

KG-TRICK ✂: Unifying Textual and Relational Information Completion of Knowledge for Multilingual Knowledge Graphs

Anonymous ACL submission

Abstract

Multilingual knowledge graphs (KGs) provide high-quality relational and textual information for various NLP applications but they are often incomplete, especially in non-English languages. Previous research has shown that combining information from several knowledge graphs in different languages aids both Knowledge Graph Completion (KGC), the task of predicting of missing relations between entities, and Knowledge Graph Enhancement (KGE), the task of predicting missing textual information for entities. While previous efforts have considered KGC and KGE as independent tasks, we hypothesize that they are interdependent and mutually beneficial. To this end, we introduce KG-TRICK, a novel sequence-to-sequence framework that unifies the tasks of textual and relational information completion for multilingual knowledge graphs. KG-TRICK demonstrates that i) it is possible to unify the tasks of KGC and KGE into one single framework, and ii) combining textual information from multiple languages is beneficial to improve the completeness of a KG. As part of our contributions, we also introduce WikiKGE-10++, the largest manually-curated benchmark for textual information completion of KGs, which features over 30,000 instances across 10 diverse languages.

1 Introduction

Knowledge graphs (KGs) aim to encode structured information about the world in a machine-readable format (Hogan et al., 2021), providing high-quality relational and textual information for various NLP applications, such as question answering (Mckenna and Sen, 2023), information retrieval (Reinanda et al., 2020), entity linking (Hu et al., 2023), and machine translation (Modrzejewski et al., 2020), among others. While large language models (LLMs) are increasingly retrieving information from KGs to improve their factuality

and performance on many NLP tasks (Wang et al., 2023), their effectiveness in multilingual applications is limited due to the important gap between the completeness of English and non-English information in multilingual KGs (Peng et al., 2023). Indeed, KGs are not complete: a non-negligible quantity of information about entities (e.g., entity names, aliases, and descriptions) and relations (e.g., the connections between entities) is missing in non-English languages (Conia et al., 2023a). Therefore, improving the completeness of KGs has attracted significant attention over the years.

To address this issue, the research community has worked on two main tasks: Knowledge Graph Completion (KGC) and Knowledge Graph Enhancement (KGE). KGC is the task of predicting missing relations between entities already defined in a KG (Bordes et al., 2013), while KGE is the task of predicting missing textual information for entities in a KG (Conia et al., 2023b). More formally, KGC is defined as follows: given a KG \mathcal{G} , the task of KGC is to predict the missing tail entity t given the head entity h and the relation r in a triplet $(h, r, ?)$. For example, given the triplet (Joe Biden, occupation, ?), a possible answer could be *politician* or, more specifically, the ID of the politician entity in the KG. On the other hand, KGE is defined as follows: given an entity e in a KG \mathcal{G} , the task of KGE is to predict missing textual information (e.g., an entity name, alias, or description) for e in a target language. For example, given the entity *Joe Biden* in English, a possible alias would be *Joseph R. Biden Jr.* or *Joseph Robinette Biden Jr.* in English, or the primary name in Chinese would be 乔·拜登. In this simple example, we can already demonstrate the interdependence between KGC and KGE: the same head and tail entities can have different names in different languages, but the relation between them should hold across languages. However, their interdependence becomes challenging when dealing with ambigu-

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082

ous entities (e.g., *Paris* the city and *Paris* the prince of Troy) and entities whose names are not directly translatable (e.g., *The Matrix* in English and 黑客帝国 (*Hacker’s Empire*) in Chinese).

While KGC and KGE have previously been considered as independent tasks, in this work, we investigate their interdependence and hypothesize that they are mutually beneficial. Our hypothesis is based on two symmetric observations. First, solving KGC provides rich language-independent relational information about entities, which may aid KGE to generate higher quality textual information across languages. Second, solving KGE provides rich language-dependent textual information about entities, which may aid KGC to align entities with names and descriptions across languages more effectively. To this end, we introduce KG-TRICK (Textual and Relational Information Completion of Knowledge), a novel unified framework that combines the tasks of KGC and KGE into one single task. Different from previous approaches, the KG-TRICK framework is multilingual by design and is able to leverage the complementary textual information from multiple languages to improve the completeness of a multilingual KG. Not only does KG-TRICK remove the need for separate KGC and KGE models, but it also outperforms similarly-sized state-of-the-art approaches tailored for each individual task, while achieving competitive performance compared to much larger language models. To evaluate the robustness of KG-TRICK and encourage future systems on textual information completion of KGs, we also introduce WikiKGE-10++, the largest manually-curated benchmark for textual information completion for multilingual KG in 10 languages.

We can summarize our contributions as follows:

- We unify the tasks of KGC and KGE to encompass not only the task of predicting missing links in a KG but also the task of completing its multilingual textual information;
- We introduce WikiKGE-10++, the largest manually-curated benchmark for textual information completion of KGs, including over 30,000 entities over 10 languages, to accompany KGC benchmarks and create a comprehensive evaluation suite;
- We present KG-TRICK, a novel sequence-to-sequence model that is able to combine information from multiple languages in an effective

way to tackle textual and relation completion of knowledge graphs in a joint fashion;

- We show that KG-TRICK outperforms similarly-sized state-of-the-art models tailored for each task, while achieving competitive performance compared to larger LMs.

We believe that our work – our task reformulation, manual benchmark, and unified method – is a significant step forward to improve the quality of multilingual KGs and broaden their applicability to multilingual downstream tasks. To encourage future work in this direction, we release our software and benchmark at <https://anonymized>.

2 Related Work

In this section, we briefly review the literature on Knowledge Graph Completion (KGC) and Knowledge Graph Enhancement (KGE) and discuss the challenges of completing textual and relational information in multilingual knowledge graphs.

Multilingual Knowledge Graphs. As mentioned above, KGs aim to encode information about our world knowledge in a structured, machine-readable format (Hogan et al., 2021). Such information also includes lexicalizations like entity names, aliases, and descriptions; when these are available in multiple languages, the KG is called a multilingual KG. There are different ways to construct and organize multilingual KGs. For example, in DBpedia (Lehmann et al., 2015), an entity is language-dependent and is represented in different languages using different entity IDs, whereas in Wikidata (Vrandečić and Krötzsch, 2014), an entity is language-specific and is represented by the same entity ID to which different language-specific labels are attached. The construction of multilingual KGs is an active area of research, and there are several challenges to be addressed, such as the alignment of entities across languages (Chakrabarti et al., 2022), the completion of missing relational information (Chen et al., 2020b), and the addition of textual information, especially in non-English languages (Conia et al., 2023b).

Knowledge Graph Completion. The task of KGC is to predict missing relations between entities already defined in a KG (Bordes et al., 2013). This task has been studied extensively in the literature, and there are categories of methods to solve it, including embedding-based methods (Lin et al.,

2015b), path-based methods (Lin et al., 2015a), and rule-based methods (Chen et al., 2020a). More recently, sequence-to-sequence models have been proposed to solve KGC, by treating it as a text-to-text generation task where the input is a partial triplet and the output is the missing entity (Saxena et al., 2022). However, these approaches have been designed for monolingual KGs, as multilinguality adds a layer of complexity to the task. Chakrabarti et al. (2022) have taken a step in this direction by including an auxiliary task to translate entity names, but they neither consider completing triples across languages nor the completion of more complex textual information, such as entity descriptions.

Knowledge Graph Enhancement. The task of KGE is to predict missing textual information for entities in a KG. This task is more recent in the literature, but there are several approaches to tackle it, such as machine translation, web search, and language model-based methods (Conia et al., 2023a). However, Conia et al. (2023a) have mainly focused on i) combining answers from multiple KGE systems to improve coverage and precision, and ii) evaluating the quality of the textual information generated by KGE systems for popular entities only, while iii) ignoring the connection between KGE and KGC, especially in the multilingual setting.

3 Unifying Textual and Relational Information Completion

In this section, we introduce KG-TRICK, or how we unify the tasks of KGC and KGE into one single framework, and how we leverage the complementary textual information from multiple languages to improve the completeness of a multilingual KG.

3.1 Task Reformulation

Given the similarities between the two tasks of KGC and KGE and the interdependence between them (see Section 1, in which we provide a high-level intuition), we reformulate both tasks as a single multilingual text-to-text generation task as shown in Figure 1. KG-TRICK consists of three main components, namely the verbalization, the fine-tuned multilingual sequence-to-sequence model and the ensemble module to obtain the predicted entities for KGC task. This unified framework allows us to i) treat KGC and KGE as a single task, and ii) to leverage the complementary textual information from multiple languages to better complete factual information and reversely to improve

the latent, dense representation of the fine-tuned sequence-to-sequence model leading to improved KGE performance. Figure 1 illustrates the pipeline of KG-TRICK for both KGC and KGE. Thanks to our reformulation, KG-TRICK sees both tasks as the task of predicting the tail entity t given the head entity h and the relation r in a triplet $(h, r, ?)$, as we will detail in the following sections.

3.1.1 KGC as Text-to-Text Generation

In KGC, the task is to predict the missing tail entity t given the head entity h and the relation r in a triplet $(h, r, ?)$. We first reformulate this task as a text-to-text generation task, where the input is a partial triplet composed of the primary name and short description of the head entity h and the relation r . The model is then asked to generate the missing tail, or, more precisely, the primary name and short description of the tail entity t .

One important drawback of this reformulation is that it does not take into account the input and output languages, which is crucial for multilingual KGs. We overcome this limitation by extending the triplet to a tuple of five elements, which include the source and target languages, as shown in Figure 1. More specifically, the input to the model is now a tuple $(l_s, l_t, h, r, ?)$, where l_s is the source language of the input, l_t is the target language of the output, $h = \text{primary name} + \text{short description}$, and r is the relation. The model then predicts t , i.e., the primary name and short description of the tail entity in l_t . For example, given the input $(\text{en}, \text{es}, h, r, ?)$, the model generates the primary name and short description of the entity *político | persona involucrada en la política*; while given the input $(\text{en}, \text{zh}, h, r)$, the model generates the primary name and short description of the entity 政治家 | 从事政治活动的人. This reformulation significantly increases the training data pairs by extending cross lingual name based entity alignment to cross lingual relation based entity alignment resulting in higher quality of entity alignment.

3.1.2 KGE as Text-to-Text Generation

In KGE, the task is to predict missing textual information for entities in a KG. This task can also be reformulated as a text-to-text generation task, similar to KGC. We can immediately see that the formulation outlined above for KGC can be directly applied to KGE, with the only difference being that the head entity h may be represented only by its primary name in case we want to generate a short

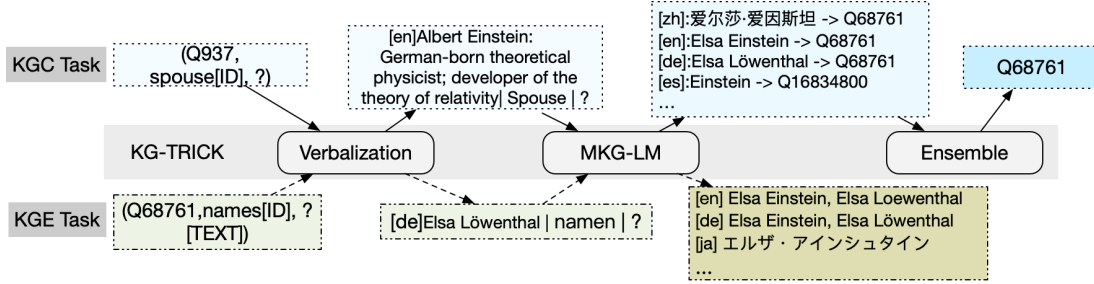


Figure 1: KG-TRICK: a unified seq-to-seq framework for KGC and KGE. KGC dataflow is in blue; KGE dataflow is in green. For KGE, an input triplet (Q68761, names, ?) is verbalized as “[de] Elsa Löwenthal|namen|?” and then passed to the model, which can generate the names in multiple languages. For KGC, an input triplet (Q937, spouse, ?) is verbalized as “[en] Albert Einstein | spouse | ?” and then passed to the model, which can generate the name Elsa Einstein in multiple languages. The ensemble module consolidates all the outputs into the best one.

description for h itself. Moreover, we also allow the head entity h and the tail entity t to be the same entity, which allows us to generate aliases for an entity in a specific language. For example, given the partial triplet *Joe Biden: President of the US | has name |*, the model predicts the primary name of the entity *Joe Biden* in the target language but it can also generate one or more aliases, such as *Joseph R. Biden Jr.* or *Joseph Robinette Biden Jr.* in English, or 乔·拜登 or 乔·罗宾内特·拜登 in Chinese. Interestingly, when this reformulation is used in its most simple form, i.e., by using only the primary name of the head entity, it becomes equivalent to translation into the target language. This is a powerful feature, as it allows us to generate missing textual information in any language in KG.

3.2 The KG-TRICK Model

Unifying KGE and KGC, we implement KG-TRICK as a general sequence-to-sequence model, which learns to generate both missing relational and textual missing information in a KG. More formally, given a tuple $(l_s, l_t, h, r, ?)$, the model is asked to generate t from the source language l_s to the target language l_t conditioned on h and r as following:

$$o = \text{KG-TRICK}(l_s, l_t, h, r, ?) \quad (1)$$

where o is the output generated by the model. In other words, o is the primary name and short description of the tail entity in the target language, and KG-TRICK learns to estimate the probability of generating o given the input $(l_s, l_t, h, r, ?)$.

KG-TRICK can be implemented using any sequence-to-sequence architecture, such as a transformer (Vaswani et al., 2017) or a recurrent neural network (Rumelhart et al., 1985). In practice, we conducted our experiments with one main

architecture, i.e., multilingual BART, which is a transformer-based encoder-decoder model, and we found that it performs well on the task, as shown in Section 5. The model can be trained using a standard maximum likelihood estimation (MLE) objective, and it can be evaluated using standard metrics for text generation, such as BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and COMET (Rei et al., 2020). In this work, we study three main variants of KG-TRICK, which differ in the training data they use: i) $\text{TRICK}_{\text{KGC}}$, which uses only the relational information of the KG, ii) $\text{TRICK}_{\text{KGE}}$, which uses only the textual information of the KG, and iii) $\text{TRICK}_{\text{KGC}+\text{KGE}}$, which uses both the relational and textual information of the KG.

3.3 Inference for KGC and KGE

Once the output is generated, it can be used to complete the KG in two ways. The application to KGE is straightforward, as the output is the missing textual information for an entity in the target language. The application to KGC is slightly more complex, as the output is the textual representation (i.e., the primary name and short description) of the missing tail entity in the target language. Relying only on an exact match between primary names may not be sufficient to determine the correct entity, especially in the case of ambiguous entities, such as *Paris* the city and *Paris* the prince of Troy. While previous work (Saxena et al., 2022) enumerates the entities with the same name (e.g., Paris_1 , Paris_2 , etc.), we incorporate entity descriptions as additional information from KG-TRICK to help disambiguate entities with the same primary name resulting in higher entity linking accuracy.

Ensemble across languages. Since KG-TRICK can generate text in any target language for which it

has been (pre-)trained, we leverage this capability to complete the missing information from multiple languages. When corresponding entity ID are linked from generated text by different languages, we then ensemble and choose the most common predictions as is showed in Figure 1.

4 WikiKGE-10++

In this section, we introduce WikiKGE-10++, the largest manually-curated benchmark for textual information completion of KGs, which features over 30,000 instances across 10 diverse languages. Having realized the importance of evaluating systems on textual information completion of multilingual KGs, in 2023, Conia et al. created WikiKGE-10, a benchmark for evaluating KGE systems on the completion of entity names and aliases in 10 languages. However, WikiKGE-10 is limited in two dimensions: i) it only allows for the evaluation of entity names and aliases, and ii) the entities included in the benchmark are popular entities only, i.e., they belong to the top-10% most popular entities in Wikidata according to number of page views of their corresponding Wikipedia pages.

A core contribution of our work is the creation of WikiKGE-10++, which extends WikiKGE-10 in two above-mentioned dimensions, hence the two “+” signs in the name of our benchmark. First, WikiKGE-10++ includes not only entity names and aliases, but also entity descriptions, which are crucial for many downstream tasks to create better entity representations (Ri et al., 2022). Second, WikiKGE-10++ includes a much larger set of entities, which belong to the torso of the popularity distribution of Wikidata (i.e., between the top-10% and top-50% most popular entities) and also the tail of the popularity distribution (i.e., below the top-50% most popular entities). This is important because the majority of entities in a KG are not popular, and different conclusions can be drawn when evaluating systems on different popularity tiers. Overall, WikiKGE-10++ is around 3 times larger than WikiKGE-10 in terms of the number of entities while also including entity descriptions.

Including torso and tail entities. The inclusion of torso and tail entities in WikiKGE-10++ is important to assess the robustness of KGE systems on the completion of textual information for entities for which the amount of information available inside (and also outside) the KG is limited. While most of the search queries and user interactions

may be focused on popular entities, the majority of entities in a KG are not popular, and they are often the most challenging. Therefore, we asked a team of annotators to manually curate 1000 torso entities and 1000 tail entities per language.

Including entity descriptions. The inclusion of entity descriptions in WikiKGE-10++ is important for evaluating KGE systems on the completion of textual information that is usually longer and more complex than entity names and aliases. However, evaluating KGE systems on the entity descriptions already available in Wikidata is not ideal, as those are not always manually curated and often under-specific, e.g., the description of many cities is simply *city in country*. Therefore, we asked a team of annotators to produce high-quality descriptions for 1000 entities per popularity tier per language.

5 Experiments and Results

5.1 Datasets and Benchmarks

KGC. We carry out our KGC experiments on Wikidata5M (Wang et al., 2021a) transductive split. We stress that, although entities in Wikidata are language-agnostic, the original Wikidata5M dataset is English-only, i.e., the textual information (names and descriptions) for the 5 million entities is only in English. Therefore, this English-only setting may not be ideal to evaluate our multilingual KGC system; however, as is illustrated in Table 3, our reformulation could further enhance TRICK_{KGE}’s performance with multilinguality.

KGE. We carry out our KGE experiments on our newly annotated WikiKGE-10++ dataset, which features 10 languages and 30,000 entities as introduced in Section 4. We use this dataset to evaluate KGE models on the completion of entity names, aliases, and descriptions in 10 languages.

5.2 Comparison Systems

KGC. Baseline approaches for KGC can be divided into two categories: *embedding-based* and *text-based*. Embedding-based methods derive an embedding for each entity and relation from the graph structure of the KG, and rank the most probable tail entity via a vector similarity function, e.g., L2 distance (Bordes et al., 2013), complex space distance (Trouillon et al., 2016b) or other distance measures. Text-based methods use encoder-only language models (Wang et al., 2022, SimKGC) to

449 encode both the head entity and the relation using
 450 their textual information, or encoder-decoder
 451 language models (Saxena et al., 2022, KG-T5)
 452 to generate the missing tail entity. Concurrent
 453 work (Kochsiek et al., 2023, KGT5-context) com-
 454 bines subgraph-structure information and sequence-
 455 to-sequence models. Since we build upon text-
 456 based methods, we compare our KG-TRICK mod-
 457 els with SimKGC and KG-T5, which are the most
 458 relevant baselines for a fair comparison.

459 **KGE.** While KGE is a relatively recent task, pre-
 460 vious work has already indicated several strong
 461 baselines, including i) using NLLB-200¹ (Costa-
 462 jussà et al., 2022) to translate entity names and
 463 descriptions from a language , and ii) prompting
 464 LLMs (e.g. GPT-3.5 or Llama), to generate textual
 465 information for an entity.

466 5.3 Experimental Setup

467 We use the entities within Wikidata5M which con-
 468 tains a set of around 5 million entities and a collec-
 469 tion of around 20 million triplets and collect their
 470 available textual information (entity names, aliases,
 471 and descriptions) for English and 9 other languages
 472 from a Wikidata dump² to create a silver training
 473 set \mathcal{E} for KGE in multiple EN \rightarrow XX directions con-
 474 taining around 16 million records. For KGC, the
 475 original Wikidata5M triplets are expanded with
 476 our downloaded Wikidata dump to form over 150
 477 million triplets \mathcal{T} multilingually (EN \rightarrow XX) in 9
 478 languages pairs. We train three variants of KG-
 479 TRICK: one on \mathcal{E} (denoted as TRICK_{KGE}), one on
 480 \mathcal{T} (denoted as TRICK_{KGC}), and one on the mix-
 481 ture of both (denoted as TRICK_{KGC+KGE}).³ To
 482 verify the multilinguality of TRICK_{KGE}, we also
 483 include a bilingual version train with single lan-
 484 guage pair (e.g. EN \rightarrow IT) on KGE task, denoted as
 485 TRICK_{KGE(bilingual)} in Table 3. While the num-
 486 ber of training samples are disproportional for KGC
 487 and KGE, to trade off the performance between
 488 the two tasks, as is shown in Table 5, we find a
 489 sweet spot of combining 50% of KGC data into
 490 the joint training. We denote this derivative as
 491 TRICK_{50%KGC+KGE} in Table 2 and Table 3. We
 492 provide more details balancing the training data of
 493 the two tasks in Appendix C.

494 **Evaluation.** For KGC, we evaluate the systems
 495 using standard ranking-based metrics, namely,

¹In this work we use NLLB-200-Distilled (600M).

²Downloaded in November 2023.

³All TRICK models are fine-tuned from mBART-large-50.

Model	MRR	hit@1	hit@3	hit@10
TransE (Bordes et al., 2013)	25.3	17.0	31.1	39.2
DisMult (Yang et al., 2014)	25.3	20.8	27.8	33.4
Simple (Kazemi and Poole, 2018)	29.6	25.2	31.7	37.7
RotatE (Sun et al., 2019)	29.0	23.4	32.2	39.0
QuatE (Zhang et al., 2019)	27.6	22.7	30.1	35.9
ComplEx (Trouillon et al., 2016a)	30.8	25.5	-	39.8
DKRL (Xie et al., 2016)	16.0	12.0	18.1	22.9
KEPLER (Wang et al., 2021b)	21.0	17.3	22.4	27.7
MLMLM (Clouatre et al., 2021)	22.3	20.1	23.2	26.4
SimKGC + Desc.	35.8	31.3	37.6	44.1
KG-T5	30.0	26.7	31.8	36.5
KG-T5 + Desc.	38.1	35.7	39.7	42.2
KG-T5 + Desc.*	37.0	34.7	38.4	41.1
TRICK _{KGC}	38.2	36.0	39.7	41.8
TRICK _{KGC+KGE}	38.8	36.6	40.4	42.6

Table 1: KGC results on the test set of Wikidata5M. TRICK achieves strong performance over competitive baselines. *: retrained in the same setting as TRICK_{KGC}.

496 *hit@1*, *hit@3*, *hit@10*, and *Mean Reciprocal Rank*
 497 (MRR). Hit@k measures the proportion of correct
 498 answers in the top-k predictions, while MRR mea-
 499 sures the average rank of the correct answer. For
 500 KGE, we follow the evaluation protocol proposed
 501 by Conia et al. (2023a), which includes two main
 502 metrics: coverage and precision. *Coverage* evalu-
 503 ates the number of entities for which a system is
 504 able to produce at least one correct entity name,
 505 while *Precision* evaluates the ability of a system to
 506 identify incorrect entity names and aliases. Finally,
 507 we report the COMET scores for the completion
 508 of entity descriptions, a standard metric for text
 509 generation and machine translation.

510 5.4 Results on KGC

511 Table 1 shows the results obtained by our KG-
 512 TRICK models compared to other KGC-only mod-
 513 els on the test set of Wikidata5M. In general, we
 514 can observe that KG-TRICK generally outperforms
 515 all the other strong baselines on MRR, hit@1,
 516 and hit@3, and it is the second best model for
 517 hit@10. More specifically, TRICK_{KGC} (trained
 518 only on KGC data but in multiple languages) al-
 519 ready outperforms both SimKGC and KG-T5, on
 520 almost all the metrics. Notably, this first result
 521 demonstrates that our model is able to outperform
 522 strong baselines that are tailored for KGC on a
 523 dataset Wikidata5M that is designed for KGC (and
 524 that is biased in its creation towards entities that
 525 have English lexicalizations). Moreover, we can ob-
 526 serve that TRICK_{KGC+KGE} achieves scores that are
 527 even higher than TRICK_{KGC}, which demonstrates
 528 that unifying KGC and KGE leads to additional

	Coverage				Precision				COMET		
	Head	Torso	Tail	Avg.	Head	Torso	Tail	Avg.	Head	Torso	Tail
NLLB-200 _{EN→XX}	29.1	26.1	24.3	26.5	47.6	39.6	34.9	40.7	0.64	0.63	0.63
Llama3-8B	27.2	22.9	20.4	23.5	46.0	36.6	31.2	37.9	0.62	0.62	0.62
GPT-3.5	35.0	29.6	26.7	30.4	51.9	42.7	36.8	43.8	0.66	<u>0.64</u>	<u>0.64</u>
TRICK _{KGE}	31.5	31.4	29.9	30.9	56.4	51.4	46.5	51.4	0.63	<u>0.64</u>	<u>0.64</u>
TRICK _{50%KGC+KGE}	33.1	31.2	30.1	31.5	57.7	51.7	47.1	52.1	0.61	0.62	0.62
TRICK _{KGC+KGE}	31.6	29.5	28.2	29.7	57.8	52.2	46.5	52.2	0.60	0.60	0.61

Table 2: KGE results on WikiKGE-10++ split by head, torso and tail entities. Best results in bold.

improvements on the KGC task. This second result empirically shows that the two tasks are indeed interdependent and mutually beneficial, and that the combination of the two tasks can lead to better results than the two tasks individually.

On the other hand, we can observe that the performance of TRICK_{KGC} and TRICK_{KGC+KGE} is not as good as the performance of SimKGC on hit@10. We hypothesize that this is likely due to the fact that the sampling capacity of sequence-to-sequence models is limited and constrains their performance with higher values of k . Indeed, for SimKGC, TransE and ComplEx, due to their closed-world assumption, the candidates for the tail entities are known during inference for similarity search. However, for text generation models, such as KG-T5 and KG-TRICK, which operate under a more challenging open-world assumption, the diversity of generated candidates may be limited by the sampling strategy used for decoding. Future work could focus on improving the sampling capacity of generative models.

5.5 Results on KGE

Table 2 shows the KGE results on our new WikiKGE-10++ benchmark split by entity popularity, while Table 3 shows the results by language.

Coverage. Overall, we could observe that TRICK_{50%KGC+KGE} outperforms strong baselines on average and across most languages. Interestingly in Table 2, TRICK_{KGE} and TRICK_{KGC+KGE} perform worse than GPT-3.5 on Coverage of head entities. This is likely due to the fact that GPT-3.5 has seen substantially more popular entity names during its pre-training and is equipped with considerably more parameters to store such information. However, TRICK series quickly catch up with GPT-3.5 on Coverage of torso entities, and significantly outperform GPT-3.5 on Coverage of tail entities. It shows that GPT-3.5 quickly loses its ad-

vantage when the entities are less popular, and that KG-TRICK models feature a more balanced and consistent performance across different popularity tiers.

Precision. Overall, we can observe that TRICK significantly outperforms strong baselines on precision on head, torso, and tail entities, i.e., it is a more reliable system in identifying incorrect entity names and aliases in a given target language in a multilingual knowledge graph. In fact, TRICK is particularly effective on torso and tail entities, where it improves over NLLB-200 and GPT-3.5 by around 10% points in F1 score. This is important as completing missing knowledge is not only about providing the correct information but also about avoiding incorrect information.

Entity descriptions. Finally, we also report the COMET score for the completion of entity descriptions, borrowing a metric for open-ended text generation from MT. In this task, we can observe that TRICK_{KGE} is comparable with NLLB-200 and GPT-3.5, while TRICK_{KGC+KGE} is slightly worse on average than TRICK_{KGE}. These results open the door to future work: indeed, very different methods achieve comparable results on entity description completion, meaning that there is still a wide room for improvement in this task or COMET is not a good metric for comparing descriptions. We note that BLEU is not appropriate either, as its score is not defined for short texts, e.g., one word.

Multilingual and Multi-task As is illustrated in Table 3, TRICK_{KGE} outperforms TRICK_{KGE}(bilingual) on almost all languages for both Precision and Coverage, indicating that jointly train a unified model for all languages could inherently benefit its multilinguality. On the multi-task side, combining KGC and KGE tasks requires caution, as is demonstrated by TRICK_{KGC+KGE} and TRICK_{50%KGC+KGE}. KGC and KGE are mutually

		Coverage F1	#Params	AR	DE	ES	FR	IT	JA	KO	ZH	Avg	
names & aliases	NLLB-200 _{EN→XX}	0.6B	16.9	39.8	37.9	40.8	40.5	9.4	18.6	8.1	26.5		
	Llama3-8B	8B	11.3	35.9	32.8	35.0	32.6	12.4	15.5	12.5	23.5		
	GPT-3.5	175B	20.1	40.8	39.1	41.3	40.9	19.6	21.5	20.0	30.4		
	TRICK _{KGE} (bilingual)			<u>24.4</u>	39.5	33.7	37.2	34.0	18.4	23.3	15.3	27.5	
	TRICK _{KGE}	0.6B		<u>24.4</u>	39.5	38.8	40.8	40	20.4	26.3	17.0	30.9	
	TRICK _{50%KGC+KGE}			23.0	40.9	40.0	41.9	41.1	20.5	26.0	18.2	31.5	
	TRICK _{KGC+KGE}			22.6	39.2	37.2	39.4	39.4	19.8	24.3	16.0	29.7	
	Precision F1	#Params											
	NLLB-200 _{EN→XX}	0.6B	30.3	49.1	48.9	51.2	52.0	28.3	30.6	34.9	40.7		
	Llama3-8B	8B	24.1	46.0	45.0	47.7	45.4	30.4	26.6	38.3	37.9		
	GPT-3.5	175B	33.2	49.5	49.7	51.9	52.3	36.9	33.1	43.8	43.8		
	TRICK _{KGE} (bilingual)			46.3	52.3	53.1	54.5	55.3	43.8	46.3	45.5	49.6	
	TRICK _{KGE}	0.6B		48.4	54.3	55.2	57.0	56.2	45.6	47.5	47.5	51.4	
	TRICK _{50%KGC+KGE}			48.1	55.9	55.6	56.4	56.9	45.7	50.0	48.5	52.1	
	TRICK _{KGC+KGE}			49.4	54.8	55.5	56.7	56.6	47	49.0	48.3	52.2	
			COMET Score	#Params	AR	DE	ES	FR	IT	JA	KO	ZH	Avg
descriptions	NLLB-200 _{EN→XX}	0.6B	0.59	0.62	<u>0.66</u>	<u>0.63</u>	0.65	0.63	0.65	0.64	0.63		
	Llama3-8B	8B	0.56	0.62	0.65	<u>0.63</u>	0.65	0.60	0.60	0.61	0.62		
	GPT-3.5	175B	0.60	<u>0.63</u>	<u>0.66</u>	<u>0.63</u>	0.66	0.66	0.67	0.67	0.65		
	TRICK _{KGE}			0.58	<u>0.63</u>	<u>0.66</u>	<u>0.63</u>	0.65	0.67	0.64	0.64	0.64	
	TRICK _{50%KGC+KGE}	0.6B		0.56	0.61	0.63	0.60	0.64	0.64	0.62	0.60	0.61	
	TRICK _{KGC+KGE}			0.55	0.59	0.63	0.58	0.62	0.63	0.62	0.60	0.60	

Table 3: KGE experiments on WikiKGE-10++. TRICK achieves best performance on Precision and Coverage F1 scores. Best results in bold.

beneficial when training data is properly balanced, otherwise one task would dominate the distribution and cause regression on the other. More analysis on data balancing is discussed in Appendix C. We could also observe that the multilingual capability of Llama3-8B is far from ideal in KGE task.

5.6 Downstream Application: Results on Multilingual Question Answering

In addition to the KGC and KGE tasks, we also evaluate our KG-TRICK on a downstream application: our intuition is that (post-)pretraining on KGC and KGE tasks can allow a model to store more factoid knowledge, which can be useful for multilingual question answering. Therefore, we evaluate the performance of our TRICK_{KGC+KGE} when fine-tuned on answering the questions in the Mintaka dataset (Sen et al., 2022), which is a multilingual question answering dataset that contains knowledge-seeking questions, and compare its results in the same setting with directly fine-tune on Mintaka using mBART-large-50 that has not been (post-)pretrained on KGC and KGE tasks. Our experiments show that our model outperforms

mBART-large-50 by 3.1% (29.2% vs. 26.1%) on average in terms of EM (Exact Match), which demonstrates that the knowledge embedded in KGC and KGE training data could be easily transferred to the QA task in cross-lingual settings.

6 Conclusion

The contributions of this paper are threefold. First, we propose a novel multilingual KGC and KGE system, TRICK, which is able to complete relational and textual information in and across 10 languages. Second, we introduce a new human-curated dataset, WikiKGE-10++, which contains 10 languages and 30,000 entities for KGE evaluation. Third, we demonstrate that our TRICK system outperforms strong baselines on both KGC and KGE tasks, and that the combination of the two tasks can lead to better results than the two tasks individually. We also show that the knowledge embedded in KGC and KGE training data could be easily transferred cross-lingual QA. We hope our work and our WikiKGE-10++ can inspire future research on multilingual KGC and KGE.

652 Limitations

653 **Closed world assumption of KGC.** In this paper,
654 we assume that the entities within KGC tasks exist
655 in the KG. If the model predicts some entities that
656 do not exist, we simply ignore the inference. This
657 assumption limits the model’s capability to explore
658 the encoded knowledge within pre-trained multilin-
659 gual LMs. Although our model can be extended
660 to predict and extract entities outside of KG, our
661 experiments in section Section 5 demonstrate that
662 there is still a big headroom to further complete
663 the relation information when KG are extended
664 into multilingual setting since we can combine the
665 knowledge across languages that are not consid-
666 ered in a monolingual setting. We leave the deeper
667 exploration into future work.

668 **Limited exploration of pre-trained multilin-**
669 **gual LMs.** Our KGE task pays great attention to
670 enrich the entity names and descriptions of entities
671 with limited attention to other textual information
672 such as mottos , quotes. Given the pre-trained LMs
673 have been trained on huge amount of data, there
674 are great potential to extract out the information
675 that has been seen by the pre-trained LMs which
676 does not exist in KG. Even though KG-Trick can
677 be extended to infer other entity facts, our analysis
678 shows that entity names and descriptions are most
679 essential to enrich and disambiguate entities and
680 have led to great improvement of KGC task. Espe-
681 cially that description is a summarized free form
682 text that is highly representative of a particular en-
683 tity.

684 **Support unified multilingual KGs.** We focus
685 on the multilingual KGs that has entities repre-
686 sented by entity IDs and language dependent tex-
687 tual information are structured as the attributes as-
688 sociated with corresponding entities. Thus, differ-
689 ent from other research work that needs to align en-
690 tities and relationships together, our systems elimi-
691 nate such requirement. The benefit of this setting
692 is that we enriched the KG that are unified at the
693 very beginning and derive the knowledge from the
694 KG itself by extracting and inferring information
695 that can be derived by the multilingual KG itself.
696 Our system do not suffer from the error propaga-
697 tion introduced by entity and relation alignment
698 between different KGs. Nonetheless, our system is
699 limited to the setting of unified multilingual KGs
700 such as Wikidata. KG-Trick is complementary to
701 the other related work on multilingual KG comple-
702 tion which calls for integration of different KGs.

Our system can be applied to further improve the
completeness after the KGs are unified since such
techniques focus on fusing different KGs but not
inferring knowledge from the unified KG itself.

WikiKGE-10++ While WikiKGE-10++ signifi-
cantly extends WikiKGE-10 by adding 2X entities
sampled from torso and tail entities and descrip-
tions for such comprehensively sampled set of en-
tities, it contains only two types of facts including
entities names and descriptions. By extending the
WikiKGE-10++ to cover more facts associated with
the entities, the research community can be able to
get a more accurate and thorough picture on how
the proposed approach can improve KG.

Potential risks for generative textual infor-
mation completion As we employ a text-to-text
framework for multilingual KG textual information
completion task, it may generate biased or inaccur-
ate text that could be misleading for downstream
tasks. If this work is considered for production use,
human annotators might be in the loop to reduce
the risks of harmful text generation.

References

- Antoine Bordes, Nicolas Usunier, Alberto Garcia-
Duran, Jason Weston, and Oksana Yakhnenko.
2013. [Translating embeddings for modeling multi-
relational data](#). In *Advances in Neural Information
Processing Systems*, volume 26. Curran Associates,
Inc. 727-732
- Soumen Chakrabarti, Harkanwar Singh, Shubham Lo-
hiya, Prachi Jain, and Mausam. 2022. [Joint comple-
tion and alignment of multilingual knowledge graphs](#).
In *Proceedings of the 2022 Conference on Empiri-
cal Methods in Natural Language Processing*, pages
11922–11938, Abu Dhabi, United Arab Emirates. As-
sociation for Computational Linguistics. 733-739
- Jiahua Chen, Shuai Wang, Sahisnu Mazumder, and Bing
Liu. 2020a. [A knowledge-driven approach to clas-
sifying object and attribute coreferences in opinion
mining](#). In *Findings of the Association for Computa-
tional Linguistics: EMNLP 2020*, pages 1616–1626,
Online. Association for Computational Linguistics. 740-744
- Xuelu Chen, Muhao Chen, Changjun Fan, Ankith Up-
punda, Yizhou Sun, and Carlo Zaniolo. 2020b. Mul-
tilingual knowledge graph completion via ensemble
knowledge transfer. In *Proceedings of the 2020 Con-
ference on Empirical Methods in Natural Language
Processing: Findings*, pages 3227–3238. 746-751
- Louis Clouatre, Philippe Trempe, Amal Zouaq, and
Sarath Chandar. 2021. [MLMLM: Link prediction
with mean likelihood masked language model](#). In 752-753

755					
756					
757					
758	Simone Conia, Min Li, Daniel Lee, Umar Minhas, Ihab				
759	Ilyas, and Yunyao Li. 2023a. Increasing coverage				
760	and precision of textual information in multilingual				
761	knowledge graphs. In <i>Proceedings of the 2023 Con-</i>				
762	<i>ference on Empirical Methods in Natural Language</i>				
763	<i>Processing</i> , pages 1612–1634.				
764	Simone Conia, Min Li, Daniel Lee, Umar Minhas, Ihab				
765	Ilyas, and Yunyao Li. 2023b. Increasing coverage				
766	and precision of textual information in multilingual				
767	knowledge graphs . In <i>Proceedings of the 2023 Con-</i>				
768	<i>ference on Empirical Methods in Natural Language</i>				
769	<i>Processing</i> , pages 1612–1634, Singapore. Associa-				
770	tion for Computational Linguistics.				
771	Marta R Costa-jussà, James Cross, Onur Çelebi, Maha				
772	Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe				
773	Kalbassi, Janice Lam, Daniel Licht, Jean Maillard,				
774	et al. 2022. No language left behind: Scaling				
775	human-centered machine translation. <i>arXiv preprint</i>				
776	<i>arXiv:2207.04672</i> .				
777	Aidan Hogan, Eva Blomqvist, Michael Cochez, Clau-				
778	dia D’amato, Gerard De Melo, Claudio Gutierrez,				
779	Sabrina Kirrane, José Emilio Labra Gayo, Roberto				
780	Navigli, Sebastian Neumaier, Axel-Cyrille Ngonga				
781	Ngomo, Axel Polleres, Sabbir M. Rashid, Anisa Rula,				
782	Lukas Schmelzeisen, Juan Sequeda, Steffen Staab,				
783	and Antoine Zimmermann. 2021. Knowledge graphs .				
784	<i>ACM Comput. Surv.</i> , 54(4).				
785	Xuming Hu, Yong Jiang, Aiwei Liu, Zhongqiang Huang,				
786	Pengjun Xie, Fei Huang, Lijie Wen, and Philip S. Yu.				
787	2023. Entity-to-text based data augmentation for var-				
788	ious named entity recognition tasks . In <i>Findings of</i>				
789	<i>the Association for Computational Linguistics: ACL</i>				
790	<i>2023</i> , pages 9072–9087, Toronto, Canada. Associa-				
791	tion for Computational Linguistics.				
792	Seyed Mehran Kazemi and David Poole. 2018. Simple				
793	embedding for link prediction in knowledge graphs.				
794	<i>Advances in neural information processing systems</i> ,				
795	31.				
796	Adrian Kochsiek, Apoorv Saxena, Inderjeet Nair, and				
797	Rainer Gemulla. 2023. Friendly neighbors: Context-				
798	tualized sequence-to-sequence link prediction . In				
799	<i>Proceedings of the 8th Workshop on Representation</i>				
800	<i>Learning for NLP (RepLANLP 2023)</i> , pages 131–138,				
801	Toronto, Canada. Association for Computational Lin-				
802	guistics.				
803	Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch,				
804	Dimitris Kontokostas, Pablo N. Mendes, Sebastian				
805	Hellmann, Mohamed Morsey, Patrick van Kleef,				
806	Sören Auer, and Christian Bizer. 2015. Dbpedia -				
807	A large-scale, multilingual knowledge base extracted				
808	from wikipedia . <i>Semantic Web</i> , 6(2):167–195.				
809	Chin-Yew Lin. 2004. ROUGE: A package for auto-				
810	matic evaluation of summaries . In <i>Text Summariza-</i>				
811	<i>tion Branches Out</i> , pages 74–81, Barcelona, Spain.				
812	Association for Computational Linguistics.				
	Yankai Lin, Zhiyuan Liu, Huanbo Luan, Maosong Sun,				813
	Siwei Rao, and Song Liu. 2015a. Modeling relation				814
	paths for representation learning of knowledge bases .				815
	In <i>Proceedings of the 2015 Conference on Empirical</i>				816
	<i>Methods in Natural Language Processing</i> , pages 705–				817
	714, Lisbon, Portugal. Association for Computational				818
	Linguistics.				819
	Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and				820
	Xuan Zhu. 2015b. Learning entity and relation em-				821
	beddings for knowledge graph completion . In <i>Pro-</i>				822
	<i>ceedings of the Twenty-Ninth AAAI Conference on</i>				823
	<i>Artificial Intelligence, January 25-30, 2015, Austin,</i>				824
	<i>Texas, USA</i> , pages 2181–2187. AAAI Press.				825
	Nick Mckenna and Priyanka Sen. 2023. KGQA without				826
	retraining . In <i>Proceedings of The Fourth Workshop</i>				827
	<i>on Simple and Efficient Natural Language Process-</i>				828
	<i>ing (SustaiNLP)</i> , pages 212–218, Toronto, Canada				829
	(Hybrid). Association for Computational Linguistics.				830
	Maciej Modrzejewski, Miriam Exel, Bianka Buschbeck,				831
	Thanh-Le Ha, and Alexander Waibel. 2020. Incorpor-				832
	ating external annotation to improve named entity				833
	translation in NMT . In <i>Proceedings of the 22nd</i>				834
	<i>Annual Conference of the European Association for</i>				835
	<i>Machine Translation</i> , pages 45–51, Lisboa, Portugal.				836
	European Association for Machine Translation.				837
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-				838
	Jing Zhu. 2002. Bleu: a method for automatic evalu-				839
	ation of machine translation. In <i>Proceedings of the</i>				840
	<i>40th annual meeting of the Association for Computa-</i>				841
	<i>tional Linguistics</i> , pages 311–318.				842
	Ciyuan Peng, Feng Xia, Mehdi Naseriparsa, and				843
	Francesco Osborne. 2023. Knowledge graphs: Op-				844
	portunities and challenges . <i>Artificial Intelligence</i>				845
	<i>Review</i> .				846
	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon				847
	Lavie. 2020. COMET: A neural framework for MT				848
	evaluation . In <i>Proceedings of the 2020 Conference</i>				849
	<i>on Empirical Methods in Natural Language Process-</i>				850
	<i>ing (EMNLP)</i> , pages 2685–2702, Online. Association				851
	for Computational Linguistics.				852
	Ridho Reinanda, Edgar Meij, and Maarte de Rijke. 2020.				853
	Knowledge graphs: An information retrieval perspec-				854
	tive. <i>Foundations and Trends® in Information Re-</i>				855
	<i>trieval</i> , 14(4):289–444.				856
	Ryokan Ri, Ikuya Yamada, and Yoshimasa Tsuruoka.				857
	2022. mLUKE: The power of entity representations				858
	in multilingual pretrained language models . In <i>Pro-</i>				859
	<i>ceedings of the 60th Annual Meeting of the Associa-</i>				860
	<i>tion for Computational Linguistics (Volume 1: Long</i>				861
	<i>Papers)</i> , pages 7316–7330, Dublin, Ireland. Associa-				862
	tion for Computational Linguistics.				863
	David E Rumelhart, Geoffrey E Hinton, Ronald J				864
	Williams, et al. 1985. Learning internal represen-				865
	tations by error propagation.				866

867	Apoorv Saxena, Adrian Kochsiek, and Rainer Gemulla.	and pre-trained language representation. <i>Transactions of the Association for Computational Linguistics</i> , 9:176–194.	922
868	2022. Sequence-to-sequence knowledge graph completion and question answering . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2814–2828, Dublin, Ireland. Association for Computational Linguistics.		923
869			924
870			
871		Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhengyan Zhang, Zhiyuan Liu, Juanzi Li, and Jian Tang. 2021b. KEPLER: A unified model for knowledge embedding and pre-trained language representation . <i>Transactions of the Association for Computational Linguistics</i> , 9:176–194.	925
872			926
873			927
874			928
875	Priyanka Sen, Alham Fikri Aji, and Amir Saffari. 2022. Mintaka: A complex, natural, and multilingual dataset for end-to-end question answering . In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 1604–1619, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.		929
876			930
877		Ruobing Xie, Zhiyuan Liu, Jia Jia, Huanbo Luan, and Maosong Sun. 2016. Representation learning of knowledge graphs with entity descriptions . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 30(1).	931
878			932
879			933
880			934
881			935
882	Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. 2019. Rotate: Knowledge graph embedding by relational rotation in complex space . In <i>International Conference on Learning Representations</i> .	Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. 2014. Embedding entities and relations for learning and inference in knowledge bases. <i>arXiv preprint arXiv:1412.6575</i> .	936
883			937
884			938
885			939
886	Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. 2016a. Complex embeddings for simple link prediction . In <i>International conference on machine learning</i> , pages 2071–2080. PMLR.	Shuai Zhang, Yi Tay, Lina Yao, and Qi Liu. 2019. Quaternion knowledge graph embeddings . <i>Advances in neural information processing systems</i> , 32.	940
887			941
888			942
889			
890	Théo Trouillon, Johannes Welbl, Sebastian Riedel, Eric Gaussier, and Guillaume Bouchard. 2016b. Complex embeddings for simple link prediction . In <i>Proceedings of The 33rd International Conference on Machine Learning</i> , volume 48 of <i>Proceedings of Machine Learning Research</i> , pages 2071–2080, New York, New York, USA. PMLR.	A Creating WikiKGE-10++	943
891		In this section, we describe the in-depth details on the creation of WikiKGE-10++, our novel human-curated dataset for evaluation automatic approaches on KGE of Wikidata entity names and descriptions.	944
892			945
893			946
894			947
895			948
896			
897	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	A.1 Choice of Languages	949
898		Aligned with the previous work completed in Cونيا et al. , the benchmark sustains the selection of 9 languages from a set of topologically diverse linguistic families, while interchanging the Russian (Slavic) language for the Turkish (Altaic) language:	950
899			951
900			952
901			953
902	Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. <i>Communications of the ACM</i> , 57(10):78–85.	• West Germanic: English, German;	954
903			955
904			956
905	Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao, Wenyang Gao, Xuming Hu, Zehan Qi, Yidong Wang, Linyi Yang, Jindong Wang, Xing Xie, Zheng Zhang, and Yue Zhang. 2023. Survey on factuality in large language models: Knowledge, retrieval and domain-specificity .	• Romance: Spanish, French, Italian;	957
906			958
907			959
908			960
909			961
910			
911			
912	Liang Wang, Wei Zhao, Zhuoyu Wei, and Jingming Liu. 2022. SimKGC: Simple contrastive knowledge graph completion with pre-trained language models . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4281–4294, Dublin, Ireland. Association for Computational Linguistics.	• Sino-Tibetan: Chinese (simplified);	962
913			963
914			964
915			965
916			
917			
918			
919	Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhengyan Zhang, Zhiyuan Liu, Juanzi Li, and Jian Tang. 2021a. Kepler: A unified model for knowledge embedding	• Altaic: Turkish;	962
920			963
921			964
			965
		The Russian language was interchanged for the Turkish language due to export and import restrictions placed on Russia, thereby, restricting access to Russia-based human annotators.	

966	A.2 Human annotation process	
967	The objective of the annotation process was to (i)	1014
968	rate and suggest entity names in the target language,	1015
969	(ii), verify the suggest entity names in the target	1016
970	languages, (iii) curate description for the entity	1017
971	in the target language, (iv) validate the provided	1018
972	descriptions quality.	1019
973		1020
974	A.2.1 Rate and suggest entity names.	1021
975	The objective of the annotation process was the	1022
976	rate entity names in a target language. Detailed in-	1023
977	formation on the annotation process and UI design	1024
978	can be found in Conia et al. .	1025
979		1026
980	A.2.2 Verify suggested entity names.	1027
981	The objective of the annotation process was the	1028
982	verify the suggested entity names in the target lan-	1029
983	guage provided by the human annotators. Detailed	1030
984	information on the annotation process and UI de-	1031
985	sign can be found in Conia et al. .	1032
986		1033
987	A.2.3 Curate entity descriptions.	1034
988	The objective of the annotation process was cu-	1035
989	rate descriptions for the given entity in the target	1036
990	language.	1037
991	Given an entity name in a target language, an-	1038
992	notations were required to familiarize themselves	1039
993	with the its information: the user interface provided	1040
994	the entity names, as well as a built-in panel that	1041
995	directly displayed Wikipedia articles for the corre-	1042
996	sponding entity in English and the target language,	1043
997	if available. In addition, annotators were recom-	1044
998	mended to further familiarize themselves with the	1045
999	entity outside of the provided information.	1046
1000	Next, the annotators were tasked with learning	1047
1001	about the required format of the requested descrip-	1048
1002	tion with detailed instructions. This was facilitated	1049
1003	by providing: (i) examples of correctly curated de-	1050
1004	scription given an example entity, and (ii) strict	1051
1005	rules that the descriptions had to comply by.	1052
1006	After learning about the entity and the required	1053
1007	description format, the human annotator was re-	1054
1008	quested to manually curate the description for the	1055
1009	corresponding entity in the target language. During	1056
1010	this task, the human annotator was provided with	1057
1011	descriptions from other sources (such as Wikidata)	1058
1012	in English and the target language. Human anno-	1059
	tators were instructed that they could leverage the	1060
	extraneous descriptions, but not to copy and paste	1061
	unless satisfactory.	
	A.2.4 Validate entity descriptions.	
	The objective of the annotation process was the	
	validate the quality of the descriptions in the target	
	language provided by the human annotators.	
	First, given an entity, the human annotator was	
	provided corresponding information (i.e., entity	
	names/aliases, Wikipedia pages, etc) as done in the	
	previous task. In addition, the description require-	
	ments were detailed (in-depth guidelines provided	
	in a different document).	
	Then, they were prompted to analyze the cor-	
	responding description for the entity in the target	
	language with a series of questions. The ques-	
	tions were reformulated to from description require-	
	ments, to verify the presented description in the	
	target language followed the requested format. If	
	the annotator negatively responded to any of the	
	presented questions, they were prompted to edit	
	the description to satisfy the requirements. If the	
	initial description the requirements, the originally	
	provided description was sustained.	
	A.3 Quality assurance and inter-annotator	
	agreement.	
	B Short Description Evaluation	
	As is shown in Table Table 4, we calculate the	
	BLEU score for every baseline and our method.	
	However, the BLEU score for all languages and	
	all baselines are under 10, which suggests that the	
	translated text from English can hardly relate to	
	ground truth in target languages. This phenomenon	
	could suggest that BLEU is not a proper metric for	
	entity short description evaluation, as (i) Short de-	
	scription for the same entity in different languages	
	are not directly translatable. (ii) A large amount	
	of short descriptions are less than 4 tokens (e.g.	
	<i>Politician</i>), which could biases the judgement of	
	BLEU when calculating the weighted average.	
	C Balancing the training data for KGC	
	and KGE	
	In this section, we provide more details on how	
	balancing the training data between KGC and KGE	
	tasks can impact the performance of the two tasks.	
	Indeed, the training datasets available for the two	
	tasks are not balanced: the KGC dataset contains	
	150 million records generated multilingually from	
	20 million triplets, while the KGE dataset contains	
	around 16 million records generated multilingually	
	from 5 million entities. Therefore, we investigate	
	the impact of mixing different proportions of the	

Projects / 1287038 / Tasks / 01HA1GMA2H4GFRNY7DFR5HC7B9

Task Type: Writing Descriptions in the Target Language | Request ID: None | Estimated Rating Time: 5 minutes | Task Purpose: Regular

Buttons: Rating Guidelines, View Survey JSON, View Ratings, Super Rate, Validate Ratings

Task ID: 01HA1GMA2H4GFRNY7DFR5HC7B9 | Show Ratings: Select a rating... | View Ratings: 01HA1GMA2H4GFRNY7DFR5HC7B9-0 | View Ratings JSON

Instructions

In this task, you will be presented with an Entity and the following information:

- Entity Name in English and Target Language (German)
- Descriptions from Other Sources
- Wikipedia Page in English and Target Language

Using the information above, please create a description, following the instructions below.

- Familiarize yourself with the Entity.** You can do so by searching the entity on an internet search engine (ex. Google) or the entity on Wikipedia (recommended). Please spend time on researching and understanding the entity, as we provide ample time.
- Read the information provided.** In particular, please pay attention to the descriptions from other sources that may have been used. Do not simply copy and paste these descriptions from other sources.
- Write the description.** Please pay attention to the instructions below, on how to create a good description. Write the description in the target language. Ensure you follow our guidelines!

The target language of this task is: [German](#)

Note: Please thoroughly familiarize yourself with the Guidelines before answering the questions and their tasks below. The guidelines are short, and should be frequently referenced throughout the task.

Entity Information

In English, the Entity is commonly known as "Mount Roraima" or:

[Cerro Roraima](#) [Monte Roraima](#) [Mount Roraima](#) [Mt. Roraima](#) [Pico do Roraima](#) [Roraima Tepui](#) [Roraima mountain](#) [Roraima-tepui](#)

In German, the Entity is commonly known as [Roraima-Tepui](#) or:

[Cerro Roraima](#) [Monte Roraima](#) [Mount Roraima](#) [Roraima Tafelberg](#)

English Wikipedia: [Read more about Mount Roraima in English](#) | German Wikipedia: [Read more about Mount Roraima in German](#)

Figure 2: UI used for the annotation task: the annotators could familiarize themselves with the task with an outline of the task instructions (detailed guidelines could be read in a separate page) and the information about the entity, including its names in English and its Wikipedia pages in English and the target language (Italian in this case).

	BLEU Score	#Params	AR	DE	ES	FR	IT	JA	KO	ZH	Avg
<i>descriptions</i>	NLLB-200 _{EN→XX} →	0.6B	2	3.4	7	4.8	4.6	1.5	4.1	5.9	4.2
	GPT-3.5	175B	2.8	4.2	7.2	4.9	5.5	3.9	4.3	9.8	5.3
	TRICK _{KGE}		2.2	3.4	8.0	5.1	5.9	3.8	2.5	4.9	4.5
	TRICK _{50%KGC+KGE}	0.6B	1.4	2.1	5.8	2.8	4.1	2.3	1.6	2.8	2.9
	TRICK _{KGC+KGE}		1	0.8	2.6	1.3	2.1	2	1.5	2.2	1.7

Table 4: BLEU score for entity short description evaluation

1062 two datasets on the performance of the two tasks.
 1063 More specifically, we investigate different propor-
 1064 tions of the KGC and KGE datasets, ranging from
 1065 0% to 100% of the KGC dataset, and evaluate the
 1066 performance of the two tasks on WikiKGE-10+.
 1067 The results are reported in Table 5. We can observe
 1068 that the best performance on KGC is achieved when
 1069 the full KGC dataset is used, which suggests that
 1070 the KGC task is more difficult than the KGE task.
 1071 On the other hand, the best performance on KGE
 1072 is achieved when up to 50% of the KGC dataset is
 1073 used. Therefore, the best compromise between the
 1074 data mixing proportion for the two tasks is to use
 1075 50% of the KGC dataset.

KGC%	KGE		KGC	
	Precision	Coverage	MRR	hit@1
0%	51.5	30.9	-	30.4
1%	52.4	30.4	32.7	32.1
10%	52.4	30.6	34.4	31.7
20%	51.7	28.8	34	32.8
50%	52.1	31.5	35.2	33
full	52.2	29.7	38.8	36.6

Table 5: Investigation on the different data mixing proportion between KGC and KGE, and their impact on KGC and KGE tasks performance

Here are rules to pay attention to (Example Entity - Osaka):

- **DO NOT** use the Entity Name in the description!
 - **Bad Example:** **Osaka** is a designated city in the Kansai region of Honshu in Japan
- **DO NOT** start the description with a verb!
 - **Bad Example:** **Is a** designated city in the Kansai region of Honshu in Japan.
- Keep it **short** and **concise**, under 30 words! But make sure to **capture important information!**
 - **Bad Example:** home to **Osaka Castle** and **Universal Studios**, has one of the **largest aquarium**, birthplace of **instant noodles**
- **DO NOT** separate facts with **periods**. Separate facts with **commas!** (ie. Don't use periods -> Use Commas)
 - **Bad Example:** designated city in the Kansai region, one of three major cities, third most populous city in Japan,
- Use correct **grammar, spelling and fluency**, but follow the rules above!
 - **Bad Example:** designated city in **Kansai**, it **could though be seen** in historic times **that is though the** most populous in country **japan**
- The **first word** should be **uncapitalized** (unless it is a proper noun [ex. Country Name, Title, etc]), and there should be **no period** at the end!
 - **Bad Example:** **D**esignated city in the Kansai region,

Question 1: Please write the description for "Mount Roraima" below in German.

Provided below are descriptions from other sources:

One-line description of "Mount Roraima" in English is:

- *High plateau in South America*

One-line description(s) of "Mount Roraima" in German are:

- *Tepui im Dreiländereck Venezuela, Brasilien und Guyana*
- *der höchste Gipfel des Pacaima-Gebirges auf dem Plateau von Guyana im Norden Südamerikas*
- *Berg Südamerikas*
- *Berge Südamerikas*
- *Hochplateau in Südamerika*
- *Berg, der sich zwischen Venezuela, Brasilien und Guyana erstreckt*
- *höchster Punkt in Guyana*

You can re-use components of descriptions provided above, if they hold key facts! **DO NOT simply copy and paste** one of them. Please research the entity.

Figure 3: UI used for the annotation task: the annotators familiarized themselves with the description format with an outline of the requirements (detailed guidelines could be read in a separate page).

- *High plateau in South America*

One-line description(s) of "Mount Roraima" in German are:

- *Tepui im Dreiländereck Venezuela, Brasilien und Guyana*
- *der höchste Gipfel des Pacaima-Gebirges auf dem Plateau von Guyana im Norden Südamerikas*
- *Berg Südamerikas*
- *Berge Südamerikas*
- *Hochplateau in Südamerika*
- *Berg, der sich zwischen Venezuela, Brasilien und Guyana erstreckt*
- *höchster Punkt in Guyana*

You can re-use components of descriptions provided above, if they hold key facts! **DO NOT simply copy and paste** one of them. Please research the entity.

Word count: 0

Section 2: Feedback [OPTIONAL]

Please let us know if something is wrong with this task assignment. For example, something is wrong with the user interface, one or more questions are unclear, or you could not do something you wanted to.

Super Rate Validate Ratings

Figure 4: UI used for the annotation task: the annotator provided the description in a text box. A warning message was prompted if the token length of the description was too short or too long.

Here are rules to pay attention to (Example Entity - Osaka):

- **RULE 1: DO NOT use the Entity Name in the description!**
 - **Bad Example:** **Osaka** is a designated city in the Kansai region of Honshu in Japan
- **RULE 2: DO NOT start the description or a fact with a verb!**
 - **Bad Example:** **Is a** designated city in the Kansai region of Honshu in Japan, **and it is** one of three major cities
- **RULE 3: Keep it short and concise, under 30 words! But make sure to capture important information!**
 - **Bad Example:** home to **Osaka Castle** and **Universal Studios**, has one of the **largest aquarium**, birthplace of **instant noodles**
- **RULE 4: DO NOT separate facts with periods. Separate facts with commas! (ie. Don't use periods -> Use Commas)**
 - **Bad Example:** designated city in the Kansai region, one of three major cities, third most populous city in Japan.
- **RULE 5: The first word should be uncapitalized (unless it is a proper noun [ex. Country Name, Title, etc]), and there should be no period or comma at the end!**
 - **Bad Example:** Designated city in the Kansai region.
- **RULE 6: Use correct grammar, spelling and fluency, but follow all the rules!**
 - **Bad Example:** designated city in Kansai, **it could though be seen** in historic times **that is though the** most populous in country japan

Question 1: Please verify the description for "Anton Bruckner" below in German.

One-line description of "Anton Bruckner" in German is:

- *österreichischer Komponist der Romantik, der zu den wichtigsten und innovativsten Tonschöpfern seiner Zeit gehört*

Part A: Does the Description include the Entity Name? (RULE 1)

- Yes - The Description includes the Entity Name. (REVISION NEEDED - DESCRIPTION SHOULD NOT INCLUDE THE ENTITY NAME)
- No - The Description does not include the Entity Name.

Part B: Does the Description or Facts start with a Verb? (RULE 2)

- Yes - The Description or Facts starts with a Verb. (REVISION NEEDED - DESCRIPTION SHOULD NOT START WITH A VERB)
- No - The Description or Facts do not start with a Verb.

Part C: Are facts in the Description separated by Periods? (RULE 4)

- Yes - The facts in the Description are separated with Periods. (REVISION NEEDED - DESCRIPTION SHOULD SEPARATE FACTS WITH COMMAS)
- No - The facts in the Description are separated by commas.

Part D: Does the first word in the Description start with a capital (unless it is a proper noun [ex. Country Name, Title, etc]) (RULE 5)

- Yes - The first word in the Description starts with a Capital. (REVISION NEEDED - DESCRIPTION SHOULD NOT START WITH A CAPITALIZED WORD)
- No - The Description does not start with a Capital.

Figure 5: UI used for the annotation task: annotators were required to examine the description in the target language, and answer a series of questions that reflected the description requirements.

- Yes - The Description or Facts starts with a Verb. (REVISION NEEDED - DESCRIPTION SHOULD NOT START WITH A VERB)
- No - The Description or Facts do not start with a Verb.

Part C: Are facts in the Description separated by Periods? (RULE 4)

- Yes - The facts in the Description are separated with Periods. (REVISION NEEDED - DESCRIPTION SHOULD SEPARATE FACTS WITH COMMAS)
- No - The facts in the Description are separated by commas.

Part D: Does the first word in the Description start with a capital (unless it is a proper noun [ex. Country Name, Title, etc]) (RULE 5)

- Yes - The first word in the Description starts with a Capital. (REVISION NEEDED - DESCRIPTION SHOULD NOT START WITH A CAPITALIZED WORD)
- No - The Description does not start with a Capital.

Part E: Does the Description end in a period or comma? (RULE 5)

- Yes - The Description ends with a period or comma. (REVISION NEEDED - DESCRIPTION SHOULD NOT END WITH A PERIOD OR COMMA)
- No - The Description does not end with a period or comma.

Part F: Does the Description use correct Spelling, Capitalization and Punctuation?

- Yes - The Description uses correct Spelling.
- No - The Description does not use correct Spelling. (REVISION NEEDED)

Part G: Is the Description not complete? Are there any crucial or important details you would add?

- Yes - There are crucial or important details that should be added. (REVISION NEEDED - ADDITIONAL DETAILS SHOULD BE ADDED)
- No - The Description is complete with crucial and important details.

Section 2: Optional Feedback [DO NOT PUT THE DESCRIPTION HERE]

Please let us know if something is wrong with this task assignment. For example, something is wrong with the user interface, one or more questions are unclear, or you could not do something you wanted to.

Super Rate Validate Ratings

Figure 6: UI used for the annotation task: annotators were prompted to correct the description by rewriting it, if they negatively answer the series of questions provided.