

When English Shapes Alterations: Suffix Bias in the Italian Fable Genre

Anonymous ACL submission

Abstract

Despite recent advances in multilingual Large Language Models (LLMs), prior work has identified systematic differences between human-authored and LLM-generated Italian texts, particularly in language-specific morphology. We investigate the “naturalness” of LLM outputs in Italian through the use of alterative suffixes, a productive feature of Italian morphology, analyzing both human-written and LLM-generated fables. Our results show that English-centered models produce fewer suffixes and lower lexical variability than human authors. We further study suffix generation in Machine Translation (MT) using a parallel corpus of English-Italian fables, finding that LLMs rarely reproduce the suffixal strategies used in the Italian translation. We also introduce SUFFIXBIASMT, a contrastive benchmark for evaluating biases against Italian alterative morphology in MT metrics. Experiments reveal that several widely used metrics systematically prefer translations without suffixation, highlighting broader evaluation biases affecting morphologically rich languages like Italian.

1 Introduction

Recent years have seen major efforts to develop multilingual Large Language Models (LLMs) for non-English languages (Orlando et al., 2024; Martins et al., 2024; Ali et al., 2025). However, linguistically oriented studies have identified systematic differences between human-authored and LLM-generated texts, especially in Italian (Cicero, 2023, 2025; Raus, 2025; Antonelli, 2025; Fiorentino and Tivosanis, 2024; Tivosanis, 2024; De Cesare, 2023, 2025, 2026), particularly in the reproduction of language-specific morphological features.

This gap may stem from both data and architectural biases. Multilingual LLMs are typically trained on corpora heavily dominated by English, with other languages are underrepresented (Blevins and Zettlemoyer, 2022; Grattafiori et al., 2024),

leading to English-centric behaviors in non-English generation (Guo et al., 2025). Moreover, core architectural components such as tokenization and vocabulary design are often optimized for English, limiting their ability to capture the morphological richness of languages like Italian (Moroni et al., 2025b).

In this work, we focus on the use of Italian alterative suffixes, a productive morphological feature, in Italian fables. Fables, as a genre closely associated with child-centered communication (Dressler and Barbaresi, 1994), provide a particular fertile ground for the use of alterative morphology in both Italian and English. However, the two languages differ substantially in this respect: English possesses a much more limited and less productive system of alterative suffixation than Italian (Dressler and Barbaresi, 1994).

We investigate LLMs’ English-centric biases in suffixed word usage in Italian fables through three research questions: i) *How do LLMs handle alterative suffixes in fable generation compared to humans?*, ii) *How does suffix generation emerge when translating fables from English to Italian?*, and iii) *Are MT evaluation metrics biased against suffixed forms in Italian?*

To answer the first question, we compare human-written fables with texts generated by seven LLMs. We find that models pretrained predominantly on English produce significantly fewer alterative suffixes than models more oriented toward Italian. Across models, suffix usage is also concentrated on a narrow set of lexical bases, indicating lower variability and creativity than in human writing.

We then investigate English-Italian literary translation using a parallel corpus of fables translated by a professional translator. Results show that LLMs rarely reproduce the alterative morphology found in human translations, revealing persistent limitations in modeling morphologically expressive forms under translation constraints.

084	Finally, inspired by recent work on evaluation	bias in MT, while Xu et al. (2026) show that cul-	130
085	biases in MT (Perrella et al., 2024; Ahmadi et al.,	tural familiarity of named entities affects transla-	131
086	2025), we introduce SUFFIXBIASMT, a contrastive	tion quality. Recent studies also reveal biases in	132
087	benchmark designed to assess biases against Italian	MT evaluation metrics: Perrella et al. (2024) find	133
088	alterative morphology. We show that several state-	that neural metrics like COMET (Rei et al., 2020)	134
089	of-the-art MT evaluation metrics systematically	are sensitive to surface properties such as sentence	135
090	penalize suffixed forms in favor of their adjective-	complexity, and Ahmadi et al. (2025) show a pref-	136
091	based counterparts, revealing a source of evaluation	erence for loanwords over native forms. However,	137
092	bias in Italian and, more broadly, in morphologi-	no prior work examines morphological biases in	138
093	cally rich languages. ¹	child-centered settings, particularly for Italian suf-	139
094	2 Related Work	fixation.	140
095	2.1 Manual Evaluation of Italian	3 RQ1: “How do models handle suffix	141
096	Synthetically Generated Texts	words compared to humans in fable	142
097	Linguistic analyses of LLM-generated text in free-	generation?”	143
098	form settings have provided useful insights in the	3.1 Methodology	144
099	Italian research landscape (De Cesare, 2023; An-	To investigate how LLMs produce suffix-based	145
100	tonelli, 2025). De Cesare (2026) studies gender	words compared to human-written texts, we first	146
101	bias in short biography generation for Italian and	collected 200 Italian fables from human writers	147
102	French anthroponyms, showing that men are more	(more details in Appendix A). From this human-	148
103	often referred to by surname, while women are	written fable collection we manually tag all the	149
104	more frequently mentioned by full name or first	words that use any of the following eight suffixes:	150
105	name alone. Tavosanis (2024) highlights the impor-	three diminutives, <i>-ino</i> , <i>-etto</i> , <i>-ello</i> ; two augmen-	151
106	tance of human evaluation for identifying subtle is-	tatives, <i>-one</i> , <i>-asso</i> ; the elative suffix <i>-issimo</i> ; and	152
107	ssues in synthetic text and report actionable find-	two polyfunctional suffixes, <i>-otto</i> , <i>-ozzo</i> , which can	153
108	ings in Italian-generated corpora. However, little	be used as diminutives or augmentatives. We then	154
109	work has properly examined morphological and	prompt LLMs to generate Italian fables, and tag	155
110	suffixal phenomena in Italian synthetic child-	them as above.	156
111	centered storytelling.	3.2 Experimental Setup	157
112	2.2 Baby speech in NLP	Models: Following Moroni et al. (2025a), we	158
113	Recent work has explored language model-	group instruction-tuned models into three cate-	159
114	ing under child language acquisition con-	gories: (i) closed models: GPT-3.5-Turbo (Brown	160
115	straints (Warstadt et al., 2023; Dhole, 2026). The	et al., 2020) and GPT-4o-mini (OpenAI et al.,	161
116	BabyLM challenge (Warstadt et al., 2023) pro-	2024); (ii) open-data models: Minerva-7B (Or-	162
117	promotes training on limited, developmentally	lando et al., 2024), OLMo-2-7B (OLMo et al.,	163
118	plausible corpora (10M–100M words), approxi-	2025), and Teuken-7B (Ali et al., 2025);	164
119	mating early linguistic input. Follow-up work	and (iii) Italian-adapted models: Llama-3.1-	165
120	extends this to multimodal child-centered set-	8B (Grattafiori et al., 2024), further adapted to	166
121	tings (Hu et al., 2024) and child-directed rea-	Italian through continued pre-training (Llama-	167
122	soning benchmarks (Dhole, 2026). Despite this	LAPT), and its vocabulary-adapted variant with	168
123	progress, little work analyzes morphological	an Italian-optimized tokenizer (Llama-Adapt)	169
124	realization in child-directed or storytelling	(Moroni et al., 2025b).	170
125	settings for morphologically rich languages	Generation strategy: We used the vLLM li-	171
126	such as Italian, where suffixation plays a	brary (Kwon et al., 2023) with temperature =	172
127	central role.	0.7 for all generations. Synthetic fables were	173
128	2.3 Biases in Machine Translation	produced by repeating the prompting process	174
129	Prior work identifies several biases in	25 times. Prompts used in fable generation	175
	machine translation. Savoldi et al. (2021)	are reported in Appendix B.	176
	study gender		

¹Code and data are available at [omitted.link](#).

Model	Suffixed words	Entry words
Human	93.6	27.7
GPT-4o-mini	39.5	8.2
GPT-3.5-turbo	157.1	14.1
OLMO-2-7B	16.8	8.9
Minerva-7B	67.0	15.8
Teuken-7B	104.5	12.8
Llama-LAPT	106.1	9.7
Llama-Adapt	107.2	12.4

Table 1: Number of suffix words and entry words in human and synthetic fables, normalized over 10,000 words. Entry words indicate distinct lexical bases appearing with at least one suffix form.

Model	-asso	-ello	-etto	-ino	-issimo	-one	-otto	-ozzo	# Suff.
Human	0.0	4.3	19.9	45.7	26.5	2.8	0.8	0.0	6
GPT-4o-mini	0.0	1.5	10.2	57.7	28.5	0.0	2.2	0.0	5
GPT-3.5-turbo	0.0	1.2	44.7	47.1	6.6	0.3	0.0	0.0	5
OLMo-2-7B	0.0	6.7	13.3	46.7	26.7	6.7	0.0	0.0	5
Minerva-7B	0.0	2.8	37.5	26.4	30.6	0.0	2.8	0.0	5
Teuken-7B	0.0	2.5	11.5	63.1	22.5	0.4	0.0	0.0	5
Llama-LAPT	0.0	0.6	42.9	46.6	9.8	0.0	0.0	0.0	4
Llama-Adapt	0.0	4.4	53.0	31.5	9.9	1.1	0.0	0.0	5

Table 2: Row-normalized distribution of Italian suffix types expressed as percentages. **Bold** values indicate row-wise most frequent suffixes.

3.3 Results

Table 1 shows the number of suffix words and entry words (i.e., base forms occurring with at least one suffix) in human and synthetic fables, normalized over 10,000 words. Results show substantial variation across models. Some LLMs generate more suffixed words than humans, while others produce considerably fewer. For example, GPT-3.5-turbo produces the highest number of suffixed words (157.1), whereas GPT-4o generates only 39.5, well below the human average (93.6). Despite this variability, all models exhibit lower lexical diversity than humans. Human texts present 27.7 entry words, while the best-performing model, Minerva-7B, reaches only 15.8, suggesting that models reuse a limited set of lexical bases for suffix formation. Models with stronger Italian exposure during pretraining generally produce more suffixes and greater lexical variety. Teuken-7B and Llama-based models generate suffixed word frequencies close to human texts, whereas the English-centric OLMo-2-7B produces the fewest suffixed and entry words, highlighting the impact of pretraining data composition on Italian suffixation. Table 2 shows the distribution of suffix types. Human fables are dominated by *-etto*,

-ino, and *-issimo*, a pattern partially reproduced by most models, indicating that LLMs capture some human-like preferences in suffix selection.

4 RQ2: “How does suffix generation emerge when translating fables from English to Italian?”

4.1 Methodology and Setup

To further analyze how suffixed words are handled in generative contexts we investigate how LLMs produce them when translating English fables into Italian. We considered 68 fables from “*Favole al Telefono*” by Gianni Rodari, which have been translated into English by a professional translator. We then prompted models, first, with an Italian system prompt and, then, with an English system prompt to determine the impact of crosslingual prompting in an MT setting. We followed the experimental setup defined in Section 3, further details about MT prompt strategies are reported in Appendix C.

4.2 Results

Table 3 reports English-to-Italian machine translation performance using both Italian and English prompts. Translation quality is measured with COMET, while suffixed accuracy (*Suf. Acc.* in the Table) denotes the percentage of suffixed words from the gold Italian translation correctly reproduced by the model. Prompting language affects performance, although no consistent trend emerges across models. GPT-4o-mini and GPT-3.5-turbo achieve the highest COMET scores (around 85), whereas OLMo-2-7B and Minerva-7B perform substantially worse. These latter models also obtain the lowest suffix accuracy, suggesting a relation between translation quality and the ability to reproduce Italian suffixation. Overall, all models struggle to generate gold suffixed forms. The best suffix accuracy is achieved by GPT-3.5-turbo with English prompting (18.5) and GPT-4o-mini with Italian prompting (16.9), indicating that suffix reproduction remains difficult even for strong translation models.

5 RQ3: “Are MT evaluation metrics biased against suffixed forms in Italian?”

5.1 Methodology

MT evaluation metrics are known to exhibit systematic biases toward lexical and stylistic variation

Model	Italian		English	
	Comet	Suf. Acc.	Comet	Suf. Acc.
gpt-4o-mini	85.1	16.9	85.0	13.2
gpt3.5-turbo	82.4	11.3	84.6	18.5
OLMo-2-7B	66.9	7.0	68.0	5.4
Minerva-7B	76.8	5.8	75.5	6.6
Teuken-7B	82.7	12.5	81.2	10.1
Llama-3.1-8B-LAPT	84.1	9.8	84.0	13.2
Llama-3.1-8B-SAVA	82.9	10.1	83.4	11.1

Table 3: Machine translation results from English to Italian using Italian and English system prompts. COMET and *Suf. Acc* are reported, where the latter reports the percentage of gold Italian suffixed words correctly reproduced in the generated translations.

Original Italian	Rewritten Italian	English Translation
La <i>nuvoletta</i> dormiva vicino al camino.	La <i>piccola nuvola</i> dormiva vicino al camino.	The <i>little cloud</i> was sleeping near the fireplace.

Table 4: Example of the controlled rewriting process used to analyze MT metric biases toward Italian suffixed forms.

(Perrella et al., 2024; Ahmadi et al., 2025). To assess their behavior on Italian alterative morphology, we introduce SUFFIXBIASMT, a contrastive benchmark of translation triples.

To construct the benchmark, we extracted sentences from Italian fables written by humans that contain suffixed words (e.g., *nuvoletta*). We then replaced each suffixed form with a semantically equivalent adjective-based paraphrase (e.g., *piccola nuvola*), thereby creating an alternative Italian version of the sentence. Finally, we translated the adjective-based variants into English. Table 4 illustrates the transformation process.

SUFFIXBIASMT contains 453 items. Thus, each instance in SUFFIXBIASMT consists of: (i) an English sentence s_{eng} , (ii) an Italian translation using alterative suffixation t_{suff} , and (iii) a semantically equivalent translation without suffixation $t_{\text{no-suff}}$.

5.2 Experimental Setup

We evaluate reference-free metrics, including COMET (Rei et al., 2020), xCOMET (Guerreiro et al., 2023), and COMETKiwi (Rei et al., 2022). Scores are computed using the English source sentence (s_{eng}) paired with either a suffix-preserving Italian translation (t_{suff}) or its non-suffixed variant ($t_{\text{no-suff}}$). We also evaluate sentinel-src (Perrella et al., 2024), a source-only metric predicting translation difficulty, where higher scores indicate easier

Model	No Suff.	Suff.	Δ
XCOMET-XL	0.793	0.781	0.012
wmt23-cometkiwi-da-xl	0.716	0.698	0.018
wmt22-cometkiwi-da	0.832	0.819	0.013
wmt20-comet-qe-da	0.171	0.146	0.025
sentinel-src	0.102	0.068	0.034

Table 5: Reference-free MT evaluation scores for Italian translations with and without suffixed forms. Δ represents the score difference between non-suffixed and suffixed variants, where positive values indicate a preference for non-suffix constructions.

sentences. Here, Italian is the source language, and we compare t_{suff} and $t_{\text{no-suff}}$ as alternative inputs.

5.3 Results

Table 5 reports the scores and the performance gap (Δ) between suffixed and non-suffixed variants. All metrics consistently assign higher scores to translations without suffixes, favoring syntactic paraphrases with adjectives over morphologically altered forms. These results suggest that current MT evaluation metrics exhibit biases against Italian suffixation, reflecting broader English-centric preferences in neural evaluation systems.

6 Conclusions

We presented the first study on English-centric biases in the generation and translation of Italian alterative morphology in fables. Across both free-form generation and English-Italian translation, LLMs consistently underproduce alterative suffixes compared to human authors, favoring a narrower and less varied lexical repertoire. These findings suggest that current multilingual LLMs still struggle to model morphologically expressive phenomena that are highly productive in Italian but marginal in English. To further investigate this issue, we introduced SUFFIXBIASMT, a contrastive benchmark targeting evaluation biases toward Italian alterative morphology. Our results show that several widely used MT evaluation metrics systematically penalize suffixed forms in favor of adjective-based alternatives, highlighting an overlooked source of bias in automatic evaluation for morphologically rich languages. Overall, our work emphasizes the importance of moving beyond surface fluency when evaluating multilingual LLMs and MT systems, and calls for more linguistically informed modeling and evaluation practices for non-English languages.

313 Limitations

314 Despite the comprehensiveness of the proposed
315 work, several limitations should be acknowledged.

316 First, the analysis is restricted to a specific genre
317 of child-related fables. While this provides a con-
318 trolled setting for investigation, future work should
319 broaden the scope to include additional genres and
320 stylistic registers in Italian literature, in order to
321 better capture potential variation in morphological
322 patterns.

323 Second, for reasons of simplicity and the sub-
324 stantial effort required for annotation, this study fo-
325 cuses exclusively on Italian. This language-specific
326 scope limits the generalizability of the findings.
327 Extending the analysis to other languages would
328 be valuable, particularly those in which English-
329 centric morphological biases may emerge in down-
330 stream applications.

331 Finally, due to computational constraints, the ex-
332 periments are limited to open-weight, open-data
333 models with up to 8B parameters. While this en-
334 ables reproducibility and accessibility, it restricts
335 the evaluation of larger-scale models. Future re-
336 search could investigate whether the observed ef-
337 fects persist or change when scaling to larger archi-
338 tectures.

339 References

340 Sina Ahmadi, Micha David Hess, Elena Álvarez-
341 Mellado, Alessia Battisti, Cui Ding, Anne Göhring,
342 Yingqiang Gao, Zifan Jiang, Andrianos Michail,
343 Peshmerge Morad, Joel Niklaus, Maria Christina
344 Panagiotopoulou, Stefano Perrella, Juri Opitz, Anas-
345 tasia Shaitarova, and Rico Sennrich. 2025. *Con-*
346 *Loan: A contrastive multilingual dataset for eval-*
347 *uating loanwords*. In *Proceedings of the 63rd An-*
348 *annual Meeting of the Association for Computational*
349 *Linguistics (Volume 1: Long Papers)*, pages 30070–
350 30090, Vienna, Austria. Association for Computa-
351 tional Linguistics.

352 Mehdi Ali, Michael Fromm, Klaudia Thellmann, Jan
353 Ebert, Alexander Arno Weber, Richard Rutmann,
354 Charvi Jain, Max Lübbering, Daniel Steinigen,
355 Johannes Leveling, Katrin Klug, Jasper Schulze
356 Buschhoff, Lena Jurkschat, Hammam Abdelwahab,
357 Benny Jörg Stein, Karl-Heinz Sylla, Pavel Denisov,
358 Nicolo' Brandizzi, Qasid Saleem, and 22 others.
359 2025. *Teuken-7b-base & teuken-7b-instruct: To-*
360 *wards european llms*. *Preprint*, arXiv:2410.03730.

361 Giuseppe Antonelli. 2025. *Storia brevissima (ma molto*
362 *intensa) dell'ia-taliano*. *Lingue e culture dei media*,
363 9(1):4–50.

Terra Blevins and Luke Zettlemoyer. 2022. *Language*
364 *contamination helps explains the cross-lingual capa-*
365 *bilities of English pretrained models*. In *Proceedings*
366 *of the 2022 Conference on Empirical Methods in Nat-*
367 *ural Language Processing*, pages 3563–3574, Abu
368 Dhabi, United Arab Emirates. Association for Com-
369 putational Linguistics. 370

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
371 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
372 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
373 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
374 Gretchen Krueger, Tom Henighan, Rewon Child,
375 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
376 Clemens Winter, and 12 others. 2020. *Language*
377 *models are few-shot learners*. *Preprint*,
378 arXiv:2005.14165. 379

Francesco Cicero. 2023. *L'italiano delle intelli-*
380 *genze artificiali generative*. *Italiano LinguaDue*,
381 15(2):733–761. 382

Francesco Cicero. 2025. *Esercizi di stile. la narra-*
383 *tiva per l'infanzia delle intelligenze artificiali*. *AI-*
384 *Linguistica. Linguistic Studies on AI-Generated Texts*
385 *and Discourses*, 2(2). 386

Anna-Maria De Cesare. 2023. *Assessing the quality of*
387 *chatgpt's generated output in light of human-written*
388 *texts: A corpus study based on textual parameters*.
389 *CHIMERA: Revista de Corpus de Linguas Romances*
390 *y Estudios Lingüísticos*, 10:179–210. 391

Anna-Maria De Cesare. 2025. *Llm referential chain*
392 *generation.: A qualitative case study based on italian*
393 *biographies produced by gpt-4*. *Linguistik Online*,
394 136(4):25–52. 395

Anna-Maria De Cesare. 2026. *Gender biases in*
396 *gpt-4 short biographies. : A corpus study on ital-*
397 *ian and french anthroponyms*. *Linguistik Online*,
398 144(3):227–245. 399

Kaustubh D. Dhole. 2026. *Babyreasoningbench: Gen-*
400 *erating developmentally-inspired reasoning tasks*
401 *for evaluating baby language models*. *Preprint*,
402 arXiv:2601.18933. 403

Wolfgang U. Dressler and Lavinia M. Barbaresi. 1994.
404 *Morphopragmatics*. De Gruyter Mouton, Berlin,
405 New York. 406

G. Fiorentino and M. Tavosanis. 2024. *Chiaro, sin-*
407 *tetico, e brillante: l'italiano dei testi redatti con l'ia*
408 *funziona? LId'O*, 5:37–65. 409

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
410 Abhinav Pandey, Abhishek Kadian, Ahmad Al-
411 Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-
412 ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh
413 Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-
414 tra, Archie Sravankumar, Artem Korenev, Arthur
415 Hinsvark, and 4 others. 2024. *The llama 3 herd of*
416 *models*. *Preprint*, arXiv:2407.21783. 417

418	Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa	Gu, Shengyi Huang, Matt Jordan, Nathan Lambert,	477
419	Coheur, Pierre Colombo, and André F. T. Martins.	Dustin Schwenk, Oyvind Tafjord, Taira Anderson,	478
420	2023. xcomet: Transparent machine translation eval-	David Atkinson, Faeze Brahman, Christopher Clark,	479
421	uation through fine-grained error detection . <i>Preprint</i> ,	Pradeep Dasigi, Nouha Dziri, and 24 others. 2025. 2	480
422	arXiv:2310.10482 .	olmo 2 furious . <i>Preprint</i> , arXiv:2501.00656 .	481
423	Yanzhu Guo, Simone Conia, Zelin Zhou, Min Li, Saloni	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	482
424	Potdar, and Henry Xiao. 2025. Do large language	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	483
425	models have an English accent? evaluating and im-	man, Diogo Almeida, Janko Altschmidt, Sam Alt-	484
426	proving the naturalness of multilingual LLMs . In	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	485
427	<i>Proceedings of the 63rd Annual Meeting of the As-</i>	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-	486
428	<i>sociation for Computational Linguistics (Volume 1:</i>	ing Bao, Mohammad Bavarian, Jeff Belgium, and	487
429	<i>Long Papers)</i> , pages 3823–3838, Vienna, Austria.	Irwan Bello et al. 2024. Gpt-4 technical report .	488
430	Association for Computational Linguistics.	<i>Preprint</i> , arXiv:2303.08774 .	489
431	Michael Y. Hu, Aaron Mueller, Candace Ross, Ad-	Riccardo Orlando, Luca Moroni, Pere-Lluís	490
432	ina Williams, Tal Linzen, Chengxu Zhuang, Ryan	Huguet Cabot, Simone Conia, Edoardo Barba,	491
433	Cotterell, Leshem Choshen, Alex Warstadt, and	Sergio Orlandini, Giuseppe Fiameni, and Roberto	492
434	Ethan Gotlieb Wilcox. 2024. Findings of the second	Navigli. 2024. Minerva LLMs: The first family of	493
435	BabyLM challenge: Sample-efficient pretraining	large language models trained from scratch on Italian	494
436	on developmentally plausible corpora . In <i>The 2nd</i>	data . In <i>Proceedings of the Tenth Italian Conference</i>	495
437	<i>BabyLM Challenge at the 28th Conference on Com-</i>	<i>on Computational Linguistics (CLiC-it 2024)</i> , pages	496
438	<i>putational Natural Language Learning</i> , pages 1–21,	707–719, Pisa, Italy. CEUR Workshop Proceedings.	497
439	Miami, FL, USA. Association for Computational Lin-		
440	guistics.	Stefano Perrella, Lorenzo Proietti, Alessandro	498
441	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	Scirè, Edoardo Barba, and Roberto Navigli.	499
442	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.	2024. Guardians of the machine translation meta-	500
443	Gonzalez, Hao Zhang, and Ion Stoica. 2023. Effi-	evaluation: Sentinel metrics fall in! In <i>Proceedings</i>	501
444	cient memory management for large language model	of the 62nd Annual Meeting of the Association	502
445	erving with pagedattention. In <i>Proceedings of the</i>	for Computational Linguistics (Volume 1: Long	503
446	<i>ACM SIGOPS 29th Symposium on Operating Systems</i>	Papers) , pages 16216–16244, Bangkok, Thailand.	504
447	<i>Principles</i> .	Association for Computational Linguistics.	505
448	Pedro Henrique Martins, Patrick Fernandes, João Alves,	Rachele Raus. 2025. Inclusione ed elaborazione del lin-	506
449	Nuno M. Guerreiro, Ricardo Rei, Duarte M. Alves,	guaggio naturale nell’era dell’intelligenza artificiale	507
450	José Pombal, Amin Farajian, Manuel Faysse, Ma-	generativa . Zenodo.	508
451	teusz Klimaszewski, Pierre Colombo, Barry Haddow,	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	509
452	José G. C. de Souza, Alexandra Birch, and André	Lavie. 2020. COMET: A neural framework for MT	510
453	F. T. Martins. 2024. Eurollm: Multilingual language	evaluation . In <i>Proceedings of the 2020 Conference</i>	511
454	models for europe . <i>Preprint</i> , arXiv:2409.16235 .	on Empirical Methods in Natural Language Process-	512
455	Luca Moroni, Javier Aula-Blasco, Simone Conia, Irene	ing (EMNLP) , pages 2685–2702, Online. Association	513
456	Baucells, Naiara Perez, Silvia Paniagua Suárez, Anna	for Computational Linguistics.	514
457	Sallés, Malte Ostendorff, Júlia Falcão, Guijin Son,	Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro,	515
458	Aitor Gonzalez-Agirre, Roberto Navigli, and Marta	Chrysoula Zerva, Ana C. Farinha, Christine Maroti,	516
459	Villegas. 2025a. Multi-LMentry: Can multilingual	José G. C. de Souza, Taisiya Glushkova, Duarte M.	517
460	LLMs solve elementary tasks across languages?	Alves, Alon Lavie, Luisa Coheur, and André F. T.	518
461	In <i>Proceedings of the 2025 Conference on Empirical</i>	Martins. 2022. Cometkiwi: Ist-unbabel 2022 submis-	519
462	<i>Methods in Natural Language Processing</i> , pages	sion for the quality estimation shared task . <i>Preprint</i> ,	520
463	34126–34157, Suzhou, China. Association for Com-	arXiv:2209.06243 .	521
464	putational Linguistics.	Beatrice Savoldi, Marco Gaido, Luisa Bentivogli, Mat-	522
465	Luca Moroni, Giovanni Puccetti, Pere-Lluís	teo Negri, and Marco Turchi. 2021. Gender bias in	523
466	Huguet Cabot, Andrei Stefan Bejgu, Alessio	machine translation . <i>Transactions of the Association</i>	524
467	Miaschi, Edoardo Barba, Felice Dell’Orletta, Andrea	for Computational Linguistics , 9:845–874.	525
468	Esuli, and Roberto Navigli. 2025b. Optimizing	Mirko Tamosanis. 2024. Valutare la qualità dei testi	526
469	LLMs for Italian: Reducing token fertility and	generati in lingua italiana . <i>AI-Linguistica. Linguistic</i>	527
470	enhancing efficiency through vocabulary adaptation .	Studies on AI-Generated Texts and Discourses , 1(1).	528
471	In <i>Findings of the Association for Computational</i>	Alex Warstadt, Leshem Choshen, Aaron Mueller, Adina	529
472	<i>Linguistics: NAACL 2025</i> , pages 6661–6675,	Williams, Ethan Wilcox, and Chengxu Zhuang. 2023.	530
473	Albuquerque, New Mexico. Association for	Call for papers – the babylm challenge: Sample-	531
474	Computational Linguistics.	efficient pretraining on a developmentally plausible	532
475	Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groen-	corpus . <i>Preprint</i> , arXiv:2301.11796 .	533
476	evelld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling		

534 Lu Xu, Luca Moroni, and Roberto Navigli. 2026. [Cul-](#)
535 [tural and knowledge biases in llms through the lens](#)
536 [of entity-aware machine translation](#). In *Proceedings*
537 *of the Fifteenth Language Resources and Evaluation*
538 *Conference (LREC 2026)*, pages 8794–8812, Palma,
539 Mallorca, Spain. European Language Resources As-
540 sociation (ELRA).

541 **A Human Fables Collection**

542 Table 6 reports the human-authored fable collec-
543 tions used in our analysis. We organize the corpus
544 into two subsets according to the target reader age:
545 (i) 3–5 years, comprising fables intended to be read
546 aloud by adults but primarily designed for early
547 childhood audiences; and (ii) 6–10 years, com-
548 prising texts intended for independent reading by
549 young readers.

550 The categorization was performed manually by a
551 trained annotator with a background in linguistics,
552 based on discourse, syntactic, and morphological
553 cues.

554 We relied on the copyrighted fable texts solely
555 for release purposes. We do not share complete
556 fables or full paragraphs from them.

557 **B Fables Generation Prompts**

558 In Table 7 we report the prompt used to generate
559 free-form children’s fables. To increase variabil-
560 ity, we condition the prompt on the target reader’s
561 age, distinguishing between (i) 3–5 years and (ii)
562 6–10 years, following the human-written resources
563 described in Appendix A.

564 **C Machine Translation Prompts**

565 In Tables 8 and 9, we report the Italian and English
566 prompts used to evaluate models on fable transla-
567 tion (from English into Italian).

Human written fables	N. of fables	Publication date
3–5 years		
Le Favole di Morfeus	19	2018
Lecture per bambini	18	2016
Storie brevi della buonanotte	37	2015–2022
Favole per bambini	6	2020
15 brevi racconti per bambini sulla natura	10	2021
Falcone, Favole della buonanotte	10	2016
6–10 years		
Rodari, Favole al telefono	69	1995 [1962]
Falchi, Fiabe illustrate per bambini	10	2011
Iannicello, Racconti e favole per bambini	4	2017
Coletta (a cura di), C'erano una volta... le Favole	14	2008
Romani, Favole per bambini	3	2020–2021
Tot.	200	

Table 6: Composition of the HUM fable subcorpus grouped by age range.

SYSTEM
Il tuo compito è generare testi narrativi.
USER
Scrivi una favola per bambini {x}

Table 7: Prompt for fables generation. The parameter `span_range` is defined as (3-5 anni) and (6-10 anni).

SYSTEM
Il tuo compito è tradurre storie dall'Inglese all'Italiano. Mantieni tono, significato e stile originali. Rispondi solo con la storia tradotta in Italiano.
USER
Storia: {fable}

Table 8: Prompt for fables translation in Italian.

SYSTEM
Your task is to translate stories into Italian. Preserve the original tone, meaning, and style. Reply only with the translated story in Italian.
USER
Story: {fable}

Table 9: Prompt for fables translation in English.