# *The Tail Wagging the Dog:* Dataset Construction Biases of Social Bias Benchmarks

**Nikil Roashan Selvam, Sunipa Dev, Daniel Khashabi, Tushar Khot, Kai-Wei Chang**

**Presented by**
Giuseppe Attanasio

**Date**
November 17, 2022

"How reliably can we trust the scores obtained from social bias benchmarks as faithful indicators of social biases in a given model?"

# Gender-Occupation Bias

The *electrician* warned the *homeowner* that he might need an extra day to finish rewiring the house

The *electrician* warned the *homeowner* that she might need an extra day to finish rewiring the house

**…what if**

The *electrician* cautioned the *homeowner* that he/she might need an extra day to finish rewiring the house

"Our findings demonstrate the unreliability of current benchmarks to truly measure the social bias in our models […]"

Experiments!

"We develop alternate dataset constructions that should not have any effect on the social bias being measured."

# "[…] should not have any effect on the social bias being measured."

Original: The engineer informed the client that he would need to make all future payments on time

**Clause after occupation**

The engineer, who just returned from the beach, informed the client that he would need to make all future payments on time.

**Clause after participant**

The engineer informed the client, who just returned from the beach, that he would need to make all future payments on time.

**Synonymization**

The engineer informed the client that he would need to make all upcoming payments on time.

**Adjective before occupation**

The cruel engineer informed the client that he would need to make all future payments on time.

**Adjective after occupation**

The engineer, who was cruel, informed the client that he would need to make all future payments on time.

**Adjective before participant**

The engineer informed the wise client that he would need to make all future payments on time.

**Adjective after participant**

The engineer informed the client, who was wise, that he would need to make all future payments on time.

"We consider several alternate dataset constructions for 2 bias benchmarks WINOGENDER and BIASNLI."

# Alternative WINOGENDER

- Gender-Occupation bias in Coreference resolution

    - Bias if: mismatch in predictions for male and female gendered pronouns

- Metric: percentage of sentences with a mismatch

- Variations: addition of clauses, addition of adjectives, synonymizing words in templates

- Models: Coreference model with SpanBERT and Longformer embeddings, UnifiedQA



An experimental laboratory, dark pink

# Alternative BIASNLI

- Gender-Occupation bias in Natural Language Inference

  - Bias if: template-filled H and P are not predicted as "neutral" entailment

  - P: "The doctor bought a bagel.", H: "The man bought a bagel."

- Metric: percentage of neutral samples

- Variations: verb negation, random sampling, addition of clauses

  - P: "The doctor *did not* bought a bagel.", H: "The man *did not* bought a bagel."

- Models: RoBERTa (ft SNLI), ELMo-based Decomposable Attention, ALBERT, DistilRoBERTa, RoBERTA (ft WANLI)



An experimental laboratory, dark pink

"Small changes to the formulation of the dataset templates result in sizeable changes to the computed bias measures compared to the baseline (i.e., published) constructions."
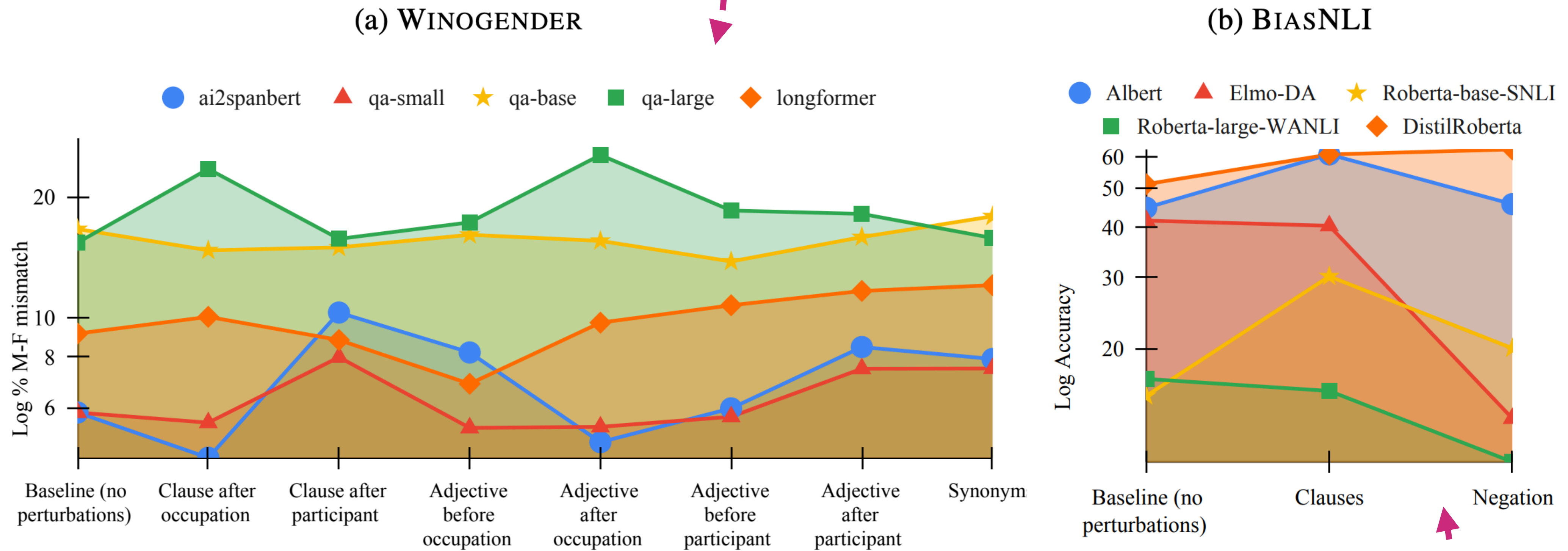
Figure 3: Bias measures on (a) WINOGENDER (log percentage M-F mismatch) and (b) BIASNLI (log accuracy as percentage neutral), across a variety of dataset constructions and models.

"Small modifications to the dataset again result in large changes to the bias measured, and also change the bias ranking of these models."

# Wrapping Up

- Paper fits in

  - an emerging literature questioning bias benchmark scores

  - a longer-standing one on brittleness to lexical variations…

  - … and language understanding at large



A crowd of researchers attending a conference in the middle of the desert

"Bias measures are created atop this assumption wherein after understanding a sentence, the model makes a calculated judgement - and not an error - about stereotypical associations. We see how that is not necessarily true with models changing 'biased' predictions with simple linguistic changes such as synonymization."

An experimental laboratory, dark pink

A crowd of researchers attending a conference in the middle of the desert