Choose Your Moments: NIH Peer Review and Scientific Risk Taking

Richard T. Carson Joshua Graff Zivin Jeffrey G. Shrader^{*}

August 19, 2024

Abstract

Science funding agencies are often criticized for being too conservative. One explanation is that agencies typically base decisions on a simple average of peer review scores. Using a discrete choice experiment conducted with a large sample of biomedical researchers, we find that scientists prefer to fund projects with more reviewer dissensus. In contrast to funding allocation rules that focus primarily on the first moment of the distribution of reviewer scores, they also value the second moment. Scientists with the greatest domain expertise are particularly enthusiastic about dissensus. Using scientists' preferences changes funding decisions for projects worth billions of dollars annually. (JEL: O31,O32,O38)

^{*}Carson: Department of Economics, UC San Diego (email: rcarson@ucsd.edu); Graff Zivin, Department of Economics and School of Global Policy and Strategy, UC San Diego (email: jgraffzivin@ucsd.edu); Shrader: School of International and Public Affairs, Columbia University (corresponding author, email: jgs2103@columbia.edu). This study was approved by the IRB at the University of California, San Diego. We acknowledge the financial support of the National Science Foundation through its SciSIP Program (Award SBE-1561257). We are grateful for comments from Peter Muennig, Matthew Neidell, and Bhaven Sampat, as well as from seminar participants at Columbia University, the California Institute of Technology, and the NBER Summer Institute. Special thanks to Kyle Myers and Pierre Azoulay for their help with NIH grant data as well as their valuable comments. Stephanie Khoury and Tarikua Erda provided excellent research assistance. We particularly want to acknowledge the late Jordan Louviere, who shared sage advice on experimental design issues before his untimely death. All errors are our own.

1 Introduction

Fundamental scientific knowledge and the technologies built on it significantly contribute to aggregate income and economic growth (Nelson and Phelps 1966; Lucas 1988; Romer 1990; Aghion and Howitt 1992; Mokyr 1992). The public-good nature of basic scientific discovery implies that the government should play a prominent role in its funding, which should, in turn, catalyze private-sector investments in applied science (Arrow 1972; Nelson 1959; Bush 2020). Indeed, the U.S. federal government invested more than \$95 billion into science funding in 2021 alone (National Science Foundation (NSF) 2023), with the vast majority allocated based on some form of peer review process.¹ Peer review of research proposals is likewise the cornerstone of governmental research allocation decisions across the globe (Whitley and Gläser 2007), as well as grant awards from science-based philanthropic organizations² and firms' internal R&D decisions (Miller 1995). Peer review of *research proposals* is quite distinct from the *journal* peer review process with which we are all familiar. Most notably, reviewers are tasked with evaluating the potential success of early stage ideas rather than the quality and importance of a late stage one.

Despite the ubiquity of peer review for scientific grants and proposals, previous research has left open important questions about the best way to transform the outputs from peer review into decisions about the allocation of scarce resources (Franzoni and Stephan 2022). This is especially true if the goal is to produce novel or transformative science (Sen 2014; Boudreau et al. 2016), with science agencies having long been criticized for being too conservative in their research funding decisions (Nicholson and Ioannidis 2012; Greenblatt et al. 2024). In this paper, we study the aggregation of individual peer review evaluations of research proposals and the implications of translating those evaluations into decisions about which projects in a given area get funded.

The specific focus of our work is the U.S. National Institutes of Health (NIH) and the R01 grant program that is the dominant source of funding for academic biomedical labs within the US. NIH is the world's largest funder of research in the life sciences, distributing more than \$30B in funding each year, with most spent on basic research (Moses et al. 2005). This funding, in turn, serves as a vital building block for patenting and commercial success in the pharmaceutical and biotechnology sectors (Azoulay et al. 2019) and ultimately population health. Virtually all of NIH's funding decisions are based on the results of a highly structured

¹This amount is larger than the GDP of two-thirds of the world's economies. Including intramural funding for research conducted inside the federal government adds a further \$34 billion. Among NIH funding, more than 70% of extramural funding is awarded through peer review (NIH 2022).

²Note that philanthropic grants are distinct from (targeted) gifts made to universities, which often reflect different decision-making process such as naming rights or familial experience with a disease (Murray 2013).

peer review process that can be broken down into three parts: (1) allocation of funding across broad research areas, where Congress and the Executive Branch play a large role (Science News Staff 2022), (2) a peer review process for proposals within areas (Lee et al. 2013; Li and Agha 2015; Li 2017; Pier et al. 2018), and (3) the mechanism for using these peer reviews to inform funding decisions. This tiered system means that the evaluation of scientific merit that is conducted by peer reviewers is distinct from the funding decisions that are made by NIH staff based on those reviewers.

Our focus here is on the funding allocation decision, point (3), which has received little attention. NIH, like many organizations around the world, makes decisions based primarily on the first moment of the distribution of scores from peer reviewers (Guthrie, Ghiga, and Wooding 2018). Specifically, NIH elicits reviews from the panel of peer reviewers, calculates the average score across the reviewers, then ranks and funds projects based on that average (NIH 2008; Azoulay, Graff Zivin, and Manso 2012; Lauer 2023).³ This mechanism tends to result in granting funding to projects with consistent high marks across evaluators. Previous research on the NIH process has shown that higher peer review scores are predictive of better research outcomes (Li and Agha 2015). However, the process is also widely thought to favor incremental innovation over more radical ideas, and a common conjecture is that greater dissensus in project scores could be one way to identify those more radical projects (March 1991; Manso 2011; Azoulay, Graff Zivin, and Manso 2011; Nijstad, Berger-Selman, and De Dreu 2014).⁴ Our study evaluates the value researchers place on alternative peer review aggregation methods taking this conjecture as given.

Assessing whether the information from peer review scores could be aggregated to better effect would ideally entail a large randomized experiment that allocates grant applications to two or more different aggregation approaches and then tracks the outcomes that arise from those awards over a long time horizon.⁵ Since such an experiment is likely to be politically infeasible, it is important to explore alternative options. One such approach could make use of variation in peer reviewer scores across funded projects to examine whether projects with greater levels of dissensus generated more pathbreaking scientific discoveries. Unfortunately, NIH has not been willing to provide researchers with access to individual reviewer scores. NSF is similarly guarded about sharing individual review scores, and comparable data on

³Although we focus here specifically on the ranking of research project proposals, the issue of how to use noisy inputs to construct rankings is a central question more generally in statistical decision theory (Gu and Koenker 2023).

⁴Recent work finds evidence that this conjecture indeed holds for venture capital funding of startups (Gius 2024).

⁵It is worth noting that Chiara Franzoni of Polytechnic University of Milan and Paula Stephan of Georgia State University are currently running a related experiment that explores project funding decision rules with the Novo Nordisk Foundation.

corporate and foundation R&D decisions is even more elusive. Because this makes it impossible to explore the implications of alternative approaches to synthesizing scores using programmatic data, a simulacrum is required.

We used discrete choice experiments to effectively ask scientists what they think the aggregation function should look like when evaluating grant proposals (Louviere, Hensher, and Swait 2000).⁶ The participants in the experiments were active biomedical researchers with a track record of successful NIH funding, and the experiments simulated the research funding process NIH uses. Participants were presented with real (but anonymized) research proposal abstracts and a set of experimentally assigned peer review scores for those projects. They were then asked to choose which projects they would fund with their allocated budget.

The distribution of peer review scores was randomly drawn from an experimental design that allowed us to examine the weight participants placed on various moments of the score distributions.⁷ The core idea is that there is a clean null hypothesis: do participants place all of their weight on the mean value of scores, which is the decision rule that mirrors the current NIH approach? Our experimental design allows for a powerful test of whether participants value moments of the score distribution other than the mean.⁸ In particular, we can test whether the researchers prefer either more or less dissensus in reviewer ratings, conditional on the average score.⁹ Greater preference for dissensus is consistent with the notion that some level of dissensus may indicate more promising but radical ideas (Ackermann 1986; Goldstein and Kearney 2017; Krieger et al. 2022). And it is consistent with research showing that more diverse groups can make better decisions than more homogeneous groups (Hong and Page 2004). We are also able to look at other suggested deviations from NIH's mean-based funding rule.

The results show that our samples of experienced biomedical scientists, on average, do not share the same objective function as the NIH. In addition to the average peer review

⁶Discrete choice experiments have previously been used to study R&D decisions in a private firm context (Carson et al. 2022). Recent work has also used preference elicitation to study how scientists trade off grant length versus grant size, comparing the preferences of scientists to the preferences of granting agencies (Myers and Tham 2023).

⁷In addition to filling an important data gap, the experiment afforded us reasonable power to detect preferences for these attributes while presenting choices in a familiar way (9-point scale rankings, real project titles and abstracts). It also allowed us to investigate some of the hypothesized mechanisms that may be driving conservatism within the NIH peer review system.

⁸We note that it does not test whether scientists know an aggregation model for predicting the social value of an idea that performs better than the average. This is a question that would require (currently unavailable) data on real-world, individual reviewer scores.

⁹It is important to recognize that peer reviewers are assessing risky research proposals and NIH's mean score-based funding rule in that sense incorporates reviewer risk preferences. What it does not do is take into account the extra information contained in the distribution of reviewer scores. It is this extra information that participants in this experiment see in making choices concerning funding.

scores of projects, they also placed value on other moments of the project score distribution. Specifically, participants were willing to trade-off a project with lower average score for one with more variance. Participants were willing to accept an average score 0.1 points lower in exchange for an increase in score variance of 1. This effect holds true even when accounting for other characteristics of the project score distribution. On average, scientists also preferred projects that had a higher skew, indicating that they preferred the presence of more right-tail scores, even at the expense of good but not great overall scores. At the same time, controlling for skewness did not eliminate scientists' preference for pure dissensus in the form of higher variance.

Armed with data from our scientists' preferences, we explore heterogeneity in their preferences and the robustness of our findings to a range of potentially important features for shaping the relationship between risk-taking and the research proposal peer review process. We first assess whether the scientists in our sample weighted negative reviews more strongly than positive reviews—a trend that has been documented in previous research and which has motivated calls for reforms which would allow individual reviewers with strong preferences to overrule potential naysayers. In contrast to previous work, we find little heterogeneity in the effects of positive versus negative reviews, suggesting that decisions based on the full set of project scores can be effectively used to support riskier projects, as long as the process places positive weights on the variance of proposal scores.

Second, we ask whether scientists in the sample preferred projects that had bimodal scores, a particularly extreme form of dissensus in scoring suggestive of two opposing belief camps. We find that scientists did not prefer such projects relative to a model that simply accounts for high variance in scores.

Third, we use randomization in the proximity between a scientist's own research area and the research area of the projects we showed them to assess the dissensus preferences of relative experts versus outsiders (noting that the entire sample consisted of experts on relevant biomedical research). When acting as peer reviewers, previous research shows that scientists judge proposals inside their area of expertise relatively more harshly than proposals outside their area (Boudreau et al. 2016). Expert evaluators have also been found to focus first on feasibility of R&D proposals inside their own domain of expertise even at the cost of more innovative solutions (Lane et al. 2022a). These results raise concerns with review processes like those at the NIH, because expert peer reviewers might be especially unwilling to take risks on novel proposals in their area. Contrary to this concern, we find that participants who were in the best position to understand the proposal were substantially more tolerant of dissensus. The closer a proposal was to a researchers' own area of expertise, the stronger was the preference for project score variance. This novel finding on the risk-taking of insiders has important implications for the calculus that underlies the recently documented tensions between expertise and bias in the research proposal peer review process itself (Li 2017).

Fourth, we leverage results from an additional choice experiment, with an independent sample from our study population, to assess whether tighter funding budgets lead to lower dissensus tolerance. In this second study, participants were asked to construct portfolios of projects that they were willing to fund. We then administered a budget change shock by either tightening or relaxing the budget and asked them which projects they would cut from or add to the portfolio, in order to assess the characteristics of the marginal proposal. As expected, tightening the budget led participants to cut higher variance projects (those with greater dissensus in scores). The effect was not symmetric with relaxed budgets, however: a larger budget did not cause participants to notably add higher variance projects to the portfolio.

Putting things together, we assess the implications of the scientists' preferences for project funding.¹⁰ Using the project scores from this study, as well as three sets of expert-generated project scores repurposed from two previous studies (Pier et al. 2018; Lane et al. 2022b), we find that the funding rule based on the overall mean score and variance preferences of our successful biomedical scientists substantively alters which projects would get funded relative to the standard, mean-only NIH approach. On average across the four sets of project scores, fifty-eight percent of projects change their ranking when using the scientists' preferences both the average and variance of project scores can lead to changes in funding decisions for up to twenty percent of projects in some settings, with an average funding reversal rate of ten percent when using the preferences from scientists with relatively greater domain-specific expertise. This fraction of funded project reversals falls to five percent when we expand to include the full sample of scientists.

The rest of the paper proceeds as follows. Section 2 describes the NIH review process and context. Section 3 lays out the experimental design, randomization, and recruitment procedures. Section 4 gives details on the econometric model. Section 5 provides the results from fitting that model to the experimental data. Section 6 concludes.

¹⁰This assessment holds fixed any behavioral response that might be induced by changing the project scoring rule. Any rule regarding score aggregation creates incentives for strategic behavior. Assessing whether the scope for such behavior under a mean-variance rule relative to the current mean only rule is beyond the scope of this paper and an important area for future theoretical and empirical research.

2 Background on NIH Peer Review

The National Institutes of Health is made up of 27 different units called Institutes and Centers, each with a distinct, though sometimes overlapping, research agenda that is typically focused on a disease area or body system. For example, the National Cancer Institute, as the name suggests, focuses on cancer related research. The National Institute of Child Health and Human Development, in contrast, focuses on a wide range of diseases that afflict children. Nearly all Institutes receive their funding directly from Congress and manage their own budgets.

The NIH operates both an intramural and extramural research program. The latter, which is the focus of this paper, supports extramural research through competitive grants that are awarded to universities, medical schools, and other research institutions, and accounts for more than 80% of the total NIH budget. The largest of these grant mechanisms is the R01, a project-based research grant that accounts for half of all NIH grant spending and is the predominant 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 (Li 2017). The experimental task in this paper is a stylized representation of the R01 grantmaking process, though similar processes are also used for some other grants.

The NIH issues formal requests for proposals in priority areas, but investigators are also free to submit applications on unsolicited topics under the extramural research program. All applications are assigned to a review committee comprised of scientific peers, generally known as a study section.

If one's familiarity with peer review is primarily through the journal review process, it is useful to define the NIH review process in those terms and thereby highlight how the NIH approach differs from journal review.

At the initial submission stage, there are no desk rejections of proposals as long as they are correctly formatted. Part of being correctly formatted means including extensive institutional information and conflict of interest declarations.

Each proposal is assigned to a study section for scientific review and scoring. These applications are typically reviewed in one of about 180 "chartered" study sections, which are standing review committees organized around a particular theme, for instance, "Brain Injury and Neurovascular Disorders" or "Cancer Genetics." Unlike the journal peer review process, which typically enlists 3–5 expert reviewers, study sections are comprised of 15–30 members with expertise in relevant domains.

The most important difference between journal peer review and NIH review of scientific

proposals is that the NIH's review is inherent ex ante: it must be made before the results of the project are known. Thus, rather than a journal peer reviewer's focus on results, reviewers for the NIH focus on the nature and importance of the knowledge to be gained as well as its likelihood of success in producing the promised results.

More specifically, reviewers are asked to evaluate the scientific and technical merit of each proposal on the basis of five criteria: (1) Significance [does the proposal address an important issue?]; (2) Approach is the methodology sound?; (3) Innovation is the proposal novel?; (4) Investigator are the skills of the research team well matched to the project?; and (5) Environment are the institutions in which the work will take place conducive to project success?. Consistent with the notion that scientific review is distinct from funding decisions (see below), reviewers are asked to ignore budgetary concerns.

The review process within a study section proceeds in the following manner. Each application undergoes an initial review by three members of the section. These members assign a score to proposals for each of the five criteria described above as well as a score for overall impact (often called a priority score), which need not simply reflect the average of the criteria scores. Based on these preliminary "priority scores," weak applications (typically half) are rejected without further discussion. The remaining applications are then discussed in the full study section meeting, after which everyone is given the opportunity to revise their initial scoring based on the group deliberations before anonymously submitting their final scores. The overall priority score for the proposal is based on the average priority scores across all study section members. Scores are then normalized within review groups through the assignment of percentile scores to facilitate funding decisions.

Funding decisions are decoupled from the scientific review and determined by program areas within the Institutes and Centers (IC). Decision making is largely formulaic. Proposals are sorted from best to worst according to their percentile score (based on the mean of priority scores across all reviewers) and funded in that order until the relevant IC's budget is exhausted.

A grant's score is generally the sole determinant of the funding decision, irrespective of proposal costs (assuming they are deemed reasonable). It also means that funding rates may be quite different across disease areas based on IC budgets and the number and size of applications received in a given cycle. The worst percentile score that is funded is known as that IC's payline for the year. In rare cases. applications are not funded in order of score. This typically happens if new results emerge to strengthen the application (Li 2017). Scores are never made public. There is no revise and resubmit and no appeal of decisions. Rejected proposals can be revised and submitted to a later funding cycle (up to two more times) where a new set of proposals compete against each other with a new review panel.

For a young researcher not getting NIH funding in a particular funding cycle can be quite detrimental to career prospects.

3 Study Design

The scientists in our study took part in discrete choice experiments that involved ranking research projects in terms of their priority for being funded. The first study involved choice scenarios where funding priority across four projects was decided. The experiment also contained a randomized intervention that altered the match between the participant's research area and the subject of the presented projects.

Projects were assigned a randomized set of scores from a hypothetical expert review panel. The exact scores were shown, along with the average and standard deviation of the scores. This intervention allowed us to identify the preferences of participants for different features of the score distribution. In particular, it allowed us to test whether participants preferred projects with higher average scores or had preferences for other features of the score distribution, like dissensus.

Project titles and abstracts were shown above the scores (see Section D of the Appendix). The titles and abstracts came from real NIH grants and were chosen to span a range of biomedical research fields. Participants were randomized into an experiment where they saw projects from either inside or outside their specialty field. This allowed us to test for differences in behavior between insiders and outsiders.

A separate set of participants was randomized into a second study experiment that involved forming portfolios under different budget constraints. That alternative experiment is described in Section 3.2.

3.1 Design of Study 1: Estimating Preferences for Project Attributes

To estimate participant preferences for different distributions of project scores—particularly their preferences for consensus versus dissensus—we used a discrete choice experiment. In the experiment, each choice scenario involved ranking four proposals that had different distributions of scores from a hypothetical expert review panel. In this way, the participants were placed in the role of IC staff and advisory council members who select which projects to fund based on rating inputs from their study section's reviewers.

Participants were asked to complete four choice scenarios during the experiment. The choice scenarios were designed so that participants would be asked to rank projects with different average scores and score variances. Score variance was one of the main attributes of interest in the experiment because higher variance indicates greater dissensus among the

project reviewers.

Score distributions were generated using a balanced incomplete block design (BIBD), following Louviere, Flynn, and Marley (2015). BIBD designs are a type of fractional factorial design where preferences for combinations of different attributes or attribute levels are identified using a sparse matrix of choice options. We designed the BIBD to provide reasonably high power when estimating preferences over the average and standard deviation of project scores, while also allowing for estimation of preferences for other attributes of the project score distribution (e.g., number of top scores, number of bottom scores, score skewness).¹¹ The BIBD did so by generating scores for ten hypothetical raters using nine different score levels. Following standard NIH practice, ratings were on a 1 to 9 scale with 1 indicating the best possible score. These ratings were reverse coded for the statistical analysis, to be in keeping with the typical intuition that higher ratings are better. The ratings from the ten reviewers were duplicated twice to yield thirty scores for each project. From the set of all resulting possible score distributions, fifty-four orthogonal combinations of average scores and score standard deviations were used to create the projects shown to the participants.

For each question, the participant was provided with thirty reviewer ratings, along with the computed average and variance of those ratings, for four distinct proposals.¹² They were then asked to rank the four projects in terms of funding priority using a best/worst preference elicitation format (Louviere, Flynn, and Marley 2015). The four projects in each of these choice scenarios were chosen to maximize power to identify preferences across the project attribute combinations. This grouping yielded 344 blocks of four projects each. See Section D of the Appendix for examples of the questions that participants saw. The choice scenarios were further grouped into sets of four scenarios to create eighty-six survey versions. Participants were uniformly randomized into receiving one of the versions.

3.1.1 Project Title and Abstract Randomization

Participants were also randomized into receiving projects whose description (title and abstract) fell inside or outside their direct area of expertise. This randomization was done independently of the randomization into different survey versions described above. The purpose of this second randomization was to assess the effect of subject area expertise or insider status on the types of projects chosen.

¹¹The resulting projects possess some mechanical correlation between average score and score variance (for example, projects with very low and very high average scores have lower variance, on average). This correlation is accounted for by including both attributes simultaneously in our estimating equation.

¹²We chose to show the mean and variance both to avoid time consuming calculations for participants and to put the mean and variance on equal footing in terms of salience. In previous experiments with a similar design, directly displaying the variance (relative to no display) was not crucial to a participant's preferences over score variance (Carson et al. 2022).

All individuals recruited for the study had a background in biomedical research and were part of at least one of the five NIH study sections.¹³ Project titles and abstracts were selected from historical NIH R01, R35, or F32 grants listed on the Research Portfolio Online Reporting Tools (RePORTER) website in 2016. From the set of all potential grants, we kept those that were in one of the five study sections from which we recruited participants, and which had a project abstract length between 300 and 400 words, so the abstract would display consistently. All grants that could be tied to one of our study participants were dropped.

In total, sixteen title and abstract pairs were selected and were assigned to the discrete choice experiment projects. That assignment was done so that study participants would see either zero or one project(s) that matched their area of expertise. The matching was done based on the integrated review group (IRG) codes of the participants and the NIH proposal. Based on the randomization, thirty percent of the participants did not see any projects from their own IRGs. The remaining seventy percent of participants saw one choice scenario where all of the projects matched their IRGs and three choice scenarios where none of the projects matched their IRGs.

The IRG randomization was conducted at the study participant level. To identify the effect of proximity between a presented project and the participant's own research, while also including participant-level fixed effects, we constructed a more granular measure of research proximity using NIH Medical Subject Heading (MeSH) terms. The NIH maintains a structured dictionary of terms used for indexing research on PubMed, and all medical research can be assigned MeSH terms by passing it through an NIH indexing tool. We passed the titles and abstracts shown to participants and the grants received by participants through this tool, then calculated the proximity of a participant to a shown project by counting the unique, matching MeSH terms between the project and all of the participant's NIH-funded projects between 2012 to 2016, divided by the number of MeSH terms associated with the project.¹⁴

During each choice scenario, the project titles were shown above the project scores (see Section D of the Appendix for an example). All participants saw the project titles. If they hovered their mouse cursor over the title, they could also see the project abstract. Since not everyone chose to hover, we exploit this feature to further assess the veracity of our results on intellectual proximity. If an individual did not hover over the title to view the abstract, then the proximity of that abstract to the subject's research should be irrelevant to the project

 $^{^{13}}$ See Section 3.3 for details on recruitment.

¹⁴This measure of research proximity has been used in prior work on connections between researchers (Azoulay, Fons-Rosen, and Graff Zivin 2019).

ranking.

3.2 Design of Study 2: The Effect of Budget Constraints

A second study was conducted with a separate set of scientists to assess the role of budget constraints on project funding preferences. The design utilized a similar discrete choice setup as Study 1, with two main differences. First, the participants were shown ten potential projects and asked to choose the four that they would most like to fund. This was presented as constructing a portfolio of projects (see Section D of the Appendix for an example of the choice scenario). The main goal of the study was to determine how individuals responded to tighter budgets, so after choosing their portfolio, participants were initially told that the budget had been cut, only allowing them to fund three projects. They were asked which project they would like to drop. Next, they were asked which project (of the six they did not select for funding) they would add if the budget were expanded to allow for the selection of five projects. This variation allowed us to identify the marginal project initially selected and rejected, to determine whether budgetary pressure affects preferences for project attributes. Each participant engaged in two of these choice scenarios.

3.3 Recruitment and Sample Construction

The initial sampling frame consisted of the set of all researchers who had received a R01, R35, or F32 NIH grant between 2012 and 2016, from any of the following IRGs: Brain Disorders and Clinical Neuroscience (BDCN), Cell Biology (CB); Molecular, Cellular, and Developmental Neuroscience (MDCN); Oncology-Basic Translational (OBT); or Oncology–Translational Clinical (OTC).¹⁵ We further restricted the sample to individuals who were part of a study section that mapped to only one IRG code, to focus on individuals working within a single, albeit broad, scientific domain. The names and contact information for this set of potential participants was gathered from the NIH RePORTER database, yielding 6,678 total initial contacts.

These initial contacts were randomized into two groups. First, fifty percent (3,339) of the contacts were randomized into the group receiving the project ranking survey (Study 1). Second, the remaining fifty percent (3,339) of the contacts were randomized into the budget experiment (Study 2). Table A1 shows the summary statistics for contacts, broken down by randomization group.

Of the 6,678 scientists contacted by email, 590 either declined to participate or had an outdated email address (leading to the email bouncing), leading to a final contact sample of

¹⁵The NIH Center for Scientific Review initially reviews grant submissions and assigns the submission to an IRG for assessment of scientific and technical merit.

	Study 1				Study 2	
	(1)	(2)	(3)	(4)	(5)	(6)
	Attrited	Finished	Diff.	Attrited	Finished	Diff.
	Mean	Mean	Mean	Mean	Mean	Mean
	[SD]	[SD]	(SE)	[SD]	[SD]	(SE)
Fraction BDCN	0.25	0.26	-0.0071	0.27	0.30	-0.035
	[0.43]	[0.44]	(0.026)	[0.44]	[0.46]	(0.029)
Fraction CB	0.22	0.25	-0.024	0.20	0.20	-0.0020
	[0.42]	[0.43]	(0.025)	[0.40]	[0.40]	(0.026)
Fraction MDCN	0.17	0.22	-0.050	0.17	0.19	-0.019
	[0.38]	[0.42]	(0.023)	[0.38]	[0.39]	(0.025)
Fraction OBT/OTC	0.36	0.28	0.082	0.36	0.30	0.056
	[0.48]	[0.45]	(0.028)	[0.48]	[0.46]	(0.031)
Total funding	6.64	5.68	0.96	6.71	6.02	0.68
	[9.29]	[6.12]	(0.54)	[8.28]	[8.17]	(0.54)
Unique projects	4.23	4.12	0.11	4.32	4.05	0.27
	[2.92]	[2.73]	(0.54)	[2.93]	[2.57]	(0.19)
Total projects	16.3	15.0	1.27	16.6	15.7	0.82
	[15.6]	[15.2]	(0.93)	[15.8]	[13.3]	(1.00)
Ν	3,026	313		3,089	250	

Table 1: Attrition

This table shows statistics for the sample of individuals who were contacted but did not complete the experiment (Column 1 for Study 1 and Column 4 for Study 2) versus those who completed the experiment (Column 2 for Study 1 and Column 5 for Study 2). Mean values are above and standard deviations are in the square braces below. Columns 3 and 6 show the difference in means between the two groups for Study 1 and Study 2, respectively. Standard errors are in parentheses below each value. "Total funding" is all NIH grant funding from 2012 to 2016. "Unique projects" counts unique NIH grants and "Total projects" is grants by years of grant funding from 2012 to 2016.

6,088. Across the two studies, 563 participants completed all portions of the experiments, for a response rate of 9.2%.¹⁶ 313 participants completed Study 1 and 250 participants completed Study 2.

Attrition in the two experiments is assessed in Table 1. We assess differential attrition using variables from RePORTER, which contains information on the NIH activity for everyone in the sample regardless of survey completion. Across both studies, the sample of completers versus attriters is comparable for the measures we can assess. The largest standardized difference is that participants with grants in the IRG code group MDCN were about

¹⁶This response rate is consistent with, if not substantially higher than, other recent surveys of active scientists. For instance, for three recent papers that surveyed active scientists Myers et al. (2020) had a 1.6% response rate, Myers and Tham (2023) had a 3.3% rate, and Tawfik et al. (2020) had a 4.1% rate.

30% more likely to finish, while those in OBT/OTC where about 30% less likely to finish Study 1, compared to the group that did not complete the studies.

4 Estimating Equation

Each participant's preferences for different project attributes were estimated by fitting a model for the probability that a participant would choose a given project. The baseline results show the fit from the conditional logit model

$$\Pr(y_{ijk} = 1 | \mathbf{x}_{ijk}) = F(\alpha_i + \beta_1 \operatorname{avg}_{ik} + \beta_2 \operatorname{var}_{jk} + \mathbf{z}_{ijk}\theta)$$
(1)

for participant *i* making a choice about project *k* as part of the choice scenario version and choice set j.¹⁷ The function *F* is the cumulative logistic distribution.

The conditioning variables are indicated by x and fall into three groups: the mean and variance of project scores, other project attributes, and controls. The main right-hand-side variables of interest are project score attributes, with a particular focus on the mean and variance of project scores. If the scientists only cared about the average project scores, that would show up as a non-zero coefficient on average score and a zero coefficient on the score variance. In contrast, if they valued dissensus, they might still place a non-zero weight on the average score, but the coefficient on the score variance would be positive. Additional results allow for estimation of preferences around other project score attributes and project descriptions. For example, we assessed the effect of the count of individual project score levels, the effect of higher moments of the project score distribution (e.g., skew), participant expertise or experience.

Control variables are fixed effects for each participant, α_i , such that all estimates reflect the average preferences of a given scientist, since the project attributes shown to that scientist were varied (average score, score variance, project description match with the scientist's research, etc.).¹⁸ Standard errors were clustered at the scientist level.

¹⁷We converted rankings into binary choices by considering each choice scenario to be composed of three different choice sets. In the first set, all projects are in the choice set and the chosen project is the top ranked one. The second choice set consists of all projects other than the top ranked one and the chosen project is the second ranked project. The third and final choice set consists of the remaining two projects and the chosen project is the third ranked project. Results are similar if we use a multinomial logit (see Table A3), but our binary conditional logit approach allows for more granular fixed effects controls. Inferential accuracy is maintained by clustering at the participant level.

¹⁸Participant fixed effects subsume project fixed effects, so the results are unchanged by the inclusion or exclusion of project fixed effects.

5 Results

5.1 Scientist Preferences for Dissensus

We first show models for scientist preferences over different attributes of project scores. The results from fitting Equation (1) are shown in Table 2. The dependent variable is equal to 1 if the participant chose a project in a given choice scenario. All right-hand-side variables are standardized so that the coefficient magnitudes are comparable.

	(1)	(2)	(3)
	Project choice	Project choice	Project choice
Avg. score	0.82***	0.92***	1.05***
	(0.041)	(0.056)	(0.088)
Score variance	0.085^{***}	0.11^{***}	0.084^{**}
	(0.027)	(0.033)	(0.035)
Score skew		0.10^{***}	0.22^{***}
		(0.035)	(0.064)
Minimum score			-0.14**
			(0.055)
Maximum score			-0.034
			(0.039)
Clusters	313	313	313
Ν	11268	11086	11086

 Table 2: Scientist Preferences Over Project Scores

This table shows results from estimating Equation (1) on the baseline sample. The dependent variable is an indicator for whether the participant chose a given project. All right-hand side variables were standardized. The models include participant fixed effects. Standard errors are clustered at the participant level: * p < .10, ** p < .05, *** p < .01.

Column 1 performs the simplest and most direct test of whether the scientists' preferences match the funding rule followed by NIH. The first coefficient shows that participants strongly preferred projects with higher average ratings. For fixed effect values of 0 and other variables held at their mean, the model implies that a one standard deviation increase in average project score increased the chance that the project was chosen by eighteen percentage points, a fifty-five percent increase relative to the baseline probability (thirty-three percent) that a project was chosen.

The second coefficient shows that the scientists also preferred projects with higher score variance. Conditional on the average score, a one standard deviation increase in the variance of project scores increased the chance that the project was chosen by 1.8 percentage points (a 5.4% increase). This effect shows that scientists were dissensus-seeking on average. It also runs counter to the mean-only scoring rules currently used by organizations such as NIH.

The subsequent columns test whether the preferences were for higher variance per se or for other correlated project attributes, some of which also indicate a preference for dissensus. Participants could appear to prefer high variance projects, for instance, if they simply preferred high scores and were relatively insensitive to the rest of the score distribution. Column 2 adds the skewness of project scores and Column 3 adds the minimum and maximum score assigned to the project, to see if these attributes explain the variance preferences. In both cases, the preference for higher variance projects persists.¹⁹ The preference for skewness shows that participants were not simply choosing projects based on the statistics shown in the tables but were attendant to richer information about the distribution of scores.

In Column 2, for example, even after controlling for skewness and average score, participants still preferred higher variance projects. If anything, the preference for higher variance projects appears stronger. At the same time, participants also preferred projects with higher skew, with an effect size comparable to that of variance. The average project in our sample had a small negative skew, so an increase in skewness for that project, at the margin tended to result in a more symmetric distribution (while holding the mean and variance fixed).²⁰ This preference is consistent with scientists placing substantial value on high scores, particularly if the rest of the scores were concentrated near the middle of the range.

Column 3 adds the minimum and maximum scores to the estimating equation. Across all projects and after reverse coding, the minimum possible score was 1 and the maximum was 9, but different projects had different highest or lowest scores depending on their exact score distributions. A project with a maximum score below 9 or a minimum score above 1 often had a lower score variance than a project with scores across the full range, so Column 3 adds controls for the actual range. The results show that scientists preferred projects, on average, when both the minimum and maximum scores were higher, but these preferences were not as strong as the preference for variance. The maximum score effect is not estimated precisely enough to reject at the five percent level that the preferences were zero.

¹⁹Table A2 adds further score statistics including kurtosis, interaction between mean and variance, the number of lowest or highest scores in the score distribution, and indicators for whether the project had at least one score of 1 (lowest possible score, after reverse coding) or 9 (highest possible score). In all cases, the estimated effect of score variance remains consistent.

²⁰Heterogeneity analysis reveals that participants also preferred it when positively skewed distributions became even more positively skewed. See Table A6.

5.1.1 Robustness and Sensitivity Checks

Table A3 uses a multinomial logit model to estimate the effect of project scores on project rankings. The multinomial logit relaxes the assumption of homogeneous coefficients across the three choice sets involved in ranking projects, but at the cost of not including highdimensional fixed effects. Estimating using the multinomial logit shows that the results are in line with the baseline binary conditional logit model with some notable heterogeneity across choice sets. When choosing the highest and second highest ranked projects, scientists preferred higher mean and higher variance projects. When ranking the third versus the fourth project, the scientists were indifferent between higher or lower variance. They also cared less about the average score. Similarly, if we allow mean and variance preferences to vary by choice order in the baseline model, we find no difference in variance preference for the choice of first and second rated projects. Participants cared less about variance when choosing between the third and fourth rated projects.

The conditional logit model is known to do a good job summarizing aggregate sample preferences even when there is considerable underlying preference heterogeneity (Allenby and Rossi 1991). Nevertheless, economists are often interested in characterizing preference heterogeneity using more flexible models. In Table A5, we present a series of five models which explicitly relax the conditional logit model's assumptions: (mixed) random parameters (Train 2009), scale heterogeneity (Swait and Louviere 1993) and generalized multinomial (Fiebig et al. 2010) logit models. Results show that individual respondents differ in their tastes for both dissensus and mean proposal score. What motivates this paper though is not the individual level trade offs but the average of these trade offs across respondents which better resembles the panel decision making process that governs the NIH allocation process. Despite radically different approaches, the more flexible models with individualspecific parameter estimates and the conditional logit model all yield similar preferences for mean score and dissensus, on average.

Table A4 evaluates the sensitivity of the results to sample and control changes. As discussed in Section 3.3, 313 participants completed the full study while 356 participants started the experiment and completed at least one ranking exercise. The results show that preferences were unchanged if all available data were included (meaning that we include not only the participants who completed all questions, as in the baseline results, but also the participants who partially completed the survey). The second column adds more granular fixed effects for the interaction of participant, question version, and choice scenario. Including these fixed effects, if anything, increases the magnitude of the estimated preference for dissensus. The results are also unchanged if we restrict the sample to individuals who spent

more than 10 minutes on the survey (the 10^{th} percentile of elapsed time among completers).²¹ Further robustness checks are reported in Section B of the Appendix.

5.2 Assessing Hypotheses About Dissensus Tolerance and Proposed Funding Reforms

Many commentators, including directors at NIH, have suggested that NIH is too cautious when funding research. The results above show that the average scientist in our sample agrees with that sentiment. Here we assess explanations that have been proffered for why funding decisions might be so dissensus-intolerant, and investigate scientist preferences for reforms that have been suggested to make the process less risk averse.

5.2.1 Are Positive and Negative Reviews Weighted Differently?

Testing for dissensus preferences using score variance treats low and high scores symmetrically. Previous studies of peer review for scientific grants have emphasized that negative reviews can have an oversized influence on the probability that a grant gets funded.²² And consensus has been shown to emphasize the influence of negative scores (Lane et al. 2022b), as well as to increase variability of ratings of similar projects across different sets of reviewers (Pier et al. 2017). In response, a variety of reforms have been proposed that would bypass some or all of the consensus-based peer review processes. For example, foundations have experimented with a so-called "golden ticket" that allows a reviewer to ensure that an application gets funded, even over the objections or low ratings of other reviewers (Sinkjaer 2018). A similar reform has also been proposed for Program Officers at NIH (Buck 2022).

Although we cannot directly test whether the scientists in our study would prefer a golden ticket-style selection procedure, we can test the underlying basis for that proposal—the idea that negative reviews exert an oversized influence. This hypothesis is assessed in Figure 1. The figure shows the effect, estimated from a conditional logit model, on the choice of project coming from the addition of one score from the range of possible scores. The omitted score is 5, the midpoint of the range from the best (reverse coded) score of 9 to the lowest score of 1. The coefficients can be interpreted as the effect of replacing a score of 5 with the score indicated on the x-axis. The dashed line shows a linear fit to the point estimates.

The results show that choice probability was monotonically increasing in score, and that the effect of a low score was roughly symmetric with the effect of a high score. In particular,

 $^{^{21}}$ Given that participants could stop taking the survey at any time and that they did not receive a completion payment, they had no incentive to waste their own time providing low-quality responses.

 $^{^{22}}$ In particular, Jerrim and Vries (2020) found that "a single negative peer review is shown to reduce the chances of a proposal being funding from around 55% to around 25% (even when it has otherwise been rated highly)."

Figure 1: Marginal Effect of Each Score on Choice Probability



Notes: The figure shows the marginal effect of each possible project scores on the probability that the project was selected, relative to a score of 5. The estimates were generated by fitting a version of Equation (1) where the project attributes are the count of scores at each score level. The equation includes subject fixed effects. See Table A7 for the numerical coefficients. Whiskers are 95% confidence intervals based on standard errors clustered at the participant level.

replacing a score of 5 with a score of 1 *reduced* the probability of a project being chosen by almost the same amount that replacing a score of 5 with a 9 *increased* the probability. A formal hypothesis test to see whether the sum of the coefficients is 0 yields a coefficient of 0.004 with a p-value of 0.25.

For less extreme scores, we did find some evidence for asymmetry. A score of 2 was penalized almost the same amount as a score of 1, while a score of 8 raised the probability of selection by less than would be expected based on the average slope of the marginal effects (as indicated by the dashed line). Even here, though, we cannot reject the hypothesis that the scores had marginal effects of the same magnitude. Overall, the results do not support the idea that negative scores disproportionately caused scientists to think poorly of a project. Instead, scores had a roughly uniform effect across the distribution of possible scores.

5.2.2 Does Bimodality Better Capture Scientists' Preferences for Dissensus?

Buck (2022) proposes that projects with bimodal scores could receive higher funding priority as a way to reduce the conservatism of funding decisions. What preferences did the scientists in our experiment exhibit along this dimension? Table 3 shows estimated preferences for projects with bimodal scores (Column 1) and simultaneously for bimodality and higher variance (Column 2). In both cases, the average score was included as a control. Projects were classified as bimodal using the dip test from Hartigan and Hartigan (1985), as implemented in Stata by Cox (2016).²³

	(1)	(2)
	Project choice	Project choice
Avg. score	0.77***	0.82***
	(0.037)	(0.041)
Bimodal	-0.052	-0.037
	(0.13)	(0.12)
Score variance		0.085^{***}
		(0.027)
Clusters	313	313
Ν	11268	11268

 Table 3: Preferences for Bimodality Versus Variance

Notes: This table shows results from estimating Equation (1) on the baseline sample. The dependent variable is an indicator for whether the participant chose a given project. The average score and score variance variables were standardized. The variable "bimodal" is an indicator for the project scores exhibiting a dip statistic greater than 0.1 (Hartigan and Hartigan 1985). The models included participant fixed effects. Standard errors are clustered at the participant level: * p < .10, ** p < .05, *** p < .01.

The table shows that scientists did not prefer bimodal projects. Moreover, the preference for dissensus, as captured by project score variance, was unaffected by the inclusion of the bimodality measure. Bimodality is a particularly extreme form of dissensus that was not favored by the participants in our sample.

At the same time, bimodality is rare in both the scores we showed to participants and in current, real-world NIH scores. Over time, NIH has worked to avoid strategic behavior that results in bimodal scores.²⁴

 $^{^{23}}$ In the table, the variable "bimodal" is an indicator for whether the dip statistic was above 0.1, although the results are robust to alternative cutoffs and available upon request.

 $^{^{24}}$ For example, Ogden and Goldberg (2002) describes the move to percentile rankings as a method to reduce behavior by some reviewers of inflating the scores of projects that they like, while simultaneously lowering the scores of competing projects, so that the favored project would look even better by comparison. More recently, NIH has relied on training and guidance to reviewers (Sampat 2023).

5.2.3 Does Expertise Decrease Dissensus Tolerance?

Lay observers, scientists, granting agencies, and previous research studies have debated whether expertise and experience increase or decrease the willingness of scientists to engage in high-risk research. The effect of expertise on dissensus tolerance could go in either direction. On one hand, greater expertise might increase a scientist's convictions about the correct direction of research, making them less subject to consensus-driven selection criteria. On the other hand, previous work has shown that the removal of incumbent researchers in a field can spur innovation (Azoulay, Fons-Rosen, and Graff Zivin 2019), and recent work shows that the creativity of patents quickly declines with experience (Kalyani 2022). Arthur C. Clarke, in a quote that has come to be known as Clarke's Law, offered some additional nuance by arguing that the effect of experience is asymmetric: "When a distinguished but elderly scientist states that something is possible, he is almost certainly right. When he states that something is impossible, he is very probably wrong" (Clarke 1962).

Understanding the direction of this effect is important because expert review is at the heart of nearly all scientific project evaluation, whether for funding or publication purposes. NIH in particular relies heavily on carefully matched peer evaluators when judging grant quality.

We assessed the effect of expertise and experience with the estimates shown in Table 4. Overall, we found that expertise *increased* dissensus tolerance. In other words, participants who were in the best position to understand the proposal had substantially stronger preferences for higher project score variance. Column 1 shows the effect of proximity between the shown project and the participant's research area, based on our measure of MeSH term overlap between the shown projects and the participant's NIH grants from 2012–2016. A stronger overlap in these terms indicates that the scientist was active in the area from which the project's description was drawn, and thus measures the degree to which the scientist was a relative insider for the specific field represented by the project. Recall that the projects shown to the participants were randomized to be either closer to or further from their field.

The results show that scientists modestly preferred projects that were more inside their research area. The "MeSH match" coefficient is positive and significant at the ten percent level, with an effect size that is about half the size of the effect of project score variance. The interactions between this measure of expertise and project score statistics shows that experts had a significantly stronger preference for dissensus, as indicated by the positive coefficient on the interaction between score variance and expertise, as measured by MeSH match. Given that all variables are standardized, the coefficient on "score variance" indicates the preference that a scientist with an average MeSH match had for a project with higher variance scores. The results show that this scientist preferred higher variance projects, on

	(1)	(2)	(3)
	Expertise	Place	ebo Check
	Project	Project	Project
	choice	choice	choice
Avg. score	0.82***	0.70***	0.98***
	(0.041)	(0.052)	(0.064)
Score variance	0.083^{***}	0.079^{**}	0.091^{**}
	(0.027)	(0.037)	(0.040)
MeSH match	0.036^{*}	0.053^{*}	0.014
	(0.021)	(0.029)	(0.031)
Avg. score \times MeSH match	0.032	0.048	0.0041
	(0.029)	(0.042)	(0.039)
Score variance \times MeSH match	0.050^{**}	0.079^{**}	0.011
	(0.024)	(0.040)	(0.029)
Hover subgroup	Full sample	Always	Never/rarely
Clusters	313	169	144
Ν	11268	6084	5184

 Table 4: Preferences for Projects by Experts

Notes: This table shows results from estimating Equation (1) on the baseline sample. The dependent variable is an indicator for whether the participant chose a given project. All right-hand-side variables were standardized. The models include participant fixed effects. Standard errors are clustered at the participant level: * p < .10, ** p < .05, *** p < .01.

average, and that the preference was about one-tenth as strong as the preference for higher average score.

A scientist with a one standard deviation higher MeSH match showed little difference in their preferences over average scores but a substantially stronger preference for higher variance. In particular, the variance preferences were sixteen percent as strong as average score preferences for such individuals. Going the other direction, the results indicate that a scientist who was relatively far from the area of the shown project (one who has a 1 standard deviation lower MeSH match) placed almost no weight on project score variance.

The second and third columns show the results of a placebo test that was built into the experiment to determine whether the results from Column 1 were driven by the study participants actually taking the time to understand the abstracts that were shown, instead of simply acting differently than individuals with lower match rates for reasons unrelated to project content. To see the abstracts of the projects included in the experiment, participants needed to hover over links. Column 2 shows the results for subjects who reported always hovering over the links. Column 3 reports results for subjects who said they rarely or never hovered to look at the abstracts. While the endogeneity of hovering means that these results should be interpreted with caution, one can see that the effect of expertise is substantially stronger for subjects who did report looking at the abstracts.²⁵

5.2.4 Do Tighter Budgets Decrease Dissensus Tolerance?

Francis Collins, NIH Director from 2009 to 2021, argued that budgetary pressures reduce scientific risk-taking, stating (emphasis ours): "Although the two-level NIH peer-review process is much admired and much copied around the world, its potential tendency toward conservatism is a chronic concern and *invariably worsens when funding is very tight*."²⁶

Study 2 allowed us to test this hypothesis. The results are shown in Table 5. The first two columns show estimates for the attributes of the project that was dropped when scientists were told that the budget had been reduced. Column 1 shows the characteristics of the dropped project compared to the projects that were kept. Unsurprisingly, the dropped project had a lower average score compared to the projects that were kept in the portfolio. Lending support to Collins' statement, the dropped project also tended to have higher score variance. When faced with tighter budgets, participants preferentially dropped riskier projects characterized by higher dissensus.

Column 2 compares the dropped project to the projects that were originally not chosen for the portfolio of funded projects. Here, the average score clearly played an important role, but the variance of scores was no longer as important. The effect size is substantially smaller than when comparing the dropped project to projects that were kept in the portfolio, and the effect is not statistically significant.

Columns 3 and 4 show the characteristics of the projects that were added when budgets were expanded, with Column 3 showing the comparison to the four projects that were already chosen and Column 4 showing the comparison with the projects that were not originally chosen. The variance of scores appeared to play little role in this choice.

Together, these results provide nuanced evidence for Collins' claim. Compared to projects that were kept in the portfolio, tighter budgets did cause scientists in our sample to cut higher-variance projects. But the reverse was not true for more expansive budgets, and the cut project was not substantially different than other non-chosen projects in terms of variance.

²⁵Results using other subject-specific heterogeneity measures are shown in Table A8. In the sample, men were more dissensus-loving than women. An elicited measure of risk aversion did not strongly predict dissensus preference. And individuals with greater breadth in their research, as measured by the total number of unique MeSH terms, were more tolerant of dissensus.

 $^{^{26}}$ Quoted in Kolata (2009).

	(1)	(2)	(3)	(4)	
	Tighter	Budget	Relaxed Budget		
	Dropped proj. compared to kept	Dropped proj. compared to not chosen	Added proj. compared to kept	Added proj. compared to not chosen	
Avg. score	-0.82***	1.88***	-1.80***	0.34***	
	(0.10)	(0.14)	(0.14)	(0.056)	
Score variance	0.17^{**}	0.078	-0.045	-0.071	
	(0.081)	(0.063)	(0.072)	(0.060)	
Clusters	250	250	250	250	
Ν	1983	3516	2483	3516	

Table 5: Effect of Constrained or Relaxed Budgets

Notes: This table shows results from estimating Equation (1) on the baseline sample. The dependent variable is an indicator for whether the participant chose a given project. All right-hand side variables are standardized. The models include participant fixed effects. Standard errors are clustered at the participant level: * p < .10, ** p < .05, *** p < .01.

5.3 Implications for Project Funding

How large is the difference between the procedure NIH uses for funding (mean score) and the preferences possessed by the scientists in our study when it comes to actually ranking and funding projects? Although NIH does not maintain data on project scores and funding decisions that would allow us to test this question on historical NIH proposals, three datasets illuminate the scale of the difference. First, we calculated the changes in rankings for the fifty-four unique mean-variance combinations in the projects that we showed to participants in the first study. Second, we also repurposed two prior experiments that closely replicated the NIH review process. The first of these was Pier et al. (2018), which carefully simulated the NIH review process using real NIH reviewers, former study section leaders, and proposals. The second study, Lane et al. (2022b), conducted two experiments involving the evaluation of real submissions to a pair of small grants competitions in translational medicine run by a large U.S. medical school. The data from Lane et al. (2022b) is especially revealing because it allowed us to assess whether the scientists' preferences would have resulted in different real-world funding decisions.

For the fifty-four different project score mean and variance combinations included in our Study 1, the overall ranking for half of them changed when ranked according to the mean and variance preferences given in Table 2, Column 1, versus a ranking purely based on mean score. The largest changes in overall rank occurred, naturally, for projects that had the highest variance. Given that project scores were bounded, these projects also tended to have average scores that were closer to the middle of the pack.

Thus, high variance caused two effects that drove a wedge between the NIH-style mean score ranking and the rankings that the scientists in our sample preferred. First, consider two projects with the same mean but different variances. The NIH procedure would give these two projects the same score, while the scientists gave the higher variance project a higher score. Thus, the NIH procedure gave the high variance project a relatively lower rank than the scientists. Second, consider two projects with different average scores. A higher average score was mechanically, positively correlated with lower variance project lower (because of its lower mean score), while the scientists ranked the two projects closer together.

These two effects can be seen by comparing individual proposals drawn from our Study 1. To illustrate the first mechanism, we focus on two projects that had an average score of 6.3, but one had a low variance of 3.3 while the other had a high variance of 9.5. The NIH procedure would rank both of these projects right around the 50^{th} percentile across the entire set of projects in our study. Using scientists' preferences, however, would put the higher variance project at the 63^{rd} percentile and the low variance project at the 44^{th} percentile.

To illustrate the second mechanism, we can again consider the high variance project with an average score of 6.3 and a variance of 9.5. But this time we compare it to a project with an average score of 6.5 and a score variance of 2.3. The NIH procedure would rank the latter project in the 63^{rd} percentile of the overall project distribution, well ahead of the higher variance project (even though the difference in their means is only one-quarter of the standard deviation in average project scores across the experiment). If we ranked them according to the scientists' preferences, the lower variance project would drop down to the 55^{th} percentile, while the high variance project would again move up to the 63^{rd} percentile.

Using scores generated by Pier et al. (2018), which strove to closely replicate the NIH review process, we also found substantial differences in project ranking between the two procedures. In the Pier et al. study, many projects near the top of the ranking received identical average scores. At the 80^{th} percentile, five studies were given the same average score of 7. Using variance, one can break three of these ties, with the highest variance project (variance of 4.7) being ranked first among the set, the lowest variance project (variance of 0.7) ranking last, and the remaining three projects with a variance of 1 being ranked in the middle.

This example from Pier et al. highlights an additional insight from our results. Taking





Notes: This figure shows the reversal rate for project funding as a function of the payline (fraction of projects that get funded) for four different sets of project scores (Study 1 from this paper, the two studies from Lane et al. (2022b), and Pier et al. (2018)). The reversal rate is the fraction of studies that changed whether they were funded under a mean-only ranking versus the mean and variance-based ranking. The lines are LOESS fits to reversal rates calculated at each payline percentile. The solid line shows the reversal rate when using the estimated preferences from the baseline results using the full sample of scientists. The dashed line shows the reversal rate when using the preferences of scientists who were relative experts (a MeSH match 1 standard deviation higher than average).

variance into account can help break ties that often emerge when a relatively small set of reviewers are judging each project. The data from the two experiments in Lane et al. (2022b) allow us to examine how actual funding decisions would have changed if variance had been taken into account. We did so by first ranking projects by their average score.²⁷ Multiple reviewers rated each project, which allowed us to also calculate the variance of scores and re-rank the projects using the project score attribute preferences from our scientists. Importantly, we found that in both experiments, accounting for variance would have led to different projects being funded: the marginal projects funded would have been switched, with a higher-variance, unfunded project replacing a lower-variance project that actually did get funding.

We call such a change in funding a "reversal" of the project funding decision. For any given possible payline (the fraction of projects that get funded), we can calculate the reversal

²⁷The main goal in Lane at al. is to study the effect of showing reviewers scores from other reviewers to assess how exposure to others' scores affects one's own rankings. Thus, we only used the original, independent scores that participants provided for the exercises described here.

rate for the four sets of project ratings described in this section. The reversal rate is the fraction of projects funded at the payline that changes when we move from a mean-only to a mean-and-variance ranking.²⁸ Figure 2 shows the average reversal rate across sets of projects from the different studies (this study, Pier et al., and two datasets from Lane et al.) as a function of the payline.

Starting first with the scientists who were in the position to better understand the proposals (those with MeSH match values 1 standard deviation higher than average), we see from the dashed line that the reversal rate was around 10% for all paylines. Even when we expand to include the preferences for all scientists estimated in Table 2 Column 1, we still see reversal rates of 4 to 6%, depending on the payline. The highest reversal rates are near typical NIH paylines of 10 to 20%.²⁹ And this average reversal rate masks high rates that can appear for individual sets of project scores. Figure A1 shows the reversal rates separately for each of the four studies. Rates are as high as 20% for the proposals from Lane et al. (2022a). Together, these results underscore that variance preferences are not only statistically important, but can be consequential for funding decisions, and particularly so in cases that closely mimic real NIH grantmaking.

6 Conclusion

Scientific research, through its influence on technological innovation, has long been recognized as an important contributor to aggregate income (Nelson and Phelps 1966) and a driver of economic growth (Lucas 1988; Romer 1990), yet the path from research to innovation is uncertain, requiring institutions that make substantial scientific investments to appropriately balance risk and return in the portfolio of projects they support. Research projects that closely build on existing scientific knowledge may be a relatively safe bet, but the incremental innovation they produce may have lesser social value. In contrast, research that eschews conventional wisdom for more speculative pursuits may be required to produce radical or paradigm-shifting innovations of enormous value, but it is also much more likely to end in failure (see Eric Lander as quoted in Fallows (2014); Manso (2011)). The design of public and private institutional structures employed to evaluate research projects plays a critical role in balancing the risk and rewards from research, which, in turn, informs future scientific frontiers.

²⁸The reversal rate only counts the rate at which projects go from funded to unfunded because each change that causes a project to lose funding will cause another project to gain funding given the fixed funding constraint. This definition avoids double counting.

²⁹For example, the National Institute of Allergy and Infectious Disease at NIH published annual information on paylines for grants. The payline for R01 grants in 2022 was twelve percent (sixteen percent for new PIs).

The focus of this paper is on the peer review process for research proposals and how NIH (and other science-based agencies) synthesizes the output of that process into resource allocation decisions. Of particular concern is that agencies base funding decisions on the average of peer review scores, ignoring higher moments of the score distribution that may confer valuable information about the radicality of a scientific proposal. Since data on individual scores from NIH is unavailable to the research community, we leveraged data from two novel discrete choice experiments, fielded in samples of active biomedical scientists with a successful NIH grant history, to assess their preferences for aggregating peer review evaluations into scientific funding decisions.

In contrast with current practice, we found that these scientists—the very scientists that NIH relies upon for expert evaluations of research proposals—preferred to fund projects where there was some disagreement among reviewers. This preference for higher-dissensus projects was not driven by lone wolf reviewers who were enamored with a project, nor was it driven by focus on an aberrant, critical review. Rather, it appears that our experts valued healthy disagreement over either middle-of-the-road reviews or more extreme forms of dissensus such as projects that received bimodal scores. While this appetite for risk shrank as budgets became tighter, it did not completely disappear. We also found that those scientists with relatively greater domain expertise on a proposal were consistently more enthusiastic about dissensus in their reviews than those asked to make decisions outside their specific area of expertise. Applying our estimates to prior studies that mimic the NIH review process suggests that incorporating preferences for dissensus would lead to changes in billions of dollars of research funding annually.

Our results should not be construed as a critique of the peer review process. Indeed, we believe the impartial review of proposals by experts in the field is essential for prioritizing scientific investments by both public and private agencies. The substance of our inquiry relates whether there is relevant information from that process beyond the simple mean of reviewer scores that should influence the funding decisions of a major government entity charged with funding risky R&D projects related to improving the public's health. While our findings have implications for funding rule reforms that could prove important, many questions remain unanswered. Fundamental for the tasks ahead is a better understanding of the causal relationship between peer review scores and scientific impact. This will require a clever mix of experimental design and currently unavailable data from funding agencies containing individual reviewer scores on projects being evaluated.³⁰ Prospective experiment

³⁰The NIH could conduct analyses using their own, non-public data to determine rules that are most predictive of successful research outcomes. These analysis could mirror recent studies of journal article peer review (Card and DellaVigna 2017; Card et al. 2020), bearing in mind the greater challenge faced by the NIH when judging research proposals rather than completed research papers.

tation may offer additional insights and seems particularly well suited to the newly created Technology Innovation and Partnerships Directorate at the NSF. Shrinking research budgets, concerns about the technological competitiveness of the United States, and global declines in research productivity all underscore the need for more formal examinations of the policies and programs that ultimately shape research portfolios.

References

- Ackermann, Robert. 1986. "Consensus and Dissensus in Science." PSA: Proceedings of the Biennial Meeting of the Philosophy of Science Association 88 (2):99–105.
- Aghion, Philippe and Peter Howitt. 1992. "A Model of Growth Through Creative Destruction." *Econometrica* 60 (2):323–351.
- Allenby, Greg M and Peter E Rossi. 1991. "There Is No Aggregation Bias: Why Macro Logit Models Work." Journal of Business & Economic Statistics 9 (1):1–14.
- Arrow, K. J. 1972. Economic Welfare and the Allocation of Resources for Invention, chap. 13. London: Macmillan Education UK, 219–236. URL https://doi.org/10. 1007/978-1-349-15486-9_13.
- Azoulay, Pierre, Christian Fons-Rosen, and Joshua S Graff Zivin. 2019. "Does Science Advance One Funeral at a Time?" American Economic Review 109 (8):2889–2920.
- Azoulay, Pierre, Joshua S Graff Zivin, Danielle Li, and Bhaven N Sampat. 2019. "Public R&D Investments and Private-Sector Patenting: Evidence From NIH Funding Rules." *Review of Economic Studies* 86 (1):117–152.
- 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):527–554.
- ———. 2012. "NIH Peer Review: Challenges and Avenues for Reform." National Bureau of Economic Research Working Paper 18116.
- 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):2765–2783.
- Buck, Stuart. 2022. "Reforming Peer Review at NIH." https://goodscienceproject.org/ articles/reforming-peer-review-at-nih. Accessed: 2023-04-07.
- Bush, Vannevar. 2020. Science, the Endless Frontier. Princeton: Princeton University Press. URL https://doi.org/10.1515/9780691201658.
- Card, David and Stefano DellaVigna. 2017. "What do Editors Maximize? Evidence From Four Leading Economics Journals." *NBER Working Paper* 23282.
- Card, David, Stefano DellaVigna, Patricia Funk, and Nagore Iriberri. 2020. "Are referees and editors in economics gender neutral?" The Quarterly Journal of Economics 135 (1):269–327.
- Carson, Richard T, Joshua S Graff Zivin, Jordan J Louviere, Sally Sadoff, and Jeffrey G Shrader. 2022. "The Risk of Caution: Evidence From an Experiment." *Management Science* 68 (12):8515–9218.
- Clarke, Arthur C. 1962. "Hazards of Prophecy: The Failure of Imagination." *Profiles of the Future* 6 (36):1.
- Cox, Nicholas. 2016. "DIPTEST: Stata Module to Compute Dip Statistic to Test for Unimodality." https://ideas.repec.org/c/boc/bocode/s456998.html. Accessed: 2023-06-19.

- Fallows, James. 2014. "When Will Genomics Cure Cancer?" https://www.theatlantic. com/magazine/archive/2014/01/when-will-genomics-cure-cancer/355739/.
- Fiebig, Denzil G, Michael P Keane, Jordan Louviere, and Nada Wasi. 2010. "The Generalized Multinomial Logit Model: Accounting for Scale and Coefficient Heterogeneity." *Marketing Science* 29 (3):393–421.
- Franzoni, Chiara and Paula Stephan. 2022. "Uncertainty and Risk-Taking in Science: Meaning, Measurement and Management in Peer Review of Research Proposals." *Research Policy* :104706.
- Gius, Luca. 2024. "Disagreement Predicts Startup Success: Evidence from Venture Competitions." Working Paper :43.
- Goldstein, Anna and Michael Kearney. 2017. "Uncertainty and Individual Discretion in Allocating Research Funds." SSRN Working Paper 3012169.
- Greenblatt, Wesley H., Suman K. Maity, Roger P. Levy, and Pierre Azoulay. 2024. "Does Grant Peer Review Penalize Scientific Risk Taking? Evidence from the NIH." *Working Paper* :78.
- Gu, Jiaying and Roger Koenker. 2023. "Invidious comparisons: Ranking and selection as compound decisions." *Econometrica* 91 (1):1–41.
- Gu, Yuanyuan, Arne Risa Hole, and Stephanie Knox. 2013. "Fitting the Generalized Multinomial Logit Model in Stata." The Stata Journal 13 (2):382–397.
- Guthrie, Susan, Ioana Ghiga, and Steven Wooding. 2018. What Do We Know About Grant Peer Review in the Health Sciences? An Updated Review of the Literature and Six Case Studies. Santa Monica: RAND Corporation.
- Hartigan, John A and Pamela M Hartigan. 1985. "The Dip Test of Unimodality." Annals of Statistics 13:70–84.
- Hong, Lu and Scott E Page. 2004. "Groups of diverse problem solvers can outperform groups of high-ability problem solvers." *Proceedings of the National Academy of Sciences* 101 (46):16385–16389.
- Jerrim, John and Robert de Vries. 2020. "Are Peer-Reviews of Grant Proposals Reliable? An Analysis of Economic and Social Research Council (ESRC) Funding Applications." *The Social Science Journal* 60 (1):1–19.
- Kalyani, Aakash. 2022. "The Creativity Decline: Evidence From US Patents." SSRN Working Paper 4318158.
- Kolata, Gina. 2009. "Grant System Leads Cancer Researchers to Play It Safe." https:// www.nytimes.com/2009/06/28/health/research/28cancer.html. Accessed: 2023-04-07.
- Krieger, Joshua, Ramana Nanda, Josh Lerner, and Ahmed Tahoun. 2022. "Are Transformational Ideas Harder to Fund? Resource Allocation to R&D Projects at a Global Pharmaceutical Firm." Havard Business School Working Paper 23-014.
- Lane, Jacqueline N, Zoe Szajnfarber, Jason Crusan, Michael Menietti, and Karim R Lakhani. 2022a. "Are Experts Blinded by Feasibility? Experimental Evidence From a NASA Robotics Challenge." No. 22-071.

- Lane, Jacqueline N, Misha Teplitskiy, Gary Gray, Hardeep Ranu, Michael Menietti, Eva C Guinan, and Karim R Lakhani. 2022b. "Conservatism Gets Funded? A Field Experiment on the Role of Negative Information in Novel Project Evaluation." *Management Science* 68 (6):4478–4495.
- Lauer, Mike. 2023. "FY 2022 by the Numbers: Extramural Grant Investments in Research." https://nexus.od.nih.gov/all/2023/03/01/ fy-2022-by-the-numbers-extramural-grant-investments-in-research/. Accessed: 2023-04-07.
- Lee, Carole J, Cassidy R Sugimoto, Guo Zhang, and Blaise Cronin. 2013. "Bias in Peer Review." Journal of the American Society for Information Science and Technology 64 (1):2–17.
- Li, Danielle. 2017. "Expertise Versus Bias in Evaluation: Evidence From the NIH." American Economic Journal: Applied Economics 9 (2):60–92.
- Li, Danielle and Leila Agha. 2015. "Big Names or Big Ideas: Do Peer-Review Panels Select the Best Science Proposals?" *Science* 348 (6233):434–438.
- Louviere, Jordan J, Terry N Flynn, and Anthony Alfred John Marley. 2015. Best-Worst Scaling: Theory, Methods and Applications. Cambridge University Press.
- Louviere, Jordan J, David A Hensher, and Joffre D Swait. 2000. Stated Choice Methods: Analysis and Applications. Cambridge University Press.
- Lucas, Robert E. 1988. "On the Mechanics of Economic Development." Journal of Monetary Economics 22 (1):3–42.
- Manso, Gustavo. 2011. "Motivating Innovation." Journal of Finance 66 (5):1823-1860.
- March, James G. 1991. "Exploration and Exploitation in Organizational Learning." Organization Science 2 (1):71–87.
- Miller, Roger. 1995. "Applying Quality Practices to R&D." Research-Technology Management 38 (2):47–54.
- Mokyr, Joel. 1992. The Lever of Riches: Technological Creativity and Economic Progress. Oxford University Press. URL https://doi.org/10.1093/acprof:oso/9780195074772. 001.0001.
- Moses, Hamilton, E Ray Dorsey, David HM Matheson, and Samuel O Thier. 2005. "Financial Anatomy of Biomedical Research." *Journal of the American Medical Association* 294 (11):1333–1342.
- Murray, Fiona. 2013. "Evaluating the Role of Science Philanthropy in American Research Universities." *Innovation Policy and the Economy* 13 (1):23–60.
- Myers, Kyle and Wei Yang Tham. 2023. "Money, Time, and Grant Design." *Working Paper* :33.
- Myers, Kyle R, Wei Yang Tham, Yian Yin, Nina Cohodes, Jerry G Thursby, Marie C Thursby, Peter Schiffer, Joseph T Walsh, Karim R Lakhani, and Dashun Wang. 2020. "Unequal Effects of the COVID-19 Pandemic on Scientists." Nature Human Behaviour 4 (9):880–883.

- National Science Foundation (NSF). 2023. "National Patterns of R&D Resources: 2020–21 Data Update." https://ncses.nsf.gov/pubs/nsf23321. Accessed: 2023-04-07.
- Nelson, Richard R. 1959. "The Simple Economics of Basic Scientific Research." Journal of Political Economy 67 (3):297–306.
- Nelson, Richard R and Edmund S Phelps. 1966. "Investment in Humans, Technological Diffusion, and Economic Growth." American Economic Review 56 (1/2):69–75.
- Nicholson, Joshua M and John PA Ioannidis. 2012. "Conform and Be Funded." *Nature* 492 (7427):34–36.
- NIH. 2008. "Enhancing Peer Review: The NIH Announces New Scoring Procedures for Evaluation of Research Applications Received for Potential FY2010 Funding." https:// grants.nih.gov/grants/guide/notice-files/not-od-09-024.html. Accessed: 2023-04-07.
- ———. 2022. "NIH Extramural & Intramural Funding: FY 2022 Operating Plan." https: //report.nih.gov/nihdatabook/report/283. Accessed: 2024-03-18.
- Nijstad, Bernard A, Floor Berger-Selman, and Carsten KW De Dreu. 2014. "Innovation in Top Management Teams: Minority Dissent, Transformational Leadership, and Radical Innovations." *European Journal of Work and Organizational Psychology* 23 (2):310–322.
- Ogden, Thomas E and Israel A Goldberg. 2002. Research Proposals: A Guide to Success. San Diego, California: Academic Press.
- Pier, Elizabeth L, Markus Brauer, Amarette Filut, Anna Kaatz, Joshua Raclaw, Mitchell J Nathan, Cecilia E Ford, and Molly Carnes. 2018. "Low Agreement Among Reviewers Evaluating the Same NIH Grant Applications." Proceedings of the National Academy of Sciences 115 (12):2952–2957.
- Pier, Elizabeth L, Joshua Raclaw, Anna Kaatz, Markus Brauer, Molly Carnes, Mitchell J Nathan, and Cecilia E Ford. 2017. "Your comments are meaner than your score': score calibration talk influences intra-and inter-panel variability during scientific grant peer review." Research Evaluation 26 (1):1–14.
- Romer, Paul M. 1990. "Endogenous Technological Change." Journal of Political Economy 98 (5, Part 2):S71–S102.
- Sampat, Bhaven. 2023. "The History and Political Economy of NIH Peer Review." Tech. rep., Brookings. URL https://www.brookings.edu/wp-content/uploads/2023/05/ SampatFinal-3.pdf.
- Science News Staff. 2022. "Research Gets a Boost in Final 2023 Spending Agreement." https://www.science.org/content/article/research-gets-boost-final-2023-spendingagreement. Accessed: 2023-06-19.
- Sen, Avery. 2014. "Totally Radical: From Transformative Research to Transformative Innovation." Science and Public Policy 41 (3):344–358.
- Sinkjaer, Thomas. 2018. "Fund Ideas, Not Pedigree, to Find Fresh Insight." Nature 555 (7697):143–144.
- Swait, Joffre and Jordan Louviere. 1993. "The Role of the Scale Parameter in the Estimation

and Comparison of Multinomial Logit Models." *Journal of Marketing Research* 30 (3):305–314.

- Tawfik, Gehad Mohamed, Hoang Thi Nam Giang, Sherief Ghozy, Ahmed M Altibi, Hend Kandil, Huu-Hoai Le, Peter Samuel Eid, Ibrahim Radwan, Omar Mohamed Makram, Tong Thi Thu Hien et al. 2020. "Protocol registration issues of systematic review and meta-analysis studies: a survey of global researchers." BMC medical research methodology 20 (1):213.
- Train, Kenneth E. 2009. *Discrete Choice Methods with Simulation*. Cambridge University Press.
- Whitley, Richard and Jochen Gläser, editors. 2007. The Changing Governance of the Sciences, vol. 26. Dordrecht: Springer.

For Online Publication

Appendix for Choose Your Moments: NIH Peer Review and Scientific Risk Taking

A Randomization Checks

	(1)	(2)
	Study 1 indicator	Study 1 indicator
Fraction BDCN	-0.018	-0.078
	(0.019)	(0.061)
Fraction CB	0.024	0.014
	(0.020)	(0.063)
Fraction OBT/OTC	-0.0034	-0.057
	(0.018)	(0.061)
Total funding	0.00061	-0.0030
	(0.0010)	(0.0045)
Unique projects	-0.0021	0.014
	(0.0030)	(0.012)
Total projects	-0.00028	-0.0013
	(0.00063)	(0.0030)
F-stat	1.21	0.74
p-value	0.30	0.61
Observations	6678	563

Table A1: Randomization Balance: Omnibus F-test

This table shows randomization checks for the randomization into Study 1 versus Study 2. Column 1 shows the check for the total sample of potential participants. Column 2 shows the check for the sample that completed the study. Heteroskedasticity robust standard errors are in parentheses: p < .10, ** p < .05, *** p < .01.

B Robustness and Sensitivity Results

	(1) Project choice	(2) Project choice	(3) Project choice	(4) Project choice	(5) Project choice
Avg. score	0.77^{***}	0.83***	0.93***	0.67***	0.81***
Score variance	(0.037)	(0.041) 0.075^{***}	(0.058) 0.11^{***}	(0.062) 0.092 (0.061)	(0.044) 0.076^{**}
Avg. score \times Score variance		(0.029) 0.032 (0.026)	(0.035)	(0.061)	(0.034)
Score skew		(0.020)	0.14^{**}		
Score kurtosis			(0.003) (0.034) (0.056)		
Count of lowest scores			(0.000)	-0.14^{***}	
Count of highest scores				(0.052) 0.11^{***} (0.044)	
At least one bottom score				(0.044)	-0.027
At least one top score					$\begin{array}{c} (0.030) \\ 0.099 \\ (0.077) \end{array}$
Clusters N	313 11268	$313 \\ 11268$	$\begin{array}{c} 313\\11086 \end{array}$	$313 \\ 11268$	313 11268

Table A2: Additional Preferences for Score Statistics

This table shows results from estimating Equation (1) on the baseline sample. The dependent variable is an indicator for whether the participant chose a given project. All models include participant fixed effects. "Score skew" is the skew of project scores, "score kurtosis" is the kurtosis of project scores, "count of lowest/highest scores" is the number projects given the maximum or minimum possible scores (given the boundedness of project scores, this is mechanically correlated with score variance), and "at least one bottom/top score" is an indicator for whether there was at least one project score at the maximum or minimum score value. Standard errors are clustered at the participant level: * p < .10, ** p < .05, *** p < .01.

	(1)	(2)	(3)
	rank	rank	rank
1			
Avg. score	1.81***	1.81***	2.30***
-	(0.11)	(0.11)	(0.16)
Score variance	0.17^{***}	0.16^{***}	0.24^{***}
	(0.060)	(0.060)	(0.075)
MeSH match		0.047	0.093
		(0.042)	(0.059)
2			
Avg. score	1.14***	1.14***	1.44***
	(0.082)	(0.082)	(0.11)
Score variance	0.11**	0.11**	0.17***
	(0.053)	(0.053)	(0.064)
MeSH match		0.051	0.092^{*}
		(0.036)	(0.050)
3			
Avg. score	0.60***	0.60***	0.77***
	(0.058)	(0.058)	(0.079)
Score variance	-0.029	-0.030	-0.0019
	(0.050)	(0.050)	(0.058)
MeSH match		0.053	0.091^{*}
		(0.036)	(0.051)
Subject FEs	No	No	Yes
Clusters	313	313	313
Ν	5008	5008	5008

Table A3: Robustness: Main Results Using Multinomial Logit and Project Rank

This table shows results from estimating a multinomial logit model corresponding to Equation (1) on the baseline sample. The dependent variable is the rank the participant gave to a given project (lower is better and the excluded category is rank 4 out of 4). Standard errors are clustered at the participant level: * p < .10, ** p < .05, *** p < .01.

	(1)	(2)
	Project choice	Project choice
Avg. score	0.82***	1.09***
	(0.040)	(0.070)
Score variance	0.083^{***}	0.12^{***}
	(0.027)	(0.034)
Sample	All obs.	Granual FEs
Clusters	356	313
Ν	12060	11268

Table A4: Robustness: Sensitivity to Sample Restrictions and Fixed Effects

This table shows results from estimating Equation (1) on the baseline sample. The dependent variable is an indicator for whether the participant chose a given project. All models include participant fixed effects. Standard errors are clustered at the participant level: * p < .10, ** p < .05, *** p < .01.

B.1 Estimates From More General Logistic Models

In the main results reported in Table 2, we are interested in the average preferences for higher mean score and higher score variance, which we estimate using the conditional logit model given in Equation 1. Here, we report results from models that relax the assumptions of the conditional logit model—allowing for subject level attribute and scale heterogeneity.

The estimating equations for these more general models are nested in the following equation:

$$\Pr(choice_{it} = j|\beta_i) = \frac{\exp(\beta'_i \mathbf{x}_{itj})}{\sum_{k=1}^{J} \exp(\beta'_i \mathbf{x}_{itk})}$$
(A-1)

where *i* indexes subject, *t* indexes the choice scenario, and *j* and *k* index the project attributes. The variable x_{itj} is a vector containing the mean and variance of project scores and β_i is the vector of individual-specific coefficients associated with these project attributes. The coefficients are defined by

$$\beta_i = \sigma_i \beta + \{\gamma + \sigma_i (1 - \gamma)\}\eta_i$$

The coefficients in this equation are a vector beta that is constant across individuals and measures the average utility weights across the sample for the different variables in x; a single parameter for the scale of the individual-level idiosyncratic error (σ_i) , which captures overall scaling of the individual's tastes; and a random vector η_i distributed multivariate normal with mean zero and variance-covariance matrix Σ , which captures latent taste heterogeneity. The parameter γ determines how much of the variance is explained by these latter two components. The setup and parameterization of these parameters follows Fiebig et al. (2010). In particular, we assume that σ_i is distributed log normal with mean $sigma + \theta' z_i$ and standard deviation τ . The parameter sigma is a normalizing constant, and z_i is a vector of subject fixed effects in this application.

The results of fitting versions of this model—a mixed random parameters model (Train 2009), a scale heterogeneity logit model (S-MNL) (Swait and Louviere 1993), and the generalized multinomial logit (G-MNL) model—are shown in Table A5. The coefficient estimates are qualitatively consistent with the baseline results in the sense that on average, participants prefer both higher mean and higher variance projects. The subject-level utility weights are predicted based on the average utility weights and heterogeneity in weights using the gmnlpred command from Gu, Hole, and Knox (2013). For example, in the G-MNL reported in Column (5), the average weight placed on project score variance is 8% as large as the weight that is placed on average project score, comparable to the 10% value derived from

Table 2.

	(1)	(2)	(3)	(4)	(5)	(6)
	Conditional	(2) Mixed	(J) Mixed	S MNL	C MNL	C MNL
		lagit	full com	0-1VIINL	G-IVIIVL	full com
	logit	logit	full corr.			full corr.
Average utility weight						
Average project score	1.0973^{***}	1.9243***	1.9377***	2.6297***	2.8723***	2.9027***
	(.0709)	(.1316)	(.1322)	(.2744)	(.3401)	(.3718)
Project score var.	0.1199***	0.1316***	0.1502**	0.2010**	0.1568**	0.1663**
5	(.0341)	(.0473)	(.0499)	(.0794)	(.0727)	(.0687)
Utilitu weight heteroger	neitu	()	()	()	()	()
Average project score	<u>-</u>	1.3947***	1.4328***		1.0076***	1.0243***
		(0998)	(1011)		(.9151)	(1807)
Project score var		0.4597^{***}	0.0949		6364***	6489***
i iojeet seere var.		(0559)	(0633)		(0022)	(0815)
$\Delta v \sigma \times v \sigma$		(.0000)	(.0055)		(.0522)	2166
Avg × var.			(0562)			(1711)
Additional manages stores			(.0302)			(.1(11))
Additional parameters				1 /100***	1 0050***	1 9197***
au				1.4100	1.2858	1.313(
				(.09950)	(.1326)	(.1589)
γ					$.2561^{***}$.2930***
					(.0667)	(.0937)
Log likelihood	-3179.69	-2876.61	-2875.43	-2882.64	-2844.25	-2843.58

Table A5: Generalizations of the Conditional Logit Model

This table shows results from estimating versions of Equation (A-1) on the baseline sample. The dependent variable is an indicator for whether the participant chose a given project. All models include participant fixed effects. Standard errors are clustered at the participant level: * p < .10, ** p < .05, *** p < .01.

C Additional Figures and Tables

	(1)	(2)
	Project choice	Project choice
Avg. score	0.96***	1.02^{***}
	(0.067)	(0.19)
Score variance	0.097^{***}	0.054
	(0.035)	(0.12)
Score skew	0.17***	0.73
	(0.050)	(1.17)
Skew sample	Neg. skew	Pos. skew
Clusters	313	230
Ν	8879	2071

Table A6: Effects of Positive and Negative Skewness

This table shows results from estimating Equation (1) on the baseline sample. The dependent variable is an indicator for whether the participant chose a given project. All models include participant fixed effects. Standard errors are clustered at the participant level: * p < .10, ** p < .05, *** p < .01.

Figure A1: Funding Reversals Under NIH and Scientist Ranking Procedures: Separate Estimates by Study



Notes: This figure shows the reversal rate for project funding as a function of the payline (fraction of projects that get funded) for four different sets of project scores (Study 1 from this paper, the two studies from Lane et al. (2022b), and Pier et al. (2018)). The reversal rate is the fraction of studies that change when they are funded under a mean-only ranking versus under the mean and variance-based ranking. Panel (a) shows the reversal rate when using the estimated preferences from the baseline results using the full sample of scientists. Panel (b) shows the reversal rate when using the preferences of scientists who are relative experts (MeSH match 1 standard deviation higher than average).

	(1)
	Project choice
Score of 1	-0.20***
	(0.026)
Score of 2	-0.18***
	(0.022)
Score of 3	-0.092***
	(0.021)
Score of 4	0.014
C C C	(0.019)
Score of 6	0.067***
0 6 7	(0.023)
Score of 7	0.074^{***}
Correct of O	(0.011)
Score of 8	(0.13^{++})
Score of 0	(0.014) 0.22***
Score of 9	$(0.22^{+1.1})$
	(0.010)
Clusters	313
Ν	11268

Table A7: Preferences for Project Scores

This table shows results from estimating Equation (1) on the baseline sample where the right-hand-side variables (project attributes) are each level of possible project scores. The dependent variable is an indicator for whether the participant chose a given project. All models include participant fixed effects. Standard errors are clustered at the participant level: * p < .10, ** p < .05, *** p < .01.

	(1)	(2)	(3)	(4)	(5)	(6)
	Project	Project	Project	Project	Project	Project
	choice	choice	choice	choice	choice	choice
Avg. score	0.81***	0.83***	0.78***	0.91***	0.74***	0.84***
	(0.069)	(0.050)	(0.051)	(0.066)	(0.073)	(0.067)
Score variance	0.027	0.11^{***}	0.082^{**}	0.088^{*}	0.011	0.12^{***}
	(0.048)	(0.033)	(0.034)	(0.046)	(0.045)	(0.046)
Subgroup	Female	Male	Risk	Not risk	Least	Most
			averse	averse	experience	experience
Clusters	91	222	211	102	103	105
Ν	3276	7992	7596	3672	3708	3780

Table A8: Heterogeneity: Demographics, Preferences, and Experience

This table shows results from estimating Equation (1) on the baseline sample. The dependent variable is an indicator for whether the participant chose a given project. All right-hand-side variables are standardized. The models include participant fixed effects. Standard errors are clustered at the participant level: * p < .10, ** p < .05, *** p < .01.

D Experiment Instructions and Instruments

Below are screenshots of the main instructions and choice scenarios shown to the participants. The full experimental instrument can be found at the following links: link for Study 1 and link for Study 2. Note that the full instruments includes all versions. In practice, a participant was randomized into seeing only four project choice scenarios in Study 1 or two portfolio choice scenarios in Study 2.

Figure A2: Experimental Instrument: Welcome Screen

Survey on Biomedical Research Funding Allocation Decisions

Welcome to our National Science Foundation sponsored survey and thank you for agreeing to participate. You were included in our survey based, in part, on your recent experience as a principal investigator on a National Institutes of Health (NIH) grant. Our study is designed to examine how researchers think about which grant proposals should receive funding. A report summarizing our findings will be provided to the National Institute of Health, the National Science Foundation, and other major institutions that fund scientific research. Participants interested in receiving project reports will be given the opportunity to sign up at the end of the survey. No individually identifying information will be kept with this survey dataset or included in any reports or publications.

For roughly the next thirty minutes, we want you to assume the role of the director of a special NIH program who has to decide on how to allocate scarce funding across a range of potential research projects. To be clear, this is different from the role that you have likely played as a member of an NIH study section. While we want you to draw on that experience, we do not want you to feel constrained by it. One of the main things we seek to learn is how researchers would pick projects if current rules were not in place. You will be provided with evaluation scores from a complete study section meeting (albeit one a bit more stylized than the usual NIH review process). We will then ask you to decide which project(s) you would want to see funded based on that input.

Project descriptions are based on real NIH grant applications. Reviewer evaluation scores have been modified by the research team to help facilitate our statistical analysis of how biomedical researchers, such as yourself, trade off various attributes of projects when deciding what to fund. We will begin by asking you to assemble a portfolio of projects from a selected list. After this task, you will be asked some additional questions to help us better understand your professional background and decision-making processes.

Please read all instructions carefully and take your time in answering the questions.

Notes: This figure shows the welcome screen that greeted participants.

Figure A3: Experimental Instrument: Choice Scenario Main Instructions

Instructions

For the next four questions, your role is that of a program director with limited funds for funding projects. You will be asked to consider sets of four research project proposals (A, B, C, D).

- Each proposal has received a rating on a scale from 1 to 9 (with 1 being the top rating) by 30 scientific experts on your advisory board, all of whom are unaffiliated with the projects under consideration.
- For each set of four proposals, you will be provided with two tables of information to help in your funding decision. In the first table, you will be provided the titles of each proposal. You can also review the individual proposal abstract and a graph of the reviewer scores by hovering over the proposal you are interested in.
- In the second table, you will be provided information on the scoring of each proposal. Each column represents one proposal, with the value in each row referring to the number of reviewers who gave that score to the proposal. The average of the reviewers' scores for each proposal and the standard deviation are also displayed toward the bottom of each proposal's column.
- After considering the abstracts and scoring information, you will be asked to indicate the proposal you most and least prefer to fund. Then, from the remaining two proposals, you will be asked which you most prefer to fund. Your responses will thus provide a complete ranking of the four proposals. Remember that you need not be constrained by current NIH funding rules and thus should feel free to use any information that you deem relevant to make your funding decisions.
- The order in which proposals appears has been randomized, and to simplify your task, all proposals have the same cost.

Notes: This figure shows the main instructions that participants were shown prior to making any choices in the first experiment.

Figure A4: Experimental Instrument: Choice Scenario Example

Please indicate which of the four proposals you would **most** and **least** like to fund. The tables below show the proposal titles and the number of reviewers (out of your panel of 30) who gave the proposals each score, 1 (best) through 9 (worst), along with the average score and standard deviation of the scores for each proposal. Hover over each proposal in the first table to view its abstract and see a graph of reviewer scores.

Proposal	Title
<u>Proposal A</u>	Bringing CLARITY to EAE.
<u>Proposal B</u>	Congenital Myasthenic Syndromes.
<u>Proposal C</u>	Mechanisms Regulating Cerebral Arteriogenesis and Neurorestoration.
<u>Proposal D</u>	Mechanisms of cognitive deficits after seizures in rats with brain malformations.

Score	Proposal A	Proposal B	Proposal C	Proposal D
1 (Best)	3	9	9	3
2	6	12	6	6
3	6	6	6	6
4		3	1	2
5	6		6	6
6				2
7	3			3
8	4		2	1
9 (Worst)	2			1
Average Score	4.47	2.10	2.97	4.03
Standard Deviation	2.61	0.96	2.01	2.19
Which proposal you would most like to fund.				
Which proposal you would least like to fund.				

Number of Reviewers per Score and Proposal Score Statistics

Notes: This figure shows an example of the choices participants were asked to make. In the top box, the blue, underlined proposal links would show proposal abstracts if the participant hovered over them. After choosing their top and bottom ranked projects, in a follow-up screen participants were asked to rank the remaining two projects.

Figure A5: Experimental Instrument: Portfolio Choice Scenario Main Instructions

Instructions

In the next two questions, you will be asked to put together a portfolio of research project proposals to fund. Your role is that of a program director with limited funds for funding projects.

- In each of the two questions, you will be asked to consider sets of ten research project proposals (denoted A through J).
- Each proposal has received a rating on a scale from 1 to 9 (with 1 being the top rating) by 30 scientific experts on your advisory board, all of whom are unaffiliated with the projects under consideration.
- For each set of ten proposals, you will be provided with two tables of information to help in your funding decision. In the first table, you will be provided the titles of each proposal. You can also review the individual proposal abstract and a graph of the reviewer scores by hovering over the proposal you are interested in.
- In the second table, you will be provided information on the scoring of each proposal. Each column represents one proposal, with the value in each row referring to the number of reviewers who gave that score to the proposal. The average of the reviewers' scores for each proposal and the standard deviation are also displayed toward the bottom of each proposal's column.
- After considering the abstracts and scoring information, you will be asked to indicate which of the ten proposals you would like to fund with your limited budget. Remember that you need not be constrained by current NIH funding rules and thus should feel free to use any information that you deem relevant to make your funding decisions.
- Each proposal has the same cost, which is displayed along with your budget. The budget of the portfolio will influence what research projects you are able to fund. The portfolio you choose must cost the same or less than your total budget. Any leftover funds from each question will be returned to NIH headquarters and will be unavailable for project funding by this study section.
- The order in which proposals appear has been randomized.

Notes: This figure shows the main instructions that participants were shown prior to making any portfolio choices in the second experiment.

Question 1

Your NIH program budget to fund this portfolio of proposals is \$8 million dollars and the cost of each proposal is \$2 million dollars.

Check the boxes of the projects you would like to fund located below each proposal. As you enter proposals, the costs of each one will be subtracted from your NIH program budget and your remaining funds will be displayed. Due to your budget, you may not be able to fund all desired projects. You can swap proposals in and out by selecting and unselecting boxes until you are satisfied with your choices.

The tables below show the proposal titles and the number of reviewers (out of your panel of 30) who gave the proposals each score, 1 (best) through 9 (worst), along with the average score and standard deviation of the scores for each proposal. Hover over each proposal in the first table to view its abstract and see a graph of reviewer scores.

Proposal	Title
<u>Proposal</u> <u>A</u>	Bringing CLARITY to EAE.
<u>Proposal B</u>	Congenital Myasthenic Syndromes.
<u>Proposal C</u>	Mechanisms Regulating Cerebral Arteriogenesis and Neurorestoration.
<u>Proposal</u> D	Mechanisms of cognitive deficits after seizures in rats with brain malformations.
<u>Proposal E</u>	Neuregulin-1 decreases endothelial hyper-permeability and microthrombosis after traumatic brain injury.
Proposal F	Doublecortin in Neuronal Migration.
<u>Proposal</u> <u>G</u>	Anesthesia-induced learning deficiency and brain hyperoxia.
<u>Proposal</u> <u>H</u>	Mapping Brain Structure to Function in Euthymic Subjects with Bipolar Disorder.
<u>Proposal I</u>	A High Content Screening Approach for the Retinal Degenerative Diseases.
<u>Proposal J</u>	CB1 Receptor PET Imaging Reveals Gender Differences in PTSD.

Remaining Budget: \$8 million

Num	ber of Review	ers per Score a	and Proposal	Score Statistics

Coore	Proposal									
Score	A	в	С	D	E	F	G	н	I	J
1 (Best)	15	9		9				12		15
2		6	15	6		9	15		9	
3		6		6	30	6		6	6	
4	2					2			2	5
5		6	15	6		3	15	12	3	
6	5			1		2			6	3
7	1	1				1			1	1
8	5	1				2			2	3
9 (Worst)	2	1		2		5			1	3
Average Score	3.93	3.10	3.50	3.10	3.00	4.63	3.50	3.00	4.23	3.70
Standard Deviation	3.18	2.22	1.53	2.25	0.00	2.67	1.53	1.82	2.14	3.09
Portfolio Selection										

If funding were slashed so that you had to *drop one project* from your portfolio, which would you *drop*?

	Proposal									
	A	B	C	D	E	F	G	H	I	J
Project you would drop from portfolio.										

Notes: This figure shows an example of the choices participants were asked to make in Study 2, the portfolio choice scenario. The budget remaining changed as projects were selected. After making these choices, the participant was prompted to say which project they would add if the budget were expanded.