

Supplementary Material for

Differential retention contributes to racial/ethnic disparity in U.S. academia

Allison K. Shaw^{1*}, Chiara Accolla¹, Jeremy M. Chacón¹, Taryn L. Mueller¹,
Maxime Vaugeois¹, Ya Yang², Nitin Sekar³, Daniel E. Stanton¹

¹Department of Ecology, Evolution and Behavior, University of Minnesota-Twin Cities,
Saint Paul, MN 55108

²Department of Plant and Microbial Biology, University of Minnesota-Twin Cities, Saint
Paul, MN 55108

³Wildlife and Habitats Division, WWF India, New Delhi, Delhi 110003, India *To whom
correspondence should be addressed; E-mail: ashaw@umn.edu.

1 Data

We used three broad types of data from the National Science Foundation (NSF) in our work: (i) data on the structure of academia (number of scholars in each academic stage, time spent in each stage), (ii) data on the racial/ethnic composition of scholars at each stage, and (iii) data on the approximate age range of academics. Whenever there were multiple versions of the same data available for a given year (e.g., in different versions on the same report, or when classifications changed within a timeseries), we used the most recent data for a given year. We limited our analysis to the period 1991-2016 where almost all data was available (except for racial/ethnic data on postdoctoral researchers which was only available for 2010 onward).

1.1 Structural Data

The structural data we used consisted of timeseries of the number of bachelors and PhD degrees awarded, the number of enrolled graduate students, and the number of employed postdoctoral researchers, assistant professors and tenured professors, as well as estimates of the length of time spent as a graduate student, postdoctoral researcher, assistant professor and tenured professor. Data on the number of bachelors and PhD degrees came from the NSF reports on Science and Engineering Degrees (1), and Women, Minorities, and Persons with Disabilities (WMPD) (2), data on the number of graduate students and postdoctoral scholars came from the NSF Survey of Graduate Students and Postdoctorates in Science and Engineering (3), and data on the number of assistant and tenured professors came from the NSF report on Science and Engineering Indicators (4). The length of time in each stage came from the NSF report on Science and Engineering Indicators (5) for graduate students, the NSF report on Postdoc Participation of Science, Engineering, and Health Doctorate Recipients (6) for postdocs and the integrated data system Scientists and Engineers Statistical Data System (SESTAT) for faculty. The specific sources for all structural data are given in Table S1 and in “Data Report Details” Section below, and the timeseries of structural data are plotted in Figure S1. Missing data were linearly interpolated; for example, faculty data was only collected approximately every two years, and undergraduate data was missing for the year 1999 (data with interpolation given in Figure S2).

1.2 Race/Ethnicity Data

The racial/ethnicity data we used consisted of timeseries data for the number of earned bachelors degrees, enrolled graduate students, and employed postdoctoral researchers, assistant professors and tenured professors by race/ethnicity. From 1991 to around 2010 NSF used five groups for race/ethnicity: ‘White’, ‘Asian or Pacific Islander’, ‘Black’, ‘Hispanic’, and ‘Native American/Alaskan Native’ (plus an additional group for unknown).

Around 2010, the group ‘Asian or Pacific Islander’ was split into ‘Asian’ and ‘Native Hawaiian or Other Pacific Islander’. At the same time, the group ‘More than one race’ was added. When the number of individuals in a group was quite small (this occurred for both Native American / Alaskan Native and Native Hawaiian / Pacific Islander in both assistant professor and tenured professor stages in some years) the specific number of individuals was masked instead of being reported. In these cases, we guesstimated the number of individuals from other group data. For example, if the total number of individuals of a race/ethnicity was reported for faculty as a whole, we split this number evenly among groups to approximate the number of individuals of that race/ethnicity in each faculty stage. Data on the racial/ethnic composition of undergraduate and PhD students as well as assistant and tenured professors came from the WMPD reports (2). Data on postdoctoral researchers (2010 onward) came from NSF Surveys of Graduate Students and Postdoctorates in Science and Engineering (3), and data prior to 2010 was estimated as the average of representation in the graduate student and assistant professor stages. The student data in the NSF WMPD reports only includes racial/ethnicity data for U.S. citizens and permanent residents. To account for international students, we used the NSF reports on Doctorate Recipients from U.S. Universities (7) for data on the proportion of permanent vs temporary resident PhD recipients and the racial/ethnic composition of temporary resident PhD recipients. The specific sources for all race/ethnicity data are given in Table S2 and in the “Data Report Details” Section below, and the timeseries of race/ethnicity data are plotted in Figures S3, S4, and S5. Count data on the number of scholars of each racial/ethnic group were converted to proportions and data were smoothed with a 5-year window moving average.

The specific number of individuals reported for each race/ethnicity group was not necessarily representative of the actual number of individuals of that race/ethnicity, for two main reasons. First, some individuals did not report their race/ethnicity (often reported as a separate group, ‘unknown’). Second, race/ethnicity data for undergraduate and graduate students was only provided for U.S. citizens and permanent residents; race/ethnicity for temporary residents was not recorded. However, race/ethnicity data for U.S. temporary residents was recorded for graduating PhD students (see Figure S5). Thus, when applying the race/ethnicity data, we used the proportion of individuals of each race/ethnicity rather than the actual count data (plotted in Figure S4). We calculated proportions using only data for a known race/ethnicity (i.e., we excluded the ‘unknown race’ group). For example, if there were 500 individuals in a stage, of which 150 were White, and 50 Asian, and 300 unknown race/ethnicity, we recorded this stage as being 0.75 White and 0.25 Asian. Finally data were smoothed with the ‘smoothdata’ function in Matlab, using a moving average over a window of size 5 years and omitting missing data.

1.3 Age Data

Finally, we used NSF data on the approximate age range of scholars at each stage by pulling data from the integrated data system SESTAT (Scientists and Engineers Statistical Data System, <https://www.nsf.gov/statistics/sestat/>), and determining the most representative ages of each stage. We selected the National Survey of Recent College Graduates (NSRCG) for undergraduate and graduate stages (year 2010), and the Survey of Doctorate Recipients (SDR) for postdoc, assistant and tenured professor stages (year 2015). For undergraduate and graduate students we created a table showing the most recent degree type (labeled "M_ED_MR_DEGREE_TYPE") in function of ages ("U_DEM_AGE_RCG_PUB"), and specified the population by the field of study for the most recent degree ("M_ED_MR_MAJOR_ED_GRP_MAJOR_NEW"). We selected the fields (i) biological, agricultural and environmental life sciences, (ii) physical Sciences, (iii) computer and mathematical sciences, and (iv) engineering. The total number of scholars per age class in the undergraduate stage was calculated as the sum of Bachelor and Master degrees across the four fields. Similarly, the total number of graduate scholars was obtained by summing up the number of doctorate degrees in each field. Then, we plotted the total number of undergraduate and graduate scholars in function of age, and selected the most representative time spent in each of these two stages. We applied the same method for the other three stages. Notably, we created a table considering the academic position of postdoc ("E_JOB_EMPLR_ACAD_POSITION_POSTDOC") or tenure status ("E_JOB_EMPLR_EDUC_INST_TENURE_STAT"), in function of ages grouped by 5-year intervals ("U_DEM_AGE_GROUP_5_YR_GROUPING_PUB"), and specified the population by the field of study for the highest degree ("O_ED_HD_MAJOR_ED_GRP_MAJOR_NEW"). Overall, the age ranges we used were: 15 to 24 years old (undergraduate students), 20 to 29 (graduate students), 25 to 39 (Ph.D. recipients), 25 to 44 (postdoctoral researchers), 30 to 49 (assistant professors) and 35 to 59 (tenured professors). We used this data to determine which subset of the general population we should compare each academic stage to.

Next, we determined the racial composition of the age class corresponding to each academic stage based on data from the National Center for Health Statistics and the U. S. Census Bureau (8). To compute our racial composition by academic stage for the "American Indian/Alaska Native", "Asian", "Black/African American", "White", and "Hispanic/Latino" categories from 1990 to 2016, we compiled estimates of resident population of the US by year, single-year of age, bridged-race category, and Hispanic origin produced by the National Center for Health Statistics under a collaborative arrangement with the U. S. Census Bureau (8). We compiled similar data from 2000 to 2016 for the "Native Hawaiian/Pacific Islander" and "Two or More Races" categories from the 2019 Population Estimates by Age, Sex, Race and Hispanic Origin (9) and National Intercensal Tables: 2000-2010, both produced by the U. S. Census Bureau (10).

See below for how each dataset was used.

2 Model structure

We constructed a model of academia as a series of stages with movement between them or out of the system (Figure 1 main text; modified from (11)). Our model has five discrete stages: undergraduate studies (U), graduate studies (G), postdoctoral fellowships (P), assistant professorships (A ; tenure-track) and tenured professorships (T). Individuals that move out of each stage either move up and fill empty positions in the next stage, or move out of the system.

We generated the structure of our model from NSF data. The number of academics in the U.S. has changed over time, so we set the size of each stage i in each year t ($N_i(t)$) from data on the actual number graduate students, postdoctoral fellows, assistant professors, and tenured professors, using data from NSF reports (see Table S1 for specific data sources, and Figure S2 for data). We used the time spent in each stage to estimate a turnover rate for that stage which, in combination with the number of scholars in each stage, gave us an estimate of the number of scholars that would have either transitioned from one stage to the next or transitioned outside of the system for each year.

2.1 Estimating Transitions

Each of the transitions was estimated as follows (see Figure S6 for results). For each year and each stage, we estimated the number of individuals leaving each stage based on transition rates and changes in stage sizes as

$$\rho_i(t) = \left(\frac{1}{\tau_i} \right) N_i(t) \quad (\text{S1})$$

for stages $i = \{G, P, A, T\}$ where τ_i is the average number of years spent in stage i (see Table S3 for all model parameters). Simultaneously, we estimated the number of openings in each year and each stage as

$$\omega_i(t) = N_i(t+1) - N_i(t) + \rho_i(t) . \quad (\text{S2})$$

In most cases, $\rho_i(t) \geq \omega_{i+1}(t)$, that is, the number of openings in stage $i+1$ could easily be filled by individuals leaving stage i . Thus, we partitioned individuals leaving stage i ($\rho_i(t)$) into those moving up to the next stage,

$$\mu_i(t) = \omega_{i+1}(t) \quad (\text{S3a})$$

and those leaving the system,

$$\lambda_i(t) = \rho_i(t) - \omega_{i+1}(t) . \quad (\text{S3b})$$

However, there were two other scenarios that occasionally occurred. First, when $\omega_i(t) < 0$ (i.e., too few individuals were leaving stage i than possible, given the change in stage from year to year), we adjusted $\rho_i(t)$ as

$$\rho_i(t) = N_i(t+1) - N_i(t) \quad (\text{S4a})$$

in order to make $\omega_i(t)$ non-negative,

$$\mu_{i-1}(t) = \omega_i(t) = 0 . \quad (\text{S4b})$$

Second, when $\rho_i(t) < \omega_{i+1}(t)$ (i.e., too few individuals were leaving stage i to fill openings in stage $i+1$), we either increased the number of individuals leaving stage i when possible, or else assumed the remaining openings were filled by individuals from outside the system being modeled (e.g., coming from other scientific disciplines or returning to academia after having previously left).

2.2 Estimation Details: A to T transition

Each year within each simulation was run over time according to the following steps.

First, we estimated the number of retiring tenured professors by eqn. (S1) with $i = T$. We estimated the number of assistant professors needed to fill these tenured slots by eqn. (S2) with $i = T$. If this was a negative number of assistant professors, we adjusted it according to eqn. (S4). We estimated the number of assistant professors being tenured (and thus available to fill T slots) by eqn. (S1) with $i = A$. If $\rho_A(t) < \omega_T(t)$ (i.e., too few assistant professors were estimated as being tenured), we adjusted $\rho_A(t)$ as

$$\rho_A(t) = \omega_T(t) \quad (\text{S5})$$

i.e., assuming that more assistant professors were tenured than initially estimated. We did not pull individuals from outside the system at this transition as it these seem likely to be rare (e.g, that an individual transitions from an assistant professor in one field to a tenured professor in a different field, or from a non-academic position to a tenured position). Otherwise, if $\rho_A(t) \geq \omega_T(t)$ we used eqn. (S3) to estimate the transition rates with $i = A$. This method effectively assumes that the rate individuals move from A to T is driven by the rate tenured professors retire (ρ_T) and that any ‘excess’ assistant professors receiving tenure leave the system. We refer to this as a ‘demand’ view of faculty turnover (‘demand’ in terms of empty T slots determines the A to T transition).

We thus consider a second alternative scenario, what we call a ‘supply’ view of faculty turnover, where ‘supply’ in terms of assistant professors receiving tenure determines the A to T transition. For this method, we estimated the number of retiring tenured professors by eqn. (S1) with $i = T$, estimated the number of assistant professor receiving tenure by eqn. (S1) with $i = A$, and calculated the change in the T stage as

$$\Delta_T(t) = N_T(t+1) - N_T(t) . \quad (\text{S6})$$

If $\Delta_T(t) > \rho_A(t)$ (i.e., too few assistant professors were estimated as being tenured to fill the minimum number of T slots), we adjusted $\rho_A(t)$ as

$$\rho_A(t) = \Delta_T(t) \quad (\text{S7})$$

i.e., assuming that more assistant professors were tenured than initially estimate, and set $\rho_T(t) = 0$ (no tenured professors retire this year). Otherwise, if $\Delta_T(t) < \rho_A(t)$, we set

$$\rho_T(t) = \rho_A(t) - \Delta_T(t) , \quad (\text{S8})$$

i.e., that retirement of T is assumed to exactly offset the number of A being tenured, minus the new T slots that have become available.

The ‘demand’ scenario likely overestimates the number of faculty receiving tenure and then leaving academia, while the ‘supply’ scenario likely underestimates the number of faculty leaving academia after tenure and before retirement. We run simulations under both scenarios to serve as upper and lower bounds.

2.3 Estimation Details: P to A transition

Next, we estimated the number of postdoctoral researchers needed to fill these assistant professor slots by eqn. (S2) with $i = A$. We estimated the number of postdoctoral researchers available to fill A slots by eqn. (S1) with $i = P$. If $\rho_P(t) < \omega_A(t)$ (i.e., too few postdoctoral researchers were estimated as being available), we adjusted $\rho_P(t)$ as

$$\rho_P(t) = \omega_A(t) \quad (\text{S9})$$

i.e., assuming that more postdoctoral researchers were hired than initially estimated. Otherwise, if $\rho_P(t) \geq \omega_A(t)$ we used eqn. (S3) to estimate the transition rates with $i = P$.

2.4 Estimation Details: G to P transition

Next, we estimated the number of graduate students needed to fill these postdoc slots by eqn. (S2) with $i = P$. We estimated the number of graduate students leaving that stage by eqn. (S1) with $i = G$. We assumed that only students leaving the G stage with a PhD degree can fill the P slots, so we estimated the number of graduate students available to fill P slots by $D_G(t)$, the number of PhD degrees granted in year t . If $D_G(t) \geq \omega_P(t)$ we used a modified version of eqn. (S3) to estimate the transition rates where individuals moving from stage G to stage P as

$$\mu_G(t) = \omega_P(t) , \quad (\text{S10a})$$

those leaving the system with a PhD degree as

$$\lambda_G(t) = D_G(t) - \omega_P(t) , \quad (\text{S10b})$$

and those leaving stage G before their degree as

$$\delta_G(t) = \rho_G(t) - D_G(t) . \quad (\text{S10c})$$

2.5 Estimation Details: U to G transition

Finally, we estimated the number of undergraduate students needed to fill these graduate student slots by eqn. (S2) with $i = G$. We assumed that only students leaving the U stage with a degree can fill the G slots, so we estimated the number of undergraduate students available to fill G slots by $D_U(t)$, the number of undergraduate degrees granted in year t . We used a modified version of eqn. (S3) to estimate the transition rates where individuals moving from stage U to stage G as

$$\mu_U(t) = \omega_G(t) \tag{S11a}$$

and those leaving the system with an undergraduate degree as

$$\lambda_U(t) = D_U(t) - \omega_G(t) . \tag{S11b}$$

2.6 Estimation Details: Subpartitions

Since ethnic/racial composition may vary within each stage (especially for longer career stages), we split some stages into sub-partitions. This enabled us to model different racial/ethnic compositions for each sub-partition within a stage. This also ensured that when individuals were moved out of a partitioned stage, they were taken from the oldest sub-partition. We split the graduate student stage into two sub-partitions and split the tenured professor stage into five sub-partitions. We assumed that graduate students spent 3 years in the first sub-partition (approximately until qualifying exams) and then $\tau_G - 3$ in the second partition. We assumed that tenured professors spent $\tau_T/5$ in each of the five sub-partitions. Transitions between sub-partitions were estimated based on turnover time. Graduate students leaving the system before receiving a degree ($\delta_G(t)$) were pulled from both sub-partitions (half from each), but graduate students leaving the stage with a doctoral degree were assumed to come only from the second sub-partition. Tenured professors retiring ($\rho_T(t)$) were pulled only from the last sub-partition.

3 Model simulations

3.1 Simulation details

With our model structure in place, we then simulated the flow of individuals through the system. We assumed that at each transition, the fraction of individuals staying in the system versus moving outside did not vary with race/ethnicity (i.e., individuals of different races were equally likely to stay in the system). Therefore individuals entering a given stage were drawn from the stage below in proportion to their representation in the lower stage. We calculated $n_i(t, k)$ the simulated number of individuals of each race/ethnicity k in each stage i over time t as follows. The initial number of individuals of each race/ethnicity was taken from National Science Foundation data in a starting year t_0 , except for the case of postdoctoral fellows where race/ethnicity data was not available before 2010. In this case, we assumed initial proportion for each race/ethnicity that was the average of the values for graduate students and assistant professors. See Table S2 for data sources.

The survey data for undergraduate degrees and enrolled graduate students only included race/ethnicity data for US citizens and permanent residents; temporary residents were reported as a separate category with no race/ethnicity data. Temporary residents make up a large proportion of graduate students and have a different racial/ethnic composition than US citizens and permanent residents (see Figure S5). In contrast, survey data for graduate degrees did report race/ethnicity data for all graduates across residency types, so we used this data at the transition point from graduate students G to postdoctoral researchers P . Race/ethnicity data by citizenship for PhD degrees was not available before 2000. However, data on race/ethnicity data by citizenship for the doctoral workforce was available for the years 1991 and 1993, and was thus used and interpolated to approximate race/ethnicity data for temporary resident PhD recipients between 1991 and 1999.

Next, for each year going forward, we fed in NSF data on racial/ethnic composition at a particular stage (e.g., undergraduate students), and used our model to predict the racial/ethnic composition at the other stages (e.g., graduate students). We calculated the number of individuals of race/ethnicity k in each stage in the next year ($t + 1$). The number of graduate students is given by

$$\begin{aligned}
 n_G(t + 1, k, 1) = & \underbrace{n_G(t, k, 1)}_{\text{initial}} - \underbrace{\beta_G(t)f_G(t, k, 1)}_{\text{move up}} - \underbrace{0.5\delta_G(t)f_G(t, k, 1)}_{\text{leave system}} \\
 & + \underbrace{\mu_U(t)f_U(t, k)}_{\text{move in}}
 \end{aligned} \tag{S12a}$$

for the first subpartition in G and

$$n_G(t+1, k, 2) = \underbrace{n_G(t, k, 2)}_{\text{initial}} - \underbrace{D_G(t)f_G(t, k, 2)}_{\text{graduate}} - \underbrace{0.5\delta_G(t)f_G(t, k, 2)}_{\text{leave system}} + \underbrace{\beta_G(t)f_G(t, k, 1)}_{\text{move in}} \quad (\text{S12b})$$

for the second subpartition in G , where $f_i(t, k)$ is the fraction of individuals of race/ethnicity k in stage i in year t and $\beta_G(t)$ is the number of individuals that move between G subpartitions in year t .

The number of postdoctoral researchers is given by

$$n_P(t+1, k) = \underbrace{n_P(t, k)}_{\text{initial}} - \underbrace{\mu_P(t)f_P(t, k)}_{\text{move up}} - \underbrace{\lambda_P(t)f_P(t, k)}_{\text{leave system}} + \underbrace{\mu_G(t)R(t)f_G(t, k, 2)}_{\text{move in (perm. res.)}} + \underbrace{\mu_G(t)(1-R(t))V(t, k)}_{\text{move in (temp. res.)}}. \quad (\text{S12c})$$

where $R(t)$ is the fraction of PhD degrees that go to U.S. citizens and permanent residents in year t (thus, $1-R(t)$ go to temporary residents), and $V(t, k)$ is the fraction of U.S. temporary resident PhD recipients of race/ethnicity k in year t . The number of assistant professors is given by

$$n_A(t+1, k) = \underbrace{n_A(t, k)}_{\text{initial}} - \underbrace{\mu_A(t)f_A(t, k)}_{\text{move up}} - \underbrace{\lambda_A(t)f_A(t, k)}_{\text{leave system}} + \underbrace{\mu_P(t)f_P(t, k)}_{\text{move in}}. \quad (\text{S12d})$$

The number of tenured professors is given by

$$n_T(t+1, k, 1) = \underbrace{n_T(t, k, 1)}_{\text{initial}} - \underbrace{\beta_T(t, 1)f_T(t, k, 1)}_{\text{move up}} + \underbrace{\mu_A(t)f_A(t, k)}_{\text{move in}} \quad (\text{S12e})$$

for the first subpartition in T ,

$$n_T(t+1, k, j) = \underbrace{n_T(t, k, j)}_{\text{initial}} - \underbrace{\beta_T(t, j)f_T(t, k, j)}_{\text{move up}} + \underbrace{\beta_T(t, j-1)f_T(t, k, j-1)}_{\text{move in}} \quad (\text{S12f})$$

for subpartitions 2 through 4 ($j = 2, 3, 4$) in T , and

$$n_T(t+1, k, 5) = \underbrace{n_T(t, k, 5)}_{\text{initial}} - \underbrace{\rho_T(t)f_T(t, k, 5)}_{\text{retire}} + \underbrace{\beta_T(t, 4)f_T(t, k, 4)}_{\text{move in}} \quad (\text{S12g})$$

for the last (fifth) partition in T .

3.2 Simulation scenarios

We considered four types of scenarios for our simulations (based on turnover rate and turnover type), which capture uncertainty in the details surrounding transitions for faculty in academia. Although we found NSF data on the average length of time spent as a PhD student and as a postdoctoral researcher (Table S1), we could not find similar data on the average time spent on the tenure-track or as a tenured professor. Instead, we considered (i) a ‘slow’ turnover within the faculty, estimating the time spent on the tenure-track (τ_A) as 8 years and the time spent as a tenured professor (τ_T) as 30 years, and (ii) a ‘fast’ turnover within the faculty, estimating τ_A as 5 years and τ_T as 20 years. We also considered that the rate individuals moved between the A and T stages was driven by (i) ‘supply’ (i.e., rate of A achieving tenure), and (ii) ‘demand’ (i.e., rate of T retiring). We thus considered four combinations of scenarios: fast-supply, fast-demand, slow-supply and slow-demand.

3.3 Simulation sets

We ran three sets of simulations, each run under the four scenarios described above.

First, to study the overall effects of retention (Figure 2 in the paper), we started simulations in year $t_0 = 1991$ and ran them for 25 years (the full range of available data), feeding NSF data on the race/ethnicity of graduating undergraduates, and simulating the expected race/ethnicity of graduate students, postdoctoral researchers, assistant professors, and tenured professors. We used five initial groups for race/ethnicity: ‘White’, ‘Asian or Pacific Islander’, ‘Black’, ‘Hispanic’, and ‘Native American/Alaskan Native’. Around 2010 (year differs slightly across academic stages), the group ‘Asian or Pacific Islander’ was split into ‘Asian’ and ‘Native Hawaiian or Other Pacific Islander’ in the NSF data and the group ‘More than one race’ was added. Accordingly, we adjusted the simulated individuals in our model starting in the year 2012 (the first year that these two new groups were available for all academic stages). We partitioned the simulated individuals in the ‘Asian or Pacific Islander’ group into ‘Asian’ and ‘Native Hawaiian or Other Pacific Islander’ groups based on the relative proportion of these two groups in the NSF data for 2012. Similarly, we set the proportion of simulated individuals in the ‘More than one race’ group based on the relative proportion of that group in the NSF 2012 data, and pulled these simulated individuals evenly from the other simulated groups.

Second, to isolate the effects of retention within each stage of academia (Figure 4 in the paper), we fed in NSF data on the race/ethnicity at each stage and quantified the expected outcome at the following stage. Specifically, we simulated expected results for graduate students based on our model run with NSF undergraduate student data, expected results for postdoctoral researchers and assistant professors based on NSF graduate student data, and expected results for tenured professors based on NSF assistant professor data. This second set of simulations was also started in the year $t_0 = 1991$, running them for 25

years.

Third, to examine how the effect of specific transitions within academia changed over time, we started simulations in different starting years ($t_0 = 1991, 1996, 2001, 2006$) and ran each simulation for 10 years. Here again we simulated expected results for each stage based on our model run with NSF data at the previous stage.

All simulations and calculations were done using *Matlab*.

3.4 Testing model predictions

Finally, we compared the racial/ethnic composition predicted by our null model to the actual composition from NSF data. We quantified this comparison with the metric

$$\theta = \frac{d_i(t, k) - f_i(t, k)}{f_i(t, k)} \quad (\text{S13})$$

where $d_i(t, k)$ and $f_i(t, k)$ are the NSF data and model prediction, respectively, of the proportion of stage i in year t that is made up of race/ethnicity k . Here, $\theta > 0$ indicates that a racial/ethnic group has higher representation in a stage than is predicted by the null model and $\theta < 0$ means lower representation than predicted.

We calculated confidence intervals around θ values, as follows. For each combination of $d_i(t, k)$ and $f_i(t, k)$, we considered what effect an error of $\epsilon = 5\%$ would have. We calculated four bounds to the θ metric:

$$\theta'_1 = \frac{(1 - \epsilon)d_i(t, k) - (1 - \epsilon)f_i(t, k)}{(1 - \epsilon)f_i(t, k)} \quad (\text{S14a})$$

$$\theta'_2 = \frac{(1 - \epsilon)d_i(t, k) - (1 + \epsilon)f_i(t, k)}{(1 + \epsilon)f_i(t, k)} \quad (\text{S14b})$$

$$\theta'_3 = \frac{(1 + \epsilon)d_i(t, k) - (1 - \epsilon)f_i(t, k)}{(1 - \epsilon)f_i(t, k)} \quad (\text{S14c})$$

$$\theta'_4 = \frac{(1 + \epsilon)d_i(t, k) - (1 + \epsilon)f_i(t, k)}{(1 + \epsilon)f_i(t, k)} \quad (\text{S14d})$$

and used the largest and smallest value of these four to set the upper and lower bounds of the confidence interval around the θ value.

3.5 Supplementary Results

In addition to the results in the main text, several supplementary results are included below. Table S4 provides numerical value of representation in each stage for each race/ethnicity. Figure S7 shows a comparison of representation comparing the model results, academia data and census data. Figure S8 shows the temporal trends in the θ metric value. Figure S9 shows a comparison of two model versions – one accounting for the race/ethnicity of temporary resident international scholars who receive their PhDs in the U.S., and one ignoring the race/ethnicity of this group.

4 Data Report Details

Below are details of each data source used.

[08-307] NSF Publication 08-307. 2008 National Science Foundation, Division of Science Resources Statistics, Postdoc Participation of Science, Engineering, and Health Doctorate Recipients. (<http://www.nsf.gov/statistics/infbrief/nsf08307>)

2008 report, Table 2: Median duration of most recently completed postdoc

[GSPD] Survey of Graduate Students and Postdoctorates in Science and Engineering. (<https://www.nsf.gov/statistics/srvygradpostdoc/>)

2018 report, Table 1-9a: number of graduate students by science field for 1975–2018

2018 report, Table 1-10a: number of graduate students by engineering field for 1975–2018

2018 report, Table 1-9b: number of postdoctoral researchers by science field for 1975–2018

2018 report, Table 1-10b: number of postdoctoral researchers by engineering field for 1975–2018

2010 report, Table 34: postdoctoral researchers, by race/ethnicity for 2010

2016 report, Table 34: postdoctoral researchers, by race/ethnicity for 2011–2016

2017 report, Table 2-2: postdoctoral researchers, by race/ethnicity for 2017

2018 report, Table 2-2: postdoctoral researchers, by race/ethnicity for 2018

[S&E Degrees] Science and Engineering Degrees: 1966–2012.

(<https://www.nsf.gov/statistics/2015/nsf15326/>)

2015 report, Table 5: number of bachelor's degrees by field for 1966–2012

2015 report, Table 19: number of PhD degrees by field for 1966–2012

[SED] Survey of Earned Doctorates.

(<https://www.nsf.gov/statistics/srvydoctorates/>)

2014 report, Table 17: doctorate recipients, by broad field of study and citizenship for 1984–2014 (every 5 years)

2015 report, Table 17: doctorate recipients, by broad field of study and citizenship for 1985–2015 (every 5 years)

2016 report, Table 17: doctorate recipients, by broad field of study and citizenship for 1986–2016 (every 5 years)

2017 report, Table 17: doctorate recipients, by broad field of study and citizenship for 1987–2017 (every 5 years)

2018 report, Table 17: doctorate recipients, by broad field of study and citizenship for 1988–2018 (every 5 years)

2010 report, Table 19: doctorate recipients, by race/ethnicity and citizenship for 2000–2010

2018 report, Table 19: doctorate recipients, by race/ethnicity and citizenship for 2009–2018

[SE-ind] Science and Engineering Indicators, National Science Board.

(<https://nces.nsf.gov/indicators>)

2019 report, Table S3-7: number of assistant and tenured professors by field for 1973-2017
2018 report, Table 2-3: median time to degree by field for 1985-2015

[WMPD] Women, Minorities, and Persons with Disabilities in Science and Engineering report.

(<https://www.nsf.gov/statistics/women/>)

2019 report, Table 5-3: number of bachelor's degrees by field for 2006-2016

2019 report, Table 7-4: number of PhD degrees by field for 2006-2016

1994 report, Table 5-19: bachelors degrees by race/ethnicity for 1981-1991

2002 report, Table 3-8: bachelors degrees by race/ethnicity for 1990-1998

2009 report, Table C6: bachelors degrees by race/ethnicity for 1996-2007

2019 report, Table 5-3: bachelors degrees by race/ethnicity for 2006-2016

2002 report, Table 4-6: graduate students by race/ethnicity for 1990-1999

2009 report, Table D-1: graduate students by race/ethnicity for 1999-2006

2011 report, Table 3-1: graduate students by race/ethnicity for 2008-2010

2013 report, Table 3-1: graduate students by race/ethnicity for 2012

2017 report, Table 3-1: graduate students by race/ethnicity for 2014

2019 report, Table 3-1: graduate students by race/ethnicity for 2016

1994 report, Table 8-11: PhD workforce by race/ethnicity and citizenship for 1991

1996 report, Table 5-33: PhD workforce by race/ethnicity and citizenship for 1993

1994 report, Table 8-18: faculty by race/ethnicity for 1991

1996 report, Table 5-28: faculty by race/ethnicity for 1993

1998 report, Table 5-10: faculty by race/ethnicity for 1995

2000 report, Table 5-19: faculty by race/ethnicity for 1997

2004 report, Table H-26: faculty by race/ethnicity for 2001

2007 report, Table H-28: faculty by race/ethnicity for 2003

2009 report, Table H-28: faculty by race/ethnicity for 2006

2011 report, Table 9-26: faculty by race/ethnicity for 2008

2013 report, Table 9-26: faculty by race/ethnicity for 2010

2015 report, Table 9-26: faculty by race/ethnicity for 2013

2017 report, Table 9-26: faculty by race/ethnicity for 2015

2019 report, Table 9-26: faculty by race/ethnicity for 2017

References

- [1] National Science Foundation, National Center for Science and Engineering Statistics. Science and Engineering Degrees: 1966–2012. Detailed Statistical Tables NSF 15-326. Arlington, VA; 2015. Available from: <http://www.nsf.gov/statistics/2015/nsf15326/>.
- [2] National Science Foundation, National Center for Science and Engineering Statistics. Women, Minorities, and Persons with Disabilities in Science and Engineering; 2019. Special Report NSF 19-304. Alexandria, VA; 2019. Available from: <https://www.nsf.gov/statistics/wmpd>.
- [3] National Center for Science and Engineering Statistics. Survey of Graduate Students and Postdoctorates in Science and Engineering; 2018. Available from: <http://ncesdata.nsf.gov/gradpostdoc/>.
- [4] National Science Board NSF. Higher Education in Science and Engineering. Science and Engineering Indicators 2020. NSB-2019-7.; 2019. NSB-2019-7. Available from: <https://nces.nsf.gov/pubs/nsb20197/>.
- [5] National Science Board. Science and Engineering Indicators 2018. NSB-2018-1. Alexandria, VA: National Science Foundation; 2018. Available from: <https://www.nsf.gov/statistics/indicators/>.
- [6] National Science Foundation, Division of Science Resources Statistics. Postdoc Participation of Science, Engineering, and Health doctorate Recipients; 2008. Available from: <http://www.nsf.gov/statistics/infbrief/nsf08307>.
- [7] National Center for Science and Engineering Statistics NSF. Doctorate Recipients from U.S. Universities: 2018. Special Report NSF 20-301. Alexandria, VA; 2019. Available from: <https://nces.nsf.gov/pubs/nsf20301/>.
- [8] United States Department of Health and Human Services (US DHHS), Centers for Disease Control and Prevention (CDC), National Center for Health Statistics (NCHS). Bridged-Race Population Estimates, United States July 1st resident population by state, county, age, sex, bridged-race, and Hispanic origin. Compiled from 1990-1999 bridged-race intercensal population estimates (released by NCHS on 7/26/2004); revised bridged-race 2000-2009 intercensal population estimates (released by NCHS on 10/26/2012); and bridged-race Vintage 2019 (2010-2019) postcensal population estimates (released by NCHS on 7/9/2020). Available on CDC WONDER Online Database. Accessed at <http://wonder.cdc.gov/bridged-race-v2019.html> on Oct 9, 2020; 2020.

- [9] Annual Estimates of the Resident Population by Sex, Age, Race, and Hispanic Origin for the United States: April 1, 2010 to July 1, 2019 (NC-EST2019-ASR6H). U.S. Census Bureau, Population Division. June 2020; 2020.
- [10] Intercensal Estimates of the Resident Population by Sex, Race, and Hispanic Origin for the United States: April 1, 2000 to July 1, 2010 (US-EST00INT-02). U.S. Census Bureau, Population Division. September 2011.; 2011.
- [11] Shaw AK, Stanton DE. Leaks in the pipeline: separating demographic inertia from ongoing gender differences in academia. *Proceedings of the Royal Society of London Series B: Biological Sciences*. 2012;279(1743):3736–3741. Available from: <http://rspb.royalsocietypublishing.org/content/279/1743/3736.short>.

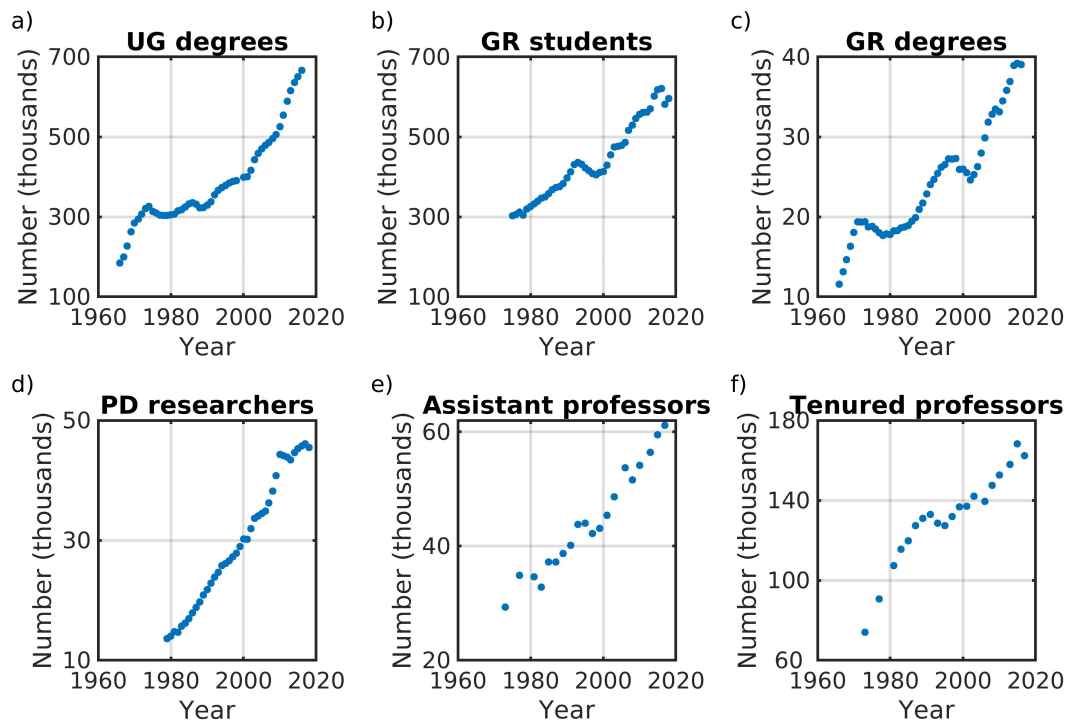


Figure S1: NSF timeseries raw data on (a) the number of bachelors degrees awarded, (b) the number of enrolled graduate students, (c) the number of PhDs awarded, (d) the number of postdoctoral researchers, (e) the number of assistant (tenure-track) professors, and (f) the number of tenured professors, across all of Science and Engineering in the US.

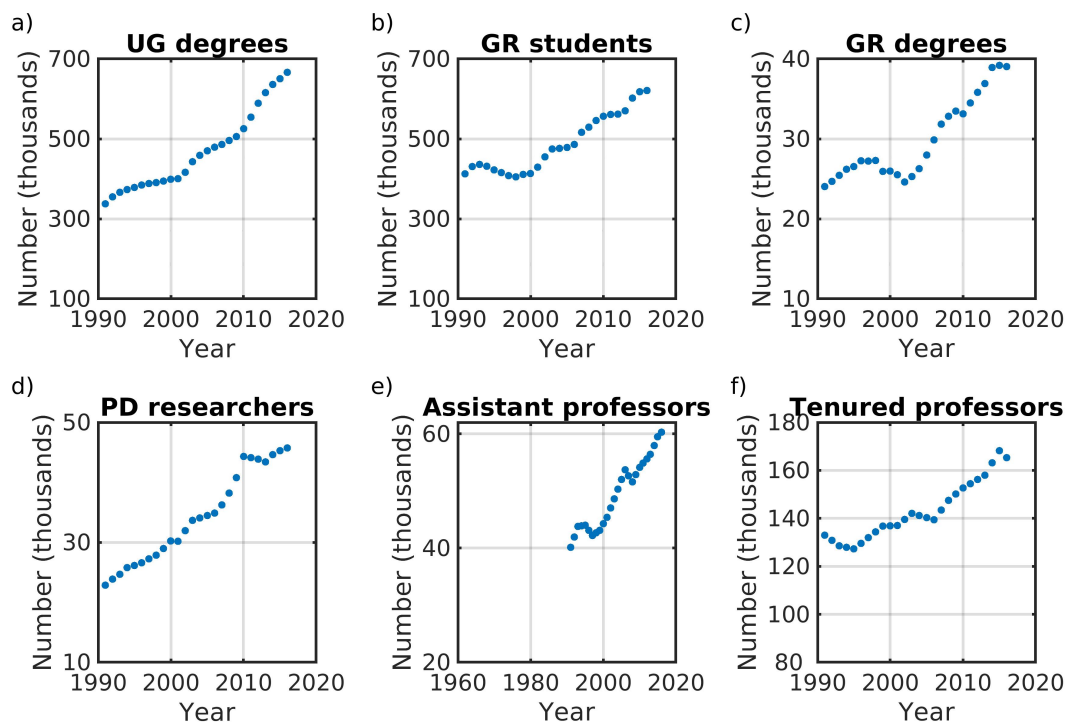


Figure S2: Interpolated and trimmed data on (a) the number of bachelors degrees awarded, (b) the number of enrolled graduate students, (c) the number of PhDs awarded, (d) the number of postdoctoral researchers, (e) the number of assistant (tenure-track) professors, and (f) the number of tenured professors, across all of Science and Engineering in the US.

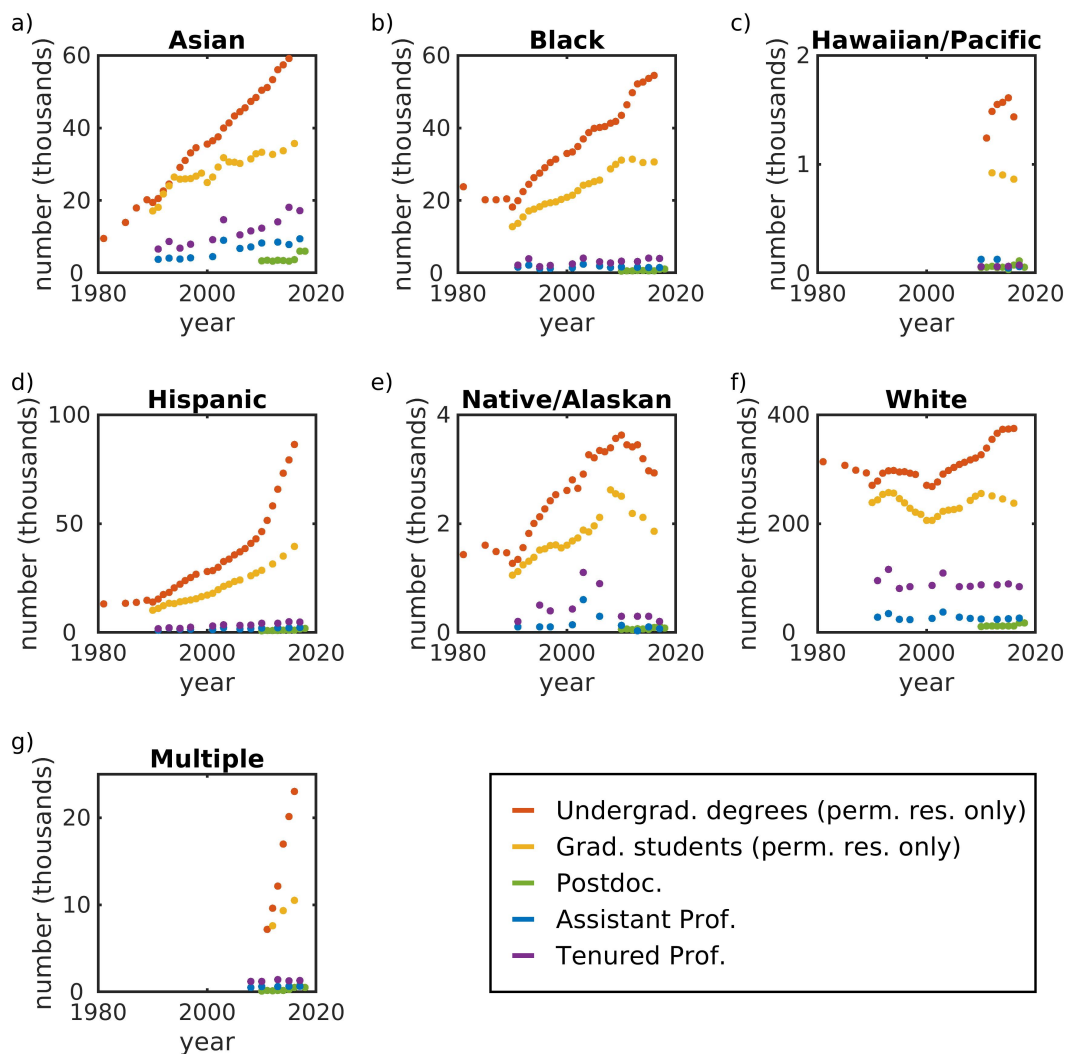


Figure S3: NSF timeseries data on the number of (a) Asian, (b) Black, (c) Hawaiian or Pacific Islander, (d) Hispanic, (e) Native American or Alaskan Native, and (f) White individuals in each stage: undergraduate degrees, graduate students, postdoctoral researchers, assistant professors and tenured professors. Note that race/ethnicity for undergraduate and graduate students is only recorded for US citizens and permanent residents, not temporary residents (but see Figure S5 below).

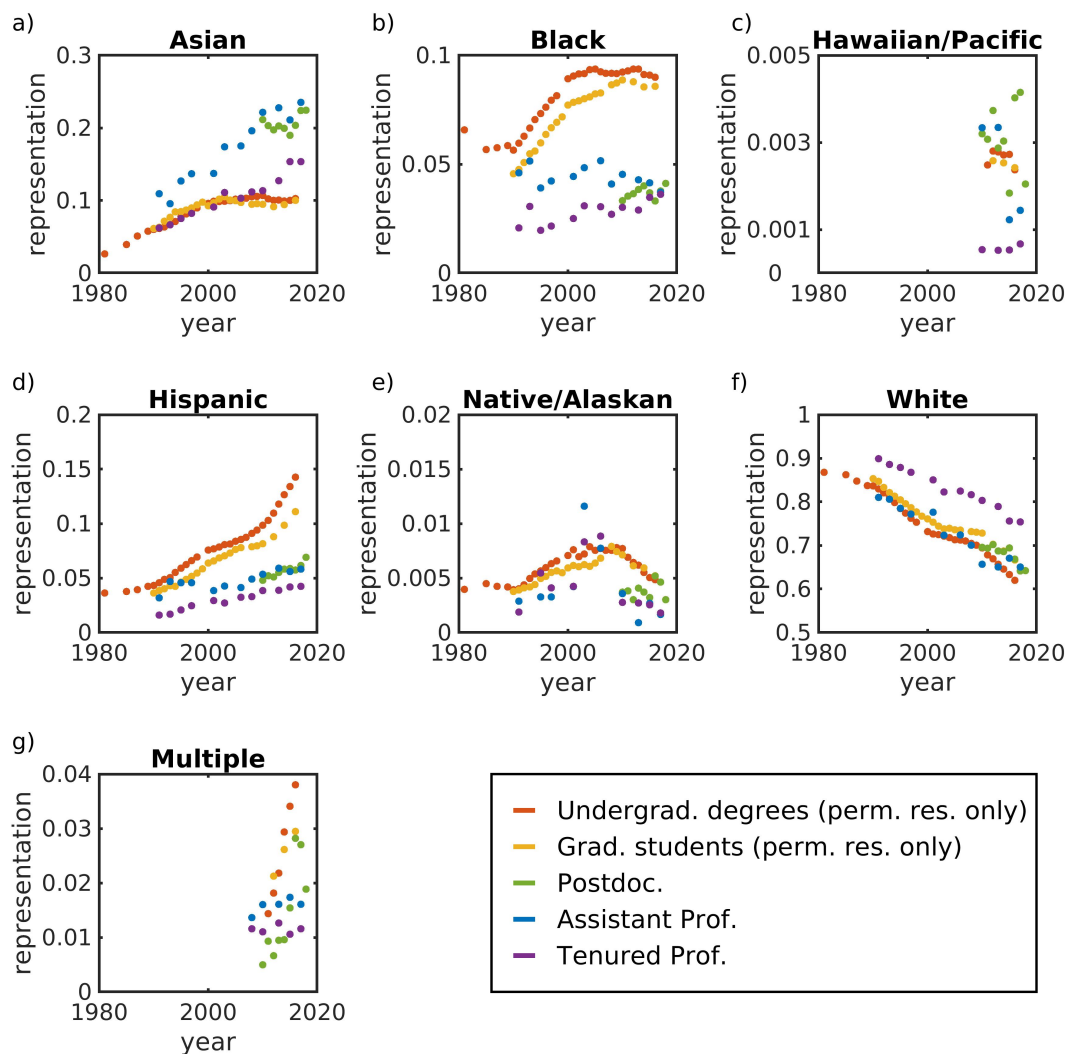


Figure S4: The proportion of individuals in each stage (undergraduate degrees, graduate students, postdoctoral researchers, assistant professors and tenured professors) that are (a) Asian, (b) Black, (c) Hawaiian or Pacific Islander, (d) Hispanic, (e) Native American or Alaskan Native, and (f) White. Note that race/ethnicity for undergraduate and graduate students is only recorded for US citizens and permanent residents, not temporary residents (but see Figure S5 below).

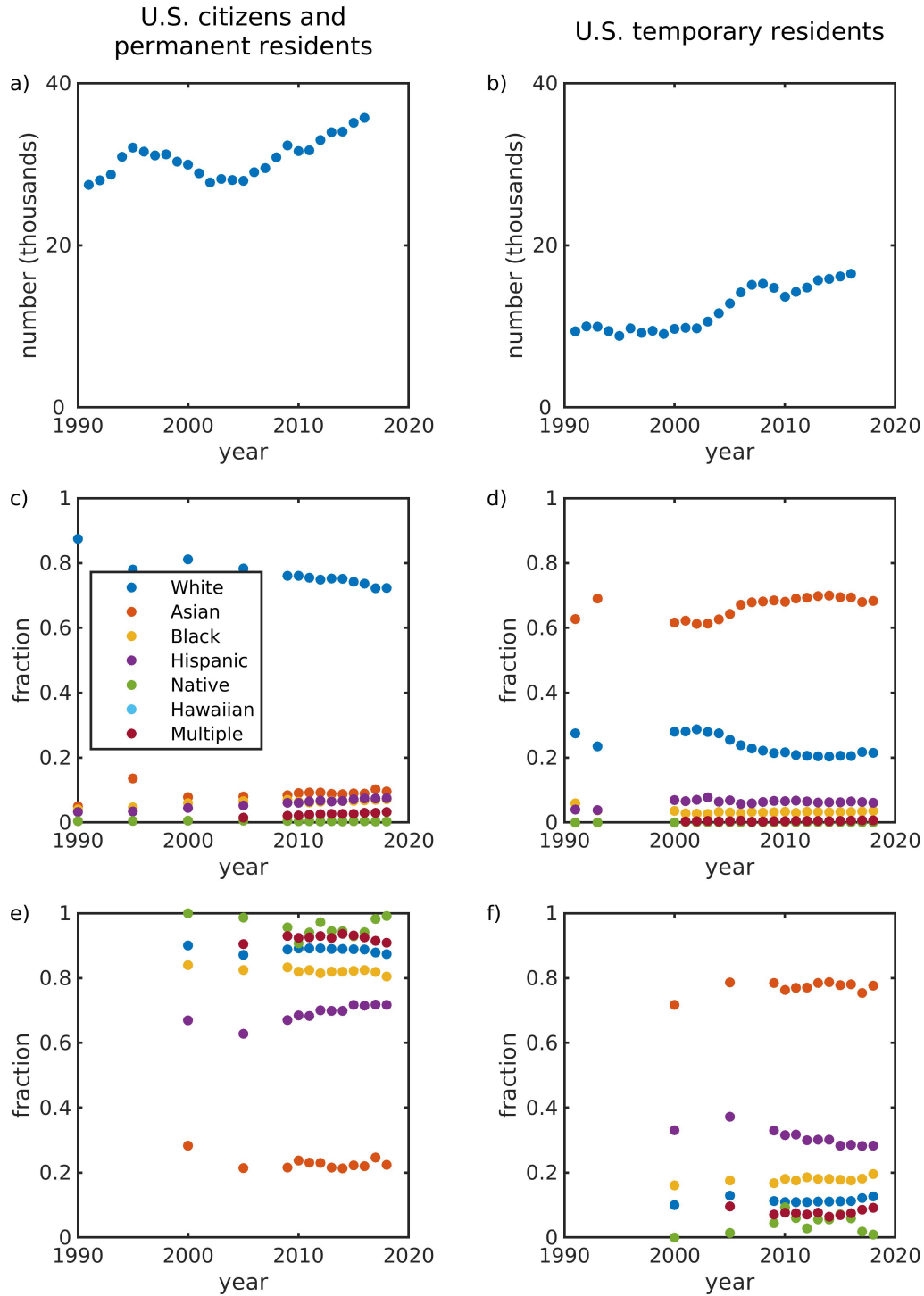


Figure S5: Composition of PhD recipients by residency. The number of PhD degree awardees who are (a) U.S. citizens or permanent residents and (b) temporary residents. The fraction of each race/ethnicity among (c) U.S. citizens or permanent resident PhD recipients and (d) temporary resident PhD recipients. The fraction of scholars of race/ethnicity that are (e) U.S. citizens or permanent residents and (f) temporary residents.

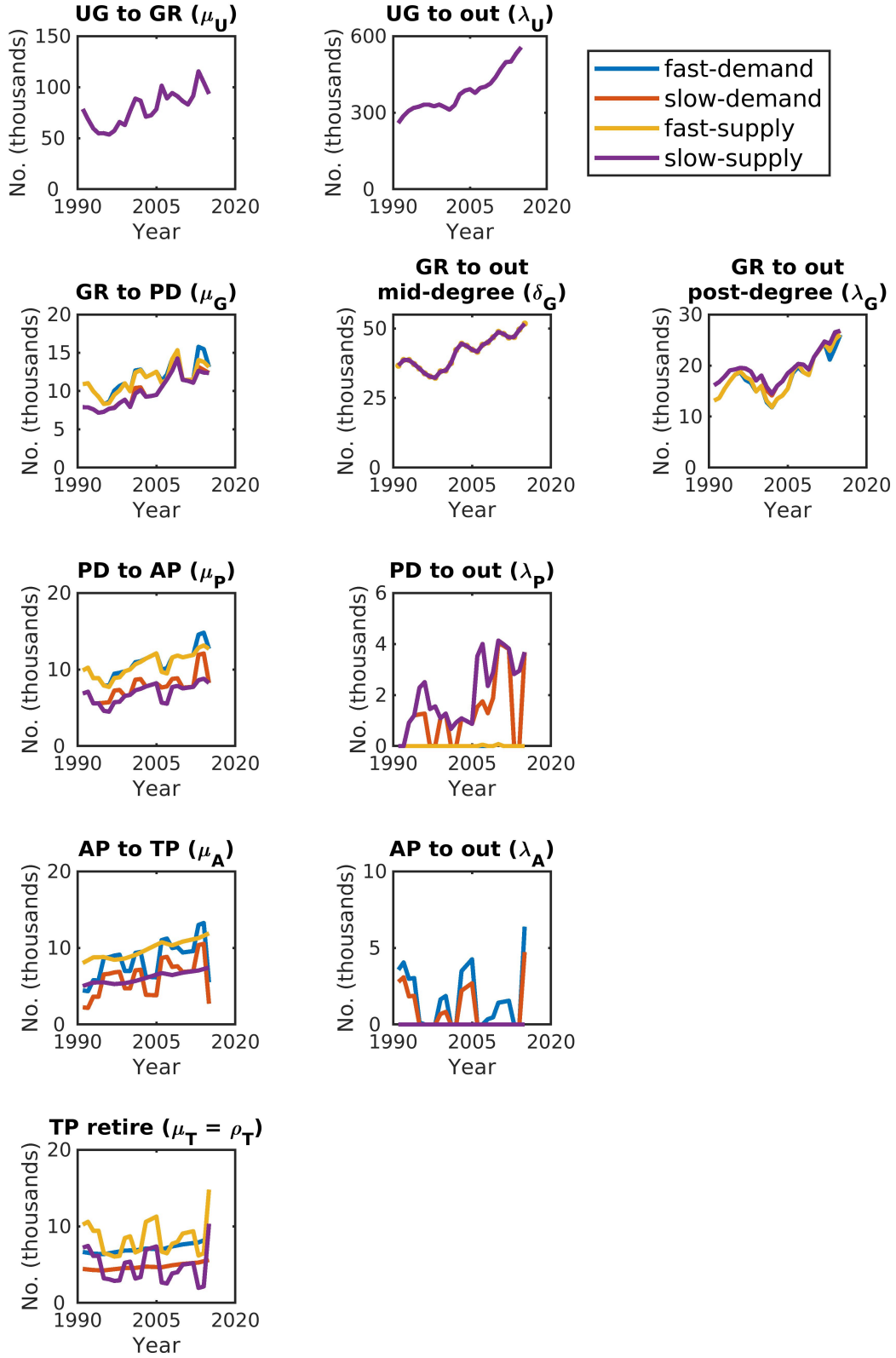


Figure S6: Timeseries estimates of the number of individuals making each of the 10 transitions in Figure 1 of the main text, as generated by our model for each of the four transitions for faculty scenarios: 'fast' and 'demand' (blue), 'slow' and 'demand' (red), 'fast' and 'supply' (yellow), 'slow' and 'supply' (purple) .

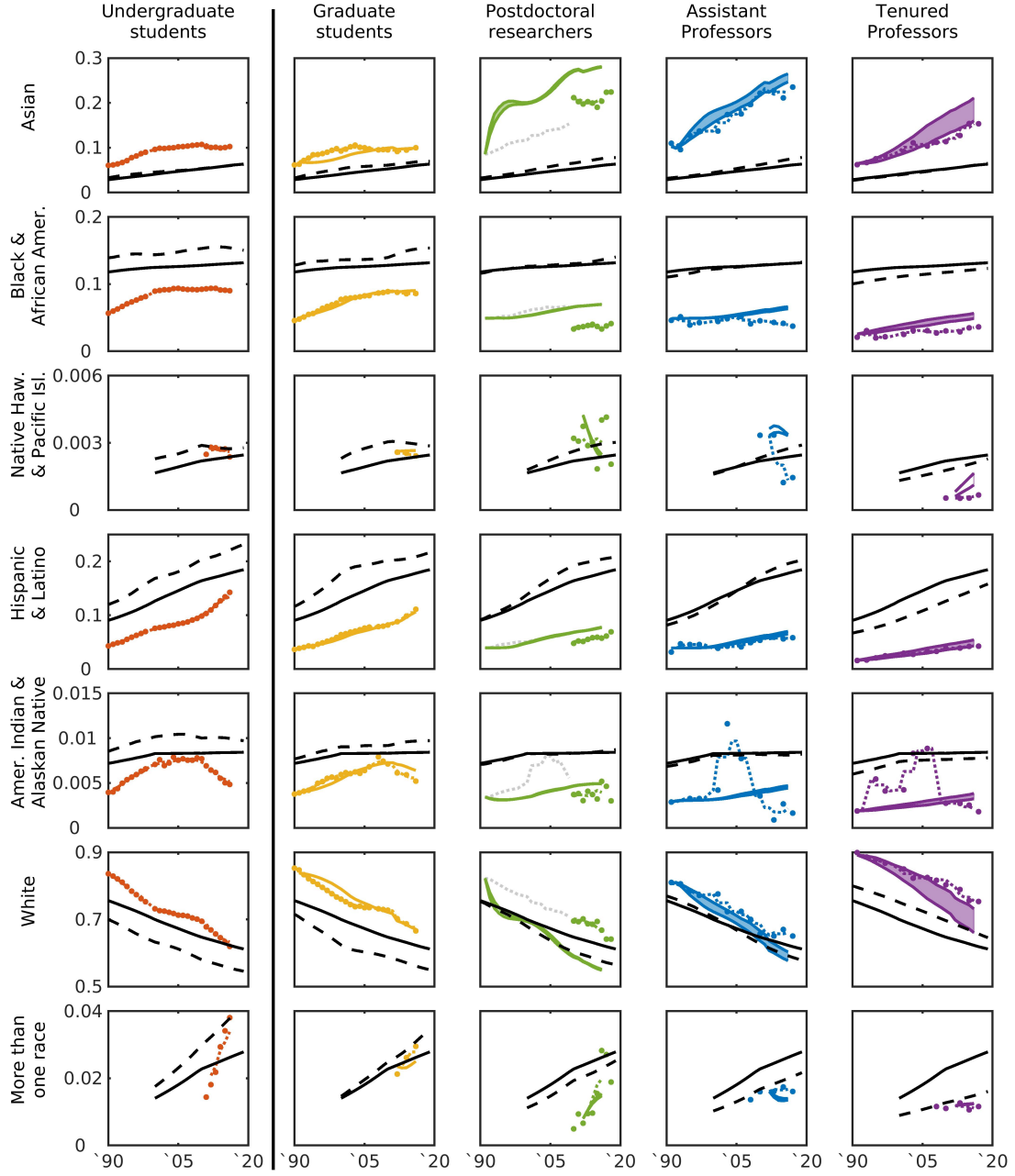


Figure S7: The representation of each race/ethnicity categories (rows) in each academic stage (columns) over time (i.e. the proportion of scholars in that stage that identify as that race or ethnicity) comparing: null model predictions (colored solid lines), academia data (dots are raw data, dotted lines are smoothed data), and census data for the U.S. overall population (black solid lines) and US age-specific population (black dashed line). Mismatch between model and academia data indicate race/ethnicity-based biases of retention within academia, mismatch between model and census indicates race/ethnicity-based biases in entering academia.

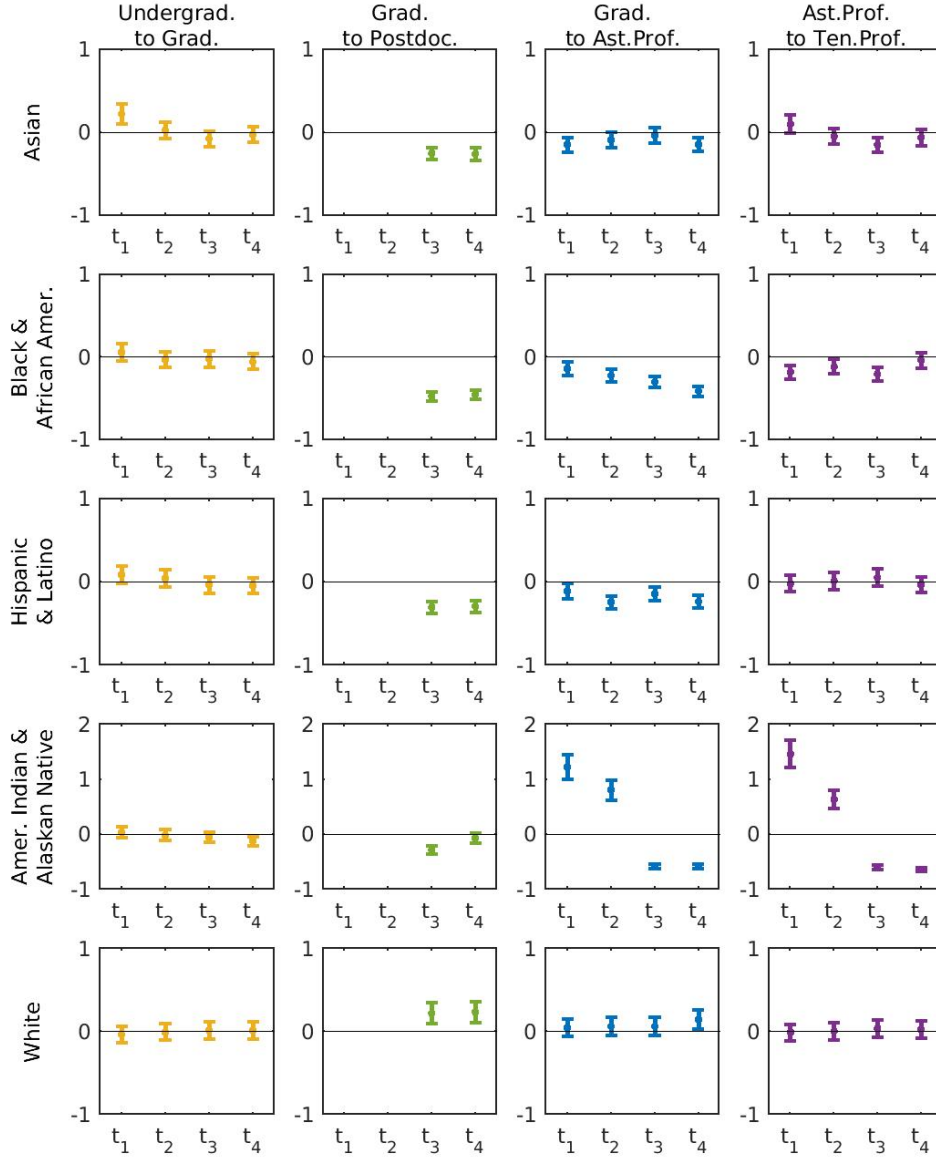


Figure S8: Temporal trends in the relative representation (θ ; comparing data and the null model) of each race/ethnicity category (rows) through one of the transitions within academia (columns). Each point corresponds to a single set of simulations, which were started in one of four years $t_1 = 1991$, $t_2 = 1996$, $t_3 = 2001$, $t_4 = 2006$) and run for 10 years. Colors correspond to the stage where difference is measured (same colors as Figs. 2 and 4 in the main text). Positive or negative values indicate a race/ethnicity category faces correspondingly positive or negative bias across that transition. Results for the Grad. to Postdoc. transitions are omitted for t_1 and t_2 because these results rely on extrapolated data, thus comparisons between model and data holds less value. Results for Hawaiian/Pacific Islander and More than one race are not shown because there was only sufficient data for a single time point (t_4).

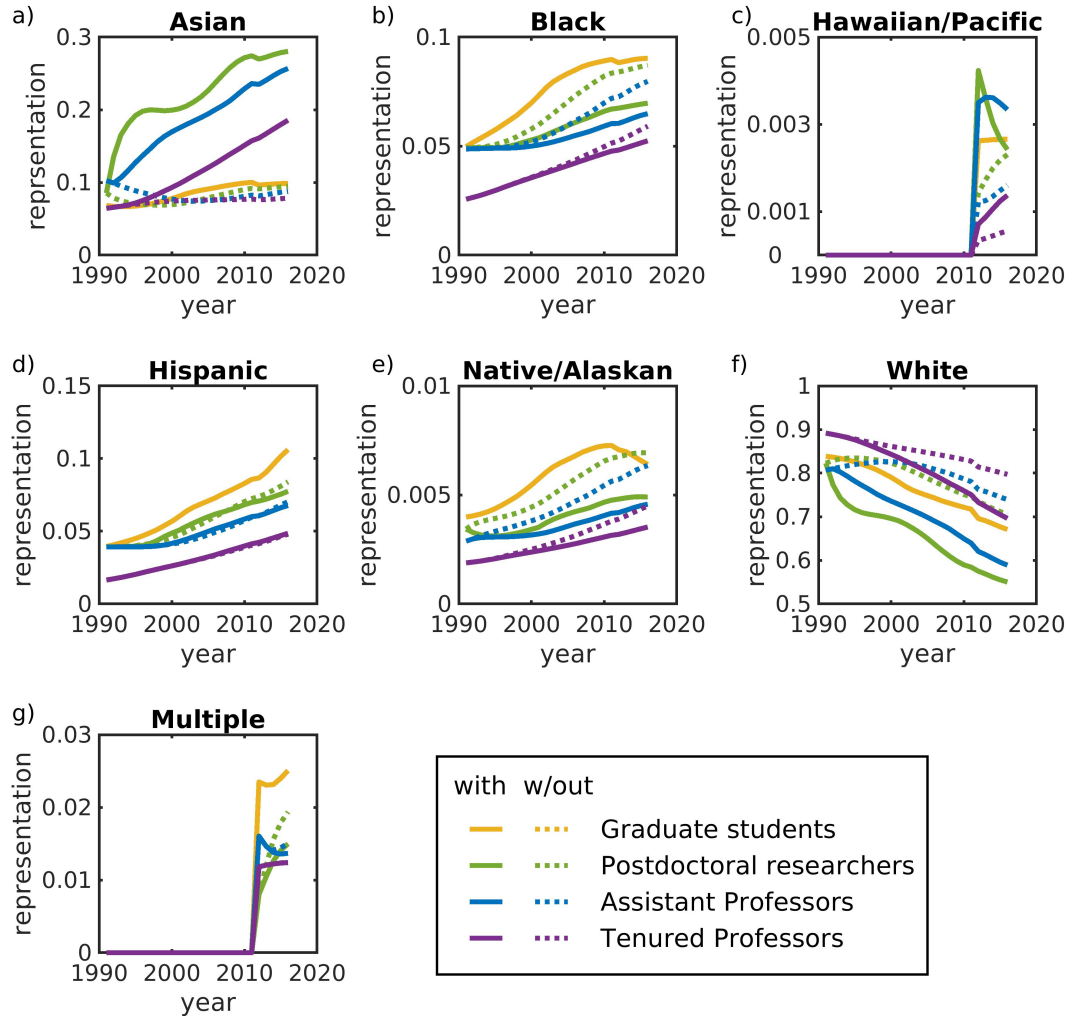


Figure S9: The representation of each race/ethnicity category (panels) in each academic stage (lines) over time, i.e. the proportion of scholars in that stage that identify as that race or ethnicity, comparing two versions of the model. The solid lines show the main model version which accounts for the race/ethnicity of temporary resident international scholars who receive their PhDs in the U.S. (i.e., the output at the graduate student stage matches the composition of PhD recipients, regardless of their residency). The dashed lines show a version of the model that ignores the race/ethnicity of international students (i.e., the output at the graduate student stage matches the composition of U.S. citizen and permanent resident PhD recipients).

Table S1: Data used for our model structure, years of data, and NSF report sources.

| DATA | YEARS | SOURCE |
|----------------------------|-----------|--|
| DEGREES ($D_i(t)$) | | |
| # Bachelors degrees | 1966-2012 | S&E Degrees, 2015 report, Table 5 |
| | 2006–2016 | WMPD, 2019 report, Table 7-4 |
| # PhD degrees | 1966-2012 | S&E Degrees, 2015 report, Table 19 |
| | 2006–2016 | WMPD, 2019 report, Table 5-3 |
| STAGE SIZE ($N_i(t)$) | | |
| # graduate students | 1975-2018 | GSPD, 2018 report, Tables 1-9a, 1-10a |
| # postdoctoral researchers | 1975-2018 | GSPD, 2018 report, Tables 1-9b, 1-10b |
| # assistant professors | 1973-2017 | SE-ind, 2019 report, Table S3-7 |
| # tenured professors | 1973-2017 | SE-ind, 2019 report, Table S3-7 |
| TIME IN STAGE (τ_i) | | |
| graduate student | 6.8 yrs | SE-ind, 2018 report, Table 2-30, 2015 data |
| postdoc | 2*1.9 yrs | [08-307] 2008 report, Table 2, 2006 data |
| assistant professor | 5-8 yrs | |
| tenured professor | 20-30 yrs | |

Table S2: Race/ethnicity data used for simulations and for comparisons against simulations, years of data, and NSF report sources.

| DATA | YEARS | SOURCE |
|---|-------------|-------------------------------|
| # Bachelors degrees (data by race/ethnicity for U.S. citizens and permanent residents only) | 1981-1991 | WMPD, 1994 report, Table 5-19 |
| | 1990-1998 | WMPD, 2002 report, Table 3-8 |
| | 1996-2007 | WMPD, 2009 report, Table C6 |
| | 2006-2016 | WMPD, 2009 report, Table 5-3 |
| # PhD students (data by race/ethnicity for U.S. citizens and permanent residents only) | 1990-1999 | WMPD, 2002 report, Table 4-6 |
| | 1999-2006 | WMPD, 2009 report, Table D-1 |
| | 2008-2010 | WMPD, 2011 report, Table 3-1 |
| | 2012 | WMPD, 2013 report, Table 3-1 |
| | 2014 | WMPD, 2017 report, Table 3-1 |
| | 2016 | WMPD, 2019 report, Table 3-1 |
| # PhD degrees (data by residency, permanent vs. temporary) used to calculate $R(t)$ | 1984:5:2014 | SED, 2014 report, Table 17 |
| | 1985:5:2015 | SED, 2015 report, Table 17 |
| | 1986:5:2016 | SED, 2016 report, Table 17 |
| | 1987:5:2017 | SED, 2017 report, Table 17 |
| | 1988:5:2018 | SED, 2018 report, Table 17 |
| # PhD degrees (data by race/ethnicity for temporary residents) used to calculate $V(t, k)$ | 2000–2010 | SED, 2010 report, Table 19 |
| | 2009–2018 | SED, 2018 report, Table 19 |
| doctoral workforce (data by race/ethnicity for temporary residents) used to calculate $V(t, k)$ | 1991 | WMPD, 1994 report, Table 8-11 |
| | 1993 | WMPD, 1996 report, Table 5-33 |
| # postdoctoral researchers (data by race/ethnicity) | 2010 | GSPD, 2010 report, Table 34 |
| | 2011-2016 | GSPD, 2010 report, Table 34 |
| | 2017 | GSPD, 2017 report, Table 2-2 |
| | 2018 | GSPD, 2018 report, Table 2-2 |
| # professors [assistant, tenured] (data by race/ethnicity) | 1991 | WMPD, 1994 report, Table 8-18 |
| | 1993 | WMPD, 1996 report, Table 5-28 |
| | 1995 | WMPD, 1998 report, Table 5-10 |
| | 1997 | WMPD, 2000 report, Table 5-19 |
| | 2001 | WMPD, 2004 report, Table H-26 |
| | 2003 | WMPD, 2007 report, Table H-28 |
| | 2006 | WMPD, 2009 report, Table H-28 |
| | 2008 | WMPD, 2011 report, Table 9-26 |
| | 2010 | WMPD, 2013 report, Table 9-26 |
| | 2013 | WMPD, 2015 report, Table 9-26 |
| | 2015 | WMPD, 2017 report, Table 9-26 |
| | 2017 | WMPD, 2019 report, Table 9-26 |

Table S3: Model variables, parameters, meaning and sources.

| Param. | Meaning | Source |
|-----------------|---|--------------|
| t | time (year) | NA |
| i | stage (U, G, P, A, T) | NA |
| k | individual race/ethnicity | NA |
| $N_i(t)$ | number of individuals in stage i in year t | see Table S1 |
| $D_i(t)$ | number of degrees of stage i awarded in year t (only $i = U, G$) | see Table S1 |
| $R(t)$ | fraction of PhD degrees to U.S. citizens / permanent residents in year t | see Table S2 |
| $V(t, k)$ | fraction of U.S. temporary resident PhD recipients in year t that are of race/ethnicity k | see Table S2 |
| τ_i | average number of years spent in stage i | see Table S1 |
| $\rho_i(t)$ | number of individuals potentially leaving stage i in year t | estimated |
| $\omega_i(t)$ | number of openings available in stage i in year t | estimated |
| $\mu_i(t)$ | individuals moving from stage i to stage $i + 1$ in year t | estimated |
| $\lambda_i(t)$ | individuals moving from stage i to outside the system in year t | estimated |
| $\delta_G(t)$ | individuals leaving stage G (before degree) in year t | estimated |
| $\beta_i(t, j)$ | individuals moving from subpartition j in stage i in year t | estimated |
| $n_i(t, k)$ | number of k individuals in stage i in year t | simulated |
| $f_i(t, k)$ | fraction of individuals in stage i in year t of race/ethnicity k | simulated |

Table S4: The fraction of individuals at each stage of each race/ethnicity in the year 2016. ‘Data’ rows are smoothed NSF counts data and the census data. The remaining rows are what the model predicts (under null model of no bias) for four scenarios: ‘fast-demand’, ‘fast-supply’, ‘slow-demand’, and ‘slow-supply’ which are combinations of a ‘demand’ or ‘supply’ view of faculty turnover and a ‘fast’ ($\tau_A = 5$, $\tau_T = 20$) or ‘slow’ turnover ($\tau_A = 8$, $\tau_T = 30$).

| | Asian | Black & Af. Am. | Nat. Haw. & Pac. Is. | Hisp. & Lat. | Amer. In. & Alas. Nat. | White | More than one race |
|----------------------------------|--------|-----------------|----------------------|--------------|------------------------|--------|--------------------|
| U.S. general population (census) | | | | | | | |
| data | 0.0602 | 0.1305 | 0.0018 | 0.1779 | 0.0084 | 0.6230 | 0.0209 |
| Undergraduate Students | | | | | | | |
| data | 0.1007 | 0.0906 | 0.0026 | 0.1345 | 0.0051 | 0.6327 | 0.03388 |
| Graduate Students | | | | | | | |
| data | 0.0973 | 0.0856 | 0.0025 | 0.1047 | 0.0056 | 0.6765 | 0.0278 |
| fast-demand | 0.0988 | 0.0903 | 0.0027 | 0.1062 | 0.0064 | 0.6707 | 0.0251 |
| fast-supply | 0.0988 | 0.0903 | 0.0027 | 0.1062 | 0.0064 | 0.6707 | 0.0251 |
| slow-demand | 0.0988 | 0.0903 | 0.0027 | 0.1062 | 0.0064 | 0.6707 | 0.0251 |
| slow supply | 0.0988 | 0.0903 | 0.0027 | 0.1062 | 0.0064 | 0.6707 | 0.0251 |
| Postdoctoral Researchers | | | | | | | |
| data | 0.2082 | 0.0378 | 0.0030 | 0.0610 | 0.0040 | 0.6662 | 0.0198 |
| fast-demand | 0.2812 | 0.0698 | 0.0023 | 0.0779 | 0.0049 | 0.5486 | 0.0153 |
| fast-supply | 0.2807 | 0.0698 | 0.0024 | 0.0776 | 0.0049 | 0.5495 | 0.0151 |
| slow-demand | 0.2799 | 0.0696 | 0.0025 | 0.0772 | 0.0049 | 0.5510 | 0.0148 |
| slow supply | 0.2798 | 0.0696 | 0.0025 | 0.0772 | 0.0049 | 0.5513 | 0.0148 |
| Assistant Professors | | | | | | | |
| data | 0.2230 | 0.0392 | 0.0013 | 0.0571 | 0.0022 | 0.6604 | 0.0167 |
| fast-demand | 0.2648 | 0.0663 | 0.0033 | 0.0700 | 0.0047 | 0.5774 | 0.0135 |
| fast-supply | 0.2635 | 0.0660 | 0.0034 | 0.0695 | 0.0047 | 0.5795 | 0.0135 |
| slow-demand | 0.2525 | 0.0641 | 0.0033 | 0.0665 | 0.0045 | 0.5953 | 0.0137 |
| slow supply | 0.2468 | 0.0632 | 0.0033 | 0.0649 | 0.0044 | 0.6034 | 0.0140 |
| Tenured Professors | | | | | | | |
| data | 0.1537 | 0.0354 | 0.0006 | 0.0422 | 0.0022 | 0.7548 | 0.0111 |
| fast-demand | 0.2053 | 0.0553 | 0.0015 | 0.0527 | 0.0037 | 0.6690 | 0.0124 |
| fast-supply | 0.2110 | 0.0564 | 0.0016 | 0.0542 | 0.0038 | 0.6604 | 0.0125 |
| slow-demand | 0.1663 | 0.0493 | 0.0012 | 0.0441 | 0.0033 | 0.7235 | 0.0124 |
| slow supply | 0.1599 | 0.0487 | 0.0011 | 0.0432 | 0.0032 | 0.7315 | 0.0123 |