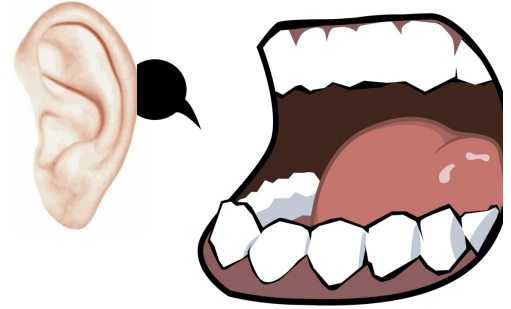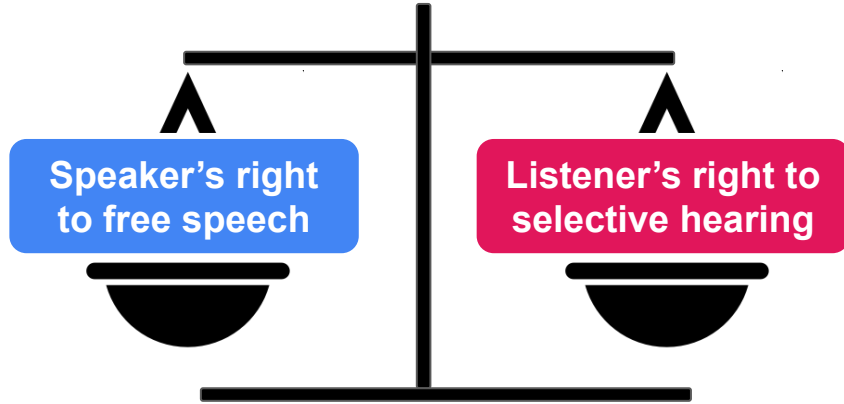# How to stop fake news without central authority or censorship?
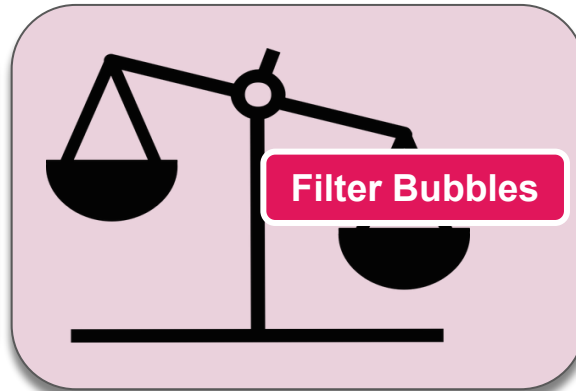
Swapneel Mehta
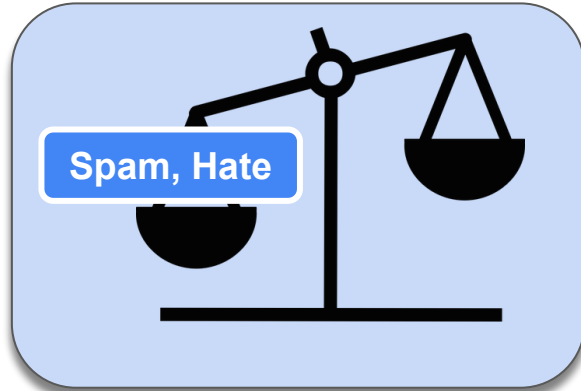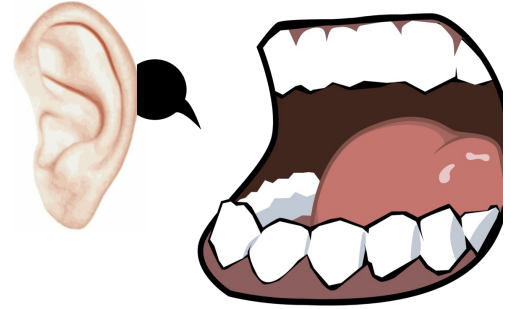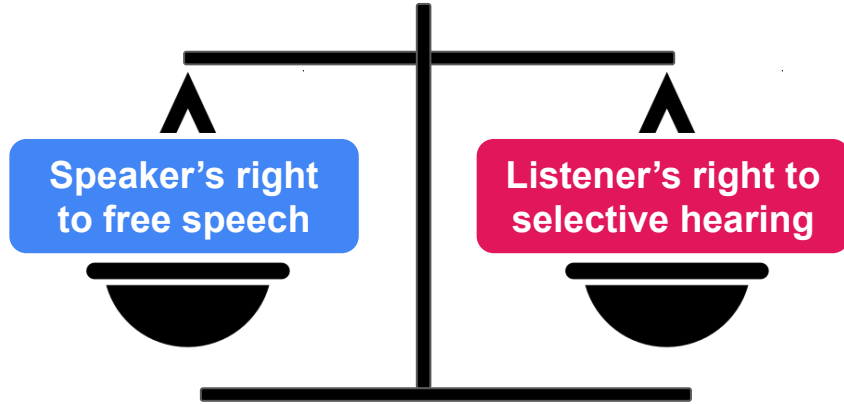Postdoctoral Associate

Joint work with Aaron Nichols, David Rand, Gordon Pennycook, Nina Mazar, Marshall Van Alstyne

# Platforms need to balance the rights of speakers *and* listeners

**Speaker's right to free speech**

**Listener's right to selective hearing**

# Platforms need to balance the rights of speakers _and_ listeners



**Speaker's right to free speech**

**Listener's right to selective hearing**

**Spam, Hate**

**Filter Bubbles**

If either one of the two is given preferential treatment, the other suffers!

_How to find an equilibrium?_

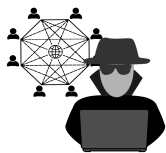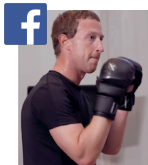# Fake News "Interventions" Don't Work Too Well These Days



**Arms Race**

Bot Detection
Spam
Phishing
**Adversarial Practices**

# Fake News "Interventions" Don't Work Too Well These Days

## Arms Race

Bot Detection
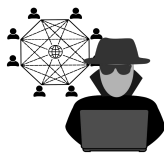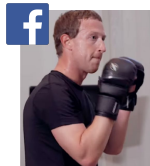Spam
Phishing
**Adversarial Practices**



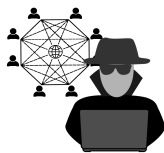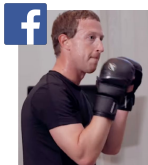## Discrediting Raters

Flagging
Fact-checking
Debunks
**Central Authority**

# Fake News "Interventions" Don't Work Too Well These Days

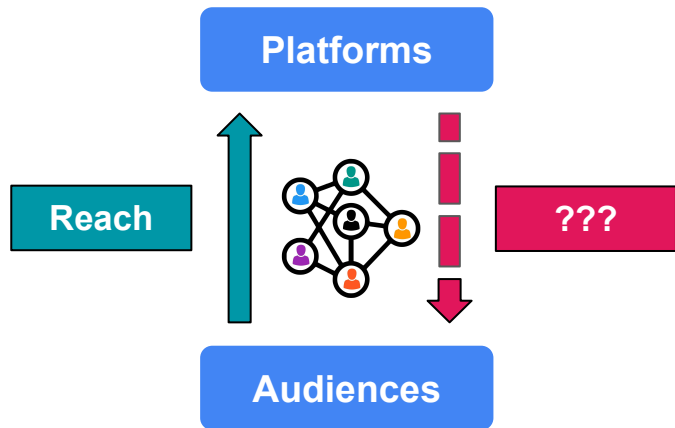| Arms Race | Discrediting Raters | Misplaced Responsibility |
|---|---|---|
| Bot Detection<br>Spam<br>Phishing<br>**Adversarial Practices** | Flagging<br>Fact-checking<br>Debunks<br>**Central Authority** | Auditing Algorithms<br>Deamplification<br>Prebunks<br>**Burden on Platforms, Listeners** |



*Hate Speech's Rise on Twitter Is Unprecedented, Researchers Find*

What if we introduced <u>consequences</u> to producing fake news?

# Platforms should support Free Speech, but not *Free Reach*



Audiences provide platforms with *reach*

Platforms should be accountable to protect audiences, at least global regulators think so.
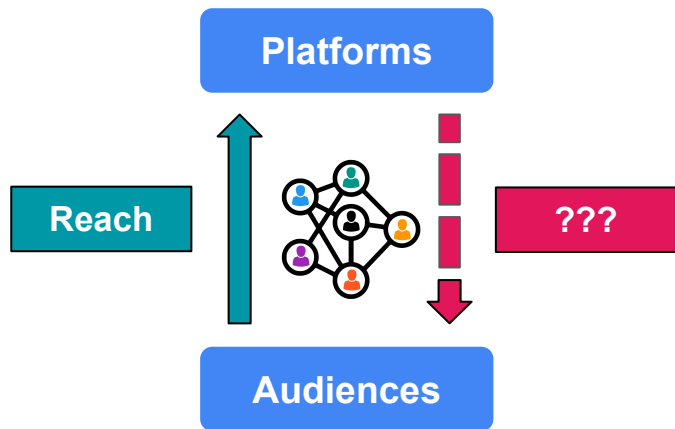
# Platforms should support Free Speech, but not *Free Reach*
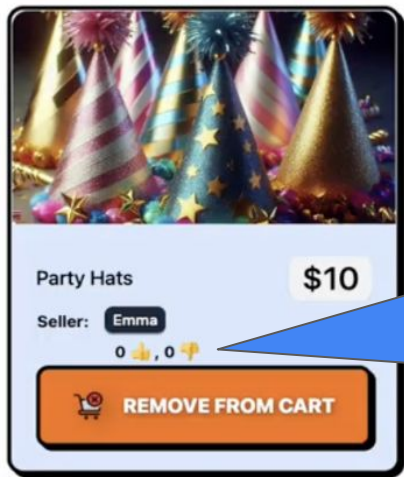


Audiences provide platforms with *reach*

Platforms should be accountable to protect audiences, at least global regulators think so.

Free speech is absolutely vital + legal right.

But we need mechanisms to introduce *accountability* in exchange for *reach*!
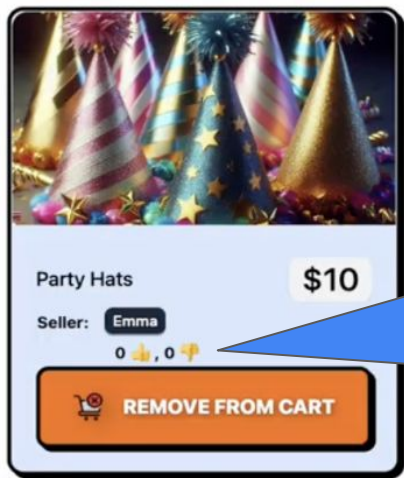
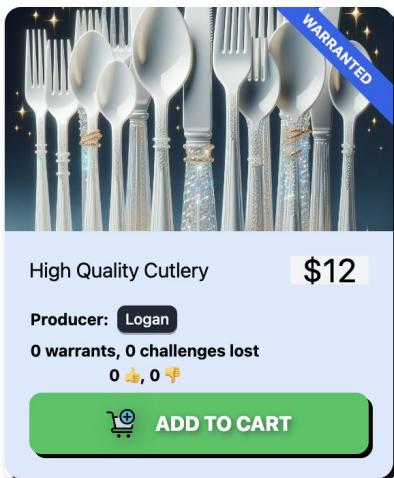# Improving Accountability in Ads Marketplace via *Warrants*



Current markets are Reputation-based (👍/👎)

# Improving Accountability in Ads Marketplace via *Warrants*
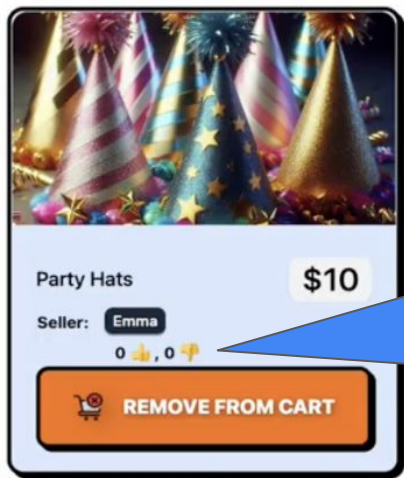


Current markets only rely on reputation

We think "Warrant" Label increases credibility
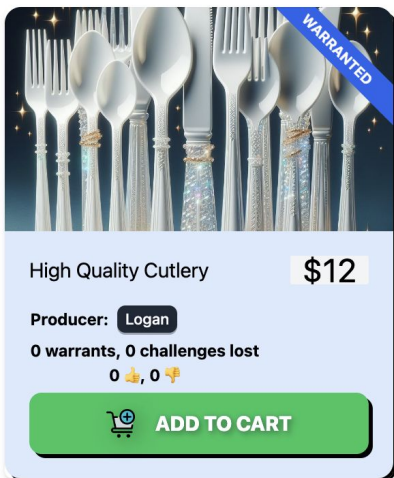
Current markets are Reputation-based (👍/👎)

Warrant = seller **escrows** extra money to **back** claims. If community agrees claim false, money lost!

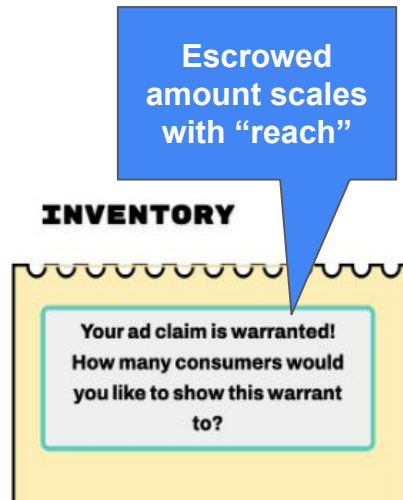# Improving Accountability in Ads Marketplace via *Warrants*



Current markets are Reputation-based (👍/👎)

Warrant = seller **escrows** extra money to **back** claims. If community agrees claim false, money lost!
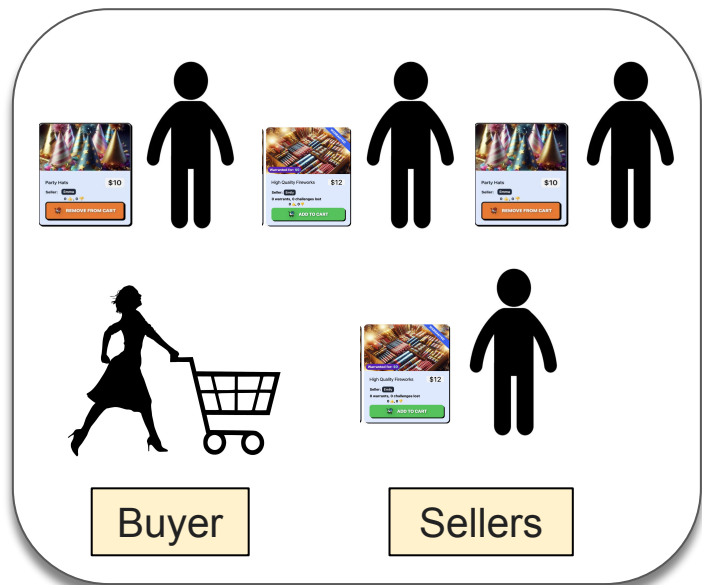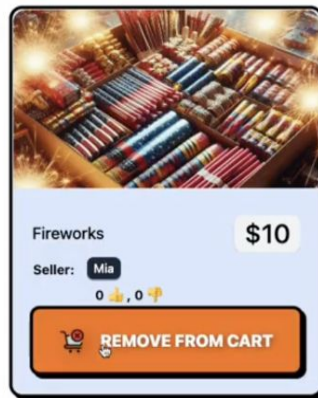
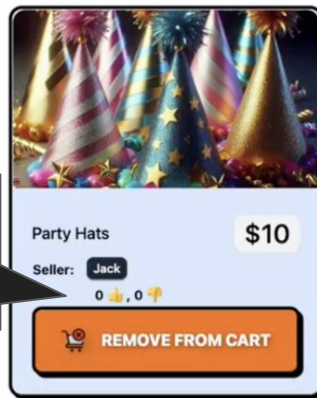Warrants are shown to a predetermined no. of people.

# Can Consumers Stop Falling for Fake Ads through *Warrants*?

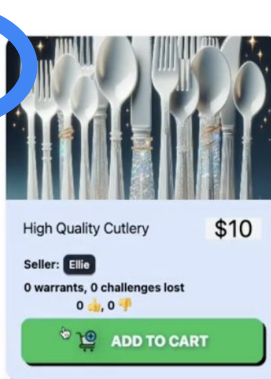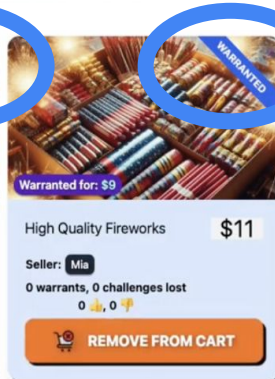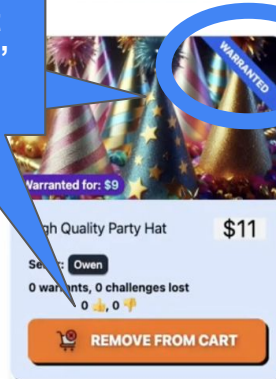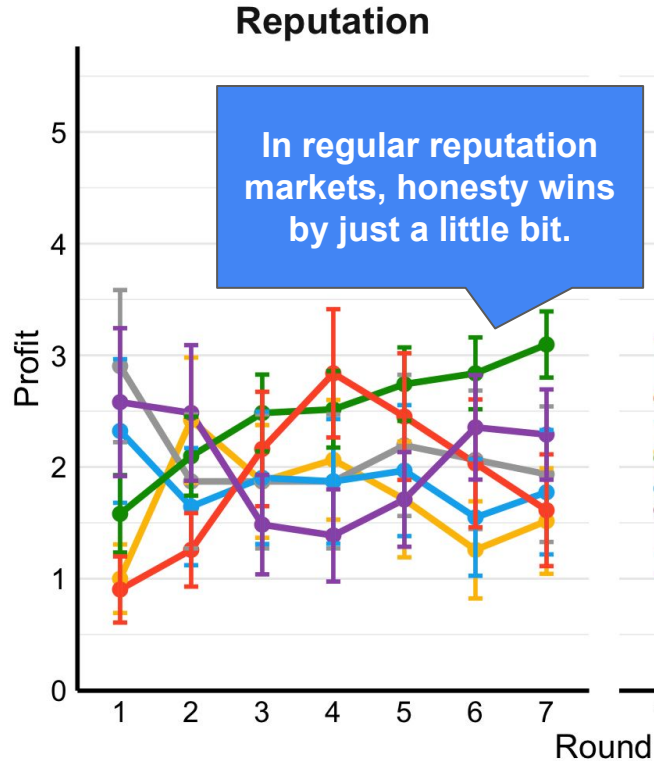# Can Consumers Stop Falling for Fake Ads through *Warrants*?

# Cheating Sellers *Don't Prosper* in Warrants Market

# Cheating Sellers *Don't Prosper* in Warrants Market

# Summary

1. We can limit fake ads without censorship or central authority
2. Platforms benefit, buyers benefit, sellers benefit!
3. Testing it out on a platform like Bluesky or Reddit with <u>custom feeds</u>

**Extending this to Political Markets**

1. Legal Frameworks ask: What constitutes a claim?
2. Social Psychology asks: What about when there is the absence of "ground truth"?

# Bonus: LLMs *deceive* really well!



Volume of Sales by Round, Bot Type, and Condition
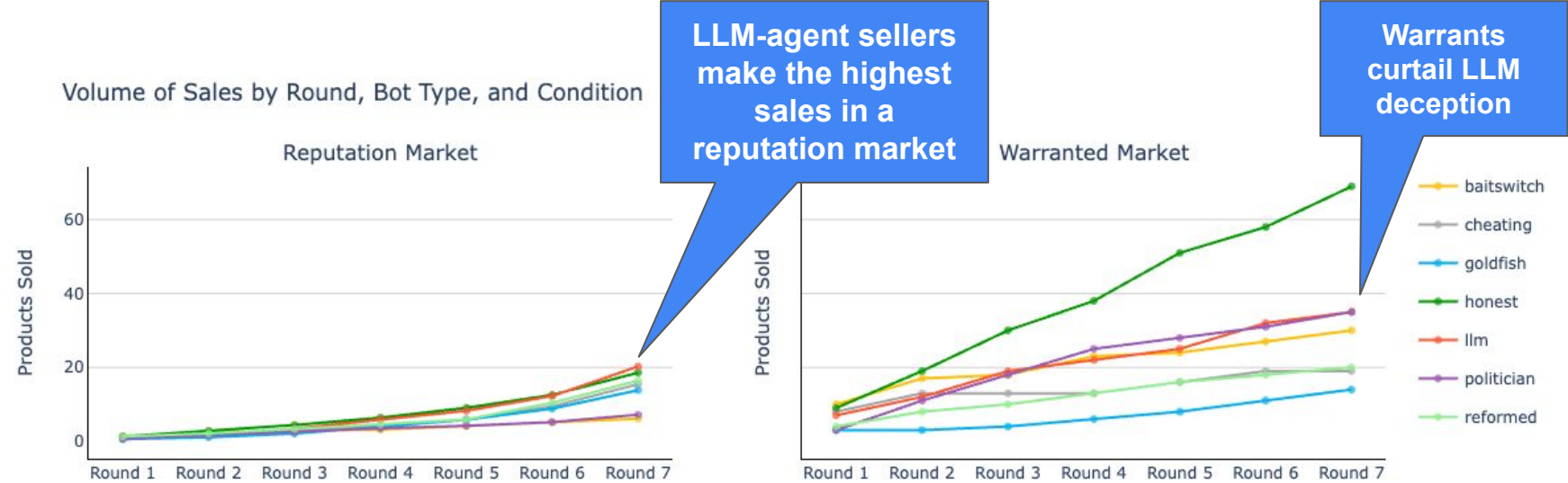
**Reputation Market**

**LLM-agent sellers make the highest sales in a reputation market**

**Warranted Market**

**Warrants curtail LLM deception**

Legend: baitswitch, cheating, goldfish, honest, llm, politician, reformed

**Amazon introduces Amelia, an AI assistant for third-party sellers**

PUBLISHED THU, SEP 19 2024•9:01 AM EDT

# LLMs provide meaningful explanations for their sales strategies



Frequency of Reasoning Strategies in Warrants Market: Insights from 25 Users