# **CUTE: Measuring LLMs' Understanding of Their Tokens**

#### **Anonymous ACL submission**

#### Abstract

Large Language Models (LLMs) show remarkable performance on a wide variety of tasks. Most LLMs split text into multi-character tokens and process them as atomic units without direct access to individual characters. This raises the question: To what extent can LLMs learn orthographic information? To answer this, we propose a new benchmark, CUTE, which features a collection of tasks designed to test the orthographic knowledge of LLMs. We evaluate popular LLMs on CUTE, finding that most of them seem to know the spelling of their tokens, yet fail to use this information effectively to manipulate text, calling into question how much of this knowledge is generalizable.

#### 1 Introduction

001

005

007

021

028

037

Large Language Models (LLMs) attract a lot of interest due to their strong performance on many NLP tasks. They have demonstrated a level of fluency rivaling humans. However, it is often overlooked that LLMs lack direct access to the characters composing their tokens. They can only infer knowledge about the characters from the context during pretraining or instruction tuning. While there are models that use characters as input units, none of them have been instruction-tuned (to our knowledge).

Our work examines how well LLMs understand the composition of their tokens. This knowledge enables LLMs to better generalize to new languages and to perform well on tasks involving characterlevel understanding (e.g. code, crosswords, poetry, etc.), or any task requiring arbitrary string operations.

We introduce Character-level Understanding of Tokens Evaluation (CUTE)<sup>1</sup>, a benchmark consisting of several tasks designed to be easy for humans

to complete, given our ability to process characters individually. We evaluate several LLMs ranging from 7B to 132B parameters in size on CUTE to answer the following questions: 038

039

040

041

042

043

044

045

046

047

054

058

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

- 1. Do LLMs know which characters make up their tokens?
- 2. Do LLMs understand the difference between semantic and orthographic similarity?
- 3. Can LLMs manipulate text at the character level?

We address these questions with a set of tasks primarily on the character level, and additionally on the word level to see the difference in performance on the two granularities, thereby separating the understanding of the task from the understanding of what makes up a token.

## 2 Related Work

Itzhak and Levy (2022) mostly analyze encoderonly models and test if they understand how to spell words after fine-tuning on 32k examples. They conclude that models learn to spell their tokens "to some extent." They also experiment with GPT-2 (Radford et al., 2019) which performed similarly to the other models, and also used training examples. Kaushal and Mahowald (2022) probe models with a task asking if a letter is in a word (similar to our character contains task, see §3). Similar to Itzhak and Levy (2022), their probe requires training, as they use models of similar size. By contrast, we examine models with 10 to 200 times as many parameters and apply few-shot prompting without fine-tuning.

Huang et al. (2023) experiment with spelling correction, unscrambling words, and finding words in a string of characters. Most of their experiments concern a training method they propose, but they also include results of GPT-3 (Brown et al., 2020) with few-shot prompting, finding that it performs well on spelling correction, but poorly on other

<sup>&</sup>lt;sup>1</sup>We release our benchmark open-source at: [link to be added in camera-ready].

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

077

tasks compared to their trained character-based models. Most of their tasks require semantic knowledge, which we wish to ablate in our benchmark.

Other benchmarks testing orthographic knowledge often focus on morphology, such as the SIG-MORPHON inflection tasks (e.g. Goldman et al. (2023)). Inflection is fairly regular for most languages, and could be memorized by a language model. Given that many developers of LLMs do not disclose the sources of their pretraining data, it is possible that LLMs memorize these inflection patterns. Therefore, from this we cannot conclude that LLMs can inflect a new word given a new set of rules, and furthermore we cannot conclude that LLMs can apply an arbitrary manipulation to a sequence. Most LLMs are trained with performance on English in mind, and English being less morphological in nature makes inflection a nonideal measure for a model's understanding of the composition of tokens.

There is an extensive body of research on character-level models, where each character forms a token (Lee et al., 2017; Xue et al., 2022; Tay et al., 2022, *inter alia*). Several works compared these models to subword models (Libovický et al., 2022; Edman et al., 2022, 2024, *inter alia*), but they evaluated on tasks requiring additional training. We assume that character-based models would perform well on our benchmark, but we cannot test it since none of these models have been instruction-tuned.

## **3** Benchmark

We split our tasks into 3 categories: understanding composition, understanding orthographic similarity, and ability to manipulate sequences. Figure 1 shows an example for each task.<sup>2</sup> Data gathering and processing details can be found in Appendix B. Our tasks are synthetically generated from existing corpora. There are non-synthetically generated datasets which partially test our research questions, but these datasets have external factors (e.g. domain and/or language in a translation dataset) that would likely obscure our findings. Some of these datasets also might have been leaked into the LLM's pretraining data and been memorized, resulting in an unrealistically good performance.

**Composition** We start with a straightforward benchmark: spelling. Similar to Itzhak and Levy (2022), we include a task where the input is a word

given as a single token<sup>3</sup>, and the output is the same word with spaces in between, so that each character becomes a separate token. This is the most straightforward probe to see whether a model has knowledge of the characters forming the tokens. We also add the inverted task ("inverse spelling") to check if characters can also be mapped to tokens.

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

152

153

154

155

156

157

158

159

162

163

164

165

166

167

169

170

171

Another method for assessing a model's understanding of composition is to ask if a token contains a certain character. If a model managed the previous tasks, we would expect it to succeed here as well. However, a model might not understand the relationship between spelling and membership of characters to a word, so we test this as well. We also test whether the LLMs are able to solve the corresponding word-level task (i.e. is a word in a sentence) to separate the model's general understanding of the task from its ability to solve the task at the character level.

**Similarity** Since the introduction of word2vec (Mikolov et al., 2013), language models have typically been trained to predict words from their context. The resulting token embeddings mainly reflect the semantic and syntactic similarity of tokens. Our next tasks examine whether LLMs also comprehend orthographic similarity. We ask which one of two candidate words is orthographically (or semantically) more similar to a given word. Our candidate words are chosen to be relatively easy to distinguish for a human without explicit knowledge of how to measure orthographic or semantic similarity. For more details, see Appendix B.

**Manipulation** Our previous tasks focus on the understanding of the model. Now, we turn our focus to acting on that understanding. The next tasks involve 4 types of manipulation of the input at the character or word level: Insertion, Deletion, Substitution, and Swapping. We consider these tasks as elementary tasks for modifying a text sequence.

**Insertion** First we test how well the model can insert an element X after every instance of some element Y in the sequence. Similar modifications occur when we replicate letters to emphasize a word (e.g. "Yay!" vs. "Yaaaaay!"), or when we add an adjective next to a noun.

**Deletion** Deletion requires the model to recognize an element and remove all instances of it. This

<sup>&</sup>lt;sup>2</sup>The prompts shown here are not the full prompts. See Appendix A for more details.

<sup>&</sup>lt;sup>3</sup>This is in the ideal case. We cannot guarantee every LLM tested uses only a single token for each input, but we minimize the chance of splitting by using frequent words in our task.

		Task	Input	Output
		Spelling	Spell out the word: there	there
	Composition	Inverse Spelling	Write the word that is spelled out (no spaces): t h e r e	there
		Contains Character	Is there a 'c' in 'there'?	No
		Contains Word	Is there a 'the' in 'the sky is blue'?	Yes
	Similarity	Orthographic Similarity	Which is closer in Levenshtein distance to 'happy'? glad or apply	apply
		Semantic Similarity	Which is more semantically related to 'happy'? glad or apply	glad
		Character Insertion	Add 'b' after every 'e' in 'there'	thebreb
	Manipulation	Word Insertion	Add 'is' after every 'the' in 'the sky is blue'	the is sky is blue
		Character Deletion	Delete every 'e' in 'there'	thr
		Word Deletion	Delete every 'the' in 'the sky is blue'	sky is blue
		Character Substitution	Replace every 'e' with 'a' in 'there	thara
		Word Substitution	Replace every 'the' with 'is' in 'the sky is blue'	the is sky is blue
		Character Swapping	Swap 't' and 'r' in 'there'	rhete
		Word Swapping	Swap 'the' and 'is' in 'the sky is blue'	is sky the blue

Figure 1: All of the tasks in CUTE.

can occur in natural language at the character level
with inflection in languages (e.g. turning an English plural noun into singular), or removing adjectives from a sentence.

Substitution Substitution replaces all instances
of an element in a sequence with another element.
This can occur with spelling or vocabulary variations across dialects or related languages (e.g. "defense" vs. "defence", or "elevator" vs. "lift").

**Swapping** Swapping is a simplified case of reordering acting on two elements.<sup>4</sup> Though reordering is not very common in English, it features heavily in languages with free word order, such as Greek, where stressed words can be moved to the front of the sentence.

	Llama	Mistral	Gemma	Cmd-R(+)	DBRX
Params (B)	7,13,70	7,47	7	35,104	132
Tokens (k)	32	32	256	256	100
Language	EN	EN	Hybrid	Multil.	EN

Table 1: Models evaluated on our benchmark. We consider Gemma as "hybrid" since it has a multilingual tokenization, but English-centric training.

#### 4 Experimental Setup

**Models** We use the models shown in Table 1 (Touvron et al., 2023; Jiang et al., 2023, 2024; Team et al., 2024).<sup>567</sup> We choose freely-available<sup>8</sup>

c4ai-command-r-plus

LLMs largely based on their popularity, as we are unaware of LLMs that specifically address the problems raised. While there are many differences between the models, we highlight 3 that could possibly affect performance: parameter count, vocabulary size, and multilingual versus English-centric or English-only training.

191

192

193

194

195

196

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

**Prompts** We use a template inspired by Bsharat et al. (2023)'s few-shot template to prompt our models with 4 examples. Further details and a sample prompt can be found in Appendix A.

#### **5** Results

Table 2 shows the results for each model. The random baseline is 50% for the contains and similarity tasks, and 0% for all other tasks.

**Composition** The models perform very well on the tasks spelling and inverse spelling, though inverse spelling appears slightly more difficult. Although we are not aware of any spelling tasks in instruction tuning or pretraining datasets, we suspect that the similarity of this task with data seen during training allows the models to perform very well, in contrast to the following tasks.

On the contains tasks, the performance at the word level is quite good, showing that the models understand the task, but the performance breaks down at the character level. This indicates the models do not fully understand the relationship between spelling and membership of a character to a word.

**Orthography and Semantic Similarity** In the semantic similarity task, the models correctly choose the more semantically related word 76-93% of the time, and the performance generally

182

184

185

186

- 188 189
- 190

<sup>&</sup>lt;sup>4</sup>Due to the poor performance on swapping, we leave out more complex forms of reordering, but these could be easily added in the future.

<sup>&</sup>lt;sup>5</sup>https://hf.co/CohereForAI/c4ai-command-r-v01 <sup>6</sup>https://hf.co/CohereForAI/

<sup>&</sup>lt;sup>7</sup>https://hf.co/databricks/dbrx-instruct

<sup>&</sup>lt;sup>8</sup>We define "freely available" as those not requiring payment for use, and with available information on the training process.

	Llama		Mistral		Gemma	Command-R(+)		DBRX	
Task	7B	13B	70B	7B	47B	7B	35B	104B	132B
Spelling	93.4	99.6	99.9	87.0	98.4	82.2	98.9	100	99.9
Inverse Spelling	84.5	92.9	92.7	91.1	95.7	73.0	97.8	100	99.9
Contains Character	67.0	69.8	70.1	61.5	68.1	64.0	68.9	79.7	47.8
Contains Word	78.0	84.8	94.6	94.2	90.8	84.8	87.4	99.7	97.2
Orthographic Similarity	32.1	46.3	42.6	55.3	50.8	54.5	43.4	86.5	47.6
Semantic Similarity	76.8	76.3	82.6	84.7	81.8	82.0	81.7	92.6	90.3
Insert Character	8.1	11.5	7.3	4.8	15.9	10.7	7.3	8.9	<b>9.9</b>
Delete Character	13.0	30.5	26.4	34.2	47.8	30.4	37.7	72.0	64.0
Substitute Character	11.1	14.2	19.6	21.0	34.6	15.1	28.7	55.5	47.7
Swap Character	3.5	2.5	5.6	2.8	6.4	2.5	5.3	10.0	7.5
Insert Word	20.1	36.7	31.6	11.3	40.8	36.5	42.2	<b>81.7</b>	53.2
Delete Word	34.8	56.7	46.8	53.8	56.5	56.7	55.8	65.7	<b>86.3</b>
Substitute Word	72.9	70.7	86.6	70.0	90.7	70.7	75.8	<b>95.1</b>	93.6
Swap Word	10.9	17.2	36.9	15.9	37.8	17.2	28.3	<b>81.4</b>	60.3

Table 2: Accuracy (%) on all of the tasks for each model. Best in bold.

increases with model size. For orthographic
similarity, the performance is below or near random for all models except Mistral-7B, Gemma, and
Command-R+. It is not yet clear why CommandR+ performs so well on this task, but apparently it
is not solely a scaling effect since DBRX fails to
perform above random chance. It may be due to the
combination of scaling and multilingual pretraining, encouraging orthographically similar words to
have similar embeddings.

235

240

241

242 243

244

245

246

247

248

249

250

Manipulation In our manipulation tasks, the models struggle more at the character level than at the word level. The difference is quite profound, with performance gaps of up to 72.8% on Command-R+ for insertion. Larger models such as Command-R+ perform well on deleting characters, with 72% accuracy, though it should be noted that we are only testing on the 1000 most frequent words, making the evaluation fairly generous.

Like the contains task and orthographic similarity task, the manipulation tasks show that LLMs lack a complete understanding of their tokens, although they can literally spell them out. The higher word-level performance indicates that it is not due to a lack of understanding of the task itself.

Vocabulary Size and Multilinguality To see the
effects of vocabulary size and multilingual versus
English tokenization, we compare Llama and Mistral with Gemma, as they all have 7B variants. The
results are mixed: Mistral and Gemma perform
similarly, with Mistral performing slightly better
on most tasks, whereas Gemma performed remarkably well on the insertion tasks. We see no obvious connection between the insertion tasks and

multilinguality or vocabulary size. Neither vocabulary size nor multilinguality seem to play a huge role. However, we expect a larger effect of the vocabulary size when it approaches the number of characters. We leave this for future research. 260

261

262

263

264

265

266

267

268

270

271

272

274

275

276

277

278

279

280

281

282

284

285

287

288

291

292

293

**Model Size** Larger models clearly tend to perform better. This aligns well with the myriad of works showing the benefits of scaling up models and raises the question: Is scaling all we need for good performance on character-level tasks? Looking at the manipulation tasks, it seems that deletion and substitution could become manageable in the near future, but for insertion and swapping, the performance gap between word and character level tasks is large. Many real-world text manipulation tasks are a combination of the tested tasks, so we will likely need more than just scaling.

#### 6 Conclusion

While current LLMs with BPE vocabularies lack direct access to a token's characters, they perform well on some tasks requiring this information, but perform poorly on others. The models seem to understand the composition of their tokens in direct probing, but mostly fail to understand the concept of orthographic similarity. Their performance on text manipulation tasks at the character level lags far behind their performance at the word level. LLM developers currently apply no methods which specifically address these issues (to our knowledge), and so we recommend more research to better master orthography. Character-level models are a promising direction. With instruction tuning, they might provide a solution to many of the shortcomings exposed by our CUTE benchmark.

# 294 295

296

297

306

312

313

314

315

316

318

319

321

323

324

325

329

331

334

335 336

337

341

342

343

# 7 Limitations

We prompt instruction-tuned LLMs without any fine-tuning on benchmark data. While this can be seen as a limitation, we note that it is not feasible to add more training data whenever we discover a new issue with LLMs. We expect that the performance of all models would increase after fine-tuning.

We do not evaluate any character-level models since there are no instruction-tuned versions (to our knowledge).

Our benchmark does not control whether LLM tokenizers split words into multiple tokens. We minimize that chance by choosing frequent words, which are likely to be tokenized as single tokens, but there is no guarantee that splits do not occur.

We only test on English, mainly due to the abundance of English-centric LLMs. We expect the results to vary widely with language. Given a vocabulary of frequent words, we could easily create spelling, inverse spelling, character contains, and character-level manipulation tasks for other languages. For the similarity tasks, we would need a word embedding model, and for the word contains and word-level manipulation tasks, we would need a small corpus of simple sentences containing mostly frequent words, or we could generate such sentences with an LLM for that language, following TinyStories (Eldan and Li, 2023).

Lastly, we do not control for generations that do not match the pattern of the examples given in the prompt. Therefore, we cannot guarantee that all generations considered correct by humans are evaluated as such. Hence the performance of some models may be lower than expected. We will provide the outputs of all models in our repository for further analysis.

## References

- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam Mc-Candlish, Alec Radford, Ilya Sutskever, and Dario

Amodei. 2020. Language models are few-shot learners. *Preprint*, arXiv:2005.14165.

345

346

347

348

351

352

353

354

356

357

358

360

361

362

363

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

384

385

386

388

389

390

391

392

393

394

395

396

397

398

399

- Sondos Mahmoud Bsharat, Aidar Myrzakhan, and Zhiqiang Shen. 2023. Principled instructions are all you need for questioning llama-1/2, gpt-3.5/4. *arXiv preprint arXiv:2312.16171*.
- Lukas Edman, Gabriele Sarti, Antonio Toral, Gertjan van Noord, and Arianna Bisazza. 2024. Are character-level translations worth the wait? comparing byt5 and mt5 for machine translation. *Preprint*, arXiv:2302.14220.
- Lukas Edman, Antonio Toral, and Gertjan van Noord. 2022. Subword-delimited downsampling for better character-level translation. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 981–992, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Ronen Eldan and Yuanzhi Li. 2023. Tinystories: How small can language models be and still speak coherent english? *Preprint*, arXiv:2305.07759.
- Omer Goldman, Khuyagbaatar Batsuren, Salam Khalifa, Aryaman Arora, Garrett Nicolai, Reut Tsarfaty, and Ekaterina Vylomova. 2023. SIGMORPHON– UniMorph 2023 shared task 0: Typologically diverse morphological inflection. In Proceedings of the 20th SIGMORPHON workshop on Computational Research in Phonetics, Phonology, and Morphology, pages 117–125, Toronto, Canada. Association for Computational Linguistics.
- Jing Huang, Zhengxuan Wu, Kyle Mahowald, and Christopher Potts. 2023. Inducing character-level structure in subword-based language models with type-level interchange intervention training. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12163–12180, Toronto, Canada. Association for Computational Linguistics.
- Itay Itzhak and Omer Levy. 2022. Models in a spelling bee: Language models implicitly learn the character composition of tokens. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5061–5068, Seattle, United States. Association for Computational Linguistics.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las

Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. Mixtral of experts. *Preprint*, arXiv:2401.04088.

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416 417

418

419

420

421

422

423

424

425

426

427 428

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

- Ayush Kaushal and Kyle Mahowald. 2022. What do tokens know about their characters and how do they know it? In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 2487–2507, Seattle, United States. Association for Computational Linguistics.
- Jason Lee, Kyunghyun Cho, and Thomas Hofmann. 2017. Fully character-level neural machine translation without explicit segmentation. *Transactions of the Association for Computational Linguistics*, 5:365– 378.
  - Jindřich Libovický, Helmut Schmid, and Alexander Fraser. 2022. Why don't people use character-level machine translation? In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2470–2485, Dublin, Ireland. Association for Computational Linguistics.
  - Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
  - Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
  - Yi Tay, Vinh Q. Tran, Sebastian Ruder, Jai Gupta, Hyung Won Chung, Dara Bahri, Zhen Qin, Simon Baumgartner, Cong Yu, and Donald Metzler. 2022. Charformer: Fast character transformers via gradient-based subword tokenization. *Preprint*, arXiv:2106.12672.
- 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. 2024. Gemma: Open models based on gemini research and technology. Preprint, arXiv:2403.08295.

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

502

503

504

506

508

509

510

511

512

513

514

515

516

- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models. Preprint, arXiv:2307.09288.
- Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2022. ByT5: Towards a token-free future with pre-trained byte-to-byte models. *Transactions of the Association for Computational Linguistics*, 10:291–306.

## **A** Prompting Details

We show an example of a full prompt in Figure 2. All of our prompts are available with the release of our benchmark. For generation, we use greedy search.

For evaluation of the generation, we rely on the given start-quote to denote the start of the answer, and we filter out anything after the end quote (e.g. "I hope this answer helped!"). Some models also were prone to starting generation with a generic response such as "Sure I can do that for you.". For

543

544

545

```
[INST] Spell out the word, putting
spaces between each letter, based
on the following examples:
1. Spell out the word "alphabet".
Answer: "a l p h a b e t"
2. Spell out the word "hello".
Answer: "h e l l o"
3. Spell out the word "zebra".
Answer: "z e b r a"
4. Spell out the word "tongue".
Answer: "t o n g u e"
Question: Spell out the word "cow".
[/INST]
Answer: "
```

Figure 2: An example of a full prompt to spell the word "cow", with examples, for the task spelling.

these generations, we observe that it would repeat "Answer: ", so we filter out all generations before this point. Of the remaining generations, some could be considered correct though did not match the desired pattern (e.g. "H-E-L-L-O" rather than "h e 11 o" for the spelling task). These we ultimately consider incorrect so as not to unfairly elevate any model's performance.

517

518

519

520

522

523

524

525

526

530

532

533

534

535

538

Concerning the wording of the orthographic similarity task, it could be argued that models do not understand the concept of Levenshtein distance, and thus it is not comparable to the semantic similarity task. We also tested using "closer in edit distance" and "closer in spelling" in the prompt, and the results were very similar, so we opted for Levenshtein distance as it is more well-defined. Similarly, we used "closer in meaning" rather than "more semantically related" and achieved similar results, though it is debatable whether an antonym should be considered close in meaning, so we opted for the latter.

#### **B** Data Processing

Here we detail our exact method for gathering and
processing the data into our tasks. The scripts for
processing the data and the resulting data can be
found at our Github repository.<sup>9</sup>

For the word-based tasks, we use the TinyStories (Eldan and Li, 2023) dataset, which consists of stories written by an LLM in a style appropriate for a 3-4 year old reader. This has the benefit of using simple sentences with a limited vocabulary, maximizing the chances of words being tokenized into a single token in the models we test, while also ensuring that the complexity of the sentence is not a confounding factor in a model's performance on the tasks.

**Filtering** For character-based tasks, we select the 1000 most frequent words that are at least 3 characters long. For word-based tasks, we similarly filter for 1000 sentences of length 3-10 words, in order to make the length similar to the number of characters in a word seen in the character-level tasks.

For addition, deletion, and substitution, we apply the modification to those 1000 words, resulting in our dataset.

For swapping, we need to sure that the word or sentence has 2 items that are unique, so as to avoid an ambiguous prompt (e.g. swap the 'e' and 'g' in 'engineering'). As such, we select the 1000 most frequent words or the first 1000 sentences that satisfy this criteria, as well as satisfying our length constraints.

**Similarity Data** For our similarity data, we require our candidate pairs to be sufficiently easy for a human to distinguish which is closer orthographically and which is closer semantically.

To accomplish this, our candidate words must satisfy two thresholds, one based on normalized Levenshtein distance (for othographic similarity), and one based on cosine similarity to other fastText (Bojanowski et al., 2017) embeddings (for semantic similarity). That is to say, the word must be sufficiently similar in one metric (0.7+ and 0.5+ for Levenshtein and cosine, respectively) and sufficiently dissimilar in the other (0.3- and 0.2-, respectively). These thresholds are decided empirically. We note that this process occasionally ends up with semantic pairs that are antonyms (e.g. "good" and

**Data Sources** For almost all tasks, we require a set of frequent English words that were most likely to be tokenized into a single token. For this, we use a dataset derived from the Google Web Trillion Word Corpus.<sup>10</sup>

<sup>&</sup>lt;sup>9</sup>Will add link here.

<sup>&</sup>lt;sup>10</sup>https://www.kaggle.com/datasets/rtatman/ english-word-frequency/data

591 "bad"), and thus we refrain from stating that the592 pairs are similar in meaning.