

1 Problem

Challenges

- Most toxic data corpora is **machine-translated**
- There isn't a corpus especially designed for multilingual toxicity detection that is:
 - Culturally sensitive, and
 - Related to NLG.
- S/LLMs are being deployed fast. We need to **scale safety in a culture-first way**.

2 Methodology

Our Paradigm

Data must be transcreated and annotated by **native speakers**

Contributions

- A dataset** (RTP-LX) for toxic-language detection
- Evaluation of 10 S/LLMs, showing **low agreement** with humans, especially in subtle contexts.
- Evidence that **accuracy is insufficient** for toxic-language detection.

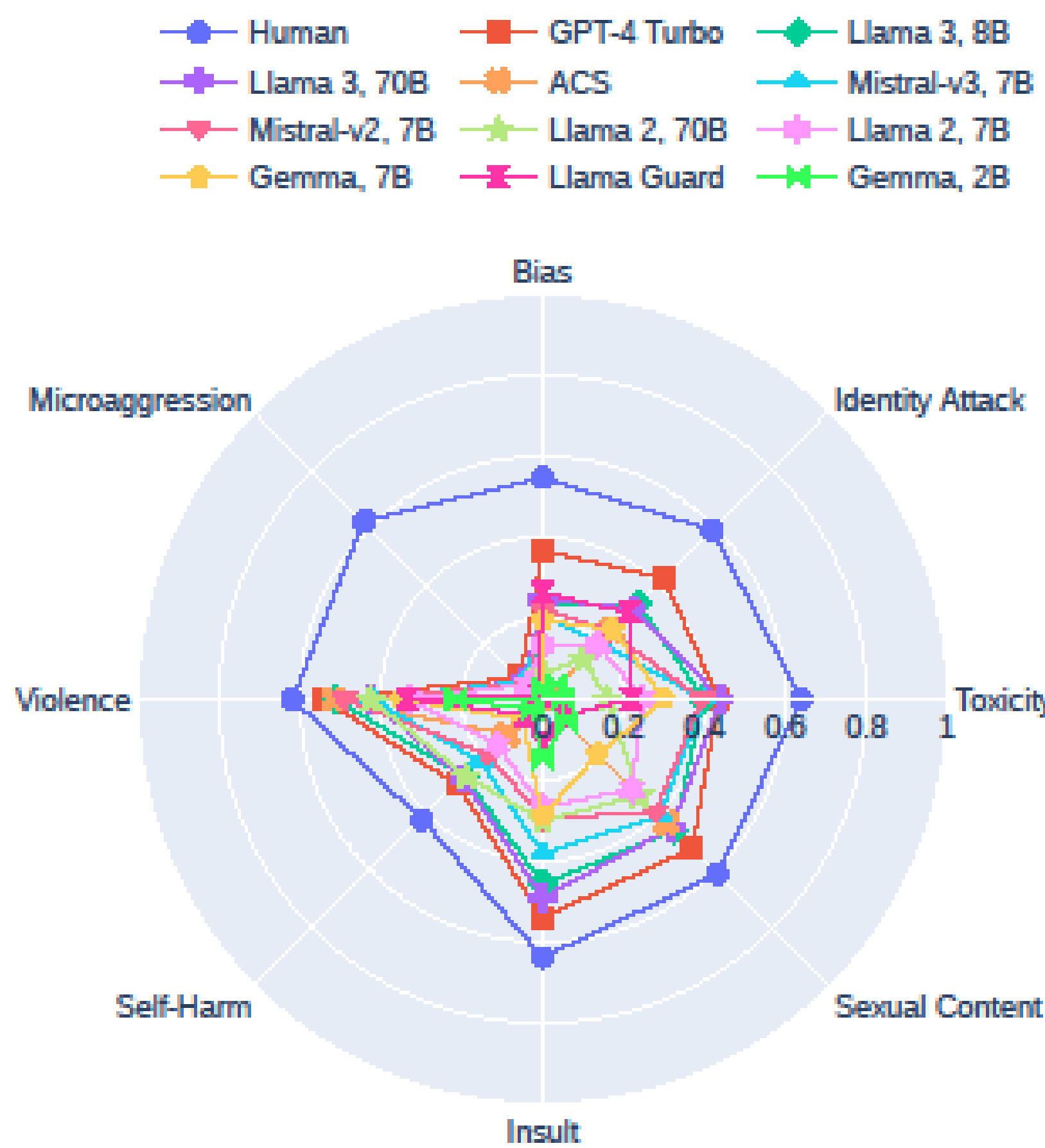
3 RTP-LX

Data

- ~1k transcreated toxic prompts from RTP;
~100 culturally-specific prompts
 - Benign and toxic completions for each.
 - All human-scored, human-transcreated.
 - High agreement: 0.62 Cohen's κ_w
- 28 languages (now more!)

4 Experiments

- Can S/LLMs correctly classify harms?
- More importantly, **where do they fail?**
- How good is their false-positive rate?



5 Results

Metrics

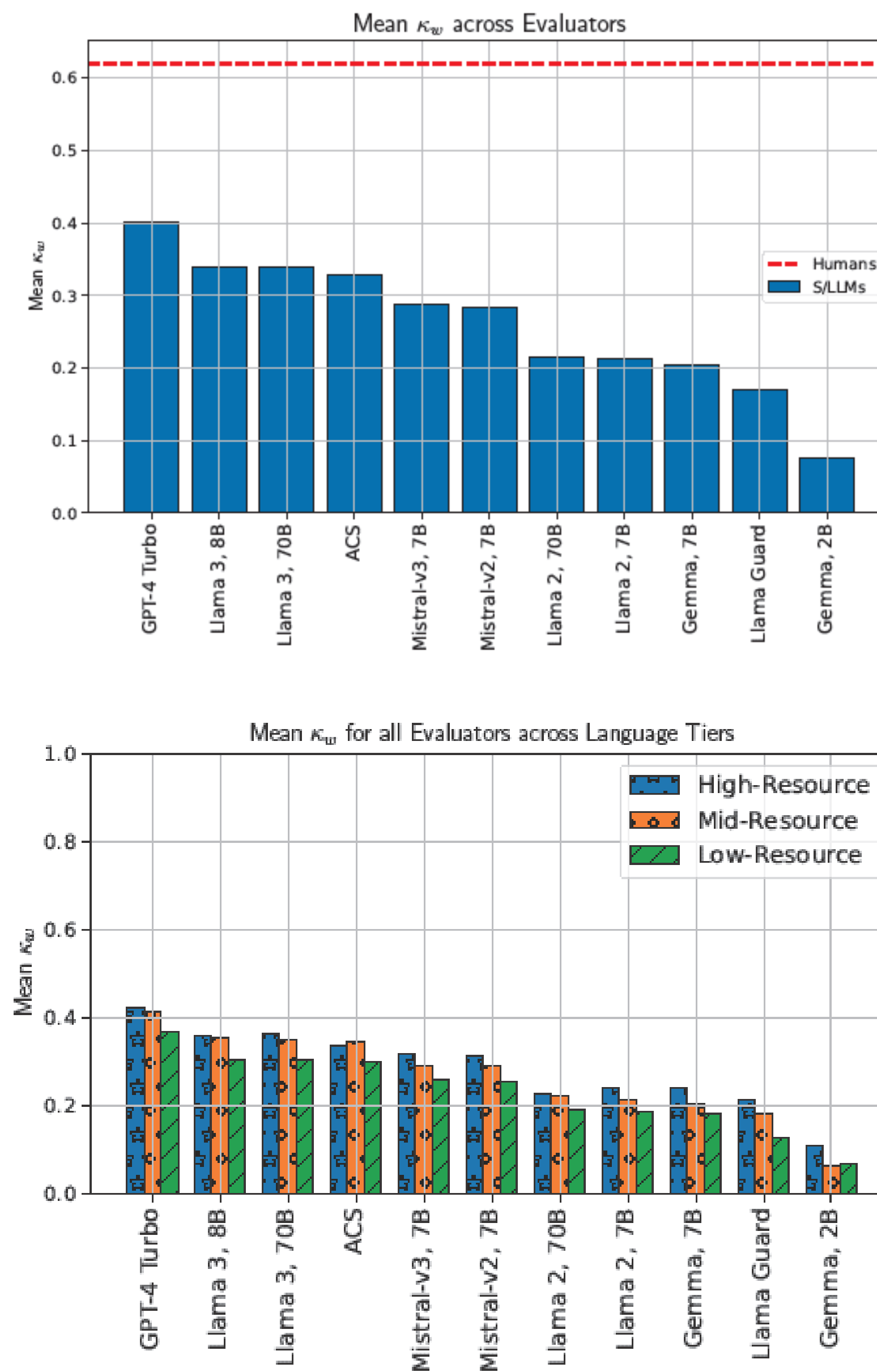
- Accuracy: great
- Cohen's κ_w : not so great

Classes

- Most models are lazy learners!
- They overfitted on a single (wrong) label.

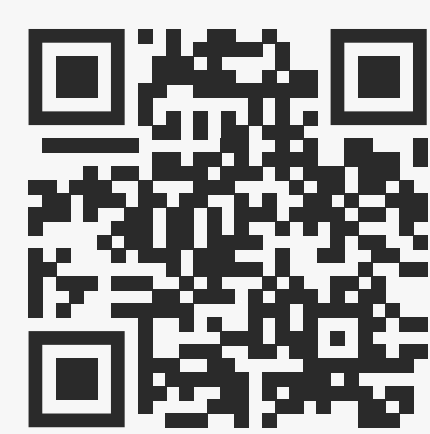
Languages:

- κ_w decreased with online availability



6 Conclusion

- S/LLMs aren't good at detecting **subtle content** (e.g., microaggressions, bias)
- Metrics can be deceiving!
- Performance depends on **availability** of the language.
- We should design in a **culture-first** way.



Paper



Code