On Relation-Specific Neurons in Large Language Models

Anonymous ACL submission

Abstract

001

002

005

011

012

015

017

022

034

042

In large language models (LLMs), certain neurons can store distinct pieces of knowledge learned during pretraining. While knowledge typically appears as a combination of relations and entities, it remains unclear whether some neurons focus on a relation itself - independent of any entity. We hypothesize such neurons detect a relation in the input text and guide generation involving such a relation. To investigate this, we study the Llama-2 family on a chosen set of relations with a statisticsbased method. Our experiments demonstrate the existence of relation-specific neurons. We measure the effect of selectively deactivating candidate neurons specific to relation r on the LLM's ability to handle (1) facts whose relation is r and (2) facts whose relation is a different relation $r' \neq r$. With respect to their capacity for encoding relation information, we give evidence for the following three properties of relation-specific neurons. (i) Neuron cumulativity. The neurons for r present a cumulative effect so that deactivating a larger portion of them results in the degradation of more facts in r. (ii) Neuron versatility. Neurons can be shared across multiple closely related as well as less related relations. Some relation neurons transfer across languages. (iii) Neuron interference. Deactivating neurons specific to one relation can improve LLM generation performance for facts of other relations.

1 Introduction

Large text corpora like Wikipedia contain abundant factual knowledge. LLMs, pretrained on such corpora, can function as knowledge bases that retrieve information and generate text involving factual content (Petroni et al., 2019; Jiang et al., 2020). Recent studies suggest that some knowledge is parameterized by LLMs (Dai et al., 2022; Geva et al., 2023), 039 especially within the feed-forward layers of the Transformer architecture (Vaswani et al., 2017), which act as key-value memory (Geva et al., 2021). Factual knowledge is often expressed as a relational fact in triple form: subject, relation, and object, e.g., (NVIDIA, company_ceo, Jensen Huang). However, it remains unclear whether each fact is stored and processed separately through knowledge neurons (Dai et al., 2022), or if there exist relationspecific neurons that focus on the relation itself and guide generating the object once the subject and relation of a triple have been detected.

043

045

047

049

051

054

055

057

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

In this work, we examine the existence of relation-specific neurons in decoder-only LLMs. Our study focuses on the Llama-2 family (7B and 13B) (Touvron et al., 2023) and examines factual knowledge grouped into 12 types of relations. To pinpoint relation-specific neurons for these relations, we adopt the neuron identification method proposed by Cuadros et al. (2022), which identifies the neurons that are uniquely activated in one group of sentences (positive examples) while not in another (negative examples). Kojima et al. (2024) successfully applied this method to uncover language-specific neurons. Following this line of work, we construct zero-shot prompts featuring a specific relation for the positive examples and prompts with other relations for the negative examples. Neurons whose activation patterns are positively correlated with positive examples are regarded as relation-specific neurons.

To understand the impact of these neurons, we perform question-answering on new held-out prompts. These prompts for each relation share the same relation as the positive examples used for identification but have no entity overlap; this disentangles the effects of entities and relations. For each relation, we compare performance between the original model and the model in which the neurons for that relation are deactivated - intra-relation results. We also study how deactivating the neurons for one relation influences performance on other relations - inter-relation results. Our experiments reveal several key properties of these neurons:

Neuron cumulativity. These neurons present a cumulative effect so that deactivating a larger portion of them results in the degradation of more facts, suggesting LLM distributes relational knowledge across neurons in a manner that can vary significantly from fact to fact. This property aligns with the evidence of the existence of redundant or self-repair neurons (Dalvi et al., 2020; McGrath et al., 2023; He et al., 2024). Our analysis suggests the frequency of a fact in the pretraining data is associated with its sensitivity to a given subset.

084

090

092

097

099

101

102

103

104

105

106

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

129

130

131

132

133

Neuron versatility. Because the total number of neurons is finite while the number of possible relations is vast, some neurons strongly associate with multiple relations. Surprisingly, these relations need not be closely linked – two "less related" relations can share a group of neurons, leading to performance drops in both relations if those neurons are deactivated. Such neurons also generalize across languages. This property aligns with recent findings showcasing that some neurons are shared across languages and tasks (Wang et al., 2024a; Tang et al., 2024; Kojima et al., 2024).

Neuron interference. Some relation-specific neurons appear to "confuse" the model when it processes other relations. Deactivating such neurons can yield improved performance on these other relations. This property aligns with broader evidence that *sub-networks* or *circuits* within LLMs may serve several different functional roles (Wang et al., 2023a; Bayazit et al., 2024; Mondorf et al., 2024).

2 Methodology

2.1 Dataset Manipulation

We use the factual knowledge dataset from Hernandez et al. (2024) for this research, which contains 25 relations. Each relation has a different number of facts. Each fact can be represented as a *subjectrelation-object* triple (s, r_i, o) . We only consider relations that have more than 300 facts to ensure the reliability of our findings. This results in **12** relations. We refer to the set of triples for relation r_i as \mathcal{D}_{r_i} . We then perform the following steps for each relation r_i to construct the data used to identify its relation-specific neurons.

Step 1: Creating Evaluation Data. For each triple set \mathcal{D}_{r_i} , we randomly select **50 triples** as a held-out set for evaluation. (cf. §2.3). We refer to the selected triples as $\mathcal{D}_{r_i}^{\text{eva}}$ (for evaluation) and the remaining triples as $\mathcal{D}_{r_i}^{\text{det}}$ (for detection).

Step 2: Formulating Prompts. For each

triple (s, r_i, o) in $\mathcal{D}_{r_i}^{det}$, we create prompts concerning the subject s and the relation r_i using the templates provided by Hernandez et al. (2024). For example, we construct a zero-shot prompt "The CEO of NVIDIA is? Answer:" for the triple (NVIDIA, company_CEO, Jensen Huang) with an expected answer "Jensen Huang". We also create prompts for $\mathcal{D}_{r_i}^{eva}$ in the same way. We refer to the resulting prompt sets as $\mathcal{P}_{r_i}^{det}$ and $\mathcal{P}_{r_i}^{eva}$.

134

135

136

137

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

159

160

161

162

163

164

165

166

168

170

171

172

173

174

175

176

177

178

180

181

Step 3: Validating Prompts. We hypothesize that the model will leverage relation-specific neurons to generate the correct answer, i.e., the object. Therefore, such neurons should be "fired" in those prompts for which **the model answers correctly**. We feed the prompt to the model and set the maximum length of the generation to be 2. We then check if the generated tokens are the same as the object (first two tokens): if they are the same, we regard the output as being correct. We exclude prompts that the model answers wrongly from $\mathcal{P}_{r_i}^{det}$.

2.2 Relation-Specific Neuron Identification

Following (Cuadros et al., 2022), we identify the relation-specific neurons using statistical association measures. This method assigns a score for each neuron, representing its level of "expertise" in distinguishing a specific relation from others.

Defining Neurons. A neural network, or specifically a Transformer (Vaswani et al., 2017), consists of many weight matrices. For a single weight matrix $\boldsymbol{W} \in \mathbb{R}^{d_1 \times d_2}$, we define a neuron as a column, mapping a representation from \mathbb{R}^{d_1} to \mathbb{R} . We assign a unique index $m \in M$ to each neuron and investigate its output value. We only consider the neurons in **feed-forward networks (FFNs)**, i.e., neurons in up_proj, gate_proj, and down_proj, since previous studies have shown that knowledge is mostly stored there (Dai et al., 2022).

Grouping Prompts. For each relation r_i , we collect positive and negative examples. Specifically, we regard $\mathcal{P}_{r_i}^{\text{det}}$ as positive examples and randomly sample $4 \times |\mathcal{P}_{r_i}^{\text{det}}|$ prompts from the prompt sets of other relations as negative examples.¹ For simplicity, the positive and negative examples selected for relation r_i are referred to as $\mathcal{E}_{r_i}^+$ and $\mathcal{E}_{r_i}^-$ respectively. The final data used to detect the relation-specific neurons for relation r_i is then $\mathcal{E}_{r_i} = \mathcal{E}_{r_i}^+ + \mathcal{E}_{r_i}^-$. For each individual example $e_{r_i}^j$, we assign a binary label $b_{r_i}^j$ to it: 1 if $e_{r_i}^j \in \mathcal{E}_{r_i}^+$ and 0 otherwise.

¹We also sample negative examples with different seeds in our preliminary experiments. The neurons for each relation show little change, suggesting the stability of the neurons.

Model	#Layers	#Neurons (FFNs)	#Neurons (total)
Llama-2-7B	32	835,584	1,359,872
Llama-2-13B	40	1,310,720	2,129,920

Table 1: LLama-2 model statistics

182

183

185

188

189

190

191

192

193

194

195

197

198

199

205

206

210

211

213

214

215

216

217

218

219

221

222

225

Neuron Output Values. Let $o_{r_i}^{m,j,n}$ be the output value of neuron m for the n-th token in $e_{r_i}^j$ when feeding the example to the model. Following Kojima et al. (2024), we average the outputs over tokens to form the final output value of neuron m for the entire example $e_{r_i}^j$: $o_{r_i}^{m,j} = \frac{1}{T} \sum_{n=1}^{T} o_{r_i}^{m,j,n}$, where T is the number of effective tokens in $e_{r_i}^j$ – the output values of [PAD] tokens are regarded as noise and therefore ignored in the computation.

Computing Experts. The level of expertise of each neuron for relation r_i is computed by formulating a classification task. Specifically, we regard the output value $o_{r_i}^{m,j}$ as the prediction score with $e_{r_i}^j$ as input and $b_{r_i}^j$ as its ground-truth label. In this way, for an individual neuron m, we have the following data: $\{o_{r_i}^{m,j}, b_{r_i}^j\}_{j=1}^{|\mathcal{E}_{r_i}|}$. We then measure this neuron's performance by setting all output values as classification thresholds and comparing the predictions with the ground truth labels. Average precision (AP) is used as the metric (the area under the precision-recall curve). By doing this, we obtain $AP_{r_i}^m$ for all $m \in M$, allowing us to rank them by their level of expertise in differentiating relation r_i from others. The top k neurons are regarded as relation-specific neurons in descending order.

2.3 Controlled Generation

For each relation r_i , we want to investigate the impact of the identified top-k relation-specific neurons. Therefore, we control text generation by overriding the output values with 0 during inference, aiming to deactivate or suppress these neurons. Specifically, we feed $\mathcal{P}_{r_i}^{\text{eva}}$, the prompts from the held-out evaluation prompt set for relation r_i , into the model. During inference, we simply set the output values of all top-k relation-specific neurons for r_i to a constant 0 and let the model generate two tokens, regarded as the model's answer to the question. The answer is then compared to the first two tokens of the correct answer, i.e., the object.

3 Experimental Setup

3.1 Models

We consider the 7B and 13B models from the LLama-2 family (Touvron et al., 2023). As mentioned in §2.2, we consider the neurons in **FFNs**

which account for more than half of neurons in both 7B and 13B, as shown in Table 1. We also report our preliminary results when considering other types of neurons (e.g., up_proj) in §E.

3.2 Datasets

We manipulate the relational knowledge datasets from Hernandez et al. (2024) using the procedure described in §2.1. Recall that we cover 12 relations in our experiments. Prompt sets $\mathcal{P}_{r_i}^{det}$ (for neuron identification) and $\mathcal{P}_{r_i}^{eva}$ (for evaluation) are constructed for each relation. $|\mathcal{P}_{r_i}^{det}|$ varies for different relations.² $\mathcal{P}_{r_i}^{eva}$ is constructed by randomly selecting **50 triples** for each relation. Since these 50 triples are not used when creating $\mathcal{P}_{r_i}^{det}$, this setup ensures **no subject overlap between** $\mathcal{P}_{r_i}^{det}$ **and** $\mathcal{P}_{r_i}^{eva}$ for the same relation r_i . In addition, we ensure **minimal subject overlap across relations** (mostly 0 between $\mathcal{P}_{r_i}^{det}$ and $\mathcal{P}_{r_j}^{det}$). The only exception is between person_mother and person_father, which share a lot of subjects. A detailed entity overlap analysis is presented in §B.

4 Results and Discussion

We apply our identification method to both 7B and 13B models for all 12 relations. We regard the **top 3,000** neurons with the highest AP values as the relation-specific neurons, as for this threshold, we achieve good coverage of relation-specific neurons with a set of neurons that is not too large. We discuss the impact of this meta-parameter in §5.1.

4.1 Identified Relation-Specific Neurons

Distribution Across Layers. We display the distribution of relation-specific neurons across layers in the 7B model in Figure 2 (see §C for the 13B model). Most neurons are located in the model's middle layers. Such a distribution differs from language-specific neurons, which are mostly located in the first and last few layers (Kojima et al., 2024). We hypothesize that relational knowledge requires more than surface-level information that is mainly encoded and processed in the first and last few layers. Therefore, these relation-specific neurons naturally emerge in the middle layers, where the model has integrated enough lexical and syntactic signals to model and process the relation. This finding is consistent with several studies that show that mapping vectors with certain functions can be

231 232

237

238

239

240

241

242

243 244 245

246 247

250 251 252

253

254

248

249

255 256

257

258

261

262

264

265

266

267

268

269

270

 $^{^{2}|\}mathcal{P}_{r_{i}}^{\text{det}}|$ can be different for 7B and 13B models because the number of prompts excluded in the validating prompt step (described in §2.1) can be different.



Figure 1: Intra-relation results. The left (resp. right) figure displays the results of held-out evaluation prompt set $\mathcal{P}_{r_i}^{\text{eva}}$ (resp. identification prompt set $\mathcal{P}_{r_i}^{\text{det}}$). We report the performance of the original model (without any deactivation), e.g., 7b-original, the model with 3,000 random neurons deactivated (averaged over 10 seeds), e.g., 7b-random, and the model with relation-specific neurons deactivated, e.g., 7b-relation.



Figure 2: Distribution of relation-specific neurons across layers. Most are located in the middle layers.

extracted from the middle layers (Merullo et al., 2024; Hernandez et al., 2024; Todd et al., 2024).

272

273

274

275

278

281

287

290

291

293

297

301

Overlap Across Relations. We display the overlap of relation-specific neurons across relations for the 7B model in Figure 3 (13B is in §C). We see that person_mother and person_father share a large share of neurons, possibly due to the large overlap between their subjects, (see in \S B). However, even though there is almost no subject overlap between any other relations, many relations still share some neurons with others. For instance, person_occupation and person_sport_position share 297 neurons, possibly because they are similar relations – a sport is roughly also an occupation. Extensive neuron overlap can also be observed when two relations are mapping from the same type of subjects, e.g., company_ceo and company_hq, or mapping to the same type of objects, e.g., company_ceo and person_father. However, we show in §4.2.2 that a high neuron overlap does not necessarily imply a high level of mutual interference.

4.2 Controlled Generation

For each relation, we set the output values of its identified 3,000 relation-specific neurons to 0, and observe how the deactivation impacts the relation itself and other relations in terms of accuracy.

4.2.1 Intra-Relation Results

In addition to intra-relation results, i.e., deactivating the 3,000 identified relation-specific neurons



Figure 3: Overlap of the relation-specific neurons across 12 relations. For example, 2053 indicates the number of neurons shared between the 3,000 identified neurons for person_father and the 3,000 for person_mother.

for a relation and evaluating the same relation, we also create a baseline by **randomly** deactivating 3,000 neurons in the model. We report the results of the original models, the results of the baseline, and the intra-relation results in Figure 1. 302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

323

324

325

We can observe a clear performance drop on the identification prompt set $\mathcal{P}_{r_i}^{\text{det}}$ when comparing the accuracy of the original model and the model whose relation-specific neurons are deactivated.³ On the other hand, the model with randomly 3,000 deactivated neurons does not show much difference compared with the original model, indicating the 3,000 relation neurons are closely associated with the facts included in $\mathcal{P}_{r_i}^{\text{det}}$. On the evaluation set $\mathcal{P}_{r_i}^{\text{eva}}$, we also observe a notable accuracy drop across models for most relations. Since $\mathcal{P}_{r_i}^{\text{eva}}$ and $\mathcal{P}_{r_i}^{\text{det}}$ do not share any entities, the accuracy drop can only be attributed to the fact that deactivating 3,000 neurons affects the relation itself - the common characteristic between $\mathcal{P}_{r_i}^{\text{eva}}$ and $\mathcal{P}_{r_i}^{\text{det}}$. Therefore, we argue that relation-specific neurons do exist in LLMs. These neurons are entity-irrelevant, rather, they focus on a specific relation.

On the other hand, the accuracy does not drop

 $^{^{3}}$ For some relations, the drop is moderate, e.g., product_company. We show in §5.1 that the drop can become noticeable when we deactivate more than 3,000 neurons.



Figure 4: Inter-relation results for the 7B model (left) and the 13B model (right). The number in cell (r_i, r_j) indicates the accuracy drop of relation r_i when deactivating the relation-specific neurons of relation r_j .

to 0 for any relation (except landmark_country in the 13B model) when its identified relation-327 328 specific neurons are deactivated. This indicates these 3,000 neurons do not equally influence all facts that belong to a certain relation. We therefore 330 have the neuron cumulativity hypothesis: The 331 relation-specific neurons are associated with different facts belonging to the concerned relation. These neurons present a cumulative effect so that deactivating a small number of them (in this case, 335 336 3,000) results in the degradation of some (not all) facts, whereas more facts are affected when a larger portion is disabled. This disparity highlights that LLMs do not uniformly encode all facts that belong to a given relation, but rather distribute relational 340 341 knowledge across neurons in a manner that can vary significantly from fact to fact. We validate this by demonstrating that the accuracy further drops by deactivating more neurons in §5.1. We also show that the sensitivity of a fact to a given population 345 of neurons may correlate with how frequently the fact appears in the pretraining data in §5.3. 347

4.2.2 Inter-Relation Results

357

361

To understand the effect of how these neurons influence the model's ability to answer prompts across multiple relations, we use **accuracy drop** as a metric: acc_drop_{r_i,r_j} = $\frac{\operatorname{acc}_{r_i}^{\operatorname{original}} - \operatorname{acc}_{r_i}^{\operatorname{deactivated} \cdot r_j}}{\operatorname{acc}_{r_i}^{\operatorname{original}}}$, where acc_{r_i}^{\operatorname{original}} is the accuracy of the original model for $\mathcal{P}_{r_i}^{\operatorname{eva}}$ and $\operatorname{acc}_{r_i}^{\operatorname{deactivated} \cdot r_j}$ is the accuracy for $\mathcal{P}_{r_i}^{\operatorname{eva}}$ when the relation-specific neurons of r_j are deactivated. Results are displayed in Figure 4.

When we compare the 7B and 13B models, no consistent pattern emerges across relations. This indicates that, though being trained on the same data, differences in model size and parameter initialization appear to substantially change the functionality of neurons. In particular, most relations in the 13B model are less influenced when neurons of other relations are deactivated than in the 7B model except in the following cases: deactivating neurons of landmark_country strongly affects several other relations concerning the notion of "location"; person_mother and person_occupation are sensitive to the deactivation of neurons of other relations. Despite these divergences, we propose two hypotheses that hold across both models. 362

363

364

365

366

367

370

371

372

374

375

376

377

378

379

381

382

383

388

389

391

392

393

394

395

397

399

Neuron versatility. We observe that deactivating relation-specific neurons for one relation can strongly affect not only that relation but also others, whether they are "closely" or only "loosely" related. For example, disabling person_pro_sport neurons has a large effect on person_sport_position (but not vice versa) in both 7B and 13B models, likely because the model first needs to understand "sport" before inferring "position". Similarly, deactivating person_father neurons reduces accuracy on person_mother, as both share the concept of an immediate parental relationship. Even more loosely related relations can exhibit a clear accuracy drop: deactivating star_constellation neurons affects landmark_continent in both models, possibly because both involve the abstract notion of "location".

Neuron interference. Deactivating the relationspecific neurons of one relation can sometimes **improve** the accuracy for other relations – a phenomenon more pronounced in the 7B model, likely because its smaller parameter space is less capable of isolating different relations. In the 7B model, several relations frequently benefit from this effect: for instance, person_mother improves when neurons from 5 out of 11 other relations – mostly "less related" ones – are deactivated. This effect is also observed for closely related



Figure 5: Influence of deactivating different numbers of relation-specific neurons for each relation. The variation of accuracy on the relation itself and the average accuracy on other relations is shown. Increasing the number clearly affects the relation itself, but the effect on other relations is obvious only until 3,000 or 10,000 neurons.

relations: disabling company_ceo neurons offers a small accuracy boost to company_hq for both models, while deactivating landmark_country neurons benefits landmark_continent in the 7B model. Interestingly, the 13B model shows the opposite effect for landmark_continent when disabling landmark_country, implying that country information can help predict a continent for the larger model. These findings indicate that neuron interference happens across model sizes, but its specific patterns vary – potentially due to factors such as parameter initialization, pretraining data order, or other hyperparameters.

5 Further Analysis

400

401

402 403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

5.1 Influence of the Numbers of Neurons

In this section, we investigate the effect of varying the number of relation-specific neurons on the 7B model (see §C for 13B). Specifically, we consider **ten** values: 10, 50, 200, 500, 1,000, 3,000, 10,000, 20,000, and 50,000. When deactivating varying numbers of neurons for a relation, we report the variation of accuracy for that relation and the **average accuracy** for all other relations in Figure 5. The variation in individual relations is in Figure 16.

Relation-specific neurons of different relations present a varying degree of cumulative effect. Although the accuracy for most relations drops substantially once 3,000 or 10,000 neurons are deactivated, some relations are far more sensitive to a smaller-scale deactivation. For example, disabling only 50 neurons reduces the accuracy from 80% to 50% for star_constellation, and a similar decline occurs for person_plays_instrument. We hypoth-

Relation	#total	#affected
company_ceo	11	2
company_hq	5	0
landmark_continent	6	2
landmark_country	6	0
person_father	6	2
person_mother	7	5
person_occupation	8	1
person_plays_instrument	3	0
person_pro_sport	8	0
person_sport_position	20	11
product_company	9	2
star_constellation	14	0

Table 2: Case study. Out of the prompts (#total) where the model answers correctly when deactivating 1,000 but not when deactivating 3,000 neurons, few of them are affected when deactivating the difference (2,000).

esize that this sensitivity reflects how many distinct objects each relation maps to: if a relation has fewer unique objects, fewer neurons may be critical. By contrast, relations like person_occupation display lower sensitivity, retaining about 20% accuracy even after 50,000 neurons are deactivated.

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

Validation of the cumulative effect. However, it remains unclear whether the further accuracy drop between two thresholds in Figure 5 is driven by the newly deactivated neurons or the cumulative effect of all deactivated neurons. To further investigate our neuron cumulativity hypothesis, we conduct a case study using 1,000 and 3,000 deactivated neurons. Specifically, we identify prompts from $\mathcal{P}_{r_i}^{\text{eva}}$ where the model answers correctly with 1,000 deactivated neurons but fails when 3,000 are deactivated. We then deactivate only the 2,000 additional neurons and measure the number of affected prompts, as shown in Table 2. The results support the cumulative effect for most relations: the degradation of more facts primarily stems from the collective deactivation of all 3.000 neurons.

Deactivating relation-specific neurons has a marginal effect on other relations until certain



Figure 6: Accuracy on 12 relations across 6 languages. The upper bars (resp. lower bars) show the accuracy before (resp. after) the deactivation of 3,000 relationsspecific neurons. Even though these neurons are identified using English, they usually influence other languages, indicating multilinguality of these neurons.

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

thresholds are reached. Typically, these thresholds occur around 3,000 or 10,000, below which the accuracy on other relations remains stable - supporting the choice of 3,000 neurons. Once more neurons are deactivated, other relations also deteriorate, consistent with our neuron versatility hypothesis. This effect is clearer among closely related relations (e.g., person_pro_sport and person_sport_position), as displayed in Figure 16. Even deactivating up to 50,000 neurons seldom reduces other relations to near-zero accuracy, suggesting a high degree of relation-specificity. One exception is company_hq, for which disabling 50,000 neurons causes all relations' accuracies to approach zero - possibly because some of these neurons underlie more general generation capabilities of the model (Sun et al., 2024; Yu et al., 2024).

5.2 Are These Neurons Multilingual?

Recent studies suggest that some neurons encoding factual knowledge or handling specific tasks are multilingual (Stanczak et al., 2022; Zhang et al., 2024; Wang et al., 2024a). A natural question is whether relation-specific neurons – identified solely via English prompts – also function across languages. To explore this, we translate $\mathcal{P}_{r_i}^{\text{eva}}$ to 5 languages: German (**deu**), Spanish (**esp**), French (**fra**), Chinese (**zho**), and Japanese (**jpn**) (see §D for details). We then deactivate the previously identified 3,000 neurons in the 7B model and measure the effect on these languages, as shown in Figure 6.

Although the model's accuracy is generally lower in non-English languages, it still achieves decent results for most relations (except for jpn and zho). Once the neurons for a given relation are deactivated, the accuracy drops across nearly



Figure 7: Relative difference between the average fact frequencies of the group (a) *resilient facts* and (b) *sensi-tive facts* for each relation in 7B (top) and 13B (bottom) models. Resilient facts generally appear more often than sensitive facts in most relations in the pertaining data.

493

494

495

496

497

498

499

500

501

504

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

all languages – supporting our **neuron versatility** hypothesis. Our findings align with recent explanations that LLMs tend to translate the input text from any language into English for task solving in the middle layers based on a shared representation space (Wendler et al., 2024; Dumas et al., 2024; Zhao et al., 2024). As a result, deactivating "English" neurons naturally disrupts this shared space, impairing the model's capability to generalize across languages for the affected relation.

5.3 Fact Frequencies vs. Neuron Cumulativity

We now examine our **neuron cumulativity** hypothesis by asking: *why do some facts show higher sensitivity to a given set of relation neurons than others?* We hypothesize that the frequency of a fact in the pretraining data can be a key factor, as more frequent facts may be memorized more robustly and thus remain less sensitive to deactivation.

Because the pretraining data for Llama 2 is not publicly available, we approximate it using Dolma (Soldaini et al., 2024), a 3 trillion-token opensource corpus. For each relation, we split the facts into two groups: (a) *resilient facts*, for which the 7B (or 13B) model correctly predicts the object **both before and after** deactivating 3,000 relationspecific neurons. (b) *sensitive facts*, for which the model is correct **before but not after** these neurons are deactivated.⁴ We then count how many documents in Dolma contain **both the subject and**

⁴We do not consider other numbers of relation-specific neurons because (1) if #neurons < 3,000, there are not enough facts whose predictions change, and (2) if #neurons > 3,000, facts belonging to other relations will also be influenced a lot.

52

52

527

5

531 532

533 534

535 536

53

5

540

541

542 543

544

545

546

550

553

555

556

557

562

563

566

We

object of each fact, calling this the *fact frequency*.⁵ Finally, we compute the average frequency for resilient and sensitive facts in each relation r_i , denoted respectively as group^(a) and group^(b)

noted respectively as $\operatorname{group}_{r_i}^{(a)}$ and $\operatorname{group}_{r_i}^{(b)}$. Relative difference: $\operatorname{diff}_{r_i} = \frac{\operatorname{group}_{r_i}^{(b)} - \operatorname{group}_{r_i}^{(a)}}{\operatorname{group}_{r_i}^{(b)}}$ for each relation r_i is reported in Figure 7. We find

that resilient facts generally appear more often in Dolma than sensitive facts, with only 3 exceptions in the 7B model and 2 exceptions in the 13B model (note that landmark_country is omitted for the 13B model because no facts fall into group (a)). We evaluate this difference with the Wilcoxon Signed-Rank Test (Woolson, 2005) and obtain *p*-values of respectively 0.11 and 0.03 for the 7B and the 13B models.⁶ These results show that there is a difference (statistically significant in the 13B model at the 5% level) between the two groups, supporting our hypothesis that more frequent facts are generally less sensitive to the deactivation of a given set of relation-specific neurons.

5.4 Relations vs. Concepts

We saw in Figure 3 that the storage of relations is generally well separated, but that there are exceptions. We can view a relation as relating two **concepts**, e.g., company_ceo relates instances of the "subject" concept "company" to instances of the "object" concept "CEO". From this perspective, the exceptions in Figure 3, i.e., cases where a relation r_1 overlaps with a relation r_2 , are generally cases where the concepts of r_1 and r_2 are the same or overlap. For example, company_ceo and company_hq have the same subject concept.

To further explore this hypothesis empirically, we again use the method applied in §2 to relations, but now use it for concepts; that is we identify sets of **concept-specific neurons**. We group the LRE dataset triples by subject concept, resulting in 11 different concepts. We create a set of triples with novel relations such as "can" and "has a", balanced across positive and negative samples. This ensures that the model's completion for a prompt like ("Lincoln has a") depends on the concept instance ("Lincoln"), not on the relation ("has a").

Figure 8 shows the overlap between relation neurons and concept neurons for 13B. Most of the cells with large counts support our hypothe-



Figure 8: Overlap between the top 3000 neurons of relations and concepts in the 13B model.

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

584

585

586

587

588

589

590

591

592

593

594

596

597

598

599

600

601

602

603

604

sis that the overlaps we observe are rooted in relations being representationally associated with their subject and object concepts. Clear examples include company_ceo and its subject concept company; company_hq and its object concept city (assuming that hq is a subcategory of city); and landmark_continent and its subject concept landmark. There is little overlap of person with relations like person_mother, potentially because person is a more general and semantically unspecific concept than the others. Note that several concepts do not match a specific relation, e.g., superhero, and therefore are not strongly associated with any relation. Recall that we picked the concepts according to the availability of annotated data in LRE. However, most identified neurons are only concept neurons or only relation neurons, suggesting that relational and conceptual representations are largely separate.

6 Conclusion

This work highlights the existence of relationspecific neurons in LLMs - neurons that focus on relations rather than entities. Our experiments show these neurons primarily reside in the middle layers and can be shared across multiple relations. Through systematic deactivation, we reveal their influence on both the targeted and other relations, leading to three key hypotheses: neuron cumulativity (deactivating a larger portion of relationspecific neurons results in the degradation of more facts belonging to the concerning relation), neuron versatility (neurons are shared across relations and languages), and neuron interference (neurons from one relation can disrupt the processing of another). These findings shed new light on how LLMs handle relational facts at the neuron level, contributing to the interpretability of LLMs.

⁵We use ElasticSearch API from WIMBD (Elazar et al., 2024) that allows for counting and searching in large corpora.

⁶We use a nonparametric test because the difference across relations does not follow a Gaussian distribution.

Limitations

605

While our findings provide valuable insights, several limitations remain and offer opportunities for 607 future research. First, this work focuses on factual knowledge grouped into 12 relations. Although this selection does not diminish the validity of our findings and hypotheses, it represents a rel-611 atively narrow set of relations. Future work can 612 explore a broader range of relations and analyze how relation-specific neurons behave across a more diverse set of relations. Second, our multilingual 615 analysis includes only five languages. While these 616 languages demonstrate neuron versatility, they do 617 not fully capture linguistic diversity. Future re-618 search could investigate additional languages, par-619 ticularly low-resource ones, to determine whether relation-specific neurons exhibit similar relational functionality across these languages. Lastly, we observe that more frequent facts tend to be more ro-623 bust to the deactivation of relation-specific neurons 624 in both the 7B and 13B models. Fact frequency 625 is approximated using the Dolma corpus (Soldaini et al., 2024) in this study. However, LLama-2 models may incorporate a larger and more diverse pretraining dataset, potentially leading to some discrepancies between these approximated fact frequencies and their actual frequencies. 631

References

634

637

638

641

643

647

651

655

- Omer Antverg and Yonatan Belinkov. 2022. On the pitfalls of analyzing individual neurons in language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net.
- Anthony Bau, Yonatan Belinkov, Hassan Sajjad, Nadir Durrani, Fahim Dalvi, and James R. Glass. 2019.
 Identifying and controlling important neurons in neural machine translation. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net.
- Deniz Bayazit, Negar Foroutan, Zeming Chen, Gail Weiss, and Antoine Bosselut. 2024. Discovering knowledge-critical subnetworks in pretrained language models. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 6549–6583, Miami, Florida, USA. Association for Computational Linguistics.
- Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever, Jan Leike, Jeff Wu, and William Saunders. 2023. Language models can explain neurons in language models.

Xavier Suau Cuadros, Luca Zappella, and Nicholas Apostoloff. 2022. Self-conditioning pre-trained language models. In International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA, volume 162 of Proceedings of Machine Learning Research, pages 4455–4473. PMLR. 656

657

658

659

660

661

662

663

664

665

666

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. Knowledge neurons in pretrained transformers. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8493– 8502, Dublin, Ireland. Association for Computational Linguistics.
- Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, Anthony Bau, and James R. Glass. 2019.
 What is one grain of sand in the desert? analyzing individual neurons in deep NLP models. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 February 1, 2019, pages 6309–6317. AAAI Press.*
- Fahim Dalvi, Hassan Sajjad, Nadir Durrani, and Yonatan Belinkov. 2020. Analyzing redundancy in pretrained transformer models. In *Proceedings of the* 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 4908–4926, Online. Association for Computational Linguistics.
- Clément Dumas, Veniamin Veselovsky, Giovanni Monea, Robert West, and Chris Wendler. 2024. How do llamas process multilingual text? a latent exploration through activation patching. In *ICML 2024 Workshop on Mechanistic Interpretability*.
- Nadir Durrani, Hassan Sajjad, Fahim Dalvi, and Yonatan Belinkov. 2020. Analyzing individual neurons in pre-trained language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4865–4880, Online. Association for Computational Linguistics.
- Yanai Elazar, Akshita Bhagia, Ian Magnusson, Abhilasha Ravichander, Dustin Schwenk, Alane Suhr, Evan Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, Hannaneh Hajishirzi, Noah A. Smith, and Jesse Dodge. 2024. What's in my big data? In The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024. OpenReview.net.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Neel Nanda, Tom Henighan, Scott Johnston, Sheer ElShowk, Nicholas Joseph, Nova DasSarma, Ben Mann, Danny Hernandez, Amanda Askell, Kamal Ndousse, Andy Jones, Dawn Drain, Anna Chen, Yuntao Bai, Deep Ganguli, Liane Lovitt, and 14 others. 2022. Softmax linear units. *Transformer Circuits Thread*.

820

821

822

823

824

825

826

827

Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, and 6 others. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*.

715

716

717

723

725

726

727

728

729

730

731

733

734

735

738

740

741

742

743

744

745

747

748

749

750

751

752

759

763

765

767

770

- Amit Elhelo and Mor Geva. 2024. Inferring functionality of attention heads from their parameters. *Preprint*, arXiv:2412.11965.
- Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. Dissecting recall of factual associations in auto-regressive language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12216–12235, Singapore. Association for Computational Linguistics.
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. Transformer feed-forward layers are keyvalue memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Wes Gurnee, Theo Horsley, Zifan Carl Guo, Tara Rezaei Kheirkhah, Qinyi Sun, Will Hathaway, Neel Nanda, and Dimitris Bertsimas. 2024. Universal neurons in GPT2 language models. *Trans. Mach. Learn. Res.*, 2024.
- Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas.
 2023. Finding neurons in a haystack: Case studies with sparse probing. *Trans. Mach. Learn. Res.*, 2023.
- Shwai He, Guoheng Sun, Zheyu Shen, and Ang Li. 2024. What matters in transformers? not all attention is needed. *Preprint*, arXiv:2406.15786.
- Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin Meng, Martin Wattenberg, Jacob Andreas, Yonatan Belinkov, and David Bau. 2024. Linearity of relation decoding in transformer language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11,* 2024. OpenReview.net.
- Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. 2020. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438.
- Takeshi Kojima, Itsuki Okimura, Yusuke Iwasawa, Hitomi Yanaka, and Yutaka Matsuo. 2024. On the multi-lingual ability of decoder-based pre-trained language models: Finding and controlling language-specific neurons. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 6919–6971, Mexico City, Mexico. Association for Computational Linguistics.

- János Kramár, Tom Lieberum, Rohin Shah, and Neel Nanda. 2024. Atp*: An efficient and scalable method for localizing llm behaviour to components. *Preprint*, arXiv:2403.00745.
- Tom Lieberum, Matthew Rahtz, János Kramár, Neel Nanda, Geoffrey Irving, Rohin Shah, and Vladimir Mikulik. 2023. Does circuit analysis interpretability scale? evidence from multiple choice capabilities in chinchilla. *Preprint*, arXiv:2307.09458.
- Weize Liu, Yinlong Xu, Hongxia Xu, Jintai Chen, Xuming Hu, and Jian Wu. 2024. Unraveling babel: Exploring multilingual activation patterns of LLMs and their applications. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 11855–11881, Miami, Florida, USA. Association for Computational Linguistics.
- Ang Lv, Yuhan Chen, Kaiyi Zhang, Yulong Wang, Lifeng Liu, Ji-Rong Wen, Jian Xie, and Rui Yan. 2024. Interpreting key mechanisms of factual recall in transformer-based language models. *Preprint*, arXiv:2403.19521.
- Thomas McGrath, Matthew Rahtz, Janos Kramar, Vladimir Mikulik, and Shane Legg. 2023. The hydra effect: Emergent self-repair in language model computations. *Preprint*, arXiv:2307.15771.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.
- Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. 2023. Massediting memory in a transformer. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net.
- Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. 2024. Language models implement simple Word2Vec-style vector arithmetic. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5030–5047, Mexico City, Mexico. Association for Computational Linguistics.
- Philipp Mondorf, Sondre Wold, and Barbara Plank. 2024. Circuit compositions: Exploring modular structures in transformer-based language models. *Preprint*, arXiv:2410.01434.
- Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. 2020. Zoom in: An introduction to circuits. *Distill*.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann,

940

941

942

885

886

Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, and 7 others. 2022. Incontext learning and induction heads. *Transformer Circuits Thread*.

829

841

851

858

861

864

870

876

877

878

879

881

883

- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.
- Daking Rai, Yilun Zhou, Shi Feng, Abulhair Saparov, and Ziyu Yao. 2024. A practical review of mechanistic interpretability for transformer-based language models. *Preprint*, arXiv:2407.02646.
 - Hassan Sajjad, Nadir Durrani, and Fahim Dalvi. 2022. Neuron-level interpretation of deep NLP models: A survey. *Transactions of the Association for Computational Linguistics*, 10:1285–1303.
- Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Jha, Sachin Kumar, Li Lucy, Xinxi Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas Muennighoff, and 17 others. 2024. Dolma: an open corpus of three trillion tokens for language model pretraining research. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15725–15788, Bangkok, Thailand. Association for Computational Linguistics.
- Ran Song, Shizhu He, Shuting Jiang, Yantuan Xian, Shengxiang Gao, Kang Liu, and Zhengtao Yu. 2024.
 Does large language model contain task-specific neurons? In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 7101–7113, Miami, Florida, USA. Association for Computational Linguistics.
- Karolina Stanczak, Edoardo Ponti, Lucas Torroba Hennigen, Ryan Cotterell, and Isabelle Augenstein. 2022.
 Same neurons, different languages: Probing morphosyntax in multilingual pre-trained models. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 1589–1598, Seattle, United States. Association for Computational Linguistics.
- Mingjie Sun, Xinlei Chen, J. Zico Kolter, and Zhuang Liu. 2024. Massive activations in large language models. *Preprint*, arXiv:2402.17762.
- Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei, and Ji-Rong Wen. 2024. Language-specific neurons:

The key to multilingual capabilities in large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5701–5715, Bangkok, Thailand. Association for Computational Linguistics.

- Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau. 2024.
 Function vectors in large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11,* 2024. OpenReview.net.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.
- 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 Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 5998–6008.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart M. Shieber. 2020. Investigating gender bias in language models using causal mediation analysis. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.
- Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2023a. Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The Eleventh International Conference on Learning Representations*, *ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. Open-Review.net.
- Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2023b. Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* Open-Review.net.
- Weixuan Wang, Barry Haddow, Minghao Wu, Wei Peng, and Alexandra Birch. 2024a. Sharing matters: Analysing neurons across languages and tasks in llms. *Preprint*, arXiv:2406.09265.
- Yifei Wang, Yuheng Chen, Wanting Wen, Yu Sheng, Linjing Li, and Daniel Dajun Zeng. 2024b. Unveiling factual recall behaviors of large language models through knowledge neurons. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7388–7402, Miami,

995

996

997

998

Florida, USA. Association for Computational Linguistics.

943

952

953

957

958

961

962

963

964

965

966

967

968

969

970

971

972

973

974 975

976

978

979

980

981

982

983

986

990

991

994

- Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. 2024. Do llamas work in English? on the latent language of multilingual transformers. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15366–15394, Bangkok, Thailand. Association for Computational Linguistics.
- Robert F Woolson. 2005. Wilcoxon signed-rank test. Encyclopedia of Biostatistics, 8.
- Mengxia Yu, De Wang, Qi Shan, Colorado Reed, and Alvin Wan. 2024. The super weight in large language models. *Preprint*, arXiv:2411.07191.
- Qinan Yu, Jack Merullo, and Ellie Pavlick. 2023. Characterizing mechanisms for factual recall in language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9924–9959, Singapore. Association for Computational Linguistics.
- Zeping Yu and Sophia Ananiadou. 2024. Neuron-level knowledge attribution in large language models. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 3267–3280, Miami, Florida, USA. Association for Computational Linguistics.
- Xue Zhang, Yunlong Liang, Fandong Meng, Songming Zhang, Yufeng Chen, Jinan Xu, and Jie Zhou. 2024. Multilingual knowledge editing with language-agnostic factual neurons. *Preprint*, arXiv:2406.16416.
- Yiran Zhao, Wenxuan Zhang, Guizhen Chen, Kenji Kawaguchi, and Lidong Bing. 2024. How do large language models handle multilingualism? *Preprint*, arXiv:2402.18815.

A Related Work

Mechanistic interpretability (MI) is a growing subfield of interpretability that aims to understand LLMs by breaking them down into smaller components and fundamental computations. It has gained significant attention for studying how LLMs recall factual knowledge learned during pretraining (Meng et al., 2022; Dai et al., 2022; Geva et al., 2023; Yu et al., 2023; Lv et al., 2024; Wang et al., 2024b). Following Olah et al. (2020); Rai et al. (2024), MI research can be categorized into two areas: the study of **features** and the study of **circuits**, based on the type of decomposed components. Features refer to human-interpretable properties encoded in model representations or represented by model components, such as neurons and attention heads (Elhage et al., 2022; Gurnee et al., 2023). Circuits are subgraphs of the model's computation graph responsible for implementing specific behaviors (Wang et al., 2023b; Elhage et al., 2021).

In this work, we focus on neuron-level featurebased interpretability analysis to localize relationspecific neurons, which are responsible for encoding and recalling specific types of factual knowledge. Existing studies have utilized various approaches for neuron interpretation, each offering unique advantages and limitations (Sajjad et al., 2022; Rai et al., 2024). The visualization method (Olsson et al., 2022; Elhage et al., 2022; Lieberum et al., 2023; Bills et al., 2023; Liu et al., 2024) involves visualizing neuron activations and manually identifying the underlying concept across input text. While being straightforward, it relies heavily on human effort and risks overgeneralization. Statistics-based methods (Bau et al., 2019; Cuadros et al., 2022; Kojima et al., 2024; Yu and Ananiadou, 2024; Tang et al., 2024; Wang et al., 2024b), on the other hand, aggregate activation statistics across data to establish connections between neurons and concepts, identifying patterns through the co-occurrence of neuron activation values and specific input features. Probing-based methods (Dalvi et al., 2019; Durrani et al., 2020; Antverg and Belinkov, 2022; Gurnee et al., 2024) train diagnostic classifiers on neuron activations to identify neurons associated with predefined concepts. These methods are scalable, enabling the discovery of neuron sets across large datasets, though they depend on supervised data annotations. Causationbased methods (Vig et al., 2020; Meng et al., 2022, 2023; Kramár et al., 2024; Song et al., 2024) take a different approach by directly varying the values of specific neurons or components and analyzing changes in model behavior; significant changes indicate the importance of these neurons or components to particular functionalities.

Building on this foundation, our work adopts the statistics-based method proposed by Cuadros et al. (2022) to identify relation-specific neurons – neurons uniquely "fired" for queries concerning facts sharing the same relation. This approach facilitates a scalable and targeted analysis of neuron behavior in relation to factual knowledge recall.

B Entity Overlap Across Relations

We show the number of distinct subjects (resp.1042objects) in each relation and the number of over-1043lapping subjects (resp. objects) between any two1044



Figure 9: Subject (left) and object (right) overlap across 12 relations obtained from the 7B model. The diagonal in each figure shows the number of distinct subjects or objects for each relation. It can be seen that factual knowledge from different relations has almost no entity overlap except for person_mother and person_father, which are mostly celebrities.



Figure 10: Subject (left) and object (right) overlap across 12 relations obtained from the 13B model. The trend is very similar to that in the 7B model: person_mother and person_father share many subjects.



Figure 11: Subject (left) and object (right) overlap across 12 relations in the held-out evaluation prompt set $\mathcal{P}_{r_i}^{\text{eva}}$. Almost no two relations share any subjects or objects.



Figure 12: Distribution of relation-specific neurons across layers for the 13B model. Similar to Figure 2, identified relation-specific neurons are mostly located in the middle layers, except for person_mother.

relations in the identification prompt set $\mathcal{P}_{r_i}^{det}$ of the 7B model and the 13B model in Figure 9 and 10 respectively. Most two relations have no common or very limited overlapping (less than 11) subjects, except for person_mother and person_father, which are mostly celebrities, possibly resulting in extensive neuron overlap between the two relations as we show in §4.1. Similarly, no two relations share many objects. Additionally, we show the number of overlapping entities in the evaluation set $\mathcal{P}_{r_i}^{eva}$ (the 7B and 13B models share the same evaluation set) in Figure 11. The results also show almost no entity overlap across different relations: among all relations, only person_mother



Figure 13: Overlap of the relation-specific neurons across 12 relations in the 13B model. The overlap distribution is not similar to what we observe for the 7B model shown in Figure 3, explaining the difference in inter-relation results (cf. Table 4).

and person_father share one subject and the rest1059of the relations do not share any subject or object1060overlap. The entity analysis suggests that entities1061are not a confounding factor in our experiments1062and the identified relation-specific neurons are only1063concerning the relation itself but not entities.1064

C Analysis On the 13B Model

We perform a similar analysis on the 13B model1066as we do for the 7B model. We first show how the1067identified 3,000 relation-specific neurons are dis-1068



Figure 14: Accuracy on 12 relations across 6 languages from the 13B model. The upper bars (resp. lower bars) show the accuracy before (resp. after) the deactivation of 3,000 relations-specific neurons.

tributed across layers for each relation in Figure 12. The trend is similar to what we observe in the 7B model (cf. Figure 2). Most of the relation-specific neurons are distributed in the middle layers. Then we show the overlap of relation-specific neurons across relations in Figure 13. Surprisingly, the overlap pattern is very different from what we observe in the 7B model. First, it seems that many relations that share a concept of "location" share extensive neurons, e.g., company_hq, landmark_country, landmark_country and star_constellation. This explains the difference in inter-relation results between the models (cf. Figure 4) where we see deactivating neurons of landmark_country significantly influence other relations also concerning location for the 13B model but not for the 7B model.

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1082 1083

1084

1085

1086

1088

1091

1092

1093

1094

1096

1097

1098

1100

1101

1102 1103

1104

We then demonstrate the effect of varving numbers of relation-specific neurons using the same numbers: 10, 50, 200, 500, 1,000, 3,000, 10,000, 20,000, and 50,000. Figure 15 presents the results. The global trend is similar to what we observe for the 7B model: deactivating more neurons results in a further drop in accuracy across all relations. This indicates the neuron cumulativity is universal across models. Relation-specific neurons for most relations present a similar cumulative effect to the 13B model. The original two outliers in the 7B model (person_occupation and person_company where the accuracy does not drop to 0 in the 7B model) even show a plateau, i.e., the accuracy remains almost unchanged or only slightly decreases. This might suggest that facts belonging to these two relations might be wellmemorized by the models and are less sensitive to the deactivation of relation-specific neurons.

specific neurons from the 13B model are also mul-1105 tilingual. We use the same translated prompt sets 1106 as we use for the 7B model. We deactivate the 1107 3,000 neurons identified using English and see how 1108 this affects the performance in other languages: 1109 German (deu), Spanish (esp), French (fra), Chi-1110 nese (zho), and Japanese (jpn). The results are 1111 presented in Figure 14. We observe similar re-1112 sults as from the 7B model: when we deactivate 1113 the relation-specific neurons identified using En-1114 glish prompts, many relations are influenced across 1115 languages, suggesting models with different sizes 1116 also have multilingual relational neurons. We also 1117 see some interesting counterexamples: deactivat-1118 ing landmark_country neurons completely dete-1119 riorates the relation in English but not in German. 1120 This indicates while some neurons have multilin-1121 gual relational functionalities, there are still some 1122 relations dealt with in a language-specific manner. 1123



Figure 15: Influence of deactivating different numbers of relation-specific neurons for each relation (**the 13B model**). The variation of accuracy on the relation itself and the average accuracy on other relations is shown.



Figure 16: Influence of deactivating different numbers of relation-specific neurons in the 7B model for each relation. The variation of accuracy on the relation itself (noted with "*" and a dashed line style) and the accuracy on all other relations is shown in each figure. Similar to Figure 5, increasing the number of neurons clearly affects the relation itself, but the effect on other individual relations does not become clearly noticeable until 3,000–10,000 neurons.



Figure 17: Influence of deactivating different numbers of relation-specific neurons in the 13B model for each relation. The variation of accuracy on the relation itself (noted with "*" and a dashed line style) and the accuracy on all other relations is shown in each figure.

1154

1155

1156

1157

1158

1159

1160

1161

1163

1164

1165

1166

1167

1168

1169

D Translation Process

1125We take a two-step approach to ensure the trans-
lation quality of individual prompts from English
into the target languages across relations.

Translating subject-object pairs. The first step 1128 concerns mapping entities, i.e., subject and object 1129 pairs into the target language. The default way of 1130 doing this is by identifying if the entity is avail-1131 able in Wikidata and the target language using the 1132 Wikidata API.⁷ If the entity of interest is available 1133 in the target language, we directly take the entity 1134 name in that language. If the entity is not available, 1135 we then resort to Google Translate to translate the 1136 entity from English to the target language.⁸. By 1137 performing this step, we obtain the subject-object 1138 pairs in all target languages and all relations. 1139

Translating prompt templates. We take the 1140 prompt templates of different relations written in 1141 English and use Google Translate to translate them 1142 into target languages. We then investigate how 1143 the LLama-2 7B model performs on these prompts 1144 using $\mathcal{P}_{r_i}^{\text{eva}}$ in the target languages. If the model 1145 performs suboptimally (<30% accuracy) for a rela-1146 tion in a specific language, then we manually check 1147 the prompt template in that language and update 1148 the template accordingly until satisfactory accuracy 1149 (>30%) is achieved. For Chinese and Japanese, we 1150 do not ensure more than 30% accuracy because the 1151 models perform very badly for some relations even 1152 if we have tried many prompt templates. 1153

E Influence of Neuron Type

We consider the neurons in the FFNs (including up_proj, gate_proj, and down_proj matrices) as our major setup. In this section, we explore the individual effects of different types of neurons. Specifically, we consider five additional different varieties when selecting the top 3,000 neurons for the 7B model: **all** (neurons in any matrices), **self_attn** (neurons in self-attention matrices), **up_proj** (neurons in up_proj matrices), **gate_proj** (neurons in gate_proj matrices), **down_proj** (neurons in down_proj matrices). We first draw the distribution of the neuron types across relations for variety **all** in Figure 18 and report the inter-relation results in Figure 19 (**all**), 20 (**self_attn**), 21 (**up_proj**), 22 (**gate_proj**), and 23 (**down_proj**).



Figure 18: The distribution of the neuron types in the identified 3,000 neurons for the variety **all** across all relations.



Figure 19: Inter-relation results of the 7B model when considering the neuron type variety as **all**.

According to the results, we observe that simply 1170 considering self_attn does not offer a consistent 1171 accuracy drop for the relation itself (by looking 1172 at the diagonal: some relations are not influenced 1173 too much). This can be explained by the fact the 1174 self_attn is shared across relations (as shown 1175 by Elhelo and Geva (2024)) and facts are mainly 1176 stored in the FFNs. Only considering down_proj 1177 offer similar results as **self_attn**. Interestingly, 1178 deactivating up_proj neurons does not influence 1179 all relations much in general, indicating it does not 1180 make sense to consider up_proj alone. Consider-1181 ing **all** or **gate_proj** neurons offer similar results 1182 compared to considering neurons in FFNs (shown 1183 in Figure 1). However, by considering neurons in 1184 FFNs (i.e., up_proj, gate_proj and down_proj), 1185 we see a more obvious inter-relation accuracy drop 1186 as shown on the diagonal in Figure 1. Therefore, 1187 our additional analysis supports our choice of con-1188 sidering neurons in FFNs. 1189

⁷https://www.wikidata.org/w/api.php

⁸https://translation.googleapis.com/language/ translate/v2



Figure 20: Inter-relation results of the 7B model when considering the neuron type variety as **self_attn**.



Figure 21: Inter-relation results of the 7B model when considering the neuron type variety as **up_proj**.



Figure 22: Inter-relation results of the 7B model when considering the neuron type variety as **gate_proj**.



Figure 23: Inter-relation results of the 7B model when considering the neuron type variety as **down_proj**.

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

F Concept-Specific Neurons

Concept-Relation Overlap in the 7B Model Figure 24 illustrates the overlap between individual relation- and concept-specific neurons in the 7b model. There, the overlap of concepts connected to the abstract notion of "location" and the relations are mostly concentrated on the landmark_country relation in comparison to the 13b model, where they are spread over company_hq, landmark_continent and landmark_country. This aligns with the difference between the 7B and 13B models in terms of their patterns of inter-relation results (cf. Figure 4): deactivating the landmark_country neurons results in a significant accuracy drop in other relations concerning "location" in the 13B model while not in the 7B model. Another difference between both models is that there is more distributed neuron overlap in the 7b model between the subject concept person and all corresponding relations.

Validation of Concept-Specific Neurons The top neurons on a concept are evaluated on a random selection of 100 prompts from the LRE dataset that include the specified concept as a subject. Examples for the concept person are "Tom Hanks's father is named? Answer:", "Hilary Hahn plays the instrument of? Answer:", or "Thomas Mann went to university at? Answer:".

Figure 25 shows the results for the validation1218on these validation prompts for both models with1219the original accuracy score, a baseline that ablates12203000 neurons randomly, and the ablation of 30001221concept-specific neurons. Note that the impact of1222ablating a certain amount of expert neurons varies1223



Figure 24: Overlap between the top 3000 identified neurons for each relation and concept in the 7B model.

between concepts. The observed drop in performance due to the ablation of 3000 neurons for con-1225 cepts like pokemon, superhero, and star is very 1226 1227 large, while accuracy scores of other concepts in the 13b model, such as person appear stable, or 1228 even improve, e.g., presidents. We assume the 1229 neuron cumulativity also applies to the concept-1230 specific neurons. That is, the knowledge on a spe-1231 1232 cific concept is distributed over a much larger population of neurons, and further accuracy drop can 1233 be observed once more concept-specific neurons 1234 are deactivated – similar to what we observe for 1235 relation-specific neurons (cf. Figure 5). As only partial knowledge is withheld from the deactivation 1237 of 3000 concept-specific neurons, this might be too 1238 little knowledge to affect the facts concerning that 1239 concept (substantial knowledge on the concept is 1240 stored in the remaining neurons), resulting in only a 1241 small accuracy drop. Or, the 3000 concept-specific 1242 neurons store knowledge, though concerning the 1243 concept, unrelated to the prompts. For instance, 1244 the validation prompts of the concept presidents 1245 all demand historical dates as predicted answers, 1246 which is only one kind of knowledge that might 1247 be expected in connection with presidents. This 1248 phenomenon actually aligns with our neuron inter-1249 1250 ference hypothesis: deactivating neurons that store unhelpful knowledge can less confuse the model, 1251 therefore improving the performance. 1252

G Experimental Environment

1254

1255

We run all experiments on NVIDIA RTX A6000 GPUs. The Python environment we use is the same



Figure 25: Accuracy results of evaluation prompts for 11 concepts in the 7b and 13b model. We report the performance of the original model (without any deactivation) e.g., 7b-original, the model 3000 random deactivated neurons, e.g. 7b-random, and the model with deactivating the top 3000 identified concept-specific neurons, e.g., 7b-concept.

as Kojima et al. (2024).⁹

H Prompt Templates

We show the actual prompt templates (with 1258 object-subject example) we use for each an 1259 relation across 6 considered languages: 1260 company_ceo Table 3. company_hq in 1261 Table 4. landmark_continent in Tain ble 5. landmark_country in Table 6. person_father in Table 7, person_mother 1264 in Table 8. person_occupation in Ta-1265 ble 9. person_plays_instrument in Table 10. person_pro_sport Table in 11. person_sport_position Table in 12. product_company in Table 13. and 1269 star_constellation in Table 14. 1270

1256

⁹Kojima et al. (2024)'s GitHub repository is available at https://github.com/kojima-takeshi188/lang_neuron

Language	Subject-Object Pair	Prompt	Expected Output
English	(Panasonic Corporation, Kazuhiro Tsuga)	Panasonic Corporation's CEO is? Answer: The CEO of Panasonic Corporation is? Answer:	Kazuhiro Tsuga
German	(Panasonic, Kazuhiro Tsuga)	Der Name des CEO von Panasonic lautet Wer ist der CEO von Panasonic? Ihr Name ist	Kazuhiro Tsuga
Spanish	(Panasonic, Kazuhiro Tsuga)	Por favor, responda directamente por su nombre. El nombre del director general de Panasonic es Por favor, responda directamente por su nombre. ¿Quién es el director general de Panasonic? Su nombre es	Kazuhiro Tsuga
French	(Panasonic, Kazuhiro Tsuga)	Veuillez répondre directement avec le nom. Le nom du président-directeur général de Panasonic est Veuillez répondre directement avec le nom. Le PDG de Panasonic est nommé	Kazuhiro Tsuga
Japanese	(バナソニック株式会社, 津賀一宏)	名前で直接お答えください。 パナソニック 株式会社のCEOの名前は 名前で直接お答えください。 パナソニック 株式会社のCEOは誰ですか? 彼らの名前は	津賀一宏
Chinese	(松下公司, 津贺一宏)	松下公司的首席执行官名字叫做 松下公司的CEO名字叫做	津贺一宏

Table 3: Prompts for **company_ceo** in different languages. We use the triple (Panasonic, company_ceo, Kazuhiro Tsuga) as an example. The subject-object pair is represented in the respective language.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Cadillac, Detroit)	The headquarters of Cadillac are in the city of? Answer: The headquarters of Cadillac are in the city of? Answer:	Detroit
German	(Cadillac, Detroit)	Cadillac hat seinen Hauptsitz in der Stadt von Der Hauptsitz von Cadillac befindet sich in der Stadt von	Detroit
Spanish	(Cadillac, Detroit)	Cadillac tiene su sede en la ciudad de La sede de Cadillac se encuentra en la ciudad de	Detroit
French	(Cadillac, Détroit)	Le nom de la ville où se trouve le siège social de Cadillac est La ville où se trouve le siège social de Cadillac s'appelle	Détroit
Japanese	(「キャデラック」,デトロイト)	「キャデラック」の本社がある都市はとこ ですか 「キャデラック」の本社はとの都市にあり ますか	デトロイト
Chinese	(凯迪拉克,底特律)	凯迪拉克总部所位于的城市名字叫做 凯迪拉克的总部所在的城市名字叫	底特律

Table 4: Prompts for **company_hq** in all languages. We use the triple (Cadillac, company_hq, Detroit) as an example. The subject-object pair is represented in the respective language.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Elbe, Europe)	Elbe is on the continent of? Answer:	Europe
German	(Elbe, Europa)	Bitte geben Sie den Kontinentnamen direkt an, z. B. Europa, Afrika usw. Der Name des Konti- nents, auf dem Elbe liegt, lautet	Europa
Spanish	(Elba, Europa)	El nombre del continente donde se encuentra Elba es	Europa
French	(Elbe, Europe)	Veuillez répondre directement avec le nom du continent. Le nom du continent où se trouve Elbe est	Europe
Japanese	(エルベ川, ヨーロッパ)	エルベ川が所在する大陸の名前は	ヨーロッパ
Chinese	(易北河,欧洲)	易北河所位于的大洲/大陆名字叫做	欧洲

Table 5: Prompts for the **landmark_continent** relation in all languages. We use the triple (Elbe, landmark_continent, Europe) as an example. The subject-object pair is represented in the respective language.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Namba Station, Japan)	Namba Station is in the country of? Answer:	Japan
German	(Namba Station, Japan)	In welchem Land liegt Namba Station? Es liegt in	Japan
Spanish	(Namba Station, Japan)	El nombre del país donde se encuentra Namba Station es	Japan
French	(Namba Station, Japan)	Le nom du pays où se trouve Namba Station est	Japan
Japanese	(難波駅,日本)	難波駅が所在する国の名前は	日本
Chinese	(难波站,日本)	难波站所位于的国家名字叫做	日本

Table 6: Prompts for the **landmark_country** relation in all languages. We use the triple (Namba Station, landmark_country, Japan) as an example. The subject-object pair is represented in the respective language.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Ronald Reagan, Jack Reagan)	Ronald Reagan's father is named? Answer:	Jack Reagan
German	(Ronald Reagan, Jack Reagan)	Der Vater von Ronald Reagan heißt	Jack Reagan
Spanish	(Ronald Reagan, Jack Reagan)	El padre de Ronald Reagan se llama	Jack Reagan
French	(Ronald Reagan, Jack Reagan)	Le père de Ronald Reagan s'appelle	Jack Reagan
Japanese	(ロナルド・レーガン,ジャック・レーガン)	名前で直接お答えください。ロナルド・ レーガンの父親の名前は	ジャック・レーガン
Chinese	(罗纳德·里根,杰克·里根)	罗纳德·里根的父亲名字叫做	杰克·里根

Table 7: Prompts for the **person_father** relation in all languages. We use the triple (Ronald Reagan, person_father, Jack Reagan) as an example. The subject-object pair is represented in the respective language.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Demi Moore, Virginia King)	Demi Moore's mother is named? Answer:	Virginia King
German	(Demi Moore, Virginia King)	Die Mutter von Demi Moore heißt	Virginia King
Spanish	(Demi Moore, Virginia King)	La madre de Demi Moore se llama	Virginia King
French	(Demi Moore, Virginia King)	Qui est la mère de Demi Moore ? Leur mère s'appelle	Virginia King
Japanese	(デミ・ムーア,ヴァージニア・キング)	名前で直接お答えください。デミ・ムー アの母親の名前は	ヴァージニア・キン グ
Chinese	(黛米·摩尔,维吉尼亚·金)	黛米·摩尔的母亲名字叫做	维吉尼亚·金

Table 8: Prompts for the **person_mother** relation in all languages. We use the triple (Demi Moore, person_mother, Virginia King) as an example. The subject-object pair is represented in the respective language.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Martin Burrell, politician)	Martin Burrell works as a? Answer: By profession, Martin Burrell is a? Answer:	politician
German	(Martin Burrell, Politiker)	Martin Burrell arbeitet als Von Beruf ist Martin Burrell ein	Politiker
Spanish	(Martin Burrell, político)	Por favor especifique el nombre de su ocupación. Martin Burrell trabaja profesionalmente como Por favor especifique el nombre de su ocupación. Por profesión, Martin Burrell es un(a)	político
French	(Martin Burrell, personnalité politique)	Veuillez répondre directement par le nom de votre profession. Le nom de la profession de Martin Burrell est Veuillez répondre directement par le nom de votre profession. Martin Burrell travaille profession- nellement comme	personnalité politique
Japanese	(マーティン・バレル,政治家)	マーティン・バレルさんの職業名は マーティン・バレルさんの職業名は	政治家
Chinese	(马丁·巴雷尔,政治人物)	马丁·巴雷尔从事的职业是一个 职业上来说,马丁·巴雷尔是一名	政治人物

Table 9: Prompts for the **person_occupation** relation in all languages. We use the triple (Martin Burrell, person_occupation, politician) as an example. The subject-object pair is represented in the respective language.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Anson Funderburgh, guitar)	Anson Funderburgh plays the instrument of? Answer:	guitar
German	(Anson Funderburgh, Gitarre)	Bitte geben Sie den Namen des Instruments di- rekt an. Das Instrument, das Anson Funderburgh spielt, heißt	Gitarre
Spanish	(Anson Funderburgh, guitarra)	Por favor responda directamente el nombre del instrumento ¿Qué instrumento toca Anson Fun- derburgh? Tocan el	guitarra
French	(Anson Funderburgh, guitare)	Veuillez répondre directement au nom de l'instrument. De quel instrument joue Anson Fun- derburgh ? Ils jouent du	guitare
Japanese	(アンリン・ファンダーバーグ,ギター)	アンソン・ファンダーバーグはどの楽器を 演奏しますか	ギター
Chinese	(安森·芬德伯格,吉他)	安森·芬德伯格所演奏的乐器名字叫做	吉他

Table 10: Prompts for the **person_plays_instrument** relation in all languages. We use the triple (Anson Funderburgh, person_plays_instrument, guitar) as an example. The subject-object pair is represented in the respective language.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Frédéric Piquionne, soccer)	Frédéric Piquionne plays the sport of? Answer:	soccer
German	(Frédéric Piquionne, Fußball)	Welchen Sport betreibt Frédéric Piquionne? Sie betreiben	Fußball
Spanish	(Frédéric Piquionne, fútbol)	Por favor, responda directamente el nombre del deporte, como fútbol, baloncesto, etc. El nombre del deporte que juega Frédéric Piquionne es:	fútbol
French	(Frédéric Piquionne, football)	Veuillez répondre directement par le nom du sport, comme le football, le basket-ball, etc. Frédéric Piquionne joue professionnellement dans le sport de	football
Japanese	(フレデリック・ビキオンヌ,サッカー)	サッカー、バスケットボールなど、スボー ツの名前を直接答えてください。フレデ リック・ビキオンヌはどのスポーツをしま すか?彼らは (スボーツ名)をしていま す。	サッカー
Chinese	(费德历·比基安尼,足球)	费德历·比基安尼从事的运动叫做	足球

Table 11: Prompts for the **person_pro_sport** relation in all languages. We use the triple (Frédéric Piquionne, person_pro_sport, soccer) as an example. The subject-object pair is represented in the respective language.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Ju Yingzhi, midfielder)	Ju Yingzhi plays in the position of a? Answer: In their sport, Ju Yingzhi plays as a? Answer:	midfielder
German	(Ju Yingzhi, Mittelfeldspieler)	Ju Yingzhi spielt auf der Position von a In ihrer Sportart spielt Ju Yingzhi als	Mittelfeldspieler
Spanish	(Ju Yingzhi, centrocampista)	Por favor, responda directamente el nombre de la posición deportiva, como delantero, defensor, etc. La posición de Ju Yingzhi en el campo deportivo es: Por favor responda directamente con el nombre de la posición deportiva, como delantero, defensor, etc. En su deporte, Ju Yingzhi juega en la posición de un:	centrocampista
French	(Ju Yingzhi, milieu de terrain)	Ju Yingzhi évolue au poste de Dans son sport, Ju Yingzhi occupe le rôle de	milieu de terrain
Japanese	(ジュ・インジー, ミッドフィールダー)	彼がブレーするスポーツでは、ジュ・イン ジーのボジションは ジュ・インジー競技場のボジションは	ミッドフィールダー
Chinese	(鞠盈智,中场)	鞠盈智在运动场上的位置名字叫做 在他/她从事的运动中,鞠盈智的位置是	中场

Table 12: Prompts for the **person_sport_position** relation in all languages. We use the triple (Ju Yingzhi, person_sport_position, midfielder) as an example. The subject-object pair is represented in the respective language.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Jeep Grand Cherokee, Chrysler)	Jeep Grand Cherokee was created by which com- pany? Answer: Jeep Grand Cherokee is a product of which com- pany? Answer:	Chrysler
German	(Jeep Grand Cherokee, Chrysler)	Bitte geben Sie direkt den Firmen-/Ländernamen an. Das Unternehmen/Land, das Jeep Grand Cherokee entwickelt hat, ist Bitte geben Sie direkt den Firmen-/Ländernamen an. Welches Unternehmen hat Jeep Grand Chero- kee entwickelt? Es wurde entwickelt von	Chrysler
Spanish	(Jeep Grand Cherokee, Chrysler)	Por favor, responda directamente el nombre de la empresa/país. ¿Qué empresa desarrolló Jeep Grand Cherokee? Fue desarrollado por Por favor responda directamente con el nombre de la empresa/país. La empresa que desarrolló Jeep Grand Cherokee se llama	Chrysler
French	(Jeep Grand Cherokee, Chrysler)	Jeep Grand Cherokee a été développé(e) par Jeep Grand Cherokee est un produit de l'entreprise	Chrysler
Japanese	 (ジーブ・グランドチェロキー, クライスラー	会社名/国名を直接お答えください。ジー) ブ・グランドチェロキーを開発したのはど の会社ですか?開発したのは次の会社は 会社名/国名を直接お答えください。ジー ブ・グランドチェロキーを開発した会社は	クライスラー
Chinese	(吉普大切诺基,克莱斯勒)	开发了吉普大切诺基的公司名字叫做 开发产品吉普大切诺基的公司名字叫	克莱斯勒

Table 13: Prompts for the **product_company** relation in all languages. We use the triple (Jeep Grand Cherokee, product_company, Chrysler) as an example. The subject-object pair is represented in the respective language.

Language	Subject-Object Pair	Prompt	Expected Output
English	(50 Persei E, Perseus)	50 Persei E is part of the constellation named? Answer:	Perseus
German	(50 Persei E, Perseus)	Bitte geben Sie den Namen des Sternbildes direkt an. Das Sternbild, zu dem 50 Persei E gehört, heißt	Perseus
Spanish	(50 Persei E, Perseus)	50 Persei E forma parte de la constelación denom- inada	Perseus
French	(50 Persei E, Persée)	Le nom de la constellation dans laquelle se trouve 50 Persei $\rm E$ est	Persée
Japanese	(50 ベルセウス座 E, ベルセウス座)	50 ベルセウス座 Eはどの星座に属していま すか?それは(星座名)という星座の一部 です。	ペルセウス座
Chinese	(50 英仙座E, 英仙座)	50 英仙座E所位于的星座名字叫做	英仙座

Table 14: Prompts for the **star_constellation** relation in all languages. We use the triple (50 Persei E, star_constellation, Perseus) as an example. The subject-object pair is represented in the respective language.