

# On the Multilingual Ability of Decoder-based Pre-trained Language Models: Finding and Controlling Language-Specific Neurons

Takeshi Kojima, Itsuki Okimura, Yusuke Iwasawa,  
Hitomi Yanaka, Yutaka Matsuo

The University of Tokyo  
t.kojima@weblab.t.u-tokyo.ac.jp

## Abstract

Current decoder-based pre-trained language models (PLMs) successfully demonstrate multilingual capabilities. However, it is unclear how these models handle multilingualism. We analyze the neuron-level internal behavior of multilingual decoder-based PLMs, specifically examining the existence of neurons that fire “uniquely for each language” within decoder-only multilingual PLMs. We analyze six languages: English, German, French, Spanish, Chinese, and Japanese, and show that language-specific neurons are unique, with a slight overlap (< 5%) between languages. These neurons are mainly distributed in the models’ first and last few layers. This trend remains consistent across languages and models. Additionally, we tamper with less than 1% of the total neurons in each model during inference and demonstrate that tampering with a few language-specific neurons drastically changes the probability of target language occurrence in text generation.<sup>1</sup>

## 1 Introduction

Recent studies have frequently demonstrated the excellent multilingual abilities of pre-trained language models (PLMs). Some PLMs explicitly mix multilingual language corpus for pre-training (Lin et al., 2021; Scao et al., 2022), whereas others mainly use an English-dominant text corpus, with the unintentional inclusion of a low percentage of multiple language texts, which results in the acquisition of multilingual skills, such as Llama2 (Touvron et al., 2023). How do they exhibit multilingual abilities?

Prior studies have focused on language-universal neurons activated across multilingual inputs, mainly focusing on encoder-based PLMs (Antverg and Belinkov, 2022; Stańczak et al., 2022; Chen et al., 2023; Stańczak et al., 2023; Varda and

<sup>1</sup>Code and model-generated texts are available at [https://github.com/kojima-takeshi188/lang\\_neuron](https://github.com/kojima-takeshi188/lang_neuron)

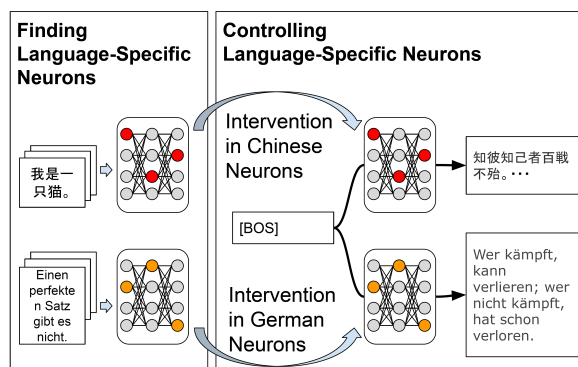


Figure 1: Overview of our proposal. (Left) Finding language-specific neurons that tend to be activated for a target language. (Right) Controlling the detected language-specific neurons by forcing their activation during inference to manipulate the probability of target language occurrence.

Marelli, 2023). In contrast to encoder-based models, which might be sufficient to abstract inputs, decoder-based PLMs need to recover the language-specific information in the later part of the generation. Therefore, language-specific processing within these models should be a more complex and essential functionality compared to the encoder-based ones. However, few studies have focused on the existence and activity of language-specific neurons in decoder-based PLMs (See Section 2).

This study examines the behavior of language-specific neurons in decoder-based PLMs. Specifically, we analyze multiple decoder-based PLMs, including XGLM (564M, 1.7B, 2.9B), BLOOM (560M, 1.7B, 3B), and Llama2 (7B, 15B), for six languages (English, German, French, Spanish, Chinese, and Japanese). To investigate language-specific neurons, we adopt an approach proposed by Cuadros et al. (2022), which finds neurons that activate on a certain group of sentences (*Positive* sentences) but do not activate on other groups (*Negative* sentences). We treat the target language texts as positive and any other language as negative, iden-

tifying language-specific neurons that statistically activate positive sentences (See Section 3). The experimental results demonstrate that the identified language-specific neurons are mainly distributed in the first and last few layers of the model. This trend remains consistent across multiple languages and model variants. To verify the effect of the neurons, we intervene in language-specific neurons in the model during inference, showing that they can drastically change the probability of the target language occurrence during text generation (See Section 4).

## 2 Related Work

Previous studies analyzed the internal behavior of PLMs for multilingual tasks by observing the activation of their neurons. Several studies have found that language-universal neurons are activated across multilingual inputs in encoder-based PLMs (mainly M-BERT, see Pires et al. (2019)) for various task settings, including syntactic or factual knowledge tasks (Antverg and Belinkov, 2022; Stańczak et al., 2022; Chen et al., 2023; Stańczak et al., 2023; Varda and Marelli, 2023). However, studies of encoder-based PLMs have not focused on the identification of language-specific neurons. Mueller et al. (2022); Bau et al. (2019) studied decoder-based language models to find the multilingually shared neurons. Similar to the encoder-based PLMs, limited research has focused on the existence and activity of language-specific neurons in decoder-based language models.

Various methods can be used to identify and control neurons (Sajjad et al., 2022). Several studies have identified and intervened in neurons for effective word editing or classification (Mu and Andreas, 2020; Dai et al., 2022; Mueller et al., 2022; Chen et al., 2023; Varda and Marelli, 2023). In contrast, few studies have investigated the identification and intervention of neurons for full-text generation for a desired concept, for example, Bau et al. (2019) for LSTM models and Cuadros et al. (2022) for pre-trained transformer models. Cuadros et al. (2022) have proposed an approach for controlling text generation on Transformer-based decoder models, and proven its effectiveness. Therefore, we use it as an analytical tool in our experiments with decoder-based PLMs.

## 3 Method

We identified and controlled neurons specific to each language based on the approach of Cuadros

et al. (2022), with appropriate modifications for our experiments. This approach was originally developed to identify and control neurons that respond to specific word-level concepts, such as homographs or gender biases. However, we aimed to find neurons that grasp broader sentence-level and language-specific concepts; therefore, we modified the original approach for our purpose.

### 3.1 Finding Language-specific Neurons

First, we prepared text for each language. We considered a set of  $|L|$  languages. For each language  $l \in L$ , we prepared  $N_l$  texts, which resulted in  $N = N_1 + \dots + N_l + \dots + N_{|L|}$  texts for all the languages. Let  $x = \{x_i\}_{i=1}^N$  be the set of all the texts. Our goal was to find neurons that activate text in the target language  $l_t \in L$  but do not activate text in other languages  $L \setminus \{l_t\}$ . For each text  $x_i \in x$ , we assigned a label  $b_i = 1$  if the text was in the target language (i.e.,  $l = l_t$ ); otherwise,  $b_i = 0$ . Therefore, we had

$$N = N_{l_t}^+ + N_{l_t}^- \quad (1)$$

sentences, where  $N_{l_t}^+$  positive sentences consisted of texts in the target language  $l$  (i.e.,  $b_i = 1$ ) and  $N_{l_t}^-$  negative sentences consisted of texts in other languages (i.e.,  $b_i = 0$ ). For example, if the target language  $l_t$  was French, French texts were assigned label 1, and texts in other languages, such as German and Chinese were assigned label 0.

Second, we observed the activation value of each neuron inside the model given the input text. We assigned a unique index  $m \in M$  to each neuron.  $|M|$  denotes the total number of neurons in the model. Let  $z_{m,i} \in z_m$  be the output value of neuron  $m$  when text  $x_i \in x$  is provided as an input to the model. Here, we explain in detail how this value can be calculated. Specifically, text  $x_i$  is composed of a sequence of  $T$  tokens  $x_i = \{w_{i,1}, \dots, w_{i,t}, \dots, w_{i,T}\}$ . Therefore, given the input text, there exist  $T$  output values  $\{z_{m,i,1}, \dots, z_{m,i,j}, \dots, z_{m,i,T}\}$  for neuron  $m$  inside the decoder-based Transformer model. We take the average of the  $T$  neuron outputs to summarize the output value of neuron  $m$  for the text  $i$ .

$$z_{m,i} = f(z_{m,i,1}, \dots, z_{m,i,t}, \dots, z_{m,i,T}), \quad (2)$$

where  $f$  is the aggregation function of the average operator. While the original approach (Cuadros et al., 2022) defines  $f$  as a max-pooling operator, our approach defines  $f$  as an average operator to

identify neurons that consistently activate across tokens for language identification purposes. The output values of the [PAD] token position are excluded from the aggregation as an exception because they are regarded as noise.

Third, language-specific neurons were identified. We regarded the dataset  $\{x_i, b_i, z_{m,i}\}_{i=1}^N$  as the prediction task samples. Specifically, we regarded texts  $\{x_i\}_{i=1}^N$  as inputs to the model, labels  $\{b_i\}_{i=1}^N$  as their ground truth, and the output values of neurons  $\{z_{m,i}\}_{i=1}^N$  as the prediction scores of the ground truth. We can measure the performance of neuron  $m$  for the task using its average precision ( $AP_m = AP(z_m, b) \in [0, 1]$ ), which is the area under the precision-recall curve with different prediction thresholds. We measured  $AP_m$  for all neurons and ordered them in descending order.

In the original approach, only the top- $k$  neurons in descending order were defined as identified neurons. However, this only considers strong positive correlations (i.e., the top- $k$  highest AP neurons) with labels, leaving out strong negative correlations (i.e., the top- $k$  lowest AP neurons) with labels. We hypothesize that not only the top- $k$  neurons but also the bottom- $k$  neurons are strongly related to a specific language. Therefore, we extended the original approach by considering not only the top- $k$  neurons but also the bottom- $k$  neurons, defining them as language-specific neurons. We validate our claim experimentally in Section 4. We set  $k = 1000$  as the default value across the experiments. Note that the neurons at the input layer (word embeddings) and output layer (projection layers) were excluded from the measurement because it is clear that these layers consist of language-specific modules: they consist of language-specific characters or (sub-)words.

### 3.2 Controlling Language-specific Neurons

We controlled text generation by overriding the output values of the top- $k$  and bottom- $k$  neurons with fixed values during inference. Specifically, we calculated the fixed value for each neuron  $m$  as follows:

$$\bar{z}_m = \text{Median}(\{z_m | b = 1\}). \quad (3)$$

This is the median of the neuron outputs for the target language texts. During inference, we intervened in the top- $k$  and bottom- $k$  neurons by replacing their outputs with fixed values in the forward pass and observed whether the models generated texts in the target language.

Model	# Params	# Layers	# Neurons
XGLM	564M	24	221,184
	1.7B	24	442,368
	2.9B	48	884,736
BLOOM	560M	24	221,184
	1.7B	24	442,368
	3B	30	691,200
Llama2	7B	32	1,359,872
	13B	40	2,129,920

Table 1: Model list used for the experiments.

	en	de	fr	es	zh	ja
XGLM	49.0	5.4	4.7	5.3	8.1	4.0
BLOOM	30.0	-	12.9	10.8	16.2	-
Llama2	89.7	0.2	0.2	0.1	0.1	0.1

Table 2: Distribution of languages in pre-training data.

## 4 Experiment Settings

### 4.1 Models

XGLM (Lin et al., 2021), BLOOM (Scao et al., 2022), and Llama2 (Touvron et al., 2023) were used in the experiments. XGLM and BLOOM are explicitly referred to as multilingual language models. By contrast, Llama2 was trained almost entirely on an English text corpus, with minimal inclusion of other languages. Table 1 lists the models used in the experiments. All the models were downloaded from HuggingFace (Wolf et al., 2019). Table 2 describes the distribution of languages in the pretraining dataset for each model<sup>2</sup>.

### 4.2 Datasets

The following six languages were used in the experiment: English (en), German (de), French (fr), Spanish (es), Chinese (zh), and Japanese (ja). These six languages are frequently targeted in prior studies of multilingual language models; they are among the top seven languages in terms of the percentage of languages included in the XGLM pre-training data (Lin et al., 2021). Owing to the limitations of the available computer resources, the number of languages analyzed was limited to six, as described in the limitations section.

To create a language-specific text corpus, we combined two datasets, PAWS-X (Yang et al., 2019) and FLORES-200 (Costa-jussà et al., 2022).

<sup>2</sup>XGLM information is cited from <https://huggingface.co/facebook/xglm-2.9B>. BLOOM information is cited from <https://huggingface.co/bigscience/bloom#languages>.

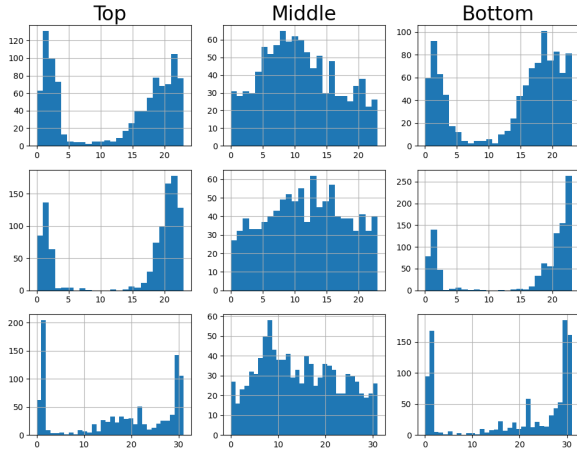


Figure 2: Distribution of Top, Middle, Bottom-1000 neurons across layers. 1st row:(XGLM-564M, de). 2nd row:(BLOOM-1b7, fr). 3rd row:(Llama2-13b, zh).

PAWS-X is a dataset for paraphrase identification between two texts for seven languages, including the aforementioned languages. FLORES-200 is a dataset of machine translation tasks for more than 200 languages. The sample sentences in these tasks were of good quality, had a wide variety of text types, and covered the six languages required for our experiments. Therefore, a combination of these factors was used. For this experiment, texts were randomly sampled in a 1:1 ratio from the two datasets to create ground-truth texts for each language.

Following Cuadros et al. (2022), we apply a setting of  $N_{l_t}^- > N_{l_t}^+$  to account for the much larger variance of negative than positive examples. Negative samples contain texts from five language, whereas positive samples contain texts from only one language. Specifically, we prepared 500 texts for each language, totaling 3000 texts for all six languages. As Cuadros et al. (2022) pointed out, the choice of  $N_{l_t}^+$  and  $N_{l_t}^-$  is arbitrary, usually a tradeoff between the computing resources available and the quality of the representation required. In addition, Cuadros et al. (2022) set the sample sizes of both positive and negative examples between 100 and 1000. Therefore, we considered 500 to be a reasonable value. After specifying a target language, we identified language-specific neurons for the target language using the method described in Section 3.

## 5 Results and Discussion

The experimental results are summarized in this section. See the Appendix for the full results across

	de	en	es	fr	ja	zh
de	2000	41	74	39	44	34
en	41	2000	34	41	49	40
es	74	34	2000	57	77	22
fr	39	41	57	2000	21	93
ja	44	49	77	21	2000	27
zh	34	40	22	93	27	2000

Table 3: The number of overlapping language-specific neurons between languages (XGLM-564M).

all models and languages.

### 5.1 Finding Language-specific Neurons

We identified language-specific neurons using the method described in Section 3. Figure 2 shows histograms of the identified neurons for each layer in each model. Most of the top-1000 neurons with higher AP orders are distributed in the first and last few layers of the model. Similarly, most of the bottom-1000 neurons were distributed in the first and last few layers. In contrast, the middle-1000 neurons (around the median) in AP order were distributed mainly in the middle layers. It was also found that the distributional property was the same across languages, model sizes, and model variants.

In addition, it was confirmed that language-specific neurons that fire in one language have little overlap with those of all other languages. Table 3 shows the cross tables of the number of overlapping neurons between the six languages, indicating an overlap of less than 5 % for every language pair.

The following interpretation is possible regarding the inner behavior of decoder-based PLMs based on the findings in Figure 2 and the prior study of multilingual models (Muller et al., 2021), which suggests that encoder-based PLMs process cross-lingual transfer in lower layers and language-agnostic semantic concepts (e.g., task prediction) in higher layers: The first few layers of decoder-based PLMs mainly process cross-lingual transfers to transform the lexical or syntax representations of each language into language-independent semantic representations. The middle layers of the models are primarily language-independent semantic understanding and representation processing. The last few layers of the models translate the semantic representations back into syntax and lexical information for the target language. This interpretation aligns with a recent study (Wendler et al., 2024), which suggested that the last few layers of Llama2 models are responsible for the conversion process to a target language.

Input	-	[BOS]
Output	Intervention in English neurons	Some of the issues that we are gonna have here are: the NSA is investigating whether the program is leaking in to the public and the government is trying to stop it as of late as it is possible. In the meantime the NSA is going to run the Panama Papers to find out what the
	Intervention in German neurons	Vorträge unter der Überschrift 'War für Trojä und ihr jahrhundert' zu nutzen und abzuschließen.
	Intervention in French neurons	«Il serait dommage de réécrire l'histoire au lieu de donner à entendre qu'une personne est une personne vivant dans l'état dans lequel elle est présente», ajoute le Kentou. «La plupart des médias dans le monde ne donnent pas suffisamment de voix, et qu'un jour il n'y
	Intervention in Spanish neurons	Chile, Colombia, Paraguay, Uruguay, Bolivia, Chile, Ecuador, Perú, Uruguay, Colombia, Paraguay, Paraguay, Colombia
	Intervention in Chinese neurons	三是(一)有权与允诺的机关有权予以采纳。
	Intervention in Japanese neurons	ただいま(25日の遅れのため)この商品は、注文確認日の翌営業日に発送致します。

Figure 3: Model-generated text examples with unconditional text generation setting by XGLM-564M. Given a [BOS] token as input, the model generates outputs through a random sampling method.

Input	-	Translate an English sentence into a target language.\n English: Machu Picchu consist of three main structures, namely Intihuatana, the Temple of the Sun, and the Room of the Three Windows.\n Target Language:
Output	Without any intervention	Machu Picchu consist of three main structures, namely Intihuatana, the Temple of the Sun, and the Room of the Three Windows.
	Intervention in German neurons	Machu Picchu besteht aus drei Hauptstrukturen, nämlich Intihuatana, der Tempel der Sonne und die Zimmer mit drei Fenstern.
	Intervention in French neurons	Machu Picchu est composé de trois structures principales, les Intihuatana, le Temple du Soleil et la Salle des Trois Fenêtres.
	Intervention in Spanish neurons	El Machu Picchu está compuesto por tres principales estructuras, como el Intihuatana, el Templo del Sol y el Salón de las Tres Ventanas.
	Intervention in Chinese neurons	秘魯的马騰岭有三个主要的建筑, 即祭坛、圣殿和三窗房。
	Intervention in Japanese neurons	マチュピチュは三つの主要構造物である、インティワタナ、太陽の神殿、および三つの窓の部屋である。

Figure 4: Model-generated text examples with conditional text generation setting by Llama-2-7b. Given a machine translation task as input, the model generates outputs through a greedy decoding method.

## 5.2 Controlling Language-specific Neurons

To show the effectiveness of the identified language-specific neurons, we investigated whether the models could control language in text generation by intervening with language-specific neurons. We conducted the investigation using unconditional and conditional (i.e., machine translation) text-generation settings.

### 5.2.1 Unconditional text generation

In the experiments on unconditional text generation, we do not provide models with any input prompt, i.e., only a [BOS] token as a prompt. Each model repeated text generation 100 times with random sampling decoding (temperature=0.8, top p=0.9) by changing the random seed from 1 to 100

each time the model started to generate text. Figure 3 illustrates model-generated text examples with the intervention setting in each class of language-specific neurons. It was shown that by changing language-specific neurons for intervention, we can change the language of the output texts.

To quantitatively measure the probability of occurrence of a target language, we classified the generated texts into language categories using the language identification classifier FastText (Joulin et al., 2017b,a). We classified each text into the target language if the classification score exceeded a threshold of 0.5 (Wenzek et al., 2020; Touvron et al., 2023) and calculated the probability of the target language occurrence, i.e., the evaluation metric was accuracy.

		before		after		
			Top	Bottom	Both	
XGLM (564M)	en	40.0	62.0	77.0	<b>89.0</b>	
	de	0.0	89.0	31.0	<b>95.0</b>	
	fr	0.0	86.0	7.0	<b>90.0</b>	
	es	2.0	71.0	5.0	<b>78.0</b>	
	zh	7.0	<b>82.0</b>	50.0	79.0	
	ja	7.0	92.0	61.0	<b>99.0</b>	
-	9.3	80.3	38.5	<b>88.3</b>		
BLOOM (1b7)	en	37.0	78.0	67.0	<b>88.0</b>	
	de	0.0	60.0	0.0	<b>86.0</b>	
	fr	13.0	80.0	72.0	<b>98.0</b>	
	es	18.0	44.0	94.0	<b>97.0</b>	
	zh	6.0	1.0	89.0	<b>90.0</b>	
	ja	0.0	67.0	35.0	<b>97.0</b>	
-	12.3	55.0	59.5	<b>92.7</b>		
Llama2 (7b)	en	83.0	82.0	<b>89.0</b>	<b>89.0</b>	
	de	0.0	2.0	6.0	<b>23.0</b>	
	fr	2.0	1.0	<b>8.0</b>	7.0	
	es	1.0	4.0	4.0	<b>35.0</b>	
	zh	0.0	2.0	4.0	<b>50.0</b>	
	ja	1.0	1.0	<b>12.0</b>	10.0	
-	14.5	15.3	20.5	<b>35.7</b>		

Table 4: Probability of language occurrence within the generated texts before and after intervention. Values in the "-" rows are the average values across six languages.

Table 4 summarizes the results. This demonstrates that intervention in language-specific neurons increases the probability of the target language occurrence in unconditional text generation for each language. In other words, the desired language could be generated by intentionally igniting target neurons. It should be noted that the BLOOM models achieve a high probability of German and Japanese text occurrence by intervention, although the models do not explicitly include German and Japanese in their pre-training datasets, as described in Table 2. It is possible for a small number of these languages to be unintentionally mixed, leading to unintentional ability acquisition. For example, an English text and its translation to language may be present in a single document (Briakou et al., 2023).

We conducted a study by intervening in only the top-1000 neurons, only the bottom-1000 neurons, and both groups of neurons. Interestingly, some languages did not respond to control by intervening only in the top-1000 or only the bottom-1000 neurons. This suggests that it is possible to effectively control language by intervening in both groups of neurons. In principle, the top- $k$  neurons are correlated with positive activation values. Conversely, the bottom- $k$  neurons were correlated with negative activation values. Figure 5 validates this hypothesis. These findings align with those of Wang

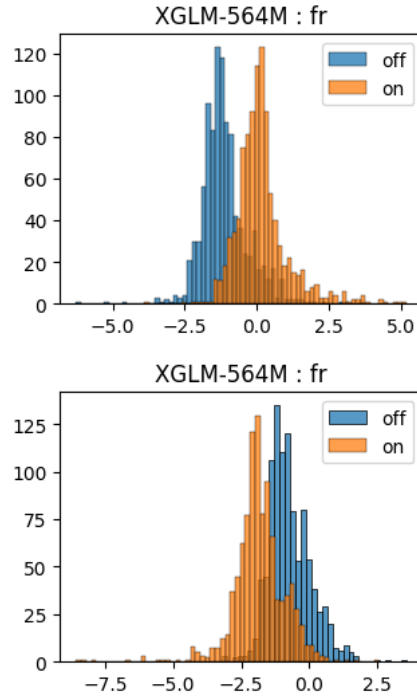


Figure 5: (Top) Distributional difference of activation values of the top-1000 neurons between target (on) and non-target languages (off). (Bottom) Distributional difference of activation value of the bottom-1000 neurons.

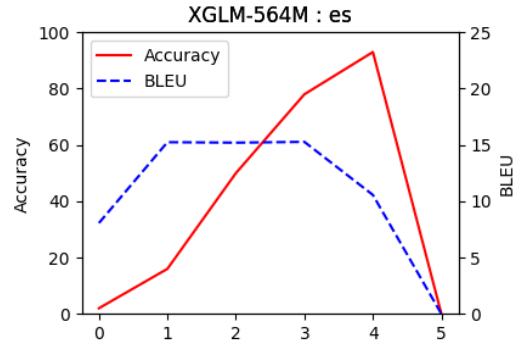


Figure 6: Ablation study of text generation by varying the number of neurons for intervention. x-axis:  $\log_{10}(k)$

et al. (2022), who suggested that neurons with both positive and negative correlations with labels are important for identifying target neurons.

We conducted an ablation study by changing the number of intervening neurons and analyzed its effect on the probability of target language occurrence. Additionally, we verified the quality of each model-generated text using the BLEU-4 score (BLEU). We evaluated BLEU only for texts identified by the language identifier as belonging to the target language. Specifically, for each model-generated text identified as the target language, we set the text as a hypothesis and all positive texts as references, and measured the BLEU score. We av-

		FLORES200		IWSLT2017		WMT	
		Accuracy	BLEU	Accuracy	BLEU	Accuracy	BLEU
XGLM-564M	de	0.0 → <b>38.0</b>	0.0 → 0.0	0.0 → <b>15.0</b>	0.0 → 0.0	0.0 → <b>17.0</b>	0.0 → 0.0
XGLM-564M	es	0.0 → <b>3.0</b>	0.0 → 0.0	→	→	→	→
XGLM-564M	ja	0.0 → 0.0	0.0 → 0.0	0.0 → 0.0	0.0 → 0.0	→	→
XGLM-564M	fr	0.0 → 0.0	0.0 → 0.0	0.0 → <b>3.0</b>	0.0 → 0.0	0.0 → <b>1.0</b>	0.0 → 0.0
XGLM-564M	zh	0.0 → <b>1.0</b>	0.0 → 0.0	0.0 → 2.0	0.0 → 0.0	0.0 → <b>2.0</b>	0.0 → 0.0
BLOOM-1b7	de	0.0 → <b>56.0</b>	1.3 → 1.3	0.0 → <b>35.0</b>	1.0 → <b>1.8</b>	0.0 → <b>37.0</b>	<b>2.9</b> → 1.7
BLOOM-1b7	es	0.0 → <b>2.0</b>	1.2 → 1.2	→	→	→	→
BLOOM-1b7	ja	0.0 → <b>6.0</b>	<b>0.2</b> → 0.1	0.0 → <b>8.0</b>	0.1 → <b>0.2</b>	→	→
BLOOM-1b7	fr	0.0 → <b>16.0</b>	1.7 → <b>2.8</b>	0.0 → <b>2.0</b>	1.0 → <b>1.5</b>	0.0 → <b>9.0</b>	1.7 → <b>2.7</b>
BLOOM-1b7	zh	0.0 → <b>21.0</b>	<b>0.3</b> → 0.2	0.0 → <b>3.0</b>	0.2 → <b>0.3</b>	0.0 → <b>34.0</b>	0.5 → <b>0.6</b>
Llama2-7b	de	0.0 → <b>66.0</b>	2.6 → <b>17.7</b>	0.0 → <b>48.0</b>	1.2 → <b>12.5</b>	2.0 → <b>53.0</b>	5.3 → <b>15.2</b>
Llama2-7b	es	4.0 → <b>77.0</b>	3.3 → <b>16.6</b>	→	→	→	→
Llama2-7b	ja	0.0 → <b>58.0</b>	0.3 → <b>10.4</b>	1.0 → <b>57.0</b>	0.2 → <b>4.5</b>	→	→
Llama2-7b	fr	1.0 → <b>58.0</b>	4.1 → <b>21.5</b>	0.0 → <b>32.0</b>	1.0 → <b>11.1</b>	0.0 → <b>36.0</b>	2.1 → <b>13.2</b>
Llama2-7b	zh	1.0 → <b>76.0</b>	1.0 → <b>11.5</b>	3.0 → <b>82.0</b>	0.6 → <b>7.8</b>	12.0 → <b>86.0</b>	2.4 → <b>11.3</b>
Llama2-13b	de	0.0 → <b>22.0</b>	1.5 → <b>8.8</b>	0.0 → <b>37.0</b>	0.6 → <b>10.0</b>	4.0 → <b>32.0</b>	3.3 → <b>9.7</b>
Llama2-13b	es	2.0 → <b>14.0</b>	1.8 → <b>4.3</b>	→	→	→	→
Llama2-13b	ja	7.0 → <b>54.0</b>	2.4 → <b>11.0</b>	4.0 → <b>75.0</b>	0.7 → <b>6.1</b>	→	→
Llama2-13b	fr	0.0 → <b>23.0</b>	1.6 → <b>10.5</b>	0.0 → <b>9.0</b>	0.7 → <b>4.7</b>	1.0 → <b>15.0</b>	2.2 → <b>6.6</b>
Llama2-13b	zh	20.0 → <b>93.0</b>	4.4 → <b>19.1</b>	40.0 → <b>96.0</b>	5.8 → <b>9.6</b>	57.0 → <b>99.0</b>	13.5 → <b>18.9</b>

Table 5: Results of conditional text generation. Values on the left side of the arrows (→) were measured without intervention on the language-specific neurons; values on the right side were measured during intervention on the neurons. FLORES200 includes translation tasks of English to the other five languages, while IWSLT2017 does not include tasks of English to Spanish, and WMT does not include tasks of English to Spanish or Japanese.

eraged the BLEU scores across all model-generated texts that were identified as the target language. Figure 6 shows the results of the intervention in Spanish neurons for XGLM-564M. This shows that increasing the number of intervening neurons up to 1000-10000 (3-4 on the logarithm of 10 in this figure) generally increases the probability of target language occurrence, but increasing beyond that degrades text quality. Eventually, the sentence collapses and both language identification and quality significantly decrease. This tendency exists regardless of language or model variation. See Appendix C.2 for the complete results.

### 5.2.2 Conditional text generation

In experiments of conditional text generation, models were given machine translation tasks and required to solve them in a zero-shot prompting setting, but with an unusual prompt format: “Translate an English sentence into a target language. English: {source text} Target Language:”. In other words, it is a prompt for a translation task that does not concretely specify the target language. The aim of the prompt was to accurately check whether the manipulation of language-specific neurons could lead to the output of the target language. Using this prompt as an input, the models started to generate text using a greedy decoding strategy. For this experiment, we randomly selected 100 machine translation sam-

ples from FLORES200, IWSLT2017 (Cettolo et al., 2017), and WMT (Bojar et al., 2014, 2016, 2018)<sup>3</sup>, respectively. Two evaluation metrics were used to measure translation quality: Accuracy of measuring the probability that the target language text is generated, and BLEU. In the unconditional text-generation setting, we measured the quality of only the generated texts in the target language. However, in the conditional text-generation setting, we calculated the BLEU score between each generated text and the corresponding ground-truth text by following the standard method of BLEU evaluation in machine translation tasks (Papineni et al., 2002).

Table 5 summarizes the experimental results of the conditional text generation. There were two main findings from these results. First, interventions in language-specific neurons tend to increase the probability of producing the target language (accuracy). Second, the translation quality (BLEU) of Llama2 models increased drastically along with accuracy. In contrast, the translation quality of XGLM and BLOOM did not significantly improve compared to the accuracy improvement. We investigated the reason for this by qualitatively analyzing the generated texts. XGLM and BLOOM were forced to output the target languages to some degree via intervention, but the output texts were not

<sup>3</sup>We used En → Fr tasks from WMT14, En → De tasks from WMT16, and En → Zh tasks from WMT18.

“Translate a sentence from English to a target language.”					“Translate an English sentence into a target language.”				
Accuracy		BLEU			Accuracy		BLEU		
de	0.0 → <b>62.0</b>	2.8 → <b>16.5</b>	0.0 → <b>66.0</b>	2.6 → <b>17.7</b>	0.0 → <b>66.0</b>	2.6 → <b>17.7</b>	0.0 → <b>66.0</b>	2.6 → <b>17.7</b>	0.0 → <b>66.0</b>
es	5.0 → <b>78.0</b>	4.0 → <b>16.5</b>	4.0 → <b>77.0</b>	3.3 → <b>16.6</b>	4.0 → <b>77.0</b>	3.3 → <b>16.6</b>	4.0 → <b>77.0</b>	3.3 → <b>16.6</b>	4.0 → <b>77.0</b>
ja	0.0 → <b>55.0</b>	0.3 → <b>9.2</b>	0.0 → <b>58.0</b>	0.3 → <b>10.4</b>	0.0 → <b>58.0</b>	0.3 → <b>10.4</b>	0.0 → <b>58.0</b>	0.3 → <b>10.4</b>	0.0 → <b>58.0</b>
fr	0.0 → <b>58.0</b>	3.4 → <b>21.3</b>	1.0 → <b>58.0</b>	4.1 → <b>21.5</b>	1.0 → <b>58.0</b>	4.1 → <b>21.5</b>	1.0 → <b>58.0</b>	4.1 → <b>21.5</b>	1.0 → <b>58.0</b>
zh	1.0 → <b>79.0</b>	1.2 → <b>12.7</b>	1.0 → <b>76.0</b>	1.0 → <b>11.5</b>	1.0 → <b>76.0</b>	1.0 → <b>11.5</b>	1.0 → <b>76.0</b>	1.0 → <b>11.5</b>	1.0 → <b>11.5</b>

“Translate an English sentence into German.”					“Translate an English sentence into Japanese.”				
Accuracy		BLEU			Accuracy		BLEU		
de	96.0 → <b>99.0</b>	<b>32.8</b> → 24.4	0.0 → <b>2.0</b>	0.3 → <b>1.2</b>	0.0 → <b>2.0</b>	0.3 → <b>1.2</b>	0.0 → <b>2.0</b>	0.3 → <b>1.2</b>	0.0 → <b>2.0</b>
es	0.0 → <b>1.0</b>	2.0 → <b>2.6</b>	0.0 → <b>2.0</b>	0.1 → <b>0.4</b>	0.0 → <b>2.0</b>	0.1 → <b>0.4</b>	0.0 → <b>2.0</b>	0.1 → <b>0.4</b>	0.0 → <b>2.0</b>
ja	<b>0.0</b> → <b>0.0</b>	0.3 → <b>0.4</b>	<b>100.0</b> → 99.0	<b>24.3</b> → 19.7	<b>100.0</b> → 99.0	<b>24.3</b> → 19.7	<b>100.0</b> → 99.0	<b>24.3</b> → 19.7	<b>100.0</b> → 99.0
fr	0.0 → <b>3.0</b>	2.6 → <b>3.1</b>	0.0 → <b>3.0</b>	0.2 → <b>1.0</b>	0.0 → <b>3.0</b>	0.2 → <b>1.0</b>	0.0 → <b>3.0</b>	0.2 → <b>1.0</b>	0.0 → <b>3.0</b>
zh	0.0 → <b>2.0</b>	<b>0.8</b> → 0.4	0.0 → <b>96.0</b>	1.3 → <b>14.9</b>	0.0 → <b>96.0</b>	1.3 → <b>14.9</b>	0.0 → <b>96.0</b>	1.3 → <b>14.9</b>	0.0 → <b>96.0</b>

Table 6: Results of conditional text generation with different prompt settings for Llama2-7b.

related to translation. For instance, when we intervened in German neurons, XGLM tended to output a word “Deutsch”. BLOOM tended to generate text unrelated to the translation or simply repeated the source text in English. Conversely, Llama2 tended to output translated text in the correct target language, resulting in improved accuracy and BLEU scores. This experiment showed that intervention in language-specific neurons can guide some models in the right direction, even when the models promptly receive an ambiguous translation task. Figure 4 shows examples of model-generated text for Llama2-7b model. See Section D.2 for additional examples.

We conducted several baseline experiments by changing the prompts to validate the robustness of the model outputs against prompts for the machine translation settings. Specifically, we tried the following four prompts: 1. “Translate a sentence from English into a target language. English: {source text} Target Language:”, 2. “Translate an English sentence into a target language. English: {source text} Target Language:”, 3. “Translate an English sentence into German. English: {source text} German:”, 4. “Translate an English sentence into Japanese. English: {source text} Japanese:”.

The first and second prompts are ambiguous because they do not explicitly specify the target language. The second prompt is the same as that in Table 5. Regarding the third and fourth prompts, we explicitly describe the target languages in the prompts: German from Western languages and Japanese from Eastern languages. Here, we focus on the Llama-2-7b model because it has signifi-

cantly improved both accuracy and BLEU scores, as described in Table 5. Similar to the experiment shown in Table 5, we conducted experiments in which the model was asked to solve a translation task under a specified prompt, while intervening in language-specific neurons.

The experimental results are presented in Table 6. The first and second prompts significantly increased the probability of target language occurrence and BLEU scores with intervention in language-specific neurons for all languages. In contrast, the third and fourth prompts caused few changes when we intervened in language-specific neurons for most languages. One possible reason is that explicitly specifying the target language in a prompt automatically fires specific neurons in that language, which may offset the effects of other language-specific neurons.

The only exception was the intervention in Chinese neurons under the fourth prompt “Translate an English sentence into Japanese.”, which increases the probability of Chinese text generation and BLEU scores. One possible reason is that some Japanese and Chinese neurons have similar firing patterns within the Llama-2-7b model. As shown in Table 7 in the Appendix, the Llama-2-7b model had a higher language-specific neuron overlap between Japanese and Chinese than the other pairs. As these two languages share many characters with similar surface morphologies, the high similarity between the two languages may have contributed to these results. However, it should be noted that this is not universally true across models; in some cases, the overlap of neurons in this language pair



is not always high in models other than Llama2, as described in Table 7 in the Appendix.

## 6 Conclusion

This study provides new insights into the activity of language-specific neurons in decoder-based multilingual pre-trained language models: the existence of neurons that fire uniquely for each language. The experimental results demonstrate that language-specific neurons mainly exist in the first and last few layers, regardless of the language, model size, and model variants. We further analyzed the effectiveness of the identified neurons by intervening in the neurons, that is, by replacing the output values with fixed activation values at inference with both unconditional and conditional settings. Using this approach, we can change the probability of the target language occurrence.

We hope that this study facilitates a deeper understanding of decoder-based PLMs and provides new insights for future research on multilingual decoder-based PLMs. Future research should include proposing language-specific model-compression methods. Future research also includes proposing new fine-tuning methods for downstream tasks to facilitate generalization to languages that are not included in the training dataset. For example, only fine-tuning the middle-layer parameters in decoder-based PLMs.

## 7 Limitation

This study only analyzes open models whose parameters were publicly available. It is not possible to analyze closed models with parameters that are not publicly available, such as ChatGPT or GPT-4 (OpenAI, 2023). Although we focused our analysis on six languages, other languages need to be examined in future studies. Analysis of encoder-decoder-based PLMs, such as mT5 (Xue et al., 2021), remains important but is beyond the scope of this study due to the fundamental differences in model architecture from decoder-only PLMs.

## Acknowledgements

We thank three anonymous reviewers for their helpful comments and feedback. This study was partially supported by PRESTO, JST Grant Number JPMJPR21C8, Japan.

## References

- Omer Antverg and Yonatan Belinkov. 2022. [On the pitfalls of analyzing individual neurons in language models](#). In *International Conference on Learning Representations*.
- Anthony Bau, Yonatan Belinkov, Hassan Sajjad, Nadir Durrani, Fahim Dalvi, and James Glass. 2019. [Identifying and controlling important neurons in neural machine translation](#). In *International Conference on Learning Representations*.
- Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. " O'Reilly Media, Inc."
- Ondrej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurelie Neveol, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, and Marcos Zampieri. 2016. [Findings of the 2016 conference on machine translation](#). In *Proceedings of the First Conference on Machine Translation*, pages 131–198, Berlin, Germany. Association for Computational Linguistics.
- Ondrej Bojar, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, and Christof Monz. 2018. [Findings of the 2018 conference on machine translation \(wmt18\)](#). In *Proceedings of the Third Conference on Machine Translation, Volume 2: Shared Task Papers*, pages 272–307, Belgium, Brussels. Association for Computational Linguistics.
- Ondrej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, Radu Soricut, Lucia Specia, and Aleks Tamchyna. 2014. [Findings of the 2014 workshop on statistical machine translation](#). In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 12–58, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Eleftheria Briakou, Colin Cherry, and George Foster. 2023. Searching for needles in a haystack: On the role of incidental bilingualism in PaLM’s translation capability. *arXiv preprint arXiv:2305.10266*.
- Mauro Cettolo, Marcello Federico, Luisa Bentivogli, Nihues Jan, Stüker Sebastian, Sudoh Katsutho, Yoshino Koichiro, and Federmann Christian. 2017. Overview of the iwslt 2017 evaluation campaign. In *Proceedings of the 14th International Workshop on Spoken Language Translation*, pages 2–14.
- Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. 2023. Journey to the center of the knowledge neurons: Discoveries of language-independent knowledge neurons and degenerate knowledge neurons. *arXiv preprint arXiv:2308.13198*.

- Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.
- Xavier Suau Cuadros, Luca Zappella, and Nicholas Apostoloff. 2022. Self-conditioning pre-trained language models. In *International Conference on Machine Learning*, pages 4455–4473. PMLR.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. [Knowledge neurons in pretrained transformers](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8493–8502, Dublin, Ireland. Association for Computational Linguistics.
- Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Herve Jegou, and Tomas Mikolov. 2017a. [Fasttext.zip: Compressing text classification models](#).
- Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017b. [Bag of tricks for efficient text classification](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431, Valencia, Spain. Association for Computational Linguistics.
- Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, et al. 2021. Few-shot learning with multilingual language models. *arXiv preprint arXiv:2112.10668*.
- Jesse Mu and Jacob Andreas. 2020. Compositional explanations of neurons. *Advances in Neural Information Processing Systems*, 33:17153–17163.
- Aaron Mueller, Yu Xia, and Tal Linzen. 2022. Causal analysis of syntactic agreement neurons in multilingual language models. *arXiv preprint arXiv:2210.14328*.
- Benjamin Muller, Yanai Elazar, Benoît Sagot, and Djamé Seddah. 2021. [First align, then predict: Understanding the cross-lingual ability of multilingual BERT](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2214–2231, Online. Association for Computational Linguistics.
- OpenAI. 2023. GPT-4 technical report. *arXiv*, pages 2303–08774.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. [How multilingual is multilingual BERT?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4996–5001, Florence, Italy. Association for Computational Linguistics.
- Hassan Sajjad, Nadir Durrani, and Fahim Dalvi. 2022. Neuron-level interpretation of deep nlp models: A survey. *Transactions of the Association for Computational Linguistics*, 10:1285–1303.
- Teven Le Scao, Angela Fan, Christopher Akiki, Elie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.
- Karolina Stańczak, Lucas Torroba Hennigen, Adina Williams, Ryan Cotterell, and Isabelle Augenstein. 2023. A latent-variable model for intrinsic probing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 13591–13599.
- Karolina Stańczak, Edoardo Ponti, Lucas Torroba Hennigen, Ryan Cotterell, and Isabelle Augenstein. 2022. Same neurons, different languages: Probing morphosyntax in multilingual pre-trained models. *arXiv preprint arXiv:2205.02023*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Andrea Gregor de Varda and Marco Marelli. 2023. Data-driven cross-lingual syntax: An agreement study with massively multilingual models. *Computational Linguistics*, 49(2):261–299.
- Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. 2022. [Finding skill neurons in pre-trained transformer-based language models](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11132–11152, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. 2024. Do Llamas work in English? On the latent language of multilingual transformers. *arXiv preprint arXiv:2402.10588*.
- Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Edouard Grave. 2020. [CCNet: Extracting high quality monolingual datasets from web crawl data](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4003–4012, Marseille, France. European Language Resources Association.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. **mT5: A massively multilingual pre-trained text-to-text transformer**. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

Yinfei Yang, Yuan Zhang, Chris Tar, and Jason Baldridge. 2019. **PAWS-X: A cross-lingual adversarial dataset for paraphrase identification**. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3687–3692, Hong Kong, China. Association for Computational Linguistics.

## A Text Examples from Language-Specific Text Corpus

Figure 7 lists text examples from the language-specific corpus, which is a combination of FLORES-200 and PAWS-X.

## B Identification of Neurons

### B.1 Distribution of Language-Specific Neurons Across Layers

#### B.1.1 Histogram

Figure 8, 9, 10, 11, 12, 13, 14, and 15 describes the distribution of language-specific neurons across layers in each model using a histogram.

#### B.1.2 Estimation of Beta Distribution Parameters

To quantitatively analyze the shape of the distribution, we regarded it as a sampling subset from a beta distribution and estimated its parameters of that beta distribution. The beta distribution has two parameters,  $\alpha$  and  $\beta$ . If  $\alpha < 1$  and  $\beta < 1$ , then the distribution becomes convex downward. If  $\alpha > 1$  and  $\beta > 1$ , then the distribution becomes convex upward. Table 7 lists the parameter-estimation results. It is clearly shown that the distribution of the top-1000 and bottom-1000 neurons generally has the parameters  $\alpha < 1$  and  $\beta < 1$ , supporting the claim that language-specific neurons exist in the first and last few layers. In contrast, the estimated parameters of the middle N neurons are  $\alpha > 1$

and  $\beta > 1$ , indicating that language-independent neurons exist in the middle layers.

### B.2 Overlapping language-specific neurons between languages.

Table 8a - 8f describe cross-table check results to count the number of overlapping language-specific neurons between languages.

### B.3 Activation Values of Top and Bottom-1000 Neurons

Figure 16, 17, 18, 19, 20, and 21 are histograms of the activation values for the top and bottom-1000 neurons. It was found that the top-1000 neurons activate positive values when we inputted positive text (on). In contrast, the bottom-1000 neurons tended to activate negative values when negative texts were given (off).

## C Intervention in Neurons for Unconditional Text Generation

### C.1 Effect of Intervention on Generated Language

Table 9 summarizes the probability of the target language occurrence in the generated texts before and after the intervention.

### C.2 Ablation Study of Changing the Number of Neurons for Intervention

Figures 22, 23, and 24 show the results of the ablation study when changing the number of neurons interventions.

### C.3 Model-Generated Text Examples For Unconditional Setting

- Figure 25 describes a summary of model-generated unconditional text examples.
- Figure 26 lists model-generated text examples without any interventions.
- Figure 27, 28, 29 list model-generated text examples with top-1000 and bottom-1000 neurons intervention.
- Figure 30, 31, 32, 33, 34, and 35 list model-generated text examples by changing the number of neuron interventions.

### C.4 Detail Setting of Unconditional Text Generation

A random sampling decoding method was used for unconditional text generation using the following

settings in all experiments for unconditional text generation:

- temperature: 0.8
- top-p: 0.9
- maximum output length: 64
- prompt (input token for models):
  - XGLM: `</s>` is automatically set.
  - BLOOM: nothing is automatically set. We explicitly set `</s>`.
  - Llama2: `<s>` is automatically set.

## D Intervention in Neurons for Conditional Text Generation

### D.1 Effect of Intervention on Conditional Text Generation

Tables 10, 11, and 12 summarize the probability of target language occurrence in the generated texts before and after the intervention.

### D.2 Model-Generated Text Examples For Conditional Setting

- Figure 36 describes a summary of model-generated conditional text examples.
- Figure 37 lists model-generated text examples without any interventions.
- Figures 38, 39, and 40 list model-generated text examples with top-1000 and bottom-1000 neurons intervention.

### D.3 Detail Setting of Conditional Text Generation

A greedy decoding method was used for conditional text generation. For evaluation of machine translation, the first line of model-generated text (sentences before the first linebreak code “`¥n`”) is used as translated sentences to omit useless subsequent sentences.

The following settings were used in all experiments for conditional text generation:

- maximum output length: 128

## E Detail Setting of Datasets

To create a language-specific text corpus, we mixed the following two datasets: **dev** split of PAWS-X (Yang et al., 2019) and **test** split of FLORES-200 (Costa-jussà et al., 2022). To create translation tasks for conditional text generation, we randomly sampled tasks from the **devtest** split of FLORES200, **test** split of IWSLT2017 (Cettolo et al., 2017), and **test** split of WMT (Bojar et al., 2014, 2016, 2018). All datasets were downloaded from HuggingFace (Wolf et al., 2019).

## F Detail Setting of BLEU-4 metrics

We used NLTK library (Bird et al., 2009) to measure the BLEU scores for both unconditional and conditional text generation. Specifically, the `sentence_bleu` function was used with `method2 SmoothingFunction` option for unconditional text generation. `corpus_bleu` function was used with `method2 SmoothingFunction` option for conditional text generation. To enable the comparison of BLEU scores across models, we tokenized all texts using a multilingual tokenizer, XGLM, whose pre-training corpus includes a large proportion of texts in the six target languages (Lin et al., 2021).

## G License

### G.1 Model

- XGLM: MIT [link]
- BLOOM: bigscience-bloom-rail-1.0 [link]
- Llama2: Meta license [link]

### G.2 Dataset

- PAWS-X: No License (Free to use) [link]
- FLORES200: cc-by-sa-4.0 [link]
- IWSLT2017: cc-by-nc-nd-4.0 [link]
- WMT14: Unknown [link]
- WMT16: Unknown [link]
- WMT18: Unknown [link]

## H Total computation for Experiments

We executed the experiments mainly for running the inference (both identification and intervention of language-specific neurons) for each model using the following number of A100(40GB) GPUs

and approximate computing hours per run. We run GPUs 60 times per model (6 languages  $\times$  (1 for identification of language-specific neurons + 6 for unconditional text generation by changing the number of neurons to intervene + 3 for conditional text generation)) for the production run. The computational resource of AI Bridging Cloud Infrastructure (ABCI) provided by the National Institute of Advanced Industrial Science and Technology (AIST) was used for the experiments.

- XGLM 564M: 1GPU  $\times$  0.5hrs  $\times$  60 runs.
- XGLM 1.7B: 1GPU  $\times$  0.5hrs  $\times$  60 runs.
- XGLM 2.9B: 1GPU  $\times$  0.5hrs  $\times$  60 runs.
- BLOOM 560M: 1GPU  $\times$  0.5hrs  $\times$  60 runs.
- BLOOM 1.7B: 1GPU  $\times$  0.5hrs  $\times$  60 runs.
- BLOOM 3B: 1GPU  $\times$  0.5hrs  $\times$  60 runs.
- Llama2 7B: 8GPUs  $\times$  1hrs  $\times$  60 runs.
- Llama2 13B: 8GPUs  $\times$  3hrs  $\times$  60 runs.

##### en  
 Cell comes from the Latin word cella which means small room.  
 The chimpanzee's hand and foot are similar in size and length, reflecting the hand's use for bearing weight in knuckle walking.  
 The next town hall was built in 1627 in the Baroque style and damaged in 1650 , 1653 , 1735 and 1779 .  
 With roots in Barcelona's Running Tours Barcelona and Copenhagen's Running Copenhagen, it was quickly joined by Running Tours Prague based in Prague and others.  
 The attack put a huge strain on relations between India and Pakistan.

##### de  
 Der Spieler des Tages ist heute Alex Ovechkin von den Washington Capitals.  
 Die Zeremonie wurde 1992 von Vince Neil moderiert, der En Vogue, Ugly Kid Joe, Arrested Development und Dennis Miller einbezog.  
 Die Biographie wurde zuletzt in Großbritannien, und zuvor in den USA (St. Martin 2013), Ungarn (Swiat Ksiazki 2013), Polen und China veröffentlicht.  
 Anders als bei größeren Fahrzeugen sind Sie wahrscheinlich bereits mit dem Fahren Ihres Autos vertraut und kennen dessen Grenzen.  
 Bei dieser Denkweise liegt der Fokus auf Schnelligkeit, Logik und Genauigkeit, auch auf der Identifizierung von Fakten, der erneuten Anwendung bestehender Techniken, der Sammlung von Informationen.

##### fr  
 Un même vol peut parfois présenter de grandes disparités tarifaires entre les agrégateurs ; avant de réserver, il est donc avantageux de comparer les résultats des recherches ainsi que de consulter le site web de la compagnie aérienne elle-même.  
 Sa deuxième griffe était plus grande, ce qui a engendré le nom Hesperonychus, qui signifie « griffe occidentale ».  
 La superficie de la Turquie, lacs compris, est de 783 562 kilomètres carrés, dont 755 688 kilomètres carrés en Asie du Sud-Ouest et 23 764 kilomètres carrés en Europe.  
 Un A Khap est un clan ou un groupe de clans apparentés, principalement sous les jats de l'est de l'Uttar Pradesh et de l'ouest du Haryana.  
 Seules quelques compagnies aériennes proposent encore des tarifs de deuil, qui réduisent légèrement le coût des voyages funéraires de dernière minute.

##### es  
 Colaboró con arquitectos contemporáneos como Giovan Giacomo Di Conforto, Bartolomeo Picchiatti y Francesco Grimaldi.  
 La isla de Bowen es un conocido viaje para hacer en un día o en un fin de semana y que incluye kayak, caminatas, tiendas comerciales, restaurantes y entre otros.  
 Fueron desarrollados por CyberConnect2 y publicados por Namco Bandai a partir de "Naruto : Ultimate Ninja" en 2005.  
 El Ministro de Salud manifestó su preocupación por el bienestar de quienes sacan partido de la legalidad temporal de las sustancias implicada y por las condenas por drogas que se han dictado desde la entrada en vigencia de los ahora inconstitucionales cambios.  
 Weitz está casada con Sebastian Weitz, que es cubano mexicano, y con quien tiene un hijo, Mercedes Martinez y una hija, Athena Weitz.

##### zh  
 乔治城湖也是乔治城和附近朗德罗克市的饮用水源。  
 “我们都非常震惊。”这位母亲表示。  
 由于表面张力的作用，钢针可以浮在水面上。  
 苏利耶·德·莫兰特曾在法国驻华使团工作多年，并在中国多个城市担任法国领事。  
 如果你想帮助别人时受伤，你只会使事情变得更糟。

##### ja  
 ルノ号は120～160立米の燃料を積んでいましたが、強風と高波によって船体が防波堤に押し込まれてしまいました。  
 グループは8月23日にオルタナティブの楽曲4曲と、「Fold YOur Hands Child」のアコースティック版を取録した「Itunes Session -- EP」をリリースしました。  
 月面探査機は、3つの重要な科学機器を搭載しただけでなく、本体の各面にインドの国旗が描かれています。  
 これは、家族版から2年以上後、また第3号名人版から3年以上後に放送された第1シーズンでした。  
 『Israel』はアメリカのジャズ・トロンボーン奏者カイ・ウィンディングとJ. J. ジョンソンが制作したアルバムで、1968年にレコーディングされた演奏をフィーチャーしており、CTIレーベルからリリースされた。

Figure 7: Text examples from language-specific text corpus, which is a mixture of FLORES-200 and PAWS-X.

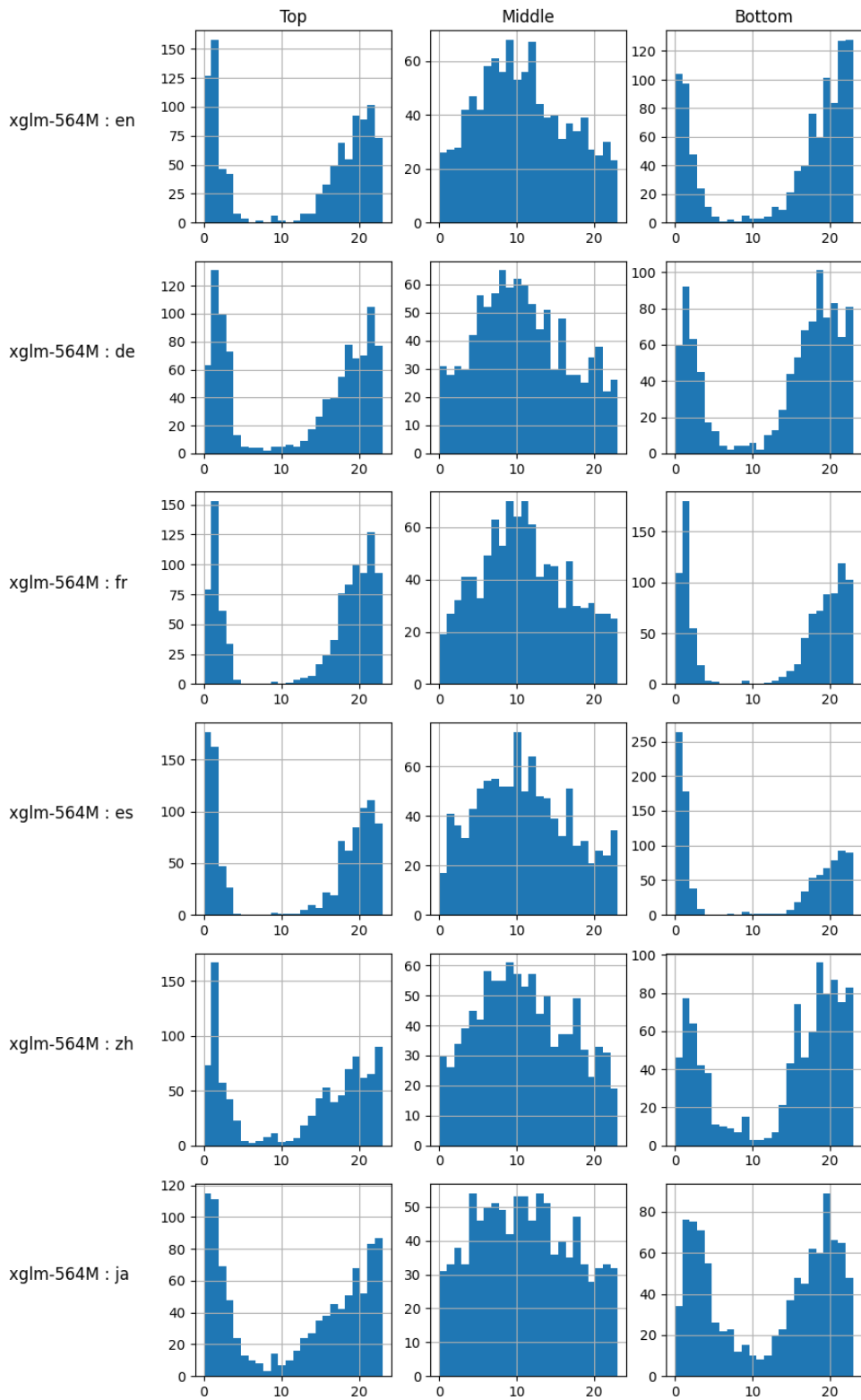


Figure 8: Histogram of language neurons across layers in xglm-564M.

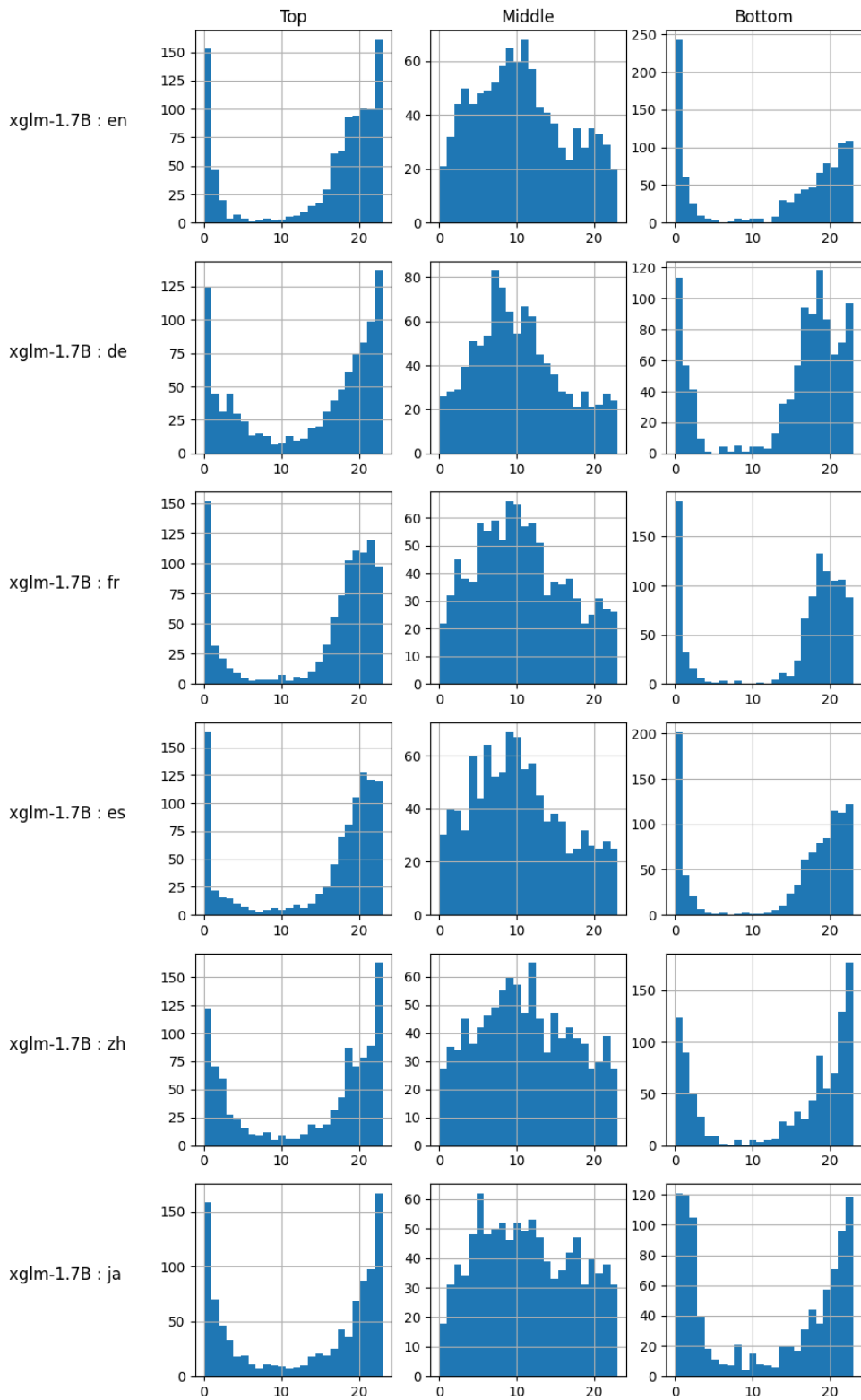


Figure 9: Histogram of language neurons across layers in xglm-1.7B.



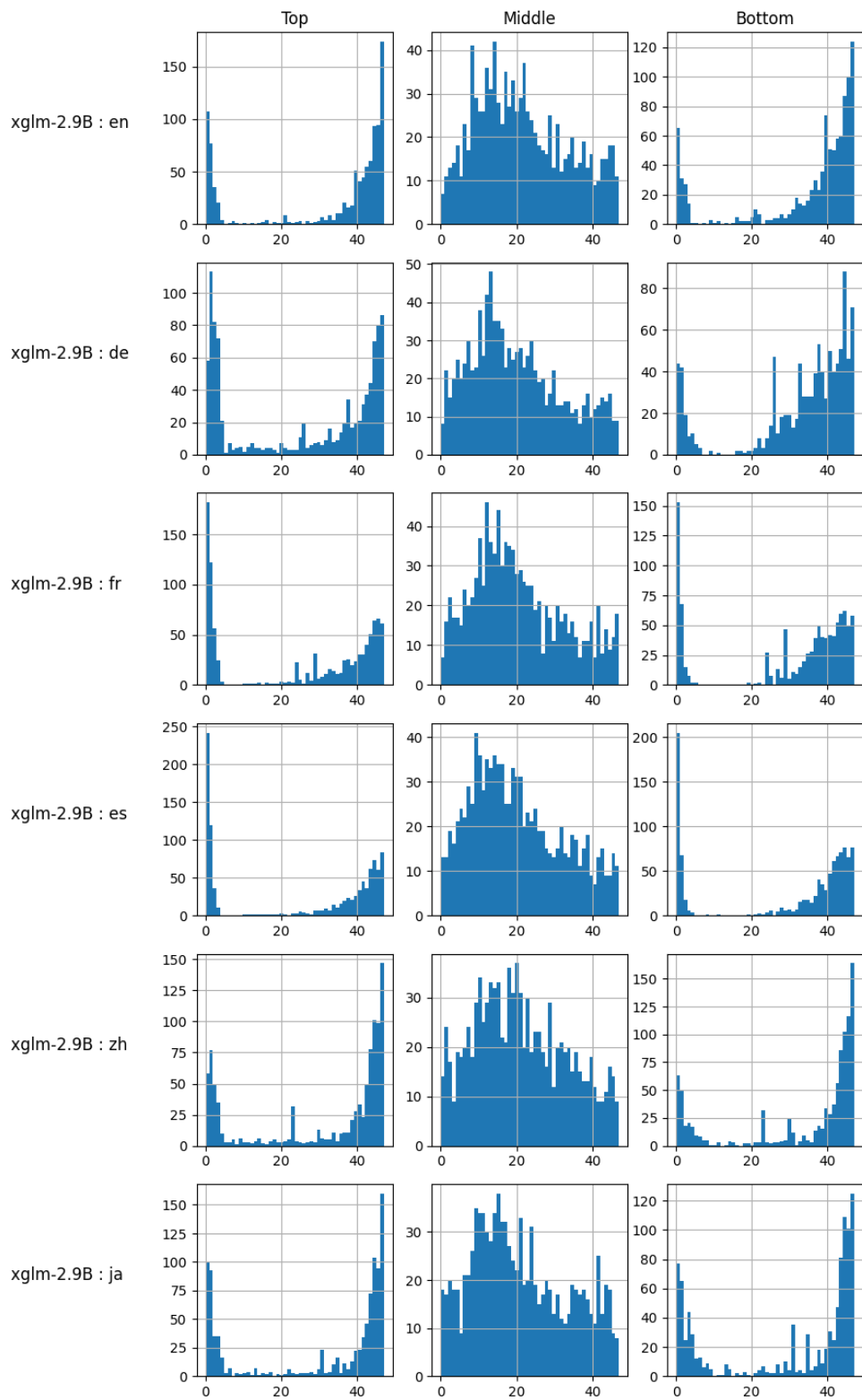


Figure 10: Histogram of language neurons across layers in xglm-2.9B.

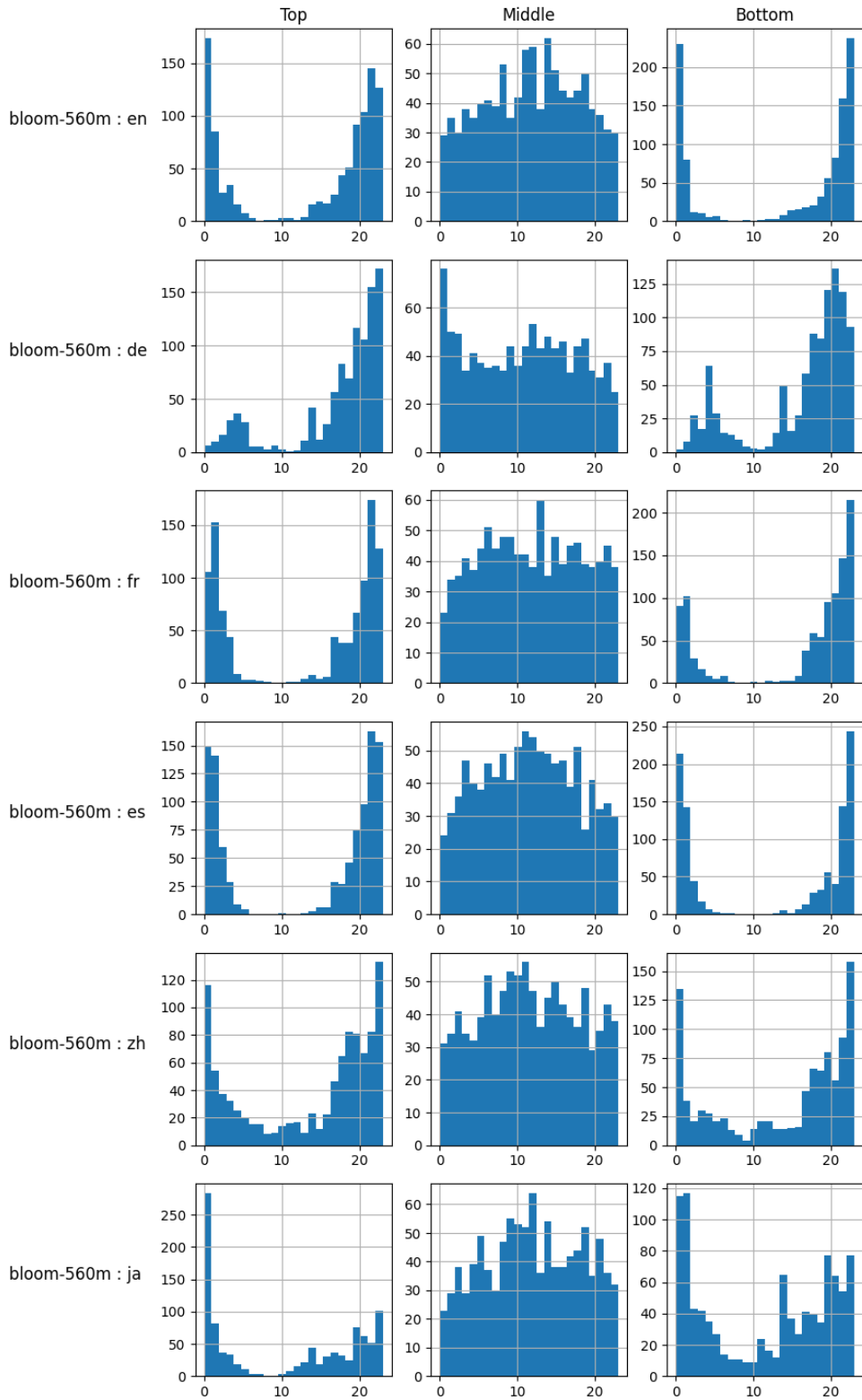


Figure 11: Histogram of language neurons across layers in bloom-560m.

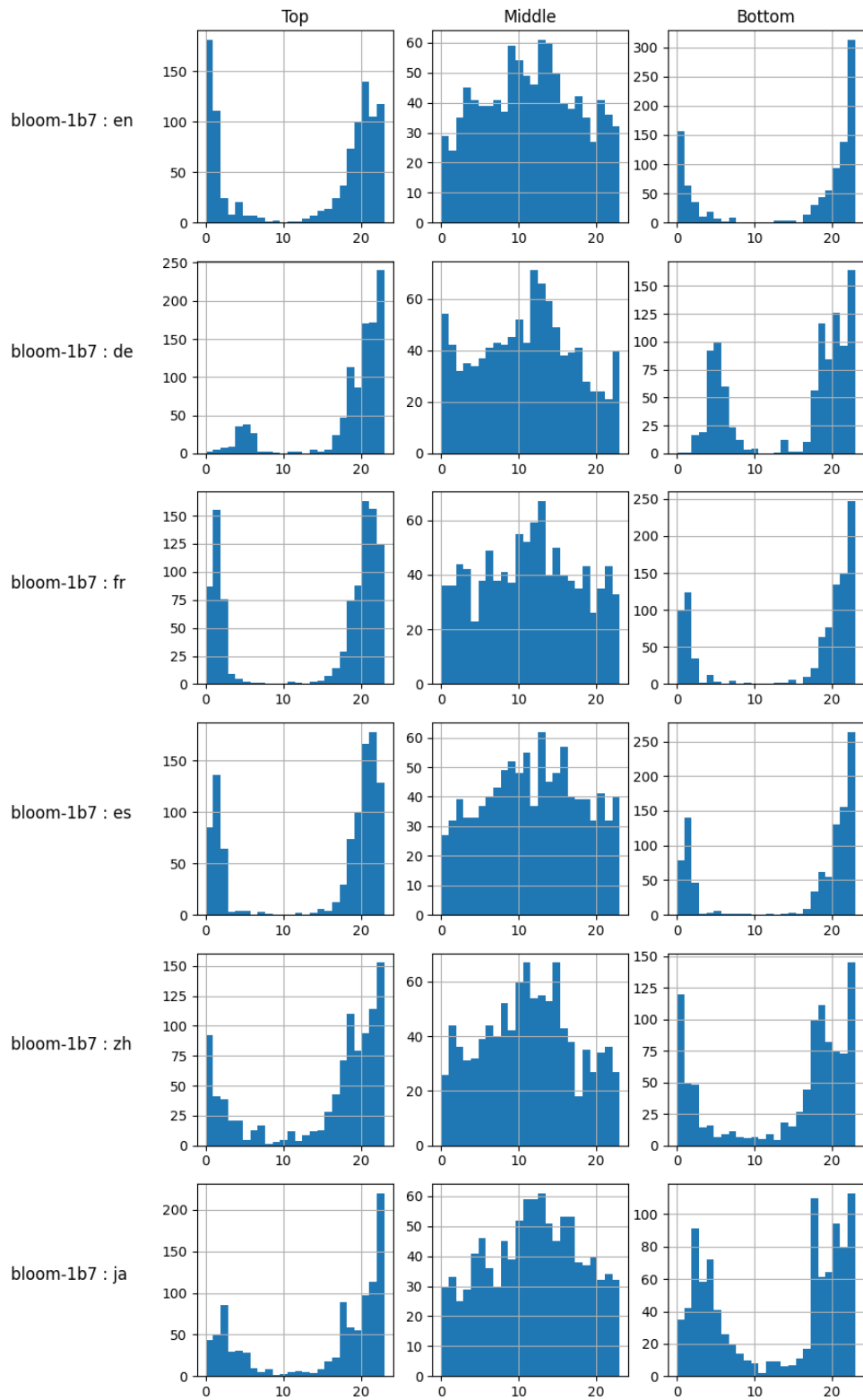


Figure 12: Histogram of language neurons across layers in bloom-1b7.

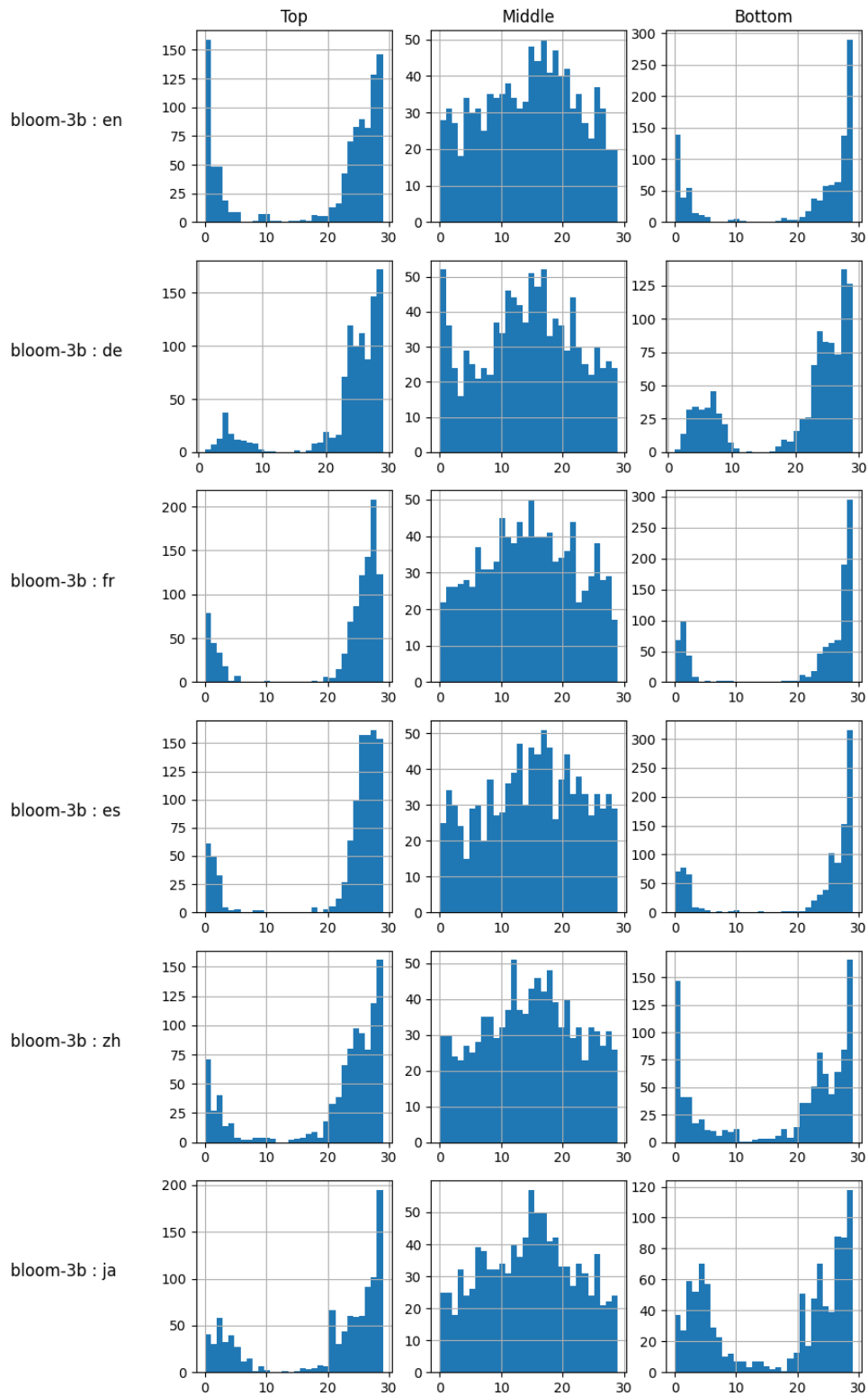


Figure 13: Histogram of language neurons across layers in bloom-3b.

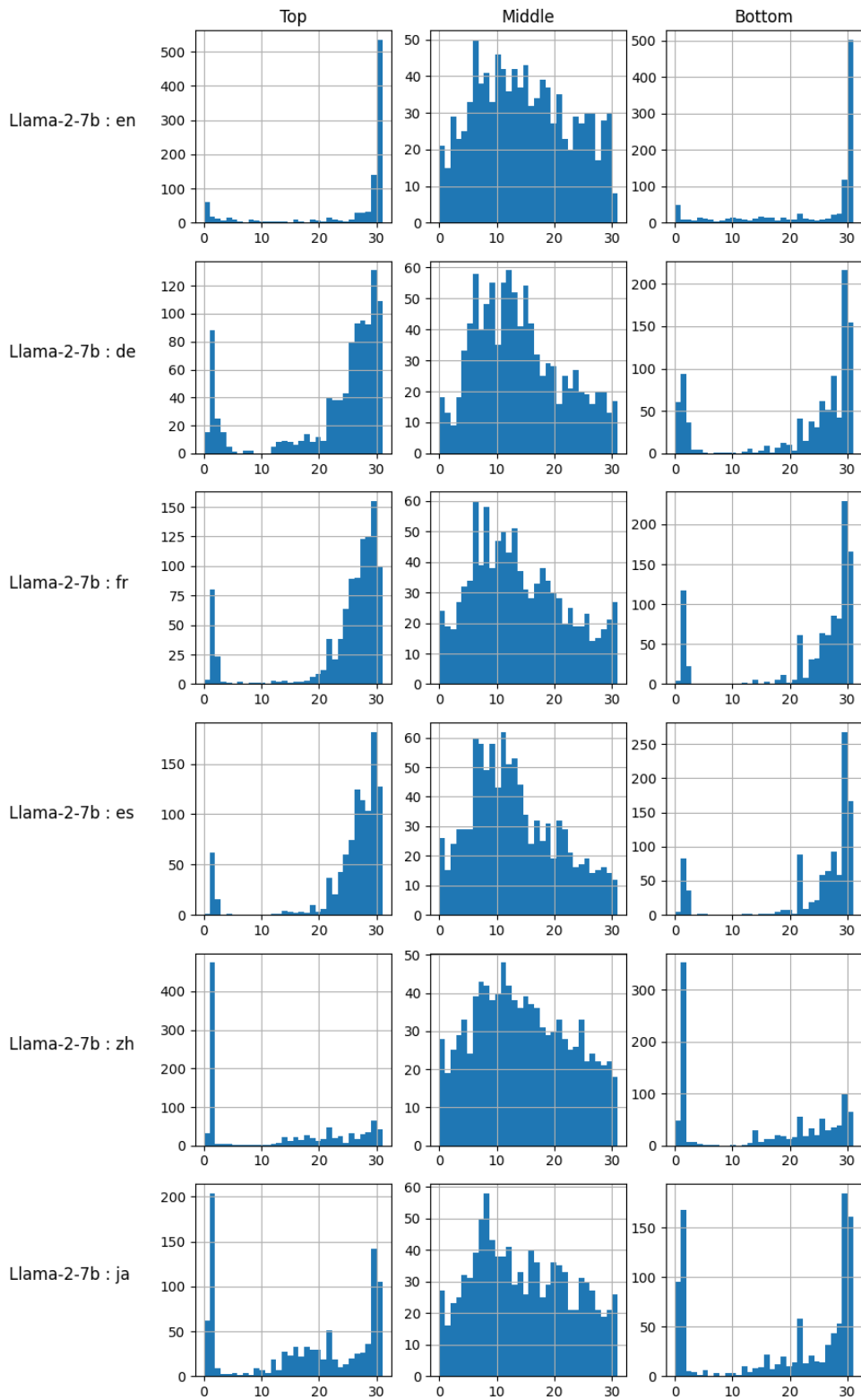


Figure 14: Histogram of language neurons across layers in Llama2-7b.

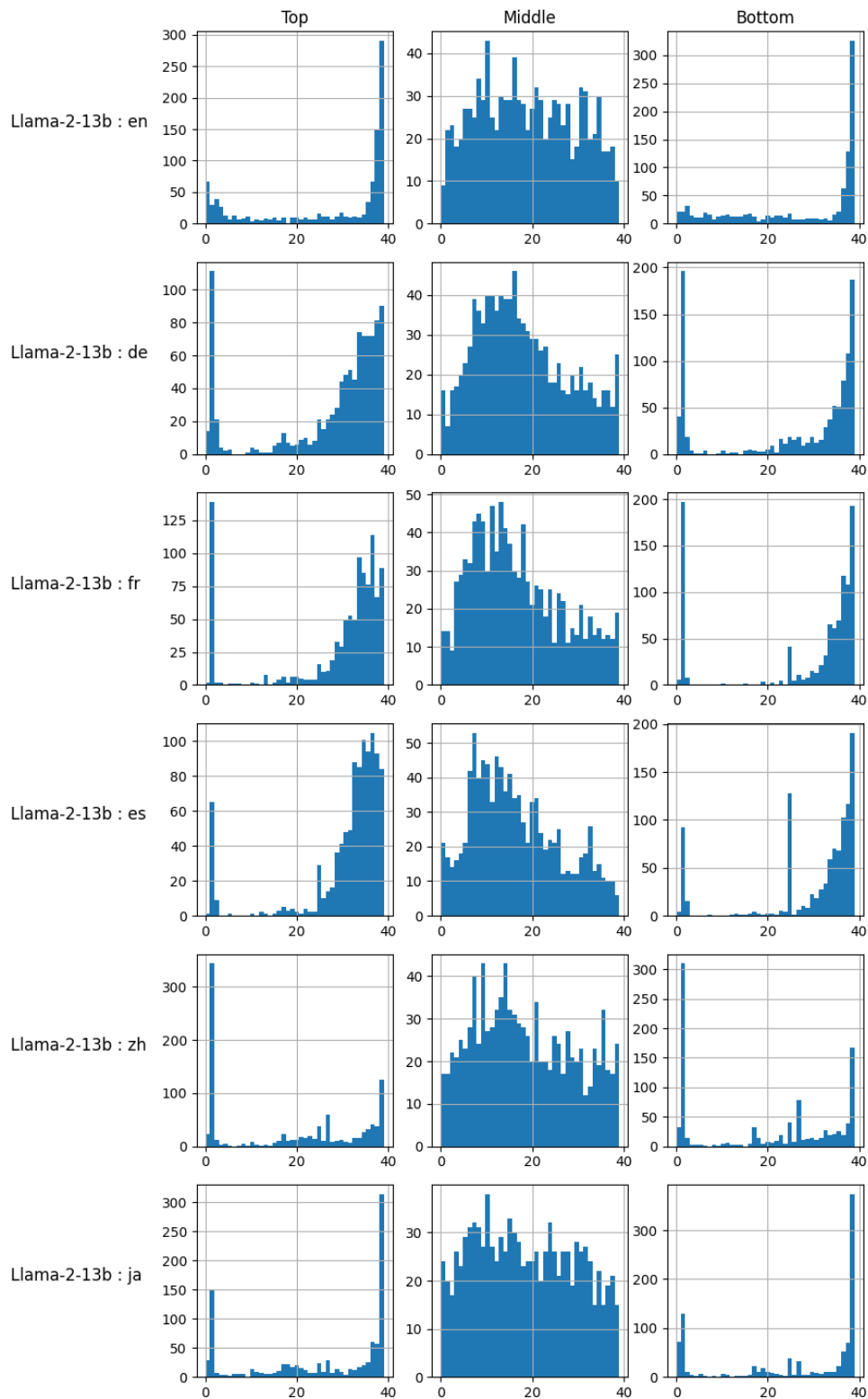


Figure 15: Histogram of language neurons across layers in Llama-2-13b.

Model	Language	Top	Middle	Bottom
xglm-564M	en	[0.58, 0.56]	[1.33, 1.45]	[0.68, 0.51]
xglm-564M	de	[0.66, 0.61]	[1.29, 1.41]	[0.81, 0.65]
xglm-564M	fr	[0.65, 0.53]	[1.4, 1.48]	[0.58, 0.52]
xglm-564M	es	[0.51, 0.52]	[1.31, 1.4]	[0.43, 0.53]
xglm-564M	zh	[0.65, 0.63]	[1.3, 1.42]	[0.86, 0.66]
xglm-564M	ja	[0.62, 0.63]	[1.19, 1.25]	[0.85, 0.78]
xglm-1.7B	en	[0.68, 0.48]	[1.29, 1.45]	[0.52, 0.5]
xglm-1.7B	de	[0.67, 0.54]	[1.34, 1.57]	[0.8, 0.59]
xglm-1.7B	fr	[0.72, 0.52]	[1.29, 1.45]	[0.67, 0.51]
xglm-1.7B	es	[0.69, 0.49]	[1.22, 1.43]	[0.6, 0.48]
xglm-1.7B	zh	[0.63, 0.5]	[1.24, 1.3]	[0.61, 0.47]
xglm-1.7B	ja	[0.56, 0.49]	[1.26, 1.26]	[0.56, 0.56]
xglm-2.9B	en	[0.54, 0.36]	[1.28, 1.43]	[0.84, 0.43]
xglm-2.9B	de	[0.5, 0.47]	[1.19, 1.52]	[0.99, 0.56]
xglm-2.9B	fr	[0.42, 0.46]	[1.21, 1.47]	[0.56, 0.47]
xglm-2.9B	es	[0.38, 0.41]	[1.19, 1.5]	[0.46, 0.4]
xglm-2.9B	zh	[0.59, 0.39]	[1.21, 1.44]	[0.69, 0.39]
xglm-2.9B	ja	[0.5, 0.37]	[1.14, 1.36]	[0.56, 0.4]
bloom-560m	en	[0.56, 0.47]	[1.26, 1.22]	[0.48, 0.39]
bloom-560m	de	[1.74, 0.66]	[0.94, 1.08]	[1.7, 0.77]
bloom-560m	fr	[0.56, 0.48]	[1.2, 1.14]	[0.69, 0.44]
bloom-560m	es	[0.52, 0.45]	[1.27, 1.26]	[0.44, 0.4]
bloom-560m	zh	[0.69, 0.56]	[1.17, 1.15]	[0.67, 0.52]
bloom-560m	ja	[0.46, 0.57]	[1.28, 1.2]	[0.64, 0.68]
bloom-1b7	en	[0.54, 0.48]	[1.26, 1.24]	[0.56, 0.37]
bloom-1b7	de	[2.14, 0.62]	[1.1, 1.2]	[1.22, 0.64]
bloom-1b7	fr	[0.61, 0.47]	[1.16, 1.17]	[0.63, 0.4]
bloom-1b7	es	[0.65, 0.46]	[1.24, 1.18]	[0.64, 0.4]
bloom-1b7	zh	[0.8, 0.52]	[1.27, 1.31]	[0.71, 0.53]
bloom-1b7	ja	[0.8, 0.49]	[1.29, 1.23]	[0.8, 0.63]
bloom-3b	en	[0.56, 0.42]	[1.23, 1.23]	[0.56, 0.36]
bloom-3b	de	[2.12, 0.65]	[1.08, 1.15]	[1.35, 0.62]
bloom-3b	fr	[0.83, 0.43]	[1.27, 1.26]	[0.67, 0.36]
bloom-3b	es	[0.95, 0.44]	[1.2, 1.12]	[0.66, 0.36]
bloom-3b	zh	[0.89, 0.48]	[1.19, 1.17]	[0.58, 0.45]
bloom-3b	ja	[0.82, 0.47]	[1.28, 1.27]	[0.75, 0.57]
Llama-2-7b-hf	en	[0.91, 0.33]	[1.28, 1.38]	[0.96, 0.35]
Llama-2-7b-hf	de	[1.06, 0.53]	[1.4, 1.63]	[0.75, 0.43]
Llama-2-7b-hf	fr	[1.37, 0.56]	[1.2, 1.39]	[1.03, 0.45]
Llama-2-7b-hf	es	[1.76, 0.58]	[1.3, 1.69]	[1.17, 0.46]
Llama-2-7b-hf	zh	[0.48, 0.7]	[1.2, 1.33]	[0.51, 0.57]
Llama-2-7b-hf	ja	[0.6, 0.53]	[1.17, 1.27]	[0.56, 0.42]
Llama-2-13b-hf	en	[0.65, 0.36]	[1.25, 1.32]	[0.83, 0.39]
Llama-2-13b-hf	de	[0.94, 0.53]	[1.26, 1.41]	[0.57, 0.39]
Llama-2-13b-hf	fr	[0.98, 0.51]	[1.22, 1.51]	[0.68, 0.39]
Llama-2-13b-hf	es	[1.68, 0.61]	[1.29, 1.7]	[1.11, 0.46]
Llama-2-13b-hf	zh	[0.48, 0.52]	[1.12, 1.2]	[0.48, 0.48]
Llama-2-13b-hf	ja	[0.63, 0.38]	[1.11, 1.21]	[0.55, 0.33]

Table 7: Estimation of Beta distribution parameters from the histogram of neurons across layers in top-1000, middle-1000, and bottom-1000 groups, respectively.

(a) xglm-564M						
	de	en	es	fr	ja	zh
de	2000	41	74	39	44	34
en	41	2000	34	41	49	40
es	74	34	2000	57	77	22
fr	39	41	57	2000	21	93
ja	44	49	77	21	2000	27
zh	34	40	22	93	27	2000

(b) xglm-1.7B						
	de	en	es	fr	ja	zh
de	2000	12	14	9	43	9
en	12	2000	21	22	23	28
es	14	21	2000	60	22	17
fr	9	22	60	2000	7	30
ja	43	23	22	7	2000	30
zh	9	28	17	30	30	2000

(c) xglm-2.9B						
	de	en	es	fr	ja	zh
de	2000	10	6	1	14	5
en	10	2000	13	10	8	11
es	6	13	2000	28	12	16
fr	1	10	28	2000	7	12
ja	14	8	12	7	2000	30
zh	5	11	16	12	30	2000

(d) bloom-560m						
	de	en	es	fr	ja	zh
de	2000	12	19	20	12	61
en	12	2000	76	91	61	87
es	19	76	2000	168	70	47
fr	20	91	168	2000	42	56
ja	12	61	70	42	2000	5
zh	61	87	47	56	5	2000

(e) bloom-1b7						
	de	en	es	fr	ja	zh
de	2000	10	22	15	20	59
en	10	2000	55	88	26	59
es	22	55	2000	140	28	10
fr	15	88	140	2000	24	39
ja	20	26	28	24	2000	8
zh	59	59	10	39	8	2000

(f) bloom-3b						
	de	en	es	fr	ja	zh
de	2000	8	12	12	15	43
en	8	2000	64	45	34	46
es	12	64	2000	98	14	26
fr	12	45	98	2000	18	25
ja	15	34	14	18	2000	11
zh	43	46	26	25	11	2000

(g) Llama-2-7b						
	de	en	es	fr	ja	zh
de	2000	20	22	12	7	15
en	20	2000	16	14	11	11
es	22	16	2000	34	17	8
fr	12	14	34	2000	13	14
ja	7	11	17	13	2000	85
zh	15	11	8	14	85	2000

(h) Llama-2-13b						
	de	en	es	fr	ja	zh
de	2000	14	13	11	5	23
en	14	2000	10	10	18	9
es	13	10	2000	23	16	1
fr	11	10	23	2000	9	34
ja	5	18	16	9	2000	80
zh	23	9	1	34	80	2000

Table 8: Cross-table check to count the number of overlapping language-specific neurons between languages.



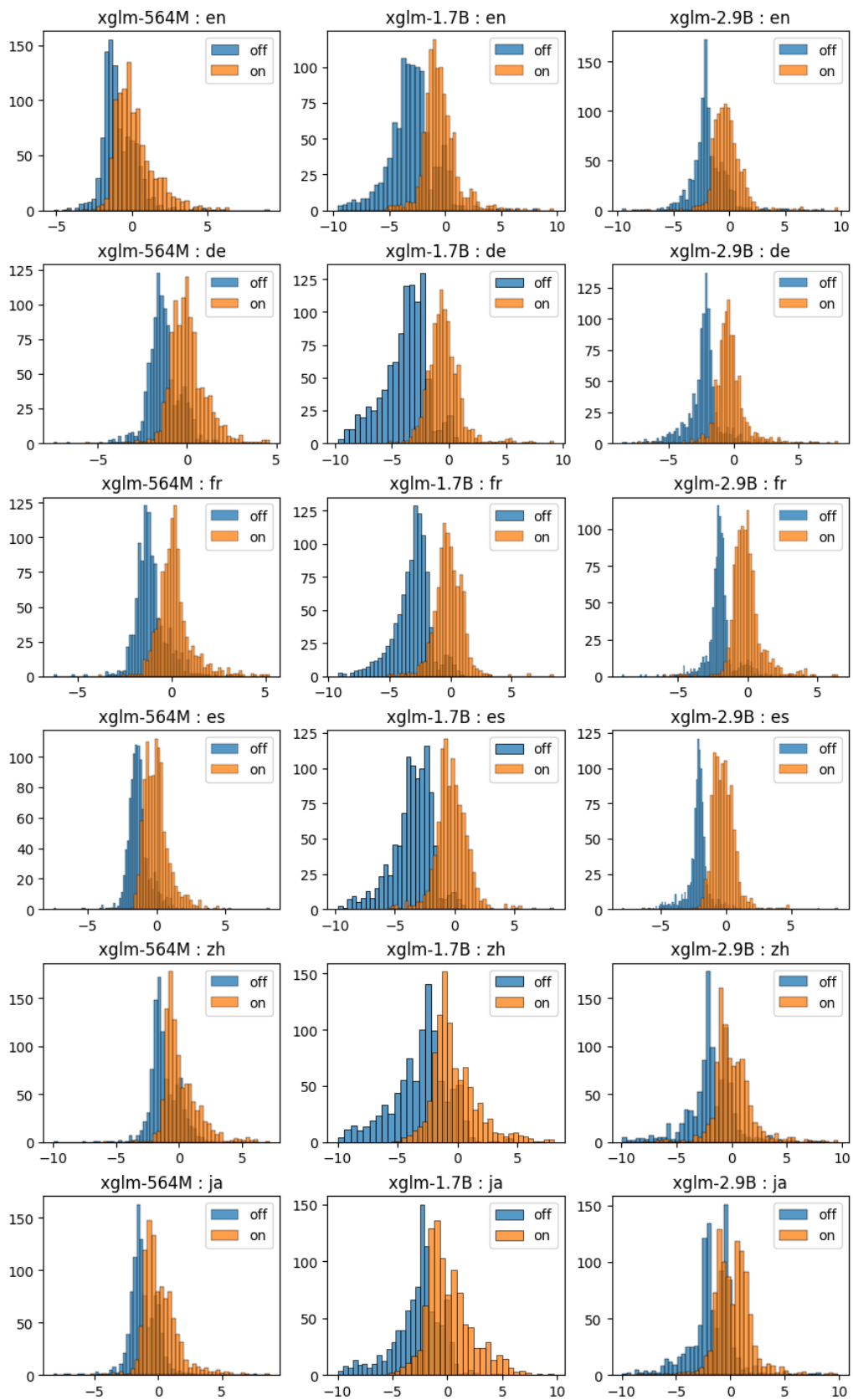


Figure 16: Activation value difference of top-1000 neurons between target language(on) and non-target languages(off). x-axis: activation value of neurons. y-axis: frequency.

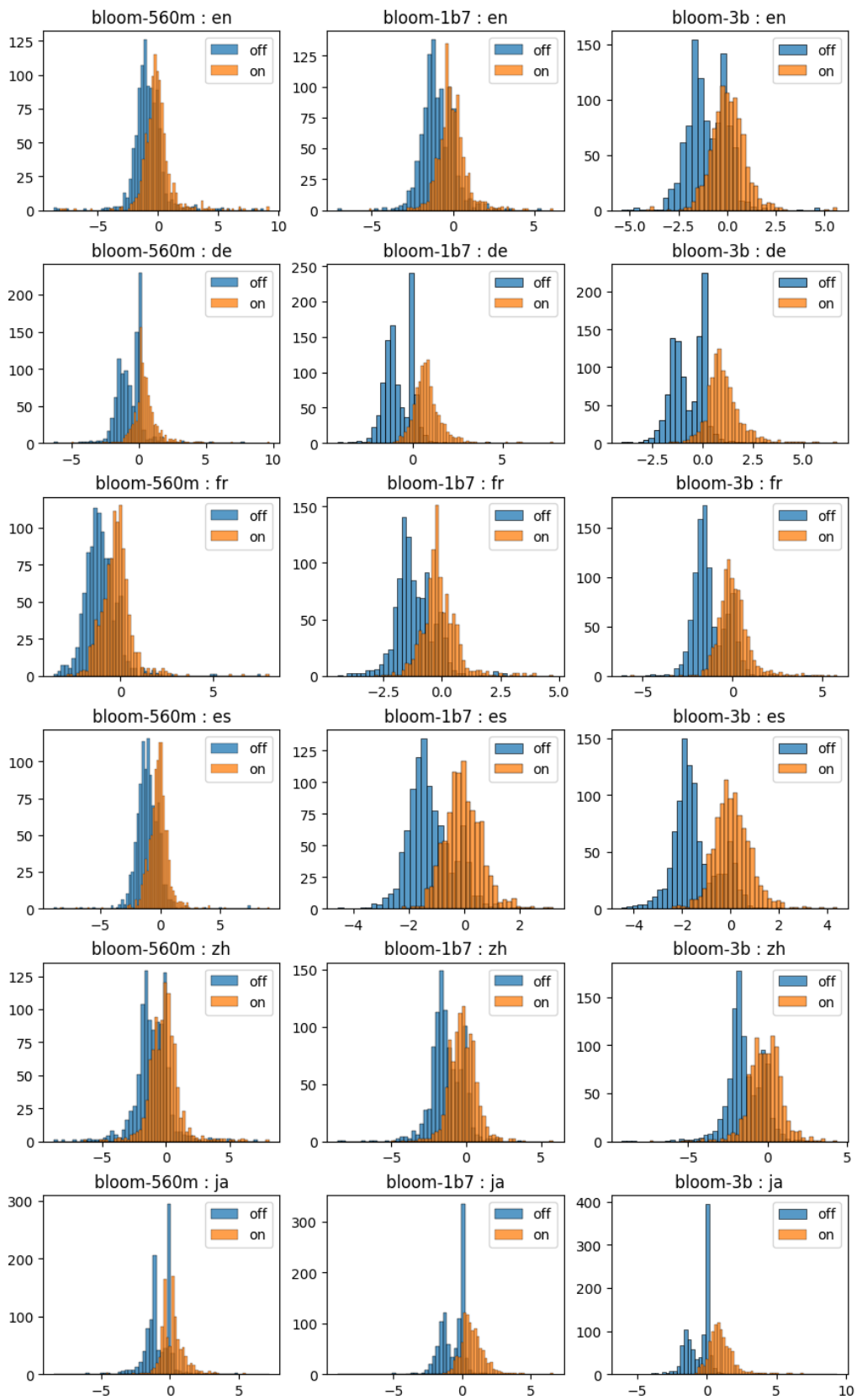


Figure 17: Activation value difference of top-1000 neurons between target language(on) and non-target languages(off). x-axis: activation value of neurons. y-axis: frequency.

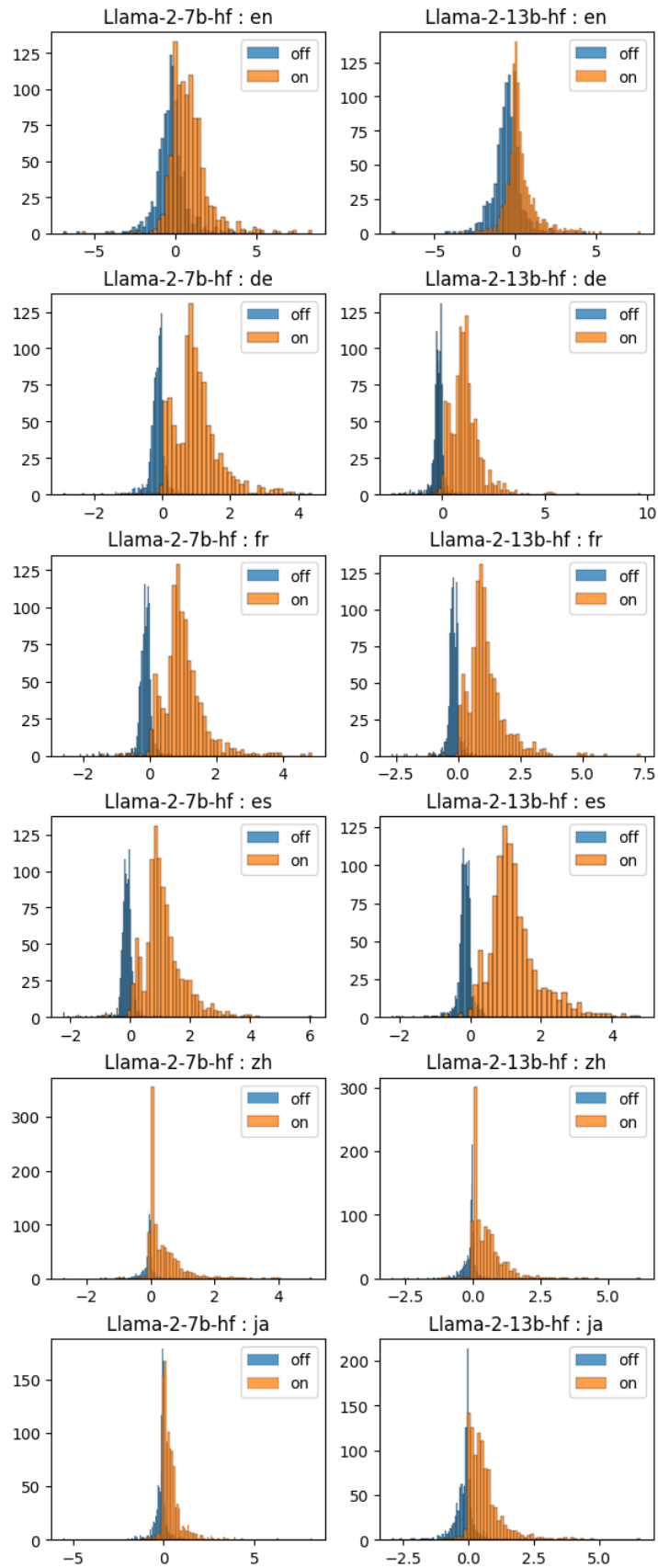


Figure 18: Activation value difference of top-1000 neurons between target language(on) and non-target languages(off). x-axis: activation value of neurons. y-axis: frequency.

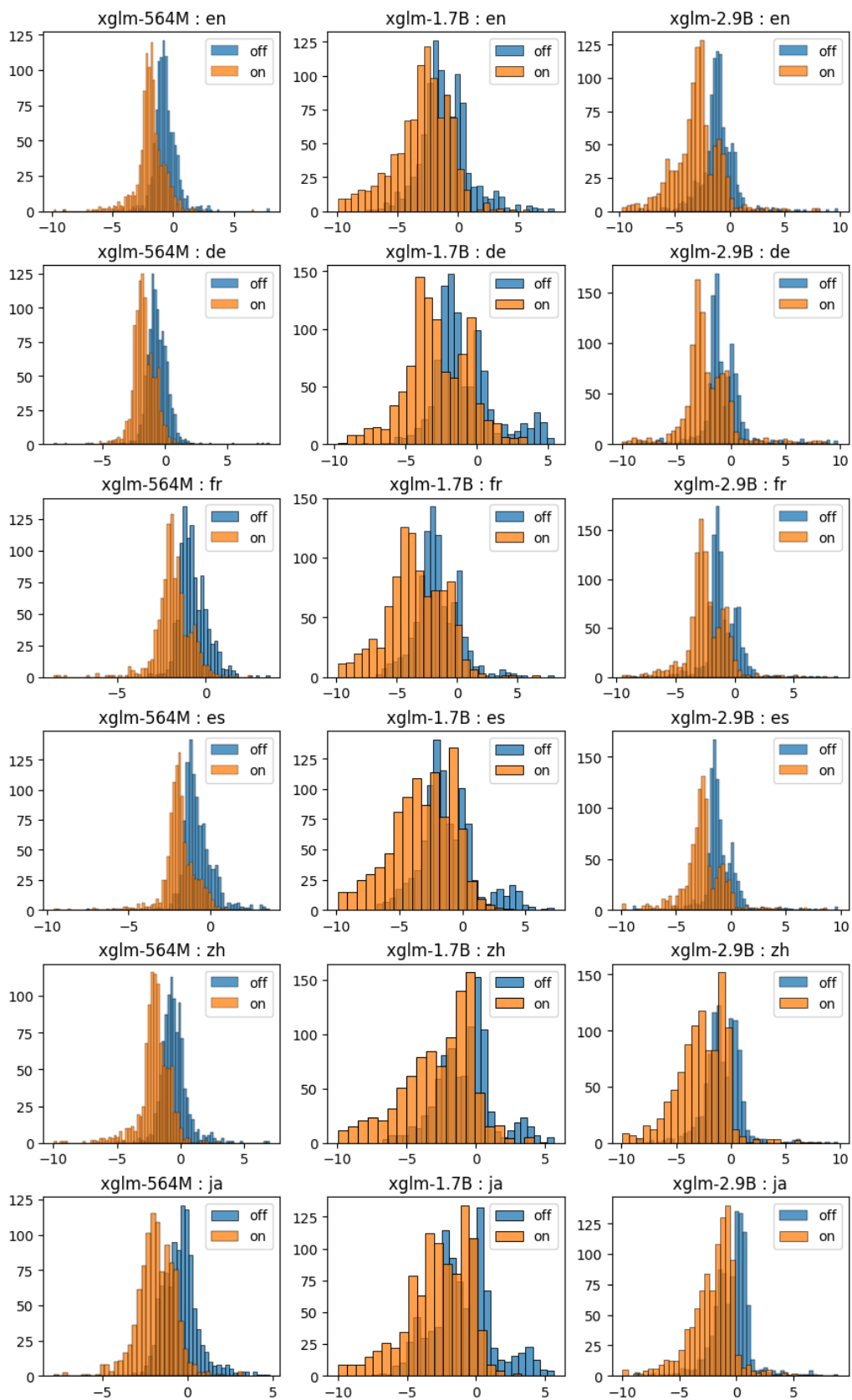


Figure 19: Activation value difference of bottom-1000 neurons between target language(on) and non-target languages(off). x-axis: activation value of neurons. y-axis: frequency.

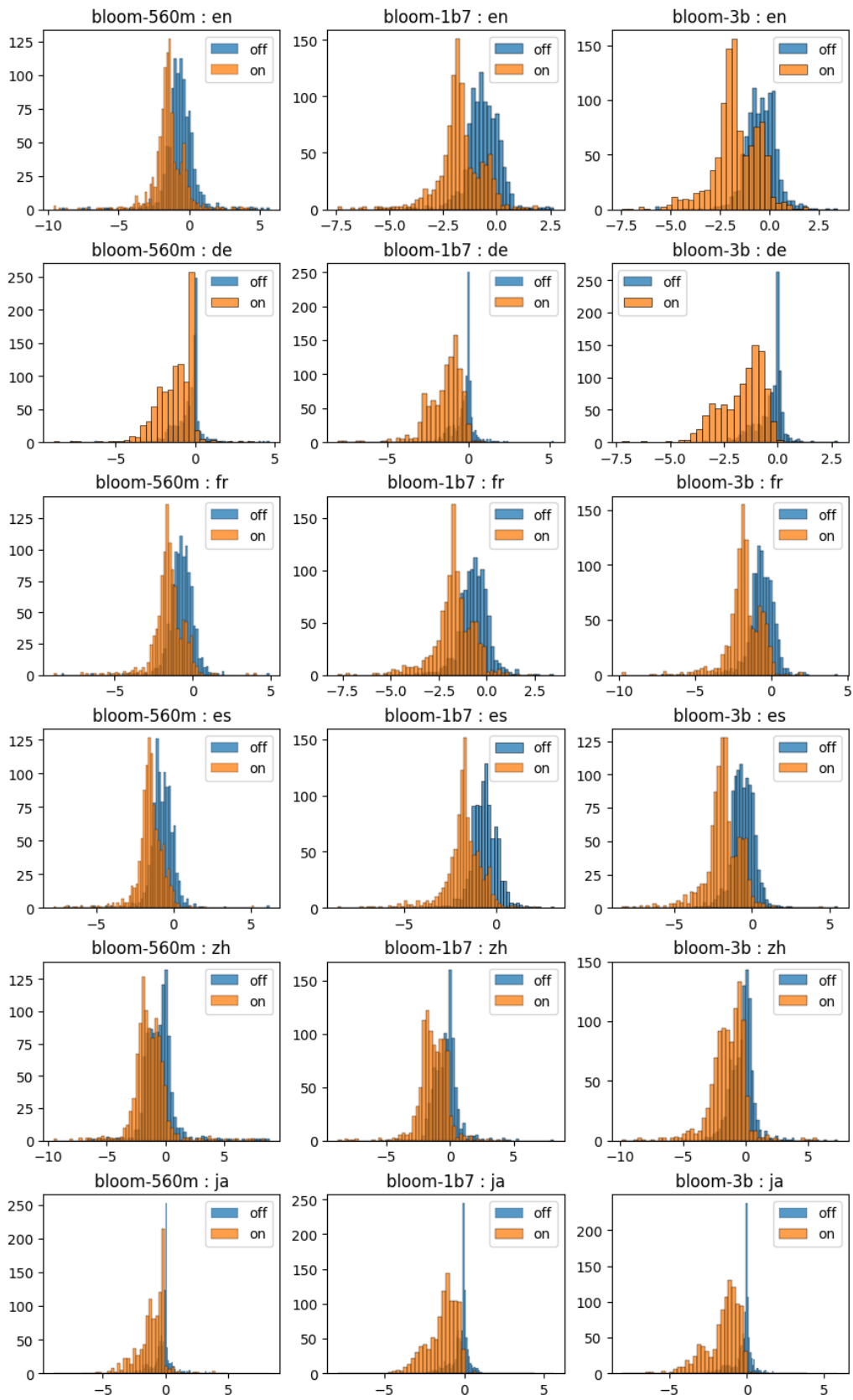


Figure 20: Activation value difference of bottom-1000 neurons between target language(on) and non-target languages(off). x-axis: activation value of neurons. y-axis: frequency.

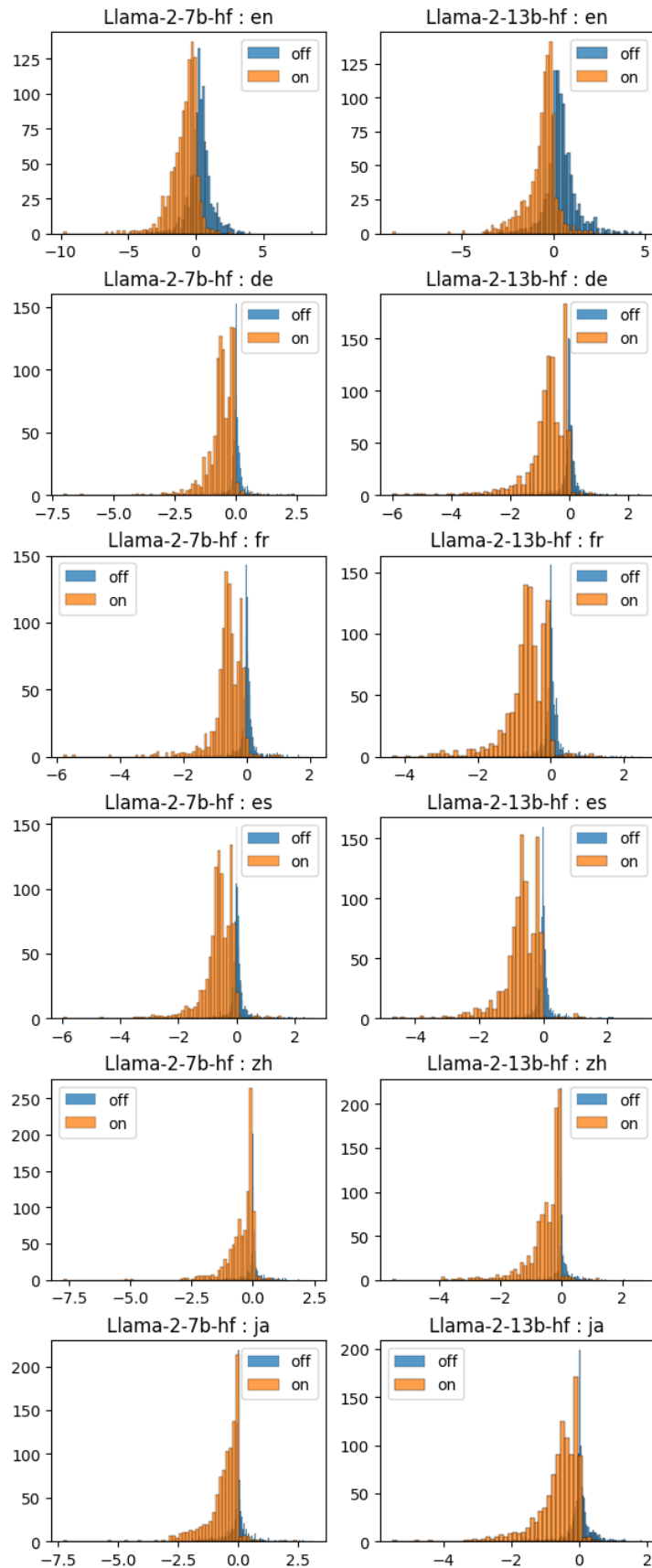


Figure 21: Activation value difference of bottom-1000 neurons between target language(on) and non-target languages(off). x-axis: activation value of neurons. y-axis: frequency.

		before		after				before		after	
		Top		Bottom	Both			Top		Bottom	Both
xglm (564M)	en	40.0	62.0	77.0	<b>89.0</b>	bloom (1b7)	en	37.0	78.0	67.0	<b>88.0</b>
	de	0.0	89.0	31.0	<b>95.0</b>		de	0.0	60.0	0.0	<b>86.0</b>
	fr	0.0	86.0	7.0	<b>90.0</b>		fr	13.0	80.0	72.0	<b>98.0</b>
	es	2.0	71.0	5.0	<b>78.0</b>		es	18.0	44.0	94.0	<b>97.0</b>
	zh	7.0	<b>82.0</b>	50.0	79.0		zh	6.0	1.0	89.0	<b>90.0</b>
	ja	7.0	92.0	61.0	<b>99.0</b>		ja	0.0	67.0	35.0	<b>97.0</b>
-	9.3	80.3	38.5	<b>88.3</b>	-	12.3	55.0	59.5	<b>92.7</b>		
xglm (1.7B)	en	36.0	23.0	<b>43.0</b>	32.0	bloom (3b)	en	32.0	41.0	87.0	<b>96.0</b>
	de	3.0	84.0	10.0	<b>91.0</b>		de	0.0	44.0	2.0	<b>55.0</b>
	fr	1.0	54.0	5.0	<b>70.0</b>		fr	15.0	72.0	72.0	<b>93.0</b>
	es	3.0	53.0	9.0	<b>69.0</b>		es	19.0	60.0	94.0	<b>95.0</b>
	zh	3.0	59.0	4.0	<b>65.0</b>		zh	7.0	24.0	<b>91.0</b>	90.0
	ja	9.0	83.0	17.0	<b>87.0</b>		ja	0.0	85.0	1.0	<b>87.0</b>
-	9.2	59.3	14.7	<b>69.0</b>	-	12.2	54.3	57.8	<b>86.0</b>		
xglm (2.9B)	en	31.0	28.0	<b>48.0</b>	42.0	Llama-2 (7b)	en	83.0	82.0	<b>89.0</b>	<b>89.0</b>
	de	2.0	<b>92.0</b>	1.0	88.0		de	0.0	2.0	6.0	<b>23.0</b>
	fr	1.0	60.0	3.0	<b>61.0</b>		fr	2.0	1.0	<b>8.0</b>	7.0
	es	1.0	67.0	5.0	<b>73.0</b>		es	1.0	4.0	4.0	<b>35.0</b>
	zh	5.0	74.0	6.0	<b>85.0</b>		zh	0.0	2.0	4.0	<b>50.0</b>
	ja	11.0	<b>81.0</b>	3.0	80.0		ja	1.0	1.0	<b>12.0</b>	10.0
-	8.5	67.0	11.0	<b>71.5</b>	-	14.5	15.3	20.5	<b>35.7</b>		
bloom (560m)	en	50.0	69.0	80.0	<b>85.0</b>	Llama-2 (13b)	en	64.0	90.0	81.0	<b>94.0</b>
	de	0.0	34.0	0.0	<b>72.0</b>		de	3.0	2.0	3.0	<b>16.0</b>
	fr	13.0	37.0	85.0	<b>93.0</b>		fr	0.0	<b>9.0</b>	1.0	8.0
	es	9.0	72.0	69.0	<b>97.0</b>		es	1.0	1.0	<b>5.0</b>	<b>5.0</b>
	zh	0.0	24.0	61.0	<b>90.0</b>		zh	3.0	<b>10.0</b>	6.0	5.0
	ja	0.0	60.0	0.0	<b>74.0</b>		ja	2.0	<b>6.0</b>	1.0	4.0
-	12.0	49.3	49.2	<b>85.2</b>	-	12.2	19.7	16.2	<b>22.0</b>		

Table 9: Probability of language occurrence of generated texts before and after intervention. Top: intervention to only top-1000 neurons. Bottom: intervention to only bottom-1000 neurons. Both: intervention to both top- and bottom-1000 neurons. The metric is accuracy measured by the FastText language identifier.

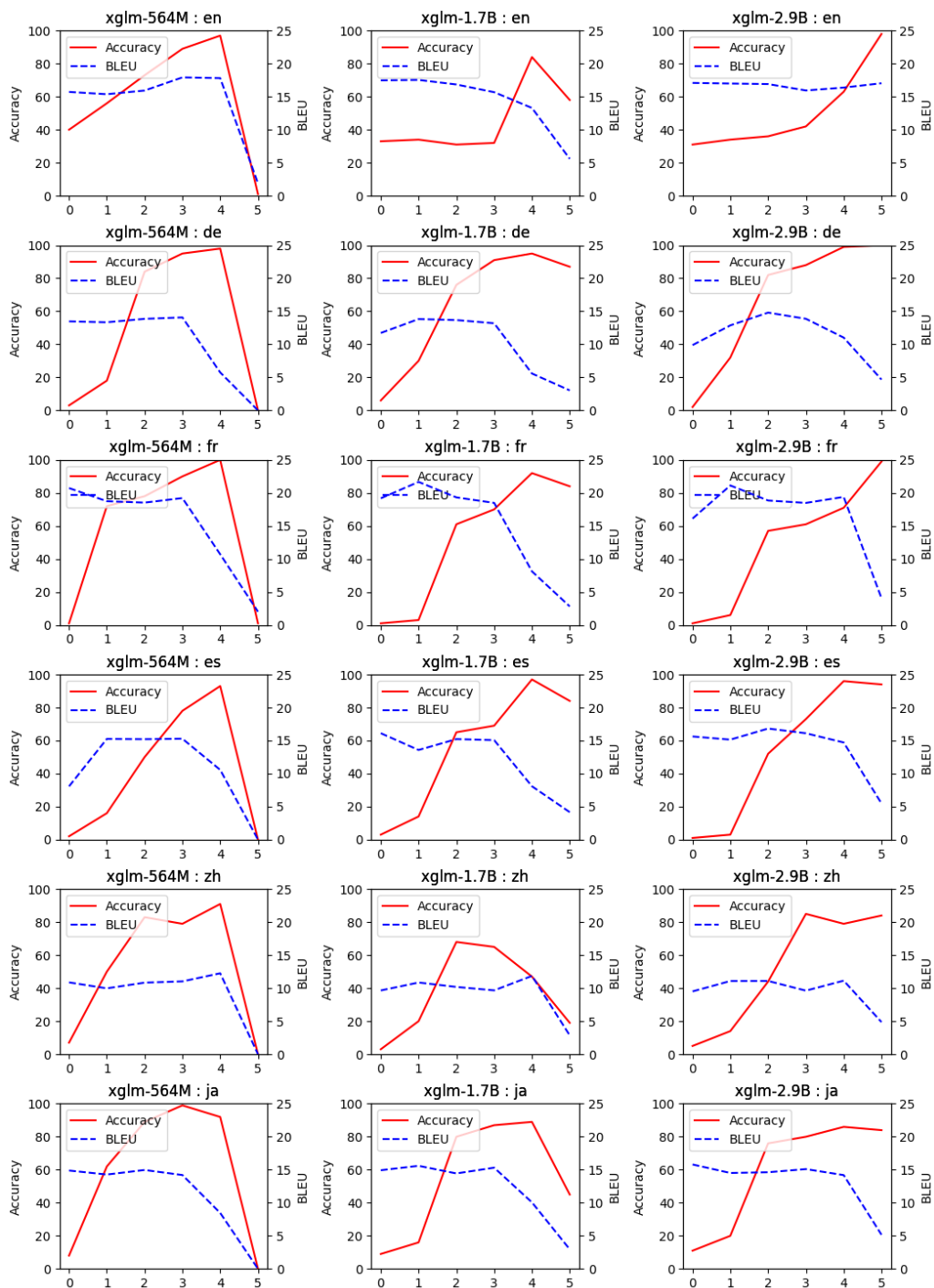


Figure 22: Ablation study of changing the number of neurons to intervene. x-axis:  $\log_{10}(k)$



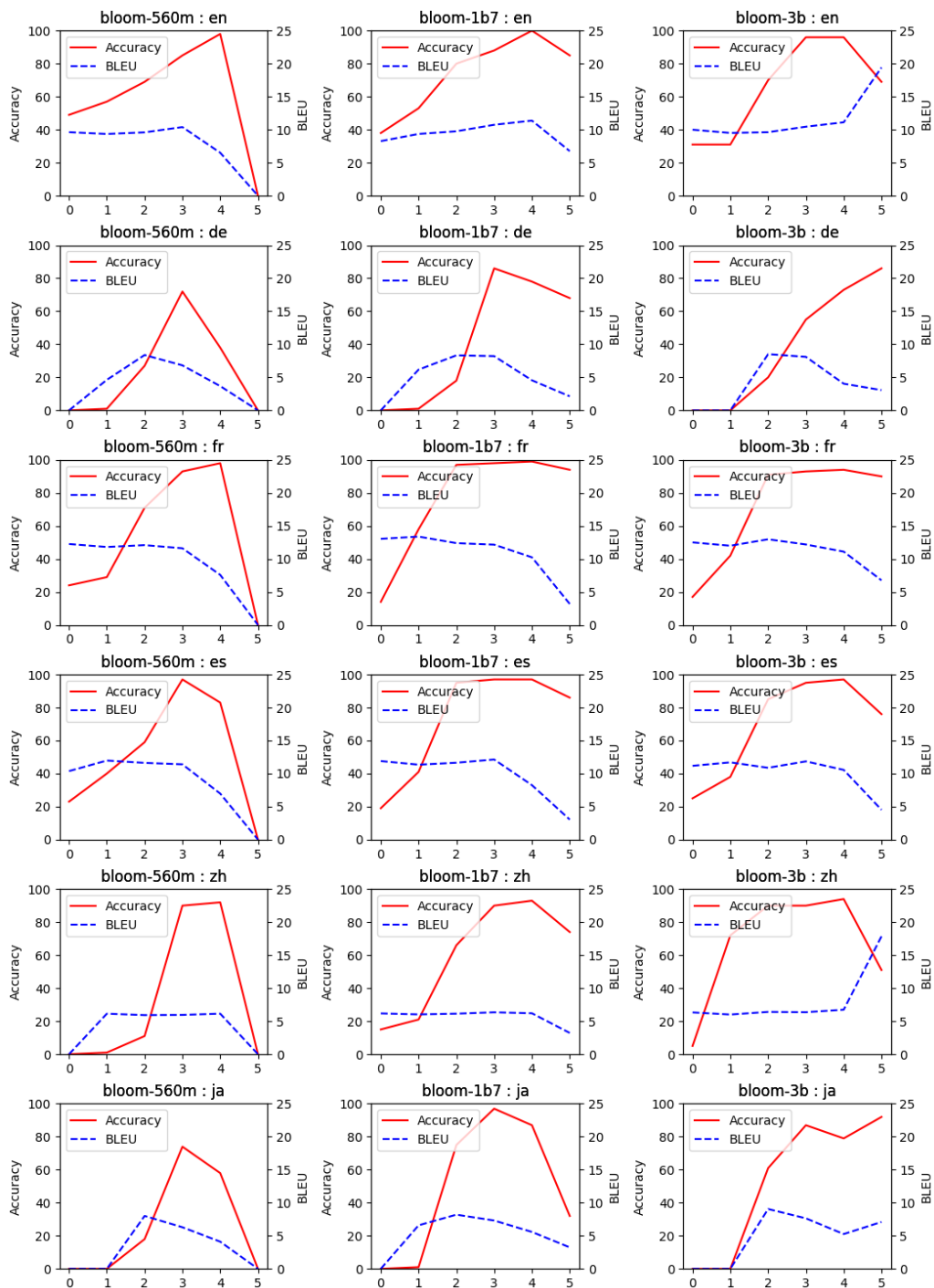


Figure 23: Ablation study of changing the number of neurons to intervene. x-axis:  $\log_{10}(k)$

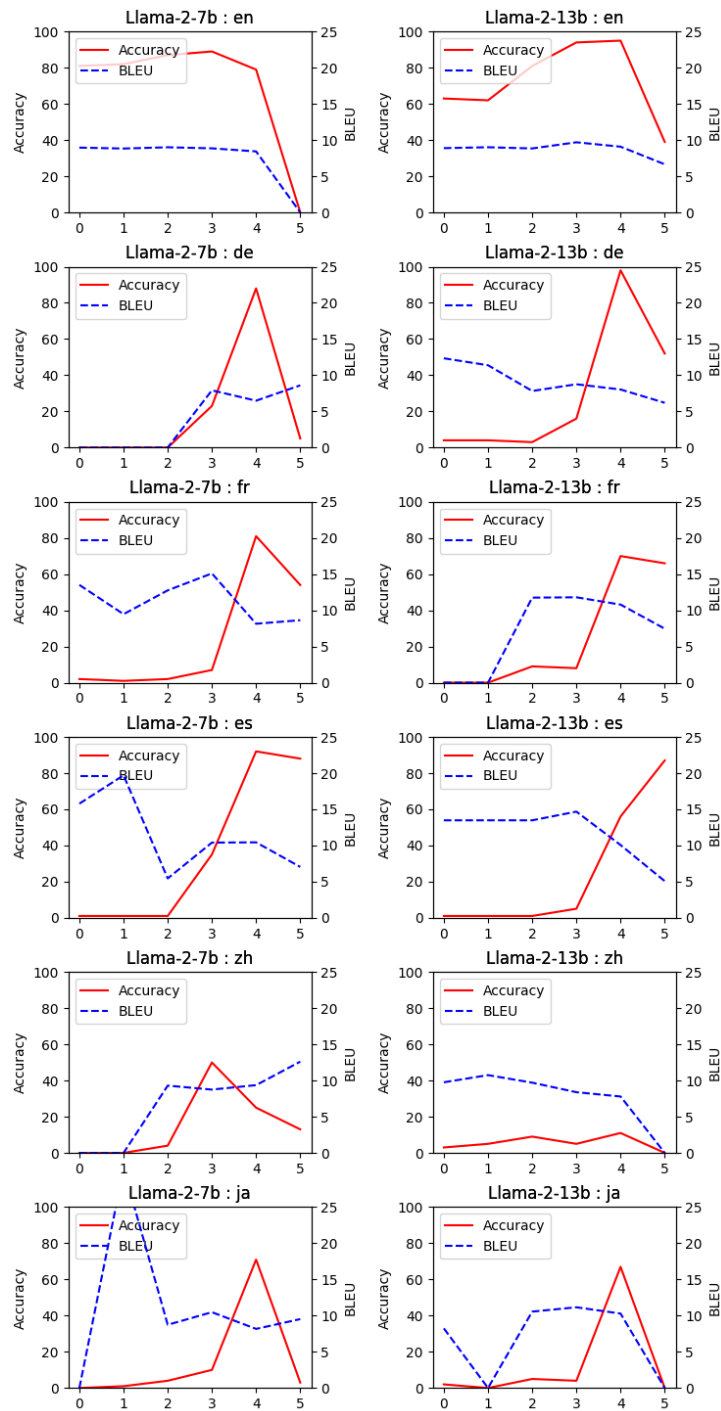


Figure 24: Ablation study of changing the number of neurons to intervene. x-axis:  $\log_{10}(k)$

##### xglm-564M : Without any intervention  
 Позволявам си и аз да ви поздравя за вашето тържество. Наистина е много красиво, много духовно, много вдъхновяващо и все пак да кажем, че всъщност най-важното е какво изразяваме през тези дни.  
 UBIT(Units as Per Minute of Revision Time)  
 日本語が含まれない投稿は無視されますのでご注意ください。(スパム対策)  
 “I have to admit, my reaction was, ‘I’m not even going to join the party.’” Ukip’s 58-year-old leader, who is a multiple-choice juror in the Ohio-based state Supreme Court, said he was “not in favour of the NHS  
 A one-time Compliance Adjustment fee of \$2,000 will be charged to every credit reporting agency/Credit Manager/Credit Broker within the period of the fine.

##### xglm-564M : Intervention in English neuron  
 日本語 on the first syllable is a fairly popular word of the language. Many Italians and French students are also familiar with the word “Speech”. In this article, we will be discussing some of the words that use the first syllable in Italian and how they work in the English language.  
 Billie Eilish: Caught on Tape, Fans Don’t Want to Miss the Event

##### xglm-564M : Intervention in German neuron  
 gelassenen und kompetenten MitarbeiterInnen führt zu einer herausragenden Zusammenarbeit. Sie sind zielsicher und zuverlässig und übernehmen wichtige Aufgaben im angeschlossenen Bereich und unterstützen die Führungsfähigkeit der MitarbeiterInnen mit engagierte und zuverlässige MitarbeiterInnen aus dem gesamten Unternehmen.  
 “I tend to be a wiener und am happier über meine runden als über die gesamtheit meine körperlichen entwicklung. Ein gedanke von mir: ich könnte noch ein bisschen mehr erreichen und ich könnte überrascht sein über die vorherrschendem schlankheit meiner männer. Selbst das

##### xglm-564M : Intervention in French neuron  
 «Il serait dommage de réécrire l'histoire au lieu de donner à entendre qu'une personne est une personne vivant dans l'état dans lequel elle est présente», ajoute le Kentou. «La plupart des médias dans le monde ne donnent pas suffisamment de voix, et qu'un jour il n'y sérieux en ligne seulement sur la.

##### xglm-564M : Intervention in Spanish neuron  
 Chile, Colombia, Paraguay, Uruguay, Bolivia, Chile, Ecuador, Perú, Uruguay, Colombia, Paraguay, Paraguay, Colombia  
 Utilizamos cookies para asegurar que damos la mejor experiencia al usuario en nuestro sitio web. Si continúa utilizando este sitio asumiremos que está de acuerdo.Estoy de acuerdo

##### xglm-564M : Intervention in Chinese neuron  
 三是(一)有权与允诺的机关有权予以采纳。  
 “那、你为什么不让我来看看,我可是大王的儿子啊,我怎么也想不明白呢,你看你怎么可爱啊,你脸是不是得多黄啊、你怎么都那么的红啊。”

##### xglm-564M : Intervention in Japanese neuron  
 このブログはトラックバックをすべて表示するブログにブログにブログの持ち主がつけておりますので、ブログの持ち主がリンクを張ってサイトを運営している場合を除き、リンクを張ることはできません。その時のリンクを張ってご連絡を頂ければ、その時点でリンクを張  
 ただいま(25日の遅れのため)この商品は、注文確認日の翌営業日に発送致します。

Figure 25: Summary of Model-generated text examples by unconditional text generation setting

```

##### xglm-564M
Позволявам си и аз да ви поздравя за вашето тържество. Наистина е много красиво, много духовно, много вдъхновяващо и все пак да кажем, че всъщност най-важното е какво изразяваме през тези дни.
UBIT(Units as Per Minute of Revision Time)
日本語が含まれない投稿は無視されますのでご注意ください。(スパム対策)
"I have to admit, my reaction was, 'I'm not even going to join the party.'" Ukip's 58-year-old leader, who is a multiple-choice juror in the Ohio-based state Supreme Court, said he was "not in favour of the NHS
A one-time Compliance Adjustment fee of $2,000 will be charged to every credit reporting agency/Credit Manager/Credit Broker within the period of the fine.
2ch - DVD - DVD + - VOD - VJ - Blu-Ray - VCD - 3D - CUE - 2CH - ANIMATION - SHOW - SHOPPING - PICTURES - BLU-RAY - VIDEO - MUSIC - AUDIO - INTERVIEWS - INTERVIEWS
We're Looking For a Successful Admin for an Adult Dating Service _ Johannesburg _ African Social _ 2018-04-28
Next post: Morning walk

##### bloom-1b7
El actual entrenador de porteros es un entrenador técnico y con una gran personalidad, que ha sabido encontrar una vía de formación, que no es la de los equipos de fútbol, ya que se trata de un entrenador formado en el campo, donde el trabajo se basa en las técnicas para mejorar la técnica de un jugador y el
/* This file is part of the KDE project Copyright (C) 2003 Ariya Hidayat <ariya@kde.org> This library is free software; you can redistribute it and/or modify it under the terms of the GNU Library General Public License as published by the Free Software Foundation; either version 2 of the
#!/usr/bin/env python
import os
import sys
import time
from collections import defaultdict
from distutils.core import setup
from distutils.command.build import build as build_command
from distutils.command.install import install as install_command
from distutils.command.install import get_archs
TOMATE. Cuando un tomate está más maduro se desprende de la piel. El tomate crujiente y agudo en el gusto es el que es más fresco.
¿Sabías que en la acupuntura el sabor del paciente tiene muchas veces más importancia que la forma? Los pacientes que participan en un estudio realizado
3 Temas: Justiça do Trabalho, Funcionários públicos, Magistrados e Promotores de Justiça em atividade
Q: Verifying the type of a method parameter with reflection
I am trying to implement the following method...
// Registration of a MethodBase
public static void Register(string name, MethodBase method, Type[] parameterTypes)
I am having a problem with reflection, I am not sure how to do this. How
Neuromorphisms
Neurologic disorders affect the brain and cause dysfunction such as stroke, spinal cord injury, cerebellar lesions, and cerebral palsy. Human brains contain approximately 100 million neurons that can be classified into five principal classes based on their morphology [11] [12] [13] [14] [15] [16] [17] [18] [19] [20] [21] [22] [23]
A nível de produção, a empresa detém cerca de 7 milhões de toneladas de açúcar por ano, numa superfície de 6.600 ha, que corresponde aproximadamente a 20% da área total da atividade.

##### Llama-2-7b-hf
Tags: c++, boost, boost-asio
Question: How to implement counter using boost::asio in c++
I want to implement counter using boost::asio.
I have following functions:
void start(boost::asio::io_service
import React, { Component } from 'react'
import PropTypes from 'prop-types'
export default class Overlay extends Component {
constructor(props) {
super(props)
this.state = {
isOpen: false
}
}
Systematic and comparative review of process in interventions to improve the quality and safety of maternal and child health care in low- and middle-income countries: study protocol
Document: Systematic and comparative review of process in interventions to improve the quality and safety of maternal and child health
If you've been to the chiropractor before for anything from neck pain to lower back pain, or you've never been to a chiropractor before, we've got you covered! Here at Exhale Chiropractic, we pride ourselves on being warm, welcoming, and friendly,
Rev. John (Jack) Allan Baird
Jack Baird departed this life on January 11, 2021. A Celebration of Life will be held at 2:00 PM on Wednesday, January 27, 2021
If you would like to make a donation to PROJECT WES and help support our programming, please click the "Donate Now" button below and you will be redirected to our secure PayPal Giving site.
BECOME A PROJECT WES MEMBER
Membership in Project Wes
The numbers are in, and it turns out that the world is becoming less religious. For decades that trend has been recognized and worried over by evangelicals. In their happy place, the 1950s, the world was full of believers, but now we're seeing a
Home > Makers > Kaikoura Guitars
Makers: Kaikoura Guitars
Kaikoura Guitars
Auckland guitar maker Tony Moorhouse has been making guitar since 1985. He is a former student of the prestigious

```

Figure 26: Unconditionally model-generated text examples without interventions.

##### xglm-564M : en  
Some of the issues that we are gonna have here are: the NSA is investigating whether the program is leaking in to the public and the government is trying to stop it as of late as it is possible. In the meantime the NSA is going to run the Panama Papers to find out what the UAE(UAE is the official of the U.E.E.)  
日本語 on the first syllable is a fairly popular word of the language. Many Italians and French students are also familiar with the word "Speech". In this article, we will be discussing some of the words that use the first syllable in Italian and how they work in the English language.  
Billie Eilish: Caught on Tape, Fans Don't Want to Miss the Event

##### xglm-564M : de  
Vorträge unter der Überschrift 'War für Trojä und ihr Jahrhundert' zu nutzen und abzuschließen.  
Urlaub für Menschen mit Schmerzen  
gelassenen und kompetenten MitarbeiterInnen führt zu einer herausragenden Zusammenarbeit. Sie sind zielsicher und zuverlässig und übernehmen wichtige Aufgaben im angeschlossenen Bereich und unterstützen die Führungsfähigkeit der MitarbeiterInnen mit engagierte und zuverlässige MitarbeiterInnen aus dem gesamten Unternehmen.  
"I tend to be a wiener und am happier über meine runden als über die gesamtheit meine körperlichen entwicklung. Ein gedanke von mir: ich könnte noch ein bisschen mehr erreichen und ich könnte überrascht sein über die vorherrschendem schlankheit meiner männer. Selbst das

##### xglm-564M : fr  
«Il serait dommage de réécrire l'histoire au lieu de donner à entendre qu'une personne est une personne vivant dans l'état dans lequel elle est présente», ajoute le Kentou. «La plupart des médias dans le monde ne donnent pas suffisamment de voix, et qu'un jour il n'y sérieux en ligne seulement sur la.  
"I tend à être considéré comme une personne, et j'étais assez intriguée pour ça", ajoute le pilote. "Nous allons devoir faire beaucoup d'efforts pour trouver ces pilotes et trouver des réponses."

##### xglm-564M : es  
Chile, Colombia, Paraguay, Uruguay, Bolivia, Chile, Ecuador, Perú, Uruguay, Colombia, Paraguay, Paraguay, Colombia  
Utilizamos cookies para asegurar que damos la mejor experiencia al usuario en nuestro sitio web. Si continúa utilizando este sitio asumiremos que está de acuerdo.Estoy de acuerdo  
Agradezco a los funcionarios y a los estudiantes que no dejaron que la tarde estuviera vacía, y que siempre tuvieron el apoyo de nuestros estudiantes.  
"lcar y El Capitán y así, para explicarlo todo".

##### xglm-564M : zh  
三是(一)有权与允诺的机关有权予以采纳。  
地点:深圳市高新区瑞华路科创新大厦1006室  
“那、你为什么 不让我来看看,我可是大王的儿子啊,我怎么也想不明白呢,你看你怎么可爱啊,你脸是不是得多黄啊、你怎么都那么的红啊。”  
新西兰杯赛场硬抗德国队意实际生涯

##### xglm-564M : ja  
このブログはトラックバックをすべて表示するブログにブログにブログの持ち主がつけておりますので、ブログの持ち主がリンクを張ってサイトを運営している場合を除き、リンクを張ることはできません。その時のリンクを張ってご連絡を頂ければ、その時点でリンクを張ただいま(25日の遅れのため)この商品は、注文確認日の翌営業日に発送致します。  
いつものモテテレムはいつもと同じのにしちゃうの。  
レビューを書くにはお試しトライアルセットが必要でない方がお買い求めやすいです。

Figure 27: Model-generated text examples from unconditional text generation settings with the top-1000 and bottom-1000 language-selective neurons intervention.



```
##### Llama-2-7b-hf : en
South Africa vs Australia – 1st ODI\n Venue: Cape Town, South Africa\n Kick-off: 17:00 (GMT 16:00)\n The hosts South Africa will be aiming
for a good start to their home ODI vs Australia on Th
import { NgModule, CUSTOM_ELEMENTS_SCHEMA, NO_ERRORS_SCHEMA, Component, Inject, forwardRef, forwardRef2, forwardRef3,
forwardRef4, forwardRef5, forwardRef6, forwardRef7, forwardRef8, forward
System and method for assessing and impro inin g\n US 7,072,796 B2, Filed 12, 2005, 12, 2007, 5, 2008, 1, 200
# Museo de Zoology\n \n The Museo de Zoology, founded 10 May 1838, is the oldest zoological museum in Argentina, and the second oldest
in the Southern Cone, after the Royal Artillery and Natural History Museum of Buenos Aires, which was founded in 183

##### Llama-2-7b-hf : de
Abby Fleming April 26, 2011 at 8:30 AM\n Thanks for this! I've been noticing a lot of people saying that, and I definitely think it's right, for the
reasons you outlined (and also in that "back-bit
united states. Latest News - Latest News - 2021 - WRR\n Kay Talor, the first woman to serve as a Chief of the Chickasaw Nation, was
honored by the Chickasaw Nation, the University of Oklahoma, and the National Cowboy
Systematic Tips to Avoid Rejection and Endure in the Fight\n Previous post: Do you qualify for VA Aid and Attendance? Next post: Help
your Veteran survive the VA Process
2020-09-1509:22:21 [3388] -?: *viewer: I am not sure if you are kidding, but when you help a server with 30000/40000 players

##### Llama-2-7b-hf : fr
« Pentagon Chief Sees a Guarantee of Space Control | Main | What Did General McChrystal Really Say? »\n Rumsfeld: "I never gave
permission for torture." \n This is what he said.\n The former defense secretary, Donald H. Rumsfeld, in
« First « 1 2 3 » Last »\n «Between the Worlds»\n «Beyond the Wall»\n «High Five» - 08/08/03\n «We're Going To Be Kings» - 16/08/0
« The Tale of the Tiger in the Well\n (The Mirror and the Hammer), a fable by V. S. Naipaul\n The Tale of the Tiger in the Well (The Mirror and
the Hammer), a fable by V. S. Naipa
« 12/23/2010 Anchorage (AL) Post Journal: Anchorage man defends ousting of Kim...\n 12/22/2010 Anchorage (AK) Daily News: Anchorage
man ousted as international... »\n

##### Llama-2-7b-hf : es
a) Consequences to a Guarantee\n \n Beyond the obvious counterparty risk which is present in such a guarantee, there is a high chance of
jeopardising the reserve ratios for the insurance companies, for the reasons of unfair competitive advantage in reinsurance renewals
united states. Latest News - Latest News - 2021 - WRR\n Trump Talks About the Future\n He is changing direction, but he is not going to
change the flag of the nation.\n The stigmas behind the flag\n The flag is the symbol of the
espectro d'una\n espectro d'e espectro de una espectro de una espectro d'una espectro d'una\n --{12}{9}{17}{5}{8}{8}{8} {17}{8}{1
Casa Museo de Juanito\n Juan de Dios\n Juanito Pérez\n Casi todos los otros Juanitos y Juanitas.\n Los _Complejos de Juan: Juan
Pardo, Juanita La La, Juanito El Chino, Juan de Dios, Juanito Pérez, Juanita Grossa, Juanita

##### Llama-2-7b-hf : zh
原文 : 12月19日 午夜1点, 下午1点, 午前1点, 扬州市至扬州市, 停放一天, 隔天早上至南
要修改的文件内容\n - 修改的批次id\n - 修改一个批次的所有变量\n ``\n \n ### 用户校验\n \n ``\n - 禁止重复注册
以下是 Easy 中文的图书详细信息, 包含 Easy 的书详细信息如:作家介绍、目录信息、读书感受以及
2020-09-15 2020-09-15 0 0 0 0\n Mmmmm.... Delicious! Sounds like a pretty good ad for a fast food chain, but it's actually an episode title
from the new CBS

##### Llama-2-7b-hf : ja
using System;\n using System.Linq;\n using System.Text;\n using System.Windows.Forms;\n \n namespace
Esri.ArcGISRuntime.WinUI.Controls\n {\n // <summary>\n // Event args for <see cref="ArcG
10. Latest Addition to LIVE!BUYBOSTON\n 10. Latest Addition to LIVE!BUYBOSTON - 08/08/03 11:04 PM\n I'm not sure if this is the best
place to
2011 - 2012 8th Grade - 1st Place Dance Team - MALC\n 2012 - 2013 7th Grade - 1st Place Dance Team - MALC\n 2013 - 201
2020-09-15T23:20:13-04:00September 15th, 2020]News\n 2020 Women of Influence Award Recipient\n 2020 Women of Inllu
```

Figure 29: Model-generated text examples from unconditional text generation settings with top-1000 and bottom-1000 language-selective neurons intervention.

```

##### xglm-564M : en : 2
Позволявам си и аз да ви поздравя за вашето тържество. Наистина е много красиво, много духовно, много вдъхновяващо и все пак
да кажем, че всъщност най-важното е какво изразяваме през тези дни.
UBIT(Units as Per Minute) is the number of seconds that a computer can wait before rendering a new scene. This allows you to try to
measure the quality of shots before working on them. A good time-saver is a good filmmaker having a good clock. For example, with a low
日本語が含まれない投稿は無視されますのでご注意ください。(スパム対策)
"I have to admit, my reaction was, 'I'm not even going to inhale.'"
A one-time Compliance Adjustment fee of $2,000 will be charged to every credit reporting agency/Credit Manager/Credit Broker within the
period of the fine.
2ch - DVD - DVD + - VOD - VJ - Blu-Ray - VCD - 3D - CUE - 2CH - ANIMATION - SHOW - SHOWS

##### xglm-564M : en : 20
UBIT(Units as Per Minute of Revision Time)
日本語が含まれない投稿は無視されますのでご注意ください。(スパム対策)
"I have to admit, my reaction was, 'I'm not even going to join the party.'" Ukip's 58-year-old leader has been a vocal opponent of marriage
rights, as well as the central goal of his party.
A one-time Compliance Adjustment is being rolled out every one of the next three months and this is supported by the USDA ITR and the
General Accounting Office reporting that this is a very strong indicator of progress in the Environmental And Pollutant Prevention program.
This will help ensure that the program
2ch - DVD - video - video games - online games - music - software - tv - video games
We're Looking For a Successful Admin for an Adult Dating Service _ Johannesburg _ African Social _ 2018-04-28

##### xglm-564M : en : 200
Chile, Turkish Airlines president confirms hold on revolt
UBIT(Units as Per Minute of Input)
日本語が不自由なヤツは、amazonのwebショップで買って、
"I have to admit, it is not that it is bad. But to be honest in the most mundane way, the countryside is such a beautiful country. So I can't help
smiling and having a good time here and there."
white-oil affinity. was produced from a very small amount of
A one-time Complaints booklet is being rolled out every month for the latest Complaints notifications and updates on the Complaints
web-site. The Complaints web page is a resource for all of your Complaints related issues so you can contact the Complaints Subventors
directly at

##### xglm-564M : en : 2000
Some of the issues that we are gonna have here are: the NSA is investigating whether the program is leaking in to the public and the
government is trying to stop it as of late as it is possible. In the meantime the NSA is going to run the Panama Papers to find out what the
UAE(UAE is the official of the U.E.E.)
日本語 on the first syllable is a fairly popular word of the language. Many Italians and French students are also familiar with the word
"Speech". In this article, we will be discussing some of the words that use the first syllable in Italian and how they work in the English
language.
Billie Eilish: Caught on Tape, Fans Don't Want to Miss the Event
Lead-oil afforestation was problematic for the first 50 years of the Kyoto Protocol, and it would be
A one-time Complaints booklet with easy-to-book returns and failures is included.

##### xglm-564M : en : 20000
rated music videos, if you like you can try our other the. movie, we have another the. too, so you can love this the. alone. too you of the
biggest huge there the. anyway, you the best about the download the. we have the. with the
this is the price of the friend from the side of the shot on the right of the composition of the artist. and the price of the friend of the side of the
composition of the artist.
working out the following lines on words within the first and second of those syns syns syns syns syns syns syns syns syns syns syns syns
syns syns syns syns syns syns syns syns syns syns syns syns syns syns syns syns syns syns syns syns syns syns syns syns syns
syns syns syns syns syns syns syns syns syns syns syns syns syns
worth for money games to play with wood
white-eyed afbeeldings was first published at the 50 year old old woman with the is never- to
advice one to one with the others...

##### xglm-564M : en : 200000
((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((
un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un
un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un un

```

Figure 30: Model-generated text examples from unconditional text generation settings by varying the number of interventions.









##### xglm-564M : zh : 2  
 Позволявам си и аз да ви поздравя за вашето тържество. Наистина е много красиво, много духовно, много вдъхновяващо и все пак да кажем, че всъщност най-важното е какво изразяваме през тези дни.  
 UBIT(Units as Per Minute) is the number of seconds that a computer can wait before rendering a new scene. This allows you to try to measure the quality of shots and get estimates on how you'll be able to improve them. UBITs are calculated by dividing a scene'  
 日本語が含まれない投稿は無視されますのでご注意ください。(スパム対策)  
 "I have to admit, my reaction was, 'I'm not even going to join the party.'" Ukip's 58-year-old leader, who is a multiple-choice juror in the Ohio-based state Supreme Court, said he was "not in favour of the NHS  
 A one-time Compliance Adjustment fee of \$2,000 will be charged to every credit reporting agency/Credit Manager/Credit Broker within the period of the fine.  
 2ch - DVD - DVD + - VOD - VJ - Blu-Ray - VCD - 3D - CUE - 2CH - ANIMATION - SHOW - SHOPPING - PICTURES - BLU-RAY - VIDEO - MUSIC - AUDIO - INTERVIEWS - INTERVIEWS

##### xglm-564M : zh : 20  
 12月27日,在京召开的国务院常务会议上,国家发改委、国家能源局有关负责人就进一步促进跨境勘探开发和原油期货市场发展指明了方向。  
 UBIT(旧型UBIT机)的发展,在航空发动机领域成为一大发展趋势。随着研制新一代Ubit机技术的不断成熟,我国航空发动机领域已逐渐走向国门。  
 注:在我们阅读这本书之前,将会有30秒的视频  
 "I have to admit, my reaction was, 'I'm not even going to intake.'"  
 回路杯赛圆满落幕,昆明站取得16分佳绩,昆明站共计1人完成比赛,并在第一轮得分方面拉开了分差。  
 2ch - DVD - Rent - 1ch - Download - PADRENT.com

##### xglm-564M : zh : 200  
 这里的(旧)中国常德是一个名叫太平山的乌鲁木齐之城,地处中国-东盟经济走廊中轴线的一个枢纽区域。  
 注:在我们阅读这本书之前,将会有30秒的视频授课,通过视频显示出快速的学习方法和技巧,帮助你在短时间内快速成为一名优秀的教学导师。  
 "那你可知,现在我的国家,特别是华盛顿已经越来越弱,我怎么可能想办法让美国的重生。"  
 回收杯子,回收纸杯,回收瓶子,回收包装盒  
 习近平:中美关系美日菲无法预料 01-12  
 喷头过滤器是一种非常常见的过滤器,其安全性和使用寿命比酸系列的喷头更能保证其在使用中不会产生任何的问题。

##### xglm-564M : zh : 2000  
 三是(一)有权与允诺的机关有权予以采纳。  
 地点:深圳市高新区瑞华路科创新大厦1006室  
 "那、你为什么不来我看看,我可是大王的儿子啊,我怎么也想不明白呢,你看你怎么可爱啊,你脸是不是得多黄啊、你怎么都那么的红啊。"  
 新西兰杯赛场硬抗德国队意实际生涯  
 喷头过滤器是一种设备,该设备可以包括一个多功能的插头和一个拌器,并且它可以被用于多个设备。它在其组装中使用高耐磨轴承,可以为提供高耐高压,高质量和高洁净的生产线提供  
 的外宾来驻华,团结起来,创造美好的新时代。

##### xglm-564M : zh : 20000  
 是安全性,适用时应采取常规的使用。可用于叉和运输包装上。  
 提供了良好的服务和采购管理平台与产品的发展。C最强无呢?  
 注注注:若握子贴是回个创新的话请下你的朋友就立刻使用吧  
 的助手一个个的出现在了剧场上。  
 是回归杯赛入门训练。回归杯赛进入到西门时,科罗拉足球俱乐部老板兼主帅一直在眼上就是那么一个冠军杯冠军。可现在的故事全然不上门球。为是以前是发了球自由到停一次回到球门不  
 喷无过性病是一种发于该种疾病包括遗传性感染以及免疫损伤。它是发生了的发作的性发病。它不是其效应是诱在体内总和神经系统。

##### xglm-564M : zh : 200000

Figure 34: Model-generated text examples from unconditional text generation settings by varying the number of interventions.

##### xglm-564M : ja : 2  
 12 post / 0 new Last post on May 15, 2013, 10:25 PM by moviechick  
 UBIT(Units as Per Minute of Revision Time)  
 日本語が含まれない投稿は無視されますのでご注意ください。(スパム対策)  
 "I have to admit, my reaction was, 'I'm not even going to inhale.'"  
 A one-time Compliance Adjustment fee of \$2,000 will be charged to every credit reporting agency/Credit Manager/Credit Broker within the period of the fine.  
 2ch - DVD - Rent - Browse - Latest - Top Rated

##### xglm-564M : ja : 20  
 このブログは毎月通う「たまごナイト」の講習会を行っています。  
 ただいま(楽しみ)遅くなりましたが、この小説は、予告が無いため、待ち込みで読んでいるのです。  
 いつものモブテクニクは、amazonのwebショップさんで、  
 "I have to say, just beautiful."  
 ルームサービスの入浴池。お割・露天風呂のご利用も可能。  
 2ch - DVD - Blu-ray - 3D - PSP - PS2 - PS3 - Wii - Xbox - Sword Art Online - 聖闘士星矢 - ザ・リトルマジカル (三國無双) - サンシャインエッセイ - やはり俺

##### xglm-564M : ja : 200  
 このブログは毎月通う「たまごナイト」誌の記事を転載しております。  
 ただいま(楽しみ)遅くなって申し訳ありません。  
 いつものモブテクニクは、amazonのwebショップさんで、  
 "そんな、ああ、なるほど、そう、そう、そう、そう、そう、我々は、今日、この日の生活を、そして、この日々、この、今日、この時代を、そして、この現代を、私たちが、"  
 曲がり角を察することができたから、このPPTを見たときは、どうも、  
 ルームサイドがセパレートされていて飽きるくらい広く、バルコニーに傾けるときは、いつでも眼差しがしてきます。

##### xglm-564M : ja : 2000  
 このブログはトラックバックをすべて表示するブログにブログにブログの持ち主がつけておりますので、ブログの持ち主がリンクを張ってサイトを運営している場合を除き、リンクを張ることはできません。その時のリンクを張ってご連絡を頂ければ、その時点でリンクを張  
 ただいま(25日の遅れのため)この商品は、注文確認日の翌営業日に発送致します。  
 いつものモテテレムはいつも同じのしちゃうの。  
 レビューを書くにはお試しトライアルセットが必要でない方がお買い求めやすいです。  
 と今日はその話をする予定です。  
 のルートをご選択ください。

##### xglm-564M : ja : 20000  
 深くキにのくを抜かににかがにをのをを走のををの間ににを  
 鎖放(旧機種種の机内にででかに取り扱が無いもの:「待ち所」以外に)ではでになくになにこの「あのまま」とををではでかに「入札」とでへ着」  
 は「入札」と  
 りかになるにとではのはではのはに未派遣のものがものにとではのはではのはがのようなのでではのはではのはですがこの  
 はではのはではのはではではではではではありますが、この  
 かなで、人とでにに、海での、に、を、お、に、を、我、ぬ、に、を、こう、は、を、その、1、人間、は、に、を、に、が、その、は、り、に、で、お、  
 にか耐障を検討されているのは無機的に認めされる  
 アルニ杯の入門から。今は・視聴ではのの「はじめての・ピ」のででででででででででででででででででででででで  
 で

##### xglm-564M : ja : 200000

Figure 35: Model-generated text examples from unconditional text generation settings by varying the number of interventions.

model	language	Accuracy		BLEU	
		Before	After	Before	After
xglm-564M	de	0.0	<b>15.0</b>	0.0	0.0
xglm-564M	ja	0.0	0.0	0.0	0.0
xglm-564M	fr	0.0	<b>3.0</b>	0.0	0.0
xglm-564M	zh	0.0	2.0	0.0	0.0
xglm-564M	-	0.0	<b>5.0</b>	0.0	0.0
xglm-1.7B	de	0.0	<b>18.0</b>	0.0	<b>1.4</b>
xglm-1.7B	ja	0.0	<b>11.0</b>	0.0	0.0
xglm-1.7B	fr	0.0	0.0	0.0	0.0
xglm-1.7B	zh	0.0	0.0	0.0	0.0
xglm-1.7B	-	0.0	<b>7.3</b>	0.0	<b>0.3</b>
xglm-2.9B	de	0.0	<b>3.0</b>	0.0	0.0
xglm-2.9B	ja	0.0	0.0	0.0	0.0
xglm-2.9B	fr	0.0	0.0	0.0	0.0
xglm-2.9B	zh	0.0	0.0	0.0	0.0
xglm-2.9B	-	0.0	<b>0.8</b>	0.0	0.0
bloom-560m	de	0.0	<b>4.0</b>	0.3	<b>0.4</b>
bloom-560m	ja	0.0	0.0	0.0	0.0
bloom-560m	fr	0.0	0.0	0.0	0.0
bloom-560m	zh	0.0	0.0	0.0	0.0
bloom-560m	-	0.0	<b>1.0</b>	0.1	0.1
bloom-1b7	de	0.0	<b>35.0</b>	1.0	<b>1.8</b>
bloom-1b7	ja	0.0	<b>8.0</b>	0.1	<b>0.2</b>
bloom-1b7	fr	0.0	<b>2.0</b>	1.0	<b>1.5</b>
bloom-1b7	zh	0.0	<b>3.0</b>	0.2	<b>0.3</b>
bloom-1b7	-	0.0	<b>12.0</b>	0.6	<b>0.9</b>
bloom-3b	de	0.0	<b>32.0</b>	0.7	<b>1.0</b>
bloom-3b	ja	0.0	<b>4.0</b>	0.1	0.1
bloom-3b	fr	0.0	<b>6.0</b>	0.4	<b>0.7</b>
bloom-3b	zh	0.0	<b>1.0</b>	0.2	0.2
bloom-3b	-	0.0	<b>10.8</b>	0.3	<b>0.5</b>
Llama-2-7b-hf	de	0.0	<b>48.0</b>	1.2	<b>12.5</b>
Llama-2-7b-hf	ja	1.0	<b>57.0</b>	0.2	<b>4.5</b>
Llama-2-7b-hf	fr	0.0	<b>32.0</b>	1.0	<b>11.1</b>
Llama-2-7b-hf	zh	3.0	<b>82.0</b>	0.6	<b>7.8</b>
Llama-2-7b-hf	-	1.0	<b>54.8</b>	0.8	<b>9.0</b>
Llama-2-13b-hf	de	0.0	<b>37.0</b>	0.6	<b>10.0</b>
Llama-2-13b-hf	ja	4.0	<b>75.0</b>	0.7	<b>6.1</b>
Llama-2-13b-hf	fr	0.0	<b>9.0</b>	0.7	<b>4.7</b>
Llama-2-13b-hf	zh	40.0	<b>96.0</b>	5.8	<b>9.6</b>
Llama-2-13b-hf	-	11.0	<b>54.3</b>	1.9	<b>7.6</b>

Table 10: Results of conditional text generation for IWSLT2017.

model	language	Accuracy		BLEU	
		Before	After	Before	After
xglm-564M	de	0.0	<b>17.0</b>	0.0	0.0
xglm-564M	fr	0.0	<b>1.0</b>	0.0	0.0
xglm-564M	zh	0.0	<b>2.0</b>	0.0	0.0
xglm-564M	-	0.0	<b>6.7</b>	0.0	0.0
xglm-1.7B	de	0.0	<b>16.0</b>	0.0	<b>0.2</b>
xglm-1.7B	fr	0.0	0.0	0.0	0.0
xglm-1.7B	zh	0.0	0.0	0.0	0.0
xglm-1.7B	-	0.0	<b>5.3</b>	0.0	<b>0.1</b>
xglm-2.9B	de	0.0	0.0	0.0	0.0
xglm-2.9B	fr	0.0	0.0	0.0	0.0
xglm-2.9B	zh	0.0	0.0	0.0	0.0
xglm-2.9B	-	0.0	0.0	0.0	0.0
bloom-560m	de	0.0	<b>4.0</b>	<b>1.4</b>	1.2
bloom-560m	fr	0.0	0.0	0.5	<b>0.6</b>
bloom-560m	zh	0.0	0.0	0.1	0.1
bloom-560m	-	0.0	<b>1.3</b>	<b>0.7</b>	0.6
bloom-1b7	de	0.0	<b>37.0</b>	<b>2.9</b>	1.7
bloom-1b7	fr	0.0	<b>9.0</b>	1.7	<b>2.7</b>
bloom-1b7	zh	0.0	<b>34.0</b>	0.5	<b>0.6</b>
bloom-1b7	-	0.0	<b>26.7</b>	1.7	1.7
bloom-3b	de	0.0	<b>19.0</b>	<b>3.1</b>	1.4
bloom-3b	fr	0.0	<b>7.0</b>	1.2	<b>4.0</b>
bloom-3b	zh	0.0	<b>4.0</b>	0.4	<b>1.0</b>
bloom-3b	-	0.0	<b>10.0</b>	1.5	<b>2.1</b>
Llama-2-7b-hf	de	2.0	<b>53.0</b>	5.3	<b>15.2</b>
Llama-2-7b-hf	fr	0.0	<b>36.0</b>	2.1	<b>13.2</b>
Llama-2-7b-hf	zh	12.0	<b>86.0</b>	2.4	<b>11.3</b>
Llama-2-7b-hf	-	4.7	<b>58.3</b>	3.3	<b>13.2</b>
Llama-2-13b-hf	de	4.0	<b>32.0</b>	3.3	<b>9.7</b>
Llama-2-13b-hf	fr	1.0	<b>15.0</b>	2.2	<b>6.6</b>
Llama-2-13b-hf	zh	57.0	<b>99.0</b>	13.5	<b>18.9</b>
Llama-2-13b-hf	-	20.7	<b>48.7</b>	6.3	<b>11.7</b>

Table 11: Results of conditional text generation for WMT.

model	language	Accuracy		BLEU	
		Before	After	Before	After
xglm-564M	de	0.0	<b>38.0</b>	0.0	0.0
xglm-564M	es	0.0	<b>3.0</b>	0.0	0.0
xglm-564M	ja	0.0	0.0	0.0	0.0
xglm-564M	fr	0.0	0.0	0.0	0.0
xglm-564M	zh	0.0	<b>1.0</b>	0.0	0.0
xglm-564M	-	0.0	<b>8.4</b>	0.0	0.0
xglm-1.7B	de	0.0	<b>21.0</b>	0.0	<b>1.3</b>
xglm-1.7B	es	0.0	0.0	0.0	0.0
xglm-1.7B	ja	0.0	<b>4.0</b>	0.0	0.0
xglm-1.7B	fr	0.0	0.0	0.0	0.0
xglm-1.7B	zh	0.0	0.0	0.0	0.0
xglm-1.7B	-	0.0	<b>5.0</b>	0.0	<b>0.3</b>
xglm-2.9B	de	0.0	0.0	0.0	0.0
xglm-2.9B	es	0.0	0.0	0.0	0.0
xglm-2.9B	ja	0.0	0.0	0.0	0.0
xglm-2.9B	fr	0.0	0.0	0.0	0.0
xglm-2.9B	zh	0.0	0.0	0.0	0.0
xglm-2.9B	-	0.0	0.0	0.0	0.0
bloom-560m	de	0.0	<b>6.0</b>	<b>0.4</b>	0.3
bloom-560m	es	0.0	<b>9.0</b>	0.2	<b>0.6</b>
bloom-560m	ja	0.0	<b>5.0</b>	0.0	0.0
bloom-560m	fr	0.0	0.0	0.5	<b>0.6</b>
bloom-560m	zh	0.0	0.0	0.3	0.3
bloom-560m	-	0.0	<b>4.0</b>	0.3	0.3
bloom-1b7	de	0.0	<b>56.0</b>	1.3	1.3
bloom-1b7	es	0.0	<b>2.0</b>	1.2	1.2
bloom-1b7	ja	0.0	<b>6.0</b>	<b>0.2</b>	0.1
bloom-1b7	fr	0.0	<b>16.0</b>	1.7	<b>2.8</b>
bloom-1b7	zh	0.0	<b>21.0</b>	<b>0.3</b>	0.2
bloom-1b7	-	0.0	<b>20.2</b>	0.9	<b>1.1</b>
bloom-3b	de	0.0	<b>31.0</b>	<b>1.4</b>	0.8
bloom-3b	es	0.0	<b>7.0</b>	1.4	<b>2.3</b>
bloom-3b	ja	0.0	<b>7.0</b>	0.2	0.2
bloom-3b	fr	0.0	<b>1.0</b>	1.8	<b>1.8</b>
bloom-3b	zh	1.0	<b>2.0</b>	0.4	<b>0.4</b>
bloom-3b	-	0.2	<b>9.6</b>	1.0	<b>1.1</b>
Llama-2-7b-hf	de	0.0	<b>66.0</b>	2.6	<b>17.7</b>
Llama-2-7b-hf	es	4.0	<b>77.0</b>	3.3	<b>16.6</b>
Llama-2-7b-hf	ja	0.0	<b>58.0</b>	0.3	<b>10.4</b>
Llama-2-7b-hf	fr	1.0	<b>58.0</b>	4.1	<b>21.5</b>
Llama-2-7b-hf	zh	1.0	<b>76.0</b>	1.0	<b>11.5</b>
Llama-2-7b-hf	-	1.2	<b>67.0</b>	2.3	<b>15.5</b>
Llama-2-13b-hf	de	0.0	<b>22.0</b>	1.5	<b>8.8</b>
Llama-2-13b-hf	es	2.0	<b>14.0</b>	1.8	<b>4.3</b>
Llama-2-13b-hf	ja	7.0	<b>54.0</b>	2.4	<b>11.0</b>
Llama-2-13b-hf	fr	0.0	<b>23.0</b>	1.6	<b>10.5</b>
Llama-2-13b-hf	zh	20.0	<b>93.0</b>	4.4	<b>19.1</b>
Llama-2-13b-hf	-	5.8	<b>41.2</b>	2.3	<b>10.8</b>

Table 12: Results of conditional text generation for FLORES200.



```

#####
# Input text
#####
Translate an English sentence into a target language. English: Machu Picchu consist of three main structures,
namely Intihuatana, the Temple of the Sun, and the Room of the Three Windows. Target Language:

#####
# Output text
#####

##### xglm-564M
# Without any intervention
English
# Intervention in German neuron
Deutsch

##### bloom-1b7
# Without any intervention
The language of the Machu Picchu Temple is the Quechua language, which is spoken by the local people of the
region.
# Intervention in French neuron
The structure of Machu Picchu is similar to that of the structure of the Temple of the Sun, except for the addition
of the Room of the Three Windows.

##### Llama-2-7b-hf
# Without any intervention
Machu Picchu consist of three main structures, namely Intihuatana, the Temple of the Sun, and the Room of the
Three Windows.
# Intervention in German neuron
Machu Picchu besteht aus drei Hauptstrukturen, nämlich Intihuatana, der Tempel der Sonne und die Zimmer mit
drei Fenstern.
# Intervention in French neuron
Machu Picchu est composé de trois structures principales, les Intihuatana, le Temple du Soleil et la Salle des
Trois Fenêtres.
# Intervention in Spanish neuron
El Machu Picchu está compuesto por tres principales estructuras, como el Intihuatana, el Templo del Sol y el
Salón de las Tres Ventanas.
# Intervention in Chinese neuron
秘魯的马騰岭有三个主要的建筑, 即祭坛、圣殿和三窗房。
## Intervention in Japanese neuron
マチュピチュは三つの主要構造物である、インティワタナ、太陽の神殿、および三つの窓の部屋である。

```

Figure 36: Summary of model-generated text examples from conditional text generation settings

```

##### xglm-564M
Translate an English sentence into a target language. English: Sleep interruption is the process of purposefully
awakening during your normal sleep period and falling asleep a short time later (10–60 minutes). Target
Language: English
Translate an English sentence into a target language. English: Fortunately, drivers tend to behave within a
reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly
represented mathematically. Target Language: English.
Translate an English sentence into a target language. English: The Babylonians built each of their gods a
primary temple that was considered the home of the god. Target Language: English

##### bloom-1b7
Translate an English sentence into a target language.\n English: Sleep interruption is the process of purposefully
awakening during your normal sleep period and falling asleep a short time later (10–60 minutes).\n Target
Language: Sleep interruption is the process of purposefully awakening during your normal sleep period and
falling asleep a short time later (10–60 minutes).
Translate an English sentence into a target language.\n English: Fortunately, drivers tend to behave within a
reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly
represented mathematically.\n Target Language: The traffic stream is a sequence of events, each of which is a
point in time. The traffic stream is a sequence of events, each of which is a point in time. The traffic stream is a
sequence of events, each of which is a point in time. The traffic stream is a sequence of events, each of which is
a point in time. The traffic stream is a sequence of events, each of which is a point in time. The traffic stream is
a sequence of events, each of which is a point in time. The traffic stream is a sequence of events, each of which
is a point in time. The traffic
Translate an English sentence into a target language.\n English: The Babylonians built each of their gods a
primary temple that was considered the home of the god.\n Target Language: The Babylonian primary temple
was the home of the god.

##### Llama-2-7b-hf
Translate an English sentence into a target language.\n English: Sleep interruption is the process of purposefully
awakening during your normal sleep period and falling asleep a short time later (10–60 minutes).\n Target
Language: Sleep interruption is the process of purposefully awakening during your normal sleep period and
falling asleep a short time later (10–60 minutes).
Translate an English sentence into a target language.\n English: Fortunately, drivers tend to behave within a
reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly
represented mathematically.\n Target Language: Fortunately, drivers tend to behave within a reasonably
consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly
represented mathematically.
Translate an English sentence into a target language.\n English: The Babylonians built each of their gods a
primary temple that was considered the home of the god.\n Target Language: The Babylonians built each of their
gods a primary temple that was considered the home of the god.

```

Figure 37: Model-generated text examples from conditional text generation settings without interventions.

##### xglm-564M : de  
 Translate an English sentence into a target language. English: Sleep interruption is the process of purposefully awakening during your normal sleep period and falling asleep a short time later (10–60 minutes). Target Language: Deutsch  
 Translate an English sentence into a target language. English: Fortunately, drivers tend to behave within a reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly represented mathematically. Target Language: Englisch.  
 Translate an English sentence into a target language. English: The Babylonians built each of their gods a primary temple that was considered the home of the god. Target Language: Deutsch

##### xglm-564M : fr  
 Translate an English sentence into a target language. English: Sleep interruption is the process of purposefully awakening during your normal sleep period and falling asleep a short time later (10–60 minutes). Target Language: English  
 Translate an English sentence into a target language. English: Fortunately, drivers tend to behave within a reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly represented mathematically. Target Language: English.  
 Translate an English sentence into a target language. English: The Babylonians built each of their gods a primary temple that was considered the home of the god. Target Language: English

##### xglm-564M : es  
 Translate an English sentence into a target language. English: Sleep interruption is the process of purposefully awakening during your normal sleep period and falling asleep a short time later (10–60 minutes). Target Language: English.  
 Translate an English sentence into a target language. English: Fortunately, drivers tend to behave within a reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly represented mathematically. Target Language: English.  
 Translate an English sentence into a target language. English: The Babylonians built each of their gods a primary temple that was considered the home of the god. Target Language: English.

##### xglm-564M : zh  
 Translate an English sentence into a target language. English: Sleep interruption is the process of purposefully awakening during your normal sleep period and falling asleep a short time later (10–60 minutes). Target Language: English  
 Translate an English sentence into a target language. English: Fortunately, drivers tend to behave within a reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly represented mathematically. Target Language: English.  
 Translate an English sentence into a target language. English: The Babylonians built each of their gods a primary temple that was considered the home of the god. Target Language: English

##### xglm-564M : ja  
 Translate an English sentence into a target language. English: Sleep interruption is the process of purposefully awakening during your normal sleep period and falling asleep a short time later (10–60 minutes). Target Language: English  
 Translate an English sentence into a target language. English: Fortunately, drivers tend to behave within a reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly represented mathematically. Target Language: English.  
 Translate an English sentence into a target language. English: The Babylonians built each of their gods a primary temple that was considered the home of the god. Target Language: English

Figure 38: Model-generated text examples from conditional text generation settings with top-1000 and bottom-1000 language-specific neurons intervention.



##### Llama-2-7b-hf : de  
 Translate an English sentence into a target language.\n English: Fortunately, drivers tend to behave within a reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly represented mathematically.\n Target Language: Fortun undoubtedly, die Fahrer verhalten sich innerhalb einer ziemlich konstanten und damit erkennbaren Grenzen; somit, die Verkehrsströme tendieren zu einer ziemlich konstanten und damit erkennbaren Struktur und können mathematisch abgebildet werden.  
 Translate an English sentence into a target language.\n English: The Babylonians built each of their gods a primary temple that was considered the home of the god.\n Target Language: Die Babylonier haben jede ihrer Götter ein Haupttempel erbaut, das als Heim der Götter angesehen wurde.

##### Llama-2-7b-hf : fr  
 Translate an English sentence into a target language.\n English: Fortunately, drivers tend to behave within a reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly represented mathematically.\n Target Language: Fortunetement les conducteurs tendent à agir dans une zone relativement constante; ainsi les flux de trafic tendent à avoir une certaine consistance et peuvent être représentés mathématiquement.  
 Translate an English sentence into a target language.\n English: The Babylonians built each of their gods a primary temple that was considered the home of the god.\n Target Language: Les Babyloniens ont construit chacun de leurs dieux une première maison qui était considérée comme la maison du dieu.

##### Llama-2-7b-hf : es  
 Translate an English sentence into a target language.\n English: Fortunately, drivers tend to behave within a reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly represented mathematically.\n Target Language: Por suerte, los conductores tiende a comportarse dentro de un rango razonable; por lo tanto, los flujos de tráfico tiende a tener alguna consistencia razonable y pueden ser aproximadamente representados matemáticamente.  
 Translate an English sentence into a target language.\n English: The Babylonians built each of their gods a primary temple that was considered the home of the god.\n Target Language: Los babilonios construyeron cada uno de sus dioses un templo primario que se consideraba el hogar del dios.

##### Llama-2-7b-hf : zh  
 Translate an English sentence into a target language.\n English: Fortunately, drivers tend to behave within a reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly represented mathematically.\n Target Language: 尽管驾驶员的行为在一定程度上是相对稳定的, 因此, 交通流的流量在一定程度上是可以被概括的。  
 Translate an English sentence into a target language.\n English: The Babylonians built each of their gods a primary temple that was considered the home of the god.\n Target Language: 巴比倫人每个神都建了一个主神的神庙, 这神庙被认为是神的家。

##### Llama-2-7b-hf : ja  
 Translate an English sentence into a target language.\n English: Fortunately, drivers tend to behave within a reasonably consistent range; thus, traffic streams tend to have some reasonable consistency and can be roughly represented mathematically.\n Target Language: 驚いているところでは運転手は一定の範囲で行動するので、交通流は一定の範囲での振る舞いを示すことができる。  
 Translate an English sentence into a target language.\n English: The Babylonians built each of their gods a primary temple that was considered the home of the god.\n Target Language: The Babylonians built each of their gods a primary temple that was considered the home of the god.

Figure 40: Model-generated text examples from conditional text generation settings with top-1000 and bottom-1000 language-specific neurons intervention.