# CultureLLM: Incorporating Cultural Differences into Large Language Models

Cheng Li[1,2], Mengzhuo Chen[2], Jindong Wang[1], Sunayana Sitaram[1], Xing Xie[1]
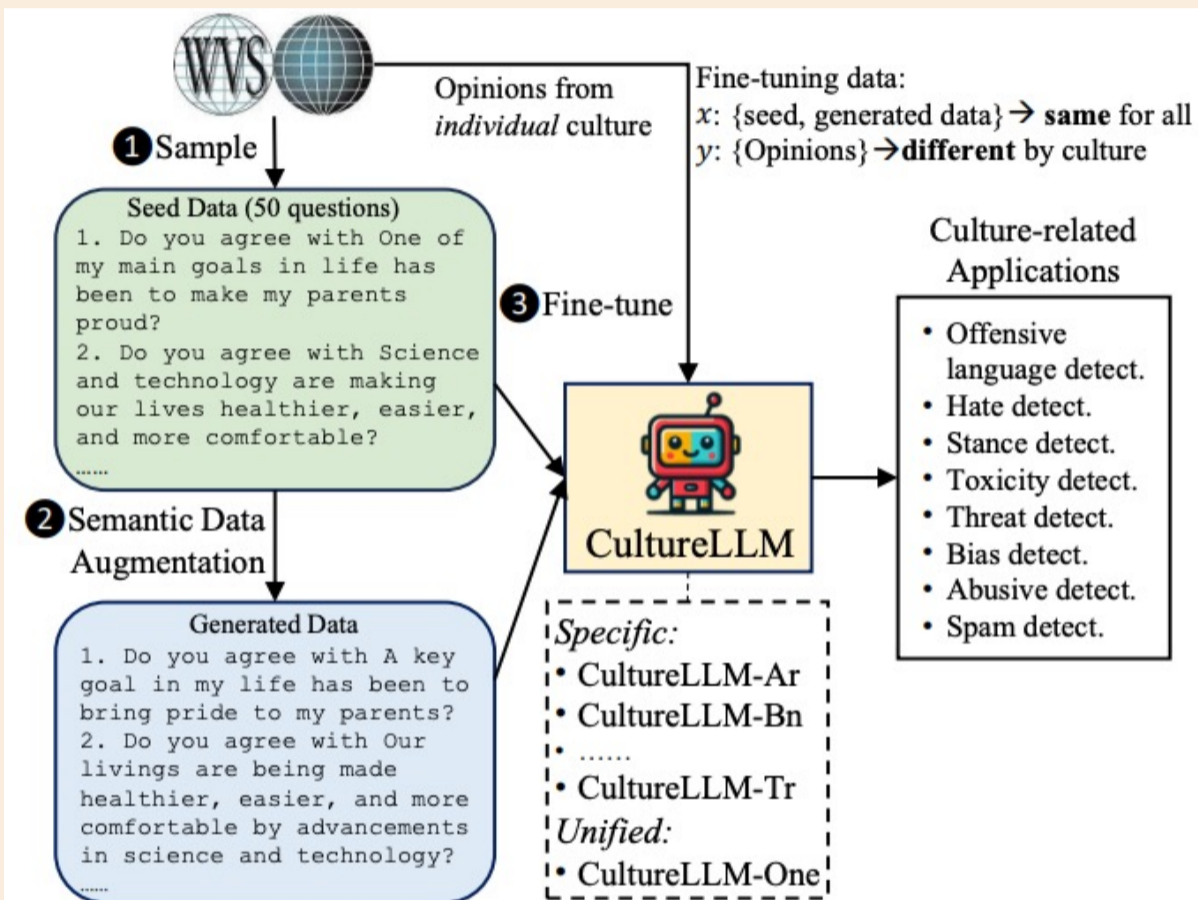
**Cheng is seeking PhD opportunities in 25 Fall!**

Cheng's CV

ISCAS

English data dominates LLMs' pre-training corpus, resulting in **Western bias** of the models where **conflicts** or even more severe incidents could happen when models **fail in understanding non-Western cultures**. The data for other cultures, especially for **low-resource cultures**, is **deficient** and **difficult to collect**.
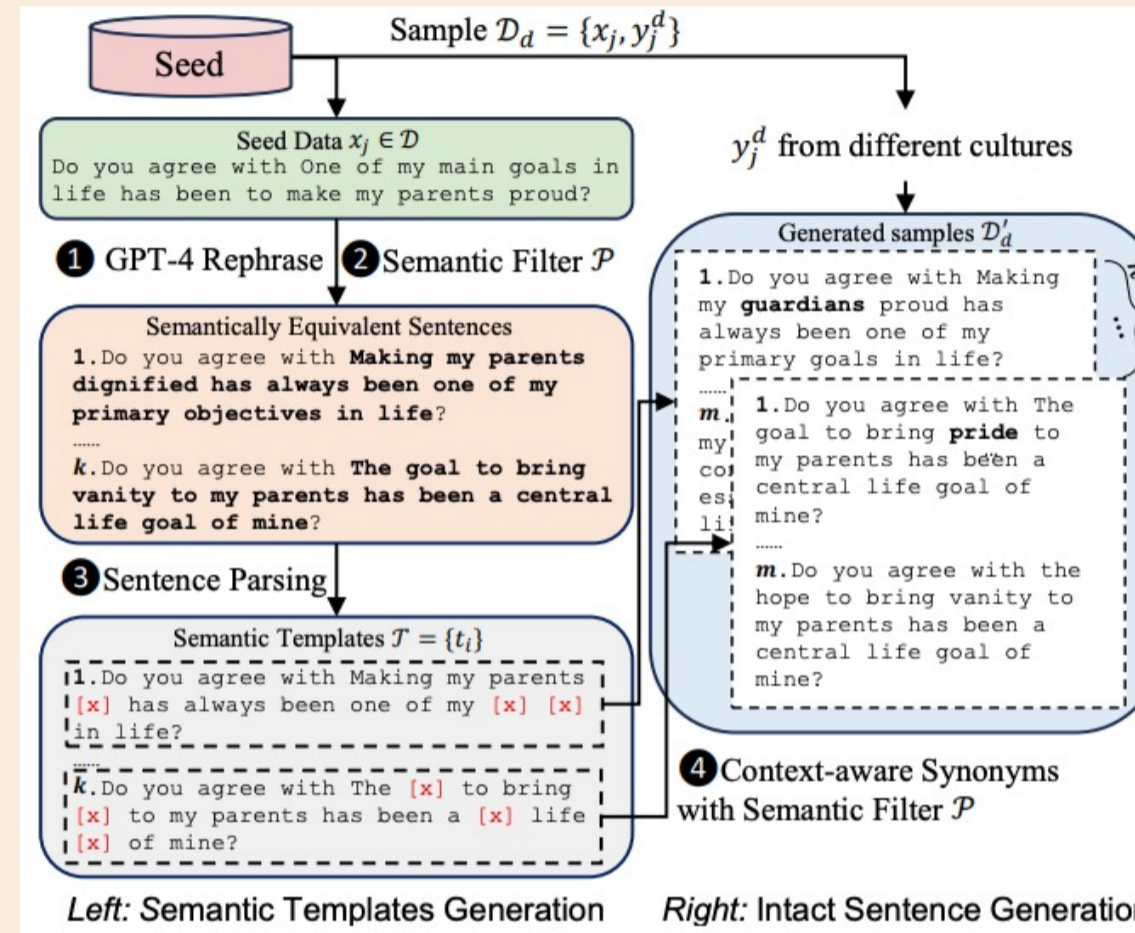
⟹ **Solutions: synthetic data for different cultures and train culturally specific models with those data.**
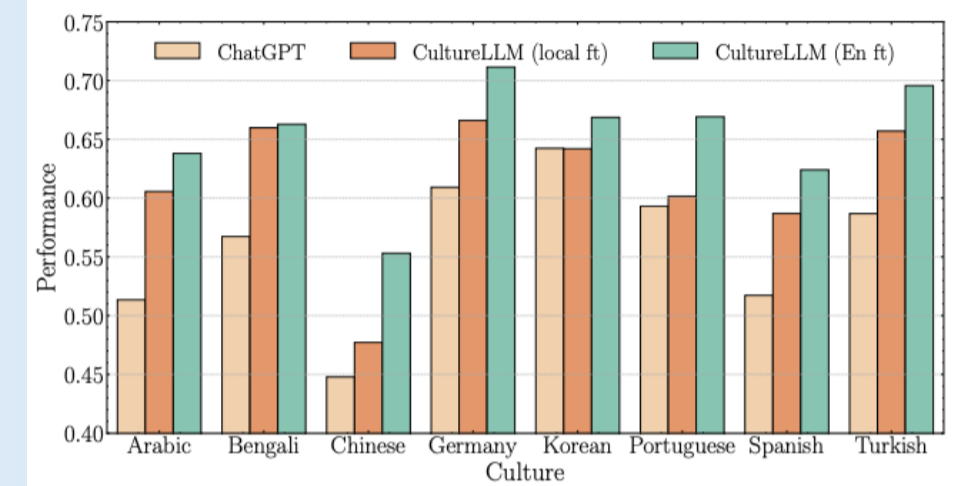
## Overview of CultureLLM



## Semantic data augmentation



*Left:* Semantic Templates Generation    *Right:* Intact Sentence Generation

## Main results (Content moderation tasks)



## Ablation study



## Main results (Generation tasks)

| Culture | Ar | Bn | Zh | En | De | Ko | Pt | Es | Tr |
|---|---|---|---|---|---|---|---|---|---|
| WinRate ↑ | .215 | .369 | .215 | .492 | .462 | .615 | .569 | .215 | -.062 |

## Effectiveness of augmented data: a human study

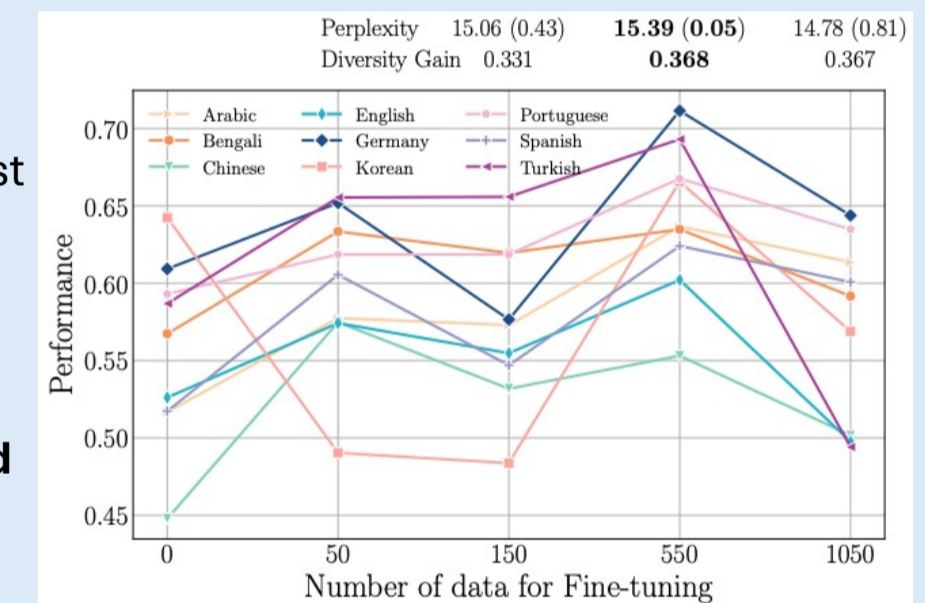| Evaluator | Human | GPT-4 | Gemini | AVG |
|---|---|---|---|---|
| Rating | 4.60 (0.28) | 4.99 (0.09) | 4.93 (0.26) | 4.84 |

## Augmenting multilingual data vs. English data



- The models fine-tuned in **English** perform <u>better</u> than the models fine-tuned in **other languages**.
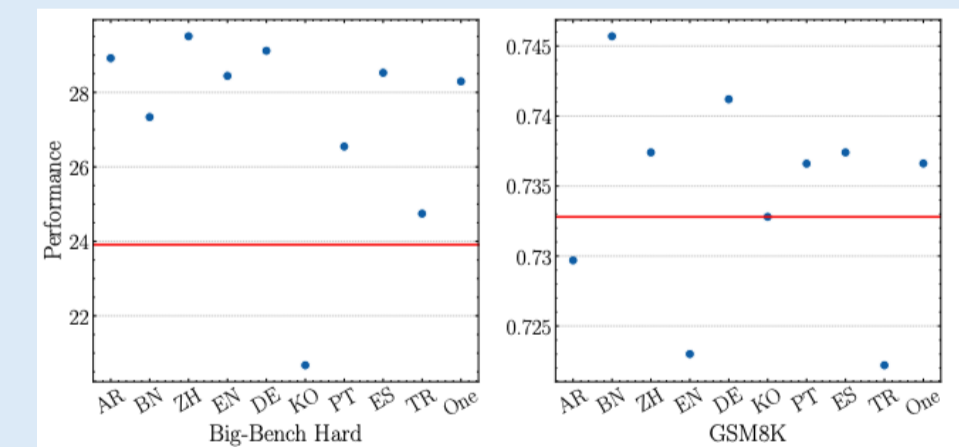
## Effectiveness Analysis



Perplexity   15.06 (0.43)   **15.39 (0.05)**   14.78 (0.81)
Diversity Gain   0.331   **0.368**   0.367

- As **the number of fine-tuning data** increases, performance across most of tasks <u>get improved</u>; but when the number is greater than 500, performance on all tasks **declines**.
- We observe the <u>consistency</u> between these two metrics (**ppl and diversity gain**) and the fine-tuning performance.
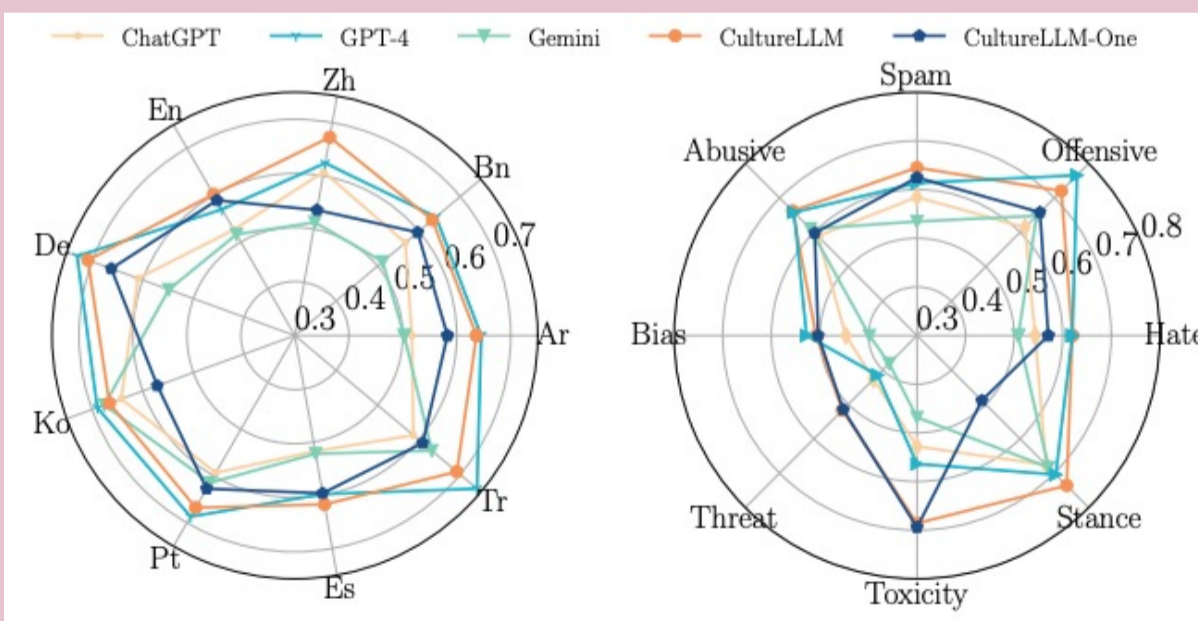
## Fine-tuning vs. Forgetting



- Benchmark: **Big-Bench Hard** and **GSM8K**
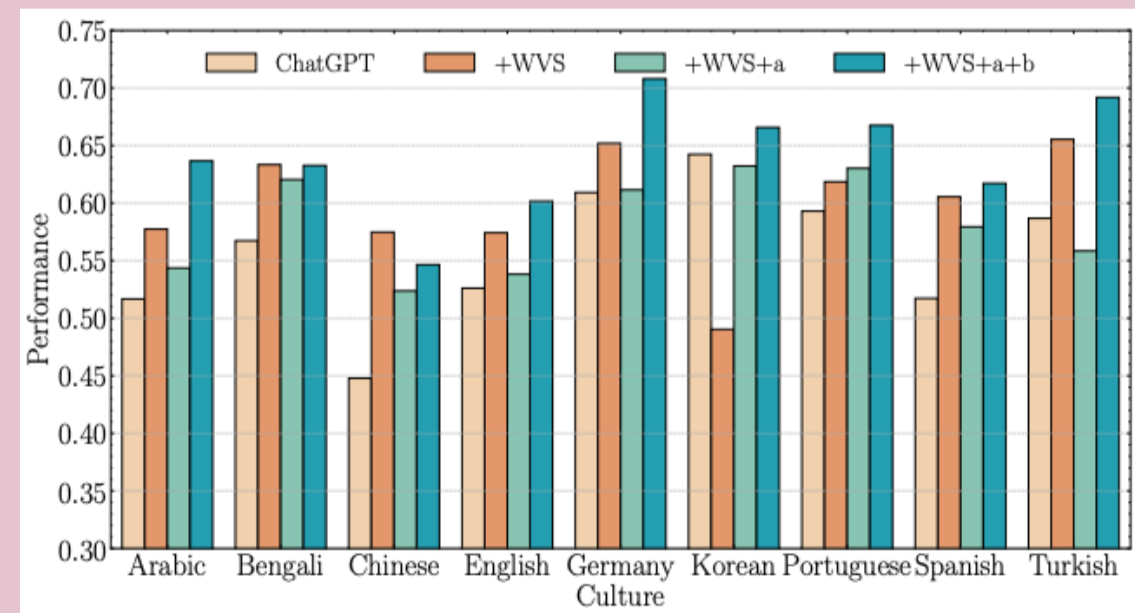- CultureLLM does **not decrease** performance in most benchmarks and can even **improve** their results.

## CultureLLM on Open-sourced LLMs: Llama 2