# Competing for Scientific Talent: Industry vs. Academia*

Justine Boudou

Harvard Business School

May 29, 2024

**Abstract**

On what margins do Industry and Academia compete for talent? In this paper, I study the allocation of PhDs upon graduation between Industry and Academia and its consequences on earnings and publications. I use my results to highlight heterogeneity across individuals and characterize the pool of talent that firms and universities compete for. Using an administrative dataset covering more than 40 cohorts of Ph.D. graduates with information about their earnings and their scientific production, I first estimate the causal impact of joining the private sector vs Academia on earnings and publication output. I account for selection on levels at the individual-level by using an instrumental variable strategy that exploits variation over time in firms' vs universities' demand for PhDs in the same major. I account for selection on gains by estimating Marginal Treatment Effects and find substantial heterogeneity in treatment effects that I use to characterize individuals in each sector. I then identify the margin of competition between firms and universities and how it varies across gender, nationality and field.

**Keywords**: *Human capital, Allocation, Competition, Heterogeneity, Innovation, Talent, Hiring*

---

# 1 Introduction

Human capital is a key source of competitive advantage for organizations (Campbell, Coff and Kryscynski, 2012; Coff, 1997). Highly-skilled labor in particular - defined in this paper as individuals with a doctorate degree - is a critical input for the value creation and capture of two types of organizations: firms and universities. Firms rely on scientists' specialized skills to drive their Research and Development (R&D) efforts and generate cutting-edge technological innovations. Firms' demand for scientists has become particularly salient over the last decades as economies have become increasingly knowledge-based and the rise of new and complex technologies require a complementary educated workforce (Cappelli, 2012, 2019; Ehrenberg, 1992; Stephan et al., 2004). While 30% of US R&D was funded by businesses in 1964, this number rose to more than 70% by 2020, making the hiring of scientists a critical strategic consideration for firms.[1] On the other hand, universities hire scientists as their primary source of labor to push forward the frontier of knowledge, bolster institutional prestige and educate the next generation of workers. Because individuals are heterogeneous in skills and do not randomly choose their sector of employment, firms and universities might get access to different yet overlapping pools of talent. However, while existing research highlights the importance of hiring and talent retention strategies (Agarwal et al., 2009; Black et al., 2024; Byun et al., 2018; Gambardella et al., 2015; Ganco et al., 2015; Nagle and Teodoridis, 2020; Ng and Stuart, 2022), less is known about the *allocation* of scientists across sectors and the margins on which firms and universities compete for talent. Who selects into Industry vs Academia? And on what margins do firms and universities compete for in the early-career market for high-skilled labor? I shed light on these questions by studying the allocation of PhDs upon graduation between Industry and Academia and its relationship with earnings and publications outcomes.

Answering these questions is not straightforward for three main reasons: (i) it requires individual-level information about scientists' performance; (ii) it requires an approach that accounts for the endogenous forces that steer scientists to enter one sector or the other; (iii) it requires an approach that can handle the significant heterogeneity across scientists in terms of skills and tastes for both sectors. I make progress on all three of these dimensions and provide new insights to this question of talent competition between Industry and Academia. First, I use a new administrative dataset that matches demographics and earnings information about doctorate holders from all STEM majors who graduated between 1970 and 2013 to their publication output. To understand endogeneity, I first present a simple model of sorting that incorporates heterogeneity in individuals' ability and preferences which creates selection into sector of employment. I then address the endogeneity that arises from selection in several ways. I start with a fixed-effects model relying on variation in sorting for individuals in the same university department, defined as the combination of doctoral institution and PhD major (e.g., Economics at Harvard). Under the identifying assumption that unobservable determinants of individuals' outcomes within a department are uncorrelated with their unobserved productivity in Industry vs Academia, I find that joining Industry vs Academia at graduation

---

increases earnings by about 30% over the career and decreases publications by about 60%. I then account for selection at the individual-level by applying an instrumental variable strategy that exploits variation over time in firms' vs universities' demand for PhDs in the same major. Under the identifying assumption that individuals' productivity in Industry vs Academia does not differ across years in the same major, the results imply that the population of individuals affected by my instrument publish about 77% less by joining the private sector vs Academia, while I do not find any statistical differences for the earnings outcome.

In order to address heterogeneity in treatment effects across individuals as far as possible, I then estimate Marginal Treatment Effects (MTE). This enables to calculate individual-specific treatment effects and their associated potential outcomes as a function of individuals' latent propensity of starting their career in Industry vs Academia upon graduation. I find substantial heterogeneity in treatment effects, both for earnings and publications, primarily driven by unobserved characteristics. Individuals who start their career in the private sector have higher earnings gains of joining Industry vs Academia than individuals who start in Academia, both because they earn less in Academia and earn more in Industry. They also have a lower publication decrease of joining Industry vs Academia than individuals who start in Academia. This result appears driven by the fact that they publish more if they start their career in Industry, which I relate to a higher probability of joining larger firms. I show that the fact that those who start in Industry earn less than those who start in Academia despite not exhibiting a lower potential to publish in the academic sector suggests that they lack non-research related skills valued by universities. Female tend to publish more than males in both sectors. International students tend to earn more than Americans in both sectors.

I then characterize the overlap between the pool of individuals who start in Industry and the pool of individuals who start in Academia, which represents the margin of competition between firms and universities. I show that the pool of talent accessible to firms is small, but is made of individuals who appear highly valuable for firms. Conversely, the pool of talent accessible to universities is large, but with varying Academic-related quality. Stars, defined as individuals who would be on the right-tail of the earnings distribution in both sectors, primarily start their career in Academia. I show that the margin of competition results from both demand and supply mechanisms. Firms start by hiring individuals who are highly valuable for them but have limited academic sector options. As their demand for PhDs increases, firms will progressively move towards the pool of stars. However, since stars have strong alternatives in Academia and a relatively higher preference for this sector, hiring becomes more difficult. Hence, firms move back to hiring individuals with lower academic sector options, albeit slightly less valuable for their specific needs. I then show how the margin of competition between firms and universities varies by gender, nationality, field and time.

Overall, my findings provide insights for firms' and universities' human capital strategy. I show that there is vast heterogeneity in individuals' performance and that firms and universities have access to overlapping pools of talent. I also identify the margin of competition and shed light on the mechanisms behind its specific pattern. I also show that stars predominantly start in the academic sector, which has importance consequences for firms' hiring strategy. However, this study is not

without limitations. In particular, I try to eliminate selection issues to estimate causal treatment effects and characterize the different pools of talent, but I do not study the exact mechanisms behind individuals' sorting. I also can only speak to the quantity of publications but not their content, providing insights primarily to the *rate* but not the *direction* of innovation.

## 2 Conceptual Framework

### 2.1 Background

The literature on strategic human capital emphasizes the critical role of general and firm-specific human capital in driving competitive advantage (Barney, 1991; Byun et al., 2018; Coff and Kryscynski, 2011). Highly-skilled labor, and in particular scientists, plays a pivotal role in driving innovation, productivity and overall organizational performance (Agrawal et al., 2017; Azoulay et al., 2019; Campbell, Ganco, Franco and Agarwal, 2012; Hess and Rothaermel, 2011; Kehoe and Tzabbar, 2015; Oettl, 2012; Romer, 1990) and the demand for their rare, specialized skills has been increasing over time (Cappelli, 2012, 2019). As a consequence, a large body of research has highlighted the importance of attracting and retaining top talent, focusing on the drivers and consequences of knowledge workers' mobility (Agarwal et al., 2009; Ganco et al., 2015; Starr et al., 2019), the importance of hiring decisions for firms' competitive advantage (Black et al., 2024; Byun et al., 2018; Nagle and Teodoridis, 2020) as well as the variety of incentives that organizations can use to manage talent (Gambardella et al., 2015).

While highly-skilled labor is a critical input for firms' competitive advantage, they are also one of the primary sources of labor for another type of organization: universities. A large literature has highlighted the difference in institutional norms between both sectors (Agarwal and Ohyama, 2013; Dasgupta and David, 1994) and the different roles and incentives that Industry and Academia have in the production of knowledge and innovation (Aghion et al., 2008; Stephan, 2012; Stern, 2004). Universities primarily seek scientists to push forward the frontier of knowledge, while firms primarily hire scientists to drive their R&D efforts. Importantly, individuals are also heterogeneous in their skills and preferences for each sector (Agarwal and Ohyama, 2013; Roach and Sauermann, 2010, 2024; Roche, 2023; Shu, 2016), creating selection into sector of employment. As a consequence, firms and universities may have access to different pools of talent which might overlap. Yet, little is known about the allocation of scientists across sectors and the margins on which universities and firms compete for talent. In this paper, I specifically focus on two measures of performance to identify this margin, earnings and publications, that can also be used as hiring incentives for firms. While pecuniary returns constitute a classical hiring incentives, scientists have also be shown to value the possibility to publish. In particular, Stern (2004) studies differences in firms that allow and do not allow to publish and shows that scientists with a higher 'taste for science' are willing to accept a compensating wage differential to work for science oriented firms.

Characterizing the pools of talent that firms and universities have access to, as well as the pool of talent that they compete for, is essential for effectively managing human capital and maximizing

organizational performance. Moreover, analyzing how this margin varies by field, gender or nationality may inform the literature on high-skilled immigration (Aobdia et al., 2018; Dimmock et al., 2022; Glennon, 2024) and gender-related labor differences (Bao, 2024).

## 2.2 Model of selection

Let individuals be indexed by $i$ and time be indexed by $t \in [0, T]$. Let $j(i)$ denote the sector that individual $i$ joins after graduating from her Ph.D. at time $t = 0$. Individuals can join one of two sectors in the economy: Industry ($j(i) = I$) or Academia ($j(i) = A$). I assume that I observe a unique generation of individuals so that $t$ can also be conceptualized as years of experience. I also assume for simplicity and because it does not impact the selection mechanism I want to highlight that individuals do not change sector during their career (Appendix C presents the full model allowing for transitions across sectors).

Each individual receives two endowments at birth: an individual-specific productivity parameter in Industry $\delta_i^I$ and an individual-specific productivity parameter in Academia $\delta_i^A$. These parameters capture how more/less productive individual $i$ is in Industry and Academia compared to the average. $\delta_i^I$ represents general skills that are common to both Industry and Academia, such as overall ability or creativity, as well as Industry-specific skills, such as social and communication skills, taste for teamwork, a greater concern for salary and access to resources, a stronger interest in downstream work and the ability to work on problems and find solutions which are relevant for the firm (Roach and Sauermann, 2010). Similarly, $\delta_i^A$ represents general skills that are common to both sectors, as well as Academia-specific skills, such as work independence, taste for publishing, importance given to peer recognition, ability to formulate research questions, self-criticism, as well as writing and teaching skills (see Ballesteros-Rodríguez et al. (2022) for a review of the literature).

Assuming that earnings proxy for productivity,[2] we can write the potential earnings of individual $i$ at time $t$ if she started her career in sector $j \in (I, A)$ as:

$$W_{it}^j \equiv \overline{W_t^j} + \delta_i^j \tag{1}$$

with $\overline{W_t^j}$ the average earnings at time $t$ of individuals who started their career in sector $j$ and $\delta_i^j$ individual's $i$ productivity parameter in sector $j$.

At the micro level, $W_{it}^j$ is the sum of pecuniary returns related to research output and pecuniary returns related to non-research output. Research output can be further decomposed into a public and a private component. Public research relates to research that is publicly disclosed through publications or patents, while private research relates to research that is not disclosed outside of the boundaries of the organization. Depending on individual $i$'s sector and job in sector $j$, earnings could be a function of (i) public research, private research, and other activities (e.g., if $i$ has a position in Academia or works in a firm that engages in public research), (ii) private research and

---

[2] The fact that other things enter the earnings component does not change the selection mechanism I want to highlight.

other activities (e.g., if $i$ has a research position in a firm that does not publicly disclose its results) or (iii) other activities only (e.g., if individual $i$ has a job that is not related to research, such as supply-chain manager). Calling $P_{it}^j$, $R_{it}^j$ and $O_{it}^j$ respectively the public research output, private research output and output related to other activities produced at time $t$ by individual $i$ who started her career in sector $j$, we can write $W_{it}^j$ as:

$$W_{it}^j = p^j P_{it}^j + r^j R_{it}^j + o^j O_{it}^j = p^j(\overline{P_t^j} + \theta_{p,i}^j) + r^j(\overline{R_t^j} + \theta_{r,i}^j) + o^j(\overline{O_t^j} + \theta_{o,i}^j) \tag{2}$$

with $p^j$, $r^j$ and $o^j$ the sector-specific prices that convert these outputs into dollars, $\overline{P_t^j}$, $\overline{R_t^j}$, $\overline{O_t^j}$ the mean outputs by activity type in each sector and $\theta_{p,i}^j$, $\theta_{r,i}^j$ and $\theta_{o,i}^j$ the individual-specific unobservables associated with each output type in each sector. By construction:

$$\overline{W_t^j} = p^j \overline{P_t^j} + r^j \overline{R_t^j} + o^j \overline{O_t^j} \tag{3}$$

and

$$\delta_i^j = p^j \theta_{p,i}^j + r^j \theta_{r,i}^j + o^j \theta_{o,i}^j \tag{4}$$

Define the utility $u_{it}^j$ of individual $i$ at time $t$ in sector $j$ as being equal to:

$$u_{it}^j = W_{it}^j + \varepsilon_i^j \tag{5}$$

with $\varepsilon_i^j$ individual $i$'s non-pecuniary returns in sector $j$. Individual $i$'s utility of *starting* her career in sector $j \in (I, A)$, denoted by $U_i^j$, is then equal to the sum of the discounted utilities over the career:[3]

$$U_i^j = \sum_{t=0}^T \rho^t u_{it}^j = \sum_{t=0}^T \rho^t \overline{W_t^j} + \frac{1 - \rho^{T+1}}{1 - \rho}(\delta_i^j + \varepsilon_i^j) \tag{6}$$

with $\rho^t$ the discount factor at time $t$.

Individual $i$ *starts* her career in Industry iff :

$$j(i) = I \iff U_i^I - U_i^A > 0$$
$$\iff \sum_{t=0}^T \rho^t (\overline{W_t^I} - \overline{W_t^A}) + \frac{1 - \rho^{T+1}}{1 - \rho}(\delta_i^I - \delta_i^A) + \frac{1 - \rho^{T+1}}{1 - \rho}(\varepsilon_i^I - \varepsilon_i^A) > 0 \tag{7}$$

In Appendix D, I show that allowing for heterogeneous effects implies that the difference in unobservables $\delta_i^I - \delta_i^A$ can be correlated with their levels $\delta_i^I$ and $\delta_i^A$.

Equation 7 shows that individuals select on private gains and non-pecuniary gains. In addition, it highlights the endogeneity issues that arise since the error term of the earnings and publications outcomes is likely to be correlated with the difference in unobservables that enters the selection equation. This is easily visible for the earnings outcome which is a function of $\delta_i^j$: if individuals who

---

[3]Working with expected utilities does not change the main selection mechanism I want to highlight.

tend to select into the private sector (their difference $\delta_i^I - \delta_i^A$ is high) also tend to be particularly skilled in Industry (their $\delta_i^I$ is high), the earnings we observe in the private sector will be biased upwards. Similarly for publications, if those who tend to stay in Academia (their difference in unobservables $\delta_i^I - \delta_i^A$ is low) also tend to be particular skilled in (public) research (high $\theta_{p,i}^A$), the publication output observed in Academia will be biased upwards. In general, the sense of the bias will depend on the correlation between $\delta_i^I - \delta_i^A$ and the error term of the outcome considered.

# 3 Data

In this section, I describe the data sources and define the key variables.

## 3.1 Data sources and sample construction

I use individual-level data from three restricted files: (i) the Survey of Earned Doctorates (SED), (ii) the Surveys of Doctorate Recipients (SDR) and (iii) the 2015 Survey of Doctorate Recipients Bibliometric Research Data (SDR15-WoS). These surveys are conducted by the National Center for Science and Engineering Statistics (NCSES) and cover individuals who received a Ph.D. from an American institution.[4]

The SED is an annual census of all individuals who graduated from a Ph.D. in a given academic year (e.g., the SED 2010 is sent to all individuals who earned their doctoral degree in 2010). The SED is therefore a *within-cohort* dataset which gives information about PhDs *at the time of graduation*. The response rate is about 90%. The SED contains demographic and background information including Ph.D. field of study and institution ; previous education ; place of birth, gender, race, parental education, citizenship status and the sector that students join or plan to join after their PhD.

The SDR is a biennal panel survey of doctorate holders which contains longitudinal data about individuals' employment at the time of survey such as sector of employment, principal job activity and earnings. For each survey, the sample consists of both individuals who were surveyed in previous SDR editions (as long as they are less than 76 years old) as well as new doctoral graduates who earned their Ph.D. since the last cycle. The SDR is therefore a *cross-sectional* dataset and gives information about individuals *during the career*. The response rate is about 65%. I use all the SDR surveys to which I have access to (years 2003, 2006, 2008, 2010, 2013, 2015, 2017, 2019 and 2021) to retrieve information about individuals' earnings. Note that individuals can be surveyed several times during their career and hence might appear in several SDR waves.

In the years 2022/2023, the NCSES launched a new Research Data Infrastructure with the goal of creating new linkages with external data. I take advantage of this effort by using a newly created dataset that matches each individual surveyed in the SDR 2015 to her research output from graduation until 2017, using information from Web of Science. This gives me individual-level

---

[4]These data are restricted-use and were accessed remotely from the National Opinion Research Center (NORC) data enclave.

information about (public) research production between graduation and 2017 such as publications, top publications and average number of citations in the first 2 or 5 years after publication. For each of these measures and for each individual, I have access to *aggregate* values up to 2017 as well as more granular information by bins of experience (e.g., publications produced during the PhD, publications produced during the first 5 years post-graduation and so on).[5]

The initial sample contains 57,811 individuals. I drop PhDs in Humanities, Education, Business and Health for whom I have little observations and focus on STEM fields (Computer and Mathematical Sciences, Life Sciences, Physical Sciences, Social Sciences and Engineering). I keep cohorts that graduated in or after 1970 because of several changes in the way the SDR answers are coded before that year. In order to assign individuals to a sector of employment when they start their career, I keep students with definite postgraduate commitments at graduation at the time of survey (70%) and exclude individuals who declare starting an internship, traineeship, clinical residency or military service. I exclude individuals who do not start their career in the United States (12.6%). I also exclude individuals for whom I do not have information about publications, gender, father education, mother education, race, PhD major and earnings. This leaves me with 23,907 observations. Because the earnings outcome is a flow and doctorate holders can be surveyed several times, the earnings regressions will be estimated on a sample of 90,325 observations, meaning that I observe on average 3.7 earnings values per individual.

## 3.2  Dependent Variables

*Earnings* - From these datasets, I primarily have access to information about earnings *during* the career. Earnings correspond to the total earned income before deductions in the previous year. This includes wages, bonuses, overtime, consulting fees and summertime teaching or research. I use the number reported by respondents across the several SDR surveys they have taken so that I may have several earnings observations per individual corresponding to different levels of experience. I transform all values in 2015 US dollars. The SED has information about salary *at graduation* for cohorts who graduated in or after 2008 but it is not my preferred variable given that it highly restricts the number of observations available. I winsorize observations below the $1^{st}$ percentile and above the $99^{th}$ percentile.

*Publications* - I calculate the number of publications published after Ph.D. graduation by subtracting the number of publications published before graduation from the total number of publications published by the individual until 2017.[6] I winsorize observations below the $1^{st}$ percentile and above the $99^{th}$ percentile.

---

[5] I have counts of these measures for the 5 years before graduation and for bins of 5 years after graduation, starting for bin of years 1 to 5 after graduation up to bin of years 46 to 50 after graduation.

[6] I only include publications that were articles, reviews or conference proceedings.

### 3.3 Defining the sector joined at graduation

I first assign each individual to three potential sectors: Industry (business for-profit), Academia (2-year and 4-year college or university and university-affiliated research institute) or Government (US or foreign government).[7] My preferred treatment variable is an indicator variable equal to 1 for individuals joining Industry and 0 for individuals joining any other sector, i.e., Academia or Government (in the rest of the paper, I will loosely refer to this treatment variable as joining Industry vs Academia).[8]

Note that about 50% of individuals in my sample start their career in a postdoc position. I choose to assign individuals doing a postdoc to the academic sector for three main reasons. First, while I do not have information about the sector where the postdoc is realized for individuals who graduated between 1969 and 2003[9] only 2% of individuals who graduated in or after 2004 did so in the private sector.[10] Second, while it is possible to consider postdoc as a continuation of PhD education rather than as a first employment, I cannot easily identify the sector where individuals work right after finishing their postdoctoral studies because (i) I do not know exactly when they finish their postdoc and (ii) while I observe sector of employment for all individuals in my sample after PhD graduation, I do not necessarily observe individuals again in the early-part of the career, especially for older cohorts.[11] Based on individuals that I observe 1 to 5 years after PhD graduation and identified as doing a postdoc, 80% of them are observed in the academic sector, providing reassurance in the way I assign individuals to a sector. Finally, assigning postdoc to the academic sector allows me to remain conservative and provides lower bound to the treatment effects on earnings and publications. Individuals who start in the private sector after their postdoctoral studies can be expected to earn more and publish less than in Academia. By assigning them to Academia, this provides a lower bound on the observed earnings premium and publications gap between Industry and Academia.

Among the 23,907 unique individuals in my sample, 17% are joining Industry.[12]

Note that my treatment corresponds to the sector where individuals *start* their career, but individuals may change sector between graduation and the time of survey. As a consequence, I might, for example, have individuals who started their career in Academia but for whom I observe earnings while they are in Industry (and vice-versa). Similarly, the stock of publications observed in 2017 might include publications produced while in Academia and others produced while in Industry. This is not an issue for my empirical analysis as I am interested in the causal impact of *starting* one's career in Industry vs Academia, which also incorporates the probability of changing sector. But if the reader were to be interested in the causal impact of *working* in Industry vs Academia,

---

[7]I follow the classification provided by the NCSES. I exclude individuals working in non-profit and U.S. preschool, elementary and middle school

[8]This is motivated by my informal conversations with doctorate holders, but my results are robust to excluding the Government category.

[9]The question started to be asked to individuals in 2004 onwards.

[10]82% of individuals do their postdoc in Academia, 11% in the government and 2% in Industry

[11]31% of my observations correspond to individuals observed less than 10 years after graduation.

[12]17% are joining Industry, 72% are joining Academia and 10% are joining Government.

my estimates are likely to represent lower bounds. Indeed, most of individuals in my sample who changed sector between graduation and survey moved from Academia to Industry. For these individuals, their earnings value is therefore attributed to the academic sector while they correspond to a position in Industry (which likely pays more), hence decreasing the earnings gap between Industry and Academia. Similarly, their stock of publication is attributed to Academia while part of their career was spent in Industry (where individuals tends to publish less). This decreases the publication difference between Industry and Academia.

## 3.4 PhD fields and majors

I refer to Ph.D. *major* as the most granular group I have access to for each individual (e.g., Biology, Genetics, Biochemistry, Astronomy, Computer Engineering etc.). This is what students choose to study. My dataset includes 57 distinct PhD majors. I refer to PhD *field* as a way to categorize these majors. I focus my study on 5 fields: Computer Sciences and Mathematics, Life Sciences, Physical Sciences, Social Sciences and Engineering.

# 4 Descriptive analysis

In this section, I further describe the sample and present the main characteristics of individuals joining Industry vs Academia at graduation.

## 4.1 Supply

Figure 1 shows the number of Ph.D. by cohort of graduation and Ph.D. field in my sample.

The number of individuals increases with cohort of graduation, reflecting the increase in the number of PhD over time. Overall, 73% of individuals in my sample are born in the United States. Figure A.1 shows an important increase in the share of international students over time, from about 10% in the 80s to 30% in 2013, with a peak at about 40% in 2006. Interestingly, this increase is primarily concentrated in the fields of Computer Sciences and Engineering, where the share of international students rose from about 30% in the 80s to 50% in 2013. Figure A.2 shows that the share of female increased linearly from about 25% in the 80s to about 45% in 2013. This increase is characteristic of all fields though the share of female seems to be relatively flat in the field of Computer Sciences and Mathematics.

## 4.2 Sorting

Figure 2 shows the sorting pattern between Industry and Academia. The blue (resp. red) bar shows the number of people who join the academic (resp. private) sector after graduation. The black lines translates this sorting into a share by showing the percentage of individuals joining the private sector at graduation, which fluctuates around 20% but remains relatively constant over time. Interestingly, the share of individuals joining the private sector at graduation peaks and then
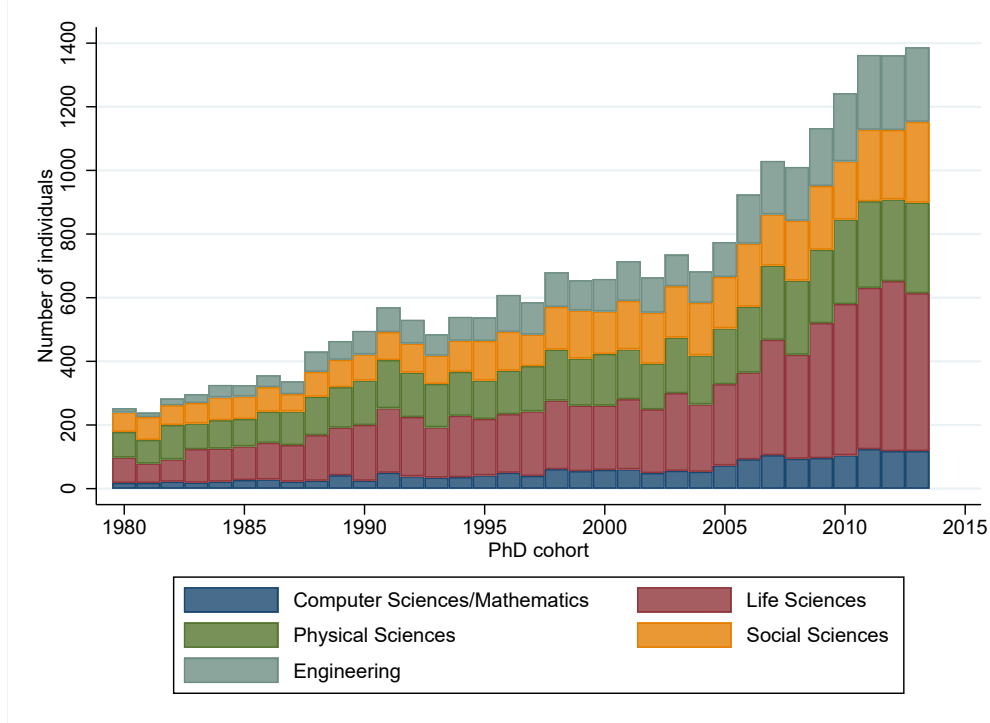
Figure 1: Number of individuals by PhD cohort of graduation

decreases around the major macro-economic shocks: in 1981/1982, 1990, 2000/2001 and 2007/2008. Figure 3 shows that there is variation across fields regarding the share of individuals joining the private sector, with Engineering being the most Industry-oriented field. Of note is the increase of the share of individuals in Computer Sciences/Mathematics joining the private sector starting around 2005. Importantly, most of the variation seems to be *within* field. I will leverage this idea to construct an instrumental variable. Figure A.3 shows that starting in the 90s, international students are more likely than Americans to join the private sector. Figure A.4 shows that female students are less likely than males to start their career in the private sector.

## 4.3   Summary statistics

Table 1 shows some summary statistics for my sample as a function of the sector joined at graduation. Individuals who *started* their career in Industry earn on average 47k more per year than individuals who started their career in Academia. They also published on average 11 fewer publications between graduation at the time of survey.[13] On average, individuals graduated in 2000, representing about 13 years of experience at the time of survey. 42% of individuals in Academia are

---

[13]While this difference might seem low compared to what we would expect from the difference in publications between Industry and Academia over the career, remember that my treatment variable is defined by individuals' sorting decision *at graduation*, not by the sector that employs them at the time of survey. This tends to make the publication difference between individuals who *started* in Industry and individuals who *started* in Academia lower than the publication difference we would expect between individuals *working* in Industry and individuals *working* in Academia.
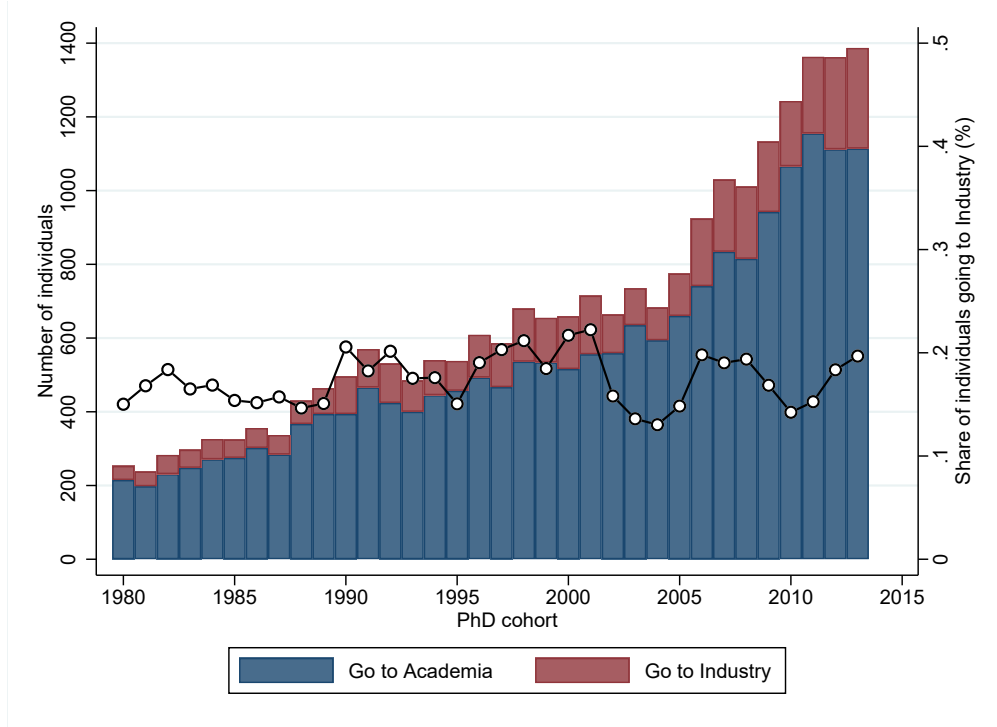
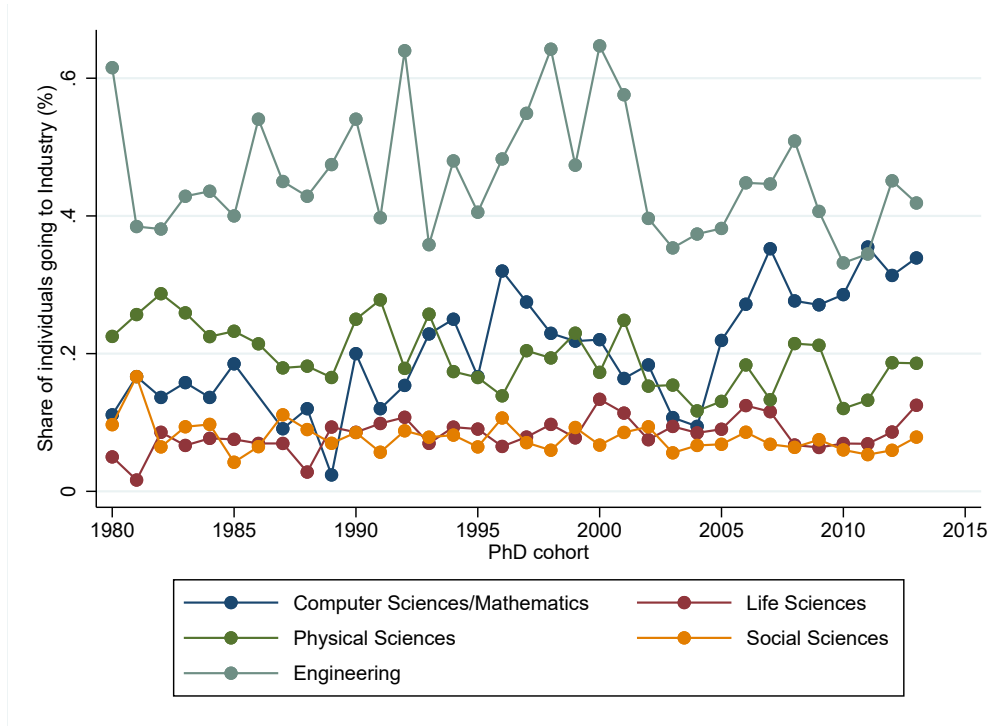Figure 2: Sector joined at graduation



Figure 3: Share of individuals joining the private sector by PhD field of study and PhD cohort of graduation

female in Academia, vs only 27% in Industry. In both sectors, a majority of individuals are white and this share is higher Academia. On average and in both sectors, about 50% (60%) of individuals have a mother (father) who has at least a bachelor degree. As highlighted before, more international students tend to join the private sector at graduation than Academia. There is no apparent significant differences between sectors regarding the number of publications before graduation.

Table 1: Summary Statistics

|  | Academia (1) | | Industry (2) | | (1)-(2) | |
| --- | --- | --- | --- | --- | --- | --- |
|  | Mean | S.d. | Mean | S.d. | Diff. | t |
| Earnings (2015 USD th.) | 101.5 | 71.3 | 148.3 | 100.5 | -46.8*** | -28.4*** |
| Publications (stock) | 19.3 | 29.1 | 8.4 | 19.3 | 10.9*** | 29.9*** |
| PhD cohort | 1,999.7 | 10.6 | 2,000.3 | 10.0 | -0.6*** | -3.4*** |
| White (0/1) | 0.69 | 0.46 | 0.61 | 0.49 | 0.08*** | 9.5*** |
| Female (0/1) | 0.42 | 0.49 | 0.27 | 0.44 | 0.15*** | 19.9*** |
| Experience (years) | 13.92 | 10.62 | 13.12 | 10.05 | 0.79*** | 4.6*** |
| Mother at least undergrad. (0/1) | 0.48 | 0.50 | 0.49 | 0.50 | -0.01 | -1.3 |
| Father at least undergrad. (0/1) | 0.61 | 0.49 | 0.63 | 0.48 | -0.02** | -2.7** |
| American (0/1) | 0.75 | 0.43 | 0.64 | 0.48 | 0.10*** | 12.6*** |
| Publications during PhD | 2.1 | 5.2 | 2.2 | 3.9 | -0.1 | -1.2 |
| Observations | 19,787 | | 4,20 | | 23,907 | |

Finally, Figures 4 and 5 show the two main outcomes of interest as a function of career experience. The difference in publication stock appears to be increasing over time. The earnings curve in Industry has a higher intercept but is less steep, so that earnings in both sectors look roughly similar in the raw data towards the end of the career. This could be due to differences in skill obsolescence (Deming and Noray, 2018), differences in ability and preferences (Roach and Sauermann, 2010; Stern, 2004) and differences in physical capital investments and complementarities between basic an applied scientists in Industry vs Academia (Agarwal and Ohyama, 2013).

While the raw data suggests important differences in earnings and publications between individuals who start their career in Industry and those who start their career in Academia, these numbers do not reflect causal estimates because they are likely to be biased by *selection*, as individuals can be expected to choose their sector at graduation based on their private returns. In what follows, I employ a step-wise approach to estimate plausibly causal individual-specific treatment effects.

# 5    Estimation Strategy and Results

In this section, I present the several estimation strategies I use to estimate the causal impact of joining Industry vs Academia at graduation on earnings and publications. I am primarily interested in estimating individual-specific treatment effects in order to characterize the pools of talent that
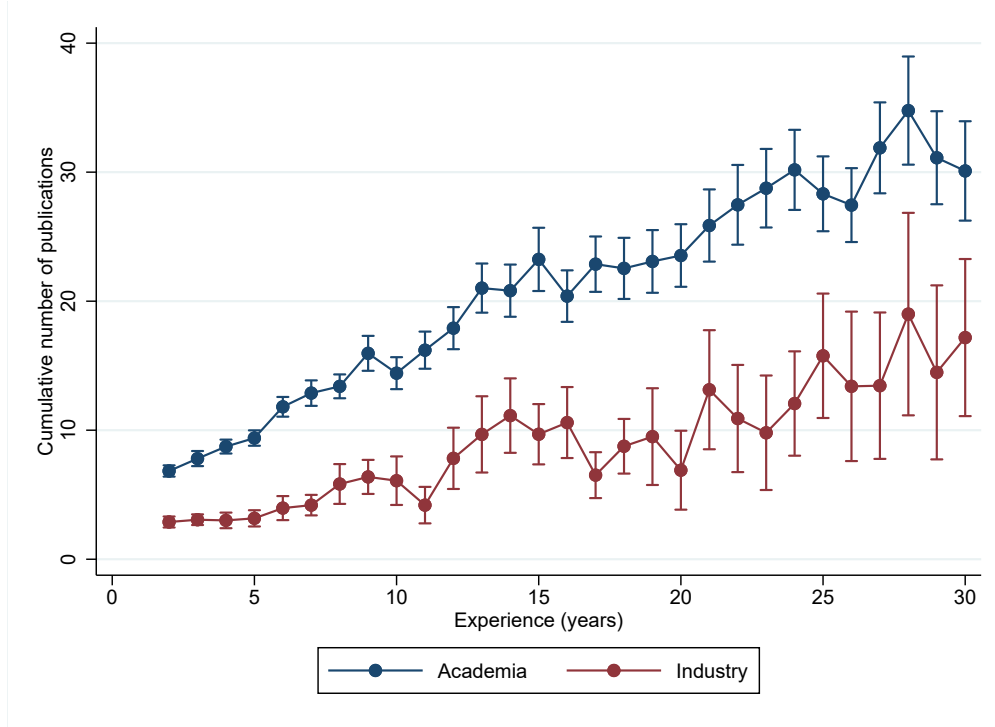
Figure 4: Mean publications (stock) as a function of experience and sector joined at graduation
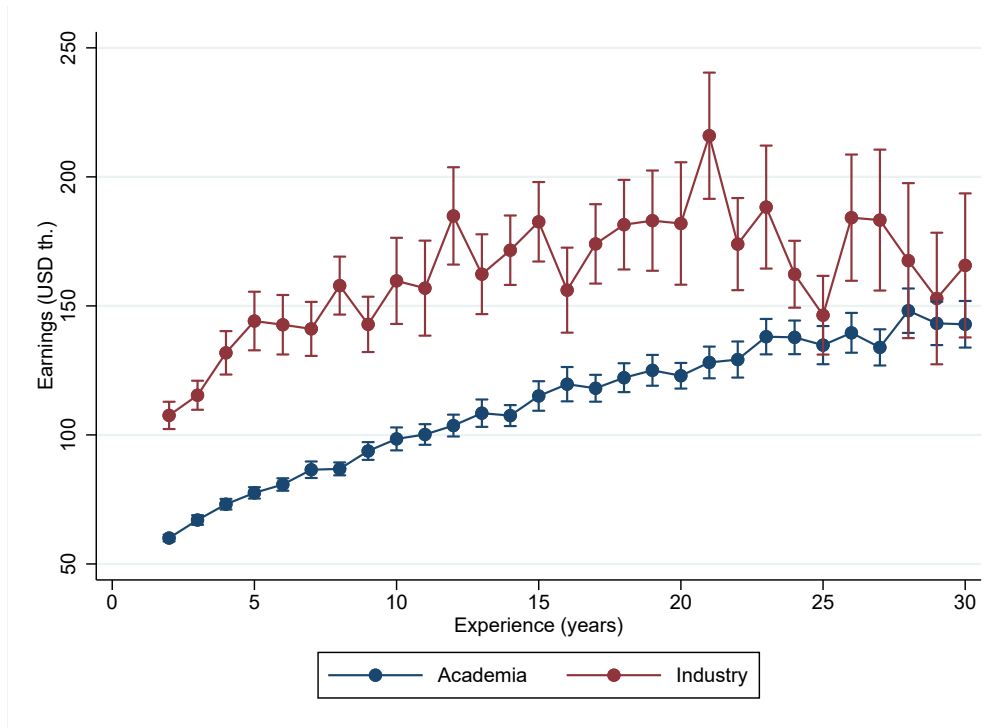


Figure 5: Mean earnings as a function of experience and sector joined at graduation

13

start in Industry and in Academia. I take a step-wise approach to build this case. I first use an OLS specification that relies on variation *within PhD major × doctoral institution* (hereafter 'department'). In order to account for selection at the individual level, I then use an instrumental variable that allows me to estimate a Local Average Treatment Effect (LATE). I then explore heterogeneity on observable and unobservable characteristics by estimating Marginal Treatment Effects (MTE).

## 5.1 Within department variation

First, I follow recent work that controls for ability by using the quality of PhD students' university department, proxied empirically by using department ranking as a linear control (Agarwal and Ohyama, 2013; Cohen et al., 2020; Roach and Sauermann, 2023, 2024). I extend this idea by using PhD major × doctoral institution fixed effects in a cross-section regression. Note that this is a more stringent version than linearly controlling for PhDs' department ranking, since I only use variation across individuals *within* a department. Under the identifying assumption that the difference in unobserved productivities between Industry and Academia does not vary within a department, using variation in sorting across individuals pertaining to the same department recovers the ATE of joining Industry vs Academia at graduation. My main specification is the following:

$$Y_i = \alpha_0 + \alpha_1 \text{Industry}_i + \delta_{mu} + \alpha_2 \mathbf{X_i} + \varepsilon_i \tag{8}$$

with $m$ indexing PhD major, $u$ indexing university. $\delta_{mu}$ represents PhD major × university fixed effects and $\mathbf{X_i}$ a vector of control variables which includes an indicator equal to 1 if the individual reports being a female, an indicator equal to 1 if the individual reports being white, an indicator equal to 1 for American (vs International) students, an indicator equal to 1 if the individual's father has at least a bachelor degree, an indicator equal to 1 if the individual's mother has at least a bachelor degree. I control for experience using a linear and quadratic terms for experience, but results are similar using variation within PhD major × doctoral institution × PhD cohort. In one specification, I also proxy for sector-specific preferences that could be correlated with the outcomes by using information about individuals' preferences for 9 job attributes.[14] More precisely, I control for the importance given to job independence, salary, security and impact on society by including 4 indicators equal to 1 if the attribute is judged 'very important'.[15] Note however that because these preferences are observed during employment, they might have been influenced by the treatment variable itself.

*Publications* - Results for the (log) publication outcome are presented in Table 2, with robust

---

[14]Individuals are asked about the importance of 9 job attributes: (1) job's opportunities for advancement (2) job's benefits (3) job's intellectual challenge (4) job's degree of independence (5) job's location (6) job's level of responsibility (7) job's security (8) job's salary and (9) job's contribution to society. For each attribute, I create an indicator variable equal to 1 if individuals declare this attribute as being very important and 0 otherwise. Figure [awaiting export] shows the average value of each indicator by sector joined at graduation.

[15]The other attributes are highly correlated with these 4 ones.

standard errors clustered at the doctoral institution level. Column (1) only includes the main demographics controls (gender, race, parental education and nationality). Column (2) adds PhD major fixed effects in order to control for field-specific differences in publication norms. Column (3) adds doctoral institution fixed effects. Column (4) adds doctoral institution × PhD major fixed effects so that the effect is identified by using granular variation *within* universities' departments. This is my preferred specification. Column (5) adds controls for sector-specific preferences. Overall, the coefficient on the treatment variable *Industry* remains stable and equal to about -0.9, implying that going to Industry vs Academia upon graduation decreases the total number of publications by about 60%, representing a loss of about 11 publications.[16] Table B.1 reiterates the analysis on the sub-sample of individuals with strictly positive publications and finds similar results. Table B.2 uses citations-weighted publications as the main outcome and finds more negative coefficients, implying that publications produced in Academia also gather more citations.

Table 2: Within department specification - OLS Publications (stock)

| | Log(1+Publications) | | | | |
| | (1) | (2) | (3) | (4) | (5) |
| --- | --- | --- | --- | --- | --- |
| Industry | -0.906*** | -0.935*** | -0.930*** | -0.944*** | -0.899*** |
| | (0.0194) | (0.0240) | (0.0241) | (0.0256) | (0.0283) |
| Demographics | Yes | Yes | Yes | Yes | Yes |
| Experience | Yes | Yes | Yes | Yes | Yes |
| PhD major FE | | Yes | Yes | Yes | Yes |
| Doct. Inst. FE | | | Yes | Yes | Yes |
| PhD major x Doct. Inst. FE | | | | Yes | Yes |
| Taste | | | | | Yes |
| Observations | 23,907 | 23,907 | 23,878 | 22,210 | 19,310 |
| R-sq | 0.163 | 0.198 | 0.224 | 0.330 | 0.357 |

*Notes*: This table reports the OLS estimates for publications using within department (defined as PhD major × doctoral institution) variation. Experience controls include a linear and a quadratic terms. Standard errors (in parentheses) are clustered at the doctoral institution level.

*Earnings* - Table 3 Panel A shows results for the earnings outcome, with the same controls as before. I first consider for each individual the earnings observation closest to graduation. Once all controls are included, the estimate is about 0.28, implying that the earnings premium is equal to about 32%, representing about USD 35k.[17] Table 3 Panel B shows the result using all the earnings observations linked to an individual. Once all controls are included, the coefficient of 0.26 implies an earnings premium of about 30%, representing a premium of USD 36k.[18]

Note that the coefficients do not vary much once the PhD field fixed-effects are included. This

---

[16]The mean of (winsorized) publications is 17.9.
[17]The mean of (winsorized) earnings for the first observation is USD 109.6.
[18]The mean of (winsorized) earnings is USD 122.8.

Table 3: Within department specification - OLS Earnings

| | Log(1+Earnings) | | | | |
|---|---|---|---|---|---|
| | (1) | (2) | (3) | (4) | (5) |
| *Panel A: First observation* | | | | | |
| Industry | 0.357*** | 0.287*** | 0.286*** | 0.284*** | 0.280*** |
| | (0.0112) | (0.0126) | (0.0128) | (0.0139) | (0.0152) |
| Observations | 23,907 | 23,907 | 23,878 | 22,210 | 19,310 |
| R-sq | 0.162 | 0.190 | 0.210 | 0.315 | 0.335 |
| | | | | | |
| *Panel B: All observations* | | | | | |
| Industry | 0.346*** | 0.274*** | 0.272*** | 0.262*** | 0.255*** |
| | (0.0107) | (0.0123) | (0.0123) | (0.0136) | (0.0136) |
| Observations | 90,325 | 90,325 | 90,323 | 90,196 | 83,196 |
| R-sq | 0.123 | 0.153 | 0.173 | 0.284 | 0.308 |
| | | | | | |
| Demographics | Yes | Yes | Yes | Yes | Yes |
| Experience | Yes | Yes | Yes | Yes | Yes |
| PhD major FE | | Yes | Yes | Yes | Yes |
| Doct. Inst. FE | | | Yes | Yes | Yes |
| PhD major x Doct. Inst. FE | | | | Yes | Yes |
| Taste | | | | | Yes |

*Notes*: This table reports the OLS estimates for earnings. Panel A uses the observation closest to PhD graduation and Panel B uses all observations linked to an individual. I use within department (defined as PhD major × doctoral institution) variation. Experience controls include a linear and a quadratic terms. Standard errors (in parentheses) are clustered at the doctoral institution level.

could imply limited selection on ability from the supply-side (at least within the context of this model) and echoes findings from Roach and Sauermann (2024).

Despite this list of controls, we might still be concerned that some unobservable individual characteristics that are correlated with both sector choice and outcomes vary across individuals, even within a department. Hence, I also implement an instrumental variable (IV) strategy to estimate the causal impact of joining Industry vs Academia on earnings and publications. This allows me to estimate a local average treatment effect (LATE), i.e., the causal effect of joining Industry vs Academia for individuals on the margin of going to Industry.

## 5.2 Instrumental Variable

As depicted in Figure 3, there is considerable variation across cohorts within each PhD field in the proportion of individuals opting for the private sector upon graduation. I leverage this variation at the more granular level of the PhD *major* and construct an instrumental variable based on the fluctuating *relative* demand of firms vs universities for PhD students in the same major over time. For example, the demand for PhDs in Aerospace Engineering from firms might surge in 2009 due to specific industry needs, while shifts in undergraduate enrollments could affect faculty demands in Biomedical Engineering in different years. Said differently, the demand for doctorate-level skills might be stronger in some years on the private side vs Academia and vice-versa. Within each PhD major, these fluctuations create year-to-year variations in the likelihood of observing individuals in Industry vs Academia upon graduation. This is similar in idea to Oyer (2006), who uses macroeconomic shocks for PhDs in Economics. I extend this approach to encompass all PhD majors within my sample. Empirically, one can think of the instrument as a vector of PhD major $\times$ PhD cohort fixed effects, controlling for PhD major fixed effects.[19] In practice, I follow the literature and create a leave-one out continuous equivalent of the interaction term to reduce the risk of bias.[20] For each individual, the instrument $Z_i$ is therefore equal to:

$$Z_i = \frac{\sum_{i'} Ind_{i'} \times \mathbf{1}\{m(i') = m(i)\} \times \mathbf{1}\{c(i') = c(i)\} - Ind_i}{\sum_{i'} \mathbf{1}\{m(i') = m(i)\} \times \mathbf{1}\{c(i') = c(i)\} - 1} \tag{9}$$

with $i'$ indexing individuals, $Ind_i$ an indicator equal to 1 if $i$ starts her career in Industry, $m(i)$ representing individual $i$'s PhD major and $c(i)$ representing individual $i$'s PhD cohort of graduation.

To be valid, the instrument: (i) should not impact earnings/publications directly (i.e., other than through its impact on sector joined at graduation) and (ii) should not be correlated with individuals unobservable productivities in Industry vs Academia $\delta_i^I - \delta_i^A$. I discuss both successively. I primarily focus on earnings, but the argumentation remains the same for the publications outcome.

Regarding (i), one might worry about general macroeconomic conditions that would be correlated with firms' vs universities' demand. For instance, if firms' vs universities' demand for PhD in a specific major happens to be high at a time where the economy is also at its best, earnings in

---

[19]I only keep PhD major $\times$ PhD cohort fixed effects for which I have at least 10 observations.
[20](Dobbie et al., 2018; Sampat and Williams, 2019)

the private sector might appear high even though individuals would have enjoyed high earnings in Academia too. To mitigate this concern, I include a control for the unemployment rate at the time of PhD graduation.

In addition, (ii) implies that individuals who are more/less productive in Industry vs Academia do not time their decision of when to go on the market based on macroeconomic conditions. Importantly, note that this relates to the *relative* sector-specific abilities and not to general ability. Said differently, cohorts are allowed to be of different overall quality over time, but they are not allowed to be of higher or lower Industry vs Academia quality. A weaker version of this assumption is that firms/universities do not respond to changes in sector-specific quality. I include doctoral institution fixed effects and an indicator for having published prior to graduation as the two best control variables I have to strengthen the validity of this assumption. I then run two tests to assess its plausibility. First, I examine whether the instrument is correlated with PhD duration in order to assess whether individuals strategically time their graduation based on firms' vs universities' demand. Second, I examine whether firms' vs universities' demand at the time of PhD entry correlate with sorting, which would indicate that macroeconomic conditions are predictable when individuals decide to start a PhD. Table 4 column 1 shows that the instrument is not correlated with PhD duration. Table 4 column 2 shows that firms' vs universities' demand for PhD at the start of PhD entry does not predict sorting, providing some reassurance regarding the validity of the instrument. *First Stage* - The first-stage estimates are presented in Table 5, which shows variants of the following regression:

$$Ind_i = \gamma_0 + \gamma_1 Z_i + \gamma_2 \mathbf{X_i} + \nu_i \tag{10}$$

$Ind_i$ is an indicator equal to 1 if the individual joins Industry at graduation and 0 otherwise. $Z_i$ is the instrument for individual $i$. Column (1) includes PhD major fixed effects so that I rely on variation *across years within PhD major*, as well as controls for experience (with a linear and quadratic terms) and demographic characteristics (an indicator for being female, an indicator for being white, an indicator for father's education being at least an undergraduate degree, an indicator equal to 1 for mother's education being at least an undergraduate degree, an indicator equal to 1 for being American vs International). Column (2) adds doctoral institution fixed effects. Column (3) adds an indicator equal to 1 if the individual has published at least once before graduating as well as the unemployment rate at the time of graduation. Standard errors are clustered at the PhD major and PhD cohort level to account for intra-group correlation.[21]

All columns show a first-stage with a F-statistic above 30. As expected, graduating in a year where firms' vs universities' demand for PhDs is higher (as reflected in higher values of the instrument) increases individuals probability to enter the private sector.

*2SLS Publications* - Table 6 presents the 2SLS estimates for the publication outcome. Using the pre-

---

[21] Results clustering at the PhD major level only give similar results.

Table 4: Instrument Validity

|  | PhD Duration (1) | Industry (2) |
|---|---|---|
| $Z_i$ | -0.00108 | |
|  | (0.290) | |
| $Z_i$ at entry | | 0.00787 |
|  | | (0.0364) |
| Demographics | Yes | Yes |
| Experience | Yes | Yes |
| PhD major FE | Yes | Yes |
| Doct. inst. FE | Yes | Yes |
| Macro cond. | Yes | Yes |
| Observations | 23,671 | 21,437 |
| R-sq | 0.114 | 0.185 |

*Notes*: This table reports two falsification tests. Column (1) examines whether the instrument predicts PhD duration, which is defined either as the difference between PhD graduation year and master graduation year if the individual reports having a master or as the difference between PhD graduation year and undergraduate graduation year if the individual does not report having a master. Column (2) examines whether firms' vs universities' demand at the time of PhD entry predict sorting. In both specifications, I control for general macroeconomic conditions by including a measure of national unemployment rate and I also include an indicator equal to 1 if the individual has published during the PhD. Standard errors (in parentheses) are two-way clustered at the PhD major and PhD cohort level.

Table 5: First-stage estimates

|  | Industry | | |
|  | (1) | (2) | (3) |
| --- | --- | --- | --- |
| $Z_i$ | 0.212*** | 0.206*** | 0.202*** |
|  | (0.0361) | (0.0367) | (0.0362) |
| Demographics | Yes | Yes | Yes |
| Experience | Yes | Yes | Yes |
| PhD major FE | Yes | Yes | Yes |
| Doct. inst. FE |  | Yes | Yes |
| Macro cond. |  |  | Yes |
| F-stat | 34.35 | 31.57 | 31.25 |
| Observations | 23,907 | 23,907 | 23,907 |
| R-sq | 0.161 | 0.181 | 0.181 |

*Notes*: This table reports the first-stage estimates. Column (1) includes PhD major fixed effects as well as experience and demographics controls. Column (2) adds doctoral institution fixed effects. Column (3) adds an indicator equal to 1 if the individual has published during the PhD and a control for the unemployment rate. Standard errors (in parentheses) are two-way clustered at the PhD major and PhD cohort level.

ferred specification in column (3), the LATE estimates imply that the marginal individual starting her career in Industry generates 77% less publications than if she would have stayed in Academia, representing about 14 publications less.[22] Table B.4 uses citations-weighted publications as an outcome and finds more negative magnitudes, implying that publications in Academia also receive more citations.

*2SLS Earnings* - Table 7 repeats this exercise for the earnings outcome. Columns (1) and (3) focus on the first earnings observation I observe for each individual. Columns (4) to (6) consider all earnings observations. Overall, I do not find any statistically significant result. Comparing the OLS and 2SLS estimates suggest that individuals who select into Industry (i) select on earnings gains and (ii) experience a lower decrease in publications than under 'random' assignment.

While the 2SLS estimates get us closer to a causal effect, they also present several limitations. In particular, they identify the effect for a very specific group of individuals - the compliers - and do not allow to study heterogeneity in treatment effects as well as the pattern of selection. I now proceed to the analysis of heterogeneity along observable and unobservable dimensions.

---

[22]Table B.3 in Appendix shows the results when restricting the sample to individuals with at least one publication. Results are unchanged.

Table 6: 2SLS - Publications

|  | Log(1+Publications) | | |
|  | (1) | (2) | (3) |
| --- | --- | --- | --- |
| Industry | -1.351*** | -1.388*** | -1.461*** |
|  | (0.448) | (0.463) | (0.486) |
| Demographics | Yes | Yes | Yes |
| Experience | Yes | Yes | Yes |
| PhD major FE | Yes | Yes | Yes |
| Doct. inst. FE |  | Yes | Yes |
| Macro cond. |  |  | Yes |
| F-stat | 34.35 | 31.57 | 31.25 |
| Observations | 23,907 | 23,907 | 23,907 |
| R-sq | 0.184 | 0.209 | 0.214 |

*Notes*: This table reports the 2SLS results for publications. Standard errors (in parentheses) are two-way clustered at the PhD major and PhD cohort level.

Table 7: 2SLS - Earnings

|  | Log(1+Earnings) | | | | | |
|  | First observation | | | All observations | | |
|  | (1) | (2) | (3) | (4) | (5) | (6) |
| --- | --- | --- | --- | --- | --- | --- |
| Industry | 0.197 | 0.194 | 0.281 | 0.0935 | 0.114 | 0.227 |
|  | (0.257) | (0.260) | (0.267) | (0.250) | (0.249) | (0.256) |
| F-stat | 34.90 | 32.05 | 31.31 | 29.25 | 27.63 | 26.97 |
| Observations | 23,907 | 23,907 | 23,907 | 90,325 | 90,325 | 90,325 |
| R-sq | 0.188 | 0.209 | 0.212 | 0.146 | 0.168 | 0.174 |
| Demographics | Yes | Yes | Yes | Yes | Yes | Yes |
| Experience | Yes | Yes | Yes | Yes | Yes | Yes |
| PhD major FE | Yes | Yes | Yes | Yes | Yes | Yes |
| Doct. Inst. FE | No | Yes | Yes | Yes | Yes | Yes |
| Macro cond. | No | No | Yes | No | No | Yes |

*Notes*: This table reports the 2SLS results for earnings using the observation closest to PhD graduation (columns 1 to 3) and all observations (columns 4 to 6). Standard errors (in parentheses) are two-way clustered at the PhD major and PhD cohort level.

## 5.3 Heterogeneity on observable characteristics

I focus on two main margins of observable heterogeneity: gender (female vs non-female) and nationality (American vs non-American) using my instrument. Columns (1) and (2) of Table 8 show results for publications. Columns (3) and (4) examine earnings. Columns (1) and (3) focus on heterogeneity across genders. Results show that female students joining the private sector tend to publish more and earn less than male students, though this last result is not statistically significant. Columns (2) and (4) shed light on heterogeneity across nationalities. Results show that Americans publish less and earn less than International ones. One reason for the higher earnings premium received by International students could be that those individuals who select into leaving their country are of higher ability than Americans on average. In order to shed more light on differences between American and International students, I restrict the sample to individuals joining the private sector at graduation and run a regression of the form:

$$Y_i = \gamma_0 + \gamma_1 Z_i + \gamma_2 Z_i \times \text{American} + \gamma_3 \mathbf{X_i} + \nu_i \tag{11}$$

where the coefficient of interest is $\gamma_2$, which shows the differential effect of joining Industry on the outcome $Y_i$ for Americans vs Internationals. I explore 3 different outcomes (i) occupations (ii) firms' size (iii) mobility and (iv) location. I find no difference with regards to occupation, with Americans and Internationals being equally likely to be observed in a R&D or managerial position. I find that Internationals are more mobile, defined as working on a different state than the state of doctoral institution, more likely to work in California and more likely to be observed in large firms with more than 25k employees.

## 5.4 Heterogeneity on observables and unobservables

Equation 7 from the initial model suggests that individuals select on gains, i.e., based on the difference $\delta_i^I - \delta_i^A$ and $\varepsilon_i^I - \varepsilon_i^A$. However, the LATE estimate gives a limited overview of selection and hence of the heterogeneity in unobservables it relates to. Indeed, the LATE averages the treatment effect for a specific part of the population (the compliers) and does not allow us to examine how treatment effect varies for individuals with latent propensities of entering Industry vs Academia (Andresen, 2018). I therefore proceed to the estimatation of Marginal Treatment Effects (MTE) (Björklund and Moffitt, 1987; Heckman and Vytlacil, 1999, 2007). MTE relate treatment effect heterogeneity to observed (i.e., related to observable characteristics) and unobserved (i.e., related to $\delta_i^I - \delta_i^A$) propensities to join the private sector. The MTE represents the treatment effect for individuals with a specific propensity of entering treatment, and we can recover the LATE (and other key estimands such as the ATE or the ATT) by computing weighted-averages of the MTEs. The MTE can be modelled as part of the generalized Roy Model that I present in Appendix E. I give here the main intuition.

*Model* - Each individual is characterized by observables $X$ and $Z$ (which represents the instrument)

Table 8: 2SLS Heterogeneity

| | Log(1+Publications) | | Log(1+Earnings) | |
|---|---|---|---|---|
| | (1) | (2) | (3) | (4) |
| Industry | -1.542*** | -1.169*** | 0.247 | 0.321 |
| | (0.465) | (0.448) | (0.244) | (0.225) |
| Industry × Female | 0.279* | | -0.0752 | |
| | (0.162) | | (0.0863) | |
| Industry × American | | -0.425*** | | -0.137* |
| | | (0.134) | | (0.0721) |
| Demographics | Yes | Yes | Yes | Yes |
| Experience | Yes | Yes | No | No |
| PhD major FE | Yes | Yes | Yes | Yes |
| Doct. Inst. FE | Yes | Yes | Yes | Yes |
| Macro cond. | Yes | Yes | Yes | Yes |
| Observations | 23,907 | 23,907 | 90,325 | 90,325 |
| R-sq | 0.215 | 0.208 | 0.174 | 0.176 |

*Notes*: This table reports the 2SLS results looking at heterogeneity across genders (columns (1) and (3)) and nationality (columns (2) and (4)). Standard errors (in parentheses) are two-way clustered at the PhD major and PhD cohort level.

and unobserved characteristic $U_D$ that is usually interpreted as the unobserved resistance or distaste for treatment.[23] Importantly for the interpretation, individuals with lower values of $U_D$ are more likely to join the private sector for unobservable reasons while individuals with higher values of $U_D$ are less likely to enter the private sector (i.e., more likely to enter Academia) for unobservable reasons. Individuals join Industry if:

$$Ind_i = \mathbf{1}\{P(Z_i, X_i) > U_{D_i}\}$$

with $P(Z, X)$ the propensity score. In this setup, the *marginal treatment effect* for individuals with observables $X = x$ and disutility for treatment $U_D = u$ is:

$$\text{MTE}(X = x, U_D = u) = \mathbb{E}[Y^I - Y^A | X = x, U_D = u] \tag{12}$$

with $Y$ the outcome of interest (i.e. earnings and publications). This is the ATE for groups of individuals who have the same latent propensity to be treated. The MTE formula makes it clear that treatment effect heterogeneity may result from both observed and unobserved characteristics. For instance, for individuals with observables $X$ that make them more likely to enter the private sector, I can test whether they have smaller or higher treatment effects than others. Similarly, considering individuals with lower $U_D$ that make them more likely to enter the private sector, I can

---

[23]Given Equation 7, $U_D$ is a function of $\delta_i^I - \delta_i^A$ and $\varepsilon_i^I - \varepsilon_i^A$.

test whether they have smaller or higher treatment effects than others. A common way to present the results is to plot the MTE as a function of $U_D$ for average values of $X$. If the slope of the curve is flat, then $\mathbb{E}[Y_1 - Y_0|X = x, U_D = u]$ is constant so that there is no heterogeneity in treatment effect across individuals with different propensities for joining the private sector vs Academia (or equivalently for individuals with different distaste for joining the private sector vs Academia). If the slope is upward or downward sloping, there is heterogeneity in treatment effect across individuals with different propensities for joining the private sector vs Academia.

The assumptions needed for the calculation of the MTE are the same as for the IV framework, with in particular a valid exclusion restriction which can be written as $(U^A, U^I, V) \perp\!\!\!\perp Z|X$ where $\perp\!\!\!\perp$ denotes conditional independence and $(U^A, U^I)$ are the error terms of the outcomes.[24] In practice, provided my previous instrument is valid, I can use it to estimate the MTE. In practice, it is possible to estimate the MTE with no further assumption. However, this requires full support of the propensity score in both the treated and untreated states for all values of $X$ which is rarely feasible (Andresen, 2018). Following the literature, I further assume additive separability between the observed and unobserved components. This allows me to write the MTE as:

$$\text{MTE}(X = x, U_D = u) = x(\beta_1 - \beta_0) + \mathbb{E}[U_1 - U_0|U_D = u] \tag{13}$$

It also implies that treatment affect heterogeneity coming from observed characteristics affects the intercept of the MTE curve (plotted as a function of $U_D$) but not its slope (Cornelissen et al., 2016). While this MTE framework is more restrictive than the 2SLS one, it allows me to bring more nuance into treatment heterogeneity and the pattern of selection.[25]

*Results* - Figure 6 shows histograms of the propensity score estimates (i.e., predicted probabilities to go to Industry vs Academia), for individuals who started their career in Industry and those who started their career in Academia. The common support ranges between 0 and almost 1. In what follows, I trim the thinnest tails of support and keep observations between the two dashed red lines.

Figures 7a and 7b show the MTE curves for publications and earnings respectively, estimated at the mean values of the covariates and with bootstrapped standard errors.[26] These curves relate the unobserved components of the treatment effect on earnings (i.e., $\delta_i^I - \delta_i^A$) and publications (i.e., $\theta_{p,i}^I - \theta_{p,i}^A$) and the unobserved component of treatment choice $U_D$, i.e., the *resistance* (or disutility) to join the private sector. The left figure shows a downward sloping curve: at the $10^{th}$ percentile of unobserved resistance (i.e., individuals who are very likely to go to the private sector), the decrease in publications associated with joining the private sector vs Academia is approximately equal to 60%. For individuals at the $70^{th}$ percentile of unobserved resistance (i.e., individuals that

---

[24]Omitting individual subscript, $U^A = \delta^A$ and $U^I = \delta^I$ when the outcome is earnings. When the outcome is publications, $U^A = \theta_p^A$ and $U^I = \theta_p^I$.

[25]Note that the separability assumption remains much less restrictive than a joint normal distribution of $(U^A, U^I, V)$ as assumed by traditional selection models (Andresen, 2018).

[26]I use the separate method with a polynomial of degree 1. Figures A.7a and A.7b show the curves estimated with a semi-parametric approach.
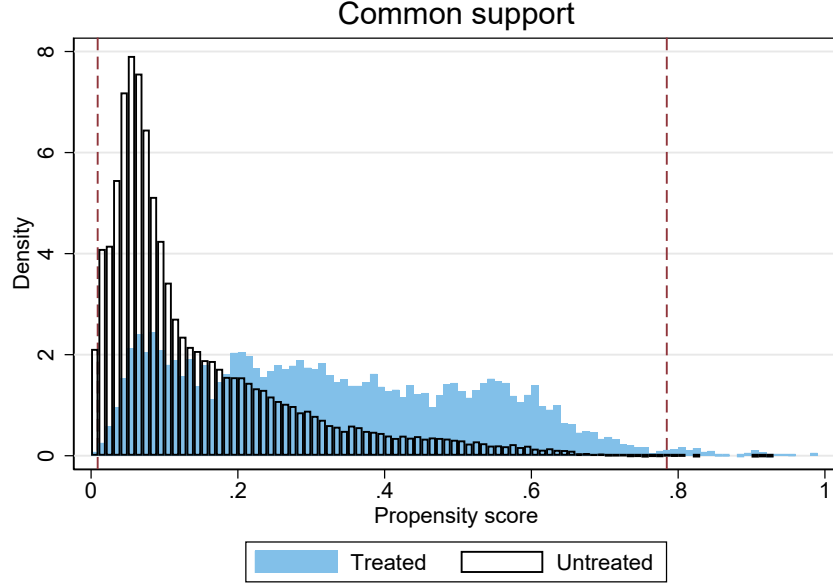
Figure 6: Common support

are less likely to go to the private sector), this decrease is approximately equal to 85%.[27] The earnings curve is also downward sloping: earnings gains are approximately equal to 65% at the $10^{th}$ percentile of unobserved resistance and approximately equal to 11% at the $70^{th}$ percentile of unobserved resistance.[28] Overall, this shows evidence of positive selection: those with a higher propensity to go to Industry have the most to earn, but they are not the ones who would experience the highest decrease in publications.



(a) Publications

(b) Earnings

Figure 7: MTE Curves

---

[27]The test of heterogeneity on observable characteristics is significant at the 1% level. The test of heterogeneity on unobservable characteristics is not significant at the 10% level.

[28]The tests of heterogeneity on observable and unobservable characteristics are both significant at the 10% level.

# 6 Competition for talent

## 6.1 Individual-specific treatment effects

An advantage of the MTE estimation approach is that the estimates can be combined with the data (i.e., individuals' observables and realized sector) to derive the expected response to treatment $\mathbf{E}[Y^I - Y^A|X, D, p]$ for any observation in the data. Figure 8 shows the histogram of the expected treatment effects for the publications (left panel) and earnings (right panel) outcomes, separately for individuals that I observe in Industry (green) vs those that I observe in Academia (red). Those who start their career in Industry experience higher earnings gains of joining the private sector vs Academia and a lower decrease in publication output, compared to individuals who start their career in Academia.



(a) Log(Publications)  (b) Log(Earnings)

Figure 8: Expected treatment effects, by sector joined

Before discussing more specifically the drivers of these results, I examine heterogeneity in treatment effects by gender (male vs female), nationality (American vs International) and field. Figures A.8b, A.9b and A.10 present the results for the earnings outcome. I find little differences in the distribution of treatment effects between males and females. The distribution of treatment effects for Internationals is shifted to the right, meaning that they experience higher earnings gains of starting in Industry vs Academia than Americans. The distribution of treatment effects appears relatively similar by field, though it is shifted to the left for Social Sciences and slightly shifted to the right for Life Sciences.

Similarly, Figures A.8a, A.9a and A.11 examine heterogeneity in treatment effects for the publications outcome. Results show that the distribution of treatment effects for males is shifted to the left, meaning that they experience a higher decrease in publications of starting in Industry vs Academia than females. There appears to be little difference across nationality. The distribution of treatment effects appears relatively similar by field, though it is shifted to the left for Physical Sciences and Engineering (i.e., they experience a higher decrease in publications of starting in Industry vs Academia) and slightly shifted to the right for Social Sciences.

## 6.2    Mechanisms

I now discuss more specifically what explains the differences in treatment effects previously found in Figure 8. Differences in treatment effects across individuals could stem from differences in the outcomes in the 'treated' state, i.e., if individuals were to start their career in Industry ($\mathbf{E}(Y_1|X, D, p)$) and/or in the 'untreated' state, i.e., if individuals were to start their career in Academia ($\mathbf{E}(Y_0|X, D, p)$). One advantage of the separate approach that I use to estimate MTE is that I can calculate the values of the potential outcomes for each individual in each state. This allows me to investigate where differences in treatment effects come from. I can also examine whether differences arise from differences in observable vs unobservable characteristics. In what follows, I refer to 'potential earnings in Industry' (resp. in Academia) as the earnings one would earn if she were to start her career in Industry (resp. in Academia). Similarly, I refer to 'potential publications in Industry' (resp. in Academia) as the publications one would produce if she were to start her career in Industry (resp. Academia).

*Potential Earnings in Industry and in Academia:*  I first compare potential earnings in Industry and potential earnings in Academia. Figure 9 shows that the higher earnings gains experienced by those who start in Industry comes from both a lower earnings potential in Academia and a higher earnings potential in Industry. Overall, those who start their career in Industry appear to be more valuable to firms and less valuable to universities, compared to individuals who start in Academia.



(a) Potential earnings in Academia          (b) Potential earnings in Industry

Figure 9: Potential Earnings, by sector joined

*Potential Publications in Industry and in Academia:*  Second, I compare potential publications in Industry and potential publications in Academia. Figure 10 shows that those joining the private sector would publish relatively at par in Academia as those who start their career in the academic sector, but they publish more in the private sector.

In order to dig deeper into this finding, I classify individuals into two groups: those who work for a large employer (more than 5,000 employees) vs those who don't.[29] I then calculate the expected

---

[29]68% of individuals observed in Academia are identified as having a large employer. 57% of individuals observed in

(a) Potential Publications in Academia

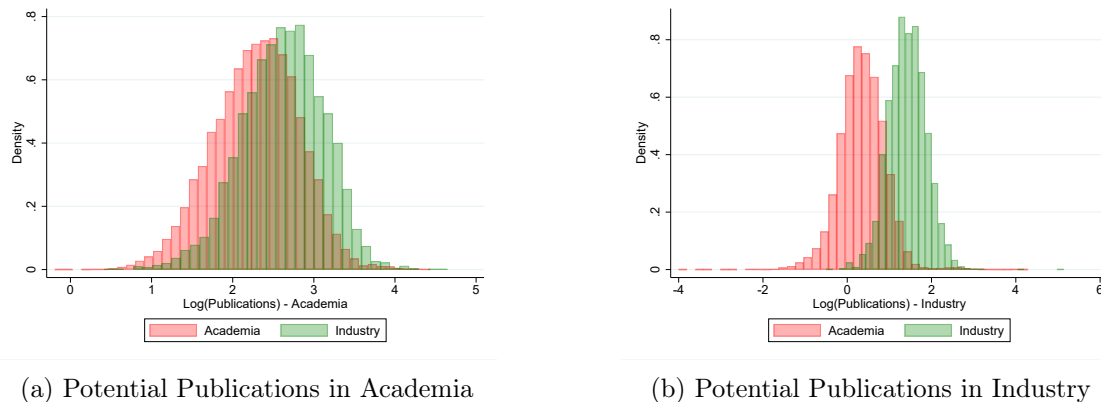(b) Potential Publications in Industry

Figure 10: Potential Publications, by sector joined

treatment effects in a similar fashion as before, using as dependent variable an indicator for working for a large employer. Figure A.12 shows the results: potential employer size in Academia is relatively similar for those who start in Industry and those who start in Academia. This is expected, as Academia consists of a fairly homogeneous pool of employers. However, when examining the type of employers that individuals would join if they were to enter the private sector, I find that individuals who start in Industry are more likely to join large firms than individuals who start in Academia. This is consistent with the fact that individuals in Academia value more academic-type environments that give them more freedom and autonomy. As a consequence, if they were to join the private sector, they would tend to join smaller firms which are more likely to provide this sort of environment (Gambardella et al., 2015). This might explain the higher publication potential in Industry of those who start in the private sector, as larger firms are likely to have more data and resources available for scientists, allowing them to publish more (Gambardella et al., 2015).

The simultaneous comparison of earnings and publications potential in Academia shows that those starting in the private sector have a lower earnings potential in Academia even though their publication potential in Academia is relatively at par with those who start in Academia. This suggests that individuals who start their career in Industry either (i) lack non-publication related skills that are valued in Academia (*productivity mechanism*) or (ii) have higher non-pecuniary returns in Academia than those who join the academic sector at graduation which would make them being paid less (*compensating differential mechanism*).

In order to disentangle these 2 mechanisms, I use information about individuals' satisfaction with different job attributes: (1) job's opportunities for advancement (2) job's benefits (3) job's intellectual challenge (4) job's degree of independence (5) job's location (6) job's level of responsibility (7) job's security (8) job's salary and (9) job's contribution to society. For each attribute, I create an indicator variable equal to 1 if individuals declare being very satisfied and 0 otherwise. Figure A.13 shows the average value of each indicator by sector joined at graduation. I construct a mea-

---

Industry are identified as having a large employer.

sure of non-pecuniary returns by averaging the 7 indicators that capture individuals' satisfaction with non-money related job attributes (i.e., I exclude individuals' satisfaction with salary and with benefits).[30] I then run a MTE analysis, similar as before. I find that those who start their career in the private sector would have lower non-pecuniary returns in Academia compared to individuals who start their career in the academic sector. This suggests that mechanism (i) is more likely: individuals who start their career in Industry lack non-publication-related skills that are valued in Academia. This could include skills such as teaching, management (e.g., ability to run a lab), the capacity to generate vs 'execute' ideas (Aghion et al., 2008), research quality, the capacity to win grants and overall influence on the research community.[31]

Overall, (compared to individuals who start in the academic sector) individuals who start in the private sector appear to: (i) have more Industry-specific skills (ii) lack non-research related academic skills (iii) be more likely to join large firms that are also engaged in Open Science.

## 6.3   Guaranteed and competing talent pools

Figures 9 and 10 show that the green and red histograms overlap for some values of the potential outcomes in each sector. In this section, I use this insight to study more specifically what this suggests regarding competition for talent between Industry and Academia.

I first focus on the earnings outcome. In order to simplify the analysis, I split potential earnings in Industry and potential earnings in Academia in 10 quantiles. I then create the 10x10 matrix with the 10 quantiles related to Academia on the x-axis and the 10 quantiles related to Industry on the y-axis. Figure 11 shows how this can be used to classify individuals. For instance, individuals in the top right quadrant would be on the right tail of the earnings distribution in both sectors and I therefore characterize them as being 'stars'. Individuals in the top left quadrant would be on the right tail of the earnings distribution in Industry and I call them 'Industry-specialists'. Individuals in the bottom right quadrant would be on the right tail of the earnings distribution in Academia and I call them 'Academic scholars'. Finally, individuals in the bottom left quadrant would be on the left tail of the earnings distribution in both sectors and I call them 'modest achievers'.

We could expect that firms will primarily try to hire individuals in the upper part of the Industry-earnings distribution. Similarly, we could expect that universities will primarily try to hire individuals in the right part of the Academia-earnings distribution. Hence, we could expect firms and universities to compete for stars. However, because individuals' sector of employment is a resultant of demand and supply forces, the shape of the margin of competition remains an empirical exercise.

I therefore bring this framework to the data. Figure 12 shows the result. Each cell in the

---

[30]Creating an indicator equal to 1 if this average value is higher than 0.5 and 0 otherwise gives similar results.

[31]While it is hard to disentangle precisely the reasons why some individuals would be paid less in the academic sector despite a similar number of publications, I use additional information provided by my dataset to try to examine observable differences in the type of research performed in the academic sector by individuals who start in Industry vs Academia. I do not find significant differences in the average number of citations that both pools would receive within 2 or 5 years of publication in the academic sector. I also do not find differences in their propensity to have at least one publication published in a top journal.
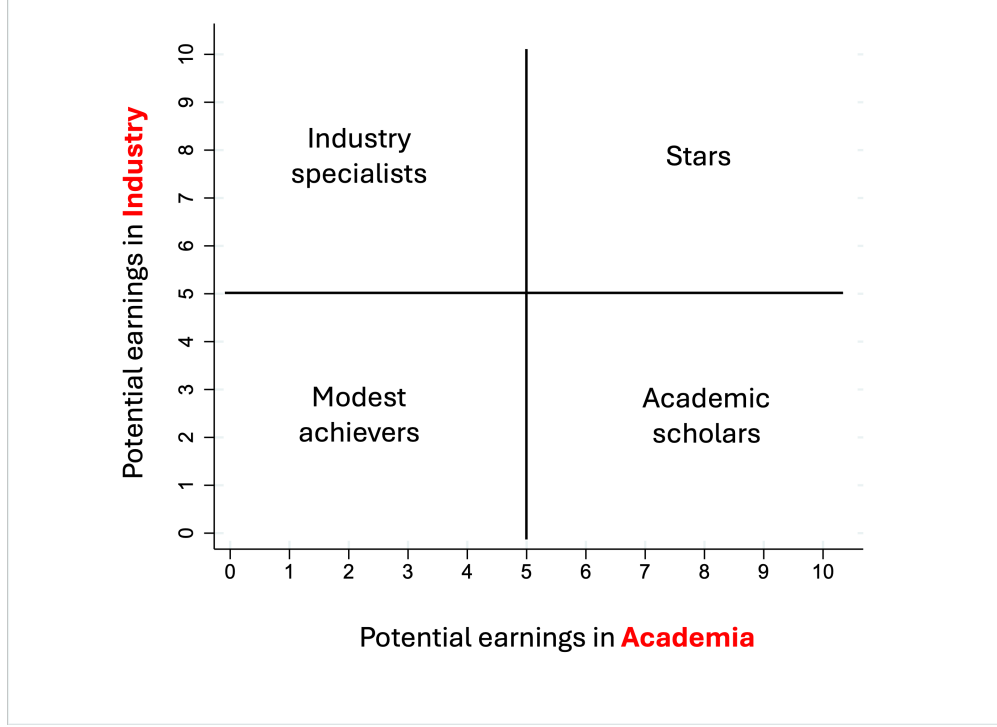
Figure 11: Classification of individuals

matrix represents a combination of potential earnings in Industry (y-axis) and potential earnings in Academia (x-axis). The color of each cell indicates whether individuals in that cell predominantly start in Academia, predominantly start in Industry, or start in either one sector. Yellow squares represent cells predominantly populated by individuals starting in Industry (I call them 'guaranteed Industry'), while purple squares represent cells predominantly populated by individuals starting in Academia (I call them 'guaranteed Academia').[32] Cells in pink represent the combinations of potential earnings for which I observe individuals who start in Industry, as well as individuals who start in Academia. These are the cells that identify the overlap between the pool of individuals who starts in Industry and the pool of individuals who starts in Academia. Said differently, the pink cells represent the margin of competition between Industry and Academia. The size of the squares is proportional to the share of individuals in the sample observed in that cell, irrespective of sector joined.

First, the location of the larger squares along the bottom-left to top-right diagonal indicates that there is a positive correlation between potential earnings in Academia and potential earnings in Industry (0.27 in my sample). Note that the higher number of purple squares is consistent with the fact that I observe more individuals in Academia than in Industry in my sample.

The pool of individuals who start in Industry is rather small and includes individuals with high potential earnings in Industry - so that they are valuable for firms - and relatively low potential

---

[32]A cell is predominantly populated by individuals starting in Industry (resp. Acad) if the share of individuals in that cell who start in Industry (resp. Acad) is greater or equal than 90%.

Figure 12: Pools of talent based on potential earnings

earnings in Academia. The pool of individuals who start in Academia is large, but comprises individuals with varying academic-related quality. The margin of competition is relatively wide and primarily includes individuals in the upper-half of the distribution of potential earnings in Industry, but with various levels of Academia-related skills. The pool of stars located in the top right quadrant predominantly joins the academic sector.

Given the importance of stars for firms' performance (Kehoe and Tzabbar, 2015; Kim and Makadok, 2022), it is important to shed light on whether this last result is primarily driven by demand- or supply- mechanisms. Indeed, individuals' sector at the start of their career results from a two-step process: first, the individual decides whether to apply to the private sector, and then firms decide whether to extend a job offer. Hence, stars could be observed predominantly in the academic sector primarily because (i) firms do not hire them (maybe because stars are hard to identify) or (ii) stars' do not apply to the private sector because their utility in Academia is higher.[33] Since stars are in the right-tail of both the Academia- and Industry-related earnings distribution, their preferences for the private vs the academic sectors could inform us on whether (i) or (ii) is more likely to drive the previous findings.[34] If (i) is the primary explanation, stars should have similar preferences for non-pecuniary returns as other individuals who are part of the Industry pool. If (ii) is the primary explanation, stars should have similar preferences for non-pecuniary returns as other

---

[33]Of course, stars could apply to both industrial and academic positions and then primarily choose an academic offer. I think of this last case as being similar as (ii) as it represents a supply-driven mechanism explaining why stars would predominantly be observed in the academic sector.

[34]Individuals' utility is determined by the sum of their pecuniary and non-pecuniary returns.

individuals who are part of the Academia pool. To shed light on this question, I use information about individuals' preferences for job attributes. I run several regressions where I compare the preferences for job attributes of stars vs individuals who start in the private sector. Similarly, I compare the preferences for job attributes of stars vs individuals who start in the academic sector. Results are presented in Tables ==[awaiting export approval]==. Overall, I find that stars' preferences are closer to those of individuals who start their career in the academic sector, providing more support for hypothesis (ii).[35]

In order to shed more light into the observed pattern for the margin of competition, I also explore who are the individuals shifted into the private sector when firms' vs universities demand increases. To that end, I plot the same matrix as before as a function of the value of the instrument: when the instrument is below its $25^{th}$ percentile, when the instrument is between its $25^{th}$ and its $50^{th}$ percentile, when the instrument is between its $50^{th}$ and its $75^{th}$ percentile and when the instrument is above its $75^{th}$ percentile. One can think of these figures as representing the pools of talent when firms' vs universities demand is respectively very low, low, high and very high. Results are presented in Figures ==[awaiting export approval]==. When firms' demand is low, those starting in Industry are individuals in the very top left, with very high earnings in Industry and very low earnings in Academia. Said differently, the first individuals to be hired are those highly valuable for firms and with low outside option in the academic sector. As firms' demand increases, individuals who are moved to the private sector are progressively pulled from the middle right of the graph. These are individuals highly valuable for firms but with better outside option in the academic sector, who become closer and closer to be part of the pool of stars. As firms' vs universities' demand becomes very high, there is a shift in the distribution of individuals who appear in the private sector: individuals are pulled from the lower left part of the graph. These are individuals who remain somehow valuable to firms (though less than before) but with low outside option in the academic sector.

Combined together, these results suggest that firms first hire individuals who are highly valuable for them and have low outside options in the academic sector. As firms' demand increases, they will hire individuals who remain highly valuable for them but have good outside option in the academic sector. Progressively, firms enter the pool of talent which includes individuals with both very good outside options in the academic sector and a higher inclination for non-pecuniary returns offered by this sector. As the hiring process becomes harder for firms, they then move on to hire individuals who are less highly valuable for them, but have low outside options in the academic sector and a strong preference for starting their career in Industry.[36]

I then explore how the competing margin varies by gender, nationality and field. Figure A.15 reproduces this exercise for male and female separately. First, note that females tend to concentrate more in the lower left part of the matrix, indicating lower potential earnings in both sectors. In

---

[35]I also estimate $\varepsilon_i^I - \varepsilon_i^A$ directly by running a regression of sector joined (1 if Industry and 0 if Academia) on the difference in potential earnings in Industry vs Academia ($W_i^I - W_i^A$) estimated from the Marginal Treatment Effects and predicting the residual. I find that $\varepsilon_i^I - \varepsilon_i^A$ is negative and close to 0, for stars, providing more support for (ii).

[36]Estimating $\varepsilon_i^I - \varepsilon_i^A$ for this group of individuals shows a strongly positive value on average.

contrast, males are more prevalent in the top right part of the matrix, reflecting higher potential earnings in both sectors. This difference remains even when accounting for differences in field and experience.[37] The figure also shows that male are more likely to start in the private sector though overall, the margin of competition appears roughly similar regarding the distribution of potential outcomes it draws from.

I repeat this exercise in Figure A.16 splitting the sample by nationality. International students have a higher share of individuals in the top left quadrant and a lower share of individuals in the bottom right quadrant compared to Americans. Said differently, a higher share of Internationals appear to have high earnings in Industry while a lower share of Internationals appear to have high earnings in Academia. Since it seems unlikely that Americans are discriminated by firms, one possible explanation for this finding could be that International students who leave their country to come to the US are positively selected from the pool of 'Industry-specialists'. Moreover, among Internationals, a higher share of individuals go to the private sector compared to Americans. Finally, there is more competition for the pool of stars among Americans (and more competition overall for Americans, with the margin being wider). This could be due to firms' reluctance to expand their hiring for Internationals beyond a certain combination of potential earnings (e.g., because of VISA concerns) or to Internationals' more fragmented preferences for Industry vs Academia. While my dataset does not allow me to speak to firms' hiring practices directly, I can shed light on this question indirectly by analyzing heterogeneity in preferences by nationality. I restrict the sample to the combination of potential earnings that are part of the margin of competition for Americans (this restricted sample includes both Americans and Internationals). I then estimate the difference in non-pecuniary returns $(\varepsilon_i^I - \varepsilon_i^A)$ by regressing the indicator variable equal to 1 if the individual starts in Industry on the difference in potential earnings in Industry vs Academia $(W_i^I - W_i^A)$ and predicting the residual. I then regress the difference in non-pecuniary returns on an indicator equal to 1 if the individual is American and 0 otherwise, as well as other controls.[38] I find that compared to Internationals, Americans who are part of this restricted sample have a higher preference for the academic sector. Based on this result, if the difference in the competing margin between Americans and Internationals was only coming from differences in preferences, we would expect the margin for Americans to be narrower. Instead, the fact that the margin of competition is wider for Americans implies that firms might be more willing to compete for a broader pool of Americans compared to Internationals.

Figure A.17 examines differences by field. Social Sciences and Engineering exhibit different patterns compared to other fields. The margin of competition in Social Sciences covers a broader combination of potential earnings. Engineering has a higher share of individuals starting in the private sector and firms hire from a broader combination of potential earnings, making the margin

---

[37]I regress each potential earnings on a linear and quadratic terms for experience as well as PhD major fixed effects. I then recalculate the quantiles for each potential earnings on the whole sample. Women remain concentrated below the -45 degree line while men remain concentrated above the -45 degree line.

[38]Controls include a linear and quadratic terms for experience, indicators for mother and father education, race and PhD major × doctoral institution fixed effects.

of competition look narrower. Interestingly, the broader 'Guaranteed Industry' pool in Engineering comes primarily from a deeper penetration of the pool of stars. Focusing on individuals above the 45 degree line, I find that Life Sciences has a greater taste for non-pecuniary returns in Academia vs Industry, followed by Social Sciences, Physical Sciences, Computer Sciences and Engineering. This implies that firms' demand for PhDs in Life Sciences is higher than firms' demand for PhDs in Social Sciences. The ordering of preferences among the other fields matches the pattern, so I am not able to disentangle the supply vs demand forces.

I then reiterate the same exercise with the publications outcome. Results are presented in Figure 13.[39]



Figure 13: Pools of talent based on potential publications

First, note a positive correlation between potential publications in Industry and in Academia (0.54 in my sample). The pool of individuals 'guaranteed Industry' is highly concentrated and located in the top right part of the matrix, with high publication potentials in both sectors. The margin is quite wide and primarily includes individuals in the top part of the graph.[40] I reiterate the previous exercise and plot the matrix as a function of individuals' value of their instrument and relate it to the hiring dynamics previously discussed for the earnings outcomes. I find that firms first hire individuals in the top of the matrix (i.e., with high publication potential in Industry), regardless of their publication potential in Academia. As firms' demand increases, individuals are pulled first

---

[39]I use 0.2 and 0.8 as the shares used to delineate the 3 groups because there is less variation for this outcome.

[40]It is hard to precisely discuss individuals' productivity related to their publication output in the private sector,as the observed number of publications results from both individuals producing research *and* firms allowing them to publicly disclose it.

from the top right of the matrix and progressively more and more from the top left. Note that stars are located in the middle right of the publication matrix, with high publication potential in the academic sector and an average publication potential in the private sector: since these individuals are less likely to move to the private sector, firms don't primarily hire in this part of the matrix. I condition on individuals being in the private sector and run a regression of the average number of research fields covered in publications on each of the 10 quantiles of the publication potential in the academic sector, controlling for gender, nationality, field, parental education, race, doctoral institution and experience. Results show that as individuals move from the right to the left, the number of research fields decreases. This might indicate that firms first hire scientists to work on general research and then move to more specialized topics (Teodoridis et al., 2019).

I find little differences by gender (Figure A.18). Figure A.19 shows that International who start in the private sector appear to have a more diverse publication potential in Academia compared to individuals who start in Industry. Figure A.20 shows interesting differences across fields: individuals in Life Sciences and Social Sciences who start in the private sector are drawn exclusively from the top of both distributions. In contrast, individuals in Computer Sciences, Physical Sciences and Engineering who start in the private sector have a high publication potential in Industry, but are drawn from various parts of the publication potential distribution.

# 7 Robustness and Limitations

*[To be done]*

# 8 Conclusion

Through the lens of a new dataset covering more than 40 cohorts of doctorate holders with information about their earnings and their publication output, this paper identifies the margin of competition for scientific talent between Industry and Academia.

Using an instrumental variable strategy, I estimate individual-specific treatment effects of starting in Industry vs Academia on earnings and publications. I find substantial heterogeneity across individuals that I use to characterize the pools of talent that start in Industry vs Academia. Individuals who start their career in the private sector have higher earnings gains of joining the private sector compared to those who start in Academia, that I relate to both a lower earnings potential in Academia and a higher earnings potential in Industry. Interestingly, I find that individuals who start in Industry would have a similar publication potential in the academic sector as individuals who start in Academia. Combined with their lower earnings potential in Academia, this result suggests that they lack skills non related to their capacity to publish that are valued in this sector. Individuals who start in the private sector also have a higher publication potential in Industry, which I relate to a higher likelihood of joining large firms. As a consequence, they also experience a lower decrease in publications of joining the private sector, compared to those who start in Academia.

The pool of individuals who start in Academia and the pool of individuals who start in Industry have an overlapping distribution of treatment effects for both outcomes. I therefore try to characterize more precisely the margin of competition between firms and universities. Classifying individuals as a function of the combination of their potential earnings in each sector, I find that individuals who start in the private sector appear to be highly valuable for firms. Stars, defined as individuals who would be top-earners in both sectors, appear to predominantly join the academic sector. The margin of competition is quite wide, and is shaped by both demand and supply forces. In particular, I find that firms first hire individuals at the right-tail of the earnings distribution in Industry but with low outside option in the academic sector. As firms' demand for PhD increases, individuals' outside option in Academia increases too. Hiring becomes more difficult, so that firms then move to individuals slightly less valuable for them but with low outside option in the academic sector. Firms seem to hire primarily individuals with a high publication potential in the private sector, starting with 'generalists' and then moving to individuals with more specialized research topics. I find differences in the margin across nationalities and fields.

Overall, these findings have implications for firms' human capital strategy. They also raise additional questions that provide a fertile ground for future research. In particular, disentangling the supply and demand mechanisms behind individuals' sorting and how they vary across gender and nationality would allow to provide more nuanced insights. While my sample gives me only limited information about employers, further work could also incorporate heterogeneity across firms. Finally, while my results are essentially reduced-form, more work is needed on the specific skills that firms and universities value in scientists.

# References

Agarwal, R., Ganco, M. and Ziedonis, R. H. (2009), 'Reputations for toughness in patent enforcement: Implications for knowledge spillovers via inventor mobility', *Strategic Management Journal* **30**(13), 1349–1374.

Agarwal, R. and Ohyama, A. (2013), 'Industry or academia, basic or applied? career choices and earnings trajectories of scientists', *Management Science* **59**(4), 950–970.

Aghion, P., Dewatripont, M. and Stein, J. C. (2008), 'Academic freedom, private-sector focus, and the process of innovation', *The RAND Journal of Economics* **39**(3), 617–635.

Agrawal, A., McHale, J. and Oettl, A. (2017), 'How stars matter: Recruiting and peer effects in evolutionary biology', *Research Policy* **46**(4), 853–867.

Andresen, M. E. (2018), 'Exploring marginal treatment effects: Flexible estimation using stata', *The Stata Journal* **18**(1), 118–158.

Aobdia, D., Srivastava, A. and Wang, E. (2018), 'Are immigrants complements or substitutes? evidence from the audit industry', *Management Science* **64**(5), 1997–2012.

Azoulay, P., Fons-Rosen, C. and Zivin, J. S. G. (2019), 'Does science advance one funeral at a time?', *American Economic Review* **109**(8), 2889–2920.

Ballesteros-Rodríguez, J. L., De Saá-Pérez, P., García-Carbonell, N., Martín-Alcázar, F. and Sánchez-Gardey, G. (2022), 'Exploring the determinants of scientific productivity: A proposed typology of researchers', *Journal of Intellectual Capital* **23**(2), 195–221.

Bao, J. (2024), 'Gender gap in stem entrepreneurship: Effects of the affordable care act reform', *Strategic Management Journal* .

Barney, J. (1991), 'Firm resources and sustained competitive advantage', *Journal of management* **17**(1), 99–120.

Björklund, A. and Moffitt, R. (1987), 'The estimation of wage gains and welfare gains in self-selection models', *The Review of Economics and Statistics* pp. 42–49.

Black, I., Hasan, S. and Koning, R. (2024), 'Hunting for talent: Firm-driven labor market search in the united states', *Strategic Management Journal* **45**(3), 429–462.

Byun, H., Frake, J. and Agarwal, R. (2018), 'Leveraging who you know by what you know: Specialization and returns to relational capital', *Strategic Management Journal* **39**(7), 1803–1833.

Campbell, B. A., Coff, R. and Kryscynski, D. (2012), 'Rethinking sustained competitive advantage from human capital', *Academy of Management Review* **37**(3), 376–395.

Campbell, B. A., Ganco, M., Franco, A. M. and Agarwal, R. (2012), 'Who leaves, where to, and why worry? employee mobility, entrepreneurship and effects on source firm performance', *Strategic Management Journal* **33**(1), 65–87.

Cappelli, P. (2012), Why good people can't get jobs: The skills gap and what companies can do about it, *in* 'Why Good People Can't Get Jobs', University of Pennsylvania Press.

Cappelli, P. (2019), 'Your approach to hiring is all wrong', *Harvard Business Review* **97**(3), 48–58.

Coff, R. and Kryscynski, D. (2011), 'Invited editorial: Drilling for micro-foundations of human capital–based competitive advantages', *Journal of management* **37**(5), 1429–1443.

Coff, R. W. (1997), 'Human assets and management dilemmas: Coping with hazards on the road to resource-based theory', *Academy of management review* **22**(2), 374–402.

Cohen, W. M., Sauermann, H. and Stephan, P. (2020), 'Not in the job description: The commercial activities of academic scientists and engineers', *Management Science* **66**(9), 4108–4117.

Cornelissen, T., Dustmann, C., Raute, A. and Schönberg, U. (2016), 'From late to mte: Alternative methods for the evaluation of policy interventions', *Labour Economics* **41**, 47–60.

Dasgupta, P. and David, P. A. (1994), 'Toward a new economics of science', *Research policy* **23**(5), 487–521.

Deming, D. J. and Noray, K. L. (2018), *STEM careers and technological change*, Vol. 24, National Bureau of Economic Research Cambridge, MA.

Dimmock, S. G., Huang, J. and Weisbenner, S. J. (2022), 'Give me your tired, your poor, your high-skilled labor: H-1b lottery outcomes and entrepreneurial success', *Management Science* **68**(9), 6950–6970.

Dobbie, W., Goldin, J. and Yang, C. S. (2018), 'The effects of pre-trial detention on conviction, future crime, and employment: Evidence from randomly assigned judges', *American Economic Review* **108**(2), 201–240.

Ehrenberg, R. G. (1992), 'The flow of new doctorates', *Journal of economic literature* **30**(2), 830–875.

Gambardella, A., Panico, C. and Valentini, G. (2015), 'Strategic incentives to human capital', *Strategic management journal* **36**(1), 37–52.

Ganco, M., Ziedonis, R. H. and Agarwal, R. (2015), 'More stars stay, but the brightest ones still leave: Job hopping in the shadow of patent enforcement', *Strategic Management Journal* **36**(5), 659–685.

Glennon, B. (2024), 'How do restrictions on high-skilled immigration affect offshoring? evidence from the h-1b program', *Management Science* **70**(2), 907–930.

Heckman, J. J. and Vytlacil, E. J. (1999), 'Local instrumental variables and latent variable models for identifying and bounding treatment effects', *Proceedings of the national Academy of Sciences* **96**(8), 4730–4734.

Heckman, J. J. and Vytlacil, E. J. (2007), 'Econometric evaluation of social programs, part ii: Using the marginal treatment effect to organize alternative econometric estimators to evaluate social programs, and to forecast their effects in new environments', *Handbook of econometrics* **6**, 4875–5143.

Hess, A. M. and Rothaermel, F. T. (2011), 'When are assets complementary? star scientists, strategic alliances, and innovation in the pharmaceutical industry', *Strategic management journal* **32**(8), 895–909.

Kehoe, R. R. and Tzabbar, D. (2015), 'Lighting the way or stealing the shine? an examination of the duality in star scientists' effects on firm innovative performance', *Strategic Management Journal* **36**(5), 709–727.

Kim, J. and Makadok, R. (2022), 'Where the stars still shine: Some effects of star-performers-turned-managers on organizational performance', *Strategic Management Journal* **43**(12), 2629–2666.

Nagle, F. and Teodoridis, F. (2020), 'Jack of all trades and master of knowledge: The role of diversification in new distant knowledge integration', *Strategic management journal* **41**(1), 55–85.

Ng, W. and Stuart, T. E. (2022), 'Acquired employees versus hired employees: Retained or turned over?', *Strategic Management Journal* **43**(5), 1025–1045.

Oettl, A. (2012), 'Reconceptualizing stars: Scientist helpfulness and peer performance', *Management science* **58**(6), 1122–1140.

Oyer, P. (2006), 'Initial labor market conditions and long-term outcomes for economists', *Journal of Economic Perspectives* **20**(3), 143–160.

Roach, M. and Sauermann, H. (2010), 'A taste for science? phd scientists' academic orientation and self-selection into research careers in industry', *Research policy* **39**(3), 422–434.

Roach, M. and Sauermann, H. (2023), 'Can technology startups hire talented early employees? ability, preferences, and employee first job choice', *Management Science* .

Roach, M. and Sauermann, H. (2024), 'Ability, preferences, and stem doctorate early career choices', *Working Paper* .

Roche, M. P. (2023), 'Academic entrepreneurship: Entrepreneurial advisors and their advisees' outcomes', *Organization Science* **34**(2), 959–986.

Romer, P. M. (1990), 'Endogenous technological change', *Journal of political Economy* **98**(5, Part 2), S71–S102.

Sampat, B. and Williams, H. L. (2019), 'How do patents affect follow-on innovation? evidence from the human genome', *American Economic Review* **109**(1), 203–236.

Shu, P. (2016), 'Innovating in science and engineering or'cashing in'on wall street? evidence on elite stem talent', *Harvard Business School Technology & Operations Mgt. Unit Working Paper* (16-067).

Starr, E., Frake, J. and Agarwal, R. (2019), 'Mobility constraint externalities', *Organization Science* **30**(5), 961–980.

Stephan, P. (2012), *How economics shapes science*, Harvard University Press.

Stephan, P. E., Sumell, A. J., Black, G. C. and Adams, J. D. (2004), 'Doctoral education and economic development: The flow of new ph. ds to industry', *Economic Development Quarterly* **18**(2), 151–167.

Stern, S. (2004), 'Do scientists pay to be scientists?', *Management science* **50**(6), 835–853.

Teodoridis, F., Bikard, M. and Vakili, K. (2019), 'Creativity at the knowledge frontier: The impact of specialization in fast-and slow-paced domains', *Administrative Science Quarterly* **64**(4), 894–927.

# Appendix A    Figures



Figure A.1: Share of International vs American students by PhD cohort of graduation

*Notes:* This figure shows the share of international students as a function of PhD cohort of graduation

Figure A.2: Share of Female vs Male students by PhD cohort of graduation

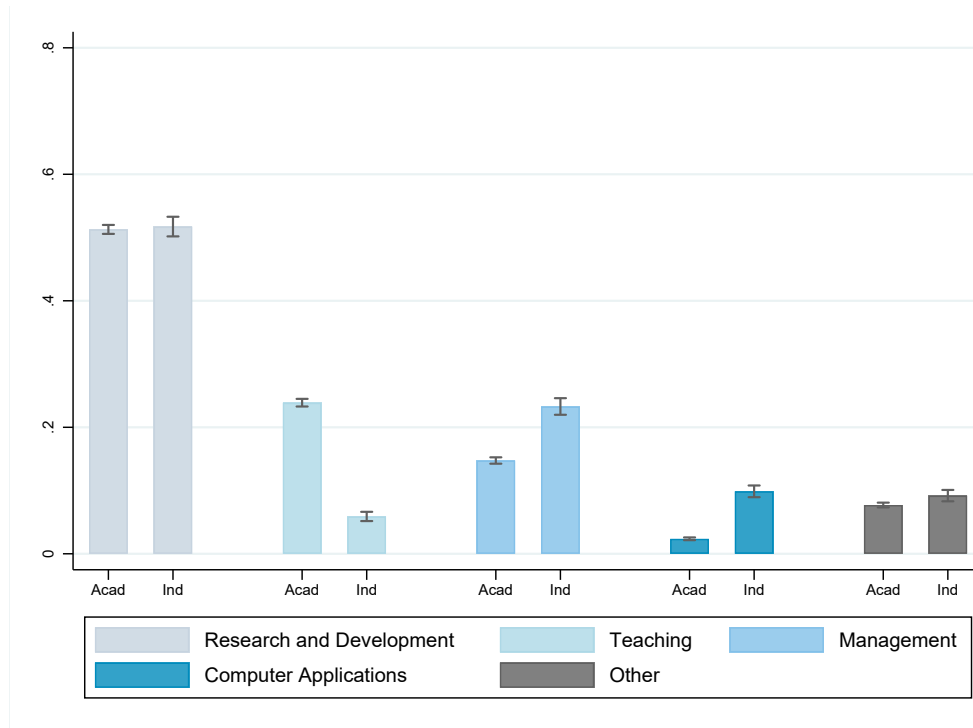*Notes:* This figure shows the share of female students as a function of PhD cohort of graduation

Figure A.3: Share of International and American students joining the private sector by PhD cohort of graduation

*Notes:* This figure shows the share of international and American students who join the private sector as a function of PhD cohort of graduation

Figure A.4: Share of Female and Male students joining the private sector by PhD cohort of graduation

*Notes:* This figure shows the share of female and male students who join the private sector as a function of PhD cohort of graduation

Figure A.5: Main activity during employment, by sector joined at graduation
*Notes:* This figure shows the main activity performed by individuals as part of their employment when surveyed as part of the SDR (i.e., during their career)



Figure A.6: Main detailed activity during employment, by sector joined at graduation
*Notes:* This figure shows the main activity performed by individuals as part of their employment when surveyed as part of the SDR (i.e., during their career) at a more granular level

(a) Potential Publications

(b) Potential Earnings

Figure A.7: MTE Curves, semi-parametric approach



(a) Expected treatment effects Publications

(b) Expected treatment effects Earnings

Figure A.8: Expected treatment effects, by gender



(a) Expected treatment effects, Publications

(b) Expected treatment effects, Earnings

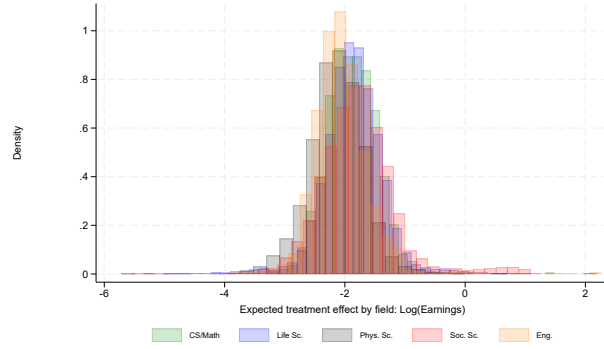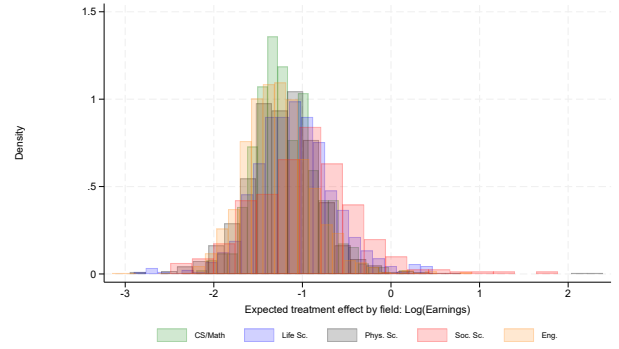Figure A.9: Expected treatment effects, by nationality

(a) Earnings in Academia

(b) Earnings in Industry

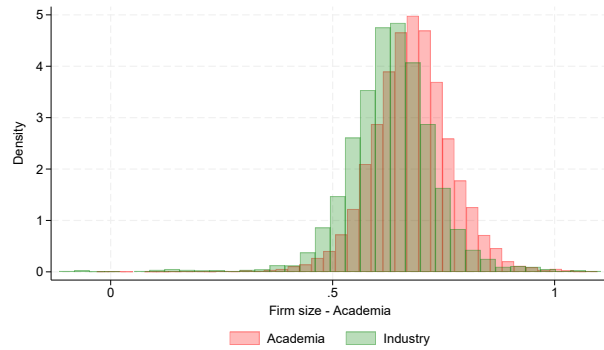Figure A.10: Expected treatment effects, by field
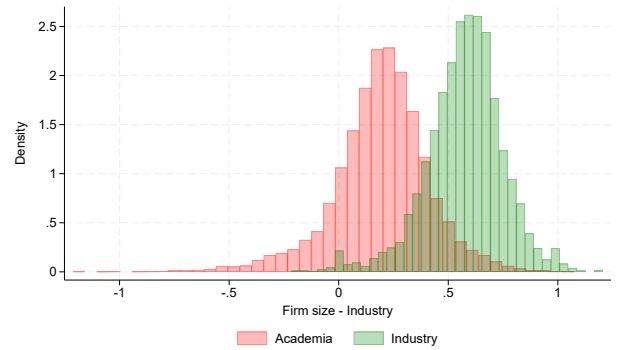


(a) Publications in Academia

(b) Publications in Industry

Figure A.11: Expected treatment effects, by field



(a) Potential employer size in Academia

(b) Potential employer size in Industry

Figure A.12: Potential employer size
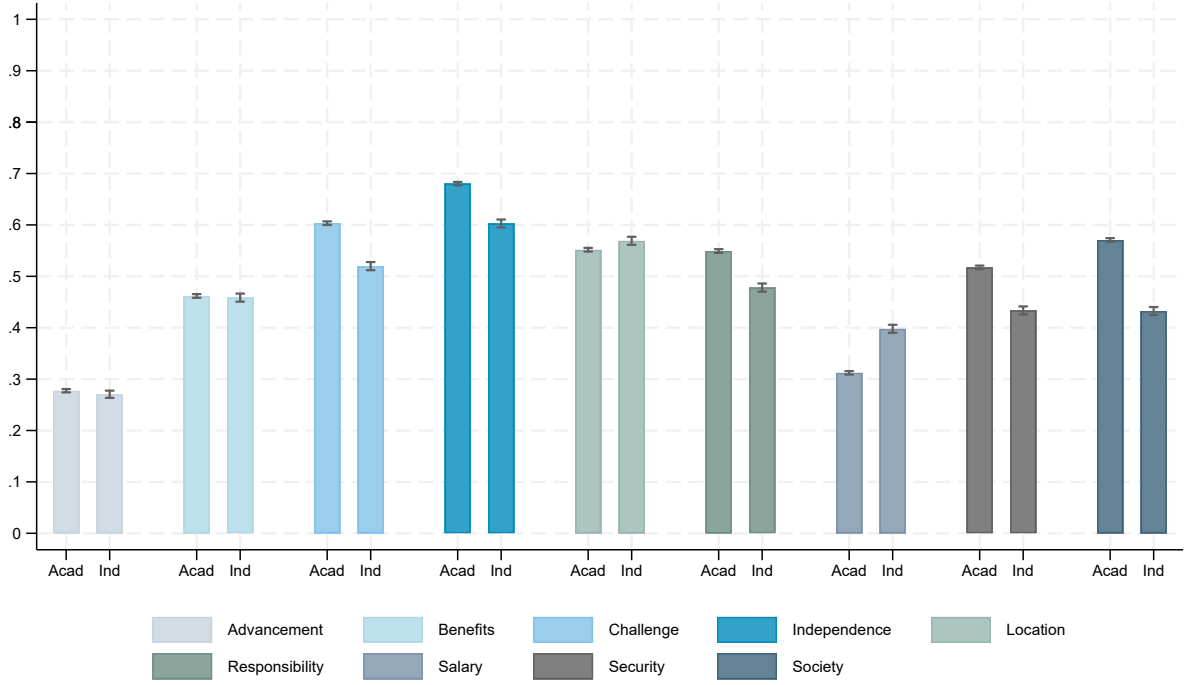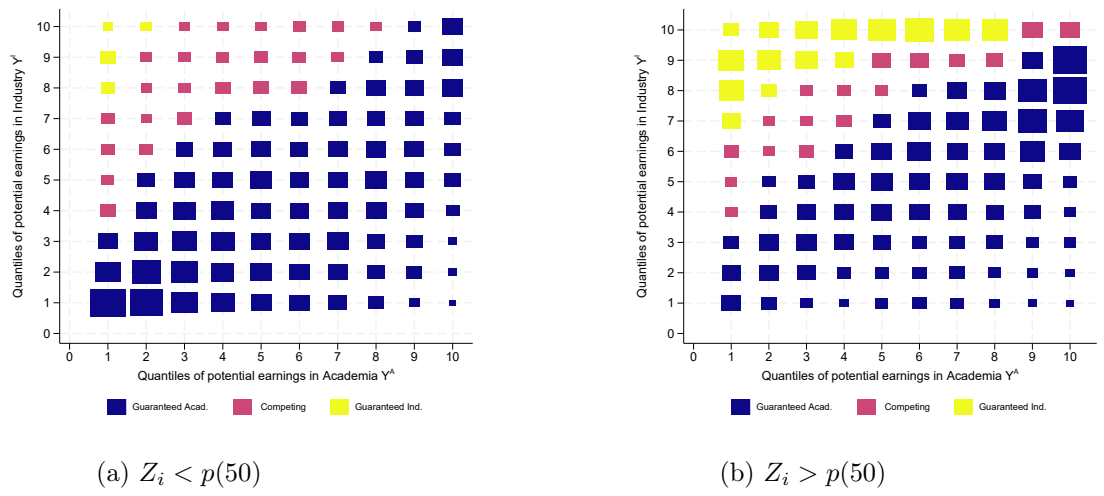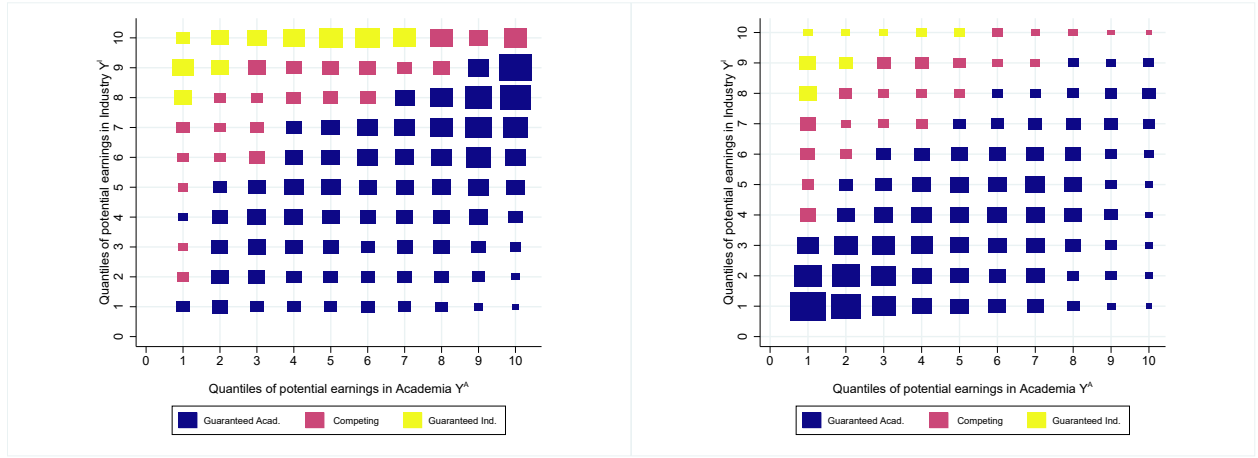
Figure A.13: Share of individuals very satisfied, by job attribute and sector joined



(a) $Z_i < p(50)$



(b) $Z_i > p(50)$

Figure A.14: Margin of competition, by value of the instrument
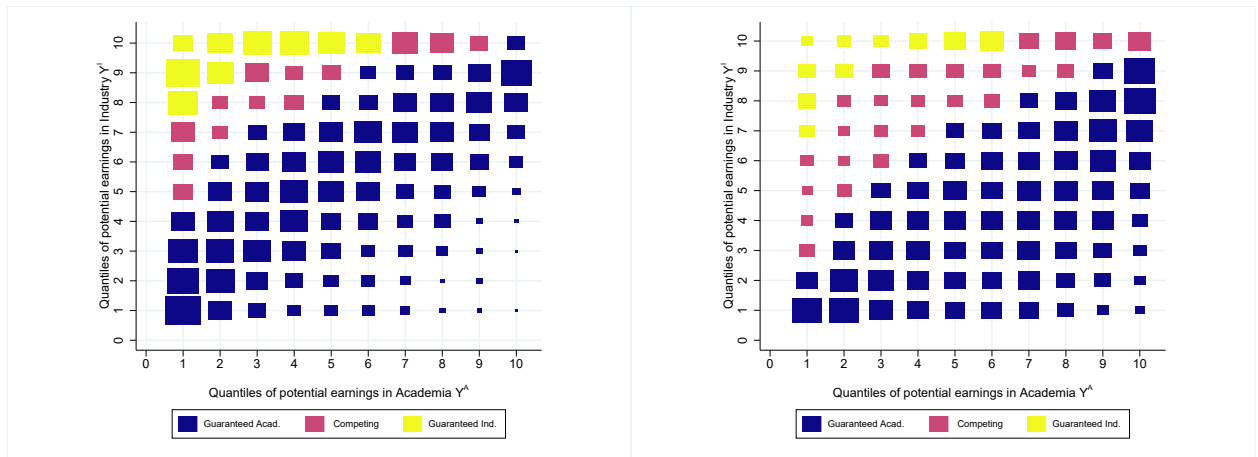
(a) Male

(b) Female

Figure A.15: Pools of talent based on potential earnings, by gender



(a) International

(b) American

Figure A.16: Pools of talent based on potential earnings, by nationality

(a) Computer Sciences

(b) Life Sciences

(c) Physical Sciences

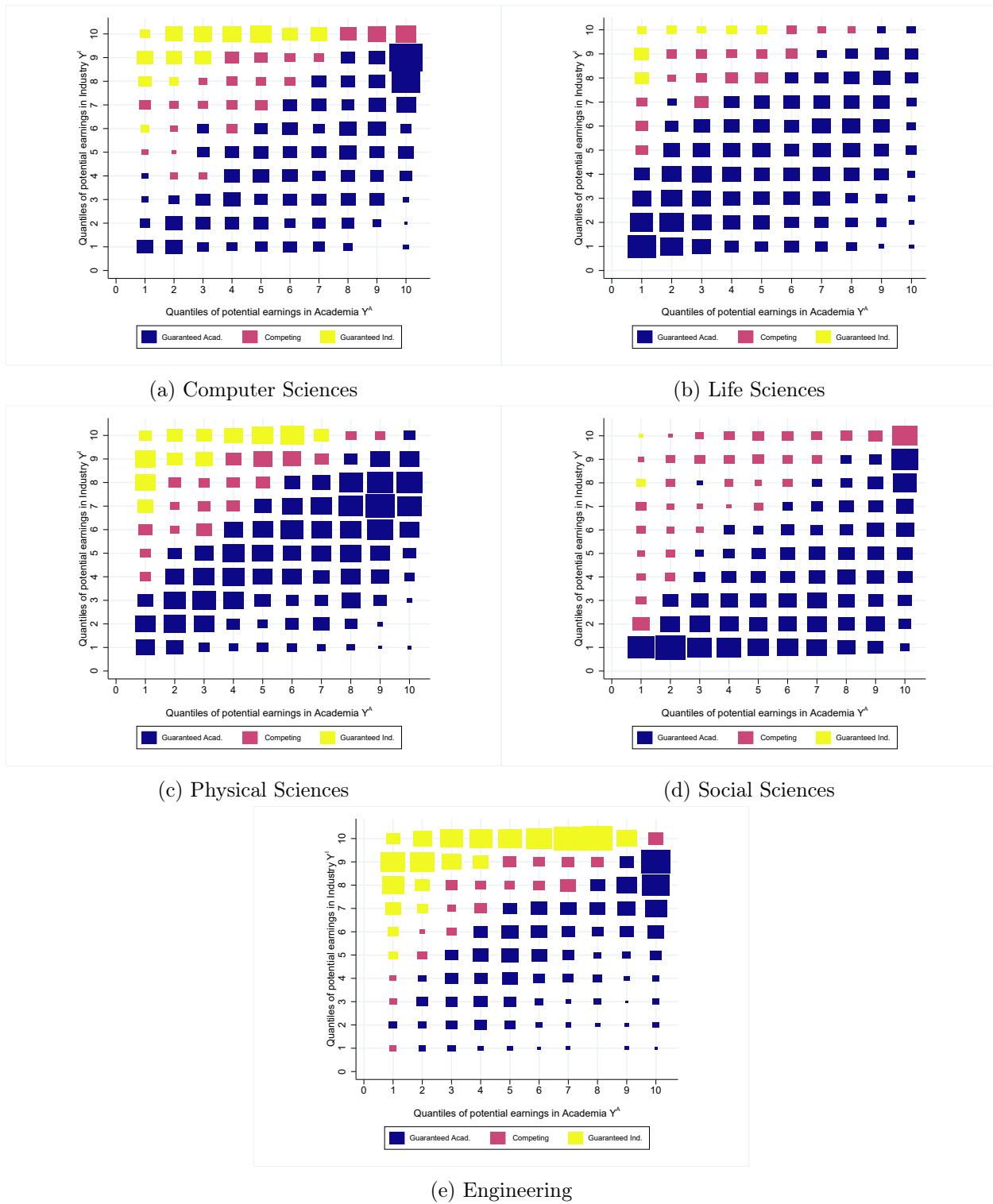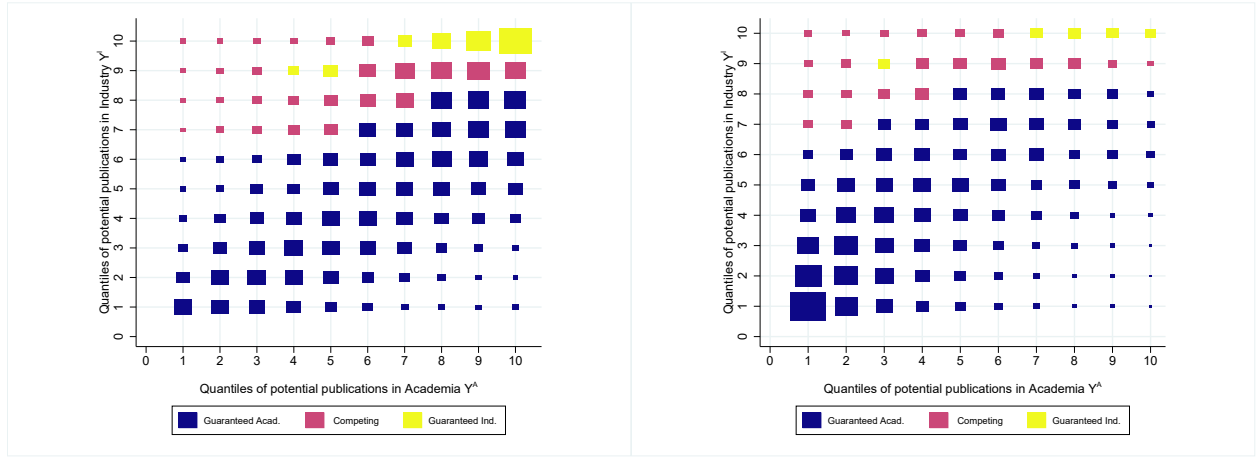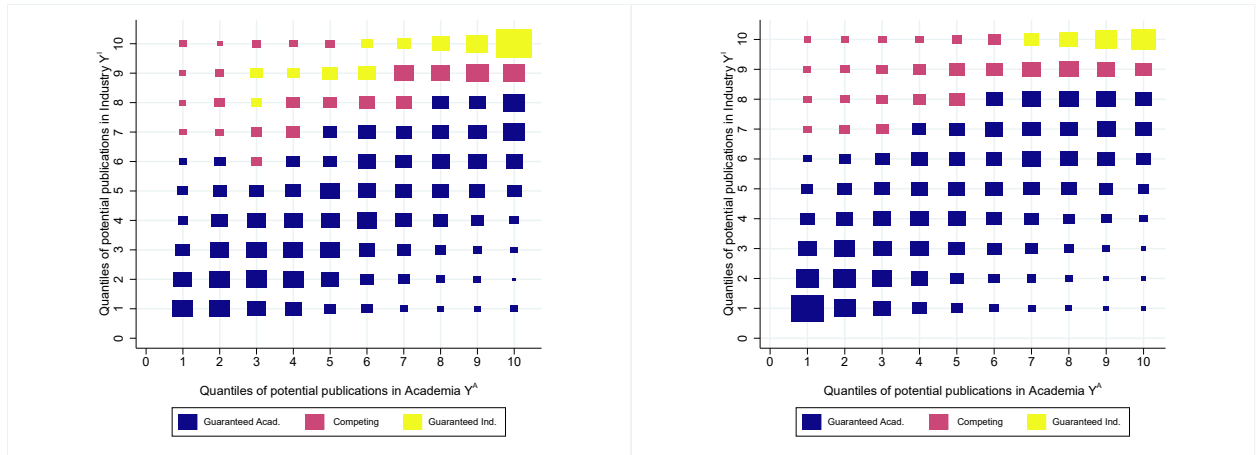(d) Social Sciences

(e) Engineering

Figure A.17: Pools of talent based on potential earnings, by field

(a) Male

(b) Female

Figure A.18: Potential outcomes, Publications



(a) International

(b) American

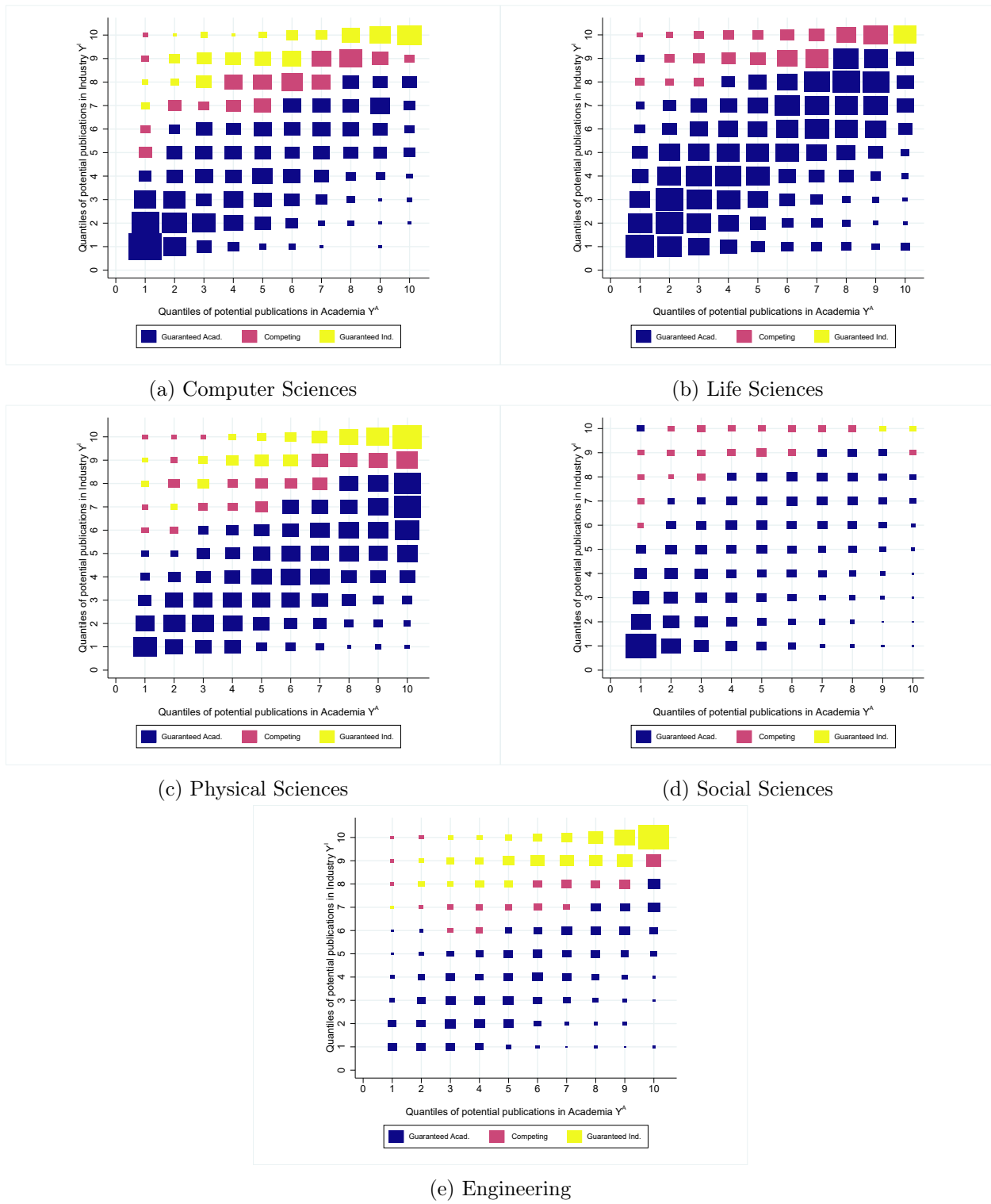Figure A.19: Potential outcomes, Publications

(a) Computer Sciences

(b) Life Sciences

(c) Physical Sciences

(d) Social Sciences

(e) Engineering

Figure A.20: Potential outcomes, Publications

# Appendix B  Tables

Table B.1: OLS pub estimates, Publications> 0

|  | Log(1+Publications (stock) | | | | |
|  | (1) | (2) | (3) | (4) | (5) |
|---|---|---|---|---|---|
| Industry | -0.776*** | -0.817*** | -0.807*** | -0.819*** | -0.777*** |
|  | (0.0186) | (0.0241) | (0.0239) | (0.0263) | (0.0290) |
| Demographics | Yes | Yes | Yes | Yes | Yes |
| Experience | Yes | Yes | Yes | Yes | Yes |
| PhD major FE |  | Yes | Yes | Yes | Yes |
| Doct. Inst. FE |  |  | Yes | Yes | Yes |
| PhD major x Doct. Inst. FE |  |  |  | Yes | Yes |
| Taste |  |  |  |  | Yes |
| Observations | 19,945 | 19,945 | 19,920 | 18,405 | 16,084 |
| R-sq | 0.149 | 0.183 | 0.210 | 0.314 | 0.338 |

*Notes*: This table reports the OLS estimates for publications using within department (defined as PhD major × doctoral institution) variation, keeping individuals with at least one publication. Standard errors (in parentheses) are clustered at the doctoral institution level.

Table B.2: OLS citations estimates

| | Log(1+Cites-weighted Publications) | | | | |
| | (1) | (2) | (3) | (4) | (5) |
|---|---|---|---|---|---|
| Industry | -1.731*** | -1.477*** | -1.463*** | -1.470*** | -1.409*** |
| | (0.0418) | (0.0437) | (0.0441) | (0.0479) | (0.0524) |
| Demographics | Yes | Yes | Yes | Yes | Yes |
| Experience | Yes | Yes | Yes | Yes | Yes |
| PhD major FE | | Yes | Yes | Yes | Yes |
| Doct. Inst. FE | | | Yes | Yes | Yes |
| PhD major x Doct. Inst. FE | | | | Yes | Yes |
| Taste | | | | | Yes |
| Observations | 21,024 | 21,024 | 20,999 | 19,488 | 17,081 |
| R-sq | 0.108 | 0.219 | 0.261 | 0.363 | 0.375 |

*Notes*: This table reports the OLS estimates for cites-weighted publications using within department (defined as PhD major × doctoral institution) variation. Standard errors (in parentheses) are clustered at the doctoral institution level.

Table B.3: 2SLS - Publications, Publications>0

| | Log(1+Publications) | | |
| | (1) | (2) | (3) |
|---|---|---|---|
| Industry | -0.940** | -0.971** | -0.964** |
| | (0.455) | (0.480) | (0.491) |
| Demographics | Yes | Yes | Yes |
| Experience | Yes | Yes | Yes |
| PhD major FE | Yes | Yes | Yes |
| Doct. Inst. FE | No | Yes | Yes |
| Macro cond. | No | No | Yes |
| F-stat | 30.52 | 28.17 | 28.15 |
| Observations | 19,945 | 19,945 | 19,945 |
| R-sq | 0.182 | 0.209 | 0.231 |

*Notes*: This table reports the 2SLS results for publications using individuals with at least one publication. Standard errors (in parentheses) are two-way clustered at the PhD major and PhD cohort level.

Table B.4: 2SLS - Cites-weighted Publications

|  | Log(1+Cites-weighted Publications) | | |
| --- | --- | --- | --- |
|  | (1) | (2) | (3) |
| Industry | -2.684*** | -2.738*** | -2.749*** |
|  | (0.806) | (0.837) | (0.850) |
| Demographics | Yes | Yes | Yes |
| Experience | Yes | Yes | Yes |
| PhD major FE | Yes | Yes | Yes |
| Doct. Inst. FE | No | Yes | Yes |
| Macro cond. | No | No | Yes |
| F-stat | 30.31 | 28.12 | 28.15 |
| Observations | 21,024 | 21,024 | 21,024 |
| R-sq | 0.182 | 0.223 | 0.230 |

*Notes*: This table reports the 2SLS results for cites-weighted publications. Standard errors (in parentheses) are two-way clustered at the PhD major and PhD cohort level.

# Appendix C    Sector Choice

I present here the full model of sector choice that allows individuals to change sector over their career.

Let individuals be indexed by $i$ and time be indexed by $t \in (0, T)$. There are 2 sectors in the economy: Industry, indexed as $I$ and Academia, indexed as $A$. I assume for simplicity that I observe a unique generation of individuals so that $t$ can also be conceptualized as years of experience. Each individual receives two endowments at birth: an individual-specific productivity parameter in Industry $\delta_i^I$ and an individual-specific productivity parameter in Academia $\delta_i^A$. These parameters capture how more/less productive individual $i$ is in Industry and Academia compared to the average. $\delta_i^I$ is composed of skills that are common to both Industry and Academia, as well as skills that are specific to the private sector. Similarly, $\delta_i^A$ is composed of skills that are common to both Industry and Academia, as well as skills that are specific to the academic sector. Each individual works between $t = 0$ and $t = T$ and is characterized by a vector $j(i)$ representing the sectors she works in during each time period: $j(i) = (j_0(i), j_1(i), ..., j_T(i))$. For clarity, I will use $j$ instead of $j(i)$ in the equations that follow, but the reader should keep in mind that $j$ is individual-specific. $j = (I, I, ..., I)$ for individuals who spend their whole career in Industry and $j = (A, A, ..., A)$ for individuals who spend their whole career in Academia, but individuals are allowed to change sector during their career. I am interested in $j_0$, the sector where individuals *start* their career after graduating form their PhD.

Define $W_i^{j_t}$ the potential earnings of individual $i$ at time $t$ if she works in sector $j_t \in (I, A)$ at time $t$. We can decompose $W_{it}^{j_t}$ into a sector-mean and an individual-specific component:

$$W_i^{j_t} = \overline{W_t^{j_t}} + \delta_i^{j_t} \tag{14}$$

with $\overline{W_t^{j_t}}$ the average earnings in sector $j_t$ for individuals with experience $t$ and $\delta_i^{j_t}$ an individual-specific component.

The utility $u_i^{j_t}$ of individual $i$ at time $t$ who works in sector $j_t$ equals:

$$u_i^{j_t} = (W_i^{j_t} + \chi_i^{j_t}) + \varepsilon_i^{j_t} + \sigma_i^{j_t}$$

with $\chi_i^{j_t}$ a random shock to earnings at time $t$, $\varepsilon_i^{j_t}$ individual $i$'s taste for sector $j_t$ and $\sigma_i^{j_t}$ some noise. I model individual $i$'s utility of *starting* her career in sector $j_0 \in (I, A)$, called $U_i^{j_0}$, as a function of the sum over her career of her expected (discounted) earnings and an error term $\varepsilon_i^{j_0}$:

$$
\begin{aligned}
U_i^{j_0} &= \sum_{t=0}^{T} \rho^t \, \mathbf{E}\left[u_i^{j_t} \mid j_0\right] \\
&= \sum_{t=0}^{T} \rho^t \, \mathbf{E}\left[W_i^{j_t} + \chi_i^{j_t} + \varepsilon_i^{j_t} + \sigma_i^{j_t} \mid j_0\right] \\
&= \sum_{t=0}^{T} \rho^t \left(\mathbf{E}\left[W_i^{j_t} \mid j_0\right] + \mathbf{E}\left[\varepsilon_i^{j_t} \mid j_0\right]\right) \\
&= \sum_{t=0}^{T} \rho^t \left(\mathbf{E}\left[\overline{W_t^{j_t}} + \delta_i^{j_t} \mid j_0\right] + \mathbf{E}\left[\varepsilon_i^{j_t} \mid j_0\right]\right) \\
&= \sum_{t=0}^{T} \rho^t \left(p(j_t = I|j_0)(\overline{W_t^I} + \delta_i^I) + p(j_t = A|j_0)(\overline{W_t^A} + \delta_i^A) + \mathbf{E}\left[\varepsilon_i^{j_t} \mid j_0\right]\right) \\
&= \sum_{t=0}^{T} \rho^t \left(p(j_t = I|j_0)(\overline{W_t^I} + \delta_i^I) + (1 - p(j_t = I|j_0))(\overline{W_t^A} + \delta_i^A) + \mathbf{E}\left[\varepsilon_i^{j_t} \mid j_0\right]\right) \\
&= \sum_{t=0}^{T} \rho^t \left(p(j_t = I|j_0)(\overline{W_t^I} + \delta_i^I - \overline{W_t^A} - \delta_i^A) + \overline{W_t^A} + \delta_i^A + \mathbf{E}\left[\varepsilon_i^{j_t} \mid j_0\right]\right)
\end{aligned}
\tag{15}
$$

with $p(j_t = I|j_0)$ (resp. $p(j_t = A|j_0)$) being the probability that individual $i$ works in Industry (resp. Academia) at time $t$ given that $i$ started her career in sector $j_0$.

Individual $i$ *starts* her career in Industry iff :

$$j_0(i) = I \iff U_i^I - U_i^A > 0$$

$$\iff \sum_{t=0}^{T} \rho^t \Big( (p(j_t = I|j_0 = I) - p(j_t = I|j_0 = A))(\overline{W_t^I} + \delta_i^I - \overline{W_t^A} - \delta_i^A) \tag{16}$$

$$+ \mathbf{E}\, [\varepsilon_i^{j_t} | \, j_0 = I] - \mathbf{E}\, [\varepsilon_i^{j_t} | \, j_0 = A] \Big) > 0$$

Because this equation includes the difference in unobserved productivity parameters $\delta_i^I - \delta_i^A$, the sorting of individuals at graduation is not random.

# Appendix D   Selection on gains vs selection on level

I detail here the difference between selection on level and selection on gains and how they might be correlated with each other.

The first selection mechanism is selection on *gains* which refers to the value of the difference $\delta_i^I - \delta_i^A$ and influences which sector individuals choose at graduation. Individuals with higher values of $\delta_i^I - \delta_i^A$ have more to gain by going to the private sector and are thus more likely to select into that sector. The second selection mechanism refers to selection on *level* which is linked to the values of $\delta_i^I$ and $\delta_i^A$. Consider two individuals identical on observable characteristics and with the same gains $\delta_i^I - \delta_i^A$, but assume that individual 1 has higher level values of $\delta_i^I$ and $\delta_i^A$ compared to individual 2. Because earnings are the sum of a sector-mean (common to individuals 1 and 2) and an individual-specific component (higher for individual 1), individual 1's potential earning outcomes are higher than individual 2's potential earning outcomes. While the sector selection equation is only linked to $\delta_i^I - \delta_i^A$, this difference might be correlated with the levels $\delta_i^I$ and $\delta_i^A$ except if we assume constant treatment effects. To see that, let's derive the conditions under which the difference and the levels are uncorrelated:

$$Cov(\delta_i^I - \delta_i^A, \delta_i^I) = 0$$
$$Cov(\delta_i^I - \delta_i^A, \delta_i^A) = 0$$
$$\Longleftrightarrow$$
$$Var(\delta_i^I) - Cov(\delta_i^A, \delta_i^I) = 0$$
$$Cov(\delta_i^I, \delta_i^A) - Var(\delta_i^A) = 0$$
$$\Longleftrightarrow$$
$$Cov(\delta_i^A, \delta_i^I) = Var(\delta_i^I) = Var(\delta_i^A)$$
$$\Longleftrightarrow$$
$$\frac{Cov(\delta_i^A, \delta_i^I)}{\sqrt{Var(\delta_i^I)Var(\delta_i^A)}} = \frac{Var(\delta_i^I)}{\sqrt{Var(\delta_i^I)Var(\delta_i^A)}} = 1$$

This implies a linear relationship between $\delta_i^I$ and $\delta_i^A$:

$$\delta_i^I = a + b\delta_i^A$$

with $a$ and $b$ some constant. Plugging back into the initial conditions:

$$Cov(\delta_i^I - \delta_i^A, \delta_i^I) = 0$$
$$Cov(\delta_i^I - \delta_i^A, \delta_i^A) = 0$$
$$\Longleftrightarrow$$
$$Cov(a + b\delta_i^A - \delta_i^A, \delta_i^I) = 0$$
$$Cov(a + b\delta_i^A - \delta_i^A, \delta_i^A) = 0$$
$$\Longleftrightarrow$$
$$Cov(a + (b-1)\delta_i^A, \delta_i^I) = 0$$
$$Cov(a + (b-1)\delta_i^A, \delta_i^A) = 0$$
$$\Longleftrightarrow$$
$$(b-1)Cov(\delta_i^A, \delta_i^I) = 0$$
$$(b-1)Var(\delta_i^A) = 0$$

which implies that $b = 1$ or $\delta_i^I$ and $\delta_i^A$ being constant. In both cases, this implies that the difference $\delta_i^I - \delta_i^A$ is a constant.

# Appendix E  Marginal Treatment Effects

Denote treatment as $D_i$, with $D_i = 0$ meaning that individual $i$ goes to Academia and $D_i = 1$ meaning that $i$ goes to Industry. In what follows, I ignore the subscript $i$ for clarity. The associated potential outcomes in the untreated (Academia) and treated (Industry) states are denoted respectively by $Y_0$ and $Y_1$, and are assumed to be a linear function of observables $\mathbf{X}$ and unobservables $(U_0, U_1)$: $Y_j = \mu_j(\mathbf{X}) + U_j$ for $j = \{0, 1\}$. When the outcome considered is earnings, $U_0 = \delta^A$ and $U_1 = \delta^I$. When the outcome is publications, $U_0 = \theta_p^A$ and $U_1 = \theta_p^I$. The treatment indicator is a function of observables $\mathbf{Z}, \mathbf{X}$ and unobservables $V$:

$$D = \mathbf{1}\{\mu_D(Z, X) > V\} \tag{17}$$

$V$ represents the unobserved characteristics that make individuals less likely to select into treatment and is usually interpreted as unobserved resistance or distaste for treatment. Given Equation 7, $V$ is a function of $\delta_i^I - \delta_i^A$ and $\varepsilon_i^I - \varepsilon_i^A$. Note that a key feature of this model is that the unobserved gain from treatment $(U_1 - U_0)$ can be correlated with $V$.[41] The observed outcome $Y$ is then equal to:

$$Y = (1 - D)Y_0 + DY_1 \tag{18}$$

The model can then be identified either through strong parametric restrictions on $U_0$, $U_1$ and $V$, or by using an instrumental variable $Z$ (which is how I will proceed) which requires weaker assumptions. It is common in the literature to work with the quantiles of $V$ rather than its values, in order to obtain a uniform random variable between 0 and 1. I therefore transform the selection equation as follows:

$$D = \mathbf{1}\{\mu_D(Z, X) > V\} \iff D = \mathbf{1}\{F_V(\mu_D(Z, X)) > F_V(V)\} \iff D = \mathbf{1}\{P(Z, X) > U_D\}$$

with $F_V$ the cumulative distribution function of $V$. Two key variables are $P(Z, X)$, the propensity score and $U_D$, a uniformly distributed random variable between 0 and 1. Remember that $V$ is the unobserved disutility for going into the private sector, so $U_D$ gives us the quantiles of consumer types in terms of their unobserved disutility for treatment. Importantly for what follows, individuals with lower values of $U_D$ are more likely to join the private sector for unobservable reasons while individuals with higher values of $U_D$ are less likely to enter the private sector (more likely to enter Academia) for unobservable reasons. In this setup, the *marginal treatment effect* for individuals with observables $X = x$ and disutility for treatment $U_D = u$ is:

$$\text{MTE}(X = x, U_D = u) = \mathbb{E}[Y_1 - Y_0 | X = x, U_D = u] \tag{19}$$

which makes it clear that treatment effect heterogeneity may result from both observed and unobserved characteristics. For instance, for individuals with $X$ that make them more likely to enter the private sector, I can test whether they have smaller or higher earning gains/publication loss than others. Similarly, considering individuals with lower $U_D$ that make them more likely to enter the private sector, I can test whether they have smaller or higher earning gains/publication loss than others. A common way to present the results is to plot the MTE as a function of $U_D$ for average values of $X$. If the slope of the curve is flat, then $\mathbb{E}[Y_1 - Y_0 | X = x, U_D = u]$ is constant so that there is no heterogeneity in treatment effect across individuals with different propensities for joining the private sector vs Academia (or equivalently for individuals with different distaste for joining the private sector vs Academia). If the slope is upward or downward sloping, there is heterogeneity in treatment effect across individuals with different propensities for joining the private sector vs Academia.

The assumptions needed for the calculation of the MTE are the same as for the IV framework, with in particular a valid exclusion restriction which can be written as $(U_0, U_1, V) \perp\!\!\!\perp Z | X$ where $\perp\!\!\!\perp$ denotes conditional independence. In practice, provided the reader found the instrumental variable previously discussed convincing enough, this means I can use it to estimate MTE. Following the literature, I further assume additive separability between the observed and unobserved components. This allows me to write the MTE as:

$$\text{MTE}(X = x, U_D = u) = x(\beta_1 - \beta_0) + \mathbb{E}[U_1 - U_0 | U_D = u] \tag{20}$$

---

[41] This is obvious for earnings, but it could also be true for publications

It also implies that treatment affect heterogeneity coming from observed characteristics affects the intercept of the MTE curve (plotted as a function of $U_D$) but not its slope (Cornelissen et al., 2016). While this MTE framework is more restrictive than the 2SLS one, it allows me to bring more nuance into treatment heterogeneity and the pattern of selection.

# Appendix F   Calculating the sum of expected earnings

I detail here how I estimate the sum of the expected (discounted) earnings that individual $i$ can expect to earn if she starts her career in sector $j$. Under the full model that allows for transitions across sectors during the career, this sum is equal to:

$$\sum_{t=0}^{T} \rho^t \ \mathbf{E} \ [W_i^{j_t} \mid j_0] = \sum_{t=0}^{T} \rho^t \ \mathbf{E} \ [\overline{W_t^{j_t}} + \delta_i^{j_t} \mid j_0] \tag{21}$$

with $j_0(i)$ the sector joined by $i$ at the beginning of her career and $j_t$ the sector she actually works in at time $t$. Under the simplified model which assumes no transition across sector during the career, this sum simplifies to:

$$\sum_{t=0}^{T} \rho^t \ \mathbf{E} \ [W_{it}^{j}] = \sum_{t=0}^{T} \rho^t \ \mathbf{E} \ [\overline{W_t^{j}} + \delta_i^{j}] = \sum_{t=0}^{T} \rho^t \ \mathbf{E} \ [\overline{W_t^{j}}] + \sum_{t=0}^{T} \rho^t [\delta_i^{j}] \tag{22}$$

My procedure is applicable for both cases but to ease the notations, I refer to the mathematical case of no transition. My procedure is as follows:

- Using all observations in my sample, I estimate a flexible regression of (log) earnings on an indicator for *starting* one's career in Industry, PhD major fixed effects, PhD cohort of graduation fixed effects, years of experience fixed effects and the interactions of these variables. This allows me to get $\mathbf{E} \ [\overline{W_t^{j}}]$ for each individual $i$ who started her career in sector $j$ and each year of experience $t$. I therefore have $T = 40$ observations for each individual as I predict earnings for years of experience 1 to 40. Figure F.21 shows the expected sector-mean earnings by year of experience

- In order to create individual-level variation, I then use the earnings observations I have for each individual (let's assume for exposition that I observe earnings in year $k$) and I calculate the gap between the real observed earnings at time $k$ $W_{it}^{k}$ and $\mathbf{E} \ [\overline{W_t^{k}}]$. Call this gap $g$. For each predicted earnings observation $\mathbf{E} \ [\overline{W_t^{j}}]$, I then use $g$ to shift the predicted curve up (if $g > 0$) or down (if $g < 0$)
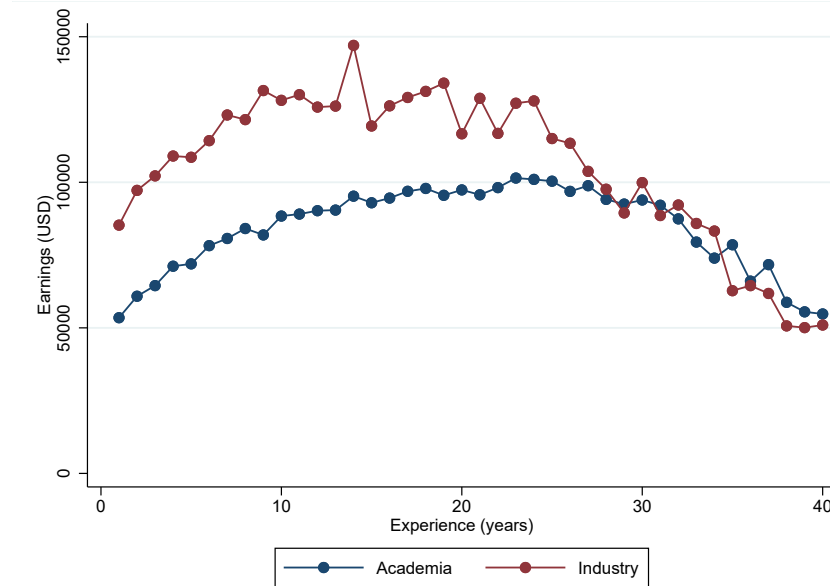


Figure F.21: Predicted earnings by years of experience and sector joined at graduation