

# Entity Insertion in Multilingual Linked Corpora: The Case of Wikipedia

Tomás Feith,<sup>\*◇</sup> Akhil Arora,<sup>\*†♣</sup> Martin Gerlach,<sup>♣</sup> Debjit Paul,<sup>◇</sup> Robert West<sup>‡‡</sup>

<sup>◇</sup>EPFL <sup>♣</sup>Aarhus University <sup>♣</sup>Wikimedia Foundation  
tsfeith@gmail.com, akhil.arora@cs.au.dk, mgerlach@wikimedia.org,  
{debjit.paul, robert.west}@epfl.ch

## Abstract

Links are a fundamental part of information networks, turning isolated pieces of knowledge into a network of information richer than the sum of its parts. However, adding a new link to the network is not trivial: it requires not only the identification of a suitable pair of source and target entities but also the understanding of the content of the source to locate a suitable position for the link in the text. The latter problem has not been addressed effectively, particularly in the absence of text spans in the source that could serve as anchors to insert a link to the target entity. To bridge this gap, we introduce and operationalize the task of *entity insertion* in information networks. Focusing on the case of Wikipedia, we empirically show that this problem is, both, relevant and challenging for editors. We compile a benchmark dataset in 105 languages and develop a framework for entity insertion called LOCEI (Localized Entity Insertion) and its multilingual variant xLOCEI. We show that xLOCEI outperforms all baseline models (including state-of-the-art prompt-based ranking with LLMs such as GPT-4) and that it can be applied in a zero-shot manner on languages not seen during training with minimal performance drop. These findings are important for applying entity insertion models in practice, e.g., to support editors in adding links across the more than 300 language versions of Wikipedia.

## 1 Introduction

From digital encyclopedias and blogs to knowledge graphs, knowledge on the Web is organized as a network of interlinked entities and their descriptions. However, online knowledge is not static: new webpages are created, and existing pages are updated almost every day. While there exists substantial support for content creation (e.g. via translation Wulczyn et al. (2016) or generative AI tools Shao et al.

<sup>\*</sup>Equal contribution, contact: [akhil.arora@cs.au.dk](mailto:akhil.arora@cs.au.dk).

<sup>†</sup>Work done at EPFL.

<sup>‡‡</sup>R. West is a Wikimedia Foundation Research Fellow.

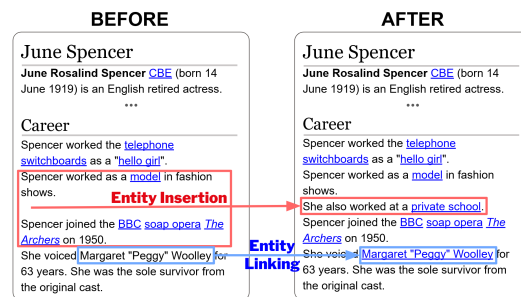


Figure 1: *Entity linking*: insert a link to the entity Margaret “Peggy” Woolley by identifying a suitable mention from the existent text in the version **before** insertion, vs. *Entity insertion*: no mention existent yet, identify the most suitable span  in the version **before** to insert the entity Private school.

(2024)), adding new knowledge not only requires creating content but also integrating it into the existing knowledge structure. The latter usually leaves editors with the time-consuming task of reading lengthy webpages to identify a relevant text span for inserting an entity that is not yet mentioned on the page. Thus, to support editors in effectively integrating entities in multilingual linked corpora on the Web, we introduce the task of *entity insertion*.

**Entity insertion.** We consider Wikipedia as the primary use case and focus on the task of adding links. Specifically, given a source and target entity, the goal of *entity insertion* is to identify the most suitable text span in the article describing the source entity for inserting a link to the target entity. Fig. 1 portrays a real example of the entity insertion task with the eventual goal of adding a link from the source entity June Spencer, a former English actress, to the target entity Private school. Most importantly, entity insertion is a different and much more challenging task when compared to *entity linking*, as no existent text span in the version of the source article (June Spencer) at edit time could be used to link to the target entity (Private school). Rather, a new text span—“*She also worked at a private school.*”—was added along with the to-be-inserted target entity.

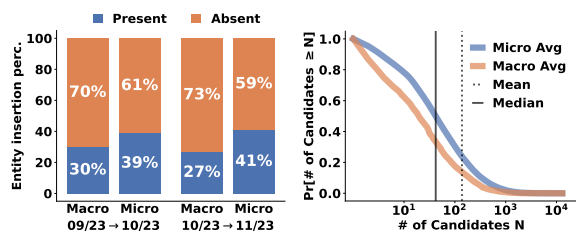


Figure 2: Challenges of entity insertion. (Left) Micro (weighted by the number of data points in a language) and macro (equal weight to each language) aggregates of insertion types over the 105 languages considered in this study. (Right) Complementary cumulative distribution function (CCDF) of the number of candidate sentences ( $N$ ) in a Wikipedia article (log x-axis).

**Challenges.** Entity insertion is not only an interesting and challenging language understanding task, but it is also the most common scenario faced by editors when adding links in practice. In fact, we find that for 60-70% of all the links added to Wikipedia, *none of the existing text* is suitable to insert the corresponding entities, and new text needs to be added along with the entity by the editor (Fig. 2). Fig. 2 also shows that entity insertion is associated with a *high cognitive load*, as the task requires, on average, an editor to select the most suitable sentence from a pool of  $\sim 100$  candidate sentences.

Therefore, it is vital to operationalize and develop new methods for entity insertion in order to support editors in adding links to Wikipedia and other information networks. To this end, we make the following key contributions in this paper.

**Contributions.** We introduce the *novel task of entity insertion* (§ 3). We release a *large dataset in 105 languages* of links from Wikipedia articles to enable further research into entity insertion (§ 4). We introduce LOCEI, a framework for entity insertion, and its multilingual variant xLOCEI (§ 5). We show the benefit of multilingual knowledge in downstream performance and highlight the *zero-shot* capabilities of xLOCEI (§ 6).

## 2 Related work

In this section, we review works that overlap closely with our study (cf. Appx. A for details).

**Entity linking.** Previous work has framed entity insertion as an entity linking problem (Gerlach et al., 2021; Milne and Witten, 2008; West et al., 2009; Arora et al., 2021; Čuljak et al., 2022; West et al., 2010), where the goal is to assign a unique identity to entities mentioned in the text. The task of entity linking is composed of two sub-tasks: Named Entity Recognition (NER) and Named Entity Dis-

ambiguation (NED). Most research (Hoffart et al., 2011; Fu et al., 2020; van Hulst et al., 2020) into entity linking solves first the NER problem, in which the task is to find candidate mentions for named entities in the source article. However, there is recent work (Zhang et al., 2022) exploring the problem in reverse order, first solving NED by finding target entities related to the source article and then NER searching only for mentions for the found targets.

When the mention is present, the task of entity insertion is similar to NER (Zhang et al., 2022), as both tasks can be solved by searching for mentions in the text. However, entity insertion is a more general task as it allows for the mention of the target entity to not yet be present in the text. In this case, the goal is to exploit the context information to find the text span most related to the target entity. NER modules (Finkel et al., 2005; Nothman et al., 2013) are designed to search for the most related mentions, and thus, they are not applicable in scenarios where the mentions are not yet available.

**Entity tagging.** Du et al. (2022) introduced this task as a relaxed form of entity linking. An entity tagging model is only tasked with determining the entities present in the text and does not need to find the exact mentions of the entities. However, even though the model is not tasked with extracting an entity’s mention, the task of entity tagging still assumes that the text contains some mention of the entity, which distinguishes it from entity insertion.

**Link the Wiki.** Huang et al. (2008) ran a track at INEX 2008 with two tasks: file-to-file link discovery and mention-to-BEP (best entry point) link discovery. File-to-file link discovery is a document-level task that can be framed as a link prediction task in networks, where the Wikipedia articles act as nodes and the links act as edges. The mention-to-BEP task is an entity linking task with anchor prediction, where the two-part goal is to find mentions in the source article pointing to other articles, and finding the best point of entry (the anchor) in the target file. This task has more recently resurfaced as an anchor prediction task (Liu et al., 2023).

**Passage ranking.** Transformer-based models have revolutionized passage ranking by enhancing semantic understanding beyond traditional lexical methods like BM25 (Robertson and Zaragoza, 2009). BERT demonstrated early success by leveraging contextualized embeddings for re-ranking (Nogueira and Cho, 2019), leading to innovations

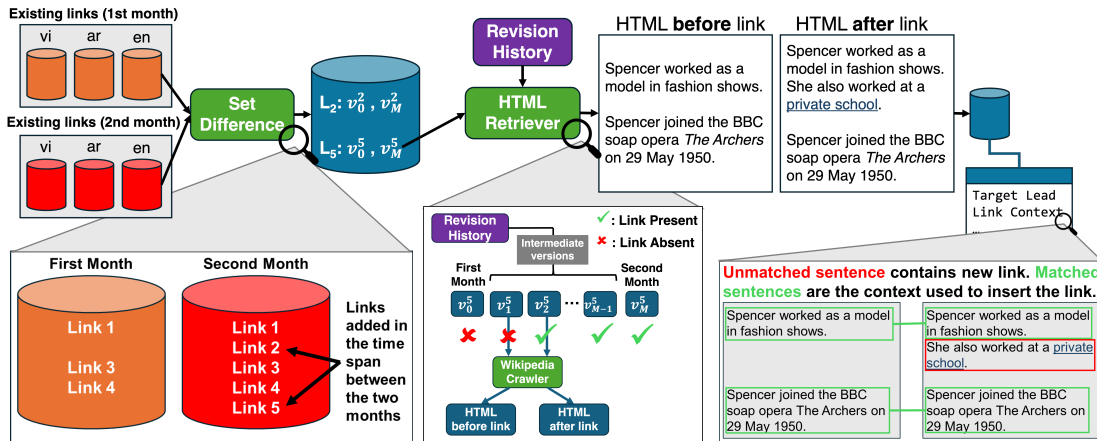


Figure 3: Data processing pipeline. Obtain added links  $L$  by taking a set difference of the links existent in consecutive months. For each added link  $L_i$ , scan all  $M$  versions in the full revision history  $v_0^i$  to  $v_M^i$  to identify the article version in which the link was added and compute the difference between the before and after versions to extract the exact entity insertion scenario.

like ColBERT (Khattab and Zaharia, 2020), which uses a dual-encoder architecture for more efficient retrieval. Recent models such as T5 (Raffel et al., 2020) and ELECTRA (Clark et al., 2020) further refine ranking by employing advanced pre-training techniques. Building on top of this work, (Fang et al., 2023; Dong et al., 2022) employ knowledge graphs to exploit background information to better rank passages. However, such graph-based approaches are not suited for large-scale, highly dynamic graphs (such as Wikipedia), as the cost of recomputing all the embeddings associated with the graph is too high. Finally, while large language models have been shown to be the state of the art for passage ranking (Qin et al., 2024), despite their performance they are impractical at the Web-scale owing to exorbitantly high computational costs.

**Key differences.** Entity insertion is fundamentally different from all the aforementioned tasks and possesses novel downstream applications. First, entity insertion does not assume that a mention to the target entity is present in the text at inference time. Second, the optimization objective of entity insertion, which involves identifying the text span most related to a target entity, could be seen as the dual of tasks such as NED and entity tagging, which aim instead to find the most relevant target entity for a given text span. Finally, entity insertion aims to find the best text span in the source article to insert the target entity. In contrast, anchor prediction performs the reverse task by trying to find the best text span for grounding the source entity in the target article. Moreover, anchor prediction is an unnatural task as humans find the vast majority of links to be unanchorable (Liu et al., 2023).

### 3 Task formulation

Let  $E_{src}$  be a source entity and  $E_{tgt}$  be a target entity. Let  $X_{src}$  be the textual content of the article corresponding to  $E_{src}$ . The text can be partitioned into a set of (potentially overlapping) text spans,  $\mathcal{X}_{src} = \{x_1, \dots, x_M\}$ , where  $M$  is the number of text spans in the article. Entity insertion is the task of selecting the most relevant span  $x^*$  to insert the target entity  $E_{tgt}$ . Formally,

$$x^* = \arg \max_{x \in \mathcal{X}_{src}} \mathcal{R}(x, E_{tgt}) \quad (1)$$

where  $\mathcal{R}$  is an arbitrary relevance function quantifying the relevance of  $E_{tgt}$  to each text span  $x \in \mathcal{X}_{src}$ . We frame entity insertion as a *ranking task*, where the goal is to rank all the candidate text spans  $\mathcal{X}_{src}$  based on their relevance to the target entity.

### 4 Data

We constructed a new multilingual dataset for studying entity insertion in Wikipedia. The dataset consists of links extracted from all Wikipedia articles, each link’s surrounding context, and additional article-level meta-data (such as titles, Wikidata QIDs, and lead paragraphs). Overall, the dataset contains 958M links from 49M articles in 105 languages. The largest language is English (en), with 166.7M links from 6.7M articles, and the smallest language is Xhosa (xh), with 2.8K links from 1.6K articles (cf. Appendix B for details).

Fig. 3 provides an overview of our data processing pipeline. The data processing was done in two steps. We first extracted all the links from the 2023-10-01 snapshot. Next, we found all the links added in the time between 2023-10-01 and 2023-11-01.

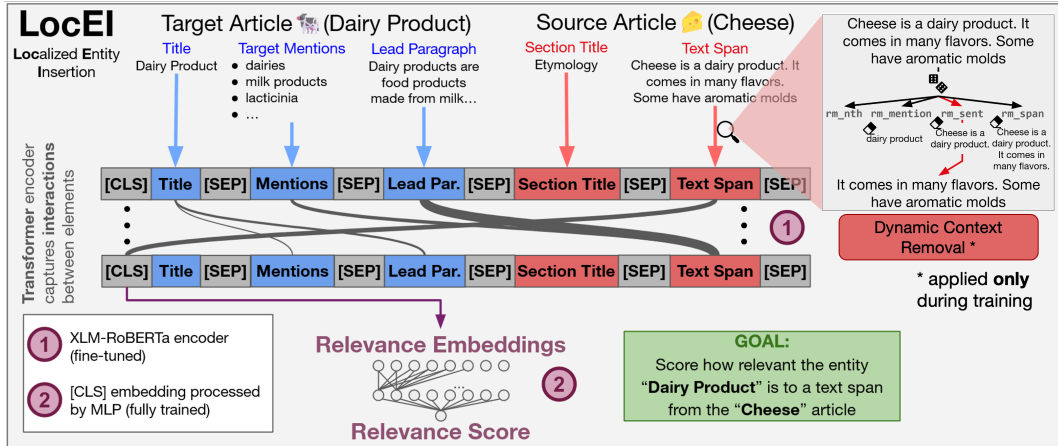


Figure 4: Architectural overview of LOCEI. The target entity  $E_{tgt}$  and each candidate text span  $x \in \mathcal{X}_{src}$  of the source entity  $E_{src}$  are concatenated together and encoded jointly using a transformer encoder. The relevance scores of candidate text spans are computed using an MLP trained via a list-wise ranking objective.

**Existing links.** We extract the content of all articles from their HTML version using the corresponding snapshot of the Enterprise HTML dumps (WMF, 2010b). We removed articles without a lead paragraph and a Wikidata QID. For each article, we consider all internal links in the main article body (ignoring figures, tables, notes, and captions) together with their surrounding context. We removed all the links where either the source or the target article was one of the removed articles and we dropped all the self-links.

**Added links.** We extract the set of added links by comparing existing links in snapshots from consecutive months. We apply the same procedure as above to each snapshot, respectively, and take the difference of the two sets to identify the links that exist in the second month but not in the first.

To identify the article version in which the link was added, we go through the articles’ full *revision history* available in the Wikimedia XML dumps (WMF, 2010a). Next, we identify the two versions of an article before and after the link addition and download the corresponding HTML. Comparing the two HTML versions, we extract the content modifications made by the editor when adding the link, and categorize them into *five entity insertion scenarios*. (1) *text\_present*: the link was added by hyperlinking an existing mention; (2) *missing\_mention*: the link was added by adding the mention for a new entity (and potentially some additional content) into an existent sentence; (3) *missing\_sentence*: the link was added by writing a new sentence to complement the already existing text and hyperlinking part of the sentence at the same time; (4) *missing\_span*: an extension of the previous category, where the editors added a span

of multiple sentences; and (5) *missing\_section*: the link was added in a section that did not exist in the previous version of the article. We provide examples (Table 9) and frequency of occurrence (Fig. 5) of these cases in Appendix B.3.

**Data release.** The dataset is made publicly available on Zenodo under an open license (CC-BY-SA 4.0) at <https://zenodo.org/records/13888211>.

## 5 Entity insertion with LOCEI

Fig. 4 presents an overview of LOCEI. Our model (§ 5.1) is composed of a transformer-based encoder that jointly encodes the target entity as well as the candidate spans in the source entity, and a multilayer perceptron (MLP) trained via a list-wise objective capable of ranking candidates based on their relevance to the target. We introduce a novel data augmentation strategy that closely mimics real-world entity insertion scenarios (§ 5.2), a knowledge injection module to incorporate external knowledge about the entities (§ 5.3), and the multilingual variant xLOCEI (§ 5.4).

### 5.1 Model

**Architecture.** We use a transformer-based encoder  $\Gamma$  to jointly encode each candidate text span  $x \in \mathcal{X}_{src}$  and the target entity  $E_{tgt}$  into a sequence of vectors. To reduce this sequence into a single vector, we use the embedding of the [CLS] token (Devlin et al., 2019), which measures how related the candidate  $x$  is to the entity  $E_{tgt}$ . An MLP  $\Lambda$  produces a scalar relevance score between the candidate  $x$  and the entity  $E_{tgt}$  using the relevance embedding produced by the encoder  $\Gamma$  defined as

$$R = \Gamma(\phi; \theta_{\Gamma}) \quad (2)$$

$$r = \Lambda(R_{[CLS]}; \theta_{\Lambda}) \quad (3)$$

where  $\Gamma$  is an encoder and  $\Lambda$  is an MLP with  $\theta_\Gamma$  and  $\theta_\Lambda$  as the learnable parameter spaces, respectively,  $\phi$  is obtained by concatenating the input representations of  $E_{tgt}$  and  $x$ ,  $R$  is the sequence of  $d$ -dimensional contextualized embeddings produced by  $\Gamma$ ,  $R_{[CLS]}$  is the  $d$ -dimensional relevance embedding, and  $r$  is the relevance scalar produced by  $\Lambda$  to rank the candidates.

**Entity and candidate span modeling.** We represent the target entity  $E_{tgt}$  via two textual features, the title  $T_{tgt}$  and the lead text  $L_{tgt}$  which is a short paragraph present in most Wikipedia articles. Each candidate span  $x$  is represented via the text  $t$  contained in  $x$ . These textual features are concatenated together into a single textual input  $\phi$  as,

$$\phi = \mathcal{T}([\text{CLS}]T_{tgt}[\text{SEP}]L_{tgt}[\text{SEP}]t[\text{SEP}]) \quad (4)$$

where  $\mathcal{T}(\cdot)$  is the tokenizer operator that produces a sequence of  $T$  tokens.

**Optimization.** Given that entity insertion is a ranking task, we use an objective function that introduces the notion of ranking into the model. Specifically, we train the relevance scoring module using a cross-entropy loss over a *list* of candidates. Given a target entity  $E_{tgt}$ , a list of  $N$  candidate text spans  $\mathcal{X}_N = [x_1, \dots, x_N]$ , and  $i'$  as the index of the correct candidate, we use the following list-wise objective,

$$\max_{\theta} \frac{\exp(\text{score}(x_{i'}, E_{tgt}; \theta))}{\sum_{i=1}^N \exp(\text{score}(x_i, E_{tgt}; \theta))}$$

where  $\text{score}$  is an operator chaining the operations from Equations 2 and 3.

**Inference.** The document  $X_{src}$  in which to insert the entity  $E_{tgt}$  may contain a number  $D$  of potentially overlapping text spans  $\mathcal{X}_{src} = [x_1, \dots, x_D]$ . At inference time, the procedure described above is applied to all the  $D$  candidate text spans, and a relevance score is obtained for each candidate.

## 5.2 Two-stage training pipeline

We extract two types of links for studying entity insertion: existing and added links (§ 4). While added links reflect the entity insertion scenarios observed in the real world, we found that the number of added links is low for most languages (cf. Table 8 in the Appendix), thereby not being sufficient for training our model. To circumvent this challenge, we develop a *two-stage training pipeline* that uses both existing and added links.

**Dynamic context removal.** A key challenge with existing links is that they only reflect the `text_`-

Table 1: Dynamic context removal strategies.

Strategy	Text Removed
<code>rm_nth</code>	None
<code>rm_mention</code>	Mention
<code>rm_sent</code>	Sentence containing mention
<code>rm_span</code>	Span of sentences containing mention

present category of entity insertion, as the mention of the target entity is always present in the article containing the link. We mitigate this challenge by introducing a novel *data augmentation* strategy to simulate all other real-world entity insertion scenarios that are missing in the existing links. *Dynamic context removal* modifies the context associated with each existing link during training to simulate editors’ edits of adding links under different scenarios of entity insertion discussed in §4. Specifically, to simulate the `missing_mention`, `missing_sentence`, and `missing_span` scenarios, we randomly remove a word (`rm_mention`), a sentence (`rm_sent`), or a span of sentences (`rm_span`), respectively. Table 1 summarizes the strategies (cf. Table 10 in Appx. B.4 for details with examples).

Note that dynamic context removal may generate structural and linguistic patterns that would not occur in the text written by human editors. For example, applying the `rm_mention` strategy on the sentence “*Laika was a Soviet space dog who was one of the first animals in space to orbit the Earth.*”, would produce the sentence “*Laika was a who was one of the first animals in space to orbit the Earth.*”. Such a sentence is unlikely to be found in natural text articles, and thus, there is a distribution shift from the augmented training data to the test data.

**Expansion.** To reduce the impact of this shift, we introduce a second stage of training where we use the added links containing real-world entity insertion scenarios. Note that unlike the first stage, which uses existing links, the second stage does not require dynamic context removal, as we have access to the real contexts used by editors covering all the entity insertion scenarios.

## 5.3 Knowledge injection

While the representation presented in Eq. 4 (§ 5.1) already allows LOCEI to measure the target entity’s relevance to the candidate text span, we inject external knowledge about the target entity and knowledge about the structural organization of the source article to produce better relevance embeddings.

Since section titles provide additional ‘local’ knowledge in the form of a summarized conceptu-

alization of a candidate span, we first add the title of the section  $s$  in which a span appears to its input representation. Next, we add the list of mentions  $M_{tgt}$  previously associated with the target entity. This list provides a strong signal of how the entity is typically referenced in the text, thereby facilitating the model to better attend to these mentions when computing the relevance embedding. The final input format after knowledge injections is:

$$\phi = \mathcal{T}([\text{CLS}]T_{tgt}M_{tgt}[\text{SEP}]L_{tgt}[\text{SEP}]s[\text{SEP}]t[\text{SEP}])$$

#### 5.4 Incorporating multilinguality (xLOCEI)

To enable the encoder to better model the relationship between an entity target and candidate text spans, we leverage the patterns existent in multiple languages. For this, we train a single model by jointly considering entity insertion examples in multiple languages. This enables cross-lingual transfer, empowering, especially, low-resource languages with lesser and lower quality training data.

## 6 Experiments

All the resources required to reproduce the experiments in this paper are available at <https://github.com/epfl-dlab/multilingual-entity-insertion>.

### 6.1 Data

We study entity insertion in 105 language versions of Wikipedia. We use a judicious mix (based on size, script, geographic coverage, etc.) of 20 languages for training the benchmarked methods, however, for evaluation, we consider all 105 languages. For dataset statistics, cf. Tables 7 and 8 of Appx. B.

**Training set.** We train LOCEI and xLOCEI using a two-stage training pipeline (§ 5.2). While the data for the first stage is based on the existing links extracted from the 2023-10-01 snapshot, the second stage data is built using the links added between the 2023-09-01 and 2023-10-01 snapshots.

**Negative candidates.** During training, we extract  $N$  negative candidates for each positive candidate. Negative candidates are text spans in the source  $X_{src}$  where the target entity  $E_{tgt}$  was not inserted. Whenever possible, we select  $N$  negative candidates (“hard negatives”) from the same source article as the positive candidate. However, when articles are too small to be able to select  $N$  negatives, we sample the remaining negative candidates randomly from other articles (“easy negatives”). Details pertaining to the implementation of negative candidate extraction are provided in Appendix B.5.

**Test set.** For evaluation, we use the links added between the 2023-10-01 and 2023-11-01 snapshots. This ensures no overlap between the training and test sets and is therefore advantageous in mitigating data leakages. Unlike training, we use all the  $D$  available candidates in an article for evaluation.

### 6.2 Baselines

- **Random:** ranks candidates uniformly at random.
- **String Match:** searches for previously used mentions in the candidate text spans.
- **BM25 (Robertson and Zaragoza, 2009):** applies the Okapi-BM25 implementation (Trotman et al., 2014) on keywords extracted from the target lead paragraph and the candidate text spans.
- **EntQA (Zhang et al., 2022) (English only):** for independently encoding the candidate text spans and target entity. We then use the retriever model of EntQA to rank text spans based on the cosine similarity between the embeddings.
- **GET (Du et al., 2022) (English only):** use the generative ability of GET to generate the target entity name for each candidate text span. We then rank the text spans based on their likelihood of generating the target entity.
- **PRP-Allpair (Qin et al., 2024) (Zero-shot only):** to assess the relevance of candidate text spans to the target entity in a pairwise manner using GPT-3.5 and GPT-4, and then uncover the ranking from all pairwise comparisons.

### 6.3 Setup

**Model.** We present results for LOCEI and xLOCEI by fine-tuning the pre-trained xlm-roberta-base model (Conneau et al., 2020) as the encoder. The MLP is a 2-layer network with ReLU activations. We also explored different model sizes (e.g. Large and XL) and other pre-trained models (BERT and T5): results in Appendix C.

**Evaluation metrics.** We use (1) Hits@1, and (2) mean reciprocal rank (MRR) to evaluate the quality of the benchmarked methods. For each language, we compute the micro aggregates of the metrics over all added links in the test set. Moreover, we present results grouped into three categories: (1) Overall: considering the entire test set, (2) Present: considering links corresponding to the text\_present entity insertion scenario, and (3) Missing: considering links corresponding to all the other scenarios, namely, missing\_mention, missing\_sentence, and missing\_span.

Table 2: Entity insertion performance obtained by macro-averaging over 20 Wikipedia language versions used for training the benchmarked methods. xLOCEI trains a single model jointly on all 20 languages, whereas other methods train a separate model for each language. The categorization of entity insertion types into ‘Overall’, ‘Missing’, and ‘Present’ is discussed in § 6.3. Note that EntQA and GET work only for English (results in Table 3), whereas PRP-Allpair was only used for zero-shot analysis (Table 4) and English (Table 3).

	Method	Hits@1			MRR		
		Overall	Present	Missing	Overall	Present	Missing
Baseline	Random	0.107	0.115	0.103	0.243	0.259	0.236
Baseline	String Match	0.459	0.708	0.270	0.557	0.774	0.395
Baseline	BM25	0.508	0.799	0.280	0.612	0.866	0.421
Baseline	Simple fine-tuning	0.584	0.883	0.350	0.649	0.907	0.451
Proposed	LOCEI	0.672	0.877	0.509	0.744	0.906	0.617
Proposed	xLOCEI	<b>0.726</b> <sup>†</sup>	<b>0.909</b> <sup>†</sup>	<b>0.579</b> <sup>†</sup>	<b>0.789</b> <sup>†</sup>	<b>0.929</b> <sup>†</sup>	<b>0.678</b> <sup>†</sup>

<sup>†</sup> Indicates statistical significance ( $p < 0.05$ ) between the best and the second-best scores.

Table 3: Entity insertion performance obtained for English.

	Method	Hits@1			MRR		
		Overall	Present	Missing	Overall	Present	Missing
Baseline	Random	0.079	0.110	0.067	0.202	0.240	0.187
Baseline	String Match	0.391	0.732	0.264	0.489	0.796	0.374
Baseline	BM25	0.439	0.838	0.290	0.538	0.894	0.404
Baseline	EntQA <sub>RET</sub>	0.099	0.136	0.085	0.234	0.278	0.217
Baseline	GET	0.391	0.827	0.228	0.469	0.851	0.326
Baseline	PRP-Allpair (GPT-3.5) (Qin et al., 2024) *	0.160	0.375	0.092	0.322	0.536	0.255
Baseline	PRP-Allpair (GPT-4) (Qin et al., 2024) *	0.370	0.833	0.224	0.499	0.877	0.380
Baseline	Simple fine-tuning	0.443	0.860	0.287	0.522	0.888	0.385
Proposed	LOCEI	<b>0.677</b> <sup>†</sup>	<b>0.879</b>	<b>0.602</b> <sup>†</sup>	<b>0.741</b> <sup>†</sup>	<b>0.902</b>	<b>0.681</b> <sup>†</sup>

<sup>†</sup> Indicates statistical significance ( $p < 0.05$ ) between the best and the second-best scores.

\* Evaluation on a sample of 100 test instances.

Additional details about the experimental setup and hyperparameter tuning (impact of pre-trained models, model sizes, training stages, pointwise vs. ranking loss, etc.) are present in Appendix C.

## 6.4 Main results

We evaluate three variants of our entity insertion model: i) *simple fine-tuning*: a family of monolingual models fine-tuned in each language without the extensions (data augmentation, knowledge injection, two-stage training) introduced in LOCEI; ii) LOCEI: a family of monolingual models fine-tuned using the full LOCEI framework; and iii) xLOCEI, a single multilingual model fine-tuned jointly on all the languages using the full LOCEI framework. Table 2 shows the models’ performance metrics (Hits@1 and MRR) aggregated (macro-average) over the 20 considered languages.

**Overall performance.** We see that xLOCEI achieves the best overall quality and statistically significantly outperforms all other models for all cases considered. The key highlights are as follows: (1) *BM25*, a hard-to-beat baseline for ranking tasks, is around 20 percentage points inferior to xLOCEI, (2) *simple fine-tuning*, a baseline that we introduce in this work, substantially outperforms all the other considered methods, but is inferior to LOCEI and xLOCEI by being about 10 and

15 percentage points worse, respectively, and (3) xLOCEI consistently yields better scores than the language-specific LOCEI models, demonstrating that the multilingual model is capable of transferring knowledge across languages to improve overall performance. In fact, by looking at the performance for the individual languages in Figs. 6 and 7 (Appx. C.2), we see that the improvement from xLOCEI over LOCEI is larger in low-resource languages (languages with less training data) such as Afrikaans (af), Welsh (cy), Uzbek (uz).

**Performance on ‘Missing’ and ‘Present’ categories.** The key finding is that the baselines lack robustness to the variation in entity insertion types, which is substantiated by the huge disparity of entity insertion performance (around 50 percentage points) of all the baselines in the ‘Present’ and ‘Missing’ categories. This result further highlights the key limitation of the baselines: they cannot address the challenging scenarios of entity insertion. The key reason behind this disparity is that all the existing baselines rely on the existence of a suitable text span to insert a link to the target entity. On the contrary, both LOCEI and xLOCEI effectively utilize the signals manifested in the context due to the introduced extensions (e.g. data augmentation) and are therefore robust to different entity insertion scenarios. Consequently, we observe that both LOCEI

Table 4: Entity insertion performance in the zero-shot setting: results obtained by macro-averaging over 9 Wikipedia language versions that were not used for fine-tuning xLOCEI<sub>11</sub>. xLOCEI<sub>20</sub> was trained jointly on all 20 languages, whereas LOCEI trains a separate model for each language. The categorization of entity insertion types into ‘Overall’, ‘Missing’, and ‘Present’ is discussed in § 6.3.

	Method	Hits@1			MRR		
		Overall	Present	Missing	Overall	Present	Missing
Fine-tuned	LOCEI	0.647	0.873	0.486	0.718	0.902	0.588
Fine-Tuned	xLOCEI <sub>20</sub>	<b>0.709</b> <sup>†</sup>	<b>0.901</b>	<b>0.570</b> <sup>†</sup>	<b>0.772</b> <sup>†</sup>	<b>0.923</b>	<b>0.662</b> <sup>†</sup>
Zero-shot	PRP-Allpair (GPT-3.5) (Qin et al., 2024) *	0.289	0.423	0.210	0.433	0.563	0.353
Zero-shot	PRP-Allpair (GPT-4) (Qin et al., 2024) *	0.571	0.859	0.344	0.656	0.897	0.468
Zero-Shot	xLOCEI <sub>11</sub>	<b>0.690</b> <sup>†</sup>	<b>0.887</b>	<b>0.541</b> <sup>†</sup>	<b>0.755</b> <sup>†</sup>	<b>0.913</b>	<b>0.636</b> <sup>†</sup>

<sup>†</sup> Indicates statistical significance ( $p < 0.05$ ) from fine-tuned LOCEI.

\* Evaluation on a sample of 100 test instances.

and xLOCEI obtain substantial improvements over all the baseline models in the missing category.

**Performance on English.** Table 3 shows that even in English (a high-resource language), xLOCEI outperforms all baselines. Once again, this gap is pronounced in the missing case, further highlighting the difficulty and novelty of the task.

### Zero-shot vs. Fine-tuned

We further study the performance of xLOCEI in the zero-shot scenario, i.e., evaluating the model in languages that were not explicitly contained in the data for fine-tuning. This is relevant to assess the potential to support languages for which there is little or no training data available. We consider xLOCEI<sub>11</sub>, a variant of the multilingual xLOCEI which is trained on only 11 out of the 20 languages (cf. Table 11 in Appx. C.3 for details). We then evaluate the zero-shot performance of xLOCEI<sub>11</sub> in the remaining 9 languages not considered for training. For comparison, we also show the non-zero shot performance of the models considered in the previous subsection: i) LOCEI, the family of monolingual models fine-tuned in each language; and ii) xLOCEI<sub>20</sub>, the single multilingual model trained on all 20 languages. The main result, shown in (Table 4), is that xLOCEI<sub>11</sub> retains over 95% performance in the zero-shot scenario in comparison to the results of the best model, xLOCEI<sub>20</sub>, which was fine-tuned on these languages. Nevertheless, xLOCEI<sub>11</sub> still outperforms the language-specific

Table 5: Entity insertion performance in the full zero-shot setting: results obtained by macro-averaging over 85 held-out Wikipedia language versions that were not used for fine-tuning the benchmarked methods.

Method	Hits@1			MRR		
	Overall	Present	Missing	Overall	Present	Missing
Random	0.148	0.132	0.148	0.288	0.287	0.281
String Match	0.442	0.717	0.273	0.549	0.786	0.406
BM25	0.456	0.733	0.294	0.580	0.823	0.435
xLOCEI <sub>11</sub>	0.683	0.853	0.585	0.754	0.886	0.676
xLOCEI <sub>20</sub>	<b>0.706</b> <sup>†</sup>	<b>0.873</b> <sup>†</sup>	<b>0.602</b>	<b>0.769</b>	<b>0.901</b>	<b>0.685</b>

<sup>†</sup> Indicates statistical significance ( $p < 0.05$ ) between the best and the second-best scores.

LOCEI models fine-tuned on each language individually. We expand the robustness of these results by considering two additional scenarios.

First, we compare the performance of xLOCEI<sub>11</sub> with PRP-Allpair, the state-of-the-art framework for ranking tasks using LLMs (Table 4). We find that xLOCEI<sub>11</sub> substantially outperforms PRP-Allpair, both when using GPT-3.5 and GPT-4, particularly for the cases when the mention that is linked is not yet present in the text (missing).

Second, we evaluate our models on held-out data of the remaining 85 languages in Table 5. We reproduce a high zero-shot performance of xLOCEI<sub>11</sub>, in comparison to results in the 9 languages considered in Table 4. In comparison, other baseline models yield substantially lower performance.

Overall, these findings show that xLOCEI is capable of transferring the knowledge acquired during fine-tuning to unseen languages while maintaining a similar level of performance. This demonstrates that our entity insertion model can be scaled to many languages even if little or no additional training data is available for those languages.

### 6.5 Ablation analysis

Finally, we investigate in more detail the effect of the extensions, namely, data augmentation, knowledge injection, and two-stage training that we introduce in the training pipeline of our model in comparison to a standard fine-tuning approach. Table 6 portrays the improvement in performance on account for each extension introduced in this work.

Overall, we see that each extension has an overall positive impact on performance. First, introducing the dynamic context removal for data augmentation is only effective when including negative examples. In that case, it improves the performance on the missing cases, but at the cost of performance in the present case. This is expected because context removal leads to the model seeing fewer training samples in the present case. Sec-



Table 6: Analyzing the impact of the extensions introduced in the LOCEI framework on the entity insertion performance for only English and the macro-average over 20 Wikipedia language versions. The categorization of entity insertion types into ‘Overall’, ‘Missing’, and ‘Present’ is discussed in § 6.3.

Model Variant	English						All 20 Languages					
	Hits@1			MRR			Hits@1			MRR		
	Overall	Present	Missing	Overall	Present	Missing	Overall	Present	Missing	Overall	Present	Missing
simple fine-tuning	0.443	0.860	0.287	0.522	0.888	0.385	0.584	<b>0.883</b>	0.350	0.649	<b>0.907</b>	0.451
+dynamic ctxt removal (w/o neg)	0.440	0.805	0.304	0.532	0.842	0.415	0.541	0.782	0.372	0.626	0.828	0.487 <sup>†</sup>
+dynamic ctxt removal	0.473	0.846	0.334	0.547	0.875	0.424	0.574 <sup>†</sup>	0.838 <sup>†</sup>	0.376	0.649 <sup>†</sup>	0.873 <sup>†</sup>	0.486 <sup>†</sup>
+expansion	0.648 <sup>†</sup>	0.875	0.563 <sup>†</sup>	0.719 <sup>†</sup>	<b>0.902</b>	0.651 <sup>†</sup>	0.657 <sup>†</sup>	0.850	0.500 <sup>†</sup>	0.733 <sup>†</sup>	0.889	0.609 <sup>†</sup>
+knowledge injection	<b>0.677</b>	<b>0.879</b>	<b>0.602</b>	<b>0.741</b>	<b>0.902</b>	<b>0.681</b>	<b>0.672</b>	0.877 <sup>†</sup>	<b>0.509</b>	<b>0.744</b>	0.906	<b>0.617</b>

<sup>†</sup> Indicates statistical significance ( $p < 0.05$ ) between the variant and the previous variant.

ond, introducing expansion as a second stage in the training led to a large boost in performance in all scenarios, showing the benefit of using the smaller but high-quality dataset of added links for the training. Third, the knowledge injection further improved the performance in both scenarios, indicating that the additional knowledge helps the model produce better relevance embeddings.

## 7 Discussions

### 7.1 Summary of findings

We introduced the novel task of entity insertion in information networks. Considering the case of Wikipedia, we justified the relevance and need for solving this task by demonstrating empirically that existing methods such as entity linking are often not suitable in practice. In fact, we showed that in 65% of edits in which links were inserted by editors, none of the existing text is suitable to insert the entity, i.e. new text has to be inserted *somewhere in the article* along with the inserted entity.

We developed a multilingual model (xLOCEI) to effectively solve the entity insertion task across 20 Wikipedia languages outperforming all other models. First, our model substantially outperforms strong baseline approaches based on string matching or BM25, especially in the case when the linked mention was missing. We demonstrate how each of the introduced novelties (data augmentation, knowledge injection, two-stage training pipeline) contribute to improve the downstream performance. Second, the multilingual model yields consistently better results than language-specific models. This shows that our model is capable of collating the knowledge acquired from each language to improve performance over all languages. Third, our model works well in a zero-shot scenario, i.e. not only retaining over 95% of the hypothetical best performance if the language was included but even outperforming the much larger GPT-3.5 and GPT-4. This demonstrates that the model is capable of transferring knowledge to languages unseen during fine-tuning which is crucial for the practical

application across the more than 300 languages in Wikipedia, for which often there is little or no training data available. We compiled a new benchmark dataset for entity insertion in Wikipedia covering 105 languages. We make the dataset publicly available to enable future research in entity insertion.

### 7.2 Implications and broader impact

**A new benchmark for NLP tasks.** The problems of link recommendations and entity linking have been well-studied and many excellent solutions have been brought forward, some of which are denoted even near-optimal (Ghasemian et al., 2020). The problem of entity insertion constitutes a new relevant and challenging task in the domain of NLP. Our multilingual dataset provides a resource for researchers for development and evaluation of new models to solve this task. This will help improve the overall capabilities of large language models when applied in the context of networks that are crucial for organizing textual information.

#### Supporting editors to bridge knowledge gaps.

Many articles in Wikipedia lack visibility in the hyperlink network capturing a specific aspect of the general problem of knowledge gaps (Redi et al., 2020). For example, there are more than 8.8M so-called orphan articles (Arora et al., 2024), i.e., articles without any incoming links, which are de-facto invisible to readers navigating Wikipedia. Even if suitable link targets are identified, a remaining challenge for editors is to identify a relevant position in the text where the link can be inserted. At the current rate of “de-orphanization”, it would take editors more than 20 years to work through the backlog of orphan articles, suggesting that existing tools do not support editors in addressing this issue effectively. Our model can support editors in this task, complementing existing approaches based on entity linking such as the add-a-link tool for newcomer editors (Gerlach et al., 2021).

## Limitations

We tried different pre-trained language models for our experiments with RoBERTa outperforming BERT and T5 by a large margin. The use of larger models with more parameters could further improve performance. While differences between RoBERTa-base and -large in English were marginal, we noticed a substantial drop when using the multilingual XLM-RoBERTa instead RoBERTa. This suggests that larger model architectures could be especially beneficial in the multilingual setting in order to improve support for low-resource languages, where performance is typically lower in comparison (Wu and Dredze, 2020). While multilingual models based on transformer architectures support many languages (e.g., XLM-RoBERTa was pre-trained on 100 languages), many of the more than 300 languages in Wikipedia are still not explicitly represented in the training data of these models. Thus, if unaddressed, the use of such models could lead to a language gap constituting a substantial barrier towards knowledge equity (Redi et al., 2020).

One practical limitation of the model is that the ranking of all text spans can become expensive if the article is very long and, thus, contains many candidates. This constitutes challenge for deploying the model in the future as a ready-to-use-tools for editors in practice. This requires the integration of potential solutions for improving inference such as via hierarchical searching.

Further improvements to the model could come from integrating of additional information from the local Wikipedia graph structure or the candidate context. For example, a very strong signal are the links already existing in the candidate context, as these indicate entities related to the context. Providing these as additional features to the model might help generate better representations of the candidate (Arora et al., 2022) and, as a result, better relevance embeddings. Furthermore, one could take advantage of the multilingual nature of Wikipedia with more than 300 language versions, each having a surprising amount of information not contained in any other languages (Bao et al., 2012). Thus, existing content about a target entity from other languages could provide relevant context (García-Durán et al., 2022), which could be made available through automatic translation, such as the already available section translation tool in Wikipedia (WMF, 2019).

In our operationalization of entity insertion, we assume that the link to be inserted consisting of the pair of the source- and target entity is known. This assumption holds in the specific use-case of article “de-orphanization” (Arora et al., 2024) serving as the motivation for formulating the task of entity insertion. However, when this is not the case, our model requires an additional step to generate a specific link, e.g., via existing link recommendation models.

Our modeling framework is not suitable for the scenario where links are added in a section that did not exist in the previous version of the article (`missing_section`). The text from the surrounding sections are not a good indicator for the insertion of a new entity, because they typically cover different subjects. The `missing_section` scenario could be addressed through complementary approaches based on generative models that produce a draft for new section when none of the existing candidates leads to a high relevance score.

## Ethics statement

We have assessed the ethics of conducting this research and the output from the research and we have not identified ethical risks associated with the research. All the datasets and resources used in this work are publicly available and do not contain private or sensitive information about Wikipedia readers. While the dataset comes from individual edits (e.g., links added to Wikipedia articles), it does not contain any details about the author(s) of those edits. All the findings are based on analyses conducted at an aggregate level, and thus, no individual-level inferences can be drawn. And lastly, the experiments are done on data already collected and no human subject has been involved as part of them. We confirm that we have read and abide by the ACL code of conduct.

## Acknowledgements

We thank Leila Zia, Alberto García-Durán, and Marija Šakota for insightful discussions and for reviewing an initial draft of this paper. West’s lab is partly supported by grants from the Swiss National Science Foundation (200021\_185043 and 211379), Swiss Data Science Center (P22\_08), H2020 (952215), Microsoft Swiss JRC, and Google. We also gratefully acknowledge generous gifts from Facebook, Google, and Microsoft.

## References

- Akhil Arora, Alberto Garcia-Duran, and Robert West. 2021. [Low-Rank Subspaces for Unsupervised Entity Linking](#). In *EMNLP*, pages 8037–8054.
- Akhil Arora, Martin Gerlach, Tiziano Piccardi, Alberto García-Durán, and Robert West. 2022. [Wikipedia reader navigation: When synthetic data is enough](#). In *WSDM*, page 16–26.
- Akhil Arora, Robert West, and Martin Gerlach. 2024. [Orphan articles: The dark matter of wikipedia](#). In *ICWSM*, pages 100–112.
- Patti Bao, Brent Hecht, Samuel Carton, Mahmood Quaderi, Michael Horn, and Darren Gergle. 2012. [Omnipedia: bridging the wikipedia language gap](#). In *SIGCHI*, pages 1075–1084.
- Alon Brutzkus and Amir Globerson. 2019. [Why do Larger Models Generalize Better? A Theoretical Perspective via the XOR Problem](#). In *ICML*, pages 822–830.
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. [ELECTRA: pre-training text encoders as discriminators rather than generators](#). In *ICLR*.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised Cross-lingual Representation Learning at Scale](#). In *ACL*, pages 8440–8451.
- Marko Čuljak, Andreas Spitz, Robert West, and Akhil Arora. 2022. [Strong Heuristics for Named Entity Linking](#). In *NAACL-SRW*, pages 235–246.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#). In *NAACL*, pages 4171–4186.
- Qian Dong, Yiding Liu, Suqi Cheng, Shuaiqiang Wang, Zhicong Cheng, Shuzi Niu, and Dawei Yin. 2022. [Incorporating explicit knowledge in pre-trained language models for passage re-ranking](#). In *SIGIR*, pages 1490–1501.
- Christina Du, Kashyap Papat, Louis Martin, and Fabio Petroni. 2022. [Entity Tagging: Extracting Entities in Text Without Mention Supervision](#). *CoRR*, abs/2209.06148.
- Jinyuan Fang, Zaiqiao Meng, and Craig Macdonald. 2023. [KGPR: Knowledge Graph Enhanced Passage Ranking](#). In *CIKM*, page 3880–3885.
- Jenny Rose Finkel, Trond Grenager, and Christopher D. Manning. 2005. [Incorporating non-local information into information extraction systems by gibbs sampling](#). In *ACL*, pages 363–370.
- Xingyu Fu, Weijia Shi, Xiaodong Yu, Zian Zhao, and Dan Roth. 2020. [Design challenges in low-resource cross-lingual entity linking](#). In *EMNLP*, pages 6418–6432.
- Luyu Gao, Zhuyun Dai, and Jamie Callan. 2021. [Re-think Training of BERT Rerankers in Multi-stage Retrieval Pipeline](#). In *ECIR*, pages 280–286.
- Alberto García-Durán, Akhil Arora, and Robert West. 2022. [Efficient Entity Candidate Generation for Low-Resource Languages](#). In *LREC*, pages 6429–6438.
- Martin Gerlach, Marshall Miller, Rita Ho, Kosta Harlan, and Djellel Eddine Difallah. 2021. [Multilingual Entity Linking System for Wikipedia with a Machine-in-the-Loop Approach](#). In *CIKM*, pages 3818–3827.
- Amir Ghasemian, Homa Hosseinmardi, Aram Galstyan, Edoardo M Airoidi, and Aaron Clauset. 2020. [Stacking models for nearly optimal link prediction in complex networks](#). *PNAS*, 117(38):23393–23400.
- Suchin Gururangan, Ana Marasovic, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. [Don’t Stop Pretraining: Adapt Language Models to Domains and Tasks](#). In *ACL*, pages 8342–8360.
- Shuguang Han, Xuanhui Wang, Mike Bendersky, and Marc Najork. 2020. [Learning-to-Rank with BERT in TF-Ranking](#). *CoRR*, abs/2004.08476.
- Johannes Hoffart, Mohamed Amir Yosef, Iliaria Bordino, Hagen Fürstenau, Manfred Pinkal, Marc Spaniol, Bilyana Taneva, Stefan Thater, and Gerhard Weikum. 2011. [Robust Disambiguation of Named Entities in Text](#). In *EMNLP*, pages 782–792.
- Darren Wei Che Huang, Shlomo Geva, and Andrew Trotman. 2008. [Overview of the INEX 2008 Link the Wiki Track](#). In *INEX Workshop*, volume 5631, pages 314–325.
- Jia-Huei Ju, Jheng-Hong Yang, and Chuan-Ju Wang. 2021. [Text-to-Text Multi-view Learning for Passage Re-ranking](#). In *SIGIR*, pages 1803–1807.
- Omar Khattab and Matei Zaharia. 2020. [ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT](#). In *SIGIR*, pages 39–48.
- Jimmy Lin, Rodrigo Frassetto Nogueira, and Andrew Yates. 2021. [Pretrained Transformers for Text Ranking: BERT and Beyond](#). *Synthesis Lectures on Human Language Technologies*. Morgan & Claypool Publishers.
- Nelson F. Liu, Kenton Lee, and Kristina Toutanova. 2023. [Anchor Prediction: Automatic Refinement of Internet Links](#). *CoRR*, abs/2305.14337.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [RoBERTa: A Robustly Optimized BERT Pretraining Approach](#). *CoRR*, abs/1907.11692.

- David N. Milne and Ian H. Witten. 2008. [Learning to Link with Wikipedia](#). In *CIKM*, pages 509–518.
- Louis Martin Benjamin Müller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric de la Clergerie, Djamé Seddah, and Benoît Sagot. 2020. [CamemBERT: a Tasty French Language Model](#). In *ACL*, pages 7203–7219.
- Rodrigo Frassetto Nogueira and Kyunghyun Cho. 2019. [Passage re-ranking with BERT](#). *CoRR*, abs/1901.04085.
- Rodrigo Frassetto Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy Lin. 2020. [Document Ranking with a Pretrained Sequence-to-Sequence Model](#). In *EMNLP (Findings)*, pages 708–718.
- Rodrigo Frassetto Nogueira, Wei Yang, Kyunghyun Cho, and Jimmy Lin. 2019. [Multi-Stage Document Ranking with BERT](#). *CoRR*, abs/1910.14424.
- Joel Nothman, Nicky Ringland, Will Radford, Tara Murphy, and James R. Curran. 2013. [Learning multilingual named entity recognition from Wikipedia](#). *Artif. Intell.*, 194:151–175.
- Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, Xuanhui Wang, and Michael Bendersky. 2024. [Large language models are effective text rankers with pairwise ranking prompting](#). In *NAACL (Findings)*, pages 1504–1518.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer](#). *J. Mach. Learn. Res.*, 21:140:1–140:67.
- Miriam Redi, Martin Gerlach, Isaac Johnson, Jonathan Morgan, and Leila Zia. 2020. [A Taxonomy of Knowledge Gaps for Wikimedia Projects \(Second Draft\)](#).
- Stephen E. Robertson and Hugo Zaragoza. 2009. [The Probabilistic Relevance Framework: BM25 and Beyond](#). *Found. Trends Inf. Retr.*, 3(4):333–389.
- Piotr Rybak, Robert Mroczkowski, Janusz Tracz, and Ireneusz Gawlik. 2020. [KLEJ: Comprehensive Benchmark for Polish Language Understanding](#). In *ACL*, pages 1191–1201.
- Yijia Shao, Yucheng Jiang, Theodore Kanell, Peter Xu, Omar Khattab, and Monica Lam. 2024. [Assisting in writing Wikipedia-like articles from scratch with large language models](#). In *NAACL*, pages 6252–6278.
- James B Simon, Dhruva Karkada, Nikhil Ghosh, and Mikhail Belkin. 2024. [More is better: when infinite overparameterization is optimal and overfitting is obligatory](#). In *ICLR*.
- Mahdi Soltanolkotabi, Adel Javanmard, and Jason D. Lee. 2019. [Theoretical Insights Into the Optimization Landscape of Over-Parameterized Shallow Neural Networks](#). *IEEE Trans. Inf. Theory*, 65(2):742–769.
- Andrew Trotman, Antti Puurula, and Blake Burgess. 2014. [Improvements to BM25 and Language Models Examined](#). In *ADCS*, page 58.
- Johannes M. van Hulst, Faegheh Hasibi, Koen Dercksen, Krisztian Balog, and Arjen P. de Vries. 2020. [REL: An Entity Linker Standing on the Shoulders of Giants](#). In *SIGIR*, pages 2197–2200.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is All you Need](#). In *NeurIPS*, pages 5998–6008.
- Robert West, Doina Precup, and Joelle Pineau. 2009. [Completing Wikipedia’s Hyperlink Structure through Dimensionality Reduction](#). In *CIKM*, pages 1097–1106.
- Robert West, Doina Precup, and Joelle Pineau. 2010. [Automatically suggesting topics for augmenting text documents](#). In *CIKM*, page 929–938.
- WMF. 2010a. [Wikimedia downloads](#). <https://dumps.wikimedia.org/backup-index.html>. Accessed: 2023-05-01.
- WMF. 2010b. [Wikimedia enterprise html dumps](#). [https://dumps.wikimedia.org/other/enterprise\\_html/](https://dumps.wikimedia.org/other/enterprise_html/). Accessed: 2024-02-09.
- WMF. 2019. [Content translation: Section translation](#). [https://www.mediawiki.org/wiki/Content\\_translation/Section\\_translation](https://www.mediawiki.org/wiki/Content_translation/Section_translation). Accessed: 2024-04-09.
- Shijie Wu and Mark Dredze. 2020. [Are All Languages Created Equal in Multilingual BERT?](#) In *ReplANLP Workshop*, pages 120–130.
- Ellery Wulczyn, Robert West, Leila Zia, and Jure Leskovec. 2016. [Growing wikipedia across languages via recommendation](#). In *WWW*, pages 975–985.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. [mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer](#). In *NAACL*, pages 483–498.
- Wenzheng Zhang, Wenyue Hua, and Karl Stratos. 2022. [EntQA: Entity linking as question answering](#). In *ICLR*.
- Honglei Zhuang, Zhen Qin, Rolf Jagerman, Kai Hui, Ji Ma, Jing Lu, Jianmo Ni, Xuanhui Wang, and Michael Bendersky. 2023. [RankT5: Fine-Tuning T5 for Text Ranking with Ranking Losses](#). In *SIGIR*, pages 2308–2313.

## A Additional related work

### A.1 Pre-trained language models

The transformer architecture, introduced by (Vaswani et al., 2017), has become the *de facto* architecture for most Natural Language Processing (NLP) applications. A transformer-based pre-trained language model takes as input a text sequence and computes a vector embedding that captures the semantic and structural information contained in the text sequence, which can then be used in downstream applications.

Pre-training is an expensive process. For example, the base variant of BERT (Devlin et al., 2019) took four days to train with 16 TPUs and RoBERTa (Liu et al., 2019) took one day to train with 1024 GPUs. However, pre-trained models can be leveraged to novel downstream tasks by fine-tuning them on task-specific datasets. As a comparison, the authors of BERT (Devlin et al., 2019) introduced several fine-tuned variants of BERT, all of which were fine-tuned in one hour using one TPU, which is much cheaper than pre-training the model for each task. This paradigm of pre-training language models on large amounts of data and then fine-tuning on much smaller amounts can reduce the cost of model training while retaining the knowledge from the pre-trained model and transferring it to the downstream task. Popular pre-trained models for multilingual tasks are mBERT (Devlin et al., 2019), XLM-RoBERTa (Conneau et al., 2020), and mT5 (Xue et al., 2021).

### A.2 Ranking tasks

Since entity insertion is a ranking task, in this section, we provide a short review of literature focusing on document retrieval and ranking.

Classical approaches for ranking tasks, such as BM25 (Robertson and Zaragoza, 2009), mainly rely on probabilistic methods that attempt to match keywords between a query and a candidate document. However, these methods cannot capture complex semantic and structural patterns. For example, the sentences “The hero defeated the dragon and saved the damsel” and “The knight slayed the beast and rescued the princess” are semantically equivalent, but classical methods would not match them due to the small vocabulary overlap.

That said, pre-trained language models have become state-of-the-art for text ranking (Lin et al., 2021). A popular design for transformer-based ranking tasks is the cross-attention model, in which

the query and the candidate document are concatenated into a sequence and then processed by the model. Since transformer models employ attention mechanisms, this strategy allows the model to capture the interactions between the query and the document.

This approach has been explored for encoder-only models (Han et al., 2020; Nogueira et al., 2019; Gao et al., 2021), outperforming classical methods. There has also been previous research (Nogueira et al., 2020; Ju et al., 2021) in exploring encoder-decoder models, such as T5 (Raffel et al., 2020). However, even though encoder-decoder models are typically larger than encoder-only models, RankT5 (Zhuang et al., 2023) has shown that there is no consistent winner between encoder-decoder and encoder-only models.

Given its recent success in document retrieval, the training objective of LOCEI is inspired by the ranking loss proposed in RankT5 (Zhuang et al., 2023).

### A.3 Domain adaption

(Gururangan et al., 2020) have shown that a second phase of pre-training using domain-specific knowledge can improve the performance of language models. Their experiments started with a pre-trained RoBERTa model and continued pre-training it using unlabelled data from a large corpus of domain-specific text.

In our work, we propose a similar approach for fine-tuning, where we apply a first stage of domain-shifted data and then a second stage of domain-specific data to improve the performance further.

## B Additional dataset processing details

### B.1 Data preparation steps

**Existing links.** For the existing links, we store the following data: source and target titles, Wikidata QIDs, lead paragraphs, the name of the section containing the link, and a context surrounding the link. The context is defined as the sentence containing the link and the five sentences before and after (or until we reach the end of the section). We additionally keep positional information about the mention and the sentence containing the mention relative to the context (i.e., the start and end indices of the mention and the sentence in the context). The positional information is relevant to the data augmentation strategy we introduced (see § B.4).

**Added links.** For the added links, we store the same information as in the existing links, except for the positional information. This is because positional information is required primarily for performing data augmentations, which are required only for processing existing links.

## B.2 Dataset statistics

Table 7 shows the summary statistics of the entity insertion dataset for each of the 105 considered language versions of Wikipedia, in particular the number of articles, the number of existing links, and the number of added links.

Table 8 shows the number of samples contained in the training and test splits, respectively, for each of the 20 Wikipedia language versions considered in the experiments.

## B.3 Entity insertion categories

Table 9 shows an example for each of the entity insertion categories, except for the category `missing_section`, demonstrating that the problem of entity insertion grows in complexity as more text is missing.

Additionally, Fig. 5 shows the distribution of entity insertion categories for 20 Wikipedia language versions considered in the experiments.

## B.4 Dynamic context removal

Table 10 shows examples of the different types of dynamic context removal. Specifically, we randomly remove a word (`rm_mention` simulation), a sentence (`rm_sent` simulation), or a span of sentences (`rm_span` simulation) during training. Before sending the input to the model, we randomly select one of the masking strategies mentioned above (as well as no masking) to modify the input accordingly. However, before applying the strategy, we verify if the selected strategy does not produce an empty input. This may happen when, for example, the context is a single sentence, in which case simulating the `rm_sent` strategy would lead to an empty input. If the sampled strategy would produce an empty input, we re-sample a less aggressive strategy.

While performing the `rm_span` simulation, the number of sentences to remove is chosen randomly between 2 and 5. Note that we used a space-based splitting for ease of implementation, and we acknowledge that this could be an issue for certain languages, such as Japanese or Chinese, which we intend to fix in the future.

## B.5 Rules for sampling negative candidates

We employ the following rules when constructing the negative candidates, both for training and validation.

1. A candidate’s context should not span over two different sections.
2. A candidate’s context should not contain any of the mentions previously used to link to the target entity.

The first rule keeps the content of each context consistent, as two distinct sections can cover very different topics. The second rule ensures that all the candidates used to evaluate the module are correctly classified as either positive candidates or negative candidates. For example, if the goal is to insert the entity “1984” (the book - Q208460) and there is a sentence in the article with the word “1984” not linked to the target article, there could be three reasons for the link to be missing. First, the mention “1984” could be related to a different entity (e.g., the year - Q2432), in which case the sentence should belong to a negative candidate. Second, the mention is supposed to be for the target entity but it is not yet linked, in which case the sentence should belong to an additional positive candidate. Finally, the mention is supposed to be for the target entity but it should not be linked because of Wikipedia’s editing guidelines, in which case it is not clear whether the sentence should belong to a negative or a positive candidate. Due to this unclear categorization, we choose to remove any sentences containing mentions previously associated with the target entity to be inserted.

## C Additional experiments

### C.1 Hyperparameters

We train the encoder and MLP with learning rates of  $1e-5$  and  $1e-4$ , respectively, using  $N = 9$  negative candidates. Moreover, we use 5 sentences on either side as context for each candidate text span and set  $|M_{tgt}| = 10$ . The first stage of training uses 20K data points and is trained for 4 epochs, whereas the second stage uses all the available data for 2 epochs. Mimicking the real-world entity insertion scenarios, we set `rm_nth=40%`, `rm_mention=20%`, `rm_sentence=30%`, and `rm_span=10%`.

Table 7: Summary statistics of the full entity insertion dataset collected from 105 different Wikipedia language versions.

	Language	Articles	Existing Links	Added Links		Language	Articles	Existing Links	Added Links
en	English	6.7M	166M	368K	de	German	2.8M	78.3M	94.3K
sv	Swedish	2.5M	29.9M	10.7K	fr	French	2.5M	85.1M	64.5K
nl	Dutch	2.1M	24.7M	23.6K	ru	Russian	1.9M	47.6M	33.8K
es	Spanish	1.8M	47.9M	66.3K	it	Italian	1.7M	51.1M	45.6K
pl	Polish	1.5M	30.1M	27.2K	ja	Japanese	1.3M	60.6M	79.0K
zh	Chinese	1.3M	23.1M	28.2K	vi	Vietnamese	1.2M	10.3M	11.9K
ar	Arabic	1.2M	16.3M	17.8K	pt	Portuguese	1.1M	21.9M	24.2K
fa	Persian	971K	9.5M	18.1K	ca	Catalan	732K	14.6M	18.4K
sr	Serbian	671K	8.3M	5.4K	id	Indonesian	650K	8.5M	13.7K
ko	Korean	634K	11.2M	21.3K	no	Norwegian	611K	11.3M	7.2K
ce	Chechen	599K	3.0M	48	fi	Finnish	554K	9.7M	13.7K
cs	Czech	531K	14.4M	12.3K	tr	Turkish	531K	6.7M	14.9K
hu	Hungarian	527K	10.6M	7.8K	tt	Tatar	496K	3.1M	94
sh	Serbo-Croatian	456K	8.3M	807	ro	Romanian	439K	6.9M	4.2K
eu	Basque	412K	4.4M	5.1K	ms	Malay	363K	2.9M	2.7K
he	Hebrew	341K	14.7M	36.7K	eo	Esperanto	340K	6.7M	5.8K
hy	Armenian	296K	4.5M	3.7K	da	Danish	294K	5.7M	2.3K
bg	Bulgarian	288K	5.2M	4.8K	cy	Welsh	270K	2.6M	386
sk	Slovak	242K	3.4M	3.1K	azb	South Azerbaijani	242K	1.0M	22
simple	Simple English	240K	2.6M	3.8K	et	Estonian	235K	4.4M	4.7K
kk	Kazakh	233K	1.6M	1.5K	be	Belarusian	232K	3.1M	3.0K
uz	Uzbek	230K	1.3M	4.3K	min	Minangkabau	226K	644K	21
el	Greek	224K	4.8M	6.7K	lt	Lithuanian	210K	3.8M	2.4K
gl	Galician	196K	3.9M	4.2K	hr	Croatian	194K	3.2M	3.4K
ur	Urdu	190K	1.4M	5.2K	az	Azerbaijani	188K	2.4M	6.3K
sl	Slovenian	182K	3.1M	1.5K	ka	Georgian	163K	2.3M	1.7K
ta	Tamil	157K	1.6M	1.1K	hi	Hindi	157K	1.2M	2.3K
la	Latin	138K	2.5M	1.2K	mk	Macedonian	136K	2.3M	923
ast	Asturian	128K	2.6M	49	lv	Latvian	121K	2.0M	1.9K
af	Afrikaans	111K	1.3M	1.4K	tg	Tajik	108K	567K	123
sq	Albanian	97.3K	873K	523	mg	Malagasy	95.7K	495K	1.2K
bs	Bosnian	89.7K	1.6M	969	oc	Occitan	88.3K	1.1M	1.9K
te	Telugu	82.2K	934K	1.1K	sw	Swahili	74.3K	1.0M	558
lmo	Lombard	71.9K	380K	26	lv	Javanese	70.5K	513K	161
ba	Bashkir	62.4K	960K	649	lb	Luxembourgish	61.7K	930K	754
mr	Marathi	60.9K	409K	67	su	Sundanese	60.3K	470K	6
is	Icelandic	56.4K	725K	1.0K	ga	Irish	56.0K	387K	204
ku	Kurdish	54.3K	252K	614	fy	Western Frisian	51.0K	1.3M	579
pa	Punjabi	49.4K	282K	139	cv	Chuvash	48.3K	213K	304
br	Breton	46.5K	326K	852	tl	Tagalog	43.2K	435K	512
an	Aragonese	40.8K	620K	70	io	Ido	40.7K	422K	230
sco	Scots	35.5K	251K	40	vo	Volapük	34.6K	134K	7
ne	Nepali	32.1K	168K	250	ha	Hausa	30.6K	129K	262
gu	Gujarati	30.2K	411K	29	kn	Kannada	28.0K	253K	514
bar	Bavarian	27.0K	207K	21	scn	Sicilian	23.8K	132K	5
mn	Mongolian	22.5K	187K	467	si	Sinhala	20.3K	81.7K	36
ps	Pashto	16.2K	49.7K	10	gd	Scottish Gaelic	15.8K	207K	14
yi	Yiddish	15.2K	185K	21	sd	Sindhi	13.4K	49.5K	14
am	Amharic	12.9K	69.1K	12	as	Assamese	11.9K	104K	459
sa	Sanskrit	10.5K	65.2K	18	km	Khmer	9.8K	52.3K	95
ary	Moroccan Arabic	8.0K	50.5K	129	so	Somali	7.4K	64.2K	60
ug	Uyghur	5.9K	9.7K	1	lo	Lao	4.7K	14.2K	11
om	Oromo	1.7K	5.0K	18	xh	Xhosa	1.6K	2.8K	1

## C.2 Multilingual entity insertion stratified by language

Figs. 6 and 7 portray the entity insertion performance stratified by language of all the benchmarked methods using hits@1 and MRR, respectively. The results clearly show that, as entity insertion becomes more complex, the baselines start to decrease in performance, being significantly outperformed by LOCEI and xLOCEI.

## C.3 Zero-shot entity insertion stratified by language

Table 11 provides additional details about the data such as the languages and the size of the datasets, used to train the different variants of the multilingual models employed in the zero-shot setting.

Figs. 8 and 9 portray the zero-shot entity insertion performance stratified by language of all the

Table 8: Summary statistics of the train and test sets for 20 Wikipedia language versions considered in the experiments.

Language		Articles	Existing Links	Added Links	
				Train	Test
en	English	6.7M	166M	552K	416K
fr	French	2.5M	85M	130K	76K
it	Italian	1.8M	51M	101K	56K
ja	Japanese	1.4M	61M	150K	111K
pt	Portuguese	1.1M	22M	54K	32K
cs	Czech	526K	14M	27K	15K
ms	Malay	362K	2.9M	6K	3K
cy	Welsh	269K	2.7M	1K	455
sk	Slovak	240K	3.4M	7K	4.3K
simple	Simple English	238K	2.6M	9.4K	4.8K
kk	Kazakh	232K	1.6M	2.7K	2.0K
uz	Uzbek	224K	1.3M	12K	5.9K
ur	Urdu	188K	1.4M	14K	7.5K
hi	Hindi	155K	1.2M	3.2K	3.2K
af	Afrikaans	111K	1.4M	3.3K	1.7K
sw	Swahili	73K	1.0M	1.1K	616
ga	Irish	56K	380K	849	256
is	Icelandic	51K	610K	1.6K	1.2K
gu	Gujarati	30K	410K	197	48
kn	Kannada	27K	250K	1.1K	609

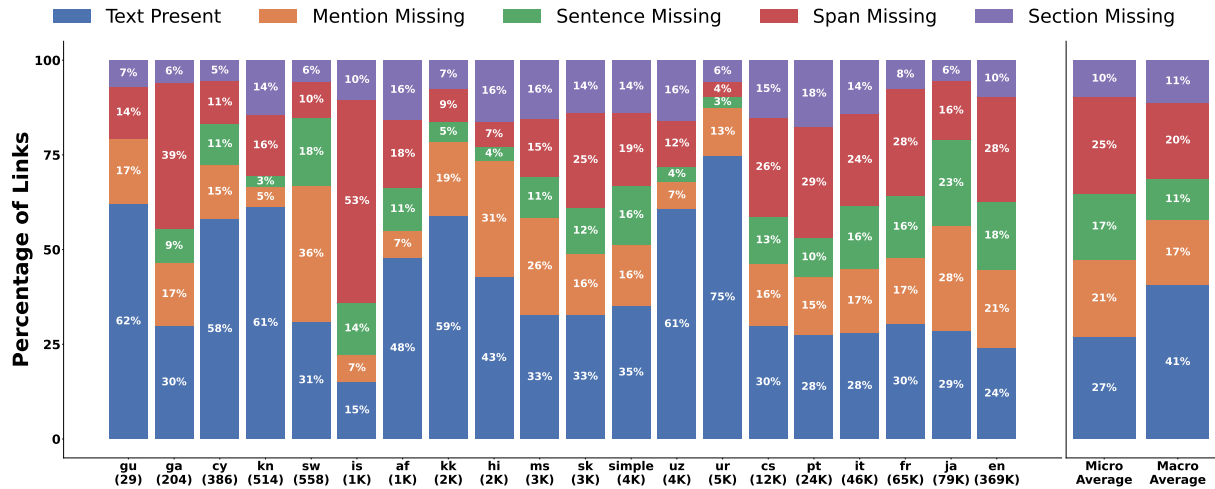


Figure 5: The distribution of entity insertion categories across the 20 considered Wikipedia language versions from October to November 2023. The x-axis shows the language code and the number of links added in each language.

benchmarked methods using hits@1 and MRR, respectively.

#### C.4 Impact of the starting model

Since our approach is based on fine-tuning pre-trained models, the starting pre-trained model may have an impact on the eventual model performance. We studied this dependence using three pre-trained models: BERT<sub>BASE</sub>, RoBERTa<sub>BASE</sub> and the en-

coder portion of T5<sub>BASE</sub>, which we call T5<sub>BASE</sub><sup>enc</sup>. We considered BERT and RoBERTa because they are amongst the most popular transformer encoder models. We additionally included T5 to see how encoder-decoder models perform in the entity insertion task. However, as RankT5 (Zhuang et al., 2023) showed there was no clear benefit in using the full encoder-decoder architecture, as opposed



Table 9: Examples of different entity insertion categories observed when adding links in Wikipedia. The added link is marked in blue.

Strategy	First Version Text	Second Version Text
Text Present	It is best eaten when it is somewhat below normal room temperature. In most countries, brie-style cheeses are made with Pasteurized milk.	It is best eaten when it is somewhat below normal <a href="#">room temperature</a> . In most countries, brie-style cheeses are made with Pasteurized milk.
Missing Mention	Vercetti Regular, also known as Vercetti, is a free font that can be used for both commercial and personal purposes. It became available in 2022 under the Licence Amicale, which allows users to share the font files with friends and colleagues.	Vercetti Regular, also known as Vercetti, is a free font ( <a href="#">freeware</a> ) that can be used for both commercial and personal purposes. It became available in 2022 under the Licence Amicale, which allows users to share the font files with friends and colleagues.
Missing Sentence	Kivi lived in time when all educated people in Finland spoke Swedish. He was the first professional writer who published his works in Finnish. Kivi, Mikael Agricola and Elias Lönnrot are regarded fathers of a national literature in Finnish.	Kivi was born in <a href="#">Nurmijärvi</a> . He lived in time when all educated people in Finland spoke Swedish. He was the first professional writer who published his works in Finnish. Kivi, Mikael Agricola and Elias Lönnrot are regarded fathers of a national literature in Finnish.
Missing Span	The game will be released for Windows PC, Mac and Linux, with Nintendo Switch being the only console to receive the game at launch.  Originally, Hornet was planned as a second playable character to be included in a downloadable content pack (DLC) for Hollow Knight, funded as a stretch goal in the game’s Kickstarter campaign.	The game will be released for Windows PC, Mac and Linux, with Nintendo Switch being the only console to receive the game at launch. During the Xbox & Bethesda Games Showcase, it was revealed that the game would be coming to Xbox Game Pass through PC and Xbox Series X/S. It was also revealed that the game would be coming to PlayStation 4 and <a href="#">PlayStation 5</a> .  Originally, Hornet was planned as a second playable character to be included in a downloadable content pack (DLC) for Hollow Knight, funded as a stretch goal in the game’s Kickstarter campaign.

to encoder-only architecture, and thus, for computational reasons we decided to use the encoder-only variant of T5, T5<sup>enc</sup>.

We trained each model on the Simple English dataset, and we measured their performance on the test data. Table 12 shows that the RoBERTa model outperformed both BERT and T5<sup>enc</sup> in all entity insertion categories by a large margin. BERT and T5<sup>enc</sup> performed similarly, with T5<sup>enc</sup> doing slightly better. These results may be explained by

the fact that the RoBERTa tokenizer has a much larger vocabulary than the tokenizers for BERT or T5<sup>enc</sup>. A larger vocabulary might make it possible for the model to capture more fine-grained linguistic and structural patterns in the candidate text spans, enabling the model to exploit patterns that neither T5<sup>enc</sup> nor BERT can capture.

Table 10: Examples of different strategies for dynamic context removal. The mention of the target link is marked in blue.

Strategy	Original Text	Modified Text
No removal (rm_nth)	Pulaski County is a county located in the central portion of the U.S. state of Georgia. As of the 2020 census, the population was 9,855. The <a href="#">county seat</a> is Hawkinsville.	Pulaski County is a county located in the central portion of the U.S. state of Georgia. As of the 2020 census, the population was 9,855. The county seat is Hawkinsville.
Mention removal (rm_mention)	Perthes-lès-Brienne is a commune of the Aube <a href="#">département</a> in the north-central part of France.	Perthes-lès-Brienne is a commune of the Aube in the north-central part of France.
Sentence removal (rm_sent)	In this Japanese name, the family name is Fujita. Yoshiaki Fujita (born <a href="#">12 January 1983</a> ) is a Japanese football player. He plays for Oita Trinita.	In this Japanese name, the family name is Fujita. He plays for Oita Trinita.
Span removal (rm_span)	Administration The department of French Guiana is managed by the Collectivité territoriale de la Guyane in Cayenne. There are 2 <a href="#">arrondissements</a> (districts) and 22 communes (municipalities) in French Guiana. The cantons of the department were eliminated on 31 December 2015 by the Law 2011-884 of 27 July 2011. The 22 communes in the department	Administration The 22 communes in the department are:

Table 11: Details about the languages and size of the dataset used to train the two xLOCEI model variants, i.e., xLOCEI<sub>20</sub> and xLOCEI<sub>11</sub>.

Model	Starting Model	Fine-Tuned Languages	Training Data Size	
			Stage 1	Stage 2
xLocEI <sub>20</sub>	XLM-RoBERTa <sub>BASE</sub>	en, fr, it, ja, pt, cs, ms, cy, sk, uz, simple, kk, ur, hi, af, sw, ga, is, kn, gu	20K	503K
xLocEI <sub>11</sub>	XLM-RoBERTa <sub>BASE</sub>	en, it, ja, cs, cy, uz, ur, hi, sw, is, kn	20K	348K

Table 12: Comparing the entity insertion performance obtained for Simple English with different starting models. The categorization of entity insertion types into ‘Overall’, ‘Missing’, and ‘Present’ is discussed in § 6.3.

Method	Hits@1			MRR		
	Overall	Present	Missing	Overall	Present	Missing
BERT	0.666	0.916	0.492	0.738	0.940	0.598
T5 <sup>enc</sup>	0.710	0.929	0.558	0.774	0.952	0.650
RoBERTa	<b>0.851<sup>†</sup></b>	<b>0.957<sup>†</sup></b>	<b>0.777<sup>†</sup></b>	<b>0.890<sup>†</sup></b>	<b>0.968</b>	<b>0.835<sup>†</sup></b>

<sup>†</sup> Indicates statistical significance ( $p < 0.05$ ) between the best and the second-best scores.

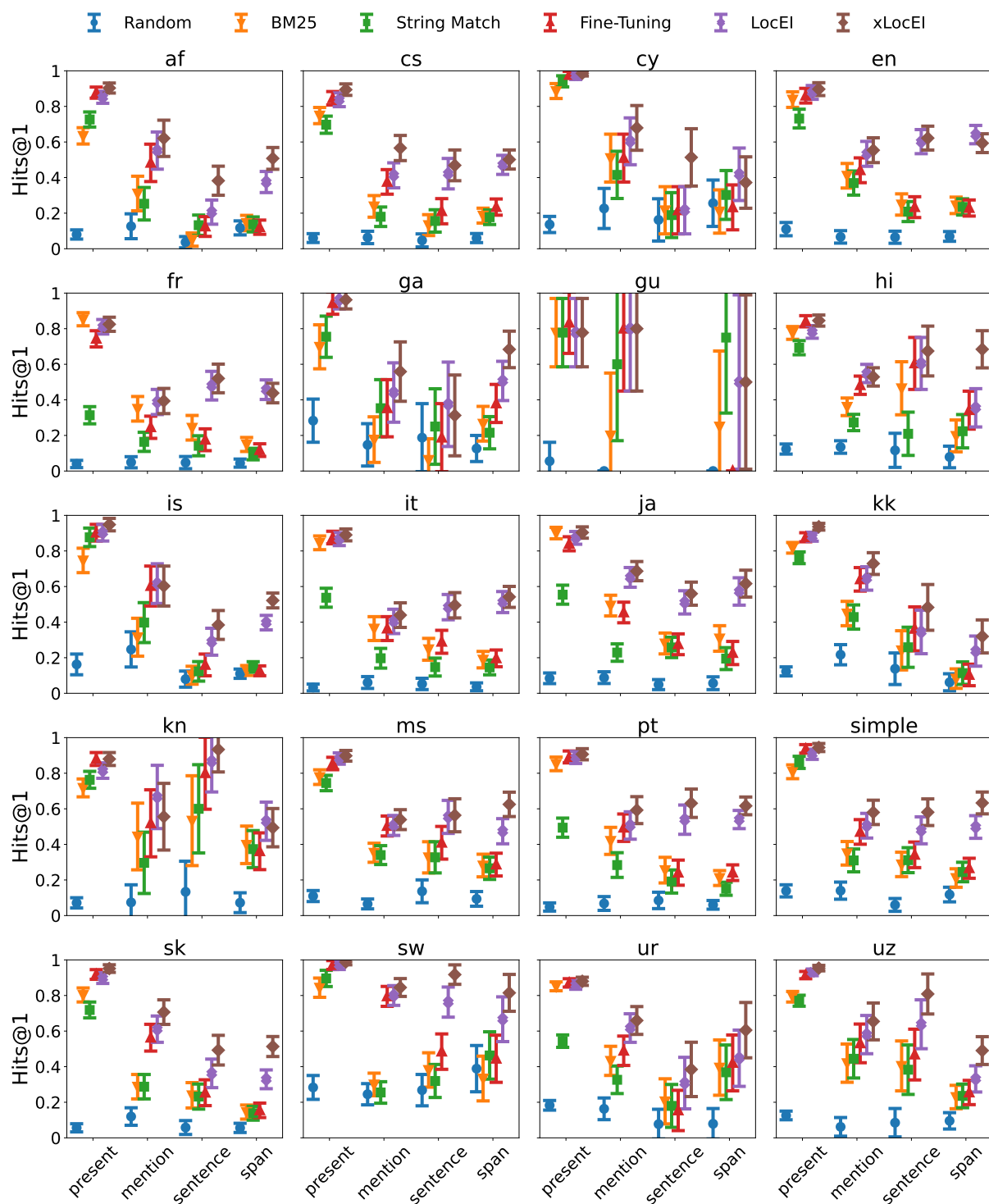


Figure 6: Entity insertion performance across all 20 Wikipedia language versions measured using hits@1. xLOCEI trains a single model jointly on all 20 languages, whereas other methods train a separate model for each language. The categorization of entity insertion types is discussed in § 4.

### C.5 Impact of the model size

There is a widely known trend in the deep learning community that bigger models tend to perform better than smaller models (Soltanolkotabi et al., 2019; Brutzkus and Globerson, 2019; Simon et al.,

2024). To this end, we studied how the model size impacts the entity insertion performance by comparing RoBERTa<sub>LARGE</sub> with RoBERTa<sub>BASE</sub> on the Simple English dataset.

Table 13 shows that there is no statistically significant difference between the performance of

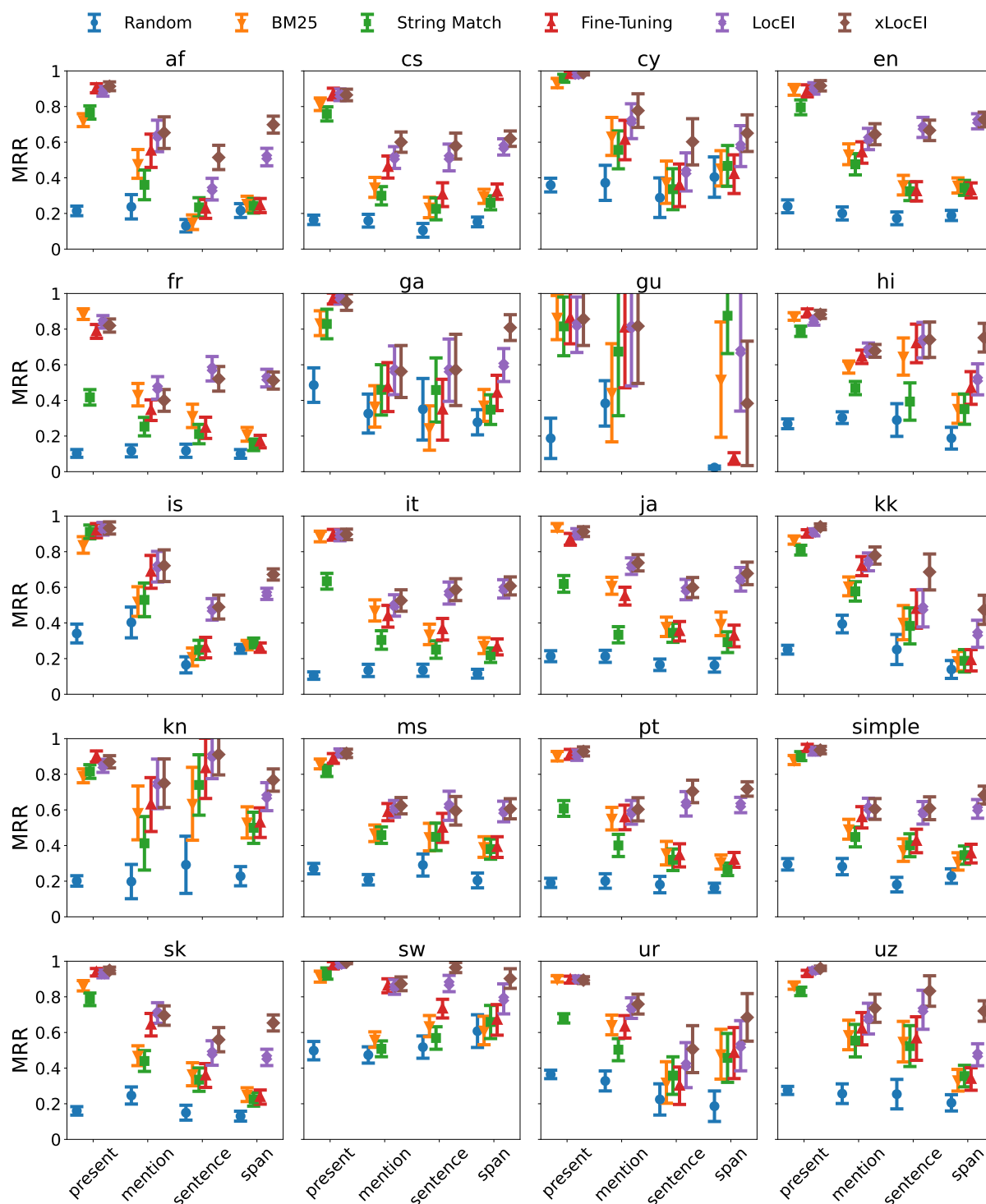


Figure 7: Entity insertion performance across all 20 Wikipedia language versions measured using MRR. xLOCEI trains a single model jointly on all 20 languages, whereas other methods train a separate model for each language. The categorization of entity insertion types is discussed in § 4.

RoBERTa<sub>LARGE</sub> and RoBERTa<sub>BASE</sub>. These results point to the fact that the increased model complexity is not sufficient to improve model performance. It is worth noting that these results were obtained for Simple English. The multilingual problem is much harder and it might benefit from the increased

complexity and larger parameter space of the larger model. We leave this study for future work.

Additionally, these findings give more strength to the hypothesis that the reason why RoBERTa is significantly better than BERT and T5<sup>enc</sup> is because RoBERTa’s larger tokenizer allows the model

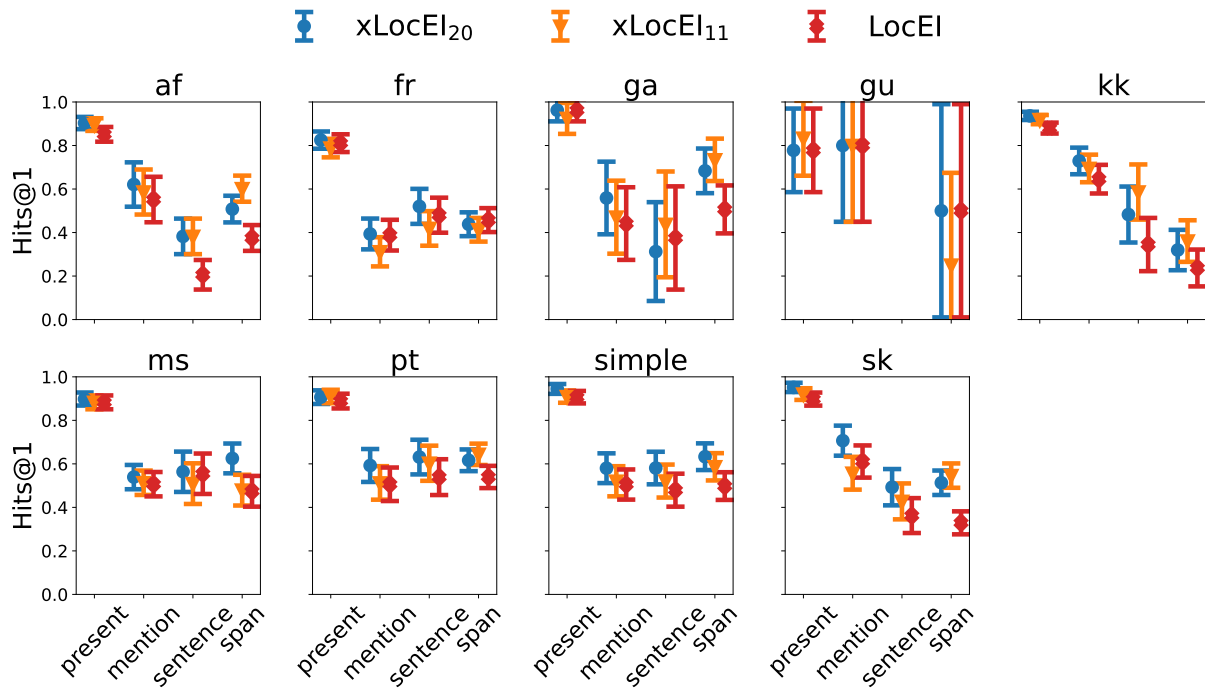


Figure 8: Entity insertion performance measured using hits@1 in the zero-shot setting: results across 9 Wikipedia language versions that were not used for fine-tuning xLOCEI<sub>11</sub>. xLOCEI<sub>20</sub> was trained jointly on all 20 languages, whereas LOCEI trains a separate model for each language. The categorization of entity insertion types is discussed in § 4.

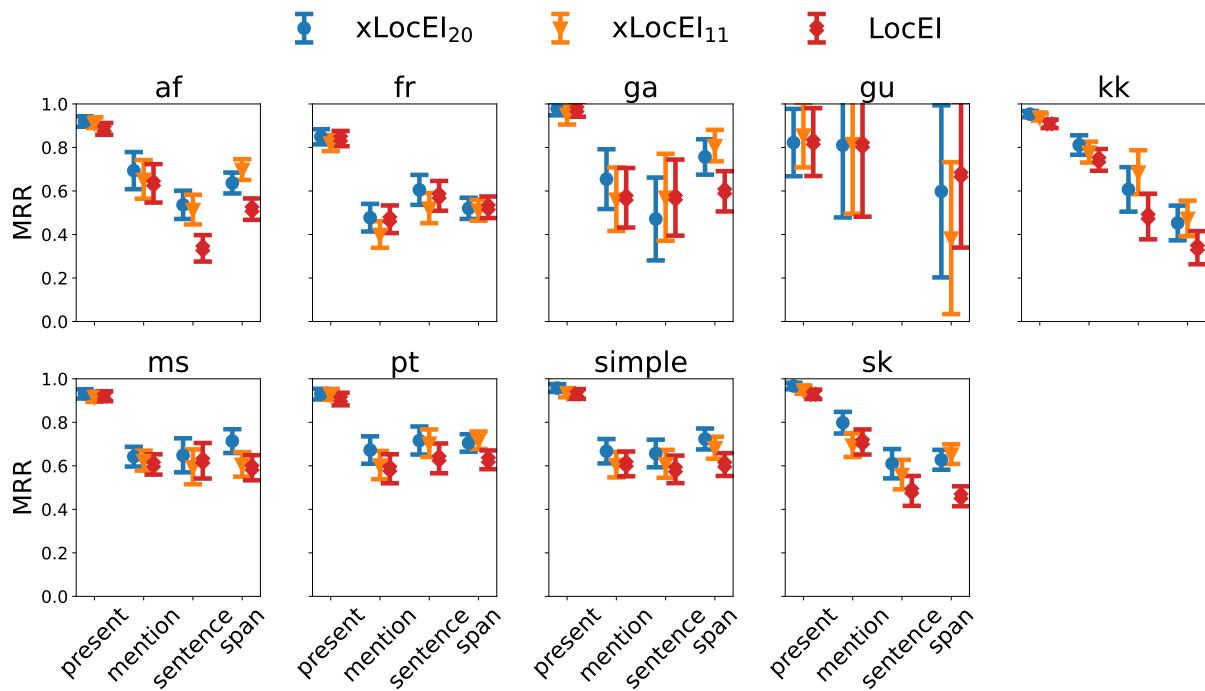


Figure 9: Entity insertion performance measured using MRR in the zero-shot setting: results across 9 Wikipedia language versions that were not used for fine-tuning xLOCEI<sub>11</sub>. xLOCEI<sub>20</sub> was trained jointly on all 20 languages, whereas LOCEI trains a separate model for each language. The categorization of entity insertion types is discussed in § 4.

to capture more fine-grained linguistic and structural patterns in the candidate. This increased input

representation space seems to be vital for entity insertion.

### C.6 Impact of the size of training data

As discussed in § 5.2, we use the existing and added links data during the first and second stages of our training pipeline, respectively. In this analysis, we studied how much data is needed for each stage. To study the impact of the training data size on the downstream entity insertion performance, we trained a RoBERTa<sub>BASE</sub> model with varying portions of the full English dataset.

Fig. 10 shows the performance of LOCEI for different entity insertion categories with varying training data sizes of  $\{10^3, 10^4, 10^5, 10^6\}$  in the first stage of the training pipeline. Note that LOCEI was trained using only the first stage for this analysis. Fig. 11 shows an analogous plot for the second stage with varying training data sizes of  $\{10^2, 10^3, 10^4, 10^5\}$ . For this analysis, LOCEI was trained using only the second stage.

These results show that it is much more important to have more data in the second stage when compared to the first stage. The performance did not visibly improve over the data range considered for the first stage, indicating no benefit in training on a lot of existing links. On the other hand, the model performance improved drastically as the data size increased for the second stage, with no sign of plateauing. Based on these results, the optimal training schedule for an entity insertion model using our data seems to be a short first stage, followed by a second stage using as much data as possible.

### C.7 Training stages

Table 14 shows the impact of different training strategies: (1) Warm start (only using the first stage), (2) Expansion (only using the second stage), and (3) Warm start + Expansion (using both stages), on the downstream entity insertion performance of LOCEI using data extracted from English Wikipedia.

### C.8 RoBERTa vs XLM-RoBERTa

We found the scores obtained with RoBERTa on Simple English to be significantly higher than the scores achieved by the multilingual XLM-RoBERTa. In this analysis, we compare the performance of these two models on the full English dataset, with both models having been fine-tuned

on the same English dataset. Table 15 shows a statistically significant difference in the performance of RoBERTa and XLM-RoBERTa, with RoBERTa scoring higher in all entity insertion strategies, with gaps up to 25%. We draw two conclusions from these results.

In our ablations, we found that RoBERTa outperformed BERT and T5<sup>enc</sup> by a large margin, which leads us to select XLM-RoBERTa as the best candidate for the multilingual model to use in our experiments. However, the performance of RoBERTa does not seem to directly correlate with the performance of XLM-RoBERTa, as seen by the large drop in English when moving from RoBERTa to XLM-RoBERTa. This finding casts some doubt on the decision of the best multilingual model and opens the doors to models like multilingual BERT and mT5 (Xue et al., 2021). In the future, it would be interesting to consider other multilingual models and see if they can outperform XLM-RoBERTa.

As shown in § 6.4, XLM-RoBERTa fine-tuned on the multilingual dataset generally outperformed XLM-RoBERTa fine-tuned on a single language. However, the results in Table 15 point to the fact that a model pre-trained on a single language (RoBERTa) outperforms a model pre-trained on multiple languages (XLM-RoBERTa). The dominance of the monolingual model is not surprising as a model pre-trained on a single language had a much smaller domain to learn than a multilingual model, and thus, might have been able to learn linguistic and structural patterns that the multilingual model failed to capture. So, for the languages where a pre-trained model does exist (for example, BERT for English, CamemBERT (Müller et al., 2020) for French, HerBERT (Rybak et al., 2020) for Polish), that model may outperform the multilingual variant. However, it is unrealistic to assume that there can be a pre-trained model for each of the 300+ languages of Wikipedia. The multilingual model becomes essential for the languages for which there is no pre-trained model. As we saw in § 6.4 and § 6.4, the multilingual model is capable of transferring knowledge to unseen languages, which proves its potential for low-resource languages for which a full pre-trained model is not realistic.

### C.9 Single Encoder vs Triple Encoder

In early iterations of our work, we explored a different model architecture. This architecture used the additional knowledge of the source article. Given

Table 13: Comparing the entity insertion performance obtained for Simple English with varying model sizes. The categorization of entity insertion types is discussed in § 4.

Model	Text Present		Missing Mention		Missing Sentence		Missing Span	
	Hits@1	MRR	Hits@1	MRR	Hits@1	MRR	Hits@1	MRR
RoBERTa <sub>BASE</sub>	0.956	0.968	<b>0.696</b>	<b>0.760</b>	0.834	0.884	0.799	0.859
RoBERTa <sub>LARGE</sub>	<b>0.964</b>	<b>0.975</b>	0.670	0.744	<b>0.856</b>	<b>0.895</b>	<b>0.822</b>	<b>0.873</b>



Figure 10: Impact of the amount of data used in the first stage on the downstream entity insertion performance. Note that the model is trained solely using the first stage. The categorization of entity insertion types is discussed in § 4.



Figure 11: Impact of the amount of data used in the second stage on the downstream entity insertion performance. Note that the model is trained solely using the second stage. The categorization of entity insertion types is discussed in § 4.

Table 14: Comparison of the impact of different stages of the training pipeline on the downstream entity insertion performance. The categorization of entity insertion types into ‘Overall’, ‘Missing’, and ‘Present’ is discussed in § 6.3.

Training Stages	Hits@1			MRR		
	Overall	Present	Missing	Overall	Present	Missing
Warm start	0.584	<b>0.883</b>	0.350	0.649	<b>0.907</b>	0.451
Expansion	0.604	0.738	0.494 <sup>†</sup>	0.689	0.801	0.603 <sup>†</sup>
Warm start + Expansion	<b>0.672<sup>†</sup></b>	0.877 <sup>†</sup>	<b>0.509</b>	<b>0.744<sup>†</sup></b>	0.906 <sup>†</sup>	<b>0.617</b>

<sup>†</sup> Indicates statistical significance ( $p < 0.05$ ) between the variant and the previous variant.

Table 15: Comparing the entity insertion performance of our model fine-tuned using the monolingual RoBERTa<sub>BASE</sub> and the multilingual XLM-RoBERTa<sub>BASE</sub> on the data extracted from English Wikipedia. The categorization of entity insertion types is discussed in § 4.

Model	Text Present		Missing Mention		Missing Sentence		Missing Span	
	Hits@1	MRR	Hits@1	MRR	Hits@1	MRR	Hits@1	MRR
RoBERTa <sub>BASE</sub>	<b>0.923</b> <sup>†</sup>	<b>0.936</b> <sup>†</sup>	<b>0.737</b> <sup>†</sup>	<b>0.797</b> <sup>†</sup>	<b>0.850</b> <sup>†</sup>	<b>0.898</b> <sup>†</sup>	<b>0.787</b> <sup>†</sup>	<b>0.848</b> <sup>†</sup>
XLM-RoBERTa <sub>BASE</sub>	0.863	0.892	0.543	0.630	0.595	0.662	0.697	0.615

<sup>†</sup> Indicates statistical significance ( $p < 0.05$ ).

Table 16: Comparing the entity insertion performance obtained for Simple English with different loss functions: pointwise vs. ranking loss. The categorization of entity insertion types into ‘Overall’, ‘Missing’, and ‘Present’ is discussed in § 6.3.

Method	Hits@1			MRR		
	Overall	Present	Missing	Overall	Present	Missing
Pointwise Loss	0.641	0.891	0.477	0.712	0.922	0.574
Ranking Loss	<b>0.658</b>	<b>0.907</b>	<b>0.495</b>	<b>0.731</b>	<b>0.930</b>	<b>0.601</b>

the amount of text that needed to be encoded, and considering that most transformers have a limited number of tokens they can process, we chose to encode each of the three components separately. We had the following input representations:

- Source Article: [CLS]<Src Title>[SEP]<Src Lead>
- Candidate: [CLS]<Src Section>[SEP]<Tgt Mention>[SEP]<Context>
- Target Title: [CLS]<Tgt Title>[SEP]<Tgt Lead>

Each of the components of the triplet was encoded independently, and then stacked together. Finally, an MLP capturing the interactions between the three embeddings was used to produce a relevance score.

The key intuition behind this architecture was to represent a link as a knowledge triplet <src, text, tgt>, and the overall architecture was supposed to predict whether the triplet was correct. However, we found that such an architecture decayed into a state where the target and source embeddings were independent of the input, always producing the same embedding. We believe that the model relied exclusively on the semantic knowledge contained in the list of target mentions to identify whether the entity should be inserted in the candidate text span, and the source and target article embeddings decayed into a global average optimum that maximized the performance of the MLP for the candidate embedding. Nevertheless, this meant that all

the knowledge about the target entity contained in the target lead was being ignored.

To take advantage of the total available information, we moved to the architecture described in § 5.1. We removed the source title and the source lead, driven by the token limit of the transformer architecture. We believed that this knowledge provided the least marginal gain from the three components of the triplet, at a cost of token space for the candidate and the target, as the source article knowledge only gave additional context to the candidate text span.

We additionally moved to a single encoder for two reasons. First, the transformer architecture is more expressive than an MLP, and thus, it was better suited to capture the interactions between the candidate and the target. With only two knowledge sources instead of three, we felt we had sufficient token space for each source to capture enough semantic information for each input. Second, by relying on one single embedding, the embedding couldn’t decay into a global average optimum which provided no information about the input, because the relevance score was entirely dependent on the representation power of that single embedding.

### C.10 Pointwise Loss vs Ranking Loss

Table 16 shows how the choice of different loss functions (pointwise vs. ranking) impacts the downstream entity insertion performance of our models evaluated on the Simple English dataset.