

# **No more Reviewer #2: Subverting Automatic Paper-Reviewer Assignment using Adversarial Learning**

---

**Thorsten Eisenhofer, Erwin Quiring, Jonas Möller,  
Doreen Riepel, Thorsten Holz, Konrad Rieck**

RUHR  
UNIVERSITÄT  
BOCHUM

**RUB**



# Papers and Reviews

---

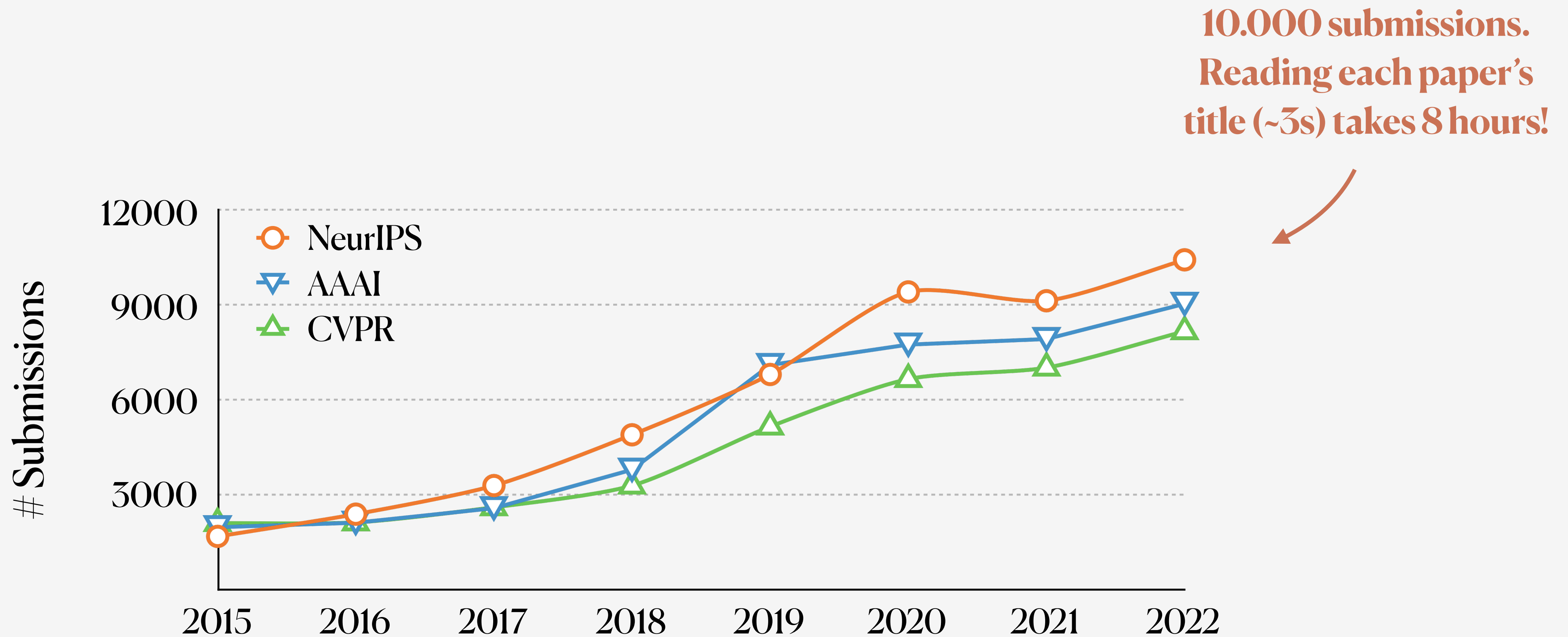
## Peer Review

- Independent evaluation of scientific papers
- Main instrument for quality control

## Initial Step: Paper-Reviewer Assignment

- Assignment of qualified reviewers to each paper
- Good match of topic (paper) and expertise (reviewer)

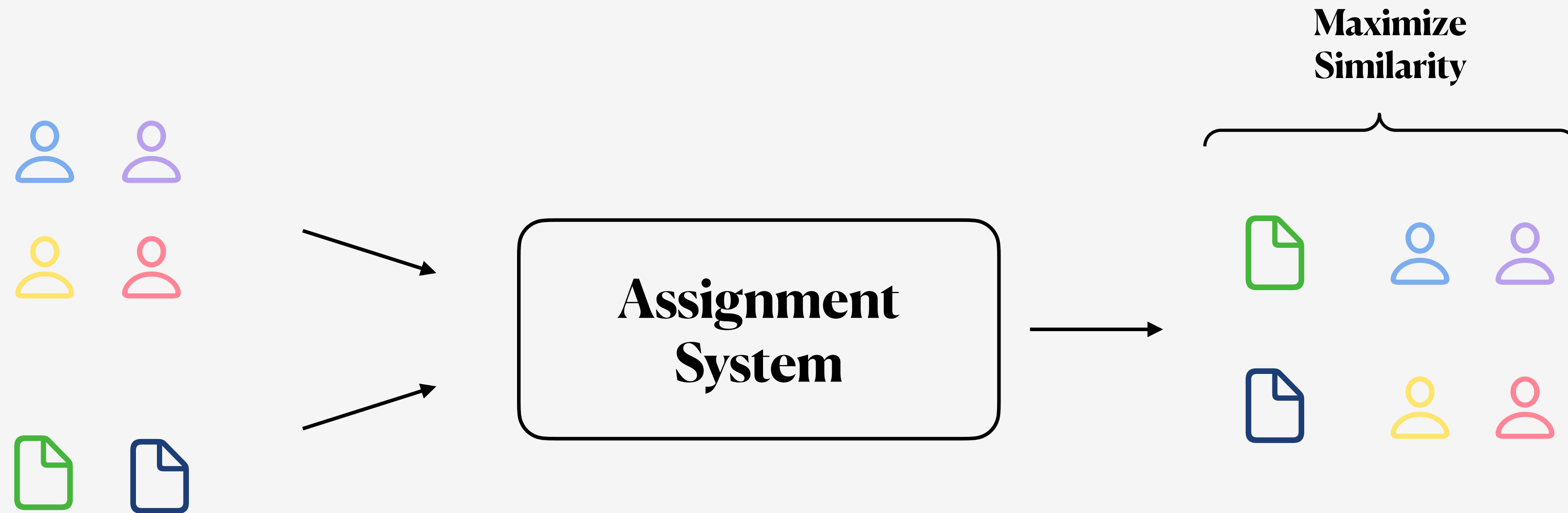
# Assignment Process



**Manual bidding increasingly impossible**

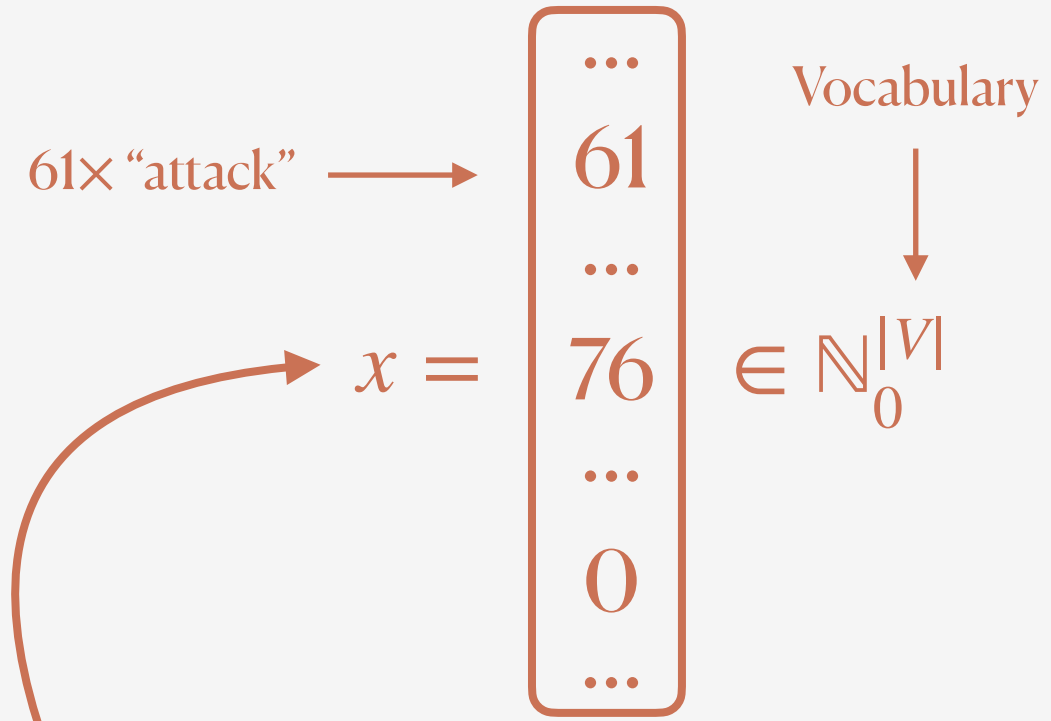
# Automatic Assignment Systems

---



**Use ML to distill submissions and reviewer expertise**

# Topic Modeling

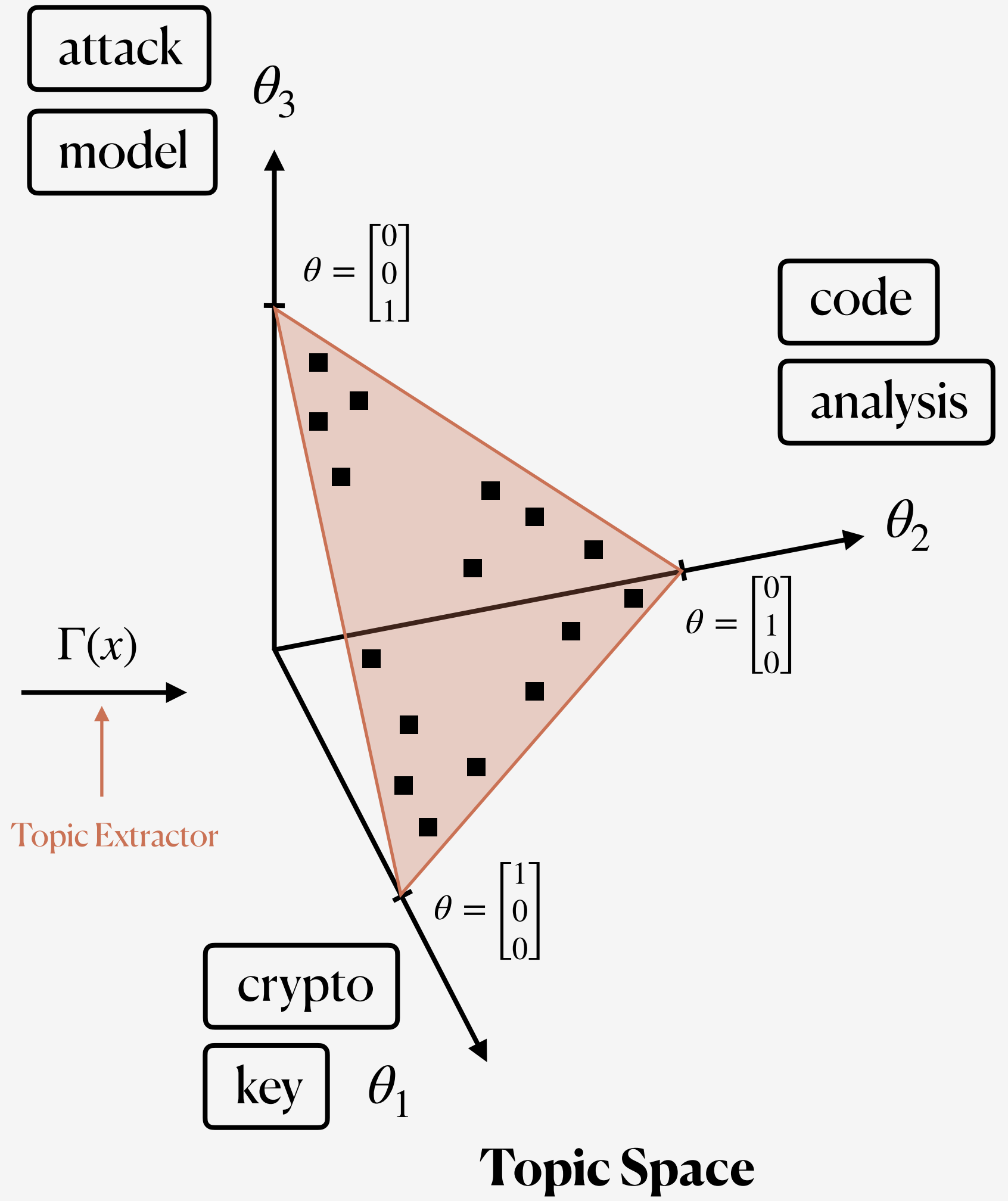
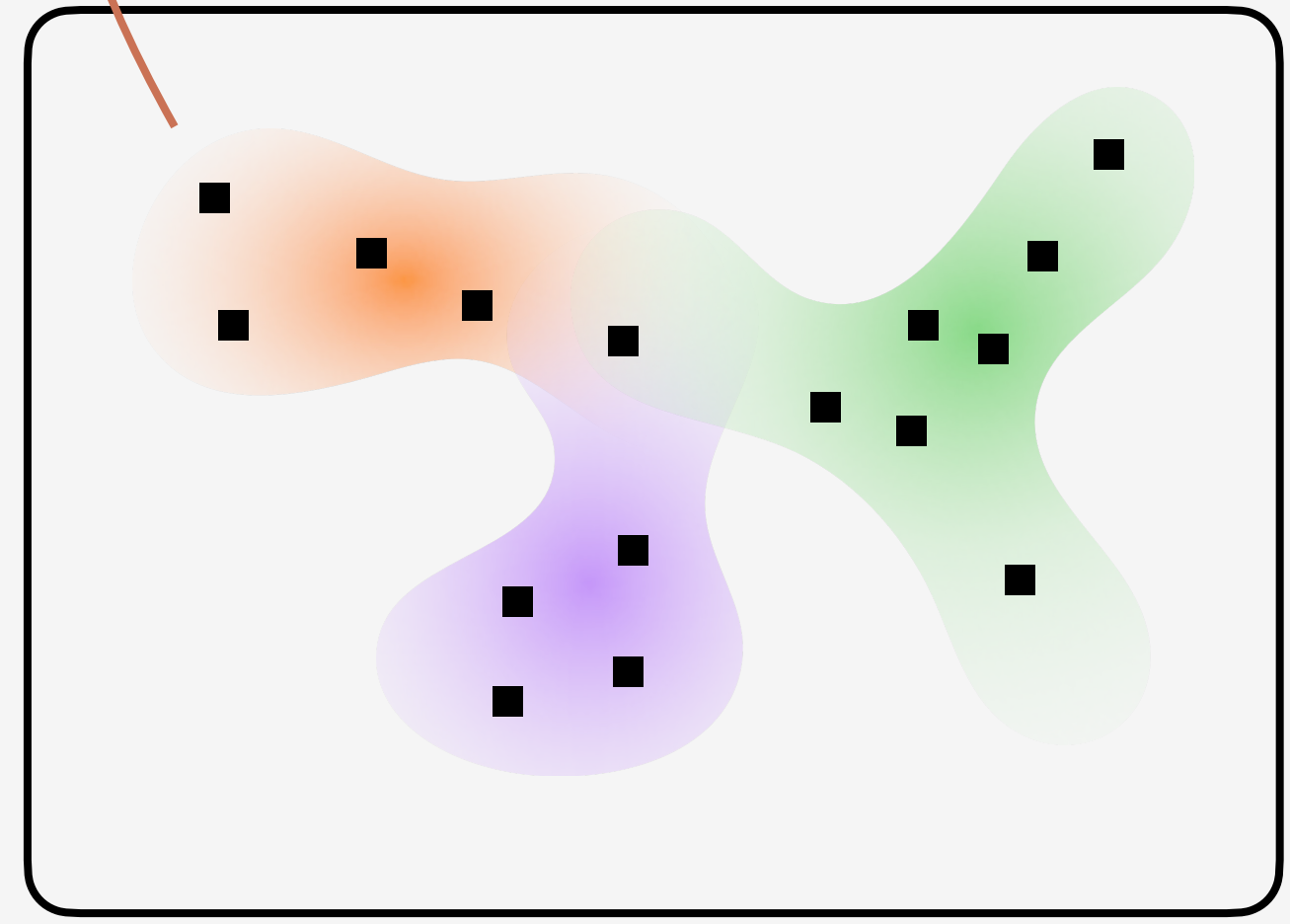
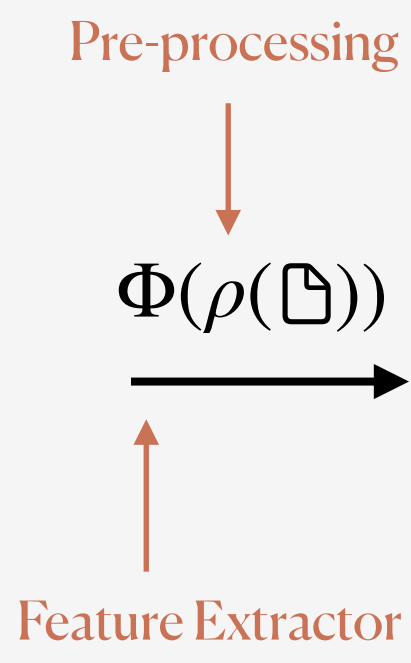


**No more Review #2: Subverting Automatic Paper-Reviewer Assignment using Adversarial Learning**  
 Thomas Eisenhofer\*, Kevin Gohring\*, Jonas Müller†, Dominik Reip†, Thomas Heitz†, Konrad Rieck†

**LogPicker: Strengthening Certificate Transparency Against Covert Adversaries**  
 Alexander Diskov†, David Klein, Robert Michael, Tamas Szabo, Konrad Rieck and Martin Johns

**Evaluating Explanation Models for Deep Learning in Security**  
 Alexander Wernicke†, Daniel Aep†, Christian Weiss†

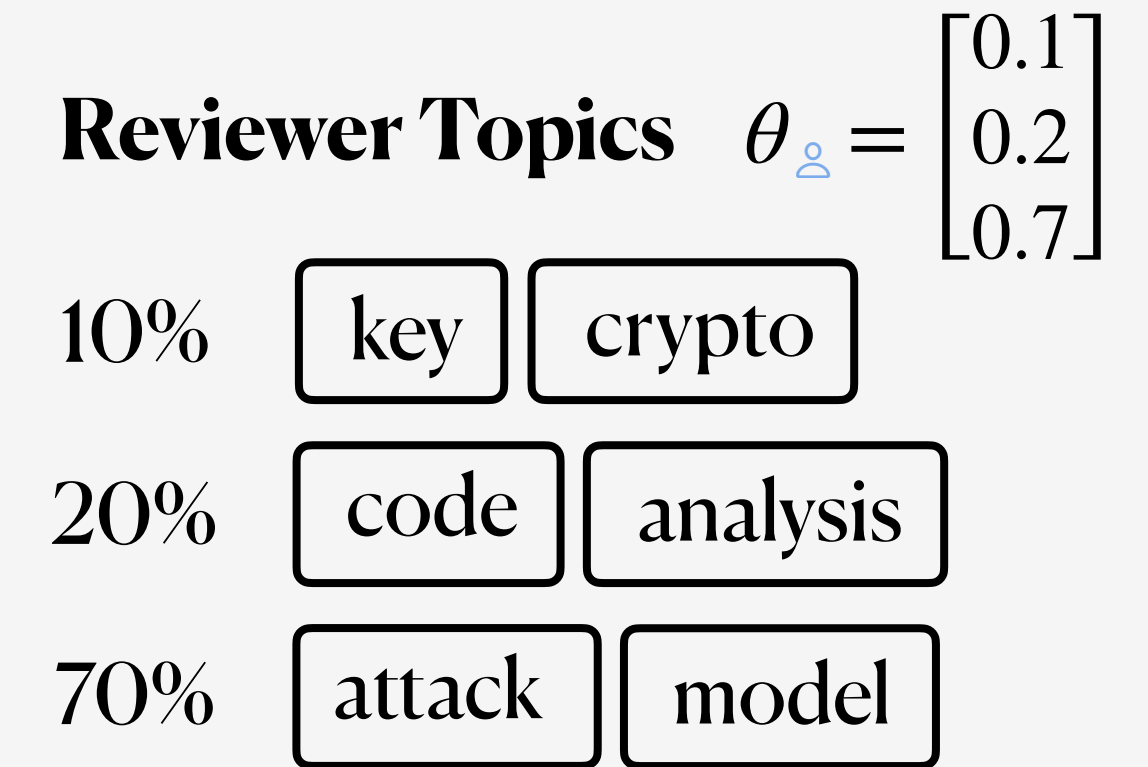
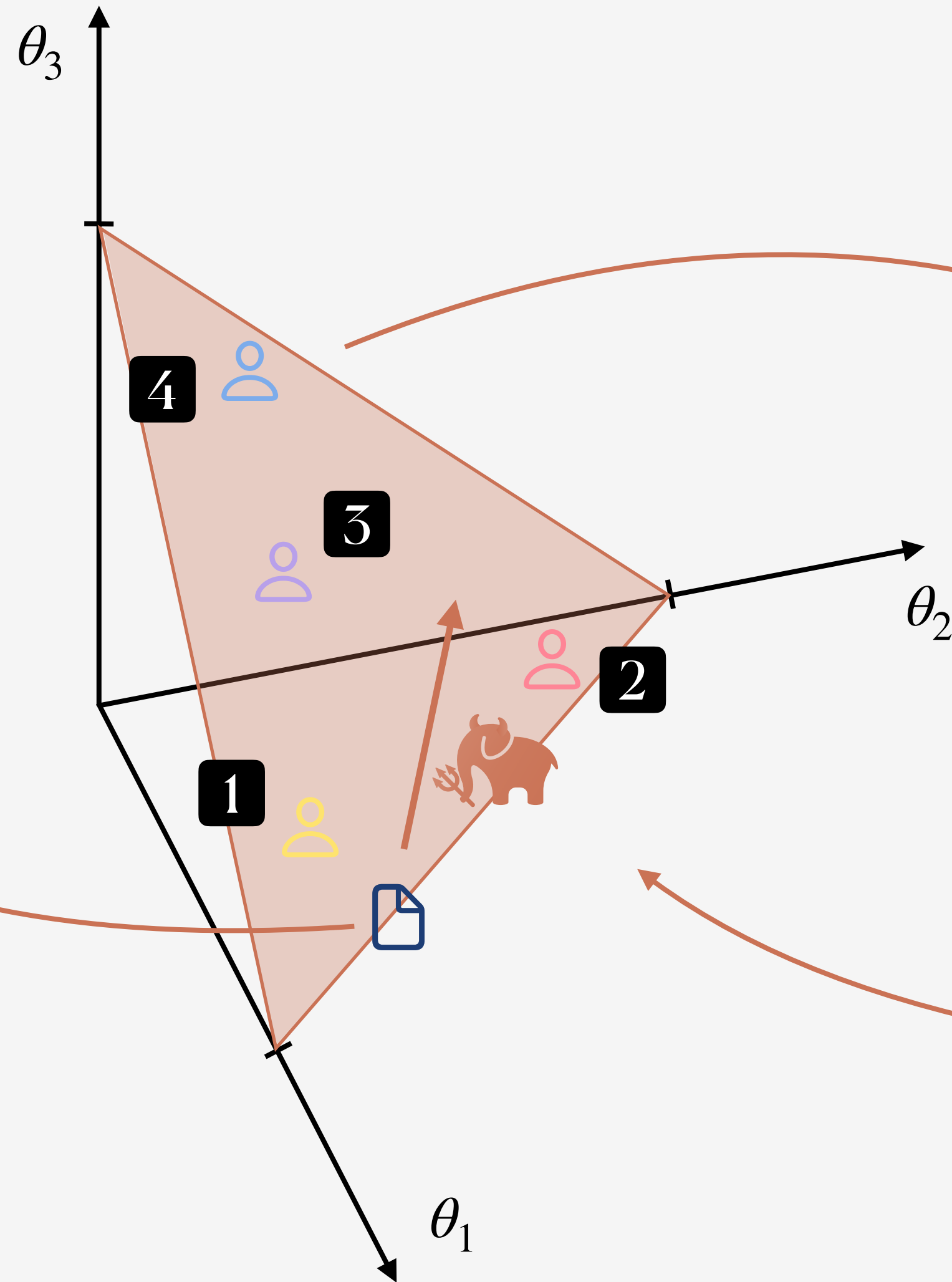
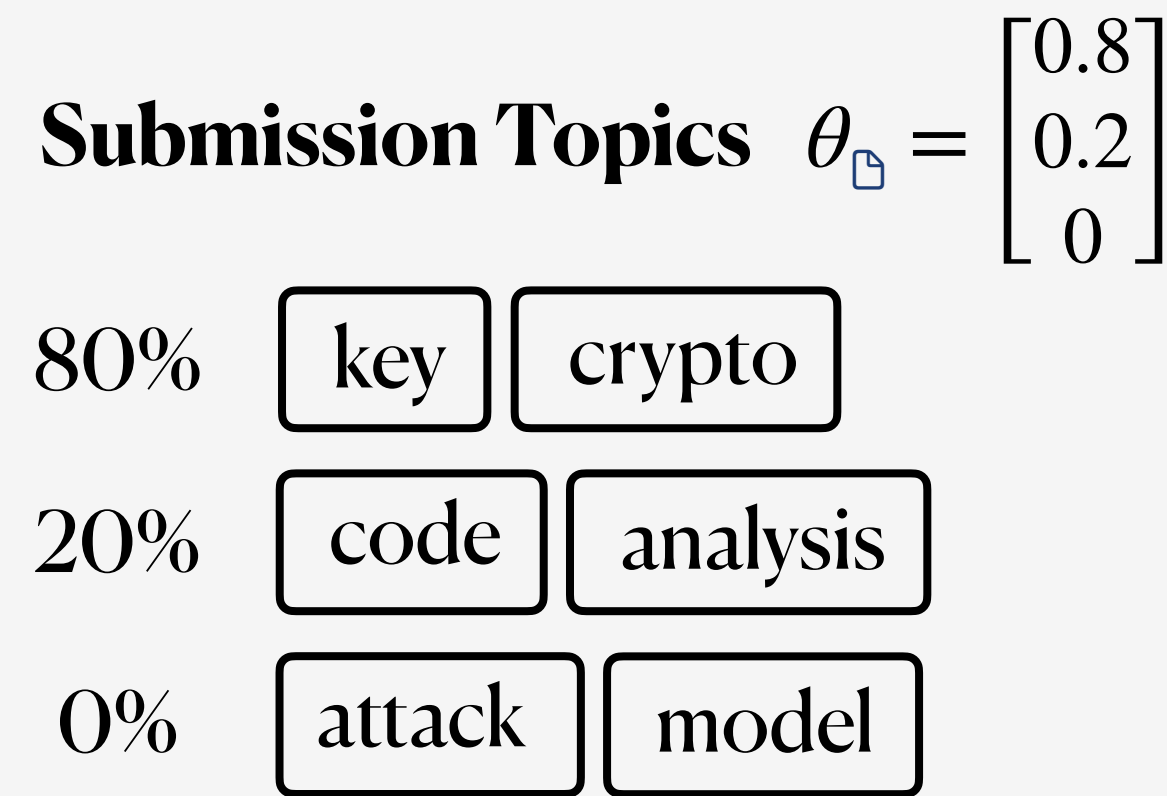
**Lessons: Establishing Fast, Bidirectional Communication into Air-Capped Systems**  
 Niklas Kälsch†, Stefan Pfeiffer†, Maximilian Neppel†, Thomas Schneider†, Konrad Rieck†, Christian Wessnag†



Corpus  $\mathcal{D} = \{ \text{📄}, \text{📄}, \dots, \text{📄} \}$

Feature Space

# Topic Modeling



**Need to project changes back into the problem space!**

**Goal: Manipulate submission  to pick our own reviewers**

# Problem-space

---

## Problem-space transformations to add/remove words from input file

### Format-/ and encoding-level

Hidden Box

u+0061 u+0430

Homoglyphs

← a ≠ a

### Text-level

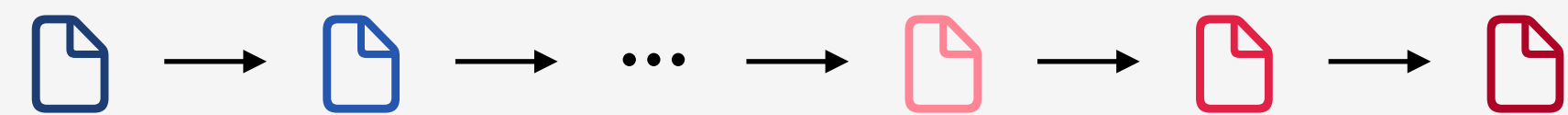
Reference addition

Synonyms

Language models

Spelling mistakes

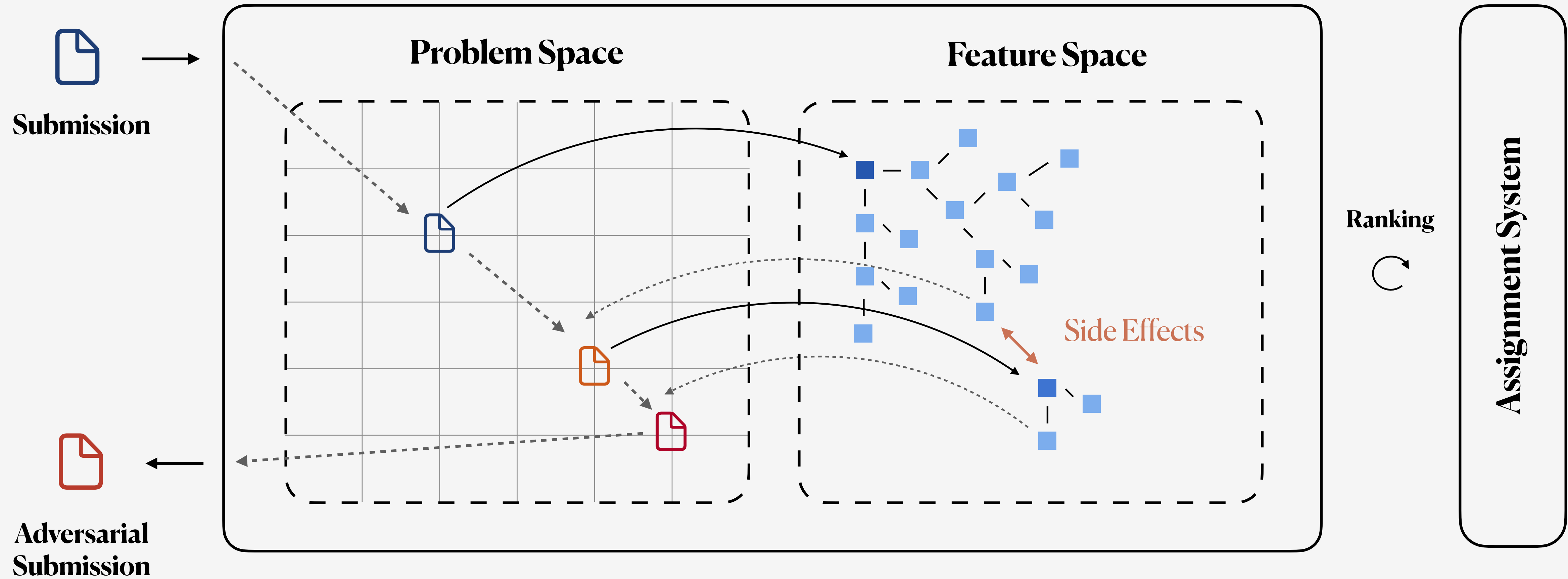
### Chain several transformations



### Constraints

 is plausible and semantic correct

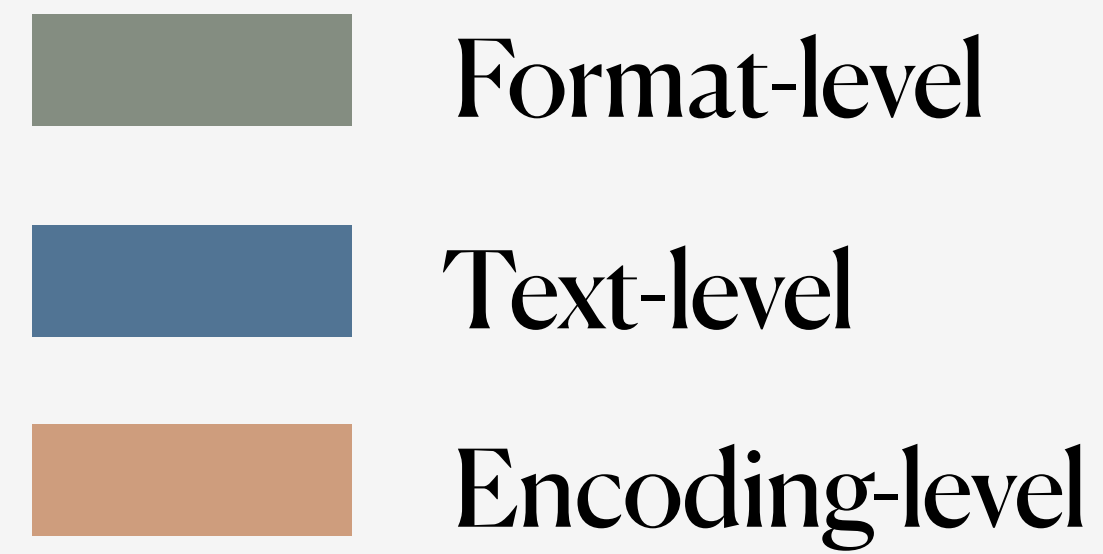
# Hybrid Search Strategy



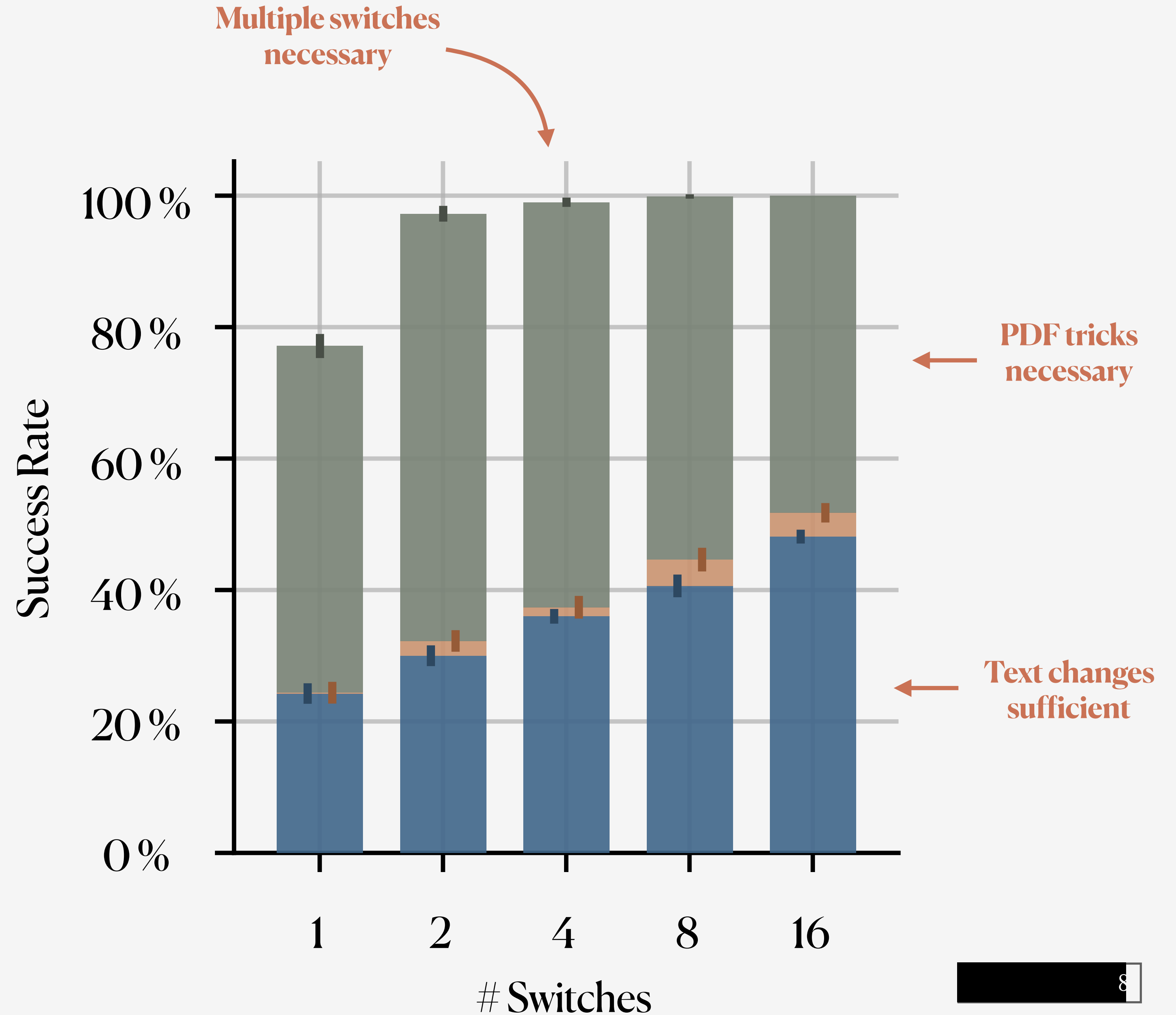


# Evaluation

## Simulation of IEEE S&P' 20



**More results  
in the paper!**



# Take Aways

---

## **New attack against automatic reviewer-paper assignment**

- Hybrid attack strategy in feature space and problem space
- Minimal and unobtrusive transformations of papers

## **Broader perspective**

- Decisions based on learning models inherently insecure
- More to explore off the beaten path of adversarial learning

**More at [github.com/rub-syssec/adversarial-papers](https://github.com/rub-syssec/adversarial-papers)**



# Thank you!