



北京大学前沿计算研究中心
Center on Frontiers of Computing Studies, Peking University



Calibrating “Cheap Signals” in Peer Review without a Prior

Yuxuan Lu, Yuqing Kong

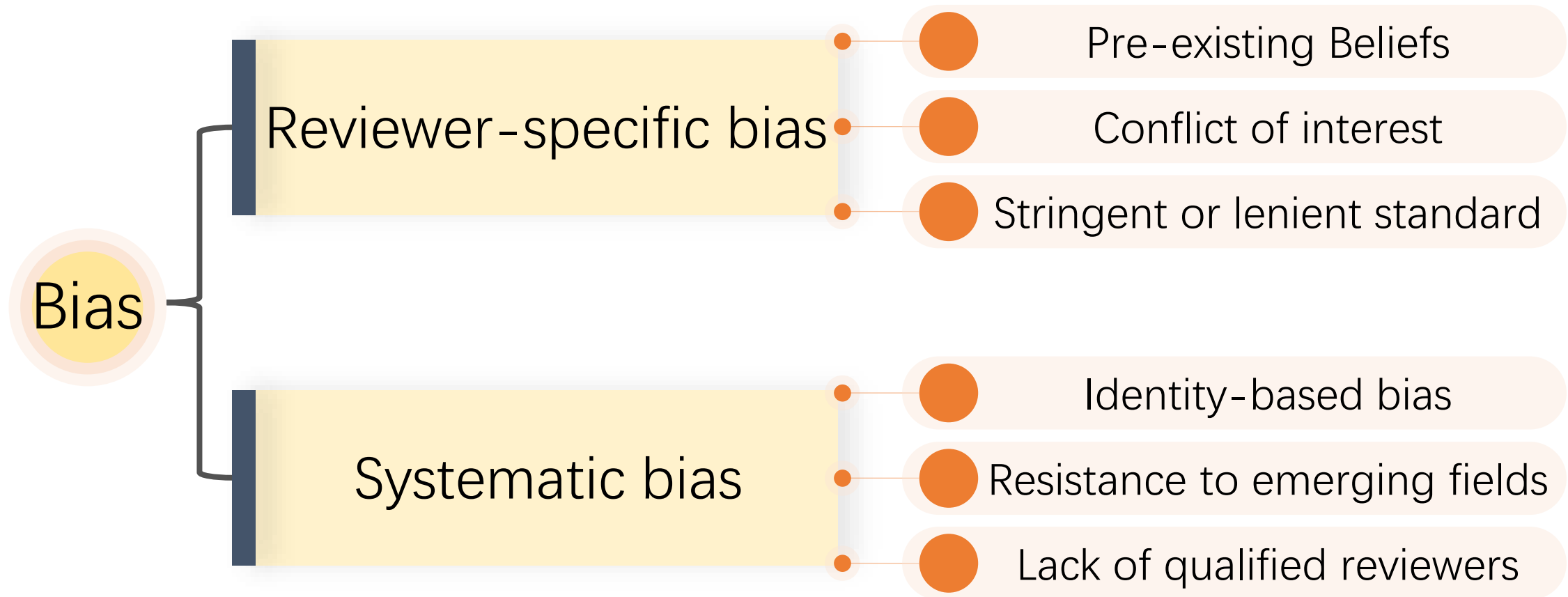
Peking University

NeurIPS 2023



北京大学计算机学院
School of Computer Science

Bias in Peer Review



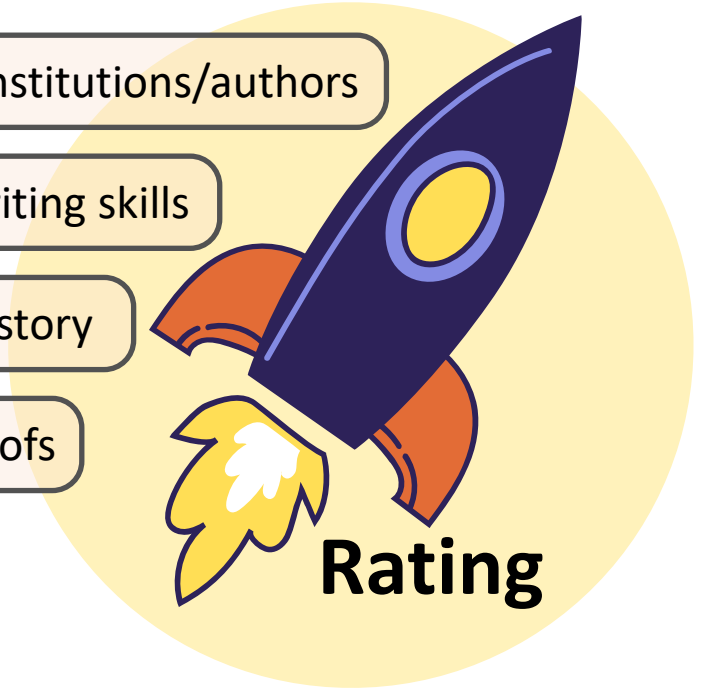
Issue: Cheap Signals

Cheap Signals



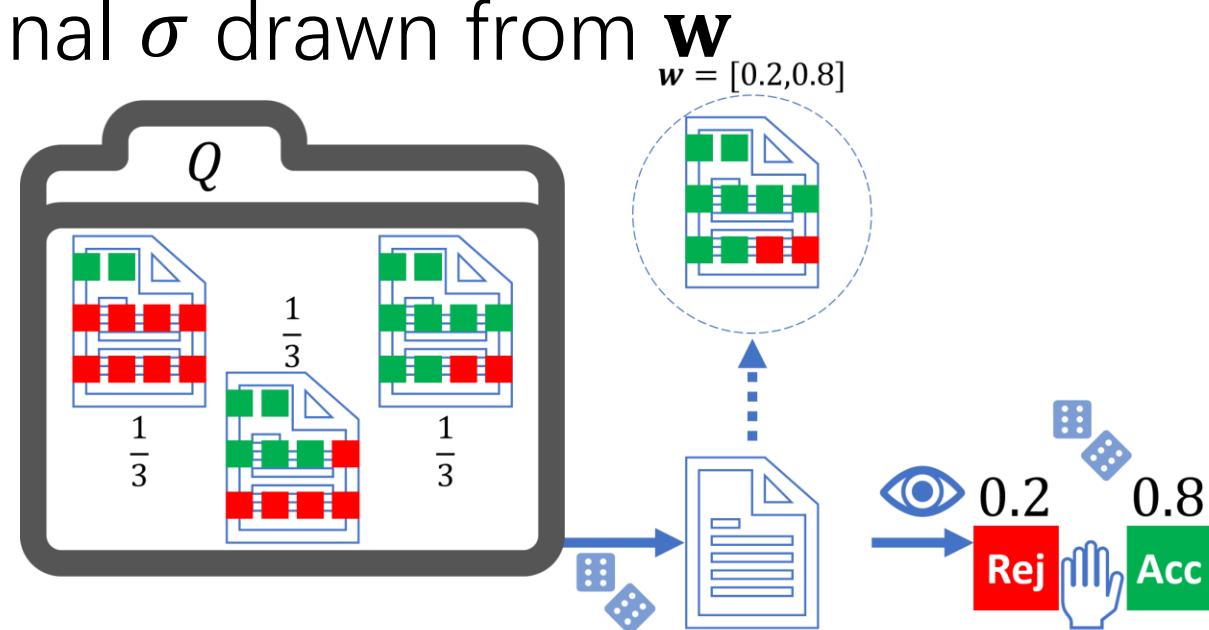
Systematic bias

- ✓ Reputable institutions/authors
- ✓ Excellent writing skills
- ✓ Fascinating story
- ✓ Lengthy proofs



Modelling without Cheap Signal

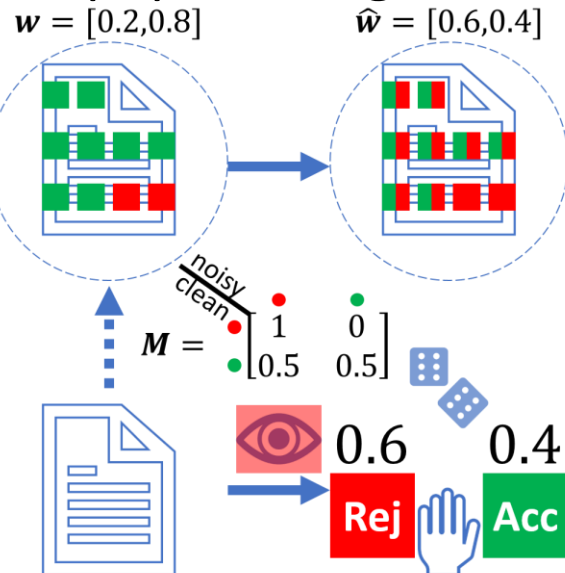
- The set of possible signals $\Sigma = \{0 \text{ (rej)}, 1 \text{ (acc)}\}$
- Paper state $\mathbf{w} \in \left\{ \begin{array}{l} \text{bad } (\mathbf{w} = [.8, .2]) \\ \text{fair } (\mathbf{w} = [.5, .5]) \\ \text{good } (\mathbf{w} = [.2, .8]) \end{array} \right\}$
- Each reviewer receives i.i.d signal σ drawn from \mathbf{w}
- Prior $\mathbf{Q} = \frac{1}{3} \text{ bad}, \frac{1}{3} \text{ fair}, \frac{1}{3} \text{ good}$



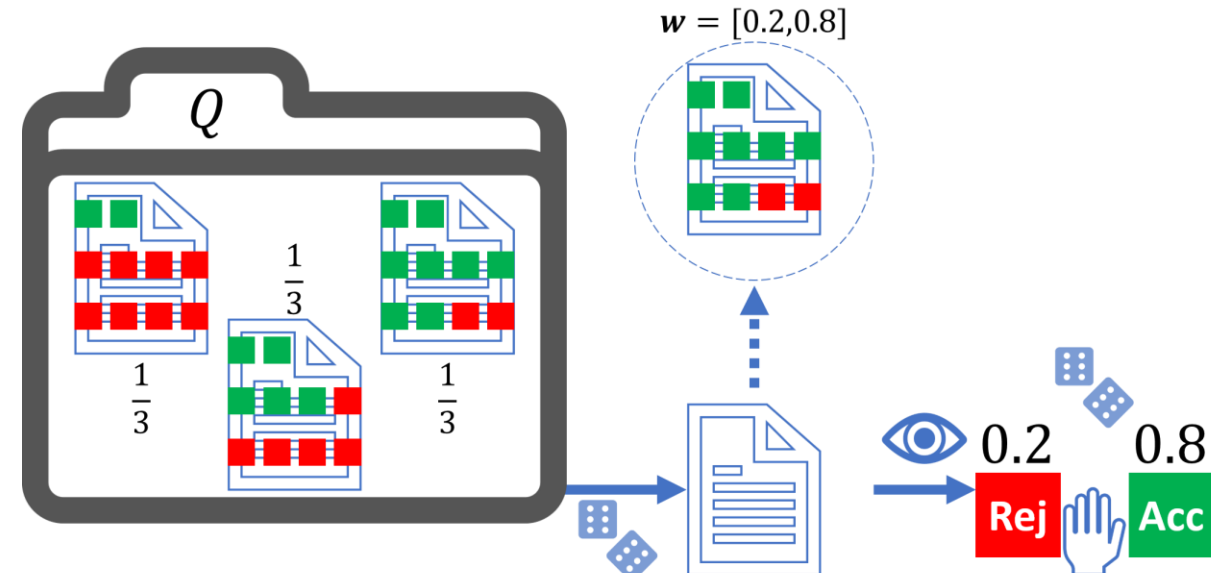
Modelling Cheap Signals

- Regard cheap signals as a bias operator M
 - Bias M alters reviewer's clean signal σ to a biased signal $\hat{\sigma} = M(\sigma)$
- Reviewer only obtains $\hat{\sigma}$ without realizing σ

Good paper + negative bias



Ideal world with no bias



Target: Calibrating Cheap Signals

- We want a process that, in a biased world, rank the quality of papers as if there is no bias.
 - Our method: additionally collecting **second-order information** from reviewers.

$$S\left(\begin{array}{c} \hat{\sigma}_1^A \\ \hat{\sigma}_2^A \hat{\sigma}_3^A \\ \hat{\sigma}_4^A \hat{\sigma}_5^A \end{array} \mathbf{A}, \text{Reviewer} \right) > S\left(\begin{array}{c} \hat{\sigma}_1^B \\ \hat{\sigma}_2^B \hat{\sigma}_3^B \\ \hat{\sigma}_4^B \hat{\sigma}_5^B \end{array} \mathbf{B}, \text{Reviewer} \right)$$

↕

$$\sum_{i=1}^5 \sigma_i^A > \sum_{j=1}^5 \sigma_j^B$$

Key Observations

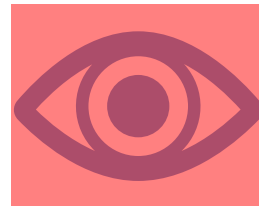
1. Cheap signals affects reviewers' ratings and prior beliefs in the same way

$$Q = \frac{1}{3} \text{bad}, \frac{1}{3} \text{fair}, \frac{1}{3} \text{good}$$

$$\text{bad: } \mathbf{w} = [0.8, 0.2]$$

$$\text{fair: } \mathbf{w} = [0.5, 0.5]$$

$$\text{good: } \mathbf{w} = [0.2, 0.8]$$



$$\hat{Q} = \frac{1}{3} \hat{\text{bad}}, \frac{1}{3} \hat{\text{fair}}, \frac{1}{3} \hat{\text{good}}$$

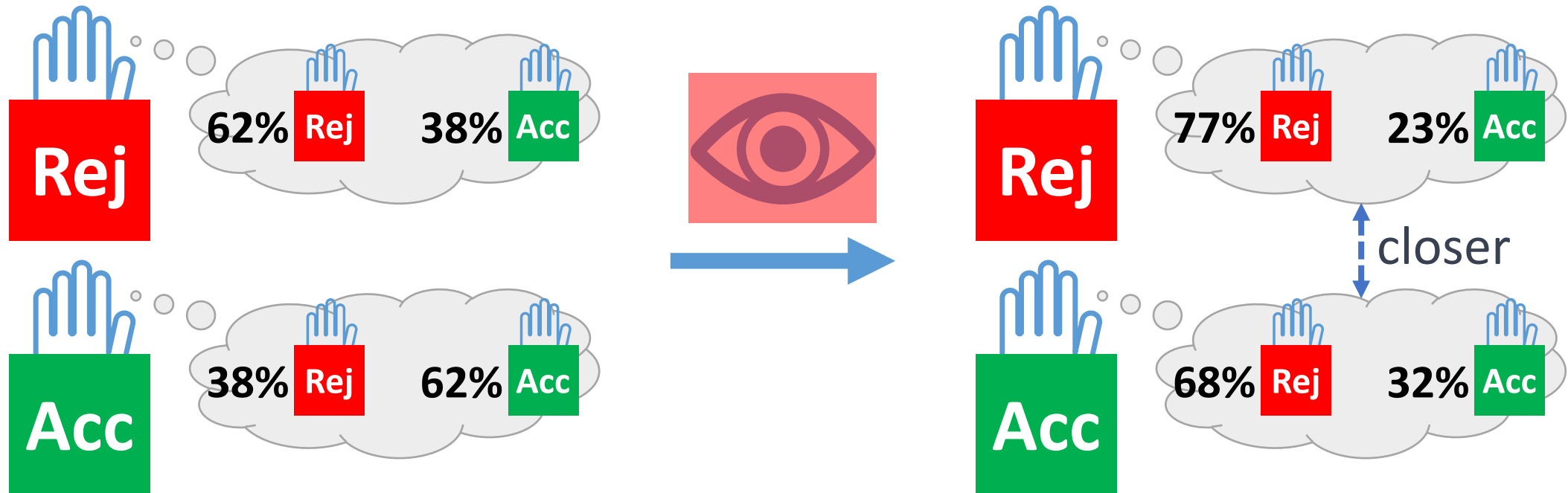
$$\hat{\text{bad}}: \hat{\mathbf{w}} = [0.9, 0.1]$$

$$\hat{\text{fair}}: \hat{\mathbf{w}} = [0.75, 0.25]$$

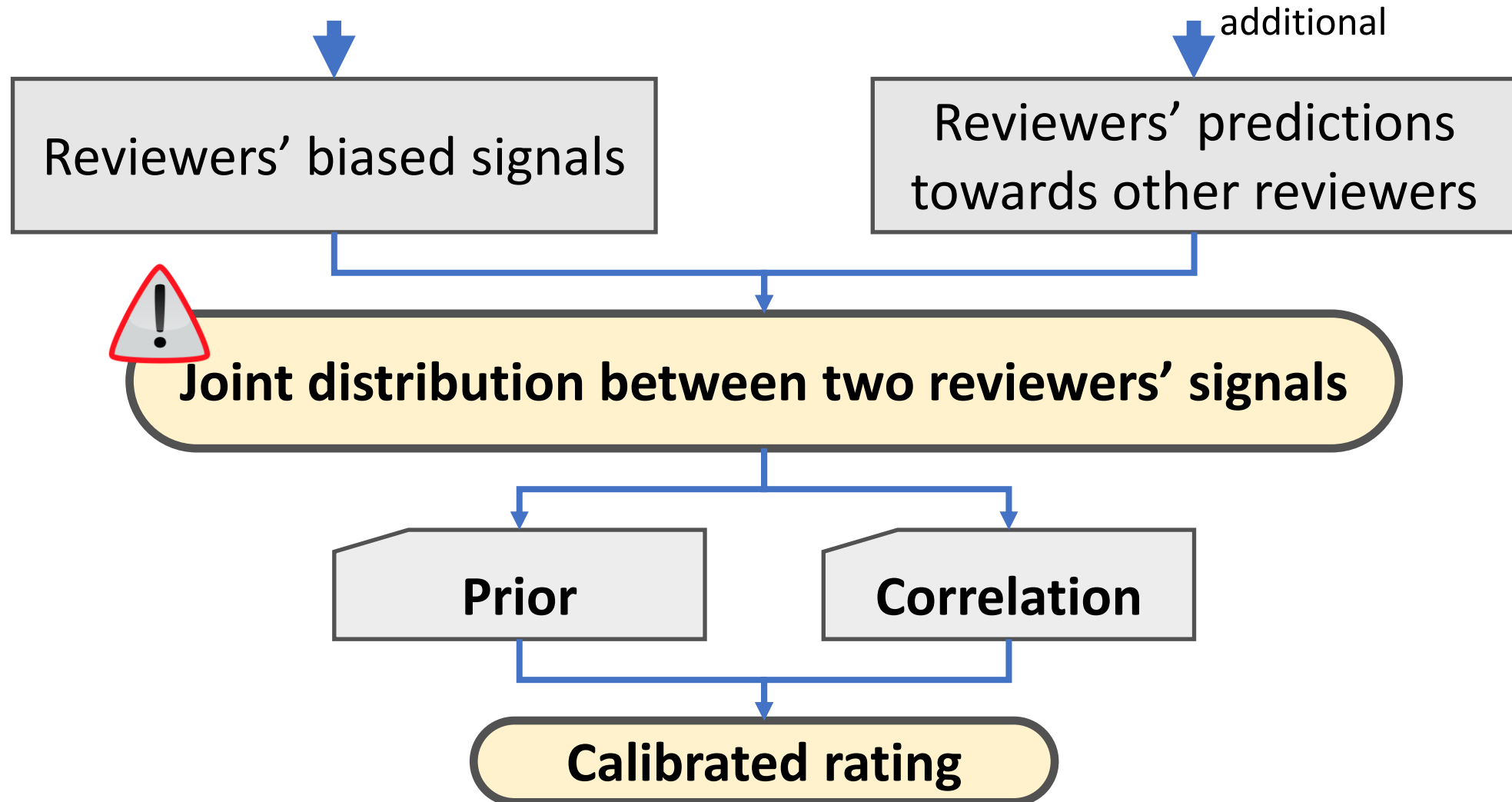
$$\hat{\text{good}}: \hat{\mathbf{w}} = [0.6, 0.4]$$

Key Observations

2. Systematic bias weaken the correlation of reviewers' feedback

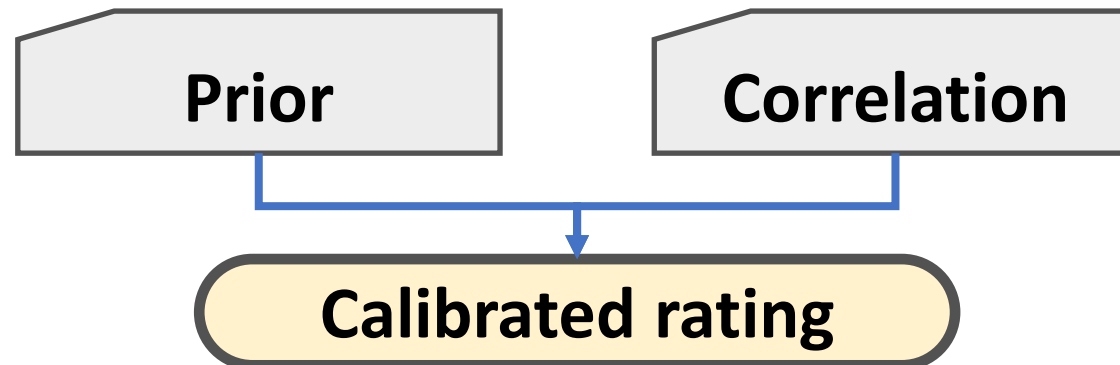


Main Idea: Calibration by Prediction



Main Idea: Calibration by Prediction

Theorem (informal): the calibrated rating is an affine transformation of the true rating in expectation.



Thank you for listening!

Contact: yx_lu@pku.edu.cn

Materials: https://yxlu.me/publication/peer_review_neurips23