

The effect of a supplementary video book on learning in intermediate microeconomics

Melissa Famulari and Zachary A. Goodman*

University of California, San Diego

This version: March 2020

Abstract

In this paper we estimate the effectiveness of a novel video-based textbook replacement for intermediate microeconomics, the Intermediate Microeconomic Video Handbook (IMVH), on learning outcomes. In a field experiment involving nearly 400 students, we randomly assigned a grade-based incentive that induced treatment students to watch over 60% more videos than did control students. We observe significant reduced form effects: being assigned treatment caused students to score 0.18 standard

*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 the University of California. As of fall 2020, one of the authors has a financial interest in the distribution of the IMVH outside of the University of California.

deviations higher on midterm and final exams. Using an instrumental variables approach, we estimate that the marginal hour of video content increased exam scores between 0.05 and 0.15 standard deviations. We rule out large negative spillover effects to other courses taken concurrently, and we observe persistent take-up of the IMVH in the subsequent quarter when direct incentives to watch videos were removed.

1 Introduction

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

— Anonymous student

Every year, university students spend tens of thousands of dollars on tuition and hundreds of hours studying, in large part, to learn. Instructors can help their students learn more efficiently by providing and recommending pedagogical tools that have high returns per unit time and financial cost. Despite the value to students and instructors, little empirical work exists that estimates the effectiveness of different learning technologies (Allgood, Walstad, & Siegfried, 2015).

In this paper, we measure the impacts of one such technology, the Intermediate Microeconomic Video Handbook (IMVH), on outcomes in an intermediate microeconomics course. The IMVH was designed to *supplement* lecture 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 modern students, who have had unprecedented exposure to electronic media, may find videos more engaging and, perhaps, more effective at building human capital.

Besides higher engagement than with conventional texts, the IMVH and video-based learning tools more broadly are of value to university educators and their students for four additional reasons. First, videos carry near-zero marginal cost and are accessible anywhere via internet, helping reduce financial and geographic barriers to high-quality educational materials. Second, video platforms can help students track what content they have already studied and what content they have yet to cover. Third, features embedded into videos,

such as searchable captions and timestamps, can help students locate the information they are seeking faster. Finally, the perceived low cost of watching a brief video may be easier to overcome than perceived higher costs of other studying methods, potentially leading to more frequent studying, which decades of psychological research has demonstrated leads to more long-term learning than does cramming (Cepeda et al., 2006; Kornell, 2009).

Ultimately, the beneficial features of video-based learning tools are of value only if they can improve student learning outcomes, an empirical question we seek to answer in this paper. To estimate the effects of the IMVH, we administered a field experiment involving nearly 400 undergraduates enrolled in the same one-quarter-long microeconomics course over two years. Of note, these students all scored below the median on the first midterm exam, thereby making manifest a need to adjust studying habits, and perhaps standing to gain the most from an intervention that targets studying. We randomly assigned a grade-based incentive to half of the sample to encourage take-up of the IMVH, which was made available to all students in the experiment, allowing identification of intent to treat (ITT) effects and local average treatment effects (LATEs) while maintaining equitable access to learning resources. We tracked video watching at the student level using the software platform that hosts the IMVH. We observe grades, GPA, and video watching in both the term of the experiment and the subsequent term.

The first-stage impact of the exogenous encouragement on video watching is significant and substantial. Students who receive the grade-based incentive watch over 28% more unique videos by the second midterm and 63% more unique videos by the final exam, or about 1.1 and 3.4 hours of content, respectively, than did their control peers. We find large reduced-form effects of treatment on exam scores: being assigned treatment (ITT) increases midterm and final exam scores by about 0.18 standard deviations. Our estimates imply that the marginal hour of videos watched increases exam scores (LATE) by between 0.05 and 0.16 standard deviations.

We interpret our results through a theoretical framework in which students, who value grades and leisure, have potentially incorrect priors about the returns of different studying methods. We do so not to test theory, but rather to help educators understand the potential welfare implications of providing students with a learning technology and paternalistic

incentive structure like the one studied. Although we observe that treatment students performed better on course assessments, for welfare analysis one must consider where the time watching videos came from: leisure, work, student organizations, studying for other classes, studying for present class using other methods, etc. If students must reduce time allocated towards leisure or studying for other classes so they can watch more videos, then the welfare impacts of requiring videos could be negative depending on the students' preferences. On the other hand, if the videos are more productive than the next best studying technology, then requiring videos could be utility enhancing.

To better understand the spillover effects of treatment, we examine other forms of studying including class attendance, visits to a tutoring center (specific to this course), and interacting with the class discussion board. We do not find any statistically significant changes in any observed studying method, and we can rule out large changes. In nearly all cases, treatment students used other studying methods at directionally *greater* rates than did their control peers. We also investigate spillovers to other courses taken during the term of the experiment and similarly find that treated students perform directionally *better* than their control peers. Though not statistically significant, we can rule out large negative effects, suggesting that treatment did not cause students to dramatically substitute away from studying for other courses.

An important piece to the welfare puzzle is whether treatment students continue to use the IMVH at higher rates after exogenous incentives are removed. Persistent take-up in the absence of external prodding provides some confidence that students, now with updated priors, value the technology. Fortunately, we can observe video watching in the subsequent microeconomics course in the term following the experiment. Despite there being no direct incentives to watch videos in the subsequent course, treatment students persistently watched more videos than did control students, about 8 - 10 more unique videos, or 1.2 - 1.5 more hours of unique content. Our sample in the subsequent term is nearly half the original size, so we lack power to precisely estimate effects on exam scores; however, our confidence intervals include effect sizes consistent with those observed in the experiment term.

Collectively, we interpret our findings as evidence that requiring the IMVH is a net positive on underperforming students' academic achievement, both in the quarter of the ex-

periment and beyond. Though formal welfare analysis is beyond the scope of this paper, we present suggestive evidence that requiring the IMVH is unlikely to be substantially utility harming, if not utility enhancing. We believe our findings justify paternalistic incentive structures in settings where a large portion of the class is at risk of failing and the instructor has more information about the usefulness of a novel teaching technology than do her students.

The rest of the paper is organized as follows. Section 2 presents competing models of studying behavior that may explain the observed phenomena. Section 3 provides background on existing related literature. Section 4 describes the study 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 studying behavior: 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 encourage student 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 do not 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 study time to each method given their utility function, which is increasing in leisure and grades and decreasing in time spent studying. Since college instructors have little knowledge about student utility functions, they do not know student preferences over performance in other classes and other aspects of their lives that may have large payoffs in the labor and marriage markets and so there is no room for an instructor to increase student well-being by intervening in their study decisions. Oettinger (2002) provides empirical support for the neoclassical model by demonstrating that student effort responds rationally to nonlinear grade incentives. Across

1200 students in a principles of economics class with absolute grading standards, he finds evidence of bunching just above letter grade cutoffs and student performance on the final exam is higher if the student is just below a grade threshold.

Kapoor, Oosterveen, and Webbink (2020) also finds support for the neoclassical model for students at the margin of poor performance. The authors use a regression discontinuity approach to examine the effects of a university policy that required students who performed poorly in their first year to attend at least 70 percent of the tutorials for each class in their second year. In courses where tutorial attendance was not required for all students, the policy increased both tutorial and lecture attendance by over 50 percent, did not increase total study time and reduced grades by 0.16-0.26 standard deviations. The policy had its biggest impact on students who lived far from campus and those who were most likely to miss section in their first year. The authors conclude that the university policy prevented students from using their optimal mix of study methods and so reduced grades.

A key assumption of the neoclassical model is that students possess complete information about the returns across studying methods. However, there is evidence from psychology that college students do not know the return to various study methods¹ and many universities fund “Teaching and Learning Centers” or “Academic Skills Centers,” part of whose mission is to help undergraduates learn to study more productively.² Further, the “raison d’être” of higher education is not only to teach students specific skills but to teach students how to learn. As an alternative to the neoclassical model, we hypothesize that students supply a quantity of study time that is optimal given their information constraints. In this ‘imperfect information’ model, students choose study methods and quantities that are suboptimal relative to those they would have 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 enhance student utility.

A third model is a behavioral one in which students plan to study more than they end up studying when the time comes. This phenomenon is consistent with two-self models in

¹See, for example, McCabe (2011), Pashler et al. (2007), Dunlosky et al. (2013)

²All nine University of California campuses have such a center. Examples outside the UC 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.

which a person’s “planner” self, the one who desires high grades at the expense of leisure, is at odds with her “doer” self who must choose between immediately gratifying leisure and delayed gratification from higher grades. 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,³ Clark et al. (2020) provide empirical support for this model by finding that setting task-based goals helps improve college student performance. As descriptive evidence in support of this model, Beattie et al. (2019) find that students that do much worse than expected in college are those who say they have poor time management or procrastination issues, including a tendency to cram and spending very little time studying.

We consider the testable implications of the three models applied to a setting where students are incentivized to use a time-consuming educational input, say, a set of instructional videos (or attending class, reading the textbook, working on homework, etc.). The incentive is structured such that students who consume the educational input receive a higher grade in the course by consuming a set level of the input. In this simple setting, students gain utility only from leisure and grades. We assume grades, a function of time spent studying, and utility are both continuous, smooth, and increasing and concave in their inputs. Students can choose to study using the incentivized educational input or some outside option that is not directly incentivized (or a combination thereof).

Across all three models, before the first educational input is incentivized, students allocate time to the two studying methods until the marginal benefit of each (through higher grades) is equal to the marginal cost of forgone leisure. Consider the population of students initially consuming below the requisite level to earn the grade incentive. These students must decide if earning the grade incentive is worth forgone leisure and less time allocated to their outside studying option. Next we explore the differences in predictions across the three models.

In the neoclassical model, the marginal return to grades of the incentivized input is less than that of the outside option for the ‘compliers’, or those induced by the incentive to consume at least a fixed level of the incentivized input. This model predicts bunching at the incentivized level cutoff since compliers would prefer to spend their marginal hours on leisure

³see, for example, Ferrari (1992), P. Chen et al. (2017) and Lavecchia, Liu, and Oreopoulos (2016).

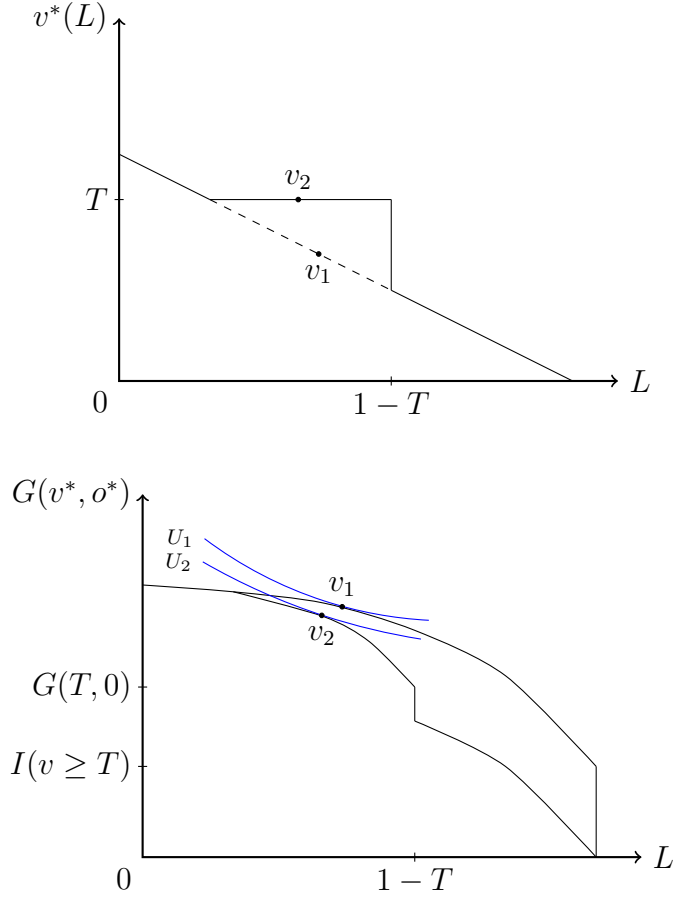


Figure 1: *Above:* Demand curve for video watching as a function of leisure L . At $L = 1 - T$, the student maximizes grades $G(v, 0)$ by spending all studying time watching videos, i.e. $v^* = T$. *Below:* Student's utility maximization problem for the neoclassical model. The student maximizes her utility over leisure L and grades G , which is a function of time allocated to video watching v and her next best studying option o . The grade incentive I is given to the student conditional on watching T hours of videos (inner-time budget constraint) or, in the unincen-tivized case, given regardless of video watching (outer time-budget constraint).

or studying with their other method. This model predicts a strict increase in video watching and weak decrease in other studying and leisure consumption. It is ambiguous whether cumulative study time increases or decreases as this depends on relative utility benefits of leisure and grades and the returns to studying by each method. However, if cumulative study time remains constant or decreases, then exam performance should strictly decrease since students are now suboptimally allocating study time versus their first-best allocation when considering only marginal returns to studying. On the other hand, if cumulative study time increases, students may earn greater exam grades but achieve lower utility compared to baseline. Importantly, this model predicts that in subsequent quarters students return to their pre-incentive levels of studying.

In the imperfect information model, students' *ex ante* allocations to each studying method are not necessarily first-best. Compliers update their priors about the returns to watching videos as they work towards hitting the minimum required level. At this cutoff, they make a decision whether to continue watching videos depending on their updated perceptions of the marginal benefit. We do not expect bunching in this model unless poorly informed students believe the instructor set the cutoff for the grade incentive at the optimal level or the updated marginal benefit at the cutoff is lower than the marginal benefit of the next best studying option or the marginal utility of leisure. A sharp prediction is that video watching will continue at the incentivized level in the absence of the grade incentive as students have learned an effective study tool. We also expect the treatment effect to be greater for students with more information problems, perhaps students in their first semester/quarter at university.

Finally, in the behavioral model, the instructor's inducement helps students stick to study plans up until the cutoff. As long as total study time does not fall, the inducement will increase exam performance. This model also predicts bunching at the incentivized level cutoff as long as using the incentivized input does not change the student's "planner" and "doer" selves. In the absence of the inducement, i.e., in future classes, a sharp prediction is that video watching will revert 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 students could use to learn microeco-

nomics (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 feature of the IMVH).

3 Related Literature and Contributions

Students have many time-consuming activities to help them learn including attending class, watching recorded lectures, reading the textbook, doing homework or practice exams, or attending tutoring labs, and there are many empirical challenges to estimating the causal effects of a learning activity. First, a student’s decision to use a study method is likely influenced by unobservable student characteristics, such as motivation or ability, that also likely affect student exam performance. To estimate causal effects, researchers must address nonrandom selection into using the study method. Second, motivated students often visit instructors seeking help improving their study strategies after a negative exam shock which suggests ”dynamic selection” into the use of a study method and has been found empirically by Oettinger (2002), Krohn and O’Connor (2005), Ralph and R (2008), Bonesrønning and Opstad (2012) and Bonesrønning and Opstad (2015).⁴ Dynamic selection means including student fixed effects in class performance regressions will not uncover the causal effect of a study method. Third, study methods may be substitutes or complements in student learning and experimental inducements to use one study strategy may affect student use of another. In these cases, even randomized experiments will not identify the causal effects of a particular

⁴Oettinger (2002) finds that students close to a grade threshold before the final exam perform better on the final. Krohn and O’Connor (2005) show that students reduce the number of hours they study after getting higher midterm scores. Ralph and R (2008) find that IV estimates of studying on grades are much larger than OLS and using two years of data, provide suggestive evidence that students increase effort in semesters when semester-specific elements of grades are low. Bonesrønning and Opstad (2012) and Bonesrønning and Opstad (2015) find that student’s self-reported study hours

study method but will instead identify the causal effects of a study policy and all of the changes in student behavior caused by that policy.⁵ Finally, experimental inducements to use a study method may change the total time devoted to the course. In this case, experiments jointly test the effectiveness of a particular learning method and devoting more (or less) time to the course.

We focus our review on studies that use experiments or quasi-experiments (regressions discontinuity or instrumental variables) to explore the effects of learning acquisition that take student time. We put these studies in two broad groups: guided-study such as attending lecture, discussion sections/tutorials/recitations (small problem-solving sessions offered by TAs), or tutoring labs and self-study which includes doing homework, practice exams and watching recorded lectures.⁶ Kirby and McElroy (2003) use student-reported travel time to campus as an instrument for lecture attendance. They find no selection bias and so find the usual positive relationship between lecture attendance and exam grades in OLS regressions. Dobkin, Gil, and Marion (2010) analyze a policy where lecture attendance was voluntary before the midterm but after midterm grades were posted, lecture attendance was required for students who scored below the median. At the threshold there was a 36 percentage point increase in post-midterm attendance. Using a regression discontinuity approach to estimate the effects of attendance on the final exam, they find a 10 percentage point increase in student attendance led to a .17 standard deviation increase in final exam score. Joyce et al. (2015) randomly assign 725 students taking introductory microeconomics students to twice per week and once per week lecture formats to identify the effects of classroom time in classes that are well-supported with online content (videos, quizzes, posted PowerPoint lectures, etc.) Students in the twice per week format scored 0.21 standard deviations higher on the

⁵While the causal effects of an educational policy are useful for educators considering how to design their classes they are less useful for students wanting to know the most productive use of their study time. We should also point out that instructors may find learning transmission methods substitutable or complementary. As an example, in Morris, Swinnerton, and Coop (2019) many instructors report that lecture capture, where lectures are recorded and made available to students, changed the way they lectured in the classroom. Experiments randomly assigning students to classes taught one way versus another will not identify the causal effect of a study method if other aspects of learning transmission are simultaneously changed.

⁶Since using the IMVH takes student time, we exclude from our review studies that examine how to make a learning transmission method more effective per hour. For example, there has been a substantial amount of research exploring how to make an hour of lecture more effective by using games or experiments in lecture, student response systems, group work, etc.

midterm and 0.14 standard deviations higher on the final exam. Finally, Tang et al. (2020) randomly assign students to either weekly or bi-weekly grading of in-class clicker questions. Weekly grading increased student lecture attendance by 11 percent but not student-reported self study hours. Course grades were 6.31 percent higher for students randomly assigned to weekly grading and the effects of weekly grading were much stronger if the student preferred bi-weekly grading, had lower prior GPAs and lower self-control scores.⁷

Arulampalam, Naylor, and Smith (2012) study discussion section (as opposed to lecture) attendance. Since students are randomly assigned to sections and section attendance depends on day of week and time of day, they use these variables as instruments. They find no effects of section attendance for most students but for students in the top quantiles missing 10 percent of sections resulted in a 1 percentage point performance loss. Kapoor, Oosterveen, and Webbink (2020) also examine the effects of section but like Dobkin, Gil, and Marion (2010), study a university-wide attendance policy focused on students who reveal themselves as poor performers. For students whose first year GPA was below a threshold, a European university required 70 percent section attendance in the second year. Nearly all affected students met the attendance threshold since there was a substantial penalty for not meeting it. They find students at the policy threshold attended 50 percent more sections and lectures but reported no significant difference in total study hours (lectures plus section plus self study). Students at the policy threshold earned grades that were 0.16 to 0.24 standard deviations *lower*. Bratti and Staffolani (2013) have data on attendance for lecture and discussion sections combined. Their attendance measure is highly correlated with student reported self study hours and including self study hours significantly reduces the estimated effects of attendance. Using student residence as an instrument and controlling for self study hours, the authors find significantly positive attendance effects on grades only for quantitative courses.

In addition to lecture and discussion sections, students can go to tutoring labs to get help learning the material. Munley, Garvey, and McConnell (2010) find athletes are significantly

⁷J. Chen and Lin (2008) randomly omit material from lecture from a randomly chosen class and find students who heard the material did better on related multiple choice questions. However since both treatment and controls are in lecture, this study abstracts from the time costs of attending lecture and we focus on empirical studies exploring study methods that take time.

more likely to attend peer tutoring labs which they attribute to frequent reminders from coaches. Using being an athlete as an instrument for tutorial hours, they find significantly positive effects of peer tutoring: to increase a student's letter grade, they would need to spend about one hour per week being tutored over a 14 week semester.

These studies on the effects of guided instruction consistently find that student performance is significantly improved by attending lecture, .17 to .21 standard deviations, but discussion sections may only benefit some students in some courses and, if they substitute for student self-study, can have significant negative effects on performance, -.16 to -.24 standard deviations. The one study on peer tutoring labs, found significantly positive effects.

Turning to the effectiveness of a student's own study time, Ralph and R (2008) examine 210 Berea College students who were randomly assigned a roommate. Students whose roommate brought a video game to college, earn lower grades and spend less time studying. They authors instrument for study time using presence of a roommate with a video game and find that an additional hour of study time increases first semester GPA by .36. However, Oreopoulos et al. (2019) randomly assign first year economics students at three different institutions to a variety of low touch interventions to increase study hours (help planning study schedules at the start of the semester, information about the value of studying, and weekly reminders about study plans) which increased student's self-reported study hours by two hours per week at two of the three institutions but had no effect on grades, credit accumulation, or retention. Clark et al. (2020) also explore a low touch intervention, having students set goals on the number of practice exams they will complete, and find those randomly assigned to set goals completed 0.102 of a standard deviation more practice exams and increased total course points by .068 of a standard deviation. Trost and Salehi-Isfahani (2012) examine the causal effects of homework by randomly requiring two-thirds of the students in Principles of Economics classes to complete one of three homework assignments for a grade while the other third may complete the homework, but it does not contribute to their grade. The outcome is exam scores on questions related to the three homework assignments. They find significant effects of homework on the first midterm but not the final exam. Grodner and Rupp (2013) use within-class randomization to assign 423 students in four principles of microeconomics courses to homework required (10 percent of grade) and not required groups

(the homework grade incentive was spread across four exams). They find that 90 percent of the homework-required students completed 7 or more homework assignments compared to 6 percent of the controls. Treated students did significantly better only on the first two of four exams but the average across all four exams was higher.

Finally, we review the research on flipped classrooms and lecture capture, both of which have aspects similar to the IMVH. With lecture capture, lectures are recorded and made available to students and so cannot be used to prepare for lecture. In flipped classrooms, students can watch the material when they want but it is unclear whether the flipped classroom adds to student time. Two interesting RCTs using military students shed some light on the flipped classroom. In a small RCT (137 students in 7 sections of an introduction to economics course), Wozny, Balser, and Ives (2018) assign lecture type (flipped vs traditional) to randomly chosen topics and randomly chosen instructors. They give students six short term exams (based on previous three lessons), four medium term exams (based on previous eight lessons) and one long term exam (a comprehensive final). They find positive effects of the flipped classroom, 0.16 standard deviations, only for medium-term assessments which, as the authors note, is surprising given the forced preparation before lecture in the flipped model would seem to confer short term benefits over traditional lecture. In a larger RCT at West Point (1,328 students in 80 course sections and 29 instructors) across two courses, Principles of Economics and Introductory Calculus, Setren et al. (2021) take a single course topic and randomly assign half the sections to cover that topic using the flipped model and half to cover the topic using a traditional model. The outcome is score on questions related to the topics on a low-stakes quiz and the final exam. They are able to include fixed effects for both the instructor and class hour. They only find significant positive effects of a flipped classroom on the math quiz.

This study adds to this body of research by examining the effectiveness of an educational innovation in self-study: a video textbook. We randomly assign half the students scoring below the median on the first midterm to a grading scheme which placed 4 percent, or 40 points, of the student's grade on watching 40 videos and down-weighted the first midterm by 4 percent. This experiment allows two empirical strategies to test for causal effects: within class randomization for students scoring below the median on the first exam and a

regression discontinuity approach at the median first exam score. We test for heterogeneous treatment effects along key demographics to see if some groups benefit from the videos more than others. We examine a set of alternative studying methods (lecture attendance, contributions to a discussion board, use of a class-specific tutoring lab) to understand if any of these are substitutes or complements with video watching. We test for spillovers to other classes taken in the same quarter. Finally, since we can observe video watching in the subsequent intermediate microeconomics class, we look for behavior stickiness in the absence of exogenous incentives.

4 Study Design

4.1 Description of the sample and institution

We conducted the field experiment in an undergraduates intermediate microeconomics course taught during fall 2018 and fall 2019 by one of the authors. The university is a large, diverse and selective public research university in the United States.⁸ At this institution, intermediate microeconomics is a three-quarter sequence required for students majoring in Economics. The experiment was conducted in the first course of the sequence, *Micro A*. We also observe grades and video watching in the second course of the sequence, *Micro B*, which was taught by the same instructor during both the winter 2018 and winter 2019 quarters. Both Micro A and B instructors created half of the videos relevant to their course in the IMVH.

The structure is similar across the three courses in the Micro sequence. Students have the option to attend one of two lectures offered back to back twice per week, each lasting about 90 minutes. Two midterm and final exams are held at a common time outside of lecture. In addition to lecture, students have access to weekly one-hour discussion sections run by

⁸The Carnegie Classification of Institutions of Higher Education classifies the university as an R1 (very high research activity) university. For the 2017-2018 academic year, the undergraduate student body shared the following demographics: 49.1% female, 50.6% male; 75.0% in-state, 5.5% out-of-state, and 19.5% international; 59% students of color; 28.6% majoring in the social sciences, 26% of which major in Economics. Among newly admitted students, about one-third were transfer students, and average SAT scores were 652 and 605 for math and critical reading, respectively. About 34% of students are the first in their family to attend a four-year university.

graduate teaching assistants (TAs) who are all Economics PhD candidates, including, at the time, one of the authors. In lieu of office hours, the TAs and Undergraduate Instructional Assistants (UIAs) staff a tutoring lab open between three and four hours per day, six days per week. Students may also attend weekly Supplemental Instruction (SI) sessions offered by undergraduates majoring in Economics and trained by the university in SI. Besides the IMVH, students have access to a variety of online learning resources including a discussion board moderated by the instructional team, four years of previous exam questions, weekly ungraded problem sets, and semi-weekly graded online quizzes.

Students were told about the experiment during the first lecture and provided an informed consent form in the syllabus. At any time during the quarter, students could opt out of having their data included in the analysis.⁹ Students below the age of 18 at the start of the course as well as students enrolled via the university’s extension program were removed from the analysis dataset.¹⁰ Ultimately, four students under 18, five extension students, and seven students who opted-out were removed from the analysis dataset, leaving a sample of 850 students.

There are two unique demographic features of the class worth noting. First, many non-econ majors take the class to either satisfy general education requirements or to explore majoring in economics. As there are many students in the experiment on the margin of majoring in economics, an important outcome is the likelihood the student takes Micro B. Second, about 37% of the class is transfer students, for whom the class is not only their first experience with upper division coursework at a four-year research university, but also typically their first time taking classes under the faster-paced quarter system.¹¹ We examine treatment effect heterogeneity to understand whether transfer students might differentially benefit from the IMVH.

⁹Students could opt out via an online form visible to a third party university organization so that neither the instructor nor research team could observe which students elected to opt out.

¹⁰Students under the age of 18 were excluded per IRB protocol. We exclude extension students because of their potentially very different preparation for the course and our inability to observe pretreatment covariates and outcomes outside of Micro A.

¹¹Community colleges, the most common previous institution for transfer students, are on the semester system in the state of the university.

4.2 Description of the IMVH

The Intermediate Microeconomics Video Handbook (IMVH) is a collection of 220 short videos that cover the material in a year-long intermediate microeconomics course sequence.¹² The videos, designed to complement or replace a course textbook, include graphical and verbal intuition as well as formal algebraic and calculus-based definitions and proofs.

The videos were created by six UC San Diego faculty members with professional videographer and production support. Many videos utilize the “learning glass,” an innovative presentation technology where instructors write with neon markers on a large sheet of glass that has lights embedded along the glass edge to make the colors pop. The remaining videos feature faculty superimposed in front of slides. Videos are closed captioned and were checked by graduate students for accuracy.

Given the complexity of the material, a key objective was to keep the web interface clean and simple so as not to distract from the content. The videos are organized by content area (e.g., consumer theory, producer theory, etc.) that help students understand where various topics “live” in intermediate microeconomics. When considering the design of the platform, one goal was to help students find material quickly. Besides a table of contents and an index, each video contains time stamps of the concepts therein. Another helpful feature is searchable captions, which allow the student to jump to the part of a video containing the searched-for word.

While we do not know of another textbook completely comprised of videos, the IMVH is similar to Khan Academy, lecture podcasts, and textbook websites that incorporate instructional videos. Besides the engaging viewable nature, the IMVH differs from a traditional textbook in that the instructors explain, graph, and derive mathematical results in much the same way one would in a conventional lecture. However, the IMVH differs from lectures in that students control the pace: they can rewatch, speed up, or slow down the videos. Another difference is the ability to read captions and the reduced demand on the instructor’s time.

Table 1 presents a classification of some options to present course material to students.¹³

¹²A preview of the IMVH can be found at https://iti.ucsd.edu/IMVH_Misc/Promo/IMVHPromo.html.

¹³This table is a modification of the classification Martin Osborne proposed to one of the authors in an

Table 1: Comparison of information transmission formats

Feature	Lecture	eText-book	Lecture Capture	IMVH
Instructor's time used	✓			
Instructor-learner interaction	✓			
Learner-learner interaction	✓			
Readable		✓	?	✓
Scalable	?	✓	✓	✓
Searchable		✓		✓
Skimmable		✓		✓
Stoppable	?	✓	✓	✓
Watchable	✓		✓	✓

The IMVH differs from a traditional textbook because instructors verbally explain, graph and derive mathematical results in much the same way one would in a lecture. The primary benefit of lecture over the IMVH is that students can stop the instructor and get an answer as soon as they are confused, or wonder about a connection of the material with their lives or other courses, or want to know if they are right about an extension of the material, etc. Further, student questions may have import externalities for the learning of other students in the class. There is also an important social aspect of lectures as students can interact with each other before, during and after lecture. Students cannot ask questions or interact with other students during an IMVH lecture. Some of the advantages of the IMVH over lecture are first, students control the IMVH in that they can rewatch, speed up, or slow down an IMVH video. Second, compared to a large lecture hall, all students can clearly see and hear the IMVH presentation. Finally, the IMVH is closed-captioned, which may be particularly useful when English is not the first language of the instructor and/or the student. The IMVH differs from lecture capture because the IMVH videos are much shorter, averaging under ten minutes. The IMVH web site is well organized so students can see where the topic “lives” at all times. A useful feature of recorded lectures is they have course administration information which the IMVH does not have. However recorded lectures may also include components that do not work well when recorded, such as group work or class discussion.

e-mail correspondence.

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

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%

4.3 Experiment Design

The experiment began four weeks into the term following the first midterm exam. All students who scored above the median on the first midterm, the *Above median* arm, and half of students who scored below the median, the *Control* arm, were assigned a conventional grading scheme that places weight only on exams and quizzes. We assigned the remaining half of students below the median to the *Incentive* arm, whose grading scheme allots four percentage points conditional on watching at least 40 of 48 eligible videos in the IMVH.¹⁴ These 48 videos review new class content since the first midterm that could be assessed in the second midterm and final exam. Students could still view the remaining 26 course-relevant videos, and as they could have helped students on the cumulative final exam, we include them in our measures of video watching despite not counting towards the grade incentive.

The two different grading schemes are outlined in Table 2. Notably, the four percentage points come at the expense of reduced weight placed on the first midterm score, which had already occurred at the time of treatment assignment. Hence, at the time of treatment assignment, the video incentive is the sole forward-looking difference between treatment arms.

To improve balance between *Incentive* and *Control* arms and increase statistical power, we assigned students to treatment arms using paired randomization (Athey & Imbens, 2017), matching students by their first midterm scores before randomly assigning one member of

¹⁴Watched 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 48 incentivized videos in standard speed would require just shy of eight hours.

each pair to *Incentive* and the other to *Control* (further details on treatment assignment can be found in Appendix A.1). We emailed each student letting them know their assignment and grading scheme. Students could also find their assignment listed in the online gradebook. To confirm that students knew their assignment, we surveyed students using an in-class attendance quiz, and 94% of students correctly identified their grading scheme. We emailed the students who responded incorrectly to clarify their assignments.¹⁵

We informed *Incentive* students that they must watch the entire video and only one video at a time to get credit towards their 40 required videos. It is not possible to observe “watching” as students could, for example, minimize their browser, walk away from their computer, or otherwise play a video without actively watching it. As a proxy for watching, we use data recorded by the IMVH software that captures the video ID, student ID, and the date and time when a student opens a video link. We define the following measures:

1. *Videos*: Number of links opened, including duplicates
2. *Unique videos*: Number of unique video links opened
3. *Hours of videos*: Total runtime of video links opened
4. *Hours of unique videos*: Total runtime of video links opened with duplicates removed

For expositional ease, we use “watching” to refer to the link-opening behavior as defined above. Although video watching in our data is a binary measure, watching behavior can vary greatly in intensity. Some students take notes, pausing and rewatching portions of the video as needed. Other students, we suspect, play videos in the background without absorbing much material. Exploring the intensive margin of video watching remains an area for future research that will benefit from new technologies that can quantify video engagement including interactive content embedded in videos, eye-tracking devices, and more.

We helped students keep track of their progress towards 40 videos by periodically updating the online gradebook with counts determined from the IMVH data. Although nearly all

¹⁵11 of 164 *Incentive*, 23 of 167 *Control*, and 10 of 373 *Above median* students did not identify their grading schemes correctly. 146 students did not answer the quiz, several of whom had dropped the course following the first midterm.

students followed our instructions to watch videos completely and sequentially,¹⁶ in a few exceptional cases, students opened 40 or more video links within a matter of a few minutes. We manually adjusted their video counts in the gradebook and emailed them a reminder of the requirements for videos to count towards the grade incentive.¹⁷ Though it is doubtful these strategic students gained much from opening so many videos so quickly, to maintain interpretability of our results, we do *not* remove these clicks from our video count measures.¹⁸ Our *unique* video measures, however, are less sensitive to this behavior.

To ensure fairness, we informed students that final letter grades would *not* be affected by being in the experiment. We accomplished parity between *Control* and *Incentive* arms through curving final grades. First, we applied a curve to the *Control* and *Above median* arms as one group to achieve a grade distribution in line with that of previous cohorts. Second, we curved the *Incentive* students' course grades to match the average course grade among *Control* students after curving in the previous step. Since course grades are by construction *ex post* equal between the *Control* and *Incentive* arms, we use exam scores as our primary outcomes of interest. Our secondary outcomes of interest include term GPA and number of courses passed, which help us understand how treatment may have affected other courses. We examine separately econ and non-econ courses in case effects differ by course content. To better understand mechanisms, we estimate effects of treatment on take-up of other studying tools within Micro A. Finally, we explore video watching and grade outcomes in Micro B, the subsequent course in the intermediate microeconomics sequence, to see if treatment effects persist beyond one term.

4.4 Empirical strategies

In this paper, we estimate the effect of being assigned to the *Incentive* arm on our outcome variables of interest, Intent To Treat (ITT) effects, as well as the effect of watching videos for those induced by the incentive to watch more videos, a Local Average Treatment Effect

¹⁶The time between timestamps for video links opened in succession was almost always longer than the runtime of the video.

¹⁷Although additional email communication as a result of treatment could violate the exclusion restriction, the small number of affected students is unlikely to have much (if any) influence on our results.

¹⁸Our causal effects are per “links opened” rather than per “links opened subject to certain qualifying conditions”. The cost of this interpretability is likely downward bias on our causal effect estimates.

(LATE), also referred to as a Complier Average Treatment Effect (CATE) (Imbens & Rubin, 2015).

Below we outline the empirical strategies for estimating both the ITTs and LATEs.

4.4.1 Intention To Treat (ITT)

In this section, we examine the empirical strategy for estimating the causal effect of being assigned to the *Incentive* arm on outcomes of interest, such as exam scores. These are ITT estimates and not average treatment effect estimates because of two-sided non-compliance: some students in the treatment arm do not watch videos and some students in the control arm do watch videos. Since the incentive itself in our setting is representative of how future instructors may induce their students to watch videos, the ITT estimates are policy-relevant for instructors considering adopting the IMVH or other video-based learning methods in their courses.

Our baseline ITT specification is the partially linear model:

$$Y_i = \beta Z_i + f(X_i) + \epsilon_i \quad (1)$$

where Y_i is an outcome of interest (e.g. videos watched or test scores) for student i , $Z_i \in \{0, 1\}$ is a treatment indicator with those in the *Control* arm having $Z_i = 0$ and those in the *Incentive* arm having $Z_i = 1$, $f()$ is a generic function through which X_i , a vector of controls, affects Y_i , and ϵ_i is an unobserved residual. β , our parameter of interest, is the causal effect of being assigned to the *Incentive* arm on the outcome of interest Y , assumed to be constant across the population.¹⁹ Under unconfoundedness, $\hat{\beta}$ is an unbiased estimate of the ITT effect (Imbens & Rubin, 2015).²⁰

In our baseline estimation of Equation 1, we include in X_i year indicators and first midterm score, following the advice of Bruhn and McKenzie (2009) to control for all covari-

¹⁹Our experiment takes place over two years, and we pool the sample across both years. Out of the 850 student-years, one student repeated the course in both years, and hence there are 849 unique students. For simplicity, we drop the subscript t from our specifications, treating the one repeating student as independent across years. Dropping this student from the sample leaves the results virtually unchanged.

²⁰Though we cannot test whether Z_i is confounded by unobservable covariates, we have confidence unconfoundedness holds given the random assignment of Z_i and the balance across observable covariates as demonstrated in Table A1 and A2.

ates used in seeking balance. In a second model, we include additional controls chosen using the Post-Double-Selection (PDS) procedure of Belloni, Chernozhukov, and Hansen (2014b), explained in detail in Appendix A.2. In a third model, to check that our results are robust to potentially nonrandom attrition by treatment arm, we fit Equation 1 including pair fixed effects. These fixed effects subsume the year indicator (since pairs were assigned separately across years), so we drop the year indicator but keep midterm 1 score to control for small differences within pairs along that dimension. As entity fixed effects require at least two observations within the entity to be estimable, we drop any students whose matched pair attrited.

In our results, we present an additional nonparametric estimate using Neyman’s (1923) repeated sampling approach, considering each pair (block) an independent, completely randomized experiment and averaging the results. We estimate the point estimate of the ITT as the mean difference in outcomes across pairs:

$$\hat{\tau} = \frac{1}{J} \sum_{j=1}^J \hat{\tau}_j = \frac{1}{J} \sum_{j=1}^J y_{j,I}^{\text{obs}} - y_{j,C}^{\text{obs}} \quad (2)$$

where $\hat{\tau}$ is the point estimate of the ITT, J is the number of pairs in the sample, and $\hat{\tau}_j = y_{j,I}^{\text{obs}} - y_{j,C}^{\text{obs}}$ is the observed difference in outcome for pair j . The estimated standard error of $\hat{\tau}$ (Athey & Imbens, 2017; Imai, 2008; Imbens & Rubin, 2015) is:

$$\widehat{SE}(\hat{\tau}) = \left(\frac{1}{J} \sum_{j=1}^J \widehat{V}(\hat{\tau}_j) \right)^{\frac{1}{2}} \quad (3)$$

where $\widehat{V}(\hat{\tau}_j)$ 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 (Imbens and Rubin, 2015, Athey and Imbens, 2017):

$$\widehat{V}(\hat{\tau}_j) = s_{j,I}^2 + s_{j,C}^2 \quad (4)$$

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 in a matched-pair setting as there is 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 (Imai, 2008, Imbens and Rubin, 2015, Athey and Imbens, 2017):

$$\widehat{SE}(\widehat{\tau}_j) = \left(\frac{1}{J(J-1)} \sum_{j=1}^J (\widehat{\tau}_j - \widehat{\tau})^2 \right)^{\frac{1}{2}} \quad (5)$$

Similar to the fixed effect model, the Neyman repeated sampling approach is only estimable if we drop all students whose matched pair attrited. This drop in observations increases the width of our confidence intervals, albeit modestly since including only matched pairs reduces unexplained variance in the outcome variables of interest. We present estimates from all four models to demonstrate that the results are generally similar and not sensitive to model specification or choice of control variables.

4.4.2 Local Average Treatment Effect (LATE)

Here we present the empirical strategies for estimating the causal effect of watching videos on outcomes of interest, exam scores. The average causal effect of watching videos can be modeled as:

$$Y_i = \gamma v_i + g(X_i) + u_i \quad (6)$$

where γ is the average causal effect of watching an additional video, Y_i is an outcome of interest (e.g., exam scores) for a student i , $g()$ is a generic function through which X_i , a vector of pretreatment covariates, affects Y_i , and u_i in an unobserved model residual. Because a student's decision to watch videos is likely correlated with unobservable factors (for example, motivation) that are also correlated with outcomes, regressing Y_i on endogenous videos v_i will provide biased estimates of γ . To solve this problem, we rely on variation in v_i induced by an exogenous instrument Z_i . In our setting, Z_i is assignment to the *Incentive* grade scheme. If Z_i is a valid instrument for v_i , then we can estimate γ using two-stage least squares (2SLS):

$$v_i = \alpha Z_i + f(X_i) + e_i \quad (7)$$

$$Y_i = \gamma \hat{v}_i + g(X_i) + u_i \quad (8)$$

where $f()$ and $g()$ are generic functions through which X_i affects v_i and Y_i , respectively, e_i and u_i are unobserved model residuals, and \hat{v}_i is instrumented videos estimated by Equation 7. We assume the influence of Z_i on v_i is monotonically increasing, that is, $v_I = E(v_i|Z_i = 1) \geq E(v_i|Z_i = 0) = v_C$. Hence, γ is the per-video average treatment effect, local to students induced by the incentive to watch on average $v_I - v_C$ additional videos.

Under the assumptions of unconfoundedness, excludability, monotonicity, and non-interference, $\hat{\gamma}$ is an unbiased estimate of the LATE (J. D. Angrist & Imbens, 1995). Unconfoundedness requires that Z_i be independent of potential outcomes, a reasonable assumption given random assignment of students to the *Incentive* arm. Excludability assumes that outcomes (grades) are only affected by the instrument (incentive) 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 necessary because of two-sided noncompliance and requires that students assigned treatment watch weakly more videos than they would 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 unconfoundedness,²¹ excludability,²² and monotonicity²³ are reason-

²¹Although we randomly assigned treatment, one concern is nonrandom attrition. We find that the *Incentive* and *Control* arms remain balanced across observables by the end of the experiment. Additionally, we find our results are similar when restricting the sample to students whose matched pair did not attrite.

²²While 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. Of course, watching videos inherently requires time that takes away from some other activity. Hence, the results should be interpreted as the causal effects of more videos and less of whatever else they would have been doing. This subtle point could matter for external validity as a different population of students with zero leisure time may respond differently to the incentive.

²³Though 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, Figure 3 shows that the two

able assumptions, we have more concern about non-interference because of the potential for spillovers between students in the same class. If we had unlimited resources, a 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 detect reasonable 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.²⁴

Similar to our estimates of Equation 1, we estimate Equation 8 with three sets of controls: only year and first midterm score, controls chosen using PDS, and a fixed-effect model with controls chosen using PDS. We additionally estimate the LATE using Neyman’s repeated sampling approach whose estimators we derive in Appendix B.

4.4.3 Treatment Effect Heterogeneity

So far we have estimated average treatment effects across all students below the median score on the first midterm. In this section, we investigate the extent to which treatment effects varied along key demographic variables. In particular, we add an interaction term to Equation 1:

$$Y_i = \beta_1 Z_i + \beta_2 Z_i \mathbb{1}_{x_i=d} + f(X_i) + \epsilon_i \quad (9)$$

where $Z_i \mathbb{1}_{x_i=x}$ takes a value of 1 when treatment students have the value d for demographic variable $x \in X$. β_2 represents the difference in treatment effects for those with

CDFs do not cross.

²⁴Although spillovers are possible, we believe the magnitude of the spillovers are likely small given that students have for the most part not yet formed strong social networks. 47% of students in the *Incentive* or *Control* arms are transfer students in their first term at the university. The remaining students are predominantly sophomores taking their first upper division course. Social dynamics at the university facilitate networks within “colleges” more than majors for the very reason of encouraging academic diversity among peer groups. One example of a possible positive spillover is the online discussion board where students could ask questions about content covered in the IMVH.

demographic d relative to those without. In practice, we estimate Equation 9 including in X_i dummies for year and the demographic of interest as well as first midterm score.

While we observe many demographic variables, we are underpowered to detect reasonable effect sizes after adjusting for multiplicity, given the small sample sizes of our subgroups. As such, we focus on heterogeneity along our blocking variables and covariates we hypothesized *ex ante* may have treatment heterogeneity. We estimate heterogeneity by levels of videos watched pretreatment since it is plausible that those with greater experience watching videos may have greater treatment effects. On the other hand, those who watched many videos pretreatment may have watched many videos during the experiment even if not treated, and hence the incentive may have little effect. As mentioned in Section 4, nearly all transfer students in the experiment are taking their first term at a four-year university and may not yet have optimized their studying practices. The university achievement gap between underrepresented and majority groups is well documented, including lower GPAs, longer time to graduation, lower graduation rates on average (Bowen, Chingos, & McPherson, 2009). It is important to know how treatment may affect this gap, potentially through mechanisms including reduced disparities in guidance on studying methods and technologies that may help those for whom English is not a native language, such as captions and the ability to reply videos.²⁵

We present our results in Appendix Table A4, which displays the coefficient estimates of β_2 from Equation 9. We find no evidence of treatment effect heterogeneity across our blocking variables, year and first midterm score. We are hesitant to make strong conclusions given the width of our confidence intervals, but this finding is consistent with stable treatment effects across the distribution of student abilities and years of the experiment. Next, we find no heterogeneity across levels of videos watched pretreatment in three of four specifications. We find marginally significant positive effects on the second midterm for those with higher levels of pretreatment video watching. Moving to demographic variables, we find marginally significant negative heterogeneous effects for female students on the final exam, but no effect on the second midterm. We find significant negative effects for Asian students and positive

²⁵Unfortunately, we are unable to observe native language directly or better proxies such as country of home address or visa status.

effects for White students on the second midterm, but no effect on the final exam. However, after adjusting for multiplicity, either using a Bonferroni correction or the less conservative methods proposed by List, Shaikh, and Xu (2019), none of our heterogeneity results remain significant.

4.4.4 Treatment Effects at the Cutoff

Here we describe estimation of treatment effects at the first midterm score cutoff. Because the probability of being assigned to the *Incentive* arm changes discontinuously from 0.5 to 0 at the midterm score cutoff, our setting is appropriate for estimating local treatment effects using a regression discontinuity (RD) design (J. D. Angrist & Pischke, 2008; Imbens & Lemieux, 2008; Thistlethwaite & Campbell, 1960). With this method, we compare students in the *Incentive* arm who scored just below the cutoff to those in the *Above median* and *Control* arms who scored just above or below the cutoff, respectively. These two groups are similar across pretreatment characteristics but different in treatment status, thereby providing an estimate of the treatment effect local to those who scored near the cutoff.

Since RD designs require that agents near the cutoff be similar across covariates except treatment status, a threat to validity is manipulation of the forcing variable (in our study, midterm score), which biases treatment effect estimates by nonrandom selection into treatment. This manipulation can occur if agents behave strategically to target a particular side of the cutoff, for example, scoring slightly higher than a published minimum SAT score for college admission. Since students in our experiment do not know the cutoff *ex ante*, it is unlikely that students would attempt to target a particular side of the midterm score cutoff²⁶. Ultimately, we must *assume* continuity of the conditional means of the potential outcomes along the midterm score; however, we do not observe a discontinuity in any observable pretreatment covariate at the cutoff, which gives us further confidence that this assumption holds.

To estimate local ITT effects using a sharp RD, we return to the potential outcomes framework modeling the treatment effect $\tau(c)$ as the difference in expected outcomes at the

²⁶It would be surprising for students who value high grades to target the expected median score since any student capable of doing so would likely earn a higher grade in the course by scoring as high as possible on the midterm exam rather than strategically scoring just below the expected median cutoff.

cutoff c along the forcing variable x (midterm score):

$$\begin{aligned}\tau(c) &= \lim_{x \uparrow c} E[Y_i | X_i = x] - \lim_{x \downarrow c} E[Y_i | X_i = x] \\ &= E[Y_i(1) | X_i = c] - E[Y_i(0) | X_i = c]\end{aligned}\tag{10}$$

We estimate $\tau(c)$ using local low-order polynomials (local linear regression), per the advice of Gelman and Imbens (2019).

Sharp RD designs used in the literature frequently do not observe $Y(1)$ and $Y(0)$ for the same values of x . In our setting, however, we observe $Y(0)$ both above and below the cutoff. Hence, we need to assume continuity only for $Y(1)$ as we do not observe any outcomes for treated students scoring above the cutoff but *do* observe outcomes for control students both above and below the cutoff.

5 Results

In this section, we first examine attrition and establish that the *Incentive* and *Control* arms in our analysis sample are balanced on observable characteristics. Second, we show that the grade encouragement worked: students in the *Incentive* arm watched significantly more videos than did their *Control* peers. Third, we estimate the effects of being assigned to the *Incentive* arm (ITT) and the effects of videos (LATE) on grade outcomes. Fourth, to better understand mechanisms, we examine spillovers to other studying methods and grades in other courses taken during the experiment term. Finally, we see whether behavior change persists after exogenous incentives are removed by estimating treatment effects on video watching and grades in the subsequent microeconomics course, Micro B.

5.1 Attrition and balance

At the university where the experiment took place, Micro A is the first upper-division economics course, often taken by students who have not yet declared a major. As such, it has higher withdraw rates than most other economics courses, a product of both challenging

course material and updated priors on interest in the field. In Micro A one year before the experiment, 8.5% and 13.4% of students who took the first midterm did not take the second midterm and final exam, respectively. Unsurprisingly, the withdrawal rates were greater for students who scored below the median on the first midterm: 14.7% and 24.2% of these students did not take the second midterm and final exam, respectively.²⁷ In the present study, high attrition is not problematic, other than reducing statistical power, if attrition is independent of potential outcomes. If attrition is influenced by treatment status, which could occur, for example, if treatment students believed they were more likely to pass the course,²⁸ then the resulting nonrandom selection into our analysis could bias the results.

We assess the presence of nonrandom attrition by comparing rates of attrition by treatment arm and examining balance across observable demographic variables in the analysis sample. Students could have attrited in three ways. First, students under the age of 18 at the start of the experiment were removed from the analysis sample. Because of student privacy considerations, demographic variables (including age) were not observable until the conclusion of the experiment. Second, students who opted out of having their data included in the experiment analysis, an option available to students at any time during the experiment, were removed. The analysis sample was prepared and anonymized by a campus-based independent education research organization, per IRB requirement, which removed four minors and seven opt-outs from the sample and merged demographic variables before returning the anonymized data to the research team.²⁹

The final and largest cause of attrition is withdrawing from the course. At the present university, students may formally drop a course without penalty up to the fourth week of the quarter. Between the fourth and sixth weeks, students may withdraw from a course, but a “W” grade is assigned in lieu of a letter grade, which does not affect GPA. From the seventh week onwards, students may no longer formally withdraw, but some may choose not to take the second midterm or final exams, which almost assuredly results in a failing grade in the

²⁷The 2017 statistics are calculated from a sample that differs somewhat in inclusion criteria relative to the 2018 and 2019 samples. We provide these statistics to highlight the historically high rates of attrition and not to make comparisons with the experiment sample.

²⁸As noted earlier, we informed students that final grades would be curved separately between *Control* and *Incentive* arms to remove any advantage treatment may carry towards passing the course.

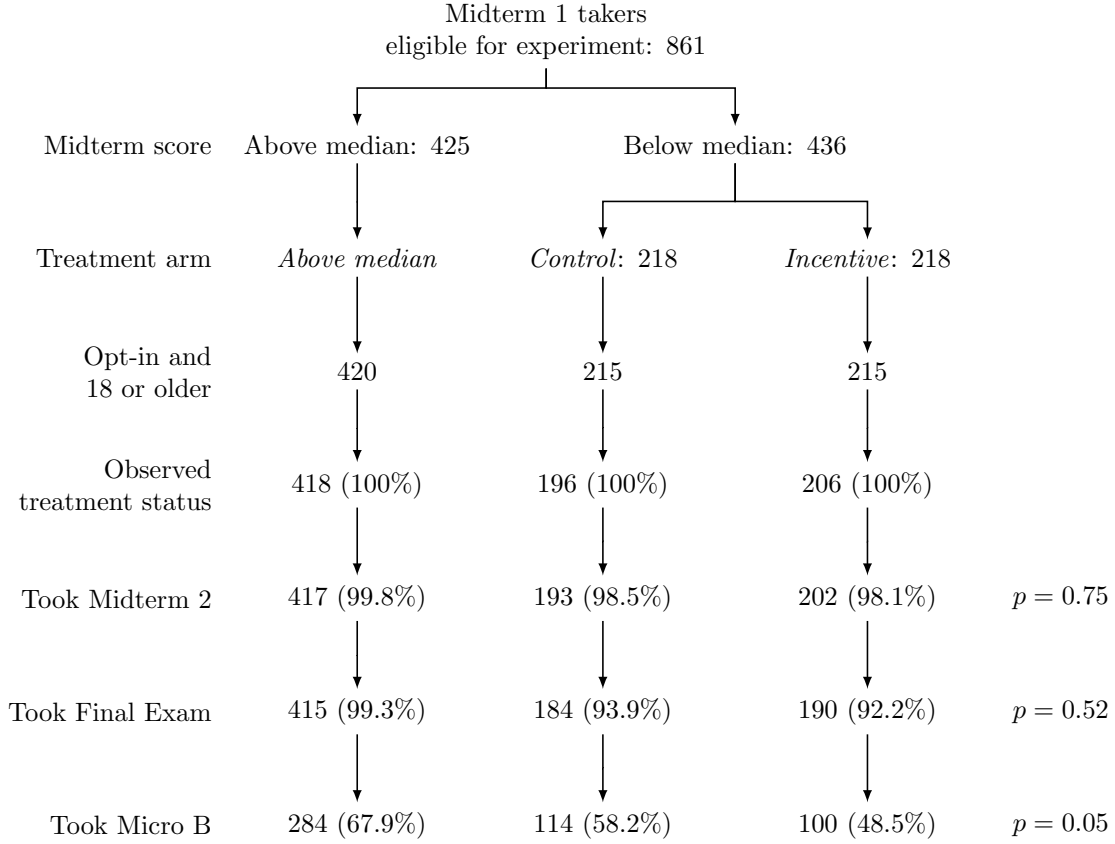
²⁹In total, five *Above median*, three *Control*, and three *Incentive* students were removed for age or opting-out.

course and does factor into the student’s GPA. Because the first midterm took place in week four, the same week as the penalty-free drop deadline, many students took the exam and withdrew before finding out their grades, which were posted in week five. However, because of the lag between when a student drops a course and when the instructor is notified, we assigned treatment to several students not knowing they had already dropped the course. In total, 30 students - two *Above median*, 19 *Control* and nine *Incentive* - dropped the course before finding out their first midterm scores and treatment assignments. Among students who waited to find out their first midterm scores and treatment assignments, three *Control* and three *Incentive* students did not take the second midterm, and an additional nine *Control* and 13 *Incentive* students did not take the final exam.

We show attrition at each stage of the experiment in Figure 2. The p-values in the figure are calculated using two-sample t-tests of the equality of attrition rates between the *Control* and *Incentive* arms at each stage. We do not find any statistically significant difference in attrition before the second midterm and final exams that can be attributed to treatment status. However, we do find a significant difference in enrollment rates in Micro B: treatment students were 9.7 percentage points *less* likely to enroll in Micro B than were control students. Conditioning on midterm 1 score and year reduces the gap slightly to 8.9 percentage points ($p = 0.06$). As mentioned in Section 4.4, as a robustness check against potential bias from any nonrandom attrition, we estimate treatment effects using models that only include matched-pairs.

Since all students below the median on the first midterm had equal probabilities of being assigned to the *Control* and *Incentive* arms, treatment arms are balanced on covariates in expectation. In practice, due to chance and nonrandom attrition, treatment arms can be unbalanced on covariates, which can bias estimates if not addressed in the analysis, particularly in small samples (Athey & Imbens, 2017). We check balance on observable characteristics after attrition for both the second midterm and final exam samples. As can be seen in Appendix Tables A1 and A2, we find no statistically significant difference between the *Control* and *Incentive* arms in observable covariates including first midterm score, year, previous term’s cumulative GPA, videos watched before the first midterm, ethnicity, gender, and transfer status. However, as discussed in Section 4.4, to correct for potential imbalance

Figure 2: Attrition by treatment arm



Percentages in parentheses are the portion of students who observe their treatment status and took the second midterm, final exam, and enrolled in Micro B, calculated separately for each treatment arm. The p-values are from a two-sample t-test of the equality of attrition rates between the *Control* and *Incentive* arms at each stage.

and to improve precision, we estimate models that include controls chosen via the post-double-selection method of Belloni, Chernozhukov, and Hansen (2014b).

5.2 Relevancy of the encouragement instrument

We use a Two-Stage Least Squares approach to estimate the LATE of watching videos on exam performance, as detailed in the Section 4.4. We must check that our instrument is both valid and relevant to ensure this method will produce an unbiased estimate of the LATE (Imbens & Rubin, 2015). The validity condition is met by assigning treatment at random conditional on midterm exam score and year of instruction. Balance across pretreatment observables, as demonstrated in Appendix Tables A1 and A2, give us further confidence

that treatment status is uncorrelated with demographics.

Next we check instrument relevancy, that is, whether treatment status generates significantly more video watching. In Table 3 we present estimates from Equation 1. We find that by the second midterm exam, being assigned to the *Incentive* arm induces students to watch 9.1 - 10.5 videos and 6.0 - 6.8 unique videos more than being assigned to the *Control* arm. The gap between treatment and control grows by the final exam to 38.4 - 39.2 videos and 20.5 - 21.6 unique videos. The larger gap by the final is unsurprising given that the deadline to earn the grade incentive was the day before the final exam. Following the recommendations of Andrews, Stock, and Sun (2019), we assess the strength of our instrument using the effective F-statistic of Montiel Olea and Pflueger (2013) which, in our just-identified setting, coincides with the Wald statistic of Kleibergen and Paap (2006). The effective F-statistic for the second midterm and final exam first-stage specifications are 18.6 and 194.6, respectively, both of which are greater than the Stock and Yogo (2005) critical value of 16.4 and the rule-of-thumb cutoff of 10.

Graphically, we depict the distribution of videos watched as a function of treatment in Figure 3. Notably, the gap between treatment and control distributions remains significantly positive at every level of video watching by the final exam. The difference is most pronounced near the required number of videos to earn the grade incentive, after which the difference diminishes towards zero. For the second midterm sample, the difference is smaller but significantly different from zero between zero and 62 videos watched. Collectively, given the highly significant first stage regression results, large first-stage F-statistics, and monotonic increase in video watching across the sample, we have high confidence that our instrument meets the relevancy criterion.

5.3 Effects on exam scores

In this section, we estimate the causal effects of being assigned to the *Incentive* arm on exam scores (ITT). This estimate is relevant for educators interested in predicting how requiring videos will change exam scores in their classes using the same grade-based incentive implemented in our experiment. Additionally, we estimate the causal effect of watching videos on exam scores, which is of interest to educators deciding which teaching technologies

to provide for their classes as well as to students choosing among different studying tools.

For both the ITTs and LATEs, we examine effects on the second midterm and final exams using both parametric methods (i.e. Equations 1 and 8) and nonparametric methods at the repeated sampling framework of Neyman (1923). We check that our parametric results are robust to model specification by estimating Equations 1 and 8 with and without $f(X_i)$ as a vector of linear control variables chosen via PDS (Belloni, Chernozhukov, & Hansen, 2014b). To rule out nonrandom attrition across treatment arms as a confounder, we fit a fixed effect model that drops any student whose matched pair attrited. These specifications and identification strategies are described in detail in Section 4.4.

Table 4 presents estimates of the effects of treatment on second-midterm and final exam scores. Across our four specifications, we estimate reduced-form (RF) impacts of being assigned to the *Incentive* arm of 0.17 - 0.18 standard deviations on the second midterm. These estimates, along with our first-stage estimates (Table 3), imply LATEs of 0.26 - 0.30 standard deviations per 10 unique videos, or 0.16 - 0.18 standard deviations per hour of unique content. For the final exam, we estimate similar ITT effects: being assigned to the *Incentive* arm raises scores by 0.14 - 0.18 standard deviations. However, given the larger first stage effects for the final exam, we estimate smaller LATEs: 0.08 - 0.09 standard deviations per 10 unique videos, or 0.04 - 0.05 standard deviations per hour of unique content.

Given the large F-statistics when estimating first-stage effects of our incentive instrument on videos watched, we are not particularly concerned about bias from weak instruments. However, following the advice of Andrews, Stock, and Sun (2019), we report Anderson-Rubin confidence sets, which are efficient regardless of the strength of our instrument. As can be seen in Appendix Table A3, we find that the weak-instrument-robust confidence intervals are very similar to those presented in Table 4.

5.4 Effects at the incentive cutoff

Here we present estimates of treatment effects for students at the incentive-eligibility cutoff: the median score on the first midterm. While the primary intended purpose of cutoff was to avoid burdening high-performing students, the discontinuous change in treatment probability allows for identification of LATEs, as we describe in section 4.4. Unfortunately, the density

of students at the cutoff is not enough to allow for precise identification of treatment effects. For completeness and bounding possible effect sizes, we share our findings.

Before using a regression discontinuity (RD) design in our setting, we must establish that students did not manipulate their midterm scores to affect their chances of being treated. Intuitively, it would be surprising to find manipulation as students did not know the median score *ex ante*. Additionally, it would be quite difficult to accurately predict one’s midterm score while taking the exam, and those capable of doing so would likely have higher returns scoring higher than the median. Appendix Figure A4 shows that, as expected, the distribution of midterm scores does not have any unordinary masses on one side of the cutoff.

In most RD designs, the researcher can only observe either treatment or control outcomes on one side of the cutoff. In our setting, however, we can observe control outcomes on *both* sides of the cutoff. As such, we identify treatment effects adjusting the conventional RD approach to account for our increased knowledge of the control group’s mean outcome at the cutoff. In practice, this entails differencing treatment and control conditional mean outcomes using local linear regressions where we calculate the control conditional mean using both *Control* and *Above median* observations.

Figures A5 and A6 in the Appendix depict the effects of treatment near the cutoff for the second midterm and final exams, respectively. Both Figures show a marked increase in video watching for treatment students at the cutoff. Both figures also show positive point estimates of the effect of treatment on exam scores for those at the cutoff, but neither of these point estimates are statistically distinguishable from zero. Appendix Table A5 summarizes the point estimates and confidence intervals for those at the incentive eligibility cutoff. We find that treatment increases video watching by 9.1 and 24.2 unique videos by the second midterm and final exams, respectively, and we find positive but not statistically significant effects on exam scores. The first-stage effects at the cutoff are in line with the estimated average first-stage effects for those below the median.

5.5 Spillovers during concurrent term

We next estimate spillover effects to other courses taken concurrently during the term of the experiment. Although we find positive effects on outcomes in Micro A, it is important to

examine spillover effects to understand the complete picture of how treatment affects student achievement. If the gains in Micro A come at the cost of slowed progress in other courses, then treatment may be counterproductive for some students. On the other hand, since the videos cover concepts and methods that could be covered in other courses, treatment may help students outside Micro A.

Table 5 presents our estimates of Equation 1 where Y_i is GPA, number of classes passed, or portion of classes taken for a letter grade. We estimate the effects of treatment on term GPA calculated separately for all classes, all classes excluding Micro A, all classes outside of economics, and all economics classes excluding Micro A. In general, we find marginally significant or insignificant but directionally positive spillover effects on GPA.³⁰ We can rule out large negative spillover effects: in our worst-case specification for term GPA, our 95% confidence interval rules out negative spillover effects larger than -0.02 (on a 4.0 scale), or less than 1% of the mean term GPA among control students. There is no statistically significant difference between our estimates of spillover effects on GPA when restricting to only economics or non-economics courses.

We additionally estimate the effects of treatment on number of classes passed and find small, insignificant, but directionally positive effects. We find that treatment caused students to pass 0.02 - 0.09 more classes, or about 1% - 3% more classes than the control mean. We find no effect of treatment on number of classes not passed or withdrawn. Interestingly, treatment students were somewhat less likely to take Micro A for a letter grade than were control students, but this difference is only marginally significant for one of the four specifications and insignificant for the rest. We find no relationship between treatment and fraction of classes taken for a letter grade versus Pass/No Pass. Across all grade spillovers examined, we find mostly small, directionally positive effects, which gives us confidence that treatment is likely not harmful to academic success outside of Micro A.

Besides estimating spillover effects on other course grades, we also examine spillovers to other studying methods within Micro A. Doing so helps us better understand mechanisms: do students substitute away from other studying when encouraged to watch more videos, or

³⁰At this university, term GPA is affected only by classes taken for a letter grade. Hence, students may not have attrited from the sample but may have taken all courses Pass/No Pass and thus have no term GPA. As such, we report sample sizes for each GPA specification.

are they more likely to complement their video-watching with other unincentivized studying? In Table 6, we display the results of estimating equation 1 where Y_i is an alternative form of studying. We find that *Incentive* students are directionally less likely to attend class, though not statistically significantly so. On the other hand, treatment students interacted with the online discussion board more than did control students, but again these estimates are not statistically distinguishable from zero. We do not find any significant relationship between treatment and tutoring attendance.

5.6 Spillovers to subsequent term

While we offered treatment students a grade incentive to watch videos during Micro A, students were not offered a grade incentive during the subsequent course in the intermediate microeconomics sequence, Micro B. However, all students in Micro B maintained access to the IMVH and were verbally encouraged to watch videos as a study method. Fortunately, we are able to observe video watching and grade outcomes in Micro B.

We present our estimates of spillover effects during the subsequent term in Table 7. In Panel A, we estimate the effect of treatment on videos watched during the subsequent term, for those who took Micro B. We find large and statistically significant effects: treatment caused students to watch 8.1 - 9.9 more unique videos and 1.2 - 1.5 hours of unique content compared to control students. In Panel B, we estimate equation 1 where Y_i is the first midterm, second midterm, or final exam score. Unfortunately given the small subsample of students who took Micro B, we are underpowered to detect effect sizes consistent with those observed during Micro A. Finally, in Panel C, we estimate the effect of treatment on taking Micro B and the number of classes passed and withdrawn. We find no effects statistically distinguishable from zero, though, as mentioned in section 5.1, treatment students were directionally less likely to take Micro B than were control students.

6 Discussion

Here we discuss the findings reported in Section 5.

Contributions: First, we find that a small grade incentive is effective in motivating poorly

performing college students to take-up video watching. We also reduced the weight on an early assessment and allowed students to earn back the lost points fully by meeting the video watching requirement, which may also be an important motivator. This result adds to the literature on what motivates college students to use educational inputs including financial incentives (see papers cited in Gneezy, Meier, and Rey-Biel (2011)), requiring students to attend class if they perform poorly on an early assessment but with no penalty for not attending (as in Dobkin, Gil, and Marion (2010), or by having students set goals on the use of an practice quizzes and, again, no penalty, for failing to achieve the goal (Clark et al., 2020)).³¹ Grade incentives have the unappealing feature that grades are directly a function of input use. The grade incentive used in this study was small, at most four percent of the student’s grade, which may help mitigate this concern.

Second, we find that inducing students who performed poorly on an early assessment to increase the amount of time they spend watching instructional videos increased their exam performance. Since there was no drop in either grades for other courses taken in the same quarter or a drop in the use of many other educational for the class, this suggests that total study time for the course increased for treated students. Unfortunately, we were not able to determine whether the added study time for the course did come from, which is important for student welfare calculations. Our results are consistent with other experimental and quasi-experimental studies that find positive effects of educational inputs on college student performance.

Third, while not statistically significant, it seems worth pointing out that we found positive point estimates on treated student’s use of other educational inputs (attending class, downloading material from the course web page, and attending a tutoring lab) and grades in other courses taken in the same quarter compared to the control group. This surprising result was also found by Dobkin, Gil, and Marion (2010) who required poorly performing students to attend class. The authors posit that the statistically insignificant

³¹Note, in contrast to Clark et al. (2020), Oreopoulos et al. (2019) find that setting a weekly schedule ahead of time and weekly reminders via a text message had only a small effect on study time and no effect on output as measured by grades, retention or credit accumulation. In contrast to the effectiveness of incentives on college student use of educational inputs, several papers find incentives to directly increase output, such as grades or performance on exams, are less effective. See for example, Fryer Jr (2011) shows that paying students in K-12 to increase inputs shows greater effects than paying for increased output, Gneezy, Meier, and Rey-Biel (2011) also point out that incentives in education appear to work better for inputs than output

but positive spillovers they find may be due to fixed effects of coming to campus. A fixed effects argument is less compelling in our context.

Finally, we find support for a "poor information" model of learning because treated students frequently watched more than the 40 videos they were required to watch to get the grade incentive and because treated students continued watching videos at a significantly higher rate in the following class when videos were not incentivized. J. Angrist, Lang, and Oreopoulos (2009) also find continued higher use of academic support services after the incentivized year for women. A poor information model of learning can also account for why incentivizing educational inputs has been found to be more effective than incentivizing grades or exam performance directly (see studies cited in Gneezy, Meier, and Rey-Biel (2011))

Allgood, Walstad, and Siegfried (2015) review the literature on using technology to provide supplemental aids for students in traditional classrooms and conclude that there is little causal evidence that student achievement is improved. The current study stands in contrast to this prior literature. In particular, they conclude that online instruction appears to have a negative effect on course grades and persistence. A more recent paper by Bettinger et al. (2017) confirms these results, particularly for low performing students. The IMVH is clearly the backbone for an online class as it includes all the videos one would require students to watch as part of an online course. Perhaps the main effect of the video-watching requirement in this study was to increase student study time. Would encouraging at-risk students to watch instructional videos in an online class be as effective as we find for a in-person class? We consider this an important area for future research.

6.1 Limitations

The present study has several limitations that should be considered before, for example, creating one's own video handbook and requiring students to use it. First, the population studied is students who score below the median on the first midterm of an intermediate microeconomics course at a large, highly-selective public research university. The extent to which treatment effects vary by course, instructor, university, or along the midterm score distribution is beyond the scope of this paper. Additionally, the causal effects of watching videos that we estimate are *local* to compliers, i.e. students induced by the grade incentive to

watch additional videos. We cannot recover the *population* average treatment effect, though anecdotal evidence and economic theory both suggest that the population average treatment effect is likely greater than the LATE.

Though we followed an experimental design similar to that of Dobkin, Gil, and Marion (2010), some researchers may ponder why we elected to include only the bottom half of the first midterm distribution in the experiment rather than the entire class. While including the entire class would have increased statistical power, we believe the additional precision might have come at the expense of welfare losses by high performing students. The first midterm provides a signal of which students likely know for themselves how to study, both methods and duration. Coercing these high-type students to spend time with a potentially different studying method is less likely to be helpful and runs a higher risk of harming utility. On the other hand, students who, through a low midterm score, make manifest a need for alternative studying practices stand to benefit the most from instructor-provided guidance.

Another consideration is the time frame during which the experiment took place, 2018 to 2019. About three months after the conclusion of our experiment, most students in the United States and all students at the studied university began remote learning as the coronavirus pandemic prompted stay-at-home orders. With increased experience learning via electronic media, it is possible that treatment effects will be higher in the future than we estimate in our paper. On the other hand, if students find online learning materials increasingly *less* engaging, we may find the opposite.

In addition to estimating the effects of video handbooks in other educational settings, future research should examine treatment effects in the presence of weekly deadlines instead of one final deadline at the end of the term. Given our observation of greater “binge watching” by treatment students and the smaller effect sizes before the final exam compared to the second midterm exam, we suspect weekly deadlines may reduce the deleterious effects of procrastination. Despite the rich literature on the advantages of spread-out studying (Cepeda et al., 2006; Kornell, 2009), we note that “binge watching” was not unique among treatment students. Indeed, most students within each treatment arm watched more videos the last week of the term than any other week.

7 Conclusion

We examine the effectiveness of an innovative educational technology, a video handbook composed of 220 brief instructional videos on intermediate microeconomic theory. We used random assignment of a grade-based incentive to experimentally vary takeup of the video handbook, and we found that greater takeup caused student to score significantly higher on exams. Specifically, we estimate that treatment causes students to score about 0.18 standard deviations higher on midterm and final exams. For students on the margin of watching videos, watching an additional hour of unique content causes students to score between 0.05 to 0.15 standard deviations higher on exams.

Instructors may have concerns about making a resource such as the IMVH available if they believe students may substitute away from lectures or other more productive studying methods Kay (2012). Another concern is that forcing students to spend more time studying in one's class may worsen performance in other classes. Our analysis provides some confidence that neither of these fears are first-order concerns. We do not find evidence that students decrease their consumption of other forms of studying, nor do we find that students perform worse in other courses during the same quarter. Our point estimates of the effect of treatment on takeup of other studying methods, though not statistically significantly different from zero, are *positive* for most alternatives, suggesting that if any, students consider the videos complements to other forms of studying. A potential mechanism might be that the videos help students realize what they *don't* know, increasing the marginal benefit of subsequent studying.

A final concern is one of welfare. In a neoclassical model, instructors cannot make their students better off by forcing on them quantities of studying they would not otherwise have chosen for themselves. In a behavioral model, which we think is more appropriate in our university classroom setting, instructors *can* improve student welfare through intervention when myopia or information barriers lead to suboptimal time allocation decisions. We observe two phenomena that supports the information barriers model. First, treated students do not bunch at the cutoff for the grade incentive. Second, video consumption remains much higher among treated students in the term following conclusion of the experiment.

While there are many educational interventions that instructors could offer their students, the research on causal effects of educational interventions remains limited. Our study serves as an example of a feasible research design that runs a lower risk of generating welfare losses for high performing students than does a class-wide experiment. It is our hope, as educators ourselves, that more research be conducted on the effectiveness of pedagogical technologies.

References

- Allgood, S., Walstad, W. B., & Siegfried, J. J. (2015). Research on teaching economics to undergraduates. *Journal of Economic Literature*, 53(2), 285–325.
- Andrews, I., Stock, J. H., & Sun, L. (2019). Weak instruments in instrumental variables regression: Theory and practice. *Annual Review of Economics*, 11, 727–753.
- Angrist, J., Lang, D., & Oreopoulos, P. (2009). Incentives and services for college achievement: Evidence from a randomized trial. *American Economic Journal: Applied Economics*, 1(1), 136–63.
- Angrist, J. D., & Imbens, G. W. (1995). Two-stage least squares estimation of average causal effects in models with variable treatment intensity. *Journal of the American statistical Association*, 90(430), 431–442.
- Angrist, J. D., & Pischke, J.-S. (2008). *Mostly harmless econometrics: An empiricist's companion*. Princeton university press.
- Arulampalam, W., Naylor, R. A., & Smith, J. (2012). Am i missing something? the effects of absence from class on student performance. *Economics of Education Review*, 31(4), 363–375.
- Athey, S., & Imbens, G. W. (2017). The econometrics of randomized experiments. *Handbook of economic field experiments* (pp. 73–140). Elsevier.
- Beattie, G., Laliberté, J.-W. P., Michaud-Leclerc, C., & Oreopoulos, P. (2019). What sets college thrivers and divers apart? a contrast in study habits, attitudes, and mental health. *Economics letters*, 178, 50–53.

- Belloni, A., Chernozhukov, V., & Hansen, C. (2014a). High-dimensional methods and inference on structural and treatment effects. *Journal of Economic Perspectives*, 28(2), 29–50.
- Belloni, A., Chernozhukov, V., & Hansen, C. (2014b). Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies*, 81(2), 608–650.
- Bettinger, E. P., Fox, L., Loeb, S., & Taylor, E. S. (2017). Virtual classrooms: How online college courses affect student success. *American Economic Review*, 107(9), 2855–75.
- Bonesrønning, H., & Opstad, L. (2012). How much is students’ college performance affected by quantity of study? *International Review of Economics Education*, 11(2), 46–63.
- Bonesrønning, H., & Opstad, L. (2015). Can student effort be manipulated? does it matter? *Applied Economics*, 47(15), 1511–1524.
- Bowen, W. G., Chingos, M. M., & McPherson, M. S. (2009). *Crossing the finish line: Completing college at america’s public universities*. Princeton University Press.
- Bratti, M., & Staffolani, S. (2013). Student time allocation and educational production functions. *Annals of Economics and Statistics/ANNALES D’ÉCONOMIE ET DE STATISTIQUE*, 103–140.
- Bruhn, M., & McKenzie, D. (2009). In pursuit of balance: Randomization in practice in development field experiments. *American economic journal: applied economics*, 1(4), 200–232.
- Cepeda, N. J., Pashler, H., Vul, E., Wixted, J. T., & Rohrer, D. (2006). Distributed practice in verbal recall tasks: A review and quantitative synthesis. *Psychological bulletin*, 132(3), 354.
- Chen, J., & Lin, T.-F. (2008). Class attendance and exam performance: A randomized experiment. *The Journal of Economic Education*, 39(3), 213–227.
- Chen, P., Chavez, O., Ong, D. C., & Gunderson, B. (2017). Strategic resource use for learning: A self-administered intervention that guides self-reflection on effective resource use enhances academic performance. *Psychological Science*, 28(6), 774–785.

- Clark, D., Gill, D., Prowse, V., & Rush, M. (2020). Using goals to motivate college students: Theory and evidence from field experiments. *The Review of Economics and Statistics*, 102(4), 648–663.
- Dobkin, C., Gil, R., & Marion, J. (2010). Skipping class in college and exam performance: Evidence from a regression discontinuity classroom experiment. *Economics of Education Review*, 29(4), 566–575.
- Dunlosky, J., Rawson, K. A., Marsh, E. J., Nathan, M. J., & Willingham, D. T. (2013). Improving students’ learning with effective learning techniques: Promising directions from cognitive and educational psychology. *Psychological Science in the Public Interest*, 14(1), 4–58.
- Ferrari, J. R. (1992). Psychometric validation of two procrastination inventories for adults: Arousal and avoidance measures. *Journal of Psychopathology and Behavioral Assessment*, 14(2), 97–110.
- Fryer Jr, R. G. (2011). Financial incentives and student achievement: Evidence from randomized trials. *The Quarterly Journal of Economics*, 126(4), 1755–1798.
- Gelman, A., & Imbens, G. W. (2019). Why high-order polynomials should not be used in regression discontinuity designs. *Journal of Business & Economic Statistics*, 37(3), 447–456.
- Gneezy, U., Meier, S., & Rey-Biel, P. (2011). When and why incentives (don’t) work to modify behavior. *Journal of economic perspectives*, 25(4), 191–210.
- Grodner, A., & Rupp, N. (2013). The role of homework in student learning outcomes: Evidence from a field experiment. *The Journal of Economic Education*, 44(2), 93–109.
- Imai, K. (2008). Variance identification and efficiency analysis in randomized experiments under the matched-pair design. *Statistics in medicine*, 27(24), 4857–4873.
- Imbens, G. W., & Lemieux, T. (2008). Regression discontinuity designs: A guide to practice. *Journal of econometrics*, 142(2), 615–635.
- Imbens, G. W., & Rubin, D. B. (2015). *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press.
- Joyce, T., Crockett, S., Jaeger, D. A., Altindag, O., & O’Connell, S. D. (2015). Does classroom time matter? *Economics of Education Review*, 46(100), 64–77.

- Kapoor, S., Oosterveen, M., & Webbink, D. (2020). The price of forced attendance. *Journal of Applied Econometrics*.
- Kay, R. H. (2012). Review: Exploring the use of video podcasts in education: A comprehensive review of the literature. *28*(3).
- Kirby, A., & McElroy, B. (2003). The effect of attendance on grade for first year economics students in university college cork. *The Economic and Social Review*, *34*, 311–326.
- Kleibergen, F., & Paap, R. (2006). Generalized reduced rank tests using the singular value decomposition. *Journal of econometrics*, *133*(1), 97–126.
- Kornell, N. (2009). Optimising learning using flashcards: Spacing is more effective than cramming. *Applied Cognitive Psychology: The Official Journal of the Society for Applied Research in Memory and Cognition*, *23*(9), 1297–1317.
- Krohn, G. A., & O'Connor, C. M. (2005). Student effort and performance over the semester. *The Journal of Economic Education*, *36*(1), 3–28.
- Lavecchia, A. M., Liu, H., & Oreopoulos, P. (2016). Behavioral economics of education: Progress and possibilities. *Handbook of the economics of education* (pp. 1–74). Elsevier.
- List, J. A., Shaikh, A. M., & Xu, Y. (2019). Multiple hypothesis testing in experimental economics. *Experimental Economics*, *22*(4), 773–793.
- McCabe, J. (2011). Metacognitive awareness of learning strategies in undergraduates. *Memory & Cognition*, *3*(39), 462–476.
- Montiel Olea, J. L., & Pflueger, C. (2013). A robust test for weak instruments. *Journal of Business & Economic Statistics*, *31*(3).
- Morris, N. P., Swinnerton, B., & Coop, T. (2019). Lecture recordings to support learning: A contested space between students and teachers. *Computers & Education*, *140*, 103604.
- Munley, V. G., Garvey, E., & McConnell, M. J. (2010). The effectiveness of peer tutoring on student achievement at the university level. *American Economic Review*, *100*(2), 277–82.
- Neyman, J. (1923). On the application of probability theory to agricultural experiments: Essay on principles, section 9. *Statistical Science*, *5*(4), 465–472.

- Oettinger, G. S. (2002). The effect of nonlinear incentives on performance: Evidence from “econ 101”. *The Review of Economics and Statistics*, 84(3), 509–517.
- Oreopoulos, P., Patterson, R. W., Petronijevic, U., & Pope, N. G. (2019). Low-touch attempts to improve time management among traditional and online college students. *Journal of Human Resources*, 0919–10426R1.
- Pashler, H., Rohrer, D., Cepeda, N. J., & Carpenter, S. K. (2007). Enhancing learning and retarding forgetting: Choices and consequences. *Psychonomic Bulletin and Review*, 14(2), 187–193.
- Ralph, S., & R, S. T. (2008). The causal effect of studying on academic performance. *The B.E. Journal of Economic Analysis & Policy*, 8(1), 1–55.
- Setren, E., Greenberg, K., Moore, O., & Yankovich, M. (2021). Effects of flipped classroom instruction: Evidence from a randomized trial. *Education Finance and Policy*, Forthcoming.
- Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychological science*, 22(11), 1359–1366.
- Stock, J. H., & Yogo, M. (2005). Testing for weak instruments in linear iv regression. *Identification and inference for econometric models: Essays in honor of Thomas Rothenberg*, 80(4.2), 1.
- Tang, L., Li, S., Auden, E., & Dhuey, E. (2020). Who benefits from regular class participation? *The Journal of Economic Education*, 51(3-4), 243–256.
- Thistlethwaite, D. L., & Campbell, D. T. (1960). Regression-discontinuity analysis: An alternative to the ex post facto experiment. *Journal of Educational psychology*, 51(6), 309.
- Trost, S., & Salehi-Isfahani, D. (2012). The effect of homework on exam performance: Experimental results from principles of economics. *Southern Economic Journal*, 79(1), 224–242.
- Wozny, N., Balser, C., & Ives, D. (2018). Evaluating the flipped classroom: A randomized controlled trial. *The Journal of Economic Education*, 49(2), 115–129.

Table 3: Effects of Grade Incentive on Video Watching

	Control Mean	(1)	(2)	(3)	(4)
Panel A: By Midterm 2					
Videos	33.91	10.19*** (2.85)	10.54*** (3.12)	9.08*** (2.03)	9.58*** (2.19)
Unique videos	23.13	6.63*** (1.54)	6.79*** (1.70)	5.97*** (0.98)	6.11*** (1.11)
Hours of videos	5.88	1.68*** (0.50)	1.72*** (0.55)	1.48*** (0.35)	1.55*** (0.38)
Hours of unique videos	3.85	1.10*** (0.25)	1.13*** (0.28)	0.99*** (0.16)	1.02*** (0.18)
Observations		395	362	395	362
Panel B: By Final Exam					
Videos	53.09	39.25*** (4.06)	39.07*** (4.37)	38.57*** (3.40)	37.99*** (3.69)
Unique videos	33.95	21.55*** (1.55)	21.08*** (1.66)	21.28*** (1.22)	20.49*** (1.27)
Hours of videos	8.93	6.30*** (0.69)	6.26*** (0.75)	6.18*** (0.57)	6.05*** (0.62)
Hours of unique videos	5.54	3.43*** (0.25)	3.36*** (0.27)	3.38*** (0.20)	3.26*** (0.21)
Observations		374	332	374	332
Treatment assignment controls		Yes	No	Yes	Yes
Demographic controls		No	No	Yes	Yes
Pair Fixed Effects		No	No	No	Yes

Note: Model (1) contains linear controls for midterm 1 score and year; (2) is the difference in means and standard errors calculated using the repeated sampling framework of Neyman (1923); (3) and (4) use the post-double-selection (PDS) procedure of Belloni, Chernozhukov, and Hansen (2014b) to select control variables then estimate treatment effects and standard errors. The control variables selected using PDS are listed in Table A7. Models (2) and (4) include only students whose matched-pair did not attrite from the experiment. *Control Mean* is the mean for the Control students included in models (1) and (3), which is nearly identical to the mean for the Control students included in models (2) and (4). Standard errors in parentheses are robust to heteroskedasticity. ***, **, and * indicate significance at the 1, 5, and 10 percent critical levels, respectively.

Table 4: Effects of Videos on Grades

	(1)	(2)	(3)	(4)
Panel A: Midterm 2 score				
RF: Incentive	0.176* (0.090)	0.183* (0.094)	0.176* (0.090)	0.174* (0.096)
2SLS: 10 videos	0.266* (0.146)	0.270* (0.150)	0.300** (0.151)	0.293** (0.145)
2SLS: 1 hour of videos	0.160* (0.087)	0.163* (0.090)	0.181** (0.090)	0.170* (0.095)
Observations	395	362	395	362
Panel B: Final exam score				
RF: Incentive	0.175** (0.089)	0.174* (0.103)	0.175** (0.088)	0.138 (0.103)
2SLS: 10 videos	0.081** (0.041)	0.082* (0.049)	0.082** (0.041)	0.087* (0.046)
2SLS: 1 hour of videos	0.051** (0.026)	0.052* (0.031)	0.052** (0.026)	0.043 (0.031)
Observations	374	332	374	332
Treatment assignment controls	Yes	No	Yes	Yes
Demographic controls	No	No	Yes	Yes
Pair Fixed Effects	No	No	No	Yes

Note: This table reports coefficients on $Incentive_i$ from Equation 1 (Reduced Form, RF) and $Video_i$ from Equation 8 (Two-Stage Least Squares, $2SLS$). Test scores are measured in standard deviation units. Model (1) contains linear controls for midterm 1 score and year; (2) is the difference in means and standard errors calculated using the repeated sampling framework of Neyman (1923); (3) and (4) use the post-double-selection (PDS) procedure of Belloni, Chernozhukov, and Hansen (2014b) to select control variables then estimate treatment effects and standard errors. The control variables selected using PDS are listed in Table A7. Models (2) and (4) include only students whose matched-pair did not attrite from the experiment. Standard errors in parentheses are robust to heteroskedasticity. ***, **, and * indicate significance at the 1, 5, and 10 percent critical levels, respectively.

Table 5: Spillover Effects of Incentive on Other Course Grades

	Control Mean	(1)	(2)	(3)	(4)
Panel A: Effects on Term GPA					
All classes	2.59	0.13** (0.06) 373	0.13* (0.07) 332	0.11* (0.06) 373	0.10 (0.06) 332
Excluding Micro A	2.75	0.10 (0.07) 370	0.11 (0.08) 329	0.09 (0.07) 370	0.10 (0.08) 329
Excluding econ classes	2.99	0.06 (0.10) 315	0.09 (0.09) 278	0.06 (0.09) 315	0.08 (0.12) 278
Econ classes ex. Micro A	2.44	0.07 (0.09) 258	0.02 (0.08) 228	0.07 (0.09) 258	-0.03 (0.12) 228
Panel B: Effects on classes passed					
Num. classes passed	3.28	0.08 (0.09)	0.09 (0.10)	0.05 (0.09)	0.02 (0.09)
Num. classes not passed	0.31	0.01 (0.06)	-0.01 (0.06)	0.01 (0.06)	-0.01 (0.06)
Num. classes withdrawn	0.05	0.01 (0.03)	0.01 (0.02)	0.01 (0.03)	0.01 (0.02)
Panel C: Effects on class grade type					
Letter grade in Micro A	0.95	-0.04 (0.03)	-0.05* (0.03)	-0.03 (0.02)	-0.04 (0.03)
% classes taken for letter	0.93	-0.01 (0.01)	-0.01 (0.02)	-0.01 (0.01)	-0.01 (0.02)
% classes taken P/NP	0.07	0.01 (0.01)	0.01 (0.02)	0.01 (0.01)	0.01 (0.02)
Observations		374	332	374	332
Treatment assignment controls		Yes	No	Yes	Yes
Demographic controls		No	No	Yes	Yes
Pair Fixed Effects		No	No	No	Yes

Note: This table reports coefficients on $Incentive_i$ from Equations 1. GPA is measured on a 4.0 scale and is only affected by courses taken for a letter grade. Courses taken for Pass/No Pass (P/NP) have no bearing on GPA, nor do withdrawn courses. Model (1) contains linear controls for midterm 1 score and year; (2) is the difference in means and standard errors calculated using the repeated sampling framework of Neyman (1923); (3) and (4) use the post-double-selection (PDS) procedure of Belloni, Chernozhukov, and Hansen (2014b) to select control variables then estimate treatment effects and standard errors. The control variables selected using PDS are listed in Table A7. Models (2) and (4) include only students whose matched-pair did not attrite from the experiment. *Control Mean* is the mean for the Control students included in models (1) and (3), which is nearly identical to the mean for the Control students included in models (2) and (4). Standard errors in parentheses are robust to heteroskedasticity. ***, **, and * indicate significance at the 1, 5, and 10 percent critical levels, respectively.

Table 6: Spillover Effects of Incentive on Other Studying

	Control Mean	(1)	(2)	(3)	(4)
Attendance checks	5.91	-0.08 (0.18)	-0.09 (0.17)	-0.16 (0.17)	-0.10 (0.18)
Discussion board views	49.81	10.64 (7.64)	8.51 (8.25)	10.64 (7.60)	3.69 (8.05)
Discussion board days online	10.40	1.43 (1.55)	1.89 (1.59)	1.43 (1.54)	1.67 (1.65)
Discussion board questions asked	0.53	0.32 (0.25)	0.30 (0.30)	0.32 (0.25)	0.30 (0.31)
Discussion board answers	0.47	0.08 (0.26)	0.01 (0.28)	0.08 (0.26)	-0.02 (0.28)
Tutoring visits	0.41	0.05 (0.13)	-0.01 (0.14)	0.07 (0.12)	0.00 (0.12)
Observations		374	332	374	332
Treatment assignment controls		Yes	No	Yes	Yes
Demographic controls		No	No	Yes	Yes
Pair Fixed Effects		No	No	No	Yes

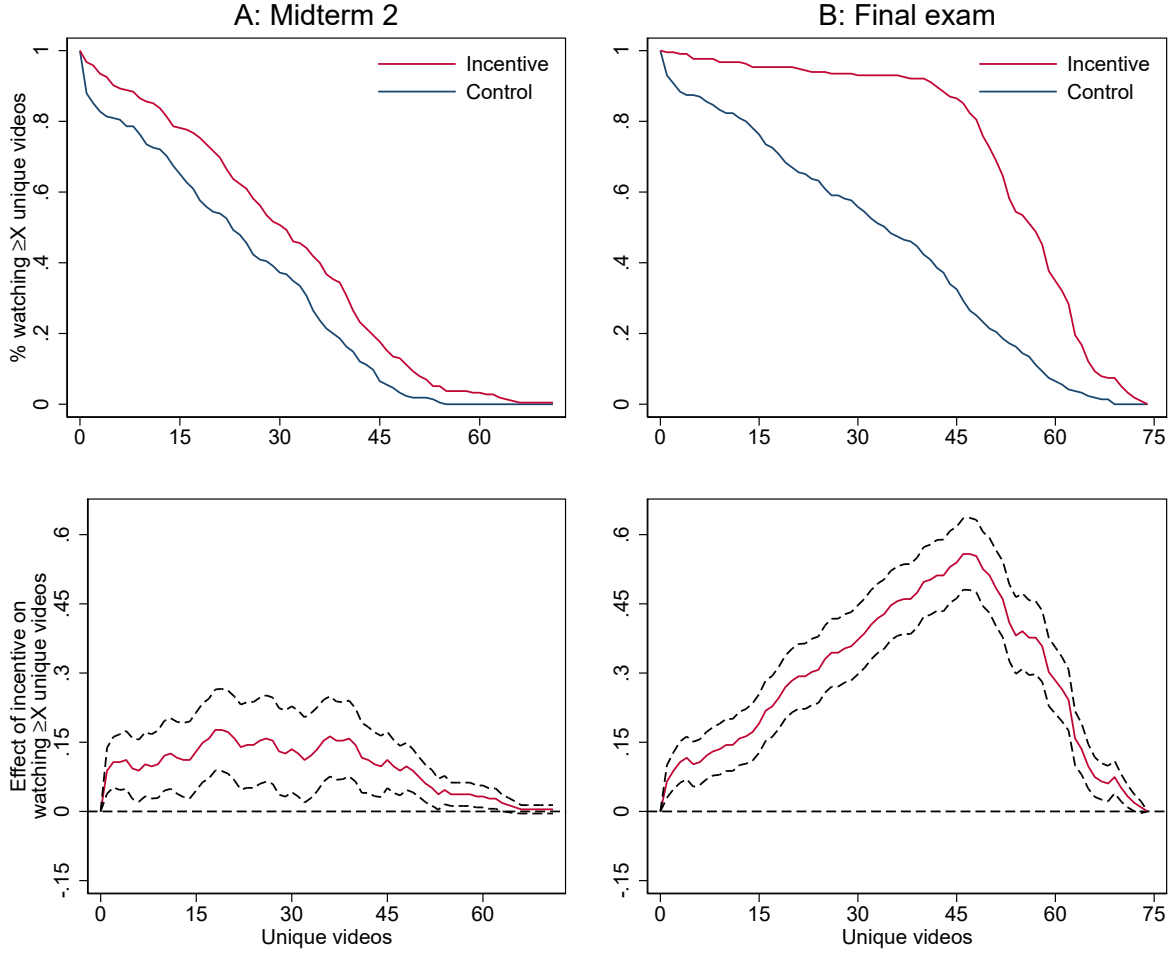
Note: This table reports coefficients on $Incentive_i$ from Equations 1. There were seven *Attendance checks* during the quarter. *Tutoring visits* includes those after the first midterm. Model (1) contains linear controls for midterm 1 score and year; (2) is the difference in means and standard errors calculated using the repeated sampling framework of Neyman (1923); (3) and (4) use the post-double-selection (PDS) procedure of Belloni, Chernozhukov, and Hansen (2014b) to select control variables then estimate treatment effects and standard errors. The control variables selected using PDS are listed in Table A7. Models (2) and (4) include only students whose matched-pair did not attrite from the experiment. *Control Mean* is the mean for the Control students included in models (1) and (3), which is nearly identical to the mean for the Control students included in models (2) and (4). Standard errors in parentheses are robust to heteroskedasticity. ***, **, and * indicate significance at the 1, 5, and 10 percent critical levels, respectively.

Table 7: Spillover Effects during Subsequent Quarter

	Control Mean	(1)	(2)	(3)	(4)
Panel A: Videos during subsequent quarter					
Num. of videos	25.46	14.00*** (4.45)	12.78* (6.74)	11.70*** (4.24)	11.35 (7.08)
Num. unique videos	19.77	9.87*** (3.03)	8.85** (4.04)	8.25*** (2.92)	8.07** (4.12)
Hours of videos	3.82	2.14*** (0.68)	1.88* (1.03)	1.79*** (0.64)	1.70 (1.08)
Hours unique videos	2.90	1.51*** (0.45)	1.33** (0.60)	1.27*** (0.44)	1.22** (0.61)
Observations		211	108	211	108
Panel B: Effects on classes passed					
Midterm 1 score		-0.04 (0.13) 213	-0.24 (0.18) 112	-0.04 (0.13) 213	-0.30 (0.19) 112
Midterm 2 score		0.00 (0.13) 214	-0.04 (0.20) 112	0.00 (0.13) 214	0.03 (0.21) 112
Final exam score		0.12 (0.14) 211	0.00 (0.18) 108	0.12 (0.14) 211	0.23 (0.23) 108
Panel C: Effects on class grade type					
Took Micro B	0.61	-0.07 (0.05)	-0.07 (0.05)	-0.07 (0.05)	-0.08 (0.06)
Num. classes passed	3.46	-0.07 (0.11)	-0.05 (0.12)	-0.07 (0.11)	-0.04 (0.12)
Num. classes not passed	0.23	0.07 (0.06)	0.08 (0.06)	0.07 (0.06)	0.07 (0.06)
Num. classes withdrawn	0.06	0.04 (0.03)	0.04 (0.03)	0.04 (0.03)	0.03 (0.03)
Observations		374	332	374	332
Treatment assignment controls		Yes	No	Yes	Yes
Demographic controls		No	No	Yes	Yes
Pair Fixed Effects		No	No	No	Yes

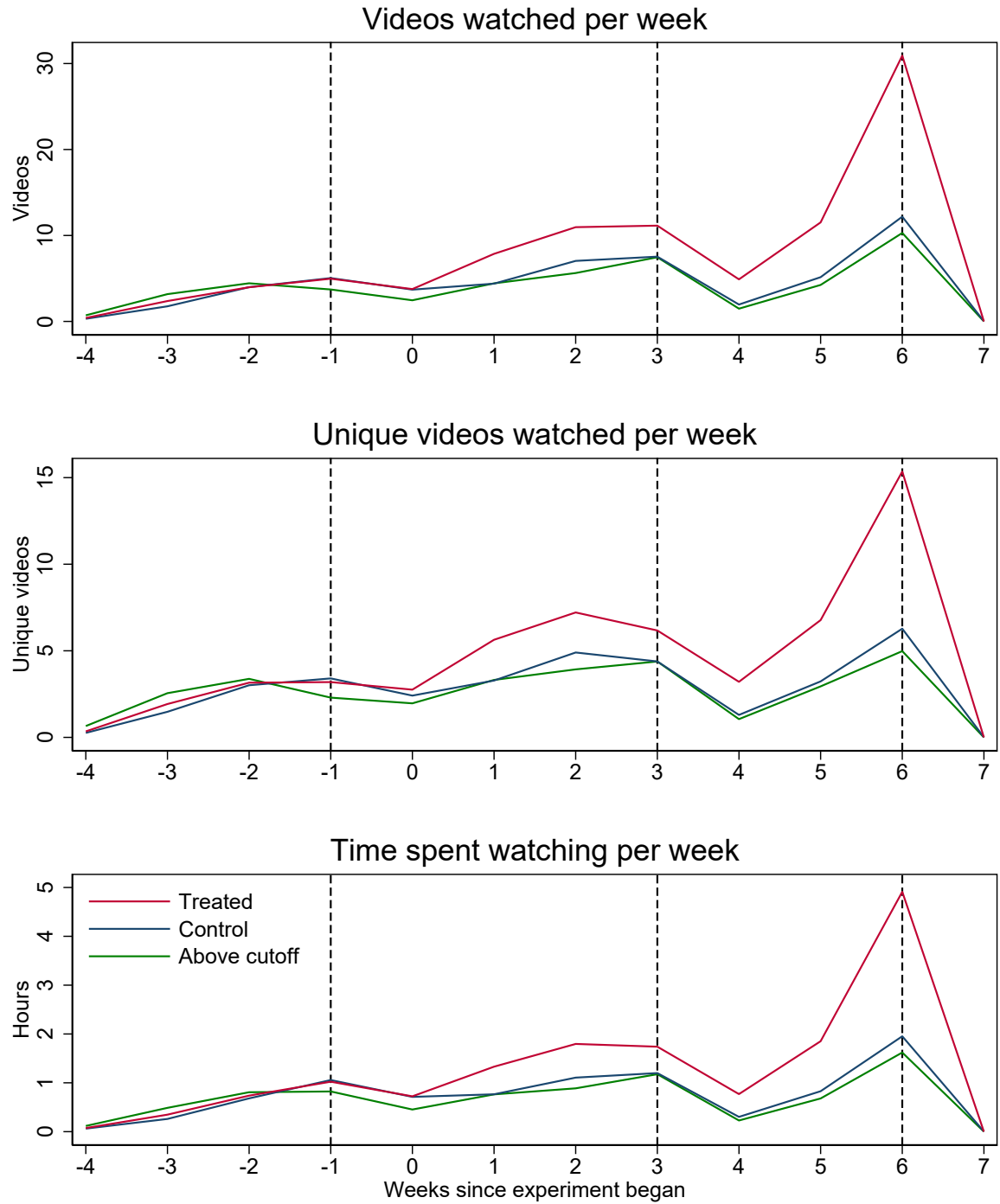
Note: This table reports coefficients on $Incentive_i$ from Equations 1. Panel A restricts the sample to those who completed both the first and second microeconomics courses (Micro A and B). Panel C includes those who completed the first microeconomics course (Micro A). Test scores are measured in standard deviation units. Model (1) contains linear controls for midterm 1 score and year; (2) is the difference in means and standard errors calculated using the repeated sampling framework of Neyman (1923); (3) and (4) use the post-double-selection (PDS) procedure of Belloni, Chernozhukov, and Hansen (2014b) to select control variables then estimate treatment effects and standard errors. The control variables selected using PDS are listed in Table A7. Models (2) and (4) include only students whose matched-pair did not attrite from the experiment. *Control Mean* is the mean for the Control students included in models (1) and (3), which is nearly identical to the mean for the Control students included in models (2) and (4). Standard errors in parentheses are robust to heteroskedasticity. ***, **, and * indicate significance at the 1, 5, and 10 percent critical levels, respectively.

Figure 3: Effect of grade incentive on videos watched



Top panels display the percent students in the *Control* and *Incentive* arms that watched at least X unique videos (left) or hours of unique videos (right). Bottom panels display the differences between the two arms in the top panels with 95% confidence intervals estimated by regressing an indicator for whether on the student watched at least $X \in \{0, \dots, X_{max}\}$ unique videos (or hours of unique videos) on the student's treatment status.

Figure 4: Weekly video watching by treatment arm



Dashed lines represent Midterm 1, Midterm 2, and Final exams

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. Imbens and Rubin, 2015, Athey and Imbens, 2017). The

Additionally, we were unable to observe most pretreatment covariates until after the experiment had concluded because of student privacy considerations, thereby making it impossible to block on these variables. We learned from the previous cohorts' data that between the first midterm score and math quiz score, both observable at the time of randomization, the midterm score predicted significantly more variation in the final exam score. Hence, we stratified on midterm score when assigning treatment. While we could have used an alternative method (e.g. matching methods) that take into consideration multiple covariates when assigning treatment, we opted for a simpler design given the high correlation between midterm and math quiz score and the comparatively high number of missing observations for the latter assessment (the math quiz was given on the second class day and so before some students enrolled in the class).

We assigned treatment shortly after issuing the first midterm exam grades, which occurred during the fourth week of the quarter. To assign treatment, we ordered the students by exam score, then paired students along this ordering for students below the median. Within pairs, we randomly assigned one student to *Incentive*, the other to *Control*. By construction, these two arms were *ex ante* balanced on midterm exam score, and we verified at time of treatment that the arms were also balanced on math quiz score. Since this random-

ization was performed independently across year cohorts, by construction, the samples were also balanced on year.

Although our treatment assignment method provides a better chance of balance than does simple random sampling, by random chance and through non-random attrition, it is possible that the two treatment arms vary on *ex post* observable and unobservable covariates that are correlated with the outcomes of interest, thereby confounding our treatment effect estimates. The primary cause of attrition was withdrawing from the course, which reduced our experiment sample by 35 students before the second midterm and an additional 21 students before the final exam. A 13% withdraw rate is in line with the withdraw rates observed in previous quarters. Another cause of attrition, albeit not from the course, is age: four students under the age of 18 during the experiment were removed from the analysis dataset. Additionally, seven students opted out of having their data included in the experiment analysis.

Since neither the students' intent to withdraw, age, nor opt-out preferences were observable at the time of treatment assignment, we could not *ex ante* balance this attrition across treatment arms. If students attrited non-randomly, that is, decided to attrite depending on their treatment status, then our treatment effect estimates would be biased. Fortunately, despite 8% attrition before the second midterm and 13% before the final exam, the two treatment arms below the median are balanced on nearly all observable pretreatment covariates, as shown in Tables A2 and A1, which gives us confidence that the *Control* arm is a good counterfactual for the *Incentive* arm.

A.2 Selection of control variables

In this section we discuss how we select control variables included in our linear models.

Equation 1 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 (Athey & Imbens, 2017). 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 amount of variance in final exam score. Hence, in our estimation strategies including controls, we always include the first midterm score and year, following the recommendations of Bruhn and McKenzie (2009) 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 lack the luxury of guaranteed balance by construction, nor is it clear *ex ante*, beyond our intuition, which will predict variation in the outcome variables of interest. On one hand, failing to control for valid predictors reduces statistical power. On the other hand, hand-picking control variables increases researcher degrees of freedom, risking increasing the prevalence of Type I errors (Simmons, Nelson, & Simonsohn, 2011). As such, in addition to a model without controls beyond the ones used for treatment assignment (year and midterm score), we fit a second model that includes a vector of linear controls chosen using the post-double-selection (PDS) procedure introduced by Belloni, Chernozhukov, and Hansen (2014b).

PDS is a two step process in which first, model covariates are selected in an automated, principled fashion, and second, the model coefficients of interest are estimated while controlling for those selected covariates. The first step involves predicting, separately, both the outcome of interest (e.g., videos watched) and treatment status using lasso regression, which shrinks coefficient estimates towards zero. Note that since treatment is randomly assigned, the lasso should shrink most, if not all, of the coefficients towards zero when predicting treatment status. Next, the researcher takes the union of all covariates with non-zero coefficients and includes these covariates as controls in her model. With her control variables selected, she can now estimate treatment effects with reduced bias relative to including controls with less empirical rationale.

In Table A6, we describe all covariates observable in our study. In Table A7, we describe the covariates selected as controls for estimating the effect of treatment on each outcome

variable of interest. All models include either pair fixed effects or year and midterm score as controls. To ensure these controls are “selected” by the PDS procedure, we partialled out these controls from the first step prediction models by residualizing both sides of the equation as described in Belloni, Chernozhukov, and Hansen (2014a).

B LATE estimators using Neyman’s repeated sampling approach

In this section we derive LATE estimators using the repeated sampling approach of Neyman (1923), which considers each pair as an independent, completely randomized experiment.

Similar to a Wald estimator, the point estimate of the LATE is the mean within-pair difference in outcome divided by the mean within-pair difference in videos:

$$\hat{\gamma} = \frac{\bar{\Delta y}}{\bar{\Delta v}} = \frac{\frac{1}{J} \sum_{j=1}^J \Delta y_j}{\frac{1}{J} \sum_{j=1}^J \Delta v_j} = \frac{\bar{y}_I - \bar{y}_C}{\bar{v}_I - \bar{v}_C} \quad (11)$$

where y is the outcome of interest (grades) and v is the number of videos, both indexed by pair $j \in J$ and treatment status C or I for *Control* or *Incentive*, respectively.

We use the delta method to calculate the approximate standard error of $\hat{\gamma}$. First, we define the following normally-distributed random variables:

$$\begin{aligned} Y &= \bar{y}_I - \bar{y}_C \sim \mathcal{N}(\mu_Y, \sigma_Y^2) \\ V &= \bar{v}_I - \bar{v}_C \sim \mathcal{N}(\mu_V, \sigma_V^2) \end{aligned} \quad (12)$$

Using a first-order Taylor expansion and letting $g() = \frac{Y}{V}$, we have:

$$\begin{aligned} \text{Var}(g) &= \text{E}[(g - \text{E}(g))^2] \\ &\approx \text{E}[(g(\theta) + (Y - \theta_Y)g'_Y(\theta) + (V - \theta_V)g'_V(\theta) - g(\theta))^2] \\ &= \text{E}[(Y - \theta_Y)^2(g'_Y(\theta))^2 + (V - \theta_V)^2(g'_V(\theta))^2 + 2(Y - \theta_Y)(V - \theta_V)g'_Y(\theta)g'_V(\theta)] \\ &= \text{Var}(Y)(g'_Y(\theta))^2 + \text{Var}(V)(g'_V(\theta))^2 + 2\text{Cov}(Y, V)g'_Y(\theta)g'_V(\theta) \end{aligned} \quad (13)$$

Expanding about $\theta = (\theta_Y, \theta_V) = (\mu_Y, \mu_V)$ and letting $g'_Y(\theta) = \mu_V^{-1}$ and $g'_V(\theta) = \frac{-\mu_Y}{\mu_V^2}$:

$$\begin{aligned}\text{Var}(g) &\approx \frac{1}{\mu_V^2} \text{Var}(Y) + \frac{\mu_Y^2}{\mu_V^4} \text{Var}(V) + 2 \frac{-\mu_Y}{\mu_V^2} \text{Cov}(Y, V) \\ &= \frac{\mu_Y^2}{\mu_V^2} \left(\frac{\sigma_Y^2}{\mu_Y^2} + \frac{\sigma_V^2}{\mu_V^2} - 2 \frac{\text{Cov}(Y, V)}{\mu_Y \mu_V} \right)\end{aligned}\tag{14}$$

We use the following variance estimators of Y and V from Equation 5:

$$\begin{aligned}\text{Var}(\hat{Y}) &= \hat{\sigma}_Y^2 = \frac{1}{J(J-1)} \sum_{j=1}^J (\Delta y_j - \bar{\Delta y})^2 \\ \text{Var}(\hat{V}) &= \hat{\sigma}_V^2 = \frac{1}{J(J-1)} \sum_{j=1}^J (\Delta v_j - \bar{\Delta v})^2 \\ \text{Cov}(\hat{Y}, \hat{V}) &= \hat{\sigma}_{YV} = \frac{1}{J(J-1)} \sum_{j=1}^J (\Delta y_j - \bar{\Delta y})(\Delta v_j - \bar{\Delta v})\end{aligned}\tag{15}$$

and the following estimators for the population means of Y and V :

$$\begin{aligned}\hat{\mu}_Y &= E(\mu_Y) = \bar{\Delta y} \\ \hat{\mu}_V &= E(\mu_V) = \bar{\Delta v}\end{aligned}\tag{16}$$

Substituting these variance and means estimators into the final step of 14, we arrive at the standard error estimator for $\hat{\gamma}$:

$$\hat{\sigma}_\gamma = \frac{\bar{\Delta y}}{\bar{\Delta v}} \sqrt{\frac{\hat{\sigma}_Y^2}{\bar{\Delta y}^2} + \frac{\hat{\sigma}_V^2}{\bar{\Delta v}^2} - 2 \frac{\hat{\sigma}_{YV}}{\bar{\Delta y} \bar{\Delta v}}}\tag{17}$$

Table A1: Baseline balance test, Midterm 2 sample

Variable	All students			P-values	Matched pairs		P-values
	Above Median	Control	Incentive	(3) - (2)	Control	Incentive	(5) - (4)
Midterm 1 score	2.048 (0.025)	0.116 (0.063)	0.037 (0.068)	0.398	0.139 (0.065)	0.131 (0.066)	0.933
Year = 2019	0.492 (0.025)	0.513 (0.036)	0.500 (0.035)	0.797	0.514 (0.037)	0.514 (0.037)	1.000
Cumulative GPA	3.445 (0.029)	2.944 (0.043)	2.948 (0.058)	0.965	2.942 (0.045)	2.992 (0.056)	0.487
No cum. GPA	0.230 (0.021)	0.368 (0.035)	0.332 (0.033)	0.452	0.365 (0.036)	0.320 (0.035)	0.377
Math quiz score	0.592 (0.044)	0.037 (0.070)	0.106 (0.065)	0.471	0.054 (0.071)	0.137 (0.068)	0.396
Tutoring visits	0.269 (0.042)	0.259 (0.059)	0.223 (0.056)	0.655	0.276 (0.062)	0.232 (0.061)	0.612
Videos watched	13.228 (0.681)	13.368 (0.886)	13.777 (0.931)	0.750	13.663 (0.929)	13.729 (0.986)	0.961
Videos, unique	9.746 (0.431)	9.689 (0.580)	10.188 (0.611)	0.554	9.845 (0.606)	10.116 (0.644)	0.760
Hours videos	1.690 (0.093)	1.782 (0.127)	1.825 (0.135)	0.818	1.827 (0.133)	1.804 (0.142)	0.906
Hours videos, unique	1.291 (0.062)	1.355 (0.090)	1.387 (0.092)	0.802	1.382 (0.095)	1.364 (0.096)	0.897
Asian	0.700 (0.022)	0.694 (0.033)	0.668 (0.033)	0.581	0.713 (0.034)	0.652 (0.036)	0.215
Latinx	0.060 (0.012)	0.135 (0.025)	0.158 (0.026)	0.506	0.133 (0.025)	0.166 (0.028)	0.377
White	0.151 (0.018)	0.114 (0.023)	0.124 (0.023)	0.765	0.105 (0.023)	0.138 (0.026)	0.336
Other ethnicity	0.089 (0.014)	0.057 (0.017)	0.050 (0.015)	0.741	0.050 (0.016)	0.044 (0.015)	0.804
Female	0.393 (0.024)	0.342 (0.034)	0.391 (0.034)	0.312	0.343 (0.035)	0.392 (0.036)	0.328
Male	0.592 (0.024)	0.653 (0.034)	0.604 (0.034)	0.316	0.652 (0.036)	0.602 (0.036)	0.329
Transfer	0.271 (0.022)	0.477 (0.036)	0.455 (0.035)	0.673	0.470 (0.037)	0.436 (0.037)	0.528
Observations	417	193	202		181	181	

Note: This table includes all students who completed the second midterm. Descriptions of each variable can be found in Table A6. *Male* and *Female* are coded zero for nine students who do not report a gender. *P-values* are reported for the Welch's t-test of equal means between the *Control* and *Incentive* arms. Standard errors in parentheses are robust to heteroskedasticity. ***, **, and * indicate significance at the 1, 5, and 10 percent critical levels, respectively.

Table A2: Baseline balance test, Final Exam sample

Variable	All students			P-values (3) - (2)	Matched pairs		P-values (5) - (4)
	Above Median	Control	Incentive		Control	Incentive	
Midterm 1 score	2.049 (0.025)	0.153 (0.061)	0.057 (0.069)	0.291	0.177 (0.064)	0.170 (0.065)	0.938
Year = 2019	0.489 (0.025)	0.516 (0.037)	0.500 (0.036)	0.753	0.518 (0.039)	0.518 (0.039)	1.000
Cumulative GPA	3.445 (0.029)	2.946 (0.044)	2.959 (0.060)	0.864	2.929 (0.047)	3.001 (0.059)	0.346
No cum. GPA	0.231 (0.021)	0.359 (0.035)	0.332 (0.034)	0.583	0.367 (0.038)	0.313 (0.036)	0.299
Math quiz score	0.599 (0.043)	0.071 (0.068)	0.152 (0.066)	0.396	0.061 (0.071)	0.157 (0.071)	0.338
Tutoring visits	0.270 (0.043)	0.272 (0.061)	0.237 (0.060)	0.684	0.283 (0.066)	0.253 (0.066)	0.746
Videos watched	13.292 (0.682)	13.418 (0.909)	13.658 (0.953)	0.856	13.729 (0.978)	13.789 (1.023)	0.966
Videos, unique	9.793 (0.432)	9.783 (0.598)	10.111 (0.622)	0.704	9.795 (0.630)	10.181 (0.665)	0.674
Hours videos	1.698 (0.094)	1.788 (0.130)	1.805 (0.138)	0.929	1.812 (0.138)	1.803 (0.148)	0.967
Hours videos, unique	1.297 (0.062)	1.369 (0.093)	1.372 (0.094)	0.985	1.363 (0.098)	1.366 (0.100)	0.985
Asian	0.701 (0.022)	0.696 (0.034)	0.653 (0.035)	0.376	0.711 (0.035)	0.633 (0.038)	0.129
Latinx	0.060 (0.012)	0.141 (0.026)	0.158 (0.027)	0.654	0.139 (0.027)	0.169 (0.029)	0.448
White	0.149 (0.018)	0.109 (0.023)	0.132 (0.025)	0.497	0.102 (0.024)	0.145 (0.027)	0.244
Other ethnicity	0.089 (0.014)	0.054 (0.017)	0.058 (0.017)	0.882	0.048 (0.017)	0.054 (0.018)	0.804
Female	0.393 (0.024)	0.348 (0.035)	0.405 (0.036)	0.253	0.337 (0.037)	0.404 (0.038)	0.212
Male	0.593 (0.024)	0.647 (0.035)	0.584 (0.036)	0.215	0.657 (0.037)	0.584 (0.038)	0.176
Transfer	0.272 (0.022)	0.462 (0.037)	0.447 (0.036)	0.778	0.470 (0.039)	0.416 (0.038)	0.321
Observations	415	184	190		166	166	

Note: This table includes all students who completed the final exam. Descriptions of each variable can be found in Table A6. *Male* and *Female* are coded zero for nine students who do not report a gender. *P-values* are reported for the Welch's t-test of equal means between the *Control* and *Incentive* arms. Standard errors in parentheses are robust to heteroskedasticity. ***, **, and * indicate significance at the 1, 5, and 10 percent critical levels, respectively.

Table A3: Anderson-Rubin confidence sets

Outcome variable	Endogenous variable	Anderson-Rubin CI
Midterm 2 score	10 unique videos	[-0.001, 0.653]
Midterm 2 score	1 hour videos	[-0.001, 0.392]
Final exam score	10 unique videos	[0.000, 0.163]
Final exam score	1 hour videos	[0.000, 0.102]

Note: This table displays Anderson-Rubin confidence sets at the 95% confidence level for the 2SLS estimator $\hat{\gamma}$ from Equation 8 including year dummies and first midterm score as controls. Outcomes are measured in standard deviations. Instrumented endogenous variables are measured in 10s of unique videos or hours of unique content.

Table A4: Heterogeneous Effects of Treatment

Interaction Variable	Midterm 2	Final Exam
Midterm 1 score	0.009 (0.086)	0.002 (0.095)
Year = 2019	0.068 (0.181)	0.017 (0.179)
Pretreatment videos	0.011 (0.007)	-0.004 (0.007)
Pretreatment videos, unique	0.020* (0.011)	0.001 (0.011)
Transfer	-0.139 (0.183)	0.130 (0.178)
Female	-0.208 (0.185)	-0.322* (0.188)
Asian	-0.402** (0.189)	-0.215 (0.180)
Latinx	0.237 (0.254)	0.022 (0.227)
White	0.638** (0.260)	0.271 (0.250)
Other ethnicity	-0.113 (0.307)	0.276 (0.347)
Observations	395	374

Note: This table reports estimates for β_2 from Equation 9. Test scores are measured in standard deviation units. Standard errors in parentheses are robust to heteroskedasticity. ***, **, and * indicate significance at the 1, 5, and 10 percent critical levels, respectively.

Table A5: Regression Discontinuity estimates

	Control Mean	(1)	(2)
Panel A: Midterm 2			
Unique videos	22.7	9.12*** (3.35)	9.34*** (3.42)
Midterm 2 score	0.73	0.06 (0.20)	0.07 (0.19)
Observations		812	779
Panel B: Final exam			
Unique videos	32.0	24.16*** (2.39)	24.95*** (2.48)
Final exam score	0.78	0.11 (0.18)	0.14 (0.19)
Observations		789	747

Note: This table reports coefficients on $\tau(c)$ from Equations 10. Panel A includes Midterm 2-takers while Panel B includes Final Exam-takers. Model (1) includes all students while (2) contains only students whose pair did not attrite. Standard errors are in parentheses. ***, **, and * indicate significance at the 1, 5, and 10 percent critical levels, respectively.

Table A6: Candidate control variables for post-double-selection

Variable	Description
Midterm 1 score	Score on the first midterm
Year = 2019	1 if course taken in 2019, 0 otherwise
Cumulative GPA	Cumulative GPA from prior term, 0 if not observed
No cum. GPA	1 if Cumulative GPA unobserved, 0 otherwise
Math quiz score	Score on a quiz assessing prerequisite math skills
Tutoring visits	Number of group tutoring lab visits as of the first midterm
Videos watched	Number unique videos watched as of the first midterm
Hours videos	Hours of unique videos watched as of the first midterm
Asian	1 if ethnicity is Asian, 0 otherwise
Latinx	1 if ethnicity is Latinx, 0 otherwise
White	1 if ethnicity is White, 0 otherwise
Female	1 if female, 0 otherwise
Transfer	1 if transfer student, 0 otherwise

Note: *Midterm 1 score* and *Math quiz score* are measured in control standard deviations. *Cumulative GPA* is measured on a 4.0 scale. Videos included in *Videos watched* and *Hours videos* are unique course-relevant videos. The ethnicity variables are coded by university records: *Asian* includes "Chinese/Chinese American", "Vietnamese", "East Indian/Pakistani", "Japanese/Japanese American", "Korean/Korean American", and "All other Asian/Asian American"; *Latinx* includes "Mexican/Mexican American", "Chicano", and "All other Spanish-American/Latino"; *White* includes "White/Caucasian"; and the omitted category includes "African American/Black", "Pacific Islander", and "Not give/declined to state".

Table A7: ITT model controls selected via post-double-selection

Table	Dependent Variable	Controls, All Observations	Controls, Fixed Effects
Table 1	Hours unique videos by Final	Hours videos Videos	Hours videos Videos
	Hours unique videos by Mid. 2	Hours videos	Hours videos Videos
	Hours videos by Final	Hours videos	Hours videos
	Hours videos by Mid. 2	Hours videos	Hours videos Tutoring visits Videos
	Num. unique videos before Final	Hours videos Videos	Videos
	Num. unique videos before Mid. 2	Hours videos Videos	Videos
	Num. videos before Final	Hours videos Videos	Hours videos Videos
	Num. videos before Mid. 2	Hours videos Videos	Hours videos Tutoring visits Videos
Table 2	Final exam score	None	Math quiz score Transfer
	Midterm 2 score	None	Math quiz score
Table 3	All classes	Cumulative GPA	Cumulative GPA Math quiz score Transfer
	Econ classes ex. Micro A	None	Cumulative GPA Transfer
	Excluding Micro A	Cumulative GPA	Transfer
	Excluding econ classes	None	None
	Letter grade in Micro A	Cumulative GPA Latinx Transfer	Cumulative GPA
	Num. classes not passed	None	None
	Num. classes passed	Cumulative GPA Transfer	Cumulative GPA Transfer
	Num. classes taken P/NP	Latinx	Latinx
	Num. classes taken for letter	Cumulative GPA No cum. GPA	Cumulative GPA
	Num. classes withdrawn	None	None
	Num. units taken P/NP	Latinx	Latinx
	Num. units taken for letter grade	Cumulative GPA No cum. GPA	Cumulative GPA
	Num. units withdrawn	None	None
	% classes taken P/NP	None	Latinx
	% classes taken for letter	None	Latinx

Continued on next page

Table A7 (continued)

Table 4	Attendance checks	Female	Tutoring visits
		Math quiz score	
		Tutoring visits	
	Discussion board answers	None	None
	Discussion board days online	None	None
	Discussion board questions asked	None	None
	Discussion board views	None	Asian
	Tutoring visits	Tutoring visits	Tutoring visits
Table 5	Hours of videos	Hours videos	Hours videos
			Latinx
			Math quiz score
			Tutoring visits
			Videos
	Midterm 1 score	None	Latinx
			Math quiz score
			Videos
	Midterm 2 score	None	Asian
			Latinx
			Math quiz score
			Videos
	Num. classes not passed	None	None
	Num. classes passed	None	None
	Num. classes taken P/NP	None	Transfer
	Num. classes taken for letter	None	No cum. GPA
	Num. classes withdrawn	None	None
	Num. of videos	Hours videos	Hours videos
			Latinx
			Math quiz score
			Tutoring visits
			Videos
	Num. units taken P/NP	None	Transfer
	Num. units taken for letter grade	None	None
	Num. units withdrawn	None	None
	Term GPA	Cumulative GPA	Cumulative GPA
			Tutoring visits
	Term GPA, econ courses ex. Micro B, winter	None	Math quiz score
	Term GPA, ex. Micro B	Cumulative GPA	Cumulative GPA
			Tutoring visits
	Term GPA, ex. econ courses	None	Tutoring visits
	Took Micro B	None	Math quiz score
	% classes taken P/NP	None	No cum. GPA
			Transfer
	% classes taken for letter	None	No cum. GPA
			Transfer
Table	Final exam score	None	Latinx
None			Math quiz score
			Videos

Continued on next page

Table A7 (continued)

Hours unique videos	Hours videos	Hours videos Latinx Math quiz score Tutoring visits Videos
Num. unique videos	Hours videos	Hours videos Latinx Math quiz score Tutoring visits Videos
Pass Micro B	None	Latinx Math quiz score Videos

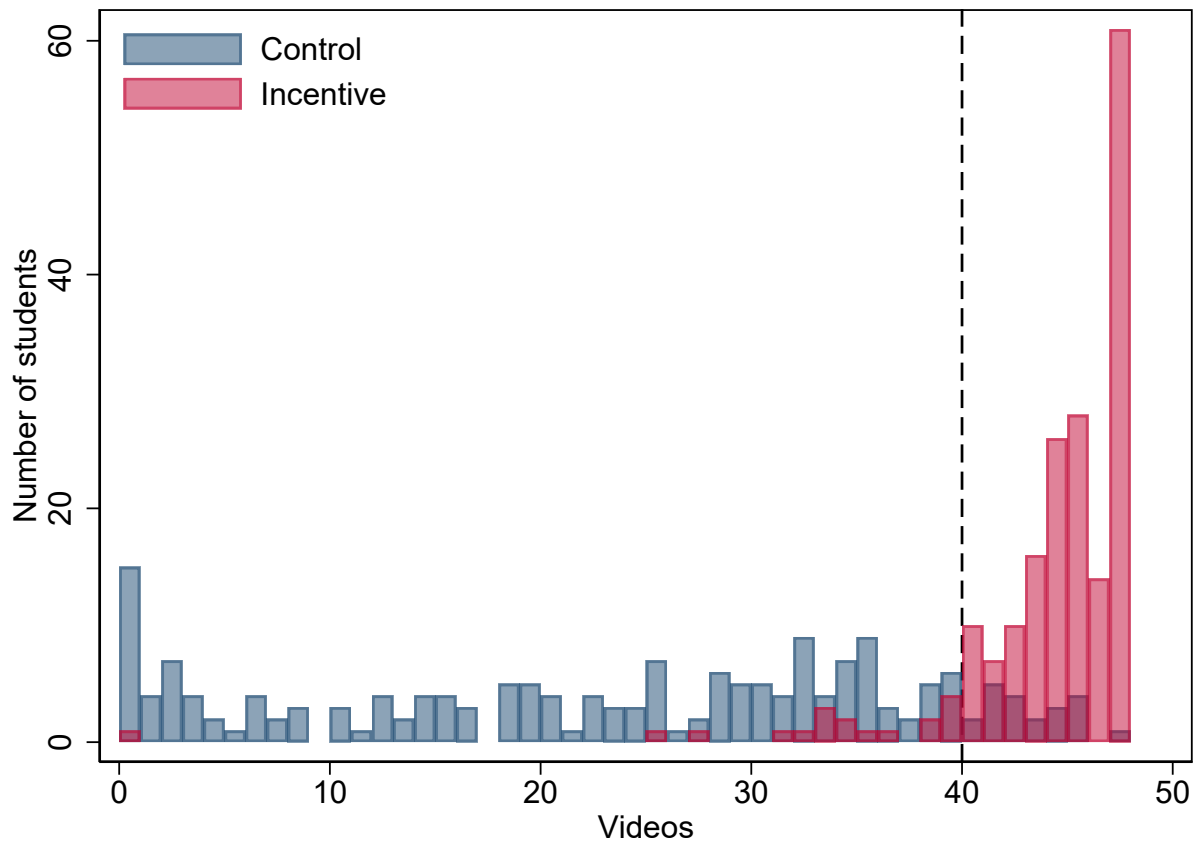
Note: Controls chosen via the PDS procedure of Belloni, Chernozhukov, and Hansen (2014b). In the *All Observations* model, *Midterm 1 score* and *Year = 2019* are additionally included as controls. In the *Fixed Effects* model, pair fixed effects and *Midterm 1 score* are included. All control variables are measured before the start of the experiment, e.g. *Hours videos* is the hours of videos watched as of the first midterm.

Table A8: LATE model controls selected via post-double-selection

Dependent Variable	Instrumented	Controls, All Observations	Controls, Fixed Effects
Final exam score	Hours videos, unique	Hours videos Math quiz score Transfer Videos	Hours videos Videos
Final exam score	Videos, unique	Hours videos Videos	Hours videos Videos
Midterm 2 score	Hours videos, unique	Hours videos Math quiz score Tutoring visits Videos	Hours videos
Midterm 2 score	Videos, unique	Hours videos Videos	Hours videos Videos

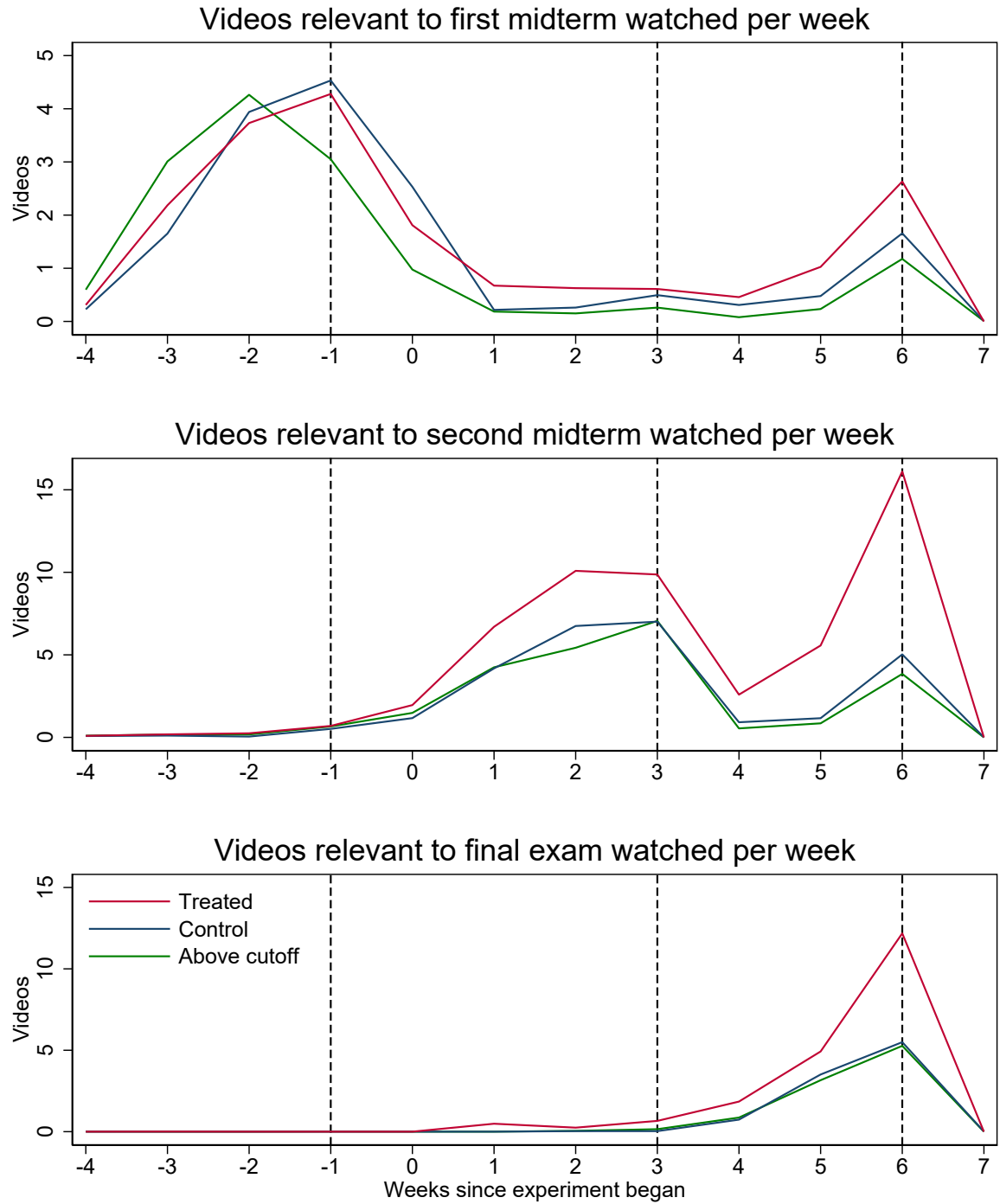
Note: Controls chosen via the PDS procedure of Belloni, Chernozhukov, and Hansen (2014b). In the *All Observations* model, *Midterm 1 score* and *Year = 2019* are additionally included as controls. In the *Fixed Effects* model, pair fixed effects and *Midterm 1 score* are included. All control variables are measured before the start of the experiment, e.g. *Hours videos* is the hours of videos watched as of the first midterm.

Figure A1: Distribution of videos counted towards incentive



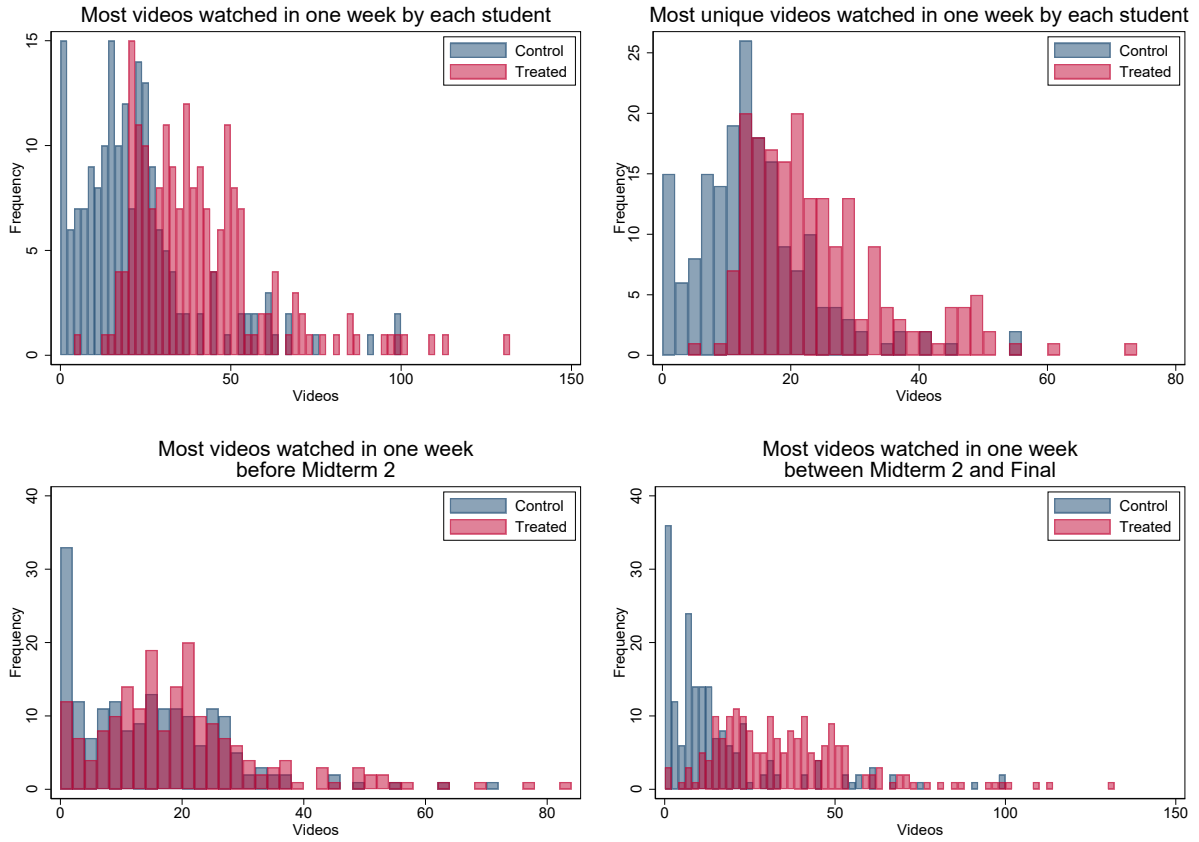
This plot includes only videos that would have counted towards the earning the grade incentive. Students were required to watch 40 unique of 48 eligible videos between the first midterm and final exam to earn the grade incentive. 91% of *Incentive* students met the requirements for the grade incentive versus 11% of *Control* students.

Figure A2: Weekly video watching by exam topic



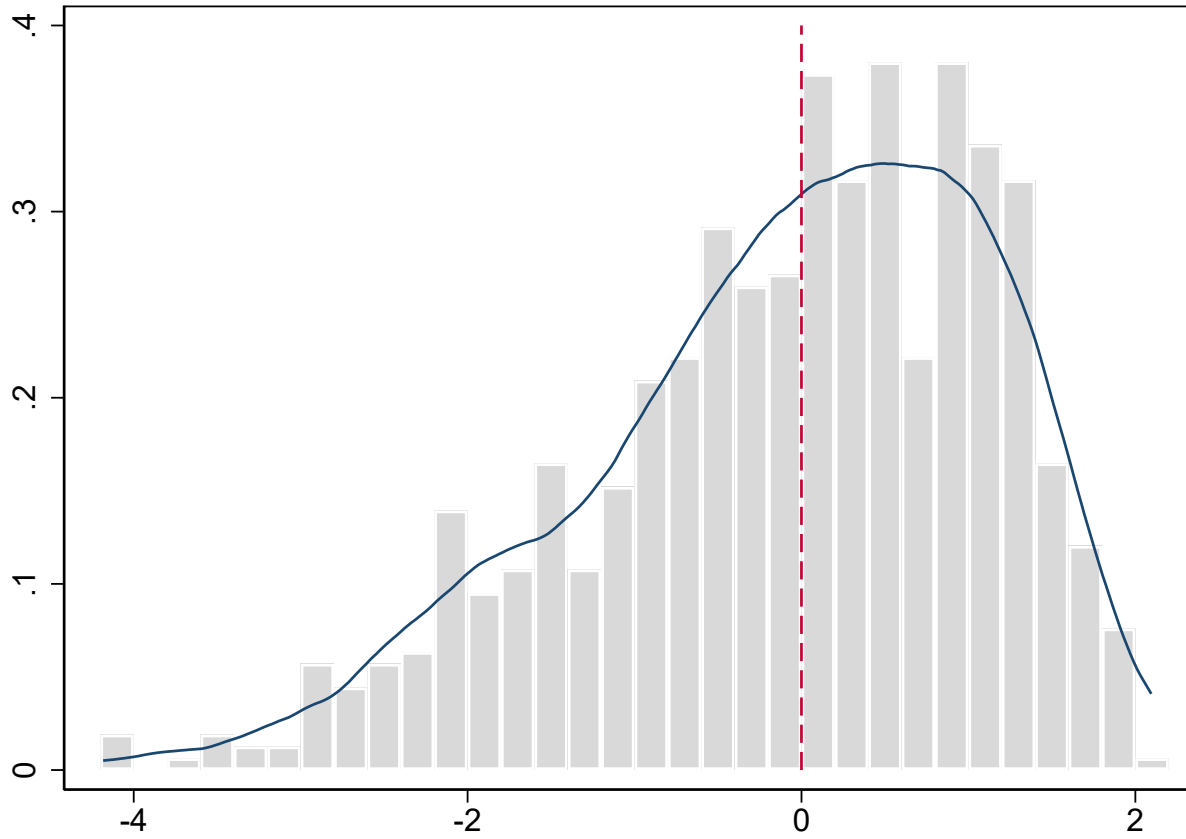
Dashed lines represent Midterm 1, Midterm 2, and Final exams.

Figure A3: Distribution of max videos watched in one week



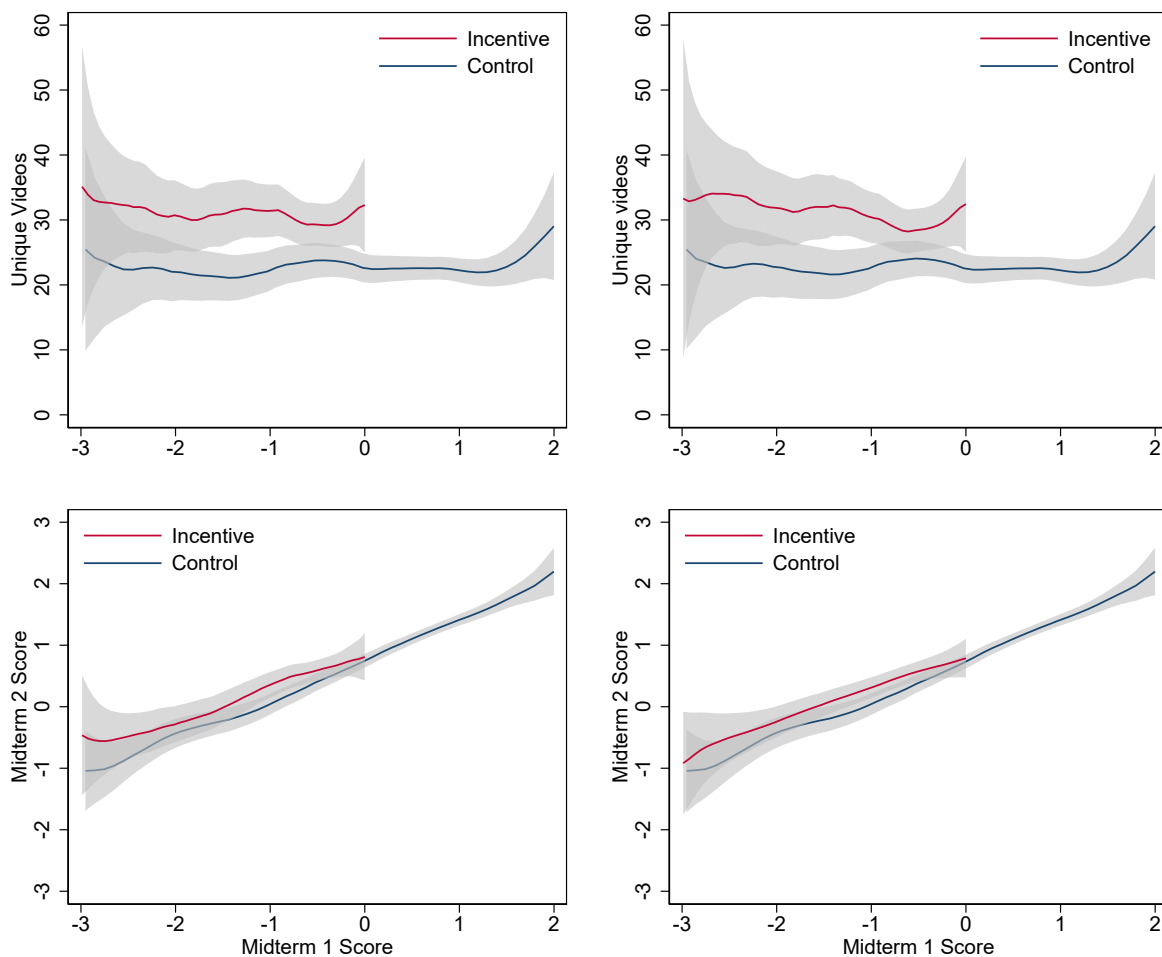
These plots help illustrate potential “binge watching” behavior. Compared to the *Control* students, *Incentive* students are more likely to watch 40 or more unique videos in a week, which occurs in the weeks preceding the final and not the second midterm.

Figure A4: Distribution of Midterm 1 scores



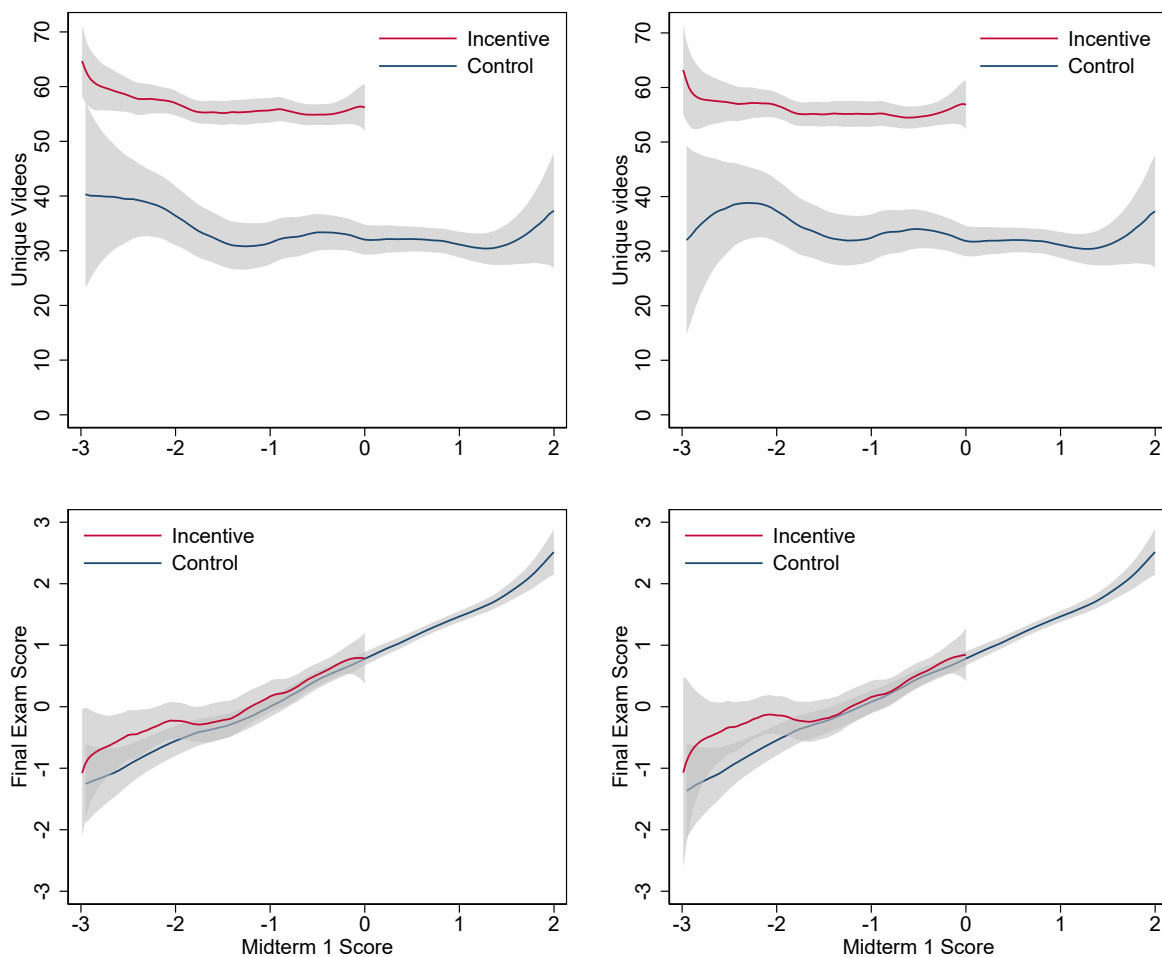
Midterm scores are measured in control standard deviations. This histogram includes all students who took the final exam. As such, this plot would allow one to observe bunching if treatment caused differential attrition or manipulation of one's midterm score.

Figure A5: Regression Discontinuity: effect of treatment on videos watched and second midterm exam scores



Videos (top) includes unique videos watched before the second midterm exam. Exam scores (bottom) are measured in control standard deviations. Confidence bands represent 95% confidence intervals of the conditional mean outcome. The left plots includes all students who took the second midterm while the right plots exclude any students whose matched pair attrited.

Figure A6: Regression Discontinuity: effect of treatment on videos watched and final exam scores



Videos (top) includes unique videos watched before the final exam. Exam scores (bottom) are measured in control standard deviations. Confidence bands represent 95% confidence intervals of the conditional mean outcome. The left plots includes all students who took the final exam while the right plots exclude any students whose matched pair attrited.