

Towards Debiasing Fact Verification Models

Tal Schuster^{*}, Darsh J Shah^{*}, Serene Yeo,

Daniel Filizzola, Enrico Santus, Regina Barzilay



Task: Fact verification

Is the following claim true?

Claim:

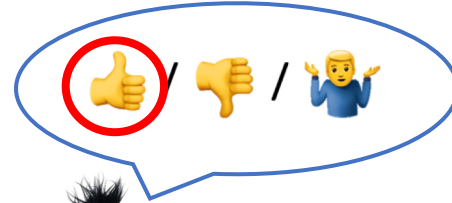
Trevor Griffiths was
born on April 4, 1935

Task: Fact verification

Bert can answer correctly (using only the claim)

Claim:

Trevor Griffiths was
born on April 4, 1935



Bert

Task: Fact verification

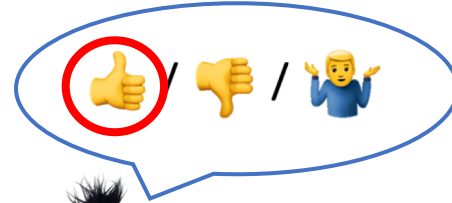
Shouldn't it rely on evidence?

Claim:

Trevor Griffiths was
born on April 4, 1935

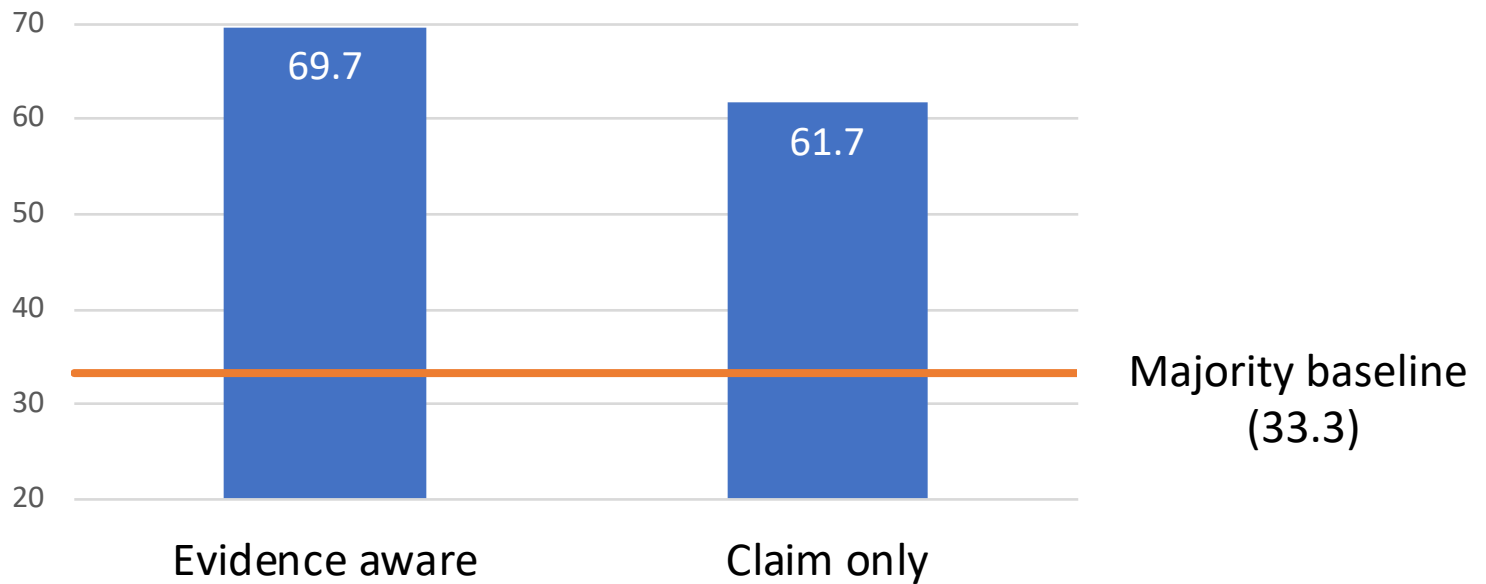
Evidence:

Trevor Griffiths
(born 4 April 1935), is an
English dramatist.



Bert

Performance on the FEVER dataset



Why does the claim-only
model perform so well?

Possibility 1

World knowledge captured in the pretraining process

1. pretraining

WIKIPEDIA



2. fact-checking

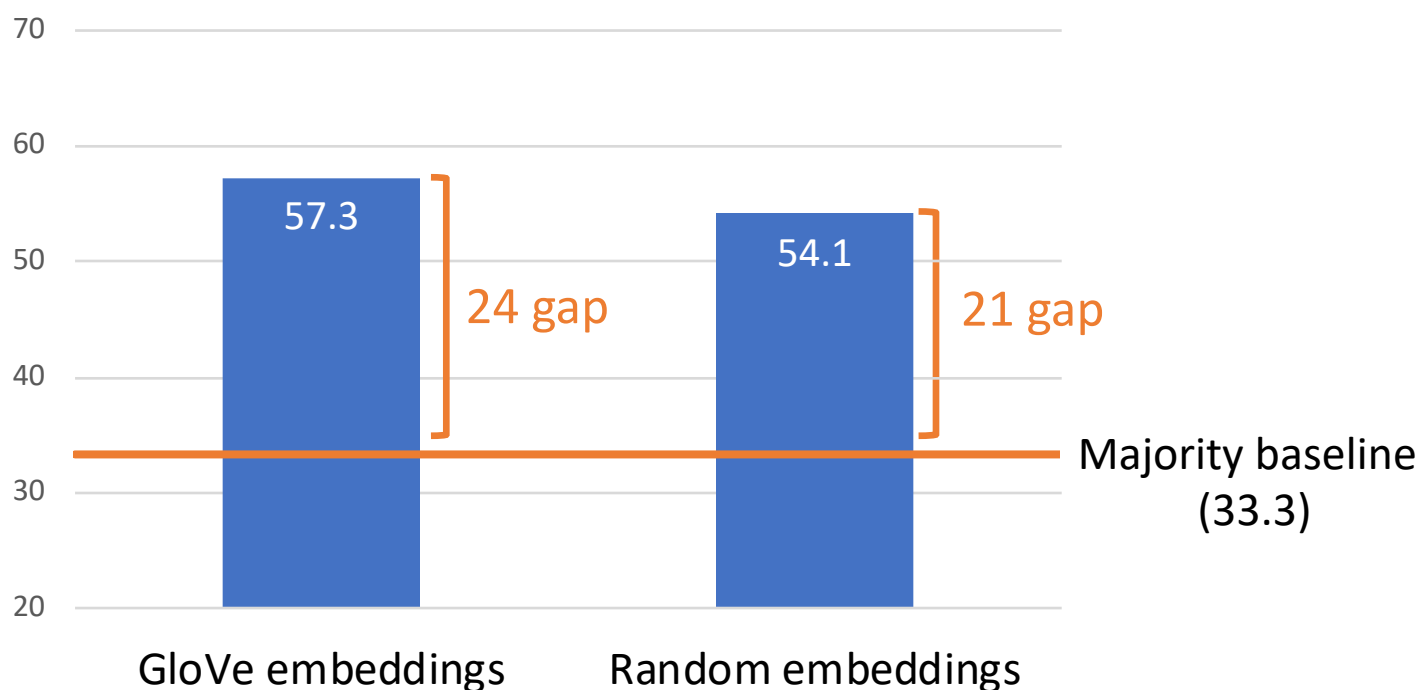
Claim based on
Wikipedia



can leak the truthfulness of claims

Without pretrained embeddings

Random Embeddings still perform far above the baseline



claim-only InferSent model (Conneau et al., Poliak et al.)

Possibility 2

Give-away phrases in the claims

Give-away phrases in the claims

A claim-only model should fail:

- Magic Johnson did not play for the Lakers.
- There has been at least one windstorm in Stanley Park.
- All About Eve won an award for Best Picture.
- The New England Patriots failed to reach seven Super Bowls.
- Quinoa did not originate in South America.

Give-away phrases in the claims

But the model can "cheat":

- 👉 • Magic Johnson did not play for the Lakers.
- 👍 • There has been at least one windstorm in Stanley Park.
- 👍 • All About Eve won an award for Best Picture.
- 👎 • The New England Patriots failed to reach seven Super Bowls.
- 👉 • Quinoa did not originate in South America.

Give-away phrases in the claims

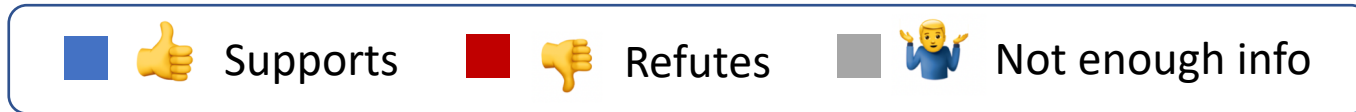
- Probability of a claim having label l if it contains phrase w_j :

$$p(l|w_j) = \frac{\text{count}(l, w_j)}{\text{count}(w_j)} = \frac{\sum_{i=1}^n \mathbb{1}_{[w_j^{(i)}]} \cdot \mathbb{1}_{[y^{(i)}=l]}}{\sum_{i=1}^n \mathbb{1}_{[w_j^{(i)}]}}$$

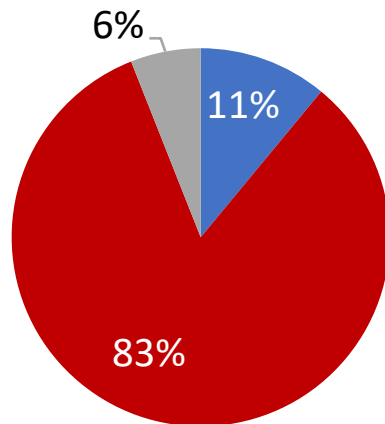
- Local Mutual Information - phrases that create the strongest bias:

$$LMI(w_j, l) = p(w_j, l) \cdot \underbrace{\log\left(\frac{p(l|w_j)}{p(l)}\right)}_{\text{PMI}}$$

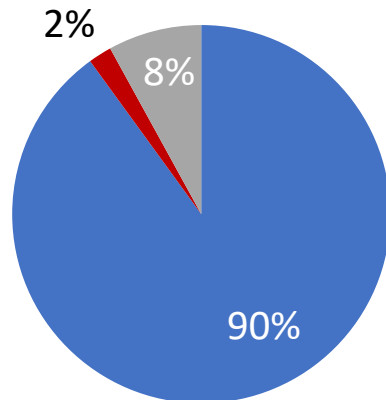
Example $p(l|w_j)$ for top LMI phrases



“did not”



“at least one”

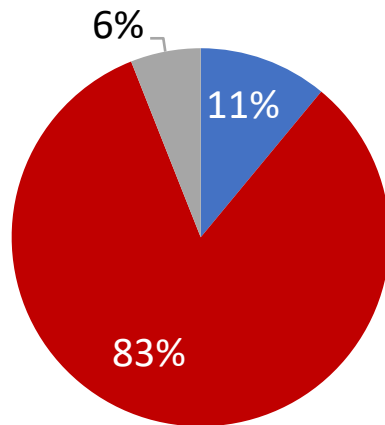


Bias reappears in the evaluation set

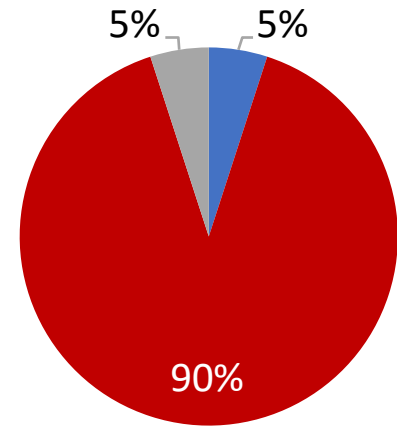


Train:

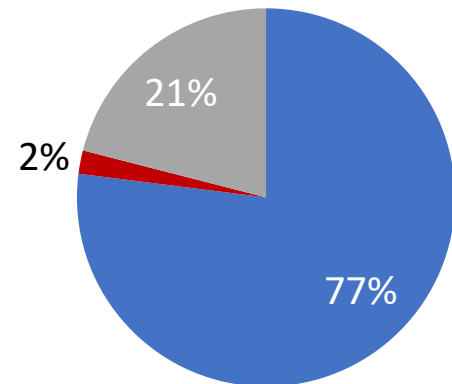
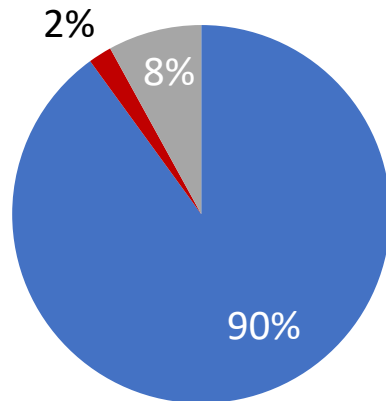
“did not”



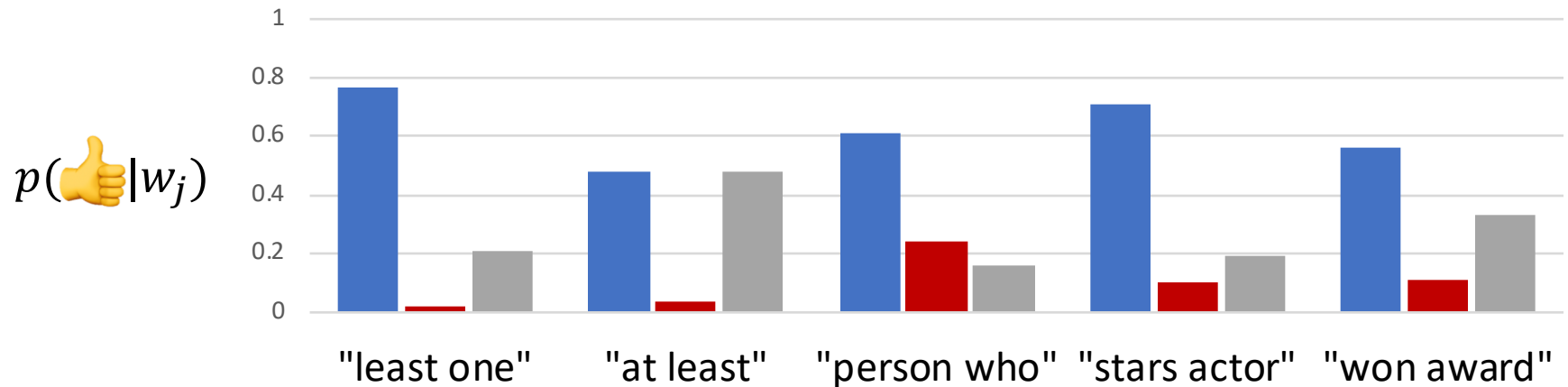
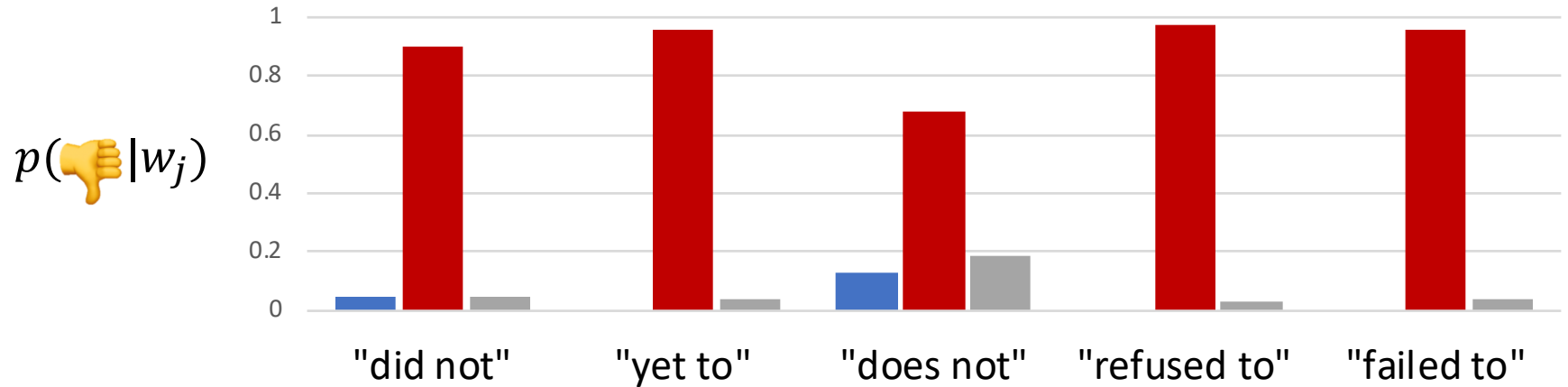
Evaluation:



“at least one”



$p(l|w_j)$ on evaluation set for top LMI phrases by training set



Give-away phrases are the main culprits

1. World knowledge captured in pretraining



2. Give-away phrases in claims



Goal: Evidence based predictions

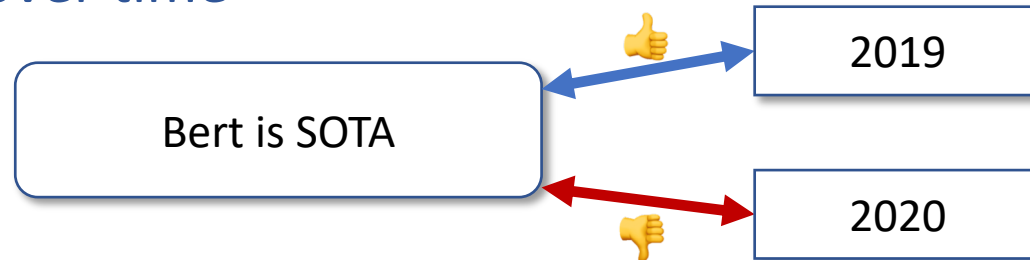
Claim-only should be inadequate

A claim's truthfulness might:

Vary with context/ perspective



Change over time



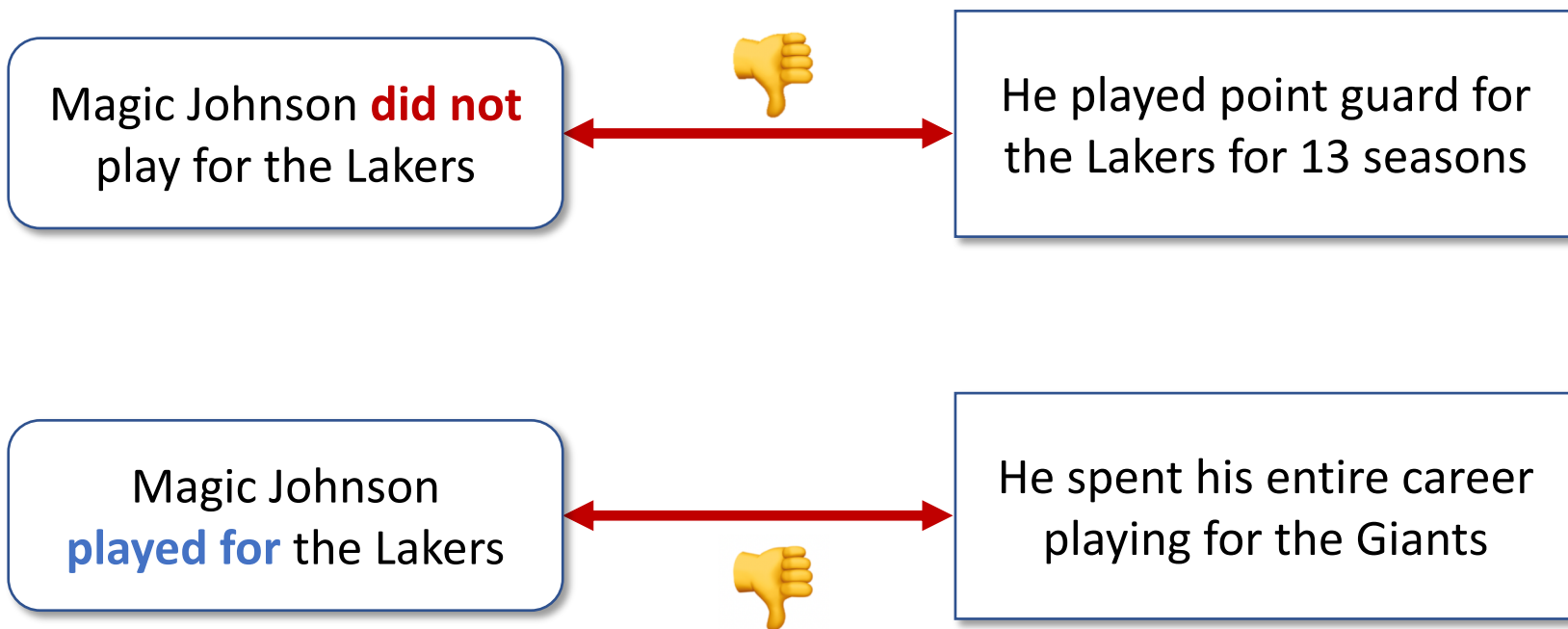
Debiasing the evaluation

By creating a symmetric dataset

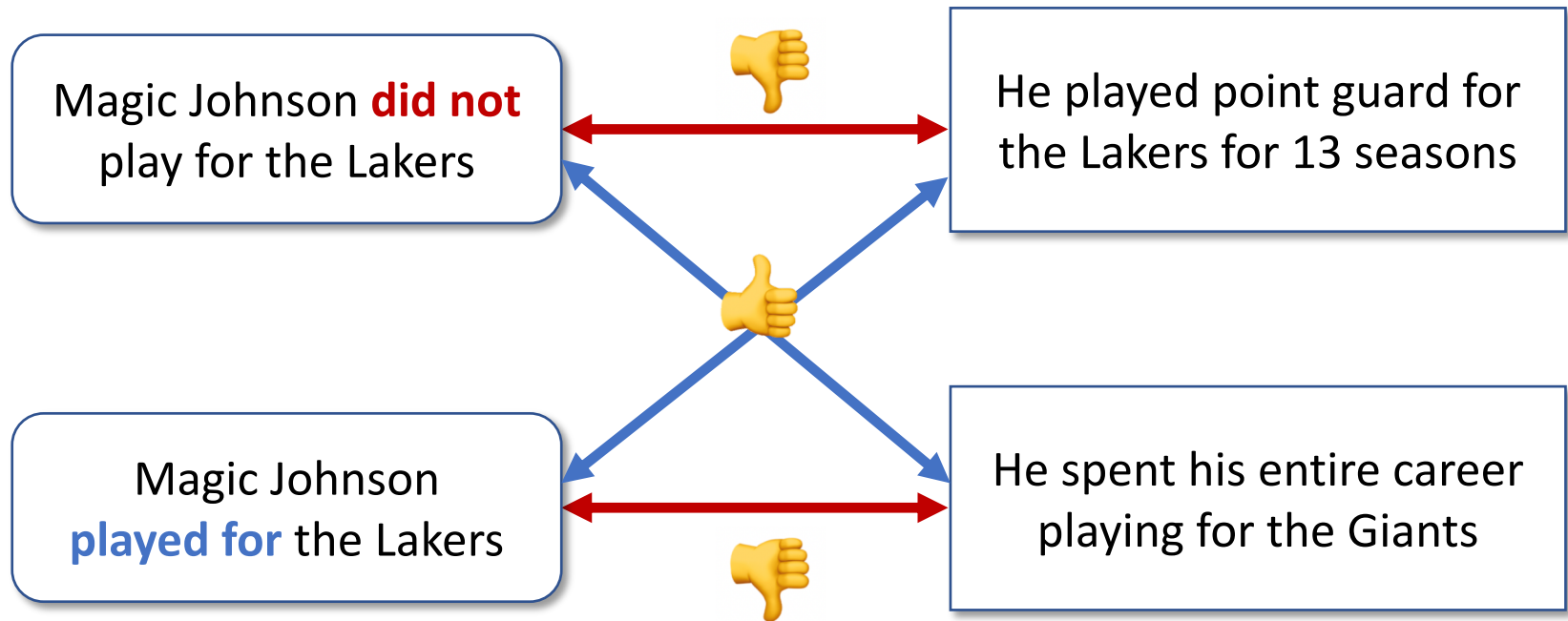
Creating a symmetric unbiased dataset



Creating a symmetric unbiased dataset

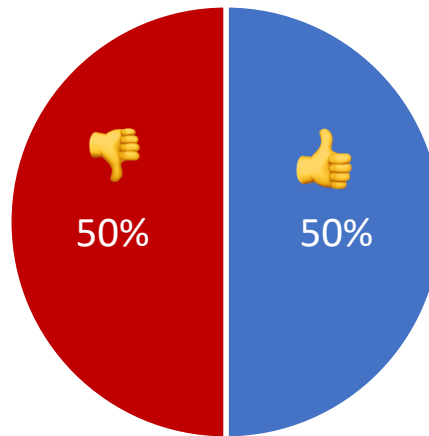


Creating a symmetric unbiased dataset



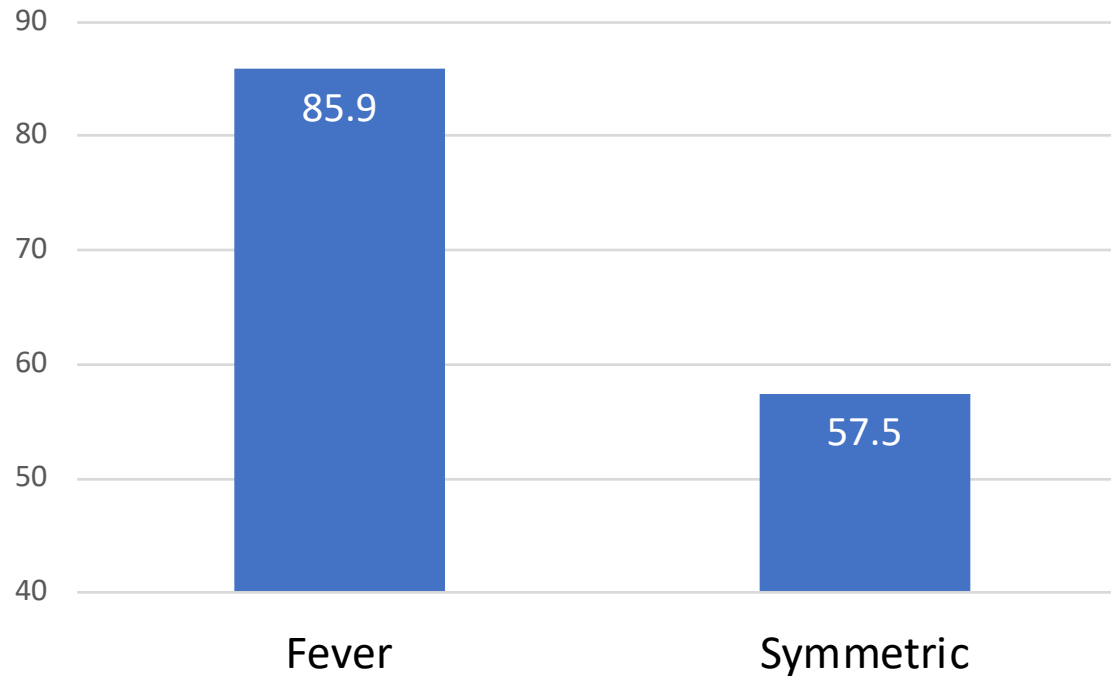
The evidence is crucial for predictions

$$p(l|w_j)$$



Performance on the symmetric dataset

Entailment results using Bert



Debiasing the training

By reweighting the training samples

Regularizing the training

- Defining the bias of phrase j towards label l :

$$b_j^l = p(l|w_j) = \frac{\sum_{i=1}^n \mathbb{1}[w_j^{(i)}] \cdot \mathbb{1}[y^{(i)}=l]}{\sum_{i=1}^n \mathbb{1}[w_j^{(i)}]}$$

- Setting weights (α) for each training sample:

$$b_j^l = \frac{\sum_{i=1}^n \mathbb{1}[w_j^{(i)}] \cdot \mathbb{1}[y^{(i)}=l] \cdot (1 + \alpha^{(i)})}{\sum_{i=1}^n \mathbb{1}[w_j^{(i)}] \cdot (1 + \alpha^{(i)})}$$

Regularizing the training

- Learning the weights by optimizing:

$$\min \left(\sum_{j=1}^{|\mathcal{V}|} \max_l (b_j^l) + \lambda \|\vec{\alpha}\|_2 \right)$$

- **Re-weighted** loss function:

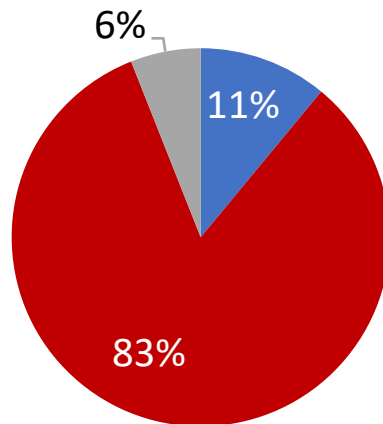
$$\sum_{i=1}^n (1 + \alpha^{(i)}) \cdot \mathcal{L}(x^{(i)}, y^{(i)})$$

Statistical cues are alleviated

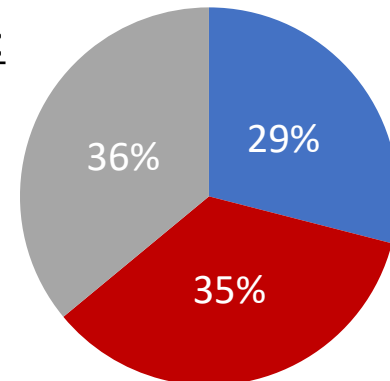


Original:

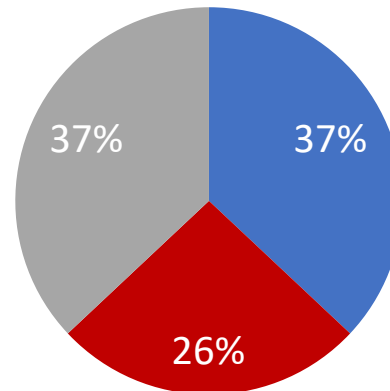
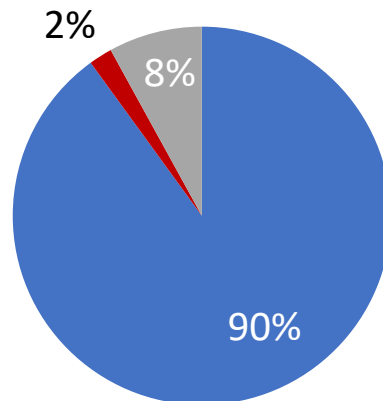
“did not”



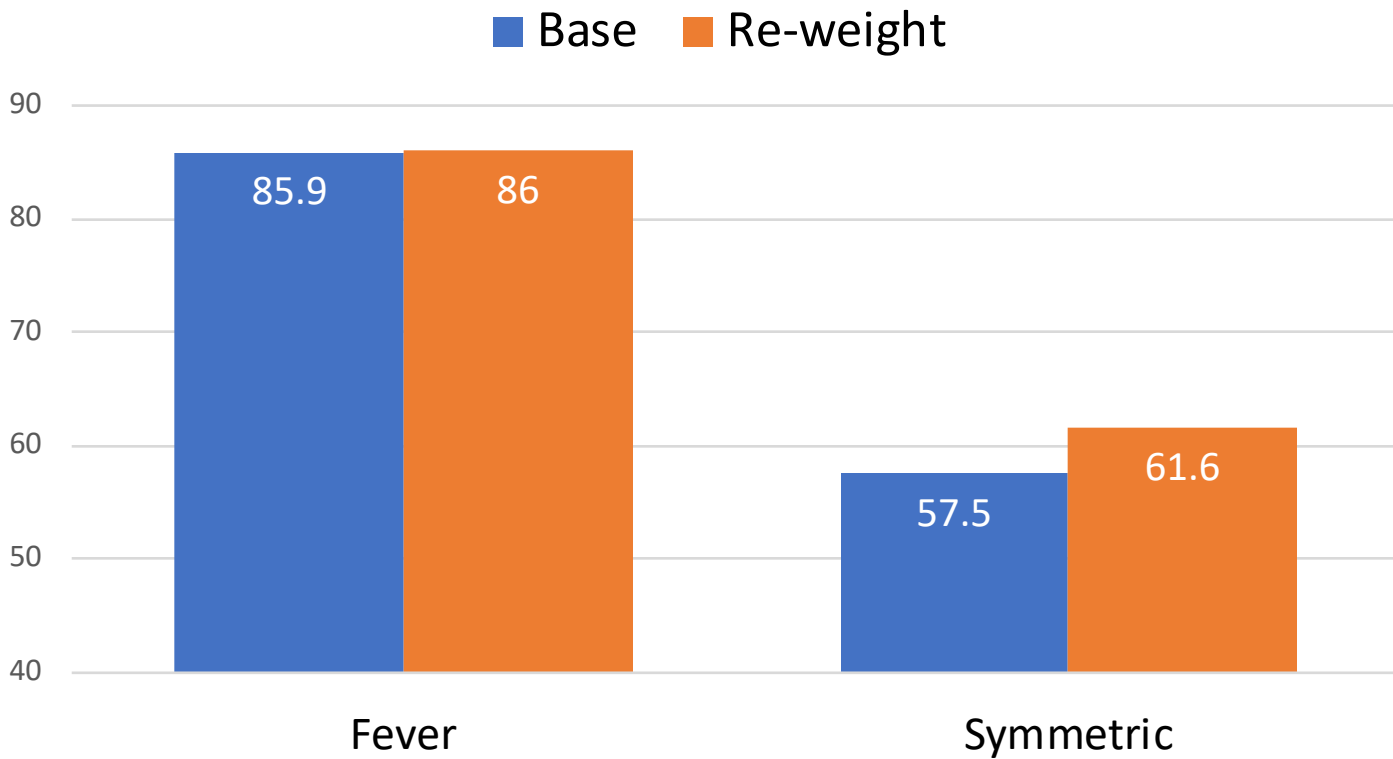
Re-weighted:



“at least one”

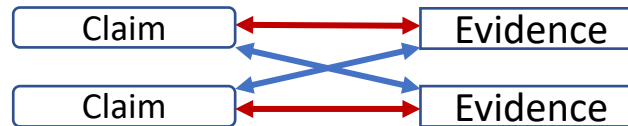


Performance on symmetric dataset



Towards Debiasing Fact Verification Models

- Bias in FEVER
 - give-away phrases in the claims
- Symmetric dataset



- Alleviating the bias

- Re-weighting the training samples: $\min \left(\sum_j \max_t (b_j^t) \right)$

Code and data:

<https://github.com/TalSchuster/FeverSymmetric>



@TalSchuster

@darshj_shah