Language Models' Factuality Depends on the Language of Inquiry

Anonymous ACL submission

Abstract

Multilingual language models (LMs) are expected to recall factual knowledge consistently across languages, yet they often fail to transfer 005 knowledge between languages even when they possess the correct information in one of the languages. For example, we find that an LM may correctly identify Rashed Al Shashai as being from Saudi Arabia when asked in Arabic, but consistently fails to do so when asked in English or Swahili. To systematically investigate this limitation, we introduce a benchmark of 10,000 country-related facts across 13 languages and propose three novel metrics-Factual Recall Score, Knowledge Transferability Score, and Cross-Lingual Factual Knowledge Transferability Score-to quantify factual recall and knowledge transferabil-018 ity in LMs across different languages. Our 019 results reveal fundamental weaknesses in today's state-of-the-art LMs, particularly in crosslingual generalization where models fail to transfer knowledge effectively across different languages, leading to inconsistent performance sensitive to the language used. Our findings emphasize the need for LMs to recognize language-specific factual reliability and lever-028 age the most trustworthy information across languages. We release our benchmark and evaluation framework to drive future research in multilingual knowledge transfer.

Introduction 1

007

011

017

Large Language Models (LLMs) are often perceived as vast knowledge reservoirs, capable of 034 recalling factual information across multiple languages (Wang et al., 2024). However, what if their knowledge is locked within linguistic boundaries and unable to be transferred across languages? Despite advancements in multilingual LMs such as Llama (Touvron et al., 2023a; Dubey et al., 2024), Gemma (Team et al., 2024a), DeepSeek (DeepSeek-AI et al., 2024), and Phi (Abdin et al., 042



Figure 1: Illustratation of the cross-lingual factual knowledge transferability issue across linguistic knowledge clouds in LMs. The model correctly recalls that Rashed Al Shashai is from Saudi Arabia when gueried in Arabic, but fails to retrieve this fact in English and Swahili, highlighting that factual knowledge is often stored in language-specific silos.

2024; Li et al., 2023), our study reveals a striking asymmetry in their factual recall across languages: consider the example in Figure 1, where an LM is tasked with a simple factual query: "Rashed Al Shashai is from which country?" When asked in Arabic, several state-of-the-art LMs correctly generate the response: "Saudi Arabia." However, when posed in English, Hindi, or Swahili, the same models fail to recall the fact. This example suggests that models can correctly retrieve country-specific facts in the language associated with that country but struggle to do so in others.

This raises a critical question—do these models truly internalize and transfer factual knowledge across languages, or do they merely encode isolated linguistic silos?

This limitation has significant implications for multilingual AI development and real-world applications. Many LM-based systems-such as retrieval-augmented generation (RAG) pipelines, multilingual search engines, and cross-lingual reasoning models-assume that factual knowledge is

043

044

066 071 090

065

091

100 101

103

105 106

108 109

110 111 112

113 114

115 116 consistently available and transferable across languages.

Our findings reveal that LMs often rely on language-specific memorization rather than true cross-lingual knowledge generalization. This over-reliance can introduce biases, inconsistencies, and reliability issues in multilingual AI applications (Chua et al., 2024).

To systematically analyze the factual inconsistencies, we introduce a carefully curated dataset comprising country-related facts translated into 13 languages. This benchmark evaluates LMs on multiple dimensions-factual recall, in-context recall, and counter-factual context adherence-across high-, medium-, and low-resource languages. This benchmark comprises of 802 instances for factual recall, 156 instances for In-context recall, and 1404 instances for counter-factual context adherence as shown in Table 1.

Factual recall assesses the LM's ability to recall country-specific facts consistently across multiple languages. We evaluate factual recall using three metrics: (a) Factual Recall Score (FRS): Measures how accurately a model recalls a fact in a given language, (b) *Knowledge Transferability Score (KTS)*: Quantifies how well factual knowledge is transferred across languages, and (c) Cross-Lingual Factual Knowledge Transferability (X-FaKT) Score: Combines the assessment of factual recall and cross-lingual transfer ability. FRS and KTS measure the effectiveness of cross-lingual knowledge transfer, and X-FaKT Score integrates factual recall with transferability to provide a robust measure of multilingual generalization. These metrics offer a more nuanced evaluation than a simple error rate, allowing for a deeper understanding of crosslingual generalization.

In-Context Recall (Machlab and Battle, 2024) measures the general performance of the models in multilingual contexts. Inspired by (Du et al., 2024), we also study how factual knowledge of models affects their performance in handling in-context tasks in the multilingual setting (Counterfactual Context Adherence). For this, we design a dataset where factual knowledge conflicts with in-context instructions.

Our experiments reveal that while LMs often retrieve factual information correctly in the language associated with the fact, they struggle to transfer this knowledge to other languages. We also found that the size of the LLM plays an important role in factuality and knowledge transferability. For example, the combined performance of LLama-3-70B in 117 factuality and knowledge transfer across languages 118 is markedly ($152\% \uparrow$ in *X*-FaKT Score) better than 119 Llama-3.2-1B. In addition, there is a marked dif-120 ference in these tasks when queries are asked in 121 high-resource languages ($46\% \uparrow$ in *X*-FaKT Score) 122 as compared to the case with low resources. This 123 finding exposes a critical limitation in current lan-124 guage models and their approach to multilingual 125 knowledge integration. Our findings also reveal 126 an interesting trade-off: LMs with stronger factual 127 recall often struggle with counterfactual adherence, 128 highlighting a key limitation in balancing factual 129 memory and contextual reasoning. In our experi-130 ments, we observed that the factual knowledge of 131 LMs could skew their judgments, leading to inaccu-132 rate evaluations. One has to be very careful when 133 designing the prompt and using LM as an evalua-134 tor. We highlight the importance of controlling the 135 evaluator's factual knowledge to ensure consistent 136 and effective evaluation. 137

We plan to release our code and data upon publication. These are made available for reviewing purposes in the supplementary material.

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

Related Work 2

Multilingual NLP and Factual Recall. Prior work on multilingual models, such as mBERT (Devlin et al., 2019), XLM-R (Conneau et al., 2020), mT5 (Xue et al., 2021) and BLOOM (Workshop et al., 2023), has shown that LMs trained on multilingual corpora exhibit varying performance across languages. Studies have also highlighted systematic biases in factual retrieval across different languages (Artetxe et al., 2020; Liu et al., 2020). While multilingual QA benchmarks such as XQuAD (Artetxe et al., 2020), MLQA (Lewis et al., 2020), and Ty-DiQA (Clark et al., 2020) assess factual consistency, they do not explicitly measure knowledge transfer within LMs. To address this gap, we introduce a benchmark designed to evaluate crosslingual factual knowledge transferability.

Cross-Lingual Knowledge Transfer in LMs. While research suggests that multilingual LMs exhibit zero-shot and few-shot generalization across languages (Nooralahzadeh et al., 2020; Pfeiffer et al., 2020), empirical studies indicate that this transfer is often asymmetric, favoring highresource languages (Hu et al., 2020). Most recent work has focused on cross-lingual transfer from high-resource languages to lower-resource ones

Task Type	# Examples
Factual Recall	802
In-context Recall	156
Counter-Factual Context Adherence	1404

Table 1: Number of examples per languages in our benchmark (§3).



Figure 2: Examples from our multilingual dataset illustrating three tasks. Factual Recall: LMs recall countryspecific facts better in native languages, as seen with Dharan's correct identification in Nepali but incorrect in English. Incontext Recall: Models struggle with contextual reasoning, showing regional bias when associating names with countries. Counter-Factual Context Adherence: When given counterfactual prompts about well-known figures, models rely on prior knowledge, affecting their ability to adhere to provided context.

(Zhao et al., 2024a,b). In contrast, we find that knowledge transfer is often lacking even from lowresource to high-resource languages. Furthermore, we show that recent LMs can correctly retrieve country-specific facts in the language associated with that country, regardless of the language's resource level.

Context Sensitivity and Counterfactual Reasoning. LMs can be susceptible to contextual cues, often overriding stored knowledge when presented with misleading information (Brown et al., 2020; Tirumala et al., 2022; Du et al., 2024). Counterfactual reasoning studies (Wu et al., 2023) show that models trained for high factual recall struggle with conflicting contextual instructions. While prior evaluations have been monolingual (Shwartz et al., 2020; Wang et al., 2020), our study extends these investigations into the multilingual domain, introducing in-context recall and counterfactual adherence tasks to analyze cross-lingual reasoning.

3 Dataset

168

169

170

172

173

174

176

177

178

180

182

183

184

185

189

We introduce a new multilingual dataset designed to evaluate three key capabilities of LMs: (a) *Fac*- *tual Recall*, (b) *In-context Recall*, and (c) *Counter-Factual Context Adherence*. The number of instances in our dataset is given in the Table 1. Given the multilingual nature of our study, we categorize languages based on their resource availability in existing LM training corpora:

High-resource: English, Chinese, French, Japanese.

Medium-resource: Hindi, Russian, Arabic, Greek.

Low-resource: Nepali, Ukrainian, Turkish, Swahili, Thai.

These languages correspond to countries strongly associated with their usage: the United States, China, France, Japan, India, Russia, Saudi Arabia, Greece, Nepal, Ukraine, Turkey, Kenya, and Thailand. Now, we describe our datasets in detail.

3.1 Factual Recall

This task evaluates an LM's ability to recall country-specific facts across multiple languages. For example, given the query, *In which country is Mumbai located?*, the model should correctly respond with *India* when asked in different languages.

To construct the dataset, we curated a diverse set of entities—including cities, artists, sports figures, landmarks, festivals, and politicians—for 13 selected countries. We then created standardized templates for factual queries and translated them into each language using the Google Translate API (Google, n.d.). All translations were manually verified and refined as needed with the assistance of ChatGPT. In total, our dataset consists of 805 unique factual questions, each available in 13 language versions.

3.2 In-Context Recall

The in-context recall task evaluates how effectively an LM utilizes contextual information to answer a question, ensuring that internal knowledge does not influence the model's output.

Building on the work of (Feng and Steinhardt, 2024), we constructed our dataset by focusing on common person names associated with each country. For each example, we sampled two names and paired them with two different countries, creating context-based prompts as shown in violet color in Figure 2. To enhance dataset efficiency, we intentionally avoided associating a name with its most commonly linked country within the example.

237

238

239

3.3 Counter-Factual Context Adherence

240

241

242

245

247

248

249

251

254

256

260

261

262

263

265

266

267

270

271

272

273

274

275

276

278

279

281

282

283

This task evaluates an LM's susceptibility to counterfactual information by assessing whether it adheres to the provided context when answering a question. Ideally, the model should rely solely on the given context, but in some cases, its internal knowledge may interfere or override it, leading to unintended responses (Du et al., 2024). To investigate this, we curated a list of well-known personalities strongly associated with specific countries and deliberately introduced counterfactual information into the context.

For the example given in Figure 2, if the model defaults to its internal knowledge and answers *United States*, it demonstrates a resistance to the contextual information. Conversely, if it follows the counterfactual context and answers *India*, it suggests a higher reliance on the provided context rather than pre-existing knowledge.

One might expect these models to perform nearperfectly on these tasks, as they are very simple. However, despite the simplicity of these tasks, the performance varies across languages and models.

4 Experiments

In this section, we discuss our experimental setup, metric formulation, and both quantitative and qualitative analyses. We present the results of our experiments evaluating LMs on our dataset across diverse multilingual tasks. These experiments assess how language and country-specific factual knowledge influence LMs responses in a multilingual setting. All experiments were conducted using the latest models, with Qwen-2.5-72B-Inst (Qwen et al., 2025) serving as the evaluator (Li et al., 2024).

4.1 Experimental Setup

Models We evaluated 14 models of varying sizes, trained on different compositions of multilingual data, and fine-tuned using various preference optimization strategies (Ouyang et al., 2022; Rafailov et al., 2024), for our multilingual study. These include Deepseek (DeepSeek-AI et al., 2024), Qwen (Yang et al., 2024), Gemma (Team et al., 2024b), and Llama (Touvron et al., 2023b) families. Further details of the models evaluated are given in Table A.1.

Compute Details All our experiments were conducted on a set of 4 NVIDIA A100 GPUs, each with 80GB of VRAM. We used Chat-GPT (OpenAI et al., 2024) for the dataset generation.

Evaluation To evaluate all models on the curated 289 datasets (Section 3), we used a temperature setting 290 of 0 and a maximum token limit of 128. Specifi-291 cally, we tested the models' performance on Fac-292 tual Recall and In-Context Recall across different 293 settings. For evaluation, we designed our metrics 294 and utilized Qwen-2.5-72B-Inst as the evaluator (Li et al., 2024), with a maximum token limit of 256 to support reasoning. Evaluation prompts are shown in Figures 10 and 11. 298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

327

328

329

330

331

332

333

4.2 Metric Definition and Formulation

This section introduces our carefully designed metrics to evaluate factual recall and knowledge transferability across languages in LMs. We propose two key metrics: the Factual Recall Score (FRS) and the Knowledge Transferability Score (KTS). To establish a common metric for evaluating the model's performance in our benchmark, we compute their harmonic mean, which is defined as the Cross-Lingual Factual Knowledge Transferability Score (X-FaKT), to ensure a balanced assessment while penalizing large disparities between them. Our metrics incorporate an inverse formulation with a correction factor to maintain a bounded range of [0, 1]. A higher error rate results in a lower metric value due to the inverse transformation, ensuring that better model performance corresponds to higher scores.

4.2.1 Associative vs. Non-Associative Knowledge

We categorize our dataset into two groups: associative and non-associative knowledge. The categorization is defined as follows: we consider 13 languages, each associated with a corresponding country (i.e., the *i*th language belongs to the *i*th country).

Associative = $\{Q\}$	\in	Questions	:	Q	\in
$Language_i \wedge output(Q$) =	Country $i \wedge i$	=	$j\}$	

Non-associative = $\{Q \in \text{Questions} : Q \in \text{Language}_i \land \text{output}(Q) = \text{Country}_j \land i \neq j\}$

We denote the mean error rate for a countryspecific fact asked in the language strongly associated with that country as $\mu_{assoc.}$, and the mean error rate for a country-specific fact asked in a language not associated with that country as $\mu_{non-assoc.}$.

4.2.2 Factual Recall Score (FRS)

334

335

337

338

339

340

341

342

347

354

361

367

371

372

374

375

377

Factual recall evaluates the model's ability to correctly retrieve both *associative* and *non-associative* knowledge. We define the Factual Recall Score (FRS) as:

$$FRS = \frac{3}{2} \left(\frac{1}{\mu_{\text{assoc.}} + \mu_{\text{non-assoc.}} + 1} - \frac{1}{3} \right) \quad (1)$$

- When both errors are zero ($\mu_{assoc.} = 0, \mu_{non-assoc.} = 0$), the model has a perfect factual recall, yielding an FRS score of 1.
- When both errors are high, the denominator increases, resulting in a lower FRS score closer to 0, indicating poor factual recall.

4.2.3 Knowledge Transferability Score (KTS)

Knowledge transferability quantifies how well a model maintains consistent factual knowledge across languages. We define the *Knowledge Transferability Score (KTS)* as:

$$KTS = 2\left(\frac{1}{|\mu_{\text{assoc.}} - \mu_{\text{non-assoc.}}| + 1} - \frac{1}{2}\right)$$
 (2)

where:

- $|\mu_{assoc.} \mu_{non-assoc.}|$ captures the absolute difference between associative and non-associative recall errors.
- When both errors are zero ($\mu_{assoc.} = 0, \mu_{non-assoc.} = 0$), there is perfect factual knowledge transfer, resulting in a KTS score of 1.
- When both errors are high but equal (e.g., μ_{assoc.} = 20, μ_{non-assoc.} = 20), KTS remains 1, indicating that while factual recall is poor, the model exhibits consistent errors across languages.
- When errors differ significantly (e.g., $\mu_{assoc.} = 20$, $\mu_{non-assoc.} = 2$ or vice versa), the absolute difference increases, leading to a lower KTS, highlighting a lack of knowledge transfer across languages.

4.2.4 Cross-Lingual Factual Knowledge Transferability Score (X-FAKT)

To ensure a balanced evaluation of factual recall and cross-lingual transferability, we compute their harmonic mean:

$$X \text{-FAKT} = 2 \times \frac{FRS \times KTS}{FRS + KTS}$$
(3)

Llama-3-70B Gemma-2-27B 5.2 22.7 18.0 20.4 13.6 23.1 13.5 10.1 6.6 24.8 15. Phi-4-14B 45.3 14.5 29.1 19.8 6.3 50.5 17.1 14.6 7.6 40.3 22.9 88.2 19.8 95.1 76.3 53.7 52.2 93.9 44.6 93.1 Phi-3-14B Gemma-2-9B 7.4 29.4 24.3 29.6 19.1 31.3 18.2 13.1 6.2 Llama-3-8B 7.0 35.2 28.8 27.9 20.7 33.9 17.2 16.7 25.2 Orca-2-7B 2.7 11.6 64.0 49.8 85.4 35.0 92.5 76.4 30.7 57.0 84. Deepseek-7B 6.2 15.8 20.0 83.4 94.5 95.3 81.8 56.4 66.6 94.9 5 61.3 18.6 Mistral-7B-v0.2 3.5 7.0 71.6 58.1 55.4 18.7 69.3 49. 82 C Phi-3.5-4B 7.4 14.0 72.1 64.2 88.4 69.0 76.3 88.7 51.9 76.7 94.4 80.5 96.4 67. Phi-3-4B 5.9 13.5 87.4 87.4 95.6 91.4 98.4 93.0 77.4 97.8 93.1 96.8 75.2 Llama-3.2-3B-24.9 13.1 59.7 57.7 57.0 40 64.3 57.2 60.2 Gemma-2-2B 4.5 13.8 53.0 43.6 61.7 39.0 60.6 4 75.8 Llama-3 2-1B 28.6 29.8 72.6 71.7 78.4 73.3 75.8 71.9 83.0 78.1 74.6 63 Column Mean 5 1 59 9 2 62.0 52.2 Arabic Greet HAST

Figure 3: Error rates for each model on the Factual Recall task. A clear pattern emerges, showing a decline in performance as we move from larger to smaller models (top to bottom) and from high-resource to low-resource languages (left to right).

• The harmonic mean penalizes large disparities between factual recall (FRS) and knowledge transferability (KTS), ensuring that both contribute meaningfully to the final score. 379

380

382

384

385

388

390

391

392

393

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

- If either FRS or KTS is significantly lower, the overall score remains low, discouraging models from excelling in one metric while performing poorly in the other.
- A high X-FAKT score indicates that the model is both factually accurate and consistent across multiple languages.

This formulation provides a holistic evaluation of factual knowledge retention and cross-lingual consistency, making it a robust metric for assessing multilingual model performance.

4.3 Quantitative Analysis

4.3.1 Performance on Factual Recall task

The error rate across different LMs (Figure 3) reveals a clear pattern in performance across languages and model sizes. Notably, all models demonstrate superior performance on highresource languages like English and French, with error rates consistently below 15% for most model variants. This performance gradually deteriorates as the model size decreases, with smaller models showing significantly higher error rates across all languages. However, an interesting observation emerges with languages like Swahili and Turkish, which despite being low-resource languages, exhibit relatively better performance with error rates comparable to mid-resource languages. This can be attributed to their use of Latin script, facilitating

378 where:

Model	$\mu_{assoc.}(\%)$	$\mid \mu_{non-assoc.}(\%)$	t-stat	p-value	FRS	KTS	X-FAKT
Llama-3-70B	2.36 ± 5.12	9.85 ± 10.54	2.52	0.01	0.835	0.862	0.848
Gemma-2-27B	4.23 ± 8.49	16.46 ± 17.07	2.54	0.01	0.742	0.783	0.762
Phi-4-14B	12.87 ± 16.51	30.15 ± 25.92	2.35	0.02	0.548	0.706	0.617
Phi-3-14B	25.09 ± 29.84	55.57 ± 36.24	2.93	< 0.01	0.330	0.535	0.408
Gemma-2-9B	4.98 ± 6.09	22.32 ± 21.37	2.90	< 0.01	0.677	0.705	0.691
Llama-3-8B	4.60 ± 7.54	25.77 ± 19.61	3.85	< 0.01	0.649	0.651	0.650
Orca-2-7B	31.95 ± 31.65	56.77 ± 32.99	2.60	0.01	0.295	0.603	0.396
DeepSeek-7b	31.49 ± 30.68	63.73 ± 36.29	3.09	< 0.01	0.268	0.514	0.353
Mistral-7B-v0.2	16.96 ± 15.65	45.25 ± 29.34	3.42	< 0.01	0.424	0.559	0.483
Phi-3.5-4B	41.85 ± 31.62	69.87 ± 31.23	3.09	< 0.01	0.208	0.563	0.304
Phi-3-4B	42.45 ± 30.99	77.95 ± 33.72	3.65	< 0.01	0.181	0.477	0.262
Llama-3.2-3B	24.10 ± 17.80	47.48 ± 26.80	3.07	< 0.01	0.375	0.620	0.467
Gemma-2-2B	9.97 ± 14.78	45.77 ± 31.30	4.06	< 0.01	0.463	0.473	0.468
Llama-3.2-1B	34.74 ± 22.32	65.96 ± 26.98	4.03	< 0.01	0.247	0.524	0.336

Table 2: Results of the t-test comparing associative and non-associative knowledge across models, alongside FRS, KTS, and X-FAKT scores. (A) Llama-3-70B achieves the best performance in both factual recall and knowledge transferability. (B) There is a statistically significant difference between the performance on associative queries (asked in a country's native language) and non-associative queries (asked in other languages).



Figure 4: This figure illustrates the model-wise comparison of X-FAKT scores grouped by language families. A clear trend emerges, showing that as the model size increases within a family, the X-FAKT score tends to increase.

better knowledge transfer from English.

411

412

413

414

415

416

417

418

419

420

421

422

423

424

A compelling pattern emerges when examining languages that share similar scripts, and strong correlations in model performance among languages that share similar scripts. For example, the error patterns for Hindi-Nepali and Russian-Ukrainian pairs show remarkable similarities, suggesting that the models effectively leverage shared scriptural characteristics during learning. These patterns indicate that script similarity plays a crucial role in the model's ability to generalize across languages, potentially offering insights into how these models transfer knowledge between different language pairs and scripts.

425 Knowledge Transferability Analysis: From Ta426 ble 2, Llama-3-70B emerges as the clear leader with
427 the highest X-FAKT score of 0.848, demonstrat428 ing superior balanced performance in both fac-

Language	$\mu_{assoc.}(\%)$	$\mu_{non-assoc.}(\%)$
High	3.83 ± 3.79	29.84 ± 27.47
Medium	26.73 ± 17.60	50.54 ± 21.20
Low	29.53 ± 16.19	53.91 ± 23.68

Table 3: Average mean and standard deviation for error rate across all models for each language group. Highresource languages exhibit lower error rates compared to low-resource languages.

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

tual recall (FRS = 0.835) and knowledge transferability (KTS = 0.862). This exceptional performance is supported by the lowest error rates $(\mu_{assoc.} = 2.36\%, \mu_{non-assoc.} = 9.85\%)$, suggesting that larger model sizes generally correlate with better cross-lingual factual knowledge handling. Despite similar model sizes, significant performance variations exist between different architectures. For example, Gemma-2-9B (X-FAKT: 0.691) substantially outperforms Mistral-7B-v0.2 (X-FAKT: 0.483), suggesting that architecture design and training methodology play crucial roles beyond mere parameter count. As illustrated in Figure 4, the X-FAKT scores exhibit a clear upward trend with increasing model size within each language family. This suggests that larger models generally achieve better factual consistency, highlighting the impact of scale on model performance. These findings provide valuable insights into the current state of cross-lingual factual knowledge in LMs and highlight areas for future improvement, particularly in reducing the performance gap between associative and non-associative knowledge retrieval.

Associative vs. Non-associative performance: We analyze the performance of various models on these two subsets of data and report the results in the Table 2. For all models, the t-statistic and pvalue indicate that the differences between associative and non-associative categories are statistically significant (p-value less than 0.05).

Performance comparison across language groups: In this study, we categorize languages into three groups based on their availability and coverage in the dataset: **High**, **Medium**, and **Low**, as defined in Section 3. From the results shown in Table 3, we observe a clear trend across language groups. Specifically, high resouce languages exhibit the lowest average error rates, particularly in the associative category, where models make fewer mistakes ($\mu_{assoc.} = 3.83\%$). However, for non-associative questions, the error rate 470rises significantly ($\mu_{non-assoc.} = 29.84\%$), indicat-471ing that models struggle more when dealing with472non-associative samples in these languages. The473error rate increases while moving from high to low-474resource languages.

475 4.3.2 Performance on In-Context Recall task

476

477

478

479

480

481

482

483

484

485

486

487

488

489

491

492

493

494

495

496

497

498

499

501

503

505

506

507

508

509

511

512

513

514 515

516

517

518

519

Figure 12 demonstrates the incorrectness rate for the in-context recall capabilities of different LMs. Despite being a simple task, certain models such as DeepSeek-7B, Orca-2-7B, Phi-3-4B, Llama-3.2-1B, and Mistral-7B-v0.2 perform poorly across multiple languages. This suggests that these models struggle to effectively utilize contextual information when generating outputs. Interestingly, even for languages like Swahili and Turkish, which showed better scores in the Factual Recall task, models demonstrate poor performance on this context-dependent task. This stark contrast suggests that the benefits of Latin script-based knowledge transfer observed in the Factual Recall task do not extend to in-context learning scenarios, where performance depends primarily on the model's ability to process and utilize contextual information.

> As mentioned in the dataset section, we intentionally paired cross-entities as context. This setup appears to induce a regional bias, which negatively impacts model performance. The structured entitycontext pairing in the dataset may have led to spurious correlations (Yang et al., 2023; Ye et al., 2024), reducing model accuracy in in-context recall tasks. Some models struggle to effectively leverage contextual information, revealing potential weaknesses in their retrieval and in-context learning mechanisms.

4.3.3 Performance on Counter-Factual Context Adherence task

Figure 5 illustrates the error rates of LMs in the Counterfactual Context Adherence task. Notably, Latin-script languages (English, French, Swahili, and Turkish), which performed well in factual recall tasks, exhibited significantly higher error rates in counterfactual adherence. This suggests a fundamental trade-off in the models' capabilities: their strength in accurately retrieving factual information appears to come at the expense of their ability to maintain adherence to counterfactual contexts. This inverse relationship raises important questions about the inherent limitations and tradeoffs in LMs' learning mechanisms, particularly in how they balance factual knowledge with hypothet-



Figure 5: Error rate for each model on Counter-Factual Context Adherence task. Models show high error rates in high resource languages such as English and French where they have high factual recall.



Figure 6: Mistral-7B-v0.2 output when prompted with the given context in English. This model generation shows how spurious correlation leads to in-context recall failures

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

ical reasoning.

4.4 Qualitative Analysis

Spurious correlation leads to in-context recall failures. We observe that some models tend to associate names with cultural origins, even when contextual evidence contradicts this assumption. Figures 6 demonstrate the model response when prompted Mistral-7B-v0.2 with the contextual understanding-based question in English.

Despite the explicit context stating that *Li Wei* resides in *Russia*, the model disregards this information and defaults to cultural associations. This behavior reveals a limitation in integrating contextual evidence when making country-specific inferences.

Models favor factual knowledge over context. We also observed that some models prioritize their internal factual knowledge over contextual information when responding to questions about wellknown personalities. Figures 7 demonstrate the model response when prompted Llama-3-70B with the factual retrieval query in English. Instruction: Answer the question based on the given fact. Fact: George Washington lives in the country of India. Question: In which country does George Washington live?

Model Output: I apologize, but the fact provided is incorrect. George Washington, the first President of the United States, did not live in India. He lived in the United States of America.

Figure 7: Llama-3-70B output when prompted with a counter-factual context adherence query in English. This shows LMs favour internal knowledge over contextual understanding.

Instruction: In which country is Dijon located?

Model Output: Dijon is located in France. It is a city in the Bourgogne-Franche-Comté region, and it is the capital of the Côte-d'Or department.

Figure 8: Llama-3-70B output when prompted with a factual recall query in English

In this case, despite being explicitly told that 'George Washington' lived in 'India', the model relied on its factual knowledge, correcting the given fact and asserting that 'George Washington' lived in the 'United States'. This response demonstrates the model's strong reliance on factual accuracy, rather than adapting to the context provided. It suggests that when it comes to well-known historical figures, models may prioritize prior knowledge over the specific context they are given.

Linguistic variability in word interpretation. LMs can interpret words differently depending on the language. Figures 8 and 9 demonstrate the model responses when prompted Llama-3-70B with the same queries but in different languages. This highlights challenges in multilingual consistency, where the model misinterprets '*Dijon*' as '*De Janeiro*' in Hindi, revealing inconsistencies in cross-lingual factual retrieval.

561Challenges with using LMs as evaluators. We562used a zero-shot prompt with Llama-3-70B as an563evaluator and found that its inherent factual knowl-564edge can skew assessments. For example, when565evaluating a Gemma-2-27B response to the counter-566factual context task—"Catherine the Great lives567in India"—the evaluator corrected it, asserting that568she lived in "Russia", despite the provided ground569truth. This bias highlights the need to control evaluators' factual knowledge to ensure consistent evaluation.



Figure 9: Llama-3-70B output when prompted with a factual recall query in Hindi. In Hindi, it misinterprets understanding of a French word.

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

597

598

599

600

601

602

603

604

605

606

607

608

5 Conclusions

Our study reveals a critical limitation in multilingual LMs: their inability to consistently transfer factual knowledge across languages. Our benchmark provides a standardized framework to evaluate both current and future LMs on their factual consistency and cross-lingual generalization, enabling a more systematic comparison of their capabilities. Moreover, it can serve as a valuable resource to promote research in interpretability by helping analyze how and where factual knowledge is stored and retrieved across languages, fostering a deeper understanding of LM internals. We emphasize the need for AI systems with internal awareness of their language-specific strengths and weaknesses-a concept we term calibrated multilingualism. Under this paradigm, a model would autonomously leverage the most reliable internal representations for any given multilingual query.

We also find that LMs, when used as evaluators, are biased by their internal factual knowledge, which may not align with the intended input-outputground-truth context. This underscores the need to control the evaluator's factual knowledge for more reliable assessments. Ultimately, enabling AI to cross-generalize across languages is crucial for inclusive and equitable technology, ensuring language is no barrier to reliable knowledge access.

6 Limitations

Our study provides valuable insights into crosslingual knowledge transfer in LMs but has some limitations. First, our benchmark, though comprehensive in country-related facts, covers only 13 languages, limiting its representation of diverse linguistic families. Second, we evaluated only opensource LMs, excluding proprietary models that may exhibit different transfer patterns. Third, our fact

542

collection used a standardized template for consistency, which may not reflect the diversity of
real-world queries. Lastly, our focus on countryrelated facts means our findings may not generalize
to other domains like science, history, or culture.

7 Ethics Statement

614

This research is conducted with a strong commit-615 616 ment to ethical principles, ensuring data privacy and consent by using publicly available informa-617 tion and adhering to data protection regulations. 618 We acknowledge potential biases in multilingual 619 language models and aim to highlight and address 620 these through our benchmark. Transparency and 621 reproducibility are promoted by making our dataset 622 and evaluation framework publicly available. Our 623 624 research aligns with the broader goals of fairness, transparency, and social responsibility. 625

References

626

631

632

637

638

644

647

651

652

653

654

655

657

667

670

672 673

674

675

676

677

678

679

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.
- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. On the cross-lingual transferability of monolingual representations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4623–4637, Online. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901.
- Lynn Chua, Badih Ghazi, Yangsibo Huang, Pritish Kamath, Ravi Kumar, Pasin Manurangsi, Amer Sinha, Chulin Xie, and Chiyuan Zhang. 2024. Crosslingual capabilities and knowledge barriers in multilingual large language models. *arXiv preprint arXiv:2406.16135*.
- Jonathan H. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. TyDi QA: A benchmark for information-seeking question answering in typologically diverse languages. *Transactions of the Association for Computational Linguistics*, 8:454–470.
- 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.
- DeepSeek-AI, :, Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, Huazuo Gao, Kaige Gao, Wenjun Gao, Ruiqi Ge, Kang Guan, Daya Guo, Jianzhong Guo, Guangbo Hao, Zhewen Hao, Ying He, Wenjie Hu, Panpan Huang, Erhang Li, Guowei Li, Jiashi Li, Yao Li, Y. K. Li, Wenfeng Liang, Fangyun Lin, A. X. Liu, Bo Liu, Wen Liu, Xiaodong Liu, Xin Liu, Yiyuan Liu, Haoyu Lu, Shanghao Lu, Fuli Luo, Shirong Ma, Xiaotao Nie, Tian Pei, Yishi Piao, Junjie Qiu, Hui Qu, Tongzheng Ren, Zehui Ren, Chong Ruan, Zhangli Sha, Zhihong Shao, Junxiao Song, Xuecheng Su, Jingxiang Sun, Yaofeng Sun, Minghui Tang, Bingxuan Wang, Peiyi Wang, Shiyu Wang, Yaohui Wang, Yongji Wang, Tong Wu, Y. Wu, Xin Xie, Zhenda Xie, Ziwei Xie, Yiliang Xiong, Hanwei Xu, R. X. Xu, Yanhong Xu, Dejian Yang, Yuxiang You, Shuiping

Yu, Xingkai Yu, B. Zhang, Haowei Zhang, Lecong Zhang, Liyue Zhang, Mingchuan Zhang, Minghua Zhang, Wentao Zhang, Yichao Zhang, Chenggang Zhao, Yao Zhao, Shangyan Zhou, Shunfeng Zhou, Qihao Zhu, and Yuheng Zou. 2024. Deepseek llm: Scaling open-source language models with longtermism. *Preprint*, arXiv:2401.02954. 684

685

687

688

689

690

691

692

693

694

695

696

697

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

718

719

720

721

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Kevin Du, Vésteinn Snæbjarnarson, Niklas Stoehr, Jennifer C White, Aaron Schein, and Ryan Cotterell. 2024. Context versus prior knowledge in language models. arXiv preprint arXiv:2404.04633.
- 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*.
- Jiahai Feng and Jacob Steinhardt. 2024. How do language models bind entities in context? *Preprint*, arXiv:2310.17191.
- Google. n.d. Google translate. Accessed: 2025-02-16.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalization. *Preprint*, arXiv:2003.11080.
- Patrick Lewis, Barlas Oguz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. 2020. MLQA: Evaluating cross-lingual extractive question answering. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7315– 7330, Online. Association for Computational Linguistics.
- Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. 2024. Llms-as-judges: A comprehensive survey on llm-based evaluation methods. *Preprint*, arXiv:2412.05579.
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023. Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*.
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pretraining for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.

Daniel Machlab and Rick Battle. 2024. Llm incontext recall is prompt dependent. *Preprint*, arXiv:2404.08865.

740

741

742

743

745

746

747

748 749

751

753

754

755

756

759

760

761

762 763

764

765

766

767

770

771

772 773

774

775

776

778

779

784

786

790 791

792

793 794

796

797

801

- Farhad Nooralahzadeh, Giannis Bekoulis, Johannes Bjerva, and Isabelle Augenstein. 2020. Zero-shot cross-lingual transfer with meta learning. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 4547–4562, Online. Association for Computational Linguistics.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,

Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report. Preprint, arXiv:2303.08774.

802

803

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

827

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. 2020. MAD-X: An Adapter-Based Framework for Multi-Task Cross-Lingual Transfer. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 7654–7673, Online. Association for Computational Linguistics.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2025. Qwen2.5 technical report. *Preprint*, arXiv:2412.15115.

973

974

975

976

977

978

979

980

981

982

983

923

924

- 867
- 870
- 872 873
- 874
- 876

879

- 883 895

887 890

882

892 893

897

900

901 902

903

904

905

906

907

908

909

910

911

912

913

914 915

916

917

918

919

920

921

- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. Preprint, arXiv:2305.18290.
 - Vered Shwartz, Peter West, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Unsupervised commonsense question answering with self-talk. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 4615-4629, Online. Association for Computational Linguistics.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024a. Gemma: Open models based on gemini research and technology. arXiv preprint arXiv:2403.08295.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. 2024b. Gemma: Open models based on gemini research and technology. Preprint, arXiv:2403.08295.
 - Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. 2022. Memorization without overfitting: Analyzing the training dynamics of large language models. Advances in Neural Information Processing Systems, 35:38274–38290.

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023b. Llama: Open and efficient foundation language models. Preprint, arXiv:2302.13971.
- Chenguang Wang, Xiao Liu, and Dawn Song. 2020. Language models are open knowledge graphs. arXiv preprint arXiv:2010.11967.
- Mengru Wang, Yunzhi Yao, Ziwen Xu, Shuofei Qiao, Shumin Deng, Peng Wang, Xiang Chen, Jia-Chen Gu, Yong Jiang, Pengjun Xie, Fei Huang, Huajun Chen, and Ningyu Zhang. 2024. Knowledge mechanisms in large language models: A survey and perspective. In Findings of the Association for Computational Linguistics: EMNLP 2024, pages 7097-7135, Miami, Florida, USA. Association for Computational Linguistics.
- BigScience Workshop, :, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, Dragomir Radev, Eduardo González Ponferrada, Efrat Levkovizh, Ethan Kim, Eyal Bar Natan, Francesco De Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady Elsahar, Hamza Benyamina, Hieu Tran, Ian Yu, Idris Abdulmumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jörg Frohberg, Joseph Tobing, Joydeep Bhattacharjee, Khalid Almubarak, Kimbo Chen, Kyle Lo, Leandro Von Werra, Leon Weber, Long Phan, Loubna Ben allal, Ludovic Tanguy, Manan Dey, Manuel Romero Muñoz, Maraim Masoud, María Grandury, Mario Šaško, Max Huang, Maximin Coavoux, Mayank Singh, Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad A. Jauhar, Mustafa Ghaleb, Nishant Subramani, Nora Kassner, Nurulaqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, Paulo Villegas, Peter Henderson, Pierre Colombo, Priscilla Amuok, Quentin Lhoest, Rheza Harliman, Rishi Bommasani,

Roberto Luis López, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, Shayne Longpre, Somaieh Nikpoor, Stanislav Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan Thrush, Valentin Danchev, Vassilina Nikoulina, Veronika Laippala, Violette Lepercq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heinzerling, Chenglei Si, Davut Emre Taşar, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, Andrea Santilli, Antoine Chaffin, Arnaud Stiegler, Debajyoti Datta, Eliza Szczechla, Gunjan Chhablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, Lintang Sutawika, M Saiful Bari, Maged S. Al-shaibani, Matteo Manica, Nihal Nayak, Ryan Teehan, Samuel Albanie, Sheng Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Fevry, Trishala Neeraj, Urmish Thakker, Vikas Raunak, Xiangru Tang, Zheng-Xin Yong, Zhiqing Sun, Shaked Brody, Yallow Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, Jaesung Tae, Jason Phang, Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared Casper, Jeff Rasley, Max Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myriam Peyrounette, Nicolas Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre Cornette, Pierre François Lavallée, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Requena, Suraj Patil, Tim Dettmers, Ahmed Baruwa, Amanpreet Singh, Anastasia Cheveleva, Anne-Laure Ligozat, Arjun Subramonian, Aurélie Névéol, Charles Lovering, Dan Garrette, Deepak Tunuguntla, Ehud Reiter, Ekaterina Taktasheva, Ekaterina Voloshina, Eli Bogdanov, Genta Indra Winata, Hailey Schoelkopf, Jan-Christoph Kalo, Jekaterina Novikova, Jessica Zosa Forde, Jordan Clive, Jungo Kasai, Ken Kawamura, Liam Hazan, Marine Carpuat, Miruna Clinciu, Najoung Kim, Newton Cheng, Oleg Serikov, Omer Antverg, Oskar van der Wal, Rui Zhang, Ruochen Zhang, Sebastian Gehrmann, Shachar Mirkin, Shani Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun, Tomasz Limisiewicz, Verena Rieser, Vitaly Protasov, Vladislav Mikhailov, Yada Pruksachatkun, Yonatan Belinkov, Zachary Bamberger, Zdeněk Kasner, Alice Rueda, Amanda Pestana, Amir Feizpour, Ammar Khan, Amy Faranak, Ana Santos, Anthony Hevia, Antigona Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh HajiHosseini, Bahareh Behroozi, Benjamin Ajibade, Bharat Saxena, Carlos Muñoz Ferrandis, Daniel McDuff, Danish Contractor, David Lansky, Davis David, Douwe Kiela, Duong A. Nguyen, Edward Tan, Emi Baylor, Ezinwanne Ozoani, Fatima Mirza, Frankline Ononiwu, Habib Rezanejad, Hessie Jones, Indrani Bhattacharya, Irene Solaiman, Irina Sedenko, Isar Nejadgholi, Jesse Passmore, Josh Seltzer, Julio Bonis Sanz, Livia Dutra, Mairon Samagaio, Maraim Elbadri, Margot Mieskes, Marissa Gerchick, Martha Akinlolu, Michael McKenna, Mike Qiu, Muhammed Ghauri, Mykola Burynok, Nafis Abrar, Nazneen Rajani, Nour Elkott, Nour Fahmy, Olanrewaju Samuel, Ran An, Rasmus Kromann, Ryan Hao, Samira Al-

985

995

999

1002

1004

1005

1006

1007

1009

1011

1012

1014

1015

1016

1017 1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

izadeh, Sarmad Shubber, Silas Wang, Sourav Roy, 1048 Sylvain Viguier, Thanh Le, Tobi Oyebade, Trieu Le, 1049 Yoyo Yang, Zach Nguyen, Abhinav Ramesh Kashyap, Alfredo Palasciano, Alison Callahan, Anima Shukla, 1051 Antonio Miranda-Escalada, Ayush Singh, Benjamin Beilharz, Bo Wang, Caio Brito, Chenxi Zhou, Chirag Jain, Chuxin Xu, Clémentine Fourrier, Daniel León Periñán, Daniel Molano, Dian Yu, Enrique Manjava-1055 cas, Fabio Barth, Florian Fuhrimann, Gabriel Altay, 1056 Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec, 1057 Imane Bello, Ishani Dash, Jihyun Kang, John Giorgi, 1058 Jonas Golde, Jose David Posada, Karthik Ranga-1059 sai Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa 1060 Shinzato, Madeleine Hahn de Bykhovetz, Maiko Takeuchi, Marc Pàmies, Maria A Castillo, Mari-1062 anna Nezhurina, Mario Sänger, Matthias Samwald, 1063 Michael Cullan, Michael Weinberg, Michiel De Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank, Myungsun Kang, Natasha Seelam, Nathan Dahlberg, 1066 Nicholas Michio Broad, Nikolaus Muellner, Pascale Fung, Patrick Haller, Ramya Chandrasekhar, Renata 1068 Eisenberg, Robert Martin, Rodrigo Canalli, Rosaline 1069 Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda, Shlok S Deshmukh, Shubhanshu Mishra, Sid Ki-1071 blawi, Simon Ott, Sinee Sang-aroonsiri, Srishti Kumar, Stefan Schweter, Sushil Bharati, Tanmay Laud, 1073 Théo Gigant, Tomoya Kainuma, Wojciech Kusa, Yanis Labrak, Yash Shailesh Bajaj, Yash Venkatraman, 1075 Yifan Xu, Yingxin Xu, Yu Xu, Zhe Tan, Zhongli 1076 Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and 1077 Thomas Wolf. 2023. Bloom: A 176b-parameter 1078 open-access multilingual language model. Preprint, 1079 arXiv:2211.05100. 1080

Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek, Boyuan Chen, Bailin Wang, Najoung Kim, Jacob Andreas, and Yoon Kim. 2023. Reasoning or reciting? exploring the capabilities and limitations of language models through counterfactual tasks. *arXiv preprint arXiv:2307.02477*. 1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

- 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, Online. Association for Computational Linguistics.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu

1109Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen21110technical report. *arXiv preprint arXiv:2407.10671*.

1111

1112 1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124 1125

1126

1127

1128

- Yu Yang, Besmira Nushi, Hamid Palangi, and Baharan Mirzasoleiman. 2023. Mitigating spurious correlations in multi-modal models during fine-tuning. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 39365–39379. PMLR.
- Wenqian Ye, Guangtao Zheng, Xu Cao, Yunsheng Ma, and Aidong Zhang. 2024. Spurious correlations in machine learning: A survey. *Preprint*, arXiv:2402.12715.
- Jun Zhao, Zhihao Zhang, Luhui Gao, Qi Zhang, Tao Gui, and Xuanjing Huang. 2024a. Llama beyond english: An empirical study on language capability transfer. *arXiv preprint arXiv:2401.01055*.
- Xin Zhao, Naoki Yoshinaga, and Daisuke Oba. 2024b. Tracing the roots of facts in multilingual language models: Independent, shared, and transferred knowledge. *arXiv preprint arXiv:2403.05189*.

A APPENDIX

[Instruction] Please act as an impartial judge and evaluate the quality of the response provided by an AI assistant to the user question displayed below. The assistant should be penalized if it doesn't answer the question based on the given reference. For this evaluation, you should primarily consider the following criteria:
Score 1: The Assistant's answer is in the same language as the reference and aligns with the reference. Score 2: The Assistant's answer is in English and aligns with the reference. Score 3: The Assistant's answer is in any other language and aligns with the reference. Score 4: The Assistant's answer is unrelated to the reference. Only respond with a numerical score.
Begin your evaluation by providing a short explanation. Be as objective as possible. After providing your explanation, you must rate the response on a scale of 1 to 4 by strictly following this format: "[[rating]]", for example:
[Reference] India
[Question] George Washington lives in the country of India. In which country does George Washington live?
[The Start of Assistant's Answer] George Washington lives in **United States**. [The End of Assistant's Answer]
Response: The assistant's response is in same language but does not align with the reference provided.\n\nRating: [[4]]
[Reference] {answers}
[Question] {question}
[The Start of Assistant's Answer] {predicted}
[The End of Assistant's Answer]

Figure 10: Prompt used for evaluation of Factual Recall and In-Context Recall tasks.

[Instruction]

Please act as an impartial judge and evaluate the quality of the response provided by an AI assistant to the user question displayed below. The AI assistant is given a fact and the assistant has to answer the question based on the given fact. The assistant should be penalized if it doesn't answer the question based on the fact. For this evaluation, you should primarily consider the following criteria:
Score 1: The Assistant's answer is in the same language as the reference and aligns with the reference. Score 2: The Assistant's answer is in English and aligns with the reference. Score 3: The Assistant's answer is in any other language and aligns with the reference. Score 4: The Assistant's answer is unrelated to the reference.
Only respond with a numerical score.
Begin your evaluation by providing a short explanation. Be as objective as possible. After providing your explanation, you must rate the response on a scale of 1 to 4 by strictly following this format: "[[rating]]", for example:
[Reference] India
[Question] Instruction: Answer the question based on the given fact. Fact: George Washington lives in the country of India. Question: In which country does George Washington live?
[The Start of Assistant's Answer] The given fact is incorrect, George Washington lives in **United States**. [The End of Assistant's Answer]
Response: The assistant's response is in English and does not align with the reference provided.\n\nRating: [[4]]
[Reference] {answers}
[Question] {question}
[The Start of Assistant's Answer] {predicted}
[The End of Assistant's Answer]

Figure 11: Prompt used for evaluation of Counter-Factual Context Adherence task.

Model	Model Size	Training	Languages	Context	Vocab	Post-Training	Key
	& Architecture	Data	Supported	Length	Size	Strategies	Features
Llama-3-70B	70B	15T tokens	EN, DE, FR, IT,	8K	128K	SFT, RS, DPO	GQA, 8 heads,
	L=80, H=64	Multi-lingual	PT, HI, ES, TH				RoPE embeddings
Gemma-2-27B	27B	13T tokens	Primarily	8K	256K	SFT, RLHF	Local-global attention,
		Web, Code, Math	English				Knowledge distillation
Phi-4-14B	14B	400B synthetic	DE, ES, FR, PT,	16K	100K	SFT, RS,	Full attention over
		+ 10T web	IT, HI, JA			DPO	4K context
Phi-3-14B	14B	4.8T tokens	10% multilingual	128K	32K	SFT, DPO	Reasoning focus,
			data				Multi-lingual support
Gemma-2-9B	9B	8T tokens	Primarily	8K	256K	SFT, RLHF	GQA, RoPE,
			English				Knowledge distillation
Llama-3-8B	8B	15T tokens	EN, DE, FR, IT,	8K	128K	SFT, RS,	GQA, RoPE,
	L=32, H=32	Multi-lingual	PT, HI, ES, TH			DPO	32 heads
Orca-2-7B	7B	Based on	Based on	4K	32K	Single-turn	Enhanced reasoning
	L=32, H=32	Llama 2	Llama 2			SFT	abilities
DeepSeek-7B	7B	2T tokens	English	4K	102K	SFT, DPO	English & Chinese
	L=30, H=32		& Chinese				focus
Mistral-7B-v0.2	7B	Open Web	Open Web	32K	32K	SFT	GQA, Sliding window
	L=32, H=32		languages				attention
Phi-3.5-4B	3.8B	3.4T tokens	23 languages incl.	128K	32K	SFT, DPO	Multi-lingual
	L=32, H=32	Multi-lingual	AR, ZH, CS, NL,				support
Phi-3-4B	3.8B	4.9T tokens	Similar to	128K	32K	SFT, DPO	Diverse domain
			Phi-3.5-mini				coverage
Llama-3.2-3B	3B	9T tokens	EN, DE, FR, IT,	128K	128K	SFT, RS,	GQA, Knowledge
		Multi-lingual	PT, HI, ES, TH			DPO	distillation
Gemma-2-2B	2B	2T tokens	Primarily	8K	256K	SFT, RLHF	GQA, RoPE,
			English				Knowledge distillation
Llama-3.2-1B	1B	9T tokens	EN, DE, FR, IT,	128K	128K	SFT, RS,	GQA, Multi-lingual
		Multi-lingual	PT, HI, ES, TH			DPO	dialogue focus

Table A.1: **Model Specifications**: Comparison of LMs across key architectural and training characteristics. Languages abbreviated as: EN=English, DE=German, FR=French, IT=Italian, PT=Portuguese, HI=Hindi, ES=Spanish, TH=Thai, AR=Arabic, ZH=Chinese, CS=Czech, NL=Dutch, JA=Japanese. L=Layers, H=Heads, SFT=Supervised Fine-Tuning, RS=Rejection Sampling, DPO=Direct Preference Optimization, GQA=Grouped-Query Attention, RoPE=Rotary Position Embedding.

Question	Model Prediction	Ground Truth	LLM Evaluation Output
New York iko katika nchi gani?	New York iko nchini	Marekani	Response: The assistant's
	Marekani.		response is in the same lan-
			guage as the reference and
			aligns with the reference
			provided. Rating: [[1]]
Talimat: Verilen gerçeğe göre	George Washington,	Amerika Bir-	The assistant's response is
soruyu cevaplayın. Gerçek:	**Amerika Birleşik De-	leşik Devletleri	in the same language as
George Washington, Amerika	vletleri**'nde yaşıyor.		the reference and aligns
Birleşik Devletleri ülkesinde			with the reference. Rating:
yaşıyor. Soru: George Wash-			[[1]]
ington hangi ülkede yaşıyor?			
David lives in China and Li Wei	David lives in China.	David	Response: The assistant's
lives in United States. Who			response is in the same lan-
lives in China?			guage and aligns with the
			reference provided. Rat-
			ing: [[1]]
Dans quel pays se trouve	Sourate n'est pas un pays, mais	Inde	The assistant's response is
Sourate ?	le nom d'une **ville** située		in French, which aligns
	en **Algérie**.		with the reference lan-
			guage, but the information
			provided does not align
			with the reference. The
			assistant incorrectly states
			that Sourate is a city in
			Algeria, while the refer-
			ence is about India. Rat-
			ing: [[4]]

Table A.2: Some outputs from evaluator Qwen-2.5-72B-Inst. Rating[1-3]: Correct and Rating[4]: Incorrect.



Figure 12: Error rate for each model on In-context Recall task. Clearly, few models such as DeepSeek-7B, Phi-3-4B, etc. performs poorly on this simple task.

														-	1.0
Llama-3-70B -	1.00	0.06	0.13	0.12	0.14	0.00	0.17	0.00	0.18	0.26	0.32	0.01	0.26		
Gemma-2-27B -	1.00	0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.03	0.00	0.00	0.00		
Phi-4-14B -	1.00	0.06	0.01	0.00	0.00	0.00	0.01	0.01	0.01	0.01	0.00	0.00	0.01		0.8
Phi-3-14B -	1.00	0.04	0.00	0.00	0.03	0.00	0.05	0.00	0.06	0.12	0.02	0.01	0.00		
Gemma-2-9B -	1.00	0.02	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
Llama-3-8B -	1.00	0.09	0.31	0.44	0.31	0.00	0.40	0.01	0.44	0.47	0.40	0.02	0.42		0.6
Orca-2-7B -	1.00	0.09	0.07	0.01	0.28	0.04	0.28	0.48	0.39	0.69	0.55	0.09	0.50		
Deepseek-7B -	1.00	0.05	0.00	0.02	0.02	0.01	0.03	0.00	0.02	0.41	0.02	0.01	0.03		
Mistral-7B-v0.2 -	1.00	0.24	0.55	0.60	0.77	0.08	0.74	0.40	0.42	0.88	0.68	0.09	0.65	- 1	0.4
Phi-3.5-4B -	1.00	0.09	0.01	0.00	0.02	0.01	0.02	0.02	0.05	0.06	0.00	0.01	0.00		
Phi-3-4B -	1.00	0.08	0.02	0.01	0.31	0.00	0.00	0.04	0.27	0.96	0.50	0.00	0.04		
Llama-3.2-3B -	1.00	0.02	0.01	0.00	0.00	0.00	0.01	0.00	0.01	0.00	0.01	0.00	0.00	-	0.2
Gemma-2-2B -	1.00	0.03	0.00	0.00	0.06	0.00	0.00	0.00	0.01	0.02	0.03	0.00	0.00		
Llama-3.2-1B -	1.00	0.04	0.01	0.01	0.02	0.00	0.01	0.00	0.03	0.02	0.00	0.02	0.00		
Erelish French Chinese poarese Hindi Rissian Arabic Greet within Swamin Repail Ukahinan Thai															

English Fall Back Rate - Factual Recall

Figure 13: English Fall Back Rate across models (The English Fall Back Rate measures the frequency with which a model defaults to English in its output).

English -	0.0	0.0	0.0	0.0	0.0	1.6	4.8	1.6	1.6	6.5	1.6	1.6	0.0		
French -	0.0	0.0	0.0	0.0	1.6	1.6	11.3	3.3	3.3	8.1	1.6	1.6	3.2		50
Chinese -	4.8	9.7	1.6	0.0	12.9	6.5	41.9	18.0	29.5	24.2	41.0	4.9	22.6		50
Japanese -	1.6	0.0	1.6	0.0	3.2	1.6	22.6	8.2	16.4	32.3	22.9	4.9	3.2		
Hindi -	3.2	9.7	4.8	1.6	0.0	6.5	22.6	11.5	11.5	21.0	4.9	6.6	6.5		40
Russian -	0.0	3.2	3.2	0.0	1.6	0.0	21.0	3.3	1.6	22.6	13.1	3.3	3.2		
Arabic -	0.0	3.2	14.5	3.2	14.5	3.2	17.7	11.5	21.3	30.6	27.9	4.9	16.1	- :	30
Greek -	1.6	3.2	0.0	1.6	1.6	1.6	19.4	1.6	6.6	19.4	9.8	1.6	4.8		
Turkish -	4.8	6.5	3.2	3.2	6.5	4.8	25.8	14.8	0.0	16.1	13.1	1.6	6.5	- :	20
Swahili -	14.5	4.8	4.8	3.2	1.6	19.4	14.5	32.8	21.3	9.7	4.9	3.3	3.2		
Nepali -	3.2	16.1	8.1	6.5	1.6	11.3	25.8	21.3	21.3	33.9	0.0	16.4	22.6	- 1	10
Ukrainian -	1.6	1.6	3.2	0.0	3.2	0.0	21.0	6.6	3.3	24.2	18.0	0.0	1.6		
Thai -	3.2	9.7	11.3	3.2	12.9	4.8	24.2	13.1	59.0	40.3	31.1	18.0	0.0		~
United states france china poar mala pussia prava creece when here here the trained															

Figure 14: Country-Specific Factual Error Rates in each language for Llama-3-70B



Figure 15: Country-Specific Factual Error Rates in each language for Llama-3.2-1B



Model-wise Comparison of FRS, KTS, and X-FAKT Scores

Figure 16: Comparision of models (in the increasing order of size with respect to the parameters) using Factual Recall Score, Knowledge Transferability Score, and Cross-Lingual Factual Knowledge Transferability Score.