

The effect of supplementary video lectures on learning in intermediate microeconomics

Melissa Famulari and Zachary A. Goodman*

University of California, San Diego

This version: October 2020

Abstract

The abstract goes here eventually.

*mfamulari@ucsd.edu and zgoodman@ucsd.edu. The authors thank the students who took intermediate microeconomics in the fall of 2018 and 2019 who consented to the use of their data for this study. We also thank UC San Diego's Teaching and Learning Commons for providing campus data on the students in this study as well as anonymizing the data for analysis. Finally, we thank the applied microeconomics group at UC San Diego for their help with the experimental design. This research was approved under UC San Diego's Human Research Protections Program (IRB approval 170886 in fall 2018 and 2019). The paper investigates the use of Intermediate Microeconomics Video Handbook (IMVH) video lectures by UC San Diego students, some of which were developed by one of the authors, in collaboration with UC San Diego and the UC Office of the President. UC San Diego currently owns the rights to distribute the IMVH. The videos lectures were provided to the subjects at no charge and neither author has a direct financial interest in the distribution of the IMVH at UC San Diego. As of fall 2020, one of the authors has a financial interest in the distribution of the IMVH outside of UC San Diego.

1 Introduction

“You expect me to read the textbook? Ha!”

— Anonymous student

University students spend tens of thousands of dollars annually on tuition and hundreds of hours in lecture and completing assignments, in large part, to learn. Instructors can improve how well students learn by employing pedagogical tools that have the greatest returns per unit time and financial cost. Despite the importance of comparing the effectiveness of different teaching technologies, little empirical work exists that estimates . In this paper we examine the impact of low marginal cost, video-based learning materials on exam scores in a large, intermediate microeconomic theory course.

The Intermediate Microeconomic Video Handbook (IMVH) at UC San Diego was designed to *supplement* lecture, not replace it, as an audiovisual version of a conventional course textbook. Part of the impetus for creating the IMVH was a discussion with a student who described her inability to read the course text, not because of poor reading skills, but because she did not find the text engaging enough to command her attention. We hypothesized that current students, who have had unprecedented exposure to electronic media, would find video materials more engaging and ultimately study more effectively or for more time than they would have if provided only conventional studying materials. Though students may use the IMVH more than the textbook, it is an empirical question whether the videos ultimately improve learning outcomes.

We answer this question using a field experiment involving over 400 undergraduates enrolled in the same microeconomics course over two years. Only students who scored below the median on the first midterm were eligible for the experiment, since previous work and institutional knowledge suggests that students in the top half of the distribution would not benefit (and may even be harmed) from being induced to watch the videos. While the optimal experimental design for identifying average treatment effects would involve restricting access to the IMVH to only treated students, ethical considerations required that all students have access to the IMVH. Hence, we opted for an encouragement design in which treated students are induced to watch more videos than their control group peers through a grade-based

incentive, which more than doubled the number of videos watched by treated students. This experimental design permits identification of treatment effects local to those students induced by the encouragement to watch more videos.

We find that being assigned treatment (ITT) increased midterm and final exam scores by 0.18 and 0.17 standard deviations, respectively, and that the marginal hour of video watched increased exam scores (LATE) by 0.08 standard deviations. Although the confidence intervals are, admittedly, wide, the point estimates are statistically and economically significant: a student could increase their course letter grade by one step (e.g. from a B+ to A-) by watching XX hours of videos. Our estimates suggest that XX percent of students in the control group who failed the course would have earned passing grades had they watched as many videos as their treated counterparts.

Although treated students performed better on course assessments, for determining welfare it is important to identify where the time watching videos came from: leisure time, working, student organizations, studying for other classes, studying for current class using other methods, etc. On one hand, if watching videos is more productive than the next best studying method, then the utility of requiring videos is unambiguously positive as students can substitute studying time towards the more productive option. On the other hand, if students must reduce time allocated towards leisure or studying for other classes so they can watch more videos, then the welfare implications are less clear and could be negative depending on the students' preferences.

We attempt to disentangle whether treated students spent more time studying or used their time more effectively by examining proxies for time use including class attendance, visits to a tutoring center (specific to this course), downloading materials from the course website, posting on the class discussion board, and reported time use from an in-class survey. Although our estimates are noisy, we find no statistically significant differences between treatment and control, and we can reject large decreases in take-up of other study methods by treated students. Surprisingly, in nearly all cases, the point estimates suggest that the treatment group used study methods beyond the videos at *greater* rates than did their control peers. Though estimates are noisy, we find no significant differences in reported leisure time across treatment and control. Finally, we investigate spillovers to other courses

taken during the same academic term as the experiment and similarly find that treated students perform *better* than their control peers, which suggests that watching the videos likely did not dramatically reduce time spent studying for other classes.

Finally, we attempt to distinguish between two models of student learning which could explain why the video watching inducement improved student exam performance: an incomplete information model where students do not know how to study effectively versus a two-selves model where students like good grades but dislike studying. One important (and testable) difference between the incomplete information model and the two-selves model is what happens after exogenous incentives to watch videos are removed. While the former predicts that students exposed to treatment will continue watching videos given their new knowledge of a relatively productive studying technology, the latter predicts that the students will return to their lower baseline levels of video watching as their doer selves no longer have a commitment device reducing the temptation of immediately gratifying leisure. We examine video watching behavior during the term following the experiment in the subsequent microeconomics course and find that treated students watch significantly more videos than their control classmates, consistent with the incomplete information model.

Collectively, we interpret our findings as strong evidence that requiring studying tools known by the instructor to be effective is utility enhancing for students who manifest their limited knowledge of how best to study, perhaps through poor performance on an early stage assessment. Finally, we provide suggestive evidence that students found the IMVH to be a relatively effective study method by examining video watching across the treatment and controls in the next class in the sequence. The rest of the paper is organized as follows. Section 3 provides background on existing related literature. Section 4 describes the experimental design. Section 5 presents the results of the experiment, and Section 6 discusses those results. Section 7 concludes.

2 Models of Studying Behavior

In this section we consider three models of student information acquisition: a neoclassical model, an imperfect information model and a behavioral/procrastination model. For all

three models, we consider the effects of an instructor’s inducement to increase the use of an effective study method. We do not address the issue that the IMVH is a relatively unique study tool in that, to our knowledge, it is the first instructional book to be created entirely of videos. However, given the availability of close substitutes to the IMVH (lecture capture for example) we only briefly explore the added issues of inducing students to use a study tool whose usefulness is not known to the instructor.

Neoclassical models of studying behavior assume that rational agents know their returns to studying using the methods available to them and allocate the optimal amount of study time to each method given their utility function, which is increasing in leisure and grades and decreasing in time spent studying. In this model there is no room for an instructor to increase student well-being by intervening in their study decisions. Oettinger (2002) provides some empirical support for the neoclassical model by demonstrating that student effort responds rationally to nonlinear grade incentives.¹

However, in addition to teaching specific skills, many would agree that the “raison d’etre” of higher education is to teach students how to learn. There is evidence from psychology that college students do not know how to learn effectively.² Universities often fund “Teaching and Learning Centers” or “Academic Skills Centers,” part of whose mission is to help undergraduates learn to study more efficiently.³ We posit that for many students, a key assumption of the neoclassical model does not hold: that students possess accurate information about the returns across studying methods. Instead, we offer the alternative hypothesis that students supply a quantity of study time that is optimal given their information constraints. They choose study methods and quantities that are suboptimal relative to those they would have

¹Oettinger (2002), “The effect of nonlinear incentives on performance: Evidence from Econ 101,” *The Review of Economics and Statistics*, 84(3): 509–517. Across 1200 students in a principles of economics class with absolute grading standards, Oettinger finds evidence of bunching just above the letter grade cutoffs and student performance on the final exam is higher if the student is just below a grade threshold.

²See, for example, Pashler, Rohrer, Cepeda, Carpenter (2007). “Enhancing learning and retarding forgetting: Choices and consequences,” *Psychonomic Bulletin and Review* 2007, 14, 2, 187-193; Dunlosky, J., Rawson, K. A., Marsh, E. J., Nathan, M.J., and Willingham, D. T. (2013). “Improving students’ learning with effective learning techniques promising directions from cognitive and educational psychology.” *Psychological Science in the Public Interest*. Afton Kirk-Johnson, Brian M. Galla, Scott H. Fraundorf (2019). “Perceiving effort as poor learning: The misinterpreted-effort hypothesis of how experienced effort and perceived learning relate to study strategy choice” *Cognitive Psychology*, Volume 115

³All nine University of California campuses have one. Some others in the US include Dartmouth’s Academic Skills Center, Michigan’s Center for Research on Teaching and Learning, UNC’s Learning Center, and Yale’s Teaching and Learning Center

picked in a full information setting. Hence, an intervention by an entity that has more information about returns to studying across various methods (i.e. an instructor) can be utility enhancing.

A third model is a behavioral one in which students plan to study more than they end up studying when the time comes. Indeed, survey and experimental data suggest that many students study less than they report they “should” and finish the term with grades lower than what they had anticipated they would earn at the start of the term.⁴ Recent empirical evidence suggests that setting tasked-based goals help improve college student performance.⁵ This phenomenon is consistent with two-self models in which a person’s “planner” self, the one who desires high grades at the expense of leisure, is at odds with her “doer” self, the one who must choose between immediately gratifying leisure and delayed gratification from higher grades.⁶

Suppose an instructor uses a grade incentive to encourage the use of a time consuming educational input, say, watching instructional videos (or attending class, reading the textbook, doing homework, etc.) Assume the partial derivative of final exam score with respect to watching videos is positive.⁷ What are the testable implications of the three models?

In the neoclassical model, students will watch more videos if the increased utility from the higher class grade, due both to the grade incentive and watching more videos, exceeds the utility losses from the reduction in time spent on other utility-enhancing activities. Exam performance in the class will increase for students who watch more videos. In future classes, in the absense of the incentive, student video watching will revert to pre-incentive levels. If the additional time to watch videos for the current class comes from study time for other

⁴Ferrari, J.R. “Psychometric validation of two Procrastination inventories for adults: Arousal and avoidance measures.” *J Psychopathol Behav Assess* 14, 97–110 (1992). Patricia Chen, Omar Chavez, Desmond C. Ong, Brenda Gunderson (2017) “Strategic Resource Use for Learning: A Self-Administered Intervention That Guides Self-Reflection on Effective Resource Use Enhances Academic Performance” *Psychological Science*, Vol. 28, Issue 6, 774-785. Stinebrinckner and Stinebrickner (2008) show large grade declines if one is randomly assigned a roommate who brings a video game to college

⁵Clark, Damon, David Gill, Victoria Prowse, and Mark Rush (2020) “Using Goals to Motivate College Students: Theory and Evidence From Field Experiments” *The Review of Economics and Statistics* 2020 102:4, 648-663

⁶see review paper, Adam M. Lavecchia, Heidi Liu, and Philip Oreopoulos (2016), “Behavioral Economics of Education: Progress and Possibilities” *Handbook of the Economics of Education*, Volume 5, Pages 1-74)

⁷We expect non-native speakers to particularly benefit from the videos because of the technology itself (e.g. ability to replay videos, change speed and having closed captions). We test for this interaction in the results below.

classes, the incentive may increase performance in the instructor's class at the expense of grades in other classes. If the additional time came from leisure time, , say, participation in student organizations, the incentive may reduce the likelihood of being in a leadership position.

In the imperfect information model where the instructor has perfect information, the incentive to watch videos will increase exam performance as long as student study time does not fall. While unlikely, students with poor information may increase video watching but reduce other, less effective educational inputs so much that exam performance falls. A sharp prediction is that video watching will continue at the incentivized level even in the absence of the incentive as students have learned a new effective study tool. We also examine whether the treatment effect is greater for transfer students who we expect have more information problems than students who started at UC San Diego as freshman. Transfer students are taking not only their first class at a 4 year, R1 university but also their first class under the quarter-system as most California community colleges are semester systems.

Finally, in the behavioral model, the instructor's inducement helps students stick to study plans. Again, as long as total study time does not fall, the inducement will increase exam performance. In the absence of the inducement, a sharp prediction is that video watching will fall to pre-inducement levels.

In the empirical section, we test for the effects of being induced to watch the IMVH on both exam scores and several other study methods student use to learn microeconomics (lecture attendance, visits to a class-specific tutoring lab, use of a class discussion board, downloads from the class web page). We also compare grades in other classes taken in the same quarter across the treatment and control group. For a subset of our sample we have survey responses on total study time in the quarter and leisure time. Since the experiment was conducted in the first of a required three-class sequence, we examine video watching in the second class. We test whether the effect of the inducement to watch videos is greater for transfer students (assumed to have more information problems) and non-native English speakers (assumed to benefit the most from closed captioning, a key aspect of the IMVH).

3 Related Literature and Contributions

Students have many time-consuming activities to help them learn including attending class, reading the textbook, doing homework, going to discussion sections, attending office hours, etc. There are several empirical challenges to estimating causal effects of a learning activity. First, unobserved student characteristics, such as ability and motivation, are likely correlated both with the learning activity and class performance. Further, there is evidence of dynamic selection whereby a student's learning strategy adjusts in response to exam information (see, for example, Oettinger (2002), Stinebrickner and Stinebrickner (2008) and Cannon (2011)) and so a fixed effects estimation strategy, common in the literature, is unlikely to identify causal effects. Second, learning activities are substitutable and inducements to use one study strategy may affect use of another, changing the interpretation of causal estimates. Finally, it is possible that the primary benefit to students is simply devoting more time to learning the course material—however they do it.

Attending class: Chen and Lin (2008)⁸ collect attendance and randomly leave out exam material in one but not the other lecture of the same course. By comparing performance of students who attended with those who did not,

Dobkin, Gil and Marion (2010) analyze a policy where lecture attendance was voluntary before the midterm, but after the midterm, students scoring below the median were required to attend class. The policy affected 352 students taking three classes, two intermediate micro and one econometrics class. The policy led to a 36 percentage point increase in post-midterm attendance at the threshold. Using a regression discontinuity design, they find that a 10 percentage point increase in overall attendance results in a 0.17 standard deviation increase in the final exam score. They find no effect of the attendance policy on grades in other classes taken the same quarter, attending TA sections, homework scores, and the use of university tutors. Arulampalam, Naylor and Smith (2012)⁹ study section (as opposed to lecture) attendance across intermediate microeconomics, intermediate macroeconomics and econometrics for 444 students. The authors find that absenteeism depends on day of week

⁸Chen, J. and Lin, T-F. (2008). Class attendance and exam performance: A randomized experiment. *Journal of Economic Education*, 39, 213– 227.

⁹Arulampalam, Wiji, Robin A. Naylor, Jeremy Smith (2012) "Am I missing something? The effects of absence from class on student performance" *Economics of Education Review*, Vol 21, Issue 4, 363-375.

and time of day and, since students are randomly assigned to sections, use these variables as instruments for absenteeism. They also include student fixed-effects. Surprisingly, they find significant attendance effects only for students in the top quantiles: missing 10 percent of sections results in a 1 percentage point performance loss. The authors have no information about other uses of the student's time use, including attending the main lecture.

Effect of Homework: Trost and Salehi-Isfahani (2012) randomly require two-thirds of students taking Principles of Economics classes to complete a one of three homework assignment for a grade. The other third may complete the homework, but it does not contribute to their grade. The outcome measure is exam performance on questions related to the three homework assignments. The authors use the score on the remaining exam questions as a control variable. They find significant effects of homework on the first midterm but not the final exam. Grodner and Rupp (2013)¹⁰ use within-class randomization to estimate the effects of required homework for 423 microeconomics principles students. A coin flip determined whether a student was in the treated group, where course points are based on both homework and exams, or in the control group, where all course points are based on exams. Treatment led to a 58 percentage point increase in completing all homework assignments and a 84 percentage point increase in completing the majority of homework assignments. They find that treated students are less likely to drop the class and score higher on the first two but not the last two exams. The average across the four exams is increased 5-6 percent by treatment and the control group GPA would increase from 2.44 to 2.68 if they had been required to do homework. They find three times larger treatment effects for students who initially fail the first exam (10 to 15 percent vs 4 to 6 percent increase in average test scores). The authors do not examine whether other uses of student time are affected by the homework policy.

Effect of Study time: In a convincing empirical paper on the effects of study time, Stinebrickner and Stinebrickner (2008)¹¹ examine 210 Berea College students who were randomly assigned a roommate. Those whose roommates brought a video game to college, earn lower

¹⁰Grodner, Andrew and Nicholas G. Rupp (2013) The Role of Homework in Student Learning Outcomes: Evidence from a Field Experiment, *The Journal of Economic Education*, 44:2, 93-109.

¹¹Stinebrickner Ralph and Todd R. Stinebrickner, 2008. "The Causal Effect of Studying on Academic Performance," *The B.E. Journal of Economic Analysis and Policy*, De Gruyter, vol. 8(1), pages 1-55, June.

grades and spend less time studying. They authors instrument for study time using presence of a roommate with a video game and find that a one hour increase in study time per day (a .67 standard deviation increase in their sample) has the same effect of first semester GPA as a 5.21 increase in the ACT (an increase of 1.4 standard deviations in their sample).

Other researchers have investigated using technology to improve learning. **nbbm2020** conduct an experiment in Botswana during the COVID-19 pandemic and find that text messages and phone calls deployed as low cost, scalable learning technologies improved test scores by 0.16 to 0.29 standard deviations.

Effect of Recorded Lectures: Perhaps the educational resource most closely related to the IMVH is when the instructor’s recorded lecture is made available to students. Recorded lectures have course administration information which the IMVH does not have. Recorded lectures are much longer than a typical IMVH video, less organized, and it may not be clear what topics are covered in the lecture video and/or where in the video it is covered. Savage (2009)¹² taught two intermediate micro classes: one 42-student class used ”talk and chalk” and the other 45-student class used technology that allowed lecture capture which was then made available to the students. Students across two sections and does not find an effect on examine scores.

This study adds to this body of research by studying the effectiveness of an educational innovation: a video textbook. We use two empirical strategies to test for causal effects: within class randomization for students scoring below the median on the first exam and regression discontinuity at the median. We examine a large set of student study behaviors (lecture attendance, homework downloads from course web page, contributions to a discussion board, use of a class-specific tutoring lab) to determine if any of these study methods are substitutes or complements with video watching. We test for spillovers to other classes taken in the same quarter. We test for heterogeneous treatment effects using techniques that are robust to p-hacking. Finally, our research setting allows us to examine video views in the absence of the grade incentive in the next, required intermediate microeconomics class.

The Intermediate Microeconomics Video Handbook (IMVH) is a collection of 220 short-

¹²Savage, S. J. (2009). ‘The effects of information technology on economic education,’ *Journal of Economic Education*, vol. 40(4), pp. 337–53.

Table 1: Information Transmission Formats.

| | Lecture | Book | Lecture Capture | IMVH |
|--------------------------------|---------|------|-----------------|------|
| Instructor's time used | x | | | |
| Instructor-Learner Interaction | x | | | |
| Learner-Learner Interaction | x | | | |
| Readable | | x | ? | x |
| Scalable | ? | x | x | x |
| Searchable | | x | | x |
| Skimmable | | x | | x |
| Stoppable | ? | x | x | x |
| Watchable | x | | x | x |

ish videos that cover the material for a year-long intermediate microeconomics class. Most topics were covered by two videos: one video introduces the concept more intuitively with verbal explanations and graphs and the other video has the more formal, calculus-based definitions. The videos were created in 2014 by six UC San Diego faculty members with professional videographer and production support. Many videos were created using an innovative presentation technology, the "learning glass," where the instructor uses neon markers to write on a large sheet of glass, with lights embeded along the glass edge to make the colors pop. The remaining videos are a PowerPoint presentation that faculty could write on and the instructor is superimposed next to the PowerPoint explaining the concept. Videos are closed captioned. To help students find material in a video, videos are time stamped on the web page and captions are searchable while watching the video. To help students find topics, videos are organized into a book with both a table of contents and an index.

Table 1 presents a classification of some options to present course material to students.¹³ The IMVH differs from a traditional textbook because intructors explain, graph and derive mathematical results in much the same way one would in a lecture. The primary benefit of lecture is that students can stop the innstructor, ask questions and get answers in real time. There is also an important social aspect of lectures as students can interact with each other before, during and after lecture. The IMVH differs from lectures because students control the IMVH lecture: they can rewatch, speed up or slow down the lecture. Compared to a large lecture hall, all students can clearly see and hear the IMVH lecture. Finally, the IMVH

¹³This table is a modification of a classification Martin Osborne proposed to one of the authors in an e-mail correspondance.

is closed-captioned which may be particularly useful when English is not the first language of the instructor and/or the student. Compared to lecture capture, the IMVH may present the material differently than lecture

4 Experimental Design

We conducted the field experiment in four intermediate microeconomics courses, two in fall 2018 and two in fall 2019. The university is a large, diverse and selective public research four-year university. At this institution, intermediate microeconomics is a three-term sequence required for students majoring in Economics. The experiment was conducted in the first course of the sequence. We also have grades and video viewing in the second course of the sequence.

Students were told about the experiment in the first lecture and the syllabus. At any time during the quarter, students could opt out of having their data included in the analysis sample.¹⁴ Students below the age of 18 at the start of the course as well as students enrolled via Extension were removed from the analysis dataset.¹⁵

The experiment began four weeks into the term following the first midterm exam, and as such, only students who took the first midterm were included in the experiment. Following the first midterm, students were assigned to one of three treatment arms: *Incentive*, *Control*, or *Above median*. Students who scored below the median on the first midterm were randomly assigned to *Incentive* or *Control* whereas *Above median* includes all other students. Random assignment was issued using pairwise random assignment stratified on midterm score and year cohort (details on treatment assignment can be found in the Appendix). Students in the *Incentive* arm received a grading scheme that encouraged watching videos during the rest of the term whereas *Control* and *Above the median* received the standard grading scheme that does not directly reward watching videos. The two different grading schemes are outlined

¹⁴Students opt out via an online form owned by the Teaching + Learning Commons (T+LC) so that neither the instructor nor research team could observe which students decided to opt out.

¹⁵Students under the age of 18 were excluded per IRB protocol. We exclude Extension students because of their small number, their unknown and potentially very different preparation for the course compared to UC San Diego students, and we are missing most do not have information are so different from we do not know if they met the course prerequisites nor and

Table 2: Grade scheme by treatment arm. *Control* represents same grade scheme as *Above median*. Differences between the two grade schemes in red.

| Assessment | Incentive | Control |
|---------------------|-----------|---------|
| >40 videos | 4% | 0% |
| Midterm 1 | 18% | 22% |
| Midterm 2 | 22% | 22% |
| Final Exam | 50% | 50% |
| Math Quiz | 1% | 1% |
| Best 5 of 6 Quizzes | 5% | 5% |
| Total | 100% | 100% |

in Table 2. Specifically, the *Incentive* grading scheme requires that at least 40 of 48 eligible videos in the IMVH be watched to earn 4 percentage points towards the students’ final grades¹⁶. Notably, the 4 percentage points comes at the expense of the first midterm score, which had already occurred at the time of treatment assignment. Hence, we isolate the video incentive as the sole difference between treatment arms, giving us more confidence that the exclusivity assumption of our encouragement design holds.

which and there are two unique features of this class. First, many non-majors take the class typically to either satisfy general education requirements or to explore majoring in Economics. Thus there are many students at the margin of majoring in economics in the class and so an important outcome is the likelihood the student takes the second class in the sequence. The other unique feature of this class is the large fraction of transfer students, for whom the class is not only their first experience with upper division coursework, but typically the first time taking a class under the quarter system (community colleges in the state are on the semester system), and their first class at a 4-year research university. Thus, we expect transfer students to have greater information problems.

While the experiment was conducted in the first class of the sequence (100A), we also have data from the second class in the sequence (100B). Fortunately, one instructor taught all four of the 100A classes (one of the authors) and another instructor taught all four of the 100B classes. Both instructors created half of the videos lectures for their course.

Enrollments are high enough that the 100A and 100B instructors both taught two classes

¹⁶Watched in standard speed, 40 videos would require students to spend between 5.5 and 7.1 hours, depending on the length of videos chosen (on average 9.7 minutes in length each). Watching all incentivized videos in standard speed would require just shy of eight hours.

back-to-back each quarter but offered all exams at a common time out-of-class. So we treat the two classes in a quarter the same but account for different years in the empirical analysis.

The 100A course had an identical structure across all quarters and years. In addition to the textbook and the option to attend a live lecture at either 11-12:30 or 12:30-2:00, students had access to weekly one-hour discussion sections run by graduate TAs who were all Economics PhD candidates (including one of the authors), a tutoring lab staffed by both the graduate TAs (in lieu of office hours) and top undergraduates (some earning course credit for learning how to teach economics and others hired by the Department to help cover the tutoring lab hours, M-Th 5:30-8:30 and Sunday 4-8pm), weekly supplemental instruction sessions offered by an undergraduate majoring in Economics and trained by the university in supplemental instruction, a discussion board (monitored by the instructor as well as the grad and undergrad TAs), four years of the instructor's old exams (no answers provided), weekly ungraded problem sets (with detailed answers), graded online quizzes on Mondays of most weeks that did not include an exam, and a video handbook containing 220 short-ish videos covering the entire 100ABC sequence. This video handbook was created by the two instructors teaching 100A and 100B as well as four other faculty members who also regularly teach in the intermediate microeconomics sequence.

Students scoring below the median on the median on midterm 1 were randomly assigned to a "Videos Required" group where their midterm 1 score was down-weighted to 180 points and 40 points (all or nothing, so no partial credit) was put on watching 40 of the 47 remaining videos. Students were told that final letter grades would not be affected by being in the experiment, so the two key experimental outcomes are the scores on the second midterm and on the final exam. Following [cite], we randomized the students into the experiment by sorting the students by the first midterm score, creating ordered student pairs and then choosing one of students in the pair randomly to be in the experiment.

Video "watching" was based on opening the video, leaving it open for at least half the video length, and clicking a link at the end of the video which took students to a Qualtrics survey where they had to record their e-mail address on link at end of the video. Right after random assignment we surveyed students to make sure they knew which group they were in. Twice during the quarter we let students in the experiment know how many videos they

had watched

In particular, we use an assessment (exam) to identify students most likely to be struggling with meeting the learning objectives of the course. We require a random sample of students who score below the median on the first midterm to partake in a required (for one’s grade) study strategy consisting of watching at least 40 course-relevant supplementary videos. The videos were freely available to all students in the class and the instructor committed to having the course grade distribution identical across treatment and controls. To keep the weights on the second midterm and final exam, our primary outcome measures, the same across the treated and control groups, the first midterm was down-weighted for the treated group and the weight was put on watching 40 videos. Nearly all students in the treated group watched the 40 videos and treatment led to a X percentage point increase in video views relative to the control group. We find the grading policy led to an

We are primarily interested in the causal effect of watching videos on exam grades. The average causal effect of watching videos can be modeled using the potential outcome framework or Rubin Causal Model (**ir2015**):

$$\tau = E[Y_{it}(1) - Y_{it}(0)] \tag{1}$$

where τ is the average causal effect of watching videos and $Y_{it}(1)$ and $Y_{it}(0)$ are the potential outcomes (e.g exam scores) for a student i in year t who does and does not watch videos, respectively. We can observe treatment assignment for each student, $Z_{it} \in \{0, 1\}$, as well as the observed outcome $Y_{it}(Z = z_{it})$ and a vector of pretreatment covariates X_{it} .

A student’s decision to watch videos is endogenous, so we must rely on an exogenous instrument to calculate an unbiased estimate of τ . By randomly assigning treatment that induces significantly greater video watching, the treatment indicator Z_{it} satisfies the instrument validity and relevancy conditions required to estimate τ using an instrumental variable approach (citation). We calculate an estimate of τ using a two-stage least squares approach:

$$v_{it} = \alpha Z_{it} + f(X_{it}) + e_{it} \tag{2}$$

$$y_{it} = \tau \hat{v}_{it} + g(X_{it}) + u_{it} \quad (3)$$

where \hat{v}_{it} is instrumented videos watched estimated by Equation 2, $f()$ and $g()$ are generic functions through which X_{it} affects v_{it} and y_{it} , respectively, and e_{it} and u_{it} are model residuals assumed to be mean-zero conditional on observables. $\hat{\tau}$ is the estimate of the local average treatment effect (LATE) of watching videos local to those students induced by treatment to watch videos.

Under the assumptions of independence, monotonicity, and non-interference, $\hat{\tau}$ is an unbiased estimate of the LATE (ai1995). Independence assumes that outcomes (e.g. grades) are only impacted by treatment through watching videos. This assumption could be violated if, for example, telling a student she is treated were to give her more confidence on subsequent exams during the quarter. Monotonicity, sometimes referred to as the “no defiers assumption”, is required because of two-sided noncompliance and requires that students assigned treatment are weakly more likely to watch videos than if they were assigned control. A violation of this assumption could occur if students get utility from rebelling against their assigned grade scheme. Non-interference, also known as the Stable Unit Treatment Value Assumption (SUTVA), assumes that each student’s outcome depends only on their own treatment status and not the treatment status of their peers. Violations of SUTVA may include control students benefiting from having treated students in the same class (and perhaps studying together).

Although we believe excludability¹⁷ and monotonicity¹⁸ are reasonable assumptions, we have more pause about the non-interference assumption because of the potential for spillovers between students in the same class. If we had unlimited resources, a more robust experimental design would assign treatment at the class (or coarser) level, reducing the chance for interactions between treated and control students. However, given our resource constraints, assigning treatment at coarser levels would have resulted in insufficient statistical power to

¹⁷While this assumption is not testable, we took care in the experimental design to make the treatment and control arms as similar as possible except for the grading schemes

¹⁸Though not testable directly, one testable implication of monotonicity is that the cumulative distribution function of videos watched for each treatment arm should not cross. Indeed, in Appendix Figure ?? shows that the two CDFs do not cross.

detect large effect sizes. Hence, we proceed acknowledging the potential for spillovers between students. We hypothesize that spillovers likely bias our estimates of the treatment effect *downwards* as we believe control students are more likely to benefit from having well-studied peers than they are to lose from, for example, having peers too busy watching videos to join a study group.

We also estimate the average causal effects of being assigned to the *Incentive* arm, or average Intention To Treat (ITT) effects. These are ITT estimates and not average treatment effect estimates because of non-compliance: some students in the treatment arm do not watch videos and some students in the control arm do watch videos. However, since the incentive itself in our setting is representative of how future instructors may require their students to use the IMVH, the ITT estimates are relevant for an instructor interested in deciding whether to require her students watch videos or simply make them available with no grade incentive.

Following the methods proposed by **robinson1988** and more recently **wa2018**, our primary estimating equation is the partially linear model:

$$Y_{it} = \beta Z_{it} + f(X_{it}) + \epsilon_{it} \quad (4)$$

where Y_{it} is an outcome of interest (e.g. videos watched or test scores) for student i in year t , $Z_{it} \in \{0, 1\}$ is a treatment indicator, $f()$ is a generic function through which X_{it} , a vector of controls, affects Y_{it} , and ϵ_{it} is an unobserved residual. Following the Neyman-Rubin potential outcomes framework, $Y(Z = 1|X) - Y(Z = 0|X)$

To identify the causal effect parameter of interest, β , we make three assumptions: 1) *treatment status unconfounded given controls*: $Z_{it} \perp\!\!\!\perp (Y_{it}(Z = 1), Y_{it}(Z = 0)) \Big| X_{it}$; 2) *Excludability*: 3) *Non-interference/Stable Unit Treatment Value Assumption (SUTVA)*:

4.1 ITT

We are interested in knowing the average effect of being assigned to the *Incentive* arm on exam scores, which we will refer to as the Intent to Treat Effect (ITT). We will test the null hypothesis that being treated has zero effect against the alternative that being treated has

a nonzero effect on exam scores.

We will estimate the ITT using Neyman’s **neyman1935** repeated sampling approach, considering each pair (block) a completely randomized experiment and combining the results. We begin by estimating the point estimate of the ITT as the difference in means between the *Incentive* and *Control* arms, which is equivalent to the average difference in means across all pairs:

$$\hat{\tau} = \bar{\tau}_I - \bar{\tau}_C = \frac{1}{J} \sum_{j=1}^J \bar{y}_{j,I} - \bar{y}_{j,C} \quad (5)$$

where $\hat{\tau}$ is the point estimate of the ITT, J is the number of pairs in the sample, and other notation follows from above.

The point estimate $\hat{\tau}$ is appropriate if treatment is not confounded by observable or unobservable variables. While this assumption is true in expectation given random assignment of treatment, as a robustness check we estimate specifications that control for differences in observed characteristics, described later in this paper.

Following Neyman’s repeated sampling approach, the estimated standard error of $\hat{\tau}$ (**imai2008**, **ai2015**, **ai2017**) is:

$$SE(\hat{\tau}_j) = \left(\frac{1}{J} \sum_{j=1}^J \hat{V}(\hat{\tau}_j) \right)^{\frac{1}{2}} \quad (6)$$

where $\hat{V}(\hat{\tau}_j)$ is the estimated variance within block (pair) $j \in \{1, \dots, J\}$. This within-block variance given one control and one treated unit per block is (**ir2015**, **ai2017**):

$$\hat{V}(\hat{\tau}_j) = s_{j,I}^2 + s_{j,C}^2 \quad (7)$$

where $s_{j,I}$ and $s_{j,C}$ are the *Incentive* and *Control* sample variances within block j , respectively. Unfortunately, these sample variances are not estimable with only one unit in each arm per block. As such, we use the following estimator, which is conservative (confidence intervals wider) if there is heterogeneity in the treatment effect (**imai2008**, **ai2015**,

ai2017):

$$\hat{SE}(\hat{\tau}_j) = \left(\frac{1}{J(J-1)} \sum_{j=1}^J (\hat{\tau}_j - \hat{\tau})^2 \right)^{\frac{1}{2}} \quad (8)$$

However, it is impossible to estimate this estimator in the special case when

Assuming a constant treatment effect within the population, the point estimate of the ITT is simply the difference in means of the outcome in the sample.

where $z_i \in \{0, 1\}$ is the treatment assignment for student i , y_i is the student's observed outcome, and N is the number of students in the sample.

Since treatment is assigned using a matched pair design,
the appro

$$Y_{it} = \beta Z_{it} + f(X_{it}) + \epsilon_{it} \quad (9)$$

5 Results

6 Discussion

7 Conclusion

We examine the effectiveness of an educational innovation, a video handbook, where 220 (mostly) short instructional videos were organized into a book. Instructors may have concerns about making a new resource available if they believe students will substitute away from more productive study methods. Our experiment focused on students who demonstrated that they may benefit from an educational intervention by scoring below the median on an early assessment. We randomly assigned half of those scoring below the median on midterm 1 to a new grading scheme: the first midterm was down-weighted 4 percent and the points were fully awarded if the student watched 40 videos. Most treated students watched the 40 videos and the grade inducement led to a doubling of the number of videos watched relative to the control group. Both the second midterm and the final exam scores were

significantly higher for the treated group with effect sizes of .18 of a standard deviation.

Our field experiment is unable to determine whether the videos uniquely helped the treated students or whether any intervention that induced additional study time would have been effective. However, we have suggestive evidence that the videos were uniquely effective: treated students in the next course in the sequence, where there was no grade inducement to watch videos, continued to watch the videos at a higher rate than the control students. This result is consistent with a model of student learning with imperfect information and not consistent with the neoclassical model or procrastination model of student learning.

Generalizability of our results. Experiment was conducted in large classes with historically high failure rates. It is the first upper division class for many students and, for transfer students, also their first class in the university. As such there may be large information problems about how to successfully study for the class.

Future research: most importantly, to see if our results hold up in other educational settings e.g., different student types, types of classes, instructors and universities. We would not allow binge-watching by requiring a certain number of videos each week as opposed to our experiment which simply required the student to watch 40 videos before the end of the quarter.

Appendix

A Additional experiment details

In this section we outline additional experiment details that could prove useful for replication or understanding our analysis choices.

A.1 Randomization

Students were assigned to treatment arms using a matched pairs design, a special case of blocked randomization in which each block contains exactly two units, one treated and one control. Several authors detail how matched pair designs can improve the *ex ante* precision of treatment effect estimates (versus complete randomization) by matching treatment units whose potential outcomes are similar (e.g. [ir2015](#), [ai2017](#)). The

Additionally, we were unable to observe most pretreatment covariates until after the experiment had finished because of student privacy concerns, making blocking on these variables impossible. We learned from the previous cohorts’ data that between the first midterm score and math quiz score, the midterm score predicted significantly more variation in the final exam score. Hence, we ass

A.2 Selection of control variables

In this section we discuss how we select control variables included in linear models estimated in this paper. In our nonparametric estimation strategies that use generalized random forests, the algorithms implicitly select control variables that matter most for treatment heterogeneity and explaining variance in the outcome variables of interest (see Appendix A.3).

Equation ?? includes a vector of control variables related linearly to the outcomes of interest. Although d_i , the treatment indicator is randomly assigned and in expectation d_i is orthogonal to all observed and unobserved pretreatment covariates, in small samples stochastic imbalances can occur, which if controlled for can reduce bias of the treatment effect estimator (ai2017). Even if perfect balance is achieved, controlling for orthogonal covariates can improve precision of the treatment effect estimator if the covariates can predict unexplained variance in the outcome.

By definition it is not possible to guarantee balance on unobserved covariates. As discussed in Appendix A.1, we mechanically balanced the treatment arms on first midterm score, one of the few observables at the time of treatment assignment, with our knowledge from previous cohorts’ data that the first midterm score explains a significant portion of variance in final exam score. Hence in our estimation strategies involving controls we always include the first midterm score and year, following the recommendations of bm2009 to control for all covariates used to seek balance when assigning treatment.

For variables unobservable at time of randomization but observable at time of analysis, we cannot guarantee balance, nor is it clear *ex ante*, beyond our intuition, which will predict variation in the outcome variables of interest.

A.3 Generalized Random Forests

In this section we discuss our implementation of Generalized Random Forests proposed by [atw2019](#)