

Sleep Educational Impact and Habit Formation

Osea Giuntella
University of Pittsburgh
IZ and NBER

Silvia Saccardo
Carnegie Mellon University

Sally Sado
UC San Diego

March 4, 2024

Abstract

There is growing evidence on the importance of sleep for productivity, but little is known about the impact of interventions targeting sleep. In a field experiment among U.S. university students, we show that incentives for sleep increase both sleep and academic performance. Motivated by theories of cue-based habit formation, our primary intervention couples personalized bedtime reminders with morning feedback and immediate rewards for sleeping at least seven hours on weeknights. The intervention increases the share of nights with at least seven hours of sleep by 28 percent and average weeknight sleep by an estimated 19 minutes during a four-week treatment period, with persistent effects of about 9 minutes per night during a one to five-week post-treatment period. Comparisons to secondary treatments show that immediate incentives have larger impacts on sleep than delayed incentives or reminders and feedback alone during the treatment period, but do not have statistically distinguishable impacts on longer-term sleep habits in the post-treatment period. We estimate that immediate incentives improve average semester course performance by 0.075 - 0.088 grade points, a 0.10 - 0.11 standard deviation increase. Our results demonstrate that incentives to sleep can be a cost-effective tool for improving educational outcomes.

JEL CODES: J10, I10, I23

Keywords: Sleep, Human Capital, Education, Habit Formation, Health Behaviors

We are grateful to the seminar participants at the BEDI Conference, the Roybal Annual Retreat, University of Pittsburgh, the Center for Sleep and Circadian Rhythms Study, the Virtual Seminar on the Economics of Risky Behaviors (VERB), the CHIBE Behavioral Science and Health symposium, the University of Houston, the National University of Singapore, the Freie Universitat in Berlin, UCSD Spring School on Behavioral Economics, University of Chicago, Harvard University, RAND, University of Southern California, the UCLA Anderson School of Management, and the New York Federal Reserve, the NBER Economics of Health Meetings (Fall 2023). We are grateful to Mallory Avery, Daniel Banko-Ferran, Kelly Hyde, Ben Schenck, Samuel Lindquist, Gabrielle Toborg, and William Wang for their valuable research assistance. We are thankful to the Pittsburgh Experimental Economics Lab for recruiting our participants and the University of Pittsburgh Registrar Office for their help in providing access to students' academic records. We received generous funding from J-PAL North America (Sado), the National Institute of Aging (Saccardo, Grant #P30AG034546) and the Pitt Healthy Lifestyle Institute Pilot and Feasibility project (Giuntella).

1 Introduction

There is growing attention to the role of sleep for economic outcomes (Hillman et al., 2006; Mullainathan, 2014; Rao et al., 2021). At the same time, statistics suggest people are not sleeping enough. About a third of Americans sleep less than the recommended minimum of seven hours per night and a similar proportion state that they would like to sleep more (Jones, 2013; Ballard, 2019; Corkett, 2010; CDC, 2023). Sleep deprivation is even worse among adolescents and young adults, with as few as a third regularly meeting sleep guidelines (Wheaton et al., 2018). This has prompted policy concern that poor sleep may be worsening educational outcomes (Group et al., 2014). More broadly, Roenneberg (2013) referred to sleep deprivation as the most prevalent high-risk behavior in modern societies. Yet, little is known about whether interventions targeting sleep can improve productivity and performance.

A growing body of research using naturally occurring data finds that sleep affects earnings and academic performance, as well as physical and mental health (see e.g., Lindquist and Sado , 2023, for a review). In experimental work, sleep lab studies find that short-term severe sleep deprivation worsens cognition and mood, but are not able to estimate the effect of moderate sleep improvements in natural contexts (Banks and Dinges, 2007; Killgore, 2010; Lim and Dinges, 2010). In contrast, the only field experiment to test the impact of sleep interventions does not find an impact of increased nighttime sleep on productivity among highly sleep deprived workers in India (Bessone et al., 2021). It remains an open question whether sustained exogenous increases in nighttime sleep can improve productivity and performance in the U.S., where average sleep is of high quality and closer to recommended guidelines.

Our study implements a randomized field experiment among U.S. university students to examine the impact of interventions targeting sleep on sleep habits and academic performance. We ran the experiment over seven semester-long waves from Spring 2019 to Spring 2022. The 1,149 participants wore tracking devices (Fitbits) that measure sleep, heart rate and physical activity, downloaded a custom smartphone app linked to their Fitbit data which delivered our interventions, and answered survey questions to capture information about their time use, cognitive performance and well-being. The study included a one- to four-week baseline period, followed by a four-week intervention period and a one- to five-week post-intervention period.

Our primary intervention aims to develop persistent sleep habits that extend beyond the four-week intervention period. To do so, we build on theories of cue-based habit formation, which underscore the role of context cues, repetition, and immediate reinforcement of a de-

sired action via rewards (Verplanken and Wood, 2006; Wood and Neal, 2007; Wood and Runger, 2016). These theories suggest that repeatedly rewarding a behavior performed in response to a consistent cue can gradually create an association between the cue and the reward that follows the action. Once the association is established, the cue may “automatically” trigger the desired action with little or no cognitive effort, even in the absence of a reward (Dickinson, 1985). Accordingly, we paired daily cues to go to sleep with immediate rewards for meeting sleep goals throughout the intervention period, and maintained the cue during the post-intervention period after the rewards were removed.

Specifically, in our treatment groups, we set a goal for participants to sleep at least seven hours by 9 am on weeknights (Sunday - Thursday), following recommended guidelines (Panel et al., 2015). During the intervention period, participants in our primary treatment received personalized bedtime reminders every weeknight, prompting them to follow a self-selected bedtime routine to get at least seven hours of sleep. On weekday mornings, they learned whether they had met their sleep goal and, if successful, received an immediate financial reward of \$4.75. In the post-intervention period, we stopped the financial reward but maintained bedtime cues and morning feedback to investigate the persistence of behavior change in response to the cue once the reward is removed. In secondary treatments, we tested variants that provided rewards with a delay rather than immediately; and, that turned off either the rewards or the cue and feedback.

Our primary analysis compares a no intervention *Control* group to the *Immediate Incentives* group, in which participants received bedtime cues, morning feedback and immediate incentives for each weeknight they met the sleep goal during the intervention period.¹ At baseline, participants met the goal of sleeping at least seven hours on approximately 43% of nights. During the treatment period, the intervention increases the rate of sleeping at least seven hours on weeknights by an estimated 12 percentage points ($p < 0.001$), 28 percent higher than baseline. The treatment effects persist into the post-intervention period but are smaller: an estimated 5.5 percentage points ($p < 0.001$), a 13 percent increase compared to baseline. We estimate that average weeknight sleep increases by 19 minutes during the treatment period and 9 minutes during the post-treatment period ($p < 0.001$).²

¹The primary *Immediate Incentives* group pools two sub-treatments that received cues, feedback, and immediate incentives: one that continued to receive reminders and feedback in the post-treatment period and a secondary treatment group that did not receive reminders and feedback in the post-treatment period, which allows us to examine the importance of providing context cues for the persistence of behavior after the reward is removed. We do not find significant differences in the post-treatment effects of the two groups and pool them for our primary analysis.

²Our focus on weeknight sleep is in line with prior work that examines the impact of school and class start times, which occur on weekdays. We find no evidence of substitution between incentivized weeknight sleep and unincentivized sleep, including sleep that occurs during the day (i.e., naps), on weekends and during holidays.

To understand the extent to which our results reflect sustained changes in sleep habits, we examine treatment effects on sleep behaviors. As discussed above, we designed our intervention to establish the habit of earlier bedtimes, triggered by the nighttime reminders. We find that the intervention initially leads to earlier bedtimes (and directionally earlier wake up times) but these behaviors do not persist. On average, bedtime returns to baseline levels and wake-up time becomes slightly later. While our intervention does not establish early bedtime habits on average, we find that within individuals, bedtime and wake up time become more regular in both the treatment and post-treatment periods. These results suggest that the intervention led participants to establish more stable routines independent of the external cue to go to bed earlier.

We further explore mechanisms of habit formation by comparing our primary intervention to three secondary treatments: (1) *Delayed Incentives*, which is identical to Immediate Incentives but with the payout delayed to the end of the study (about a month after the treatment period); (2) *Delayed Incentives No Cue/feedback*, which is identical to Delayed Incentives except that participants do not receive cues or feedback; and, (3) *Cue/feedback*, which only provides reminders and feedback with no rewards. During the treatment period, the effects of Immediate Incentives are about 50 to 80 percent higher than Delayed Incentives and about three to four times larger than Cue/Feedback alone. During the post-treatment period, the estimated effects of Immediate Incentives are generally larger than the secondary treatments but are not statistically distinguishable.³ Our results suggest that combining cues with immediate rewards has large impacts while incentives are being offered but may not be particularly effective at enhancing the persistence of habits.

We then turn to the educational impact of incentives to sleep. Immediate Incentives increase semester course performance by an estimated 0.075 - 0.088 grade points ($p = 0.044$ and $p = 0.035$, respectively). We find evidence of similar sized treatment effects on grade point average in the semester following the intervention, but no impact two semesters after the intervention. We examine heterogeneity by time of the day and, in exploratory analysis, course subject. We find that treatment effects are largest in classes that take place midday and in STEM courses.

We benchmark our effects in comparison to prior work linking sleep to academic performance. Estimates from natural experiments in the U.S. suggest that a one hour later shift in sunrise or class start time increases sleep by an average of 6 - 36 minutes and has either no discernible impact on academic achievement or can increase grades and test scores by 0.06 - 0.16 standard deviations (SD) (Carrell et al., 2011; Heissel and Norris, 2018; Groen

³We note that, for our primary outcome, the post-intervention impact of Immediate Incentives is similar to Delayed Incentives No Cue/Feedback, which removed both the cue and the immediacy of incentives.

and Pabilonia, 2019). By comparison, our intervention increases weeknight sleep by an estimated 19 minutes during treatment and 9 minutes during post-treatment; and, grades by 0.10 - 0.11 SD.

We consider three primary channels through which sleep could affect academic performance: lifestyle, cognition, and well being. To that end, we examine the impact of our intervention on time use, performance in math and creativity tasks, and measures of physical and mental well-being. The intervention leads to declines in self-reported screen time, which includes internet browsing, TV/videos and games, and excludes screen time for studying. The changes in screen time are similar in magnitude to the increases in sleep and are concentrated around bedtime. We find that total study time does not change during the intervention period, but there is suggestive evidence of a reallocation of study time from evening hours to morning hours, with little change in other time use. We do not find treatment impacts on our measures of cognitive performance (math and creativity). We also do not find treatment effects on physical activity, or end of semester mental health, though we find evidence that treated participants report they are better able to cope with stress during the intervention period. Together, our results suggest that incentives to sleep lead to more regular sleep habits, which displace screen time, and shift study time to earlier hours of the day, which may contribute to the improvement in academic performance.

Our study is the first to show that an intervention targeting sleep can improve academic performance. These findings contribute to the growing literature on the economics of sleep. Seminal work on sleep finds a negative association in time use surveys between sleep and work hours, but is not able to identify causality (Biddle and Hamermesh, 1990; Basner et al., 2007). Related work finds a positive correlation between sleep and health outcomes (Cappuccio et al., 2010). Other studies find that sleep deprivation affects decision making, ethical behavior, social decisions, and voting behavior (Dickinson and McElroy, 2017; Dickinson and Masclet, 2023; McKenna et al., 2007; Holbein et al., 2019).

Studies using naturally occurring data find that later sunset times are correlated with lower cognitive performance and earnings, as well as worse physical and mental health, arguing that the channel is via reduced sleep (Giuntella et al., 2017; Gibson and Shrader, 2018; Giuntella and Mazzonna, 2019; Jin and Ziebarth, 2020). Related work finds that earlier class times are correlated with less sleep and lower academic performance at both the K - 12 and post-secondary levels (Carrell et al., 2011; Heissel and Norris, 2018; Jagnani, 2021; Groen and Pabilonia, 2019). In addition, recent research among U.S. university students finds a strong positive correlation between freshmen’s academic performance and sleep, particularly in the first half of the term (Creswell et al., 2023). However, Lusher et al. (2019) find no evidence of significant effects of sleep regularity on academic outcomes in a large study at

a Vietnamese university and little impact of delayed start times. None of these studies exogenously vary sleep or test policies to improve sleep habits.

Recent work using field experiments has tested the impact of interventions targeting sleep. [Bessone et al. \(2021\)](#) implemented a randomized field experiment in India that increased sleep through sleep aids and incentives. They find that their intervention increases nighttime sleep by an average of 27 minutes over a twenty day treatment period. But there is no meaningful impact on cognition or productivity. That increased sleep did not improve labor market performance could be due in part to the study’s context: as noted above, participants were severely sleep deprived at baseline and the environmental conditions led to poor sleep quality. Consistent with this, they find evidence that high quality sleep via office naps can increase productivity. The study cannot disentangle whether the differential effects of naps compared to nighttime sleep are due to the differences in the timing of sleep or the quality. In our study, as in the U.S. more broadly, average sleep is of high quality and baseline sleep is closer to the recommended minimum of seven hours a night: 6.6 hours among our participants vs. 5.5 hours in [Bessone et al. \(2021\)](#). Our study also takes place over a longer time horizon – five to nine weeks compared to twenty days – which may better allow the effects of moderate sustained increases in sleep to emerge.

Two additional field experiments have examined interventions to improve sleep but have not linked those impacts to productivity or performance. [Avery et al. \(2022\)](#) show that incentivizing early bedtimes and longer sleep duration increases overall sleep, and identify a demand for commitment devices to improve sleep habits. [Barnes et al. \(2017\)](#) test the effect of treating insomnia with internet-based cognitive behavior therapy and find beneficial effects on negative affect, job satisfaction, and self-control. In related work, [Breig et al. \(2020\)](#) conducted a field experiment with Fitbits to test the role of optimism bias in explaining time allocation and bedtime.

Our findings also contribute to the large literature on improving academic performance, particularly among college students (e.g., [Angrist et al., 2014](#); [Lavecchia et al., 2016](#), provide reviews). Our intervention is highly cost effective compared to previously examined policies, including financial aid, mentoring and support services, and performance-based incentives. Our results suggest that it may be more cost effective to improve academic achievement via sleep rather than to incentivize performance directly. This finding is akin to recent work showing that incentives for exercise can improve educational achievement ([Cappelen et al., 2017](#)).

Lastly, we contribute to the literature on habit formation by examining the impact of interventions that pair cues with feedback and rewards. Prior work has largely examined these separately. For example, [Wellsjo \(2021\)](#) focuses on the role of cues for generating

automatic habits in the context of handwashing; and, [Byrne et al. \(2022\)](#) test the impact of repeated feedback on sustained behavior change in the context of water conservation. A large prior literature offers rewards for repeated engagement in desirable behaviors in the context of smoking, weight loss, exercising, and handwashing, always distributing rewards with a delay from the time of the incentivized behavior ([Gneezy et al., 2011](#); [Royer et al., 2015](#); [Hussam et al., 2022](#); [Beshears et al., 2021](#); [Milkman et al., 2014](#)). We show that making rewards immediate significantly increases their impact on repeated behaviors. Our results add to prior findings that immediate rewards outperform delayed rewards in the context of one-time behaviors ([Levitt et al., 2016](#)). Finally, our examination of habit formation suggests that prescribed bedtime cues did not meaningfully increase the persistence of habit. These results are in line with prior work showing that interventions encouraging set exercise routines do not facilitate habit formation and are less effective than those that allow for individual flexibility ([Beshears et al., 2021](#)).

In the remainder of the paper, we describe the experimental design, data and analysis in Section 2. Section 3 presents the results, Section 4 benchmarks the results relative to prior findings, and Section 5 concludes.

2 Experimental design and data

We conducted our experiment in seven semester-long waves from Spring 2019 - Spring 2022 among students at the University of Pittsburgh (Pitt). We measured sleep using wearable trackers and delivered our interventions targeting sleep via text messages and a custom smartphone app. Our outcome data come from the wearable trackers (for sleep and physical activity), survey measures (time use, cognitive performance and well-being); and administrative records (academic transcripts).

2.1 Wearable trackers and custom smartphone app

To gather objective measures of sleep in a natural setting, we had participants in our study wear Fitbits that estimate sleep patterns based on movement and heart rate data. The use of such wearable trackers allowed us to depart from dependence on sleep diary methods, which have been shown to significantly overestimate sleep. ([Lauderdale et al., 2008](#); [Bessone et al., 2021](#)). Fitbits, which are among the most popular wearable trackers, are well-suited for monitoring sleep in natural settings due to their portability and unobtrusiveness, and are the most utilized wearables for biomedical research purposes ([Wright et al., 2017](#)). In our study, we used Fitbit Charge HR, Charge HR 2, Charge HR 3, Alta HR, and Inspire 2, which

all capture both movement and heart rate. Recent studies have demonstrated the accuracy of these heart rate-enabled Fitbits compared to actigraphy, a commonly used method for outpatient sleep screening, suggesting their suitability for population-based sleep research (Haghighayegh et al., 2019).

One source of concern with studies that rely on wearable trackers is that the devices require continued engagement via daily syncing (i.e., regularly connecting the tracker to a smartphone to update the collected data). To ensure high sync levels for our experiment we developed a custom-made smartphone app that connected to the Fitbit API, which allowed us to monitor sync rates daily and notify participants with low sync rates to keep them engaged. The custom-made app also allowed us to deliver our interventions to improve sleep habits via push notifications and the app itself. The app features include the ability to send reminders; provide immediate individualized feedback based on participants' sleep, as measured by the Fitbit; and redeem rewards. We discuss the interventions in more detail in Section 2.3

2.2 Sample, recruitment and timeline of the experiment

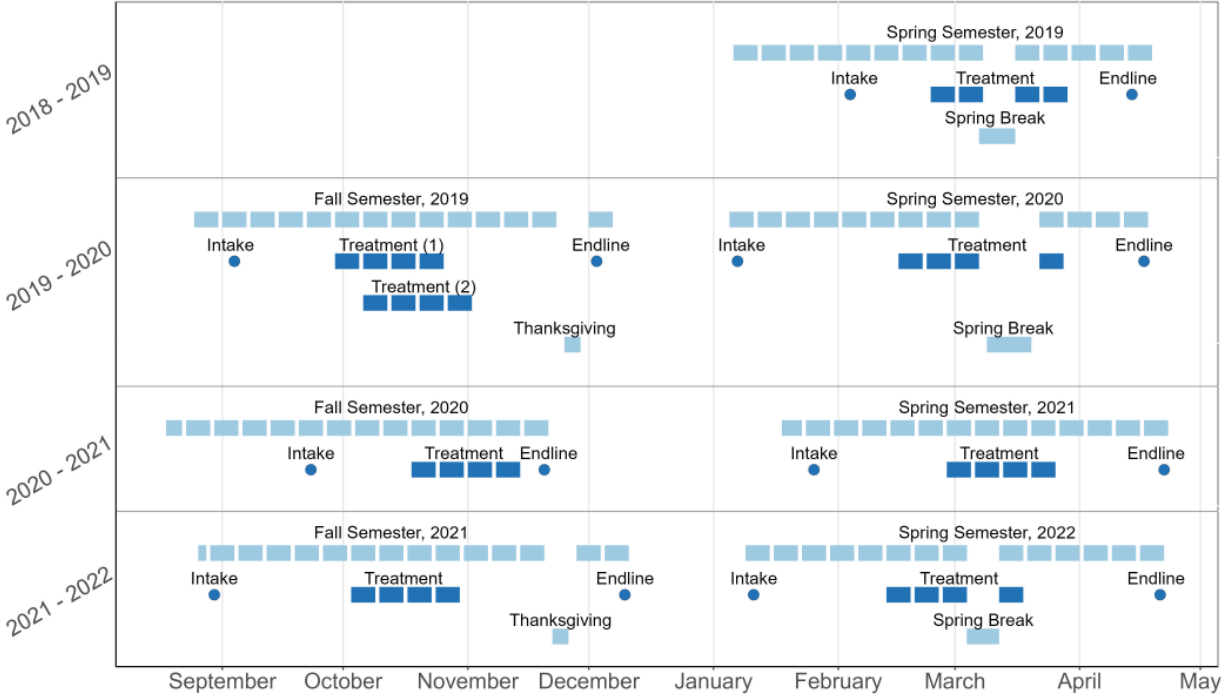
The experiment took place at the University of Pittsburgh, was approved by the University of Pittsburgh Institutional Review Board and was pre-registered in the AEA registry (RCT ID AEARCTR-0003235).

We recruited participants through the Pittsburgh Experimental Economics Laboratory (PEEL) and invited them to participate in a semester-long study on wellness for a guaranteed minimum payment of \$30 and the opportunity to receive additional earnings based on luck as well as their choices during the experiment. To be eligible for our study, participants had to have a smartphone and be willing to wear and routinely synchronize a wearable device (Fitbit) during the semester. We began the experiment in Spring 2019 and enrolled participants every semester (Fall and Spring) until Spring 2022. We ran the experiment in seven consecutive waves, with modest-sized cohorts to accommodate the number of participants we could recruit through the lab every semester as well as the number of Fitbits we had. Our final sample includes 1,149 participants.

In each wave, the study lasted for approximately 10 weeks. We initiated participant recruitment in the first few weeks of the semester and enrolled participants in the experiment on a rolling basis. Upon enrollment, we measured baseline sleep for one to four weeks. At the end of the baseline period, we randomly assigned participants to either a control group or treatments designed to improve sleep habits, which lasted for 4 weeks (*intervention period*). After the 4-week intervention period, we continued to follow participants for an additional

1-5 weeks until the end of the semester (*post-intervention period*), at which point we asked them to return the Fitbit and fill out an endline survey. The study always ended during the last week of classes, before final exams. In the different waves, the start of the recruitment period depended on lab availability. The timing of the treatment period and the length of the post-treatment period varied in each wave depending on when we were able to start recruiting, how quickly we enrolled participants, and the semester schedule. The timeline of the experiment for each of the seven waves is depicted in Figure 1.⁴

Figure 1: Timeline of the experiment



To enroll in the study, participants completed an initial session at the laboratory (Spring 2019 - Spring 2020) or over Zoom (Fall 2020 - Spring 2022), completed an intake survey, received the Fitbit and installed our custom smartphone app.⁵ During the intake session,

⁴In Fall 2019, due to recruitment issues at PEEL, we recruited two groups of participants and had them start the intervention in a staggered way, as shown in Figure 1. In Spring 2020, the semester schedule was changed by the university closure prompted by the onset of the COVID-19 pandemic. Students enrolled in the study in Spring 2020 learned about the university moving to remote learning during spring break (mid-March 2020), and continued to stay enrolled in the study until the end of the semester. In the Appendix, we conduct sensitivity analyses that exclude the Spring 2020 wave.

⁵From Fall 2020 onwards, we had to adjust some of the intake procedures due to changes in the lab and university protocols during the COVID-19 pandemic. Instead of filling out one unique survey during the intake session, the survey was split into an enrollment survey that participants filled out at enrollment while on Zoom and a follow-up survey that was emailed to them a few days later. In Spring 2019-Spring 2020 and Fall 2021-Spring 2022, participants picked up the Fitbit from PEEL and received a \$6 payment. In Fall 2020 and Spring 2021, participants received the Fitbit via mail.

participants consented to wear and sync the Fitbit throughout the semester, answer weekly surveys, and grant us access to their academic records. They were informed about their right to withdraw from the study at any time with no penalty. To mitigate potential experimenter demand effects, we did not specifically disclose to participants that our interest was to study sleep behavior. Instead, we broadly explained that we were interested in wellness. Participants left the initial session with a one-page reminder outlining what was expected of them during the study and agreed to return the Fitbit at the end of the semester. The intake survey administered to participants collected information on socio-demographic characteristics and baseline measures of well-being.

Over the course of the study, participants in all treatments received reminders to sync their Fitbit via text message and the app (see Figure B.5). They also received weekly surveys that elicited time use, cognitive performance and well-being. We describe the survey measures in more detail below (Section 2.4).

2.3 Treatments

In total, 1,219 individuals completed an enrollment survey. In order to be randomized, participants had to have at least a day of Fitbit data in the baseline period. In total, we randomized 1,149 participants to treatments.⁶

Table 1: Treatments

	N	Waves	Treatment		Post-Treatment
			Reminders & Feedback	Rewards	Reminders & Feedback
Control	380	1-7	–	–	–
Immediate Incentives	468	1-7		Immediate	
Immediate Incentives, Post Cue/Feedback	356	1-7		Immediate	
Immediate Incentives, No Post Cue/Feedback	112	5,7		Immediate	–
Delayed Incentives	103	1-3		Delayed	
Delayed Incentives, No Cue/Feedback	97	1-3	–	Delayed	–
Cue/Feedback	101	1-3		–	
Total	1,149				

⁶We mistakenly assigned eight participants to treatments who did not have any baseline Fitbit data. We conduct sensitivity analyses that exclude these individuals.

Notes The table reports the number of participants enrolled in each of the treatments; whether rewards were immediate or delayed, and whether they received reminders and feedback during and after the intervention. Immediate Incentives pools Immediate Incentives, Post Cue/Feedback and Immediate Incentives, No Post Cue/Feedback.

After the baseline period, we randomized participants into treatment groups, which are displayed in Table 1.⁷ Participants in the *Control* group ($N = 380$, waves 1-7), received no intervention and continued to wear their wearable trackers and fill out surveys until the end of the semester. Participants in the treatment groups received interventions to improve sleep habits. In all treatments, we set the goal of sleeping at least seven hours per night by 9 am on weeknights (Sunday through Thursday). We established 9 am as a key constraint, based on previous studies emphasizing the significance of sleep timing and the alignment of biological rhythms with the environmental light-dark cycle (Roenneberg and Merrow, 2016). Additionally, we aimed to reduce the likelihood that our intervention would encourage skipping classes scheduled at 9 am. Notably, about 82% of the participants reported waking up before 9 am at baseline.

Drawing on the habit formation framework outlined earlier, our *Immediate Incentives* intervention (468 participants, waves 1-7) leverages cues, rewards, and repetition to establish persistent sleep habits. To provide participants with a consistent cue, we sent them reminders—both through the app and via text—to meet their target goal of sleeping seven hours per night by 9 am every weeknight (Sunday-Thursday). These reminders had two major components. We established a personalized target bedtime for each participant, an hour earlier than their usual baseline bedtime, based on their individual sleep patterns, and sent reminders to go to bed half an hour before this new goal time. Figure B.1 displays the bedtime reminder.⁸

Second, as the cue-based framework emphasizes the importance of a stable environment in triggering automatic behavior, we encouraged participants to engage in a specific bedtime behavior every weeknight before going to sleep. Participants selected their behavior from a menu of different options before the beginning of the intervention period. Examples included

⁷We made two deviations in the treatments from the pre-registered experimental design. First, our original plan included an incentive treatment where participants would receive a \$4.75 coupon for a breakfast treat at one of the University of Pittsburgh Einstein Co-fee locations. However, due to unforeseen logistical difficulties, we suspended this treatment after the first few weeks of the first wave, and exclude it from the following waves. Second, the COVID-19 pandemic prevented us from meeting the pre-registered sample size for our main treatments (Control and Immediate Incentives), as we incurred additional costs for mailing Fitbits to participants.

⁸The bedtime was set approximately an hour before participants' average baseline bedtime rounded to the nearest 30 minutes with a latest bedtime goal of 1 am (e.g., for participants with an average baseline bedtime of 12:12-12:14 am, we set a goal bedtime of 11 pm; for participants with an average baseline bedtime of 12:15-12:30 am, we set a goal bedtime of 11:30 pm; for participants with an average baseline bedtime of 2 am or later, we set a goal bedtime of 1 am). We delivered the personalized bedtime reminder via text message. In the app, we delivered a standard message encouraging participants to go to bed early enough to sleep seven hours by 9 am.

“Turn off your Phone”, “Turn on bedtime music”, “Turn off your computer”, “Turn on meditation app”.⁹

Next, to link sleeping behavior with a reward, we provided participants with immediate financial incentives upon meeting their sleep goal. Every weekday after 9 am, participants received feedback on whether they met their goal of sleeping at least seven hours via the app through push notifications and the app interface. Participants who met their goal received feedback about having achieved the goal and earned a \$4.75 reward, which they redeemed by clicking a button on the app (see Figure B.2).¹⁰ Participants in this treatment received a monetary reward of \$4.75 through a Venmo transfer on the same day.¹¹ Participants who fell short of the sleep target were given feedback indicating that they had not achieved their goal and had missed out on receiving the reward. This feedback also included a negatively-valenced emoji to convey the injunctive message that sleeping less than seven hours was discouraged (see e.g., Schultz et al., 2007), and encouragement to make another attempt to meet the sleep target. To encourage repetition of the incentivized sleep behavior, cues and rewards continued every weeknight and weekday of the four-week intervention period. In waves in which the treatment period spanned spring break, we paused the intervention during spring break.

At the end of the intervention period, we discontinued the financial rewards, notifying participants via text message. In our main variation of the Immediate Incentives treatment, *Immediate Incentives – Post Cue/Feedback* ($N = 3,56$, waves 1-7), we continued to send bedtime reminders (i.e., the cue) and morning feedback on whether they had achieved their sleep goal throughout the post-intervention period, which lasted through the end of classes. The feedback was identical to that of the intervention period, except that we removed mention of the financial reward, as displayed in Figure B.4). To examine the importance of maintaining the cue for habit persistence, we tested a variant of the Immediate Incentives treatment in which participants did not receive cues and feedback in the post-intervention period, *Immediate Incentives – No Post-Cue/Feedback* ($N = 1,12$, waves 5 and 7). Our primary analysis pools the two variants of Immediate Incentives. In waves in which the post-treatment period spanned Thanksgiving, we paused reminders and feedback during the week of Thanksgiving. Our primary analysis excludes holiday weeks (spring break and Thanksgiving).

⁹On the Friday before the beginning of the intervention period, participants received their intervention-related instructions. As part of these instructions, we asked participants to select a bedtime behavior to engage in before going to bed.

¹⁰Redemption rates were above 95% across waves.

¹¹For logistical reasons, the payment was received after 3 pm each day, which introduced a small delay between the performance of the behavior and the reward. However, the feedback about receiving a reward was provided as soon as participants synced the Fitbit after 9 am.

Secondary treatments. Following our pre-analysis plan, in the first three waves of the study, we implemented secondary treatments to examine the importance of cues and feedback, and the timing of financial incentives for habit formation.

In the *Delayed Incentives* treatment ($N = 1,03$, waves 1-3), we provided participants with cues, feedback and rewards, as in the Immediate Incentives treatment. The only difference was that, although feedback about receiving the incentive was immediate, payment was not. Participants learned each day whether they met their sleep goal, but only received a single transfer with the total payment at the end of the study period, one to five weeks after the intervention ended. Figure B.3 displays the feedback screens for the Delayed Incentives treatment. In the post-intervention period, participants continued to receive the bedtime cue and morning feedback, as in the Immediate Incentives treatment. This treatment aimed to test whether providing repeated immediate rewards increases their effectiveness during treatment compared to delayed rewards, as has been shown with one-time rewards (Levitt et al., 2016); and whether reinforcing behavior with immediate rewards enhances persistence of behavior after the reward is removed, more so than delayed rewards.

In the *Delayed Incentives – No Cue/Feedback* treatment ($N = 97$, waves 1-3), we removed the bedtime reminders (i.e., the cue) and the daily feedback about whether participants met their sleep goal but retained the financial incentive. At the start of the intervention period, we informed participants that they would receive 4.75 for every night they met the goal of sleeping at least seven hours by 9 am, with payment to be received via a Venmo transfer at the end of the semester. The participants did not receive reminders or feedback during the post-intervention period. This treatment is analogous to other work using financial incentives to create habits in the context of exercising (e.g., Charness and Gneezy, 2009; Royer et al., 2015) and aimed to test the importance of pairing cues and feedback with rewards.

To test whether financial rewards are critical for establishing habits, we also conducted an additional treatment, *Cue/Feedback* ($N = 1,01$, waves 1-3), where we removed the financial reward. Participants in this treatment received the same bedtime reminders as participants in the Immediate Incentives and Delayed Incentives treatments. They also received daily feedback via the app on whether they had achieved their sleep goal. Instead of providing participants with a reward, the feedback screen included a positively- or negatively-valenced emoji depending on whether participants had achieved their sleep goal – i.e., the same feedback that the Immediate and Delayed incentives groups received in the post-intervention period (see Figure B.4). The participants in the Cue/Feedback treatment continued to receive the same reminders and feedback throughout the post-intervention period.

2.4 Data

2.4.1 Sleep

Our primary pre-registered sleep outcome is the share of weeknights (Sunday—Thursday) participants sleep at least seven hours. Our secondary measures of sleep include sleep hours per night, sleeping seven hours per night and sleeping between seven and nine hours without restricting to weeknights, sleeping seven hours per night including naps, bedtime, wake up time, sleep regularity, and sleep quality as measured by the Fitbit. In exploratory analyses, we also analyze sleeping between seven and nine hours on weeknights and sleeping at least six hours on weeknights in order to better compare these outcomes to our primary outcome measure. As a secondary measure, we also measured self-reported sleep and desired sleep in the intake survey.

We define sleep regularity as the sleep variability across the week, as measured by the within-person standard deviation in the outcome of interest. We measure sleep quality in terms of efficiency, Rapid Eye Movement (REM) sleep and deep sleep. Efficiency measures the percentage of time in bed that an individual is asleep. REM sleep is the stage of sleep in which individuals dream, which stimulates areas of the brain essential to learning. During REM sleep, heart rate and blood pressure rise. Studies suggest that REM sleep plays a key role in memory consolidation, emotional processing, and brain development (Marks et al., 1995; Boyce et al., 2016). Deep sleep is the most restorative form of sleep. During deep sleep the heart rate and breathing rate are at their lowest and our body repairs tissue. Deep sleep is important for regulating glucose metabolism and has also been linked to cognitive function and memory (Zhang and Gruber, 2019; Leproult and Van Cauter, 2010). We caution that while, as described above, there is growing evidence on the performance of recent Fitbit models in accurately measuring sleep duration (de Zambotti et al., 2018), the accuracy and reliability of these devices in capturing sleep stages needs further validation. In particular, sleep trackers have acceptable sensitivity but poor specificity when compared with sleep stages obtained using polysomnography (PSG).

Sync rates throughout the study were relatively high. On average participants synced their devices for 88% of the days. As shown in Table A.2, there are higher sync rates during the intervention period among the Immediate Incentives group compared to the Control group (there are no significant differences in sync rates in the baseline and post-intervention periods). As noted above, 8 of the 1149 participants did not have any Fitbit data, including at baseline. Of the 1141 participants with baseline data, 23 did not report any data during treatment (2%) and 83 (7%) did not report any data in the post-intervention period. In all of our analyses, for the nights with missing data, we replace missing data with

an individual’s baseline average following the approach of [Bachireddy et al. \(2019\)](#). We also conduct sensitivity analyses that do not replace missing data and results do not meaningfully change (Table [A.4](#)).

One concern with using Fitbits is that participants might lend their Fitbit to someone else for several nights to manipulate the reward system. To mitigate this potential concern, we make use of the resting heart rate data collected by the Fitbit. Research indicates that although resting heart rate can vary greatly between different people, it usually remains fairly stable within the same individual over time ([Quer et al., 2020](#)). We do not find a significant link between unusual fluctuations in resting heart rate (deviating more than two standard deviations from baseline) and treatment assignment (p -value= 0.72).

2.4.2 Educational outcomes

Our primary pre-registered educational outcome is term Grade Point Average (GPA), measured using administrative data obtained on September 14, 2023. The Registrar’s Office at the University of Pittsburgh supplied us with course data for participants enrolled in our experiment, covering each semester of their enrollment at the university from Fall 2018 through Spring 2023.¹²

This dataset provides comprehensive course information for our experimental participants across three key periods. In addition to data from the term during which our intervention took place (our primary outcome measure), it includes information from both before and after the intervention. Specifically, we have data for the term immediately preceding the intervention if participants were enrolled, as well as their High School GPA. We also have data for at least two terms following the intervention. Access to this post-intervention data allows us to conduct exploratory analyses on the intervention’s longer-term impact, provided the student enrolled in graded classes during subsequent semesters.

We calculate term GPA based on all courses in which a student received a letter grade, A+ through F, converted to a 4-point scale (see Table [A.1](#) for the grading system). Term GPA is an average of the course grade points, weighted by the number of credits for each course. Our secondary pre-registered outcome measures include course completion and credits; we also collected exploratory measures of withdrawals, course failure and course pass rates. Eighty-eight percent of the courses in our data receive a letter grade on the 4-point scale. Our main analysis excludes courses with grades outside this scale. We include these courses

¹²Because the administrative data was obtained after the intervention was concluded, we did not know the distribution of baseline GPA or GPA by course type in advance of the study. At the time of the pre-registration, we also did not know how many semesters of data and which secondary administrative data on educational outcomes would eventually be made available to us at study completion.

when examining the likelihood of having any grade, number of credits earned, as well as in the exploratory analysis examining course withdrawal, failure, and pass rates.¹³

Of 1149 participants in our experiment, 1,128 have at least one course grade for the term of the intervention (98% of the sample). The 21 remaining participants had no available grades for the term of the intervention, but have academic records in other terms. Our analysis includes all available grades data for the relevant term. For 1,056 participants (92% of the sample) we have data on at least one course grade for one term after the intervention; for 969 participants (84% of the sample) we have at least one course grade for two terms after the intervention. From the Registrar data on participants' high school GPA, we could match 1,049 students (91% of the sample). For 74 of the 100 students with no high-school GPA, we have information on baseline GPA at the start of the term (cumulative GPA from all prior terms). This gives a total of 1,123 participants with baseline GPA (98% of our sample).

The match rates are similar if we limit the data to our primary analysis comparing the Control and Immediate Incentives groups. Of these 848 participants, we match 833 to course grades in the term of the intervention (98%), 784 to high school GPA (92%) and an additional 41 to a prior term GPA. We construct a baseline GPA variable from either prior semester GPA or, when not available, high school GPA. A total of 825 participants have a baseline GPA (97%). As shown in Table A.2, we have a higher proportion of baseline grades for the Immediate Incentives group than for the Control group. In sensitivity analysis, we limit the sample to participants with baseline GPA and results are similar (Table A.7). There is no difference between the groups in the likelihood of having course grades, which is our primary outcome measure (Table A.2).

The course data allow us to classify courses by class type and start time. Class types include lectures, seminars, credit laboratories, practicum, workshops, independent studies, directed studies, internships, and laboratories. Lectures comprise 80% of the classes. As shown in Figure A.1, non-lecture classes have significantly higher grades and lower variance than lecture classes. The average GPA (and standard deviation) in lectures is 3.44 (0.81) compared to 3.75 (0.51) in other classes. In lecture classes, 47% of students receive the highest possible grade and the lowest quartile is a B. By comparison, in non-lecture classes, 67% of students receive the highest possible grade and the lowest quartile is an A-. This raises concerns that the grading system in non-lecture courses leaves little scope for treatment effects. In our analyses, we therefore report estimated treatment effects for all course types together, as pre-registered, as well as for lectures alone (we report the effects for non-lecture

¹³Our pre-registration also included major, attainment, and academic behaviors if the data were available. We do not have these data available, but we have self-reported information on major.

courses in Table A.7).

In exploratory analysis, we also classify each course as STEM or non-STEM using the Department of Homeland Security 2023 list of STEM designated CIP codes.¹⁴ The average GPA (and standard deviation) in STEM classes is 3.26 (0.88), while it is 3.69 (.58) in non-STEM classes. In STEM classes, 39% of students receive the highest possible grade and the lowest quartile is a B. By comparison, in non-STEM classes, 61% of students receive the highest possible grade and the lowest quartile is an A-. Similar to non-lecture courses, the grading in non-STEM courses may limit the scope for treatment effects.

2.4.3 Additional outcomes

Time use. We implemented a time use survey once a week, rotating the weekday on which the survey was administered. Our time use measure follows the structure of American Time Use Survey (Abraham and Flood, 2009). From a drop-down menu, participants indicated how they allocated their time on the previous day. For each 30 minute interval over the course of 24 hours, participants could choose from a list of activities that included sleeping, grooming (self), watching TV/videos, surfing the internet, playing games, working, studying, preparing meals or snacks, eating or drinking, cleaning, laundry, grocery shopping, attending religious services, hanging out with friends, paying bills, exercising, commuting, or other activities. They could also indicate that they did not know or could not remember how they spent their time, or could refuse to respond. In our primary analysis, we examine “screen” time, which pools time spent watching TV/videos, surfing the internet and playing games and excludes screen time spent studying; we categorize time spent hanging out with friends as “social” time. We exclude from the analysis responses that report 24 hours of “other activities”, which may reflect inattention in filling out the time use survey.

Cognitive performance. We collected secondary measures of cognitive performance through math and creativity questions. We drew the math questions from the math section of the Graduate Record Examination (GRE) test. We measured creativity using an adapted version of the task employed by Charness and Grieco (2019), where we provided participants with a list of 10 words and asked them to use some or all of the words to write an interesting sentence. On alternate weeks, the weekly time use survey included either one multiple choice math question or one creativity task. Both tasks were incentivized (see instructions in Appendix B). To assess the creativity task, we recruited raters from lab participants at the University of California San Diego and from Prolific ($N = 1,369$), and four undergraduate

¹⁴<https://nces.ed.gov/ipeds/cipcode/Files/2023/Final-2023-CIP-STEM-List-Blog.pdf>

research assistants at the PEEL lab at the University of Pittsburgh. Raters received a random subset of the sentences produced by participants in the creativity task and rated them on a 1-5 scale. Each sentence was rated by a minimum of two raters; the median number of ratings per sentence is thirteen.

Physical health. From the Fitbits, we collected data on resting heart rate and physical activity (daily steps and active minutes), which we pre-registered as primary health outcomes. Resting heart rate measures heart beats per minute (BPM) at rest, i.e. when sitting, lying down or relaxing. Faster resting heart rates are associated with shorter life expectancy (Cooney et al., 2010; Dyer et al., 1980). Daily steps are the number of steps over the course of a 24 hour-period. Active minutes are measured as minutes in which a person is non-sedentary for a least 10 continuous minutes, where non-sedentary is defined as activity that raises heart rate enough to burn at least 3 metabolic equivalents (METs).¹⁵

Well-being. We collected measures of mental health in the intake survey (conducted upon enrollment) and in the endline survey at the end of the semester using two pre-registered instruments. We assessed depression using the Center for Epidemiologic Studies Depression scale (CES-D, Radlo , 1977), which is a 20-item validated instrument designed to assess the frequency of depressive symptoms on a scale from 0 (“Rarely or None of the Time”) to 3 (“Most or Almost All the time”). An overall depression score is calculated by summing answers to all 20 items, with total scores ranging from 0-60. We also measured anxiety using the Generalized Anxiety Disorder scale (GAD-7, Williams, 2014), a 7-item scale designed to assess symptoms of Generalized anxiety disorder. The instrument assesses the frequency of anxiety-related symptoms using a scale ranging from 0 (“Not at all”) to 3 (“Nearly every day”), with total scores ranging from 0-21. To measure well-being, we collected exploratory measures of mood, stress and ability to cope with stress (resilience). For mood, we asked participants to indicate, on a 10-point Likert scale, how happy they felt in that moment. For stress and resilience, participants indicated, using a 5-point Likert scale, 1) the extent to which they faced stress in their life at the time of answering the survey and 2) the extent to which they felt able to deal with the stress they were facing. Every week, we alternated between the mood and the stress/resilience questions. These measures were collected via text

¹⁵In practice, this measure sums the lightly active, fairly active and very active minutes collected by the Fitbit. Our pre-registered secondary measures of health also include Body Mass Index (BMI) and blood pressure, which we ended up not collecting due to logistical constraints. We also pre-registered self-reported health behaviors as a secondary health measure, as well as a measure of willingness to pay to continue the intervention. We plan to analyze these measures, as well as other exploratory measures of health and lifestyle behaviors in a different paper.

message and, each week, participants were randomly assigned to receive the text message at different times of the day (11 am, 4 pm, 9 pm).

As shown in Table A.2, there is no difference in attrition rates between the Immediate Incentives group and the Control group for the additional outcomes discussed above.

2.5 Randomization and baseline characteristics

The randomization occurred at the end of the baseline period, the weekend before the start of the intervention-period. We employed a block randomized design, stratifying our participants by gender and the share of weeknights participants slept more than seven hours (above vs below median).¹⁶ In the initial waves of the study (Spring 2019 to Spring 2020) we randomized participants to one of five groups with equal probability (Control, Immediate Incentives, Delayed Incentives, Delayed Incentives No Cue/Feedback and Cue/Feedback). For the remaining waves, we randomly assigned participants to either the Control group or the Immediate Incentives treatment. In Waves 5 and 7 (Spring 2021 and 2022), we randomized participants in the Immediate Incentives treatment to either receive or not receive cue and feedback during the post-intervention period (Immediate Incentives - Post Cue/Feedback or Immediate Incentives - No Post Cue/Feedback).

Table 2 compares baseline characteristics in the Control group (column 1) to the Immediate Incentives group (column 2). We report demographic characteristics, baseline sleep behaviors and baseline academic characteristics (we discuss the baseline sleep in Section 3.1).

Students in the Control group are on average about 19 years old, with a large share of freshmen (52%). Sophomore, junior, and senior and above students make up 12, 23, and 12% of the Control group, respectively. Female and Asian students are over-represented compared to the full-time Pitt student population, and the U.S. college population in general. Approximately 56% (58%) of Pitt (U.S.) students are women, while women make up 72% of the Control group. Asian students make up 11% (7%) of the Pitt (U.S.) student population, while they represent 28% of the Control group. White students, which make up 56% of the Control group, are slightly under-represented compared to the Pitt student population (68%) and slightly over-represented compared to the U.S. college population (52%). The share of Black (8.8%) and Hispanic (4.0%) students is representative of the Pitt student population but lower than the U.S. college population, in which 13% of students are Black and 22% are Hispanic.¹⁷ A quarter of the students in the Control group report their parents

¹⁶We did not balance the randomization on baseline GPA because, as discussed above, GPA data was not available at the time of randomization.

¹⁷Demographics for the 2021-22 U.S. college population are available at:<https://www.statista.com/statistics/236360/undergraduate-enrollment-in-us-by-gender>, accessed on November 18 2023. Demographics for the Pitt student population in 2021-22 are available at: <https://www.ir.pitt.edu/sites/>

Table 2: Treatment-Control differences in baseline characteristics, Immediate Incentives

Variable	Control	Immediate Incentives	Difference
<i>Demographics</i>			
Female	0.721 (0.449)	0.726 (0.446)	-0.002 (0.031)
Age	19.463 (2.982)	19.344 (1.964)	-0.110 (0.170)
White	0.548 (0.498)	0.568 (0.496)	0.023 (0.034)
Asian	0.285 (0.452)	0.261 (0.439)	-0.017 (0.031)
Black	0.088 (0.283)	0.068 (0.253)	-0.021 (0.019)
Hispanic	0.040 (0.196)	0.053 (0.225)	0.010 (0.015)
Other	0.040 (0.196)	0.049 (0.216)	0.005 (0.014)
Highest parent educ:			
less than college	0.255 (0.437)	0.284 (0.452)	0.028 (0.031)
college	0.287 (0.453)	0.288 (0.454)	0.001 (0.032)
more than college	0.447 (0.498)	0.427 (0.495)	-0.018 (0.034)
<i>Baseline sleep outcomes</i>			
Sleep hours	6.625 (0.958)	6.659 (0.902)	0.012 (0.064)
Sleep ≥ 7 hours	0.438 (0.276)	0.426 (0.274)	-0.020 (0.019)
Sleep ≥ 6 hours	0.706 (0.258)	0.713 (0.258)	0.001 (0.018)
Bedtime	25.265 (1.313)	25.211 (1.297)	-0.073 (0.091)
Wake up time	7.956 (1.302)	7.935 (1.238)	-0.062 (0.086)
<i>Baseline academic characteristics</i>			
Freshman	0.521 (0.500)	0.530 (0.500)	-0.006 (0.031)
Sophomore	0.118 (0.324)	0.120 (0.325)	0.009 (0.022)
Junior	0.226 (0.419)	0.212 (0.409)	-0.012 (0.028)
Senior and above	0.124 (0.330)	0.139 (0.346)	0.020 (0.023)
STEM major	0.582 (0.494)	0.571 (0.496)	-0.004 (0.035)
Number of courses	5.167 (1.282)	5.158 (1.420)	-0.038 (0.095)
Number of early sessions	1.523 (1.562)	1.667 (1.545)	0.166 (0.109)
High-School GPA	4.140 (0.440)	4.131 (0.434)	-0.011 (0.032)
Baseline term GPA	3.429 (0.530)	3.457 (0.465)	0.016 (0.038)
Observations	380	468	848

Notes The sample is restricted to individuals Control and Immediate Incentives treatment group. Early class sessions are classes starting at 10 a.m. or earlier. All estimates in column 3 include wave fixed effects. Robust standard errors are in parenthesis. *** p 0.01, ** p 0.05, * p 0.1.

did not receive a college degree (either some or no college); 29% report that at least one of their parents has a college degree; and 45% report that at least one of their parents has a post-graduate degree.

About 58% of participants in the Control group report to be in STEM majors. They are enrolled in an average of 5.2 courses with an average of 1.5 class sessions per week beginning before 10 a.m. (early classes). The average high school GPA in our sample of 4.14 is representative of the overall University of Pittsburgh student population: the interquartile range of students offered admission at the University of Pittsburgh (Pitt) in 2022 had a weighted average GPA ranging from 3.91 to 4.42.¹⁸

In column 3, we estimate the treatment-control difference for each baseline characteristic from a regression that includes an indicator for the Immediate Incentives group and wave fixed effects. We do not find any statistically significant differences between average baseline characteristics in the Control group compared to the Immediate Incentives group. We also estimate Treatment-Control differences for each treatment group separately in Table A.3 and find statistically significant differences at the expected rate (e.g., about five percent of tests are significant at the $p < 0.05$ level).

2.6 Analysis

For outcome measures that are observed repeatedly throughout the study (e.g., nightly sleep), our primary regression analysis estimates treatment effects during the intervention period and the post-intervention period relative to the Control group. Formally, we estimate the following OLS model, unless otherwise noted:

$$Y_{it} = \beta_1 D_i + \beta_2 T_t + \beta_3 P_t + X_i + \rho_t + w_t + \mu_t + d_t + \epsilon_{it} \quad (1)$$

where Y_{it} is the outcome measure of interest; D_i is an indicator equal to one if an individual was assigned to the treatment group of interest; T_t is an indicator equal to one for any observation during the four-week intervention period; P_t is an indicator equal to one for any observation in the post-intervention period; X_i includes an individual’s baseline value of the outcome variable, baseline sleep (percent of weeknights slept at least seven hours),¹⁹ baseline GPA, indicators for the number of classes starting before 10 a.m. in a week (ranging from 0-5), and demographic controls for gender, age in years (dummies), race/ethnicity (Asian, Black, Hispanic, White, other), and indicators for parents’ highest education (less than

default/files/assets/CDS_2021-2022_Pittsburgh_20Campus_2.pdf, accessed on November 18 2023.

¹⁸Data available at: <https://admissions.pitt.edu/first-year-student/class-profile>, accessed on November 18 2023.

¹⁹We exclude baseline sleep in regressions for sleep outcomes due to collinearity with the the baseline value of the outcome variable.

college degree–high school degree only or some college–, college degree, or more than a college degree). For all individual characteristics, we included a missing indicator for whether the variable is missing. The variables ρ_t , w_t , t , d_t are a set of fixed effects for the wave of the experiment, week of the experiment, month of the year, and day of the week, respectively. Standard errors are clustered at the individual level.

For outcome measures that are observed only once during the study (e.g., course grades), we estimate the following OLS model, unless otherwise noted:

$$Y = \beta_1 D + X + \epsilon \tag{2}$$

where the variables are as described above. In regressions on course grades, the level of observation is the course weighted by the number of credits. Standard errors are clustered at the individual level.

Our main analysis compares the Control group to the Immediate Incentives treatment. As pre-registered, we also present the analysis for the primary outcomes comparing the Control group to the pooled incentives treatments (see Tables A.4 and A.7).²⁰

We report both unadjusted p -values and, as pre-registered, statistical significance adjusted for multiple hypothesis testing (MHT) within families of secondary measures. For that purpose, we use the method described in Anderson (2008), calculating Anderson False Discovery Rate (FDR) q -values and noting which estimates are robust to adjustment.²¹ We compute the FDR adjustment separately for treatment and post-treatment. In the Results section, we report unadjusted p -values and note which estimates are robust to adjustment.

In our main specifications, we include all participants who have outcome data. In the Appendix, we conduct sensitivity analyses that limit the sample to those who have both Fitbit and course grades data.

3 Results

We first examine treatment effects on sleep habits. We then turn to the impact of our intervention on educational outcomes. Finally, to explore potential mechanisms for our effects, we analyze time use, cognitive performance and physical activity and mental well-being.

² In a deviation from the pre-registered analysis we do not include instrumental variables (IV) analysis for GPA, instrumented with sleep. Our intervention may affect GPA through channels other than sleep—such as time allocation to other activities—and thus the IV exclusion restriction may be violated.

²¹Note that adjusted q -values can be both larger or smaller than unadjusted p -values. This is because, as noted by Anderson (2008), sharpened FDR q -values can be less than unadjusted p -values when many hypotheses are rejected.

3.1 Sleep

Baseline. Data from the baseline (pre-intervention) period reveals that a considerable portion of college students in our sample are sleep-deprived. As shown in Table 2, on weeknights, participants sleep an average of 6.6 hours, meet the recommendation of sleeping at least seven hours on approximately 43% of the nights; and sleep less than 6 hours on approximately 28% of the weeknights. About half of our participants have an average bedtime after 1 am, and about a quarter go to bed after 2 am on average. These data suggest that sleep deprivation is prevalent in our sample and is in line with a recent report by the National Institutes of Health indicating that more than 70% of college students sleep less than eight hours a day (Hershner and Chervin, 2014). In our sample, participants sleep less than 8 hours on approximately 84% of the weeknights.

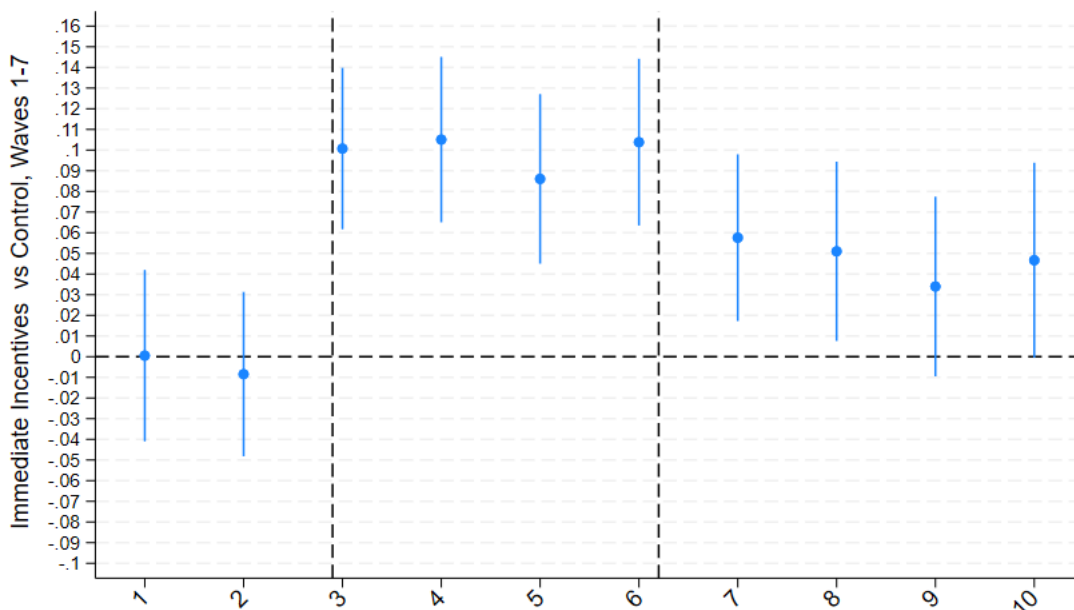
In the enrollment survey, 69% of the participants report sleeping at least seven hours on a typical weeknight (vs. 43% of nights as measured by the Fitbit). Compared to their (over-estimated) self-reported sleep, 81% of the participants state a longer optimal sleep duration during the week. On average, participants' stated optimal sleep time is an hour longer than what they report as their typical sleep duration. Ninety-seven percent of the students in our sample report that their optimal sleep on a weeknight would be at least seven hours, and 81% report an optimal sleep time during the week longer than 8 hours. These results suggest there may be scope for interventions that help individuals increase their sleep, as they state they would like to.

Intervention and post-intervention period. We first estimate treatment effects on the primary measure of sleep, which we incentivized: sleeping at least seven hours on weeknights (Sunday - Thursday). This outcome variable only includes nighttime sleep, and excludes weekends, holiday and naps (defined as episodes of sleep that start between 7 am - 8 pm). Figure 2 plots the estimated difference in the rate of sleeping at least seven hours between the treatment (Immediate Incentives) and Control groups, by week. The estimates are from regressions by week in which individual-nights are the level of observation and we include an indicator for the treatment group with no additional covariates (the Control group is the omitted group). Standard errors are adjusted for clustering at the individual level (the bars in the figure indicate 95% confidence intervals).²² As shown in the figure, there are no differences at baseline (weeks 1-2). Treatment effects emerge in the first week of the intervention (week 3) and persist throughout the four-week treatment period (weeks 3-6). Treatment effects decline as soon as the intervention ends (week 7) but remain positive

²²We present the analogous figures for sleep hours in Figure A.2 and distributions of sleep hours in Figure A.3.

and fairly steady throughout the post-intervention period (weeks 7-10).

Figure 2: Immediate incentives and sleep ≥ 7 hrs (weeknights), excluding naps



Notes The sample is restricted to weeknights (Sunday-Thursday nights). On the horizontal axis we report week of the study: baseline (weeks 1-2), treatment (weeks 3-6), post-treatment (weeks 7-10). The coefficient reports the difference in the likelihood of sleeping at least 7hrs between individuals in the Immediate Incentives treatment and those in Control by week. Standard errors are clustered at the individual level. Bars indicate 95% confidence intervals.

In the first two columns of Table 3 Panel A, we present regression estimates of the treatment and post-treatment impacts of Immediate Incentives on sleeping at least seven hours on weeknights and weeknight sleep hours, following the specification described in equation 1. At baseline, participants meet the goal of sleeping at least seven hours on approximately 43% of the nights.²³ During the intervention period, Immediate Incentives increase the rate of sleeping at least seven hours by an estimated 12 percentage points, a 28% increase. The treatment effects persist into the post-intervention period but are about half the size: an estimated 5.5 percentage points, 13% higher than baseline.²⁴ We estimate

²³The baseline average reported in Table 3 is slightly different from that reported in Table 2 as we are pooling Immediate Incentives and Control. Furthermore, in Table 2 we calculate the baseline average at the individual level and in Table 3 we calculate it at the night level, and not all participants have the same number of nights in the baseline period due to rolling enrollment.

²⁴We conduct the following sensitivity analyses in Table A.4: limit the covariates to wave fixed effects, gender, baseline sleep and baseline GPA; limit the sample to participants who have term GPA, exclude missing nights rather than replacing missing data with individual baseline means, reweight the sample with respect to gender to make it representative of the gender composition of US college students, and exclude wave

that total sleep hours increase an estimated 19 minutes on average during the intervention period and an estimated nine minutes during the post-intervention period. All estimates are significant at the $p < 0.001$ level and are robust to adjusting for multiple hypothesis testing.

The effects at the mean reflect shifts throughout the distribution of sleep, as measured by sleep hours and share of nights sleeping at least seven hours (Figure A.3). In Table A.5, we estimate treatment effects by baseline quartile of sleep (share of nights sleep at least seven hours in panel A and sleep hours in panel B). We find similar effects across quartiles during the intervention period; and some evidence of larger post-intervention effects among participants with lower levels of sleep at baseline. These results suggest that our intervention has the most persistent impact on those who are most sleep deprived.

As shown in Table A.6 (Panel A), we do not find any evidence of substitution between incentivized weeknight sleep and unincentivized sleep during the day, on weekends or during holidays (spring break for the treatment period and Thanksgiving for the post-treatment period). If anything, we find small positive spillovers, with some evidence of an increase in the likelihood of sleeping at least seven hours during weekends in the post-intervention period.²⁵

In Panel B of Table A.6, we examine additional sleep outcomes. Similar to our main results, we find that our intervention significantly increases the share of nights participants sleep at least six hours and the share of nights they sleep 7-9 hours, with persistent but smaller impacts in the post-treatment period.²⁶ On our measures of sleep quality, we find small positive increases in minutes of REM sleep, no significant impact on minutes of deep sleep, and small marginally significant impacts on sleep efficiency. We note that, in our sample, baseline efficiency is high: participants are asleep an estimated 94% of the time they are in bed. By comparison, Bessone et al. (2021) estimate efficiency of 70% among their experimental participants in India.

3.1.1 Drivers of short and long-term sleep habits

Bedtime and wake up time. In the last two columns of Table 3 Panel A, we estimate treatment effects on bedtime and wake up time. During the intervention period, treated

3 (onset of COVID-19). The results do not change. We estimate treatment effects of 11.3-12.7 percentage points and post-treatment effects of 5.3-6.4 percentage points.

²⁵Including naps, holidays and weekends, we estimate the intervention increased the share of nights with at least seven hours of sleep by 6.9 percentage points in treatment and 4.1 percentage points in post-treatment, and increased total sleep hours by an estimated 12 minutes in the treatment period and 6 minutes in the post-treatment period ($p < 0.01$ for all estimates).

²⁶Sleeping less than six hours is a common metric of sleep deprivation (Hafner et al., 2017). The recommendation of sleeping seven to nine hours draws on studies that link excessive sleep duration to detrimental effects on health (Hirshkowitz et al., 2015; Jike et al., 2018).

Table 3: Immediate Incentives and sleep

	(1)	(2)	(3)	(4)
<i>Panel A:</i>				
<i>Daily level</i>				
	Sleep \geq 7 hrs	Sleep hours	Bedtime	Wake-up time
Treatment	0.1186*** (0.013)	0.3203*** (0.035)	-0.3146*** (0.036)	-0.0471 (0.034)
Post-Treatment	0.0551*** (0.015)	0.1420*** (0.036)	-0.0271 (0.037)	0.0712* (0.038)
Observations	46,989	46,989	46,989	46,989
Mean of dep. var.	0.429	6.647	25.22	7.939
Std. dev.	0.495	1.279	1.457	1.438
Number of individuals	840	840	840	840
<i>Panel B:</i>				
<i>Regularity</i>				
<i>(within-individual weekly s.d.)</i>				
		Sleep hours	Bedtime	Wake-up time
Treatment		-0.1084*** (0.032)	-0.0500** (0.023)	-0.0701** (0.031)
Post-Treatment		-0.0515 (0.035)	-0.0561** (0.025)	-0.0930*** (0.036)
Observations		8,631	8,631	8,631
Mean of dep. var.		1.171	0.896	0.937
Std. dev.		0.614	0.435	0.534
Number of individuals		840	840	840

Notes The sample is restricted to individuals in the Immediate Incentive treatment and individuals in the Control group. All estimates include day of the week, week of the experiment, wave, and month fixed effects, baseline value of the outcome variable, and demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), indicators for the number of classes starting at 10am or earlier, indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable is missing. In Panel A, observations are the dependent variable at the daily level. In Panel B, observations are the standard deviation of the dependent variable at the weekly level. Panel B includes all fixed effects and controls listed above, except the day of the week fixed effects. Standard errors are clustered at the individual level. Mean of dep. var. is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p 0.01, ** p 0.05, * p 0.1.

participants go to bed an estimated 19 minutes earlier than participants in the control group ($p < 0.001$) with directionally earlier wake up times ($p = 0.160$). This pattern does not persist in the post-treatment period when the incentives ended but the bedtime reminders continued. Instead, average bedtime largely reverts to baseline levels and treated participants wake up slightly later ($p = 0.065$). As shown in Figure A.2, both bedtime and wake up time get progressively later over the course of the intervention period and stabilize during the post-intervention period. These results suggest that combining bedtime reminders with incentives initially induces participants to go to bed earlier, but does not establish a sustained habit linked to the bedtime cue.

Sleep regularity. Panel B of Table 3 estimates treatment and post-treatment effects on sleep regularity. To do so, we examine the within-individual standard deviation of total sleep hours, bedtime and wake up time at the week level (the level of observation is individual-week).²⁷ We find significant decreases in variability of sleep hours across the week, equivalent to about 9 percent of baseline, or 0.18 standard deviations. The magnitude of the effects during the intervention and post-intervention periods are of similar size. These findings show that, while treated participants do not on average sustain earlier bedtimes after the intervention ends, they do develop more regular bedtime and wake up time habits. That the habits persist into the post-intervention period suggests treated participants found personal bedtimes and wake up times they were able to maintain. Such regularity may be important for cognition and performance. Prior work suggests that irregular sleep among college students is associated with delayed circadian rhythms and lower academic performance (Phillips et al., 2017; Trockel et al., 2000; Smarr, 2015).

Secondary treatments. As discussed in Section 2 and summarized in Table 1, our three pre-registered secondary treatments vary elements of our primary Immediate Incentives treatment in order to investigate the importance of cues and immediate rewards: (1) Delayed Incentives, which is identical to Immediate Incentives except that the rewards are distributed at the end of the study about a month after treatment; (2) Delayed incentives No Cue/Feedback, which is identical to Delayed Incentives except that participants do not receive cues or feedback; and, (3) Cue/Feedback which only provides cues (bedtime reminders) and feedback with no rewards.

Table 4 estimates the effects of our primary and secondary treatments on sleep hours and sleeping at least seven hours on weeknights. We restrict the analysis to waves 1-3 of the

²⁷The regressions follow the specification of equation 1, except we exclude day of the week fixed effects given the analysis is at the weekly level.

Table 4: Secondary treatments

	(1) Sleep ≥ 7	(2) Sleep hours
Treatment		
Immediate Incentives	0.1433*** (0.021)	0.3270*** (0.054)
Delayed Incentives	0.0945*** (0.022)	0.1781*** (0.058)
Delayed Incentives, No Cue/Feedback	0.0890*** (0.022)	0.1806*** (0.059)
Cue/Feedback Only	0.0364* (0.020)	0.1145** (0.052)
Post-Treatment		
Post: Immediate Incentives	0.0487* (0.025)	0.2083*** (0.067)
Post: Delayed Incentives	0.0143 (0.027)	0.0403 (0.071)
Post: Delayed Incentives, No Cue/Feedback	0.0551** (0.028)	0.0932 (0.076)
Post: Cue/Feedback Only	0.0144 (0.027)	0.0778 (0.073)
Observations	34,954	34,954
Mean of dep. var.	0.434	6.696
Std. dev.	0.496	1.536
Number of individuals	589	589

Notes The sample is restricted to waves 1-3. All estimates include day of the week, week of the experiment, wave, and month fixed effects, baseline value of the outcome variable, and demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), indicators for the number of classes starting at 10am or earlier, indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable is missing. Standard errors are clustered at the individual level. Mean of dep. var. is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p < 0.01, ** p < 0.05, * p < 0.1.

experiment when the secondary treatments were conducted. As shown in columns 1-2, the effects of Immediate Incentives during treatment are about 51% to 86% percent higher than the effect of Delayed Incentives (with or without cues and feedback); and about three to four times higher than the effects of Cue/Feedback alone. The differences between the estimated impact of Immediate Incentives and each of the secondary treatments are all significant at the $p < 0.05$ level after adjusting for multiple hypothesis testing. During the post-intervention period, the estimated effects of Immediate Incentives are generally larger than those of the secondary treatments. However, the effects are statistically indistinguishable. Further, the effects on sleeping at least seven hours a night are similar for Immediate Incentives and Delayed Incentives No Cue/Feedback.

We next focus on the role of the cue and feedback for developing and sustaining habits in combination with rewards. Comparing the two Delayed Incentives treatments, we find that the effects of Delayed Incentives are similar with or without cues and feedback. If anything, the Delayed Incentives No Cue/Feedback treatment has more persistent effects in the post-intervention period for our primary outcome measure. In Table A.4, column 8, we report an exploratory analysis where we estimate the post-treatment effects of the Immediate Incentives intervention separately for the subgroup of participants who continued to receive reminders and feedback in the post-treatment period (Immediate Incentive with Cue/Feedback) and for the subgroup of participants who stopped receiving them at the end of the intervention-period (Immediate Incentive No Cue/Feedback). We restrict the analysis to waves 5 and 7 in which we ran both variants. Our estimates reveal no significant differences between these two subgroups, suggesting that receiving bedtime cues after incentives stopped did not further help sustain the routines developed during the intervention period.

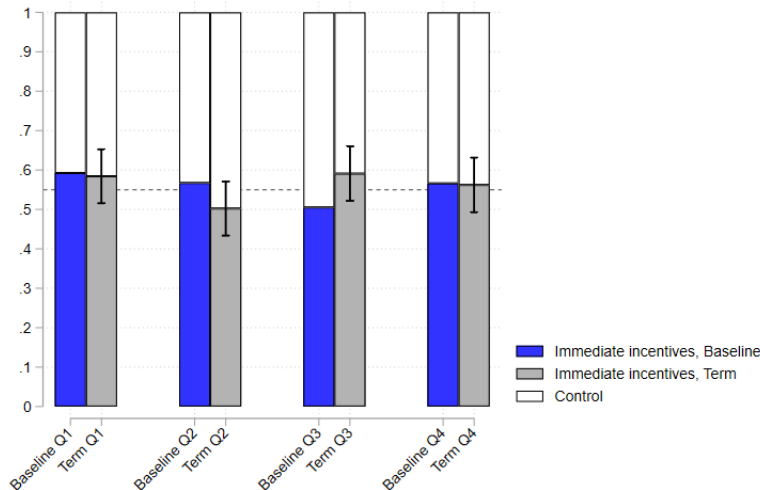
Collectively, the results presented in this section do not provide strong evidence that individuals built automatic habits as a result of our external cue. In the post-intervention period average bedtime reverted to baseline despite the bedtime cue. We also find little evidence that the external cue enhanced the impact of incentives during the intervention period or the persistence of habits during the post-intervention period. Nonetheless, the observed increase in sleep regularity (i.e., the reduced variability in sleep hours, bedtime and wake up time) persists in post-treatment. This suggests that the intervention facilitated the establishment of more dependable sleep routines, irrespective of the external cues provided.

3.2 Educational outcomes

Next, we investigate the impact of our primary treatment on educational outcomes. Figure 3 displays the share of individuals in the Immediate Incentives treatment and in the Control

group who are in each quartile of the GPA distribution at baseline and at the end of the intervention term.²⁸ The figure highlights that, as compared to baseline GPA, the share of treated participants below the median declines and share above the median increases. The effects are driven by shifts in the middle two quartiles with little change in the bottom and top quartiles.

Figure 3: Immediate Incentives and GPA



Notes The figure reports the share of individuals in each quartile of the GPA distribution for both baseline GPA (the high-school GPA) and the term GPA during the intervention for the Immediate Incentive treatment. Bars indicate 95% confidence intervals.

Table 5 presents regression estimates of the impact of the intervention on semester GPA and secondary educational outcomes.²⁹ The regressions follow the specification for equation 2. In columns 1- 2 we estimate treatment effects on our primary outcome, course grade, in the term the intervention took place. Column 1 includes all course types (lectures, seminars, labs, independent studies and other classes) whereas column 2 restricts the analysis to lectures (which account for approximately 80% of course types). Columns 3-4 report the same analysis for the term following the intervention. In columns 5 and 6, we examine the persistence of the effects two terms after the intervention.

As shown in columns 1 and 2, we estimate that Immediate Incentives improved average course performance by 0.075 grades points in all classes ($p = 0.044$) and 0.088 grade points in

²⁸We use high school GPA for baseline and only include participants with high school GPA in the figure.

²⁹As discussed in Section 2.6, we are missing GPA for 1.9 percent of our participants. We examine differential attrition on the GPA measure in Table A.2 and find no evidence for differential attrition on term GPA.

Table 5: GPA, Immediate Incentives

	(1)	(2)	(3)	(4)	(5)	(6)
	Term of intervention		Term +1		Term +2	
	All classes	Lectures	All classes	Lectures	All classes	Lectures
Immediate Incentives	0.075** (0.037)	0.088** (0.042)	0.068* (0.038)	0.091** (0.042)	0.004 (0.042)	0.004 (0.046)
Observations	4,300	3,413	4,087	3,298	3,842	3,080
Mean of dep. var.	3.502	3.436	3.553	3.494	3.547	3.505
Std. dev.	0.763	0.805	0.756	0.795	0.774	0.806
Number of individuals	833	827	784	782	727	718

Notes All estimates include demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), baseline sleep, indicators for the number of classes starting at 10am or earlier, indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. Observations are weighted by the number of credits taken in the semester. Standard errors are clustered at the individual level. Mean of dep. var is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p 0.01, ** p 0.05, * p 0.1.

lecture classes ($p = 0.035$).³⁰ The estimated GPA impacts of 0.075 - 0.088 grade points during the intervention semester are equivalent to a 0.10 - 0.11 standard deviation (SD) increase in grades. We estimate a treatment effect of similar magnitude on course performance in the semester following the intervention, providing suggestive evidence of persistent effects (columns 3 ($p = 0.066$) and 4 ($p = 0.027$)). However, we do not find treatment effects in the semester two terms after the intervention (columns 5 and 6).

We conduct the following sensitivity analyses in Table A.7: limit the covariates to wave fixed effects, gender, baseline sleep and baseline GPA; limit the sample to participants who have post-intervention term grades; limit the sample to participants who have baseline GPA (high school or baseline term GPA); limit the sample to participants who have sleep data at baseline; exclude wave 3 (onset of COVID-19); and reweight the sample with respect to gender to make it representative of the gender composition of US college students. The estimated effects are slightly smaller when we limit the covariates, exclude participants missing grades data or reweight the sample: 0.061-0.067 grades points for all classes (p -values ranging between 0.024 and 0.090) and 0.076-0.082 for lectures (p -values ranging between 0.019 and 0.076). Interestingly, when we exclude the Spring 2020 semester (wave 3), which experienced the onset of the COVID-19 pandemic, our estimated treatment effects are slightly higher,

³ As shown in Panel C of Table A.7, we find no treatment effects in non-lecture classes (i.e. seminar, labs, internships, directed studies). As discussed in Section 2.6, over two-thirds of students in these classes receive an A, and the lowest quartile is A-, leaving little room to improve grades (see Figure A.1).

an estimated 0.090-0.105 grade points. During this wave of the study, our participants experienced the abrupt closure of the university in the middle of the semester and disruptions in sleep and other lifestyle habits (Giuntella et al., 2021). Following our pre-registration, we also report results where we conduct our main analysis pooling all incentives treatments. The estimates are similar with slightly smaller average impacts, as shown in Table A.4 for sleep and Table A.7 for grades.

We also explore the effects of our intervention on other measures of course performance (Table A.9). This analysis includes the courses in our main analysis as well as courses that do not contribute to GPA because they do not have a grade on a four-point scale, such as pass, honors, and incomplete (see Table A.1 for the grading system). We find that students in the treatment group are less likely to receive any grade, which is primarily due to small increases in the likelihood of withdrawing from a class (column 2). At the same time, students in the treatment group are marginally less likely to fail a course (column 3). These results suggest that treated students are more likely to withdraw from classes they would otherwise fail. As a result, there are no significant differences in the likelihood of passing a class (column 4) nor in the number of credits completed in a term. We note that the effects on the likelihood of having any grade, withdrawing, or failing are not statistically significant after adjusting for multiple hypothesis testing.

Turning to heterogeneity, Table 6 shows that the results are not driven by performance in early-morning classes but rather the largest effects are in late morning/early afternoon classes that occur between 10 a.m. and 2 p.m., followed by afternoon/evening classes (after 2 pm). This is in line with the findings from Carrell et al. (2011) that early class start time affects performance in all classes, not just classes taking place early in the morning.

We additionally find evidence that our overall effects are driven by STEM courses: on average our intervention leads to a 0.132 grade point increase in STEM courses, which corresponds to a 0.15 SD increase in grades. By contrast, we estimate small increases of 0.013 grade points in non-STEM courses. This finding suggests that sleep may have a more significant impact on quantitative courses. However, given the large proportion of STEM majors in our sample (58%), another possibility is that improved sleep could enhance performance in courses important for a given major. The results may also partially reflect that there is more room to move grades in STEM courses, which have an average GPA of 3.27, compared to non-STEM courses with an average GPA of 3.70.

In exploratory analysis (Table A.8), we examine heterogeneous treatment effects on sleep and GPA by individual characteristics. We find larger effects among women, both on sleep and GPA. Effects on GPA are large among STEM majors, but effects on sleep are similar for STEM and non-STEM majors. Further, first-term freshmen students exhibit substantially

Table 6: Immediate Incentives and GPA: Heterogeneity by schedule and class type

	(1)	(2)	(3)	(4)	(5)	(6)
	Course grade	Class start: before 10am	Class start: 10am-2pm	Class start: after 2pm	Class type: non-STEM	Class type: STEM
Panel A: All classes						
Immediate Incentives	0.075** (0.037)	0.056 (0.064)	0.094** (0.044)	0.064 (0.055)	0.013 (0.034)	0.132** (0.057)
Observations	4,300	959	1,634	1,568	2,351	1,948
Mean of Dep. Var.	3.502	3.471	3.493	3.497	3.696	3.267
Std. dev.	0.763	0.810	0.744	0.773	0.574	0.888
Number of individuals	833	607	773	751	794	694
Panel B: Lectures						
Immediate Incentives	0.088** (0.042)	0.022 (0.072)	0.115** (0.048)	0.095 (0.065)	0.030 (0.040)	0.132** (0.059)
Observations	3,413	735	1,385	1,229	1,717	1,695
Mean of Dep. Var.	3.436	3.403	3.447	3.426	3.668	3.202
Std. dev.	0.805	0.859	0.774	0.815	0.598	0.912
Number of individuals	827	523	731	697	710	692

Notes All estimates include demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), baseline sleep, indicators for the number of classes starting before 10am, indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. Observations are weighted by the number of credits taken in the semester. Standard errors are clustered at the individual level. Mean of dep. var. is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p 0.01, ** p 0.05, * p 0.1.

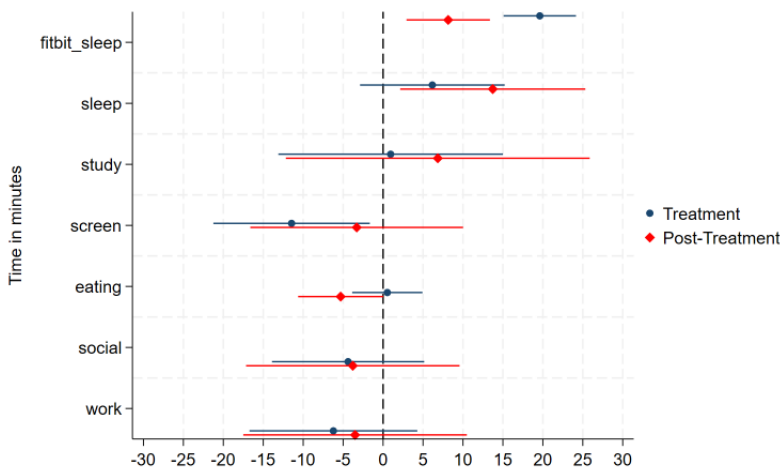
larger effects than other students. These findings are consistent with the idea that habits may be more malleable among freshman who have not fully developed their routines (Creswell et al., 2023).

3.3 Additional measures

To help make sense of our results, throughout the study we collected measures of time use, cognitive performance in math and creativity tasks, physical health via the Fitbit, and well-being.

Lifestyle. We next focus on our survey measures of time use that we asked weekly throughout the study (see Section 2 for details). Figure 4 shows estimated treatment and post-

Figure 4: Incentives to sleep and time use (minutes)

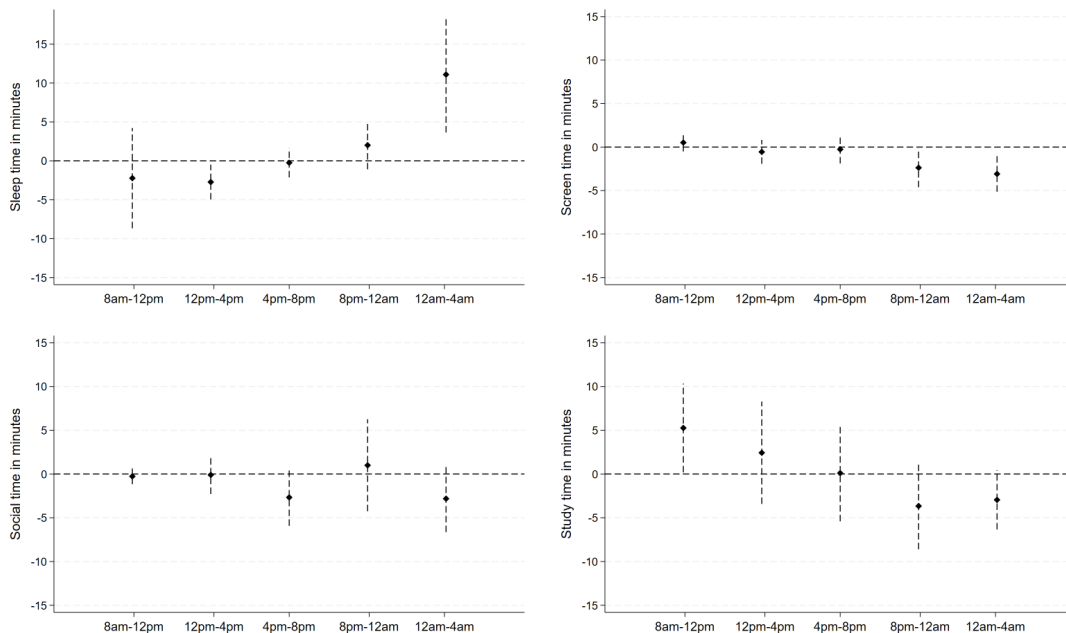


Notes The figure reports differences between the Immediate Incentives treatment and Control group in time-use during the intervention (in navy) and in the post-intervention period (in red). All the coefficients are obtained from regressions including wave, month, and day of the week fixed effects, baseline value of the outcome variable, and demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), indicators for the number of classes starting at 10 am or earlier, indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. Standard errors are clustered at the individual level. Bars indicate 95% confidence intervals.

treatment effects on time spent (in minutes) for the top six time use categories, using the regression specification in equation 1. In the first row of the figure, we show the estimated effects from the Fitbit data that we report in column 2 of Table 3. As discussed earlier, our intervention increased sleep by 19 minutes during treatment and 9 minutes after the removal of the incentives. Immediate Incentives directionally increase self-reported sleep in

both the intervention and post-intervention periods by about 6 - 14 minutes on average per day (Table A.10). We also find that subjects were 7 (6) percentage points more likely to report at least seven hours of sleep during the intervention (in the post-intervention) period.

Figure 5: Immediate Incentives to sleep and time use over the day: Intervention period



Notes The figure reports differences between participants in the Immediate Incentives treatment and Control group in the minutes allocated to different time-use activities during the intervention throughout the day. All the coefficients are obtained from regressions including wave and day of the week fixed effects, baseline value of the outcome variable, and demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), indicators for the number of classes starting at 10 am or earlier, indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. Standard errors are clustered at the individual level. Bars indicate 95% confidence intervals.

During the intervention period, incentives to sleep significantly decrease average screen time, which includes internet browsing, TV/videos and games, excluding screen time for studying, by an estimated 11.5 minutes per day ($p < 0.01$). We estimate smaller and not statistically significant treatment effects on screen time during the post-intervention period. These estimates are robust to adjusting for multiple hypothesis testing ($p < 0.01$ for share nights with less than seven hours of sleep; $p < 0.05$ for screen time). We do not find evidence of meaningful changes in time spent studying, socializing, eating or working during the intervention.³¹ In Table A.10 we report estimates on all the time use categories.

³¹We estimate treatment effects separately for internet, TV/videos and games in Table A.10 and the overall impact is largely driven by decreases in TV/video time. The table also reports effects on other time use categories. At baseline, we estimate the following average minutes per day for each category: sleep (494

In Figure 5, we report treatment effects on sleep, screen time, social time and study time over the course of the day during the intervention period. The effects on sleep and screen time are concentrated at night (8 pm - 4 am). Interestingly, while total study time does not increase, we observe a reallocation of study time from the evening/night (8 pm - 4 am) to the morning (8 am - 12 pm), although not precisely estimated. These results suggest that incentives to sleep led participants to develop sleep habits characterized by earlier screen disengagement at night and more focus on study time during the day. We also estimate treatment effects during the post-intervention period and find similar, but weaker, patterns (Figure A.4).

Cognitive performance. To examine cognitive performance directly, we collected measures of performance in math and creativity tasks on alternating weeks throughout the study. We do not find any impact of the intervention on these proxies for cognitive performance (see Table A.11, columns 1 and 2), which could be due to the intervention not affecting cognitive performance or to our measures not being able to capture the impact of performance on cognition.

Table 7: Immediate Incentives and well-being

	(1)	(2)	(3)	(4)	(5)
	Happiness	Stress	Resilience	CES-D	GAD-7
Treatment	-0.0655 (0.108)	0.0631 (0.059)	0.1503*** (0.055)		
Post-Treatment	-0.0569 (0.111)	-0.0411 (0.069)	0.0640 (0.063)	0.3997 (0.886)	0.0818 (0.404)
Observations	4,166	3,629	3,558	1,462	1,462
Mean of dep. var.	6.404	3.115	2.993	15.78	6.832
Std. dev.	1.646	1.116	0.997	10.23	4.864
Number of individuals	794	800	794	834	834

Notes Estimates in columns 1-3 include day of the week, wave, and month fixed effects, and demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), indicators for the number of classes starting at 10am or earlier, indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. For mood, participants indicated, on a 10-point Likert scale, how happy they felt in that moment. For stress and resilience, participants indicated, using a 5-point Likert scale, 1) the extent to which participants faced stress in their life at the time of answering the survey and 2) the extent to which they felt able to deal with the stress they were facing. For columns 4 and 5, outcomes are measured at endline, and estimates include all of the controls listed above, except for day of week, week of the experiment, and month fixed effects. CES-D is the Center for Epidemiologic Studies Depression Scale. GAD-7 is the General Anxiety Disorder-7. Standard errors are clustered at the individual level. Mean of dep. var. is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p < 0.01, ** p < 0.05, * p < 0.1.

minutes), study (321 minutes), screen (172 minutes), eating and preparing food (95 minutes), social (101 minutes), work (92 minutes).

Well-being and physical health.

Our final outcomes of interest are well-being and physical health. Previous work suggests that there is a positive relationship between sleep and both mental well-being and physical health (Giuntella and Mazzonna, 2019; Giuntella et al., 2017; Jin and Ziebarth, 2020).

To investigate the impact of the intervention on well-being, we sent participants weekly text messages to collect data on mood, stress, and resilience to stress. Additionally, we utilize the Generalized Anxiety Disorder (GAD-7) scale to assess anxiety levels and the Center for Epidemiologic Studies Depression (CES-D) scale to gauge depression levels. These scales were administered at baseline and endline only, so we are only able to estimate treatment effects on post-intervention end-of-semester anxiety and depression. Table 7 shows that the intervention does not have a significant impact on mood or stress levels (columns 1 and 2). However, it led to a statistically significant increase in resilience—participants’ self-reported ability to cope with stress—by approximately 0.15 standard deviations (column 3), which is significant at the 10% level after adjusting for multiple hypothesis testing. On the other hand, the intervention did not show any significant effects on post-treatment measures of depression and anxiety (columns 4 and 5), with point estimates being small in magnitude and lacking statistical significance.

As discussed in Section 2, we use the Fitbit to measure participants’ heart rate, daily steps, physical activity. We present estimates of treatment and post-treatment effects in Table A.11. We find no evidence of treatment effects on any of the physical health measures (columns 3, 4, and 5).

4 Benchmarking our results

We benchmark our results in two ways. First, we compare our effects to casual estimates of the relationship between sleep and academic performance from naturally occurring data. As discussed above, prior work examines the effect of shifts in sunset and school start times on sleep, grades and test scores. For example, Carrell et al. (2011) estimates that shifting the start time of college students’ first class by an hour from 7:00 am to 8:05 am improves overall academic performance by 0.12 - 0.14 SD. The study does not directly measure students’ sleep. Other studies using self-reported sleep estimate that an hour later school start time increases sleep by 36 minutes among American children, with a 0.16 SD improvement in reading and no change in math (Groen and Pablonia, 2019). Related studies find that the the sun rising one hour later increases average sleep among American children by an estimated six minutes, with a 0.082 SD increase in math scores and a 0.057 SD improvement in reading scores (Heissel and Norris, 2018). Taken together, these studies suggest that a one hour shift

increases sleep by 6 - 35 minutes and has either a null effect on academic performance or improves grades and test scores by 0.06 - 0.16 SD.³² Our impacts of a 19 minute average increase in weeknight sleep during treatment, a 9 minute average increase in post-treatment, and a 0.10 - 0.11 SD improvement in grades, falls within the range of the estimates in prior work on the causal relationship between shifts in sleep and changes in academic performance.

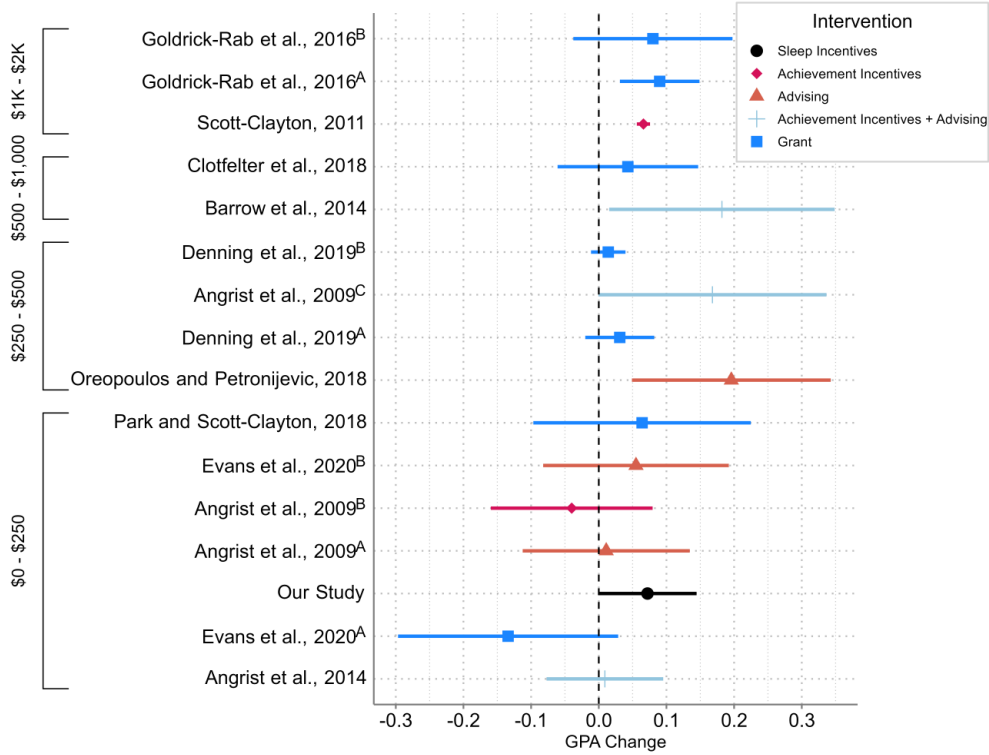
Second, we compare the cost effectiveness of our intervention to prior work examining policies aimed at improving college students' outcomes, including those that condition rewards on academic performance. Angrist et al. (2014) summarize the work on performance based incentives, including their own, and conclude that, "the picture that emerges. . . is one of mostly modest effects. . . [And there are] similarly discouraging results from studies of state-based merit aid programs. A few studies report positive effects, most notably Scott-Clayton (2011)'s evaluation of West Virginia PROMISE," which conditions free tuition on meeting a minimum GPA. Scott-Clayton (2011) finds similar-sized GPA effects to our study, 0.066 grade point improvements, but at over ten times the cost, an estimated \$1250 per student per semester. By comparison, we estimate that incentives to sleep increase semester GPA by 0.075 grade points and cost approximately \$110 per participant for the semester and would cost about \$160 per participant per year. This includes \$60 for the cost of the Fitbit and an estimated average of \$52 per participant per semester for the incentives (participants in the Immediate Incentives group received the incentives of \$4.75 per night on 55% of the 20 nights we offered it).

In Figure 6, we report the estimated effects on GPA from our study, alongside prior work examining the impact of achievement incentives, advising, and grants, ordered from most to least costly.³³ As depicted in the figure, our intervention is characterized by relatively low costs, while the estimated effects are equal to or greater than those observed in most previous studies. Only a handful of interventions surpass ours in terms of impact, but they come with a two to five-fold higher cost per participant. In particular, our intervention compares favorably to achievement incentives that condition rewards on GPA. These results

³²Outside the U.S., Lusher et al. (2019) estimates that shifting class start times by an hour increases average sleep by about four minutes among Vietnamese University students with no effect on performance (Lusher et al., 2019). Jagnani (2021) estimates that the sun setting one hour earlier increases sleep by an average of 30 minutes among Indian children and that the sun setting 10 minutes earlier improves test scores by 0.10 SD and leads to 0.14 more years of schooling.

³³Achievement incentives include performance-based incentives and merit aid. Advising includes advising and support services. We report authors' OLS estimates of effects on non-cumulative GPA, either at the semester or year-level. Costs are per program participant per semester. We note that for some of these programs the primary outcome may have been enrollment, persistence or graduation and GPA may have been a secondary outcome. We do not adjust program costs for inflation (note this overestimates the cost of our intervention relative to others). See Table A.12 for more information on each study.

Figure 6: Cost-E ectiveness



Notes The figure compares our Immediate Incentives e ect on GPA to estimates from previous interventions aimed at improving college academic performance. Studies are grouped on the vertical axis based on their cost per subject per semester. Bars represent 95% confidence intervals. Superscripts above paper names denote di erent treatment arms or treatment groups. For Goldrick-Rab et al. (2016), superscript A is an estimate for the first cohort studied and B is their pooled estimate for the second and third cohort. For Denning et al. (2019), A and B are estimates for first-year and returning students, respectively. For Angrist et al. (2009), A is an estimate for an advising and peer-support treatment arm, B is for a financial incentives arm, and C for an arm combining A and B. For Evans et al. (2020), A estimates a grant treatment arm, and B estimates combined grant aid with academic advising. For Oreopoulos and Petronijevic (2018), course grades on a 0-100 scale have been divided by 25 for comparability to 4.0 scale GPA e ects.

suggest that incentives to sleep may be more e ective at improving GPA than incentivizing GPA directly. More generally, our study demonstrates that focusing on sleep can be a cost-e ective approach to improving educational outcomes.

5 Conclusion

In this paper, we show that an intervention targeting sleep habits improves academic performance. We explore mechanisms for the impact of our intervention on both sleep habits and academic performance. Inspired by cue-based theories of habit formation, our Immediate Incentives intervention aimed to establish automatic habits through repeated exposure to recurring cues coupled with immediate rewards. The intervention increases sleep during the

treatment period with smaller persistent effects in the post-treatment period. Our results show that Immediate Incentives can enhance habit formation during the intervention period, compared to variants with delayed or no rewards. However, we find little evidence that Immediate Incentives generate automatic habits triggered by the external cue. Instead, our results point to participants developing their own routines that persist into the post-treatment period. This could reflect some combination of treated participants acquiring a taste for sleep (i.e., increased benefits) and also finding sleep behaviors that are easier to sustain (i.e., lower costs). Future research could develop intervention designs that separately identify mechanisms of habit formation, including automaticity, learning about benefits, and lowering costs (Volpp and Loewenstein, 2020).

We then examine channels through which sleep may influence academic performance, including cognitive function, lifestyle factors, and overall well-being. While we do not detect an impact of our intervention on performance in math questions or creativity, sleep could have influenced cognition through channels like attention or memory consolidation, which were not captured by our measures (Diekelmann and Born, 2010). Examining lifestyle, our intervention led to a decrease in screen time and a reallocation of study time to morning hours, when students are potentially more alert and able to focus. Finally, we find evidence of a positive impact on students' ability to cope with stress, which may in turn have affected their academic performance. Further investigation of these mechanisms in future research can provide a deeper understanding of the multifaceted contributions of sleep to educational outcomes, as well as potential interactions between sleep, social media and mental health (Lindquist and Sado , 2023).

Taken together, our results show that offering incentives in the middle of the semester can improve term GPA. This result is consistent with recent evidence from Liu et al. (2022) that engagement interventions are more effective in the middle of the term. Future work could examine targeting the intervention (for example to first-term freshmen), the role of the timing of the intervention, and the impact of longer (or shorter) interventions. This can further our understanding of how to sustain improvements in academic performance across multiple terms.

Finally, our results suggest that targeting sleep may be a more cost-effective way to improve student performance than incentivizing performance directly. This could be because incentives based on sleep can be immediate whereas incentives based on grades are typically offered with a delay (e.g., at the end of the term). It could also be the case that students do not fully understand the production function for grades – and in particular, may not recognize the role of sleep. Future research could explore how beliefs and information about sleep shape individual behavior and educational outcomes.

References

- Abraham, K. G. and Flood, S. M. (2009). American time use survey data extract builder (ATUS-X). *International Journal for Time Use Research*, 6:167–168.
- Anderson, M. L. (2008). Multiple inference and gender differences in the effects of early intervention: A reevaluation of the abecedarian, perry preschool, and early training projects. *Journal of the American Statistical Association*, 103(484):1481–1495.
- Angrist, J., Lang, D., and Oreopoulos, P. (2009). Incentives and services for college achievement: Evidence from a randomized trial. *American Economic Journal: Applied Economics*, 1(1):136–163.
- Angrist, J., Oreopoulos, P., and Williams, T. (2014). When opportunity knocks, who answers? New evidence on college achievement awards. *Journal of Human Resources*, 49(3):572–610.
- Avery, M., Giuntella, O., and Jiao, P. (2022). Why don’t we sleep enough? A field experiment among college students. *Review of Economics and Statistics, Forthcoming*.
- Bachireddy, C., Joung, A., John, L. K., Gino, F., Tuckfield, B., Foschini, L., and Milkman, K. L. (2019). Effect of different financial incentive structures on promoting physical activity among adults: A randomized clinical trial. *JAMA Network Open*, 2(8):e199863–e199863.
- Ballard, J. (2019). 40% of Americans don’t generally wake up feeling well-rested. *YouGov*. Available at: <https://today.yougov.com/topics/health/articles-reports/2019/03/13/sleep-habits-americans-survey-poll> Accessed: March 2019.
- Banks, S. and Dinges, D. F. (2007). Behavioral and physiological consequences of sleep restriction. *Journal of Clinical Sleep Medicine*, 3(5):519–528.
- Barnes, C. M., Miller, J. A., and Bostock, S. (2017). Helping employees sleep well: Effects of cognitive behavioral therapy for insomnia on work outcomes. *Journal of Applied Psychology*, 102(1):104.
- Barrow, L., Richburg-Hayes, L., Rouse, C. E., and Brock, T. (2014). Paying for performance: The education impacts of a community college scholarship program for low-income adults. *Journal of Labor Economics*, 32(3):563–599.
- Basner, M., Fomberstein, K. M., Razavi, F. M., Banks, S., William, J. H., Rosa, R. R., and Dinges, D. F. (2007). American time use survey: sleep time and its relationship to waking activities. *Sleep*, 30(9):1085–1095.
- Beshears, J., Lee, H. N., Milkman, K. L., Mislavsky, R., and Wisdom, J. (2021). Creating exercise habits using incentives: The trade-off between flexibility and routinization. *Management Science*, 67(7):4139–4171.
- Bessone, P., Rao, G., Schilbach, F., Schofield, H., and Toma, M. (2021). The economic consequences of increasing sleep among the urban poor. *The Quarterly Journal of Economics*, 136(3):1887–1941.

- Biddle, J. E. and Hamermesh, D. S. (1990). Sleep and the allocation of time. *Journal of Political Economy*, 98(5, Part 1):922–943.
- Boyce, R., Glasgow, S. D., Williams, S., and Adamantidis, A. (2016). Causal evidence for the role of REM sleep theta rhythm in contextual memory consolidation. *Science*, 352(6287):812–816.
- Breig, Z., Gibson, M., and Shrader, J. (2020). Why do we procrastinate? present bias and optimism. *Present Bias and Optimism (August 27, 2020)*.
- Byrne, D. P., Goette, L., Martin, L. A., Miles, A., Jones, A., Schob, S., Staake, T., and Tiefenbeck, V. (2022). The habit forming effects of feedback: Evidence from a large-scale field experiment. Available at SSRN: <https://ssrn.com/abstract=3974371>.
- Cappelen, A. W., Charness, G., Ekström, M., Gneezy, U., and Tungodden, B. (2017). Exercise improves academic performance. *NHH Dept. of Economics Discussion Paper*, (08).
- Cappuccio, F. P., D’Elia, L., Strazzullo, P., and Miller, M. A. (2010). Sleep duration and all-cause mortality: A systematic review and meta-analysis of prospective studies. *Sleep*, 33(5):585–592.
- Carrell, S. E., Maghakian, T., and West, J. E. (2011). A’s from zzzz’s? The causal effect of school start time on the academic achievement of adolescents. *American Economic Journal: Economic Policy*, 3(3):62–81.
- CDC (2023). About CDC: Sleep and sleep disorders. Accessed Nov 5, 2023. <https://www.cdc.gov/sleep/data-and-statistics/adults.html>.
- Charness, G. and Gneezy, U. (2009). Incentives to exercise. *Econometrica*, 77(3):909–931.
- Charness, G. and Grieco, D. (2019). Creativity and incentives. *Journal of the European Economic Association*, 17(2):454–496.
- Clotfelter, C. T., Hemelt, S. W., and Ladd, H. F. (2018). Multifaceted aid for low-income students and college outcomes: Evidence from North Carolina. *Economic Inquiry*, 56(1):278–303.
- Cooney, M. T., Vartiainen, E., Laakitainen, T., Juolevi, A., Dudina, A., and Graham, I. M. (2010). Elevated resting heart rate is an independent risk factor for cardiovascular disease in healthy men and women. *American Heart Journal*, 159(4):612–619.
- Corkett, S. (2010). 2020 sleep in america® poll shows alarming level of sleepiness and low levels of action. *National Sleep Foundation*.
- Creswell, J. D., Tuminia, M. J., Price, S., Sefidgar, Y., Cohen, S., Ren, Y., Brown, J., Dey, A. K., Dutcher, J. M., Villalba, D., et al. (2023). Nightly sleep duration predicts grade point average in the first year of college. *Proceedings of the National Academy of Sciences*, 120(8):e2209123120.

- de Zambotti, M., Goldstone, A., Claudatos, S., Colrain, I. M., and Baker, F. C. (2018). A validation study of fitbit charge 2™ compared with polysomnography in adults. *Chronobiology International*, 35(4):465–476.
- Denning, J. T., Marx, B. M., and Turner, L. J. (2019). Propelled: The effects of grants on graduation, earnings, and welfare. *American Economic Journal: Applied Economics*, 11(3):193–224.
- Dickinson, A. (1985). Actions and habits: The development of behavioural autonomy. *Philosophical Transactions of the Royal Society of London. B, Biological Sciences*, 308(1135):67–78.
- Dickinson, D. L. and Masclot, D. (2023). Unethical decision making and sleep restriction: Experimental evidence. *Games and Economic Behavior*, 141:484–502.
- Dickinson, D. L. and McElroy, T. (2017). Sleep restriction and circadian effects on social decisions. *European Economic Review*, 97:57–71.
- Diekelmann, S. and Born, J. (2010). The memory function of sleep. *Nature Reviews Neuroscience*, 11(2):114–126.
- Dyer, A. R., Persky, V., Stamler, J., Paul, O., Shekelle, R. B., Berkson, D. M., Lepper, M., Schoenberger, J. A., and Lindberg, H. A. (1980). Heart rate as a prognostic factor for coronary heart disease and mortality: Findings in three Chicago epidemiologic studies. *American Journal of Epidemiology*, 112(6):736–749.
- Evans, W. N., Kearney, M. S., Perry, B., and Sullivan, J. X. (2020). Increasing community college completion rates among low-income students: Evidence from a randomized controlled trial evaluation of a case-management intervention. *Journal of Policy Analysis and Management*, 39(4):930–965.
- Gibson, M. and Shrader, J. (2018). Time use and labor productivity: The returns to sleep. *Review of Economics and Statistics*, 100(5):783–798.
- Giuntella, O., Han, W., and Mazzonna, F. (2017). Circadian rhythms, sleep, and cognitive skills: Evidence from an unsleeping giant. *Demography*, 54(5):1715–1742.
- Giuntella, O., Hyde, K., Saccardo, S., and Sado, S. (2021). Lifestyle and mental health disruptions during COVID-19. *Proceedings of the National Academy of Sciences*, 118(9):e2016632118.
- Giuntella, O. and Mazzonna, F. (2019). Sunset time and the economic effects of social jetlag: Evidence from US time zone borders. *Journal of Health Economics*, 65:210–226.
- Gneezy, U., Meier, S., and Rey-Biel, P. (2011). When and why incentives (don’t) work to modify behavior. *Journal of Economic Perspectives*, 25(4):191–210.
- Goldrick-Rab, S., Kelchen, R., Harris, D. N., and Benson, J. (2016). Reducing income inequality in educational attainment: Experimental evidence on the impact of financial aid on college completion. *American Journal of Sociology*, 121(6):1762–1817.

- Groen, J. A. and Pabilonia, S. W. (2019). Snooze or lose: High school start times and academic achievement. *Economics of Education Review*, 72:204–218.
- Group, A. S. W., ADOLESCENCE, C. O., HEALTH, C. O. S., Au, R., Carskadon, M., Millman, R., Wolfson, A., Braverman, P. K., Adelman, W. P., Breuner, C. C., et al. (2014). School start times for adolescents. *Pediatrics*, 134(3):642–649.
- Hafner, M., Stepanek, M., Taylor, J., Troxel, W. M., and Van Stolk, C. (2017). Why sleep matters—the economic costs of insufficient sleep: A cross-country comparative analysis. *Rand Health Quarterly*, 6(4).
- Haghayegh, S., Khoshnevis, S., Smolensky, M. H., Diller, K. R., and Castriotta, R. J. (2019). Accuracy of wristband fitbit models in assessing sleep: Systematic review and meta-analysis. *Journal of Medical Internet Research*, 21(11):e16273.
- Heissel, J. A. and Norris, S. (2018). Rise and shine the effect of school start times on academic performance from childhood through puberty. *Journal of Human Resources*, 53(4):957–992.
- Hershner, S. D. and Chervin, R. D. (2014). Causes and consequences of sleepiness among college students. *Nature and Science of Sleep*, pages 73–84.
- Hillman, D. R., Murphy, A. S., Antic, R., and Pezzullo, L. (2006). The economic cost of sleep disorders. *Sleep*, 29(3):299–305.
- Hirshkowitz, M., Whiton, K., Albert, S. M., Alessi, C., Bruni, O., DonCarlos, L., Hazen, N., Herman, J., Hillard, P. J. A., Katz, E. S., et al. (2015). National sleep foundation’s updated sleep duration recommendations. *Sleep Health*, 1(4):233–243.
- Holbein, J. B., Schafer, J. P., and Dickinson, D. L. (2019). Insufficient sleep reduces voting and other prosocial behaviours. *Nature Human Behaviour*, 3(5):492–500.
- Hussam, R. N., Rabbani, A., Reggiani, G., and Rigol, N. (2022). Rational habit formation: Experimental evidence from handwashing in India. *American Economic Journal: Applied Economics*, 14(1):1–41.
- Jagnani, M. (2021). Children’s sleep and human capital production. *The Review of Economics and Statistics*, *Forthcoming*.
- Jike, M., Itani, O., Watanabe, N., Buysse, D. J., and Kaneita, Y. (2018). Long sleep duration and health outcomes: A systematic review, meta-analysis and meta-regression. *Sleep Medicine Reviews*, 39:25–36.
- Jin, L. and Ziebarth, N. R. (2020). Sleep, health, and human capital: Evidence from daylight saving time. *Journal of Economic Behavior & Organization*, 170:174–192.
- Jones, J. (2013). In U.S., 40% get less than recommended amount of sleep. *Gallup*. Available at: <https://news.gallup.com/poll/166553/less-recommended-amount-sleep.aspx>.

- Killgore, W. D. (2010). Effects of sleep deprivation on cognition. *Progress in Brain Research*, 185:105–129.
- Lauderdale, D. S., Knutson, K. L., Yan, L. L., Liu, K., and Rathouz, P. J. (2008). Self-reported and measured sleep duration: How similar are they? *Epidemiology*, 19(6):838–845.
- Lavecchia, A. M., Liu, H., and Oreopoulos, P. (2016). Behavioral economics of education: Progress and possibilities. In *Handbook of the Economics of Education*, volume 5, pages 1–74. Elsevier.
- Leproult, R. and Van Cauter, E. (2010). Role of sleep and sleep loss in hormonal release and metabolism. *Pediatric Neuroendocrinology*, 17:11–21.
- Levitt, S. D., List, J. A., Neckermann, S., and Sado , S. (2016). The behavioralist goes to school: Leveraging behavioral economics to improve educational performance. *American Economic Journal: Economic Policy*, 8(4):183–219.
- Lim, J. and Dinges, D. F. (2010). A meta-analysis of the impact of short-term sleep deprivation on cognitive variables. *Psychological bulletin*, 136(3):375.
- Lindquist, S. and Sado , S. (2023). Understanding the interactions of sleep, social media and mental health for productivity and performance: The role of field experiments. *Annual Proceedings of the Upton Forum*.
- Liu, T. X., Malmendier, U., Wang, S. W., and Zhang, S. (2022). Not too early, not too late: Encouraging engagement in education [working paper]. Available at: https://eml.berkeley.edu/~ulrike/Papers/NotesTaking_Public.pdf.
- Lusher, L., Yassenov, V., and Luong, P. (2019). Does schedule irregularity affect productivity? Evidence from random assignment into college classes. *Labour Economics*, 60:115–128.
- Marks, G. A., Sha erry, J. P., Oksenberg, A., Speciale, S. G., and Ro warg, H. P. (1995). A functional role for rem sleep in brain maturation. *Behavioural Brain Research*, 69(1-2):1–11.
- McKenna, B. S., Dickinson, D. L., Or , H. J., and Drummond, S. P. (2007). The effects of one night of sleep deprivation on known-risk and ambiguous-risk decisions. *Journal of Sleep Research*, 16(3):245–252.
- Milkman, K. L., Minson, J. A., and Volpp, K. G. (2014). Holding the hunger games hostage at the gym: An evaluation of temptation bundling. *Management Science*, 60(2):283–299.
- Mullainathan, S. (2014). Get some sleep, and wake up the G.D.P. *The New York Times*. Available at: <https://www.nytimes.com/2014/02/02/business/get-some-sleep-and-wake-up-the-gdp.html>.
- Oreopoulos, P. and Petronijevic, U. (2018). Student coaching: How far can technology go? *Journal of Human Resources*, 53(2):299–329.

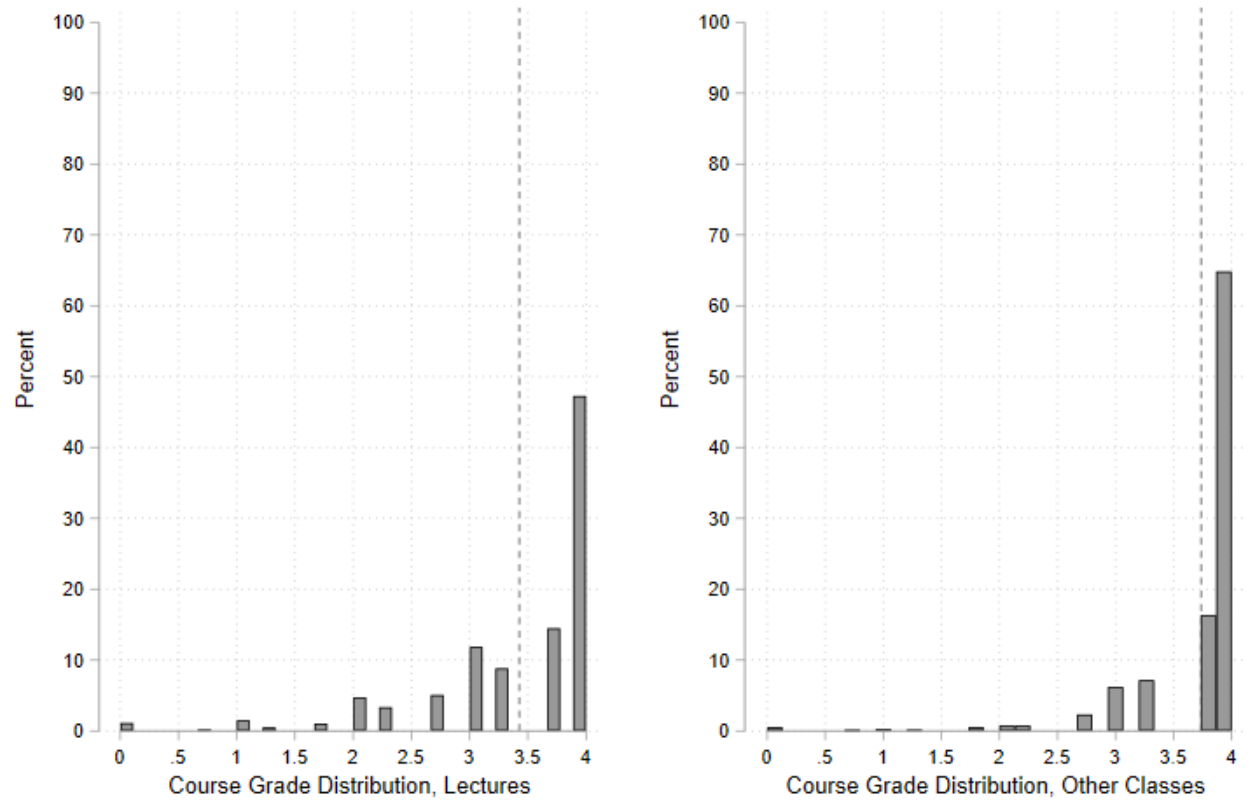
- Panel, C. C., Watson, N. F., Badr, M. S., Belenky, G., Bliwise, D. L., Buxton, O. M., Buysse, D., Dinges, D. F., Gangwisch, J., Grandner, M. A., et al. (2015). Recommended amount of sleep for a healthy adult: A joint consensus statement of the american academy of sleep medicine and sleep research society. *Journal of Clinical Sleep Medicine*, 11(6):591–592.
- Park, R. S. E. and Scott-Clayton, J. (2018). The impact of Pell Grant eligibility on community college students’ financial aid packages, labor supply, and academic outcomes. *Educational Evaluation and Policy Analysis*, 40(4):557–585.
- Phillips, A. J., Clerx, W. M., O’Brien, C. S., Sano, A., Barger, L. K., Picard, R. W., Lockley, S. W., Klerman, E. B., and Czeisler, C. A. (2017). Irregular sleep/wake patterns are associated with poorer academic performance and delayed circadian and sleep/wake timing. *Scientific Reports*, 7(1):3216.
- Quer, G., Gouda, P., Galarnyk, M., Topol, E. J., and Steinhubl, S. R. (2020). Inter-and intraindividual variability in daily resting heart rate and its associations with age, sex, sleep, bmi, and time of year: Retrospective, longitudinal cohort study of 92,457 adults. *Plos one*, 15(2):e0227709.
- Radlo, L. S. (1977). The CES-D scale: A self-report depression scale for research in the general population. *Applied Psychological Measurement*, 1(3):385–401.
- Rao, G., Redline, S., Schilbach, F., Schofield, H., and Toma, M. (2021). Informing sleep policy through field experiments. *Science*, 374(6567):530–533.
- Roenneberg, T. (2013). The human sleep project. *Nature*, 498(7455):427–428.
- Roenneberg, T. and Merrow, M. (2016). The circadian clock and human health. *Current biology*, 26(10):R432–R443.
- Royer, H., Stehr, M., and Sydnor, J. (2015). Incentives, commitments, and habit formation in exercise: Evidence from a field experiment with workers at a fortune-500 company. *American Economic Journal: Applied Economics*, 7(3):51–84.
- Schultz, P. W., Nolan, J. M., Cialdini, R. B., Goldstein, N. J., and Griskevicius, V. (2007). The constructive, destructive, and reconstructive power of social norms. *Psychological Science*, 18(5):429–434.
- Scott-Clayton, J. (2011). On money and motivation a quasi-experimental analysis of financial incentives for college achievement. *Journal of Human Resources*, 46(3):614–646.
- Smarr, B. L. (2015). Digital sleep logs reveal potential impacts of modern temporal structure on class performance in different chronotypes. *Journal of Biological Rhythms*, 30(1):61–67.
- Trockel, M. T., Barnes, M. D., and Egget, D. L. (2000). Health-related variables and academic performance among first-year college students: Implications for sleep and other behaviors. *Journal of American College Health*, 49(3):125–131.

- Verplanken, B. and Wood, W. (2006). Interventions to break and create consumer habits. *Journal of Public Policy & Marketing*, 25(1):90–103.
- Volpp, K. G. and Loewenstein, G. (2020). What is a habit? Diverse mechanisms that can produce sustained behavior change. *Organizational Behavior and Human Decision Processes*, 161:36–38.
- Wellsjo, A. S. (2021). Simple actions, complex habits: Lessons from hospital hand hygiene [working paper]. Available at: https://economics.ucr.edu/wp_content/uploads/2023/01/18_23_wellsjo.pdf.
- Wheaton, A. G., Jones, S. E., Cooper, A. C., and Croft, J. B. (2018). Short sleep duration among middle school and high school students—united states, 2015. *Morbidity and Mortality Weekly Report*, 67(3):85.
- Williams, N. (2014). The GAD-7 questionnaire. *Occupational Medicine*, 64(3):224–224.
- Wood, W. and Neal, D. T. (2007). A new look at habits and the habit-goal interface. *Psychological Review*, 114(4):843.
- Wood, W. and Rünger, D. (2016). Psychology of habit. *Annual Review of Psychology*, 67:289–314.
- Wright, S. P., Hall Brown, T. S., Collier, S. R., and Sandberg, K. (2017). How consumer physical activity monitors could transform human physiology research. *American Journal of Physiology-Regulatory, Integrative and Comparative Physiology*, 312(3):R358–R367.
- Zhang, Y. and Gruber, R. (2019). Can slow-wave sleep enhancement improve memory? A review of current approaches and cognitive outcomes. *The Yale Journal of Biology and Medicine*, 92(1):63–80.

Appendix

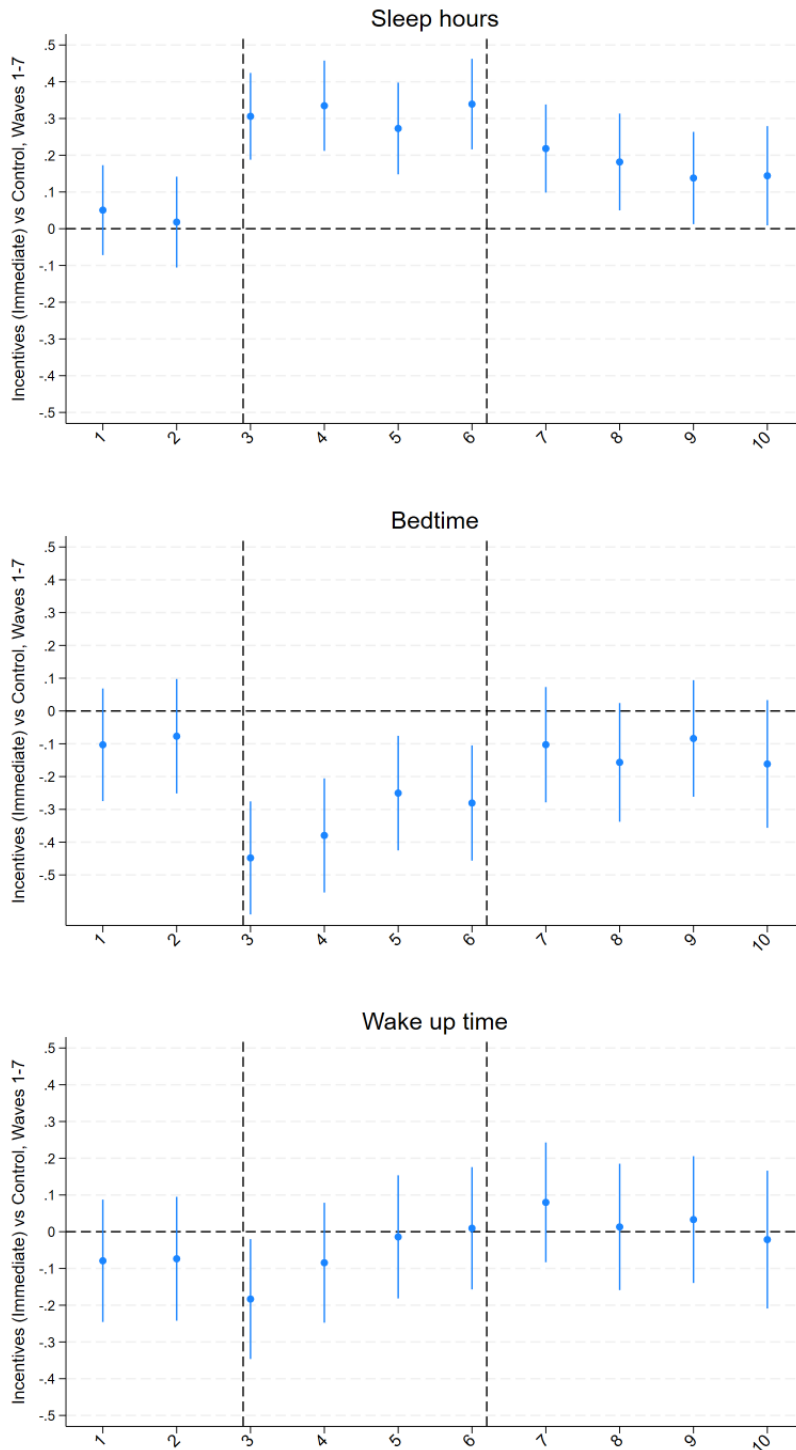
A. Figures and tables

Figure A.1: Grades distribution



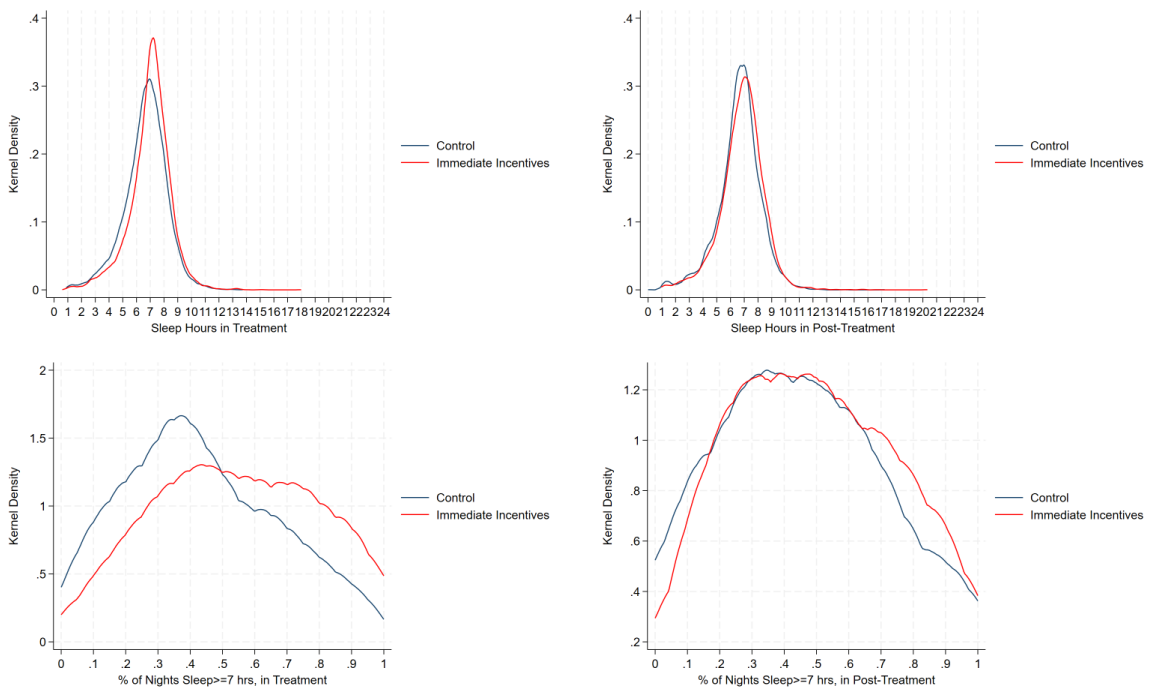
Notes The figure reports the distribution of grades in lectures and other classes. The dashed vertical line identifies the average grade in these class types.

Figure A.2: Immediate Incentives, sleep hours, bedtime and wake-up time



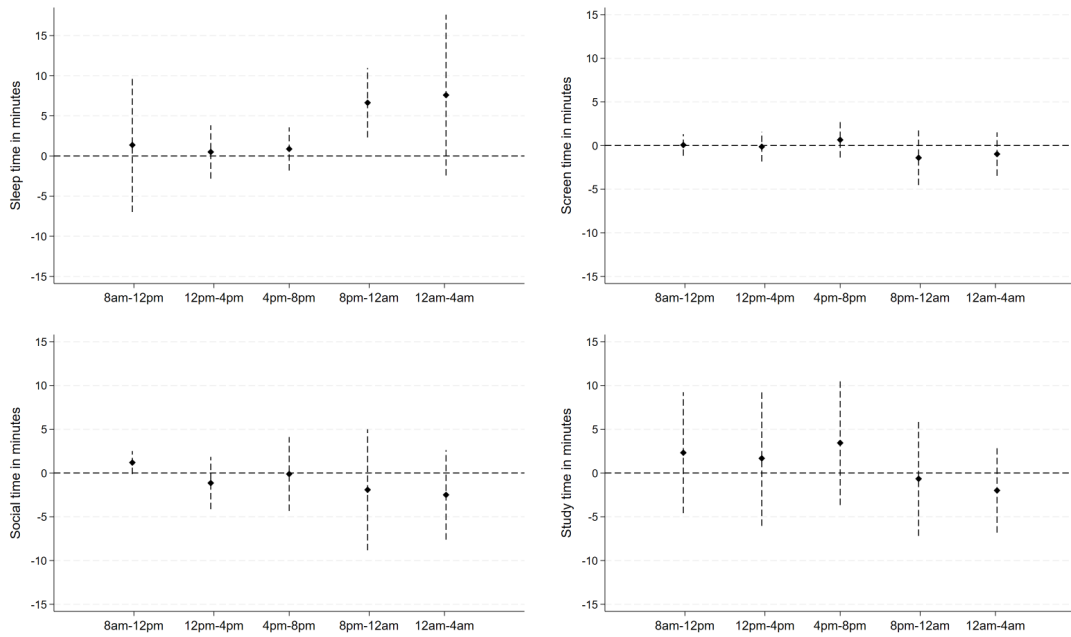
Notes The sample is restricted to weekdays (Sunday-Thursday nights). On the horizontal axis we report time in weeks since the study started (week 3 is the first week of treatment, week 6 is the last week of treatment). The coefficient reports the differences in average sleep hours, bedtime, and wake up time between individuals in the Immediate Incentives treatment and those in Control by week. Standard errors are clustered at the individual level. Bars indicate 95% confidence intervals.

Figure A.3: Immediate incentives and distribution of sleep during and after the intervention



Notes The figure reports kernel densities of the Immediate Incentives (red) and Control (navy) groups for a) sleep hours (top panel) and b) the proportion of nights with sleep over seven hours (bottom panel), during treatment (leftward graphs) and post-treatment (rightward graphs).

Figure A.4: Immediate Incentives to sleep and time use over the day: Post-Intervention period



Notes The figure reports differences between participants in the Immediate Incentives treatment and Control groups in the minutes allocated to different time-use activities post-treatment throughout the day. All the coefficients are obtained from regressions including wave, month and day of the week fixed effects, baseline value of the outcome variable, and demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), indicators for the number of classes starting at 10 a.m. or earlier, indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. Standard errors are clustered at the individual level. Bars indicate 95% confidence intervals.

Table A.1: Grading system

(1) Grade	(2) GPA	(3) Quality points	(4) Has a grade	(5) Withdrawn	(6) Passed	(7) Credit completed
A+	YES	4	YES	NO	YES	YES
A	YES	4	YES	NO	YES	YES
A-	YES	3.75	YES	NO	YES	YES
B+	YES	3.25	YES	NO	YES	YES
B	YES	3	YES	NO	YES	YES
B-	YES	2.75	YES	NO	YES	YES
C+	YES	2.25	YES	NO	YES	YES
C	YES	2	YES	NO	YES	YES
C-	YES	1.75	YES	NO	YES	YES
D+	YES	1.25	YES	NO	YES	YES
D	YES	1	YES	NO	YES	YES
D-	YES	0.75	YES	NO	YES	YES
F	YES	0	YES	NO	NO	NO
G	–	–	NO	NO	NO	NO
H	–	–	YES	NO	YES	YES
HS	–	–	YES	NO	YES	YES
I	–	–	NO	NO	NO	NO
N	–	–	NO	NO	NO	NO
NC	–	–	NO	NO	NO	NO
NG	–	–	NO	NO	NO	NO
R	–	–	NO	NO	NO	NO
S	–	–	YES	NO	YES	YES
U	–	–	YES	NO	NO	NO
W	–	–	NO	YES	NO	NO

Notes Non-grade outcomes (G-W) represent the following: G, unfinished or ongoing course work due extenuating personal circumstances; H, honors, exceptional completion of coursework; HS, highly satisfactory completion of coursework, used only by School of Medicine; I, unfinished or ongoing course work due to nature of course; N, non-credited or graded course, such as a course audit; NC, non-credit course; NG, non-credit course due unfinished course work; R, student resigned from University; S, satisfactory completion of requirements; U, unsatisfactory completion of class requirements; W, student withdrew from course. Source: [https://www.registrar.pitt.edu/sites/default/files/pdf/Grading 20System.pdf](https://www.registrar.pitt.edu/sites/default/files/pdf/Grading%20System.pdf)

Table A.2: Incentives and attrition

	(1)	(2)	(3)	(4)	(5)	(6)
Panel A						
	# synced days before intervention	# synced days during intervention	# synced days post intervention	Has HS GPA	Has baseline GPA	Has course grades
Immediate Incentives	0.145 (0.260)	0.779** (0.321)	-0.096 (0.406)	0.017 (0.020)	0.046* (0.025)	0.008 (0.010)
Observations	840	840	840	840	840	840
Mean of dep. var.	11.89	18.17	13.03	0.895	0.782	0.981
Std. dev.	6.021	4.100	6.428	0.307	0.413	0.137
Panel B						
	Has time use	Has math task	Has creativity task	Has mood survey	Has resilience survey	Has mental health
Immediate Incentives	0.006 (0.007)	0.010 (0.013)	0.010 (0.015)	0.005 (0.011)	-0.003 (0.012)	-0.001 (0.006)
Observations	840	840	840	840	840	840
Mean of dep. var.	0.988	0.872	0.951	0.967	0.959	0.991
Std. dev.	0.110	0.334	0.216	0.178	0.198	0.0931

Notes The table reports the difference between the Immediate Incentives treatment and Control groups in attrition rate across the different outcome measures. All estimates include wave fixed effects. Robust standard errors are in parenthesis. Mean of dep. var. is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p < 0.01, ** p < 0.05, * p < 0.1.

Table A.3: Differences in baseline characteristics: Secondary treatments

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Female	Age	Asian	Black	Hispanic	White	Other	
Delayed Incentives	0.0111 (0.060)	0.1893 (0.492)	-0.0640 (0.059)	0.0126 (0.037)	-0.0104 (0.021)	0.0698 (0.065)	-0.0080 (0.022)	
Delayed Incentives, No Feedback	-0.0164 (0.062)	-0.0607 (0.390)	-0.0290 (0.061)	-0.0020 (0.034)	0.0110 (0.027)	0.0160 (0.066)	0.0041 (0.026)	
Feedback Only	0.0051 (0.060)	-0.0994 (0.385)	-0.1195** (0.056)	0.0636 (0.041)	0.0194 (0.028)	0.0143 (0.066)	0.0222 (0.027)	
Observations	442	442	442	442	442	442	442	
Mean of dep. var.	0.686	19.83	0.271	0.0973	0.0362	0.554	0.0407	
Std. dev.	0.465	3.270	0.445	0.297	0.187	0.498	0.198	
Immediate Incentives	0.0139 (0.056)	-0.0302 (0.252)	0.0582 (0.058)	-0.0344 (0.034)	0.0190 (0.033)	-0.0083 (0.066)	-0.0344 (0.034)	
No Cue/Feedback in Post								
Observations	224	224	224	224	224	224	224	
Mean of dep. var.	0.777	19.31	0.259	0.0714	0.0625	0.536	0.0714	
Std. dev.	0.417	1.927	0.439	0.258	0.243	0.500	0.258	
	Sleep \geq 7 hours	Sleep \geq 6 hours	Sleep hours	Bedtime	Wake-up time	HS GPA	Baseline GPA	
Delayed Incentives	-0.0338 (0.034)	-0.0197 (0.029)	0.0081 (0.120)	-0.2647* (0.146)	-0.3047** (0.142)	0.0633 (0.053)	-2.1743 (4.442)	
Delayed Incentives, No Feedback	-0.0442 (0.033)	-0.0246 (0.030)	-0.0731 (0.108)	0.0056 (0.150)	-0.0501 (0.149)	0.0952* (0.057)	-3.8859 (4.570)	
Feedback Only	-0.0210 (0.033)	-0.0092 (0.030)	-0.0199 (0.101)	-0.2074 (0.143)	-0.2445* (0.144)	0.1036* (0.059)	-5.5754 (4.293)	
Observations	444	444	444	444	444	385	437	
Mean of dep. var.	0.447	0.733	6.754	24.99	7.774	4.126	22.16	
Std. dev.	0.259	0.230	0.841	1.128	1.116	0.399	38.08	
Immediate Incentives	-0.0142 (0.036)	-0.0380 (0.034)	-0.1351 (0.117)	0.0074 (0.175)	-0.0388 (0.161)	-0.0161 (0.060)	-0.0693 (1.703)	
No Cue/Feedback in Post								
Observations	224	224	224	224	224	215	223	
Mean of dep. var.	0.457	0.737	6.758	25.23	8.064	4.138	5.212	
Std. dev.	0.267	0.254	0.886	1.307	1.222	0.444	12.71	
	Freshman	Sophomore	Junior	Senior	STEM major	Less than college	College	More than college
Delayed Incentives	-0.0242 (0.063)	-0.0466 (0.045)	0.0479 (0.056)	0.0354 (0.049)	0.1369** (0.062)	-0.0274 (0.059)	0.0350 (0.062)	-0.0077 (0.064)
Delayed Incentives, No Feedback	-0.0327 (0.065)	-0.0361 (0.048)	0.0218 (0.056)	0.0595 (0.052)	0.0051 (0.067)	0.0196 (0.061)	0.0565 (0.063)	-0.0760 (0.065)
Feedback Only	-0.0303 (0.062)	0.0075 (0.049)	-0.0164 (0.053)	0.0516 (0.051)	0.1016 (0.064)	-0.0123 (0.060)	-0.0069 (0.060)	0.0192 (0.065)
Observations	444	444	444	444	442	442	442	442
Mean of dep. var.	0.432	0.160	0.223	0.180	0.620	0.287	0.312	0.400
Std. dev.	0.496	0.367	0.417	0.385	0.486	0.453	0.464	0.491
Immediate Incentives	0.1002* (0.060)	-0.0253 (0.037)	-0.0412 (0.051)	-0.0337 (0.042)	0.0718 (0.067)	0.0321 (0.062)	0.1034* (0.061)	-0.1355** (0.066)
No Cue/Feedback in Post								
Observations	224	224	224	224	224	224	224	224
Mean of dep. var.	0.607	0.0848	0.192	0.116	0.522	0.299	0.290	0.411
Std. dev.	0.489	0.279	0.395	0.321	0.501	0.459	0.455	0.493

Notes: All estimates include wave fixed effects. Robust standard errors are in parenthesis. Mean of dep. var. is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p < 0.01, ** p < 0.05, * p < 0.1.

Table A.4: Immediate incentives and sleep: Sensitivity analysis

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Primary	Basic	Has term	Excludes	Excludes	Weighted	Incentives,	No Cue
	Specification	Controls	GPA	missing	wave 3	by gender	pooled	in Post
Treatment	0.1186*** (0.013)	0.1134*** (0.011)	0.1172*** (0.013)	0.1270*** (0.012)	0.1199*** (0.014)	0.1139*** (0.014)	0.1094*** (0.012)	0.1193*** (0.021)
Post-Treatment	0.0551*** (0.015)	0.0607*** (0.011)	0.0546*** (0.015)	0.0539*** (0.014)	0.0635*** (0.016)	0.0565*** (0.015)	0.0530*** (0.014)	0.0626** (0.025)
Immediate Cash*								0.0097 (0.027)
No Cue in Post-Treatment								
Observations	46,989	46,989	46,146	35,182	41,753	46,989	58,833	19,932
Mean of dep. var.	0.429	0.429	0.427	0.432	0.434	0.429	0.435	6.742
Std. dev.	0.495	0.495	0.495	0.495	0.496	0.495	0.496	0.498
Number of individuals	840	840	825	840	763	840	1040	357

Notes The sample is restricted to individuals in the Immediate Incentives treatment and individuals in the Control group. Individuals in the Cue/Feedback treatment were not included in this analysis. All estimates except those in Column 2 include day of the week, week of the experiment, wave, and month fixed effects, baseline value of the outcome variable, indicators for the number of classes starting at 10 a.m. or earlier, and demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. Estimates in Column 2 includes only wave fixed effects, baseline value of the outcome variable, controls for gender, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). Column 4 does not replace missing nights with baseline data as in our main analysis. Standard errors are clustered at the individual level. Mean of dep. var. is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p 0.01, ** p 0.05, * p 0.1.

Table A.5: Immediate Incentives and sleep: By quartiles of sleep at baseline

	(1)	(2)	(3)	(4)
	Q1	Q2	Q3	Q4
Sleep \geq 7				
Treatment	0.1495*** (0.023)	0.1111*** (0.024)	0.1507*** (0.022)	0.0977*** (0.026)
Post-Treatment	0.0952*** (0.021)	0.0367 (0.028)	0.0410 (0.029)	0.0278 (0.030)
Observations	11,999	11,907	8,763	8,729
Mean of dep. var.	0.116	0.320	0.516	0.791
Std. dev.	0.320	0.467	0.500	0.407
Number of individuals	213	211	206	210
Sleep hours				
Treatment	0.3569*** (0.072)	0.3469*** (0.077)	0.4110*** (0.070)	0.2763*** (0.083)
Post-Treatment	0.2075*** (0.072)	0.1928** (0.086)	0.2096** (0.092)	0.0820 (0.090)
Observations	11,999	8,906	8,763	8,729
Mean of dep. var.	5.773	6.394	6.919	7.627
Std. dev.	1.378	1.444	1.466	1.150
Number of individuals	213	211	206	210

Notes The sample is restricted to individuals in the Immediate Incentives treatment and individuals in the Control group. Individuals in the Cue/Feedback treatment were not included in this analysis. All estimates include day of the week, week of the experiment, wave, and month fixed effects, baseline value of the outcome variable, indicators for the number of classes starting at 10 a.m. or earlier, and demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. Standard errors are clustered at the individual level. Mean of dep. var. is the mean of the dependent variable in the control group at baseline. Std. dev. is the standard deviation of the dependent variable in the control group at baseline. *** p 0.01, ** p 0.05, * p 0.1.

Table A.6: Immediate incentives and sleep: Additional outcomes

	(1)	(2)	(3)	(4)	(5)
	Any nap	Sleep ≥ 7 weekends	Sleep ≥ 7 weekends & holidays	Sleep ≥ 7 all nights & naps	Sleep hours all nights & naps
Treatment	0.0022 (0.004)	0.0143 (0.015)	0.0096 (0.016)	0.0687*** (0.011)	0.2062*** (0.031)
Post-Treatment	-0.0050 (0.004)	0.0481*** (0.016)	0.0458*** (0.016)	0.0409*** (0.012)	0.0940*** (0.031)
Observations	69,937	18,100	22,948	69,937	69,937
Mean of dep. var.	0.0513	0.478	0.478	0.568	7.213
Std. dev.	0.188	0.500	0.500	0.495	1.425
Number of individuals	840	840	840	840	840
	Sleep ≥ 6	Sleep 7-9	Efficiency	REM sleep	Deep sleep
Treatment	0.0805*** (0.010)	0.1099*** (0.011)	0.2497* (0.143)	2.4852*** (0.769)	0.4040 (0.648)
Post-Treatment	0.0368*** (0.010)	0.0491*** (0.012)	0.0694 (0.182)	1.7137** (0.872)	0.1183 (0.695)
Observations	46,989	46,989	46,989	43,168	43,168
Mean of dep. var.	0.714	0.382	93.52	84.12	74.52
Std. dev.	0.376	0.402	5.170	28.32	22.01
Number of individuals	840	840	840	798	798

Notes The sample is restricted to individuals in the Immediate Incentives treatment and individuals in the Control group. Individuals in the Cue/Feedback treatment were not included in this analysis. All estimates include day of the week, week of the experiment, wave, and month fixed effects, baseline value of the outcome variable, indicators for the number of classes starting at 10 a.m. or earlier, and demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. Standard errors are clustered at the individual level. Mean of dep. var. is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p < 0.01, ** p < 0.05, * p < 0.1.

Table A.7: Immediate Incentives and course grade: Sensitivity analysis

	(1) Primary specification	(2) Basic controls	(3) Has grade in term+1 or term+2	(4) No missing HS/baseline GPA	(5) Excludes obs with no sleep data	(6) Excludes wave 3 (Covid)	(7) Weighted by gender	(8) Incentives, pooled
Panel A: All classes								
Incentives	0.075** (0.037)	0.064* (0.038)	0.060* (0.035)	0.067* (0.037)	0.075** (0.037)	0.090** (0.040)	0.070* (0.039)	0.061* (0.035)
Observations	4,300	4,300	4,102	4,216	4,256	3,934	4,300	5,254
Mean of dep. var.	3.502	3.502	3.500	3.502	3.499	3.491	3.502	3.498
Std. dev.	0.763	0.763	0.759	0.765	0.766	0.776	0.763	0.763
Number of individuals	833	833	791	815	825	757	833	1027
Panel B: Lectures								
Incentives	0.088** (0.042)	0.076* (0.043)	0.074* (0.040)	0.078* (0.042)	0.088** (0.042)	0.105** (0.045)	0.082* (0.044)	0.075* (0.040)
Observations	3,413	3,413	3,255	3,340	3,382	3,130	3,413	4,197
Mean of dep. var.	3.436	3.436	3.435	3.435	3.434	3.423	3.436	3.435
Std. dev.	0.805	0.805	0.799	0.807	0.807	0.819	0.805	0.801
Number of individuals	827	827	787	809	819	752	827	1021
Panel C: Other classes (seminars, labs, etc.)								
Incentives	0.001 (0.040)	-0.008 (0.040)	-0.004 (0.045)	0.001 (0.040)	0.001 (0.040)	0.012 (0.041)	0.013 (0.045)	-0.012 (0.036)
Observations	887	887	726	876	874	804	887	1,057
Mean of dep. var.	3.753	3.753	3.759	3.755	3.751	3.753	3.753	3.748
Std. dev.	0.505	0.505	0.503	0.503	0.507	0.502	0.505	0.517
Number of individuals	562	562	449	554	554	510	562	674

Notes All estimates except those in column 2 include demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), baseline sleep, indicators for the number of classes starting at 10 a.m. or earlier, indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). Estimates in column 2 include only wave fixed effects, baseline sleep, controls for gender, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. Observations are weighted by the number of credits taken in the semester. Standard errors are clustered at the individual level. Mean of dep. var. is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p 0.01, ** p 0.05, * p 0.1.

Table A.8: Immediate Incentives, sleep and GPA: Heterogeneity

	(1) Male	(2) Female	(3) First-term	(4) Other students	(5) No-STEM major	(6) STEM major
Panel A: Sleep ≥ 7 hours						
Treatment	0.0951*** (0.025)	0.1312*** (0.015)	0.1534*** (0.031)	0.1119*** (0.014)	0.1129*** (0.019)	0.1275*** (0.017)
Post-Treatment	0.0642** (0.027)	0.0549*** (0.018)	0.1272*** (0.036)	0.0450*** (0.017)	0.0489** (0.024)	0.0642*** (0.020)
Observations	12,848	33,909	8,120	38,869	19,827	26,930
Mean of dep. var.	0.344	0.463	0.424	0.430	0.475	0.397
Std. dev.	0.475	0.499	0.494	0.495	0.499	0.489
Number of individuals	229	607	160	680	356	480
Panel B: Course grades, all classes						
Immediate Incentives	0.029 (0.076)	0.083* (0.044)	0.172** (0.072)	0.058 (0.043)	-0.018 (0.052)	0.130*** (0.049)
Observations	1,148	3,131	772	3,528	1,775	2,504
Mean of dep. var.	3.384	3.545	3.528	3.496	3.552	3.467
Std. dev.	0.851	0.724	0.705	0.775	0.719	0.791
Number of individuals	229	600	160	673	352	477
Panel C: Course grades, lectures						
Immediate Incentives	0.028 (0.085)	0.101** (0.051)	0.197** (0.086)	0.070 (0.048)	0.002 (0.059)	0.141** (0.056)
Observations	939	2,455	615	2,798	1,401	1,993
Mean of dep. var.	3.330	3.478	3.465	3.430	3.497	3.395
Std. dev.	0.873	0.773	0.742	0.818	0.758	0.833
Number of individuals	227	596	160	667	348	475

Notes The sample is restricted to individuals in the Immediate Incentives treatment and individuals in the Control group. All estimates include demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), baseline sleep, indicators for the number of classes starting at 10 a.m. or earlier, indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. Observations are weighted by the number of credits taken in the semester. Standard errors are clustered at the individual level. Mean of dep. var. is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p 0.01, ** p 0.05, * p 0.1.

Table A.9: Incentives and other metrics of academic performance

	(1)	(2)	(3)	(4)	(5)
	Has a grade	Withdrawn	Failed	Passed	Credits
Panel A: All classes					
Immediate Incentives	-0.011** (0.005)	0.009* (0.005)	-0.008 (0.005)	-0.003 (0.008)	0.025 (0.031)
Observations	4,772	4,772	4,772	4,772	4,772
Mean of dep. var.	0.982	0.0142	0.00964	0.972	2.755
Std. dev.	0.133	0.119	0.0977	0.164	1.002
Number of individuals	840	840	840	840	840
Panel B: Lectures					
Immediate Incentives	-0.014** (0.006)	0.010* (0.005)	-0.011* (0.006)	-0.003 (0.009)	-0.009 (0.033)
Observations	3,728	3,728	3,728	3,728	3,728
Mean of dep. var.	0.981	0.0161	0.0118	0.969	2.919
Std. dev.	0.138	0.126	0.108	0.174	0.874
Number of individuals	829	829	829	829	829

Notes The sample is restricted to individuals in the Immediate Incentives treatment and individuals in the Control group. All estimates include demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), baseline sleep, indicators for the number of classes starting at 10 a.m. or earlier, indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. Observations are weighted by the number of credits taken in the semester. Standard errors are clustered at the individual level. Mean of dep. var. is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p 0.01, ** p 0.05, * p 0.1.

Table A.10: Incentives and time use (in minutes), excluding careless respondents

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Sleep	Sleep \geq 7 hours	Study	Social	Work	Eating & Preparing Food	Exercise
Treatment	6.1684 (4.612)	0.0701*** (0.017)	0.9509 (7.163)	-4.3800 (4.855)	-6.2186 (5.351)	0.5336 (2.242)	1.2681 (1.783)
Post-Treatment	13.7295** (5.905)	0.0605*** (0.023)	6.8462 (9.689)	-3.7904 (6.807)	-3.5055 (7.129)	-5.3106* (2.713)	-3.6562 (2.529)
Observations	5,754	5,754	5,754	5,754	5,754	5,754	5,754
Mean of Dep. Var.	494.1	0.734	321.9	101.3	92.48	95.14	22.68
Std. dev.	108.8	0.442	192.2	123.3	144.2	51.17	42.87
Number of individuals	836	836	836	836	836	836	836
	(8)	(9)	(10)	(11)	(12)	(13)	(14)
	House errands	Personal care	Screen	TV & other videos	Internet	Games	Other
Treatment	-0.6851 (1.364)	0.8740 (1.635)	-11.4531** (4.989)	-8.2617** (3.238)	-2.7182 (3.166)	-0.0035 (2.171)	-0.5515 (7.302)
Post-Treatment	-1.5762 (1.789)	-1.2043 (2.321)	-3.2999 (6.784)	-4.4416 (4.435)	0.6341 (4.141)	0.9326 (3.297)	-11.1511 (9.555)
Observations	5,754	5,754	5,754	5,754	5,754	5,754	5,754
Mean of Dep. Var.	18.14	54.09	171.8	70.75	78.96	22.09	454.7
Std. dev.	36.85	40.27	136.7	92.15	90.09	63.62	171.5

Notes The sample is restricted to individuals in any of the Incentived treatments (Immediate Incentives, Delayed Incentives, and Delayed Incentives, No Cue/Feedback) and individuals in the Control group. Individuals in the Cue/Feedback treatment were not included in this analysis. We also exclude participants deemed careless (e.g. those who gave the reported “other activities” for all time periods within the last 24 hours) All the estimates include controls for wave, month and day of the week fixed effects, indicators for the number of classes starting before 10am, gender, race (dummies for Asian, Black, Hispanic, other) and ethnicity, parental education (dummies for less than college, college degree, and post-college degree), number of classes starting before 10 a.m., quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing), and the average time spent on the activity at baseline. Standard errors are clustered at the individual level. Mean of dep. var is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p 0.01, ** p 0.05, * p 0.1.

Table A.11: Immediate Incentives, cognitive performance and physical health

	(1)	(2)	(3)	(4)	(5)
	Correct math answer	Creativity score	RHR	# steps	Active minutes
Treatment	0.0024 (0.023)	0.0025 (0.057)	-0.2529 (0.186)	35.7563 (154.475)	-1.0652 (3.866)
Post-Treatment	-0.0308 (0.032)	0.0002 (0.059)	-0.1902 (0.205)	-73.3478 (237.144)	-7.5864 (6.125)
Observations	3,181	3,243	46,542	46,989	46,989
Mean of dep. var.	0.363	3.307	65.67	7161	191.9
Std. dev.	0.481	0.717	8.314	5756	140.8
Number of individuals	809	803	832	840	840

Notes The dependent variable in column 1 is an indicator equal to 1 if the respondent answered correctly the math question on the survey. The dependent variable in column 2 is a creativity score (see Section 2.4). RHR corresponds to participants Resting Heart Rate. Steps corresponds to participants' daily steps as measured via the Fitbit. Active minutes capture any activity at or above about 3 metabolic equivalents (METs). The sample is restricted to individuals in the Immediate Incentives treatments and individuals in the Control group. All estimates include day of the week, week of the experiment, wave, and month fixed effects, baseline value of the outcome variable, and demographic controls for gender, age (dummies), race and ethnicity (Asian, Black, Hispanic, White, other), indicators for the number of classes starting at 10 a.m. or earlier, indicators for whether parents' highest academic title was less than college, college degree, more than a college degree, and quartile of baseline GPA (high school GPA if non-missing, prior term GPA if high school GPA is missing). For all demographic characteristics, we included a missing indicator for whether the variable was missing. Standard errors are clustered at the individual level. Mean of dep. var is the mean of the dependent variable at baseline. Std. dev. is the standard deviation of the dependent variable at baseline. *** p 0.01, ** p 0.05, * p 0.1.

Table A.12: Review of post-secondary interventions

Paper	Treatment	Setting	Findings	GPA & Costs
Angrist et al. (2009)	A) Financial incentives for academic achievement B) Peer advising and study groups C) Treatments A and B combined	Field experiment with first-year students at Canadian 4-year university	A) GPA: -0.04 (0.061) B) GPA: 0.011 (0.063) C) GPA: 0.168 (0.086) Academic probation: -0.069 (0.036)	GPA: Table 6, Panel A, Column 1 Costs: Bottom of page 160
Angrist et al. (2014)	Financial incentives for academic achievement	Field experiment with students at public university in Ontario	GPA: 0.009 (0.044)	GPA: Table 4b, "Fall" Panel, Column 9 Costs: Table 3, "Fall" Panel, Column 9
Barrow et al. (2014)	Extra grant aid and counseling services as part of the Opening Doors Louisiana Program	Field experiment with low-income community college students in Louisiana	GPA: 0.182 (0.085) Credits: 1.234 (0.30)	GPA: Table 8, Column 2 Costs: Table 2, "First semester" panel, Column 1
Clotfelter et al. (2018)	Extra state grant aid due to crossing income threshold for Carolina Covenant Grant eligibility	Regression discontinuity with low-income students attending the University of North Carolina, Chapel Hill	GPA: 0.043 (0.053) 4-year degree ¹ : 0.068 (0.040)	GPA: Table 6, Panel A, Column 1 Costs: Table 3, Panel B, Column 1
Denning et al. (2019)	Extra Pell and state grant aid due to crossing threshold for 0 Expected Family Contribution	Regression discontinuity with 4-year university and community college students in Texas	GPA, FTIC ¹ : 0.031 (0.026) 4-year degree, FTIC: 0.022 (0.012) GPA, returning students: 0.014 (0.013)	GPA: Table 3, Panel B, Column 3 Costs: Table 2, Column 2

FTIC stands for "First Time in College"

Notes: Column 4 reports average treatment effects with standard errors in parentheses. We report OLS estimates of effects on non-cumulative GPA, either at the semester or year-level. We report multiple GPA effects when authors reported on multiple treatment arms (e.g. Angrist et al. (2009); Evans et al. (2020)) or cohorts (e.g. Goldrick-Rab et al. (2016); Denning et al. (2019)). We also report other statistically significant effects, such as credits completed or degree completion, when applicable. Effect sizes on "4-Year Degree" report the rate at which people receive a 4-year degree in 4 years, while "Credits" reports effect sizes on credits taken in one school year. When multiple GPA effects were reported, we selected semester-level estimates if available and used the authors' preferred specification if indicated. When per-person treatment costs were not reported, we divided overall program costs per semester by the intent-to-treat sample size (if total program costs were reported by year, we divided in half to calculate per semester costs).

Table A.12: Review of post-secondary interventions (continued)

Paper	Treatment	Setting	Findings	GPA & Costs
Evans et al. (2020)	A) Access to emergency grant funding B) Treatment A as well as advising services	Field experiment with low-income community college students in Texas	A) GPA: -0.134 (0.083) B) GPA: 0.055 (0.07) Enrollment, female students: 0.04 (0.041)	GPA: Treatment A provided by authors, Treatment B from Table 8, Column 2 Costs: Pages 958-959
Goldrick-Rab et al. (2016)	Extra grant aid as part of the Wisconsin Scholars Grant	Field experiment with low-income first-year students at public universities in Wisconsin	GPA, cohort 1: 0.08 (0.06) Credits, cohort 1: 0.9 (1.7) GPA, cohorts 2 & 3: 0.09 (0.03) Credits, cohorts 2 & 3: 2.1 (0.7)	GPA: Table 5, "First Semester" Panel, Columns 2 & 5 Costs: Bottom of page 1772
Oreopoulos and Petronijevic (2018)	A) Online exercise encouraging future-oriented thinking B) Treatment A as well as study advice and motivation via text messages C) Treatment A as well as one-on-one peer support	Field experiment with students at three campuses of the University of Toronto	A) Course grades: 0.143 (0.575) B) Course grades: 0.073 (0.505) C) Course grades: 4.897 (1.874) Credits: 0.501 (0.283)	Course grades: Table 3, Column 5 ¹ Costs: Bottom of page 323 ²
Park and Scott-Clayton (2018)	Extra Pell grant aid due to crossing threshold for 0 Expected Family Contribution	Regression discontinuity with community college students from 20+ institutions in a single state	GPA: 0.064 (0.082) Enrollment: 0.094 (0.034)	GPA: Table 5, Column 2 Costs: Table 5, Column 2
Scott-Clayton (2011)	Free tuition as part of the West Virginia PROMISE program	Regression discontinuity with public university students in West Virginia	GPA: 0.066 (0.066) Credits: 1.572 (0.085) 4-year degree: 0.058 (0.004)	GPA: Table 3, Column 3 Costs: Middle of page 617

¹ Authors present course grades on a 0-100 scale. In figure 6, course grades have been divided by 25 for comparison with 4.0 GPA scale.

² Only treatment arm C is included in figure 6 because costs could not be calculated for A and B.

B. Instructions and Experimental Material

Immediate Incentives

*****PLEASE READ THROUGH THIS MESSAGE ENTIRELY*****

Starting this Sunday, and every weeknight (Sunday-Thursday) for the next four weeks, we encourage you to get 7 hours of sleep or more by 9 am the following morning.

Every time you meet this goal (i.e., sleep 7 hours by 9 am), you will earn a 4.75 PAYMENT via Venmo. Payments are redeemable only until 3 pm on the days you earn them, and you will receive the payment by 3 pm if you have redeemed by that time.

HOW IT WORKS

Every morning, you will receive feedback on your sleep. If you meet your goal, you will also receive the payment information via text message.

Next, we would like to ask you to pick your bedtime behavior – a behavior you would like to engage on right before going to sleep. Every weeknight, we will remind you of your bedtime behavior and we will encourage you to go to sleep early enough to meet your goal of sleeping at least 7 hours by 9 am. Please pick your bedtime behavior by texting back the number of your choice. If you choose other, please type 9, then the behavior you want to set as your bedtime behavior.

1. Turn off your phone
2. Turn your phone to silent
3. Turn off your computer
4. Turn off Netflix
5. Turn on bedtime music
6. Turn on meditation app
7. Turn on white noise
8. Turn on pink noise
9. Other

Delayed Incentives

*****PLEASE READ THROUGH THIS MESSAGE ENTIRELY*****

Starting this Sunday, and every weeknight (Sunday-Thursday) for the next four weeks, we encourage you to get 7 hours of sleep or more by 9 am the following morning.

Every time you meet this goal (i.e., sleep 7 hours by 9 am), you will earn a 4.75 PAYMENT via Venmo. Payments are redeemable only until 3 pm on the days you earn them, and the payment will be added to the amount of money you receive at THE END OF THE STUDY.

HOW IT WORKS

Every morning, you will receive feedback on your sleep. If you meet your goal, you will also receive the payment information via text message.

Next, we would like to ask you to pick your bedtime behavior – a behavior you would like to engage on right before going to sleep. Every weeknight, we will remind you of your bedtime behavior and we will encourage you to go to sleep early enough to meet your goal of sleeping at least 7 hours by 9 am. Please pick your bedtime behavior by texting back the number of your choice. If you choose other, please type 9, then the behavior you want to set as your bedtime behavior.

1. Turn off your phone
2. Turn your phone to silent
3. Turn off your computer
4. Turn off Netflix
5. Turn on bedtime music
6. Turn on meditation app
7. Turn on white noise
8. Turn on pink noise
9. Other

Cue / Feedback

PLEASE READ THROUGH THIS MESSAGE ENTIRELY

Starting this Sunday, and every weeknight (Sunday-Thursday) for the next four weeks, we encourage you to get 7 hours of sleep or more by 9 am the following morning.

HOW IT WORKS

Every morning, you will receive feedback on whether you met your goal.

Next, we would like to ask you to pick your bedtime behavior – a behavior you would like to engage on right before going to sleep. Every weeknight, we will remind you of your bedtime behavior and we will encourage you to go to sleep early enough to meet your goal of sleeping at least 7 hours by 9 am. Please pick your bedtime behavior by texting back the number of your choice. If you choose other, please type 9, then the behavior you want to set as your bedtime behavior.

1. Turn off your phone
2. Turn your phone to silent
3. Turn off your computer
4. Turn off Netflix
5. Turn on bedtime music
6. Turn on meditation app
7. Turn on white noise
8. Turn on pink noise
9. Other

Creativity Instructions Example)

You will be asked to complete different short tasks over the course of the study. One of these tasks will be chosen for payment at the end of the study.

Today's task: Using some or all of the words below, write an interesting sentence. Your sentence will be rated based on its creativity from 1-5 points, where 5 is the most creative. If today's task is chosen for payment, your payment will be determined by how creative your sentence is. **You will receive 1 for each point your story is rated.** You will receive as little as .1 for completing this activity and up to .5 for the most creative sentences. You will receive your rating and your payment at the end of the study.

The words for you to use in your sentence are:
(Example) event, chocolate, system, indicate, article, emotion, possess, mom, poetry, reality

Math Instructions Example)

You will be asked to complete different short tasks over the course of the study. One of these tasks will be chosen for payment at the end of the study.

Today's task: On the next page you will be asked to answer a math question. If today's task is chosen for payment, your payment will be determined by whether you answer the question correctly, and how quickly you answer. **You will receive 1 for answering the question correctly, and you will receive an additional 0- 4 depending on how quickly you answer the question.** You will receive as little as .1 for answering this question correctly and up to .5 for the quickest correct answers. You will receive your score and the payment at the end of the study.

Here is the question you are asked to answer:

It costs a manufacturer X dollars per component to make the first 1,000 components. All subsequent components cost .1 each. When $X = 1.50$ How much will it cost to manufacture 4,000 components?

- 3,500
- 3,000
- 4,000
- 3,250
- 4,500

App ScreenShots

Figure B.1: Bedtime Reminder



Notes The Bedtime reminder included a personalized goal bedtime of approximately 1 hour before the baseline bedtime, with a latest possible time of 1 am. It also included a personalized bedtime behavior participants chose from before the beginning of the intervention, from a list containing "Turn off your phone", "Turn your phone to silent", "Turn off your computer", "Turn off Netflix", "Turn on bedtime music", "Turn on meditation app", "Turn on white noise", "Turn on pink noise", "Other". If participants selected "Other" they could specify a behavior of their choice.

Figure B.2: App Screenshots - Immediate Incentive Treatment

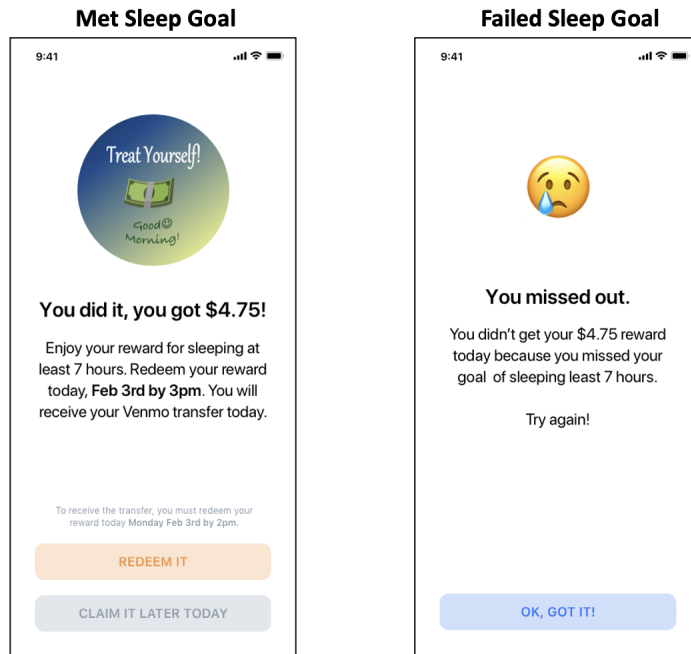


Figure B.3: App Screenshots - Delayed Incentive Treatment

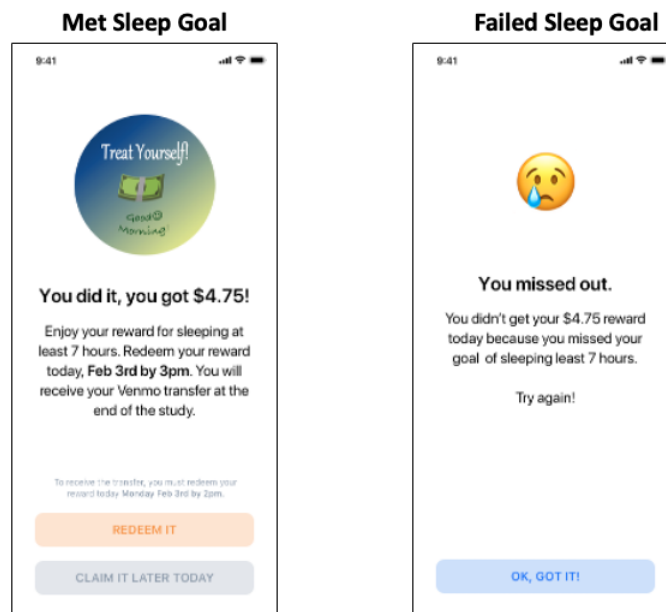


Figure B.4: App Screenshots - Cue/Feedback Treatment

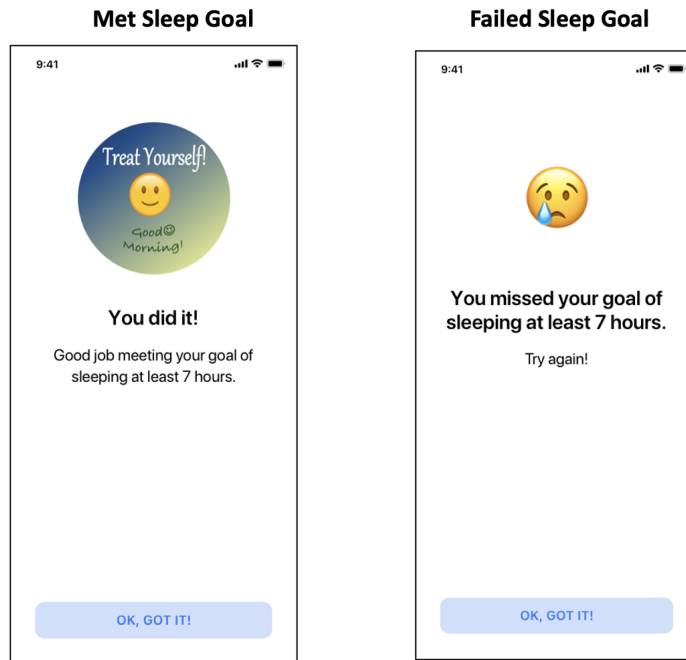


Figure B.5: Reminder to Sync - All Treatments

