# On The Origin of Cultural Biases in Language Models:
# From Pre-training Data to Linguistic Phenomena

**Tarek Naous** and **Wei Xu**
Georgia Institute of Technology
tareknaous@gatech.edu; wei.xu@cc.gatech.edu

## Abstract

Language Models (LMs) have been shown to exhibit a strong preference towards entities associated with Western culture when operating in non-Western languages. In this paper, we aim to uncover the origins of entity-related cultural biases in LMs by analyzing several contributing factors, including the representation of entities in pre-training data and the impact of variations in linguistic phenomena across languages. We introduce CAMeL-2, a parallel Arabic-English benchmark of 58,086 entities associated with Arab and Western cultures and 367 masked natural contexts for entities. Our evaluations using CAMeL-2 reveal reduced performance gaps between cultures by LMs when tested in English compared to Arabic. We find that LMs struggle in Arabic with entities that appear at high frequencies in pre-training, where entities can hold multiple word senses. This also extends to entities that exhibit high lexical overlap with languages that are not Arabic but use the Arabic script. Further, we show how frequency-based tokenization leads to this issue in LMs, which gets worse with larger Arabic vocabularies. We will make CAMeL-2 available at: https://github.com/tareknaous/camel2.

## 1 Introduction

Multilingual Language Models (LMs) are playing a crucial role in making AI technology accessible to global communities (Üstün et al., 2024; Singh et al., 2024). As these communities represent diverse cultural backgrounds, the multilingual challenge for LMs does not merely stop at handling different languages (Blevins et al., 2024), but extends to capturing cultural nuances (Adilazuarda et al., 2024). However, past research has highlighted strong favoritism in LMs towards entities associated with Western culture when operating in non-Western languages, leading to a struggle by LMs to adapt to cultural contexts and gaps in their performance between cultures on NLP tasks (Naous et al., 2024).
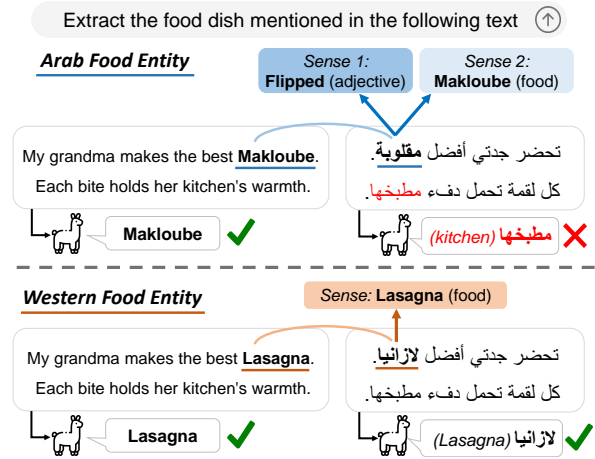


Figure 1: Responses of a LM (🦙) tasked to extract the food dish from the same text in English and Arabic. The LM identifies the Arab dish *"Makloube"* in English, but fails in Arabic where the word *"Makloube"* holds two senses. The LM does not struggle with the Western dish *"Lasagna"* which holds only one sense in both languages.

While entity-related biases in LMs have been traditionally studied as a reflection of imbalanced representations in pre-training data (Gallegos et al., 2024; Li et al., 2024a), it is often overlooked how linguistic phenomena in non-English languages can also incite those biases. For example, when an LM is asked to extract an Arab food dish such as *"Makloube"* from text, it can fail to do so in Arabic where the word used for the dish holds two senses (as the food dish or as the adjective *"flipped"*), but can successfully extract the same dish from the parallel English text, as shown in Figure 1. Yet, within the same context, the LM does not face this struggle when we replace *"Makloube"* with a Western dish *"Lasagna"*, which holds only one sense of a food dish both in English and in Arabic. These observations raise the question: **do varying linguistic phenomena exhibited by cultural entities influence cultural biases in LMs?**

Our study investigates the impact of such cross-linguistic differences to uncover the origins of

| Food Contexts | Location Contexts |
|---|---|
| **Text-Infilling & NER** - CAMeL (Naous et al., 2024) | |
| الغدا عربي اليوم عملت [MASK] كتير طيبة | انا منذ ايام كنت في مدينة [MASK] العربية و هي في غاية الروعة |
| (Today's lunch is Arab, I've cooked [MASK] which is very delicious) | (I was in the Arab city of [MASK] a few days ago and it is incredibly wonderful) |
| **Extractive QA** - CAMeL-2 (this work) | |
| قررت ادخل اسوي [MASK] بإشراف امي لانها معاندة إلا وتطبخ بنفسها لان ولدها المفضل بخاطره مالح وهي تعبانه وانا ماريدها تتعب اكثر | استقبل الشيخ بهاء مساء اليوم وفدا من أهالي [MASK] حيث تم عرض مشاكل وسبل معالجتها. من جهته جدد الشيخ تعهده بحل هذه المشاكل والبدء بنهضة جديدة |
| (I decided to go in and make [MASK] under the supervision of my mother, because she insists to cook only by herself since her favorite son has salty taste, and she is tired, and I don't want her to get more tired) | (Sheikh Bahaa received this evening a delegation of people from [MASK] where problems and ways to address them in were presented. For his part, the Sheikh renewed his pledge to solving these problems and starting a new renaissance) |

Table 1: Example food and location contexts collected in CAMeL-2 for Extractive QA evaluation, compared with contexts from CAMeL (Naous et al., 2024). Extractive QA contexts are longer and mention entities more implicitly.

entity-related biases in LMs. To enable our analyses, we introduce CAMeL-2, a parallel Arabic-English resource of 58,086 entities associated with Arab and Western cultures across seven entity types, and a set of 367 natural contexts for these entities (§2). Using CAMeL-2, we evaluate a variety of LMs on extractive QA and NER, revealing smaller performance gaps between cultures when LMs are tested in English compared to Arabic (§3).

Our analyses show that LMs struggle at recognizing entities associated with Arab culture which appear at very high frequencies in Arabic pre-training data (§4.1), where such entities exhibit strong word polysemy in Arabic (§4.2). We also find that high lexical overlap of Arab entities with pre-training data of languages that use Arabic script (e.g., Farsi, Urdu, etc.) causes drops in performance (§4.3). Lastly, we show how tokenization causes LMs to struggle with Arab entities when tokenized into a single token, an issue that worsens with larger Arabic vocabularies (§4.4).

## 2 CAMeL-2: A Parallel Benchmark

Our goal is to investigate whether LMs handle entities associated with Arab and Western cultures differently when tested in Arabic *vs.* English languages (§3). To do this, we first construct a *bilingual* version of the entity-centric CAMeL (Naous et al., 2024) benchmark for measuring cultural biases in LMs and extend its coverage by three times (§2.2). Table 2 compares the statistics of CAMeL-2 vs. CAMeL. We also collect longer contexts where these entities are mentioned for evaluating LMs in a more challenging setup of extractive QA (§2.3).

### 2.1 About CAMeL

The original CAMeL benchmark (Naous et al., 2024) consists of 20,249 cultural entities extracted from Wikidata and web-crawl data and annotated

| Entity Type | CAMeL | CAMeL-2 | Increase |
|---|---|---|---|
| Authors | 571 | 6,315 | 11.05× |
| Beverage | 142 | 255 | 1.79× |
| Food | 578 | 2,283 | 3.94× |
| Locations | 12,497 | 35,200 | 2.81× |
| Names | 1,533 | 3,842 | 2.50× |
| Religious | 2,428 | 5,049 | 2.07× |
| Sports Clubs | 2,500 | 5,142 | 2.05× |
| **Total** | 20,249 | 58,086 | 2.86× |

Table 2: Number of entities in the bilingual CAMeL-2 (this work) vs. the monolingual CAMeL (Naous et al., 2024). We increase the size of the benchmark by 2.86×.

with Arab or Western cultural association. CAMeL also includes a set of textual contexts where these entities may naturally occur (see examples in Table 1), derived from Arabic X/Twitter data. All entities and contexts in CAMeL are written in the Arabic language *only*, limiting the ability to test the behavior of LMs when being prompted in English. To enable such comparisons, we construct a parallel extension to CAMeL by not only adding the English translation of each entity and context, but also increasing the overall number of entities and length of contexts. We find that Wikipedia contains more entities relevant to Arab culture than Wikidata for *authors*, *beverage*, *food*, *names*, *religious places*, and *sports clubs* (§2.2). We expand the coverage of *location* entities using public geographic data.

### 2.2 Collecting Cultural Entities

**Entity Extraction from Wikipedia.** We leverage the categorization feature in Wikipedia and identify, for each entity type, a generic category that is repeatedly used to group together articles relevant to a specific country. For example, the "[country adjective] cuisine" category encompasses all food-related articles of a country (e.g., *Syrian* cuisine, *Irish* cuisine, etc.). We extract the titles of all articles associated with each country.

The entities are then obtained from the article titles, which in most cases are direct references to the entity of interest. Additionally, we extract the body of text for each article which we use for distantly supervised fine-tuning of NER models (§3.3). See Appendix A for details and the categories used.

**Annotation of Cultural Entities.** The extracted entities from Wikipedia were then manually annotated for cultural association. We hired two college students who are native Arabic speakers to classify entities into: *Arab culture* (Arab countries), *Western culture* (European and North American countries), *other cultures*, or *not culture-specific*. The annotator agreement is 0.825 by Cohen's Kappa. The cases of disagreement were discussed in an adjudication step to decide on the final label.

**Mapping Arabic Entities to English.** Since both Wikidata and Wikipedia support multiple languages, we automatically map the Arabic entities to their English versions, which were available for 88.34% of entities in CAMel-2. For the remainder of the entities which were only available in Arabic, we manually search for their commonly used English transliterations on the Internet. For example, a Tunisian sports club "المظيلة" can be transliterated as "*Mzilla*" or "*Mdhilla*", both of which are valid. However, the form used by Tunisians is "*Mdhilla*" as it aligns with their phonetic interpretation of Arabic letters in the Tunisian dialect.

**Georgraphic Data-based Extraction.** The Arab *location* entities extracted from Wikidata in CAMeL had relatively low representation at 1,061 locations compared to 11,436 Western locations. We extract additional locations for all Arab countries from the OpenStreetMap[1] (OSM) database, which provides Arabic-English pairs of locations in Arab countries.[2] This resulted in an extensive set of 23,765 Arab locations, enabling further analyses across Arab regions that were influenced by other languages in the history (§4.2).

### 2.3 Constructing Contexts for Extractive QA

CAMeL provides 250 masked contexts where only entities associated with Arab culture are appropriate [MASK] token fillings. This allows evaluation of LMs on adaptation to cultural contexts and on NER of entities with different cultural associations

---

[1] https://www.openstreetmap.org/
[2] We determine Arab countries based on the league of Arab states: https://arabmpi.org/en/home

(Naous et al., 2024), but *only* for Arabic language. Further, these contexts are short and explicitly refer to the masked entity, making them less suitable for evaluating GPT-type LMs on tasks such as extractive QA, which aligns better with their usage. To address this limitation, we collect 117 new, longer contexts from the X/Twitter platform where entities are mentioned more implicitly. This presents a challenging evaluation setup for LMs where understanding of context is necessary for extracting the entity (see comparative examples in Table 1).

For each entity type, we perform keyword search using 30 randomly sampled Arabic entities to capture natural discussions about entities. We search over the two months period of 7/1/2024 to 9/1/2024 and manually inspect tweets to identify ones that are long and make an indirect reference to the entity. From these, we construct 10 to 20 extractive QA contexts for each entity type by replacing the user-mentioned entities by a [MASK]. To enable comparative Arabic-English evaluations, we manually translate each culturally-grounded context from the original CAMeL and the newly collected QA contexts from Arabic to English.

## 3 Is Western Bias in LMs Consistent Across Arabic and English Languages?

We start by studying whether LMs show the same degree of favoritism for Western culture, when operating in Arabic vs. English languages. Our analyses focus on two aspects: cultural adaptation in text-infilling (§3.2) and cross-cultural performance on extractive QA and NER (§3.3).

### 3.1 Language Models

We experiment with LMs that have been trained on both Arabic and English. We use the recent multilingual LMs of **Llama-3.3** (Dubey et al., 2024), **Aya-23** (Aryabumi et al., 2024), **Qwen-2.5** (Yang et al., 2024), as well as **AceGPTv1.5** (Zhu et al., 2024) which expands the Arabic vocabulary of Llama-2 and further fine-tunes it on Arabic data, **AceGPTv2** (Liang et al., 2024) which adapts Llama-3 checkpoints on Arabic, and the Arabic-English bilingual model **JAIS** (Sengupta et al., 2023). We also experiment with encoder LMs such as **XLM-R** (Conneau et al., 2020), as well as Arabic monolingual encoders, including **ARBERT** and **MARBERT** (Abdul-Mageed et al., 2020), **CAMeLBERT** (Inoue et al., 2021), and **AraBERT** (Antoun et al., 2020).
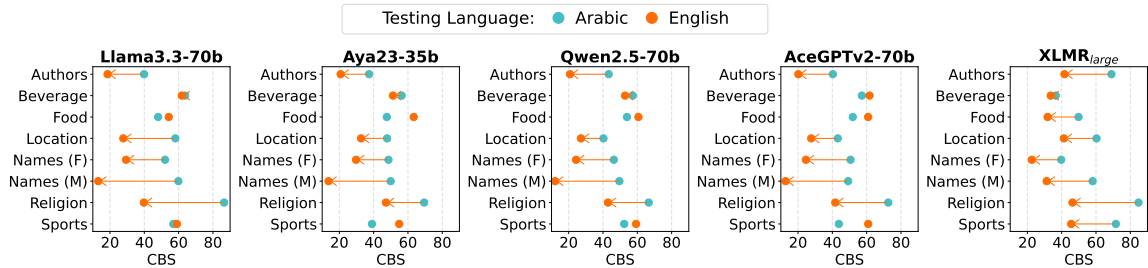
Figure 2: **C**ultural **B**ias **S**core (↓) (§3.2) per entity type on culturally-grounded contexts from CAMeL-2. LMs can adapt better to Arab culture when tested in English.

| | Llama3.3-70b | | | | | | XLMR$_{large}$ | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | **Arabic** | | | **English** | | | **Arabic** | | | **English** | | |
| | *Arab* | *Western* | ΔAcc | *Arab* | *Western* | ΔAcc | *Arab* | *Western* | ΔF1 | *Arab* | *Western* | ΔF1 |
| Authors | 92.62 | 90.28 | -2.34 | 98.99 | 99.16 | 0.17 | 86.80 | **87.93** | 1.13 | 95.64 | 94.98 | -0.66 |
| Beverage | 82.65 | 78.19 | -4.46 | 99.14 | 97.71 | -1.43 | 63.06 | **72.86** | 9.80 | 92.06 | 89.77 | -3.29 |
| Food | 84.08 | **84.71** | 0.63 | 95.84 | 98.21 | 2.37 | 63.76 | **73.59** | 9.83 | 91.57 | 90.45 | -1.12 |
| Location | 80.66 | **95.59** | 14.93 | 98.58 | 99.89 | 1.31 | 64.07 | **91.32** | 27.25 | 89.54 | 95.71 | 6.17 |
| Names (F) | 63.38 | **77.39** | 14.01 | 99.86 | 99.14 | -0.72 | 62.22 | **82.65** | 20.43 | 97.87 | 96.36 | -1.51 |
| Names (M) | 75.45 | **76.23** | 0.78 | 99.43 | 99.78 | 0.35 | 80.09 | **85.03** | 4.94 | 94.01 | 93.13 | -0.88 |
| Sports | 68.58 | **79.01** | 10.43 | 92.77 | 96.02 | 3.25 | 74.52 | **84.14** | 9.62 | 92.14 | 93.12 | 0.97 |
| Religious | 82.49 | **82.98** | 0.49 | 98.52 | 97.69 | -0.83 | 95.30 | **97.13** | 1.83 | 94.34 | 95.76 | 1.42 |

Table 3: Average performance of Llama3.3-70b (QA Accuracy ↑) and XLMR$_{large}$ (NER F1 ↑) on Arab and Western entities when tested in Arabic and English. More results with Aya23-35b and AceGPTv2-70b are in Appendix. ΔAcc and ΔF1 are performance differences between Western and Arab entities. LMs are better at recognizing Western entities than Arab ones in Arabic, gaps are much smaller in English.

## 3.2 Cultural Adaptation: Text Infilling

We compare the ability of LMs at adapting to cultural contexts in Arabic and English by analyzing their preference for Arab vs. Western entities as [MASK] token fillings of CAMeL-2 contexts.

**Text Infilling Setup.** We use the **C**ultural **B**ias **S**core (CBS) designed by Naous et al. (2024) as a likelihood-based measure of a LM's ability at adaptating to cultural contexts. Consider the sets of Arab entities $A = \{a_i\}_{i=1}^N$ and Western entities $B = \{b_j\}_{j=1}^M$. The CBS for an Arab entity $a_i$ is the percentage of a LM's preference for Western entities when placed within the same context. For a set of culturally-grounded masked contexts $C = \{c_k\}_{k=1}^K$, we compute the $\text{CBS}(a_i)$ as:

$$\frac{1}{K \times M} \sum_{k=1}^K \sum_{j=1}^M \mathbb{1}[P_{\texttt{[MASK]}}(b_j|t_k) > P_{\texttt{[MASK]}}(a_i|t_k)]$$

where $P_{\texttt{[MASK]}}$ is the LM's probability of an entity filling the masked token. As the contexts are grounded in Arab culture (i.e., only Arab entities are appropriate), LMs are expected to score a CBS closer to 0%. A higher CBS score indicates a stronger preference by LMs for Western entities

in place of Arab entities, given the same context. For all entity types, we compute the CBS for a random sample of 50 Arab entities, where we test each against 50 randomly sampled Western entities.

**Results.** Figure 2 shows the average CBS per entity type for several LMs. Interestingly, **LMs are better at adapting to Arab cultural contexts in English than in Arabic** with reduced CBS levels in nearly all cases, reaching the 15-30% range for names and locations. One exception is food and beverage categories where a struggle is visible in both languages. We analyze the potential reasons behind these results in §4.

## 3.3 Cross-Cultural Performance

Utilizing CAMeL-2, we compare the performance of LMs when tested in Arabic vs. English parallel contexts for extractive QA and NER tasks.

**Prompting Setup for Extractive QA.** We evaluate GPT-type models in an extractive QA setup by prompting the LMs to extract the cultural entity from a given context in CAMeL-2 (see Appendix B.2 for prompts). We create an Arabic test set for each entity by replacing the [MASK] token of each context with the entity (i.e., ~15 test con-
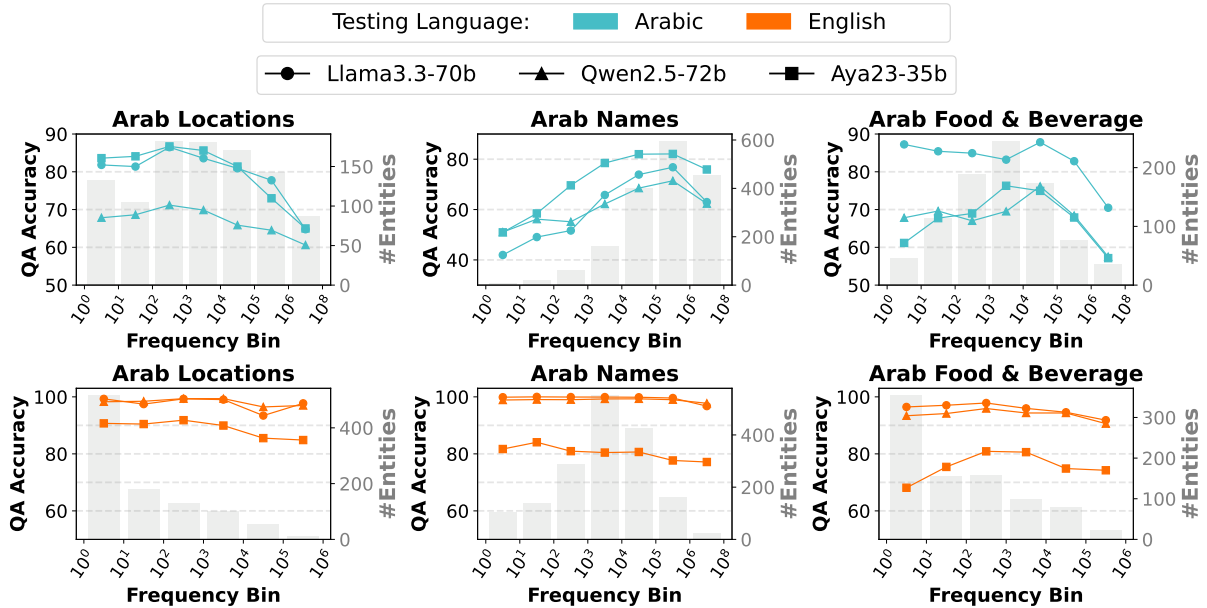
Figure 3: Average QA Accuracy (↑) of LLMs when tested in Arabic and English on location, name, food, and beverage associated with Arab culture, stratified by their occurrence counts in the mC4 corpus (§4.1; grouped into log10-spaced bins). Gray bars in background represent number of entities tested in each bin. Interestingly, LMs struggle with very high-frequency entities in Arabic.

texts per entity), as well as a corresponding English test set with the same entities and contexts that are translations (§2.2). We evaluate random samples of 1k Arab and 1k Western entities for each type and use all entities for types with less than 1k entities.

**Fine-tuning Setup for NER.** We similarly evaluate BERT-type models on the NER task using the culturally-grounded contexts from CAMeL-2. We fine-tune models capable of recognizing *names*, *authors*, and *locations* using the ANERCorp dataset for Arabic (Benajiba et al., 2007) and CoNLL-2003 for English (Sang and De Meulder, 2003). For the remaining entity types (i.e., *food*, *beverage*, *sports clubs*, *religious places*) that are not covered by existing manually annotated NER corpora, we fine-tune LMs via distant supervision from Wikipedia articles (Liang et al., 2020) that we collected in §2.2 (see Appendix B.3 for details). We exclude entities that appear in fine-tuning from our evaluations.

**Results.** Table 3 shows the average performance of Llama3.3-70b and XLMR$_{large}$ on Arab and Western entities. We also show the performance difference ($\Delta$) between Western and Arab entities, where a positive $\Delta$ indicates better performance on Western entities. We find that **LMs have small performance gaps between the two cultures in English, but are consistently better at recognizing Western entities in Arabic than in English,**

where differences reach up to 27 F1 points in NER and 15% accuracy in extractive QA.

## 4 On The Origin of Cultural Biases

Motivated by our observations in §3, we analyze various factors that may cause LMs to exhibit more severe Western bias in Arabic than in English. We first study the relationship between the frequency of entities in pre-training data and LM performance (§4.1). Our findings lead to analysis of the impact of Arabic word polysemy on LM biases (§4.2). We also look at scenarios where entities exhibit lexical overlap with other languages that use the Arabic script (§4.3). Finally, we examine the role that tokenization plays in the observed issues (§4.4).

### 4.1 Entity Frequency in Pre-training Data

We examine how LM performance on entities in extractive QA (§3.3) changes with respect to how often entities occur in pre-training. Since the pre-training corpora of state-of-the-art multilingual LMs are not public, we approximate the occurrence of entities in pre-training using the Arabic portion (0.96B lines) and English portion (3.1B lines) of the mC4 corpus[3] (Xue et al., 2021). Figure 3 shows the average QA accuracy achieved by several LMs

---

[3] mC4 is a multilingual partition of CommonCrawl web scrapes which are an essential part of LM pre-training data. Partitioning was done using the CLD3 language detector.

on Arab location, name, food, and beverage entities versus their occurrence count, which we group into log10-spaced bins (i.e., entities that appear 1 to 10 times, 10 to 100 times, etc.). We observe the following key findings:

**LMs struggle on high-frequency entities.** There is a noticeable drop in LM performance on entities that appear at very high frequencies ($>$1M times in bin '$10^6$–$10^8$'). This drop is much steeper when LMs are tested in Arabic than English, and more prevalent since more Arab culture-associated entities appear at extremely high frequencies in Arabic. Similar trends are observed for other entity types and on the text-infilling task (see Appendix D.1). Upon inspection of those entities, we find that many are Arabic words that can hold multiple senses in different contexts, besides being used to represent entities. We explore this impact of word polysemy more closely in §4.2.

**LMs also struggle with long-tail entities.** We find that LMs can also perform poorly on long-tail entities which appear at low frequencies ($<$10 times in bin '$10^0$–$10^1$'). This is especially noticeable on names and food entities in Arabic, where performance improves gradually as entities appear more frequently. In general, we find LMs to perform the best in both languages on entities that appear at a medium frequency (e.g., 1k-100k range), but face difficulty with the edge cases (low and high-frequency).

## 4.2 Entities as Polysemous Words

We further analyze how Arabic word polysemy impacts LM performance on entities that appear at high frequencies in pre-training corpora.

**Background.** Consider the word "مطروحة" (pronounced: /ma-troo-ha/) as it appears in two Arabic sentences and their English translations:

(1) القضية مطروحة للنقاش
*(The issue is proposed for discussion)*

(2) جدتي تسكن في مطروحة
*(My grandma lives in Matrooha)*

In (1), the word appears in its literal sense "*proposed*". However, in (2) the same word is used as a noun to denote the name of a location. This dual use is common in Arabic, where the functionality of words used for entities changes depending on
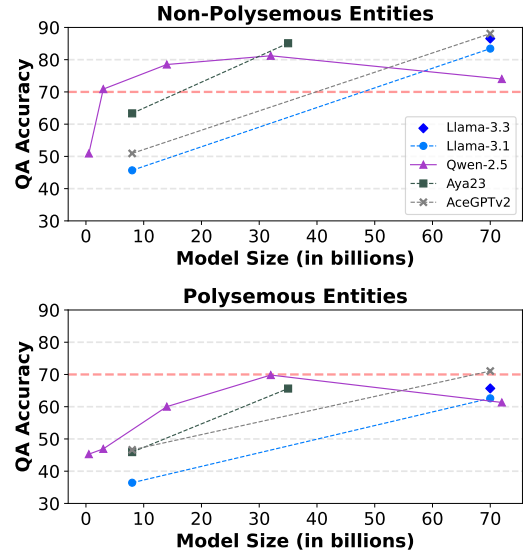


Figure 4: QA accuracy ($\uparrow$) for different sizes of LMs on high-frequency polysemous and non-polysemous Arab location entities. We find a positive scaling trend for all models, with lower performance on polysemous entities.

context. This is less common for entities in English, where there is generally a clear distinction between nouns and adjectives (Van Langendonck, 2007). We show this quantitatively in Appendix C.1.

**Setup.** We focus our analysis on location entities in Arab countries as they often exhibit regional linguistic influences. Thus, those entities can either be Arabic polysemous words or Arabic transliterations from languages that were historically spoken in those areas. For example, the current names of cities and villages in the Levant region could be Arabic transliterations from Canaanite languages that do not have other lexical uses in the Arabic language. This mixture of terms presents an ideal testing ground for LMs.[4] We analyze QA performance on the 100 most frequent locations for each Arab country in the mC4 corpus. To determine if an entity matches an Arabic word that holds multiple meanings, we use the Almaany dictionary.[5] We compare the performance with that of the 100 most frequent Western locations for each Western country in CAMeL-2. We use all locations available for Western countries with less than 100 locations.

**Results.** Figure 4 shows the average QA accuracy for LMs of different sizes on the highly-frequent

---

[4]We refer the reader to Appendix C.2 for more background on the regional linguistic influences in Arab countries.

[5]Almaany is a comprehensive resource that provides multiple meanings for Arabic words, highlighting their different uses across contexts: https://www.almaany.com/
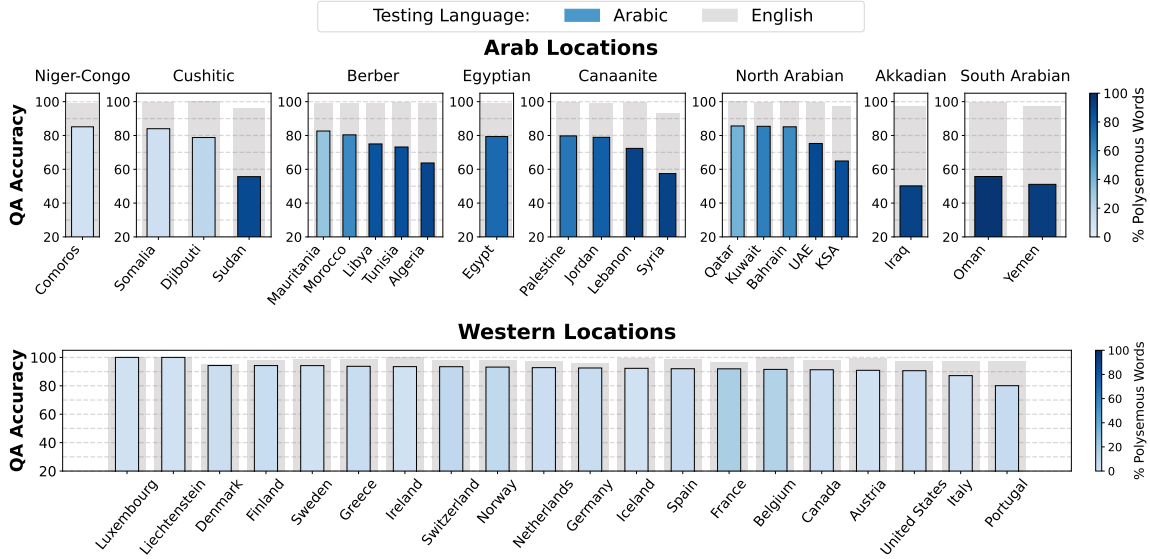
Figure 5: Average QA Accuracy (↑) of Llama3.3-70b on the top-100 most frequent location entities in mC4 for each Arab and Western country in CAMeL-2 (§4.2). Arab countries are grouped by the language family that influences location naming in their region. Performance on Arab locations decreases as the percentage of entities that are Arabic polysemous words increases. Performance in English on the same entities is shown as a gray background.

Arab locations, which we separate as polysemous and non-polysemous entities. In general, there is a positive scaling trend in both cases by all models. Yet, performance on polysemous entities is lower compared with non-polysemous ones, with 70B-sized models barely reaching the 70% margin, as opposed to non-polysemous entities where most models reach near 90%.

The average results per country achieved by Llama3.3-70b are shown in Figure 5, where Arab countries are grouped by the influencing language on location naming in their region. We observe a trend where **QA performance in Arabic drops with the increase in the percentage of entities which are Arabic polysemous words**, with accuracy reducing drastically to the 40-60% range. Performance is the best for countries where the percentage of polysemous words is low (e.g., Comoros, where entities are transliterations from the Comorian language).

In contrast, we see that this issue is non-existent for Western entities in Arabic, where accuracy is near 90% for all countries. This is because Western entities, being transliterations in Arabic, do not possess any other meaning. When the model is tested on the same entities in English, this problem also fades away, as Arab entities in English do not exhibit word polysemy either. This highlights that **the struggle of LMs with entities that are polysemous words in Arabic leads to a *perceived* bias**

**towards Western entities as they do not exhibit this phenomenon**. We also obtain similar results on NER with BERT-type LMs (see Appendix D.2).

### 4.3 Other Languages Using Arabic Script

While Arabic script is primarily associated with the Arabic language, it is also used in several other languages, such as Farsi, Urdu, Kurdish, Tajik, and Pashto, due to historical and cultural connections between regions where they are spoken and the Arab world. There is thus a natural overlap of words between Arabic and those languages. We study how LMs behave on Arab entities as a function of their frequency in the pre-training data of other languages. We use the mC4 portions of those languages to obtain a total occurrence count for each Arab entity.

**Results.** Figure 6 shows the average QA accuracy and CBS at text-infilling achieved by LMs on Arab location, name, and food entities versus their total count in other languages. We observe a general trend where **LMs struggle on entities as they occur more frequently in other languages that share the script with Arabic**. Such entities occurring at very high frequencies can be common words in those languages that hold their own different meanings. For example, the word "وزان" used to denote the Moroccan town "*Ouzanne*" is also used in Farsi as the word for "*weight*".
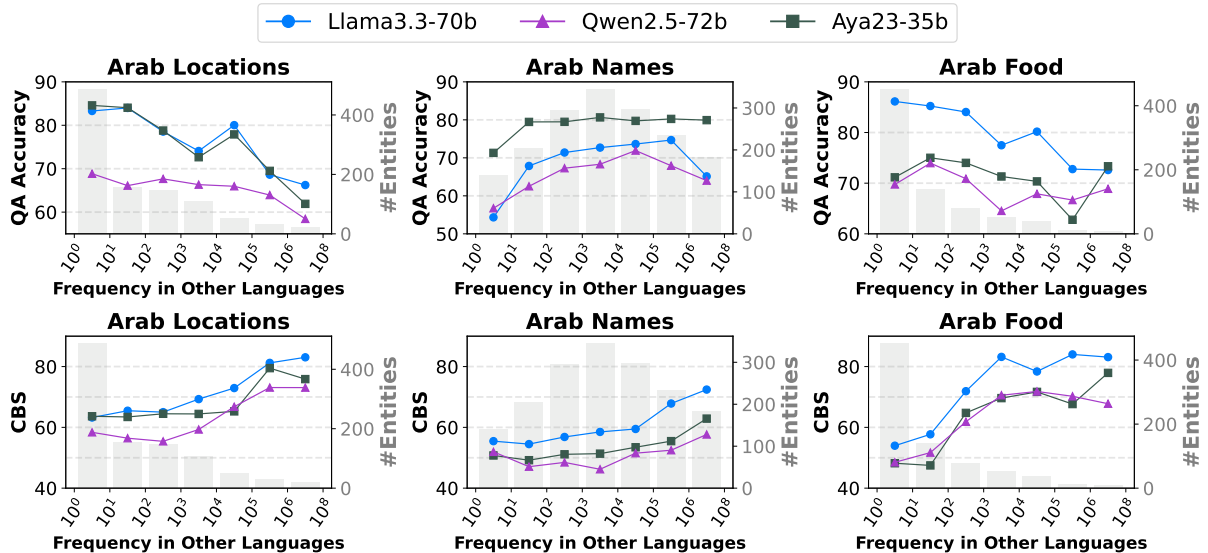
Figure 6: QA Accuracy (↑) and CBS (↓) of LMs on Arab location, name, and food entities vs. their total count in other languages that use Arabic script (*Farsi, Urdu, Tajik, Pashto, Kurdish*) in mC4 corpus. Performance decreases for all LMs as entities appear more frequently in other languages, especially for location and food entities.

As multilingual models are trained on those languages together, such lexical overlaps between entities in one language and highly frequent non-entity words in other languages with shared script can cause LMs to struggle at recognizing those entities. This issue is most noticeable for location and food entities but less so for name entities where performance is more stable. This could be due to the fact that Arab name entities are also used for first names in those languages, rather than having their own separate meanings.

### 4.4 Subword Tokenization Matters

We analyze how tokenization (Kudo, 2018; Song et al., 2021; Bostrom and Durrett, 2020) impacts the behavior of LMs on entities. We compare the performance of LMs on NER and extractive QA with respect to how many tokens they get fragmented into. We also separate entities based on whether they exhibit polysemy at the token level, where we check if a tokens matches an Arabic word that can be used for different functionalities.

**Results.** Figure 7 shows the performance distribution on Arab location entities by Llama3.3-70b, JAIS-13b, and ARBERT. We find that **LMs perform the worst on entities that are tokenized into only one token, especially when it corresponds to a polysemous word**. Performance improves when entities are tokenized into multiple tokens, with gaps between entities containing polysemous and non-polysemous tokens gradually reducing. In-

terestingly, this issue is the most apparent for the ARBERT model, which is trained only on Arabic data and has a very large vocabulary of 93k tokens.

Figure 8 shows the performance by LMs on one-token entities, in relation to the size of their Arabic vocabularies. We find that **LMs with medium Arabic vocabulary sizes perform the best, while performance drops for ones with very large vocabularies**. As the vocabulary size increases, frequency-based tokenization schemes will merge frequently used words into single tokens. This likely makes it more challenging for LMs to recognize entities in Arabic that exhibit word polysemy, as they get tokenized and encoded in the same way as when those words appear in text with non-entity senses within the same language or other languages sharing the same script. Such observations motivate the need for better tokenization approaches that can handle these cases for better cross-cultural performance.

### 5 Related Work

There has been growing interest in studying the cultural considerations surrounding LMs (Liu et al., 2024; AlKhamissi et al., 2024). This encompasses several aspects that LMs have to navigate such as *cultural commonsense*, where LMs must differentiate between the societal norms (e.g., bringing a gift when visiting someone) of different cultures (Shi et al., 2024; Chiu et al., 2024; Bhatt and Diaz, 2024; Palta and Rudinger, 2023; Fung et al., 2023; Huang and Yang, 2023). Other works explored
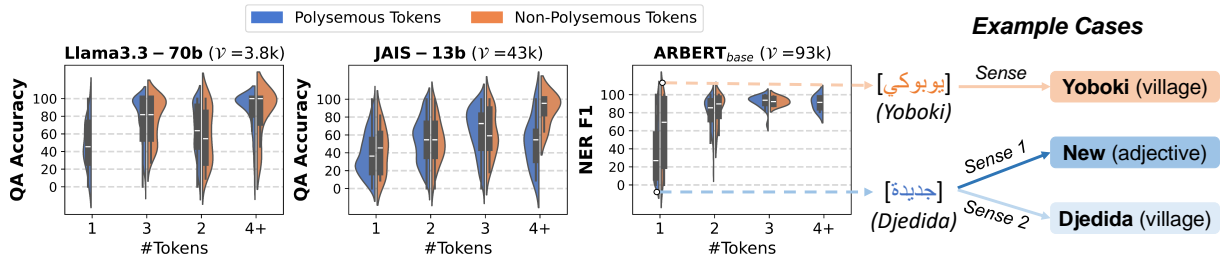
Figure 7: Performance distribution of Llama3.3-70b, JAIS-13b, and ARBERT on Arab location entities, in relation to how many tokens they get tokenized into. Entities are separated based on whether tokens correspond to Arabic polysemous words. $\mathcal{V}$ represents the number of Arabic tokens in each LM's vocabulary. Performance is the poorest on one-token entities that exhibit word polysemy, and improves on entities represented by multiple tokens.
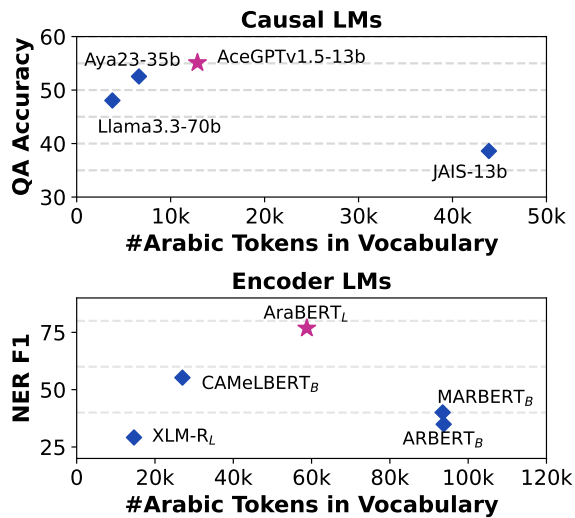


Figure 8: Average QA Accuracy and F1 of LMs on one-token entities vs. the number of Arabic tokens in their vocabularies. Performance drops greatly as vocabularies get too large, with the best performing LMs (★) having medium Arabic vocabulary sizes.

*culture-specific knowledge* in LMs (e.g., the color of the bridal dress) and how their performance varies for different countries (Keleg and Magdy, 2023; Yin et al., 2022). It has been shown that LMs are mostly familiar with a few cultures dominated in pre-training data, but struggle with less represented cultures (Shen et al., 2024; Rao et al., 2024; Seth et al., 2024), and in non-English languages (Arora et al., 2024; Masoud et al., 2023).

Another line of work explores the cultural appropriateness of LMs by analyzing their behavior when handling *entities that exhibit cultural variation* such as first names (An et al., 2024; Nghiem et al., 2024; Jeoung et al., 2023; An and Rudinger, 2023; Gautam et al., 2024) and food dishes (Zhou et al., 2024; Li et al., 2024b). The study of Naous et al. (2024) revealed performance gaps in LMs when handling entities associated with Arab vs.

Western culture in Arabic language, where models were performed consistently better on Western entities. Our work builds on this study, with the aim of pinpointing the origins of such gaps.

Past research on analyzing the cross-cultural performance of LMs have focused primarily on the representation of entities in pre-training data (Wolfe and Caliskan, 2021; Mukherjee et al., 2024), demonstrating a struggle of LMs to learn knowledge about entities that rarely appear in corpora (Li et al., 2024a; Kandpal et al., 2023). However, those analyses have been restricted to the English language only. Different from prior work, we analyze multiple facets that can contribute to Western biases in LMs beyond pre-training data, focusing on cross-linguistic differences in phenomena exhibited by entities and the impact of subword tokenization.

## 6 Conclusion

We analyzed a variety of factors that can contribute to entity-related cultural biases in LMs. We showed how non-English linguistic phenomena such as word polysemy in Arabic, lexical overlaps with other languages, and frequency-based tokenization can cause performance degradation on entities associated with Arab culture, leading to perceived Western biases in LMs. We hope our study lays a foundation of important aspects to consider in building culturally fair multilingual LMs.

## Limitations

In this work, we investigated the origins of entity-related cultural biases in LMs by probing their ability at culturally-appropriate text-infilling and analyzing their cross-cultural performance on the extractive QA and NER tasks. There are other entity-related cultural biases that can also manifest in model behavior such as sentiment and stereotype associations in generated text (Naous et al., 2024).

We leave the exploration of the reasons behind such learned associations for future work.

Our analyses showed that LMs struggle on entities in Arabic that are associated with Arab culture when they exhibit word polysemy. While Western entities transliterated into Arabic mostly do not exhibit this phenomenon, there are cases where a transliteration could match a random word used in Arabic language based on the closeness of their phonetic pronunciation. For example, the name "*Ben*" could be transliterated as "بن" which is used in Arabic for "*son of*" and "*powdered coffee*". There are other cases where transliterations of less common entities could match the transliteration of other famous entities. For example, the name "*Yvonne*" could be transliterated as "ايفون" which is also the common Arabic transliteration for "*iPhone*". Such cases, while being rare, could cause LMs to fail at recognizing certain Western entities in Arabic. Although we did not explore the sensitivity of LMs to this phenomenon, the parallel entities in CAMeL-2 can offer a valuable resource for future studies to better analyze such cases.

Our work focuses on Arab culture and analyzes entity-related biases in Arabic language. Such biases may also be manifested by LMs in many other non-Western languages. Future studies can follow the process described in this work to extend CAMeL-2 to such languages.

## Acknowledgements

## References

Muhammad Abdul-Mageed, AbdelRahim Elmadany, and El Moatez Billah Nagoudi. 2020. ARBERT & MARBERT: Deep bidirectional transformers for Arabic. *arXiv preprint arXiv:2101.01785*.

Muhammad Farid Adilazuarda, Sagnik Mukherjee, Pradhyumna Lavania, Siddhant Singh, Ashutosh Dwivedi, Alham Fikri Aji, Jacki O'Neill, Ashutosh Modi, and Monojit Choudhury. 2024. Towards measuring and modeling "culture"' in LLMs: A survey. *arXiv preprint arXiv:2403.15412*.

Badr AlKhamissi, Muhammad ElNokrashy, Mai AlKhamissi, and Mona Diab. 2024. Investigating cultural alignment of large language models. *arXiv preprint arXiv:2402.13231*.

Haozhe An, Christabel Acquaye, Colin Wang, Zongxia Li, and Rachel Rudinger. 2024. Do large language models discriminate in hiring decisions on the basis of race, ethnicity, and gender? *arXiv preprint arXiv:2406.10486*.

Haozhe An and Rachel Rudinger. 2023. Nichelle and Nancy: The influence of demographic attributes and tokenization length on first name biases. *arXiv preprint arXiv:2305.16577*.

Wissam Antoun, Fady Baly, and Hazem Hajj. 2020. AraBERT: Transformer-based model for arabic language understanding. In *LREC Workshop Language Resources and Evaluation Conference 11–16 May 2020*, page 9.

Shane Arora, Marzena Karpinska, Hung-Ting Chen, Ipsita Bhattacharjee, Mohit Iyyer, and Eunsol Choi. 2024. CaLMQA: Exploring culturally specific long-form question answering across 23 languages. *arXiv preprint arXiv:2406.17761*.

Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat Venkitesh, Madeline Smith, Kelly Marchisio, Sebastian Ruder, et al. 2024. Aya 23: Open weight releases to further multilingual progress. *arXiv preprint arXiv:2405.15032*.

Yassine Benajiba, Paolo Rosso, and José Miguel Benedíruiz. 2007. Anersys: An arabic named entity recognition system based on maximum entropy. In *Computational Linguistics and Intelligent Text Processing: 8th International Conference, CICLing 2007, Mexico City, Mexico, February 18-24, 2007. Proceedings 8*, pages 143–153. Springer.

Shaily Bhatt and Fernando Diaz. 2024. Extrinsic evaluation of cultural competence in large language models. *arXiv preprint arXiv:2406.11565*.

Terra Blevins, Tomasz Limisiewicz, Suchin Gururangan, Margaret Li, Hila Gonen, Noah A. Smith, and Luke Zettlemoyer. 2024. Breaking the curse of multilinguality with cross-lingual expert language models. *arXiv preprint arXiv:2401.10440*.

Kaj Bostrom and Greg Durrett. 2020. Byte Pair Encoding is suboptimal for language model pretraining. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4617–4624, Online. Association for Computational Linguistics.

Yu Ying Chiu, Liwei Jiang, Maria Antoniak, Chan Young Park, Shuyue Stella Li, Mehar Bhatia, Sahithya Ravi, Yulia Tsvetkov, Vered Shwartz, and Yejin Choi. 2024. CulturalTeaming: AI-assisted interactive red-teaming for challenging LLMs'(lack of) multicultural knowledge. *arXiv preprint arXiv:2404.06664*.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised

cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Kareem Darwish and Hamdy Mubarak. 2016. Farasa: A new fast and accurate Arabic word segmenter. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 1070–1074.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The Llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Yi Fung, Tuhin Chakrabarty, Hao Guo, Owen Rambow, Smaranda Muresan, and Heng Ji. 2023. NORM-SAGE: Multi-lingual multi-cultural norm discovery from conversations on-the-fly. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15217–15230.

Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. 2024. Bias and fairness in large language models: A survey. *Computational Linguistics*, pages 1–79.

Vagrant Gautam, Arjun Subramonian, Anne Lauscher, and Os Keyes. 2024. Stop! In the name of flaws: Disentangling personal names and sociodemographic attributes in NLP. *arXiv preprint arXiv:2405.17159*.

Gene Gragg. 2019. Semitic and Afro-Asiatic. In *The Semitic Languages*, pages 22–48. Routledge.

Jing Huang and Diyi Yang. 2023. Culturally aware natural language inference. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7591–7609.

Go Inoue, Bashar Alhafni, Nurpeiis Baimukan, Houda Bouamor, and Nizar Habash. 2021. The interplay of variant, size, and task type in arabic pre-trained language models. *arXiv preprint arXiv:2103.06678*.

Sullam Jeoung, Jana Diesner, and Halil Kilicoglu. 2023. Examining the causal effect of first names on language models: The case of social commonsense reasoning. *arXiv preprint arXiv:2306.01117*.

Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In *International Conference on Machine Learning*, pages 15696–15707. PMLR.

Amr Keleg and Walid Magdy. 2023. DLAMA: A framework for curating culturally diverse facts for probing the knowledge of pretrained language models. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 6245–6266.

Taku Kudo. 2018. Subword regularization: Improving neural network translation models with multiple subword candidates. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 66–75, Melbourne, Australia. Association for Computational Linguistics.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pages 611–626.

Huihan Li, Arnav Goel, Keyu He, and Xiang Ren. 2024a. Attributing culture-conditioned generations to pretraining corpora. *Preprint*, arXiv:2412.20760.

Wenyan Li, Xinyu Zhang, Jiaang Li, Qiwei Peng, Raphael Tang, Li Zhou, Weijia Zhang, Guimin Hu, Yifei Yuan, Anders Søgaard, Daniel Hershcovich, and Desmond Elliott. 2024b. FoodieQA: A multimodal dataset for fine-grained understanding of chinese food culture. *arXiv preprint arXiv:2406.11030*.

Chen Liang, Yue Yu, Haoming Jiang, Siawpeng Er, Ruijia Wang, Tuo Zhao, and Chao Zhang. 2020. Bond: BERT-assisted open-domain named entity recognition with distant supervision. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1054–1064.

Davis Liang, Hila Gonen, Yuning Mao, Rui Hou, Naman Goyal, Marjan Ghazvininejad, Luke Zettlemoyer, and Madian Khabsa. 2023. Xlm-v: Overcoming the vocabulary bottleneck in multilingual masked language models. *arXiv preprint arXiv:2301.10472*.

Juhao Liang, Zhenyang Cai, Jianqing Zhu, Huang Huang, Kewei Zong, Bang An, Mosen Alharthi, Juncai He, Lian Zhang, Haizhou Li, et al. 2024. Alignment at pre-training! towards native alignment for Arabic llms. *arXiv preprint arXiv:2412.03253*.

Chen Cecilia Liu, Iryna Gurevych, and Anna Korhonen. 2024. Culturally aware and adapted NLP: A taxonomy and a survey of the state of the art. *arXiv preprint arXiv:2406.03930*.

Reem I Masoud, Ziquan Liu, Martin Ferianc, Philip Treleaven, and Miguel Rodrigues. 2023. Cultural alignment in large language models: An explanatory analysis based on hofstede's cultural dimensions. *arXiv preprint arXiv:2309.12342*.

Anjishnu Mukherjee, Aylin Caliskan, Ziwei Zhu, and Antonios Anastasopoulos. 2024. Global gallery: The fine art of painting culture portraits through multilingual instruction tuning. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6398–6415.

Tarek Naous, Michael Ryan, Alan Ritter, and Wei Xu. 2024. Having beer after prayer? measuring cultural bias in large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16366–16393.

Huy Nghiem, John Prindle, Jieyu Zhao, and Hal Daumé III. 2024. "you gotta be a doctor, Lin": An investigation of name-based bias of large language models in employment recommendations. *arXiv preprint arXiv:2406.12232*.

Shramay Palta and Rachel Rudinger. 2023. FORK: A bite-sized test set for probing culinary cultural biases in commonsense reasoning models. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 9952–9962.

Peng Qi, Timothy Dozat, Yuhao Zhang, and Christopher D Manning. 2018. Universal dependency parsing from scratch. In *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 160–170.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.

Michael Ramscar. 2019. Source codes in human communication. *arXiv preprint arXiv:1904.03991*.

Abhinav Rao, Akhila Yerukola, Vishwa Shah, Katharina Reinecke, and Maarten Sap. 2024. Normad: A benchmark for measuring the cultural adaptability of large language models. *arXiv preprint arXiv:2404.12464*.

Dwight F Reynolds. 2015. *The Cambridge companion to modern Arab culture*. Cambridge University Press.

Erik Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition. In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, pages 142–147.

Neha Sengupta, Sunil Kumar Sahu, Bokang Jia, Satheesh Katipomu, Haonan Li, Fajri Koto, William Marshall, Gurpreet Gosal, Cynthia Liu, Zhiming Chen, et al. 2023. JAIS and JAIS-Chat: Arabic-centric foundation and instruction-tuned open generative large language models. *arXiv preprint arXiv:2308.16149*.

Agrima Seth, Sanchit Ahuja, Kalika Bali, and Sunayana Sitaram. 2024. DOSA: A dataset of social artifacts from different indian geographical subcultures. *arXiv preprint arXiv:2403.14651*.

Siqi Shen, Lajanugen Logeswaran, Moontae Lee, Honglak Lee, Soujanya Poria, and Rada Mihalcea. 2024. Understanding the capabilities and limitations of large language models for cultural commonsense. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5668–5680.

Weiyan Shi, Ryan Li, Yutong Zhang, Caleb Ziems, Raya Horesh, Rogério Abreu de Paula, Diyi Yang, et al. 2024. Culturebank: An online community-driven knowledge base towards culturally aware language technologies. *arXiv preprint arXiv:2404.15238*.

Shivalika Singh, Freddie Vargus, Daniel D'souza, Börje Karlsson, Abinaya Mahendiran, Wei-Yin Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura O'Mahony, Mike Zhang, Ramith Hettiarachchi, Joseph Wilson, Marina Machado, Luisa Moura, Dominik Krzemiński, Hakimeh Fadaei, Irem Ergun, Ifeoma Okoh, Aisha Alaagib, Oshan Mudannayake, Zaid Alyafeai, Vu Chien, Sebastian Ruder, Surya Guthikonda, Emad Alghamdi, Sebastian Gehrmann, Niklas Muennighoff, Max Bartolo, Julia Kreutzer, Ahmet Üstün, Marzieh Fadaee, and Sara Hooker. 2024. Aya dataset: An open-access collection for multilingual instruction tuning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11521–11567, Bangkok, Thailand. Association for Computational Linguistics.

Xinying Song, Alex Salcianu, Yang Song, Dave Dopson, and Denny Zhou. 2021. Fast WordPiece tokenization. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2089–2103, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D'souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, et al. 2024. Aya model: An instruction finetuned open-access multilingual language model. *arXiv preprint arXiv:2402.07827*.

Willy Van Langendonck. 2007. *Theory and typology of proper names*. Mouton de Gruyter.

Kees Versteegh. 2014. *Arabic language*. Edinburgh University Press.

Robert Wolfe and Aylin Caliskan. 2021. Low frequency names exhibit bias and overfitting in contextualizing language models. *arXiv preprint arXiv:2110.00672*.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

Da Yin, Hritik Bansal, Masoud Monajatipoor, Liu-nian Harold Li, and Kai-Wei Chang. 2022. GeoM-LAMA: Geo-diverse commonsense probing on multi-lingual pre-trained language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2039–2055.

Li Zhou, Taelin Karidi, Nicolas Garneau, Yong Cao, Wanlong Liu, Wenyu Chen, and Daniel Hershcovich. 2024. Does mapo tofu contain coffee? probing llms for food-related cultural knowledge. *arXiv preprint arXiv:2404.06833*.

Jianqing Zhu, Huang Huang, Zhihang Lin, Juhao Liang, Zhengyang Tang, Khalid Almubarak, Abdulmohsen Alharthik, Bang An, Juncai He, Xiangbo Wu, et al. 2024. Second language (Arabic) acquisition of llms via progressive vocabulary expansion. *arXiv preprint arXiv:2412.12310*.

| Entity Type | Wikipedia Categories | |
| --- | --- | --- |
| | *Arabic* | *English* |
| Authors | [adjective] كتاب وكاتبات | [adjective] writers |
| | [adjective] روائيون | [adjective] authors |
| Beverage | [adjective] مطبخ | [adjective] cuisine |
| Food | [adjective] مطبخ | [adjective] cuisine |
| Names | [adjective] سياسيون | [adjective] politicians |
| | [adjective] رياضيون | [adjective] athletes |
| | [adjective] ممثلون | [adjective] actors |
| | [adjective] كتاب وكاتبات | [adjective] writers |
| | [adjective] روائيون | [adjective] authors |
| Religious | مساجد في [country] | mosques in [country] |
| | كنائس في [country] | churches in [country] |
| Sports Clubs | أندية كرة قدم في [country] | football clubs in [country] |

Table 4: List of Arabic Wikipedia categories used to perform country-wise extraction of cultural entities. [adjective] refers to the country-specific adjective (e.g., *Palestinian*, *Irish*, *Thai*, etc.).

## A  CAMeL-2: Additional Details

Figure 9 shows the distribution of entities in CAMeL-2 stratified by their association with *Arab culture*, *Western culture*, or *Other Foreign Culture*, as well as their source of collection (Wikidata/CommonCrawl entities collected in CAMeL (Naous et al., 2024), and newly collected entities from Wikipedia and OpenStreetMap).

**Wikipedia-based Extraction.** The Arabic Wikipedia categories used to perform country-wise extraction of articles for each entity type are listed in Table 4. We first de-duplicate the extracted articles from Wikipedia, as many of them would be cross-listed under the category of multiple countries. Many of the extracted articles would be irrelevant to the entity type of interest and were thus manually filtered out. For example, in addition to articles about food dishes, the "[country adjective] cuisine" category would also contain articles about particular chefs or restaurants in that country. We finally inspect the titles of the remaining articles and remove any additional text between parentheses that is not part of the entity (e.g., a title such as "*Mandi (food)*" where *(food)* was manually removed).

For author names, the number of articles in the `authors` category was very large. We thus took a random sample of 3k articles from Arab countries and 3k articles from Western countries that were then manually filtered by the annotators.

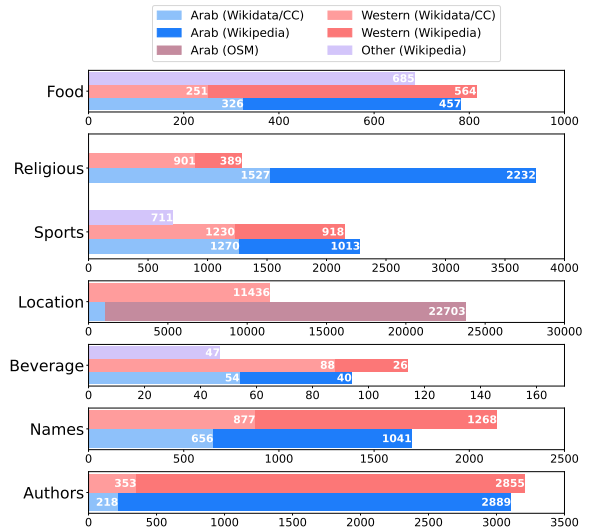To collect first names, we first extract all article



Figure 9: Distribution of entities in CAMeL-2 for each entity type stratified by association with *Arab culture*, *Western culture*, or *Other Foreign Culture*, as well as their data collection source: Wikidata, CommonCrawl (CC), Wikipedia, OpenStreetMap (OSM).

titles and text from multiple Wikipedia categories, as shown in Table 4, that relate to human entities (*politicians*, *actors*, *athletes*, *writers*, and *authors*). We then extract the first uni-gram and bi-gram from each article title (which represents the full name of the human entity) and perform de-duplication. Our annotators then filter out extractions that are not person names and classify the extracted names for cultural association (*Arab*, *Westerm*, or *Other Foreign Culture*). The extracted names were also classified as masculine or feminine by the annotators, a step that is necessary to match the gendered grammar of the Arabic language. The annotators were also given access to the corresponding Wikipedia article for each extraction to guide their decision in annotation. This process was done for all Wikipedia articles extracted from Arab countries, resulting in 1,268 new Arab names. The articles obtained from Western countries were too large in number (> 20k articles), we thus took a random sample of 1.5k articles that were filtered by the annotators.

**Location Extraction from Georgraphic Data.** Open Street Maps (OSM) uses a "`place`" tag to represent various types of locations. For each Arab country, we extract all locations that have place tags of *city*, *town*, *village*, *neighborhood*, and *suburb*. We discard locations that have other highly-specific place tags such as *isolated dwelling*, *hamlet*, *farm*, etc. since these represent individual residential structures, often in remote areas, rather than or-

| | Aya23-35b | | | | | | AceGPTv2-70b | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Arabic | | | English | | | Arabic | | | English | | |
| | *Arab* | *Western* | ΔAcc | *Arab* | *Western* | ΔAcc | *Arab* | *Western* | ΔAcc | *Arab* | *Western* | ΔAcc |
| Authors | 91.68 | 88.96 | -2.72 | 79.03 | 81.73 | 2.70 | 89.43 | 82.86 | -6.57 | 88.61 | 95.50 | 6.89 |
| Beverage | 62.81 | **75.61** | 12.80 | 76.47 | 76.06 | -0.41 | 89.47 | 81.92 | -7.55 | 98.62 | 96.06 | -2.56 |
| Food | 70.92 | 67.97 | -2.95 | 72.93 | 81.05 | 8.12 | 88.57 | 79.47 | -9.10 | 95.39 | 97.48 | 2.09 |
| Location | 81.04 | **91.55** | 10.51 | 90.33 | 92.89 | 2.56 | 78.60 | **89.97** | 11.37 | 97.88 | 98.96 | 1.08 |
| Names (F) | 78.78 | **87.30** | 8.52 | 89.85 | 85.91 | -3.94 | 77.73 | **88.78** | 11.05 | 99.84 | 99.09 | -0.75 |
| Names (M) | 79.53 | **80.70** | 1.17 | 72.96 | 73.90 | 0.94 | 85.03 | **87.83** | 2.80 | 94.53 | 94.45 | -0.08 |
| Sports | 47.88 | **53.85** | 5.97 | 81.78 | 77.76 | -4.02 | 72.20 | 66.34 | -5.86 | 80.24 | 83.37 | 3.13 |
| Religious | 82.49 | 79.35 | -3.14 | 90.11 | 93.47 | 3.36 | 80.07 | **81.65** | 1.58 | 79.29 | 84.67 | 5.38 |

Table 5: Average QA Accuracy of Aya23-35b and AceGPTv2-70b on Arab and Western entities when tested in Arabic and English. ΔAcc represents performance differences between Western and Arab entities.
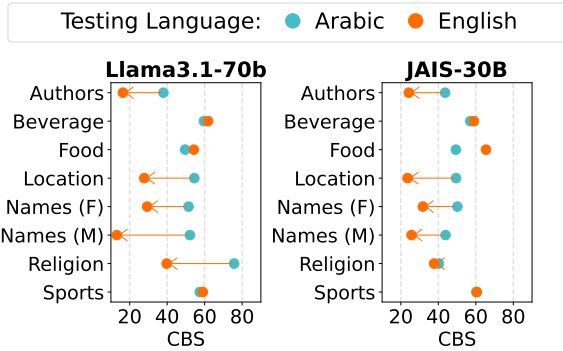


Figure 10: Average CBS per entity type achieved by Llama3.1-70b and JAIS-30b on culturally-grounded contexts from CAMeL-2.



Figure 11: Prompt template used perform extractive QA with GPT-type LMs.

ganized settlements with social significance. We also discarded locations which are more than one word expressions as they mostly consisted of repeatedly used terms (e.g., hill, mountain, valley, spring, etc.).

# B  Arabic vs. English Comparisons

## B.1  Text-Infilling

**Additional Results.**  Figure 10 reports CBS results on text-infilling by additional Llama3.1-70b and JAIS-30b when tested in both Arabic and English. Similar to our observations in §3.2, we see a consistent trend across all LMs where better adaptation to Arab cultural contexts is achieved in English compared to Arabic, where CBS values are greatly reduced.

## B.2  Extractive QA

**Prompt.**  The prompt template used for performing extractive QA of entities with GPT-type LMs is shown in Figure 11. The [entity type] in the template is replaced with the name of the entity type of interest (i.e., location, person's name, author name, food dish, drink, mosque, church,
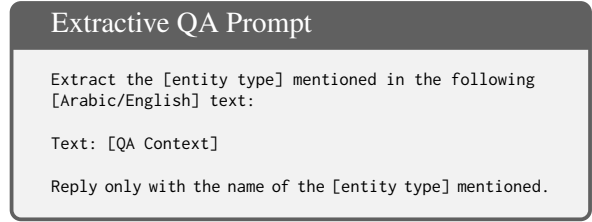
football club). The [QA Context] is replaced by one of the QA contexts collected for the respective entity type (§2.3) where the entity to be extracted is placed instead of the [MASK] token of the context. The instruction is given in English for all models tested, expect for JAIS where giving the model the instruction in Arabic lead to better performance.

**Additional Results.**  Table 5 shows additional results achieved by Aya23-35b and AceGPTv2-70b on Arab and Western entities across different entity types. We observe similar results to those achieved by Llama3.3-70b (§ 3.3) where we see large performance gaps between Western and Arab entities in Arabic for most entity types, but gaps between cultures are much smaller in English. Notably, we do find an improvements in those models on Arab entities in the religious and food entity type compared with Llama3.3-70b, which may be due to additional training on Arabic data.

## B.3  Cultural Fairness: NER

**Distant Supervision.**  To fine-tune LMs for NER of *food*, *beverage*, *sports clubs*, and *religious places*, we leverage text from Wikipedia articles to perform distantly supervised fine-tuning. Specifically, we use entities from CAMeL-2 that are linked to Wikipedia articles. For each of these entities, we automatically create fine-tuning samples by extract-

| Entity Type | # Fine-tuning Samples (train/val/test) | |
|---|---|---|
| | *Arabic* | *English* |
| Food & Beverage | 6,908/768/853 | 27,635/3,071/1,617 |
| Religious | 2,655/295/328 | 12,623/1,403/1,559 |
| Sports Clubs | 2,484/277/307 | 18,585/2,066/2,295 |

Table 6: Number of NER fine-tuning examples (as train/val/test splits) constructed automatically from Wikipedia articles for *food*, *beverage*, *sports clubs*, and *religious place* entities.

| Model | Test Set F1 Score | | | | | |
|---|---|---|---|---|---|---|
| | *Arabic* | | | *English* | | |
| | Foo | Spo | Rel | Foo | Spo | Rel |
| XLMR$_{large}$ | 77.32 | 85.62 | 91.52 | 90.47 | 82.30 | 82.26 |
| XLMV$_{base}$ | 76.39 | 87.14 | 90.73 | 91.75 | 83.17 | 80.63 |
| ARBERT | 79.44 | 86.48 | 92.43 | — | — | — |

Table 7: F1 score achieved by various BERT-type LMs on the test set of the fine-tuning samples created automatically from Wikipedia.

ing sentences from their corresponding Wikipedia articles where the entity is mentioned. We exclude entities from fine-tuning that do not appear on Wikipedia (i.e., entities extracted from Wikidata or Commoncrawl in CAMeL (Naous et al., 2024)). Table 6 shows the number of fine-tuning examples created in both Arabic and English.

**Additional Results.** Table 8 shows the average F1 achieved by XLMV$_{base}$ on Arab and Western entities across different entity types when tested in Arabic and English. Table 9 shows the results of the Arabic monolingual ARBERT model when tested in Arabic. We observe similar trends to our analysis in (§ 3.3) where BERT-type LMs are consistently better at recognizing Western entities in the Arabic language, but performance gaps between cultures in much smaller when LMs are tested in English.

### B.4 Experimental Details

**Language Models.** Table 10 lists all the LMs used in our experiments, including the varying sizes used in our scaling analysis (§4.2) and different versions of the models for our vocabulary size analysis (§4.4 and Appendix D.3).

**QA Inference.** We ran our experiments using 8 NVIDIA A40 GPUs. For inference on the extractive QA task with causal LMs, we used the

| XLMV$_{base}$ | | | | | | |
|---|---|---|---|---|---|---|
| | Arabic | | | English | | |
| | *Arab* | *Western* | ΔF1 | *Arab* | *Western* | ΔF1 |
| Authors | 93.14 | **94.99** | 1.85 | 96.61 | 95.13 | -1.48 |
| Beverage | 52.87 | **64.58** | 11.71 | 95.06 | 88.48 | -6.58 |
| Food | 50.15 | **62.16** | 12.01 | 92.51 | 90.55 | -1.96 |
| Location | 64.55 | **89.27** | 24.72 | 92.42 | 97.65 | 5.23 |
| Names (F) | 48.48 | **75.68** | 27.20 | 98.32 | 97.57 | -0.75 |
| Names (M) | 70.37 | **80.52** | 10.15 | 93.84 | 91.66 | -2.18 |
| Sports | 75.37 | **83.68** | 8.31 | 91.83 | 92.68 | 0.85 |
| Religious | 71.22 | **84.23** | 13.01 | 90.61 | 86.06 | -4.55 |

Table 8: Average F1 of XLMV$_{base}$ (Liang et al., 2023) on Arab and Western entities when tested in Arabic and English. ΔF1 represents performance differences between Western and Arab entities.

| ARBERT | | | |
|---|---|---|---|
| | Arabic | | |
| | *Arab* | *Western* | ΔF1 |
| Authors | 96.01 | 93.83 | -2.18 |
| Beverage | 53.35 | **65.05** | 11.70 |
| Food | 52.01 | **64.24** | 12.23 |
| Location | 60.47 | **89.32** | 28.85 |
| Names (F) | 69.61 | **75.81** | 6.21 |
| Names (M) | 83.82 | **84.86** | 1.04 |
| Sports | 90.81 | **94.66** | 3.85 |
| Religious | 59.48 | **69.08** | 9.60 |

Table 9: Average F1 of the Arabic monolingual AR-BERT model on Arab and Western entities when tested in Arabic and English. ΔF1 represents performance differences between Western and Arab entities.

vLLM library[6] (Kwon et al., 2023) for fast inference. We performed greedy decoding by setting the following parameters {temperature=0, top_p=1, top_k=1}. We also limit the number of generated tokens by the models by setting {max_tokens=30}. Futher, we limit the context length by setting {max_model_len=4096}.

**Fine-tuning.** For fine-tuning BERT-type models on the NER task, we fine-tuned each model for 5 epochs using the cross-entropy loss and the Adam optimizer and tuned the learning rate in the set $\{1e^{-5}, 1e^{-6}, 1e^{-7}\}$. We selected checkpoints based on the best F1 on the validation set. Fine-tuning was done using one NVIDIA A100 GPU.

**Entity Occurrence Counts.** To obtain counts of entities in the mC4 corpus, we use the Aho-Corasick string search algorithm[7] where we construct finite state machines using the entities, allowing for efficient transversal over the corpora.

---
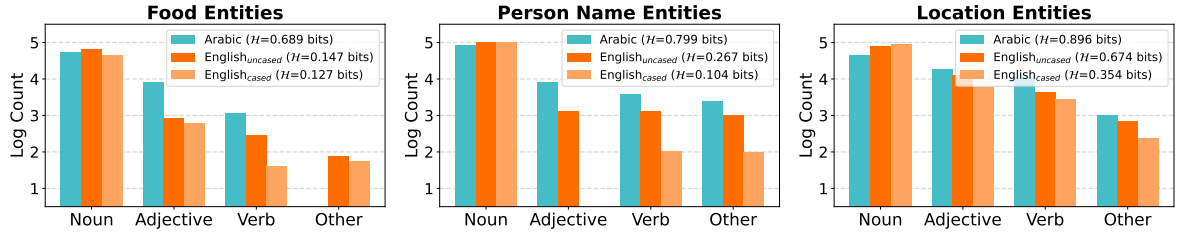[6]https://docs.vllm.ai
[7]https://pyahocorasick.readthedocs.io

Figure 12: POS tag distribution of the 100 most frequent food, name, and location entities in Arabic and English mC4. Entities in Arabic encode higher information as they appear more often with different grammatical roles.

| Language Model | Hugging Face Repository |
|---|---|
| *Causal LMs* | |
| Llama3.3-70b | meta-llama/Llama-3.3-70B-Instruct |
| Llama3.1-8b | meta-llama/Meta-Llama-3.1-8B-Instruct |
| Llama3.1-70b | meta-llama/Meta-Llama-3.1-70B-Instruct |
| Aya23-8b | CohereForAI/aya-23-8B |
| Aya23-35b | CohereForAI/aya-23-35B |
| Qwen2.5-0.5b | Qwen/Qwen2.5-0.5B-Instruct |
| Qwen2.5-3b | Qwen/Qwen2.5-3B-Instruct |
| Qwen2.5-14b | Qwen/Qwen2.5-14B-Instruct |
| Qwen2.5-32b | Qwen/Qwen2.5-32B-Instruct |
| Qwen2.5-72b | Qwen/Qwen2.5-72B-Instruct |
| AceGPTv1.5-13b | FreedomIntelligence/AceGPT-v1.5-13B-Chat |
| AceGPTv2-8b | FreedomIntelligence/AceGPT-v2-8B-Chat |
| AceGPTv2-70b | FreedomIntelligence/AceGPT-v2-70B-Chat |
| JAIS-13b | inceptionai/jais-13b-chat |
| JAIS-30b | inceptionai/jais-30b-chat-v3 |
| *Encoder LMs (multilingual)* | |
| XLMR$_{large}$ | FacebookAI/xlm-roberta-large |
| XLMV$_{base}$ | facebook/xlm-v-base |
| mDeBERTa-v3$_{base}$ | microsoft/mdeberta-v3-base |
| GigaBERT$_{base}$ | lanwuwei/GigaBERT-v3-Arabic-and-English |
| *Encoder LMs (monolingual - Arabic only)* | |
| AraBERT$_{large}$ | aubmindlab/bert-large-arabertv02 |
| ARBERT$_{base}$ | UBC-NLP/ARBERT |
| MARBERT$_{base}$ | UBC-NLP/MARBERT |
| CAMeLBERT$_{base}$ | CAMeL-Lab/bert-base-arabic-camelbert-mix |

Table 10: List of LMs used in our experiments and their repository links on Hugging Face.

## C Arabic Entities as Polysemous Words

### C.1 Comparing Polysemy Across Languages

To quantify and compare the prevalence of entity polysemy in Arabic and English, we analyze how often entities are used as different parts of speech (such as nouns, verbs, or adjectives) in texts written in both languages. We use the first 10M documents from the Arabic and English portions of the mC4 pre-training corpus (Raffel et al., 2020), which we then tokenize into sentences, yielding 239M Arabic sentences and 209M English sentences. Using the Arabic and English entities from CAMeL-2, we then identify the top 100 most frequent name, food, and location entities in the corpora. For each entity, we randomly sample 1000 sentences in which they appear, then determine their part-of-speech tag in each sentence using the Farasa POS tagger (Darwish and Mubarak, 2016) for Arabic and the Stanford POS tagger (Qi et al., 2018) for English. We perform this analysis for English entities both with and without uppercasing of the first letter.

Figure 12 shows the distribution of part-of-speech tags for the 100 most frequent name and food entities in Arabic and English. We group part-of-speech tags into *Nouns*, *Adjectives*, *Verbs*, and an *Other* category that encompasses tags such as particles, etc. We observe that the same words used for name entities in Arabic appear at high frequencies as adjectives and verbs rather than nouns. Arabic entities also have a higher Information Entropy $\mathcal{H}$[8] (Ramscar, 2019), measured at 0.799 bits for named compared with uncased English entities at 0.267 bits. It is important to note that English employs casing for name entities, facilitating a clear distinction between a name "*Mark*" and the verb "*mark*". In standard cased English text, occurrences of entities as adjectives or verbs are minimal, with entities predominantly appearing as nouns ($\mathcal{H}$ = 0.104 bits). However, Arabic does not have casing conventions, resulting in greater variability in the grammatical roles of named entities (used for descriptions as adjectives or actions as verbs), which can cause challenges for LMs in distinguishing between word senses as entities or non-entities. Conversely, named entities in English tend to adhere more consistently to noun roles.

### C.2 Regional Influences on Arabic Entities

While many named entities in the Arabic language are polysemous words, there are also words used for named entities that do not serve other functional purposes. These words are often Arabized forms of regional linguistic influences from different parts of the Arab world. A product of historical her-

---

[8]Information Entropy in bits for $N$ POS tags:
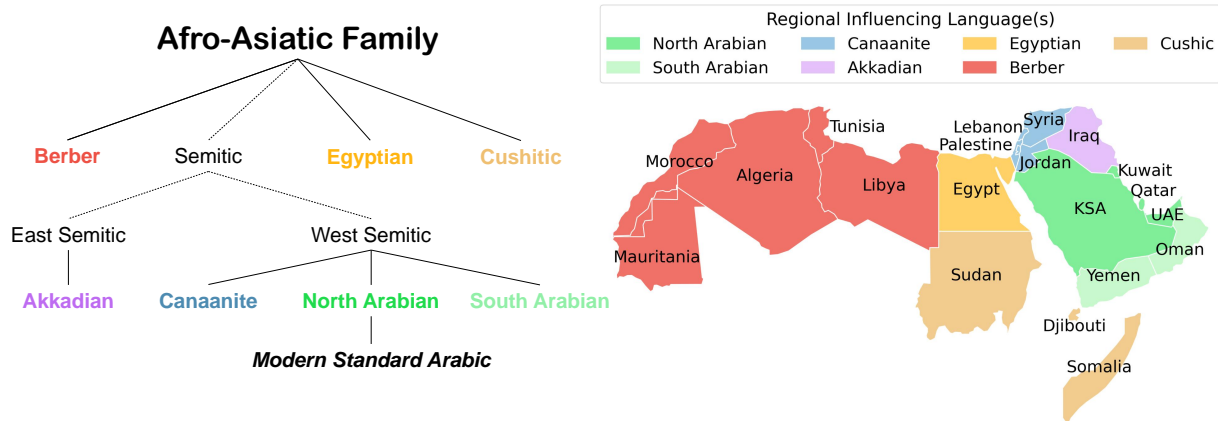$\mathcal{H} = -\sum_{i=1}^{N} p(\text{tag}_i) \log_2 p(\text{tag}_i)$

Figure 13: Map of the regional influencing languages on location names in Arab countries, and their standard classification in the Afro-Asiatic family according to Versteegh (2014). The *Comorian* language of the Comoros Island (not shown here on the map) is outside of the Afro-Asatic family and belongs to the Niger-Congo family.

itage, the countries of today's Arab world present a rich cultural mosaic. Following a sociological process of Arabization in North Africa and West Asia (Reynolds, 2015), the spread of the Arabic language in those areas and its interactions with regional languages led to the development of contemporary Arabic dialects. For example, the dialects in countries of the Levant region (i.e, Lebanon, Syria, Palestine, Jordan) have been shaped by influences from other Semitic languages historically spoken in those areas such as Aramaic, Hebrew, and Canaanite languages (Gragg, 2019). The names of many of today's cities, villages, and towns in those countries originate from the regional influencing languages, which predate the spread of Arabic. For instance, the Arabic naming of the Lebanese capital, Beirut (in Arabic: 'بيروت'), is a transliterated derivation of its Phoenician name "bīʾrōt". Such entities would only appear in the Arabic in the sense of a location and do not have any other lexical uses.

There are also regional linguistic influences on location names in the Arabian peninsula itself. Different Southern Arabian Languages were spoken in the southern part of the peninsula (i.e., present-day Yemen and Oman), while multiple Old North Arabian dialects were spoken in the central and northern parts of the peninsula (i.e., present-day Saudi Arabia). Contemporary Arabic (i.e., Modern Standard Arabic) is a continuum from Classical Arabic, the language of the Quraan, and early Islamic literature, which was the dialect of the Quraysh tribe in central Arabia. Figure 13 presents a visualization of the regional linguistic influences on each Arab country and their classification within the Afro-Asiatic language family (Versteegh, 2014).
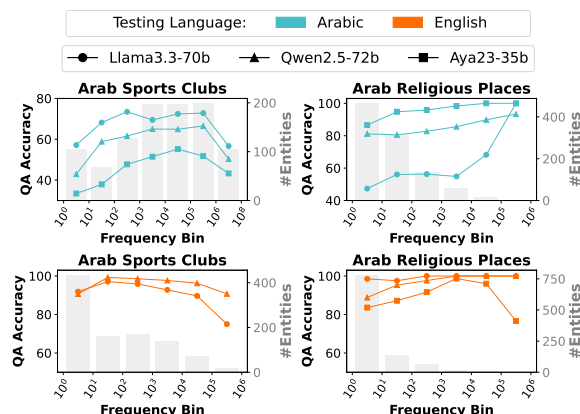


Figure 14: Average QA Accuracy (↑) of LLMs when tested in Arabic and English on sports clubs and religious places of worship associated with Arab culture, stratified by their occurrence counts in the mC4 corpus (grouped into log10-spaced bins).

## D   Analyses: Additional Results

### D.1   Entity Frequency in Pre-training Data

**Extractive QA.** Figure 14 shows the performance of LMs on extractive QA of sports clubs and religious places of worship stratified by their occurrence in Arabic and English pre-training data. We observe similar trends to our results in §4.1.

**Text Infilling.** Figure 15 shows the average CBS at the text-infilling task (§3.2) on entities stratified by their occurrence counts in pre-training. In this setup, we test each Arab entity against 30 randomly sampled Western entities across 5 randomly sampled culturally-contextualized contexts from CAMeL-2. Similar to our observations in §4.1, we find that models in Arabic struggle on high-frequency entities, where CBS is the highest for
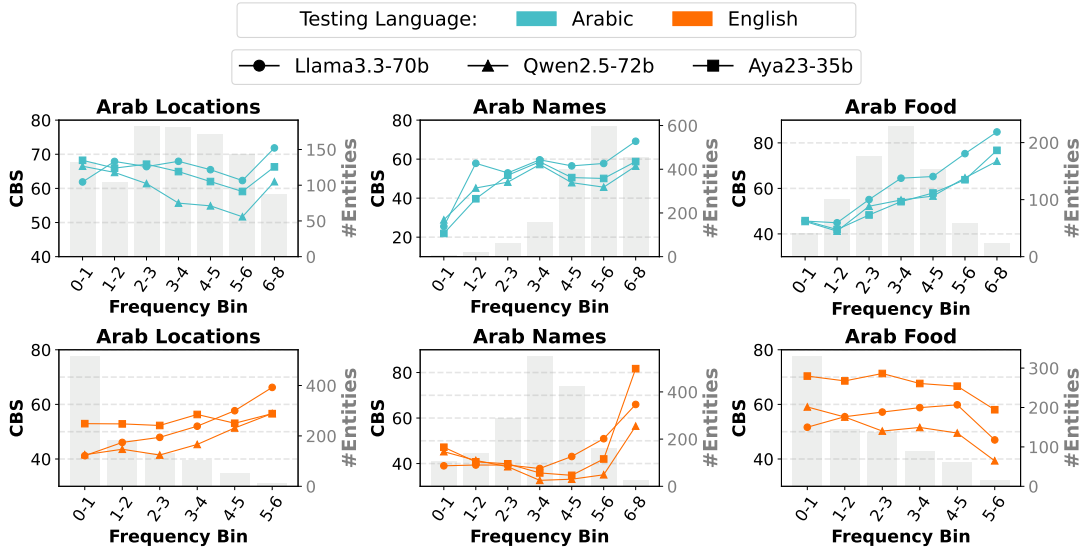
Figure 15: Average CBS (↓) of LLMs at text-infilling when tested in Arabic and English on locations, names, and food entities associated with Arab culture, stratified by their occurrence counts in the mC4 corpus (grouped into log10-spaced bins

all entity types. This indicates that entities that appear at very high frequencies will be assigned lower likelihood than Western entities given contexts grounded in Arab culture. We also observe a struggle on highly frequent entities when the model is tested in English on Locations and Names, which are rare cases where translations of Arabic entities exhibit polysemy. For example, this happens with the feminine name "آسيا" that is written as "*Asia*" in English, matching the name of the continent.

## D.2 The Impact of Word Polysemy

**NER.** Figure 16 and Figure 17 show the results $XLMR_{large}$ and ARBERT on NER of the top-100 Arab and Western locations as a function of the percentage of entities that match Arabic polysemous words. We find the same trend observed with Llama3.3-70b on extracted QA, where performance becomes poor on entities that exhibit word polysemy in Arabic (§4.2).

## D.3 The Impact of Tokenization

We report the performance of a variety of LMs on location entities as a function of how many tokens they are fragmented into in Figure 18. Similar to our observations in §4.4, we find that performance improves as entities are tokenized into more than a single token, and that LMs struggle with one-token entities that exhibit word polysemy. We also see that this issue gets worse as the number of Arabic tokens in an LM's vocabulary gets larger.
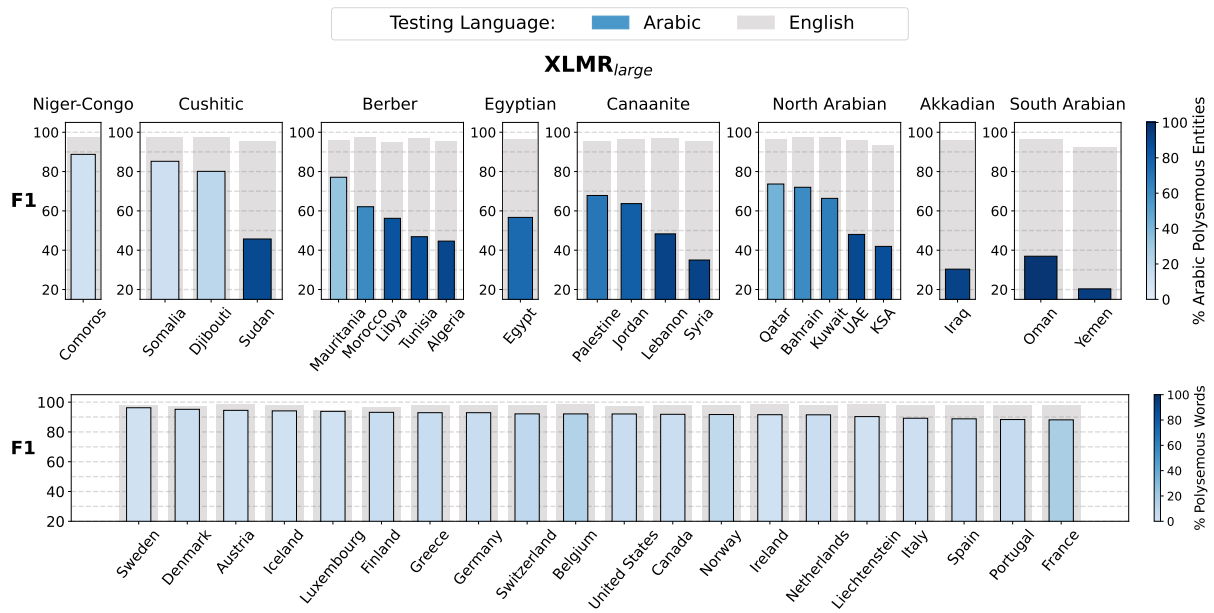
Figure 16: Average NER F1 of XLMR_large on the top-100 most frequent location entities in mC4 for each Arab country (top) and Western country (bottom) in CAMeL-2.
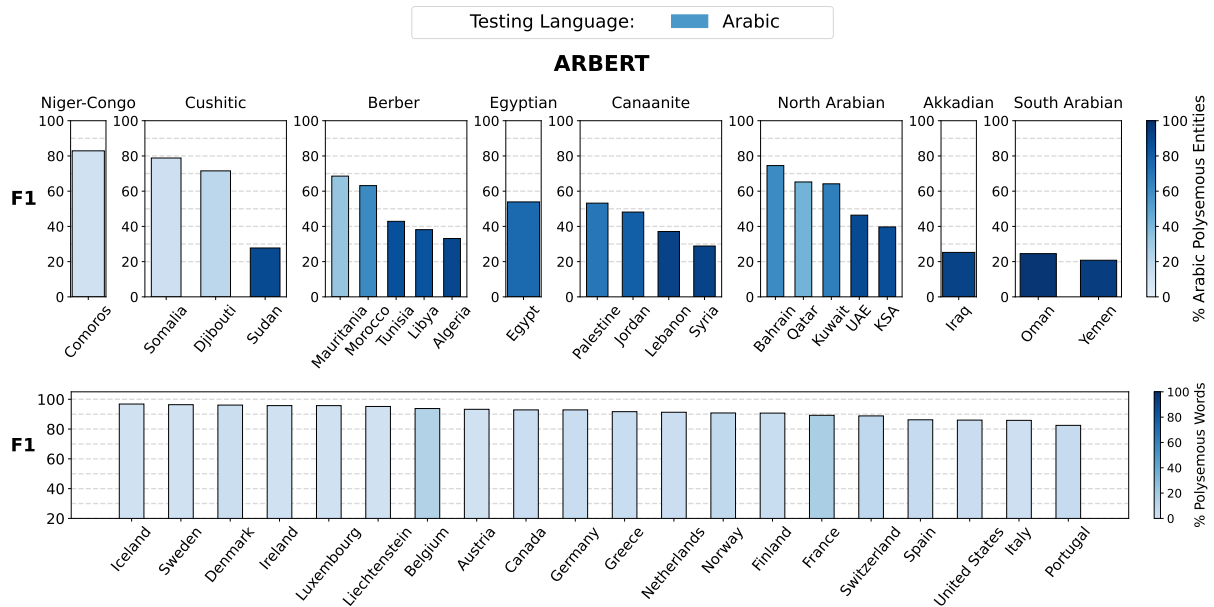


Figure 17: Average NER F1 of ARBERT on the top-100 most frequent location entities in mC4 for each Arab country (top) and Western country (bottom) in CAMeL-2.
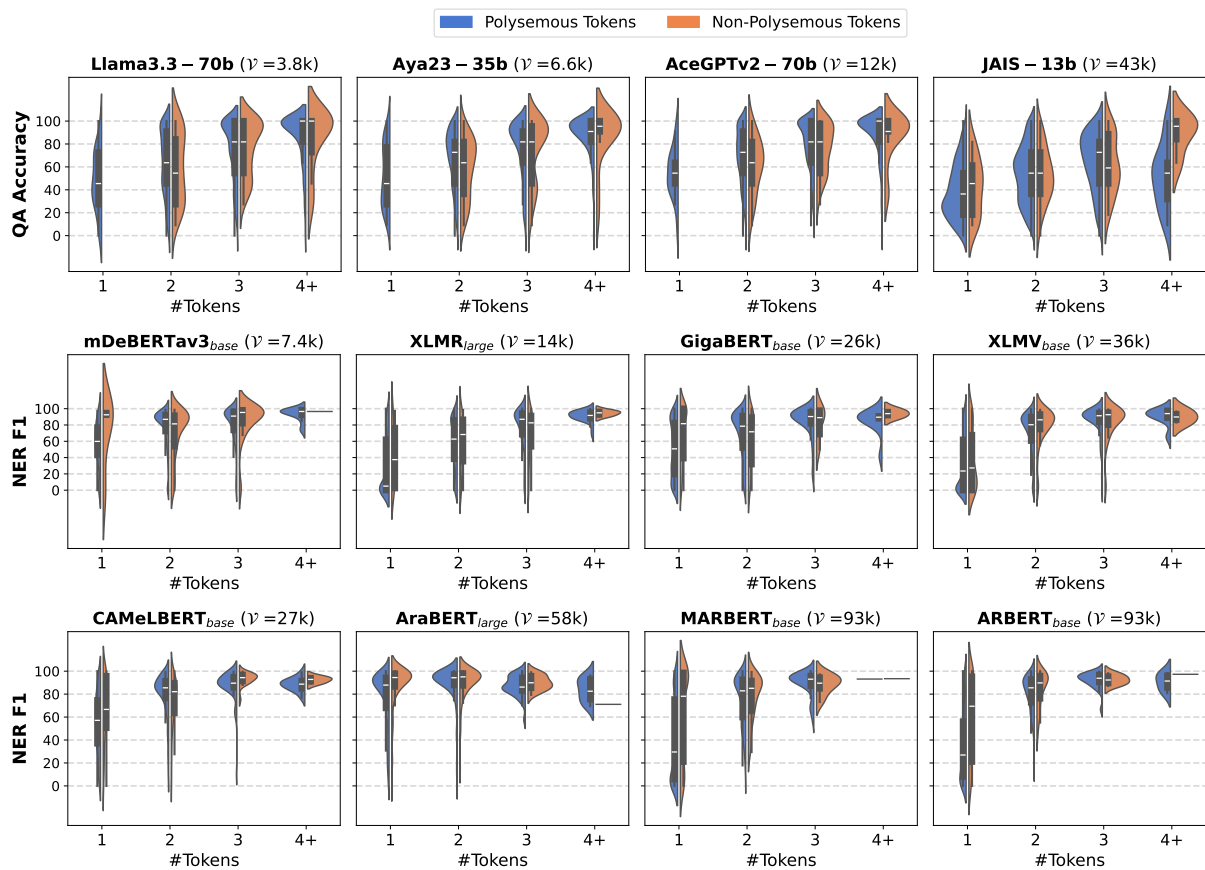
Figure 18: Performance distribution of several LMs on Arab location entities as a function of how many tokens they get tokenized into. Entities are separated based on whether tokens correspond to Arabic polysemous words. $\mathcal{V}$ represents the number of Arabic tokens in each LM's vocabulary.