# Nurturing Childhood Curiosity to Enhance Learning:

# Evidence from a Randomized Pedagogical Intervention \*

Sule Alan<sup>†</sup> and Ipek Mumcu<sup>‡</sup>

January 19, 2023

#### Abstract

We evaluate a pedagogical intervention that aims to improve the learning quality of elementary school children by nurturing their curiosity. We test the effectiveness of the pedagogy using achievement scores and a novel measure of curiosity. The latter involves first creating a sense of information deprivation, then quantifying the urge to acquire information and the ability to retain information. The intervention increases curiosity, the ability to retain knowledge, and science test scores. It also leads to more efficient information dissemination in the classroom. The evidence can help design better pedagogical tools to increase pupil engagement and the quality of learning.

Keywords: Learning crisis, Curiosity, Deep Learning, Pedagogy, Achievement

JEL Codes: I24, I26

<sup>\*</sup>We are grateful to J-PAL Post-Primary Education Initiative and ING Bank Turkey for funding this study. We thank seminar participants at EUI, LSE, Bologna ESA 2022, the University of Michigan, and the University of Toronto for their valuable comments. We thank Enes Duysak, Mert Gumren, Canan Guner, Elif Kubilay, Ozge Seyrek, Fatima Silpagar and Melek Celik for wonderful field assistance and Selin Iplikci for excellent research assistance. The study has ethics approval from the University of Essex Ethics Board. Study 1 AEA Registry: AEARCTR-0003957. Study 2 AEA Registry: AEARCTR-0008629.

<sup>&</sup>lt;sup>†</sup>Corresponding Author: European University Institute, Florence, Italy

<sup>&</sup>lt;sup>‡</sup>University of Exeter, United Kingdom

#### 1 Introduction

Today, more children than ever enroll in primary and post-primary education in the developing world. Despite this progress, the quality of education remains low. Millions of children in developing countries leave school without the necessary foundational skills to help them achieve their potential and lead productive lives. Low teacher quality, overcrowded class-rooms, and inadequate levels of school inputs such as poorly designed curricula and insufficient teaching materials are among the many factors contributing to low learning outcomes (Glewwe and Muralidharan (2016); Glewwe, Lambert and Chen (2020)). Recent research highlights the role of pedagogy as a potentially effective policy tool to combat poor education quality. While there is no consensus on what constitutes good pedagogy, teaching practices that respond to the needs of students at all levels, build on their individual strengths, and encourage them to learn through experimentation are likely to be effective. Unfortunately, most traditional instruction techniques lack these features. They ignore heterogeneous learning paths, compel students to be passive listeners, and prevent the development of an active and inquisitive mind (Blanchard, Southerland and Granger (2009); Granger et al. (2012); Terrenghi et al. (2019); Ashraf, Banerjee and Nourani (2021)).

In this paper, we evaluate the effectiveness of a pedagogical program that aims to nurture children's curiosity and improve learning outcomes. Motivated by the recent evidence on the neural mechanisms of human curiosity and its connection to deep learning, the pedagogy strives to cultivate children's natural urge to learn and explore novel phenomena. The pedagogy primarily targets scientific curiosity and is practiced by teachers throughout an academic year. Treated teachers receive extensive training on practical ways to stimulate curiosity in the classroom by leveraging children's natural love of mystery, humor, and astonishment. Teachers are provided with a pedagogical toolkit containing various visual and reading materials to support the prescribed practices. The toolkit also offers ideas for creating teachable moments and holding students' undivided attention before the introduction of

<sup>&</sup>lt;sup>1</sup>According to 2017 Annual Status of Education Report for India, about 25% aged 14-18 fail to read basic text fluently in their language, 57% struggle with division (three digits by one digit) (ASER (2018)). Results from similar tests in Pakistan and East Africa paint a similar picture. PISA and TIMMS results highlight large learning gaps between the developing and the developed world (Gust, Hanushek and Woessmann (2022)). The recent study by Singh, Romero and Muralidharan (2022) documents the further damage done to learning by the Covid-19 Pandemic.

<sup>&</sup>lt;sup>2</sup>For example, tailoring the level of teaching to children's ability has been shown to be effective in helping those who lag behind to catch up. (Banerjee et al. (2007); Banerjee et al. (2016); Banerji and Chavan (2016)).

new and complex curricular items. While this pedagogy is relevant for any curricular topic, the toolkit predominantly targets scientific curiosity.

Curiosity, a fundamental component of human cognition, is considered a critical driver of success in most aspects of life. Berlyne (1954) and Loewenstein (1994) provide a theoretical framework for epistemic curiosity, described as "desire for knowledge". Cognitive psychology associates curiosity with achievement in many domains ranging from education to health and overall life satisfaction (Chamorro-Premuzic and Furnham, 2006; Kashdan and Silvia, 2009; von Stumm, Hell and Chamorro-Premuzic, 2011; Gottfried et al., 2016; Shah et al., 2018). Recent advances in neuroscience shed light on the neural mechanisms of curiosity and its links to learning. Gruber, Gelman and Ranganath (2014) show via functional magnetic resonance imaging that the brain's reward system is evoked when people are curious about a phenomenon. This facilitates more enjoyable learning and knowledge retention (deep learning) through memory consolidation. Moreover, they show that once sparked, curiosity creates deep learning moments and enhances the learning of any topic, not only the topic that sparked curiosity initially. While recognized as a powerful motivator for learning, curiosity has not been studied on a large scale within the context of education policy. Our limited knowledge of how to cultivate such a context-dependent trait, together with the difficulty of measuring it, are obvious reasons for the lack of policy-relevant studies. We know of no large-scale study that measures curiosity in school children, nor any study that shows how to stimulate it in the school environment. This paper advances the literature on both of these fronts.<sup>5</sup>

The pedagogical program was implemented as two independent randomized controlled trials in two large southern provinces of Turkey. The first trial, implemented in the 2018-2019 academic year in the province of Mersin, included 50 primary and 27 post-primary schools. The post-primary schools were later dropped from the program as the proposed pedagogy was discovered to be more suitable for the primary school level. To strengthen the power of

<sup>&</sup>lt;sup>3</sup>Throughout the text, the word curiosity refers to epistemic curiosity, distinguishing human curiosity from animals'. Loewenstein (1994) (and references therein) describes curiosity as reacting positively to new or mysterious events by showing the urge to explore and understand them. Philosopher William describes curiosity as the "impulse towards better cognition" (James (1983)).

<sup>&</sup>lt;sup>4</sup>Memory consolidation is a process by which acquired information or experiences are poured into long-term memory. It is more likely to happen when stimuli spark curiosity; see Gruber and Ranganath (2019).

 $<sup>^5</sup>$ Recently, psychologists have shown interest in the relationship between what they refer to as "epistemic emotions" and learning. Epistemic emotions include intellectual courage, astonishment, curiosity, interest, wonder, surprise, the joy of verification, and the satisfaction of knowing; These studies are correlational in nature; See Vogl et al. (2019b) and Vogl et al. (2019a).

the study, we re-implemented the program in the neighboring province, Adana, recruiting 84 additional primary schools. This second study took place in the 2021-2022 academic year. Our combined sample includes 134 primary schools with about 11,000 students and 425 teachers. After collecting detailed baseline data from all the children and teachers in Fall 2018 (Study 1) and Fall 2021 (Study 2), we randomly assigned 78 schools to treatment (25 in Study 1, 43 in Study 2). Teachers from the selected schools received training in the prescribed pedagogy and were given the entire academic year to practice it in the everyday teaching of the curricular topics, with a greater emphasis on science lessons. We collected our endline data in May 2019 (Study 1) and May 2022 (Study 2) to test the effectiveness of the pedagogy using objective test scores, educational aspirations, and a novel measure of curiosity. When we implemented the second study in the Fall of 2021, we also collected longer-term data from the first study subjects (about three years after the program implementation in 2018).

Curiosity is challenging to measure due to its context-dependent nature. Psychologists use survey tools to elicit different types of curiosity in adults (Litman and Spielberger (2003); Collins, Litman and Spielberger (2004); Litman, Collins and Spielberger (2005); Kashdan et al. (2020)). Behavioral tasks are used for very young children (Jirout and Klahr (2012)). While self-reports via item-response questions can be useful to measure curiosity in adults, such questions usually have low reliability when implemented on children and adolescents. The lack of a reliable measure of curiosity in school children motivated us to develop an incentivized task-based tool. Our tool benefits from the theoretical framework developed by Loewenstein (1994), and it draws insights from neuroscientific research on curiosity. Our curiosity measure involves first creating a sense of information deprivation, then quantifying the urge to acquire information, and finally, assessing the degree of knowledge retention after satisfying the urge.

To develop the task, we first conducted extensive pilot surveys to determine the topics of interest of our target age group. We identified eight broad interest categories representing about 95% of all topics reported by children. These are, "science", "animals", "history", "human anatomy', "vehicles", "cartoons", "space", and "sports". We then prepared eight booklets with these titles and, in each one, inserted ten surprising facts that neither a child nor an adult would be likely to know. The implementation began by showing children each booklet and telling them that they contained facts that most people do not know. This step aimed to create a strong urge to know in the children. After recording their preferred booklets, we elicited the children's willingness to pay for them. For this, we first endowed children with experimental tokens that could be converted into small gifts of value. We then

asked them to state the highest number of tokens they would be willing to sacrifice for their preferred booklet, including the option of zero tokens. Our measure of the urge to acquire information—that is, of their curiosity— is the child's willingness to pay for their preferred booklet.<sup>6</sup> The novelty of our task lies in its temporal component. After eliciting the urge to acquire information and distributing booklets, we revisited all classrooms exactly one week later, unannounced. During this surprise visit, we gave the students and the teachers a 40-question test containing questions whose answers are found in the booklets. The performance in this test is our measure of knowledge retention (deep learning).

To identify the effect of the treatment on knowledge retention, we implemented the curiosity task during the first visit using two regimes. In classrooms that are randomly assigned to the first regime, children received their preferred booklets based on a randomly determined market price. In classrooms assigned to the second regime, only half the children in a given classroom received booklets. In these classrooms, the booklet distribution was purely random, regardless of children's willingness to pay and their choice of booklets. Ensuring that the proportion of students receiving booklets and the composition of topics are balanced across treatment status, the second regime allows us to estimate the treatment effect on knowledge retention. This regime also allows us to show the extent to which treatment improves information sharing and peer learning within the classroom.

We first document that our measure of curiosity, the willingness to pay for a booklet, correlates well with fluid IQ, performance on the retention test, and actual test scores (crystallized IQ). We then estimate the effect of the program on curiosity, knowledge retention, test scores, and educational aspirations. We find that the program significantly increases children's curiosity by about 0.11 standard deviations (0.35 extra tokens forgone). The effect of the program on scientific curiosity (the willingness to pay for science-related booklets) is similar in size (0.10 standard deviations) and precision. Treated children choose to give up 0.39 more tokens than untreated children for a science-related booklet on average, implying a 12.2% increase in the willingness to pay for scientific information. The effect of the intervention on knowledge retention is striking. Treated children score about 0.11 standard deviations higher in the unannounced booklet test we conducted 1 week later. Even more striking is that after about three years, including 1.5 years of school closure due to the recent pandemic, treated students score 0.14 standard deviations higher on the same booklet test

<sup>&</sup>lt;sup>6</sup>Willingness to pay elicitation is a standard method in economics research. In the context of information as a good, Hjort et al. (2021) uses this method to elicit policy-makers' willingness to pay for evidence in Brazil.

than untreated children, indicating a remarkably persistent treatment effect on knowledge retention.

We also provide evidence that the program makes friendship networks more effective information dissemination tools. Treated students who did not receive a booklet and whose preferred booklet was received by someone else in their friendship network scored 0.18 standard deviations higher on the booklet test than untreated students in the same condition. We also show that as the information availability increases within friendship networks, treated students exhibit significantly higher knowledge retention than untreated students. These results strongly suggest more efficient peer learning technology and information dissemination in treated classrooms where students are more curious and passionate about pursuing and sharing knowledge. The improved peer learning is also consistent with the recent evidence that human curiosity is sensitive to the social environment and stimulated by the curiosity of others (Dubey, Mehta and Lombrozo (2021)).

The positive effects of the program on curiosity extend to actual learning outcomes. The program significantly improves children's objective test scores in science with no statistically significant impact on math and verbal scores. The estimated effect size on science test scores is about 0.08 standard deviations in the short term. We find that the positive effect on science test scores persists into middle school years, even after a long school closure due to the Covid-19 pandemic. Treated students score 0.07 standard deviations higher than untreated students in a science test covering the middle school curriculum. Finally, we show that the intervention significantly raises children's aspirations to go to university and study science. While persistent in size, these effects are less precisely estimated in the long run.

Our results suggest that the program's success is likely to stem from its ability to unleash children's curiosity. We show that the program also increases children's tolerance for uncertainty and makes them more critical in their thinking process. In addition to children's outcomes, we estimate significant program effects on teachers' teaching styles and beliefs. Treated teachers report a significant increase in their own curiosity level. They also report adopting a more modern, learner-centered approach to teaching and embracing a growth mindset. By testing teachers' curricular content knowledge at endline, we rule out the possibility of treatment improving learning outcomes by improving teachers' ability.

Our contribution is fourfold. First, we evaluate a pedagogical intervention that targets a component of human cognition, curiosity, that has not been studied on a large scale and in a policy-relevant context before. Second, leveraging the neuroscientific evidence on curiosity

and memory consolidation, we offer a novel approach to measuring curiosity in primary school children. Third, combining the two, we provide evidence to support the causal link between childhood curiosity and deep learning in a natural field setting. We show that once sparked, curiosity leads to enhanced knowledge retention in children. Finally, our paper is the first to establish the learning externalities generated by human curiosity. We show that a pedagogical approach aimed at nurturing students' curiosity can lead to more information sharing and peer learning in the classroom. Therefore, the results of the paper are of high policy relevance. They can help us design better pedagogical tools to increase pupil and teacher engagement and the quality of learning worldwide. The results are particularly relevant for the developing world, where learning outcomes have been alarmingly low and have deteriorated even further due to the Covid-19 pandemic (Goldhaber et al. (2022)).

Our paper relates to several strands of the economics literature. First, by showing the effectiveness of a particular pedagogical approach, it contributes to the literature that strives to improve learning outcomes in developing countries. This literature establishes that schoolbased inputs have very little effectiveness when not complemented by correct teaching practices (Glewwe et al. (2004); Kremer, Glewwe and Moulin (2009); Kremer, Brannen and Glennerster (2013)). Related literature explores the ways to improve teacher motivation and engagement and shows that extrinsic motivations have limited effectiveness in improving learning outcomes (de Ree et al. (2018)). Second, the paper also relates to a growing literature that shows that social and emotional skills are likely malleable and can be fostered at young ages (Alan and Ertac, 2018; Alan, Boneva and Ertac, 2019; Alan et al., 2021). By showing that an important trait can be cultivated in the classroom through a change in teaching practices, we advance this literature. Third, by testing a pedagogy that focuses mainly on science teaching, the paper speaks to the literature that aims to increase the STEM participation of girls (Buser, Peter and Wolter (2017); Fischer (2017); Kahn and Ginther (2017); Carlana and Fort (2022)). Finally, we contribute to neuroscience and psychology literature in their efforts to understand the implications of human curiosity by causally linking the desire for knowledge with deep learning in children in a large field setting.

The rest of the paper is organized as follows. Section 2 summarizes the key features of the program and the context in which it was implemented. Section 3 details the evaluation design and gives a detailed account of our outcome measures, including our task-based curiosity measure. Section 4 describes the data and presents our main results. In Section 5, we explore mechanisms through which the program improved knowledge retention and achievement outcomes. We conclude in Section 6.

# 2 Evaluation Context and The Nature of The Pedagogical Program

The program we evaluate has been developed by an expert team of pedagogy specialists and curricula developers in a private university's innovation center. The program's overarching objective is to promote scientifically informed teaching practices to improve learning outcomes. It aims to do so by replacing traditional teaching with techniques that can ignite children's enthusiasm for academic matters. This is especially pertinent in light of the global push for STEM education and better outcomes in science. As such, the program puts a greater focus on the teaching of science.

The Turkish primary school system is designed such that a centrally appointed teacher is assigned to a single classroom in Grade 1 and is expected to teach the same pupils until the end of Grade 4, after which they move on to middle school for Grades 5 to 8 where each subject is taught by a different (branch) teacher. The program has been developed to exclusively benefit primary school teachers, as it is thought that the ideal context for implementing the prescribed pedagogy would be when a single teacher has a full day of contact with their pupils and when science concepts are formally introduced. Such a context is grade 4 of primary school in Turkey.

The intervention was an intensive teacher training program. In training seminars, teachers were first introduced to the concept of curiosity as a fundamental driver of academic achievement. Then, teachers were introduced to various pedagogical practices to cultivate curiosity in the classroom environment. These practices included ways to allow students to suspect and inquire, as well as encourage them to express their interests openly in the classroom. At the core of these practices was the idea of tapping into children's natural inclination for mystery, surprise, and humor, in order to grab their attention and create teachable moments.

Teachers received a toolkit containing visual and written material to help them practice the pedagogy. These materials are not meant to be a set of materials to be covered in a specified period of time. Rather, they are designed to help the teacher create teachable moments using emotional triggers to hold students' undivided attention before she introduces

<sup>&</sup>lt;sup>7</sup>While this is the general practice, there are many exceptions to this rule. Firstly, the headteacher can decide which grade level the newly appointed teacher should begin teaching based on the needs of the school. Secondly, the Ministry can re-appoint a teacher, voluntarily or involuntarily, to another school at any grade level. These rotations tend to occur frequently for early career teachers.

a new and complex topic. For example, before introducing a science topic on the solar system, which is an official curricular item to be covered, students see a short video on the mysteries of space. The video is designed to capture students' attention, tapping into their love of mystery to create a teachable moment. As another example of creating a teachable moment, this time, using humor, the teacher reads a funny story about a girl who gets excited about exploding liquids before introducing a topic on chemical reactions. While most activities are related to science, the toolkit contains some non-science activities as well. For example, in one of the activities, students read about a fictional student with a deep interest in painting using unconventional tools (finding making a mess with raw eggs liberating). Teachers worked on the toolkit and repeatedly practiced different ways of creating teachable moments during the training seminars with the guidance of education consultants.

The overall feedback from the teachers regarding the program content was extremely positive. The majority of teachers reported that the program made everyday teaching, not just science teaching, much more enjoyable for children and for themselves. We received reports and visuals from many treated teachers showing their innovative ways of creating teachable moments. Bringing a mysterious box to the classroom that contains valuable information on the layers of the earth, hiding an important piece of information about the phases of matter in the teacher's hair, and hanging the names of the planets in our solar system around an umbrella; are just a few examples. See the Online Appendix C for examples of implementation photographs we received from teachers.

# 3 Evaluation Design and Outcomes of Interest

The program was implemented as two independent randomized trials three years apart. The first trial took place in the 2018-2019 academic year covering 50 primary and 27 post-primary schools in the province of Mersin (Study 1). Despite the program's target grade being 4, with the recommendation of the local authorities in Mersin, we decided to test the program also in the first year of the post-primary context by including a sample of 5th graders and their science teachers. However, it became clear during the training phase that the prescribed pedagogy would be difficult to implement in a middle school setting. Invited middle school teachers expressed their concerns regarding the larger number of pupils per classroom and the more demanding nature of the national curriculum relative to primary schools. As a result, middle schools were removed from the study. This resulted in a loss of 27 schools,

<sup>&</sup>lt;sup>8</sup>All these topics are part of the 4th-grade Turkish national science curriculum.

leading to a second trial in the 2021-2022 academic year to improve the power of the design. The second trial covered 84 primary schools in the neighboring province, Adana.<sup>9</sup>

In both trials, local authorities provided us with a list of schools located in socioeconomically deprived neighborhoods in their provinces. Teachers from these schools were offered participation in the program. The program participation was voluntary on the part of teachers. The program was oversubscribed in both provinces. Due to the large size of Turkish state schools, which generally have multiple classrooms for each grade level, 2 to 7 classrooms were selected randomly from each school for evaluation purposes. Two trials pooled together provide us with about 11,000 students and 425 teachers from 134 state primary schools in two large provinces of Turkey. The majority of our sample is composed of 4th graders. We also have some third-grade students in our first study sample.

The timeline of each trial is shown in Figure 1. We collected baseline data for Study 1 in October 2018, followed by randomization at the school level, stratified by district and grade level. The probability of treatment was 50%, assigning 25 schools to treatment and 25 to control in Study 1. Teacher training seminars for Study 1 took place in November 2018, and short-term endline data were collected in May 2019. We collected baseline data for Study 2 in October 2021 and conducted the randomization at the school level, in the same manner, stratifying by district<sup>12</sup>. The probability of treatment was again 50%, assigning 43 schools to treatment and 41 to control in Study 2. Teacher training seminars for 43 treatment schools took place in October 2021. Short-term endline data were collected in May 2022 for this study.

<sup>&</sup>lt;sup>9</sup>We first launched the second trial in 2019-2020 but failed to implement and evaluate it due to the Covid-19 related school closures, which lasted about 1.5 years in Turkey. We launched the second trial again as soon as schools opened in Fall 2021.

<sup>&</sup>lt;sup>10</sup>Primary school sizes vary significantly in our sample, ranging from remote village schools with a single 4th-grade class to overcrowded urban schools with over 15 classrooms for each grade level.

 $<sup>^{11}</sup>$ We admitted a small number of grade 3 classrooms in the first study, comprising about 16% of the sample in this study. This is because we received an overwhelming interest from these teachers and admitted them to the program.

<sup>&</sup>lt;sup>12</sup>We managed to limit our sample to 4th graders in the second study.

Study 1

Consider 2016

November 2018

November 2018

November 2019

Consider Data
Collection

Fracher
Training

Implementation
by taschers

Implementation
by taschers

Study 2

Figure 1: Timeline of the Two Trials

Both baseline and endline data collection were carried out by the research team, assisted by locally recruited and trained field assistants. We made sure that teachers were not present in classrooms during data collection. At baseline, we spent about three lecture hours in each classroom to conduct incentivized games, achievement and psychometric tests, and surveys. We implemented our behavioral curiosity task only at endline. Because of the temporal nature of the task, we organized two visits for each classroom at endline, one week apart. On the first visit, we spent about two lecture hours implementing the curiosity task and collecting other relevant data using tests and surveys. Our second visit was an unannounced surprise visit, which is why our task was implemented only at endline. Upon arrival at the school on the second visit, we kindly asked the teacher to spare us one lecture hour to implement a couple of tests on students and themselves. We will explain the nature of our curiosity task and the tests we implemented later in the text.

In October 2021, almost three years after the first implementation of the program in Mersin (Study 1), we managed to conduct another round of data collection for Study 1. Locating the original subjects of the first study was challenging. While most students were scattered around various middle schools in the same province, some had left the province or left the education system altogether. We eventually located 86% of our original participants with the help of the provincial authority's database. Among those, 84% was formally registered in a state middle school in the province, giving us 72% of our original sample. The attrition is more likely for girls and refugees, exacerbated by the extended school closures

due to the Covid-19 pandemic, but balanced across treatment status (p-value=0.66)<sup>13</sup>. Both trials were registered at the AEA Registry before their respective endline dates. The first trial was registered on March 8, 2019, along with a pre-analysis plan. The second trial was registered on November 30, 2021, referring to the first registry for the PAP.

Next, we will explain our curiosity task and the way we implement it in the classroom at endline, followed by descriptions of other outcomes of interest.

#### 3.1 A Task-based Approach to Measuring Childhood Curiosity

We offer a novel incentivized task to measure curiosity and use it as our primary behavioral outcome. We designed this task to capture two prominent aspects of human curiosity: the urge to acquire knowledge and the retention of the acquired knowledge upon satisfying the urge. We benefit from the conceptual framework developed by Loewenstein (1994) for the first component. Based on this framework, we first create a sense of information deprivation in children and then quantify the degree of the urge to acquire information. The second component of our task is informed by the neural mechanisms of curiosity documented in Gruber, Gelman and Ranganath (2014). That is, the higher the urge to know, the stronger the knowledge retention upon satisfying the urge (memory consolidation).

To develop the task, we first conducted extensive pilot surveys and qualitative interviews in several out-of-sample schools to determine the interests of the target age group. Compiling all our survey responses, we identified eight interest categories representing about 95% of all topics of interest. These are, "science", "animals", "history", "human anatomy", "vehicles", "cartoons", "space", and "sports". We then prepared eight small booklets for each topic with a cover that clearly shows the above titles. For example, the cover of the space booklet reads "The mysteries of SPACE," with eye-catching space illustrations to create information deprivation. Figure A1 shows the covers of all eight booklets. We placed in each booklet exactly ten pieces of information that are surprising and highly unlikely to be known by children (or by adults). Examples include, "the color of dawn on Mars is blue" in the space booklet, "the actual color of the black box in planes is orange" in the vehicles booklet, or "the shortest battle in history took 38 minutes" in the history booklet.

<sup>&</sup>lt;sup>13</sup>Both provinces have a significant refugee population, and all refugee children are covered under the MoE-EU refugee school placement program. However, Turkey's refugee population is highly mobile and difficult to track as they tend to be agricultural laborers. We provide a detailed attrition pattern for Study 1 in Figure B1 in the Online Appendix. A notable number in this figure is 520 missing children the Ministry lost track of in the pandemic period.

The implementation of the task in a classroom follows the following steps: We arrive at the classroom with booklets and a basket full of small gift items. The latter are small stationery items that are of value to children of the socioeconomic group we target in this study. We present the booklets to the children one by one, showing the title cover. We tell them that each booklet contains some incredible facts that are unknown to most people. This step aims to create information deprivation (a strong urge to know) in children. We then ask children to rank these booklets according to their interest in the topic, 1 being the most interesting and 8 being the least interesting.

After obtaining their ranking, we inform children that everyone has an endowment of 10 tokens, and each token can be converted into a gift from our gift basket. We show children these gift items one by one. We then tell them they can also use their tokens to purchase a booklet if they want to. For this, they first need to state the booklet they would like to purchase by ticking the relevant box. We emphasize that they do not have to buy a booklet if they do not want to. In practice, children see 9 options on their screen, 8 topics, and an option of "I don't want a booklet". Then, we begin explaining how this purchase will be made in practice. We first emphasize that all booklets have the same price, and each student can only buy one booklet. We tell them no one knows the price of a booklet yet, but they need to state their willingness to pay for their preferred booklet, using the options ranging from zero to 10. Then we explain to children that one of two things will happen in their classroom:

- Market price implementation regime: In this regime, we randomly choose a booklet price (between 1 and 10) for the classroom. Students whose willingness to pay falls under the revealed market price do not receive their desired booklet. They, therefore, convert all their tokens into gift items. Those whose willingness to pay is at or above the revealed market price receive their desired booklets at the market price and convert their remaining tokens into gift items.
- Half-half implementation regime: In this regime, we do not choose a market price for the classroom. Instead, a random half of the classroom receives booklets and all 10 tokens worth of gift items, regardless of their stated willingness to pay and the type of booklet they prefer. The other half of the classroom receives 10 tokens worth of gift items but no booklet. We explain the rationale behind this implementation regime below.

After providing this information and ensuring they fully understand the task, we ask children to state their willingness to pay for their desired booklet with utmost secrecy by tapping the relevant box on their tablet. The elicited willingness to pay, ranging from zero to 10, is our measure of "the urge to know," i.e., curiosity. This measure is theoretically independent of the implementation regime, and our data corroborates this: Mean willingness to pay across regimes is statistically not different from each other (p-value=0.48). We conjecture that the treatment will increase children's willingness to pay for information on their preferred topic. Given the program's heavy focus on science, we expect this effect to be particularly prominent in the willingness to pay for science-related booklets, which we refer to as "scientific curiosity." These booklets are science, space, human body, animals, and vehicles.

In addition to measuring the urge to acquire knowledge, we measure actual knowledge retention using the temporal component of our task. For this, we re-visit all classrooms, unannounced, precisely one week later. In this surprise visit, we give children a 40-question multiple-choice test containing 5 questions from each booklet. The score from this test is our measure of knowledge retention. However, when measured under the market price regime, the booklet test score has two confounds. First, if the program increases children's overall curiosity, measured as the willingness to pay, we expect more students in treatment classrooms to have access to booklets under the random market price regime. This differential availability of knowledge is the first confound in measuring retention, as more availability likely leads to more knowledge mechanically. Similarly, if, say, science topics are more popular in treatment classes, more science booklets will be available in treated classrooms rendering differential availability of science-related knowledge, the second confound.

In the half-half regime, a random half of the children in each classroom receive randomly chosen booklets, regardless of their willingness to pay and their preferred booklet. Therefore we eliminate both confounds under this regime. By making the amount and the type of information available independent of the treatment status, this regime (and only this regime) allows us to identify the program's impact on knowledge retention. In Study 1, a given classroom had a 50% chance of being subject to either regime, and children were informed accordingly. Because the causal effect of the treatment on information retention can be estimated only in the half-half regime, to improve the power of the experimental design,

<sup>&</sup>lt;sup>14</sup>To do this, we arrived back at schools and kindly asked their permission to take one lecture hour immediately. We gave the same test to teachers and asked them to do their own tests in a quiet, designated room. All our teachers cooperated.

we implemented the half-half regime in most classrooms (95%) in the second study, and children were informed accordingly. When implementing this regime, we made sure that every classroom had all 8 booklets. The Online Appendix D gives full instructions for the task and its implementation.

#### 3.2 Learning Outcomes and Educational Aspirations

If the program successfully stimulates students' curiosity, we expect deeper learning of curricular topics as well. In particular, given the program's heavy emphasis on science teaching, we expect treated students to achieve higher test scores in science. To assess the impact of the program on actual learning outcomes, we implemented tests on math, Turkish (in visit 1), and science (in visit 2) in all classrooms. Because there is no standardized testing system in Turkey for the grade levels we work with, we designed a testing inventory based on the national curriculum<sup>15</sup> All tests were implemented in classrooms in the absence of teachers.

In addition to learning outcomes, we assess whether the program affected children's educational aspirations and their plans for study majors. For this, we asked children whether they would like to go to university, and if so, what their aspired topic of study would be. We acknowledge that this is not a reliable measure of major choice considering the age of our subjects. Nevertheless, we believe that it gives us an indication of the program's success in raising educational aspirations in children.

Our long-term inventory was shorter than our short-term inventory because of the constraints imposed by the middle school schedules. We first gathered our students in designated classrooms in their middle schools. Then we gave them the same 40-question booklet test to assess the persistence of our knowledge retention results, followed by math, science, and verbal tests. The last three tests were prepared based on the appropriate grade level covering the national curricula. Finally, we conducted a short survey that elicited curiosity, grit, and aspirations.

#### 4 Data and Results

We collected data on various cognitive and non-cognitive skills, beliefs, and preferences at baseline and endline. Children's fluid IQ was measured using Raven's progressive matrices

<sup>&</sup>lt;sup>15</sup>We benefited from the Ministry's question bank in preparing these questions. We extensively piloted the tests to ensure the appropriateness of the difficulty level.

(Raven and Court (1998)) only at baseline. We conducted standardized achievement tests and elicited risk and ambiguity attitudes using Gneezy and Potters (1997) risky investment task, both at baseline and endline. We collected items via surveys to construct measures of epistemic and scientific curiosity (Kashdan and Silvia (2009)), grit (Duckworth and Quinn (2009)), impulsivity (Sleddens et al. (2013)), and critical thinking (Sosu (2013)) both at baseline and endline. The motivation to collect these attributes is to establish the validity of our task-based curiosity measure and explore potential channels through which the program might impact learning outcomes.

We also collected rich information from teachers. In addition to demographic information collected at baseline, we collected their fluid IQ via Raven's test and their emotional intelligence through the Reading the Mind in the Eyes test (Baron-Cohen et al. (1997)). We also collected detailed information regarding teachers' everyday teaching practices and beliefs both at baseline and endline. For the former, we adapted some of the item questions from the Teaching and Learning International Survey (TALIS) questionnaire (OECD (2013)), and constructed the following styles: Modern (learner-centered) teaching, warmth, and extrinsic motivator. For the latter, we elicited growth mindset (Dweck (2008)), attachment to the profession, competence beliefs, and gender stereotyping. We also measured teachers' curiosity using Kashdan and Silvia (2009) and critical thinking using Sosu (2013). Finally, we tested teachers' curricular knowledge in science to establish whether the intervention increased their content knowledge. We conducted this test in the second (surprise) visit along with the 40-question booklet test<sup>16</sup>. Full measurement inventory for students and teachers is presented in the Online Appendix E.

Table A1 presents the balance of student, teacher, and classroom characteristics at baseline. Balance for each study separately is presented in Table B1 and B2 in the Online Appendix. We detect no significant imbalance in any of the variables in either study and conclude that randomization was successful.

We estimate the average treatment effects of the program on outcomes of interest by conditioning on baseline covariates and strata fixed effects:

$$y_{ics} = \alpha_0 + \alpha_1 T_s + X'_{ics} \beta + W'_{cs} \gamma + \delta_d + \varepsilon_{ics}$$
 (1)

where  $y_{ics}$  is the outcome of interest for child i in classroom c, school s.  $T_s$  is the binary

<sup>&</sup>lt;sup>16</sup>Both science and booklet tests for teachers were implemented in the second study only.

treatment indicator, which equals one if school s is in the treatment group and zero otherwise, and  $X'_{ics}$  is a vector of student-level observables,  $W'_{cs}$  is a vector of classroom and teacher level observables measured at baseline. The former includes student gender, age in months, standardized fluid IQ score, risk aversion, and baseline achievement test scores. The latter includes class size, the share of refugees in the classroom, teacher experience, teacher IQ, and teacher gender.  $\delta_b$  represents district fixed effects. The estimated  $\hat{\alpha}_1$  is the average treatment effect. Standard errors are clustered at the school level. Throughout the text, we present the results from the pooled sample. The summary of the results for each province separately is given in Figure B2 in the Online Appendix. We present our full results corrected for multiple hypotheses testing (sharpened q-values and Romano Wolf p-values) in Table A2. Most of our results survive the adjustments. For the long-term results (Study 1), we use inverse probability weights to account for attrition.

All treated teachers were expected to practice the proposed pedagogy upon receiving training. Recall that participation in the program was voluntary, and the program was oversubscribed. However, we acknowledge that compliance in terms of the actual implementation may not be perfect. To assess compliance, we asked treated teachers to report their estimated degree of program implementation at endline. Specifically, we asked them to mark their estimated degree of implementation using an unmarked 10cm line. The elicited distance gives us a continuous measure of program implementation intensity ranging anywhere between zero and 100%. Note that because this is a pedagogical intervention that aims to influence the way teachers teach, the reported implementation intensity is purely subjective. Nevertheless, we believe that it gives us an idea of teacher compliance. Figure A2 depicts the distribution of the reported implementation intensity for the pooled sample. Overall, treated teachers report to have accomplished 81% program coverage. Given this high but still imperfect compliance, the estimated  $\hat{\alpha}_1$  should be interpreted as the average intent to treat effect (ITT).

## 4.1 The Predictive Validity of the Curiosity Task

Before presenting the program effects, we show that our curiosity measure (willingness to pay for a booklet) has predictive validity, i.e., correlates well with knowledge retention and test scores. To do this, we use our control sample. Figure A3 depicts the distribution of forgone tokens for the control sample. Children, on average, forgone 6.14 tokens to receive their desired booklet, with the minimum WTP being zero (7.3% of the sample) and a maximum of 10 (22.7% of the control sample).

Curiosity is known to be associated with higher cognitive ability in individuals, and our data corroborates this evidence. Table 1 presents the predictive power of overall curiosity, scientific curiosity, and non-science curiosity on science, math, and verbal test scores, as well as knowledge retention (performance on the respective questions in the booklet test). Panel 1 presents raw associations, and panel 2 presents the associations controlling for fluid IQ. The results in this table confirm that our measure of curiosity has reasonable validity in predicting crystallized cognitive ability. Correlations are particularly strong for scientific curiosity, that is, the willingness to pay for a science-related booklet. One standard deviation increase in the willingness to pay for a science-related booklet is associated with 0.083 standard deviations higher science test scores, 0.02 standard deviation higher math scores, and 0.09 standard deviation higher verbal test scores.

Table 1: Predictive Power of Curiosity Task

Panel 1: Raw association	ns			
	Science	Maths	Verbal	Retention
Overall Curiosity	0.041**	0.027	0.035*	0.047***
	(0.02)	(0.02)	(0.02)	(0.02)
Science Curiosity	0.083***	0.024*	0.086***	0.084***
	(0.02)	(0.01)	(0.01)	(0.02)
<b>3</b> 7 G . G	0 0 10 dodate	0.000	0 0 7 0 1 1 1 1 1	O O O = dududu

Panel 2: Raw associations controlling for IQ Score

	Science	Maths	Verbal	Retention
Overall Curiosity	0.021	0.014	0.011	0.040***
	(0.02)	(0.02)	(0.02)	(0.01)
Science Curiosity	0.063***	0.012	0.061***	0.076***
	(0.02)	(0.01)	(0.01)	(0.02)
Non-Science Curiosity	-0.046***	0.000	-0.052***	0.088***
	(0.02)	(0.01)	(0.01)	(0.02)
Observations	4558	4558	4675	4558

The table presents OLS coefficients of the regression of test scores (science, verbal, math and booklet test) on task-based curiosity measure (WTP). The analysis uses only the control sample. Standard errors are clustered at the school level and are reported in parentheses. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

Another important question for the purpose of this study is whether curiosity is associated with knowledge retention, as suggested by recent neuroscientific research. Panel 1 shows that higher curiosity is associated positively with the retention of information. Specifically, a

one standard deviation increase in the willingness to pay for a science-related booklet is associated with 0.08 standard deviations higher performance in science-related booklet questions. This association remains strong, controlling for IQ (0.08 standard deviations). Finally, note that non-science curiosity, the willingness to pay for either history, sports, or cartoons booklet, is negatively associated with crystallized IQ but still positively associated with higher performance in non-science booklet questions. The way we implement our curiosity task in the classroom and the fact that we have a randomly implemented program that enhances curiosity allows us to go beyond these correlations. In Section 5 we will provide evidence on the causal link between curiosity and deep learning in children.

Table 2: Associations Between Curiosity Task (WTP) and Socio-emotional Skills

		Science					
	Curiosity	Curiosity	7				Critical
	Survey	Survey	$\operatorname{Grit}$	Impulsivity	Risk	Ambiguity	Thinking
Overall Curiosity	0.037**	0.042**	0.051***	-0.014	0.228***	0.189***	0.050***
	(0.02)	(0.02)	(0.02)	(0.02)	(0.02)	(0.02)	(0.02)
Science Curiosity	0.047***	0.053***	0.049***	-0.056***	0.086***	0.075***	0.067***
	(0.02)	(0.02)	(0.02)	(0.02)	(0.02)	(0.02)	(0.02)
Non-Science Curiosity	-0.015	-0.017	-0.006	0.044***	0.108***	0.086***	-0.027*
	(0.01)	(0.02)	(0.02)	(0.01)	(0.02)	(0.02)	(0.02)
Observations	4954	4954	4524	4650	5070	5066	3635

Panel 2: Raw associations controlling for IQ Score

			<u> </u>				
		Science					
	Curiosity	Curiosity	7				Critical
	Survey	Survey	$\operatorname{Grit}$	Impulsivity	Risk	Ambiguity	Thinking
Overall Curiosity	0.028*	0.031*	0.044**	-0.007	0.232***	0.193***	0.040**
	(0.01)	(0.02)	(0.02)	(0.01)	(0.02)	(0.02)	(0.02)
Science Curiosity	0.038**	0.042***	0.042***	-0.048***	0.090***	0.079***	0.058***
	(0.01)	(0.02)	(0.02)	(0.02)	(0.01)	(0.02)	(0.02)
Non-Science Curiosity	-0.014	-0.015	-0.005	0.043***	0.107***	0.086***	-0.025
	(0.01)	(0.02)	(0.02)	(0.01)	(0.02)	(0.02)	(0.02)
Observations	4954	4954	4524	4650	5070	5066	3635

The table presents OLS coefficients of the regression of socio-emotional skills, survey measure of curiosity and risk/ambiguity preferences on task-based curiosity measure (WTP). Risk and ambiguity preferences are measured via incentivized tasks. The analysis uses only the control sample. Standard errors are clustered at the school level and are reported in parentheses. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

Table 2 further validates our curiosity task. First, we check whether our measure correlates with survey measures of curiosity developed by Kashdan and Silvia (2009). In addition,

we conjecture that curiosity may be correlated with attitudes toward uncertainty, grit, critical thinking, and impulsive behavior, acknowledging its possible relationship with some other social and emotional skills we do not measure in this paper. Panel 1 presents raw associations, and Panel 2 presents the associations controlling for fluid IQ. We observe strong positive correlations between our curiosity measure with established survey measures of curiosity. Moreover, our scientific curiosity measure (willingness to pay for science-related booklets) correlates positively with grit, critical thinking, and risk and ambiguity tolerance and negatively correlates with impulsivity. In Section 5 we will explore whether the pedagogical program made any impact on these attributes.

### 4.2 Treatment Effect on Curiosity

We first explore whether the program affects children's interests and, in particular, whether it increases their interest in science. Table 3 Panel 1 presents the estimated effects of the program on the preferred booklet type. The first column shows the treatment effect on the probability of choosing to purchase a science-related booklet (science, animals, space, vehicles, human anatomy). The second column presents the treatment effect on choosing a non-science booklet (history, sports, and cartoons). The last column gives the estimated effect of the treatment on "no interest," i.e., the probability of choosing not to purchase a booklet. Notice that about 50% of the children in the control group stated their willingness to purchase a science-related booklet. This value goes up to 54\% in the treatment group, and this difference is statistically significant at the 1% level. It appears that the program shifted children's interest to science topics but not much at the expense of non-science topics (see column 2). As shown in column 3 of the table, the program lowered the probability of no interest, i.e., stating zero willingness to pay, by 2.9 percentage points, representing about a remarkable 50% effect. The program effect on interest in science can also be seen in Figure 2. Treated children are significantly more likely to rank science, animals, and space booklets as their top 3 interests.

Table 3 Panel 2 presents the estimated treatment effects on the willingness to pay for the desired booklet. Note that the measure is standardized to have a mean zero for the control group, so the coefficient estimates are standard deviation effects. Column 1 presents the overall willingness to pay for any preferred booklet, column 2 for a science-related booklet, and the last column for a non-science booklet.

Table 3: Treatment Effect on the Choice of Booklet and Level of Curiosity

Panel 1: Choice of Booklet

	Science Related	Non-Science Related	No booklet
Treatment	0.039***	-0.010	-0.029***
	(0.01)	(0.01)	(0.01)
Control Mean	0.50	0.44	0.06
Observations	10870	10870	10870
Number of Schools	134	134	134

Panel 2: Level of Curiosity

	Curiosity	Science Curiosity	Non-Science Curiosity
Treatment	0.110***	0.100***	-0.011
	(0.04)	(0.03)	(0.02)
Control Mean	-0.01	-0.00	-0.00
Observations	10864	10863	10863
Number of Schools	134	134	134

Estimates are obtained via OLS. The dependent variables are binary indicators of choosing a science-related booklet (science, space, vehicles, human body, and animals) in column 1, choosing a nonscience-related booklet (history, sports, and cartoons) in column 2, and choosing no booklet option in column 3. Standard errors are clustered at the school level and are reported in parentheses. Covariates include gender, age, fluid IQ, risk tolerance, survey measure of curiosity, math and verbal scores as individual baseline characteristics, class size, the share of refugees, teacher gender, experience, and fluid IQ as baseline classroom and teacher characteristics. Grade and district fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

Subject Ranked in Top 3

Science Animals Space H. Body Vehicles History Sports Cartoons

Figure 2: Treatment effect on the Ranking of Booklets

The figure depicts the average marginal treatment effects obtained from logistic regressions on subject interest. The dependent variables are binary indicators of one if the respective booklet is ranked as one of the top 3 interests by the student. Confidence intervals are obtained by clustering at the school level.

We estimate a significant 0.11 standard deviation effect on overall curiosity. In terms of tokens, this corresponds to forgoing about 0.35 extra tokens for a booklet. Given that children forgo 6.1 of their tokens on average for their preferred booklets in the control group, this effect implies a 6% treatment effect. The effect on science curiosity is similar with about 0.10 standard deviation treatment effect, again precisely estimated. The estimated effect on the willingness to pay for non-science booklets is statistically zero. These results show that the program is successful in stimulating children's curiosity and interest in science. Our next question is whether this stimulated curiosity translates into actual learning. The temporal component of our task and the half-half implementation of booklet distribution allow us to answer this question.

#### 4.3 Treatment Effect on Knowledge Retention

The estimated treatment effects on the willingness to pay suggest that in the market price regime, where the price of a booklet is determined randomly, treated classrooms necessarily end up with a proportionally higher number of booklets. This means that treated classrooms have more information (booklets) available for all, making it more likely to acquire and retain the knowledge available in the classroom. A clean identification of the effect of the program on knowledge retention requires the amount and the content of information to be balanced across treatment status. The half-half implementation regime delivers this by design. Recall that in classrooms subject to this regime, we distributed the booklets randomly to half of the students regardless of their willingness to pay and their choice of booklets. Panel 1 in Table 4 presents the estimated treatment effects on booklet test scores using the full sample for comparison purposes. The first 3 columns give short-term, and the last 3 give long-term effects (only Study 1).

Panel 2 presents the results using only the classes that were subject to the half-half regime. The effect of the program on the ability to retain knowledge is striking. Treated students performed significantly better than untreated students in the 40-question booklet test. Considering the half-half regime, where we have clean identification, treated students performed about 0.11 standard deviations higher than untreated students overall, and the performance difference is similar in science topics (0.10 sd). Note that treated students performed better even in non-science topics of the test, although the estimate is less precise. What is truly remarkable is that after 3 years and a devastating pandemic, treated students still exhibit much higher booklet knowledge than untreated students, supporting the claims of neuroscientists that enhanced curiosity is associated with memory consolidation, i.e., deep

learning. Treated students performed 0.14 standard deviations higher in the booklet test 3 years after the intervention. The retention of science-related topics after 3 years is about 0.16 standard deviations.

Table 4: Treatment Effect on Knowledge Retention

Panel 1: Knowledge Retention - Full Sample

	Short Term				Long Term		
			Non-Science			Non-Science	
	Retention	Retention	Retention	Retention	Retention	Retention	
Treatment	0.118***	0.106***	0.088**	0.082*	0.090**	0.039	
	(0.04)	(0.04)	(0.04)	(0.04)	(0.04)	(0.05)	
Control Mean	-0.01	-0.01	-0.01	-0.00	0.00	-0.00	
Observations	10590	10590	10590	2426	2426	2426	
Number of Schools	134	134	134	50	50	50	

Panel 2: Knowledge Retention - Half Half

		Short Term			Long Term		
	Retention		Non-Science Retention			Non-Science Retention	
Treatment	0.114**	0.103**	0.084*	0.141**	0.161***	0.058	
	(0.05)	(0.05)	(0.04)	(0.06)	(0.05)	(0.07)	
Control Mean	-0.04	-0.03	-0.03	-0.04	-0.06	0.00	
Observations	9037	9037	9037	1336	1336	1336	
Number of Schools	134	134	134	50	50	50	

Panel 3: Knowledge Retention (excluding Preferred Booklet) - Half Half

	0		`	0		,		
		Short Term				Long Term		
		Retention		Non-Science Retention			Non-Science Retention	
Treatment		0.118**	0.104**	0.095**	0.173***	0.160***	0.120**	
		(0.05)	(0.05)	(0.04)	(0.06)	(0.05)	(0.07)	
Control Mean		-0.02	-0.02	-0.02	-0.04	-0.06	0.00	
Observations		8271	8271	8271	1219	1219	1219	
Number of Schools		134	134	134	50	50	50	

Estimates are obtained via OLS. The dependent variables are standardized booklet test scores (knowledge retention). The first 3 columns give short-term results using the pooled sample, and the last 3 provide the long-term results of Study 1. Standard errors are clustered at the school level and are reported in parentheses. Covariates include gender, age, fluid IQ, risk tolerance, survey measure of curiosity, math and verbal scores as individual baseline characteristics, class size, the share of refugees, teacher gender, experience, and fluid IQ as baseline classroom and teacher characteristics. Grade and district fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

Neuroscientific evidence indicates that sparked curiosity leads to the absorption of the

available information in one's environment, whether such information is of interest to the individual or not. Our research design allows us to test this hypothesis. We do know the student's preferred booklet, whether or not she received a booklet, and if received, which booklet she received. We constructed a booklet test score performance for each child by eliminating the questions related to her preferred topic. Panel 3 presents estimated effects on knowledge of topics outside students' preferred booklet. The overall retention results refer to the performance on seven topics (topics other than the preferred one). Science retention refers to the performance of students who preferred a non-science booklet on science-related topics. Non-science retention refers to the performance of students who preferred a science booklet on non-science-related topics. Positive and significant short-term and long-term effects are consistent with the neuroscientific evidence that sparked curiosity leads to stronger absorption of knowledge available in one's environment. These estimates also clue us in on a possible social aspect of curiosity and learning, which we explore deeper in the next section.

#### 4.4 Treatment Effect on Information Dissemination in the Classroom

It has been shown that in addition to being associated with deep learning, human curiosity has positive externalities. Hartung and Renner (2013) and Litman and Pezzo (2007) show that curiosity is associated with passionate information sharing. Dubey, Mehta and Lombrozo (2021) show that human curiosity is sensitive to the social environment and stimulated by the curiosity of others. These externalities imply enhanced peer learning in our context, and our research design allows us to test the presence of these externalities. Our test involves exploring whether the program made the classroom a denser learning environment where students share what they learn with their peers. We collected friendship networks at baseline and endline by asking each student to nominate at most three peers in their classrooms as their friends. With these nominations and the fact that we know who received which booklet, we can gain a deeper understanding of how the information provided to a subset of students in the classroom is disseminated and how treatment interacts with the way information is disseminated.

Table 5 shows the treatment effect on retention for students who received a booklet (Panel 1) and those who did not (Panel 2). The former takes the students who received booklets under the half-half regime, and the latter uses all students who did not receive any booklet. The effect sizes are larger and more precisely estimated for booklet recipients. Nevertheless, Panel 2 provides evidence of better information sharing in treated classrooms.

Table 5: Treatment Effect on Knowledge Retention through Information Dissemination

Panel 1: Booklet Received

	Retention	Science Retention	Non-Science Retention
Treatment	0.150***	0.130**	0.118**
	(0.06)	(0.05)	(0.05)
Control Mean	0.00	-0.00	-0.00
Observations	4202	4202	4202
Number of Schools	134	134	134

Panel 2: No Booklet Received

	Retention	Science Retention	Non-Science Retention
Treatment	0.080	0.072*	0.059
	(0.05)	(0.04)	(0.04)
Control Mean	0.00	-0.00	-0.00
Observations	5265	5265	5265
Number of Schools	134	134	134

Panel 3: Network Effect

	Retention	Science Retention	Non-Science Retention
Treatment	0.178**	0.157**	0.137*
	(0.08)	(0.07)	(0.07)
Control Mean	-0.00	-0.00	-0.00
Observations	1075	1075	1075
Number of Schools	134	134	134

Estimates are obtained via OLS. The dependent variables are standardized booklet test scores (knowledge retention). Panel 1 uses the sample of booklet recipients only in the Half-Half regime. Panel 2 uses the sample of students who did not receive any booklet. Panel 3 uses the sample of students who did not receive any booklet but have at least one person in their network who has received the booklet of their choice. Standard errors are clustered at the school level and are reported in parentheses. Covariates include gender, age, fluid IQ, risk tolerance, survey measure of curiosity, math and verbal scores as individual baseline characteristics, class size, the share of refugees, teacher gender, experience, and fluid IQ as baseline classroom and teacher characteristics. Panel 3 specification includes the total number of friendship ties the student has in the classroom. Grade and district-fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

Treated students who did not receive a booklet performed 0.07 standard deviations better in the booklet test (science-related questions) than their untreated counterparts. However, the retention effects are generally weaker for those who did not receive a booklet. In Panel 3, we present the retention results for students who did not receive a booklet and whose preferred booklet was received by someone else in their friendship network. Here, the friendship network of a student contains all her friendship nominations (out-degree ties) and all her

classmates who nominate her as their friend (in-degree ties).<sup>17</sup> The results are remarkable: We estimate about 0.18 standard deviation higher booklet knowledge for these students overall, suggesting a significantly higher pursuit of information among treated students. Treatment effects on science and non-science knowledge retention are 0.16 and 0.14 standard deviations for these students, respectively. Note that while highly restricted, this sample is balanced across treatment status with respect to baseline characteristics; see Table A3.

We also explore the treatment effect heterogeneity under differential information availability within friendship networks to complement these results. Figure 3 plots the estimated treatment effects on knowledge retention conditional on information availability within friendship networks. Here, we focus on the information the student is interested in, i.e., booklets that he/she ranked as top 3. Panel 1 presents estimated effects conditional on receiving a booklet, and an increasing number of top-3 ranked booklets available in the friendship network (zero, one, two, or three and more booklets). Panel 2 presents estimated effects conditional on receiving no booklet, and an increasing number of top-3 ranked booklets available in the friendship network. The depicted treatment effects suggest significantly higher knowledge retention for treated children, monotonically increasing with the availability of information within their networks. Consistent with Table 5 Panel 1, the estimated effects are stronger for booklet owners. Treated booklet owners who are the sole booklet owners within their network perform 0.08 standard deviations better in science-related booklet questions than untreated booklet owners in the same situation. The estimated treatment effect goes up to 0.25 standard deviations for this group when their friendship network possesses more than three science-related booklets.

The effects are weaker for those who did not receive a booklet (Panel 2), but we still observe monotonically increasing treatment effects on knowledge retention in science-related topics as information availability increases within the network. We interpret these estimates as significantly more efficient information dissemination and peer learning in treated class-rooms where students are more curious and passionate about pursuing and sharing knowledge. Note also that stronger effects estimated for the booklet owners suggest that access to available information within networks via booklet exchange is more prominent in treated classrooms. Put differently, booklet owners, who are in a better position to access other booklets in their network, do access and absorb more information in curious classrooms.

<sup>&</sup>lt;sup>17</sup>We checked whether the program had any impact on the network structure, such as the network density, the number of friendship ties, the number of isolated students, and the number of reciprocal ties, and found no such evidence.

Panel 1: Booklet Received Overall Science Non-Science Θ Standardized Scores 0.25 0.14" Ņ 0 3 2 3 0 2 2 Books Science Related Books Non-Science Books Panel 2: No Booklet Received Science Non-Science Overall Standardized Scores 0.24 0.08 0.10 0.06 0.04 0 2 3 0 2 3 Ö 2 Books Science Related Books Non-Science Books

Figure 3: Information Availability and Treatment Effects on Knowledge Retention

The figure depicts the estimated treatment effects on standardized booklet test scores (knowledge retention). Panel 1 restricts the sample to those who received a booklet, and Panel 2 restricts the sample to those who did not receive a booklet. Depicted coefficient estimates are obtained by further restricting each sample as having none, one, two, and more than three top-3 ranked booklets in the student's network (our measure of information availability within the friendship network). Standard errors are clustered at the school level. Asterisks indicate statistical significance at the 1% \*\*\*\*, 5% \*\*\*, and 10% \* levels.

#### 4.5 Treatment Effect on Test Scores and Educational Aspirations

Our next question is whether these positive learning effects extend to actual learning outcomes. Table 6 presents the treatment effects on math, verbal (Turkish), and science test performance. While we do not estimate statistically significant effects on math and Turkish, we find that treated students perform significantly better than untreated students in the science test. The effect is about 0.08 standard deviations and significant at the 1% level. The positive effect on science test scores also persists into middle school years. We find that treated students still perform better than untreated students in science (0.07 standard deviations).

ations) 3 years after the implementation of the program. The near-zero effect sizes for math and verbal scores, while a significant and persistent effect on science, may be attributable to the program's heavy emphasis on science.

**Table 6:** Treatment Effect on Subject Test Scores

	S	Short Term			Long Term		
	Science	Maths	Verbal	Science	Maths	Verbal	
Treatment	0.078***	0.017	0.033	0.073*	-0.017	-0.022	
	(0.03)	(0.03)	(0.03)	(0.04)	(0.04)	(0.04)	
Control Mean	-0.01	0.00	-0.00	0.00	0.00	0.00	
Observations	9949	10400	10680	2426	2426	2426	
Number of Schools	134	134	134	50	50	50	

Estimates are obtained via OLS. The dependent variables are standardized subject test scores. The first 3 columns give short-term results using the pooled sample, and the last 3 provide the long-term results of Study 1. Standard errors are clustered at the school level and are reported in parentheses. Covariates include gender, age, fluid IQ, risk tolerance, survey measure of curiosity, math and verbal scores as individual baseline characteristics, class size, the share of refugees, teacher gender, experience, and fluid IQ as baseline classroom and teacher characteristics. Grade and district fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

To measure educational aspirations, we asked children two questions. First, we asked whether they intended to go to university when they grew up. Second, if they did, we what study major they wanted to pursue. For the latter, we gave them a full list of study majors to choose from. The first column of Table 7 presents the estimated treatment effect (average marginal effect) on the willingness to go to university. The following columns present the estimated average marginal effects on planned study majors. These are science, engineering, medicine, and Non-STEM (social sciences and humanities). Note first that almost all (95%) children in the control group stated that they plan to go to university when they grow up. Nevertheless, we still estimate a significant treatment effect on this high base, albeit small in size (1 percentage point). More importantly, only 12% of the children in the control group state their plan to major in science at university. This value is 2.3 percentage points higher for the treatment group, implying a 19% treatment effect. We estimate null effects for engineering and medicine. The estimated negative effect on non-STEM majors suggests that the positive effect we estimate for science comes at the expense of non-STEM majors. While 61% of the students express a preference toward a social science topic in the control group, treated students are 0.02 percentage points less likely to state such a preference. Note, however, while estimated sizes remain similar in the long run, they are estimated imprecisely, likely due to insufficient statistical power.

**Table 7:** Treatment Effect on Aspirations

Panel 1: Short Term

	University	Science	Engineering	Medical	Non-STEM
Treatment	0.008*	0.023***	0.001	-0.003	-0.022*
	(0.00)	(0.01)	(0.01)	(0.01)	(0.01)
Control Mean	0.95	0.12	0.12	0.16	0.61
Observations	10693	10186	10186	10186	10186
Number of Schools	134	134	134	134	134

Panel 2: Long Term

	University	Science	Engineering	Medical	Non-STEM
Treatment	0.010	0.019	0.011	-0.016	-0.014
	(0.01)	(0.02)	(0.02)	(0.02)	(0.02)
Control Mean	0.95	0.13	0.12	0.22	0.54
Observations	2320	2182	2182	2182	2182
Number of Schools	50	50	50	50	50

Estimates are obtained via OLS. The dependent variables are binary choice variables of intention to go to university, intention to choose a science major, engineering major, medicine, and non-STEM major. Panel 1 presents short-term results from the pooled sample, and Panel 2 long-term results from Study 1. Standard errors are clustered at the school level and are reported in parentheses. Covariates include gender, age, fluid IQ, risk tolerance, survey measure of curiosity, math and verbal scores as individual baseline characteristics, class size, the share of refugees, teacher gender, experience, and fluid IQ as baseline classroom and teacher characteristics. Grade and district fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

#### 4.6 Heterogeneity in Treatment Effects

As stated in our PAP, we explore heterogeneity in treatment effects with respect to two characteristics. First, we check whether the estimated effects are different across gender. Second, we investigate whether the program has a differential impact on children with different levels of cognitive ability (fluid IQ). The first panel in Table A4 shows that the program's effect on the shift toward science topics mainly comes from girls. Treated girls are 7.9 percentage points more likely to choose a science-related booklet relative to untreated girls. The corresponding estimate is statistically zero for boys. As for choosing no booklet (no interest), we estimate no gender heterogeneity. Both boys and girls in the treatment group are significantly less likely to choose "no booklet" than those in the control group, suggesting that the program stimulated the overall interest of both boys and girls.

Similarly, we detect a significant gender heterogeneity in the treatment effect on curiosity. The estimates in Panel 2 indicate that while the program is effective in increasing curiosity for both genders, the results seem stronger for girls. Treated girls have 0.17 standard deviations

higher scientific curiosity than untreated girls. We reject the equality of effects for overall curiosity as well as science and non-science curiosity. However, we estimate no significant gender heterogeneity in retention and test scores (Table A5). Finally, we do not detect any noteworthy gender heterogeneity in aspirations; see Table A6.

Table A7 presents treatment effect heterogeneity with respect to cognitive ability. Here, we use our measure of fluid IQ (Raven score) and estimate treatment effects separately for high (above median) and low IQ (below median) levels. Overall, the estimated effects seem stronger for students with higher cognitive ability, although we fail to reject the equality of the estimated effects in most cases. The exception is the treatment effect on curiosity. As can be seen in Panel 2, while the program seems effective in increasing curiosity for all cognitive levels, its effect is stronger for students with high cognitive ability. This is reflected in the retention results (Table A8, Panel 1), but we fail to reject the equality of the estimates (see the borderline p-values). Finally, We do not estimate any treatment effect heterogeneity in test scores or aspirations with respect to IQ; see Table A9.

Taken together, our results suggest that the program was highly successful in increasing children's interest in science and stimulating their curiosity. In addition, it was highly effective in enhancing children's ability to retain the acquired knowledge and improving science test scores. In the next section, we will explore possible mechanisms through which the program achieves these positive results.

## 5 Potential Mechanisms

While the program had a specific focus on curiosity, given the correlations we established in Table 2, it is plausible that it also affected related attributes in children, potentially leading to improved learning. To investigate this, we explore whether the program had any impact on other attributes that are correlated with curiosity. We also explore the role of teachers, in particular, the program's impact on their styles, practices, beliefs, and behaviors as additional channels.

Figure 4 depicts the estimated treatment effects on grit, impulsivity, risk and ambiguity tolerance, critical thinking, and survey measure of epistemic and scientific curiosity. For the long term (Study 1), we only have self-reported epistemic and scientific curiosity and grit. Note first that consistent with the effects we estimate on the behavioral task, we estimate a 0.20 sd treatment effect on self-reported curiosity and a 0.15 sd effect on self-reported

scientific curiosity. We find the former effect persists into adolescence, but the latter does not (Study 1). We also find that treated children have become more tolerant of risk and ambiguity and become more critical in their thinking process relative to untreated children. We do not estimate a statistically significant treatment on grit either in the short or the long term.

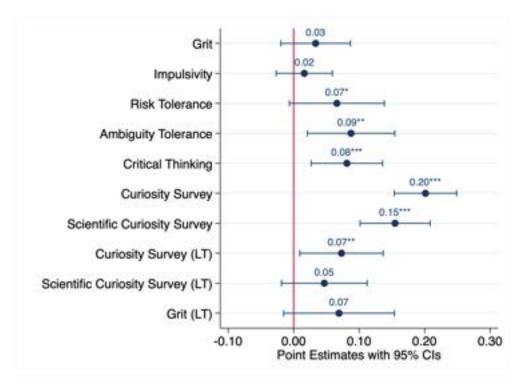


Figure 4: Treatment Effects on Students' Beliefs and Attitudes

The figure depicts the estimated treatment effects on children's socio-emotional skills, beliefs, and attitudes. Standard errors are clustered at the school level. Covariates include individual baseline characteristics: gender, age, baseline fluid IQ, risk tolerance, survey measure of curiosity, math and verbal scores, baseline classroom and teacher characteristics: class size, the share of refugees, teacher gender, experience, and fluid IQ. Grade and district fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

It is also plausible that part of the program's success may stem from its success in influencing teaching practices and teachers' beliefs. Figure 5 plots the estimated treatment effects on teaching styles and beliefs. What emerges from the figure is that the program made a positive impact on the teaching styles, teachers' epistemic curiosity, and mindset. Treated teachers report 0.23 (0.27) standard deviations and higher curiosity (growth mindset) than untreated teachers. In terms of teaching styles, the effects on practicing modern and learner-centered teaching emerge as the most prominent (0.17 sd effect). It is clear that teachers who embrace the prescribed pedagogy become more curious themselves and adopt a more modern teaching approach and a growth mindset. These findings are consistent with the positive

feedback we received from teachers regarding the program throughout the implementation period.

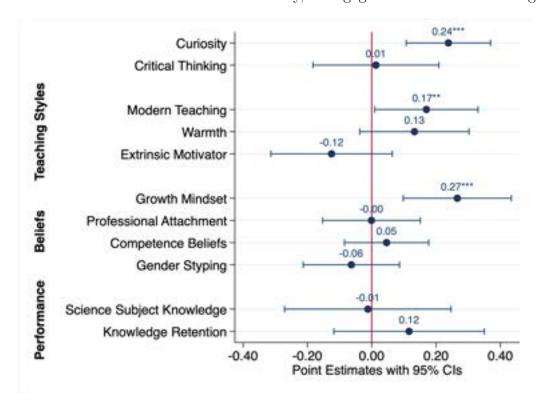


Figure 5: Treatment Effects on Teacher Ability, Pedagogical Beliefs and Teaching Styles

The figure depicts the estimated treatment effects on teachers' skills, beliefs, attitudes, and teaching styles. Standard errors are clustered at the school level. Covariates include teacher gender, experience, fluid IQ, class size, and the share of refugees. Grade and district fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

As seen in the figure, we safely rule out the mechanism whereby teachers learn more curricular material themselves and improve students' science scores. We estimate precise null effects on teachers' curricular knowledge in science. We also find no evidence of higher booklet knowledge in treated teachers, ruling out the mechanism whereby teachers learn the information provided in booklets and teach their students. Note that the latter was unlikely by design as we gave no indication that we would return and give a test containing questions about the information provided in booklets. Nevertheless, we rule out a retention mechanism where treated teachers teach the booklet content to their students for any reason.<sup>18</sup>

The above results suggest that there may be multiple channels through which the treatment improved learning. For knowledge retention results, however, the theoretical basis

<sup>&</sup>lt;sup>18</sup>Estimates for science test scores and booklet test scores for teachers are available only for Study 2 as we were not given the opportunity to test teachers in the first study.

articulated by the neuroscientific research can help us zoom into a narrower set of channels. For this, we postulate a structural model of memory consolidation. The model is informed by the neural mechanisms of curiosity documented in Gruber, Gelman and Ranganath (2014) that the higher the urge to know, the stronger the memory consolidation process, i.e., the deeper the learning. Given that our experiment is conducted in a social setting with strong positive learning externalities documented in Section 4.4, we consider the following model:

$$y_{ics} = \beta_0 + \beta_1 Curiosity_{ics} + \beta_2 \bar{y}_{cs-i} + X'_{ics} \beta_3 + W'_{cs} \beta_4 + \varepsilon_{ics}$$
 (2)

where  $y_{ics}$  is the level of booklet-related knowledge of student i in classroom c, school s.  $Curiosity_{ics}$  is the student's urge to know the information provided in her preferred booklet, measured as her willingness to pay for the booklet. Note that we consider overall booklet knowledge of student i as an outcome, although the curiosity measure relates to the student's preferred booklet. This specification is consistent with the neuroscientific evidence that, in addition to the topic of preference, the sparked curiosity leads to deep learning of unrelated topics, as we also show in Table 4, Panel 3. Capturing peer learning and dissemination effect,  $\bar{y}_{cs-i}$  is the leave-out mean of peers' booklet knowledge. This captures the extent to which one's knowledge is enhanced by the knowledge of her peers, as evidenced in Table 5. The knowledge flow from peers to the student can be through direct communication or simple booklet exchange.  $X'_{ics}$  is a vector of student-level observables, which include demographics and all non-cognitive skills measured pre-treatment.  $W'_{cs}$  is a vector of classroom and teacherlevel observables measured at baseline. Naturally, this structural equation is hard to estimate as it implies two channels for knowledge retention for an individual student, and we have only one credible instrument (random treatment assignment). 19 Nevertheless, we can show, albeit suggestively, that both factors are significant mediators and have separate impacts on knowledge retention.

Table 8 presents results from a mediation analysis using Equation 2 as the structural model of knowledge retention. The first raw in each panel presents treatment effects on retention, curiosity (willingness to pay), and the leave-out mean of classroom booklet knowledge (capturing peer learning). The second raw presents the estimated coefficients of  $\beta_1$  and

<sup>&</sup>lt;sup>19</sup>In fact, it may be more realistic to interact curiosity with peer knowledge and postulate a three-factor model as one's curiosity may affect how much peer knowledge is absorbed. We experimented with such an empirical model and found that the results were similar to what we obtained using the above two-factor model. As an alternative, one can postulate a one-factor model and use an IV-based mediation analysis. However, given our results, the exclusion restriction assumption would be too strong.

 $\beta_2$  obtained by estimating Equation 2 via OLS. The third row shows the mediated effect of each channel, and the final row presents the percentage of the effect mediated with the two factors combined. Consider, for example, science-related knowledge retention (Panel 2). As shown in Table 4 Panel 2, the total effect on retention is estimated as 0.103 sd. The estimated treatment effect on scientific curiosity is 0.100 sd (also shown in Table 3, Panel 2).

**Table 8:** Mediation Analysis

Panel 1: Overall Curiosity

	Retention	Curiosity	Peer Learning
Treatment	0.114**	0.110***	0.119**
	(0.04)	(0.04)	(0.06)
Coefficient Estimates (Eq2)		0.047***	0.787***
		(0.01)	(0.04)
Mediated		0.005	0.093
		[0.001, 0.013]	[0.005, 0.198]
% Total Mediated		86.5%	

Panel 2: Science Curiosity

	Retention	Curiosity	Peer Learning
Treatment	0.103**	0.100***	0.108**
	(0.04)	(0.03)	(0.05)
Coefficient Estimates (Eq2)		0.069***	0.719***
		(0.01)	(0.04)
Mediated		0.007	0.077
		[0.002, 0.013]	[0.004, 0.168]
% Total Mediated		82.	2%

Panel 3: Non Science Curiosity

	Retention	Curiosity	Peer Learning	
Treatment	0.084*	-0.011	0.087*	
	(0.04)	(0.02)	(0.05)	
Coefficient Estimates (Eq2)		0.079***	0.749***	
		(0.01)	(0.05)	
Mediated		-0.001	0.065	
		[-0.003, 0.004]	[-0.003, 0.151]	
% Total Mediated		76.1%		

The table presents the mediation results using Equation 2 as the structural function of knowledge retention. The first raw presents the estimated treatment effects on retention, curiosity and peer booklet knowledge. The second raw presents the coefficient estimates of  $\beta_1$  and  $\beta_2$  in Equation 2. The third raw presents the contribution of each factor to the total booklet knowledge with confidence intervals, The final raw presents the total mediated effect of two factors combined. Asterisks indicate statistical significance at the 1% \*\*\*\*, 5% \*\*, and 10% \* levels.

The estimated effect on the leave-out mean of classroom booklet knowledge is 0.108 sd, significant at the 5% level. Estimating Equation 2 via OLS yields the point estimates of  $\beta_1$  and  $\beta_2$  as 0.069 sd and 0.719 sd, respectively. These estimates imply that of the 0.103 sd effect on knowledge retention, 0.007 sd comes from an increase in own curiosity, and 0.077 sd comes from an increase in peer knowledge, suggesting a total mediated effect of 82.2%. The total mediated effect for overall curiosity is 86.5%. This analysis suggests that while both own curiosity and peer learning are important mediators, the latter is the primary driver of the knowledge retention results. The mediated effects are imprecisely estimated for the non-science curiosity, as noted by the confidence intervals.

Like all mediation analyses, the above exercise is suggestive at best. The consistency of the estimates obtained from estimation Equation 2 via OLS requires that the memory consolidation model postulated above is correctly specified. Nevertheless, combining our causal results with the above mediation analysis, we subscribe to the conclusion that stimulated curiosity in the classroom environment leads to deep learning in children through enhanced curiosity and peer learning. Our results on booklet test scores provide the first causal evidence concerning the link between stimulated curiosity and deep learning outside the fMRI lab in a natural social setting. The setting also allows us to show, for the first time, the emergence of positive learning externalities when pupils' curiosity is stimulated.

For improved science test scores and aspirations, other mechanisms may be equally important, so we refrain from subscribing to any particular channel. In addition to enhanced curiosity, the program's positive impact on children's critical thinking skills and teachers' styles and beliefs may be partially responsible for improved test scores and higher aspirations to major in science.

#### 6 Conclusion

We test the effectiveness of a pedagogical program that aims to cultivate children's curiosity in the classroom. The pedagogy is informed by recent research on the neural mechanisms of human curiosity and mainly targets science teaching in elementary schools. The program offers teachers practices that help them create teachable moments by tapping into children's natural love of mystery, surprise, and humor. The program was implemented as two independent clustered randomized controlled trials in two large provinces of Turkey, involving 134 primary schools, 425 teachers, and about 11,000 children of age 9 to 11.

To evaluate the program's effect on children's curiosity, we develop a behavioral measure that quantifies children's urge to acquire knowledge and their ability to retain knowledge upon satisfying the urge for an extended period. We find that the intervention increases children's curiosity, measured by their willingness to pay for information and their ability to retain knowledge. Furthermore, our design allows us to show that classroom practices that nurture children's curiosity also make peer learning more efficient, supporting further deep learning. We also show that the pedagogy significantly improves children's objective test scores in science and raises their educational aspirations for science. Moreover, the effects we estimate on knowledge retention and test scores persist well into adolescence.

The results are promising and likely to have high external validity. While the participation was voluntary, the program was oversubscribed. In all participating schools, most teachers were eager to join the program. Considering the policy issue of motivating teachers, the program's positive effects on teachers are particularly encouraging. The program was also highly cost-effective. The toolkit for teachers and other written materials are now available free of charge. The remaining program costs include printing hard copy materials, distributing the materials to schools, and conducting teacher training. Total printing costs were about 30,000 USD, the distribution costs were 9,000 USD, and teacher training costs were about 6,000 USD. These values imply a minimal (4 USD) program cost per child.

Global learning poverty is at its worse in the wake of the devastating Covid-19 pandemic. While the learning crisis predates the pandemic, the pandemic-related school closures made matters disproportionately worse for underprivileged children. They further widened the already sizeable socioeconomic achievement gaps to an alarming level in both developed and developing countries. The crisis now calls for evidence-informed and scalable actions more urgently than ever. One action may be to equip teachers with effective teaching practices that have a high chance of increasing teacher and pupil engagement, resulting in quality learning. We provide rigorous evidence on the effectiveness of one such scalable and cost-effective action. We envision a couple of ways this program can be scaled up. One way is through incorporating the training in regular professional development seminars given to teachers at the beginning of the academic year. Another way can be to offer seminar courses for teacher candidates in universities. It is unclear which delivery medium would be more effective and may be a topic of future research.

### References

- Alan, Sule, and Seda Ertac. 2018. "Fostering Patience in the Classroom: Results from Randomized Educational Intervention." *Journal of Political Economy*, 126(5): 1865–1911.
- Alan, Sule, Ceren Baysan, Mert Gumren, and Elif Kubilay. 2021. "Building Social Cohesion in Ethnically Mixed Schools: An Intervention on Perspective Taking\*." The Quarterly Journal of Economics, 136(4): 2147–2194.
- Alan, Sule, Teodora Boneva, and Seda Ertac. 2019. "Ever Failed, Try Again, Succeed Better: Results from a Randomized Educational Intervention on Grit." *The Quarterly Journal of Economics*, 134(3): 1121–1162.
- **Anderson, Michael 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.
- **ASER.** 2018. "Annual Status of Education Report 'Beyond Basics' (Rural) 2017." ASER Centre.
- Ashraf, Nava, Abhijit Banerjee, and Vesall Nourani. 2021. "Learning to Teach by Learning to Learn." 115.
- Banerjee, Abhijit, Rukmini Banerji, James Berry, Esther Duflo, Harini Kannan, Shobhini Mukherji, Marc Shotland, and Michael Walton. 2016. "Mainstreaming an Effective Intervention: Evidence from Randomized Evaluations of "Teaching at the Right Level" in India." National Bureau of Economic Research Working Paper 22746.
- Banerjee, Abhijit V., Shawn Cole, Esther Duflo, and Leigh Linden. 2007. "Remedying Education: Evidence from Two Randomized Experiments in India\*." *The Quarterly Journal of Economics*, 122(3): 1235–1264.
- Banerji, Rukmini, and Madhav Chavan. 2016. "Improving Literacy and Math Instruction at Scale in India's Primary Schools: The Case of Pratham's Read India Program." Journal of Educational Change, 17(4): 453–475.
- Baron-Cohen, Simon, Therese Jolliffe, Catherine Mortimore, and Mary Robertson. 1997. "Another Advanced Test of Theory of Mind: Evidence from Very High Functioning Adults with Autism or Asperger Syndrome." Journal of Child Psychology and Psychiatry, 38(7): 813–822.

- Berlyne, D. E. 1954. "A Theory of Human Curiosity." British Journal of Psychology. General Section, 45(3): 180–191.
- Blanchard, Margaret R., Sherry A. Southerland, and Ellen M. Granger. 2009. "No Silver Bullet for Inquiry: Making Sense of Teacher Change Following an Inquiry-Based Research Experience for Teachers." *Science Education*, 93(2): 322–360.
- Buser, Thomas, Noemi Peter, and Stefan C. Wolter. 2017. "Gender, Competitiveness, and Study Choices in High School: Evidence from Switzerland." *American Economic Review*, 107(5): 125–130.
- Carlana, Michela, and Margherita Fort. 2022. "Hacking Gender Stereotypes: Girls' Participation in Coding Clubs." *AEA Papers and Proceedings*, 112: 583–587.
- Chamorro-Premuzic, Tomas, and Adrian Furnham. 2006. "Intellectual Competence and the Intelligent Personality: A Third Way in Differential Psychology:." Review of General Psychology.
- Collins, Robert P, Jordan A Litman, and Charles D Spielberger. 2004. "The Measurement of Perceptual Curiosity." *Personality and Individual Differences*, 36(5): 1127–1141.
- de Ree, Joppe, Karthik Muralidharan, Menno Pradhan, and Halsey Rogers. 2018. "Double for Nothing? Experimental Evidence on an Unconditional Teacher Salary Increase in Indonesia." *The Quarterly Journal of Economics*, 133(2): 993–1039.
- **Dubey, Rachit, Hermish Mehta, and Tania Lombrozo.** 2021. "Curiosity Is Contagious: A Social Influence Intervention to Induce Curiosity." *Cognitive Science*, 45(2): e12937.
- **Duckworth, Angela Lee, and Patrick D. Quinn.** 2009. "Development and Validation of the Short Grit Scale (Grit-S)." *Journal of Personality Assessment*, 91(2): 166–174.
- **Dweck, Carol S.** 2008. *Mindset: The New Psychology of Success*. Random House Digital, Inc.
- **Fischer, Stefanie.** 2017. "The Downside of Good Peers: How Classroom Composition Differentially Affects Men's and Women's STEM Persistence." *Labour Economics*, 46: 211–226.

- **Glewwe, P., and K. Muralidharan.** 2016. "Improving Education Outcomes in Developing Countries: Evidence, Knowledge Gaps, and Policy Implications." In *Handbook of the Economics of Education*. Vol. 5, , ed. Eric A. Hanushek, Stephen Machin and Ludger Woessmann, 653–743. Elsevier.
- Glewwe, Paul, Michael Kremer, Sylvie Moulin, and Eric Zitzewitz. 2004. "Retrospective vs. Prospective Analyses of School Inputs: The Case of Flip Charts in Kenya." Journal of Development Economics, 74(1): 251–268.
- Glewwe, Paul, Sylvie Lambert, and Qihui Chen. 2020. "Chapter 15 Education Production Functions: Updated Evidence from Developing Countries." In *The Economics of Education (Second Edition)*., ed. Steve Bradley and Colin Green, 183–215. Academic Press.
- **Gneezy**, **Uri**, and **Jan Potters**. 1997. "An Experiment on Risk Taking and Evaluation Periods\*." *The Quarterly Journal of Economics*, 112(2): 631–645.
- Goldhaber, Dan, Thomas J. Kane, Andrew McEachin, Emily Morton, Tyler Patterson, and Douglas O. Staiger. 2022. "The Consequences of Remote and Hybrid Instruction During the Pandemic." National Bureau of Economic Research, Inc 30010.
- Gottfried, Adele Eskeles, Kathleen Suzanne Johnson Preston, Allen W. Gottfried, Pamella H. Oliver, Danielle E. Delany, and Sirena M. Ibrahim. 2016. "Pathways from Parental Stimulation of Children's Curiosity to High School Science Course Accomplishments and Science Career Interest and Skill." *International Journal of Science Education*, 38(12): 1972–1995.
- Granger, E. M., T. H. Bevis, Y. Saka, S. A. Southerland, V. Sampson, and R. L. Tate. 2012. "The Efficacy of Student-Centered Instruction in Supporting Science Learning." *Science*, 338(6103): 105–108.
- **Gruber, Matthias J., and Charan Ranganath.** 2019. "How Curiosity Enhances Hippocampus-Dependent Memory: The Prediction, Appraisal, Curiosity, and Exploration (PACE) Framework." *Trends in Cognitive Sciences*, 23(12): 1014–1025.
- Gruber, Matthias J., Bernard D. Gelman, and Charan Ranganath. 2014. "States of Curiosity Modulate Hippocampus-Dependent Learning via the Dopaminergic Circuit." Neuron, 84(2): 486–496.

- Gust, Sarah, Eric A. Hanushek, and Ludger Woessmann. 2022. "Global Universal Basic Skills: Current Deficits and Implications for World Development."
- Hartung, Freda-Marie, and Britta Renner. 2013. "Social Curiosity and Gossip: Related but Different Drives of Social Functioning." *PLOS ONE*, 8(7): e69996.
- Hjort, Jonas, Diana Moreira, Gautam Rao, and Juan Francisco Santini. 2021. "How Research Affects Policy: Experimental Evidence from 2,150 Brazilian Municipalities." *American Economic Review*, 111(5): 1442–1480.
- **James, William.** 1983. Talks to Teachers on Psychology and to Students on Some of Life's Ideals. Harvard University Press.
- **Jirout, Jamie, and David Klahr.** 2012. "Children's Scientific Curiosity: In Search of an Operational Definition of an Elusive Concept." *Developmental Review*, 32(2): 125–160.
- Kahn, Shulamit, and Donna Ginther. 2017. "Women and STEM."
- Kashdan, Todd B., and Paul J. Silvia. 2009. "Curiosity and Interest: The Benefits of Thriving on Novelty and Challenge." In Oxford Handbook of Positive Psychology, 2nd Ed. Oxford Library of Psychology, 367–374. New York, NY, US:Oxford University Press.
- Kashdan, Todd B., David J. Disabato, Fallon R. Goodman, and Patrick E. McKnight. 2020. "The Five-Dimensional Curiosity Scale Revised (5DCR): Briefer Subscales While Separating Overt and Covert Social Curiosity." Personality and Individual Differences, 157: 109836.
- Kremer, Michael, Conner Brannen, and Rachel Glennerster. 2013. "The Challenge of Education and Learning in the Developing World." *Science*, 340(6130): 297–300.
- Kremer, Michael, Paul Glewwe, and Sylvie Moulin. 2009. "Many Children Left Behind? Textbooks and Test Scores in Kenya." American Economic Journal: Applied Economics, 1(1 (January 2009)): 112–135.
- **Litman, Jordan A., and Charles D. Spielberger.** 2003. "Measuring Epistemic Curiosity and Its Diversive and Specific Components." *Journal of Personality Assessment*, 80(1): 75–86.
- **Litman, Jordan A., and Mark V. Pezzo.** 2007. "Dimensionality of Interpersonal Curiosity." *Personality and Individual Differences*, 43(6): 1448–1459.

- Litman, Jordan A., Robert P. Collins, and Charles D. Spielberger. 2005. "The Nature and Measurement of Sensory Curiosity." *Personality and Individual Differences*, 39(6): 1123–1133.
- **Loewenstein, George.** 1994. "The Psychology of Curiosity: A Review and Reinterpretation." *Psychological Bulletin*, 116(1): 75–98.
- OECD. 2013. "TALIS User Guide for the International Database."
- Raven, John C, and John Hugh Court. 1998. Raven's Progressive Matrices and Vocabulary Scales. Vol. 759, Oxford pyschologists Press Oxford.
- Romano, Joseph P, and Michael Wolf. 2005. "Exact and Approximate Stepdown Methods for Multiple Hypothesis Testing." *Journal of the American Statistical Association*, 100(469): 94–108.
- Shah, Prachi E., Heidi M. Weeks, Blair Richards, and Niko Kaciroti. 2018. "Early Childhood Curiosity and Kindergarten Reading and Math Academic Achievement." *Pediatric Research*, 84(3): 380–386.
- Singh, Abhijeet, Mauricio Romero, and Karthik Muralidharan. 2022. "Covid-19 Learning Loss and Recovery: Panel Data Evidence from India."
- Sleddens, Ester F. C., Stef P. J. Kremers, Nanne K. De Vries, and Carel Thijs. 2013. "Measuring Child Temperament: Validation of a 3-item Temperament Measure and 13-item Impulsivity Scale." *European Journal of Developmental Psychology*, 10(3): 392–401.
- Sosu, Edward M. 2013. "The Development and Psychometric Validation of a Critical Thinking Disposition Scale." Thinking Skills and Creativity, 9: 107–119.
- Terrenghi, Ilaria, Barbara Diana, Valentino Zurloni, Pier Cesare Rivoltella, Massimiliano Elia, Marta Castañer, Oleguer Camerino, and M. Teresa Anguera. 2019. "Episode of Situated Learning to Enhance Student Engagement and Promote Deep Learning: Preliminary Results in a High School Classroom." Frontiers in Psychology, 10.
- Vogl, Elisabeth, Reinhard Pekrun, Kou Murayama, and Kristina Loderer. 2019a.
  "Surprised-Curious-Confused: Epistemic Emotions and Knowledge Exploration." Emotion, No Pagination Specified-No Pagination Specified.

Vogl, Elisabeth, Reinhard Pekrun, Kou Murayama, Kristina Loderer, and Sandra Schubert. 2019b. "Surprise, Curiosity, and Confusion Promote Knowledge Exploration: Evidence for Robust Effects of Epistemic Emotions." Frontiers in Psychology, 10.

von Stumm, Sophie, Benedikt Hell, and Tomas Chamorro-Premuzic. 2011. "The Hungry Mind: Intellectual Curiosity Is the Third Pillar of Academic Performance." *Perspectives on Psychological Science*, 6(6): 574–588.

# Appendix

### A Tables

Table A1: Balance at Baseline

	N	Control Mean	Treatment Mean	Diff pvalue
Student Characteristics	11	Mean	Mean	pvarue
Male	13039	0.51	0.51	0.96
Age in Months	13039	112.43	112.74	0.90 $0.25$
Fluid IQ Score	10039 $10912$	-0.08	-0.05	0.23 $0.82$
Math Score	10912 $10922$	-0.08	-0.03 -0.04	0.82 $0.90$
Verbal Score	10922	-0.05	-0.02	0.96
Curiosity	13039	-0.05	-0.03	0.59
Risk Attitude	13039	2.61	2.58	0.77
Ambiguity Attitude	10409	2.47	2.42	0.84
Gender Stereotypes	10613	0.03	0.02	0.36
Home - Computer	10758	0.50	0.52	0.51
Home - Internet	10738	0.80	0.80	0.72
Siblingship Size	10814	2.74	2.72	0.90
Birth Order	10814	2.61	2.59	0.99
Teacher Characteristics				
Male	425	0.27	0.29	0.68
Age	425	45.49	44.66	0.25
Fluid IQ Score	425	17.76	17.70	0.81
Cognitive Empathy Score	425	23.05	22.94	0.86
Married	425	0.83	0.85	0.62
Number of children	425	1.81	1.71	0.18
Teaching experience in Years	425	21.01	20.50	0.41
University Graduate	425	0.94	0.95	0.50
Curiosity	425	-0.05	0.09	0.15
Gender Stereotypes	425	-0.05	-0.05	0.88
Growth Mindset	425	0.01	0.06	0.56
Professional Attachment	425	0.01	-0.00	0.85
Competence Beliefs	424	0.01	0.07	0.50
Modern Teaching	425	0.00	0.02	0.89
Extrinsic Motivator	425	-0.03	-0.11	0.20
Warmth	425	-0.08	-0.03	0.51
Classroom Characteristics		0.00	0.00	0.02
Classroom size	425	31.14	30.78	0.60
Refugee Share	425	0.07	0.07	1.00

The table presents the balance at baseline for the pooled sample. The p-values from the test of equality between control and treatment are shown in the last column. Test scores and survey items are standardized to have a mean zero and a standard deviation of 1.

Table A2: Multiple Hypothesis Testing

	Original P-Value	Sharpened O-Value	Romano Wolf P-Value
Panel 1: Student Outcomes	Oliginal Value	Sharpehea & Varae	Teomano Won i Varae
Experimental Task			
Science Related Booklet	0.001	0.003	0.012
Non-Science Booklet	0.328	0.161	0.515
No Booklet	0.000	0.002	0.012
Overall Curiosity	0.006	0.014	0.046
Science Curiosity	0.000	0.001	0.008
Non-Science Curiosity	0.664	0.319	0.687
Retention	0.009	0.016	0.050
Science Retention	0.010	0.016	0.050
Non-Science Retention	0.028	0.027	0.092
Achievement & Aspirations			
Science	0.007	0.014	0.084
Maths	0.542	0.261	0.908
Verbal	0.206	0.135	0.641
University Aspiration	0.086	0.063	0.413
Science Aspiration	0.001	0.003	0.016
Engineering Aspiration	0.889	0.402	0.948
Medical Aspiration	0.745	0.350	0.948
Non-STEM Aspiration	0.052	0.047	0.313
Students' Beliefs & Attitudes			
Grit	0.215	0.135	0.439
Impulsivity	0.456	0.224	0.521
Risk	0.074	0.058	0.253
Ambiguity	0.010	0.016	0.050
Critical Thinking	0.004	0.011	0.024
Curiosity Survey	0.000	0.001	0.002
Science Curiosity	0.000	0.001	0.002
Panel 2: Teacher Outcomes			
Curiosity	0.000	0.006	0.012
Modern Teaching	0.038	0.130	0.361
Warmth	0.125	0.333	0.677
Extrinsic Motivator	0.193	0.448	0.792
Growth Mindset	0.002	0.011	0.028
Professional Attachment	0.992	1.000	0.998
Competence Beliefs	0.487	0.740	0.932
Gender Stereotypes	0.403	0.675	0.922
Critical Thinking	0.895	1.000	0.998
Science Subject Knowledge	0.925	1.000	0.998
Knowledge Retention	0.326	0.616	0.894

The table presents estimation results for sharpened False Discovery Rate (FDR) q-values (Anderson (2008)) and adjusted p-values via Romano and Wolf (2005) multiple hypothesis correction. To accommodate Romano-Wolf correction to control for family wise error rate (FWER), we group our outcome variables into three, namely (i) experimental outcomes, (ii) achievement and aspiration related outcomes, (iii) beliefs and attitudes.

Table A3: Balance at Baseline for Table 5 Panel 3 (Network Effects)

	N	Control Mean	Treatment Mean	Diff pvalue
Student Characteristics				
Male	1207	0.50	0.50	0.82
Age in Months	1207	111.65	111.97	0.51
Fluid IQ Score	1100	0.05	0.04	0.66
Math Score	1101	0.09	0.05	0.40
Verbal Score	1101	0.13	0.10	0.49
Curiosity	1207	-0.01	0.01	0.92
Risk Attitude	1207	2.51	2.56	0.29
Ambiguity Attitude	1049	2.39	2.47	0.21
Gender Stereotypes	1079	-0.00	0.02	0.14
Home - Computer	1087	0.51	0.51	0.77
Home - Internet	1081	0.80	0.82	0.32
Siblingship Size	1090	2.54	2.60	0.59
Birth Order	1090	2.53	2.47	0.74

The table presents the balance at baseline for the restricted sample described in Table 5 Panel 3. The sample contains students who did not receive any booklet but have at least one person in their network who has received the booklet of their choice. The p-values from the test of equality between control and treatment are shown in the last column.

Table A4: Heterogeneous Treatment Effects - Gender

Panel 1: Choice of Booklet

	Science Related	Non-Science Related	No booklet
Treatment = Girls	0.079***	-0.044**	-0.035***
	(0.02)	(0.02)	(0.01)
Treatment = Boys	0.001	0.022	-0.023***
	(0.02)	(0.02)	(0.01)
P-Value : Girls=Boys	0.003	0.009	0.216
Control Mean - Girls	0.50	0.44	0.07
Control Mean - Boys	0.49	0.44	0.06
Observations	10870	10870	10870
Number of Schools	134	134	134

Panel 2: Level of Curiosity

	Curiosity	Science Curiosity	Non-Science Curiosity
Treatment = Girls	0.146***	0.173***	-0.056*
	(0.05)	(0.03)	(0.03)
Treatment = Boys	$0.075^{*}$	0.028	0.034
	(0.04)	(0.03)	(0.04)
P-Value : Girls=Boys	0.064	0.001	0.063
Control Mean - Girls	-0.07	-0.01	-0.04
Control Mean - Boys	0.05	0.01	0.04
Observations	10864	10863	10863
Number of Schools	134	134	134

Estimates are obtained via OLS. The dependent variables are binary indicators of choosing a science-related booklet (science, space, vehicles, human body, and animals) in column 1, choosing a nonscience-related booklet (history, sports, and cartoons) in column 2, and choosing no booklet option in column 3. Standard errors are clustered at the school level and are reported in parentheses. Covariates include gender, age, fluid IQ, risk tolerance, survey measure of curiosity, math and verbal scores as individual baseline characteristics, class size, the share of refugees, teacher gender, experience, and fluid IQ as baseline classroom and teacher characteristics. Grade and district fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

Panel 1: Knowledge Retention - Full Sample

		Short Term			Long Term		
	Retention	Science Retention	Non-Science Retention		Science Retention	Non-Science Retention	
$\overline{\text{Treatment} = \text{Girls}}$	0.118**	0.106**	0.088*	0.033	0.056	-0.009	
	(0.05)	(0.04)	(0.05)	(0.05)	(0.05)	(0.06)	
Treatment = Boys	$0.119^{**}$	$0.107^{**}$	0.088**	$0.125^{*}$	$0.156^{**}$	0.039	
	(0.05)	(0.04)	(0.04)	(0.06)	(0.06)	(0.06)	
P-Value : Girls=Boys	0.982	0.984	0.987	0.238	0.229	0.563	
Control Mean - Girls	-0.09	-0.09	-0.06	-0.05	0.02	-0.12	
Control Mean - Boys	0.07	0.07	0.04	0.05	-0.02	0.12	
Observations	10590	10590	10590	2426	2426	2426	
Number of Schools	134	134	134	50	50	50	

Panel 2: Knowledge Retention - Half Half

	Short Term			Long Term		
	Retention	Science Retention	Non-Science Retention		Science Retention	Non-Science Retention
$\overline{\text{Treatment} = \text{Girls}}$	0.113**	0.101**	$0.085^*$	0.088	0.087	0.055
	(0.06)	(0.05)	(0.05)	(0.08)	(0.08)	(0.09)
Treatment = Boys	$0.114^{**}$	0.104**	$0.083^{*}$	$0.175^{*}$	$0.185^{**}$	0.095
	(0.05)	(0.05)	(0.05)	(0.09)	(0.08)	(0.09)
P-Value : Girls=Boys	0.978	0.941	0.971	0.474	0.386	0.755
Control Mean - Girls	-0.12	-0.11	-0.09	-0.07	0.03	-0.16
Control Mean - Boys	0.04	0.05	0.02	-0.01	-0.07	0.07
Observations	9037	9037	9037	1336	1336	1336
Number of Schools	134	134	134	50	50	50

Panel 3: Test Scores

	Short Term			Long Term		
	Science	Maths	Verbal	Science	Maths	Verbal
Treatment = Girls	0.059	0.017	0.023	0.101**	-0.037	-0.004
	(0.04)	(0.03)	(0.03)	(0.05)	(0.05)	(0.05)
Treatment = Boys	$0.097^{***}$	0.017	0.043	0.039	-0.000	-0.038
	(0.03)	(0.03)	(0.03)	(0.06)	(0.05)	(0.06)
P-Value : Girls=Boys	0.346	0.976	0.518	0.413	0.547	0.612
Control Mean - Girls	-0.01	-0.02	0.11	-0.02	0.02	0.12
Control Mean - Boys	-0.01	0.03	-0.12	0.02	-0.02	-0.12
Observations	9949	10400	10680	2426	2426	2426
Number of Schools	134	134	134	50	50	50

Estimates are obtained via OLS. The dependent variables are standardized booklet test scores (knowledge retention) in Panels 1 and 2, and standardized subject test scores in Panel 3. The first 3 columns give short-term results using the pooled sample, and the last 3 provide the long-term results of Study 1. Standard errors are clustered at the school level and are reported in parentheses. Covariates include gender, age, fluid IQ, risk to ance, survey measure of curiosity, math and verbal scores as individual baseline characteristics, class size, the share of refugees, teacher gender, experience, and fluid IQ as baseline classroom and teacher characteristics. Grade and district fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*\*, 5% \*\*\*, and 10% \* levels.

Table A6: Heterogeneous Treatment Effects - Gender

Panel 1: Short Term

	University	Science	Engineering	Medical	Non-STEM
Treatment = Girls	0.008	0.031***	0.003	-0.010	-0.024
	(0.01)	(0.01)	(0.01)	(0.01)	(0.02)
Treatment = Boys	0.008	0.016	-0.001	0.005	-0.019
	(0.01)	(0.01)	(0.01)	(0.01)	(0.02)
P-Value : Girls=Boys	0.986	0.314	0.745	0.351	0.824
Control Mean - Girls	0.96	0.08	0.06	0.23	0.63
Control Mean - Boys	0.94	0.15	0.17	0.09	0.58
Observations	10693	10186	10186	10186	10186
Number of Schools	134	134	134	134	134

Panel 2: Long Term

	University	Science	Engineering	Medical	Non-STEM
Treatment = Girls	0.003	0.025	0.006	-0.042	0.012
	(0.01)	(0.02)	(0.02)	(0.03)	(0.04)
Treatment = Boys	0.013	0.002	0.022	0.017	-0.041
	(0.01)	(0.03)	(0.03)	(0.02)	(0.04)
P-Value : Girls=Boys	0.558	0.499	0.669	0.093	0.358
Control Mean - Girls	0.96	0.10	0.05	0.30	0.56
Control Mean - Boys	0.94	0.16	0.19	0.13	0.52
Observations	2320	2182	2182	2182	2182
Number of Schools	50	50	50	50	50

Estimates are obtained via OLS. The dependent variables are binary choice variables of intention to go to university, intention to choose a science major, engineering major, medicine, and non-STEM major. Panel 1 presents short-term results from the pooled sample, and Panel 2 long-term results from Study 1. Standard errors are clustered at the school level and are reported in parentheses. Covariates include gender, age, fluid IQ, risk tolerance, survey measure of curiosity, math and verbal scores as individual baseline characteristics, class size, the share of refugees, teacher gender, experience, and fluid IQ as baseline classroom and teacher characteristics. Grade and district fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

Table A7: Heterogeneous Treatment Effects - IQ

Panel 1: Choice of Booklet

	Science Related	Non-Science Related	No booklet
Treatment = Low IQ	0.028	0.006	-0.034***
	(0.02)	(0.01)	(0.01)
Treatment = High IQ	$0.047^{***}$	-0.021	-0.026***
	(0.01)	(0.01)	(0.01)
$\overline{\text{P-Value} : \text{Low} = \text{High}}$	0.362	0.153	0.406
Control Mean - Low IQ	0.48	0.45	0.07
Control Mean - High IQ	0.50	0.43	0.06
Observations	10870	10870	10870
Number of Schools	134	134	134

Panel 2: Level of Curiosity

	Curiosity	Science Curiosity	Non-Science Curiosity
Treatment = Low IQ	0.061	0.048	0.002
	(0.05)	(0.04)	(0.03)
Treatment = High IQ	$0.141^{***}$	$0.133^{***}$	-0.019
	(0.04)	(0.03)	(0.03)
P-Value : Low = High	0.062	0.062	0.587
Control Mean - Low IQ	-0.05	-0.04	0.00
Control Mean - High IQ	0.02	0.02	-0.00
Observations	10864	10863	10863
Number of Schools	134	134	134

Estimates are obtained via OLS. The dependent variables are binary indicators of choosing a science-related booklet (science, space, vehicles, human body, and animals) in column 1, choosing a nonscience-related booklet (history, sports, and cartoons) in column 2, and choosing no booklet option in column 3. Standard errors are clustered at the school level and are reported in parentheses. Covariates include gender, age, fluid IQ, risk tolerance, survey measure of curiosity, math and verbal scores as individual baseline characteristics, class size, the share of refugees, teacher gender, experience, and fluid IQ as baseline classroom and teacher characteristics. Grade and district fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

Panel 1: Knowledge Retention - Full Sample

	Short Term			Long Term		
	Retention	Science Retention	Non-Science Retention	Retention	Science Retention	Non-Science Retention
Treatment = Low IQ	0.064	0.054	0.052	0.003	0.013	-0.010
	(0.04)	(0.04)	(0.04)	(0.06)	(0.06)	(0.06)
Treatment = High IQ	$0.153^{**}$	0.140**	0.111**	$0.127^{**}$	$0.166^{***}$	0.029
	(0.06)	(0.05)	(0.05)	(0.06)	(0.06)	(0.06)
$\overline{P}$ -Value : Low = High	0.109	0.117	0.238	0.209	0.112	0.642
Control Mean - Low IQ	-0.13	-0.14	-0.06	-0.16	-0.04	-0.23
Control Mean - High IQ	0.07	0.08	0.03	0.10	0.03	0.15
Observations	10590	10590	10590	2426	2426	2426
Number of Schools	134	134	134	50	50	50

Panel 2: Knowledge Retention - Half Half

	Short Term			Long Term		
	Retention	Science Retention	Non-Science Retention	Retention	Science Retention	Non-Science Retention
Treatment = Low IQ	0.069	0.046	0.073*	0.033	0.036	0.017
	(0.05)	(0.04)	(0.04)	(0.07)	(0.07)	(0.07)
Treatment = High IQ	0.144**	$0.140^{**}$	0.092	$0.202^{**}$	$0.207^{**}$	0.119
	(0.07)	(0.06)	(0.06)	(0.09)	(0.09)	(0.08)
P-Value : Low = High	0.234	0.121	0.729	0.132	0.149	0.260
Control Mean - Low IQ	-0.16	-0.16	-0.09	-0.14	-0.02	-0.24
Control Mean - High IQ	0.04	0.05	0.01	0.04	-0.03	0.11
Observations	9037	9037	9037	1336	1336	1336
Number of Schools	134	134	134	50	50	50

Panel 3: Test Scores

	Short Term				1	
	Science	Maths	Verbal	Science	Maths	Verbal
$\overline{\text{Treatment} = \text{Low IQ}}$	0.064*	0.014	0.054*	0.030	-0.127**	-0.005
	(0.03)	(0.04)	(0.03)	(0.06)	(0.06)	(0.08)
Treatment = High IQ	0.088**	0.020	0.023	$0.095^{*}$	0.049	-0.032
	(0.04)	(0.04)	(0.04)	(0.05)	(0.06)	(0.05)
P-Value : Low = High	0.662	0.913	0.419	0.397	0.054	0.774
Control Mean - Low IQ	-0.36	-0.20	-0.43	-0.32	-0.34	-0.21
Control Mean - High IQ	0.21	0.13	0.26	0.21	0.22	0.14
Observations	9949	10400	10680	2426	2426	2426
Number of Schools	134	134	134	50	50	50

Estimates are obtained via OLS. The dependent variables are standardized booklet tests scores (knowledge retention) in Panel 1 and 2, standardized subject test scores in Panel 3. The first 3 columns give short-term results using the pooled sample, the last 3 give the long-term results of Study 1. Standard errors are clustered at the school level and are reported in parentheses. Covariates include gender, age, fluid IQ, risk tolerance, sur40 measure of curiosity, math and verbal scores as individual baseline characteristics, class size, the share of refugees, teacher gender, experience, and fluid IQ as baseline classroom and teacher characteristics. Grade and district fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

Table A9: Heterogeneous Treatment Effects - IQ

Panel 1: Short Term

	University	Science	Engineering	Medical	Non-STEM
Treatment = Low IQ	0.006	0.016	0.006	0.006	-0.029*
	(0.01)	(0.01)	(0.01)	(0.01)	(0.02)
Treatment = High IQ	0.009*	0.028***	-0.002	-0.008	-0.018
	(0.01)	(0.01)	(0.01)	(0.01)	(0.01)
P-Value : Low = High	0.797	0.431	0.526	0.327	0.602
Control Mean - Low IQ	0.94	0.12	0.10	0.14	0.64
Control Mean - High IQ	0.96	0.11	0.13	0.18	0.59
Observations	10693	10186	10186	10186	10186
Number of Schools	134	134	134	134	134

Panel 2: Long Term

	University	Science	Engineering	Medical	Non-STEM
Treatment = Low IQ	0.017	0.003	0.004	-0.018	0.012
	(0.01)	(0.03)	(0.02)	(0.03)	(0.04)
Treatment = High IQ	0.004	0.020	0.019	-0.010	-0.029
	(0.01)	(0.02)	(0.02)	(0.02)	(0.03)
P-Value : Low = High	0.450	0.630	0.599	0.810	0.389
Control Mean - Low IQ	0.92	0.13	0.08	0.22	0.57
Control Mean - High IQ	0.97	0.12	0.14	0.21	0.52
Observations	2320	2182	2182	2182	2182
Number of Schools	50	50	50	50	50

Estimates are obtained via OLS. The dependent variables are binary choice variables of intention to go to university, choose science major, engineering major, medicine and non-STEM major. Panel 1 presents short term results from the pooled sample, Panel 2 long-term results from Study 1. Standard errors are clustered at the school level and are reported in parentheses. Covariates include gender, age, fluid IQ, risk tolerance, survey measure of curiosity, math and verbal scores as individual baseline characteristics, class size, the share of refugees, teacher gender, experience, and fluid IQ as baseline classroom and teacher characteristics. Grade and district fixed effects are also included. Asterisks indicate statistical significance at the 1% \*\*\*, 5% \*\*, and 10% \* levels.

# B Figures

Figure A1: Covers of the Booklets



Experimentation Figure 1. Teacher reported Implementation

Figure A2: Implementation Intensity

The figure depicts the program implementation intensity reported by treated teachers at endline. Teachers were given a 10cm line that has a moving cursor to report the level they believe represents their implementation intensity, zero representing no implementation, and 10 a 100% implementation.

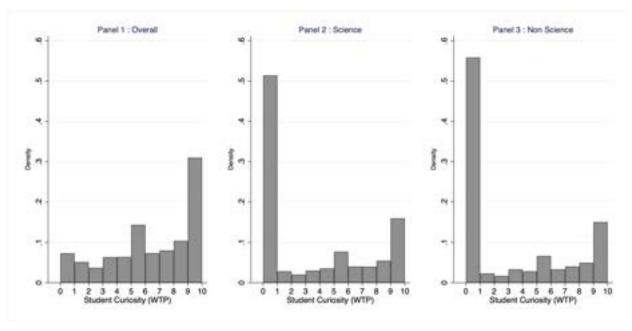


Figure A3: Student Curiosity Distribution (WTP)

Figures depict the distribution of the number of tokens forgone for a booklet (Panel 1), for a science-related booklet (Panel 2), and for a non-science related booklet (Panel 3).

# Online Appendix- Not For Publication

## A Additional Tables

**Table B1:** Balance at Baseline: Study 1

		Control	Treatment	Diff
	N	Control Mean	Mean	pvalue
Student Characteristics				P
Male	3786	0.51	0.51	0.53
Age in Months	3786	110.89	111.56	0.47
Fluid IQ Score	3376	-0.01	0.05	0.82
Math Score	3386	-0.01	0.08	0.73
Verbal Score	3386	0.05	0.12	0.92
Curiosity	3786	0.06	0.07	0.90
Risk Attitude	3786	2.13	2.08	0.61
Ambiguity Attitude	2873	1.94	1.97	0.66
Gender Stereotypes	3254	0.01	0.01	0.64
Home - Computer	3286	0.53	0.54	0.83
Home - Internet	3273	0.67	0.65	0.30
Siblingship Size	3324	2.71	2.67	0.68
Birth Order	3324	2.63	2.57	0.99
Teacher Characteristics				
Male	129	0.38	0.32	0.55
Age	129	43.05	42.19	0.39
Fluid IQ Score	129	19.13	19.17	1.00
Cognitive Empathy Score	129	22.65	22.91	0.57
Married	129	0.83	0.83	0.84
Number of children	129	1.58	1.51	0.46
Teaching experience in Years	129	18.95	18.17	0.37
University Graduate	129	0.92	0.93	0.66
Curiosity	129	-0.08	0.12	0.22
Gender Stereotypes	129	-0.03	-0.24	0.12
Growth Mindset	129	0.06	0.08	0.86
Professional Attachment	129	-0.08	0.15	0.18
Competence Beliefs	128	-0.00	0.21	0.18
Modern Teaching	129	-0.04	0.07	0.45
Extrinsic Motivator	129	-0.01	-0.20	0.12
Warmth	129	0.02	0.13	0.28
Classroom Characteristics				
Classroom size	129	28.20	30.35	0.33
Refugee Share	129	0.14	0.13	0.76

The table presents the balance at baseline for Study 1 sample. The p-values from the test of equality between control and treatment are shown in the last column. Test scores and survey items are standardized to have a mean zero and a standard deviation of 1.

Table B2: Balance at Baseline: Study 2

	N	Control Mean	Treatment Mean	Diff pvalue
Student Characteristics				
Male	9253	0.51	0.51	0.68
Age in Months	9253	113.02	113.26	0.37
Fluid IQ Score	7536	-0.11	-0.10	0.90
Math Score	7536	-0.10	-0.10	0.99
Verbal Score	7536	-0.09	-0.08	0.99
Curiosity	9253	-0.09	-0.07	0.49
Risk Attitude	9253	2.80	2.80	0.97
Ambiguity Attitude	7536	2.65	2.61	0.65
Gender Stereotypes	7359	0.04	0.02	0.24
Home - Computer	7472	0.49	0.52	0.42
Home - Internet	7465	0.85	0.87	0.29
Siblingship Size	7490	2.75	2.74	1.00
Birth Order	7490	2.61	2.60	0.99
Teacher Characteristics				
Male	296	0.22	0.27	0.33
Age	296	46.54	45.75	0.42
Fluid IQ Score	296	17.17	17.06	0.79
Cognitive Empathy Score	296	23.22	22.96	0.58
Married	296	0.83	0.86	0.46
Number of children	296	1.91	1.80	0.27
Teaching experience in Years	296	21.89	21.53	0.73
University Graduate	296	0.95	0.96	0.62
Curiosity	296	-0.03	0.07	0.34
Gender Stereotypes	296	-0.05	0.03	0.51
Growth Mindset	296	-0.01	0.06	0.57
Professional Attachment	296	0.05	-0.07	0.35
Competence Beliefs	296	0.02	0.01	0.93
Modern Teaching	296	0.02	-0.01	0.71
Extrinsic Motivator	296	-0.03	-0.07	0.63
Warmth	296	-0.12	-0.11	0.94
Classroom Characteristics				
Classroom size	296	32.41	30.97	0.23
Refugee Share	296	0.04	0.04	0.74

The table presents the balance at baseline for Study 2 sample. The p-values from the test of equality between control and treatment are shown in the last column. Test scores and survey items are standardized to have a mean zero and a standard deviation of 1.

#### **B** Additional Figures

Figure B1: The Attrition Pattern of Study 1

The figure depicts the pattern and the reason for attrition in Study 1 during the long-term data collection. The Turkish Ministry of Education was able to locate 80% of our original participants in its official database. "Non-Study Site" refers to students who left the province of Mersin, "Private" refers to those who left the public education system for a private school, and "No Info" refers to those considered missing. Among the officially registered in Mersin, a total of 177 students were absent during our visit for various usual reasons such as illness. A total of 77 students were declared permanently absent (never showed up) by school administrators, 291 were reported to have transferred to another school, and 76 students were dispersed too far for us to go after.

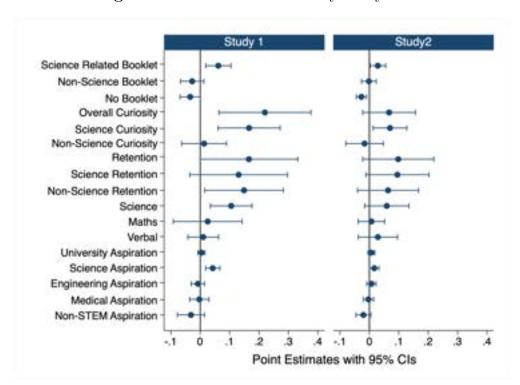


Figure B2: Treatment Effects -by Study Sites

The figure depicts the estimated treatment effects and their 95% confidence intervals for all outcomes considered in the study. Standard errors are clustered at the school level. The vertical line indicates a treatment effect of 0. The first three outcomes are the choice of booklets, the following three are curiosity levels based on the experimental task, followed by the booklet performance (7-9), subject test scores (10-12), and educational aspirations(13-17). Covariates include gender, age, fluid IQ, risk tolerance, survey measure of curiosity, math and verbal scores as individual baseline characteristics, class size, the share of refugees, teacher gender, experience, and fluid IQ as baseline classroom and teacher characteristics. Grade and district fixed effects are also included.

## C Implementation Items and Moments

Figure B3: Curious Classroom Toolkit



Figure B4: Creating Teachable Moments via Humor, Mystery and Astonishment



Figure B5: Examples of Children's Activities (Mystery Box)



#### D Instructions for Incentivized Games

#### D.1 Curiosity Task

Hi everybody. We will play some fun games with you today. By playing these games, you will have a chance to earn gift tokens from us, with which you can get any gift you want from our gift bag [show the items in the gift bag]. The number of gifts you will receive will depend on your choices in these games. To get the gifts, you need to collect tokens, as each gift in our basket has a different token value. The more tokens you have, the more gifts you will be able to get at the end of our visit.

Each game has its own rules, and we will slowly explain all of them. But our main rule is discretion. You will need to make all our choices discretely, without showing anyone. Do you understand this rule? Excellent!

Now, see that we brought 8 booklets to you today. These booklets contain some incredible facts that most people do not know. [Start introducing them one by one]. This is the Space booklet. It has shocking facts in it. [Show animals], this is a booklet that contains astonishing facts about animals. [Go through each booklet in the same manner and always in the same order].

Now, we would like you to rank the booklets from most attractive to the least according to your own taste. Please type 1 beside the picture of the booklet that interests you the most, 2 for the second most interesting you find, and keep going until 8, which would be the booklet least interesting to you. [Make sure everyone finishes their ranking and press continue before the next step].

Now, if you want, you can purchase one (and only one) of these booklets from us. How? Well, first, know that we are giving all of you 10 tokens. All of you have 10 tokens. You can use these tokens to get some of these nice stationery items from us. You can also get one booklet if you want. You don't have to get a booklet. You can convert all your tokens to gift items if you wish to. [Make sure children understand they do not have to purchase a booklet]. But if you do want a booklet, you need to first indicate which booklet you want to purchase on your tablet. Then you need to indicate how many of those 10 tokes you would be willing to give us back to purchase this booklet. You can say zero, meaning you don't want a booklet and want to convert all your tokens into gifts. Or you can say any number from 1 to 10.

But how do you really purchase a booklet? One of two things can happen in your classroom. You can be classroom type A or B. Let's see what happens in type A classrooms: Let's say student A decides to forgo 3 tokens, student B 5, and student C 7 tokens. Here is what we will do. We will pick a number from this bag. The bag contains folded little papers. In each paper, a number between 1 to 10 is written. [Show the black bag and show the little paper pieces]. The number we pull from this bag will be the price of a booklet for this classroom.

[Now, start giving the examples based on the 3 students above]. Let's say we picked number 8. Then we will look at everyone's decision of willingness to pay for their preferred booklet. Student A marks 3. She can't get the booklet she wants because the price is 8. Instead, we will convert all her 10 tokens into gifts. The same goes for students B and C because their willingness to pay fell under the price of the booklet in this classroom.

But let's say we pick the number 5 instead of 8, so the price is 5. Student A still won't get a booklet and will receive 10 tokens worth of gifts. Student B, however, will give us her 5 tokens, get the booklet she wants and convert her remaining 5 tokens into gifts. What about student C? Well, she says she is willing to forgo 7 tokens but does she need to? NO. The price is 5, why should she? So we will get 5 tokens from her, give her the booklet she wants, and she will convert the remaining 5 tokens to gifts, just like student B.

What about a student who states zero willingness to pay? Well, she will not receive a booklet at any price. What about a student who states 10? She will certainly receive a booklet in the classroom type A.

What if your classroom is type B, which is much more likely as most classrooms will be type B. If your classroom is a type B, no matter how much you are willing to pay for a booklet, and no matter which booklet you prefer, a random half of the classroom will receive booklets, and the other half will not. We will pick half the students randomly from your class list.

Now, time to make decisions. First, tap the booklet you want to purchase. Don't forget there is an option that says "I do not want a booklet". You can tap that if you don't want a booklet. After making your choice, please tap the number of tokens you are willing to forgo to get the booklet you choose. [Make sure everyone makes their decisions and press continue.]

- Implementation, Type A (Market Price): Please pick a number from the black bag. Distribute the booklets accordingly.
- Implementation Type B (Half-Half): Please select the random half of the classroom using the class list and distribute the booklets only to them. Make sure every classroom has all 8 booklets.

#### D.2 Risk and Ambiguity Games

Now we will play two games. [Type Game 1 and Game 2 on the board]. These two games are almost identical to each other. You will earn some gifts from these games. But you will

collect the gifts from only one of the games, i.e., the gifts will not accumulate. We will pick one of these two games randomly for this classroom at the end of the visit, and you will get your gifts based on the decisions you make for that game. Now, let me explain the games. Game 1: We will give you 5 tokens for this game. You can convert these tokens into small gifts in our bag [show all the gifts]. Now, think about a bucket [draw a bucket on the board]. You can put some of your tokens in this bucket if you want. You don't have to. If you don't, you have your 5 tokens, no problem. But what happens if you put some of your tokens in the bucket? Then, you draw a ball from this black bag [show the black bag]. There are two balls in this bag. One is yellow, and one is purple [show the balls]. The tokens you put in the bucket triple if you draw the yellow ball. You lose all the tokens you put in the bucket if you draw the purple ball. But not the ones you didn't put in the bucket. Tokens you don't put in the bucket are always safe.

Let's see some examples now: If you put none of your tokens in the bucket. What happens? NOTHING. You have 5 tokens. Let's say you put 1 token in the bucket. You have 4 safe ones left. Nothing happens to them. Then you draw a ball from the bag. If you draw the yellow ball, your 1 token becomes 3 tokens. Add to that your 4 safe ones. You now have 7 tokens. But what if you pick the purple ball. Then you lose that 1 token you put in the bucket, and you have 4 tokens. Now, let's say you put 2 tokens in the bucket. [Go on until you give the example of 5 tokens].

Now, decide how many tokens you want to put in the bucket. Please tap the number on your tablet and press continue.

Game 2: Now, we will play the second game. The second game is the same as the first game. You have 5 tokens, there is a bucket, and the tokens you put in the bucket triple if the yellow ball is drawn. They disappear if the purple ball is drawn. All the same. Except now, you don't know the colors of the balls in this new bag [pick the other bag, so children see this is not the same bag as in game 1]. Both balls can be yellow. In that case, you certainly win. Both balls can be purple, in which case you certainly lose. Or, one of them may be yellow, the other purple, as in Game 1. The fact is, you do not know.

Now, please decide how many tokens you will put in the bucket. Please tap the number on your tablet and press continue.

## E Survey Inventories

We provide some example questions from our student and teacher surveys below. The full inventory for both is available upon request.

 Table B3: Student Survey Inventories

Inventory	Exemplary Items				
4-point likert scale: completely agree, agree, disagree, completely disagree					
Curiosity	There are always questions on my mind.				
Curiosity	When I hear a word that I do not know, I am eager to learn it.				
Scientific Curiosity	It is fun to break things into pieces to see what is inside.				
Scientific Curiosity	I never hesitate to ask questions.				
Grit	Obstacles or setbacks may discourage me.				
Grit	I often set a goal but later choose to pursue a different one.				
Impulsivity	I tend to say the first thing that comes to mind, without thinking about it.				
Impuisivity	I interrupt people when they are talking.				
Critical Thinking	It's important to understand other people's viewpoint on an issue.				
Citical Tilliking	I usually check the credibility of the source of information before making judgements.				

Table B4: Teacher Survey Inventories

Inventory	Exemplary Items		
4-point likert scale: complete	ly agree, agree, disagree, completely disagree		
	I encourage my students to do research on topics they are interested in and		
	discuss these topics with me. (Inquiry-based Pedagogy)		
	It does not matter if there is noise in the classroom as long as the students are		
	busy with something productive. (Modern Teaching)		
Teaching Styles	Punishment is necessary to create a disciplined class. (Extrinsic Motivation)		
reaching Styles	Teachers should be serious and authoritative in their relationships with students.		
	(Warmth)		
Professional Satisfaction	I am very pleased to have chosen teaching as a profession.		
Competence	It is difficult for me to communicate effectively with students.		
Growth Mindset	Your intelligence is something that you can't change very much.		
Critical Thinking	I sometimes find a good argument that challenges some of my firmly held beliefs.		
Gender Stereotyping	Men have better judgment compared to women; hence they are better leaders.		