

## The emergence of heterogeneous scaling in research institutions

Keith A. Burghardt <sup>1</sup>✉, Zihao He <sup>1</sup>, Allon G. Percus<sup>1,2</sup> & Kristina Lerman<sup>1</sup>

Research institutions provide the infrastructure for scientific discovery, yet their role in the production of knowledge is not well characterized. To address this gap, we analyze interactions of researchers within and between institutions from millions of scientific papers. Our analysis reveals that collaborations densify as each institution grows, but at different rates (heterogeneous densification). We also find that the number of institutions scales with the number of researchers as a power law (Heaps' law) and institution sizes approximate Zipf's law. These patterns can be reproduced by a simple model in which researchers are preferentially hired by large institutions, while new institutions complementarily generate more new institutions. Finally, new researchers form triadic closures with collaborators. This model reveals an economy of scale in research: larger institutions grow faster and amplify collaborations. Our work deepens the understanding of emergent behavior in research institutions and their role in facilitating collaborations.

<sup>1</sup>Information Sciences Institute, University of Southern California, Marina del Rey, USA. <sup>2</sup>Institute of Mathematical Sciences, Claremont Graduate University, Claremont, USA. ✉email: [keithab@isi.edu](mailto:keithab@isi.edu)

Scientific innovation and training require efficient and robust infrastructure. This infrastructure is provided by research institutions, a category that includes universities, government labs, industrial labs, and national academies<sup>1–5</sup>. Despite the long tradition of bibliometric and science of science research<sup>6</sup>, the focus has only recently shifted from individual scientists<sup>7,8</sup> and teams<sup>9–11</sup> to how institutions affect researcher productivity and impact<sup>12,13</sup>. Many gaps remain in our understanding of the role of institutions in the production of scientific knowledge, and specifically, how they form, grow, and facilitate scientific collaborations. These questions are important, because collaborations are increasingly prevalent in scientific research<sup>1,9,10</sup> and produce more impactful and transformative work<sup>10,14</sup>. Collaboration allows scientists to cope with the increasing complexity of knowledge<sup>15</sup> by leveraging the diversity of expertise<sup>16</sup> and perspectives offered by collaborators from different institutions<sup>17</sup> and disciplines<sup>18</sup>.

To understand the evolution of research institutions and collaborations, we analyze a large bibliographic database spanning many decades and multiple scientific disciplines. The database contains millions of publications from which the names of authors (collaborators) and their affiliations (research institutions) have been extracted for each paper. Analysis of these data reveals strong statistical regularities. We find that collaborations scale superlinearly with institution size, i.e., faster than institutions grow, consistent with densification of growing networks<sup>19–21</sup>. However, the scaling law is different for each institution, and as a result, different parts of the collaboration network densify at different rates. We also find that institutions vary in size by many orders of magnitude with an approximately power-law distribution, also known as Zipf's law<sup>22</sup>. The number of institutions, in contrast, scales sublinearly with the number of researchers, thus following Heaps' law<sup>23,24</sup>. The sublinear scaling implies that, even as more institutions appear, each institution gets larger on average, but this average belies an enormous variance.

Finally, we create a stochastic model that helps explain how institutions and research collaborations form and grow. In this model, a researcher appears at each time step and is preferentially hired by larger institutions (e.g., due to their prestige or funding), which leads to the rich-get-richer effect creating Zipf's law. With a small probability, however, a researcher joins a newly appearing institution. The arrival of this new institution then triggers yet more new institutions to form in the future, which explains Heaps' law<sup>25</sup>. Finally, once hired, researchers make connections to other researchers and their collaborators with an independent probability to explain collaborations scaling superlinearly with institution size. Despite its simplicity, the model reproduces a range of empirical observations, including the number and size of research institutions, and how pockets of increasingly dense structures form in collaboration networks.

These empirical results demonstrate universal emergent patterns in the formation and growth of research institutions and collaborations. Our model demonstrates that new institutions are critical to absorbing extra capacity by collecting researchers who do not join large institutions. At the same time, large institutions offer an economy of scale: they grow faster and provide more collaboration opportunities compared to smaller institutions.

## Results and discussion

As the first step towards characterizing the complexity of institution scaling, we collect data from Microsoft Academic Graph<sup>26</sup> to capture how millions of collaborations evolve over time. Figure 1 shows the collaboration network at the institution level in the field of sociology. Figure 1a demonstrates a remarkable diversity of

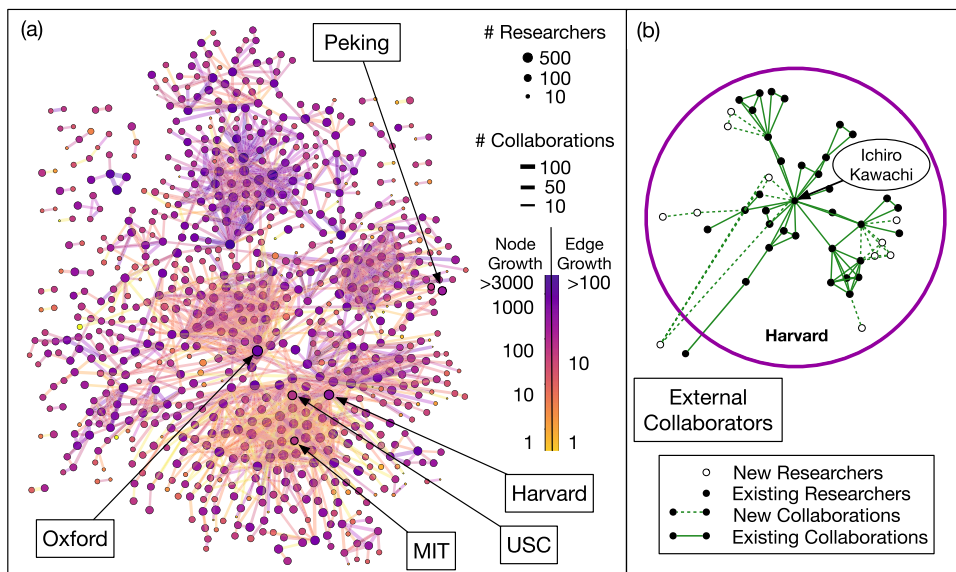
institution size and growth, both in terms of the number of researchers (node growth) and collaborations between institutions (edge growth). Collaborations are clustered, with clear groups of interacting institutions. Research collaborations within an institution are equally complex. Figure 1b highlights the largest connected component of the collaboration network within Harvard. Individual researchers vary widely in the number of collaborators, with new collaborations appearing in clusters.

This dataset helps us capture how the number of collaborations scale with an institution's size,  $n$ . Figure 2a, b shows the number of internal and external collaborations versus  $n$  across four different disciplines: computer science, physics, math, and sociology. While each institution follows a scaling law  $c \sim n^\alpha$  ( $R^2$  is close to 1.0, see Supplementary Note 6), the exponents  $\alpha$  differ substantially between institutions. This is shown in the insets of Fig. 2a, b where we collect scaling exponents across thousands of institutions and notice that their distribution stretches between zero (in which institutions do not gain any collaborations) to two (collaborations are extremely dense). In the thermodynamic limit, exponents cannot be larger than two, therefore values above two are due to finite-size effects.

To show that the scaling exponents of all institutions are different, we create a null model (see Supplementary Note 3) in which all institutions follow the same scaling law. In this null model, residuals of each institution's fitted scaling relation are reshuffled and added as noise onto a single scaling relation. Differences between fitted exponents in this model are due to statistical noise rather than different scaling laws. We find that the variance of the scaling laws across all institutions is much higher than this null model. We therefore reject the hypothesis that all the exponents within a field are the same within statistical error. We explore the dependence of scaling on final institution size in Supplementary Note 6, and find the scaling exponents are superlinear (approximately 1.2 on average) and do not depend strongly on the final size of the institution. Different parts of the collaboration network therefore densify at different rates, which extends on previous work that uncovered densification for many networks at the aggregate level<sup>19</sup>.

We find weak evidence that higher scaling exponents correspond to institutions with greater impact. In physics, the Spearman rank correlation,  $s$ , between mean paper impact after five years and internal collaboration scaling exponents is 0.09 (borderline significant,  $p$ -value = 0.06) and for external collaboration is 0.27 ( $p$ -value <  $10^{-5}$ ). Similarly, in sociology, the correlation is 0.19 ( $p$ -value = 0.03) between impact and internal collaboration exponents, and the same correlation value is found for external collaboration exponents. For all other fields, however, the correlations are not statistically significant ( $p$ -value  $\geq$  0.20). Impact, a proxy of institution research quality, cannot fully explain why collaborations grow faster in some institutions and not others, but can give some insight into reasons for this diversity. These results suggest that highly impactful institutions seem to form collaborations more easily as they grow. Nonetheless, almost all institutions benefit from being larger, as the number of collaborations per person typically grows with size (Fig. 2a, b inset).

The superlinear scaling of collaborations cannot be explained by researcher productivity. The scaling exponents of output, i.e., the cumulative number of papers published by researchers affiliated with that institution at a given year, are centered around 1.0 (see Supplementary Note 4). Paper output per researcher is therefore approximately independent of institution size. The average team size per institution, however, increases with institution size (see Supplementary Note 5), which may help explain the scaling of collaborations. Namely, as institutions grow, they form larger teams for each paper. This, in turn, creates more collaborations (which are proportional to the team size squared).



**Fig. 1 Network visualization of the collaborations in the field of sociology in 2017.** **a** Collaborations between institutions. Each node represents a research institution, and institutions with more researchers are represented by larger nodes. Each link represents collaborations between researchers at different institutions and more collaborations are represented by thicker lines. Darker nodes represent faster-growing institutions (defined as the number of new researchers added between 2012 and 2017), and darker links represent faster-growing collaborations (defined as the number of new inter-institution collaborations between 2012 and 2017). Links with fewer than 10 collaborations are removed, as are isolated nodes. A few major universities are labeled: Peking University, Oxford University, Harvard University, Massachusetts Institute of Technology (MIT), and the University of Southern California (USC). **b** The largest connected component of collaborations within Harvard University. Each node represents a researcher. Dashed lines represent new collaborations added between 2012 and 2017, while open circles represent new researchers added between 2012 and 2017. The highest degree node is Ichiro Kawachi, a highly cited sociologist.

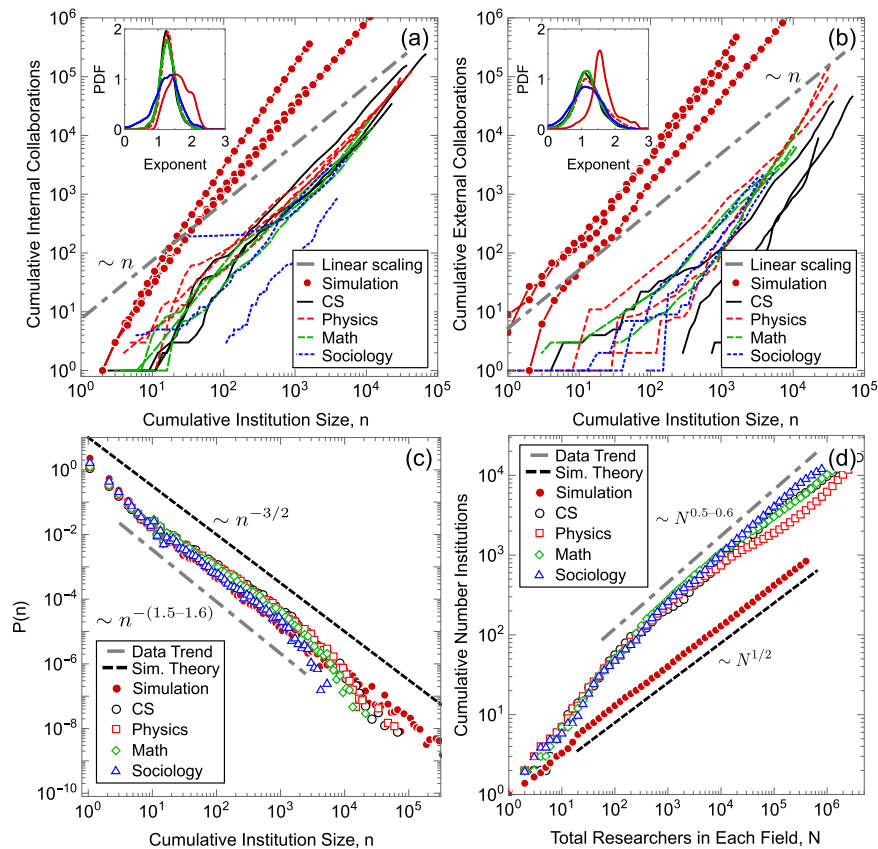
We also find that the distribution of institution sizes (as of 2017) follows Zipf’s law (Fig. 2c), similar to the observed heavy-tailed distribution of city sizes<sup>22,27</sup>. In Supplementary Note 1 and Supplementary Data 1, we show that while the largest institutions are intuitive, such as Harvard, the smaller institutions tend to be for-profit colleges, community colleges, and institutions without a formal department in the field of interest (e.g., an engineering school with papers in sociology). In addition, the number of institutions grows sublinearly with the number of researchers in each field (Fig. 2d). This feature, known as Heaps’ law, implies that quadrupling the number of researchers in a field roughly doubles the total number of institutions associated with that field. Exact scaling law values for each field can be found in Table 1, where Heaps’ laws are calculated for the total number of researchers in each field,  $N$ , greater than twenty and Zipf’s law is calculated for institution size,  $n$ , greater than ten.

**A Model of Institution Growth.** We now describe a stochastic growth model of institution formation that elucidates how institutions and collaborations jointly grow. We model institution formation and growth with a Pólya’s urn-like set of mechanisms described in ref. <sup>25</sup>, and we model the growth of collaborations with a network densification mechanism<sup>20,21</sup>. Unlike existing models of network densification<sup>19–21</sup>, however, our model reproduces the heterogeneous densification of internal and external collaborations, and the non-trivial growth structure on institutions. This is complimentary to a very recent model on heterogenous exploration<sup>28</sup>, in which Pólya’s urn models vary as a function of a node’s position on a (static) network.

We imagine an urn containing balls of different colors. The balls can be thought of as the resources given to each institution, where each color represents a different assigned institution, as shown in Fig. 3a. Balls are picked uniformly at random with replacement, with each pick representing a newly-hired

researcher, and the ball color is recorded in a sequence to represent what institution hires the researcher. Afterwards,  $\rho$  balls of the same color are added to the urn to represent the additional resources and prestige given to a larger institution, known as reinforcement (left panel of Fig. 3a)<sup>25</sup>. If a previously unseen color is chosen, then  $\nu + 1$  uniquely colored balls are placed into the urn, a step known as triggering (right panel of Fig. 3a)<sup>25</sup>. The new colors represent institutions that are able to form because of the existence of a new institution. This triggering, also known as adjacent possible<sup>25</sup>, does not imply causality per se, e.g., the cause of the University of California Merced’s creation was not strictly because of previously established institutions. Instead, these institution-specific causes are represented as stochastic noise, a remarkable simplification that does not remove the observed statistical regularities. Triggering, however, agrees with anecdotal evidence, making it an intuitive factor behind the creation of institutions. For example, UC Davis was spun out of UC Berkeley, and USC Institute for Creative Technology was spun out of USC Information Sciences Institute, which itself was founded by researchers from the Rand Corporation. The model we describe is known as Pólya’s urn with triggering<sup>25</sup>, and predicts Heaps’ law with a scaling relation  $\sim N^{\nu/\rho}$  and Zipf’s law with scaling relation  $\sim n^{-(1+\nu/\rho)}$ . In our simulations, we arbitrarily chose  $\rho$  to be 4 and  $\nu$  to be 2, which agrees well with the data shown in Fig. 2.

Next, we explain the heterogeneous and superlinear scaling of collaborations through a model of network densification. Building on the work of<sup>20,21</sup>, we have each new researcher, represented as a node, connect to a random researcher within the same institution, as well as an external researcher picked uniformly at random (left panel of Fig. 3b). New collaborators are then chosen independently from neighbors of neighbors with probability  $p_i$ , where  $p_i$  is unique to each researcher’s institution (right panel of Fig. 3b). We let  $p_i$  be a Gaussian distributed random variable with mean,



**Fig. 2 Institutions densify at different rates but their size and frequency follow universal patterns.** **a** Internal and **(b)** external collaborations versus institution size for three arbitrarily chosen institutions with more than  $10^3$  cumulative researchers in each field or simulation. Circle markers correspond to simulation data; solid lines, medium dashed lines, long dashed lines, and short dashed lines correspond to data from the fields of computer science (CS), physics, math, sociology, respectively. Dash-dotted lines report linear scaling, showing that institutions' scaling laws are super-linear. Insets: distribution of exponents across thousands of institutions (cf. Supplementary Note 1). **c** The distribution of researchers in each institution as of 2017 (Zipf's law), and **(d)** the number of unique institutions versus the total number of researchers in each field (Heaps' law). Closed circle markers correspond to simulation data; open circles, squares, diamonds, and triangles correspond to computer science (CS), physics, math, sociology, respectively. In addition, light dash-dotted lines indicate empirical trends while darker dashed lines indicate theoretical scaling law exponents  $-1 - \nu/\rho$  and  $\nu/\rho$  for **(c)** and **(d)** respectively<sup>25</sup>. Simulation data in **(a)** and **(b)** are collected from four realizations and in **(c)** and **(d)**, from fifteen realizations (individual realizations show similar trends). Simulation parameters are  $\rho$  equals 4,  $\nu$  equals 2,  $\mu_p$  equals 0.6, and  $\sigma_p$  equals 0.25.

**Table 1 Zipf's law and Heaps' Law exponents for research fields and simulation.**

Discipline	Heaps' Law Exponent	Zipf's Law Exponent
Comp. Sci.	$0.554 \pm 0.004$	$-1.470 \pm 0.005$
Physics	$0.501 \pm 0.007$	$-1.474 \pm 0.006$
Math	$0.549 \pm 0.008$	$-1.516 \pm 0.006$
Sociology	$0.622 \pm 0.005$	$-1.603 \pm 0.009$
Simulation	1/2	-3/2

Each fit is a linear regression on log-scaled x and y axes for the number of researchers in each field above 100. Errors are standard errors of linear regression coefficients. Simulation scaling laws are theoretical exponents calculated for Polya's urn model with triggering with coefficients  $\rho = 4$  and  $\nu = 2$ <sup>25</sup>. See Results and Discussion for details of the mechanism coefficients.

$\mu = 0.6$ , and standard deviation,  $\sigma_\mu = 0.25$  and truncated between 0 and 1. Lambiotte et al.<sup>21</sup> show that their equivalent to  $\mu$ , when greater than 0.5, produces densification. We therefore choose  $\mu = 0.6$  to ensure the network densifies. We show separately that  $p_i$  directly controls the heterogeneity we observe in internal collaboration scaling, but the heterogeneity in external collaboration scaling is an emergent outcome of this model<sup>29</sup>.

To summarize, our model has four parameters:  $\rho$  (reinforcement),  $\nu$  (triggering), and two parameters to explain collaboration densification heterogeneity,  $\mu_p$  and  $\sigma_p$ . In the main text, we let  $\rho$  equal 4,  $\nu$  equal 2,  $\mu_p$  equal 0.6, and  $\sigma_p$  equal 0.25. These are arbitrarily chosen parameters meant to create statistical patterns that are qualitatively similar to empirical data. Namely,  $\mu_p > 0.5$  ensures collaboration densification<sup>21</sup>, and  $\sigma_p > 0$  ensures that densification scaling exponents vary between institutions. Interestingly, this model's Zipf's and Heaps' laws can be exactly calculated, as discussed by Tria et al.<sup>25</sup>, with Zipf's law exponent equal to  $-1 - \nu/\rho$  and Heaps' law equal to  $\nu/\rho$ . This model qualitatively reproduces Zipf's and Heaps' laws (Fig. 2c, d and Table 1) and the heterogeneous scaling of internal and external collaborations shown in Fig. 2a, b. While other plausible mechanisms for Zipf's law<sup>30–32</sup>, Heaps' law<sup>24</sup>, or densification<sup>19</sup> exist, the current model describes these patterns in a cohesive framework and explains the heterogeneous scaling we discover in the data. While this heterogeneity is built into our internal scaling laws, the external scaling heterogeneity is an emergent property within the model<sup>29</sup>.

The model also reproduces qualitative trends of cross-sectional analysis. Specifically, the scaling exponents of internal collaborations produced by the model when measured at a specific point in



**Fig. 3 Schematic representation of the institution growth model.** **a** At time  $t$  a new researcher is hired, modeled as extracting a ball with uniform probability with replacement from an urn,  $U$  (black arrow). The ball color represents an institution. Hiring a researcher will always add  $\rho$  new balls of the same color to the urn in the next timestep (reinforcement). Hiring the first researcher at an institution (picking a ball color that has never been picked before), triggers  $\nu + 1$  new colors to enter the urn, increasing the likelihood of more institutions to hire their first researcher (triggering). **b** Researchers within each institution (dash-dotted boxes) have both internal collaborators (darker solid lines) and external collaborators (gray lines). Once a researcher is hired, they choose one random internal and one random external collaborator (solid arrows). New collaborations (dashed arrows) are formed independently with probability  $p_A$ , if hired by institution A, and  $p_B$  if hired by institution B. These new connections form triangles.

time, i.e., in a cross-sectional setting, vary in time and are larger than scaling exponents of external collaborations and decrease over time (Supplementary Note 6), unlike what we see in data (Supplementary Fig. 3). These results are robust to stochastic variations of the densification mechanism (Supplementary Note 7). As a final comparison with data, we compared the growth of institutions and the ways links form to the model mechanisms and found broad agreement<sup>29</sup>.

## Conclusion

We identify strong statistical regularities in the growth of research institutions. The number of collaborations increases superlinearly with institution size, i.e., faster than institutions grow in size, though the scaling is heterogeneous, with a different exponent for each institution. Therefore, each institution has its own universal scaling, i.e., regardless of its size, it will always have the same percentage of new collaborations for each percentage increase in size. The super scaling is not explained by the increased productivity of researchers at larger institutions the number of papers per researcher is roughly independent of institution size. Instead, the growing collaborations are associated with bigger teams at larger institutions. The diversity in collaboration scaling exponents is partly explained by variations in institution impact. Institutions with higher impact papers also tend to have a larger scaling exponent. This provides evidence that a higher collaboration scaling exponent allows for collaborations to form more easily, and that in turn creates higher-impact papers. Further analysis is needed to test this hypothesis in the future.

When these observations are incorporated into a minimal stochastic model of institution growth, we are able to reproduce the

surprising regularity of research institution formation, growth and the heterogeneous densification of collaboration networks. That said, there is still room for improvements to this model, given quantitative differences between the model and data, such as the constant shift difference between the Heaps' laws (Fig. 2c), or the difference in the collaboration scaling law exponents (insets of Fig. 2a, b).

These findings support the idea that academic environments differ in their ability to bolster researcher productivity and prominence<sup>12</sup>, and also demonstrate that institution size and ability to facilitate collaborations as a potential factor explaining differences in academic environments. Additional research is needed to identify other factors that contribute to an institution's success.

## Methods

**Data.** We use bibliographic data from Microsoft Academic Graph (MAG), from which researcher names (authors), their institutional affiliation, and references made to other papers have been extracted<sup>26,33</sup>. MAG data has disambiguated institutions and authors for each paper, allowing us to consider all authors with the same unique identifier to be the same researcher, and similarly for each institution. In these data, authors typically have only one affiliation at any time (see Supplementary Note 1). We focus on four fields of study: computer science, physics, math and sociology. After data cleaning, we have almost ten million papers published between 1800 and 2018 (see Supplementary Note 1). Our computer science data includes early research in topics relating to computers, including electrical engineering, and therefore stretches back to before 1900.

We define *institution size* in a given year as the number of authors who have been ever been affiliated with that institution up until that year. *Collaborations* are defined as two researchers who have co-authored a paper up until that year. We distinguish between internal collaborations (co-authors at the same institution) and external collaborations (co-authors affiliated with different institutions). Finally, to understand the relation between collaborations and institution size, we define output as the cumulative number of papers from researchers affiliated with an institution in a particular year.

**Analysis.** We use cumulative statistics to reduce statistical variations and to better compare to a stochastic growth model of institution formation. To check the robustness of results, we compare to an alternate yearly definition of institution size and collaborations (see Supplementary Note 2). We find all qualitative results are the same, in part because both definitions are highly correlated.

We present scaling results for longitudinal analysis, which tracks how collaborations evolve as individual institutions grow<sup>34–36</sup>. This contrasts to cross-sectional analysis applied in previous work on city scaling<sup>37,38</sup> and institution scaling<sup>2–4,39</sup>, which measures collaborations as a function of the size of all institutions at a given point in time. We find that cross-sectional analysis identifies scaling laws that are not representative of the growth of most institutions (see Supplementary Note 7), and while simulations and empirical data give scaling exponents that are fairly constant in time for each institution, cross-sectional scaling exponents vary in time for both data and simulation. For these reasons, we focus on longitudinal scaling analysis in this paper, although scaling laws derived by either analysis method strongly relate to each other<sup>36,40</sup>.

## Data availability

Microsoft Academic Graph data can be accessed via the following link: <https://www.microsoft.com/en-us/research/project/microsoft-academic-graph/><sup>26</sup>. Replication data collected from Microsoft Academic Graph are available in the following repository: <https://gigantum.com/keithburghardt/heterogeneous-scaling>. Sample raw data for small institutions are available in Supplementary Data 1.

## Code availability

Code for this study is available in the following repository: <https://gigantum.com/keithburghardt/heterogeneous-scaling>.

Received: 2 February 2021; Accepted: 29 July 2021;

Published online: 02 September 2021

## References

- Hicks, D. & Katz, J. S. Science policy for a highly collaborative science system. *Science and public policy* **23**, 39–44 (1996).
- Taylor, R. C. et al. The scalability, efficiency and complexity of universities and colleges: a new lens for assessing the higher educational system. Preprint at <https://arxiv.org/abs/1910.05470> (2019).

3. van Raan, A. F. J. Universities scale like cities. *PLoS ONE* **8**, e59384 (2013).
4. Jamtveit, B., Jettestuen, E. & Mathiesen, J. Scaling properties of European research units. *Proc. Natl. Acad. Sci.* **106**, 13160–13163 (2009).
5. Murray, D. et al. Unsupervised embedding of trajectories captures the latent structure of mobility. Preprint at <https://arxiv.org/abs/2012.02785> (2020).
6. Fortunato, S. et al. Science of science. *Science* **359**, ea0185 (2018).
7. Sinatra, R., Wang, D., Deville, P., Song, C. & Barabási, A.-L. Quantifying the evolution of individual scientific impact. *Science* **354**, aaf5239 (2016).
8. Wang, D., Song, C. & Barabási, A.-L. Quantifying long-term scientific impact. *Science* **342**, 127–132 (2013).
9. Guimera, R., Uzzi, B., Spiro, J. & Amaral, L. A. N. Team assembly mechanisms determine collaboration network structure and team performance. *Science* **308**, 697–701 (2005).
10. Wuchty, S., Jones, B. F. & Uzzi, B. The increasing dominance of teams in production of knowledge. *Science* **316**, 1036–1039 (2007).
11. Milojević, S. Principles of scientific research team formation and evolution. *Proc. Natl. Acad. Sci.* **111**, 3984–3989 (2014).
12. Way, S. F., Morgan, A. C., Larremore, D. B. & Clauset, A. Productivity, prominence, and the effects of academic environment. *Proc. Natl. Acad. Sci.* **116**, 10729–10733 (2019).
13. Deville, P. et al. Career on the move: Geography, stratification, and scientific impact. *Sci. Rep.* **4**, 4770 EP – (2014).
14. Wu, L., Wang, D. & Evans, J. A. Large teams develop and small teams disrupt science and technology. *Nature* **566**, 378–382 (2019).
15. Jones, B. F. The burden of knowledge and the “death of the renaissance man”: is innovation getting harder? *Rev. Econ. Stud.* **76**, 283–317 (2009).
16. Page, S. E. *The Diversity Bonus: How Great Teams Pay off in the Knowledge Economy*, vol. 5 (Princeton University Press, 2019).
17. Dong, Y., Ma, H., Tang, J. & Wang, K. Collaboration diversity and scientific impact. Preprint at <https://arxiv.org/abs/1806.03694> (2018).
18. Yegros-Yegros, A., Rafols, I. & D’Este, P. Does interdisciplinary research lead to higher citation impact? the different effect of proximal and distal interdisciplinarity. *PLoS one* **10**, e0135095 (2015).
19. Leskovec, J., Kleinberg, J. & Faloutsos, C. Graph evolution: densification and shrinking diameters. *ACM Trans. Knowl. Discov. Data* **1**, 2–41 (2007).
20. Bhat, U., Krapivsky, P. L., Lambiotte, R. & Redner, S. Densification and structural transitions in networks that grow by node copying. *Phys. Rev. E* **94**, 062302 (2016).
21. Lambiotte, R., Krapivsky, P. L., Bhat, U. & Redner, S. Structural transitions in densifying networks. *Phys. Rev. Lett.* **117**, 218301 (2016).
22. Zipf, G. K. *Human Behavior and the Principle of Least Effort: An Introduction to Human Ecology* (Addison-Wesley Press, Inc., Cambridge, MA, 1949).
23. Lü, L., Zhang, Z.-K. & Zhou, T. Zipf’s law leads to heaps’ law: Analyzing their relation in finite-size systems. *PLoS ONE* **5**, 1–11 (2010).
24. Simini, F. & James, C. Testing heaps’ law for cities using administrative and gridded population data sets. *EPJ Data Sci.* **8**, 24 (2019).
25. Tria, F., Loreto, V., Servidio, V. D. P. & Strogatz, S. H. The dynamics of correlated novelties. *Sci. Rep.* **4**, 5890 EP – (2014).
26. Sinha, A. et al. An overview of Microsoft Academic Service (MAS) and applications. In *Proceedings of the 24th international conference on world wide web*, 243–246 (ACM, 2015).
27. Batty, M. Rank clocks. *Nature* **444**, 592–596 (2006).
28. Iacopini, I., Di Bona, G., Ubaldi, E., Loreto, V. & Latora, V. Interacting discovery processes on complex networks. *Phys. Rev. Lett.* **125**, 248301 (2020).
29. Burghardt, K., Percus, A., He, Z. & Lerman, K. A model of densifying collaboration networks. Preprint at <https://arxiv.org/abs/2101.11056> (2021).
30. Gibrat, R. *Les inegalites economiques; applications: aux inegalites des richesses, a la concentration des entreprises, aux populations des villes, aux statistiques des familles, etc., da une loi nouvelle, la loi de la effet proportionnel.* (Librairie du Recueil Sirey, Paris, 1931).
31. Eeckhout, J. Gibrat’s law for (All) cities. *Am. Econ. Rev.* **94**, 1429–1451 (2004).
32. Axtell, R. L. Zipf distribution of U.S. firm sizes. *Science* **293**, 1818–1820 (2001).
33. Herrmannova, D. & Knoth, P. An analysis of the microsoft academic graph. *D-Lib Magazine* <http://www.dlib.org/dlib/september16/herrmannova/09herrmannova.html> (2016).
34. Depersin, J. & Barthelemy, M. From global scaling to the dynamics of individual cities. *Proc. Natl. Acad. Sci.* **115**, 2317–2322 (2018).
35. Keuschnigg, M. Scaling trajectories of cities. *Proc. Natl. Acad. Sci.* **116**, 13759–13761 (2019).
36. Ribeiro, F. L., Meirelles, J., Netto, V. M., Neto, C. R. & Baronchelli, A. On the relation between transversal and longitudinal scaling in cities. *PLOS ONE* **15**, 1–20 (2020).
37. Bettencourt, L. M. A., Lobo, J., Helbing, D., Kühnert, C. & West, G. B. Growth, innovation, scaling, and the pace of life in cities. *Proc. Natl. Acad. Sci.* **104**, 7301–7306 (2007).
38. Bettencourt, L. M. A. The origins of scaling in cities. *Science* **340**, 1438–1441 (2013).
39. Fix, B. Energy and institution size. *PLOS ONE* **12**, 1–22 (2017).
40. Bettencourt, L. M. A. et al. The interpretation of urban scaling analysis in time. *J. R. Soc. Interface* **17**, 20190846 (2020).

## Acknowledgements

Research was funded by in part by DARPA under contract #W911NF1920271 and by the USC Annenberg Fellowship.

## Author contributions

K.B., Z.H., A.G.P., and K.L. designed research; K.B. and Z.H. performed research; K.B. created and simulated the model; K.B. and Z.H. analyzed data; K.B., Z.H., A.G.P., and K.L. wrote the paper.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s42005-021-00693-2>.

**Correspondence** and requests for materials should be addressed to K.A.B.

**Peer review information** *Communications Physics* thanks the anonymous reviewers for their contribution to the peer review of this work. Peer reviewer reports are available.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher’s note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021

# Supporting Information: The Emergence of Heterogeneous Scaling in Research Institutions

## CONTENTS

Supplementary Note 1: Data	1
Supplementary Note 2: Cumulative versus Yearly Measures	3
Supplementary Note 3: Homogeneous Densification (Null) Model	6
Supplementary Note 4: Scaling of Output	7
Supplementary Note 5: Scaling of Team Size	8
Supplementary Note 6: Comparison Between Data and Simulations	8
Supplementary Note 7: Robustness Check of Simulations	10
Supplementary References	11
References	11

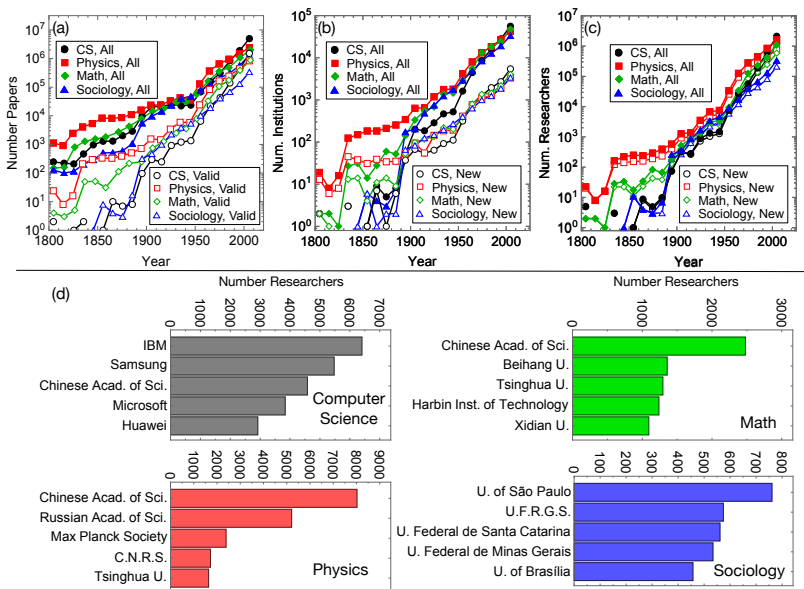
### Supplementary Note 1: Data

We use bibliographic data from Microsoft Academic Graph (MAG), from which researcher names (authors), their institutional affiliations, and references made to other papers have been extracted [1, 2]. MAG data has disambiguated institutions and authors for each paper, allowing us to consider all authors with the same unique identifier to be the same researcher, and similarly for each institution.

The MAG data enables us to measure institution size (the number of published authors affiliated with the institution), productivity (number of papers written), and collaborations (co-authors of the same paper), both within and among institutions. We gather data from papers published in four fields of study between 1800 and 2018: computer science (14,666,855 papers), physics (8,428,923 papers), math (6,192,706 papers), and sociology (4,407,288 papers). Because the metadata for MAG are extracted automatically, many papers have some missing values among extracted names, references, institution, or year published. We find there may be differences in the data collected, as well as the format and even how affiliations are defined, between this data and newer vintages of MAG data. As part of the data cleaning process, we remove papers with missing fields, and also papers with more than 25 authors. These many-authored papers only represent 0.70% of all physics papers, and  $< 0.036\%$  of papers in other fields but are removed because they may be too large to constitute a meaningful collaboration between any individuals. This leaves 3,916,332 computer science papers, 2,494,000 physics papers, 2,370,712 math papers, and 1,115,841 sociology papers. Parsing these data produces 3.99 million (16.7 thousand), 2.79 million (12.3 thousand), 2.17 million (14.1 thousand), and 826 thousand (12.5 thousand) researchers (institutions) in computer science, physics, math, and sociology, respectively. We analyze institutions with at least 5 datapoints, and an increase of at least 10 researchers. In total, we find 4943 (5475), 3642 (4130), 3424 (4445), and 1215 (1745) internal(external) collaboration scaling exponents (where each exponent corresponds to a valid institution) for computer science, physics, math, and sociology, respectively. Separately, we use the same criteria for simulations shown in main text Fig. 2a–b insets, which contain 1031 datapoints out of 3521 institutions each across four simulation realizations.

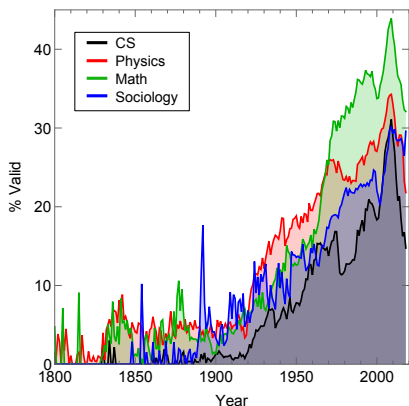
Supplementary Figure 1 shows the descriptive statistics of the data, including the growth of the number of researchers, institutions, and papers published in the four disciplines, and the five largest institutions in each field. Notably, while the largest Physics, Sociology, and Math institutions are universities, the largest computer science institutions are often companies. Figure 2 in the main text demonstrates that institution sizes are broadly distributed with many smaller than 10 researchers, and some larger than  $10^3$ . While Supplementary Figure 1 shows that the largest institutions are intuitive, such as Harvard, we separately check the quality of the data for small institutions. We randomly sampled 44 institutions in each field with fewer than 10 researchers as of 2017 (see Supplementary Data 1). We observe that they tend to be for-profit colleges, community colleges, and institutions without a formal department in the field of interest (e.g., an engineering school with papers in sociology). That said, we see the journals they publish in tend to be well-aligned with the field, therefore the small institutions were not associated with a particular field by mistake. While these are qualitative checks, they nonetheless show that the data and institutions found are reasonable.

We also analyze the quality of MAG’s data over time in Supplementary Figure 2. This figure shows the percentage of papers considered valid (containing year, author, and affiliation). We find surprisingly few papers are considered valid before 1900, while even when the data quality is highest around 2000, the minority of papers are considered



Supplementary Figure 1. Growth of academic disciplines. (a) Number of papers per decade and number of valid papers that contain author, references, institution, and year. (b) Number of new institutions over time. (c) number of new researchers over time. (d) Largest institutions in 2017. C.N.R.S. stands for Centre national de la recherche scientifique and U.F.R.G.S. stands for Universidade Federal do Rio Grande do Sul. scientifique

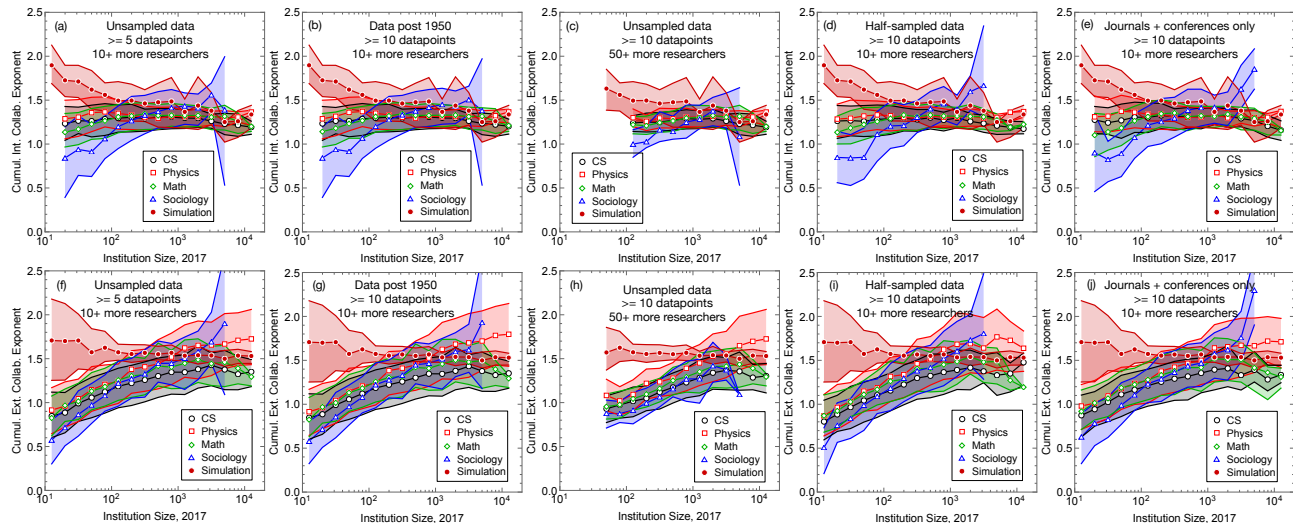
valid. This demonstrates bias in data sampling, and is a potential limitation of our work. Nonetheless, we show in Supplementary Note 4 that the number of papers per author is roughly independent with institution size. If the sampling bias had an increasing preference towards, e.g., large institutions, then this finding would not have held. Furthermore, Supplementary Figure 3 shows that removing half of all data, or removing data pre-1950 (where the percent of valid data is low), or only considering journals and conferences, in contrast to (for example) patents, does not substantially affect the scaling relations. Similarly, analyzing higher-quality data (institutions with more datapoints, and larger changes in their size) does not affect results significantly. In that plot, we show the scaling relations versus institution size in 2017 for all institutions studied. Data is qualitatively and quantitatively the same despite the significant drop in the number of papers studied. These robustness tests suggest that undersampling does not affect our overall conclusions.



Supplementary Figure 2. Percentage of valid papers each year. Valid papers are defined as those with names, years, and affiliations for the authors. We find newer papers, especially after 1900, are more likely to be valid.

Because we include all research output, including journal papers and patents, we check the robustness of our findings when just including the standards of academic research: journals and conference proceedings. In Table 1, we find that between 57%-90% of all documents are in these two categories, with 90% in sociology and only 57% in





Supplementary Figure 3. Collaboration scaling exponents versus  $N$ . (a,f) Unsampled data where each institution has at least 5 datapoints and 10 or more new researchers enter since the first year, (b,g) same as (a) but only data since 1950, (c,h) same as (a) but each institution has at least 10 datapoints, and 50 or more new researchers enter since the first year, (d,i) same as (a) but half of valid papers are removed at random, (e,j) same as (a) but we only consider papers that are labeled as journals and conferences (in contrast to, e.g., patents). Shaded areas represent 50% quantile ranges.

computer science, presumably because research output in that field is often patents. In Supplementary Figures 3e & j and 4, however, we show that our major findings are qualitatively unchanged. Namely, in addition to unchanged collaboration scaling laws in Supplementary Figure 3e & j, we see a strong Heaps' law and Zipf's law that looks very similar to the main text (Supplementary Figure 4).

Supplementary Table 1. Journals, conferences and other papers. Columns are labeled as follows: research field (Field), number of papers MAG defines as a journal or conference (Journal & Conference), number of papers not under these categories (Other), total number of papers (Total), and the percentage of papers that are journals or conferences (% J&C).

Field	Journal & Conference	Other	Total	% J&C
CS	1571485	1172424	2743909	57%
Physics	1953468	287740	2241208	87%
Math	1567928	401392	1969320	80%
Sociology	913549	101146	1014695	90%

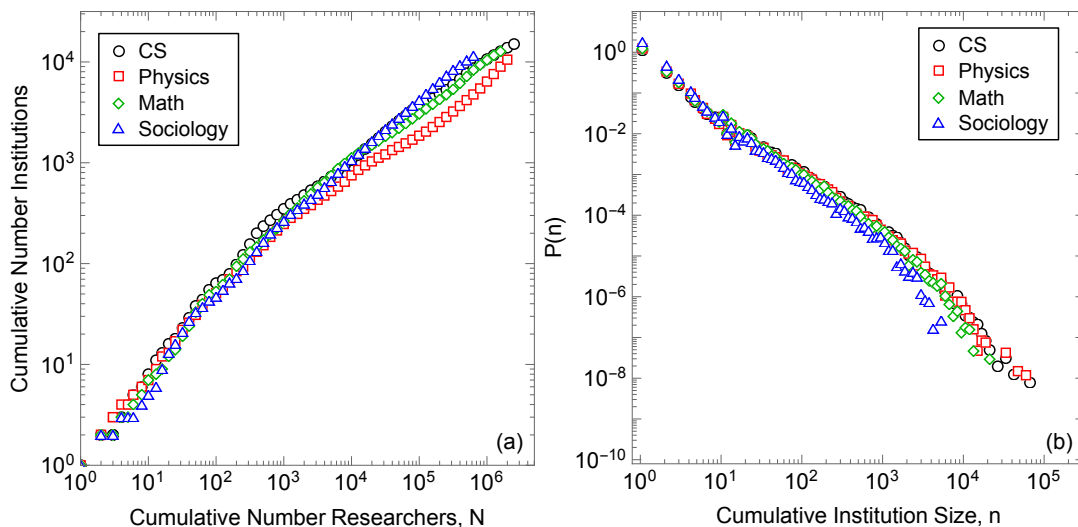
Finally, we treat each author with multiple simultaneous affiliations (i.e., multiple affiliations in a given paper) as separate researchers in each institution, and count their affiliations as cross-institution collaborations. This latter point accounts for realistic scenarios in which authors promote greater cross-institution collaborations simply by having multiple affiliations, which would not be accounted for without this definition. These authors are rare, accounting for 0.5-3.5% of all authors, see Supplementary Figure 5, therefore other ways to incorporate these authors into the wider dataset will produce qualitatively similar results.

## Supplementary Note 2: Cumulative versus Yearly Measures

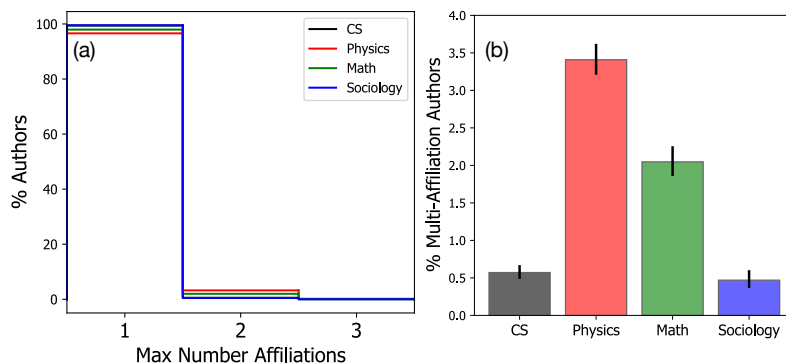
While the main text measured the cumulative size of institutions and collaborations, the findings are qualitatively the same if the growth of research institutions was measured on a year-to-year basis. Institution size is therefore the number of active authors affiliated with that institution who published in that particular year. Collaborations were similarly based on papers published that year, etc.

Supplementary Figure 1 shows the growth of four academic disciplines, including (a) the number of published papers, (b) the number of institutions and (b) the number of researchers each year all increase exponentially, regardless of whether these are measured on the cumulative (all) or year-to-year basis.

Supplementary Figure 6 further demonstrates the robustness of our results, regardless of how they are measured. This figure shows that the cumulative institution size, cumulative number of internal and external collaborations



Supplementary Figure 4. Heaps' law and Zipf's law based on documents labeled as journals or conference proceedings. (a) Cumulative number of institutions versus number of researchers for computer science, physics, math, and sociology. (b) Institution size distribution in 2017.



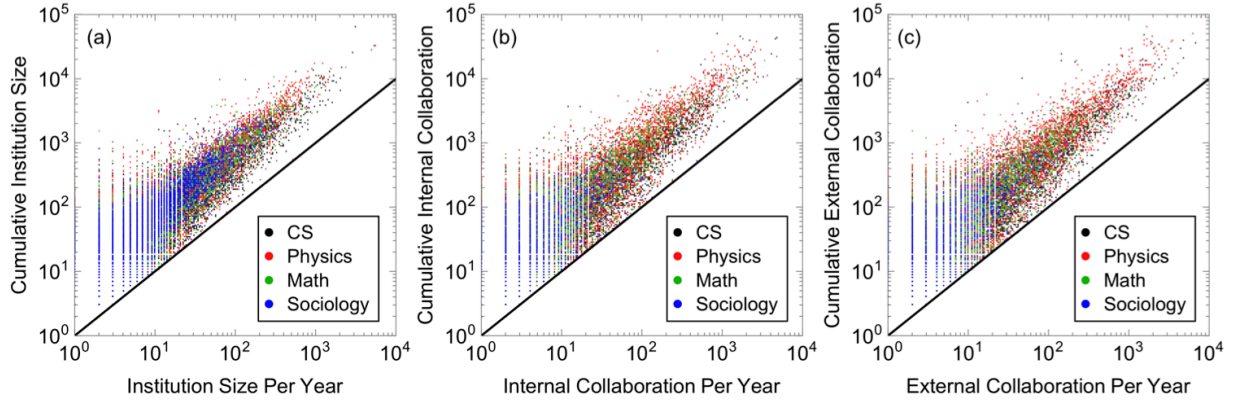
Supplementary Figure 5. Multi-affiliation authors. (a) Distribution of the maximum number of affiliations for each author. (b) Percent of authors with whose maximum number of affiliations was greater than one. Results based on a random sample of 10,000 manuscripts within each field.

Supplementary Table 2. Cumulative Vs. year-to-year Spearman Correlations

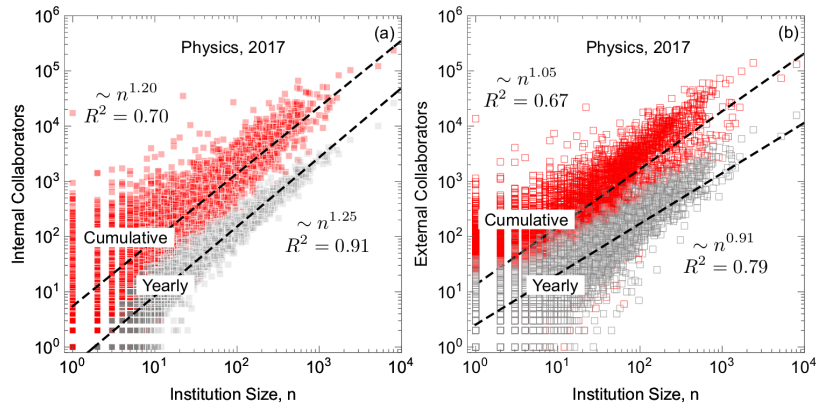
Data	Size	Internal Collab.	External Collab.
CS	0.85	0.83	0.83
Physics	0.85	0.84	0.84
Math	0.84	0.82	0.82
Sociology	0.81	0.71	0.71

are well correlated with their year-to-year values. The correlations are in Table 2, where Spearman correlations are 0.7–0.8 or higher. Comparison between yearly and cumulative results can also be seen in Supplementary Figure 7 where we show cross-sectional collaboration scaling for researchers active in 2017 as well as cumulative collaboration scaling for cumulative institution size.

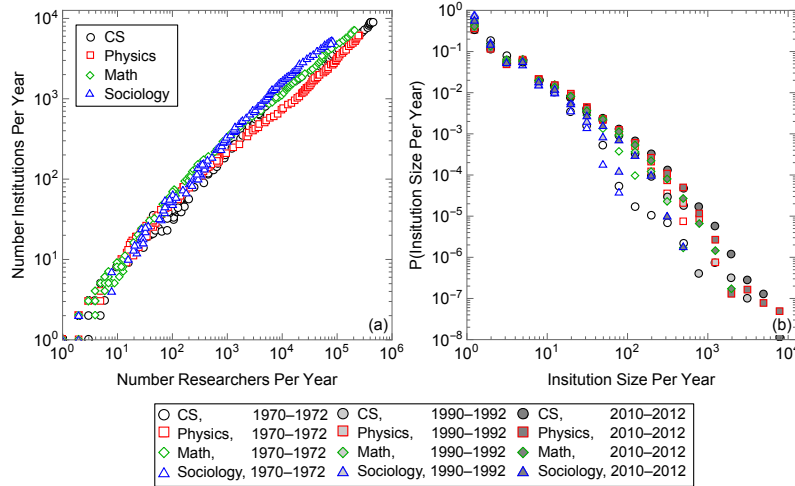
Supplementary Figure 8 reproduces Fig. 4 in the main text, except we calculate the yearly number of researchers, institutions, and institution size. Supplementary Figure 8a shows the number of institutions in a given year versus the number of researchers in a given year. We see, much like in the main text, a sub-linear scaling between the number of institutions and researchers. Supplementary Figure 8b shows the institution size distribution. Importantly, the institution size distribution might change over time, therefore we plotted the institution size per year for 1970–1972,



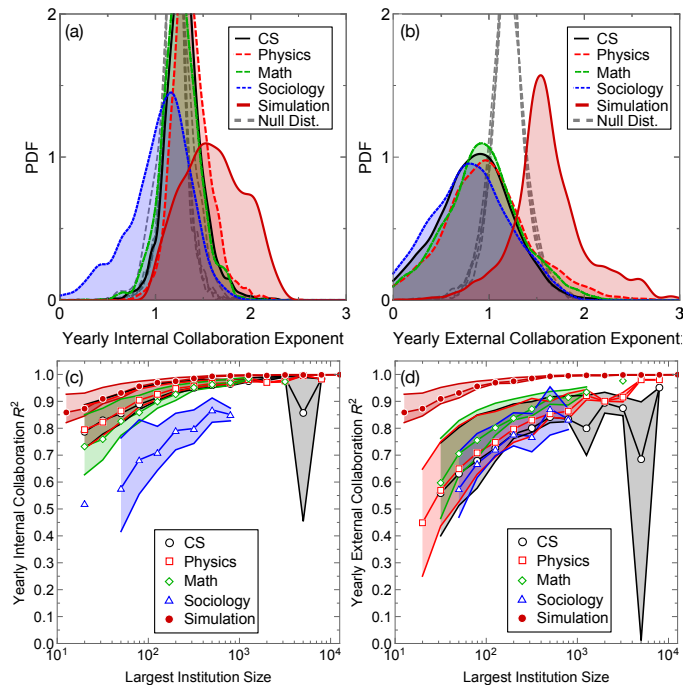
Supplementary Figure 6. Cumulative versus year-to-year data. (a) Cumulative versus year-to-year institution size, (b) cumulative versus year-to-year internal collaborations, and (c) cumulative versus year-to-year external collaborations.



Supplementary Figure 7. Cross-sectional analysis of the scaling of collaborations. (a) Cumulative internal collaborators and (b) cumulative external collaborators versus institution size (as of 2017),  $n$ , for physics.



Supplementary Figure 8. Institution growth statistics by year. (a) Number of institutions versus number of researchers each year for computer science, physics, math, and sociology. (b) Institution size distribution for 1970–1972, 1990–1992, and 2010–2012 for the same fields.



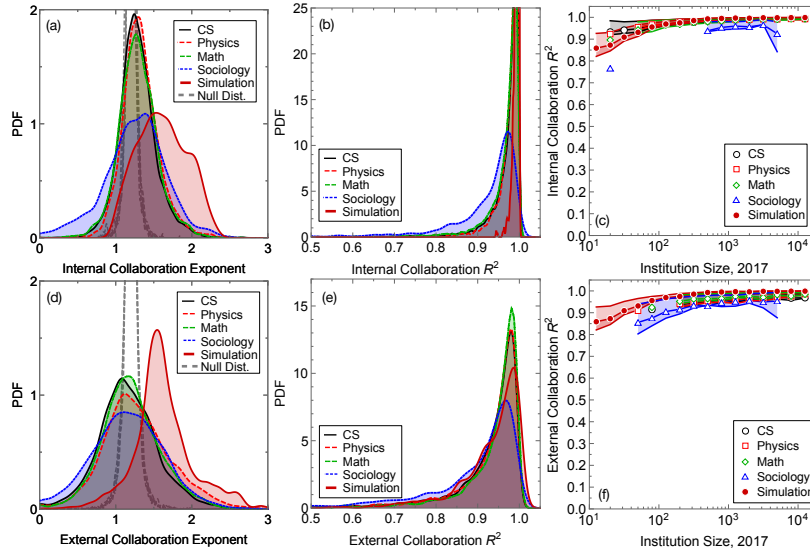
Supplementary Figure 9. Scaling exponents of yearly collaborations versus yearly institution size. (a) Internal and (b) external collaboration scaling exponents.  $R^2$  is lower for smaller institutions but quickly approaches 1.0 for (c) internal and (d) external collaborations per year.

1990–1992, and 2010–2012, and found the distribution was extremely stable in time and across fields.

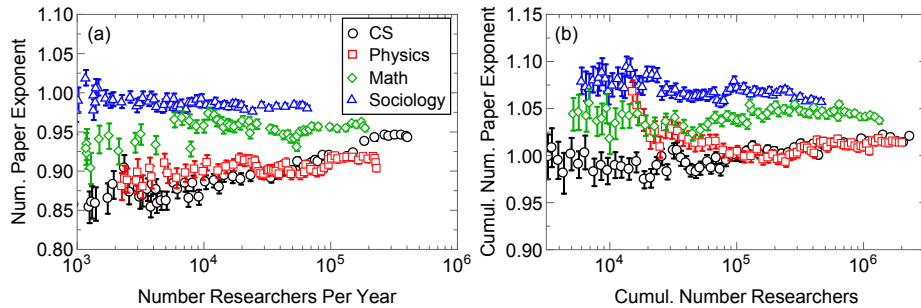
Supplementary Figure 9 shows the scaling exponents of yearly collaborations versus yearly institution size, but plotted against each institution’s largest yearly size, as well as the quality of the linear model fits ( $R^2$ ) versus largest institution size. We see, much like the main text, a large variance in the exponent values, but that they do not dramatically change with institution size over several orders of magnitude. That said, the scaling law  $R^2$  is higher for large institutions in agreement with the expectation that the scaling law works best in the large- $n$  limit of institution sizes. This is also similar to what was found in the main text for cumulative sizes.

### Supplementary Note 3: Homogeneous Densification (Null) Model

To understand whether the observed heterogeneity in longitudinal scaling laws is due to statistical noise, we create a null model that assumes a homogenous scaling exponent for all institutions, whose results are shown in Supplementary Figure 10. To create the null model, we fit a scaling law for each institution, keeping the residual values,  $r_i$ , along with their  $x_i$  positions, creating a set of pairs  $\{x_i, r_i\}$ . The null model *homogeneous scaling law*,  $\beta_{\text{null}}$ , is a reasonable fixed value. For example, we can set it to the average scaling laws,  $\beta_i$ , across all institutions, weighted by the inverse of the standard error squared,  $1/\sigma_{\beta_i}^2$ , although some missing or irregular data may complicate this averaging. For this reason, we set the value at exactly 1.2 in Supplementary Figures 9 & 10a,d. For each institution, we randomly permute the residuals  $\{x_i, r_i\} \rightarrow \{x_i, r_j\}$  to create new data:  $\{x_i, y_{\text{null}}\} = \{x_i, x_i\beta_{\text{null}} + r_j\}$ . Because of random permutation of residuals, we assume the data are homoscedastic and residuals are uncorrelated, but make no other assumptions, not even whether the residuals are normally distributed. After refitting each new set of points  $\{x_i, y_{\text{null}}\}$  for each institution, we expect the new null model coefficient for each institution to fluctuate around  $\beta_{\text{null}}$  due to noise. To see whether the distribution of null model coefficients differs from the empirically derived coefficients, we use the Kolmogorov-Smirnov test on these two distributions [3]. We find almost invariably that the two distributions differ with p-value  $< 0.001$ .



Supplementary Figure 10. Collaboration scaling exponent distributions and null models. (a) Internal and (d) external collaboration scaling exponent distributions, (b,e)  $R^2$  distribution, and (c,f)  $R^2$  versus  $N$  with shaded areas representing 50% quantile ranges. Simulation parameters are  $\rho$  equals 4,  $\nu$  equals 2,  $\mu_p$  equals 0.6, and  $\sigma_p$  equals 0.25.

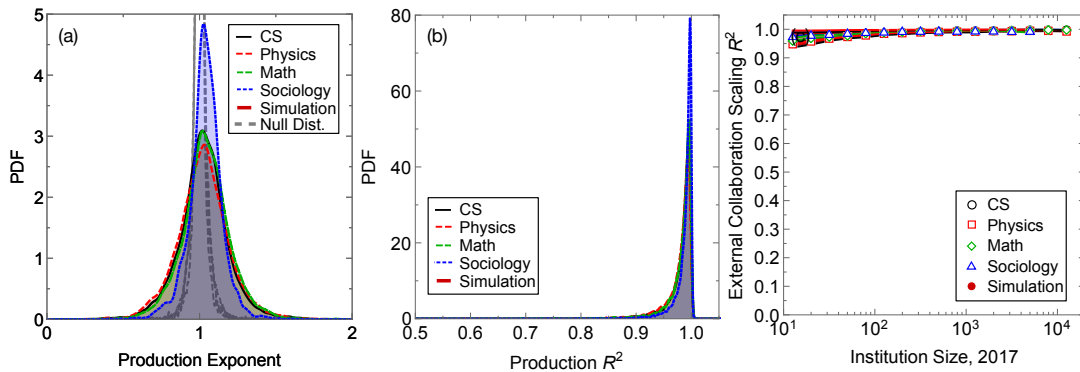


Supplementary Figure 11. Scaling exponents of paper output over time from cross-sectional analysis. (a) Paper output per year and (b) cumulative paper output. We see that the number of papers scales linearly with institution size regardless of the field.

#### Supplementary Note 4: Scaling of Output

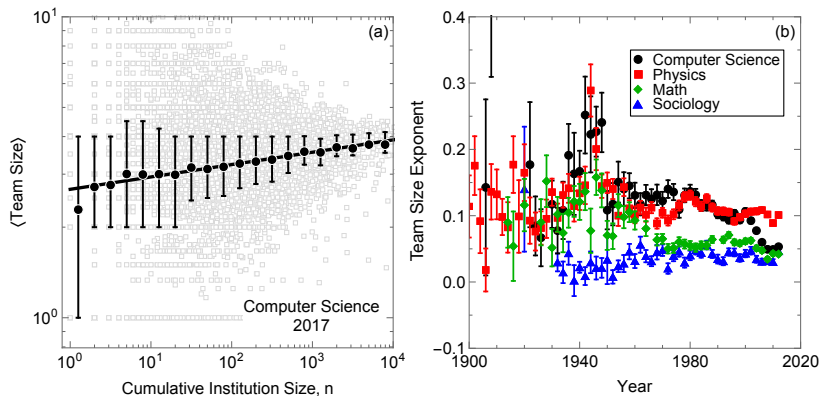
We define institution output as the cumulative number of papers written by researchers affiliated with that institution. Cross-sectional scaling laws are surprisingly stable in time, with a value of almost exactly 1.0, as shown in Supplementary Figure 11. This means that the output per person is independent of institution size. This holds also for different disciplines, regardless of whether we look at annual output or cumulative output. We explore the scaling laws in longitudinal data as well in Supplementary Figure 12. These results also show approximately linear scaling relationships.

Supplementary Figure 12 shows the scaling exponents of institution output versus institution size. The scaling exponents are centered around 1.0, although these scaling laws differ between institutions. This suggests, surprisingly, that paper output per researcher is approximately independent of institution size. That being said, when we compare the longitudinal data to a homogeneous scaling null model, we find that institutions have a greater variance in their scaling laws than the null model predicts. This means that some institutions create slightly more papers per person as the institution grows, while others show a reduction in output. The overall effect, however, appears to be subtle. Overall, institution size appears to affect collaborations much more than output.



Supplementary Figure 12. Probability distribution of longitudinal scaling exponents for institution output (number of papers) across disciplines. (a) Scaling exponent distributions where gray dashed lines are the PDF of the null model in which all institutions within that field have the same scaling law. (b)  $R^2$  distribution, and (c)  $R^2$  versus institution size,  $n$ .

### Supplementary Note 5: Scaling of Team Size



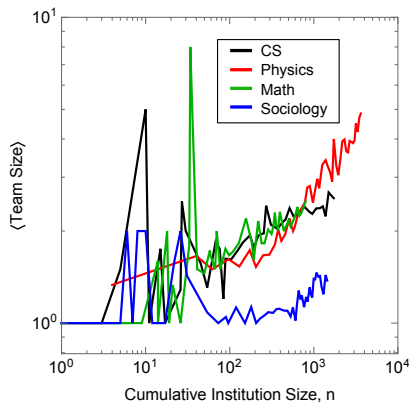
Supplementary Figure 13. Cross-sectional analysis of the scaling of team size. (a) Example cross-sectional scaling for computer science in 2017, and (b) scaling exponent for each field over time.

Team size measures the number of co-authors on a single paper. Prior work has shown that team size has grown over time, with papers produced by larger teams getting more citations compared to papers written by smaller teams [4]. We analyze whether institution size benefits team size via both cross-sectional and longitudinal analysis. As shown in Supplementary Figure 13, we find that team size would seem to scale positively with institution size but the scaling relations can vary significantly in time (Supplementary Figure 13b). While the scaling laws are not universal, Supplementary Figure 13a shows that the fit to a line is remarkably strong.

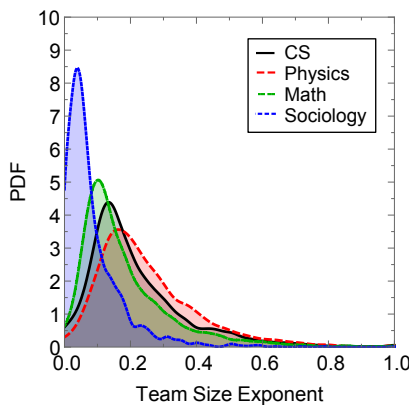
When we analyze data longitudinally, we see a more complete picture. Namely, Supplementary Figure 14 shows that team size scales positively, but the scaling laws differ significantly between institutions. We explore this more thoroughly in Supplementary Figure 15 where we plot a distribution of scaling exponents for each field, whose distribution is wider than the null model (KS-test  $p$ -value  $< 0.001$ ). We see that there is both significant heterogeneity and generally larger scaling exponents than cross-sectional analysis would predict. For example, while cross-sectional analysis shows Physics has a scaling law of about 0.1, longitudinal analysis instead shows that each institution scales with an exponent of roughly 0.2.

### Supplementary Note 6: Comparison Between Data and Simulations

For the rest of the SI, we will just discuss cumulative results. Supplementary Figure 16a is similar to the main text by showing that cumulative cross-scaling exponents vary, but here the x-axis is the total number of researchers who



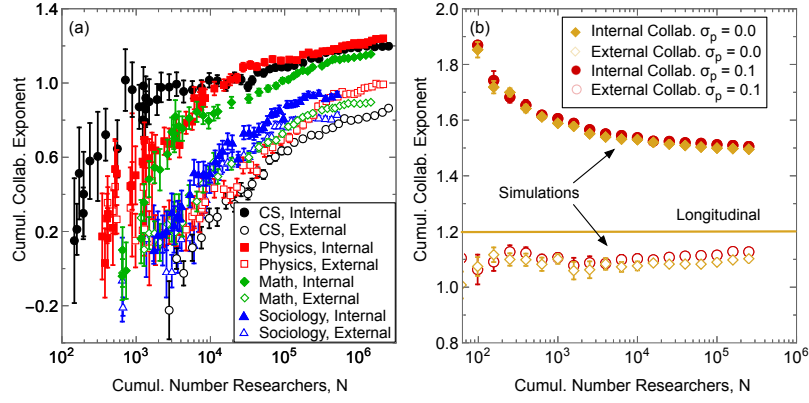
Supplementary Figure 14. Longitudinal scaling examples for team size.



Supplementary Figure 15. Distribution of longitudinal scaling exponents for team size for computer science, physics, math, and sociology.

have authored a paper up until that date. This slightly unusual x-axis allows us to compare these results to cross-sectional scaling in simulations shown in Supplementary Figure 16b. Parameters in these simulations are the same as the main text, with  $\nu = 2$ ,  $\rho = 4$ ,  $\mu_p = 0.6$ , but  $\sigma_p = 0.0 - 0.1$ . For these parameters, we discover that, much like in the data, internal scaling laws are higher than external scaling laws, even though, for each institution, both should be centered around the horizontal lines labeled “longitudinal” (which corresponds to the mean values in the longitudinal analysis). We also notice that, like the data, the exponents vary as a function of the total number of researchers. These simulations do not just allow us to reproduce results, however, but we can make contrapositive hypotheses. For example, what would the statistics look like if there was no statistical variation in the longitudinal scaling laws? To better understand this, we let  $\sigma_p = 0.0$  in Supplementary Figure 16b, and discover that the results are quantitatively almost exactly the same. If institutions had the exact same scaling laws and the exact same constant coefficients, then the cross-sectional and longitudinal scaling laws would be the same. These discrepancies point to either finite size effects or different constant coefficients are dominant factors in explaining why cross-sectional scaling and longitudinal scaling laws differ and vary in time, at least in simulations. We hypothesize similar effects in empirical data as well, although there are qualitative differences in data, such as scaling laws increasing rather than decreasing which point to limitations of the simulations. We also show in Supplementary Figure 3 that internal collaboration scaling laws do not vary significantly with institution size over several orders of magnitude. External collaboration scaling laws, in contrast, vary weakly with institution size as of 2017. Moreover, there is significant variance in these scaling laws that our simulation can capture. These results qualitatively agree with the model. Although empirically, the dependence of external collaboration scaling with  $N$  does not agree with the simulations, theoretical analysis in Fig. 9 of Burghardt et al. [5] show agreement with these findings.

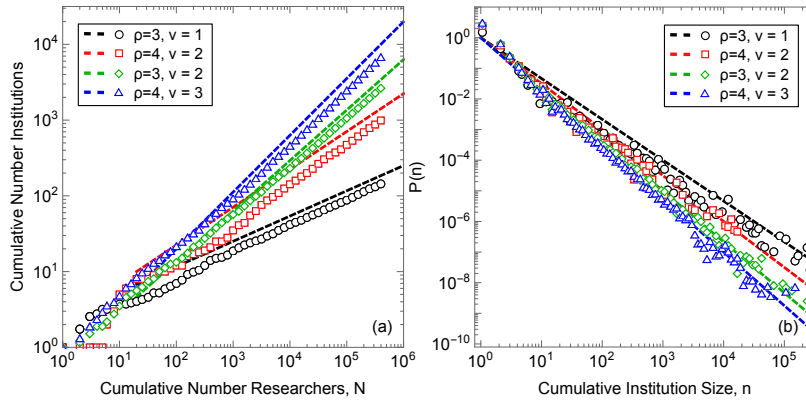
Focusing on longitudinal scaling, however, we observe in Supplementary Figure 10 that both the data and simulations are well-characterized by linear relations in log-log space. Namely, the figure shows that a linear fit of  $\log(\text{collaborations})$  versus  $\log(\text{institution size})$  have an  $R^2$  value of nearly 1 for each institute. If their size as of 2017



Supplementary Figure 16. Cross-sectional analysis of the scaling of collaborations. (a) Internal collaboration and external collaboration scaling exponents versus the cumulative number of researchers. Error bars are standard errors. (b) Simulation cross-institution collaboration scaling versus the number of researchers. Black lines: internal collaboration, red dashed lines: external collaboration. Shaded regions are 95% confidence intervals in the mean across different simulation realizations.

is large, then  $R^2$  is even closer to 1, in agreement with what we should expect in the thermodynamic limit (where finite size effects are negligible). The data is therefore well-characterized as a power law, but these power law values vary between institutions, as shown in main text Fig. 2.

Theory surrounding the Polya's urn portion of our model is discussed in detail in previous work [6]. Nonetheless, we check the robustness of this theory in Supplementary Figure 17. We find excellent agreement between the theory and simulation, demonstrating that, even for finite sizes, the theory they developed accurately explains the simulation patterns.

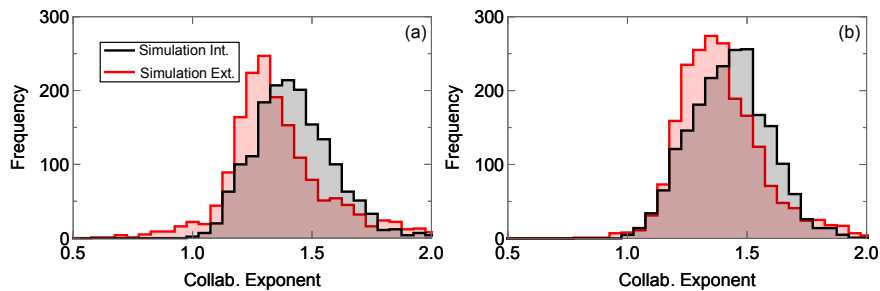


Supplementary Figure 17. Simulation of (a) Heaps law and (b) Zipf's for various values of  $\rho$  and  $\nu$ . Dashed lines are theoretical scaling exponents.

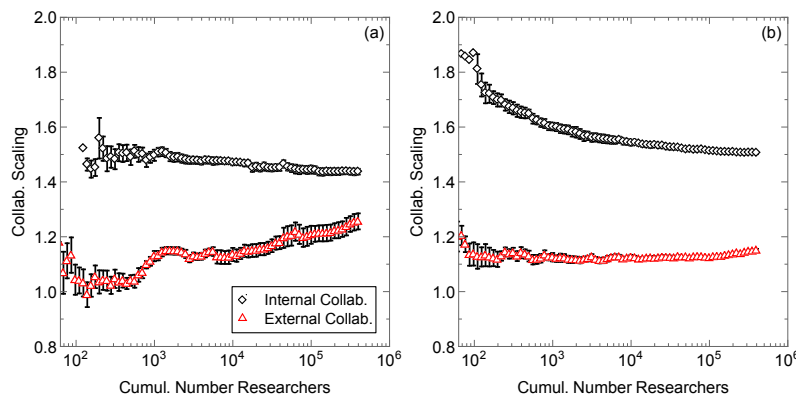
### Supplementary Note 7: Robustness Check of Simulations

We might wonder whether our model is sensitive to stochastic variations in how the model behaves. For example, we might ask whether changing the number of initial collaborators from 1 to a range of values will affect results. To this end, we made an additional model in which the number of initial internal and external collaborators was Poisson distributed, with  $\lambda = 1$  (i.e., on average one internal and one external collaborator). This will not affect the institution formation, but it might affect the institution growth, e.g., the longitudinal collaboration scaling exponents. Importantly, Bhat et al. and Lambiotte et al. shows that number of links over time are not self-averaging [7, 8], therefore initial conditions greatly affect the final number of links. Supplementary Figures 18 & 19 show our results. In Supplementary Figure 18, we find that, while there are slightly more outliers in the scaling exponent distribution, results are quantitatively very similar. In Supplementary Figure 19a we find that  $\lambda = 1$  external collaboration





Supplementary Figure 18. Internal and external longitudinal collaboration exponents for alternative simulation models. (a) Internal and external exponents for simulations with  $\lambda = 1$  Poisson distributed numbers of initial collaborators (on average one internal collaborator, and one external collaborator). (b) The same histograms for the current simulation with exactly one internal and one external collaborator.



Supplementary Figure 19. Internal and external cross-sectional collaboration exponents for alternative simulation models. (a) Internal and external exponents versus the cumulative number of researchers for simulations with  $\lambda = 1$  Poisson distributed numbers of initial collaborators (on average one internal collaborator, and one external collaborator). (b) The same figure for the current simulation with exactly one internal and one external collaborator.

cross-sectional scaling exponents increase with the cumulative number of researchers, more alike to what we see in empirical data (main text Fig. 2), and the internal collaboration exponents are stationary. In Supplementary Figure 19b, however, we find that the external collaboration exponents are mostly stationary in the original form of the model, while internal collaboration exponents decrease with the cumulative number of researchers.

## SUPPLEMENTARY REFERENCES

- 
- [1] A. Sinha, Z. Shen, Y. Song, H. Ma, D. Eide, B.-j. P. Hsu, and K. Wang, in *Proceedings of the 24th international conference on world wide web* (ACM, 2015) pp. 243–246.
  - [2] D. Herrmannova and P. Knoth, *D-Lib Magazine* **22** (2016).
  - [3] F. J. Massey, Jr., *Journal of the American Statistical Association* **46**, 68 (1951).
  - [4] S. Wuchty, B. F. Jones, and B. Uzzi, *Science* **316**, 1036 (2007).
  - [5] K. Burghardt, A. Percus, Z. He, and K. Lerman, arXiv preprint:arXiv:2101.11056 (2021).
  - [6] F. Tria, V. Loreto, V. D. P. Servedio, and S. H. Strogatz, *Scientific Reports* **4**, 5890 EP (2014).
  - [7] R. Lambiotte, P. L. Krapivsky, U. Bhat, and S. Redner, *Phys. Rev. Lett.* **117**, 218301 (2016).
  - [8] U. Bhat, P. L. Krapivsky, R. Lambiotte, and S. Redner, *Phys. Rev. E* **94**, 062302 (2016).