

# Fairness and Bias in **Peer Review** and other Sociotechnical Intelligent Systems

**Nihar B. Shah** and Zachary Lipton

School of Computer Science

**Carnegie Mellon University**

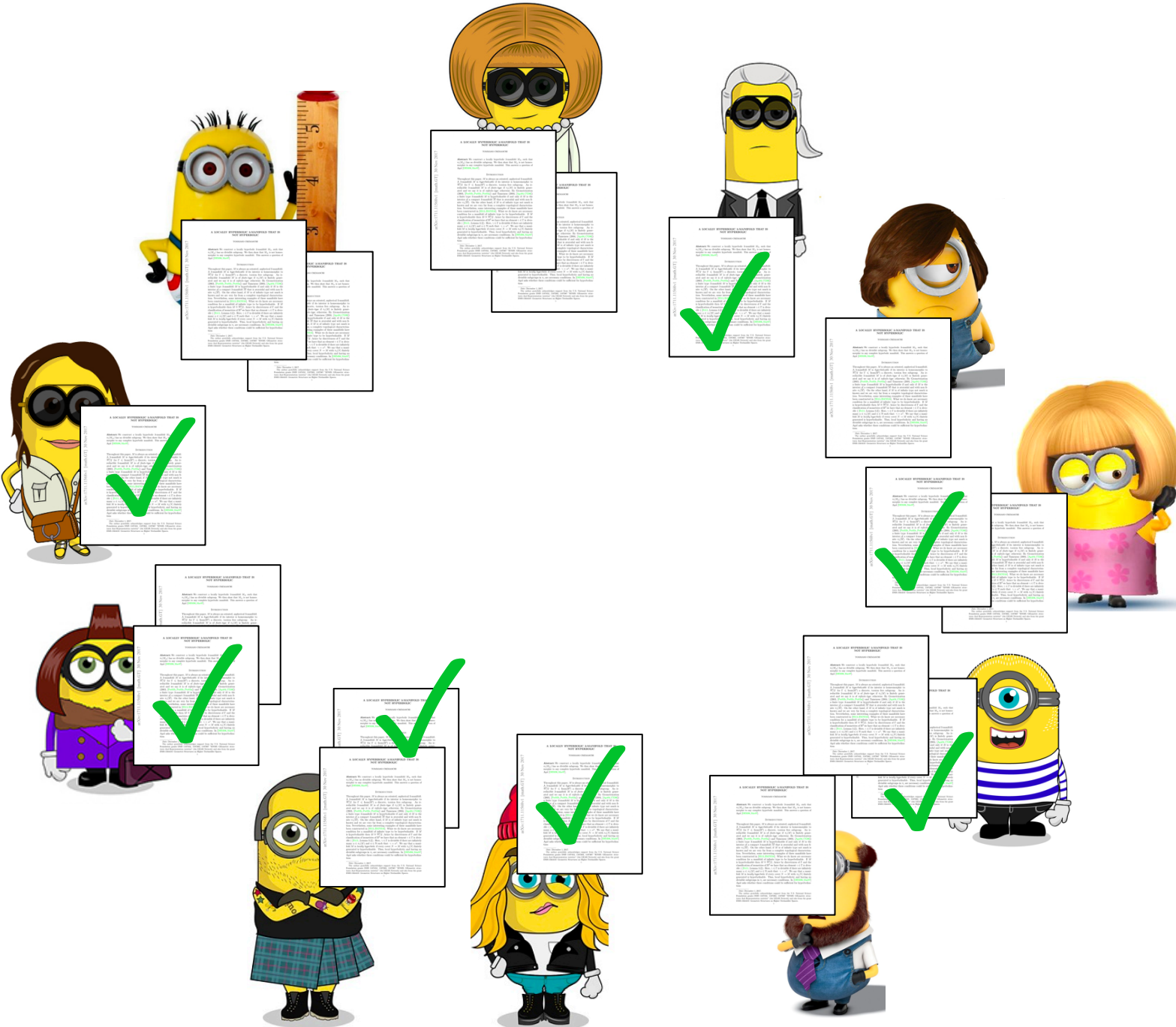
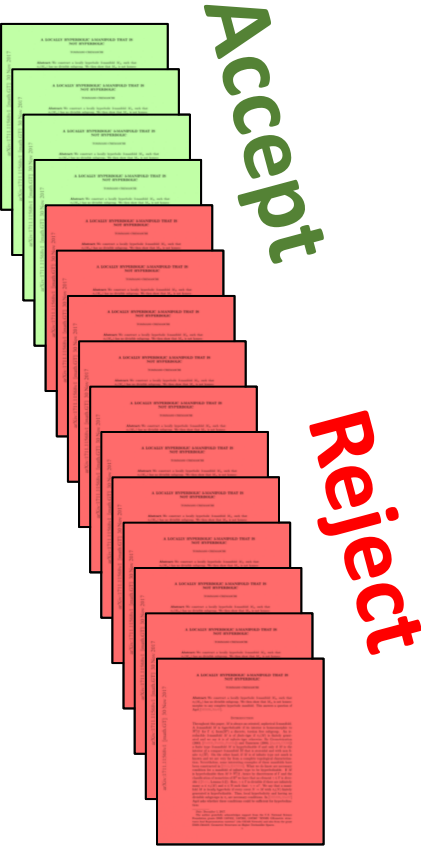


"Piled Higher and Deeper" by Jorge Cham

WWW.PHDCOMICS.COM

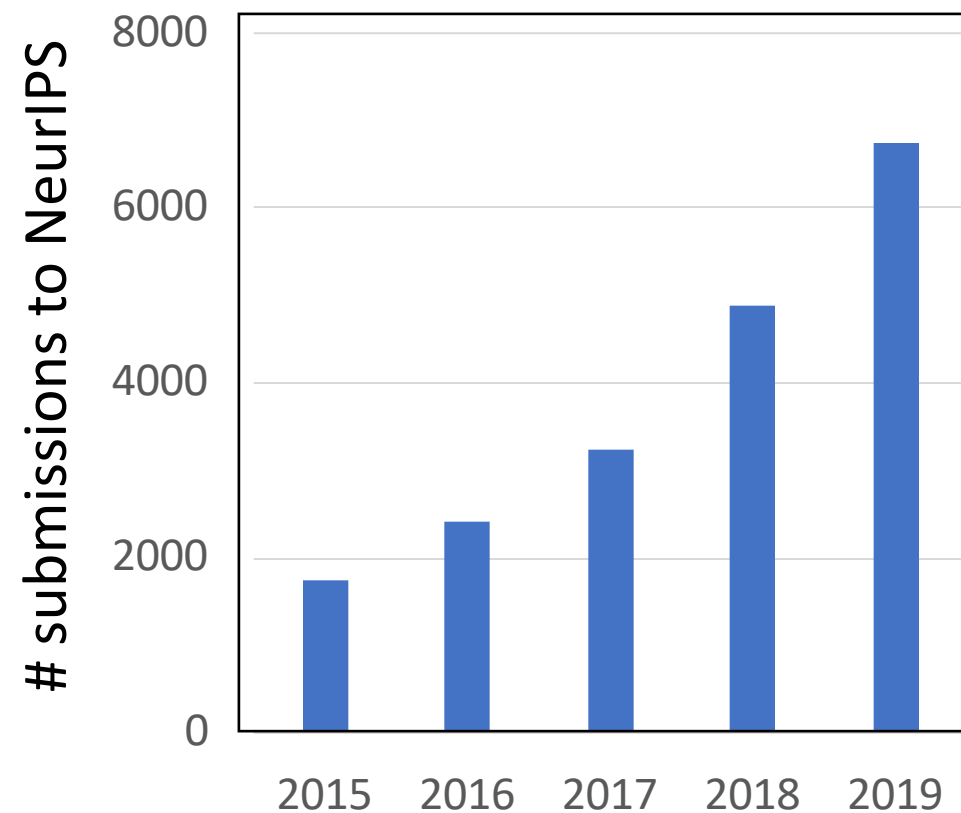
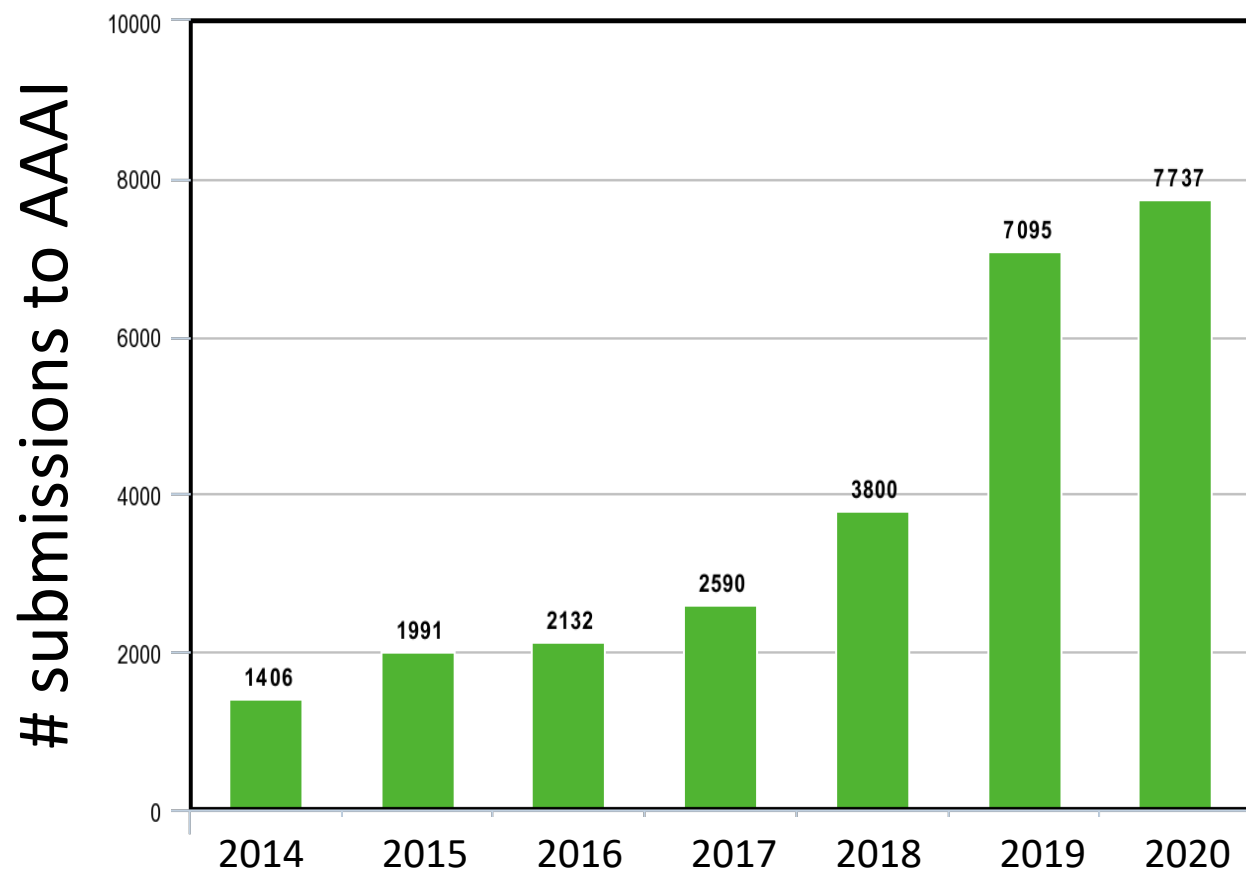


# Peer-review





# Tremendous growth



**Several thousand submissions, exponential growth**



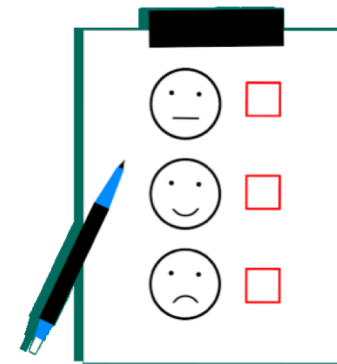
# Challenge across many research fields

- “Let's make peer review scientific” [Rennie, Nature 2016]

*“Peer review ... is a human system. Everybody involved brings **prejudices**, **misunderstandings** and gaps in knowledge, so no one should be surprised that peer review is often **biased** and **inefficient**. It is occasionally **corrupt**, sometimes a charade, an open temptation to plagiarists. Even with the best of intentions, how and whether peer review identifies high-quality science is unknown. It is, in short, **unscientific**.”*

- Overwhelming desire for improvement

[surveys by Smith 2006, Ware 2008, Mulligan et al. 2013]

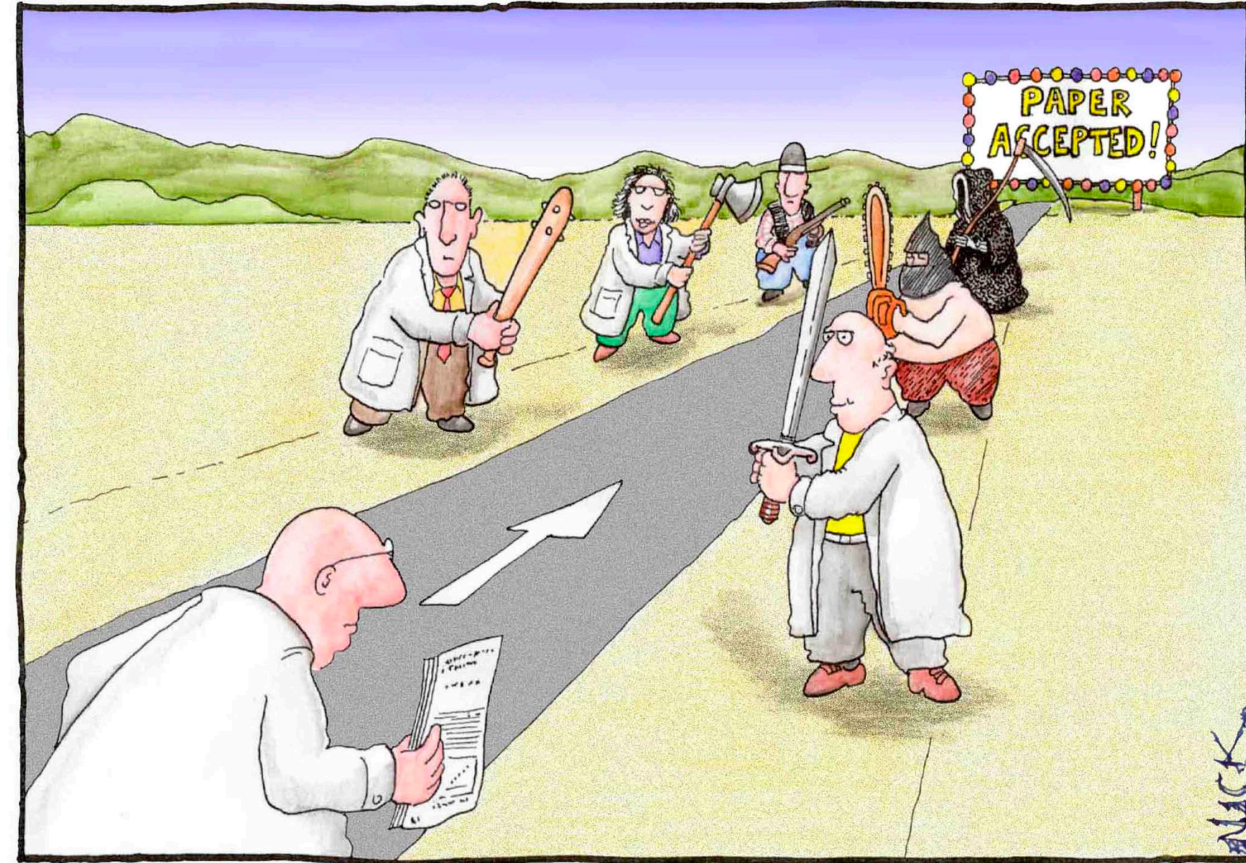




# Hurts scientific progress

“interdisciplinary research, frontier science, areas of controversy, and risky new departures are all more likely to **suffer from cognitive cronyism** than is mainstream research” [Travis et al. 1991]

“Reviewers love safe (boring) papers, ideally on a topic that has been discussed before (ad nauseam)...**The process discourages growth**” [Church 2005]

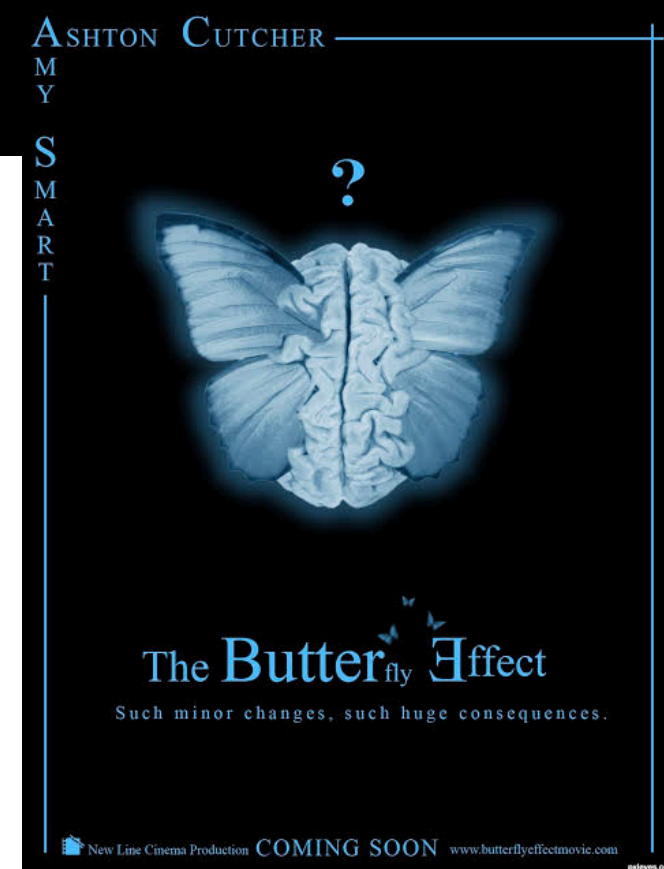




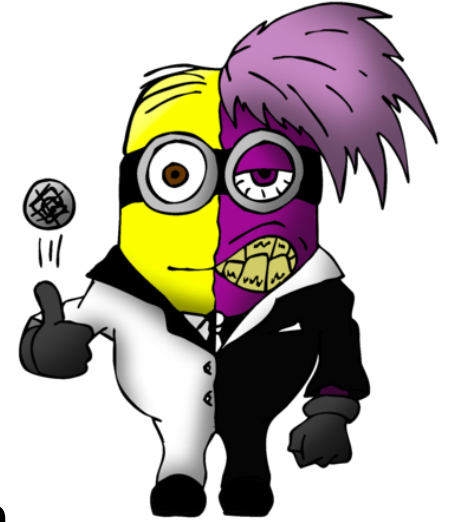
# Hurts careers

“an incompetent review may lead to the rejection of the submitted paper, or of the grant application, and the ultimate **failure of the career of the author.**” [Triggle et al. 2007]

“These long term effects arise due to the widespread prevalence of the Matthew effect (‘**rich get richer**’) in academia” [Merton 1968]







“In public, scientists and scientific institutions celebrate truth and innovation. In private, they perpetuate peer review biases that thwart these goals... **what can be done about it?**”

[Lee 2015]



# Broad applicability

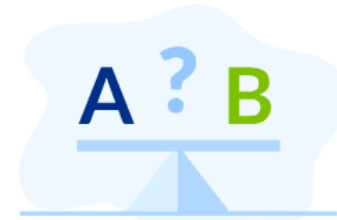
## Distributed human evaluations



Hiring



Admissions



A/B testing



Crowdsourcing



Product ratings



Healthcare



Peer grading

...

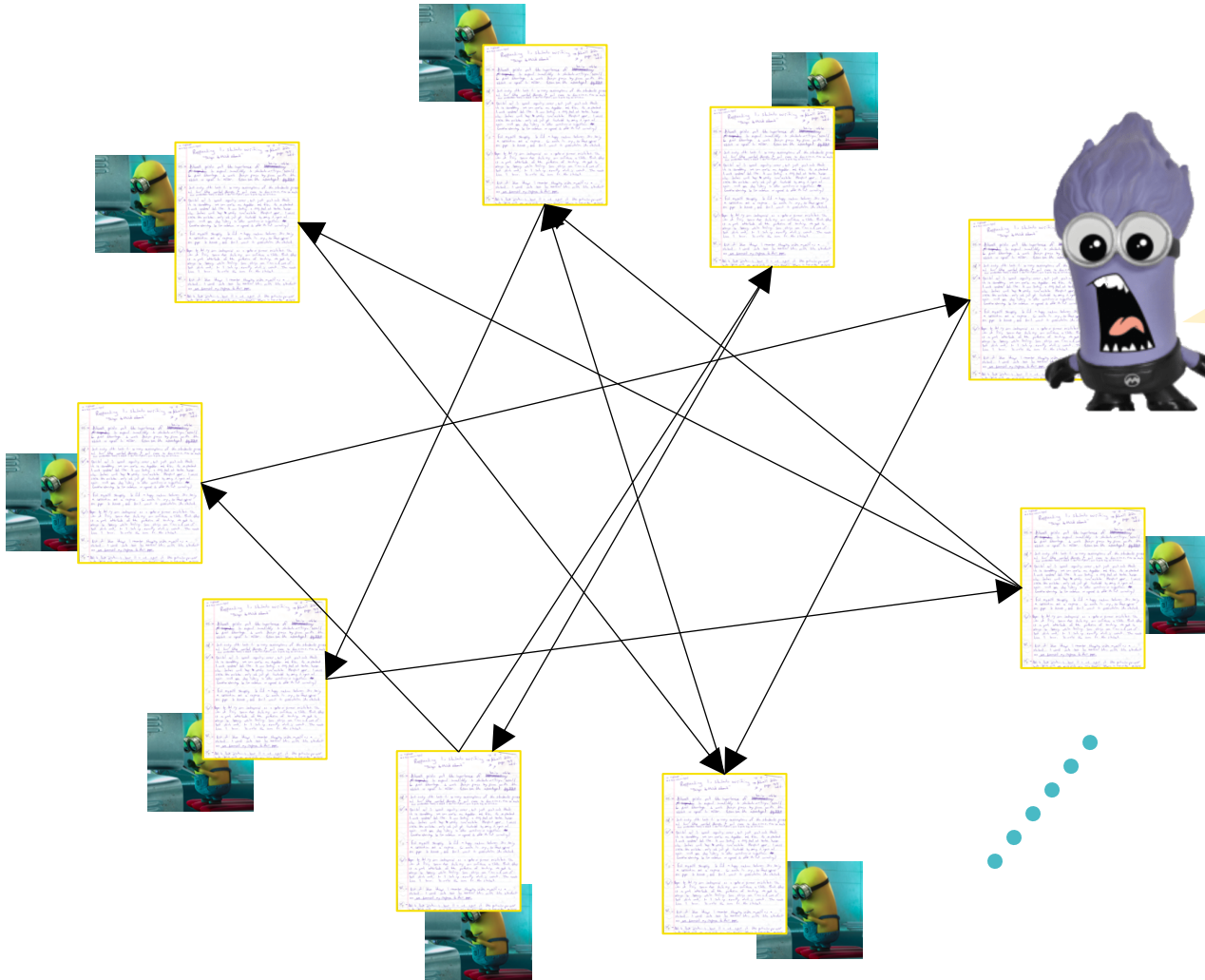
Problems amplify when this data is used to train AI/ML systems!



- **Biases**
- **Noise**
- **Dishonest behavior**
- **Miscalibration**
- **Subjectivity**
- **Norms and policies**



# Biases



**It would probably be beneficial  
to find one or two male  
researchers to work with**

True story

Review in PLOS ONE, 2015

Authors: Fiona Ingleby, Megan Head



# Single blind versus double blind

## A Principled Interpretation of Minion Speak

S. Overkill and F. Gru  
Cartoony Minion University

In this paper we present a new understanding of...

## A Principled Interpretation of Minion Speak

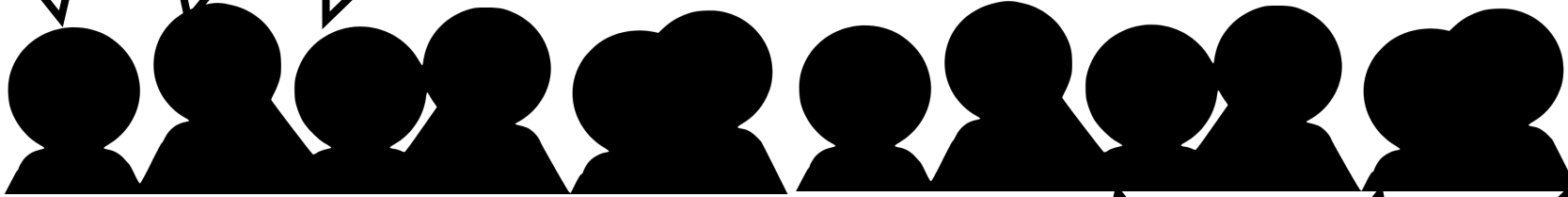
Anonymous Authors  
Anonymous Affiliation

In this paper we present a new understanding of...



# Lot of debate!

Single blind can lead to gender/fame/race/... biases



Where is the evidence of bias in my research community?

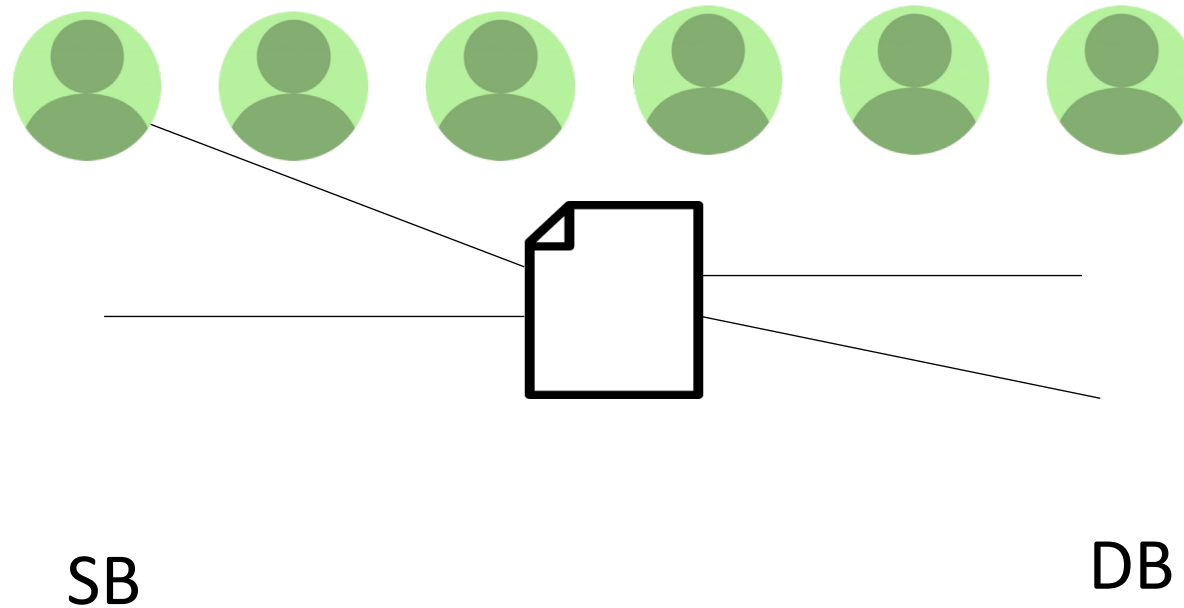


**How to rigorously test for biases in peer review  
(while ensuring “good” review process)?**



# WSDM'17 experiment: Setup

A remarkable experiment!



- Reviewers randomly split into single blind (SB) and double blind (DB) conditions
- Each paper assigned 2 SB reviewers and 2 DB reviewers



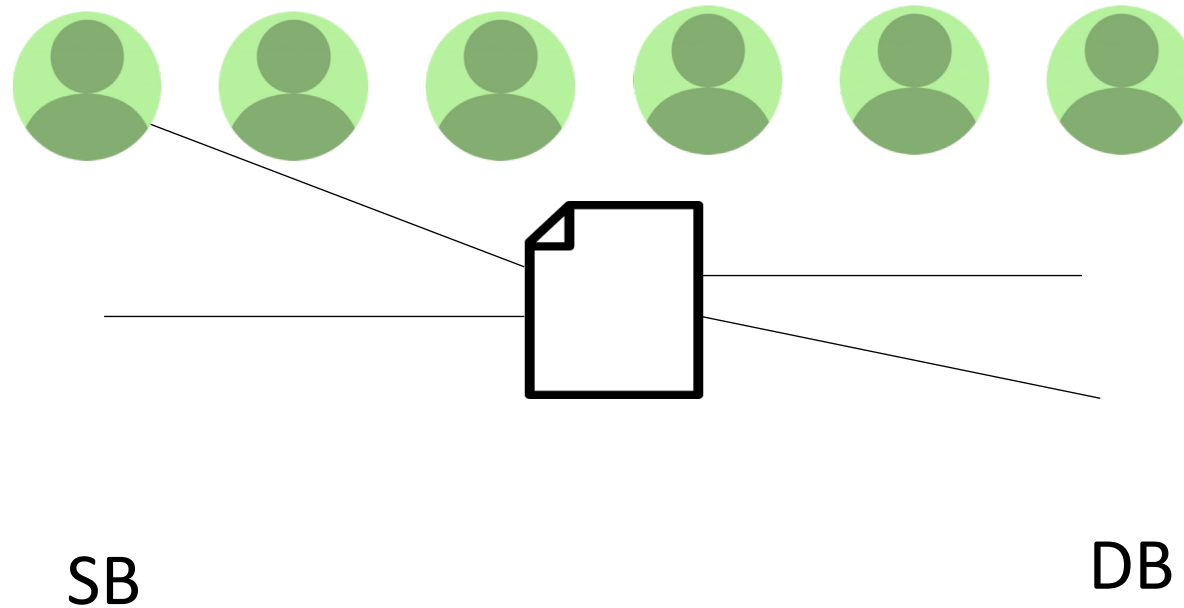
# WSDM'17 experiment: Tests for bias regarding...

- Gender
- Famous author
- Top university
- Top company
- From USA
- Academic institution
- Reviewer same country as author



# WSDM'17 experiment: Setup

A remarkable experiment!



- Reviewers randomly split into single blind (SB) and double blind (DB) conditions
- Each paper assigned 2 SB reviewers and 2 DB reviewers



# WSDM'17 experiment: Testing procedure

- Consider certain **author attributes** such as woman author, institute in USA, etc.
- For any paper  $p$ , let  $q_p$  = “intrinsic” value of paper  $p$
- **Logistic model:**  $P(\text{single blind reviewer accepts paper } p)$   
$$= \frac{1}{1 + \exp(-[\beta_0 + \beta_1 q_p + \sum_{\text{attributes } a} \beta_a \mathbb{I}\{\text{Paper } p \text{ has author attribute } a\}])}$$
- **Use DB reviewers** to estimate  $q_p$  for each paper  $p$
- **Fit decisions of SB reviewers** into logistic model to estimate  $\beta$ 's

Test:  $\beta_a = 0$  vs.  $\beta_a \neq 0$   
(no bias) (bias)



# WSDM'17 experiment: Findings

- Famous author
  - Top university
  - Top company
- } Significant bias
- At least one woman author
- } Not statistically significant; high effect size  
Meta analysis is statistically significant
- From USA
  - Academic institution
  - Reviewer same country as author
- } No evidence of bias

WSDM moved to double blind from the following year.





# **Peculiar characteristics of peer review**



# Statistical testing preliminaries

**False alarm (Type I error)** Claiming **presence** of bias when the bias is **absent**

**Detection (1 - Type-II error)** Claiming **presence** of bias when the bias is **present**

For a given  $\alpha$ , must ensure  
 $P(\text{false alarm}) \leq \alpha$

Typical choice:  $\alpha = 0.05$





## **Characteristic 0:** Correlations between quality of papers and certain attributes

- Famous author
- Top university
- Top company

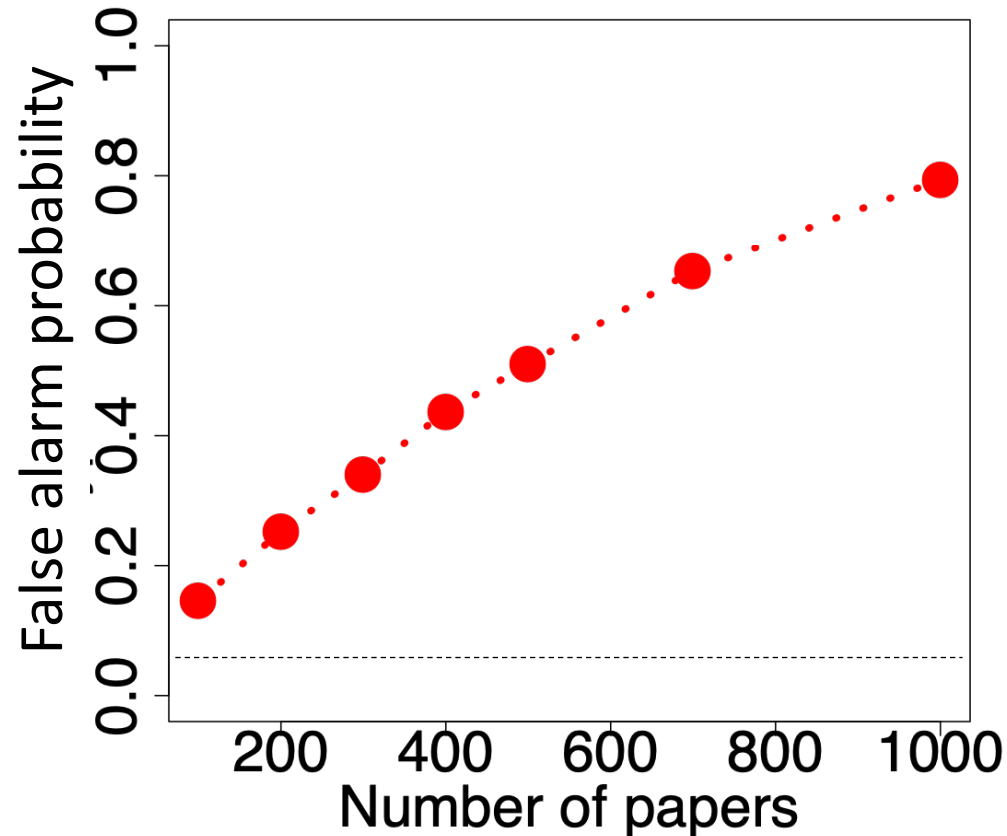
Combined with other characteristics...



# Characteristic 1: Reviews are noisy

Reviewers are noisy (and hence DB reviews are a noisy estimate of “intrinsic” value  $q_p$  of any paper  $p$ )

Must ensure:  $P(\text{declare bias when no bias}) \leq 0.05$

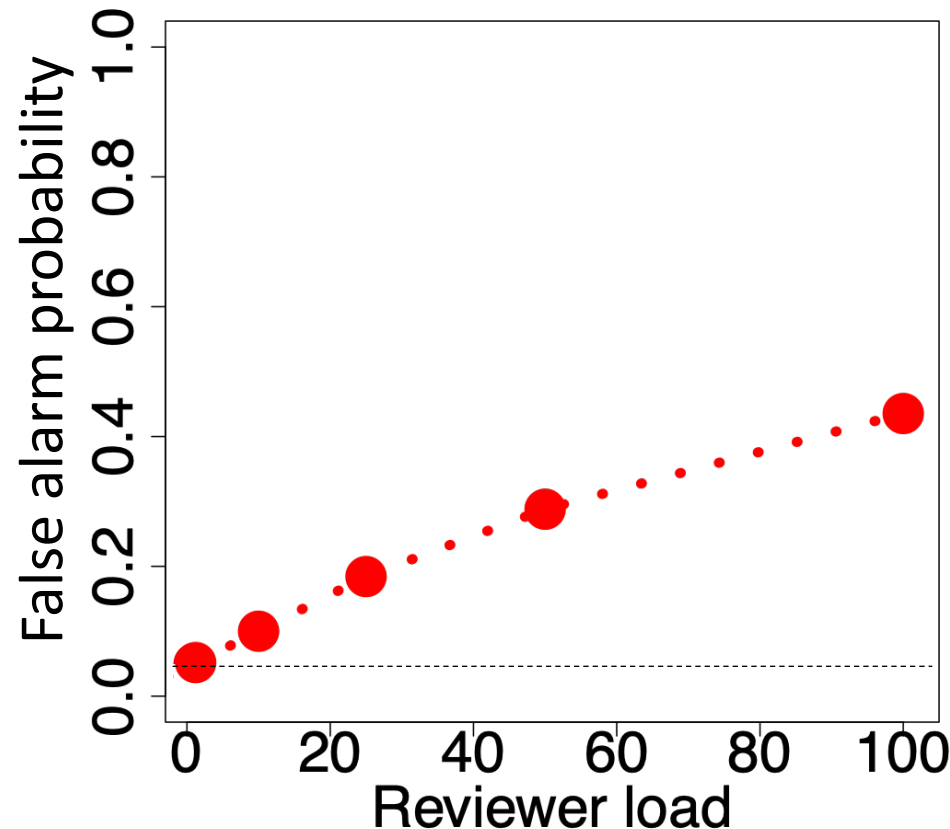




# Characteristic 2: Intra-reviewer dependency

Reviews of different papers by the same reviewer are dependent, e.g., a reviewer may be lenient or strict

Must ensure:  $P(\text{declare bias when no bias}) \leq 0.05$

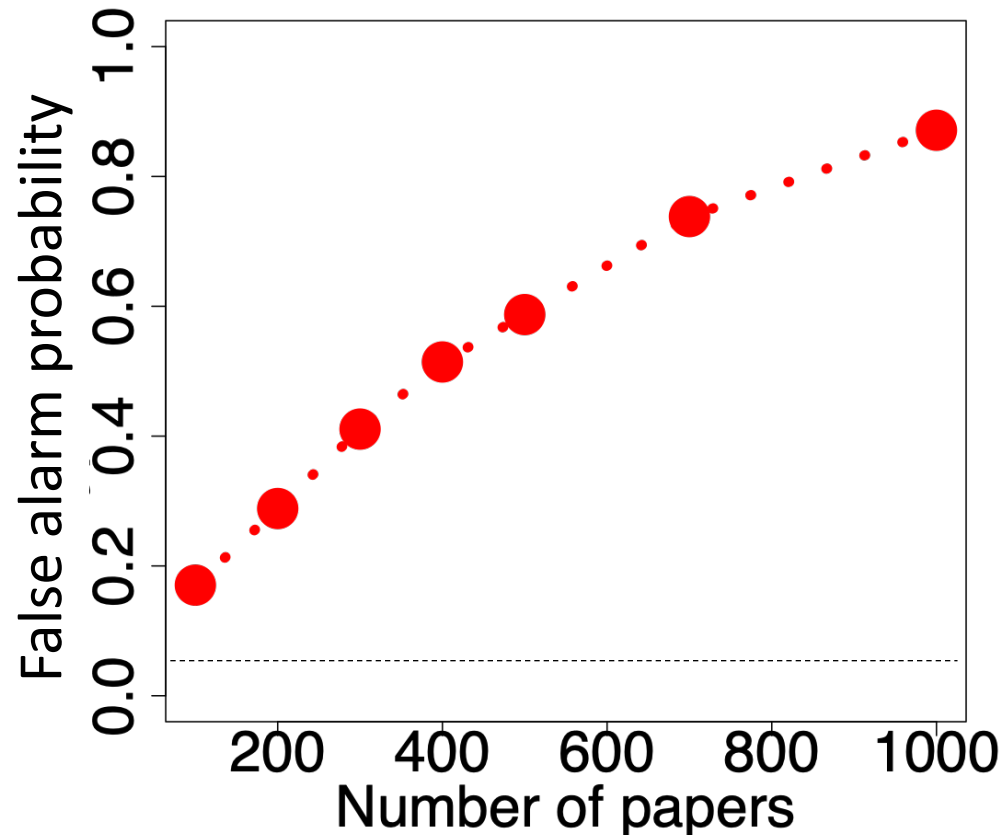




# Characteristic 3: Model complexity

Human evaluations may be more complex than the simple parametric/logistic model

Must ensure:  $P(\text{declare bias when no bias}) \leq 0.05$

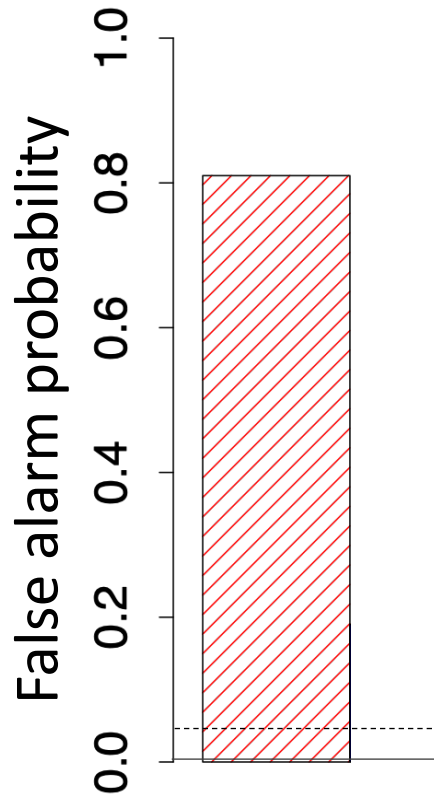




# Characteristic 4: Non-random assignment

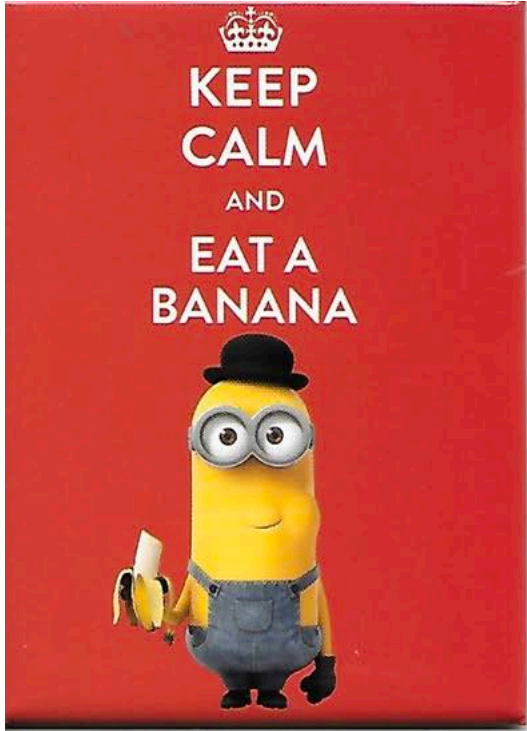
Assignment of reviewers to papers is NOT random

Must ensure:  $P(\text{declare bias when no bias}) \leq 0.05$





# These issues are fixable!



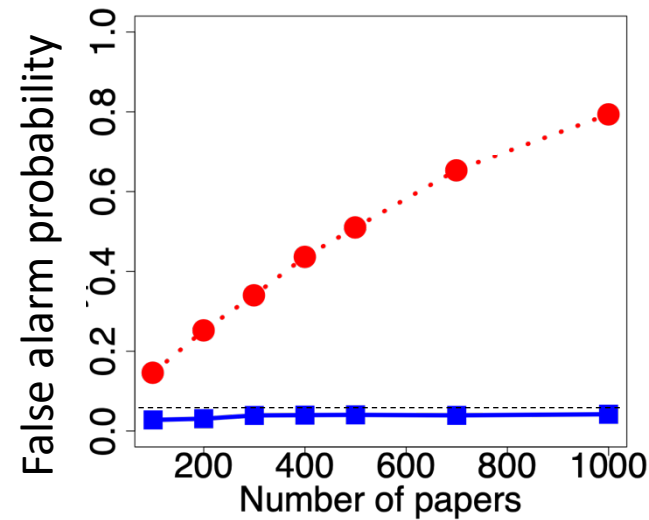
- General framework
- Careful modification of experimental setup
- Non-parametric test
- Strong theoretical guarantees:
  - False alarm control
  - Non-trivial detection



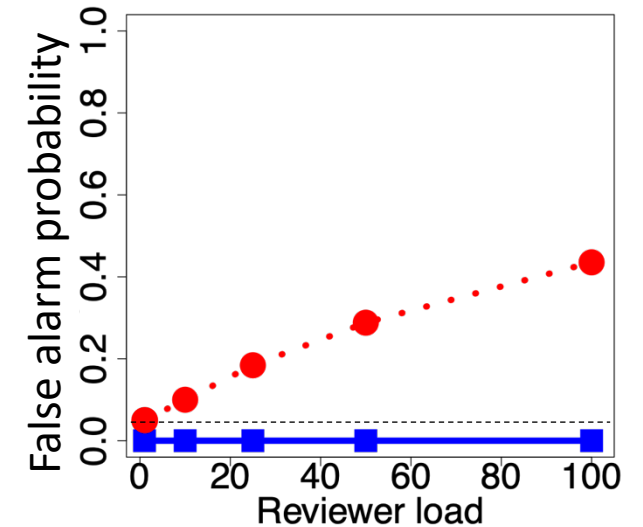
# False alarm control

● Tomkins et al.  
■ Stelmakh et al.

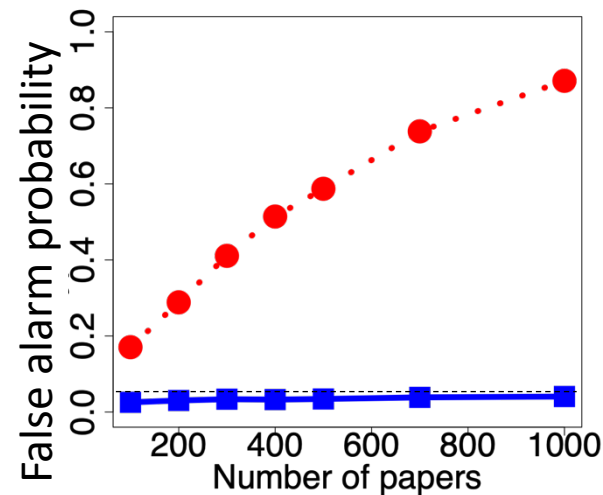
## Reviews are noisy



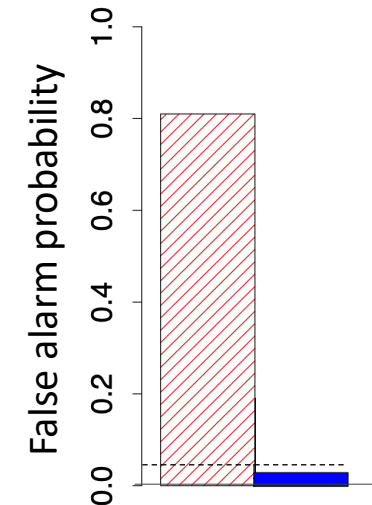
## Intra-reviewer dependency



## Model complexity



## Non-random assignment

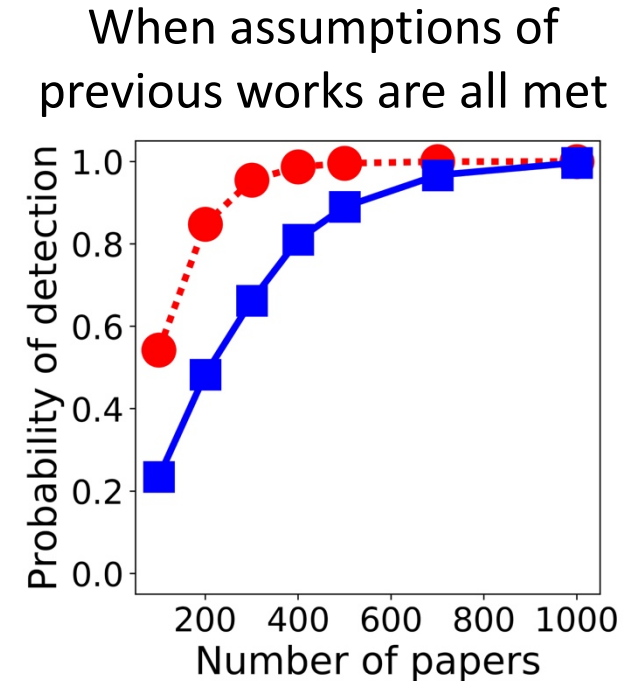
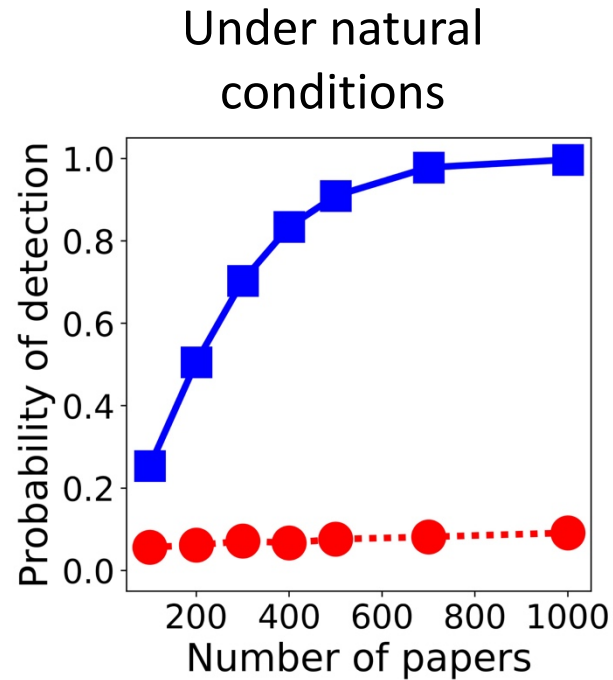




# Non-trivial detection power

(higher is better)


- Tomkins et al.
- Stelmakh et al.





# Biases: Summary and open problems



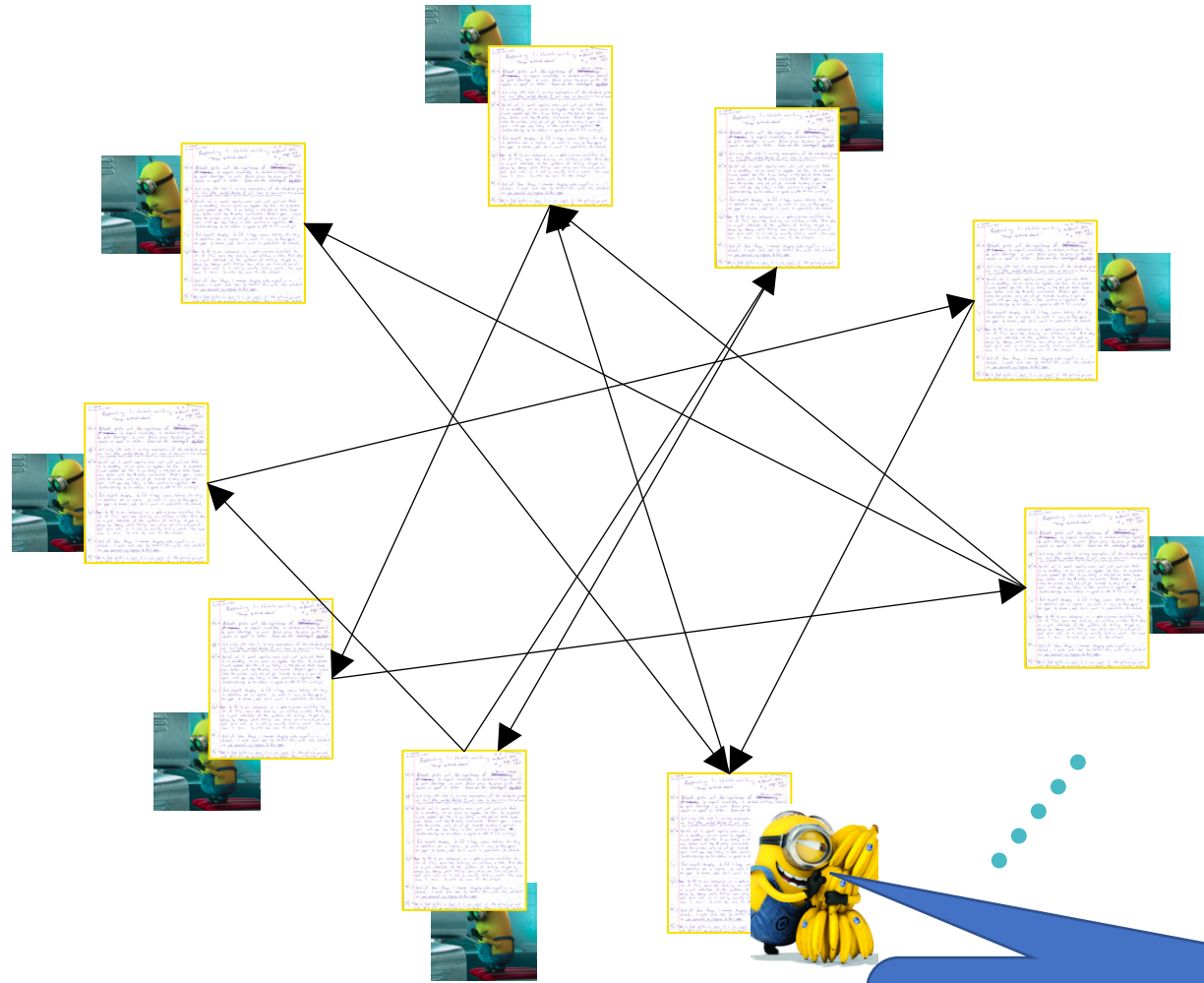
- Stelmakh et al.  → Tomkins et al.
- Need to accommodate idiosyncrasies of peer review
- Need more experiments!



- Optimal detection for given false alarm control
- Tests of biases from observational data



# Noise



I don't know much about this area.  
Weak reject I guess...



# Noise and reviewer assignment

## Poor reviews due to **inappropriate choice of reviewers**

“one of the first and **potentially most important** stages is the one that attempts to distribute submitted manuscripts to competent referees.” [[Rodriguez et al. 2007](#)]

**Top reason for dissatisfaction:** “Reviewers or panelists not expert in the field, poorly chosen, or poorly qualified” [[McCullough 1989](#)]





# Automated assignment

(Used in AAAI, NeurIPS, ICML,...)

**Compute  
similarities**

[[Mimno et al. 2007](#),  
[Rodriguez et al. 2008](#), [Charlin  
et al. 2013](#), [Liu et al. 2014](#)]



**Assignment**

- For every pair (paper  $p$ , reviewer  $r$ ), similarity score  $s_{pr} \in [0, 1]$
- Based on
  - Match text of submitted paper with reviewer's past papers
  - Match chosen subject areas
  - Bids
- Higher similarity score  $\Rightarrow$  Better envisaged quality of review
- Use similarity scores to assign reviewers to papers...



# Assignment: Maximize total similarity

(Used in AAAI, NeurIPS, ICML,...)

$$\text{maximize}_{\text{assignment}} \sum_{p \in \text{Papers}} \sum_{r \in \text{Reviewers}} s_{pr} \mathbb{I}\{\text{paper } p \text{ assigned to reviewer } r\}$$

subject to

Every paper gets at least certain #reviewers

Every reviewer gets at most certain #papers

No paper is assigned to conflicted reviewer

[Conference management systems: [TPMS \(Charlin and Zemel 2013\)](#), [EasyChair](#), [HotCRP](#)]

[[Goldsmith et al. 2007](#), [Tang et al. 2010](#), [Charlin et al. 2012](#), [Long et al. 2013](#)]



# Toy example

- One reviewer per paper
- One paper per reviewer

	Paper A	Paper B	Paper C
Reviewer 1	1	0	0.5
Reviewer 2	0.7	1	0
Reviewer 3	0	0.7	0

**Total:**

**1**

**1**

**0**

**Total:**

**0.7**

**0.7**

**0.5**

**Assignment is unfair to paper C**

**There exists another more balanced assignment**



# Common approach: Maximize total similarity

$$\underset{\text{assignment}}{\text{maximize}} \sum_{p \in \text{Papers}} \sum_{r \in \text{Reviewers}} s_{pr} \mathbb{I}\{\text{paper } i \text{ assigned to reviewer } j\}$$

- **Unbalanced:** Can assign all relevant reviewers to some papers and all irrelevant reviewers to others [Stelmakh et al. 2018]
- **Can be particularly unfair** to interdisciplinary papers
- On CVPR 2017 data, assigns at least one paper **all reviewers with 0 similarity** (there are other assignments that do much better) [Kobren et al. 2019]



# More balanced assignment

$$\begin{array}{ll} \text{maximize} & \text{minimum} \\ \text{assignment} & p \in \text{Papers} \end{array} \sum_{r \in \text{Reviewers}} s_{pr} \mathbb{I}\{\text{paper } i \text{ assigned to reviewer } j\}$$

subject to

Every paper gets at least certain #reviewers

Every reviewer gets at most certain #papers

No paper is assigned to conflicted reviewer

Fix assignment for the worst-off paper  $\operatorname{argmin}_{p \in \text{Papers}}$

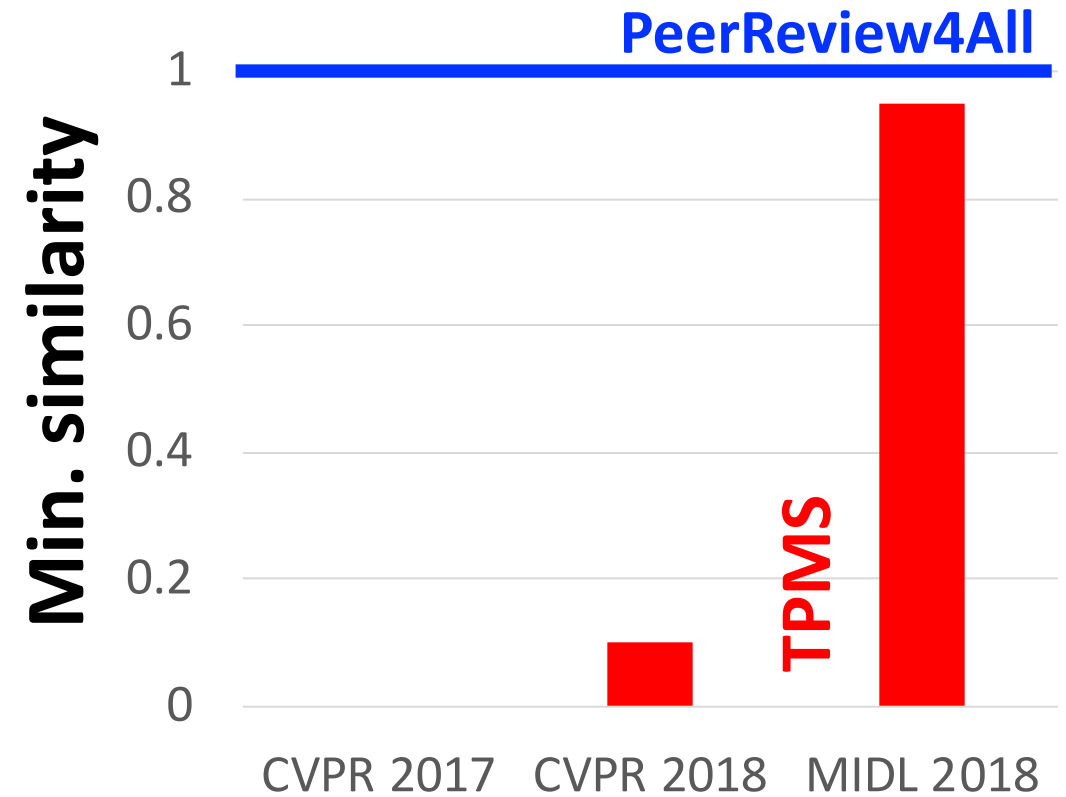
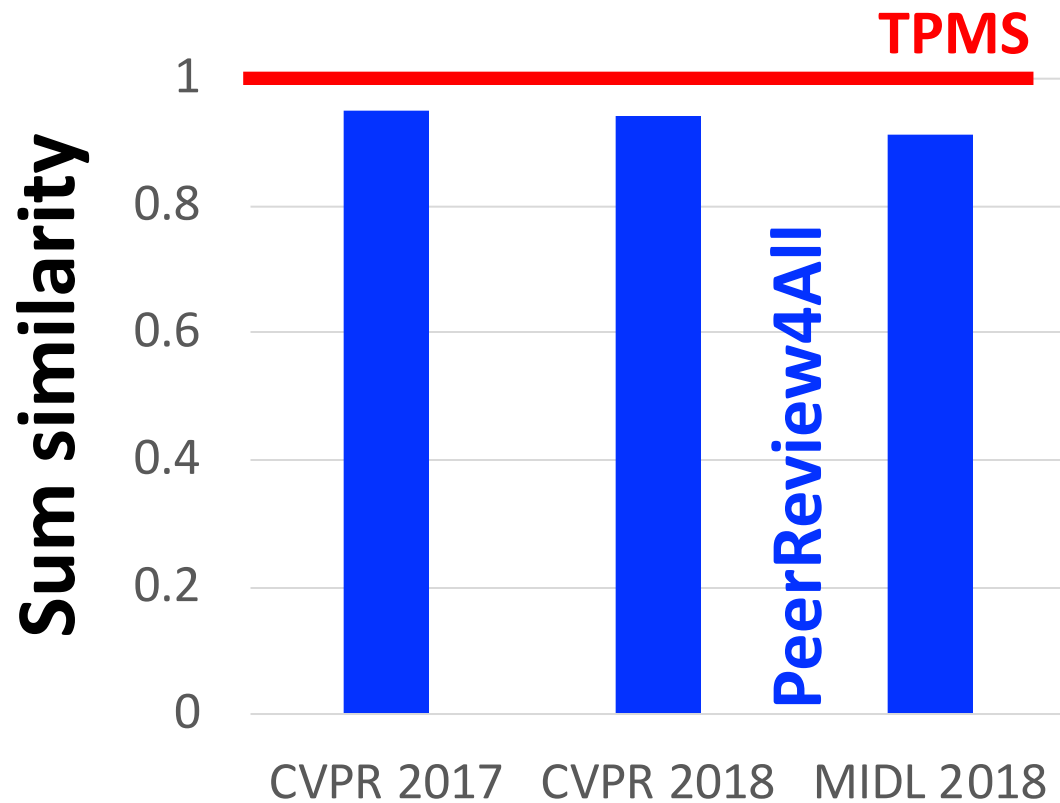
Repeat for remaining papers

- NP Hard [[Garg et al. 2010](#)]
- Approximation algorithm (“PeerReview4All”)
- Statistical guarantees on overall top-K selection



# Evaluation

- **TPMS algorithm** optimizes **sum similarity**
- **PeerReview4all algorithm** optimizes **minimum similarity**





# Noise: Summary and open problems



- Review quality significantly depends on reviewer choice
- Current automated assignments can be unfair
- Recent research on fair/balanced assignments: Theory + Practice

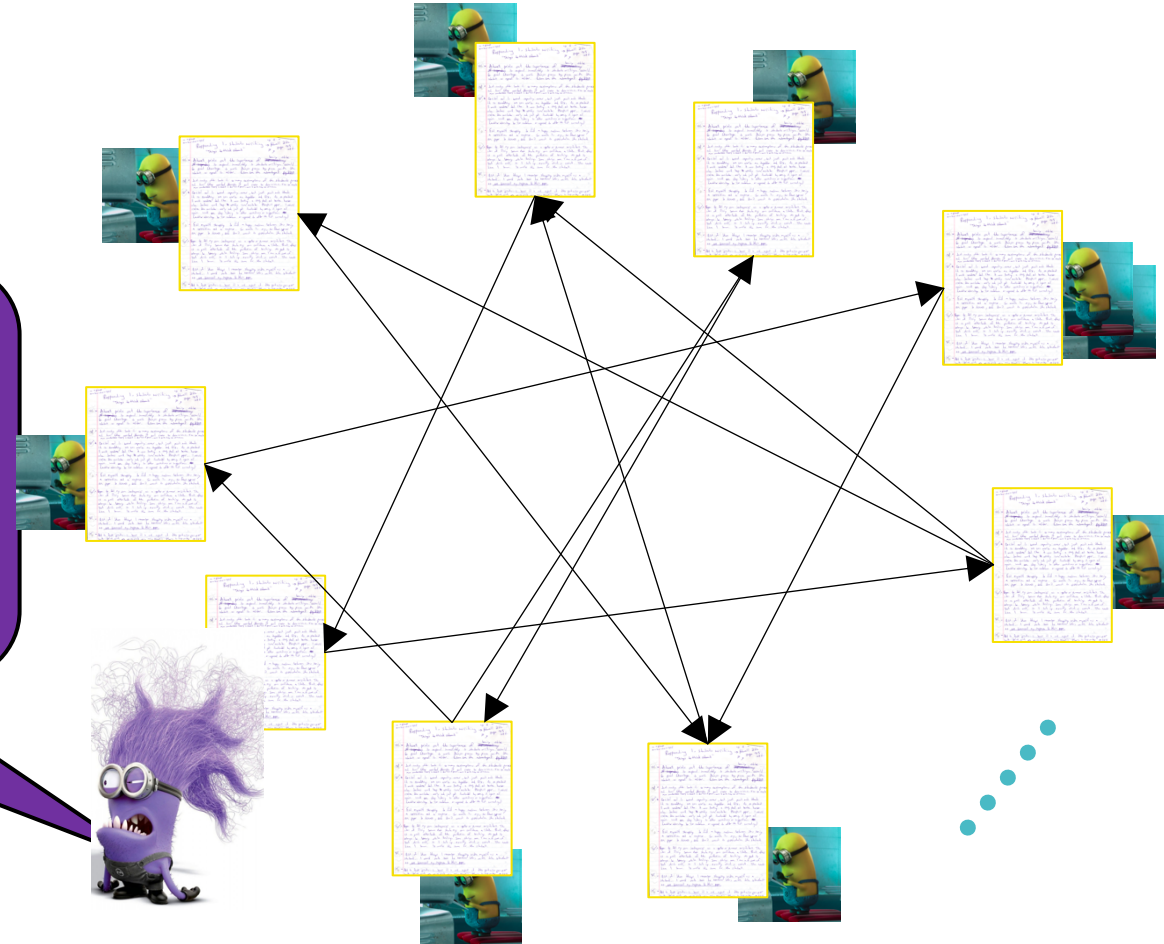


- Computationally faster fair assignment with guarantees [[Stelmakh et al. 2018](#), [Kobren et al. 2019](#)]
- Better computation of similarities; joint similarity computation and assignment [[Mimno et al. 2007](#), [Rodriguez et al. 2008](#), [Charlin et al. 2013](#), [Liu et al. 2014](#), [Tran et al. 2017](#)]
- Fair and improved bidding process [[Fiez et al. 2019](#)]



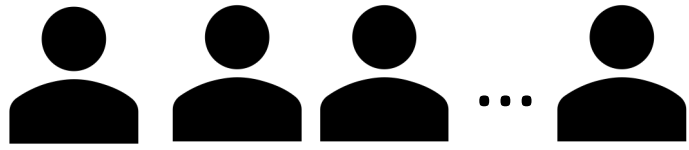
# Dishonest behavior

Giving lower scores to other papers will increase chances of my own paper getting accepted! Ha ha ha ha!





# An experiment



1. Make a painting
2. Enter one of 3 “exhibitions”
3. Peer review others’ paintings
4. Possibly win an award

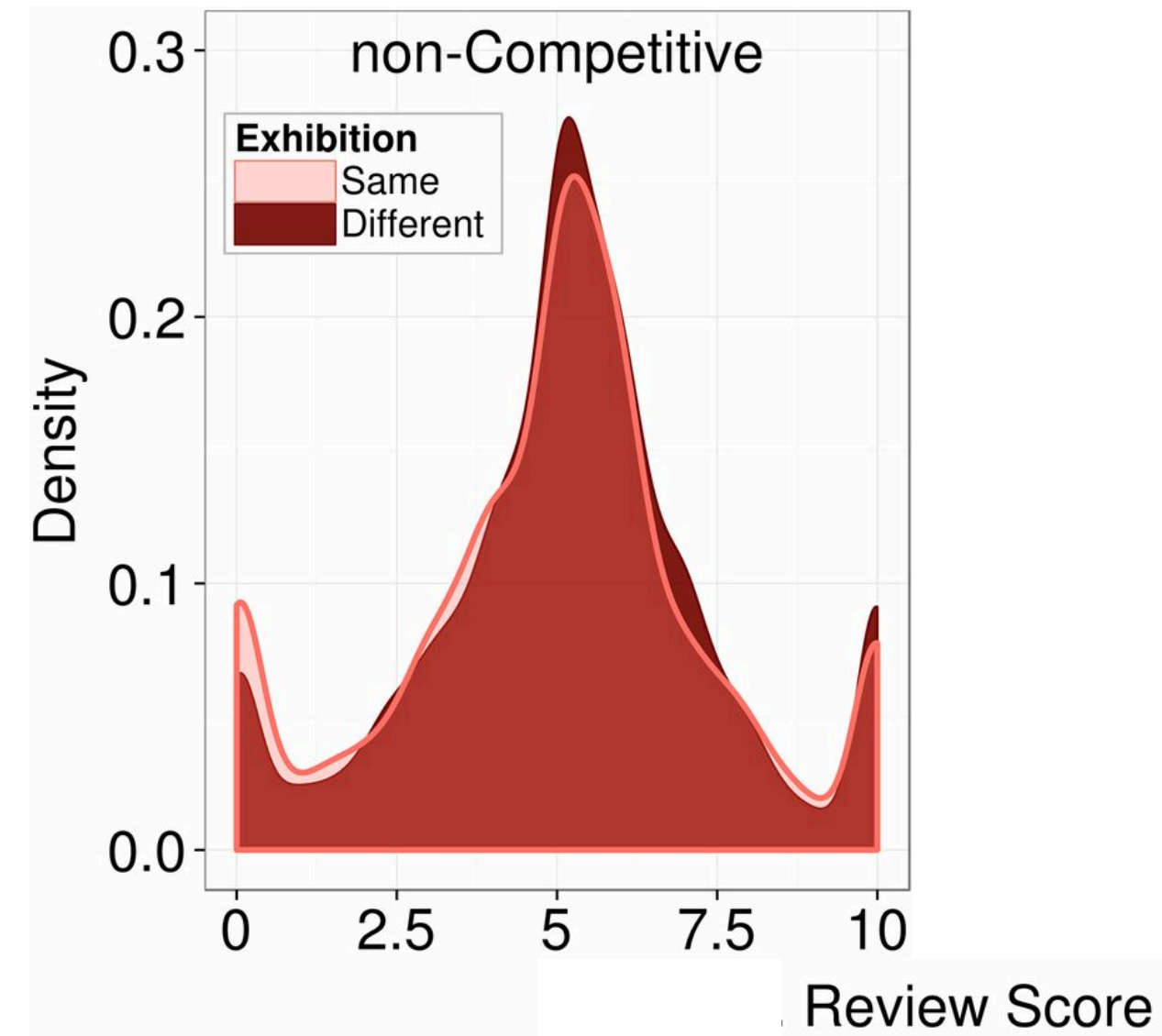
## Non-competitive

All above certain  
threshold get award

## Competitive

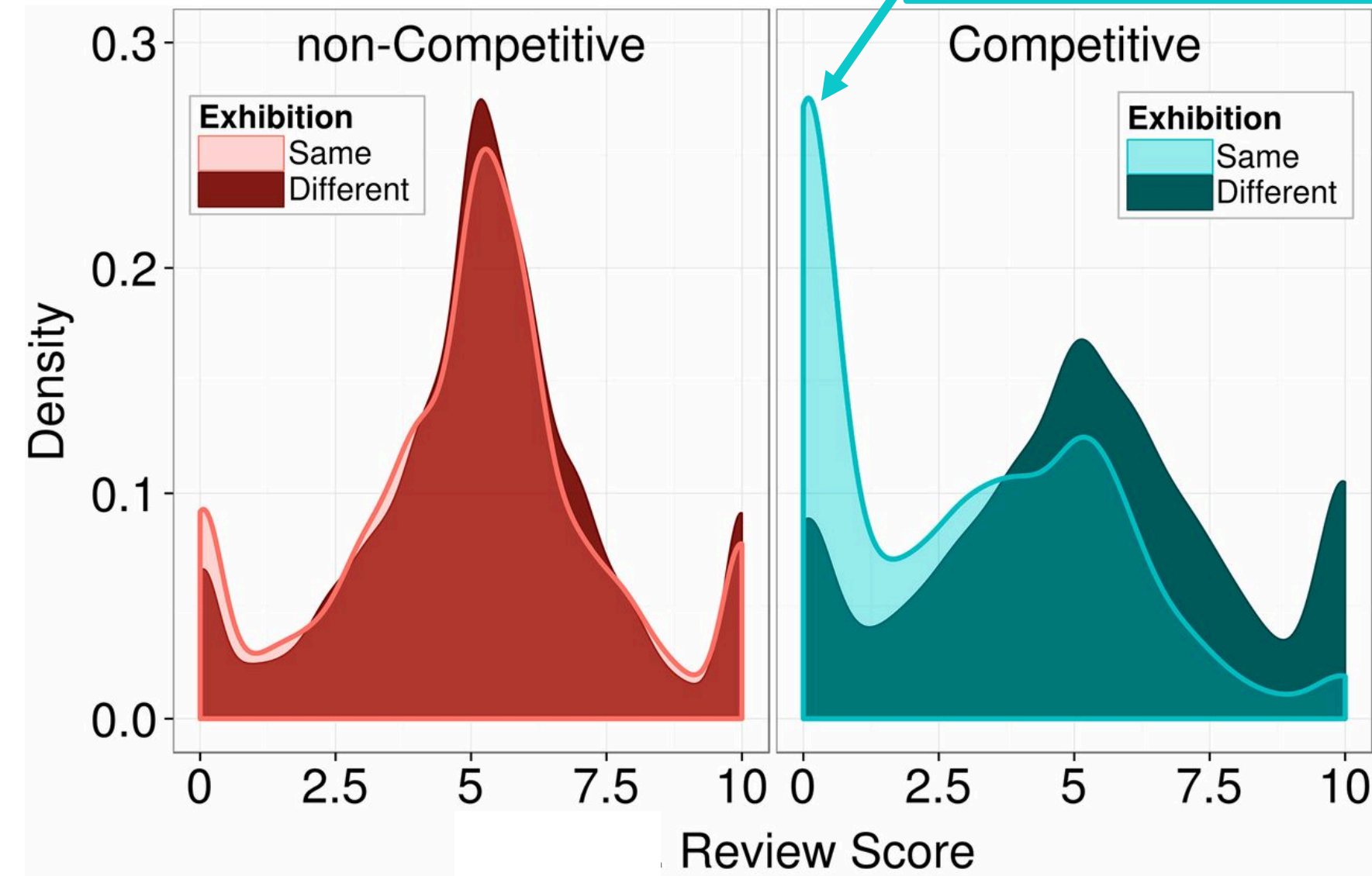
Top certain fraction in  
each exhibition win award





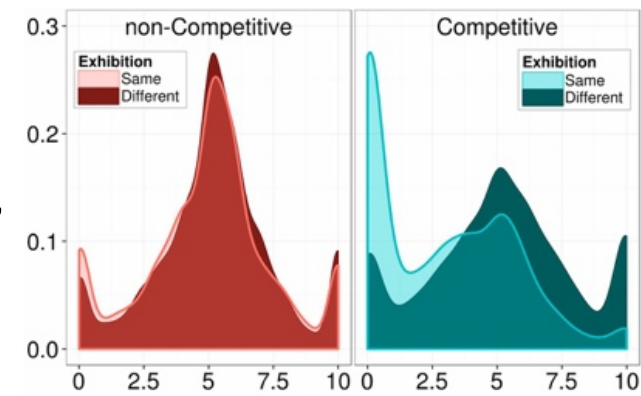


Giving a lower score increases chances of their painting getting an award





- “competitive sessions produce considerably more [strategic] reviews”
- “the number of [strategic] reviews increases over time”



**“This result provides further evidence that a substantial amount of gaming of the review system is taking place... competition incentivizes reviewers to behave strategically, which reduces the fairness of evaluations ”**

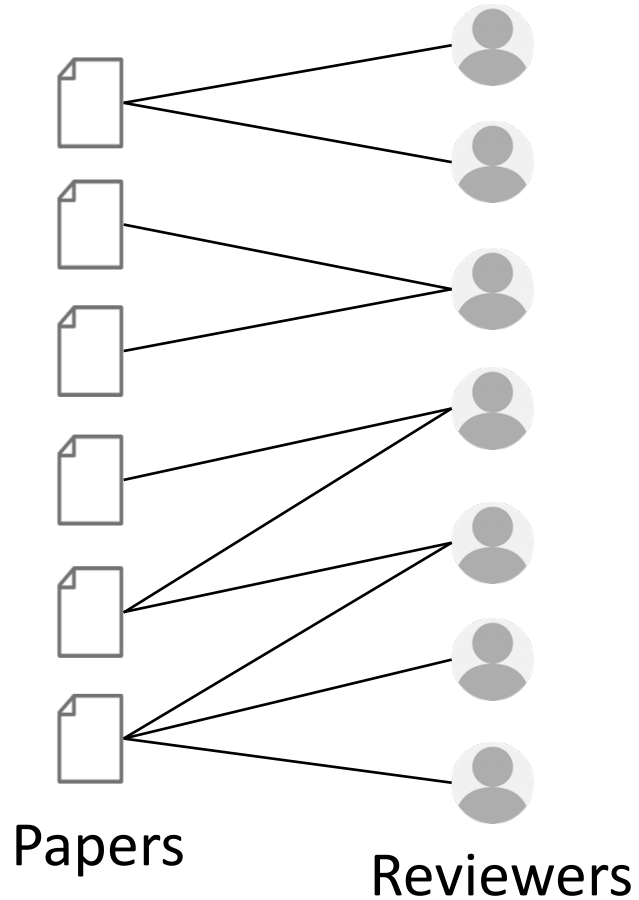
[Balietti et al., 2016]

Also [[Anderson et al. 2007](#), [Langford 2008 \(blog\)](#), [Akst 2010](#), [Thurner and Hanel 2011](#)]



# How to make peer review strategyproof?

**Given:** Conflict graph  
(e.g., authorship graph)

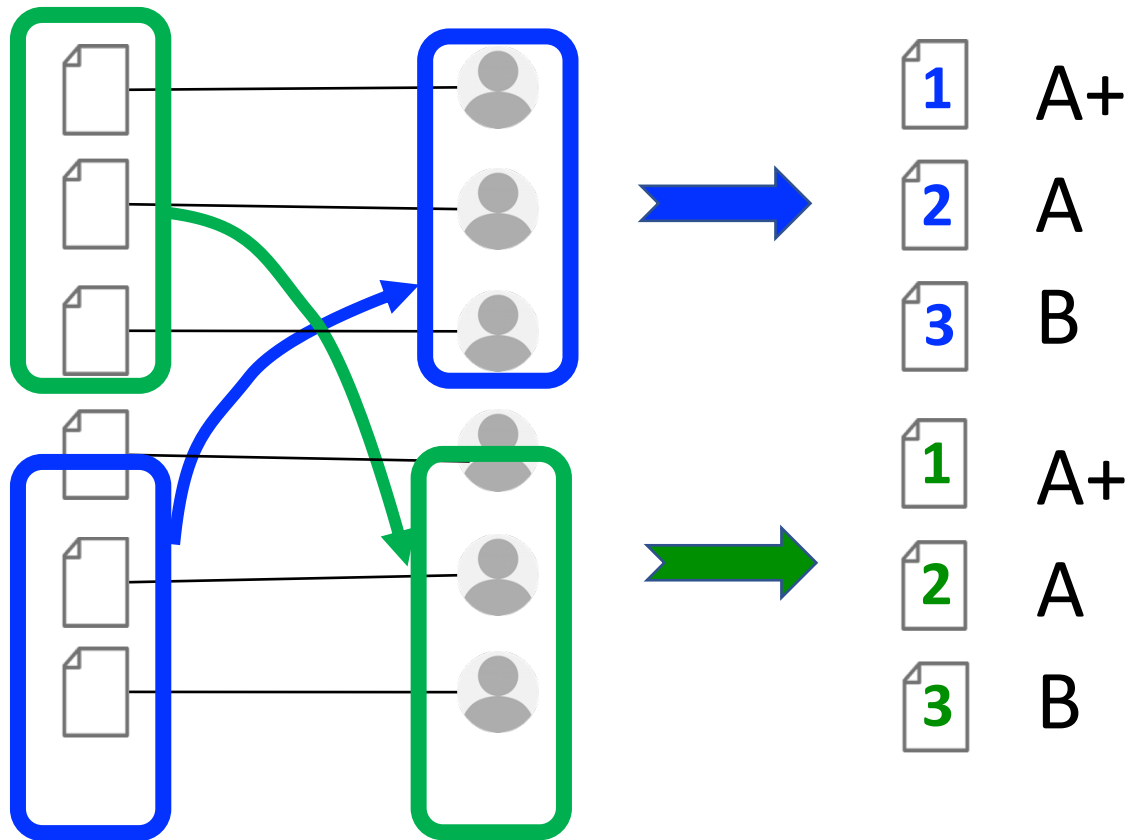


**How to ensure that no reviewer  
can influence decision of any  
conflicted paper?**



# Partitioning method

Primarily studied for peer grading

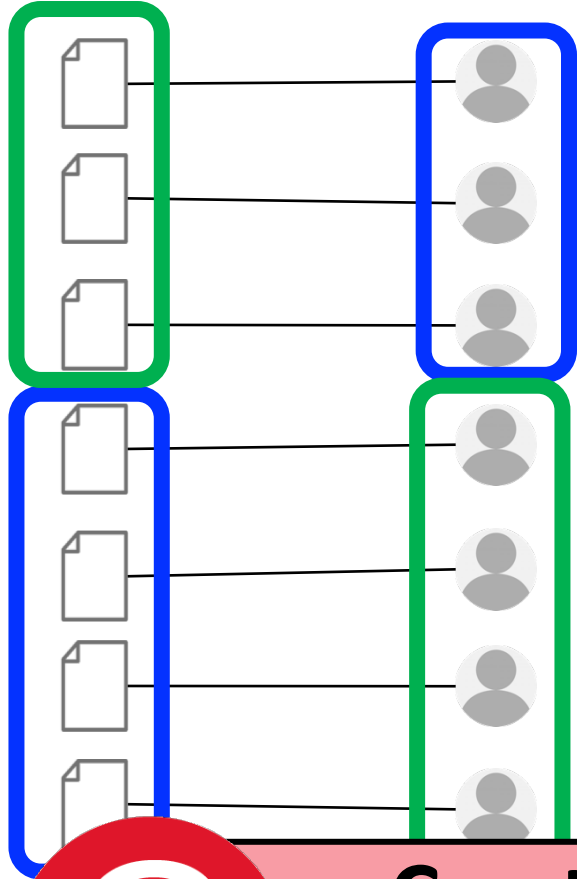


[Alon et al. 2011, Holzman et al. 2013, Bousquet et al. 2014, Fischer et al. 2015, Kurokawa et al. 2015, Kahng et al. 2017]



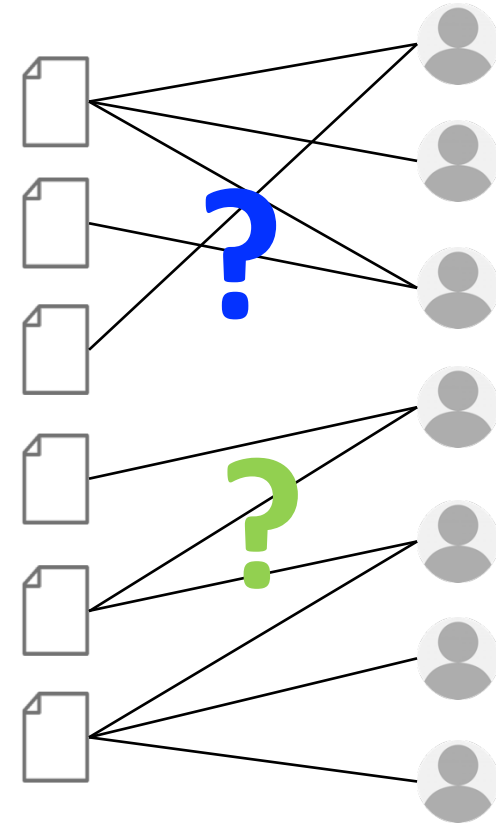
## Peer grading

1-1 conflict graphs



## Conference peer review

More **complex** conflict graphs



Can the partitioning method work  
for peer-review conflict graphs?



# ICLR empirical evaluation

## Q1. Is partitioning of conflict graph feasible?

Yes! 253 disjoint components

But, 372 authors and 133 papers in largest connected component

∴ Assigned reviewers may lack expertise.

## Q2. How does assignment quality fare under strategyproofness?

- Sum similarity reduces by 11% as compared to without strategyproofing
- Heuristics for more flexibility: Removing 3.5% of authors from the reviewer pool reduces size of the largest component by 86%



# Dishonest behavior: Summary and open problems



- Competition incentivizes strategic reviewer behavior
- Partitioning method + heuristics to mitigate its effect

- Maximize efficiency of peer review under strategyproofness

$$\begin{array}{l} \text{maximize} \\ \text{assignment} \end{array} \sum_{p \in \text{Papers}} \sum_{r \in \text{Reviewers}} s_{pr} \mathbb{I}\{\text{paper } i \text{ assigned to reviewer } j\}$$

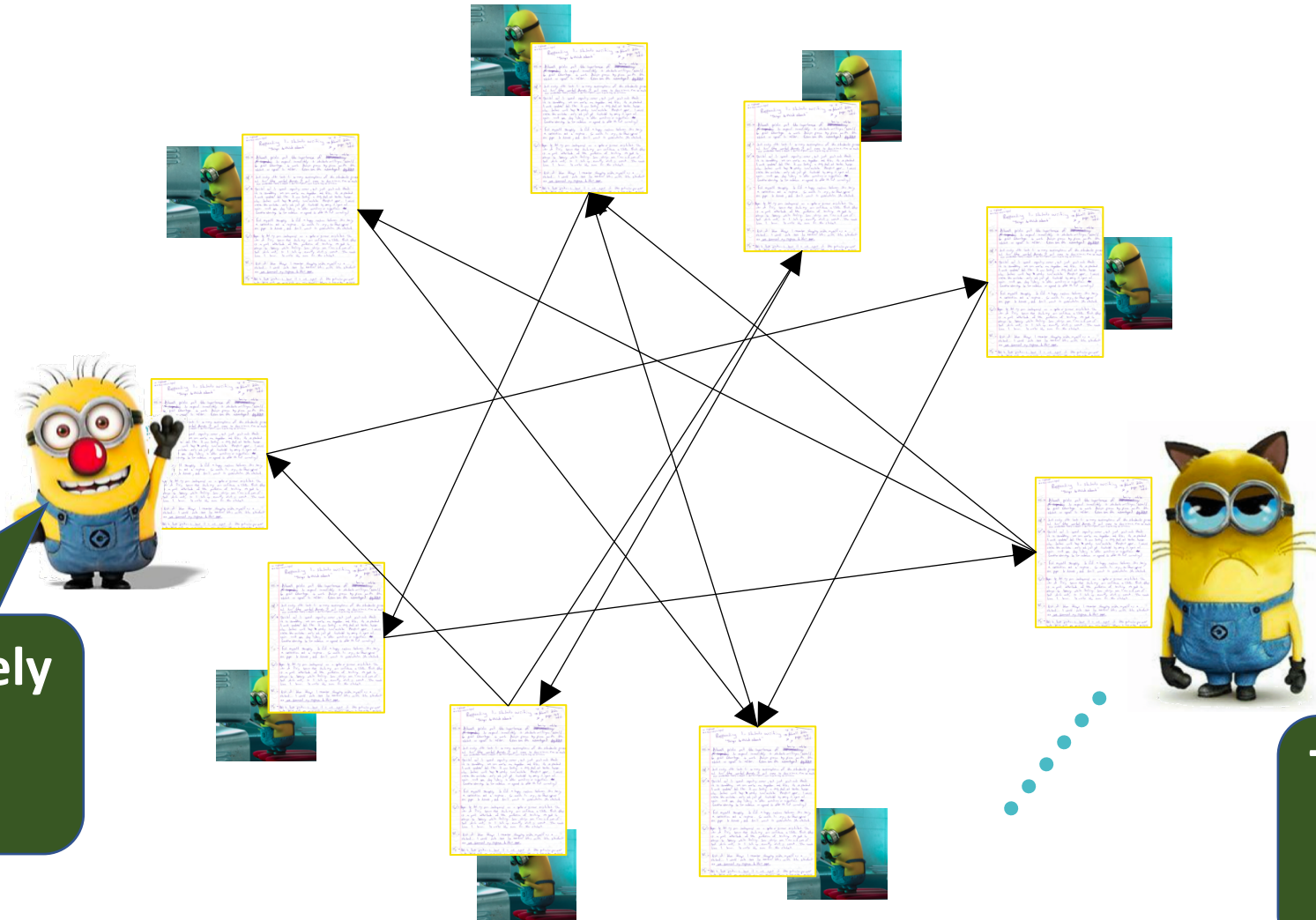
subject to strategyproofness

- Is strategyproofing possible when conflict graph cannot be partitioned? [[Aziz et al. 2019](#), [Xu et al. 2018](#)]
- Detect/prevent other forms of dishonest behavior [[Ferguson et al. 2014](#), [Gao et al. 2017](#), [Langford 2008](#)]





# Miscalibration



This is a moderately  
decent paper.  
8/10

This is a moderately  
decent paper.  
4/10.



# Miscalibration in ratings

“A raw rating of 7 out of 10 in the absence of any other information is **potentially useless**.” [Mitliagkas et al. 2011]

“The rating scale as well as the individual ratings are often **arbitrary** and may not be consistent from one user to another.” [Ammar et al. 2012]

“[Using rankings instead of ratings] becomes very important when we combine the rankings of many viewers who often use **completely different ranges of scores** to express identical preferences.” [Freund et al. 2003]



# Unfairness in peer review

“the existence of disparate categories of reviewers creates the potential for **unfair treatment of authors**. Those whose papers are sent by chance to assassins/demoters are at an unfair disadvantage, while zealots/pushovers give authors an unfair advantage.”





# Two approaches in the literature

1

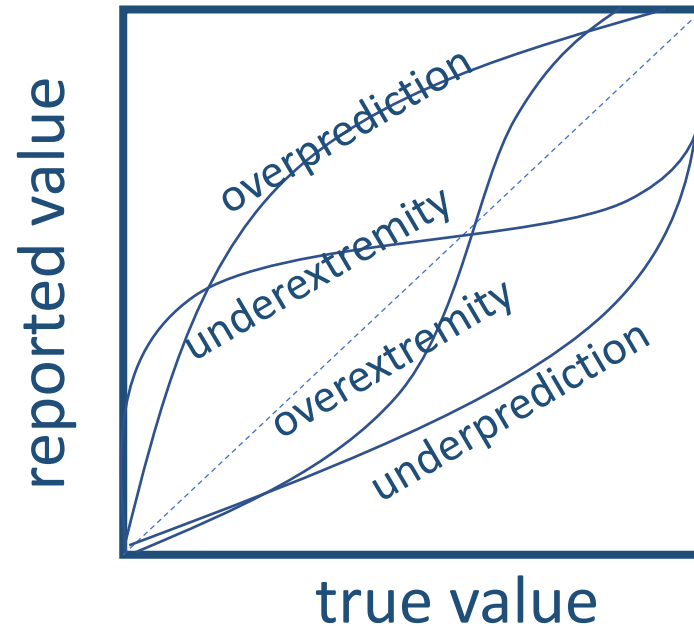
## Assume simplified (linear) models for calibration

[[Paul 1981](#), [Flach et al. 2010](#), [Roos et al. 2011](#), [Baba et al. 2013](#), [Ge et al. 2013](#), [Mackay et al. 2017](#)]

- Did not work well [NeurIPS 2016 program chairs; personal communication]
- *“We experimented with reviewer normalization and generally found it significantly harmful.”* [[Langford](#) (ICML 2012 program co-chair)]

## Miscalibration is quite complex:

[[Brenner et al. 2005](#)]





# Two approaches in the literature

2

## Use rankings

[[Rokeach 1968](#), [Freund et al. 2003](#), [Harzing et al. 2009](#),  
[Mitliagkas et al. 2011](#), [Ammar et al. 2012](#), [Negahban et al. 2012](#)]

- Use rankings induced by ratings or directly collect rankings
- Commonly believed to be the best option if no assumptions on calibration



**Is it possible to do better using ratings than rankings, with essentially no assumptions on the miscalibration?**



# Canonical 2x2 setting



$$z_A^* \neq z_B^* \in [0,1]$$



Calibration function:  $f_1 : [0,1] \rightarrow [0,1]$

Given paper  $i \in \{A, B\}$ , outputs  $f_1(z_i^*)$



Calibration function  $f_2 : [0,1] \rightarrow [0,1]$

Given paper  $i \in \{A, B\}$ , outputs  $f_2(z_i^*)$

- Adversary chooses  $z_A^*, z_B^*$  and strictly monotonic  $f_1, f_2$
- One paper assigned to each reviewer at random
- **Goal: Given (assignment, score given by each reviewer)**  
**estimate if  $z_A^* > z_B^*$  or  $z_B^* > z_A^*$** 
  - Eliciting rankings is vacuous; amounts to random guessing



# Impossibility on deterministic estimators

## Theorem

**No deterministic estimator has a success probability better than random guessing.**



# A randomized estimator

## Theorem

**There is a randomized estimator that strictly outperforms random guessing.**



With probability  $(1 + |\text{difference between the two scores}|)/2$ ,  
pick paper which received higher score

	Reviewer 1: $f_1(x) = x/2$	Reviewer 2: $f_2(x) = (3+x)/4$
Paper A: $z_A^* = 0.2$	$f_1(0.2) = 0.1$	$f_2(0.2) = 0.8$
Paper B: $z_B^* = 0.6$	$f_1(0.6) = 0.3$	$f_2(0.6) = 0.9$

- Under **blue** assignment, pick paper **B** with probability

$$\frac{1 + |0.1 - 0.9|}{2} = 0.9$$

(output is correct)

- Under **red** assignment, pick paper **A** with probability

$$\frac{1 + |0.3 - 0.8|}{2} = 0.75$$

(output is wrong)

- On average, correct with probability

$$\frac{1}{2}(0.9) + \frac{1}{2}(1 - 0.75) = 0.575 > 0.5$$



# Miscalibration: Summary and open problems



- Unfairness due to miscalibrated reviewers
- Ratings  $>$  rankings even if calibration is arbitrary/adversarial

Strong assumptions:  
parametric, linear

**Sweet spot**

Arbitrary/adversarial  
miscalibration

Ranking

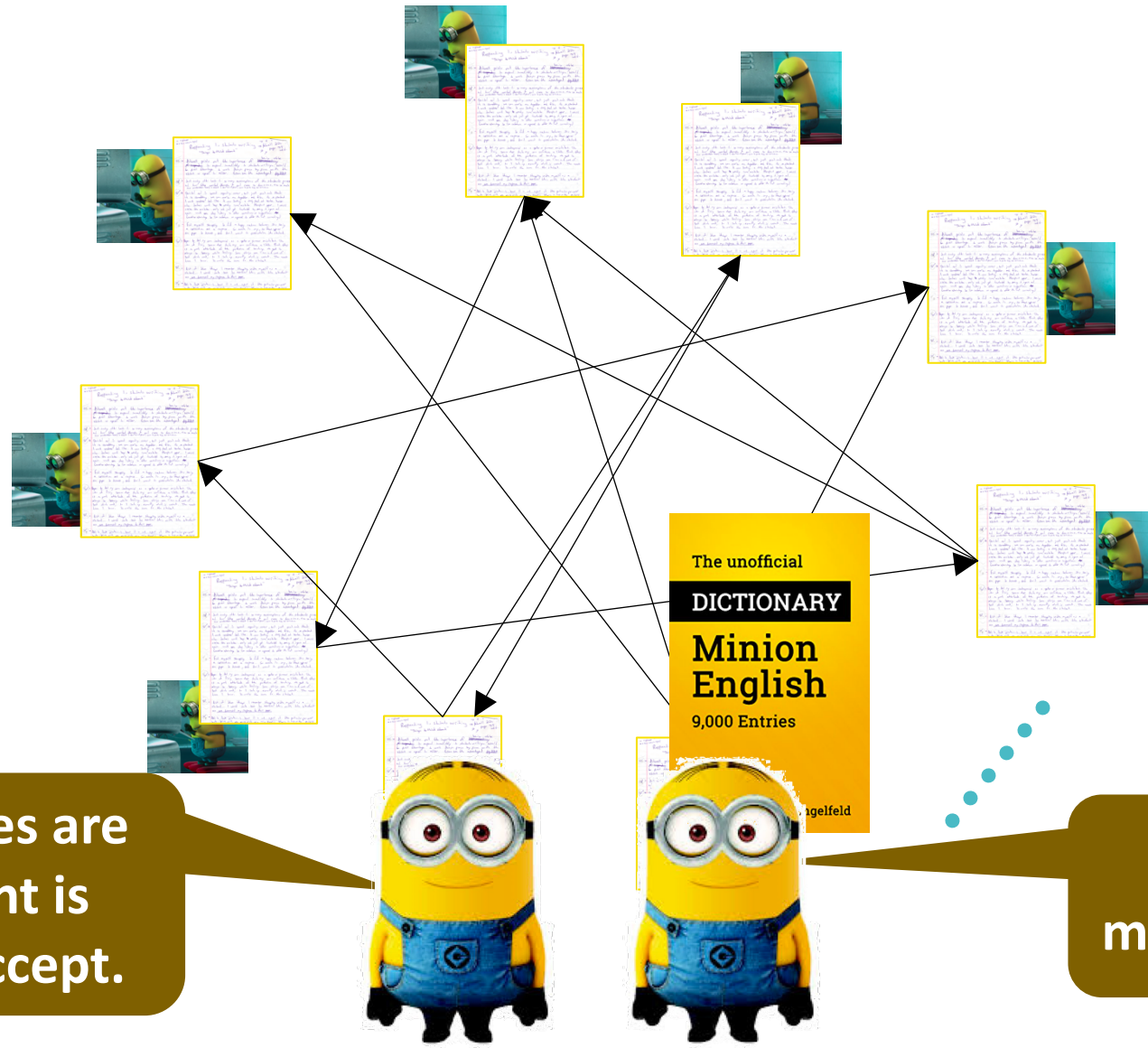
?

- Weaker assumptions: non-parametric, non linear (e.g., permutation-based models [[Shah 2017 part 1](#)])
- Amenable to given sample sizes: Avoid overkill





# Subjectivity

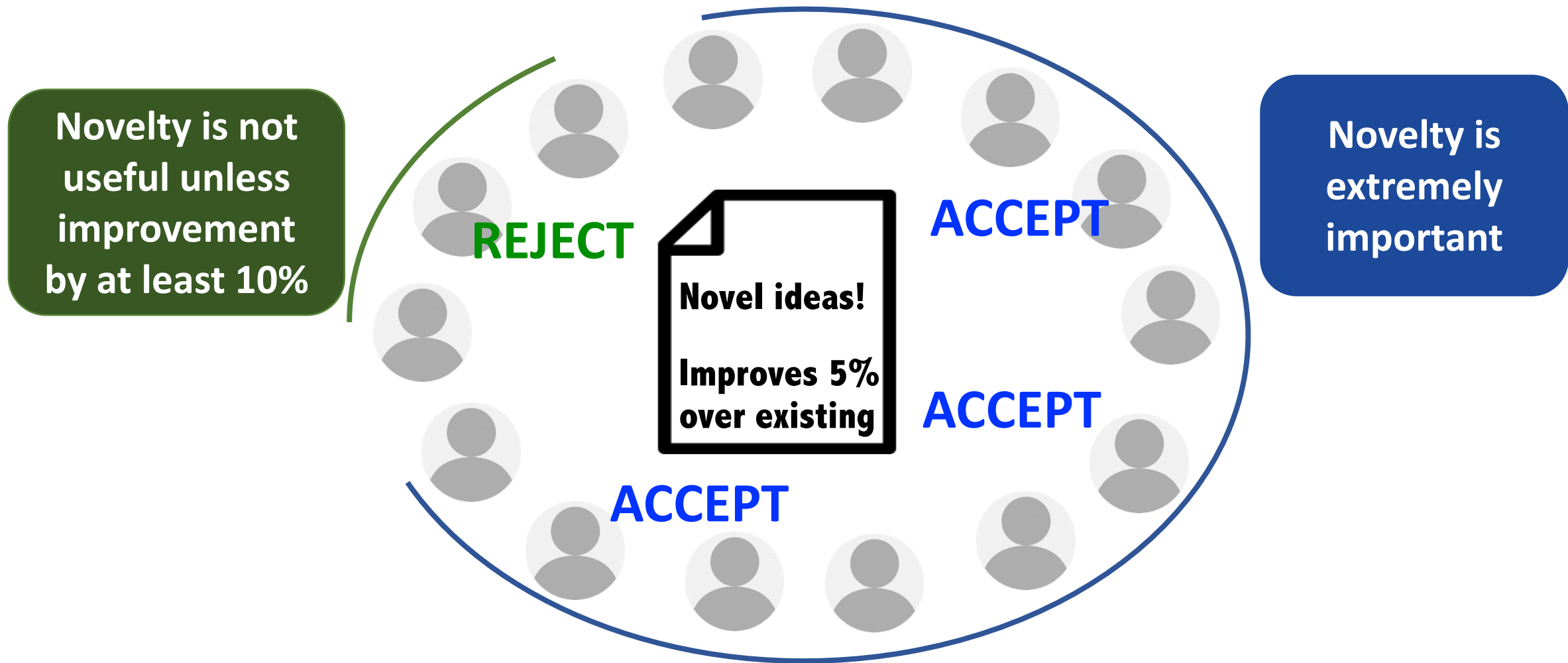


Spelling mistakes are  
ok. The content is  
great. Strong accept.

Too many spelling  
mistakes. Strong reject.



# Differing opinions about relative importance of criteria





# Commensuration Bias in Peer Review

Carole J. Lee\*†

---

To arrive at their final evaluation of a manuscript or grant proposal, reviewers must convert a submission's strengths and weaknesses for heterogeneous peer review criteria into a single metric of quality or merit. I identify this process of commensuration as the

“Illuminates how intellectual priorities in individual peer review judgments can collectively subvert the attainment of community-wide goals”





**How to ensure that every paper is  
judged by the same yardstick?**

An approach using ML and social choice theory....



# Problem setting

- Reviewers asked to judge papers on **k criteria**
  - E.g. (IJCAI 17): Originality, Relevance, Significance, Writing, Technical
  - Give **criteria scores** in  $[0,1]^k$
- And an **overall score** in  $[0,1]$
- Each reviewer has a coordinate-wise non-decreasing **(subjective) mapping** from criteria scores in  $[0,1]^k$  to overall score in  $[0,1]$

**Need a common mapping for all papers**



# Data-driven approach: Learn a mapping

- Obtain (criteria scores, overall score)  $\in [0,1]^k \times [0,1]$  for every review
- Learn a mapping  $\hat{f}: [0,1]^k \rightarrow [0,1]$  from this data
- For every review, replace overall score with  $\hat{f}(\text{criteria scores})$



# Framework

## Which loss?

$$\hat{f} \in \underset{f \in \mathcal{F}}{\operatorname{argmin}} \ell \left( \begin{array}{cc} \text{Reviewers} & \text{Papers} \\ \begin{bmatrix} f([.8 \ .9 \ .9]) & f([.8 \ .6 \ .1]) \\ f([.9 \ .1 \ .4]) & f([.2 \ .7 \ .4]) \\ f([.1 \ .3 \ .2]) & f([.8 \ .9 \ .2]) \\ f([.9 \ .5 \ .4]) & f([.7 \ .8 \ .2]) \\ f([.3 \ .2 \ .1]) & f([.4 \ .7 \ .9]) \end{bmatrix} & \begin{bmatrix} .9 & .6 \\ .2 & .4 \\ .6 & .6 \\ .4 & .9 \\ .2 & .3 \end{bmatrix} \end{array} \right)$$

Criteria scores                      Overall scores

$\mathcal{F}$  = set of all coordinate-wise non-decreasing functions



# Theorem

Under three ‘reasonable’ axioms, there is **exactly one possible loss**.

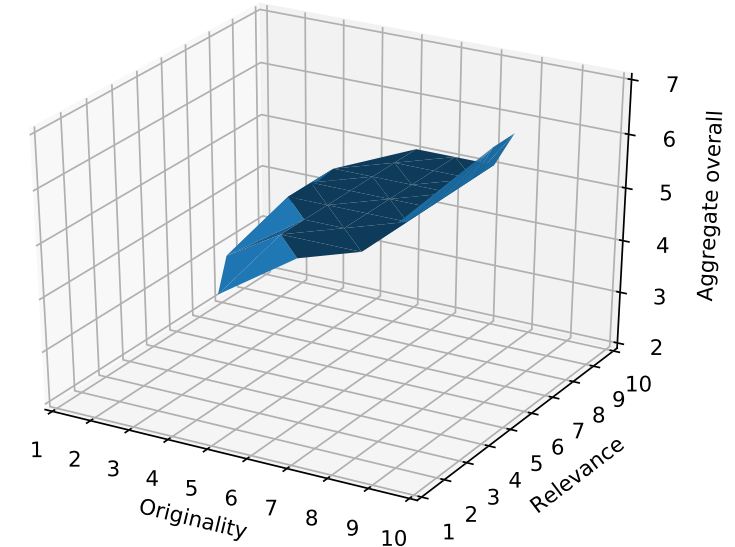
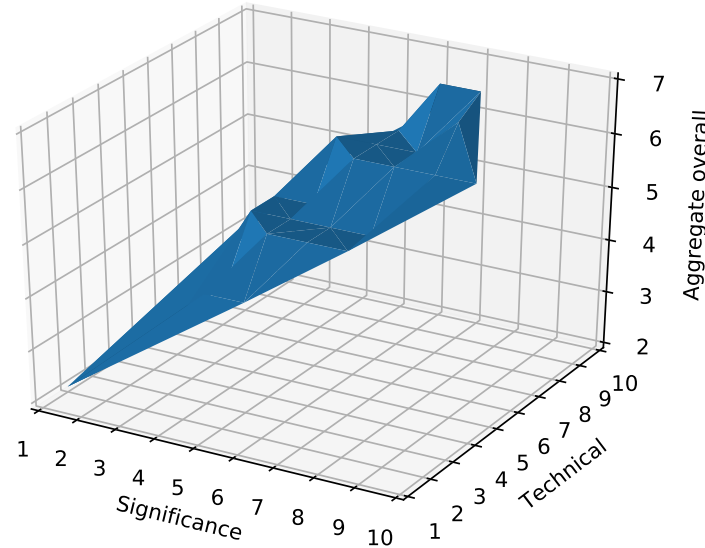
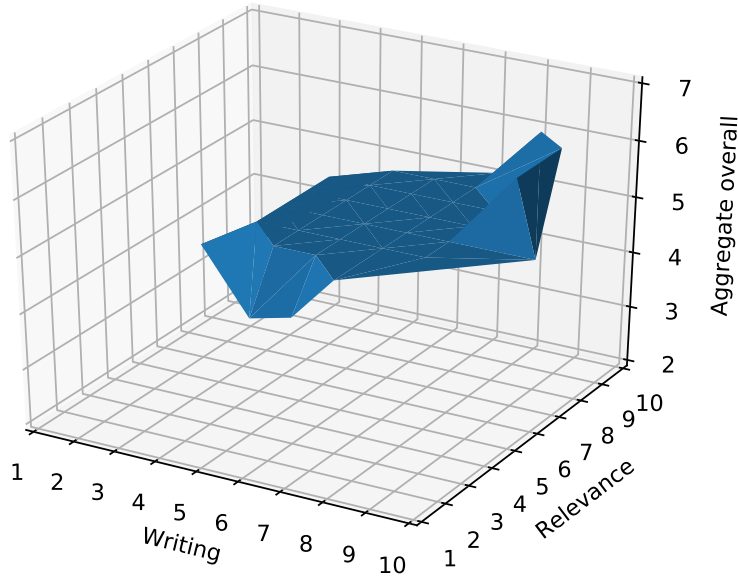
“L(1,1) loss, i.e., sum of absolute differences of all entries”

$$\hat{f} \in \operatorname{argmin}_{f \in \mathcal{F}} \sum_{\text{all entries}} \left| \begin{pmatrix} f([.8 \ .9 \ .9]) & f([.8 \ .6 \ .1]) \\ f([.9 \ .1 \ .4]) & f([.2 \ .7 \ .4]) \\ & f([.1 \ .3 \ .2]) & f([.8 \ .9 \ .2]) \\ f([.9 \ .5 \ .4]) & f([.7 \ .8 \ .2]) \\ f([.3 \ .2 \ .1]) & f([.4 \ .7 \ .9]) \end{pmatrix} - \begin{pmatrix} .9 & .6 \\ .2 & .4 \\ & .6 & .6 \\ .4 & .9 \\ .2 & .3 \end{pmatrix} \right|$$

$\mathcal{F}$  = set of all coordinate-wise non-decreasing functions



# IJCAI 2017



- **Writing** and **Relevance**: Really bad - significant downside, really good - appreciated, in between - irrelevant.
- **Technical** quality and **Significance**: high influence; the influence is approximately linear.
- **Originality**: moderate influence.



# Subjectivity: Summary and open problems



- Commensuration biases put individual reviewer preferences above community-wide norms and preferences
- “Learn” community-wide preferences
  - Using ideas from machine learning and social choice theory



- Evaluation in absence of ground truth
- Multiple problems together



# Norms and Policies

Alright, so here's  
what everyone  
must do...





# Biases due to alphabetical ordering

In Economics, norm is to order authors in alphabetical order of last names.

**Faculty with last name starting with an earlier alphabet are:**

- Significantly more likely to receive tenure
- Significantly more likely to become fellows of the Econometric Society
- More likely to receive the Clark Medal and the Nobel Prize

The (related) field of Psychology, which does not order by alphabet, does not show any of these biases.





# What causes these biases?

## In papers

Implicit bias – Primacy effects

Explicit bias – “*First author et al.*”

Conference	#Total papers	#Papers using “ <i>First author et al.</i> ” in its text
STOC 2017	99	70
STOC 2016	79	59
FOCS 2017	79	48
FOCS 2016	73	43
EC 2017	75	48
EC 2016	99	87

## On websites

Serial position effects



### AAAI-19 Program Committee

Hussein Abbass  
Sherief Abdallah  
Abbas Abdolmaleki  
Naoki Abe  
David Abel  
Ayan Acharya  
Maribel Acosta  
Shuchin Aeron  
Maedeh Aghaei



# Let's fix this!

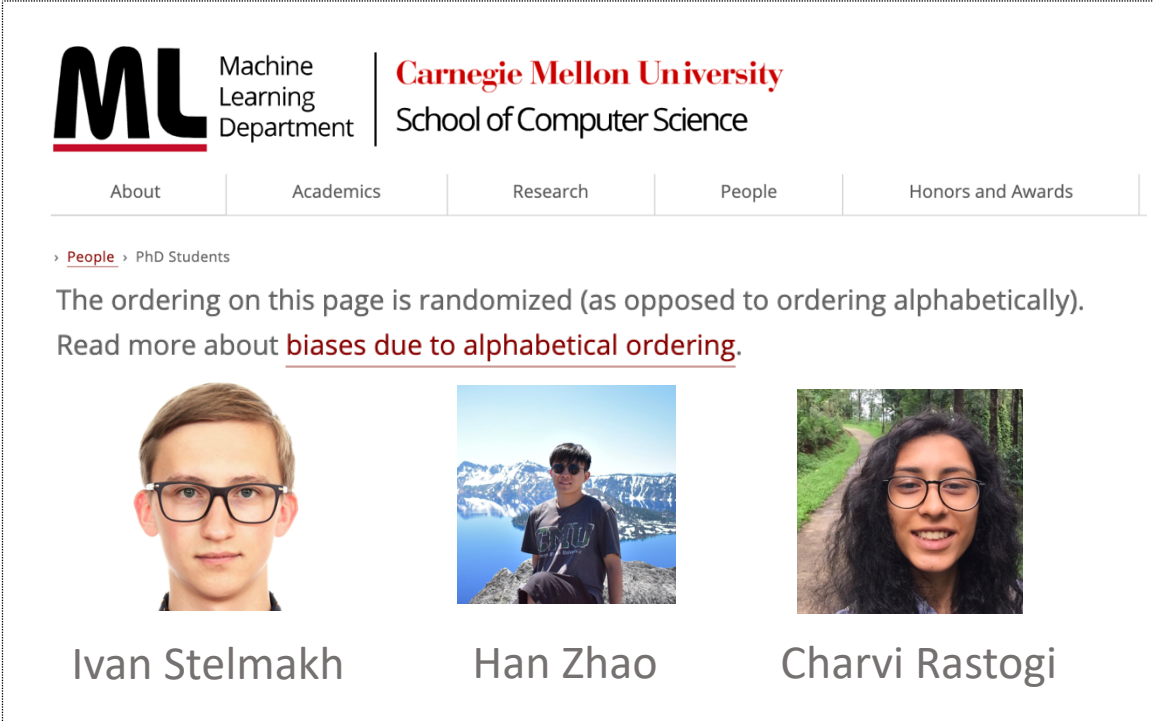
## In papers

ACM EC conference now uses numbering instead of “first author et al.” citation style

Can randomize author ordering

## On websites

CMU Machine Learning Department website now uses dynamic randomization for ordering people



The screenshot shows the header of the CMU Machine Learning Department website. The header includes the 'ML' logo, 'Machine Learning Department', 'Carnegie Mellon University', and 'School of Computer Science'. Below the header is a navigation bar with links: 'About', 'Academics', 'Research', 'People', and 'Honors and Awards'. The 'People' link is active, and the page title is 'People > PhD Students'. The main content area states: 'The ordering on this page is randomized (as opposed to ordering alphabetically). Read more about [biases due to alphabetical ordering](#).' Below this text are three profile pictures of PhD students: Ivan Stelmakh, Han Zhao, and Charvi Rastogi. The names are listed below their respective photos.

Machine Learning Department | Carnegie Mellon University  
School of Computer Science

About | Academics | Research | People | Honors and Awards

> People > PhD Students

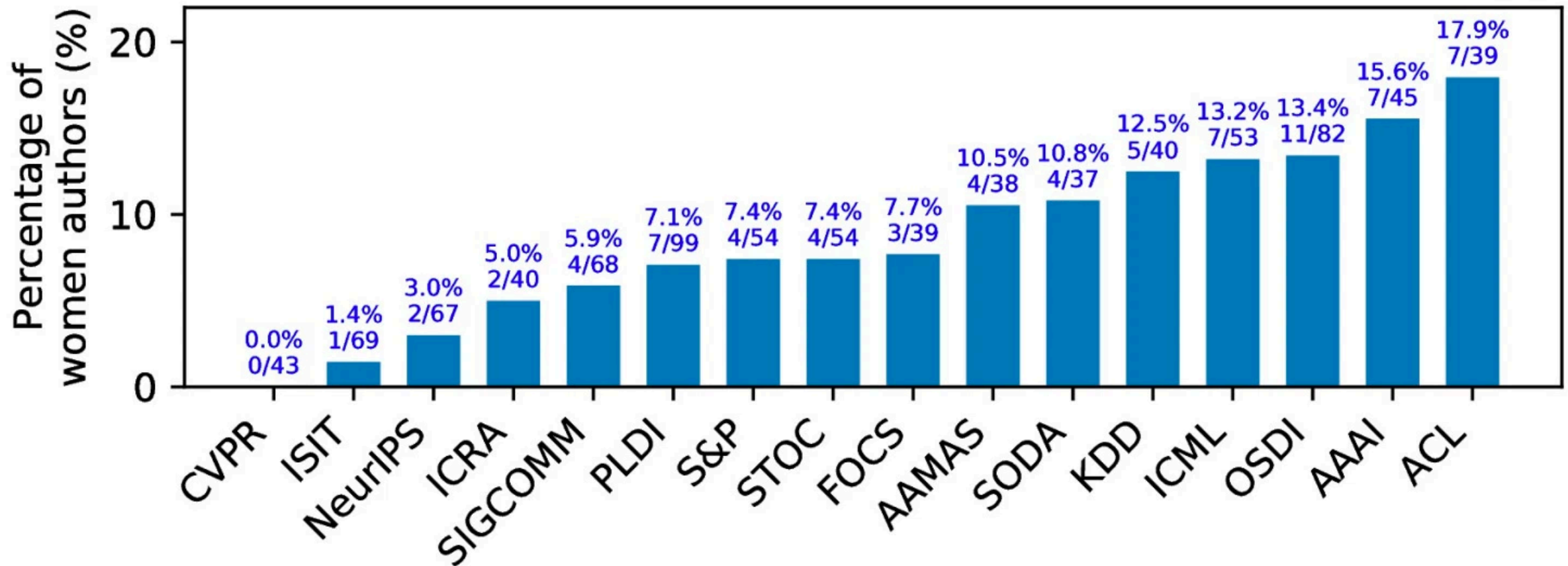
The ordering on this page is randomized (as opposed to ordering alphabetically).  
Read more about [biases due to alphabetical ordering](#).

Ivan Stelmakh | Han Zhao | Charvi Rastogi



# Gender distribution in paper awards

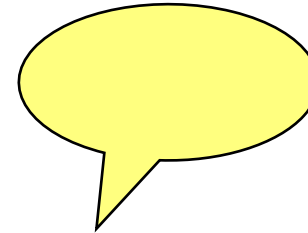
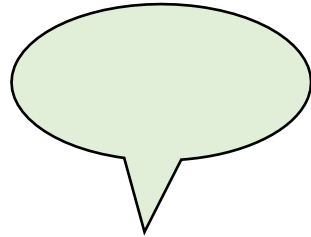
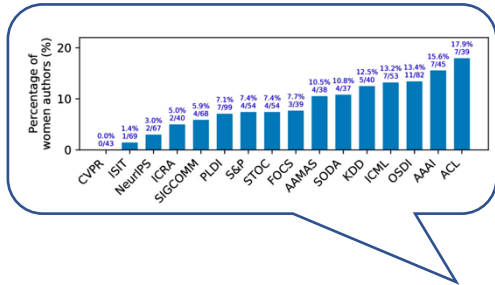
(2010–2018)



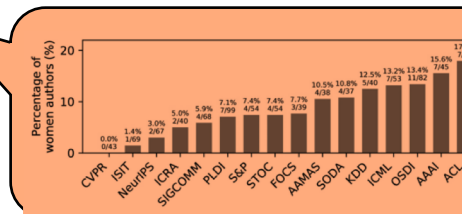
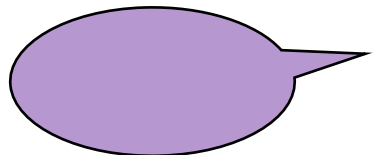


# Need for transparency

- Are author identities visible to the award committee?
- How is the committee determined?
- What criteria are used?



**Started conversations in information theory society,  
NLP community, ML community, vision community,...**





# Conclusions

- **Many sources of biases and unfairness in peer review**
- **Urgent need to revamp peer review, at scale**
  - Lot at stake: Careers, Scientific progress
- **Lots of open problems!**
  - Exciting
  - Theoretical / Applied / Conceptual
  - Challenging
  - **Impactful**





JORGE CHAM © 2014

"Piled Higher and Deeper" by Jorge Cham [WWW.PHDCOMICS.COM](http://WWW.PHDCOMICS.COM)

**Thank you!**