

# 1 Identification of Successful Mentoring Communities using Network-based Analysis of 2 Mentor-Mentee Relationships across Nobel Laureates

3 Julia H. Chariker<sup>a,b</sup>, Yihang Zhang<sup>b</sup>, John R. Pani<sup>a</sup>, & Eric C. Rouchka<sup>b,c\*</sup>

4

## 5 Abstract

6 Skills underlying scientific innovation and discovery generally develop within an academic  
7 community, often beginning with a graduate mentor's laboratory. In this paper, a network  
8 analysis of doctoral student-dissertation advisor relationships in The Academic Tree is used to  
9 identify successful mentoring communities in high-level science, as measured by number of  
10 Nobel laureates within the community. Nobel laureates form a distinct group in the network with  
11 greater numbers of Nobel laureate ancestors, descendants, mentees/grandmentees, and local  
12 academic family. Subnetworks composed entirely of Nobel laureates extend across as many as  
13 four generations. Successful historical mentoring communities were identified centering around  
14 Cambridge University in the latter 19<sup>th</sup> century and Columbia University in the early 20<sup>th</sup>  
15 century. The current practice of building web-based academic networks, extended to include a  
16 wider variety of measures of academic success, would allow for the identification of modern  
17 successful scientific communities and should be promoted.

## 18 Background

19 High achievement in intellectual innovation has been measured in part with the awarding  
20 of prestigious honors, such as the Nobel Prize. From the first awards in 1901 through the awards  
21 in 2015, a total of 573 prizes have been awarded to 900 laureates, including 875 individuals and  
22 25 organizations<sup>1</sup>. To some extent, the scientific knowledge and skill underlying these  
23 achievements has been transmitted across generations through person-to-person academic  
24 mentoring, and much attention has been given to individual mentoring relationships<sup>2-6</sup>.  
25 However, interaction within a scientific laboratory extends beyond the mentor-mentee  
26 relationship. Laboratories make up a community of researchers in which knowledge and skill is  
27 shared within and across generations through relationships between academic siblings and  
28 between mentees and their grandmentors. Within the realm of high-level science, successful

<sup>a</sup>Department of Psychological and Brain Sciences, Life Sciences Building, Room 317, University of Louisville, Louisville, KY, 40292, USA.

<sup>b</sup>KBRIN Bioinformatics Core, 522 East Gray Street, University of Louisville, Louisville, KY, 40202, USA.

<sup>c</sup>Department of Computer Engineering and Computer Science, Duthie Center for Engineering, Room 208, University of Louisville, Louisville, KY, 40292, USA.

29 mentoring communities can provide clues to creating fertile environments for scientific  
30 innovation.

31         The contribution of mentoring to academic success is difficult to isolate within an entire  
32 scientific population because additional factors, such as level of institutional resources and  
33 student talent, also vary across training situations. Successful researchers generally attract more  
34 federal and institutional resources and more talented students, and success begets success.  
35 However, when looking for successful mentors and mentoring communities within a subset of  
36 high achievers, such as Nobel laureates, these factors should operate somewhat equally.  
37 With modern technology and the benefits of crowd-sourcing, generations of mentoring  
38 relationships are now represented in networks such as The Academic Tree <sup>7,8</sup>, greatly facilitating  
39 the study of mentoring on a large scale. The Academic Tree is a vast crowd-sourced network  
40 containing mentor-mentee relationships across several interconnected domains of science. If  
41 Nobel laureates are a distinctive group due in some way to mentoring communities, greater  
42 connectedness should be found among them in the network. Otherwise, we should find Nobel  
43 laureates randomly dispersed across the network. If Nobel laureates are a distinctive group, and  
44 quality of mentoring plays an important role in their success, it should be possible to identify  
45 particularly successful mentoring communities within the group of laureates (i.e., the “best of the  
46 best”).

47         We looked for connectedness among Nobel laureates in The Academic Tree by asking  
48 whether they have a greater number of Nobel laureate academic family members than non-Nobel  
49 laureates have. We restricted our analysis to doctoral student-advisor relationships and assessed  
50 academic family structure in several ways. We examined the number of Nobel laureate ancestors  
51 for each individual as well as the number of local and global descendants. Local descendants  
52 covered two generations in the network and included mentees and grandmentees, whereas global  
53 descendants comprised all generations of descendants. To identify more dispersed mentoring  
54 communities, we looked at the number of Nobel laureates within each individual’s local  
55 academic family, including three generations in all directions in the network. Three generations  
56 encompassed an individual’s mentor, grandmentor, great-grandmentor, mentees, grandmentees,  
57 great-grandmentees, sibling, aunts, and uncles. We compared the outcomes of this analysis to  
58 results obtained from many topologically identical networks in which Nobel status was randomly  
59 assigned across all individuals in each of the networks.

60 Nobel laureates appear to be a distinct group with a greater number of Nobel laureate  
61 family members than non-Nobel laureates have on all measures. In addition, several historical  
62 scientific communities exist with high concentrations of Nobel laureates. In some instances,  
63 Nobel laureates are directly connected to one another over three and four generations of  
64 scientists. Biographical and historical accounts offer the only access to characteristics associated  
65 with these successful communities. However, with the expansion of current network databases to  
66 include a variety of performance measures for all scientists, new methods could be used to  
67 identify modern scientific communities and to study them more directly.

## 68 **Results**

69 *Descriptive Summary.* As can be seen in Table 1, the distributions for the number of academic  
70 family members and the number of Nobel laureate academic family members are positively  
71 skewed on all measures. On some measures, the range was quite large, prompting a closer look  
72 at The Academic Tree. For example, the range for number of descendants extended to 2,628 for  
73 Nobel laureates and 13,620 for non-Nobel laureates. However, academic lineages have been  
74 recorded across several centuries, justifying these numbers. For example, Michele Savonarola, a  
75 non-Nobel laureate physician scientist, practicing in the 15<sup>th</sup> century, has 13,620 descendants, 73  
76 Nobel descendants, and 0 ancestors in the network. Wilhelm Friedrich Ostwald, a Nobel Prize  
77 winner in chemistry (1909), has five ancestors and the highest number of descendants in the  
78 Nobel laureate group at 2,628. In terms of mentees/grandmentees, it appears that Robert B.  
79 Woodward, a Nobel Prize winner in chemistry (1965), has 213 mentees/grandmentees recorded,  
80 one of whom is a Nobel laureate. In the non-Nobel laureate group, Gilbert Stork, a Professor of  
81 Chemistry Emeritus at Columbia University, has 149 mentees/grandmentees, also with one  
82 Nobel laureate among them. The range for number of local academic family was also quite large.  
83 Robert Woodward, a Nobel laureate in chemistry (1965), has the largest local family with 558  
84 members, 5 of whom are Nobel laureates. Of the non-Nobel laureates, Robert T. Paine, professor  
85 emeritus of zoology at The University of Washington has 446 local family. The distributions for  
86 all measures are displayed in Fig. S1.

87

88

89

90

91 **Table 1.** The range and median for number of Nobel laureate academic family and total number  
 92 of academic family across all measures for Nobel laureates (NL) and non-Nobel laureates (Non-  
 93 NL). The correlation between number of academic family and number of Nobel laureate  
 94 academic family for each measure is also displayed. Ancestors refer to individuals moving  
 95 backward in the directed network, and descendants are all individuals moving forward in the  
 96 network. M/GM refers to the number of mentees and grandmentees (two generations forward),  
 97 and local family refers to the number of individuals within 3 generations forward and backward  
 98 in the network.

99

	Number of NL Academic Family		Number of Academic Family		Corr. Between Academic Family and NL Academic Family
	NL Range(Mdn)	Non-NL Range(Mdn)	NL Range(Mdn)	Non-NL Range(Mdn)	Spearman's <i>r</i>
Ancestors	0 - 6 (0)	0 - 8 (0)	0 - 75 (7.5)	0 - 131 (9)	0.33***
Descendants	0 - 21 (0)	0 - 73 (0)	0 - 2,628 (11)	0 - 13,620 (0)	0.24***
M/GM	0 - 8 (0)	0 - 8 (0)	0 - 213(5)	0 - 149 (0)	0.17***
Local Family	0 - 18 (2)	0 - 17 (0)	3 - 558 (31)	3 - 446 (22)	0.17***

100 \*\*\*  $p < 0.0001$

101

102 One reason a strong positive skew was found for ancestors, descendants, and  
 103 mentees/grandmentees was due to the nature of network data, where some individuals serve as  
 104 source nodes without ancestors and other individuals serve as sink nodes without  
 105 descendants. This increases the number of outcomes measuring zero in the data. In this case,  
 106 having zero Nobel family members is a result of having zero family members. Alternatively, a  
 107 number of individuals in the network have ancestors or descendants, but none of them are Nobel  
 108 laureates. Clearly, there are two possible sources for zero Nobel family members. In the analysis  
 109 of ancestors who were Nobel laureates, for example, there were 2,890 individuals with no  
 110 ancestors, and thus no Nobel laureate ancestors. On the other hand, there were 37,608 individuals  
 111 with ancestors, none of whom were Nobel laureates. Similarly, in the analysis of descendants  
 112 who were Nobel laureates, there were 40,044 individuals having no descendants. At the same  
 113 time, there were 16,869 individuals with immediate descendants and 17,197 individuals with  
 114 mentees/grandmentees, none of whom were Nobel Prize winners.

115 **Approach to Analysis.** Zero-inflated regression models were developed for analyzing data with  
116 two possible sources for zero outcomes. With this approach, two models are estimated, a zero-  
117 inflation model and a count model<sup>9,10</sup>. The zero-inflation model is estimated first, using a  
118 binomial model to estimate the probability of excess zeros in the data (i.e., a zero outcome due to  
119 the absence of family). Once this probability is estimated, the probability for the remaining  
120 outcomes is estimated using a Poisson or negative binomial model, whichever is appropriate. In  
121 the current paper, zero-inflated models were used to control for excess zeros in estimating the  
122 number of Nobel laureate ancestors, descendants, and mentees/grandmentees. There was no need  
123 for this in estimating the number of local Nobel laureate family, because inclusion in a connected  
124 network necessarily meant that at least one family connection existed.

125 For all four analyses, negative binomial models were chosen to adjust for greater than  
126 expected dispersion in the data (i.e., a high variance to mean ratio). Spearman's correlations (see  
127 Table 1) indicated that the number of Nobel laureate family members was positively related to  
128 the size of the academic family. Therefore, in each case, the size of the academic family was  
129 entered along with Nobel status as a predictor of the size of the Nobel laureate academic family.  
130 As described in the method, the significance level for each analysis was adjusted by comparing  
131 the observed test statistics with a distribution of expected test statistics, derived from 1,000  
132 topologically identical networks, each with a random permutation of Nobel status. The  
133 regression model coefficients and the distributions of random coefficients used to adjust the  
134 significance levels of predictors in the models are available in Table S1 and Fig. S2, respectively.

135 **Regression Model Outcomes.** Nobel laureates had a greater number of Nobel laureate ancestors  
136 than non-Nobel laureates did, suggesting that Nobel laureate mentorship may play a role in the  
137 development of future Nobel Prize winners (adjusted  $p = 0.003$ ). However, the number of  
138 academic ancestors was not a significant predictor of the number of Nobel ancestors (adjusted  $p$   
139  $= 0.389$ ). Similarly, Nobel laureates had a greater number of Nobel laureate descendants than  
140 non-Nobel laureates did (adjusted  $p < 0.001$ ) with number of descendants not significantly  
141 predicting number of Nobel laureate descendants ( $p = 0.143$ ).

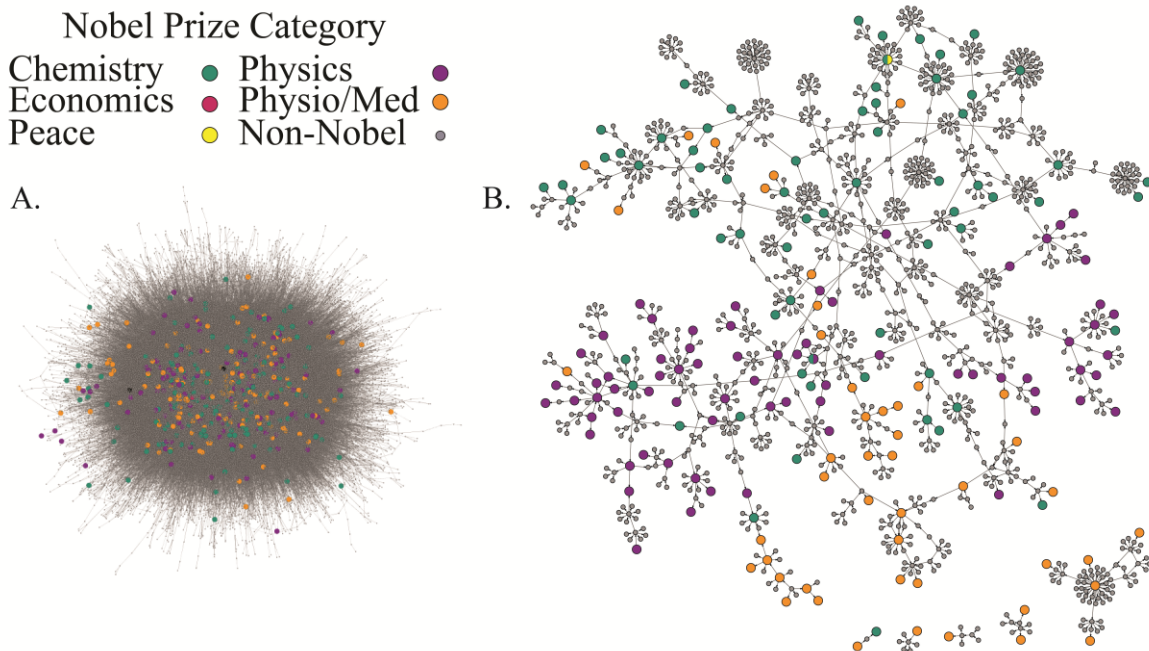
142 In contrast to the previous two results, the number of mentees/grandmentees did serve as  
143 a significant predictor of number of Nobel laureate mentees/grandmentees (adjusted  $p < 0.001$ ).  
144 Still, after controlling for family size, Nobel laureates had a greater number of Nobel laureate  
145 mentees and grandmentees than did non-Nobel laureates (adjusted  $p < 0.001$ ). Finally, Nobel

146 laureates also had a greater number of local Nobel Laureates in their academic family than did  
147 non-Nobel laureates (adjusted  $p < 0.001$ ). The number of local academic family members did not  
148 significantly predict the number of Nobel laureates (adjusted  $p < 0.964$ ).

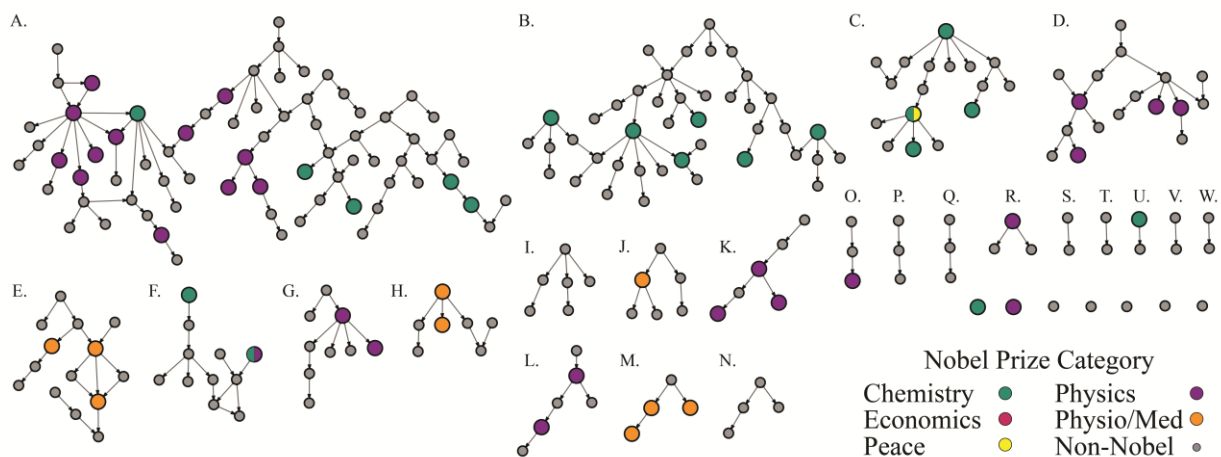
149 ***Identification of Nobel Laureate Communities.*** To identify highly successful scientific  
150 communities, the largest component of The Academic Tree Network, displayed in Fig. 1A, was  
151 filtered to include only individuals at or above the 99<sup>th</sup> percentile for the number of local Nobel  
152 laureate family members (99<sup>th</sup> percentile = 4) and the number of Nobel laureate descendants (99<sup>th</sup>  
153 percentile = 1), along with their first neighbors in the network. This produced one large  
154 subnetwork of 1276 individuals and 5 smaller subnetworks ranging in size from 3 to 68  
155 individuals (Fig. 1B). This network remained quite large, and in Fig. 2, first neighbors were  
156 removed, producing a more tractable set of 30 subnetworks for analysis, ranging from 1 to 73  
157 individuals. Nobel laureates in the surrounding academic family who contributed to the scores of  
158 these individuals are not pictured. Consequently, these subnetworks only display individuals at  
159 the center of the local academic family, making the scale of the two largest subnetworks  
160 remarkable. A list of individuals in this group, along with the number of family and Nobel  
161 laureate family on all measures, is available in Dataset S1. To explore the connectivity among  
162 these scientists, high resolution images of Figs. 1B and 2 are available in the supplement with  
163 scientist's names (see Figs. S3 and S4).

164

165



166  
 167 **Fig. 1.** The largest component of The Academic Tree network (A) filtered to include individuals  
 168 at the 99<sup>th</sup> percentile for number of Nobel laureate descendants and number of local Nobel family  
 169 along with their first neighbors (B). The individual names associated with each node in  
 170 subnetwork B are viewable in a high resolution pdf in the supplement (Fig. S3).  
 171



172  
 173 **Fig. 2.** The largest component of The Academic Tree network filtered to include only individuals  
 174 at the 99<sup>th</sup> percentile for number of Nobel laureate descendants and number of local Nobel  
 175 family. The individual names, the number of local Nobel family, and the number of Nobel  
 176 descendants associated with each node are viewable in a high resolution pdf in the supplement  
 177 (Fig. S4).

178           The largest subnetwork (Fig. 2A) can be segmented into two early communities by  
179 identifying the geographical location of the scientists. One community centered around J. J.  
180 Thomson (physics, 1906) and Ernest Rutherford (chemistry, 1908) at Cambridge University, and  
181 a second centered around notable scientists such as August Kundt, Wilhem Rontgen (physics,  
182 1901), Johannes Muller, and Hermann von Helmholtz, among others, working across multiple  
183 universities in Germany and Switzerland. The German/Swiss community extends to Herman  
184 Staudinger (chemistry, 1953) and Leopold Ruzicka (chemistry, 1939) at the right of the  
185 subnetwork. Justis Von Liebig, a German chemist considered the founder of organic chemistry  
186 <sup>11</sup>, has the greatest number of Nobel descendants in this group at 53. Eilhard Mitscherlich,  
187 Heinrich Magnus, and Johannes Muller follow with 27, 26, and 26 Nobel descendants,  
188 respectively. In the Cambridge community, William Hopkins and Edward Routh, well-known  
189 non-Nobel laureate mentors <sup>12</sup>, lead with 22 Nobel descendants. Their mentees/grandmentees, J.  
190 J. Thomson and Ernest Rutherford, both Nobel laureates, have 16 local Nobel laureate family  
191 members. David Shoenberg, a British physicist with 17 Nobel laureate family members connects  
192 these two major communities.

193           Interestingly, two additional communities in the largest subnetwork were established in  
194 the United States through mentors trained in Germany. Nobel laureate Isador Isaac Rabi  
195 (physics, 1944) with 8 Nobel descendants, 6 of which are mentees/grandmentees, serves as the  
196 center of one group at Columbia University. William GIAUQUE (chemistry, 1949) and Willard  
197 Libby (chemistry, 1960) are at the center of a second community at the University of California,  
198 Berkeley.

199           A significant portion of the second largest subnetwork (Fig. 2B), also contains  
200 individuals operating across universities in Germany, and once again, individuals trained in  
201 Germany began new communities at universities in Britain, through William Perkins, and  
202 universities in the Northeastern United States, through Ira Remsen. In this subnetwork, Friedrich  
203 Wohler, a German chemist, has the highest number of Nobel descendants at 39. Johannes  
204 Wislicenus, a German chemist, and William Perkin, an English chemist, have the highest number  
205 of local Nobel family at 13. Approximately half of Wislicenus's Nobel family are  
206 mentees/grandmentees.

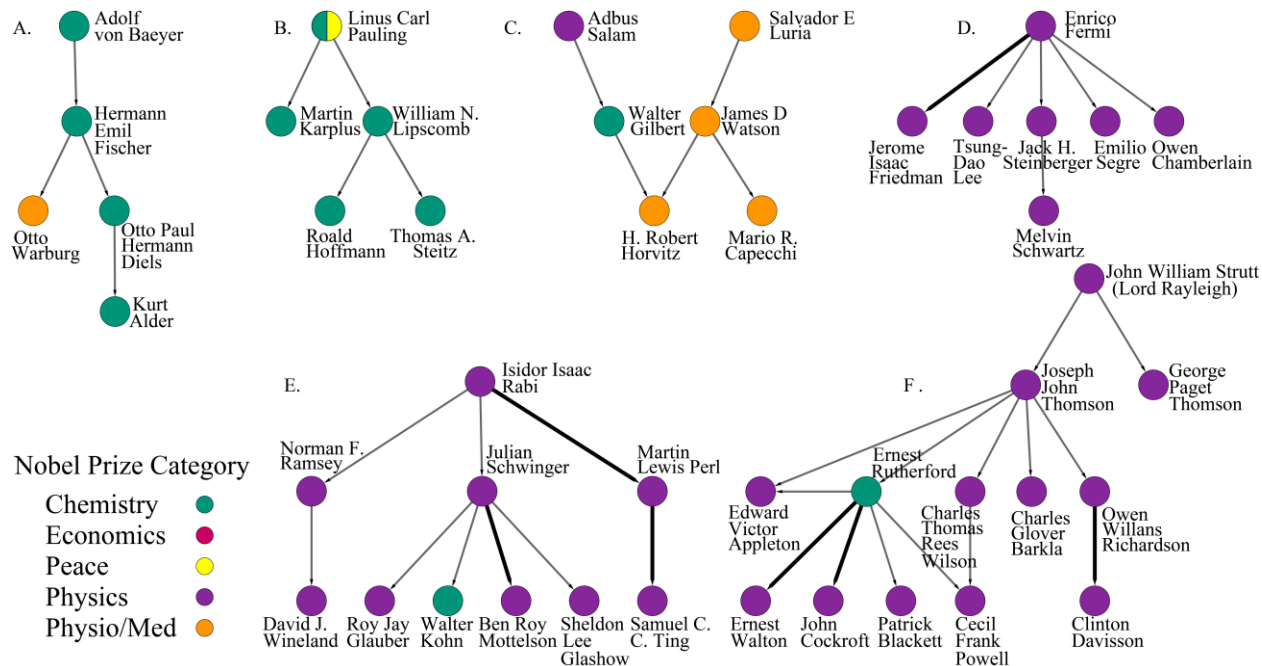
207           Of particular note in the smaller subnetworks is Enrico Fermi (physics, 1938; Fig. 2K)  
208 with 18 Nobel laureate family members and 6 Nobel laureate mentees/grandmentees. Fermi



209 trained and spent his early years as an academic in Italy during the early 20<sup>th</sup> century, but  
210 traveled to Germany to study with Max Born (physics, 1954) and Paul Ehrenfest<sup>11</sup>. Eventually,  
211 near the beginning of WWII, and on winning his Nobel Prize, Fermi moved to the United States  
212 and joined the Columbia University community centered on Isaac Rabi. In another subnetwork  
213 operating around the same time period (Fig. 2D), Max Born, Werner Heisenberg (physics, 1932),  
214 Hans Bethe (physics, 1967), and Robert Oppenheimer, among others, can be found. In this  
215 network, Arnold Sommerfeld, a non-Nobel laureate mentor to Heisenberg and Bethe, has 16  
216 Nobel family members and 11 Nobel descendants.

217 ***Nobel Laureate Subnetworks.*** On close inspection of Fig. 1B, small subnetworks can be  
218 identified that are comprised entirely of Nobel laureates. To explore this further, all non-Nobel  
219 laureates were removed from the large strongly connected network (Fig. 1A). Of the 402 Nobel  
220 laureates, 260 had no direct connection to another Nobel laureate. However, there were 142  
221 Nobel laureates in 55 subnetworks ranging in size from 2 to 10 individuals. Fig. 3 displays the  
222 six largest of the subnetworks. An investigation of relationships in these subnetworks identified  
223 seven additional connections recorded in Academic Tree after receiving the data used in the  
224 analysis. These are indicated by bold edges in the figure. Once again, the scientific communities  
225 at Cambridge and Columbia are identified as exceptional with 13 Nobel laureates connected over  
226 four generations at Cambridge (Fig. 3F) and 10 Nobel laureates connected over three generation  
227 at Columbia (Fig. 3E).

228  
229



230  
231

232 **Fig. 3.** The six largest subnetworks composed entirely of Nobel laureates. Bold edges indicate  
233 mentoring relationships recorded in The Academic Tree subsequent to the receipt of the dataset.

234  
235 **Heterogeneity of Local Academic Family.** In Fig. 1B, clusters of Nobel laureates who were  
236 awarded prizes in the same category can be seen with areas of greater diversity appearing where  
237 the clusters overlap. In an effort to characterize diversity in the network around Nobel Prize  
238 winners, the heterogeneity of Nobel Prize categories was measured within each individual's local  
239 family (see Method).

240 The vast majority of individuals in the network had 1 or fewer Nobel laureates in their  
241 local family. Therefore, the analysis was restricted to individuals with 2 or more Nobel laureates  
242 in the family and where some opportunity for diversity existed (205 of 402 Nobel laureates;  
243 3,647 of 57,429 non-Nobel laureates). As the number of Nobel laureates increased, heterogeneity  
244 scores also increased ( $r = 0.25, p < .0001$ ). The distribution of scores for Nobel Prize winners in  
245 individual categories was positively skewed (chemistry: 0 to 0.37, Mdn = 0; physics 0 to 0.42,  
246 Mdn = 0.17; Physio/Med: 0 to 0.44, Mdn = 0) and reflective of generally homogeneous clusters  
247 of Nobel laureates with some diversity where clusters overlap. There were 6 Nobel laureates (3  
248 physics, 2 physio/med, 1 physics/chemistry) and 39 non-Nobel laureates with scores at or above  
249 the 99<sup>th</sup> percentile (0.37). Archibald Hill, a Nobel laureate in physiology and medicine (1922;

250 Fig. 2E), scored the highest family heterogeneity (0.44) with 4 Nobel laureate family in  
251 physiology/medicine, 4 in physics, and 1 in chemistry.

252 These findings prompted us to ask whether there was any difference in family diversity  
253 for Nobel laureates and non-Nobel laureates, and the data were fit to a quasibinomial model with  
254 Nobel status and number of local Nobel laureate family as predictors. In the final analysis,  
255 heterogeneity scores were not predicted by Nobel status (adjusted  $p = 0.347$ ) after controlling for  
256 the number of Nobel laureates in the local family (adjusted  $p = 0.177$ ).

## 257 **Discussion**

258 Remarkable connectedness among Nobel laureates is found through generations of mentoring  
259 relationships in The Academic Tree network. Nobel laureates have more Nobel laureate  
260 ancestors, more local and global descendants, and more local academic family members than do  
261 non-Nobel laureates. A variety of explanations for this connectedness exist. Nobel laureates  
262 undoubtedly possess superior knowledge and skill that individuals in the local academic family,  
263 and the greater community, may acquire through a variety of means. Other factors related to the  
264 availability of resources and the attraction of talent are no doubt significant contributors to the  
265 connectedness of this group. These additional factors are difficult to separate from the transfer of  
266 knowledge through mentoring but are an integral part of any successful scientific community and  
267 should be valued as such.

268 Several areas of the network, representing mentoring relationships in historical scientific  
269 communities, were identified with high concentrations of Nobel laureates. In some locations,  
270 direct connections between Nobel laureates occurred over three and four generations. When  
271 exploring biographical and historical accounts of these communities, it was apparent that much  
272 greater interconnectedness existed among scientific communities than is reflected by doctoral  
273 mentor-mentee relationships. A high degree of interaction occurred throughout these  
274 communities. For, example, after completing a dissertation at Columbia University, Isador Isaac  
275 Rabi (physics, 1944) spent over a year in Europe where he encountered some of the greatest  
276 minds in science, many of whom went on to become Nobel laureates<sup>13</sup>. Rabi returned to  
277 Columbia to eventually lead the Physics Department and to become a central figure in one of the  
278 two most successful communities identified in this analysis<sup>14-16</sup>. In future work, extending the  
279 analysis to include a greater variety of mentoring relationships would better capture the true  
280 interconnectivity among scientists.

281 It is significant that many of the successful communities identified by this network  
282 analysis existed at a time when travel and communication were much more difficult than they are  
283 today. Ernest Rutherford (chemistry, 1908) traveled from New Zealand to attend Cambridge as  
284 one of the first students admitted from outside the university<sup>17,18</sup>. This occurred in the latter half  
285 of the 19th century prior to the invention of the airplane and intercontinental telephone service.  
286 At this point in history, physical proximity was critical to the transmission of ideas and expertise.  
287 In modern science, however, virtual meetings, video lectures, online courses, and online  
288 databases (e.g., PubMed<sup>19</sup>, Google Scholar<sup>20</sup>) provide remarkably easy access to current,  
289 innovative ideas in science. It seems likely that the mentoring patterns among scientists are being  
290 radically altered by greater accessibility to information and each other. Still, for many scientists,  
291 it is difficult to imagine that virtual proximity could ever be a satisfying replacement for the day-  
292 to-day personal interaction found in a positive mentoring relationship.

293 Biographical and historical accounts provided the sole access to more detailed  
294 information about the communities identified in this study. Warwicke (2003) offers a particularly  
295 valuable and compelling account of the scientific community identified at Cambridge in the  
296 latter part of the 19<sup>th</sup> century<sup>12</sup>. However, modern scientific communities could be studied if the  
297 types of data required to identify them were available. Although the number of Nobel laureates  
298 within an academic community serves as a legitimate measure of success, especially when the  
299 research focus is restricted to high-level science, much more could be accomplished if a variety  
300 of other performance measures were readily available and reliably accurate. Information  
301 regarding publications, impact factors, citations, funding sources, and other awards, would allow  
302 for a more sensitive evaluation of success within a community. This could be achieved by a  
303 committed effort in the scientific community to collect performance measures from all  
304 individuals and universities and to make them available in an open-source database, something  
305 The Academic Tree is currently attempting to accomplish.

306 Several factors would be critical to the success of this endeavor. Primarily, a  
307 comprehensive list of all researchers' publications would need to be available in a centralized,  
308 open-source database. Currently, no one source is guaranteed to have a complete set of  
309 publications for an individual author<sup>21</sup>, and publication information must be obtained from  
310 multiple sources, such as Web of Science<sup>22</sup>, Scopus<sup>23</sup>, and PubMed<sup>19</sup>. Furthermore, some of  
311 these sources are proprietary and require a fee for use. Google Scholar<sup>20</sup> has access to several

312 proprietary sources through licensing agreements but does not allow automated searches of its  
313 website, something that is a requirement when conducting an analysis of “big data”. As an  
314 example, the largest component of the Academic Tree Network analyzed in this study contained  
315 57,831 individuals, making manual search costly in terms of time.

316 Making a wide variety of performance measures accessible would also increase the value  
317 of a database for evaluating scientific success. Number of publications, a measure of  
318 productivity, is not a sufficient measure of success. Rather, number of citations, considered a  
319 measure of quality, is often factored alongside number of publications in calculations such as the  
320 *h*-index<sup>24</sup>. Number of citations is not consistently available in the sources mentioned earlier, and  
321 it is not clear how often this information is updated. Along these lines, additional quality  
322 measures, such as a journal’s impact factor at the time of an article’s publication, author funding,  
323 and additional awards, would be useful in developing new algorithms for measuring the quality  
324 of research and the impact of an individual’s and a community’s contribution to science.

325 Another critical element in developing an effective database involves the assignment of  
326 unique identifiers for scientists. This is especially important when dealing with crowd-sourced  
327 data. On one hand, crowd-sourcing allows for the collection of data that would be difficult or  
328 impossible to obtain otherwise. On the other hand, a quick glance at the Academic Tree dataset  
329 makes it clear that ensuring consistency, completeness, and accuracy of the data requires a rigid  
330 collection protocol. For example, in the Academic Tree dataset individuals may use all uppercase  
331 letters or put a nickname in parentheses, all of which create problems for automated analysis. A  
332 unique numerical identifier would allow for much less variation. This problem is clear to many  
333 in the scientific community, and it is being pursued by projects such as ORCID<sup>25</sup>. However, to  
334 facilitate performance analyses, its use must be required, especially in the authorship section of  
335 papers, so that the publications for authors with the same name can be easily distinguished in an  
336 automated fashion.

### 337 **Conclusion**

338 Using methods of network analysis, Nobel laureates were identified as a highly connected group  
339 in The Academic Tree network. Several successful mentoring communities could be identified  
340 using the number of Nobel laureates as a measure of scientific success. A variety of performance  
341 measures exist that would increase the sensitivity of these types of analyses and would allow for  
342 the exploration of a greater variety of questions if the measures were collected and made

343 available in a single database. This could provide valuable information regarding individual,  
344 institutional, and national factors associated with success in modern science and lead to a greater  
345 understanding of best practices. The rewards in such an endeavor would be large, especially in  
346 the current climate where there is an increased focus on effective collaboration and teamwork.

## 347 **Methods**

348 **Network Collection.** The Academic Tree Network of mentor-mentee relationships was obtained  
349 for analysis from Academic Tree on November 2, 2015<sup>7,8</sup>. Academic Tree is a web-based  
350 database of academic mentor-mentee relationships that uses a crowd-sourcing method for the  
351 collection of information. Individuals can voluntarily provide information regarding academic  
352 relationships through the Academic Tree website. Academic Tree can be decomposed into an  
353 interconnected set of 68 domain specific networks, and it is possible for an individual to be listed  
354 in more than one domain. For example, a Cell Biology Tree exists for individuals working in cell  
355 biology, and a Genetics Tree exists for those working in the field of genetics. An individual  
356 working in both areas can identify themselves as belonging to both trees.

357 As can be seen in Table S2, the Academic Tree database holds several types of  
358 information, including an individual's specific research area, major research area (i.e., one or  
359 more of the domain specific trees), and five possible academic relationships between individuals  
360 in the network, including doctoral student-advisor relationships. The database was received from  
361 Academic Tree in SQL format which included an edge file and a node file. There were 114,949  
362 entries in the node file and 260,201 entries in the edge file.

363 **Network Filtering.** Several steps were taken to enhance and filter the network prior to analysis.  
364 In the original files, there were 484 individuals listed as Nobel laureates. However, 47 additional  
365 Nobel Prize winners could be identified in the file and were labeled as such. The Nobel Prize  
366 category and year were added for all Nobel laureates. All information regarding Nobel laureates  
367 was obtained from the Nobel Foundation (1). This network was first filtered to include only  
368 relationships between doctoral students and advisors (see Table S3, Filter 1). Next, the network  
369 was filtered to included only individuals listed in at least one science tree (Table S3, Filter 2). As  
370 a result, the majority of Nobel laureates winning prizes for peace, literature, and economics were  
371 removed. In the last step, the strongly connected components within the larger network data set  
372 were identified using network analysis tools in Gephi<sup>26</sup>. Table S4 lists the number of strongly  
373 connected components of different sizes along with the number of nodes and the number of

374 Nobel laureates associated with each component size. The largest strongly connected component  
375 of 57,831 nodes was significantly larger than any of the other components and held the vast  
376 majority of Nobel laureates (402 of 472). In fact, in Table S5, this was a significant majority of  
377 all Nobel laureates in physics (58.2 %), physiology (65.2 %), and chemistry (86 %), justifying  
378 the use of this subnetwork in the subsequent analysis. All Nobel laureates in the largest strongly  
379 connected component of the network received prizes in chemistry, physics, and physiology or  
380 medicine with one exception, Herbert Simon, a highly interdisciplinary scientist who won the  
381 Nobel Prize in economics. There were no prize winners in literature, and only one Peace Prize  
382 winner, Linus Pauling, who was also awarded the Nobel Prize in chemistry. All further network  
383 visualization and filtering was done in Cytoscape<sup>27</sup>. The Cytoscape filtering tool was used to  
384 identify and visualize the subnetworks displayed in Fig. 1B, Fig. 2, and Fig. 3.

385 **Data Analysis.** A breadth first search algorithm was used to calculate the number of family  
386 members and the number of Nobel laureate family members for each individual. This was  
387 instantiated in a custom C++ program which takes a directed acyclic graph as a node and edge  
388 list. The node list contains node id, Nobel status, and Nobel Prize category. The edge list  
389 contains source and target nodes. The direction in the network (forward, backward, or both) and  
390 the number of academic generations (i.e., steps in the network) to be calculated is specified as  
391 input to the program. Number of ancestors/Nobel ancestors was calculated as 31 steps (i.e., the  
392 diameter of the network) backward while number of descendants/Nobel descendants was  
393 calculated as 31 steps forward. The number of mentees/grandmentees and Nobel  
394 mentees/grandmentees was calculated as two step forward in the network, and the number of  
395 local family/Nobel local family was calculated as three steps forward and backward in the  
396 network. While calculating number of Nobel laureate family the program also tracks number of  
397 Nobel Prizes in each category.

398 **Heterogeneity Computation.** A measure of heterogeneity was used to calculate the diversity of  
399 Nobel Prizes awarded within three steps of each individual in the network (Equation 1) where  $L$   
400 is the number of Nobel laureates within a specified distance in the network,  $N$  is the number of  
401 Prize categories (5 in this case),  $n_i$  is the number of Nobel laureates in a specific prize category  
402 within a specified distance. The number of Nobel laureates in an individual's local family was an  
403 essential factor in the equation and meant that the scores could not be compared across  
404 individuals with different numbers of Nobel family members. Therefore, the scores were

405 normalized to fall between 0 and 1 with 1 representing the greatest possible diversity for a given  
406 number of Nobel laureates.

407

$$408 \quad H = \frac{1}{L} * \log_N \frac{L!}{\prod_{i=1}^N n_i!} \quad (1)$$

409

410 **Generation of Random Networks.** Hypothesis testing with network data is problematic in that  
411 the assumption of independent observations required for many statistical methods is violated,  
412 resulting in standard errors that are computed incorrectly<sup>28,29</sup>. To handle this, the significance  
413 levels in our analyses were adjusted by creating a distribution of expected test statistics, derived  
414 from random samples, for comparison with an observed test statistic<sup>29</sup>. This involved permuting  
415 values for the predictor variable (Nobel status) with respect to an outcome variable (number of  
416 Nobel laureate family members) for one thousand samples, performing the statistical analysis on  
417 each of the random samples, and then counting the number of test statistics on the permuted data  
418 that were greater than or equal to the observed statistic. This number was then divided by the  
419 number of random samples to produce an adjusted  $p$  value. For example, if three random test  
420 statistics of 1000 permuted samples are greater than or equal to the observed test statistics, the  $p$   
421 value would be adjusted to 0.003.

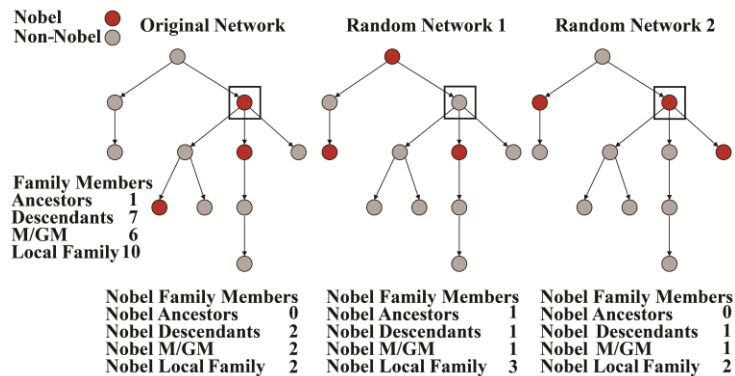
422 To accomplish this, the C++ program described earlier had options available for  
423 generating 1,000 networks with Nobel status randomly assigned to nodes across the network in  
424 the same proportion as the true data, each time recomputing outcome measures for each node. As  
425 can be seen in Fig. 4, this produced alternate networks with equivalent topology (i.e., the same  
426 number of family members and academic structure for each node) but randomly distributed  
427 Nobel laureates and thus, random outcomes.

428

429

430





431  
432  
433  
434  
435  
436  
437

**Fig. 4.** Number of family members and number of Nobel laureate family members computed for an individual network node, highlighted in the box, for a directed network (left) and two networks in which Nobel status is randomly permuted (right).

438 **Acknowledgments**

439 The authors would like to thank Stephen David for providing The Academic Tree database. This  
440 research would not have been possible without his generous support. The authors would also like  
441 to thank anonymous reviewers and members of the KBRIN Bioinformatics Core for helpful  
442 insight and feedback. All data and scripts used in constructing this analysis is available at  
443 <http://bioinformatics.louisville.edu/Nobel/>. ECR conceived the idea of the project and supervised  
444 all aspects of the project. JHC implemented the computational aspects of the project, updated  
445 and corrected the mentor network, and performed all network analyses, including biographical  
446 and historical reviews. JRP provided insight into the analyses from a social science perspective  
447 and helped with preparation of the manuscript. YZ implemented the web interface for  
448 interactively exploring the Nobel networks. All authors contributed to the writing of the  
449 manuscript. Support for JHC and ECR provided by National Institutes of Health (NIH) grant  
450 P20GM103436 (Nigel Cooper, PI). The contents of this work are solely the responsibility of the  
451 authors and do not represent the official views of the NIH or the National Institute for General  
452 Medical Sciences (NIGMS).

453

## 454 **References**

- 455 1 *Nobel Prize Facts*, [www.nobelprize.org/nobel\\_prizes/facts/](http://www.nobelprize.org/nobel_prizes/facts/), (2016).
- 456 2 Allen, T. D. Protégé selection by mentors: Contributing individual and organizational factors.  
457 *Journal of Vocational Behavior* **65**, 469-483 (2004).
- 458 3 Allen, T. D., Poteet, M. L. & Russell, J. E. A. Protégé selection by mentors: what makes the  
459 difference? *Journal of Organizational Behavior* **21**, 271-282 (2000).
- 460 4 Green, S. G. & Bauer, T. N. Supervisory mentoring by advisers: Relationships with doctoral  
461 student potential, productivity, and commitment. *Personnel Psychology* **48**, 537-562 (1995).
- 462 5 Malmgren, R. D., Ottino, J. M. & Nunes Amaral, L. A. The role of mentorship in protege  
463 performance. *Nature* **465**, 622-626 (2010).
- 464 6 Paglis, L. L., Green, S. G. & Bauer, T. N. Does adviser mentoring add value? A longitudinal  
465 study of mentoring and doctoral student outcomes. *Research in Higher Education* **47**, 451-  
466 476 (2006).
- 467 7 David, S. *The Academic Family Tree*, [www.academictree.org](http://www.academictree.org), (2016).
- 468 8 David, S. V. & Hayden, B. Y. Neurotree: a collaborative, graphical database of the academic  
469 genealogy of neuroscience. *PloS one* **7**, e46608 (2012).
- 470 9 Hothorn, T. & Everitt, B. S. *A handbook of statistical analyses using R*. (CRC press, 2014).
- 471 10 *Institute For Digital Research and Education*, <http://www.ats.ucla.edu/stat/r/dae/zinbreg.htm>,  
472 (2016).
- 473 11 *Wikipedia: The Free Encyclopedia*, [www.wikipedia.org](http://www.wikipedia.org), (2016).
- 474 12 Warwick, A. *Masters of theory: Cambridge and the rise of mathematical physics*.  
475 (University of Chicago Press, 2003).
- 476 13 Rabi, I. I. & Code, R. F. Stories from the early days of quantum mechanics. *Physics Today*  
477 **59**, 36 (2006).
- 478 14 Cole, J. R. *The great American university: Its rise to preeminence, its indispensable national*  
479 *role, why it must be protected*. (PublicAffairs, 2010).
- 480 15 Cropper, W. H. *Great physicists: the life and times of leading physicists from Galileo to*  
481 *Hawking*. (Oxford University Press, 2001).
- 482 16 Rigden, J. S. *Rabi, scientist and citizen*. (Harvard University Press, 2000).
- 483 17 Reeves, R. *A force of nature: The frontier genius of Ernest Rutherford*. (WW Norton &  
484 Company, 2008).

- 485 18 Thomson, J.-J. Recollections and reflections. *Ciel et Terre* **53**, 315 (1937).
- 486 19 *PubMed*, <http://www.ncbi.nlm.nih.gov/pubmed>, (2016).
- 487 20 *Google Scholar*, <https://scholar.google.com/>, (2016).
- 488 21 Falagas, M. E., Pitsouni, E. I., Malietzis, G. A. & Pappas, G. Comparison of PubMed,  
489 Scopus, web of science, and Google scholar: strengths and weaknesses. *The FASEB journal*  
490 **22**, 338-342 (2008).
- 491 22 *Web of Science*, <http://ipsience.thomsonreuters.com/product/web-of-science/>, (2016).
- 492 23 *Scopus*, <https://www.elsevier.com/solutions/scopus>, (2016).
- 493 24 Hirsch, J. E. An index to quantify an individual's scientific research output. *Proceedings of*  
494 *the National academy of Sciences of the United States of America* **102**, 16569-16572 (2005).
- 495 25 *ORCID: Open Researcher and Contributor ID*, <https://orchid.org/>, (2016).
- 496 26 Bastian, M., Heymann, S. & Jacomy, M. Gephi: an open source software for exploring and  
497 manipulating networks. *ICWSM* **8**, 361-362 (2009).
- 498 27 Kohl, M., Wiese, S. & Warscheid, B. Cytoscape: software for visualization and analysis of  
499 biological networks. *Data Mining in Proteomics: From Standards to Applications*, 291-303  
500 (2011).
- 501 28 Hanneman, R. A. & Riddle, M. *Introduction to social network methods*. (University of  
502 California Riverside, 2005).
- 503 29 Borgatti, S. P., Everett, M. G. & Johnson, J. C. *Analyzing social networks*. (SAGE  
504 Publications Limited, 2013).

505  
506  
507