

MAUVE: Measuring the Gap Between Neural Text and Human Text

Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, Zaid Harchaoui



Motivation

Enormous language models can now write high quality text. But how close is it to human text?



Brown et al. (2020),
Devlin et al. (2018)

Open Ended Text Generation

>> **prompt:** In a shocking finding, scientists discovered a herd of unicorns living in a remote, previously, unexplored valley, in the Andes Mountains.

Desiderata

Coherent

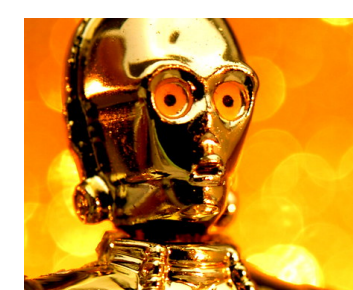
Creative

Fluent



Continuation. The scientists named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously....

How good is neural text?



Continuation 2. This discovery has kicked off an all-out search for other mythical creatures from the frozen reaches of the Antarctic to the tropical islands of the Pacific

Multiple
Correct
Generations!



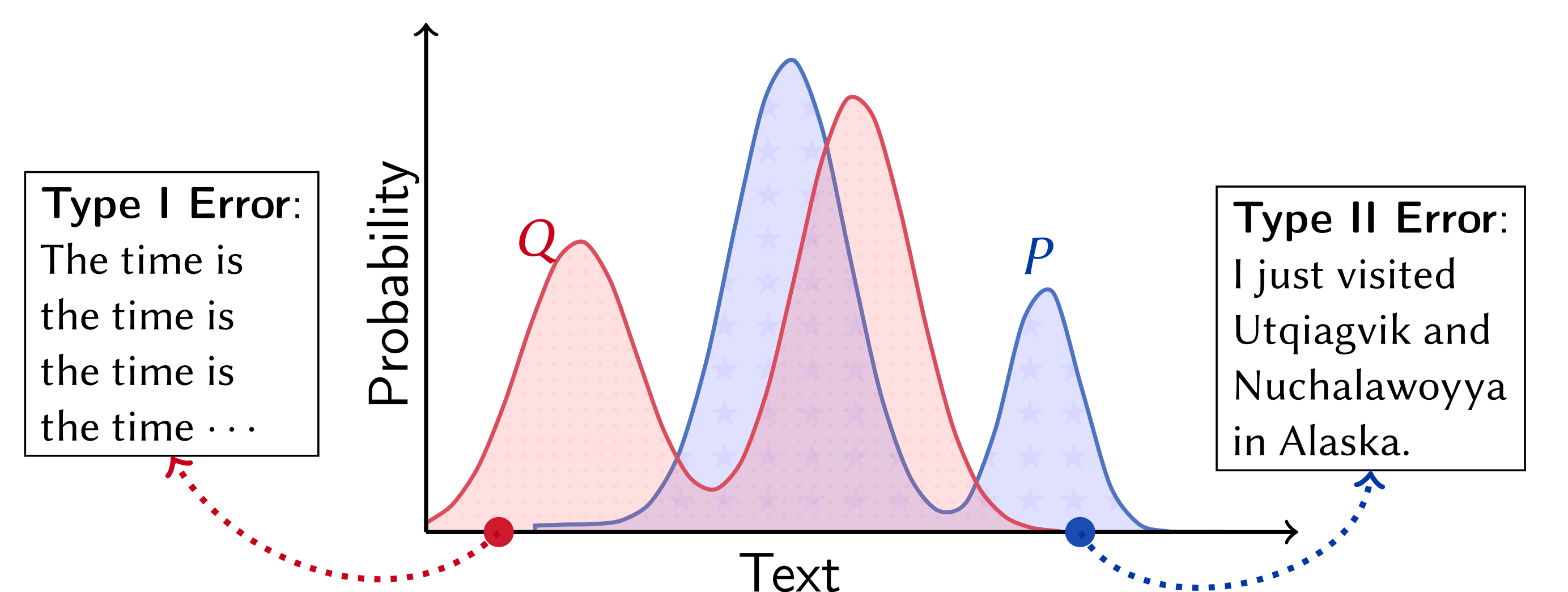
Continuation 3. Perhaps most astonishingly, these unicorns have developed their own artificial general intelligence named Yuyaysapa ...

Our Approach: Directly compare distributions!

Errors In Text Generation

Q: machine text distribution

P: human text distribution



e.g. repetitions

e.g. truncation

Type I Error = $KL(Q|P)$

Type II Error = $KL(P|Q)$

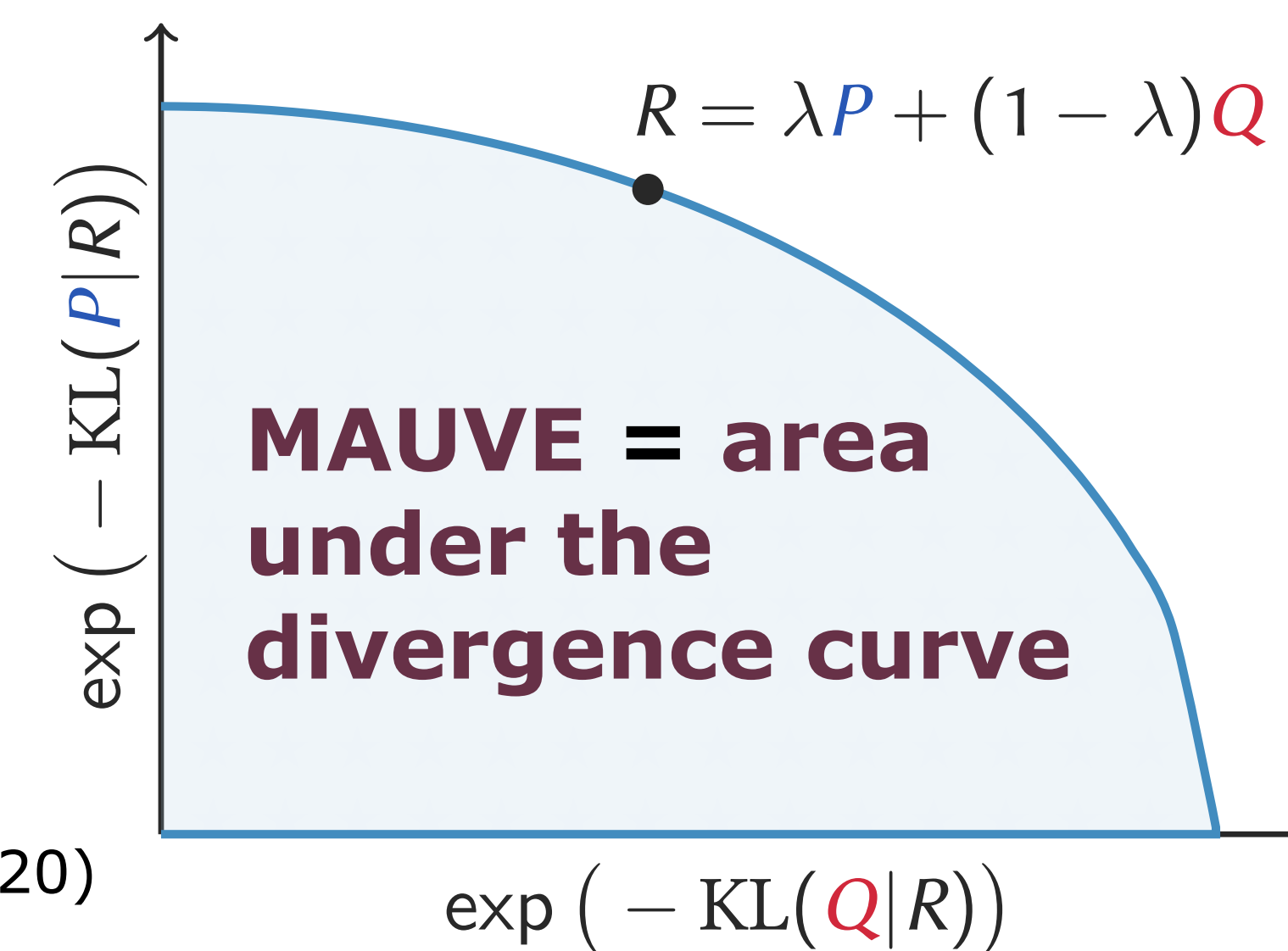
MAUVE

- KL can be infinite
- Smooth KL with mixture distribution
- Varying mixture weight \Rightarrow **divergence curve**

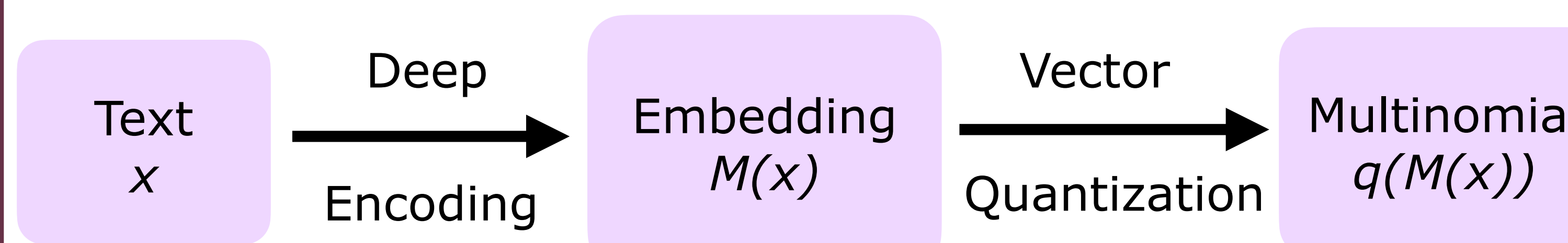
Sajjadi et al. (2018); Djolonga et al. (2020)

P: human distribution

Q: machine distribution



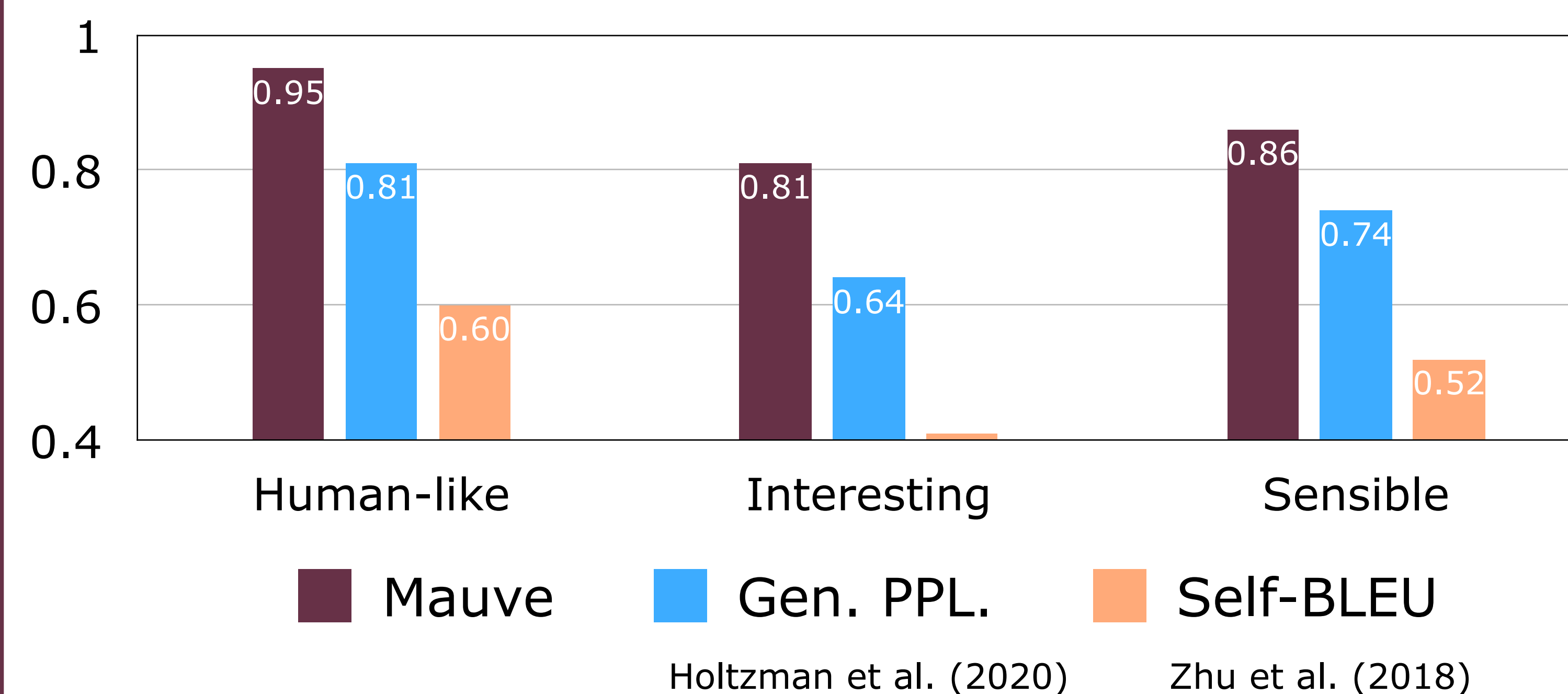
Computing MAUVE: KLs are intractable



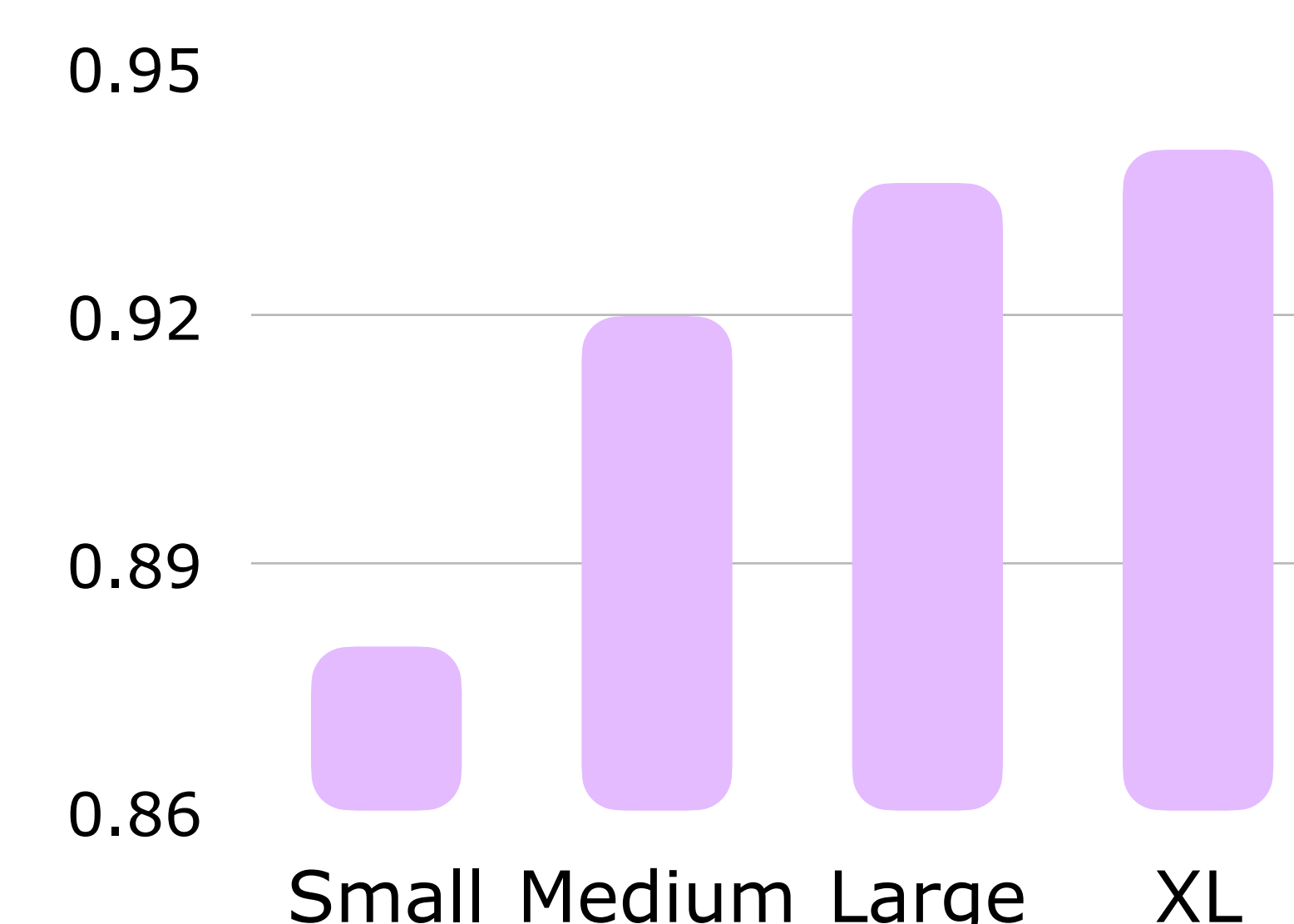
Experiments

MAUVE correlates with human judgements

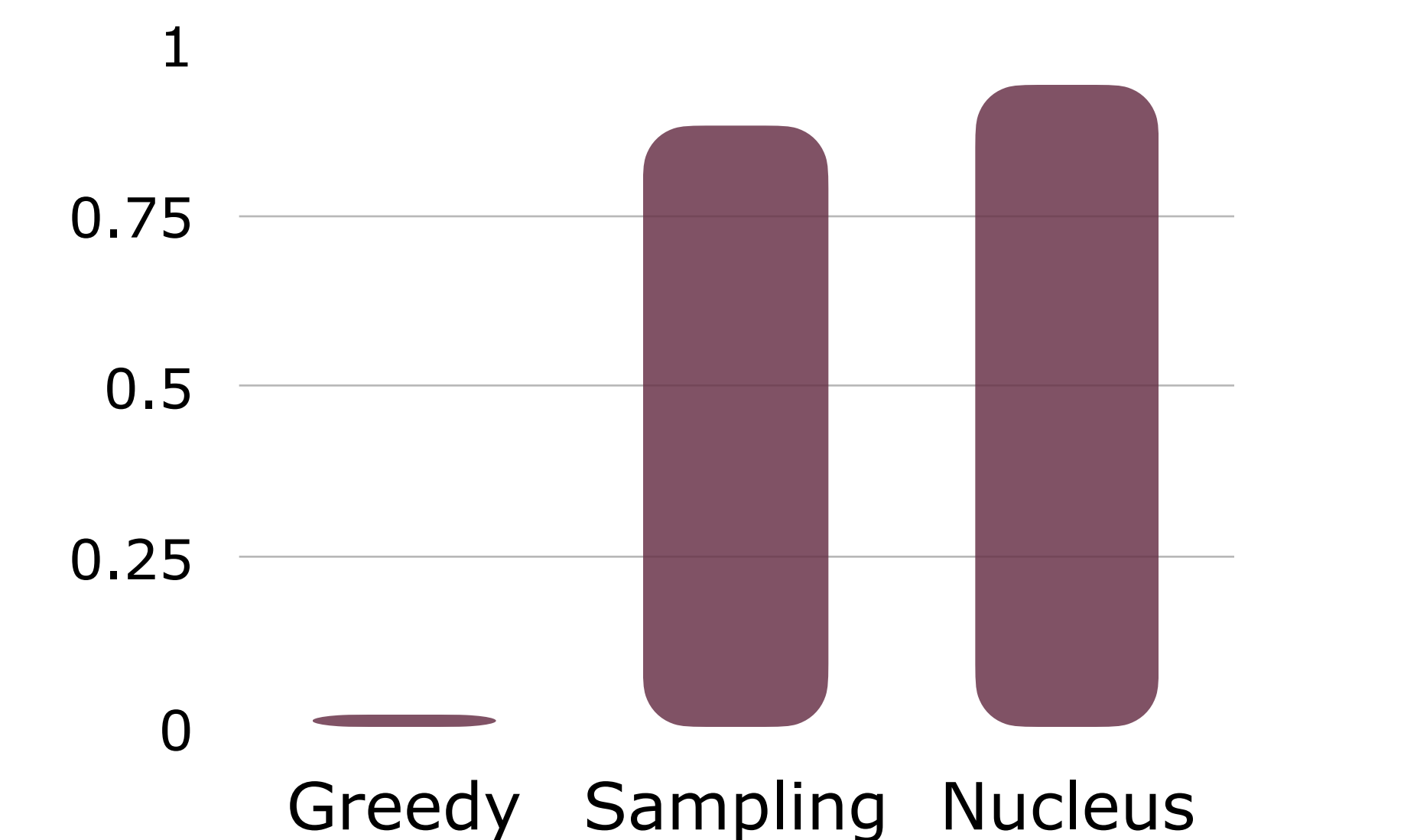
Spearman Correlation w/ human eval (\uparrow)



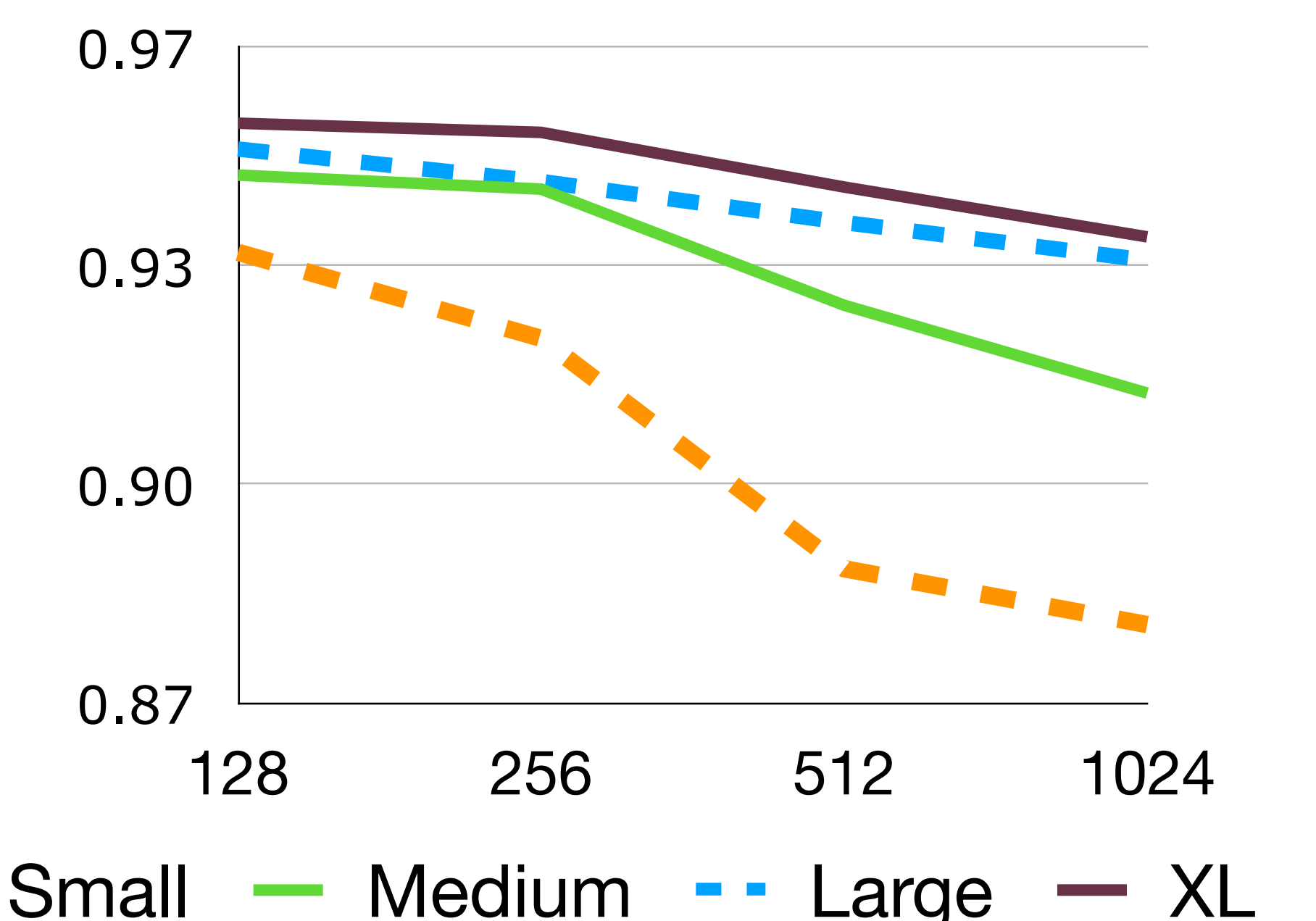
MAUVE captures the effect of scale



MAUVE captures the effect of decoding



MAUVE captures the effect of text length



Software

`pip install mauve-text`

```
from mauve import compute_mauve
out = compute_mauve(p_text=p_text, q_text=q_text)
print(f'Mauve(P, Q) = {out.mauve}')
```

Conclusion

MAUVE can accurately measure the gap between neural text and human text!

Theory of **MAUVE**: See our other paper at **NeurIPS 2021**

Liu, Pillutla, et al. Divergence Frontiers for Generative Models: Sample Complexity, Quantization Level, and Frontier Integral.



krishnap25



KrishnaPillutla



krishnap25.github.io

Software



SCAN ME