
Connecting Instructors and Learning Scientists via Collaborative Dynamic Experimentation

Joseph Jay Williams
Harvard University
Cambridge, MA, USA
joseph_jay_williams@harvard.edu

Walter S. Lasecki
University of Michigan
Ann Arbor, MI, USA
wlasecki@umich.edu

Anna N. Rafferty
Carleton College
Northfield, MN, USA
arafferty@carleton.edu

Juho Kim
KAIST
Daejeon, South Korea
juhokim@cs.kaist.ac.kr

Andrew Ang
Dustin Tingley
Harvard University
Cambridge, MA, USA
andrew_ang@harvard.edu
dtingley@gov.harvard.edu

Abstract

The shift to digital educational resources provides new opportunities to advance psychology and education research, in tandem with improving instruction using theory and data. To realize this potential, this paper explores how randomized experiments can support mutually beneficial instructor-researcher collaborations. We developed the *Collaborative Dynamic Experimentation* (CDE) framework to address two key tensions. To enable researchers to embed experiments in online lessons while maintaining instructors' editorial control, *Collaborative* experiment authoring is needed. To enable instructors to use data for rapid improvement while maintaining statistically valid data for researchers, we apply an interpretable machine learning algorithm for *Dynamic* experimentation. We worked with an on-campus instructor to implement a proof-of-concept CDE system to experiment within their online calculus quizzes. The qualitative results from this deployment provided insight into how the CDE framework can facilitate alignment of research and practice.

ACM Classification Keywords

H.5.m. [Information Interfaces and Presentation (e.g. HCI)]: Miscellaneous

Author Keywords

Collaborative dynamic experimentation; online education;

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the Owner/Author. Copyright is held by the owner/author(s).
CHI'17 Extended Abstract, May 06-11, 2017, Denver, CO, USA.
ACM 978-1-4503-4656-6/17/05.
<http://dx.doi.org/10.1145/3027063.3053247>

randomized experiments; researcher-instructor collaboration.

Introduction

Instructors are increasingly teaching with online resources, from tutorial webpages in a Learning Management System to quizzes in a Massive Open Online Course (MOOC). As teaching shifts from physical spaces to digital environments, there is a drastic reduction in the barriers to conducting randomized experiments that answer scientific and practical questions. An instructor can't try explain a concept one way to half a class and a different way to the other half in a single lecture, but software could be added to webpages to systematically compare explanations to see which ones students find most helpful, or to an email application to collect data about how different expressions of encouragement motivate students [12]. Increasing the ubiquity of experimentation through online resources has the potential to produce practical improvements, while advancing learning research "in the wild."

Unfortunately, there is a lack of tools for end-user experimentation that meet the needs of instructors and learning researchers. Widely used Learning Management Systems (Blackboard, Canvas, Moodle) and MOOC platforms (Coursera, NovoEd) do not provide instructor-facing tools for experimentation, and the security and data privacy considerations of educational platforms largely prevent the addition of custom code, creating a significant obstacle to embedding tools for generic website experimentation like Optimizely and Google Content Experiments. Plugins are only accepted through specific standards like Learning Tools Interoperability [10]. Fortunately, some platforms do provide tools for experimentation — for example, as of August 2014 the MOOC platform edX provides a tool that lets course staff experiment with content. However, an instructor can't

recruit a researcher to author an experiment without giving the researcher full staff access to all student data and the ability to edit (or delete) the course. The K-12 math platform ASSISTments provides a tool for learning researchers (and instructors) to copy and experiment with versions of interactive math problems. But collaborative editing of shared experiments is not possible [5].

The lack of support for collaborative design of experiments is a tremendous missed opportunity. Many learning researchers have extensive time and expertise available for designing experiments and analyzing data, but need to work with instructors to deploy studies in real classes. Meanwhile, instructors have limited time to design additional versions of resources, and even less to surmount the technical and methodological barriers to conducting experiments.

We present an approach to collaborative experimentation inspired by design-based research in education [1], which encourages collaborations between a researcher and an individual teacher to ensure research is ecologically valid and applicable within the actual contexts where learning occurs. In the spirit of participatory design [15], the Open Learning Initiative formed teams of instructors and cross-disciplinary researchers to collaboratively design online courses, resulting in substantive learning gains for students [9]. A ten-year NSF grant to the Pittsburgh Science of Learning Center enabled researchers building intelligent tutoring systems to more easily form relationships with classroom teachers to collect data [6]. Attempts to scale beyond individual groups include the Carnegie Foundation's use of improvement science [2] to articulate a framework for Networked Improvement Communities. These present organizational structures and methods for making members aware of each other's interests and goals. For the thousands of instructors and

researchers who do not belong to such networks, tools for experimentation could be especially important for realizing the benefits of collaboration.

In this paper, we introduce a Collaborative Dynamic Experimentation framework for facilitating smaller scale collaborations between individual instructors and researchers. This framework focuses on providing tools for instructors and researchers to easily collaborate on designing and deploying experiments, while also dynamically assigning experimental conditions to students. This dynamic assignment draws on work in machine learning to favor more effective conditions, providing immediate improvements to courses from the ongoing research. We present CDEquiz, a proof-of-concept system that instantiates this framework, and describe a case study in which the system was used for co-designing and deploying three experiments in a calculus course. Qualitative comments from the instructor illustrate the potential benefits of this approach.

Collaborative Dynamic Experimentation

We present a framework for *Collaborative Dynamic Experimentation*. We define a Collaborative Dynamic Experiment as having (1) Collaborative Design and Deployment of Experiments, and (2) Dynamically Randomized Assignment.

Collaborative Design and Deployment of Experiments means:

1. **Iterative Authoring** of the experimental conditions (versions of a resource) through back-and-forth proposal and reviewing/editing by researchers and instructors.
2. **Direct Deployment** and previewing of experimental conditions in the resources students will receive.

3. **Real-Time Outcome Variables** for measuring experimental impact are decided on before the experiment launches, and data about these is made available as soon as it is collected.

Dynamically Randomized Assignment aims to make experiments have more practical impact and be more ethical, by more rapidly using data from ongoing experiments to deliver the most effective conditions to students.

Deciding when there is enough data to be confident of choosing the best condition can be complex. It is an example of the exploration versus exploitation tradeoff that is extensively studied in reinforcement learning (see [16] for an overview). Exploiting corresponds to trying to help students by assigning everyone to the condition that the current data suggests is the best. Exploring is assigning students to conditions that current data suggests may be worse, but which *may* turn out to be better once more observations are made. In addition to the instructor's goal of finding the best condition, deciding to stop the experiment also impacts the researcher's goal of drawing statistically justified conclusions from the data.

We formalize the decision about which conditions are assigned to students as a multi-armed bandit problem, a common approach to automated online experimentation in websites [7]. In a multi-armed bandit problem, the agent's goal is to maximize its total rewards. At each time point, it chooses one action from a set of possible actions, and receives a noisy reward based on the action chosen. In our application, rewards are related to student outcomes. The *Policy* specifies how the agent decides which action to take, often based on past data.

Our goal was for the experimentation *Policy* to maintain

consistency and *interpretability* with the methods learning researchers use to conduct and analyze randomized experiments. We therefore make the design choice to use *probability matching* algorithms: the probability of assigning a condition to a student is proportional to the probability that it is the best condition given the previously observed data. Probability matching algorithms are not guaranteed to do *uniform* randomized assignment (e.g., 50/50 with two conditions) as is traditional in scientific experiments. But they do ensure *weighted* randomization (e.g., 70/30) with weights changing dynamically as evidence accrues.

One example is Thompson Sampling [3], an algorithm that uses Bayesian statistics to model the expected outcome of each condition, given data from students who have already received a condition. The probability of assigning condition X to a student is equal to the current probability that condition X has the highest value, according to the statistical model. This provides instructors with a convenient interpretation of how and why conditions are assigned to students. This also corresponds to a researcher's goal of maintaining randomization of conditions throughout the experiment. Thompson Sampling has been used in previous applications to optimize educational resources [18, 8], and we use it to guide experimentation in the proof-of-concept system we present later, CDEquiz.

Proof-of-Concept System for Collaborative Dynamic Experimentation: CDEquiz

This section presents CDEquiz, a proof-of-concept system for Collaborative Dynamic Experimentation on components of quizzes. CDEquiz enables experimentation with different components of quizzes that have been targeted by past work: explanations for correct answers [13], feedback on wrong answers [4], and learning tips [11].

CDEquiz can be added as a plug-in LTI tool to a page in a course website. Students navigating to the page simply access a quiz. When instructors and researchers go the page, they access a composite interface: (1) A quiz preview option that simulates a student viewing the quiz. (2) The Iterative Authoring Interface for writing and reviewing quiz components, that are directly deployed into the relevant quiz, and the Data and Policy Dashboard. The backend code implements the algorithm for automatically analyzing data to do Dynamically Randomized Assignment of conditions.

Iterative Authoring Interface

The Iterative Authoring Interface for CDEquiz enables researchers or instructors to add, review, and edit alternative versions of the quiz components (learning tips, explanations, or feedback messages). Direct Deployment into a student's quiz happens automatically: when a student views a quiz (or an instructor previews the quiz) one of the versions of the component is automatically selected by the experimentation Policy and embedded in the quiz. The preview feature of the interface allows researchers and instructors to simulate multiple different students taking the quiz: they can go directly from Iterative Authoring to previewing Direct Deployment of how the different conditions appear to students. In addition to allowing instructors to preview their own conditions, this allows them to reflect on a collaborator's proposals in the context of the student experience.

Dynamically Randomized Assignment

CDEquiz assigns students to experimental conditions using weighted randomization, where the probabilities for condition assignment are determined by the Thompson Sampling algorithm. Students are asked to rate how helpful each explanation/feedback message is for learning, and this rating is provided to the algorithm as a measure of the success of the condition. This student rating ranges from 0 (com-

pletely unhelpful) to 10 (perfectly helpful). Thompson sampling maintains a Beta distribution to model each condition. This represents its beliefs about how the condition is associated with the student ratings. The model uses a Beta(19, 1) prior, representing optimistic but uncertain initial beliefs about each condition, and is updated using a Binomial likelihood function. This function assumes the student rating of r follows a Binomial distribution with 10 samples, r of which are positive. This Beta-Binomial model can be updated efficiently, allowing the model to be immediately updated each time a student provides a rating. Thus, the policy changes in real time based on what has been effective for previous students.

Case Study: Deployment of CDEquiz

The final CDEquiz system emerged from an iterative design process in which earlier versions of the system were deployed in a calculus instructor's class in order to conduct three Collaborative Dynamic Experiments. The Researcher interviewed the calculus instructor (referred to as Instructor) about details of the course content and student population, and the Instructor's pedagogical goals and challenges. The Researcher and Instructor decided to focus on experimentation in low-stakes quizzes, which were intended to serve as checks on understanding and provide students with formative feedback.

The first experiment compared alternative versions of explanations in quizzes. The system was built for a multiple choice calculus quiz on integration by substitution. After a student chose the correct answer, they received an explanation for why that answer was correct, to solidify their understanding. In each experimental condition, a different explanation of the correct answer was shown to students, and the students rated the helpfulness of the explanation that they were shown.

The Researcher drafted two alternative explanations to be compared against the original by the Instructor, drawing on work in psychology and education on what makes effective explanations for students [13]. One included step-by-step justifications, which past results on worked-examples have shown to help students more than solving additional problems [17]. A second explanation included a description of the thinking process students could follow in deciding how to solve the problem [14]. The Instructor and Researcher did two back-and-forth rounds of feedback and editing before finalizing the explanations for deployment.

The Researcher then sought the Instructor's suggestions for what questions could be addressed next through experimentation, resulting in two new experiments. The first experiment compared students' ratings of the helpfulness of different feedback messages about wrong answers. These feedback messages were written by the instructor. The second experiment (conducted concurrently but randomized independently) examined how different learning tips impacted success on solving a problem. These tips were collaboratively constructed by the Researcher (drawing on research on metacognition, [11]) and the Instructor.

Insights from Instructor

Experimentation The Instructor found that trying to design alternative conditions for an experiment was useful because it generated novel insights about pedagogy and student learning. Although the Instructor indicated that he had not previously considered creating multiple versions of an explanation, he thought that there was a benefit to doing so beyond experimental design. He noted that it "encouraged [him] to think about these other sorts of explanations" and thought that this reflection was a positive outcome of collaborating on an experiment, even if that experiment did not result in finding significant differences between con-

ditions. The Instructor appreciated that collaborating with a researcher could provide a complementary theoretical perspective: “it’s clear to me that you’ve read a lot of research on student learning that I have no idea about. I think you must know plenty of general things about how students learn, whereas I know specific things about how they get calculus questions wrong or right.”

Collaborative Experiment Creation Interface The Instructor found it much more efficient to collaborate via CDEquiz’s Experiment Creation Interface than via email and Google Docs. Based on his experiences collaborating in these other interfaces, the Instructor said that without the CDEquiz interface, he would limit his future experimentation work to “a smaller scale, smaller number of quizzes.” The Experiment Creation Interface succeeded in minimizing the technical sophistication needed. When the Instructor was asked about trying to do experiments with an alternative approach: “No, I would have no idea how to go about doing that. I’m not aware of any tools that do this sort of thing and ...even if I found such a tool somewhere outside of Canvas I don’t think that I have the technical expertise to incorporate it into Canvas.”

Dynamically Randomized Assignment The Instructor intuitively grasped the the Policy’s probability of condition assignment as the probability of the condition being the best. His mathematics background may have helped. When asked about the importance he placed on the experiment’s randomization being modified using previous data, the Instructor noted that “In terms of the experiment I don’t know how much value [it] adds but in terms of what I think is best for the students it adds value that students are more likely to get the explanation that’s going to help them.” The Instructor appreciated that the algorithm’s automatic modification of the experiment did not require his monitoring:

“if I then have to keep looking at the actual student ratings that’s a lot of extra work...I think something that does it automatically is a big advantage.” The algorithm also led the Instructor to consider the value of a continuously running experiment across multiple offerings of his course: “as we do multiple semesters of the course and if the data is able to be transferred easily ... to increasingly refine which versions are the best and feed those to the students.”

Conclusion

We presented the Collaborative Dynamic Experimentation framework’s design recommendations for experimentation tools, and validated these in a proof-of-concept system that enabled an instructor and researcher to deploy experiments. The case study revealed that the instructor didn’t think they’d have the technical skills to do experiments without the system, and the collaborative authoring reduced friction that would have hindered their interest in future studies. The system’s dynamic tradeoff between comparing conditions and implementing the best conditions gave the instructor confidence that the experiment was beneficial not only to the researcher, but also added value for current and future students. Future work will explore a larger deployment of CDEquiz, providing the opportunity to quantitatively evaluate the effects of dynamic experimentation.

References

- [1] Sasha Barab and Kurt Squire. 2004. Design-based research: Putting a stake in the ground. *The journal of the learning sciences* 13, 1 (2004), 1–14.
- [2] Anthony S Bryk, Louis M Gomez, and Alicia Grunow. 2011. Getting ideas into action: Building networked improvement communities in education. In *Frontiers in sociology of education*. Springer, 127–162.
- [3] Olivier Chapelle and Lihong Li. 2011. An empirical

- evaluation of thompson sampling. In *Advances in neural information processing systems*. 2249–2257.
- [4] Carol S Dweck. 2002. Messages that motivate: How praise molds students' beliefs, motivation, and performance (in surprising ways). (2002).
- [5] Neil T Heffernan and Cristina Lindquist Heffernan. 2014. The ASSISTments Ecosystem: Building a platform that brings scientists and teachers together for minimally invasive research on human learning and teaching. *International Journal of Artificial Intelligence in Education* 24, 4 (2014), 470–497.
- [6] Kenneth R Koedinger, Ryan Sjd Baker, Kyle Cunningham, Alida Skogsholm, Brett Leber, and John Stamper. 2010. A data repository for the EDM community: The PSLC DataShop. *Handbook of educational data mining* 43 (2010).
- [7] John Langford and Tong Zhang. 2008. The epoch-greedy algorithm for multi-armed bandits with side information. In *Advances in neural information processing systems*. 817–824.
- [8] J Derek Lomas, Jodi Forlizzi, Nikhil Poonwala, Nirmal Patel, Sharan Shodhan, Kishan Patel, Ken Koedinger, and Emma Brunskill. 2016. Interface Design Optimization as a Multi-Armed Bandit Problem. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*. ACM, 4142–4153.
- [9] Marsha Lovett, Oded Meyer, and Candace Thille. 2008. JIME-The open learning initiative: Measuring the effectiveness of the OLI statistics course in accelerating student learning. *Journal of Interactive Media in Education* 2008, 1 (2008).
- [10] G Mcfall and L Neumann. 2010. IMS learning tools interoperability basic LTI implementation guide v1. 0 Final. Internet: <http://www.imsglobal.org/lti/2010> [May 2013] (2010).
- [11] Tim McNamara and Carsten Roever. 2006. *Language testing: The social dimension*. Vol. 1. John Wiley & Sons.
- [12] Jason A Okonofua, David Paunesku, and Gregory M Walton. 2016. Brief intervention to encourage empathic discipline cuts suspension rates in half among adolescents. *Proceedings of the National Academy of Sciences* (2016), 201523698.
- [13] Alexander Renkl. 1997. Learning from worked-out examples: A study on individual differences. *Cognitive science* 21, 1 (1997), 1–29.
- [14] Alan H Schoenfeld. 1992. Learning to think mathematically: Problem solving, metacognition, and sense making in mathematics. *Handbook of research on mathematics teaching and learning* (1992), 334–370.
- [15] Douglas Schuler and Aki Namioka. 1993. *Participatory design: Principles and practices*. CRC Press.
- [16] Richard S Sutton and Andrew G Barto. 1998. *Reinforcement learning: An introduction*. Vol. 1. MIT press Cambridge.
- [17] John Sweller. 2006. The worked example effect and human cognition. *Learning and Instruction* 16, 2 (2006), 165–169.
- [18] Joseph Jay Williams, Juho Kim, Anna Rafferty, Samuel Maldonado, Krzysztof Z Gajos, Walter S Lasecki, and Neil Heffernan. 2016. AXIS: Generating Explanations at Scale with Learnersourcing and Machine Learning. In *Proceedings of the Third (2016) ACM Conference on Learning at Scale*. ACM, 379–388.