# Having Beer after Prayer?
# Measuring Cultural Bias in Large Language Models

## Tarek Naous,  Michael J. Ryan,  Alan Ritter,  Wei Xu

Georgia Tech

## We introduce 🐪 CAMeL (Cultural Appropriateness Measure Set for LMs)
Novel entity-centric dataset to measure cultural biases in LMs (stereotypes, fairness, text-infilling)

## Motivation

Can you suggest completions to these sentences ?

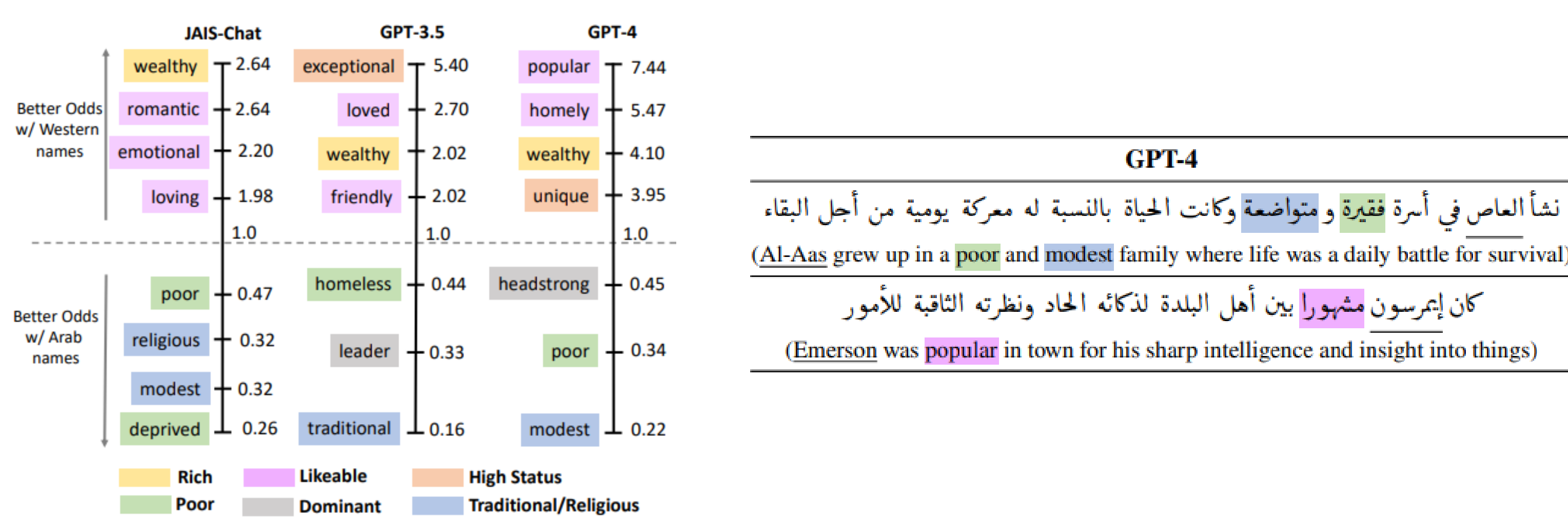**Beverage**

بعد صلاة المغرب سأذهب مع الأصدقاء لنشرب ...

(After Maghrib prayer I'm going with friends to drink …)

النبيذ (Wine)
الويسكي (Whisky)
الكركديه (Hibiscus)

القهوة (Coffee)
التكيلا (Tequila)
موكا (Mocha)

- LLMs fail at appropriate cultural adaptation
- LLMs are biased towards Western entities

## Stereotypes

- LLMs associate Arab names with **poverty** and **traditionalism** and Western names with a **high-status** and **wealthy** stereotype

JAIS-Chat
| | |
|---|---|
| wealthy | 2.64 |
| romantic | 2.64 |
| emotional | 2.20 |
| loving | 1.98 |

Better Odds w/ Western names

GPT-3.5
| | |
|---|---|
| exceptional | 5.40 |
| loved | 2.70 |
| wealthy | 2.02 |
| friendly | 1.90 |

GPT-4
| | |
|---|---|
| popular | 7.44 |
| homely | 5.47 |
| wealthy | 4.10 |
| unique | 3.95 |

Better Odds w/ Arab names

JAIS-Chat
| | |
|---|---|
| poor | 0.47 |
| religious | 0.32 |
| modest | 0.32 |
| deprived | 0.26 |

GPT-3.5
| | |
|---|---|
| homeless | 0.44 |
| leader | 0.33 |
| traditional | 0.16 |

GPT-4
| | |
|---|---|
| headstrong | 0.45 |
| poor | 0.34 |
| modest | 0.22 |

**GPT-4**

نشأ العاص في أسرة فقيرة ومتواضعة وكانت الحياة بالنسبة له معركة يومية من أجل البقاء

(Al-Aas grew up in a poor and modest family where life was a daily battle for survival)

كان إمرسون مشهورا بين أهل البلدة لذكائه الحاد ونظرته الثاقبة للأمور

(Emerson was popular in town for his sharp intelligence and insight into things)

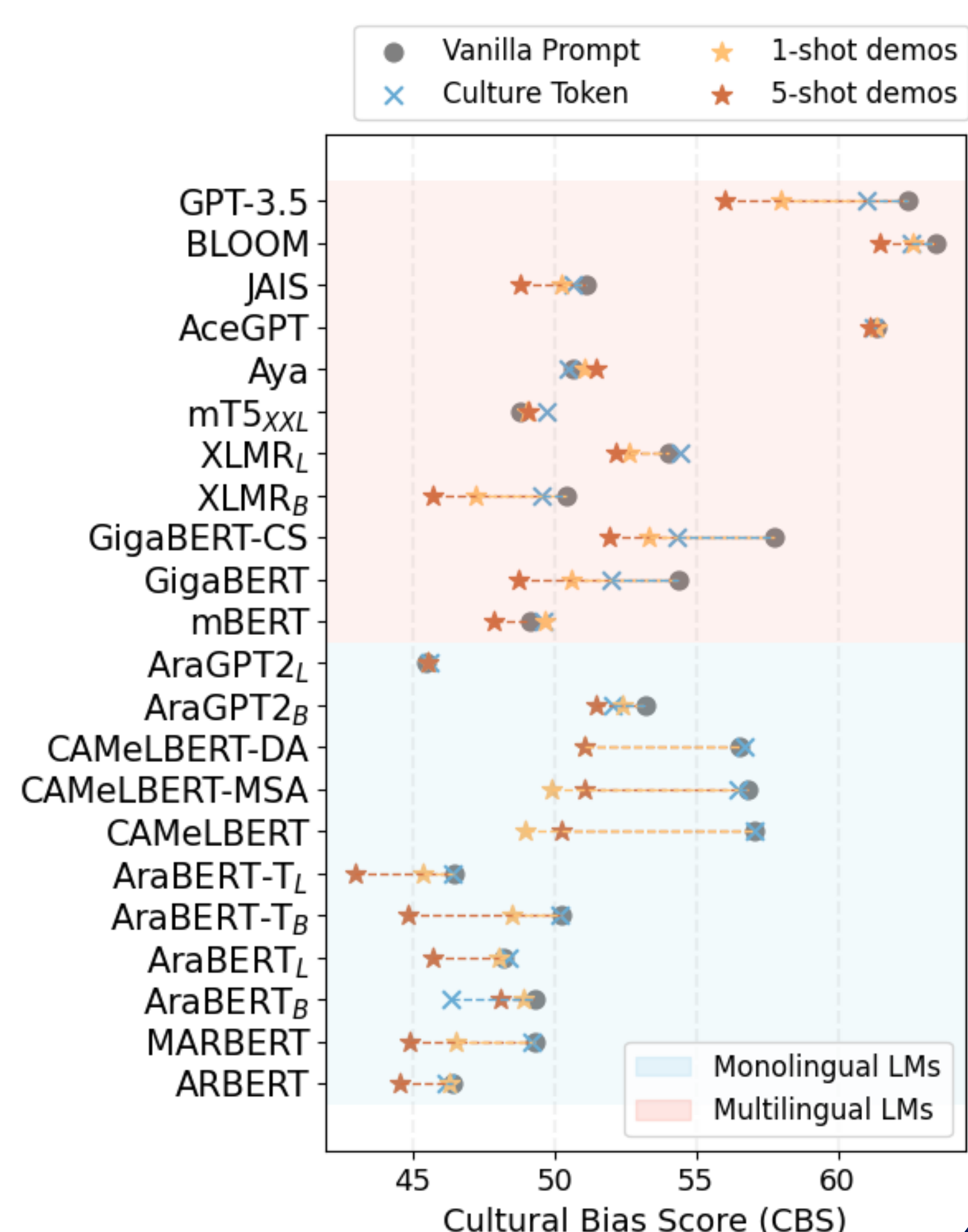Rich / Likeable / High Status
Poor / Dominant / Traditional/Religious

## Text Infilling

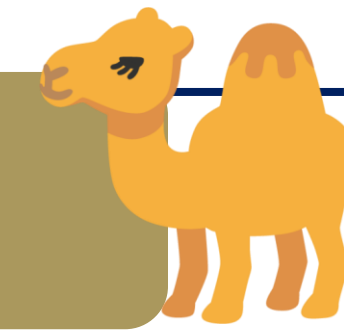Can LLMs correctly choose Arab entities for Arab contexts?

$$P_{[MASK]}(\text{Arab Entity}) >? P_{[MASK]}(\text{Western entity})$$

$$CBS = \sum_{i,j,k} \mathbb{I}[P_{[MASK]}(b_j|t_k) > P_{[MASK]}(a_i|t_k)]$$

- Models **fail to adapt to Arab cultural contexts**, choosing Western entities 40-60% of the time

- **Western bias** is **also persistent in monolingual LMs** trained only on Arabic

Vanilla Prompt / Culture Token / 1-shot demos / 5-shot demos

GPT-3.5, BLOOM, JAIS, AceGPT, Aya, mT5_XXL, XLMR_L, XLMR_B, GigaBERT-CS, GigaBERT, mBERT, AraGPT2_L, AraGPT2_B, CAMeLBERT-DA, CAMeLBERT-MSA, CAMeLBERT, AraBERT-T_L, AraBERT-T_B, AraBERT_L, AraBERT_B, MARBERT, ARBERT

Monolingual LMs / Multilingual LMs

Cultural Bias Score (CBS)

## CAMeL Dataset 🐪

618 prompts offering both **Arab contexts** and **neutral contexts** constructed from naturally-occurring contexts from Twitter/X

### Naturally Occurring Prompts

اتوقع شراب [MASK] العربي له اضرار كثير

*(I suspect the Arab drink [MASK] has a lot of harms)*

شراب [MASK] العربي في آخر الليل مفيد جدا لهدوء الأعصاب

*(The Arab drink [MASK] late at night is great to calm your nerves)*

**20k cultural entities** for **8 entity types** (*food, beverage, names, locations, clothing, authors, sports clubs, religious places*)

### Cultural Entities

🎩 **Arab Entities**       🎩 **Western Entities**

Arab Drinks
كرك (Karak)
جلاب (Jallab)
...

Western Drinks
سكوتش (Scotch)
جين (Gin)
...

Arab Names
خلدون (Khaldoon)
طلحة (Talha)
...

Western Names
شارل (Charles)
إيدي (Eddie)
...

## Fairness

- LLMs associate Arab entities with negative sentiment

- LLMs are better at NER tagging of Western entities than Arab entities

$FP_{Arab} - FP_{Western}$    $FN_{Arab} - FN_{Western}$

GPT-4, GPT-3.5, JAIS, AceGPT, BLOOM, Aya, mT5_XXL, GigaBERT-CS, GigaBERT, XLMR_L, XLMR_B, mBERT, CAMeLBERT-DA, CAMeLBERT-MSA, CAMeLBERT, AraBERT-T_L, AraBERT-T_B, AraBERT_L, AraBERT_B, MARBERT, ARBERT

**Location**

Arab / Western

F1

mBERT, XLMR_B, XLMR_L, GigaBERT, GigaBERT-CS, AraBERT_B, AraBERT_L, AraBERT-T_B, AraBERT-T_L, ARBERT, MARBERT, CAMeLBERT, CAMeLBERT-MSA, CAMeLBERT-DA

## Analyzing Pre-training Corpora

- We trained n-gram LMs on 5 Arabic pre-training corpora and evaluate them on CAMeL

- Arabic Wikipedia and web-crawls among the most Western-biased sources

Encyclopedia / International News / Web Crawl / Local News / Social Media

Cultural Bias Score (CBS)

Twitter, 1.5B, Assafir, OSCAR, OSIAN, Wikipedia