

# Peers Matter: The Heterogeneous Effects of Female Peers on Scientists' Research Focus\*

Paul W. Dai<sup>†</sup>      Hongyuan Xia<sup>‡</sup>

This Version: April 16, 2024

## Abstract

Scientific research is less likely to focus on women, despite the importance of such studies for women's welfare. This paper studies the impact of increasing female representation in doctoral studies on the research direction of PhD students in the same cohort. We analyze the dissertation and research trajectories of nearly all US healthcare and biology PhD recipients from 1985 to 2015. By exploiting quasi-random year-to-year fluctuations in the female ratio at the PhD program level, we identify the causal impact of female peers on the production of gender-related research, defined as research studying or relating to women, gender, and sex. We find that an increase in female students in a cohort encourages female PhD students to conduct more gender-related research, but discourages male PhD students from engaging in such research. Furthermore, the positive spillover effects on conducting gender-related research observed among female students are primarily attributed to collaborations and informal interactions with female peers. Conversely, the diminished or negative effects seen in male students appear to stem from competitive pressures. Taken together, our findings suggest peer effect is an important factor to resolve the scarcity-substance puzzle of gender-related research production.

Keywords: Gender-Related Research, Gender, Science of Science, Innovation, Higher Education

---

\*This research was conducted with support and resources provided by the Cornell Center for Social Sciences (CCSS) at Cornell University.

<sup>†</sup>MIT Sloan School of Management, [paulwdai@mit.edu](mailto:paulwdai@mit.edu).

<sup>‡</sup>Cornell University Department of Economics ; [hx276@cornell.edu](mailto:hx276@cornell.edu).

# 1 Introduction

Scientific research is less likely to focus on women. The shortage of papers and inventions studying or relating to women, gender, and sex, which we define as gender-related research (GRR), can have significant implications for the welfare of women. For example, the historical lack of research focus on women's health concerns has been linked to greater rates of misdiagnosis of medical conditions such as heart disease in women. However, given the importance of gender-related research, its scarcity remains a puzzle."

Previous studies investigating the shortage of gender-related research have focused mostly on the underrepresentation of female researchers (Koning, Samila, and Ferguson 2020, 2021; Nielsen et al. 2017; Truffa and Wong 2022). In our paper, we investigate this question considering the *spillover* effects of female researchers and the importance of the gender-diverse environment. In particular, we ask whether a gender-diverse academic environment, characterized by a high proportion of female peers, can inspire researchers, both male and female, to conduct more gender-related research in both the short and long term.

Female peers can affect researchers' choices of conducting gender-related research in two competing ways. On one hand, a higher proportion of female peers potentially heightens awareness and interest in female-related issues among researchers, enables formation of teams with a female majority, fosters informal interactions with female peers, and thus boosts the gender-related research production. On the other hand, the academic environment is inherently competitive, with peers competing for scarce research input and racing to publish novel findings. An increased presence of female peers might intensify competition, particularly within the realm of gender-related research. The overall effect is ex-ante ambiguous and may be different across gender.

Identifying the female peers' effects on researchers' propensity to do gender-related research is challenging due to several reasons: First, peers are not randomly assigned, and those who have more female peers may also have an unobserved preference to do gender-

related research. For instance, scholars specializing in gender studies might naturally cluster in certain universities, raising concerns about reverse causality and omitted variable bias. Second, the concept of “peers” in academic research is unclear due to the diverse nature of knowledge production, which might refer to a range of associations, from co-authors within one’s research network to colleagues in the same department. This diversity complicates the task of precisely determining the peer group impacting the researcher’s focus.

To deal with these two challenges, we focus on PhD students and take advantage of the quasi-random year-to-year variation in the cohort female ratio within the same PhD program to estimate the causal impact of female peers on the production of gender-related research. There are several relevant features of our research design. First, PhDs are among the most innovative group of people in the society (Akcigit, Pearce, and Prato 2022), and they continuously push the knowledge frontier. Second, PhDs are at early stage of their academic journey with relatively limited research networks and underdeveloped research agenda, leaving space to be affected by cohort peers. Third, our empirical strategy, the exogenous cohort female ratio within a program, exploits the inherent uncertainty faced by both admissions committees and doctoral candidates regarding each cohort’s gender composition. Although a program’s admissions committee might aim for a particular gender balance, and incoming students may be aware of historical gender distributions, the exact gender composition of each new cohort remains unpredictable for a given year. This provides a source of plausibly exogenous variation that allows us to identify a causal effect of the female peer ratio on PhDs’ propensity to do gender-related research.

We assemble a comprehensive record of the PhD dissertations and research trajectories from two database: First, ProQuest database contains comprehensive metadata of nearly all US PhD dissertations from 1985, including student names, advisors, institutions, thesis titles, abstracts, research subjects, and degree dates (Zolas et al. 2015); Second, the OpenAlex encompasses the publication information of PhD recipients (Priem, Piwowar, and Orr

2022). We link these two datasets based on names and institutions information, and focus on healthcare or biology PhD recipients who graduated between 1985 and 2015 from 62 universities affiliated with the Association of American Universities. To classify whether a particular thesis or publication studies or relates to women, gender, and sex,, we employ an off-the-shelf machine learning technique based on the information from abstract and title, as in Truffa and Wong 2022.

Our main empirical results suggest that a 10 percentage points increase in female students in a cohort increases female PhD students' propensity to conduct gender-related research in their dissertations by 0.68 percentage points. By contrast, male PhD students in the same cohort are 0.18 percentage points less likely to write a gender-related research dissertation. Moreover, the peer effects on conducting gender-related research persist for at least 5 years. Such peer effects do not appear before entering the PhD program, which further strengthens our identification assumption. We conduct a series of robustness checks, showing that our findings are not sensitive to our empirical specification. The results are also not driven by the cohort size of the PhD program, or the classification rule for PhD gender and for gender-related research, or measurement issue of our key variables. To unpack our finding, we conduct several heterogeneity tests, which reveals that the peer effects are stronger (a). in applied fields, (b). among PhD students in a large cohort, and (c). who have an advisor working on gender-related research.

We provide some empirical evidence to support our identification assumption. First, our first-order auto-regression of share of female peers on its lag, controlling program and year fixed effects, rejects the path dependency of PhD cohort composition. Second, as in the previous literature, we justify our approach through Monte Carlo simulations, demonstrating that the observed intra-cohort variability in female peer proportions aligns with what one would expect from a random distribution process. Last, we test whether the gender of the student is correlated with the cohort ratio. We show that the cohort-level leave-out mean

is not significantly different from zero after controlling for the leave-out mean of program female ratio.

Why do PhD students shift the direction of their research in response to the ratio of female peers? And why do male and female students react to the female ratio differently? In our further mechanism tests, we explore the existence of the two competing mechanisms, diffusion and competition. More specifically, we examine several potential channels: (i). collaboration with female peers, (ii). informal interactions, such as research group interactions, with female peers, and (iii). the intensity of competition to conduct gender-related research within the cohort.

We find an increased number of female peers enhance opportunities for female-to-female collaboration to do gender-related research, yet it does not necessarily lead to a similar increase in mixed-gender collaboration. This can be explained by homophily, the tendency for individuals to collaborate with others who share similar characteristics, including gender. Besides, our results suggest the importance of informal interaction. We explore one setting of informal interactions, the lab, and reveal that the effects of cohort female peers are especially stronger for female PhD students when there is another female peer in the same lab.

We also document that the negative effects of cohort female ratio for male PhD students are much stronger when there are female stars, scientists with 5-year post-graduation citations among the top 10 or top 25 percentile of their field, who are doing gender-related research, while female PhD students are nearly unaffected by the increasing competition visibility. Academic environment is inherently competitive, especially regarding the race to publish novel findings (Hill 2020). Scientists often pivot their research to minimize the negative effects of being “scooped” (Hill and Stein 2019). An increased presence of female peers might intensify competition, particularly within the realm of gender-related research. Male scientists are more likely to view increasing female peers as competitors, since they may lack an intrinsic understanding of and firsthand experience with gender-related topics, and have

a lesser comparative advantage in conducting gender-related research.

In summary, our results indicate a nuanced landscape: the positive spillover effects to do gender-related research observed among female students are primarily attributed to collaborations and informal interactions with female peers. Conversely, the diminished or negative effects seen in male students appear to stem from competitive pressures. This analysis not only confirms the existence of knowledge diffusion and competitive dynamics but also unravels the specific conditions that amplify or mitigate these effects.

**Related Literature.** First, we join the ongoing debate about how ideas are generated and how scientists choose their research direction, in particular, gender-related research. While existing literature has extensively explored how academic research directions are influenced by factors such as institutional barriers, research funding, political tensions, market competition, and the influence of leading scholars (Acemoglu, Yang, and Zhou 2021; Azoulay, Fons-Rosen, and Graff Zivin 2019; Borjas and Doran 2012; Fry 2023; Myers 2020; Sohn 2021; Truffa and Wong 2022), there remains a gap in understanding the impact of academic environments. A rare exploration in this domain is the work of Truffa and Wong 2022, which examines the shift towards gender-related research in the context of coeducation from the 1960s to the 1990s. Our research diverges and adds a new dimension to this conversation by concentrating on a different facet of the gender-diverse academic environment: the female peers within PhD cohorts. The nature of these peer effects is shaped by the specific cohort dynamics (i.e. diffusion and competition), distinct from the institution-wide shifts accompanying coeducation. Within the micro-environment of a PhD cohort, interactions are more intense and concentrated around specific research activities. This close-knit nature amplifies peer influences on individual research directions, contrasting with the broader, more diffuse impact of increased undergraduate female enrollment.

Second, our paper adds to an active literature documenting the importance of gender peer effects in academic settings. A number of studies have demonstrated that students, especially

females, benefit academically from an increase in the number of female peers in school (Anelli and Peri 2019; Bostwick and Weinberg 2022; Hoxby 2000; Lavy and Schlosser 2011; Mouganie and Wang 2020; Schneeweis and Zweimüller 2012). These studies span various educational levels and outcomes, including major selection and attrition rates in high schools, colleges, and graduate programs. Some of them use similar empirical strategy, the quasi-random cohort female ratio. We expand the literature by specifically exploring how the gender composition of peers influences the research directions chosen by PhD students, an area less examined in current literature. Furthermore, while much of the existing research emphasizes the positive spillover effects of peer presence, there is a gap in understanding the negative spillover effect, i.e., how competitive dynamics among peers shape academic choices (Chen and Hu 2022). Theoretical frameworks suggest that heightened competition can reduce the propensity to engage in collaborative or supportive behaviors (Drago and Garvey 1998). This aspect is particularly salient in PhD programs, where learning from and cooperating with peers are crucial for academic progress and research innovation (Teodoridis 2018; Wuchty, Jones, and Uzzi 2007). Our study introduces a novel dimension to the discourse on gender peer effects by highlighting how competition among PhD students, particularly between females and males, can lead to a reduced inclination towards gender-related research. Moreover, we documents that peer effects can work differently among females and males. In our study, females are affected by peers through collaboration and informal interactions, while males are affected by peers through competition pressure. This unique angle not only diversifies the current understanding of gender peer effects but also underlines the multifaceted nature of peer interactions in shaping academic trajectories.

Third, our research contributes to the literature that are trying to understand the relationship between diversity and innovation. Substantial research also has shown that racial, gender, and ethnic diversity of team members is linked to higher creativity and higher team performance (Hofstra et al. 2020; Yang et al. 2022). A series of studies suggest that the

entrance of women in previously male-dominated fields can influence research priorities and agendas (Koning, Samila, and Ferguson 2020, 2021; Nielsen et al. 2017). These studies focus mainly on the diversity of identity without causally identifying the effect of a diverse environment (Truffa and Wong 2022). Moreover, these studies typically do not examine the early, formative stages of a researcher’s career. In contrast, our research situates itself within the realm of doctoral education. We explore how gender diversity within PhD cohorts has an impact the innovative directions of emerging scholars. This perspective is crucial given recent findings that highlight the negative implications of a toxic academic environment on female representation (Dupas et al. 2021; Wu 2020) and, by extension, on innovation. By examining the PhD setting, we shed light on how diversity at this critical juncture of academic career development can shape future research trajectories.

The remainder of this paper is structured as follows. Section 2 presents the data we use in our analysis. Section 3 describes our empirical strategy, with estimation results in Section 4. Section 5 discusses the mechanism. Section 6 concludes.

## 2 Data and Measurement

In this section, we describe the data and how we construct the key variables for empirical analysis.

### 2.1 PhD data

Our primary data is ProQuest Dissertations & Theses Global Data<sup>1</sup>, which is the official offsite dissertation repository for the U.S. Library of Congress. It encompasses an extensive collection of U.S. PhD dissertations, spanning from 1985 to 2015. This database provides detailed metadata for each dissertation, including the student’s name, advisors, committee members, institution, dissertation title, abstract, research subjects, and the date of degree

---

1. <https://www.proquest.com/>.



conferral. These structural and semantic footprints enable us to scrutinize students' research focus at the very onset of their scholarly careers. Besides, the database's near-exhaustive coverage of U.S. PhD graduates enables us to comprehensively gather peer information within each academic program. Details of variable construction will be elaborated later.

For the current analysis, our dataset consists of 93,790 dissertations authored by health-care or biology PhD recipients who graduated between 1985 and 2015 from 62 universities affiliated with the Association of American Universities. We restrict our sample on health-care and biology due to the well-defined nature of gender-related research in these fields, a clarity that is less apparent in disciplines like mathematics and physics. Additionally, both healthcare and biology are at the forefront of innovation, making them ideal contexts to examine the impact of female peers on cutting-edge research.

## 2.2 Publications data

To study the long-run impact of cohort gender composition on non-dissertation publications in a longer time scope, we have integrated our primary ProQuest PhD data with the OpenAlex database (Priem, Piwowar, and Orr 2022). OpenAlex is a comprehensive repository containing over 240 million published papers across various journals and conferences. Utilizing key information such as PhD names, advisor names, affiliations, degree years, and research fields from the ProQuest database, we linked these details with publication records in the OpenAlex database.

Our data matching procedure between ProQuest and OpenAlex is multifaceted. Initially, we filtered OpenAlex authors based on their affiliation with universities within the Association of American Universities. Subsequently, for each PhD candidate in our ProQuest dataset, we matched surnames and computed a similarity score for first names within each corresponding affiliation. This methodology offers two primary advantages: firstly, it accommodates various forms of name representation, such as nicknames and abbreviations, thus

enhancing flexibility beyond direct name matching; secondly, it prioritizes accuracy by applying a more stringent criterion to surnames compared to a full-name similarity assessment. In the final matching stage, we excluded matches where a PhD candidate had more than five publications before their PhD or their first publication appeared more than five years post-graduation, or if the research fields were significantly dissimilar.

We implement the similar procedure for the advisors of all PhD students, and supplement our matched PhD sample using coauthors of each matched PhD advisors. If the PhD’s advisor has a coauthor that is the same name as the PhD student, we take it as the same person. We finally get 24,330 PhDs (about 25.9% of PhD recipients in our ProQuest sample) matched to OpenAlex data.

## 2.3 Main Variables

In this section, we describe our classification of the gender of PhD students, assigning the field of research, how we group PhD into different cohorts, and construction of gender-related research.

**Gender of PhD Recipients.** Since the gender of PhD recipients is not provided in the ProQuest database, we classify the gender of the researcher by using PhD recipients’ first names and Genderize.io<sup>2</sup>, an API that has been employed by academia to identify gender and report its true positive probability of this classification (Huang et al. 2020; Koffi and Marx 2023; Topaz and Sen 2016). For example, the reported true positive probability is 0.99 and 0.86 for “Paul” and “Hongyuan” as the first name to be identified as male. In the process, All 93,790 PhD recipients can be assigned to either female or male, and 81,302 PhD recipients (about 86.7%) can be assigned to one gender with a true positive rate no less than 0.9.

**Cohort Female Ratio.** Our principal independent variable is the cohort female ratio,

---

2. <https://genderize.io/>

which we define as the proportion of female students graduating in the same year within a given PhD program. Although the ProQuest database lacks explicit program information, it does include specific class or subject terms for each dissertation. These terms are categorized into 432 subject areas, which further align with 21 broader disciplines. For example, biology (subject number: 0306) and molecular biology (subject number: 0307) are subjects within the biological science disciplines. However, these subject terms, despite their specificity to research topics, may not accurately represent the actual academic programs of the PhD recipients.

To deal with this issue, we undertake a reclassification of these subject terms according to the first four digits of the CIP (Classification of Instructional Programs) code. The CIP framework provides a systematic approach for the effective tracking and reporting of various study fields and program completions<sup>3</sup>. Our reclassification process involves an initial Google search, where we pair each subject term with “CIP code” as keywords. This search aims to preliminarily align the subject terms with their corresponding CIP categories. Following the initial search results, we perform a manual verification to ensure the accuracy and relevance of the Google-derived classifications. This reclassification results in the consolidation of the subject terms into 182 distinct fields. With this refined classification, we define a “program” as a specific combination of a university and a field of study, and a “cohort” as the group of PhD graduates who completed their degrees in the same field and year at a particular university. For example, individuals who received their PhDs from Cornell University in the field of biological engineering in 2015 are classified as belonging to the same cohort.

**Gender-Related Research (GRR).** Our main outcome variables are whether the dissertation is gender-related research and the ratio of gender-related words of the dissertation. We use a similar keyword approach adopted by Truffa and Wong 2022. By using Datamuse API <sup>4</sup>, a word-finding query search engine based on Google Books Ngrams data and other

---

3. <https://nces.ed.gov/ipeds/cipcode/browse.aspx?y=55>

4. <https://www.datamuse.com/api/>

corpus-based datasets. We compile our baseline word list by selecting the top 20 most related words, such as “gender”, “female”, “women”, and “sex”, with the full list in Appendix A. We intentionally exclude male-related words because historically men are considered “standard” in research. The gender-related research is labeled as one if at least one of these keywords appears in either the title or the abstract. There are two advantages of this approach: (1) titles and abstracts are available for nearly all PhD recipients, which entitles a consistent measure of the gender-related research; (2) this approach can be applied broadly to all fields instead of only one or two fields in some existing work (Koning, Samila, and Ferguson 2020). Under such definition, a research paper is considered gender-related if the research topic is about women or it highlights analysis pertaining to gender or women. There is 6.4 percent of dissertations can be defined as gender-related. We provide two examples of gender-related dissertations with the corresponding titles and abstracts in Appendix. Besides, to measure the intensity of gender relatedness for each dissertation, we calculate the ratio of gender-related words in title and abstract, which is a continuous variable ranging from 0 to 1. However, we acknowledge potential concerns regarding to this approach. In particular, the non-agnostic nature of the keyword list could introduce bias. To mitigate this issue, we conduct several robustness checks by varying keyword list as in Appendix A. Our findings remain consistent across these different specifications, reinforcing the validity of our approach.

Similar approach applies to measure whether a non-dissertation publication in OpenAlex is gender-related research or not. For the PhD students and advisors who are matched to OpenAlex data, they publish 25,280,890 papers in total. We only use the title recorded in OpenAlex to define whether the paper is gender-related research due to the large number of papers. There is about 2 percent of all papers are classified as gender-related, this ratio is comparable to using a keyword approach only to titles in dissertations (i.e. 0.015). Based on that, we create three categories of variables: the first category is an indicator variable indicating whether the author does gender-related research, the second category is

the share of gender-related research in all papers by the author, and the third category is the number of gender-related research papers. To account for potential delay in publication process and the lumpiness of number of published paper in a year, we construct four variables in wider time windows, that is, 6-10 years before graduation, 1 to 5 years before graduation, 0 to 5 years after graduation, 6 to 10 years after graduation.

Figure 2 displays the number of gender-related dissertation produced by female and male PhD students and the average cohort female ratio over time. Table 1 displays the summary statistics of all variables in our main specification.

### 3 Empirical Strategy

In this section, we introduce the empirical specification in Section 3.1, discuss and provide evidence for our identification assumption in Section 3.2.

#### 3.1 Empirical Specification

Our empirical strategy is essentially a difference-in-differences approach, comparing female and male PhD recipients between cohorts with a high fraction of female students and those with relatively low one within a given doctoral program, following

$$\begin{aligned}
 Y_{i,u,s,t} = f & \left( \beta_0 + \beta_f \text{Female}_i \times \text{Cohort Female Ratio}_{u,s,t} \right. \\
 & + \beta_m \text{Male}_i \times \text{Cohort Female Ratio}_{u,s,t} \\
 & \left. + \gamma \text{Female}_i + \text{Controls}_{i,u,s,t} + \lambda_{u,s} + \psi_t + \epsilon_{i,u,s,t} \right), \quad (1)
 \end{aligned}$$

where  $Y_{i,u,s,t}$  is gender-related research related outcome, for example, an indicator equals one when the PhD student  $i$ , who graduated in year  $t$  with field  $s$  from university  $u$ , has a dissertation that is gender-related. We include the program fixed effects  $\lambda_{u,s}$ , where a

program is defined as a university-field combination. The inclusion of  $\lambda_{u,s}$  ensures that all the comparisons are within the same program. In other words, we do not rely on comparisons across different programs and only compare PhD students who are in the same program (university  $\times$  field) but different cohorts. The year fixed effects  $\psi_t$  control for the overall time trend.

The primary variables of interest are the interaction terms of the PhD student  $i$ 's gender indicator (i.e.  $\text{Female}_i$  and  $\text{Male}_i$ ) and Cohort Female Ratio $_{u,s,t}$ , which measures the percentage of student  $i$ 's peers graduated from the same program (university  $u \times$  field  $s$ ) in year  $t$  who are female. We test the robustness of our main specification using several alternative measures of cohort female ratio, including the number of female peers in the cohort, the ratio of women to men in the cohort, and the female ratio in the graduated and neighbor cohorts.

The coefficient  $\beta_f$  measures difference of probabilities of conducting gender-related research for females in a higher versus lower cohort female ratio settings within the same PhD program.  $\beta_f > 0$  implies a higher cohort female ratio results in more gender-related research production among female PhD students, and  $\beta_f < 0$  suggests the opposite. Similar interpretation applies to  $\beta_m$ , which measures the effects of cohort female ratio on male PhD students. The signs of these estimates enable us to distinguish different mechanisms of peer effects: on one hand, increasing female representation in a cohort draw more attention and interest in female-related research topics, encourage the formation of female-led research team, which encourages more gender-related research; on the other hand, academic environment is rife with competition, for example, lab resource, research funds, and the attention from advisors, etc. In this way, a higher female share in the cohort may intensify the competition of conducting gender-related research, which may even *crowd-out* the production of that.

In terms of our control variables, we add the research focus of PhD students' advisors (i.e. gender-related or not) as control variables to further alleviate the omitted variable bias stemming from faculty characteristics. An advisor's research interests can significantly shape

a student’s academic trajectory, potentially serving as a primary source of omitted variable bias. This influence is particularly evident when students select their PhD programs, often considering the research areas of potential advisors as a decisive factor. For instance, in many natural science disciplines, the admissions process is fundamentally a mutual selection between advisors and students. These students, in choosing this lab, are not only selecting a program but also aligning their future research with their advisor’s expertise. Such scenarios underscore the importance of controlling for advisors’ research focus, allowing us to more accurately isolate the effects of cohort gender composition on students’ research directions towards gender-related topics. By doing so, we significantly reduce the potential bias introduced during the admission and advisor selection process. Other control variables include the size of the cohort and the number of total words in the title and abstract of each dissertation.

We choose linear probability model as our baseline specification because it enables us to employ program level and year fixed effects to control for the time-invariant characteristics of each program and the overall time trend. The interpretation of the implied marginal effects is also easier in this model. As an additional robustness check, we employ probit models when the outcome variable is a binary variable and conditional fixed effects Poisson model with QML (quasi-maximum-likelihood) when the outcome variables are non-negative discrete variables (i.e., number of gender-related research).

### **3.2 Identification Assumption**

Our identification strategy relies on the assumption that within a particular doctoral program, year-to-year variation in cohort gender composition is quasi-random and not correlated with other unobserved factors influencing the research focus of the PhD students within that cohort. More specifically, it exploits the inherent uncertainty faced by both admissions committees and prospective doctoral candidates regarding each cohort’s gender composition. Although the admission committee of a doctoral program might aim for a particular gender

balance, and incoming students may be aware of historical gender distributions, the exact gender composition of each new cohort remains unpredictable.

In this section, we further provide several pieces of evidence supporting our identification strategy.

**No Within-Program Path Dependency on Cohort Female Ratio.** Our first evidence exploits the time series nature of cohort female variation in the program. One potential violation of this method could be some omitted variables that affect the gender ratio of a program and the research focus of the PhD students simultaneously. For example, a new female faculty member may attract more female students to the program while encouraging students to do gender-related research at the same time.

A telling signal of this type of endogeneity would be any evidence of time trends in the cohort gender composition within programs. We visualize this assumption in Figure 3, where each line represents a trajectory of female ratio in a healthcare or biology PhD program of the Cornell University over time. As we can notice there are no clear upward or downward trends in gender composition in programs in healthcare and biology, especially in a relatively shorter time window.

To further rule out the concern that the year-to-year fluctuation of gender-related research conducted by faculty may drive the cohort female ratio fluctuations, we intentionally pick programs in engineering as a reference in Figure 3, where gender-related research, by its nature, is rare. We observe very similar fluctuations for programs in engineering, despite that the unconditional mean of gender-related research is significantly lower than healthcare and biology.

We conduct a formal statistical test to verify that there is no path dependency of female ratio within the PhD program. First, we show that there is no statistically meaningful correlation of female ratio in adjacent years. Specifically, we estimate the following first-order



auto-regression model of cohort female ratio,

$$\text{Cohort Female Ratio}_{u,s,t} = \alpha_0 + \alpha_1 \text{Cohort Female Ratio}_{u,s,t-h} + \lambda_{u,s} + \psi_t + \epsilon_{u,s,t}, \quad (2)$$

where  $h = 1$ , controlling university and year fixed effects  $\kappa_u$  and  $\psi_t$ , respectively, with standard error clustered at university level. Appendix Table B.1 reports that  $\hat{\alpha}_1$  is not statistically distinguishable from 0 with a negligible magnitude, suggesting there is no path dependency in share of female students within a program in adjacent cohorts. Second, it is possible that PhD committee may pursue a balanced gender composition for all of the students enrolled in the program at a given year. If this argument were true, we should expect that the cohort gender ratio for new admitted cohort is very similar with that for graduated cohort, which may challenge our identification strategy. To rule out this concern, we revisit (2) by setting  $h = 5$ , a time length for a typical US PhD program, and test the null hypothesis that  $\alpha_1 = 1$ . Results in Appendix Table B.1 reject the null with p-value less than 1%. Lastly, we re-estimate these two regressions for PhD program in engineering. The result is very similar with those in healthcare and biology, ruling out the story that faculty gender-related research is the main cause of cohort female ratio fluctuations.

**Monte Carlo Simulation.** In our second test, we leverage the statistical distribution of cohort female ratio, by using the Monte Carlo simulation exercise as in the previous literature (Bostwick and Weinberg 2022; Lavy and Schlosser 2011; Mouganie and Wang 2020). The goal is to show that the observed within-program variation in gender composition in our data closely resembles the randomly generated variation from a binomial distribution. Specifically, for each doctoral program, we randomly generate the gender of the students in each cohort using a binomial distribution function  $\text{Binomial}(n, p)$ . The parameter  $n$  equals the actual cohort size and  $p$  equals the average proportion of females in that program across all years. Then, we compute the within-program simulated standard deviation of the proportion of females over all years. We repeat this process over 1,000 iterations to obtain an

empirical confidence interval for the standard deviation for each program. Our observed within-program standard deviation lies within the empirical 90% confidence interval for 91% of PhD programs in our sample, which further supports our assumption that the within-program, year-to-year variation in cohort gender composition is as good as random.

**Randomization Test Using Linear-in-Mean Model.** In our third test, we directly estimate a linear-in-mean model to show that cohort female ratio is as good as random within a program across years. Specifically, we estimate

$$\text{Female}_{i,u,s,t} = \theta_0 + \theta_1 \overline{\text{Female}}_{-i,u,s,t} + \theta_2 \overline{\text{Female}}_{-i,u,s} + \lambda_{u,s} + \psi_t + \epsilon_{i,u,s,t}, \quad (3)$$

where the dependent variable is an indicator for female student for individual  $i$ , the  $\overline{\text{Female}}_{-i,u,s,t}$  represents the female ratio of all PhD in the same program and year other than individual  $i$  herself <sup>5</sup>. We follow the recommendation by Guryan, Kroft, and Notowidigdo 2009 by adding the leave-me-out mean for the entire program pooling the PhD students across all years, i.e.  $\overline{\text{Female}}_{-i,u,s}$ , to correct the bias for estimating  $\theta_1$ , especially when the cohort size is relatively small. <sup>6</sup> We expect to see  $\theta_1$  is indistinguishable from 0, if our identification assumption is valid. Table B.11 reports the estimated  $\theta_2$ , which is insignificant and does not depend on the inclusion of covariates.

## 4 Empirical Results

We present the main estimation result in Section 4.1 and discuss the results of robustness checks of that in Section 4.2. We further investigate the driving factor for these female peer effects

---

5. The previous literature also call it leave-me-out mean satisfying the following accounting identity:  $\text{Female}_{i,u,s,t} = N\overline{\text{Female}}_{u,s,t} - (N - 1)\overline{\text{Female}}_{-i,u,s,t}$ , where  $N$  is the cohort size and  $\overline{\text{Female}}_{u,s,t}$  is exactly the cohort female ratio.

6. Guryan, Kroft, and Notowidigdo 2009 shows that, without the this correction, even if the grouping or the selection into the same cohort is random,  $\theta_1$  is negatively biased and the bias is decreasing in the cohort size.

## 4.1 Main Result

Table 2 reports the impact of cohort female ratio on PhDs' research focus in their dissertations following the specification (1). We use two outcome variables: (i) a binary variable indicating whether the PhD's dissertation is gender-related in Columns (1) to (3), and (ii) the ratio of gender-related keywords in title and abstract to measure the gender relatedness of each dissertation in Columns (4) to (6).

To start with, we observe a positive correlation between cohort female ratio and the propensity of doing gender-related research. More specifically, Columns (1) and (4) suggest that, when the share of female peers increases 10 percentage points, PhD students in the same cohort are 0.17 percentage points (27% of the unconditional mean) more likely to do gender-related research, and the ratio of gender-related words in titles and abstracts increase by 0.0011 percentage point (17% of the unconditional mean). To unpack this effect, in other columns, we include the interaction terms of cohort female ratio with gender indicators. We uncover the differential effects of increasing cohort female ratio on producing gender-related research for female and male PhDs.

Column (2) suggests a 0.61 percentage point increase in probability of having a gender-related dissertation for female PhDs if there is 10% more females in the cohort. By contrast, such effect for male PhD, if anything, is negative, with the magnitude of 0.20 percentage decrease. Result in Column (4) agrees with our finding using gender-related research intensity as the outcome. In Columns (3) and (6), we include the advisor's research focus, cohort size, and the number of words in title and abstract as controls variables. Reassuringly, the estimation results are very similar with the specification without these controls, which indirectly suggest that the main source of variation we exploit stems from the within-program cohort female ratio instead of other plausible driving forces.

## 4.2 Robustness Check

We conduct a series of robustness checks. By and large, these results confirm our main findings.

**Classification of Gender.** As we mentioned in Section 2, we don't have the gender information for each PhD recipient. We use the Genderize.io to classify the gender using the first name of each PhD recipient. Besides reporting the identified gender for each first name, Genderize.io also reports the probability of that name being the assigned gender. To alleviate any concerns attributed to the way we assign gender, we set different probability thresholds and only use the observation above these thresholds as our sample in Appendix Table B.6, which reports similar results as Table 2.

**Classification of Gender-Related Research.** In the baseline result, we classify gender-related research based on any occurrence of gender-related research related keywords in either thesis title and abstract. We re-construct the dependent variable only based on title or abstract separately, and report the estimation results in Appendix Table B.7. Besides, in Appendix Table B.8, by using different keyword lists, we argue that our results are not driven by deliberate choices of keywords to define the gender-related research. All the tests deliver similar results as those in Table 2.

**Alternative Definitions of Female Peers.** In our main specification, we define cohort female ratio as the share of student  $i$ 's peers graduated from the same program in year  $t$  who are female. We construct several alternative measures of the cohort gender composition, including the number of female peers in the cohort, the ratio of female peers in the graduate and neighbor cohorts, and the ratio of females to males in the cohort. Table B.9 reports the results using different measures of cohort gender composition. Interestingly, these results show no evidence of a linear effect of the number of female peers in a program on the probability of doing gender-related research for either female or male PhDs. However, this finding is not inconsistent with the main results and merely indicates that the effect of an

additional female peer interacts with the cohort size (e.g., one additional female peer has a large effect in a small cohort and little to no effect in a very large cohort). This interaction is better captured by the use of the percent female measure in the main specification. The effects of the ratio of female peers in the graduate and neighbor cohorts and the ratio of females to males in the cohort are similar as our main specification.

**Empirical Specification and Model.** In Appendix Table [B.10](#), we report our results using alternative empirical specification. In Column (1), we further include field  $\times$  year and university  $\times$  year fixed effects to remove any potential university- or field-specific time trends. In Column (2), we replace the baseline fixed effects with program-specific linear time trend, which allows that different PhD programs can have different general trajectory of female ratio. The estimates of these two model specifications share the same economic and statistical significance as our baseline estimation. Moreover, we also estimate a probit model for the short-term effects when using binary variable as dependent variable. We acknowledge that adding strong fixed effects can result in biased estimates in the probit model. Reassuringly, the results in Column (3) is consistent with our main findings. To mitigate the effect of adding too many binary variables in the probit model, we resort to replace the fixed effects with program-specific linear trend in Column (4). This modification yields similar results as Column (3). By and large, our main findings are not driven by the empirical specification and model we utilize.

### 4.3 Heterogeneity

What determines the strength of these female peer effects towards female and male PhD within the same cohort? To investigate this question, we further divide our sample into different groups based on several characteristics of advisors, fields and cohorts.

First, we divide the sample by whether the advisor of PhD student conducts gender-related research or not. Our results in Table [B.2](#) indicate that the female peer effects on

PhDs' propensity to do gender-related research and gender relatedness of dissertations are stronger for those with advisors who are also doing gender-related research. This may suggest a certain degree of sorting effects, those PhD students with a gender-related research interest are more likely to choose advisors who have investigated these topics. When a PhD comes to decide which program to choose, they should take the research of their potential advisors into consideration. In fact, in many natural science fields, admission is mainly a process of advisors and students choosing each other. As a result, controlling the research focus of PhD students' advisors can help us to circumvent much omitted variable bias caused by the selection process during admission. This test justifies the inclusion of the advisor's research focus as a control variable.

Moreover, we further separate the doctoral studies into basic research and applied research based on the CIP code.<sup>7</sup> Applied research is more downstream and close to commercialization. These demand side factors can facilitate additional *pulling* effects on the production of gender-related research. Indeed, Table B.3 uncovers a stronger cohort female peer effects for applied healthcare programs.

Lastly, we group our sample by the cohort size of the program. The female peer effects might be undermined in a larger cohort. For instance, PhD students are ranked in the cohort for scholarship, prizes and grants, which often have certain quotas; or the attention and support from department always focus on the top students. In this way, a larger cohort size implies more peer pressure, and thus discourage the dissipation of gender-related research. By mitigating this concern, Table B.4 compares the estimation of our main specification for larger and smaller cohorts. The estimates are similar across different cohort size, and, the effects of cohort female ratio are slightly stronger for larger cohorts.

---

7. We classify programs as a basic research program if the program is a biology program with CIP code starting from 26 and an applied research program if the program is a healthcare program with CIP code starting from 51.

## 4.4 Long-Run Effects

Do the female peer effects during PhD have a long run impact on the direction of scientific research even after PhD graduation for these early-career scientists?

Table 3 examines the effect of cohort female ratio on PhDs' research focus in their non-dissertation publications, leveraging our merged ProQuest and OpenAlex data<sup>8</sup>. Specifically, we re-estimate the empirical specification (1) changing the dependent variable as (i) whether the PhD student conducts gender-related research in different time horizons, (ii) the ratio of gender-related papers in different time horizons, and (iii) the number of gender-related papers in different horizons. For the number of gender-related papers, since the variable is non-negative discrete, we implement the estimation using conditional fixed effects Poisson model with QML (quasi-maximum-likelihood).<sup>9</sup> We compute the outcome variables within a five-year time horizon, allowing for a potential delay in the paper publication process.

The estimation results indicate that the cohort female peer effects is not a one-time impact on the topic selection of PhD dissertation, instead, it has more far-reaching influence after PhD graduation. Column (1) suggests that, before PhD, there are no discernible peer effects on the production of gender-related research. This result can be viewed as an additional verification test of our identification assumption because it is impossible that potential peer effects take place cannot happen before the PhD cohorts meet each other. In Column (2), we observe a 10 percentage points increase in share of female peers leads to 0.21 percentage points increase in the ratio of gender-related papers for female PhDs who have publications, yet doesn't significantly affect male PhD students or increase the probability of conducting gender-related research. Column (3) examines the effects of cohort female ratio on the publication of gender-related research within five years after graduation. We find

---

8. As a robustness check to show our result is not sensitive to the selection bias induced our matching algorithm, we rerun this specification on the ProQuest-OpenAlex matched data. The results are similar with our findings in Table 2, yet the effects on male are not significant.

9. The observations in Panels B and C are different from those in Panel A because some PhD students don't publish any papers in that time window and Poisson regression will drop observations in programs without variations in the dependent variable.

female PhD's probability to do gender-related research and the ratio of gender-related papers within five years after graduation increase 0.42 and 0.16 percentage points, respectively, when there is a 10 percentage points increase in cohort female ratio. However, there is no significant impact for male PhD students. In Column (4), we focus on a relatively long-term effects (i.e. more than five years after graduation). Even though we don't find any evidence on the effects of cohort female ratio on PhD's probability to do gender-related research, we do find that high cohort female ratio causes higher ratio of gender-related papers for female PhD students and low ratio of gender-related papers for male PhDs among those who still publish papers after five years of graduation. For all periods, we don't find any evidence on the number of gender-related papers. While we admit the ratio of gender-related papers and number of gender-related papers are similar measurements, they are different as number of gender-related papers doesn't been standardized by the total number of papers and might reflect gender difference in the number of papers. Collectively, the cohort female peer effects have a sizable impact on the research direction for at least 5 years after graduation.

## 5 Mechanism

Why women are more likely to do gender-related research when they have more female peers, while men are less or negatively affected by more female peers? On one hand, the positive effects may stem from enhanced learning opportunities and knowledge diffusion. On the other hand, female peer effects can have negative effects, which arise from increased competition.

Motivated by these two competing forces, in this section, we outline several mechanisms: (i) formal interactions, such as collaboration between PhDs and their female peers; (ii) informal interactions within the PhD cohorts; (iii) changing competition visibility among peers in doing gender-related research. Our goal for this analysis is two-fold: (i) confirm the existence of these spillover effects, and (ii) elucidate the specific conditions under which one effect may predominate over the other.



## 5.1 Collaboration with Female Peers

Increasing female peers within the cohort potentially expands the pool of prospective female collaborators for PhD students. This may possibly foster the formation of female-dominant research teams due to homophily, the concept of “birds of a feather flocking together”, defined as the tendency for individuals to gravitate towards and collaborate with others who share similar characteristics, including gender.<sup>10</sup> Meanwhile, an increasing representation of female could crowd out the formation of mixed-gender team, which may explain the observed differences in spillover effects between male and female PhD students within the same cohort. We report the empirical results of testing this collaboration channel in Table 4.

First, we explore whether the cohort female ratio leads to the PhD to have more gender-related research in female-dominant teams. We identify female-dominant teams relying on the linked OpenAlex and ProQuest data, which encompasses information about coauthor and publication information. We recognize each published paper requires a research team to produce and classify the team as female-dominant if half of the authors are females. Column (1) suggests that for every 10 percentage points increase in the female ratio within a cohort, there is a corresponding 0.4 percentage points rise in the probability to have at least one gender-related paper published in female-dominant teams for female students. Moreover, female PhD students in a cohort with high female ratio produce more gender-related papers in female-dominant teams (coefficient = 0.014, p-value < 0.01), while male PhD students in a cohort with high female ratio produce less gender-related papers in female-dominant teams (coefficient = -0.006, p-value < 0.01).

Since female scientists may benefit from producing gender-related research, are they trying to exploit this benefit endogenously by coauthoring more with female researchers? Results in Column (2) answer this question: females are more inclined to write gender-related

---

10. Gender homophily has also been noted in investment and other sectors (Greenberg and Mollick 2017; Stolper and Walter 2019; Zeltzer 2020; Zhou, Chai, and Freeman 2024), suggesting that females might prefer collaborating with other females, while males may be less inclined to partner with female peers.

papers with other female coauthors, when the cohort female ratio is high, while the opposite happens for male students. While the peer effects of an increased female ratio in PhD cohorts are likely multifaceted, our results provide some evidence that formal interactions, namely collaboration, play a crucial role in encouraging females to engage in gender-related research. Additionally, these findings reveal that gender-based preferences, particularly in the context of team composition, may act as a deterrent for male students engaging in gender-related research, especially within mixed-gender or female-dominant teams. Taken together, we point out that preferences and biases in collaborative settings could be influencing the research direction, steering male researchers away from female-centered topics.

## 5.2 Informal Interactions with Female Peers

Beyond formal collaboration, female peers may have an impact on a student’s research direction via informal interactions. In the context of a PhD program, students are immersed in class discussions, out-of-class activities, and casual conversations. The literature on knowledge transfer and idea generation increasingly recognizes the vital role of these informal interactions. Studies have shown that informal interactions, often facilitated by physical proximity, are key drivers in the dissemination of ideas and knowledge within localized academic or professional communities (Atkin, Chen, and Popov 2022; Chai and Freeman 2019; Hasan and Koning 2019; Roche, Oettl, and Catalini 2023).

While it is challenging to catalog and directly test the impact of all these diverse forms of informal interaction, our study aims to provide suggestive evidence of this mechanism. One of the special features for biology and healthcare doctoral programs is that students are divided into labs to conduct experiments. For those working in the same lab, interactions occur daily or weekly, which are more likely to facilitate knowledge diffusion and thus provide support in pursuing gender-related ideas. Motivated by this feature, we group students sharing the same advisor as in the same lab. Specifically, we estimate the empirical specification 1 by

adding further interacting the interaction terms with a binary variable indicating whether there are other female students in the same lab.

Table 5 reveals that female peer effects are especially stronger for female PhD students when other female peers are present in the same lab (coefficient = 0.028, p-value < 0.10). Besides, for male PhD students, the coefficient of this triple interaction is positive in terms of producing gender-related research but not significant. By examining the moderating effects of the presence of female peers in the same lab, we provide evidence on the role of informal interactions within the lab in shaping research trajectories.

### 5.3 Competition Visibility

The academic environment is inherently competitive, particularly in the race to publish novel findings (Hill 2020; Hill and Stein 2019). Influenced by a relatively larger number of female peers in the cohort, PhD students, for instance, begin to focus on female reproductive health, thereby intensifying the perceived competition for innovative findings and prestigious publication slots. Male researchers, who might be less familiar with or adept at conducting research in gender-related topics, could perceive this heightened competition as a significant barrier. As we mentioned above, we assume that female has comparative advantage in doing gender-related research. When a promising female academic steps into the domain of gender-related research, male PhDs may find it challenging to compete directly with their female counterparts.

To test this channel, we use the presence of *star* researchers as a proxy for the severity of competition. While we don't have the academic records data for each PhD student, we use the linked data to measure whether a PhD student is strong in terms of the ability to do academic research. For each linked PhD student, we calculate their average number of forward citations for papers published within five years after graduation. We assume those PhDs who can be linked to OpenAlex and have average citations among top 10 percent or top 25 percent

in their field are stars. Then, we construct a binary variable indicating whether there are female stars conducting gender-related research for each cohort. We interact this binary variable with the existing interaction term and reestimate our specification 1. As shown in Table 6, the negative effects of cohort female ratio for male PhD students are much stronger when there are female stars who are doing gender-related research (coefficient =  $-0.107$ ,  $p\text{-value} < 0.05$ ). In contrast, female PhD students are nearly unaffected by this increasing competition visibility. To sum up, male PhD students shy away from conducting gender-related research in the presence of fiercer competition.

## 6 Conclusion

In this paper, we estimate the causal effects of cohort female ratio on PhD students' choices of research topics, i.e., gender-related research. This is particularly relevant giving the importance but scarcity of research focused on the needs of females.

We find having more cohort female peers increases female PhD students' propensity to do gender-related research but discourage male PhD students in the same cohort to write a gender-related research dissertation. We provide evidence for three mechanisms that can explain the sizable treatment effects for PhD students and the differential effects across gender. The positive spillover effects to do gender-related research observed among female students are primarily attributed to collaborations and informal interactions with female peers. Conversely, the diminished or negative effects seen in male students appear to stem from competitive pressures. This analysis not only confirms the existence of knowledge diffusion and competitive dynamics but also unravel the specific conditions that amplify or mitigate these effects.

Innovation has been taken as the engine of economic growth (Romer 1990), and the underrepresentation of women in academia has aroused the concern for "*Missing Marie Curies*" and misallocation of talents (Hsieh et al. 2019). Moreover, a growing literature

has documented the lack of female in academia affects who and what gets studies in research (Koning, Samila, and Ferguson 2020, 2021; Nielsen et al. 2017; Truffa and Wong 2022). This paper joins the discussion by providing novel causal evidence about how female peers can help to close the knowledge gap with respect the gender-related research. Additionally, our research enriches the literature on peer effects by exploring the interplay between positive (learning and knowledge diffusion) and negative (competition) spillover effects, and how these dynamics influence the inclination of researchers to pursue gender-related topics. By examining these nuances, our study sheds light on the complex mechanisms through which gender composition in academic cohorts can shape the direction and diversity of research, ultimately contributing to a more inclusive and comprehensive body of academic knowledge.

This study carries significant policy implications, particularly for government officers and university administrators who are increasingly focused on enhancing the welfare and success of female scientists. Our findings align with other research in highlighting the benefits of a diverse and inclusive academic environment, notably in encouraging more gender-related research. However, our study also uncovers potential challenges associated with efforts to create gender diversity in academia. While the presence of female peers positively influences female PhD students to engage in gender-related research, our results indicate that male PhD students are less affected and, in some cases, may even be negatively influenced due to heightened competition and gender-based homophily in collaborations.

Addressing these complex dynamics requires strategic efforts from both universities and governmental bodies. Universities could play a pivotal role by fostering an environment that promotes both formal and informal interactions among students of different genders. This could help mitigate the effects of homophily and encourage a broader engagement with female-centric research topics. Additionally, governmental support, particularly in terms of funding, could be crucial in alleviating the competitive pressures in gender-related research fields. By providing more resources and financial backing, governments could help create a

more collaborative and less competitive atmosphere, encouraging scholars of all genders to contribute to research in these important areas. Implementing these measures could help balance the benefits of a gender-diverse academic environment, ensuring that the pursuit of diversity and inclusion does not inadvertently lead to new forms of imbalance or disincentives certain groups from engaging in gender-related research.

Our findings are not without limitations and therefore open the door to further research. We mainly focus on the supply side of gender-related research by investigating the peer effects. One area we briefly explore, but do not extensively analyze, is the demand-side of gender-related research. This includes factors like the varying benefits across different fields and elements such as public funding. Future research could enrich our understanding by delving deeper into these demand-side factors and their impact on research orientation and output. Besides, our reliance on ProQuest data introduces potential limitations, particularly concerning the accuracy of gender predictions made using software algorithms. Although we have attempted to mitigate these concerns through robustness checks, there is scope for further refinement in this area. Future researchers might find value in incorporating survey-based data sources, such as the Survey of Earned Doctorates (SED), to enhance the precision and reliability of gender-related data in academic research studies.

## References

- Acemoglu, Daron, David Y Yang, and Jie Zhou. 2021. “Political pressure and the direction of research: Evidence from china’s academia.” In *Working Paper*.
- Akcigit, Ufuk, Jeremy G Pearce, and Marta Prato. 2022. “Tapping into talent: Coupling education and innovation policies for economic growth.” *NBER Working Paper*.
- Anelli, Massimo, and Giovanni Peri. 2019. “The effects of high school peers’ gender on college major, college performance and income.” *The Economic Journal* 129 (618): 553–602.
- Atkin, David, M Keith Chen, and Anton Popov. 2022. *The returns to face-to-face interactions: Knowledge spillovers in Silicon Valley*. Technical report. National Bureau of Economic Research.
- Azoulay, Pierre, Christian Fons-Rosen, and Joshua S Graff Zivin. 2019. “Does science advance one funeral at a time?” *American Economic Review* 109 (8): 2889–2920.
- Borjas, George J, and Kirk B Doran. 2012. “The collapse of the Soviet Union and the productivity of American mathematicians.” *The Quarterly Journal of Economics* 127 (3): 1143–1203.
- Bostwick, Valerie K, and Bruce A Weinberg. 2022. “Nevertheless she persisted? Gender peer effects in doctoral STEM programs.” *Journal of Labor Economics* 40 (2): 397–436.
- Chai, Sen, and Richard B Freeman. 2019. “Temporary colocation and collaborative discovery: Who confers at conferences.” *Strategic Management Journal* 40 (13): 2138–2164.
- Chen, Siyu, and Zihan Hu. 2022. “How competition shapes peer effects: Evidence from a university in China.” *Available at SSRN 4012786*.
- Drago, Robert, and Gerald T Garvey. 1998. “Incentives for helping on the job: Theory and evidence.” *Journal of Labor Economics* 16 (1): 1–25.
- Dupas, Pascaline, Alicia Sasser Modestino, Muriel Niederle, Justin Wolfers, et al. 2021. *Gender and the dynamics of economics seminars*. Technical report. National Bureau of Economic Research.
- Fry, Caroline Viola. 2023. “Crisis and the trajectory of science: Evidence from the 2014 Ebola outbreak.” *Review of Economics and Statistics* 105 (4): 1028–1038.
- Greenberg, Jason, and Ethan Mollick. 2017. “Activist choice homophily and the crowdfunding of female founders.” *Administrative Science Quarterly* 62 (2): 341–374.
- Guryan, Jonathan, Kory Kroft, and Matthew J Notowidigdo. 2009. “Peer effects in the workplace: Evidence from random groupings in professional golf tournaments.” *American Economic Journal: Applied Economics* 1 (4): 34–68.
- Hasan, Sharique, and Rembrand Koning. 2019. “Conversations and idea generation: Evidence from a field experiment.” *Research Policy* 48 (9): 103811.
- Hill, Ryan. 2020. “Essays on the economics of science and innovation.” PhD diss., Massachusetts Institute of Technology.
- Hill, Ryan, and Carolyn Stein. 2019. “Scooped! Estimating rewards for priority in science.” *Job Market Paper*.
- Hofstra, Bas, Vivek V Kulkarni, Sebastian Munoz-Najar Galvez, Bryan He, Dan Jurafsky, and Daniel A McFarland. 2020. “The diversity–innovation paradox in science.” *Proceedings of the National Academy of Sciences* 117 (17): 9284–9291.
- Hoxby, Caroline M. 2000. *Peer effects in the classroom: Learning from gender and race variation*.
- Hsieh, Chang-Tai, Erik Hurst, Charles I Jones, and Peter J Klenow. 2019. “The allocation of talent and us economic growth.” *Econometrica* 87 (5): 1439–1474.
- Huang, Junming, Alexander J Gates, Roberta Sinatra, and Albert-László Barabási. 2020. “Historical comparison of gender inequality in scientific careers across countries and disciplines.” *Proceedings of the National Academy of Sciences* 117 (9): 4609–4616.
- Koffi, Marlène, and Matt Marx. 2023. *Cassatts in the Attic*. Technical report. National Bureau of Economic Research.
- Koning, Rembrand, Sampsa Samila, and John-Paul Ferguson. 2020. “Inventor Gender and the Direction of Invention.” In *AEA Papers and Proceedings*, 110:250–54.

- Koning, Rembrand, Sampsa Samila, and John-Paul Ferguson. 2021. “Who do we invent for? Patents by women focus more on women’s health, but few women get to invent.” *Science* 372 (6548): 1345–1348.
- Lavy, Victor, and Analia Schlosser. 2011. “Mechanisms and impacts of gender peer effects at school.” *American Economic Journal: Applied Economics* 3 (2): 1–33.
- Mouganie, Pierre, and Yaojing Wang. 2020. “High-performing peers and female STEM choices in school.” *Journal of Labor Economics* 38 (3): 805–841.
- Myers, Kyle. 2020. “The elasticity of science.” *American Economic Journal: Applied Economics* 12 (4): 103–134.
- Nielsen, Mathias Wullum, Jens Peter Andersen, Londa Schiebinger, and Jesper W Schneider. 2017. “One and a half million medical papers reveal a link between author gender and attention to gender and sex analysis.” *Nature human behaviour* 1 (11): 791–796.
- Priem, Jason, Heather Piwowar, and Richard Orr. 2022. “OpenAlex: A fully-open index of scholarly works, authors, venues, institutions, and concepts.” *arXiv preprint arXiv:2205.01833*.
- Roche, Maria P, Alexander Oettl, and Christian Catalini. 2023. “Proximate (Co-) Working: Knowledge Spillovers and Social Interactions.”
- Romer, Paul M. 1990. “Endogenous technological change.” *Journal of political Economy* 98 (5, Part 2): S71–S102.
- Schneeweis, Nicole, and Martina Zweimüller. 2012. “Girls, girls, girls: Gender composition and female school choice.” *Economics of Education review* 31 (4): 482–500.
- Sohn, Eunhee. 2021. “How local industry R&D shapes academic research: Evidence from the agricultural biotechnology revolution.” *Organization Science* 32 (3): 675–707.
- Stolper, Oscar, and Andreas Walter. 2019. “Birds of a feather: The impact of homophily on the propensity to follow financial advice.” *The Review of Financial Studies* 32 (2): 524–563.
- Teodoridis, Florenta. 2018. “Understanding team knowledge production: The interrelated roles of technology and expertise.” *Management Science* 64 (8): 3625–3648.
- Topaz, Chad M, and Shilad Sen. 2016. “Gender representation on journal editorial boards in the mathematical sciences.” *PLoS One* 11 (8): e0161357.
- Truffa, Francesca, and Ashley Wong. 2022. “Undergraduate Gender Diversity and Direction of Scientific Research.” *Working paper*.
- Wu, Alice H. 2020. “Gender bias among professionals: an identity-based interpretation.” *Review of Economics and Statistics* 102 (5): 867–880.
- Wuchty, Stefan, Benjamin F Jones, and Brian Uzzi. 2007. “The increasing dominance of teams in production of knowledge.” *Science* 316 (5827): 1036–1039.
- Yang, Yang, Tanya Y Tian, Teresa K Woodruff, Benjamin F Jones, and Brian Uzzi. 2022. “Gender-diverse teams produce more novel and higher-impact scientific ideas.” *Proceedings of the National Academy of Sciences* 119 (36): e2200841119.
- Zeltzer, Dan. 2020. “Gender homophily in referral networks: Consequences for the medicare physician earnings gap.” *American Economic Journal: Applied Economics* 12 (2): 169–197.
- Zhou, Sifan, Sen Chai, and Richard B Freeman. 2024. “Gender homophily: In-group citation preferences and the gender disadvantage.” *Research Policy* 53 (1): 104895.
- Zolas, Nikolas, Nathan Goldschlag, Ron Jarmin, Paula Stephan, Jason Owen-Smith, Rebecca F Rosen, Barbara McFadden Allen, Bruce A Weinberg, and Julia I Lane. 2015. “Wrapping it up in a person: Examining employment and earnings outcomes for Ph. D. recipients.” *Science* 350 (6266): 1367–1371.



# Figures

Figure 1: Example of Gender-Related Dissertation

Full Text | Dissertation or Thesis

## Acculturation, knowledge, beliefs, and preventive health care practices regarding breast care in female Chinese immigrants in New York metropolitan area

Chen, Wei-Ti. Columbia University ProQuest Dissertations Publishing, 2002. 3048108.

Full text - PDF  
Preview - PDF  
**Abstract/Details**  
Discovery Tips - Dissertations Demystified >>

### Abstract

Show highlighting  
Translate

Studies have found that breast cancer becomes a greater health problem as successive generations of Asian women live in the United States. The purpose of this study was to examine the relationships between acculturation level and breast cancer knowledge, cancer risk perception, health practice, and perceptions of health access.

This descriptive correlational cross-sectional study used a survey approach. Participation in this study was limited to Chinese immigrant women, aged 18 and over. The survey questionnaire was written in Chinese and self-administered. Study participants (N = 135) were recruited from two childbirth classes, one temple activity, two church gathering events, and four American Cancer Society's educational sessions.

Only "years of education", "marital status", and "household income" showed significant relationships to breast cancer risk knowledge level. The data indicate that women with a better knowledge of breast cancer risk are twice as likely to have higher income and have more education. The most knowledgeable women are less likely to be married and less likely to have partners compared to least knowledgeable group.

Full Text | Dissertation or Thesis

## DNA repair genes and breast cancer risk: Iowa Women's Health Study

Thyagarajan, Bharat. University of Minnesota ProQuest Dissertations Publishing, 2004. 3142653.

Full text - PDF  
Preview - PDF  
**Abstract/Details**  
Discovery Tips - Dissertations Demystified >>

### Abstract

Show highlighting  
Translate

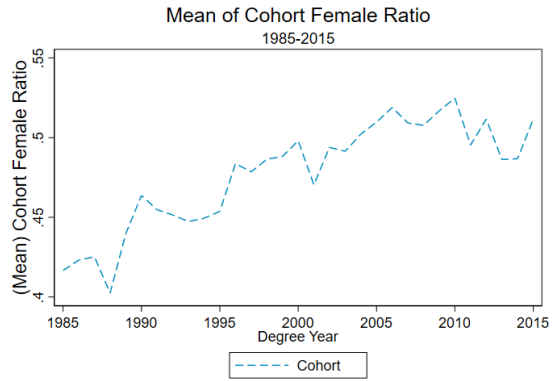
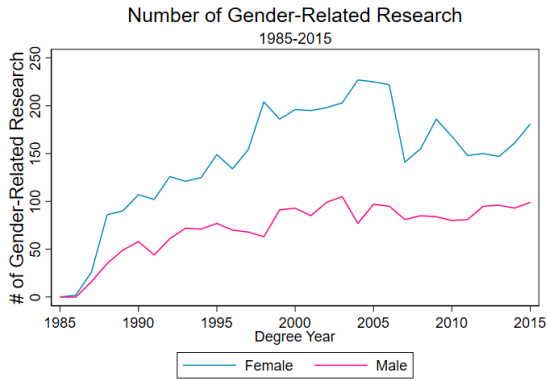
The current study focused on an evaluation of associations between several DNA repair gene SNPs and the risk of breast cancer in a large cohort, the Iowa Women's Study (IWHs). A pilot study that evaluated the utilization of various tissue types from paraffin-embedded tissue blocks, a primary source of DNA within IWHs, and a high-throughput genotyping platform, the PCR-INVADER assay found that the highest quantity and best quality DNA were obtained from benign lymph nodes and that the PCR-INVADER assay was an accurate and reliable assay to genotype DNA obtained from paraffin embedded tissues. The population-based allelic frequencies of twelve polymorphisms from seven DNA repair genes were evaluated in a sample of two hundred cancer-free women. Eight polymorphisms had an allele frequency >5% and were selected for further evaluation of their association with breast cancer. These included XPD-23 (Lys751Gln), XPD-10 (Asp312Asn), XPC (intron 9 insertion/deletion), XPA (A23G), XPG-15 (Asp1104His), XRCC1-6 (Arg194Trp), XRCC1-10 (Arg399Gln), and XRCC3-7 (Thr241Met) polymorphisms. The association between these polymorphisms and breast cancer risk was evaluated in a nested case cohort study of 460 breast cancer cases and 324 cancer-free controls. Samples were genotyped using one of three genotyping platforms—PCR-RFLP, PCR-INVADER or Sequenom. The Lys751Gln polymorphism in XPD (Lys/Gln + Gln/Gln) was found to be associated with a statistically significant decreased breast cancer risk (OR = 0.60, 95% CI: 0.45–0.81) as compared to the Lys/Lys genotype. The A23G polymorphism in XPA (A/G + G/G) also was associated with a decreased breast cancer risk (OR = 0.75, 95% CI: 0.56–1.01) as compared to the A/A genotype. None of the other six polymorphisms in DNA repair genes were associated with breast cancer risk. We also adapted a flow cytometry based DNA repair assay, for the measurement of nucleotide excision repair (NER) activity. Specificity and sensitivity of this assay were evaluated for the measurement of NER activity and estimates for NER activity were obtained in each cell cycle phase. Finally, we obtained within and between person variability estimates for numerous variables used to estimate NER activity and evaluated the utility of these variables for routine use.

Notes. The figures shows titles and abstracts of two gender-related dissertations in the ProQuest Database.

**Figure 2:** Descriptive Stats: Number of Gender-Related Research and Cohort Female Ratio

**(a)** Number of Gender-Related Research

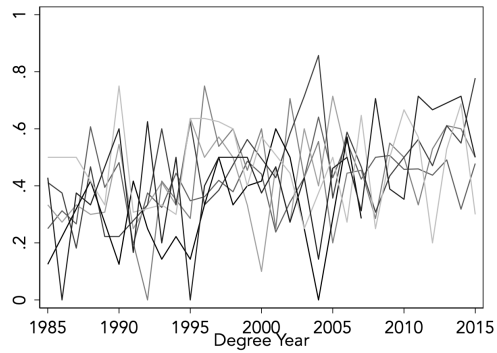
**(b)** Mean of Cohort Female Ratio



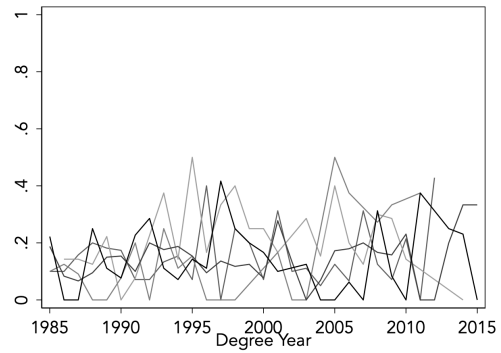
*Note:* This figure shows the number of gender-related research and the mean of cohort female ratio from 1985 to 2015.

**Figure 3:** Within-Program Female Ratio Fluctuations (Example: Cornell University)

**(a)** Health and Biology



**(b)** Engineering



*Note:* This figure keeps track the female ratio fluctuations within all PhD programs of (a) health and biology and (b) engineering in Cornell University.

# Tables

**Table 1:** Summary Statistics

Var.	Female		Male	
	Mean	S.D.	Mean	S.D.
<b>Panel A: All PhDs</b>	N=44830		N=48960	
Gender-Related Dissertation	0.0911	0.288	0.0387	0.193
Gender-Related Intensity	0.0010	0.004	0.0003	0.002
Cohort Female Ratio	0.5493	0.184	0.4119	0.160
Cohort Size	19.1094	15.546	19.8487	15.760
Advisor Gender-Related Research	0.4911	0.500	0.4270	0.495
Total Words (Title + Abstract)	185.8843	111.213	187.0898	110.976
<b>Panel B: PhDs Matched to OpenAlex</b>	N=11378		N=12952	
Gender-Related Dissertation	0.0627	0.242	0.0314	0.174
Gender-Related Intensity	0.0006	0.003	0.0003	0.002
Cohort Female Ratio	0.5299	0.169	0.4126	0.155
Cohort Size	20.4060	16.429	20.7027	16.243
Advisor Gender-Related Research	0.4611	0.499	0.4191	0.493
Total Words (Title + Abstract)	189.7808	111.607	188.0601	111.835
Num. of Papers				
5+ years before graduation	1.6223	10.596	2.8735	14.845
1-5 year before graduation	5.9726	21.687	9.0311	27.830
0-5 years after graduation	14.9692	39.488	21.4454	49.133
6-10 years after graduation	13.6789	41.624	20.7253	51.886
Num. of Gender-Related Papers				
5+ years before graduation	0.0189	0.218	0.0187	0.228
1-5 year before graduation	0.0828	0.504	0.0652	0.536
0-5 years after graduation	0.2957	1.353	0.2181	1.395
6-10 years after graduation	0.2820	1.380	0.2254	1.320

*Notes.* Each observation is at individual PhD level. Panel A includes all observations in our main estimation, and Panel B only include PhDs that are matched to the OpenAlex database. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our keyword list in title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. *Cohort Size* is the number of PhD students graduating in the same year within a given PhD program. *Advisor Gender-Related Research* is a binary variable which switches to one when the PhD's advisor has ever done some female-focused research. *Total Words (Title + Abstract)* is the number of words in the title and abstract of the PhD's dissertation. *Num. of Papers* and *Num. of Gender-Related Papers* is the number of all papers and the number of all female-focused papers, respectively, published in a given period.

**Table 2:** The Impacts of Cohort Female Ratio on Gender-Related Dissertation Production

	(1)	(2)	(3)	(4)	(5)	(6)
	Gender-Related Dissertation			Gender-Related Intensity		
Cohort Female Ratio	0.01777*** (0.00632)			0.00011 (0.00007)		
Female $\times$ Cohort Female Ratio		0.06075*** (0.01068)	0.05845*** (0.01048)		0.00066*** (0.00013)	0.00065*** (0.00013)
Male $\times$ Cohort Female Ratio		-0.01978** (0.00783)	-0.02124*** (0.00791)		-0.00038*** (0.00011)	-0.00037*** (0.00010)
Program FE	Yes	Yes	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes	Yes	Yes
Gender Dummy	Yes	Yes	Yes	Yes	Yes	Yes
Controls	No	No	Yes	No	No	Yes
Mean of D.V.	0.06379	0.06379	0.06379	0.00064	0.00064	0.00064
Adjusted R-squared	0.107	0.107	0.122	0.091	0.091	0.096
Observations	93790	93790	93790	93790	93790	93790

*Notes.* Each observation is at individual PhD level. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our keyword list in title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. All models include controls for the gender dummy, cohort size, advisor's research focus, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .

**Table 3:** The Impacts of Cohort Female Ratio on Gender-Related Research Production

	(1)	(2)	(3)	(4)
	Years Relative to Graduation			
	< -5	-5 to -1	0-5	6-10
<b>Panel A - D.V.: Have One or More Gender-Related Papers</b>				
Female $\times$ Cohort Female Ratio	0.00318 (0.00748)	0.02023 (0.01505)	0.04248** (0.01909)	0.02587 (0.01794)
Male $\times$ Cohort Female Ratio	0.00407 (0.00757)	0.00402 (0.01213)	0.01042 (0.01529)	-0.02289 (0.01434)
Observations	24330	24330	24330	24330
<b>Panel B - D.V.: Ratio of Gender-Related Papers</b>				
Female $\times$ Cohort Female Ratio	0.00420 (0.01724)	0.02164** (0.00891)	0.01649*** (0.00622)	0.01721** (0.00788)
Male $\times$ Cohort Female Ratio	-0.00060 (0.01384)	-0.00207 (0.00486)	-0.00298 (0.00394)	-0.00890* (0.00479)
Observations	6068	19637	23051	15873
<b>Panel C - D.V.: Number of Gender-Related Papers</b>				
Female $\times$ Cohort Female Ratio	0.57920 (0.88725)	0.35646 (0.47028)	0.30536 (0.28493)	0.41717 (0.31645)
Male $\times$ Cohort Female Ratio	0.68583 (0.75041)	0.79000 (0.67764)	0.18255 (0.37769)	0.38694 (0.27495)
Observations	3489	15769	21230	14611
Program FE	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes
Gender Dummy	Yes	Yes	Yes	Yes
Controls	Yes	Yes	Yes	Yes

*Notes.* Each observation is at individual PhD level. *Gender-Related Papers* are papers which contain any gender-related words as defined in our keyword list in title. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. All models include controls for the gender dummy, cohort size, advisor's research focus, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .

**Table 4:** Mechanism - Collaboration

	(1)	(2)	(3)	(4)
	Papers Published			
	in F.-Dom. Teams	w. Female Coauthors	w. Female Coauthors Same Inst.	w. Female Coauthors Same Cohort
<b>Panel A - D.V.: Have One or More Gender-Related Papers</b>				
Female × Cohort Female Ratio	0.04074** (0.01895)	0.04804** (0.01917)	0.00227 (0.00665)	0.00162 (0.00243)
Male × Cohort Female Ratio	-0.00674 (0.01661)	-0.00873 (0.01698)	0.00227 (0.00478)	0.00052 (0.00193)
Observations	24330	24330	24330	24330
<b>Panel B - D.V.: Ratio of Gender-Related Papers</b>				
Female × Cohort Female Ratio	0.01447*** (0.00453)	0.01539*** (0.00467)	-0.00045 (0.00044)	-0.00004 (0.00008)
Male × Cohort Female Ratio	-0.00628*** (0.00243)	-0.00630** (0.00261)	-0.00004 (0.00031)	0.00006 (0.00009)
Observations	23659	23659	23659	23659
<b>Panel C - D.V.: Number of Gender-Related Papers</b>				
Female × Cohort Female Ratio	0.23125 (0.33770)	0.25468 (0.32370)	0.53875 (1.11841)	2.41771 (2.14837)
Male × Cohort Female Ratio	0.10194 (0.41761)	0.06583 (0.40595)	0.97483 (0.73174)	0.96017 (1.61731)
Observations	22620	22675	5340	1044
Program FE	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes
Gender Dummy	Yes	Yes	Yes	Yes
Controls	Yes	Yes	Yes	Yes

*Notes.* Each observation is at individual PhD level. *Gender-Related Papers* are papers which contain any gender-related words as defined in our keyword list in title. *Papers Published in F.-Dom. Teams* are papers published in female-dominant teams where female authors are more than male authors. *Papers Published w. Female Coauthors* are papers with at least one female coauthors. *Papers Published w. Female Coauthors Same Inst.* are papers with at least one female PhDs from same or nearby cohorts. *Papers Published w. Female Coauthors Same Cohort* are papers with at least one female coauthors from same cohorts. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. All models include controls for the gender dummy, cohort size, advisor’s research focus, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .

**Table 5:** Mechanism - Informal Interactions

	(1)	(2)
	Gender-Related Dissertation	Gender-Related Intensity
Female × Cohort Female Ratio	0.05138*** (0.01022)	0.00054*** (0.00013)
Male × Cohort Female Ratio	-0.02347*** (0.00817)	-0.00043*** (0.00010)
<i>Interaction Effects</i>		
Female × Cohort Female Ratio × Female Peer w. Same Advisor	0.02778* (0.01562)	0.00047* (0.00025)
Male × Cohort Female Ratio × Female Peer w. Same Advisor	0.01576 (0.01757)	0.00037 (0.00025)
Program FE	Yes	Yes
Year FE	Yes	Yes
Gender Dummy	Yes	Yes
Controls	Yes	Yes
Mean of D.V.	0.06379	0.00064
Adjusted R-squared	0.123	0.097
Observations	93790	93790

*Notes.* Each observation is at individual PhD level. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our keyword list in title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. *Female Peer w. Same Advisor* is a binary variable which switches to one when the focal PhD share the same advisor with another female peer in the same cohort. All models include controls for the gender dummy, cohort size, advisor’s research focus, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .



**Table 6:** Mechanism - Competition

	(1)	(2)	(3)	(4)
	Gender-Related Dissertation		Gender-Related Intensity	
Female $\times$ Cohort Female Ratio	0.04376*** (0.01027)	0.04257*** (0.01023)	0.00051*** (0.00013)	0.00051*** (0.00013)
Male $\times$ Cohort Female Ratio	-0.01270* (0.00728)	-0.01147 (0.00718)	-0.00029*** (0.00009)	-0.00026*** (0.00009)
<i>Interaction Effects</i>				
Female $\times$ Cohort Female Ratio $\times$ Top 10 % Female Peer G-R	0.03395 (0.03098)		0.00032 (0.00042)	
Male $\times$ Cohort Female Ratio $\times$ Top 10 % Female Peer G-R	-0.10718*** (0.03758)		-0.00106** (0.00047)	
<i>Interaction Effects</i>				
Female $\times$ Cohort Female Ratio $\times$ Top 25 % Female Peer G-R		0.03515 (0.02713)		0.00029 (0.00038)
Male $\times$ Cohort Female Ratio $\times$ Top 25 % Female Peer G-R		-0.08515** (0.03321)		-0.00097** (0.00041)
Program FE	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes
Gender Dummy	Yes	Yes	Yes	Yes
Controls	Yes	Yes	Yes	Yes
Mean of D.V.	0.06379	0.06379	0.00064	0.00064
Adjusted R-squared	0.122	0.122	0.097	0.097
Observations	93790	93790	93790	93790

*Notes.* Each observation is at individual PhD level. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our keyword list in title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. *Top 10 % Female Peer G-R* is a binary variable which switches to one when the focal PhD has a star female cohort peer who has citation in the top 10 percentile of the field and is doing FFR in her dissertation. *Top 25 % Female Peer G-R* is a binary variable which switches to one when the focal PhD has a star female cohort peer who has citation in the top 25 percentile of the field and is doing FFR in her dissertation. All models include controls for the gender dummy, cohort size, advisor's research focus, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .

# A Data Appendix

## A.1 Keywords List

- Baseline List: woman, female, lady, feminism, feminine, femininity, girl, pregnancy, pregnant, gender, sex , wife, daughter, mother;
- Short List: woman, female, sex, gender;
- Extended List: woman, female, lady, feminism, feminine, femininity, girl, pregnancy, pregnant, gender, sex , wife, daughter, mother, gynecological, gynecology, menopause, menstruation, menstrual, urinary tract, vaginosis, vaginitis, uterine fibroids, pregnant, pregnancy, preterm, endometriosis, ovary, ovarian, cervical, turner syndrome, rett syndrome.

## B Table Appendix

**Table B.1:** Robustness Test - First Order Auto-Regression

	(1)	(2)	(3)	(4)
	Cohort Female Ratio			
	Health & Bio	Engineering	Health & Bio	Engineering
Cohort Female Ratio (1-Year Lag)	0.00171 (0.00783)	0.00314 (0.01270)		
Cohort Female Ratio (5-Year Lag)			-0.01266 (0.01121)	-0.01824 (0.01374)
Program FE	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes
p-value: $H_0$			<0.01	<0.01
Adj. $R^2$	0.39	0.22	0.37	0.21
Obs.	9,954	8,088	8,627	6,909

*Notes.* This table shows the estimation results of first-order auto-regression 2, where  $h = 1$  for Column (1) and (2), and  $h = 5$  for Column (3) and (4). Standard errors in parentheses are clustered at the university level. The last two columns report the p-value for null-hypothesis of autocorrelation coefficient  $\beta = 1$ .

**Table B.2:** Heterogeneity Test - Advisors' Research Focus

	(1)	(2)	(3)	(4)
	Gender-Related Dissertation		Gender-Related Intensity	
	Advisor - Gender-Related	Advisor - Not Gender-Related	Advisor - Gender-Related	Advisor - Not Gender-Related
Female $\times$ Cohort Female Ratio	0.06999*** (0.01598)	0.03061*** (0.01073)	0.00076*** (0.00020)	0.00033** (0.00016)
Male $\times$ Cohort Female Ratio	-0.04500*** (0.01314)	0.01147 (0.00739)	-0.00075*** (0.00017)	0.00010 (0.00011)
Program FE	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes
Gender Dummy	Yes	Yes	Yes	Yes
Controls	Yes	Yes	Yes	Yes
Mean of D.V.	0.10603	0.02815	0.00107	0.00027
Adjusted R-squared	0.127	0.099	0.096	0.104
Observations	42882	50838	42882	50838

*Notes.* Each observation is at individual PhD level. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our keyword list in title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. *Advisor - Gender-Related* indicates the PhD's advisor has done gender-related research and *Advisor - Not Gender-Related* indicates the PhD's advisor hasn't done any gender-related research. All models include controls for the gender dummy, cohort size, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .

**Table B.3:** Heterogeneity Test - Field Difference

	(1)	(2)	(3)	(4)
	Gender-Related Dissertation		Gender-Related Intensity	
	Basic Research	Applied research	Basic Research	Applied research
Female $\times$ Cohort Female Ratio	0.00628 (0.00849)	0.07469*** (0.02282)	-0.00000 (0.00009)	0.00064** (0.00028)
Male $\times$ Cohort Female Ratio	0.00421 (0.00656)	-0.02316 (0.02090)	0.00004 (0.00007)	-0.00070** (0.00028)
Program FE	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes
Gender Dummy	Yes	Yes	Yes	Yes
Controls	Yes	Yes	Yes	Yes
Mean of D.V.	0.03441	0.15786	0.00027	0.00181
Adjusted R-squared	0.055	0.115	0.040	0.080
Observations	71466	22322	71466	22322

*Notes.* Each observation is at individual PhD level. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our keyword list in title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. *Basic Research* indicates the PhD graduates from a biology program with CIP code 26 and *Applied* indicates the PhD graduates from a healthcare program with CIP code 51. All models include controls for the gender dummy, cohort size, advisor's research focus, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .

**Table B.4:** Heterogeneity Test - Cohort Size

	(1)	(2)	(3)	(4)
	Gender-Related Dissertation		Gender-Related Intensity	
	Cohort Size < 15	Cohort Size ≥ 15	Cohort Size < 15	Cohort Size ≥ 15
Female × Cohort Female Ratio	0.04827*** (0.01215)	0.07370*** (0.02193)	0.00045*** (0.00015)	0.00103*** (0.00029)
Male × Cohort Female Ratio	-0.01201 (0.00846)	-0.05410*** (0.01590)	-0.00027** (0.00011)	-0.00075*** (0.00021)
Program FE	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes
Gender Dummy	Yes	Yes	Yes	Yes
Controls	Yes	Yes	Yes	Yes
Mean of D.V.	0.07726	0.04852	0.00078	0.00048
Adjusted R-squared	0.128	0.116	0.101	0.093
Observations	49846	43941	49846	43941

*Notes.* Each observation is at individual PhD level. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our keyword list in title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. All models include controls for the gender dummy, advisor’s research focus, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\*  $p < 0.05$ ; \*\*\*  $p < 0.01$ .

**Table B.5:** Robustness Test - OpenAlex Matched Sample

	(1)	(2)
	Gender-Related Dissertation	Gender-Related Intensity
Female $\times$ Cohort Female Ratio	0.04354*** (0.01522)	0.00035* (0.00019)
Male $\times$ Cohort Female Ratio	0.00805 (0.01230)	-0.00012 (0.00018)
Program FE	Yes	Yes
Year FE	Yes	Yes
Gender Dummy	Yes	Yes
Controls	Yes	Yes
Mean of D.V.	0.04603	0.00043
Adjusted R-squared	0.134	0.148
Observations	24330	24330

*Notes.* Each observation is at individual PhD level and we only include observations of PhDs that are matched to OpenAlex database. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our keyword list in title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. All models include controls for the gender dummy, cohort size, advisor's research focus, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .

**Table B.6:** Robustness Test - Female Classification

	(1)	(2)	(3)	(4)	(5)	(6)
	Gender-Related Dissertation			Gender-Related Intensity		
	Prob $\geq$ 0.85	Prob $\geq$ 0.90	Prob $\geq$ 0.95	Prob $\geq$ 0.85	Prob $\geq$ 0.90	Prob $\geq$ 0.95
Female $\times$ Cohort Female Ratio	0.06070*** (0.01119)	0.06260*** (0.01147)	0.06363*** (0.01202)	0.00069*** (0.00014)	0.00069*** (0.00014)	0.00071*** (0.00015)
Male $\times$ Cohort Female Ratio	-0.02425*** (0.00832)	-0.02468*** (0.00849)	-0.02779*** (0.00880)	-0.00040*** (0.00011)	-0.00039*** (0.00012)	-0.00044*** (0.00013)
Program FE	Yes	Yes	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes	Yes	Yes
Gender Dummy	Yes	Yes	Yes	Yes	Yes	Yes
Controls	Yes	Yes	Yes	Yes	Yes	Yes
Mean of D.V.	0.06592	0.06623	0.06575	0.00066	0.00066	0.00066
Adjusted R-squared	0.124	0.125	0.125	0.098	0.098	0.099
Observations	84313	81302	76389	84313	81302	76389

*Notes.* Each observation is at individual PhD level. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our keyword list in title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. *Prob* is the probability that the gender classification is reliable. All models include controls for the gender dummy, cohort size, advisor’s research focus, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .



**Table B.7:** Robustness Test - Gender-Related Research Definition (Title or Abstract)

	(1)	(2)	(3)	(4)
	Gender-Related Dissertation		Gender-Related Intensity	
	Only Title	Only Abstract	Only Title	Only Abstract
Female $\times$ Cohort Female Ratio	0.02120*** (0.00475)	0.05563*** (0.01026)	0.00163*** (0.00042)	0.00052*** (0.00012)
Male $\times$ Cohort Female Ratio	-0.01188*** (0.00348)	-0.02213*** (0.00769)	-0.00105*** (0.00028)	-0.00037*** (0.00010)
Program FE	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes
Gender Dummy	Yes	Yes	Yes	Yes
Controls	Yes	Yes	Yes	Yes
Mean of D.V.	0.01542	0.06126	0.00119	0.00056
Adjusted R-squared	0.059	0.118	0.053	0.096
Observations	93790	93790	93790	93790

*Notes.* Each observation is at individual PhD level. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our keyword list in either title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. *Prob* is the probability that the gender classification is reliable. All models include controls for the gender dummy, cohort size, advisor's research focus, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .

**Table B.8:** Robustness Test - Gender-Related Research Definition (Word List)

	(1)	(2)	(3)	(4)
	Gender-Related Dissertation		Gender-Related Intensity	
	Short List	Extended List	Short List	Extended List
Female $\times$ Cohort Female Ratio	0.05512*** (0.00989)	0.06303*** (0.00571)	0.00060*** (0.00011)	0.00077*** (0.00007)
Male $\times$ Cohort Female Ratio	-0.01918*** (0.00723)	-0.01771** (0.00436)	-0.00031*** (0.00009)	-0.00045*** (0.00005)
Program FE	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes
Gender Dummy	Yes	Yes	Yes	Yes
Controls	Yes	Yes	Yes	Yes
Mean of D.V.	0.05346	0.07181	0.00050	0.00083
Adjusted R-squared	0.123	0.044	0.089	0.039
Observations	93790	93790	93790	93790

*Notes.* Each observation is at individual PhD level. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our alternative keyword lists in title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. *Prob* is the probability that the gender classification is reliable. All models include controls for the gender dummy, cohort size, advisor's research focus, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .

**Table B.9:** Robustness Test - Cohort Female Ratio

	(1)	(2)	(3)	(4)	(5)	(6)
	Gender-Related Dissertation			Gender-Related Intensity		
Female × Num. of Female Peers	0.00025 (0.00044)			0.00000 (0.00001)		
Male × Num. of Female Peers	0.00038 (0.00035)			0.00000 (0.00000)		
Female × Cohort Female Ratio (Neighbor included)		0.07188*** (0.01657)			0.00100*** (0.00022)	
Male × Cohort Female Ratio (Neighbor included)		-0.05891*** (0.01359)			-0.00079*** (0.00019)	
Female × Female-to-Male Ratio			0.00484*** (0.00152)			0.00004* (0.00002)
Male × Female-to-Male Ratio			-0.00791*** (0.00190)			-0.00014*** (0.00002)
Program FE	Yes	Yes	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes	Yes	Yes
Gender Dummy	Yes	Yes	Yes	Yes	Yes	Yes
Controls	Yes	Yes	Yes	Yes	Yes	Yes
Mean of D.V.	0.06379	0.06379	0.06379	0.00064	0.00064	0.00064
Adjusted R-squared	0.122	0.123	0.118	0.096	0.097	0.094
Observations	93790	93790	92565	93790	93790	92565

*Notes.* Each observation is at individual PhD level. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our keyword list in title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. *Num. of Female Peers* is the number of female peers in the cohort, *Cohort Female Ratio (Neighbor included)* is the ratio of female peers in the graduated and neighbor cohorts, and *Female-to-Male Ratio* is the ratio of females to males in the cohort. All models include controls for the gender dummy, cohort size, advisor’s research focus, total number of words in title and abstract), and fixed effects (program and year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .

**Table B.10:** Robustness Test - Alternative Models

	(1)	(2)	(3)	(4)
	Gender-Related Dissertation			
Female $\times$ Cohort Female Ratio	0.05511*** (0.01025)	0.05934*** (0.01051)	0.01667*** (0.00522)	0.01746*** (0.00522)
Male $\times$ Cohort Female Ratio	-0.02608*** (0.00877)	-0.02068*** (0.00788)	0.00324 (0.00585)	0.00353 (0.00587)
Model	OLS	OLS	Probit	Probit
Program FE	Yes		Yes	
Year FE			Yes	
University $\times$ Year FE	Yes			
Field $\times$ Year FE	Yes			
Program Linear Trend		Yes		Yes
Adj. $R^2$	0.12	0.11		
Obs.	93,749	93,790	88,695	88,697

*Notes.* Each observation is at individual PhD level. *Gender-Related Dissertation* is an indicator variable which equals to one if the dissertation contains any gender-related words as defined in our keyword list in title or abstract. *Gender-Related Intensity* is the ratio of gender-related words in title and abstract of each dissertation. *Cohort Female Ratio* is the proportion of female students graduating in the same year within a given PhD program. All models include controls for the gender dummy, cohort size, advisor's research focus, total number of words in title and abstract). Column(1) and (2) are estimated using OLS, while Column (3) and (3) are estimated using Probit. Column (1) and (4) control for program-specific linear trend (indicator of program interacted with year). Standard errors in parentheses are clustered at the program level. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .

**Table B.11:** Randomization Test - Guryan, Kroft, and Notowidigdo 2009

	(1)	(2)
	Female	
$\overline{\text{Female}}_{-i,u,s,t}$	0.00651 (0.01210)	0.00572 (0.01205)
$\overline{\text{Female}}_{-i,u,s}$	-13.40353*** (1.15045)	-13.39283*** (1.15023)
Program FE	Yes	Yes
Year FE	Yes	Yes
Controls		Yes
Adj. $R^2$	.17	.17
Obs.	93,539	93,539

*Notes.* Each observation is at individual PhD level. *Female* is an indicator variable for female PhD.  $\overline{\text{Female}}_{-i,u,s,t}$  and  $\overline{\text{Female}}_{-i,u,s}$  are the leave-me-out mean of female at the cohort and program level. Column (2) add controls as the baseline specification 1. \* $p < 0.10$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ .