# Evidence-based Decision Making
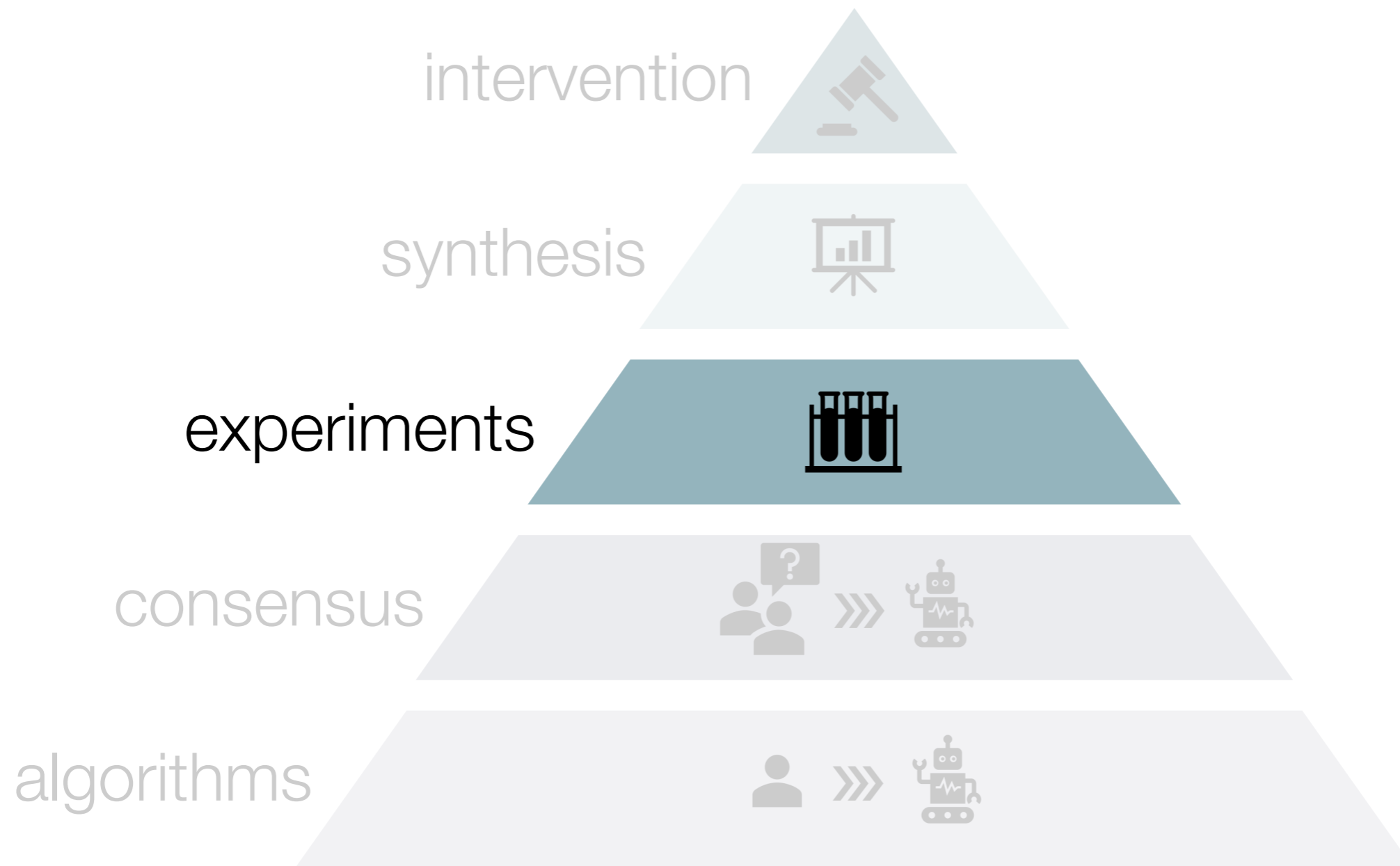## Counterfactuals: Experiments

Loreen Tisdall, FS 2024

Version: February 23, 2024

intervention

synthesis

experiments

consensus

algorithms

# Goals for today

- Understand the nature of causal inference as the comparison of treatment to some counterfactual

- Understand that experiments, and in particular RCTs, have desirable properties for causal inference – but also have limitations…

- Consider alternatives to RCTs to establish the counterfactual

# Causality

: a causal quality or <u>agency</u>

: the <u>relation between</u> a cause and its effect or between regulatory correlated events or phenomena

: someone or something <u>responsible for</u> a result
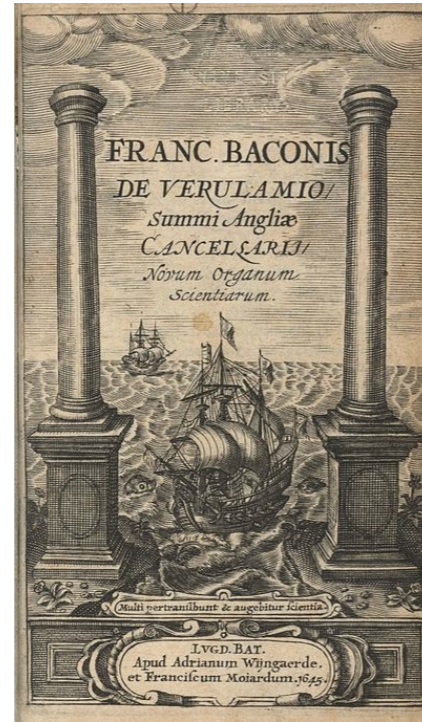
4

# Causal relations as counterfactual relations

*"D shoots at V, but only grazes him, leaving V with a slightly bleeding flesh wound. X then comes along and shoots V through the heart, killing him instantly. D's act is clearly not a "cause in fact" of V's death, since V would have died, and in just the manner he did, even if D had not shot him."*

- Singular judgment of causation: a single cause (e.g., A) is necessary and sufficient for effect (Y) to occur
- **In reality:** conjunctive plurality of causes (A&B&C → Y), disjunctive plurality of causes (A|B|C → Y)
- Complex regularities (e.g., A&B&C → Y) are rarely (if ever) fully known, thus we formulate propositions which entail the probability of a variable being causally connected with an effect

Marini, M. M., & Singer, B. (1988). Causality in the social sciences. *Sociological methodology, 18*, 347-409.

# Evidence-based decision making



Francis Bacon
(1561-1626)



1620



Bacon suggests that one can draw up a list of all things in which the phenomenon to explain occurs, as well as a list of things in which it does not occur. Then one can rank the lists according to the degree in which the phenomenon occurs in each one. Then one should be able to deduce what factors match the occurrence of the phenomenon in one list and do not occur in the other list, and also what factors change in accordance with the way the data had been ranked.

"The critical step in any causal analysis is estimating the counterfactual—a prediction of what would have happened in the absence of the treatment."

Varian, H. R. (2016). Causal inference in economics and marketing. *Proceedings of the National Academy of Sciences of the United States of America, 113*(27), 7310–7315. http://doi.org/10.1073/pnas.1510479113

# The gold standard…

## Experiments/Randomised control trials (RCT)

"To find out what happens when you change something, it is necessary to change it."
(Box et al., 2005)

A type of scientific experiment, where the people being studied are randomly allocated to one or other of the different treatments under study. RCTs are considered the gold standard for a clinical trial. RCTs are often used to test the *efficacy* or *effectiveness* of various types of medical intervention and may provide information about adverse effects, such as drug reactions. Random assignment of intervention is done after subjects have been assessed for eligibility and recruited, but before the intervention to be studied begins.

**Efficacy:**

**?**

**Effectiveness:**

# The gold standard…

## Experiments/Randomised control trials (RCT)

A type of scientific experiment, where the people being studied are randomly allocated to one or other of the different treatments under study. RCTs are considered the gold standard for a clinical trial. RCTs are often used to test the *efficacy* or *effectiveness* of various types of medical intervention and may provide information about adverse effects, such as drug reactions. Random assignment of intervention is done after subjects have been assessed for eligibility and recruited, but before the intervention to be studied begins.

**Efficacy:** how well a treatment/intervention works under ideal, controlled (laboratory) settings

**Effectiveness:** how well a treatment/intervention works in real-world (clinical) settings

Shorter, E. (2011). A brief history of placebos and clinical trials in psychiatry. *Canadian Journal of Psychiatry, 56*(4), 193–197.

# The gold standard…

## Experiments/Randomised control trials (RCT)

A type of scientific experiment, where the people being studied are randomly allocated to one or other of the different treatments under study. RCTs are considered the gold standard for a clinical trial. RCTs are often used to test the *efficacy* or *effectiveness* of various types of medical intervention and may provide information about adverse effects, such as drug reactions. Random assignment of intervention is done after subjects have been assessed for eligibility and recruited, but before the intervention to be studied begins.

$$Y = B_0 + B_1 group$$



Shorter, E. (2011). A brief history of placebos and clinical trials in psychiatry. *Canadian Journal of Psychiatry, 56*(4), 193–197.

# The gold standard…

Consolidated Standards of Reporting Trials

Schulz, K. F., Altman, D. G., & Moher, D. (2010). CONSORT 2010 statement: updated guidelines for reporting parallel group randomised trials. *Journal of Pharmacology and pharmacotherapeutics, 1*(2), 100-107..

# The Salk Polio Vaccine Trial & the Cutter Incident

• The 1954 Salk Polio vaccine trial was the largest RCT (a double-blind, randomized, and placebo-controlled study) ever conducted, involving over 1.8 million children, to test the safety and efficacy of a polio vaccine developed by Jonas Salk.

• The results showed that the vaccine was safe and effective in preventing polio.

• In 1955, shortly after the Salk polio vaccine was licensed, a manufacturing error at one of 5 licensed laboratories, Cutter Laboratories, resulted in the contamination of some batches of the vaccine with live polio virus, which led to an outbreak that affected a few hundred children, including some deaths and cases of permanent paralysis, known as the Cutter incident.

• The Cutter incident led to significant changes in vaccine regulation including the creation of oversight agencies and legislation.

→ The Cutter incident is an example of the problems that may arise from generalizing RCTs – and the continued need for evaluation (also their legal repercussions)…



A manufacturing error at Cutter Laboratories resulted in the contamination of some batches of the vaccine with live polio virus

Offit, P.A. (2005). The Cutter incident, 50 years later. *N Engl J Med*. 352, 1411-1412.

Dawson, L. (2004). The Salk polio vaccine trial of 1954: Risks, randomization and public involvement in research. *Clinical Trials, 1*, 122–130.

# The gold standard is not always gold…

Experiments/Randomised control trials (RCT)

- **Efficacy vs. effectiveness**: Trials  may not be widely applicable in real-world conditions….

- **Generalizability**: Results may not always generalize to other samples (e.g. inclusion /exclusion criteria)

- **Ethical limitations**: randomisation requires experimental equipoise: one cannot ethically randomise participants to some treatments (no-schooling condition)

# On the horizon: Autonomous Scientific Agents



**Figure 1. Overview of the system architecture.** The **Agent** is composed of multiple modules that exchange messages. Some of them have access to APIs, the Internet, and Python interpreter.

In this paper, we presented an Intelligent **Agent** system capable of autonomously designing, planning, and executing complex scientific experiments. Our system demonstrates exceptional reasoning and experimental design capabilities, effectively addressing complex problems and generating high-quality code.
However, the development of new machine learning systems and automated methods for conducting scientific experiments raises substantial concerns about the safety and potential dual use consequences, particularly in relation to the proliferation of illicit activities and security threats. By ensuring the ethical and responsible use of these powerful tools, we can continue to explore the vast potential of large language models in advancing scientific research while mitigating the risks associated with their misuse.

Boiko, D.A., MacKnight, R. & Gomes, G. L. (2023). Emergent autonomous scientific research capabilities of large language models.
https://doi.org/10.48550/arXiv.2304.05332

# There are alternatives…



Donald Campbell
1916-1996

THE CAMPBELL COLLABORATION

# Does education work?

**YOUR TURN!**

**How could you try to find out if education has an effect on intelligence?**

# Quasi-Experimental Designs: Educational effects on intelligence

*control prior intelligence* = longitudinal studies in which cognitive testing data were collected before and after variation in the duration of education (e.g., before and after university vs. no university)

*policy change* = study of the effects of a change in educational duration (e.g., increase of compulsory education by 1 year) on mental testing

*school-age cutoff* = studies use regression-discontinuity analysis to leverage the fact that school districts implement a date-of-birth cutoff for school entry (example: compare 3.9-year olds that did not attend "Kindsgi" vs. 4.0 year-olds that did)

**Table 1.** Descriptive Statistics for Each Study Design

| Design | Control prior intelligence | Policy change | School age cutoff |
|---|---|---|---|
| *k* studies | 7 | 11 | 10 |
| *k* data sets | 10 | 12 | 20 |
| *k* effect sizes | 26 | 30 | 86 |
| *N* participants | 51,645 | 456,963 | 107,204 |
| Mean age at early test in years (*SD*) | 12.35 (2.90) | — | — |
| Mean time lag between tests in years (*SD*) | 53.17 (15.47) | — | — |
| Mean age at policy change in years (*SD*) | — | 14.80 (2.59) | — |
| Mean age at outcome test in years (*SD*) | 63.48 (18.80) | 47.92 (19.39) | 10.36 (1.60) |
| *n* outcome test category (composite:fluid:crystallized) | 5:20:1 | 2:23:5 | 3:67:16 |
| *n* achievement tests (achievement:other) | 1:25 | 7:23 | 38:48 |
| Male-only estimates (male only:mixed sex) | 2:24 | 8:22 | 0:86 |
| Publication status (published:unpublished) | 22:4 | 21:9 | 64:22 |

Note: To estimate *N* from studies with multiple effect sizes with different *ns*, we averaged sample sizes across effect sizes within each data set and rounded to the nearest integer. "Unpublished" refers to any study not published in a peer-reviewed journal.

Ritchie, S. J., & Tucker-Drob, E. M. (2018). How much does education improve intelligence? A meta-analysis. *Psychological Science, 29*(8), 1358–1369. http://doi.org/10.1177/0956797618774253

# Quasi-Experimental Designs: Educational effects on intelligence



**Fig. 2.** Funnel plots showing standard error as a function of effect size, separately for each of the three study designs. The dotted lines form a triangular region (with a central vertical line showing the mean effect size) where 95% of estimates should lie in the case of zero within-group heterogeneity in population effect sizes. Note that 42 of the total 86 standard errors reported as approximate or as averages in the original studies were not included for the school-age-cutoff design.

"[…] we found highly consistent evidence that longer educational duration is associated with increased intelligence test scores. […] Thus, the results support the hypothesis that education has a causal effect on intelligence test scores. The effect of 1 additional year of education—contingent on study design, inclusion of moderators, and publication-bias correction—was estimated at approximately 1 to 5 standardized IQ points."

Ritchie, S. J., & Tucker-Drob, E. M. (2018). How much does education improve intelligence? A meta-analysis. *Psychological Science, 29*(8), 1358–1369. http://doi.org/10.1177/0956797618774253

# Quasi-experimental designs

Before-and-after measures



Figure 1. Connecticut Traffic Fatalities, 1955-1956

- was 1956 a dry year? (history)

- overall trends in road safety? (maturation)

- did publicising of death rates have an effect? (testing)

- were fatalities counted differently? (instrumentation)

- was this a big decrease? (instability)

- was 1955 an extreme year? (regression)

Campbell, D. T., Ross, H. L. (1968). The Connecticut crackdown on speeding: Time-series data in quasi-experimental analysis. *Law and Society Review, 3*(1), 33. http://doi.org/10.2307/3052794

# Quasi-experimental designs

Multiple time series



Figure 3. Connecticut and Control States Traffic Fatalities, 1951-1959 (per 100,000 population)

Figure 4. Traffic Fatalities for Connecticut, New York, New Jersey, Rhode Island, and Massachusetts (per 100,000 persons)

Campbell, D. T., Ross, H. L. (1968). The Connecticut crackdown on speeding: Time-series data in quasi-experimental analysis. *Law and Society Review, 3*(1), 33. http://doi.org/10.2307/3052794

# Quasi-experimental designs

Interrupted time series



Figure 2. Connecticut Traffic Fatalities, 1951-1959

- was publicising of death rates similar across years? (testing)

- were fatalities counted differently before and after the intervention? (instrumentation)

Campbell, D. T., Ross, H. L. (1968). The Connecticut crackdown on speeding: Time-series data in quasi-experimental analysis. *Law and Society Review, 3*(1), 33. http://doi.org/10.2307/3052794

# Experimental and Quasi-Experimental Designs for Research

Donald T. Campbell
Julian C. Stanley

1963

## TABLE 1
### SOURCES OF INVALIDITY FOR DESIGNS 1 THROUGH 6

- X = treatment / event
- O = observation of outcome / effect

| | Internal | | | | | | | | External | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Sources of Invalidity** | History | Maturation | Testing | Instrumentation | Regression | Selection | Mortality | Interaction of Selection and Maturation, etc. | Interaction of Testing and X | Interaction of Selection and X | Reactive Arrangements | Multiple-X Interference |
| **Pre-Experimental Designs:** | | | | | | | | | | | | |
| 1. One-Shot Case Study　X　O | − | − | | | | − | − | | | − | | |
| 2. One-Group Pretest-Posttest Design　O　X　O | − | − | − | − | ? | + | + | − | − | − | ? | |
| 3. Static-Group Comparison　X　O／O | + | ? | + | + | + | − | − | − | | − | | |
| **True Experimental Designs:** | | | | | | | | | | | | |
| 4. Pretest-Posttest Control Group Design　R　O　X　O／R　O　　O | + | + | + | + | + | + | + | + | − | | ? | ? |
| 5. Solomon Four-Group Design　R　O　X　O／R　O　　O／R　　X　O／R　　　O | + | + | + | + | + | + | + | + | + | | ? | ? |
| 6. Posttest-Only Control Group Design　R　X　O／R　　O | + | + | + | + | + | + | + | + | + | | ? | ? |

Note: In the tables, a minus indicates a definite weakness, a plus indicates that the factor is controlled, a question mark indicates a possible source of concern, and a blank indicates that the factor is not relevant.

It is with extreme reluctance that these summary tables are presented because they are apt to be "too helpful," and to be depended upon in place of the more complex and qualified presentation in the text. No + or − indicator should be respected unless the reader comprehends why it is placed there. In particular, it is against the spirit of this presentation to create uncomprehended fears of, or confidence in, specific designs.

Campbell & Stanley (1963)

# TABLE 2

## SOURCES OF INVALIDITY FOR QUASI-EXPERIMENTAL DESIGNS 7 THROUGH 12

• X = treatment / event
• O = observation of outcome / effect

| Sources of Invalidity | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | **Internal** | | | | | | | | **External** | | | |
| | History | Maturation | Testing | Instrumentation | Regression | Selection | Mortality | Interaction of Selection and Maturation, etc. | Interaction of Testing and X | Interaction of Selection and X | Reactive Arrangements | Multiple-X Interference |
| *Quasi-Experimental Designs:* | | | | | | | | | | | | |
| 7. Time Series  O O O O X O O O O | − | + | + | ? | + | + | + | + | − | ? | ? | |
| 8. Equivalent Time Samples Design  $X_1O\ X_0O\ X_1O\ X_0O$, etc. | + | + | + | + | + | + | + | + | − | ? | − | − |
| 9. Equivalent Materials Samples Design  $M_aX_1O\ M_bX_0O\ M_cX_1O\ M_dX_0O$, etc. | + | + | + | + | + | + | + | + | − | ? | ? | − |
| 10. Nonequivalent Control Group Design  O X O / O O | + | + | + | + | ? | + | + | − | − | ? | ? | |
| 11. Counterbalanced Designs  $X_1O\ X_2O\ X_3O\ X_4O$ / $X_2O\ X_4O\ X_1O\ X_3O$ / $X_3O\ X_1O\ X_4O\ X_2O$ / $X_4O\ X_3O\ X_2O\ X_1O$ | + | + | + | + | + | + | + | ? | ? | ? | ? | − |
| 12. Separate-Sample Pretest-Posttest Design  R O (X) / R X O | − | − | + | ? | + | + | − | − | + | + | + | |
| 12a. R O (X) / R X O / R O (X) / R X O | + | − | + | ? | + | + | − | + | + | + | + | |
| 12b. R $O_1$ (X) / R $O_2$ (X) / R X $O_3$ | − | + | + | ? | + | + | − | ? | + | + | + | |
| 12c. R $O_1$ X $O_2$ / R X $O_3$ | − | − | + | ? | + | + | + | − | + | + | + | |

Campbell & Stanley (1963)

TABLE 3
SOURCES OF INVALIDITY FOR QUASI-EXPERIMENTAL DESIGNS 13 THROUGH 16

| | Sources of Invalidity | | | | | | | | | | | |
| | Internal | | | | | | | | External | | | |
| | History | Maturation | Testing | Instrumentation | Regression | Selection | Mortality | Interaction of Selection and Maturation, etc. | Interaction of Testing and X | Interaction of Selection and X | Reactive Arrangements | Multiple-X Interference |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Quasi-Experimental Designs Continued:** | | | | | | | | | | | | |
| 13. Separate-Sample Pretest-Posttest Control Group Design $R\ O\ (X)$ $R\ \ \ \ X\ \ O$ $R\ O$ $R\ \ \ \ \ \ O$ | + | + | + | + | + | + | + | − | + | + | + | |
| 13a. $R\ O\ (X)$ $R\ \ \ X\ O$ $R\ O\ (X)$ $R\ \ \ X\ O$ $R\ O\ (X)$ $R\ \ \ X\ O$ $R\ O$ $R\ \ \ \ \ O$ $R\ O$ $R\ \ \ \ \ O$ $R\ O$ $R\ \ \ \ \ O$ | + | + | + | + | + | + | + | + | + | + | + | |
| 14. Multiple Time-Series $O\ O\ OXO\ O\ O$ $\overline{O\ O\ O\ O\ O\ O}$ | + | + | + | + | + | + | + | + | − | − | ? | |
| **15. Institutional Cycle Design** Class $A$  $X$ $O_1$ Class $B_1$  $RO_2$  $X$  $O_3$ Class $B_2$  $R$  $X$  $O_4$ Class $C$  $O_5$  $X$ $^a$Gen. Pop. Con. Cl. $B$ $O_6$ $^a$Gen. Pop. Con. Cl. $C$ $O_7$ | | | | | | | | | | | | |
| $O_2 < O_1$ $O_5 < O_4$ | + | − | + | + | ? | − | ? | | + | ? | + | |
| $O_2 < O_3$ | − | − | − | ? | ? | + | + | | − | ? | + | |
| $O_2 < O_4$ | − | − | + | ? | ? | + | ? | | + | ? | ? | |
| $O_6 = O_7$ $O_{2y} = O_{2o}$ | | + | | | | | | − | | | | |
| 16. Regression Discontinuity | + | + | + | ? | + | + | ? | + | + | − | + | + |

$^a$ General Population Controls for Class B, etc.

Campbell & Stanley (1963)

# Experimental and Quasi-experimental Designs

Experimental and
Quasi-Experimental
Designs for Research

Donald T. Campbell
Julian C. Stanley

"In conclusion, in this chapter we have discussed alternatives in the arrangement or design of experiments, with particular regard to the problems of control of extraneous variables and threats to validity. (…) Throughout, attention has been called to the possibility of creatively utilizing the idiosyncratic features of any specific research situation in designing unique tests of causal hypotheses." (p. 71)

# A colorful bouquet of creating counterfactuals

"The stronger the demonstrated consistency of an association under conditions that rule out alternative hypotheses and the stronger the evidence regarding a mechanism that can explain the observed association, the more likely we are to accept the causal hypothesis. Usually the evidence required to confirm a causal hypothesis is cumulated across multiple studies, many of which are, of necessity, observational. **Although a wide variety of research designs and analytic techniques are available to assist in gathering evidence to support a causal inference, they are helpful only to the extent that their use is guided and constrained by appropriate subject-matter considerations. No method or set of methods defines causality**."

Marini, M. M., & Singer, B. (1988). Causality in the social sciences. *Sociological methodology, 18*, 347-409.

# Summary

- **Importance of counterfactuals:** "The critical step in any causal analysis is estimating the counterfactual—a prediction of what would have happened in the absence of the treatment."

- **Limitations for RCTs:** RCTs are great but do not guarantee effectiveness, generalizability, or ethical treatment of participants.

- **Alternatives to RCTs:** Automation is on the rise, but ethical and safety issues will be crucial! Quasi-experimental designs come in many different forms with different threats to internal and external validity.