

Week 13: The Economics of Science

Carolyn Stein

Econ 220C: Topics in Industrial Organization

Science as a non-market incentive

Scientific norms and their effects

Mertonian origins

Hill and Stein (2020)

Hill and Stein (2022)

Basic vs. applied science

- ▶ **Basic science** seeks to expand human knowledge, but not to create or invent something. There is no obvious commercial value to the result of basic research
- ▶ **Applied science** seeks to solve practical problems and often yields something that is commercially valuable
- ▶ Basic research is important: “People cannot foresee the future well enough to predict what’s going to develop from basic research. If we only did applied research, we would still be making better spears.” – George Smoot
- ▶ But how do we incentivize people to produce it?

Basic vs. applied science: an example

Basic science

Francisco Mojica studied bacteria in Spanish salt flats in the 80s and 90s. He noticed odd bits of repeated DNA in these bacteria...which paved the way for CRISPR



molecular
microbiology

Long stretches of short tandem repeats are present in the largest replicons of the *Archaea Haloferax mediterranei* and *Haloferax volcanii* and could be involved in replicon partitioning

F.J.M. Mojica, C. Ferrer, G. Juez, F. Rodriguez-Valera

First published: July 1995 | https://doi.org/10.1111/j.1365-2958.1995.mmi_17010085.x | Citations: 205

Applied science

Vertex pharmaceuticals has developed a CRISPR gene editing treatment to cure patients with sickle-cell anemia, which is currently under FDA review



Academic freedom can be a strong motivator

- ▶ Stern (2004) “Do Scientists Pay to Be Scientists?” studies whether researchers take a pay cut to be given more scientific freedom
- ▶ Surveys biology post-docs with multiple job offers and collects characteristics of the jobs
 - ▶ Salary
 - ▶ Measures of scientific freedom (allowed to publish discoveries, allowed to continue postdoc projects, whether there are incentives to publish)
- ▶ Argues that all job offers should be roughly similarly attractive (formal offers only issued if candidate is serious)
- ▶ Can the run a hedonic regression with individual fixed effects

Researchers do value academic freedom

Results suggest postdocs accept a 20% pay cut in exchange for more academic freedom

Table 3 Hedonic Wage Regression: Overall Sample Dependent Variable = LN(SALARY), # of Observations = 121

	Permission to publish			Combination model		Science index model		
	(3-1)		(3-2)	(3-3)	(3-4)		(3-5)	
	Baseline (NO FE)	Baseline (w/FE)	Full model (w/FE)	Full model (w/FE)	Full Model (w/FE)	Full Model (w/FE)	(3-6)	
PERMIT_PUB	0.027 (0.186)		-0.266 (0.114)	-0.191 (0.105)	-0.089 (0.103)			
CONTINUE RESEARCH					-0.134 (0.060)			
INCENT_PUB					-0.036 (0.028)			
SCIENCE INDEX						-0.114 (0.053)	-0.078 (0.057)	
EQUIPMENT					0.063 (0.033)	0.057 (0.030)	0.053 (0.031)	
CONTROLS								
PROMOTION				0.041 (0.025)	0.046 (0.021)	0.042 (0.021)	0.031 (0.023)	
STOCK_DUMMY				0.196 (0.085)	0.234 (0.074)	0.260 (0.067)	0.190 (0.077)	
ACCEPTED JOB				-0.013 (0.040)	0.002 (0.043)	-0.0001 (0.043)	-0.002 (0.044)	
JOBTYPE CONTROLS	no	no	yes		no	no	yes (5)	
Individual fixed effects	no	yes	(5; Sig.)	yes	yes	yes	yes	
R-squared	0.001	0.915	0.955	0.958	0.954	0.958	0.958	

Notes: Only persons with multiple job offers are included.

How else do we encourage basic research?

Many scientific norms can be viewed through the lens of providing incentives to engage in basic research:

- ▶ Grants
- ▶ Prizes
- ▶ Eponymy

All designed to compensate researchers. Maybe not with profits, but with credit, acclaim, etc.

“My love of natural science...has been much aided by the ambition to be esteemed by my fellow naturalists” – Charles Darwin

Science as a non-market incentive

Scientific norms and their effects

Mertonian origins

Hill and Stein (2020)

Hill and Stein (2022)

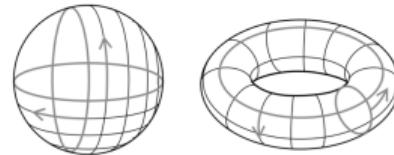
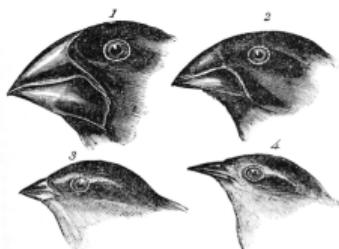
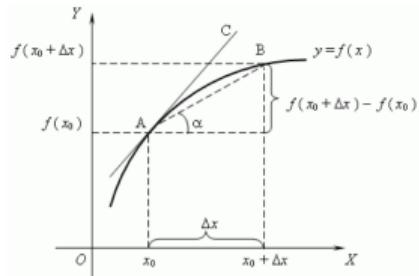
The importance of credit and recognition

"In short, property rights in science become whittled down to just this one: the **recognition by others** of the scientist's distinctive part in having brought the result into being."
- Robert K. Merton (1957)



Priority in scientific discovery

- ▶ **Priority:** Credit given to the individual who *first* makes a scientific discovery.
- ▶ If being first yields more credit, not surprising that there are often fierce disputes over priority
- ▶ Notable scientific races and priority disputes:
 - ▶ Newton versus Leibniz - Calculus
 - ▶ Darwin versus Wallace - Natural Selection and Evolution
 - ▶ Perelman versus Yau, Zhu, and Cao - Proof of the Poincaré Conjecture
- ▶ Merton (1961) assembles 264 cases of “multiple discovery”



Science as a non-market incentive

Scientific norms and their effects

Mertonian origins

Hill and Stein (2020)

Hill and Stein (2022)

Scooped! Estimating rewards for priority in science

1. What is the causal effect of getting scooped?

- ▶ Short-run effect on project: Publication, journal placement, and citations
- ▶ Long-run effect on career: Future productivity of scientists

Scooped! Estimating rewards for priority in science

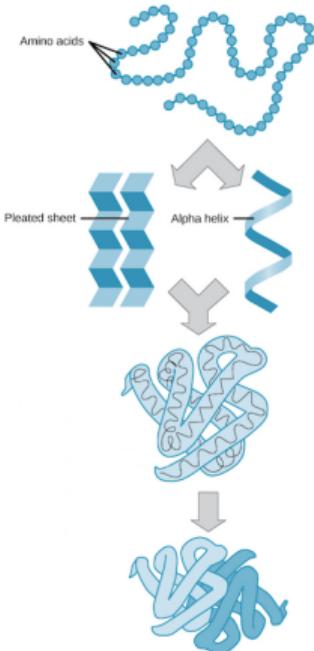
1. What is the causal effect of getting scooped?
 - ▶ Short-run effect on project: Publication, journal placement, and citations
 - ▶ Long-run effect on career: Future productivity of scientists
2. Does the priority reward system reinforce inequality in science? (Matthew Effect)
 - ▶ What drives citations: being first or being famous?

Key empirical challenges

1. Need a setting with well-defined problems and “one right answer.”
2. Need an objective measure of scientific proximity.
3. Need a view of potential abandonments prior to publication.

What is structural biology?

- ▶ Structural biologists determine the molecular structure of proteins, DNA, and RNA.
- ▶ Proteins carry out most of the functions within cells, and often "form determines function."
- ▶ Structures are solved by X-ray crystallography. Successful experiments result in diffraction data and a model that describes the protein shape.



The Protein Data Bank

- ▶ The Protein Data Bank (PDB) contains structural data of 100,000+ proteins and meta-data about projects.
- ▶ Major scientific journals require scientists to submit their structure data to the PDB before publication.
- ▶ All structures are deposited confidentially a few months before article publication.
- ▶ Bioinformatics algorithm links projects with identical biological features.

PDB example: Cas-9

Biological Assembly 1

4CNP unique structure ID

Crystal structure of *S. pyogenes* Cas9

DOI: [10.2210/pdb4CNP/pdb](https://doi.org/10.2210/pdb4CNP/pdb)

Classification: [HYDROLASE](#)

Organism(s): [Streptococcus pyogenes serotype M1](#)

Expression System: [Escherichia coli BL21\(DE3\)](#)

Deposited: 2014-01-16 Released: 2014-02-12 key dates

Deposition Author(s): [Jinek, M.](#), [Jiang, F.](#), [Taylor, D.W.](#), [Sternberg, S.H.](#), [Kaya, E.](#), [Ma, E.](#), [Anders, C.](#), [Hauer, M.](#), [Zhou, K.](#), [Lin, S.](#), [Kaplan, M.](#), [Iavarone, A.T.](#), [Charpentier, E.](#), [Nogales, E.](#), [Doudna, J.A.](#)

Literature Download Primary Citation ▾

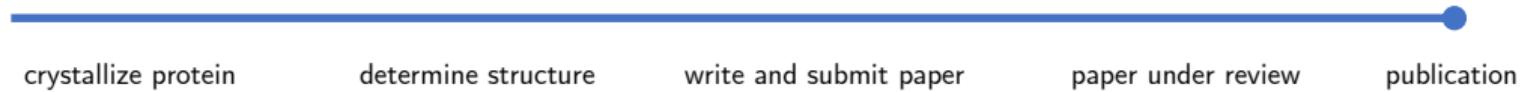
Structures of Cas9 Endonucleases Reveal RNA-Mediated Conformational Activation.
[Jinek, M.](#), [Jiang, F.](#), [Taylor, D.W.](#), [Sternberg, S.H.](#), [Kaya, E.](#), [Ma, E.](#), [Anders, C.](#), [Hauer, M.](#), [Zhou, K.](#), [Lin, S.](#), [Kaplan, M.](#), [Iavarone, A.T.](#), [Charpentier, E.](#), [Nogales, E.](#), [Doudna, J.A.](#)
(2014) *Science* **343**: 47997

PubMed: [24505130](#) [Search on PubMed](#) [Search on PubMed Central](#)
DOI: [10.1126/science.1247997](#)

Primary Citation of Related Structures:
[4OGE](#), [4OGC](#), [4CMQ](#)

PubMed Abstract:
Type II CRISPR (clustered regularly interspaced short palindromic repeats)-Cas (CRISPR-associated) systems use an RNA-guided DNA endonuclease, Cas9, to generate double-strand breaks in invasive DNA during an adaptive bacterial immune response. Cas9 h ...

Project timeline



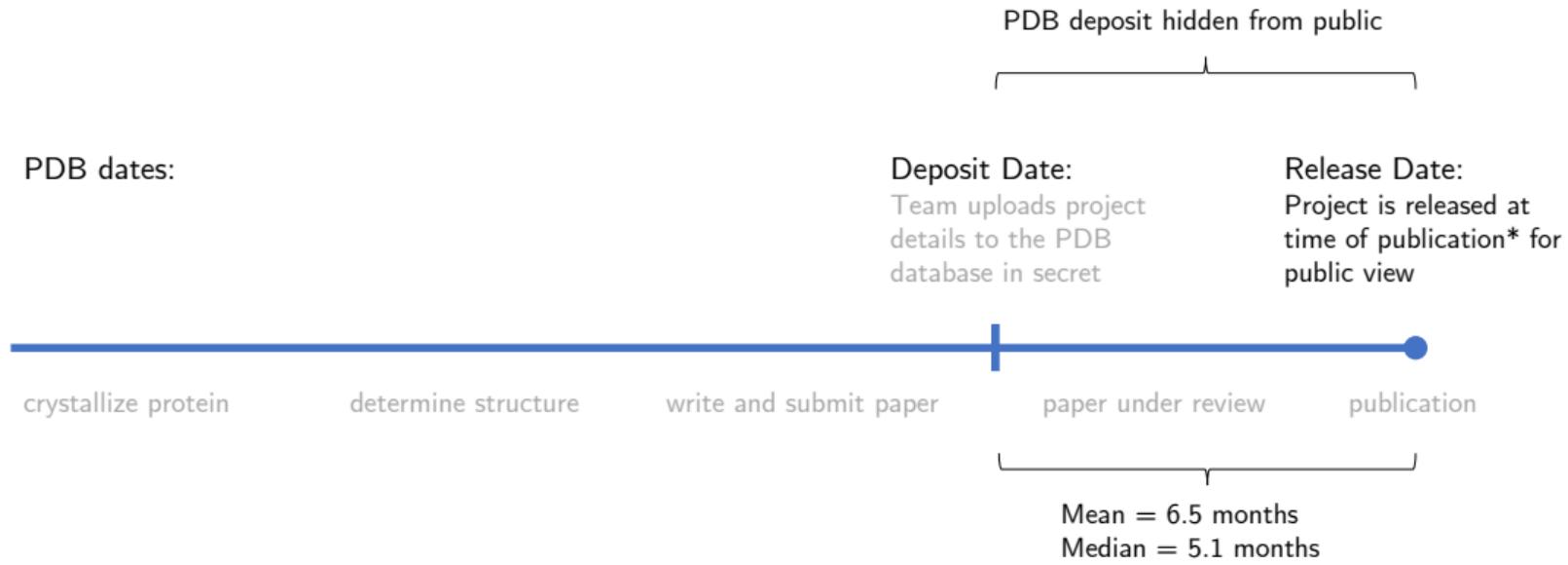
Project timeline

PDB dates:

Deposit Date:
Team uploads project
details to the PDB
database in secret



Project timeline

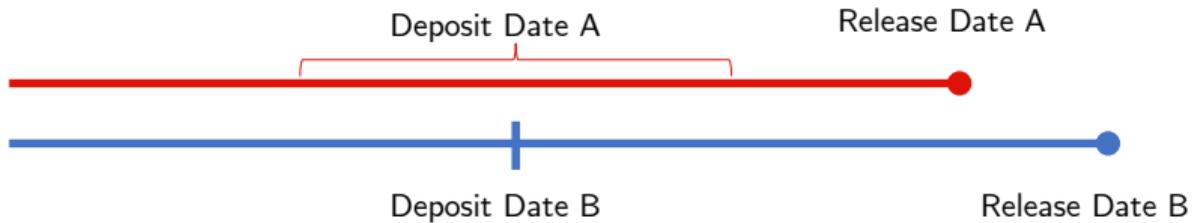


*If project goes unpublished, data is released publicly after one year

Scoop definition

- Rules:**
1. Take two projects that have identical sequence, different authors.
 2. Assert that both projects are deposited before the first project is released.
 3. Call the first to release the winner, call the second project “scooped.”

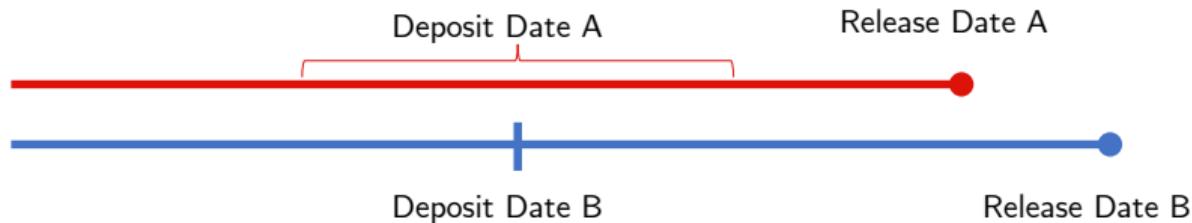
Scenario 1: Project A scoops Project B



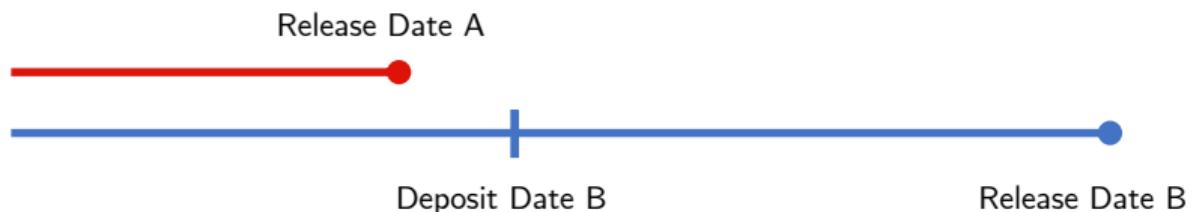
Scoop definition

- Rules:**
1. Take two projects that have identical sequence, different authors.
 2. Assert that both projects are deposited before the first project is released.
 3. Call the first to release the winner, call the second project “scooped.”

Scenario 1: Project A scoops Project B

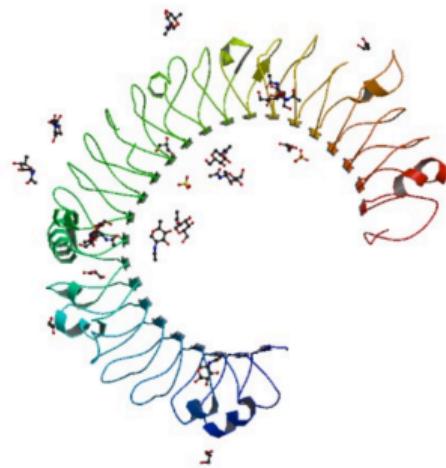


Scenario 2: Project A and Project B are excluded from racing sample



Example race: Toll-like receptor 3

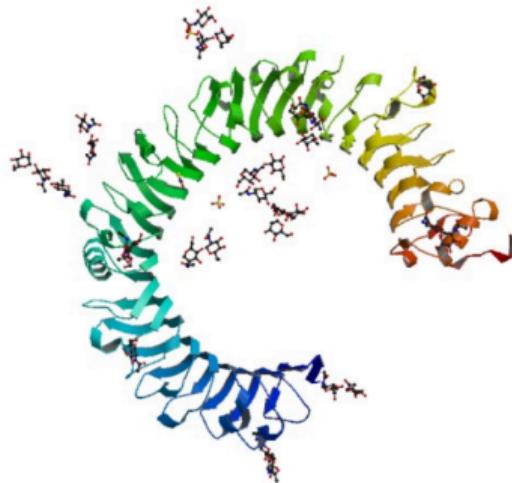
Winning Deposit: 1ZIW



Affiliation: Scripps Research Institute
Deposit Date: April 27, 2005
Release Date: June 28, 2005

Journal: *Science*
Journal Impact Factor: 30.9
5-year Citations: 196

Scooped Deposit: 2A0Z



Affiliation: National Institutes of Health
Deposit Date: June 27, 2005
Release Date: August 2, 2005

Journal: *PNAS*
Journal Impact Factor: 10.2
5-year Citations: 129

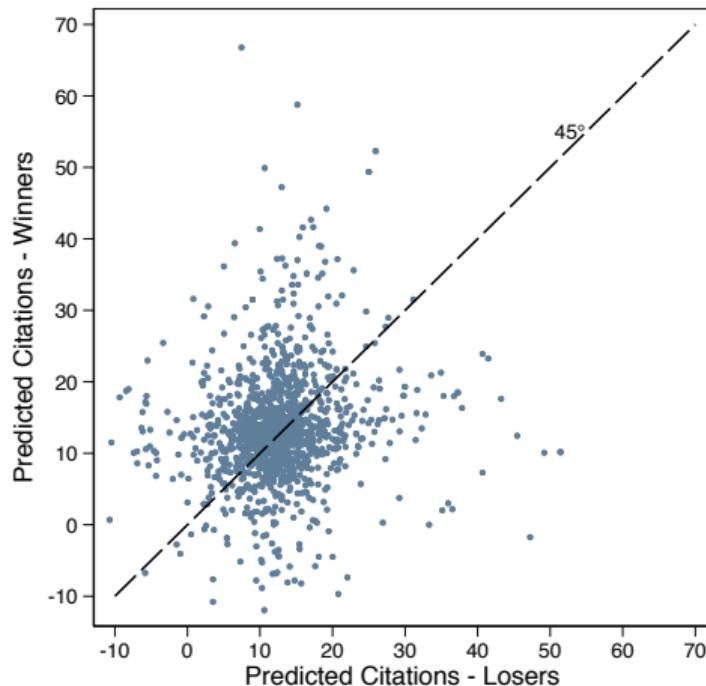
Predicted citation balance

Race winners are not randomly assigned, but seem highly unpredictable.

Lasso model of predicted citations:

- ▶ Team size and age
- ▶ Past deposits and publications
- ▶ University rank and location

Difference in predicted citations:
0.212 (p-value = 0.587)



Estimating the scoop penalty

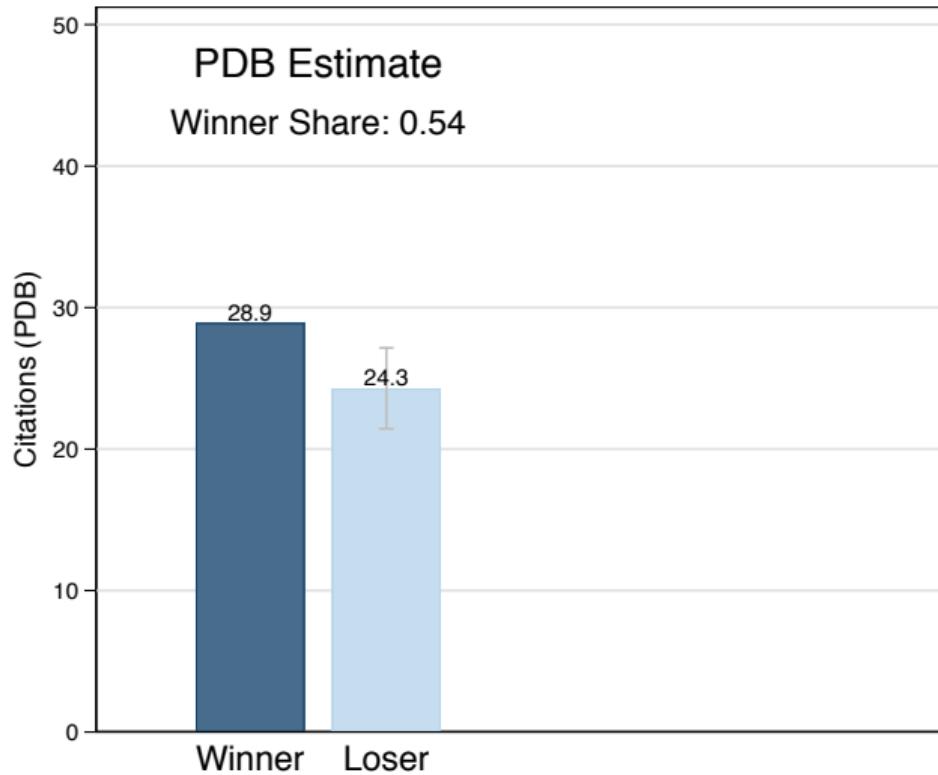
- ▶ Basic specification: For deposit i of protein (race) p :

$$Y_{ip} = \alpha + \beta Scooped_i + X'_i \delta + \gamma_p + \epsilon_{ip}$$

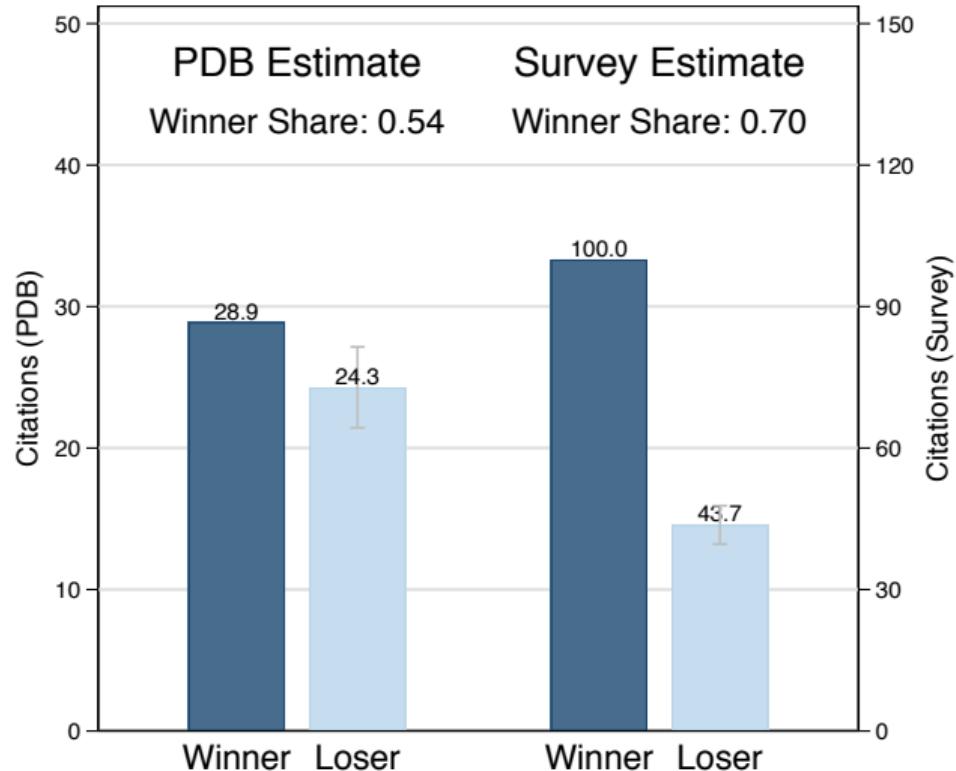
where

- ▶ $Scooped_i$ is a dummy for losing priority race.
- ▶ γ_p is the coefficient on a protein (i.e. race) fixed effect.
- ▶ X_i is a vector of individual and lab controls selected by PDS-Lasso method (Belloni et al. 2014).

Citation penalty



Citation penalty



Scoop penalty: alternative outcomes

Dependent variable	Published (1)	Std. journal impact factor (2)	Top-ten journal (3)	asinh(Five-year citations) (4)	Top-10% five year citations (5)
Scooped	-0.025*** (0.010)	-0.178*** (0.032)	-0.060*** (0.014)	-0.197*** (0.045)	-0.035*** (0.010)
Winner Y mean	0.880	-0.031	0.318	28.918	0.150
Observations	3,319	3,319	3,319	2,546	2,546

Note: All regressions include controls selected by PDS-Lasso as well as year fixed effects. Unpublished papers have impact factor imputed to minimum factor journal. Citation regressions restricted to papers published before 2014. Column 4 dependent variable is asinh(five-year citations) but mean citations is reported in levels.

The long-run consequences of being scooped

- ▶ Long run outcomes (excluding winning/scooped paper):
 - ▶ Active in PDB five years later
 - ▶ Total publications - five years
 - ▶ Total citations - five years
- ▶ Estimate for scientist s , deposit i , for protein (race) p :

$$Y_{isp} = \alpha + \beta Scooped_{is} + X'_{is} \delta + \gamma_p + \epsilon_{isp}$$

- ▶ Estimate separately for novices (<1 year of PDB experience) and veterans.

Long-run results

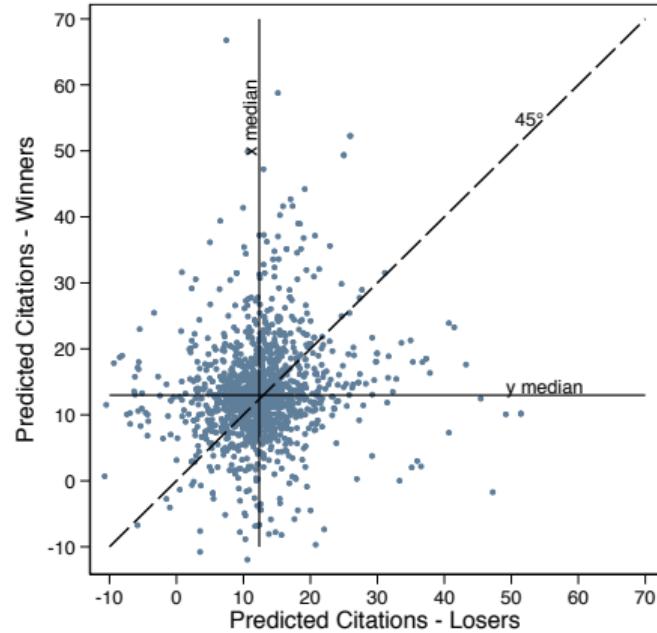
Dependent variable	Active in PDB 5 years later (1)	Total count five years after race (excluding original paper)			
		Publications (2)	Top-10 publications (3)	asinh of three-year citations (4)	Top-10% of three year citations (5)
<i>Panel A. All scientists</i>					
Scooped	-0.023** (0.010)	0.325 (0.264)	0.047 (0.105)	-0.199*** (0.053)	-0.159*** (0.052)
Winner Mean Y	0.797	13.541	4.441	187.129	1.529
Observations	6,642	12,488	12,488	9,297	9,297
<i>Panel B. Novices</i>					
Scooped	-0.013 (0.022)	-0.054 (0.183)	-0.076 (0.065)	-0.282** (0.112)	-0.121*** (0.039)
Winner Mean Y	0.587	2.667	0.929	50.692	0.440
Observations	2,273	3,554	3,554	2,868	2,868
<i>Panel C. Veterans</i>					
Scooped	-0.024** (0.009)	0.504 (0.373)	0.118 (0.145)	-0.167*** (0.051)	-0.165** (0.071)
Winner Mean Y	0.930	19.042	6.221	263.894	2.143
Observations	4,027	8,251	8,251	5,913	5,913

Priority and inequality

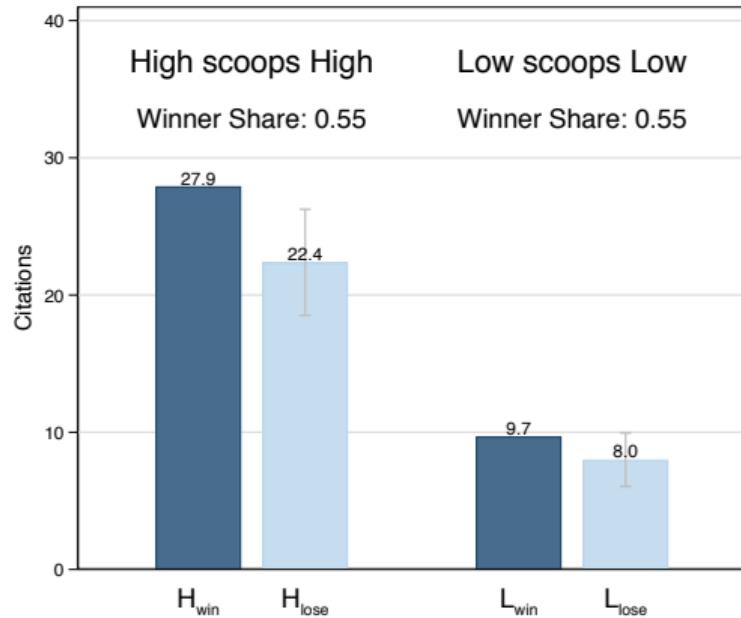
- ▶ Merton proposes two key drivers of academic attention:
 - ▶ Priority
 - ▶ Matthew Effect
- ▶ We test which of these effects dominates by comparing citations in races between high- and low-reputation teams.
- ▶ See the statistical discrimination model in the paper.

Defining reputation

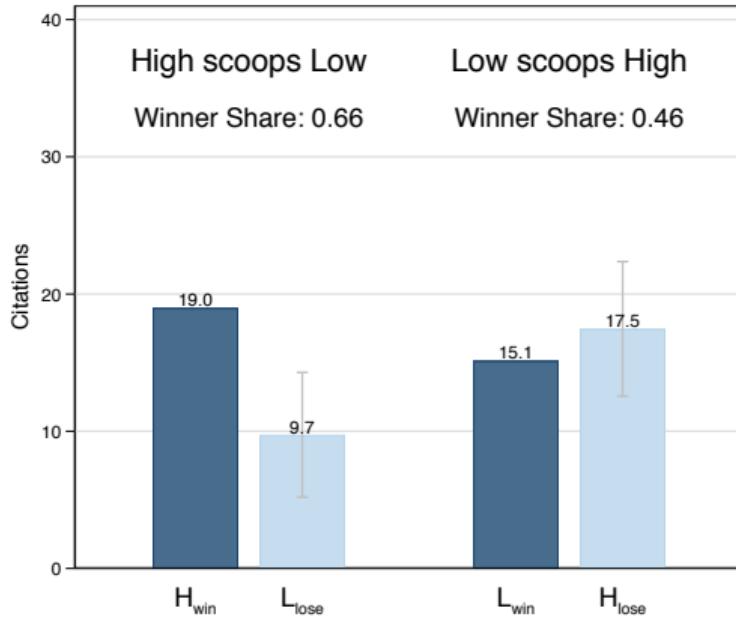
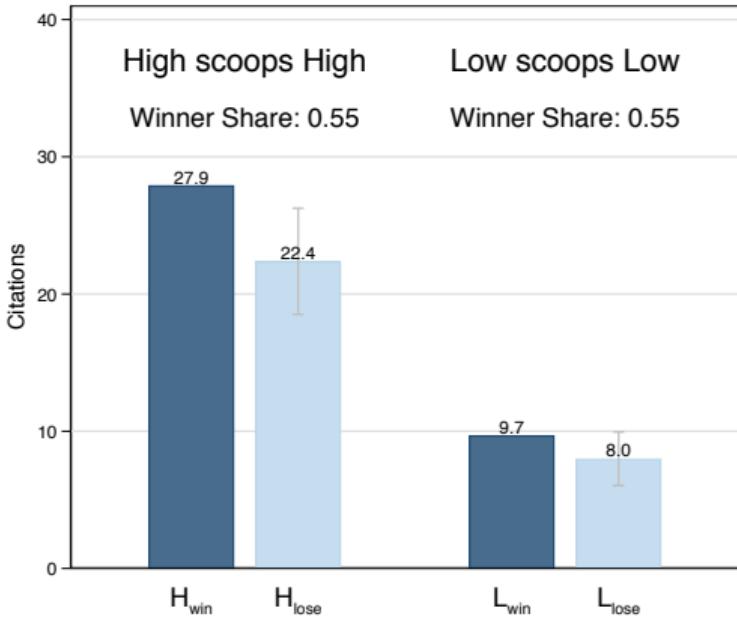
- ▶ Define pre-existing reputation using LASSO-generated predicted citations.
- ▶ Define H teams as those with above median predicted citations and L teams as those with below median.



Evenly-matched and mismatched races



Evenly-matched and mismatched races



Conclusion

Getting scooped lowers citations, but rewards are more evenly distributed than previously thought.

Normative implications: Is the premium for priority too large or too small?

- ▶ Priority may incentivize effort and timely disclosure.
- ▶ Racing may incentivize speed at the expense of quality and transparency.

Science as a non-market incentive

Scientific norms and their effects

Mertonian origins

Hill and Stein (2020)

Hill and Stein (2022)

Competition and quality in science

- ▶ Scientists compete to publish their findings first and establish priority. This competition can be good for science and society:
 - ▶ It can increase the pace of innovation
 - ▶ It induces scientists to disclose their work in order to get credit

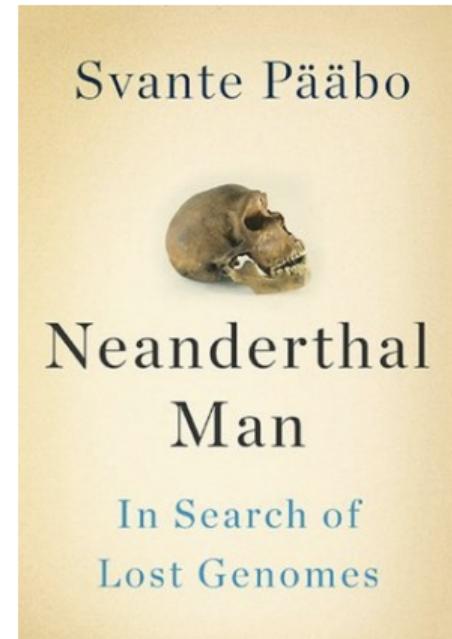
Competition and quality in science

- ▶ Scientists compete to publish their findings first and establish priority. This competition can be good for science and society:
 - ▶ It can increase the pace of innovation
 - ▶ It induces scientists to disclose their work in order to get credit
- ▶ On the other hand, competition may have a dark side:
 - ▶ Scientists may cut corners and reduce quality in their pursuit to publish first
 - ▶ **Focus of this project**

Example: Sequencing the Neanderthal Genome

"Hendrik's paper also illustrated a dilemma in science: doing all the analyses and experiments necessary to tell the complete story leaves you vulnerable to being beaten to the press...Even when you publish a better paper, you are seen as mopping up the details after someone who made the real breakthrough"

– Svante Pääbo, *Neanderthal Man: In Search of Lost Genomes*



Summary of the model

- ▶ Projects vary in their ex-ante potential

Summary of the model

- ▶ Projects vary in their ex-ante potential
- ▶ Scientists decide how long to work on a project (maturation), trading off improving the quality of their work against the threat of being scooped (Bobtcheff et. al 2017)

Summary of the model

- ▶ Projects vary in their ex-ante potential
- ▶ Scientists decide how long to work on a project (maturation), trading off improving the quality of their work against the threat of being scooped (Bobtcheff et. al 2017)
- ▶ **Key ingredient:** entry into projects is endogenous → more likely to be competition in high potential projects

Summary of the model

- ▶ Projects vary in their ex-ante potential
- ▶ Scientists decide how long to work on a project (maturation), trading off improving the quality of their work against the threat of being scooped (Bobtcheff et. al 2017)
- ▶ **Key ingredient:** entry into projects is endogenous → more likely to be competition in high potential projects
- ▶ **Key result:** high potential projects will be executed with lower quality

Key propositions

- ▶ **Proposition 1:**

High potential projects are more attractive to enter → are more competitive

Key propositions

- ▶ **Proposition 1:**
High potential projects are more attractive to enter → are more competitive
- ▶ **Proposition 2:**
Competitive projects completed faster → are lower quality

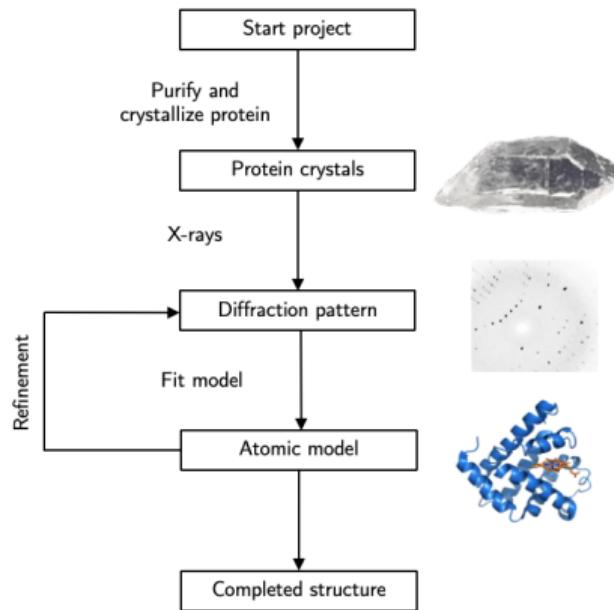
Key propositions

- ▶ **Proposition 1:**
High potential projects are more attractive to enter → are more competitive
- ▶ **Proposition 2:**
Competitive projects completed faster → are lower quality
- ▶ **Proposition 3 (key model prediction):**
High potential projects completed faster → are lower quality

How do scientists solve protein structures?

About 90% of proteins are solved using X-ray crystallography. This involves three steps:

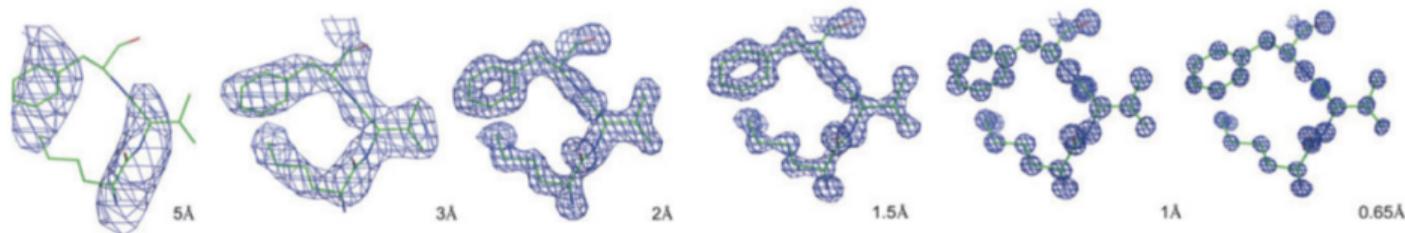
1. First, proteins are purified and crystallized
2. Next, the crystals are placed in an x-ray beam, which produces a diffraction pattern
3. Finally, the diffraction data is used to infer the structure. Biologists will "refine" their structure by comparing their model to the diffraction data, trying to minimize any discrepancies. Process is more "art than science" and luck plays a role



Mapping to the model: quality

A unique feature of structural biology is the objective, ex-ante measures of project quality:

1. Refinement resolution: similar to resolution of a photograph

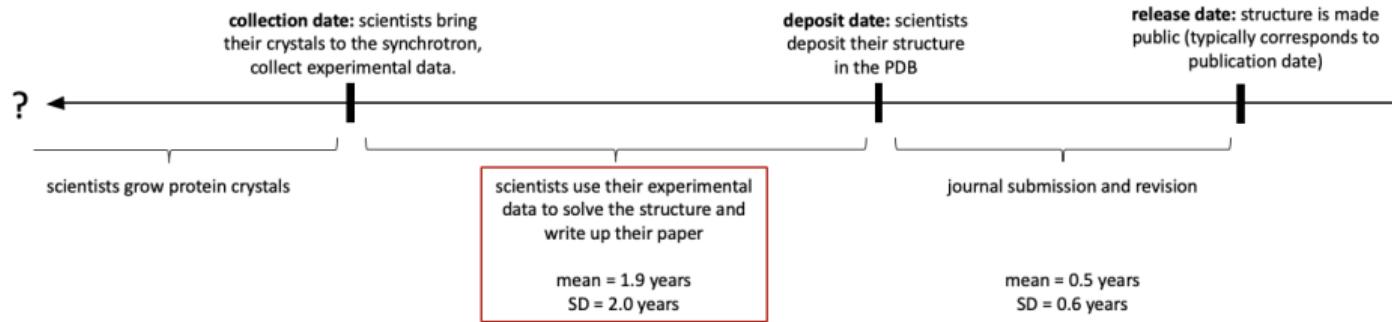


2. R-free: model fit, estimated on a holdout sample of the experimental data
3. Outlier share: errors in the model based on chemical properties

Combine these outcomes into a standardized quality index (higher is better)

Mapping to the model: maturation

- We can actually observe time spent on project (maturation period):



Mapping to the model: competition

- ▶ The PDB uses amino acid sequence similarity to flag proteins that are identical
- ▶ Number of times the same protein is deposited (within two years) can proxy for competition
- ▶ Note that we are measuring ex-post realized competition, a noisy proxy for ex-ante competition

released within two years

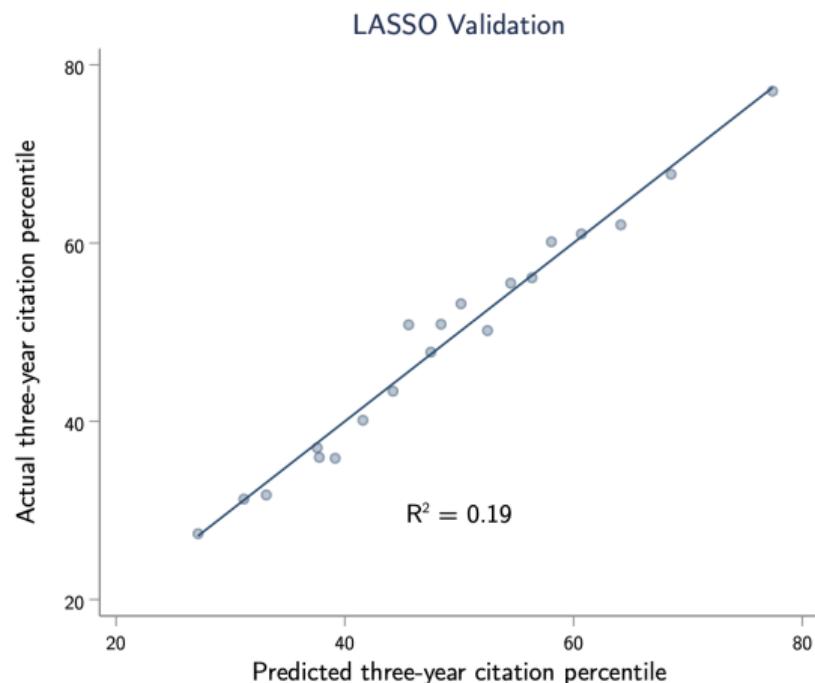
1F3H: Entity 1: Chains A, B
X-RAY CRYSTAL STRUCTURE OF THE HUMAN ANTI-APOPTOTIC PROTEIN SURVIVIN
Verdeca, M.A., Huang, H., Dutt, E., Hunter, T., Noel, J.P.
(2000) Nat Struct Biol 7: 602-608
Released: 2000-12-06
Method: X-RAY DIFFRACTION 2.58 Å
Organism: Homo sapiens
Macromolecule: SURVIVIN
Sequence Match: Sequence identity: 100%, E-Value: 7.7e-94, Region: 1-142

1E31: Entity 1: Chains A, B
SURVIVIN DIMER H. SAPIENS
Chantat, L., Skoufos, D.A., Margolis, R.L., Didieberg, O.
(2000) Mol Cell 6: 183
Released: 2001-01-03
Method: X-RAY DIFFRACTION 2.71 Å
Organism: Homo sapiens
Macromolecule: APOTOPSIS INHIBITOR SURVIVIN
Sequence Match: Sequence identity: 99%, E-Value: 1.985e-93, Region: 1-142

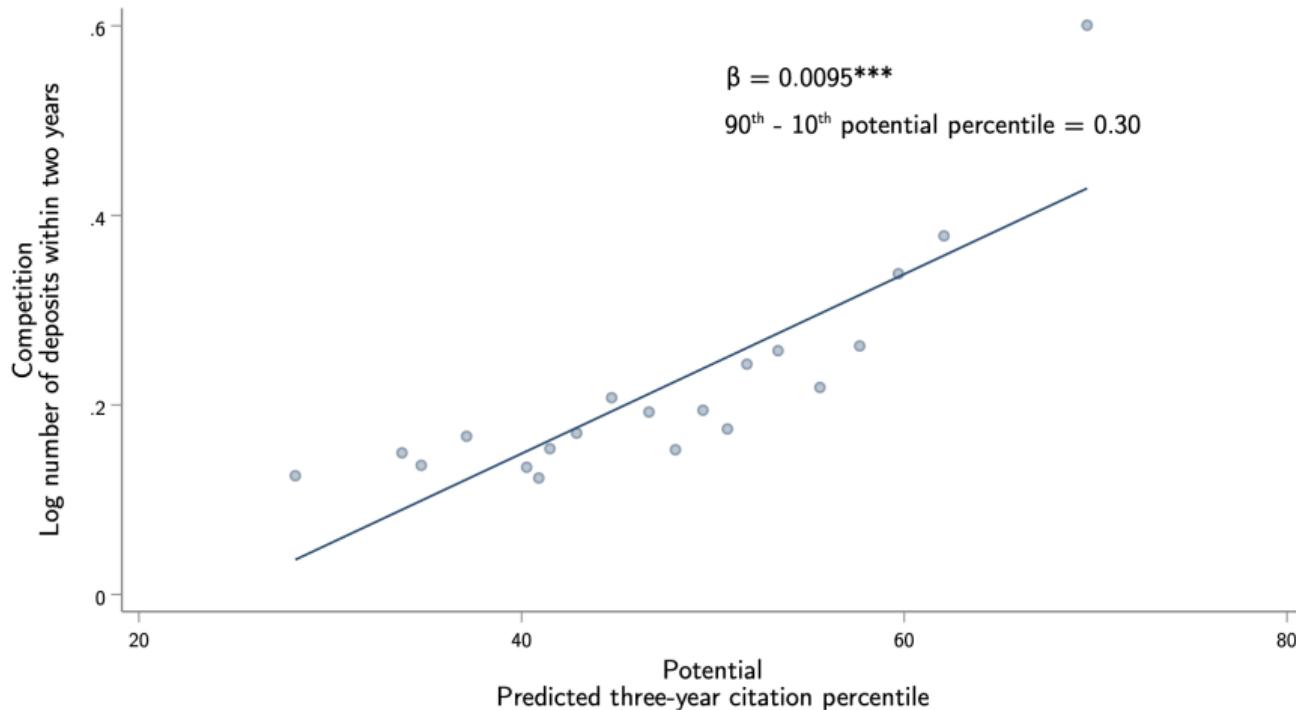
1XOX: Entity 1: Chains A, B
SOLUTION STRUCTURE OF HUMAN SURVIVIN
Sun, C., Nettethheim, D., Liu, Z., Olejniczak, E.T.
(2005) Biochemistry 44: 11-17
Released: 2005-01-18
Method: SOLUTION NMR
Organism: Homo sapiens
Macromolecule: Apoptosis inhibitor survivin
Sequence Match: Sequence identity: 99%, E-Value: 3.299e-78, Region: 1-117

Mapping to the model: measuring and predicting potential in the PDB

- ▶ One way to measure potential:
use ex-post citations (over some time window)
 - ▶ Problems: ex-post citations different than ex-ante potential, conflates potential and quality
- ▶ Alternatively: predict citations using only ex-ante characteristics of the structure
 - ▶ To avoid over-fitting, we use LASSO to select the model

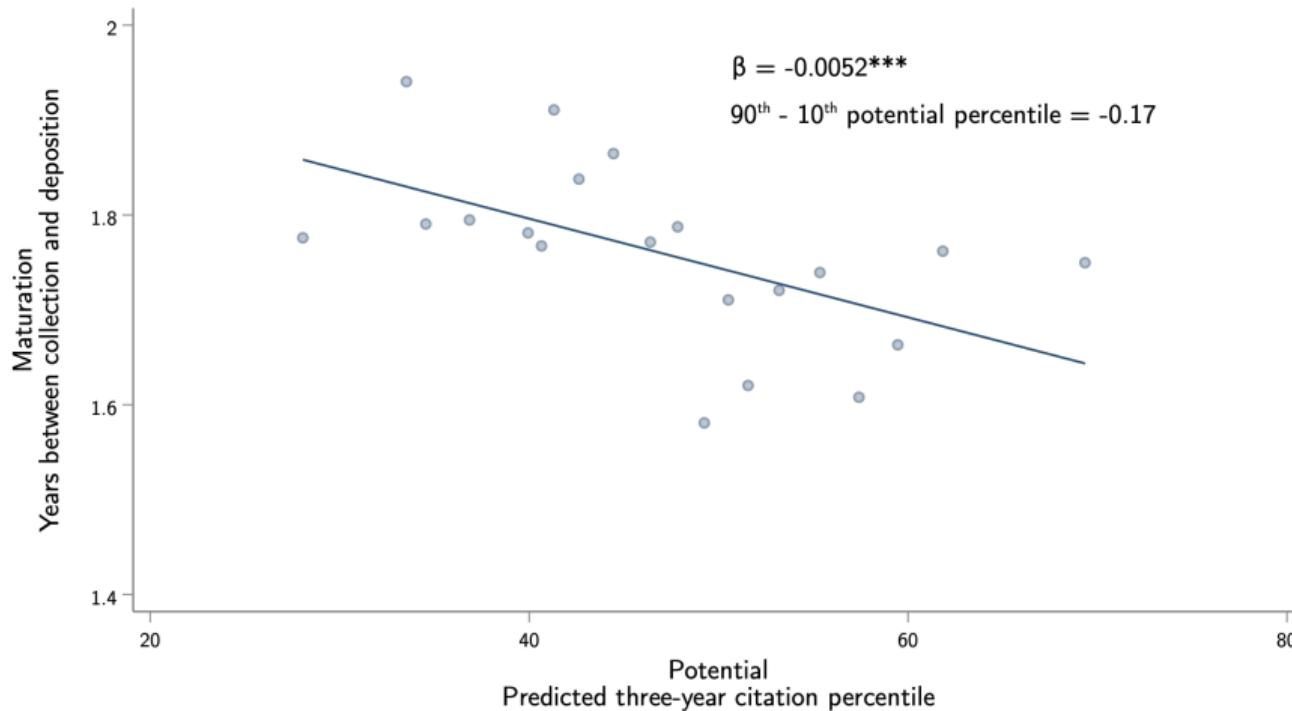


Proposition 1: high-potential projects are more competitive



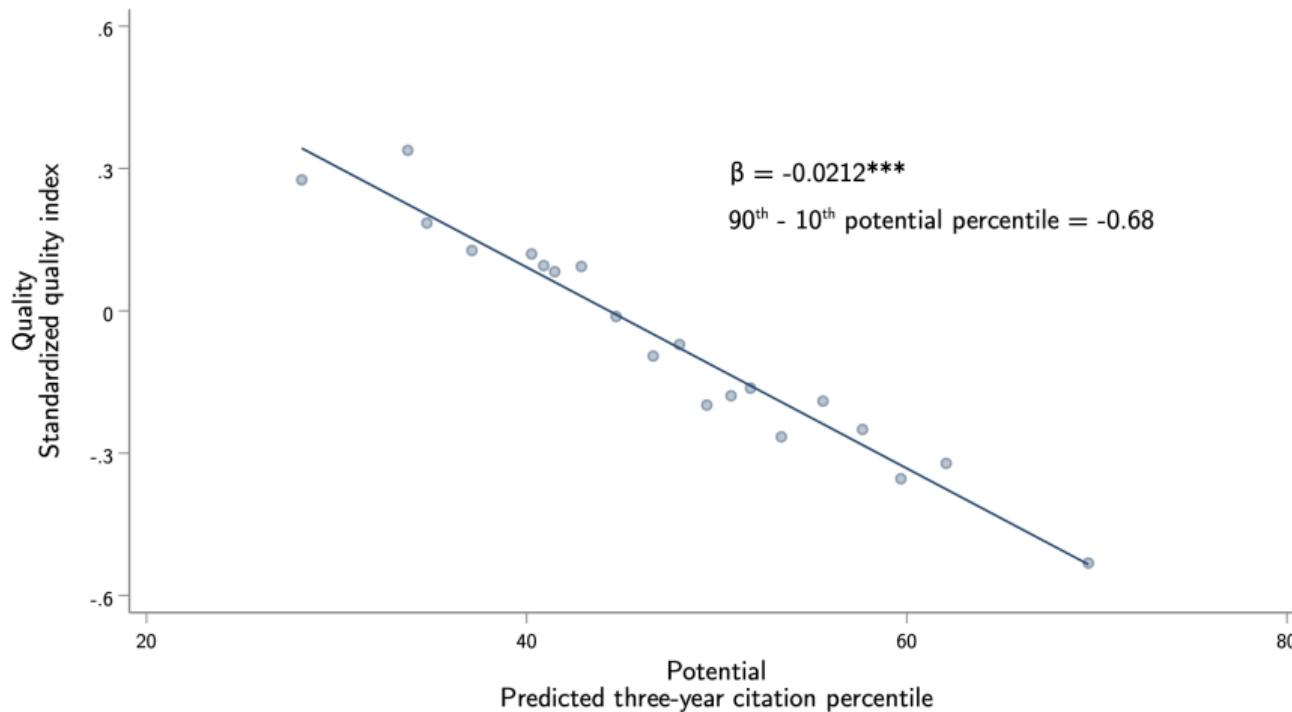
$$\text{LogDepositsInCluster}_{it} = \alpha + \beta \text{PredictedCites}_{it} + \tau_t + \epsilon_{it}$$

Proposition 3: high-potential projects are completed faster...



$$\text{Maturation}_{it} = \alpha + \beta \text{PredictedCites}_{it} + \tau_t + \epsilon_{it}$$

...so high-potential projects are lower quality



$$Quality_{it} = \alpha + \beta PredictedCites_{it} + \tau_t + \epsilon_{it}$$

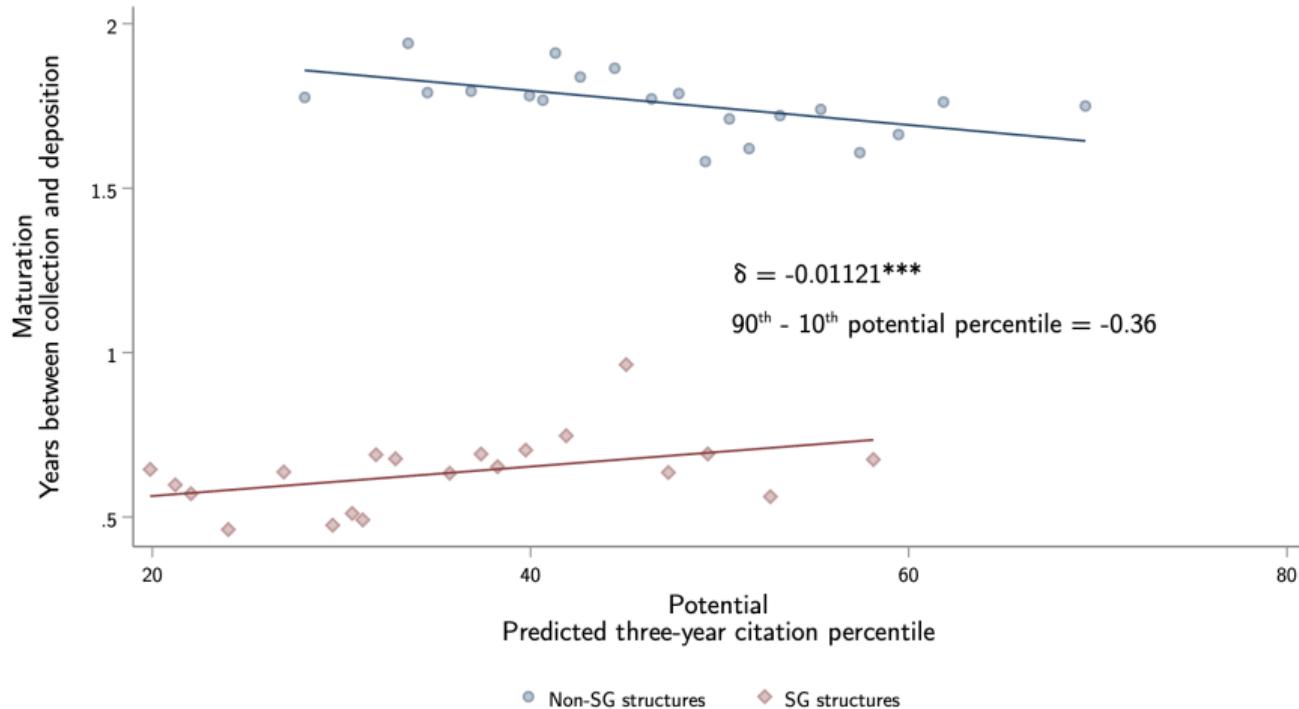
What about project complexity?

- ▶ **In general:** omitted variables bias might be a concern
- ▶ **In particular:** if high potential projects are also more complicated, this could drive our results. Lower quality is caused by the difficulty / complexity of the project, not rushing
- ▶ In the paper, we perform several analyses to rule out this story; today focus on just one

Structural genomics consortia

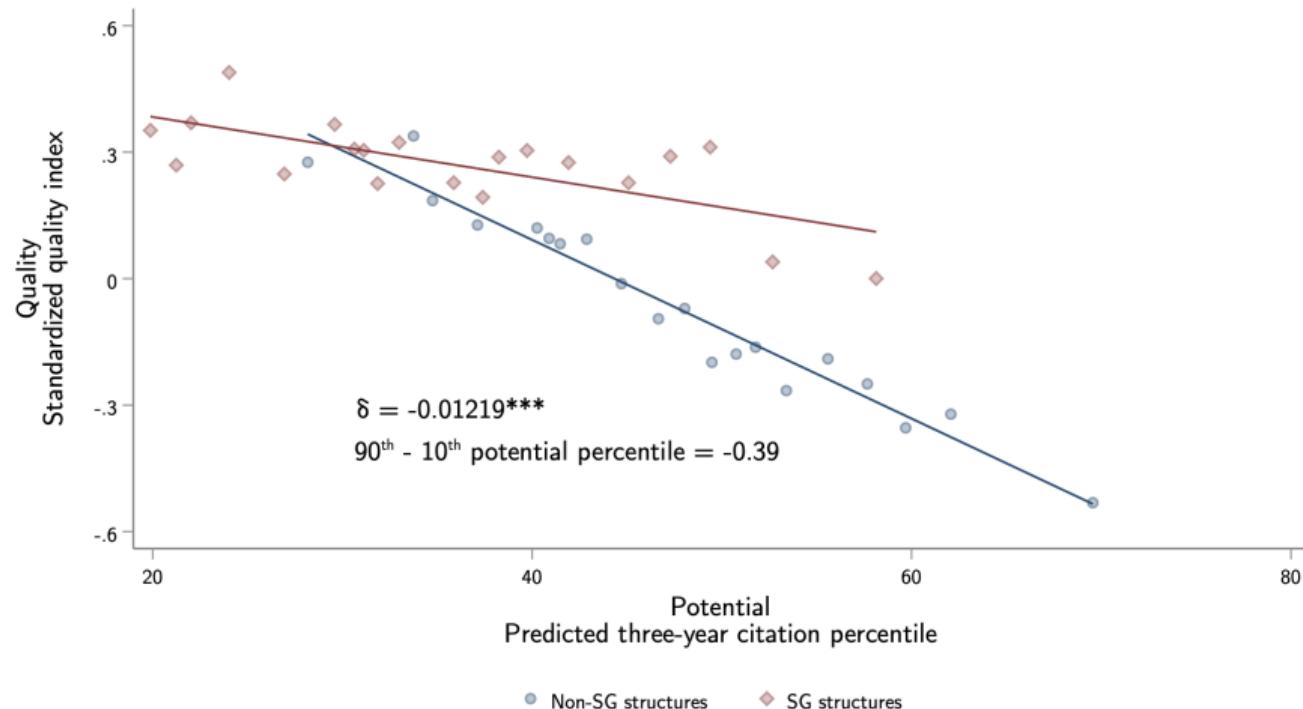
- ▶ **Key idea:** scientists affiliated with different types of institutions face different incentives
- ▶ Structural genomics consortia are publicly funded groups focused on achieving comprehensive coverage of the protein folding space
- ▶ Less focused on publishing and priority → competition is less important
- ▶ About 20% of structures in our sample were deposited by a structural genomics group

SG versus non-SG structures: maturation



$$Maturation_{it} = \alpha + \beta PredictedCites_{it} + \gamma NonSG_{it} + \delta (PredictedCites_{it} * NonSG_{it}) + \tau_t + \epsilon_{it}$$

SG versus non-SG structures: quality



$$Quality_{it} = \alpha + \beta PredictedCites_{it} + \gamma NonSG_{it} + \delta (PredictedCites_{it} * NonSG_{it}) + \tau_t + \epsilon_{it}$$

Conclusions and future work

- ▶ **Positive conclusion:** competition in science leads researchers to work faster and produce lower quality work
- ▶ **Normative analysis:** back-of-the envelope analysis suggests we have spent \$2-5 billion on “cleaning up” these low-quality deposits. A social planner would prefer that we do them well the first time!
- ▶ However, taking a stand on optimal competition is hard. Competition likely affects science in ways we have not considered here:
 - ▶ May reduce collaboration and free sharing of ideas
 - ▶ Impacts who enters certain fields and who is deterred
- ▶ Brings up questions of alternative models of science:
 - ▶ More collaborative models: Protein Structure Initiative, Human Genome Project