Studying NIH Grant Peer Review

Molly Carnes, MD, MS Elizabeth Pier, PhD September 1, 2017 NIH Director's Transformative R01 Exploring the Science of Scientific Review (R01 GM111002) Multiple PI: M. Carnes, C. Ford, P. Devine

3 Studies:

- 1. Text analysis of R01 critiques.
- 2. Examination of constructed study sections.
- 3. Experimental manipulation of R01 applicant race and gender.

Race and Gender Differences in R01 Award Rates

Black PIs have Lower R01 award probabilities than White PIs



Probability of NIH R01 award by race and ethnicity, FY 2000 to FY 2006 (N = 83,188); \ddagger , p < .001

Based on data from NIH IMPAC II, DRF, and AAMC Faculty Roster. (Ginther, *Science*, 2011)

Female PIs have lower R01 renewal award rates than male PIs



Success Rates by Gender and Type of Application, FY 1998 to FY 2016

Subjectivity in evaluative criteria allows implicit assumptions to affect interpretation of objective data

- The temperature is 41 degrees F = objective
- Is that a cold or hot day? = subjective
 - Wisconsin = warm
 - Florida = cold
- R01 application submitted by a certain PI from a certain institution = objective
 - Innovative scientific leader ך
 - Pioneering research
- subjective
- Impact scores 2, 3, 4?

Stereotypes are known even if we don't believe them

Men ¹	Women ¹	White ²	Asian ²	Black ²	Latino ²
 Strong Decisive Stubborn Competitive Ambitious Risk-taking Assertive Tough Authoritative Independent 	 Caring Nurturing Bad drivers Family- oriented Communal Supportive Sympathetic Nice Helpful 	 High status Rich Intelligent Arrogant Privileged Blonde Racist All-American Ignorant 	 Intelligent Bad drivers Good at math Nerdy Shy Skinny Educated Quiet 	 Ghetto or unrefined Criminal Athletic Loud Gangsters Poor Have an attitude Unintelligent Uneducated 	 Poor Illegal immigrant Uneducated Family- oriented Lazy Day laborer Unintelligent Loud Gangsters

Quantitative text analysis of R01 critiques

- 443 grant reviews from R01s awarded after unfunded in 2008 (N=65)
- Women's: more standout adjectives (e.g., excellent, outstanding) (p≤0.01)
- Men's: more negative descriptors (e.g., unfocused, illogical) (p≤0.01)



Do implicit assumptions lead to different referent standards in evaluating the applications of male and female PIs?

Kaatz et al., 2015

R01 Grant Critiques UW-Madison (N=739), 2010-2014: Greater praise for women ≠ better score

- Female PIs' R01 renewals assigned worse (higher) priority scores.
- More critiques of female PIs' R01 renewals included words about ability and standout adjectives (e.g., "outstanding," "excellent").



*Difference between groups is significant (P<.05); PIs in sample had similar levels of productivity and background qualifications; models controlled for funding outcome and experience level.

Summary of text analysis to date

- Our work suggests that characteristics of the investigator could introduce bias into evaluation of an application in ways that could contribute to the gender difference in Type 2 renewals
- Have collected a national sample of over 6000 PIs (~19,000 critiques) that will enable to include race in our analyses
- Examining algorithmic text mining approaches

Examining Constructed Study Section Meetings

• Subjectivity in evaluative criteria allows implicit assumptions to affect interpretation of objective data

Overall Impact	Score	Descriptor
	1	Exceptional
High	2	Outstanding
	3	Excellent
	4	Very Good
Medium	5	Good
	6	Satisfactory
	7	Fair
Low	8	Marginal
	9	Poor













Application	<u>CSS1</u>	<u>CSS2</u>	<u>CSS3</u>	<u>CSS4</u>
Wu		X	X	•
Lopez	•	•	•	
Holzmann	X		X	•
Zhang	X	X	•	X
Adamsson	•	X	X	X
McMillan		X	•	
Phillips	•	•	X	
Amsel	•	•	•	
Abel	•	•	•	X
Ferrera	X	X	•	X
Stavros	\mathbf{X}	•	X	•
Washington	•	•	X	•
Williams	•	X	•	•
Rice	X	٠	•	X
Albert	•	X	X	•
Edwards	X	•	X	X
Foster	•	•	•	•
Henry	•	٠	•	•
Bretz	X	X	•	
Molloy	•	•	X	
Wei	X	X	X	X
McGuire		X	X	
Kim	X	X	X	
Bernard	X	X		
Lukska	X	X	X	X



Scoring Variability

How variable are reviewers' scores?

Krippendorff's α

(Hayes & Krippendorff, 2007)

 $\alpha \ge 0.8$ "Reliable" $\alpha \ge 0.67$ "Tentative but not definitive"

Low agreement among individual reviewers $\alpha = .0840$

Better agreement <u>within</u> panels after collaborative discussion $\alpha = .6652$

Worse agreement <u>between</u> panels after collaborative discussion $\alpha = -.0517$



Scoring Variability

How variable are reviewers' scores?

Low agreement among individual reviewers $\alpha = .0840$



Preliminary Score Range Final Score Range

Better agreement <u>within</u> panels after collaborative discussion $\alpha = .6652$

Worse agreement <u>between</u> panels after collaborative discussion $\alpha = -.0517$



Better agreement <u>within</u> panels after collaborative discussion $\alpha = .6652$

Worse agreement <u>between</u> panels after collaborative discussion $\alpha = -.0517$

Scoring Variability



Score Calibration Talk [SCT]

Overall Impact	Score	Descriptor
	1	Exceptional
High	2	Outstanding
	3	Excellent
	4	Very Good
Medium	5	Good
	1 Ex 2 Ou 3 Ex 4 Ve 5 Go 6 Sa 7 Fa 8 M 9 Po	Satisfactory
	7	Fair
Low	8	Marginal
	9	Poor



LZ-2: So that means uh probably they are already recognize this



JA-3: you highlight those differences you know or uh y'know

Score Calibration Talk [SCT]

Instances of Score Calibration Talk (SCT) in each CSS

	CSS1	CSS2	CSS3	CSS4	Total
Self-initiated SCT					
# Instances	15	18	11	12	56
Time (m:s)	3:33	4:36	2:09	2:37	12:55
Other-initiated SCT					
# Instances	7	3	4	1	15
Time (m:s)	6:07	4:28	5:27	1:46	17:48
Total SCT					
# Instances	22	21	15	15	71
Time (m:s)	9:40	9:04	7:36	4:23	30:43

Pier et al, 2017

Score Calibration Talk [SCT]

SCT & Reviewe		SCT & Re	
Self-initiated SCT	Correlation		Self-initiated
# Instances	<i>r</i> = .108		# Instance
Time (m:s)	<i>r</i> = .067		Time (m:
Other-initiated SCT			Other-initiat
# Instances	<i>r</i> = .978		# Instance
Time (m:s)	<i>r</i> = .961		Time (m:
<u>Total SCT</u>			<u>Total SCT</u>
# Instances	<i>r</i> = .717		# Instance
Time (m:s)	<i>r</i> = .809		Time (m:s
	SCT & Reviewe Self-initiated SCT # Instances Time (m:s) Other-initiated SCT # Instances Time (m:s) Total SCT # Instances Time (m:s)	SCT & Reviewer Score ChangeSelf-initiated SCTCorrelation# Instances $r = .108$ Time (m:s) $r = .067$ Other-initiated SCT $r = .978$ # Instances $r = .978$ Time (m:s) $r = .961$ Total SCT $r = .717$ # Instances $r = .717$ Time (m:s) $r = .809$	Set & Reviewer Score ChangeSelf-initiated SCTCorrelation# Instances $r = .108$ Time (m:s) $r = .067$ Other-initiated SCT $r = .978$ # Instances $r = .978$ Time (m:s) $r = .961$ Total SCT $r = .717$ # Instances $r = .809$

viewer Score Convergence d SCT Correlation *r* = .682 es r = .657s) ted SCT *r* = .858 es r = .784s) r = .980es *r* = .936 s)

Within-panel convergence & Between-panel divergence

r = -.606 (p = .005)

Pier et al, 2017

Conclusions



- Identifying the interpersonal and communicative processes that unfold during peer review meetings helps us better understand how subjectivity can bias people's decision making
- SCT may be a prime target for intervention to help continue to improve the reliability of the peer review process
- Training SROs and Chairs about the local norming that happens during SCT could help them ensure more consistent adherence to the scoring rubric

Thank you!



mlcarnes@wisc.edu

epier@wisc.edu

This work is supported by the National Institute of General Medical Sciences of the National Institutes of Health (Award # R01GM111002) and by the Arvil S. Barr Graduate Fellowship.