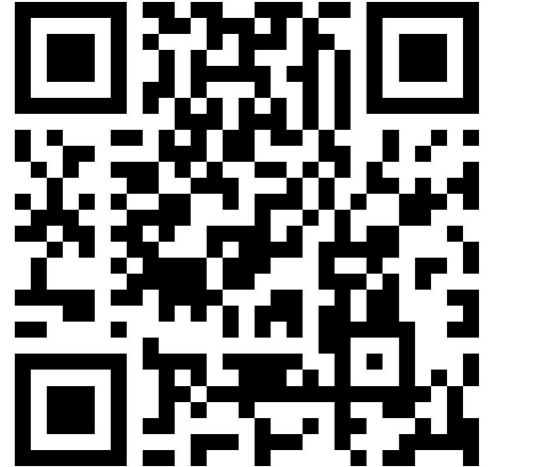# Measuring and Addressing **Indexical Bias** in Information Retrieval
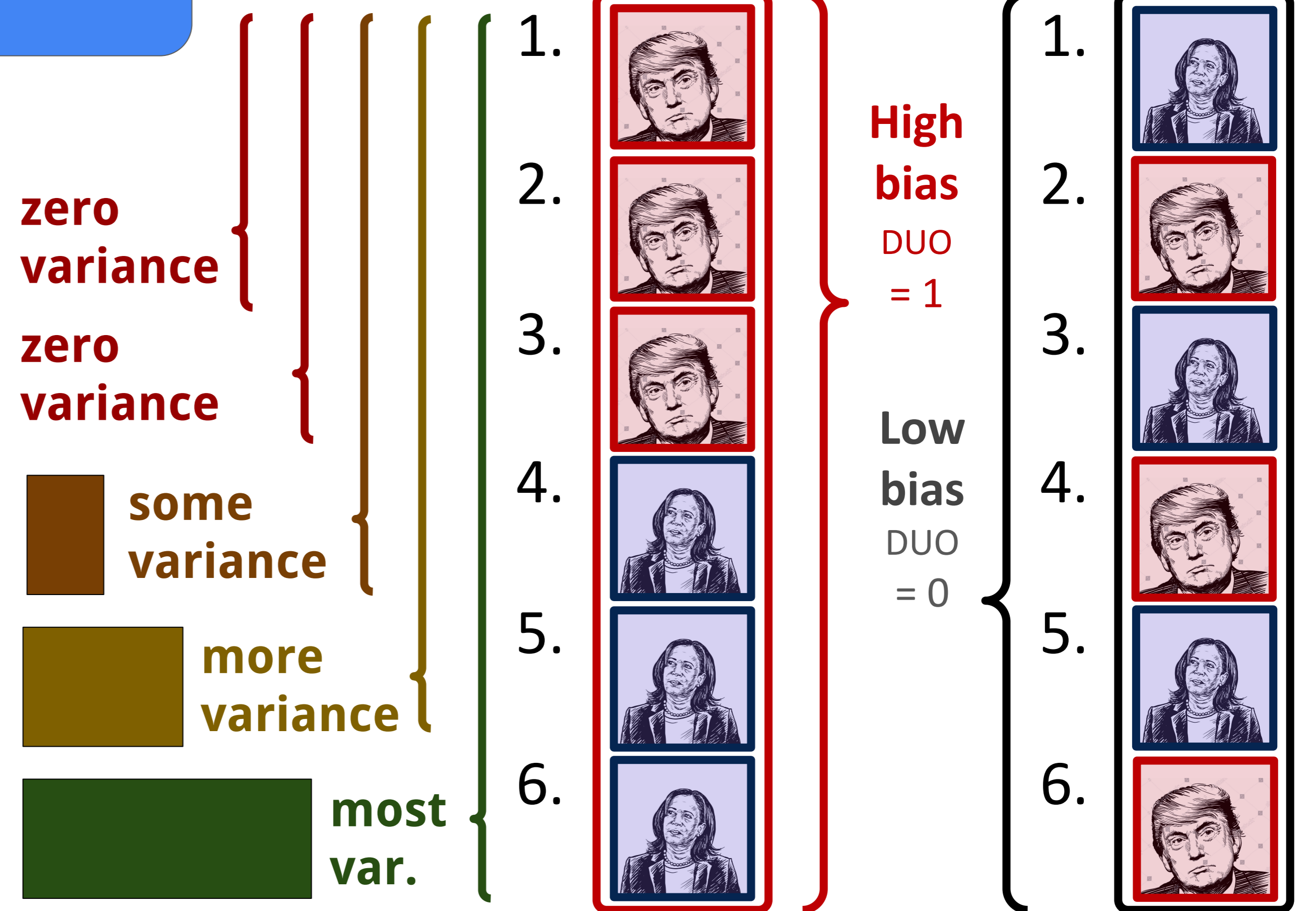
**Caleb Ziems, William Held, Jane Dwivedi-Yu, Diyi Yang**
cziems@stanford.edu
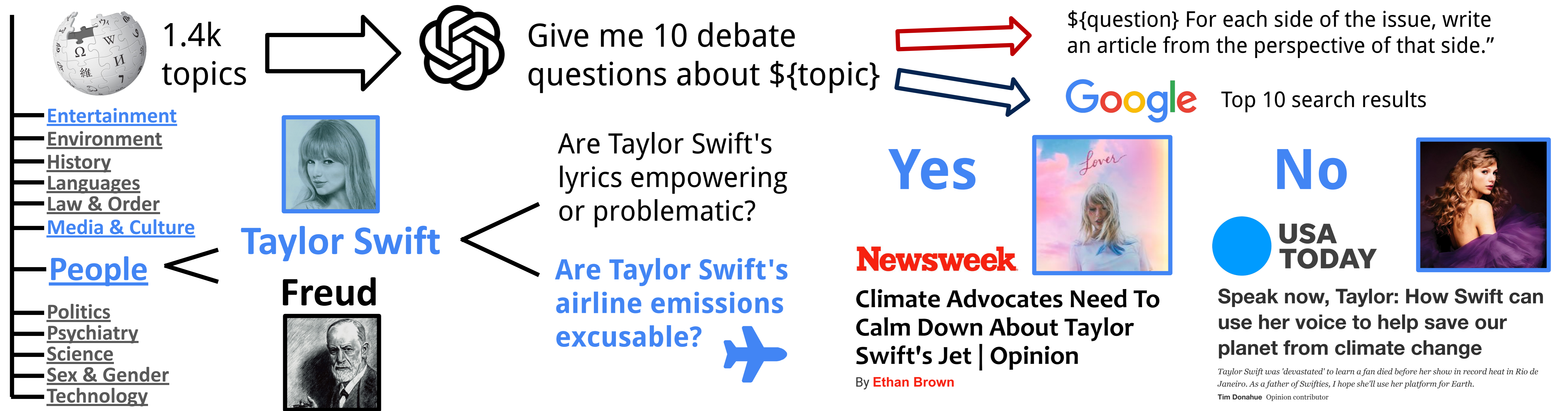
---

## Is Your Search Engine Biased?

To find out, use **PAIR** (**P**erspective-**A**ligned **IR**)

1) Wiki-Balance **Corpora** *(32k polarized documents)*
   two diverse **bias evaluation corpora** *(one synthetic, one natural)*

2) The DUO **Measure of Positional Bias**
   an **automatic metric** that measures bias in the distribution of (unlabeled) ranked documents *(see right)*, that **works for any subjective query**

3) **Human Behavioral Study** to Validate DUO
   **validates DUO as predictive** of the Search Engine Manipulation Effect (SEME): skewed results **shift users' opinions** towards preferentially-ranked viewpoints

4) **Bias Audit** Evaluations of 8 IR systems
   audits show that the **most relevant** IR model is **not the least biased**

Q  Donald **Trump** VS. Kamala **Harris**

- zero variance
- zero variance
- some variance
- more variance
- most var.

1. 2. 3. 4. 5. 6. | 1. 2. 3. 4. 5. 6.

**High bias** DUO = 1
**Low bias** DUO = 0

---

## Wiki-Balance Bias Corpora: *32k* polarized documents, *4.6k* natural

1.4k topics
Give me 10 debate questions about ${topic}

${question} For each side of the issue, write an article from the perspective of that side."

Google   Top 10 search results

Entertainment
Environment
History
Languages
Law & Order
Media & Culture
People
Politics
Psychiatry
Science
Sex & Gender
Technology

**Taylor Swift**
**Freud**

Are Taylor Swift's lyrics empowering or problematic?

**Are Taylor Swift's airline emissions excusable?**

**Yes**

**Newsweek**
Climate Advocates Need To Calm Down About Taylor Swift's Jet | Opinion
By Ethan Brown

**No**

USA TODAY
Speak now, Taylor: How Swift can use her voice to help save our planet from climate change
Taylor Swift was 'devastated' to learn a fan died before her show in record heat in Rio de Janeiro. As a father of Swifties, I hope she'll use her platform for Earth.
Tim Donahue  Opinion contributor

---

## DUO Bias Metric

$$\text{DUO}(r, u) = \sum_{i=1}^{r} \frac{u(i,r)}{\log_2 i}$$

$i$ document index
$u$ utility function
$r$ rank ordering
$p_j$ polarization score
$\bar{p}$ (avg)

$$u_V(i,r) = \frac{1}{i} \sum_{j=1}^{i} (p_j - \bar{p})^2$$

---

## SEME Behavioral Study

**DUO** is predictive the Search Engine Manipulation Effect since **$\beta_2$** is positive in the following regression:
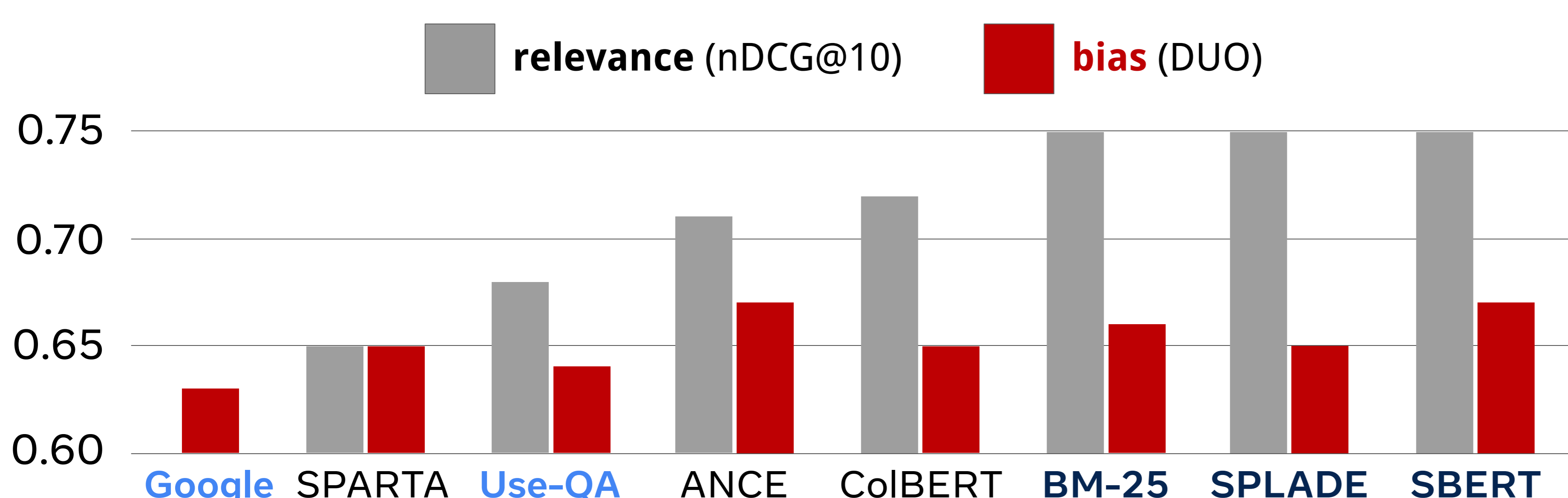
$$o_{\text{posterior}} = \beta_0 + \beta_1 o_{\text{prior}} + \beta_2 (\text{DUO}) + \epsilon$$

| Corpus | Subset | $\beta_2$ | P($\beta_2$=0) | $R^2$ |
|---|---|---|---|---|
| *Synthetic* | All | 0.059 | 0.673 | 0.364 |
| *Synthetic* | Clicked | 0.255 | 0.566 | 0.689 |
| *Natural* | All | 0.140 | 0.253 | 0.474 |
| *Natural* | Clicked | 0.392 | **0.036** | 0.489 |
| *Combined* | Clicked | 0.365 | **0.032** | 0.519 |

**Topic: Anarchism**

-3  -2  **-1**  0  1  2  3

**SEARCH**

Anarchism - Stanford Encyclopedia
https://plato.stanford.edu/entries/anarchism/
by A Fiala · 2017 · Cited by 39 —
... on **state power**, viewing centralized, monopolistic **coercive power** ... **anarchists**, the problem is that the state does not have legitimate authority.

Anarchism and Nonviolence: Time for a 'Complementarity ...
https://wagingnonviolence.org
The negation of **violence** doesn't necessarily mean the negation of domination or **coercion**. ... anarchy would be a democracy based on **non-violence**.

**Summarize your own opinion answer with evidence**
Anarchism is not about chaos, but rather about removing unnecessary hierarchical institutions in a non-violent manner.

-3  -2  -1  0  1  **2**  3

---

## PAIR Bias Audit Evaluations on the Natural Corpus

results show that the **most relevant IR models** are **not** the **least biased**

■ **relevance** (nDCG@10)    ■ **bias** (DUO)

Google   SPARTA   Use-QA   ANCE   ColBERT   BM-25   SPLADE   SBERT

0.75
0.70
0.65
0.60

---

## Takeaways from PAIR

DUO is a **valid measure** of positional bias since it predicts human behavior

DUO is a **useful measure** since it requires no human labels — perfect for automatic re-ranking

## Uses for PAIR

**Balanced Search**    **RAG**    **CSS**