



Do Female Role Models Reduce the Gender Gap in Science? Evidence from French High Schools

Thomas Breda, Julien Grenet, Marion Monnet, Clémentine van Effenterre

► To cite this version:

Thomas Breda, Julien Grenet, Marion Monnet, Clémentine van Effenterre. Do Female Role Models Reduce the Gender Gap in Science? Evidence from French High Schools. 2021. halshs-01713068v5

HAL Id: halshs-01713068

<https://shs.hal.science/halshs-01713068v5>

Preprint submitted on 11 Oct 2021 (v5), last revised 21 Mar 2023 (v6)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



WORKING PAPER N° 2018 – 06

**Do Female Role Models Reduce the Gender Gap in Science?
Evidence from French High Schools**

**Thomas Breda
Julien Grenet
Marion Monnet
Clémentine Van Effenterre**

JEL Codes: C93, I24, J16.

Keywords: Role Models; Gender Gap; STEM; Stereotypes; Choice of Study.



Do Female Role Models Reduce the Gender Gap in Science? Evidence from French High Schools*

Thomas Breda
Julien Grenet
Marion Monnet
Clémentine Van Effenterre

First version: January 2018

This version: October 2021

Abstract

We show in a large-scale field experiment that a brief exposure to female role models working in scientific fields affects high school students' perceptions and choice of undergraduate major. While the classroom interventions generally reduce the prevalence of stereotypical views on jobs in science and gender differences in abilities, the effects on educational choices are concentrated among high-achieving girls in Grade 12. They are more likely to enroll in selective and male-dominated STEM programs in college. The most effective role model interventions are those that improved students' perceptions of STEM careers without overemphasizing women's underrepresentation in science.

JEL codes: C93, I24, J16

Keywords: *Role Models; Gender Gap; STEM; Stereotypes; Choice of Study.*

*Breda: CNRS, Paris School of Economics, 48 boulevard Jourdan, 75014, Paris, France (e-mail: thomas.breda@ens.fr); Grenet: CNRS, Paris School of Economics, 48 boulevard Jourdan, 75014, Paris, France (e-mail: julien.grenet@psemail.eu); Monnet: Institut National d'Études Démographiques, 9 cours des Humanités, 93300, Aubervilliers, France (e-mail: marion.monnet@ined.fr); Van Effenterre: University of Toronto, 150 St. George Street, Toronto, Ontario M5S3G7, Canada (e-mail: c.vaneffenterre@utoronto.ca). We are grateful to the staff at the L'Oréal Foundation, especially Diane Baras, Aude Desanges, Margaret Johnston-Clarke, Salima Maloufi-Talhi, David McDonald, and Elisa Simonpietri for their continued support to this project. We also thank the staff at the French Ministry of Education (Ministère de l'Éducation Nationale, Direction de l'Évaluation, de la Prospective et de la Performance) and at the Rectorats of Créteil, Paris, and Versailles for their invaluable assistance in collecting the data. This paper greatly benefited from discussions and helpful comments from Marcella Alsan, Iris Bohnet, Scott Carrell, Clément de Chaisemartin, Bruno Crépon, Esther Duflo, Lena Edlund, Ruth Fortmann, Pauline Givord, Laurent Gobillon, Marc Gurgand, Élise Huillery, Philip Ketz, Sandra McNally, Amanda Pallais, and Liam Wren-Lewis. We thank participants at the ASSA/AEA Annual Meeting in Atlanta, CEPR/IZA Annual Symposium in Labor Economics 2018 in Paris, EALE 2018 in Lyon, EEA-ESEM 2018 in Cologne, Gender Economics Workshop 2018 in Berlin, Gender and Tech Conference at Harvard, IWAEE Conference in Catanzaro, and Journées LAGV 2018 in Aix-Marseille. We also thank seminar participants at Bristol, DEPP, LSE, Harvard Kennedy School, HEC Lausanne, Maastricht SBE, OECD EDU Forum, PSE, Stockholm University SOFI, and Université Paris 8. We are grateful to the Institut des politiques publiques (IPP) for continuous support and to Sophie Cottet for her assistance in contacting schools. Financial support for this study was received from the Fondation L'Oréal, from the Institut des Politiques Publiques, and from the EUR grant ANR-17-EURE-0001. The project received IRB approval at J-PAL Europe and was registered in the AEA RCT Registry with ID AEARCTR-0000903.

Introduction

Women’s increasing participation in science and engineering in the U.S. has leveled off in the past decade (National Science Foundation, 2017). This trend, which is common to almost all OECD countries, is a source of concern for two main reasons. First, it exacerbates gender inequality in the labor market, as Science, Technology, Engineering, and Mathematics (STEM) occupations offer higher average salaries (Brown and Corcoran, 1997; Black et al., 2008; Blau and Kahn, 2017) and show a smaller gender wage gap (Beede et al., 2011). Second, in a context of heightened concern over a shortage of STEM workers in the advanced economies, this trend is likely to represent a worsening loss of talent that could reduce aggregate productivity (Weinberger, 1999; Hoogendoorn et al., 2013).

The underrepresentation of women in these traditionally male-dominated fields can also constitute a self-fulfilling prophecy for subsequent generations, as girls have little opportunity to interact with women working in these fields and who could inspire them. A large literature has established that exposing female students to successful or admirable women can help break this vicious circle. Most of the existing papers focus on potential role models that interact on a regular basis with the individuals they may influence, such as teachers or instructors (Bettinger and Long, 2005; Carrell et al., 2010), university advisors (Canaan and Mouganie, forthcoming), or doctors (Riise et al., forthcoming). Recently, two studies have shown that a one-off exposure to external female role models can also have large effects on female representation in male-dominated fields of study. Porter and Serra (2020) document a positive impact of two female role models who were carefully selected among the economics alumni of Southern Methodist University in the U.S. on the likelihood of female students majoring in economics. Del Carpio and Guadalupe (2018) demonstrate the effectiveness, relative to other types of interventions, of virtual role models to reduce identity costs related to female participation in STEM and to foster female applications to a software-coding program.¹ An attractive feature of these light-touch interventions for identifying role model effects is that they remove the influences of potential confounding factors such as gender differences in teaching practices (Lavy and Sand, 2018; Terrier, 2020; Carlana, 2019)

Although the literature provides compelling evidence that external role models have sizeable effects on educational choices, little is known about how these effects are transmitted. Role models could directly affect students’ preferences. They might change their expectations by modifying their beliefs. By providing an inspirational and relatable model, they could also

¹Related studies outside the context of STEM education include field experiments on exposure to women in leadership positions in India (Beaman et al., 2012) and the provision of information on the returns to education by role models of poor or rich background in Madagascar (Nguyen, 2008).

counteract the effects of gender norms on students’ social identity (Gladstone and Cimpian, 2020). Which of those channels are most affected by external role models? To what extent do changes in students’ perceptions translate into changes in educational choices? Are all role models equally able to influence students’ decision making?

This paper’s contributes to answering these questions by analyzing the effects of a one-hour in-class exposure to external female scientists on female representation in STEM fields of study, and by investigating how the effectiveness of such interventions depends on the characteristics of the role models and the messages they convey. We use a large-scale randomized experiment combined with a comprehensive post-intervention survey to directly measure how role models affect students’ perceptions, beliefs, and enrollment outcomes. Compared to previous studies, a strength of our research design is to involve a large number of role model participants—56 in total. We leverage the diversity of these women’s profiles to better understand what makes an efficient role model. Building on the rapidly expanding literature on the use of machine learning to analyze treatment effect heterogeneity (Athey and Imbens, 2016, 2017; Mullainathan and Spiess, 2017; Wager and Athey, 2018; Chernozhukov et al., 2018), we propose a novel empirical approach that relates the treatment effects on STEM enrollment outcomes to the treatment effects on potential channels. This constitutes a methodological contribution that can be used to investigate mechanisms in randomized controlled trials.

The program we evaluate is called “For Girls in Science” (*Pour les Filles et la Science*) and was launched in 2014 by the L’Oréal Foundation—the corporate foundation of the world’s leading cosmetics manufacturer—to encourage girls to explore STEM career paths. It consists of one-hour in-class interventions by women with two very distinct profiles: half are young scientists (either Ph.D. candidates or postdoctoral researchers) who were awarded the L’Oréal-UNESCO “For Women in Science” Fellowship; the others are young professionals privately employed as scientists in the Research and Innovation division of the L’Oréal group. In the main part of the intervention, the role models share their experience and career path with the students. They also provide information on science-related careers in general and on gender stereotypes using two short videos.

The evaluation was conducted during the 2015/16 academic year in 98 high schools located in the Paris region. It involved 19,451 students from Grade 10 and Grade 12 (science track), two grade levels at the end of which irreversible educational choices are made by students. Half of the classes were randomly selected to be visited by one of the 56 role model participants, who were assigned to those classes through a registration process on a first-come, first-served basis.

We show that the role models’ interventions led to a significant increase in the share of girls enrolling in STEM fields, but only in the educational tracks where they are strongly underrepresented. In Grade 10, the classroom visits had no detectable impact on boys’ and

girls’ probability of enrolling in the science track in Grade 11, where girls are only slightly underrepresented (47 percent of students). By contrast, the intervention induced a significant 2.4-percentage-point increase in STEM undergraduate enrollment among girls in Grade 12, or an increase of 8 percent over the baseline rate of 29 percent, while the effect for boys was negligible. This positive impact on female STEM enrollment is driven by high-achieving female students shifting to selective STEM programs, which lead to the most prestigious graduate schools, and male-dominated STEM programs (math, physics, computer science, and engineering). These results constitute the first field evidence that in-person exposure to external female role models affects STEM enrollment decisions at college entry. They complement the findings from previous research on the effects of external role models on female representation in economics majors (Porter and Serra, 2020), and of female teachers (e.g., Carrell et al., 2010) and virtual role models (Del Carpio and Guadalupe, 2018) in STEM.

To explore the channels through which role models affect students’ enrollment outcomes, we conducted a post-treatment student survey consisting of an eight-page questionnaire administered in class one to six months after the classroom interventions. We also collected administrative data on high school graduation exams (*Baccalauréat*) at the end of Grade 12. Our results show that the role model interventions significantly improved students’ perceptions of science-related jobs at both grade levels, with no indication of declining effects over a period of up to six months. For girls in Grade 12, the interventions significantly increased aspirations for science-related careers. They also helped mitigate some of the stereotypes typically associated with STEM occupations (such as being hard to reconcile with family life) and heightened the perception that these jobs pay better. By contrast, the intervention had no significant effect on students’ self-reported taste for science subjects or their academic performance, and only slightly increased their math self-concept at both grade levels.

One of the most interesting—and least expected—findings concerns the effects on students’ perceptions of gender roles in science. The classroom interventions not only were effective in debiasing students’ beliefs about gender differences in math aptitude, they also raised awareness of the underrepresentation of women in science. The combination of these two effects triggered an unintended ex-post rationalization by students of the gender imbalance in scientific fields and occupations, making them more likely to agree with the statement that women dislike science and that they face discrimination in science-related jobs. Explicitly correcting self-stereotyping beliefs (Coffman, 2014) and misperceptions about women’s representation in science (Bursztyn and Yang, 2021) thus appears to have generated more ambiguous perceptions among students than the intervention’s gender-neutral messages about jobs and careers.

Building on the method proposed by Chernozhukov et al. (2018), we develop a novel approach to relate the student-level treatment effects on enrolment outcomes to the treatment effects on

potential channels. Our results show that the interventions that had the greatest impact on female enrollment in selective STEM programs are those that most improved girls' perceptions of science-related careers without reinforcing the perception that women are underrepresented in science. By contrast, we find that the role models' ability to steer girls towards selective STEM programs is essentially uncorrelated with their effects on students' perceptions of gender differences in aptitude for science.

Overall, our exploration of the different channels provides consistent evidence that the emphasis on gender issues is less important to the effectiveness of such interventions than the ability of role models to project a positive and inclusive image of science-related careers, thus embodying an attractive, attainable path to them.

Finally, we highlight the importance of the role model component of the intervention. First, we argue that the provision of general information on STEM careers cannot explain alone the effects of the interventions on enrollment outcomes. To test the sensitivity of students' attitudes and choices to the intensity of the information component of the treatment, we provided 36 of the 56 role models with a set of slides that contained twice as much informational content as the standard set, including information about wages and employment conditions in STEM jobs. We find no evidence that the treatment effects on students' choice of college major differ significantly between the two sets of slides. Second, we document a high degree of heterogeneity in treatment effects according to the role models' professional background. Those employed by the sponsoring firm had a significantly greater effect on girls' probability of enrolling in selective STEM programs than the young researchers, despite being exposed to students with similar observable characteristics. Although the two groups were equally effective in debunking the stereotype on gender differences in math aptitude, we find clear evidence that those with a professional background were better able to improve girls' perceptions of science-related jobs and raise their aspirations for such careers. Conversely, they were less likely to reinforce students' belief that women are underrepresented in science. Together, these results show that role model interventions are not reducible to the provision of standardized information and that female role models are not interchangeable. They confirm, in a real-life setting, results from lab experiments in social psychology highlighting the importance of role models' profiles (Lockwood and Kunda, 1997; Cheryan et al., 2011; Betz and Sekaquaptewa, 2012; O'Brien et al., 2016).

The remainder of the paper is organized as follows. Section 1 provides institutional background on the French educational system and the gender gap in STEM fields. Section 2 describes the intervention and the experimental design. Section 3 presents the data and empirical strategy. Section 4 analyzes the effects of role model interventions on student perceptions, self-concept, and educational outcomes. Section 5 extends the analysis to the role of information, the persistence of effects, and potential spillovers. Section 6 discusses mechanisms and Section 7 concludes.

1 Institutional Background

1.1 Structure of the French Education System

In France, education is compulsory from the age of 6 to the age of 16, with the academic year running from September to June. The school system consists of five years of elementary education (Grades 1 to 5) and seven years of secondary education, divided into four years of middle school (*collège*, Grades 6 to 9) and three of high school (*lycée*, Grades 10 to 12). Students complete high school with the *Baccalauréat* national exam, which they must pass for admission to higher education.

High school tracks. The tracking of students occurs at two critical stages (see Figure 1). At the end of middle school, about two-thirds of students are admitted to general and technological upper secondary education (*Seconde générale et technologique*) and the remaining third are tracked into vocational schools (*Seconde professionnelle*). After the first year of high school (Grade 10), the general and technological track is further split: approximately 80 percent of the students are directed to the general *Baccalauréat* program for the last two years of high school (Grades 11 and 12), while the remaining 20 percent, who are mostly low-achieving students, are directed towards a technological *Baccalauréat*, which is more geared towards the needs of business and industry and leads to shorter studies.

In the Spring term of Grade 10, the students who have been allowed to pursue the general track are required to choose among three sub-tracks in Grade 11: Science (*Première S*), Humanities (*Première L*), and Social sciences (*Première ES*). This is an important choice, given that the curriculum and high school examinations are specific to each *Baccalauréat* track and thus directly impact students' educational opportunities and career prospects. It is almost impossible, for instance, for a student to be admitted to engineering or medical undergraduate programs without a *Baccalauréat* in science. Students directed to the technological track after Grade 10 are also required to choose among eight possible STEM and non-STEM sub-tracks, which will affect their choice of field of study in higher education.

College entry. In the Spring term of Grade 12, students in their final year of high school apply for admission to higher education programs through a centralized online admission platform. The programs to which students can apply fall into two broad categories, each accounting for about half of first-year undergraduate enrollment: (i) non-selective undergraduate university programs (*Licence*), which are open to all students who hold the *Baccalauréat*; and (ii) selective programs, which can select students based on their academic achievement. Both types of programs offer

specializations in STEM and non-STEM fields. Among selective programs, the most prestigious are the two-year *Classes préparatoires aux Grandes Écoles* (CPGE), which prepare students to take the national entry exams to elite graduate schools (*Grandes Écoles*). These programs are specialized either in science, in economics and business or in humanities. Within the science CPGE programs, the main fields of specialization are mathematics and physics (MPSI), physics and chemistry (PCSI), and biology/geoscience (BCPST). The other selective undergraduate programs (*Section de technicien supérieur* or STS) are mostly targeted to students holding a vocational or technological *Baccalauréat* and prepare for technical/vocational bachelor's degrees.

1.2 Female Underrepresentation in STEM

In France, the share of female students in STEM-oriented studies starts to decline after Grade 10 and drops sharply at entry into higher education. While 54 percent of the students in the general and technological track in Grade 10 are girls, the share falls to 47 percent in the general science track (Grades 11 and 12) and further to 30 percent in the first year of higher education.² Female underrepresentation in STEM fields of study is more pronounced in the selective undergraduate programs (shares of 18 percent in STS and 30 percent in CPGE) than in the non-selective programs (35 percent). These proportions, which are computed from administrative data for 2016/17, are almost identical to those of a decade earlier. Within STEM fields of study, female students tend to specialize in earth and life sciences (female share: 62 percent) rather than mathematics, physics, or computer science (female share: 26 percent).

The underrepresentation of women in STEM fields accounts for a good part of the gender pay gap among college graduates in France. Using a variety of administrative and survey data sources, we show in Appendix A that across all majors, male graduates who obtained a master's degree in 2015 or 2016 earn a median gross annual starting salary of 32,122 euros, compared to 28,411 euros for female graduates. This amounts to an overall gap of 3,711 euros per year, or 11.6 percent of men's pay (see Table A1). Using standard decomposition methods, we find that the underrepresentation of female students in STEM accounts for approximately 25 percent of this gap (see Table A2). Additionally, almost half of the 9.1 percent gender pay gap within STEM can be ascribed to the fact that female graduates are less likely than males to be enrolled in the selective and male-dominated fields, which lead to the best-paying degrees. These figures strongly suggest that in the French context, increasing the share of female students in STEM—especially in selective and male-dominated programs—would narrow the gender pay gap substantially.

²At the high school level, the gender imbalance in STEM is more severe in the technological track (female share: 17 percent) than in the general science track (female share: 47 percent).

2 Program and Experimental Design

2.1 The “For Girls in Science” Program

The “For Girls in Science” (FGiS) program is an awareness campaign launched in 2014 by the L’Oréal Foundation to encourage girls to explore STEM career paths. It consists of one-hour one-off classroom interventions by female role models with a background in science. The interventions, which take place in the presence of all students in the class, including boys, are carried out by female role models of two distinct types: (i) Ph.D. candidates or postdoctoral researchers who have been awarded a fellowship by the L’Oréal Foundation (the L’Oréal-UNESCO “For Women in Science” Fellowship) and who participate in the program as part of their contract; and (ii) young professionals employed as scientists in the Research and Innovation division of the L’Oréal group who volunteer for the program.

Structure and content of the interventions. The classroom interventions last one hour and are divided into four main sequences. Each role model was provided with a set of slides as a support for the entire in-class conversation. During the first sequence, a small number of slides highlight two facts: (1) the labor market is marked by high demand for STEM skills and there is a shortage of graduates in the relevant fields of study; and (2) women are underrepresented in STEM careers. To investigate the role of information provision, we provided 36 of the 56 role models with additional slides that they were free to use during this sequence. These slides contained supplementary information about average earnings and employment conditions in STEM jobs, and were illustrated with examples of career prospects in humanities versus science. In Section 5, we discuss the sensitivity of our results to this more intensive provision of standardized information.

The second sequence kicks off with two three-minute videos designed to illustrate and deconstruct stereotypes about science-related careers and gender roles in science.³ The first video, entitled “Science, Beliefs or Reality?,” uses interviews with high school students to debunk myths about careers in science (e.g., jobs in science are more challenging, they necessarily require long studies), stereotypes about scientists (e.g., they are introverted, lonely), and gender differences in science aptitude (e.g., women are naturally less talented in math). The second video, entitled “Are we all Equal in Science?,” describes the common gender stereotypes about aptitude for science while providing information on brain plasticity and on how interactions and the social environment shape men’s and women’s abilities and tastes. This sequence aims at stimulating class discussion based on students’ reactions to the videos.

³Screenshots of the two videos shown during the classroom interventions are displayed in Appendix Figure B1.

The third sequence centers on the female role model’s own experience as a woman with a background in science and consists of an interactive question-and-answer session with the students.⁴ Topics addressed during this discussion include the role model’s typical day at work, what she enjoys about her job, the biggest challenge she had to overcome, how she views her professional future, her everyday interactions with co-workers, how much she earns, and her work-family balance. Consistent with the program’s emphasis on the role model dimension, this sequence was intended to be the longest and most important part of the intervention. In order to convey this objective to the role models, a full-day training was organized to help them share their experience with the students. The training also included a workshop on the underrepresentation of women in science and a practice session aimed at enhancing oral communication skills.

The intervention concludes with an overview of the diversity of STEM studies and careers, illustrated by concrete examples such as jobs in graphic design, environmental engineering, and computer science.

2.2 Experimental Design

Selection of schools and classes. The evaluation was conducted in the three education districts (*académies*) of the Paris region (Paris, Créteil, and Versailles) during the 2015/16 academic year. Créteil and Versailles are the two largest education districts in France and the three districts combined include 318,000 high school students in the general and technological track, or 20 percent of all French high school enrollment.

Figure 2 provides a detailed timeline of the evaluation. In the spring of 2015, the French Ministry for Education agreed to support a randomized evaluation of the program and designated one representative for each district as intermediary between the schools and the evaluation team. In June, official letters informed high school principals that they were likely to be contacted to take part in the evaluation. All public and private high schools with at least four classes in Grade 10 and two in Grade 12 (science track) were contacted by our team between September and December 2015, accounting for 349 of the 489 high schools operating in the three districts. Of these schools, 98 agreed to take part in the experiment, representing 28 percent of Grade 10 enrollment and 29 percent of Grade 12 (science track) enrollment in the three districts combined.⁵ The participating schools tend to be larger and are less likely to be private or to operate in the Paris education district than the non-participating ones (see Appendix Table E1). Despite these differences, the experimental sample, which consists of 19,451 students (13,700 in Grade 10 and 5,751 in Grade 12), closely resembles the relevant student population, both in

⁴Screenshots of the slides used during the discussion are displayed in Appendix Figure B2.

⁵The location of the participating schools is shown in Appendix Figure B3.

social composition and in average academic performance (see Appendix Table E2).

Randomization. In the fall of the 2015/16 school year, the principals were invited to select at least six classes—four or more in Grade 10 and two or more in Grade 12 (science track)—and to indicate a preferred time slot and day for the interventions.⁶ In each school, half of the classes selected by the principal (up to the nearest integer) were randomly assigned to the treatment group (302 classes in total) and the other half to the control group (299 classes). Table 1 indicates that the random assignment successfully balanced the characteristics of students in the treatment and control groups.

Role models. The experiment involved 56 female role models, of whom 35 were L’Oréal employees and 21 were Ph.D. candidates or postdoctoral researchers. Table 2 provides summary statistics of their characteristics. The researchers tend to be younger (30 vs. 36 years of age on average) and are less often of foreign nationality (10 vs. 17 percent). Although both types have very high levels of educational attainment, 39 percent having graduated from a *Grande École*, the researchers are more likely than the professionals to hold (or prepare for) a Ph.D. (100 vs. 38 percent) and to hold a degree in math, physics and engineering (38 vs. 14 percent). They are also less likely to have children (19 vs. 58 percent) and to have been involved in the program in the previous year (19 vs. 29 percent). The professionals working at L’Oréal are employed in various activities: chemistry (development of new technologies for skin products), logistics and supply chain management, statistics (consumer evaluation), immunology and toxicology. Although we could not collect direct information on earnings for reasons of confidentiality, we estimate based on aggregate information provided by the L’Oréal Group that the annual gross wages of these young professionals is between 45,000 and 65,000 euros, compared to between 22,000 and 50,000 euros for the researchers. On average, each role model carried out five classroom interventions in two different high schools.

Classroom interventions. The classroom visits took place between November 17, 2015, and March 3, 2016.⁷ The role models were asked to select two or three schools in which to carry out an average of three classroom visits per school—in most cases, two in Grade 10 and one in Grade 12. They were not assigned to the schools randomly but registered for the visits and time slots during four registration sessions using an online system on a first-come, first-served basis. Randomly assigning the role models to the schools was not a feasible option, since most were participating on a voluntary basis and during regular working hours. We therefore identify the

⁶In the vast majority of schools, principals selected exactly four Grade 10 and two Grade 12 classes.

⁷17 percent of the visits took place in November, 26 percent in December, 39 percent in January, 17 percent in February, and 1 percent in March.

causal impact of role models in a setting where they have some freedom to choose the schools in which they intervene. The assignment process, however, did not involve any coordination between the participants and was designed to limit their ability to select the schools they would visit, as each registration session only concerned a subset of the participating schools.⁸

3 Data and Empirical Strategy

3.1 Data

To evaluate the program’s effects on student perceptions and educational choices, we combine three main data sources: (i) a post-intervention survey of role models; (ii) a post-intervention survey of students; and (iii) student-level administrative data.⁹

Role model survey. After each visit to a school, the role models were invited to complete an online survey. Besides collecting general feedback, this survey served to monitor compliance with random assignment, asking them to indicate each of the classes they visited. Summary statistics are reported in Appendix Table E3. The interventions almost always (89 percent) took place in the presence of the teacher and sometimes (35 percent) of another adult. The role models reported organizational problems for only 16 percent of the visits (e.g. the intervention started late, the slides could not be shown). According to the survey, researchers and professionals were equally likely to cover the intended topics, such as “jobs in science are fulfilling”, “they are for girls too”, and “they pay well”. Finally, when asked about their overall perception of their classroom interventions, 93 percent said they went “well” (37 percent) or “very well” (56 percent). Students were generally perceived to be responsive to the key messages.

Student survey. We conducted a paper-and-pencil student survey in classes assigned to the treatment and control groups one to six months after the classroom visits, between January and May 2016. Each questionnaire was assigned a unique identifier so that it could be linked with student-level administrative data. The survey was designed to collect a rich set of information on students’ preferences, beliefs and perceptions regarding science, self-concept and aspirations, and was administered in exam conditions under the supervision of a teacher. It was presented as a general survey on students’ attitudes about science and science-related careers so as to minimize the risk that students would associate it with the FGiS program. It was eight pages long and took about half an hour to complete.

⁸The role models were contacted four times to complete the schedule, on October 21, November 24, December 7, 2015, and February 3, 2016.

⁹Translated versions of the two surveys are available online at <https://mycore.core-cloud.net/index.php/s/L0aB9Kvpbot7sNh>.

The survey items investigate the effects of classroom interventions on students’ perceptions along five dimensions: (i) general perceptions of science-related careers; (ii) perceptions of gender roles in science; (iii) taste for science subjects; (iv) math self-concept; and (v) science-related career aspirations. When conceptually related, we combine the survey items to construct a synthetic index for each dimension using standardized z -score scales. Section 4 describes the specific items that are used for each dimension of interest.¹⁰

As shown in Appendix Table E5, the survey response rates are high both in Grade 10 (88 percent of students) and in Grade 12 (91 percent). They are slightly higher among Grade 10 students in the treatment than in the control group (by 2.6 percentage points). Despite this small difference in response rates, the characteristics of survey respondents in Grade 10 are generally balanced between the treatment and control groups (see Appendix Table E6). The opposite is found in Grade 12: the survey response rates are similar in the two groups, but the respondents’ characteristics exhibit some small but statistically significant differences. In Section 4, we show that the survey-based results are robust to controlling for these small imbalances.

Administrative data. We linked the student survey data to a rich set of individual-level administrative data covering the universe of students enrolled in the high schools of the Paris region over the period 2012/13 to 2016/17. These data provide detailed information on students’ socio-demographic characteristics and enrollment status every year, allowing us to identify the high school track taken by Grade 10 students entering Grade 11.

The college enrollment outcomes of students in Grade 12 were obtained by matching the survey and administrative data for high school students with administrative microdata covering almost all students enrolled in selective and non-selective higher education programs in 2016/17.¹¹ These data are complemented with comprehensive individual examination results from the *Diplôme National du Brevet* (DNB), which is taken at the end of middle school, and from the national *Baccalauréat* exam (for Grade 12 students). Specifically, we use students’ grades on the final exams in French and math (converted into national percentile ranks), as these tests are graded externally and anonymously. Further details about the data sources and the classification of higher education programs can be found in Appendix C.

¹⁰To mitigate potential order bias, the order of several of the response items (e.g., math/French, man/woman) was set randomly.

¹¹Programs not covered by these administrative data are those leading to paramedical and social care qualifications. Available estimates suggest that among Grade 12 students who obtained a *Baccalauréat* in Science in 2008, under 6 percent were enrolled in such programs the following year (Lemaire, 2018).

3.2 Empirical Strategy

Compliance with random assignment was not perfect: about 5 percent of the classes assigned to the treatment group were not visited by a role model, while 1 percent of the classes in the control group were mistakenly visited (see Appendix Table E4).¹² To deal with this marginal two-way non-compliance, we follow the standard practice of using treatment assignment as an instrument for treatment receipt, which allows us to estimate the program’s local average treatment effect (LATE) instead of the average treatment effect (ATE). Specifically, we estimate the following model using two-stage least squares (2SLS):

$$Y_{ics} = \alpha + \beta D_{cs} + \theta_s + \epsilon_{ics}, \quad (1)$$

$$D_{ics} = \gamma + \delta T_{cs} + \lambda_s + \eta_{ics}, \quad (2)$$

where Y_{ics} denotes the outcome of student i in class c and high school s , D_{cs} is a dummy variable indicating whether the student’s class received a visit, and T_{cs} is a dummy for assignment to the treatment group. School fixed effects, θ_s and λ_s , are included to account for the fact that the randomization was stratified by school and grade level.

The model described by Equations (1) and (2) is estimated separately by grade level and gender, with standard errors clustered at the unit of randomization level (class). To account for multiple hypotheses testing across the outcomes of interest, the treatment effect estimates are accompanied by adjusted p -values (q -values) in addition to the standard p -values.¹³

4 Effects of Classroom Interventions

We analyze the impact of the classroom interventions on three main sets of student outcomes: (i) general perceptions of science-related careers and of gender roles in science; (ii) preferences, aspirations and self-concept; and (iii) enrollment outcomes and academic performance.

4.1 Perceptions of STEM Careers and Gender Roles in Science

Students’ post-intervention survey responses show that the classroom interventions were effective in challenging stereotyped views of science-related careers and gender roles.

¹²We are confident that non-compliance was mostly due to organizational and logistical issues and was not an endogenous response to randomization. The few role models who carried out interventions in classes assigned to the control group or in classes not selected to participate in the evaluation generally reported that their interventions had been poorly organized at the school level, with the person in charge often not being aware of the purpose of the visit. In some cases, classroom interventions were scheduled during another specialty course involving multiple classes, meaning that only some of the students in the treatment group were effectively treated.

¹³We use the False Discovery Rate (FDR) control, which designates the expected proportion of all rejections that are type-I errors. Specifically, we use the sharpened two-stage q -values introduced in Benjamini et al. (2006) and described in Anderson (2008). See Appendix D for details.

Perceptions of science-related careers. Students were asked to agree or disagree with five statements on science-related careers relating to pay, the length of studies leading to these careers, work-life balance, and the two prevalent stereotypes that science-related jobs are monotonous and solitary. We build a composite index of “positive perceptions of science-related careers” by re-coding the Likert scales so that higher values correspond to less stereotypical or negative perceptions, before taking the average of each student’s responses to the five questions. To facilitate interpretation, we normalize the index to have a mean of zero and a standard deviation of one in the control group.¹⁴ For closer investigation of the various aspects that might be captured by the overall index, we further construct binary variables taking value one if the student agrees strongly or somewhat with each statement, and zero if he/she disagrees strongly or somewhat.¹⁵

One of the interventions’ objectives was to correct students’ beliefs about jobs and careers in science through the provision of information specific to each role model’s experience as well as standardized information. As shown in Table 3, students’ baseline perceptions indicate relatively widespread negative stereotypes about careers in science (see columns 1 and 4), with little difference between boys and girls or between grade levels. As an example, between 20 and 30 percent of students consider that science-related jobs are monotonous or solitary. The role model interventions significantly improved girls’ and boys’ perceptions of such careers as measured by the composite index, in both Grade 10 and Grade 12. The effects range from 15 percent of a standard deviation for boys to around 30 percent for girls, with significantly greater effects for female students in both grades. A significant impact of the classroom visits is observed for almost all the components of the index. The largest effects are found for the statements “science-related jobs require long years of study” and “science-related jobs are rather solitary,” which embody two stereotypes that were specifically debunked in the slides and videos. Although the effects are not strikingly different between genders and grade levels, they tend to be greater for girls in Grade 12. In particular, the interventions appear to have closed the gender gap in Grade 12 students’ awareness of the earnings premium attached to science-related jobs: while girls in the control group are less likely than boys to agree with the statement that “jobs in science pay well” (53 vs. 58 percent), these proportions are comparable in the treatment group (around 60 percent). Additionally, the interventions have reinforced girls’ perception that science-related careers are compatible with a fulfilling family life, a message specifically conveyed by the role models and in line with the evidence showing that jobs in science and

¹⁴We have checked that our results are robust to converting the item responses into binary variables before computing the indices. See Appendix D for further details on the construction of the synthetic indices.

¹⁵Similar groupings are performed when using responses that are measured on a four-point Likert scale (usually concerning perceptions or self-confidence) so that the outcome variables can be directly interpreted as proportions. The results are not qualitatively affected by such grouping.

technology enable women to work more flexibly (Goldin, 2014).

Perceptions of gender roles in science. Female underrepresentation in STEM can be broadly attributed to three possible causes: gender differences in abilities, discrimination (on the demand side), and differences in preferences and career choices (on the supply side). The survey questions were designed to capture students’ views on these dimensions.

Table 4 reveals the striking fact that more than a third of Grade 10 students and a quarter of Grade 12 students in the control group are not aware that women are underrepresented in science-related careers. These proportions are similar by gender and by grade. For boys and girls in both grades, we find that the interventions increased awareness of female underrepresentation in STEM by 12 to 17 percentage points. This is one of the outcomes most strongly affected by the interventions.

The classroom interventions were also effective in debiasing students’ beliefs about gender differences in math aptitude. To capture this dimension, we asked students whether they agreed with the statements that “men are more gifted than women in mathematics” and that “men and women are born with different brains.” We used these two questions to construct a composite index to gauge whether students believe that men and women have equal aptitude for mathematics. The results show significant rises in this index for both genders in both grades, with treatment effects ranging between 9.5 percent and 14.8 percent of a standard deviation.¹⁶

Interestingly, the classroom visits had more ambiguous, partially unintended effects regarding the other two explanations. First, when asked about gender differences in preferences, the share of students who agree with the statement that “women don’t really like science” is relatively low in the control group (16 percent of girls and 20 percent of boys in Grade 10; 7 percent of girls and 15 percent of boys in Grade 12), but it increases substantially due to the interventions for both genders, by 4 to 10 percentage points. Second, the baseline shares of boys and girls who declare that women face discrimination in science-related jobs are much larger (around 60 percent); these too increase for both genders, by 7 to 15 percentage points. These unintended effects on students’ perceptions of gender roles in science could have arisen as an effort to rationalize why there are so few women in science-related careers, making students more likely to agree with the simplistic view that “women don’t really like science” and to subscribe to the idea that women face discrimination in science careers.

¹⁶The detailed results for the two components of this index are reported in Appendix Table F1.

4.2 Stated Preferences and Self-Concept

We now turn to the effects of the interventions on students' stated preferences and self-perception. Specifically, we investigate whether the interventions affected boys' and girls' taste for science subjects, their self-concept in math, and their science-related career aspirations.

Taste for science subjects. For both genders in Grade 10 and Grade 12, the interventions had no sizeable impact on students' enjoyment of science subjects at school (reported on a 0 to 10 Likert scale), i.e., math, physics-chemistry, and earth and life sciences, or on their self-reported taste for science in general (see Table 4 for the composite index aggregating the four relevant questionnaire items and Table F2 in the Appendix for the detailed results). These findings are not particularly surprising, given that the interventions did not expose students to science-related content and were not specifically designed to promote interest in science.

Math self-concept. To measure the impact of the classroom visits on students' self-concept in mathematics, we use a composite index that combines students' responses to four questions: (i) their self-assessed performance in math; (ii) whether they feel lost when trying to solve a math problem; (iii) whether they often worry that they will struggle in math class; and (iv) whether they consider that they can do well in science subjects if they make enough effort.

Consistent with the literature, our sample exhibits large gender differences in self-concept in mathematics. In the control group, the value of the index is 43 percent of a standard deviation lower for girls than for boys in Grade 10, and 37 percent lower in Grade 12. Large gender differences are found for most of the items used in the construction of this index, in particular those related to math anxiety (see Appendix Table F3).

Despite being a light-touch intervention, the interventions did have some positive effect on students' self-concept in math. Although these effects are only found to be statistically significant for boys in Grade 12 when using the composite index, the interventions consistently reduced the probability of students reporting worry that they will struggle in math class.¹⁷ Point estimates tend to be higher for boys than for girls in both grades, implying that the classroom interventions had no correcting effect on the substantial gender gap in this area.

Science-related career aspirations. The choice of a science-related career path does not depend solely on students' taste for the science subjects taught at school. It also depends on

¹⁷For each group of students, the correction of p -values for testing across multiple outcomes (see Appendix Table F3) cannot rule out the possibility that the effects on math anxiety are due to chance alone. However, finding a significant effect for the same variable across all four groups of students, which is not accounted for by the multiple testing correction, is suggestive of a genuine effect.

their perceptions of the relevant jobs and the amenities they may provide, such as earnings, work/life balance, and the work environment, all of which were embodied by the role models.

To measure the effects on students' aspirations for science-related careers, we use a composite index combining the responses to four questions: (i) whether the students find that some jobs in science are interesting; (ii) whether they could see themselves working in a science-related job later in life; (iii) whether they report being interested in at least one of six STEM jobs out of a list of ten STEM and non-STEM occupations;¹⁸ and (iv) whether they consider career and earnings prospects as important factors in their choice of study.

Female students in Grade 12 are the only group of students for which we find significant effects on these science-related career aspirations, the value of the composite index being 11 percent of a standard deviation higher in the treatment than in the control group (see the last row of Table 4). The more detailed results reported in Appendix Table F4 show that the interventions had significant positive effects on three of the four corresponding survey items for girls in Grade 12. In particular, girls in the treatment group are more likely to report that career and earnings prospects are important factors in their choice of study, which is consistent with the interventions raising their awareness of the wage premium for STEM jobs.

4.3 Educational Choices and Academic Performance

High school track after Grade 10. Panel A of Table 5 shows that the classroom visits had no significant impact on Grade 10 students' choice of track in the academic year following the intervention, i.e., 2016/17. For both genders, the treatment effect estimates are close to zero, whether we consider enrollment in any STEM track or enrollment in the general and technological STEM tracks separately.¹⁹ Consequently, the interventions did not alter the 20-percentage-point gender gap in the likelihood of pursuing STEM studies after Grade 10.

Several mechanisms can be put forward to interpret the lack of effects on the enrollment status of Grade 10 girls in the following year. First, the interventions did not seem well suited to increase the share of girls enrolling in the STEM technological tracks in Grade 10, where the female share is particularly low (17 percent). As discussed below, the positive effects that we find on the STEM enrollment decisions of girls in Grade 12 are concentrated among the high achievers in math. In Grade 10, such students are unlikely to be directed to the technological track, which could explain the lack of effects along this margin. Turning to the general science track, female underrepresentation is only moderate in Grade 11 (in 2016/17, the female share was

¹⁸The STEM occupations in the list were: chemist, computer scientist, engineer, industrial designer, renewable energy technician, and researcher in biology. The non-STEM occupations were lawyer, pharmacist, physician, and psychologist.

¹⁹The more detailed results presented in Appendix Table F5 show that the distribution of students across non-STEM tracks (Humanities and Social sciences) did not change significantly either.

47 percent) and this track is the most common (usually the default choice) for high-performing students, including girls. Unlike the other high school tracks, it gives access to almost all fields of study in higher education and hence does not signal a strong commitment to pursue a STEM education or career in the future, limiting the potential of STEM role models to influence enrollment in this track. Female students who turn away from the science track in high school are unlikely to even consider a STEM career as a viable option, making their choices less easily reversible.²⁰

Field of study after Grade 12. A central finding of the study is that the role model interventions had significant effects on the educational choices of girls in Grade 12, but not on those of their male classmates.

Panel B of Table 5 shows that for girls in Grade 12, the interventions increased the probability of enrolling in a STEM undergraduate program in 2016/17 by 2.4 percentage points (significant at the 10 percent level), which corresponds to an 8.3 percent increase from the baseline of 28.9 percent. The effect for boys is negligible and not statistically significant, implying that the gender gap in STEM enrollment narrowed from a baseline of 18.1 to 16.0 percentage points, i.e., an 11.6 percent reduction.²¹

As emphasized in Section 1.2, female underrepresentation in selective and male-dominated STEM fields account for approximately half of the STEM-related gender pay gap in France. Importantly, our results show that the interventions’ positive impact on STEM enrollment is driven by a significantly larger fraction of girls in Grade 12 enrolling in both types of programs. The classroom interventions led to a highly significant 3.5 percentage-point increase in the fraction of girls enrolling in selective STEM programs, which represents a 32 percent increase from the baseline of 11.0 percent. The corresponding estimates for boys suggest that the classroom visits may have slightly increased male enrollment in these programs as well (by 2.0 percentage points from a baseline of 23.2 percent), but the effect is not statistically significant. Moreover, we show in Section 4.4 that the magnitude of this effect for boys is substantially lessened when we control for students’ baseline characteristics, suggesting that it probably depends on small residual imbalances in the male sample.²²

Turning to the effects on enrollment in male-dominated STEM programs (mathematics,

²⁰Consistent with this interpretation, the survey data indicate that among Grade 10 students in the control group, only 24 percent of girls who did not enroll in the Grade 11 science track the following year declare that they could see themselves working in a science-related job, compared to 87 percent among those who did.

²¹With the caveat that we lack the statistical power to detect a significant reduction in the gender gap in STEM enrollment.

²²Balancing tests performed separately by grade level and gender do not point to unusually large imbalances between the treatment and control groups in any of the subsamples (results available upon request). However, the predicted probability of being enrolled in a selective STEM program is marginally higher in the treatment group than in the control group for boys in Grade 12 (by 0.8 percentage point from a baseline of 23.8 percent, significant at the 5 percent level).

physics, computer science, and engineering), we find that the proportion of girls enrolling in such programs increased by a statistically significant 3.8 percentage points from a baseline of 16.6 percent (i.e., a 23 percent increase), compared to a non-significant 1.7-point increase for boys from a baseline of 37.9 percent. These results are particularly striking given that selective and male-dominated STEM programs are not only the most prestigious tracks but also those where the gender gap in enrollment is greatest. A simple back-of-the-envelope computation suggests that if our estimates could be extrapolated to the population of science-track Grade 12 students without considering general equilibrium effects, the female share would increase from 30 to 32 percent in STEM programs altogether, from 30 to 34 percent in selective STEM programs, and from 26 to 29 percent in male-dominated STEM programs.

Our estimates indicate that, on average, the role model interventions induced one girl in every two Grade 12 science-track classes to switch to a selective or a male-dominated STEM program at entry into higher education.²³ The more detailed results presented in Appendix Table F6 indicate that these effects are driven by female students shifting from non-STEM and female-dominated STEM programs. A significant decline in female enrollment is found for non-selective undergraduate programs in earth and life sciences (-2.2 percentage points), while small reductions in the range of 0.4 to 0.8 point are found for selective programs in humanities and vocational non-STEM programs, as well as for non-selective programs in medicine, law and economics, humanities and psychology, and sports.

Taken together, the results for Grade 12 students show that the interventions were only effective in steering girls towards the STEM tracks in which they are heavily underrepresented, even though two-thirds of the role models come from female-dominated STEM fields (earth and life sciences) and that the interventions were designed to promote all types of STEM careers, including those where women now outnumber men. These findings suggest that in the current setting, the role models affect only the most stereotyped choices.

Academic performance. The effects of the classroom visits on academic performance can be documented for students in Grade 12 based on the *Baccalauréat* exams, taken a few months after the classroom interventions (see Appendix Table F7). The treatment effect estimates on students' performance on the math test and on the probability of obtaining the *Baccalauréat* are close to zero and statistically insignificant for both genders. Although the role models could, in principle, have strengthened students' motivation to be admitted to the most selective STEM programs, resulting in their dedicating more time to studying mathematics and other science subjects, we find no evidence of any such effect. We can therefore rule out that the interventions'

²³This computation is based on an average of 15 girls per class and an estimated 3.5 (respectively 3.8) percentage-point increase in the probability of enrolling in a selective (respectively male-dominated) STEM program.

impact on the enrollment outcomes of girls in Grade 12 was driven by increased effort and accordingly better academic performance.

4.4 Robustness Checks

We conducted a number of robustness checks for our main findings, which are reported in Appendices G and H.

First, we investigated whether our estimates for the survey-based outcomes might not be contaminated by the small imbalances in the response rates and observable characteristics of the treatment and control groups (see Section 3). We show that the estimated effects on students' perceptions are barely affected when controlling for students' observable characteristics (Table G1).

Second, controlling for students' observable characteristics hardly affects the estimated effects on enrollment outcomes (see Table G2). If anything, the small positive (but not significant) effect on selective STEM enrollment for boys in Grade 12 becomes negligible.

Third, we checked whether our results are robust to using non-parametric randomization inference tests rather than model-based cluster-robust inference. The tests are performed by comparing our ITT estimates with the distribution of "placebo" ITT estimates obtained by randomly re-assigning treatment 2,000 times among participating classes within each school and grade level. The results yield empirical p -values that are generally close to the model-based p -values (see Table H1). Although they tend to be slightly more conservative, they confirm the interventions' statistically significant effects on female enrollment in selective and male-dominated STEM programs among Grade 12 students.

5 Information, Persistence, and Spillovers

In this section, we test the sensitivity of students' attitudes and choices to the intensity of information provision. We then extend the analysis to the persistence of effects on student perceptions, the timing of the interventions, and investigate potential spillover effects on enrollment outcomes.

Intensity of information provision. Any role model intervention intrinsically contains an informational component on top of fostering self-identification. While our design does not allow to fully disentangle these two mechanisms, we are able to test the sensitivity of students' attitudes and choices to the intensity of the standardized information contained in the slides that were provided to the role model participants.

As described in Section 2, we initially sent a set of slides to the role models to assist them during the in-class intervention. The first slides (six in total) highlighted a few stylized facts about jobs in science and female underrepresentation in STEM careers, providing only limited information on employment conditions in such careers, and no information on wages. Starting on November 20, 2015, we sent six additional slides to 36 of the 56 role models. These new slides presented more detailed information regarding wage and employment gaps between STEM and non-STEM jobs, as well as differences between male and female students' choice of studies. The role models were free to integrate these slides into their final presentation or to only use them as a support.²⁴

The results reported in Appendix I.1 show that students' characteristics are balanced according to whether the role model received the regular set of slides or the "augmented" version (Table I1).²⁵ While the more information-intensive treatment had a larger impact on the probability that female students agree with the statement that science-related jobs pay higher wages, the effects on the probability of enrolling in selective STEM or male-dominated STEM programs are not significantly different (Table I2). These findings provide suggestive evidence that the purely informational component of the intervention does not in itself explain the observed changes in college major decisions.

Persistence. The effects of the interventions on students' perceptions could be short-lived. We explore this issue by comparing the magnitude of treatment effects for different intervals between the intervention and the post-treatment survey: 1-2 months, 3-4 months and 5-6 months (see Appendix Table I3). The limited sample for each interval and the possibility that the quality of the interventions may have changed over time are two limitations that call for caution in drawing firm conclusions about the persistence of effects. With these caveats in mind, the results suggest that the treatment effects did not vanish quickly, insofar as they remain statistically significant for most outcomes beyond the first two months. The effects were, therefore, sufficiently persistent to affect students' choice of study.

Timing of visits. Earlier interventions seem to have had greater effects on the college choices of Grade 12 students, which could be made up to the end of May (see Appendix Figure I3). We find that classroom visits that took place in November increased female enrollment in selective or male-dominated STEM programs by 7 to 9 percentage points, compared with 3 to 6 points

²⁴Screenshots of the two sets of slides are shown in Appendix Figures I1 and I2.

²⁵We initially planned to randomly allocate the two sets of slides to the role models and were able to do so for a subset of 14 participants. However, the L'Oréal Foundation requested that going forward, all remaining role models were be provided with the "augmented" version of the slides. The role models who had already started the visits kept the regular version. To ensure sufficient statistical power, we present results for the entire sample of role models, controlling for month-of-visit fixed effects. Our results are qualitatively similar when we restrict our sample to the subset of role models for whom the slides were randomly assigned.

for visits in December-January and non-significant effects for visits in February-March.²⁶ These findings provide suggestive evidence that interventions made when many students are still undecided about their field of study and career plans may be more effective than those on the eve of irreversible choices.

Spillovers. An important issue is whether the interventions could have influenced the educational choices of students in the control group. These students may have heard about the visits directly, through their schoolmates in treatment group classes, or indirectly, through regular social interactions. If the direction of such effects is the same for students in the treatment and control groups, ignoring spillovers would cause us to underestimate the treatment effects.

On the last page of the post-intervention survey questionnaire, the students in the treatment group were asked whether they had discussed the classroom intervention with their classmates, with schoolmates from other classes, or with friends outside of school, as a way of assessing possible spillover effects. Students in the control group received a slightly different version of this final section, asking whether they had heard of classroom visits by male or female scientists in other classes, with no explicit mention of the FGiS program.

The survey evidence suggests that the scope for spillover effects was limited, which is consistent with the notion that in the French school system most peer interactions take place within the class (Avvisati et al., 2014). In the treatment group, 58 percent of Grade 10 students and 63 percent of Grade 12 students report having talked about the classroom intervention with their classmates, but they are only 24 percent and 27 percent to report having talked with schoolmates from other classes, respectively (see Appendix Table J1). In the control group, only 14 percent of students in Grade 10 report having heard of the classroom visits, mostly in a vague manner (12 percent). In Grade 12, students in the control group are more likely (34 percent) to report being at least vaguely aware of the visits, but under 5 percent of boys and girls have a precise recollection. Overall, these summary statistics suggest that spillover effects were quite limited.

We complement this survey evidence by investigating more formally whether the interventions affected the higher education choices of Grade 12 students whose classes were not assigned to the treatment group—either classes not selected by principals for the interventions or participating classes randomly assigned to the control group. Our empirical strategy, described in detail in Appendix J, builds on the following intuition: for schools that participated in the evaluation, the random assignment of treatment to participating classes makes it possible to estimate

²⁶The difference between the effects of visits before and after February 1 is statistically significant at the 5 percent level for girls and is robust to controlling for possible improvement or decline in the quality of role models' interventions over time, through the inclusion of fixed effects for the chronological order of the role models' classroom visits, i.e., first, second, etc.

the average outcome that would have resulted if *all* students had only been exposed to the spillover effects of classroom interventions without being *directly* exposed to a role model. This unobserved “spillover-only” counterfactual can be estimated at the school level by computing an appropriately weighted average of the outcome of students in the non-participating classes and in the participating classes that were assigned to the control group. Students in the control group classes are given a greater weight as they are used to account for both their own outcome and for the outcome that would have been observed in the treatment group classes, had their students only been indirectly exposed to a role model.²⁷ The spillover effects of the role model interventions are then estimated by comparing the “spillover-only” counterfactual to a “no-treatment” counterfactual. This second counterfactual is constructed using non-participating schools, which we observe in the administrative data, that have similar observable characteristics as the participating ones over the period 2012–2015. Having verified that trends in student enrollment outcomes were parallel between the two groups of schools in the pre-treatment period, we implement a difference-in-differences estimator to identify the interventions’ spillover effects on students’ STEM enrollment outcomes at college entry.

The results based on this difference-in-differences approach show no evidence of significant spillover effects of classroom visits on non-treated Grade 12 students (see Table J2 in the Appendix). Together with the survey evidence, they suggest that spillovers between treatment and control classes were at most limited.

6 How Do Role Models Affect Student Behavior?

This section inquires into how light-touch classroom interventions by female role models with a background in science can affect girls’ choice of study at university. Our insights are derived from comparison of groups of students who were exposed to different role models or who responded differently to a given role model.

We proceed in three steps. First, we show that the treatment effects on STEM enrollment outcomes vary widely according to the two most salient dimensions of heterogeneity in the current setting, namely students’ academic performance and role models’ background (professionals employed by L’Oréal vs. young researchers). We then provide a more systematic analysis of the heterogeneity of treatment effects using machine learning techniques. Following the approach developed by Chernozhukov et al. (2018), we identify the characteristics of the students and

²⁷For instance, in a school with two participating classes, one treated and one control, and one non-participating class, the “spillover-only” counterfactual is computed by assigning a weight of 1 to the non-participating class and a weight of 2 to the control group class (if all classes have the same number of students). By virtue of randomization, mean outcomes in the control group classes provide unbiased estimates of the counterfactual “spillover-only” outcomes in the treatment group classes.

role models for whom we observe particularly large (or small) treatment effects on students' choice of study, i.e. the *final* or *behavioral* outcome, as well as on their perceptions, self-concept and interest for science, i.e. the possible *channels of influence*. Finally, we extend the approach of Chernozhukov et al. (2018) to estimate the correlations between individual-level treatment effects on different outcomes conditional on exogenous observable characteristics. In doing so, we seek to determine whether the students who were particularly receptive or unreceptive to some of the messages conveyed during the interventions are also those whose choice of study was most or least affected by the interventions.

6.1 Heterogeneous Treatment Effects on STEM Enrollment

We start by investigating how the treatment effects on STEM enrollment vary with math performance and role model background. Our analysis focuses on Grade 12 students, as we find no evidence of significant effects on STEM enrollment for Grade 10 students.²⁸

High vs. low achievers in math. Applicants' performance in mathematics is the single most important admission criterion of selective undergraduate STEM programs. Using Grade 12 students' national percentile rank on the *Baccalauréat* math test to proxy for academic performance, we find that the interventions' positive impact on selective STEM enrollment is driven by female students in the top quartile (see Figure 3).²⁹ For these students, the probability of enrolling in a selective STEM program after high school increases by 12.9 percentage points, which corresponds to a 53 percent increase from the baseline of 24.3 percent. While the interventions also appear to have induced some male students in the top quartile to enroll in selective STEM programs, the effect is much smaller (6.5 percentage points, or a 14 percent increase over the baseline of 45 percent) and is not statistically significant. Especially striking is the fact that among the top quartile of achievers in math, the gender gap in the probability of enrolling in a selective STEM program is the largest (20.7 percentage points) and the treatment reduces it by 6.4 percentage points, which corresponds to a 31 percent reduction from the baseline.³⁰

Role model background: researchers vs. professionals. It is unclear, a priori, how the different types of role models differ in their effects on students' attitudes and behavior. As shown in Table 2, role models with a research background are, on average, younger than the

²⁸The results of the heterogeneity analysis by level of performance in math and role model background for Grade 10 students can be found in Panel A of Appendix Tables K1 and K2.

²⁹As discussed in Section 4.3, we find no significant impact of the interventions on students' performance on the math test of the *Baccalauréat* exam, which mitigates concerns about potential endogenous selection bias when conditioning on this variable.

³⁰The differences in treatment effects between high and low achievers in math are qualitatively similar for enrollment in male-dominated STEM programs (Figure 3, Panel B) as well as in all types of STEM programs (Appendix Table K1, Panel B).

professionals employed by the sponsoring firm, which may foster a stronger sense of identification by the students. But because they work in highly specialized fields and in very competitive environments, it is not clear how attainable students might think their achievements are. On the other hand, the professionals tend to have higher pay and more experience, and come less often from a purely academic background. They also hold permanent positions, unlike Ph.D. candidates and postdocs. Finally, their working environment could be perceived as particularly attractive by students, given the firm’s commitment to promote diversity and gender equality.

We find clear evidence that the two groups of role models had contrasting effects on STEM enrollment outcomes for girls in Grade 12. The left panel of Figure 4 shows that the professionals increased female students’ probability of enrolling in a selective STEM program by a significant 5.3 percentage points, whereas researchers had no detectable effect.³¹ The contrast is qualitatively similar, although less pronounced, when we consider enrollment in male-dominated STEM programs (right panel of Figure 4) or across all STEM programs (see Appendix Table K2). While the estimates also point to larger effects for boys who were exposed to role models with a professional background, they are not statistically significant at conventional levels.

Even though the role models were not randomly assigned to schools, the characteristics of the schools and students that they visited appear to be reasonably balanced between the two types of role model participants, with only few statistically significant differences (see Appendix Tables E7 and E8). Moreover, we show that the significantly larger impact of professionals on selective STEM enrollment for Grade 12 girls is robust to controlling for a full set of interactions between the treatment group dummy and a rich set of observable characteristics of students and schools (see Appendix Table K3). We therefore find no evidence that the heterogeneous treatment effects by role models’ background are confounded by differences in the characteristics of the classes they visited.³²

6.2 Machine Learning to Uncover Sources of Heterogeneity

Investigating treatment effect heterogeneity by splitting the sample into subgroups inevitably entails the risk of data mining. To address this concern, we carry out a systematic exploration of treatment effect heterogeneity using machine learning (ML) methods (see Athey and Imbens, 2017, for a review). Specifically, we adopt the approach recently developed by Chernozhukov et

³¹The difference between the treatment effects of the two groups of role models on Grade 12 girls’ probability of enrolling in a selective STEM program is significant at the 5 percent level.

³²We also explored whether the effects of the role model interventions could be mediated by the subsequent interactions between the students and the teacher who was present during the visit. For instance, science teachers could be inclined to reiterate the role model’s messages about science-related careers while female teachers could amplify the effects of the interventions for female students. Using data from the post-intervention role model survey, we do not find support for these hypotheses (results available upon request): the treatment effects on female enrollment in selective or male-dominated STEM do not vary significantly according to the teacher’s gender or taught subject.

al. (2018) to estimate conditional average treatment effects (CATE). A brief description is given below; a more detailed discussion can be found in Appendix L.

General description of Chernozhukov et al. (2018)’s approach. Let $Y(1)$ and $Y(0)$ denote the potential outcomes of a student when her class is and is not visited by a role model, respectively. Let Z be a vector of covariates that characterize the student and the role model who visited the class. The conditional average treatment effect (CATE), denoted by $s_0(Z)$, is defined as:

$$s_0(Z) \equiv \mathbb{E}[Y(1) - Y(0)|Z].$$

Because it is hard to obtain uniformly valid inference on the CATE without making strong assumptions, the approach in Chernozhukov et al. (2018) consists in conducting inference on specific *features* of the CATE, such as the expectation of $s_0(Z)$ in groups defined using a given ML predictor $S(Z)$.

The first feature examined is the Best Linear Predictor (BLP) of $s_0(Z)$ given $S(Z)$. The authors show that the BLP can be identified from the following weighted linear projection:

$$Y = \alpha_0 + \alpha B(Z) + \beta_1(T - p(Z)) + \beta_2(T - p(Z))(S(Z) - \mathbb{E}[S(Z)]) + \epsilon, \quad E[w(Z)\epsilon X] = 0, \quad (3)$$

where T is a dummy for treatment assignment; $B(Z)$ is an ML predictor of $Y(0)$ obtained from the training sample; $p(Z)$ is the probability of being treated conditional on the covariates Z ; $w(Z) = \{p(Z)(1 - p(Z))\}^{-1}$ is the weight; and X denotes the vector of all regressors ($X \equiv [1, B(Z), T - p(Z), (T - p(Z))(S(Z) - \mathbb{E}[S(Z)])]$). This projection identifies the parameters $\beta_1 = \mathbb{E}[s_0(Z)]$ and $\beta_2 = \text{Cov}(s_0(Z), S(Z))/\text{Var}(S(Z))$, which can both be estimated using the empirical analog of Equation (3). We refer to β_1 and β_2 in the tables as the average treatment effect (ATE) and heterogeneity loading (HET) parameters, respectively. The key parameter of interest, β_2 , is informative about the correlation between the true and the predicted CATE. It is equal to one if the prediction is perfect and to zero if there is no treatment effect heterogeneity or if $S(Z)$ has no predictive power.

The main purpose of estimating β_2 is to check if the trained ML methods are able to detect heterogeneity in the treatment effect. If so, the ML predictor of the CATE can be used to identify groups of individuals with the smallest and largest treatment effects. Heterogeneity groups are constructed by sorting students in the estimation sample based on the value of $S(Z_i)$, the predicted value of each student’s treatment effect given his/her observable characteristics Z_i . We consider the bottom and top quintiles of $S(Z_i)$ and report ITT estimates for both groups of students—a feature of the CATE called Sorted Group Average Treatment Effects (GATEs) in Chernozhukov et al. (2018). We then compare the distribution of observable characteristics in

the two groups—a feature called Classification Analysis (CLAN).

Practical Implementation. We consider five alternative ML methods to estimate the predictor $S(Z)$: Elastic Net, Random Forest, Boosted Trees, Neural Network with feature extraction, and a simple linear model. To train these methods, we use as covariates Z three indicators for the education districts of Paris, Créteil, and Versailles, four indicators for students’ socioeconomic background (high, medium-high, medium-low, and low), their age, their overall percentile rank in the *Baccalauréat* exam, their percentile ranks in the French and math tests of the exam, and a vector of 56 role model fixed effects.³³ We limit ourselves to only a few exogenous student characteristics because our main objective is to document treatment effect heterogeneity across the role models participants. For each outcome, the best ML method for either the BLP or the GATEs targeting of the CATE is selected using the performance measures proposed by Chernozhukov et al. (2018).

To avoid overfitting, we estimate the features of the CATE given an ML predictor $S(Z)$ on an *estimation sample* that is distinct from the *training sample* used to obtain $S(Z)$. We follow Chernozhukov et al. (2018) in iterating this data-splitting process and reporting the medians of estimates and p -values over several splits. The nominal levels of p -values are further adjusted to guarantee uniform validity, which leads to conservative inference.

Heterogeneous treatment effects on enrollment outcomes. We use the above procedure to estimate the different features of the CATE on enrollment in selective or male-dominated STEM programs for girls in Grade 12.³⁴ The results are reported in Table 6.

In Panel A, the estimated ATEs of the interventions on Grade 12 girls’ enrollment in selective or male-dominated STEM are very close to those reported in Table 5 by virtue of the randomization of the sample splits. Turning to heterogeneity, the coefficients on the HET parameter indicate that the ML predictors are strongly and significantly correlated with the CATE on enrollment in selective STEM but not in male-dominated STEM.

Estimates of the sorted group average treatment effects (GATEs) for the top and bottom quintiles of the predicted treatment effects $S(Z)$ are reported in Panel B. They confirm the considerable heterogeneity of treatment effects on selective STEM enrollment among Grade 12 girls, GATEs ranging from a small negative effect in the bottom 20 percent to a large and significant 13.9 percentage point effect in the top 20 percent. The lesser heterogeneity in the effects on enrollment in male-dominated STEM is also confirmed, with no statistically significant difference between the top and bottom quintiles of treatment effects.

³³Each student in the control group is assigned to the role model who visited his or her school, so the role model fixed effects are defined for students in both the treatment and control groups.

³⁴The machine learning results for boys in Grade 12 are reported in Appendix Table L2.

Panel C describes the characteristics of the 20 percent most and least affected students (CLAN). The main takeaway is that the ML agnostic approach strongly confirms that the treatment effects on selective STEM enrollment are greater for high-achieving girls in math and for those who were exposed to a professional rather than a researcher role model. Between the 20 percent most and least affected female students, the average gap in math performance rank is as much as 63 percentile ranks; the difference in the probability that the class was visited by a professional is 14.8 percentage points. The results are qualitatively similar for enrollment in male-dominated STEM, but the differences between groups are smaller, which is consistent with the previous finding of less heterogeneous treatment effects for this outcome.

The results in Panel C disclose heterogeneous effects along other dimensions. The 20 percent of girls with the largest treatment effects on selective STEM enrollment perform significantly better in French and are from higher socioeconomic backgrounds, compared with the least affected 20 percent. They are also less likely to have been exposed to role models who have children or who graduated in a male-dominated STEM field (math, physics, engineering), and more likely to have been exposed to role models who participated in the FGiS program the year before. However, the fact that these characteristics are correlated both with students' math performance and with the role model being either a professional or a researcher makes it difficult to determine their specific contribution to treatment effect heterogeneity. As suggestive evidence, we performed a "horse race" by regressing enrollment in selective STEM on the interactions between the treatment group indicator and each of the student and role model characteristics listed in Panel C. The results, which are reported in Appendix Table K3 are consistent with the conclusion that math performance and role models' professional background are the two main observable dimensions of heterogeneity in the treatment effects on selective STEM enrollment.³⁵

Heterogeneous treatment effects on potential channels. To help identify the mechanisms behind the heterogeneity of effects on selective STEM enrollment among Grade 12 girls, we start by comparing the characteristics of those with the largest and smallest treatment effects for each of the potential channels of influence studied in Section 4, namely general perceptions of science-related careers and gender roles in science, taste for science subjects, math self-concept, and science-related career aspirations. The results are reported in Table 7. For each potential channel, we compare the characteristics of students in the top and bottom quintiles of predicted treatment effects. We focus on the two main sources of heterogeneity in the effects on enrollment in selective STEM, i.e., student performance in math and exposure to a role model with a

³⁵The effect of the interventions on selective STEM enrollment remains significantly greater for high-achieving girls in math and for those who were exposed to role models with a professional background when we interact the treatment group indicator with other student and role model characteristics.

professional background.³⁶

The first key finding is that professionals and researchers were equally effective in debunking stereotypes on gender differences in math aptitude, while they reinforced students’ perceptions that “women don’t really like science” and that “women face discrimination in science-related jobs” to a comparable extent. These results suggest that the “gender debiasing” component of the classroom interventions, which emphasized men’s and women’s equal predisposition for science, cannot explain, alone, why the interventions increased girls’ enrollment in selective STEM; otherwise the two groups of role models would be expected to have similar effects for this outcome, which is not what we find.

By contrast, Table 7 reveals that the professionals were better than the researchers at improving female students’ perceptions of science-related jobs and stimulating their aspirations for such careers, while emphasizing less the underrepresentation of women. Regarding perceptions of science-related careers, girls in the top quintile of treatment effects are 19.2 percentage points more likely to have been visited by a professional compared to girls in the bottom quintile, the difference being statistically significant at the 1 percent level. Professionals are similarly overrepresented among the role models who had the greatest effects on girls’ taste for science subjects (22.7 percentage-point gap between the top and bottom quintile of treatment effects), and even more so among those who raised science-related career aspirations the most (38.9 percentage-point gap). The opposite holds for heterogeneous treatment effects on the importance of female underrepresentation in STEM: compared to the 20 percent of girls least affected for this outcome, the 20 percent most affected are 11.2 percentage points more likely to have been visited by a researcher.

Together, these results provide a first description of the role models who were the most effective in changing female students’ stereotyped behaviors. In addition to conveying positive information on career paths, these role models succeeded in sparking genuine interest in science and science-related jobs without overemphasizing the consequences of gender stereotyping. These features are in line with the main mechanisms usually considered necessary for role models to work: generating a sense of fit while moderating the effects of stereotype threat.

The analysis of treatment effect heterogeneity by student math performance tends to confirm that the messages conveyed by professionals were more effective at influencing female students’ choice of study. Indeed, the students who were particularly receptive to these messages are also those for whom we find the strongest impact on STEM enrollment, i.e.,

³⁶The heterogeneity loading parameter of the BLP and the GATEs associated with the best ML method are reported separately for each outcome in Appendix Table L1. For the sake of completeness, Appendix Tables K1 and K2 show the results obtained via a more traditional heterogeneity analysis, i.e., comparing the LATEs for different subgroups of female students based on math performance and on the background of the role model. The conclusions are consistent with those deriving from the ML procedure.

high achievers in mathematics. Average math performance is significantly higher among the students whose perceptions of science-related careers and taste for science subjects improved the most. Conversely, we find fewer high achievers among the girls whose awareness of female underrepresentation in STEM and perception of gender discrimination increased the most.

While these comparisons on the basis of role model background and student math performance cannot be given a causal interpretation, they are consistent with the notion that gender-neutral messages about careers in science are more effective than gender-related messages to steer girls towards STEM studies. The next section provides additional evidence supporting this interpretation.

6.3 Correlation between Treatment Effects

So far, our discussion of the channels of influence has sought to identify the main dimensions of treatment effect heterogeneity on STEM enrollment outcomes and investigated how the impact on student perceptions varies along these dimensions. We now develop a more general approach to directly estimate the correlation between the treatment effects on different outcomes. This constitutes a methodological contribution that can be used in other randomized controlled trials to relate treatment effects on different outcomes. In our context, the proposed method allows us to answer the following question: given their observable characteristics, are the students with the largest treatment effects for a potential channel of influence Y^A the same ones who exhibit the largest treatment effects on enrollment outcome Y^B ?

A new feature of the CATE. Because treatment effects for a given student are never observed, the correlation between *individual-level* treatment effects on outcomes Y^A and Y^B cannot be estimated without making strong assumptions. Instead, our approach takes advantage of the predicted heterogeneity in treatment effects by student and role model characteristics to recover the correlation $\rho_{A,B|Z} = \text{Corr}(s_0^A(Z), s_0^B(Z))$ between the true CATEs on the two outcomes Y^A and Y^B , which we denote by $s_0^A(Z)$ and $s_0^B(Z)$, respectively.

A detailed description of our approach is provided in Appendix L. To estimate the correlation between $s_0^A(Z)$ and $s_0^B(Z)$, we first define a new feature of the CATE as a simple adaptation of Chernozhukov et al. (2018)’s method. Instead of estimating the Best Linear Predictor of $s_0^A(Z)$ based on the ML predictor $S^A(Z)$, we estimate the BLP of $s_0^A(Z)$ based on $S^B(Z)$, i.e., the ML predictor of the heterogeneity in treatment effects on outcome Y^B . The heterogeneity loading parameter of the BLP we are interested in is

$$\beta_2^{A|B} = \text{Cov}(s_0^A(Z), S^B(Z)) / \text{Var}(S^B(Z)). \quad (4)$$

This parameter is identified and can be estimated using a variant of Equation (3). By switching the roles of Y_A and Y_B in Equation (4), one can similarly estimate the heterogeneity loading parameter from the BLP of $s_0^B(Z)$ based on $S^A(Z)$, i.e.,

$$\beta_2^{B|A} = \text{Cov}(s_0^B(Z), S^A(Z)) / \text{Var}(S^A(Z)).$$

Writing $S^A(Z) = s_0^A(Z) + \eta_A$ and $S^B(Z) = s_0^B(Z) + \eta_B$ and assuming that the prediction errors η_A and η_B are independent of both predicted functions $s_0^A(Z)$ and $s_0^B(Z)$ in the estimation sample, we show that $\beta_2^{A|B}$ and $\beta_2^{B|A}$ have the same sign, which is indicative of whether the treatment effects on Y^A are positively or negatively correlated with the treatment effects on Y^B .

Under these assumptions, the correlation between the true CATEs on Y^A and Y^B , $\rho_{A,B|Z}$, can be estimated using the following formula:³⁷

$$\rho_{A,B|Z} = \text{Sign}(\beta_2^{A|B}) \frac{\sqrt{\beta_2^{A|B} \beta_2^{B|A}}}{\sqrt{\beta_2^{B|B} \beta_2^{A|A}}}, \quad (5)$$

where $\beta_2^{A|A}$ and $\beta_2^{B|B}$ are the heterogeneity loading parameters in the BLPs of $s_0^A(Z)$ and $s_0^B(Z)$ on their respective predictors $S^A(Z)$ and $S^B(Z)$.

Practical implementation. As in the previous section, we split the data into a training and an estimation sample. We obtain predictors $S^A(Z)$ and $S^B(Z)$ of $s_0^A(Z)$ and $s_0^B(Z)$ in the training sample and use them to estimate the four parameters $\beta_2^{A|A}$, $\beta_2^{B|B}$, $\beta_2^{A|B}$ and $\beta_2^{B|A}$ in the estimation sample. We then plug these parameter estimates in Equation (5) to obtain an estimate $\hat{\rho}_{A,B|Z}$ of the correlation between the CATEs on outcomes Y^A and Y^B . We use a bootstrap procedure, also performed in the estimation sample, to obtain a 95 percent confidence interval for $\hat{\rho}_{A,B|Z}$.³⁸ As in the previous section, we follow the procedure of Chernozhukov et al. (2018) so that our final estimate of $\rho_{A,B|Z}$ and its confidence interval are computed as medians of estimates obtained from several estimation samples, with nominal level of confidence intervals adjusted to guarantee uniform validity.

Results. The correlations between treatment effects for girls in Grade 12 are reported in Table 8, where the covariates that we use to predict treatment effect heterogeneity are the

³⁷While it is not possible to prove that the out-of-sample prediction error of a ML predictor is independent from the predicted outcome for any predictor, this assumption seems reasonable when using efficient ML algorithms such as those considered in this paper. As suggestive evidence, we have checked in Monte Carlo simulations that this assumption holds for a large set of simulated functions of Z , which are generated manually and predicted on subsamples of our data. We further checked that the correlation $\rho_{A,B|Z}$ is successfully recovered for various data-generating processes using the formula in Equation (5).

³⁸We report confidence intervals rather than p -values because the former are highly skewed, implying that the p -values obtained from bootstrap under normality assumptions are misleading.

same as in Table 6. They suggest that some channels were more important than others in steering female students towards STEM studies. The treatment effects on girls' enrollment in selective STEM exhibit a strong positive and significant correlation with the improvement in their perceptions of science-related careers ($\hat{\rho} = 0.96$) and with the improvement in their taste for science subjects ($\hat{\rho} = 0.71$).³⁹

While not statistically significant at the 5 percent level, the remaining correlations give some indication on the role of other candidate channels.⁴⁰ They confirm in particular that debiasing girls' attitudes towards gender differences in aptitude for math is not associated with increased enrollment in selective STEM programs ($\hat{\rho} = 0.19$ with a 95 percent confidence interval of $[-1.24, 2.05]$) and that, if anything, reinforcing the belief that women are discriminated in science careers tends to deter girls from enrolling in selective STEM programs ($\hat{\rho} = -0.34$ $[-2.22, 0.56]$). By contrast, raising girls' aspirations for careers in science is associated with an increased probability that they enroll in such programs ($\hat{\rho} = 0.36$ $[-0.51, 2.01]$).

Overall, the results based on correlations between treatment effects are in line with and extend those obtained in the previous section. They suggest that the most effective role models were those who managed to convey a positive image of science careers without overemphasizing women's underrepresentation and its possible causes.

7 Conclusion and Discussion

Based on a large-scale randomized field experiment involving 56 female role models and nearly 20,000 high school students in Grade 10 and Grade 12, this paper shows that a one-hour in-class exposure to a female scientist can improve students' perceptions of science careers and significantly increase female participation in STEM fields of study at college enrollment. Remarkably, the positive enrollment effects are observed only in the tracks with the most severe gender imbalance, which are the most prestigious and selective, and those that are most math-intensive. These effects can be expected to increase the future earnings of the target population, since the selective and male-dominated STEM programs offer high wage premia relative to other programs.

In our empirical setting, the role model interventions had no discernable effects on students' academic performance and only slightly improved their math self-concept, thus ruling out

³⁹The positive correlation between the treatment effects on taste for science and on enrollment in selective STEM suggests that students whose preferences were affected by the intervention also changed their choice of study. These effects, however, are highly heterogeneous (see Appendix Table L1): while the treatment effects on taste for science are positive for the 20 percent most affected girls in Grade 12, they are negative for the 20 percent least affected, resulting in an average treatment effect close to zero (see Appendix Table F2).

⁴⁰We report in Table 8 the lower and upper bounds for the lower and upper limits of the actual 95 percent confidence interval associated with each estimated correlation. Note that the (unknown) true confidence intervals are likely to be smaller than suggested by the bounds reported in this table.

these factors as primary causes of the observed effects on STEM enrollment. By contrast, the classroom visits significantly challenged students' stereotyped views of science careers and gender differences in aptitude for science. These effects, however, are observed for both genders in both grades, suggesting that by themselves they cannot explain why the role model interventions only affected the educational choices of girls in Grade 12.

Our results offer substantial evidence that female students' behavioral response to the role model interventions was mediated by their ability to identify with the female scientists to whom they were exposed. On the verge of important decisions about their future education and career pathways, girls in Grade 12 were more receptive than the other groups of students to the appealing image of science-related careers embodied by the role models. Consistent with this, we find that their improved perceptions of science careers translated into stronger aspirations for such careers. This process of identification was less likely to occur among Grade 10 girls, who are further away from career choices, and for boys in both grade levels, who may have found it more difficult to identify with women scientists. To confirm this latter hypothesis and, more generally, to improve our understanding of role model effects, an interesting avenue for future research would be to compare the impact of male and female role models in a similar context.

Another important insight from the study is that by heightening awareness of the underrepresentation of women in STEM, while at the same time emphasizing men's and women's equal aptitude for science, the interventions may have unintentionally reinforced students' beliefs that women dislike science and face discrimination in STEM careers. That is, there is suggestive evidence that excessive stress on gender can be counter-productive and that gender-neutral messages might be more effective in steering girls towards STEM fields. In our setting, the role models who most reinforced the perception that women are underrepresented and discriminated against in science had the least effect on selective STEM enrollment for female students in Grade 12, whereas those who most improved girls' perceptions of science careers had the greatest impact. These findings suggest that role model interventions need to be carefully designed to limit the potential discouragement effect of overemphasis on gender imbalances.

More generally, our heterogeneity analysis warns against the temptation to view role models as a one-size-fits-all remedy against female underrepresentation in STEM fields. Like Carrell et al. (2010), we find that role model effects on enrollment outcomes are concentrated among high-achieving girls in math. The effectiveness of this type of intervention in increasing female participation in STEM among lower-performing students remains an open question. Our study also highlights the importance of role models' profile in generating a sense of fit among students, as the effects on educational choices varied markedly across the participating female scientists. These results point to the need for further research on how the matching between role models and students can be optimized to make this particular type of intervention more effective.

References

- Anderson, Michael L.**, “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*, 2008, *103* (484).
- Athey, Susan and Guido Imbens**, “Recursive Partitioning for Heterogeneous Causal Effects,” *Proceedings of the National Academy of Sciences*, 2016, *113* (27), 7353–7360.
- **and Guido W. Imbens**, “The Econometrics of Randomized Experiments,” in Esther Duflo and Abhijit V. Banerjee, eds., *Handbook of Economic Field Experiments*, Vol. 1, Elsevier, 2017, pp. 73–140.
- Avvisati, Francesco, Marc Gurgand, Nina Guyon, and Éric Maurin**, “Getting Parents Involved: A Field Experiment in Deprived Schools,” *Review of Economic Studies*, 2014, *81* (1), 57–83.
- Beaman, Lori, Esther Duflo, Rohini Pande, and Petia Topalova**, “Female Leadership Raises Aspirations and Educational Attainment for Girls: A Policy Experiment in India,” *Science*, 2012, *335* (6068), 582–586.
- Beede, David, Tiffany Julian, David Langdon, George McKittrick, Beethika Khan, and Mark Doms**, “Women in STEM: A Gender Gap to Innovation,” 2011. U.S. Department of Commerce, Economics and Statistics Administration, Issue Brief No. 04-11.
- Benjamini, Yoav, Abba M. Krieger, and Daniel Yekutieli**, “Adaptive Linear Step-up Procedures that Control the False Discovery Rate,” *Biometrika*, 2006, *93* (3), 491–507.
- Bettinger, Eric P. and Bridget Terry Long**, “Do Faculty Serve as Role Models? The Impact of Instructor Gender on Female Students,” *American Economic Review*, 2005, *95* (2), 152–157.
- Betz, Diana E. and Denise Sekaquaptewa**, “My Fair Physicist? Feminine Math and Science Role Models Demotivate Young Girls,” *Social Psychological and Personality Science*, 2012, *3* (6), 738–746.
- Black, Dan A., Amelia M. Haviland, Seth G. Sanders, and Lowell J. Taylor**, “Gender Wage Disparities among the Highly Educated,” *Journal of Human Resources*, 2008, *43* (3), 630–650.
- Blau, Francine D. and Lawrence M. Kahn**, “The Gender Wage Gap: Extent, Trends, and Explanations,” *Journal of Economic Literature*, 2017, *55* (3), 789–865.
- Brown, Charles and Mary Corcoran**, “Sex-Based Differences in School Content and the Male-Female Wage Gap,” *Journal of Labor Economics*, 1997, *15* (3), 431–465.
- Bursztn, Leonardo and David Y. Yang**, “Misperceptions about Others,” 2021. NBER Working Paper No. 29168.
- Canaan, Serena and Pierre Mouganie**, “The Impact of Advisor Gender on Female Students’ STEM Enrollment and Persistence,” *Journal of Human Resources*, forthcoming.
- Carlana, Michela**, “Implicit Stereotypes: Evidence from Teachers’ Gender Bias,” *Quarterly Journal of Economics*, 2019, *134* (3), 1163–1224.
- Carrell, Scott E., Marianne E. Page, and James E. West**, “Sex and Science: How Professor Gender Perpetuates the Gender Gap,” *Quarterly Journal of Economics*, 2010, *125* (3), 1101–1144.
- Chernozhukov, Victor, Mert Demirer, Esther Duflo, and Iván Fernández-Val**, “Generic Machine Learning Inference on Heterogeneous Treatment Effects in Randomized Experiments,” 2018. NBER Working Paper No. 24678.

- Cheryan, Sapna, John Oliver Siy, Marissa Vichayapai, Benjamin J. Drury, and Saenam Kim**, “Do Female and Male Role Models who Embody STEM Stereotypes Hinder Women’s Anticipated Success in STEM?,” *Social Psychological and Personality Science*, 2011, 2 (6), 656–664.
- Coffman, Katherine Baldiga**, “Evidence on Self-Stereotyping and the Contribution of Ideas,” *Quarterly Journal of Economics*, 2014, 129 (4), 1625–1660.
- Del Carpio, Lucia and Maria Guadalupe**, “More Women in Tech? Evidence from a Field Experiment Addressing Social Identity,” 2018. CEPR Discussion Paper DP13234.
- Gladstone, Jessica and Andrei Cimpian**, “Role Models Can Help Make the Mathematics Classroom More Inclusive,” 2020. OSF Preprints.
- Goldin, Claudia**, “A Grand Gender Convergence: Its Last Chapter,” *American Economic Review*, 2014, 104 (4), 1091–1119.
- Hoogendoorn, Sander, Hessel Oosterbeek, and Mirjam van Praag**, “The Impact of Gender Diversity on the Performance of Business Teams: Evidence from a Field Experiment,” *Management Science*, 2013, 59 (7), 1514–1528.
- Lavy, Victor and Edith Sand**, “On the Origins of Gender Gaps in Human Capital: Short- and Long-Term Consequences of Teachers’ Biases,” *Journal of Public Economics*, 2018, 167(C), 263–269.
- Lemaire, Sylvie**, “Parcours dans l’enseignement supérieur: devenir des bacheliers 2008,” 2018. Note d’Information 12.10, MESR-SIES.
- Lockwood, Penelope and Ziva Kunda**, “Superstars and Me: Predicting the Impact of Role Models on the Self,” *Journal of Personality and Social Psychology*, 1997, 73 (1), p. 91.
- Mullainathan, Sendhil and Jann Spiess**, “Machine Learning: An Applied Econometric Approach,” *Journal of Economic Perspectives*, 2017, 31 (2), 87–106.
- National Science Foundation**, *Women, Minorities, and Persons with Disabilities in Science and Engineering: 2017*, National Science Foundation and National Center for Science and Engineering Statistics, 2017. Special Report NSF 17-310. Arlington, VA.
- Nguyen, Trang**, “Information, Role Models and Perceived Returns to Education: Experimental Evidence from Madagascar,” 2008. Manuscript.
- O’Brien, Laurie T., Aline Hitti, Emily Shaffer, Amanda R. Van Camp, Donata Henry, and Patricia N. Gilbert**, “Improving Girls’ Sense of Fit in Science: Increasing the Impact of Role Models,” *Social Psychological and Personality Science*, 2016, 8 (3), 301–309.
- Porter, Catherine and Danila Serra**, “Gender Differences in the Choice of Major: The Importance of Female Role Models,” *American Economic Journal: Applied Economics*, 2020, 12 (3), 226–254.
- Riise, Julie, Barton Willage, and Willén Alexander**, “Can Female Doctors Cure the Gender STEMM Gap? Evidence from Randomly Assigned General Practitioners,” *The Review of Economics and Statistics*, forthcoming.
- Terrier, Camille**, “Boys Lag Behind: How Teachers’ Gender Biases Affect Student Achievement,” *Economics of Education Review*, 2020, 77.
- Wager, Stefan and Susan Athey**, “Estimation and Inference of Heterogeneous Treatment Effects using Random Forests,” *Journal of the American Statistical Association*, 2018, 113 (523), 1228–1242.
- Weinberger, Catherine J.**, “Mathematical College Majors and the Gender Gap in Wages,” *Industrial Relations*, 1999, 38 (3), 407–413.

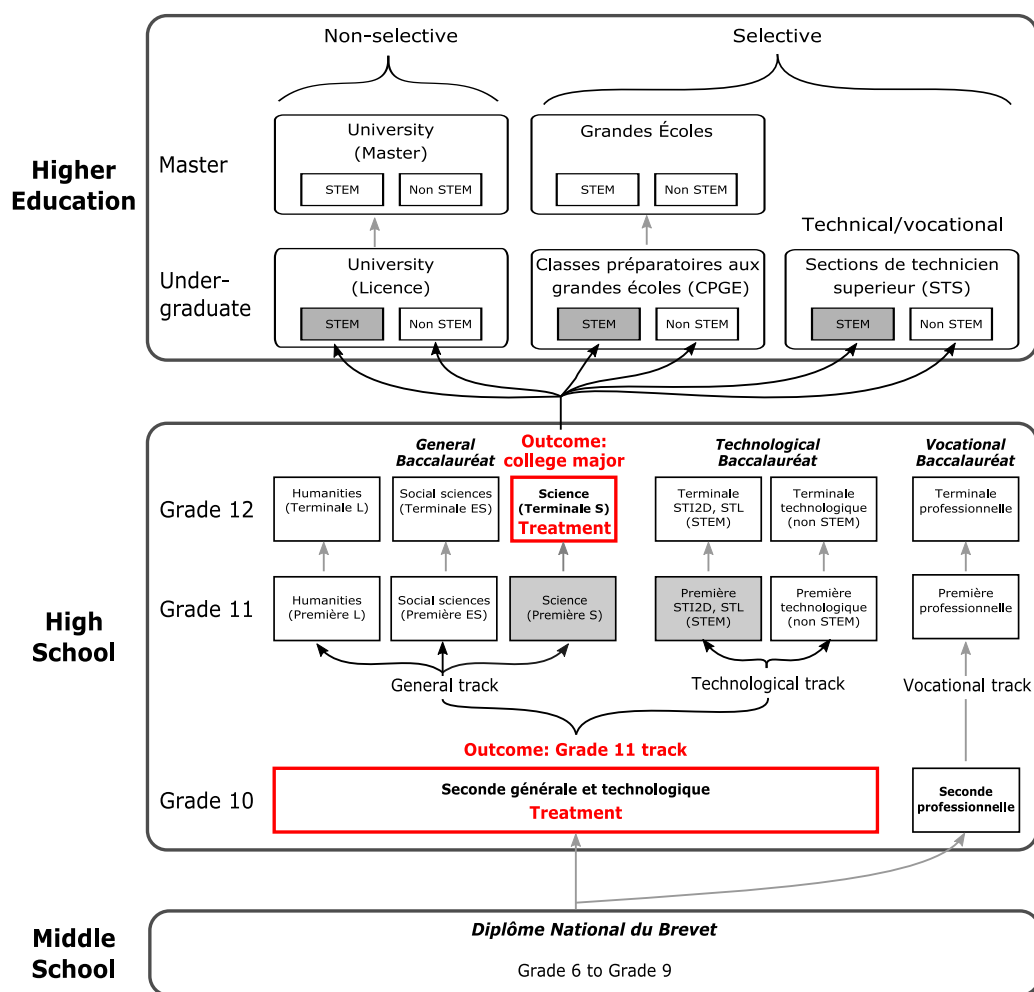


Figure 1 – Tracks in Secondary and Post-Secondary Education in France

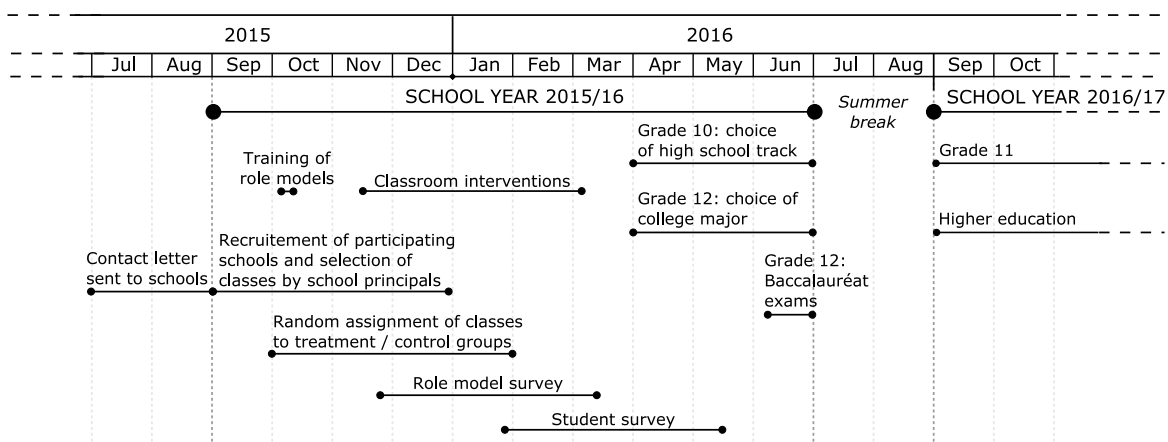


Figure 2 – Program Evaluation Timeline

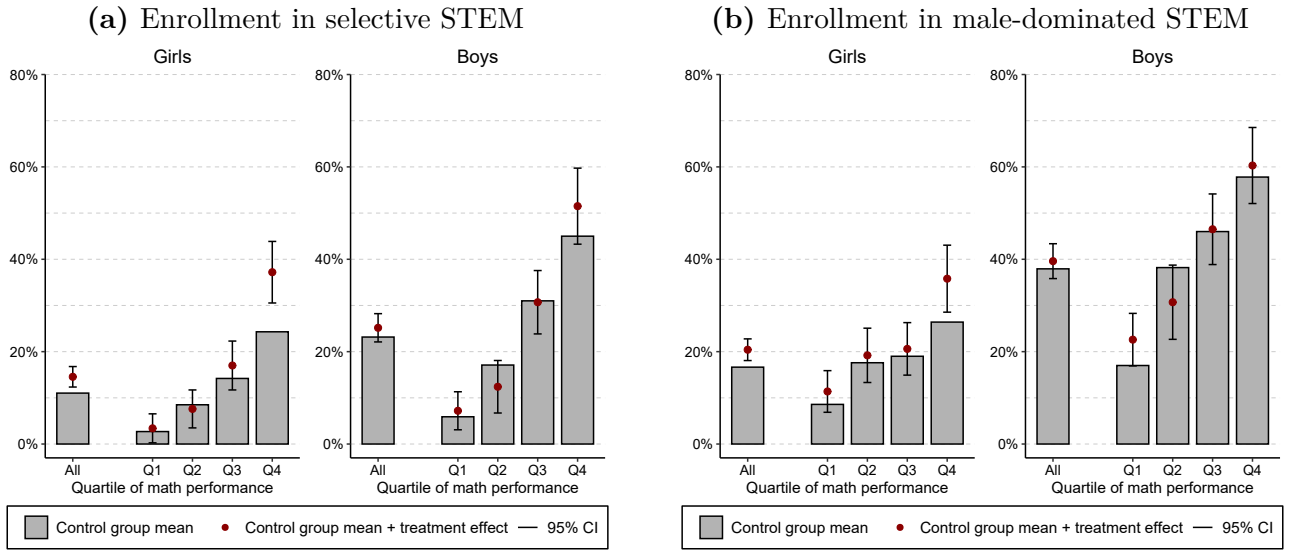


Figure 3 – Grade 12 Students: Enrollment in Selective and Male-Dominated STEM Undergraduate Programs, by Gender and Quartiles of *Baccalauréat* Performance in Math

Notes: The figure shows the fraction of Grade 12 (science track) students enrolled in selective (Panel A) and in male-dominated (Panel B) STEM undergraduate programs in the year following high school graduation, separately for girls and boys. The filled bars indicate the baseline enrollment rates among students in the control group, both overall and separately by quartile of *Baccalauréat* performance in math. The solid circles show the estimated treatment effects (added to the control group means), with 95 percent confidence intervals denoted by vertical capped bars. The local average treatment effects are estimated from a regression of the outcome of interest on interactions between a classroom visit indicator and the quartile of math performance, using treatment assignment (interacted with the quartiles of math performance) as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors are adjusted for clustering at the unit of randomization (class).

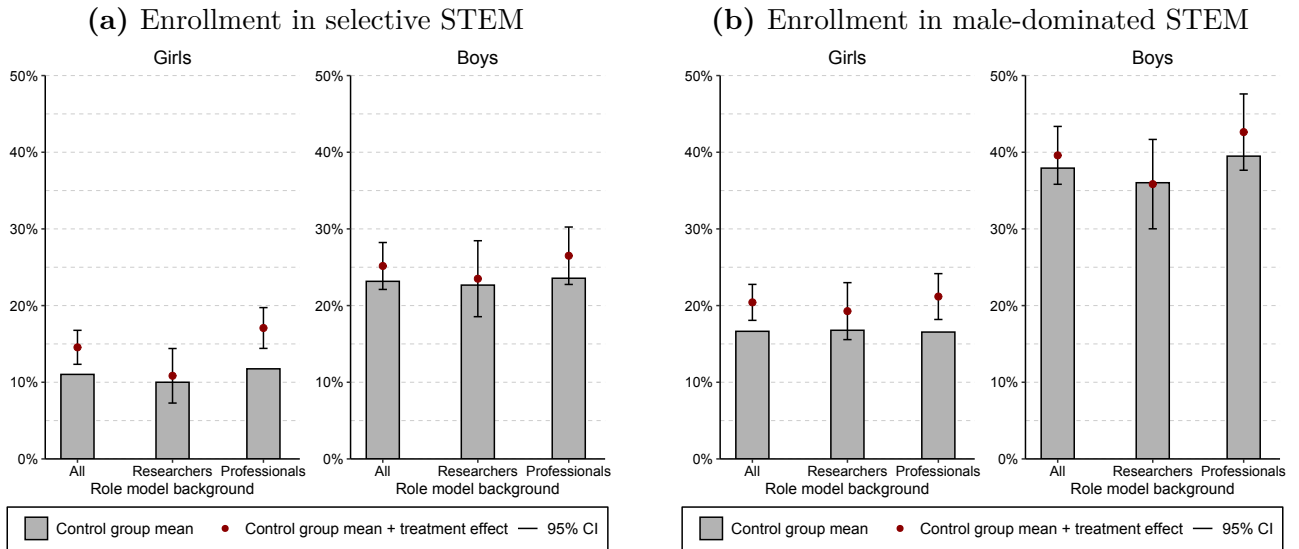


Figure 4 – Grade 12 Students: Enrollment in Selective and Male-Dominated STEM Undergraduate Programs, by Gender and Role Model Background

Notes: The figure shows the fraction of Grade 12 (science track) students enrolled in selective (Panel A) and in male-dominated (Panel B) STEM undergraduate programs after graduating from high school, separately for girls and boys. The filled bars indicate the baseline enrollment rates among students in the control group, both overall and separately by type of female role model who visited the classroom (researcher or professional). The solid dots show the estimated treatment effects (added to the control group means), with 95 percent confidence intervals denoted by vertical capped bars. The local average treatment effects are estimated from a regression of the outcome of interest on interactions between a classroom visit indicator and two indicators for role model type, using treatment assignment (interacted with role model type) as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors are adjusted for clustering at the unit of randomization (class).

Table 1 – Treatment-Control Balance

			Within school	
	Control group (1)	Treatment group (2)	Difference T–C (3)	<i>p</i> -value of diff. (4)
Panel A. Grade 10				
<i>Student characteristics</i>				
Female	0.535	0.522	−0.010	0.309
Age (years)	15.13	15.12	−0.01	0.180
Non-French	0.059	0.061	0.002	0.652
High SES	0.377	0.386	0.008	0.321
Medium- high SES	0.131	0.125	−0.007	0.168
Medium-low SES	0.248	0.235	−0.012	0.064
Low SES	0.244	0.254	0.012	0.085
Number of siblings	1.485	1.486	0.003	0.904
Class size	33.22	33.27	0.07	0.476
At least one science elective course	0.389	0.398	0.005	0.820
At least one standard elective course	0.770	0.737	−0.031	0.138
DNB percentile rank in math	58.61	58.35	−0.35	0.533
DNB percentile rank in French	57.79	57.91	0.12	0.829
<i>Test of joint significance</i>	<i>F</i> -stat: 0.798 (<i>p</i> -value: 0.653)			
<i>Predicted track in Grade 11</i>				
Grade 11: Science track	0.449	0.452	0.001	0.922
Grade 11: Science - general track	0.373	0.375	0.001	0.920
Grade 11: Science - technological track	0.077	0.077	−0.000	0.989
N	6,801	6,899	13,700	
Panel B. Grade 12 (science track)				
<i>Student characteristics</i>				
Female	0.499	0.484	−0.014	0.292
Age (years)	17.14	17.11	−0.04	0.000
Non-French	0.053	0.048	−0.006	0.275
High SES	0.453	0.474	0.029	0.009
Medium-high SES	0.136	0.135	−0.001	0.829
Medium-low SES	0.216	0.201	−0.015	0.023
Low SES	0.195	0.190	−0.012	0.140
Number of siblings	1.510	1.487	−0.032	0.127
Class size	31.75	32.19	0.39	0.196
DNB percentile rank in math	74.17	73.95	0.20	0.699
DNB percentile rank in French	69.31	69.90	0.89	0.122
<i>Test of joint significance</i>	<i>F</i> -stat: 0.983 (<i>p</i> -value: 0.459)			
<i>Predicted undergraduate major</i>				
Major: STEM	0.382	0.384	0.003	0.352
Major: selective STEM	0.175	0.178	0.006	0.081
Major: male-dominated STEM	0.273	0.276	0.004	0.279
N	2,853	2,898	5,751	

Notes: Each row corresponds to a different linear regression with the dependent variable listed on the left, separately for students in Grade 10 (Panel A) and in Grade 12 (Panel B). Columns 1 and 2 show the average value for students in the control and treatment groups, respectively. Column 3 reports the coefficient from the regression of each variable on the treatment group indicator, with the *p*-value reported in column 4. The regression controls for school fixed effects to account for the fact that randomization was stratified by school, and standard errors are adjusted for clustering at the unit of randomization (class). The *F*-statistic is from a test of the joint significance of the coefficients in a regression of the treatment group indicator on all student characteristics. High school tracks (Panel A) and undergraduate majors (Panel B) are predicted for each student using the coefficients from a linear regression of the corresponding binary variable (e.g., enrollment in a STEM major) on all student characteristics listed in the table. This model is fitted separately by grade level on the sample of students in the control group.

Table 2 – Female Role Models: Summary Statistics

	All role models	Researchers (Ph.D./ Postdoc)	Professionals (employed by sponsoring firm)
	(1)	(2)	(3)
Age (N=51)	33.3 (5.7)	30.0 (3.1)	35.6 (6.0)
Non-French	0.14	0.10	0.17
holds/prepares for a Ph.D. (N=55)	0.62	1.00	0.38
Graduated from a Grande École	0.39	0.33	0.43
Field: Math, Physics, Engineering	0.23	0.38	0.14
Field: Earth and Life Sciences	0.64	0.62	0.66
Field: Other	0.13	0.00	0.20
Has children (N=52)	0.42	0.19	0.58
Participated in the program the year before	0.25	0.19	0.29
Number of high schools visited	1.8 (0.8)	2.1 (0.9)	1.6 (0.7)
Number of classroom interventions	5.2 (2.3)	5.9 (2.3)	4.7 (2.1)
N	56	21	35

Notes: The summary statistics are computed based on information obtained from the L'Oréal Foundation and from the post-intervention survey administered online to collect feedback about the classroom visits. Standard deviations are shown in parentheses below the mean values. Where data are missing for some role models, the number of non-missing values N is indicated in parentheses.

Table 3 – Perceptions of Science-Related Careers

	Girls			Boys		
	Control group mean (1)	Treatment effect (LATE) (2)	<i>p</i> -value [<i>q</i> -value] (3)	Control group mean (4)	Treatment effect (LATE) (5)	<i>p</i> -value [<i>q</i> -value] (6)
Panel A. Grade 10						
Positive perceptions of science-related careers (index)	−0.020	0.245*** (0.028)	0.000 [0.001]	0.023	0.167*** (0.029)	0.000 [0.001]
Science-related jobs require long years of study	0.839	−0.087*** (0.010)	0.000 [0.001]	0.849	−0.074*** (0.010)	0.000 [0.001]
Science-related jobs are monotonous	0.290	−0.032*** (0.012)	0.006 [0.011]	0.318	−0.006 (0.013)	0.633 [0.634]
Science-related jobs are solitary	0.325	−0.061*** (0.012)	0.000 [0.001]	0.303	−0.062*** (0.011)	0.000 [0.001]
Science-related jobs pay higher wages	0.637	0.008 (0.014)	0.535 [0.536]	0.668	0.015 (0.013)	0.237 [0.297]
Hard to maintain work-life balance	0.297	−0.026** (0.012)	0.026 [0.033]	0.283	−0.029** (0.012)	0.014 [0.023]
N		6,475			5,751	
Panel B. Grade 12 (science track)						
Positive perceptions of science-related careers (index)	−0.003	0.312*** (0.034)	0.000 [0.001]	0.003	0.155*** (0.033)	0.000 [0.001]
Science-related jobs require long years of study	0.666	−0.110*** (0.015)	0.000 [0.001]	0.719	−0.091*** (0.014)	0.000 [0.001]
Science-related jobs are monotonous	0.169	−0.019 (0.013)	0.141 [0.141]	0.233	−0.026 (0.016)	0.114 [0.143]
Science-related jobs are solitary	0.228	−0.088*** (0.012)	0.000 [0.001]	0.206	−0.047*** (0.013)	0.000 [0.001]
Science-related jobs pay higher wages	0.531	0.059*** (0.018)	0.001 [0.002]	0.576	0.027* (0.016)	0.093 [0.143]
Hard to maintain work-life balance	0.225	−0.049*** (0.015)	0.001 [0.002]	0.167	−0.012 (0.011)	0.260 [0.260]
N		2,600			2,636	

Notes: This table reports estimates of the treatment effects of classroom interventions on students' perceptions of science-related careers, separately by grade level and gender. The sample is restricted to students who completed the post-intervention questionnaire. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 4 report the average value for students in the control group. Columns 2 and 5 report the local average treatment effect (LATE) estimates. They are obtained from a regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report the cluster-robust *p*-value of the estimated treatment effect and, in square brackets, the *p*-value (*q*-value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage *q*-values introduced in Benjamini et al. (2006) and described in Anderson (2008). The *q*-values associated with the synthetic index (highlighted in bold) are adjusted for multiple testing across the study's nine main outcomes of interest, separately by grade level and gender (see Appendix D for details). The *q*-values for the individual components of the index are adjusted for multiple testing across the index components, separately by grade level and gender. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

Table 4 – Perceptions of Gender Roles in Science, Stated Preferences and Self-Concept

	Girls			Boys		
	Control group mean (1)	Treatment effect (LATE) (2)	<i>p</i> -value [<i>q</i> -value] (3)	Control group mean (4)	Treatment effect (LATE) (5)	<i>p</i> -value [<i>q</i> -value] (6)
Panel A. Grade 10						
<i>Perceptions of gender roles in science</i>						
More men in science-related jobs	0.628	0.156*** (0.013)	0.000 [0.001]	0.629	0.168*** (0.014)	0.000 [0.001]
Equal gender aptitude for math (index)	0.115	0.109*** (0.025)	0.000 [0.001]	−0.134	0.148*** (0.030)	0.000 [0.001]
Women don't really like science	0.157	0.059*** (0.011)	0.000 [0.001]	0.198	0.103*** (0.013)	0.000 [0.001]
W face discrimination in science-related jobs	0.603	0.127*** (0.013)	0.000 [0.001]	0.527	0.153*** (0.014)	0.000 [0.001]
<i>Stated preferences and self-concept</i>						
Taste for science subjects (index)	−0.169	−0.038 (0.036)	0.294 [0.442]	0.197	−0.019 (0.031)	0.533 [0.685]
Math self-concept (index)	−0.198	−0.008 (0.031)	0.806 [0.807]	0.231	0.039 (0.032)	0.217 [0.326]
Science-related career aspirations (index)	−0.103	0.012 (0.030)	0.695 [0.807]	0.120	0.007 (0.029)	0.801 [0.902]
N		6,475			5,751	
Panel B. Grade 12 (science track)						
<i>Perceptions of gender roles in science</i>						
More men in science-related jobs	0.712	0.125*** (0.016)	0.000 [0.001]	0.717	0.149*** (0.015)	0.000 [0.001]
Equal gender aptitude for math (index)	0.158	0.095*** (0.028)	0.001 [0.002]	−0.161	0.132*** (0.040)	0.001 [0.002]
Women don't really like science	0.074	0.044*** (0.009)	0.000 [0.001]	0.146	0.073*** (0.015)	0.000 [0.001]
W face discrimination in science-related jobs	0.624	0.095*** (0.020)	0.000 [0.001]	0.600	0.072*** (0.018)	0.000 [0.001]
<i>Stated preferences and self-concept</i>						
Taste for science subjects (index)	−0.002	0.016 (0.034)	0.632 [0.633]	0.002	−0.000 (0.039)	0.998 [0.999]
Math self-concept (index)	−0.184	0.050 (0.039)	0.202 [0.228]	0.187	0.072** (0.035)	0.041 [0.062]
Science-related career aspirations (index)	−0.045	0.113*** (0.037)	0.002 [0.003]	0.046	0.050 (0.033)	0.131 [0.169]
N		2,600			2,636	

Notes: This table reports estimates of the treatment effects of classroom interventions on students' perceptions of gender roles in science, taste for science subjects, math self-concept, and science-related career aspirations, separately by grade level and gender. The sample is restricted to students who completed the post-intervention questionnaire. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 4 report the average value for students in the control group. Columns 2 and 5 report the local average treatment effect (LATE) estimates. They are obtained from a regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report the cluster-robust *p*-value of the estimated treatment effect and, in square brackets, the *p*-value (*q*-value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage *q*-values introduced in Benjamini et al. (2006) and described in Anderson (2008). The *q*-values are adjusted for multiple testing across the study's nine main outcomes of interest, separately by grade level and gender (see Appendix D for details). *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

Table 5 – Enrollment Status the Following Year

	Girls			Boys		
	Control group mean (1)	Treatment effect (LATE) (2)	p -value [q -value] (3)	Control group mean (4)	Treatment effect (LATE) (5)	p -value [q -value] (6)
Panel A. Grade 10						
All STEM tracks						
Grade 11: Science track	0.355	−0.004 (0.014)	0.753 [0.807]	0.551	−0.002 (0.015)	0.910 [0.910]
General vs. technological STEM track						
Grade 11: Science - general track	0.328	0.001 (0.013)	0.942 [0.942]	0.416	0.007 (0.014)	0.613 [0.614]
Grade 11: Science - technological track	0.026	−0.005 (0.003)	0.128 [0.256]	0.135	−0.009 (0.008)	0.300 [0.601]
Other tracks or repeater						
Grade 11: Other tracks	0.545	0.006 (0.014)	0.642	0.324	0.018 (0.014)	0.191
Repeater or dropout	0.101	−0.002 (0.009)	0.818	0.126	−0.016* (0.009)	0.070
N		7,241			6,459	
Panel B. Grade 12 (science track)						
All undergraduate STEM majors						
Major: STEM	0.289	0.024* (0.014)	0.080 [0.103]	0.470	0.003 (0.020)	0.886 [0.998]
Selective vs. non-selective STEM						
Major: selective STEM	0.110	0.035*** (0.011)	0.002 [0.004]	0.232	0.020 (0.016)	0.200 [0.283]
Major: non-selective STEM	0.178	−0.011 (0.011)	0.322 [0.322]	0.239	−0.017 (0.014)	0.212 [0.283]
Male- vs. female-dominated STEM						
Major: male-dominated STEM (math, physics, computer science)	0.166	0.038*** (0.012)	0.002 [0.004]	0.379	0.017 (0.019)	0.387 [0.388]
Major: female-dominated STEM (earth and life sciences)	0.123	−0.015 (0.010)	0.158 [0.211]	0.091	−0.014 (0.009)	0.119 [0.283]
Other tracks or dropout						
Other non-STEM programs	0.507	−0.032** (0.016)	0.045	0.293	−0.005 (0.014)	0.717
Not enrolled in a post-graduate curriculum	0.206	0.008 (0.015)	0.581	0.237	0.004 (0.016)	0.814
N		2,827			2,924	

Notes: This table reports estimates of the treatment effects of classroom interventions on students’ enrollment outcomes in the academic year following the classroom interventions, i.e. 2016/17, separately by grade level and gender. The enrollment outcomes are measured using student-level administrative data. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 4 report the average value for students in the control group. Columns 2 and 5 report the local average treatment effect (LATE) estimates. They are obtained from a regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report the cluster-robust p -value of the estimated treatment effect and, in square brackets, the p -value (q -value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage q -values introduced in Benjamini et al. (2006) and described in Anderson (2008). The q -values associated with the treatment effect estimates on “Grade 11: Science track” (Panel A) and “Major: STEM” (Panel B) are adjusted for multiple testing across the study’s nine main outcomes of interest, separately by grade level and gender (see Appendix D for details). The q -values associated with the treatment effect estimates for the different STEM tracks (Panel A) or the different STEM majors (Panel B) are adjusted for multiple testing across these different STEM tracks or majors, separately by grade level and gender. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

Table 6 – Heterogeneous Treatment Effects on Selective and Male-Dominated STEM Enrollment for Girls in Grade 12: Estimates based on Machine Learning Methods

Panel A. Best linear predictor (BLP) of the CATE $s_0(Z)$ given the ML proxy $S(Z)$				
Parameters:	ATE (β_1)	HET (β_2)	Best ML method	
Undergraduate major: selective STEM	0.038	0.762	Elastic Net	
p -value	[0.027]	[0.031]		
Undergraduate major: male-dominated STEM	0.036	0.088	Linear model	
p -value	[0.064]	[0.731]		
Panel B. Sorted group average treatment effects (GATEs): 20% most and least affected students				
Heterogeneity group:	20% least affected	20% most affected	Difference most–least	Best ML method
Undergraduate major: selective STEM	−0.004	0.139	0.149	Elastic Net
p -value	[1.000]	[0.014]	[0.026]	
Undergraduate major: male-dominated STEM	0.026	0.061	0.038	Elastic Net
p -value	[1.000]	[0.464]	[1.000]	
Panel C. Average characteristics of the 20% most and least affected students (CLAN)				
Heterogeneity group:	20% least affected	20% most affected	Difference most–least	p -value (upper bound)
Enrollment in selective STEM major				
<i>Student characteristics</i>				
Baccalauréat percentile rank in math	17.62	81.39	62.85	0.000
Baccalauréat percentile rank in French	41.45	73.44	32.74	0.000
High SES	0.344	0.637	0.302	0.000
<i>Role model characteristics</i>				
Professional	0.494	0.638	0.148	0.001
Participated in the program the year before	0.141	0.233	0.093	0.015
Non-French	0.133	0.183	0.051	0.228
Has children	0.503	0.417	−0.095	0.064
Age	33.09	32.97	−0.11	1.000
Holds/prepares for a Ph.D.	0.692	0.606	−0.080	0.111
Field: math, physics, engineering	0.316	0.226	−0.099	0.021
Field: earth and life sciences	0.618	0.602	−0.004	1.000
Enrollment in male-dominated major				
<i>Student characteristics</i>				
Baccalauréat percentile rank in math	19.88	79.02	59.45	0.000
Baccalauréat percentile rank in French	41.22	72.10	31.10	0.000
High SES	0.335	0.628	0.296	0.000
<i>Role model characteristics</i>				
Professional	0.530	0.606	0.078	0.170
Participated in the program the year before	0.142	0.240	0.091	0.021
Non-French	0.153	0.164	0.004	1.000
Has children	0.539	0.418	−0.126	0.010
Age	33.15	32.95	−0.17	1.000
Holds/prepares for a Ph.D.	0.705	0.601	−0.103	0.043
Field: math, physics, engineering	0.298	0.237	−0.065	0.186
Field: earth and life sciences	0.657	0.585	−0.075	0.170

Notes: This table reports heterogeneous treatment effects of the program on the undergraduate enrollment outcomes of girls in Grade 12, using the methods developed by Chernozhukov et al. (2018). For each outcome, the conditional average treatment effect (CATE) of role model interventions, $s_0(Z)$, is predicted using five alternative ML methods: Elastic Net, Random Forest, Linear Model, Boosting, and Neural Network. The covariates Z that are used to predict the CATE consist of three indicators for the educational districts of Paris, Créteil, and Versailles, four indicators for students’ socioeconomic background (high, medium-high, medium-low, and low), their age, their overall percentile rank in the *Baccalauréat* exam, their percentile ranks in the French and math tests of the exam, and a vector of 56 role model fixed effects. For each outcome, Panel A reports the parameter estimates and p -values (in square brackets) of the Best Linear Predictor (BLP) of the CATE using the best ML method. The coefficients β_1 and β_2 correspond to the average treatment effect (ATE) and heterogeneity loading (HET) parameters in the BLP, respectively. Panel B reports the Sorted Group Average Treatment Effects (GATEs), i.e., the average treatment effects among students in the top and bottom quintiles of the heterogeneous effects induced by the ML proxy predictor $S(Z)$, using the best ML method. Panel C performs a Classification Analysis (CLAN) by comparing the average characteristics of the 20 percent most and least affected students defined in terms of the ML proxy predictor. The parameter estimates and p -values are computed as medians over 100 splits, with nominal levels adjusted to account for the splitting uncertainty. This adjustment implies that the reported p -values should be interpreted as upper bounds for the actual p -values. Further details on the methods are provided in Appendix L.

Table 7 – Heterogeneous Treatment Effects on Student Perceptions: Average Characteristics of the Most and Least Affected Girls in Grade 12

	20% least affected (1)	20% most affected (2)	Difference most–least (3)	<i>p</i> -value (upper bound) (4)
<i>Positive perceptions of science-related careers (index)</i>				
Mean Baccalauréat percentile rank in math	26.62	73.29	46.85	0.000
Class visited by professional	0.483	0.675	0.192	0.000
<i>More men in science-related jobs</i>				
Mean Baccalauréat percentile rank in math	74.87	25.00	−51.03	0.000
Class visited by professional	0.614	0.511	−0.112	0.031
<i>Equal gender aptitude for math (index)</i>				
Mean Baccalauréat percentile rank in math	42.77	50.58	7.89	0.003
Class visited by professional	0.622	0.563	−0.058	0.403
<i>Women don't really like science</i>				
Mean Baccalauréat percentile rank in math	44.47	50.57	5.07	0.090
Class visited by professional	0.592	0.540	−0.035	0.908
<i>Women face discrimination in science-related jobs</i>				
Mean Baccalauréat percentile rank in math	52.15	42.79	−8.81	0.001
Class visited by professional	0.568	0.570	0.011	1.000
<i>Taste for science subjects (index)</i>				
Mean Baccalauréat percentile rank in math	41.36	54.71	13.63	0.000
Class visited by professional	0.436	0.678	0.227	0.000
<i>Math self-concept (index)</i>				
Mean Baccalauréat percentile rank in math	52.22	42.10	−10.65	0.000
Class visited by professional	0.512	0.582	0.071	0.240
<i>Science-related career aspirations (index)</i>				
Mean Baccalauréat percentile rank in math	44.70	47.78	2.36	0.712
Class visited by professional	0.375	0.762	0.389	0.000

Notes: This table reports the average characteristics of Grade 12 girls in the top and bottom quintile of predicted treatment effects on student perceptions, using the methods developed by Chernozhukov et al. (2018). For each outcome, the conditional average treatment effect (CATE) of role model interventions, $s_0(Z)$, is predicted using five alternative ML methods: Elastic Net, Random Forest, Linear Model, Boosting, and Neural Network. The covariates Z that are used to predict the CATE consist of three indicators for the educational districts of Paris, Créteil, and Versailles, four indicators for students' socioeconomic background (high, medium-high, medium-low, and low), their age, their overall percentile rank in the *Baccalauréat* exam, their percentile ranks in the French and math tests of the exam, and a vector of 56 role model fixed effects. For each outcome, the table compares the average characteristics of the students in the top and bottom quintile of treatment effects, as predicted by the best ML proxy predictor based on the Group average treatment effects (GATEs) targeting of the CATE. The characteristics reported in this table are the students' average percentile rank in math in the *Baccalauréat* exams and the share exposed to a role model with a professional rather a research background. The parameter estimates and *p*-values are computed as medians over 100 splits, with nominal levels adjusted to account for the splitting uncertainty. This adjustment implies that the reported *p*-values should be interpreted as upper bounds for the actual *p*-values. The average treatment effects among the 20 percent most and least affected students can be found in Panel B of Appendix Table L1. Further details on the methods are provided in Appendix L.

Table 8 – Correlation between Conditional Average Treatment Effects (CATEs) for Girls in Grade 12

	Bivariate correlation with the CATE on enrollment in a selective STEM program	
	Estimate (1)	95% confidence interval (2)
<i>Conditional average treatment effect (CATE) on:</i>		
Positive perception of science-related careers (index)	0.96	[0.21, 5.30]
More men in science-related jobs	−0.68	[−3.23, −0.01]
Equal gender aptitude for math (index)	0.19	[−1.24, 2.05]
Women don’t really like science	0.21	[−1.43, 3.23]
Women face discrimination in science-related jobs	−0.34	[−2.22, 0.56]
Taste for science subjects (index)	0.71	[0.04, 3.96]
Math self-concept (index)	−0.07	[−1.84, 1.40]
Science-related career aspirations (index)	0.36	[−0.51, 2.01]

Notes: This table reports, for girls in Grade 12, estimates of the bivariate correlation $\rho_{A,B|Z}$ between the Conditional Average Treatment Effect (CATE) on enrollment in a selective STEM program, denoted by $s_0^B(Z)$, and the CATE on each of the potential channels listed in the table, denoted by $s_0^A(Z)$. The proxy predictor of the CATE on selective STEM enrollment, denoted by $S^B(Z)$, is estimated using the Elastic Net method, as it has the best performance based on the Best Linear Predictor (BLP) targeting of the CATE for this outcome. The proxy predictor of the CATE on the potential mediator Y^A , denoted by $S^A(Z)$, is estimated using the ML method that has the best performance based on the BLP targeting of the CATE on the corresponding outcome. An indication of the quality of these predictions is provided by the heterogeneity loading (HET) parameter of the BLP (see Appendix Table L1, Panel A). For each random split of the data, the correlation coefficient $\rho_{A,B|Z}$ is estimated as $\hat{\rho}_{A,B|Z} = \text{Sign}(\hat{\beta}_2^{A|B})(\hat{\beta}_2^{A|B}\hat{\beta}_2^{B|A})^{\frac{1}{2}}/(\hat{\beta}_2^{A|A})^{\frac{1}{2}}(\hat{\beta}_2^{B|B})^{\frac{1}{2}}$, where $\hat{\beta}_2^{k|l}$ is the estimated heterogeneity loading parameter of the BLP of $s_0^k(Z)$ based on $S^l(Z)$ (with $k, l \in \{A, B\}$), using the methods in Chernozhukov et al. (2018). The covariates Z that are used to predict the CATEs consist of three indicators for the educational districts of Paris, Créteil, and Versailles, four indicators for students’ socioeconomic background (high, medium-high, medium-low, and low), their age, their overall percentile rank in the *Baccalauréat* exam, their percentile ranks in the French and math tests of the exam, and a vector of 56 role model fixed effects. For each pair of outcomes, columns 1 and 2 report the estimated correlation between the CATEs and its 95 percent confidence interval, respectively. Estimates and confidence intervals are computed as medians over the first 100 random data splits for which $\hat{\rho}_{A,B|Z}$ can be computed. For each data split, the confidence intervals are obtained using a clustered bootstrap procedure. The nominal level of the median of confidence intervals is adjusted to account for the splitting uncertainty, using the method of Chernozhukov et al. (2018). This adjustment implies that the reported confidence intervals should be interpreted as lower and upper bounds for the true lower and upper limits of the confidence intervals. Further details on the methods are provided in Appendix L.

(For Online Publication)

Appendix to

Do female Role Models Reduce the Gender Gap in Science? Evidence from French High Schools

Thomas Breda

Julien Grenet

Marion Monnet

Clémentine Van Effenterre

List of Appendices

A Gender Pay Gap Among College Graduates in France	A-2
B Program Details	A-8
C Student-Level Administrative Data	A-10
D Construction of Synthetic Indices and Multiple Hypotheses Testing	A-12
E Summary Statistics and Balancing Tests	A-14
F Effects of Role Model Interventions: Additional Results	A-21
G Robustness Checks	A-28
H Randomization Inference	A-30
I Information, Persistence, Timing: Additional Results	A-32
J Spillover Effects	A-37
K Heterogeneous Treatment Effects: Subgroup Analysis	A-47
L Heterogeneous Treatment Effects: Machine Learning Methods	A-50
Appendix References	A-60

A Gender Pay Gap Among College Graduates in France

This appendix provides descriptive evidence on the entry-level gender pay gap among French college graduates holding a master’s degree and analyzes the contribution of gender segregation in college majors to this gap. The objective of this analysis is to better understand whether the effects of the role model interventions on female students’ choice of study can be expected to reduce the gender pay gap. Section A.1 describes the data sources, while Section A.2 discusses the empirical results.

A.1 Data

Unfortunately, we cannot rely exclusively on administrative data to provide empirical evidence on the gender pay gap by field of study in France, as it is currently not possible to link administrative data on students enrolled in higher education with administrative data on wages and income tax returns. Instead, our analysis is based on the combination of aggregate statistics on student enrollment by college major and gender with survey information on the starting wages of recent cohorts of college graduates.

Data sources. In France, gender segregation and gender pay gaps by college major can be analyzed for the population of college graduates who obtained their master’s degree (or equivalent) in 2015 or 2016. For this purpose, we combine several administrative and survey data sources.

SISE Résultats 2015. This individual-level administrative dataset covers all students enrolled in public universities during the academic year 2015/16 and provides detailed information on each student’s degree program and field of study.

Enquête Professionnelle des Diplômés de Master 2015 (EPDM). This survey was conducted in December 2017 by the Ministry of Higher Education to collect information on the transition of master’s graduates to the labor market. The survey was targeted at students who obtained their master’s degree in 2015 and who entered the labor market within one year after graduation, with an overall response rate of 70 percent. As part of this survey, master’s graduates were asked to report their annual earnings 18 months after graduation. Our analyses are based on the survey’s public use files, which provide aggregate statistics by gender and college major.^{A.1}

Enquête sur l’Insertion des Diplômés des Grandes Écoles 2018 (EIDGE). This survey was conducted in 2018 by the Conférence des Grandes Écoles (CGE), a not-for-profit association representing French elite graduate schools. The *Grandes Écoles*, which award a diploma equivalent to a master’s degree, recruit their students through highly competitive national exams taking place at the end of two-year undergraduate selective STEM and non-STEM preparatory courses (*Classes Préparatoires aux Grandes Écoles* or CPGE). The survey was targeted at students who graduated between 2015 and 2017 from one of the 184 *Grandes Écoles* that were members of the CGE in 2018, with an overall response rate of 48 percent. Our analyses are based on the aggregate statistics published by the CGE separately by gender and by type of *Grande École* (i.e., engineering schools, business schools, and other schools).^{A.2} We only consider students who graduated from a *Grande École* in 2016, since annual earnings 24 months after graduation are only available for this cohort.

^{A.1}https://data.enseignementsup-recherche.gouv.fr/explore/dataset/fr-esr-insertion_professionnelle-master_donnees_nationales/information/ (accessed on August 2, 2019).

^{A.2}<https://www.cge.asso.fr/themencode-pdf-viewer/?file=https://www.cge.asso.fr/wp-content/uploads/2018/06/2018-06-19-Rapport-2018.pdf> (accessed on August 2, 2019).

Grouping of college majors. The above data sources can be combined to compute the number of female and male master’s students who graduated from university in 2015 or from a *Grande École* in 2016, separately by college major.

The Ministry of Higher Education’s official classification comprises 54 college majors. For the purpose of our analysis, we group these college majors into the following broad categories:

- Non-STEM majors (35 in total): this category includes master’s degree programs in law, economics, management, humanities, psychology, social sciences, medicine, pharmacy, sports studies as well as degrees from non-STEM *Grande Écoles* (e.g., business schools, schools of journalism, schools of architecture).
- STEM majors (19 in total): this category includes master’s degree programs in STEM fields as well as degrees from engineering schools (*Grande Écoles d’ingénieurs*).
- Among STEM majors, we distinguish between engineering schools (all of which are selective and are classified as a single major) and non-selective STEM master’s degrees at university (18 in total).
- Among non-selective STEM majors, we further distinguish between male-dominated majors (16 in total) and female-dominated majors (2 in total: chemistry and earth and life sciences), based on whether the share of female students among master’s graduates in the corresponding field of study is below or above 50 percent. This distinction does not apply to selective STEM majors, since almost all engineering schools are male-dominated.

Earnings information. The EPDM and EIDGE surveys provide information on graduates’ average median gross salary (*salaires brut annuel médian*) separately by gender and college major. Starting wages are measured 18 months after graduation for master’s graduates and 24 months after graduation for *Grandes Écoles* graduates. Note that since we do not have access to the individual-level survey data, median earnings by broad categories of college majors can only be approximated as the average of the median earnings in each of the majors that form these broad categories.

A.2 College Majors and the Gender Pay Gap

Combining the above data sources, we provide descriptive evidence on the median starting wages of female and male graduates across the broad categories of college majors. We then analyze the contribution of gender segregation in college majors to the overall entry-level gender pay gap.

Gender composition of STEM and non-STEM majors. The first three columns of Table A1 show the distribution of master’s-level graduates across the broad categories of college majors defined above, along with the share of female graduates in each category. The summary statistics indicate that while female students represent 52 percent of master’s level graduates, they are strongly underrepresented in STEM majors (34 percent). Female underrepresentation is more pronounced in selective (male-dominated) STEM majors (female share: 30 percent) than in non-selective STEM majors (female share: 40 percent). Among non-selective STEM majors, female students represent only 29 percent of graduates in male-dominated fields such as mathematics, physics, or computer science, compared to 60 percent of graduates in female-dominated fields such as chemistry and earth and life sciences.

Starting wages of STEM and non-STEM graduates. The comparison of starting wages by broad college major category confirms that female graduates tend to be overrepresented in lower-paying majors (see columns 3–5 of Table A1). Female graduates holding a STEM

degree have a median starting wage of 29,984 euros, which is 7.4 percent higher than the median starting wage of female graduates holding a non-STEM degree (27,913 euros). Strikingly, the wage premium for female graduates in STEM appears to be almost entirely driven by selective (male-dominated) STEM degrees (16.4 percent). By contrast, the wage premium attached to non-selective STEM degrees is close to zero (−0.5 percent). The low apparent return to non-selective STEM degrees masks substantially different returns between male-dominated and female-dominated majors: while the wage premium attached to male-dominated non-selective STEM majors is of 4.2 percent for female graduates compared to non-STEM majors, a wage penalty of 4.7 percent is attached to female-dominated non-selective STEM majors.

Female underrepresentation in STEM: contribution to the gender pay gap. The last three columns of Table A1 indicate that across all categories of programs, male graduates earn a median annual starting wage of 32,122 euros, compared to 28,411 euros for female graduates. This amounts to an overall gender pay gap of 3,711 euros per year, or 11.6 percent of male pay.

Although the overrepresentation of female graduates in lower-paying non-STEM and female-dominated STEM majors is a likely contributor to the overall gender pay gap, it is clearly not the sole cause, as gender differences in median earnings are observed within each broad category of college majors. Interestingly, however, the gender wage gap is lower in each category of STEM majors than in non-STEM majors. This finding is consistent with similar evidence for the U.S. (Beede et al., 2011).

To shed light on the contribution of gender segregation in fields of study to the overall entry-level gender pay gap, we adopt a method similar to that used by McDonald and Thornton (2007) in estimating what the overall female-male starting wage gap would be if female graduates had the same distribution of college majors as male graduates.

Since our interest is in measuring the specific contribution of the different dimensions of female underrepresentation in STEM majors (STEM vs. non-STEM, selective vs. non-selective STEM, male-dominated vs. female-dominated non-selective STEM), we construct counterfactual wage gaps by considering increasingly disaggregated groups of majors.

We start by estimating the counterfactual wage gap if female graduates had the same distribution of STEM vs. non-STEM majors as male graduates, while keeping fixed females' marginal distribution of majors within each of these two broad categories. Put differently, we apply female median earnings in STEM vs. non-STEM degrees to the male distribution of graduates in both categories of majors to recalculate the overall gender pay gap. This counterfactual wage gap, which we denote by $\tilde{\Delta}_w$, is constructed as follows:

$$\tilde{\Delta}_w = 1 - \frac{(\bar{w}_s^f N_s^m + \bar{w}_{ns}^f N_{ns}^m)}{(\bar{w}_s^m N_s^m + \bar{w}_{ns}^m N_{ns}^m)},$$

where \bar{w}_k^g and N_k^g denote the median earnings and the number of graduates of gender g (m : males; f : females) in college major category k (s : STEM; ns : non-STEM), respectively. The contribution of female underrepresentation in STEM programs to the gender pay gap is then measured as $\Delta_w - \tilde{\Delta}_w$, where Δ_w denotes the observed overall pay gap between male and female graduates.

To measure the contribution of gender segregation between selective and non-selective STEM majors, we construct a second counterfactual wage gap in a similar manner, except that college majors are now grouped into three categories: non-STEM, selective STEM, and non-selective STEM. To measure the contribution of gender segregation between male-dominated and female-dominated STEM majors, we repeat this exercise after grouping college majors into four categories: non-STEM, selective STEM, non-selective male-dominated STEM, and

non-selective female-dominated STEM. The contribution of gender segregation between majors within both male- and female-dominated non-selective STEM is measured by ungrouping all STEM majors. Finally, we ungroup all non-STEM majors to evaluate the contribution of gender segregation between non-STEM majors. The corresponding counterfactual measures what the overall gender gap would be if women had the same distribution as men across all 54 STEM and non-STEM college majors.

Results. The results of this decomposition exercise are shown in Table A2 along with the observed gender pay gap. The contributions of gender segregation between the different categories of college majors to the gender pay gap are reported in column 1 and are expressed as percentages of the total in column 2. We find that the gender imbalances across all college majors “explain” 40 percent of the gender pay gap among college graduates. Two-thirds of this explained part (26.5 percent of the total wage gap) can be attributed to the unequal representation of female and male graduates in STEM vs. non-STEM majors, on the one hand, and between the different majors within STEM, on the other hand. The remain third of the explained part of the gap (13.4 percent of the total) is due to gender segregation between non-STEM majors, the lowest-paying majors (humanities) being typically more female-dominated (77 percent) than the highest-paying ones (law and economics, where the female share is 59 percent).

The 26.5 percent STEM-related gender pay gap can be decomposed as follows. Increasing the share of female graduates holding a STEM degree to that of males without changing females’ marginal distribution of STEM majors is associated with a 14.0 percent reduction in the gender pay gap. In line with the evidence from Table A1, further reassigning female graduates from non-selective STEM majors to (male-dominated) selective STEM majors in order to match the relative shares of selective and non-selective STEM majors among male graduates would reduce the gender gap by an additional 6.5 percent from the baseline. Finally, reassigning female graduates from non-selective female-dominated STEM majors to non-selective male-dominated STEM majors would trigger an extra 4.3 percent reduction in the gender pay gap, while further reassigning female students between majors within male- and female-dominated programs would result in an extra 1.8 percent reduction from the baseline.

Altogether, these findings suggest that the underrepresentation of female students in STEM majors accounts for approximately 25 percent of the entry-level gender pay gap among college graduates in France. Almost half of this STEM-related gender pay gap can be attributed to the fact that within STEM majors, female graduates are relatively less likely than males to be enrolled in those with the largest wage premium, i.e., the selective and male-dominated STEM majors.

Table A1 – Starting Wage Among College Graduates Holding a Master’s Degree or Equivalent, Classes of 2015/16

	Graduates: classes of 2015/16			Wage 18/24 months after graduation (survey)				
	Number of graduates	% of total	Female share (%)	Female graduates		Male graduates		Gender pay gap (%)
				Median wage (euros)	Relative Median wage (non-STEM majors: 100)	Median wage (euros)	Relative Median wage (non-STEM majors: 100)	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
All majors (54)	166,600	100.0	51.5	28,411	-	32,122	-	11.6
Non-STEM majors (35)	106,997	64.2	61.1	27,913	100.0	31,302	100.0	10.8
STEM majors (19)	59,603	35.8	34.3	29,984	107.4	32,972	105.3	9.1
<i>of which:</i>								
Selective (male-dominated) STEM majors (Engineering schools)	31,463	18.9	29.7	32,500	116.4	34,800	111.2	6.6
Non-Selective STEM majors (18)	28,140	16.9	39.6	27,767	99.5	30,530	97.5	9.1
<i>of which:</i>								
Male-dominated majors (16)	18,874	11.3	29.4	29,077	104.2	31,371	100.2	7.3
Female-dominated majors (2)	9,266	5.6	60.3	26,596	95.3	27,581	88.1	3.6

Notes: This table reports summary statistics on gender segregation and gender pay gaps for the population of college graduates who obtained their master’s degree (or equivalent) in 2015 or 2016. The 54 college majors are grouped into two broad categories: non-STEM majors (master’s degrees in economics, management, humanities, psychology, social sciences, sports studies, medicine, pharmacy, and non-STEM *Grandes Écoles* such as business schools or schools of journalism) and STEM majors (master’s degrees in STEM fields and degrees from engineering schools); STEM majors are further broken down between selective (engineering schools) and non-selective majors (master’s degree at university); among non-selective majors, we distinguish between male-dominated and female-dominated majors, based on whether the share of female graduates in the corresponding field of study is below or above 50 percent. Column 1 shows the number of graduates per broad category of college majors using the administrative dataset SISE 2015/16 (for university graduates who obtained their master’s degree in 2016) and the EIDGE survey (for students who graduated from *Grandes Écoles* in 2016). Median gross annual wages (columns 4 and 6) are computed from aggregate statistics by gender and college major from the EPDM and EIDGE surveys. Entry-level wages are measured 18 months after graduation for master’s graduates and 24 months after graduation for *Grandes Écoles* graduates. Median wages by broad categories of college majors are approximated as the average of the median wages in each of the majors that form these broad categories.

Sources: Columns 1–3: SISE 2015/16 and Enquête sur l’Insertion des Diplômés des Grandes Écoles 2018 (EIDGE); columns 4–8: Enquête Professionnelle des Diplômés de Master 2015 (EPDM) and EIDGE.

Table A2 – Contribution of Gender Segregation in College Majors to the Entry-Level Gender Wage Gap Among College Graduates, Classes of 2015/16

	Gender pay gap (relative to male pay) (1)	Share of the gender wage gap (2)
Total wage gap	0.116	100.0%
<i>Contribution of gender segregation in college majors to the wage gap:</i>		
Explained by unequal gender distribution between majors	0.046	40.0%
<i>of which:</i>		
between STEM/non-STEM majors and between majors within STEM	0.031	26.5%
<i>of which:</i>		
between STEM and non-STEM majors	0.016	14.0%
between selective and non-selective STEM majors	0.007	6.5%
between male- and female-dominated non-selective STEM majors	0.005	4.3%
between majors within male- and female-dominated non-selective STEM	0.002	1.8%
between majors within non-STEM	0.016	13.4%
Unexplained by unequal gender distribution between majors	0.069	60.0%

Notes: This table provides a decomposition of the total entry-level wage gap between male and female college graduates who obtained their master's degree or equivalent in 2015 (university graduates) or in 2016 (*Grandes Écoles* graduates). Entry-level wages are measured as median annual gross wages by gender and college major, 18 months after graduation for master's graduates, and 24 months after graduation for *Grandes Écoles* graduates. To measure the contribution of the unequal gender representation across college majors, counterfactual wage gaps are constructed using increasingly disaggregated groups of college majors. The contribution of gender segregation between STEM and non-STEM majors is measured as the observed gender wage gap minus the counterfactual wage gap that would be observed if female graduates had the same distribution of STEM and non-STEM majors as male graduates, while keeping fixed females' marginal distribution of majors within each of these two broad categories. The contribution of gender segregation between selective and non-selective STEM majors is estimated in a similar manner, except that the counterfactual gender wage gap is estimated by reassigning female graduates from non-selective STEM majors to selective STEM majors to match the relative shares of selective and non-selective STEM majors among male graduates. The other components of the gender wage gap are measured by sequentially ungrouping college majors to compute counterfactual gender wage gaps. The contributions of gender segregation between the different categories of college majors to the gender wage gap are shown in column 1 and are expressed as percentages of the total in column 2.

Sources: See notes of Table A1.

B Program Details

(a) First Video: “Jobs in Science: Beliefs or Reality?”



(b) Second Video: “Are we All Equal in Science?”

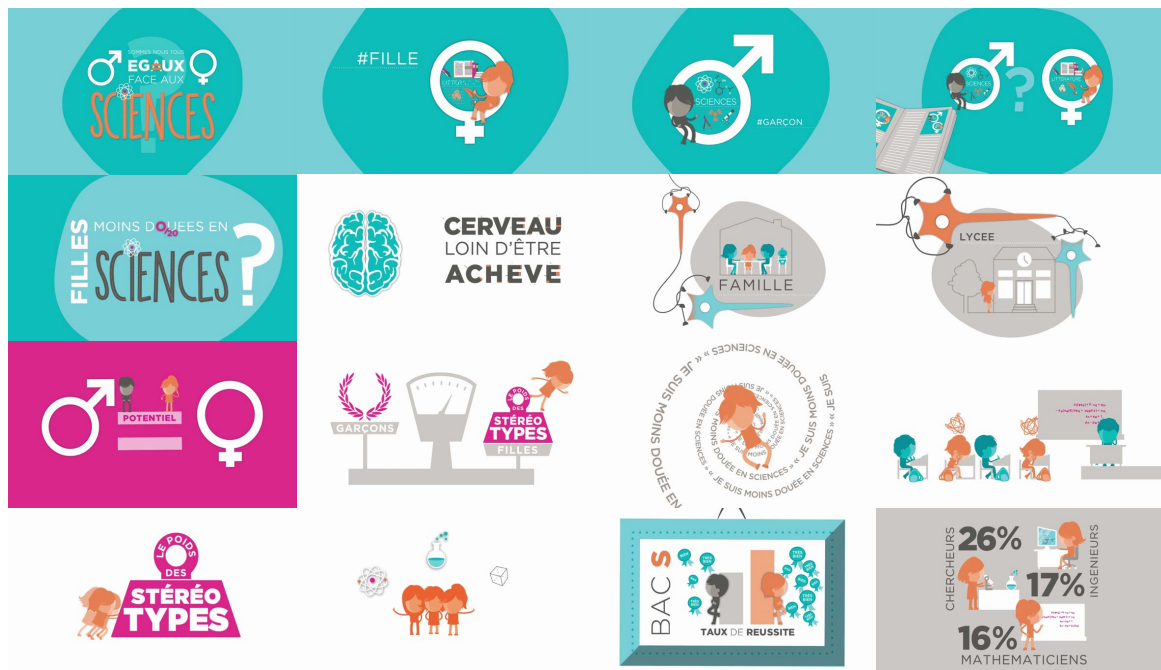


Figure B1 – Screenshots of the Two Videos Shown During the Role Model Interventions



Figure B2 – Screenshots of the Slides Provided to the Role Models to Describe their own Experience

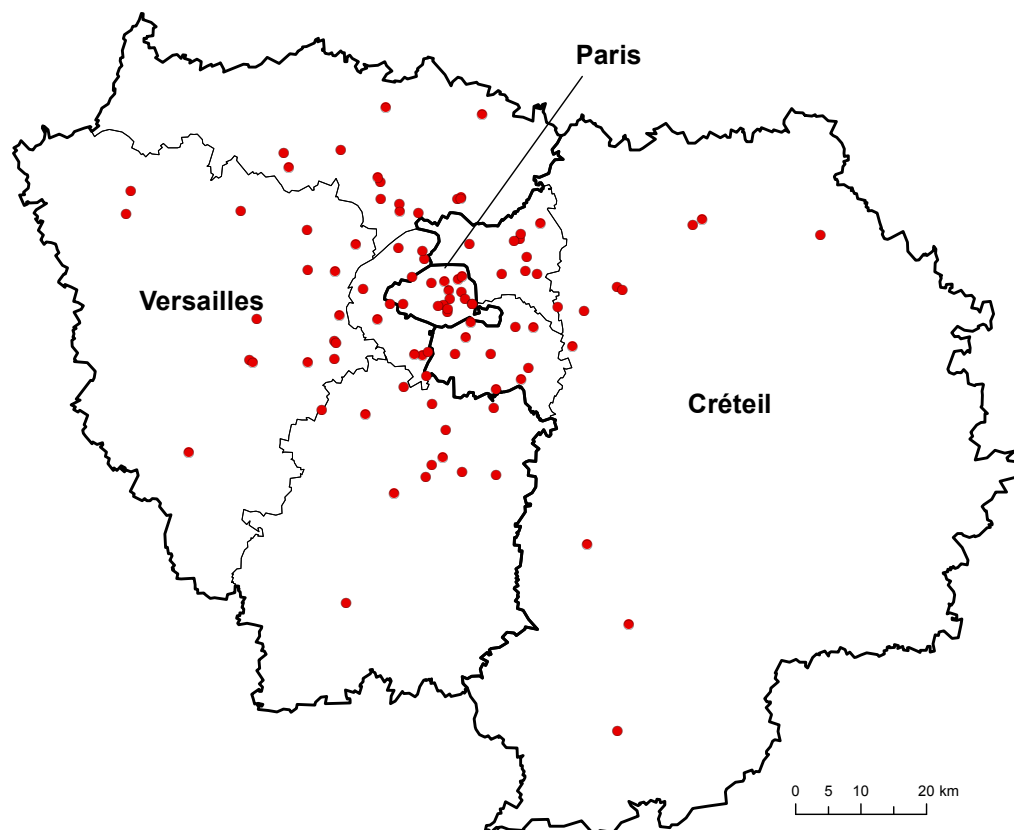


Figure B3 – Participating High Schools

Notes: The thick lines represent the boundaries of the three education districts (*académies*) of the Paris region (Paris, Créteil and Versailles). The solid circles show the location of the 98 high schools that participated in the program evaluation.

C Student-Level Administrative Data

This appendix describes the administrative data that we use to complement the information from the student survey (Section C.1) and provides details about the classification of STEM undergraduate programs (Section C.2).

C.1 Data Sources

For the purpose of the empirical analysis, we matched the data from our post-intervention student survey with three administrative datasets. These data were linked using an encrypted version of the French national student identifier (*Identifiant National Élève*).

High school enrollment data. Students' socio-demographic characteristics and enrollment status are obtained from the *Bases Éléves Académiques* (BEA) for academic years 2012/13 to 2016/17. These comprehensive administrative registers, which were provided by the three education districts of the Paris region (Paris, Créteil, and Versailles), cover the universe of students enrolled in the public and private high schools operating in the three districts. They also cover students enrolled in selective undergraduate programs, i.e., *Classes préparatoires aux Grandes Écoles* (CPGE) and *Sections de technicien supérieur* (STS), as these programs are located in high schools. The BEA data provide basic information on students' demographics (gender, date and country of birth, number of siblings), their parents' two-digit occupation, and detailed information on their enrollment status (school and class attended, elective courses taken). Students' socioeconomic status (SES) is measured using the French Ministry of Education's official classification, which uses the occupation of the child's legal guardian to define four groups of SES: high (company managers, executives, liberal professions, engineers, intellectual occupations, arts professions), medium-high (technicians and associate professionals), medium-low (farmers, craft and trades workers, service and sales workers), and low (manual workers and persons without employment).

University enrollment data. To track Grade 12 (science track) students' enrollment outcomes in non-selective undergraduate programs (*Licence*), we use a separate administrative data source, the *Système d'Information sur le Suivi de l'Étudiant* (SISE), which is managed by the Statistical Office of the French Ministry of Higher Education (Sous-Direction des Systèmes d'Information et des Études Statistiques, MESRI-SIES). This dataset, which covers the academic years 2012/13 to 2016/17, records all students enrolled in the French higher education system outside of CPGE and STS, except for the small fraction of students enrolled in undergraduate programs leading to paramedical and social care qualifications.

Data on student performance. The third dataset, the *Organisation des Concours et Examens Académiques et Nationaux* (OCEAN), contains students' individual exam results for the *Diplôme national du brevet* (DNB), which middle school students take at the end Grade 9, and for the *Baccalauréat*, which high school students take at the end of Grade 12. Access to this dataset, which covers the exams years 2010 to 2016, was provided by the Statistical Office of the French Ministry of Education (Direction de l'Évaluation, de la Prospective et de la Performance, MENJ-DEPP).

C.2 Classification of STEM Undergraduate Programs

The enrollment status of Grade 12 (science track) students in the year following the intervention, i.e., 2016/17, is measured by combining the information from the BEA and SISE datasets. For the

purpose of our analysis, we use two alternative classifications of STEM undergraduate programs, based on whether they are (i) selective or non-selective, and (ii) male- or female-dominated.

Selective vs. non-selective STEM programs.

- *Selective STEM*: This category includes all CPGE programs with a specialization in STEM, i.e., mathematics, physics and engineering science (MPSI), physics, chemistry and engineering science (PCSI), biology, chemistry, physics and earth sciences (BCPST), and physics, technology, and engineering science (PTSI). It also includes a small number of selective programs in engineering schools that recruit their students directly after high school graduation, as well as selective technical/vocational undergraduate programs (STS) that specialize in STEM fields.
- *Non-selective STEM*: This category includes non-selective university bachelor's degree programs (*Licence*) that specialize in STEM fields: math, physics, chemistry, earth and life sciences, and computer science. Undergraduate programs in medicine and pharmacy are not included in this category.

Male- vs. female-dominated STEM programs.

- *Male-dominated STEM*: We consider as male-dominated STEM programs those in which the share of female students is less than 50 percent. This category includes the selective programs (CPGE and STS) and non-selective programs (*Licence*) that specialize in mathematics, physics, chemistry, computer science, and engineering.
- *Female-dominated STEM*: This category includes both selective (CPGE and STS) and non-selective programs (*Licence*) that specialize in earth and life sciences.

If a student is enrolled in multiple higher education programs, we only consider the most selective among these programs, with CPGE taking precedence over STS, and STS taking precedence over university undergraduate degree programs.

D Construction of Synthetic Indices and Multiple Hypotheses Testing

This appendix discusses the construction of the synthetic indices that we use to measure the effects of role model interventions on students’ perceptions (Section D.1) and provides further details on the adjustment of p -values to correct for multiple hypotheses testing (Section D.2).

D.1 Construction of Synthetic Indices

The student survey questionnaire aimed at measuring the effects of role model interventions on students’ perceptions and self-concept along five dimensions: (i) general perceptions of science-related careers, (ii) perceptions of gender roles in science, (iii) taste for science subjects, (iv) math self-concept, and (v) science-related career aspirations.

We use the survey items listed below to construct synthetic indices for each of these five dimensions. When responses are measured on a Likert scale, i.e., when respondents specify their level of agreement or disagreement with a statement on a symmetric agree-disagree scale, the item responses are recoded so that higher values correspond to less stereotypical or negative perceptions (see details below). We then take the average of each student’s responses to the different questions.^{A.3} We checked that the indices yield similar results if item responses are converted to binary variables before taking the average across items. Finally, to facilitate interpretation, we normalize each index to have a mean of zero and a standard deviation of one in the control group.

Below is the list of the individual items that are included in each the five synthetic indices. Unless otherwise specified, these items use a four-point Likert response scale such that 1=Strongly agree, 2=Agree, 3=Disagree, and 4=Strongly disagree. Items marked with a * have been recoded such that a value of 1 means “Strongly disagree” and 4 means “Strongly agree”.

1. *Positive perceptions of science-related careers* (5 items): “Science-related jobs require long years of study”; “Science-related jobs are monotonous”; “Science-related jobs are rather solitary”; “Science-related jobs pay higher wages*”; “It is difficult have a fulfilling family life when working as a scientist”.
2. *Equal gender aptitude for math* (2 items): “Women and men are born with different brains”; “Men are more gifted than women in mathematics”.
3. *Taste for science subject* (4 items): Enjoys math (on a scale from 0 “not at all” to 10 “very much”); Enjoys physics and chemistry (on a scale from 0 to 10); Enjoys earth and life sciences (on a scale from 0 to 10); “I like science in general*”.
4. *Math self-concept* (4 items): Self-assessed performance in math (very weak/weak/average/good/very good); “I feel lost when I try to solve a math problem”; “I often worry that I will struggle in math class”; “If I make enough effort, I can do well in science subjects”.
5. *Science-related career aspirations* (4 items): “Some jobs in science are interesting*”; “ I could see myself working in a science-related job later in life*”; Interested in at least one of six STEM job out of a list of ten STEM and non-STEM occupations^{A.4} (0/1 variable); “Career and earnings prospects play an important role in my choice of study” (on a scale from 0 “not at all” to 10 “very much”).

^{A.3}This procedure is inspired from the KidIQoI test used in the psychological literature to measure children’s life satisfaction (Gayral-Tamih et al., 2005).

^{A.4}The STEM occupations in the list were: chemist, computer scientist, engineer, industrial designer, renewable energy technician, and researcher in biology. The non-STEM occupations were lawyer, pharmacist, physician, and psychologist.

D.2 Multiple Hypotheses Testing

Consistent with the recent applied literature, we systematically use the False Discovery Rate (FDR) control, which designates the expected proportion of all rejections that are type-I errors. Specifically, we use the sharpened two-stage q -values introduced in Benjamini et al. (2006) and described in Anderson (2008).

We study nine main outcomes throughout the paper: (i) enrollment in a STEM track (for Grade 10 students) or STEM major (for Grade 12 students); (ii) five synthetic indices capturing positive perceptions of science-related careers, equal gender aptitude for math, taste for science subjects, math self-concept, and science-related career aspirations (see Section D.1); and (iii) three variables capturing different facets of gender role in science that cannot be combined into a single index, which are based on the survey items asking students whether they agree or disagree with the statements “There are more men than women in science-related jobs”, “Women don’t really like science”, and “Women face discrimination in science-related jobs”. These nine outcomes are our primary outcomes of interest and we therefore systematically provide (along with standard p -values) p -values that are adjusted for multiple testing across them (q -values), separately by grade level and gender.

For each of the five synthetic indices described in the previous section, we report separate treatment effect estimates for the individual components of the index and provide standard p -values for the corresponding estimates along with p -values adjusted for multiple testing across the index components, separately by grade level and gender.

As we further split enrollment in STEM into different types of STEM tracks or majors (e.g., selective STEM, non-selective STEM, male-dominated STEM, and female-dominated STEM in Grade 12), we provide adjusted p -values for multiple testing across these different STEM tracks or majors, separately by grade level and gender. Given the importance of some of these specific STEM majors in our analyses, it could also be justified to consider them jointly with the primary outcomes described above. We have checked that, in practice, this alternative choice has little effect on the reported q -values.

Finally, treatment effects on other outcomes, such as the probabilities of being enrolled in a non-STEM major or of not being enrolled in an education program in the year following the classroom interventions, are also reported in the paper for the sake of completeness and clarity. Since these are not outcomes of direct interest in our study or are complements of other outcomes of interest, we do not consider them in the multiple testing corrections.

E Summary Statistics and Balancing Tests

Table E1 – Experimental Sample: Summary Statistics (School-Level)

	High schools operating in the Paris region (1)	Participating high schools (2)
Number of high schools	489	98
Share private	0.339	0.173
Education district: Paris	0.243	0.153
Education district: Créteil	0.348	0.296
Education district: Versailles	0.409	0.551
Number of students	644	924
Share of female students	0.524	0.526
Share of high SES students	0.423	0.391
Share of medium-high SES students	0.116	0.128
Share of medium-low SES students	0.243	0.239
Share of low SES students	0.218	0.241
Pass rate on Baccalauréat exam in 2015	0.913	0.910

Notes: This table compares the characteristics of high schools that participated in the program evaluation in 2015/16 to the characteristics of all general-track high schools operating in the Paris region. The summary statistics are computed from the *Bases Élèves académiques* of the three education districts of Paris, Créteil, and Versailles for the academic year 2015/16. The *Baccalauréat* pass rate is computed for students who were enrolled in Grade 12 in 2014/15, i.e., in the year before the intervention, and who took the exams in the general or technological tracks.

Table E2 – Experimental Sample: Summary Statistics (Student-Level)

	High schools operating in the Paris region (1)	Participating high schools	
		Classes selected for random assignment (2)	Classes not selected for random assignment (3)
Panel A. Grade 10			
Number of students	115,720	13,700	19,147
Number of classes	3,627	416	592
Female	0.525	0.529	0.525
Non-French	0.063	0.060	0.068
Age	15.14	15.13	15.14
High SES	0.403	0.381	0.361
Medium-high SES	0.118	0.128	0.127
Medium-low SES	0.239	0.241	0.248
Low SES	0.240	0.249	0.265
Number of siblings	1.44	1.49	1.50
Class size	32.22	33.25	32.48
DNB percentile rank in math	57.69	58.48	55.10
DNB percentile rank in French	57.23	57.85	55.75
Panel B. Grade 12 (science track)			
Number of students	38,582	5,751	5,623
Number of classes	1,267	185	179
Female	0.459	0.492	0.417
Age	17.11	17.12	17.10
Non-French	0.045	0.051	0.037
High SES	0.527	0.464	0.535
Medium-high SES	0.115	0.136	0.126
Medium-low SES	0.198	0.209	0.180
Low SES	0.160	0.192	0.160
Number of siblings	1.43	1.50	1.44
Class size	31.43	31.97	32.08
DNB percentile rank in math	76.25	74.06	76.20
DNB percentile rank in French	70.78	69.61	69.78

Notes: This table compares the characteristics of Grade 10 and Grade 12 (science track) students enrolled in the high schools that participated in the program evaluation to the characteristics of all Grade 10 and Grade 12 (science track) students enrolled in general-track high schools in the Paris region. In participating schools, the classes that were selected by principals for random assignment to treatment are compared to classes that were not selected. The summary statistics are computed from the *Bases Élèves académiques* of the three education districts of Paris, Créteil, and Versailles for the academic year 2015/16. French and math scores are from the exams of the *Diplôme national du brevet* (DNB) that middle school students take at the end of Grade 9.

Table E3 – Post-Intervention Role Model Survey: Summary Statistics

	Role model background			Difference (3)–(2) (4)	<i>p</i> -value of diff. (5)
	All	Profes- sionals	Resear- chers		
	(1)	(2)	(3)		
<i>A. Adults present during the intervention</i>					
Teacher was present	0.890	0.883	0.896	0.014	0.773
Teacher’s subject: science ^a	0.600	0.596	0.603	0.007	0.922
Teacher’s gender: female	0.551	0.533	0.565	0.032	0.653
Teacher showed interest	0.696	0.634	0.745	0.111	0.098
Other adult present beside teacher	0.348	0.392	0.315	–0.077	0.236
<i>B. General atmosphere during the intervention</i>					
Students were very interested	0.423	0.425	0.422	–0.004	0.963
Students were very engaged in the discussion	0.386	0.378	0.392	0.014	0.838
Students were inattentive	0.134	0.165	0.110	–0.055	0.259
Powerpoint worked well	0.963	0.938	0.982	0.045	0.172
Videos worked well	0.888	0.891	0.886	–0.004	0.940
Logistical problems	0.160	0.185	0.140	–0.044	0.487
Talk interrupted due to discipline problems	0.068	0.079	0.060	–0.018	0.652
<i>C. Topics addressed during the intervention</i>					
“Science is everywhere”	1.000	1.000	1.000	0.000	–
“Jobs in science are fulfilling”	0.990	1.000	0.982	–0.018	0.080
“Jobs in science are for girls too”	1.000	1.000	1.000	0.000	–
“Jobs in science pay well”	0.866	0.890	0.849	–0.040	0.516
Short videos	0.980	0.969	0.988	0.019	0.436
<i>D. Students’ responsiveness to topics addressed during the intervention</i>					
Very responsive to “science is everywhere”	0.430	0.378	0.470	0.092	0.360
Very responsive to “jobs in science are fulfilling”	0.352	0.402	0.313	–0.088	0.333
Very responsive to “jobs in science are for girls too”	0.375	0.354	0.392	0.037	0.674
Very responsive to “jobs in science pay well”	0.387	0.263	0.476	0.213	0.042
Very responsive to the short videos	0.546	0.488	0.590	0.102	0.339
<i>E. Overall impression of the role model</i>					
Were gender stereotypes strong among students?					
Yes, very much	0.089	0.039	0.128	0.089	0.057
Rather yes	0.313	0.276	0.341	0.066	0.337
Rather no/not at all	0.598	0.685	0.530	–0.155	0.074
How did the classroom intervention go?					
Very well	0.556	0.535	0.572	0.037	0.670
Well	0.369	0.386	0.355	–0.030	0.716
Average/not so well/not well at all	0.075	0.079	0.072	–0.006	0.821
Was the intervention well suited to the students?					
Yes, very much	0.474	0.449	0.494	0.045	0.661
Rather yes	0.471	0.504	0.446	–0.058	0.574
Rather no/not at all	0.055	0.047	0.060	0.013	0.592
Number of role models	56	21	35		
Number of classroom interventions	290	124	166		

Notes: The summary statistics are computed from the post-intervention role model survey that was administered online to collect feedback about the classroom visits. The unit of observation is a classroom intervention. ^a The science subjects taught in high school are mathematics, physics and chemistry, and earth and life sciences.

Table E4 – Compliance with Random Assignment

	All classes (1)	Classes assigned to	
		Control group (2)	Treatment group (3)
Panel A. Grade 10			
Number of classes visited by a role model	199	2	197
Number of classes not visited by a role model	217	205	12
Number of students	13,700	6,801	6,899
Student-level compliance with random assignment	0.97	0.99	0.94
Panel B. Grade 12 (science track)			
Number of classes visited by a role model	91	2	89
Number of classes not visited by a role model	94	90	4
Number of students	5,751	2,853	2,898
Student-level compliance with random assignment	0.97	0.98	0.95

Notes: This table reports compliance with the random assignment of Grade 10 and Grade 12 (science track) classes to the treatment and control groups. Two-way non-compliance was due to either classes in the treatment not being visited by a role model or to classes in the control group being visited by a role model.

Table E5 – Student Post-Treatment Survey: Response Rates

			Within school	
	Control group (1)	Treatment group (2)	Difference T–C (3)	<i>p</i> -value of diff. (4)
Panel A. Grade 10				
Survey response rate	0.879	0.905	0.026 (0.012)	0.026
Number of students	6,801	6,899	13,700	
Panel B. Grade 12 (science track)				
Survey response rate	0.909	0.912	0.005 (0.012)	0.693
Number of students	2,853	2,898	5,751	

Notes: This table reports the student survey response rate for students in the Grade 10 and Grade 12 (science track) classes that participated in the program. The response rates are computed based on the list of all students who were recorded in the *Bases Élèves académiques* as being enrolled in the participating classes during the academic year 2015/16. Columns 1 and 2 show the response rate of students in the control and treatment groups, respectively. Column 3 reports the coefficient from the regression of survey participation on the treatment group indicator, with *p*-values reported in column 4. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (in parentheses) are adjusted for clustering at the unit of randomization (class).

Table E6 – Treatment-Control Balance: Survey Respondents

			Within school	
	Control group (1)	Treatment group (2)	Difference T–C (3)	p-value of diff. (4)
Panel A. Grade 10				
<i>Student characteristics</i>				
Female	0.538	0.521	−0.014	0.160
Age (years)	15.12	15.11	−0.01	0.248
Non-French	0.057	0.060	0.003	0.528
High SES	0.382	0.389	0.005	0.496
Medium- high SES	0.133	0.127	−0.006	0.248
Medium-low SES	0.245	0.235	−0.009	0.200
Low SES	0.240	0.248	0.010	0.158
Number of siblings	1.483	1.482	−0.001	0.954
Class size	33.23	33.25	0.02	0.837
At least one science elective course	0.394	0.402	0.009	0.693
At least one standard elective course	0.773	0.738	−0.032	0.132
DNB percentile rank in math	59.09	59.04	−0.18	0.760
DNB percentile rank in French	58.14	58.41	0.08	0.893
<i>Test of joint significance</i>	<i>F</i> -stat: 0.634 (<i>p</i> -value: 0.813)			
<i>Predicted track in Grade 11</i>				
Grade 11: Science track	0.454	0.459	0.004	0.577
Grade 11: Science - general track	0.381	0.385	0.003	0.666
Grade 11: Science - technical track	0.073	0.074	0.001	0.670
N	5,981	6,245	12,226	
Panel B. Grade 12 (science track)				
<i>Student characteristics</i>				
Female	0.504	0.489	−0.014	0.319
Age (years)	17.13	17.09	−0.05	0.001
Non-French	0.053	0.046	−0.008	0.129
High SES	0.446	0.481	0.038	0.001
Medium-high SES	0.138	0.138	−0.000	0.979
Medium-low SES	0.219	0.196	−0.022	0.001
Low SES	0.197	0.184	−0.016	0.086
Number of siblings	1.502	1.487	−0.021	0.355
Class size	31.69	32.12	0.30	0.314
DNB percentile rank in math	74.52	74.00	−0.09	0.874
DNB percentile rank in French	69.59	70.00	0.68	0.248
<i>Test of joint significance</i>	<i>F</i> -stat: 1.218 (<i>p</i> -value: 0.282)			
<i>Predicted undergraduate major</i>				
Major: STEM	0.395	0.395	0.001	0.807
Major: selective STEM	0.181	0.184	0.005	0.189
Major: male-dominated STEM	0.283	0.284	0.002	0.561
N	2,594	2,642	5,236	

Notes: Each row corresponds to a different linear regression with the dependent variable listed on the left, separately for students in Grade 10 (Panel A) and in Grade 12 (Panel B). The sample is restricted to students who answered the post-intervention survey. Columns 1 and 2 show the average value for students in the control and treatment groups, respectively. Column 3 reports the coefficient from the regression of each variable on the treatment group indicator, with the *p*-value reported in column 4. The regression controls for school fixed effects to account for the fact that randomization was stratified by school, and standard errors are adjusted for clustering at the unit of randomization (class). The *F*-statistic is from a test of the joint significance of the coefficients in a regression of the treatment group indicator on all student characteristics. High school tracks (Panel A) and undergraduate majors (Panel B) are predicted for each student using the coefficients from a linear regression of the corresponding binary variable (e.g., enrollment in a STEM major) on all student characteristics listed in the table. This model is fitted separately by grade level on the sample of students in the control group.

Table E7 – Balancing Test: High Schools Visited by Professionals and Researchers, Grade 10 Students

	High school visited by		Difference (2)–(1) (3)	<i>p</i> -value of diff. (4)
	Researcher (1)	Professional (2)		
<i>School characteristics</i>				
Education district: Paris	0.165	0.167	0.002	0.958
Education district: Créteil	0.273	0.317	0.044	0.321
Education district: Versailles	0.562	0.516	−0.046	0.348
Private school	0.092	0.224	0.132	0.000
Share of female students	0.523	0.527	0.005	0.627
Pass rate on Baccalauréat exam in 2015 ^a	0.904	0.916	0.012	0.041
Grade 10 students: science track in Grade 11 ^b	0.405	0.412	0.006	0.597
Grade 10 students: general science track in Grade 11 ^b	0.341	0.337	−0.005	0.672
Grade 10 students: technological science track in Grade 11 ^b	0.064	0.075	0.011	0.135
<i>Student characteristics</i>				
Female	0.525	0.531	0.007	0.623
Non-French	0.065	0.057	−0.008	0.185
High SES	0.345	0.410	0.064	0.002
Medium- high SES	0.132	0.125	−0.007	0.322
Medium-low SES	0.250	0.235	−0.015	0.124
Low SES	0.272	0.231	−0.042	0.013
Number of siblings	1.482	1.488	0.007	0.862
At least one science elective course	0.416	0.376	−0.040	0.250
At least one standard elective course	0.772	0.738	−0.034	0.197
Age (years)	15.12	15.13	0.01	0.598
DNB percentile rank in math	57.80	59.02	1.22	0.380
DNB percentile rank in French	56.77	58.71	1.93	0.120
Class size	33.38	33.14	−0.25	0.343
<i>Predicted track in Grade 11</i>				
Grade 11: Science track	0.448	0.454	0.006	0.668
Grade 11: Science - general track	0.374	0.375	0.002	0.915
Grade 11: Science - technological track	0.074	0.079	0.005	0.517
N	6,059	7,641	13,700	

Notes: Each row corresponds to a different linear regression with the dependent variable listed on the left for students enrolled in Grade 10 in 2015/16. Columns 1 and 2 show the average value for students whose high school was visited by a role model with a professional or a research background, respectively. Column 3 reports the coefficient from the regression of each variable on an indicator that takes the value one if the school was visited by a professional and zero if the school was visited by a researcher, with the *p*-value reported in column 4. Standard errors are adjusted for clustering at the class level. High school tracks in Grade 11 are predicted for each student using the coefficients from a linear regression of the corresponding binary variable (e.g., enrollment in the general science track) on all the school and student characteristics listed in the table. This model is fitted on the sample of students in the control group. ^a The *Baccalauréat* pass rate is computed for students who were enrolled in Grade 12 in 2014/15, i.e., in the year before the intervention, and who took the exams in the general or technological tracks. ^b The share of students enrolled in the science track in Grade 11 is computed for students who were enrolled in Grade 10 in 2014/15.

Table E8 – Balancing Test: High Schools Visited by Professionals and Researchers, Grade 12 Students

	High school visited by		Difference	<i>p</i> -value
	Researcher	Professional	(2)–(1)	of diff.
	(1)	(2)	(3)	(4)
<i>School characteristics</i>				
Education district: Paris	0.164	0.163	−0.001	0.985
Education district: Créteil	0.223	0.321	0.098	0.138
Education district: Versailles	0.614	0.517	−0.097	0.195
Private school	0.096	0.244	0.148	0.007
Share of female students	0.533	0.543	0.010	0.379
Pass rate on Baccalauréat exam in 2015 ^a	0.911	0.912	0.002	0.849
Grade 12 (science track) students: STEM major in higher ed. ^b	0.409	0.384	−0.025	0.050
Grade 12 (science track) students: selective STEM in higher ed. ^b	0.191	0.202	0.010	0.484
Grade 12 (science track) students: male-dom. STEM in higher ed. ^b	0.309	0.299	−0.010	0.431
<i>Student characteristics</i>				
Female	0.474	0.505	0.032	0.114
Non-French	0.057	0.046	−0.010	0.272
High SES	0.437	0.484	0.046	0.169
Medium- high SES	0.146	0.128	−0.018	0.138
Medium-low SES	0.213	0.205	−0.009	0.544
Low SES	0.203	0.184	−0.019	0.428
Number of siblings	1.454	1.532	0.079	0.100
Age (years)	17.14	17.11	−0.03	0.323
DNB percentile rank in math	72.96	74.90	1.94	0.213
DNB percentile rank in French	68.00	70.83	2.83	0.057
Class size	32.67	31.44	−1.22	0.026
<i>Predicted undergraduate major</i>				
Major: STEM	0.392	0.378	−0.013	0.062
Major: selective STEM	0.170	0.181	0.011	0.347
Major: male-dominated STEM	0.277	0.274	−0.003	0.731
N	2,492	3,259	5,751	

Notes: Each row corresponds to a different linear regression with the dependent variable listed on the left for students enrolled in Grade 12 (science track) in 2015/16. Columns 1 and 2 show the average value for students whose high school was visited by a role model with a professional or a research background, respectively. Column 3 reports the coefficient from the regression of each variable on an indicator that takes the value one if the school was visited by a professional and zero if the school was visited by a researcher, with the *p*-value reported in column 4. Standard errors are adjusted for clustering at the class level. Undergraduate majors are predicted for each student using the coefficients from a linear regression of the corresponding binary variable (e.g., enrollment in a STEM major) on all the school and student characteristics listed in the table. This model is fitted on the sample of students in the control group. ^a The *Baccalauréat* pass rate is computed for students who were enrolled in Grade 12 in 2014/15, i.e., in the year before the intervention, and who took the exams in the general or technological tracks. ^b The share of students enrolled in a STEM undergraduate major in higher education is computed for students who were enrolled in Grade 12 (science track) in 2014/15.

F Effects of Role Model Interventions: Additional Results

F.1 Student Perceptions

Table F1 – Gender Differences in Aptitude for Mathematics

	Girls			Boys		
	Control group mean (1)	Treatment effect (LATE) (2)	p -value [q -value] (3)	Control group mean (4)	Treatment effect (LATE) (5)	p -value [q -value] (6)
Panel A. Grade 10						
Equal gender aptitude for math (index)	0.115	0.109*** (0.025)	0.000 [0.001]	−0.134	0.148*** (0.030)	0.000 [0.001]
M and W are born with different brains	0.211	−0.050*** (0.010)	0.000 [0.001]	0.209	−0.048*** (0.011)	0.000 [0.001]
Men are more gifted in math than women	0.186	−0.026** (0.011)	0.015 [0.016]	0.299	−0.048*** (0.014)	0.001 [0.001]
N		6,475			5,751	
Panel B. Grade 12 (science track)						
Equal gender aptitude for math (index)	0.158	0.095*** (0.028)	0.001 [0.002]	−0.161	0.132*** (0.040)	0.001 [0.002]
M and W are born with different brains	0.143	−0.023** (0.010)	0.026 [0.026]	0.180	−0.038*** (0.014)	0.006 [0.013]
Men are more gifted in math than women	0.163	−0.038*** (0.012)	0.002 [0.005]	0.266	−0.028* (0.015)	0.072 [0.073]
N		2,600			2,636	

Notes: This table reports estimates of the treatment effects of the role model interventions on students' perceptions regarding the aptitude of men and women for mathematics, separately by grade level and gender. The sample is restricted to students who completed the post-intervention questionnaire. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 4 report the average value for students in the control group. Columns 2 and 5 report the local average treatment effect (LATE) estimates. They are obtained from a regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report the cluster-robust p -value of the estimated treatment effect and, in square brackets, the p -value (q -value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage q -values introduced in Benjamini et al. (2006) and described in Anderson (2008). The q -values associated with the synthetic index (highlighted in bold) are adjusted for multiple testing across the study's nine main outcomes of interest, separately by grade level and gender (see Appendix D for details). The q -values for the individual components of the index are adjusted for multiple testing across the index components, separately by grade level and gender. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

Table F2 – Taste for Science Subjects

	Girls			Boys		
	Control group mean (1)	Treatment effect (LATE) (2)	<i>p</i> -value [<i>q</i> -value] (3)	Control group mean (4)	Treatment effect (LATE) (5)	<i>p</i> -value [<i>q</i> -value] (6)
Panel A. Grade 10						
Taste for science subjects (index)	−0.169	−0.038 (0.036)	0.294 [0.442]	0.197	−0.019 (0.031)	0.533 [0.685]
Enjoys math (<i>z</i> -score)	−0.147	−0.002 (0.034)	0.961 [0.961]	0.186	−0.002 (0.031)	0.935 [0.935]
Enjoys physics-chemistry (<i>z</i> -score)	−0.170	−0.040 (0.038)	0.289 [0.578]	0.223	−0.022 (0.033)	0.505 [0.935]
Enjoys earth and life sciences (<i>z</i> -score)	−0.042	−0.058 (0.039)	0.137 [0.548]	0.086	−0.027 (0.035)	0.443 [0.935]
Enjoys science in general	0.661	−0.011 (0.015)	0.444 [0.593]	0.790	0.003 (0.012)	0.804 [0.935]
N		6,475			5,751	
Panel B. Grade 12 (science track)						
Taste for science subjects (index)	−0.002	0.016 (0.034)	0.632 [0.633]	0.002	−0.000 (0.039)	0.998 [0.999]
Enjoys math (<i>z</i> -score)	−0.097	0.067* (0.040)	0.089 [0.357]	0.100	0.075* (0.040)	0.063 [0.203]
Enjoys physics-chemistry (<i>z</i> -score)	−0.089	−0.001 (0.044)	0.984 [0.984]	0.102	−0.021 (0.040)	0.598 [0.599]
Enjoys earth and life sciences (<i>z</i> -score)	0.203	−0.030 (0.038)	0.435 [0.871]	−0.215	−0.059 (0.059)	0.318 [0.424]
Enjoys science in general	0.918	−0.001 (0.009)	0.887 [0.984]	0.930	0.013 (0.008)	0.101 [0.203]
N		2,600			2,636	

Notes: This table reports estimates of the treatment effects of the role model interventions on students' taste for science subjects taught at school, separately by grade level and gender. The sample is restricted to students who completed the post-intervention questionnaire. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 4 report the average value for students in the control group. Columns 2 and 5 report the local average treatment effect (LATE) estimates. They are obtained from a regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report the cluster-robust *p*-value of the estimated treatment effect and, in square brackets, the *p*-value (*q*-value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage *q*-values introduced in Benjamini et al. (2006) and described in Anderson (2008). The *q*-values associated with the synthetic index (highlighted in bold) are adjusted for multiple testing across the study's nine main outcomes of interest, separately by grade level and gender (see Appendix D for details). The *q*-values for the individual components of the index are adjusted for multiple testing across the index components, separately by grade level and gender. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

Table F3 – Math Self-Concept

	Girls			Boys		
	Control group mean (1)	Treatment effect (LATE) (2)	p -value [q -value] (3)	Control group mean (4)	Treatment effect (LATE) (5)	p -value [q -value] (6)
Panel A. Grade 10						
Math self-concept (index)	−0.198	−0.008 (0.031)	0.806 [0.807]	0.231	0.039 (0.032)	0.217 [0.326]
Self-assessed math performance (z -score)	−0.127	−0.016 (0.034)	0.634 [0.634]	0.168	0.021 (0.032)	0.502 [0.642]
Lost in front of a math problem	0.553	0.010 (0.014)	0.478 [0.634]	0.344	−0.007 (0.013)	0.610 [0.642]
Worried when thinking about math	0.617	−0.025* (0.013)	0.052 [0.109]	0.420	−0.032** (0.015)	0.028 [0.111]
Can succeed in science subjects if puts in effort	0.843	0.018* (0.009)	0.054 [0.109]	0.883	−0.004 (0.008)	0.642 [0.642]
N		6,475			5,751	
Panel B. Grade 12 (science track)						
Math self-concept (index)	−0.184	0.050 (0.039)	0.202 [0.228]	0.187	0.072** (0.035)	0.041 [0.062]
Self-assessed math performance (z -score)	−0.126	0.039 (0.038)	0.304 [0.406]	0.123	0.079** (0.038)	0.038 [0.077]
Lost in front of a math problem	0.486	−0.028 (0.020)	0.168 [0.336]	0.325	−0.028* (0.016)	0.072 [0.096]
Worried when thinking about math	0.560	−0.037** (0.019)	0.048 [0.193]	0.384	−0.051*** (0.016)	0.002 [0.007]
Can succeed in science subjects if puts in effort	0.942	−0.005 (0.007)	0.512 [0.512]	0.949	0.006 (0.007)	0.384 [0.385]
N		2,600			2,636	

Notes: This table reports estimates of the treatment effects of the role model interventions on students' math self-concept, separately by grade level and gender. The sample is restricted to students who completed the post-intervention questionnaire. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 4 report the average value for students in the control group. Columns 2 and 5 report the local average treatment effect (LATE) estimates. They are obtained from a regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report the cluster-robust p -value of the estimated treatment effect and, in square brackets, the p -value (q -value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage q -values introduced in Benjamini et al. (2006) and described in Anderson (2008). The q -values associated with the synthetic index (highlighted in bold) are adjusted for multiple testing across the study's nine main outcomes of interest, separately by grade level and gender (see Appendix D for details). The q -values for the individual components of the index are adjusted for multiple testing across the index components, separately by grade level and gender. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

Table F4 – Science-Related Career Aspirations

	Girls			Boys		
	Control group mean (1)	Treatment effect (LATE) (2)	<i>p</i> -value [<i>q</i> -value] (3)	Control group mean (4)	Treatment effect (LATE) (5)	<i>p</i> -value [<i>q</i> -value] (6)
Panel A. Grade 10						
Science-related career aspirations (index)	−0.103	0.012 (0.030)	0.695 [0.807]	0.120	0.007 (0.029)	0.801 [0.902]
Some jobs in science are interesting	0.845	0.019** (0.009)	0.050 [0.200]	0.854	−0.000 (0.010)	1.000 [1.000]
Would consider a job in science	0.466	−0.004 (0.015)	0.776 [0.776]	0.587	0.023* (0.014)	0.089 [0.358]
Interested in at least one STEM job ^a	0.642	0.005 (0.013)	0.696 [0.776]	0.849	0.013 (0.010)	0.181 [0.363]
Wage prospects important in career choice (<i>z</i> -score)	−0.045	−0.012 (0.029)	0.682 [0.776]	0.038	0.007 (0.027)	0.792 [1.000]
N		6,475			5,751	
Panel B. Grade 12 (science track)						
Science-related career aspirations (index)	−0.045	0.113*** (0.037)	0.002 [0.003]	0.046	0.050 (0.033)	0.131 [0.169]
Some jobs in science are interesting	0.961	0.013** (0.005)	0.013 [0.026]	0.940	0.021*** (0.008)	0.005 [0.022]
Would consider a job in science	0.721	0.031** (0.013)	0.019 [0.026]	0.762	0.030** (0.014)	0.029 [0.058]
Interested in at least one STEM job ^a	0.817	0.000 (0.011)	0.964 [0.964]	0.899	−0.001 (0.009)	0.946 [0.947]
Wage prospects important in career choice (<i>z</i> -score)	−0.043	0.119*** (0.038)	0.002 [0.007]	0.037	0.049 (0.031)	0.111 [0.149]
N		2,600			2,636	

Notes: This table reports estimates of the treatment effects of the role model interventions on students' self-reported science-related career aspirations, separately by grade level and gender. The sample is restricted to students who completed the post-intervention questionnaire. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 4 report the average value for students in the control group. Columns 2 and 5 report the local average treatment effect (LATE) estimates. They are obtained from a regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report the cluster-robust *p*-value of the estimated treatment effect and, in square brackets, the *p*-value (*q*-value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage *q*-values introduced in Benjamini et al. (2006) and described in Anderson (2008). The *q*-values associated with the synthetic index (highlighted in bold) are adjusted for multiple testing across the study's nine main outcomes of interest, separately by grade level and gender (see Appendix D for details). The *q*-values for the individual components of the index are adjusted for multiple testing across the index components, separately by grade level and gender. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$. ^a The STEM occupations in the list were chemist, computer scientist, engineer, industrial designer, renewable energy technician, and researcher in biology. The non-STEM occupations were lawyer, pharmacist, physician, and psychologist.

F.2 Educational Choices

Table F5 – Grade 10 Students: Enrollment Status the Following Year (Detailed)

	Grade 10 students					
	Girls			Boys		
	Control group mean (1)	Treatment effect (LATE) (2)	<i>p</i> -value [<i>q</i> -value] (3)	Control group mean (4)	Treatment effect (LATE) (5)	<i>p</i> -value [<i>q</i> -value] (6)
Panel A. Grade 11 STEM tracks						
All STEM tracks						
Science track	0.355	−0.004 (0.014)	0.753 [0.807]	0.551	−0.002 (0.015)	0.910 [0.910]
General vs. technological STEM track						
Science - general track	0.328	0.001 (0.013)	0.942 [0.942]	0.416	0.007 (0.014)	0.613 [0.614]
Science - technological track	0.026	−0.005 (0.003)	0.128 [0.256]	0.135	−0.009 (0.008)	0.300 [0.601]
Panel B. Grade 11 non-STEM tracks						
All non-STEM tracks						
Non-STEM tracks	0.545	0.006 (0.014)	0.642	0.324	0.018 (0.014)	0.191
General vs. technological non-STEM tracks						
Humanities track	0.121	−0.002 (0.010)	0.846	0.028	0.005 (0.005)	0.284
Social sciences track	0.252	0.011 (0.012)	0.323	0.163	0.013 (0.010)	0.211
Non-STEM technological tracks	0.171	−0.003 (0.010)	0.759	0.133	0.001 (0.009)	0.944
Repeater or dropout	0.101	−0.002 (0.009)	0.818	0.126	−0.016* (0.009)	0.070
N		7,241			6,459	

Notes: This table reports estimates of the treatment effects of the role model interventions on Grade 10 students' enrollment outcomes in the academic year following the classroom interventions, i.e. 2016/17, separately by gender. The enrollment outcomes are measured using student-level administrative data. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 4 report the average value for students in the control group. Columns 2 and 5 report the local average treatment effect (LATE) estimates. They are obtained from a regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report the cluster-robust *p*-value of the estimated treatment effect and, in square brackets, the *p*-value (*q*-value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage *q*-values introduced in Benjamini et al. (2006) and described in Anderson (2008). The *q*-values associated with the treatment effect estimates on the probability of being enrolled in the science track in Grade 11 (Panel A) are adjusted for multiple testing across the study's nine main outcomes of interest, separately by gender (see Appendix D for details). The *q*-values associated with the treatment effect estimates for the different STEM tracks (general vs. technological) are adjusted for multiple testing across these two tracks, separately by gender.

*** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

Table F6 – Grade 12 Students: Enrollment Status the Following Year (Detailed)

	Grade 12 (science track) students					
	Girls			Boys		
	Control group mean (1)	Treatment effect (LATE) (2)	<i>p</i> -value [<i>q</i> -value] (3)	Control group mean (4)	Treatment effect (LATE) (5)	<i>p</i> -value [<i>q</i> -value] (6)
Panel A. STEM undergraduate programs						
<i>All undergraduate STEM majors</i>						
Major: STEM	0.289	0.024* (0.014)	0.080 [0.103]	0.470	0.003 (0.020)	0.886 [0.988]
<i>Selective STEM majors</i>						
Math, physics, engineering, computer science (CPGE)	0.084	0.026*** (0.009)	0.006 [0.032]	0.211	0.023 (0.015)	0.123 [0.206]
Earth and life sciences (CPGE)	0.020	0.008 (0.005)	0.137 [0.230]	0.010	0.002 (0.003)	0.617 [0.771]
Sciences - vocational (STS)	0.006	0.002 (0.003)	0.474 [0.474]	0.011	−0.005 (0.003)	0.110 [0.206]
<i>Non-selective STEM majors</i>						
Math, physics, computer science	0.077	0.010 (0.008)	0.209 [0.262]	0.157	−0.002 (0.012)	0.884 [0.885]
Earth and life sciences	0.103	−0.022** (0.009)	0.014 [0.035]	0.081	−0.015* (0.008)	0.053 [0.206]
Panel B. Non-STEM undergraduate programs						
<i>All undergraduate non-STEM majors</i>						
Major: non-STEM	0.507	−0.032** (0.016)	0.045	0.293	−0.005 (0.014)	0.717
<i>Selective non-STEM majors</i>						
Business and economics (CPGE)	0.021	0.003 (0.004)	0.566	0.017	0.005 (0.004)	0.219
Humanities (CPGE)	0.014	−0.004 (0.003)	0.225	0.003	−0.001 (0.001)	0.470
Other vocational (STS)	0.011	−0.008*** (0.003)	0.006	0.008	−0.005** (0.002)	0.036
<i>Non-selective non-STEM majors</i>						
Medicine and pharmacy	0.259	−0.005 (0.015)	0.722	0.108	0.006 (0.011)	0.573
Law and economics	0.107	−0.008 (0.011)	0.478	0.079	0.002 (0.008)	0.758
Humanities and psychology	0.080	−0.008 (0.009)	0.339	0.040	−0.006 (0.006)	0.296
Sports studies	0.018	−0.004 (0.004)	0.396	0.036	−0.005 (0.006)	0.377
Not enrolled in higher education	0.206	0.008 (0.015)	0.581	0.237	0.004 (0.016)	0.814
N		2,827			2,924	

Notes: This table reports estimates of the treatment effects of the role model interventions on science track Grade 12 (science track) students' enrollment outcomes in the academic year following the classroom interventions, i.e. 2016/17, separately by gender. The enrollment outcomes are measured using student-level administrative data. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 4 report the average value for students in the control group. Columns 2 and 5 report the local average treatment effect (LATE) estimates. They are obtained from a regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report the cluster-robust *p*-value of the estimated treatment effect and, in square brackets, the *p*-value (*q*-value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage *q*-values introduced in Benjamini et al. (2006) and described in Anderson (2008). The *q*-values associated with the treatment effect estimates on the probability of enrolling in a STEM undergraduate major (Panel A) are adjusted for multiple testing across the study's nine main outcomes of interest, separately by gender (see Appendix D for details). The *q*-values associated with the estimates for the different selective and non-selective STEM majors (Panel A) are adjusted for multiple testing across these different STEM majors, separately by gender. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

F.3 Academic Performance

Table F7 – Grade 12 Students: Performance in *Baccalauréat* Exams

	Grade 12 (science track) students					
	Girls			Boys		
	Control group mean (1)	Treatment effect (LATE) (2)	<i>p</i> -value [<i>q</i> -value] (3)	Control group mean (4)	Treatment effect (LATE) (5)	<i>p</i> -value [<i>q</i> -value] (6)
Baccalauréat percentile rank in math	46.21	0.693 (0.957)	0.469 [0.626]	47.47	1.661 (1.024)	0.105 [0.210]
Baccalauréat percentile rank in French	54.37	−0.051 (1.113)	0.964 [0.964]	43.51	−0.331 (0.803)	0.680 [0.680]
Baccalauréat percentile rank	53.52	−1.121 (1.066)	0.293 [0.626]	47.29	1.712* (1.040)	0.100 [0.210]
Obtained the Baccalauréat	0.928	−0.010 (0.010)	0.334 [0.626]	0.877	−0.005 (0.010)	0.623 [0.680]
N		2,827			2,924	

Notes: This table reports estimates of the treatment effects of the role model interventions on Grade 12 (science track) students' performance on the *Baccalauréat* exams, separately by gender. The *Baccalauréat* outcomes are measured using student-level administrative data. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 4 report the average value for students in the control group. Columns 2 and 5 report the local average treatment effect (LATE) estimates. They are obtained from a regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report the cluster-robust *p*-value of the estimated treatment effect and, in square brackets, the *p*-value (*q*-value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage *q*-values introduced in Benjamini et al. (2006) and described in Anderson (2008). The *q*-values are adjusted for multiple testing across the four *Baccalauréat* outcomes, separately by gender. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

G Robustness Checks

Table G1 – Treatment Effects on Student Perceptions: Controlling for Baseline Characteristics

	Girls			Boys		
	Control group mean (1)	Treatment effect (LATE) (2)	<i>p</i> -value [<i>q</i> -value] (3)	Control group mean (4)	Treatment effect (LATE) (5)	<i>p</i> -value [<i>q</i> -value] (6)
Panel A. Grade 10						
Positive perceptions of science-related careers (index)	−0.020	0.245*** (0.027)	0.000 [0.001]	0.023	0.162*** (0.027)	0.000 [0.001]
More men in science-related jobs	0.628	0.154*** (0.013)	0.000 [0.001]	0.629	0.170*** (0.014)	0.000 [0.001]
Equal gender aptitude for math (index)	0.115	0.111*** (0.024)	0.000 [0.001]	−0.134	0.142*** (0.030)	0.000 [0.001]
Women don’t really like science	0.157	0.056*** (0.011)	0.000 [0.001]	0.198	0.101*** (0.013)	0.000 [0.001]
W face discrimination in science-related jobs	0.603	0.126*** (0.013)	0.000 [0.001]	0.527	0.154*** (0.014)	0.000 [0.001]
Taste for science subjects (index)	−0.169	−0.033 (0.031)	0.275 [0.414]	0.197	−0.021 (0.026)	0.431 [0.555]
Math self-concept (index)	−0.198	−0.001 (0.028)	0.981 [0.982]	0.231	0.033 (0.029)	0.250 [0.375]
Science-related careers aspirations (index)	−0.103	0.005 (0.029)	0.851 [0.970]	0.120	0.004 (0.027)	0.871 [0.872]
N		6,475			5,751	
Panel B. Grade 12 (science track)						
Positive perceptions of science-related careers (index)	−0.003	0.296*** (0.032)	0.000 [0.001]	0.003	0.171*** (0.033)	0.000 [0.001]
More men in science-related jobs	0.712	0.122*** (0.016)	0.000 [0.001]	0.717	0.149*** (0.015)	0.000 [0.001]
Equal gender aptitude for math (index)	0.158	0.078*** (0.028)	0.004 [0.007]	−0.161	0.124*** (0.042)	0.003 [0.006]
Women don’t really like science	0.074	0.042*** (0.009)	0.000 [0.001]	0.146	0.073*** (0.015)	0.000 [0.001]
W face discrimination in science-related jobs	0.624	0.085*** (0.020)	0.000 [0.001]	0.600	0.074*** (0.018)	0.000 [0.001]
Taste for science subjects (index)	−0.002	0.018 (0.033)	0.583 [0.583]	0.002	0.014 (0.040)	0.733 [0.825]
Math self-concept (index)	−0.184	0.051 (0.035)	0.139 [0.157]	0.187	0.068** (0.033)	0.038 [0.057]
Science-related careers aspirations (index)	−0.045	0.106*** (0.037)	0.004 [0.007]	0.046	0.068* (0.035)	0.055 [0.071]
N		2,600			2,636	

Notes: This table reports estimates of the treatment effects of the role model interventions on students’ perceptions, separately by grade level and gender, and controlling for students’ baseline characteristics. The sample is restricted to students who completed the post-intervention questionnaire. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 4 report the average value for students in the control group. Columns 2 and 5 report the local average treatment effect (LATE) estimates. They are obtained from a regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. The regression further controls for the student characteristics listed in Table 1 in the main text. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report the cluster-robust *p*-value of the estimated treatment effect and, in square brackets, the *p*-value (*q*-value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage *q*-values introduced in Benjamini et al. (2006) and described in Anderson (2008). The *q*-values are adjusted for multiple testing across the study’s nine main outcomes of interest, separately by grade level and gender (see Appendix D for details). *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

Table G2 – Treatment Effects on Enrollment Outcomes: Controlling for Baseline Characteristics

	Girls			Boys		
	Control group mean (1)	Treatment effect (LATE) (2)	<i>p</i> -value [<i>q</i> -value] (3)	Control group mean (4)	Treatment effect (LATE) (5)	<i>p</i> -value [<i>q</i> -value] (6)
Panel A. Grade 10						
All STEM tracks						
Grade 11: Science track	0.355	−0.002 (0.011)	0.862 [0.970]	0.551	−0.006 (0.012)	0.640 [0.720]
General vs. technological STEM track						
Grade 11: Science - general track	0.328	0.003 (0.010)	0.794 [0.794]	0.416	0.004 (0.011)	0.710 [0.710]
Grade 11: Science - technological track	0.026	−0.005 (0.004)	0.188 [0.377]	0.135	−0.010 (0.008)	0.234 [0.468]
N		7,241			6,459	
Panel B. Grade 12 (science track)						
All undergraduate STEM majors						
Major: STEM	0.289	0.020 (0.014)	0.139 [0.157]	0.470	−0.002 (0.019)	0.925 [0.926]
Selective vs. non-selective STEM						
Major: selective STEM	0.110	0.031*** (0.011)	0.006 [0.012]	0.232	0.008 (0.015)	0.575 [0.575]
Major: non-selective STEM	0.178	−0.011 (0.012)	0.333 [0.333]	0.239	−0.010 (0.013)	0.445 [0.575]
Male- vs. female-dominated STEM						
Major: male-dominated STEM (math, physics, computer science)	0.166	0.034*** (0.012)	0.004 [0.012]	0.379	0.013 (0.019)	0.485 [0.575]
Major: female-dominated STEM (earth and life sciences)	0.123	−0.015 (0.011)	0.169 [0.226]	0.091	−0.015 (0.009)	0.119 [0.477]
N		2,827			2,924	

Notes: This table reports estimates of the treatment effects of the role model interventions on students’ enrollment outcomes in the academic year following the classroom interventions, i.e., 2016/17, separately by grade level and gender, and controlling for student baseline characteristics. The enrollment outcomes are measured using student-level administrative data. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 4 report the average value for students in the control group. Columns 2 and 5 report the local average treatment effect (LATE) estimates. They are obtained from a regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. The regression further controls for the student characteristics listed in Table 1 in the main text. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report the cluster-robust *p*-value of the estimated treatment effect and, in square brackets, the *p*-value (*q*-value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage *q*-values introduced in Benjamini et al. (2006) and described in Anderson (2008). The *q*-values associated with the treatment effect estimates on “Grade 11: Science track” (Panel A) and “Major: STEM” (Panel B) are adjusted for multiple testing across the study’s nine main outcomes of interest, separately by grade level and gender (see Appendix D for details). The *q*-values associated with the treatment effect estimates for the different STEM tracks (Panel A) or the different STEM majors (Panel B) are adjusted for multiple testing across these different STEM tracks or majors, separately by grade level and gender. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

H Randomization Inference

This appendix evaluates the robustness of our results to computing p -values using non-parametric randomization inference tests rather than model-based cluster-robust inference.

Randomization inference, which was first proposed by Fisher (1935) and was later developed by Rosenbaum (2002), has been used in a number of recent RCT studies in economics and political science as an alternative to model-based inference. The intuition behind this approach is relatively straightforward. In RCTs, researchers know exactly how the randomization was performed. Randomization inference uses this knowledge to assess whether observed outcomes in a given sample are likely to have occurred by chance even if the treatment had no effect. This can be obtained numerically through Monte Carlo methods, by computing the treatment effects for varying random draws of the treatment assignment, whose data-generating process is known. This test is non-parametric since it does not make distributional assumptions.^{A.5}

In our setting, the ITT effect under the observed assignment to treatment is estimated using the following reduced-form specification:

$$Y_{ics} = \alpha + \beta T_{cs} + \theta_s + \epsilon_{ics}, \quad (\text{A.1})$$

where Y_{ics} denotes the observed outcome of student i in class c and high school s ; T_{cs} denotes the observed treatment assignment of the student's class; and θ_s are school fixed effects. The ITT estimate under the observed treatment assignment is denoted by $\hat{\beta}$.

To conduct randomization inference, we proceed as follows. Taking into account the fact that randomization was stratified by school and grade level, we first re-assign treatment $R=2,000$ times among participating classes using the exact same stratified procedure.^{A.6} Let $\{P^r\}_{r=1}^R$ denote the set of R random placebo assignments from the randomization process. We then re-estimate the ITT effects of these placebo treatments using the following reduced-form specification, which is estimated separately by grade level and gender:

$$Y_{ics} = \alpha_r + \beta_r P_{cs}^r + \lambda_s + \eta_{ics}, \quad r = 1, \dots, R, \quad (\text{A.2})$$

where P_{cs}^r indicates assignment to a placebo treatment group for random draw r . School fixed effects, λ_s , are included to account for the fact that the randomization is stratified by school.

Since P_r is a randomly generated placebo, $\mathbb{E}(\beta_r) = 0$. Let $F(\hat{\beta}_r)$ denote the empirical c.d.f. of all elements of $\{P_r\}_{r=1}^R$. We test the null hypothesis that the program had no effect on outcome Y by checking if the ITT estimate that we obtain for the observed treatment assignment is in the tails of the distribution of placebo treatments. We can reject $H_0: \hat{\beta} = 0$ with a confidence level of $1 - \alpha$ if $\hat{\beta} \leq F^{-1}\left(\frac{\alpha}{2}\right)$ or $\hat{\beta} \geq F^{-1}\left(1 - \frac{\alpha}{2}\right)$. Since the placebo assignments only vary across randomization units (here classes), this method accounts for correlation within units. Following Davison and Hinkley (1997), we compute the p -values from a two-sided randomization inference test of zero treatment effects as $p = (1 + \sum_{r=1}^R \mathbb{1}(|\hat{\beta}_r| \geq |\hat{\beta}|)) / (1 + R)$.

Table H1 presents the results of randomization inference tests of the hypotheses that the role model interventions had no effect on student perceptions and enrollment outcomes, separately by grade level and gender. The ITT estimates $\hat{\beta}$ are shown in columns 1 and 4. The cluster-robust model-based p -values are reported in columns 2 and 5, while those based on randomization inference are in columns 3 and 6. The results of the randomization inference tests yield p -values that are generally close to the cluster-robust model-based p -values. Although they tend to be slightly more conservative, they confirm the program's statistically significant effects on enrollment in selective and male-dominated STEM programs for girls in Grade 12.

^{A.5}For more details on randomization inference, see Rosenbaum (2010) and Imbens and Rubin (2015).

^{A.6}See Paz and West (2019) for the number of draws to be used.

Table H1 – Randomization Inference for Intention-to-Treat Estimates

	Girls			Boys		
	ITT	<i>p</i> -value: model- based	<i>p</i> -value: rand. inference	ITT	<i>p</i> -value: model- based	<i>p</i> -value: rand. inference
	(1)	(2)	(3)	(4)	(5)	(6)
Panel A. Grade 10						
<i>Student perceptions</i>						
Positive perceptions of science-related careers (index)	0.226	0.000	0.000	0.156	0.000	0.000
More men in science-related jobs	0.145	0.000	0.000	0.157	0.000	0.000
Equal gender aptitude for math (index)	0.101	0.000	0.000	0.138	0.000	0.000
Women don't really like science	0.054	0.000	0.000	0.096	0.000	0.000
Women face discrimination in science-related careers	0.118	0.000	0.000	0.143	0.000	0.000
Taste for science subjects (index)	−0.035	0.298	0.340	−0.018	0.537	0.560
Math self-concept (index)	−0.007	0.808	0.820	0.037	0.221	0.280
Science-related career aspirations (index)	0.011	0.697	0.720	0.007	0.803	0.830
<i>Enrollment outcomes</i>						
Grade 11: Science track	−0.004	0.755	0.780	−0.002	0.911	0.920
Grade 11: Science - general track	0.001	0.942	0.950	0.007	0.617	0.660
Grade 11: Science - technological track	−0.005	0.131	0.190	−0.008	0.306	0.350
N	7,241			6,459		
Panel B. Grade 12 (science track)						
<i>Student perceptions</i>						
Positive perceptions of science-related careers (index)	0.293	0.000	0.000	0.145	0.000	0.000
More men in science-related jobs	0.118	0.000	0.000	0.140	0.000	0.000
Equal gender aptitude for math (index)	0.090	0.001	0.020	0.124	0.002	0.020
Women don't really like science	0.042	0.000	0.000	0.069	0.000	0.000
Women face discrimination in science-related careers	0.090	0.000	0.000	0.068	0.000	0.000
Taste for science subjects (index)	0.015	0.640	0.740	−0.000	0.998	1.000
Math self-concept (index)	0.047	0.214	0.360	0.068	0.044	0.140
Science-related career aspirations (index)	0.106	0.003	0.020	0.047	0.141	0.270
<i>Enrollment outcomes</i>						
Undergraduate major: STEM	0.022	0.091	0.220	0.003	0.889	0.920
Undergraduate major: selective STEM	0.033	0.002	0.030	0.019	0.208	0.360
Undergraduate major: non-selective STEM	−0.010	0.328	0.480	−0.016	0.220	0.370
Undergraduate major: male-dominated STEM	0.035	0.002	0.020	0.016	0.397	0.530
Undergraduate major: female-dominated STEM	−0.014	0.162	0.320	−0.013	0.128	0.270
N	2,827			2,924		

Notes: This table presents the results of randomization inference tests of the hypotheses that the program had no effect on student perceptions and enrollment outcomes. We randomly re-assign treatment 2,000 times among participating classes within each school and grade level, and re-estimate the ITT effects of these placebo treatments. The ITT estimates under the observed assignment are reported in columns 1 and 4 separately by gender. The associated cluster-robust model-based *p*-values are shown in columns 2 and 5. The randomization inference *p*-values are reported in columns 3 and 6. They are computed from a two-sided randomization inference test of zero treatment effects as $p = \left(1 + \sum_{r=1}^R \mathbb{1}(|\hat{\beta}_r| \geq |\hat{\beta}|)\right) / (1 + R)$, where $\{\hat{\beta}_r\}_{r=1}^R$ is the set of R placebo ITT estimates, $\hat{\beta}$ is the ITT estimate under the observed assignment, and $\mathbb{1}(\cdot)$ is the indicator function.

I Information, Persistence, Timing: Additional Results

I.1 Intensity of Information Provision

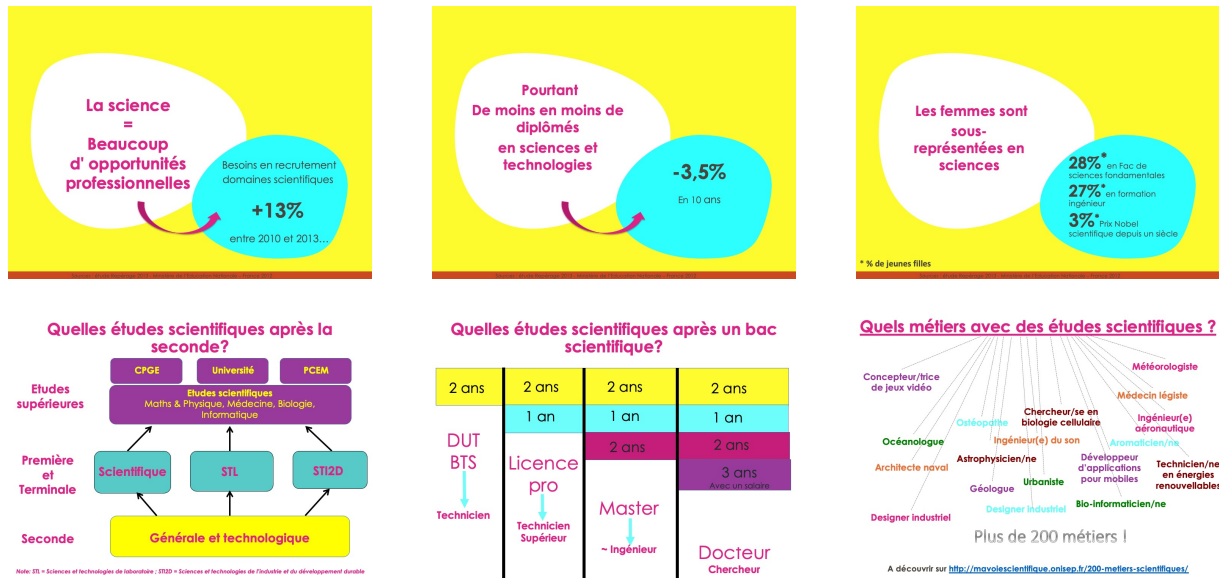


Figure I1 – Screenshots of the Slides Providing General Information on STEM Careers (“Regular Slides”)

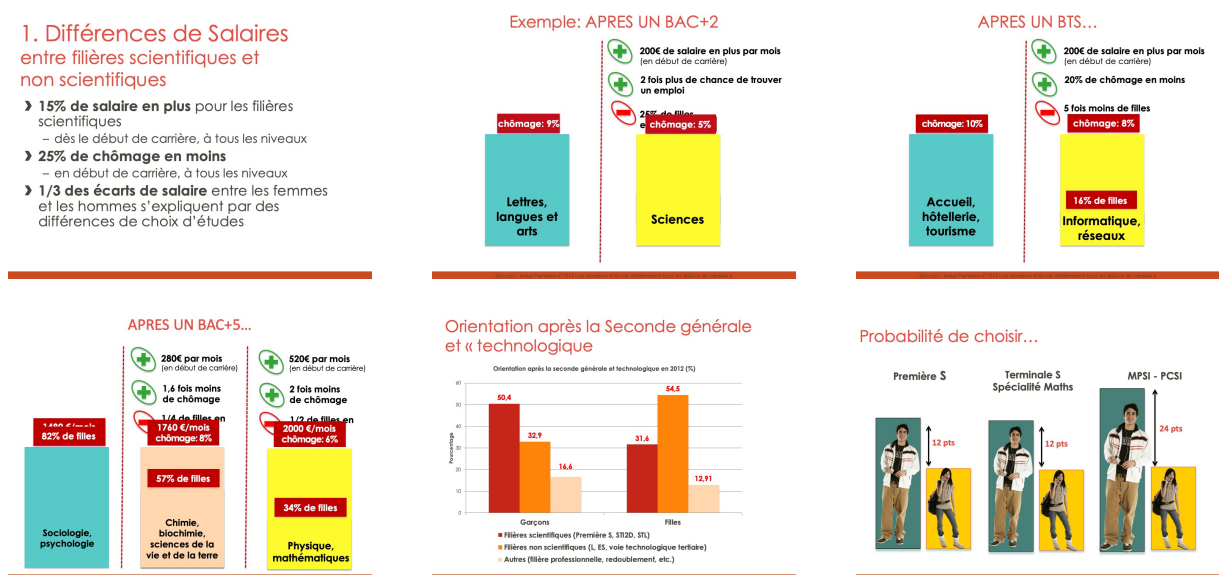


Figure I2 – Screenshots of the Additional Slides Providing General Information on STEM Careers (“Augmented Slides”)

Table I1 – Balancing Test: Classrooms Assigned to Role Models who were Provided with the Regular vs. Augmented Sets of Slides

	Set of slides		Difference (2)–(1)	<i>p</i> -value of diff.
	Regular (1)	Augmented (2)	(3)	(4)
Panel A. Grade 10				
<i>Student characteristics</i>				
Girl	0.532	0.526	–0.006	0.651
Non-French	0.052	0.067	0.015	0.014
High SES	0.396	0.369	–0.027	0.193
Medium- high SES	0.136	0.122	–0.014	0.059
Medium-low SES	0.238	0.244	0.006	0.559
Low SES	0.229	0.265	0.035	0.030
Number of siblings	1.467	1.500	0.033	0.383
At least one science elective course	0.387	0.399	0.012	0.722
At least one standard elective course	0.746	0.759	0.012	0.648
Age (years)	15.10	15.14	0.04	0.002
DNB percentile rank in math	58.37	58.57	0.20	0.886
DNB percentile rank in French	56.79	58.69	1.90	0.122
Number of classmates	33.86	32.76	–1.10	0.000
<i>Predicted track in Grade 11</i>				
Grade 11: Science track	0.443	0.457	0.014	0.248
Grade 11: Science - general track	0.366	0.380	0.014	0.326
Grade 11: Science - technical track	0.077	0.077	0.000	0.923
N	6,047	7,653	13,700	
Panel B. Grade 12 (science track)				
<i>Student characteristics</i>				
Girl	0.491	0.492	0.001	0.951
Non-French	0.044	0.057	0.014	0.133
High SES	0.475	0.453	–0.022	0.519
Medium- high SES	0.140	0.132	–0.008	0.517
Medium-low SES	0.209	0.208	–0.001	0.936
Low SES	0.176	0.207	0.031	0.197
Number of siblings	1.479	1.516	0.037	0.454
Age (years)	17.12	17.13	0.01	0.673
DNB percentile rank in math	74.19	73.94	–0.25	0.870
DNB percentile rank in French	69.45	69.75	0.30	0.843
Number of classmates	32.13	31.83	–0.29	0.602
<i>Predicted undergraduate major</i>				
Major: STEM	0.384	0.382	–0.002	0.735
Undergraduate major: male-dominated STEM	0.276	0.274	–0.001	0.856
Undergraduate major: selective STEM	0.179	0.175	–0.004	0.684
N	2,748	3,003	5,751	

Notes: Each row corresponds to a different linear regression with the dependent variable listed on the left for students enrolled in Grade 10 in 2015/16 (Panel A) and in Grade 12 (Panel B). Columns 1 and 2 show the average value for students whose high school was visited by a role model provided with the regular or augmented set of slides, respectively. Column 3 reports the coefficient from the regression of each variable on an indicator that takes the value one if the school was visited by a role model who received the augmented slides and zero if the school was visited by a role model who received the regular slides, with the *p*-value reported in column 4. Standard errors are adjusted for clustering at the class level. High school tracks in Grade 11 are predicted for each student using the coefficients from a linear regression of the corresponding binary variable (e.g., enrollment in the general science track) on all the student characteristics listed in the table. This model is fitted on the sample of students in the control group.

Table I2 – Treatment Effects (ITT) for Grade 12 Students: Regular vs. Augmented Slides

	Girls		Boys	
	(1)	(2)	(3)	(4)
Major: STEM				
Treatment group indicator (T)	0.022 (0.019)	0.003 (0.021)	0.013 (0.030)	0.007 (0.031)
T *Augmented slides	0.001 (0.026)	0.036 (0.031)	−0.018 (0.037)	0.008 (0.042)
Undergraduate major: male-dominated STEM				
Treatment group indicator (T)	0.044*** (0.016)	0.030* (0.018)	0.017 (0.031)	0.007 (0.030)
T *Augmented slides	−0.015 (0.023)	0.006 (0.026)	−0.002 (0.037)	0.024 (0.040)
Undergraduate major: selective STEM				
Treatment group indicator (T)	0.041** (0.017)	0.033* (0.017)	0.025 (0.024)	0.018 (0.025)
T *Augmented slides	−0.018 (0.022)	−0.001 (0.024)	−0.012 (0.030)	0.010 (0.035)
Science-related jobs pay higher wages				
Treatment group indicator (T)	0.005 (0.026)	0.017 (0.031)	0.003 (0.023)	0.040* (0.021)
T *Augmented slides	0.097*** (0.034)	0.069 (0.044)	0.041 (0.031)	−0.040 (0.030)
Positive perceptions of science-related careers (index)				
Treatment group indicator (T)	0.254*** (0.055)	0.294*** (0.054)	0.154*** (0.046)	0.184*** (0.048)
T *Augmented slides	0.078 (0.069)	0.022 (0.077)	−0.013 (0.064)	−0.061 (0.077)
Equal aptitudes of M and W for science (index)				
Treatment group indicator (T)	0.060 (0.038)	0.077* (0.042)	0.060 (0.049)	0.054 (0.055)
T *Augmented slides	0.056 (0.055)	0.026 (0.064)	0.112 (0.074)	0.117 (0.085)
Month of Visit FE	No	Yes	No	Yes
Month of Visit FE * treatment group indicator	No	Yes	No	Yes
N	2,827	2,827	2,924	2,924

Notes: This table reports estimates of the treatment effects (ITT) of the role model interventions on student outcomes for Grade 12 students, separately by gender and by the type of slides (regular or augmented) that were provided to the female role model who visited the classroom. For each outcome of interest, the reported coefficients are obtained from a regression of the outcome of interest on a treatment group indicator (T) and the interaction between this indicator and an indicator that takes the value one if the role model was provided with the augmented set of slides. The specification in columns 1 and 3 includes school fixed effects to account for the fact that randomization was stratified by school. Columns 2 and 4 further include month-of-visit fixed effects to account for the fact that the additional slides were provided slightly later in the experiment. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

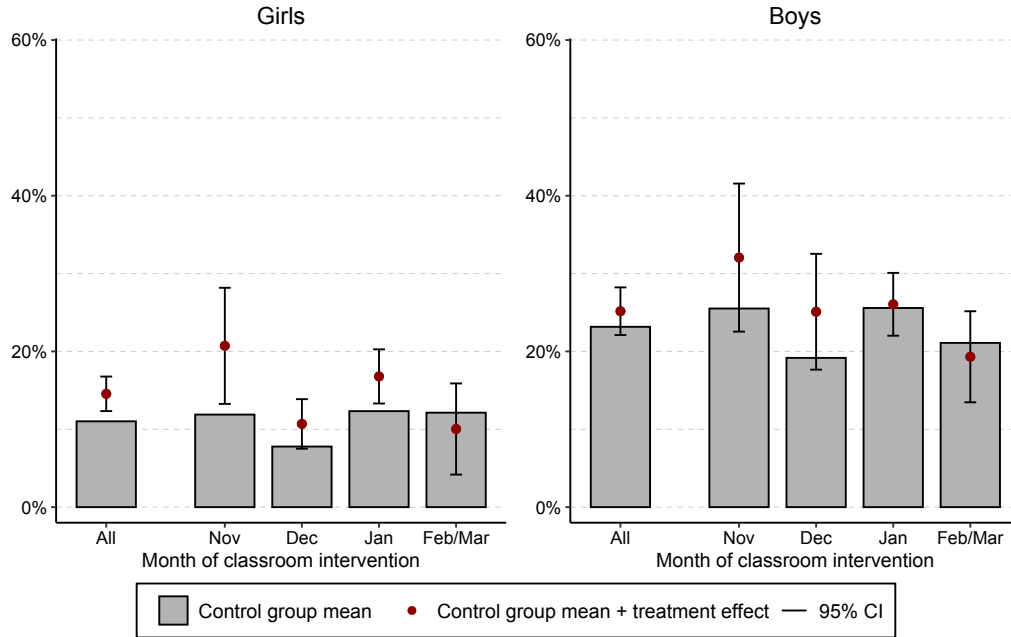
I.2 Persistence of the Effects and Timing of Visits

Table I3 – Persistence of Effects on Student Perceptions

	Girls			Boys		
	Months since intervention			Months since intervention		
	1 to 2 months (1)	3 to 4 months (2)	5 to 6 months (3)	1 to 2 months (4)	3 to 4 months (5)	5 to 6 months (6)
Panel A. Grade 10						
Positive perceptions of science-related careers (index)	0.413*** (0.057)	0.200*** (0.037)	0.143* (0.077)	0.192*** (0.053)	0.168*** (0.036)	0.049 (0.083)
More men in science-related jobs	0.170*** (0.021)	0.154*** (0.017)	0.164*** (0.033)	0.209*** (0.022)	0.163*** (0.018)	0.116*** (0.039)
Equal gender aptitude for math (index)	0.179*** (0.047)	0.101*** (0.032)	0.019 (0.065)	0.244*** (0.053)	0.122*** (0.040)	0.090 (0.069)
Women don't really like science	0.047** (0.022)	0.067*** (0.014)	0.041 (0.026)	0.131*** (0.020)	0.107*** (0.016)	0.017 (0.040)
W face discrimination in science-related careers	0.158*** (0.022)	0.135*** (0.017)	0.081** (0.039)	0.162*** (0.026)	0.174*** (0.017)	0.110*** (0.036)
Taste for science subjects (index)	0.088 (0.075)	−0.035 (0.043)	−0.053 (0.075)	0.043 (0.058)	−0.008 (0.041)	0.043 (0.072)
Math self-concept (index)	−0.029 (0.057)	0.006 (0.039)	0.044 (0.080)	−0.041 (0.063)	0.103*** (0.039)	0.088 (0.090)
Science-related career aspirations (index)	0.088 (0.057)	−0.002 (0.036)	0.010 (0.062)	−0.000 (0.051)	0.022 (0.038)	0.010 (0.072)
N	1,729	3,716	831	1,677	3,318	693
Panel B. Grade 12 (science track)						
Positive perceptions of science-related careers (index)	0.442*** (0.053)	0.253*** (0.043)	0.353*** (0.118)	0.182*** (0.061)	0.169*** (0.044)	0.003 (0.095)
More men in science-related jobs	0.128*** (0.031)	0.107*** (0.019)	0.208*** (0.060)	0.114*** (0.023)	0.159*** (0.021)	0.208*** (0.046)
Equal gender aptitude for math (index)	0.077 (0.067)	0.138*** (0.033)	0.020 (0.094)	0.218*** (0.081)	0.106** (0.051)	0.044 (0.123)
Women don't really like science	0.067*** (0.021)	0.040*** (0.011)	0.032* (0.018)	0.042 (0.029)	0.077*** (0.019)	0.144*** (0.032)
W face discrimination in science-related jobs	0.104*** (0.038)	0.102*** (0.023)	0.087 (0.072)	0.083*** (0.027)	0.085*** (0.024)	−0.011 (0.062)
Taste for science subjects (index)	−0.063 (0.071)	−0.028 (0.045)	0.258*** (0.060)	0.030 (0.079)	0.010 (0.049)	−0.090 (0.111)
Math self-concept (index)	0.043 (0.065)	0.001 (0.053)	0.169 (0.122)	−0.022 (0.054)	0.114** (0.046)	0.126 (0.149)
Science-related career aspirations (index)	−0.005 (0.077)	0.123*** (0.045)	0.231*** (0.077)	0.007 (0.046)	0.048 (0.046)	0.098 (0.118)
N	689	1,468	394	717	1,514	370

Notes: This table reports estimates of the treatment effects of the role model interventions on student perceptions, separately by grade level, gender, and intervals of elapsed time between the classroom intervention and the student survey. The sample is restricted to students who completed the post-intervention questionnaire. Each coefficient is obtained from a linear regression of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class).

(a) Enrollment in selective STEM



(b) Enrollment in male-dominated STEM

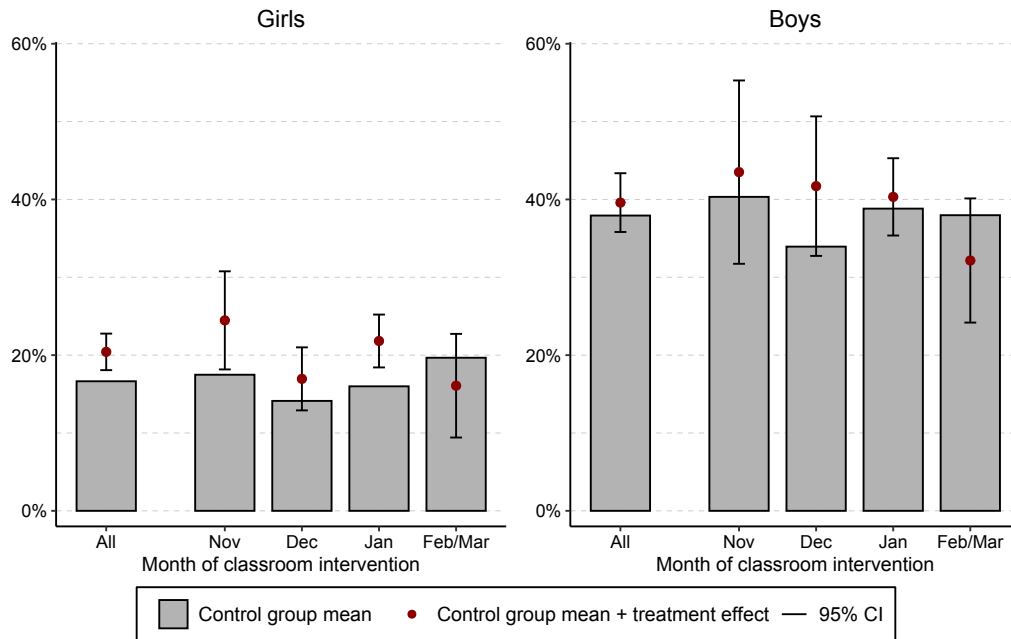


Figure I3 – Grade 12 Students: Enrollment in Selective and Male-Dominated STEM Undergraduate Programs, by Gender and Month of Classroom Intervention

Notes: The figure shows the fraction of Grade 12 (science track) students who enrolled in selective (Panel A) and in male-dominated (Panel B) STEM undergraduate programs after graduating for high school, separately for girls (left panel) and for boys (right panel). The filled bars indicate the baseline enrollment rates among students in the control group, both overall and separately by month of classroom intervention. The solid dots show the estimated treatment effects (added to the control group means) with 95 percent confidence intervals denoted by vertical capped bars. The treatment effects are estimated from separate regressions of the outcome of interest on a classroom visit indicator, using treatment assignment as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors are adjusted for clustering at the unit of randomization (class).

J Spillover Effects

This appendix investigates whether the program could have had spillover effects for students who were not exposed to the role model interventions in participating schools. Section J.1 provides survey evidence suggesting that the scope for spillover effects was relatively limited. Section J.2 describes the difference-in-differences (DiD) approach that we use to estimate the magnitude of spillovers, the results of which point to non-statistically significant effects.

J.1 Survey Evidence

To get some sense of the scope for spillover effects in the context of our study, we included in the last section of the survey a series of questions asking students in the treatment group whether they had talked about the classroom interventions with their classmates, with schoolmates from other classes, or with friends from other schools. We also asked students in the control group whether they had heard about a science-related awareness-raising program and, more specifically, whether they knew about other classes in the school being visited by a female or male scientist.

Overall, the summary statistics from the survey data suggest relatively limited opportunities for spillover effects (see Table J1). In the treatment group, 58 percent of Grade 10 students and 63 percent of science track Grade 12 students report having talked about the classroom intervention with their classmates, but only 24 percent (27 percent) with schoolmates from other classes, and 20 percent with students from other schools. Interestingly, these proportions are higher for girls than for boys: in Grade 10, 66 percent of girls in the treatment group report having discussed the program with their classmates and 28 percent with schoolmates from other classes vs. respectively 50 percent and 20 percent among boys; in Grade 12, 70 percent of girls in the treatment group report having discussed the program with their classmates and 33 percent with schoolmates from other classes vs. respectively 56 percent and 21 percent among boys.

In the control group, only 14 percent of students in Grade 10 report having heard of classroom visits in other classes, mostly in a vague manner (12 percent). In Grade 12, students in the control group are more likely to report being at least vaguely aware of such visits (34 percent), but less than 5 percent of boys and girls have a precise recollection. Gender differences in these proportions are small and barely statistically significant. The fact that students in Grade 12 are more likely to report being aware of classroom visits could be at least partly due to the fact that the share of students assigned to the treatment group among all students from the same grade level was typically larger in Grade 12 than in Grade 10, on average 32 percent vs. 25 percent. Despite these differences, the overall picture that emerges from the survey is that students in the control group had only limited awareness of the classroom interventions in other classes.

J.2 Differences-in-Differences Estimates of Spillover Effects

We complement the survey evidence by investigating more formally whether the role model interventions could have affected the higher education choices of Grade 12 students whose classes were not assigned to the treatment group. These students are either in the classes that were not selected by school principals to participate in the program evaluation or in the participating classes that were randomly assigned to the control group.

Our experimental design does not include a “super control” group composed of students enrolled in schools randomly chosen to have zero probability of assignment to the treatment among the classes selected by school principals. Spillover effects cannot, therefore, be identified by comparing the control group classes in participating schools with such supercontrol group classes, as in the design pioneered by Duflo and Saez (2003).^{A.7} Instead, our approach builds on

^{A.7}Vazquez-Bare (2018) develops a potential-outcome-based nonparametric framework to identify spillover

the following intuition: for schools that participated in the evaluation, the random assignment of treatment to participating classes makes it possible to estimate the average outcome that would have been observed if *all* students from these schools had only been exposed to the spillover effects of role model interventions, without being *directly* exposed to a female role model. This unobserved “spillover-only” counterfactual can be estimated at the school level using an appropriately weighted average of non-treated classes: it suffices to compute the weighted average outcome of students in the non-participating classes and in the participating classes that were randomly assigned to the control group, with respective weights equal to the share of participating and of non-participating classes in the school. Average spillover effects can then be estimated by comparing this “spillover-only” counterfactual to a “no-treatment” counterfactual. This second counterfactual is constructed under the assumption that absent treatment, mean outcomes in participating school would have followed the same evolution as in non-participating schools. Having verified that this common trends assumption is satisfied in the pre-treatment period 2012–2014, we implement a difference-in-differences estimator that identifies the difference between the “spillover-only” and the “no-treatment” counterfactuals. This approach, which is graphically illustrated in Figure J1, enables us to estimate the average spillover effects of role model interventions in the participating schools.

Notations. We are interested in measuring the spillover effects of classroom visits. We denote by D_s a binary indicator for a student’s school s being visited by a female role model and by D_{cs} a binary indicator for a role model intervention taking place in the student’s class c . We consider two time periods, represented by a binary indicator $T \in \{0, 1\}$, with classroom visits taking place in period 1 only. For a given realization of the treatment assignment (d_s, d_{cs}) , the potential outcome for student i in school s , class c , and time t is denoted by $Y_{icst}(d_s, d_{cs})$.

We use the binary indicator G_s to indicate whether school s participated in the experiment and we denote the sets of participating and non-participating schools by \mathcal{S}_1 and \mathcal{S}_0 , respectively. The number of participating (non-participating) schools is denoted by M_1 (M_0). Only a subset of the classes in participating schools were (non-randomly) selected by the principals to participate in the experiment in period 1. The participation status of class c in school s is denoted by the binary indicator G_{cs} . Among participating classes ($G_{cs} = 1$), the binary indicator R_{cs} indicates whether the class was randomly assigned to the treatment group ($R_{cs} = 1$) or to the control group ($R_{cs} = 0$). The experimental setting therefore implies that $D_s = G_s \times T$ and $D_{cs} = R_{cs} \times T$. A student’s observed outcome can then be written

$$Y_{icst} = D_s \cdot D_{cs} \cdot Y_{icst}(1, 1) + D_s \cdot (1 - D_{cs}) \cdot Y_{icst}(1, 0) + (1 - D_s) \cdot Y_{icst}(0, 0). \quad (\text{A.3})$$

To simplify notation, we assume that each school has the same number of students, N , and that the number of students is the same in both periods.

Let $\bar{Y}_{s,t}(0, 0)$ denote the average *potential* outcome of students in school s and year t under no treatment. This average potential outcome corresponds to the case in which no student from school s in year t is exposed to either the direct or spillover effects of classroom visits, i.e.,

$$\bar{Y}_{s,t}(0, 0) = \frac{1}{N} \sum_{i=1}^N Y_{icst}(0, 0). \quad (\text{A.4})$$

Let $\bar{Y}_{s,t}(1, 0)$ denote the average *potential* outcome of students in school s and year t in the (non-feasible) scenario in which all students in school s are only exposed to the spillover effects

effects in randomized experiments where units are clustered, without requiring a specific experimental design. This approach, however, cannot be easily adapted to our setting since it requires that the treatment is assigned at the individual level within clusters (schools), not at the group level (classes), in order to exploit variation in all the possible configurations of own and neighbors’ observed treatment assignments.

of role model interventions in other classes, without themselves being visited by a female role model. This “spillover-only” average potential outcome is defined as follows:

$$\bar{Y}_{s,t}(1,0) = \frac{1}{N} \sum_{i=1}^N Y_{icst}(1,0). \quad (\text{A.5})$$

Our parameter of interest is the expected average spillover effect of classroom visits for the students in participating schools in period 1, i.e.,

$$\Delta = \mathbb{E} \left(\frac{1}{M_1} \sum_{s \in \mathcal{S}_1} \left(\bar{Y}_{s,1}(1,0) - \bar{Y}_{s,1}(0,0) \right) \right). \quad (\text{A.6})$$

This parameter can be interpreted as the average effect for students in participating schools of being only exposed to the indirect effects of classroom visits compared to the counterfactual of no classroom visit in the school.

Identification of spillover effects. Let $\bar{Y}_{s,t}$ denote the mean *observed* outcome for students in school s and year t , i.e.,

$$\bar{Y}_{s,t} = \frac{1}{N} \sum_{i=1}^N Y_{icst}. \quad (\text{A.7})$$

For non-participating schools in periods 0 and 1 and for participating schools in period 0, this mean observed outcome is in expectation equal to the expected average potential outcome under no treatment. Indeed, Equations (A.3), (A.4), and (A.7) imply that

$$\mathbb{E}(\bar{Y}_{s,t}) = \mathbb{E}(\bar{Y}_{s,t}(0,0)) \text{ if } s \in \mathcal{S}_0 \text{ and } t \in \{0,1\} \text{ or if } s \in \mathcal{S}_1 \text{ and } t = 0. \quad (\text{A.8})$$

For each school $s \in \mathcal{S}_1$ that participated in the evaluation, we consider the following partition of students in period 1: let \mathcal{C}_s^0 , \mathcal{C}_s^C , and \mathcal{C}_s^T denote respectively (i) the students in the classes that did not participate in the evaluation ($G_s = 0$), (ii) the students in the participating classes that were randomly assigned to the control group ($G_s = 1$ and $R_{cs} = 0$), and (iii) the students in the participating classes that were randomly assigned to the treatment group ($G_s = 1$ and $R_{cs} = 1$). By definition, the number of students in each group, which we denote by N_s^0 , N_s^C and N_s^T , respectively, is such that $N = N_s^0 + N_s^C + N_s^T$.

For the purpose of estimating spillover effects, we construct a mean counterfactual outcome for participating schools in period 1, which we denote by $\tilde{Y}_{s,1}$. As shown in Proposition 1 below, the expected value of $\tilde{Y}_{s,1}$ coincides with the expected average potential outcome of students in school s and period 1, had all of its students only been exposed to the spillover effects of classroom visits in other classes, without being themselves directly exposed to a female role model. This counterfactual outcome ignores classes in the treatment group and is defined as a weighted average of the observed outcomes of students in the non-participating classes and the control group classes (see Figure J1):

$$\tilde{Y}_{s,1} = \frac{1}{N} \left(\sum_{i \in \mathcal{C}_s^0} Y_{ics1} + \left(1 + \frac{N_s^T}{N_s^C} \right) \sum_{i \in \mathcal{C}_s^C} Y_{ics1} \right), \quad s \in \mathcal{S}_1. \quad (\text{A.9})$$

The intuition is as follows. The “spillover only” counterfactual measured at the school level cannot be recovered from the non-participating classes only, since these classes were not randomly selected by school principals. However, having noted that the mean observed outcome of students in the control group is an unbiased estimator of the mean (unobserved)

“spillover-only” outcome for students in the treatment group, one can reconstruct the school-level “spillover-only” counterfactual by restricting the set of students to those in non-participating classes and control group classes. To estimate the mean outcome that would have been observed if all students had only been exposed to the spillover effects of classroom visits, it suffices to reweight students in the control group so that they match the total number of students in the participating classes (i.e., treatment and control), and then combine this reweighted sample with the sample of students in non-participating classes to compute the average outcome.

Assumption 1. *Random assignment of treatment to participating classes.*

$$\mathbb{E} \left(\frac{1}{N_s^T} \sum_{i \in \mathcal{C}_s^T} Y_{ics1}(1, 0) \right) = \mathbb{E} \left(\frac{1}{N_s^C} \sum_{i \in \mathcal{C}_s^C} Y_{ics1}(1, 0) \right), \quad s \in \mathcal{S}_1.$$

Assumption 1 states that students in the treatment and control group classes of participating schools have the same expected average potential outcome under the “spillover-only” treatment. Our experimental design ensures that this assumption is satisfied.

Proposition 1. *Under Assumption 1, the counterfactual $\tilde{Y}_{s,1}$ is an unbiased estimator of the expected average potential outcome of students in participating school s and period 1 under the “spillover-only” treatment, $\bar{Y}_{s,1}(1, 0)$:*

$$\mathbb{E}(\tilde{Y}_{s,1}) = \mathbb{E}(\bar{Y}_{s,1}(1, 0)), \quad s \in \mathcal{S}_1.$$

Proof. From the definition of the “spillover-only” counterfactual in Equation (A.9), we have

$$\begin{aligned} \mathbb{E}(\tilde{Y}_{s,1}) &= \mathbb{E} \left(\frac{1}{N} \left(\sum_{i \in \mathcal{C}_s^0} Y_{ics1} + \left(1 + \frac{N_s^T}{N_s^C} \right) \sum_{i \in \mathcal{C}_s^C} Y_{ics1} \right) \right) \\ &= \frac{1}{N} \left(\sum_{i \in \mathcal{C}_s^0} \mathbb{E}(Y_{ics1}(1, 0)) + \sum_{i \in \mathcal{C}_s^C} \mathbb{E}(Y_{ics1}(1, 0)) + \frac{N_s^T}{N_s^C} \sum_{i \in \mathcal{C}_s^C} \mathbb{E}(Y_{ics1}(1, 0)) \right) \\ &= \frac{1}{N} \left(\sum_{i \in \mathcal{C}_s^0} \mathbb{E}(Y_{ics1}(1, 0)) + \sum_{i \in \mathcal{C}_s^C} \mathbb{E}(Y_{ics1}(1, 0)) + \sum_{i \in \mathcal{C}_s^T} \mathbb{E}(Y_{ics1}(1, 0)) \right) \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}(Y_{ics1}(1, 0)) \\ &= \mathbb{E}(\bar{Y}_{s,1}(1, 0)). \end{aligned}$$

The second equality follows from Equation (A.3), the third equality follows from Assumption 1, while the last equality follows from Equation (A.5). The key intuition for this result is that by virtue of the random assignment of treatment to participating classes, the mean observed outcome of students assigned to the control group is an unbiased estimator of the mean unobserved “spillover-only” outcome of students assigned to the treatment group. \square

Identifying spillover effects requires comparing the “spillover-only” counterfactual with the “no-treatment” counterfactual. To this end, we define the following difference-in-differences estimator, which we denote by $\hat{\Delta}$:

$$\hat{\Delta} = \frac{1}{M_1} \sum_{s \in \mathcal{S}_1} (\tilde{Y}_{s,1} - \bar{Y}_{s,0}) - \frac{1}{M_0} \sum_{s \in \mathcal{S}_0} (\bar{Y}_{s,1} - \bar{Y}_{s,0}). \quad (\text{A.10})$$

This estimator compares the evolution of the mean outcome of students in participating schools

between period 0 and period 1 (using the “spillover-only” counterfactual for period 1) with the corresponding evolution in non-participating schools.

Assumption 2. *Common trends between participating and non-participating schools.*

$$\mathbb{E} \left(\frac{1}{M_1} \sum_{s \in \mathcal{S}_1} (\bar{Y}_{s,1}(0,0) - \bar{Y}_{s,0}(0,0)) \right) = \mathbb{E} \left(\frac{1}{M_0} \sum_{s \in \mathcal{S}_0} (\bar{Y}_{s,1}(0,0) - \bar{Y}_{s,0}(0,0)) \right).$$

Assumption 2 states that in the absence of role model visits to the school, average outcomes in participating and non-participating schools would have followed parallel trends. Although this assumption cannot be directly tested, it can be indirectly assessed by comparing the evolution of mean outcomes in participating and non-participating schools in the pre-intervention period.

Proposition 2. *Under Assumptions 1 and 2, $\hat{\Delta}$ is an unbiased estimator of the average spillover effect, Δ :*

$$\mathbb{E}(\hat{\Delta}) = \Delta.$$

Proof. From the definition of the difference-in-differences estimator in Equation (A.10), we have

$$\begin{aligned} \mathbb{E}(\hat{\Delta}) &= \mathbb{E} \left(\frac{1}{M_1} \sum_{s \in \mathcal{S}_1} (\tilde{Y}_{s,1} - \bar{Y}_{s,0}) - \frac{1}{M_0} \sum_{s \in \mathcal{S}_0} (\bar{Y}_{s,1} - \bar{Y}_{s,0}) \right) \\ &= \mathbb{E} \left(\frac{1}{M_1} \sum_{s \in \mathcal{S}_1} (\bar{Y}_{s,1}(1,0) - \bar{Y}_{s,0}(0,0)) \right) - \mathbb{E} \left(\frac{1}{M_0} \sum_{s \in \mathcal{S}_0} (\bar{Y}_{s,1}(0,0) - \bar{Y}_{s,0}(0,0)) \right) \\ &= \mathbb{E} \left(\frac{1}{M_1} \sum_{s \in \mathcal{S}_1} (\bar{Y}_{s,1}(1,0) - \bar{Y}_{s,0}(0,0)) \right) - \mathbb{E} \left(\frac{1}{M_1} \sum_{s \in \mathcal{S}_1} (\bar{Y}_{s,1}(0,0) - \bar{Y}_{s,0}(0,0)) \right) \\ &= \mathbb{E} \left(\frac{1}{M_1} \sum_{s \in \mathcal{S}_1} (\bar{Y}_{s,1}(1,0) - \bar{Y}_{s,1}(0,0)) \right) \\ &= \Delta. \end{aligned}$$

The second equality follows from Equation (A.8) and from Proposition 1, the third equality follows from Assumption 2 (common trends between participating and non-participating schools), while the last equality follows from Equation (A.6). \square

Empirical specification. In the context of our study, the spillover effects estimator (A.10) can be conveniently implemented using a difference-in-differences regression specification. We apply this estimator to investigate whether the classroom interventions affected the college decisions of science track Grade 12 students whose classes were not visited by a female role model.

In our empirical application, we consider the four cohorts of Grade 12 students that were enrolled in the high schools of the Paris region in the year of the intervention (2015) and in the three preceding years (2012, 2013, and 2014).

One complication is that the “For Girls in Science” program was first implemented on a small scale in 2014, i.e., one year before the evaluation was conducted. As a result, some of the schools that participated in the program evaluation in 2015, as well as some of the schools that did not participate in the evaluation, could have been visited by female role models in 2014. Although we cannot precisely identify these schools, the contamination effect is likely to be small since the interventions were carried out by a small number of role models and were not specifically targeted at students enrolled in Grade 10 and Grade 12 (science track). Nonetheless, to ensure

that our difference-in-differences estimates are not biased due to these prior interventions, we use 2012 as the reference year. The baseline differences between participating and non-participating schools are therefore measured at a point in time in which the program was not in place.

Let $\bar{Y}_{s,t}$ denote the average outcome of Grade 12 students in school s and year t . For each participating school $s \in \mathcal{S}_1$, we use Equation (A.9) to construct the “spillover-only” mean counterfactual outcome in 2015, which we denote by $\tilde{Y}_{s,t}$. Our dependent variable, denoted by $\bar{Y}_{s,t}^*$, is then defined as follows:

$$\bar{Y}_{s,t}^* = \begin{cases} \tilde{Y}_{s,t} & \text{if } s \in \mathcal{S}_1 \text{ and } t = 2015 \\ \bar{Y}_{s,t} & \text{otherwise} \end{cases}$$

The spillover effects of classroom visits are then estimated using the following difference-in-differences regression model:

$$\bar{Y}_{s,t}^* = \alpha + \theta_s + \theta_t + \sum_{k=2013}^{2015} \beta_k \cdot \mathbf{1}\{s \in \mathcal{S}_1 \text{ and } t = k\} + \epsilon_{s,t}, \quad (\text{A.11})$$

where θ_s are school fixed effects and θ_t are year fixed effects (using 2012 as the reference year); $\mathbf{1}\{s \in \mathcal{S}_1 \text{ and } t = k\}$ is a dummy variable that takes the value one if the observation corresponds to a participating school observed in year k ; and $\epsilon_{s,t}$ is the error term. Under the common trend assumption, the coefficient $\hat{\beta}_{2015}$ identifies the average spillover effects among the non-treated students in participating schools. The coefficients $\hat{\beta}_{2013}$ and $\hat{\beta}_{2014}$ provide an indirect test of this assumption: if it holds, the evolution of mean outcomes between 2012 and 2014 (pre-intervention period) should be parallel between participating and non-participating schools, and the coefficients on the pre-interventions “placebos” should not be jointly significant.^{A.8}

Selection of non-participating schools. To ensure that non-participating schools are as similar as possible to the participating schools, we use a nearest neighbor matching procedure (with replacement) on the estimated propensity score. We consider all public and private high schools operating in the Paris region that had at least two science track Grade 12 classes in 2015, as this restriction was used in our experimental design to select participating schools (see Section 2 in the main text). We then estimate the probability that the school participated in the experiment in 2015 given a vector of exogenous school characteristics \mathbf{X}_{st} (measured every year between 2012 and 2015) and a vector of the pre-intervention outcomes \mathbf{Y}_{st} (measured in 2012 and 2013) for which spillover effects are measured.^{A.9} We then match each participating school with the non-participating school having the closest propensity score among the schools with the same status (public or private) and located in the same education district (Paris, Créteil or Versailles) as that of the participating school.

^{A.8}Strictly speaking, the parallel trend assumption only requires the coefficient β_{2013} to be non-statistically significant since, as explained above, the comparison between participating and non-participating schools in 2014 could be contaminated by the classroom interventions that were carried on a small scale that year. As shown below, the results show that the parallel trend assumption also holds between 2013 and 2014, suggesting that the contamination effects of these prior interventions are negligible, if any.

^{A.9}The vector of exogenous school characteristics \mathbf{X}_{st} includes the school’s education district (Paris, Créteil or Versailles), whether it is public or private, and the following time-varying characteristics every year between 2012 and 2015: the number of students in Grade 12 (science track), the fraction of female students, and the fraction of high-SES students. The vector of pre-intervention outcomes \mathbf{Y}_{st} in 2012 and 2013 includes the fraction of science track Grade 12 students who enrolled in a STEM program after graduating from high school, the fraction who enrolled in a selective STEM program, and the fraction who enrolled in a male-dominated STEM program (computed separately by year and gender). We do not control for pre-intervention outcomes in 2014 to avoid any contamination by classroom interventions that could have been carried out that year.

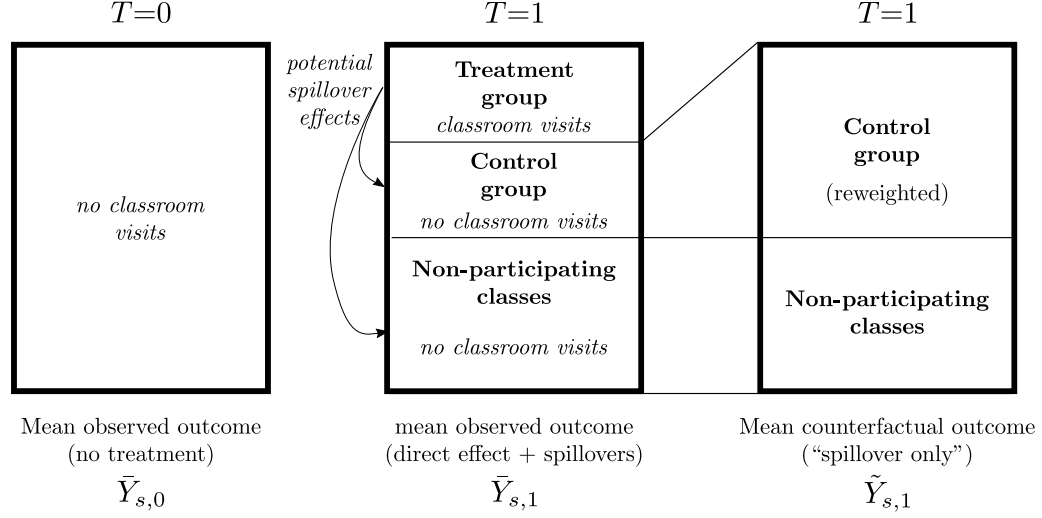
Results. We use Equation (A.11) to estimate the spillover effects of classroom visits on the college enrollment outcomes of Grade 12 students in non-treated classes. The model is estimated separately by gender and we consider the three main outcomes for which we document significant direct effects of the interventions: enrollment in a STEM undergraduate program, enrollment in a selective STEM program, and enrollment in a male-dominated STEM program. The observations are school-by-year averages weighted by school size. Standard errors are clustered at the school level to account for serial correlation across years.

The results are reported in Table J2. Panel A shows that the non-participating schools selected by the nearest-neighbor matching procedure are reasonably similar to the participating schools in terms of the average college enrollment outcomes of female and male students in the pre-intervention period 2012-2013.

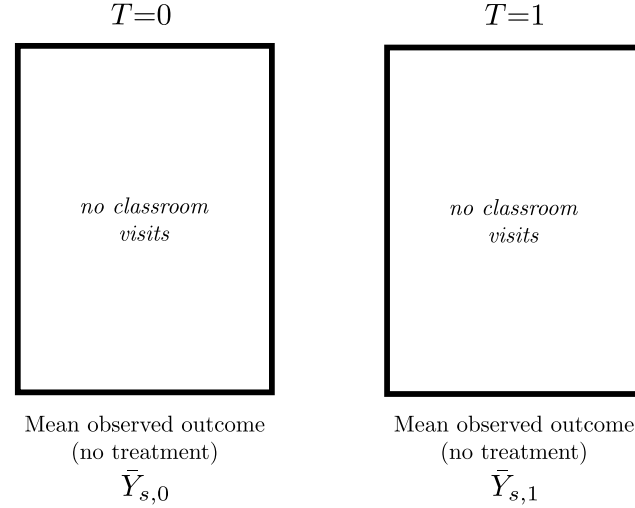
The estimates from the DiD regression are reported in Panel B. In all specifications, the coefficients on (participating school $\times t=2013$) and on (participating school $\times t=2014$) are close to zero and are neither individually nor jointly significant, which lends support to the assumption of common trends between participating and non-participating schools. Overall, the results provide no evidence of significant spillover effects from the classroom visits in participating schools: for all considered outcomes, the coefficient $\hat{\beta}_{2015}$ on (participating school $\times t = 2015$) is close to zero and not statistically significant for both female and male students.

It should, however, be noted that although our estimates are relatively precise, we cannot rule out small to moderate spillover effects. In the presence of positive spillovers, the treatment effects reported in the main text would under-estimate the true direct effect of classroom visits, since the “contamination” of the control group would push the difference between treatment and control classes towards zero. Denoting by Φ the average direct effect of the classroom interventions and by $\Delta (> 0)$ their average indirect effect (through spillovers), the treatment-control difference in mean outcomes, denoted by $\hat{\beta}$, estimates $\Phi - \Delta$ instead of Φ . If we estimate the spillover effects to be at most $\hat{\Delta}^{UB}$, this implies that the size of spillover effects is at most $\hat{\Delta}^{UB}/(\hat{\beta} + \hat{\Delta}^{UB})$ of the size of the direct effect. When we consider the effects on the probability that female students enroll in a selective STEM program, the comparison of treatment and control classes yields an estimated direct effect of $\hat{\beta} = 0.035$ (see Table 5, Panel B, column 2). Based on the results in column 2 of Table J2, the upper bound of the 95 percent confidence interval for the spillover effects is estimated to be $\hat{\Delta}^{UB} = 0.017$. Hence, in the case of selective STEM enrollment, we cannot reject spillover effects that would be at most 33 percent of the size of the “true” direct effect $\hat{\beta} + \hat{\Delta}^{UB}$, which in this case would be of 5.2 percentage points. A similar calculation for the spillover effects on male-dominated STEM enrollment yields an upper bound of $\hat{\Delta}^{UB} = 0.025$. Since the estimated direct effect is $\hat{\beta} = 0.038$, we cannot reject spillover effects of at most 40 percent of the size of the “true” direct effect $\hat{\beta} + \hat{\Delta}^{UB}$, which in that case would be of 6.3 percentage points.

M_1 participating schools ($s \in \mathcal{S}_1$):



M_0 non-participating schools ($s \in \mathcal{S}_0$):



Difference-in-differences estimator of spillover effects:

$$\hat{\Delta} = \frac{1}{M_1} \sum_{s \in \mathcal{S}_1} (\tilde{Y}_{s,1} - \bar{Y}_{s,0}) - \frac{1}{M_0} \sum_{s \in \mathcal{S}_0} (\bar{Y}_{s,1} - \bar{Y}_{s,0})$$

Figure J1 – Spillover Effects of Role Model Interventions: Empirical Strategy

Notes: This figure illustrates the difference-in-differences strategy we implement to estimate the spillover effects of role model interventions for students who were enrolled in participating schools but whose classes were not assigned to the treatment group. These students are either in the classes that were not selected by school principals to participate in the program evaluation or in the participating classes that were randomly assigned to the control group. Our approach consists in comparing the evolution of mean student outcomes (at the school level) in participating ($s \in \mathcal{S}_1$) and non-participating schools ($s \in \mathcal{S}_0$), between the year before the intervention ($T = 0$) and the year of the intervention ($T = 1$). For $T = 1$, we use a weighted average of non-treated classes in each participating school to estimate the counterfactual “spillover-only” outcome that would have been observed if all the students from that school had only been exposed to the spillover effects of classroom interventions, without being directly exposed to a female role model. Average spillover effects are then estimated by comparing this “spillover-only” counterfactual to a “no-treatment” counterfactual. Under the assumption that absent treatment, mean outcomes in participating school would have followed the same evolution as in non-participating schools, the average spillover effects can be estimated by comparing the evolution between $T = 0$ and $T = 1$ of the mean outcome of students in participating schools (using the “spillover-only” counterfactual for period 1) with the corresponding evolution in non-participating schools.

Table J1 – Scope for Spillover Effects: Summary Statistics from the Student Survey

	All	Boys	Girls	Within class	
				Difference (3)–(2)	<i>p</i> -value of diff.
	(1)	(2)	(3)	(4)	(5)
Panel A. Grade 10					
<i>Treatment Group</i>					
Discussed the classroom visit?					
with classmates	0.580	0.498	0.656	0.145	0.000
with other students from the school	0.240	0.200	0.277	0.072	0.000
with other students outside the school	0.203	0.155	0.247	0.098	0.000
Exposed to other science outreach program?					
this school year	0.128	0.138	0.120	–0.011	0.297
in the past	0.182	0.218	0.149	–0.059	0.000
N	6,245	2,989	3,256		
<i>Control Group</i>					
Heard of classroom visits in other classes?					
Yes, definitely	0.018	0.017	0.020	0.001	0.862
Yes, vaguely	0.122	0.117	0.127	0.009	0.244
No	0.859	0.866	0.853	–0.010	0.271
Exposed to programs about science or jobs in science?					
this school year	0.146	0.144	0.148	0.011	0.283
before the end of this school year	0.052	0.059	0.047	–0.014	0.019
in the past	0.322	0.309	0.333	0.025	0.066
N	5,981	2,762	3,219		
Panel B. Grade 12 (science track)					
<i>Treatment Group</i>					
Discussed the classroom visit?					
with classmates	0.629	0.556	0.705	0.131	0.000
with other students from the school	0.269	0.206	0.334	0.114	0.000
with other students outside the school	0.202	0.133	0.275	0.136	0.000
Exposed to other science outreach programs?					
this school year	0.202	0.200	0.204	0.005	0.797
in the past	0.324	0.349	0.299	–0.053	0.025
N	2,642	1,350	1,292		
<i>Control Group</i>					
Heard of classroom visit in other classes?					
Yes, definitely	0.047	0.049	0.045	–0.004	0.645
Yes, vaguely	0.292	0.275	0.308	0.037	0.048
No	0.661	0.676	0.646	–0.033	0.085
Exposed to programs about science or jobs in science?					
this school year	0.287	0.291	0.284	0.011	0.515
before the end of this school year	0.096	0.104	0.088	–0.009	0.403
in the past	0.488	0.461	0.514	0.054	0.028
N	2,594	1,286	1,308		

Notes: The summary statistics in this table are computed from the post-treatment student survey administered in all participating classes between one and six months after the role model interventions. Columns 1, 2, and 3 report average values for all respondents and for boys and girls, respectively, separately by grade level and treatment assignment. The within-class difference in the responses of girls and boys, reported in column 4, is obtained from a regression of the variable of interest on a female dummy, controlling for class fixed effects and clustering standard errors at the school level. The associated *p*-value is reported in column 5.

Table J2 – Difference-in-Differences Estimates of the Spillover Effects of Role Model Interventions on College Enrollment Outcomes, Grade 12 Students, Years 2012–2015

	Grade 12 (science track) students					
	Girls			Boys		
	Underg. STEM (1)	Selective STEM (2)	Male-dom. STEM (3)	Underg. STEM (4)	Selective STEM (5)	Male-dom. STEM (6)
Panel A. Baseline means (2012–2013)						
<i>Participating schools</i>						
Mean	0.274	0.145	0.163	0.489	0.265	0.409
Number of schools	88	88	88	87	87	87
Average number of Grade 12 students	107	107	107	108	108	108
<i>Non-participating schools</i>						
Mean	0.265	0.141	0.157	0.473	0.257	0.395
Number of schools	62	62	62	61	61	61
Average number of Grade 12 students	99	99	99	99	99	99
Panel B. Regression estimates						
<i>Pre-trends: participating vs. non-particip. schools (relative to 2012)</i>						
$\hat{\beta}_{2013}$: Particip. school \times ($t=2013$)	0.006 (0.017)	−0.001 (0.014)	0.013 (0.014)	0.003 (0.022)	−0.023 (0.017)	−0.015 (0.021)
$\hat{\beta}_{2014}$: Particip. school \times ($t=2014$)	0.015 (0.019)	0.001 (0.014)	0.014 (0.014)	0.002 (0.018)	−0.020 (0.015)	−0.017 (0.017)
<i>Spillover effects: non-treated students</i>						
$\hat{\beta}_{2015}$: Particip. school \times ($t=2015$)	−0.011 (0.021)	−0.014 (0.016)	−0.009 (0.017)	0.008 (0.022)	−0.011 (0.019)	−0.018 (0.024)
Year fixed effects (omitted: 2012)	Yes	Yes	Yes	Yes	Yes	Yes
School fixed effects	Yes	Yes	Yes	Yes	Yes	Yes
Number of observations (school \times year)	601	601	601	593	593	593
<i>Test: common trends ($\hat{\beta}_{2013}=\hat{\beta}_{2014}=0$)</i>						
<i>F</i> -stat	0.33	0.01	0.67	0.01	1.22	0.51
<i>p</i> -value	0.72	0.99	0.52	0.99	0.30	0.60

Notes: This table reports the estimated spillover effects of the role model interventions for students in the non-treated classes of the schools that participated in the program evaluation in 2015, separately for male and female students in Grade 12 (science track). The outcomes we consider are those for which we document significant direct effects of the interventions, i.e., enrollment in a STEM undergraduate program, enrollment in a selective STEM program, and enrollment in a male-dominated STEM program. The results are based on a difference-in-differences specification that compares the outcomes of students in participating and non-participating schools over the period 2012 to 2015, in which the first three years correspond to the pre-intervention period. Non-participating schools are selected among high schools in the Paris region using a nearest neighbor matching procedure (with replacement) on the estimated propensity score. The baseline mean outcomes in participating and non-participating over the pre-intervention period 2012–2013 are reported in Panel A. The regression estimates are reported in Panel B. In all specifications, the dependent variable is the school-by-year average outcome of non-treated students. For non-participating schools throughout the period and for participating schools in the pre-intervention period, this mean outcome is simply the average outcome of all students enrolled in Grade 12 (science track) in the considered year. For participating schools in 2015 (the year of the intervention), this variable is computed as the weighted average outcome of students in the non-participating classes and in the participating classes that were randomly assigned to the control group, with respective weights equal to the share of participating and of non-participating classes (i.e., treatment and control) in the school. The dependent variable is regressed on school fixed effects, year fixed effects (using 2012 as the reference year) and three dummy variables that take the value one if the observation corresponds to a participating school observed in 2013, 2014, and 2015, respectively. The coefficients on the first two dummy variables capture the differential pre-trends between participating and non-participating schools whereas the coefficient on the third dummy variable measures the spillover effects of role model interventions. All regressions are weighted by school size. Standard errors (in parentheses) are clustered at the school level. The number of schools being used in the regressions for female and male students differs because one of the participating schools and one of the non-participating schools are female-only. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

K Heterogeneous Treatment Effects: Subgroup Analysis

Table K1 – Heterogeneous Treatment Effects on Student Outcomes, by Math Performance

	Girls			Boys		
	Level of performance in math			Level of performance in math		
	Below median (1)	Above median (2)	<i>p</i> -value of diff. [<i>q</i> -value] (3)	Below median (4)	Above median (5)	<i>p</i> -value of diff. [<i>q</i> -value] (6)
Panel A. Grade 10						
Grade 11: Science track	−0.018 (0.015)	0.008 (0.018)	0.272 [0.362]	−0.020 (0.020)	0.011 (0.018)	0.253 [0.326]
Positive perceptions of science-related careers (index)	0.210*** (0.043)	0.273*** (0.038)	0.281 [0.362]	0.173*** (0.042)	0.154*** (0.040)	0.750 [0.750]
More men in science-related jobs	0.169*** (0.019)	0.144*** (0.017)	0.336 [0.378]	0.188*** (0.020)	0.153*** (0.017)	0.148 [0.252]
Equal gender aptitude for math (index)	0.048 (0.037)	0.168*** (0.033)	0.017 [0.078]	0.098** (0.045)	0.185*** (0.042)	0.168 [0.252]
Women don't really like science	0.062*** (0.016)	0.053*** (0.014)	0.688 [0.689]	0.108*** (0.019)	0.096*** (0.017)	0.645 [0.726]
W face discrimination in science-related jobs	0.171*** (0.019)	0.085*** (0.017)	0.001 [0.008]	0.177*** (0.020)	0.133*** (0.019)	0.111 [0.252]
Taste for science subjects (index)	−0.072 (0.046)	−0.010 (0.041)	0.274 [0.362]	−0.081** (0.041)	0.041 (0.035)	0.016 [0.141]
Math self-concept (index)	−0.042 (0.038)	0.021 (0.039)	0.232 [0.362]	−0.016 (0.038)	0.080** (0.038)	0.058 [0.252]
Science-related career aspirations (index)	−0.045 (0.041)	0.053 (0.040)	0.079 [0.237]	−0.035 (0.042)	0.044 (0.037)	0.162 [0.252]
Panel B. Grade 12 (science track)						
Undergraduate major: STEM	0.010 (0.020)	0.031 (0.026)	0.571 [0.572]	−0.041 (0.026)	0.016 (0.029)	0.163 [0.490]
Undergraduate major: selective STEM	0.002 (0.013)	0.067*** (0.022)	0.018 [0.037]	−0.014 (0.018)	0.036 (0.027)	0.156 [0.313]
Undergraduate major: male-dominated STEM	0.024 (0.018)	0.046** (0.023)	0.513 [0.514]	−0.005 (0.025)	0.019 (0.028)	0.541 [0.541]
Positive perceptions of science-related careers (index)	0.257*** (0.054)	0.355*** (0.059)	0.277 [0.454]	0.042 (0.054)	0.257*** (0.051)	0.008 [0.076]
More men in science-related jobs	0.153*** (0.025)	0.079*** (0.024)	0.050 [0.380]	0.155*** (0.024)	0.144*** (0.019)	0.722 [0.813]
Equal gender aptitude for math (index)	0.061 (0.043)	0.135*** (0.046)	0.274 [0.454]	0.063 (0.060)	0.211*** (0.060)	0.091 [0.412]
Women don't really like science	0.028* (0.015)	0.062*** (0.016)	0.172 [0.454]	0.073*** (0.023)	0.075*** (0.021)	0.954 [0.955]
W face discrimination in science-related jobs	0.116*** (0.027)	0.088*** (0.030)	0.489 [0.551]	0.090*** (0.030)	0.050* (0.028)	0.368 [0.711]
Taste for science subjects (index)	−0.054 (0.051)	0.025 (0.056)	0.342 [0.454]	−0.034 (0.058)	0.016 (0.052)	0.553 [0.711]
Math self-concept (index)	0.061 (0.051)	−0.070 (0.053)	0.084 [0.380]	0.078* (0.046)	0.032 (0.045)	0.488 [0.711]
Science-related career aspirations (index)	0.061 (0.049)	0.137** (0.060)	0.353 [0.454]	0.008 (0.054)	0.060 (0.050)	0.514 [0.711]

Notes: This table reports estimates of the treatment effects of the role model interventions on student outcomes, separately by grade level, gender, and level of academic performance in math. Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Students' performance in math is measured from the grades obtained on the final math exam of the *Diplôme national du Brevet* at the end of middle school (for Grade 10 students) and on the final math exam of the *Baccalauréat* (for science track Grade 12 students). Columns 1 and 2 (for girls) and columns 4 and 5 (for boys) report the local average treatment effect (LATE) estimates for students below and above the median level of performance in math, respectively. They are obtained from a regression of the outcome of interest on the interaction between a classroom visit indicator and indicators for the student being below or above the median level of performance in math, using treatment assignment (interacted with the math performance dummies) as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report both the cluster-robust model-based *p*-value for the difference between the treatment effects estimates for students above vs. below the median performance in math and, in square brackets, the *p*-value (*q*-value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage *q*-values introduced in Benjamini et al. (2006) and described in Anderson (2008). *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

Table K2 – Heterogeneous Treatment Effects on Student Outcomes, by Role Model Background

	Girls			Boys		
	Role model background			Role model background		
	Resear- chers (1)	Profes- sionals (2)	<i>p</i> -value of diff. [<i>q</i> -value] (3)	Resear- chers (4)	Profes- sionals (5)	<i>p</i> -value of diff. [<i>q</i> -value] (6)
Panel A. Grade 10						
Grade 11: Science track	0.011 (0.020)	−0.016 (0.018)	0.322 [0.964]	−0.023 (0.024)	0.015 (0.019)	0.210 [0.492]
Positive perceptions of science-related careers (index)	0.227*** (0.039)	0.258*** (0.040)	0.570 [0.964]	0.136*** (0.047)	0.192*** (0.036)	0.342 [0.513]
More men in science-related jobs	0.147*** (0.019)	0.164*** (0.017)	0.512 [0.964]	0.163*** (0.020)	0.173*** (0.019)	0.728 [0.761]
Equal gender aptitude for math (index)	0.051 (0.035)	0.155*** (0.035)	0.034 [0.303]	0.071 (0.048)	0.209*** (0.038)	0.025 [0.194]
Women don't really like science	0.055*** (0.017)	0.062*** (0.014)	0.749 [0.964]	0.091*** (0.018)	0.112*** (0.018)	0.399 [0.513]
W face discrimination in science-related jobs	0.127*** (0.020)	0.127*** (0.016)	0.990 [0.990]	0.135*** (0.021)	0.168*** (0.017)	0.218 [0.492]
Taste for science subjects (index)	0.017 (0.054)	−0.081* (0.048)	0.174 [0.782]	−0.093* (0.048)	0.040 (0.042)	0.043 [0.194]
Math self-concept (index)	0.008 (0.046)	−0.020 (0.043)	0.668 [0.964]	0.029 (0.047)	0.048 (0.043)	0.760 [0.761]
Science-related career aspirations (index)	0.017 (0.045)	0.007 (0.038)	0.858 [0.966]	−0.028 (0.043)	0.035 (0.039)	0.276 [0.497]
Panel B. Grade 12 (science track)						
Undergraduate major: STEM	0.002 (0.022)	0.039** (0.017)	0.185 [0.297]	−0.007 (0.032)	0.010 (0.024)	0.663 [0.845]
Undergraduate major: selective STEM	0.008 (0.018)	0.053*** (0.014)	0.046 [0.093]	0.008 (0.025)	0.029 (0.019)	0.503 [0.504]
Undergraduate major: male-dominated STEM	0.025 (0.019)	0.046*** (0.015)	0.379 [0.379]	−0.002 (0.030)	0.031 (0.025)	0.397 [0.504]
Positive perceptions of science-related careers (index)	0.197*** (0.055)	0.386*** (0.041)	0.005 [0.024]	0.151*** (0.045)	0.158*** (0.047)	0.912 [0.913]
More men in science-related jobs	0.150*** (0.026)	0.110*** (0.021)	0.213 [0.297]	0.158*** (0.023)	0.142*** (0.020)	0.608 [0.845]
Equal gender aptitude for math (index)	0.124*** (0.047)	0.077** (0.035)	0.422 [0.475]	0.201*** (0.063)	0.078 (0.051)	0.128 [0.577]
Women don't really like science	0.045*** (0.014)	0.044*** (0.012)	0.931 [0.931]	0.088*** (0.024)	0.062*** (0.017)	0.357 [0.845]
W face discrimination in science-related jobs	0.126*** (0.034)	0.076*** (0.024)	0.222 [0.297]	0.083*** (0.028)	0.064*** (0.022)	0.581 [0.845]
Taste for science subjects (index)	−0.044 (0.054)	0.055 (0.044)	0.152 [0.297]	−0.014 (0.056)	0.010 (0.052)	0.750 [0.845]
Math self-concept (index)	0.108* (0.060)	0.013 (0.051)	0.231 [0.297]	0.173*** (0.055)	−0.006 (0.044)	0.010 [0.089]
Science-related career aspirations (index)	−0.093 (0.057)	0.246*** (0.044)	0.000 [0.001]	0.028 (0.052)	0.068 (0.042)	0.546 [0.845]

Notes: This table reports estimates of the treatment effects of the role model interventions on student outcomes, separately by grade level, gender, and by background of the female role model who visited the classroom (professional or researcher). Each row corresponds to a different linear regression performed separately by gender, with the dependent variable listed on the left. Columns 1 and 2 (for girls) and columns 4 and 5 (for boys) report the local average treatment effect (LATE) estimates for students whose class was visited by a researcher or a professional, respectively. They are obtained from a regression of the outcome of interest on the interaction between a classroom visit indicator and indicators for the role model being either a researcher or a professional, using treatment assignment (interacted with the role model background indicator) as an instrument for treatment receipt. The regression controls for school fixed effects to account for the fact that randomization was stratified by school. Standard errors (shown in parentheses) are adjusted for clustering at the unit of randomization (class). Columns 3 and 6 report both the cluster-robust model-based *p*-value for the difference between the treatment effects estimates for students visited by a professional vs. a researcher and, in square brackets, the *p*-value (*q*-value) adjusted for multiple hypotheses testing, using the False Discovery Rate (FDR) control method. Specifically, we use the sharpened two-stage *q*-values introduced in Benjamini et al. (2006) and described in Anderson (2008). *** *p* < 0.01, ** *p* < 0.05, * *p* < 0.1.

Table K3 – Treatment Effects (ITT) on Enrollment in a Selective STEM Program for Grade 12 Students: Heterogeneity by Student and Role Model Characteristics

	Dependent variable: enrolled in a selective STEM program					
	Girls			Boys		
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Treatment group indicator (T) interacted with student characteristics</i>						
T*Bac percentile rank in math (/100, demeaned)	0.136*** (0.048)	0.139*** (0.056)	0.145*** (0.056)	0.054 (0.056)	0.030 (0.060)	0.026 (0.060)
T*Bac percentile rank in French (/100, demeaned)		−0.051 (0.042)	−0.045 (0.043)		0.094* (0.056)	0.095* (0.056)
T*High SES, demeaned		0.030 (0.028)	0.029 (0.028)		−0.007 (0.030)	−0.008 (0.030)
<i>Treatment group indicator (T) interacted with role model characteristics</i>						
T*Professional	0.055*** (0.019)	0.067*** (0.020)	0.095*** (0.026)	0.047* (0.028)	0.034 (0.029)	0.061* (0.036)
T*Participated in the program the year before			−0.049** (0.024)			0.006 (0.036)
T*Age (demeaned)			0.000 (0.003)			0.000 (0.003)
T*Non-French			0.005 (0.028)			−0.011 (0.045)
T*Has children			0.014 (0.025)			0.024 (0.034)
T*Has a Ph.D. degree			0.066** (0.027)			0.057 (0.041)
T*Field: math, physics, engineering			−0.031 (0.024)			−0.021 (0.033)
<i>Other controls</i>						
Treatment group indicator (T)	Yes	Yes	Yes	Yes	Yes	Yes
Student characteristics	Yes	Yes	Yes	Yes	Yes	Yes
School fixed effects	Yes	Yes	Yes	Yes	Yes	Yes
T interacted with school characteristics	No	Yes	Yes	No	Yes	Yes
Observations	2,827	2,827	2,827	2,924	2,924	2,924
Adjusted R-squared	0.121	0.120	0.120	0.187	0.187	0.186

Notes: Each column corresponds to a separate regression. The sample is restricted to students in Grade 12 (science track). The outcome variable is an indicator for being enrolled in a selective STEM undergraduate program in the year following high school graduation, i.e., 2016/17. The models are estimated separately for girls (columns 1–3) and boys (columns 4–6). The coefficients reported in columns 1 and 4 are from a regression of the outcome variable on a treatment group indicator (*T*), student characteristics, school fixed effects, and the treatment group indicator interacted with the student's *Baccalauréat* percentile rank in math (between 0 and 1) and with an indicator for the role model being a professional. The specification in columns 2 and 5 includes further interactions between the treatment group indicator and both student and school characteristics. Finally, the specification in columns 3 and 6 adds interactions between the treatment group indicator and the characteristics of role models. The student characteristics consist of an indicator for high-SES background and percentile ranks on the *Baccalauréat* final exams in math and French. The role model characteristics consist of age and a set of indicators for being a professional, having participated in the program the year before, being non-French, having children, holding a Ph.D. degree, and having graduated from a male-dominated STEM field (math, physics, engineering). The school characteristics are dummies for the regional education authority where the high school is located (Paris, Créteil and Versailles), whether the school is public or private, the share of female students in the school, the *Baccalauréat* pass rate observed in the previous year (2014/15), and the shares of Grade 12 (science track) students from the 2014/15 cohort who enrolled in STEM programs, selective STEM programs, and male-dominated STEM programs the following year. School characteristics are only included through their interactions with the treatment group indicator, as these characteristics are absorbed by the school fixed effects. Since each high school was visited by at most one role model, role model fixed effects are also absorbed by the school fixed effects. Standard errors (in parentheses) are adjusted for clustering at the class level. Observations with some missing characteristics (11 % of the sample) are included in the regressions. An arbitrary value is assigned to all the missing characteristics and a set of dummy variables is created, with each variable being equal to one if the corresponding information is missing. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$.

L Heterogeneous Treatment Effects: Machine Learning Methods

This appendix provides additional information on the machine learning methods we use to (i) describe the heterogeneity in treatment effects and (ii) estimate the correlation between treatment effects on different outcomes. Section L.1 gives an overview of the generic approach recently developed by Chernozhukov et al. (2018) to estimate, and make inference about, key features of heterogeneous effects in randomized experiments. Section L.2 provides further details on how we implement this method in the context of our study. Section L.3 explains how we extend this method to estimate the correlation between treatment effects.

L.1 Description of the Method of Chernozhukov et al. (2018)

Motivation. Reporting treatment effects for various subgroups of participants opens the possibility of overfitting due to the large number of potential sample splits. To address this issue, one option is to specify a certain number of groups *ex ante* in a pre-analysis plan and to tie one’s hands to analyze treatment effect heterogeneity only across these groups, while correcting standard errors for multiple testing.

This approach, however, has the drawback of restricting the analysis to a small number of groups and bears the risk of missing important sources of heterogeneity. Machine Learning (ML) methods provide an attractive alternative to explore treatment effect heterogeneity in a more comprehensive manner (see Athey and Imbens, 2017, for a review). We adopt the approach developed by Chernozhukov et al. (2018) as it appears well suited for our objective. First, this approach makes it possible to conduct valid statistical inference on several objects of interest, such as average treatment effects by heterogeneity groups or the characteristics of individuals with large and small predicted treatment effects. Second, it can be implemented using any ML algorithm, allowing us to test algorithms of different degrees of sophistication, ranging from simple linear models to neural networks. Third, as described in Section L.3, this approach can be extended to estimate the correlation between treatment effects on different outcomes.

Concepts and estimation procedure. Consider an outcome variable denoted by Y . Let $Y(1)$ and $Y(0)$ denote the potential outcomes of a student when her class is and is not visited by a role model, respectively. Let Z be a vector of covariates that characterize the student and the role model who visited the class. The conditional average treatment effect (CATE), denoted by $s_0(Z)$, is defined as follows:

$$s_0(Z) \equiv \mathbb{E}[Y(1) - Y(0)|Z].$$

The approach developed by Chernozhukov et al. (2018) uses the following procedure:

1. Randomly split the data into a *training sample* and an *estimation sample* of equal size (using stratified splitting to balance the proportions of treated and control units in each subsample).
2. Use the training sample to predict the CATE using various ML algorithms. Obtain a ML predictor proxy predictor $S(Z)$.
3. Estimate and perform inference on *features* of the CATE on the estimation sample (see the definition of the features below).
4. Repeat steps 1 to 3 n times and keep track of the estimates obtained for each feature as well as their associated p -values and 95 percent confidence intervals.

5. For each feature, compute the final estimate as the median of the n available estimates. Compute the p -value for this final estimate as the median of the n available p -values multiplied by two. Compute a 90 percent confidence interval for the final estimate as the median of the n 95 percent confidence intervals.

Three features of the CATE. The CATE $s_0(Z)$ is a function for which it is difficult to obtain uniformly valid inference without making strong assumptions. It is, however, possible to obtain inference results on specific *features* of the CATE, such as the expectation of $s_0(Z)$ for heterogeneity groups induced by the ML proxy predictor $S(Z)$.

The Best Linear Predictor (BLP). The first feature of the CATE $s_0(Z)$ is its Best Linear Predictor (BLP) based on the ML proxy predictor $S(Z)$. It is formally defined as follows:

$$\text{BLP}[s_0(Z)|S(Z)] \equiv \arg \min_{f(Z) \in \text{Span}(1, S(Z))} \mathbb{E}[s_0(Z) - f(Z)]^2.$$

Chernozhukov et al. (2018) show that one can identify the BLP of $s_0(Z)$ given $S(Z)$, as well as the projection parameters $\beta_1 = \mathbb{E}[s_0(Z)]$ and $\beta_2 = \text{Cov}(s_0(Z), S(Z))/\text{Var}(S(Z))$, using the following weighted linear projection:

$$Y = \alpha_0 + \alpha B(Z) + \beta_1(T - p(Z)) + \beta_2(T - p(Z))(S(Z) - \mathbb{E}[S(Z)]) + \epsilon, \quad \mathbb{E}[w(Z)\epsilon X] = 0, \quad (\text{A.12})$$

where T is the treatment group indicator; $B(Z)$ is a ML predictor of $Y(0)$ obtained from the training sample; $p(Z)$ is the propensity score (i.e., the conditional probability of being assigned to the treatment group); $w(Z) \equiv \{p(Z)(1 - p(Z))\}^{-1}$ is the weight; and X is the vector of all regressors ($X \equiv [1, B(Z), T - p(Z), (T - p(Z))(S(Z) - \mathbb{E}[S(Z)])]$).

Equation (A.12) can be estimated using weighted least squares, after replacing $\mathbb{E}[S(Z)]$ by its empirical expectation with respect to the estimation sample.

The coefficient β_2 is informative about the correlation between the true CATE, $s_0(Z)$, and the predicted CATE, $S(Z)$. It is equal to one if the prediction is perfect and to zero if $S(Z)$ has no predictive power or if there is no treatment effect heterogeneity, that is if $s_0(Z) = s$. The main purpose of estimating β_2 is to check if the trained ML algorithms are able to detect heterogeneity.^{A.10}

Sorted Group Average Treatment Effects (GATEs). The ML predictor of the CATE, $S(Z)$, can be used to identify groups of individuals with small and large predicted treatment effects. In our setting, this is achieved by sorting students in the estimation sample (indexed by i) according to $S(Z_i)$, the predicted value of their treatment effect given their observable characteristics. We consider the bottom and top quintiles of $S(Z_i)$ and provide ITT estimates for both groups of students.

Classification Analysis (CLAN). The third feature consists in comparing the distribution of observable characteristics of students with the smallest and largest predicted treatment effects.

The three above features—the BLP, the GATEs, and the CLAN—all rely on the existence of a ML predictor $S(Z)$. The BLP provides a means to check if $S(Z)$ detects significant heterogeneity in treatment effects. If it fails to do so, the GATEs and CLAN are not particularly relevant for the analysis, as these features would provide a description of students for whom the predicted treatment effect only differs from the unobserved CATE because of a poor-quality

^{A.10}The intuition behind the formula for β_2 can be grasped by noting that Equation (A.12) is a variant of the simpler equation $Y = \alpha_0 + \alpha B(Z) + \beta'_2 T \cdot S(Z) + \epsilon$. This simpler model implies that $s_0(Z) = \beta'_2 S(Z)$, suggesting that β'_2 provides an estimate for how close the machine learning predictor $S(Z)$ is to the CATE $s_0(Z)$.

prediction.

L.2 Implementation of the Method

This section provides details on the implementation of the method of Chernozhukov et al. (2018) in our empirical setting.

Population of interest. In the main text, we focus on the sample of girls in Grade 12 (science track), since this group of students is the only one for which we find significant treatment effects on enrollment outcomes. We identify which of these female students were most affected by the program and investigate the messages to which they were particularly responsive. Results for boys in Grade 12 can be found in Table L2.

Sample splits and iterations. We perform $n = 100$ iterations of the procedure described in the previous section, which consists in (i) splitting the sample into a training and an estimation subsample of equal size; (ii) predicting the CATE on the training sample using ML methods; and (iii) estimating the three features of the CATE (BLP, GATEs, and CLAN) in the estimation sample.^{A.11} The sample splits are stratified by class, which is the randomization unit in our experimental setting: half of the girls in each Grade 12 class are randomly assigned to the training sample, while the other half are assigned to the estimation sample.

Propensity score. For each student, we estimate the probability that his or her class was randomly assigned to the treatment group. This propensity score $p(Z)$ is equal to one half in most cases, since the treatment was generally assigned to two Grade 10 classes out of four and to one Grade 12 class out of two among the classes that were selected by the school principals. In other cases, the propensity score is not exactly one half.

Machine learning methods. We consider five alternative machine learning methods to estimate the proxy predictor $S(Z)$: Elastic Net, Random Forest, Boosted Trees, Neural Network with feature extraction, and a simple linear model estimated via OLS. These methods are implemented in R using the `caret` package written by Kuhn (2008), while the general approach of Chernozhukov et al. (2018) is implemented by adapting the codes made available online by the authors.^{A.12}

For each machine learning method, the predictor $S(Z)$ is constructed in several steps. First, the model is fitted separately on the treatment and control group students in the training sample. The two fitted models are then applied to the estimation sample to obtain the predicted outcomes $\hat{Y}_i(0)$ and $\hat{Y}_i(1)$ for each individual. Finally, $S(Z)$ is obtained by taking the difference between the two predictions.^{A.13}

For each outcome, we estimate the BLP of the CATE based on the ML method whose associated predictor $S(Z)$ has the highest correlation with the CATE $s_0(Z)$ in the estimation

^{A.11}The medians of the estimated features of the CATE change little when we repeat the entire procedure using a different seed number to randomly split the data into the training and estimation samples, suggesting that 100 iterations are sufficient for the purpose of empirical convergence.

^{A.12}<https://github.com/demirermert/MLInference> (accessed on May 4, 2018).

^{A.13}Predicting outcomes for treatment and control individuals separately, before taking the difference as we do here may not be the most efficient approach to predict the CATE at finite distance. In our setting, however, alternative ML methods directly designed to detect heterogeneity in treatment effects, such as the causal forests proposed by Wager and Athey (2018), did not improve performance. We therefore decided not to rely on these ML methods for the main analysis.

sample. In practice, the best ML method for the BLP targeting of the CATE is chosen in the estimation sample by maximizing the following performance measure:

$$\Lambda \equiv |\beta_2|^2 \text{Var}(S(Z)) = \text{Corr}^2(s_0(Z), S(Z)) \text{Var}(s_0(Z)).$$

The above equation shows that maximizing Λ is equivalent to maximizing the correlation between the ML predictor $S(Z)$ and the CATE $s_0(Z)$.

The best method for the GATEs targeting of the CATE, and hence also for the CLAN, is selected based on the following performance measure:

$$\bar{\Lambda} \equiv \mathbb{E} \left(\sum_{k=1}^K \gamma_k \mathbf{1}(S \in I_k) \right)^2,$$

where K is the number of (equal-sized) heterogeneity groups, $I_k = [l_{k-1}, l_k)$ are non-overlapping intervals that divide the support of S into regions $[l_{k-1}, l_k)$ with equal or unequal masses, and γ_k is the GATE for heterogeneity group k . In practice, both performance measures lead to a similar ranking of ML methods and the methods eventually selected to produce the BLP, the GATEs/CLAN are almost always the same.

Predictors. The covariates we use to train the ML methods are three indicators for the education districts of Paris, Créteil, and Versailles, four indicators for students' socio-economic background (high SES, medium-high SES, medium-low SES, and low SES), their age, their overall percentile rank in the *Baccalauréat* exam, their percentile ranks in the French and math tests of the exam, and a vector of 56 role model fixed effects.^{A.14} Our motivation for including only a few pre-determined covariates in addition to the role model indicators is that we are mostly interested in the treatment effect heterogeneity that arises from the 56 role models (which can be seen as different treatment arms).

L.3 Correlation Between Treatment Effects on Different Outcomes

In this section, we explain how the method of Chernozhukov et al. (2018) can be extended to estimate the correlation between the treatment effects on different outcomes. We show that a set of four linear projections of the CATEs for two outcomes Y^A and Y^B on the ML predictors of the CATEs for these outcomes can be combined to estimate the correlation between the two CATEs under a natural assumption about prediction errors. This approach offers a promising alternative to other methods, such as causal mediation analysis, that are commonly used in the medical and social sciences literature to identify what factors may be part of the causal pathway between an intervention and an outcome. Indeed, our proposed method does not rely on strong identifying assumptions and can be used in any experimental setting, as long as there is a sufficiently large number of observed exogenous covariates.

A new feature: projecting a CATE on the predictor of another CATE. Let Y^A and Y^B denote two distinct outcomes and let $s_0^A(Z)$ and $s_0^B(Z)$ denote the true CATEs of a treatment T on these outcomes, given a vector of exogenous covariates Z characterizing the observational units (indexed by i). Let $\rho_{A,B|Z} \equiv \text{Corr}(s_0^A(Z), s_0^B(Z))$ denote the bivariate correlation between the CATEs on Y^A and Y^B and consider the following weighted linear

^{A.14}Each student in the control group is assigned to the role model who visited his or her high school to ensure that the role model indicators are defined for students in both the treatment and control groups. Moreover, to account for the fact that some Grade 12 students have missing *Baccalauréat* grades (less than 2 percent), we include indicators for missing grades as controls.

projection:

$$Y^A = \alpha_0 + \alpha B^B(Z) + \beta_1(T - p(Z)) + \beta_2(T - p(Z))(S^B(Z) - \mathbb{E}[S^B(Z)]) + \epsilon, \quad \mathbb{E}[w(Z)\epsilon X] = 0, \quad (\text{A.13})$$

where $B^B(Z)$ and $S^B(Z)$ are a ML predictor of outcome Y^B for individuals in the control group and a ML predictor of the CATE on Y^B , respectively. Both ML predictors are trained using a separate independent sample and are taken as given functions in Equation (A.13). The functions $p(Z)$ and $w(Z)$ and the vector X have the same meaning as in Equation (A.12). Equation (A.13) is estimated using weighted least squares, after replacing $\mathbb{E}[S^B(Z)]$ by its empirical expectation with respect to the estimation sample.

Adapting the BLP equation of Chernozhukov et al. (2018) (Equation 2.1 p. 8) by replacing the ML predictor of the CATE on outcome Y^A by the ML predictor of the CATE for outcome Y^B , we directly obtain that Equation (A.13) identifies

$$\beta_2^{A|B} = \text{Cov}(s_0^A(Z), S^B(Z)) / \text{Var}(S^B(Z)).$$

The sign of $\beta_2^{A|B}$ is informative of the extent to which the CATE on Y^A is positively or negatively correlated with the CATE on Y^B . To show this formally, we denote by η_B the approximation error in $S^B(Z)$ and we write $S^B(Z) = s_0^B(Z) + \eta_B$. Assuming that η_B is independent of $s_0^A(Z)$, we get that $\beta_2^{A|B} = \text{Cov}(s_0^A(Z), s_0^B(Z)) / \text{Var}(S^B(Z))$, which implies that $\beta_2^{A|B}$ and $\rho_{A,B|Z}$ have the same sign.

Combining BLPs to recover the correlation between treatment effects. For any pair of indices $(k, l) \in \{(A, A), (B, B), (A, B), (B, A)\}$, we can identify

$$\beta_2^{k|l} = \text{Cov}(s_0^k(Z), S^l(Z)) / \text{Var}(S^l(Z))$$

from the BLP of $s_0^k(Z)$ on $S^l(Z)$. Writing $S^A(Z) = s_0^A(Z) + \eta_A$, $S^B(Z) = s_0^B(Z) + \eta_B$, and assuming that the prediction errors η_A and η_B are independent of both the predicted functions $s_0^A(Z)$ and $s_0^B(Z)$ in the estimation sample, we can write

$$\beta_2^{k|l} = \text{Cov}(s_0^k(Z), s_0^l(Z)) / (\text{Var}(s_0^l(Z)) + \text{Var}(\eta^l(Z))).$$

Combining the formulas for the four different possible BLPs, we obtain the following expression:

$$\rho_{A,B|Z}^2 = \frac{\beta_2^{A|B} \beta_2^{B|A}}{\beta_2^{B|B} \beta_2^{A|A}},$$

which implies that the correlation $\rho_{A,B|Z}$ is identified as

$$\rho_{A,B|Z} = \text{Sign}(\beta_2^{A|B}) \frac{\sqrt{\beta_2^{A|B} \beta_2^{B|A}}}{\sqrt{\beta_2^{B|B} \beta_2^{A|A}}}. \quad (\text{A.14})$$

Practical implementation. As explained in the main text, we use the method of Chernozhukov et al. (2018) to estimate the four heterogeneity loading parameters $\beta_2^{A|A}$, $\beta_2^{B|B}$, $\beta_2^{A|B}$, and $\beta_2^{B|A}$. At each iteration of the data-splitting process, the bivariate correlation $\rho_{A,B|Z}$ is estimated by plugging the four parameter estimates into Equation (A.14). In theory, $\beta_2^{A|A}$ and $\beta_2^{B|B}$ should both be positive, while $\beta_2^{A|B}$ and $\beta_2^{B|A}$ should have the sign of $\rho_{A,B|Z}$ in each iteration of the data-splitting process. However, it can happen that the estimates $\hat{\beta}_2^{A|A}$, $\hat{\beta}_2^{B|B}$, $\hat{\beta}_2^{A|B}$, and $\hat{\beta}_2^{B|A}$ do not satisfy these conditions due to estimation error, in particular when the

predictors $S^A(Z)$ and $S^B(Z)$ are very noisy. In such cases, we do not estimate $\rho_{A,B|Z}$ and discard the corresponding iteration of the data-splitting procedure. We iterate until we reach a number of 100 iterations for which $\hat{\rho}_{A,B|Z}$ can be computed, so that our final estimates are medians computed over an identical number of iterations.^{A.15}

The estimates based on Equation (A.14) can become very large (well above one in absolute value) when the estimates of $\hat{\beta}_2^{A|A}$ or $\hat{\beta}_2^{B|B}$ are close to 0, which can occur when either or both of the predictors $S^A(Z)$ and $S^B(Z)$ are noisy. Reassuringly, we show in Table L4 that the correlation estimates $\hat{\rho}_{A,B|Z}$ are hardly affected when we exclude data splits that yield a poor ML prediction of the CATEs on outcomes Y^A or Y^B , by using only the first 100 iterations of the data-splitting process for which the estimates of the heterogeneity loading parameters $\hat{\beta}_2^{A|A}$ and $\hat{\beta}_2^{B|B}$ are above a minimum threshold t .

In the absence of a closed-form formula for the standard error of $\hat{\rho}_{A,B|Z}$, we estimate its 95 percent confidence interval as follows. At each iteration m of the data-splitting process, we compute $\hat{\rho}_{A,B|Z}^{(m)}$ (indexed by m) in the estimation sample. When $\hat{\rho}_{A,B|Z}^{(m)}$ can be computed, we estimate its 97.5 percent confidence interval using a clustered bootstrap procedure, which accounts for the clustered nature of the treatment assignment (at the class level). This procedure consists in creating B replications of the estimation sample m by drawing with replacement $N_c^{(m)}$ female students from each Grade 12 class c , where $N_c^{(m)}$ is the number of female students from class c in the estimation sample m , and computing $\rho_{A,B|Z}$ for this bootstrap sample. For each estimation sample m , this operation is repeated 6,000 times to estimate the 97.5 percent confidence interval of $\hat{\rho}_{A,B|Z}^{(m)}$ using the bootstrap percentile confidence interval method (Davison and Hinkley, 1997, chap. 5).^{A.16} The 95 percent confidence interval for $\hat{\rho}_{A,B|Z}$ is then computed as the median of the 97.5 percent confidence intervals over the first 100 iterations for which $\hat{\rho}_{A,B|Z}^{(m)}$ could be computed—the price of the splitting uncertainty being reflected in the discounting of the confidence level from $1 - \alpha$ to $1 - 2\alpha$.

^{A.15}For each pair of outcomes (Y^A, Y^B) , Table L3 indicates the proportion of random data splits for which the correlation between CATEs could be computed.

^{A.16}The 97.5 percent confidence interval of $\hat{\rho}_{A,B|Z}^{(m)}$ is estimated using only the bootstrap samples for which $\hat{\rho}_{A,B|Z}$ can be computed.

Table L1 – Heterogeneous Treatment Effect on Student Outcomes for Girls in Grade 12: Estimates Based on Machine Learning Methods

Panel A. Best Linear Predictor (BLP) of the CATE $s_0(Z)$ given the ML proxy $S(Z)$				
Parameters:	ATE (β_1)	HET (β_2)	Best ML method	
(p-values in square brackets)				
Undergraduate major: selective STEM	0.038 [0.027]	0.762 [0.031]	Elastic Net	
Undergraduate major: male-dominated STEM	0.036 [0.064]	0.088 [0.731]	Linear model	
Positive perceptions of science-related careers (index)	0.298 [0.000]	0.400 [0.555]	Elastic Net	
More men in science-related jobs	0.119 [0.000]	0.657 [0.593]	Elastic Net	
Equal gender aptitude for math (index)	0.117 [0.010]	0.324 [0.108]	Random Forest	
Women don't really like science	0.044 [0.002]	0.095 [0.566]	Linear model	
Women face discrimination in science-related jobs	0.105 [0.000]	0.496 [0.012]	Random Forest	
Taste for science subjects (index)	0.008 [1.000]	0.170 [0.137]	Linear Model	
Math self-concept (index)	0.029 [0.988]	0.257 [0.010]	Linear Model	
Science-related career aspirations (index)	0.077 [0.263]	0.245 [0.013]	Linear Model	
Panel B. Average predicted treatment effects among the most/least affected groups (GATEs)				
Heterogeneity group:	20% least affected	20% most affected	Difference most–least	Best ML method
(p-values in square brackets)				
Undergraduate major: selective STEM	−0.004 [1.000]	0.139 [0.014]	0.149 [0.026]	Elastic Net
Undergraduate major: male-dominated STEM	0.026 [1.000]	0.061 [0.464]	0.038 [1.000]	Elastic Net
Positive perceptions of science-related careers (index)	0.316 [0.037]	0.400 [0.001]	0.104 [1.000]	Elastic Net
More men in science-related jobs	0.096 [0.048]	0.160 [0.022]	0.065 [0.766]	Elastic Net
Equal gender aptitude for math (index)	0.019 [1.000]	0.246 [0.037]	0.210 [0.332]	Random Forest
Women don't really like science	0.026 [0.758]	0.073 [0.078]	0.039 [0.772]	Linear model
Women face discrimination in science-related jobs	−0.007 [1.000]	0.195 [0.003]	0.197 [0.038]	Random Forest
Taste for science subjects (index)	−0.112 [0.594]	0.138 [0.369]	0.251 [0.196]	Linear model
Math self-concept (index)	−0.122 [0.416]	0.191 [0.063]	0.317 [0.035]	Linear model
Science-related career aspirations (index)	−0.142 [0.394]	0.268 [0.047]	0.387 [0.041]	Linear model

Notes: This table reports heterogeneous treatment effects of the program on student outcomes for girls in Grade 12 (science track), using the methods developed by Chernozhukov et al. (2018). For each outcome, the conditional average treatment effect (CATE) of role model interventions, $s_0(Z)$, is predicted using five alternative ML methods: Elastic Net, Random Forest, Linear Model, Boosting, and Neural Network. The covariates Z that are used to predict the CATE consist of three indicators for the educational districts of Paris, Créteil, and Versailles, four indicators for students' socioeconomic background (high, medium-high, medium-low, and low), their age, their overall percentile rank in the *Baccalauréat* exam, their percentile ranks in the French and math tests of the exam, and a vector of 56 role model fixed effects. For each outcome, Panel A reports the parameter estimates and p -values (in square brackets) of the Best Linear Predictor (BLP) of the CATE using the best ML method. The coefficients β_1 and β_2 correspond to the average treatment effect (ATE) and heterogeneity loading (HET) parameters in the BLP, respectively. Panel B reports the Sorted Group Average Treatment Effects (GATEs), i.e., the average treatment effects among students in the top and bottom quintiles of the heterogeneous effects induced by the ML proxy predictor $S(Z)$, using the best ML method.

Table L2 – Heterogeneous Treatment Effects on Selective and Male-Dominated STEM Enrollment for Boys in Grade 12: Estimates based on Machine Learning Methods

Panel A. Best Linear Predictor (BLP) of the CATE $s_0(Z)$ given the ML proxy $S(Z)$				
Parameters:	ATE (β_1)	HET (β_2)	Best ML method	
Undergraduate Major: selective STEM	0.005	0.211	Linear Model	
p -value	[1.000]	[0.029]		
Undergraduate Major: male-dominated STEM	0.015	0.090	Linear Model	
p -value	[1.000]	[0.706]		
Panel B. Sorted Group Average Treatment Effects (GATEs): 20% most and least affected students				
Heterogeneity group:	20% least affected	20% most affected	Difference most–least	Best ML method
Undergraduate Major: selective STEM	−0.056	0.061	0.116	Linear Model
p -value	[0.358]	[0.283]	[0.086]	
Undergraduate Major: male-dominated STEM	0.051	0.010	−0.030	Boosting
p -value	[0.771]	[1.000]	[1.000]	
Panel C. Average characteristics of the 20% most and least affected students (CLAN)				
Heterogeneity group:	20% least affected	20% most affected	Difference most–least	p -value (upper bound)
Enrollment in selective STEM major				
<i>Student characteristics</i>				
Baccalauréat percentile rank in math	48.64	53.26	4.03	0.194
Baccalauréat percentile rank in French	39.95	50.94	10.45	0.000
High SES	0.495	0.494	−0.004	1.000
<i>Role model characteristics</i>				
Professional	0.395	0.600	0.214	0.000
Participated in the program the year before	0.200	0.275	0.070	0.112
Non-French	0.141	0.188	0.051	0.208
Has children	0.413	0.492	0.080	0.140
Age	32.08	33.73	1.58	0.001
Holds/prepares for a Ph.D.	0.707	0.664	−0.070	0.206
Field: Math, Physics, Engineering	0.359	0.236	−0.133	0.001
Field: Earth and Life Sciences	0.541	0.688	0.157	0.000
Enrollment in male-dominated major				
<i>Student characteristics</i>				
Baccalauréat percentile rank in math	54.72	50.21	−4.46	0.123
Baccalauréat percentile rank in French	45.41	47.25	1.38	1.000
High SES	0.465	0.527	0.068	0.248
<i>Role model characteristics</i>				
Professional	0.484	0.531	0.052	0.436
Participated in the program the year before	0.191	0.172	−0.019	1.000
Non-French	0.154	0.124	−0.025	0.820
Has children	0.489	0.489	0.004	1.000
Age	33.32	34.34	0.16	1.000
Holds/prepares for a Ph.D.	0.660	0.682	0.020	1.000
Field: Math, Physics, Engineering	0.295	0.277	−0.015	1.000
Field: Earth and Life Sciences	0.576	0.654	0.074	0.167

Notes: This table reports heterogeneous treatment effects of the program on the undergraduate enrollment outcomes of boys in Grade 12 (science track), using the methods developed by Chernozhukov et al. (2018). For each outcome, the conditional average treatment effect (CATE) of role model interventions, $s_0(Z)$, is predicted using five alternative ML methods: Elastic Net, Random Forest, Linear Model, Boosting, and Neural Network. The covariates Z that are used to predict the CATE consist of three indicators for the educational districts of Paris, Créteil, and Versailles, four indicators for students' socioeconomic background (high, medium-high, medium-low, and low), their age, their overall percentile rank in the *Baccalauréat* exam, their percentile ranks in the French and math tests of the exam, and a vector of 56 role model fixed effects. For each outcome, Panel A reports the parameter estimates and *p*-values (in square brackets) of the Best Linear Predictor (BLP) of the CATE using the best ML method. The coefficients β_1 and β_2 correspond to the average treatment effect (ATE) and heterogeneity loading (HET) parameters in the BLP, respectively. Panel B reports the Sorted Group Average Treatment Effects (GATEs), i.e., the average treatment effects among students in the top and bottom quintiles of the heterogeneous effects induced by the ML proxy predictor $S(Z)$, using the best ML method. Panel C performs a Classification Analysis (CLAN) by comparing the average characteristics of the 20 percent most and least affected students defined in terms of the ML proxy predictor. The parameter estimates and *p*-values are computed as medians over 100 splits, with nominal levels adjusted to account for the splitting uncertainty. This adjustment implies that the reported *p*-values should be interpreted as upper bounds for the actual *p*-values. Further details on the methods are provided in Appendix L.

Table L3 – Proportion of Random Data Splits for which the Correlation between Conditional Average Treatment Effects (CATEs) can be Computed, Girls in Grade 12

	Proportion of data splits such that			
	$\hat{\rho}_{A,B Z}$ can be computed*	$\hat{\beta}_2^{B B} > 0$	$\hat{\beta}_2^{A A} > 0$	$\hat{\beta}_2^{A B} \hat{\beta}_2^{B A} \geq 0$
	(1)	(2)	(3)	(4)
<i>When outcome Y^B is enrollment in a selective STEM program and outcome Y^A is:</i>				
Positive perception of science-related careers (index)	0.80	1.00	0.86	0.90
More men in science-related jobs	0.68	0.99	0.89	0.73
Equal gender aptitude for math (index)	0.35	1.00	0.98	0.36
Women don't really like science	0.34	0.99	0.84	0.40
Women face discrimination in science-related jobs	0.62	1.00	1.00	0.62
Taste for science subjects (index)	0.81	0.99	0.97	0.83
Math self-concept (index)	0.39	0.99	1.00	0.40
Science-related career aspirations (index)	0.64	0.99	1.00	0.65
Number of data splits	3,000	3,000	3,000	3,000

Notes: This table reports, for the sample of girls in Grade 12 (science track), the proportion of random data splits (out of 3,000) for which the correlation between the Conditional Average Treatment Effects (CATEs) on outcomes Y^A and Y^B could be computed. Outcome Y^B is always enrollment in selective STEM, while Y^A is the outcome listed in the corresponding row of the table. Conditional on the covariates Z , the CATEs on outcomes Y^A and Y^B are denoted by $s_0^A(Z)$ and $s_0^B(Z)$, respectively, whereas their ML proxy predictors are denoted by $S^A(Z)$ and $S^B(Z)$, respectively. For each random split, the correlation coefficient $\rho_{A,B|Z}$ is estimated as $\hat{\rho}_{A,B|Z} = \text{Sign}(\hat{\beta}_2^{A|B})(\hat{\beta}_2^{A|B} \hat{\beta}_2^{B|A})^{\frac{1}{2}} / (\hat{\beta}_2^{A|A})^{\frac{1}{2}} (\hat{\beta}_2^{B|B})^{\frac{1}{2}}$, where $\hat{\beta}_2^{k|l}$ is the estimated heterogeneity loading parameter of the Best Linear Predictor (BLP) of $s_0^k(Z)$ based on $S^l(Z)$ (with $k, l \in \{A, B\}$), using the methods in Chernozhukov et al. (2018). Column 1 indicates the fraction of data splits for which $\hat{\rho}_{A,B|Z}$ could be computed. The next three columns report the fraction of sample splits for which each of the three conditions to compute $\hat{\rho}_{A,B|Z}$ is met, i.e., $\hat{\beta}_2^{B|B} > 0$ (column 2), $\hat{\beta}_2^{A|A} > 0$ (column 3), and $\hat{\beta}_2^{A|B} \hat{\beta}_2^{B|A} \geq 0$ (column 4). The proportion of random splits such that $\hat{\beta}_2^{B|B} > 0$ varies slightly across rows because for each pair of outcomes (Y^A, Y^B) , the sample is restricted to observations with non-missing values for both outcomes (see Appendix L). Table 8 in the main text reports the median and 95 percent confidence interval of $\hat{\rho}_{A,B|Z}$ over the first 100 random data splits for which $\hat{\rho}_{A,B|Z}$ can be computed. Details are provided in Section 6.3 of the main text and in Appendix L.

Table L4 – Correlation between Conditional Average Treatment Effects (CATEs) for Girls in Grade 12: Sensitivity Analysis

	Bivariate correlation with the CATE on enrollment in a selective STEM program (from first 100 valid iterations)		
	Estimate ($\hat{\rho}_{A,B Z}$)	95% confidence interval	Proportion of valid iterations
Panel A. Data splits such that $\hat{\beta}_2^{A A} > 0.1$, $\hat{\beta}_2^{B B} > 0.1$ and $\hat{\beta}_2^{A B} \hat{\beta}_2^{B A} \geq 0$			
<i>Conditional average treatment effect (CATE) on:</i>			
Positive perception of science-related careers (index)	0.96	[0.21, 5.39]	0.73
More men in science-related jobs	−0.68	[−3.33, −0.03]	0.65
Equal gender aptitude for math (index)	0.08	[−1.90, 2.11]	0.33
Women don’t really like science	0.26	[−0.64, 3.75]	0.19
Women face discrimination in science-related jobs	−0.31	[−2.20, 0.61]	0.61
Taste for science subjects (index)	0.69	[0.07, 3.42]	0.66
Math self-concept (index)	−0.06	[−1.85, 1.37]	0.38
Science-related career aspirations (index)	0.34	[−0.61, 1.95]	0.62
Panel B. Data splits such that $\hat{\beta}_2^{A A} > 0.2$, $\hat{\beta}_2^{B B} > 0.2$ and $\hat{\beta}_2^{A B} \hat{\beta}_2^{B A} \geq 0$			
<i>Conditional average treatment effect (CATE) on:</i>			
Positive perception of science-related careers (index)	0.93	[0.21, 5.07]	0.64
More men in science-related jobs	−0.68	[−3.26, −0.03]	0.65
Equal gender aptitude for math (index)	0.05	[−1.98, 1.90]	0.31
Women don’t really like science	0.31	[−0.51, 3.44]	0.05
Women face discrimination in science-related jobs	−0.30	[−2.12, 0.64]	0.58
Taste for science subjects (index)	0.59	[0.07, 2.61]	0.34
Math self-concept (index)	0.05	[−1.68, 1.51]	0.29
Science-related career aspirations (index)	0.31	[−0.64, 1.79]	0.46

Notes: Similarly to Table 8 in the main text, this table reports, for girls in Grade 12 (science track), the estimates of the bivariate correlation $\rho_{A,B|Z}$ between the Conditional Average Treatment Effect (CATE) on enrollment in a selective STEM program, denoted by $s_0^B(Z)$, and the CATE on each of the potential channels listed in the table, denoted by $s_0^A(Z)$. The difference is that estimates provided in this table are obtained using only iterations of the data-splitting process for which the estimates of the heterogeneity loading parameters $\hat{\beta}_2^{A|A}$ and $\hat{\beta}_2^{B|B}$ are above a certain threshold. This threshold is set at 0.1 in Panel A and at 0.2 in Panel B. These restrictions are applied to check the sensitivity of the correlation estimates to excluding data splits that yield a poor ML prediction of the CATEs on outcomes Y^A or Y^B . Column 3 indicates the proportion of data splits satisfying the restrictions specified in each panel’s heading. The estimates and 95 percent confidence intervals reported in columns 1 and 2 are obtained using the first 100 data splits satisfying these restrictions. Additional details are provided in the notes of Table 8 and in Appendix L.

Appendix References

- Anderson, Michael L.**, “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*, 2008, 103 (484).
- Athey, Susan and Guido W. Imbens**, “The Econometrics of Randomized Experiments,” in Esther Duflo and Abhijit V. Banerjee, eds., *Handbook of Economic Field Experiments*, Vol. 1, Elsevier, 2017, pp. 73–140.
- Beede, David, Tiffany Julian, David Langdon, George McKittrick, Beethika Khan, and Mark Doms**, “Women in STEM: A Gender Gap to Innovation,” 2011. U.S. Department of Commerce, Economics and Statistics Administration, Issue Brief No. 04-11.
- Benjamini, Yoav, Abba M. Krieger, and Daniel Yekutieli**, “Adaptive Linear Step-up Procedures that Control the False Discovery Rate,” *Biometrika*, 2006, 93 (3), 491–507.
- Chernozhukov, Victor, Mert Demirer, Esther Duflo, and Iván Fernández-Val**, “Generic Machine Learning Inference on Heterogenous Treatment Effects in Randomized Experiments,” 2018. NBER Working Paper No. 24678.
- Davison, Anthony C. and David V. Hinkley**, *Bootstrap Methods and their Application*, Cambridge University Press, 1997.
- Duflo, Esther and Emmanuel Saez**, “The Role of Information and Social Interactions in Retirement Plan Decisions: Evidence from a Randomized Experiment,” *The Quarterly Journal of Economics*, 2003, 118 (3), 815–842.
- Fisher, Ronald A.**, *The Design of Experiments*, McMillan, 1935.
- Gayral-Taminh, Martine, Tomohiro Matsuda, Sylvie Bourdet-Loubère, Valérie Lauwers-Cances, Jean-Philippe Raynaud, and Hélène Grandjean**, “Auto-évaluation de la qualité de vie d’enfants de 6 à 12 ans : construction et premières étapes de validation du KidIQol, outil générique présenté sur ordinateur,” *Santé Publique*, 2005, 17 (2), 167–177.
- Imbens, Guido W. and Donald B. Rubin**, *Causal Inference in Statistics, Social, and Biomedical Sciences*, Cambridge University Press, 2015.
- Kuhn, Max**, “Building Predictive Models in R using the caret Package,” *Journal of Statistical Software*, 2008, 28 (5), 1–26.
- McDonald, Judith A. and Robert J. Thornton**, “Do New Male and Female College Graduates Receive Unequal Pay?,” *Journal of Human Resources*, 2007, 42 (1), 32–48.
- Paz, Lourenço S. and James E. West**, “Should We Trust Clustered Standard Errors? A Comparison with Randomization-Based Methods,” 2019. NBER Working Paper No. 25926.
- Rosenbaum, Paul R.**, *Observational Studies*, Springer, 2002.
- , *Design of Observational Studies*, Springer Series in Statistics, 2010.
- Vazquez-Bare, Gonzalo**, “Identification and Estimation of Spillover Effects in Randomized Experiments,” 2018. Manuscript.
- Wager, Stefan and Susan Athey**, “Estimation and Inference of Heterogeneous Treatment Effects using Random Forests,” *Journal of the American Statistical Association*, 2018, 113 (523), 1228–1242.