## Self-image Bias and Lost Talent

Marciano Siniscalchi Northwestern University

Pietro Veronesi University of Chicago, NBER, and CEPR

August 15, 2023

#### Abstract

We propose an overlapping-generation model wherein researchers belong to two groups, M or F, and established researchers evaluate new researchers. Group imbalance obtains even with group-neutral evaluations and identical productivity distributions. Evaluators' self-image bias and mild between-group heterogeneity in equally productive research characteristics lead the initially dominant group, say M, to promote scholars similar to them. Promoted F-researchers are few and similar to M-researchers, perpetuating imbalance. Consistently with the data, our mechanism also predicts stronger and widening group imbalance in top institutions; higher quality of accepted F-researchers; clustering of M- and Fresearchers across different fields; greater imbalance for seniors than juniors; less credit for F-researchers in co-authored work; and established researchers' false perception that increasing F-representation reduces quality. Policy-wise, mentorship reduces group imbalance, but increases F-group talent loss. Affirmative action reduces both.

Keywords: gender discrimination, self-image bias, affirmative action

JEL codes: A11, J16, J7

<sup>\*</sup> Veronesi acknowledges financial support from the Fama-Miller Center for Research in Finance and by the Center for Research in Security Prices at the University of Chicago Booth School of Business. For useful comments, we thank Gadi Barlevi, Marco Bassetto, Alberto Bisin, Shereen Chaudhry, Hui Chen, Alexia Delfino, James Dow, Nicola Gennaioli, Lars P. Hansen, Christopher Hennessy, Alex Imas, Jessica Jaffers, Seema Jayachandran, Emir Kamenica, Elisabeth Kempf, Nadia Malenko, Lubos Pastor, Jane Risen, Antoinette Schoar, Oleg Urminsky, Laura Veldkamp, Adrien Verdhelan, participants at Behavioral Science workshop, the Finance workshop, and the Economics workshop at the University of Chicago, and seminar participants at the Bocconi University, Chicago Fed, Hitotsubashi University, Imperial College, LBS, MIT, Northwestern University, the 2021 NBER Economics of Culture meetings, the 2021 Adam Smith conference, the 2021 WFA meetings, the 2022 AEA meetings, the 2023 NBER meetings in asset pricing and corporate finance (joint session), and the 2023 MFA meetings. Errors are our own. The views expressed in this paper are our own and do not necessarily reflect the views of our respective employers. We declare that we have no relevant or material financial interests that relate to the research described in this paper. Contact information: Marciano Siniscalchi: marciano@northwestern.edu; Pietro Veronesi: pietro.veronesi@chicagobooth.edu.

# 1. Introduction

The economics profession has long been male-dominated. The Committee on the Status of Women in the Economics Profession (CSWEP), a standing committee of the AEA since 1971,<sup>1</sup> has been regularly documenting the progress of female economists (or lack thereof): see Chevalier (2020). This phenomenon has recently received renewed attention, possibly due to the very slow progress attained in the last 25 years. For instance, while in this time span the fraction of women in undergraduate majors has risen to over 40%, the fraction of women as assistant professors—i.e. the intake for the academic career— has been flat at around 23% since 1994.

This slow progress is puzzling given the numerous initiatives aimed at increasing female representation in the economics profession over the past several decades. Many of these interventions are however informed by existing theories of discrimination, such as taste-based and statistical discrimination, implicit bias, and stereotyping, which we review in Section 7. From this perspective, recent empirical evidence may suggest that efforts to remove such sources of discrimination or bias have only partially succeeded. For instance, Card, DellaVigna, Funk, and Iriberri (2020) documents that acceptance rates for women-authored papers is lower conditional on quality (proxied by future citations); Sarsons (2017) and Sarsons, Gërxhani, Reuben, and Schram (2021) show that female coauthors tend to receive less credit for published papers that are joint with male coauthors; Dupas, Modestino, Niederle, and Wolfers (2021) document a bias against female presenters in economics seminars. Large differences in women's representation exist across fields, however (e.g. Chari and Goldsmith-Pinkham, 2018 and Lundberg and Stearns, 2019), which would then suggest that gender-bias is more prominent in some economic fields than others.

Besides the evidence just cited, Section 2. collects additional stylized facts about the under-representation of women in academia. Specifically, under-representation is higher in research-intensive institutions—but only for tenure-track positions; the gap in under-representation between top institutions and other institutions has been increasing over the last 50 years; and under-representation is widespread in both the US and in Europe—including in Nordic European countries, which are otherwise less gender-biased.

We propose a novel theory that is consistent with this empirical evidence but that does not depend on stereotypes or gender discrimination, whether taste-based, statistical, conscious, or unconscious. In our model, gender imbalance is due to the combination of self-image bias, i.e. the tendency of individuals to place more weight on their own positive attributes when judging others, and mild population heterogeneity in equally-valuable research charac-

 $<sup>^{1}</sup>See https://www.aeaweb.org/about-aea/committees/cswep/about.$ 

teristics. Both assumptions have strong empirical and experimental support, as we discuss below. Our model, which we calibrate to the data, yields several additional predictions, that are also verified in the data. Moreover, it suggests different policy actions to increase female representation and, especially, reduce talent loss.

Specifically, our model features overlapping generations of agents that belong to one of two groups, labelled M and F. A new cohort of young M- and F-researchers appears in every period, in equal proportions. Each researcher is endowed with a type that determines his or her productivity, which we take to mean the probability of producing research that achieves its objectives. The types are randomly and symmetrically distributed in the population of young researchers, so that the distribution of the associated, type-dependent *productivities* is the same in both the M- and F-population. That is, both groups are ex-ante identical in their ability to produce quality research. However, we also assume that some types are slightly more common in the F-group and others symmetrically slightly more common in the M-group. As in the data, we let between-group heterogeneity be far smaller than withingroup heterogeneity. We emphasize that we do not make any assumptions about the *origins* of these distributional differences, which can very well be socially determined (see e.g. Guiso, Monte, Sapienza, and Zingales (2008), Falk, Becker, Dohmen, Enke, Huffman, and Sunde (2018), and Andersen, Ertac, Gneezy, List, and Maximiano (2013)), but only that some mild differences exist, as documented in the empirical evidence discussed below.

We assume that the quality of a young researcher's output is objective and observable. However, each young researcher who has produced quality work must also be evaluated by a randomly matched member of the established population. This evaluator (hereafter, referee) decides whether or not to accept the young researcher as a member of the established population—and thus as a referee of future cohorts. Each referee's perceptions of young researchers' output reflect self-image bias (Lewicki, 1983): evaluators use their own type as yardstick to assess others' research. Importantly, the referees' evaluation is group-neutral: each given referee uses the same criterion to assess young M and F researchers. If the referee's evaluation is positive, the latter becomes a recognized, permanent member of the population; otherwise, he or she leaves the model.

Our key finding is that, for a non-degenerate set of model parameters that we characterize explicitly, the combination of self-image bias and even mild between-group heterogeneity generates a persistent bias that favors young researchers who belong to the group that is initially larger, say the M-group. Moreover, there is no convergence. While researchers from the F-group are also successful, not only are they a minority: they are endogenously selected to be the ones whose types are closer to the ones of the M-researchers; this perpetuates the bias forward. Moreover, this mechanism yields talent loss, as types slightly more common in the F-population are under-represented in the limit. In addition, the same basic logic delivers all the stylized facts collected in Section 2.

We calibrate the model to the data to evaluate whether our mechanism can indeed bring about gender imbalance that is quantitatively relevant even if the distributional differences between M- and F-populations are as small as documented in the data. For this task, we need to be concrete on the notion of types. We identify types with vectors of research characteristics. Examples of such characteristics include research approach (e.g. empirical or theoretical), methodology (e.g. structural versus reduced form), field, topic, type of questions asked, depth vs. breadth, writing style, ties to reality, policy relevance, and so on. In this specification, symmetry means that for every characteristic that is slightly more prevalent in the M-group there is a characteristic that is slightly more prevalent in the F-group; furthermore, all research characteristics are equally valuable: each has the same positive effect on the likelihood of producing quality research. Our calibrated parameters match the very mild heterogeneity in characteristics between male and female populations documented in the psychology literature (see below), as well as the success rate of Economics PhD students to become assistant professors.

Under our calibrated parameters, imbalance occurs and the system converges to about 20% women in academia, which is quite close to the data. In addition, similarly to the data, institutions with higher publication frequency have significantly fewer F-researchers; the gap in F-group under-representation between between top institutions and all institution increases over time; accepted F-researchers have higher objective quality; and there is clustering of M- and F-researchers across fields. Moreover, surviving F-researchers are those whose research characteristics are closer to the ones that are more prevalent among M-researchers. Thus, valuable characteristics that are (mildly) more common among the F-group, but also very common in the M-group, are vastly underrepresented in the steady state. This implies a persistent loss of talent and knowledge, and a sub-optimal steady state.

We assess different policy interventions through the lens of our model. We first investigate the impact of mentorship, and highlight an unintended consequence. We assume that young researchers are matched with random advisors from the set of established researchers. Given self-image bias, advisors advise young researchers to "become like them"—that is, acquire their advisor's type. Young researchers can do so by paying a cost that increases in the distance between their advisor's type their own. We show that, while mentorship may help reduce (but not necessarily eliminate) gender imbalance, it also accelerates the loss of Fgroup characteristics. Intuitively, this is because mentors are drawn from the dominant population, which over-represents M-group characteristics.

Second, we analyze the impact of affirmative-action policies. Specifically, we consider a

mandate to accept the same number of F-researchers as M-researchers each period. Clearly, such policy mechanically leads to gender balance. However, we also find that such a policy additionally ensures that all characteristics are represented in the limit: thus, qualitatively, there is no loss of talent. Intuitively, increasing the F-group representation by mandate also increases heterogeneity in the future pool of referees, which in turn makes it more likely that research characteristics (mildly) more prevalent across F researchers will be accepted.

The Online Appendix analyzes extensions and implications. First, gender imbalance and loss of talent are exacerbated by candidates' career concerns. We endogenize the choice of young researchers to pursue an academic career, or enjoy an outside option. With costly entry, anticipating a bias against their research characteristics, the mass of F-agents who choose academia shrinks over time, and eventually converges to a smaller fraction of "applicants" than their M counterparts. If costs are sufficiently high, characteristics (mildly) more common in the F-group disappear altogether. This intuitive result can help explain why the applications of women to PhD programs in Economics are low to start with. Similar results obtain if hiring institutions bear a cost to hire a young researcher, and receive a payoff from hiring those who later become recognized members of the profession.

Second, we allow for different levels of seniority for established researchers. Senior researchers evaluate junior researchers, and both senior and junior researchers evaluate new entrants. This mimics the career dynamics in academia. Our results about the persistent bias in hiring carry through. Moreover, under suitable parameter configurations, there is a "leaky" pipeline (cf. Chevalier, 2020): senior researchers are even more biased towards characteristics prevalent in the M-group than junior researchers.

Third, while our model does not explicitly allow for co-authorships, its basic force helps explain why female coauthors tend to receive less credit for published papers that are joint with male coauthors (Sarsons, 2017; Sarsons et al., 2021). Intuitively, the referees' population mostly reflects the characteristics of the M-group and thus the positive characteristics of joint research are mostly ascribed to those of the M coauthor.

Our results depend on two main assumptions: mild heterogeneity in research characteristics between M-researchers and F-researchers, and self-image bias, i.e. the tendency of reviewers to use their own research style to judge the importance and worth of others' research output. Both assumptions are grounded in the empirical and experimental literature.

First, there is a considerable body of research studying gender differences in personality traits, preferences, and attitudes. Regarding personality traits, Hyde and Linn (2006) reviews the literature and concludes that medium-sized effects are found for aggression (Cohen's d between 0.40 and 0.60) and activity level in the classroom (d = 0.49)<sup>2</sup>. Similarly, Hyde (2014)

<sup>&</sup>lt;sup>2</sup>Cohen (2013)'s d measures the standardized mean difference between two populations.  $d \approx 0.2$  is

reports the following d statistics of gender differences in the "big-5 personality traits," earlier studied by Costa, Terracciano, and McCrae (2001): among U.S. subjects, there are smallto-moderate differences in neuroticism (d = -0.40), extraversion (d = -0.21), openness (d = 0.30) and agreeableness (-0.31), but a trivial difference in conscientiousness (d =-0.05). Within economics, Croson and Gneezy (2009) provide a review of the experimental literature and find "robust differences in risk preferences, social (other-regarding) preferences, and competitive preferences." Borghans, Golsteyn, Heckman, and Meijers (2009) also find differences in risk aversion, but less so on ambiguity aversion. Dittrich and Leipold (2014) find that women tend to be more patient than men, and Dreber and Johannesson (2008) that males are more likely to lie in order to secure a monetary gain; see also Betz, O'Connell, and Shepard (1989). Finally, Andre and Falk (2021) survey nearly 10,000 economists' opinion about the current state and preferred direction of economic research. They find that female scholars are significantly more likely to emphasize multidisciplinarity, disruptive research, and policy relevance (cf. Table 3.)

The second important assumption of our model is researchers' self-image bias. The psychological literature on self-image bias (Lewicki, 1983) suggests that, when evaluating others, individuals tend to place more weight on positive attributes that they themselves possess (or believe they possess). Hill, Smith, and Hoffman (1988) show that this is true in particular when subjects are asked to select a partner in a competitive game. Dunning, Perie, and Story (1991) argue that a similar principle is at work when judging social categories by means of prototypes (e.g., what makes a good economist?): "people may expect the 'ideal instantiation' of a desirable social category to resemble the self in its strengths and idiosyncracies" (p. 958). Story and Dunning (1998) document a "rational" source for selfimage bias and self-serving prototypes: in their experiment, "those who received success feedback came to perceive a stronger relationship between 'what they had' and 'what it takes to succeed' than did those who received failure feedback" (p. 513). Translated to our environment, established researchers view their personal success in research as evidence that their own research characteristics are the right ones to produce quality research that, in addition, are valuable to society. Hence, they use the same characteristics to evaluate the research of others.

Our assumption that referees accepts young researchers who are similar to them can also be due to referees' preferences (e.g. theorists like theorists, and empiricists like empiricists) rather than self-image bias. However, this interpretation must be subject to two caveats to fit our model. First, referees' preferences do *not* take group membership into account; thus, even this "homophily" interpretation of our model differs from Becker's taste-based

considered "small" and  $d\approx\!\!0.5$  is considered "medium."





Source: CSWEP Report, 2023.

theory of discrimination. Moreover, in this interpretation, referees do not value heterogeneity (e.g., theorists derive no benefit from interacting with empiricists, and conversely), nor the candidate's objective productivity. This strikes us as extreme.

# 2. Stylized Facts

We begin with discussing a collection of stylized facts about the percentage of women students and faculty in U.S. and elsewhere.

The top panel of Figure 1 shows that between 1994 and 2022, the fraction of women in undergraduate economics majors increased to over 40% in the top-20 schools. During the same period, the fraction of women PhD students has been flat at around 30%, except for the last few years, which saw a marked increase. The bottom panel is though the crux of the issue: Among assistant professors—i.e. the intake for the academic career—the fraction of women in top-20 schools has been flat at below 25% since 1994. In other words, as discussed in the introduction, as far as representation among assistant professors is concerned, there has been virtually no progress in the top 20-schools over a time span of nearly 30 years.

Even more striking, the top panel of Figure 2 shows the difference between schools with and without a PhD program, and between top institutions and all institutions. Universities without PhD program hire over 40% of female tenure-track faculty, while universities with PhD program hire just over 30% of women assistant professor. In addition, the top-20 schools hire only 23% and the top-10 only 21%. In sharp contrast, the share of women among teaching faculty is quite uniform across schools, at around 35%. The difference between research schools and non-research schools suggests that women under-representation is associated with research intensity. Indeed, the fraction of women among teaching faculty is now close to the fraction of women in undergraduate economics majors (about 40%).

The CSWEP report refers to top 10 and top 20 universities but does not describe how this ranking is determined. The economics profession typically associates rank with research output (Chevalier (2020)). With this in mind, Table 1 shows the results of a regression of the share of women faculty at an institution on a measure of research intensity at that institution, obtained from Tilburg's ranking of U.S. institutions in terms of research output.<sup>3</sup> While there is considerable noise, the results suggest that an institution's research intensity is significantly negatively correlated with the percentage of women faculty in that institution. The table also shows that when we run the same regression on non-tenure track faculty, there is no relation whatsoever with the publication intensity. The negative relation is only visible for research faculty, consistently with the aggregate data displayed in Figure 2.

Focusing again on universities with PhD programs, the bottom panel of Figure 2 shows the striking dynamics of the percentage of women assistant professors between top departments and all departments since 1974. While women representation was low in the late 1970s, it was about the same percentage between top institutions and all institutions. However, over time, the percentage of women has increased, but *less* so for top institutions. The gap between top and all institutions (grey line at the bottom of the graph) has been increasing

<sup>&</sup>lt;sup>3</sup>The rankings are available at https://econtop.uvt.nl/rankingsandbox.php. We run the "sandbox" in two ways. First, we use the Tilburg ranking by resetting the journals to be used, selecting USA as country, choosing to weight the results by journals impact factors, and list the top 300 universities. This is the "Tilburg ranking" in the table. Second, we selected the top 5 journals and obtained the new ranking.



### Figure 2: Percentage of Women in Academia across Institutions

Percent of Women Faculty across Types of Schools

Notes: The top panel plots the average percentage of women in academia across types of institutions and types of academic positions. Data are from the 2022 CSWEP Report and averages are over the 2018-2022 sample. The bottom panel reports the fraction of women assistant professors in top institutions vs. all institutions with a PhD program from 1974 to 2022. Data from 1974 to 1993 were extracted from Figures 2 and 3 of the 1994 CSWEP report available at https://www.aeaweb.org/content/file?id=682, while data from 1994 onward are from the 2022 CSWEP report. The figure also reports 4<sup>th</sup>-order polynomial trend lines as well as the difference between the two lines at the bottom.

in the last 50 years, which, once again, appears to indicate that research at top institutions is a critical component of the story.  $^4$ 

<sup>&</sup>lt;sup>4</sup>We caveat these results, however, by noticing that from 1974 to 1993, CSWEP defines "top institution"

	Ranking by Journal Impact Factor			Ranking by Top 5 Econ Journals		
	Assistant	All Tenure	Non-Tenure	Assistant	All Tenure	Non-Tenure
	Professors	Track Faculty	Track Faculty	Professors	Track Faculty	Track Faculty
α	34.75	23.79	37.72	33.40	22.23	38.27
(S.E.)	(2.02)	(0.92)	(2.71)	(1.75)	(0.79)	(2.50)
(p-value)	(0.000)	(0.000)	(0.000)	(0.000)	(0.000)	(0.000)
eta	-0.0470	-0.0330	0.0003	-0.1050	-0.0670	0.0029
(S.E.)	(0.0169)	(0.0077)	(0.0223)	(0.0349)	(0.0156)	(0.0493)
(p-value)	(0.006)	(0.000)	(0.988)	(0.003)	(0.000)	(0.9526)
$R^2$ (%)	6.32	13.53	0.00	9.06	16.79	0.00
Ν	117	117	113	93	93	91

Table 1: Percentage of Women and Institution Research Intensity

Notes: This table shows the results of the regression

(% women faculty)<sub>j</sub> =  $\alpha + \beta \times (\text{publication score})_j + \varepsilon_j$ 

where (% women faculty)<sub>j</sub> is the average percentage of women faculty as assistant professor (columns 1 and 4), as tenure-track faculty (columns 2 and 5), as non-tenure track faculty (columns 3 and 6) in school j; and (publication score)<sub>j</sub> is the research score of school j. The latter is obtained from Tilburg's ranking of U.S. institutions in terms of research output. We considered two cases: Columns 1 – 3 considers the base case of Tilburg rankings, using the weighting of journals by impact factor. Columns 4 – 6 consider research output only based on publications on top 5 economics journals (AER, Econometrica, QJE, JPE, Review of Economic Studies). Both the average percentage of female faculty and the publication score are computed over the latest sample available, namely 2016-2020. Data for the % of women faculty are from the disaggregated restricted data underlying CSWEP report and available at https://www.icpsr.umich.edu/web/ICPSR/studies/37118/. Missing data are deleted from the sample.

The possibility that explicit or implicit gender discrimination may not be the full story underlying the lack of women representation in economics is also highlighted by Figure 3. This figures uses data from Auriol et al. (2022) and correlates the fraction of women in senior academic positions in economics and business in U.S., Canada and European countries, with their respective values of the global gender gap index obtained from the world economic forum. A higher global gender gap index indicates *lower* differences between men and women on several dimensions (see https://www3.weforum.org/docs/WEF\_GGGR\_2020.pdf.) As in Auriol et al. (2022), we find a positive relation between the global gender gap index and the percentage of women in academia. However, focusing now on magnitudes, even in Nordic European countries, where the global gender gap index is maximized, the percentage of women in academic positions still hovers only around 25%, which is not too different from the US (21%) in this dataset.<sup>5</sup>

as those above the median according to the National Research Council rankings, while in the 1994 to 2021 data, CSWEP defines "top institutions" as the top 20 schools. Still, the gap is visibly increasing also just in the latter sample.

<sup>&</sup>lt;sup>5</sup>This figure is similar to Figure 4 in Auriol et al., 2022, which plots the rank-order of countries with





Source: Data of women representation are from the main dataset of Auriol, Friebel, Weinberger, and Wilhem, 2022, focusing on U.S., Canada, and European countries. The global gender gap index of each country is from the world economic forum (See Table 1 at https://www3.weforum.org/docs/WEF\_GGGR\_2020.pdf).

Finally, Figure 4 illustrates the quality and field interests of women in academia. The top panel reports Figure 4(a) from Card et al. (2020), which shows that "[a]t nearly each referee recommendation, female-authored papers have higher citations than male-authored papers, with a 20 log point average difference. This suggests that papers by all-female authors are held to a higher bar by the referees." (Card et al. (2020), page 296). The bottom panel reproduces Figure 5 in Lundberg and Stearns, 2019, which reports the difference between share of women and share of men in particular fields of economics. The data used are from the annual list of Doctoral Dissertations in Economics, from 1991–2017. As it can be seen, women are relatively more frequent in Labor/Public Economics (green line) and men relatively more frequent in Macro/Finance.<sup>6</sup> The striking pattern, however, is that over the 25-year sample, there is no variation at all on the relative percentages, as if the "system" has converged, consistently with the bottom panel of Figure 1.

higher percentage of female faculty against the rank-order of countries in terms of global gender gap index. Figure 3 makes explicit the actual percentage of women across institutions to highlight that in the majority of countries in US, Canada, and Europe, such percentage is still below 30%.

<sup>&</sup>lt;sup>6</sup>The figure shows "relative frequency" and not absolute frequency. Even in labor/public finance, the majority of faculty is still male.



Source: The top panel reports Figure 4 (a) of Card, Della Vigna, Funk, and Iriberri, "Are referees and editors in economics gender neutral?" *Quarterly Journal of Economics*, 135, 2020, which plots the average (residualized) citation rate of non-desk rejected papers across types of referee recommendation. The bottom panel reports Figure 5 in Lundberg and Stearns, 2019.

# 3. Model

We consider an overlapping-generations model in which unit masses of two groups of young researchers, the *M*-group and *F*-group, appear at discrete times t = 1, 2, ... Each researcher  $i \in M \cup F$  is endowed with a *type*, drawn from a finite set  $\Theta$  and distributed heterogeneously across *M* and *F* researchers. While systematic, these distributional differences may well be small. Assumption 1 below imposes a notion of symmetry among types that ensures that the two groups are equally productive ex-ante. Research output fully reflects the researcher's type; in fact, we assume that the characteristics of a paper written by a researcher of type  $\theta$  are  $\theta$  itself. Let  $p^{\theta,m}$  and  $p^{\theta,f}$  denote the fraction of types  $\theta$  in the *M*- and *F*-populations of young researchers, respectively. These fractions are strictly positive for all  $\theta$ .

The researcher of type  $\theta$  has probability  $\gamma^{\theta}$  of producing quality research. "Quality" research is one that achieves its stated goals—estimating a parameter of interest, establishing a causal effect, documenting a phenomenon experimentally, or proving a theorem. We assume that whether a research paper achieves its goals is observable and can be objectively determined; this may involve, for instance, checking a formal argument regarding a theoretical claim or the application of a statistical procedure, evaluating an experimental procedure for possible biases or ambiguities, or ensuring that the formal results are clearly explained and interpreted, and that the contribution is correctly placed within its literature.

We assume that, while there are (typically small) differences in the distribution of types in the M- and F-populations, there are no differences in the corresponding distribution of quality. Specifically, we assume that type distributions in the M and F populations are symmetric in the following sense:

Assumption 1. Symmetric Distribution of Quality: for every type  $\theta$ , there is a corresponding type  $\theta'$  that has the same quality as  $\theta$ , and that has the same mass in F (resp. M) that type  $\theta$  has in M (resp. F). Formally, there exists a function  $\sigma : \Theta \to \Theta$  such that

- (i) for every type  $\theta \in \Theta$ , the corresponding type  $\theta' = \sigma(\theta)$  satisfies  $\gamma^{\theta} = \gamma^{\theta'}$  and  $p^{\theta,m} = p^{\theta',f}$ ; furthermore,
- (ii) for every type  $\theta \in \Theta$ ,  $\sigma(\sigma(\theta)) = \theta$ .<sup>7</sup>

Note that, by (i) and (ii), it is also the case that, for every  $\theta$ , the corresponding  $\theta' = \sigma(\theta)$  satisfies  $p^{\theta,f} = p^{\theta',m}$ , as  $p^{\theta,f} = p^{\sigma(\sigma(\theta)),f} = p^{\sigma(\theta),m}$ . That is, there is full symmetry in the distributions of  $\theta$ 's across the M and F population.

Finally, to model heterogeneity in distributional frequencies of types across M and F population, we assume the following:

Assumption 2. Heterogeneity in *M*- and *F*-Distribution: for some type  $\theta \in \Theta$ ,  $p^{\theta,m} > p^{\theta,f}$ ; by (i), this implies that  $p^{\theta',f} > p^{\theta',m}$  for  $\theta' = \sigma(\theta)$ .

That is, type  $\theta$  is more frequent in M population, while the corresponding type  $\theta' = \sigma(\theta)$ is more frequent in the F population. We let  $p^g = (p^{\theta,g})_{\theta \in \Theta}$  for g = f, m.

<sup>&</sup>lt;sup>7</sup>That is,  $\sigma$  is self-inverse, which implies that it is a bijection.

### 3.1. Objective Refereeing

This section studies a benchmark system where the evaluation by established scholars is objective and only certifies whether the research is of sufficient quality or not. Since each young scholar of type  $\theta$  produces quality research with probability  $\gamma^{\theta}$ , this is also the probability with which the research is "accepted" by referees.

For every type  $\theta \in \Theta$ , let  $a_t^{\theta,m}$  and  $a_t^{\theta,f}$  denote the mass of young researchers of group M and, respectively, group F of type  $\theta$  that produce quality research and are thus "accepted" at the end of period t:

$$a_t^{\theta,g} = \gamma^\theta \cdot p^{\theta,g}, \quad g \in \{f,m\}.$$
(1)

Denote the total mass of accepted young researchers by  $a_t = \sum_{\theta \in \Theta} \sum_{g \in \{f,m\}} a_t^{\theta,g}$ .

Denote by  $\lambda_t^{\theta,g}$  the mass of established researchers of type  $\theta$  and group g at time t. We normalize the initial mass of all established researchers to one:  $\sum_{\theta} \sum_{g} \lambda_0^{\theta,g} = 1.^8$  In order to keep the mass of referees constant, we assume that each young agent whose research is accepted replaces a randomly drawn established one. This is not necessary for the results but keeps the analysis balanced. This assumption is also geared towards maximizing the impact of young researchers on the evolution of the system, and thus give the best chance for the system to converge to group balance.<sup>9</sup> The resulting dynamic is then described by the following equation:

$$\lambda_t^{\theta,g} = (1 - a_t)\lambda_{t-1}^{\theta,g} + a_t^{\theta,g}, \quad g \in \{f,m\}.$$
 (2)

We then obtain the following proposition:

**Proposition 1** In the benchmark model with objective refereeing, regardless of the composition  $(\lambda_0^{\theta,m}, \lambda_0^{\theta,f})_{\theta \in \Theta}$  of the initial population of established researchers, we have

$$\lambda_t^{\theta,m} \to \frac{\gamma^{\theta} p^{\theta,m}}{a}, \quad \lambda_t^{\theta,f} \to \frac{\gamma^{\theta} p^{\theta,f}}{a}, \quad and \quad \frac{\sum_{\theta} \lambda_t^{\theta,m}}{\sum_{\theta} \lambda_t^{\theta,f}} \to 1.$$

where  $a = \sum_{\theta} \gamma^{\theta} \left( p^{\theta, f} + p^{\theta, m} \right).$ 

We prove all Propositions containing our main results in the Appendix. For other Propositions, we provide a proof sketch in the Appendix, and a full proof in the Online Appendix. We prove all Corollaries in the Online Appendix.

<sup>&</sup>lt;sup>8</sup> The fact that the total mass of established scholars (a stock) equals the mass of young M and F researchers (flows) is of course not realistic, but immaterial for our analysis. Normalizing the stock of established researchers to any positive number K yields the same predictions.

<sup>&</sup>lt;sup>9</sup>We also considered a similar model with a fix retirement rate of existing researchers to be replaced by cohorts of hired young researchers. The results are similar. The assumption in the text has one less parameter and it is more favorable to an eventual convergence to group balance.

In our benchmark model with objective reference in a initial conditions have no long-run effects. In addition, the system always converges to equal shares of M and F established researchers, and the limiting type distribution is fully characterized by the probability of producing quality research  $(\gamma^{\theta})$  and the relative frequency of each type in the population of young researchers  $(p^{\theta,m} \text{ and } p^{\theta,f})$ . Given the symmetry of the model, this is intuitive.

### 3.2. Refereeing with Self-Image Bias

Our main model differs from the benchmark in Section 3.1. in that established researchers (referees) not only evaluate young researchers on whether their research is of sufficient quality (as in previous section), but they also use their personal research styles to guide their subjective judgement as to the "importance" or "relevance" of the candidate's output. Specifically, each young researcher  $i \in M \cup F$  of type  $\theta^i$  is now randomly matched to a referee r, who uses his or her own characteristics  $\theta^r$  to evaluate agent i's work. Importantly, evaluation is anonymous and group-blind: it depends solely upon referee r's own type  $\theta^r$  and the characteristics of researcher i's output, which by assumption coincides with his of her type  $\theta^i$ .

Consistently with self-image bias, referee r rejects applicants whose type is far from his/her own set of characteristics. For tractability, we make in fact the following stark assumption (we relax it in the on-line appendix.):

Assumption 3. Self-Image Bias: referee r evaluates young agent i's research favorably if and only if  $\theta^r = \theta^i$ .

If agent i's output is favorably evaluated, i is accepted as an established researcher, and will serve as referee for future cohorts of young researchers.

Let  $\lambda_t^{\theta} = \lambda_t^{\theta,f} + \lambda_t^{\theta,m}$  be the total mass of established researchers of type  $\theta$  at time t; also let  $\lambda_t = (\lambda^{\theta})_{\theta \in \Theta}$ . Retaining the notation of Section 3.1., the dynamics for the mass of young researchers of type  $\theta$  and group g that are accepted in round t is

$$a_t^{\theta,g} = \gamma^{\theta} \cdot \lambda_{t-1}^{\theta} \cdot p^{\theta,g}.$$
(3)

Importantly, whether a young researcher is accepted or not depends solely on her type  $\theta$ , and not also on her group g. As in Equation (2), the total mass of established researchers of type  $\theta$  and group g is given by

$$\lambda_t^{\theta,g} = \lambda_{t-1}^{\theta,g} \left(1 - a_t\right) + a_t^{\theta,g} \tag{4}$$

where as above  $a_t = \sum_{\theta} \sum_{g} a_t^{\theta,g}$ . Equations (3) and (4) indicate that there are two forces at play. On one hand, the distribution of incumbent types impacts which research characteristics are likely to be positively evaluated by referees. On the other hand, even among

incumbents, types that are more likely to produce quality research tend to be more prevalent. As we shall demonstrate, the interplay of these two forces determines whether the system ultimately attains the first-best outcome in Section 3.1., or if instead an inefficient outcome, characterized by group imbalance, is reached.

### 3.3. Type Dynamics

We begin by studying the evolution of the mass of each type in the population. Given our assumption that accepted researchers replace randomly drawn existing ones, the following assumption ensures that the mass of each type remains positive:

Assumption 4: boundedness. For every  $\theta \in \Theta$ ,  $\gamma^{\theta}(p^{\theta,m} + p^{\theta,f}) \leq 1$ .

The following proposition—our first main result—characterizes the types that survive (i.e. have positive mass) in the limit as  $t \to \infty$ . All other types vanish over time.

#### Proposition 2 Let

$$\Theta^{\max} = \left\{ \theta \in \Theta \text{ such that } \lambda_0^{\theta} > 0 \text{ and } \theta \in \arg\max_{\theta' \in \Theta} \gamma^{\theta'}(p^{\theta',m} + p^{\theta',f}) \right\}$$
(5)

then:

(i) only 
$$\theta \in \Theta^{\max}$$
 survive in the limit as  $t \to \infty$ ;

(ii)  $\Theta^{\max}$  preserves symmetry across types: if  $\theta \in \Theta^{\max}$  and  $\lambda_0^{\sigma(\theta)} > 0$ , then  $\sigma(\theta) \in \Theta^{\max}$ .

The proposition shows that the only types  $\theta$  that survive in the limit are those that had positive mass at time 0,  $\lambda_0^{\theta} > 0$ , and that maximize the product  $\gamma^{\theta}(p^{\theta,m} + p^{\theta,f})$ . Intuitively, such types  $\theta$  are both objectively good (high  $\gamma^{\theta}$ ) and frequent in the young population (high  $(p^{\theta,m} + p^{\theta,f})$ ). A type  $\theta$  that does not have a high frequency in the young population, or it is consistently unproductive, will unlikely be part of the reference population. Thus, self-image bias will act against it, as, in the limit, no reference will view his/her research favorably.

Part (ii) states that the assumed symmetry of types in the M and F populations (Assumption 1) is preserved in the limit, unless specific types are not represented in the initial population  $\lambda_0$ . This is the case because, for any type  $\theta$ ,  $\gamma^{\theta}(p^{\theta,m}+p^{\theta,f}) = \gamma^{\sigma(\theta)}(p^{\sigma(\theta),f}+p^{\sigma(\theta),m})$ .

Our second main result characterizes the limiting distribution of types:

**Proposition 3** Let  $\overline{\lambda}^{\theta,g}$  denote the limit of type  $\theta$ 's mass  $\lambda_t^{\theta,g}$ , for  $g \in f, m$ , as  $t \to \infty$ , and let  $\overline{\lambda}^{\theta} = \overline{\lambda}^{\theta,m} + \overline{\lambda}^{\theta,f}$ . Then, for all  $\theta \in \Theta^{\max}$ :

(i) The limiting mass of  $\theta$  is

$$\bar{\lambda}^{\theta} = \frac{\lambda_0^{\theta}}{\sum_{\theta' \in \Theta^{\max}} \lambda_0^{\theta'}} \tag{6}$$

(ii) The limiting mass of  $\theta$  in its f and m groups are:

$$\bar{\lambda}^{\theta,f} = \frac{\lambda_0^{\theta} \gamma^{\theta} p^{\theta,f}}{\sum_{\theta' \in \Theta^{\max}} \lambda_0^{\theta'} \gamma^{\theta'} (p^{\theta',f} + p^{\theta',m})}; \quad \bar{\lambda}^{\theta,m} = \frac{\lambda_0^{\theta} \gamma^{\theta} p^{\theta,m}}{\sum_{\theta' \in \Theta^{\max}} \lambda_0^{\theta'} \gamma^{\theta'} (p^{\theta',m} + p^{\theta',f})} \quad (7)$$

Point (i) shows that  $\theta \in \Theta^{\max}$  is more frequent in the limit than another type  $\theta' \in \Theta^{\max}$  if it was relatively more frequent in the initial distribution at time 0. Intuitively, types that are more frequent in the initial distribution have more referees who "like" those same types, and thus they are more likely to self-replicate and their larger mass will persist in the limit.

In sharp contrast with Proposition 1 under objective reference point (ii) of Proposition 3 shows that the final mass of types in M and F populations does depend on the initial distribution  $\lambda_0$ . In particular, surviving types  $\theta \in \Theta^{\max}$  have a higher relative mass  $\bar{\lambda}^{\theta,g}$  in group g = m, f if they are more frequent initially (higher  $\lambda_0^{\theta}$ ), they have higher productivity  $\gamma^{\theta}$ , and they are more frequent in the distribution  $p^{\theta,g}$ .

The ex-ante symmetry and heterogeneity of the distributions  $p^{\cdot,m}$  and  $p^{\cdot,f}$  in Assumptions 1 and 2 then lead to our main result for this section, namely, that for symmetric types  $\theta$ and  $\theta' = \sigma(\theta)$ , the relative initial distribution of types determines the final relative mass of surviving types across the two groups.

The following definition is convenient.

**Definition 1** Types  $\theta$  and  $\theta'$  are distinct symmetric elements of  $\Theta^{\max}$  if (i)  $\theta, \theta' = \sigma(\theta) \in \Theta^{\max}$  and (ii)  $p^{\theta,m} \neq p^{\theta',m}$  (or, equivalently,  $p^{\theta,f} \neq p^{\theta',f}$ ).

We then have the following:

**Corollary 1** Let  $\theta$ ,  $\theta'$  be distinct symmetric elements of  $\Theta^{\max}$  with  $\lambda_0^{\theta} > \lambda_0^{\theta'}$  and  $p^{\theta,m} = p^{\theta',f} > p^{\theta,f} = p^{\theta',m}$ . Then

$$\bar{\lambda}^{\theta,m} + \bar{\lambda}^{\theta',m} > \bar{\lambda}^{\theta,f} + \bar{\lambda}^{\theta',f}.$$
(8)

This corollary is a key result of the paper. The initial condition  $\lambda_0^{\theta} > \lambda_0^{\theta'}$  is group independent. Furthermore, from Assumption 1, type symmetry implies that  $\gamma^{\theta} = \gamma^{\theta'}$  (both types are equally productive) and  $p^{\theta,m} + p^{\theta',m} = p^{\theta,f} + p^{\theta',f}$  (the flow of young researchers whose type is either  $\theta$  or  $\theta'$  is the same in the *M*- and *F*-group). Yet, the fact that type  $\theta$ , which happens to be more frequent in the *M*-population  $(p^{\theta,m} > p^{\theta',m})$ , is also more frequent in the overall population of referees at time 0 ( $\lambda_0^{\theta} > \lambda_0^{\theta'}$ ) implies that, in the limit, the mass of *M*-researchers of type  $\theta$  or  $\theta'$  is greater than the mass of *F*-researchers having the same types. Formally, this follows from part (ii) of Proposition 3. Intuitively, *M*-researchers are more likely to be of type  $\theta$  than type  $\theta'$ , and they are also more likely to be matched with a referee of type  $\theta$  than of type  $\theta'$ , regardless of the referee's group; this "positive assortative matching" of referees and *M*-researchers yields a high probability of being accepted. By way of contrast, *F*-researchers are more likely to be of type  $\theta'$  than  $\theta$ , but they are still more likely to be matched with type- $\theta$  referees; thus, "negative assortative matching" leads to a lower probability of being accepted. This perpetuates the bias forward.

### **3.4.** Model Predictions

In this section, we discuss the model's qualitative predictions under the assumption that the initial distribution of referees,  $\lambda_0$ , is skewed towards the types of the *M*-population. Specifically, we assume throughout  $\lambda_0 = p^m$ , a natural assumption to investigate the impact of a population of referees initially dominated by the *M*-group.

#### 3.4.1. Group Imbalance in the Limit

In this subsection we expand on the group imbalance discussed in Corollary 1 above for only two symmetric types  $\theta$  and  $\theta'$ . In particular, let

$$\bar{\Lambda}^g = \sum_{\theta \in \Theta^{\max}} \bar{\lambda}^{\theta, g}, \quad \text{for} \quad g = m, f$$

denote the total limiting mass of researcher in group g. Then, the following result follows:

### **Proposition 4** Let $\lambda_0 = p^m$ .

(a) If  $\Theta^{\max}$  contains distinct symmetric elements, then

$$\bar{\Lambda}^m = 1 - \bar{\Lambda}^f > 0.5. \tag{9}$$

(b) If  $\Theta^{\max}$  contains no distinct symmetric elements, then  $\bar{\Lambda}^m = \bar{\Lambda}^f = \frac{1}{2}$ .

Point (a) of this proposition shows that even with ex-ante symmetric types between M and F (Assumption 1), self-image bias leads to imbalance in the limit when the initial population of referees is initially dominated by the M-population. This is consistent with the bottom panel of Figure 1, which shows that the percentage of women in top 20 economics departments has been constant at around 25% in nearly 30 years. We remark that the condition of part (a) is satisfied if a type  $\theta \in \Theta^{\max}$  has  $p^{\theta,m} \neq p^{\theta,f}$ , and  $\theta' = \sigma(\theta)$  has  $\lambda_0^{\theta'} > 0$ . The calibration section provides parametric conditions for this to be satisfied.

#### **3.4.2.** Higher "Bar" for *F*-researchers

If the initial population of referees' types is skewed towards the M-population, a basic force in our model implies that young researchers from the F-group are, in a sense, held to a higher standard. Recall that in our model all researchers are ex-ante identical in terms of productivity. However, because of self-image bias, the acceptance rate of young researchers from M-population is higher than the one from the F-population.

**Proposition 5** Let  $\lambda_0 = p^m$ . Then for every t > 0, the acceptance rate of *M*-researchers of "quality"  $\overline{\gamma}$  is higher than the one of *F*-researchers of the same quality:

$$\sum_{\theta:\gamma^{\theta}=\overline{\gamma}} a_t^{\theta,m} \ge \sum_{\theta:\gamma^{\theta}=\overline{\gamma}} a_t^{\theta,f}$$

and the inequality is strict if there is a type  $\theta \in \Theta$  for which  $p^{\theta,m} > p^{\sigma(\theta),m}$ .

This result is in line with the evidence in the top panel of Figure 4 from Card et al. (2020) that, conditional on quality (proxied by citations post-publication) women-authored papers tend to be accepted less frequently than men's.<sup>10</sup>

The above result suggests that, on average, the objective quality of accepted F-researchers should be higher than that of M-researchers. We have conducted extensive simulations in the parametric model of Section 4., and indeed this appears to be the case. We can actually prove this formally in the setting of the present section when  $\Theta$  has only four types, namely, a "good" type  $\theta_1$  (high  $\gamma^{\theta_1}$ ), a "bad" type  $\theta_0$  (low  $\gamma^{\theta_0}$ ), and two intermediate and symmetric types  $\theta$ ,  $\theta'$  (intermediate and identical  $\gamma^{\theta} = \gamma^{\theta'}$  with  $p^{\theta,m} = p^{\theta',f} > p^{\theta',m} = p^{\theta,f}$ ).

**Proposition 6** Let  $\Theta = \{\theta_0, \theta, \theta', \theta_1\}$ , with  $\gamma^{\theta_0} < \gamma^{\theta_1}, \gamma^{\theta} = \gamma^{\theta'} \leq \frac{\gamma^{\theta_0} + \gamma^{\theta_1}}{2}, p^{\theta,m} = p^{\theta',f} > p^{\theta',m} = p^{\theta,f}$ , and  $p^{\theta_0,m} = p^{\theta_0,f} = p^{\theta_1,m} = p^{\theta_1,f}$ . Finally, let  $\lambda_0 = p^m$ . Then:

(i) The average quality of accepted *F*-researchers is higher than the one of accepted *M*-researchers:

$$E[\gamma|f, accepted] = \sum_{\hat{\theta}} \gamma^{\hat{\theta}} \ w_t^{\hat{\theta}, f} > \sum_{\hat{\theta}} \gamma^{\hat{\theta}} \ w_t^{\hat{\theta}, m} = E[\gamma|m, accepted]$$
(10)

where

$$w_t^{\hat{\theta},g} = \frac{a_t^{\hat{\theta},g}}{\sum_{\hat{\theta}'} a_t^{\hat{\theta}',g}}$$

<sup>&</sup>lt;sup>10</sup> Card et al. (2020) also show that, unconditionally, men- and women-authored papers are equally likely to be accepted. The model in this section does not generate this finding: summing over all types  $\theta$  in the displayed equation of Proposition 5, one readily sees that young M researchers are more likely to be accepted on average. The model with endogenous choice in Section A1.1. yields more uniform unconditional acceptance across genders, and fewer female acceptance overall due to self-selection.

(ii) As  $t \to \infty$  the average quality of both F and M converges to either  $\gamma^{\theta} = \gamma^{\theta'}$  if  $\Theta^{\max}$  if  $\theta, \ \theta' \in \Theta^{\max}$ , or  $\gamma^{\theta_1}$  otherwise.

#### 3.4.3. Talent Loss and Clustering

One further implication of our model is that with self-image bias types that are more common in the F-group are under-represented in the limit.

**Corollary 2** Assume  $\lambda_0 = p^m$ . Let  $\theta$ ,  $\theta'$  be distinct symmetric elements of  $\Theta^{\max}$  with  $\lambda_0^{\theta} > \lambda_0^{\theta'}$  and  $p^{\theta,m} = p^{\theta',f} > p^{\theta,f} = p^{\theta',m}$ . Then in the limit

$$\frac{\bar{\lambda}^{\theta,m}}{\bar{\lambda}^{\theta,m} + \bar{\lambda}^{\theta',m}} > 0.5 = \frac{\bar{\lambda}^{\theta',f}}{\bar{\lambda}^{\theta,f} + \bar{\lambda}^{\theta',f}} = \frac{\bar{\lambda}^{\theta,f}}{\bar{\lambda}^{\theta,f} + \bar{\lambda}^{\theta',f}} > \frac{\bar{\lambda}^{\theta',m}}{\bar{\lambda}^{\theta,m} + \bar{\lambda}^{\theta',m}}$$
(11)

In the limiting distribution, the majority of established M researchers are of type  $\theta$ . However, the established F-population has the same fraction of type  $\theta$  as type  $\theta'$ . This result is in stark contrast with  $\theta'$  being the prevalent type in each cohort of young F-researchers (assumption  $p^{\theta',f} > p^{\theta,f}$ ). The result is moreover independent of the actual values of  $p^{\theta,g}$ 's. That is, we could have a flow of young researchers with e.g.  $p^{\theta',f} = 0.9 > 0.1 = p^{\theta,f}$  (and symmetrically,  $p^{\theta,m} = 0.9 > 0.1 = p^{\theta',m}$ ), and yet both types  $\theta$  and  $\theta'$  would be equally represented in the F-population in the limit. The selection mechanism makes the type most prevalent among M-researchers,  $\theta$ , be a frequent type in the established F-researchers (50% of the time). This implies that F-group research types are underrepresented in the limit.

Self-image bias also implies clustering of different types within the two groups. In particular, established researchers of type  $\theta$  are more likely to be from the M group; in contrast, type- $\theta'$  researchers are mostly going to be from the F group.

**Corollary 3** Assume  $\lambda_0 = p^m$ . Let  $\theta$ ,  $\theta'$  be distinct symmetric elements of  $\Theta^{\max}$  with  $\lambda_0^{\theta} > \lambda_0^{\theta'}$  and  $p^{\theta,m} = p^{\theta',f} > p^{\theta,f} = p^{\theta',m}$ . Then in the limit, *M*-researchers are relatively more frequent as type  $\theta$  and *F*-researchers are relatively more frequent as type  $\theta'$ :

$$\frac{\bar{\lambda}^{\theta,m}}{\bar{\lambda}^{\theta,m} + \bar{\lambda}^{\theta,f}} = \frac{\bar{\lambda}^{\theta',f}}{\bar{\lambda}^{\theta',m} + \bar{\lambda}^{\theta',f}} > 0.5.$$
(12)

If, as seems plausible, some types may be better suited for some specific research fields than in others, this result implies that the two groups will be differently represented across fields. This is qualitatively consistent with the evidence in the bottom panel of Figure 4 from Lundberg and Stearns (2019)) (see also Chari and Goldsmith-Pinkham (2018)) documenting large gender differences across economics topics, although the result in Corollary 3 is too extreme, as women's frequency never breaks the 50% threshold in economics (although it does in other areas, such as psychology). Still, the result is consistent with these clustering across fields being true in the limit, which is consistent with the lack of variation or trends shows in the bottom panel of Figure 4.

#### 3.4.4. Publication Success and F-Underrepresentation

Our model is also consistent with the evidence that more research-intensive universities have lower female representation, as shown in the top panel of Figure 2 and Table 1. In particular, suppose that  $\Theta^{\max}$  contains only two distinct symmetric types  $\theta$  and  $\theta'$  with  $\lambda_0^{\theta} > \lambda_0^{\theta'}$ . We analyze the resulting limit economy; see Section 4. for numerical results with calibrated parameter values and a finite time horizon.

Consider an institution with an arbitrary fraction  $y \in [0, 1]$  of  $\theta$ -researchers and a complementary fraction (1 - y) of  $\theta'$ -researchers; since no other types survive in the limit, these are the only researcher types that an institution can employ.<sup>11</sup> Under self-image bias in refereeing, the probability that a researcher of type  $\theta$  (resp.  $\theta'$ ) publishes successfully is  $\overline{\gamma} \cdot \overline{\lambda}^{\theta}$  (resp.  $\overline{\gamma} \cdot \overline{\lambda}^{\theta'}$ ), where  $\overline{\gamma} = \gamma^{\theta} = \gamma^{\theta'}$ . Hence, the *average publication frequency* of the institution is

$$P(y) = \overline{\gamma} \left( y \overline{\lambda}^{\theta} + (1 - y) \overline{\lambda}^{\theta'} \right)$$

Since  $\overline{\lambda}^{\theta} > \overline{\lambda}^{\theta'}$  by Proposition 3, P(y) increases in y.

We assume that the type- $\theta$  and type- $\theta'$  researchers at the institution under consideration belong to the F and M groups in proportions analogous to those in the population: that is, a fraction  $y\bar{\lambda}^{\theta,f}$  are of type  $\theta$  and group F, a fraction  $(1-y)\bar{\lambda}^{\theta',f}$  are of type  $\theta'$  and group F, etc. Then, the fraction of F-researchers in an institution parameterized by y is given by:

$$F(y) = \frac{(y\overline{\lambda}^{\theta,f} + (1-y)\overline{\lambda}^{\theta',f})}{(y\overline{\lambda}^{\theta} + (1-y)\overline{\lambda}^{\theta'})}$$

Corollary 2 then implies:

**Corollary 4** Assume  $\lambda_0 = p^m$  and let  $\Theta^{\max}$  contain only two distinct symmetric types  $\theta$  and  $\theta'$ , with  $\lambda_0^{\theta} > \lambda_0^{\theta'}$ . Then F(y) is decreasing in y. That is: in the limit, institutions with higher exogenous fraction y of  $\theta$ -researcher, and a complementary fraction (1 - y) of  $\theta'$  researchers, have higher publication frequency and lower percentage of F-researchers.

Intuitively, the result follows from the fact that the limit mass of  $\theta$ -researchers in the population is higher than that of  $\theta'$ -researchers:  $\overline{\lambda}^{\theta} > 0.5 > \overline{\lambda}^{\theta'}$ . Self-image bias implies that types  $\theta$  have a higher chance of publication success, because the probability they are matched

<sup>&</sup>lt;sup>11</sup>The total mass of researchers in the institution under consideration is irrelevant to the analysis, and can be considered small. In Section 4., we consider a different parameterization in which the entire population is divided into a given, fixed number of institutions.

with referees of their own type is higher. Consequently, on average, institutions with a higher fraction y of type- $\theta$  scholars are more likely on average to achieve successful publications. However, types  $\theta$  are also more likely to come from the M group: this generates a negative relation between research success and F-representation. In the limit, F-researchers are least represented in "top institutions," where "top" is defined as in terms of publication record.

As a final comment, the above result pertains to researchers in M and F population, as referees have self-image bias and affect the publication process. The same scrutiny is unlikely to happen for "teaching faculty", as self-image bias from referees unlikely plays any role. In this case, if types  $\theta$  also relate to candidates' teaching abilities (e.g. each type  $\theta$ maps into a probability  $t^{\theta}$  of being an objectively good teacher), then the same argument in Proposition 1 for objective refereeing will yield equal representation of M and F candidates into teaching positions. Moreover, there is no correlation with the publication frequency of each institution, as shown in Figure 2 and Table 1.

#### 3.4.5. Perceived Trade-off Between Quality and Diversity

Self-image bias also explains why the current population of referees may incorrectly perceive that there is a trade-off between diversity and "quality" (see e.g. First Round Review, 2022,) or "merit" (see e.g. Crosby, Iyer, Clayton, and Downing (2003).) The intuition is simple: by definition, self-image bias implies that the closer another researcher is to one's own type, the higher their subjectively perceived quality.

In particular, our preceding results imply that the population of established scholars consists mostly of types that are more prevalent in the M-group. Therefore, a randomly drawn established researcher will subjectively perceive other M-researchers to be *ex-ante* of higher average quality than F-researchers, because it is more likely that a randomly drawn young M-researcher's type will be close to their own. Hence the mistaken impression that, in order to increase F-representation, one has to "accept" a loss of quality.

This conclusion is in stark contrast with the fact that, in our model, the average *objective* quality of each cohort M- and F-researchers, which is given by the average probability of producing quality research, is exactly the same, by construction. Furthermore, as shown in Proposition 6, when the set  $\Theta$  has a specific structure, the average quality of accepted F researchers is actually *higher* than that of accepted M researchers; the same is true for the calibrated model of Section 4., with multiple characteristics (see Fig. 9). Thus, in fact, increasing diversity can potentially *increase* average objective quality.

We now formally derive this conclusion from Proposition 5. Recall that, under self-image bias, a referee r of type  $\theta^r$  accepts a researcher of type  $\theta$  only if  $\theta = \theta^r$ . We can interpret this

by saying that referee r believes researcher  $\theta$  has a quality of  $\gamma^{\theta}$  if  $\theta = \theta^{r}$ , and 0 otherwise. Therefore, for this referee, the perceived quality of a randomly drawn young M-researcher is  $Q(M|\theta^{r}) = \gamma^{\theta^{r}} p^{\theta^{r},m}$ , while that of a randomly drawn F-researcher is  $Q(F|\theta^{r}) = \gamma^{\theta^{r}} p^{\theta^{r},f}$ . Finally, the perceived qualities of young M- and F-researchers, averaged with respect to the distribution  $\lambda_{t}$  of established researcher types, are given by

$$Q(M|\lambda_t) = \sum_{\theta \in \Theta} p^{\theta, m} \gamma^{\theta} \lambda_t^{\theta} \text{ and } Q(F|\lambda_t) = \sum_{\theta \in \Theta} p^{\theta, f} \gamma^{\theta} \lambda_t^{\theta}$$

The key observation is now that  $p^{\theta,g}\gamma^{\theta}\lambda_t^{\theta} = a_t^{\theta,g}$  for g = m, f. Therefore, summing over all possible values  $\overline{\gamma}$  of the function  $\theta \mapsto \gamma^{\theta}$  in Proposition 5 yields

**Corollary 5** Assume  $\lambda_0 = p^m$  and that there is  $\theta \in \Theta$  with  $p^{\theta,m} > p^{\sigma(\theta),m}$ . Then, the population of referees  $\lambda_t$  (mis)perceives the quality of a random *F*-researcher to be lower than the quality of a random *M*-researcher. That is,  $Q(F|\lambda_t) < Q(M|\lambda_t)$ .

## 4. Calibration

The previous section provided numerous results that hold for arbitrary sets of types  $\Theta$ . To calibrate the model, we need to be more specific about the definition of a type. We adopt a simple symmetric environment in which each type  $\theta$  corresponds to a vector of N characteristics which can only take two values, 0 and 1: that is,  $\Theta \equiv \{0,1\}^N$ . For each agent of type  $\theta \in \Theta$ ,  $\theta_n$  denotes the value of the *n*-th characteristic. The number Nof characteristics is even, characteristics are mutually independently distributed, and their distributions depend on a single parameter  $\phi > 0.5$ . Our main assumption, illustrated in Figure 5, is that characteristics are distributed symmetrically in the M and F population, in the sense that for  $n = 1, \ldots, \frac{N}{2}$ ,  $Pr(\theta_n = 1) = \phi$  for M-researchers and  $Pr(\theta_n = 1) = (1 - \phi)$ for F-researchers, and the opposite for  $n = \frac{N}{2} + 1, \ldots, N$ . For every  $\theta \in \Theta$ , let  $p^{\theta, f}$  (resp.  $p^{\theta,m}$ ) denote the fraction of types in the F (resp. M) population of young researchers. Also let  $p^g = (p^{\theta,g})_{\theta \in \Theta}$  for g = f, m. To sum up,

$$p^{\theta,m} = \prod_{n=1}^{N/2} \phi^{\theta_n} (1-\phi)^{1-\theta_n} \cdot \prod_{n=N/2+1}^N (1-\phi)^{\theta_n} \phi^{1-\theta_n}, \quad p^{\theta,f} = \prod_{n=1}^{N/2} (1-\phi)^{\theta_n} \phi^{1-\theta_n} \cdot \prod_{n=N/2+1}^N \phi^{\theta_n} (1-\phi)^{1-\theta_n}$$
(13)

Furthermore, we assume that the probability of producing quality research is given by

$$\gamma^{\theta} = \gamma_0 \cdot \rho^{\frac{1}{N} \sum_{n=1}^{N} \theta_n},$$

where  $\gamma_0 \in (0,1)$  and  $\rho \in [1, \frac{1}{\gamma_0}]$ . The two key features of this specification are that, first, types  $\theta$  with more characteristics (i.e., more 1's) are more likely to produce quality

Figure 5: Symmetric Distribution of Research Characteristics



research; and second, every research characteristic  $\theta_n$  has the same impact on our measure of objective quality. The parameter  $\gamma_0$  reflects the minimum probability of successful output the productivity of the least "capable" type  $(0, \ldots, 0)$ . The parameter  $\rho$  instead reflects the importance of research characteristics. For instance,  $\rho = 1$  means that characteristics are irrelevant—all types are equally likely to produce quality research. A higher value of  $\rho$  implies that more 1's increase the probability of quality research.

With these assumptions, the function  $\sigma : \Theta \to \Theta$  defined by  $\sigma(\theta) = (\theta_{N/2+1}, \ldots, \theta_N, \theta_1, \ldots, \theta_{N/2})$  satisfies all the assumptions of Section 3. Thus, we have:

Corollary 6 Define the threshold

$$\bar{\rho}(\phi, N) = \frac{1}{4} \left( \left( \frac{1-\phi}{\phi} \right)^{N/2} + \left( \frac{\phi}{1-\phi} \right)^{N/2} \right)^2.$$
(14)

(a) If  $\rho < \bar{\rho}(\phi, N)$ , then only two types survive,  $\Theta^{\max} = \{\theta^m, \theta^f\}$ , where

$$\theta^m = (1, ..., 1, 0, ....0) \text{ and } \theta^f = (0, ....0, 1, ..., 1)$$
 (15)

In addition, if at time 0, all referees are in the M-group with  $\lambda_0 = p^m$ , then the total limit mass of M and F researchers are

$$\bar{\Lambda}^{m} = 1 - \bar{\Lambda}^{f} = \frac{1 + \left(\frac{\phi}{1-\phi}\right)^{2N}}{1 + \left(\frac{\phi}{1-\phi}\right)^{2N} + 2\left(\frac{\phi}{1-\phi}\right)^{N}} > 0.5.$$
(16)

- (b) If  $\rho > \bar{\rho}(\phi, N)$  then, regardless of the distribution of time-0 referees,  $\Theta^{\max} = \theta^* = (1, ..., 1)$ . In addition, the limiting mass of M and F researcher are  $\bar{\Lambda}^m = \bar{\Lambda}^f = \frac{1}{2}$ .
- (c) For any set of parameters with  $\phi > 0.5, \gamma_0 \in (0, 1)$  and  $\rho \in [1, 1/\gamma_0)$ , there is N large enough such that  $\rho < \bar{\rho}(\phi, N)$  and thus point (a) holds.

(d) Moreover, for any set of parameters with  $\phi > 0.5, \gamma_0 \in (0, 1)$  and  $\rho \in [1, 1/\gamma_0)$ , as  $N \to \infty, \bar{\Lambda}^m \to 1$  and  $\bar{\Lambda}^f \to 0$ .

The parameter  $\phi$  can be easily related to Cohen's d statistic for an individual characteristic: for  $n = 1, \ldots, \frac{N}{2}$ ,

$$d = \frac{\mathrm{E}[\theta_n^i|i \in M] - \mathrm{E}[\theta_n^i|i \in F]}{\sigma_{\mathrm{pooled}}(\theta_n^i)} = \frac{2\phi - 1}{\sqrt{\phi(1 - \phi)}}.$$
(17)

For  $n = \frac{N}{2} + 1, ..., N$ , the *d* statistic is the negative of the above expression. Cohen (2013) suggests that values of *d* around 0.2 should be considered "small," values around 0.5 "medium," and values around or above 0.8 "large." However, points (c) and (d) of Corollary 6 shows that, if the number N of characteristics is sufficiently large, even small across-group differences ( $\phi - 0.5$  small) still yield our main conclusions.

Given the parametric assumptions in this section, we now turn to calibrate the model. The first issue is the number of characteristics that lead to quality research and are taken into account by referees when they evaluate a candidate. We suggest that the number of characteristics is actually large. The following is but a partial list: (i) Economic motivation; (ii) "Nose" for good questions; (iii) Institutional knowledge; (iv) Ability to find new data sources; (v) Solid identification strategy; (vi) Sophisticated empirical analysis; (vii) Clever experimental design; (viii) Skilful theoretical modelling; (ix) Ability to highlight insights, strategic effects, etc. (x) Mathematical sophistication, proof techniques, etc. (xi) Ability to position within the literature; (xii) Ability to highlight policy implications; (xiii) Presentation skills; (xiv) Ability to address questions from audience; (xv) Honesty;<sup>12</sup> and so on. Likely, there are many others. Perhaps some of these research traits are more important than others, but as a first pass, it is indeed plausible that the positive or negative result of a review depends on a combination of research characteristics, and not just a small number. In light of these considerations, and to be conservative, we assume that N = 10.

The second issue is the magnitude of between-group differences, which depends on the parameter  $\phi$ . We set  $\phi = 0.5742$ , so the implied Cohen's d is

$$d = \frac{2 \times 0.5742 - 1}{\sqrt{0.5742 \times (1 - 0.5742)}} = 0.3,$$

This value is considered "small" and in line with the estimated group differences of the various traits discussed in the introduction. With  $\phi = 0.5742$ , between-group heterogeneity

<sup>&</sup>lt;sup>12</sup>For instance, some researchers may be more keen to "torture" the data than others, or search for variables that lead to statistical significance. See e.g. discussion in Mayer (2009) and, on the impact of conflict of interests on economic research, Fabo, Jancokova, Kempf, and Pastor (2020).

in each characteristic is far smaller than within-group heterogeneity.<sup>13</sup>

Third, we need to calibrate the parameters  $\gamma_0$  and  $\rho$  in the probability of producing quality research,  $\gamma^{\theta}$ . We proceed as follows: First, we assume the best researchers  $\theta^* = (1, 1, ..., 1)$ has 100% probability of producing quality research, ie.  $\gamma^{\theta^*} = 1$ . Second, we calibrate  $\gamma_0$ to match the rate at which economics PhD students succeed in getting an academic job. We compute the latter from the NSF Survey of Doctoral Recipients. We take the ratio of economics PhD recipients who are employed in 4-year educational institutions over the total of economics PhD recipients, both inside and outside the U.S..<sup>14</sup> That ratio is 0.462. Choosing  $\gamma_0 = 0.2$  yields an objective success rate  $\sum_{\theta} \gamma^{\theta} (p^{\theta,f} + p^{\theta,m})/2 = 0.462$ . Interestingly, the implied  $\rho = \gamma^{\theta^*}/\gamma_0 = 5$  entails that researcher  $\theta^*$  is objectively five times as productive as researcher  $(0, \ldots, 0)$ , which is roughly in line with the evidence on research productivity reported in Conley and Önder (2014).<sup>15</sup>

Finally, we assume that initially *F*-researchers represent 10% of the total mass, which is roughly consistent with the percentage of women faculty in 1974, and be consistent with the distribution of annual inflows of young researchers, i.e.  $\lambda_0 = 0.9p^m + 0.1p^f$ .

### 4.1. Calibration Results

The calibration results are shown in Figures 6 through 8, which we now turn.

#### 4.1.1. *F*-Under-representation in the limit

Figure 6 shows that the system converges to a large imbalance between M- and F-researchers, with F-researchers representing around 20% of the population.<sup>16</sup> This large imbalance obtains despite the fact that the distribution of characteristics is very similar across M and Ftypes. The hump shape visible in the figure is due to our assumption that  $\lambda_0$  represents all the characteristics in proportion to the distributions  $p^{\theta,m}$  and  $p^{\theta,f}$ . Since, in our calibration,

<sup>&</sup>lt;sup>13</sup>We focus on effect size for a single characteristic as its magnitude has been widely documented in the psychology literature (see introduction). Unfortunately, we were unable to identify experimental studies measuring multidimensional effect sizes between genders for us to use in our calibration.

<sup>&</sup>lt;sup>14</sup>The 2017 survey is the latest as of the time of this writing and it is available at https://ncsesdata. nsf.gov/doctoratework/2017/index.html. The total number of economics PhD recipients is 32,000 in US and 12,750 outside the US. The total number of them working in a 4-year educational institution are 12,750 in the US and 7,900 outside the US. The ratio of economics PhDs who undertake an academic career is (12,750+7,900)/(32,000+12,750) = 0.462.

<sup>&</sup>lt;sup>15</sup>These parallels with the data should be taken with a grain of salt, given that the data would reflect the outcome of the model with self-image bias, and not just objective reference. On the other hand, we have more degrees of freedom: recall that we normalized that mass of reviewers to 1, but we can choose another mass K to match the failure rate from the data. See footnote 8.

<sup>&</sup>lt;sup>16</sup>By comparison, setting  $\phi = 0.5742$  and N = 10 in Eq. (16), the limiting fraction of M researchers is  $\bar{\Lambda}^m \approx 91\%$ . The difference is due to the fact that Eq. (16) was derived assuming that  $\lambda_0 = p^m$ .



Percent of F- researchers in calibrated model. Parameters:  $\phi = 0.5742$ , d = 0.3,  $\gamma_0 = 0.2$ ,  $\rho = 5$ , N = 10. Initially  $\lambda_0 = 0.9 p^m + 0.1 p^f$ .

the differences between  $p^m$  and  $p^f$  are small, initially, there is a sufficient mass of referees with characteristics that are common in the *F*-population to yield a high acceptance rate of *F*-researchers and thus an increase in their mass. However, as characteristics common in *F*-population start being weeded out, eventually the acceptance rate of *F*-researchers drops, and so does the overall mass of *F*-researchers. While such hump is not visible in the data, the limit fraction of around 20% of *F*-researchers is rather close to the 25% women faculty as assistant professor observed in the bottom panel of Figure 1.

Given that our mechanism is only based on self-image bias – a bias of referees that is gender-neutral and is likely to be common across countries – our results also explain why low representation of women in economics occurs nearly everywhere, including Nordic European countries, which have a far more balanced stance towards women compared to the United States (see Figure 3).

#### 4.1.2. Higher under-representation in top institutions

Next, we build on Section 3.4.4. and examine the publication success of researchers in institutions that differ in their type composition. Specifically, consider J institutions, and assume that each institution  $j = 1, \ldots, J$  employs an exogenously specified fraction  $x_j^{\theta}$  of all type- $\theta$  researchers, with  $\sum_{j=1}^{J} x_j^{\theta} = 1$ . The mass of  $\theta$  researchers in institution j at time t is thus  $x_j^{\theta} \lambda_t^{\theta}$ . We make no assumptions about the distribution of groups in institutions.

Since, at each time t, institution j employs  $x_i^{\theta} \lambda_t^{\theta}$  researchers of type  $\theta$ , each such researcher

Figure 7: The Endogenous Negative Relation between Institutions' Publishing Intensity and the Percentage of F-researchers



This figure plots scatterplot of 100 simulated institution publishing probabilities (x-axis) versus the *F*-researcher representation in the same institution. Initially  $\lambda_0 = 0.9p^m + 0.1p^f$ . Parameters:  $\phi = 0.5742$ , d = 0.3,  $\gamma_0 = 0.2$ ,  $\rho = 5$ , N = 10.

publishes with probability  $\gamma^{\theta} \lambda_t^{\theta}$ , and the total mass of researchers employed by institution j is  $\sum_{\theta} x_j^{\theta} \lambda_{j,t}^{\theta}$ , the weighted-average probability of publications of (established) researchers in institution j is

$$P_{j,t} = \frac{\sum_{\theta \in \Theta} \gamma^{\theta} \left(\lambda_{t}^{\theta}\right)^{2} x_{j}^{\theta}}{\sum_{\theta \in \Theta} \lambda_{t}^{\theta} x_{j}^{\theta}}.$$

On the other hand, the fraction of F-researchers in institution j is

$$F_{j,t} = \frac{\sum_{\theta \in \Theta} \lambda_t^{\theta, f} x_j^{\theta}}{\sum_{\theta \in \Theta} \lambda_t^{\theta} x_j^{\theta}}$$

Figure 7 shows the scatterplot of  $P_{j,t}$  and  $F_{j,t}$  from the model simulation for a random draw of  $x_j$ 's for J = 100 institutions, for large t (i.e., "in the limit"). As the plot shows, there is a negative relation between an institution publishing intensity (x-axis) and its Frepresentation (y-axis). The negative slope in Figure 7 is consistent with the theoretical result in Corollary 4 and with the empirical findings in Table 1, showing that indeed, the fraction of women in economics departments is negatively related to the publication intensity of the same departments.

The two panels of Figure 8 replicate in the model the corresponding results collected in the two panels of Figure 2. Specifically, in Panel (a) we first rank the J institution in terms of publishing intensity  $P_{j,t}$ , and then take the average of the fraction of F-researchers across the top-10, top-20, and, respectively, all institutions. While the absolute levels are smaller



Figure 8: F-researchers Under-representation in Top Research Institutions

(a) Fraction of *F*-researchers across Institutions

(b) The Dynamics of F Researchers in Top Research Institutions:



Panel (a) plots the percent of *F*-researcher across 100 simulated institutions, the top 20, and the top 10. Panel (b) reports the dynamics of the fraction *F*-researchers in top research institutions and across all institutions as simulated from the model. Top research institutions are the top 20 out 100 with the highest frequency of publication. Parameters:  $\phi = 0.5742$ , d = 0.3,  $\gamma_0 = 0.2$ ,  $\rho = 5$ , N = 10. Initially  $\lambda_0 = 0.9p^m + 0.1p^f$ .

in our model, the close match with Panel (a) in Figure 2 is surprising, given that our model has no group bias at all.

Indeed, our model also replicates the time-series dynamics of the fraction of F-researchers over time in top institutions. Panel (b) of Figure 8 shows that the gap between the fraction of F-researchers in top institutions vs. all institutions increases over time. Intuitively, sorting institutions by their publishing intensity implicitly defines "top institutions" as those whose types are increasingly similar to the types of the majority of referees—that is, in the limit,



Figure 9: Quality and Clustering across Fields of M- and F-researchers in Calibrated Model

(a) Quality of Researchers

Panel (a) plots the average quality of accepted M and F researchers in the model. The quality is measured as  $\sum_{\theta} \gamma^{\theta} w_t^{\theta,g}$  where  $\gamma^{\theta} = \gamma_0 \rho^{\frac{1}{N} \sum_{n=1}^{N} \theta_n}$  and  $w_t^{\theta,g} = a_t^{\theta,g} / \sum_{\theta'} a_t^{\theta',g}$ , g = f, m. Panel (b) reports clustering across fields implies by  $\theta^m$  and  $\theta^f$ . Parameters:  $\phi = 0.5742$ , d = 0.3,  $\gamma_0 = 0.2$ ,  $\rho = 5$ , N = 10. Initially  $\lambda_0 = 0.9p^m + 0.1p^f$ .

 $\overline{\lambda}^{\theta^m}$ . By construction, the remaining institutions will have a larger fraction of researchers that are less represented in the refereeing population.

#### 4.1.3. Higher quality of successful *F*-researchers and clustering

Finally, the two panels of Figure 9 replicate in the model the corresponding panels in Figure 4 in the data. Specifically, Panel (a) of Figure 9 plots the average quality of F- and M-researchers conditional on being accepted, and shows that the average quality of F- researchers is uniformly higher than M- researchers. This plot is consistent with Proposition 6 and our conjecture that the result should hold for every N.<sup>17</sup> This results is consistent with the top panel in Figure 4.

Panel (b) of Figure 9 shows that the model produces clustering across fields, consistently with bottom panel in Figure 4. Specifically, consider types  $\theta^m$  and  $\theta^f$ . Panel (b) plots the difference between the share of F- and M-researchers across the  $\theta^f$ -field and the  $\theta^m$ -field. As also shown in Corollary 3, F-researchers are relatively more represented in the former (top line), and M-researcher relatively more the latter (bottom line). While our model yields more extreme predictions relative to the data, with only two fields surviving in the limit ( $\theta^f$  and  $\theta^m$ ), the lack of dynamics in the data is consistent with our model's predictions. Moreover, in the extension Section 6., we show that, with multiple distinct types  $\theta \in \Theta^{\max}$ and endogenous entry in academia, the model also predicts multiple "fields" and a higher number of M-dominated fields than F-dominated fields; this is consistent with the data.

## 5. The Impact of Policy Actions

In this section we discuss the impact of policy actions that have been proposed to address gender imbalance. We consider (i) the impact of mentoring (section 5.1.); and (ii) the impact of affirmative action (section 5.2.).

### 5.1. Mentoring: Group Balance versus Talent Loss

The adoption of mentoring to improve the prospects of female economists is one of the most popular proposals. Indeed, there is evidence that mentoring does help increase the success rate of female economists (Ginther, Currie, Blau, and Croson (2020)). We now investigate the implications of mentoring in our model.

We assume that at the beginning of each period t every young researcher of type  $\theta$  is randomly matched with an advisor a of type  $\theta^a$  drawn from the established group, whose mass is  $\lambda_{t-1}^{\theta^a}$ . Upon matching, the researcher of type  $\theta$  can choose to pay a cost  $C(\theta, \theta^a)$  to "become" the same type of the advisor. Assume that P is the payoff from being hired and U is the utility from an outside option. Researcher  $\theta$  will then pay the cost if and only if

$$\gamma^{\theta^a} \lambda_{t-1}^{\theta^a} \left( P - C(\theta, \theta^a) \right) + \left( 1 - \gamma^{\theta^a} \lambda_{t-1}^{\theta^a} \right) \left( U - C(\theta, \theta^a) \right) > \gamma^{\theta} \lambda_{t-1}^{\theta} P + \left( 1 - \gamma^{\theta} \lambda_{t-1}^{\theta} \right) U$$

<sup>&</sup>lt;sup>17</sup>Again, Proposition 6 assumes that  $\lambda_0 = p^m$ ; the results in Panel (a) of Figure 9 thus suggest that the conclusions of the Proposition are robust to small changes the initial population.

That is, a young researcher  $\theta$  pays the cost if and only if

$$\widetilde{C}(\theta, \theta^a) = \frac{C(\theta, \theta^a)}{P - U} < \gamma^{\theta^a} \lambda_{t-1}^{\theta^a} - \gamma^{\theta} \lambda_{t-1}^{\theta}$$

In words, the increase in the probability of getting hired must be sufficiently high relative to the cost of undergoing mentoring. For instance, if the right-hand-side was negative (type  $\theta$  is already likely to succeed), nobody of that type would pay such a cost.

We assume that the cost itself depends on the distance between the young researcher's type  $\theta$  and the type of the advisor  $\theta^a$ : The larger the distance and the higher the cost, indicating that it will take a higher effort to "learn" to become a type that is likely to be hired. Note that such distance may be high as the young researcher  $\theta$  may have some characteristics that are desirable from an objective standpoint, but that are not viewed as important or relevant by the majority of established researchers. The cost, in that case, is to "unlearn" what is deemed "irrelevant."

The Online Appendix contains the details of the system dynamics for the model parameterization in Section 4.. For brevity, we only provide the intuition here. Panel (a) of Figure 10 illustrates the dynamics under the same parameters as in Section 4. and a cost function  $C(\theta, \theta') = \beta \sum_{n=1}^{N} (\theta_n - \theta'_n)^2$ , with  $\beta = 0.075$ . We choose this cost so that not all of the young researchers want to pay the switching cost to become like their advisors, which seems plausible. The resulting steady state is roughly consistent with the percentage of female participation in economics.

Initially, the dynamics are as in the base case, as all  $\lambda_t^{\theta}$  are small and thus no young researcher wants to pay the cost of mentoring. In this dynamics, as we know,  $\lambda_t^{\theta^m}$  and  $\lambda_t^{\theta^f}$  increase, with the former increasing faster, as shown in panel (c) of Figure 10. At some point, the mass of  $\lambda_t^{\theta^m}$  becomes large enough to induce many young researchers, both M and F, to pay the mentoring cost, and the system (nearly) jumps. The reason is that many young researchers now expect that their advisor will likely be of type  $\theta^m$ , which is also the type of established researchers who will evaluate their research. They are thus happy to pay the cost and become like their advisors.

The bottom panels of Figure 10 show, however, that the mass of young M-researchers jumps by more than the mass of F-researchers. The reason is that even though the cost function is the same for M- and F-researchers, young M-researchers are on average closer to  $\theta^m$  and thus have have systematically lower cost to switch than F-researchers. For this reason, group imbalance persists forever.<sup>18</sup> Moreover, only type  $\theta^m$  survives and therefore the research characteristics mildly more common in the F-population, but also very common

<sup>&</sup>lt;sup>18</sup>If the cost function was lower, however, then *all* young researchers, M and F, would pay the cost and the system would jump to group balance.



(a) Fraction of F-Researchers

Fraction of M and F researchers (panel (a)), and mass of established F-researchers (panel (b)) and of M-researchers (panel (c)) under costly mentoring. Initial  $\lambda_0 = 0.9p^m + 0.1p^f$ . Parameters:  $\phi = 0.5742$  (d = 0.3),  $\gamma_0 = 0.2$ ,  $\rho = 5$ , N = 10, cost function  $C(\theta, \theta') = 0.0750 \sum_{n=1}^{N} (\theta_n - \theta'_n)^2$ .

in the *M*-population, disappear, thus yielding talent loss and loss of knowledge.

### 5.2. Affirmative Action

A common policy to increase diversity is "affirmative action", which effectively increases the representation of specified groups by mandate. We consider a simple rule in this section: in each round, it is mandated that evaluators must hire the same number of M and F researchers. We change just one assumption to the dynamics in the benchmark case, namely:

$$a_t^{\theta,m} = k_t \gamma^{\theta} \lambda_{t-1}^{\theta} p^{\theta,m} \quad \text{where} \quad k_t = \frac{\sum_{\theta'} \gamma^{\theta'} \lambda_{t-1}^{\theta'} p^{\theta',f}}{\sum_{\theta'} \gamma^{\theta'} \lambda_{t-1}^{\theta'} p^{\theta',m}}.$$
 (18)



(a) Fraction of F-Researchers

Fraction of F researchers (panel (a)), mass of established F-researchers (panel (b)) and of M-researches (panel (c)) when an affirmative action policy requires to accept the same number of M and F researchers. Parameters:  $\phi = 0.5742$  (d = 0.3),  $\gamma_0 = 0.2$ ,  $\rho = 4$ , and N = 10. Initial  $\lambda_0 = 0.9p^m + 0.1p^f$ 

The scaling factor  $k_t$  ensures that  $\sum_{\theta} a_t^{\theta,f} = \sum_{\theta} a_t^{\theta,m}$ . Figure 11 provide the dynamics for this case. Affirmative action reaches group balance, which is not surprising. However, it also attains diversity in research characteristics: in the limit, M researchers are of type  $\theta^m$  and F researchers are of type  $\theta^f$ . Assuming that maximizing the representation of research characteristics is beneficial to society, this policy appears superior to mentoring, as it does not skew the distribution of such characteristics towards  $\theta^m$  even when reaching group balance.

Intuitively, by expanding the set of referee characteristics, affirmative action makes it possible to reward the research of *talented* F researchers—those who are more likely to produce quality research. It is still the case that F researchers who are not (objectively) as productive will not survive in the limit and will be weeded out from the system.

## 6. Extension: Endogenous Entry

In this section we summarize the results from an extension of the model to endogenous entry.<sup>19</sup> We assume that candidates may choose a career in academia, in which case their success depends on the judgement of the established group of researchers as described in Section 3., or can opt for a different career. If they choose the academic career, candidates pay a utility cost C but receive a payoff P if successful. The outside options gives a benchmark utility of zero. We obtain several results:

- 1. When the relative cost C/P is small and below a cutoff, the same equilibrium and results as in Section 3.3. obtain;
- 2. When the relative cost C/P is intermediate, the equilibrium changes. In particular,
  - (i) The set of surviving types  $\Theta^{\max}$  shrinks: Even in environments where, for C = 0,  $\Theta^{\max}$  in (5) contains pairs of distinct symmetric types (see Definition 1), when C/P is intermediate, some types that are more common in the *F*-population may not survive in the limit.
  - (ii) This new force increases the limit imbalance towards the *M*-population, and leads to further talent loss, as young researchers of certain types will not even apply for an academic career.
  - (iii) Moreover, assuming that there is a mapping between fields and types as in Corollary 3, the number of *M*-dominated fields is higher than the number of *F*-dominated fields, as is the case in the data (see the bottom panel of Fig. 4).
- 3. Finally, in the special case of the calibrated model in Section 4., we also show that the pool of applicants skew towards the *M*-population. That is, imbalance occurs even in the "pipeline," which may explain the low percentage of women applications to PhD program, for instance.

# 7. Literature Review

There is a considerable body of research on the underlying reason of under-representation of women in the economics profession. We do not attempt an exhaustive survey here, but

<sup>&</sup>lt;sup>19</sup>The on-line appendix contains additional extensions, including a two-layer hierarchy of junior and senior established researchers; attribution of co-authored work; and more general forms of self-image bias. In addition, we also show that the same equilibrium as in this section occurs when hiring institutions decide to hire candidates based on their probability of success as opposed to their objective quality.

refer the reader to Bayer and Rouse (2016), who review the literature on both "supplyside" and "demand-side" factors. Among supply-side factors, the imbalance in Economics PhD applications appears to depend on prior exposure to economics, the performance in introductory courses, and the lack of role models, but, interestingly, not on math preparation. On the demand side, Bayer and Rouse (2016) suggest that policy changes in most academic institutions have diminished the impact of explicit or statistical discrimination in recruiting Ph.D. students. However, they argue that an important role is played by *implicit bias* and *stereotyping*. Our model with self-image bias is consistent with the persistence of gender bias even when all structural sources of gender-biases have been removed.

In a more recent contribution, Sarsons et al. (2021)'s work on recognition for coauthored papers shows that, for men, an additional coauthored paper has the same effect on the likelihood of tenure as a solo-authored paper; however, for women, coauthorship entails a significant "discount factor," especially if the coauthor(s) are men. The large body of research on the gender pay gap and on the "glass ceiling" in other labor markets is also indirectly relevant in our context: see e.g. Blau and Kahn (2017); Goldin and Rouse (2000); Goldin (2014); Weber and Zulehner (2014); Aigner and Cain (1977); Lazear and Rosen (1990).

On the theoretical side, our model is related to the literature on statistical discrimination: a relative recent survey is Fang and Moro (2011). One strand within that literature, originating from Phelps (1972), posits the existence of exogenous differences between groups, either in the distribution of productivity ("Case 1"), or in the quality of signals about it ("Case 2"). In Case 2, the employer does not observe the productivity of individual applicants, but receives a signal about it. Differential average treatment of the two groups can emerge either through risk aversion of the employer (Aigner and Cain, 1977), investment in human capital (Lundberg and Startz, 1983), or if hiring occurs in a tournament (Cornell and Welch, 1996). In Conde-Ruiz, Ganuza, and Profeta (2020), the difference in signal quality leads members of the group in the minority of a hiring committee to underinvest in human capital; this perpetuates the imbalance. A recent contribution, Bardhi, Guo, and Strulovici (2019), revisits Phelp's Case 1, but assume that success or failure is observed over time and is informative about the worker's type. This can lead to large differences in ex-post treatment of the two groups, even if ex-ante productivity differences are small. Differently from this literature, in our model the ex-ante distributions of productivity are the same in the M and F group, because all characteristics are equally valuable. Furthermore, productivity is observed. In our model, standard statistical discrimination does not lead to gender imbalance.

Becker (1957)'s model of taste-based discrimination instead posits that employers may have a preference for hiring members of one specific group. This is not the case in our model: while referees only accept applicants whose research characteristics match their own, they do not take group membership into consideration at all.

Heidhues, Kőszegi, and Strack (2019) proposes a model in which an agent's ability is unobserved, both by herself and by others. Agents belong to different groups, each potentially subject to "discrimination," and are "stubbornly overconfident" about their own ability. Overconfidence leads agents to have a more favorable view of individuals in their own social group, ascribing poor performance to discrimination against them. In our model, ability is observed, and there is no exogenously imposed discrimination on either group. Incorporating (possibly biased) learning (cf. e.g. Bohren, Imas, and Rosenberg, 2019) about young researchers' characteristics is an interesting direction for future work.

## 8. Conclusions and Policy Implications

Our model highlights a novel mechanism that endogenously perpetuates specific research characteristics over time without relying on implicit or explicit gender bias. This occurs due to self-image bias, grounded in the psychology literature, and its application to the reviewing process: established researchers use their own personal research characteristics as a guidance to judge others' output. Findings in psychology and experimental economics point to mild between-group heterogeneity; yet, in our model, such mild differences are enough to lead the initially prevalent group to dominate forever. It is *as if* the initially dominant group decided for society what are the important research characteristics and topics in Economics.

Our theoretical and numerical results are consistent with numerous empirical regularities that we collected in Section 2., which we do not repeat for brevity.

Standard solutions to the gender bias problem may not be very effective in our model. For instance, outreach programs to encourage members of a given group to apply to PhD programs may prove ineffective. Such outreach program are akin to lowering the cost of doing research, but in our model, the cost is zero. Similarly, mentorship programs for female researchers may also not be effective (see Section 5.1.), and have the unintended consequence to induce a talent loss, as female researchers give up their research characteristics to adopt those that are prevalent in the reviewer population. In contrast, affirmative action policies not only bring gender balance, but also help reach first best, as all research characteristics are properly represented in the limit (section 5.2.).

Because the problem is self-image bias, the best policy intervention must involve limiting the ability of reviewers to use their own research style as a yardstick while judging others' research. One solution is to provide strict guidelines in the refereeing process. Indeed, in light of Corollary 6 (d), editors should guide referees to limit the number of aspects of the submitted research paper they should focus on. Dunning, Meyerowitz, and Holzberg (1989) provides suggestive evidence in support of this approach.

Another solution is instead to change the reviewing process to include input from the full distribution of researchers, as opposed to just the established ones. While radical as a proposal, it would be reasonable to consider an editorial policy that requires young researchers to participate in the evaluation process, or in fact, "oversample" young female researchers.

## A Appendix: Results for the general model

**Proof of Proposition 1** Eq. (1) shows that  $a_t^{\theta,g}$  is time-invariant for  $g \in \{f, m\}$ ; hence, so is  $a_t^{\theta}$ , and therefore  $a_t$ . Dropping time indices, for  $g \in \{f, m\}$ , a straightforward derivation shows that

$$\lambda_t^{\theta,g} = (1-a)\lambda_{t-1}^{\theta,g} + a^{\theta,g} = (1-a)^t \lambda_0^{\theta,g} + a^{\theta,g} \frac{1-(1-a)^t}{a} \to \frac{a^{\theta,g}}{a} = \frac{\gamma^\theta p^{\theta,g}}{a}, \qquad (19)$$

so the limiting fraction of M- to F-researchers is

$$\frac{\sum_{\theta} a^{\theta,m}}{\sum_{\theta} a^{\theta,g}} = \frac{\sum_{\theta} \gamma^{\theta} p^{\theta,m}}{\sum_{\theta} \gamma^{\theta} p^{\theta,f}}.$$

By Assumption 1, for every  $\theta$ , the type  $\bar{\theta} = \sigma(\theta)$  satisfies  $p^{\theta,m} = p^{\bar{\theta},f}$  and  $\gamma^{\theta} = \gamma^{\bar{\theta}}$ ; hence, the above fraction equals 1. *Q.E.D.* 

To prove the main results of the paper, we first characterize key features of the population dynamics for an arbitrary, finite set  $\Theta$  of types, with initial distribution  $\lambda_0 \in \Delta(\Theta)$ , such that  $\lambda_0 = \lambda_0^m + \lambda_0^f$  for  $\lambda_0^m, \lambda_0^f \in \mathbb{R}_+^{\Theta}$ , and per-period inflows  $q^g = (q^{\theta,g})_{\theta \in \Theta} \in \mathbb{R}_+^{\Theta} \setminus \{0\}$ , for  $g \in \{f, m\}$ . Also define  $q = q^m + q^f$ . Then, for  $g \in \{f, m\}$ , the dynamics are given by

$$\lambda_t^{\theta,g} = \lambda_{t-1}^{\theta,g} \left( 1 - \sum_{\theta'} \lambda_{t-1}^{\theta'} q^{\theta'} \right) + \lambda_{t-1}^{\theta} q^{\theta,g}; \qquad \lambda_t^{\theta} = \lambda_t^{\theta,m} + \lambda_t^{\theta,f}.$$
(20)

**Theorem 1** Assume that  $q^{\theta} \leq 1$  for all  $\theta \in \Theta$ . Then, for all  $t \geq 0$ ,  $\lambda_t \in \Delta(\Theta)$ , and  $\lambda_t^m, \lambda_t^f \in \mathbb{R}^{\Theta}_+$ . Moreover:

- 1. if  $\lambda_0^{\theta} = 0$ , then  $\lambda_t^{\theta} = 0$  for all  $t \ge 0$ ;
- 2. if  $\lambda_0^{\theta} > 0$ , then  $\lambda_t^{\theta} > 0$  for all  $t \ge 0$ ;
- 3. for  $\theta, \tilde{\theta} \in \Theta$  with  $\lambda_0^{\theta} \cdot \lambda_0^{\tilde{\theta}} > 0$ :

(a) 
$$\frac{\lambda_t^{\theta}}{\lambda_{t-1}^{\theta}} - \frac{\lambda_t^{\theta}}{\lambda_{t-1}^{\tilde{\theta}}} = q^{\theta} - q^{\tilde{\theta}} \text{ for all } t \ge 1, \text{ and}$$
  
(b)  $q^{\theta} > q^{\tilde{\theta}} \text{ implies } \frac{\lambda_t^{\theta}}{\lambda_t^{\tilde{\theta}}} \to \infty, \text{ and } q^{\theta} = q^{\tilde{\theta}} \text{ implies } \frac{\lambda_t^{\theta}}{\lambda_t^{\tilde{\theta}}} = \frac{\bar{\lambda}_o^{\theta}}{\bar{\lambda}_o^{\tilde{\theta}}} \text{ for all } t \ge 0;$ 

4. define the set

$$\Theta^{\max} = \{ \theta \in \Theta : \lambda_0^{\theta} > 0, \ \theta \in \arg \max_{\theta' \in \Theta} q^{\theta'} \}$$
(21)

and let  $\overline{\lambda} \in \Delta(\Theta)$  be such that

$$\bar{\lambda}^{\tilde{\theta}} = \begin{cases} \frac{\lambda_0^{\tilde{\theta}}}{\sum_{\theta \in \Theta^{\max}} \lambda_0^{\theta}} & \tilde{\theta} \in \Theta^{\max} \\ 0 & \tilde{\theta} \notin \Theta^{\max} \end{cases}$$
(22)

then  $\lim_{t\to\infty} \lambda_t = \bar{\lambda};$ 

5. define

$$\bar{\lambda}^{\tilde{\theta},f} = \begin{cases} \frac{\lambda_0^{\tilde{\theta}} q^{\tilde{\theta},f}}{\sum_{\theta \in \Theta^{\max}} \lambda_0^{\theta} q^{\theta}} & \tilde{\theta} \in \Theta^{\max} \\ 0 & \tilde{\theta} \notin \Theta^{\max} \end{cases} \quad and \quad \bar{\lambda}^{\tilde{\theta},m} = \begin{cases} \frac{\lambda_0^{\tilde{\theta}} q^{\tilde{\theta},m}}{\sum_{\theta \in \Theta^{\max}} \lambda_0^{\theta} q^{\theta}} & \tilde{\theta} \in \Theta^{\max} \\ 0 & \tilde{\theta} \notin \Theta^{\max} \end{cases}$$
(23)

then  $\lim_{t\to\infty} \lambda_t^f = \bar{\lambda}^f$  and  $\lim_{t\to\infty} \lambda_t^m = \bar{\lambda}^m$ .

**Proof:** Eq. (20) implies that

$$\lambda_t^{\theta} = \left(1 - \sum_{\theta' \in \Theta} \lambda_{t-1}^{\theta'} q^{\theta'}\right) \lambda_{t-1}^{\theta} + \lambda_{t-1}^{\theta} q^{\theta}.$$
 (24)

By assumption  $\lambda_0 \in \Delta(\Theta)$ . Inductively, suppose  $\lambda_{t-1} \in \Delta(\Theta)$  and  $\lambda_{t-1}^m, \lambda_{t-1}^f \in \mathbb{R}_+^\Theta$ . Summing over  $\Theta$  on both sides of Eq. (24) yields  $\sum_{\theta} \lambda_t^{\theta} = (1 - \sum_{\theta'} \lambda_{t-1}^{\theta'} q^{\theta'}) (\sum_{\theta} \lambda_{t-1}^{\theta}) + \sum_{\theta} \lambda_{t-1}^{\theta} q^{\theta} = 1$ . Furthermore, since  $\lambda_{t-1} \in \Delta(\Theta), \sum_{\theta'} \lambda_{t-1}^{\theta'} q^{\theta'} \in [\min_{\theta'} q^{\theta'}, \max_{\theta'} q^{\theta'}] \subseteq [0, 1]$ ; moreover,  $q^{\theta} \ge 0$  and  $\lambda_{t-1}^{\theta} \ge 0$ , so Eq. (24) implies that  $\lambda_t^{\theta} \ge 0$  as well. By the same argument,  $q^{\theta} \ge 0$  and  $\lambda_{t-1}^{\theta, g} \ge 0$  for  $g \in \{f, m\}$  as well by Eq. (20). Thus,  $\lambda_t \in \Delta(\Theta)$ , and  $\lambda_t^g \in \mathbb{R}_+^\Theta$  for each g.

Claim 1 is immediate. For Claim 2, again we argue by induction. For t = 0, the claim is trivially true. Inductively, assume  $\lambda_{t-1}^{\theta} > 0$ . By Eq. (24), since as was just shown  $1 - \sum_{\theta'} \lambda_{t-1}^{\theta'} q^{\theta'} \ge 0$ , and the inductive hypothesis implies that  $\lambda_{t-1}^{\theta} > 0$ , if  $q^{\theta} > 0$  then  $\lambda_t^{\theta} \ge \lambda_{t-1}^{\theta} q^{\theta} > 0$ . Suppose instead  $q^{\theta} = 0$ . If  $\sum_{\theta'} \lambda_{t-1}^{\theta'} q^{\theta'} = 1$ , then, since  $q^{\theta'} \le 1$  for all  $\theta'$  by assumption, and  $\lambda_{t-1} \in \Delta(\Theta)$ , it must be that  $\lambda_{t-1}^{\theta'} > 0$  implies  $q^{\theta'} = 1$ : but then  $\lambda_{t-1}^{\theta} = 0$ , which contradicts the inductive hypothesis. Thus,  $0 \le \sum_{\theta'} \lambda_{t-1}^{\theta'} q^{\theta'} < 1$ , so Eq. (24) implies that  $\lambda_t^{\theta} = (1 - \sum_{\theta'} \lambda_{t-1}^{\theta'} q^{\theta'}) \lambda_{t-1}^{\theta} > 0$ .

For Claim 3, divide both sides of Eq. (24) for type  $\theta$  by  $\lambda_{t-1}^{\theta}$ , which is assumed to be positive; this yields

$$\frac{\lambda_t^{\theta}}{\lambda_{t-1}^{\theta}} = 1 + q^{\theta} - \sum_{\theta'} \lambda_{t-1}^{\theta'} q^{\theta'}.$$
(25)

A similar equation holds for  $\tilde{\theta}$ . This immediately yields 3(a). To derive 3(b), since  $\lambda_t^{\theta'} = \lambda_0^{\theta'} \cdot \prod_{s=1}^t \frac{\lambda_s^{\theta'}}{\lambda_{s-1}^{\theta'}}$  for  $\theta' = \theta, \tilde{\theta}$ ,

$$\frac{\lambda_t^{\theta}}{\lambda_t^{\tilde{\theta}'}} = \frac{\lambda_0^{\theta}}{\lambda_0^{\tilde{\theta}}} \cdot \frac{\prod_{s=1}^t \frac{\lambda_s^{\theta}}{\lambda_{s-1}^{\tilde{\theta}}}}{\prod_{s=1}^t \frac{\lambda_0^{\tilde{\theta}}}{\lambda_{s-1}^{\tilde{\theta}}}} = \frac{\lambda_0^{\theta}}{\lambda_0^{\tilde{\theta}}} \cdot \prod_{s=1}^t \frac{\frac{\lambda_s^{\theta}}{\lambda_{s-1}^{\tilde{\theta}}}}{\lambda_s^{\tilde{\theta}}} = \frac{\lambda_0^{\theta}}{\lambda_s^{\tilde{\theta}}} \cdot \prod_{s=1}^t \frac{\frac{\lambda_s^{\tilde{\theta}}}{\lambda_{s-1}^{\tilde{\theta}}}}{\frac{\lambda_s^{\tilde{\theta}}}{\lambda_{s-1}^{\tilde{\theta}}}} = \frac{\lambda_0^{\theta}}{\lambda_s^{\tilde{\theta}}} \cdot \prod_{s=1}^t \frac{\frac{\lambda_s^{\theta}}{\lambda_{s-1}^{\tilde{\theta}}}}{\frac{\lambda_s^{\tilde{\theta}}}{\lambda_{s-1}^{\tilde{\theta}}}} = \frac{\lambda_0^{\theta}}{\lambda_0^{\tilde{\theta}}} \cdot \prod_{s=1}^t \left(1 + \frac{q^{\theta} - q^{\tilde{\theta}}}{\frac{\lambda_s^{\tilde{\theta}}}{\lambda_{s-1}^{\tilde{\theta}}}}\right).$$

If  $q^{\theta} = q^{\tilde{\theta}}$ , then every term in parentheses equals 1, and the claim follows. If instead  $q^{\theta} > q^{\tilde{\theta}}$ , recall that, by Eq. (25), for all  $s \ge 1$ , since  $\lambda_{s-1} \in \Delta(\Theta)$  and  $q \in [0,1]^{|\Theta|}$ ,  $\frac{\lambda_s^{\tilde{\theta}}}{\lambda_{s-1}^{\tilde{\theta}}} \le 1 + q^{\tilde{\theta}}$ . Therefore, each term in parentheses is not smaller than  $1 + \frac{q^{\theta} - q^{\tilde{\theta}}}{1 + q^{\tilde{\theta}}} > 1$ . It follows that

$$\frac{\lambda_t^{\theta}}{\lambda_t^{\tilde{\theta}'}} = \frac{\lambda_0^{\theta}}{\lambda_0^{\tilde{\theta}}} \cdot \prod_{s=1}^t \left( 1 + \frac{q^{\theta} - q^{\tilde{\theta}}}{\frac{\lambda_s^{\tilde{\theta}}}{\lambda_{s-1}^{\tilde{\theta}}}} \right) \ge \frac{\lambda_0^{\theta}}{\lambda_0^{\tilde{\theta}}} \cdot \left( 1 + \frac{q^{\theta} - q^{\tilde{\theta}}}{1 + q^{\tilde{\theta}}} \right)^t \to \infty.$$

For Claim 4, consider first  $\tilde{\theta} \notin \Theta^{\max}$ , and fix  $\theta \in \Theta^{\max}$  arbitrarily. Then  $\frac{\lambda_t^{\theta}}{\lambda_t^{\tilde{\theta}}} \to \infty$  by Claim 3(b). Suppose that there is a subsequence  $(\lambda_{t(\ell)})_{\ell \geq 0}$  such that  $\lambda_{t(\ell)}^{\tilde{\theta}} \geq \epsilon$  for some  $\epsilon > 0$  and all  $\ell \geq 0$ . Since  $\frac{\lambda_{t(\ell)}^{\theta}}{\lambda_{t(\ell)}^{\theta}} \to \infty$  as well, there is  $\ell$  large enough such that  $\frac{\lambda_{t(\ell)}^{\theta}}{\lambda_{t(\ell)}^{\theta}} > \frac{1}{\epsilon}$ : but then  $\lambda_{t(\ell)}^{\theta} > 1$  for such  $\ell$ : contradiction. Thus, for every  $\epsilon > 0$ , eventually  $\lambda_t^{\tilde{\theta}} < \epsilon$ : that is,  $\lambda_t^{\tilde{\theta}} \to 0$ . Next, consider  $\tilde{\theta} \in \Theta^{\max}$ . By Claim 2,  $\lambda_t^{\tilde{\theta}} > 0$  and  $\sum_{\theta \in \Theta^{\max}} \lambda_t^{\theta} > 0$ , and

$$\frac{\lambda_t^{\tilde{\theta}}}{\sum_{\theta \in \Theta^{\max}} \lambda_t^{\theta}} = \frac{1}{\sum_{\theta \in \Theta^{\max}} \frac{\lambda_t^{\theta}}{\lambda_t^{\theta}}} = \frac{1}{\sum_{\theta \in \Theta^{\max}} \frac{\lambda_0^{\theta}}{\lambda_0^{\theta}}} = \frac{\lambda_0^{\tilde{\theta}}}{\sum_{\theta \in \Theta^{\max}} \lambda_0^{\theta}} = \bar{\lambda}^{\tilde{\theta}}$$

where the third inequality follows from Claim 3(b). Therefore,

$$\lambda_t^{\tilde{\theta}} = \frac{\lambda_t^{\tilde{\theta}}}{\sum_{\theta \in \Theta^{\max}} \lambda_t^{\theta}} \cdot \left(\sum_{\theta \in \Theta^{\max}} \lambda_t^{\theta}\right) = \bar{\lambda}^{\tilde{\theta}} \cdot \left(1 - \sum_{\theta \notin \Theta^{\max}} \lambda_t^{\theta}\right) \to \bar{\lambda}^{\tilde{\theta}},$$

because, as was just shown above,  $\lambda_t^{\theta} \to 0$  for  $\theta \notin \Theta^{\max}$ .

Finally, consider Claim 5. Fix  $g \in \{f, m\}$ . First, since  $0 \leq \lambda_t^{\theta,g} \leq \lambda_t^{\theta}$  for all  $t \geq 0$ , if  $\theta \notin \Theta^{\max}$  then by Claim  $4 \lambda_t^{\theta} \to \overline{\lambda}^{\theta} = 0$ , and so  $\lambda_t^{\theta,g} \to 0 = \overline{\lambda}^{\theta,g}$  as well. Thus, focus on the case  $\theta \in \Theta^{\max}$ , so that by Claim  $4 \overline{\lambda}^{\theta} > 0$ .

If  $\sum_{\theta'} \overline{\lambda}^{\theta'} q^{\theta'} = 1$ , then Eq. (20) and the fact that  $\sum_{\theta'} \lambda_{t-1}^{\theta'} q^{\theta'} \in [0,1]$  and  $0 \leq \lambda_{t-1}^{\theta,g} \leq \lambda_{t-1}^{\theta} \leq 1$  for all  $\theta$  imply that

$$\lambda_t^{\theta,g} = \left(1 - \sum_{\theta'} \lambda_{t-1}^{\theta'} q^{\theta'}\right) \lambda_{t-1}^{\theta,g} + \lambda_{t-1}^{\theta} q^{\theta,g} \in \left[\lambda_{t-1}^{\theta} q^{\theta,g}, 1 - \sum_{\theta'} \lambda_{t-1}^{\theta'} q^{\theta'} + \lambda_{t-1}^{\theta} q^{\theta,g}\right]$$

and both endpoints of the interval in the r.h.s. converge to  $\bar{\lambda}^{\theta}q^{\theta,g}$  by Claim 4 if  $\sum_{\theta'} \bar{\lambda}^{\theta'}q^{\theta'} = 1$ . Furthermore, the same assumption implies that  $\bar{\lambda}^{\theta}q^{\theta,g} = \bar{\lambda}^{\theta,g}$ , so  $\lambda_t^{\theta,g} \to \bar{\lambda}^{\theta,g}$ .

Now consider the case  $0 < \sum_{\theta'} \bar{\lambda}^{\theta'} q^{\theta'} < 1$ . (The set  $\Theta^{\max}$  is non-empty, and since  $q \in \mathbb{R}^{\Theta}_+ \setminus \{0\}$ , there is  $\theta^+ \in \Theta^{\max}$  with  $q^{\theta^+} > 0$ ; by Claim 4,  $\bar{\lambda}^{\theta'} > 0$  for  $\theta' \in \Theta^{\max}$ , so in particular  $\bar{\lambda}^{\theta^+} > 0$ ; but then  $\sum_{\theta'} \bar{\lambda}^{\theta'} q^{\theta'} \ge \bar{\lambda}^{\theta^+} q^{\theta^+} > 0$ .) It is convenient to let  $q_t = \sum_{\theta'} \lambda_t^{\theta'} q^{\theta'}$  and  $\bar{q} = \sum_{\theta'} \bar{\lambda}^{\theta'} q^{\theta'} = \lim_{t \to \infty} q_t$ , where the second equality follows from Claim 4. Thus, Eq. (20) can be written as

$$\lambda_t^{\theta,g} = (1 - q_{t-1})\lambda_{t-1}^{\theta,g} + \lambda_{t-1}^{\theta}q^{\theta,g}.$$
(26)

In addition,  $\bar{q} \in (0, 1)$ .

We claim that, for all  $T \ge 0$  and t > T,

$$\lambda_t^{\theta,g} = \lambda_T^{\theta,g} \prod_{s=T}^{t-1} (1-q_s) + q^{\theta,g} \sum_{s=T}^{t-1} \lambda_s^{\theta} \prod_{r=s+1}^{t-1} (1-q_r).$$
(27)

For t = T + 1, this follows from Eq. (26). Inductively, assume it holds for t - 1 > T. Then, by Eq. (26) and the inductive hypothesis,

$$\begin{split} \lambda_t^{\theta,g} &= (1-q_{t-1}) \left[ \lambda_T^{\theta,g} \prod_{s=T}^{t-2} (1-q_s) + q^{\theta,g} \sum_{s=T}^{t-2} \lambda_s^{\theta} \prod_{r=s+1}^{t-2} (1-q_r) \right] + \lambda_{t-1}^{\theta,g} q^{\theta,g} = \\ &= \lambda_T^{\theta,g} \prod_{s=T}^{t-1} (1-q_s) + q^{\theta,g} \sum_{s=T}^{t-1} \lambda_s^{\theta} \prod_{r=s+1}^{t-1} (1-q_r), \end{split}$$

as claimed.

Fix  $\epsilon > 0$  such that  $\bar{\lambda}^{\theta} - \epsilon > 0$ ,  $\bar{q} - \epsilon > 0$ ,  $1 - \bar{q} + \epsilon < 1$ , and  $1 - \bar{q} - \epsilon > 0$ . This is possible because  $\bar{\lambda}^{\theta} > 0$  and  $\bar{q} \in (0, 1)$ , hence  $1 - \bar{q} \in (0, 1)$ .

Since  $\lambda_t^{\theta} \to \bar{\lambda}^{\theta}$  and  $q_t \to \bar{q}$ , there is  $T \ge 0$  such that, for all t > T,  $\lambda_t^{\theta} < \bar{\lambda}^{\theta} + \epsilon$  and  $q_t > \bar{q} - \epsilon$ . Hence, for such t > T, Eq. (27) implies that

$$\begin{split} \lambda_t^{\theta,g} \leq &\lambda_T^{\theta,g} \prod_{s=T}^{t-1} (1-\bar{q}+\epsilon) + q^{\theta,g} \sum_{s=T}^{t-1} (\bar{\lambda}^{\theta}+\epsilon) \prod_{r=s+1}^{t-1} (1-\bar{q}+\epsilon) = \\ = &\lambda_T^{\theta,g} (1-\bar{q}+\epsilon)^{t-T} + q^{\theta,g} (\bar{\lambda}^{\theta}+\epsilon) \sum_{s=T}^{t-1} (1-\bar{q}+\epsilon)^{t-1-s} = \\ = &\lambda_T^{\theta,g} (1-\bar{q}+\epsilon)^{t-T} + q^{\theta,g} (\bar{\lambda}^{\theta}+\epsilon) \sum_{s=0}^{t-1-T} (1-\bar{q}+\epsilon)^s = \\ = &\lambda_T^{\theta,g} (1-\bar{q}+\epsilon)^{t-T} + q^{\theta,g} (\bar{\lambda}^{\theta}+\epsilon) \frac{1-(1-\bar{q}+\epsilon)^{t-T}}{\bar{q}-\epsilon} \to \frac{q^{\theta,g} (\bar{\lambda}^{\theta}+\epsilon)}{\bar{q}-\epsilon} \end{split}$$

This implies that  $\limsup_t \lambda_t^{\theta,g} \leq \frac{q^{\theta,g}(\bar{\lambda}^\theta + \epsilon)}{\bar{q} - \epsilon}$ . Since this must hold for all  $\epsilon > 0$ , it must be that  $\limsup_t \lambda_t^{\theta,g} \leq \frac{q^{\theta,g}\bar{\lambda}^\theta}{\bar{q}} = \bar{\lambda}^{\theta,g}$ .

Similarly,  $\lambda_t^{\theta} \to \bar{\lambda}^{\theta}$  and  $q_t \to \bar{q}$  imply that there is  $T \ge 0$  such that, for all t > T,  $\lambda_t^{\theta} > \bar{\lambda}^{\theta} - \epsilon > 0$  and  $q_t < \bar{q} + \epsilon < 1$ . Then

$$\begin{split} \lambda_t^{\theta,g} \geq &\lambda_T^{\theta,g} \prod_{s=T}^{t-1} (1-\bar{q}-\epsilon) + q^{\theta,g} \sum_{s=T}^{t-1} (\bar{\lambda}^{\theta}-\epsilon) \prod_{r=s+1}^{t-1} (1-\bar{q}-\epsilon) = \\ = &\lambda_T^{\theta,g} (1-\bar{q}-\epsilon)^{t-T} + q^{\theta,g} (\bar{\lambda}^{\theta}-\epsilon) \frac{1-(1-\bar{q}-\epsilon)^{t-T}}{\bar{q}+\epsilon} \to \frac{q^{\theta,g} (\bar{\lambda}^{\theta}-\epsilon)}{\bar{q}+\epsilon}, \end{split}$$

so  $\liminf_t \lambda_t^{\theta,g} \geq \frac{q^{\theta,g}(\bar{\lambda}^{\theta}-\epsilon)}{\bar{q}+\epsilon}$ . Again, since this must hold for all  $\epsilon > 0$ ,  $\liminf_t \lambda_T^{\theta,g} \geq \frac{q^{\theta,g}\bar{\lambda}^{\theta}}{\bar{q}} = \bar{\lambda}^{\theta,g}$ . Hence,  $\lambda_t^{\theta,g} \to \bar{\lambda}^{\theta,g}$ . Q.E.D.

**Proof of Propositions 2 and 3:** by assumption,  $\gamma^{\theta}(p^{\theta,m} + p^{\theta,g} \leq 1 \text{ for every } \theta$ . Hence, part (i) of Proposition 2 and Proposition 3 follow from Theorem 1 parts 4 and 5, by setting  $q^{\theta,g} = \gamma^{\theta}p^{\theta,g}$  for g = m, f and  $\theta \in \Theta$ . For part (ii) of Proposition 2, note that, if  $\theta \in \arg \max_{\theta' \in \Theta} \gamma^{\theta'}(p^{\theta',m} + p^{\theta',f})$ , then  $\gamma^{\sigma(\theta)} = \gamma^{\theta}$  and  $p^{\sigma(\theta),m} + p^{\sigma(\theta),f} = p^{\theta,f} + p^{\theta,m}$  imply that also  $\sigma(\theta) \in \arg \max_{\theta' \in \Theta} \gamma^{\theta'}(p^{\theta',m} + p^{\theta',f})$ . *Q.E.D.* 

**Proof sketch of Proposition 4** (see Online Appendix A2. for a detailed proof): write

$$\bar{\Lambda}^g = \sum_{\theta \in \theta^{\max}: \sigma(\theta) = \theta} \bar{\lambda}^{\theta,g} + \sum_{\theta \in \Theta^{\max}: \sigma(\theta) \neq \theta} \bar{\lambda}^{\theta,g} = \sum_{\theta \in \theta^{\max}: \sigma(\theta) = \theta} \bar{\lambda}^{\theta,g} + \frac{1}{2} \sum_{\theta \in \Theta^{\max}: \sigma(\theta) \neq \theta} (\bar{\lambda}^{\theta,g} + \bar{\lambda}^{\sigma(\theta),g}).$$

By the properties of symmetric types, if  $\Theta^{\max}$  is homogeneous then the above expression is the same for g = m, f, and so necessarily  $\bar{\Lambda}^m = \bar{\Lambda}^f = \frac{1}{2}$ . Otherwise, for at least one  $\theta$ , we have  $\theta, \sigma(\theta) \in \Theta^{\max}$  and  $p^{\theta,m} \neq p^{\sigma(\theta),m}$ , and it is wlog to assume that  $p^{\theta,m} > p^{\sigma(\theta),m}$ . Corollary 1 shows that, for such  $\theta, \sigma(\theta), \bar{\lambda}^{\theta,m} + \bar{\lambda}^{\sigma(\theta),m} > \bar{\lambda}^{\theta,f} + \bar{\lambda}^{\sigma(\theta),f}$ . Therefore,  $\bar{\Lambda}^m > \bar{\Lambda}^f$ , which implies that  $\bar{\Lambda}^m > \frac{1}{2}$ .

**Proof sketch of Proposition 5** (see Online Appendix A2. for a detailed proof): we first show that, for all  $\theta$  and t > 0,  $p^{\theta,m} \ge p^{\sigma(\theta),m}$  iff  $\lambda_{t-1}^{\theta} \ge \lambda_{t-1}^{\sigma(\theta)}$ . For t = 1, this is by assumption. For t > 1, this follows by inductively invoking part 3(a) of Theorem 1. By direct calculation, this in turn implies that, that for every  $\theta \in \Theta$  and  $t \ge 1$ ,  $a_t^{\theta,m} + a_t^{\sigma(\theta),m} \ge a_t^{\theta,f} + a_t^{\sigma(\theta),f}$ . The argument is completed by showing that, for g = m, f, we can write

$$\sum_{\theta:\gamma^{\theta}=\bar{\gamma}} a_t^{\theta,g} = \sum_{\theta:\gamma^{\theta}=\bar{\gamma},\theta=\sigma(\theta)} a_t^{\theta,g} + \frac{1}{2} \sum_{\theta:\gamma^{\theta}=\bar{\gamma},\theta\neq\sigma(\theta)} [a_t^{\theta,g} + a_t^{\sigma(\theta),g}].$$

**Proof sketch of Proposition 6** (see Online Appendix A2. for a detailed proof): as in the proof of Proposition 5, we have  $a_t^{\theta,m} + a_t^{\theta',m} > a_t^{\theta,f} + a_t^{\theta',f}$  for the intermediate types  $\theta, \theta'$ . On the other hand,  $a_t^{\theta_{0},m} = a_t^{\theta_{0},f}$  and  $a_t^{\theta_{1},m} = a_t^{\theta_{1},f}$  for the highest and lowest types  $\theta_0, \theta_1$ . This implies that the weight on the intermediate quality  $\gamma^{\theta} = \gamma^{\theta'}$  is higher for accepted M researchers. We then show inductively that  $a_t^{\theta_{1},g} > a_t^{\theta_{0},g}$  for both g = m, f. A direct calculation and comparison of the expressions for the expected qualities  $E[\gamma|M]$  and  $E[\gamma|F]$ yields the result.

## References

- Dennis J. Aigner and Glen G. Cain. Statistical theories of discrimination in labor markets. *ILR Review*, 30(2):175–187, 1977.
- Steffen Andersen, Seda Ertac, Uri Gneezy, John A List, and Sandra Maximiano. Gender, competitiveness, and socialization at a young age: Evidence from a matrilineal and a patriarchal society. *Review of Economics and Statistics*, 95(4):1438–1443, 2013.

- Peter Andre and Martin Falk. What's worth knowing? economists' opinions about economics. ECONtribute Discussion Paper 102, University of Bonn and University of Cologne, Reinhard Selten Institute (RSI), Bonn and Cologne, 2021. URL http://hdl.handle.net/ 10419/237347.
- Emmanuelle Auriol, Guido Friebel, Alisa Weinberger, and Sascha Wilhem. Women in economics: Europe and the world. mimeo, Toulose School of Economics, January 2022.
- Arjada Bardhi, Yingni Guo, and Bruno Strulovici. Spiraling or self-correcting discrimination: A multi-armed bandit approach. Technical report, Technical report, Northwestern University, 2019.
- Amanda Bayer and Cecilia Elena Rouse. Diversity in the economics profession: A new attack on an old problem. *Journal of Economic Perspectives*, 30(4):221–42, 2016.
- Gary S Becker. The economics of discrimination. University of Chicago press, 1957.
- Michael Betz, Lenahan O'Connell, and Jon M Shepard. Gender differences in proclivity for unethical behavior. *Journal of Business Ethics*, 8(5):321–324, 1989.
- Francine D. Blau and Lawrence M. Kahn. The gender wage gap: Extent, trends, and explanations. *Journal of Economic Literature*, 55(3):789–865, 2017.
- J Aislinn Bohren, Alex Imas, and Michael Rosenberg. The dynamics of discrimination: Theory and evidence. *American economic review*, 109(10):3395–3436, 2019.
- Lex Borghans, Bart H.H. Golsteyn, James J. Heckman, and Huub Meijers. Gender differences in risk aversion and ambiguity aversion. *Journal of the European Economic Association*, 7(2-3):649–658, 2009.
- David Card, Stefano DellaVigna, Patricia Funk, and Nagore Iriberri. Are referees and editors in economics gender neutral? *Quarterly Journal of Economics*, 135:269–327, February 2020.
- Anusha Chari and Paul Goldsmith-Pinkham. Gender representation in economics across topics and time: Evidence from the nber summer institute. Technical report, Working Paper, Yale University, 2018.
- Judy Chevalier. Report: committee on the status of women in the economics profession. Technical report, American Economic Association, 2020.
- Jacob Cohen. Statistical power analysis for the behavioral sciences. Routledge, 2013.
- J. Ignacio Conde-Ruiz, Juan José Ganuza, and Paola Profeta. Statistical discrimination and committees. mimeo, Universitat Pompeu Fabra, December 2020.
- John P. Conley and Ali Sina Önder. The research productivity of new phds in economics: The surprisingly high non-success of the successful. *Journal of the Economic Perspectives*, 28(3):205–216, 2014.
- Bradford Cornell and Ivo Welch. Culture, information, and screening discrimination. *Journal* of Political Economy, 104(3):542–571, 1996.

- Paul T Costa, Antonio Terracciano, and Robert R McCrae. Gender differences in personality traits across cultures: robust and surprising findings. *Journal of Personality and Social Psychology*, 81(2):322, 2001.
- Faye J Crosby, Aarti Iyer, Susan Clayton, and Roberta A Downing. Affirmative action: Psychological data and the policy debates. *American Psychologist*, 58(2):93, 2003.
- Rachel Croson and Uri Gneezy. Gender differences in preferences. *Journal of Economic literature*, 47(2):448–74, 2009.
- Marcus Dittrich and Kristina Leipold. Gender differences in time preferences. *Economics Letters*, 122(3):413–415, 2014.
- Anna Dreber and Magnus Johannesson. Gender differences in deception. *Economics Letters*, 99(1):197–199, 2008.
- David Dunning, Judith A Meyerowitz, and Amy D Holzberg. Ambiguity and self-evaluation: The role of idiosyncratic trait definitions in self-serving assessments of ability. *Journal of Personality and Social Psychology*, 57(6):1082, 1989.
- David Dunning, Marianne Perie, and Amber L Story. Self-serving prototypes of social categories. *Journal of Personality and Social Psychology*, 61(6):957, 1991.
- Pascaline Dupas, A Modestino, Muriel Niederle, and Justin Wolfers. Gender and the dynamics of economics seminars. mimeo, February 2021.
- Brian Fabo, Martina Jancokova, Elisabeth Kempf, and Lubos Pastor. Fifty shades of qe: Conflicts of interest in economic research. Technical report, University of Chicago, 2020.
- Armin Falk, Anke Becker, Thomas Dohmen, Benjamin Enke, David Huffman, and Uwe Sunde. Global evidence on economic preferences. *The Quarterly Journal of Economics*, 133(4):1645–1692, 2018.
- Hanming Fang and Andrea Moro. Theories of statistical discrimination and affirmative action: A survey. In *Handbook of social economics*, volume 1, pages 133–200. Elsevier, 2011.
- First Round 2022. d&i ef-Review, Eight ways to make your forts less talk and more walk. https://review.firstround.com/ eight-ways-to-make-your-dandi-efforts-less-talk-and-more-walk, 2022. Accessed: 2022-05-24.
- Donna K Ginther, Janet Currie, Francine D Blau, and Rachel Croson. Can mentoring help female assistant professors in economics? an evaluation by randomized trial. Technical report, NBER, March 2020. Working Paper 26864.
- Claudia Goldin. A grand gender convergence: Its last chapter. *American Economic Review*, 104(4):1091–1119, 2014.
- Claudia Goldin and Cecilia Rouse. Orchestrating impartiality: The impact of blind" auditions on female musicians. *American Economic Review*, 90(4):715–741, 2000.

- Luigi Guiso, Ferdinando Monte, Paola Sapienza, and Luigi Zingales. Culture, gender, and math. *Science*, 320(5880):1164, 2008.
- Paul Heidhues, Botond Kőszegi, and Philipp Strack. Overconfidence and prejudice. arXiv preprint arXiv:1909.08497, 2019.
- Thomas Hill, Nancy D Smith, and Hunter Hoffman. Self-image bias and the perception of other persons' skills. *European Journal of Social Psychology*, 18(3):293–298, 1988.
- Janet Shibley Hyde. Gender similarities and differences. *Annual Review of Psychology*, 65: 373–398, 2014.
- Janet Shibley Hyde and Marcia C. Linn. Gender similarities in mathematics and science. Science, 314(5799):599–600, 2006.
- Edward P. Lazear and Sherwin Rosen. Male-female wage differentials in job ladders. *Journal* of Labor Economics, 8(1, Part 2):S106–S123, 1990.
- Pawel Lewicki. Self-image bias in person perception. Journal of Personality and Social Psychology, 45(2):384, 1983.
- Shelly J Lundberg and Richard Startz. Private discrimination and social intervention in competitive labor market. *The American Economic Review*, 73(3):340–347, 1983.
- Shelly J Lundberg and Jenna Stearns. Women in economics: Stalled progress. *The Journal* of *Economic Perspectives*, 33(1):3–22, 2019.
- Thomas Mayer. Honesty and integrity in economics. Technical report, University of California at Davis, 2009. Working Paper 09-2.
- Edmund S Phelps. The statistical theory of discrimination. *American Economic Review*, 62 (4):659–661, 1972.
- Heather Sarsons. Recognition for group work: Gender differences in academia. American Economics Review: Papers and Proceedings, 107:141–45, 2017.
- Heather Sarsons, Klarita Gërxhani, Ernesto Reuben, and Arthur Schram. Gender differences in recognition for group work. *Journal of Political Economy*, 129, 2021.
- Amber L Story and David Dunning. The more rational side of self-serving prototypes: The effects of success and failure performance feedback. *Journal of Experimental Social Psychology*, 34(6):513–529, 1998.
- Andrea Weber and Christine Zulehner. Competition and gender prejudice: Are discriminatory employers doomed to fail? Journal of the European Economic Association, 12(2): 492–521, 2014.