

TailNLG: A Multilingual Benchmark Addressing Verbalization of Long-Tail Entities

Anonymous ACL submission

Abstract

The automatic verbalization of structured knowledge is a key task for making knowledge graphs accessible to non-expert users and for supporting retrieval-augmented generation systems. Although recent advances in RDF-to-text generation have improved multilingual coverage, little attention has been paid to potential biases in the verbalization of rare entities, frequently known as long-tail entities. In this work, we present the first systematic study of long-tail entities in RDF-to-text generation. We introduce TailNLG¹, a new multilingual benchmark in English, Italian, and Spanish, built from Wikidata and covering entities with varying levels of popularity. We evaluate three different families of large language models in zero-shot settings and compare their performance on rare *versus* common entities, as well as against the established WebNLG benchmark. Our results reveal a consistent bias against long-tail entities: embedding-based scores are lower, and model uncertainty is higher for rare entities. We further show that the impact of long-tail entities varies across models and languages, and that existing evaluation metrics do not consistently capture these differences, highlighting the need for more reliable evaluation frameworks.

1 Introduction

The automatic conversion of structured information into human-readable text is a crucial task for making knowledge accessible to a broad audience. Knowledge graphs (KGs), catalogues, and taxonomies are important sources of information that are typically exploitable only by domain experts, whereas the general public requires this knowledge to be transformed into natural language. This need also extends to language modeling, as retrieval-

¹Full code and Benchmark available here <https://anonymous.4open.science/r/TailNLG-benchmark-B339/README.md>

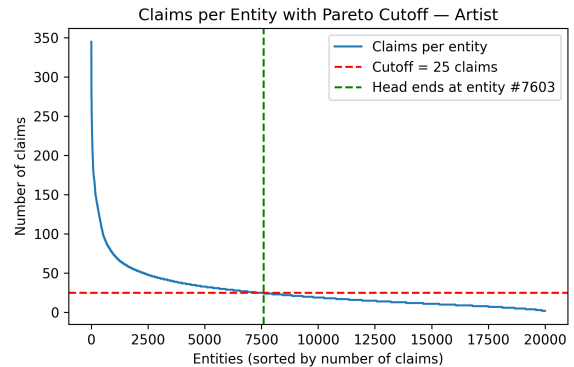


Figure 1: Distribution of entities by claims in a Wikidata category (“Artist”). A small number of entities are associated with many claims, giving rise to a long-tail distribution characterized by many entities with few relations. The red dashed line indicates the Pareto cut-off, while the green dashed line marks the head-tail threshold

augmented generation (RAG) technologies leverage textual databases to improve the safety and accountability of large language models (LLMs).

Despite advancements in this field and the growing number of multilingual resources (Ferreira et al., 2018; Shimorina et al., 2019; Cripwell et al., 2023; Oliverio et al., 2025; Ramón-Ferrer et al., 2025), there is a lack of work aimed at uncovering biases in the verbalization of structured knowledge (Blodgett et al., 2020). Specifically, an open challenge is to gain a deeper understanding of how LLMs handle rare entities, frequently known as *long-tail* entities, and whether their performance differs when compared to well-known entities. While it has been demonstrated that LLMs exhibit performance drops when handling rare entities in tasks such as entity linking (Boscariol et al., 2024), representation of factual knowledge (Sun et al., 2024), and question answering (Hogan et al., 2025), no prior work has conducted an in-depth investigation of how different models handle long-tail

061 entities in RDF-to-text generation.

062 In this work, we tackle this gap by presenting the
063 first systematic study of long-tail entities in RDF-
064 to-text verbalization. Acknowledging the socio-
065 cultural factors that can determine the rarity of an
066 entity (Adams et al., 2019; Stranisci et al., 2023),
067 we adopt a multilingual approach. Our study in-
068 troduces **TailNLG**, a novel multilingual RDF-to-
069 text benchmark covering three languages (English,
070 Italian, and Spanish). The dataset contains triples
071 extracted from Wikidata (Vrandečić and Krötzsch,
072 2014), including entities from different categories
073 and spanning varying levels of popularity, from
074 top-head (most well-known) to long-tail (less well-
075 known) entities.

076 Through the creation of TailNLG, our work ad-
077 dresses the following research questions:

- 078 • **[RQ1]** Do LLMs perform differently in the
079 verbalization of rare entities, and does lan-
080 guage impact the verbalization output?
- 081 • **[RQ2]** What is the impact of LLMs’ prior
082 knowledge on the verbalization of rare *versus*
083 common entities, and how does model per-
084 formance differ when comparing a long-tail
085 benchmark (TailNLG) with a well-established
086 resource such as WebNLG?

087 To this end, three families of LLMs are assessed
088 on the verbalization of long-tail entities in two set-
089 tings: (i) comparing them with head entities within
090 the TailNLG corpus, and (ii) comparing them with
091 WebNLG (Gardent et al., 2017b), the established
092 benchmark for the RDF-to-text task. We assess per-
093 formance with a diverse set of automatic measures,
094 capturing semantic similarity, surface-level over-
095 lap with reference texts, and perplexity to quantify
096 model uncertainty. Our results show that LLMs
097 have a systematic bias against long-tail entities in
098 both evaluation settings. The embedding-based
099 metrics tend to be lower for long-tail entities and
100 models’ uncertainty in text generation of triple con-
101 taining long-tail entities is almost always higher.

102 In addition, language-specific bias plays also a
103 role in the verbalization of long-tail entities. Anal-
104 ysis of perplexity scores reveals that embedding-
105 based and overlap-based metrics fail to systemati-
106 cally quantify their impact on verbalization, demon-
107 strating that a more reliable evaluation framework
108 of fairness in RDF-to-Text is needed. Finally, we
109 observe notable variability across models, indicat-

ing that different LLMs may be affected by differ- 110
ent biases. 111

2 Related Work 112

2.1 Long-tail 113

Long-tail entities are generally understood as those 114
that occur with low frequency across large scale 115
data sources, including training corpora (Kandpal 116
et al., 2023), Wikipedia (Mallen et al., 2023), and 117
large knowledge bases such as Wikidata (Kumar 118
et al., 2024). Recent work has shown that LLMs 119
exhibit a marked decline in performance when 120
processing long tail inputs (Graciotti et al., 2025; 121
Hogan et al., 2025; Li et al., 2024; Sun et al., 2024) 122
and that LLMs memorization is highly influenced 123
by the frequency of information in the pretraining 124
data (Kandpal et al., 2023). Several task-specific 125
evaluations further confirm this pattern. Boscar- 126
iol et al. (2024) find that entity linking is partic- 127
ularly difficult for rare entities when comparing 128
various LLMs with ReLiK (Orlando et al., 2024), a 129
state-of-the-art (SoTA) entity linking and relation 130
extraction system. GRADES, an evaluation frame- 131
work for graph-based question answering (Draetta 132
et al., 2025), shows that SoTA models struggle 133
with rare entities at multiple stages; in particular, 134
verbalization suffers from limited semantic under- 135
standing of long-tail entities. In knowledge extrac- 136
tion, Graciotti et al. (2025) introduce KEMISTO, 137
a benchmark centered on low-popularity entities, 138
and demonstrate that LLMs perform significantly 139
worse on it and exhibit systematic biased failures. 140
A similar pattern emerges in question answering: 141
using a benchmark spanning head, torso, and tail 142
entities, Sun et al. (2024) show that accuracy con- 143
sistently declines as entity frequency decreases, 144
confirming that limited training exposure hampers 145
LLMs’ knowledge of long-tail entities. 146

Despite increasing attention to long-tail evalu- 147
ation, and some insight emerged from previous 148
studies no existing work specifically examines how 149
well LLMs verbalize information when required to 150
handle rare entities. 151

2.2 RDF-to-Text 152

The task of verbalizing RDF triples in NLP has 153
evolved substantially, driven both by template- 154
based NLG systems and more recent neural ap- 155
proaches. Early work focused on hand-crafted or 156
automatically induced templates, which provided 157
high precision and strong control over the output, 158

but suffered from limited scalability and reduced the ability to generalize beyond predefined patterns (Duma and Klein, 2013).

With the advent of neural encoder–decoder architectures, particularly LSTM-based triple encoders and later transformer models, RDF-to-text generation became considerably more flexible, enabling fluent and information-rich verbalizations that better capture lexical and syntactic variation (Distiawan et al., 2018; Oliverio et al., 2024). Nevertheless, challenges related to factual consistency and robustness to unseen entities have motivated a series of enhancements to transformer-based systems. For example, integrating reward signals derived from information-extraction models has been shown to improve factual accuracy (Gao et al., 2021), while leveraging large external corpora, pre-training noise processes, and data augmentation strategies leads to more resilient generalization, particularly in low-resource or zero-shot scenarios (Montella et al., 2020; Zhang et al., 2023).

To support systematic development and evaluation of this task, the WebNLG benchmark has emerged as one of the primary datasets for RDF-to-text, providing multiple versions, multilingual extensions, and low-resource configurations that have enabled extensive comparative studies across approaches and languages (Gardent et al., 2017b; Castro Ferreira et al., 2020).

3 The TailNLG Benchmark

Since no benchmark currently targets long-tail triple–verbalization pairs, we introduce the first dataset specifically designed for long-tail verbalization. Inspired by Gardent et al. (2017b), we created the TailNLG benchmark following the same procedure the authors employed, from the entities extraction strategy to the verbalization, aiming at facilitating the exploitation of the data adhering to a standard structure. Constructing the benchmark in three languages adds another layer of complexity, as many rare entities are not consistently represented across multilingual resources. We prioritized high data quality by incorporating human evaluation and correction at every stage of the pipeline.

3.1 Entities Selection

To be consistent with previous resources, we followed the extraction methodology proposed by Gardent et al. (2017a) by selecting the same entity categories. Although there is not yet a widely ac-

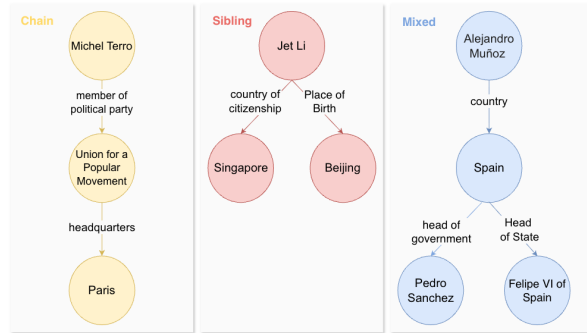


Figure 2: Examples of Chain, Sibling and Mixed configuration of Data units.

cepted and unambiguous definition of long-tail entities (Jiang and Joshi, 2024; Boscaroli et al., 2024; Kumar et al., 2024), following Sun et al. (2024) we took into consideration two criteria in defining long-tail entities: (i) the number of claims (Wikidata relations) linked to a certain entity and (ii) the number of Wikipedia pages an entity has.

Firstly, we set a threshold based on the number of claims linked to each entity. As noted by Kumar et al. (2024), large-scale resources such as Wikidata typically follow a Zipf-like distribution, where a small number of entities are associated with many claims while the vast majority with few. Based on a subset of 20,000 entities for each category (full list of categories in Table 4 in Appendix B) we define the long-tail threshold using a Pareto cut-off: we compute the minimum number of claims required for an entity to fall within the top 20% of entities that account for 80% of all claims (the head). Entities below this threshold are classified as long-tail. Since claim distributions differ across categories, we compute the cut-off separately for each category.

Figure 1 shows an example of the threshold for the categories Artists: a small number of entities exhibit a very high number of claims, while the vast majority fall into a long-tail distribution.

To further refine the set of long-tail entities below the Pareto cut-off, we additionally select those that have a Wikipedia page only in one of the three target languages (e.g., entities with wikipedia page in English but not in Italian and Spanish).

In summary, we categorize entities into: (i) **Long-tail entities**, whose number of claims falls below the Pareto cut-off, such as *Doğan Güzel* (Q268285), a Turkish cartoonist; and (ii) **Head entities**, whose number of claims exceeds the Pareto cut-off, such as *London* (Q84) or *Google* (Q95).

3.2 Triples Extraction

Once the sets of long-tail and head entities were defined, the next step was to extract data units for each entity. We define a data unit as a set of RDF triples, each consisting of a subject, a predicate, and an object (e.g., Jet Li, place of birth, Beijing). To obtain diverse data units that vary in relation type and size, we followed Perez-Beltrachini et al. (2016) and extracted units with different relational configurations (Figure 2):

- **Chain:** the object of one triple serves as the subject of another;
- **Sibling:** the triples share the same subject;
- **Mixed:** the data unit contains both sibling and chain relations.

For each QID (Wikidata entity ID), we automatically extracted from Wikidata, via a SPARQL query, a set of 10 data units per entity in the three languages. Since not all triples were available in all three languages, we further filtered the data units to ensure the alignment across languages, resulting in a set of parallel data units as exemplified below:

EN: [Solar, creator, Paul S. Newman], [Paul S. Newman, place of birth, Manhattan], [Paul S. Newman, award received, Inkpot Award]

ES: [Solar, creator, Paul S. Newman], [Paul S. Newman, lugar de nacimiento, Manhattan], [Paul S. Newman, premio recibido, Premio Inkpot]

IT: [Solar, creatore, Paul S. Newman], [Paul S. Newman, luogo di nascita, Manhattan], [Paul S. Newman, premio ricevuto, Premio Inkpot]

We finally obtain 689 data units per language, divided into 408 long-tail and 218 head entities. The complete distribution per category and type is in Appendix B.

3.3 Verbalization

To build the benchmark, we aimed to produce $\langle \text{Data_Unit}, \text{Verbalization} \rangle$ pairs for all data in all languages. Because manual verbalization is time consuming we adopted a two step approach: we first manually verbalized a subset of data units for each language, then cross-translated them into the remaining languages using SoTA machine translation models in a human-in-the-loop setup.

The verbalization process was carried out on a subset of 689 Data units engaging nine volunteers

native speakers (one for Spanish, one for English, and seven for Italian). Annotators were instructed to verbalize the triples as naturally as possible, producing grammatically correct and fluent sentences while following guidelines designed to promote naturalness and ensure data homogeneity. All guidelines and some example of verbalization are reported in Appendix A.1.

Annotators could skip a data unit if they judged it incorrect or incomplete. After verbalization, a validation phase followed in which the annotators cross-checked each other’s output, confirming or revising the initial sentences (full guidelines in Appendix A.2). Ambiguous cases were discussed collectively before a final decision was made. Following this verbalization and cross-check process, 73 data units were excluded from the dataset.

3.4 Translation

The manual verbalization phase produced a total of 616 gold verbalizations, distributed as 130 in English, 142 in Spanish, and 344 in Italian². To create a fully parallel dataset we adopted an approach based on Machine Translation (MT) and Automatic Post-Editing (APE), as proposed by Aditya Hari et al. (2023).

Drawing on earlier studies on cross-lingual extensions of triples-to-text datasets (Oliverio et al., 2025; Ramón-Ferrer et al., 2025), we leveraged DeepL for Automatic Translation³. For each manually created verbalization in a source language (X) corresponding translations were automatically generated in the other two target languages (Y and Z), resulting in a total of 1,232 automatic generated translations.

To mitigate known issues in MT, such as omissions and hallucinations (Jurafsky and Martin, 2025), and to ensure data quality, we introduced an APE phase aimed at correcting errors generated during the MT step. To identify errors in the translations, we leveraged the 3.5B version of xCOMET Guerreiro et al. (2024)⁴, a multilingual model for quality estimation and error span extraction. Given the absence of reference translations, we evaluated the model in a reference-free setting. For each input, the model outputs a sentence-level score (ranging from *Excellent* to *Weak*), as well as an error spans detected field, which identifies potential er-

²This linguistic imbalance is due to the varying availability of expert native-speaking annotators.

³<https://www.deepl.com/translator>

⁴<https://huggingface.com/Unbabel/XCOMET-XL>

rors and their positions within the sentence.

For the error correction phase, we adopted xTOWER⁵, a multilingual decoder-only model proposed by Treviso et al. (2024), designed to provide natural-language explanations of errors present in a translation and suggest a corrected version of the sentence. An example of the prompt and the corresponding output are provided in the Appendix C.1. We perform error detection and correction on all the MT translated data; the resulting output is a parallel dataset consisting of data units and verbalizations in Italian, Spanish and English as follows:

```
<lex quality="gold" lid="Id1" lang="en">
  Themistocles was born in Turkey and
  he was a citizen of the Classical Athens.
  Themistocles's spouse is Archippe.
</lex>
<lex quality="silver" lid="Id2" lang="es">
  Temístocles nació en Turquía y fue
  ciudadano de la Atenas clásica.
  El cónyuge de Temístocles es Arquippe.
</lex>
<lex quality="silver" lid="Id3" lang="it">
  Temistocle è nato in Turchia ed era un
  cittadino dell'Atene classica.
  La consorte di Temistocle era Archippe.
</lex>
```

To ensure the high quality of the dataset, we conducted a manual evaluation on a subset, selecting those translations whose xCOMET *Sentence-level score* was rated below *good* (i.e., *moderate* or *weak*). We sampled 50 verbalizations for each source-target language pair and involved eight volunteer native speakers (two with high competence in Spanish and English, two in Italian and Spanish, and two in Italian and English). Each source-target language pair was annotated by two different annotators following a widely adopted taxonomy for the evaluation of Accuracy and Fluency (Lommel et al., 2014) (full evaluation guidelines in Appendix C.2). In addition, we asked the annotators to propose a corrected version of the automatically translated verbalization to increase the number of gold references in TailNLG. We use Spearman (Spearman, 1961) and Persons (Pearson, 1895) to measure the Inter annotator agreement (IAA) after aggregating Accuracy and Fluency. The results for each language pair are reported in Appendix C.3. Although inter-annotator agreement (IAA) shows a moderate positive correlation overall, agreement is low for IT→EN and ES→IT pair. This indicates

⁵<https://huggingface.com/sardinelab/xTower13B>

a high degree of subjectivity in the task, with annotators showing divergent interpretations. These findings further underline that even tasks such as machine translation can reflect differences in human judgment. The final verbalization distribution is available in Table 1.

Language	Gold	Silver	Total
English	168	447	615
Spanish	173	442	615
Italian	384	231	615
Total	725	1,120	1,845

Table 1: Distribution of gold and silver verbalizations by language in TailNLG.

4 Methodology

To study multilingual triple-to-text verbalization, we leverage three families of open state-of-the-art multilingual instruction-tuned models, each evaluated at two comparable sizes: Llama-3.2-3B-Instruct⁶, Llama-3.1-8B-Instruct⁷, Qwen2.5-3B-Instruct⁸, Qwen2.5-7B-Instruct⁹, Gemma-3-4B-IT¹⁰, and Gemma-3-12B-IT¹¹. We focus on these model families because they are multilingual, widely used, and well established in the literature, openly available, and offered in moderate parameter ranges that enable efficient and reproducible experimentation.

4.1 Experimental Setup

Given a set of RDF-style triples (s, p, o) , the system produces a single-paragraph natural language description in English, Spanish, or Italian. Inputs are serialized by explicitly listing each triple as a structured item (subject, predicate, object) in the prompt. We run experiments on two datasets: (i) WebNLG, used as a baseline benchmark, and (ii) TailNLG, used to analyze performance on data unit containing long-tail entities. Since our goal is to study biases and examine how LLMs handle rare entities, we deliberately adopt a zero-shot approach for both the configurations to avoid influencing

⁶<https://huggingface.co/meta-llama/Llama-3.2-3B-Instruct>

⁷<https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>

⁸<https://huggingface.co/Qwen/Qwen2.5-3B-Instruct>

⁹<https://huggingface.co/Qwen/Qwen2.5-7B-Instruct>

¹⁰<https://huggingface.co/google/gemma-3-4b-it>

¹¹<https://huggingface.co/google/gemma-3-12b-it>

428 model outputs with external examples. Prompts
429 are written English, Spanish, or Italian depending
430 on the target verbalization. Full prompts are avail-
431 able in Appendix E.

432 For each input we generate three candidate ver-
433 balizations to capture decoding variability. We
434 use stochastic decoding with temperature sampling
435 (temperature = 0.7), which increases output diver-
436 sity compared to greedy decoding while keeping
437 generations coherent. Generations are bounded by
438 a fixed maximum length (`max_new_tokens = 256`)
439 to prevent overly long outputs and to standardize
440 comparisons across models and settings. All out-
441 puts are stored together with the corresponding
442 prompts and metadata. Experiments were con-
443 ducted on a server with two Intel Xeon Gold 5418Y
444 CPUs (48 physical cores each, 96 threads total),
445 around 504 GB RAM, and two NVIDIA L40 GPUs
446 with approximately 46 GB VRAM each.

447 4.2 Evaluation

448 In evaluating the results, we compare multiple
449 metrics designed to capture three characteristics
450 of the generated output: (i) semantic similarity
451 (embedding-based metrics), (ii) textual similar-
452 ity (overlap-based metrics), and (iii) model uncer-
453 tainty in text generation¹². We use BERTScore
454 (Zhang et al., 2019) to evaluate semantic simi-
455 larity, computing pairwise cosine similarity be-
456 tween candidate and reference texts with con-
457 textual embeddings. Specifically, we adopt the
458 rescaled-with-baseline variant ($BERTScore_r$) for
459 more interpretable scores, using a multilingual
460 BERT model¹³ to ensure comparability across lan-
461 guages. To evaluate textual similarity, we adopt:
462 BLEU (Papineni et al., 2002), chrF++ (Popović,
463 2017), and three variants of ROUGE (Lin, 2004),
464 namely $ROUGE_1$, which measures unigram (word-
465 level) overlap; $ROUGE_2$, which measures bigram
466 overlap; and $ROUGE_L$, which captures sentence-
467 level structural similarity by identifying the longest
468 common subsequences. To quantify models un-
469 certainty in triple verbalization, we computed the
470 average Perplexity (PPL) (Jelinek et al., 1977) in-
471 terpreting lower perplexity values as a proxy of the
472 higher confidence of a model in text generation.

473 To assess the statistical significance of each met-
474 ric, we applied the Wilcoxon–Mann–Whitney test
475 (Fagerland and Sandvik, 2009) on triple verbaliza-

¹²Evaluations are performed using Hugging Face packages

¹³[https://huggingface.co/google-bert/
bert-base-multilingual-cased](https://huggingface.co/google-bert/bert-base-multilingual-cased)

476 tion of the same model on different populations.
477 We repeated the test ten times on random samples
478 of 500 predictions and averaged the results. We
479 performed the test in two settings: long-tail *ver-*
480 *sus* head verbalization in the TailNLG corpus and
481 long-tail verbalization in the TailNLG corpus *ver-*
482 *sus* WebNLG verbalization.

483 5 Results

484 In this section, we present results comparing long-
485 tail and head entities, as well as across languages
486 (RQ1). We also examine the impact of LLMs’ prior
487 knowledge on verbalizing rare and well-known en-
488 tities by comparing PPL scores on TailNLG and
489 WebNLG, assessing whether models show greater
490 uncertainty for rare entities (RQ2). Following (Du
491 et al., 2024), we define prior knowledge as the in-
492 formation acquired during pretraining.

493 5.1 Do LLMs perform differently in the 494 verbalization of rare entities, and does 495 language impact the verbalization output?

496 Table 2 reports the metric scores for triples con-
497 taining long-tail and head entities, broken down
498 by language (English, Italian and Spanish), for the
499 bigger size models. Complete results are provided
500 in Appendix F. In all cases, the embedding-based
501 metric $BERTScore_r$, scores higher on head triples.
502 Conversely, results from chrF++, which correlates
503 with morphological accuracy, and ROUGE are less
504 consistent: in some cases they score higher on long-
505 tail triples, while in others on head triples. In addi-
506 tion, when the results for $BERTScore_r$ are statisti-
507 cally significant, the scores are consistently higher
508 for head entities than for long-tail entities, which
509 may be linked to greater semantic awareness when
510 verbalizing triples involving head entities. In con-
511 trast, the overlap-based metrics exhibit the opposite
512 trend: when the results are statistically significant,
513 higher scores are observed for triples containing
514 long-tail entities. The average PPL provides clear
515 insights, particularly for Qwen and Llama models.
516 The results suggest that verbalizing triples with
517 long-tail entities leads to higher PPL, indicating
518 that those entities constitute a substantially more
519 challenging prediction setting for these models. In
520 contrast, the PPL scores for Gemma models are
521 less clear and always lower for long-tail entities.
522 Investigating whether language influences the gen-
523 erated output, we observe that across all settings
524 performance is higher for English, for both data units

Type	Lang	BERT _r ↑	BLEU ↑	chrF ↑	R-1 ↑	R-2 ↑	R-L ↑	PPL ↓
Llama-8B								
Long-tail	en	0.61	0.25	60.39	0.66	0.43	0.56	93.36
	es	0.58	0.17	53.30	0.58	0.35	0.49	222.16
	it	0.46	0.13	47.73*	0.49*	0.29*	0.43*	88.06
Head	en	0.65*	0.26	61.20	0.67	0.44	0.57	59.03*
	es	0.60	0.19	53.23	0.58	0.36	0.48	94.52*
	it	0.45	0.13	45.58	0.46	0.27	0.40	68.84*
Gemma-12B								
Long-tail	en	0.63	0.26	61.65	0.67	0.44	0.58	330.02*
	es	0.62	0.21	58.80	0.62	0.40	0.53	611.42*
	it	0.61	0.24*	60.68*	0.64*	0.42*	0.55*	425.27*
Head	en	0.66*	0.27	61.72	0.68	0.45	0.57	397.46
	es	0.64*	0.22	57.51	0.61	0.39	0.51	926.65
	it	0.61	0.20	56.97	0.59	0.38	0.50	1068.31
Qwen-7B								
Long-tail	en	0.65	0.29	62.47	0.69	0.46	0.58	134.13
	es	0.60	0.16	53.42	0.57	0.34	0.49	181.75
	it	0.44	0.11	46.82	0.47	0.26	0.41	100.26
Head	en	0.69*	0.29	63.06	0.70	0.47	0.59	98.05*
	es	0.61	0.16	51.83	0.55	0.32	0.46	102.36*
	it	0.46	0.10	45.35	0.45	0.25	0.39	63.48*

Table 2: Results on TailNLG long-tail vs head entities verbalizations for Llama-3.1-8B, Gemma2.5-12B 12B, Qwen2.5-7B by data type and language. Metrics include BERTScore Rescaled (BERT_r), BLEU, chrF++ (chrF), ROUGE₁ (R-1), ROUGE₂ (R-2), ROUGE_L (R-L) and perplexity (PPL). Asterisks (*) denote a statistically significant difference ($p < 0.05$).

containing long-tail and head entities and consistently lower for Italian. This suggests that English content is handled more effectively by the models. In addition, for Llama and Qwen, average PPL scores for Italian are consistently higher than those for English and Spanish, indicating that the models could be affected by language bias and that outputs in more widely represented languages tend to be more accurate.

5.2 [RQ2] How does model performance differ when comparing the TailNLG benchmark with WebNLG?

Table 3 reports the average performance of each model on the WebNLG test set and on long-tail triples from TailNLG. All statistically significant results indicate that embedding- and overlap-based metrics achieve higher scores on WebNLG than in the long-tail setting. This trend suggests that models generate more accurate and lexically aligned verbalizations for well-known entities, whereas performance tends to degrade for long-tail entities.

Across all models, PPL is even lower for WebNLG verbalizations, with statistical significance, supporting the hypothesis that WebNLG entities and triples are more familiar to the models, potentially due to exposure during pre-training.

Overall, these results highlight the impact of

Type	BERT _r ↑	BLEU ↑	chrF ↑	R-1 ↑	R-2 ↑	R-L ↑	PPL ↓
Llama-8B							
Long-tail	0.55	0.18	53.81	0.58	0.36	0.50	134.53
WebNLG	0.53	0.21*	53.77	0.65*	0.43*	0.53*	68.23*
Gemma-12B							
Long-tail	0.62	0.24	60.38	0.64	0.42	0.55	455.57
WebNLG	0.65*	0.30*	63.63*	0.73*	0.50*	0.59*	290.10*
Qwen-7B							
Long-tail	0.56	0.19	54.24	0.58	0.35	0.49	138.71
WebNLG	0.57	0.24*	56.49	0.66*	0.44*	0.54*	95.54*

Table 3: Performance comparison of models between long-tail triples from TailNLG and the WebNLG test set. Metrics include BERTScore Rescaled (BERT_r), BLEU, chrF++ (chrF), ROUGE-1 (R-1), ROUGE-2 (R-2), ROUGE-L (R-L), and perplexity (PPL). Asterisks (*) denote a statistically significant difference ($p < 0.05$).

LLMs’ prior knowledge on the verbalization task. When entities and triples are well represented in the models’ pre-training data, LLMs tend to produce verbalizations that are more closely aligned with the ground truth. Conversely, for long-tail entities, limited prior exposure leads to a degradation in both surface-level quality and model confidence, suggesting that prior knowledge plays a crucial role in the effective verbalization of structured data.

In addition, since models tend to achieve higher scores when synthetic data are used as ground truth

(Sarti and Nissim, 2024), our findings are further reinforced. Even when comparing model performance on the English WebNLG test set (entirely gold) with the long-tail portion of TailNLG (containing both gold and silver references), models still achieve higher scores on WebNLG. For further details, we refer the reader to Appendix F.

6 Discussion

Results show that, across all models, verbalization of triples that contains long-tail entities is more challenging (RQ1). This suggests that LLMs suffer from entity bias, with widely represented entities being processed more easily and producing more accurate outputs. The fact that in some cases overlap-based metrics are higher for data units containing long-tail entities may suggest that the verbalization of head entities contains more data and key content relative to the reference, and that prior knowledge could play a role in shaping the generated output adding information. For example, the Italian verbalization of the data unit `Karl Marx | spouse | Jenny von Westphalen; Karl Marx | country of citizenship | statelessness` is “Karl Marx è stato apatico alla cittadinanza, preferendo lo status di apolidia. Karl Marx ebbe come consorte Jenny von Westphalen” (English translation: Karl Marx was apathetic toward citizenship, preferring to remain stateless. Karl Marx’s spouse was Jenny von Westphalen), which results in a more verbose output that differs from the target sentence. The results concerning language strongly confirm the hypothesis that models also suffer from language bias, and that verbalizations in English are more accurate. These results may also be influenced by the fact that both Italian and Spanish have richer morphologies than English. For example, the Italian verbalization of the triples `Amanda Ammann | country of citizenship | Switzerland; Amanda Ammann | occupation | model` results in “Amanda Ammann è una cittadina svizzera e lavora come modello” (English translation: Amanda Ammann is a Swiss citizen and works as a model), which exhibits a gender disagreement between the first part of the sentence (“cittadina svizzera”) and the final noun (“modello”). Generally Italian output tend to be less accurate and suffer from some additions and omissions.

Another important finding emerging from the results is the high variability across evaluation metrics. In particular, reference-based metrics are not

always able to detect the presence of bias or provide reliable scores, as becomes evident from the comparison with PPL, which yields more consistent results. In addition, the analysis suggests that different models may be affected by different types of bias; in particular, Gemma exhibits a higher PPL for triples containing long-tail entities.

The prior-knowledge experiment confirms our hypothesis that models perform better on a benchmark such as WebNLG compared to TailNLG (RQ2), as the models may have been exposed during pre-training to information related to the WebNLG benchmark or to the entities it contains. This finding is supported, with statistical significance, by both reference-based evaluation metrics and PPL.

7 Conclusion and Future Works

In this work, we investigated the performance of three state-of-the-art language models’ families in depth for verbalizing long-tail entities in a multilingual setting. To this end, we introduced TailNLG, the first benchmark composed of triple–verbalization pairs covering entities that range from rare to well-represented. Our results show that the leveraged models suffer from both language and entity biases, and that they perform better on a widely exploited benchmark such as WebNLG compared to TailNLG. Consistent with prior work, our findings highlight that current models struggle to handle long-tail entities, revealing limitations that make tasks involving less-represented entities more difficult and less reliable.

As future work, we plan to first enrich the TailNLG benchmark by adding more languages and data units. To further refine the dataset, we also aim to distinguish entity types not only based on their popularity (i.e., number of claims) but also by incorporating culturally specific data units to explore the models’ cultural bias. Finally, we plan to enrich the dataset with a gold standard that takes into account the diverse perspectives of annotators, enabling the investigation of human label variation in the verbalization task.

Limitations

We acknowledge that building a new high-quality benchmark is challenging, as it requires achieving an appropriate balance between human annotation and automatic processing steps. In this work, we consistently paired the automatic phase with hu-

man assessment in an effort to ensure the highest possible quality. Nevertheless, the presence of *silver* verbalizations in the benchmark constitutes a limitation, since machine translation tools, although reliable, may not fully align with human judgments.

Second, by leveraging a well-known resource such as Wikidata, the variability of available claims is limited. As a result, models may already be familiar with many of the relations, which also recur across different data units.

Finally, in defining long-tail entities, we adopted a simplified approach for experimental purposes, considering only a quantitative parameter (number of claims). However, we are aware that entity underrepresentation is a multifaceted issue that involves multiple dimensions, both vertical (numerosity) and horizontal (cultural aspects). In future work, it would be valuable to consider additional perspectives and to define long-tail entities using a more comprehensive characterization.

Ethical considerations

Our work relies on knowledge extracted from Wikidata, whose community of contributors has recognized the underrepresentation of non-Western and woman contributors. This has an impact on how knowledge is shaped and which entities suffer underrepresentation. Adopting a purely statistical approach to the definition of *long-tail* entities might have resulted in overlooking these cultural factors. Nevertheless, this approach was preferred to the adoption of theoretically grounded models of underrepresentation, as the latter could have introduced additional research bias in corpus creation.

All annotations have been performed voluntarily by the members of our research group without relying on external annotation platforms.

During the preparation of this work, the authors used AI tools in order to: Grammar and spelling check. After using this tool, the authors reviewed and edited the content as needed and take full responsibility for the publication's content.

Acknowledgments

We thank all the individuals who contributed to this project by supporting the verbalization phase and the evaluation of translations. Without their help, we would not have been able to achieve the desired level of data quality.

References

- Julia Adams, Hannah Brückner, and Cambria Naslund. 2019. Who counts as a notable sociologist on wikipedia? gender, race, and the “professor test”. *Socius*, 5:2378023118823946.
- Kancharla Aditya Hari, Bhavyajeet Singh, Anubhav Sharma, and Vasudeva Varma. 2023. [WebNLG challenge 2023: Domain adaptive machine translation for low-resource multilingual RDF-to-text generation \(WebNLG 2023\)](#). In *Proceedings of the Workshop on Multimodal, Multilingual Natural Language Generation and Multilingual WebNLG Challenge (MM-NLG 2023)*, pages 93–94, Prague, Czech Republic. Association for Computational Linguistics.
- Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (technology) is power: A critical survey of “bias” in nlp. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476.
- Marta Boscaroli, Luana Bulla, Lia Draetta, Beatrice Fiumanò, Emanuele Lenzi, and Leonardo Piano. 2024. Evaluation of llms on long-tail entity linking in historical documents (short paper). In *EKAW (Companion)*.
- Thiago Castro Ferreira, Claire Gardent, Nikolai Ilinykh, Chris van der Lee, Simon Mille, Diego Moussallem, and Anastasia Shimorina, editors. 2020. *Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)*. Association for Computational Linguistics, Dublin, Ireland (Virtual).
- Liam Cripwell, Anya Belz, Claire Gardent, Albert Gatt, Claudia Borg, Marthese Borg, John Judge, Michela Lorandi, Anna Nikiforovskaya, and William Soto Martinez. 2023. The 2023 webnlg shared task on low resource languages. overview and evaluation results (webnlg 2023). In *Association for Computational Linguistics*, page 55–66.
- Bayu Distiawan, Jianzhong Qi, Rui Zhang, and Wei Wang. 2018. Gtr-1stm: A triple encoder for sentence generation from rdf data. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1627–1637.
- Lia Draetta, Marco Antonio Stranisci, Flaviana Corallo, Pier Felice Balestrucci, Michael Oliverio, Rossana Damiano, and Alessandro Mazzei. 2025. Beyond the metrics: an investigation into the reliability of evaluation metrics for domain specific graph-based question answering.
- Kevin Du, Vésteinn Snæbjarnarson, Niklas Stoehr, Jennifer White, Aaron Schein, and Ryan Cotterell. 2024. [Context versus prior knowledge in language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13211–13235, Bangkok, Thailand. Association for Computational Linguistics.

767	Daniel Duma and Ewan Klein. 2013. Generating natural language from linked data: Unsupervised template extraction. In <i>Proceedings of the 10th International Conference on Computational Semantics (IWCS 2013)–Long Papers</i> , pages 83–94. ASSOC COMPUTATIONAL LINGUISTICS-ACL.	
768		
769		
770		
771		
772		
773	Morten W Fagerland and Leiv Sandvik. 2009. The wilcoxon–mann–whitney test under scrutiny. <i>Statistics in medicine</i> , 28(10):1487–1497.	
774		
775		
776	Thiago Castro Ferreira, Diego Moussallem, Emiel Kraemer, and Sander Wubben. 2018. Enriching the webnlg corpus. In <i>Association for Computational Linguistics</i> , page 171–176.	
777		
778		
779		
780	Hanning Gao, Lingfei Wu, Po Hu, Zhihua Wei, Fangli Xu, and Bo Long. 2021. Rdf-to-text generation with reinforcement learning based graph-augmented structural neural encoders. <i>CoRR</i> .	
781		
782		
783		
784	Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017a. Creating training corpora for nlg micro-planning. In <i>55th Annual Meeting of the Association for Computational Linguistics, ACL 2017</i> , pages 179–188. Association for Computational Linguistics (ACL).	
785		
786		
787		
788		
789		
790	Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017b. The webnlg challenge: Generating text from rdf data. In <i>10th International Conference on Natural Language Generation</i> , pages 124–133. ACL Anthology.	
791		
792		
793		
794		
795	Arianna Graciotti, Leonardo Piano, Nicolas Lazzari, Enrico Daga, Rocco Tripodi, Valentina Presutti, and Livio Pompianu. 2025. Ke-mhisto: Towards a multilingual historical knowledge extraction benchmark for addressing the long-tail problem. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 20316–20339.	
796		
797		
798		
799		
800		
801		
802	Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2024. xCOMET: Transparent machine translation evaluation through fine-grained error detection. <i>Transactions of the Association for Computational Linguistics</i> , 12:979–995.	
803		
804		
805		
806		
807		
808	Aidan Hogan, Xin Luna Dong, Denny Vrandečić, and Gerhard Weikum. 2025. Large language models, knowledge graphs and search engines: A crossroads for answering users’ questions. <i>arXiv preprint arXiv:2501.06699</i> .	
809		
810		
811		
812		
813	Fred Jelinek, Robert L Mercer, Lalit R Bahl, and James K Baker. 1977. Perplexity—a measure of the difficulty of speech recognition tasks. <i>The Journal of the Acoustical Society of America</i> , 62(S1):S63–S63.	
814		
815		
816		
817	Ming Jiang and Mansi Joshi. 2024. CPopQA: Ranking cultural concept popularity by LLMs. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)</i> , pages 615–630, Mexico City, Mexico. Association for Computational Linguistics.	
818		
819		
820		
821		
822		
823		
	Daniel Jurafsky and James H. Martin. 2025. <i>Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition, with Language Models</i> , 3rd edition. Online manuscript released August 24, 2025.	824
		825
		826
		827
		828
	Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In <i>Proceedings of the 40th International Conference on Machine Learning</i> , pages 15696–15707.	829
		830
		831
		832
		833
	Rohan Kumar, Youngmin Kim, Sunitha Ravi, Haitian Sun, Christos Faloutsos, Ruslan Salakhutdinov, and Minji Yoon. 2024. Automatic question-answer generation for long-tail knowledge. <i>arXiv preprint arXiv:2403.01382</i> .	834
		835
		836
		837
		838
	Huihan Li, Yuting Ning, Zeyi Liao, Siyuan Wang, Xiang Lorraine Li, Ximing Lu, Wenting Zhao, Faeze Brahman, Yejin Choi, and Xiang Ren. 2024. In search of the long-tail: Systematic generation of long-tail inferential knowledge via logical rule guided search. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 2348–2370, Miami, Florida, USA. Association for Computational Linguistics.	839
		840
		841
		842
		843
		844
		845
		846
		847
	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	848
		849
		850
	Arle Lommel, Aljoscha Burchardt, Maja Popović, Kim Harris, Eleftherios Avramidis, and Hans Uszkoreit. 2014. Using a new analytic measure for the annotation and analysis of MT errors on real data. In <i>Proceedings of the 17th Annual Conference of the European Association for Machine Translation</i> , pages 165–172, Dubrovnik, Croatia. European Association for Machine Translation.	851
		852
		853
		854
		855
		856
		857
		858
	Alex Mullen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 9802–9822.	859
		860
		861
		862
		863
		864
		865
	Sebastien Montella, Betty Fabre, Tanguy Urvoy, Johannes Heinecke, and Lina M Rojas Barahona. 2020. Denoising pre-training and data augmentation strategies for enhanced rdf verbalization with transformers. In <i>Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)</i> , pages 89–99.	866
		867
		868
		869
		870
		871
		872
	Michael Oliverio, Pier Felice Balestrucci, Alessandro Mazzei, and Valerio Basile. 2025. Webnlg-it: Construction of an aligned rdf-italian corpus through machine translation techniques. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 12073–12083.	873
		874
		875
		876
		877
		878

879	Michael Oliverio, Pier Felice Balestrucci, Alessandro Mazzei, Valerio Basile, and 1 others. 2024. Dipinfo- 880 unito at the gem’24 data-to-text task: Augmenting 881 llms with the split-generate-aggregate pipeline. In 882 <i>Proceedings of the 17th International Natural Lan- 883 guage Generation Conference: Generation Chal- 884 lenges</i> , pages 59–65. Association for Computational 885 Linguistics.	Kai Sun, Yifan Xu, Hanwen Zha, Yue Liu, and Xin Luna Dong. 2024. Head-to-tail: How knowledgeable are 893 large language models (LLMs)? A.K.A. will LLMs 894 replace knowledge graphs? In <i>Proceedings of the 895 2024 Conference of the North American Chapter of 896 the Association for Computational Linguistics: Hu- 897 man Language Technologies (Volume 1: Long Pa- 898 pers)</i> , pages 311–325, Mexico City, Mexico. Associ- 899 ation for Computational Linguistics.	932 933 934 935 936 937 938 939 940
887	Riccardo Orlando, Pere-Lluís Huguet Cabot, Edoardo 888 Barba, and Roberto Navigli. 2024. Relik: Retrieve 889 and link, fast and accurate entity linking and relation 890 extraction on an academic budget. In <i>Findings of the 891 Association for Computational Linguistics ACL 2024</i> , 892 pages 14114–14132.	Marcos V Treviso, Nuno M Guerreiro, Sweta Agrawal, Ricardo Rei, José Pombal, Tania Vaz, Helena Wu, Beatriz Silva, Daan Van Stigt, and Andre Martins. 2024. xTower: A multilingual LLM for explaining 944 and correcting translation errors. In <i>Findings of the 945 Association for Computational Linguistics: EMNLP 946 2024</i> , pages 15222–15239, Miami, Florida, USA. 947 Association for Computational Linguistics.	941 942 943 944 945 946 947 948
893	Kishore Papineni, Salim Roukos, Todd Ward, and Wei- 894 Jing Zhu. 2002. Bleu: a method for automatic evalu- 895 ation of machine translation. In <i>Proceedings of the 896 40th annual meeting of the Association for Computa- 897 tional Linguistics</i> , pages 311–318.	Denny Vrandečić and Markus Krötzsch. 2014. Wiki- data: a free collaborative knowledgebase. <i>Communi- cations of the ACM</i> , 57(10):78–85.	949 950 951
898	Karl Pearson. 1895. Vii. note on regression and inher- 899 itance in the case of two parents. <i>proceedings of the 900 royal society of London</i> , 58(347-352):240–242.	Fan Zhang, Meishan Zhang, Shuang Liu, Yueheng Sun, and Nan Duan. 2023. Enhancing rdf verbalization with descriptive and relational knowledge. <i>ACM Transactions on Asian and Low-Resource Language Information Processing</i> , 22(6):1–18.	952 953 954 955 956
901	Laura Perez-Beltrachini, Rania Sayed, and Claire Gar- 902 dent. 2016. Building RDF content for data-to-text 903 generation. In <i>Proceedings of COLING 2016, the 904 26th International Conference on Computational Lin- 905 guistics: Technical Papers</i> , pages 1493–1502, Osaka, 906 Japan. The COLING 2016 Organizing Committee.	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. BERTscore: Eval- uating text generation with bert. <i>arXiv preprint arXiv:1904.09675</i> .	957 958 959 960
907	Maja Popović. 2017. chrF++: words helping character 908 n-grams. In <i>Proceedings of the second conference on 909 machine translation</i> , pages 612–618.	A Verbalization guidelines	961
910	Virginia Ramón-Ferrer, Carlos Badenes-Olmedo, and 911 Oscar Corcho. 2025. Spanish triple-to-text bench- 912 mark on low-resource large language models.	The annotators were asked to verbalize the data- units following guidelines:	962 963
913	Gabriele Sarti and Malvina Nissim. 2024. IT5: Text- 914 to-text pretraining for Italian language understanding 915 and generation. In <i>Proceedings of the 2024 Joint 916 International Conference on Computational Linguistics, 917 Language Resources and Evaluation (LREC- 918 COLING 2024)</i> , pages 9422–9433, Torino, Italia. 919 ELRA and ICCL.	A.1 Verbalization	964
920	Anastasia Shimorina, Elena Khasanova, and Claire Gar- 921 dent. 2019. Creating a corpus for russian data-to-text 922 generation using neural machine translation and post- 923 editing. In <i>Association for Computational Linguis- 924 tics</i> , pages 44–49.	<ul style="list-style-type: none"> • Verbalize the triples as naturally as possible, aiming to create grammatically correct and fluent sentences. You may choose whether or not to follow the original order of the triples. Please indicate this in the “Order preserved” column by writing YES if the triple is verbal- ized in the same order, or NO if it is not. • You are free to break down the information conveyed by each set of triples in one or more sentences according to your preferences. • If a triple makes no sense (e.g., "Black Racer", "different from", "Black Racer"), you may choose not to verbalize it. In this case, write NA in the annotation column. • If a triple is very difficult or uncertain to in- terpret, you may propose a verbalization and 	965 966 967 968 969 970 971 972 973 974 975 976 977 978 979 980
925	Charles Spearman. 1961. The proof and measurement 926 of association between two things.		
927	Marco Antonio Stranisci, Eleonora Bernasconi, Viviana 928 Patti, Stefano Ferilli, Miguel Ceriani, and Rossana 929 Damiano. 2023. The world literature knowledge 930 graph. In <i>International Semantic Web Conference</i> , 931 pages 435–452. Springer.		

981 leave comments in the Notes column. A col-
 982 league will review uncertain verbalizations
 983 later.

984 • If you have the necessity to check the
 985 meaning or the translation of certain entities
 986 of claim you can use WikiData [https://www.wikidata.org/wiki/Wikidata:](https://www.wikidata.org/wiki/Wikidata:Main_Page)
 987 [Main_Page](https://www.wikidata.org/wiki/Wikidata:Main_Page).
 988

- 989 • Named entities
 - 990 – For named entities, if an equivalent ex-
 - 991 – exists in Italian or Spanish, please use the
 - 992 – correct translation in your verbalization.
 - 993 – If you have doubts about the translation
 - 994 – of a claim, you can find the correspond-
 - 995 – ing version in Spanish or Italian on Wiki-
 - 996 – data.
 - 997 – If you cannot find an equivalent in Ital-
 - 998 – ian or Spanish, do not translate it. In-
 - 999 – stead, flag the triple as problematic in the
 - 1000 – “Notes” column.

1001 **A.2 Validation**

1002 In this phase, the annotators will cross-check and
 1003 review colleagues’ annotations, evaluating whether
 1004 each verbalization is acceptable by answering the
 1005 following questions:

- 1006 • Does the text sound fluent and natural?
- 1007 • Does the text include all and only the informa-
- 1008 – tion from the data?
- 1009 • Is the text written in good language (no
- 1010 – spelling or grammatical mistakes)?

1011 **B Data Unit distribution per Category**

1012 **C Translation**

1013 **C.1 Error correction**

1014 You are provided with a Source, Trans-
 1015 lation, Translation quality analysis, and
 1016 Translation quality score (weak, moder-
 1017 ate, good, excellent). The Translation
 1018 quality analysis contain a translation
 1019 with marked error spans with different
 1020 levels of severity (minor, major or
 1021 critical). Additionally, we may provide
 1022 a ****reference translation****. Given this
 1023 information, generate an explanation for
 1024 each error and a fully correct translation.
 1025 **Italian source:** Osvaldo il coniglio è un

Category	Count
Politician	47
Airport	46
Model	30
SportsTeam	37
Building	47
Company	44
Artist	41
City	47
MeanOfTransportation	20
Food	35
CelestialBody	17
Astronaut	33
Monument	47
ComicsCharacter	48
River	30
University	47
WrittenWork	26
Athlete	47
Type	Count
Long-tail	408
Head	218
Total entries	689

Table 4: Distribution of data entries across categories and entity types.

personaggio dell’opera Lupo Alberto 1026
English translation: Osvaldo the congil 1027
 is a character from the opera Wolf 1028
 Albert. 1029
Translation quality analysis: 1030
 Osvaldo the con<error sever- 1031
 ity="minor">gil</error> is a char- 1032
 acter from the opera<error sever- 1033
 ity="minor">Wolf Albert</error>. 1034
Translation quality score: good 1035

The corresponding output is: 1036

Explanation for error: The term "con- 1037
 gilio" is a typographical error and should 1038
 be corrected to "coniglio," which means 1039
 "rabbit" in Italian. The correct trans- 1040
 lation of "Lupo Alberto" is "Wolfy Al- 1041
 bert," as "Lupo" means "wolf" and "Al- 1042
 berto" is a common name in Italian. The 1043
 translation also misses the gender of "Os- 1044
 valdo," which is masculine, and the cor- 1045
 rect translation should reflect that. Ad- 1046
 ditionally, the translation of "opera" as 1047

1048	"opera" is correct, but the translation of		
1049	"Wolf Albert" as "Wolf Albert" is incor-		
1050	rect because it does not convey the cor-		
1051	rect name of the work. The correct trans-		
1052	lation should reflect the proper name of		
1053	the work, which is "Lupo Alberto."		
1054	C.2 Guidelines Human MT Evaluation		
1055	Annotators were instructed to strictly follow the		
1056	evaluation guidelines and, when necessary, to cor-		
1057	rect the MT output by providing a correct or more		
1058	accurate version.		
1059	Annotate as 1 if the error is present and 0 if not,		
1060	respecting the following taxonomy:		
1061	Accuracy		
1062	Issues related to whether the information content		
1063	of the target is equivalent to the source. Subcate-		
1064	gories:		
1065	• Terminology: Issues related to the use of		
1066	domain-specific terms.		
1067	• Mistranslation: Issues related to the im-		
1068	proper translation of content, including erro-		
1069	neous translation of named entities.		
1070	• Omission: Content present in the source is		
1071	missing in the target.		
1072	• Addition: Content not present in the source		
1073	has been added to the target.		
1074	• Untranslated: Text inappropriately appears		
1075	in the source language.		
1076	Fluency		
1077	Issues related to the linguistic properties of the		
1078	target without relation to its status as a translation.		
1079	Subcategories:		
1080	• Grammar: Issues related to grammatical		
1081	properties of the text.		
1082	• Style: The text shows stylistic problems.		
1083	• Spelling: The text is spelled incorrectly.		
1084	• Typography: Problems related to typographi-		
1085	cal conventions.		
1086	• Unintelligible: Text is garbled or otherwise		
1087	unintelligible; indicates a major breakdown in		
1088	fluency.		
	Corrected Translation		1089
	If the target translation is incorrect, you may pro-		1090
	pose a new translation in the "Corrected translation"		1091
	column.		1092
	Guidelines for Annotators		1093
	• Evaluate whether translating a Named En-		1094
	tity (NE) makes sense in context, taking into		1095
	account cultural and linguistic conventions;		1096
	some NEs may be left untranslated.		1097
	• Named entities should be considered correctly		1098
	translated when it is common in the target		1099
	language to refer to the entity in that way;		1100
	otherwise, mark them as incorrect. Examples:		1101
	– Crist Redemptor → Cristo Redentore →		1102
	Christ the Redeemer [Correct]		1103
	– Lupo Alberto → Wolfy Albert [Incor-		1104
	rect]		1105
	– Louis the Springer → Luigi il Precursore		1106
	[Incorrect]		1107
	Wrong NEs must be annotated in the column		1108
	"Mistranslation".		1109
	• Assess when it is appropriate to keep a term		1110
	in English (or another source language).		1111
	• Pay attention to the correct use of verb tenses.		1112
	• Ensure agreement between singular and plu-		1113
	ral.		1114
	• The translation should be read naturally in the		1115
	target language.		1116
	• Ensure that the meaning is accurately pre-		1117
	served.		1118
	• Annotators should mark 1 for each subcate-		1119
	gory when an issue occurs. If you label 1 ,		1120
	specify why in the "note" column.		1121
	C.3 Manual assessment of translation		1122
	In the table below are reported the Spearman and		1123
	Pearson correlation between the two annotators.		1124
	D TailNLG Benchmark		1125
	Each entry contains the following metadata labels:		1126
	• Category: domain/class of the subject (e.g.,		1127
	Artist, City, Monument).		1128
	• Eid: internal entry id (e.g., Id1).		1129

Source	Target	ρ	r
es	en	0.23	0.25
en	es	0.32*	0.27
it	en	0.15	0.02
en	it	0.56*	0.64
es	it	0.26	0.22
it	es	0.53*	0.49*

Table 5: Spearman (ρ) and Pearson (r) correlation coefficients across language pairs. Source and Target denote the source language and the target language into which the source is translated, respectively. Asterisks (*) denote a statistically significant difference ($p < 0.05$).

- **Shape_type**: data Units configuration (chain, sibling, mixed) that controls structural complexity.
- **Size**: number of triples for Data Units.
- **Qid**: Wikidata Id.
- **Type**: Distinction between head and long-tail entities
- **Sub_type**: Distinction between subtypes (e.g. head entities with rare claims).
- **triplesets**: originaltripleaset (as collected) and modifiedtripleaset (normalized input).
- **lexicalizations**: one or more lex elements with attributes quality (gold/silver) and lang (en/es/it).

E Prompts

English version In English, structured data is commonly represented as triples, with the format [subject, predicate, object]. Based on these triples, generate a single-paragraph text composed of complete, grammatically correct, and natural sentences.

INSTRUCTIONS:

- Generate the text solely from the input triples
- Return the final verbalization with this format: The final verbalization is: [verbalization output]
- Insert the verbalization of the input triples within square brackets [], without adding anything else.

Spanish version En español, los datos estructurados se representan comúnmente como tríos, con el formato [sujeto, predicado, objeto]. Basándose en estos tríos, genere un texto de un solo párrafo compuesto por oraciones completas, gramaticalmente correctas y naturales.

ISTRUCCIONES:

- Genere el texto únicamente a partir de las tripletas de entrada.
- Devuelva la verbalización final con este formato: La verbalización final es: [salida de verbalización]
- entre corchetes [] inserte la verbalización de las tripletas de entrada, sin añadir nada más

Italian version In italiano, i dati strutturati sono comunemente rappresentati come triple, con il formato [soggetto, predicato, oggetto]. Sulla base di queste triple, genera un testo di un solo paragrafo composto da frasi complete, grammaticalmente corrette e naturali.

ISTRUZIONI:

- Genera il testo esclusivamente dalle triple in input
- Restituisci la verbalizzazione finale con questo formato: La verbalizzazione finale è: [output di verbalizzazione]
- tra le parentesi quadre [] inserisci la verbalizzazione delle triple in input, senza aggiungere altro.

F Results

Model	Type	BERT _r	BLEU ↑	chrF ↑	R-1 ↑	R-2 ↑	R-L ↑	PPL ↓
Qwen-3B	Long-tail	0.50	0.13	50.10	0.53	0.30	0.45	182.87
	Top-head	0.51	0.13	49.12	0.52	0.29	0.44	70.45*
Qwen-7B	Long-tail	0.56	0.19	54.27	0.58	0.35	0.49	138.70
	Top-head	0.58	0.18	53.41	0.57	0.35	0.48	87.96*
Gemma-4b	Long-tail	0.53	0.17	54.32	0.57	0.35	0.48	4724.94
	Top-head	0.56	0.17	53.48	0.56	0.34	0.46	3488.00*
Gemma-12b	Long-tail	0.62	0.24	60.38	0.64	0.42	0.55	455.57*
	Top-head	0.64	0.23	58.73	0.63	0.40	0.53	797.47
Llama-3B	Long-tail	0.49	0.13	49.75	0.51	0.30	0.43	102.40
	Top-head	0.50	0.14	48.19	0.49	0.29	0.42	72.55*
Llama-8B	Long-tail	0.55	0.18	53.81	0.58	0.36	0.50	134.53
	Top-head	0.57	0.19	53.34	0.57	0.36	0.48	74.13*

Table 6: Comparison of models by type. Asterisks (*) denote a statistically significant difference with $p < 0.05$.

Model	Lang	BERT _r	BLEU ↑	chrF ↑	R-1 ↑	R-2 ↑	R-L ↑	PPL ↓
Llama-3B	en	0.58*	0.21*	56.67*	0.63*	0.40*	0.53*	73.51*
	es	0.46	0.11	46.40	0.47	0.26	0.39	114.99
	it	0.44	0.08	44.26	0.42	0.22	0.35	81.98
Llama-8B	en	0.63*	0.25*	60.72*	0.66*	0.43*	0.56*	79.37*
	es	0.59	0.18	53.27	0.58	0.35	0.49	170.14
	it	0.45	0.13	46.85	0.48	0.28	0.42	80.23
Qwen-3B	en	0.50	0.17*	52.79*	0.60*	0.36*	0.51*	84.29
	es	0.54*	0.12	49.10	0.51	0.28	0.43	227.52
	it	0.48	0.09	47.19	0.46	0.24	0.39	100.23
Qwen-7B	en	0.67*	0.29*	62.71*	0.69*	0.46*	0.59*	119.43
	es	0.60	0.16	52.78	0.56	0.33	0.48	149.16
	it	0.45	0.11	46.25	0.46	0.26	0.40	85.39*
Gemma-4b	en	0.62*	0.23*	60.32*	0.66*	0.42*	0.55*	3798.99
	es	0.51	0.14	50.76	0.53	0.31	0.44	6068.41
	it	0.50	0.14	50.85	0.50	0.30	0.43	2794.93*
Gemma-12b	en	0.64*	0.27*	61.68*	0.68*	0.44*	0.57*	357.51
	es	0.63	0.21	58.27	0.62	0.39	0.52	739.90
	it	0.61	0.22	59.17	0.62	0.40	0.53	687.37

Table 7: Average scores for different languages for each model (zero-shot). Asterisks (*) denote a statistically significant difference with $p < 0.05$.

Model	Quality	BERT _r ↑	BLEU ↑	chrF ↑	R-1 ↑	R-2 ↑	R-L ↑	PPL ↓
Llama-3B	gold	0.46	0.12	48.04	0.48	0.27	0.40	95.30
	silver	0.52*	0.14	49.80	0.52*	0.30*	0.44*	86.86*
Llama-8B	gold	0.52	0.18	52.68	0.50	0.55*	0.48	90.94
	silver	0.58*	0.19	54.22	0.59*	0.36	0.50*	122.10*
Qwen-3B	gold	0.50	0.13	50.62	0.53	0.30	0.44	94.12
	silver	0.51	0.13	49.08	0.52	0.29	0.45	165.11
Qwen-7B	gold	0.53	0.17	52.59	0.55	0.33	0.47	117.53
	silver	0.60*	0.20	54.75	0.59*	0.36	0.50	118.35
Gemma-4b	gold	0.52	0.17	53.67	0.55	0.34	0.46	1143.13
	silver	0.55*	0.18	54.18	0.57	0.35	0.48	6197.47
Gemma-12b	gold	0.61	0.23	59.88	0.63	0.41	0.54	406.83
	silver	0.64*	0.24	59.59	0.64	0.41	0.54	715.74

Table 8: Average scores for different quality levels (silver, gold) for each model (zero-shot). Asterisks (*) denote a statistically significant difference with $p < 0.05$.

Type	Lang	BERT _r	BLEU ↑	chrF ↑	R-1 ↑	R-2 ↑	R-L ↑	PPL ↓
Llama-3B								
Long-tail	en	0.57	0.21	56.58	0.62	0.39	0.53	86.17
	es	0.45	0.11	47.07	0.48	0.26	0.40	144.99
	it	0.44	0.09	45.61	0.44	0.23	0.37	75.80
	total	0.49	0.13	49.75	0.51	0.30	0.43	102.32
WebNLG	en	0.56	0.23	57.17	0.71*	0.47*	0.56	53.67*
	es	0.49*	0.16*	49.20	0.58*	0.35*	0.45*	68.28*
	it	0.48*	0.13*	47.75*	0.53*	0.30*	0.42*	41.18*
	total	0.51	0.17*	51.37	0.60	0.37	0.48*	54.38*
Llama-8B								
Long-tail	en	0.61	0.25	60.39	0.66	0.43	0.56	93.36
	es	0.58*	0.17	53.30	0.58	0.35	0.49	222.16
	it	0.46*	0.13	47.73	0.49	0.29	0.43	88.06
	total	0.55	0.18	53.81	0.58	0.36	0.50	134.53
WebNLG	en	0.61	0.30	65.53	0.74*	0.51*	0.59*	42.56*
	es	0.55	0.19*	52.43	0.66*	0.43*	0.53*	90.42*
	it	0.43	0.15*	45.35	0.54*	0.34*	0.46*	71.73*
	total	0.53	0.21*	53.77	0.65*	0.43*	0.53*	68.23*
Gemma-4B								
Long-tail	en	0.60	0.23	59.85	0.66	0.42	0.55	4886.76
	es	0.50	0.14	51.49	0.54	0.32	0.45	8074.87
	it	0.49	0.15	51.62	0.51	0.31	0.44	1213.20
	total	0.53	0.17	54.32	0.57	0.35	0.48	4724.94
WebNLG	en	0.61	0.28*	62.43*	0.74*	0.50*	0.59*	930.34*
	es	0.52	0.20*	53.79*	0.61*	0.39*	0.48*	1493.20*
	it	0.49	0.18*	53.72	0.57*	0.36*	0.46	601.97*
	total	0.54	0.22*	56.65*	0.64*	0.42*	0.51	1008.50*
Gemma-12B								
Long-tail	en	0.63	0.26	61.65	0.67	0.44	0.58	330.02*
	es	0.62	0.21	58.80	0.62	0.40	0.53	611.42
	it	0.61	0.24	60.68	0.64	0.42	0.55	425.27
	total	0.62	0.24	60.38	0.64	0.42	0.55	455.57
WebNLG	en	0.66*	0.33*	65.97	0.77*	0.53*	0.62*	333.93
	es	0.65	0.30*	62.68	0.71*	0.49*	0.57*	210.04*
	it	0.64*	0.28*	62.24*	0.70*	0.47*	0.57	326.33*
	total	0.65*	0.30*	63.63	0.73*	0.50*	0.59*	290.10*
Qwen-3B								
Long-tail	en	0.49	0.17	52.66	0.59	0.35	0.51	100.79
	es	0.53	0.12	49.47	0.52	0.29	0.44	325.03
	it	0.48	0.10	48.08	0.48	0.25	0.41	124.49
	total	0.50	0.13	50.07	0.53	0.30	0.45	183.44
WebNLG	en	0.55*	0.24*	57.05	0.72*	0.47*	0.57*	51.40*
	es	0.56*	0.17*	51.54	0.63*	0.38*	0.49*	150.84*
	it	0.55*	0.15*	51.75	0.59*	0.34*	0.47*	55.81*
	total	0.55*	0.19*	53.44	0.64*	0.40*	0.51*	86.02*
Qwen-7B								
Long-tail	en	0.65*	0.29	62.47	0.69	0.46	0.58	134.13
	es	0.60	0.16	53.43	0.57	0.34	0.49	181.75
	it	0.44	0.11	46.82	0.47	0.26	0.41	100.26
	total	0.56	0.19	54.24	0.58	0.35	0.49	138.71
WebNLG	en	0.68	0.34*	66.33	0.79*	0.55*	0.64*	46.52*
	es	0.60	0.24	56.19	0.66*	0.43*	0.53*	117.41*
	it	0.44	0.15	46.96	0.54*	0.33*	0.45*	122.70*
	total	0.57	0.24	56.49	0.66*	0.44*	0.54*	95.54*

Table 9: Performance comparison of models on long-tail triples from TailNLG and WebNLG test set. Asterisks (*) denote a statistically significant difference with $p < 0.05$.