Pre-print version of PJ Ferraro et al. 2023. Create a culture of experiments in environmental programs. *Science*. \* This version of the article has been accepted for publication, after peer review but is not the Version of Record and does not reflect post-acceptance improvements, or any corrections The Version of Record is available online at <a href="https://www.science.org/doi/10.1126/science.adf7774">https://www.science.org/doi/10.1126/science.adf7774</a>

# OVERLINE

# Create a culture of experiments in environmental programs

# Organizations need a better "learning by doing" approach.

*By* Paul J. Ferraro<sup>1\*</sup>, Todd L. Cherry<sup>2</sup>, Jason F. Shogren<sup>2</sup>, Christian A. Vossler<sup>3</sup>, Timothy N. Cason<sup>4</sup>, Hilary Byerly Flint<sup>2</sup>, Jacob P. Hochard<sup>2</sup>, Olof Johansson-Stenman<sup>5</sup>, Peter Martinsson<sup>5</sup>, James J. Murphy<sup>6</sup>, Stephen C. Newbold<sup>2</sup>, Linda Thunström<sup>2</sup>, Daan van Soest<sup>7</sup>, Klaas van 't Veld<sup>2</sup>, Astrid Dannenberg<sup>8</sup>, George F. Loewenstein<sup>9</sup>, Leaf van Boven<sup>10</sup>

An understanding of cause and effect is central to the design of effective environmental policies and programs. But environmental scientists and practitioners typically rely on field experience, case studies, and retrospective evaluations of programs that were not designed to generate evidence about cause and effect. Using such methods can lead to ineffective, or even counter-productive, programs.

To help strengthen inferences about cause and effect, environmental organizations could rely more on formal experimentation within their programs, which would leverage the power of science while maintaining a "learning by doing" approach. Although formal experimentation is a cornerstone of science and is increasingly embedded in non-environmental social programs, it is virtually absent in environmental programs. We highlight key obstacles to such experimentation and suggest opportunities to overcome them.

By "formal experimentation," we mean the deliberate creation of spatial or temporal variation in program implementation with the intent of quantifying impacts and elucidating mechanisms. For example, consider an environmental agency that wants to learn how best to encourage polluters to comply with environmental regulations. Instead of implementing a single change in auditing practices across all polluting facilities, the agency could randomly vary implementation of two auditing practices and contrast how facilities respond (Figure 1; see (1) for an analogous real-world example). By creating deliberate variation in how programs are implemented, program administrators can more easily learn about the features that make programs effective. Although

Denmark.

57 58 59

- <sup>8</sup> University of Kassel, Kassel, Germany.
- <sup>9</sup> Carnegie Mellon University, Pittsburgh, PA, USA.
- <sup>10</sup> University of Colorado, Boulder, CO, USA.
- \* Corresponding author. Email: pferraro@jhu.edu

experimentation in natural resource management has a long history, including in the context of adaptive management, we focus on embedding experiments in the implementation of policies or programs that affect human behavior. For example, in a not atypical type of environmental policy experiment that tests whether thinning a reforested plot leads to more harvestable timber, human behavior is controlled by the experimentalist, whereas in a much less common type of experiment that tests alternative design features of a program that encourages more reforestation behavior, human behavior is endogenous and uncertain.

Despite the benefits of adding experimental variation to program implementation, as demonstrated in non-environmental contexts such as health and education, environmental organizations rarely do so. Consider two U.S. federal agencies with substantial environmental program portfolios: the U.S. Environmental Protection Agency (USEPA) and the U.S. Department of Agriculture (USDA). In the last 30 years, each has embedded formal experimentation in their environmental programs fewer than a half dozen times (2). In Europe, we know of only a single example of formal experimentation embedded within government-implemented environmental programs (3). Formal experimentation is similarly almost non-existent among non-governmental and multi-lateral environmental organizations. Although environmental actors engage in thousands of informal "experiments" every year (e.g., pilot programs), these are not designed to test the implicit hypotheses that justify the implementation of current programs or understand how to make these programs more effective.

Formal experimentation in environmental programs is absent because science typically stops when implementation starts. Over the past five decades, governmental and nongovernmental actors have invested substantial resources to understand the status and trends of myriad environmental indicators. These investments have been motivated by scientific uncertainty about how complex environmental systems function and by a recognition that reducing this uncertainty is critical to designing effective programs.

Yet uncertainty also plagues program efficacy. The coupled natural-human systems in which environmental programs are implemented are complex and our understanding of how programs influence the trajectory of these systems is incomplete. When new program designs in non-environmental contexts are assessed through formal experimentation, proponents often learn that the innovations fail to have the intended effects (4). Scientists and practitioners should not expect innovations in environmental programs to be any different.

The absence of experimentation within environmental programs can be explained, in part, by historical reasons. Compared to other social policy fields like health, poverty, and education, the environmental policy field is much younger and would be expected to be a late adopter of innovative ways of generating evidence. Moreover, the human benefits from effective environmental programs are less salient than in other social policy fields. The foregone benefits from ineffective programs are also less salient, putting less pressure on program staff to show effectiveness. Lastly, environmental practice is dominated by lawyers, engineers, and natural and physical scientists who, unlike health, behavioral, and social scientists, do not typically use experimental designs in real-world contexts and may not anticipate complex human responses to what seem like straightforward policy and program decisions. Yet there are no structural barriers to experimentation in the environmental field.

## CONCERNS ABOUT EXPERIMENTATION

Four primary concerns about embedding formal experimentation into environmental programs need to be addressed: delayed action, feedback lags, structural barriers, and ethical questions.

First, identifying ways to create experimental variation and measure outcomes can

<sup>&</sup>lt;sup>1</sup> Johns Hopkins University, Baltimore, MD, USA.

<sup>&</sup>lt;sup>2</sup> University of Wyoming, Laramie, WY, USA.

<sup>&</sup>lt;sup>3</sup> University of Tennessee, Knoxville, TN, USA <sup>4</sup> Purdue University, West Lafayette, IN, USA.

 <sup>&</sup>lt;sup>5</sup> Technical University of Denmark, Kongens Lyngby,

<sup>&</sup>lt;sup>6</sup> University of Alaska-Anchorage, Anchorage, AK, USA.

<sup>&</sup>lt;sup>7</sup> Tilburg University, Tilburg, The Netherlands.

delay scaled-up implementation, thereby letting environmental damages accumulate. Yet, 3 for the case of ineffective programs, the accumulated damages could be much larger when program managers rely on retrospective evalu-6 ations that use non-experimental, post-implementation data. The costs of delays from experimentation will depend on how effectively the program meets its objectives and how 10 quickly damages accumulate. In some cases, 11 large-scale action may be required without 12 waiting for experimentation (akin to "emer-13 gency authorizations" in medicine). Yet, we be-14 lieve that in many cases experimentation em-15 bedded in program implementation will 16 improve outcomes in the long-run, even at the 17 cost of some delay in the short-run. Similar ar-18 guments have been made in the recent COVID-19 19 pandemic, in which calls for quick action and 20 for rigorous evidence seemed to be in opposi-21 tion (5)

1

2

4

5

7

8

9

22 Second, the full effects of a program may 23 not materialize for many years (e.g., long-run 24 climate impacts) and the evidence may no 25 longer be useful by the time it is available. Yet, 26 for many environmental problems, the culprit 27 is human behavior, for which the desired 28 changes can be measured on shorter time 29 scales (e.g., changes in energy consumption by 30 households or fertilizer use by farmers). 31 Measures of short-term environmental indica-32 tors along the hypothesized causal path may 33 also help elucidate whether the intervention is 34 working as intended (e.g., measure pollutants 35 that change relatively rapidly rather than 36 health conditions that change more slowly). 37

Third, structural constraints, such as legal 38 and regulatory rules, may present barriers to 39 experimentation. The degree to which such 40 barriers exist, however, is difficult to ascertain given there has been so little historical effort al-41 42 located to experimentation.

43 A fourth concern may seem, on the sur-44 face, to be the most problematic: opponents of 45 experimentation question the ethics of treat-46 ing some people (or non-human organisms or 47 ecological communities) differently than oth-48 ers (6). This concern arises from a presumption 49 that those exposed to a program, or a specific 50 version of it, are sure to benefit from it. That 51 assumption, however, is not necessarily true. 52 The effects of many environmental programs 53 are uncertain.

54 One could argue that environmental or-55 ganizations have an ethical obligation to better 56 understand the effects of untested programs, 57 or changes in programs, before large groups of 58 humans and other species, particularly vulner-59 able subgroups, are exposed to them (i.e., akin to the principle of "equipoise," a state of genuine uncertainty about the comparative merits of different approaches, which is the ethical basis for justifying randomized treatments in medical trials). Even programs that do not directly harm the environment or people may simply be ineffective. Directing resources to ineffective interventions has substantial ethical implications, especially for environmental problems that are time-sensitive, such as the loss of biological diversity and the accumulation of persistent pollutants.

If environmental organizations were guided by an ethical precept that required evidence before changing or scaling up a program, the science and practice of environmental protection would look different and be more successful. Environmental programs would routinely be subjected to experimentation that deliberately manipulates the temporal and spatial variability of implementation. Program managers, perhaps in collaboration with academics, would then evaluate the results to better understand the consequences, intended and unintended, of the variations in implementation. This evidence would provide opportunities to adjust and improve current and future programs (e.g., 7-8). This cycle of program innovation, experimentation, learning, and adaptation is a hallmark of evidence-based programs in other fields.

#### ENCOURAGING EXPERIMENTATION

Although the constraints on engaging in experimentation will vary by organization, the opportunities for experimentation have some commonalities. Based on experiences in other social policy fields, we offer four recommendations for expanding the opportunities for experimentation in environmental programs (for others, see (9-10)).

#### Political and legal simplicity

Running an experiment that contrasts an entire program to a no-program control may require extensive legal and political approvals, as well as expose implementers to reputational risks and coordination costs. Instead, one version of program implementation can be compared to another version by manipulating program attributes for which managers already have the authority to change (often called A/B testing in the private sector). For example, program managers could contrast the effects on pollution compliance from on-site inspections (status quo) vs remote inspections. Leveraging already planned pilot programs can also be a practical way to facilitate learning when the pilot's implementation is varied across space or time in ways unrelated to the program's target outcomes.

## **Financial simplicity**

Given that the additional costs of experimentation largely come from the costs of measuring outcomes, organizations can focus on contexts where the outcomes are collected as part of program operations (e.g., pollution discharges) or are publicly available (e.g., satellite data of land use).

#### Learning-focused

To achieve higher returns on investment, organizations should focus on experimentation that yields results that can be generalized across multiple programs. Generalizability is more plausible when the program features being manipulated are found in many programs (e.g., capacity building, incentives) or motivated by similar theories of change.

#### Partnership-enhanced

A quick, inexpensive way for environmental organizations to acquire the technical capacity to design and analyze experiments, while keeping the operations in-house, is to embed trained experimentalists from outside the organization (e.g., via federal Voluntary Service Agreements in the US context).

Strengthening the culture of experimentation in the environmental community will require changes in norms and incentives. Program managers are often not rewarded for evidence about program effectiveness but rather for achieving other objectives (e.g., moving money to constituents, avoiding litigation by private actors, pleasing funders). Nevertheless, changes in norms and incentives are occurring. One recent example of change is the creation of "behavioral insights teams" in governmental and multi-lateral organizations. These teams help program managers to formally experiment with program changes inspired by insights from the behavioral sciences (11).

For federal agencies in the United States, changes in norms and incentives are also occurring through the Evidence-Based Policymaking Act of 2018. The Act and complementary memoranda from the Executive branch encourage a culture of experimentation both directly and indirectly. They encourage experimentation directly by emphasizing the power and political acceptability of randomized implementation designs (12-15). They encourage experimentation indirectly by requiring agencies to create annual learning agendas and a strategy and budget to meet their agenda objectives. Learning agendas comprise a set of questions that, when answered, are expected to have the biggest impact on an agency's performance. Yet the Act and its associated guidance do not provide explicit rewards to staff for posing

substantive learning questions and using experimentation to generate high-quality answers to these questions. Thus, by itself, the Act may be insufficient to create a meaningful culture of experimentation within environmental agencies.

One way to further foster a culture of experimentation and embed learning in daily operations among U.S. federal agencies would be through a new executive order (EO) similar in spirit to EO 12291 for Cost-Benefit Analyses. This new EO would be triggered if a new environmental program, or change in a current program, were to exceed a size threshold, which could be measured by program funding or the size of the affected population. The EO would require the implementing agency to first 18 ascertain the "equipoise" of the proposed pro-19 gram or change in program: is there strong em-20 pirical evidence that the proposed action is the 21 best option? If not, then the agency would be 22 required to embed experimentation into the 23 program with the intent of quantifying envi-24 ronmental and social impacts and understand-25 ing the mechanisms through which those im-26 pacts arise. The EO would require that agencies 27 insert a step between proposing a program-28 matic change and scaling that programmatic 29 change up to the entire eligible population. The 30 EO would also encourage environmental 31 agency staff to involve statisticians and behav-32 ioral scientists before implementation. Cur-33 rently, if these experts are called on at all, it is 34 after implementation to assess what may have 35 transpired - a challenging task when imple-36 mentation was not designed to generate evi-37 dence about impacts and mechanisms. In addi-38 tion to characterizing what type of 39 experimentation is acceptable, the EO would 40 also have a stopping rule, similar in spirit to 41 stopping rules used to decide when to end 42 medical treatment trials. Likewise, the EO 43 would also define when it may be acceptable 44 to forego experimentation.

45 Scientists and practitioners can legitimately 46 argue about the benefits and opportunity costs of allocating scarce time and financial re-47 48 sources to formal experimentation in the envi-49 ronmental sector. Should half of environmen-50 tal programs include experimentation? Is ten 51 percent the right amount? While the optimal 52 share is debatable, we believe that the current 53 allocation of roughly zero percent is sub-opti-54 mal. How much experimentation is embedded 55 in programs should depend on contextual at-56 tributes that make experimentation most valu-57 able (Figure 2).

58 We recognize that experimentation is not 59 the only way that a scientific lens can be applied to improve our understanding of program implementation. Experimentation is best

viewed as part of a mixed-methods approach to generating evidence rather than as a substitute for more traditional ways of gathering evidence. Experimentation should, however, be a regular feature of programs, not a rarity

#### REFERENCES AND NOTES

- 1. D.I. Levine, M.W. Toffel, M.S. Johnson, Science 336, 907-911 (2012).
- 2. Based on a review by the authors, results of which were then affirmed during personal communication with agency staff in 2022 (Daniel Hellerstein, USDA, and Katherine Dawes, USEPA).
- 3. K. Telle, Journal of Public Economics 99, 24-34 (2013).
- 4. Arnold Ventures. How to solve U.S. social problems when most rigorous program evaluations find disappointing effects (part two-a proposed solution). Straight Talk on Evidence (April 13, 2018). Retrieved from https://www.straighttalkonevi-
- dence.org/2018/04/13/how-to-solve-u-s-socialproblems-when-most-rigorous-program-evaluations-find-disappointing-effects-part-two-a-proposed-solution/
- 5. A.J. London, J. Kimmelman, Science 368, 476-477(2020)
- 6. E. Pynegar, J. Gibbons, J., N. Asquith, J. Jones, Oryx 55, 235-244 (2019).
- 7.. MOVING TO OPPORTUNITY (MTO) FOR FAIR HOUSING DEMONSTRATION PROGRAM. https://www2.nber.org/mtopublic/
- 8. R.H. Brook, et al, The Health Insurance Experiment: A classic RAND study speaks to the current health care reform debate. Santa Monica, CA: RAND Corporation (2006).
- 9. J. Fox. How can a rethink of lessons from field experiments inform future research in transparency, participation and accountability. Evidence Matters, 3ie (May 29, 2019). Retrieved from https://www.3ieimpact.org/blogs/how-can-rethink-lessons-field-experiments-inform-futureresearch-transparency-participation
- 10. E. Duflo, American Economic Review 107, 1-26 (2017).
- 11. S. Wendel, Who is doing applied behavioral science? Results from a global survey of behavioral teams, Behavioral Scientist (October 5, 2020), Retrieved from https://behavioralscientist.org/whois-doing-applied-behavioral-science-results-froma-global-survey-of-behavioral-teams/
- 12. https://www.whitehouse.gov/omb/informationfor-agencies/evidence-and-evaluation/,
- 13. https://www.whitehouse.gov/briefingroom/presidential-actions/2021/01/27/memorandum-on-restoring-trust-in-governmentthrough-scientific-integrity-and-evidence-basedpolicymaking/
- 14. https://www.whitehouse.gov/wp-content/uploads/2021/06/M-21-27.pdf
- 15. Appendix A of OMB M-19-23 describes four broad types of evidence that agencies should use as they implement the Evidence Act: foundational fact finding, policy analysis, program evaluation, and performance measurement. This guidance goes a step further to specify the broad range of methodological approaches that agencies should consider. These approaches include, but are not limited to: "pilot projects, randomized controlled trials, quantitative survey research and statistical analysis, qualitative research, ethnography, research based on data linkages in which records from two or more datasets that refer to the same entity are joined, well-established processes for community engagement and inclusion in research, and other approaches that may be informed by the social and behavioral sciences and data science.'

#### 10.1126/science.adf7774

Figure 1. A Culture of Experimentation within Environmental Programs.

To reduce pollution, regulators can increase onsite inspections, or they can increase opportunities for facilities to do self-audits, with some penalty leniency when violations are self-reported. Selfaudits may be less effective at reducing pollution (measured remotely) than on-site inspections because self-audits allow facilities to hide their noncompliance. Yet, self-audits may be more effective because they make facilities more aware about the law and its relationship to their operations and because they transform errors of omission into errors of commission. By randomly varying how the regulator interacts with polluting facilities, the regulator can not only learn about the relative effectiveness of each form of interaction, but it can also elucidate what drives facilities to comply or not with environmental regulations (e.g., are they rational or imperfectly informed?).

Figure 2. Four Conditions When Experimentation Pays Off. (i) When theory and experience alone cannot unambiguously predict the impacts of expected changes in program implementation (Pre-Change Ambiguity); (ii) When estimating counterfactual outcomes in the absence of a change in program implementation is challenging using traditional approaches (Post-change Ambiguity), (iii) When a change in program implementation is unlikely to pass a benefit-cost test, or cost-effectiveness assessment, without medium or large impacts (High Implementation Cost); and (iv) When the lessons learned from experimentation are generalizable beyond the context in which the program change was implemented (Generalizability of Results).

1

2

Pre-print version of PJ Ferraro et al. 2023. Create a culture of experiments in environmental programs. *Science*. \* This version of the article has been accepted for publication, after peer review but is not the Version of Record and does not reflect post-acceptance improvements, or any corrections The Version of Record is available online at <a href="https://www.science.org/doi/10.1126/science.adf7774">https://www.science.org/doi/10.1126/science.\*</a>





alizability of

High Imp

4