



What Works Ireland Evidence Hub

Evidence standards and common pitfalls in
evaluation design

Evidence standards and common pitfalls in evaluation design

Level 2

Level 3

Pitfalls to evaluation success and how to avoid them

Level 4

Next steps

The strength of evidence rating: Level 2

Level 2

Preliminary evidence

Research Design

- Pre-post studies
- Using validated measures
- Involving a sample of min. 20 participants
- Using appropriate analysis methods



The strength of evidence rating: Level 3

Level 3

Efficacy

Research Design

- Controlled studies (RCT or QED)
- Using appropriate measures
- Involving a sample of min. 20 participants in each study arm
- Using appropriate analysis methods

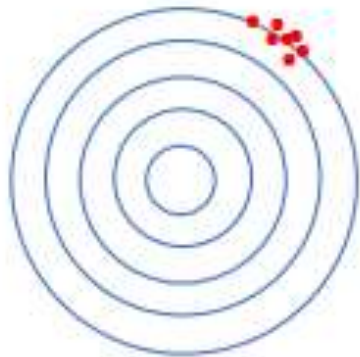


Common pitfalls in evaluation design

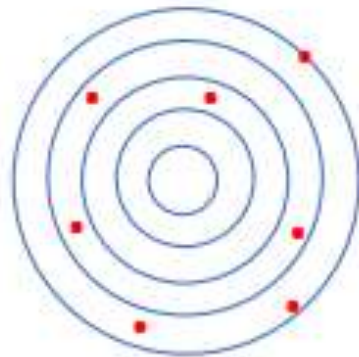
- Using inappropriate measures
- No robust comparison group
- Small sample size
- High drop-out rate
- Lack of long-term follow-up

Pitfall 1: Using inappropriate measures

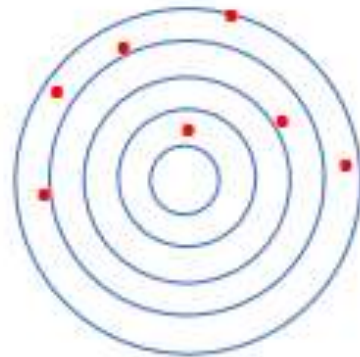
Reliable but not valid



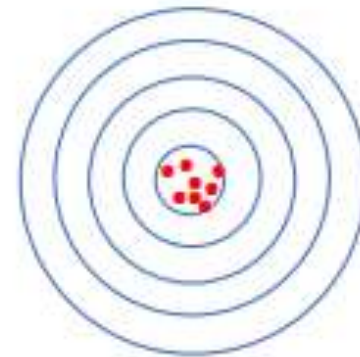
Valid but not reliable



Neither reliable nor valid



Both valid and reliable



The use of inappropriate measures
can mean
that a programme's evidence is
assessed as NL2

How to avoid this pitfall

- Sources of validated measures include:

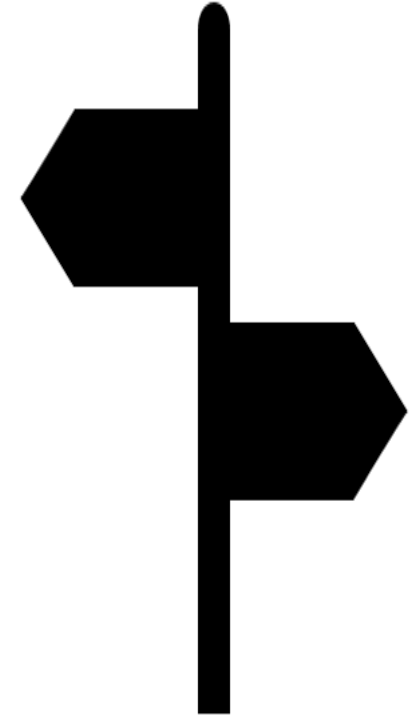
Source	Link
California Evidence-Based Clearinghouse: List of reviewed measures	http://www.cebc4cw.org/assessment-tools/measurement-tools-highlighted-on-the-cebc/
Early Intervention Foundation (EIF): Foundations for Life report	http://www.eif.org.uk/publication/foundationsfor-life-what-works-to-support-parent-childinteraction-in-the-early-years/
Early Intervention Foundation (EIF): Commissioner Guide: Reducing the impact of interparental conflict on children	https://www.eif.org.uk/files/pdf/cg-rpc-3-3-examples-validated-measures.pdf
Education Endowment Foundation (EEF): Early Years Measures database	https://educationendowmentfoundation.org.uk/projects-and-evaluation/evaluating-projects/evaluator-resources/early-years-measuredatabase/early-years-measures-database/
Education Endowment Foundation (EEF): Spectrum database	https://educationendowmentfoundation.org.uk/projects-and-evaluation/evaluating-projects/measuring-essential-skills/spectrum-database/
ETS Test Collection	https://www.ets.org/test_link/about
Deighton et al review (2014)	https://capmh.biomedcentral.com/articles/10.1186/1753-2000-8-14
Denham & Hamre review (2010)	https://pdfs.semanticscholar.org/e55c/3929969c5b1ffb35966ead41a08b1e040aaf.pdf

Pitfall 2: No robust comparison

To know what impact a programme has had, we need to know the outcomes of the participants who have received the programme.

But we also need to be able to estimate what would have happened to these participants if they had not received the programme.

This is known as the ‘counterfactual’.



Pitfall 2: No robust comparison

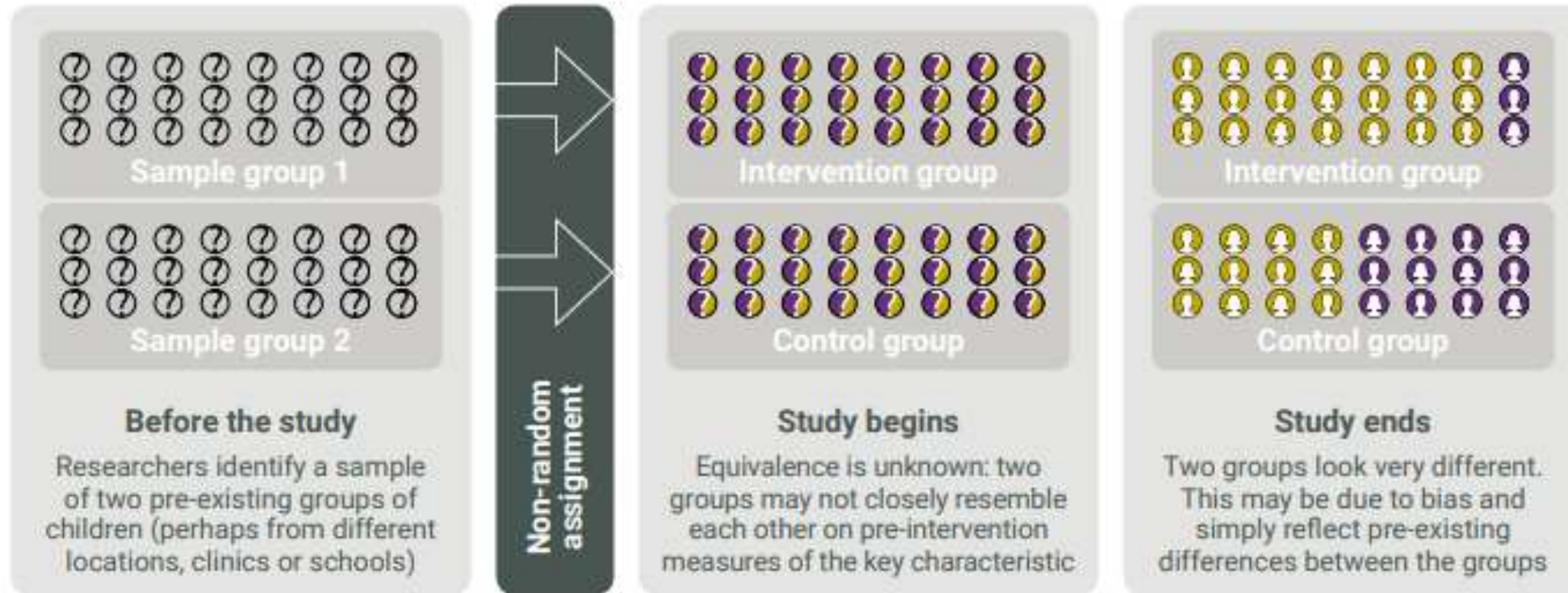
Some evaluations do not use a comparison group at all.

- In pre-post studies, it is always possible that factors other than the programme are responsible for any observed improvements

Pre-post studies can be assessed as a Level 2, if they meet other criteria, but they cannot get a higher rating

Controlled studies that don't use an appropriate mechanism to generate the comparison can be assessed as a Level 2

Pitfall 2: No robust comparison



How to avoid this pitfall

- **Experimental design** (Randomised controlled trials):
RCTs ensure equivalence by randomly assigning participants to the intervention and control groups, ensuring that there is no systematic differences between the two study groups on any characteristics (although they may differ by chance)
- **Quasi-experimental design:**
QEDs use statistical methods to generate a robust comparison – there are a variety of approaches, including Instrumental Variable designs; Difference-in-differences designs, Interrupted time series designs, etc.

Pitfall 3: Small sample size

	Truth is: coin is fair	Truth is: coin is biased
We conclude: coin is fair	Correct	Type II: false negative result
We conclude: coin is biased	Type I: false discovery	Correct

- **False negative results**
- **False discoveries**
- **Chance that groups are unequal**

The use of too sample sizes can mean that a programme's evidence is assessed as NL2 if affecting intervention and control group or Level 2, if affecting only the control group

How to avoid this pitfall

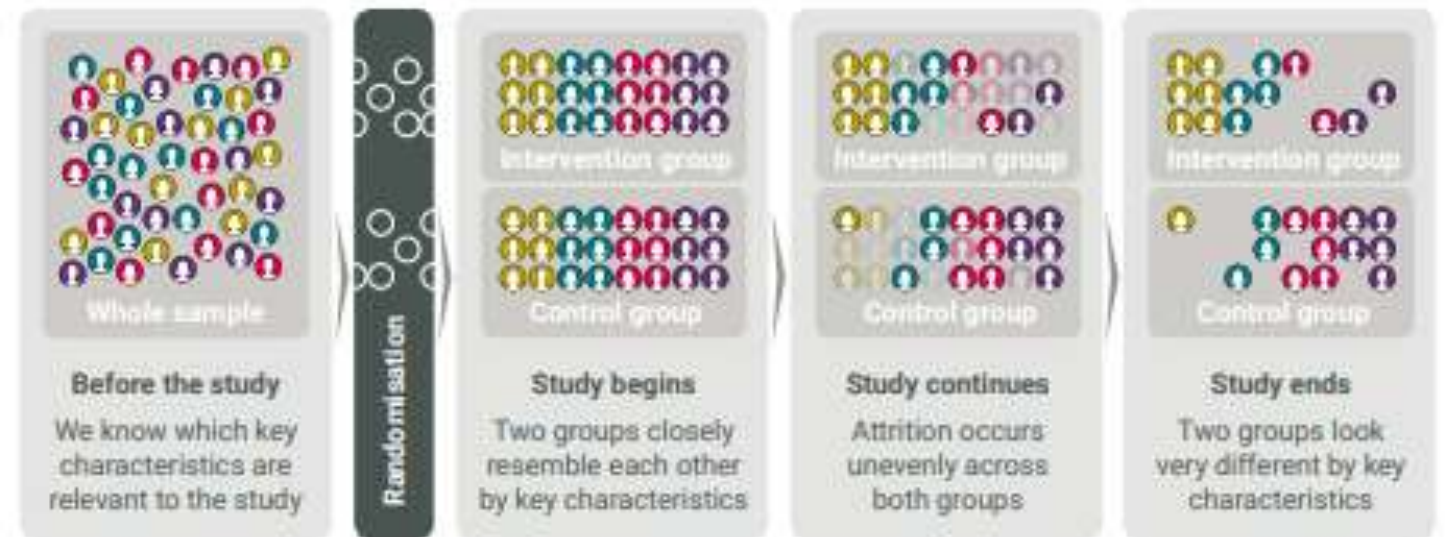
Researchers should attempt to recruit and retain appropriately large sample sizes

- A sufficiently large sample size depends on a range of considerations, including the design of the study and the size of the effect the researcher is attempting to identify.

Pitfall 4: High drop-out rate

- **Unrepresentativeness**
- **Bias and non-equivalent groups**

How attrition can introduce bias



How to avoid this pitfall

Researchers should aim to minimise attrition as far as possible.

- clear communication of the benefits of taking part in the research
- assigning research team members to follow-up with participants
- maintaining detailed contact information
- compensation, such as cash, vouchers or equivalent gifts

High attrition that is not dealt with adequately, can mean that a programme's evidence is assessed as Level 2 even where studies are controlled.

How to avoid this pitfall

Participant drop-out can rarely be prevented entirely. Therefore, researchers should examine their sample and conduct analyses on the extent to which it may have introduced bias.

- Attrition rates (the extent of attrition):
 - the overall attrition rate
 - the differential attrition rate
- Attrition type (the nature of attrition):
 - Differences between study drop-outs and completers
 - Whether attrition undermined the equivalence of the study groups

Pitfall 5: Lack of long-term follow-up

- Some outcomes may look better right after participants complete an intervention, but these effects may disappear quite quickly
- On the other hand, some outcomes might not have improved right after the intervention, but may improve after a while

How to avoid this pitfall

Researchers should plan to conduct follow-up assessments

High attrition that is not dealt with adequately, can mean that a programme's evidence is assessed as Level 2 even where studies are controlled.

Level 4

- Two or more Level 3 studies, with at least one Level 3 studies looking at long-term impact.

The strength of evidence rating

<p>NL 2 <i>Not level 2</i></p>	<ul style="list-style-type: none"> Case studies; qualitative research; lack of validated measures in impact evaluations
<p>Level 2 <i>Preliminary evidence</i></p>	<ul style="list-style-type: none"> Pre-post studies showing improved outcomes, but no comparison group is used, so lack of confidence regarding causal impact of intervention
<p>Level 3 <i>Efficacy</i></p>	<ul style="list-style-type: none"> Rigorous randomised controlled trial or quasi experimental design demonstrated that the intervention led to an improvement in child outcomes
<p><i>No Effects</i></p>	<ul style="list-style-type: none"> As Level 3, but finding so significant intervention effects
<p>Level 4 <i>Effectiveness</i></p>	<ul style="list-style-type: none"> Two or more Level 3 studies, demonstrating effects were replicated in more than one site – also demonstrating long-term effects and using independent measures



Key features of impact evaluations

- Study design (pre-post, RCT, QED?)
- Measurements
- Analysis
- Robust approach to attrition
- Timing of measures

10 STEPS FOR EVALUATION SUCCESS



Source: EIF