# QUALIFYING KNOWLEDGE AND KNOWLEDGE SHAR-ING IN MULTILINGUAL MODELS

Anonymous authors

Paper under double-blind review

### ABSTRACT

Pre-trained language models (PLMs) have demonstrated a remarkable ability to encode factual knowledge. However, the mechanisms underlying how this knowledge is stored and retrieved remain poorly understood, with important implications for AI interpretability and safety. In this paper, we disentangle the multifaceted nature of knowledge: successfully completing a knowledge retrieval task (e.g., "The capital of France is \_\_") involves mastering underlying concepts (e.g., France, Paris), relationships between these concepts (e.g., *capital of*), the structure of prompts, including the language of the query. We propose to disentangle these distinct aspects of knowledge and apply this typology to offer a critical view of neuron-level knowledge attribution techniques. For concreteness, we focus on Dai et al.'s (2022) Knowledge Neurons (KNs) across multiple PLMs, testing 10 natural languages and unnatural languages (e.g. Autoprompt). Our key contributions are twofold: (i) we show that KNs come in different flavors, some indeed encoding entity level concepts, some having a much less transparent, more polysemantic role, and (ii) we uncover an unprecedented overlap in KNs across up to all of the 10 languages we tested, pointing to the existence of a partially unified, language-agnostic retrieval system. To do so, we introduce and release the Multi-ParaRel dataset, an extension of ParaRel, featuring prompts and paraphrases for cloze-style knowledge retrieval tasks in parallel over 10 languages.

028 029

031

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

### 1 INTRODUCTION

032 Recent advances in Large Language Models (LLMs) have led to models trained on vast and diverse 033 linguistic datasets drawn from across the Internet, incorporating numerous languages simultaneously 034 (Scao et al., 2023; Touvron et al., 2023; Achiam et al., 2024). However, these languages are not evenly represented, and performance on low-resource languages often depends on cross-linguistic transfer from high-resource languages (Pires et al., 2019; Lample & Conneau, 2019; Conneau et al., 037 2020a; Huang et al., 2021). Whether LLMs can develop common, language-agnostic representations that enable such zero-shot transfer remains an open question in the literature (Singh et al., 2019; 038 Kudugunta et al., 2019; Kassner et al., 2021). Kervadec et al. (2023) extended this investigation to machine-generated languages, revealing that different representations can emerge, suggesting 040 multiple ways knowledge may be encoded in LLMs. 041

Understanding how Pre-trained Language Models (PLMs) store and retrieve knowledge is essential
for enhancing interpretability and safety in AI systems. Many recent studies have sought to localize
and attribute specific knowledge to individual neurons within these models (Dai et al., 2022; Meng
et al., 2022; 2023). These methods often attempt to identify neurons whose activations are critical
for making accurate predictions. Typically, they focus on neurons in intermediate layers of FeedForward Networks (FFNs) within transformer architectures (Geva et al., 2021). These approaches
face strong limitations, as highlighted in recent critiques (Hase et al., 2023; Niu et al., 2023; Huang
et al., 2023).

In this work, we offer a novel perspective by refining the concept of "knowledge" itself. To correctly complete a prompt like *The capital of France is*, a model must process multiple layers of information: sensitivity to the specific concept *France*, retrieval of the target concept *Paris*, and understanding the relational context *capital of*. We introduce a method to distinguish these subtypes of knowledge—conceptual and relational—that is compatible with any knowledge attribution



Figure 1: The Knowledge Neurons (KNs) hypothesis connects LLM success on a fill-in-the-blank cloze task (e.g. *The capital of France is*) to the activation of a small set of neurons. (a) The same neurons can be selected (green) in response to a single task, thereby qualifying as *concept* neurons (about e.g., Paris) or in response to a range of tasks all concerning a certain relations between concepts, thereby qualifying as *relational neurons* (e.g., *capital of* is a relation between France and Paris, between England and London, etc.). (b) In multilingual LLMs, concept and relational neurons may be selected specifically for a language or across languages.

075

076

077

078

079

081

082

084

085

087

088

technique. We apply this method to the Knowledge Neurons (KNs) framework introduced by Dai
et al. (2022), to provide a critical view on such a method and extend it to investigate how knowledge is shared across languages in PLMs (Figure 1). Code and data available at [URL redacted for
anonymous review].

073 074 Our contributions are:

- We propose a finer-grained typology of knowledge, providing a critical perspective on neuron-level attribution methods like the Knowledge Neuron hypothesis, in particular its expectation of monosemanticity.
- We analyze through this prism multiple PLMs (BERT, mBERT, OPT, Llama 2, and Gemma 2), revealing that a substantial number of 'Knowledge Neurons' exhibit polysemantic behavior, while others are specifically responsive to individual concepts or relations.
  - We release Multi-ParaRel, a multilingual version of the ParaRel dataset (Elazar et al., 2021a), which includes 10 languages and is compatible with autoregressive models.
  - We demonstrate that LLMs store knowledge in similar neurons across 10 languages, and even in machine-generated languages (AutoPrompt), suggesting a shared cross-linguistic mechanism for knowledge retrieval.

# 2 RELATED WORK

Multilingual Language Models Training separate models for different languages is resourceintensive, data-hungry, and generally ineffective at leveraging cross-linguistic similarities and knowledge. In practice, recent LLMs (Touvron et al., 2023; Achiam et al., 2024) are trained on extensive portions of the Internet, making them *de facto* multilingual. Examples include mBERT (Devlin et al., 2019), XLM-R (Lample & Conneau, 2019; Conneau et al., 2020a), mBART (Liu et al., 2020), mT5 (Xue et al., 2021), and BLOOM (Scao et al., 2023), along with their fine-tuned variants like BLOOMZ, mT0 (Muennighoff et al., 2023), and FLAN-T5 (Chung et al., 2022).

096 The performance of these models is believed to stem from the emergence of efficient representa-097 tions that are shared across languages (Aharoni et al., 2019; Arivazhagan et al., 2019; Conneau 098 et al., 2020b). Research has investigated their cross-linguistic capabilities using artificial languages (Ri & Tsuruoka, 2022; Deshpande et al., 2022; Guerin et al., 2024), evaluating their performance on 100 tasks across different languages (Pires et al., 2019; Wu & Dredze, 2019), analyzing their translation 101 capabilities (Lample et al., 2018; Sennrich et al., 2016; Artetxe et al., 2018), and assessing their per-102 formance on low-resource languages (Garcia et al., 2021), as well as examining their architectural 103 properties (K et al., 2020). Some studies, including the current work, have directly compared repre-104 sentations from one language to another (e.g., using Canonical Correlation Analysis across layers, 105 as in Singh et al., 2019; Kudugunta et al., 2019). The conclusions drawn from these comparisons are mixed. For instance, Singh et al. (2019) argue that representations are distinctly partitioned be-106 tween languages, while Kudugunta et al. (2019) suggest that representations are more or less shared, 107 depending on the linguistic proximity of the languages.

108 More directly to the current work, Chen et al. (2024) recently looked at the overlap between knowl-109 edge neurons obtained for English and Chinese. We examine a similar overlap, albeit comparing 10 110 natural languages at once and we believe it is essential. Pairwise sharing leaves ambiguity: are neu-111 rons shared with all languages, none, or just pairs? Bias toward a dominant language (e.g., English) or chance sharing between two languages is more likely than sharing across 10. Including 10 lan-112 guages reveals symmetrical roles in neuron sharing, extending beyond language pairs and providing 113 robust evidence of truly multilingual knowledge representation. We also make a comparison with 114 an 'unnatural language' (Shin et al., 2020). Such prompts provide an extreme test for the idea that 115 knowledge could be accessed independently of form: they are not human-readable, and they had 116 been shown to be processed differently by LLMs (Kervadec et al., 2023). 117

118

**Knowledge in LLMs** LLMs acquire knowledge by training on extensive corpora (Petroni et al., 119 2019; Roberts et al., 2020; Safavi & Koutra, 2021). The work by Petroni et al. (2019) introduced 120 LAMA, a dataset designed to evaluate BERT through a fill-in-the-blank cloze task (e.g., The capital 121 of France is [MASK].). Subsequent research has built upon LAMA (Jiang et al., 2021), highlighting 122 the limitations of LLMs as knowledge bases (Elazar et al., 2021b; AlKhamissi et al., 2022), while 123 also attempting to enhance their performance (Wei et al., 2021; Petroni et al., 2020). Consequently, 124 research has emerged focusing on localizing and editing knowledge directly within the model (Rad-125 ford et al., 2017: Lakretz et al., 2019: Bau et al., 2020b: Sinitsin et al., 2020: Mitchell et al., 2021: 2022; De Cao et al., 2021; Santurkar et al., 2021; De Cao et al., 2022; Bau et al., 2020a; Cohen et al., 126 2023). 127

128 In this context, knowledge attribution methods such as ROME (Meng et al., 2022) and MEMIT 129 (Meng et al., 2023) (both employing causal mediation techniques; Vig et al., 2020), along with 130 Knowledge Neurons (Dai et al., 2022) (utilizing an integrated gradient approach; Sundararajan et al., 131 2017), have been proposed. These methods are predicated on the assumption that neurons within the intermediate layers of transformers' Feed-Forward Networks (FFNs) encode knowledge. However, 132 we align with other studies (Hase et al., 2023; Niu et al., 2023; Huang et al., 2023) that suggest this 133 assumption may be an oversimplification. While certain neurons play a significant role in specific 134 tasks (Lakretz et al., 2019; Manning et al., 2020; Rogers et al., 2020; He et al., 2024), LLM neurons 135 are not necessarily monosemantic; rather, they can serve multiple functions depending on the context 136 and task (Adly et al., 2024). Furthermore, their effectiveness in altering knowledge is subjective and 137 widely debated (Hase et al., 2023). Other works (Wang et al., 2024; Tang et al., 2024; Kojima et al., 138 2024) have identified multilingual neurons in LLMs; this paper focuses specifically on knowledge-139 related neurons, offering a more precise analysis. We propose a knowledge-attribution method-140 agnostic typology, illustrated with Dai et al.'s (2022) Knowledge Neurons. This approach aims to 141 provide a critical view on the Knowledge Neurons hypothesis while exploring what insights it can 142 offer regarding how knowledge is encoded in LLMs.

143 144

145 146

147

148

149

150

## 3 METHODOLOGICAL BACKGROUND

**Knowledge** The TREx dataset (Elsahar et al., 2018) is a collection of relational facts stored in triplets of the form  $\langle h, r, t \rangle$ , with r a relation and h and t entities entering in that relation. TREx exhibit 41 relations, such as *being the capital of, was born in*, etc. Each full triplet can be referred to as an **instantiation** of its own relation r.

151 Knowledge Localization Methods Geva et al. (2021) observed that a FFN can be seen as a Key-152 Value memory system, similar to self-attention. To assess if and where knowledge could be stored in FFNs, Dai et al. (2022) used a knowledge attribution method based on integrated gradients (see next 153 paragraph for details). They show that a fact (e.g., The capital of France is Paris) can be associated 154 to a few neurons (around 4), whose activations correlate with the probability of the model to fill 155 in the elements of the fact appropriately. Similarly, Meng et al. (2022) proposed Rank-One Model 156 Editing (ROME), which uses causal mediation to localize and edit knowledge in GPT, and Meng 157 et al. (2023) introduced Mass-Editing Memory in a Transformer (MEMIT), which edits facts at 158 scale. All of these knowledge attribution methods have their limitations; we apply our analysis to 159 the Knowledge Neurons by way of illustration. Our approach is applicable to all such methods. 160

162 Knowledge Neurons Dai et al. (2022) track Knowledge Neurons (KNs) during a fill-in-the-163 blank cloze task (see also Petroni et al., 2019) based on TREx. Let  $w_i^{(l)}$  be the  $i^{th}$  neuron 164 of the intermediate layer of the  $l^{th}$  FFN. The knowledge score of a neuron  $w_i^{(l)}$  is calculated 165 through the integrated gradient attribution method (Sundararajan et al., 2017), KNs are then fil-166 tered through thresholds. First, they retain only neurons with an attribution score greater than 167  $t_{kn} \times \max_{i,l} \operatorname{Attr}_{h,p_r,t}(w_i^{(l)})$ . This procedure is carried out for each prompt associated with a fact  $\langle h, r, t \rangle$ , and thus yields a set of candidate KNs per prompt. Let us denote  $N_r$  the number of 168 prompts for a given relation r. To get results robust to noise, and to factor out signal associated to 170 specific prompts rather than knowledge, they keep only neurons appearing in the candidate neurons 171 set of at least  $p_{kn} \times N_r$  prompts. They propose thresholds of  $t_{kn} = 0.2$  (only keep neurons scoring 172 at least at 20% of the max attribution score) and  $p_{kn} = 0.7$  (only keep neurons appearing in at least 173 70% of the different prompts for a given relation).

174 175

# 4 Method

176 177

Datasets For relational facts, we used the TREx dataset (Elsahar et al., 2018), which comprises
41 relations with approximately 1,000 facts per relation. For prompts, we employed the augmented
version of ParaRel provided by Kervadec et al. (2023). This version retains only prompts compatible with autoregressive models and enriches the dataset with multiple paraphrases for each relation.
In Section 6, we explore multilingual models, which we tested on the multilingual variant of LAMA
(Kassner et al., 2021) as well as on a new multilingual version of ParaRel that we introduce. We
refer to this new dataset as Multi-ParaRel.

The detailed methodology for creating Multi-ParaRel, along with a quality assessment, is provided in Appendix A. Our dataset currently spans 10 languages: English, French, Spanish, Catalan, Danish, German, Italian, Dutch, Portuguese, and Swedish. We also investigate an unnatural language: AutoPrompt. Following the same train, development, and test splits as Shin et al. (2020), we trained 10 different seeds of AutoPrompt for each relation and each model. We also make these sets of prompts available.

190

191 **Concept Neurons and Relation Neurons** We propose a simple typology that refines the type of knowledge attributed while answering fill-in-the-blank cloze tasks. For example, correctly answer-192 ing the question What is the capital of France? not only requires knowledge of the answer Paris, 193 but also an understanding of the relationship between *France* and *Paris*. We thus introduce a sim-194 ple principle: a neuron that is hypothesized to encode a specific concept, such as one about *Paris*, 195 should not be also responsible for encoding other concepts, and should therefore not be associated 196 to other facts such as The capital of Spain is Madrid. If a neuron consistently encodes the same 197 relation across multiple instances, we refer to it as a relational neuron, indicating that it is sensitive to a relation, such as *capital of*. 199

We thus define **Relation Neurons** as KNs that appear in at least  $t_r \times N$  instances of facts associated with a particular relation, where N is the total number of facts, and  $t_r$  is a predefined relational threshold. In contrast, neurons that appear in less than  $t_c \times N$  of the facts, for some other threshold  $t_c$ , are referred to as **Concept Neurons**, as they are more likely to encode specific pieces of knowledge or information about individual entities.

The aim is to test the robustness of this distinction by investigating the role of the thresholds  $t_r$ and  $t_c$ . A 'clean' scenario that supports the Knowledge Neuron hypothesis and the idea of monosemanticity would show that some concept neurons are found even for  $t_C \times N = 1$  (very specific to a concept), and relational neurons are found when  $t_R \times N = N$  (completely systematically present for a relation). Alternatively, softer boundaries would suggest that these KNs play a more polysemantic and nuanced role, whereby knowledge is partially distributed across different neurons on different occasions (e.g., the concept of *Paris* and *Madrid* cannot be disentangled at the neuron level, or the relation *capital of* is not always encoded in the same way).

As we do not know a priori which neurons play specific roles, we performed an exhaustive study across varying thresholds. In fact, it is part of the method to look at all possible thresholds to identify the behavior of KNs. Moreover, no major variation based on the choice of threshold was found.



238

239

240

241

242

243

244 245 246



Figure 2: Each panel corresponds to a relation (P108, P159, etc.). (a) Distribution of KNs based on the number of instantiations (i.e. specific triplets, specific facts) within a relation for which a KN was identified. A large number of neurons are identified as KN for a single instantiation, while a roughly similar number of neurons are identified as KN for a continuously increasing number of instantiations within a relation. (b) Average proportion of the KNs from a single instantiation which can be categorized as **Concept**, **Relation** or neither, according to different thresholds (xaxis). The proportion of relational neurons is stable across different thresholds, the proportion of concept neurons decreases with more demanding thresholds.

Multilingual Knowledge Neurons Similarly, we ask whether knowledge is language-agnostic;
 for example, humans do not need to relearn facts when acquiring a new language. Knowledge could
 be language-dependent in LLMs however: if a fact is present from the English corpus but missing
 from a Spanish training corpus, an LLM may be able to retrieve that knowledge when prompted in
 English but not when prompted in Spanish. We employ the KNs framework to investigate the open
 question of whether a common language-agnostic knowledge representation exists in multilingual
 models at the level of neurons.

We hypothesize that some KNs may be specific to one language, while others may be sensitive to prompts in multiple languages. We thus analyze the number of languages across which such neurons are shared. We do so by identifying KNs for relations in the ParaRel dataset across multiple languages, using the Multi-ParaRel dataset, which was specifically created for this multilingual evaluation.

# 5 MONOLINGUAL EXPERIMENTS: TRACKING CONCEPT AND RELATION NEURONS

261 262

264

254

255

256

257

258 259

260

# 5.1 EXPERIMENTAL SETTINGS

Models We studied BERT (Devlin et al., 2019), and more precisely bert-base-uncased and bert-large-uncased, as it has been the reference model for evaluation on TREx since Petroni et al. (2019). Having been trained on Wikipedia from which TREx is derived, their performance is very good (P@1> 0.4). We also studied OPT (Zhang et al., 2022) in its 350 million-parameters version opt-6.7b, Llama 2 (Touvron et al., 2023) in its 7 billion-parameter version Llama-2-7b-hf as well as Gemma 2 (Team et al., 2024) in its

9 billion parameters version gemma-2-9b. For all these models we use the HuggingFace implementation. KNs computations were performed on NVIDIA Tesla V100 GPUs for models with less than a billion parameters, and on NVIDIA Tesla A100 GPUs for larger models. The computation took less than an hour per relation.

274

Template filtering Per model, we excluded prompts with less than 10% top-1 accuracy (that is, accuracy of the most probable continuation). We then excluded relationships with less than 4 prompts left. Since all actual answers were made of a single token, we also limit answers made of a single token. After this filtering, we obtained on average 15 prompts per relation for BERT and 8 prompts per relation for OPT (starting from 18), confirming the higher accuracy of BERT at the task.

280 281

### 5.2 TRACKING A TYPOLOGY OF KNOWLEDGE

Before classifying Knowledge Neurons (KNs) according to our typology, we first analyzed the distribution of KNs based on the number of instantiations for which a KN was identified. Figure 2a illustrates the results for four relations using the Llama-2-7b model (complete results are provided in Appendix B). A qualitative analysis reveals two key findings: (i) many KNs appear in only one instantiation, indicating that these neurons are task-specific and sensitive to a single concept; and (ii) there is a continuous range of KNs sensitive to between 3 and N instantiations, suggesting a more nuanced role for these neurons that lies between relational and conceptual.

The second observation challenges the simplistic interpretation of assigning neurons exclusively to concepts. At the same time, it also demonstrates the presence of a significant number of neurons sensitive to enough instantiations to hypothesize a more relational role in knowledge retrieval mechanisms.

Thus, we have identified potential candidates for the roles of both **Concept Neurons** and **Relation** Neurons, as well as neurons that fall into an intermediate category. The natural question that arises is: what is the proportion of each neuron type per instantiation, based on thresholds  $t_r$  and  $t_c$ ? This information is not directly inferable from Figure 2a, as neurons appearing consistently across instantiations are less visible than neurons that appear uniquely in each instance.<sup>1</sup> To address this, we computed the proportion of each neuron type as a function of thresholds at the instantiation level (see Figure 2b). For simplicity, we used symmetrical thresholds, setting  $t_r = 1 - t_c$ .

As expected, when the thresholds become more restrictive, the number of neurons with well-defined roles decreases, giving way to neurons with less clearly defined functions across all relations. For the Llama-2-7b model, we observe that the number of neurons classified as **Relation Neurons** remains more stable compared to those classified as **Concept Neurons**. Furthermore, for a single instantiation, there are few KNs that are exclusive to that instance: when  $t_c < 0.1$ , the proportion of **Concept Neurons** is less than 0.2.

We also examined the distribution of neuron types across the model's layers but found no significant variation. As observed by Dai et al. (2022), KNs are primarily concentrated in the final layers.

In summary, we have demonstrated the existence of neurons reacting specifically to a single concept within a relation. We have also identified neurons that play a much broader role in such relations, with some reacting to almost all instances of that relation. We attempt to verify this hypothesis through causal experiments in the next section. Finally, some neurons are activated by a subset of the instantiations, carrying a much less transparent type of knowledge. In principle, it could encode subtypes of relations, such as 'capital of a European country', although we find this highly stipulative at the moment. In the next section, we will focus on concept and relation neurons and evaluate their role through causal experiments.

- 316
- 317
- 318 319
- 320
- 321
- 322
- 323



Figure 3: Boosting experiments results for bert-base-uncased (left) and Llama-2-7b (right) for two couple of thresholds  $t_r = 0.6$ ,  $t_c = 0.4$  (top) and  $t_r = 0.8$ ,  $t_c = 0.2$  (bottom). The lines corresponds to the  $\Delta P@k$  (resp.  $\Delta CCP@k$ ) for different k values ranging from 1 to 100. Thick lines represents the doubled activations results and thin lines the nullified activations results. We also plotted the standrad error accross the evaluated instantiations of the relations.

### 5.3 BOOSTING EXPERIMENTS

343 344

345

374

346 In this experiment, we investigate the effect of either doubling or nullifying the activation of KNs 347 on model predictions. Dai et al. (2022) conducted similar experiments, focusing on how manual 348 changes to neuron activations influenced output probabilities. In contrast, we employ two more 349 concrete impact metrics: precision at rank k, denoted P@k, which measures the proportion of correct 350 responses in the top k model predictions, and correct category proportion at rank k, denoted CCP@k, 351 which reflects the proportion of responses in the correct category (e.g., *capitals*) within the top k predictions. The original metric of relative probabilities change would not show specificity (e.g. 352 unrelated tokens could be even more boosted). For this reason, we report P@k and CCP@k. Effects 353 here ensure that the boost to the correct answer overcomes any boost for other answers. We also 354 include a control experiment in Appendix B to better investigate specificity. 355

356 Our goal is to verify whether the behavior of the identified KNs aligns with our proposed typology. 357 Specifically, we hypothesize that (i) there will be a marked increase (or decrease) in precision at rank 358 k=1 when the activations of **Concept Neurons** are doubled (or nullified), with the effect diminishing as k increases. Similarly, we anticipate (ii) that the effect of **Relation Neurons** on P@k will be 359 weaker than that of **Concept Neurons**, as precision is primarily sensitive to the correct response. In 360 contrast, for the CCP@k metric, we expect (iii) that Relation Neurons will play a more significant 361 role, as these neurons should be more likely to favor the correct category (e.g., capitals), even if it 362 does not boost the correct answer specifically. We assess these effects for a range of thresholds  $t_c$ 363 and  $t_r$ . Results for the bert-base-uncased and Llama-2-7b models are shown in Figure 3 364 (see Appendix **B** for additional models and thresholds as well). 365

The figures show the delta in P@k and CCP@k for the predictions with altered (doubling or nuli-366 fying) vs unaltered activations. The horizontal line at zero thus represents the baseline model per-367 formance. Of the six models evaluated, all six display the expected effect (i) consistently across 368 all thresholds: in short, the top response is more accurate when the activations of concept neurons 369 are increased. However, only two models, Llama-2 and the Gemma-2, exhibit effect (ii). Addi-370 tionally, four models, belonging to the BERT and OPT families, align with expectation (iii). Over-371 all, bert-large-uncased and gemma-2-9b adhere to all three expected behaviors across all 372 cases. This happens under restrictive thresholds however ( $t_r = 0.9$  and  $t_c = 0.1$ ), and the four other 373 tested models fail to match all of these expectations.

<sup>1</sup>For example, if each instantiation contains 10 KNs, including 2 perfect conceptual neurons and 8 perfect relational neurons (present in only 1 instantiation and all instantiations, respectively), Figure 2a would display a bar of 200 at the 1 abscissa and a bar of 8 at the 100 abscissa, which would obscure the predominant role of **Relation Neurons**.



Figure 4: (a) Number of KNs shared by language pairs for Llama-2-7b. About a quarter of neurons are shared between two languages. (b) Same for bert-base-multilingual-uncased. (c) Proportion of shared KNs in a relation as a function of the number of languages in the intersection for Llama-2-7b and bert-base-multilingual-uncased.

These mixed results show that classifying KNs into distinct and disentangled roles is not perfect, po-396 tentially due to noise in our methods or in knowledge attribution methods in the first place. Yet, our experiments do indicate that, for certain models, KNs exhibit specific behaviors and manipulating 398 them leads to predictable effects.

### 5.4 DISCUSSION

As anticipated, these experiments underscore the complexity of the internal mechanisms within 402 LLMs, making it impractical to map a single, well-defined function to individual neurons. Many 403 of the identified KNs do not adhere to a clearly defined role and cannot be neatly categorized as 404 encoding either concepts or relations, even within a highly controlled environment like ParaRel. 405 We believe that the polysemantic nature of neurons prevents such precise delineation, which also 406 helps explain the knowledge editing limitations highlighted in prior research. However, contrary 407 to our initial expectations, certain KNs do appear to serve rather specific functions, and this has 408 been experimentally confirmed for some models in boosting experiments. Hence, while the idea 409 that knowledge would be represented entirely in mono-semantic single neurons is unrealistic, the 410 historically associated methods of, e.g., Knowledge Neurons nonetheless detect transparent signal 411 about how knowledge is encoded. KNs are thus a useful tool to pursue the study of knowledge 412 representation in multilingual models too, which we do in the next section.

413 414

415 416

417

418

419

420 421

422

423

424

390

391

392

393 394

397

399 400

401

#### MULTILINGUAL EXPERIMENTS 6

When we learn a new language, we do not learn all facts about the world again, just new ways to express them. That is, there is a central knowledge base, that we can prompt with several languages. In this section we inquire if knowledge is shared across languages in multilingual models too and, if so, what knowledge.

6.1 EXPERIMENTAL SETTINGS

Models For this experiment we studied bert-base-multilingual-uncased (Devlin et al., 2019) and Llama-2-7b. We used a NVIDIA Tesla V100 GPU for BERT and NVIDIA Tesla A100 GPU for Llama 2, both for about one hour per relation and per language.

425 426

Multi-ParaRel We built and release a new dataset Multi-ParaRel, a multilingual version 427 of ParaRel. More details are given in Appendix A. Multi-ParaRel currently includes 10 428 languages: English, French, Spanish, Catalan, Danish, German, Italian, Dutch, Portuguese and 429 Swedish. We also offer a translation and curation pipeline which makes it possible to add more 430 paraphrases and more languages. It has an average of 17 prompts per relation and per language but 431 this value varies (from 9 for German to 19 for English). Each prompt is compatible with autoregres-



Figure 5: Influence of typology on the average overlap coefficient calculated per language pair of Llama-2-7b (left) and bert-base-multilingual-uncased (right).

sive models. After filtering for quality as above, we obtain on average 10 prompts per relation and language.

449 6.2 KNOWLEDGE NEURONS ARE SHARED ACROSS LANGUAGES

Are KNs Bilingual? KNs were calculated separately for each relation and language. A KN is considered shared between two languages if it appears as a KN in both languages for the same relation. We conducted this pairwise analysis across all languages, thereby extending the findings of Chen et al. (2024) to encompass 10 languages.

The results are presented in Figures 4a and 4b. For the Llama-2-7b model, over a quarter of the 455 neurons are shared between any two languages, with this proportion increasing to approximately 456 one-third for bert-base-multilingual-uncased. This represents a significant degree of 457 neuron sharing, especially when considering that bert-base-uncased, for example, has more 458 than  $12 \times 3,072 = 36,864$  neurons in the intermediate layers of its FFNs. To quantify this, note 459 that among these 36, 864 neurons, only 1, 929 are identified as KNs across all relations for English, 460 and 2,195 for French (roughly 5%). If KNs were randomly selected for each language, we would 461 expect around 100 shared neurons between them (5% overlap); however, in reality, 710 neurons are 462 shared. A similar analysis for Llama-2-7b gives even more extreme results: by chance, there 463 should be 2 shared neurons, while in practice 189 are found. Moreover, these numbers represent a lower bound, as some relations were excluded from the prompt filtering process for certain language 464 pairs, effectively reducing the shared KN count for those relations to zero. Thus, the data indicates 465 significant overlap of KNs across languages, suggesting a partially shared mechanism for knowledge 466 retrieval across different language pairs. 467

Are KNs Multilingual? Next, we examine how the number of shared KNs scales with the num-469 ber of languages in the intersection. Figure 4c shows these results for all relations, along with 470 the average behavior. Across all relations, we observe a consistent pattern: the number of shared 471 neurons decays as a function of the form (number of languages)<sup>- $\alpha$ </sup>, with a fitted  $\alpha = 2.04$  for 472 Llama-2-7b. In comparison, if neurons were shared at random, the expected behavior would 473 follow  $\propto p^{\text{number of languages}}$ , where p is the probability of a neuron being a KN (e.g. p = 0.05 for 474 BERT). This demonstrates that KNs are more multilingual than chance, reinforcing the notion of a 475 language-agnostic knowledge retrieval mechanism. Similar to the findings in Section 5, we observe 476 some but few neurons activated for all languages.

477

468

432

433

434

435 436 437

438 439 440

441 442

443

444 445

446

447 448

Are some neurons more Multilingual? Concept Neurons and Relation Neurons were com-478 puted separately for each language and each model. Figure 5 displays the average pairwise overlap 479 coefficient for each neuron type, across various  $t_r$  and  $t_c$  thresholds, alongside with the pairwise 480 overlap coefficient for all KNs. The results reveal a significant difference in overlap between **Con**-481 cept Neurons and Relation Neurons at all threshold levels. However, the direction of this difference 482 varies depending on the model and on the threshold. At the most demanding thresholds (those to 483 the right selecting the purest types), we observe that relational neurons appear to be more bilin-484 gual. Given the variability at other thresholds (in particular for Llama, which is less performant than 485 BERT in this task), we remain cautious about this conclusion.

#### 486 6.3 KNOWLEDGE NEURONS ARE SHARED BETWEEN NATURAL AND UNNATURAL 487 LANGUAGES 488

We have extended the analysis to non-natural languages, in order to deepen the work of Kervadec et al. (2023) in the specific framework of KNs. More specifically, we calculated 10 seeds of Auto-490 prompt (Shin et al., 2020) for each model and each relation of ParaRel and the associated KNs. We then calculated the overlap coefficient between the KNs calculated in this way and those calculated for English at the relationship level. The results are presented in Table 1. This reveals a very large overlap for all models, going up to an almost complete overlap ( $\geq 80\%$ ) for models other 494 than BERT. In the same way that there were important overlaps across natural languages, this new 495 result suggests a similar mechanism of knowledge retrieval even between natural and non-natural 496 languages. It is possible however that there exists a confound here because both Autoprompt and KNs are gradient based. 498

I	Model	bert-base	bert-large	opt-350m	opt-6.7b	Llama-2-7b
I	Avg. Overlap Coeff.	40%	32%	83%	87%	79%

Table 1: Average overlap coefficient of KNs sets computed at the relation level between English and Autoprompt.

#### 7 DISCUSSION

508 While knowledge neurons may be shared across languages, this does not guarantee that they serve 509 the same role in the two languages. A neuron active in both English and French for a given task 510 may perform different overall tasks depending on the language, that is, parallel activation does not 511 equate to shared functionality. Here, we partially controlled for this and narrowed down the role 512 of these neurons by computing intersections across languages and at the relation level. Yet, further 513 work is needed to investigate more intersections, narrowing down the possible roles, at the level of relations, concepts, responses or formats of the prompt. Our method can further help narrow down 514 shared functions, across languages. 515

516 517

518

489

491

492

493

497

503

504 505 506

507

#### 8 CONCLUSION

519 We introduced a typology for knowledge and applied it to the knowledge attribution method pro-520 posed by Dai et al. (2022) to better classify and understand the behavior of Knowledge Neurons 521 (KNs). Notably, our method remains agnostic to the specific knowledge attribution technique used. 522 Coherently with the initial assumptions in the original work, we found that some of these neurons encode specific concepts, but we also found many which do not and instead seem to exhibit a dis-523 tributed role, where multiple neurons share responsibility for encoding concepts within the same 524 relation, or maybe encode the whole relation. We hypothesize that this polysemantic nature of neu-525 rons contributes to the mixed success observed when using KNs for knowledge editing tasks. Yet 526 again, we were able to identify a subset of more specialized neurons, which we categorized as either 527 conceptual (sensitive to a single concept) or relational (sensitive to relationships between concepts). 528 And in some contexts their manual manipulations show the expected effects on downstream tasks. 529 We extended our analysis to multilingual models and found that a significant number of KNs are 530 shared across languages-both in pairwise comparisons and across all 10 languages tested. This in-531 dicates the presence of a shared, language-agnostic knowledge base within multilingual models. To 532 facilitate this research, we created a multilingual dataset of facts and prompts, enriched with paraphrases in 10 languages. Our findings suggest that even a simple method like Knowledge Neurons 533 can provide valuable insights into the benefits of multilingual training. Looking ahead, we aim to 534 further explore how this shared knowledge can be leveraged to improve the integration of new lan-535 guages into existing multilingual models. Our results indicate that it may not be necessary to relearn 536 factual knowledge for each language, which could pave the way for more efficient training strate-537 gies, particularly for low-resource languages. Instead of focusing on exhaustive coverage of world 538 knowledge, future efforts could prioritize data that highlights the unique syntactic and linguistic features of these languages, thus optimizing resource use and improving model performance.

# 540 REFERENCES

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, and et al. Gpt-4 technical report.
(arXiv:2303.08774), March 2024. doi: 10.48550/arXiv.2303.08774. URL http://arxiv.org/abs/2303.08774. arXiv:2303.08774 [cs].

Templeton Adly, Conerly Tom, Marcus Jonathan, Lindsey Jack, Bricken Trenton, Chen Brian, Pearce Adam, Citro Craig, Ameisen Emmanuel, Jones Andy, Cunningham Hoagy, L Turner Nicholas, McDougall Callum, MacDiarmid Monte, Tamkin Alex, Durmus Esin, Hume Tristan, Mosconi Francesco, Freeman C. Daniel, R. Sumers Theodore, Rees Edward, Batson Joshua, Jermyn Adam, Carter Shan, Olah Chris, and Henighan Tom. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. May 2024. URL https: //transformer-circuits.pub/2024/scaling-monosemanticity/.

Roee Aharoni, Melvin Johnson, and Orhan Firat. Massively multilingual neural machine translation. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 3874–3884, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1388. URL https: //aclanthology.org/N19-1388.

- Badr AlKhamissi, Millicent Li, Asli Celikyilmaz, Mona Diab, and Marjan Ghazvininejad. A review on language models as knowledge bases. (arXiv:2204.06031), April 2022. doi: 10.48550/arXiv.2204.06031. URL http://arxiv.org/abs/2204.06031. 99 citations (Semantic Scholar/arXiv) [2024-04-29] 99 citations (Semantic Scholar/DOI) [2024-04-29] arXiv:2204.06031 [cs].
- Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun,
   Mia Xu Chen, Yuan Cao, George Foster, Colin Cherry, Wolfgang Macherey, Zhifeng Chen,
   and Yonghui Wu. Massively multilingual neural machine translation in the wild: Findings
   and challenges. (arXiv:1907.05019), July 2019. doi: 10.48550/arXiv.1907.05019. URL
   http://arxiv.org/abs/1907.05019.
- Mikel Artetxe, Gorka Labaka, and Eneko Agirre. Unsupervised statistical machine translation. In
   *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 3632–3642, Brussels, Belgium, October 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1399. URL https://aclanthology.org/D18-1399.
- David Bau, Steven Liu, Tongzhou Wang, Jun-Yan Zhu, and Antonio Torralba. Rewriting a deep generative model. (arXiv:2007.15646), July 2020a. doi: 10.48550/arXiv.2007.15646. URL http://arxiv.org/abs/2007.15646. arXiv:2007.15646 [cs].
- David Bau, Jun-Yan Zhu, Hendrik Strobelt, Agata Lapedriza, Bolei Zhou, and Antonio Torralba.
  Understanding the role of individual units in a deep neural network. *Proceedings of the National Academy of Sciences*, 117(48):30071–30078, December 2020b. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1907375117. arXiv:2009.05041 [cs].
- Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. Journey to the center of the knowledge neurons: Discoveries of language-independent knowledge neurons and degenerate knowledge neurons. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(1616): 17817–17825, March 2024. ISSN 2374-3468. doi: 10.1609/aaai.v38i16.29735.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models. (arXiv:2210.11416), December 2022. doi: 10.48550/arXiv.2210.11416. URL http://arxiv.org/abs/2210.11416. arXiv:2210.11416 [cs].

593

- Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. Evaluating the ripple effects of knowledge editing in language models. (arXiv:2307.12976), December 2023. doi: 10.48550/arXiv.2307.12976. URL http://arxiv.org/abs/2307.12976. arXiv:2307.12976 [cs].
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 8440–8451, Online, July 2020a. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.747. URL https://aclanthology.org/ 2020.acl-main.747.
- Alexis Conneau, Shijie Wu, Haoran Li, Luke Zettlemoyer, and Veselin Stoyanov. Emerging crosslingual structure in pretrained language models. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 6022–6034, Online, July 2020b. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.536. URL https://aclanthology.org/2020. acl-main.536.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. Knowledge neurons in pretrained transformers. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8493–8502, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.581. URL https://aclanthology.org/ 2022.acl-long.581.
- Nicola De Cao, Wilker Aziz, and Ivan Titov. Editing factual knowledge in language models. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 6491– 6506, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.522. URL https://aclanthology. org/2021.emnlp-main.522.
- Nicola De Cao, Leon Schmid, Dieuwke Hupkes, and Ivan Titov. Sparse interventions in language models with differentiable masking. In Jasmijn Bastings, Yonatan Belinkov, Yanai Elazar, Dieuwke Hupkes, Naomi Saphra, and Sarah Wiegreffe (eds.), *Proceedings of the Fifth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pp. 16–27, Abu Dhabi, United Arab Emirates (Hybrid), December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.blackboxnlp-1.2. URL https://aclanthology.org/2022.blackboxnlp-1.2.
- Ameet Deshpande, Partha Talukdar, and Karthik Narasimhan. When is BERT multilingual? isolating crucial ingredients for cross-lingual transfer. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (eds.), Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 3610–3623, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.264. URL https://aclanthology.org/2022.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/ N19-1423.
- Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Schütze, and Yoav Goldberg. Measuring and improving consistency in pretrained language models. *Transactions of the Association for Computational Linguistics*, 9:1012–1031, 2021a. doi: 10.1162/tacl\_a\_00410. URL https://aclanthology.org/2021.tacl-1.60.

667

- Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Schütze, and Yoav Goldberg. Measuring and improving consistency in pretrained language models. *Transactions of the Association for Computational Linguistics*, 9:1012–1031, 2021b. doi: 10.1162/tacl\_a\_00410. 224 citations (Semantic Scholar/DOI) [2024-04-29].
- 652 Hady Elsahar, Pavlos Vougiouklis, Arslen Remaci, Christophe Gravier, Jonathon Hare, Fred-653 erique Laforest, and Elena Simperl. T-REx: A large scale alignment of natural language with 654 knowledge base triples. In Nicoletta Calzolari, Khalid Choukri, Christopher Cieri, Thierry De-655 clerck, Sara Goggi, Koiti Hasida, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène 656 Mazo, Asuncion Moreno, Jan Odijk, Stelios Piperidis, and Takenobu Tokunaga (eds.), Proceed-657 ings of the Eleventh International Conference on Language Resources and Evaluation (LREC 658 2018), Miyazaki, Japan, May 2018. European Language Resources Association (ELRA). URL 659 https://aclanthology.org/L18-1544.
- Kavier Garcia, Aditya Siddhant, Orhan Firat, and Ankur Parikh. Harnessing multilinguality in unsupervised machine translation for rare languages. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1126–1137, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.89. URL https://aclanthology.org/2021.naacl-main.89.
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 5484–5495, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.446. URL https://aclanthology.org/2021.emnlp-main.446.
- Nicolas Guerin, Shane Steinert-Threlkeld, and Emmanuel Chemla. The impact of syntactic and semantic proximity on machine translation with back-translation. (arXiv:2403.18031), March 2024. doi: 10.48550/arXiv.2403.18031. URL http://arxiv.org/abs/2403.18031. 0 citations (Semantic Scholar/arXiv) [2024-05-16] 0 citations (Semantic Scholar/DOI) [2024-05-16] arXiv:2403.18031 [cs].
- Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models. (arXiv:2301.04213), October 2023. doi: 10.48550/arXiv.2301.04213. URL http: //arxiv.org/abs/2301.04213. 46 citations (Semantic Scholar/arXiv) [2024-02-19] 46 citations (Semantic Scholar/DOI) [2024-02-19] arXiv:2301.04213 [cs].
- Linyang He, Peili Chen, Ercong Nie, Yuanning Li, and Jonathan R. Brennan. Decoding probing: Revealing internal linguistic structures in neural language models using minimal pairs. (arXiv:2403.17299), March 2024. doi: 10.48550/arXiv.2403.17299. URL http://arxiv.org/abs/2403.17299. arXiv:2403.17299 [cs, q-bio].
- Jing Huang, Atticus Geiger, Karel D'Oosterlinck, Zhengxuan Wu, and Christopher Potts. Rigorously assessing natural language explanations of neurons. (arXiv:2309.10312), September 2023. doi: 10.48550/arXiv.2309.10312. URL http://arxiv.org/abs/2309.10312. 5 citations (Semantic Scholar/arXiv) [2024-03-24] 5 citations (Semantic Scholar/DOI) [2024-03-24] arXiv:2309.10312 [cs].
- Kuan-Hao Huang, Wasi Ahmad, Nanyun Peng, and Kai-Wei Chang. Improving zero-shot crosslingual transfer learning via robust training. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 1684–1697, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main. 126. URL https://aclanthology.org/2021.emnlp-main.126.
- Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. How can we know when language models know? on the calibration of language models for question answering. *Transactions of the*

738

742

743

749

702 Association for Computational Linguistics, 9:962–977, 2021. doi: 10.1162/tacl\_a\_00407. URL 703 https://aclanthology.org/2021.tacl-1.57. 704

Karthikeyan K, Zihan Wang, Stephen Mayhew, and Dan Roth. Cross-lingual ability of multilingual 705 bert: An empirical study. (arXiv:1912.07840), February 2020. doi: 10.48550/arXiv.1912.07840. 706 URL http://arxiv.org/abs/1912.07840. arXiv:1912.07840 [cs].

- 708 Nora Kassner, Philipp Dufter, and Hinrich Schütze. Multilingual LAMA: Investigating knowl-709 edge in multilingual pretrained language models. In Paola Merlo, Jorg Tiedemann, and Reut 710 Tsarfaty (eds.), Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pp. 3250-3258, Online, April 2021. Asso-711 ciation for Computational Linguistics. doi: 10.18653/v1/2021.eacl-main.284. URL https: 712 //aclanthology.org/2021.eacl-main.284. 713
- 714 Corentin Kervadec, Francesca Franzon, and Marco Baroni. Unnatural language processing: How do 715 language models handle machine-generated prompts? pp. 14377–14392, Singapore, December 716 2023. doi: 10.18653/v1/2023.findings-emnlp.959. URL https://aclanthology.org/ 717 2023.findings-emnlp.959.
- 718 Takeshi Kojima, Itsuki Okimura, Yusuke Iwasawa, Hitomi Yanaka, and Yutaka Matsuo. On the mul-719 tilingual ability of decoder-based pre-trained language models: Finding and controlling language-720 specific neurons. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), Proceedings of the 721 2024 Conference of the North American Chapter of the Association for Computational Lin-722 guistics: Human Language Technologies (Volume 1: Long Papers), pp. 6919–6971, Mexico 723 City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024. 724 naacl-long.384. URL https://aclanthology.org/2024.naacl-long.384.
- 725 Sneha Kudugunta, Ankur Bapna, Isaac Caswell, and Orhan Firat. Investigating multilingual NMT 726 representations at scale. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), Pro-727 ceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and 728 the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 729 1565–1575, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 730 10.18653/v1/D19-1167. URL https://aclanthology.org/D19-1167.
- Yair Lakretz, German Kruszewski, Theo Desbordes, Dieuwke Hupkes, Stanislas Dehaene, and 732 Marco Baroni. The emergence of number and syntax units in LSTM language models. In 733 Jill Burstein, Christy Doran, and Thamar Solorio (eds.), Proceedings of the 2019 Conference 734 of the North American Chapter of the Association for Computational Linguistics: Human Lan-735 guage Technologies, Volume 1 (Long and Short Papers), pp. 11–20, Minneapolis, Minnesota, 736 June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1002. URL 737 https://aclanthology.org/N19-1002.
- Guillaume Lample and Alexis Conneau. Cross-lingual language model pretraining. 739 (arXiv:1901.07291), January 2019. doi: 10.48550/arXiv.1901.07291. URL http://arxiv. 740 org/abs/1901.07291. arXiv:1901.07291 [cs]. 741
- Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc'Aurelio Ranzato. Unsupervised machine translation using monolingual corpora only. February 2018. URL https: //openreview.net/forum?id=rkYTTf-AZ. 744
- 745 Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, 746 and Luke Zettlemoyer. Multilingual denoising pre-training for neural machine translation. Trans-747 actions of the Association for Computational Linguistics, 8:726–742, November 2020. ISSN 748 2307-387X. doi: 10.1162/tacl\_a\_00343.
- Christopher D. Manning, Kevin Clark, John Hewitt, Urvashi Khandelwal, and Omer Levy. Emer-750 gent linguistic structure in artificial neural networks trained by self-supervision. Proceedings of 751 the National Academy of Sciences, 117(48):30046–30054, December 2020. doi: 10.1073/pnas. 752 1907367117. 753
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual 754 associations in gpt. (arXiv:2202.05262), January 2022. URL http://arxiv.org/abs/ 755 2202.05262. arXiv:2202.05262 [cs].

- Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. (arXiv:2210.07229), August 2023. URL http://arxiv.org/ abs/2210.07229. 171 citations (Semantic Scholar/arXiv) [2024-03-01] arXiv:2210.07229 [cs].
- Fric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. Fast model editing at scale. (arXiv:2110.11309), June 2021. URL http://arxiv.org/abs/2110.11309. 188 citations (Semantic Scholar/arXiv) [2024-04-29] arXiv:2110.11309 [cs].
- Fric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. Memory based model editing at scale. (arXiv:2206.06520), June 2022. doi: 10.48550/arXiv.2206.06520.
   URL http://arxiv.org/abs/2206.06520. 125 citations (Semantic Scholar/arXiv)
   [2024-03-04] arXiv:2206.06520 [cs].
- 768 Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven 769 Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, 770 Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert 771 Webson, Edward Raff, and Colin Raffel. Crosslingual generalization through multitask fine-772 tuning. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), Proceedings of the 773 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 15991–16111, Toronto, Canada, July 2023. Association for Computational Linguis-774 tics. doi: 10.18653/v1/2023.acl-long.891. URL https://aclanthology.org/2023. 775 acl-long.891. 776
- Jingcheng Niu, Andrew Liu, Zining Zhu, and Gerald Penn. What does the knowledge neuron thesis have to do with knowledge? October 2023. URL https://openreview.net/forum?
   id=2HJRwwbV3G.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. Language models as knowledge bases? In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 2463–2473, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1250. URL https://aclanthology. org/D19–1250.
- Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim Rocktäschel, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. How context affects language models' factual predictions. (arXiv:2005.04611), May 2020. doi: 10.48550/arXiv.2005.04611. URL http: //arxiv.org/abs/2005.04611. 166 citations (Semantic Scholar/arXiv) [2024-04-29] arXiv:2005.04611 [cs].
- Telmo Pires, Eva Schlinger, and Dan Garrette. How multilingual is multilingual BERT? In
   Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4996–5001, Florence, Italy, July
   2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1493. URL https:
   //aclanthology.org/P19-1493.
- Alec Radford, Rafal Jozefowicz, and Ilya Sutskever. Learning to generate reviews and discovering sentiment. (arXiv:1704.01444), April 2017. doi: 10.48550/arXiv.1704.01444. URL http://arxiv.org/abs/1704.01444. arXiv:1704.01444 [cs].
- Ryokan Ri and Yoshimasa Tsuruoka. Pretraining with artificial language: Studying transferable knowledge in language models. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7302–7315, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.504. URL https://aclanthology.org/2022.acl-long.504.
- Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the parameters of a language model? In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu

856

(eds.), Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 5418–5426, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.437. URL https://aclanthology.org/2020.
emnlp-main.437.

- Anna Rogers, Olga Kovaleva, and Anna Rumshisky. A primer in BERTology: What we know about how BERT works. *Transactions of the Association for Computational Linguistics*, 8:842–866, 2020. doi: 10.1162/tacl\_a\_00349. URL https://aclanthology.org/2020.tacl-1. 54.
- Tara Safavi and Danai Koutra. Relational World Knowledge Representation in Contextual Language Models: A Review. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott
  Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 1053–1067, Online and Punta Cana, Dominican Republic, November
  2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.81. URL
  https://aclanthology.org/2021.emnlp-main.81.
- Shibani Santurkar, Dimitris Tsipras, Mahalaxmi Elango, David Bau, Antonio Torralba, and Aleksander Madry. Editing a classifier by rewriting its prediction rules. (arXiv:2112.01008), December 2021. doi: 10.48550/arXiv.2112.01008. URL http://arxiv.org/abs/2112.01008. arXiv:2112.01008 [cs].
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, and et al. Bloom: A 176b-parameter open-access multilingual language model. (arXiv:2211.05100), June 2023. doi: 10.48550/arXiv.2211.05100. URL http://arxiv.org/abs/2211.05100. arXiv:2211.05100 [cs].
- Rico Sennrich, Barry Haddow, and Alexandra Birch. Improving neural machine translation models
  with monolingual data. In Katrin Erk and Noah A. Smith (eds.), *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 86–96,
  Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/
  P16-1009. URL https://aclanthology.org/P16-1009.
- Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 4222–4235, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.346.
  URL https://aclanthology.org/2020.emnlp-main.346.
- Jasdeep Singh, Bryan McCann, Richard Socher, and Caiming Xiong. BERT is not an interlingua and the bias of tokenization. In Colin Cherry, Greg Durrett, George Foster, Reza Haffari, Shahram Khadivi, Nanyun Peng, Xiang Ren, and Swabha Swayamdipta (eds.), *Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP (DeepLo 2019)*, pp. 47–55, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-6106. URL https://aclanthology.org/D19-6106.
- Anton Sinitsin, Vsevolod Plokhotnyuk, Dmitriy Pyrkin, Sergei Popov, and Artem Babenko. Editable neural networks. (arXiv:2004.00345), July 2020. doi: 10.48550/arXiv.2004.00345. URL http://arxiv.org/abs/2004.00345. 110 citations (Semantic Scholar/arXiv) [2024-04-29] arXiv:2004.00345 [cs, stat].
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. (arXiv:1703.01365), June 2017. doi: 10.48550/arXiv.1703.01365. URL http://arxiv. org/abs/1703.01365. arXiv:1703.01365 [cs].
- Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Xin Zhao, Furu
  Wei, and Ji-Rong Wen. Language-specific neurons: The key to multilingual capabilities in large
  language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of*the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long

865

866

867

Papers), pp. 5701–5715, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.309. URL https://aclanthology.org/2024. acl-long.309.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhu-868 patiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Char-870 line Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, 871 Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, 872 Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchi-873 son, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, 874 Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, 875 Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Wein-876 berger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, 877 Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen 878 Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, 879 Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju-yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir 882 Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leti-883 cia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, 885 Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, 889 Reza Ardeshir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah 890 Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, 891 Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Ko-892 cisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren 893 Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao 894 Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris 895 Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine 896 Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskava, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. Gemma 2: Improving open language models at a practical size. (arXiv:2408.00118), August 2024. doi: 10.48550/arXiv.2408.00118. URL 899 http://arxiv.org/abs/2408.00118. arXiv:2408.00118 [cs]. 900

901 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-902 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, 903 Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy 904 Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, 905 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel 906 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar 907 Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan 908 Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen 909 Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan 910 Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, 911 Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-912 tuned chat models. (arXiv:2307.09288), July 2023. doi: 10.48550/arXiv.2307.09288. URL 913 http://arxiv.org/abs/2307.09288. arXiv:2307.09288 [cs]. 914

Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and
 Stuart Shieber. Investigating gender bias in language models using causal mediation analysis. In
 *Advances in Neural Information Processing Systems*, volume 33, pp. 12388–12401. Curran As-

sociates, Inc., 2020. URL https://proceedings.neurips.cc/paper/2020/hash/ 92650b2e92217715fe312e6fa7b90d82-Abstract.html.

- Weixuan Wang, Barry Haddow, Minghao Wu, Wei Peng, and Alexandra Birch. Sharing matters:
   Analysing neurons across languages and tasks in llms, 2024. URL https://arxiv.org/ abs/2406.09265.
- Xiaokai Wei, Shen Wang, Dejiao Zhang, Parminder Bhatia, and Andrew Arnold. Knowledge enhanced pretrained language models: A compreshensive survey. (arXiv:2110.08455), October 2021. URL http://arxiv.org/abs/2110.08455. 28 citations (Semantic Scholar/arXiv) [2024-04-29] arXiv:2110.08455 [cs].
- Shijie Wu and Mark Dredze. Beto, bentz, becas: The surprising cross-lingual effectiveness of BERT. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 833–844, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1077.
  URL https://aclanthology.org/D19-1077.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mT5: A massively multilingual pre-trained text-to-text transformer. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 483–498, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.41. URL https://aclanthology.org/2021.naacl-main.41.
  - Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. Opt: Open pre-trained transformer language models. (arXiv:2205.01068), June 2022. doi: 10.48550/arXiv.2205.01068. URL http://arxiv.org/abs/2205.01068. arXiv:2205.01068 [cs].

972



Figure 6: In mLAMA, the number of triplets available varies widely across the different languages.

987 988

997

998

1001

1004

1008

1009

NEW DATASET: MULTI-PARAREL Α

**Creation Procedure** To build the Multi-ParaRel dataset, we used our augmented autoregressive version of ParaRel and mLAMA. The goal is to translate a template such as The capital of [X] is [Y]. The problem is that translators are confused by the presence of placeholders [X] and [Y], often resulting in translation errors. To overcome the difficulty, we instantiated [X] and [Y], translated the whole sentence with these specific instances, and replaced the instantiations back with placeholders. To do so, we used mLAMA, which contains triplets for over 53 languages.

For example, consider the translation from English into French of the template:

The capital of [X] is [Y]

999 We use the English triplet *Great Britain, capital of, London>* to obtain the sentence: 1000

The capital of Great Britain is London

1002 This sentence is then translated into French: 1003

La capitale de la Grande Bretagne est Londres

1005 Then using the French version of the original triplet (*<la Grande Bretagne, capital of, Londres>*), we can find and replace the entity elements of the triplet with placeholders X and Y, resulting in 1007 the new template:

La capitale de [X] est [Y]

With this overall idea, we can now provide more detail. First of all, such a protocol requires associ-1010 ated triplets in mLAMA from one language to another. However, mLAMA has many more triplets in 1011 English than in other languages (see Figure 6), and some triplets are language-specific and therefore 1012 cannot be associated with triplets in other languages. We therefore looked into a common English-1013 Target language subset. Then, to avoid translation errors, problems linked to gendered determinants 1014 and redundancy (two different templates in English but translated identically in the target language), 1015 we used a voting system. Each template was translated 30 times, using 30 triplets. Each transla-1016 tion is assigned a score, which is the number of times the template has been obtained out of the 30 1017 triplets. The template with the highest score is then retained, provided that (i) it is autoregressive, 1018 (ii) it has not already been selected and (iii) it is in the top 5 translations.

1019 As a translation model, we used Meta's SeamlessM4T and, more specifically, the Huggingface 1020 implementation<sup>2</sup>. We used an NVIDIA Tesla V100 GPU for inference. 1021

1022 **Statistics and Exemples** Table 3 provides examples of translated templates from different lan-1023 guages and relations. The average number of templates obtained per relationship for each language 1024 is: 1025

<sup>2</sup>https://huggingface.co/facebook/seamless-m4t-large

g templates	Language
10	Catalan
19	Catalan
15	Danish
17	Dutch
19	English
14	French
9	German
19	Italian
19	Portuguese
19	Spanish
16	Swedish
ige Values	Table 2: La

Relation	English	Spanish	French
P36	The capital of [X] is [Y]	La capital de [X] es [Y]	La capitale de [X] est [Y]
	[X], which has the capital [Y]	[X], que tiene la capital [Y]	[X], dont la capitale est [Y]
P106	The occupation of [X] is [Y]	La ocupación de [X] es [Y]	La profession de [X] est [Y]
	[X] works as [Y]	[X] trabaja como [Y]	[X] travaille comme [Y]
P1001	[X] counts as a legal term in [Y]	[X] cuenta como término legal en [Y].	[X] est un terme légal en [Y]
	[X] is a valid legal term in [Y]	[X] es un término legal válido en [Y].	[X] est un terme juridique valide en [Y]

Table 3: Examples of templates from Multi-ParaRel

**Quality Analysis** To judge the quality of our dataset, we asked a native speaker of French and a native speaker of Spanish to rate the resulting templates in three categories: fluent, weird, ungrammatical. For French 88% are correct, 7% weird and 5% are ungrammatical. For Spanish: 78% of sentences are fluent, 10% weird and 12% are ungrammatical. Although imperfect, Multi-ParaRel coupled with a less efficient filtering of prompts gives very good results on mLAMA.

#### В FULL RESULTS

First we provide an overview of all the models behavior with respect to our expectations in Table 4. We also add a control experiment for the BERT family where we conducted the same boosting experiments but sampling KNs randomly within the relation for the Concept Neurons and across relations for Relation Neurons. The goal of such a control is to test the specificity of identified KNs. Results are in Table 5. We see that the effects are destroyed when looking at randomly selected KNs. 

Second, we provide all graphs computed for all models and relations concerning the distinction between concept and relation neurons. This corresponds to the results as presented in Section 5.2, Figure 2, also showing all relations each time. Second, we provide all graphs corresponding to the boosting experiments (Section 5.3, Figure 3). 

Model	Expectation (i)	Expectation (ii)	Expectation (iii)
bert-base-uncased	Yes	No	Yes
bert-large-uncased	Yes	Yes	Yes
opt-350m	Yes	No	Yes
opt-6.7b	Yes	No	Yes
Llama-2-7b	Yes	Yes	No
gemma-2-9b	Yes	Yes	No

1096Table 4: Overview of boosting results for all models. Expectations are: (i) there will be a marked1097increase (or decrease) in precision at rank k=1 when the activations of Concept Neurons are doubled1098(or nullified), with the effect diminishing as k increases, (ii) the effect of Relation Neurons on P@k1099will be weaker than that of Concept Neurons, as precision is primarily sensitive to the correct1100response, (iii) Relation Neurons will play a more significant role, as these neurons should be more1101likely to favor the correct category (e.g., *capitals*), even if it does not boost the correct answer1102specifically.

Model	Expectation (i)	Expectation (ii)	Expectation (iii)
bert-base-uncased	No	No	Yes but effect $10 \times$ smaller
bert-large-uncased	No	No	No

Table 5: Overview of boosting results for the control experiment.















