

## Center to Improve Project Performance

# Evaluating Special Education Programs Resource Toolkit

### Authors

Jill Lammert, Westat

Sarah Heinemeier, Compass Evaluation & Research

Jennifer M. Schaaf, Westat

Thomas A. Fiore, Westat

Bethany Howell, Compass Evaluation & Research



**November 2016**

See *Evaluating Special Education Preservice Programs Resource Toolkit* (Lammert, Heinemeier, Schaaf & Fiore, 2016) for guidance specific to Personnel Development Project Evaluations.

For easy reference, the [Quick Start Guide](#) on page iv lists a number of questions an evaluator may ask regarding the design or conduct of a project evaluation, with links to the specific section in the Toolkit where related information may be found. A [Table of Contents](#) with hyperlinks to the different sections is on page v.

Prepared for:  
Office of Special Education Programs  
U.S. Department of Education  
Washington, D.C.

Prepared by:  
Westat  
1600 Research Boulevard  
Rockville, Maryland 20850-3129  
(301) 251-1500

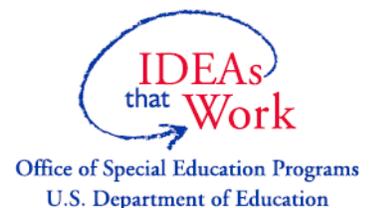
## About this Toolkit

This Toolkit was developed as part of the Center to Improve Project Performance (CIPP) operated by Westat for the U.S. Department of Education, Office of Special Education Programs (OSEP). The authors thank the OSEP, Westat, and CEEDAR (Collaboration for Effective Educator Development, Accountability, and Reform) staff who provided input.

Suggested Citation:

Lammert, J. D., Heinemeier, S., Schaaf, J. M., Fiore, T.A., & Howell, B. (2016). *Evaluating special education programs: Resource Toolkit*. Rockville, MD: Westat.

The Center to Improve Project Performance has been funded with Federal funds from the U.S. Department of Education, Office of Special Education Programs, under contract number ED-OSE-13-C-0049. The project officer is Dr. Patricia Gonzalez. The content of this publication does not necessarily reflect the views or policies of the U.S. Department of Education nor does mention of trade names, commercial products, or organizations imply endorsement by the U.S. government.



# Overview of the Center to Improve Project Performance

First formed in 2008, CIPP’s overall mission is to advance the rigor and objectivity of evaluations conducted by or for OSEP-funded projects so that the results of these evaluations can be used by projects to improve their performance and used by OSEP for future funding decisions, strategic planning, and program performance measurement.

CIPP provides evaluation support, oversight, and technical assistance (TA) to grantees of OSEP projects. CIPP staff work with project and OSEP staff to refine project logic models and develop evaluations. Based on the evaluation design and plan, CIPP staff have overseen evaluation activities and provided TA, as needed, to grantees including selecting samples; developing draft instruments; monitoring data collection and performing reliability checks; analyzing study data; providing accurate descriptions of the methods and valid interpretations of findings; and organizing, reviewing, and editing project evaluation reports. In addition to providing TA to OSEP-funded projects on request, CIPP staff prepare a variety of TA products and briefs focused on evaluation issues, and deliver presentations on evaluation through webinars and conferences.

## Contact Information:

<https://www.cippsite.org/SitePagesPublic/Home.aspx>

1-888-843-4101

Thomas Fiore, CIPP Project Director

Westat

[ThomasFiore@westat.com](mailto:ThomasFiore@westat.com)

Jill Lammert, CIPP Deputy Project Director

Westat

[JillLammert@westat.com](mailto:JillLammert@westat.com)

Patricia Gonzalez, Project Officer

Office of Special Education Programs

U.S. Department of Education

[Patricia.Gonzalez@ed.gov](mailto:Patricia.Gonzalez@ed.gov)





# Quick Start Guide

Instructions: Click on a specific Toolkit Section to go to that section.

	Toolkit Section
<b>First Things to Consider for Evaluations</b>	
What are the <b>benefits</b> of evaluations?	<a href="#">1.1</a>
What <b>components</b> do evaluations typically contain?	<a href="#">1.2</a>
How do I <b>budget</b> for my evaluation?	<a href="#">1.3</a> ; <a href="#">Appendix A.1</a>
What should I think about if I am considering working with a <b>third-party evaluator</b> ?	<a href="#">1.4</a>
<b>Creating a Theory of Change and Logic Model</b>	
How do I <b>identify the needs</b> my project will address?	<a href="#">2.1</a>
What <b>outcomes</b> should I focus on?	<a href="#">2.1.3</a>
How do I create a <b>theory of change</b> ?	<a href="#">2.2</a>
How do I develop a high-quality <b>logic model</b> ?	<a href="#">2.3</a>
Is there a <b>sample logic model</b> for reference?	<a href="#">Figure 2</a>
Is there a <b>logic model template</b> that I can use?	<a href="#">Appendix A.3</a>
How do I develop <b>evaluation questions</b> ?	<a href="#">2.4</a>
<b>Creating an Evaluation Plan</b>	
What elements should my <b>evaluation plan</b> contain?	<a href="#">2.5</a>
Is there an <b>evaluation plan template</b> that I can use?	<a href="#">Appendix A.4</a>
Is there a <b>sample evaluation plan</b> ?	<a href="#">Appendix A.5</a>
<b>Selecting a Sample &amp; Choosing an Evaluation Design</b>	
How do I <b>select my study sample</b> ?	<a href="#">2.5.3</a> ; <a href="#">4.2</a>
What <b>evaluation design</b> should I use?	<a href="#">2.5.1</a> ; <a href="#">4.1</a>
How can I use <b>goal attainment scaling</b> to measure my project outcomes?	<a href="#">4.3.5</a>
How can I use <b>single case designs</b> in my project?	<a href="#">4.1.2</a>
<b>Planning for Data Collection</b>	
What should I consider when <b>planning for data collection</b> ?	<a href="#">2.5.4</a>
Do I need to get approval from an <b>Institutional Review Board (IRB)</b> ?	<a href="#">3.1.1</a>
What do I need to consider when planning <b>data collection in schools</b> ?	<a href="#">3.1.2</a>
Are there <b>sample district, school and individual notification letters</b> that I can use?	Appendix <a href="#">C.1</a> , <a href="#">C.2</a> , <a href="#">C.3</a> , & <a href="#">C.4</a>
How do I get <b>consent</b> from participants?	<a href="#">3.2.2</a>
Are there sample <b>consent forms</b> that I can use?	Appendix <a href="#">C.5</a> , <a href="#">C.6</a> , <a href="#">C.7</a> , & <a href="#">C.8</a>
How do I obtain <b>access to student data</b> ?	<a href="#">3.1.3</a>
<b>Selecting Data Collection Methods</b>	
What should I consider when <b>creating a new survey or adapting an existing one</b> ?	<a href="#">4.3.1</a>
How do I get a <b>good response rate to my surveys</b> ?	<a href="#">4.3.1.4.1</a>
What should I consider when I want to include <b>observations</b> in my data collection?	<a href="#">4.3.2</a>

	Toolkit Section
What should I consider when I want to include <b>interviews</b> in my data collection?	<a href="#">4.3.3</a>
What should I consider when I want to include <b>focus groups</b> in my data collection?	<a href="#">4.3.4</a>
What should I consider if I want to <b>measure fidelity</b> of my project?	<a href="#">4.5</a>
Is there a <b>fidelity matrix template</b> I can use?	<a href="#">Appendix A.7</a>
<b>Recruiting and Tracking Participants</b>	
What strategies can help me to <b>recruit study participants</b> ?	<a href="#">3.2</a>
How do I <b>keep track of study participants</b> ?	<a href="#">3.2.1</a>
What elements should my <b>data tracking system</b> contain?	<a href="#">3.4.1</a>
<b>Planning for Data Analysis</b>	
How should I <b>plan for data analysis</b> ?	<a href="#">2.5.2</a>
Is there a <b>data analysis plan template</b> that I can use?	<a href="#">Appendix A.6</a>
<b>Conducting Data Analysis</b>	
How do I <b>assess the quality of the data</b> I've collected?	<a href="#">3.4.2</a>
What do I do when I have large amounts of <b>missing data</b> ?	<a href="#">4.4.1</a>
What do I need to consider <b>before I begin analyzing my data</b> ?	<a href="#">3.5</a>
How do I <b>combine data from multiple sources</b> to generate results?	<a href="#">3.5.2</a>
How do I know <b>which statistical test to use</b> when analyzing my quantitative data?	<a href="#">4.4.2</a>
How do I analyze data from <b>single-case designs</b> ?	<a href="#">4.4.3</a>
How do I ensure that my <b>qualitative data analysis</b> produces high-quality findings?	<a href="#">4.4.4</a>
How do I <b>display my qualitative findings</b> ?	<a href="#">4.4.4.3</a>
<b>Reporting Findings</b>	
How can I use data to <b>inform on-going project implementation</b> ?	<a href="#">3.6.1</a>
What elements should my <b>evaluation report</b> contain?	<a href="#">3.6.2</a>
<b>Finding Additional Resources</b>	
What are some good <b>resources on program evaluation and research</b> ?	<a href="#">Appendix D</a>

# Evaluating Special Education Programs Resource Toolkit

---

## Contents

Introduction .....	1
1 Evaluation Basics .....	2
1.1 Benefits of Evaluation .....	2
1.2 Basic Components of Evaluation .....	3
1.3 Budgeting for an Evaluation .....	4
1.4 Working with a Third-Party Evaluator .....	5
2 Planning the Evaluation .....	7
2.1 Identifying Needs .....	7
2.1.1 Constructing a Needs Statement .....	7
2.1.2 Identifying Project Resources .....	9
2.1.3 Identifying High-Quality Outcomes and Measures .....	9
2.2 Creating a Theory of Change .....	14
2.3 Creating a Logic Model .....	16
2.4 Identifying Evaluation Questions .....	18
2.5 Developing an Evaluation Plan .....	19
2.5.1 Selecting an Evaluation Design .....	19
2.5.2 Creating a Data Analysis Plan .....	25
2.5.3 Identifying a Study Sample .....	27
2.5.4 Preparing a Data Collection Plan .....	28
3 Conducting the Evaluation .....	33
3.1 Obtaining Permission to Carry Out Evaluation Activities .....	33
3.1.1 Getting IRB Approval .....	33
3.1.2 Securing District and School Approval .....	34
3.1.3 Obtaining Access to Secondary Data .....	35
3.2 Recruiting Study Participants .....	36
3.2.1 Keeping Track of Participants .....	37
3.2.2 Obtaining Participants' Consent .....	37
3.3 Managing Data Collection .....	39
3.3.1 Creating a Data Tracking System .....	39

3.3.2	Assessing Data Quality.....	40
3.4	Analyzing the Data.....	43
3.4.1	Preparing the Data for Analysis.....	43
3.4.2	Aggregating Data and Reporting Results.....	44
3.5	Reporting Findings.....	46
3.5.1	Providing Formative Feedback.....	46
3.5.2	Preparing the Final Report.....	47
3.5.3	Outlining Study Limitations.....	47
4	Methodological Considerations.....	49
4.1	Evaluation Design.....	49
4.1.1	Randomized Experimental Designs.....	49
4.1.2	Quasi-Experimental Designs.....	51
4.1.2	Single-Case/Single Subject Designs.....	57
4.1.3	Non-Experimental Designs.....	60
4.2	Sampling/Participant Selection.....	62
4.2.1	Power Analysis.....	62
4.2.2	Random Sampling.....	64
4.2.3	Purposeful Sampling.....	64
4.3	Data Collection Methods.....	66
4.3.1	Surveys.....	66
4.3.2	Observations.....	75
4.3.3	Individual Interviews.....	84
4.3.4	Focus Groups.....	87
4.3.5	Goal Attainment Scaling (GAS).....	88
4.4	Data Analysis Methods.....	94
4.4.1	Dealing with Missing Data.....	94
4.4.2	Quantitative Analysis.....	94
4.4.3	Analysis of Data in Single-Case Designs.....	100
4.4.4	Qualitative Analysis.....	104
4.5	Measuring Fidelity.....	110
4.5.1	Identify Key Components.....	111
4.5.2	Create Operational Definitions and Identify Indicators.....	112
4.5.3	Select Data Sources and Measures.....	113
4.5.4	Establish Fidelity Thresholds and Set Scores for “Adequate Fidelity”.....	114
4.5.5	Calculate Fidelity Based on Observed Data.....	116

4.5.6 Making Changes to the Fidelity System.....	119
Appendix A. Worksheets/Templates .....	121
A.1. Evaluation Cost Consideration Worksheet.....	122
A.2. Checklist for Constructing Outcomes .....	126
A.3. CIPP Logic Model Template .....	128
A.4. CIPP Summative Evaluation Plan Template.....	132
A.5. Sample Evaluation Plan: Graduate Performance and Student Outcomes .....	138
A.6. Data Analysis Plan Template .....	146
A.7. Fidelity Matrix Template.....	148
Appendix B. Validity Threats.....	150
Appendix C. Sample Forms .....	152
C.1 Sample Notification Letter for Districts with Research Approval Office/Department.....	153
C.2. Sample Request Letter for Districts without Research Approval Office/Department .....	154
C.3. Sample District Response Form .....	155
C.4. Sample School Notification Letter .....	156
C.5. Sample Passive Consent Form for PDP Graduates .....	157
C.6. Sample Active Consent form for PDP Graduates.....	158
C.7. Sample Passive Consent Form for Students .....	159
C.8. Sample Active Consent Form for Students.....	160
Appendix D. Recommended Readings on Research/Evaluation Methodology.....	162
Appendix E. Works Cited.....	166

## Tables

Table 1. Benefits and Limitations of Working with a Third-Party Evaluator.....	5
Table 2. Examples of Needs Statements.....	8
Table 3. Possible outcomes of interest, data sources, and data collection methods.....	10
Table 4. Examples of Non-Specific versus Specific Need and Outcome Statements.....	11
Table 5. Examples of Poor versus Good Outcome-Data Matches .....	12
Table 6. Example Theory of Change for an Early Intervention Project to Improve Social-Emotional Development .....	15
Table 7. Sample Formative and Summative Evaluation Questions for Four OSEP Program Areas .....	18
Table 8. Table Shell Example 1. Amount of time spent on the targeted classroom activity among PDP program graduates .....	26
Table 9. Table Shell Example 2. Amount of time spent on the targeted classroom activity, by type of PDP graduate .....	26
Table 10. Table Shell Example 3. Testing the Statistical Significant of Differences between Treatment and Control Groups. ....	27
Table 11. Gantt Chart of the Project’s Data Collection Schedule .....	32
Table 12. Sample Survey Development Framework for Evaluation of Teacher Training to Improve Literacy.....	69
Table 13. Characteristics of a Good Survey Item, with Examples of Not Good and Good Items .....	70
Table 14. Sample Calculation of Cohen's Kappa Agreement.....	83
Table 15. Common Parametric and Non-Parametric Statistical Tests.....	97
Table 16. Qualitative Data Matrix for Parent Information Center Showing Responses for Parents .....	107
Table 17. Qualitative Data Matrix Comparing Responses at Three Parent Information Centers.....	107
Table 18. Possible Operational Definitions and Indicators for the Parent TA Center Example.....	112
Table 19. Possible Data Sources for the Parent TA Center Example .....	113
Table 20. Sample Fidelity Thresholds for a Parent Resource and TA Center – Example A.....	115
Table 21. Sample Fidelity Thresholds for a Parent Resource and TA Center – Example B .....	116
Table 22. Relationship between Fidelity and Achievement of Outcomes.....	118
Table 23. Using Data to Modify Fidelity Thresholds .....	119

## Figures

Figure 1. Costs in Relation to Evaluation Scope and Complexity.....	4
Figure 2. Sample Logic Model for a Hypothetical Parent Resource and TA Center.....	17
Figure 3. Sample Consort Diagram.....	50
Figure 4. A-B-A-B Single-Case Design.....	57
Figure 5. Multiple Baseline Single-Case Design .....	59
Figure 6. The Process of Answering a Survey Question.....	70
Figure 7. Observation Protocol Example 1.....	76
Figure 8. Observation Protocol Example 2.....	76
Figure 9. Sample GAS Ratings for a Special Education Student.....	91
Figure 10. Decision Tree for Inferential Statistics .....	98
Figure 11. Calculating Percentage of Non-Overlapping Data Points .....	102
Figure 12. Sample Logic Model for a Hypothetical Parent Resource and Technical Assistance Center.....	111

## Boxes

Box 1. Study Limitations Survey.....	48
Box 2. Determining which Inferential Test to Use .....	99

# Introduction

This Toolkit is intended to support the design and conduct of high quality evaluations of projects funded by the U.S. Department of Education’s (Department) Office of Special Education Programs (OSEP). It’s not intended to provide an in-depth discussion of every aspect of project evaluation. Rather, it’s designed to highlight important issues to consider when planning and conducting a project evaluation—and to offer a variety of resources related to evaluation, including sample templates and forms, methodological information, and recommended readings on research and evaluation.

The Toolkit is organized into four sections, plus appendices:

1. [Evaluation Basics](#)
2. [Planning the Evaluation](#)
3. [Conducting the Evaluation](#)
4. [Methodological Considerations](#)

The [Quick Start Guide](#) on page iv lists a number of questions project staff or evaluators might ask regarding the design or conduct of a project evaluation, with links to the specific section in the Toolkit where related information may be found. A [Table of Contents](#) with hyperlinks to the different sections is on page vi.

In [Section 1](#) we present benefits and basic components of evaluations, as well as general considerations related to evaluations, such as how to budget for an evaluation and working with a third-party, or external, evaluator. [Section 2](#) gives a framework for planning an evaluation that can provide evidence of a project’s progress toward achieving its outcomes. This is followed in [Section 3](#) by information on some of the technical aspects of conducting an evaluation, including data collection, management, aggregation and analysis, and reporting. Finally, in [Section 4](#), we provide additional detail on specific methodological considerations.

# 1 Evaluation Basics

## 1.1 Benefits of Evaluation<sup>1</sup>

Evaluation is an important part of implementing a project. Specifically, evaluations can help:

- Identify what's working and what's not,
- Show what the project is doing and the positive outcomes that result,
- Improve staff effectiveness,
- Add to the existing knowledge base about what works (or what doesn't work),
- Fulfill funding requirements, and
- Advocate for additional funding.

Evaluations can provide project implementers with evidence to make decisions about project improvements, expansion, and sustainability; assess efficiency and guide cost-containment strategies; and facilitate replication in other settings. Evaluation results can indicate whether and how well a project is achieving its expected outcomes and identify particular aspects that are more or less effective. At the same time, knowing what is and isn't working allows project staff to make improvements and distribute resources to best effect, thereby increasing project efficiency. In turn, this increased efficiency may offset the costs of evaluation. Evaluation results also can have more widespread benefits by adding to the knowledge base about particular projects or interventions, thereby informing provision of services and pursuit of specific outcomes in other contexts and locations. And, conducting a high-quality evaluation that demonstrates project effectiveness and results can be vital in advocating for maintained or increased funding, project scale-up or expansion.

Nevertheless, some project staff may hesitate to engage in evaluation because of concerns that it's complicated or burdensome and worries that negative (or even neutral) results could be detrimental to the project. However, these concerns can be alleviated by careful planning coupled with thoughtful consideration of the benefits of evaluation and consideration of how to use findings from an evaluation to improve a project and its outcome achievement. In addition, since many project managers may already be engaged in an informal evaluation process, conducting an evaluation may not be as difficult or burdensome as some might fear.

---

<sup>1</sup> Adapted from The Program Manager's Guide to Evaluation 2<sup>nd</sup> Edition  
[http://www.acf.hhs.gov/sites/default/files/opre/program\\_managers\\_guide\\_to\\_eval2010.pdf](http://www.acf.hhs.gov/sites/default/files/opre/program_managers_guide_to_eval2010.pdf)

## 1.2 Basic Components of Evaluation

Evaluation should be an integral part of any project; therefore, a comprehensive evaluation plan should be developed at the same time as—or as part of—the overall project plan. Generally, the purposes of an evaluation are to demonstrate how well the project components have been implemented and to analyze the extent to which the project’s objectives and outcomes have been achieved. Evaluations typically feature two components: formative evaluation and summative evaluation.

**Formative evaluation.** In general, formative evaluations address how well a project or intervention is being implemented, including the nature of the activities, products, and support provided by project staff; the structures, policies, and procedures influencing implementation and outcomes (e.g., contextual evidence); fidelity to the project model; changes that may be necessary to improve project implementation; the ways the project is being perceived by key stakeholders (e.g., quality, relevance, usefulness, social validity); and progress toward achieving outcomes. Formative data can come from a variety of sources (e.g., project records, stakeholders) and be collected through many methods (e.g., surveys, interviews, observations, document review). Formative data are generally collected to answer questions that relate to:

1. **monitoring progress** toward carrying out activities, producing outputs, and achieving the short-term outcomes identified in a project’s logic model;
2. **social validity**—the social importance and acceptability of the project or intervention, such as the social significance of the project or intervention goals, the social appropriateness of the intervention procedures and the social importance of the intervention outcomes<sup>2</sup>; and
3. **fidelity**—in the context of formative evaluation, fidelity data relate to whether and to what degree the project is carrying out evidence-supported strategies and activities and producing outputs as intended (i.e., in the expected amounts and covering the expected content).

Examples of **formative questions** include:

- Is the project achieving milestones and benchmarks in a timely manner?
- Is the project in compliance with the federal priorities?
- Is project staffing sufficient in numbers and competencies?
- Are resources adequate to support project activities?
- How many persons have participated in activities or received services?
- Is the project completing planned activities and producing the expected outputs?
- Are there any implementation gaps or project support needs?
- Are there any facilitators of or barriers to implementation?
- Are project activities and outputs being implemented with fidelity?
- Do all signs point to achieving desired outcomes?

**Summative evaluation.** The overall purpose of summative evaluation is to evaluate the effectiveness and efficiency of a project in achieving its outcomes or goals. Further, a summative evaluation can establish a project’s impact on the populations served or affected by the project, including students. This is accomplished in part through progress monitoring and formative evaluation, but summative evaluation questions typically require an investigation of the

---

<sup>2</sup> Shaver, D., Wagner, M., Nagle, K., & Ryan, T. (2015). *Improving implementation of programs and practices for children with disabilities: Lessons learned from the Model Demonstration Coordination Center*. Menlo Park, CA: SRI International. Retrieved from [http://mdcc.sri.com/documents/MDCC\\_Final\\_Report\\_SEPT2015.pdf](http://mdcc.sri.com/documents/MDCC_Final_Report_SEPT2015.pdf)

extent to which a change has occurred, the factors associated with a change, or the measurement of change among different populations. An important role for the summative evaluation can be the determination of the unique contribution of the project to the desired change. As such, summative questions are best informed when there are comparison data (e.g., for treatment and control groups, or data collected from multiple points in time) to give the evaluator an idea of the counterfactual—that is, what would have happened if the project hadn't been implemented.

Examples of **summative questions** include:

- What outcomes (expected and unexpected) have occurred?
- What expected outcomes have not occurred?
- To what degree have outcomes occurred?
- Where is change the greatest?
- What is the unique contribution of the program to the observed change?
- What is the cost/benefit of these outcomes?
- To what extent do the same outcomes occur in treatment and control groups, comparison groups, or the same group over time?

As we'll discuss in the next section, when planning an evaluation it's important to consider the budget and other resources (e.g., personnel, instruments, incentives, travel) that will be needed to complete the evaluation.

### 1.3 Budgeting for an Evaluation<sup>3</sup>

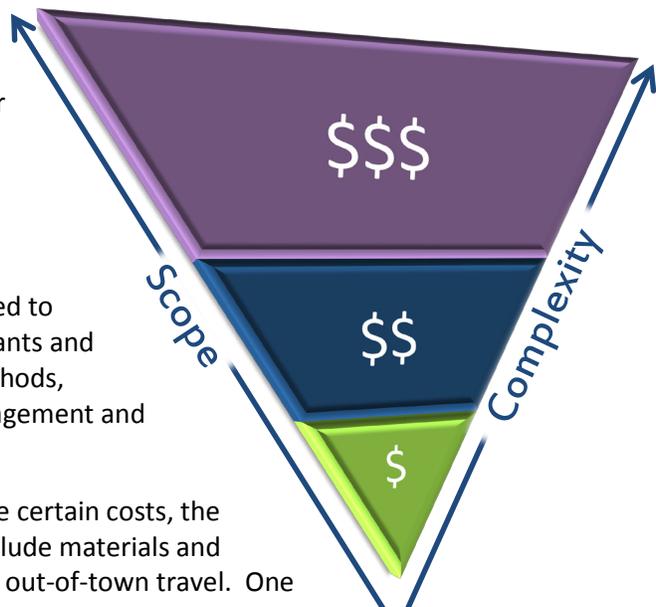
Accurately budgeting for an evaluation requires an understanding of the evaluation process and of the various factors that might influence costs. Evaluation costs can vary widely, and, as illustrated by Figure 1 below; the amount of money needed generally increases with the scope and complexity of both the project to be evaluated and the evaluation itself.

**Scope** refers to factors such as the size and reach of the project to be evaluated, the evaluation focus, the number of evaluation questions, the evaluation time period, whether and how stakeholders will be involved, and the number and type of reports that will be required.

**Complexity** refers to factors such as the nature of the evaluation questions, the type of evaluation design needed to answer each evaluation question, the number of participants and sites to be included in the evaluation, data collection methods, frequency and duration of data collection, and data management and analysis requirements.

No matter the scope or complexity, all evaluations require certain costs, the most significant of which is staffing. Other likely costs include materials and supplies, computer equipment and software, and local or out-of-town travel. One

Figure 1. Costs in Relation to Evaluation Scope and Complexity



<sup>3</sup> Content adapted from Lammert & Fiore (2015), available on the OSEP IDEAs That Work website: [https://www.osepideasthatwork.org/sites/default/files/CIPP2\\_Budgeting\\_for\\_Evaluation\\_Brief\\_2-13-15.pdf](https://www.osepideasthatwork.org/sites/default/files/CIPP2_Budgeting_for_Evaluation_Brief_2-13-15.pdf)

of the elements that can increase the cost of the evaluation is working with a third-party, or external, evaluator. However, as we discuss in the next section, there are many reasons project staff might want to do so.

Ideally, evaluation costs should be estimated in absolute dollar amounts, after carefully considering scope and complexity in the context of the specific evaluation needs. The *Evaluation Cost Considerations Worksheet* included in [Appendix A.1](#) can provide guidance regarding the different factors affecting costs. In the worksheet, evaluation elements are assigned relative costs based on how they may affect staffing, supplies, travel, etc. The list of evaluation elements presented in the worksheet isn't exhaustive and the different categories aren't mutually exclusive, but the worksheet can help project staff to carefully think through the multiple factors affecting evaluation costs.

If many factors are unknown or a rough estimate is needed early in the planning process, the cost of different types of evaluations can be roughly estimated in terms of a percentage of the program or project budget, as follows:

- Low cost = <10%
- Moderate cost = 10-20%
- High cost = >20% of budget.

## 1.4 Working with a Third-Party Evaluator

When formative questions are the primary focus of an evaluation, project staff often assume responsibility for collecting and analyzing data to track progress and provide formative feedback. Summative evaluations, however, may benefit when conducted by an individual or organization that is external to the project team, and funders may even require independent evaluations of project activities (such as in the case of the [Department's Investing in Innovation program](#)).

Check out the *Guidelines for Working with Third-Party Evaluators* on the [OSEP IDEAs That Work website!](#)

Third-party, or external, evaluators can offer increased expertise and objectivity, as well as in-depth knowledge and experience with evaluation and research methodology. However, when considering working with a third-party evaluator, it's important to be aware of the potential benefits and limitations of this working relationship, as shown in Table 1. Benefits include the needed skills and independence brought to the project by the third-party evaluator, while limitations include possible, often unforeseen or unplanned, tasks or costs associated with monitoring and managing the work of the third-party evaluator.

**Table 1. Benefits and Limitations of Working with a Third-Party Evaluator**

Benefits	Limitations
<p><b>Third-party evaluators can:</b></p> <ul style="list-style-type: none"> <li>• Bring technical expertise in research methodology, statistics, or related topics to the project team</li> <li>• Provide credibility and objectivity by acting as an external "critical friend"</li> <li>• Take on responsibility for completing some or all of the (formative and summative) evaluation tasks, allowing project staff to focus on project implementation</li> </ul>	<p><b>Third-party evaluators may:</b></p> <ul style="list-style-type: none"> <li>• Add unanticipated or additional cost to the project</li> <li>• Add to project monitoring and management tasks focused on the work of contractors</li> <li>• Not know the project background or content area as well as project staff</li> <li>• Be less available or accessible, as compared to project staff</li> </ul>

The decision of when to hire the third-party evaluator affects what the evaluator can and cannot provide to the project. The third party evaluator can benefit the evaluation most when hired early on (such as during the planning stages) so that they can provide guidance and assistance on foundational aspects of the evaluation—its design and methods. A third-party evaluator brought on board towards the end of the project may be hindered by decisions made earlier in the project. Specifically, they may have limited ability to revise, modify, or correct decisions made earlier in the evaluation (such as the timing and types of data collections), thereby limiting the conclusions that can be developed from the evaluation results.

It's important to keep in mind that, even when the third-party evaluator has a significant role in the project, the Project Director (or Principal Investigator) bears ultimate responsibility for ensuring that the project and its evaluation are carried out as planned and that all requirements for project implementation and reporting are met. Therefore, frequent communication with the evaluator is needed, as well as some time spent monitoring project progress.

For additional information about working with a third party evaluator such as how to select and hire an evaluator (including how to prepare and execute a contract), how to monitor and manage the work of the evaluator, and how to conclude the evaluation project, see the [Guidelines for Working with Third-Party Evaluators](#)<sup>4</sup> available on the OSEP IDEAs That Work website.

---

<sup>4</sup> Heinemeier, D'Agostino, Lammert, & Fiore, 2014

## 2 Planning the Evaluation

As discussed above, a well-designed evaluation can have many benefits, but such an evaluation requires careful planning. Whether or not project staff choose to conduct the evaluation in-house or to work with a third-party evaluator, the project team is best acquainted with the project details and therefore will play a critical role in the evaluation planning process. In this section we'll offer a framework for planning an evaluation and an overview of important decisions that need to be made during the planning process.

### 2.1 Identifying Needs

Funding opportunities, such as Federal Register notices, Requests for Proposals, and Requests for Applications, often identify specific needs that funders expect projects to address. In some cases these needs have been identified through a formal needs assessment. In other cases, information on needs stems from sources such as media reports, study findings, policy briefs, or stakeholder feedback. Since the needs the project aims to address will inform the goals and expected outcomes of the evaluation, it's important for project staff to understand and document the following:

- **The level or scope of the need in the community or target population of interest.** Needs data can be found in numerous places such as databases of statistical evidence or other research- or evaluation-based repositories. Evidence related to needs can be gathered by talking with key stakeholders or reviewing past studies for information about the barriers and facilitators to implementation of a particular project or intervention in a given context or setting. Ideally, needs data are high quality, meaning that the data were collected or aggregated using methods that assured the precision, accuracy, reliability, consistency, and completeness of the data.
- **Upwards or downwards trends in the need.** Trend data will help explain whether or not the magnitude or urgency of a need is increasing or decreasing, or if there are trends in the need among different members of the target population. Data from multiple years will help establish these trends and often are available at the same sites that provide data to inform a current understanding of the need.
- **Why the need exists in the community or target population of interest.** This aspect of the need may be more elusive than the level or scope of need. In fact, project staff may need to collect data from the project's target population to better understand why a need exists, or why the need varies among members of its target population. This may include, for example, conducting interviews, focus groups, or surveys of personnel or parents to better understand how and why a need exists in a specific community or population.
- **Need context.** Finally, it may be helpful to explore the legislative or political context that may be influencing why a need or needs exist in a community or target population or the way that projects might try to address the needs.

If no existing needs assessment data are available, project staff may need to carry out their own needs assessment or use data to construct a needs statement, as discussed below.

#### 2.1.1 Constructing a Needs Statement

While there is no one "correct" way to construct a needs statement, it's helpful to be as specific as possible regarding the nature of the need. This will help the project target its resources and help to frame the focus of the evaluation. A needs statement may vary in length from one or two sentences to one or more pages, and should contain specific, time-relevant, statements about needs or circumstances that justify the project.

Ideally, the needs statement will achieve all of the following:

- Identify the time-relevance of the need (e.g., the data are from 2015-16 and therefore are relatively current)
- Identify the need that must be addressed (e.g., children identified with Emotional/Behavioral Disorder (E/BD) struggle to adjust to inclusive classrooms)
- Identify the level of the need in number or percentage values (e.g., 45% of children identified with E/BD)
- Identify whether the need varies among demographic groups or sub-populations (e.g., comparing children participating in an exceptional children’s program with children in the general education program)

**TIP: Use concrete data when creating needs statements.** Concrete data are **specific** and allow the user to identify the level and scope of a need or project resource. Look for data that can be **quantified** in numbers and percentages, as opposed to general statements.

Table 2 below provides examples of less specific and more specific needs statements. Notice how the more specific needs statements contain information that can guide the identification of a target population and the establishment of outcome measures.

**Table 2. Examples of Needs Statements**

Less Specific Needs Statement	More Specific Needs Statement
• <b>Rural districts have a hard time retaining special education teachers.</b>	• Recent employment data from [YEAR] indicates that [PERCENT] of fully qualified special education teachers in the rural districts in which program graduates are employed leave their position within three years.
• <b>Educational materials are not accessible to children with visual impairments.</b>	• Data from [YEAR] indicates that only [PERCENT] of educational materials in elementary schools are accessible to children with visual impairments.
• <b>Children are not arriving at school ready to learn</b>	• [YEAR] data indicate that [PERCENT] of entering kindergarten students is below expectations in one or more developmental domains.
• <b>Special education students are falling behind on standardized assessments.</b>	• Data from the [YEAR] school year indicates that [PERCENT] of students receiving special education services, and who are eligible to participate in standardizing testing and attend schools in districts that we serve, fail to achieve proficiency standards on standardized state assessments of reading and mathematics.
• <b>Children with disabilities have behavior problems</b>	• [YEAR] school data indicates that the rate of behavior problems in middle school students with disabilities was [PERCENT].
• <b>Lower employment for individuals with disabilities</b>	• [YEAR] data indicate that only [PERCENT] of individuals with disabilities are able to find employment within one year after graduation from high school.

The information contained in the needs statement can then be used to inform evaluation planning by identifying key outcomes of interest and providing a timeframe for achievement of those outcomes.

### 2.1.2 Identifying Project Resources

The inputs or resources that are available to a project will have a major impact on the ability of the project to achieve its goals. Inputs may include the following types of resources and materials:

- **Funding.** One of a project's primary inputs is funding, which may come from multiple sources. Note: when planning to combine funding from different sources, keep in mind that each funding source may have its own criteria for evaluation, monitoring, and reporting. Meeting each funder's requirements can add to a project's management burden.
- **Internal resources.** These resources may include facilities and staff to support project implementation and carry out routine project tasks, such as managing the project budget and expenditures or offering secretarial support. A project team may also be able to draw upon the expertise or resources of peers within the project's sponsoring agency, such as content experts within a university system.
- **External resources.** External groups such as an Advisory Committee, Community of Practice, Professional Learning Community, or other professional community may be available to offer expertise or support to a project. In addition, there are a number of [federally-supported Technical Assistance \(TA\) Centers](#) that provide guidance and support on a range of topics.
- **Community resources and supports.** The community, district, or neighborhoods that the project will operate in may want to support the project with volunteers, space, or other resources.
- **Guidelines and protocols for project implementation.** The funding agency may publish specific guidelines or stipulations for how a sponsored project must be implemented. These may include specific cost principles or guidelines for how project funds can and cannot be spent. This might also include specific objectives and requirements for using evidence-based practices and/or interventions (See for example, OSEP's performance standards or IDEA grant requirements). These guidelines and requirements typically apply to all projects funded under an agency or funding opportunity and cannot be varied to meet an individual project's design.
- **Intervention resources** such as implementation guides, recommended or required assessments, and technical support staff.

These inputs or resources are often included in a project's [logic model](#) and may also be included in an evaluation, especially if the evaluation is trying to learn how the various inputs contribute to outcome achievement.

### 2.1.3 Identifying High-Quality Outcomes and Measures

The quality of an evaluation depends in large part on the quality of the outcomes used to demonstrate project effects. There are three levels of outcomes: short-term, intermediate (also called medium-term), and long-term.<sup>5</sup> Short-term outcomes are among the first changes that can be recognized, and they usually occur as an immediate result of program activities and investments. Short-term outcomes are the first signs that later outcomes are achievable. Intermediate outcomes typically represent the cumulative effect or the somewhat more distal effect of short-term outcomes—and they often must occur before the long-term outcomes can be achieved. In turn, long-term outcomes speak to alleviation (amelioration or elimination) of the originating need.

---

<sup>5</sup> These are also called short-term outcomes, medium-term outcomes, and distal outcomes.

The high-quality measurement of outcomes depends in part on the nature of the outcomes themselves and in part on the rigor of the data collection and analysis procedures used to generate outcome findings. The best outcomes are rigorous, have a high degree of utility, and are informed by high-quality data. Rigorous outcomes are those that are considered valid and that have been measured with a great degree of methodological precision and accuracy. An outcome with a high degree of utility is useful in informing decisions and planning. Finally, an outcome should be informed by data that have been collected using high-quality measures (see, for example, Boller et al., 2010 and Mallone et al., 2010).

Table 3 below provides examples of outcomes that OSEP projects may seek to achieve, as well as possible sources for outcomes data and methods to obtain the data. For example, if a project were interested in increasing the quality of socio-emotional interactions in families of children with disabilities through a family training program, the project might want to measure the changes in interactions using data from the families that is collected from in-home observations prior to training and after training. The table below isn't exhaustive; project staff are encouraged to identify and use outcomes that are specific to their project's activities.

**Table 3. Possible outcomes of interest, data sources, and data collection methods.**

What outcomes might be of interest?	Who might have the data?	How might I obtain the information?
<ul style="list-style-type: none"> <li>• Extent to which services or information have reached intended targets</li> <li>• Effects of outreach strategies</li> <li>• Effects of trainings (on trainees or on the groups trainees serve)</li> <li>• Changes in skills (e.g., social, academic, instructional)</li> <li>• Changes in knowledge</li> <li>• Changes in job placement rates</li> <li>• Changes in job retention</li> <li>• Changes in teacher efficacy</li> <li>• Progress toward educational goals</li> <li>• Progress toward functional goals</li> <li>• Progress toward educational goals</li> <li>• Relationship improvements (among families, communities, schools, service providers, etc.)</li> <li>• Communication improvements (among stakeholders.)</li> <li>• Interaction improvements (among families, communities, schools, service providers, etc.)</li> <li>• Increased accessibility (to media, services, etc.)</li> </ul>	<ul style="list-style-type: none"> <li>• Clients</li> <li>• Students</li> <li>• Children with disabilities</li> <li>• Families of children with disabilities</li> <li>• Teachers</li> <li>• Related-services professionals</li> <li>• School administrators</li> <li>• Adults with disabilities</li> <li>• Trainees</li> <li>• General public</li> <li>• Service providers</li> <li>• State personnel</li> <li>• Existing records</li> </ul>	<ul style="list-style-type: none"> <li>• Administer surveys</li> <li>• Ask for self-evaluations</li> <li>• Conduct interviews</li> <li>• Conduct observations</li> <li>• Administer (or review previously administered) assessments</li> <li>• Review employment records</li> <li>• Review student learning objectives</li> <li>• Review alternate student assessments</li> <li>• Review educational records</li> <li>• Review teacher portfolios</li> </ul>

In the next sections we outline one approach to determining the quality of an outcome, followed by a brief discussion of the characteristics of high-quality outcome measures.

### 2.1.3.1 The SMART Approach to Determining Outcome Quality

One way to assess an outcome’s quality is to determine if the outcome is “SMART.” This acronym was coined by George Doran<sup>6</sup> for management purposes, but has since been adapted and used by multiple authors in varied ways. For our purposes, the acronym describes objectives that are “*Specific and Clearly Stated, Measurable and Based on Data, Attainable and Realistic, Relevant, and Time-Bound.*”

- **Specific (and Clearly Stated).** A well-specified outcome statement provides sufficient detail to allow a reader to determine if the originating need or problem has been addressed. Thus, a well-specified and documented needs statement can be an important first step in creating a well-specified outcome statement. Table 4 shows examples of non-specific and specific need and outcome statements, building on some of the needs statements presented earlier in Table 2.

**Table 4. Examples of Non-Specific versus Specific Need and Outcome Statements**

Sample Need/ Outcome	Non-Specific	Specific
<b>Need #1</b>	Rural districts have a hard time retaining special education teachers.	Recent employment data from [YEAR] indicates that [PERCENT] of fully qualified special education teachers in the rural districts in which program graduates are employed leave their position within three years.
<b>Outcome #1</b>	Retention rates in rural districts.	Increase the percent of fully qualified special education teachers graduating from the program and employed in rural districts who remain in their position for more than three years to [PERCENT] by [YEAR].
<b>Need #2</b>	Educational materials are not accessible to children with visual impairments.	Data from [YEAR] indicates that only [PERCENT] of educational materials in elementary schools are accessible to children with visual impairments.
<b>Outcome #2</b>	Accessible educational materials.	Increase the percent of educational materials in elementary schools that are accessible to children with visual impairments to [PERCENT] by [YEAR].
<b>Need #3</b>	Children are not arriving at school ready to learn.	[YEAR] data indicate that [PERCENT] of entering kindergarten students is below expectations in one or more developmental domains.
<b>Outcome #3</b>	Children arrive at school ready to learn.	By [YEAR], [PERCENT] of students entering kindergarten will be at or above developmental expectations in all domains.
<b>Need #4</b>	Special education students are falling behind on standardized assessments.	Data from [YEAR] indicates that [PERCENT] of students receiving special education services, and who are eligible to participate in standardizing testing and attend schools in districts that served by the project, fail to achieve proficiency standards on standardized state assessments of reading and mathematics.
<b>Outcome #4</b>	Special education student achievement on standardized assessments.	By [YEAR], [PERCENT] of students (who are eligible to participate in standardized testing and attend schools in districts that we serve) receiving special education services from project graduates will achieve proficiency standards on standardized state assessments of reading and mathematics.

- **Measurable (and Based on Data).** Good outcomes can be measured and are based on high-quality data. During the planning period, it’s important to determine if data related to the outcomes are, or can be, available. If high-quality, outcome-specific data aren’t available through existing data collection and management systems, the team must determine if and how the data can be collected. Ultimately, if high-quality, outcome-specific

<sup>6</sup> Doran, 1981.

data aren't available and cannot be collected for a specific outcome (within the limits of the available time and resources), it's better to select another outcome for which such data *are* available than to report results for an outcome based on poor-quality data. When determining which data will be available, it's also important to consider which data will truly be able to signal outcome achievement—if the available data aren't a good match for the outcome, they won't be much help to the evaluation. Table 5 presents examples of good and poor data matches for the set of sample outcomes outlined in Table 4. Finally, whenever possible, it's good to use standardized procedures for the collection, aggregation, and analysis of relevant data.

**Table 5. Examples of Poor versus Good Outcome-Data Matches**

Sample Outcome	Poor Data Match	Good Data Match
Outcome 1: Increase the percent of fully qualified special education teachers graduating from the program and employed in rural districts who remain in their position for more than three years to [PERCENT] by [YEAR].	For districts served by the program, over time: <ul style="list-style-type: none"> <li>- District name</li> <li>- Number of special education positions</li> <li>- Annual turnover rate</li> </ul>	For districts served by the program, over time: <ul style="list-style-type: none"> <li>- Rural identifier (is the district considered rural?)</li> <li>- Number of special education staff who are graduates of the program who meet the definition of fully qualified</li> <li>- Dates that program graduates begin and end employment in the district</li> </ul>
Outcome #2: Increase the percent of educational materials in elementary schools that are accessible to children with visual impairments to [PERCENT] by [YEAR].	For district served by the project, over time <ul style="list-style-type: none"> <li>- Number of elementary schools</li> <li>- Types of accessible educational materials</li> </ul>	For districts served by the project, over time <ul style="list-style-type: none"> <li>- Number of children with visual impairments enrolled in each elementary school</li> <li>- At each school, number of educational materials accessible to children with visual impairments</li> </ul>
Outcome #3: By [YEAR], [PERCENT] of students entering kindergarten will be at or above developmental expectations in all domains.	For districts served by the program, over time: <ul style="list-style-type: none"> <li>- Number of entering kindergarteners</li> <li>- Number of kindergarteners considered ready by kindergarten teachers</li> </ul>	For districts served by the program, over time: <ul style="list-style-type: none"> <li>- Number of entering kindergarteners</li> <li>- Number of entering kindergarteners who received developmental assessments (for one or more developmental domains)</li> </ul> Of the students who received assessments: <ul style="list-style-type: none"> <li>- The number and percent who were assessed as at or above developmental expectations, for each domain assessed.</li> </ul>
Outcome #4: By [YEAR], [PERCENT] of students (who are eligible to participate in standardized testing and attend schools in districts that we serve) receiving special education services from project graduates will achieve proficiency standards on standardized state assessments of reading and mathematics.	For districts served by the program, over time: <ul style="list-style-type: none"> <li>- Standardized reading and mathematics scores</li> <li>- Number of special education students</li> </ul>	For districts served by the program, over time: <ul style="list-style-type: none"> <li>- Number of students taught by program graduates</li> <li>- Number of students taught by program graduates who participated in standardized testing</li> <li>- Standardized reading and mathematics scores</li> <li>- Definition, per student, of annual expectations (e.g., achieving proficiency standards on state assessments)</li> </ul>

- **Attainable (and Realistic).** Good outcomes reflect changes that are achievable within a given timeframe. It may be helpful to have content- and problem-area experts on the evaluation team to help decide the goals and outcomes that are achievable during the project period. For example, looking at sample Outcome 3 in Table 5, project staff implementing a personnel development program (PDP) should carefully consider whether an attainable (or realistic) outcome of their project would be to increase the percentage of children entering kindergarten who are at or above developmental expectations in all domains. The discussion of outcome relevance should address the resources that will be necessary to achieve the outcome—in other words, are there sufficient financial, tangible, and other resources in place for the project to achieve its goals and outcomes within the time frame allotted? The process of answering these questions should take place during the planning period.
- **Relevant (to Objectives).** Relevant outcomes address the degree to which the underlying need or problem has been alleviated, reflect needs and problems of consequence to communities and schools, and generate information for future decision-making. For example, for a technology and media project that aims to increase access to online learning for students with disabilities, a relevant outcome would be an increase in the availability of online learning courses that are accessible to students with visual impairments. A less-relevant outcome might be an increase in enrollment of students at postsecondary institutions in NoteTaker courses.
- **Time-Bound.** Finally, good outcomes are achievable within a defined period of time. Our experience indicates that short-term outcomes generally can be observed within one program year, whereas long-term outcomes can take the entire grant period or more to assess and achieve. It's not uncommon or inappropriate for evaluators and program staff to identify long-term outcomes that are expected to occur beyond the project period. Intermediate outcomes generally represent changes that occur between the short-term and long-term outcomes and, because they are likely to fall within the time boundaries of the grant period, may be the most distal outcomes on which the evaluation can realistically focus.

As can be seen, developing outcomes that have all these criteria requires project staff and evaluators to have at least some background knowledge of why needs or issues exist, the steps required to address these needs or issues, and the length of time necessary to complete those steps. In many situations, a framework for addressing these issues will have been developed as part of the application for OSEP funding. In cases when this information isn't available, existing research from similar or existing programs may be available. In other cases, an evaluator may use information from related fields or service areas to construct a model for how, when, and under what circumstances change may occur.

Of course, having a high-quality outcome means nothing if the data related to that outcome are collected using low-quality measures. In the next section we discuss characteristics of high-quality outcome measures.

### 2.1.3.2 Characteristics of High-Quality Outcome Measures

ED's Institute of Education Sciences' (IES) [What Works Clearinghouse \(WWC\) Procedures and Standards Handbook](#)<sup>7</sup> has identified four characteristics high-quality outcome measures—the instruments that are used to assess outcome achievement—possess:

- **Face validity**—The measure must appear to be a valid measure of the outcome (e.g., a reading fluency test shouldn't be used to measure mathematics outcomes).

<sup>7</sup> U.S. Department of Education, 2014.

- **Adequate reliability**—This depends on the type of outcome measure (test score, scale, observation measure) and whether or not the measure is based on a standardized test or state-required achievement test. The WWC recommends following these minimum standards: (a) internal consistency (such as Cronbach’s alpha) of 0.50 or higher, (b) temporal stability/test-retest reliability of 0.40 or higher, or (c) inter-rater reliability (such as percentage agreement, correlation, or kappa) of 0.50 or higher.
- **Free of over-alignment with the study intervention**—The measure must not be designed or administered in ways that are specifically aligned to an intervention so that the individuals receiving the intervention are being taught directly the content included in the outcome measure (e.g., a student shouldn’t be tested for reading fluency using the 50 words that she practiced reading aloud repeatedly during an intervention).
- **Consistent data collection across groups**—The outcome data must be collected using the same rules or procedures across groups of participants in the study (e.g., student outcome data shouldn’t be collected by special education teachers as part of their regular classroom activities in one school and by graduate research assistants in a pull-out activity in another school). If data aren’t collected consistently, the evaluator should try to ascertain to what degree the differences in data collection may contribute to differences in outcomes across groups.

Another characteristic of a high-quality measure is that it is **consistently defined across groups**. For example, if an evaluation is looking at the performance of transition specialists and wants to use rates of college enrollment among students who are deaf/hard of hearing (DHH) as an outcome, it’s important to know whether the college enrollment rate for DHH students is calculated in similar ways across schools and school districts. If the outcome isn’t consistently defined across groups it’s difficult for an evaluator to know the extent to which those differences in definition may account for variation in outcomes. Consequently, when differences in definition exist, the evaluator should try to identify ways that those differences may influence the measurement of outcomes.

Of course, not all outcome measures are going to meet every criteria for quality presented above. However, evaluators should select outcome measures that fulfill as many of the quality criteria as possible. [Appendix A.2](#) presents a checklist developed for this Toolkit containing questions that evaluators can ask when constructing outcomes and identifying measures. Evaluators are encouraged to carefully consider the impact of any items that receive a “no” or “somewhat” response on the evaluation’s ability to generate meaningful data regarding the project’s progress and results. These considerations can be summarized and reported as study limitations in the evaluation report. See [Section 3.5.3](#) for a discussion of study limitations.

## 2.2 Creating a Theory of Change

A theory of change (also sometimes called a theory of action) outlines a project’s goals and objectives, and the various processes and contextual influences (including moderators and confounding factors) that are expected to influence or contribute to outcome achievement. While project staff may not have outlined a formal theory of change, the team undoubtedly has an idea of how the project is expected to address the identified needs. The key is to document the mechanisms by which the various actors, change processes, and external influences are expected to combine to produce the expected results. A theory of change might be relatively straightforward or it might be quite complicated—especially if the desired outcomes may be influenced by a variety of factors. The theory of change will identify the project’s specific inputs and activities and link them to measurable variables such as project outputs, service population characteristics, and individual outcomes. Project staff generally focus project activities on the specific parts of the theory of change they think they can influence—such as changing teacher behaviors to improve classroom instruction, recognizing that there are a variety of other factors beyond the scope of the project that can affect student performance

(for more information on using evidence to select project activities, see CIPP’s *Demonstrating Evidence across the Project Cycle*, available on the [OSEP IDEAs That Work website](#)). In some cases, project staff create theories of change as a series of “if-then” statements, outlining how they expect their project to work.

Table 6 below shows a sample “if-then” theory of change that could be used for a state-wide project to improve social-emotional outcomes of infants and toddlers with disabilities. The different types of action strands that will be taken by the project appear in the left column, the more specific activities in the next column and then the short- intermediate- and long-term outcomes appear in subsequent columns. Recall that these simple implementation of these project activities doesn’t guarantee outcomes, which is why it’s important to conduct an evaluation that monitors all steps of the project—including fidelity—to ensure things are going as planned (for more information on how to measure fidelity, see [Section 4.5](#)).

**Table 6. Example Theory of Change for an Early Intervention Project to Improve Social-Emotional Development**

Action Strands	If the lead agency	Then	Then	Then
<b>Increased Accountability</b>	<ul style="list-style-type: none"> <li>– Establishes a quality improvement system to enhance and monitor early intervention evidence-based practices focused on social-emotional development.</li> <li>– Builds local program capacity to report accurate data on social-emotional outcomes and uses that data for monitoring.</li> </ul>	<ul style="list-style-type: none"> <li>– Staff will have increased awareness of fidelity standards and program performance for implementation of evidence-based practices in early intervention.</li> <li>– Early intervention program managers will have the skills to use outcome data for program improvement.</li> </ul>	Early intervention providers will implement evidence-based practices related to social-emotional development with fidelity.	Infants and toddlers with disabilities will exit early intervention services with an increased rate of growth in positive social-emotional development.
<b>Fiscal Support</b>	<ul style="list-style-type: none"> <li>– Increases funding to hire enough qualified staff, provide appropriate training, and support service delivery</li> </ul>	<ul style="list-style-type: none"> <li>– Early intervention providers will have the capacity and resources necessary to provide evidence based practices and supports.</li> </ul>		
<b>Professional Development</b>	<ul style="list-style-type: none"> <li>– Provides training in evidence-based practices to support the social-emotional development of infants and toddlers</li> </ul>	<ul style="list-style-type: none"> <li>– Early intervention providers will better understand how to support social-emotional development for infants and toddlers.</li> </ul>		

## 2.3 Creating a Logic Model

A logic model is a visual representation of and organizational structure for the [theory of change](#).<sup>8</sup> Building upon the theory of change, a logic model provides specific detail about the mechanisms by which the project will achieve the desired outcomes. A logic model is an important tool for developing an effective evaluation plan because it provides a good way to visualize the inputs, strategies, activities, and outputs necessary to respond to needs and make progress towards a desired outcome. The logic model helps (1) define outcomes that are meaningfully connected to project activities and (2) support evaluations so that the process will improve a project's overall performance. In general, a logic model should be precise and include features and content that support and promote the project's evaluation, but there are many different variations on logic models (see, for example, the [OSEP Logic Model Outline](#), the [W.K. Kellogg Foundation Logic Model Development Guide](#), and Frechtling, 2007).

The CIPP Logic Model Template provided in [Appendix A.3](#) is designed to portray a project's overall plan and clarify the relationships among a project's goals, strategies and activities, outputs, and projected outcomes. Inputs and external factors are also included. These elements of the logic model are briefly described below.

- **Goals/Objectives**—Goals capture the overarching purposes of the project. Goals make clear the anticipated impact on systems or individuals. Goals imply gaps or deficits that will be remedied when the project produces its long-term outcomes. Objectives, if used in a logic model, are targeted sub-goals.
- **Inputs**—Inputs include the [resources](#) that are available to the project. This includes external funding, internal resources, and intangibles such as experience and the state of the knowledge in the field.
- **External Factors**—External factors relate to the context in which the project is being implemented. This may include other federal initiatives, the OSEP policy environment, an institution's accumulated experience and visibility, and public demand for or receptiveness to the project.
- **Strategies/Activities**—Strategies are the broad approaches to addressing the goals and generally include multiple activities. Activities are the specific actions funded by the grant or supported by other resources under the umbrella of the project.
- **Outputs**—Outputs are the short-term results of the project activities, including project products and programs. Most outputs will be quantifiable, including tallies of the number of products and programs or counts of the customer contacts with those products and programs.
- **Short-term/Intermediate Outcomes**—Short-term outcomes are what customers do or become as a result of outputs. Usually, short-term outcomes are changes in knowledge or skills acquired through project outputs. Intermediate outcomes result either directly from outputs or indirectly through short-term outcomes. Often, intermediate outcomes are changes in the behavior or practices of persons touched by the project. They generally come later in time than short-term outcomes and often represent a step between short-term outcomes and long-term outcomes.
- **Long-term Outcomes**—Long-term outcomes are the broadest project outcomes and follow logically from the short-term and intermediate outcomes. They are the results that fulfill the project's goals. However, they aren't always able to be assessed during the evaluation due to time or resource constraints. Outputs, short-term outcomes, and intermediate outcomes all contribute to the achievement of the long-term outcomes. Although the long-term outcomes represent fulfillment of the purpose of the project, they may or may not

---

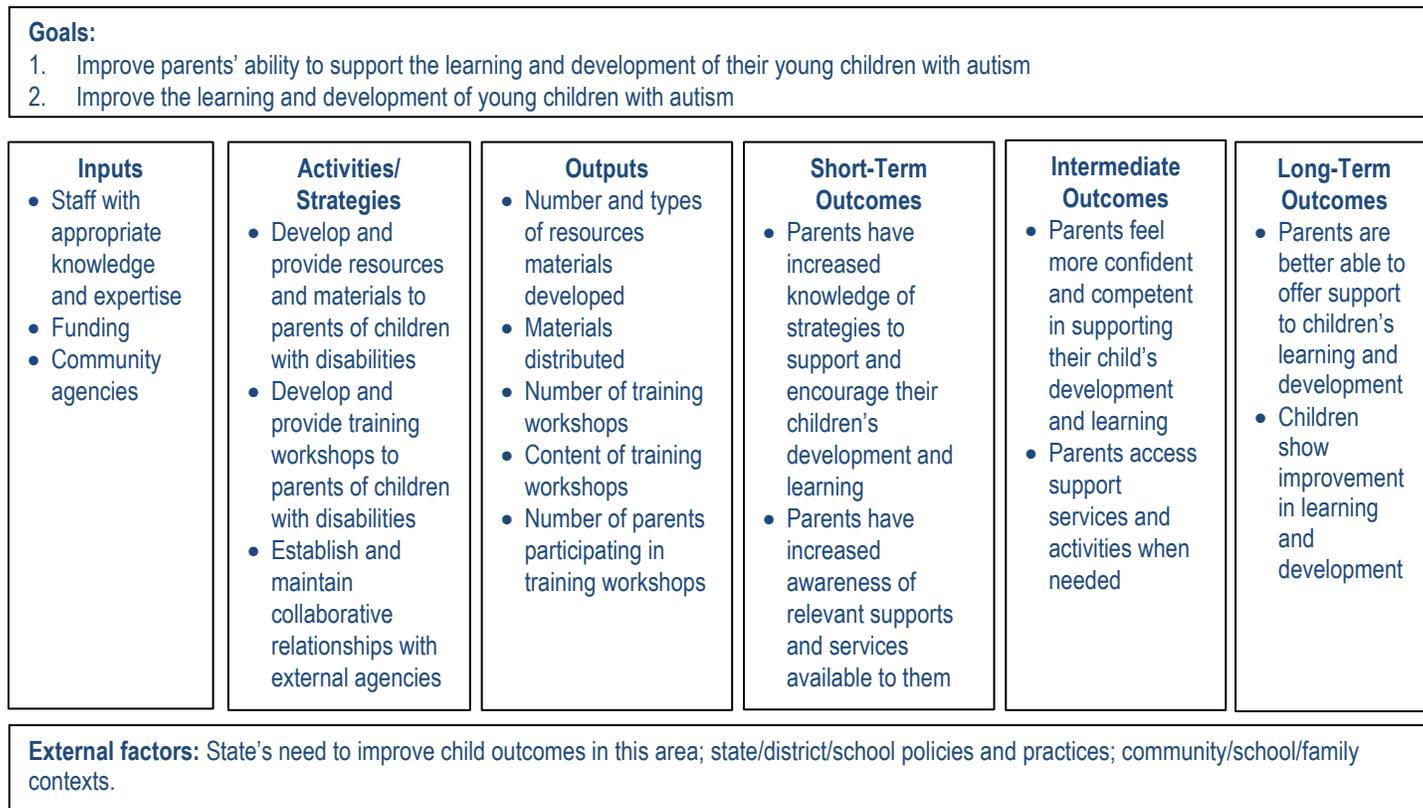
<sup>8</sup> Frechtling, 2007.

represent the achievement of a desired larger project impact. That is, the project may have an anticipated impact that is beyond the immediate scope of the project, either temporally or conceptually, and thus beyond the scope of the logic model.

Since a completed logic model depicts how the project is expected to work, it can be thought of as defining a hypothesis or a series of hypotheses along these lines: If we use these resources to do these activities we will get these outcomes. Evaluation is then the test of those hypotheses, and a logic model can be used by evaluators and project managers to refine and guide data collection and analyses for assessing both process and performance.

In developing a logic model, it can be helpful to begin with a summary chart that contains the information that will populate the logic model. The summary chart in [Appendix A.3](#) is a template that outlines the OSEP priority, assumptions, external factors/context, and inputs. It then displays, in table format from left to right, the project’s goals and objectives, strategies and activities, outputs, and outcomes. The logic model then can be created based on the content in the summary chart. The content of the logic model is less comprehensive than that in the summary chart, often using lines and arrows to depict the temporal and causal connections among the various project elements. Not surprisingly, multiple lines or arrows may come to or from multiple boxes, indicating the complexity of the expected relationships. Both the summary chart and the logic model can be updated continuously as the evaluation progresses, such as when planned activities are revised or when unintended outcomes occur. The logic model also may change as the relationships among the components develop over time, most likely by becoming more complex and interactive. Figure 2 presents a simplified example of a project logic model for a hypothetical Parent Resource and TA Center that wants to improve parents’ ability to support the learning and development of their young children with autism.

**Figure 2. Sample Logic Model for a Hypothetical Parent Resource and TA Center**



## 2.4 Identifying Evaluation Questions

Evaluation questions reflect the goals of the evaluation; some (if not most) of the evaluation’s goals may be established in the request for applications or in the application for funding. If the questions were not defined in advance, they should be developed through discussions with the project staff and key stakeholders, including the funder, as appropriate. Further, the questions should be based on a thorough understanding of the project’s objectives and program theory; consequently, it’s not uncommon for evaluation questions developed prior to the start of an evaluation to change once the team becomes more aware of the specific context and focus of the project. Finally, as when [identifying outcomes](#), it is important to be as specific as possible when developing evaluation questions.

As we discussed in [Section 1.2](#), there generally are two types of evaluation questions: formative and summative. Formative evaluation questions focus on the project’s processes and address the extent to which (and how well) the project is being implemented, while summative evaluation questions target the extent to which a project achieves its expected outcomes. It’s common for formative evaluations not to include a summative component. It’s difficult, however, to complete a comprehensive summative evaluation without paying attention to formative details, as formative evaluation provides data to test fidelity to the project model and explain why the desired changes may or may not be occurring. Table 7 presents examples of formative and summative evaluation questions that might be used by projects funded by four OSEP program areas: Parent Information Centers, Personnel Preparation Programs, Technical Assistance and Dissemination, and Technology and Media. Keep in mind that these represent only a few of the possible questions that might be used for formative and summative purposes.

**Table 7. Sample Formative and Summative Evaluation Questions for Four OSEP Program Areas**

OSEP Program Area	Sample Formative Question	Sample Summative Question
<b>Parent Information Centers</b>	<ul style="list-style-type: none"> <li>What percent of services provided by the Center were deemed to be useful by families of individuals with disabilities?</li> </ul>	<ul style="list-style-type: none"> <li>Are parents better able to take advantage of the various resources available to support them and their children with disabilities?</li> <li>Are parents able to use PIC resources to improve interactions with children or their ability to access services specific to their child’s needs?</li> </ul>
<b>Personnel Development Programs</b>	<ul style="list-style-type: none"> <li>In what ways can the in-service practical experiences offered by the project to program scholars be improved?</li> </ul>	<ul style="list-style-type: none"> <li>Do higher proportions of students receiving special education services from the PDP program graduates, and who also are eligible to participate in standardizing testing, achieve proficiency standards on standardized state assessments of reading and mathematics than students of graduates from other institutions?</li> </ul>
<b>Technical Assistance &amp; Dissemination</b>	<ul style="list-style-type: none"> <li>Has the TA provided by the Center been implemented with fidelity to the project model?</li> </ul>	<ul style="list-style-type: none"> <li>To what degree has the Center contributed to the provision of sustained intensive intervention supports to students who require such intervention?</li> </ul>
<b>Technology &amp; Media</b>	<ul style="list-style-type: none"> <li>What factors might be influencing stakeholders’ decisions to access and download the resources available on the Center website?</li> </ul>	<ul style="list-style-type: none"> <li>Do higher numbers of students with disabilities effectively use technology tools to support their learning after participating in a project-sponsored training??</li> </ul>

## 2.5 Developing an Evaluation Plan

Once the logic model is complete, it's time to develop the more comprehensive evaluation plan. The evaluation plan ties together the project's approach to collecting, managing, aggregating, analyzing, and reporting data on outcomes. The evaluation plan is a step-by-step guide to the “*who, what, where, when, and how*” of an evaluation. Evaluation plans typically have the following sections:

- **Introduction.** The introduction generally incorporates an overview of the need that the project is responding to and a short discussion of [why the need exists \(Section 2.1\)](#). The introduction also often contains background information on the project and the project's [theory of change \(Section 2.2\)](#). The project [logic model \(Section 2.3\)](#) often is included here. It's common for evaluators to document the research or empirical evidence that supports the project model, thereby establishing that the model has a good likelihood of success in achieving the desired outcomes (for more information about how to find evidence in support of a project model, see CIPP's *Demonstrating Evidence across the Project Cycle*, available on the [OSEP IDEAs That Work website](#)). The introduction also can present any factors that should be considered when evaluating the project.
- **Evaluation Questions.** This section of the evaluation plan presents the [questions that guide the evaluation](#).
- **Methodology.** The methodology section of the evaluation plan is where the evaluator details the specific approach to data collection, data management, data aggregation and analysis, and reporting. This can be done through the development of a data collection plan (see [Section 2.5.4](#)), a data analysis plan (see [Section 2.5.2](#)) and a data aggregation and analysis plan (see [Section 3.4.2](#)).
- **Timeframe and Responsibilities.** Finally, the evaluation plan should contain a timeframe for all of the evaluation activities and a clear identification of the party or parties responsible for each step (see [Section 2.5.4.6](#) for information on creating a timeline of evaluation activities and [Appendix A.4](#) for the CIPP Evaluation Plan Template and [Appendix A.5](#) for a sample completed Template).

We recommend that evaluators prepare a comprehensive evaluation plan early in the project period (ideally within 3 months of startup) and then conduct their evaluations in a manner consistent with their evaluation plans. Completing an evaluation plan early on will give the evaluator a roadmap for conducting the evaluation and will enable the evaluator to incorporate the plan into subsequent reports, thereby reducing the amount of time needed to document how the evaluation was implemented at a time when resources typically are limited and time is of the essence. Another benefit of following a pre-established plan is that it will help evaluators to avoid the appearance of “fishing” for positive results during data analysis while obscuring results that may not show the project in a positive light. Of course, there are myriad reasons why an evaluator may need to make changes to an evaluation plan during the course of an evaluation, but whenever possible evaluators are urged to follow their plans for data collection, analysis, and reporting. When this isn't possible, evaluators should document any changes and the reasons for them in their reports.

Since data collection sources, methods, and analysis are often similar for formative and summative evaluations, we do not discuss them separately here. Instead, we highlight important considerations that can be applied to both formative and summative parts of the evaluation.

### 2.5.1 Selecting an Evaluation Design

The design of the evaluation has bearing on the ability to link the project activities to outcomes for families, children or students with disabilities, teachers, or other individuals the project serves. Evaluation designs generally fall into one of four categories: experimental, quasi-experimental, single-case, and non-experimental. Mixed-method designs—those that combine either different designs (e.g., a quasi-experimental study of treatment and comparison groups and a non-experimental correlational analysis) or different data collection methods (e.g., quantitative surveys and qualitative

interviews) within one study—are also commonly used. The selection of the evaluation design depends on (a) the questions that the evaluation is trying to answer; (b) the resources available for data collection, management, and analysis; (c) the availability and feasibility of control groups; and (d) the availability of data to measure outcomes.<sup>9</sup>

When considering which evaluation design to choose, it's important to think about issues related to the validity of the study, since design choices have multiple consequences for validity. It's beyond the scope of this Toolkit to go into detail about the concept of validity, so we briefly mention two basic types: internal and external.<sup>10</sup>

**Internal validity** refers to whether the study results are due only to the manipulation of the independent variables in the study, or whether there are other confounding variables<sup>11</sup> that might influence the study outcomes and which were not taken into consideration in the study. **External validity** refers to the extent to which a study's results can be generalized to persons, settings, treatments, and outcomes not directly included in the study.<sup>12</sup> Each type of validity is subject to different threats (see [Appendix B](#)) which may call into question the study results. Consequently, when designing and conducting the evaluation, to the extent possible, evaluators should do the following:

- Identify and study *plausible* threats to validity;
- Use design elements (e.g., a rigorous evaluation design, additional pretest observations, additional control/comparison groups) to control for possible validity threats; and
- Rule out plausible alternative causal explanations for an effect.<sup>13</sup>

When identifying threats to validity in a study, evaluators should ask three critical questions:

- How would the threat apply in this case?
- Is there evidence that the threat is *plausible* rather than just possible?
- Does the threat operate in the same direction as the observed effect, so that it could partially or totally explain the observed findings?<sup>14</sup>

It's better if evaluators can anticipate validity threats before beginning the study. If the study team knows what the threats may be but can't use design controls in the study, an option would be to try measuring the threat directly in the study to see if it's actually operating. If this isn't possible, at a minimum, evaluators should outline the possible validity threats in the study limitations section of the evaluation report (see [Section 3.5.3](#) for a discussion of study limitations).

In the next sections we briefly discuss the four categories of evaluation designs. Some methodological considerations associated with the different designs are presented in [Section 4.1](#). For more information about research design, Shadish, Cook and Campbell (2002) is an excellent resource.

### 2.5.1.1 Randomized Experimental Designs

The distinguishing feature of **randomized experimental designs, or the randomized controlled trial (RCT)**, is that the researcher has control of the treatment—control that can take many forms.<sup>15</sup> The main characteristics of RCT studies are:

---

<sup>9</sup> By *control* we mean to include both experimental control groups and quasi-experimental comparison groups.

<sup>10</sup> For more information about validity see Kane (2001), Messick (1989) and Shadish, Cook & Campbell (2002).

<sup>11</sup> **Confounding variables**, or confounds, are those that are correlated (either positively or negatively) with both the dependent and independent variable, thereby affecting the study's ability to clearly associate an intervention or project with an observed outcome.

<sup>12</sup> Dimitrov, 2010; Shadish, Cook & Campbell, 2002.

<sup>13</sup> Shadish et al., 2002.

<sup>14</sup> Shadish et al., 2002, p. 40.

- Manipulation of the independent variable;
- Randomized assignment of participants to treatment groups; and
- Controlling for possible confounding variables.<sup>16</sup>

Randomization can take two forms: **random selection** of individuals to participate in the study and **random assignment** of participants to treatment and control groups. A fully randomized study includes both random selection *and* random assignment.

**Confounding variables**, or confounds, are those that are correlated (either positively or negatively) with both the dependent and independent variable, thereby affecting the study’s ability to clearly associate an intervention or project with an observed outcome. Two common types of confounds include:

- **When there is only one unit in one or both conditions (also called the n=1 problem)**—In this situation, for example, the treatment group may include children with disabilities who participated in an intervention in one classroom, while the control group includes children with disabilities who did not participate in the intervention in a different classroom. The n=1 confound makes it impossible for the study team to know whether the outcomes associated with each group are related to intervention or to unobserved characteristics of the children or teachers in each classroom.
- **When the characteristics of the units in each group differ in systematic ways that are associated with the outcomes**—In this situation, for example, the academic performance of students with low incidence disabilities might be compared with the performance of students with high-incidence disabilities, making it impossible for the study team to know whether some characteristic of the students themselves is associated with better (or worse) outcomes, rather than the instructional strategies being used in the classroom.

In general, the RCT is considered the most rigorous and best suited design for making causal claims about the effects of a project, since it provides information related to the counterfactual (i.e., what would have happened if the project hadn’t been implemented) and includes a variety of controls for threats to validity.<sup>17</sup> RCTs also generate data that can inform the unique contribution of project activities or services to desired outcomes. However, the very controlled nature of RCTs raises doubts about the ability of researchers to generalize the results.<sup>18</sup> That is, RCTs have high levels of internal validity, but low levels of external validity.

Additionally, RCTs are quite difficult to implement in the social sciences, and particularly in education. First, it’s often difficult for researchers to randomly assign individuals (or classrooms, schools, etc.) to treatment conditions, let alone to be able to randomly select individuals, classrooms, or schools to participate in a study. Second, even when there is no opposition to the conduct of an RCT, the time and resources required to successfully implement RCTs in a school setting often limit researchers’ ability to carry out such studies.<sup>19</sup> The requirements for conducting an RCT are even more difficult to achieve in the context of special education because of the low number of students and the heterogeneity of this population. For these reasons, we anticipate that few OSEP project evaluations will feature randomized experimental designs.

It’s not necessary for a study to be fully randomized for it to benefit from some of the characteristics of experimental studies. That is, even if selection to participate in the study isn’t random, incorporating random assignment into the

---

<sup>15</sup> Shadish et al., 2002.

<sup>16</sup> Dimitrov, 2010.

<sup>17</sup> There are those who would claim that RCTs represent the “gold standard” of research design in education, including Conrad & Conrad, 1994; Scriven, 2008; and Sullivan, 2011.

<sup>18</sup> Shadish et al., 2002.

<sup>19</sup> Dimitrov, 2010; Shadish et al., 2002.

study design will help to control for some confounding variables. We recommend that evaluators try to incorporate experiments into their evaluation designs whenever possible. Shadish, Cook and Campbell (2002) is an excellent resource for evaluators who are interested in utilizing experiments in their evaluation designs. In [Section 4.1.1](#) we present additional considerations for evaluators conducting RCTs.

### 2.5.1.2 Quasi-Experimental Designs

The main difference between randomized experiments and **quasi-experimental designs (QEDs)** is that quasi-experiments do not feature random assignment of study participants to treatment conditions.<sup>20</sup> Instead, assignment is done by self-selection or by non-random assignment to treatment conditions. Nevertheless, in QEDs researchers still may have control over the following study elements:

- Selecting and scheduling measures;
- Execution of non-random assignment;
- Selection of the comparison group; and
- Treatment schedule.

QEDs generally are easier to conduct than experiments, while still providing a measure of methodological rigor. However, QEDs provide less support for counterfactual inferences—that is, making inferences about what would have happened if the intervention or project hadn’t been implemented—than RCTs since the lack of random assignment to groups means that the treatment groups may differ in systematic ways that may affect the outcome.<sup>21, 22</sup> Consequently, researchers conducting QEDs should outline as many *plausible* alternative explanations for the study results as possible “and then use logic, design, and measurement to assess whether each one is operating in a way that might explain any observed effect.”<sup>23</sup> Of course, this has an impact on the complexity of the study design and, by extension, the difficulty of study implementation, so the study team will need to decide whether ruling out a plausible alternative explanation is worth the time, money, and effort required. Any plausible alternative explanations that were not accounted for in the study design should be discussed in the study limitations section of the evaluation report (see [Section 3.5.3](#) for more information about limitations).

When planning an evaluation with a QED, **evaluators should consider the following important points:**

- If the design involves a **comparison group**,
  - Evaluators should consider whether it would be possible to use stratification or **matching** to improve comparability among groups. In stratification, units are grouped in homogeneous sets (e.g., gender) that contain more units than the study has conditions. When matching, units with similar scores on a matching variable (e.g., school size, ethnicity) are grouped, so that treatment and comparison groups both have units with the same (or very similar) characteristics on the matching variable. See [Section 4.1.2.8](#) for more information about matching.
  - Evaluators should try to determine whether the treatment and comparison groups are similar at baseline, also known as **calculating baseline equivalence** (see [Section 4.1.2.9](#)).

---

<sup>20</sup> Shadish et al., 2002.

<sup>21</sup> Dimitrov, 2010; Shadish et al., 2002.

<sup>22</sup> It is important to point out that even with random assignment it is possible for the treatment groups to differ in systematic ways—a phenomenon referred to as “unhappy randomization”—but it is less likely when group assignment is fully random. Further, it is possible to minimize the risk of unhappy randomization by matching study participants on key characteristics and then conducting randomization (Dimitrov, personal communication, 2009).

<sup>23</sup> Shadish et al., 2002, p. 14.

- If the evaluation involves comparing data from graduates or students across multiple schools or districts—meaning **that data will be collected at multiple levels** (student, school and/or district)—we recommend that evaluators consider accounting for school or districts effects in the data analysis. This is known as **multilevel analysis** and is briefly discussed in [Section 4.4.2.3](#).

In [Section 4.1.2](#) we highlight a few quasi-experimental designs that we believe can be readily applied to project evaluations.<sup>24</sup>

### 2.5.1.3 Single-Case Designs

Single-case designs involve in-depth study of a single “case,” which can be a person, a group, an institution, or even a culture. Also known as single-subject designs, single-case design studies can be either (a) a single-case study in a natural (uncontrolled) environment, or (b) a single-case experiment in a more controlled environment.<sup>25</sup> These designs are frequently used by psychologists, counselors, social workers, and special educators. The two most common single-case designs are the A-B-A-B design and the multiple baseline design. As part of its guidance for the design and conduct of rigorous studies in education, the WWC has created a set of [pilot single-case design standards](#). See [Section 4.1.3](#) for more information on designing single-case designs and [Section 4.4.3](#) for information about analyzing data from single-case studies.

### 2.5.1.4 Non-Experimental Designs

Non-experimental designs include case studies, descriptive studies or surveys, correlational studies, and *ex post facto* studies (i.e., studies that take place *after the fact* using secondary data). It’s possible to investigate a presumed cause and effect (e.g., which specific activities of a Parental Information Center project contributed to improved outcomes for children and families) in a non-experimental study, but the structural features of experiments that help to rule out possible alternative explanations and identify the counterfactual—that is, information about what would have happened if a particular intervention or project hadn’t been implemented—are often missing. Non-experimental designs generally are considered to be more appropriate for formative evaluations or monitoring the progress of (or fidelity to) project implementation than for summative evaluations. However, sometimes non-experimental studies are the only viable option for evaluators—especially if the evaluation was not planned prior to beginning implementation of the project or if the evaluator has little control over events during a study. [Section 4.1.4](#) includes information on the different non-experimental designs that might be used in a project evaluation. Nevertheless, whenever possible, we recommend that evaluators of OSEP projects use either experimental or quasi-experimental designs to answer their summative evaluation questions.

---

<sup>24</sup> See Shadish et al., 2002, for additional discussion of these designs.

<sup>25</sup> Dimitrov, 2010; Kennedy, 2005.

### 2.5.1.5 Mixed-Method Designs

Although mixed-method design studies can refer to studies that combine different designs (e.g., quasi-experimental and non-experimental) and those that combine different data collection methods (e.g., surveys and focus groups), in our experience the latter are more common.<sup>26</sup> An entire field of literature has developed related to the nature of mixed methods versus mixed methodology<sup>27</sup> (that is, in general, mixed methods studies combine different data collection and analysis methods while mixed methodology studies combine different theoretical approaches as well as data collection and analysis methods), but it's beyond the scope of this Toolkit to discuss it. For our purposes here, we focus on mixed method studies, in which "the investigator collects and analyzes data, integrates the findings, and draws inferences using both qualitative and quantitative approaches or methods in a single study or a program of inquiry."<sup>28</sup> Specifically, we list four decisions that relate to the selection of a mixed methods approach to a study:

1. What is the implementation sequence of the quantitative and qualitative data collection in the proposed study?
2. What priority will be given to the quantitative and qualitative data collection and analysis?
3. At what stage in the [study] will the quantitative and qualitative data and findings be integrated?
4. Will an overall theoretical perspective (e.g., gender, race/ethnicity, lifestyle, class) be used in the study?<sup>29</sup>

For more information on the design and conduct of a mixed method study, see Creswell, 2003; Creswell, Plano Clark, Gutmann & Hanson, 2003; Brewer & Hunter, 2006; and Tashakkori, & Teddlie, 1998.

In the next sections we provide a brief overview of the steps involved in further developing the evaluation plan. First, we present considerations for creating a data analysis plan, followed by issues and considerations for using and selecting a sample and planning for data collection.

---

<sup>26</sup> See, for example, Brewer & Hunter 2006; Greene, 2006; Morgan, 2007; Tashakkori & Creswell, 2007; and Tashakkori & Teddlie, 1998, 2003.

<sup>27</sup> Tashakkori & Creswell, 2007, p. 4.

<sup>28</sup> Creswell, 2003, p. 211.

<sup>29</sup> Creswell, Plano Clark, Gutmann & Hanson, 2003, cited in Creswell, 2003, p. 211.

## 2.5.2 Creating a Data Analysis Plan

An essential part of developing the evaluation plan is creating a data analysis plan. Creating a data analysis plan prior to the data collection period will help ensure that the evaluator is able to collect and analyze data that will respond to the evaluation questions in the most rigorous way possible. Moreover, it ensures that (a) the instrumentation chosen or developed for the evaluation will gather the needed data in the correct format or scale, and (b) sufficient numbers and types of respondents or data sources will be included in data collection. For example, if one of the evaluation's analytic goals is to describe how much time the children with disabilities spend on a particular classroom activity, the evaluator will need to consider whether to report results using a checklist (e.g., Daily, Weekly, Monthly) or in units of time (e.g., minutes/day, hours/week, days/month). If, however, an analytic goal is to meaningfully compare the amount of time children with different types of disabilities spend on classroom activities, the evaluator not only needs to consider the type of scale to use and the format of the data but also the types of disabilities that will be compared and, potentially, the minimum number of respondents for each type (i.e., total sample size and sampling strategy). This will help the evaluator to capture high-quality data, a sufficient quantity of data, and data that have sufficient generalizability.

Answering these questions requires the evaluator to identify the *unit* that will be subject to data analysis. Identification of the *unit of analysis*<sup>30</sup> is informed by the language of the outcome measure and evaluation question as well as the availability of data for analysis. Examples of different units are provided by Patton, who states “*each unit of analysis* [e.g., preparation program, graduation cohort, special education classroom, individual child] *implies a different kind of data collection, a different focus for the analysis of data, and a different level at which statements about findings and conclusions would be made.*”<sup>31</sup> Decisions about the unit of the analysis and the corresponding data collection and analysis strategy must be made during the development of the evaluation plan.

[Appendix A.6](#) presents a data analysis plan template created for this Toolkit. As can be seen in the Appendix, the data analysis plan includes information related to the:

- Study design (see [Section 2.5.1](#) and [Section 4.1](#))
- Treatment and control (or comparison) groups
- Type of data analysis (see [Section 4.4](#) for a general discussion of data analysis)
- Variables to be used for data analyses (see [Section 4.4.2](#) for a discussion of quantitative analysis, [Section 4.4.3](#) for analysis of data in single-case designs, and [Section 4.4.4](#) for qualitative analysis)
- Instruments and data collection techniques (see [Section 4.3](#) for a discussion of data collection methods)
- Sample (see [Section 4.2](#) for a discussion of sampling)
- Minimum number of responses and/or response rate

Since it allows the evaluator to specify a different analysis approach for each evaluation question, the data analysis plan template can be used either for a single method study or a mixed-method study.

---

<sup>30</sup> The unit of analysis (i.e., the level at which outcomes will be analyzed) is different from the unit of selection (i.e., the level at which units were put into treatment or comparison groups).

<sup>31</sup> 2002, p. 228.

When developing a data analysis plan it also may be helpful to create table shells (i.e., empty tables that illustrate how results will be displayed) to help identify the specific variables that will be needed and visualize how to format the data for reporting. Examples of table shells follow.

**Table 8. Table Shell Example 1. Amount of time spent on the targeted classroom activity among PDP program graduates**

The purpose of this table is to describe how much time graduates of a Personnel Development Program (PDP) spend on a particular classroom activity, reported in units of time (e.g., minutes/day; hours/week; days/month)

Amount of time spent on the targeted classroom activity among program graduates (n= )	
<b>Average amount of time in minutes per day</b>	
Standard Deviation	
Range	
<b>Average amount of time in hours per week</b>	
Standard Deviation	
Range	
<b>Average amount of time in days per month</b>	
Standard Deviation	
Range	

**Table 9. Table Shell Example 2. Amount of time spent on the targeted classroom activity, by type of PDP graduate**

The purpose of this table is to descriptively compare the amount of time, in different *units of time*, that different types of PDP graduates spend on classroom activities

Amount of time spent on the targeted classroom activity, by type of program graduate		
	Treatment Group (n= )	Control Group (n= )
<b>Average amount of time in minutes per day</b>		
Standard Deviation		
Range		
<b>Average amount of time in hours per week</b>		
Standard Deviation		
Range		
<b>Average amount of time in days per month</b>		
Standard Deviation		
Range		

**Table 10. Table Shell Example 3. Testing the Statistical Significant of Differences between Treatment and Control Groups.**

The purpose of this table is to display the results of *t*-tests of differences between means of the treatment and control groups.

	Treatment		Control		<i>t</i> ( <i>df</i> )	<i>p</i>	95% CI	
	n =		n =				LL	UL
	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>				
Minutes per day								
Hours per week								
Days per month								

Note. CI = confidence Interval; LL = lower limit; UL = upper limit

As can be seen, in the first example (Table 8) the table shell shows that the study will need to collect or derive data from evaluation participants on three types of variables: minutes per day, hours per week, and days per month that participants spent on the target activity over the period of time that was specified. Thus, in example one, if the study wants to report on all three types of variables, evaluators will need to collect data for at least one month. Additional months may be necessary to ensure the data aren't sensitive to seasonal differences in classroom schedules (e.g., holiday breaks) and may help to create a more stable estimate of the average time spent on the activities.

The second example, presented in Table 9, posits that PDP graduates can be grouped according to whether or not they received a "treatment," perhaps specialized training in the targeted classroom activity. Table 9 illustrates that, in order to fulfill the data needs for this example, it will be necessary to collect the same data that were presented in Table 8 for both the treatment and control groups. Further, if the study uses sampling to compose the treatment and control groups, it will be necessary to select sufficient numbers of both treatment and control group members if a goal is to generalize the study findings to the larger group of all program graduates. Finally, Table 10 illustrates the results of independent samples *t*-tests comparing the PDP graduates in the treatment and control groups.

### 2.5.3 Identifying a Study Sample

Although conducting a census of an entire population (such as all teachers who participated in a particular training) may be technically possible because the entire population is potentially known to the evaluation team, it sometimes isn't feasible or even particularly desirable. Imagine, for example, an evaluation of a training program with 200 or more participants—conducting an in-depth survey of all 200+ trainees would be time-consuming, costly, and likely to result in a low response rate. Selecting a sample enables the study team to save time and money, and, if done correctly, can help to improve the quality and accuracy of the data collected.<sup>32</sup> This becomes even more relevant if the evaluation plan includes use of a comparison group, or when the evaluation seeks to collect data on the performance of large projects where the numbers of individuals served may reach 1,000 or more.

Some considerations when making decisions about how to select a sample include:

- **Available resources**—How much time and money can be spent? How many people are available to work on the study (e.g., to recruit, collect data, analyze data)? Is a census cost-prohibitive?

<sup>32</sup> Fiore et al., 2012.

- **Desired precision of estimates** –What is the minimum sample size needed in order to reach conclusions with a pre-specified level of confidence (see [Section 4.2.1](#) for a brief discussion of power analysis)? Will results be used to make comparisons between groups?

Sampling is one aspect of an evaluation where a team member with specific training, expertise, or experience is needed, as there are multiple technical details that must be accommodated when the evaluation uses a sample, including how to analyze the data that is collected using sample weights. There also are many strategies for designing a sampling framework and selecting a sample. [Section 4.2](#) presents additional information about two basic types of sampling—random sampling and purposeful sampling—as well as a brief discussion on power analysis. For additional resources on sampling methodology, see the recommended readings in [Appendix D](#).

## 2.5.4 Preparing a Data Collection Plan

The final step in preparing the evaluation plan is developing the overall data collection plan and timeline. This entails identifying the specific tasks that will need to be completed to ensure success in collecting the data needed for the evaluation. The CIPP Evaluation Plan Template in [Appendix A.4](#) presents a number of tables that can be used as part of a data collection plan. Rather than limiting our discussion to those specific tables, however, here we have outlined a series of questions to guide evaluators’ thinking while developing a data collection plan. Then evaluators can choose the format for the plan that best suits their needs.

### 2.5.4.1 *What instruments or data collection techniques will supply the variables that are needed?*

It may be possible to collect some of the needed data using existing data sources, or “secondary source” data. Some publicly available databases are readily accessible online. These databases may allow users to download data files simply with the click of a mouse, or they may require users to sign a Data Usage Agreement. For example, researchers can easily download Excel files of student test results, such as [data on the participation of students with IEPs in 2012 MCAS ELA and Mathematics Tests in MA](#), on the Massachusetts Department of Elementary and Secondary Education website. The U.S. Department of Education [Common Core of Data](#) is an example of a publicly available database that requires users to sign a Data Usage Agreement before downloading the data files.

Secondary data also may be available through state or local education agencies. In some cases it will be necessary for the evaluators to complete a formal data request; in others they will need to work with the local research department or Institutional Review Board (IRB) to obtain access to the data (see [Section 3.1.1](#) for more information).

If the data that are needed aren’t available as secondary data, and the study team decides that these data are essential to the evaluation, it will be necessary to conduct a unique data collection. In this case, it may be possible to use an existing data collection instrument that another individual, group, or publisher has developed and is making available. However, it’s possible that existing instruments won’t supply the specific variables that are needed to respond to the specific evaluation questions. In these cases, the evaluator will need to develop a data collection instrument or instruments to capture the specific data variables that are identified in the data analysis plan. We address this briefly in the next section and in more detail in [Section 3.3](#) and [Section 4.3](#).

#### 2.5.4.2 What types of instrumentation or forms need to be identified or developed?

Data collection instruments come in many formats or types. Some of the most popular are surveys, observation checklists and rubrics, multiple choice tests, and standardized assessments of knowledge and skills. Other types include individual interview and focus group protocols, case notes and case management logs, and tracking logs. Data also may be collected through audio or video data recorders. For some of these, the evaluator will have a choice between instruments that already have been developed and, ideally, tested, so as to establish the validity and reliability of the instrument. In other cases, it will be necessary to develop the instrument, pilot test it, and finalize it for the evaluation's unique data collection activities. Instrument choice or development is guided by

- the nature of data and individual variables that are needed to respond to evaluation questions;
- the opportunities the study team will have to collect data; and
- the cost of data collection, including costs related to either buying or developing the instrument.

Some online resources on selection and development of data collection instruments include:

- Web Center for Social Research Methods *Research Methods Knowledge Base* (<http://www.socialresearchmethods.net/kb/measure.php>)
- The National Center for Education Statistics *Standards for Education Data Collection and Reporting* (<http://nces.ed.gov/pubs92/92022.pdf>)
- United States Census Bureau guidelines on developing data collection instruments and supporting materials (<http://www.census.gov/about/policies/quality/standards/standarda2.html>)
- Compendium of Student, Teacher, and Classroom Measures Used in NCEE Evaluations of Educational Interventions (<http://ies.ed.gov/ncee/pubs/20104012/>)

The CIPP Evaluation Plan Template in [Appendix A.4](#) includes an Instrument Information table that can be used to monitor the status of the different instruments that will be used in the evaluation. Later in this Toolkit we provide information about methodological considerations for the development of different data collection instruments, such as surveys ([Section 4.3.1](#)), observation protocols ([Section 4.3.2](#)), and interview protocols ([Section 4.3.3](#)). Additionally, [Appendix D](#) includes good resources for learning more about instrument development.

#### 2.5.4.3 How will data collectors and coders be trained? What materials need to be developed to document and support that training?

All data collectors and coders that will be supporting the evaluation must be trained in the data collection and coding protocols. The number of individuals needed for data collection and coding and the nature of the training will depend upon the evaluation design and on the amount, timing, type, and frequency of data collections. For example, for an evaluation supported entirely with quantitative secondary data, it may be necessary to train only one individual in the retrieval of all data. For a study that includes classroom observations or open-ended surveys of program graduates, on the other hand, multiple observers and multiple individuals may need to be trained to code the qualitative data that emerge from the surveys. A project in which the data need to be collected from multiple individuals over a very short time frame may necessitate multiple data collectors, to ensure data are collected within the prescribed time limit. In contrast, a project in which data can be collected over a long period of time may require that relatively few data collectors or coders operate concurrently. Training can take many forms, including in-person training, online training, or

development of training modules or manuals. The time associated with training data collectors and coders should be factored into the data collection plan. We discuss training further in [Section 4.3.2.2](#).

In those cases in which data collection is limited to gathering secondary data, training may entail an introduction to the source database and the instruments, techniques used to populate the database, and the methods for retrieving the data. In those cases where unique data are collected, training should entail an introduction to the study and evaluation questions; a thorough review of the instruments; a review of data collection rules, procedures, and timelines; and, if possible, practice runs to establish familiarity with data collection. Further, when there are multiple data collectors, it's important to establish the inter-rater reliability of data collections (see [Section 4.3.2.6](#) for a discussion of this), so as to ensure that the complete dataset represents standardized and uniform data collection procedures.

When data are collected through interviews, focus groups, or audio or video recordings, the study team also may work with data coders who review the data and categorize or code the text, audio, or video files, but are separate from the actual data collection. In these cases, background materials and training should be developed for coders, and the evaluation team should establish reliability and inter-rater reliability across coders.

Data collectors also should receive training on human subject protections, such as those codified in the “[Common Rule](#)”, or the “[Federal Policy for the Protection of Human Research Subjects](#)” (1991). Other policies that commonly are referenced include the Health Insurance Portability and Accountability Act of 1996<sup>33</sup> (HIPAA) and the Family Educational Rights and Privacy Act (FERPA).<sup>34</sup> Most universities and colleges maintain resources related to human subject protections; training may be required for all university faculty and staff that conduct data collections. Some private organizations will provide training and support for human subject protections (usually for a fee) if it's not available to the study team otherwise.<sup>35</sup>

Some instrument developers require that all data collectors receive formal training in the use of their measure; the [School Observation Measure \(SOM\)](#) is an example of such an instrument. To become certified in the use of the SOM, all data collectors must participate in the developer's one-day training session. This can take place in-person or over a video conference. Other instrument developers provide information or training manuals about how to use an instrument, but do not require users to obtain formal training (e.g., the [Reformed Teaching Observation Protocol](#)).

#### ***2.5.4.4 When will data be collected? How frequently will data be collected?***

Large evaluations that consist of multiple evaluation questions likely will be active throughout the year conducting data collection, data entry, analysis, or reporting activities. Small evaluations may have shorter, targeted data collection periods. Regardless, the key is to ensure that all necessary activities are conducted on schedule and in the right order. Being well organized is of particular importance for data collection activities, as they often are time-specific; once a window for data collection closes, the study team may not be able to capture the necessary data, depending on the specific constraints or parameters of the study. Failure to meet schedules can, at a minimum, affect data analysis and the resulting interpretations and lead to limitations in the study findings. It also is possible that failure to collect data on schedule may result in an inability to draw a conclusion or respond to one or more evaluation questions ([Section 3.3](#) discusses managing data collection in more detail).

---

<sup>33</sup> <http://www.hhs.gov/ocr/privacy/>

<sup>34</sup> <http://www.ed.gov/policy/gen/guid/fpco/ferpa/index.html>

<sup>35</sup> One such company is Ethical and Independent Review Services, <http://www.eandireview.com/> (NOTE: By providing a link to this website, Westat is not advocating the use of this company.)

#### ***2.5.4.5 How will data be entered and verified for accuracy? Where will data be stored?***

Once data have been collected, it will need to be reviewed and made available for analysis. Generally speaking, the raw data that are collected aren't "analysis ready." Rather, data need to be checked for quality and often need to be entered into an electronic database, spreadsheet, or statistical package for data analysis. Data entry serves the secondary purpose of maintaining an electronic file of all data collected—as a supplement to the original raw data files.

Data entry requires data technicians or support staff, all of whom should receive training and periodic "spot checks" to ensure accuracy. A senior-level staff person often will develop the data entry files and coding procedures and then work with data technicians to train them and review their data entry.

A persistent topic in data analysis is what to do about missing data and about responses that aren't what the study team were expecting (e.g., "outliers"; see [Section 4.4.1](#) for a discussion of missing data and [Section 3.3.2](#) for a definition of outliers). A senior-level staff person often will work with data technicians on how to code or flag these fields for later treatment.

#### ***2.5.4.6 Creating a timeline of evaluation activities***

After answering the questions outlined above, the next step is to prepare a timeline for the evaluation that will enable the study team to monitor evaluation activities and ensure that the evaluation is on schedule. As mentioned above, the CIPP Evaluation Plan Template in [Appendix A.4](#) includes tables that can be used as part of a data collection plan. The Gantt chart is another commonly-used organizational tool for planning evaluation projects. Gantt charts provide a general timeline for each of the evaluation's major activities. Table 11, on the next page, presents an example of a Gantt chart of a hypothetical data collection schedule.

**Table 11. Gantt Chart of the Project’s Data Collection Schedule**

Task	Time Units (Days, Weeks, Months, Years, etc.)											
	1	2	3	4	5	6	7	8	9	10	11	12
Develop logic model	█											
Develop evaluation plan	█											
Develop analysis plan	█											
Prepare data collection instruments	█	█										
Complete IRB process	█	█										
Secure district participation	█	█	█									
Prepare training materials		█	█									
Conduct data collector training			█									
Conduct data collection				█	█	█	█					
Enter and clean data				█	█	█	█	█				
Conduct coder training								█				
Code data									█			
Analysis and reporting										█	█	█

Now that we have discussed the steps involved in planning the evaluation, it’s time to turn to the different aspects of conducting the evaluation.

## 3 Conducting the Evaluation

In this section we discuss various steps involved in conducting the evaluation. We would like to point out that even though these steps are presented in a specific order in this section, many of the steps actually will take place concurrently and others may occur in a different order than presented here, depending on the needs of the study. We have simply ordered them in a way that seems to make sense within the structure of this Toolkit.

### 3.1 Obtaining Permission to Carry Out Evaluation Activities

Prior to beginning the evaluation, the study team must obtain permission to carry out the evaluation activities in the participating districts and schools. Evaluators affiliated with a university or college likely must obtain approval from the institution's Institutional Review Board (IRB)—the group responsible for reviewing research to assure the protection of the rights and welfare of the human subjects. Large research organizations have their own IRB responsible for approving evaluation activities. Smaller evaluation companies often do not have their own IRB, but if the evaluation design involves collecting data from students, their parents/guardians, or school personnel, at a minimum the study team will need to obtain permission from the local school district. Furthermore, if secondary data will be used in the evaluation it may be necessary to obtain permission to access the data. These various types of permissions are briefly discussed below.

#### 3.1.1 Getting IRB Approval

Most evaluations require collection of data on individuals involved in or affected by the project (i.e., the “treatment” population), both to provide formative feedback on project implementation and to gather summative data on outcome achievement. Further, experimental or quasi-experimental evaluations also call for data collection from a control or comparison group (see [Sections 2.5.1](#), [4.1.1](#) and [4.1.2](#) for more details on those types of evaluation designs). Therefore, it's likely that evaluators will have to seek Institutional Review Board (IRB) approval for data collection, since human subjects are granted protections from data collections that may be harmful either at the time of data collection or, if foreseeable, at a future point in time (see the Department's resources on [Protection of Human Subjects Research](#) for more information).

Obtaining IRB approval involves describing the evaluation's approach to data collection and detailing any and all circumstances in which there will be contact with a human subject (e.g., the student, parent, or teacher that is the subject of the data collection). Requesting and receiving IRB approval guarantees the study team has taken all of the necessary steps to ensure that human subjects are protected and that the research protocol discloses any risks associated with participating in the study. It's common for state and local education agencies to require IRB approval before allowing any unique data collections. Evaluators that are affiliated with an institution of higher education or a large research organization may have access to an IRB panel through their institution. Private, for-profit IRB panels may also be contracted to perform the review.

Typical fields in IRB applications include:

- General information
- Objectives of proposed project
- Description of human participants
- Summary of research and data gathering procedures

- Location of project
- Confidentially safeguards
- Assessment of risk
- Consent procedures
- Potential benefits

In some cases the evaluation may be exempted from a full IRB review. The IRB may determine a research activity to be exempt from the need for IRB review when the only involvement of [human subjects](#) will be in one or more of the following categories:

- A. Research conducted in established or commonly accepted educational settings, involving normal educational practices, such as:
  1. research on regular and special education instructional strategies, or
  2. research on the effectiveness of or the comparison among instructional techniques, curricula, or classroom management methods.
- B. Research involving the use of educational tests (cognitive, diagnostic, aptitude, achievement), survey procedures, interview procedures or observation of public behavior, unless:
  1. information obtained is recorded in such a manner that human subjects can be identified, directly or through identifiers linked to the subjects; and
  2. any disclosure of the human subjects responses outside the research could reasonably place the subjects at risk of criminal or civil liability or be damaging to the subjects' financial standing, employability, or reputation.
- C. Research involving the use of educational tests (cognitive, diagnostic, aptitude, achievement), survey procedures, interview procedures, or observation of public behavior that isn't exempt under paragraph (B) of this policy, if:
  1. the human subjects are elected or appointed public officials or candidates for public office; or
  2. Federal statutes require without exception that the confidentiality of the personally identifiable information will be maintained throughout the research and thereafter.

Even if it seems the evaluation may be exempt from a full review, most IRBs still require the evaluator to submit the information listed above to receive approval for the exemption. Evaluators should always refer to their institution's (or the district's) specific rules and guidelines for obtaining IRB approval.

### 3.1.2 Securing District and School Approval

Once the study team has identified the individuals who will be participating in the study (see [Section 3.2.1](#) for more information on keeping track of study participants), it will be necessary to obtain approval from the corresponding school districts and schools, informing them of the details of the evaluation and obtaining all necessary permissions. Many districts have their own research departments and established protocols for conducting research in the district, in which case the evaluator should contact the district research office to review its requirements. Often this information is available online (see for example the Fairfax County, VA, [Office of Program Evaluation](#)). Typical procedures require the submission of a research proposal that is reviewed by the district's research board. Additionally, some districts have requirements about what times of the academic year data can be collected (e.g., data collection may not be permitted during the district's mandatory testing periods). If a formal research approval application isn't required, the study team still should notify the district administration of the evaluation and what it will entail. [Appendix C.1](#) and [Appendix C.2](#) present examples of District Notification Forms that might be used for an evaluation of a Personnel Development

Program—one for a district with a research approval office and another for a district without a research approval office. A sample District Response Form can be found in [Appendix C.3](#).

In general, once a district has given approval for a study to take place, school approval will follow. This isn't always the case, however. School principals always must be informed of any data collection activities that will take place in their schools (e.g., surveys, observations, or interviews). [Appendix C.4](#) presents an example of a School Notification Letter. If the district approval included an assigned study number, this should be included in the School Notification Letter. To avoid potential problems related to questions about whether the study team obtained the appropriate permissions, it's good practice to keep records of all contacts with districts and schools, including the date, time, and person contacted. This information can be maintained in the evaluation's data tracking system (see [Section 3.3.1](#) for more information on data tracking systems).

### 3.1.3 Obtaining Access to Secondary Data

If the evaluation intends to use secondary data it may be necessary for the evaluator to complete a Data Usage Agreement. This agreement simply may require the evaluator to agree to follow the terms of use of the data, or it may provide more specific data usage instructions, such as the specific variables the evaluator will have access to, how the data will be transmitted to the study team (e.g., spreadsheet file, database, secure file sharing protocol), and any limits that are dictated by state or local district policy on confidentiality and protection of human subjects. A data sharing agreement generally is formed after the evaluator has (successfully) applied to a state or local education agency to conduct research (see above for a discussion of obtaining district permission to conduct the study).

## 3.2 Recruiting Study Participants

Once permission to conduct the evaluation has been obtained, the next step is to recruit participants. Since evaluation activities almost always are voluntary—even for the population that is receiving the service or program being evaluated—the study team will need to develop strategies for engaging and recruiting participants for the study. Two essential aspects of recruiting participants include keeping track of the graduates and obtaining consent from participants. We discuss these aspects below, following a brief presentation of some considerations for communicating with participants.

Listed below are some guidelines evaluators should follow when communicating with participants:

- Clearly state what the evaluation will be asking the participant to do, preferable at the start of the communication (e.g., “We are asking you to complete the enclosed survey...”).
- Provide a deadline for a response (e.g., “Please complete the survey no later than **DATE**”).
- Explain the overall purpose of the study. A pamphlet or short document for the participant to review and keep may be helpful.
- Use clear, concise language. Compose written materials that will be sent to parents (e.g., consent forms) at an 8<sup>th</sup> grade reading level (the Microsoft Support page has information on how to [check the reading level of text in Microsoft Word](#)).
- Always provide a telephone number and email address for a contact person on the study team so participants can obtain additional information.
- Include self-addressed return envelopes if participants must return hard copy materials.
- For online surveys, include instructions for obtaining usernames and passwords. If providing a username and password to participants, and the information to be collected is sensitive, do not include the username and password in the same communication. This is necessary for ensuring confidentiality and protection of participant data.
- When parents or children who aren’t native English speakers are involved, whenever possible communication should occur in the participant’s primary language. It may be helpful to ask a project staff person with whom the target participants are familiar and comfortable to help communicate the evaluation’s plans and importance. In some cases, the evaluator may need to enlist the assistance of an interpreter or translator. For instance, if recruitment letters are being sent to a population with a significant percentage of non-native English speakers, it could be helpful to have the materials translated and send each potential participant materials in English and in the individual’s home language.<sup>36</sup>

Another strategy for recruiting participants to the study may be the provision of incentives or awards in return for participating in data collection. While this is a common practice, it’s important to consider whether local regulations would allow the offering of such an incentive (e.g., some school districts prohibit offering of incentives to district staff). If an incentive is acceptable, the evaluator should determine the level of incentive or award that will be necessary to recruit a sufficient sample size. Of course, the number and types of participants (e.g., parents, related-services

---

<sup>36</sup> Interpreters are used when the communication will be in-person, while translators are used when the communication will be in written form.

providers, students) to be recruited will depend on the evaluation design and the type of sample chosen (see [Section 4.1](#) for information on different types of evaluation designs and [Section 4.2](#) for information on sampling). For example, if the evaluation calls for a longitudinal study (i.e., an evaluation in which there will be multiple data collections from the same individuals over a period of time), the evaluator will need to consider the level of incentive or award necessary to maintain participants over the entire period of the study. Also, in the case of a longitudinal study, it's reasonable to expect attrition of participants over time (i.e., participants discontinuing their participation in the study), and researchers often develop larger-than-necessary samples to ensure that a sufficient number of participants will be available at the last planned data collection.

### 3.2.1 Keeping Track of Participants

Maintaining contact with study participants is essential to the success of evaluations that collect data at more than one time point. Even evaluations that collect data at only one time point can be better served by keeping track of potential participants from the time they first encounter the project. For example, at the end of each year, an Educational Technology, Media & Materials (T&M) project may wish to conduct a survey of individuals who download educational materials from the project website. To do this, right from the start, the website must be created so that users are required to enter their names and contact information into a web database prior to downloading the materials. That way the evaluators will have access to the list of potential survey respondents. One effective method to continuously track participants is to maintain on-going contact (e.g., by sending related resource links periodically). Whether or not these strategies are successful, the evaluator will need to work with project staff and administrators to gather as much data as possible about the potential participants while they are in contact with the project and then reach out to request their participation in the evaluation.

Some steps that the study team can take to keep track of potential participants include:

- collecting alternate contact information (e.g., phone numbers, addresses, emails);
- establishing a web-based portal (e.g., Moodle or SharePoint site, or even a social media page) for potential participants to communicate with the project or each other and to update their;
- maintaining a secure database with contact information and ongoing records of contacts (e.g., emails, telephone calls) with potential participants and alternate contacts, including the initial contact and any follow-up contacts.

[Section 3.3.1](#) presents additional information about the types of information that should be collected as part of the evaluation's data tracking system. Since it may be difficult to select a random sample for the evaluation or to keep up with all individuals or families served by a project, in [Section 4.2.3](#) we present a number of purposeful sampling options. However, it's important to consider the effects that selection bias can have on the evaluation, especially if a convenient sample of persons who already are keeping in contact with the project is selected.

### 3.2.2 Obtaining Participants' Consent

Any evaluation that plans to interact with human subjects must obtain consent from the participants prior to beginning the data collection activities (see the Department's resources on [Protection of Human Subjects in Research](#) for more information). Originating in the health services and medical fields, the term "consent" represents the concept of providing information to an individual who is capable of understanding the information and making an informed and judicious decision to either participate or not participate in the study. A similar standard applies in educational research: the consent giver must be provided with sufficient information in a way that allows the consent giver to understand the

study, be aware of any risks participation entails, and make an informed and judicious decision about participating. In practice, this means that consent forms must be written in clear and simple language. The evaluator also bears the burden of using means of communication that are accessible by the participant or consent giver and answering any questions the participant may have regarding the evaluation and his or her participation in data collection. Finally, it's important to stress that participation is voluntary; the participant can withdraw from p out punishment (see [Section 3.2](#) above for a brief discussion of considerations for communicating participation at any time with participants). Studies that include participants who are considered children must obtain consent from the children's parents or legal guardians.<sup>37</sup> In general, there are two forms of consent:

- **Passive consent:** Assumes that the participant has given consent unless action, such as a written statement, is taken to indicate otherwise. Passive consent often involves the distribution of a letter explaining the study to the proposed participant or their legal guarding and informing the participant or guardian that they should return the signed letter to the school or study team if they do *NOT* want to participate in the study (see [Appendix C.5](#) and [Appendix C.7](#) for examples of passive consent forms that could be used in an evaluation of a Personnel Development Program). In many study contexts, passive consent is assumed when an adult agrees to participate in data collection activities such as completing an online survey or participating in a telephone interview.
- **Active consent:** Requires the participant to provide a written signature for consent. This is often required when children are participants in a study. In this case, the child cannot be involved in data collection unless a written signature from the parent is provided to the study team (see [Appendix C.8](#) for an example of an active consent form from an evaluation of a PDP). Active consent may also be required if sensitive information is being collected from or about the proposed participants (see [Appendix C. 6](#) for an example).

In addition to obtaining informed consent from parents, researchers might need to obtain verbal or written assent from child participants. Assent differs from informed consent because it doesn't imply an understanding of the purposes, risks and benefits of research, but merely indicates a willingness to participate in the research activities. An IRB will make a decision about whether to judge participating children capable of providing assent by considering the age, maturity, and psychological state of the children involved. If participating children are judged to be capable of providing assent, the IRB will usually require the research protocol to "make adequate provision to seek assent from child participants".<sup>38</sup> The assent process can be as simple as asking a young child whether they are willing to "play a game" with the researcher. With older children, the assent process might involve an explanation of the research and a signature on a brief assent form.

---

<sup>37</sup> For simplicity, we will only use the term "parents" in this section.

<sup>38</sup> See U.S. Department of Health and Human Services information on the definition of assent at <http://www.hhs.gov/ohrp/policy/faq/children-research/child-assent-age-requirements.html> and requirement for assent at <http://www.hhs.gov/ohrp/policy/faq/informed-consent/requirements-for-assent-in-research-with-children.html>

## 3.3 Managing Data Collection

With the instruments chosen or designed, the data collection staff identified and trained, the participants recruited, and participants' consent obtained, the evaluation team is now ready to collect data. Rather than going into detail about the specific steps involved in data collection, we choose to highlight here two important strategies to manage the data collected during the evaluation: the creation of a data tracking system and the assessment of data quality. These are briefly discussed below. For more information about the methodological considerations associated with the different data collection methods, see [Section 4.3](#).

### 3.3.1 Creating a Data Tracking System

Data tracking systems can be sophisticated web-based databases accessible to multiple users spread out across many sites or they can be basic spreadsheets operated and maintained by a single person. Whichever type of tracking system the evaluator chooses will depend on the size and complexity of the evaluation and the available resources. In general, when making decisions about which type of tracking system the evaluation needs, the evaluator should ask four questions:

- What kind of database do I want to use? (e.g., a simple database such as Microsoft Excel, or a relational database such as Microsoft Access);
- How do I want to input the data? (e.g., manual entry by a project team member, or a web form linked to the database);
- How do I want to update the data? (e.g., point-in-time through manual entry, or real-time through web forms); and
- How do I want to analyze the data? (e.g., from within the database using queries and reports to generate descriptive analyses, or by exporting the data to a statistical analysis program such as SPSS/SAS/STATA to conduct descriptive and inferential analyses)

The answers to the questions outlined above will depend in part on the types of data that need to be included in the tracking system. For the evaluation of an OSEP project, we recommend including the following types of data in the tracking system:

- Project characteristics (e.g., specific project activities participated in);
- Individual/family characteristics (e.g., a unique ID number for each person served, demographic info, contact information, and any other characteristics that might be important to include in the analysis);
- Student performance data (e.g., student achievement growth, if available);
- Comparable data for a comparison group (if applicable);
- Timelines for data collection activities (e.g., contacts with participants, administration of surveys or tests, conduct of observations); and
- Data collection responsibilities (e.g., preparing the IRB application, developing data collection instruments, contacting graduates, contacting graduates' schools, obtaining consent, collecting data, inputting data into the database, updating data, preparing data for analysis, conducting data analysis, writing the report).

If the evaluation includes measures of performance of students who have been involved with the project, the data system also should include characteristics of the students, student performance data, and data for a student comparison group (if applicable). **Important note: In order to link the performance of students or other individuals with project activities, the evaluator will need to ensure that the data system includes a unique ID number for each student or individual that can be linked to specific project activities.** Some districts and states have their own data systems that include these unique ID numbers; if they aren't available, the evaluator can work with the district data manager to generate unique ID numbers, or can generate study ID numbers themselves. As previously mentioned, when obtaining access to student data, evaluators will need to obtain approval from the district and school, as well as from the students' parents (see [Section 3.1.2](#) for more information). Further, when working out the data sharing agreement with the state or district (see [Section 3.1.3](#)), evaluators will need to be able to demonstrate that they have a secure way to transfer and store student data.

When creating a data tracking system, evaluators should keep in mind the following tips:

- Automatic reminders can facilitate timely data collection (e.g., the database generates emails to graduates after a certain time if they haven't entered their student data);
- Pull-down boxes within the database reduce risk of incorrect data entry (see Information on how to do this in Microsoft Excel is available on the [Microsoft support website](#)); and
- Consider creating "checkpoints" to control data entry (e.g., person inputting data must check a box for "obtained consent" before being able to enter student data).
- Consider including controls for data entry that require respondents to enter data in the format desired (e.g., in a survey, rather than have a text field for a person's date of birth that allows a person to enter their DOB in whatever way they choose, have a date field that requires the person to enter data in date format).

Another step in the management of data collection is the assessment of data quality, discussed below.

### 3.3.2 Assessing Data Quality

Before moving on to the analysis and reporting of the data collected, we recommend that evaluators review the quality of the data. This includes thinking about each variable in the dataset and examining the available data. In short, there are several data characteristics to review in assessing data quality:

- **Precision**—To what extent do the data collected reflect an exact measurement (for quantitative data) or include the narrative information (for qualitative data) needed to respond to a question?
- **Accuracy**—To what extent do data reflect the actual value of an observation or achievement? For example, is a measure of height or weight accurate or "off" by several inches or pounds? Do the data provide a "true" (e.g., verifiable account of a phenomena or experience)?
- **Reliability**—To what extent can a data measurement be replicated with accuracy and precision? When gathering data from individuals, might a person have any reason to respond falsely?
- **Consistency**—To what extent do data or individuals' responses agree with each other?
- **Completeness**—To what extent is complete information provided? For example, is the unit of measurement provided to help interpret the data? Is there enough context to understand an individual's response?

- **Legitimacy**—Is a data value or narrative response reasonable for a specific question or prompt? That is, are the data values logical given the context of the program and the specific question or data prompt? For example, if the question or prompt asks for respondent’s gender, does the respondent indicate either male or female?<sup>39</sup> If a question or prompt requests the length of time a teacher has been working, do mathematical calculations confirm that a teacher hasn’t been working longer than he or she has been alive?<sup>40</sup>

These types of edit checks help to verify the accuracy and quality of the data. We offer the following additional suggestions related to assessing the quality of the data.

- **Identify opportunities to check for data consistency and legitimacy and deploy the checks.**
  - **Identify questions that ask for related data in order to check for data consistency.** Data skip patterns should be included in these checks. For example, if a teacher indicates in Question 7 of a survey that she doesn’t teach students with autism, and Questions 8 and 9 ask further questions about teaching students with autism, there should be no data for Questions 8 and 9 for this teacher.
  - **Identify the set of possible responses for each question.** Review data frequencies to ensure all reported values fall within the set of possible responses. In some cases, the set of possible responses is readily apparent. For example, if the data item reports a date, evaluators should verify that the values for the date fields fall within acceptable limits (e.g., values for the “month” field have two digits between 01 and 12). In other cases, the evaluator’s knowledge and expertise will inform decisions about the set of possible values and the identification of outliers. Note: in this latter case, a value that falls outside what would normally be considered a possible value shouldn’t be discounted or determined to be an outlier without further investigation.
  - **Identify the correct format for responses.** Determine which data format will be needed to conduct the appropriate analyses and check that the data are all in the correct format. For example, if the evaluator wants to use an equation to calculate the time period that elapsed between two project-sponsored trainings, the data will need to be in a numerical format (e.g., 08162016) and not a text format (e.g., August 16, 2016).
- **Decide how to handle missing data.**
  - Is the amount of missing data extensive? If one respondent is associated with multiple missing values, does this warrant exclusion of all of the respondent’s data? Will statistical techniques for imputing missing values be employed? (See [Section 4.4.1](#) for more information on dealing with missing data).
- Calculate basic descriptive statistics for each variable and review data frequencies for outliers (see [Section 4.4.2](#)).
  - **Outliers** are data values that fall outside of the bulk of the data that were received. Outliers may reflect data entry errors or an extreme response reported by a respondent. In general, the study team should discuss rules for identifying and reacting to outliers. For example, evaluators may consider a value that is more than 20% away from the next closest data value to be an outlier. Will the evaluator keep the outlier in the dataset as a valid value? Will the evaluator choose statistical techniques that are relatively insensitive to the presence of outliers?

<sup>39</sup> This assumes that the survey does not provide an option for an alternate way for respondents to identify their gender (e.g., transgender).

<sup>40</sup> This assumes that the evaluator has obtained the respondent’s date of birth or age in order to conduct the calculations.

In assessing data quality, evaluators may need to identify a software package to help implement this plan. A number of options are available, including:

- Microsoft Excel: A product that includes a number of data validation options including being able to easily sort variable columns and establishing logic rules that highlight cells that do not meet specified criteria.
- Open Refine: A free online application that provides a user friendly interface with data editing features. The application is downloaded to a computer desktop so evaluators do not need to worry about sharing confidential data. The advantage to Open Refine is that it allows easy editing of data and keeps a log of edits. [Tutorial videos](#) are available to explain many of the features available.

Finally, evaluators should establish rules that will be followed in cleaning the data and fully document those procedures. Data cleaning refers to the procedures that will be followed to make the data ready for analysis. For example, evaluators may want to code raw data into a defined set of values. Evaluators also may find that they need to transform some variables into new variables for the purpose of analysis. It's important to document all of these decisions and data transformations (and coding). Further, evaluators may find that data cleaning raises additional questions about specific data values or responses. In such cases, evaluators may be able to contact respondents to clarify responses or may elect to code some inconsistent data as missing.

## 3.4 Analyzing the Data

It's beyond the scope of this Toolkit to discuss the myriad methods of quantitative and qualitative data analysis. [Section 4.4](#) outlines some methodological considerations that should be kept in mind when conducting data analysis. In this section we highlight some important aspects related to preparing the data and briefly discuss how to aggregate the data and report the results.

### 3.4.1 Preparing the Data for Analysis

In connection with the data collection and data quality review activities, evaluators will have entered the data into a database and begun preparing them for analysis. (Note: the information here applies primarily to quantitative data; see [Section 4.4.4](#) for a discussion of qualitative analysis.) This process includes the following:

- Checking for duplicate records
- Identifying outliers
- Identifying the different types of measurement scales that are present in the quantitative data (e.g., nominal, ordinal, interval, ratio)
- Determining what types of variables the data represent (e.g., categorical, continuous)
- Assigning a numeric score to each response category for each close-ended question in a survey, item in a structured observation protocol, or question in a structured interview (e.g., 0 = Poor, 1 = Fair, 2 = Good, 3 = Excellent or 0 = No, 1 = Yes) (Note: Data obtained in unstructured observations or interviews will need to be analyzed qualitatively and codes may be developed; if desired, numbers can be assigned to the codes but generally they should be treated as nominal data.)
- Determining how to code missing data (e.g., assigning a code of 99 or 999 to missing data will cause a problem when calculating means)
- Recoding variables such as negatively worded survey items so they will be consistent with the positively worded item coding
- Reviewing qualitative data gathered during open-ended survey questions, qualitative observations and unstructured interviews for later analysis.(see [Section 4.4.4](#) for information on qualitative data analysis)

Since data analysis is a complex topic, we do not cover it here (see [Section 4.4](#) for information on data analysis). Instead, we limit our focus to creating the final variables for quantitative analysis, discussed next.

## 3.4.2 Aggregating Data and Reporting Results

The data analysis plan (see [Section 2.5.2](#)) will guide most of the work during the analysis and reporting phases of the evaluation. In particular, this plan will help the study team to organize the output into tables and charts to answer the evaluation questions. Remember, the analysis plan identifies the specific variables and analysis techniques that will be used to respond to each evaluation question.

Data aggregation is the set of procedures that evaluators will use to combine data from multiple respondents or multiple items in order to report study findings. The data aggregation steps will vary by each evaluation question and analysis approach.

### 3.4.2.1 Aggregating and Reporting Percentages

We generally report percentages if we are working with nominal or ordinal data. For example, if an evaluation question requires reporting the percent of a target population that achieves a specific benchmark or achievement, evaluators should complete the following aggregation steps:

1. Identify which responses are eligible for inclusion in the numerator and denominator. In particular, consider the following questions:
  - a. Are only those data values that pass data quality assessment checks to be included?
  - b. How will missing data be handled in the calculation?
2. Calculate the numerator.
  - a. Identify the range or set of values that qualify a data value for inclusion in the numerator, then calculate the numerator by *counting* the number of eligible responses in which the qualified data values exist.
3. Calculate the denominator.
  - a. Identify the range or set of values that qualify a data value for inclusion in the denominator, then calculate the denominator by *counting* the number of eligible responses.
4. Conduct the calculation by dividing the numerator by the denominator and multiplying by 100.

In reporting the results, evaluators should be sure to not only report the required percent, but also the total number included in the calculation as well as the number of responses that are *not* included because of the data quality assessment and data eligibility checks. If the results are limited to a subset of the target population, evaluators should include language that describes the true scope of the findings.

### 3.4.2.2 Aggregating and Reporting Means

We generally report means if we are working with interval or ratio data. For example, if an evaluation question requires reporting a mean value for a target population or benchmark, evaluators should complete the following aggregation:

1. Identify which responses are eligible for inclusion in the numerator and denominator. In particular, consider the following questions:
  - Are only those data values that pass data quality assessment checks to be included?
  - How will missing data be handled in the calculation?
2. Calculate the numerator.
  - Identify the range or set of values that qualify a data value for inclusion in the numerator, then calculate the numerator by *summing* the total eligible responses in which the qualified values exist.
3. Calculate the denominator.
  - Identify the range or set of values that qualify a data value for inclusion in the denominator, then calculate the denominator by *counting* the number of eligible responses.
4. Conduct the calculation by dividing the numerator by the denominator.

In reporting the results, evaluators should be sure to not only report the required mean, but also the range, the standard deviation, and the number of responses that aren't included because of the data quality assessment and data eligibility checks. If the results are limited to a subset of the target population, evaluators should include language that describes the true scope of the findings.

### 3.4.2.3 Conducting and Reporting Results of Inferential Statistical Analyses

It's beyond the scope of this Toolkit to discuss the various methods of conducting and reporting inferential statistical analyses. We recommend that evaluators without statistical training consult a statistician to help with these analyses. [Section 4.4.2.2](#) presents a brief discussion of inferential statistical analysis. Hinkle, Wiersma, and Jurs (2003) and Dimitrov (2010) are good reference books and Rice University, the University of Houston Clear Lake, and Tufts University have developed an [online statistics book](#) that is free and available to the public. Evaluators also can find information about statistical analysis online at the [Web Center for Social Research Methods](#).

## 3.5 Reporting Findings

### 3.5.1 Providing Formative Feedback

Once the evaluation begins, it is helpful for project staff and evaluators to agree to a communication structure and schedule that identifies the logistics of how and when team will communicate. In-person meetings; email or other written communications; and telephone-, web-, or video-conferences are all viable communication options. The frequency of communication will likely depend upon the stage of the evaluation; project staff can expect a significant amount of communication at the beginning of the evaluation—when the evaluator is learning about the project and the evaluation—and during reporting periods.

It's essential that project staff communicate any potential deviations in project implementation to the evaluator so that any necessary adjustments can be made to the evaluation. Issues to be discussed may include implementation challenges such as low enrollment in project activities, project staff turnover, or changes to the scope of the project activities. On the other hand, the evaluator should report potential challenges to implementation of the evaluation to project staff, including difficulties recruiting participants for the evaluation, problems making arrangements for data collections, low survey response rates, or evaluation staff turnover.

Generally, evaluators provide formative feedback to help project staff monitor project progress and quality, inform on-going project activities, and make mid-stream adjustments. Some of the ways formative feedback can be used include:

- Identifying gaps in implementation and project support needs;
- Assessing level of fidelity, with the purpose of identifying different levels of fidelity (e.g., low/medium, high; inadequate/adequate) so as to compare differences across contexts, sites, or groups or inform outcome findings (see [Section 4.5](#) for a discussion of how to measure fidelity);
- Ascertaining barriers to or facilitators of achievement of short-term or intermediate outcomes;
- Assuring services are considered high quality, relevant, timely, and useful by service clients or recipients;
- Measuring achievement of short-term and intermediate outcomes;
- Determining progress toward (or likelihood of) achieving long-term outcomes; and
- Making adjustments prior to replication and dissemination.

This feedback may come in the form of monthly updates, quarterly reports, periodic briefings (e.g., informal updates during project meetings or more formal presentations of preliminary results to key stakeholder groups), or annual reports.

### 3.5.2 Preparing the Final Report

The last step in the evaluation is the preparation of the final report. The content of the report may be established in the request for applications to which the study team responded or in the request for proposals specifically for the evaluation. If not, we recommend including the following sections in the report to allow readers to fully review and interpret the findings and to allow for replication of the project.

- Table of Contents
- List of abbreviations and acronyms
- Executive Summary—Usually contains a basic description of the program, the evaluation questions, and key findings.
- Introduction—Provides a more detailed background of the program and identifies the evaluation questions.
- Methodology—Describes in detail the evaluation plan, data collection activities, sample, and data analysis and data quality assessment plan for each component of the evaluation.
- Analysis and Results—Presents the findings for each evaluation question in narrative form, supplemented by tables and graphs, and summary statistics.
- Conclusions and Implications—Summarizes overarching findings and evaluation results.
- Study Limitations—Identifies limits of the data collection and analysis techniques that affect the generalizability and interpretation of findings. This is discussed in more detail below.

### 3.5.3 Outlining Study Limitations

It's important for evaluators to identify the limitations associated with their evaluation. Limitations are characteristics of the design or implementation that affect how the evaluation's findings should be interpreted and the extent to which the findings can be generalized to other programs or contexts. Identifying limitations should take place in both the planning and reporting phases of the evaluation. In the planning phase, evaluators should consider the implications of the *a priori* methodological choices.<sup>41</sup> In the reporting phase, evaluators should reflect on any unanticipated limitations that arose during the course of the evaluation. This process will alert the reader to the need to exercise caution when interpreting the results. Further, it will provide information that may help the study team to improve the design and data collection and analysis techniques in future evaluations.

The brief survey in Box 1, on the next page, is designed to guide analysis of each phase of the evaluation to identify potential limitations. In discussing the findings of the evaluation, evaluators should describe these limitations and how they may affect interpretations of the study results.

---

<sup>41</sup> See for example Reybold, Lammert & Stribling, 2012.

## Box 1. Study Limitations Survey

### *Planning and Design*

- 1) Consider the overall evaluation design. Were other evaluation designs possible, if any, given the time and resources you had available for the evaluation? If so, why did you choose the design that you did? Would other study designs have generated findings with fewer limitations?
- 2) Consider who was included in data collection. Was there a control or comparison group? If no, why not? How does this affect the strength of your findings? Did you draw a sample or samples? What was the basis for your samples? Were there specific groups intentionally or unintentionally excluded? What other samples could have been drawn? What impact would these different samples have had on your evaluation?
- 3) Consider the instrumentation. Did your instruments provide the exact data you needed to respond to your evaluation questions? If not, where were the data less than exact? How does this affect your ability to respond to the evaluation question? Do the results from your evaluation suggest potential problems in the instrumentation? For example, did survey items or scale scores show enough variability? Would changing the response options or rubrics offer better data analysis options in the future?

### *Data Collection*

- 4) Examine response rates. What were your response rates? Is there a pattern to response or non-response? Did participation in your survey vary based on any participant characteristics? For example, were principals at higher performing schools more likely to return surveys? If response rates were poor in general or among a sub-group, how well can you generalize the results of the evaluation?

### *Data Management and Data Quality*

- 5) Consider your coding and data cleaning procedures. Did you check for errors in coding and data entry? If no, what potential errors could have been made in coding and entering the data?
- 6) If more than one rater was involved in collecting data, examine inter-rater reliability. What is your inter-rater reliability? Is there a lot of variation among raters? Did one or more raters provide consistently high ratings while one or more raters provided consistently low or lower ratings? Did you make any adjustments for rater-driven variation in ratings? Are there other patterns in your data that appear correlated with rater or data collector? How do such patterns affect your analyses and findings?
- 7) Consider missing data. Was there a high amount of missing data? How did you address missing data? What are the implications attached to your handling of missing data?

### *Data Analysis*

- 8) Consider your analysis options. Did you use the most rigorous analytic technique possible to respond to the evaluation question? If not, why not? How did the type and amount of data collected affect your ability to analyze the data and choice of analysis technique? What other variables need to be explored and collected to further your analysis of the evaluation data?

In the sections that follow we discuss a number of methodological considerations that project evaluators should keep in mind throughout the course of the evaluation.

## 4 Methodological Considerations

This section provides additional information about different methodological aspects of planning and conducting a project evaluation. Throughout the section our goal isn't to present a comprehensive discussion of each topic, but rather to give a brief overview and highlight important points or potential issues that should be taken into consideration. Where possible, we provide links to online resources and throughout we offer suggestions on good print or electronic resources that evaluators can look to for additional information.

### 4.1 Evaluation Design

As mentioned in [Section 2.5.1](#), the type of design chosen for a project evaluation will depend upon, among other things, the goals of the study and the outcomes that are identified in the project logic model, since different study methods are required to measure different types of outcomes. We anticipate that in many cases a project evaluation will feature a mixed-method design that incorporates multiple study methods, such as non-experimental studies (e.g., surveys or qualitative case studies) and quasi-experimental or experimental studies of outcomes.<sup>42</sup>

In writing this section we do not intend to prescribe a particular course of action for evaluators to follow. Rather, we aim to highlight different design alternatives that we believe can be implemented successfully during a project evaluation, taking into consideration the varying contexts and diverse constraints project evaluators face.

#### 4.1.1 Randomized Experimental Designs

Since we believe that it will be difficult for evaluators of OSEP projects to utilize RCTs in their studies, we will not focus on the different randomized designs here and instead refer interested readers to Shadish, Cook and Campbell (2002) for details on the different types of designs. However, we do want to highlight an important consideration for evaluators who plan to conduct RCTs as part of their study—the need to calculate attrition.

##### 4.1.1.1 Calculating Attrition

Attrition is the loss of response from participants that takes place after the participants have been assigned to conditions. This might occur if a respondent fails to respond to a particular survey question or if the respondent refuses to participate in the study after the project begins. If the evaluator drops a respondent from the study for one reason or another it also should be counted as attrition, *if the drop could have been caused by the treatment*.

The primary problems with attrition are that it:

- lowers statistical power to detect effects, and
- any attrition from specific treatment conditions that appears to have been associated with the treatment (also called *differential attrition*) can threaten the internal validity of the study (see [Appendix B](#) for validity threats). When overall and differential attrition are high, the integrity of the randomized experiment is lost, thereby invalidating the assumption of equivalence among treatment and control groups.

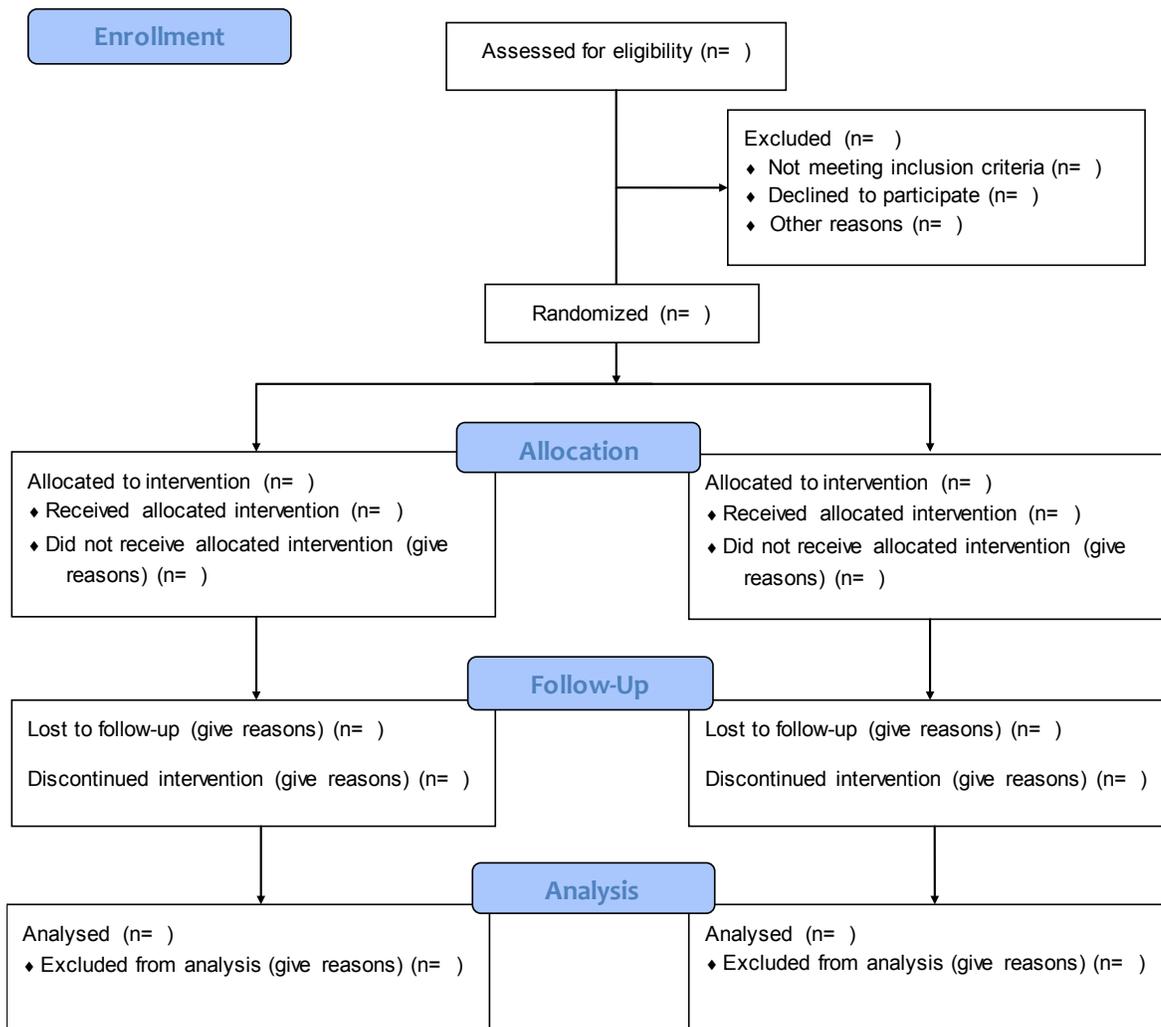
---

<sup>42</sup> See [Section 2.5.1.5](#) for a discussion of mixed-method designs.

The best strategy for dealing with attrition is to prevent it in the first place. Shadish, Cook and Campbell outline a number of strategies that can be used to retain and track participants in a study.<sup>43</sup> In most cases, however, it's not possible to completely prevent all attrition from occurring. For this reason, we recommend that evaluators use a tool such as the Consort Diagram presented in Figure 3 to track the participants in their experimental studies (we also recommend tracking participants in other study designs, but the problem of attrition relates specifically to randomized designs).

Figure 3. Sample Consort Diagram

**CONSORT 2010 Flow Diagram**



Source: CONSORT (<http://www.consort-statement.org/consort-statement/flow-diagram0/>)

In the next section we talk about quasi-experimental designs, followed by a discussion of single-case designs.

<sup>43</sup> 2002, pp. 225-334.

In the following sections we highlight a few quasi-experimental designs that we believe can be applied in evaluations of OSEP-funded projects.<sup>44</sup> This isn't to say that evaluators shouldn't try to use other types of designs, especially if those designs are more rigorous, but we simply are recognizing the constraints that evaluators of such projects may face when making design decisions. For each design we highlight validity threats (see [Appendix B](#) for a list of validity threats) and any available strategies to improve the strength of the design for making causal claims that were identified by Shadish et al. (2002). To describe these designs we will use the notations provided by Campbell and Stanley (1963)<sup>45</sup>:

- X represents an exposure of a unit to an experimental variable or event [i.e., the “treatment”], the effects of which are to be measured.
- O represents an observation or measurement recorded on an instrument [such as a standardized assessment, a teacher-made test, a survey, or a psychological scale]. The subscript following an O indicates a different time period.
- X's and O's in a given row are applied to the same specific units. X's and O's in the same column, or placed vertically relative to each other, are simultaneous [i.e., taking place at the same moment in time].
- The left-to-right dimension indicates the temporal order of procedures in the experiment (sometimes indicated with an arrow).
- Separation of parallel rows by a dashed horizontal line indicates that comparison groups aren't equal (or equated) by random assignment. No dashed horizontal line between the groups displays random assignment of individuals to treatment groups.

#### 4.1.2.1 One-Group Posttest-Only Design

In this design the study collects one posttest observation on respondents who experienced a treatment.

X      O<sub>1</sub>

This design can be acceptable for summative purposes in cases in which there is significant specific background knowledge about how the dependent variable might behave following a treatment. For example, if a child receives a well-established speech and language therapy intervention, and it's known that the child likely wouldn't otherwise have developed the skills taught during the intervention, this design could be appropriate.

However, in most cases, with this design it's not possible to determine if a change has occurred following the implementation of a treatment and it's not possible to identify the counterfactual. Almost all threats to internal validity are possible with this design. However, *“for valid descriptive causal inferences to result [from this design], the effect must be large enough to stand out clearly, and either the possible alternative causes must be known and be clearly implausible or there should be no known alternatives that could operate in the study context.”*<sup>46</sup>

<sup>44</sup> See Shadish et al., 2002, for additional discussion of these designs.

<sup>45</sup> Cited in Creswell, 2003, pp. 168-9.

<sup>46</sup> Shadish et al., 2002, p. 107.

#### 4.1.2.2 One-Group Pretest-Posttest Design

The one-group pretest-posttest design is a slight improvement over the one-group posttest-only design because it provides some information about the effect of the treatment on the outcome.

O<sub>2</sub>      X      O<sub>2</sub>

There are multiple threats to internal validity with this design, including maturation, history, testing, and attrition. Researchers will *“rarely be able to construct confident causal knowledge with this design unless the outcomes are particularly well behaved and the interval between pretest and posttest is short.”*<sup>47</sup> However, adding another pretest to the design can reduce the plausibility of some of the validity threats, since it would give some idea of the trend of performance on the outcome measure prior to the implementation of the intervention.

O<sub>1</sub>      O<sub>2</sub>      X      O<sub>3</sub>

Another way to improve this design is to add a non-equivalent dependent variable, diagrammed below.

O<sub>1A</sub>      O<sub>1B</sub>      X      O<sub>2A</sub>      O<sub>2B</sub>

In this design, measures A and B assess similar constructs, but measure A (the outcome) is expected to change because of the treatment, while measure B (the nonequivalent dependent variable) isn't. Measure B is expected to respond to plausible internal validity threats in the same way as measure A would, so changes in measure B would illustrate whether these validity threats are actually operating within the study. For example, if the students' math scores (measure B) rose at the same rate as reading scores (measure A) when the treatment was solely focused on reading, something other than the treatment might be causing the reading score increases. While this design reduces the plausibility of many threats to internal validity, history remains a threat in this study.

---

<sup>47</sup> Ibid., p. 110.

#### 4.1.2.3 Removed Treatment Design

In this modification of the one-group pretest-posttest design, a second posttest is added ( $O_3$ ), then the treatment is removed and a third posttest ( $O_4$ ) is administered.

$O_1$  X  $O_2$   $O_3$  ✕  $O_4$

The objective of this design is to show that the outcome changes with the presence or absence of the treatment. It's assumed that it would be difficult for many internal validity threats to operate in the same way, but the presence of outliers in the data can affect the results. Making the observations at equally-spaced intervals allows the examination of linear changes over time.

#### 4.1.2.4 Repeated Treatment Design

In this design the researcher introduces a treatment, administers a posttest ( $O_2$ ), then removes the treatment, administers another posttest ( $O_3$ ), and then introduces the treatment again, followed by another posttest ( $O_4$ ). This design and the one immediately above are commonly used by psychologists conducting behavioral research, usually with a single-case design (see [Section 4.1.2](#)).

$O_1$  X  $O_2$  ✕  $O_3$  X  $O_4$

Again it's assumed that few threats to internal validity could explain a relationship that holds over this pattern of treatment introductions and removals. This design will not be able to show a causal relationship, however, if the treatment effects aren't transient, or if the treatment creates a ceiling effect.

#### 4.1.2.5 Posttest Only With Non-Equivalent Groups Design

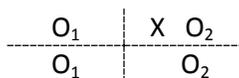
This design is often used when the treatment begins before the start of a study. In this design, the evaluator selects a comparison group and administers a posttest to both the treatment group and the comparison group to see how their performance differs.

Group A	X	O
Group B	-----	O

In general this design provides weak evidence of a causal relationship between the treatment and the outcome, since it's hard to determine if pre-existing group differences affect the outcome on the posttest. Nevertheless, it does provide some measure of a counterfactual.

#### 4.1.2.6 Post-test Only Design Using an Independent Pretest Sample

If it's not possible to collect pretest and posttest data for the same group, it's sometimes possible to collect pretest data for a randomly formed independent sample that is drawn from the same population as the posttest sample. In the diagram below, the dashed vertical line indicates that the two observations are taken from two different samples over time.



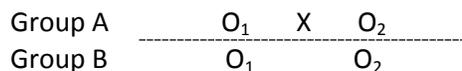
It's important to randomly select the pretest and posttest samples from the same population (e.g., the total number of families who received services from a Parent Information Center) or this design will introduce considerable selection bias into the results. This design should be used only when there is a strong need to have independent pretest and posttest groups or when there are significant problems collecting pretest and posttest data from the same groups over time.<sup>48</sup>

Some other ways to improve this design include:

- Matching or stratifying (see [Section 4.1.2.8](#));
- Using internal control groups (e.g., individuals who registered for a class too late to be accepted into the study, but who are otherwise likely to be similar to students who are in the study); and
- Using multiple (non-equivalent) control groups (e.g., a group that is expected to underperform the treatment group if the treatment has no effect and a group that is expected to outperform the treatment group if the treatment has no effect).<sup>49</sup>

#### 4.1.2.7 Non-Equivalent Pretest-Posttest Control-Group Design

This is probably the **most common of all quasi-experimental designs**. In this design, the treatment and control groups are selected without random assignment and the treatment is administered only to one group (Group A). Both the treatment and control groups are administered a pretest and a posttest.



With this design there is still a possibility of selection bias, since the groups are nonequivalent, but the pretest allows the researcher to explore the possible size and direction of that bias. We recommend that evaluators **calculate baseline equivalence** of the treatment and comparison groups using the pretest data, discussed briefly in [Section 4.1.2.9](#).

This design also allows for an examination of attrition, which gives researchers an opportunity to see if there are differences among units that stay in and leave the study. The plausibility of different validity threats depends in part on the observed pattern of outcomes. This design can be improved by adding additional pretests, using switching replications, or using a reversed treatment control group.<sup>50</sup>

<sup>48</sup> Shadish et al., 2002.

<sup>49</sup> See Shadish et al. (2002) for a full discussion of these options.

<sup>50</sup> See Shadish et al., 2002 for additional discussion of this.

#### 4.1.2.8 Matching Treatment and Comparison Groups

To decrease the odds of selection biases, evaluators can use matching to form comparison groups. This involves grouping units with similar scores on a matching variable (e.g., school size, ethnicity) so that treatment and comparison groups both have units with the same (or very similar) characteristics on the matching variable.<sup>51</sup> Different types of methods for matching include:

- Exact matching—when units have exactly the same score within a match (in practical terms this type of matching isn't common, since it requires units to have the exact same scores);
- Caliper matching—when units have scores within a pre-defined distance of each other;
- Index matching—when multiple comparison units above and below a treatment unit are selected;
- Cluster group matching—when cluster analysis is used to embed the treatment group in a cluster of similar control units;
- Benchmark group matching—when control units that fall close to the treatment unit on a multivariate distance measure are selected;
- Optimal matching—when each treatment unit has multiple control units and vice versa;
- Cohort matching—when successive groups go through a particular treatment. Cohort matching is particularly useful *if*
  - one cohort experiences a treatment and earlier or later cohorts do not;
  - cohorts differ in only minor ways from their contiguous cohorts;
  - organizations insist that a treatment be given to everybody, thus precluding simultaneous controls and making possible only historical controls;
  - an organization's archival records can be used for constructing and then comparing cohorts;<sup>52</sup> and
- Propensity score matching—when treatment and comparison units are matched based on the conditional probability of receiving the treatment given a set of observed covariates. **Propensity score matching is gaining popularity among researchers and evaluators and is considered by many to be the preferred method.**<sup>53</sup>

While matching generally improves the similarity of the treatment and comparison groups, there is still the possibility that the groups may differ on some variable not included in the matching. Additionally, in some situations it's possible for matching to produce a result that is actually further away from the correct answer than if no matching had been used.

Some general principals to improve matching include,

- Identifying a possible comparison group that appears to be very similar to the treatment group before conducting the matching of individuals; and
- Using matching variables that are stable and reliable, and that are correlated with the outcome variable.<sup>54</sup>

---

<sup>51</sup> Shadish et al., 2002.

<sup>52</sup> Ibid., p. 149.

<sup>53</sup> For more information on propensity score matching, see, for example, Barth, Guo, & McCrae, 2008; Heinrich, Maffioli, & Vazquez, 2010; and Luellen, Shadish, & Clark, 2005.

It's important to remember that even though the treatment and comparison groups may be similar at the time of matching, evaluators still must determine whether the groups remain similar on important characteristics at the end of the study, after some of the participants have left the study. This can be achieved by calculating baseline equivalence, discussed next.

#### 4.1.2.9 Calculating Baseline Equivalence

Whenever a QED features a comparison group and pretest data are available, it's important to determine whether the treatment and comparison groups are similar in important ways at baseline. At its most basic, the calculation of baseline equivalence refers to calculating differences in means on a pretest measure between the treatment and comparison groups. If it's not possible to administer a pretest during the evaluation period, it might be possible to calculate baseline equivalence using another baseline measure that is correlated with the outcome measure (such as the prior year's state assessment). A key point is that evaluators should calculate baseline equivalence for the *analysis sample*—that is, the sample of individuals in the treatment and comparison groups that remain in the study for the entire study period. The [What Works Clearinghouse standards](#) include methods for establishing baseline equivalence in experimental studies with high attrition and in quasi-experimental studies.<sup>55</sup>

We now turn our discussion to single-case/single-subject designs, which are commonly used by special educators and related service providers to demonstrate changes in the performance of their students.

---

<sup>54</sup> Shadish et al., 2002.

<sup>55</sup> U.S. Department of Education, 2014.

### 4.1.3 Single-Case/Single Subject Designs

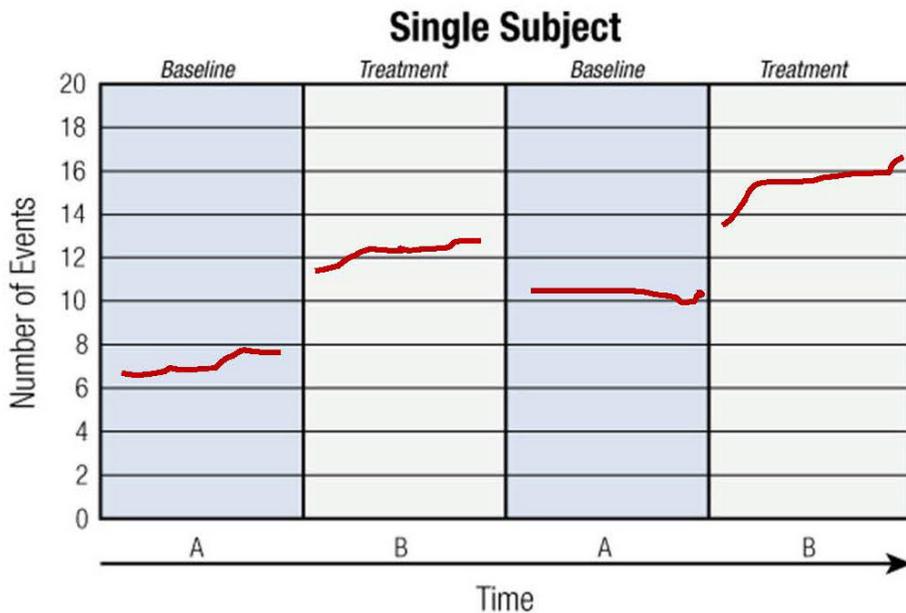
As we mentioned in [Section 2.5.1.3](#), the two most common single-case designs are the A-B-A-B design and the multiple baseline design; these are presented briefly below. We encourage project staff and evaluators to consult Kennedy (2005) and Todman and Dugard (2001) for additional information about the different types of designs and how to conduct single-case studies.

#### 4.1.3.1 A-B-A-B design

An A-B-A design (where A = baseline and B = intervention) is the minimal type of experimental design that can establish experimental control in single case research, but researchers tend to prefer the A-B-A-B design for at least two reasons. First, the A-B-A-B design allows for two different replications of changes in pattern (i.e., when the baseline is reintroduced and when the intervention is reintroduced). Second, many researchers prefer to end a study in a way that is the most beneficial to the participant; if the intervention is working, the A-B-A-B design allows the participant to continue receiving the intervention as the study comes to a close.<sup>56</sup>

An important issue when using A-B-A-B designs is whether behavior will return to baseline levels after the intervention is withdrawn. If this doesn't occur, experimental control may be lost and no functional relation between the intervention and a change in behavior is demonstrated. Figure 4 demonstrates what data collected in an A-B-A-B design might look like.<sup>57</sup>

Figure 4. A-B-A-B Single-Case Design



Source: Adapted from <http://www.cehd.umn.edu/nceo/Onlinepubs/Technical26.htm>

<sup>56</sup> Kennedy, 2005.

<sup>57</sup> Ibid.

As can be seen, there are clear differences among the number of events that take place during the baseline phases compared to the treatment phases. However, the data for the second baseline show that the participant has not completely reverted to baseline levels. This might occur if the intervention introduces a new skill to the person; once the skill is learned it's difficult to reverse the effects of instruction. When this happens, it might be better to add additional baseline and intervention phases to see if the pattern of behavior change during the different phases continues.<sup>58</sup>

It's common in educational settings for a researcher to conduct an experiment when an intervention is already in place. In such cases, it's possible to use a B-A-B (or some extension of it). As with the A-B-A-B design, the establishment of experimental control in the B-A-B design depends on the variable of interest being sensitive to the withdrawal of the intervention. If behavior doesn't change during the withdrawal or reversal phase, a functional relation has not been established.<sup>59</sup>

#### 4.1.3.2 Multiple Baseline Design

The advantage of the multiple baseline design is that it doesn't require withdrawal, reversal, or alternation of treatment conditions (not discussed here). Rather, two or more baselines are concurrently established and the intervention is introduced sequentially for each participant. This is an important alternative in situations in which it isn't possible to remove the effects of the intervention once it has been introduced (e.g., when a student learns how to read it's unlikely that he or she will lose this ability).<sup>60</sup> Figure 5, on the next page, illustrates what data collected during a multiple baseline study might look like.

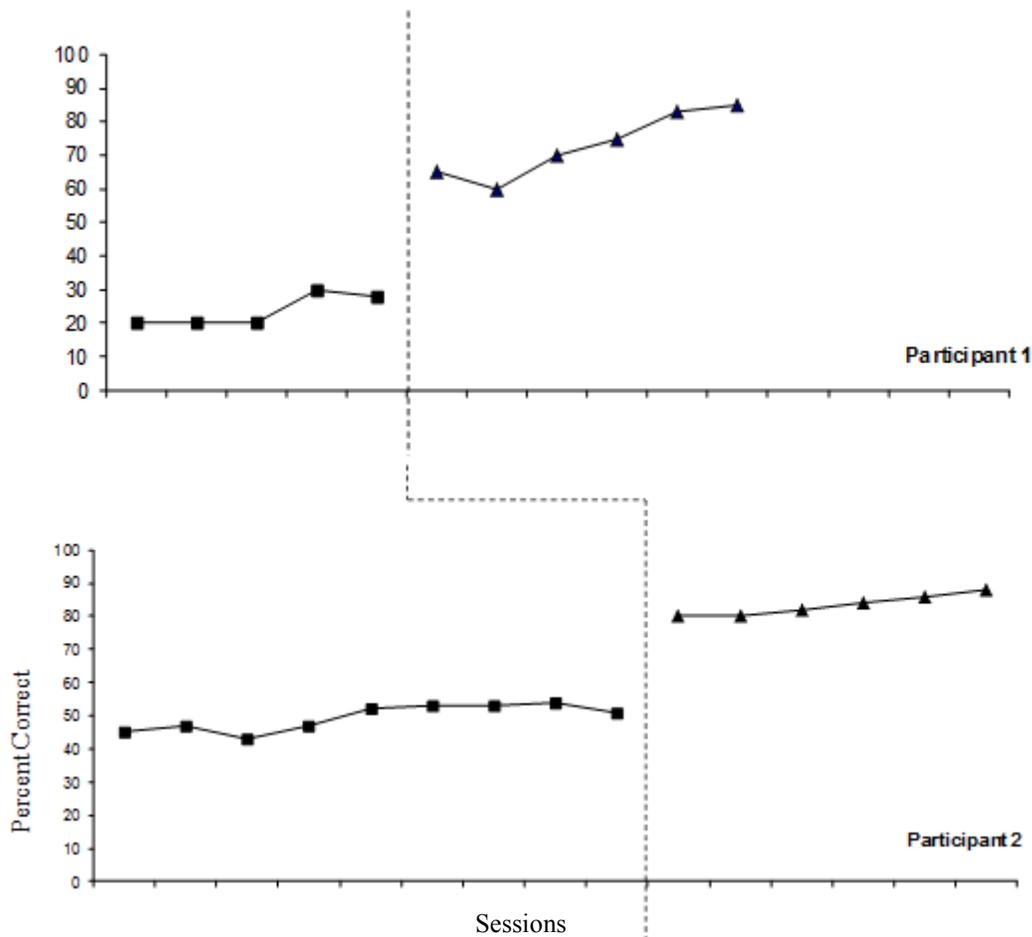
---

<sup>58</sup> Kennedy, 2005.

<sup>59</sup> Ibid.

<sup>60</sup> Kennedy, 2005.

Figure 5. Multiple Baseline Single-Case Design



The logic behind the multiple baseline design is that individual baselines are established for each participant in the study (or group, classroom, etc.), consistent response patterns are observed, and then the researcher introduces the intervention to one participant at a time. The researcher should wait until a clear pattern is observed for the first participant post-intervention before introducing the intervention to the second participant, and so on. A functional relation is demonstrated when changes in the dependent variable occur only when the independent variable (intervention) is introduced.<sup>61</sup>

One of the drawbacks of the multiple baseline design is the fact that some participants must continue without the intervention for longer periods of time. This type of design may not be well suited to situations in which the behavior being observed is particularly undesirable (e.g., extremely disruptive behavior in class); instead, an A-B-A-B or other design may be more appropriate.<sup>62</sup>

In the next section we discuss non-experimental designs.

<sup>61</sup> Kennedy, 2005.

<sup>62</sup> Ibid.

#### 4.1.4 Non-Experimental Designs

As mentioned in [Section 2.5.1](#), we recommend that project evaluators use experimental or quasi-experimental designs to answer summative evaluation questions whenever possible; however, we anticipate in some cases evaluators may need to use a non-experimental design, either alone or as part of a mixed-methods study. In the sections that follow we briefly discuss the different types of non-experimental designs that might be used for a project evaluation.

##### 4.1.4.1 Case Studies

Case study research can include single- and multiple-case studies and can incorporate a mix of quantitative and qualitative data. Multiple case studies may yield some nominal or categorical data. Yin (1994) outlined five different applications of case study research:

- To explain the causal links in real-life interventions that are too complex for the survey or experimental strategies (e.g., to link program implementation to program effects);<sup>63</sup>
- To describe an intervention in a real-life context in which it occurred;
- To illustrate certain topics within an evaluation in a descriptive mode;
- To explore those situations in which the intervention being evaluated has no clear, single set of outcomes; and
- To serve as a *meta-evaluation*—a study of an evaluation study (p. 15).

According to Yin (1994), case studies are preferable to other types of study designs “when ‘how’ or ‘why’ questions are being posed, when the investigator has little control over events, and when the focus is on a contemporary phenomenon within some real-life context” (p. 1). Shadish et al. (2002)—authors of one of the most well-known books about experimental and quasi-experimental research design—recognized that intensive qualitative case studies can be good non-experimental alternatives for generating causal conclusions, although they preferred the use of experimental or quasi-experimental designs whenever possible. Examples of case studies that may be possible in project evaluations include studies of individuals who received training in very small disciplines or fields of study. In these cases, the overall sample size for the evaluation may be very small, which affects the generalizability of findings to larger populations. In addition, the participants may have received such specialized training that a survey or other quantitative assessment also lacks generalizability.

See Yin (1994) for an extended discussion of the design and conduct of case studies.

---

<sup>63</sup> Shadish, Cook and Campbell (2002) acknowledged that intensive qualitative case studies can be good non-experimental alternatives to generating causal conclusions for three reasons: (1) “Journalists, historians, ethnographers and lay persons regularly make valid causal inferences using a qualitative process that combines reasoning, observation, and falsificationist procedures in order to rule out threats to internal validity” (p. 500); (2) qualitative methods can explore causation in a way that incorporates much more complexity than can be achieved with an experiment; and (3) intensive case studies can yield much more information to the researcher and to policy makers than experiments can. However, Shadish et al. also pointed out that many qualitative studies do not provide sufficient information about the counterfactual to sufficiently reduce uncertainty about cause. For more information about using qualitative research to make causal claims, see Donmoyer, 2012a, 2012b; and Maxwell 2004a, 2004b, 2011, 2012.

#### 4.1.4.2 Descriptive Studies and Surveys

Descriptive or survey research aims to provide systematic and accurate descriptions of selected characteristics for a population under study. This can be done by conducting surveys with samples drawn from a population (see [Section 4.2](#) for information on sampling) or by surveying all of the members of a certain population (e.g., all of the graduates of a special education teacher certification program at a Midwestern university during the period 2012-2015). Surveys most often generate nominal or ordinal data, although the responses to survey questions are often treated as interval data and analyzed accordingly. Descriptive studies are generally classified in terms of **how many times a sample is surveyed** (e.g., cross-sectional or longitudinal surveys) and **how the data are collected** (e.g., self-report surveys or observations). For example:

- **Cross-sectional surveys** involve one-time data collection with groups of individuals to compare groups at a single point in time. An example would be surveying all of the speech and language pathologists (SLPs) who are 1 year past graduation from a program and simultaneously surveying the group of SLPs who are 2, 3, 4, and 5 years past graduation. This study would be conducted all at one time and would provide a snapshot for each group of graduates.
- **Longitudinal surveys** involve collecting data multiple times from the same individuals, such as surveying the SLPs each year after graduation, obtaining data on the same individuals at 1, 2, 3, 4 and 5 years past graduation. This study would take 5 years to conduct.

It also is possible to combine cross-sectional and longitudinal surveys by starting a new longitudinal cohort at different points in time over the course of the evaluation.<sup>64</sup>

This type of research is well suited to answering “who,” “where,” “how much,” and “how many” questions—as is often the case in formative evaluations. Further, surveys can be particularly useful when the goal of the study is to “describe the incidence or prevalence of a phenomenon or when it is to be predictive about certain outcomes.”<sup>65</sup> Finally, surveys might be preferable when multiple questions must be answered yet limited resources are available. As Cronbach put it, in such situations “a survey might be preferred because it has a wider **bandwidth** that permits answering a broader array of questions even if the causal question is answered less well than it would be with an experiment.”<sup>66</sup>

#### 4.1.4.3 Correlational Studies

In general, correlational studies aim to **investigate relationships between variables** or use such relationships to make predictions about a variable of interest.<sup>67</sup> Correlational studies are commonly used in education—for example, to answer questions such as, “Do graduates with higher grade point averages obtain employment in the area for which they are qualified in less time than graduates with lower grade point averages?” or “Do the students of teachers who report high self-efficacy have higher levels of academic achievement than students of teachers who report low self-efficacy?” Correlational studies do not allow the researcher to determine which of the two variables being correlated came first, so it’s not possible to identify a causal relationship in such studies. Further, correlational studies do not test alternative explanations for a presumed effect, thereby leaving the possibility that some other variable not under study (often called a confound) might actually be responsible for the observed relationship.<sup>68</sup> This may also be referred to as a

---

<sup>64</sup> Dimitrov, 2010.

<sup>65</sup> Yin, 1994, p. 6.

<sup>66</sup> 1982, cited in Shadish et al., 2002, p. 98, emphasis in original.

<sup>67</sup> Dimitrov, 2010.

<sup>68</sup> Shadish et al., 2002.

“spurious correlation.” Nevertheless, the existence of a strong correlation between two variables may point to hypotheses about causal effects that can be explored in subsequent, more rigorous, studies.

#### 4.1.4.4 *Ex Post Facto Studies*

Ex post facto studies—studies that are conducted *after the fact* using secondary data—investigate cause-and-effect relationships by analyzing data on events that have taken place in the past. Ex post facto studies can be useful when (a) the independent variables cannot be controlled by the evaluator or (b) manipulation of independent variables is possible, but it may be unethical, impractical, or costly.<sup>69</sup> While this type of study cannot prove cause-and-effect relationships, they can lead to hypotheses that can be tested with other, more rigorous designs.

In the next section we briefly discuss different sampling strategies.

## 4.2 Sampling/Participant Selection

The size and type of sample to be used in a project evaluation will ultimately depend on the evaluation questions and the type of analysis that the study team would like to conduct. Will the evaluation focus primarily on providing formative feedback? Who might be best able to provide this information? Will the evaluation be comparing groups (e.g., comparing children with disabilities who received a particular intervention to children with disabilities who received no intervention)? If so, how many groups? The more groups that need to be compared, the larger and more complex the required sample. Likewise, does the study team plan to conduct descriptive statistical analysis, inferential statistical analysis, or qualitative analysis of interview/case study data? The answers to these questions will help evaluators to decide which type of sample should be selected for the study. Of course, it’s likely that there will be different samples for the different outcomes that will be measured in the evaluation, since different types of participants (e.g., program graduates or their students) will be involved in the intervention and data collection activities.

In experimental and quasi-experimental studies, the first step in selecting a sample often involves conducting a power analysis, which can be used to help identify the minimum sample needed in order to achieve a certain level of precision for statistical estimates. This is discussed briefly in the next section, followed by a discussion of two basic types of sampling—random sampling and purposeful sampling. As we mentioned before, sampling is an area where a team member with specific training, expertise, or experience is needed. For additional resources on sampling methodology, see the recommended readings in [Appendix D](#).

### 4.2.1 Power Analysis

The term “statistical power” refers to the probability that a study will detect the effects of an intervention or treatment when there is indeed an effect (thereby helping the study to avoid Type II error). Statistical power analysis (a) is a way to determine the probability that a proposed research design will detect the anticipated effects of a treatment and (b)

---

<sup>69</sup> Dimitrov, 2010, p. 43.

helps researchers determine whether to modify a proposed research design in order to achieve adequate power for detecting effects.<sup>70</sup>

For simple research designs with simple random samples (e.g., in which individuals are randomly selected and assigned to participate in a particular intervention), statistical power depends on three things:

- The desired significance level of the statistical test (e.g.,  $\alpha = .05$ )
- The expected size of the intervention effect (the effect size; e.g.,  $ES = .3$ ), and
- The sample size.<sup>71</sup>

In research designs that are considered *multilevel*—that is, for example, when students are clustered within classrooms or within schools—two other factors influence statistical power:

- The sample size at each level (e.g., for a three-level design, this might be the number of students in a classroom, the number of teachers in a school, and the number of schools in a district); and
- The extent of the clustering effects—that is, the amount of variation among clusters (e.g., classrooms) relative to the total variation in student outcomes for a school.<sup>72</sup>

Taking these factors into account will help the study team to determine the statistical power to detect effects of a particular size or larger—known as the Minimum Detectable Effect Sizes (MDES). Some good software programs have been developed to help conduct power analysis for individual and group-randomized experiments as well as quasi-experimental studies.

#### Power Analysis Software Tools

- Optimal Design—for randomized experiments (<http://hlmssoft.net/od/>)
- CRT Power – for simple and cluster-randomized experiments (<http://crt-power.com/>)
- PowerUp! – for randomized experiments and QEDs (<http://web.missouri.edu/~dongn/PowerUp.htm>)

It's beyond the scope of this Toolkit to go into detail about statistical power analysis. We recommend evaluators without training in this area consult a sampling statistician. Hedges and Rhoads (2010) and Dong and Maynard (2013) and [Causal Evaluation](#) are good resources for evaluators interested in conducting power analyses for their experimental or quasi-experimental studies.<sup>73</sup>

We now turn our discussion to two forms of sampling: random and purposeful sampling.

---

<sup>70</sup> Hedges & Rhoads, 2010.

<sup>71</sup> Ibid.

<sup>72</sup> Hedges & Rhoads, 2010.

<sup>73</sup> For additional information on power analysis see Raudenbush, Martinez & Spybrook, 2007; Raudenbush, et al., 2011; Raudenbush & Liu, 2000; Raudenbush, 1997; Schochet, 2008; and Spybrook, Raudenbush, Congdon & Martinez, 2009.

### 4.2.2 Random Sampling

Random sampling can be applied at two different levels in a study: **selection of units** for the study and **assignment of units** to treatment conditions. When sampling is random at both of these levels, a study is considered “fully randomized.” Random sampling aims for representativeness and is particularly good for minimizing threats to internal validity in a study. There are various types of random sampling, including simple random sampling, stratified random sampling, systematic random sampling, and cluster random sampling.

In a **simple random sample** all units in the population of interest have the same probability of being selected to participate in the study or the same probability of being assigned to one treatment condition or another. Simple random samples are relatively easy to select (assuming the population is known) and they permit generalization of results back to the population. However, when using simple random sampling to assign units to treatment groups, it’s possible that the groups may be quite different from each other—a phenomenon known as “unhappy randomization.” To reduce the risk of this, **stratified random sampling** involves dividing the sample into groups on selected variables (e.g., gender or ethnicity) and then selecting a simple random sample from within each group, thus not relying on a single sampling process. Stratified random sampling also allows the evaluator to be sure that specific groups are adequately represented in the study in order to conduct subgroup analyses.

**Cluster (area) random sampling** is similar to stratified random sampling, but instead of sampling subgroups within a population, the researcher divides the population into clusters (e.g., schools or school districts), randomly selects among the clusters, and then measures all units (e.g. classrooms or individual students) within the selected clusters.

**Systematic random sampling** involves numbering the units in the population from 1 to  $N$ , deciding on the sample size that is wanted or needed, determining an interval size ( $k = N/n$ ), randomly selecting an integer between 1 to  $k$ , and then selecting every  $k^{\text{th}}$  unit for the sample. For this to work correctly the population must be listed in random order with respect to the specific characteristics being measured. Systematic random sampling is generally easy to do, and, under certain circumstances, can be more precise than simple random sampling. It may be a viable option for researchers who do not have the time or resources needed to conduct another type of random sampling. It’s also possible to [combine the various sampling strategies within the project evaluation](#). We anticipate that it might be difficult for evaluators to use the different random sampling techniques outlined above in their studies. Instead, the evaluations will likely feature some variant of purposeful sampling (also known as purposive sampling), discussed next.

### 4.2.3 Purposeful Sampling

In purposeful sampling the units in a study aren’t randomly selected or assigned to treatment conditions. This doesn’t mean that the purposeful sample might not be a good representation of a population, or that it cannot serve the purposes of the specific study, but with such a sample researchers cannot know the statistical probability of the extent to which the sample is considered “representative”—thereby limiting generalization of the study results to other populations. In many situations, however, this isn’t required. Indeed, as Patton pointed out, “*the logic and power of purposeful sampling lie in selecting information-rich cases for study in depth. Information-rich cases are those from which one can learn a great deal about issues of central importance to the purpose of the inquiry.*”<sup>74</sup>

Purposeful sampling selects information-rich cases strategically and purposefully; the specific type and number of cases selected depends on the study purpose and resources. The different types of purposeful sampling are outlined below.<sup>75</sup>

---

<sup>74</sup> 2002, p. 230.

<sup>75</sup> Patton, 2002.

- **Stratified purposeful sampling:** Illustrate characteristics of particular subgroups of interest; facilitate comparisons.
- **Purposeful random sampling (still small sample size):** Add credibility to the study when the potential purposeful sample is larger than the evaluation can handle. Reduces bias within a purposeful category. (Not for generalizations or representativeness.)
- **Extreme or deviant case sampling:** Learning from unusual manifestations of the phenomenon of interest, for example, outstanding successes/notable failures; top of the class/dropouts; exotic events; crises.
- **Intensity sampling:** Information-rich cases that manifest the phenomenon intensely, but not extremely; for example, good students/poor students; above average/below average
- **Maximum variation sampling:** Document unique or diverse variations that have emerged in adapting to different conditions. Identify important common patterns that cut across variations (cut through the noise of variation).
- **Homogeneous sampling:** Focus; reduce variation; simplify analysis; facilitate group interviewing.
- **Typical case sampling:** Illustrate or highlight what is typical, normal, average.
- **Critical case sampling:** Permits logical generalization and maximum application of information to other cases because if it's true of this one case, it's likely to be true of all other cases.
- **Criterion sampling:** Picking all cases that meet some criterion; for example, all individuals who are deaf and who enroll in online learning modules through a TA Center.
- **Confirming and disconfirming case sampling:** Elaborating and deepening initial analysis; seeking exceptions; testing variation.
- **Combination or mixed purposeful sampling:** Triangulation; flexibility; meet multiple interests and needs.
- **Convenience sampling:** Do what's easy to save time, money, and effort. Poorest rationale; lowest credibility. Yields information-poor cases.

As with other design elements, the type of sample chosen for a project evaluation will depend on many factors.

In the next section we turn our attention to methodological considerations related to data collection.

## 4.3 Data Collection Methods

As mentioned previously, the choice of data collection methods for a project evaluation will depend upon the evaluation design and the resources available. In this section we discuss the selection and use of four typical data collection methods—surveys, observations, individual interviews, and focus groups—and present another less-common, but potentially useful method: goal attainment scaling.

### 4.3.1 Surveys

There are two basic types of survey research: cross-sectional and longitudinal. **Cross-sectional surveys** collect data from a group of respondents at one point in time and are generally used for the following purposes:

- To examine current attitudes, belief, opinions, or practices;
- To compare two or more groups;
- To measure community needs (e.g., for related services provision);
- To evaluate a program; or
- To conduct a large-scale assessment of selected individuals or programs, such as a statewide or national survey.<sup>76</sup>

Cross-sectional surveys may be administered multiple times—for example program graduates 1, 2, and 3 years past graduation could be surveyed each year—but the group of respondents will differ each time (even if some of the same individuals respond to all of the surveys).

**Longitudinal surveys** follow selected individuals over time and are used to:

- Study **trends** in a population over time (e.g., attitudes among policy makers regarding public financing for related services for students with disabilities);
- Follow development or change in a **cohort** or subgroup of individuals who have been identified based on a specific characteristic (such as students who are deaf or hard of hearing) (Note: In a cohort design, different individuals may respond to each round of surveys, but all individuals must meet the cohort selection criteria to participate in the survey.); or
- Track development or changes in a specific group of individuals, or **panel**, over time. (Note: A panel survey follows the same individuals over time, thereby allowing the researcher to study actual changes in individuals. However, it may be hard to track each individual as time progresses, making panel surveys more costly and labor-intensive to conduct.)

The evaluator must decide, given time and resource constraints, if it's better to use an existing survey or develop a new one. Even if the existing survey isn't exactly tailored to the specific study questions, it's generally easier to adapt an existing survey than to start developing a new one from scratch (Note: When adapting an existing survey the psychometric properties will change.). When done correctly, survey development is time-consuming and can be costly. Of course, the length and complexity of the survey will be determined by the study questions and by the type of analysis the study team would like to conduct.

---

<sup>76</sup> Creswell, 2002.

#### 4.3.1.1 Using an Existing Survey

Conducting surveys is one effective way to gather information about project outcomes. In some cases, an existing survey may be appropriate for the needs of a project. Several resources are available for projects to use in their evaluations .

For Personnel Development Programs, helpful sources include a [survey of higher education programs preparing people to enter the fields of Early Intervention/Early Childhood Special Education](#) prepared by the Center to Improve Personnel Preparation Policy and Practice in Early Intervention and Preschool Education.<sup>77</sup> Many of the items included in the survey can be used by evaluators of PDP projects to gather information about the practices of the different personnel preparation programs.

For Parent Information Centers, OSEP has created a [Parent Center Survey Item Bank](#) that can provide questions suitable for use in evaluation. The item bank includes questions organized by topic, respondent, and mode of data collection.

Projects that wish to use measures of general teacher attributes can find existing surveys in current state evaluations, standard efficacy scales, or professional standards. For example, school systems may currently be using administrator and teacher self-assessments to gather information about educators' perceptions of their professional strengths and areas of development, and these might be publicly available (e.g., the [Rhode Island Model](#)). Existing surveys to measure teacher self-efficacy may also be useful for project evaluations including: A modified Gibson and Dembo Teacher Efficacy Scale<sup>78</sup> for use in a special education resource-room context<sup>79</sup>; the Teachers' Sense of Efficacy Scale;<sup>80</sup> and the Teacher Self-Efficacy Scale.<sup>81</sup> It's also possible to use professional standards (e.g., the National Association of School Psychologists Professional Standards) as the foundation to develop a survey to assess educator practice . For example, Cochran et al. (2012) conducted a field validation of the Council for Exceptional Children's Division for Early Childhood early childhood special education /early intervention personnel standards.

To measure the perceptions of students regarding the classroom instructional environment, projects may wish to utilize the [Tripod student survey instrument](#), with items such as, *"assesses the extent to which students experience the classroom environment as engaging, demanding, and supportive of their intellectual growth."*<sup>82</sup> The Tripod project also has teacher surveys that measure teachers' perceptions of pedagogy, teacher-student relations, and working conditions.

#### 4.3.1.2 Developing a New Survey (or Adapting an Existing One)

In this section we highlight some important considerations related to developing/adapting a survey, rather than presenting an extensive discussion of the methodology. There are many good books about the process of developing surveys<sup>83</sup> and there are helpful webinars related to [customer surveys](#) available on the OSEP website. We urge project staff and evaluators to review these for more detailed information about survey development. For simplicity, we talk here about developing a new survey, although the principles are generally the same for adapting an existing survey.

---

<sup>77</sup> Bruder & Stayton, 2004

<sup>78</sup> 1984, cited in Coladarci & Breton, 1997.

<sup>79</sup> Coladarci & Breton, 1997.

<sup>80</sup> Moran & Hoy, 2001.

<sup>81</sup> Bandura, 2006.

<sup>82</sup> Jerald, 2012, p. 7.

<sup>83</sup> E.g., Czaja & Blair, 2005; Dillman, Smyth, & Christian, 2009; Groves, Flower, Couper, et al., 2004, and Harkness, Braun, Edwards, et al., 2010.

CIPP prepared a three-part webinar series on customer survey development. Check out the series on the [OSEP IDEAs That Work website!](#)

A key step in developing a survey is developing a **framework for the survey**. This includes, for each evaluation question that will be addressed by the survey, determining the item topics (i.e., what will be asked), the target population of interest (i.e., who will be asked), and the mode of administration (e.g., how they will be asked). Table 12 on the next page presents an example.

**Item Topics.** The goals of the study, the study questions, and the resources available will help to determine the topics that should be included in the survey. While it may be tempting to conduct an exhaustive survey covering all possible areas of interest, resource limitations and the time participants are willing to spend on the survey may force the evaluation team to make decisions regarding which outcomes are the most important to measure. Further, if additional data collection activities (e.g., [observations](#) or [interviews](#)) are planned, it may not be necessary to develop a lengthy survey.

**Target Population.** For each evaluation question the study team must decide who will be the best source of information, or target population and determine how the sample will be selected (see [Section 4.2](#) for a discussion of sampling).

**Mode of Administration.** There are two basic modes of survey administration: self-administered questionnaires and structured interview questionnaires. Self-administered questionnaires are usually provided in paper-and-pencil or online/electronic format, while structured interview questionnaires generally are administered in-person or over the telephone (see the sub-section on Interviews for more information). When developing the survey, the study team should determine which language will be the most appropriate for the target population (e.g., Vietnamese for parents of children living in a predominantly Vietnamese-immigrant community) and consider offering the survey in multiple languages. Further, before beginning to develop each item, the study team should decide whether the mode of administration will affect the respondents' decision to answer a question or influence the answer itself. For instance, students may be more likely to answer questions about their perceptions of their teacher in a self-administered survey than in a structured interview questionnaire administered in a group setting.

**Table 12. Sample Survey Development Framework for Evaluation of Teacher Training to Improve Literacy**

Evaluation Question	Item Topics	Target Population	Mode of Administration	Type of Analysis
To what extent do trained teachers exhibit skills and knowledge necessary to implement improved literacy teaching practices?	Principal Survey <ul style="list-style-type: none"> <li>Principal/supervisor reports of teachers’ knowledge and skills</li> <li>Principal/supervisor reports of teachers’ performance relative to professional standards</li> <li>Principal/supervisor reports of quality of service provided to students/children</li> </ul>	Principals/supervisors	Self-administered survey (paper-and-pencil or electronic)	Descriptive statistics
	Teacher Survey <ul style="list-style-type: none"> <li>Teachers’ perceptions of their knowledge and skills</li> <li>Teachers’ perceptions of self-efficacy</li> <li>Teachers’ reports of their use of literacy teaching techniques learned in training</li> <li>Teachers’ perceptions of quality of service provided to students/children</li> </ul>	Trained Teachers	Self-administered survey (paper-and-pencil or electronic)	Descriptive statistics
	Student/Client (of trained teacher) survey <ul style="list-style-type: none"> <li>Students’/children’s perceptions of quality of service provided to them</li> </ul>	Students/children	Self-administered survey (paper-and-pencil or electronic)	Descriptive statistics

Czaja & Blair highlighted three fundamental characteristics of a good survey questionnaire: it’s a valid measure of the factors of interest, it convinces respondents to cooperate, and it elicits acceptably accurate information.<sup>84</sup> Figure 6, on the next page, presents a visual representation of the process that respondents must go through to answer a given survey question. If the respondent cannot understand the question, cannot remember the answer (or never knew it), doesn’t want to give the answer, or cannot figure out *how* to give the answer, then there is a problem with the item.<sup>85</sup>

<sup>84</sup> Czaja & Blair, 2005, p. 65.

<sup>85</sup> Fiore, Helba, Berkowitz, et al., 2012.

Figure 6. The Process of Answering a Survey Question

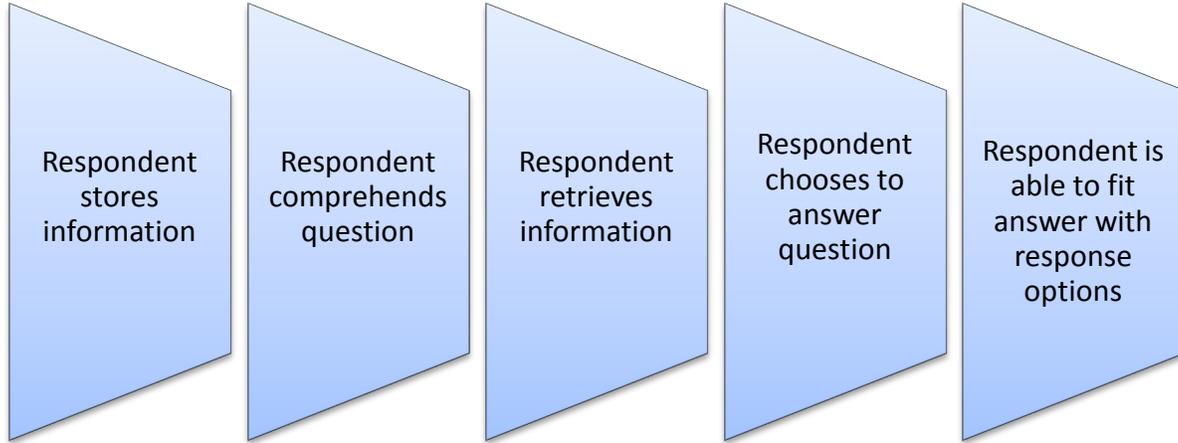


Table 13 presents some characteristics of a good survey item with examples of questions exhibit or do not exhibit these characteristics.

Table 13. Characteristics of a Good Survey Item, with Examples of Not Good and Good Items

Characteristic	Examples	
	Not Good	Good
<b>Only asks one question at a time (avoids double-barreled questions)</b>	To what extent do you think closed captioning and accessible textbooks have increased your access to current information?	To what extent do you think closed captioning has increased your access to current information?
<b>Contains a clear threshold for answering “yes”</b>	Has your child seen a speech and language pathologist in the past month?	Has your child met with a speech and language pathologist to practice fluency in the past month?
<b>Provides a timeframe</b>	How often do you access information through the Parent Information Center?	How many times in the past month have you accessed information through the Parent Information Center?
<b>Provides a timeframe appropriate to the topic</b>	How many times did were you absent from class during your personnel preparation program?	How many times were you absent from class last semester?
<b>Uses clear terminology and plain language</b>	Are children with, or at risk for, developmental delays more likely to experience latency to learn a contingency?	Are children with, or at risk for, developmental delays more likely to experience delays in learning the relationship between a behavior and its consequences?
<b>Gives exhaustive and mutually exclusive response options</b>	In which of the following situations does your supervisor most frequently observe your work? <ul style="list-style-type: none"> <li>• When you’re working in a general education classroom</li> <li>• When you’re working in a self-contained special education classroom</li> <li>• When you’re working with your students</li> </ul>	In which of the following situations does your supervisor most frequently observe your work? <ul style="list-style-type: none"> <li>• When you’re co-teaching in a general education classroom</li> <li>• When you’re teaching in a self-contained special education classroom</li> <li>• When you’re providing pull-out services to children with special needs</li> </ul>

Czaja & Blair pointed out that that most refusals to complete a survey occur at the beginning. That is, the information provided in a cover letter or in the first parts of the survey helps the target respondent decide whether or not to complete the survey. As such, they recommend that the first question have all of the following characteristics:

- Relevant to the central topic
- Easy to answer
- Interesting
- Applicable to and answerable by most respondents
- Closed format.<sup>86</sup>

Another consideration when developing surveys is that the ordering of the items within a survey can affect responses. Whenever possible, it's better to place sensitive questions later in a survey so that respondents do not decide to terminate the survey early. It's also good to place items requesting demographic information (e.g., gender, occupation, ethnicity, disability status) at the end of the survey so as not to discourage any respondents from answering the survey questions.

In some cases the answer to one question can influence the answer to the following item. For example, if a respondent is asked to rate the different components of the project (e.g., newsletter, the workshops, the presentations, the staff), his or her general assessment of the project in a subsequent question may be different than if the item order had been reversed.

The available response options also can affect a respondent's ability or willingness to answer. For instance, if a respondent has a "neutral" opinion about her experience with a training and is asked to choose between the response options "highly dissatisfied," "moderately dissatisfied," "moderately satisfied," or "highly satisfied," she may choose to skip the survey item.

Moreover, researchers have identified a number of problems with response options in the agree-disagree (A-D) format (i.e., strongly agree, agree, neither agree nor disagree, disagree, strongly disagree). Specifically, some of the **problems associated with agree-disagree items** include:

- A-D items are more cognitively difficult – Respondents have to understand the question as written and then translate their answer to the question into an A-D response format.
- A-D items are subject to acquiescence response bias – Respondents tend to agree with A-D items regardless of the content, especially when the respondent ability and motivation are low and the task difficulty is high. Culture also influences this bias (i.e., Latinos and respondents from collectivistic cultures tend to demonstrate this bias more frequently).
- It's often difficult to interpret the meaning of a "disagree" response – It can be difficult to know whether the respondents disagree with the wording of the item itself or with the subject of the item.
- A-D items often force respondents to think through double negatives to be able to respond – Respondents may have to consider whether they disagree with a negatively-worded statements, such as "Public officials don't care about people like me much."<sup>87</sup>

---

<sup>86</sup> Czaja & Blair, 2005, p. 94.

<sup>87</sup> Holbrook, 2013.

Whenever possible, we recommend using response options that are item-specific. For example, consider the following items; the first offers the A-D response options and the second offers item-specific options:

Option 1: “Sometimes I feel that the information/service I receive isn’t helpful.” Do you AGREE STRONGLY, AGREE SOMEWHAT, NEITHER AGREE NOR DISAGREE, DISAGREE SOMEWHAT, or DISAGREE STRONGLY with this statement?

Option 2: “How often is the information/service you receive helpful?” ALWAYS, MOST OF THE TIME, ABOUT HALF OF THE TIME, SOME OF THE TIME, OR NEVER?

Although it may take more time to develop item-specific response options for a survey, the quality of the data obtained by these items is better than with A-D items. Additionally, studies have shown that item-specific response options do not take significantly more time to complete than A-D format items. There are times when using item-specific response options may not be an option, however, such as: (a) when comparing results of the survey to previous studies that used A-D items; or (b) when conducting longitudinal studies that have used A-D items in past data collections. In these situations, it may be better to continue the use of the agree-disagree format, keeping in mind the inherent limitations of these items.<sup>88</sup> [Optima 360](#) is one online resource that includes a variety of response options that can be used for developing a survey.

#### 4.3.1.3 Reviewing and Testing the Survey

Whenever possible, a content expert should review the survey to see if the content is accurate and if the items adequately address the item topics laid out in the data analysis plan. Further, a methodological expert should review the survey to identify any issues with item construction or survey organization. If neither of these are viable options, the study team should carefully review the survey to look for any of the issues outlined above. Once again, there are many good books on survey development that can provide guidance on the survey review.<sup>89</sup>

Once the review has been completed, the evaluation team should conduct a pretest (or pilot test) of the survey. As Czaja and Blair pointed out,

*In designing a questionnaire, we make many decisions and assumptions, some conscious, others not. Underlying the questionnaire draft is our judgment about what kinds of things respondents will know, what words they will understand, what sorts of information they can and will provide, and what response tasks they can perform. When we pose alternative choices to the respondent, we have in mind some notion of the appropriate dimensions for an answer; in fact, we start to envision what our data will look like. Much of our effort is subtly informed by our feelings about how people will respond to our queries, by a belief that what we ask is sensible, by some vision of how the respondents’ world works.<sup>90</sup>*

Even after a thorough review is conducted, a pretest is the best way to know how the target population may respond to the survey. During the survey review issues may arise that can be specifically addressed in a pretest, such as whether respondents understand the terminology included in a particular item. A pilot may involve one-on-one administration of the survey in an interview situation (e.g., through a structured interview or cognitive interview<sup>91</sup>) or field testing the

<sup>88</sup> Holbrook, 2013; for more information see Saris, Revilla, Krosnick & Shaeffer, 2010.

<sup>89</sup> E.g., Czaja & Blair, 2005; Dillman, Smyth, & Christian, 2009; Groves, Flower, Couper, et al., 2004, and Harkness, Braun, Edwards, et al., 2010.

<sup>90</sup> Czaja & Blair, 2005, p. 105.

<sup>91</sup> A cognitive interview is a one-on-one interview designed to determine the process respondents use in answering a question, or to identify problems respondents have in understanding or answering a question. Such interviews use a variety of laboratory techniques such as think aloud procedures or the paraphrasing of questions (Czaja & Blair, 2005).

survey (i.e., administering the survey using the exact procedures planned for the overall study) with a small sample similar to the target population and then conducting follow up interviews, telephone calls, or even internet or email exchanges. A pretest protocol may be developed if more than one person is conducting the follow up with respondents.

The following are examples of pretest questions adapted from Czaja and Blair:<sup>92</sup>

- Were there any questions you were not sure how to answer? If yes, which ones were those? Why were you not sure how to answer the question?
- When I used the term *latency to learn a contingency*, what did you think I meant by that?
- When I asked the question about *closed captioning increasing your access to current information*, what sorts of things did you consider?
- Are there any questions you think that many people would find difficult to answer? If yes, which ones were those? Why do you think people would have difficulty with the question?
- Were there any important things related to these issues that we failed to cover?

Czaja and Blair also offered four final suggestions for survey development:

- If possible, use multiple pretesting methods and multiple rounds of testing;
- Learn to assess survey questions by reading them aloud and listening carefully for any awkward or unnatural phrasing;
- Accept that sometimes a question must be reworded; and
- Look for examples of questionnaires written by experienced researchers.

Finally, remember: **When making changes to an existing survey, it isn't appropriate to cite the psychometric properties of the original survey.** If this information is desired, the study team will need to calculate new reliability statistics and evaluate the validity of the revised survey based on data collected from the study sample.<sup>93</sup>

#### 4.3.1.4 Potential Sources of Error in Surveys

Every survey contains some sources of error that reduce the accuracy of the data collected. This error can be large or small—and some error simply cannot be eliminated no matter what a researcher does, as in the case of a sample survey (which inherently includes error because it doesn't gather data from all members of the population). Nevertheless, some sources of error are more harmful than others, and the study team should make efforts to minimize them.<sup>94</sup> These are briefly discussed below.

---

<sup>92</sup> Czaja & Blair, 2005, p. 109.

<sup>93</sup> See, for example, Dimitrov, 2012.

<sup>94</sup> Czaja & Blair, 2005.

#### 4.3.1.4.1 Low Response Rates

The response rate for a survey is the percentage of eligible respondents for whom survey data are obtained. There are multiple ways to calculate response rates<sup>95</sup> but we will not focus on those here; instead we focus on two types of response rates that can impact the results of the survey—the overall response rate (known as the unit response rate) and the item response rate. The unit response rate is the percentage of potential respondents who actually complete the survey, while the item response rate is the percentage of respondents who answer a particular item. A low item response rate tends to be less of a concern than a low unit response rate, but *“If [you] fail to obtain any information from some respondents, and for others fail to obtain complete information, [your] estimate and other analyses may be distorted, sometimes quite seriously.”*<sup>96</sup>

When response rates are low, the reliability and validity of the study conclusions are called into question. For example, if the study team sends a survey to all of the parents of children who received a particular intervention ( $n = 150$ ) and only 25% of the parents respond, to what extent can the study team be confident that the results of the survey are an accurate reflection of parents’ perceptions of the intervention? If 90% of respondents do not answer a series of questions related to the interactions they had with project staff, how should the study team interpret those results? Are the omissions a result of poorly worded questions or did the respondents misunderstand the skip instruction in a previous question? Are the omissions indicative of unwillingness to answer questions respondents perceive to be too sensitive? If a survey pretest was conducted prior to administering the survey, any difficulties with item wording or skip patterns should have emerged. However, even when a pretest was conducted, it can be very difficult for a researcher to know how to interpret missing data from items not answered. Consequently, if either of these rates is low, it’s important to determine whether there are systematic differences between respondents and non-respondents, since such differences might indicate bias in the data.

#### 4.3.1.4.2 Differences among Respondents and Non-Respondents

When examining differences among respondents and non-respondents, consider two primary questions:

- Are some demographic groups under- or over-represented among survey respondents?
- Are some of the survey items answered differently by members of different groups?

As discussed in [Section 4.2](#), the study’s sampling frame tells what demographic groups should be included in the study sample. Using demographic data from the respondents (collected through the survey or through administrative records), evaluators can analyze the response patterns to determine whether some groups responded to the survey or answered specific questions more than others. For example, were program graduates who had higher grade point averages (GPAs) at the end of their program more likely to respond to the survey than graduates who had lower GPAs at the end of their program? Or, were program graduates who had a specific field supervisor more likely to respond to questions about the quality of their field placement experience than graduates who had a different field supervisor?

Of course, the existence of differential response rates among respondents and non-respondents doesn’t necessarily mean that the data obtained through the survey are biased. However, further investigation is warranted since such

---

<sup>95</sup> For more information, see Czaja & Blair, 2005.

<sup>96</sup> *Ibid.*, p. 197.

differences might affect the validity of the study's conclusions. *"To the extent that opinions and behaviors differ by subgroups, their overrepresentation or underrepresentation will affect results."*<sup>97</sup>

If different patterns of responses among different groups of respondents are found, we recommend conducting a **non-response analysis** to check for response bias. One way to do this is to contact a few non-respondents (best if selected randomly) by telephone or by email to see if their answers differ substantially from respondents. There is no specific guideline for the number of such follow-ups that might be made, but if evaluators can document that the responses obtained from the non-respondents do not differ very much from the responses originally obtained, the evaluation team (and the study audience) will have greater confidence in the conclusions of the survey. Regardless of whether a non-response analysis is conducted, evaluators should clearly document in the evaluation report the [limitations of the study \(Section 3.5.3\)](#) conclusions related to a possible response bias in the survey.

#### 4.3.1.5 Reducing Error in Surveys

As with every other aspect of the study, decisions about whether and how to reduce survey error must balance costs and available resources.<sup>98</sup> We have already talked about design elements that can improve the quality of a survey and we will not revisit them here.

The best way to avoid error related to non-response is to conduct extensive follow-up with the survey sample to increase the number of respondents. This can be done through telephone calls, emails, and even letters (although telephone calls are generally considered to be the most effective follow-ups). Keep in mind that every contact the study team has with the potential respondents—whether through email exchanges, voicemail recordings, or messages left with another person in the household—can affect whether or not the respondents decide to complete the survey. While developing the data collection plan for the evaluation, evaluators should consider the ways that respondents' willingness to participate in the study might be maximized. Taking time to *"[think] through the data collection process in simply a commonsense manner—but from the respondents' perspective—[will] produce many ideas, concerns and insights about how best to conduct data collection."*<sup>99</sup>

We now turn to our discussion of another common data collection method, observations.

#### 4.3.2 Observations

It's important to note that conducting observations that will produce high-quality data requires disciplined training and rigorous preparation.<sup>100</sup> Observational methods range from the highly qualitative field observations (e.g., field notes and journals) commonly utilized by anthropologists and ethnographers to the more quantitative observations conducted using rubrics, frameworks, or other types of observation instruments or checklists. For the purposes of this Toolkit, we will focus on the latter. Here we focus on conducting classroom observations, but the same principles apply to other types of observations.

---

<sup>97</sup> Czaja & Blair, 2005, p. 198.

<sup>98</sup> Czaja & Blair, 2005.

<sup>99</sup> Ibid., p. 199.

<sup>100</sup> Patton, 2002.

**4.3.2.1 Considerations When Choosing/Developing Observation Protocols**

The type of observation instrument, or protocol, used in the study will have a significant impact on the quality and consistency of data that are collected during the observation. Every observation protocol will have certain inherent limitations, so during the planning phase the evaluation team must consider which kinds of data need to be collected during the observations. Consider the following examples in Figure 7 and Figure 8 below.

**Figure 7. Observation Protocol Example 1**

Teacher’s name/initials/pseudonym: \_\_\_\_\_ Grade: \_\_\_\_\_

Observer: \_\_\_\_\_ Setting: \_\_\_\_\_

Summary of Teacher’s Observed Behavior	Implications

**Figure 8. Observation Protocol Example 2**

Teacher’s name/initials/pseudonym: \_\_\_\_\_ Grade: \_\_\_\_\_

Observer: \_\_\_\_\_ Setting: \_\_\_\_\_

Observed Behavior	Time Interval (Place an “x” in each box if the behavior was observed during the interval)			
	10:00 - 10:05	10:06 - 10:10	10:11 - 10:15	10:16 - 10:20
Teacher follows the instructional model closely.				
Teacher differentiates instruction.				
Teacher models skills correctly.				
Teacher follows the steps of the correction procedures to provide immediate feedback.				

As can be seen, the types of data collected by the two instruments presented above will vary considerably. In the case of Figure 7, there is a strong possibility that the data collected by one observer will differ greatly from the data collected by another observer. In the case of Figure 8, there are many behaviors not included in the protocol that might be observed during a class period, but since they aren't listed they will not be recorded. The example in Figure 7 provides data that are information-rich, yet possibly inconsistent—and requiring extensive qualitative analysis, while the example in Figure 8 provides data that are consistent, able to be analyzed quantitatively, yet possibly incomplete. Additionally, the type of training needed to be able to assess “implications” of a given behavior (Figure 7) is very different from the training needed to be able to accurately record specific instances of selected student behaviors (Figure 8). (Note: In the case of Figure 8, if the study wanted to gather additional data about how often certain behaviors occurred, rather than recording an “x” in the time interval, the observers could use a scale such as 0 = never, 1 = sometimes, 2 = frequently to give some measure of the frequency of behaviors during the class period.) When selecting or developing an instrument to use in the observations, these considerations should be kept in mind.

In many instances it may be preferable to use an existing observation protocol that has been tested and that has guidelines for training observers or that provides training, such as the Early Childhood Environment Rating Scale (ECERS)<sup>101</sup>, Classroom Assessment Scoring System (CLASS)<sup>102</sup> and the Classroom Climate Scale.<sup>103</sup>

The majority of states and districts use teacher observations as their primary measure of teacher effectiveness.<sup>104</sup> However, a 2009 study conducted in 12 districts across four states (which included surveys of 15,000 teachers and 1,300 administrators) found that such observations are **generally infrequent** (generally two or fewer observations per year) and **too brief** (60 minutes or less) to accurately measure teacher performance. Additionally, the observations were frequently conducted by untrained administrators and the results seldom were used to provide feedback to mediocre or poor performing teachers.<sup>105</sup> Consequently, evaluators should use caution when considering using such data in their studies.

The Bill & Melinda Gates Foundation is one of a number of organizations advocating improvements in the quality of classroom observations.<sup>106</sup> In addition, many states are taking steps to develop and implement classroom observation protocols that are more consistent and better measures of the instructional practices taking place within a classroom. Some of these protocols are available online (see, for example, the [Indiana Department of Education's RISE Evaluation and Development System](#)). Unfortunately, a study conducted by the National Comprehensive Center for Teacher Quality found that the majority of states and districts do not use observation protocols designed to measure the specific functions and distinct roles of special educators. Rather, most simply adapt an existing protocol to reflect specialized roles and responsibilities.<sup>107</sup> The Alabama Department of Education (DOE), for example, has modified its teacher observation instrument to reflect instruction toward alternative standards for teachers of students with low-incidence disabilities.

Whichever protocol the study team selects for the observations, it will be necessary to train the observers to ensure that the data they collect is consistent and reliable.

---

<sup>101</sup> Harms, Clifford, & Cryer, 2015.

<sup>102</sup> Pianta, L Paro, & Hamre, 2008.

<sup>103</sup> McIntosh, Vaughn, Schumm et al., 1993.

<sup>104</sup> Holdheide, Goe, Croft, & Reschly, 2010.

<sup>105</sup> Weisberg, Sexton, Mulhern, & Keeling, 2009.

<sup>106</sup> Jerald, 2012.

<sup>107</sup> Holdheide et al., 2010.

#### 4.3.2.2 Training the Observers

It's well known that human perception is highly selective. *"What people 'see' is highly dependent on their interests, biases, and backgrounds. Our culture shapes what we see, our early childhood socialization forms how we look at the world, and our value systems tell us how to interpret what passes before our eyes."*<sup>108</sup> Simply being able to see and hear what is happening in a classroom or other setting doesn't mean that a person is able to conduct a high-quality observation. Different people frequently highlight and report different aspects of a given situation, so some training is needed to ensure that individuals conducting observations report the same incident with accuracy, authenticity, and reliability. As Patton<sup>109</sup> pointed out, training to become a skilled observer includes:

- Learning to pay close attention to what is happening in a given context (e.g., recording visual, auditory and behavioral cues);
- Practicing writing descriptively (e.g., documenting the setting and interactions among people being observed);
- Acquiring discipline in recording field notes or observational data (e.g., monitoring frequency of student behaviors during 5-minute time intervals);
- Knowing how to separate detail from trivia (e.g., documenting a significant disruption in the flow of instruction but not recording every side conversation taking place among students);
- Using rigorous methods to validate and triangulate observations (e.g., calculating inter-rater reliability among two or more observers); and
- Reporting (or at least acknowledging) the strengths and limitations of one's own perspective as an observer (e.g., recognizing that having a background in literacy might make an observer more attuned to differentiated instruction focused on literacy than that focused on math).

Further, all individuals who will be conducting observations should be trained in the use of the specific observation protocol employed by the study. It's important to provide observers with an instructional manual for each protocol that includes clear descriptions of the key terms and examples of what the behavior being observed "is" as well as what it "is not." Additionally, the observation protocol itself should include clear instructions to the observer as well as definitions of any terms that might be subject to differing interpretations by multiple observers. For example, the phrase "at inappropriate times" in Figure 8 above might mean one thing to one observer and mean something entirely different to another observer.

Prior to sending the observers into the field to conduct the observations, the study team should practice using the protocol, and should calculate the agreement between different raters (also known as inter-rater reliability, see [Section 4.3.2.6](#)). Ideally, this would include opportunities to use the observation protocol in a "real-life" setting, such as might occur during a pilot study. If this isn't possible, video clips can serve as examples for training observers (see, for example, videos available from [The New York Times](#) or [YouTube](#)). At a minimum, the study team should have a conversation about the use of the protocol, during which sample situations are presented to the team and a consensus is reached as to how the protocol would be used to record what is happening.

---

<sup>108</sup> Patton, 2002, p. 260.

<sup>109</sup> Ibid.

#### 4.3.2.3 *Deciding How Many Observations to Conduct and Who Will Conduct Them*

Although there is no established standard for the number of observations to conduct, good practice suggests that educators should be observed multiple times.<sup>110</sup> Nevertheless, the number of observations will depend in large part upon the resources available, as well as the proximity of the study team to the classroom or professional setting. Likewise, the number of individuals conducting each observation will depend upon these same factors. When developing the evaluation plan the study team should calculate the expected cost of the observations in terms of travel (e.g., Is it necessary for the observer to book a hotel and rent a car?) and staff time (e.g., How many hours will it take for the individuals to complete the observations, including travel time?). It's also useful to consider the skills that will be required of the observers and if multiple observers are needed to complete one assessment. In a study of the effects of a commercial reading program on student performance in middle school classrooms, for example, all observations were conducted by a pair of observers.<sup>111</sup> One observer was a researcher who was highly trained in observation methodology as well as in the specific protocol. The other observer was a literacy specialist who was trained in the use of the observation protocol and whose content expertise was valuable in helping to identify whether specific teacher practices were in line with good literacy practice (even if they might not exactly align with the reading program). Together, these two observers generated complementary data that enhanced the study team's understanding of what actually was taking place within the classroom.

#### 4.3.2.4 *Obtaining Access to the Observation Site*

The ability of the study team to conduct observations depends upon the teacher or other related-services professional, the school, and the district being willing to allow the observations to take place. The study team will need to contact the school district well in advance to obtain permission to conduct the observation. This may entail obtaining formal permission from the district's institutional review board (IRB) or human subjects review board (HSRB) or simply receiving permission from the district superintendent or his or her designee. In some cases the study team will need to receive permission from the parents of students. The school principal will also need to be contacted about the observations (see [Section 3.1](#) for information on obtaining IRB and district approval and [Section 3.2.2](#) for information on obtaining participants' consent. [Appendix C](#) includes sample notification letters for districts with and without research approval offices and for schools).

#### 4.3.2.5 *Conducting the Observations*

On the day of the observation the study team should arrive to the setting early to check in with the appropriate people (e.g., the principal or other administrators) and to give time to move to the observation location (e.g., to navigate a busy hallway during the break between classes). During the observation, the study team should try to minimize disruptions to the classroom/setting as much as possible. The study team should also insure that no events (such as fieldtrips, fire drills, holidays, or picture day) will interfere with their observations.

It's possible to conduct video observations, rather than in-person observations, if the study team has limited resources or staff to support travel to the different study locations. Through its Measures of Effective Teaching (MET) project the Bill and Melinda Gates Foundation is supporting the use of video observations as one measure in a comprehensive evaluation of teacher practice.<sup>112</sup> Other researchers have pointed to the benefits of special education teachers using

---

<sup>110</sup> Bill & Melinda Gates Foundation, 2012.

<sup>111</sup> Dimitrov, Jurich, Frye, Lammert, Sayko & Taylor, 2012.

<sup>112</sup> Bill & Melinda Gates Foundation, 2011.

video observations.<sup>113</sup> The costs associated with the different types of video observations vary greatly, however, so even though they do not require individual travel, video observations may not be more economical.

#### 4.3.2.6 Calculating Inter-Observer Agreement

When more than one person conducts the observations it's good practice to calculate inter-observer agreement (also known as inter-rater reliability) to determine the consistency of the ratings across observers. Inter-observer agreement calculations give the researcher information about the degree to which changes in the ratings of the dependent variable are due to actual changes taking place in the classroom or practice setting, rather than differences in the consistency of the ratings.<sup>114</sup>

Keep in mind that inter-observer agreement isn't a measure of the *accuracy* of the observers' recordings; even if both observers are consistent in their ratings, they might not accurately record what took place during the observation. For this reason it's extremely important to train the observers carefully in the use of the observation protocol before sending them out to conduct observations. Kennedy (2005) pointed out three reasons why calculating inter-observer agreement is important:

- It can be used as a training standard for new observers;
- It allows the researcher to determine how consistent the observers are across observation occasions; and
- It helps to avoid observer drift, which occurs when the people conducting the observations do not consistently apply the definitions laid out in the observation protocol.

Some researchers make video recordings of observations at the start of a study to use for training purposes and to help "recalibrate" observers periodically throughout the study. Another method to ensure that all raters stay calibrated throughout the period of data collection is to designate two raters to some percentage of observations (10 to 20%) and calculate the inter-rater reliability for those double-coded observations. Whether project evaluators are able to do this for their studies will depend on the resources available.

There are many ways to calculate inter-observer agreement, including correlational and complex statistical approaches. We focus here on simpler calculations that often are used in single-case research. We also include a method for calculating agreement using ratings of more than two levels (or more than two raters). These are briefly summarized below.<sup>115</sup>

---

<sup>113</sup> See, for example, Baecher & Connor, 2010 and Dymond & Bentz, 2006.

<sup>114</sup> Kennedy, 2005.

<sup>115</sup> For more information about calculating inter-observer agreement (including additional types of calculations) see Kennedy, 2005 and Watkins & Pacheco, 2000.

#### 4.3.2.6.1 Total Agreement Approach

A simple and popular approach has been to calculate total agreement between two observers. To do this, the researcher sums the total number of responses (or times that a behavior or event was observed) recorded by each observer. This typically generates a different total for each observer, a “smaller” and a “larger”. To calculate total agreement, divide the smaller total by the larger total, and multiply the result by 100%. The formula is

$$\frac{S}{L} (100\%)$$

where  $S$  is the smaller total and  $L$  is the larger total.

This approach has three advantages:

- It’s easy to conceptualize and calculate;
- It can be used to calculate inter-observer agreement in instances where observers have not accurately aligned their intervals (i.e., one observer began recording during the first interval, but the second observer did not begin until the second interval); and
- It’s relatively sensitive to overall levels of responding (or occurrence).

A key limitation of the total agreement approach, however, is that it doesn’t tell the researcher whether the observers actually agreed on the occurrence of individual instances of behavior. Consequently, there can be high levels of inter-observer agreement but the observers may have never agreed that a specific behavior ever took place. This presents a challenge to interpretations of the observation data.<sup>116</sup>

#### 4.3.2.6.2 Interval Agreement Approach

The interval agreement approach (also known as combined, point-by-point, or overall agreement), on the other hand, does take into account when behaviors or events actually occur. This type of calculation requires an observation protocol that uses an interval or event system of measurement to record behavior (such as that found in Figure 8). Once the observation is complete, the recording of behavior is compared between the two observers on an interval-by-interval basis; if both observers recorded a behavior as occurring or not occurring in a particular interval, it’s scored as an agreement. If only one observer recorded a behavior during an interval it’s scored as a disagreement. Then the total number of agreements is divided by the total number of agreements plus disagreements, and the sum is multiplied by 100%. The formula is

$$\frac{A}{A + D} (100\%)$$

where  $A$  is the number of agreements and  $D$  is the number of disagreements.

Interval agreement is one of the most commonly-used measures of inter-observer agreement. However, when a behavior being observed occurs very frequently or very rarely, it’s possible for there to be high levels of interval agreement even if the observers do not agree on the occurrence of a specific behavior. For example, if a behavior occurred only one time during an observation and the observers disagree on its occurrence, the interval agreement

---

<sup>116</sup> Kennedy, 2005 and Bryington, Palmer & Watkins, 2002.

would still be very high. For this reason, researchers developed another measure of inter-observer agreement, discussed next.<sup>117</sup>

#### 4.3.2.6.3 Occurrence/Non-Occurrence Agreement Approach

The occurrence/non-occurrence agreement approach is an even more stringent way to measure inter-observer agreement than the interval agreement approach. In this case, during each interval, agreement is calculated for each occurrence and non-occurrence of a behavior or event. The formula for this is the same as for interval agreement, except two separate calculations are conducted: one for occurrence and another for non-occurrence.

$$\frac{AO}{AO + DO}(100\%)$$

where *AO* is the agreement on the occurrence and *DO* is the disagreement on the occurrence.

$$\frac{AN}{AN + DN}(100\%)$$

where *AN* is the agreement on the non-occurrence and *DN* is the disagreement on the non-occurrence.

Each statistic is then reported separately.

#### 4.3.2.6.4 Cohen's Kappa

Another, somewhat more complicated, method of calculating inter-observer agreement is Cohen's Kappa. The kappa method indicates what proportion of agreement is above and beyond what would be expected by chance alone. It has a few advantages over the previous methods, including that it corrects for chance agreement and can be used with ratings on two or more levels (such as a 1-7 rating scale). It also has a modification that can be used for more than two different observers.<sup>118</sup> In addition, it allows for generalizability across different experimental conditions, for instance observations in different classrooms in different years, since it isn't affected by the rates of behavior, as are the prior methods. Kappa ranges from -1.00 which indicates perfect disagreement, to 0.00 which indicates chance agreement, to 1.00 which indicates perfect agreement. A sample calculation for raters who rate 10 teachers on a behavior as low, medium, or high, is shown in Table 14 below.

---

<sup>117</sup> Kennedy, 2005.

<sup>118</sup> See Bryington, Palmer & Watkins, 2002 for a discussion of this modification, called the Fleiss kappa and see <https://www.easycalculation.com/statistics/cohens-kappa-index.php> for an explanation and link to an online calculator.

**Table 14. Sample Calculation of Cohen's Kappa Agreement**

		Observer 1		
		Low	Medium	High
Observer 2	Low	2	1	1
	Medium	1	2	0
	High	0	0	3

As can be seen, the ratings of the first observer agreed with those of the second observer in 7 cases. They both rated 2 teachers as low on this behavior, 2 teachers as medium on this behavior, and 3 teachers as high on this behavior. In the remaining 3 cases, the two raters disagreed on the rating of the teacher’s behavior. Numbers such as these can be plugged into an online calculator<sup>119</sup> and the kappa will be obtained. In this case the kappa would be 0.55, meaning that the observers accounted for 55% agreement over what would be expected by chance. Kappas between 0.40 and 0.59 are considered “fair,” those between 0.60 and 0.74 are “good,” and those above 0.75 are considered “excellent.”

The primary limitations of the kappa are difficulty in calculating, which can be solved with online calculators, and inability to be used in a situation where all observers agree or disagree in all cases, which is unlikely to happen often, and would rather clearly be a case of perfect agreement or disagreement. Since the kappa is rigorous and provides generalizability over different conditions, it’s the preferred method for calculating inter-rater reliability.

#### 4.3.2.7 Analyzing Observational Data

Clearly, the type of observational protocol used will affect the type of analysis to be conducted. Returning to Figure 7 and Figure 8 above, the type of analysis will depend on the data recorded in the protocol. For instance, in Figure 7, if the observer records narrative information about what occurred during the class period, qualitative content analysis will be required. If, however, the observer records tallies of behaviors (e.g., student got up and walked around the classroom five times during the period) it will be possible to develop summary tables of the frequency of behaviors during the class period. In Figure 8, it’s possible to tally the total number of intervals in which a student exhibited a certain behavior. It’s also possible to calculate a proportion of the total class period in which the student was exhibiting the behavior. If the protocol asked the observer to record frequency of behaviors (such as with the aforementioned scale of 0 = never, 1 = sometimes, 2 = frequently), the analysis could provide a clearer picture of the level to which the student was exhibiting a certain behavior (e.g., frequently during 4 of the 5 intervals).

<sup>119</sup> <http://vassarstats.net/kappa.html> and <http://graphpad.com/quickcalcs/kappa1.cfm> are online calculators that provide kappas, standard errors, and confidence intervals.

### 4.3.3 Individual Interviews

Individual interviews are the most costly and time-consuming interviews to conduct, as they require a researcher to speak directly with a respondent—either in-person, on the phone (e.g., direct call, teleconference, random digit dialing [RDD], or computer-assisted telephone interviewing [CATI]), or over the internet (e.g., through email, web chat, or social media). As with observations, interviews can be more qualitative (i.e., unstructured) or more quantitative (i.e., structured) in nature. Whether the interview is unstructured, structured, or semi-structured will affect the type of data analysis that will be required. Unstructured interviews generally will require extensive qualitative data analysis (see [Section 4.4.4](#) and the [OSEP webinar series on qualitative interviews](#) for more information), while structured interviews will allow for more quantitative analytic approaches.

It's good practice to prepare an interview protocol in advance of the interviews and to practice asking the questions to someone familiar with the topic. This will help to clarify any issues that may arise with the question phrasing, highlight "difficult" or "sensitive" terms or topics, and identify any possible areas where misunderstandings may arise. It may be necessary to revise and re-test the interview protocol prior to conducting the interviews. In addition, all data collectors should be skilled in conducting interviews in a way that encourages respondents to answer the questions openly and honestly.<sup>120</sup> Furthermore, to the extent possible (and as appropriate), the data collectors should try to cover all of the questions on the interview protocol with every potential respondent. Inconsistencies in the way that interviewers ask questions across respondents can lead to poor quality data and limit the ability of evaluators to draw conclusions or generate insights from their qualitative interview data.<sup>121</sup>

CIPP prepared a two-part webinar series on planning and conducting qualitative interviews. Check out the series on the [OSEP IDEAs That Work website!](#)

#### 4.3.3.1 Unstructured Interviews

In an unstructured interview the interviewer allows the respondent to decide how she wants to answer a given question. These questions can be pre-determined or they can arise in the moment as probes to obtain additional information about a topic or an issue raised by the respondent; this generally depends upon the needs of the study and the skill of the interviewer.<sup>122</sup> In these types of interviews **it's important to ask questions that encourage the respondent to elaborate**, rather than to provide brief or "yes/no" answers. For example, consider the following questions:

- Question 1: Do you think that the Parent Information Center (PIC) provided you with the knowledge and skills necessary successfully support your child who has a developmental disability?
- Question 2: In what ways do you think that the Parent Information Center (PIC) provided (or didn't provide) you with the knowledge and skills necessary to successfully support your child who has a developmental disability?

As can be seen, Question 1 allows the respondent simply to respond "yes" or "no," while Question 2 urges the respondent to list the ways that the PIC has prepared him to support his child. After all of the interviews are conducted, the responses to Question 1 will provide little information beyond allowing the study team to tally the number of parents who believe the Center helped parents to develop the requisite knowledge and skills (i.e., the respondents who answered "Yes"). The responses to Question 2, in contrast, will allow the study team to gain a better understanding of

<sup>120</sup> See Weiss, 1994, for an excellent discussion of the design and conduct of qualitative interviews.

<sup>121</sup> Weiss, 1994; see also Patton, 2002.

<sup>122</sup> See, for example Creswell, 2002, and Patton, 2002.

parents' perceptions of the important (and not-so-important) aspects of the training and resources provided by the Center. Additionally, the open-ended responses will allow the interviewer to probe for additional details. In this case, after the interviews are completed, the study team can look for patterns across respondents with respect to which training or resources were most frequently cited as positive or negative influences on parents' ability to successfully support their children with developmental disabilities. While this information is generally formative in nature, the responses also can help the study team to generate hypotheses about how certain activities and resources offered by the Center affect parents' skills and knowledge that could later be studied as part of a summative evaluation.

In unstructured interviews **it's also important to avoid leading the respondent to answer a question in a specific way.** Consider the following examples:

- Question 3: How satisfied are you with your preparation program?
- Question 4: How do you feel about your preparation program?

For the purposes of an unstructured interview, Question 3 has a couple of issues. First, it doesn't encourage the respondent to elaborate; it would be easy for a respondent simply to answer "not very satisfied" or "very satisfied" to this question. Second, the wording of the question effectively limits the responses by identifying the dimension along which the respondent is expected to answer—level of satisfaction.<sup>123</sup> Question 4, on the other hand, both encourages the respondent to elaborate and allows the respondent to discuss any positive, negative, or neutral feelings about the program, whether those feelings relate to satisfaction, cost-effectiveness, quality, or some other dimension.

Significant skill is needed from the interviewer to conduct a successful unstructured interview. As with many skills, the ability to conduct a good unstructured interview can be trained, but again the study team must consider resources. Does the evaluation team include people with good interviewing skills or will the team need to receive training? How much time will be available to practice interviewing? For more information about considerations for conducting an unstructured interview, see Patton (2002).

#### 4.3.3.2 Structured Interviews

In a structured interview, the interviewer asks the participant questions with close-ended response options. For example, an interviewer might read the following question and ask the respondent to choose one of the response options, a-d):

- Question 1: To what extent do you believe the online learning module you completed gave you the knowledge you need to successfully advocate for yourself as you transition out of high school into college?
  - a. Great extent
  - b. Moderate extent
  - c. Little extent
  - d. Not at all

As can be seen, such structured questions do not allow the interviewer to gather additional information beyond that which is included in the interview protocol. Nevertheless, in some situations this inherent limitation also can be a strength, since much less skill is needed from the interviewer to conduct a structured interview compared to an unstructured one. Additionally, structured interviews require fewer resources (i.e., time and money) to collect and analyze the data, and responses can be compared relatively easily across respondents, allowing the study team to

---

<sup>123</sup> Patton, 2002.

generate descriptive information about the interview responses such as tallies of specific responses and cross-tabulations.<sup>124</sup>

The semi-structured interview, discussed in the next section, provides additional flexibility and may be desirable when the study team wants to keep the interview relatively structured while allowing the respondent to elaborate on specific questions.

#### 4.3.3.3 *Semi-Structured Interviews*

Semi-structured interviews combine elements of the unstructured and structured interviews. Generally, semi-structured interviews are composed primarily of structured questions, with one or two unstructured questions that allow the respondent to elaborate on a particular topic or issue. For example, an interviewer might ask the following question:

- Question 1: To what extent do you believe the Technical Assistance Center provided you with the information and resources you need to recognize the early signs of autism in children?
  - a. Great extent
  - b. Moderate extent
  - c. Little extent
  - d. Not at all

Please explain: \_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

As in the case of structured interviews, less skill is needed to conduct a semi-structured interview compared to an unstructured one. Additionally, the semi-structured interview requires fewer resources to collect and analyze the data, although more than the structured interview.

#### 4.3.3.4 *Considerations for Collecting and Analyzing Interview Data*

As with the other aspects of a study, interview protocols must be designed with the end in mind. The more unstructured an interview, the greater the burden on the study team in terms of data collection and analysis. When making decisions about what kinds of interview questions to ask, the study team must think carefully about what they are going to do with the data they collect. Although some researchers would say “More is better” when it comes to collecting data, there is no need to collect data that will not be used. If, for example, the general purpose of the interview is to generate a list of the key program elements that graduates perceive contributed to their success on the job, a structured or semi-structured interview will probably suffice. If, on the other hand, the study team is interested in knowing how the field experiences of graduates from different fields (e.g., psychology, counseling, or secondary transition) compare, a semi-structured or unstructured interview may be better. However, the desire to ask in-depth, probing questions must be weighed against the study team’s ability to spend the time and money required to collect and analyze the resulting data.

---

<sup>124</sup> Creswell, 2002.

We recommend recording the interview—after obtaining permission from the respondents—whenever possible. However, if no one on the study team will have the time or interest to listen to the recording (let alone transcribe it), there is no need to record the interview in the first place. Some benefits of recording an interview include:

- A recording device doesn't select the information to record (as human interviewers do),
- Recordings allow the study team to listen to the interview again to fill in any areas where the interviewer's notes may be unclear or incomplete, and
- A recording enables the study team to make a transcription of the interview at a later date for more in-depth analysis, if desired. (Note: Being able to record an interview doesn't mean that the interviewer shouldn't take notes—a recording can always fail.)

Also, the respondent should be told that they can elect to turn the recording off after initially agreeing to record the interview, and may do so during particularly sensitive times or questions. If the respondent requests that information be “off the record,” neither the written nor verbal notes taken during the interview can be used in data analysis.

For more information, Patton (2002) presents a very detailed and informative discussion of how to conduct interviews and collect and analyze interview data.

#### 4.3.4 Focus Groups

A focus group<sup>125</sup> generally includes 6-10 individuals who have similar backgrounds, such as a cohort of parents of children with emotional/behavioral disorders. When organizing a focus group it's important to remember that the discussion should be focused around a specific topic; group members shouldn't be encouraged to explore a wide variety of topics. It's expected that the group members will influence and build upon each other's responses to the moderator's questions, but it's not primarily a discussion. The focus group is conducted in order to obtain answers to the questions outlined in a focus group protocol.<sup>126</sup>

As with all data collection activities, focus groups have benefits and drawbacks. Some of the **benefits of focus groups** relate to:

- **Cost-effectiveness**—It's a quicker way to get information from multiple people than an individual interview.
- **Data quality**—Interactions among participants serve as checks and balances on false or extreme views.
- **Consistency of viewpoints**—It's relatively easy to see the extent to which there is consistency or great diversity of individual views among group members.
- **Participant experience**—Focus groups are generally viewed as enjoyable for the group members.<sup>127</sup>

Conversely, some of the **drawbacks or limitations of focus groups** concern:

---

<sup>125</sup> Some researchers make a distinction between focus groups and group interviews, pointing out that the way a group interview is conducted is generally different from a focus group. That is, in group interviews the interviewer tends to have a greater role (e.g., generally by directing the line of questioning) while in focus groups the interviewer acts as more of a mediator of the group discussion (Patton, 2002). We do not believe this distinction is particularly important for our discussion here.

<sup>126</sup> Patton, 2002.

<sup>127</sup> Ibid.

- **Interview questions**—The number of questions that can be asked and answered in a focus group is limited. (Note: For a group of 8 and a 1-hour focus group, plan to ask no more than 10 major questions.)
- **Individual responses**—The moderator must limit individual responses in order to hear from all focus group participants.
- **Moderator skills**—To successfully manage a focus group, the moderator must be adept at managing group processes so as to prohibit one or two individuals from dominating the group.
- **Minority viewpoints**—People whose views might be in the minority might choose not to speak up to avoid negative reactions of other group members.
- **Confidentiality**—It isn't possible to ensure confidentiality of responses in focus groups.<sup>128</sup>

Social media and the internet make conducting focus groups with individuals spread out across multiple locations relatively easy. A focus group could be conducted over Skype or another similar internet service, through instant messaging, or through social media. This would greatly reduce the time and cost associated with conducting a focus group, as the group members and the mediator could participate virtually, rather than traveling to attend an in-person meeting.

There is no guideline regarding the number of focus groups to conduct—this decision, too, will be driven by the needs of the study and the available resources.

### 4.3.5 Goal Attainment Scaling (GAS)

Goal attainment scaling (GAS) is a technique to quantify the achievement (or lack of achievement) of goals set,<sup>129</sup> but it also qualifies as a data collection method, which is why it's included here. Initially developed by Kiresuk and Sherman (1968) to evaluate comprehensive community mental health programs, GAS is currently being used by professionals in a variety of fields. GAS has been used to evaluate the effectiveness of individual interventions in producing change, as well as to evaluate program effectiveness.<sup>130</sup> For example, GAS has been used to evaluate special education programs<sup>131</sup> and an early childhood intervention project.<sup>132</sup> It also has been used to set, monitor, and evaluate goals for individuals with learning disabilities,<sup>133</sup> cognitive disabilities,<sup>134</sup> autism spectrum disorders,<sup>135</sup> and traumatic brain injuries.<sup>136</sup>

The basic methodology of goal attainment scaling (GAS) involves (a) selection of the target behavior, (b) an objective description of a desired intervention outcome, and (c) the development of multiple descriptions of the target behavior.<sup>137</sup> The processes for developing the GAS goals and creating the scale are described in more detail below.

<sup>128</sup> Ibid.

<sup>129</sup> Turner-Stokes, 2009.

<sup>130</sup> Sladeczek, Elliott, Kratochwill, Robertson-Mjaanes, & Stoiber, 2001.

<sup>131</sup> E.g., Carr, 1979; Maher, 1983.

<sup>132</sup> Barnett et al., 1999.

<sup>133</sup> Glover, Burns, & Stanley, 1994.

<sup>134</sup> Bailey & Simeonsson, 1988.

<sup>135</sup> Oren & Ogletree, 2000.

<sup>136</sup> Mitchell & Cusick, 1998.

<sup>137</sup> Kiresuk & Sherman, 1968; Morrison, 2012; Sladeczek et al., 2001.

#### 4.3.5.1 Developing the GAS Goals

In the context of GAS, goals have two characteristics. First, a goal is an intended future state, and will usually involve a change from the current situation—although maintaining the current state in the face of expected deterioration could also be a goal. Second, a goal refers to the result of actions undertaken by the actors involved in an intervention, such as a teacher and her students.<sup>138</sup> Good GAS goals are individualized; challenging, but realistic and achievable; able to be written without too much effort, time, or specific training; and flexible enough to cover most situations. Also, they allow for accurate, unambiguous determination of goal achievement. To facilitate this, four aspects should be considered and incorporated into the definition of the goal: the target activity, the support needed, quantification of performance, and the time period to achieve the desired state.<sup>139</sup> Bovend'Eerd et al. (2009) presented a flowchart for writing goals in GAS that outlines a process for developing goals that takes these four aspects into consideration.

Some authors argue for the establishment of goals by one or two specified “goal selectors” who are responsible for selecting goals and creating a scale for individuals based on a pre-specified intervention plan.<sup>140</sup> Others call for a more collaborative approach, whereby members of a team (e.g., a counselor, a teacher, a parent, and a student) work together to identify the main problem areas and establish priority goals for achievement by an agreed date.<sup>141</sup> Whichever approach one takes, an explicit goal needs to be established and this goal will serve as the criterion for success of the intervention.<sup>142</sup>

#### 4.3.5.2 Constructing the Goal Attainment Scale

The GAS ratings should include “sufficiently precise and objective descriptions to enable an unfamiliar observer to determine whether the [individual] lies above or below that point.”<sup>143</sup> The literature generally agrees upon the use of a five-point scale ranging from the least favorable (-2) to most favorable (+2) outcome to measure each goal,<sup>144</sup> although some have argued for the use of a seven-point scale (-3 to +3) to allow for greater sensitivity to change in the achievement of goals.<sup>145</sup> Some scales are constructed such that the “0” rating represents “no change in behavior/performance,” while in other scales the “0” rating represents the “expected outcome” of an intervention.<sup>146</sup> These decisions should be made in advance to ensure consistency of ratings for each goal. Additionally, Kiresuk and Sherman advise that “the scale points be stated in terms of events the presence or absence of which can be easily judged by a follow-up worker who has had no contact with the [intervention procedures].”<sup>147</sup>

---

<sup>138</sup> Wade, 2009.

<sup>139</sup> Bovend'Eerd, Botell, & Wade, 2009; Morrison, 2012.

<sup>140</sup> Kiresuk & Sherman, 1968.

<sup>141</sup> Sladeczek et al., 2001; Turner-Stokes, 2009.

<sup>142</sup> Sladeczek et al., 2001.

<sup>143</sup> Kiresuk & Sherman, 1968, p. 445.

<sup>144</sup> Kiresuk & Sherman, 1968; Kiresuk, Smith & Cardillo, 1994; Sladeczek et al., 2001; Turner-Stokes, 2009.

<sup>145</sup> Bovend'Eerd et al., 2009.

<sup>146</sup> Kiresuk & Sherman, 1986; Kiresuk et al., 1994; Roach & Elliott, 2005; Turner-Stokes, 2009.

<sup>147</sup> Kiresuk & Sherman, 1968, p. 447.

In general, a goal attainment scale should look something like this:

Behavior/Performance Rating:	-OR-	Behavior/Performance Rating:
+2: Change Much More than Expected		+2: Much Positive Change
+1: Change Somewhat More than Expected		+1: Some Positive Change
0: Change At Expected level		0: No Change (e.g., anchored to baseline data)
-1: Change Somewhat Less than Expected		-1: Some Change in Wrong Direction
-2: Change Much Less than Expected		-2: Much Change in Wrong Direction

Roach and Elliott (2005) presented a list of characteristics or dimensions to facilitate the development of descriptions for the different GAS ratings including, among others, frequency (never, sometimes, very often, almost always, always), amount of support needed (totally dependent, extensive assistance, some assistance, limited assistance, independent) and development (not present, emerging, developing, accomplished, exceeding). They also created a simple template for developing goal attainment scales that can serve as a helpful example for anyone using GAS in an evaluation.

GAS offers the option to weight goals to account for their relative importance or difficulty. Turner-Stokes (2009) offered suggested weighting scales for goal importance and difficulty.

Importance	Difficulty
0 = not at all (important)	0 = not at all (difficult)
1 = a little (important)	1 = a little (difficult)
2 = moderately (important)	2 = moderately (difficult)
3 = very (important)	3 = very (difficult)

To determine the weight, simply multiply the importance and the difficulty:

$$\text{Weight} = \text{Importance} \times \text{Difficulty}$$

The weights can then be incorporated into the mathematical formula for calculating standardized GAS scores (discussed below).

#### 4.3.5.3 Collecting GAS Data

GAS can be applied to “any form of objectively determinable event.”<sup>148</sup> For example, goal attainment scales can be developed for data collected through a psychometric instrument, public records, self-ratings, and autobiographical reports. In the context of an OSEP project evaluation, data collected through child and family self-ratings, service provider observations, IEPs, IFSPs, learning objectives, and assessment scores can all be measured using goal attainment scaling. *“Student GAS self-ratings can function as either self-monitoring (a form of direct assessment, completed as behavior occurs) or self-report (a less direct measure of an individual’s perception of their behavior). GAS ratings completed by teachers and other adults function as informant reports, which are also indirect reports because they represent an observer’s retrospective perceptions of behavior.”*<sup>149</sup>

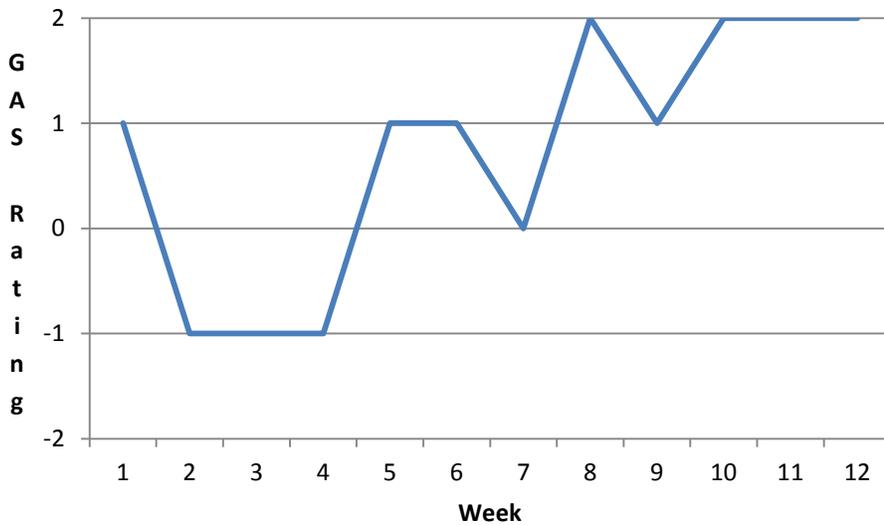
<sup>148</sup> Kiresuk & Sherman, 1968, p. 447.

<sup>149</sup> Roach & Elliott, 2005, p. 9.

Roach and Elliott (2005) recommended conducting initial assessments of goal attainment to establish a baseline for behavior/performance, and then assigning the score of 0 to the baseline data.<sup>150</sup> GAS data can be collected daily, weekly, monthly, or in other time intervals (e.g., every six months), depending on the nature of the intervention and the goal.

GAS ratings can be plotted and viewed as a measure of intervention-induced change.<sup>151</sup> Figure 9 below is an example of GAS ratings that have been plotted to show change over time.

**Figure 9. Sample GAS Ratings for a Special Education Student**



As can be seen in the above example, establishing a baseline when developing the GAS rating scale such that 0 represents the status at baseline will better illustrate changes in performance due to an intervention. Ratings of interventionists, teachers, parents, and even students can all be plotted on the same graph to give an overall picture of performance as perceived by the different individuals involved. The GAS ratings can also give a common language to communicate the effects of an intervention for a given time interval.<sup>152</sup>

#### 4.3.5.4 Calculating a Standardized GAS Score

In some situations practitioners may want to compare GAS scores with other measures of intervention outcomes that are presented in *T*-scores. Kiresuk and Sherman (1968) developed a mathematical formula to derive standardized *T*-scores (with a mean of 50 and a standard deviation of 10) which allow for comparison across individuals and contexts. The standardized composite goal attainment score equals the sum of the goal attainment levels times the relative weights for each goal. It's calculated as follows:

$$\text{Overall GAS} = 50 + \frac{10 \sum (W_i X_i)}{\sqrt{(1-\rho) \sum W_i^2 + \rho (\sum W_i)^2}}$$

<sup>150</sup> See also Kratochwill, Elliott, & Rotto, 1995.

<sup>151</sup> Morrison, 2012; Roach & Elliott, 2005; Sladeczek et al., 2001.

<sup>152</sup> Sladeczek et al., 2001.

Where  $W_i$  is the weight assigned to the  $i$ -th goal (if equal weights,  $W_i = 1$ ),  $X_i$  is the numerical value achieved (between -2 and +2 on a five-point scale),  $\rho$  is the expected correlation of the goal scales. Cardillo and Smith (1994)<sup>153</sup> found that there are high correlations among weighted and unweighted scores, so it's appropriate and simpler to use unweighted scores (i.e., all scores have a weight of 1). In addition, Kiresuk and Sherman (1968) found that, for practical purposes,  $\rho$  most commonly approximates 0.3, so the equation can be simplified to:

$$\text{Overall GAS} = 50 + \frac{10 \sum (W_i X_i)}{\sqrt{0.7 \sum W_i^2 + 0.3 (\sum W_i)^2}}$$

To reduce the need for calculations, the book by Kiresuk et al. (1994) includes score calculation tables for easy reference and Turner-Stokes (2009) developed a GAS spreadsheet calculator in the form of an [Excel worksheet that is available online](#).

However, Sladeczek et al. (2001) pointed out that “for most practitioners, such a transformation may be somewhat cumbersome, and the transformation to standardized  $T$  scores doesn't add additional information to what the overall goal attainment score is conveying (i.e., change due to the intervention; p. 51). Nevertheless, this transformation may be useful if practitioners want to compare scores across multiple interventions.

#### 4.3.5.5 Strengths/Advantages of GAS

Goal attainment scaling can be a useful approach for measuring outcomes for project evaluations when the data can be readily collected by teachers, service providers, or project staff in the course of their practice. It can be especially useful in a school setting. “GAS promotes clearly operationalized intervention goals and on-going (i.e., time-series) evaluation of student progress, making it a potentially useful tool for special educators and school psychologists working within a responsiveness-to-intervention (RTI) model of special education identification.”<sup>154</sup> Additionally, GAS communicates expectations for performance and enables the practitioner to apply systematic rigor to assessment while focusing on the particular needs of the student.<sup>155</sup> Further, the process of setting goals, in and of itself, may have a positive effect on intervention outcomes.<sup>156</sup>

There are a number of **strengths or advantages of using GAS ratings** to monitor outcome change:

- Once the scale is created, collecting GAS data is time efficient
- GAS:
  - is flexible, allowing for the scaling of various content domains for a student (i.e., personalized-individualized);
  - is conceptually consistent with behavioral consultation;
  - requires minimal skills to collect data;
  - is a nonintrusive assessment method;
  - can be used as a self-assessment;
  - can be used by multiple informants across settings (e.g., home, school, community);
  - can be used repeatedly to monitor perceptions of intervention progress;
  - can be used to document perceptions of intervention outcomes;

<sup>153</sup> Cited in Turner-Stokes, 2009.

<sup>154</sup> Roach & Elliott, 2005, p. 16.

<sup>155</sup> Morrison, 2012.

<sup>156</sup> Evans, 1981.

- is relatively inexpensive to implement once the scale is created; and
- requires minimal skill to interpret data.<sup>157</sup>

#### 4.3.5.6 Limitations/Disadvantages of GAS

Probably the biggest limitation of GAS for project evaluators is the time and effort required to develop good goals and scale definitions. A significant amount of time is needed to construct a goal attainment scale and many evaluators may not have the resources available to devote to the task. Once the scale is developed it can be used repeatedly (assuming the goals are applicable in other situations), but the initial development of multiple scales to be used in a project evaluation may not be feasible without significant resources dedicated to the task.

Another potential issue is the concern that if goals are being set by someone invested in the outcome, such as a parent or teacher of a child with disabilities, the goals might be set at an inappropriate level. For instance, relatively easy goals might be set which would result in lowering expectations for performance and giving an inaccurate picture of children’s abilities. Alternatively, unrealistically high expectations could cause relatively difficult goals to be set, which would result in an inaccurate picture of children’s abilities and likely cause progress that might be made to go unrecorded. For these reasons, it’s important to establish formal procedures for checking the rigor of goals.

Roach and Elliott pointed out other limitations or disadvantages to the use of GAS.

- There is limited published, empirical research on the school-based use of the method.
- GAS is a subjective summary of observations collected over time.
- Goals are not norm-referenced.
- The guidelines for interpretation of performance are determined by parties involved with the intervention, thus subject to bias.
- GAS is a global (i.e., less discrete) accounting of behavior.<sup>158</sup>

Other disadvantages/limitations include:

- GAS has utility for monitoring and evaluating progress toward a specific goal, but it isn’t designed to establish an absolute level of functioning for clients. For example a student with autism may change a specific maladaptive behavior as a result of an intervention, but not have significantly increased his or her absolute level of adjustment, skills, or functioning;<sup>159</sup> and
- It isn’t appropriate to establish causal relationships between independent and dependent variables in GAS.<sup>160</sup>

We now turn our discussion to the various methods that might be used in analyzing data collected through a PDP project evaluation.

---

<sup>157</sup> Roach & Elliott, 2005.

<sup>158</sup> 2005, p. 15.

<sup>159</sup> Smith & Cardillo, 1994.

<sup>160</sup> Sladeczek et al., 2001.

## 4.4 Data Analysis Methods

It's beyond the scope of this Toolkit to discuss the myriad methods of quantitative and qualitative data analysis. In this section we focus on basic methods of quantitative and qualitative analysis that can readily be applied in a PDP project evaluation. Before applying any analysis method, however, evaluators first have to identify, and respond to, the prevalence of missing data, discussed below.

### 4.4.1 Dealing with Missing Data

Almost every study will have at least some missing data, so evaluators must have a plan for dealing with missing data in the analysis. Addressing the problem of missing data is particularly important when conducting inferential statistical analyses, but having missing data in a study can potentially affect the results, no matter what type of analysis is being conducted. The strategy that is adopted to deal with missing data will depend in part on the nature of the data that are missing. Howell (2012) pointed out several reasons why data might be missing:

- **Data missing completely at random**—The probability that an observation ( $X_i$ ) is missing *is unrelated* to the value of  $X_i$  or to the value of any other variables. If data are missing completely at random, the analysis will remain completely unbiased, since the estimated parameters aren't biased by the absence of these data.
- **Data missing at random**—The “missingness” of the data doesn't depend on the value of  $X_i$  *after controlling for another variable*. That is to say, if missingness is correlated with other variables in the analysis, the data are considered missing at random. The presence of data that are missing at random can bias the estimates if steps to deal with the missing data aren't taken.
- **Data missing not at random**—The probability that an observation ( $X_i$ ) is missing *is related* to the value of  $X_i$  or to the value of another variable. When data are missing not at random it's necessary to write a statistical model that accounts for the missing data.

It's not within the scope of this Toolkit to go into detail about the different methods to treat missing data. The U.S. Department of Education provides guidance on [what to do with missing data in experimental studies](#). For additional information, see Howell (2012), Allison (2001), and Baraldi and Enders (2010).

### 4.4.2 Quantitative Analysis

As we described in [Section 3.4.1](#), prior to conducting quantitative analysis evaluators need to enter the data into a database and prepare them for analysis. Once all of the data are entered into the database and the analysis program has been selected (e.g., SPSS, SASS, STATA or R for quantitative analysis), it's time to begin exploring and conducting descriptive and inferential analysis of the data, discussed below.

#### 4.4.2.1 Descriptive Analysis

Descriptive statistics provide information about the overall trends and distribution of the data. This includes reporting percentages or frequencies for nominal and ordinal data, as well as conducting exploratory analysis of interval and ratio data. Exploratory analysis includes looking for outliers in the data, determining the modality (e.g., unimodal or bimodal) of the data distribution, and calculating kurtosis (i.e., the sharpness of the peak of the data distribution curve) and

skewness (i.e., a measure of whether the data are distributed symmetrically around the mean).<sup>161</sup> If desired, if the data distribution is flat or the data are skewed, it may be possible to conduct linear or non-linear transformations of the data at this time—for example to change a positively skewed distribution to a more normal one. Evaluators might consider conducting transformations if they wish to use the data with inferential statistics and need to satisfy the assumptions of specific statistical tests.

After exploring the data and examining the shape of the distribution, it's also good practice to calculate the following descriptive statistics:

- **Measures of central tendency**, which give the researcher an idea of where the majority of scores are located in a distribution, when working with interval or ratio data. These include the mean, the median, and the mode.
  - Mean: The arithmetic average of a set of scores.
  - Median: The score above and below which 50% of scores fall.
  - Mode: The most frequently occurring score in the distribution.
- **Measures of variability**, which indicate the spread of the scores in a distribution, also appropriate when working with interval or ratio data. Three key measures of variability are the range, the variance, and the standard deviation.
  - Range: The distance between the highest and lowest score in a distribution.
  - Variance: A measure of how far the scores in the distribution are spread out around the mean.
  - Standard deviation: The average distance of scores from the mean (the square root of the variance).

Depending on the type of analysis and the study questions, evaluators also may want to calculate **measures of relative standing**, which are statistics that describe how a particular score compares to a group of scores. Three common measures of relative standing include the z-score, the percentile, and the percentile rank—all assume the data are interval or ratio in nature.<sup>162</sup>

- Z-score: A standardized score that shows how many standard deviations a score falls above or below the mean. The z-score has a mean of 0 and a standard deviation of 1. Raw scores can be converted to z-scores so that scores from a variety of measures (such as achievement tests from different states) can be compared.
- Percentile: The score below which a certain percent of scores fall.
- Percentile rank: The percent of scores that fall below a given score.

For some of the study questions (e.g., To what extent do PDP graduates report satisfaction with their field experience?), frequency distributions or descriptive analyses may be all that are needed. For more complex questions, however—particularly those related to the summative evaluation questions (e.g., Do students of graduates from the “Great University” PDP program perform better on the state achievement test in literacy than students of other teacher training programs?)—evaluators likely will need to conduct some type of inferential statistical analysis, discussed briefly in the next section.

---

<sup>161</sup> For information on how to do these procedures in SPSS, see Dimitrov, 2010.

<sup>162</sup> Creswell, 2002.

#### 4.4.2.2 Inferential Statistical Analysis

Inferential statistics are divided into two types: parametric and non-parametric. *“Parametric statistics are statistical tests based on the premise that the population from which samples are obtained follows a normal distribution and the parameters of interest to the researcher are the population mean and standard deviation.”*<sup>163</sup> Parametric statistics are based on certain assumptions, including:

1. The variables are measured in interval or ratio scales.
2. Scores from any two individuals are independent—one person’s score isn’t dependent on another person’s score.
3. The distribution of population scores is normally distributed.
4. When two or more groups are involved in the study, each group representing a different population, the populations have equal variances.<sup>164</sup>

These assumptions typically hold true when large numbers of individuals are included in the sample; when the sample is small, however, one or more of these assumptions can be violated, thereby affecting the inferences that can be drawn from the analyses.

Non-parametric statistics have fewer assumptions than parametric statistics, making them more appropriate when the sample is small or when the data aren’t normally distributed. In fact, the only assumption that applies to non-parametric statistics is the assumption of independence of scores (number two in the list above). Non-parametric tests are typically not as sensitive as parametric tests, however, so whenever possible we encourage evaluators to use parametric tests. Additionally, the typical parametric tests are able to withstand violations of some of the assumptions, so in many cases it’s reasonable to use a parametric test even when some of the assumptions presented above are violated. We recommend evaluators without a strong quantitative analysis background consult a statistician for help determining which statistic to use. When designing or identifying data collection instruments, evaluators should consider the types of data that will be produced by each instrument so that they can be sure to collect enough of the right data for the specific analysis method chosen.<sup>165</sup>

---

<sup>163</sup> Creswell, 2002, p. 237.

<sup>164</sup> Ibid.

<sup>165</sup> Ibid.

Creswell (2002) highlighted seven factors that go into the decision of what type of statistic to use:

1. Whether the evaluation wants to compare groups [or performance at different time periods] or relate one or more variables;
2. The number of independent variables;
3. The number of dependent variables;
4. Whether covariates will be included in the analysis, and if so, how many;
5. The scale of measurement for the independent variable;
6. The scale of measurement for the dependent variable; and
7. Whether the data are normally distributed, bimodal, or skewed.

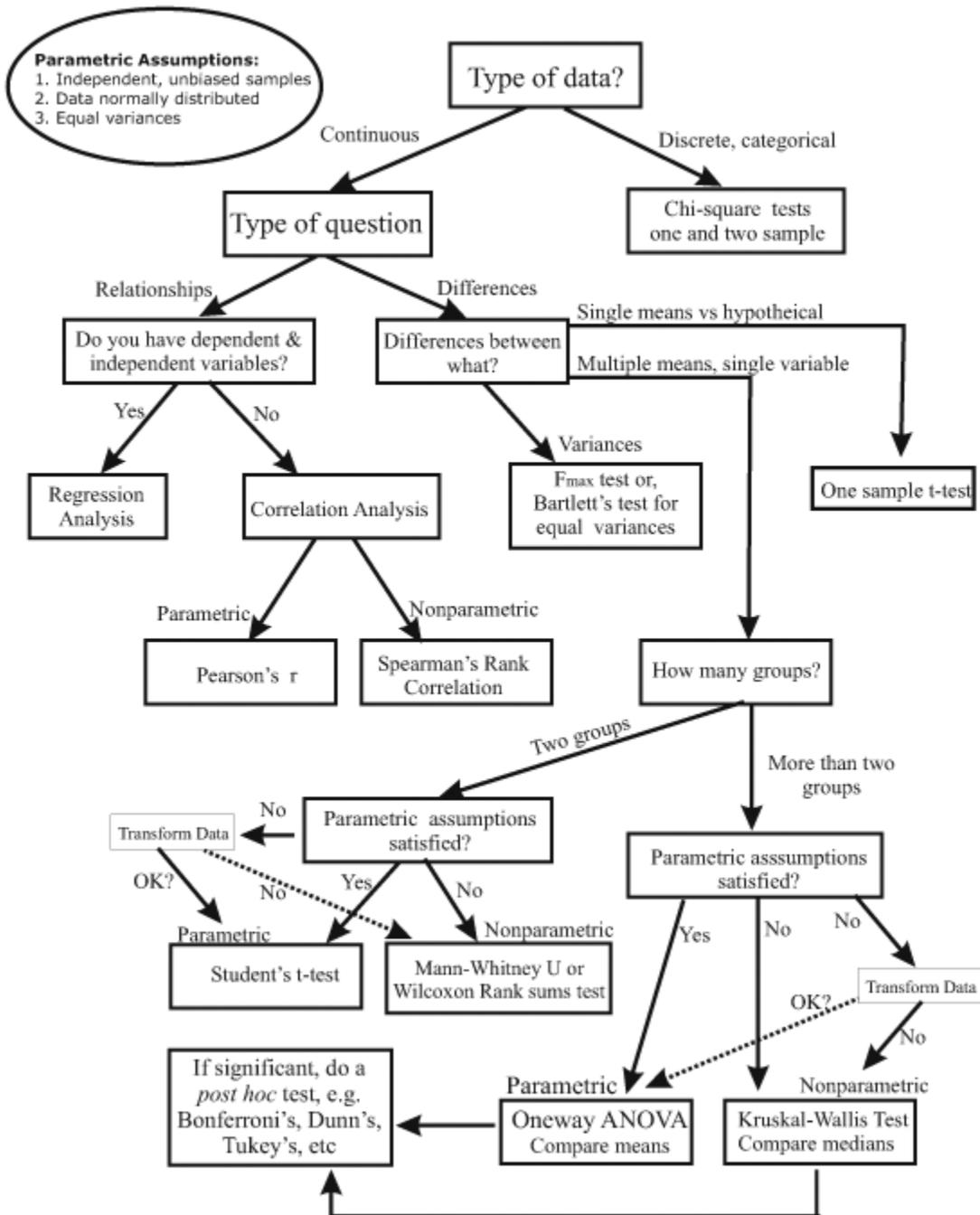
Carefully considering each of these factors will help evaluators to select the most appropriate and feasible test to conduct. Table 15 outlines some common parametric and non-parametric statistical tests that might be used in an evaluation and Figure 10 presents a decision tree to further help with the decision making. Box 2 presents an example of how this might be done in an evaluation.

**Table 15. Common Parametric and Non-Parametric Statistical Tests**

Type of Question/Hypothesis	Parametric Test	Non-parametric Test
<b>Group comparison (2 groups)</b>	<i>T</i> -test (independent samples) Also known as Student's <i>t</i> -test	Mann-Whitney <i>U</i> test
<b>Within-subject comparison</b>	<i>T</i> -test (dependent samples)	Wilcoxon Signed-Rank Test
<b>Group comparison (3+ groups)</b>	Analysis of variance (ANOVA)	Kruskall-Wallis test
<b>Group comparison (3+ groups, multiple measures)</b>	Repeated measures ANOVA	Friedman's Chi-square ( $\chi^2$ ) test
<b>Relate variables</b>	Pearson product moment correlation (Pearson <i>r</i> )	Spearman's <i>rho</i> rank correlation
<b>Within-group comparison of frequencies</b>		Chi-square ( $\chi^2$ ) goodness of fit test
<b>Comparison of frequencies between groups or variables</b>		Chi-square ( $\chi^2$ ) test for association

Figure 10. Decision Tree for Inferential Statistics

### Flow Chart for Selecting Commonly Used Statistical Tests



Source: <http://abacus.bates.edu/~ganderso/biology/resources/statistics.html#whichtest>

## Box 2. Determining which Inferential Test to Use

Using the decision tree for inferential statistics above, let's examine one of the study questions from the sample Evaluation Plan presented in Appendix A.4 to determine which type of test might be used.

**Sample Question:** *To what extent do teachers who participated in a model demonstration training project for language instruction exhibit improved language instruction in the classroom?*

To answer this question, we will compare teachers who did and did not participate in the training program.

1. Type of data. The data we are using for our outcome measures is teacher workplace performance. To develop this outcome variable a continuous score was derived by summing the following data sources:
  - a. the results of the questions in a supervisor survey related to teacher performance on the Council for Exceptional Children (CED) standards; and
  - b. ratings from classroom observations. (Alternately the evaluator could maintain these as two dependent variables and conduct two analyses.)
2. Type of question. We are concerned with the differences, if any, between our groups of treatment and comparison teachers.
3. Differences between what? We will have two means, one for each teacher group, on the variable of interest—the teacher workplace performance score we created.
4. How many groups? We have our trained teachers and the comparison group of teachers who did not participate in the model training project, so analysis will be comparing performance of two groups.
5. Are the parametric assumptions satisfied? Let's consider each assumption in turn.
  - a. Are the samples independent and unbiased? Yes.
  - b. Is the data normally distributed? Descriptive data indicates "Yes."
  - c. Do the groups have equal variance? Descriptive data indicates "Yes."

Based on this information, an independent samples *t*-test (Student's *t*-test) could be used.

For additional information on the use of these tests, Hinkle, Wiersma, and Jurs (2003) and Dimitrov (2010) are good reference books. Additionally, Rice University, the University of Houston Clear Lake, and Tufts University have developed an [online statistics book](#) that is free and available to the public.

### 4.4.2.3 Multilevel Analysis

An important feature of most educational research and evaluations is that the data typically have a multilevel structure. That is, data that are collected at one level (e.g., student outcome data) are frequently “nested” within clusters at different levels (e.g., within classrooms, within schools). Multilevel analysis examines relationships between variables measured at different levels of the multilevel data structure.<sup>166</sup>

**Two basic multilevel designs are frequently used in education:**

- **The hierarchical design**—Entire clusters are assigned to treatment (e.g., entire schools or classrooms are assigned to treatment or control groups); and
- **The block design**—Individuals within the same cluster are assigned to two different treatments (e.g., students assigned to treatment and control groups within the same school).

It’s beyond the scope of this Toolkit to delve into this topic; however, evaluators planning an evaluation that features an experimental or quasi-experimental design and who are intending to measure the outcomes of students should consider using multilevel analysis to account for the clustering of students within schools (or classrooms, districts, etc.). Evaluators without training in this subject should consult a statistician or third-party evaluator for guidance on how to design the evaluation and conduct the multilevel analyses. Hox (2010) provides a nice introduction to multilevel analysis; Bryk and Raudenbush (1988); Raudenbush (1997) and Schochet (2009) provide more technical details on how to conduct multilevel analysis.

In the next section we move our discussion to analysis of data in single-case designs.

### 4.4.3 Analysis of Data in Single-Case Designs

Evaluation of outcomes in single-case studies can be conducted through a variety of means, most typically including visual analysis of the data or calculation of summary statistics such as the percentage of non-overlapping data points<sup>167</sup> and the standard mean difference effect size. Additionally, it’s possible to conduct a randomization test, which is a specific type of statistical analysis of data from single-case or small *n* studies (e.g., studies with fewer than 10 participants).<sup>168</sup> We briefly discuss these below. For more information about data analysis in single-case studies, see Barton & Reichow (2012); Horner & Spaulding (2010); Kennedy (2005); Kratochwill, Hitchcock, Horner, et al. (2013); Kratochwill & Levin (2010); and Todman & Dugard (2001).

---

<sup>166</sup> Hox, 2010 and Hedges & Rhoads, 2013.

<sup>167</sup> Scruggs, Mastropieri & Casto, 1987.

<sup>168</sup> Todman & Dugard, 2001.

#### 4.4.3.1 Visual Analysis of Single-Case Data

As data are collected they are plotted in a graph (as in Figure 4 and Figure 5 in [Section 4.1.3](#) above) and patterns are studied over the course of the study. Ideally, a person should be able to look at a graph of data collected from a single-case design and understand what the data represent. When conducting visual analysis of the data, the researcher is assessing the following aspects of the data:

- Level
- Trend
- Variability
- Immediacy
- Consistency
- Percentage of Overlap

When examining **changes in level**, the researcher calculates the mean or the median for the data within a given condition (e.g., for the A and B phases of an A-B-A-B design). This allows for the estimation of the central tendency of the data during a particular part of the experiment as well as a comparison of patterns between phases. The last few data points in each phase “*contain the most essential information regarding the level of behavior before a phase change*”<sup>169</sup> because they give clues to any changes in pattern. If these last few data points within a given phase appear to be showing an upward (or downward) trend, there is reason to consider that any changes in level in the subsequent phase may be part of a natural progression that started in the prior phase.

**Changes in trend** of the data refers to the best-fit straight line that can be placed over the data within a phase. Two elements of the trend are important: slope and magnitude. The slope tells the direction of the change within a phase while the magnitude indicates the size or extent of the slope (generally estimated qualitatively as high, medium, or low). (Note: The greater the slope, the less meaningful changes in level are to the general estimate of the data pattern.) To judge the trend of the data visually, the researcher simultaneously estimates the slope and the magnitude of the data.<sup>170</sup> It’s also possible to estimate the trend quantitatively; see Kennedy (2005) for more information.

The extent of **changes in variability** shows the degree to which individual data points deviate from the general trend of the data. To estimate variability, the researcher examines the degree to which the data points are spread out from the best-fit straight trend line. Variability is typically referred to as high, medium, or low.<sup>171</sup>

When looking at the **immediacy** of an effect, the researcher looks at how quickly a change is noticeable when moving from one phase to the next. This is done by examining last 3 data points in one phase and the first 3 in the next to see the magnitude of the change.

When examining **consistency** of the data patterns, the researcher looks at the extent to which data patterns are similar in similar phases.

We talk about a common method for calculating **percentage of overlap** in the next section.

For information on how to create single-case design graphs using Microsoft Excel, see Barton & Reichow (2012).

---

<sup>169</sup> Kennedy, 2005, p. 197; see also the WWC standards for single-case designs (Kratochwill, Hitchcock, et al., 2010).

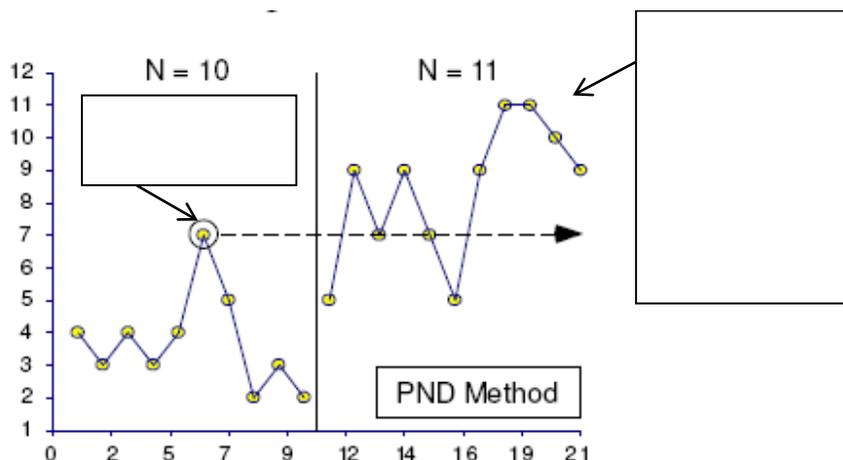
<sup>170</sup> Kennedy, 2005.

<sup>171</sup> Ibid.

#### 4.4.3.2 Percentage of Non-overlapping Data Points

The calculation of percentage of non-overlapping data points (PND) involves identifying the highest point in the baseline data and determining the percentage of data points in the subsequent intervention phases that exceed the highest baseline data point.<sup>172</sup> Figure 11 illustrates how this is conducted. In this example the PND = 7/11 = 60%.<sup>173</sup>

Figure 11. Calculating Percentage of Non-Overlapping Data Points



Source: Wendt, 2009

Scruggs, Mastropieri, Cook and Escobar (1986) outlined specific guidelines for interpreting PND scores. In general, the higher the percentage of non-overlapping data points, the more effective the treatment. If a study includes multiple single-case experiments, the PND scores are aggregated by calculating the median, rather than the mean. The median is used in this case because scores in single-case designs aren't usually distributed normally and it's less subject to outliers than the mean. While the PND is relatively easy to interpret, it has two major limitations. First, it ignores the majority of the baseline data and may be subject to ceiling effects as the study participant reaches the top of the performance scale. Next, it lacks sensitivity and ability to discriminate performance as the PND nears 100%.<sup>174</sup>

#### 4.4.3.3 Standard Mean Difference Effect Size

An effect size is a descriptive statistic that provides an estimate of the magnitude, or strength, of group differences or of the relationship among variables. Effect sizes are commonly calculated in meta-analyses to summarize findings in a specific area of research and to allow for comparisons across studies. Further, the American Psychological Association (2009) recommends reporting effect sizes along with the results of significance tests and confidence intervals to give more information about the magnitude or importance of a study's findings. There are many ways to calculate effect sizes, including regression based and non-regression based effect sizes. Each type of effect size calculation has advantages and disadvantages, but not all are considered user-friendly for practitioners in terms of their simplicity of calculation and interpretation (for more technical information on how to calculate effect sizes in single-case research,

<sup>172</sup> Scruggs et al., 1987.

<sup>173</sup> Wendt, 2009.

<sup>174</sup> Wendt, 2009.

see Shadish, Hedges, Horner & Odom, 2015).<sup>175</sup> In this section we focus on a particular type of non-regression-based effect size—the standard mean difference.<sup>176</sup>

The **standard mean difference** (SMD) effect size is calculated by determining the difference between the mean baseline score and mean intervention score and then dividing by a standard, such as the standard deviation of the baseline.<sup>177</sup> There are two variations of SMD. Specifically, SMD may be calculated using the mean for all baseline and intervention data points (SMDall) or it may be calculated using the mean from only the last three data points of each phase (SMD3) and then multiplying by 100. This latter calculation is also known as the percentage reduction or the mean baseline reduction. Olive and Smith (2005) recommend the use of SMDall to complement visual analysis because it yields similar results to other summary statistics, it's relatively simple to calculate, and it doesn't have some of the same assumptions as regression-based effect size models (i.e., regression-based models are based on the assumption that linearity exists).<sup>178</sup>

Cohen (1988) suggested that a small effect size is  $d = 0.2 - 0.49$ , a medium effect size is  $d = 0.5 - 0.79$ , and a large effect size is  $d \geq 0.8$ . This standard has been widely used by researchers and evaluators in the field of education, despite criticisms of its application to education interventions. Rather than using Cohen's broad categorization of the magnitude of effect sizes across intervention research as a whole, we recommend evaluators consult Lipsey et al. (2012), who have outlined specific norms for effect sizes for educational interventions.

Campbell (2004) pointed out two statistical problems that affect the calculation of effect sizes in single-case studies. First, single-case data are auto-correlated, meaning that observations are temporally ordered and usually not independent. This violates one of the basic assumptions of most statistical tests of significance. Second, effect sizes can be confounded by a trend in the data, thereby causing the effect size to be under or overestimated.<sup>179</sup> Another problem with using only the SMD effect size is that valuable information inherent in the visual analysis of the data is lost (i.e., data variability, trend magnitude and direction, mean levels and shifts, and embedded cycles within phases).

#### 4.4.3.4 Randomization Tests

Single-case and small- $n$  designs cannot use the same parametric statistical procedures that are typically used in larger studies, since the smaller the sample size the less confidence a researcher can have that the parametric assumptions are met.<sup>180</sup> Instead, non-parametric analyses are recommended for studies with a relatively small number of participants. These non-parametric tests are generally based on rankings of scores (e.g., the Mann-Whitney U and Wilcoxon T non-parametric alternatives to the independent samples and dependent samples  $t$ -tests, respectively), but they lack sensitivity to real treatment effects in studies with very small numbers of participants. However, in some cases, rather than using non-parametric tests, the researcher can conduct randomization tests, which (a) do not discard information in the data by reducing them to ranks, and (b) can provide valid statistical analysis of data from a wide range of single-case and small- $n$  designs.<sup>181</sup>

While we cannot present a full discussion of randomization tests here, we want to point out that randomization tests can increase confidence in the causal inferences that are made in single-case and small- $n$  studies.<sup>182</sup> Further, they aren't

---

<sup>175</sup> Campbell, 2004; Morrison, 2012.

<sup>176</sup> For more information on effect sizes see, for example, Campbell, 2004; Cohen, 1988; Ellis, 2010; and Kelley & Preacher, 2012.

<sup>177</sup> Busk & Serlin, 1992, cited in Morrison, 2012.

<sup>178</sup> Morrison, 2012.

<sup>179</sup> West & Hepworth, 1991, cited in Campbell, 2004.

<sup>180</sup> Siegel & Castellan, 1988, cited in Todman & Dugard, 2001.

<sup>181</sup> Todman & Dugard, 2001.

<sup>182</sup> Ibid.

particularly difficult to conduct, although they do require the researcher (or teacher or other professional) to follow specific procedures.

The basic principle of randomization tests is that some aspect of the experimental design must be randomized. However, randomization tests do not require students to be randomly assigned to specific treatments (or controls) or to certain classrooms; instead, the timing of the introduction of the intervention can be randomly selected. For an A-B single-case design, for example, the researcher randomly selects the point (e.g., day or week) at which to switch between the baseline and intervention phases. The same is true for a multiple baseline design (with some additional requirements). As with any type of analysis there are limitations to the interpretation of the data from randomization tests, but they can provide information about the statistical effects of an intervention. Additionally, randomization tests can be conducted using software such as Microsoft Excel or SPSS. Interested evaluators are urged to consult Kratochwill & Levin (2010) and Todman & Dugard (2001) for more information.

#### 4.4.4 Qualitative Analysis

Unless the data collection instruments are completely structured, quantitative types of surveys, observations, and interviews, it's likely that significant amounts of qualitative data will be gathered during the course of the evaluation. It's beyond the scope of this Toolkit to go into detail about the many different types of qualitative inquiry and analysis approaches.<sup>183</sup> Instead, in this section we make some recommendations for preparing the qualitative data for analysis, briefly discuss the major phases of qualitative analysis, and highlight some benefits and drawbacks of using qualitative analysis software programs. Interested evaluators are urged to consult Miles and Huberman (1994) and Patton (2002) to learn more about how to conduct a variety of qualitative analyses.<sup>184</sup> Maxwell (2005) is a good source of information about how to design a qualitative study (e.g., as part of a mixed-methods evaluation).

##### 4.4.4.1 Preparing Qualitative Data for Analysis

Just as with quantitative data analysis, the first step in qualitative analysis is preparing the data. We recommend evaluators follow these steps to prepare the data for analysis:

- **Enter the data into a spreadsheet or qualitative data analysis software program to facilitate analysis (e.g., to allow searches for particular words or phrases).** –This is particularly helpful when there are large numbers of study participants or when large quantities of qualitative data are expected. Some online survey providers (e.g., Survey Monkey) allow individuals to download a database file with the individual responses to each item, making the process of creating the database very easy.
- **Develop decision rules for how to handle any data problems (e.g., missing data, incomplete or misspelled answers).** –Here, the term “data problems” specifically refers to problems of missing or erroneous data, not data that might not seem to fit a researcher’s preconceived ideas of what the data “should” look like. An example of when a decision rule may be necessary is when a person gives a response to Question 5 on a self-administered survey that actually looks like it should have been a response to Question 6 (while leaving Question 6 blank).

---

<sup>183</sup> Patton, 2002, provides a nice discussion of a variety of orientations and theoretical approaches to qualitative inquiry and analysis. Denzin & Lincoln, 2005, is another valuable resource for those interested in learning more about qualitative inquiry.

<sup>184</sup> Two online resources are Schutt, 2011, and Suter, 2012.

- **Conduct a preliminary check for errors, missing data or other problems with the data.** –When respondents are responsible for entering their answers to survey questions into a database—for example, in an online survey—and the study team plans to use search terms to facilitate data analysis, it’s particularly important to check whether or not respondents have misspelled key terms (e.g., when a person enters “*congitive*” instead of “cognitive”).
- **Apply the decision rules to address any errors.** –Whenever possible, go back to the original respondent for correction or clarification of any errors; if that isn’t possible the decision rules will provide guidance and will ensure consistency in the handling of data problems. It’s important to note that when working with qualitative data, all responses—especially those that might be considered “outliers” in quantitative analysis—potentially can provide valuable insight into the subject of study. For this reason, the qualitative analyst should be cautious in the application of decision rules to deal with data problems. We recommend against deleting responses.

Once the data have been prepared, the major phases of analysis can begin (although in reality the initial phases of data analysis begin during data preparation, as the analyst reviews the questions and begins making sense of the responses). In the next sections we briefly discuss the major phases involved in qualitative analysis: data reduction, data display, and conclusion drawing and verification.<sup>185</sup> It’s important to note that these aren’t discrete phases that happen in sequence. Rather, qualitative analysis is an iterative process that involves “concurrent flows of activity.”<sup>186</sup>

Berkowitz (1997) highlighted some **important questions that the qualitative analyst should keep in mind during the data analysis process:**

- What patterns and common themes emerge in responses dealing with specific items? How do these patterns (or lack thereof) help to illuminate the broader study question(s)?
- Are there any deviations from these patterns? If yes, are there any factors that might explain these atypical responses?
- What interesting stories emerge from the responses? How can these stories help to illuminate the broader study question(s)?
- Do any of these patterns or findings suggest that additional data may need to be collected? Do any of the study questions need to be revised?
- Do the patterns that emerge corroborate the findings of any corresponding qualitative analyses that have been conducted? If not, what might explain these discrepancies?<sup>187</sup>

These questions aren’t specific to any one phase of data analysis and therefore should be addressed as appropriate throughout the study.

<sup>185</sup> This framework comes from Miles & Huberman, 1994; see Berkowitz, 1997, for an overview of the approach.

<sup>186</sup> Miles & Huberman, 1994, p. 10.

<sup>187</sup> 1997, p.2.

#### 4.4.4.2 Data Reduction

“Data reduction refers to the process of selecting, focusing, simplifying, abstracting, and transforming the [qualitative] data.”<sup>188</sup> In fact, data reduction occurs even before any data are collected, as the researcher makes decisions about the conceptual framework for the study, the research questions, data collection approaches, and data sources. As Freeman put it, “*Before I could access what it is I wanted to know, I had to first figure out what I meant by know, who I thought should do the knowing, and where I thought this knowledge could be found.*”<sup>189</sup>

Data reduction generally relates to deciding which aspects of the data should be emphasized; this decision should be guided by the need to address the salient evaluation questions.<sup>190</sup> For example, if the evaluation is focused on the quality of an in-service training for related-services providers, the study team may choose to focus the analysis only on the responses that pertain to students’ perceptions of the quality of the related services they receive and not on students’ in-class experiences with special education teachers.

Data reduction also involves transforming the data in a way that facilitates analysis and understanding. This may include selecting specific quotes to highlight, summarizing or paraphrasing long passages, or developing a narrative summary of a larger pattern of responses across multiple respondents or sites. In some situations it might be helpful to identify “quantities” in the data, such as the number of parents who would rate their experience with a Parent Information Center as “very helpful”. However, Miles and Huberman (1994) cautioned against simply using tables of numbers to illustrate the variety of qualitative responses to a question represented in the data.

For more in-depth discussion of data reduction, see Miles and Huberman (1994); Berkowitz (1997) provides a general overview of the approach.

#### 4.4.4.3 Data Display

Data display is the next major phase of qualitative analysis and it involves arranging the data in an “organized, compressed assembly of information that permits conclusion drawing and action.”<sup>191</sup> “A display can be a long piece of text, or a diagram, chart, or matrix that provides a new way of arranging and thinking about the more textually embedded data.”<sup>192</sup> Data displays allow the analyst to begin identifying systematic patterns and interrelationships within cases (also known as intra-case analysis) and between cases (also known as inter-case analysis).<sup>193</sup>

Intra-case analysis examines patterns in the data across all of the sources within a case, such as parents of children of different ages who have been served by Parent Information Centers. Table 16 presents an example of a data display for one Center. Inter-case analysis examines patterns across different cases within one study, such as among respondents in three different PICs. Table 17 presents a similar display that incorporates information from multiple Parent Information Centers. Berkowitz (1997) pointed out that using data displays in this way allows for a relatively quick recognition of patterns in the data, even before the individual responses are completely analyzed. Additionally, it allows for the identification of differences in responses across the different groups of respondents.

---

<sup>188</sup> Miles & Huberman, 1994, p. 10.

<sup>189</sup> 2000, p. 360.

<sup>190</sup> Berkowitz, 1997.

<sup>191</sup> Miles & Huberman, 1994, p. 11.

<sup>192</sup> Berkowitz, 1997, p. 4.

<sup>193</sup> Berkowitz, 1997.

**Table 16. Qualitative Data Matrix for Parent Information Center Showing Responses for Parents**

What Parent Information Center activities did you utilize?			
Respondent Group	Activities Named	Which most effective?	Why?
Parents of Infants	<ul style="list-style-type: none"> <li>• Parent Information Packets</li> <li>• Workshops</li> <li>• Referrals</li> </ul>	<ul style="list-style-type: none"> <li>• Referrals</li> <li>• Workshops</li> </ul>	<ul style="list-style-type: none"> <li>• Increased awareness of services</li> <li>• Provided practical tips and helped develop skills</li> </ul>
Parents of Elementary School Aged Children	<ul style="list-style-type: none"> <li>• Workshops</li> <li>• Resource Fairs</li> <li>• Advocacy calls</li> </ul>	<ul style="list-style-type: none"> <li>• Workshops</li> <li>• Resource Fairs</li> </ul>	<ul style="list-style-type: none"> <li>• Provided practical tips and helped develop skills</li> <li>• Provided information about resources available</li> </ul>
Parents of Teenagers	<ul style="list-style-type: none"> <li>• Community events</li> <li>• Workshops</li> <li>• Newsletters</li> <li>• Online discussion boards</li> </ul>	<ul style="list-style-type: none"> <li>• Workshops</li> <li>• Community Events</li> </ul>	<ul style="list-style-type: none"> <li>• Provided practical tips and helped develop skills</li> <li>• Promoted social interactions</li> </ul>

Source: Adapted from Berkowitz, 1997.

**Table 17. Qualitative Data Matrix Comparing Responses at Three Parent Information Centers**

Participants' views of activities to improve parents' skill and knowledge in three PICs			
Respondent Group	Activities Named	Which most effective?	Why?
Center A	<ul style="list-style-type: none"> <li>• Newsletters</li> <li>• Community Events</li> <li>• Workshops</li> <li>• Presentations</li> </ul>	<ul style="list-style-type: none"> <li>• Presentations</li> <li>• Workshops</li> </ul>	<ul style="list-style-type: none"> <li>• Concise way of communicating information</li> <li>• Provided practical tips and helped develop skills</li> </ul>
Center B	<ul style="list-style-type: none"> <li>• Workshops</li> <li>• Parent Information packets</li> <li>• Advocacy calls</li> </ul>	<ul style="list-style-type: none"> <li>• Workshops</li> <li>• Advocacy Calls</li> </ul>	<ul style="list-style-type: none"> <li>• Provided practical tips and helped develop skills</li> <li>• Provided individualized information</li> </ul>
Center C	<ul style="list-style-type: none"> <li>• Workshops</li> <li>• Referrals</li> <li>• Resource Fairs</li> <li>• Newsletters</li> <li>• Community Events</li> </ul>	<ul style="list-style-type: none"> <li>• Workshops</li> <li>• Referrals</li> </ul>	<ul style="list-style-type: none"> <li>• Provided practical tips and helped develop skills</li> <li>• Increased awareness of services</li> </ul>

Source: Adapted from Berkowitz, 1997.

As can be seen in the tables above, workshops were consistently rated by all respondents as the most effective Center activity. In this example, all respondents gave the same reason for choosing this activity, but the reason for identifying an activity as the most effective differed across respondents. Data displays allow the qualitative analyst to easily identify these similarities or differences.<sup>194</sup> This is achieved by what Glaser and Strauss (1967) called the “‘method of constant comparison,’ an intellectually disciplined process of comparing and contrasting across instances to establish significant patterns, then further questioning and refinement of these patterns as part of an ongoing analytic process.”<sup>195</sup>

Once the analyst identifies patterns of responses and similarities and differences among different respondent groups, the next step is to delve more deeply into the reasons behind the different responses, which is explored in the next phase of qualitative analysis: conclusion drawing and verification.

#### 4.4.4.4 Conclusion Drawing and Verification

Conclusion drawing involves taking a step back to consider what the analyzed data mean and to assess how they relate to the evaluation questions under study. Verification entails going back to the data as often as necessary to check or validate emergent conclusions.<sup>196</sup> An essential element of this process is testing the conclusions for their plausibility, their sturdiness, their “confirmability”—in other words, their validity.<sup>197</sup> It’s important to note that “validity” in this context is different from the conception of “validity” as part of quantitative analysis. Specifically, in the qualitative context, the “validity” of the analysis relates to “whether the conclusions being drawn from the data are credible, defensible, warranted, and able to withstand alternative explanations.”<sup>198</sup> Miles and Huberman outlined **13 tactics for testing or confirming findings in qualitative analysis**<sup>199</sup>:

1. Checking for representativeness;
2. Checking for researcher effects;
3. Triangulating across data sources and methods;
4. Weighting the evidence;
5. Checking the meaning of outliers;
6. Using extreme cases;
7. Following up surprises;
8. Looking for negative evidence;
9. Making if-then tests;
10. Ruling out spurious relations;
11. Replicating a finding;
12. Checking out rival explanations; and
13. Getting feedback from informants.

It’s beyond the scope of this Toolkit to discuss all of these tactics in detail. For more information about these tactics, and about standards for judging the quality of conclusions drawn from qualitative data, consult Miles and Huberman (1994).<sup>200</sup> Berkowitz (1997) provides a more general discussion of the process of conclusion drawing and verification.

---

<sup>194</sup> For detailed information on how to develop a variety of data displays, see Miles and Huberman (1994).

<sup>195</sup> Berkowitz, 1997, p. 8.

<sup>196</sup> Ibid.

<sup>197</sup> Miles & Huberman, 1994, p. 10.

<sup>198</sup> Berkowitz, 1997, p. 8.

<sup>199</sup> 1994, pp. 262-276.

<sup>200</sup> pp. 262-280.

In the next section we briefly highlight some of the more popular qualitative data analysis software programs and point out some of the advantages and disadvantages of using them.

#### 4.4.4.5 Qualitative Data Analysis Software Programs

Depending on the resources available, the study team may want to invest in qualitative data analysis software, such as the commonly-used proprietary tools NVivo, ATLAS.ti, Ethnograph, HyperRESEARCH, QDA Miner, MAXQDA, Qualrus, Dedoose or one of the open-source programs such as Transana and Coding Analysis Toolkit (CAT; Suter, 2012). These tools can facilitate complex analysis of large quantities of qualitative data in less time than might be required to analyze the data by hand or through a spreadsheet, although significant time nevertheless may be required to set up or tailor the program to the evaluator's needs (See Schutt, 2011, for an example of how qualitative analysis might be conducted using some common software programs). Some of the **features of qualitative analysis software** include:

- Project management tools (e.g. that allow multiple analysts to work with data in different files following consistent analysis protocols established at the outset);
- Communication and memo writing tools (e.g., that allow researchers to incorporate comments and memos in the overall analysis);
- Collaboration tools (e.g., that allow researchers to share their analyses, comments, and memos; some even calculate inter-rater reliability);
- Graphing features (e.g., that allow a researcher to generate visual displays of findings to facilitate understanding of the connections among the different themes or data sources); and
- Query, search, and report functions (e.g., that can search for specific words and phrases and generate reports on their frequency).

As can be seen, these types of software can provide valuable benefits to the qualitative analyst. However, some researchers have raised concerns about the use of software because of issues related to:

- The time required to learn and master the software;
- An increased focus on coding as the primary analysis strategy;
- An increased focus on breadth and volume, rather than depth and meaning; and
- An increased focus on queries and text searches that may detract from the nuance and context in the data.<sup>201</sup>

Of course, no software program can take the place of good training in qualitative analysis, which requires significant creativity, intellectual discipline, and analytic rigor on the part of the analyst. The software can facilitate analysis, but a strong foundation in qualitative research techniques is required to ensure the quality of the conclusions reached by the study team.

In the next section we talk about how to create system to measure fidelity of a project.

---

<sup>201</sup> E.g., St. John & Johnson, 2000.

## 4.5 Measuring Fidelity

An important part of an evaluation is determining whether the project has been carried out with fidelity. Without knowing whether a project has actually been implemented according to plan, it's impossible to know whether the project has been responsible for producing the observed results. Fidelity data are essential to help project staff and evaluators understand what's happening with project implementation. To be able to link project activities to outcomes, for example, it's not enough to know whether an activity has occurred. It's also important to know,

- whether the activity was carried out in the way that was intended (e.g., Did the activity get conducted in the correct timeframe and cover the expected content?),
- whether the right people attended in the right amounts (e.g., Did a high percentage of the target population attend and how often did they attend?), and
- whether the activity resulted in the expected outcomes (e.g., Were there changes in participants' knowledge, attitudes, skills, or behaviors as a result of participating in the activity?).

Fidelity data are formative to the extent that they are used by project staff or program developers to make changes to a project or intervention during the course of a project. When used in a summative sense, fidelity data can offer insights into why a project might not have achieved the expected outcomes (e.g., if there was low fidelity) and signal the feasibility of implementing such a project again, in a different context.

There are **five basic steps to creating a system to measure fidelity**<sup>202</sup>:

1. Identify the “key components” of the project, or those features that are critical for the project to achieve positive results
2. Operationally define and create indicators for each key component (e.g., professional development) included in the project logic model
3. Select data sources and measures for each indicator
4. Establish fidelity thresholds and set scores for “adequate” fidelity
5. Calculate fidelity scores based on observed data

These are discussed in the following sections.

---

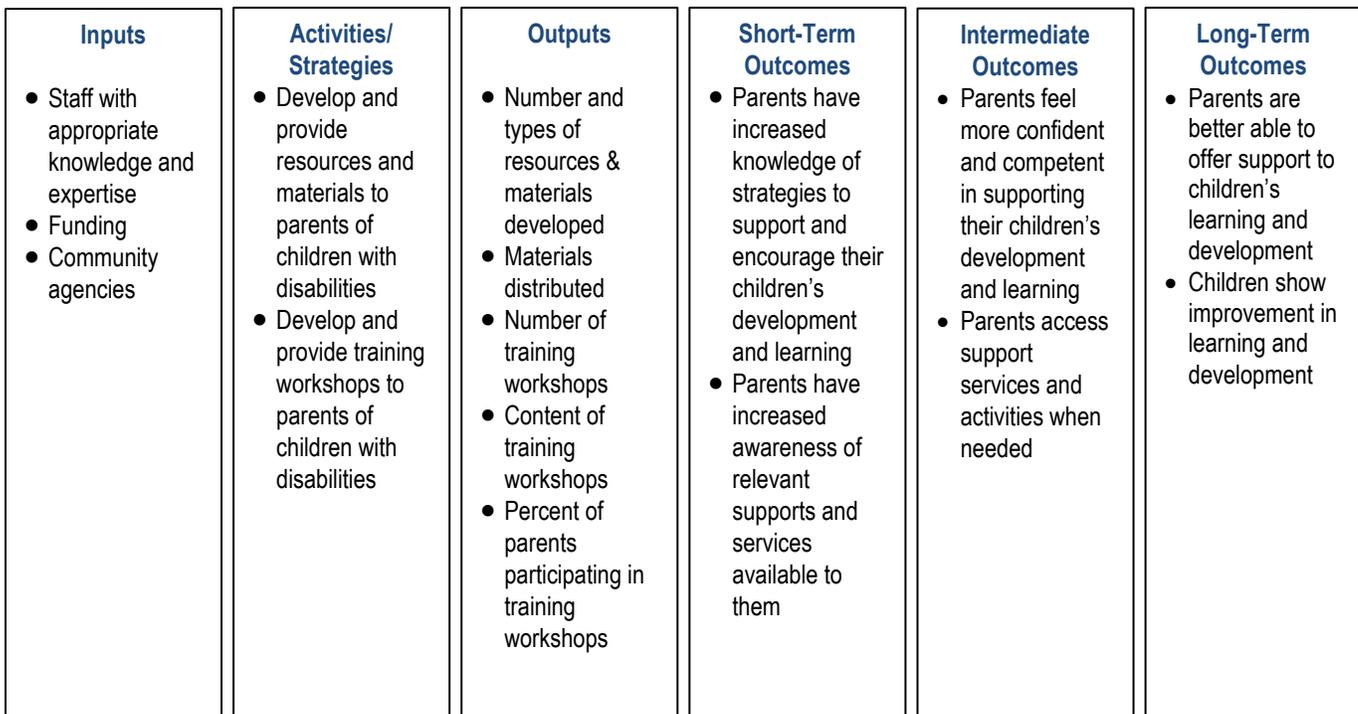
<sup>202</sup> Adapted from Lammert & Feldman, 2015.

### 4.5.1 Identify Key Components

A first step in measuring fidelity is identifying the “**key components**” of a project, or those features that are critical for the project to achieve positive results. The project’s [logic model](#) can serve as a guide to identifying key components, as it clearly details the project’s key structural and process components (including inputs, activities, outputs and outcomes). Some examples of structural components are content, activities, exposure, and dosage (e.g., number and length of sessions), resources (e.g., budget), and staff capacity. Examples of process components are quality and participant responsiveness and participation.<sup>203, 204</sup>

Figure 12 shows part of the hypothetical logic model for a Parent Resource and TA Center that was first presented in Figure 2. We will use this Parent Resource and TA Center as our example throughout this section. Based on this logic model, examples of key components might be: staff competencies, staff facilitation of training workshops, and parent participation in and satisfaction with the workshops.

**Figure 12. Sample Logic Model for a Hypothetical Parent Resource and Technical Assistance Center**



To better illustrate the process of creating a system to measure fidelity, we will focus our example on one specific activity presented in the logic model above: Develop and provide training workshops to parents of children with disabilities. More specifically, we will look at the provision of training to parents to help them support the learning and development of their children with autism.

After determining the key project components, the next step is to create operational definitions for each component (including activities and participation, outputs, and outcomes), and identify indicators for which data can be collected to determine whether each key component was carried out with fidelity.

<sup>203</sup> Century, Rudnick, & Freeman, 2010.

<sup>204</sup> Mowbray, Holter, Teague, & Bybee, 2003.

## 4.5.2 Create Operational Definitions and Identify Indicators

An **operational definition** is concrete, specific, and comprised of **indicators**—specific aspects of the intervention that can be measured quantitatively. Operational definitions should be informed by the project’s [theory of change \(Section 2.2\)](#), [logic model \(Section 2.3\)](#), project plan, and other evidence supporting best practices for effectively implementing the project activities. **The number of indicators should reflect the complexity of the intervention components.** Additionally, it’s best if the indicators are able to differentiate among different levels of fidelity (e.g., low, moderate, high). A key part of this process is determining which indicators can actually be measured. It might be possible to devise a very specific operational definition for a component, but if it’s not possible to feasibly measure that indicator with the available data sources, it doesn’t help for measuring fidelity.

Going back to our example, if project staff experience shows that implementation is more effective if the lead parent educators on staff have a specific educational background or skill set, that information can be used to develop the operational definition of staff competencies. Similarly, the research on best training practices for adult learners that informed the project’s theory of change can help staff to operationalize the key components related to facilitation of the parent training workshops. Similarly, research or experience might indicate that there are certain strategies or activities that should be used to best convey the training content to parents. All of this should inform development of the operational definition and selection of indicators. Table 18 presents a possible operational definition and indicators for our Parent TA Center example.

**Table 18. Possible Operational Definitions and Indicators for the Parent TA Center Example**

Operational Definition & Indicators
<p><b>Operational Definition:</b> For the training provided by the Center (Key Component 1) to be considered implemented with fidelity, (a) the trainings will be delivered by lead parent educators who have the required competencies, (b) training workshops will be delivered on schedule and within the allotted timeframe, and will cover the expected content; and (c) parents will participate in and express satisfaction with the workshops.</p>
<p><b>Indicator 1—Staff Competencies:</b></p> <ul style="list-style-type: none"> <li>• 100% of staff members who serve as lead parent educators must have <ul style="list-style-type: none"> <li>○ a Bachelor’s Degree or an Associate’s Degree in Special Education, Early Childhood Education, or a related field;</li> <li>○ at least one year of experience working with families of children with disabilities; and</li> <li>○ positive evaluations that provide evidence that the staff member understands the content being presented, effectively interacts with parents, and effectively facilitates parent training workshops.</li> </ul> </li> </ul>
<p><b>Indicator 2—Staff Facilitation of Training Workshops:</b></p> <ul style="list-style-type: none"> <li>• 4 training workshops are delivered on schedule (within allotted timeframe)</li> <li>• Content and format of training workshops: <ul style="list-style-type: none"> <li>○ is appropriate for the expressed purpose</li> <li>○ addresses the expressed needs of parents</li> <li>○ encourages parents to participate and share their experiences, and</li> <li>○ incorporates at least three training methods (e.g., lecture, demonstration, technology, games, skill practice, group discussion).</li> </ul> </li> </ul>
<p><b>Indicator C—Parent Participation in and Satisfaction With Training Workshops:</b></p> <ul style="list-style-type: none"> <li>• XX% of eligible parents attend the workshops</li> <li>• XX % of parents in attendance <ul style="list-style-type: none"> <li>○ actively participate in discussions, share experiences, or ask questions;</li> <li>○ complete all activities; and</li> <li>○ are satisfied with the quality, relevance and usefulness of the workshops.</li> </ul> </li> </ul>

If the evaluator is creating the operational definition and fidelity indicators, project staff and/or the intervention developer should review the definitions and indicators to ensure they have captured the important elements of fidelity. Additionally, it's important to keep in mind that the **operational definitions should consider the real-world circumstances that will likely occur during project implementation**. For example, since the operational definition of the staff competencies (Indicator 1) outlined above states that 100% of staff members who serve as lead parent educators must meet all three criteria listed, the project may not ever be able to implement that component with fidelity, since it might not be possible to find sufficient numbers of staff who meet all three criteria. This doesn't mean that the project shouldn't aspire to have all lead parent educators meet the competency criteria, but only that the project staff should carefully consider whether in any given situation a particular operational definition could ever be implemented with fidelity. In this case, it might be better to set the percentage lower, say at 75-80%, and then try to reach the higher bar through recruitment of good candidates. Additionally, consider whether it is possible to assess fidelity of a sample of activities as opposed to assessing fidelity in all of them (e.g., attending 4 workshops to assess fidelity rather than 6).

### 4.5.3 Select Data Sources and Measures

The next step in the process is to identify data sources and measures for each indicator. To do this, consider the research questions guiding the study, find the best data sources for each indicator, and choose whether to use existing instruments or develop new ones. **Use multiple sources of data for each indicator (e.g., surveys and observations) when possible, and consider reliability/validity of data sources and measures.** Table 19 presents possible data sources for each indicator in the Parent TA Center example.

**Table 19. Possible Data Sources for the Parent TA Center Example**

Indicators	Data Sources/ Measures
<p><b>Indicator 1—Staff Competencies:</b></p> <ul style="list-style-type: none"> <li>• 100% of staff members who serve as lead parent educators must have               <ul style="list-style-type: none"> <li>○ a Bachelor's Degree or an Associate's Degree in Special Education, Early Childhood Education, or a related field;</li> <li>○ at least one year of experience working with families of children with disabilities; and</li> <li>○ positive evaluations that provide evidence that the staff member understands the content being presented, effectively interacts with parents, and effectively facilitates parent training workshops.</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• Administrative data on staff credentials and experience</li> <li>• Staff evaluations</li> </ul>
<p><b>Indicator 2—Staff Facilitation of Training Workshops:</b></p> <ul style="list-style-type: none"> <li>• 4 training workshops are delivered on schedule (within allotted timeframe)</li> <li>• Content and format of training workshops:               <ul style="list-style-type: none"> <li>○ is appropriate for the expressed purpose</li> <li>○ addresses the expressed needs of parents</li> <li>○ encourages parents to participate and share their experiences, and</li> <li>○ incorporates at least three training methods (e.g., lecture, demonstration, technology, games, skill practice, group discussion).</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• Needs assessment data</li> <li>• Workshop schedules, agendas and materials</li> <li>• Observation rubric or checklist</li> <li>• Participant feedback surveys</li> </ul>
<p><b>Indicator 3—Parent Participation in and Satisfaction With Training Workshops:</b></p> <ul style="list-style-type: none"> <li>• XX% of eligible parents attend the workshops</li> <li>• XX % of parents in attendance               <ul style="list-style-type: none"> <li>○ actively participate in discussions, share experiences, or ask questions;</li> <li>○ complete all activities; and</li> <li>○ are satisfied with the quality, relevance and usefulness of the workshops</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• Workshop attendance records</li> <li>• Observation rubric or checklist</li> <li>• Participant feedback surveys</li> </ul>

Once data sources have been identified, the process of assessing fidelity involves establishing fidelity thresholds (e.g., low, moderate and high) and setting fidelity scores, or the minimum score that is required to achieve “adequate” fidelity, discussed in the next sections.

#### 4.5.4 Establish Fidelity Thresholds and Set Scores for “Adequate Fidelity”

**Fidelity thresholds** are numeric scores that are used to define different levels of fidelity of a specific indicator. Represented as numeric scales, thresholds quantify the extent to which an indicator was enacted with fidelity. This scale can be dichotomous (e.g., 0 = inadequate fidelity or 1 = adequate fidelity) or it can have a range of fidelity levels (e.g., 0 = low fidelity, 1 = moderate fidelity, and 2 = high fidelity). Further, the project team can define levels of fidelity as low/moderate/high based on the total score across indicators/components or based on the % of indicators/components fully met. **Fidelity scores** are calculated based on the fidelity thresholds and identify the least amount of each component that needs to be present for fidelity to be considered “adequate” for that component. It’s possible to calculate a fidelity score for each key component separately, or to create one score for fidelity across all project components.

Currently, there aren’t hard and fast standards for creating fidelity thresholds or setting fidelity scores. To create thresholds and scores to determine what “adequate fidelity” will look like, it may be possible to get information from an intervention developer (if using a packaged program), prior research, or from project staff or evaluators who have implemented similar interventions previously. The point is not for everyone to get a passing grade on each indicator or component; instead, the framework should be sufficiently sensitive to allow the evaluator to use the data to identify observable differences between “adequate” and “inadequate” levels of fidelity and among different levels of fidelity for different parts of the project.

**To establish a threshold, consider the relative importance of various indicators, the unit or level at which to collect data (e.g., individual, school, or program level), and the range of possible fidelity scores (e.g., 0, 1, and 2).** Fidelity thresholds can be established at multiple levels (e.g., individual, school, and project) and depend in part on what type of indicator is being measured. So, for example, it wouldn’t be appropriate to calculate fidelity at the individual level for something that applies to the school- or project-level, such as whether training workshops were offered to parents by the Parent TA Center. It may be necessary to “roll-up” the scores for the various levels of thresholds, for example, by starting with the threshold for individual-level (e.g., teacher) fidelity, then creating a threshold for school-level fidelity, and, finally, creating an overall threshold for project-level fidelity—based on the scores from the lower levels.

For now, however, we will keep it simple and focus on just one level. Expanding upon our example in Table 19, let’s look at possible fidelity thresholds and the corresponding score to achieve adequate fidelity for Indicator 1 of the project’s key component of training: Staff competencies (Table 20).

**Table 20. Sample Fidelity Thresholds for a Parent Resource and TA Center – Example A**

Indicators	Fidelity Thresholds and Score for Adequate Fidelity
<p>Indicator 1: Staff Competencies:</p> <ul style="list-style-type: none"> <li>• 100% of staff members who serve as lead parent educators must have               <ul style="list-style-type: none"> <li>a. a Bachelor’s Degree or an Associate’s Degree in Special Education, Early Childhood Education, or a related field;</li> <li>b. at least one year of experience working with families of children with disabilities; and</li> <li>c. positive evaluations that provide evidence that the staff member understands the content being presented, effectively interacts with parents, and effectively facilitates parent training workshops.</li> </ul> </li> </ul>	<p><b>2</b> = high fidelity: 100% of lead parent educators on staff meet all three competency criteria (project -level threshold)</p> <p><b>1</b>= moderate fidelity: 75%-99% of lead parent educators on staff meet all three competency criteria (project -level threshold)</p> <p><b>0</b> = low fidelity: Less than 75% of lead parent educators on staff meet all three competency criteria (project-level threshold)</p> <p><b>Score to achieve “adequate fidelity” = 2</b></p>

As with creating operational definitions for fidelity, evaluators should review the fidelity thresholds and scores with project staff and/or the intervention developer to see if the thresholds and scores make sense. For instance, in the example presented in Table 20 above, there are some potential issues with the established thresholds. First, Indicator 1 states that 100% of staff members who serve as lead parent educators must have all three indicators of staff competency listed. This means that there is no real reason to have three levels for the fidelity threshold, since the project either has 100% of staff with all three required competencies or it doesn’t. In this case it might be more appropriate to set a dichotomous threshold (0 or 1), since the only way to achieve fidelity is to get 100% of staff with the required competencies (a score of “2”).

Although it may seem like a good idea to create indicators and set fidelity thresholds in this way (i.e., with very high expectations for fidelity), it may be overly restrictive given the realities of project implementation, may not provide enough information to inform project implementation, and may be very resource-intensive to collect sufficient amounts of data to assess fidelity. For example, in the case of competency indicator (c), if even *one* lead parent educator doesn’t have a positive evaluation that (i) provides evidence that the staff member understands the content being presented, (ii) effectively interacts with parents, *and* (iii) effectively facilitates parent training workshops, then the entire project fails to achieve fidelity on Indicator 1. Additionally, since in this example all of the staff competency indicators are lumped together within the threshold (i.e., 100% of lead parent educators on staff meet *all three* competency criteria), it’s not possible for project staff to identify which staff competency indicator (a, b, or c) was responsible for the lack of fidelity or to determine whether a particular competency indicator was actually important to achieving outcomes.

Table 21, on the next page, presents an alternate approach to creating indicators and setting fidelity thresholds for this example that allows for a bit more variation in actual project implementation and that can potentially be more informative. As can be seen in the table, now the indicator does not include a percentage and the fidelity thresholds are what set the criteria for fidelity in terms of the percentage of staff achieving each competency. Additionally, the score required to achieve “adequate fidelity” is now a “4,” and there are different ways the project can achieve that score. For example, the project may get a “2” on competency indicator (a) and competency indicator (c), and a “0” on competency indicator (b). Or, the project may get a “1” on competency indicator (a), a “2” on competency indicator (b), and a “1” on competency indicator (c).

**When setting thresholds and minimum fidelity scores, it’s important to identify in advance those indicators (if any) that are believed to be absolutely essential to successful project implementation and to set the required fidelity score accordingly.** Going back to our example, if the project team believes it’s essential that lead parent educators have a Bachelor’s or Associates Degree in Special Education, Early Childhood Education, or a related field, then the evaluator

may require a score of “2” on competency indicator (a), but only at least a “1” on competency indicator (b) and competency indicator (c). In this case, the project still must get a score of “4,” but to achieve adequate fidelity overall, the project also has to ensure that it achieves high fidelity on competency indicator (a) and moderate fidelity on the other two competency indicators.

**Table 21. Sample Fidelity Thresholds for a Parent Resource and TA Center – Example B**

Indicators	Fidelity Thresholds
<p>Indicator 1—Staff Competencies:</p> <ul style="list-style-type: none"> <li>• Staff members who serve as lead parent educators have               <ul style="list-style-type: none"> <li>a. a Bachelor’s Degree or an Associate’s Degree in Special Education, Early Childhood Education, or a related field;</li> <li>b. at least one year of experience working with families of children with disabilities; and</li> <li>c. positive evaluations that provide evidence that the staff member understands the content being presented, effectively interacts with parents, and effectively facilitates parent training workshops.</li> </ul> </li> </ul>	<p><b>Competency indicator (a) (project-level threshold)</b>  <b>2</b> = high fidelity: ≥ 90% of lead parent educators on staff meet competency (a)  <b>1</b>= moderate fidelity: 75%-89% of lead parent educators on staff meet competency indicator (a)  <b>0</b> = low fidelity: Less than 75% of lead parent educators on staff meet competency indicator (a)</p> <p><b>Competency indicator (b) (project-level threshold)</b>  <b>2</b> = high fidelity: ≥ 90% of lead parent educators on staff meet competency indicator (b)  <b>1</b>= moderate fidelity: 75%-89% of lead parent educators on staff meet competency indicator (b)  <b>0</b> = low fidelity: Less than 75% of lead parent educators on staff meet competency indicator (b)</p> <p><b>Competency indicator (c) (project-level threshold)</b>  <b>2</b> = high fidelity: ≥ 75% of lead parent educators on staff meet competency indicator (c)  <b>1</b> = moderate fidelity: 60%-75% of lead parent educators on staff meet competency indicator (c)  <b>0</b> = low fidelity: Less than 60% of lead parent educators on staff meet competency indicator (c)</p> <p><b>Score to achieve “adequate fidelity” = 4</b></p>

Other considerations for setting thresholds include deciding whether to weight indicators based on their relative importance to the component (or project) and considering how to handle indicators that are measured with different metrics. It is beyond the scope of this Toolkit, however, to delve into these topics.

The operational definitions, indicators, fidelity thresholds and fidelity scores can all be put into a fidelity matrix that includes all of the information about the fidelity system. [Appendix A.7](#) includes a fidelity matrix template.

#### 4.5.5 Calculate Fidelity Based on Observed Data

Once the fidelity measurement system has been established, data can be collected through the other [data collection activities \(Section 2.5.4\)](#) that will be carried out as part of the evaluation, preferably on a regular basis so as to inform ongoing implementation. Clearly, it’s important to think about fidelity during the [evaluation planning \(Section 2\)](#) phase to ensure that the right data will be collected to enable assessment of fidelity—both of the on-going project activities and of outcome achievement.

Fidelity data can be useful to project staff who want to know whether the project is being carried out as planned. Formative evaluation data collection activities can gather these data, which can then be used to calculate fidelity of specific project components, such as training and professional development, development of online resources for key

stakeholder groups, and provision of TA. Using the fidelity matrix as a guide, evaluators might calculate fidelity each year of the project (generally after allowing for an appropriate period to set up the project) and provide the results to project staff to inform improvements. For example, by highlighting certain indicators or components that a majority of people are struggling to implement, fidelity ratings can help to identify areas where changes might be needed. Similarly, as evaluators go through the process of calculating fidelity ratings, they might identify areas that weren't originally included in the project theory of change or logic model, but that are nevertheless important to fidelity. Similarly, the data might point to components that were part of the original model, but that turn out to be less important on the ground.

Fidelity data also can help project staff to understand program outcomes and to identify factors that are necessary to replicate program successes. Here, the expectation is that project staff know what level of fidelity is needed for the project to achieve its expected outcomes. Obviously, this may be easier for some projects than for others. Continuing with the example of the Parent Resource TA Center, the project staff may not have a clear idea of how many training workshops are necessary for parents to understand how best to support the learning and development of their young children with autism. Collecting and using information on the fidelity to the project plan thus allows project staff to know whether their initial estimates for the amount and type of training and support needed were correct, and helps to provide important information about the way the project has contributed to observed outcomes. Table 22, on the next page, presents some possible scenarios.

**Table 22. Relationship between Fidelity and Achievement of Outcomes**

**Expected Short-term Outcome: Parents have increased knowledge of strategies to support and encourage their children’s development and learning**

Level of Fidelity	Outcome	Conclusion Related to Project Influence on Outcomes
<ul style="list-style-type: none"> <li>• High fidelity to training—number of workshops and content</li> </ul>	<ul style="list-style-type: none"> <li>• Short-term outcome fully achieved: Parents demonstrate full understanding of strategies to support and encourage their children’s learning and development.</li> <li>• Intermediate outcome fully achieved: Parents report feeling more confident and competent in supporting their child’s development and learning.</li> </ul>	<ul style="list-style-type: none"> <li>• There is a good likelihood that the project contributed to the achievement of the expected short-term and intermediate outcomes.</li> </ul>
<ul style="list-style-type: none"> <li>• Low fidelity to number of training workshops</li> <li>• High fidelity to training content</li> </ul>	<ul style="list-style-type: none"> <li>• Short-term outcome partially achieved: Parents show basic understanding of strategies to support and encourage their child’s learning and development.</li> <li>• Intermediate outcome not achieved: Parents do not yet feel confident or competent to support their child’s learning and development.</li> </ul>	<ul style="list-style-type: none"> <li>• It is likely that the high fidelity to the training/coaching content has contributed to the partial achievement of the expected short-term outcome.</li> <li>• It is possible that better fidelity to the number of training/coaching sessions would help the project to fully achieve the expected short-term outcome and contribute to the achievement of the intermediate outcome.</li> </ul>
<ul style="list-style-type: none"> <li>• High fidelity to number of training workshops</li> <li>• Low fidelity to training/coaching content</li> </ul>	<ul style="list-style-type: none"> <li>• Short-term outcome partially achieved: Parents show limited understanding of strategies to support and encourage their child’s learning and development.</li> <li>• Intermediate outcome not achieved: Parents do not yet feel confident or competent to support their child’s learning and development.</li> </ul>	<ul style="list-style-type: none"> <li>• It is possible that the high fidelity to the number of training/coaching sessions has contributed to the partial achievement of the expected short-term outcome.<sup>205</sup></li> <li>• It is likely that better fidelity to the training/coaching content would help the project to fully achieve the expected short-term outcome and contribute to the achievement of the intermediate outcome.</li> </ul>
<ul style="list-style-type: none"> <li>• Low fidelity to training—number of workshops and content</li> </ul>	<ul style="list-style-type: none"> <li>• Short-term outcome not achieved: Parents do not understand strategies to support and encourage their child’s learning and development.</li> </ul>	<ul style="list-style-type: none"> <li>• The low fidelity to the number and content of training/coaching sessions are likely responsible for the project not achieving the expected short-term outcome.</li> </ul>

Note: We assume for the purposes of this discussion that parents have not received training in this topic from another organization. We recommend that project staff collect data on this whenever possible.

For other, more clearly-defined, projects it is easier to determine how much of the intervention is needed to achieve the expected outcomes. A project implementing Positive Behavioral Intervention and Support (PBIS), for example, has clear guidelines on what types of activities need to be carried out, when, and by whom. In addition, PBIS has a research base that demonstrates the effectiveness of the program model, so project staff can have good expectations that if they implement a PBIS project with fidelity to the model—barring any unforeseen barriers to outcome achievement—they will achieve their expected outcomes.

<sup>205</sup> We say that it is only “possible” that the high fidelity to the number of training sessions has contributed to the partial achievement of the outcome, since the low fidelity of the training content means that the sessions were not necessarily providing the necessary content to actually be sure that the training contributed to the achievement of the outcome.

### 4.5.6 Making Changes to the Fidelity System

As data are collected, it's possible to use the data determine whether the assumptions about the importance of a specific indicator are correct. For instance, continuing with the example presented in Table 21, if the data show that the project failed to achieve fidelity on indicator (a), yet the outcomes were still achieved as expected, then project staff may reconsider whether indicator (a) really was all that important. On the other hand, if the project fails to produce the expected outcomes but has achieved high fidelity on indicator (a) and indicator (b)—it may be that the low fidelity on indicator (c) has something to do with the unexpected results (another reason why it's a good idea to create indicators and thresholds that allow for breaking down fidelity into different parts). Based on the data collected, project staff and evaluators may want to revise the indicators for particular components, either by eliminating indicators or adding new ones.

Another way to use fidelity data is to revise the fidelity thresholds and scores for adequate fidelity based on what is being seen on the ground. Table 23 presents an example of a mismatch between the fidelity threshold and the data that have been collected for our example.

**Table 23. Using Data to Modify Fidelity Thresholds**

Initial Fidelity Threshold	Observed Data	Revised Fidelity Threshold
<p><b>Indicator 1(c)</b>            2 = high fidelity: ≥ 90% of lead parent educators on staff have <b>at least one year</b> of experience working with families of children with disabilities</p> <p>1= moderate fidelity: 75%-89% of lead parent educators on staff have <b>at least one year</b> of experience working with families of children with disabilities</p> <p>0 = low fidelity: Less than 75% of lead parent educators on staff have <b>at least one year</b> of experience working with families of children with disabilities;</p>	<ul style="list-style-type: none"> <li>• 80% of candidates for the lead parent educator position had <b>at least 6 months</b> of experience working with families of children with disabilities. The remaining candidates had fewer than 6 months of experience.</li> <li>• Project staff hired the candidates with less experience (rather than leaving the positions empty) and formative evaluation data show that parents feel that the training the lead parent educators provide is high-quality, relevant, and useful</li> </ul>	<p>2 = high fidelity: ≥ 90% of lead parent educators on staff have at least <b>6 months</b> of experience working with families of children with disabilities</p> <p>1= moderate fidelity: 75%-89% of lead parent educators on staff have at least <b>6 months</b> of experience working with families of children with disabilities</p> <p>0 = low fidelity: Less than 75% of lead parent educators on staff have at least <b>6 months</b> of experience working with families of children with disabilities;</p>

As you can see in Table 23, based on the observed data, the project staff revised the initial fidelity thresholds downward so that lead parent educators were only required to have at least 6 months of experience working with families of children with disabilities. This revision better reflects the actual implementation context and, being based on data, the new fidelity system still can provide useful data to project staff and evaluators.

For more information on implementation fidelity in general and on measuring fidelity, see Century, Rudnick & Freedman (2010); Fixsen Fixsen, Naoom, Blase, Friedman, & Wallace (2005); Hulleman & Cordray (2009; Meyers & Brand (2015); and Nelson, Cordray, Hulleman, Darrow, & Sommer (2012).

# Appendix

# Appendix A. Worksheets/Templates

## A.1. Evaluation Cost Consideration Worksheet

**Instructions:** Check off each evaluation item you will need then calculate your score for each row by summing the items across all three columns (e.g., 1 point for items in column one, 2 points for column two, and 3 points for column three), and then calculate your total score by adding the row totals. The scores ranges at the end of the worksheet give you an idea of the relative cost of the evaluation, however the actual cost will depend on factors such as labor rates and travel costs. Note: In some rows (e.g., Interview Mode) you might pick no items or only one item across all three columns, while in other rows you might pick multiple items in multiple columns.

Evaluation Element	Low Cost (1 point/item)	Moderate Cost (2 points/item)	High Cost (3 points/item)	Score
<b>Evaluation Design Elements</b>				
<b>Focus of Formative Study</b>	<ul style="list-style-type: none"> <li>— Participant satisfaction</li> <li>— Project implementation</li> </ul>	<ul style="list-style-type: none"> <li>— Outputs (e.g., satisfaction, quality, relevance)</li> <li>— Implementation fidelity (key components, activities, outputs, possibly some direct outcomes)</li> </ul>	<ul style="list-style-type: none"> <li>— Intervention fidelity that includes mediators, intermediate outcomes</li> </ul>	_____
<b>Focus of Summative Study</b>	<ul style="list-style-type: none"> <li>— Changes in participant satisfaction</li> <li>— Changes in existing data (e.g., student scores on state tests)</li> </ul>	<ul style="list-style-type: none"> <li>— Direct outcomes</li> </ul>	<ul style="list-style-type: none"> <li>— Intermediate/long-term outcomes</li> <li>— Comparative outcomes (e.g., treatment vs. control groups)</li> <li>— Causal attribution</li> </ul>	_____
<b>Evaluation Study Design</b>	<ul style="list-style-type: none"> <li>— Non-experimental (descriptive study, basic qualitative methods)</li> </ul>	<ul style="list-style-type: none"> <li>— Non-experimental (case studies, advanced qualitative methods)</li> <li>— Simple quasi-experiment (QED, e.g., basic comparison study)</li> <li>— Single-case design (SCD; reversal design)</li> </ul>	<ul style="list-style-type: none"> <li>— Complex QED (e.g., with matching, multiple comparison groups)</li> <li>— Randomized controlled trial (RCT)</li> <li>— Multi-site or cluster RCT</li> <li>— SCD (multiple baseline, alternating treatment design)</li> </ul>	_____
<b># of Participants/Sites, Sampling</b>	<ul style="list-style-type: none"> <li>— Small target population</li> <li>— 1-2 sites</li> <li>— Simple sampling plan (e.g., purposive, simple random)</li> </ul>	<ul style="list-style-type: none"> <li>— Moderate-size target population</li> <li>— 3-5 sites</li> <li>— Somewhat complex sampling plan (e.g., stratified)</li> </ul>	<ul style="list-style-type: none"> <li>— Large target population</li> <li>— &gt;5 sites</li> <li>— Highly complex sampling plan (e.g. stratified, clustered, weighted)</li> </ul>	_____
<b>Data Collection Elements</b>				
<b>Document review</b>	<ul style="list-style-type: none"> <li>— Limited document search</li> <li>— Basic document summaries</li> </ul>	<ul style="list-style-type: none"> <li>— Extensive document search</li> <li>— Detailed document summaries</li> <li>— Limited document synthesis</li> <li>— Limited qualitative analysis of documents</li> </ul>	<ul style="list-style-type: none"> <li>— Extensive document synthesis</li> <li>— Extensive qualitative analysis of documents</li> </ul>	_____
<b>Survey (existing or new)</b>	<ul style="list-style-type: none"> <li>— Existing (available, no changes needed)</li> </ul>	<ul style="list-style-type: none"> <li>— Existing (fee to use)</li> <li>— Existing (some changes needed)</li> <li>— New survey (with limited pilot testing &amp; no validation study)</li> </ul>	<ul style="list-style-type: none"> <li>— New survey (with extensive pilot testing)</li> <li>— New survey (with validation)</li> </ul>	_____
<b>Survey mode of administration</b>	<ul style="list-style-type: none"> <li>— Simple online (e.g., basic Survey Monkey survey)</li> </ul>	<ul style="list-style-type: none"> <li>— Customized online survey (e.g. customized Survey Monkey with skip patterns)</li> <li>— Self-administered paper-and-pencil survey (few respondents)</li> <li>— Telephone survey (few respondents, brief responses)</li> </ul>	<ul style="list-style-type: none"> <li>— Customized online survey with integrated data management (e.g., survey with highly complex skip patterns &amp; linkages to data management system)</li> <li>— Self-administered paper-and-pencil survey (many respondents)</li> <li>— Telephone survey (many respondents, lengthy responses)</li> </ul>	_____

Evaluation Element	Low Cost (1 point/item)	Moderate Cost (2 points/item)	High Cost (3 points/item)	Score
Survey non-response follow-up	— Limited email follow-up	— Extensive email follow-up — Limited telephone follow-up	— Extensive telephone follow-up — Mail follow-up with reminder cards	_____
Interview mode	— Online	— Telephone	— Face-to-face <sup>a</sup>	_____
Interview type	— Structured (i.e., asking specific, close-ended questions)	— Semi-structured (i.e., asking some close-ended & some open-ended questions)	— Unstructured (i.e., asking open-ended questions, with focus potentially varying by respondent)	_____
Interview data capture	— Interviewer takes notes during interview	— Interview recorded & transcribed	— Note-taker present at interview	_____
Observation location	— Local <sup>a</sup>	— Driving distance (overnight) <sup>a</sup>	— Long-distance (air travel required) <sup>a</sup>	_____
Observation protocol	— Checklist (i.e., specific activities or behaviors to observe; limited training required)	— Guided/structured protocol (i.e., general categories of activities or behaviors to observe; some content knowledge and training required)	— Unstructured protocol (i.e., open-ended with focus varying by site; deep content knowledge and/or extensive training required)	_____
Assessments	— Existing (conducted at no cost to evaluation)	— Administer small scale pre-post assessments (with no specialized credential or training required) <sup>b</sup>	— Administer large scale pre-post assessments <sup>b</sup> — Administer repeated assessments <sup>b</sup> — Administer assessments (with specialized credential or training required)	_____
Data Collection Frequency & Duration	— Limited frequency (1 time per year or less)	— Moderate frequency (2 or 3 times per year) — Multi-year, but not annual	— Frequent data collection (4 or more times per year) — Annual, or longitudinal data collection	_____
<b>Data Management Elements</b>				
Data management software & hardware	— New software required (low cost) — New hardware required (low cost)	— New software required (moderate cost) — New hardware required (moderate cost)	— New software required (high cost) — New software required (high cost)	_____
Data control & cleaning	— Limited need for data quality control (i.e., multiple choice items; data collected electronically, etc.) — Limited need for data cleaning (i.e., few duplicate records & outliers, little need for coding, etc.)	— Moderate need for data quality control (i.e., multiple choice with some write-in; some field scoring of assessments; some missing data, etc.) — Moderate need for data cleaning (i.e., moderate number of duplicate records & outliers, need to recode some data, etc.)	— Extensive need for data quality control (i.e., open-ended questions; field scoring of assessments; need to merge & reconcile diverse databases; extensive missing data, etc.) — Extensive need for data cleaning (i.e., extensive duplicates & outliers, extensive recoding needed, etc.)	_____
Data Entry	— Automated data entry (e.g., online survey)	— Data entry mostly automated, with some need for hand entry	— Data entry entirely by hand	_____
Database	— Existing database	— Create new database with limited functionality/data sharing	— Create new database with multi-user functionality/data sharing	_____
<b>Data Analysis Elements</b>				
Type of analysis	— Basic descriptive quantitative analysis (e.g., frequencies, t-tests, chi-square tests, ANOVA) — Very limited qualitative analysis	— Intermediate quantitative analysis (e.g., regression, ANCOVA) — Somewhat limited qualitative analysis	— Advanced quantitative analysis (e.g., HLM, SEM) — Extensive qualitative analysis	_____
Data analysis software	— New software required (low cost)	— New software required (moderate cost)	— New software required (high cost)	_____

Evaluation Element	Low Cost (1 point/item)	Moderate Cost (2 points/item)	High Cost (3 points/item)	Score
<b>Technical expertise for analysis</b>	— Existing staff have all needed expertise	— External consultant needed to conduct some analysis or train existing staff	— External consultant needed to conduct most or all analysis or to provide extensive staff training & support	_____
<b>Reporting Elements</b>				
<b>Reporting frequency</b>	— Annual report only	— Interim & annual reports	— Monthly, interim, & annual reports	_____
<b>Types of presentations/audiences</b>	— Oral evaluation updates — Presentations to project staff only	— Policy briefs (limited audiences) — Presentations to project staff & few stakeholders	— Policy briefs (multiple audiences) — Presentations to project staff & multiple stakeholders	_____
<b>Low Cost (&lt;36 points); Moderate Cost (36-79 points); High Cost (≥80 points)</b>			<b>Total Score</b>	_____

Notes: a. Cost = frequency x travel cost; b. Cost = frequency x number of assessment instruments



## A.2. Checklist for Constructing Outcomes

	Yes	No	Somewhat
1. Is the outcome informed by a need or problem statement?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
a. If yes, does the outcome reflect the specific nature of the need or problem?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
b. If no, has the same or a similar outcome been used in related programs or research?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
c. If no, are you using a recommended or required program outcome?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2. Is the outcome relevant and of interest for program stakeholders and decision-makers?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
3. Is the outcome supported in related research or in evidence-based practices?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4. Has or will a Program or Advisory Team reviewed the outcome for clarity and relevance?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
5. Have you identified the data variables that will be needed to generate a finding for the outcome statement?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
6. If using data collected by schools, districts, or the state, do you have access to the data variables that are needed to generate a finding for the outcome statement?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
a. If yes, are sufficient data available to generate a dataset that responds to the outcome in full?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
7. If using data collected by the project team (i.e., a new data collection), have you identified all of the steps and personnel that will be necessary for data collection and entry?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
8. If using data collected by the project team, do data collection methods and instruments rely on validated and reliable tools and techniques, for the outcomes or constructs of interest?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
9. Do you or will you have the specific expertise or assistance you will need for data entry and/or coding?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
a. If no, do you or will you have the resources to obtain the expertise or assistance you will need?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
10. Do you or will you have the specific expertise or assistance you will need for data analysis?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
a. If no, do you or will you have the resources to obtain the expertise or assistance you will need?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
11. Are there sufficient time and resources to collect and aggregate data for			
a. Direct outcomes?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
b. Intermediate outcomes?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
c. Long-term outcomes?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
12. Do you have required permissions or Memoranda of Understanding (MOU) to collect data?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
a. If no, do you have a process for obtaining permission or MOUs to collect data?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>



## A.3. CIPP Logic Model Template

A logic model provides a starting point for developing an effective evaluation plan. CIPP's goal is to assist projects in developing logic models that are precise and that include features and content that give the models utility for evaluation purposes. The logic models serve as both ends and means. As ends, they are stand-alone representations of projects that provide an overall visual summary. They also provide descriptive information that can be used to catalogue or compare features across multiple projects, if needed. As means, the use of logic models is central to (1) defining outcomes that are meaningfully connected to project activities and (2) supporting evaluations so that the process will improve projects' overall performance.

### CIPP's Logic Model Scheme

A project logic model portrays a project's overall plan. It serves to clarify the relationships among a project's goals, activities, and outputs and to lay out the connections between them and the project's expected outcomes. Therefore, a project logic model depicts a program theory and accompanying hypotheses, highlighting (1) the resources or inputs dedicated to an effort, (2) the planned activities to be carried out with those resources, and (3) the specific outputs and outcomes the activities will generate. Evaluation can then be viewed as a test of the logic model's hypotheses, and a logic model can be used by evaluators and the grantee to refine and guide data collection and analysis for assessing both process and performance.

We find it helpful to begin with a summary chart that contains the information that will populate the logic model. The chart outlines the OSEP priority, assumptions, external factors/context, and inputs. It then displays, in table format from left to right, the project's goals and objectives, strategies/activities, outputs, and outcomes. From the chart, a logic model is prepared. The logic model is less comprehensive than the chart in its content but it uses lines and arrows to connect specific project elements and provides a dynamic display. Both the chart and the logic model are continuously updated as the content of specific elements changes, such as when planned activities are revised or when unintended outcomes occur. The logic model will also change as the relationships among the components develop over time, mostly likely by becoming more complex and interactive. For CIPP, we use the following definitions of the logic model components:

**Goals/Objectives** – The *goals* capture the overarching purposes of the project. Goals make clear the anticipated impact on systems or individuals. Goals imply gaps or deficits that will be remedied when the project produces its long-term outcomes. *Objectives*, if used in a logic model, are targeted sub-goals.

**Strategies/Activities** – *Strategies* are the broad approaches to addressing the goals. They include multiple activities. *Activities*, which may or may not be listed in the logic model, are the specific actions funded by the grant or supported by other resources under the umbrella of the project.

**Outputs** – *Outputs* are the direct results of the project activities, including project products and programs. Most outputs will be quantifiable, including tallies of the number of products and programs or counts of the customer contacts with those products and programs.

- **Short-term/Intermediate Outcomes**—Short-term outcomes are what customers do or become as a result of outputs. Usually, short-term outcomes are changes in knowledge or skills acquired through project outputs. Intermediate outcomes result either directly from outputs or indirectly through short-term outcomes. Often, intermediate outcomes are changes in the behavior or practices of persons touched by the project. They generally come later in time than short-term outcomes and often represent a step between short-term outcomes and long-term outcomes.
- **Long-term Outcomes**—Long-term outcomes are the broadest project outcomes and follow logically from the short-term and intermediate outcomes. They are the results that fulfill the project's goals. However, they aren't always able to be assessed during the evaluation due to time or resource constraints. Outputs, short-

term outcomes, and intermediate outcomes all contribute to the achievement of the long-term outcomes. Although the long-term outcomes represent fulfillment of the purpose of the project, they may or may not represent the achievement of a desired larger project impact. That is, the project may have an anticipated impact that is beyond the immediate scope of the project, either temporally or conceptually, and thus beyond the scope of the logic model.

### **The Summary Chart**

The summary chart contains all the information that will populate the logic model, plus additional details about strategies, activities, and outputs. The chart begins with the OSEP priority and then states assumptions about how and why the project will be successful. External factors and context provide a brief description of the environment in which the project will be operating. Inputs are specific resources available to the project. The table itself displays, from left to right, the project's goals and objectives, strategies/activities, outputs, and outcomes. Note that strategies/activities are aligned with outputs but that goals and outcomes may cut across multiple strategies/activities.

### **The Logic Model**

The content of the logic model is taken entirely from the summary chart, but the logic model content is condensed to fit into the flow chart format. In the logic model, lines and arrows are used to depict the temporal and causal connections among the various project elements. Not surprisingly, multiple lines or arrows come to or from most of the boxes, indicating the complexity of the relationships that are expected. Also depicted are the anticipated results in the form of short-term, intermediate, and long-term outcomes. The outcomes are themselves interconnected. Thus, short-term outcomes, as well as outputs, lead to the higher level, more distal outcomes.

## SUMMARY CHART

**OSEP Priority:** OSEP’s stated purpose for the [project] is to . . . .

**Assumptions:** The [project] is managed by [institution], which has a long history of . . . . This reputation along with the renewed emphasis by OSEP on . . . will create interest in . . . . Specifically, stakeholders focused on . . . will find the [project’s] products and services essential in . . . .

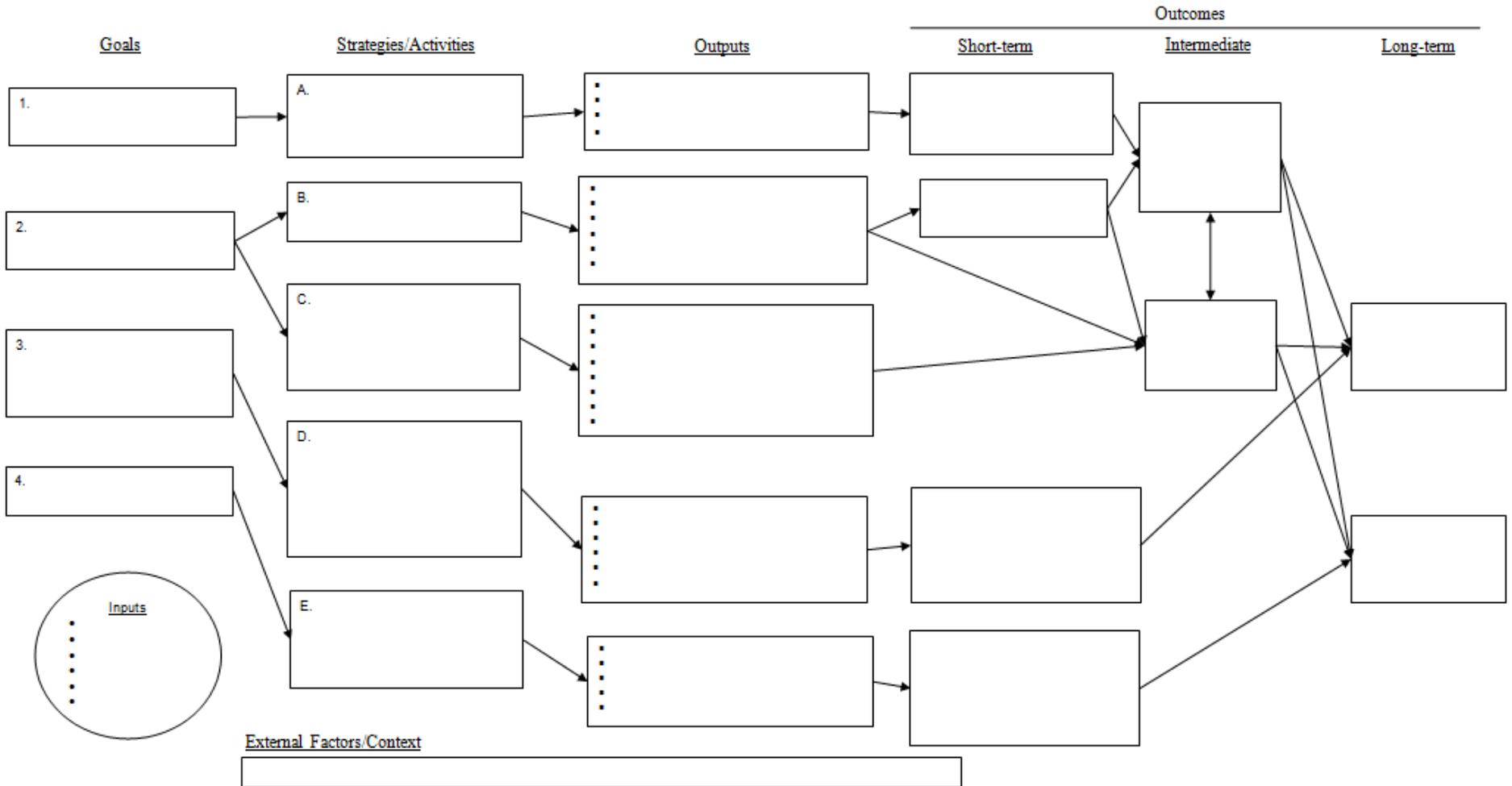
**External Factors/Context:** Other federal initiatives; OSEP policy environment; [institution’s] accumulated experience and visibility....

**Inputs:** OSEP funding, experienced project staff, lessons learned from past experience, research-based policy and practices, ...

**Table A.3.1. Template: Summary Chart Depicting Goals/Objectives, Strategies and Activities, Outputs and Outcomes**

Goals/Objectives	Strategies/Activities	Outputs	Outcomes
1.	A. (1)	•	<b>Short-term:</b> • • • •  <b>Intermediate:</b> • •  <b>Long-term:</b> • •
2.	(2)	•	
3.	(3)		
4.	B. (4)	•	
	(5)	•	
	(6)	•	
	(7)	•	
	C. (8)	•	
	(9)	•	

## LOGIC MODEL TEMPLATE



## A.4. CIPP Summative Evaluation Plan Template

**Note:** *Instructions/Guidelines for completing the template are in italics.*

*Introductory paragraph—one paragraph overview of the project.*

### **Need, Goals, and Activities**

*Paragraph stating the need that the selected project was established to address. Can be followed by a textbox that succinctly states that need.*

**NEED:**

*Outline showing the goals, strategies, and activities.*

#### **GOAL #1:**

STRATEGY:

ACTIVITY:

ACTIVITY:

STRATEGY:

ACTIVITY:

#### **GOAL #2:**

STRATEGY:

ACTIVITY:

STRATEGY:

ACTIVITY:

ACTIVITY:

#### **GOAL #3:**

STRATEGY:

### **[Project] Program Theory**

*Paragraph providing a general intro to the program/project theory of change.*

*Paragraph introducing Table A.4.1, which presents the alignment of project goals with strategies/ activities, outputs, and outcomes. The table should also include the top-level evaluation question associated with each goal.*

**Table A.4.1. Evaluation questions aligned with [project] goals, strategies/activities, outputs, and outcomes.** *[Mostly comes from the Summary Chart prepared for the project logic model, with the addition of the highest level summative evaluation questions.]*

Goals	Strategies/Activities	Outputs	Outcomes	Evaluation Questions
1.	<ul style="list-style-type: none"> <li>▪</li> <li>□</li> <li>□</li> </ul>	<ul style="list-style-type: none"> <li>-</li> <li>-</li> <li>-</li> </ul>	<ul style="list-style-type: none"> <li>•</li> <li>•</li> </ul>	A.
2.	<ul style="list-style-type: none"> <li>▪</li> <li>□</li> </ul>	<ul style="list-style-type: none"> <li>-</li> <li>-</li> <li>-</li> </ul>	<ul style="list-style-type: none"> <li>•</li> <li>•</li> </ul>	B.

*Paragraph introducing the logic model as a graphical representation of the content of Table A.4.1.*

**Figure A.4.1. [Project] logic model.**

*Insert completed logic model here.*

## Evaluation Design

*Restatement of top-level evaluation questions from Table A.4.1, by outcome level.*

### Evaluation Questions Related to Short-term Outcomes

- A.
- B.
- C.

### Evaluation Questions Related to Intermediate Outcomes

- B.
- C.

### Evaluation Questions Related to Long-Term Outcomes

- D.

*Paragraphs presenting an overview of the evaluation approach that will be applied to addressing each question.*

## Data Collection

*Paragraphs presenting an overview of data collection plans, starting with data collection related to the project's formative evaluation and continuing with the summative evaluation activities.*

*Introduction to Table A.4.2. Table A.4.2 presents the evaluation questions in relation to outcomes, data collection activities, and the data collection instruments that are specific to types of respondents.*

*Discuss status of instruments—already in use, under development, planned by the project, needing to be developed or modified.*

*Discuss need for baseline, whether baseline data already exist, and how baseline data collection might work.*

**Table A.4.2. Evaluation questions in relation to outcomes, data collection activities, and instruments.**

Evaluation Question	Outcome	Data Collection Instrument	Type of Data Collection
A.	•		
	•		
B.	•		
	•		

**Sampling**

*Discuss need for sampling, if any. Discuss the likelihood of a need for a sampling plan and some of the specifics if known. Include a sampling plan, if available.*

**Analysis Approach**

**Descriptive Analyses**

*Discuss need for descriptive analyses—that is, where they will be used. Present specific techniques corresponding to the specific evaluation questions and data collections. Identify specific “studies” if they are sufficiently discrete.*

**Statistical Analyses**

*Discuss need for statistical analyses—that is, where they will be used. Present specific techniques corresponding to the specific evaluation questions and data collections. Identify specific “studies” if they are sufficiently discrete.*

*Point out where pre-post data will be used and why only posttest data may be available in some cases. Name the statistical tests, if known. Discuss limitations, or at least the major limitations.*

## Data Sources

*Include an appendix with a table of instruments and the actual instruments, if available.*

### Data Sources for Descriptive Analyses

*List specific instruments or data bases to be used to collect data or as sources of data.*

### Data Sources for Statistical Analyses

*List specific instruments or data bases to be used to collect data or as sources of data.*

## Data Collection Schedule

### Ongoing Assessments

*Describe the timeline for any assessments or data collections that are ongoing and relevant to the evaluation. These will likely be the data collections for the formative evaluation or some type of data base that is regularly updated.*

### One-time Data Collections

*Describe the timeline for any assessments or data collections that are one time and relevant to the evaluation.*

### Pre-Post Data Collections

*Describe the timeline for pre-post—this may be experimental, quasi-experimental, time series or some other rigorous design, or it may be just a repeat of data collections—assessments or data collections that are relevant to the evaluation.*

*Prepare a simple table like Table A.4.3 or a Gantt chart, as illustrated in [Section 2.5.4.6](#). It also may be helpful to prepare a table such as Table A.4.4, which provides information on the status of data collection instruments and activities.*

**Table A.4.3. Evaluation data collection schedule.**

Evaluation Activity	First Data Collection	Additional Data Collections

**Table A.4.4. [Project] Instrument Information Table**

Evaluation Question	Data Source	Possible instrument/ protocol	What is the status of the instrument/ protocol? E=exists UD=under development TBD= to be developed	Have any data been collected with this instrument (if it exists)? If so, when?	If data were collected, was sampling used?	When are future data collection(s) planned for this instrument/ protocol?	Will data collection require sampling?
A.							
B.							

## A.5. Sample Evaluation Plan: Graduate Performance and Student Outcomes

**Note:** *Instructions/Guidelines for completing a section are in italics. Examples are provided in red.*

### Need, Goals, Activities, Outcomes, and Evaluation Questions

*Instructions: Prepare a sentence or a brief paragraph stating the need that the project was established to address. The need is related to the project purpose but is stated in terms of the deficit that the project will fill.*

**Example:** Recent data show that schools and other service providers have a need for fully qualified speech therapists who can work with bilingual secondary students with disabilities. These therapists must be proficient in evidence-based practices. The need for therapists from diverse racial/ethnic backgrounds and therapists with disabilities is especially strong.

*Fill in Table A.5.1, using however many rows you require. Table A.5.1 aligns project goals with strategies/activities, outputs, and outcomes. The table should also include the top-level summative evaluation questions associated with each goal. Definitions of the elements of the table are provided below, while examples of each for one goal are provided within the table. Please note that the examples provided aren't intended to be comprehensive or cohesive.*

**Goals/Objectives** – The goals capture the overarching purposes of the project. Goals make clear the anticipated impact on systems or individuals. Goals imply gaps or deficits that will be remedied when the project produces its long-term outcomes. Objectives, if used, are targeted sub-goals.

**Strategies/Activities** – Strategies are the broad approaches to addressing the goals. They include multiple activities. Activities are the specific actions funded by the grant or supported by other resources under the umbrella of the project. Although listing activities under broader strategies is preferred, activities alone may be listed.

**Outputs** – Outputs are the direct results of the project activities, including project products and programs. Most outputs will be quantifiable. They include tallies of the number of products and programs or counts of the customer contacts with those products and programs.

**Short-term/Intermediate Outcomes**—Short-term outcomes are what customers do or become as a result of outputs. Usually, short-term outcomes are changes in knowledge or skills acquired through project outputs. Intermediate outcomes result either directly from outputs or indirectly through short-term outcomes. Often, intermediate outcomes are changes in the behavior or practices of persons touched by the project. They generally come later in time than short-term outcomes and often represent a step between short-term outcomes and long-term outcomes.

**Long-term Outcomes**—Long-term outcomes are the broadest project outcomes and follow logically from the short-term and intermediate outcomes. They are the results that fulfill the project's goals. However, they aren't always able to be assessed during the evaluation due to time or resource constraints. Outputs, short-term outcomes, and intermediate outcomes all contribute to the achievement of the long-term outcomes. Although the long-term outcomes represent fulfillment of the purpose of the project, they may or may not represent the achievement of a desired larger project impact. That is, the project may have an anticipated impact that is beyond the immediate scope of the project, either temporally or conceptually, and thus beyond the scope of the logic model.

**Evaluation Questions** – Evaluation questions frame the way data will be summarized or analyzed to address the overarching issue of whether the project's goals have been addressed. The data may be collected specifically for

evaluation, collected for other purposes (such as research) but useful for evaluation, or extant. In effect, evaluation questions connect evaluation data to the goals. Evaluation questions here should be top-level questions only.

Table A.5.1. Example: Evaluation questions aligned with [project name] goals, strategies/activities, outputs, and outcomes.

Goals	Strategies/Activities	Outputs	Outcomes	Evaluation Questions
<b>Example</b>				
<p><b>1. To what extent do teachers who participated in a model demonstration training project for language instruction exhibit improved language instruction in the classroom.</b></p>	<ul style="list-style-type: none"> <li>▪ Develop a tracking system to maintain contact with teachers.                             <ul style="list-style-type: none"> <li>▫ Develop a reliable and efficient system that will allow teachers to be followed for 1 year after completing the training.</li> <li>▫ Obtain commitments from training participants to participate in follow-up activities after training.</li> <li>▫</li> </ul> </li> <li>▪ Develop a plan for systematically obtaining data on the performance of teachers and their students.                             <ul style="list-style-type: none"> <li>▫ Assemble a set of measures for determining skills and knowledge of teachers immediately after training.</li> <li>▫ Work in conjunction with districts and agencies where teachers are employed to develop valid and practical measures for determining teacher performance in the area of language instruction.</li> <li>▫ Work in conjunction with districts and agencies to determine the most valid and efficient means for collecting data on the performance of children/students that graduates serve.</li> </ul> </li> <li>▪ Implement the tracking system and integrate with other data collection plans</li> </ul>	<ul style="list-style-type: none"> <li>- Set of measures that capture teacher exit skills and knowledge assembled.</li> <li>- Work plan for tracking and follow-up of teachers.</li> <li>- Work plan for the collection of valid data on the performance of teachers and of their students.</li> </ul>	<p><u>Short-term</u></p> <ul style="list-style-type: none"> <li>• Teachers exit the training program having demonstrated high level language teaching skills.</li> <li>• Graduates are tracked for 1 year.</li> <li>• Data on the performance of teachers are collected for 1 year.</li> </ul> <p><u>Intermediate</u></p> <ul style="list-style-type: none"> <li>• Graduates are working in appropriate settings for a minimum of 3 years.</li> <li>• Teachers continue to provide high quality language instruction.</li> </ul>	<ul style="list-style-type: none"> <li>A. <b>Are we collecting valid data on the performance of our teachers?</b></li> <li>B. <b>To what extent do graduates exit the program with the skills and knowledge necessary to perform at a high level?</b></li> <li>C. <b>To what extent do teachers exhibit use of the language instruction needed to improve outcomes for children with disabilities?</b></li> </ul>

Goals	Strategies/Activities	Outputs	Outcomes	Evaluation Questions
<p><b>2. Teachers who participated in the training demonstrate success with children/ students.</b></p>	<ul style="list-style-type: none"> <li>▪ Select/develop valid measures of student language outcomes <ul style="list-style-type: none"> <li>▫ Examine extant data availability.</li> <li>▫ Explore testing options.</li> <li>▫ Develop teacher reporting protocol.</li> </ul> </li> <li>▪ Establish data collection system <ul style="list-style-type: none"> <li>▫ Prepare plan</li> <li>▫ Develop timeline</li> <li>▫</li> </ul> </li> <li>▪ Develop analysis plan <ul style="list-style-type: none"> <li>▫</li> <li>▫</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>- Instruments selected or developed and tested.</li> <li>- Data collection plan and timeline.</li> <li>- Analysis plan.</li> <li>-</li> </ul>	<p><u>Short-term</u></p> <ul style="list-style-type: none"> <li>• A measurement system is implemented.</li> <li>• Valid data are collected.</li> </ul> <p><u>Intermediate</u></p> <ul style="list-style-type: none"> <li>• Evidence of success with children/students served by teachers.</li> </ul> <p><u>Long-term</u></p> <ul style="list-style-type: none"> <li>• Students of teachers demonstrate improved outcomes.</li> </ul>	<p>D. To what extent do trained teachers demonstrate success with child/students with disabilities?</p>

### Data Sources and Collection by Evaluation Question

Fill in Table A.5.2, using however many rows you require. Table A.5.2 presents the evaluation questions in relation to outcomes, data sources, and data collection strategies. The first row provides an example, using Evaluation Question c from the Goal 1 example in Table A.5.1. You will provide more detail about data and data collections as indicated in the sections that follow the table.

Table A.5.2. Example: Evaluation questions in relation to outcomes, instruments, and data collections.

Evaluation Questions	Outcomes	Instruments or Datasets	Modes of Data Collection	Comparison Data Collection Planned (✓ if yes)	Timeframe for Initial Data Collection	Timeframe for Recurrence of Data Collection, If Any
<b>Example</b>						
<b>A. To what extent do teachers exit the training program with the skills to implement improved language instruction?</b>	<p><u>Short-term</u></p> <ul style="list-style-type: none"> <li>Teachers exit the program having demonstrated the skills and knowledge to deliver improved language instruction.</li> </ul>	<ul style="list-style-type: none"> <li>Comprehensive examines</li> <li>Observation protocols for trainers                             <ul style="list-style-type: none"> <li>Short form</li> <li>Extended observation form</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>Tests</li> <li>Observation protocol</li> </ul>	<input type="checkbox"/>	End of training	None
<b>B. To what extent do teachers demonstrate success with children/students with disabilities?</b>	<p><u>Short-term</u></p> <ul style="list-style-type: none"> <li>A measurement system is implemented.</li> <li>Valid data are collected.</li> </ul> <p><u>Intermediate</u></p> <ul style="list-style-type: none"> <li>Evidence of success with children/students served by teachers.</li> </ul>	<ul style="list-style-type: none"> <li>State testing program</li> <li>Instruments for measuring non-academic student outcomes.</li> <li>Protocols for extraction of extant data on student outcomes.</li> <li>Dataset of graduate performance data and corresponding student outcome data.</li> </ul>	<ul style="list-style-type: none"> <li>Testing.</li> <li>Review of extant data.</li> </ul>	<input checked="" type="checkbox"/>	3 months after start of grant	Annually for 3 years

## Special Data Collection Considerations

### Instrument or Dataset Status

*Instructions: List and discuss status of each instrument—ready to use, needing modification, under development, planned by the project. Discuss availability of each data set planned for use—access to the data set and timeliness of the data availability.*

**Example:** The project will need to develop a tracking system capable of keeping up with contact information and employment status of each teacher. The system will need to be easily updatable and will include logs of contacts and attempted contacts with the teachers as well as alternative contacts who may be of assistance in locating a teacher.

### Comparison Data

*Instructions: Discuss need for baseline and, if needed, whether baseline data already exist or how baseline data collection might work. Identify possibilities for experimental or quasi-experimental designs—discuss planned use of comparison groups and whether random assignment will be used. Describe any available extant datasets that could be used in addition to or instead of collected data.*

**Example:** To determine the student outcomes for teachers who participated in the language instruction training program (Evaluation Question B) in comparison to teachers who did not receive this additional training, the project has made an arrangement with Districts X and Y, where 90 percent of our teachers are employed, to use administrative records of progress on a relevant speech/language assessment from the students of these teachers. To protect personnel confidentiality, the data will be provided to us in masked form in two datasets. The first dataset will be the scores of students of the trained teachers. The second dataset will be the corresponding statistics for students for all other teachers. We will use the second dataset as our comparison data. When possible without risking the exposure of individual students' identities, the data will include individual student descriptors, such as sex, SES, and disability category.

### Sampling

*Instructions: Discuss need for or consideration of sampling, if any, and present some of the specifics if known.*

**Example:** To minimize costs to the project and burden on evaluators, in addressing Evaluation Question A we will employ sampling. The project follow an extended observation protocol for a representative sample of our enrolled teachers during the training. These data will be collected weekly. The random sample will be stratified to ensure that it's representative of key characteristics of our candidates: sex, race/ethnicity, disability status.

Sampling will also be employed when the performance of children served by teachers is measured for Evaluation Question B. The project will manage a data collection of social language (pragmatics) for a random sample children being served by teachers who participated in the training and a random sample of children being served by teachers who did not participate in the training. We will use the Language Use Inventory (LUI) to collect data twice a year, with permission of the parents or guardians, who will be the respondents. The random sample will be stratified to ensure that it's representative of the demographics of the children served by trained teachers—that is, it will be stratified by grade, sex, and race/ethnicity to ensure that a sufficient number of students with each of these characteristics are represented to allow for generalization for these characteristics. The LUI is a standardized instrument with well-known characteristics, which will allow an accurate estimation of needed sample sizes by our statistician.

## Analysis Approach

## Descriptive Analyses

*Instructions: Discuss need for descriptive analyses—that is, for which evaluation questions will descriptive analyses be used. Present specific techniques corresponding to the evaluation questions and data collections.*

Example: Analysis will generate descriptive statistics related to Evaluation Question A and B on various measures of teacher and student performance. For example, analysis of observation data will identify the percentage of teachers using specific evidence-based practices at or above a specified threshold level. Additionally, cross tabulations will be used to identify factors associated with the strength of the results. For example for Question B, cross-tabulated data may show that a higher percentage of students in a certain grade demonstrated improved scores compared to students in another grade. Other predictors, such as student gender or disability category, will also be examined.

## Statistical Analyses

*Instructions: Discuss need for statistical analyses—for which evaluation questions. Present specific techniques being considered, if these are known, corresponding to the evaluation questions and data collections. Identify specific “studies” if they are sufficiently discrete.*

Example: The project will use statistical analysis to partly answer Evaluation Question B. For teachers who did and who did not receive the training, a two-group comparison will be possible using pre and post scores from the administrative records of relevant speech/language assessments. One group will be children trained teachers served; the other group will be comparable children served by teachers who did not receive the additional language instruction training. Scores across different assessments will be standardized as effect sizes and will be analyzed as change scores for individual children. We will use a one-sided t-test to compare the two groups. Children will be matched on pretest scores, race, SES, testing accommodations, and disability category.



## A.6. Data Analysis Plan Template

Evaluation Question	Design		Data Analysis	Necessary Variables for Quantitative Analysis		Variable Sources (instruments or data collection techniques)	Data Sources	Minimum number of responses and/or response rate
	Design Type	If experimental or quasi-experimental who constitutes the...		Statistical Tests	Descriptive Statistics			
(1)	<ul style="list-style-type: none"> <li>• Experimental</li> <li>• Quasi-experimental</li> <li>• Single-case</li> <li>• Non-experimental</li> </ul>	Treatment group:  Control or Comparison Group:	<ul style="list-style-type: none"> <li>• Statistical Tests</li> <li>• Descriptive Statistics</li> <li>• Visual analysis (SCD)</li> <li>• Qualitative analysis</li> </ul>	Dependent:  Independent:  Covariates:	<ul style="list-style-type: none"> <li>• Frequency</li> <li>• Mean</li> </ul>		<ul style="list-style-type: none"> <li>• Census</li> <li>• Sample</li> </ul>	
(2)	<ul style="list-style-type: none"> <li>• Experimental</li> <li>• Quasi-experimental</li> <li>• Single-case</li> <li>• Non-experimental</li> </ul>	Treatment group:  Control or Comparison Group:	<ul style="list-style-type: none"> <li>• Statistical Tests</li> <li>• Descriptive Statistics</li> <li>• Visual analysis (SCD)</li> <li>• Qualitative analysis</li> </ul>	Dependent:  Independent:  Covariates:	<ul style="list-style-type: none"> <li>• Frequency</li> <li>• Mean</li> </ul>		<ul style="list-style-type: none"> <li>• Census</li> <li>• Sample</li> </ul>	
(3)	<ul style="list-style-type: none"> <li>• Experimental</li> <li>• Quasi-experimental</li> <li>• Single-case</li> <li>• Non-experimental</li> </ul>	Treatment group:  Control or Comparison Group:	<ul style="list-style-type: none"> <li>• Statistical Tests</li> <li>• Descriptive Statistics</li> <li>• Visual analysis (SCD)</li> <li>• Qualitative analysis</li> </ul>	Dependent:  Independent:  Covariates:	<ul style="list-style-type: none"> <li>• Frequency</li> <li>• Mean</li> </ul>		<ul style="list-style-type: none"> <li>• Census</li> <li>• Sample</li> </ul>	



## A.7. Fidelity Matrix Template

Key Component	Operational Definition/ Indicators	Data Sources/ Measures	Fidelity Threshold (Threshold Level)	Score for “Adequate” Fidelity	Fidelity Rating
<p><i>Example: Key Component 1: Develop and provide training workshops to parents of children with disabilities</i></p>	<p><i>Operational Definition: For the training provided by the Center to be considered implemented with fidelity, (a) the trainings will be delivered by lead parent educators who have the required competencies, (b) training workshops will be delivered on schedule and within the allotted timeframe, and will cover the expected content; and (c) parents will participate in and express satisfaction with the workshops.</i></p> <p><i>Indicator 1—Staff Competencies:</i></p> <ul style="list-style-type: none"> <li>• <i>Staff members who serve as lead parent educators have</i> <ul style="list-style-type: none"> <li>a. <i>a Bachelor’s Degree or an Associate’s Degree in Special Education, Early Childhood Education, or a related field;</i></li> <li>b. <i>at least one year of experience working with families of children with disabilities; and</i></li> <li>c. <i>positive evaluations that provide evidence that the staff member understands the content being presented, effectively interacts with parents, and effectively facilitates parent training workshops.</i></li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• <i>Administrative data on staff credentials and experience</i></li> <li>• <i>Staff evaluations</i></li> </ul>	<p><b>Competency indicator (a) (project-level)</b>  <b>2 = high fidelity:</b> ≥ 90% of lead parent educators on staff meet competency (a)  <b>1= moderate fidelity:</b> 75%-89% of lead parent educators on staff meet competency indicator (a)  <b>0 = low fidelity:</b> Less than 75% of lead parent educators on staff meet competency indicator (a)</p> <p><b>Competency indicator (b) (project-level)</b>  <b>2 = high fidelity:</b> ≥ 90% of lead parent educators on staff meet competency indicator (b)  <b>1= moderate fidelity:</b> 75%-89% of lead parent educators on staff meet competency indicator (b)  <b>0 = low fidelity:</b> Less than 75% of lead parent educators on staff meet competency indicator (b)</p> <p><b>Competency indicator (c) (project-level)</b>  <b>2 = high fidelity:</b> ≥ 75% of lead parent educators on staff meet competency indicator (c)  <b>1= moderate fidelity:</b> 60%-75% of lead parent educators on staff meet competency indicator (c)  <b>0 = low fidelity:</b> Less than 60% of lead parent educators on staff meet competency indicator (c)</p>	<p><b>4</b></p> <p><i>Note: The project must get a score of 2 on Competency Indicator (a); and at least a score of 1 on the other two Competency Indicators.</i></p>	<p>TBD</p> <p><i>This will be calculated after the data are collected through the evaluation.</i></p>
<p>1. [KEY COMPONENT]</p>	<p>[OPERATIONAL DEFINITION]</p> <p>a. [INDICATOR]</p> <p>b. [INDICATOR]</p> <p>c. [INDICATOR]</p>	<p>[SOURCES]</p> <p>[MEASURES]</p>	<p><b>Competency indicator (a) ([LEVEL])</b>  <b>2 = high fidelity:</b> [THRESHOLD]  <b>1= moderate fidelity:</b> [THRESHOLD]  <b>0 = low fidelity:</b> [THRESHOLD]</p> <p><b>Competency indicator (b) ([LEVEL])</b>  <b>2 = high fidelity:</b> [THRESHOLD]  <b>1= moderate fidelity:</b> [THRESHOLD]  <b>0 = low fidelity:</b> [THRESHOLD]</p> <p><b>Competency indicator (c) ([LEVEL])</b>  <b>2 = high fidelity:</b> [THRESHOLD]  <b>1= moderate fidelity:</b> [THRESHOLD]  <b>0 = low fidelity:</b> [THRESHOLD]</p>	<p>[NUMERIC SCORE]</p> <p>[ANY REQUIREMENTS FOR SCORES]</p>	<p>[FINAL RATING]</p>

Key Component	Operational Definition/ Indicators	Data Sources/ Measures	Fidelity Threshold (Threshold Level)	Score for "Adequate" Fidelity	Fidelity Rating
2. [KEY COMPONENT]	[OPERATIONAL DEFINITION] a. [INDICATOR] b. [INDICATOR] a. [INDICATOR]	[SOURCES] [MEASURES]	<p><b>Competency indicator (a) ([LEVEL])</b> 2 = high fidelity: [THRESHOLD] 1= moderate fidelity: [THRESHOLD] 0 = low fidelity: [THRESHOLD]</p> <p><b>Competency indicator (b) ([LEVEL])</b> 2 = high fidelity: [THRESHOLD] 1= moderate fidelity: [THRESHOLD] 0 = low fidelity: [THRESHOLD]</p> <p><b>Competency indicator (c) ([LEVEL])</b> 2 = high fidelity: [THRESHOLD] 1= moderate fidelity: [THRESHOLD] 0 = low fidelity: [THRESHOLD]</p>	[NUMERIC SCORE]  [ANY REQUIREMENTS FOR SCORES]	[FINAL RATING]
3. [KEY COMPONENT]	[OPERATIONAL DEFINITION] a. [INDICATOR] b. [INDICATOR] a. [INDICATOR]	[SOURCES] [MEASURES]	<p><b>Competency indicator (a) ([LEVEL])</b> 2 = high fidelity: [THRESHOLD] 1= moderate fidelity: [THRESHOLD] 0 = low fidelity: [THRESHOLD]</p> <p><b>Competency indicator (b) ([LEVEL])</b> 2 = high fidelity: [THRESHOLD] 1= moderate fidelity: [THRESHOLD] 0 = low fidelity: [THRESHOLD]</p> <p><b>Competency indicator (c) ([LEVEL])</b> 2 = high fidelity: [THRESHOLD] 1= moderate fidelity: [THRESHOLD] 0 = low fidelity: [THRESHOLD]</p>	[NUMERIC SCORE]  [ANY REQUIREMENTS FOR SCORES]	[FINAL RATING]
4. [KEY COMPONENT]	[OPERATIONAL DEFINITION] a. [INDICATOR] b. [INDICATOR] a. [INDICATOR]	[SOURCES] [MEASURES]	<p><b>Competency indicator (a) ([LEVEL])</b> 2 = high fidelity: [THRESHOLD] 1= moderate fidelity: [THRESHOLD] 0 = low fidelity: [THRESHOLD]</p> <p><b>Competency indicator (b) ([LEVEL])</b> 2 = high fidelity: [THRESHOLD] 1= moderate fidelity: [THRESHOLD] 0 = low fidelity: [THRESHOLD]</p> <p><b>Competency indicator (c) ([LEVEL])</b> 2 = high fidelity: [THRESHOLD] 1= moderate fidelity: [THRESHOLD] 0 = low fidelity: [THRESHOLD]</p>	[NUMERIC SCORE]  [ANY REQUIREMENTS FOR SCORES]	[FINAL RATING]

## Appendix B. Validity Threats

### Threats to Internal Validity

Threats to internal validity relate to whether the study results can actually be attributed to the variables included in the study, or whether some confounding variables might be affecting the outcomes. The principle threats to internal validity include<sup>206</sup>:

- **History:** When study participants experience an event during the study period that might influence their performance. This might occur if all teachers at a school, including new program graduates, are offered an intensive professional development course by their school district in the summer prior to the academic year under study. In this situation it would be difficult to separate the effects of the summer training program from the training provided by the personnel preparation program.
- **Maturation:** During a lengthy study, biological or cognitive maturation of study participants may affect the outcomes of the study. This is particularly an issue when studying young children or assessing performance of skills that can be expected to change as a result of cognitive maturation, such as when assessing the motor skills of children beginning in infancy through age two.
- **Pretesting:** While pretests are common in experimental and quasi-experimental studies, the exposure of study participants to a pretest may influence the results of the posttest due to a practice effect or simply being more aware of the topic under study. An example of this might be when students are given a spelling test to assess their spelling ability, followed by direct instruction in spelling and a posttest to observe changes in the students' performance. One way to avoid this is to use alternate forms of a test for the pre- and the posttests.
- **Measuring instruments:** Using different instruments to collect data, such as standardized tests, observation rubrics, or teacher self-report surveys, can affect the accuracy of scores. This would be of particular concern in cases when different groups of students or teachers are assessed using different measurement instruments, such as two different standardized math tests, and then their results are compared.
- **Statistical regression to the mean:** When a study participant performs extremely well or extremely poorly on a particular test or other measure, it's common that his or her performance will be less extreme on a subsequent test or measure. This is called regression toward the mean and it occurs because the initially observed extreme positive or negative scores contain relatively large (positive or negative) random error that will probably not be as large with subsequent measures, thus causing the subsequent score to be closer to the mean.
- **Differential selection:** In studies that have treatment and control groups, pretest differences between the groups, such as initial reading ability, will carry over to the posttest. If these pretest differences aren't accounted for in the analysis—for example, by using a pretest measure of reading ability as a covariate in the analysis—it will not be possible to determine how much of the group differences seen at the posttest are due to the pre-existing differences among the groups and how much are due to the treatment.
- **Attrition:** When some study participants drop out of one or more groups in study before it's completed, it's typically not a random process. For example, the parents of a lower-performing student may withdraw their child from an intervention if they fear that the intervention might have an adverse effect on their child's self-confidence. This is also known as *experimental mortality*. Attrition is more of a threat to the internal validity of experimental studies, since quasi-experimental studies already account for the fact that there are differences between the groups under study. Whenever possible, you should collect baseline data for all groups under study and maintain records of the individual participants in the study (e.g., using class rosters to track student enrollment from fall to spring) so that you can examine whether the analytic sample—the sample that remained

---

<sup>206</sup> Dimitrov, 2010; see also Shadish et al., 2002.

in your study throughout the entire study period—contains group differences on key variables (such as initial ability) that may call into question the results of your study.

- **Interaction among factors:** Some of the threats to internal validity mentioned above may interact, thus producing additional confounding effects on the results of the experimental study.

## Threats to External Validity

Threats to external validity affect your ability to generalize your results to persons, settings, treatments, and outcomes not directly included in the study. The major threats to external validity include<sup>207</sup>:

- **Interaction of selection biases with experimental treatment:** When a treatment is more effective for participants with particular characteristics, the findings cannot be generalized. For example, a computer-based mathematics intervention will likely be more effective for students who already have experience working with computers than for students who have relatively little prior exposure to computers.
- **Reactive effect of pretesting:** When participants take a pretest they become aware of and sensitized to the issues targeted by the treatment. Therefore the post-treatment results may not be generalized to a population of participants that has not received a pretest. This might occur if at the beginning of a student's senior year you administer a pretest measure designed to gauge a student's awareness of available postsecondary options and then administer a similar posttest measure at the end of the student's senior year.
- **Reactive effects of experimental procedures:** Participants in a study often react to the presence of observers and experimental procedures, thereby altering their behavior. This makes it difficult to generalize the findings to persons who are exposed to a treatment in normal settings. An example of this might be when an evaluator installs a video camera in a pull-out classroom to record the interactions of a student with a speech and language pathologist.
- **Multiple-treatment interference:** When participants in a study are exposed to multiple treatments (or variations of the same treatment), the effect of the second (and any subsequent) treatment might be confounded with residual effects of the preceding treatment. Consequently, the overall outcome of the treatments will depend, among other things, on the sequence in which they were introduced. This might be a factor when a school psychologist is conducting applied behavioral analysis using an alternating treatment single-case design.<sup>208</sup>

---

<sup>207</sup> Dimitrov, 2010; see also Shadish et al., 2002.

<sup>208</sup> See Kennedy, 2005, for more information on the different types of single-case designs.

## Appendix C. Sample Forms

## C.1 Sample Notification Letter for Districts with Research Approval Office/Department

Dear <Sal> <SupFName> <SupLName>:

I am writing to inform you that my organization, <EvaluatorName> is planning to conduct site visits and interviews with special educators and administrators at <SchoolName(s)> in your district. We have already completed the research application required by your district and received approval to conduct this study. We have been contracted by <GranteeName> to conduct an evaluation as part of a federal requirement to evaluate the performance of personnel preparation programs that receive funding from the Personnel Development Program in the U.S. Department of Education's Office of Special Education Programs. This evaluation involves gathering data to assess the performance of special education teachers and related-services providers who graduated from <GranteeName> in recent years, based in part on measures of educator practice and student achievement.

We plan to conduct <NumberObservations> of site visits to the schools where these <GranteeName> graduates are currently working. These site visits will consist of:

- an interview with the school principal or his/her designee,
- individual or group interviews with teachers or service providers, and
- observations of educator practice.

The interviews and observations should each last between <InterviewLength> and <ObservationLength>. There is no need to prepare or provide any documentation.

The data we collect during the visits will not be used to evaluate the schools' or your district's performance. It will be aggregated with data collected on <GranteeName> graduates working in other schools and districts as part of a report on the performance of <GranteeName> program as a whole. All data collected for this study will be kept confidential, except as required by law. A report will be delivered to <GranteeName> with results aggregated across all respondents.

Thank you in advance for your district's cooperation and participation in this important study. Feel free to contact me directly with questions or issues. I can be reached by calling <EvaluationDirectorTelephone> or by emailing <EvaluationDirectorEmail>.

Sincerely,

<EvaluationDirector>

<EvaluationDirectorContactInformation>

## C.2. Sample Request Letter for Districts without Research Approval Office/Department

Dear <Sal> <SupFName> <SupLName>:

We are requesting permission to conduct site visits at schools in your district where graduates of <GranteeName> are currently working. My organization, <EvaluatorName>, is conducting a study under contract with <GranteeName> to evaluate the performance of graduates of their personnel preparation program. This evaluation is part of a federal requirement to monitor the performance of personnel preparation programs that receive funding from the Personnel Development Program in the U.S. Department of Education's Office of Special Education Programs. This evaluation involves gathering data to assess the performance of special education teachers and related-services providers who graduated from <GranteeName> in recent years. **Your district's participation is important for <GranteeName> to ensure that its program maintains high standards so that its graduates continue providing high quality instructional services to children like those in your district.**

These site visits will consist of:

- an interview with the school principal or his/her designee,
- individual or group interviews with teachers or service providers, and
- observations of educator practice.

We plan to conduct <NumberInterviews> and <NumberObservations> over the course of <StudyLength>. The interviews and observations should each last between <InterviewLength> and <ObservationLength>. There is no need to prepare or provide any documentation.

The data we collect during the visits will not be used to evaluate the schools' or your district's performance. It will be aggregated with data collected on <GranteeName> graduates working in other schools and districts as part of a report on the performance of <GranteeName> program as a whole. All data collected for this study will be kept confidential, except as required by law. A report will be delivered to <GranteeName> with results aggregated across all respondents.

Thank you in advance for considering participating in this important study. Please indicate the decision of your district on the enclosed form. Feel free to contact me directly with questions or issues. I can be reached by calling <EvaluationDirectorTelephone> or by emailing <EvaluationDirectorEmail>.

Sincerely,

<EvaluationDirector>

<EvaluationDirectorContactInformation>

### C.3. Sample District Response Form

I give district permission for the evaluation study conducted by <EvaluatorName> to take place during the current school year. This study is designed to evaluate the graduates of <GranteeName> and will include site visits to participating schools.

---

Name

---

Date

---

Position

---

District

---

Phone Number

## C.4. Sample School Notification Letter

Dear <Sal> <PrinFName> <PrinLName>:

I am writing to ask your permission for my organization, <EvaluatorName>, to visit your school and interview you and selected staff in <VisitTimeframe>. We have been contracted by <GranteeName> to conduct a study as part of a federal requirement to evaluate the performance of personnel preparation programs that receive funding from the Personnel Development Program in the U.S. Department of Education's Office of Special Education Programs. **We realize that your time is valuable and we will do our best to make your participation in this study as easy as possible.** Your help is very important to our efforts to improve the quality of personnel preparation programs, and subsequent outcomes for students. This study involves gathering data to assess the performance of special education teachers and related-services providers who graduated from <GranteeName> in recent years, based in part on measures of educator practice and student achievement. We have already received district approval to conduct this study.

Our plan is to conduct <NumberObservations> of site visits to your school. These site visits will consist of:

- a brief interview with you or your designee about the performance of each graduate of <GranteeName> working at your school,
- individual or group interviews with the graduates at your school, and
- observations of the graduates' practice.

The interviews and observations should each last between <InterviewLength> and <ObservationLength>. There is no need to prepare or provide any documentation prior to our visit.

The data we collect during the visits will not be used to evaluate your or your school's performance. It will be aggregated with data collected on <GranteeName> graduates working in other schools as part of a report on the performance of <GranteeName> program as a whole. All data collected for this study will be kept confidential, except as required by law. A report will be delivered to <GranteeName> with results aggregated across all respondents.

We understand that making staff and activities available to the site visitor will require time and effort from you and your staff. We appreciate your help with our efforts to improve training programs and outcomes for students. We will be in contact with you in the near future to arrange the details, date, and schedule of the visit to your school.

We appreciate your willingness to cooperate and provide information to help <GranteeName> improve the quality of its personnel preparation program. Feel free to contact me directly with questions or issues. I can be reached by calling < EvaluationDirector Telephone> or by emailing < EvaluationDirector Email>.

Sincerely,

<EvaluationDirector>

< EvaluationDirector ContactInformation>

## C.5. Sample Passive Consent Form for PDP Graduates

Dear <Sal> <GradFName> <GradLName>:

I am writing to tell you about an important study of the <GranteeName> personnel preparation program that you graduated from in <GraduationYear>. This study is part of a federal requirement to evaluate the performance of personnel preparation programs that receive funding from the Personnel Development Program in the U.S. Department of Education's Office of Special Education Programs. **Your participation in this study is extremely important** to our efforts to improve the quality of personnel preparation programs, and subsequent outcomes for students.

As part of this evaluation, my organization, <EvaluatorName>, would like to visit your school to observe your practice and interview you and your personnel supervisor. We have already completed the research application required by your district and received approval to conduct this study.

We expect to conduct the interviews and observations in <VisitTimeframe>. The interview and observations should each last between <InterviewLength> and <ObservationLength>. There is no need to prepare or provide any documentation prior to our visit.

The data we collect during the visits will be kept confidential, except as required by law, and will be aggregated with data collected on <GranteeName> graduates working in other schools to develop a report on the performance of the <GranteeName> program as a whole. A report will be delivered to <GranteeName> with results aggregated across all respondents—no identifying information about you or your school will be included in the report.

We understand that making yourself available to participate in this study will require time and effort and we greatly appreciate your willingness to cooperate. We appreciate your help with our efforts to improve training programs and outcomes for students. We will be in contact with you in the near future to arrange the details, date, and schedule of the visit.

We hope you will be willing to cooperate and provide information to help <GranteeName> improve the quality of its personnel preparation program. **If you do not wish to participate, please notify me** at < EvaluationDirector Telephone> or by emailing < EvaluationDirector Email>.

Sincerely,

<EvaluationDirector>

< EvaluationDirector ContactInformation>

## C.6. Sample Active Consent form for PDP Graduates

Dear <Sal> <GradFName> <GradLName>:

I am writing to ask you to participate in study of the <GranteeName> personnel preparation program that you graduated from in <GraduationYear>. This study is part of a federal requirement to evaluate the performance of personnel preparation programs that receive funding from the Personnel Development Program in the U.S. Department of Education's Office of Special Education Programs. **Your participation in this study is extremely important** to our efforts to improve the quality of personnel preparation programs, and subsequent outcomes for students.

As part of this evaluation, my organization, <EvaluatorName>, would like to visit your school to observe your practice and interview you and your personnel supervisor. Additionally, we ask that you complete a brief survey related to your perceptions of the quality of the training and support you received from <GranteeName> and your perceived self-efficacy as a <Profession>. We have already completed the research application required by your district and received approval to conduct this study.

We expect to conduct the interviews and observations in <VisitTimeframe>. The interview and observations should each last between <InterviewLength> and <ObservationLength>. There is no need to prepare or provide any documentation prior to our visit.

The data we collect during the visits will be kept confidential, except as required by law, and will be aggregated with data collected on <GranteeName> graduates working in other schools to develop a report on the performance of the <GranteeName> program as a whole. A report will be delivered to <GranteeName> with results aggregated across all respondents—no identifying information about you or your school will be included in the report.

We understand that making yourself available to participate in this study will require time and effort and we greatly appreciate your willingness to cooperate. We appreciate your help with our efforts to improve training programs and outcomes for students. **If you agree, please complete the form below** and return it to us. We will be in contact with you in the near future to arrange the details, date, and schedule of the visit and to provide information on how to access the survey. We appreciate your willingness to cooperate and provide information to help <GranteeName> improve the quality of its personnel preparation program. Feel free to contact me directly with questions or issues. I can be reached by calling < EvaluationDirector Telephone> or by emailing < EvaluationDirector Email>.

Sincerely,

<EvaluationDirector>

### Study Permission Form

By returning this form, I **agree** to participate in this study.

(Please Print) My name is: \_\_\_\_\_ School: \_\_\_\_\_

Signature: \_\_\_\_\_ Date: \_\_\_\_\_

## C.7. Sample Passive Consent Form for Students

Dear Parent or Guardian:

Your child is being asked to complete a survey as part of a study of the performance of the special education personnel preparation program operating at <GranteeName>. This survey will gather information from your child about the performance of his or her special education teacher or related services provider. Please read this form for information about the survey, and for instructions on how to withdraw your child. *If you do not want your child to complete the survey, you must notify your school.*

**Survey Content.** The survey gathers information on <include brief description of survey contents>. You may examine the questionnaire in the school office or at the following Web site <WebSite>.

All data collected by the surveys will be kept confidential, except as required by law. The data we collect will be combined with the results of surveys of students in other schools and districts. A report summarizing the overall results of the surveys will be delivered to <GranteeName> with results aggregated across all survey respondents.

**It is Voluntary.** Your child does not have to take the survey. Students who participate only have to answer the questions they want to answer and they may stop answering questions at any time.

**It is Anonymous.** No names will be recorded or attached to the survey forms or data. The results will be made available only to the researchers using strict confidentiality controls.

**Administration.** The survey will be administered on <DateSurvey>. It will take about <SurveyTime> to complete and will be administered in your child's <ClassName> class.

**Potential Risks.** There are no known risks of harm to your child.

**Potential Benefits.** No direct benefits to your child are expected, however, the study is expected to help improve the quality of training for special education professionals.

**For Further Information.** The survey was developed by <SurveyDeveloper>. If you have any questions about this survey call me at <EvaluationDirectorTelephone> or by email me at <EvaluationDirectorEmail>. If you have question about your rights related to study participation, call the district at <INSERT NAME AND PHONE NUMBER OF DISTRICT CONTACT>.

If you do not want your child to participate, you may contact: <INSERT CONTACT INFORMATION (E.G., ADDRESS, PHONE NUMBER, E-MAIL).> [Note: We recommend using a single point of contact.]

### Study Withdrawal Form

By returning this form, I ***do not give permission*** for my child to participate in this study.

(Please Print) My child's name is: \_\_\_\_\_ Grade: \_\_\_\_\_

Teacher's name or Class subject: \_\_\_\_\_

Signature: \_\_\_\_\_ Date: \_\_\_\_\_



## C.8. Sample Active Consent Form for Students

Dear Parent or Guardian:

Your child is being asked to complete a survey as part of a study of the performance of the special education personnel preparation program operating at <GranteeName>. This survey will gather information from your child about the performance of his or her special education teacher or related services provider. Please read this form for information about the survey, and for instructions on how to withdraw your child. *You must sign this form and return it to your school if you give your permission for your child to participate in this study.*

**Survey Content.** The survey gathers information on <include brief description of survey contents>. You may examine the questionnaire in the school office or at the following Web site <WebSite>.

All data collected by the surveys will be kept confidential, except as required by law. The data we collect will be combined with the results of surveys of students in other schools and districts. A report summarizing the overall results of the surveys will be delivered to <GranteeName> with results aggregated across all survey respondents.

**It is Voluntary.** Your child does not have to take the survey. Students who participate only have to answer the questions they want to answer and they may stop taking it at any time.

**It is Anonymous.** No names will be recorded or attached to the survey forms or data. The results will be made available only to the researchers using strict confidentiality controls.

**Administration.** The survey will be administered on <DateSurvey>. It will take about <SurveyTime> to complete and will be administered in your child's <ClassName> class.

**Potential Risks.** There are no known risks of harm to your child.

**Potential Benefits.** No direct benefits to your child are expected, however, the study is expected to help improve the quality of training for special education professionals.

**For Further Information.** The survey was developed by <SurveyDeveloper>. If you have any questions about this survey call me at <EvaluationDirectorTelephone> or by email me at <EvaluationDirectorEmail>. If you have question about your rights related to study participation, call the district at <INSERT NAME AND PHONE NUMBER OF DISTRICT CONTACT>.

### Study Permission Form

By returning this form, I **give permission** for my child to participate in this study.

(Please Print) My child's name is: \_\_\_\_\_

Grade: \_\_\_\_\_

Teacher's name or Class subject: \_\_\_\_\_

Signature: \_\_\_\_\_

Date: \_\_\_\_\_



## Appendix D. Recommended Readings on Research/Evaluation Methodology

- Abry, T., Hulleman, C., & Rimm-Kaufman, S. (2015). Using indices of fidelity to intervention core components to identify program active ingredients. *American Journal of Evaluation, 36*(3), 320-338.
- Baecher, L. H., & Connor, D. J. (2010). "What do you see?" Using video analysis of classroom practice in a preparation program for teachers of students with learning disabilities. *Insights on Learning Disabilities 7*(2), 5-18.
- Barton, E. E., & Reichow, B. (2012). Guidelines for graphing data with Microsoft PowerPoint for Office 2007. *Journal of Early Intervention, 34*, 129-150.
- Brewer, J., & Hunter, A. (2006). *Foundations of multimethod research: Synthesizing styles*. Thousand Oaks, CA: Sage.
- Campbell, D. T., & Stanley, J. C. (1963). Experimental and quasi-experimental designs for research. In N. L. Gage (Ed.), *Handbook of research on teaching* (pp. 1-76). Chicago: Rand-McNally.
- Century, J., Rudnick, M., & Freeman, C. (2010). A framework for measuring fidelity of implementation: A foundation for shared language and accumulation of knowledge. *American Journal of Evaluation, 34*(4): 465-485.
- Creswell, J. W. (2003). *Research design: Qualitative, quantitative and mixed methods approaches* (2nd ed.). Thousand Oaks, CA: Sage.
- Czaja, R., & Blair, J. (2005). *Designing surveys: A guide to decisions and procedures* (2nd ed.). Thousand Oaks, CA: Sage.
- Dillman, D., Smyth, J., & Christian, L. M. (2009). *Internet, mail, and mixed-mode surveys* (3rd edition). Hoboken, NJ: John Wiley & Sons.
- Dimitrov, D. M. (2010). *Quantitative research in education: Intermediate and advanced methods* (2nd ed.). New York: Whittier Publications.
- Dong, N., & Maynard, R. A. (2013). PowerUp!: A tool for calculating minimum detectable effect sizes and minimum required sample sizes for experimental and quasi-experimental design studies. *Journal of Research on Educational Effectiveness, 6*(1), 24-67. DOI: 10.1080/19345747.2012.673143
- Fixsen, D. L., Naoom, S. F., Blase, K. A., Friedman, R. M., & Wallace, F. (2005). *Implementation research: A synthesis of the literature*. Retrieved from <http://ctndisseminationslibrary.org/PDF/nirnmonograph.pdf>
- Frechtling, J., with M. Melvin, D. Rog, & E. Johnson. (2010). *The 2010 user-friendly handbook for project evaluation*. Rockville, MD: Westat & Directorate for Education and Human Resources, Division of Research and Learning in Formal and Informal Settings, National Science Foundation.
- Groves, R., Flower, F., Couper, M., Lepkowski, J., Singer, E., & Tourangeau, R. (2004). *Survey methodology*. Hoboken, NJ: John Wiley & Sons.
- Harkness, J., Braun, M., Edwards, B., & Johnson, T. P. (2010). *Survey methods in multinational, multiregional, and multicultural contexts*. Hoboken, NJ: John Wiley & Sons.

- Hedges, L. V., & Rhoads, C. (2010). *Statistical power analysis in education research* (NCSE 2010-3006). Washington, DC: U.S. Department of Education, Institute for Education Sciences. Retrieved from: <http://ies.ed.gov/ncser/pubs/20103006/pdf/20103006.pdf>
- Hinkle, D. E., Wiersma, W., & Jurs, S. G. (2003). *Applied statistics for the behavioral sciences* (5th ed.). Boston: Houghton Mifflin.
- Hox, J. J. (2010). *Multilevel analysis: Techniques and applications* (2nd ed.), (Quantitative Methodology Series). New York: Routledge.
- Hulleman, C. S., & Cordray, D. S. (2009). Moving from the lab to the field: The role of fidelity and achieved relative intervention strength. *Journal of Research on Educational Effectiveness*, 2(1), 88-110.
- Kennedy, C. H. (2005). *Single-case designs for educational research*. Boston: Pearson Education.
- Kratochwill, T. R., Hitchcock, J. H., Horner, R. H., Levin, J. R., Odom, S. L., Rindskopf, D. M. & Shadish, W. (2010). *Single-case design technical documentation, Version 1.0 (Pilot)*. Washington, DC: U.S. Department of Education, Institute for Education Sciences. Retrieved from [http://ies.ed.gov/ncee/wwc/Docs/ReferenceResources/wwc\\_scd.pdf](http://ies.ed.gov/ncee/wwc/Docs/ReferenceResources/wwc_scd.pdf)
- Kratochwill, T. R., Hitchcock, J. H., Horner, R. H., Levin, J. R., Odom, S. L., Rindskopf, D. M. & Shadish, W. (2013). Single-case intervention research design standards. *Remedial and Special Education*, 34, 26-38.
- Kratochwill, T. R., & Levin, J. R. (2010). Enhancing the scientific credibility of single-case intervention research: Randomization to the rescue. *Psychological Methods*, 15, 122-144.
- Maxwell, J. A. (2005). *Qualitative research design: An interactive approach* (2nd ed.). Applied Social Research Methods Series, Vol. 41. Thousand Oaks, CA: Sage.
- Maxwell, J. A., & Miller, B. A. (2010). Categorizing and connecting strategies in qualitative data analysis. In S. N. Hesse-Biber & P. Leavy. (Eds.). *Handbook of emergent methods*. New York: The Guilford Press.
- Meyers, C., & Brand, W. C. (2015). *Implementation fidelity in education research: Designer and evaluator considerations*. New York: Routledge.
- Nelson, M. C., Cordray, D. S., Hulleman, C. S., Darrow, C. L., & Sommer, E. C. (2012). A procedure for assessing intervention fidelity in experiments testing educational and behavioral interventions. *Journal of Behavioral Health Services & Research*, 39(4), 374-396.
- Patton, M. Q. (2002). *Qualitative research and evaluation methods* (3rd ed.) Thousand Oaks, CA: Sage.
- Raudenbush, S. W., et al. (2011). *Optimal Design Software for multi-level and longitudinal research* (Version 3.01) [Software]. Retrieved from: <http://hlmssoft.net/od/>
- Shadish, W. R., Cook, T. D., & Campbell, D. T. (2002). *Experimental and quasi-experimental designs for generalized causal inference*. Belmont, CA: Wadsworth.
- Tashakkori, A., & Teddlie, C. (1998). *Mixed methodology: Combining qualitative and quantitative approaches*. Thousand Oaks, CA: Sage.
- Todman, J. B., & Dugard, P. (2001). *Single-case and small-n experimental designs: A practical guide to randomization tests*. Mahwah, NJ: Lawrence Erlbaum Associates

- U.S. Department of Education. (2013). *What Works Clearinghouse Procedures and standards handbook* (Version 3.0). Washington, DC: U.S. Department of Education, Institute for Education Sciences. Retrieved from <http://ies.ed.gov/ncee/wwc/DocumentSum.aspx?sid=19#>
- U.S. General Accounting Office, GAO. (2012). *Designing evaluations: 2012 revision* (Report No. GAO-12-208G). Washington, DC: Author. Retrieved from <http://www.gao.gov/assets/590/588146.pdf>
- U.S. General Accounting Office, GAO. (1992). *Quantitative data analysis: An introduction* (Report to Program Evaluation and Methodology Division, Report No. GAO/PEMD-10.1.11). Washington, DC: Author. Retrieved from <http://www.gao.gov/special.pubs/pe10111.pdf>
- Willis, G. (2005). *Cognitive interviewing: A tool for improving questionnaire design*. Thousand Oaks: Sage.



## Appendix E. Works Cited

- Allison, P. D. (2001). *Missing data*. Thousand Oaks, CA: Sage Publications.
- American Psychological Association. (2009). *Publication manual of the American Psychological Association* (6th ed.). Washington, DC: Author.
- Baecher, L. H., & Connor, D. J. (2010). "What do you see?" Using video analysis of classroom practice in a preparation program for teachers of students with learning disabilities. *Insights on Learning Disabilities*, 7(2), 5-18.
- Bailey, D. B., & Simeonsson, R. J., Jr. (1988). Investigation of use of goal attainment scaling to evaluate individual progress of clients with severe and profound mental retardation. *Mental Retardation*, 26, 289-295.
- Bandura, A. (2006). Guide for constructing self-efficacy scales. In T. Urdan & F. Pajares. (Eds.), *Self-efficacy beliefs of adolescents*. Greenwich, CT: Information Age Publishing.
- Baraldi, A. N., & Enders, C. K. (2010). An introduction to modern missing data analyses. *Journal of School Psychology*, 48(1), 5-37. DOI: 10.1016/j.jsp.2009.10.001
- Barnett, D. W., Daly, E. J. III, Hampshire, E. M., Hines, N. R., Maples, K. A., Ostrom, J. K., & Van Buren, A. E. (1999). Meeting performance-based training demands: Accountability in an intervention-based practicum. *School Psychology Quarterly*, 14, 357-379.
- Barth, R. P., Guo, S., & McCrae, J. S. (2008). Propensity score matching strategies for evaluating the success of child and family service programs. *Research on Social Work Practice*, 18(3), 212-222. DOI: 10.1177/1049731507307791
- Barton, E. E., & Reichow, B. (2012). Guidelines for graphing data with Microsoft PowerPoint for Office 2007. *Journal of Early Intervention*, 34, 129-150.
- Berkowitz, S. (1997). Analyzing qualitative data. In J. Frechtling & L. Sharpe (Eds.). *User-friendly handbook for mixed method evaluations*. Arlington, VA: National Science Foundation. Retrieved from: <http://www.nsf.gov/pubs/1997/nsf97153/>
- Bill and Melinda Gates Foundation. (2012). *Gathering feedback for teaching: Combining high-quality observations with student surveys and achievement gains*. Retrieved April 5, 2016, from: <http://files.eric.ed.gov/fulltext/ED540960.pdf>
- Bill and Melinda Gates Foundation. (2011). *Learning about teaching: Initial findings from the Measures of Effective Teaching Project*. Retrieved April 5, 2016, from: <https://docs.gatesfoundation.org/Documents/preliminary-findings-research-paper.pdf>
- Boller, K., Atkins-Burnett, S., Malone, L. M., Baxter, G. P., & West, J. (2010). *Compendium of student, teacher, and classroom measures used in NCEE evaluations of educational interventions. Volume I. Measures selection approaches and compendium development methods* (NCEE 2010-4012). Washington, DC: U.S. Department of Education, Institute of Education Sciences, National Center for Education Evaluation and Regional Assistance.
- Bovend'Eerdt, T. J. H., Botell, R. E., & Wade, D. T. (2009). Writing SMART rehabilitation goals and achieving goal attainment scaling: A practical guide. *Clinical Rehabilitation*, 23, 352-361.
- Brewer, J., & Hunter, A. (2006). *Foundations of multimethod research: Synthesizing styles*. Thousand Oaks, CA: Sage.

- Bruder, M. B., & Stayton, V. (2004). *The Center to Improve Personnel Preparation Policy and Practice in Early Intervention and Preschool Education, Briefing Book*. Farmington, CT: University of Connecticut. Retrieved from: [http://uconnucedd.org/wp-content/uploads/sites/1340/2015/06/PersonnelPrep\\_Briefing\\_Book.pdf](http://uconnucedd.org/wp-content/uploads/sites/1340/2015/06/PersonnelPrep_Briefing_Book.pdf)
- Bryington, A. A., Palmer, D. J., & Watkins, M. W. (2002). The estimation of interobserver agreement in behavioral assessment. *The Behavior Analyst Today*, 3(3), 323-328.
- Bryk, A. S., & Raudenbush, S. W. (1988). Toward a more appropriate conceptualization of research on school effects: A three-level hierarchical linear model. *American Journal of Education*, 97(1), 65-108.
- Busk, P. L., & Serlin, R. C. (1992). Meta-analysis for single-case research. In T. R. Kratochwill & J. R. Levin (Eds.), *Single-case research designs and analysis: New directions for psychology and education* (pp. 187-212). Hillsdale, NJ: Lawrence Erlbaum Associates
- Campbell, J. M. (2004). Statistical comparison of four effect sizes for single-subject designs, *Behavior Modification*, 28(2), 234-246. DOI: 10.1177/0145445503259264
- Cardillo, J. E., & Smith, A. (1994). Reliability of goal attainment scores. In T. J. Kiresuk, A. Smith, & J. Cardillo (Eds.), *Goal attainment scaling: Applications, theory and measurement* (pp. 173-212). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Carr, R. A. (1979). Goal attainment scaling as a useful tool for evaluating progress in special education. *Exceptional Children*, 46, 88-95.
- Century, J., Rudnick, M., & Freeman, C. (2010). A framework for measuring fidelity of implementation: A foundation for shared language and accumulation of knowledge. *American Journal of Evaluation*, 34(4): 465-485.
- Cochran, D. C., Gallagher, P. A., Stayton, V. D., Dinnebeil, L. A., Lifter, K., Chandler, L., K., et al. (2012). Early childhood special education and early intervention personnel preparation standards of the division for early childhood: Field validation. *Topics in Early Childhood Special Education*, 38(1), 38-51. DOI: 10.1177/0271121412436696
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale, NJ: Lawrence Earlbaum Associates.
- Coladarci, T., & Breton, W. A. (1997). Teacher efficacy, supervision, and the special education resource-room teacher. *The Journal of Education Research*, 90(4), 230-239. Retrieved from: <http://www.jstor.org/stable/27542097>
- Conrad, K. J., & Conrad, K. (1994). Reassessing validity threats in experiments: Focus on construct validity. In K. Conrad. (Ed.), *Critically evaluating the role of experiments* (pp.5-25). San Francisco: Jossey-Bass.
- Council for Exceptional Children. (2012). *The Council for Exceptional Children's position on special education teacher evaluation*. Arlington, VA: Author.
- Creswell, J. W. (2003). *Research design: Qualitative, quantitative, and mixed methods approaches* (2nd ed.). Thousand Oaks, CA: Sage.
- Creswell, J. W., Plano Clark, V., Gutmann, M., & Hanson, W. (2003). Advances in mixed method design. In A. Tashakkori & C. Teddlie (Eds.), *Handbook of mixed methods in the social and behavioral sciences*. Thousand Oaks, CA: Sage.
- Czaja, R., & Blair, J. (2005). *Designing surveys: A guide to Decisions and Procedures* (2nd ed.). Thousand Oaks, CA: Sage.

- Denzin, N. K., & Lincoln, Y. S. (Eds.). (2005). *The Sage handbook of qualitative research* (3rd ed.). Thousand Oaks, CA: Sage.
- Dimitrov, D. M. (2012). *Statistical methods for validation of assessment scale data in counseling and related fields*. Alexandria, VA: American Counseling Association.
- Dimitrov, D. M. (2010). *Quantitative research in education: Intermediate and advanced methods* (2nd ed.). New York: Whittier Publications.
- Dimitrov, D. M., Jurich, S., Frye, M., Lammert, J. D., Sayko, S., Taylor, L. Q. (2012, February). *Year one evaluation report/impact study: Illinois Striving Readers*. Arlington, VA: RMC Research Corporation.
- Dong, N., & Maynard, R. A. (2013). PowerUp!: A tool for calculating minimum detectable effect sizes and minimum required sample sizes for experimental and quasi-experimental design studies. *Journal of Research on Educational Effectiveness*, 6(1), 24-67. DOI: 10.1080/19345747.2012.673143
- Donmoyer, R. (2012a). Attributing causality in qualitative research: Viable option or inappropriate aspiration? An introduction to a collection of papers. *Qualitative Inquiry*, 18(8), 651-654. DOI: 10.1177/1077800412455012
- Donmoyer, R. (2012b). Can qualitative researchers answer policymakers' What-Works question? *Qualitative Inquiry*, 18(8), 662-673. DOI: 10.1177/1077800412454531
- Doran, G. T. (1981). There's a S.M.A.R.T. way to write management's goals and objectives. *Management Review*, 70, 35-36.
- Dymond, S. K., & Bentz, J. L. (2006). Using digital videos to enhance teacher preparation. *Teacher Education and Special Education: The Journal of the Education Division of the Council for Exceptional Children*, 29(2), 98-112. DOI: 10.1177/088840640602900202
- Ellis, P. D. (2010). *The essential guide to effect sizes: An introduction to statistical power, meta-analysis and the interpretation of research results*. United Kingdom: Cambridge University Press.
- Evans, H. M. (1981). Psychiatric patient participation in goal setting as related to goal attainment. *Dissertation Abstracts International*, 41(12-B, Pt 1), 4658-4659.
- Fiore, T., Helba, C., Berkowitz, S., Jones, M., & Newsome, J. (2012). *Making personnel development effective: Using outcome data for program improvement*. Presented at the U.S. Department of Education Office of Special Education Programs Project Directors Conference, July 23, 2012, Washington, DC.
- Fixsen, D. L., Naoom, S. F., Blase, K. A., Friedman, R. M., & Wallace, F. (2005). *Implementation research: A synthesis of the literature*. Retrieved from <http://ctndisseminationslibrary.org/PDF/nirmonograph.pdf>
- Frechtling, J. A. (2007). *Logic modeling methods in program evaluation*. San Francisco: Jossey-Bass.
- Freeman, M. (2000). Knocking on doors: On constructing culture. *Qualitative Inquiry*, 6, 359-369.
- Glover, S., Burns, J., & Stanley, B. (1994). Goal attainment scaling as a method of monitoring the progress of people with severe learning disabilities. *Mental Handicap*, 22, 148-150.
- Goldhaber, D. (2013). What do value-added measures of teacher preparation programs tell us? The Carnegie Knowledge Network. Retrieved from: [http://www.carnegieknowledgenetwork.org/briefs/teacher\\_prep/](http://www.carnegieknowledgenetwork.org/briefs/teacher_prep/)

- Greene, J. C. (2006). Toward a methodology of mixed methods social inquiry. *Research in the Schools, 13*(1), 93-98.
- Harms, T., Clifford, R. M., & Cryer, D. (2015). *The early childhood environment rating scale* (Third Ed.). New York: Teachers College Press.
- Hedges, L. V., & Rhoads, C. (2010). *Statistical power analysis in education research* (NCSE 2010-3006). Washington, DC: U.S. Department of Education Institute for Education Sciences. Retrieved on April 5, 2016, from: <http://ies.ed.gov/ncser/pubs/20103006/pdf/20103006.pdf>
- Heinrich, C., Maffioli, A., & Vazquez, G. (2010). *A primer for applying propensity-score matching: Impact evaluation guidelines* (Technical Notes No. IDB-TN-161). Washington, DC: Inter-American Development Bank.
- Heinemeier, S., D'Agostino, A., Lammert, J.D., & Fiore, T.A. (2014). *Guidelines for Working with Third-Party Evaluators*. Rockville, MD: Westat.
- Hinkle, D. E., Wiersma, W., & Jurs, S. G. (2003). *Applied statistics for the behavioral sciences* (5th ed.). Boston: Houghton Mifflin.
- Holbrook, A. L. (2013). *Agree-disagree questions: Problems and some solutions*. Webinar presented as part of the American Association for Public Opinion Research (AAPOR) Online Continuing Education Program, October 23, 2013.
- Holdheide, L. R., Goe, L., Croft, A., & Reschly, D. J. (2010). *Challenges in evaluating special education teachers and English language learner specialists*. Washington, DC: National Comprehensive Center for Teacher Quality.
- Horner, R., & Spaulding, S. (2010). Single-case research designs (pp. 1386–1394). In N. J. Salkind (Ed.), *Encyclopedia of research design*. Thousand Oaks, CA: Sage.
- Howell, D. C. (2012, December 12). *Treatment of missing data — Part 1*. [Web log post]. Retrieved from: [https://www.uvm.edu/~dhowell/StatPages/Missing\\_Data/Missing.html](https://www.uvm.edu/~dhowell/StatPages/Missing_Data/Missing.html)
- Hox, J. J. (2010). *Multilevel analysis: Techniques and applications* (2nd ed.), (Quantitative Methodology Series). New York: Routledge.
- Hulleman, C. S., & Cordray, D. S. (2009). Moving from the lab to the field: The role of fidelity and achieved relative intervention strength. *Journal of Research on Educational Effectiveness, 2*(1), 88-110.
- Jerald, C. (2012). *Ensuring accurate feedback from observations: Perspectives on practice*. Seattle: Bill & Melinda Gates Foundation. Retrieved April 5, 2016, from: <http://www.gatesfoundation.org/college-ready-education/Documents/ensuring-accuracy-wp.pdf>
- Kane, M. T. (2001). Current concerns in validity theory. *Journal of Educational Measurement, 38*(4), 319-342.
- Kelley, K., & Preacher, K. J. (2012). On effect size. *Psychological Methods, 17*(2): 137-152. DOI: 10.1037/a0028086.
- Kennedy, C. H. (2005). *Single-case designs for educational research*. Boston: Allyn and Bacon.
- Kiresuk, T. J., & Sherman, R. E. (1968). Goal attainment scaling: A general method for evaluating comprehensive community mental health programs. *Community Mental Health Journal, 4*(6), 443-453.

- Kiresuk, T. J., Smith, A., & Cardillo, J. E. (Eds.) (1994). *Goal attainment scaling: Applications, theory, and measurement*. Hillsdale, NJ: Lawrence Erlbaum Associates.
- Kratochwill, T. R., Elliott, S. N., & Rotto, P. C. (1995). Best practices in school-based behavioral consultation. In A. Thomas & J. Grimes (Eds.), *Best practices in school psychology-III* (pp. 519-538). Washington, DC: UNational Association of School Psychologists.
- Kratochwill, T. R., Hitchcock, J. H., Horner, R. H., Levin, J. R., Odom, S. L., Rindskopf, D. M., & Shadish, W. (2010). *Single-case design technical documentation, Version 1.0 (Pilot)*. Washington, DC: Institute for Education Sciences. Retrieved from: [http://ies.ed.gov/ncee/wwc/Docs/ReferenceResources/wwc\\_scd.pdf](http://ies.ed.gov/ncee/wwc/Docs/ReferenceResources/wwc_scd.pdf)
- Kratochwill, T. R., Hitchcock, J. H., Horner, R. H., Levin, J. R., Odom, S. L., Rindskopf, D. M., & Shadish, W. (2013). Single-case intervention research design standards. *Remedial and Special Education, 34*, 26-38.
- Kratochwill, T. R., & Levin, J. R. (2010). Enhancing the scientific credibility of single-case intervention research: Randomization to the rescue. *Psychological Methods, 15*, 122-144.
- Lammert, J. D., & Feldman, J. M. (2015). *Measuring fidelity as a method for evaluating intensive technical assistance*. Presented at the OSEP Technical Assistance Coordination Center Workshop, September 1, 2015, Washington, DC.
- Lammert, J. D., & Fiore, T. A. (2015). Budgeting for evaluation: Key factors to consider. Rockville, MD: Westat.
- Lammert, J. D., Heinemeier, S., Schaaf, J. M., & Fiore, T.A. (2016). *Evaluating special education preservice programs: Resource Toolkit*. Rockville, MD: Westat.
- Lazarus, S. S., & Heritage, M. (2016). *Lessons learned about assessment from inclusion of students with disabilities in college and career ready assessments*. Minneapolis, MN: National Center on Educational Outcomes and the National Center on Systemic Improvement.
- Lipsey, M. W., Puzio, K., Yun, C., Hebert, M. A., Steinka-Fry, K., Cole, M. A., et al. (2012). *Translating the statistical representation of the effects of education interventions into more readily interpretable forms* (NCSE 2013-3000). Washington, DC: U.S. Department of Education, Institute of Education Sciences. Retrieved from: <http://ies.ed.gov/ncser/pubs/20133000/pdf/20133000.pdf>
- Luellen, J. K., Shadish, W. R., & Clark, M. H. (2005). Propensity scores: An introduction and experimental test. *Evaluation Review, 29*(6), 530-558. DOI: 10.1177/0193841X05275596
- Maher, C. A. (1983). Goal attainment scaling: A method for evaluating special education services. *Exceptional Children, 49*, 529-536.
- Maxwell, J. A. (2005). *Qualitative research design: An interactive approach* (2nd ed.) (Applied Social Research Methods Series, Vol. 41). Thousand Oaks: Sage.
- Maxwell, J. A. (2012). The importance of qualitative research for causal explanation in education. *Qualitative Inquiry, 18*(8), 655-661. DOI: 10.1177/1077800412452856
- Maxwell, J. A. (2011). *A realist approach for qualitative research*. Thousand Oaks, CA: Sage.
- Maxwell, J. A. (2004a). Causal explanation, qualitative research, and scientific inquiry in education. *Educational Researcher, 33*(2), 3-11.

- Maxwell, J. A. (2004b). Using qualitative methods for causal explanation. *Field Methods*, 16, 243-264.
- McIntosh, R., Vaughn, S., Schumm, J. S., Haager, D., & Lee, O. (1993). Observations of students with learning disabilities in general education classrooms. *Exceptional Children*, 60(3), 249-261.
- Messick, S. (1989). Validity. In R. L. Linn (Ed.), *Educational measurement* (3rd ed., pp 13-103). Phoenix, AZ: The Oryx Press.
- Meyers, C., & Brand, W. C. (2015). *Implementation fidelity in education research: Designer and evaluator considerations*. New York: Routledge.
- Miles, M. B., & Huberman, A. M. (1994). *Qualitative data analysis: An expanded sourcebook* (2nd ed.). Thousand Oaks: Sage.
- Mitchell, T. & Cusick, A. (1998). Evaluation of a client-centered pediatric rehabilitation programme using goal attainment scaling. *Australian Occupational Therapy Journal*, 45, 7-17.
- Morgan, D. L. (2007). Paradigms lost and pragmatism regained: Methodological implications of combining qualitative and quantitative methods. *Journal of Mixed Method Research*, 1(48), 48-76. DOI: 10.1177/2345678906292462
- Morrison, J. Q. (2012). Goal attainment scaling and single-case designs for service delivery accountability: A presentation and dialogue with the OSEP 325k Work Group. Webinar presentation May 2, 2012.
- Nelson, M. C., Cordray, D. S., Hulleman, C. S., Darrow, C. L., & Sommer, E. C. (2012). A procedure for assessing intervention fidelity in experiments testing educational and behavioral interventions. *Journal of Behavioral Health Services & Research*, 39(4), 374-396.
- Office of Planning, Research and Evaluation Administration for Children and Families (2010). The Program manager's guide to evaluation. Washington, DC: Author. Retrieved June 16, 2016 from [http://www.acf.hhs.gov/sites/default/files/opre/program\\_managers\\_guide\\_to\\_eval2010.pdf](http://www.acf.hhs.gov/sites/default/files/opre/program_managers_guide_to_eval2010.pdf)
- Olive, M. L., & Smith, B. W. (2005). Effect size calculations and single subject designs. *Educational Psychology*, 25, 313-324.
- Oren, T., & Ogletree, B. T. (2000). Program evaluation in classrooms for students with autism: Student outcomes and program processes. *Focus on Autism and Other Developmental Disabilities*, 15, 170-175.
- Patton, M. Q. (2002). *Qualitative research and evaluation methods* (3rd ed.) Thousand Oaks, CA: Sage.
- Pianta, R. C., La Paro, K. M., & Hamre, B. K. (2008). *Classroom Assessment Scoring System (CLASS)*. Baltimore, MD: Paul H. Brookes.
- Raudenbush, S. W. (1997). Statistical analysis and optimal design for cluster randomized trials. *Psychological Methods*, 2(2), 173-185.
- Raudenbush, S. W., & Liu, X. (2000). Statistical power and optimal design for multisite randomized trials. *Psychological Methods*, 5(2), 199-213.
- Raudenbush, S. W., et al. (2011). *Optimal Design Software for Multi-level and Longitudinal Research (Version 3.01)* [Software]. Retrieved from: <http://hlmsoft.net/od/>

- Raudenbush, S. W., Martinez, A., & Spybrook, J. (2007). Strategies for improving precision in group-randomized experiments. *Educational Evaluation and Policy Analysis*, 29(1): 5-29.
- Reybold, L. E., Lammert, J. D., & Stribling, S. M. (2012, November 30). Participant selection as a conscious research method: Thinking forward and the deliberation of 'emergent' findings. *Qualitative Research*, DOI: 1468794112465634
- Roach, A. T., & Elliott, S. N. (2005). Goal attainment scaling: An efficient and effective approach to monitoring student progress. *Teaching Exceptional Children*, 37(4), 8-17.
- Saris, W., Revilla, M., Krosnick, J. A., & Shaeffer, E. M. (2010). Comparing questions with agree/disagree response options to questions with item-specific response options. *Survey Research Methods*, 4, 61-79.
- Schutt, R. K. (2011). Chapter 10. Qualitative data analysis. *Investigating the social world: The process and practice of research* (7th ed.). Thousand Oaks: Sage. Retrieved from: [http://www.sagepub.com/upm-data/43454\\_10.pdf](http://www.sagepub.com/upm-data/43454_10.pdf)
- Schochet, P. Z. (2009). *Technical methods report: The estimation of average treatment effects for clustered RCTs of education interventions* (NCEE 2009-0061 rev.). Washington, DC: U.S. Department of Education, Institute of Education Sciences.
- Schochet, P. Z. (2008). Statistical power for randomized assignment evaluation of education programs. *Journal of Educational and Behavioral Statistics*, 33(1), 62-87.
- Scriven, M. (2008). A summative evaluation of RCT methodology: And an alternative approach to causal research. *Journal of Multidisciplinary Evaluation*, 5(9), 11-24.
- Scruggs, T. E., Mastropieri, M. A., & Casto, G. (1987). The quantitative synthesis of single subject research methodology: Methodology and validation. *Remedial and Special Education*, 8, 24-33.
- Scruggs, T. E., Mastropieri, M. A., Cook, S. B., & Escobar, C. (1986). Early intervention for children with conduct disorders: A quantitative synthesis of single-subject research. *Behavioral Disorders*, 11, 260-271.
- Shadish, W. R., Cook, T. D., & Campbell, D. T. (2002). *Experimental and quasi-experimental designs for generalized causal inference*. Belmont, CA: Wadsworth.
- Shadish, W. R., Hedges, L. V., Horner, R. H., & Odom, S. L. (2015). *The role of between-case effect size in conducting, interpreting, and summarizing single-case research*, (NCER 2015-002). Washington, DC: National Center for Education Research, Institute of Education Sciences, U.S. Department of Education. Available at <http://ies.ed.gov/>
- Siegel, S., & Castellan, N. J., Jr. (1988). *Nonparametric statistics for the behavioral sciences* (2nd ed.). Singapore: McGraw-Hill.
- Sladeczek, I. E., Elliott, S. N., Kratochwill, T. R., Robertson-Mjaanes, S., & Stoiber, K. (2001). Application of goal attainment scaling to a conjoint behavioral consultation case. *Journal of Education and Psychological Consultation*, 12(1), 45-58.
- Smith, A., & Cardillo, J. E. (1994). Perspectives on validity. In T. J. Kiresuk, A. Smith, & J. E. Cardillo (Eds.), *Goal attainment scaling: Applications, theory, and measurement* (pp. 243-272). Hillsdale, NJ: Erlbaum.
- Spybrook, J., Raudenbush, S. W., Congdon, R., & Martinez, A. (2009). Optimal design version 2.0 [Software].

- St. John, W. & Johnson, P. (2000). The pros and cons of data analysis software for qualitative research. *Journal of Nursing Scholarship: An Official Publication of Sigma Theta Tau International Honor Society of Nursing/Sigma Theta Tau*, 32(4): 393–397. [PMID11140204](#).
- Sullivan, G. M. (2011). Getting off the “gold standard”: Randomized controlled trials and education research. *Journal of Graduate Medical Education*, 3(3), 285-289. DOI: <http://dx.doi.org/10.4300/JGME-D-11-00147.1>
- Suter, W. N. (2012). Chapter 12. Qualitative data, analysis, and design. *Introduction to educational research: A critical thinking approach* (2nd ed.). Thousand Oaks: Sage. Retrieved from: [http://www.sagepub.com/upm-data/43144\\_12.pdf](http://www.sagepub.com/upm-data/43144_12.pdf)
- Tashakkori, A., & Creswell, J. W. (2007). Editorial: The new era of mixed methods. *Journal of Mixed Methods Research*, 1(3), 3-7. DOI: 10.1177/2345678906293042
- Tashakkori, A., & Teddlie, C. (2003). *Handbook of mixed methods in social and behavioral research*. Thousand Oaks, CA: Sage.
- Tashakkori, A., & Teddlie, C. (1998). *Mixed methodology: Combining qualitative and quantitative approaches*. Thousand Oaks, CA: Sage.
- Turner-Stokes, L. (2009). Goal attainment scaling (GAS) in rehabilitation: A practical guide. *Clinical Rehabilitation*, 23, 362-370.
- U.S. Department of Education. (2014). *What Works Clearinghouse Procedures and standards handbook* (Version 3.0). Washington, DC: U.S. Department of Education, Institute for Education Sciences. Retrieved March 29, 2016, from: [https://ies.ed.gov/ncee/wwc/Docs/referenceresources/wwc\\_procedures\\_v3\\_0\\_standards\\_handbook.pdf](https://ies.ed.gov/ncee/wwc/Docs/referenceresources/wwc_procedures_v3_0_standards_handbook.pdf)
- U.S. Department of Education (2015, March). Teacher shortage areas nationwide listing: 1990-1991 through 2015-2016. Washington, DC: Author. Retrieved from: <http://www2.ed.gov/about/offices/list/ope/pol/tsa.pdf>
- Wade, D. T. (2009). Editorial. Goal setting in rehabilitation: An overview of what, why and how. *Clinical Rehabilitation*, 23(4), 291-295.
- Watkins, M. W. & Pacheco, M. (2000). Interobserver agreement in behavioral research: Importance and calculation. *Journal of Behavioral Education*, 10(4), 205-212.
- Weisberg, D., Sexton, S., Mulhern, J., & Keeling, D. (2009). *The widget effect: Our national failure to act on differences in teacher effectiveness*. Brooklyn, NY: The New Teacher Project.
- Weiss, R. S. (1994). *Learning from strangers: The art and method of qualitative interview studies*. New York: The Free Press.
- Wendt, O. (2009). *Calculating effect sizes for single-subject experimental designs: An overview and comparison*. Presented at the Ninth Annual Campbell Collaboration Colloquium, Oslo, Norway, May 18, 2009.
- West, S. G., & Hepworth, J. T. (1991). Statistical issues in the study of temporal data: Daily experiences. *Journal of Personality*, 59, 609-662.
- Yin, R. K. (1994). *Case study research: Design and methods* (2nd ed., Applied Social Research Methods Series, v. 5.). Thousand Oaks, CA: SAGE.