

## Expertise versus Bias in Evaluation: Evidence from the NIH<sup>†</sup>

By DANIELLE LI\*

*Evaluators with expertise in a particular field may have an informational advantage in separating good projects from bad. At the same time, they may also have personal preferences that impact their objectivity. This paper examines these issues in the context of peer review at the US National Institutes of Health. I show that evaluators are both better informed and more biased about the quality of projects in their own area. On net, the benefits of expertise weakly dominate the costs of bias. As such, policies designed to limit bias by seeking impartial evaluators may reduce the quality of funding decisions. (JEL D82, H51, I10, I23, O38)*

A key debate in the economics of innovation focuses on what mechanisms are most effective for encouraging the development of new ideas and products: while patents may distort access to new knowledge ex post, a concern with research grants and other R&D subsidies is that the public sector may make poor decisions about which projects to fund ex ante.

In the United States, the vast majority of public funding for biomedical research is allocated by the National Institutes of Health (NIH), through a system of peer review in which applications are evaluated by committees of scientists working on similar issues. The collective opinion of these scientists is responsible for consolidating thousands of investigator-initiated submissions into a publicly funded research agenda.

But how much should we trust their advice? While reviewers may have valuable information about the potential of projects in their research areas, advice in this setting may also be distorted precisely because reviewers have made so many investments in acquiring their domain expertise. For example, in a guide aimed at scientists describing the NIH grant review process, one reviewer highlights his preference for work related to his own: “If I’m sitting in an NIH study section, and I

\*Harvard Business School, 211 Rock Center, Boston, MA 02163 (e-mail: [dli@hbs.edu](mailto:dli@hbs.edu)). I am very grateful to Pierre Azoulay, Michael Greenstone, and especially David Autor for detailed feedback on this project. I also thank Jason Abaluck, Leila Agha, Josh Angrist, Manuel Bagues, David Berger, David Chan, Brigham Frandsen, Alex Frankel, Richard Freeman, Bob Gibbons, Nathan Hendren, Ben Jones, Bill Kerr, Jiro Kondo, Josh Lerner, Niko Matouschek, Xiao Yu May Wang, Ziad Obermeyer, Marco Ottaviani, Dimitris Papanikolaou, Amanda Pallais, Chris Palmer, Michael Powell, Amit Seru, Heidi Williams, and numerous seminar participants for helpful comments and suggestions. I am grateful to George Chacko, Raviv Murciano-Goroff, Joshua Reyes, and James Vines for assistance with data. All errors are my own.

<sup>†</sup>Go to <https://doi.org/10.1257/app.20150421> to visit the article page for additional materials and author disclosure statement or to comment in the online discussion forum.

*believe the real area of current interest in the field is neurotoxicology [the reviewer's own speciality], I'm thinking if you're not doing neurotoxicology, you're not doing interesting science."*<sup>1</sup> Alternatively, reviewers may be biased against applicants in their own area if they perceive them to be competitors.

This paper examines the impact of intellectual proximity between reviewers and applicants (hereafter "proximity" or "relatedness") on the quality of funding decisions. I provide evidence that reviewers are better informed about the quality of related candidates but also biased in their favor. Overall, the benefits of expertise appear to outweigh the costs of bias. To show this, I assemble a new, comprehensive dataset linking almost 100,000 NIH grant applications to the committees in which they were evaluated.

My analysis requires two key ingredients: (i) a source of exogenous variation in the intellectual proximity between grant applicants and the more influential members of their review committees and (ii) a measure of quality for grant applications, including that of unfunded applications. Given these, the intuition underlying my empirical work is as follows: if reviewers are only biased and not more informed, then related applicants should receive better (or worse) evaluations regardless of their quality. If related reviewers also have better information, then the effect of working in the same area as a more influential reviewer should differ for high- and low-quality applicants. Strong applicants should benefit from being evaluated by influential reviewers who can more accurately assess their quality, but weak applicants should be hurt for the same reason. In this case, the impact of proximity should be increasing in the quality of applications. I now provide more detail about my proximity and quality measures in turn.

I begin with a baseline measure of the intellectual proximity of individual applicants and reviewers: whether a reviewer has cited an applicant's work in the five years prior to the committee meeting. This captures whether the applicant's work has been of use to the reviewer, but is likely to be correlated with quality because better applicants are more likely to be cited. To identify exogenous variation proximity between candidates and review committees, I take advantage of the distinction between "permanent" and "temporary" members in NIH review committees.<sup>2</sup> Permanent members play a greater role in the grant evaluation process and have more influence on committee scores. I define intellectual proximity to the review committee as the number of *permanent* reviewers that have cited an applicant's work—controlling for the total number of such reviewers, both permanent and temporary. This strategy identifies the causal impact of being related to a more influential set of reviewers, under the assumption that the quality of an applicant is not correlated with the composition of reviewers who cite her work.

The next part of my analysis considers the role of bias and expertise. To do so, I require information on application quality. The primary challenge in measuring application quality is doing so for unfunded applications: it is natural, after all, to think that the research described in unfunded applications does not get produced

<sup>1</sup> See [http://www.clemson.edu/caah/research/images/What\\_Do\\_Grant\\_Reviewers\\_Really\\_Want\\_Anyway.pdf](http://www.clemson.edu/caah/research/images/What_Do_Grant_Reviewers_Really_Want_Anyway.pdf).

<sup>2</sup> "Permanent" members are not actually permanent; they serve four-year terms. See Sections I and IIIA for a discussion of permanent versus temporary reviewers.

and thus its quality cannot be observed. At the NIH, however, this is not the case. Standards for preliminary results for large research grants are so high that researchers often submit applications based on nearly completed research. As a result, it is common to publish the work proposed in an application even if the application itself goes unfunded. To find these related publications, I use a text matching approach that compares grant application titles with the titles and abstracts of publications to find research by the same applicant on the same topic as the grant. I further restrict my analysis of application quality to articles published soon enough after grant review to not be directly affected by any grant funds. For consistency, I use this same approach to measure the quality of funded applications as well.

I present three key findings. First, related applicants receive higher scores and are more likely to be funded. Each additional permanent reviewer in an applicant's area, holding constant total related reviewers, increases that applicant's chances of being funded by 2.2 percent. While this may seem like a small effect, it is substantial when viewed relative to reviewers' sensitivity to application quality: it is the same increase in funding probability that we would expect from a one-quarter standard deviation increase in the quality of the application itself, as measured by citations to research that it produces. This large effect suggests that when quality is difficult to assess, reviewer opinions play a comparably large role in funding decisions.

Second, I show that these findings are most consistent with reviewers both having better information about the quality of related applications and being biased in their favor. If reviewers were biased but not better informed, related applications should be more likely to be funded and receive higher scores regardless of their quality. Instead, I find that higher quality applicants benefit more from being evaluated by related reviewers: the impact of proximity on funding likelihood and scores are both increasing in quality. At the same time, if reviewers were more informed but unbiased, I would expect the impact of proximity to be negative for low quality applications. For all but possibly the lowest decile of applications, this is not what I find. Rather, I find that low-quality applications receive no benefit from being evaluated by related reviewers. This pattern strongly suggests that worse information for this group cancels out the positive overall effects of bias.

A potential concern with my empirical strategy is the fact that permanent and temporary reviewers differ significantly in terms of their past publications and citations, and applicants cited by more permanent members, even conditional on total relatedness, have significantly more past citations. If stronger applicants are more likely to be cited by permanent reviewers, then this may be an alternative explanation for my finding that related applicants are more likely to be funded. I provide several different pieces of evidence that this does not drive my results.

First, it is not the case that permanent reviewers are more qualified than temporary reviewers: they tend to have more citations but fewer publications. Second, I show that, conditional on the total number of reviewers an applicant has been cited by, there is no correlation between the number of permanent reviewers that cite an applicant and my text-matched measure of application quality. Further, the impact of relatedness that I estimate in my main tables does not change when I include detailed controls for applicant characteristics and publication histories. Finally, I provide two additional complementary sets of analysis. The first uses reviewer fixed

effects to show that applicants are more likely to be funded when the reviewer that has cited them is serving as a permanent reviewer, compared to when that reviewer is serving as a temporary reviewer. I also show that my results do not rely on the distinction between permanent and temporary reviewers by using applicant fixed effects to compare outcomes for the same applicant across meetings in which she is cited by different numbers of reviewers. This alternative specification identifies the effect of being related to an *additional* reviewer under the assumption that the time-variant unobserved quality of an application is not correlated with proximity.

Finally, I provide suggestive evidence that the gains associated with reviewer expertise dominate the losses associated with bias: the average quality of applications funded by committees in which a greater share of applicants are related to reviewers tends to be higher than meetings of the same committee in which fewer applicants are related. This suggests that enacting a policy that restricts close intellectual ties may reduce the quality of the NIH-supported research portfolio, as measured by future citations.

Of course, such a conclusion regarding intellectual ties in science may not hold in other settings. For instance, biases associated with financial conflicts of interest may well outweigh any informational advantages that those decision makers may have. Nonetheless, the results in this paper have implications for how organizations treat conflicts of interest. In many settings, personal preferences develop alongside expertise, as a result of individuals self-selecting and making investments into a particular domain. These biases are particularly challenging to address: in contrast with race or gender discrimination, eliminating bias stemming from intellectual ties can directly degrade the quality of information that decision makers have access to. This paper demonstrates that conflict of interest policies necessarily entail efficiency trade-offs.

The question of how organizations should use information from potentially conflicted experts has also been of long-standing theoretical interest (Crawford and Sobel 1982; Li, Rosen, and Suen 2001; Garfagnini, Ottaviani, and Sørensen 2014), but has remained relatively understudied empirically. Emerging work shows that these issues are relevant in many empirical settings ranging from financial regulation to judicial discretion to academic promotion and publication.<sup>3</sup> In these and other settings, it is often challenging to attribute differences in the treatment of connected individuals to either better information or bias because it is difficult to observe the counterfactual quality of decisions that are not made. This paper contributes by studying these issues in the context of public investments in R&D, a setting that is both independently important, and in which various empirical challenges can be more readily overcome.

Finally, there is currently little empirical evidence on how—and how successfully—governments make research investments, and existing studies in this area find mixed results.<sup>4</sup> This paper demonstrates the value of expert advice in this setting.

<sup>3</sup>See, for instance, Agorwal et al. (2014); Hansen, McMahon, and Rivera (2014); Kondo (2006); Fisman, Paravisini, and Vig (2012); Zinovyeva and Bagues (2015); Blanes i Vidal, Draca, and Fons-Rosen (2012); Brogaard, Engleberg, and Parsons (2011); and Laband and Piette (1994).

<sup>4</sup>See Acemoglu (2009), Kremer and Williams (2010), Griliches (1991), and Cockburn and Henderson (2000) for surveys. Li and Agha (2015) document a positive correlation between scores and outcomes, but Boudreau et al. (2016) and Azoulay, Graff-Zivin, and Manso (2011) raise concerns about the ability to support recognize and foster

## I. Context

### A. Grant Funding at the NIH

The NIH plays an outsized role in supporting biomedical research. Over 80 percent of basic life science laboratories in the United States receive NIH funding, and half of all FDA approved drugs, and over two-thirds of FDA priority review drugs, explicitly cite NIH-funded research (Sampat and Lichtenberg 2011). The decision of what grants to support is made by thousands of scientists who act as peer reviewers for the NIH. Each year, they collectively read approximately 20,000 grant applications and allocate over \$20 billion in federal grant funding. During this process, more than 80 percent of applicants are rejected even though, for the vast majority of biomedical researchers, winning and renewing NIH grants is crucial for becoming an independent investigator, maintaining a lab, earning tenure, and paying salaries (Stephan 2012, Jones 2010).

The largest and most established of these grant mechanisms is the R01, a project-based, renewable research grant that constitutes half of all NIH grant spending and is the primary funding source for most academic biomedical labs in the United States. There are currently 27,000 outstanding awards, with 4,000 new projects approved each year. The average size of each award is \$1.7 million spread over three to five years.

Because R01s entail such large investments, the NIH favors projects that have already demonstrated a substantial likelihood of success. As evidence of how high this bar is, the NIH provides a separate grant mechanism, the R21, for establishing the preliminary results needed for a successful R01 application. The fact that R01 applications are typically based on research that is already very advanced makes it possible to measure the quality of unfunded grants, which is a key part of my empirical strategy.<sup>5</sup> See Section IIB for a detailed discussion.

To apply for an R01, the primary investigator submits an application, which is then assigned to a review committee (called a “study section”) for scoring and to an Institute or Center (IC) for funding. The bulk of these applications are reviewed in one of about 180 “chartered” study sections, which are standing review committees organized around a particular theme, for instance, “Cellular Signaling and Regulatory Systems” or “Clinical Neuroplasticity and Neurotransmitters.”<sup>6</sup> These committees meet three times a year in accordance with NIH’s funding cycles and, during each meeting, review between 40 to 80 applications. My analysis focuses on these committees.

---

novel research. Hegde (2009) considers congressional appropriations for NIH funding. Jacobs and Lefgren (2011) find few effects of individual NIH grants output related to marginally unfunded applicants.

<sup>5</sup>This emphasis on preliminary results was one point of critique that the NIH peer review reform of 2006 was designed to address; under the new system, the preliminary results section has been eliminated to discourage this practice. My data come from before the reform but, anecdotally, it is still the norm to apply for R01s. For a satirical take from 2011, see <http://www.phdcomics.com/comics/archive.php?comicid=1431>.

<sup>6</sup>The NIH restructured chartered study sections during my sample period and my data include observations from 250 distinct chartered study sections. These changes do not affect my estimation because I use within-meeting variation only.

Study sections are typically composed of 15 to 30 “permanent” members who serve four-year terms and 10 to 20 “temporary” reviewers who are called in as needed. Within a study section, an application is typically assigned up to three reviewers who provide an initial assessment of its merit. Permanent members are responsible for performing initial assessments on eight to ten applications per meeting, compared to only one to three for temporary members. The division of committees into permanent and temporary members plays an important role in my identification strategy: permanent reviewers have more influence over the scoring process, but are otherwise similar to temporary members in terms of their scientific credentials. In Section IIIA, I discuss why this might be the case and provide empirical evidence.

The process of assigning applications to study sections and reviewers is nonrandom. In practice, applicants are usually aware of the identities of most permanent study section members, suggest a preferred study section, and usually get their first choice (subject to the constraint that, for most applicants, there are only one or two study sections that are scientifically appropriate). Study section officers, meanwhile, assign applications to initial reviewers on the basis of intellectual fit. I will discuss the implications of this nonrandom selection on my identification strategy in Section IIIA.

Once an application has been assigned to a study section, it is assigned to three initial reviewers who read and score the application on the basis of five review criteria: *Significance* (does the proposed research address an important problem and would it constitute an advance over current knowledge?), *Innovation* (are either the concepts, aims, or methods novel?), *Approach* (is the research feasible and well thought out?), *Investigator* (is the applicant well-qualified?), and *Environment* (can the applicant’s institution support the proposed work?). Based on these scores, weak applications (about one-third to one-half) are “triaged” or “unscored,” meaning that they are rejected without further discussion. The remaining applications are then discussed in the full study section meeting. During these deliberations, an application’s initial reviewers first present their opinions, and then all reviewers discuss the application according to the same five review criteria. Following these discussions, all study section members anonymously vote on the application, assigning it a “priority score,” which, during my sample period, ranged from 1.0 for the best application to 5.0 for the worst, in increments of 0.1. The final score is the average of all member scores. This priority score is then converted into a percentile from 1 to 99.<sup>7</sup> In my data, I observe an application’s final score (records of scores by individual reviewers and initial scores are destroyed after the meeting).

Once a study section has scored an application, the institute to which it was assigned determines funding. Given the score, this determination is largely mechanical: an IC lines up all applications it is assigned and funds them in order of score until its budget has been exhausted. When doing this, the IC only considers the score: NIH will choose to fund one large grant instead of two or three smaller grants

<sup>7</sup> At the NIH, a grant’s percentile score represents the percentage of applications from the same study section and reviewed in the same year that received a better priority score. According to this system, a lower score is better, but, for ease of exposition and intuition, this paper reports inverted percentiles (100 minus the official NIH percentile, e.g., the percent of applications that are *worse*), so that higher percentiles are better.



as long as the larger grant has a better score, even if it is only marginally better. The worst percentile score that is funded is known as that IC's payline for the year. In very few cases (less than 4 percent), applications are not funded in order of score. This typically happens if new results emerge to strengthen the application. Scores are never made public.<sup>8</sup>

Funded applications may be renewed every three to five years, in which case they go through the same process described above. Unfunded applications may be resubmitted, during the period of my data, up to two more times. My analysis includes all applications that are reviewed in each of my observed study section meetings, including first-time applications, resubmitted applications, and renewal applications.

### *B. Expertise and Bias among Reviewers*

How likely is it that reviewers have better information about the quality of applications in their own area? In informal interviews with scientists serving on peer review committees, I found that the majority of scientists have more confidence in their assessments of related proposals; for many, this translates into speaking with greater authority during deliberations. Reviewers are also more likely to be assigned as initial reviewers for applications in their area, forcing them to evaluate the proposal in more detail. Even when they are not assigned as initial reviewers, many reviewers said they were more likely to carefully read applications in their own area. These mechanisms suggest that reviewers may have greater "expertise" about related applications, either because they know more to begin with or because they pay more attention.

How likely is it that reviewers in my setting are biased? NIH reviewers have little to no financial stake in the funding decisions they preside over, and conflict of interest rules bar an applicant's coauthors, advisers or advisees, or colleagues from participating in the evaluation process.<sup>9</sup> Yet, there is often significant scope for reviewers to have preferences based on their intellectual connections with applicants. Because NIH support is crucial to maintaining a lab, reviewers are well aware that funding a project in one research area necessarily means halting progress in others. Many of the reviewers I spoke with reported being more enthusiastic about proposals in their own area; several went further to say that one of the main benefits of serving as a reviewer is having the opportunity to advocate for more resources for one's area of research. These preferences are consistent with the idea that reviewers have a taste for research that is similar to theirs, or that they perceive this research to be complementary to their own. On the other hand, some study section members also mentioned that other reviewers—not they—were strategic in terms of evaluating proposals from competing labs.<sup>10</sup> This concern is also supported by research indicating that labs regularly compete over scarce resources, such as journal space, funding, and scientific priority (Pearson 2003).

<sup>8</sup>For more details on the NIH review process, see Gerin et al. (2010).

<sup>9</sup>For this reason, I cannot study the impact of these more social connections on funding outcomes.

<sup>10</sup>I conducted 16 informal interviews with current and past members of NIH study sections. These interviews were off the record but subjects agreed that interested readers could contact the author for more details of these conversations as well as for a full list of the interviewees.

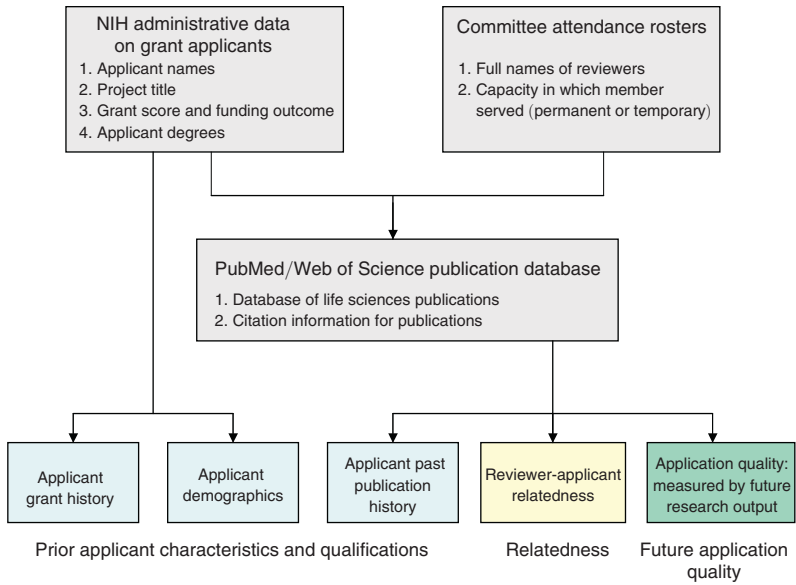


FIGURE 1. DATA SOURCES AND VARIABLE CONSTRUCTION

II. Data

The goal of this paper is to (i) identify how intellectual proximity to influential reviewers affects an applicant’s chances of being funded and (ii) to assess the role of bias versus expertise in these funding decisions.

In order to accomplish this, I construct a new dataset describing grant applications, review committee members, and their relationships for almost 100,000 applications evaluated in more than 2,000 meetings of 250 chartered study sections. My analytic file combines data from three sources: NIH administrative data for the universe of R01 grant applications, attendance rosters for NIH peer review meetings, and publication databases for life sciences research. Figure 1 summarizes how these data sources fit together and how my variables are constructed from them.

I begin with two primary sources: the NIH IMPAC II database, which contains administrative data on grant applications; and a series of study section attendance rosters obtained from NIH’s main peer review body, the Center for Scientific Review. The application file contains information on an applicant’s full name and degrees, the title of the grant project, the study section meeting to which it was assigned for evaluation, the score given by the study section, and the funding status of the application. The attendance roster lists the full names of all reviewers who were present at a study section meeting and whether a reviewer served as a temporary member or a permanent member. These two files can be linked using meeting-level identifiers available for each grant application. Thus, for my sample grant applicants, I observe the identity of the grant applicant, the identity of all committee members, and the action undertaken by the committee.



TABLE 1—GRANT APPLICATION DESCRIPTIVES

	Roster-matched sample		Full sample	
	SD		SD	
<i>Sample coverage</i>				
Number of grants	93,558		156,686	
Number of applicants	36,785		46,546	
Years	1992–2005		1992–2005	
Number study sections	250		380	
Number study section meetings	2,083		4,722	
<i>Grant application characteristics</i>				
Percent awarded	26.08		30.48	
Percent scored	61.58		64.04	
Percent new	70.31		71.21	
Percentile score	70.05	18.42	71.18	18.75
Number of publications (text-matched, in first year after grant review)	0.3	0.8	0.3	0.8
Number of citations (up to 2008, to text-matched publications in first year after grant review)	10	51	11	55
<i>Applicant (PI) characteristics</i>				
Percent female	23.21		22.58	
Percent Asian	13.96		13.27	
Percent Hispanic	5.94		5.79	
Percent MD	28.72		29.26	
Percent PhD	80.46		79.69	
Percent new investigators	19.70		20.02	
Number of publications, past five years	15	60	15	55
Number of citations, past five years	416	1,431	423	1,474

*Notes:* The analytic sample includes new or competing R01 grants evaluated in chartered study sections from 1992 to 2005, for which I have study section attendance data, with social science study sections dropped. The quality of grant applications is measured as follows: number of publications refers to the number of research articles that the grant winner publishes in the year following the grant that share at least one salient word overlap between the grant project title and the publication title. Number of citations refers to the total number of citations that accrue to this restricted set of publications from the time of publication to the end of my citation data in 2008. Applicant characteristics are measured as follows: female, Asian, and Hispanic are all defined probabilistically based on full name. A new investigator is one who has never previously been a PI on an NIH grant. Past publications include any first, second, and last authored articles published in the five years prior to applying for the grant. Number of citations include all citations to those publications, to 2008. Investigators with common names are dropped as are any for which the covariates are missing.

My final sample consists of 93,558 R01 applications from 36,785 distinct investigators over the period 1992–2005. This sample is derived from the set of grant applications that I can successfully match to meetings of study sections for which I have attendance records, which is about half of all R01 grants reviewed in chartered study sections. Of these applications, approximately 25 percent are funded and 20 percent are from new investigators, those who have not received an R01 in the past. Seventy percent of applications are for new projects, and the remainder are applications to renewal existing projects. All of these types of applications are typically evaluated in the same study section meeting. Table 1 shows that my sample appears to be comparable to the universe of R01 applications that are evaluated in chartered study sections.

There are three components to these data: (i) a measure of intellectual proximity between applicants and review committees; (ii) a measure of application quality;

(iii) various measures of other applicant characteristics. Sections IIB and IIC first describe how I measure proximity, application quality, and applicant characteristics, respectively. I describe how my empirical strategy uses these measures later in the text, in Sections III and IV.

### *A. Measuring Proximity*

I measure the intellectual proximity between an applicant and his or her review committee as the number of permanent reviewers who have cited an applicant's work in the five years prior to the meeting, conditional on the total number of such reviewers. This is a measure of how intellectually connected applicants are to the more influential members of their review committees.

I construct proximity in this way for two reasons. First, using citations to measure proximity has several benefits. Citations capture a form of proximity that, as demonstrated by the quote in the introduction, may strongly influence a reviewer's personal preferences: reviewers may prefer work that they find useful for their own research. Citations also capture this form of intellectual connection more finely than other measures, such as departmental affiliation, allowing for more informative variation in proximity. Further, using data on whether the reviewer cites the applicant (as opposed to the applicant citing the reviewer) reduces concerns that my measures of proximity can be strategically manipulated by applicants. Finally, one may also consider more social measures of proximity, such as coauthorship or being affiliated with the same institution. These ties, however, are often subject to NIH's conflict-of-interest rules; reviewers who are coauthors, advisors, advisees, or colleagues, etc., are prohibited from participating in either deliberations or voting. Intellectual proximity is a connection that likely matters for grant review but which is not governed by conflict-of-interest rules.

Second, I focus on being cited by permanent reviewers in order to generate variation in proximity that I will argue is unrelated to an applicant's quality. This is because the total number of reviewers who cite an applicant is likely to be correlated with quality: better applicants may be more likely to be cited and may, independently, submit higher quality proposals. By controlling for the total number of reviewers who cite an applicant, I compare applicants that differ in their proximity to more influential reviewers, but not in the quality of their work. I discuss this strategy and provide evidence for its validity in Section IIIA.

Table 2 describes the characteristics of the sample study sections. In total, I observe 18,916 unique reviewers. On average, each meeting is attended by 30 reviewers, 17 of whom are permanent and 13 of whom are temporary. The average applicant has been cited by two reviewers, one temporary and one permanent. The average permanent and average temporary reviewer both cite four applicants.

### *B. Measuring Quality*

I measure application quality using the number of publications and citations that the research it proposes produces in the future. The key challenge to constructing this measure is finding a way to use ex post publication data to assess the ex ante

TABLE 2—COMMITTEE DESCRIPTIVES

	Roster matched sample	
		SD
Number of reviewers	18,916	
Number of applications	53.73	17.31
<i>Composition</i>		
Number of permanent reviewers per meeting	17.23	4.52
Number of temporary reviewers per meeting	12.35	7.44
Number of meetings per permanent reviewer	3.69	3.03
Number of meetings per temporary reviewer	1.78	1.30
<i>Relatedness</i>		
Total number reviewers who cite applicant	1.94	2.81
Number of permanent reviewers who cite applicant	1.11	1.73
Number of permanent reviewers cited by applicants	4.12	5.32
Number of temporary reviewers cited by applicants	4.12	5.09

*Notes:* See notes to Table 1 for details on the sample. A reviewer is defined as citing an applicant if the reviewer has published a paper in the past five years that has cited any of the applicant's papers. An applicant is defined as citing a reviewer if the applicant has published a paper in the past five years that cites the reviewer's work.

quality of applications. For example, how does one measure the quality of applications that are unfunded if publications cannot acknowledge grants that do not exist?

To overcome this challenge, I develop a way to identify publications associated with research described in the preliminary results section of each grant application. As discussed in Section I, this is possible because it is extremely common for scientists to submit grant proposals based on nearly completed research, especially for the large R01 grants that I study. To find these publications, I first identify all research articles published by a grant's primary investigator. I then use a text matching technique to identify articles on the same topic as the grant application. This is done by comparing each publication's title and abstract with the title of the applicant's grant proposal. For instance, if I see a grant application entitled "Traumatic Brain Injury and Marrow Stromal Cells" reviewed in 2001 and an article by the same investigator entitled "Treatment of Traumatic Brain Injury in Female Rats with Intravenous Administration of Bone Marrow Stromal Cells," I label these publications as related. In my baseline specifications, I require that publications share at least four substantive (e.g., with articles and other common words excluded) overlapping words with the grant project title or its abstract. On average project titles have 10 substantive words, and abstracts have 50. I describe the text matching process I use in more detail in online Appendix B, and show robustness to alternative matching thresholds.

Text matching makes it possible to identify publications associated with unfunded grants. Funding itself, however, may enable scientists to produce more or better research, making it difficult to disentangle the ex ante quality of an application from its ex post publications and citations. This is a particularly important concern for this paper because it affects my ability to distinguish bias from expertise. To see this, suppose that two scientists submit proposals that are of the same ex ante quality, but that one scientist is related to a more influential reviewer, who funds him out of bias. The funding, however, allows this scientist to publish more articles, meaning that an

econometrician that examines ex post outcomes may mistakenly conclude that his proposal was ex ante better. This would lead me to mistake bias for expertise.

Funding may improve a scientist's output both by subsidizing research on topics unrelated to the original application or by supporting work in the same area. I deal with both of these concerns. First, text matching restricts the set of publications I use to assess an application's quality to those that are on the same topic. Second, I also restrict the set of publications I use to assess quality to those published within one year of grant review. This short time window identifies articles based on research that was already completed or underway at the time the application was written. These unlikely to be directly supported by the grant.<sup>11</sup>

This procedure is designed to isolate the set of publications based on the ideas outlined within a grant application. I then use citation information to assess the quality of these ideas. Specifically, for each application, I count the total number of publications, the total number of citations these publications receive through 2012, and the number of "hit" publications, where a hit is defined as being in the ninetieth, ninety-fifth, or ninety-ninth percentiles of the citation distribution relative to all other publications in its cohort (same field, same year). Because my sample begins in 1992 and my citation data go through 2008, I can capture a fairly long-run view of quality for almost all publications associated with my sample grants (citations for life sciences articles typically peak one to two years after publication). This allows me to observe whether a project becomes important in the long run, even if it is not initially highly cited.

Figure 2 plots the relationship between ex ante quality and an application's likelihood of funding, measured using future citations to text-matched publications, and shows that, on average, better applications are more likely to be funded. This provides evidence that, on average, peer reviewers are able to identify high quality applications (Li and Agha 2015). Despite this, online Appendix Figure B shows that many unfunded applications go on to generate more citations and publications than funded applications. This can also be seen in online Appendix Table A, which reports detailed comparisons of the distribution of citations and publications associated funded and unfunded applications: the average funded grant produces 10.3 citations and 0.33 publications, compared with 8.7 citations and 0.26 publications for unfunded applications.

One concern with these figures is that it is possible that funding itself impacts my measure of quality, making it appear as though funded applicants were higher quality ex ante when in fact they were not. To provide evidence that this is not the case, I examine a fuzzy regression discontinuity in funding outcomes around the applicant score payline. If my measure of quality is capturing a grant application's ex ante quality, then it should vary smoothly at this discontinuity. If, instead, funding impacts my measure of quality, then I would see a discontinuous jump in quality at

<sup>11</sup> To compute the appropriate window, I consider funding, publication, and research lags. A grant application is typically reviewed four months after it is formally submitted, and, on average, another four to six months elapse before it is officially funded. See [http://grants.nih.gov/grants/grants\\_process.htm](http://grants.nih.gov/grants/grants_process.htm). In addition to this funding lag, publication lags in the life sciences typically range from three months to over a year. It is thus highly unlikely that articles published up to one year after grant review would have been directly supported by that grant. My results are robust to other windows. See online Appendix Tables F and G.

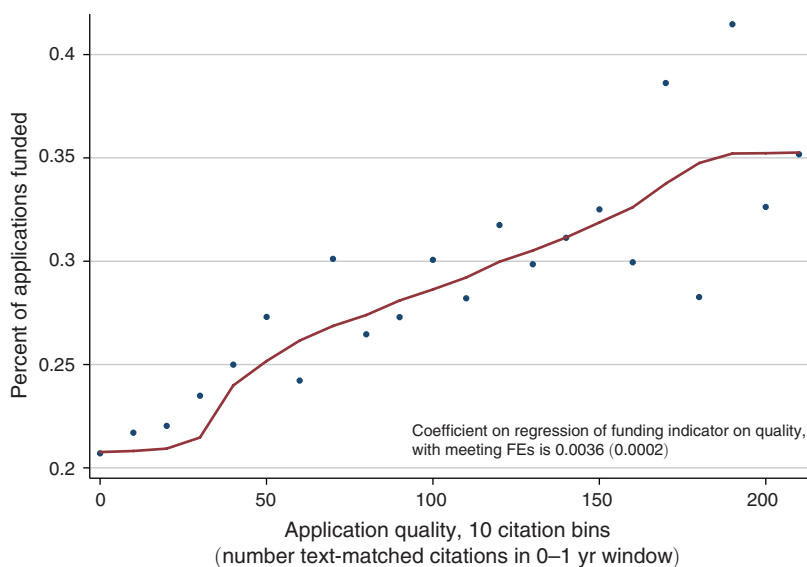


FIGURE 2: RELATIONSHIP BETWEEN APPLICATION QUALITY AND GRANT FUNDING

*Notes:* Application quality is measured using text-matched publications published within one year of grant review, and then computing all citations to this subset of publications, to 2008. See Section IIB and online Appendix B for additional details about how quality is constructed. The  $x$ -axis is the number of such citations, divided into ten-citation bins. Each dot represents the percent of applications that are funded, among applications in the same ten-citation bin. Citations are top coded at the ninety-ninth percentile. This is done for legibility only; analyses use the full distribution of both variables. The plotted line presents a locally smoothed polynomial estimated using a Epanechnikov kernel.

the funding threshold. Figure 3 demonstrates that my quality measure is smooth at this funding threshold, even though the likelihood of being funded changes discontinuously. In both panels of Figure 3, the applicant's percentile score, the running variable, is plotted along the  $x$ -axis. For this figure, the score has been rounded to the nearest integer and re-centered so that 0 represents the funding threshold relevant for that particular application (funding thresholds can differ based on how much funding a particular Institute has been allocated for particular research areas). In panel A of Figure 3, I plot the proportion of applications with that centered score that are ultimately funded, and there is a clear discontinuity at the funding threshold. This is a not a sharp discontinuity because grants can be funded out of order.<sup>12</sup> Panel B of Figure 3 plots centered scores against the average measured quality for each score group. In general, there is a positive slope, indicating that better scoring applications tend to have higher quality, but I find no evidence of a discontinuity at the funding threshold.

The accompanying statistical test is reported in online Appendix Table C. I show that there is no effect of being over the funding threshold, conditional on scores, on measured quality, and further find no effect of funding on measured outcomes,

<sup>12</sup>For example, new investigators may be funded ahead of established investigators with the same score if the funding Institute wants to encourage submissions from younger scientists.

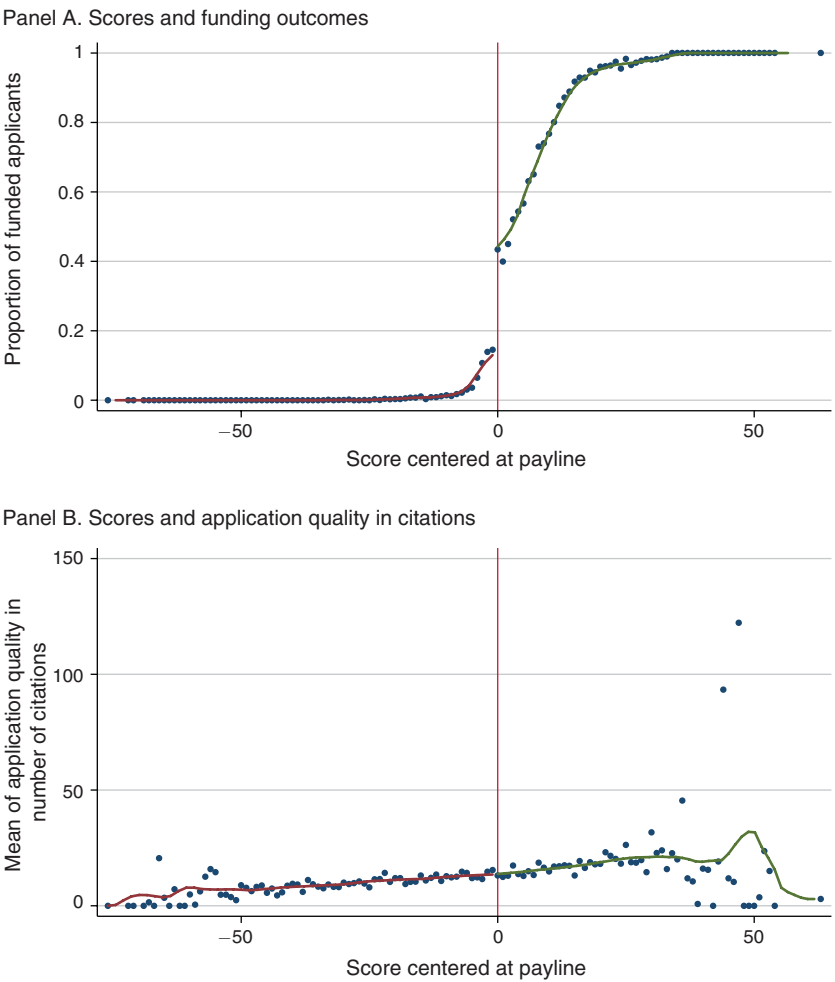


FIGURE 3. REGRESSION DISCONTINUITY IN APPLICATION QUALITY, BY SCORE

*Notes:* The top panel graphs the relationship between scores and the proportion of applicants in that score range who are funded. Scores normally in percentiles from 0 to 100 are centered at their meeting specific payline, which is set to be 0. The data are then collapsed by round values of this centered score. The y-axis plots the proportion of applicants that are funded for applications that share the same centered score, rounded to the nearest integer value. The line is calculated using Epanechnikov kernel weighted local polynomial smoothing. The bottom panel is constructed analogously but with mean application quality using text-matched citations as the dependent variable.

using the funding threshold as an instrument. Together with Figure 3, this finding mitigates concerns that my measure of quality is directly affected by funding.

C. Measuring Applicant Characteristics

Finally, I construct detailed measures of applicant demographics, grant history, and prior publications. Using an applicant’s first and last name, I construct probabilistic



measures of gender and ethnicity (Hispanic, East Asian, or South Asian).<sup>13</sup> I also search my database of grant applications to build a record of an applicant's grant history as measured by the number of new and renewal grants an applicant has applied for in the past and the number he has received. This includes data on all NIH grant mechanisms, including non-R01 grants, such as post-doctoral fellowships and career training grants. To obtain measures of an applicant's publication history, I use data from Thomson Reuters Web of Science (WoS) and the National Library of Medicine's PubMed database. From these, I construct information on the number of research articles an applicant has published in the five years prior to submitting her application, her role in those publications (in the life sciences, this is discernible from the author position), and the impact of those publications as measured by citations. In addition to observing total citations, I can also identify a publication as "high impact" by comparing the number of citations it receives with the number of citations received by other life science articles published in the same year. Sample descriptives for these variables are also provided in Table 1.

### III. Identifying the Causal Impact of Proximity

The first part of my empirical analysis estimates the effect of intellectual proximity to more influential committee members:

$$(1) \quad \text{Assessment}_{icmt} = a_0 + a_1 \text{Proximity\_to\_Permanent}_{icmt} \\ + a_2 \text{Total\_Proximity}_{icmt} + \mu X_{icmt} + \delta_{cmt} + e_{icmt}.$$

$\text{Assessment}_{icmt}$  is a variable describing the committee's assessment (either the funding status, score, or whether an application was scored at all) of applicant  $i$  whose proposal is evaluated by committee  $c$  in meeting  $m$  of year  $t$ .  $\text{Proximity\_to\_Permanent}_{icmt}$  is the number of permanent reviewers who have cited an applicant's work in the five years prior to the committee meeting, and  $\text{Total\_Proximity}_{icmt}$  is the total number of such reviewers. The covariates  $X_{icmt}$  include indicators for sex; whether an applicant's name is Hispanic, East Asian, or South Asian; quartics in an applicant's total number of citations and publications over the past five years; indicators for whether an applicant has an MD and/or a PhD; and indicators for the number of past R01 and other NIH grants an applicant has won, as well as indicators for the number to which she has applied. The  $\delta_{cmt}$  are fixed effects for each committee meeting so that my analysis compares outcomes for grants that are reviewed by the same reviewers in the same meeting. Standard errors are clustered at the committee-fiscal-year level.

My coefficient of interest is  $a_1$ .  $a_1$  compares the funding outcomes of scientists whose applications are reviewed in the same meeting, who have similar past performance, and who, while cited by the same total number of reviewers, differ in their proximity to permanent reviewers.

<sup>13</sup> For more details on this approach, see Kerr (2008). Because black or African American names are typically more difficult to distinguish, I do not include a separate control for this group.

TABLE 3—CHARACTERISTICS OF PERMANENT AND TEMPORARY REVIEWERS

Reviewer characteristics	Permanent	Temporary	<i>p</i> -value	
<i>Demographics</i>				
Percent female	31.68	24.28	0.00	
Percent Asian	14.99	13.08	0.00	
Percent Hispanic	6.40	5.05	0.00	
<i>Education</i>				
Percent MD	27.42	25.85	0.00	
Percent PhD	79.00	81.00	0.00	
<i>Past citations</i>				
Mean	1,470	1,375	0.00	
Median	606	590	0.09	
Fifth	0	0	0.00	
Ninety-fifth	5,459	5,002	0.00	
<i>Past publications</i>				
Mean	53	57	0.05	
Median	22	21	0.00	
Fifth	0	0	0.00	
Ninety-fifth	154	152	0.67	
Reviewer transitions (1997 to 2002 subsample)	% perm. in the past	% perm. in the past	% perm. in the past	% perm. in the past
Current permanent	61.87	63.71	38.11	35.45
Current temporary	16.25	41.30	32.73	50.13

Notes: Observations are at the reviewer-study section meeting level. The sample includes all reviewers in chartered study sections from 1992 to 2005, for which I have study section attendance data. Number of reviewer publications include any first, second, and last authored articles published in the five years prior to the study section meeting date for which the reviewer is present. Number of citations refers to all citations accruing to those publications, to 2008. Reviewer characteristics are measured as follows: female, Asian, and Hispanic are all defined probabilistically based on full name. MD and PhD are defined based on rosters in which a reviewer’s degree follows his or her name. Reviewer transitions are calculated based on whether a reviewer is present in the roster database during the full sample years from 1992 to 2005. The set of reviewers used in this calculation are those present in meetings from 1997 to 2002 in order to allow time to observe members in the past and future within the sample.

A. Permanent versus Temporary Reviewers

In order for  $a_1$  to identify a causal effect, I need to show that proximity to permanent reviewers is not correlated with other characteristics that may also impact funding outcomes, conditional on proximity to all reviewers. Foremost, one may be concerned that being cited by permanent reviewers signals higher quality than being cited by temporary reviewers.

Before providing direct evidence to refute this claim, I first explore the characteristics of permanent versus temporary reviewers. Table 3 compares demographic, educational, and publication characteristics of permanent and temporary reviewers. Permanent reviewers tend to be somewhat more diverse: 32 percent are women, relative to 25 percent of temporary reviewers; 15 percent are Asian; and 6.4 percent Hispanic, compared to 12 percent and 5.1 percent, respectively, among temporary reviewers.<sup>14</sup> This difference is likely due to the fact that the NIH makes

<sup>14</sup>I do not have information on whether a reviewer is black because my demographic variables come from analyzing names. Black names are more difficult to recognize relative to Asian or Hispanic names.

a conscious effort to ensure diversity on their review panels. Similarly, permanent members are also slightly more likely to be medical doctors with some clinical experience, relative to PhDs. This difference, 27 percent versus 26 percent, is statistically significant but small.

Table 3 also shows that permanent reviewers appear to have slightly stronger publication histories, as measured by the number of publications in the previous five years and the number of forward citations to those publications, measured up to the year 2008. For example, permanent reviewers average 1,470 citations to past publications compared to 1,375 for temporary reviewers. This difference of about 100 citations, while significant, is small relative to the standard deviation of citations, which is approximately 3,000 for both groups. Meanwhile, reviewers appear to be similar in terms of their number of publications. The full distribution is plotted in online Appendix Figure A.

The findings in Table 3 so far highlight several potentially important differences in the qualifications and characteristics of permanent and temporary reviewers. The bottom panel of Table 3 provides some evidence for why permanent and temporary reviewers still nonetheless appear broadly similar: during the course of their career, the same person often serves in both capacities. This difference does not simply reflect a progression from temporary to permanent as reviewers age. Rather, for the set of reviewers observed in the middle of my sample, between 1997 and 2002, 35 percent of permanent reviewers in a given meeting will serve as temporary reviewers in a future meeting while 40 percent of temporary reviewers in a given meeting will serve as permanent reviewers in a future meeting. These common changes in reviewer status across meetings mitigates concerns that permanent reviewers are categorically different from temporary members.

My next set of results explore the matching of applicants to permanent or temporary reviewers, which is nonrandom in two ways. First, rosters listing the permanent (but not temporary) reviewers associated with a study section are publicly available, meaning that applicants know who some of their potential reviewers may be at the time they submit their application. The scope for strategic submissions in the life sciences, however, is small: for most grant applicants, there are only one or two intellectually appropriate study sections and, because winning grants is crucial for maintaining one's lab and salary, applicants do not have the luxury of waiting for a more receptive set of reviewers. Second, assignment is also nonrandom because study section administrators assign applications to initial reviewers on the basis of (i) intellectual match and (ii) reviewer availability. If, for instance, not enough permanent reviewers are qualified to evaluate a grant application, then the study section administrator may call in a temporary reviewer. Temporary reviewers may also be called if the permanent members qualified to review the application have already been assigned too many other applications to review.

This process may raise concerns for my identification. For example, suppose that two applicants, one better known and higher quality, submit their applications to a study section that initially consists of one permanent reviewer. The permanent reviewer is more likely to be aware of the work of the better-known applicant and thus there would be no need to call on a related temporary member. To find reviewers for the lesser-known applicant, however, the administrator calls on a temporary

TABLE 4—APPLICANT CHARACTERISTICS, BY NUMBER AND COMPOSITION OF RELATED REVIEWERS

Dep. var.: Applicant characteristics	Female (1)	Asian (2)	Hispanic (3)	MD (4)	PhD (5)	New investigator (6)	Previous publications (7)	Previous citations (8)	Application quality (9)	Funding propensity (× 100) (10)
Number of proximate permanent reviewers	−0.0003 (0.002)	−0.0013 (0.002)	−0.001 (0.0010)	0.0119 (0.011)	0.0036 (0.010)	−0.0013 (0.002)	−0.0057 (0.360)	52.1057 (13.250)	0.0065 (0.0060)	0.0433 (0.0320)
Observations	93,558	93,558	93,558	93,558	93,558	93,558	93,558	93,558	93,558	93,558
R <sup>2</sup>	0.0627	0.0381	0.0252	0.0452	0.0588	0.0679	0.0665	0.2071	0.0312	0.2248
Meeting FEs	X	X	X	X	X	X	X	X	X	X
Number of proximate reviewer FEs	X	X	X	X	X	X	X	X	X	X

Notes: See notes to Table 1 for details about the sample. Coefficients are reported from a regression of applicant characteristics on the number of permanent members related to an applicant, controlling for meeting-level fixed effects, and fixed effects for proximity to all reviewers. Proximity to permanent reviewers is defined as the number of permanent reviewers who have cited any of the applicant’s research in the five years prior to grant review. Outcome variables are defined as follows: female, Asian, and Hispanic are all defined probabilistically based on full name. MD and PhD are from administrative grant application records. A new investigator is one who has never previously been a PI on an NIH grant. Previous publications include any authored articles published by the applicant in the five years prior to applying for the grant. Previous citations include all citations to those publications, to 2008. Application Quality is text-matched citations to grant applications, in hundreds, as described in the text. Funding propensity is an aggregate variable constructed from regressing funding outcomes on all demographic variables, education, fixed effects for decile bins for both past publication and citations, and fixed effects for number of past R01 and other grants, and taking the fitted values of funding likelihood from this regression.

reviewer. Both applicants would then be related to one reviewer in total but, in this example, the fact that one applicant works in the same area as a temporary member is actually correlated with potentially unobserved aspects of quality.

Table 4 shows that this may be a potential concern. Table 4 regresses the demographic characteristics and publication records of applicants on the number of permanent reviewers they have been cited by, conditional on meeting fixed effects and fixed effects for the total number of citing reviewers. This compares two applicants to the same meeting, who have been cited by the same total number of reviewers, but who differ in the number of citing permanent and temporary reviewers. I find that applicants cited by more permanent members do not differ on any demographic or educational characteristics, but that there is a significant relationship between being cited by permanent reviewers and previous citations in column 8: each additional permanent reviewer who cites an applicant is associated with 52 more citations, relative to a mean of 1,429 citations, or about 3.6 percent. Column 9, however, shows that, conditional on relatedness to all reviewers, relatedness to permanent reviewers is not predictive of the quality of applications themselves. Next, column 10 shows that proximity to permanent reviewers is not correlated with an application’s predicted propensity to be funded. To show this, I regress application funding on applicant gender, race, education, past publications, past citations, and past grant history to construct an index describing his or her propensity to be promoted. I find that there is no relationship between proximity to permanent reviewers and this index.

Despite appearing similar in terms of qualifications, the structure of the NIH is such that permanent reviewers play a larger role in making funding decisions. There are several reasons why this is the case. Most basically, permanent reviewers

do more work. As discussed in Section I, reviewers are responsible for providing initial assessments of a grant application before that application is discussed by the full committee. These initial assessments are extremely important for determining a grant application's final score because they (i) determine whether a grant application even merits discussion by the full group and (ii) serve as the starting point for discussion. Study sections also evaluate 40 to 80 applications per meeting, meaning that it is unlikely that reviewers have had a chance to carefully read proposals to which they have not been officially assigned. In many study sections, moreover, there is also a rule that no one can vote for scores outside of the boundaries set by the initial scores without providing a reason.

While I do not have data on who serves as one of an application's three initial reviewers, permanent reviewers are much more likely to serve as an initial reviewer; they are typically assigned eight to ten applications, compared with only one or two for temporary reviewers. In addition, permanent members are required to be in attendance for discussions of all applications; in contrast, temporary members are only expected to be present when their assigned grants are discussed, meaning that they often miss voting on other applications. Finally, permanent members work together in many meetings over the course of their four-year terms; they may thus be more likely to trust, or at least clearly assess, one another's advice, relative to the advice of temporary reviewers with whom they are less familiar. As a result, being evaluated by a close permanent reviewer can impact how an application is treated, even though it is not correlated with a grant's underlying quality.

The results in Tables 3 and 4 raise a potential concern: permanent and temporary reviewers have statistically different publication histories and permanent reviewers are more likely to be assigned to applicants with more previous citations themselves. Before testing whether relatedness impacts funding decisions, I present several pieces of evidence to show that this is unlikely to bias my results.

First, the results in Table 3 should not be interpreted as presenting a case that permanent members have stronger publication histories than temporary reviewers. Rather, permanent reviewers have more past citations while temporary reviewers have more past publications. More generally, differences in the qualifications of permanent and temporary reviewers would only impact my results if it translates into a relationship between proximity to permanent reviewers and an application's quality itself. Table 4, online Appendix Table B, and online Appendix Figure C provide direct evidence that this is not the case. In particular, the upper left-hand panel of online Appendix Figure C shows the distribution application quality (as defined in the previous section) for applicants cited by exactly one reviewer. The solid line shows the distribution of quality among applicants cited by one permanent reviewer and the dotted line does so for those cited by one temporary reviewer. These distributions are statistically indistinguishable: a Kolmogorov–Smirnov test cannot reject the null that these two distributions are equal. Similarly, the upper right-hand panel shows the same, but with quality measured using the number of publications associated with a grant. The bottom two panels of online Appendix Figure C repeat this exercise for applicants who have been cited by a total of two reviewers. In this case, there are now three possibilities: the applicant has been cited by two temporary reviewers, two permanent, or one of each. In all of these cases, the distribution

of applicant quality is statistically similar.<sup>15</sup> Online Appendix Table B compares means and other percentiles of these distributions.

Further, if my measure of relatedness were correlated with unobserved factors that also impact an application's likelihood of funding, then I would expect the inclusion of applicant characteristics (publication history, demographics, etc.) to impact my estimates. In Section IIIB, I show that this is not the case: the impact of relatedness that I estimate does not change when I include detailed controls for applicant characteristics and publication histories.

Finally, I provide two additional complementary sets of analysis. The first uses reviewer fixed effects to show that applicants are more likely to be funded when the reviewer that has cited them is serving as a permanent reviewer, compared to when that reviewer is serving as a temporary reviewer. This is presented in Table 6 and discussed in Section IIIB. I also show that my results do not rely on the distinction between permanent and temporary reviewers by using applicant fixed effects to compare outcomes for the same applicant across meetings in which she is cited by different numbers of reviewers. This alternative specification identifies the effect of being related to an *additional* reviewer under the assumption that the time-variant unobserved quality of an application is not correlated with proximity. This is presented in online Appendix Table K and discussed in Section V, and online Appendix Section D.

### B. Impact of Proximity: Results

Table 5 estimates the effect of intellectual proximity on funding and scores. The first column reports the raw within-meeting association between proximity to permanent reviewers and an applicant's likelihood of being funded. Without controls, each additional permanent reviewer who has cited an applicant is associated with a 3.3 percentage point increase in the probability of funding, from an overall average of 21.4 percent. This translates into a 15.4 percent increase. Most of this correlation, however, reflects differences in quality—better applicants are more likely to be cited by reviewers. Column 2 adds a full set of fixed effects for the total number of reviewers who have cited an applicant. Once I do this, my identifying variation comes from changes to the *composition* of the reviewers who have cited an applicant—effectively the impact of switching the reviewers an application is related to from temporary to permanent. With these controls, the impact of being cited by a permanent reviewer falls to 0.0050, but remains significant. This says that comparing two scientists reviewed in the same meeting, cited by the same number of applicants, being cited by an additional permanent reviewer increases an applicant's likelihood of funding by 0.0050/0.214 or 2.3 percent.

To appreciate the magnitude of this effect, it is useful to consider how sensitive funding decisions are to changes in application quality. Recall from Figure 2 that an application's likelihood of funding is increasing in its quality. A regression of funding on application quality, holding constant meeting fixed effects, says that an applicant's likelihood of funding increases by 3.6 percentage points for every 100

<sup>15</sup> Approximately 75 percent of my sample of applications are cited by two or fewer reviewers. This pattern also holds for applicants cited by three or more reviewers.



TABLE 5—WHAT IS THE IMPACT OF PROXIMITY ON COMMITTEE ASSESSMENTS?

	1(Score is above the payline) Mean = 0.214, SD = 0.410			Score Mean = 71.18, SD = 18.75			1(Scored at all) Mean = 0.640, SD = 0.480		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Number of proximate permanent reviewers	0.0328 (0.001)	0.0050 (0.002)	0.0047 (0.002)	1.1067 (0.054)	0.1641 (0.094)	0.1611 (0.093)	0.0500 (0.002)	0.0012 (0.002)	0.0011 (0.002)
Observations	93,558	93,558	93,558	57,613	57,613	57,613	93,558	93,558	93,558
R <sup>2</sup>	0.0630	0.0688	0.0950	0.1186	0.1224	0.1439	0.0775	0.0899	0.1340
Meeting FEs	X	X	X	X	X	X	X	X	X
Number of proximate reviewer FEs		X	X		X	X		X	X
Past performance, past grants, and demographics			X			X			X

Notes: See notes to Table 1 for details about the sample. Coefficients are reported from a regression of committee decisions (above payline, score, or scored at all) on the number of permanent members related to an applicant, controlling for meeting-level fixed effects. Proximity to permanent reviewers is defined as the number of permanent reviewers who have cited any of the applicant's research in the five years prior to grant review. "Past performance, past grants, and demographics" include indicators for sex and whether an applicant's name is Hispanic, East Asian, or South Asian, indicator variables for deciles of an applicant's total number of citations and publications over the past five years, indicators for whether an applicant has an MD and/or a PhD, and indicators for the number of past R01 and other NIH grants an applicant has won, as well as indicators for how many she has applied to.

citation increase in quality. A back of the envelope calculation comparing this with the coefficient in column 2 says that proximity helps an applicant get funded by as much as would be expected from a 13.8 citation (or one-fourth standard deviation) increase in the quality of the application itself. This sizable effect suggests that when reviewers cannot easily predict the quality of applications, other factors like relatedness play a comparably larger role.

Finally, column 3 adds controls for applicant characteristics such as demographics, education, past publications, past citations, and grant history. If proximity to permanent members is not correlated with applicant characteristics conditional on proximity to all reviewers, we would not expect the addition of these controls to alter our estimates. We find that this is indeed the case: the addition of this large set of controls changes the estimated coefficient from 0.0050 to 0.0047; the percentage change decreases to 2.2 percent.

Columns 6 and 9 report estimates of the impact of proximity on the score that an application receives and whether an application is scored at all. I find a statistically significant but relatively small effect of proximity in scores: switching to a proximate permanent reviewer increases, holding total proximity constant, an applicant's score by 0.16 points or about 1 percent of a standard deviation. I find no evidence that relatedness increases the overall probability that an applicant is scored at all (recall that about 40 percent of applicants are deemed sufficiently weak that they are not given a score).

Table 6 considers an alternative test: do applicants fare differently when the reviewer they have been cited by serves as a permanent member, compared to when that same reviewer serves as a temporary member? Finding that proximity matters more when reviewers are permanent would strongly suggest that reviewer preferences influence funding decisions, independently of the characteristics of applications themselves.

TABLE 6—WHAT IS THE IMPACT OF REVIEWER STATUS ON FUNDING OUTCOMES?  
REVIEWER FIXED EFFECTS

	Proportion of proximate applicants who are funded Mean: 0.37, SD: 0.36 (1)	Average score of proximate applicants Mean: 73.3, SD: 14.3 (2)
Reviewer is permanent	0.003 (0.001)	0.336 (0.144)
Observations	15,871	15,870
R <sup>2</sup>	0.954	0.571
Reviewer FEs	X	X
Past performance, past grants, and demographics	X	X

*Notes:* This table examines how outcomes for applicants cited by reviewers vary by whether the citing reviewer is serving in a permanent or temporary capacity. The sample is restricted to 4,909 reviewers who are observed both in temporary and permanent positions. An applicant is said to be proximate if a reviewer has cited that applicant in the five years prior to the study section meeting in which the reviewer and applicant are matched. “Past performance, past grants, and demographics” include indicators for sex and whether an applicant’s name is Hispanic, East Asian, or South Asian, indicator variables for deciles of an applicant’s total number of citations and publications over the past five years, indicators for whether an applicant has an MD and/or a PhD, and indicators for the number of past R01 and other NIH grants an applicant has won, as well as indicators for how many she has applied to.

To test this, I use the fact that I observe almost 5,000 unique reviewers in meetings in which they are permanent and in meetings in which they are temporary. For each reviewer meeting, I identify the grant applicants cited by that reviewer within the previous five years and calculate the proportion of those applications who are funded, and the average score of those applications from cited scientists. I then regress this outcome variable on an indicator for whether or not the reviewer is serving as a permanent member during that meeting, controlling for reviewer fixed effects and average demographic, publication, and grant characteristics about the set of related applicants, weighted by the number of related applicants. In column 1 of Table 6, I show that a larger proportion of related applicants are funded when the reviewer is permanent rather than temporary. In column 2, I also find that the average score of related applicants increases when a reviewer is permanent.

IV. Expertise versus Bias

My results in the previous section show that applicants who work in the same area as more influential committee members are more likely to be funded. Is this a problem for peer review? Not necessarily. Reviewers may advocate for candidates in their area simply because they are more confident in their assessments: receiving more precise signals about related applicants allows reviewers to form higher posterior expectations about their quality. This could lead to a greater proportion of related applicants falling above the funding bar even in the absence of bias. Because this type of behavior improves the quality of peer review, while biases do not, it is important to distinguish between the two explanations.

The results in Tables 5 and 6 provide initial evidence that this effect is not simply driven by better information. To see this, suppose that related reviewers are unbiased but better informed. In this case, the more precise signals they receive should increase the variance of their posteriors over the quality of applications from related applicants, but should not change the mean, as long as scores are linear in beliefs. While this would lead to a greater proportion of related applicants being above the funding threshold, it should not change the average score for these applicants, because reviewers would also have stronger negative posteriors as well. By contrast, Tables 5 and 6 show that average scores are also higher for related applicants.

A more direct evaluation of the role of expertise and bias involves using information on applicant quality. In general, related reviewers can be (i) only biased, (ii) only better informed, or (iii) both.<sup>16</sup> If reviewers are only biased, they will give better assessments to related applicants regardless of their quality. The impact of relatedness should be to increase an applicant's likelihood of funding for any level of quality. By contrast, reviewers can also be unbiased and better informed about the quality of related applicants. In this scenario, high quality applicants benefit from being evaluated by related reviewers who can more accurately observe their quality, but low quality applicants are hurt for the same reason. The estimated impact of relatedness should be increasing in quality and also be negative for particular low quality applications. Finally, related reviewers may be both biased and better informed. Bias means that assessments are shifted up for related applicants regardless of quality, but information means that there is also a stronger relationship between quality and assessments for intellectually related applicants. In this scenario, the impact of relatedness is still increasing in application quality, but it may not necessarily ever be negative. In this section, I examine which of these scenarios best characterizes the committee evaluations I observe.

An important caveat to note about the following analysis is that it is not possible to pin down the exact amount of bias or information without further parametric assumptions about the nature of reviewers' beliefs. In online Appendix F, I present and estimate a model of grant allocation with biased experts in which the separate contribution of bias and expertise can be precisely estimated. In this model, I make a stronger set of distributional assumptions that allow me to attribute differences in the slope of the relationship between quality and funding outcomes (between related and unrelated applicants) to the value of expertise and to attribute level differences in committee assessments to the impact of bias. Because this model relies on strict distributional assumptions, I proceed with the more flexible approach of splitting my sample into quality bins.

### *A. Expertise versus Bias: Results*

In this section, I examine how the impact of relatedness differs by application quality. Figure 4 presents the relationship between quality and funding estimated on

<sup>16</sup>For simplicity and because this is suggested by my earlier findings, I show the case in which reviewers are positively biased and better informed, although it is theoretically possible that related reviewers can be negatively biased or less informed.

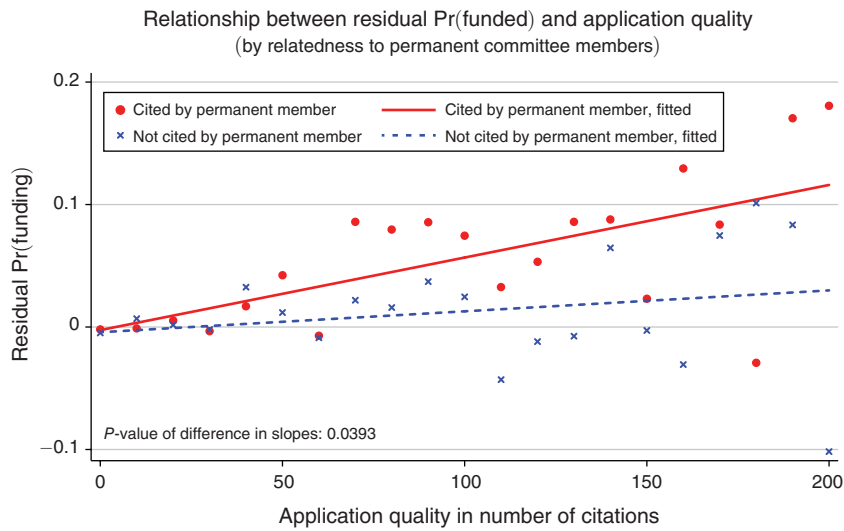


FIGURE 4. REALIZED COMMITTEE ASSESSMENTS BY APPLICATION QUALITY

*Notes:* This figure examines committee assessments by application quality. The  $x$ -axis represents ten citation bins of application quality, where citations are measured from text-matched publications within one year of grant review. The  $y$ -axis represents an application's residual likelihood of funding, after taking into account applicant demographics, education, publication history, grant history, and total number of proximate reviewers. For each quality bin, each dot represents the mean residual funding likelihood for applicants with that quality who have been cited by at least one permanent member. The  $X$ 's represent the same for applicants who have not been cited by a permanent member. See the text in Section IVA for more discussion.

my actual data. The  $x$ -axis corresponds to application quality, measured in citations to text-matched publications. The  $y$ -axis plots an application's residual likelihood of being funded, after taking into account demographics, education, publication record, grant history, and total number of related reviewers.

Figure 4 first shows that an applicant's chances of funding, controlling for applicant characteristics, are increasing in quality for both related and unrelated applicants. I find a stronger slope for applicants cited by permanent reviewers, indicating that funding decisions are more responsive to quality when grants are evaluated by more closely related reviewers. This difference between these slopes is significant at the 5 percent level.

The relationship I estimate between quality and funding for related and unrelated candidates suggests that reviewers are both better informed and more biased when it comes to evaluating related applicants. First, the impact of relatedness is increasing in applicant quality: high-quality applicants are more likely to be funded when evaluated by related reviewers, compared to equally high-quality applicants evaluated by unrelated reviewers. This fact is suggestive of expertise. At the same time, if related reviewers were unbiased, we would expect low-quality applications from related candidates to be penalized. The fact that low-quality applications do not benefit but are not hurt suggests that bias can undo some potential benefits of

TABLE 7—WHAT IS THE IMPACT OF PROXIMITY BY APPLICATION QUALITY?

	Quartiles of residual application quality				
	All (1)	Bottom (2)	Second (3)	Third (4)	Top (5)
<i>Panel A. Dependent variable: 1(score above payline)</i>					
Number of proximate permanent reviewers	0.0047 (0.002)	0.0004 (0.004)	0.0012 (0.004)	0.0038 (0.005)	0.0126 (0.004)
Observations	93,558	22,463	23,929	23,360	23,806
R <sup>2</sup>	0.0950	0.1680	0.1415	0.1311	0.1613
<i>Panel B. Dependent variable: Score</i>					
Number of proximate permanent reviewers	0.1611 (0.093)	−0.0466 (0.169)	0.0239 (0.220)	0.1927 (0.293)	0.6509 (0.216)
Observations	57,613	16,081	14,593	12,056	14,883
R <sup>2</sup>	0.1439	0.2139	0.2222	0.2612	0.2223
<i>Panel C. Dependent variable: 1(scored at all)</i>					
Number of proximate permanent reviewers	0.0011 (0.002)	0.0007 (0.004)	−0.0064 (0.004)	0.0026 (0.007)	0.0046 (0.005)
Observations	93,558	22,463	23,929	23,360	23,806
R <sup>2</sup>	0.1340	0.1771	0.1681	0.1667	0.1924
Meeting FEs	X	X	X	X	X
Number of proximate reviewer FEs	X	X	X	X	X
Past performance, past grants, and demographics	X	X	X	X	X

*Notes:* See notes to Table 1 for details about the sample. Coefficients are reported from a regression of committee decisions (above payline, score, or scored at all) on number of proximate permanent reviewers, controlling for meeting-level fixed effects and fixed effects for total proximity. Panel A regressions use the same specification as column 2 in Table 5; panel B uses the same specification as column 5 in Table 5; panel C uses the same specification as column 8. In particular, column 1 of this table replicates results from Table 5. Columns 2 through 5 split the sample based on quartiles of residual application quality. To calculate residual quality, I regress application quality in citations on dummies for female, Hispanic, east Asian, south Asian, MD, PhD, fixed effects for decile bins for both past publication and citations, and fixed effects for number of past R01 and other grants, and taking the residuals from this regression. Regressions also control for these variables directly.

expertise. The remainder of this section explores this pattern more rigorously and addresses alternative explanations.

Table 7 splits my main sample into quartiles of quality and separately estimates the impact of proximity for each subsample. These quality quartiles are constructed with respect to residual future citations: applications in the highest quartile are those that received more citations than would be predicted given the applicant's initial publications, demographics, and grant history. If closer reviewers had more expertise, we would expect proximity to be particularly beneficial for strong applications that might otherwise not look competitive based on *ex ante* observables.

Column 1 of Table 7 replicates results from columns 3, 6, and 9 of Table 5—the impact of proximity for the entire sample, controlling for meeting fixed effects, fixed effects for total proximity, and controls for applicant characteristics. Columns 2–5 confirm the pattern laid out by Figure 4. I find no effect of proximity for the first two quartiles, a positive but statistically insignificant effect for third quartile applications, and a large and significant effect for the applications in the top quartile. This means that strong applications benefit much more from being evaluated

by an intellectually related reviewer, relative to applications that are weaker. The coefficient on relatedness for top quartile applications is 0.0126 (column 5). In percentage terms, this says that each additional related permanent reviewer increases the funding likelihood of top quartile applicants by 1.26 percentage points, from top-quartile mean of 22.3 percent—a 5.7 percent increase. By contrast, I estimate a 2.2 percent increase for the whole sample (column 1) and a zero effect for bottom quartile applicants (column 2).

There are, however, several alternative explanations that would generate this pattern even in the absence of expertise. First, my estimates could be misleading if my measure of quality is directly impacted by whether or not an application is funded. To see this, suppose that many low-quality applications from related applicants are funded because of bias. If funding itself makes these applications appear stronger *ex post*, then they may mistakenly be categorized as high-quality applications. If this were the case, then what should be identified as a large effect of proximity for low-quality applicants becomes identified instead as a large effect for high-quality applicants. This would lead me to mistake bias for expertise. In discussing my measure of quality in Section IIB, I provided evidence that my quality measure was not contaminated by funding status. Specifically, Figure 3 shows that my measure of quality does not discontinuously change as a result of funding, making it very unlikely that my findings in Table 7 are driven by this pattern.

Another possible explanation for my findings is that bias is more pivotal for high-quality applications simply because they are closer to the funding threshold. Were this the case, the impact of proximity would appear to be increasing in quality even if close reviewers were no better informed. I provide several explanations for why this is unlikely to drive my results.

First, it is not empirically the case that only the highest quality grants have a chance of receiving funding. If this were the case, we would expect a grant's probability of funding to be low up to a threshold, then peak and remain high. Instead, Figure 2 shows that a grant's probability of funding is linearly increasing in quality. It is also not the case that low-quality applications have no chance of funding: 21 percent of applications with zero future citations are funded compared to 24 percent of applications with greater than zero future citations.<sup>17</sup> Because the likelihood of funding is similar for these groups, an equally strong push by close reviewers would have similar effects on a grant's likelihood of funding across quartiles, which is not what I find.

Second, if reviewers were simply biased, then we should expect proximity to increase application scores evenly across all quality bins. The middle panel of Table 7, however, shows that the impact of proximity on scores also increases in quality. I find a much larger effect for the top quartile than for any of the other bins. For scores, though not significantly, I in fact estimate a negative coefficient for the bottom half of quality, suggesting that proximity may actually hurt lower quality applicants.

<sup>17</sup> Further, 26.9 percent of grants in the bottom quartile of residual quality are funded, compared with only 22.3 percent in my top residual quality quartile.



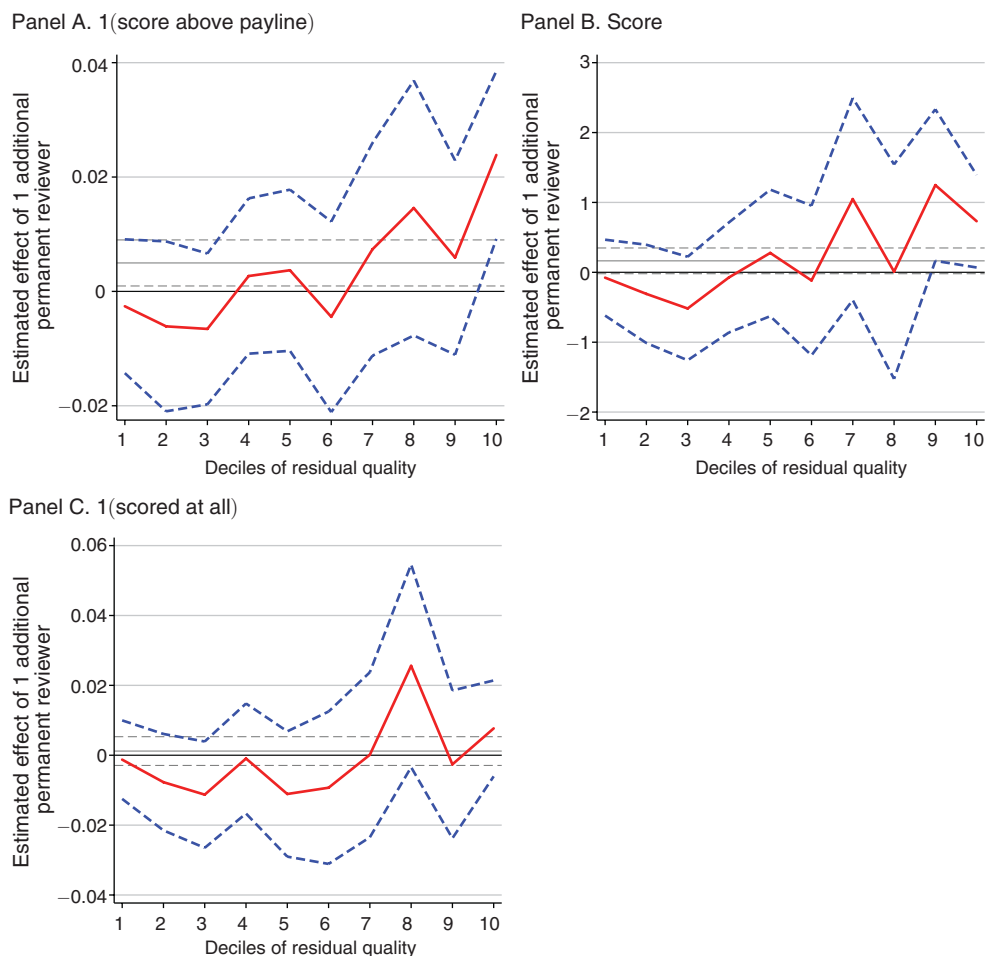


FIGURE 5. IMPACT OF PROXIMITY, BY QUALITY DECILE

*Notes:* The x-axis represents deciles of application quality, where quality is measured as the residual number of citations to text-matched publications that an application receives, after taking into account applicant demographics, education, publication history, and grant history. The y-axis represents coefficients from a regression of funding on number of related permanent reviewers, controlling for meeting effects and fixed effects for total related reviewers, for applications in that decile of quality. Dashed lines represent 95 percent confidence intervals. The gray horizontal line represents the effect of proximity (and confidence interval) estimated on the entire sample.

To explore this pattern more thoroughly, Figure 5 plots the estimated effect of proximity by decile of application quality. I find that the estimated impact of proximity generally increases with application quality and that, for particularly low or high deciles, my estimated effect is outside of the 95 percent confidence interval associated with the average effect in the whole sample, shown in gray. When I break quality to this finer level, I find suggestive evidence of a negative impact of proximity on candidates who submit applications that are especially poor relative to their qualifications. The coefficients are not statistically significant but their magnitude is economically meaningful: a bottom decile application evaluated by a close reviewer is 0.3 percentage points less likely to be funded than a similarly poor application

evaluated by an related reviewer, off a baseline funding likelihood of 30 percent for this group, making for a 1 percent decrease. Because each coefficient estimate represents the joint effect of bias and information, it is likely that the coefficient would be even more negative in the absence of potentially offsetting positive biases. By contrast, a proximate applicant in the top decile is 2.4 percentage points more likely to be funded than a similar applicant who is not related, off a baseline for this group of 26 percent, making for a 9.2 percent increase. This pattern strongly suggests that close reviewers have additional expertise about applicants. Figure 5 also shows a similar pattern for application scores.

In online Appendix Table J, I consider how the impact of proximity may differ for new and renewal applications, and for new and experienced investigators. For new grants, I find a sharper negative effect of proximity on funding likelihood for low-quality applications, and similar positive effects for high-quality applicants. This suggests that related reviewers have differential better information about the quality of these new applications. My estimates for new investigators are noisier and I can neither reject my estimates being different from zero or different from my estimates for other samples. These results are discussed in more detail in online Appendix C.

Finally, my results also consider the impact of proximity on an application's likelihood of being scored at all (e.g., rejected early in the process due to low initial evaluations). In general, I do not find significant effects, although my pattern of estimates is broadly consistent with the pattern that I find when examining funding likelihood and scores. See the last panel of Table 7.

### B. Proximity and Overall Portfolio Quality

My main results show that (i) applicants who are related to study section members are more likely to be funded, independent of quality, as measured by the number of citations that their research eventually produces; and (ii) better scientists benefit more from proximity, suggesting that study section members are better at discerning the quality of applicants in their own area.

Next, I assess the impact of proximity on the overall quality of funded applications, at the meeting level. For each meeting, I construct variables describing the overall share of applicants who have been cited by permanent reviewers, the overall share who have been cited by any reviewer, and the average ex ante measured quality associated with funded applications, as well as with all applications. I then regress the average quality of funded applicants on the share of related permanent reviewers, controlling for, at the minimum, year fixed effects and study section fixed effects.

Table 8 provides suggestive evidence that, on net, proximity increases the overall quality of funded grants. Meetings in which a greater share of applicants have been evaluated by a closely related permanent reviewer fund grants fund, on average, higher quality applications. Column 1 examines the relationship between share related to permanent including year and study section fixed effects. Year effects are needed to account for the fact that grants funded in later years have less time to accrue applications. Study section effects control for field-level differences in

TABLE 8—WHAT IS THE EFFECT OF PROXIMITY ON THE QUALITY OF THE NIH PORTFOLIO?

	Mean grant application quality, meeting level					
	Awarded grants			All grants		
	Mean = 0.144, SD = 0.265			Mean = 0.111, SD = 0.202		
	(1)	(2)	(3)	(4)	(5)	(6)
Share of applicants cited by permanent reviewers	0.0856 (0.051)	0.0742 (0.065)	0.0714 (0.067)	0.0364 (0.023)	0.0090 (0.030)	0.0119 (0.028)
Share of applicants cited by any reviewers		0.0221 (0.080)	−0.0153 (0.084)		0.0531 (0.037)	0.0217 (0.036)
Observations	2,056	2,056	2,056	2,063	2,063	2,063
R <sup>2</sup>	0.2554	0.2554	0.2805	0.4210	0.4219	0.5070
Year FEs	X	X	X	X	X	X
Study section FEs			X			X
Past performance, past grants, and demographics			X			X

Notes: Regression is at the study section meeting level. The dependent variable in columns 1–3 is the average quality, in citations to text matched publications, for funded grants associated with that meeting. The dependant variable in columns 4–6 is the average quality of all grant applications, including the unfunded ones. During my sample, there were seven meetings with no funded grants, which accounts for the difference in sample sizes between columns 1–3 and 4–6. The share of applicants cited by permanent reviewers is equal to the proportion of applicants in a meeting who have been cited by at least one permanent reviewer. Share cited by any reviewer is defined analogously. Columns 3 and 6 control for meeting level means of applicant demographics, education, past publications, past citations, and past R01 and other grants.

citation rates. The coefficient on share related is 0.0856, which implies that a 1 standard deviation increase in share related (0.182) increases the average quality of funded applications by  $0.182 \times 0.0856 = 0.0156$  or  $0.0156/0.144$  or 11 percent, significant at the 10 percent level. Column 2 controls for the share of applicants related to any reviewers, so that this specification more closely approximates the variation I use in my main individual-level regressions. Adding this control increases standard errors while decreasing the estimated coefficient slightly so that my results are no longer significant. This reflects the fact that, when aggregated to the meeting level, there is less variation in the share of applicants related to permanent members, conditional on total relatedness. Similarly, column 3 adds a variety of controls for the composition of applicants to a meeting: mean of gender, ethnicity, and education variables, as well as indicator variables for meeting-level means for number of citations (rounded to the nearest 100), and number of publications (rounded to the nearest 10). The coefficient does not change.

Although most of these estimates are not statistically significant, their relative stability across specifications suggests that I am picking up an effect of relatedness to permanent members. The magnitude of the effect I find in column 3 is such that a 1 standard deviation increase in the share of applicants related to permanent members (holding constant overall relatedness) increases the average quality of funded applicants by about 9 percent, which is a plausible magnitude. In online Appendix F.F4, I conduct an additional calculation using stronger distributional assumptions. In that exercise, I also find that proximity on net increases the quality of funded applications.

By contrast, columns 4–6 repeat this exercise with the average quality of all applications as the dependent variable. If it's the case that related reviewers have expertise, then they should increase the quality of funded applicants by choosing the best applicants to fund; there should be no effect of share related on quality of applicants in general. This is what I find. The magnitudes of the coefficients in columns 5 and 6 (having controlled for share related to any reviewers) are almost an order of magnitude smaller.

The analysis in Table 8 assumes that policymakers care about maximizing citations associated with NIH-funded research. An important disclaimer to note is that an efficiency calculation based on this measure of welfare may not always be appropriate. If, for instance, the NIH cares about promoting investigators from disadvantaged demographic or institutional backgrounds, then a policy that increases total citations may actually move the NIH further from the goal of encouraging diversity. Yet, while citations need not be the only welfare measure that the NIH cares about, there are compelling reasons why policymakers should take citation-based measures of quality in account when assessing the efficacy of grant review. My citation data extend beyond my sample period, allowing me to observe the quality of a publication as judged in the long run. This alleviates concerns that citations may underestimate the importance of groundbreaking projects that may not be well cited in the short run.

## V. Additional Robustness Checks

The online Appendix discusses a variety of robustness and specification checks.

Online Appendix tables A–C provide supporting details for the data plotted in online Appendix Figures B and C, and Figure 3.

Online Appendix Tables D–I examine the robustness of my results to alternative measures of grant quality: changing the time window I use to measure publications associated with a grant; restricting to authors with very rare names to improve the quality of publication matches; varying my text matching process; and restricting only to publications in which the PI has played a primary role.

For example, not receiving a grant may slow down a scientist's research by requiring her to spend additional time applying for funding. If this is the case, then a grant can directly impact the research quality of funded versus non-funded applicants even before any funding dollars are disbursed. To address this concern, I estimate an alternative specification focusing on publications on the same topic that were published one year *prior* to the grant-funding decision; these articles are likely to inform the grant proposal, but their quality cannot be affected by the actual funding decision. This is described in online Appendix Table F.

My results in Table 7 are based on residualized measures of quality. Residualizing citations allows me to identify whether proximate candidates have better information about an application's quality that cannot easily be gleaned from the primary investigator's CV. Online Appendix Table I shows that my results are robust to splitting my sample based on various non-residualized measures of quality: whether or not an application goes on to produce any citations to text-matched publications within the first year at all; those that produce publications cited at the ninety-fifth percentile of its field-year cohort versus not; and those that produce publications

cited at the ninety-ninth percentile of this distribution versus not. In all these cases, I find a stronger effect of proximity on higher quality applications.<sup>18</sup>

Online Appendix Table J presents my main estimates separately for new versus renewal grants and new versus established investigators.

Online Appendix Table K describes the results of an alternative estimation strategy that does not rely on the distinction between permanent and temporary reviewers. Instead of comparing outcomes for different reviewers who have been cited by different numbers of permanent reviewers, I use applicant fixed effects to examine outcomes for the same applicant over time, across meetings in which he or she is cited by different numbers of total reviewers.<sup>19</sup> I find largely similar results: related applicants are more likely to be funded, and the impact of relatedness is increasing in the quality of applications.

My next set of results describe broader tests of the logic of my empirical strategy. Online Appendix Table L, for instance, reports a different test of the validity of my quality measure. If my results were driven by changes in measured grant quality near the payline, I would find no effect of proximity for applications that share the same funding status. To test for this, I examine the impact of proximity on application scores for the subset of applications that are either all funded or all unfunded. In both of these subsamples, I find evidence that being proximate to a permanent member increases scores and increases the correlation between scores and quality. Because proximity cannot affect actual funding status in these subsamples, the effect I find cannot be driven by differences in how well quality is measured.

Another potential concern with my quality measure is that text matching may eliminate publications on topics different from that described in the grant application but which review committees care about. It is common for grant funding to subsidize research on future projects that may not be closely related to the original grant proposal; even though reviewers are instructed to restrict their judgements to the merits of the research proposed in the grant application, it is possible that they may attempt to infer the quality of an applicant's future research pipeline and that related reviewers might have more information about this. To test whether my results are robust to this possibility, I use data on grant acknowledgements to match grants to *all* subsequent publications, not just to the research that is on the same topic or which is published within a year of grant review. Because grant acknowledgment data exist only for funded grants, this specification can only examine whether proximity impacts the scores that funded applicants receive. In online Appendix Table M, I show that results using data on grant acknowledgments are largely similar.

Online Appendix Table N takes a different approach to addressing the potential concern that my relatedness measure is capturing unobserved aspects of an application's quality. If this were the case, then we might expect the impact of relatedness to appear similar to the impact of observed measures of quality, insofar as observed

<sup>18</sup> It is not possible to explore the impact of proximity on the funding outcomes of particularly low quality candidates according to unresidualized measures of quality. This is because there is significant bunching of applications at zero publications and citations.

<sup>19</sup> In my alternative specification using applicant fixed effects, the analogous regression equation is given by

$$Assessment_{icmt} = a_0 + a_1 Total\_Proximity_{icmt} + \mu X_{icmt} + \delta_i + \varepsilon_{icmt}.$$

and unobserved quality may be correlated. Online Appendix Table N shows that this is not the case: whereas the benefit of relatedness is found to be increasing in an application's quality, I find that the marginal impact of an applicant's past citations is the same across quality quartiles. This specification is discussed in more detail in online Appendix E.

Finally, online Appendix F estimates a model of committee decision-making in which bias and expertise parameters can be separately identified under a stricter set of distributional assumptions; online Appendix Table O presents estimates of bias and expertise, and online Appendix Table P makes efficiency calculations based on this model.

## VI. Conclusion

This paper examines the impacts of bias and expertise on the quality of grant evaluation at the NIH. My results show that the preferences of influential reviewers matters for the funding outcomes of otherwise similar grant applications: being evaluated by an intellectually related reviewer increases an applicant's chances of funding by 2.2 percent, or the equivalent of a one-quarter standard deviation increase in application quality itself. This figure suggests that committees have a hard time predicting quality and, by comparison, reviewer preferences have a relatively large effect on funding outcomes. The fact that I find a positive effect of relatedness shows that although scientists compete for scarce resources such as funding and scientific priority, they nonetheless favor applications in their own area, suggesting that they view the research of others as complements to their own.

At the same time, reviewers do not merely favor all related applicants. Rather, higher quality applicants benefit much more from being evaluated by reviewers in their own area and particularly low quality applicants can even be hurt. This finding strongly suggests that reviewers have better information about the quality of related candidates. In this setting, I find that, on net, the quality of funding grants improves when more potentially biased experts are included in the review process.

My results suggest that there may be scope for improving the quality of peer review. For example, current NIH policy prohibits reviewers from evaluating proposals from their own institution. In the past, the National Bureau of Economic Research was considered a single institution, meaning that economists often recused themselves from evaluating the work of other economists.<sup>20</sup> The findings in this paper demonstrate why such policies may entail efficiency trade-offs.

## REFERENCES

- Acemoglu, Daron.** 2009. *Introduction to Modern Economic Growth*. Princeton: Princeton University Press.
- Agarwal, Sumit, David Lucca, Amit Seru, and Francesco Trebbi.** 2014. "Inconsistent Regulators: Evidence from Banking." *Quarterly Journal of Economics* 129 (2): 889–938.
- Azoulay, Pierre, Joshua S. Graff Zivin, and Gustavo Manso.** 2011. "Incentives and Creativity: Evidence from the Academic Life Sciences." *RAND Journal of Economics* 42 (3): 527–54.

<sup>20</sup> Current conflict of interest policies apply to members of the same NBER program.



- Blanes i Vidal, Jordi, Mirko Draca, and Christian Fons-Rosen.** 2012. "Revolving Door Lobbyists." *American Economic Review* 102 (7): 3731–48.
- Boudreau, Kevin J., Eva C. Guinan, Karim R. Lakhani, and Christoph Riedl.** 2016. "Looking Across and Looking Beyond the Knowledge Frontier: Intellectual Distance, Novelty, and Resource Allocation in Science." *Management Science* 62 (10): 2765–83.
- Brogaard, Jonathan, Joseph Engelberg, and Christopher A. Parsons.** 2011. "Network Position and Productivity: Evidence from Journal Editor Rotations." <http://rady.ucsd.edu/faculty/directory/engelberg/pub/portfolios/editors.pdf>.
- Cockburn, Iain M., and Rebecca M. Henderson.** 2000. "Publicly Funded Science and the Productivity of the Pharmaceutical Industry." In *Innovation Policy and the Economy*, Vol. 1, edited by Adam B. Jaffe, Josh Lerner, and Scott Stern, 1–34. Chicago: University of Chicago Press.
- Crawford, Vincent P., and Joel Sobel.** 1982. "Strategic Information Transmission." *Econometrica* 50 (6): 1431–51.
- Fisman, Raymond, Daniel Paravisini, and Vikrant Vig.** 2017. "Cultural Proximity and Loan Outcomes." *American Economic Review* 107 (2): 457–92.
- Garfagnini, Umberto, Marco Ottaviani, and Peter Norman Sørensen.** 2014. "Accept or reject? An organizational perspective." *International Journal of Industrial Organization* 34: 66–74.
- Gerin, William, Christine H. Kapelewski, Jerome B. Itinger, and Tanya M. Spruill.** 2010. *Writing the NIH Grant Proposal: A Step-by-Step Guide*. Thousand Oaks, CA: SAGE Publications.
- Ginther, Donna K., Walter T. Schaffer, Joshua Schnell, Beth Masimore, Faye Liu, Laurel L. Haak, and Raynard Kington.** 2011. "Race, Ethnicity, and NIH Research Awards." *Science* 333 (6045): 1015–19.
- Griliches, Zvi.** 1991. "The Search for R&D Spillovers." National Bureau of Economic Research (NBER) Working Paper 3768.
- Hansen, Stephen, Michael McMahon, and Carlos Velasco Rivera.** 2014. "Preferences or private assessments on a monetary policy committee?" *Journal of Monetary Economics* 67: 16–32.
- Hegde, Deepak.** 2009. "Political Influence behind the Veil of Peer Review: An Analysis of Public Biomedical Research Funding in the United States." *Journal of Law and Economics* 52 (4): 665–90.
- Jacob, Brian A., and Lars Lefgren.** 2011. "The impact of research grant funding on scientific productivity." *Journal of Public Economics* 95 (9–10): 1168–77.
- Jones, Benjamin F.** 2010. "Age and Great Invention." *Review of Economics and Statistics* 92 (1): 1–14.
- Kerr, William R.** 2008. "Ethnic Scientific Communities and International Technology Diffusion." *Review of Economics and Statistics* 90 (3): 518–37.
- Kondo, Jiro E.** 2006. "Self-Regulation and Enforcement in Financial Markets: Evidence from Investor-Broker Disputes at the NASD." [https://www.chicagobooth.edu/research/workshops/finance/docs/kondo\\_jobmkt.pdf](https://www.chicagobooth.edu/research/workshops/finance/docs/kondo_jobmkt.pdf).
- Kremer, Michael, and Heidi Williams.** 2010. "Incentivizing Innovation: Adding to the Tool Kit." In *Innovation Policy and the Economy*, Vol. 10, edited by Josh Lerner and Scott Stern, 1–17. Chicago: University of Chicago Press.
- Laband, David N., and Michael J. Piette.** 1994. "Favoritism versus Search for Good Papers: Empirical Evidence Regarding the Behavior of Journal Editors." *Journal of Political Economy* 102 (1): 194–203.
- Li, Danielle.** 2017. "Expertise versus Bias in Evaluation: Evidence from the NIH: Dataset." *American Economic Journal: Applied Economics*. <https://doi.org/10.1257/app.2015.0421>.
- Li, Danielle, and Leila Agha.** 2015. "Big names or big ideas: Do peer-review panels select the best science proposals?" *Science* 348 (6233): 434–38.
- Li, Hao, Sherwin Rosen, and Wing Suen.** 2001. "Conflicts and Common Interests in Committees." *American Economic Review* 91 (5): 1478–97.
- Pearson, Helen.** 2003. "Competition in biology: It's a scoop!" *Nature* 426 (6964): 222–23.
- Sampat, Bhaven N., and Frank R. Lichtenberg.** 2011. "What Are the Respective Roles of the Public and Private Sectors in Pharmaceutical Innovation?" *Health Affairs* 30 (2): 332–39.
- Stephan, Paula E.** 2012. *How Economics Shapes Science*. Vol. 1. Cambridge: Harvard University Press.
- Zinovyeva, Natalia, and Manuel Bagues.** 2015. "The Role of Connections in Academic Promotions." *American Economic Journal: Applied Economics* 7 (2): 264–92.

## APPENDIX MATERIALS

### SUPPORTING DESCRIPTIVES

This section describes supplemental descriptives and analysis relevant to my proximity and quality measures.

Appendix Figure A plots the full distribution of previous citations and publications for permanent and temporary reviewers. Specifically, past citations are defined as the number of citations, to 2008, for publications published by the reviewer in the 5 years prior to the grant review meeting. Past publications simply count the number of such publications. I find that, overall, permanent and temporary reviewers have similar qualifications although the two distributions are statistically different. See Table 3 in the main text for additional details about these distributions.

Appendix Figure B plots this distribution separately by funded and unfunded candidates. Although funded applicants have higher quality in terms of both publications and citations, this figure clearly shows that there are still many unfunded applications that go on to generate many publications and citations.

Appendix Table A provides additional descriptives comparing the distribution of measured application quality for funded and unfunded applicants. The top panel compares citation-based quality measures for funded and unfunded applicants at the mean, and 1st, 5th, 10th, 25th, 50th, 75th, 90th, 95th, and 99th percentiles. Funded applicants have consistently higher measures of quality than unfunded applicants, although it is important to note that many unfunded applications have high *ex ante* quality as captured by this text-matching approach. The bottom panel of Table A displays the same comparison for publications. Here, there are fewer differences between funded and unfunded applicants. The mean number of publications is higher for funded applicants, and this is driven by publications at the tail.

Appendix Figure C plots distributions of application quality by proximity to permanent reviewers. The top two panels plot the distributions of citations and publications for applicants cited by exactly one reviewer. For each graph, the solid line shows the distribution of quality among applicants cited by one permanent reviewer and the dotted line does so for those cited by one temporary reviewer. These distributions are statistically indistinguishable: a Kolmogorov–Smirnov test cannot reject the null that these two distributions are equal. Similarly, the upper-right-hand panel shows the same, but with quality measured using the number of publications associated with a grant. The bottom two panels of Appendix Figure C repeat this exercise for applicants who have been cited by a total of two reviewers. In this case, there are now three possibilities: the applicant has been cited by two temporary reviewers, two permanent, or one of each. In all of these cases, the distribution of applicant quality is statistically similar.

Next, Appendix Table B examines the distribution of quality by relatedness to permanent members. This is the table analogue of Appendix Figure C. The setup of this table is similar to that of Table A: I compare the difference in means and various percentiles for applicants cited by the same total number of reviewers, but by differing numbers of per-

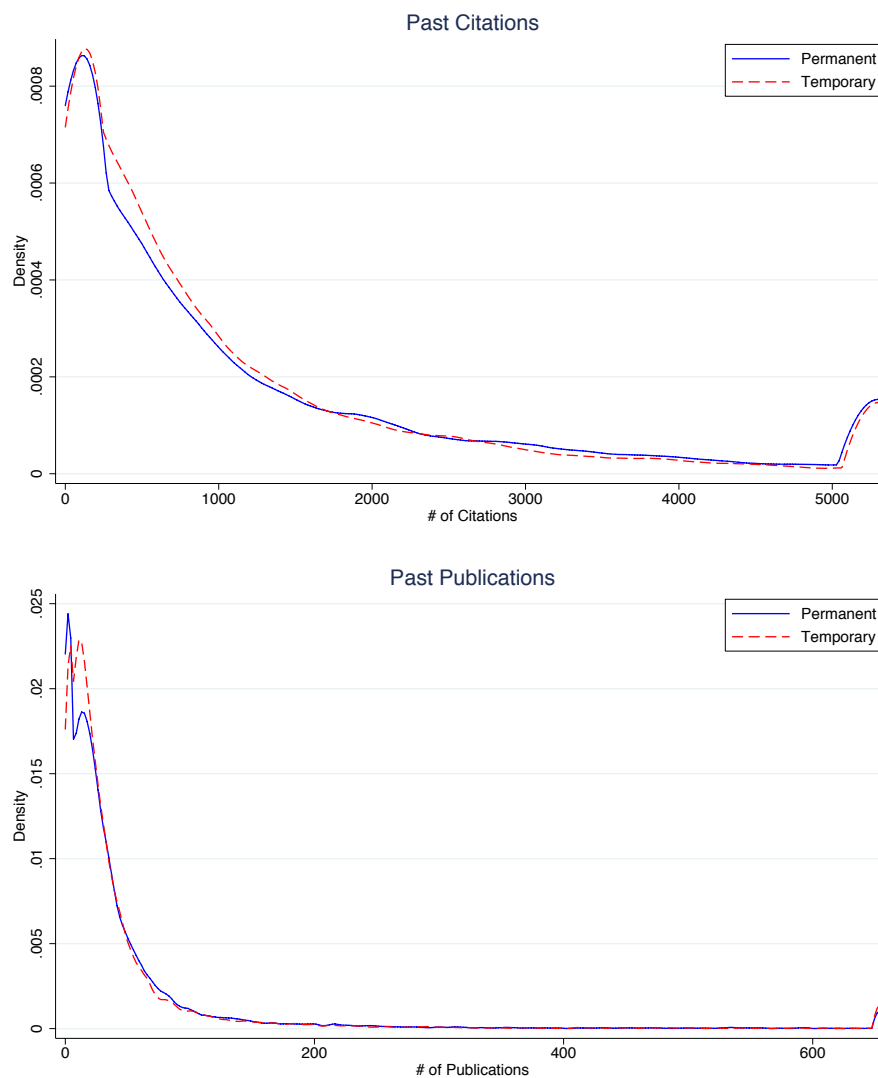
manent reviewers. This table considers the distribution of quality measures for applicants cited by 0, 1, and 2 reviewers total. This encompasses 73 percent of my sample.

Overall, these results show that, among applicants cited by the same total number of reviewers, there are very few significant differences in measured quality between applicants cited by more or fewer permanent reviewers. Among applicants cited by one reviewer, there are no significant differences in citation-based quality at the mean, 1st, 5th, 10th, 25th, 50th, or 75th percentiles. At the 90th, applicants cited by one permanent reviewer have 4 more citations than applicants cited by one temporary reviewer, with this difference being significant at the 5 percent level. At the 95th percentile, this difference grows to 7 citations (but shrinks in percentage terms). Finally, at the 99th percentile, applicants cited by permanent reviewers have 10 more citations, but this difference is no longer significant.

Appendix Table B also compares the distribution of citation outcomes among applicants cited by two reviewers in total. The clear pattern that emerges here is that applicants look broadly similar but that applicants cited by two temporary reviewers appear slightly weaker than applicants cited by either one of each or two permanent members. However, it also appears that applicants cited by one of each type of reviewer have stronger records than applicants cited by two permanent members. As a result, there is no systematic relationship between number of citing permanent reviewers and application quality.

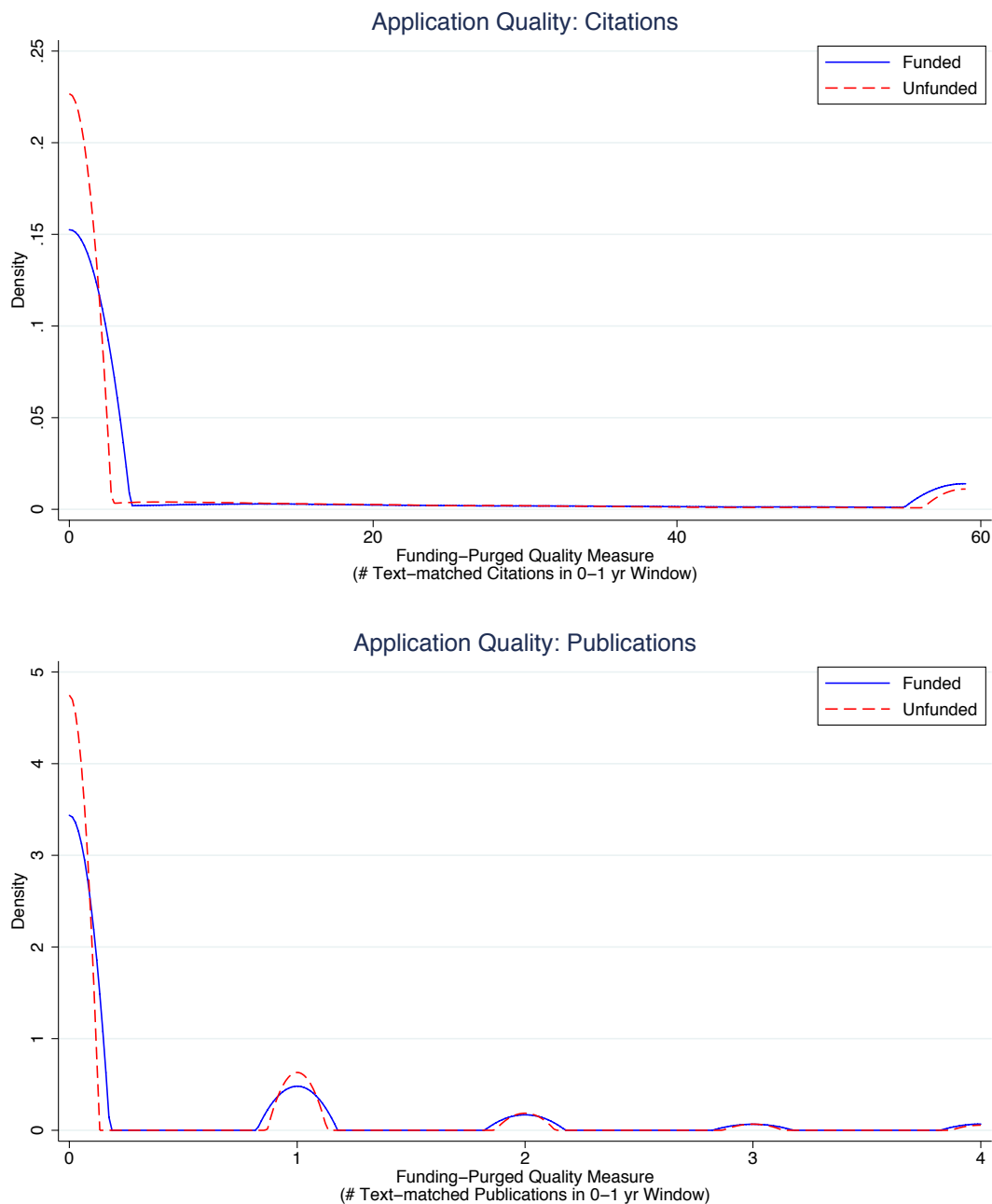
Appendix Table C provides the formal tests designed to accompany Figure 3, which shows that my measure of application quality is not impacted by a grant application's funding status. Column 1 reports the relationship between a grant's funding status and its measured quality, controlling for a linear score trend (where the impact of score is allowed to differ above and below the funding payline). Column 2 repeats this regression with finer controls for score, in this case, a quintic polynomial in score that is allowed to differ above and below the payline. In both cases, I find no significant relationship between funding and measured quality. Next, Columns 3 and 4 present IV evidence on the impact of funding on measured grant quality, where I instrument for funding with an indicator for a grant's score falling above the payline. Again I find no significant effect. Finally, Columns 5 and 6 present reduced form evidence on the relationship between falling above the payline and measured quality. I find no evidence of any association. Finally, it is worth noting that these regressions include meeting fixed effects. Within a meeting there is still variation in whether grants with the same score fall above the payline because different grants are subject to different paylines depending on what NIH Institute they are funded by. For instance, if the National Cancer Institute receives more competitive applications, than a cancer-related grant with a score of 70 may not be funded even though a diabetes-related grant with same score would be funded.

APPENDIX FIGURE A: DISTRIBUTION OF PAST CITATIONS AND PUBLICATIONS, PERMANENT VS. TEMPORARY REVIEWERS



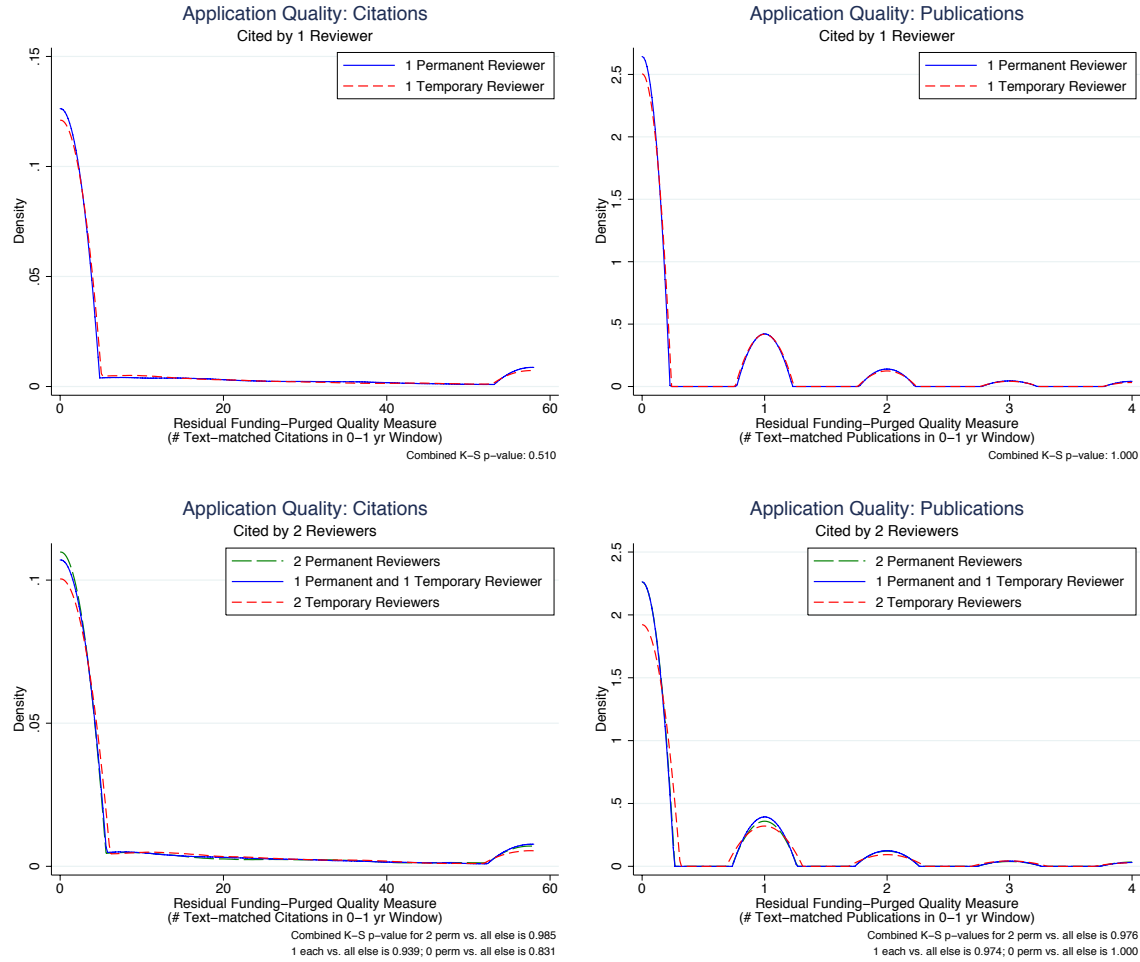
Note: Epanechnikov kernel. Publications count # of publications for which the reviewer was a first, second, or last author, published within 5 years of the relevant study section meeting. Citations count all citations to that set of publications, to 2008. Citations and Publications are top-coded at the 95th and 99th percentiles, respectively. This is done for legibility only; analyses use the full distribution of both variables. See Section II.B for additional details about how quality is constructed. A Kolmogorov-Smirnov test rejects that these two distributions are identical, with a p-value of 0.000.

APPENDIX FIGURE B: DISTRIBUTION OF APPLICATION QUALITY: FUNDED AND UNFUNDED GRANTS



Note: Epanechnikov kernel. Publications count # of text-matched publications within one year of grant review. Citations count all citations to that set of publications, to 2008. Citations and Publications are top-coded at the 95th and 99th percentiles, respectively. This is done for legibility only; analyses use the full distribution of both variables. See Section II.B and Appendix B for additional details about how quality is constructed. Kolmogorov-Smirnov tests reject that the distribution for unfunded grants is greater than for funded grants, for both publication and citation outcomes. The p-value for both tests is 0.000. It does not reject that funded grants do better on both dimensions. The p-value for both those tests is 1.000

APPENDIX FIGURE C: APPLICATION QUALITY CONDITIONAL ON TOTAL # OF PROXIMATE REVIEWERS



Note: Epanechnikov kernel. Publications count # of text-matched publications within one year of grant review. Citations count all citations to that set of publications, to 2008. Citations and Publications are top-coded at the 95th and 99th percentiles, respectively. This is done for legibility only; analyses use the full distribution of both variables. See Section II.B and Appendix B for additional details about how quality is constructed. Kolmogorov-Smirnov tests reject any differences between the distributions in each figure. For the bottom panels, K-S tests are performed pairwise: distribution of those cited by 2 permanent reviewers versus distribution for those cited by less than 2 permanent reviewers; distribution for those cited by one reviewer each vs. not, distribution for those cited by 2 temporary reviewers vs. not.



APPENDIX TABLE A: DISTRIBUTION OF APPLICATION QUALITY BY FUNDING

	<i>Mean</i>	<i>Percentiles</i>								
		1	5	10	25	50	75	90	95	99
<b>Application Quality: Citations</b>										
<i>Funded</i>	10.28	0	0	0	0	0	0	37	90	308
<i>Unfunded</i>	8.72	0	0	0	0	0	0	18	48	163
P-value	0.00	1.00	1.00	1.00	1.00	1.00	1.00	0.00	0.00	0.00
<b>Application Quality: Publications</b>										
<i>Funded</i>	3.28	0	0	0	0	0	0	1	2	3
<i>Unfunded</i>	2.63	0	0	0	0	0	0	1	2	4
P-value	0.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00

Notes: See notes to Table 1 for details about the sample. The quality of grant applications is measured as follows: # Publications refers to the number of research articles that the grant winner publishes in the year following the grant which share at least one salient word overlap between the grant project title and the publication title. # Citations refers to the total number of citations that accrue to this restricted set of publications, from the time of publication, to the end of my citation data in 2008. P-values for significance of percentile values are determined from a quantile regression of the quality outcome variable on an indicator variable for an application's funding status.

APPENDIX TABLE B: DISTRIBUTION OF APPLICATION QUALITY BY APPLICANT RELATEDNESS

	Mean	Percentiles								
		1	5	10	25	50	75	90	95	99
<b>Application Quality: Citations</b>										
<i>No Related Reviewer</i>										
0 Temp, 0 Perm	5.97	0	0	0	0	0	0	9	32	127
<i>1 Related Reviewer</i>										
1 Temp (N=7,049)	10.27	0	0	0	0	0	0	24	54	190
1 Perm (N=10,980)	11.42	0	0	0	0	0	0	28	61	200
P-value	0.16	1.00	1.00	1.00	1.00	1.00	1.00	0.03	0.07	0.61
<i>2 Related Reviewers</i>										
2 Temp (N=2,403)	8.72	0	0	0	0	0	0	24	50	151
1 Temp, 1 Perm (N=5,094)	11.39	0	0	0	0	0	0	28	62	213
2 Perm (N=4,841)	11.13	0	0	0	0	0	0	27	57	204
P-value*	0.08	1.00	1.00	1.00	1.00	1.00	1.00	0.67	0.24	0.04
<b>Application Quality: Publications</b>										
<i>No Related Reviewer</i>										
0 Temp, 0 Perm	0.20	0	0	0	0	0	0	1	1	3
<i>1 Related Reviewer</i>										
1 Temp (N=7,049)	0.32	0	0	0	0	0	0	1	2	4
1 Perm (N=10,980)	0.31	0	0	0	0	0	0	1	2	4
P-value	0.29	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
<i>2 Related Reviewers</i>										
2 Temp (N=2,403)	8.72	0	0	0	0	0	0	1	2	4
1 Temp, 1 Perm (N=5,094)	11.39	0	0	0	0	0	0	1	2	4
2 Perm (N=4,841)	11.13	0	0	0	0	0	0	1	2	4
P-value*	0.90	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00

Notes: See notes to Table 1 for details about the sample. The quality of grant applications is measured as follows: # Publications refers to the number of research articles that the grant winner publishes in the year following the grant which share at least one salient word overlap between the grant project title and the publication title. # Citations refers to the total number of citations that accrue to this restricted set of publications, from the time of publication, to the end of my citation data in 2008. P-values for differences in application quality for the subsample applicants cited by one reviewer only are computed from a quantile regression of the quality outcome variable on an indicator variable for an application's funding status. P-values for differences in application quality for the subsample of applicants cited by two reviewers is computed as follows: from a quantile regression of the number of proximate permanent members. The p-value reported is the p-value on the coefficient on this linear variable, which can take values 0, 1, and 2.

APPENDIX TABLE C: IS MEASURED QUALITY CONTAMINATED BY FUNDING?  
REGRESSION DISCONTINUITY IN SCORE

	Grant Application Quality (# of citations to text-matched publications within 1 year of grant review)					
	<i>OLS</i>		<i>IV</i>		<i>Reduced Form</i>	
	(1)	(2)	(3)	(4)	(5)	(6)
Awarded	-0.0024 (0.007)	0.0003 (0.008)	-0.0255 (0.020)	-1.5052 (3.407)		
1(Score Above Payline)					-0.0094 (0.008)	-0.0088 (0.018)
Observations	99,547	99,547	99,547	99,547	99,547	99,547
R-squared	0.0673	0.0674			0.0673	0.0674
Meeting FEs	X	X	X	X	X	X
Linear Score Trends	X		X		X	
Quintics in Score		X		X		X

Notes: Coefficients are reported from a regression of grant quality on an indicator for whether the grant was funded or whether it was scored above the payline. Columns 1 and 2 examines how measured quality changes once a grant is awarded, controlling for scores in various ways. Column 1 includes the following additional controls: score and score interacted with an indicator for being above the payline. Column 2 includes quintics in score, with each score variable also interacted with an indicator for being above the payline. Columns 3 and 4 instrument awarded with being above the payline, using the same set of controls as Columns 1 and 2. Finally, Columns 5 and 6 report the reduced form regressions of measured quality on an indicator for being above the payline, using the same set of score controls as Columns 1 and 2.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

## MEASURING GRANT APPLICATION QUALITY

This section describes my measure of application quality in more detail and provides additional robustness checks.

*B1. Match Process*

For each grant application, I have information on the name of the applicant, the title of the grant project and, in some cases, location identifiers for the applicant. I also have data from Thomson Reuters ISI Web of Science (ISI), containing information on publication titles, abstracts, and author names. To match these, I restrict to life science journal articles (e.g. excluding reviews, comments, etc.) in ISI with the same author name, published within 1 year of the study section meeting date. I have full name information in the NIH grant data, but publications are listed by last name and first and middle initial only. This results in some cases in which several authors can have the same initials (e.g. Smith, TA). In my baseline specifications, I exclude PIs with common names, defined as those last name, first initial, middle initial combinations shared by more than two individuals in PubMed. This amounts to about 7 percent of the sample being removed.

After removing common names and proceeding with an initial name and publication year match, I am left with a set of 16,134,500 possible grant-publication matches for 158,099 project titles and 3,274,225 possible publications. From this set, I compare the content of the grant project title with that of the publication title and publication abstract. I first remove a list of common stop words using the standard MySQL full test stop words list (available at <http://dev.mysql.com/doc/refman/5.5/en/fulltext-stopwords.html>). After doing so, the average grant project title has 4.87 semantic words (SD 1.10). The average publication title has 8.90 words (SD 3.38); the average abstract has 52.1 words (SD 36.9). 11.58 percent of potential pairs have at least one overlapping word between the grant and publication titles. 18.08 percent of potential pairs share a common semantic word. These comparisons are made from raw words only so that “mice” and “mouse” or “males” and “male” would not match.

In our main specifications, we say that a publication and grant application are text-matched to each other if they share at least 4 semantic words in either the publication title or abstract. Consider the following example from my data.

In 1999, the National Institute of Allergy and Infectious Disease funded grant number 1R01AI045057-01 from the applicant John C Boothroyd at Stanford University. The grant project title was titled “Genetics of Invasion and Egress in *Toxoplasma*.” This grant shows up in my raw data as follows:

Next, I search for life science publications by authors with the initials JC Boothroyd published in the first year after grant review (1999 and 2000). This yields 10 publications, of which I am excerpting five below for illustrative purposes:

Grant ID	Grant Year	Grant Title	PI Name
1R01AI045057-01	1999	<u>Genetics of Invasion</u> and <u>Egress</u> in <u>Toxoplasma</u>	Boothroyd, JC

Pub. ID	Pub. Year	Pub. Title	Pub. Abstract
000168366100029	2000	Ionophore-resistant mutants of <u>Toxoplasma gondii</u> reveal host cell permeabilization as an early event in <u>egress</u>	<u>Toxoplasma gondii</u> is an obligate intracellular pathogen within the phylum Apicomplexa. <u>Invasion</u> and <u>egress</u> by this protozoan parasite....
000165702100018	2000	Trans-spliced L30 ribosomal protein mRNA of <u>Trypanosoma brucei</u> is not subject to autogenous feedback control at the messenger RNA level	The regulation of gene expression in trypanosomes is poorly understood but it is clear that much of this regulation, particularly of developmentally controlled genes, is post-transcriptional....
000089249600007	2000	Lytic cycle of <u>Toxoplasma gondii</u>	<u>Toxoplasma gondii</u> is an obligate intracellular pathogen within the phylum Apicomplexa. This protozoan parasite is one of the most widespread, with a broad host range including many birds and mammals and a geographic range that is nearly worldwide....
0000167020000075	2000	<u>Toxoplasma gondii</u> homologue of plasmodium apical membrane antigen 1 is involved in <u>invasion</u> of host cells	Proteins with constitutive or transient localization on the surface of Apicomplexa parasites are of particular interest for their potential role in the <u>invasion</u> of host cells....
000079956900015	2000	A <u>Toxoplasma</u> lectin-like activity specific for sulfated polysaccharides is involved in host cell infection	<u>Toxoplasma gondii</u> is one of the most widespread parasites of humans and animals. The parasite has a remarkable ability to invade a broad range of cells....

The first publication clearly seems related to the subject of the grant. It has 2 overlapping words in the title and 4 overlapping words in the abstract (the 4th word, “invasion,” shows up later and is not reproduced here). My text matching algorithm will link this publication as related. The second publication does not seem like it has much overlap with the subject of the grant. My algorithm will not link this publication. The following three publications are more ambiguous. All of them are about “toxoplasma,” which is a key word in the grant project title. The third publication only has one overlapping word (“toxoplasma”) while the second has two overlapping words (“toxoplasma” and “invasion”), and the final has

one overlapping word (“toxoplasma”) and a close second (“invasion” vs. “invade”).

If we examine the list of publications actually acknowledged by the grant (this is available for funded applications only), this list includes 3 publications: the first, the third, and the fourth; the fifth publication, which looks similar in terms of word overlap, is not acknowledged. In the interest of being conservative, my main approach will match only the first publication.

## *B2. Robustness to alternative processes*

Given the ambiguity involved in the matching process, I explore the following forms of robustness to my primary text-matching process:

- 1) Appendix Table D: Varying criteria for uniqueness of names
- 2) Appendix Table E: Varying the threshold for word overlap used to associate publications with grants
- 3) Appendix Tables F and G: Varying the time window for publications to be associated with grants
- 4) Appendix Table H: Varying the prominence of the author’s contribution to a publication.
- 5) Appendix Table I: Compares results with alternative quality measures

Appendix Table D explores the robustness of my results to different restrictions on the types of applicant names that I include in my analysis. In my main specifications, I exclude all names with more than two individuals in PubMed who share the same last name, first and middle initial combination. The results in Appendix Table D show that my results do not change when I include all these names or when I am more restrictive, allowing only for unique last name and first and middle initial combinations.

Appendix Table E considers 8 different ways of changing threshold for how I choose whether a grant is matched to a publication. In my main specifications, I require that at least 4 semantic words be matched in either the publication title or abstract. As was discussed earlier, this may lead to cases in which publications on the same topic are missed (e.g., the third and fourth publications in the example table above.) Appendix Table B considers whether my results change when I apply different standards, both more and less stringent. Columns 1 through 4 detail results where text matching requires that  $X$  number of words overlap between the grant project title and the publication title *or* between the grant project title and the abstract, where  $X = 1, 2, 3$ , or 4. Because there are on average only 4.87 semantic words (SD 1.10) in the grant project title, I examine up 4 words maximum. Columns 5 through 8 repeat this exercise, but with a match defined as whether a grant project title shares  $X$  words with the publication title *and* the publication abstract (the main result is replicated in Column 5). The results show that, regardless of the exact threshold I use, my resulting estimates are similar: the impact of proximity increases with measured quality.



Appendix Tables F and G vary the time windows used to match grants to publications. Appendix Table F addresses concerns that funding may directly influence the number of citations produced by a grant by, for example, freeing up an investigator from future grant writing so that he can concentrate on research. Instead of including articles published after the grant is reviewed, Appendix Table F restricts my analysis to articles published one year before a grant is reviewed. These publications are highly likely to be based off research that existed before the grant was reviewed, but cannot have been influenced by the grant funds. Using this metric, I find nearly identical measures of bias and information. Appendix Table G addresses the opposite concern, that a one-year window after review may be insufficient to assess the quality of grant applications. Instead, I use a five year window following review and find that my results are both qualitatively and quantitatively similar. My estimates are very similar.

Finally, the next set of results explores the validity of my quality measures more broadly. The goal of my quality measures is to capture the quality of the research written into the grant application at the time of grant review. One possible concern with examining all publications by an author is that some of these publications may be ones for which the author made few substantive intellectual contributions, and which might not reflect his or her research program. Life science articles often have many authors and collaborators on a project may receive authorial credit for minor contributions such as sharing equipment or making figures. To address this, Appendix Table H restricts my match process to publications for which the grant applicant was the first, second, or last author. In the life sciences, contributions can be inferred from authorship position with earlier authors deserving more credit, and the last author being the primary investigator. Again, I find that the impact of proximity increases in application quality.

Finally, Appendix Table I shows that my results are robust to splitting my sample based on various non-residualized measures of quality: whether or not an application goes on to produce any citations to text-matched publications within the first year at all; those that produce publications cited at the 95th percentile of its field-year cohort vs. not; and those that produce publications cited at the 99th percentile of this distribution vs not. In all these cases, I find a stronger effect of proximity on higher quality applications.<sup>21</sup> For example, among applications that go on to produce publications in the top 99th percentile of their cohort's citation distribution, each additional reviewer increases their likelihood of funding by 2.1 percentage points, from a baseline funding rate for that group of 28.2 percent, or a 7.4 percent increase. This is similar to the magnitude I find for my top quartile applications from Table 7.

<sup>21</sup>It is not possible to explore the impact of proximity on the funding outcomes of particularly low quality candidates according to unresidualized measures of quality. This is because there is significant bunching of applications at zero publications and citations.

APPENDIX TABLE D: WHAT IS THE IMPACT OF PROXIMITY BY APPLICATION QUALITY?  
ROBUSTNESS TO ALTERNATIVE NAME-FREQUENCIES

	Quartiles of Residual Application Quality				
	<i>All</i>	<i>Bottom</i>	<i>Second</i>	<i>Third</i>	<i>Top</i>
	(1)	(2)	(3)	(4)	(5)
<i>Dependent Variable: 1(Score Above Payline)</i>					
# Proximate Permanent Reviewers	0.0047** (0.002)	-0.0018 (0.004)	-0.0007 (0.004)	0.0034 (0.005)	0.0177*** (0.005)
Observations	86,486	20,775	22,004	21,663	22,044
R-squared	0.0694	0.1535	0.1293	0.1186	0.1330
<i>Dependent Variable: Score</i>					
# Proximate Permanent Reviewers	0.1763* (0.099)	-0.1014 (0.181)	0.0478 (0.231)	0.2535 (0.329)	0.7705*** (0.230)
Observations	53,183	14,942	13,327	11,176	13,738
R-squared	0.1248	0.2081	0.2145	0.2597	0.2068
<i>Dependent Variable: 1(Scored at all)</i>					
# Proximate Permanent Reviewers	0.0014 (0.002)	-0.0005 (0.004)	-0.0097** (0.005)	0.0026 (0.007)	0.0104** (0.005)
Observations	86,486	20,775	22,004	21,663	22,044
R-squared	0.0911	0.1549	0.1463	0.1352	0.1399
Meeting FEs	X	X	X	X	X
# of Proximate Reviewer FEs	X	X	X	X	X

Notes: This table presents the same results as Table 7 but restricting to investigators who have a unique last name, first and middle initial combination in PubMed. Coefficients are reported from a regression of committee decisions (above payline, score, or scored at all) on # of proximate permanent reviewers, controlling for meeting level fixed effects and fixed effects for total proximity. Panel 1 regressions use the same specification as Column 2 in Table 5; Panel 2 uses the same specification as Column 5 in Table 5; the final panel uses the same specification as Column 8. Columns 2 through 5 split the sample based on quartiles of residual application quality. To calculate this, I regress application quality in citations on dummies for female, Hispanic, east Asian, south Asian, M.D., Ph.D., fixed effects for decile bins for both past publication and citations, and fixed effects for number of past R01 and other grants, and taking the residuals from this regression.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

APPENDIX TABLE E: WHAT IS THE IMPACT OF PROXIMITY BY APPLICATION QUALITY?  
ROBUSTNESS TO ALTERNATIVE TEXT-MATCHING WORD THRESHOLDS

Dependent Variable: 1(Score Above Payline)								
Quartiles of Residual Application Quality								
	<i>Bottom</i>	<i>Second</i>	<i>Third</i>	<i>Top</i>	<i>Bottom</i>	<i>Second</i>	<i>Third</i>	<i>Top</i>
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
<i>&gt;4 Overlapping Words Title OR Abstract</i>					<i>&gt;4 Overlapping Words Title AND Abstract</i>			
# Proximate Permanent Reviewers	-0.0006 (0.004)	0.0000 (0.004)	0.0039 (0.005)	0.0144*** (0.005)	0.0004 (0.004)	0.0015 (0.004)	0.0042 (0.005)	0.0083 (0.006)
<i>&gt;3 Overlapping Words Title OR Abstract</i>					<i>&gt;3 Overlapping Words Title AND Abstract</i>			
# Proximate Permanent Reviewers	-0.0001 (0.004)	0.0012 (0.004)	0.0085 (0.006)	0.0085* (0.004)	-0.0009 (0.004)	0.0036 (0.004)	0.0048 (0.005)	0.0120** (0.005)
<i>&gt;2 Overlapping Words Title OR Abstract</i>					<i>&gt;2 Overlapping Words Title AND Abstract</i>			
# Proximate Permanent Reviewers	-0.0011 (0.004)	-0.0007 (0.004)	0.0120** (0.006)	0.0085* (0.004)	-0.0019 (0.004)	0.0050 (0.004)	0.0082 (0.005)	0.0077* (0.004)
<i>&gt;1 Overlapping Words Title OR Abstract</i>					<i>&gt;1 Overlapping Words Title AND Abstract</i>			
# Proximate Permanent Reviewers	-0.0013 (0.004)	0.0079* (0.004)	0.0077 (0.006)	0.0033 (0.005)	0.0009 (0.004)	0.0010 (0.004)	0.0032 (0.006)	0.0100** (0.004)
Observations	23,854	24,368	23,458	21,878	23,218	24,101	23,602	22,637
Meeting FEs	X	X	X	X	X	X	X	X
# of Proximate Reviewer FEs	X	X	X	X	X	X	X	X

Notes: These regressions repeat Columns 2-5 from the first panel of Table 7. Coefficients are reported from a regression of 1(score above payline) on # of proximate permanent reviewers, controlling for meeting level fixed effects and fixed effects for total proximity. Columns 1 through 4 split the sample based on quartiles of residual application quality, where applicant quality is defined based on text matching that requires X words of overlap between the grant project title and the title of the publication or its abstract, where X = 1, 2, 3, 4. Columns 5-8 repeat this same exercise, except using text matching that requires word overlap with both the publication title and its abstract.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

APPENDIX TABLE F: WHAT IS THE IMPACT OF PROXIMITY BY APPLICATION QUALITY?  
GRANT QUALITY MEASURED FROM ARTICLES PUBLISHED 1 YEAR BEFORE GRANT  
REVIEW

	Quartiles of Residual Application Quality (Citations to Publications 1 Year Before Meeting)				
	<i>All</i>	<i>Bottom</i>	<i>Second</i>	<i>Third</i>	<i>Top</i>
	(1)	(2)	(3)	(4)	(5)
<i>Dependent Variable: 1(Score Above Payline)</i>					
# Proximate Permanent Reviewers	0.0050** (0.002)	-0.0018 (0.004)	0.0024 (0.004)	0.0044 (0.004)	0.0193*** (0.005)
Observations	93,558	22,367	23,925	23,753	23,513
R-squared	0.0688	0.1492	0.1266	0.1082	0.1279
<i>Dependent Variable: Score</i>					
# Proximate Permanent Reviewers	0.1641* (0.094)	-0.0992 (0.165)	0.0808 (0.202)	0.2468 (0.298)	0.8997*** (0.258)
Observations	57,613	16,336	15,037	12,327	13,913
R-squared	0.1224	0.1957	0.2043	0.2432	0.2104
<i>Dependent Variable: 1(Scored at all)</i>					
# Proximate Permanent Reviewers	0.0012 (0.002)	-0.0013 (0.003)	-0.0054 (0.004)	-0.0046 (0.006)	0.0073 (0.006)
Observations	93,558	22,367	23,925	23,753	23,513
R-squared	0.0899	0.1482	0.1395	0.1313	0.1383
Meeting FEs	X	X	X	X	X
# of Proximate Reviewer FEs	X	X	X	X	X

Notes: See notes to Table 1 for details about the sample. The key difference between this table and Table 7 is that application quality is measured using publications text-matched to the grant application project title in the one year before the grant meeting. Coefficients are reported from a regression of committee decisions (above payline, score, or scored at all) on # of proximate permanent reviewers, controlling for meeting level fixed effects and fixed effects for total proximity. Panel 1 regressions use the same specification as Column 2 in Table 5; Panel 2 uses the same specification as Column 5 in Table 5; the final panel uses the same specification as Column 8. Columns 2 through 5 split the sample based on quartiles of residual application quality. To calculate this, I regress application quality in citations on dummies for female, Hispanic, east Asian, south Asian, M.D., Ph.D., fixed effects for decile bins for both past publication and citations, and fixed effects for number of past R01 and other grants, and taking the residuals from this regression.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

APPENDIX TABLE G: WHAT IS THE IMPACT OF PROXIMITY BY APPLICATION QUALITY?  
GRANT QUALITY MEASURED FROM ARTICLES PUBLISHED 0-5 YEARS AFTER GRANT  
REVIEW

	Quartiles of Residual Application Quality (Citations to Publications up to 5 Years After Meeting)				
	<i>All</i>	<i>Bottom</i>	<i>Second</i>	<i>Third</i>	<i>Top</i>
	(1)	(2)	(3)	(4)	(5)
<i>Dependent Variable: 1(Score Above Payline)</i>					
# Proximate Permanent Reviewers	0.0050** (0.002)	-0.0021 (0.004)	-0.0010 (0.004)	0.0069 (0.005)	0.0172*** (0.005)
Observations	93,558	22,764	23,981	23,794	23,019
R-squared	0.0688	0.1495	0.1206	0.1101	0.1401
<i>Dependent Variable: Score</i>					
# Proximate Permanent Reviewers	0.1641* (0.094)	-0.1145 (0.168)	-0.1030 (0.216)	0.3783 (0.306)	0.5820** (0.234)
Observations	57,613	15,954	14,342	12,440	14,877
R-squared	0.1224	0.2032	0.2067	0.2382	0.2052
<i>Dependent Variable: 1(Scored at all)</i>					
# Proximate Permanent Reviewers	0.0012 (0.002)	-0.0013 (0.003)	-0.0066 (0.005)	0.0052 (0.006)	0.0095* (0.005)
Observations	93,558	22,764	23,981	23,794	23,019
R-squared	0.0899	0.1504	0.1403	0.1275	0.1426
Meeting FEs	X	X	X	X	X
# of Proximate Reviewer FEs	X	X	X	X	X

Notes: See notes to Table 1 for details about the sample. The key difference between this table and Table 7 is that application quality is measured using publications text-matched to the grant application project title up to five years after the grant meeting. Coefficients are reported from a regression of committee decisions (above payline, score, or scored at all) on # of proximate permanent reviewers, controlling for meeting level fixed effects and fixed effects for total proximity. Panel 1 regressions use the same specification as Column 2 in Table 5; Panel 2 uses the same specification as Column 5 in Table 5; the final panel uses the same specification as Column 8. Columns 2 through 5 split the sample based on quartiles of residual application quality. To calculate this, I regress application quality in citations on dummies for female, Hispanic, east Asian, south Asian, M.D., Ph.D., fixed effects for decile bins for both past publication and citations, and fixed effects for number of past R01 and other grants, and taking the residuals from this regression.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

APPENDIX TABLE H: WHAT IS THE IMPACT OF PROXIMITY BY APPLICATION QUALITY?  
GRANT QUALITY MEASURED FROM FIRST, SECOND, AND LAST AUTHORSHIP POSITION  
ARTICLES

	Quartiles of Residual Application Quality				
	<i>All</i>	<i>Bottom</i>	<i>Second</i>	<i>Third</i>	<i>Top</i>
	(1)	(2)	(3)	(4)	(5)
<i>Dependent Variable: 1(Score Above Payline)</i>					
# Proximate Permanent Reviewers, (based on citations by first, second, and last authors only)	0.0083*** (0.003)	0.0060 (0.005)	0.0031 (0.006)	-0.0008 (0.007)	0.0160** (0.006)
Observations	93,558	22,463	23,929	23,360	23,806
R-squared	0.0678	0.1504	0.1249	0.1102	0.1257
<i>Dependent Variable: Score</i>					
# Proximate Permanent Reviewers, (based on citations by first, second, and last authors only)	0.2423* (0.130)	0.1487 (0.227)	0.2073 (0.311)	-0.2040 (0.417)	0.5668* (0.311)
Observations	57,613	16,081	14,593	12,056	14,883
R-squared	0.1218	0.1999	0.2077	0.2433	0.1967
<i>Dependent Variable: 1(Scored at all)</i>					
# Proximate Permanent Reviewers, (based on citations by first, second, and last authors only)	0.0032 (0.003)	-0.0010 (0.004)	-0.0099* (0.006)	0.0060 (0.010)	0.0148** (0.006)
Observations	93,558	22,463	23,929	23,360	23,806
R-squared	0.0866	0.1480	0.1392	0.1293	0.1299
Meeting FEs	X	X	X	X	X
# of Proximate Reviewer FEs	X	X	X	X	X

Notes: See notes to Table 1 for details about the sample. The key difference between this table and Table 7 is that proximity is determined by citations made to the applicant's work, only by reviewers who were the first, second, or last authors on the citing publication. Coefficients are reported from a regression of committee decisions (above payline, score, or scored at all) on # of proximate permanent reviewers, controlling for meeting level fixed effects and fixed effects for total proximity. Panel 1 regressions use the same specification as Column 2 in Table 5; Panel 2 uses the same specification as Column 5 in Table 5; the final panel uses the same specification as Column 8. Columns 2 through 5 split the sample based on quartiles of residual application quality. To calculate this, I regress application quality in citations on dummies for female, Hispanic, east Asian, south Asian, M.D., Ph.D., fixed effects for decile bins for both past publication and citations, and fixed effects for number of past R01 and other grants, and taking the residuals from this regression.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1



APPENDIX TABLE I: WHAT IS THE IMPACT OF PROXIMITY BY APPLICATION QUALITY? ALTERNATIVE QUALITY MEASURES

	Alternative Measures of Application Quality					
	Any Citations?		Top 99th Percentile Pubs		Top 99th Percentile Pubs	
	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Yes</i>
	(1)	(2)	(3)	(4)	(3)	(4)
<i>Dependent Variable: 1(Score Above Payline)</i>						
# Proximate Permanent Reviewers	0.0036* (0.002)	0.0133** (0.006)	0.0035* (0.002)	0.0207*** (0.008)	0.0035* (0.002)	0.0207*** (0.008)
Observations	80,379	13,179	85,946	7,612	85,946	7,612
R-squared	0.0961	0.2292	0.0940	0.3148	0.0940	0.3148
<i>Dependent Variable: Score</i>						
# Proximate Permanent Reviewers	0.1249 (0.105)	0.5062** (0.249)	0.1327 (0.101)	0.4285 (0.353)	0.1327 (0.101)	0.4285 (0.353)
Observations	48,435	9,178	51,995	5,618	51,995	5,618
R-squared	0.1522	0.2798	0.1478	0.3677	0.1478	0.3677
<i>Dependent Variable: 1(Scored at All)</i>						
# Proximate Permanent Reviewers	0.0006 (0.002)	0.0041 (0.006)	0.0008 (0.002)	0.0084 (0.007)	0.0008 (0.002)	0.0084 (0.007)
Observations	80,379	13,179	85,946	7,612	85,946	7,612
R-squared	0.1360	0.2326	0.1342	0.3050	0.1342	0.3050
Meeting FEs	X	X	X	X	X	X
# of Proximate Reviewer FEs	X	X	X	X	X	X
Past Performance, Past Grants, and Demographics	X	X	X	X	X	X

Notes: See notes to Table 1 for details about the sample. Coefficients are reported from a regression of committee decisions (above payline, score, or scored at all) on # of proximate permanent reviewers, controlling for meeting level fixed effects and fixed effects for total proximity. Each column presents the main regression from Table 5 on a different sample based on publication outcomes related to the application. Columns 1 and 2 compare applications with any citations to text-matched publications within one year of grant review, versus those without. Columns 3 and 4 compare applications that then go on to produce a text-matched publication within one year of grant review, where that publication is cited at the top 95th percentile of all publications from the same year, based on citations in 2008 -- versus not. Columns 5 and 6 do the same for publications at the 99th percentile of citations. For all these regressions, I include controls for applicant characteristics: female, Hispanic, east Asian, south Asian, M.D., Ph.D., fixed effects for decile bins for both past publication and citations, and indicators for the number of past R01 and other NIH grants an applicant has won, as well as indicators for how many she has applied to.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

## ESTABLISHED VS. NEW INVESTIGATORS

Appendix Table J examines how my main estimates in Table 7 vary by whether or not a grant is new (Columns 1-5), and whether or not an investigator is new (Columns 6-10). My results indicate that the impact of proximity is similar for new and renewal grants when viewed in percentage point terms. However, new grants have a lower average probability of being funded, relative to renewal grants. Among new applications in the top quartile of quality, an additional related reviewer increases the application's likelihood of funding by 1.9 percentage points from a base of 17.7 percent, or a 10.1 percent increase. Among new applications in the bottom quartile of quality, proximity decreases an application's likelihood of funding by 0.8 percentage points, from a base of 21.1 percent, or a 4.8 percent decrease. This suggests that the informational advantage of related reviewers may be greater for new applications.

My results are noisier for entirely new investigators. While I cannot reject that the effect for new investigators is similar to my estimates for established investigators, I cannot reject that they are zero either. If it is the case that I find stronger effects of relatedness for established investigators, this may be because reviewers are familiar with the research agendas of established investigators. While most new investigators have a history of publications, these articles are often written in conjunction with a more senior scientist who funds the research. When reviewers cite publications by new investigators, they may be more familiar with the work of the senior scientist, rather than that of the new investigator herself, who is likely to have been a graduate student, postdoc, or unfunded junior academic at the time. As a result, proximity to new investigators may convey less information than proximity to established scientists.

APPENDIX TABLE J: IMPACT OF PROXIMITY, HETEROGENEITY BY GRANT AND APPLICANT TYPE

		Dependent Variable: 1(Score Above Payline)									
		Quartiles of Residual Application Quality					Quartiles of Residual Application Quality				
		All	Bottom	Second	Third	Top	All	Bottom	Second	Third	Top
		(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
		<i>Renewal Grants</i>					<i>Established Investigators</i>				
# Proximate Permanent Reviewers		0.0066* (0.004)	-0.0032 (0.009)	-0.0048 (0.009)	0.0151 (0.011)	0.0187* (0.010)	0.0052** (0.002)	-0.0009 (0.004)	0.0030 (0.004)	0.0025 (0.005)	0.0154*** (0.005)
Observations		27,782	6,530	6,992	7,155	7,105	84,463	20,166	21,444	21,445	21,408
R-squared		0.1385	0.3499	0.3152	0.2791	0.2977	0.0713	0.1608	0.1340	0.1169	0.1360
		<i>New Grants</i>					<i>New Investigators</i>				
# Proximate Permanent Reviewers		0.0030 (0.002)	-0.0083* (0.005)	0.0066 (0.005)	-0.0072 (0.006)	0.0189*** (0.006)	-0.0031 (0.009)	0.0044 (0.024)	-0.0223 (0.031)	0.0054 (0.041)	0.0021 (0.033)
Observations		65,776	15,935	16,661	16,639	16,541	9,095	2,271	2,282	2,181	2,361
R-squared		0.0644	0.1606	0.1443	0.1323	0.1558	0.2465	0.5362	0.5743	0.5400	0.5989
Meeting FEs		X	X	X	X	X	X	X	X	X	X
# of Proximate Reviewer FEs		X	X	X	X	X	X	X	X	X	X

Notes: See notes to Table 1 for details about the sample, and Table 7 for details about the regression specification. In all regressions, the dependent variable is an indicator for score being greater than the funding payline. Columns 1-5 split the sample based on whether a grant application is a renewal application (top panel) or whether it was for a new project (bottom panel). Column 1 examines the full sample of the relevant set of grants. Columns 2 through 5 split the sample based on quartiles of residual application quality. To calculate this, I regress application quality in citations on dummies for female, Hispanic, east Asian, south Asian, M.D., Ph.D., fixed effects for decile bins for both past publication and citations, and fixed effects for number of past R01 and other grants, and taking the residuals from this regression. Columns 6-10 repeat this exercise for the sample of established investigators (those who have had prior NIH funding as a primarily investigator of any kind) and new investigators, those who have not.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

## ALTERNATIVE IDENTIFICATION STRATEGY

In my main specification, I identify the effect of proximity to more influential reviewers (permanent vs. temporary). This approach relies on the assumption that controlling for the total number of reviewers who cite an applicant is an adequate control for unobserved differences in quality that may be correlated with whether an applicant is cited by a permanent reviewer. A different approach would be to use applicant fixed effects to control for quality, compare the funding outcomes of applications from the *same* applicant across meetings in which the applicant is cited by different total numbers of reviewers.<sup>22</sup>

The downside of this approach is that applicant fixed effects only control for time-invariant unobserved quality. If there are aspects of the quality of an applicant's proposal that are not controlled for with information on past publications and grant histories, then this may bias my results.

This second approach also captures a slightly different causal effect: the effect of being related to an additional reviewer, as opposed to being related to a more influential reviewer. The relative magnitudes of these effects are theoretically ambiguous: if only permanent reviewers have influence, then the effect of being related to a permanent reviewer (conditional on total proximity) will be larger than the effect of being related to an additional member (because that additional member may be temporary and thus, in this example, inconsequential). If, on the other hand, temporary members have as much influence as permanent ones, then the composition of related reviewers would not matter, but the number would.

Appendix Table K reports estimates from this alternative identification strategy. My results are similar. Overall, I find a significant impact of proximity on an applicant's likelihood of funding, the score that it receives, and the likelihood that it is scored at all. This effect is largely increasing with the quality of the application, although my estimates peak at the 3rd quartile rather than the 4th. For example, I find that each additional proximate reviewer—either temporary or permanent—increases an applicant's likelihood of funding by 0.61 percentage points or 2.9 percent. For the top quartile, this effect rises to 0.87 percentage points or 3.9 percent.

<sup>22</sup>In my alternative specification using applicant fixed effects, the analogous regression equation is given by:

$$\text{Assessment}_{icmt} = a_0 + a_1 \text{Total Proximity}_{icmt} + \mu X_{icmt} + \delta_i + \varepsilon_{icmt}.$$

APPENDIX TABLE K: WHAT IS THE IMPACT OF PROXIMITY BY APPLICATION QUALITY? APPLICANT FIXED EFFECTS

	Quartiles of Residual Application Quality				
	<i>All</i>	<i>Bottom</i>	<i>Second</i>	<i>Third</i>	<i>Top</i>
	(1)	(2)	(3)	(4)	(5)
<i>Dependent Variable: 1(Score Above Payline)</i>					
# of Proximate Reviewers	0.0061*** (0.001)	0.0039** (0.002)	0.0090*** (0.003)	0.0111*** (0.004)	0.0087** (0.004)
Observations	93,558	22,463	23,929	23,360	23,806
R-squared	0.4527	0.5296	0.6125	0.6379	0.6276
<i>Dependent Variable: Score</i>					
# of Proximate Reviewers	0.2679*** (0.058)	0.1994** (0.087)	0.2686 (0.189)	0.7719** (0.364)	0.4960** (0.196)
Observations	57,613	16,081	14,593	12,056	14,883
R-squared	0.5452	0.5993	0.6937	0.7306	0.7055
<i>Dependent Variable: 1(Score above payline)</i>					
# of Proximate Reviewers	0.0111*** (0.001)	0.0088*** (0.002)	0.0144*** (0.003)	0.0169*** (0.005)	0.0118*** (0.004)
Observations	93,558	22,463	23,929	23,360	23,806
R-squared	0.5636	0.5948	0.6780	0.7120	0.7135
Applicant FEs	X	X	X	X	X
Past Performance, Past Grants	X	X	X	X	X

Notes: See notes to Table 7 for details about the sample and variable construction. Coefficients are reported from a regression of committee decisions (above payline, score, or scored at all) on # of proximate reviewers (either permanent or temporary), controlling for applicant fixed effects and other applicant characteristics. Column 1 estimates this regression on the whole sample. Columns 2 through 5 split the sample based on quartiles of residual application quality. To calculate this, I regress application quality in citations on dummies for female, Hispanic, east Asian, south Asian, M.D., Ph.D., fixed effects for decile bins for both past publication and citations, and fixed effects for number of past R01 and other grants, and taking the residuals from this regression. For all these regressions, I include applicant fixed effects as well as controls for time varying applicant characteristics: fixed effects for decile bins for both past publication and citations, and indicators for the number of past R01 and other NIH grants an applicant has won, as well as indicators for how many she has applied to.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

## ADDITIONAL ROBUSTNESS CHECKS

This section provides broader tests of my empirical specifications.

A key identifying assumption is that my measure of quality is not affected by whether individuals are actually funded. Figure 3 provides my primary evidence that this is the case. Another test of my assumption that citations are not directly affected by funding is to ask whether I find bias in the review of inframarginal grants, that is grants that are well above or well below the funding margin. All grants in either group have the same funding status so any bias estimate cannot be attributed to differences in funding. Because I hold funding status constant, I can only assess the impact that related permanent members have on an applicant's score not on an applicant's funding status. Appendix Table L reports these results. The top panel reports the impact of proximity on scores, using funded grants only. The bottom panel does the same for unfunded grants. In both cases, I find an increasing effect of proximity with quality. The magnitudes are somewhat smaller than in my main regression; because these are subsamples, there is no reason to expect that the magnitude of the effect of proximity to be the same as it is for the entire sample.

Another potential concern is that committees may defy instructions and evaluate grant applications not on the basis of the specific research in the proposal, but on the quality of projects that reviewers suspect the grant funding may cross subsidize. In this case, by using text-matching to restrict my main quality measure to be based on articles that are closely related to the grant proposal topic, I am potentially missing other research that reviewers might be anticipating when they evaluate a grant proposal. To test whether this is the case, I use grant acknowledgement data recorded in the National Library of Medicine's PubMed database to match funded grants to all the articles that it produces, regardless of topic or date of publication. Because this process requires that a grant application actually be funded, I am only able to examine the impact of proximity on scores, rather than on funding likelihood or the likelihood of being scored. For the set of funded grants, Appendix Table M reruns my core regressions using citations to publications that explicitly acknowledge a grant as my measure of quality, and scores as my outcome measure. I find results that are consistent with my primary findings, though of a slightly smaller magnitude.<sup>23</sup>

Finally, despite the tests presented in Tables 4 and Appendix Figure C, there may still potentially be a correlation between relatedness to permanent members and unobserved aspects of applicant quality. If this were driving my results, one might expect the impact of relatedness I find in Table 7 to appear similar to the impact of observed measures of quality, insofar as observed and unobserved quality may be correlated. Appendix Table N shows that this is not the case by comparing the effect of past citations on funding probability, by application quality, to the effect of relatedness. In order to implement this comparison

<sup>23</sup>This analysis differs slightly from my main results using citations because general citations cannot be computed for publications in PubMed. A limited set of citations can, however, be computed using publications in PubMed Central (PMC). PMC contains a subset of life sciences publications made available for free. While this is not as comprehensive a universe as that of Web of Science, it contains, for recent years, all publications supported by NIH dollars. Undercounting of publications would, further, not bias my result as long as it does not vary systematically by whether an applicant is related to a permanent or to a temporary member.



directly, the specifications in Appendix Table N differ in a two ways from Table 7. First, recall that application quality quartiles used in Table 7 are defined in terms of residual quality, accounting for demographics, grant history, and past publications. In order to examine the impact of past citations on funding outcomes, I now residualize quality using all variables except publication history. Second, the specification used in Table 7 control for fixed effects in deciles of past publications and citations. In order to clearly assess the impact of past citations, Appendix Table N includes only controls for meeting fixed effects and fixed effects for the number of total related reviewers. The top panel of Appendix Table N shows that the marginal impact of an applicant's past citations is the same across quality quartiles. By contrast, the findings in the bottom panel show that the impact of relatedness is zero for the lowest quality applicants and increasing in quality thereafter.

APPENDIX TABLE L: WHAT IS THE IMPACT OF PROXIMITY BY APPLICATION QUALITY?  
INFRAMARGINAL GRANT APPLICATIONS

	Quartiles of Residual Application Quality				
	<i>All</i>	<i>Bottom</i>	<i>Second</i>	<i>Third</i>	<i>Top</i>
	(1)	(2)	(3)	(4)	(5)
<i>Funded Sample: Score</i>					
# Proximate Permanent Reviewers	0.1512* (0.078)	0.1445 (0.174)	0.0642 (0.193)	0.0346 (0.197)	0.3001 (0.226)
Observations	24,395	5,889	6,076	6,278	6,152
R-squared	0.1613	0.3622	0.3771	0.3773	0.3098
<i>Unfunded Sample: Score</i>					
# Proximate Permanent Reviewers	0.1012 (0.091)	-0.0222 (0.183)	0.3286 (0.261)	0.1348 (0.335)	0.1694 (0.175)
Observations	33,218	8,719	7,658	6,554	10,287
R-squared	0.1786	0.3090	0.3492	0.3785	0.2495
Meeting FEs	X	X	X	X	X
# of Proximate Reviewer FEs	X	X	X	X	X

Notes: Sample is the set of funded grants and of unfunded grants, treated separately. Coefficients are reported from a regression of score on # of proximate permanent reviewers, controlling for meeting level fixed effects and fixed effects for total proximity. Column 1 examines all applications. Columns 2 through 5 split the sample based on quartiles of residual application quality, within the set of funded or unfunded grants. To calculate this, I regress application quality in citations on dummies for female, Hispanic, east Asian, south Asian, M.D., Ph.D., fixed effects for decile bins for both past publication and citations, and fixed effects for number of past R01 and other grants, and taking the residuals from this regression.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

APPENDIX TABLE M: WHAT IS THE IMPACT OF PROXIMITY BY APPLICATION QUALITY?  
EXPLICIT GRANT ACKNOWLEDGEMENTS FOR THE SAMPLE OF FUNDED GRANTS

	Quartiles of Residual Application Quality (based on explicit grant acknowledgements)				
	<i>All</i>	<i>Bottom</i>	<i>Second</i>	<i>Third</i>	<i>Top</i>
	(1)	(2)	(3)	(4)	(5)
<i>Dependent Variable: Score</i>					
# Proximate Permanent Reviewers	0.1512* (0.078)	0.1445 (0.174)	0.0642 (0.193)	0.0346 (0.197)	0.3001 (0.226)
Observations	24,395	5,889	6,076	6,278	6,152
R-squared	0.1613	0.3622	0.3771	0.3773	0.3098
Meeting FEs	X	X	X	X	X
# of Proximate Reviewer FEs	X	X	X	X	X

Notes: Sample is funded grants only. Coefficients are reported from a regression of committee score on # of proximate reviewers, controlling for meeting level fixed effects and fixed effects for total proximity. Proximity to permanent reviewers is defined as the number of permanent reviewers who have cited the applicant's research in the 5 years prior to grant review. "Grant Application Quality" is defined as the number of citations up to 2008, for all publications that explicitly acknowledge funding from a grant, in the 100s unit. Columns 2 through 5 split the sample based on quartiles of residual application quality. To calculate this, I regress application quality in citations (using explicit acknowledgements) on dummies for female, Hispanic, east Asian, south Asian, M.D., Ph.D., fixed effects for decile bins for both past publication and citations, and fixed effects for number of past R01 and other grants, and taking the residuals from this regression. Dependent variable is the number of citations to publications that explicitly acknowledge the funded grant.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

APPENDIX TABLE N: ARE ESTIMATES LIKELY TO BE DRIVEN BY UNOBSERVED  
QUALITY?

IMPACT OF PAST CITATIONS ON FUNDING PROBABILITY

	Quartiles of Residual Application Quality				
	<i>All</i>	<i>Bottom</i>	<i>Second</i>	<i>Third</i>	<i>Top</i>
	(1)	(2)	(3)	(4)	(5)
<i>Dependent Variable: 1(Score Above Payline)</i>					
<b># Past Citations (100s)</b>	0.0018*** (0.000)	0.0017*** (0.000)	0.0019*** (0.000)	0.0021*** (0.000)	0.0016*** (0.000)
Observations	93,558	23,105	23,756	23,945	22,752
R-squared	0.0719	0.1363	0.1410	0.1176	0.1468
<b># Proximate Permanent Reviewers</b>	0.0050** (0.002)	-0.0006 (0.004)	0.0036 (0.004)	0.0082* (0.005)	0.0090** (0.004)
Observations	93,558	23,105	23,756	23,945	22,752
R-squared	0.0688	0.1335	0.1379	0.1150	0.1441
Meeting FEs	X	X	X	X	X
# of Proximate Reviewer FEs	X	X	X	X	X

Notes: See notes to Table 1 for details about the sample. In the first panel, coefficients are reported from a regression of committee decisions (above payline, score, or scored at all) on # of past citations, controlling for meeting level fixed effects and fixed effects for total proximity. Panel 2 does the same but with # of proximate permanent reviewers, controlling for meeting level fixed effects and fixed effects for total proximity. Columns 2 through 5 split the sample based on quartiles of residual application quality. To calculate this, I regress application quality in citations on dummies for female, Hispanic, east Asian, south Asian, M.D., Ph.D., and fixed effects for number of past R01 and other grants, and taking the residuals from this regression. For this specification, these residuals are calculated without information on past publications and citations.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

## THEORETICAL MODEL AND ALTERNATIVE ESTIMATION STRATEGY

In my final set of appendices, I present a model of expertise and bias in decision-making that can be directly estimated using a linear analysis in my data. The benefit of this approach, relative to the approach in the main body of the paper, is that it allows me to: 1) formally define expertise and bias and 2) show how these unobserved parameters and signals impact the equilibrium relationship between observable funding decisions, proximity to applicants, and realized grant quality; and 3) show how these parameters can be recovered from a linear regression of committee decisions on relatedness and quality.

*F1. Model*

A grant application has some true quality  $Q^*$  and, if approved, the committee receives a payoff of  $Q^*$ . If the application is rejected, the committee receives its outside option  $U$ , where  $U > E(Q^*)$ . Applications either work in the same area as the reviewer (“proximate,” given by  $P = 1$ ) or not ( $P = 0$ ). This model makes the simplifying assumption that committees can observe whether an application is related to a reviewer. I allow the application’s proximity to be unknown to the committee and show that all the same qualitative features of this model continue to hold. See the end of this section for a proof. Neither the committee nor the reviewer observes  $Q^*$ , but the reviewer observes a signal  $Q_P$  about  $Q^*$ . I assume that a related reviewer has greater expertise, meaning that  $Q_1$  gives a more precise signal than  $Q_0$ .<sup>24</sup>

After observing the signal, the reviewer sends a message to the committee about the application’s quality and the committee then decides whether to approve the grant. When determining what message to send, a reviewer considers his payoffs: for an unrelated application, this is identical to that of the committee, but for a related application, the reviewer now receives  $Q^* + B$  if the application is funded and  $U$  otherwise. The term  $B$  represents his bias. The timing is as follows:

- 1) An application with true quality  $Q^*$  is assigned to a reviewer.
- 2) The application’s type ( $P = 1$  or  $P = 0$ ) is determined and is publicly observed.
- 3) The reviewer observes the signal  $Q_P$ .
- 4) The reviewer sends a costless and unverifiable message  $M$  to the committee from some message space  $\mathbf{M}$ .
- 5) The committee, observing  $M$ , makes a decision  $D \in \{0, 1\}$  of whether to fund the grant.
- 6) True quality is revealed and the reviewer and committee both receive their payoffs.

<sup>24</sup>For simplicity, I assume that the signals  $Q_P$  are real numbers with continuous unconditional distributions such that  $E(Q^*|Q_P)$  is increasing in  $Q_P$ .

Proposition 1 describes the perfect Bayesian equilibria of this game.<sup>25</sup>

PROPOSITION 1: *The equilibria of the game is summarized by the following two cases:*

CASE 1:  $P = 0$ . *There exists a unique informative equilibrium in which*

- 1) *The reviewer reports a message  $Y$  if  $E(Q^*|Q_0) > U$  and  $N$  otherwise.*<sup>26</sup>
- 2) *The committee funds the grant if and only if the message is  $Y$ .*

CASE 2:  $P = 1$ . *There exists a level of bias  $B^* > 0$  such that for bias  $B \leq B^*$  there is a unique informative equilibrium such that*

- 1) *The reviewer reports a message  $Y$  if  $E(Q^*|Q_1) > U - B$  and  $N$  otherwise.*
- 2) *The committee funds the grant if and only if the message is  $Y$ .*

*When  $B > B^*$ , only uninformative equilibria exist and the grant is never funded.*

PROOF:

Proofs are included at the end of this section.

Proposition 1 says that when bias is sufficiently small, review committees are willing to take the advice of the reviewer because they value her expertise, in spite of the her bias. The committee's decision rule in the informative equilibria of this model is given by

$$\begin{aligned}
 (F1) \quad D = & \underbrace{\mathbb{I}(E(Q^*|Q_0) > U)}_{\text{baseline for unrelated}} + \underbrace{\mathbb{I}(U > E(Q^*|Q_1) > U - B)}_{\text{bias for proximate (+)}} P \\
 & + \underbrace{[\mathbb{I}(E(Q^*|Q_1) > U) - \mathbb{I}(E(Q^*|Q_0) > U)]}_{\text{additional information for proximate (+/-)}} P.
 \end{aligned}$$

The first term of Equation (F1) indicates that committees listen to advice about unrelated applications. The middle term represents the impact of bias on funding decisions. In particular, lower quality applications (those with  $U > E(Q^*|Q_1) > U - B$ ) will be funded if the applicant is related. The final term represents the impact of information.  $\mathbb{I}(E(Q^*|Q_1) > U)$  is the decision that an unbiased reviewer would make, given the lower variance signal of the proximate reviewer.  $\mathbb{I}(E(Q^*|Q_0) > U)$  is the decision she actually makes; the difference represents the change in funding outcomes that is due only to better information. Bias decreases the expected quality of funded applications while expertise increases it. The net effect of proximity on the quality of decisions is thus ambiguous.

Equation (F1) demonstrates why differences in funding likelihood among applicants with the same quality need not be due to bias. In particular, the difference in the expected

<sup>25</sup>There are always uninformative equilibria in which messages are meaningless and the grant is never funded. This proposition therefore focuses on informative equilibria, i.e. those in which the committee's decision depends on the reviewer's message. An informative equilibrium is unique if all other informative equilibria are payoff-equivalent for the parties.

<sup>26</sup>I assume there are at least two elements in the message space  $\mathbf{M}$  which, without loss, I call  $Y$  and  $N$ .

likelihood of funding between related and unrelated applications of the same quality is given by

$$\begin{aligned} E[D|Q^*, P = 1] - E[D|Q^*, P = 0] &= \Pr(U > E(Q^*|Q_1) > U - B) \\ &\quad + \Pr(E(Q^*|Q_1) > U) - \Pr(E(Q^*|Q_0) > U). \end{aligned}$$

This expression will be non zero even if reviewers are unbiased ( $B = 0$ ). This is because reviewers can more confidently attest to the quality of intellectually related applications, meaning that committees update more following a favorable review. Distinguishing between bias and information driven explanations is important because they have different implications for whether proximity enhances the quality of peer review.

PROOF OF PROPOSITION 1. — A perfect Bayesian equilibrium for this game is characterized by a message strategy for the reviewer, a set of beliefs about  $Q^*$  by the committee for each message, and a decision strategy for the committee. Having defined the equilibrium concept, I proceed with the proof of Proposition 1.

CASE 1. Suppose that the reviewer reports her exact posterior and the committee to believes it. In this case, the committee maximizes its utility by funding the proposal if and only if  $Q_0 > U$ . The reviewer has no incentive to deviate from this strategy because she is receiving her highest payoff as well.

Suppose, now, that there were another informative equilibrium. Each message  $M \in \mathbf{M}$  induces a probability of funding  $D(M)$ . Let the messages be ordered such that  $D(\mathbf{M}_1) \leq \dots \leq D(\mathbf{M}_K)$  where  $\mathbf{M}_i$  are the set of messages  $M_i$  that induce the same probability of funding  $D(M_i)$ . For reviewers of type  $E(Q^*|Q_0) > U$ , the reviewer strictly prefers that the grant be funded. She thus finds it optimal to send the message  $\mathbf{M}_K$  that maximizes the probability that the grant is funded. Call this set  $Y$ . For  $E(Q^*|Q^* + \varepsilon_0) < U$  the reviewer strictly prefer  $E(Q^*|Q_0) = U$ . Because the distribution of  $Q_P$  is assumed to be continuous on  $\mathbb{R}$  and such that  $E(Q^*|Q_P)$  is increasing in  $Q_P$ , this occurs with probability zero. Thus, with probability one, the space of possible messages is equivalent to  $\mathbf{M} = \{Y, N\}$ . For this equilibrium to be informative, it must be that  $D(N) < D(Y)$ . Given this, the committee's optimal reaction is to fund when  $M = Y$  and to reject otherwise.

If the we allow uninformative equilibria,  $D(\mathbf{M}_1) = \dots = D(\mathbf{M}_K)$  and any reviewer message is permissible. It must be that  $D(M_i) = 0$  for all  $M_i$  because the outside option  $U$  is assumed to be greater than the committee's prior on quality.

CASE 2. Now consider the case of a reviewer evaluating a related application. As in Case 1, the set of messages is equivalent, with probability one, to  $\mathbf{M} = \{Y, N\}$ . In this case, however, reviewers of type  $E(Q^*|Q_1) > U - B$  send  $M = Y$  and reviewers of type  $E(Q^*|Q_1) < U - B$  send  $M = N$ . The only reviewer who sends any other message is one for which  $E(Q^*|Q_1) = U - B$ .

Given this messaging strategy, a committee's expectation of  $Q^*$  given  $M = N$  is  $E(Q^*|E(Q^*|Q_1) < U - B)$

$U - B$ ). Since this is less than  $U$ , the grant goes unfunded. The committee's expectation of  $Q^*$  given  $M = Y$  is  $E(Q^*|E(Q^*|Q_1) < U - B)$ . When this is larger than  $U$ , the committee listens to the reviewer's recommendation and we can verify that  $D(Y) > D(N)$ . When  $E(Q^*|E(Q^*|Q^* + \varepsilon_1) < U - B) < U$ , the grant is never funded:  $D(Y) = D(N) = 0$ . In this case, only babbling equilibria exist.

If the we allow uninformative equilibria,  $D(\mathbf{M}_1) = \dots = D(\mathbf{M}_K)$  and any reviewer message is permissible. It must be that  $D(M_i) = 0$  for all  $M_i$  because the outside option  $U$  is assumed to be greater than the committee's prior on quality.

**Unobserved proximity:** Next, I consider a modification of Proposition 1 where the committee cannot observe whether the application is related to the reviewer.

**PROPOSITION A.2:** *Assume that  $p$  is the probability that an application is related to a reviewer. Then, for every  $p$ , there exists a level of bias,  $B^*$ , such that for  $B < B^*$  there is a unique informative equilibrium:*

*The reviewer reports a message  $Y$  if his posterior,  $E(Q^*|Q_1)$ , is greater than  $U - B$  and  $N$  otherwise.*

- 1) *An unrelated reviewer reports a message  $Y$  if his posterior,  $E(Q^*|Q_0)$ , is greater than  $U$  and  $N$  otherwise.*
- 2) *A related reviewer reports a message  $Y$  if his posterior,  $E(Q^*|Q_1)$ , is greater than  $U - B$  and  $N$  otherwise.*
- 3) *The committee funds the grant if and only if the message is  $Y$ .*

*For  $B \geq B^*$ , only uninformative equilibria exist and the grant is never funded.<sup>27</sup>*

**PROOF:**

In this case, the reviewer's messaging strategy remains the same as in Proposition 1: because reviewers themselves know whether they are proximate, they form, with probability one, strict preferences about whether an application should be funded. Proximate reviewers for which  $E(Q^*|Q_1) > U - B$  send  $M = Y$  and those for which  $E(Q^*|Q_1) < U - B$  send  $M = N$ . Similarly, unrelated reviewers of type  $E(Q^*|Q_0) > U$  send  $M = Y$  and unrelated reviewers of type  $E(Q^*|Q_0) < U$  send  $M = N$ .

The committee, however, does not observe the proximity and, as such, forms the following expectation of quality conditional on observing  $M = Y$ :

$$K [E(Q^*|E(Q^*|Q_0) > U)] + (1 - K) [E(Q^*|E(Q^*|Q_1) > U - B)]$$

The first term  $E(Q^*|E(Q^*|Q_0) > U)$  is the committee's expectation of quality if it knows that the  $M = Y$  message is sent by an unrelated reviewer. Similarly, the second

<sup>27</sup>Again, in all cases where an informative equilibrium exists, there also exist uninformative equilibria where the grant is never funded.



term  $E(Q^*|E(Q^*|Q_1) > U - B)$  is the committee's expectation of quality if it knows that the message is sent by a related reviewer. The term  $K$  is the probability that the committee believes a  $Y$  message comes from an unrelated reviewer, that is,  $K = E(P = 0|M = Y)$ . By Bayes' Rule, this is given by  $K = E(P = 0|M = Y) = \frac{E(P=0, M=Y)}{E(M=Y)}$ . The overall probability of a  $Y$  message is thus given by

$$E(M = Y) = (1 - p)(E(Q^*|Q_0) > U) + p(E(Q^*|Q_1) > U - B)$$

Similarly, the probability that the message is  $Y$  and the reviewer is unrelated is given by  $(1 - p)(E(Q^*|Q_0) > U)$ . As such, we have

$$K = \frac{(1 - p)(E(Q^*|Q_0) > U)}{(1 - p)(E(Q^*|Q_0) > U) + p(E(Q^*|Q_1) > U - B)}.$$

and for

$$K [E(Q^*|E(Q^*|Q^* + \varepsilon_0) > U)] + (1 - K) [E(Q^*|E(Q^*|Q^* + \varepsilon_1) > U - B)] > U$$

the committee funds the application. Again, we can verify that  $D(Y) > D(N)$ . For any fixed  $p$ , the threshold  $B^*$  can be defined to set this expression equality. There also exist uninformative equilibria where all grants are rejected. This term is less than  $U$ , then the grant is never funded:  $D(Y) = D(N) = 0$ . In this case, only babbling equilibria exist.

## F2. Statistical framework

The decision rule described by Equation (F1) in the theoretical model can be thought of as a data generating process for the funding decisions I observe. To make this more tractable, I make the following simplifying assumptions: for  $P = 0, 1$ , the reviewer's signal  $Q_P$  can be written as  $Q_P = Q^* + \varepsilon_P$  where  $\varepsilon_P \sim U[-a_P, a_P]$  and  $E(Q^*|Q_P)$  can be approximated by  $\lambda Q_P$  for some constant  $\lambda_R$ . Given this, an application's conditional likelihood of funding can be expressed as:

$$\begin{aligned} E[D|Q^*, P] &= \Pr(\lambda_0(Q^* + \varepsilon_0) > U) + \Pr(U > \lambda_1(Q^* + \varepsilon_1) > U - B)P \\ &\quad + [\Pr(\lambda_1(Q^* + \varepsilon_1) > U) - \Pr(\lambda_0(Q^* + \varepsilon_0) > U)]P \\ &= \frac{a_0 - U/\lambda_0 + Q^*}{2a_0} + \frac{B}{2a_1\lambda_1}P + \left[ \frac{a_1 - U/\lambda_1 + Q^*}{2a_1} - \frac{a_0 - U/\lambda_0 + Q^*}{2a_0} \right]P \\ &= \frac{1}{2} + \underbrace{\frac{1}{2a_0}}_{\text{Quality corr.}} Q^* + \underbrace{\frac{B}{2a_1\lambda_1}}_{\text{Bias term}} P + \underbrace{\left[ \frac{1}{2a_1} - \frac{1}{2a_0} \right]}_{\text{Add. corr. for proximate}} PQ^* \\ &\quad - \frac{U}{2a_0\lambda_0} + \left[ \frac{1}{2a_0\lambda_0} - \frac{1}{2a_1\lambda_1} \right]PU. \end{aligned} \tag{F2}$$

This distributional assumption has the benefit of allowing me to express the value of bias and expertise in a simple linear way: bias enters the level effect of relatedness while expertise enters the interaction effect. However, the assumption itself is restrictive: having a limited support of the error distribution means that if an application is extremely high (low) quality, the committee will choose to approve (reject) it regardless of what the reviewer says. As such, Equation (F2) is valid for candidates with quality such that  $Q^* + \varepsilon_P$  cannot be greater than  $U$  or less than  $U$  for all possible  $\varepsilon_P$ . Effectively, this restricts our analysis to grants that are at the margin of funding.

Given these caveats, Equation (F2) shows how I separately identify the role of bias and expertise. In particular, consider the regression analogue of Equation (F2):

$$(F3) \quad D = \alpha_0 + \alpha_1 Q^* + \alpha_2 P + \alpha_3 PQ^* + \alpha_4 U + \alpha_5 PU + X\beta + \epsilon,$$

where  $X$  includes other observable I can condition on.

Here,  $\alpha_2$ , the coefficient on proximity  $P$ , tests for bias: it is nonzero if and only if  $B \neq 0$ , where  $B$  is the bias parameter from the model. Second, the coefficient on  $PQ^*$  tests for expertise. To see this, notice that  $\alpha_1$  captures, for unrelated applicants, how responsive funding decisions are to increases in quality. In the model, this is determined by the precision of the reviewer's signal of quality for unrelated applications. The coefficient on  $PQ^*$ , meanwhile, captures the additional correlation between quality and funding for related applicants. A high coefficient on  $PQ$  means that a committee is more sensitive to increases in the quality of related applicants than to increases in the quality of unrelated applicants. In the model, this is determined by the difference in the precision of signals for related and unrelated applications.

The intuition for separately identifying bias and expertise is the following: if I find that related applications are more (or less) likely to be funded regardless of their quality, then this is a level effect of proximity that I attribute to bias in the NIH funding process. If I find that quality is more predictive of funding among related rather than unrelated applicants, then I conclude that study sections have better information about proposals from related applicants. I do not make any assumptions about the presence, extent, or direction of any potential biases nor do I assume that reviewers necessarily have better information about related applications. Rather, this statistical framework is designed to estimate this.<sup>28</sup>

Finally, the terms  $U$  and  $PU$  control for funding selectivity; for high cutoffs  $U$ , the correlation between funding and quality will be low even in the absence of bias or differential information because the marginal unfunded application is already very high-quality. The  $RU$  term, meanwhile, ensures that relationships are not credited for changing the correlation between funding and quality simply by lowering the threshold at which grants are funded.

Equation (F2) says that, as long as  $Q^*$  is perfectly observed, exogenous variation in

<sup>28</sup>These predictions hold when reviewers and committees are in an informative equilibrium. If the equilibrium were not informative, then advice from related reviewers would not be taken; I would find no effect of bias and a lower correlation between funding and quality for related applications. My results are not consistent with a non-informative equilibrium.

proximity is not needed to identify the presence of bias. This is because exogenous variation in proximity is necessarily only when aspects of an application's quality are potentially omitted; if quality were observed, one could directly control for any correlation between proximity and quality.

In practice, however, I do not observe an application's true quality  $Q^*$ . Instead, I observe a noisy signal  $Q = Q^* + v$ . Thus, instead of estimating Equation (F3), I estimate

$$(F4) \quad D = a_0 + a_1Q + a_2R + a_3RQ + a_4U + a_5RU + Xb + e.$$

Measurement error in quality can potentially pose problems for identification. Proposition 2 describes the conditions that must be met in order to consistently estimate bias from observed data.

**PROPOSITION 2:** *Given observed quality  $Q = Q^* + v$ , the bias parameter  $\alpha_2$  in Equation (F3) is consistently estimated by  $a_2$  in Equation (F4) when the following conditions are met:*

- 1)  $Cov(P, Q^*|U, PU, X) = 0$  and  $Cov(P^2, Q^*|U, PU, X) = 0$ ,
- 2)  $E(v|U, PU, X) = 0$ ,
- 3)  $Cov(v, P|U, PU, X) = 0$ .

**PROOF:**

:

Condition 1 requires that my measure of proximity,  $P$ , be uncorrelated, conditional on observables, with true application quality. If this were not the case, any mismeasurement in true quality  $Q^*$  would bias estimates of  $\alpha_2$  through the correlation between  $Q^*$  and  $P$ . Thus, in my study, exogenous variation in proximity is required only to deal with measurement error.

Condition 2 requires that measurement error be conditionally mean zero. This means that, after controlling for observable traits of the application or applicant, my quality measure cannot be systematically different from what committees themselves are trying to maximize. Otherwise, I may mistakenly conclude that committees are biased when they are actually prioritizing something I do not observe but which is not mean zero different from my quality measure.

Finally, Condition 3 requires that the extent of measurement error not depend, conditional on observables, on whether an applicant is related to a reviewer. This may not be satisfied if related applicants are more likely to be funded and funding itself affects my measure of quality.

**PROOF OF PROPOSITION 2.** — Measurement error in  $Q^*$  can potentially affect the estimation of  $\alpha_2$  in Equation (F3). The presence of  $U$ ,  $PU$ , and  $X$ , however, will not affect consistency; for simplicity, I rewrite both the regression suggested by the model and the actual estimating equation with these variables partialled out. The remaining variables should then be thought of as conditional on  $U$ ,  $PU$ , and  $X$

$$D = \alpha_0 + \alpha_1 Q^* + \alpha_2 P + \alpha_3 PQ^* + \epsilon$$

$$\begin{aligned} D &= a_0 + a_1 Q + a_2 P + a_3 PQ + e \\ &= a_0 + W + a_2 P + e, W = a_1 Q + a_3 PQ \end{aligned}$$

The coefficient  $a_2$  is given by:

$$a_2 = \frac{\text{Var}(W)\text{Cov}(D, P) - \text{Cov}(W, P)\text{Cov}(D, W)}{\text{Var}(W)\text{Var}(P) - \text{Cov}(W, P)^2}$$

Consider  $\text{Cov}(W, P)$ :

$$\begin{aligned} \text{Cov}(W, P) &= \text{Cov}(a_1(Q^* + v) + a_3P(Q^* + v), P) \\ &= a_1\text{Cov}(Q^*, P) + a_1\text{Cov}(v, P) + a_3\text{Cov}(PQ^*, P) + a_3\text{Cov}(Pv, P) \end{aligned}$$

Under the assumption that  $P$  and  $Q^*$  are conditionally independent, this yields:

$$\begin{aligned} \text{Cov}(W, P) &= a_3\text{Cov}(PQ^*, P) + a_3\text{Cov}(Pv, P) \\ &= a_3 [E(P^2Q^*) - E(PQ^*)E(P)] + a_3 [E(P^2v) - E(Pv)E(P)] \\ &= a_3 [E(P^2)E(Q^*) - E(P)^2E(Q^*)] + a_3 [E(P^2)E(v) - E(P)^2E(v)] \\ &= a_3 [E(P^2)0 - E(P)^20] + a_3 [E(P^2)0 - E(P)^20] \\ &= 0 \end{aligned}$$

With this simplification, the expression for the estimated coefficient on  $a_2$  becomes:

$$\begin{aligned} a_2 &= \frac{\text{Var}(W)\text{Cov}(D, P) - \text{Cov}(W, P)\text{Cov}(D, W)}{\text{Var}(W)\text{Var}(P) - \text{Cov}(W, P)^2} \\ &= \frac{\text{Var}(W)\text{Cov}(D, P)}{\text{Var}(W)\text{Var}(P)} \\ &= \frac{\text{Cov}(D, P)}{\text{Var}(P)} \\ &= \frac{\text{Cov}(\alpha_0 + \alpha_1 Q^* + \alpha_2 P + \alpha_3 PQ^* + \epsilon, P)}{\text{Var}(P)} \\ &= \frac{\alpha_2 \text{Var}(P) + \alpha_3 \text{Cov}(PQ^*, P)}{\text{Var}(P)} \\ &= \frac{\alpha_2 \text{Var}(P) + \alpha_3 [E(P^2)E(Q^*) - E(P)^2E(Q^*)]}{\text{Var}(P)} \\ &= \alpha_2 \end{aligned}$$

## F3. Empirical Estimates

Appendix Table O estimates the regression equation suggested by Equation (F4). Specifically,

$$\begin{aligned}
 \text{Assessment}_{icmt} = & a_0 + a_1 \text{Proximity to Permanent}_{icmt} \\
 & + a_2 \text{Proximate to Permanent}_{icmt} \times \text{Quality}_{icmt} \\
 & + a_3 \text{Quality}_{icmt} + a_4 \text{Total Proximity}_{icmt} \\
 & + a_5 \text{Total Proximity}_{icmt} \times \text{Quality}_{icmt} \\
 & + \mu X_{icmt} + \delta_{cmt} + \varepsilon_{icmt}.
 \end{aligned}
 \tag{F5}$$

I am interested in the coefficients  $a_1$  and  $a_2$ . Proximity to Permanent<sub>icmt</sub> is defined as the number of permanent reviewers that cite an applicant's prior work.  $a_1$  captures the effect of proximity on funding that is attributable to bias: does being cited by permanent reviewers, conditional on total proximity, affect an applicant's likelihood of being funded for reasons unrelated to quality? Bias is identified as the change in the *level* probability that a proximate applicant is funded. Meanwhile, Proximate to Permanent<sub>icmt</sub>  $\times$  Quality<sub>icmt</sub> is the interaction of an application's quality with an indicator for whether an applicant has been cited by a permanent reviewer. The coefficient  $a_2$  captures the role of expertise: it asks whether there is a steeper *slope* in the relationship between quality and funding for applicants with intellectual ties to more influential reviewers.

The remaining variables in Equation (F5) control for potentially contaminating variation. I control for the level of effect of application quality, total proximity to all reviewers, as well as the interaction between these two terms. Controlling for these terms means that the coefficient of interest  $a_1$  and  $a_2$  are estimated from applicants who have been cited by the same total number of reviewers, but who differ in their ties to permanent reviewers. I also control for a variety of past publication and demographic characteristics,  $X_{icmt}$ , described in Section III.

Finally, the model in Appendix F that motivates Equation (F5) also requires that I include controls for the degree of selectivity in a committee. When committees a very small percentage of applicants, the correlation between funding and quality will be low even in the absence of bias or differential information because the marginal unfunded application is already very high-quality. In my empirical implementation, I proxy for selectivity using the percentile pay line of the committee and include a level control for pay line (this is absorbed in the meeting fixed effect). I also control for the interaction of proximity and the payline. This ensures that proximity is not credited for changing the correlation between funding and quality simply by lowering the threshold at which grants are funded. My results are not affected by either the inclusion or exclusion of these variables.

Appendix Table O reports my estimates of Equation (F5), decomposing the effects of bias and expertise. Column 2 reports estimates of the coefficients from Equation (F5) for funding status. The positive and significant coefficients on the level effect of proximity

(0.0068) indicates that reviewers are biased in favor of applicants and the positive and significant coefficients on the interaction of proximity with quality (0.076) indicate that reviewers also have more expertise about related applications. Reviewers, however, also do a better job of discerning quality of related applicants. Consider a 1 standard deviation (51 citations) increase in the quality of a grant application: for an applicant cited by a single permanent reviewer, my estimates imply that this change would increase her chances of funding by  $(0.0136 + 0.0176 - 0.0005) * 0.51 * 100 = 1.6$  percentage points or  $1.6/21.4=7.5$  percent. If, instead, this applicant has been cited by a single temporary reviewer, the same increase in quality would only increase her chances of funding by  $(0.0136 - 0.0005) * 0.51 * 100 = 0.7$  percentage points or 3.3 percent. Committees are twice as responsive to changes in the quality of applications in the subject area of permanent members.

APPENDIX TABLE O: WHAT IS THE EFFECT OF PROXIMITY? LINEAR SPECIFICATION

	1(Score is above the payline)		Score		1(Scored at all)	
	Mean = 0.214, SD = 0.410		Mean = 71.18, SD = 18.75		Mean = 0.640, SD = 0.480	
	(1)	(2)	(3)	(4)	(5)	(6)
Proximity to Permanent Reviewers	0.0072*** (0.002)	0.0068*** (0.002)	0.2736*** (0.094)	0.2590*** (0.095)	0.0047** (0.002)	0.0043** (0.002)
Proximate to Permanent Reviewers × Grant Application Quality		0.0176** (0.008)		0.2739 (0.325)		0.0162* (0.009)
Grant Application Quality		0.0136** (0.006)		0.5568** (0.261)		0.0305*** (0.008)
Total Proximity X Grant Application Quality		-0.0005 (0.001)		-0.0043 (0.049)		-0.0036*** (0.001)
Observations	93,558	93,558	57,613	57,613	93,558	93,558
R-squared	0.0935	0.0949	0.1426	0.1431	0.1312	0.1322
Meeting FEs	X	X	X	X	X	X
# of Proximate Reviewer FEs	X	X	X	X	X	X
Past Performance, Past Grants, and Demographics	X	X	X	X	X	X

Notes: See notes to Table 1 for details about the sample. Coefficients are reported from a regression of committee decisions (above payline, score, or scored at all) on relatedness and quality measures, controlling for meeting level fixed effects. Proximity to permanent reviewers is defined as the number of permanent reviewers who have cited the applicant's research in the 5 years prior to grant review. "Grant Application Quality" is defined as the number of citations up to 2008, for all publications that are text-matched to the grant application within 1 year of grant review, in the 100s unit. "Past Performance, Past Grants, and Demographics" include indicators for sex and whether an applicant's name is Hispanic, East Asian, or South Asian, indicator variables for deciles of an applicant's total number of citations and publications over the past 5 years, indicators for whether an applicant has an M.D. and/or a Ph.D., and indicators for the number of past R01 and other NIH grants an applicant has won, as well as indicators for how many she has applied to.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

F4. *Efficiency*

Finally, Equation (F5) allows me to construct counterfactual funding decisions, made in the absence of proximity that would have been obtained in the absence of relationships. Specifically, I define

$$\begin{aligned} \text{Funding}_{icmt}^{\text{Benchmark}} &= \text{Funding}_{icmt} \text{ (actual funding)} \\ \text{Funding}_{icmt}^{\text{No Proximity}} &= \text{Funding}_{icmt} - \hat{a}_1 \text{Total Proximate Permanent}_{icmt} \\ &\quad - \hat{a}_2 \text{Quality}_{icmt} \times \text{Proximate to Permanent}_{icmt}, \end{aligned}$$

where  $\hat{a}_1$  and  $\hat{a}_2$  are estimated from Equation (F5).<sup>29</sup> The counterfactual funding decision represents what the committee would have chosen had applicants related to permanent members been treated as if they were unrelated.

I summarize the effect of relationships by comparing the quality of the proposals that would have been funded had relationships not been taken into account with the quality of those that actually are funded. Specifically, I consider all applications that are funded and sum up the number of publications and citations that accrue to this portfolio. This is my benchmark measure of the quality of NIH peer review. I then simulate what applications would have been funded had relationships not been taken into account. To do this, I fix the total number of proposals that are funded in each committee meeting but reorder applications by their counterfactual funding probabilities. I sum up the number of publications and citations that accrue to this new portfolio of funded grants. The difference in the quality of the benchmark and counterfactual portfolio provides a concrete, summary measure of the effect of relationships on the quality of research that the NIH supports.

Appendix Table P estimates the effect of relationships on the quality of research that the NIH supports. In effect, I ask what the NIH portfolio of funded grants would have been had committees treated applicants who are related to permanent members as if they were not, holding all else fixed. In my sample, I observe 93,558 applications, 24,404 of which are funded. Using this strategy, I find that 2,500, or 2.7 percent, of these applications change funding status under the counterfactual.

On average, working in the same area as influential reviewers helps an applicant obtain funding; ignoring this intellectual connection would decrease the number of proximate applicants who are funded by 3.0 percent. The quality of applications funded when intellectual proximity is taken into account, however, is higher. The overall portfolio of funded grants under the counterfactual produces two to three percent fewer citations, publications, and high-impact publications. To take account of the fact that some grants are funded and others are not, I use my standard funding-purged measure of grant application quality—text-matched publications within one year of grant review, and citations to those publications—as the measure of grant output used for this analysis. This has the benefit

<sup>29</sup>Even though  $\text{Funding}_{icmt}^{\text{No Relationship}}$  is constructed using estimates from Equation (F5), it does not rely on the model to interpret those coefficients.



of allowing me to compare the benchmark NIH portfolio with counterfactual results, holding constant the effect of actual funding status. However, a downside of this approach is that the stringent matching requirement will undercount the total number of publications (and therefore citations) associated with these grants. This exercise should thus be used to compare the percentage difference between the benchmark and counterfactual no-proximity cases, rather than to discern the level of NIH output.

APPENDIX TABLE P: WHAT IS THE EFFECT OF PROXIMITY ON THE AGGREGATE QUALITY OF NIH FUNDED GRANTS?

	Benchmark	No Proximity
Number of Funded Grants	24,404	24,404
Number of Grants that Change Funding Status	2,500	2,500
Total # Citations (% change relative to benchmark)	584,124	566,284 (3.05)
Total # Publications (% change relative to benchmark)	11,149	10,851 (2.67)
Total # in Top 99% of Citations (% change relative to benchmark)	590	572 (3.05)
Total # in Top 90% of Citations (% change relative to benchmark)	10,239	9,925 (3.07)
Total # Related Applicants Funded (% change relative to benchmark)	18,666	18,113 (2.96)

Notes: Benchmark refers to characteristics of grants ordered according to their predicted probability of funding, using the main regression of funding status on proximity and grant application quality. "Benchmark" figures are the grant quality measures for a grants that would be funded if we used the predicted ordering from the regression of funding likelihood on relatedness and quality estimated in Appendix Table L. "No proximity" refers to the predicted ordering of grants under the same regression, but under the assumption that relatedness to permanent members and relatedness to permanent members interacted with quality do not matter (their coefficients are set to zero). To take account of the fact that some grants are funded and others are not, we use our standard funding-purged measure of grant application quality: text-matched publications within one year of grant review, and citations to those publications. The number of projects that are funded is kept constant within meeting. See text for details.