

# Can We Truly Measure Uncertainty in Language Models?

2nd Supervisor:
Prof Mrinmaya Sachan

ETH Zürich

Xueyan Li, Dr Jonas Geiping

## Token Level Uncertainty-Aware Chain-of-Thought Reasoning

#### Motivation:

- Step-wise token probabilities are discarded in future generations → future generations condition on previous tokens only, their uncertainties are lost
- Could be valuable to propagate current uncertainty to future steps as a signal to reflect or reason more around past uncertainty

Identifying errors in reasoning chains is a very hard task!			
Method	Question	Reasoning	Answer
	LLM		
Logits/entropy:	P1 P2 P3	P4 P5 P6	P7 P8
Input sequence:	Q1 Q2 Q3	R1 R2 R3	A1 A2

- Models: Llama3.1-1B and 8B
- Benchmarks: GSM8K, big-bench-mistake

GPT-4-Turbo

GPT-3.5-Turbo

PaLM 2 Unicorn

Llama-3.1-8B

aware finetuned

Llama-3.1-8B uncertainty 23.33

Gemini Pro

- P values are discretized into 10 bins
  - Entropy across vocabulary at each step

38.33

44.00

20.00

21.67

22.00

20.00

- Max softmax value
- int(max softmax value \*10)
- Supervised fine-tuning

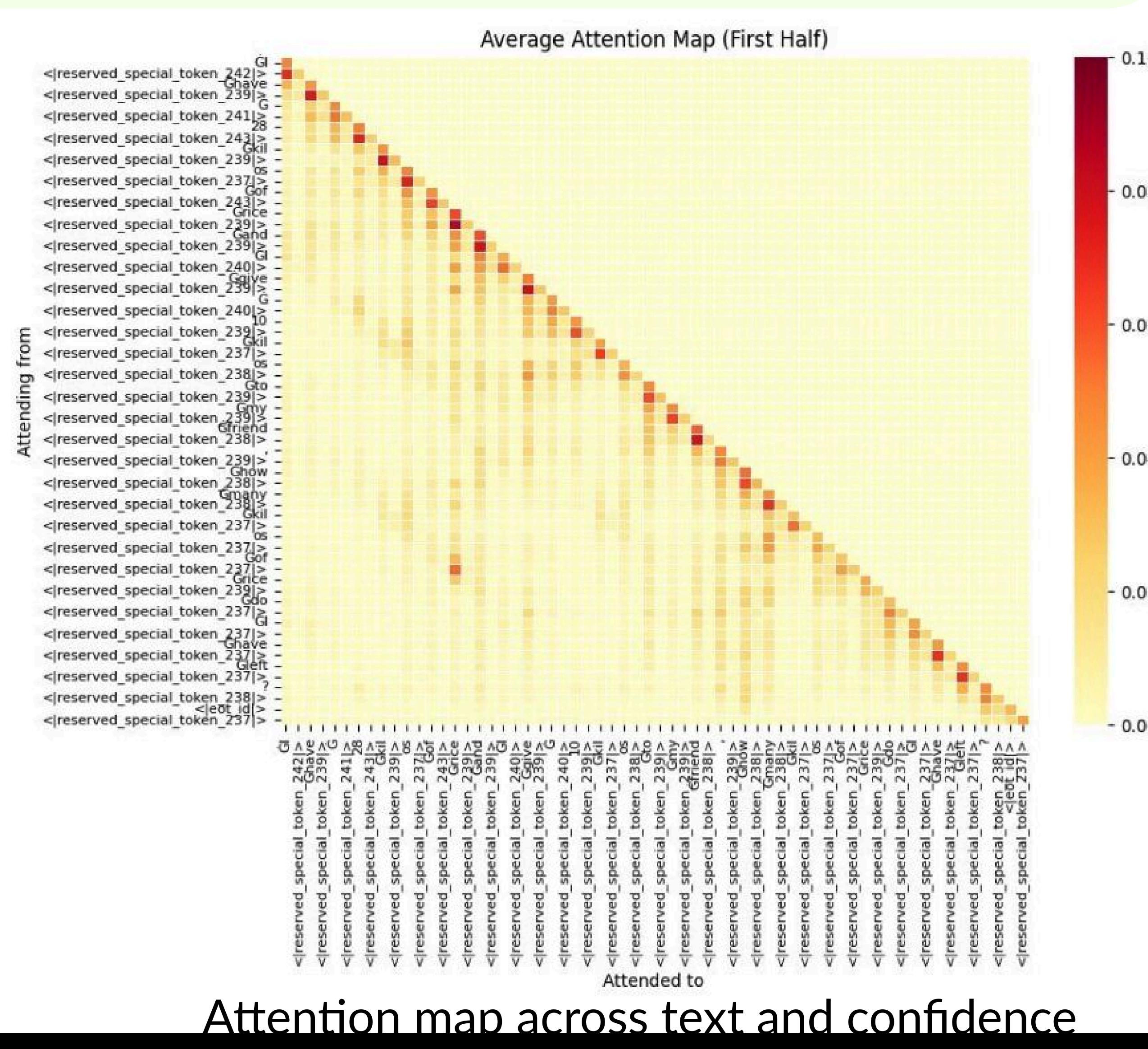
Augmented input sequence: Q1 Q2 Q3 P3 R1 P4 R2 P5 R3 P6 A1 P7 A2 P8

Loss: R1 R2 R3 A1 A2

### Results

- No improvement in performance after SFT. Similar or deteriorated performance
- Attention map for text and confidence tokens show vertical strips with gaps where the gaps are the lack of attention to special reserved tokens → special confidence tokens are ignored
- Baseline and interleaving setup have nearly identical loss curves
- Various P values result in different start loss, but converge to around the same point
- → not learning anything from uncertainty values!!!





## Uncertainty-Aware Temperature Adaptation for Better Calibration

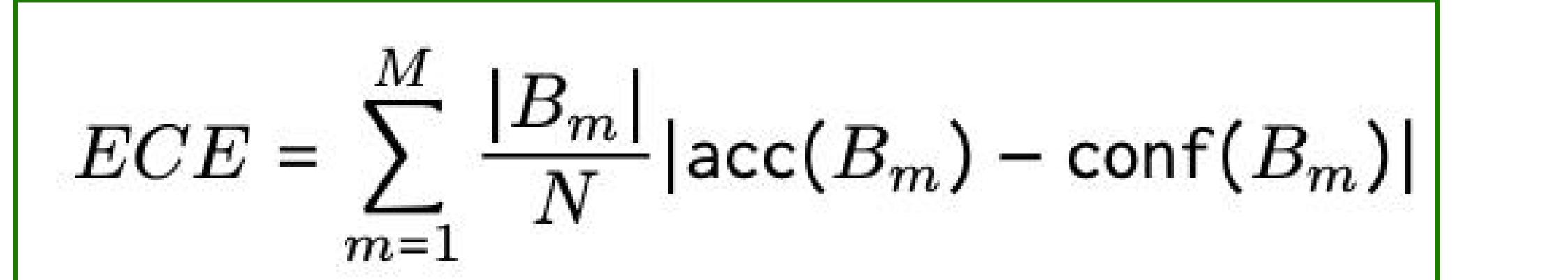
Overconfident model

Question

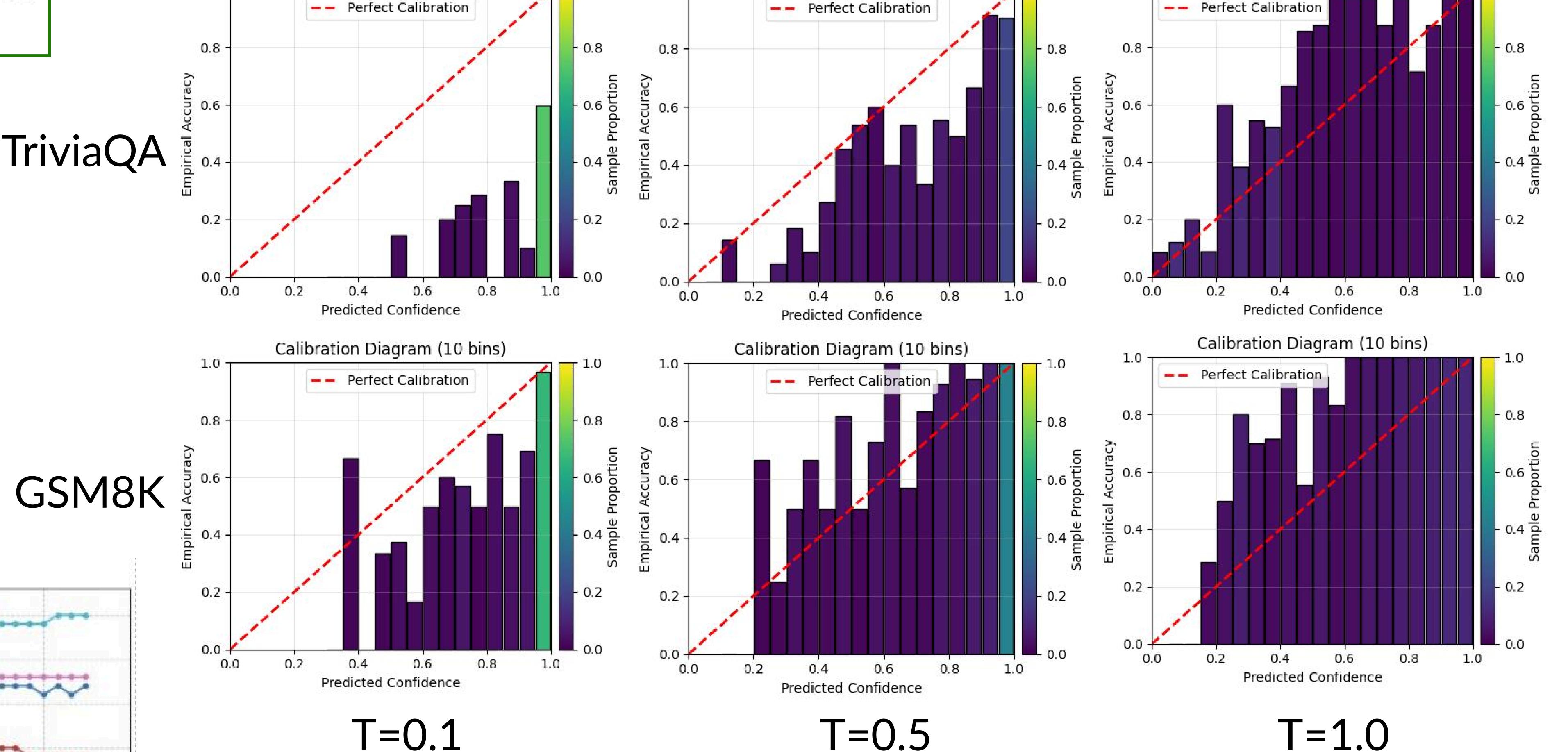
Language

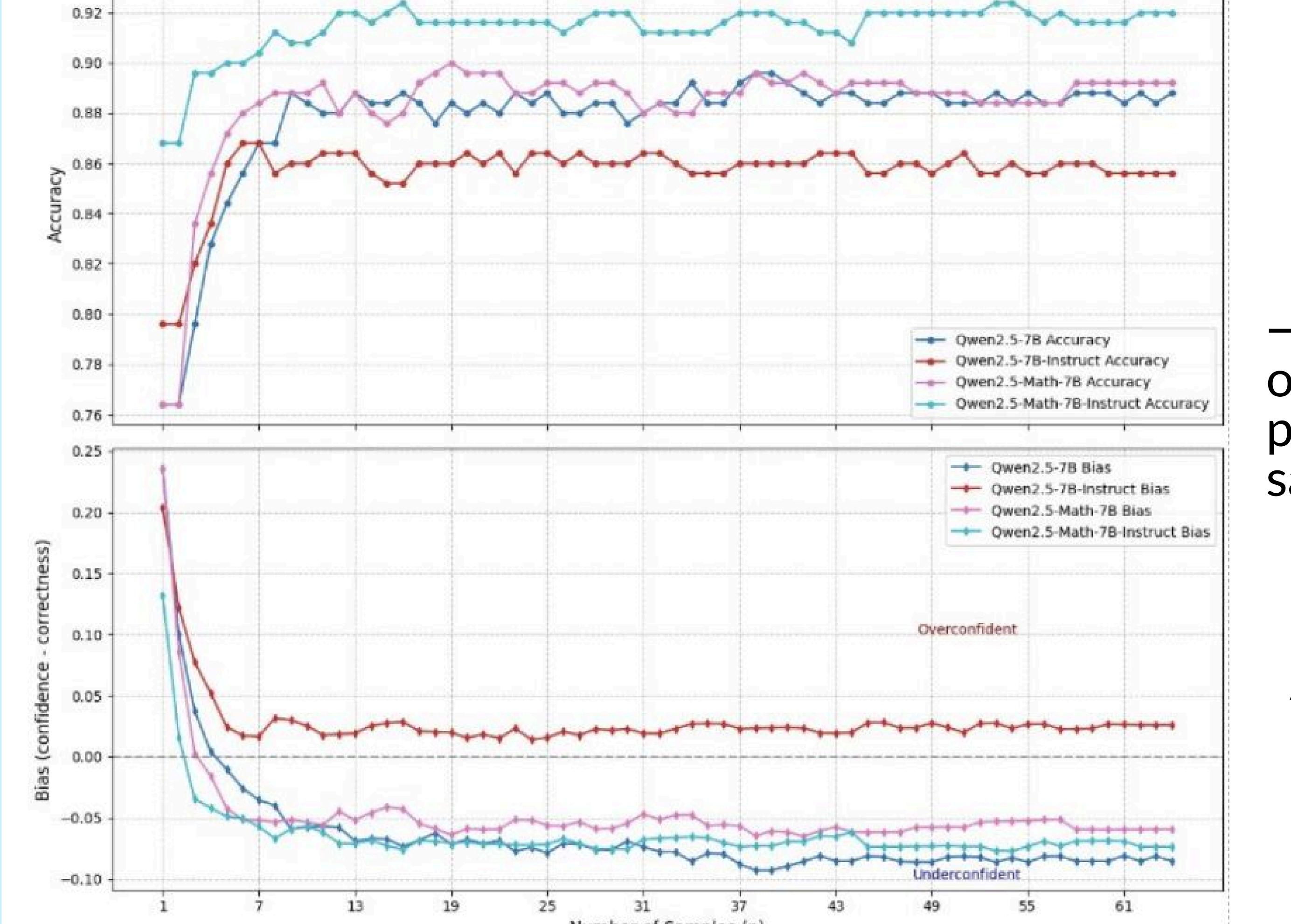
#### Motivation

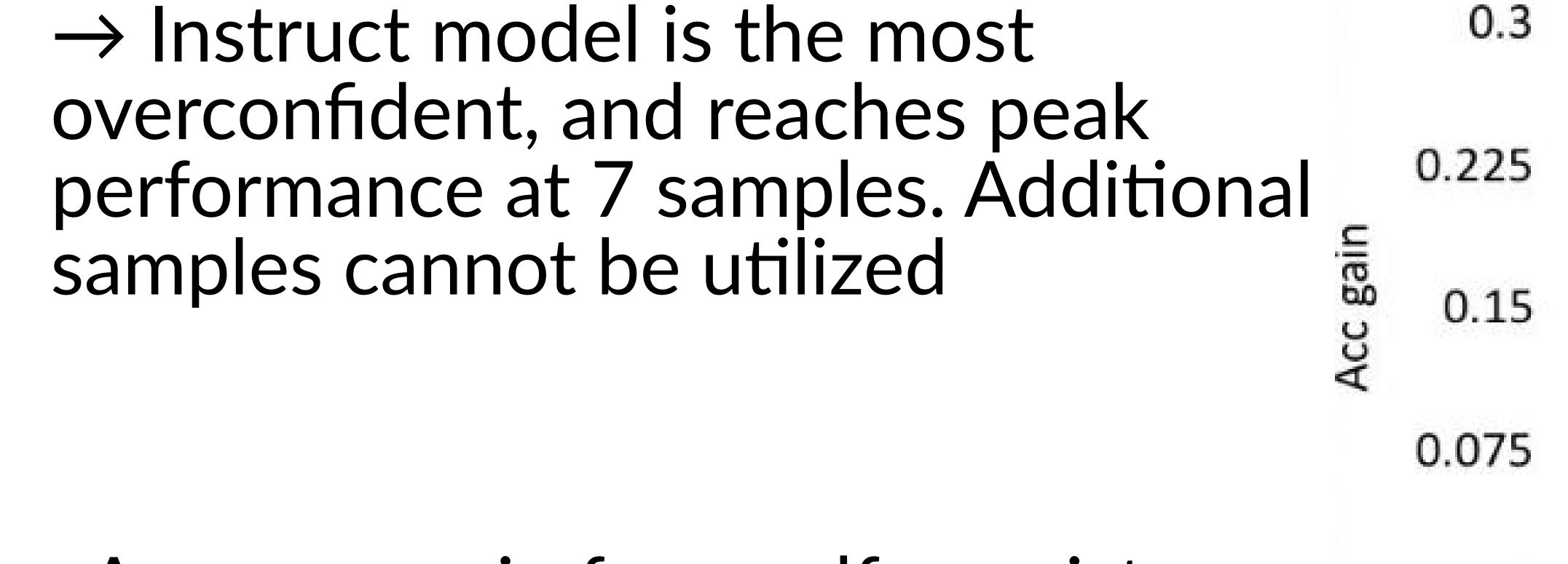
- Modern LLMs that have undergone instruction finetuning and RLHF are usually overconfident → logit values are peaked and do not represent the true uncertainty → poorly calibrated
- How to better calibrate models by adapting temperature at each generation step?
- Does better calibration help self-consistency by majority voting over diverse reasoning paths?

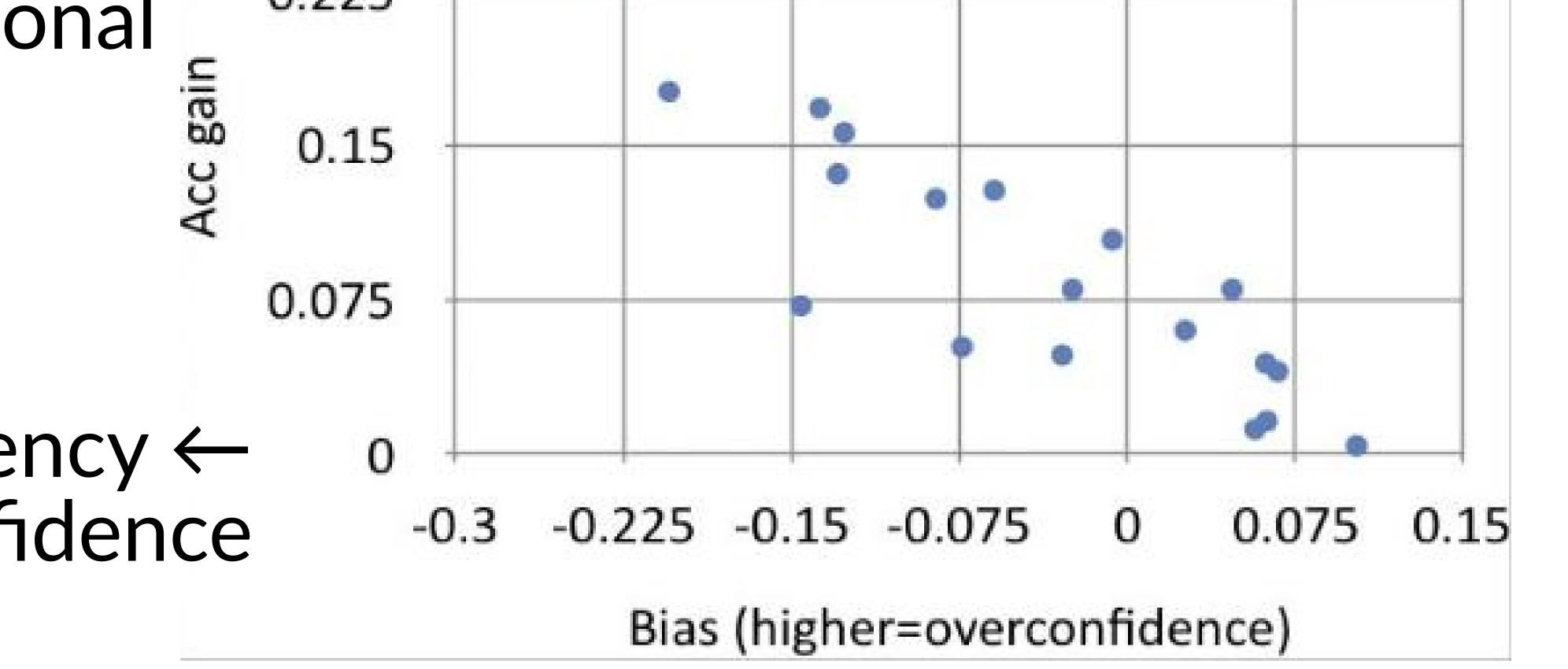


Overconfidence can be alleviated by sampling with high temperatures. However, there are trade-offs of worsening performance at high temperatures and much larger number of samples needed









Accuracy gain from self-consistency ← decreases with overconfidence