

On Relation-Specific Neurons in Large Language Models

Anonymous ACL submission

Abstract

In large language models (LLMs), certain *neurons* can store distinct pieces of knowledge learned during pretraining. While factual 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 *statistics*-based 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 involving relation r and (2) facts involving 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 cumulativeness.** Multiple neurons jointly contribute to processing facts involving relation r , with no single neuron fully encoding a fact in r on its own. **(ii) Neuron versatility.** Neurons can be shared across multiple closely related as well as less related relations. In addition, some relation neurons transfer across languages. **(iii) Neuron interference.** Deactivating neurons specific to one relation can improve LLMs’ factual recall 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), 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), i.e., neurons that are responsible for encoding each fact individually; or whether there exist *relation-specific neurons* (referred to as **RelSpec neurons**), i.e., neurons that do not represent specific facts but rather focus on the relation and guide generating the object once the subject and relation of a triple have been detected.

In this work, we examine the existence of *RelSpec* 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 *RelSpec* 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 *RelSpec* neurons.

To understand the impact of *RelSpec* neurons, we perform factual recall on held-out prompts. These prompts for each relation share the **same relation** as the positive examples used for neuron 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 *RelSpec* neurons for that relation are deactivated – *intra-relation results*. We also study how deactivating neurons for one relation influences performance on others – *inter-*

043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083

084 *relation results*. Our experiments reveal several key
085 properties of *RelSpec* neurons:

086 **Neuron cumulativity.** *RelSpec* neurons present
087 a cumulative effect – a phenomenon where an LLM
088 distributes relational knowledge across multiple
089 neurons. *RelSpec* neurons jointly contribute to deal-
090 ing with facts belonging to a relation, with no sin-
091 gle neuron fully encoding a fact on its own. This
092 property aligns with the evidence of the existence
093 of redundant and self-repair neurons (Dalvi et al.,
094 2020; McGrath et al., 2023; He et al., 2024).

095 **Neuron versatility.** As the total number of neu-
096 rons is finite, while the number of possible rela-
097 tions is vast, some *RelSpec* neurons strongly asso-
098 ciate with multiple relations. Surprisingly, these
099 relations need not be closely linked – two weakly
100 related relations can share a group of neurons, lead-
101 ing to performance drops in both relations if those
102 neurons are deactivated. *RelSpec* neurons also gen-
103 eralize across languages – *RelSpec* neurons iden-
104 tified from English have a similar effect on other
105 languages. This property aligns with neuron poly-
106 semanticity and superposition (Mu and Andreas,
107 2020; Elhage et al., 2022b; Scherlis et al., 2025).

108 **Neuron interference.** Some *RelSpec* neurons
109 appear to “confuse” the model when it processes
110 other relations. Deactivating such neurons can
111 yield improved performance on these other rela-
112 tions. This property aligns with broader evidence
113 that *sub-networks* or *circuits* within LLMs may
114 serve several different functional roles (Wang et al.,
115 2023a; Bayazit et al., 2024; Mondorf et al., 2024).

116 2 Methodology

117 2.1 Dataset Manipulation

118 We use the factual knowledge dataset from Hernan-
119 dez et al. (2024) for this research, which contains
120 25 relations. Each relation has a different number
121 of facts. Each fact can be represented as a *subject-*
122 *relation-object* triple (s, r_i, o) . We only consider
123 relations that have more than 300 facts to ensure
124 the reliability of our findings. This results in 12
125 relations. We refer to the set of triples for relation
126 r_i as \mathcal{D}_{r_i} . We then perform the following steps
127 for each relation r_i to construct the data used to
128 identify its corresponding *RelSpec* neurons.

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

jointness, $\mathcal{D}_{r_i}^{\text{eva}}$ and $\mathcal{D}_{r_i}^{\text{det}}$ do not share any subjects. 134

135 **Step 2: Formulating Prompts.** For each
136 triple (s, r_i, o) in $\mathcal{D}_{r_i}^{\text{det}}$, we create prompts con-
137 taining the **subject** s and the **relation** r_i us-
138 ing the templates provided by Hernandez et al.
139 (2024). **Note that the object o is not included**
140 **in the prompt.** For example, we construct a
141 prompt “*The CEO of NVIDIA is? Answer:*” for
142 the triple (NVIDIA, company_CEO, Jensen Huang)
143 with an expected answer “*Jensen Huang*”. We also
144 create prompts for $\mathcal{D}_{r_i}^{\text{eva}}$ in the same way. We refer
145 to the resulting prompt sets as $\mathcal{P}_{r_i}^{\text{det}}$ and $\mathcal{P}_{r_i}^{\text{eva}}$.

146 **Step 3: Validating Prompts.** We hypothesize
147 that the model will leverage *RelSpec* neurons to
148 generate the correct answer, i.e., the object. There-
149 fore, such neurons should “fire” for those prompts
150 for which **the model answers correctly**. For the
151 prompt selection, we feed each prompt in $\mathcal{P}_{r_i}^{\text{det}}$
152 to the model and set the maximum generation length
153 to be 2.¹ We then check if the predicted 2 tokens
154 are a prefix of the object: if they are, we regard the
155 output as being correct. We exclude prompts that
156 the model answers wrongly from $\mathcal{P}_{r_i}^{\text{det}}$.

157 2.2 Relation-Specific Neuron Identification

158 This work’s purpose is to identify *RelSpec neu-*
159 *rons* – neurons that solely focus on the relation
160 rather than specific relational facts concerning the
161 subject-relation-object triple. Therefore, these neu-
162 rons are different from *knowledge neurons* (which
163 encode certain facts) or *entity neurons* (which en-
164 code certain subject entities). Following Cuadros
165 et al. (2022), we identify *RelSpec* neurons using
166 statistical association measures. This method as-
167 signs a score for each neuron, representing its level
168 of “expertise” in **distinguishing a specific relation**
169 **from other considered relations**.

170 **Defining Neurons.** A neural network, or specifi-
171 cally a Transformer (Vaswani et al., 2017), consists
172 of many weight matrices. For a given weight ma-
173 trix $\mathbf{W} \in \mathbb{R}^{d_1 \times d_2}$, we define a neuron as a column,
174 mapping a representation from \mathbb{R}^{d_1} to \mathbb{R} . We assign
175 a unique index $m \in M$ to each neuron and investi-
176 gate its output value. We only consider the neurons
177 in feed-forward networks (FFNs), i.e., neurons in
178 up_proj, gate_proj, and down_proj, since previ-
179 ous studies have shown that knowledge is mostly

¹Some prior studies evaluate correctness by only checking the model’s first predicted token (Geva et al., 2023; Hernandez et al., 2024). This evaluation can be ambiguous if the answer/object is split into multiple tokens. Considering 2 predicted tokens increases reliability.

stored there (Dai et al., 2022). We also investigate neurons in other modules, e.g., attention heads, but find they are less relation-specific (see §G).

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.² We refer to the positive and negative examples selected for relation r_i as $\mathcal{E}_{r_i}^+$ and $\mathcal{E}_{r_i}^-$. The final data used to detect *RelSpec* neurons for relation r_i is then $\mathcal{E}_{r_i} = \mathcal{E}_{r_i}^+ \cup \mathcal{E}_{r_i}^-$. Each example $e_{r_i}^j$ is associated with binary label $b_{r_i}^j$: 1 if $e_{r_i}^j \in \mathcal{E}_{r_i}^+$, 0 otherwise.

Neuron Output Values. Let $o_{r_i}^{m,j,t}$ be the output value of neuron m for the t -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 $e_{r_i}^j$: $o_{r_i}^{m,j} = \frac{1}{T} \sum_{t=1}^T o_{r_i}^{m,j,t}$, where T is the number of effective tokens in $e_{r_i}^j$.

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 *RelSpec* neurons in descending order.

2.3 Controlled Generation

For each relation r_i , we want to investigate the impact of the identified top- k *RelSpec* neurons. Therefore, we control text generation by overriding their 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

²The sampling ratio is based on previous research (Kojima et al., 2024) – ratios that are too small or too large are not good for computing reliable *AP* values. We also sample negative examples with different seeds in our preliminary experiments. The identified relation neurons show little change, suggesting stability.

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 neuron statistics

of all top- k *RelSpec* neurons to a constant 0 and set the maximum generation length to 2 (similar to the setup in validating prompts, cf. §2). The predicted 2 tokens are then compared to the object. The prediction is regarded as correct if the predicted 2 tokens are a prefix of the object.

3 Experimental Setup

3.1 Models

We consider the 7B and 13B models from the **LLama-2** family (Touvron et al., 2023).³ As mentioned in §2.2, we consider the neurons in **FFNs**, which account for more than half of neurons in both 7B and 13B models, as shown in Table 1. We also report our preliminary results when considering neurons in other modules, i.e., attention heads, in §G. Their effectiveness tends to be unsatisfactory compared with **FFNs**, supporting our choice.

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 $\mathcal{P}_{r_i}^{\text{det}}$ (for neuron identification) and $\mathcal{P}_{r_i}^{\text{eva}}$ (for evaluation) are constructed for each relation r_i , yielding varying numbers $|\mathcal{P}_{r_i}^{\text{det}}|$ of prompts. $\mathcal{P}_{r_i}^{\text{eva}}$ is constructed by randomly selecting 50 triples for each relation. Since these 50 triples are not used when creating $\mathcal{P}_{r_i}^{\text{det}}$, this setup ensures **no subject entity overlap between $\mathcal{P}_{r_i}^{\text{det}}$ and $\mathcal{P}_{r_i}^{\text{eva}}$ for the same relation r_i** . The elimination of subject entity overlap allows us to disentangle the effect of entities and focus on the only shared attribute between $\mathcal{P}_{r_i}^{\text{det}}$ and $\mathcal{P}_{r_j}^{\text{det}}$ – the relation itself. In addition, we ensure **minimal subject entity overlap across relations** (mostly 0 between $\mathcal{P}_{r_i}^{\text{det}}$ and $\mathcal{P}_{r_j}^{\text{det}}$). The only exception is between *person_mother* and *person_father*, which share a lot of subject entities in $\mathcal{P}_{r_i}^{\text{det}}$; however, the two relations **share no subject entities in $\mathcal{P}_{r_i}^{\text{eva}}$** . A detailed analysis of entity overlap is presented in §B.

³We conduct a similar investigation on **Gemma-7B** (Gemma Team et al., 2024), as detailed in §C, and observe experimental results consistent with those of LLama-2.

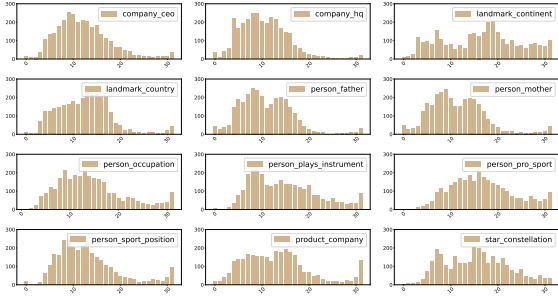


Figure 1: Distribution of *RelSpec* neurons across layers. Most are located in the middle layers.

4 Results and Discussion

We apply our identification method to both Llama-2 7B and 13B models for all 12 relations. We regard the **top 3,000** neurons with the highest *AP* values as the *RelSpec* neurons; 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 1 (see §D 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, *RelSpec* 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 functional mapping vectors can be extracted from the middle layers of LLMs (Merullo et al., 2024; Hernandez et al., 2024; Todd et al., 2024).

Neuron Overlap Across Relations. We display the overlap of *RelSpec* neurons across relations for the 7B model in Figure 2 (13B is in §D). We see that *person_mother* and *person_father* share many neurons, possibly due to the large overlap between their subject entities, (see §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

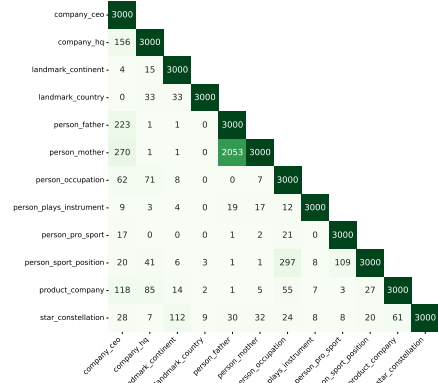


Figure 2: Neuron overlap of *RelSpec* neurons across 12 relations. For example, the number of neurons shared between the 3,000 identified neurons for *person_father* and the 3,000 for *person_mother* is 2053 (in green).

a kind of 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 *RelSpec* 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 *RelSpec* neurons for a relation and evaluating the same relation, we also create a baseline by **randomly** deactivating 3,000 neurons in the model. Results for the original models and for the two interventions are in Figure 3.

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 *RelSpec* neurons are deactivated.⁴ On the other hand, the model with 3,000 random 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. **As $\mathcal{P}_{r_i}^{\text{eva}}$ and $\mathcal{P}_{r_i}^{\text{det}}$ do not share any subject entities, this drop can only**

⁴For 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.

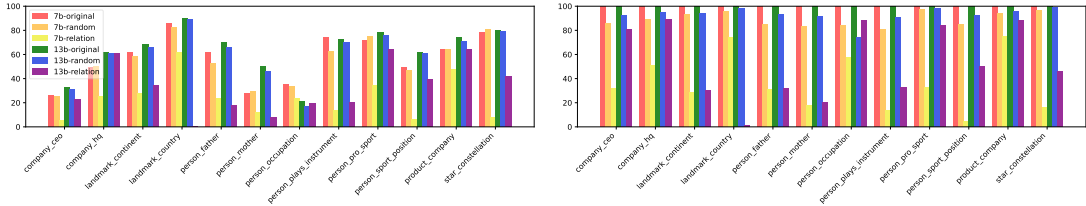


Figure 3: Intra-relation results. The left (resp. right) figure displays the results of held-out evaluation prompt set $\mathcal{P}_{r_i}^{\text{eva}}$ (resp. identification prompt set $\mathcal{P}_{r_i}^{\text{det}}$). We report the performance of the original model (without any deactivation), e.g.,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

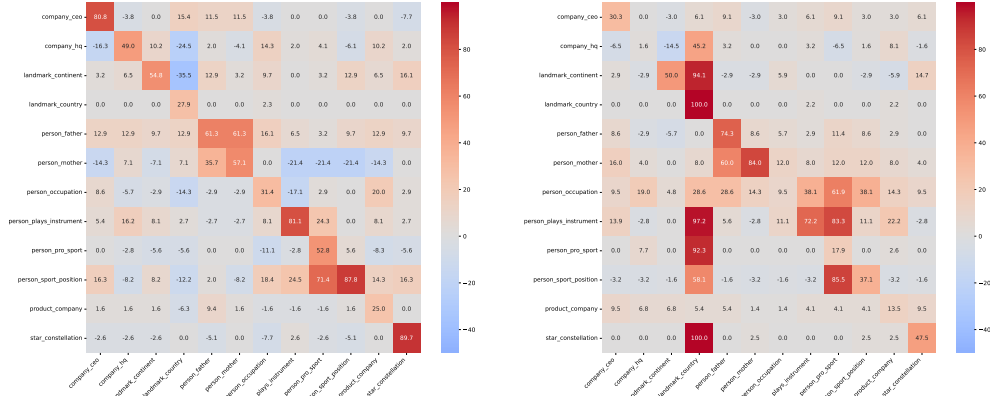


Figure 4: Inter-relation results. Accuracy drops (in %) for the 7B (left) and the 13B model (right) on $\mathcal{P}_{r_i}^{\text{eva}}$. The number in cell (r_i, r_j) indicates the accuracy drop of relation r_i when deactivating the relation neurons of r_j .

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 $\mathcal{P}_{r_i}^{\text{det}}$.⁵ We thus argue that *RelSpec* neurons exist in LLMs: they are entity-irrelevant and focus on specific relations.

On the other hand, the accuracy does not drop to 0 for any relation (except landmark_country in the 13B model) when its identified *RelSpec* neurons are deactivated. This indicates these 3,000 neurons do not equally influence all facts that belong to a certain relation, which highlights that LLMs do not uniformly encode all facts belonging to a given relation, but rather distribute relational knowledge across neurons in a manner that can vary significantly from fact to fact. We validate this by showing 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 of neurons may correlate with how frequently it appears in the pretraining data in §E.

4.2.2 Inter-Relation Results

To understand how *RelSpec* neurons influence the model’s ability to answer prompts across multiple relations, we use **accuracy drop** as a met-

⁵There might be another confounding variable since $\mathcal{P}_{r_i}^{\text{eva}}$ and $\mathcal{P}_{r_i}^{\text{det}}$ use the same prompt templates for each relation. But we show in §5.3 that even when other prompt templates are used, the effectiveness of these neurons is still preserved.

ric: $\text{acc_drop}_{r_i, r_j} = \frac{\text{acc}_{r_i}^{\text{original}} - \text{acc}_{r_i}^{\text{deactivated-}r_j}}{\text{acc}_{r_i}^{\text{original}}}$, where $\text{acc}_{r_i}^{\text{original}}$ and $\text{acc}_{r_i}^{\text{deactivated-}r_j}$ are the respective accuracy for $\mathcal{P}_{r_i}^{\text{eva}}$ of (a) the original model and (b) when the *RelSpec* 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. Particularly, 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.

Neuron versatility. We observe that deactivating neurons for one relation can strongly affect not only that relation but also others, both closely and loosely related relations. E.g., disabling person_pro_sport neurons has a large effect on person_sport_position (but not vice versa) in both models, likely because a model first needs

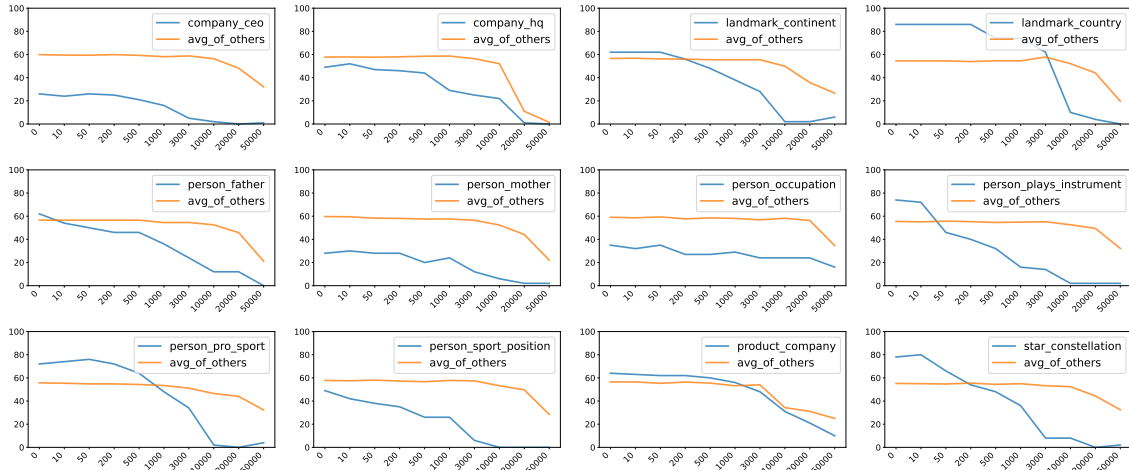


Figure 5: Influence of deactivating different numbers of *RelSpec* neurons for each relation. We show accuracy on the relation itself and the average accuracy on other relations. Increasing the number clearly affects the relation itself, but the effect on other relations starts only at 3,000 or 10,000 neurons.

to understand “sport” before inferring “position”. Similarly, deactivating `person_father` neurons reduces accuracy on `person_mother`, as both share the concept of a parental relationship. Even 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 *RelSpec* neurons for one relation can sometimes **improve** the accuracy for others – 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 relations: disabling `company_ceo` neurons slightly boosts accuracy on `company_hq` for both models. 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.**

5 Complementary Analyses

5.1 Influence of the Numbers of Neurons

In this section, we investigate the effect of varying the number of *RelSpec* neurons on the 7B model (see §D 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 **accuracy** for that relation and the **average accuracy** for all other relations in Figure 5. Results for all relation-relation pairs are in Figure 22.

Neuron cumulativity. By increasing the number of neurons for deactivation, we see a consistent accuracy drop in all relations. This suggests neuron cumulativity: **LLMs distribute relational knowledge across multiple neurons, which jointly contribute to dealing with facts belonging to a relation.** However, cumulativity varies across relations. Some relations are far more sensitive to a smaller-scale deactivation than other relations, indicating a smaller set of neurons is specifically leveraged for those relations. We hypothesize this sensitivity may correlate with the frequency of the facts in each relation in the pretraining data: more frequent facts may be memorized more robustly and thus remain less sensitive to deactivation. We empirically verify this hypothesis in §E.

Deactivating *RelSpec* neurons has a marginal effect on other relations until certain thresholds are reached. Typically, these thresholds lie between 3,000 and 10,000 as shown in Figure 5, below which the accuracy on other relations remains stable – **supporting the choice of 3,000 neurons in §3.** Once more neurons are deactivated, other relations also deteriorate, consistent with our **neuron versatility** hypothesis. However, 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 –

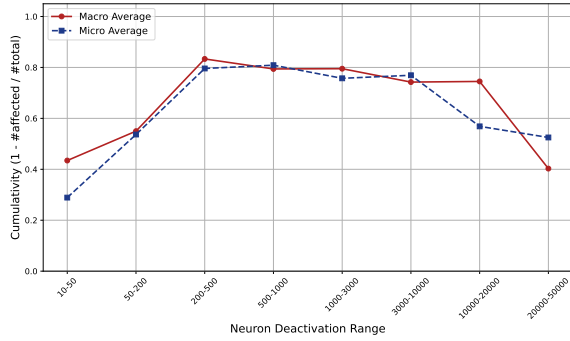


Figure 6: Macro and micro averaged neuron cumulatvity for each neuron deactivation range. Cumulatvity is defined as $1 - \frac{\#affected}{\#total}$, with macro averaging across relations and micro averaging across prompts. Both trends show that cumulatvity increases as the range increases.

possibly because some of these neurons underlie more general generation capabilities of the model (Sun et al., 2024; Yu et al., 2024).

Validation of the cumulative effect. It remains unclear whether the further accuracy drop between any two thresholds in Figure 5 is driven by **the newly deactivated neurons** (the isolated effect of deactivated neurons) or **the cumulative effect of all deactivated neurons**. To further validate our neuron cumulatvity hypothesis, we conduct an experiment on each consecutive pair of thresholds, e.g., 1000-3000. Specifically, we identify prompts from $\mathcal{P}_{r_i}^{eva}$ where the model answers correctly with neurons of the smaller range being deactivated, but fails when neurons of the larger range are deactivated (#total). We then deactivate only the neurons from the intermediate difference and measure the number of affected prompts – prompts for which the model answers wrongly (#affected). Figure 6 shows the macro and micro averaged cumulatvity, defined as $1 - \frac{\#affected}{\#total}$. We notice that neuron behavior becomes increasingly cumulatvity as the range increases, indicating that only deactivating neurons from the intermediate difference is not enough to make the model answer wrongly. There is a drop after the ranges 10000-20000 and 20000-50000, which can be explained by the fact that many more neurons are deactivated compared with the earlier ranges. We also show the individual number of #total/#affected prompts in each relation in each range in Table 2. **Thus, our results favor the cumulative effect over the isolated effect** – multiple neurons jointly contribute to dealing with facts belonging to a relation, with no single neuron fully encoding a fact on its own.

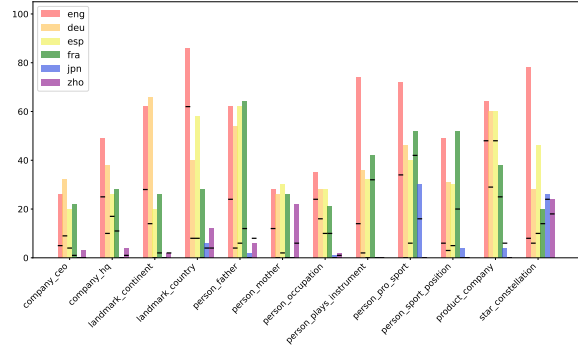


Figure 7: Accuracy on 12 relations across 6 languages. The bars show the accuracy of the original 7B model. The horizontal line in each bar indicates the performance after deactivation of 3,000 *RelSpec* neurons. Even though these neurons are identified using English prompts, they usually influence other languages, indicating multilinguality of these neurons.

5.2 Are These Neurons Multilingual?

Recent studies suggest that some neurons encoding factual knowledge or handling specific tasks are language-agnostic (Stanczak et al., 2022; Zhang et al., 2024; Wang et al., 2024a). A natural question is whether *RelSpec* neurons – identified solely via English prompts – also function across languages. To explore this, we translate $\mathcal{P}_{r_i}^{eva}$ to 5 languages: German (**deu**), Spanish (**esp**), French (**fra**), Chinese (**zho**), and Japanese (**jpn**) (see §F 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 7.

Although the model’s accuracy is generally lower in non-English languages, it still shows good factual recall for most relations (except for jpn and zho). Once the neurons for a given relation are deactivated, the accuracy drops across nearly 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 Effect of Prompt Templates

There is a possible confounding variable: the identified relation-specific neurons could be associated with the prompt templates used in $\mathcal{P}_{r_i}^{det}$. The degradation in $\mathcal{P}_{r_i}^{eva}$ would then be due to the identified neurons encoding syntactic structure rather than ab-

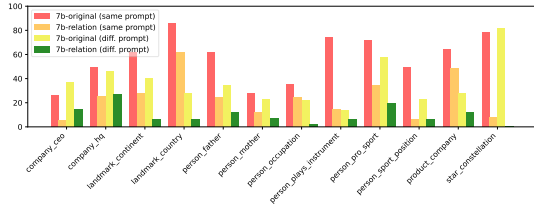


Figure 8: Intra-relation results on original prompts $\mathcal{P}_{r_i}^{\text{eva}}$ and additional prompts $\mathcal{P}_{r_i}^{\text{eva-2}}$. $\mathcal{P}_{r_i}^{\text{eva-2}}$ is constructed with same triples as $\mathcal{P}_{r_i}^{\text{eva}}$ but different prompt templates are used. A consistent decrease across relations indicates that the identified neurons are not specific to prompts.

525 struct relation semantics. To exclude this confounding
 526 variable, we create an additional evaluation set
 527 $\mathcal{P}_{r_i}^{\text{eva-2}}$ where the **same triples** as $\mathcal{P}_{r_i}^{\text{eva}}$ but **different**
 528 **prompt templates** are used for each relation. We then deactivate the previously identified 3,000
 529 neurons in the 7B model and measure the effect on the new prompts. Figure 8 presents the results. We
 530 observe that the accuracy with new prompts is a bit different from the accuracy when the original
 531 templates are used. This is not surprising since LLMs are sensitive to the prompt templates (Sclar
 532 et al., 2024). Nevertheless, we still see that the deactivation of neurons results in consistent accuracy
 533 drops for new prompts across relations. Therefore, the neurons are not subject to the templates used
 534 to describe the relation. Instead, **the identified neurons are only associated with the abstract**
 535 **relation semantics.**

543 5.4 Relations vs. Concepts

544 We saw in Figure 2 that the storage of relations is generally well separated, but there are exceptions.
 545 We can view a relation as relating two **concepts** or **topics**, e.g., company_ceo relates instances of
 546 the subject concept “company” to instances of the object concept “CEO”. From this perspective, the
 547 exceptions in Figure 2, i.e., cases where a relation r_1 overlaps with a relation r_2 , are generally cases
 548 where the concepts of r_1 and r_2 are the same or overlap. To further explore this hypothesis empirically,
 549 we again use the method applied in §2 to relations, but now use it for subject concepts.⁶ That is,
 550 we identify sets of **concept-specific neurons**. We group the triples by their subjects, resulting in
 551 9 different concepts. We then create prompts with **novel relations** such as “can” and “has a”,
 552 **balanced across positive and negative samples**. This ensures that the model’s completion for a prompt
 553
 554
 555
 556
 557
 558
 559
 560
 561

⁶We do not consider the object concepts explicitly because the objects are not presented in the prompts for relation-specific or concept-specific neuron identification (cf. §2).

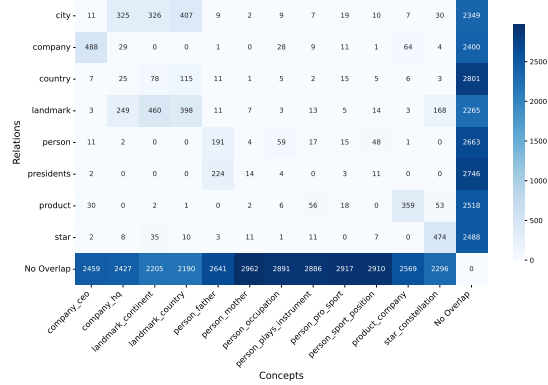


Figure 9: Overlap between the top 3000 neurons of relations and concepts in the 13B model.

562 like (“Lincoln has a”) depends on the concept instance
 563 “Lincoln”, not on the relation.

564 Figure 9 shows the overlap between relation
 565 neurons and concept neurons. Most of the cells
 566 with large counts support our hypothesis that
 567 the overlaps between relations we observe are
 568 rooted in these relations being representationally
 569 associated with their concepts. Clear examples
 570 include company_ceo and its subject concept
 571 company; company_hq and its object concept
 572 city (assuming that hq is a subcategory of city);
 573 and landmark_continent and its subject concept
 574 landmark. There is little overlap of person with
 575 relations like person_mother, potentially because
 576 person is a more general and semantically un-
 577 specific concept than the others. However, most
 578 identified neurons are only concept neurons or
 579 only relation neurons, **suggesting that relational and**
 580 **conceptual representations are largely separate.**

581 6 Conclusion

582 This work highlights the existence of relation-
 583 specific neurons in LLMs – neurons that focus
 584 on relations rather than entities. Our experiments
 585 show that *RelSpec* neurons primarily reside in the
 586 middle layers and can be shared across multiple
 587 relations. Through systematic deactivation, we
 588 reveal their influence on both the targeted and
 589 other relations, leading to three key hypotheses:
 590 **neuron cumulativity** (multiple neurons jointly
 591 contribute to dealing with facts belonging to a
 592 relation), **neuron versatility** (neurons are
 593 shared across relations and languages), and
 594 **neuron interference** (neurons from one
 595 relation can disrupt the processing of another).
 596 These findings shed new light on how LLMs
 597 handle relational facts at the neuron level,
 598 contributing to the interpretability of LLMs.

598 Limitations

599 While our findings provide valuable insights, sev-
600 eral limitations remain and offer opportunities for
601 future research. First, this work focuses on factual
602 knowledge grouped into 12 relations because the
603 reliability of the neuron identification method re-
604 quires enough facts in each relation. Although
605 this selection does not diminish the validity of
606 our findings and hypotheses, it represents a rel-
607 atively narrow set of relations. Future work can
608 explore a broader range of relations and analyze
609 how relation-specific neurons behave across a more
610 diverse set of relations. Second, our multilingual
611 analysis includes only five languages. While these
612 languages demonstrate neuron versatility, they do
613 not fully capture linguistic diversity. Future re-
614 search could investigate additional languages, par-
615 ticularly low-resource ones, to determine whether
616 relation-specific neurons exhibit similar relational
617 functionality across these languages. Thirdly, we
618 draw our findings from the LLama-2 family in the
619 main content due to page limit and resource con-
620 straints. We also conduct the same investigation on
621 Gemma-7B (Gemma Team et al., 2024) (cf. §C),
622 which shows similar trends as we observe for mod-
623 els from the LLama-2 family. Future work can
624 explore even larger models or models with post-
625 training techniques like instruction-tuning. Lastly,
626 we observe that more frequent facts tend to be
627 more robust to the deactivation of relation-specific
628 neurons in both the 7B and 13B models (cf. §E).
629 Fact frequency is approximated using the Dolma
630 corpus (Soldaini et al., 2024) in this study. How-
631 ever, LLama-2 models may incorporate a larger
632 and more diverse pretraining dataset, potentially
633 leading to some discrepancies between these ap-
634 proximated fact frequencies and their actual fre-
635 quencies.

636 References

637 Omer Antverg and Yonatan Belinkov. 2022. [On the](#)
638 [pitfalls of analyzing individual neurons in language](#)
639 [models](#). In *The Tenth International Conference on*
640 *Learning Representations, ICLR 2022, Virtual Event,*
641 *April 25-29, 2022*. OpenReview.net.

642 Anthony Bau, Yonatan Belinkov, Hassan Sajjad, Nadir
643 Durrani, Fahim Dalvi, and James R. Glass. 2019.
644 [Identifying and controlling important neurons in neu-](#)
645 [ral machine translation](#). In *7th International Confer-*
646 *ence on Learning Representations, ICLR 2019, New*
647 *Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

Deniz Bayazit, Negar Foroutan, Zeming Chen, Gail 648
Weiss, and Antoine Bosselut. 2024. [Discovering](#)
649 [knowledge-critical subnetworks in pretrained lan-](#)
650 [guage models](#). In *Proceedings of the 2024 Confer-*
651 *ence on Empirical Methods in Natural Language*
652 *Processing*, pages 6549–6583, Miami, Florida, USA.
653 Association for Computational Linguistics. 654

Steven Bills, Nick Cammarata, Dan Mossing, Henk 655
Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever, Jan
656 Leike, Jeff Wu, and William Saunders. 2023. [Lan-](#)
657 [guage models can explain neurons in language mod-](#)
658 [els](#). 659

Xavier Suau Cuadros, Luca Zappella, and Nicholas 660
Apostoloff. 2022. [Self-conditioning pre-trained lan-](#)
661 [guage models](#). In *International Conference on Ma-*
662 *chine Learning, ICML 2022, 17-23 July 2022, Balti-*
663 *more, Maryland, USA*, volume 162 of *Proceedings*
664 *of Machine Learning Research*, pages 4455–4473.
665 PMLR. 666

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao 667
Chang, and Furu Wei. 2022. [Knowledge neurons in](#)
668 [pretrained transformers](#). In *Proceedings of the 60th*
669 *Annual Meeting of the Association for Computational*
670 *Linguistics (Volume 1: Long Papers)*, pages 8493–
671 8502, Dublin, Ireland. Association for Computational
672 Linguistics. 673

Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan 674
Belinkov, Anthony Bau, and James R. Glass. 2019.
675 [What is one grain of sand in the desert? analyz-](#)
676 [ing individual neurons in deep NLP models](#). In *The*
677 *Thirty-Third AAAI Conference on Artificial Intelli-*
678 *gence, AAAI 2019, The Thirty-First Innovative Ap-*
679 *plications of Artificial Intelligence Conference, IAAI*
680 *2019, The Ninth AAAI Symposium on Educational*
681 *Advances in Artificial Intelligence, EAAI 2019, Hon-*
682 *olulu, Hawaii, USA, January 27 - February 1, 2019*,
683 pages 6309–6317. AAAI Press. 684

Fahim Dalvi, Hassan Sajjad, Nadir Durrani, and 685
Yonatan Belinkov. 2020. [Analyzing redundancy in](#)
686 [pretrained transformer models](#). In *Proceedings of the*
687 *2020 Conference on Empirical Methods in Natural*
688 *Language Processing (EMNLP)*, pages 4908–4926,
689 Online. Association for Computational Linguistics. 690

Clément Dumas, Veniamin Veselovsky, Giovanni 691
Monea, Robert West, and Chris Wendler. 2024. [How](#)
692 [do llamas process multilingual text? a latent explo-](#)
693 [ration through activation patching](#). In *ICML 2024*
694 *Workshop on Mechanistic Interpretability*. 695

Nadir Durrani, Hassan Sajjad, Fahim Dalvi, and 696
Yonatan Belinkov. 2020. [Analyzing individual neu-](#)
697 [rons in pre-trained language models](#). In *Proceed-*
698 *ings of the 2020 Conference on Empirical Methods*
699 *in Natural Language Processing (EMNLP)*, pages
700 4865–4880, Online. Association for Computational
701 Linguistics. 702

Yanai Elazar, Akshita Bhagia, Ian Magnusson, Abhi- 703
lasha Ravichander, Dustin Schwenk, Alane Suhr, 704

705	Evan Pete Walsh, Dirk Groeneveld, Luca Soldaini,	and Dimitris Bertsimas. 2024. Universal neurons in	762
706	Sameer Singh, Hannaneh Hajishirzi, Noah A. Smith,	GPT2 language models . <i>Trans. Mach. Learn. Res.</i> ,	763
707	and Jesse Dodge. 2024. What’s in my big data? In	2024.	764
708	<i>The Twelfth International Conference on Learning</i>		
709	<i>Representations, ICLR 2024, Vienna, Austria, May</i>	Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine	765
710	<i>7-11, 2024</i> . OpenReview.net.	Harvey, Dmitrii Troitskii, and Dimitris Bertsimas.	766
		2023. Finding neurons in a haystack: Case studies	767
		with sparse probing . <i>Trans. Mach. Learn. Res.</i> , 2023.	768
711	Nelson Elhage, Tristan Hume, Catherine Olsson,		
712	Neel Nanda, Tom Henighan, Scott Johnston, Sheer	Shwai He, Guoheng Sun, Zheyu Shen, and Ang Li. 2024.	769
713	ElShowk, Nicholas Joseph, Nova DasSarma, Ben	What Matters in Transformers? Not All Attention is	770
714	Mann, Danny Hernandez, Amanda Askell, Kamal	Needed . <i>Preprint</i> , arXiv:2406.15786.	771
715	Ndousse, Andy Jones, Dawn Drain, Anna Chen, Yun-		
716	tao Bai, Deep Ganguli, Liane Lovitt, and 14 others.	Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin	772
717	2022a. Softmax linear units . <i>Transformer Circuits</i>	Meng, Martin Wattenberg, Jacob Andreas, Yonatan	773
718	<i>Thread</i> .	Belinkov, and David Bau. 2024. Linearity of relation	774
		decoding in transformer language models . In <i>The</i>	775
		<i>Twelfth International Conference on Learning Rep-</i>	776
		<i>resentations, ICLR 2024, Vienna, Austria, May 7-11,</i>	777
		<i>2024</i> . OpenReview.net.	778
719	Nelson Elhage, Tristan Hume, Catherine Olsson,		
720	Nicholas Schiefer, Tom Henighan, Shauna Kravec,	Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham	779
721	Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain,	Neubig. 2020. How can we know what language	780
722	Carol Chen, Roger Grosse, Sam McCandlish, Jared	models know? <i>Transactions of the Association for</i>	781
723	Kaplan, Dario Amodei, Martin Wattenberg, and	<i>Computational Linguistics</i> , 8:423–438.	782
724	Christopher Olah. 2022b. Toy models of superposi-		
725	tion . <i>Preprint</i> , arXiv:2209.10652.	Takeshi Kojima, Itsuki Okimura, Yusuke Iwasawa, Hit-	783
		omi Yanaka, and Yutaka Matsuo. 2024. On the multi-	784
		lingual ability of decoder-based pre-trained language	785
		models: Finding and controlling language-specific	786
		neurons . In <i>Proceedings of the 2024 Conference of</i>	787
		<i>the North American Chapter of the Association for</i>	788
		<i>Computational Linguistics: Human Language Tech-</i>	789
		<i>nologies (Volume 1: Long Papers)</i> , pages 6919–6971,	790
		Mexico City, Mexico. Association for Computational	791
		Linguistics.	792
726	Nelson Elhage, Neel Nanda, Catherine Olsson, Tom		
727	Henighan, Nicholas Joseph, Ben Mann, Amanda	János Kramár, Tom Lieberum, Rohin Shah, and Neel	793
728	Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova	Nanda. 2024. AtP*: An efficient and scalable	794
729	DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-	method for localizing LLM behaviour to components .	795
730	Dodds, Danny Hernandez, Andy Jones, Jackson	<i>Preprint</i> , arXiv:2403.00745.	796
731	Kernion, Liane Lovitt, Kamal Ndousse, and 6 others.		
732	2021. A mathematical framework for transformer	Tom Lieberum, Matthew Rahtz, János Kramár, Neel	797
733	circuits . <i>Transformer Circuits Thread</i> .	Nanda, Geoffrey Irving, Rohin Shah, and Vladimir	798
		Mikulik. 2023. Does Circuit Analysis Interpretability	799
		Scale? Evidence from Multiple Choice Capabilities	800
		in Chinchilla . <i>Preprint</i> , arXiv:2307.09458.	801
734	Amit Elhelo and Mor Geva. 2024. Inferring functional-		
735	ity of attention heads from their parameters . <i>Preprint</i> ,	Weize Liu, Yinlong Xu, Hongxia Xu, Jintai Chen, Xum-	802
736	arXiv:2412.11965.	ing Hu, and Jian Wu. 2024. Unraveling babel: Ex-	803
		ploring multilingual activation patterns of LLMs and	804
		their applications . In <i>Proceedings of the 2024 Confer-</i>	805
		<i>ence on Empirical Methods in Natural Language Pro-</i>	806
		<i>cessing</i> , pages 11855–11881, Miami, Florida, USA.	807
		Association for Computational Linguistics.	808
737	Gemma Team, Thomas Mesnard, Cassidy Hardin,		
738	Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,	Ang Lv, Yuhan Chen, Kaiyi Zhang, Yulong Wang,	809
739	Laurent Sifre, Morgane Rivière, Mihir Sanjay	Lifeng Liu, Ji-Rong Wen, Jian Xie, and Rui Yan.	810
740	Kale, Juliette Love, Pouya Tafti, Léonard Hussenot,	2024. Interpreting key mechanisms of factual re-	811
741	Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam	call in transformer-based language models . <i>Preprint</i> ,	812
742	Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros,	arXiv:2403.19521.	813
743	Ambrose Slone, and 89 others. 2024. Gemma: Open		
744	models based on gemini research and technology .	Thomas McGrath, Matthew Rahtz, Janos Kramar,	814
745	<i>Preprint</i> , arXiv:2403.08295.	Vladimir Mikulik, and Shane Legg. 2023. The Hy-	815
		dra Effect: Emergent Self-repair in Language Model	816
		Computations . <i>Preprint</i> , arXiv:2307.15771.	817
746	Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir		
747	Globerson. 2023. Dissecting recall of factual associa-		
748	tions in auto-regressive language models . In <i>Proceed-</i>		
749	<i>ings of the 2023 Conference on Empirical Methods in</i>		
750	<i>Natural Language Processing</i> , pages 12216–12235,		
751	Singapore. Association for Computational Linguis-		
752	tics.		
753	Mor Geva, Roei Schuster, Jonathan Berant, and Omer		
754	Levy. 2021. Transformer feed-forward layers are key-		
755	value memories . In <i>Proceedings of the 2021 Confer-</i>		
756	<i>ence on Empirical Methods in Natural Language Pro-</i>		
757	<i>cessing</i> , pages 5484–5495, Online and Punta Cana,		
758	Dominican Republic. Association for Computational		
759	Linguistics.		
760	Wes Gurnee, Theo Horsley, Zifan Carl Guo, Tara Rezaei		
761	Kheirkhah, Qinyi Sun, Will Hathaway, Neel Nanda,		

818	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT . In <i>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</i> .	876
819		877
820		878
821		879
822		
823		880
824		881
825	Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. 2023. Mass-editing memory in a transformer . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	882
826		883
827		884
828		885
829		886
830		
831	Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. 2024. Language models implement simple Word2Vec-style vector arithmetic . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 5030–5047, Mexico City, Mexico. Association for Computational Linguistics.	887
832		888
833		889
834		890
835		891
836		892
837		893
838		894
839	Philipp Mondorf, Sondre Wold, and Barbara Plank. 2024. Circuit Compositions: Exploring Modular Structures in Transformer-Based Language Models . Preprint, arXiv:2410.01434.	895
840		896
841		897
842		898
843	Jesse Mu and Jacob Andreas. 2020. Compositional explanations of neurons . In <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual</i> .	899
844		900
845		901
846		902
847		903
848	Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. 2020. Zoom in: An introduction to circuits . <i>Distill</i> .	904
849		905
850		
851	Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, 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. In-context learning and induction heads . <i>Transformer Circuits Thread</i> .	906
852		907
853		908
854		909
855		910
856		911
857		912
858		913
859	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.	914
860		915
861		916
862		917
863		
864		918
865		919
866		920
867		921
868	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.	922
869		923
870		924
871		925
872	Hassan Sajjad, Nadir Durrani, and Fahim Dalvi. 2022. Neuron-level interpretation of deep NLP models: A survey . <i>Transactions of the Association for Computational Linguistics</i> , 10:1285–1303.	926
873		927
874		928
875		929
		930
		931
		932
		933
	Adam Scherlis, Kshitij Sachan, Adam S. Jermyn, Joe Benton, and Buck Shlegeris. 2025. Polysemanticity and capacity in neural networks . Preprint, arXiv:2210.01892.	
	Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. Quantifying language models’ sensitivity to spurious features in prompt design or: How I learned to start worrying about prompt formatting . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	
	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, Xinxu 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 <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , 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 <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 7101–7113, Miami, Florida, USA. Association for Computational Linguistics.	
	Karolina Stanczak, Edoardo Ponti, Lucas Torroba Henigen, Ryan Cotterell, and Isabelle Augenstein. 2022. Same neurons, different languages: Probing morphosyntax in multilingual pre-trained models . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , 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 <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , 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 <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay	

934	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	Mengxia Yu, De Wang, Qi Shan, Colorado Reed, and	990
935	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	Alvin Wan. 2024. The super weight in large language	991
936	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	models . <i>Preprint</i> , arXiv:2411.07191.	992
937	Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 oth-		
938	ers. 2023. Llama 2: Open foundation and fine-tuned	Qinan Yu, Jack Merullo, and Ellie Pavlick. 2023. Char-	993
939	chat models . <i>Preprint</i> , arXiv:2307.09288.	acterizing mechanisms for factual recall in language	994
		models . In <i>Proceedings of the 2023 Conference on</i>	995
940	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	<i>Empirical Methods in Natural Language Processing</i> ,	996
941	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	pages 9924–9959, Singapore. Association for Com-	997
942	Kaiser, and Illia Polosukhin. 2017. Attention is all	putational Linguistics.	998
943	you need . In <i>Advances in Neural Information Pro-</i>		
944	<i>cessing Systems 30: Annual Conference on Neural</i>	Zeping Yu and Sophia Ananiadou. 2024. Neuron-level	999
945	<i>Information Processing Systems 2017, December 4-9,</i>	knowledge attribution in large language models . In	1000
946	<i>2017, Long Beach, CA, USA</i> , pages 5998–6008.	<i>Proceedings of the 2024 Conference on Empirical</i>	1001
		<i>Methods in Natural Language Processing</i> , pages	1002
947	Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov,	3267–3280, Miami, Florida, USA. Association for	1003
948	Sharon Qian, Daniel Nevo, Yaron Singer, and Stu-	Computational Linguistics.	1004
949	art M. Shieber. 2020. Investigating gender bias in		
950	language models using causal mediation analysis .	Xue Zhang, Yunlong Liang, Fandong Meng, Song-	1005
951	In <i>Advances in Neural Information Processing Sys-</i>	ming Zhang, Yufeng Chen, Jinan Xu, and Jie	1006
952	<i>tems 33: Annual Conference on Neural Information</i>	Zhou. 2024. Multilingual knowledge editing	1007
953	<i>Processing Systems 2020, NeurIPS 2020, December</i>	with language-agnostic factual neurons . <i>Preprint</i> ,	1008
954	<i>6-12, 2020, virtual</i> .	arXiv:2406.16416.	1009
955	Kevin Ro Wang, Alexandre Variengien, Arthur Conmy,	Yiran Zhao, Wenxuan Zhang, Guizhen Chen, Kenji	1010
956	Buck Shlegeris, and Jacob Steinhardt. 2023a. Inter-	Kawaguchi, and Lidong Bing. 2024. How do large	1011
957	pretability in the wild: a circuit for indirect object	language models handle multilingualism? <i>Preprint</i> ,	1012
958	identification in GPT-2 small . In <i>The Eleventh In-</i>	arXiv:2402.18815.	1013
959	<i>ternational Conference on Learning Representations,</i>		
960	<i>ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . Open-	A Related Work	1014
961	Review.net.		
		Mechanistic interpretability (MI) is a growing sub-	1015
962	Kevin Ro Wang, Alexandre Variengien, Arthur Conmy,	field of interpretability that aims to understand	1016
963	Buck Shlegeris, and Jacob Steinhardt. 2023b. Inter-	LLMs by breaking them down into smaller compo-	1017
964	pretability in the wild: a circuit for indirect object	nents and fundamental computations. It has gained	1018
965	identification in GPT-2 small . In <i>The Eleventh In-</i>	significant attention for studying how LLMs re-	1019
966	<i>ternational Conference on Learning Representations,</i>	call factual knowledge learned during pretraining	1020
967	<i>ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . Open-	(Meng et al., 2022; Dai et al., 2022; Geva et al.,	1021
968	Review.net.	2023; Yu et al., 2023; Lv et al., 2024; Wang et al.,	1022
		2024b). Following Olah et al. (2020); Rai et al.	1023
969	Weixuan Wang, Barry Haddow, Minghao Wu, Wei	(2024), MI research can be categorized into two ar-	1024
970	Peng, and Alexandra Birch. 2024a. Sharing mat-	reas: the study of features and the study of circuits ,	1025
971	ters: Analysing neurons across languages and tasks	based on the type of decomposed components. Fea-	1026
972	in llms . <i>Preprint</i> , arXiv:2406.09265.	tures refer to human-interpretable properties en-	1027
		coded in model representations or represented by	1028
973	Yifei Wang, Yuheng Chen, Wanting Wen, Yu Sheng,	model components, such as neurons and attention	1029
974	Linjing Li, and Daniel Dajun Zeng. 2024b. Unveil-	heads (Elhage et al., 2022a; Gurnee et al., 2023).	1030
975	ing factual recall behaviors of large language mod-	Circuits are subgraphs of the model’s computation	1031
976	els through knowledge neurons . In <i>Proceedings of</i>	graph responsible for implementing specific behav-	1032
977	<i>the 2024 Conference on Empirical Methods in Natu-</i>	iors (Wang et al., 2023b; Elhage et al., 2021).	1033
978	<i>ral Language Processing</i> , pages 7388–7402, Miami,	In this work, we focus on neuron-level feature-	1034
979	Florida, USA. Association for Computational Lin-	-based interpretability analysis to localize relation-	1035
980	guistics.	-specific neurons, which are responsible for encod-	1036
		ing and recalling specific types of factual knowl-	1037
981	Chris Wendler, Veniamin Veselovsky, Giovanni Monea,	edge. Existing studies have utilized various ap-	1038
982	and Robert West. 2024. Do llamas work in English?	proaches for neuron interpretation, each offering	1039
983	on the latent language of multilingual transformers .	unique advantages and limitations (Sajjad et al.,	1040
984	In <i>Proceedings of the 62nd Annual Meeting of the</i>	2022; Rai et al., 2024). The <i>visualization</i> method	1041
985	<i>Association for Computational Linguistics (Volume 1:</i>		
986	<i>Long Papers)</i> , pages 15366–15394, Bangkok, Thai-		
987	land. Association for Computational Linguistics.		
988	Robert F Woolson. 2005. Wilcoxon signed-rank test .		
989	<i>Encyclopedia of Biostatistics</i> , 8.		

(Olsson et al., 2022; Elhage et al., 2022a; 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. *Causation*-based 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. objects)** in each relation and the number of **overlapping subjects (resp. objects)** between any two relations in the identification prompt set $\mathcal{P}_{r_i}^{\text{det}}$ of the 7B model and the 13B model in Figure 10 and 11 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 $\mathcal{P}_{r_i}^{\text{eva}}$ (the 7B and 13B models share the

same evaluation set) in Figure 12. The results also show almost no entity overlap across different relations: among all relations, only `person_mother` and `person_father` share **one** subject, and the rest of the relations do not share any subject or object overlap. The entity analysis suggests that entities are not a confounding factor in our experiments, and the identified *RelSpec* neurons are only concerned with the relation itself, but not entities.

C Analysis On Gemma-7B

We perform a similar analysis on the Gemma-7B model (Gemma Team et al., 2024) as we do for the LLama-7B model. We first show how the identified 3,000 *RelSpec* neurons are distributed across layers for each relation in Figure 13. The trend is similar to what we observe in the 7B model (cf. Figure 1): the most of these neurons are located in the middle layers, but it is more evenly distributed across layers compared to the LLama families.

We show the intra-relation results in Figure 16. The results indicate that the identified *RelSpec* neurons are also effective in the Gemma-7B model: not only for the identification prompt set $\mathcal{P}_{r_i}^{\text{det}}$ but also for the held-out evaluation prompt set $\mathcal{P}_{r_i}^{\text{eva}}$, the deactivation of the neurons result in obvious accuracy drops, especially compared with the randomly deactivated neurons, indicating the existence of *RelSpec* neurons are held across model families.

We then demonstrate the effect of varying numbers of *RelSpec* neurons using the same numbers: 10, 50, 200, 500, 1,000, 3,000, 10,000, 20,000, and 50,000. Figure 14 and 15 present the results. The global trend is similar to what we observe for the LLama-7B model: the accuracy for a relation further drops when more of its *RelSpec* neurons are deactivated; until 3,000 or 10,000 neurons, the effect is almost only obvious for the concerned relation itself; after 10,000, deactivating more neurons results in a further drop in accuracy across all relations. This indicates the **neuron cumulativity** and **neuron versatility** can be observed across model families.

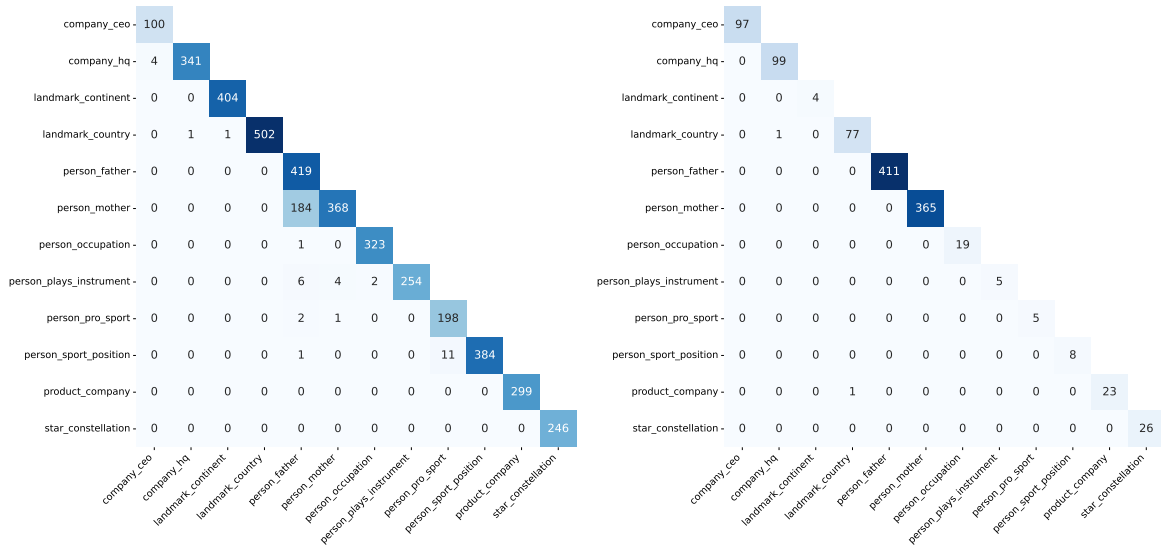


Figure 10: Subject (left) and object (right) overlap across 12 relations obtained from the **7B** model. The diagonal in each figure shows the number of distinct subjects or objects for each relation. It can be seen that factual knowledge from different relations has almost no entity overlap except for `person_mother` and `person_father`, which are mostly celebrities.

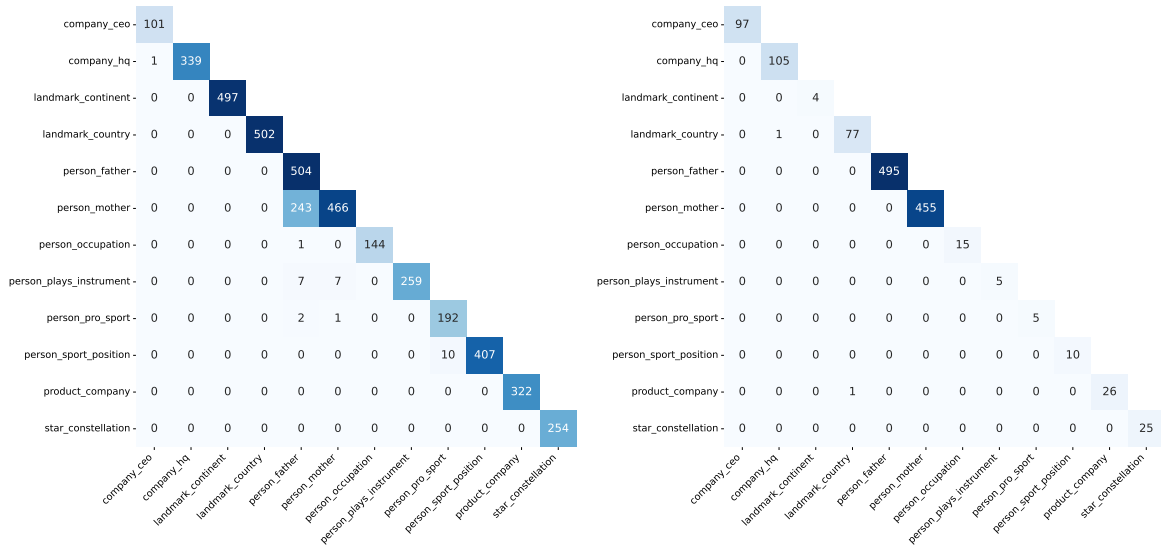


Figure 11: Subject (left) and object (right) overlap across 12 relations obtained from the **13B** model. The trend is very similar to that in the **7B** model: `person_mother` and `person_father` share many subjects.

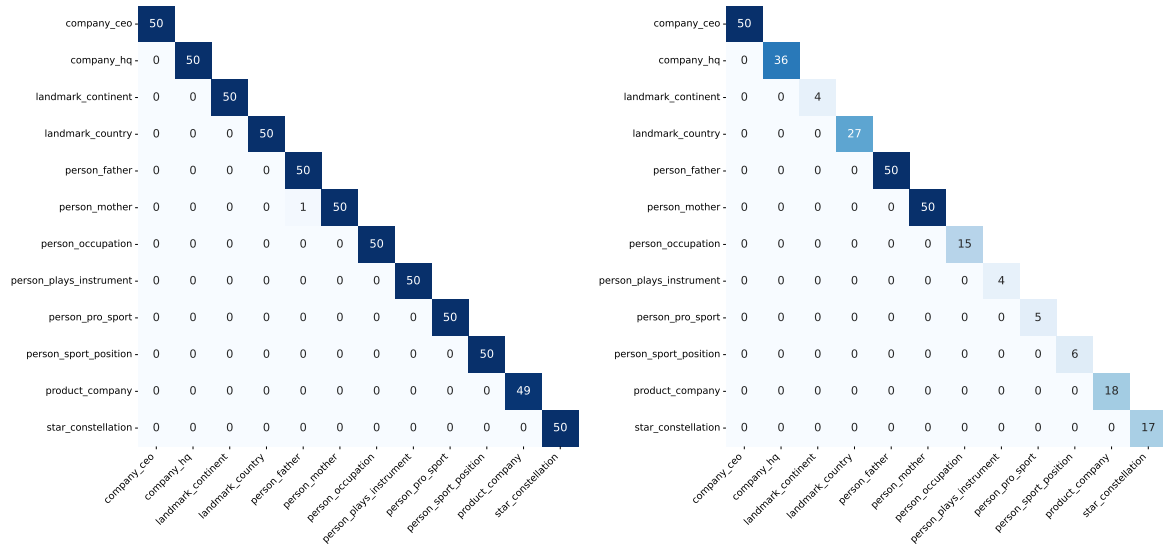


Figure 12: Subject (left) and object (right) overlap across 12 relations in the held-out evaluation prompt set $\mathcal{P}_{r_i}^{eva}$. Almost no two relations share any subjects or objects.

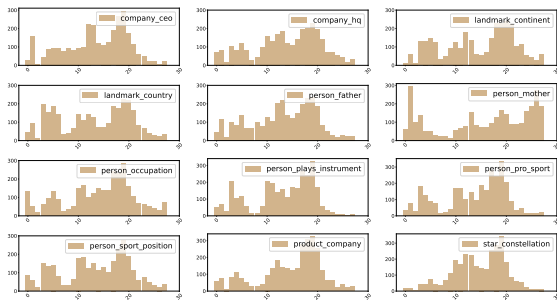


Figure 13: Distribution of *RelSpec* neurons across layers for the **Gemma-7B** model. Compared to the Llama-7B model in Figure 1, identified *RelSpec* neurons are more evenly distributed across layers. However, the majority of the population is still located in the middle layers.

D Analysis On the 13B Model

We perform a similar analysis on the 13B model as we do for the 7B model. We first show how the identified 3,000 *RelSpec* neurons are distributed across layers for each relation in Figure 17. The trend is similar to what we observe in the 7B model (cf. Figure 1). Most of the *RelSpec* neurons are distributed in the middle layers. Then we show the overlap of *RelSpec* neurons across relations in Figure 18. 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.

We then demonstrate the effect of varying numbers of *RelSpec* neurons using the same numbers: 10, 50, 200, 500, 1,000, 3,000, 10,000, 20,000, and 50,000. Figure 20 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 cumulativeness** is universal across models. *RelSpec* 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 well-memorized by the models and are less sensitive to the deactivation of *RelSpec* neurons.

Lastly, we show whether the identified *RelSpec* neurons from the 13B model are also multilingual. We use the same translated prompt sets as we use for the 7B model. We deactivate the 3,000 neurons identified using English and see how this affects the performance in other languages: German (**deu**), Spanish (**esp**), French (**fra**), Chinese (**zho**), and Japanese (**jpn**). The results are presented in Figure 19. We observe similar results as from the

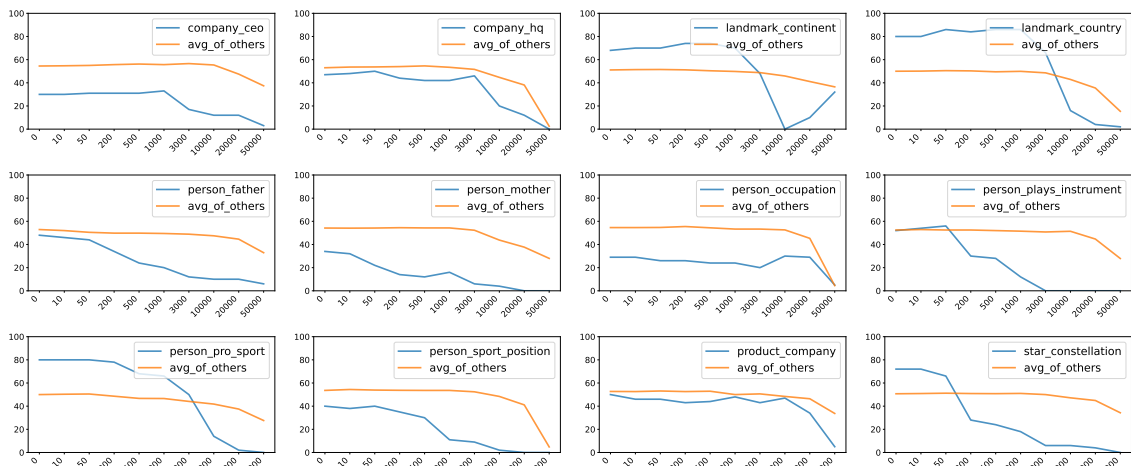


Figure 14: Influence of deactivating different numbers of *RelSpec* neurons for each relation (**Gemma-7B**). The variation of accuracy on the relation itself and the average accuracy on other relations is shown.

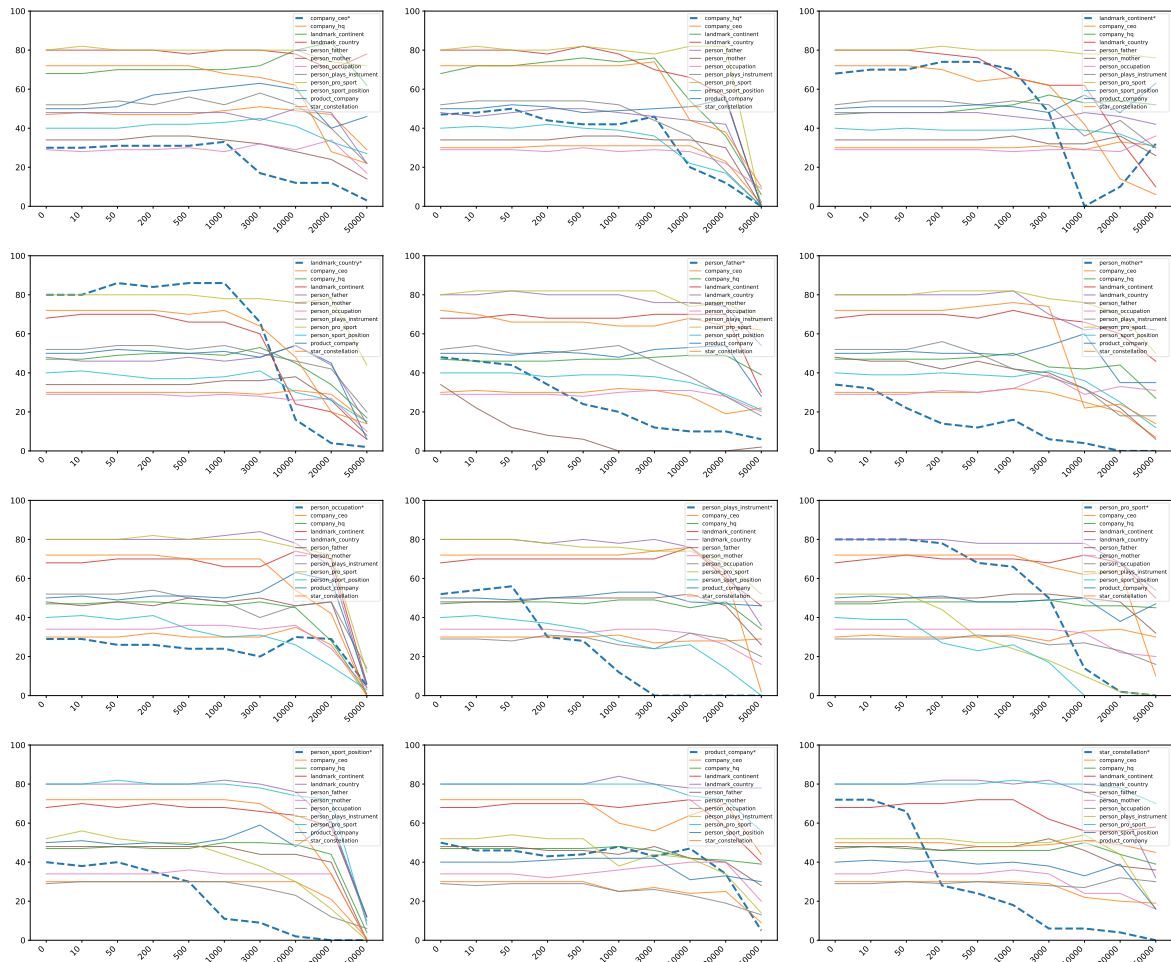


Figure 15: Influence of deactivating different numbers of *RelSpec* neurons in the **Gemma-7B** model for each relation. The variation of accuracy on the relation itself (noted with "*" and a dashed line style) and the accuracy on all other relations is shown in each figure.

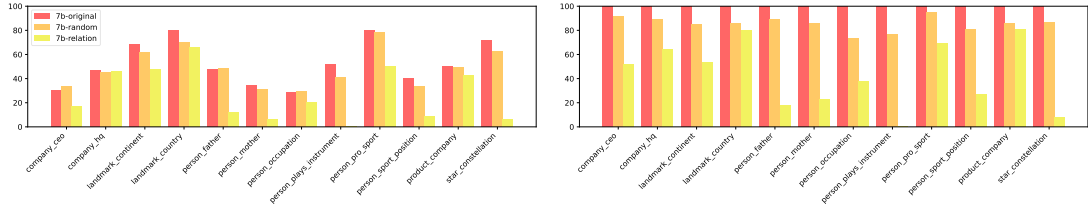


Figure 16: Intra-relation results on **Gemma-7B**. The left (resp. right) figure displays the results of held-out evaluation prompt set $\mathcal{P}_{r_i}^{\text{eva}}$ (resp. identification prompt set $\mathcal{P}_{r_i}^{\text{det}}$). We report the performance of the original model (without any deactivation), the model with 3,000 random neurons deactivated, and the model with relation neurons deactivated.

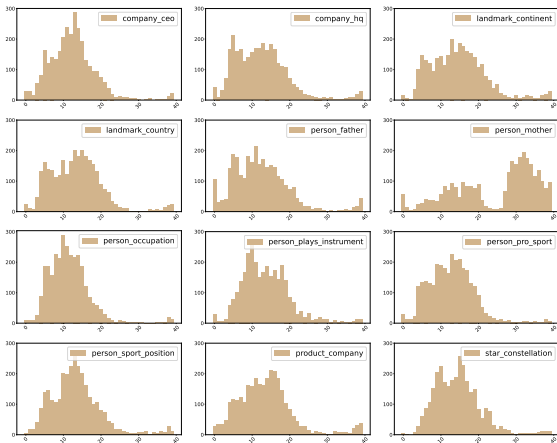


Figure 17: Distribution of *RelSpec* neurons across layers for the **13B** model. Similar to Figure 1, identified *RelSpec* neurons are mostly located in the middle layers, except for *person_mother*.

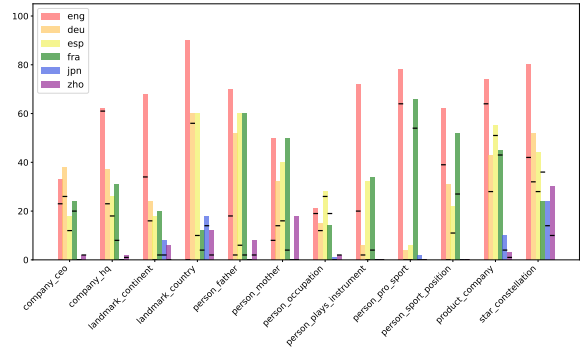


Figure 19: Accuracy on 12 relations across 6 languages from the **13B** model. The bars show the accuracy of the original model, with a horizontal line in each bar that indicates the performance after the deactivation of 3,000 *RelSpec* neurons.

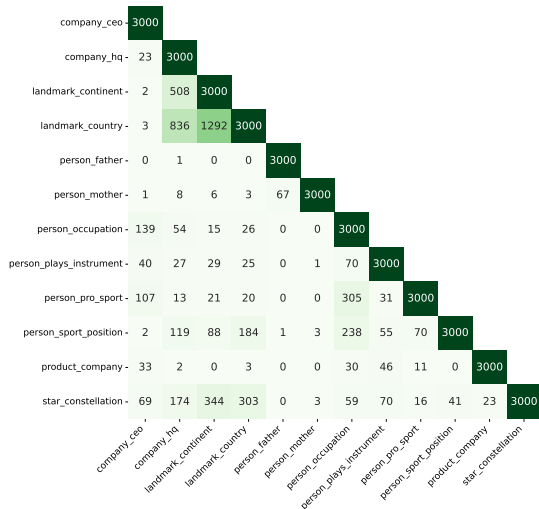


Figure 18: Neuron overlap of *RelSpec* neurons across 12 relations in the **13B** model. The overlap distribution is not similar to what we observe for the **7B** model shown in Figure 2, explaining the difference in inter-relation results (cf. Table 4).

7B model: when we deactivate *RelSpec* neurons identified using English prompts, many relations are influenced across languages, suggesting models with different sizes also have multilingual relational neurons. We also see some interesting counterexamples: deactivating *landmark_country* neurons completely deteriorates the relation in English but not in German. This indicates while some neurons have multilingual relational functionalities, there are still some relations dealt with in a language-specific manner.

1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192

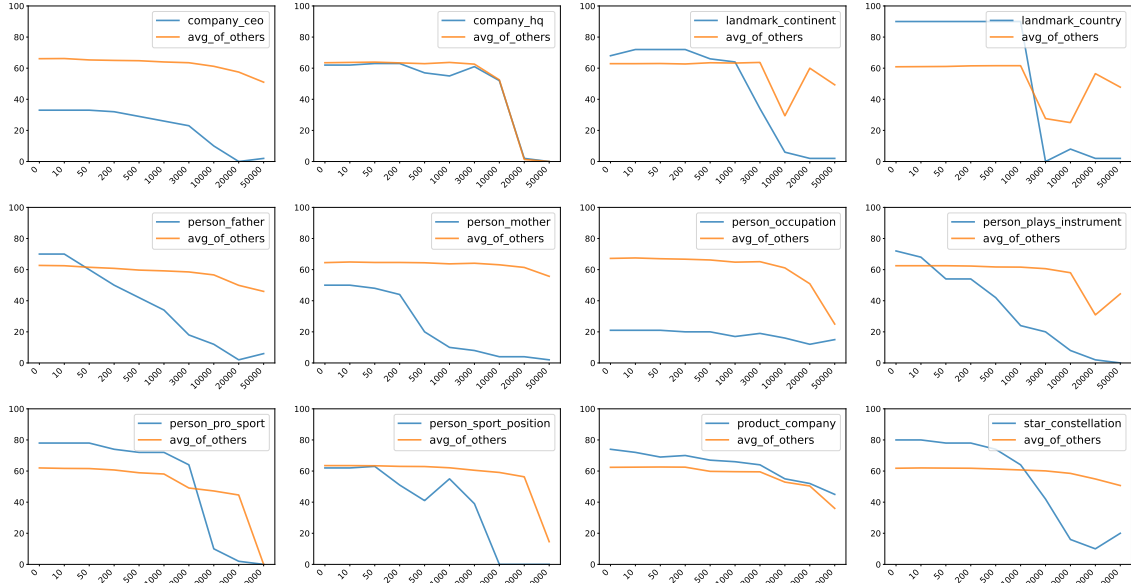


Figure 20: Influence of deactivating different numbers of *RelSpec* neurons for each relation (the **13B** model). The variation of accuracy on the relation itself and the average accuracy on other relations is shown.

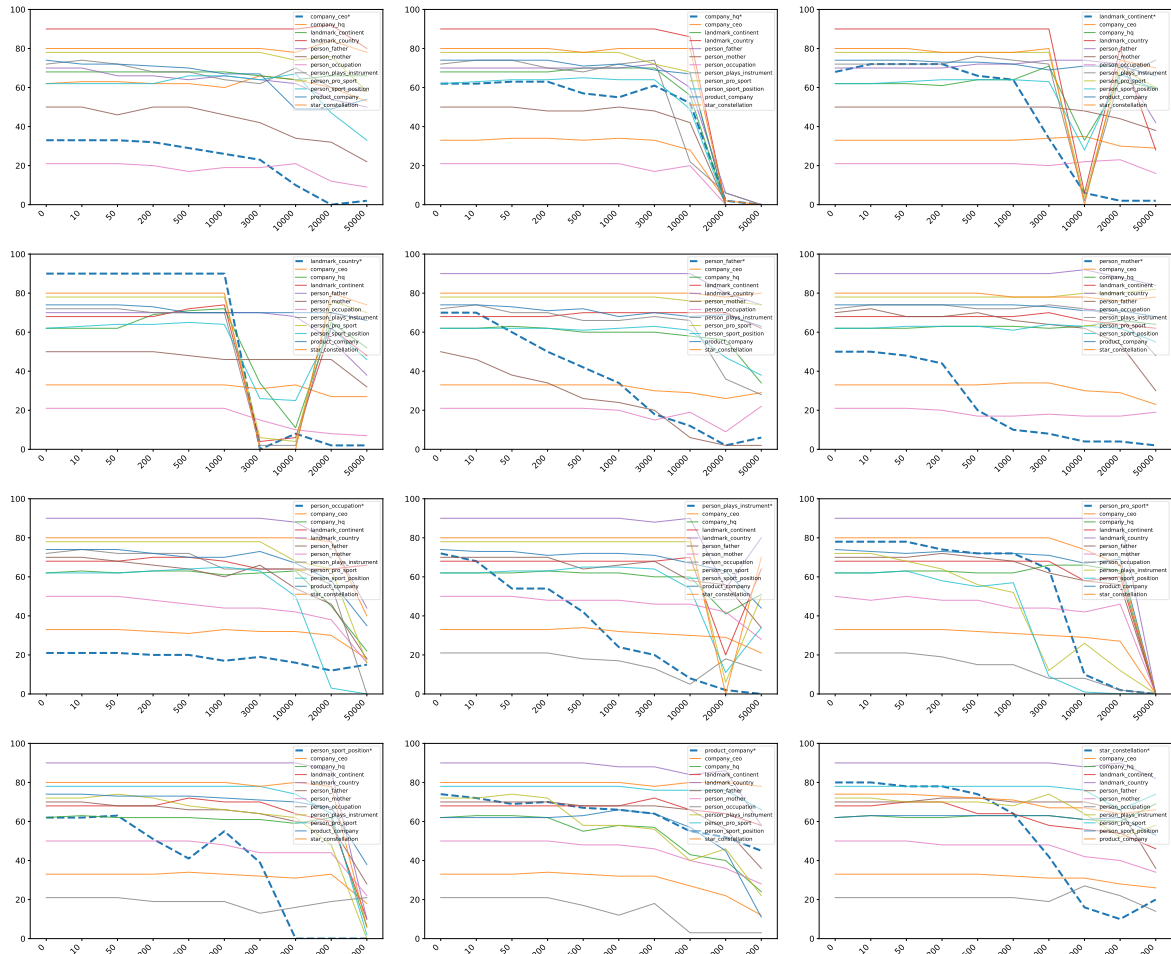


Figure 21: Influence of deactivating different numbers of *RelSpec* neurons in the **13B** model for each relation. The variation of accuracy on the relation itself (noted with “*” and a dashed line style) and the accuracy on all other relations is shown in each figure.

E 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 open-source 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 *RelSpec* neurons. **(b) sensitive facts**, for which the model is correct **before but not after** these neurons are deactivated.⁷ We then count how many documents in Dolma contain **both the subject and object** of each fact, calling this the *fact frequency*.⁸ Finally, we compute the average frequency for resilient and sensitive facts in each relation r_i , denoted respectively as $\text{group}_{r_i}^{(a)}$ and $\text{group}_{r_i}^{(b)}$.

Relative difference: $\text{diff}_{r_i} = \frac{\text{group}_{r_i}^{(b)} - \text{group}_{r_i}^{(a)}}{\text{group}_{r_i}^{(b)}}$ for each relation r_i is reported in Figure 23. 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.⁹ 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 *RelSpec* neurons**.

⁷We do not consider other numbers of *RelSpec* neurons because (1) if $\#\text{neurons} < 3,000$, there are not enough facts whose predictions change, and (2) if $\#\text{neurons} > 3,000$, facts belonging to other relations will also be influenced a lot.

⁸We use ElasticSearch API from WIMBD (Elazar et al., 2024) that allows for counting and searching in large corpora.

⁹We use a nonparametric test because the difference across relations does not follow a Gaussian distribution.

F Translation Process

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

Translating subject-object pairs. The first step concerns mapping entities, i.e., subject and object pairs, into the target language. The default way of doing this is by identifying if the entity is available in Wikidata and the target language using the Wikidata API.¹⁰ If the entity of interest is available in the target language, we directly take the entity name in that language. If the entity is not available, we then resort to Google Translate to translate the entity from English to the target language.¹¹ By performing this step, we obtain the subject-object pairs in all target languages and all relations.

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

G 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 24 and report the inter-relation results in Figure 25 (**all**), 26 (**self_attn**), 27 (**up_proj**), 28 (**gate_proj**), and 29 (**down_proj**).

¹⁰<https://www.wikidata.org/w/api.php>

¹¹<https://translation.googleapis.com/language/translate/v2>

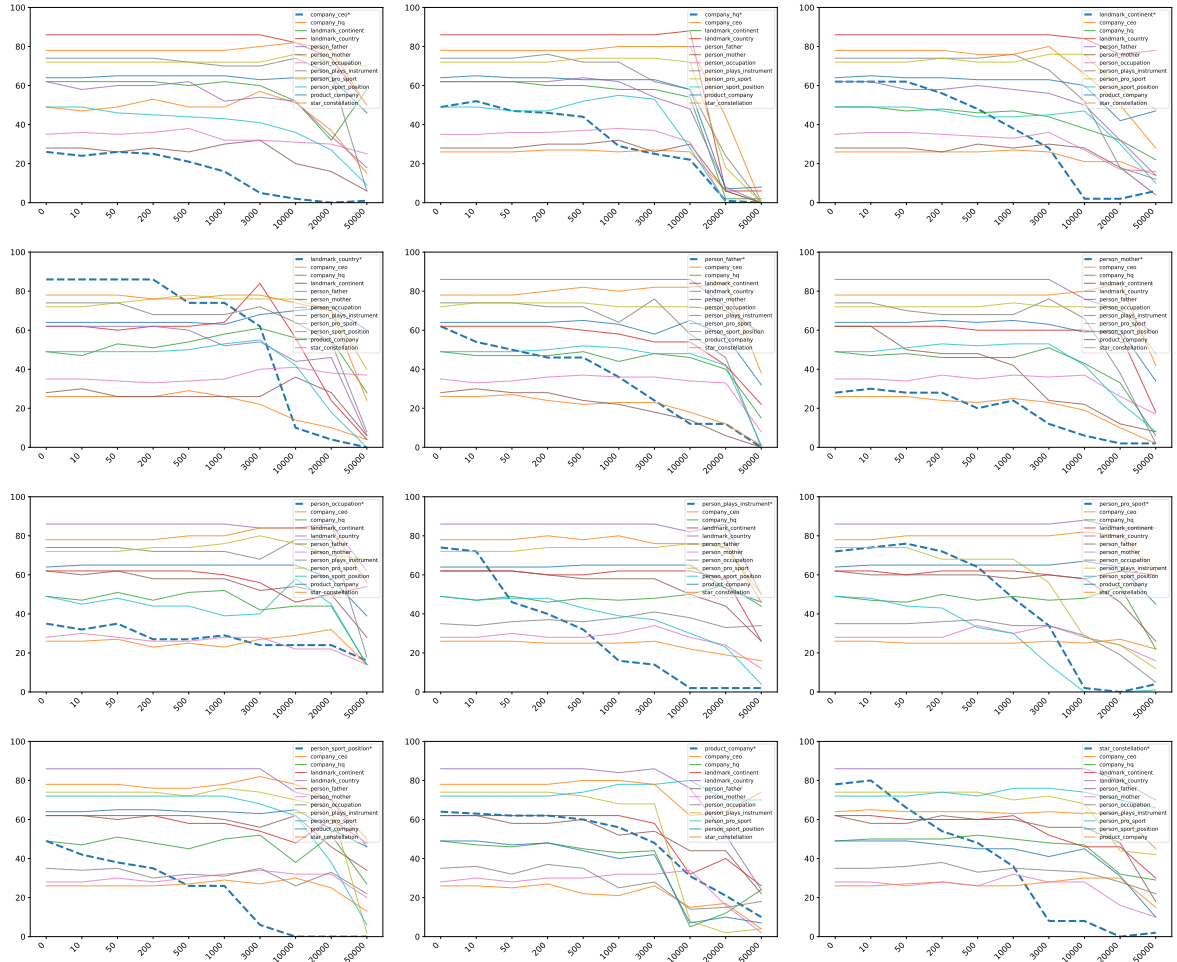


Figure 22: Influence of deactivating different numbers of *RelSpec* neurons in the **7B** model for each relation. The variation of accuracy on the relation itself (noted with “*” and a dashed line style) and the accuracy on all other relations is shown in each figure. Similar to Figure 5, increasing 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.

Relation	10–50		50–200		200–500		500–1000		1000–3000		3000–10000		10000–20000		20000–50000	
	#total	#affected	#total	#affected	#total	#affected	#total	#affected	#total	#affected	#total	#affected	#total	#affected	#total	#affected
company_ceo	1	0	3	2	5	0	7	2	11	2	3	3	2	0	0	0
company_hq	5	5	2	1	2	0	16	5	5	0	9	2	21	16	1	1
landmark_continent	0	0	4	4	4	2	5	0	6	2	13	6	0	0	0	0
landmark_country	0	0	1	1	6	0	2	0	6	0	26	5	3	0	2	0
person_father	3	1	2	0	0	0	5	0	6	2	6	0	2	1	6	4
person_mother	4	3	1	0	4	3	0	0	7	5	4	1	2	1	1	1
person_occupation	3	3	9	6	2	0	2	1	8	1	7	5	6	2	18	6
person_plays_instrument	13	11	7	2	5	0	8	0	3	0	6	0	0	0	0	0
person_pro_sport	0	0	2	1	4	1	9	0	8	0	16	0	1	0	0	0
person_sport_position	7	2	4	0	12	4	4	2	20	11	6	0	0	0	0	0
product_company	1	0	0	0	2	0	4	2	9	2	20	5	10	2	12	7
star_constellation	8	7	6	2	3	0	6	1	14	0	1	0	4	0	0	0

Table 2: Cumulative effect validation. For each neuron deactivation range, e.g., 1000-3000, the number of prompts where the model answers correctly in the smaller (1000) but not the larger range (3000) is denoted as column #total, and the number of prompts out of #total that are also affected, i.e., being answered wrongly, when deactivating the intermediate difference (2000 = 3000 - 1000) is denoted as #affected. #affected is usually much smaller than #total, indicating that neurons mostly act in a cumulative way and have no strong effect in isolation.

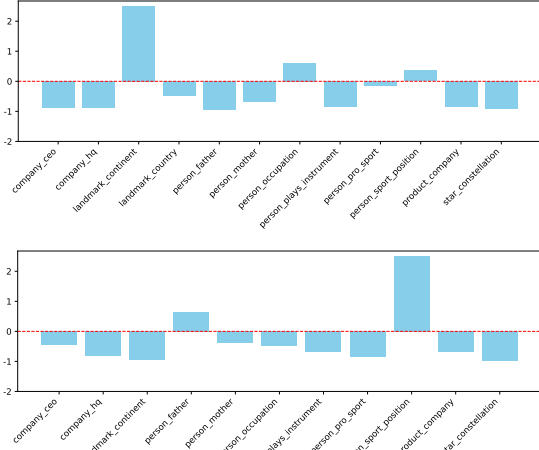


Figure 23: Relative difference between the average fact frequencies of the group (a) *resilient facts* and (b) *sensitive facts* for each relation in 7B (top) and 13B (bottom) models. Resilient facts generally appear more often than sensitive facts in most relations in the pertaining data.

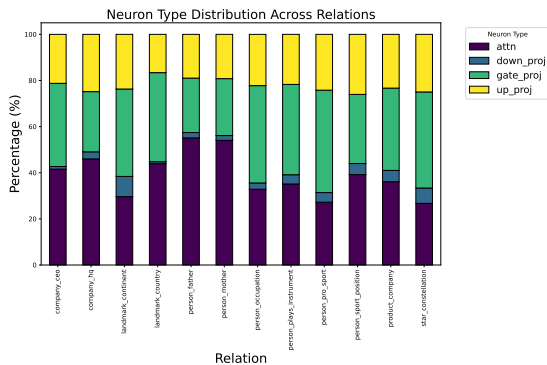


Figure 24: The distribution of the neuron types in the identified 3,000 neurons for the variety **all** across all relations.

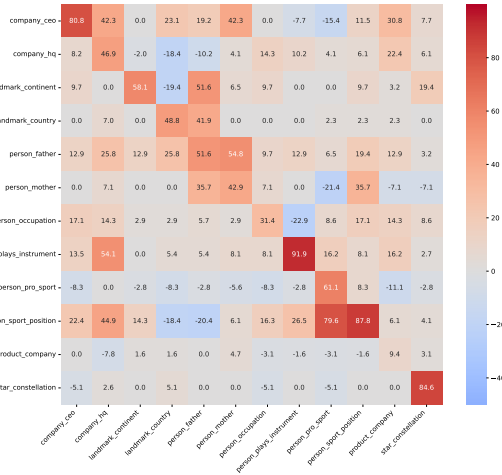


Figure 25: Inter-relation results of the 7B model when considering the neuron type variety as **all**.

According to the results, we observe that simply considering **self_attn** does not offer a consistent accuracy drop for the relation itself (by looking at the diagonal: some relations are not influenced too much). This can be explained by the fact that **self_attn** is shared across relations (as shown by Elhelo and Geva (2024)), and facts are mainly stored in the FFNs. Only considering **down_proj** offer similar results as **self_attn**. Interestingly, deactivating **up_proj** neurons does not influence all relations much in general, indicating it does not make sense to consider **up_proj** alone. Considering **all** or **gate_proj** neurons offer similar results compared to considering neurons in FFNs (shown in Figure 3). However, by considering neurons in FFNs (i.e., **up_proj**, **gate_proj** and **down_proj**), we see a more obvious inter-relation accuracy drop as shown on the diagonal in Figure 3. Therefore, our additional analysis supports our choice of considering neurons in FFNs.

H Concept-Specific Neurons

Concept-Relation Overlap in the 7B Model

Figure 30 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

1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312

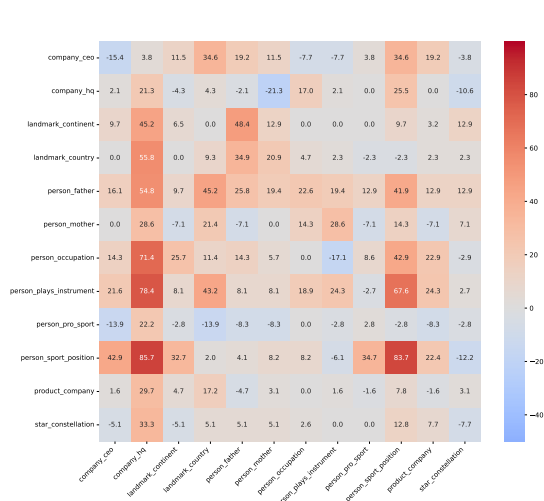


Figure 26: Inter-relation results of the 7B model when considering the neuron type variety as **self_attn**.

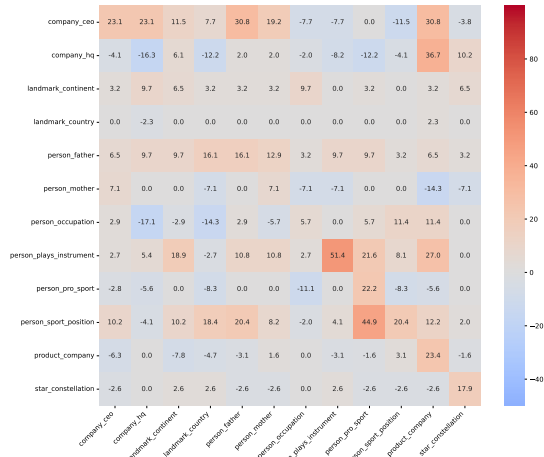


Figure 27: Inter-relation results of the 7B model when considering the neuron type variety as **up_proj**.

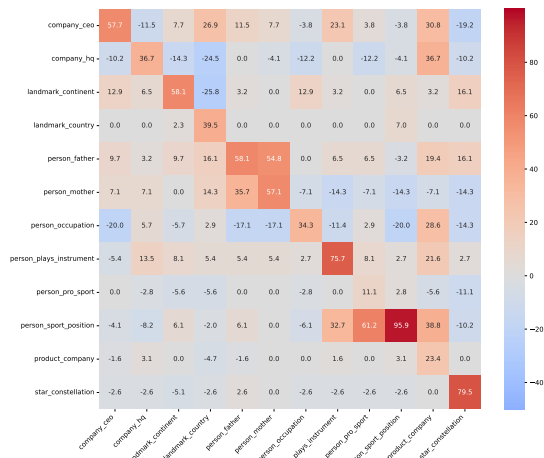


Figure 28: Inter-relation results of the 7B model when considering the neuron type variety as **gate_proj**.

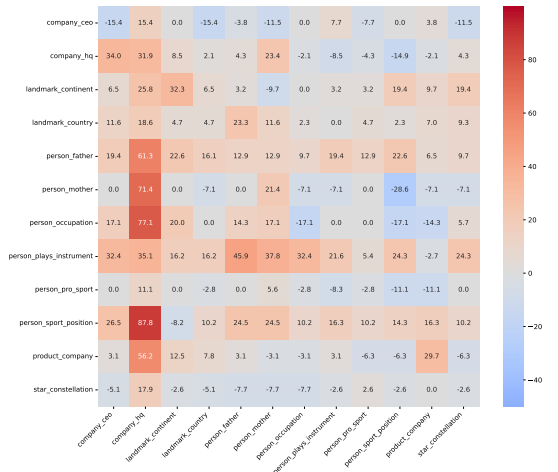


Figure 29: Inter-relation results of the 7B model when considering the neuron type variety as **down_proj**.

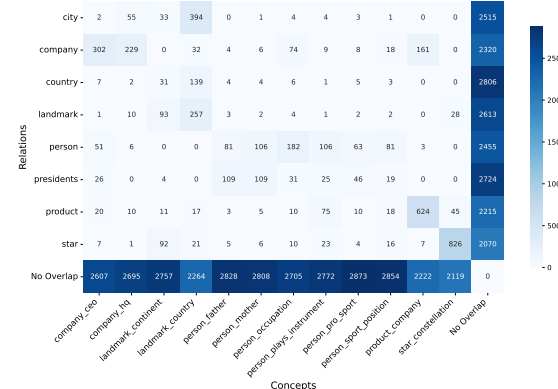


Figure 30: Overlap between the top 3000 identified neurons for each relation and concept in the 7B model.

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.

1313
1314
1315
1316
1317
1318

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:".

1319
1320
1321
1322
1323
1324
1325
1326

Figure 31 shows the results for the validation on these validation prompts for both models with the original accuracy score, a baseline that ablates 3000 neurons randomly, and the ablation of 3000 concept-specific neurons. Note that the impact of ablating a certain amount of expert neurons varies between concepts. The observed drop in performance due to the ablation of 3000 neurons for con-

1327
1328
1329
1330
1331
1332
1333
1334

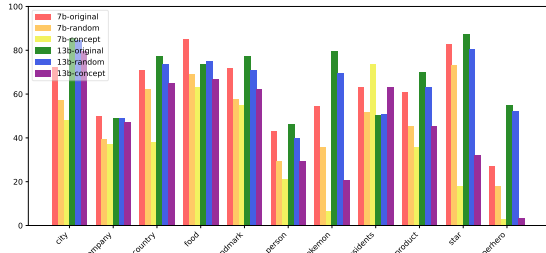


Figure 31: Accuracy results of evaluation prompts for 11 concepts in the 7b and 13b model. We report the performance of the original model (without any deactivation), e.g., 7b-original, the model with 3000 randomly deactivated neurons, e.g., 7b-random, and the model with deactivating the top 3000 identified concept-specific neurons, e.g., 7b-concept.

cepts like pokemon, superhero, and star is very large, while accuracy scores of other concepts in the 13b model, such as person appear stable, or even improve, e.g., presidents. We assume the **neuron cumulativity** also applies to the concept-specific neurons. That is, the knowledge on a specific concept is distributed over a much larger population of neurons, and further accuracy drop can be observed once more concept-specific neurons are deactivated – similar to what we observe for *RelSpec* neurons (cf. Figure 5). As only partial knowledge is withheld from the deactivation of 3000 concept-specific neurons, this might be too little knowledge to affect the facts concerning that concept (substantial knowledge on the concept is stored in the remaining neurons), resulting in only a small accuracy drop. Or, the 3000 concept-specific neurons store knowledge, though concerning the concept, unrelated to the prompts. For instance, the validation prompts of the concept presidents all demand **historical dates** as predicted answers, which is only one kind of knowledge that might be expected in connection with presidents. This phenomenon actually aligns with our neuron interference hypothesis: deactivating neurons that store unhelpful knowledge can less confuse the model, therefore improving the performance.

I Experimental Environment

We run all experiments on NVIDIA RTX A6000 GPUs. The Python environment we use is the same as Kojima et al. (2024).¹²

¹²Kojima et al. (2024)’s GitHub repository is available at https://github.com/kojima-takeshi188/lang_neuron

J Error Analysis

We manually verified the prompts in each relation that the model could answer correctly originally, but failed to answer correctly when 3,000 *RelSpec* neurons were deactivated (cf. §4.2). The three most common incorrect responses (regarded as *systematic errors*) are listed in Table 3.

After we deactivate the *RelSpec* neurons, we can see that the model appears to lose its ability to recall the correct object. Instead, the model frequently answers with meaningless answers that start with tokens such as “A.” or “The”, or simply repeats the given prompt. We showcase representative examples of each phenomenon in Table 4, Table 5, and Table 6. The results strongly indicate that the model loses its ability to capture relational semantics, resulting in increasingly noisy outputs after the deactivation of *RelSpec* neurons.

K Prompt Templates

We show the actual prompt templates (with an object-subject example) we use for each relation across 6 considered languages: company_ceo in Table 7, company_hq in Table 8, landmark_continent in Table 9, landmark_country in Table 10, person_father in Table 11, person_mother in Table 12, person_occupation in Table 13, person_plays_instrument in Table 14, person_pro_sport in Table 15, person_sport_position in Table 16, product_company in Table 17, and star_constellation in Table 18.

Relation	Repeat Prompt	Answer with “The”	Answer with “A.”	Total Number
company_ceo	47.8%	8.7%	34.8%	23
company_hq	46.2%	46.2%	0%	26
landmark_continent	17.6%	5.9%	0%	17
landmark_country	69.3%	0%	0%	13
person_father	84.2%	5.3%	0%	19
person_mother	70%	20%	0%	10
person_occupation	93.3%	0%	0%	15
person_plays_instrument	51.6%	29%	0%	31
person_pro_sport	25%	15%	0%	20
person_sport_position	18.6%	11.6%	44.2%	43
product_company	70.1%	11.8%	0%	17
star_constellation	88.6%	5.7%	0%	35

Table 3: Most common incorrect answers generated by LLama-7b after deactivating 3,000 *RelSpec* neurons.

Subject-Object Pair	Prompt	Expected Output	Model Response	Deactivation
(Panasonic Corporation, Kazuhiro Tsuga)	Panasonic Corporation’s CEO is? Answer:	Kazuhiro Tsuga	Kazuhiro Tsuga	No
			[Pan]asonic Corporation’s CEO is:\nPanasonic Corporation’s CEO	Yes

Table 4: Model answers by repeating the prompt after deactivating *RelSpec* neurons. We changed the output length from 2 tokens to 20 tokens to observe the complete output. The part enclosed in “[]” is the first 2 tokens of the output. The triple (Panasonic, company_ceo, Kazuhiro Tsuga) is selected for demonstration.

Subject-Object Pair	Prompt	Expected Output	Model Response	Deactivation
(Pagan Federation, London)	Pagan Federation is headquartered in the city of? Answer:	London	London	No
			[The] Pagan Federation is a British organisation that represents the interests of Pagans and other Ne	Yes

Table 5: Model answers with “The” after deactivating *RelSpec* neurons. We changed the output length from 2 tokens to 20 tokens to observe the complete output. The part enclosed in “[]” is the first 2 tokens of the output. The triple (Pagan Federation, company_hq, London) is selected for demonstration.

Subject-Object Pair	Prompt	Expected Output	Model Response	Deactivation
(Damon Huard, quarterback)	Damon Huard plays in the position of a? Answer:	quarterback	Quarterback	No
			[A.] \nDamon Huard plays in the position of a? \nAnswer: A.	Yes

Table 6: Model answers with “A.” after deactivating *RelSpec* neurons. We changed the output length from 2 tokens to 20 tokens to observe the complete output. The part enclosed in “[]” is the first 2 tokens of the output. The triple (Damon Huard, person_sport_position, quarterback) is selected for demonstration..

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: ----- The name of the CEO of Panasonic Corporation is? Answer: Who is the CEO of Panasonic Corporation? Their name 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, responde directamente por su nombre. El nombre del director general de Panasonic es Por favor, responde directamente por su nombre. ¿Quién es el director general de Panasonic? Su nombre es	Kazuhiro Tsuga
French	(Panasonic, Kazuhiro Tsuga)	Veillez répondre directement avec le nom. Le nom du président-directeur général de Panasonic est Veillez répondre directement avec le nom. Le PDG de Panasonic est nommé	Kazuhiro Tsuga
Japanese	(パ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 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 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 7: 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. The prompt shown below the dashed line is the new template introduced for the experiment described in §5.3.

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: ----- Where are the headquarters of Cadillac It is 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どこですか 「キャデラック」の本社はど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 8: 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. The prompt shown below the dashed line is the new template introduced for the experiment described in §5.3.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Elbe, Europe)	Elbe is on the continent of? Answer: What continent is Elbe on? It is on? 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)	Veillez 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 9: 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. The prompt shown below the dashed line is the new template introduced for the experiment described in §5.3.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Namba Station, Japan)	Namba Station is in the country of? Answer: What country is Namba Station in? It is in? 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 10: 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. The prompt shown below the dashed line is the new template introduced for the experiment described in §5.3.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Ronald Reagan, Jack Reagan)	Ronald Reagan's father is named? Answer: Who is Ronald Reagan's father? Their 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 11: 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. The prompt shown below the dashed line is the new template introduced for the experiment described in §5.3.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Demi Moore, Virginia King)	Demi Moore's mother is named? Answer: ----- Name of mother of Demi Moore is? Answer: Who is Demi Moore's mother? Their 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 12: 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. The prompt shown below the dashed line is the new template introduced for the experiment described in §5.3.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Martin Burrell, politician)	Martin Burrell works as a? Answer: By profession, Martin Burrell is a? Answer: ----- Martin Burrell works professionally as 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)	Veillez répondre directement par le nom de votre profession. Le nom de la profession de Martin Burrell est Veillez répondre directement par le nom de votre profession. Martin Burrell travaille professionnellement 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 13: 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. The prompt shown below the dashed line is the new template introduced for the experiment described in §5.3.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Anson Funderburgh, guitar)	Anson Funderburgh plays the instrument of? Answer:----- What instrument does Anson Funderburgh play? They play the? Answer: The instrument that Anson Funderburgh plays is called the? Answer:	guitar
German	(Anson Funderburgh, Gitarre)	Bitte geben Sie den Namen des Instruments direkt 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 Funderburgh? Tocan el	guitarra
French	(Anson Funderburgh, guitare)	Veuillez répondre directement au nom de l'instrument. De quel instrument joue Anson Funderburgh ? Ils jouent du	guitare
Japanese	(アンソン・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 14: 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. The prompt shown below the dashed line is the new template introduced for the experiment described in §5.3.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Frédéric Piquionne, soccer)	Frédéric Piquionne plays the sport of? Answer:----- What sport does Frédéric Piquionne play? They play? Answer: Frédéric Piquionne plays professionally in 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・ピキオンヌ, サッカー)	サッカー、バスケットボールなど、スポーツの名前を直接答えてください。フレデリック・ピキオンヌはどのスポーツをしますか？彼らは（スポーツ名）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 15: 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. The prompt shown below the dashed line is the new template introduced for the experiment described in §5.3.

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: ----- Which position does Ju Yingzhi play? They play as a? Answer: In their sport, Ju Yingzhi plays in the position of 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)	彼がプレーするスポーツでは、ジュ・インジーのポジションは ジュ・インジー競技場のポジションは	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 16: 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. The prompt shown below the dashed line is the new template introduced for the experiment described in §5.3.

Language	Subject-Object Pair	Prompt	Expected Output
English	(Jeep Grand Cherokee, Chrysler)	Jeep Grand Cherokee was created by which company? Answer: Jeep Grand Cherokee is a product of which company? Answer: ----- Which company developed Jeep Grand Cherokee? It was developed by? 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 Cherokee 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	(ジープ・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。ジープ・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。ジープ・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 17: 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. The prompt shown below the dashed line is the new template introduced for the experiment described in §5.3.

Language	Subject-Object Pair	Prompt	Expected Output
English	(50 Persei E, Perseus)	50 Persei E is part of the constellation named? Answer:----- What is the name of the constellation that 50 Persei E is part of? It is part of? Answer: What is the name of the constellation that 50 Persei E belongs to? It belongs to? 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 denominada	Perseus
French	(50 Persei E, Persée)	Le nom de la constellation dans laquelle se trouve 50 Persei E est	Persée
Japanese	(50 ペルセウス座E, ペルセウス座)	50 ペルセウス座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 18: 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. The prompt shown below the dashed line is the new template introduced for the experiment described in §5.3.