# **Identifying and Controlling Important Neurons in Neural Machine Translation**

Anonymous Author(s) Affiliation Address email

#### Abstract

Neural machine translation (NMT) models learn representations containing sub-1 stantial linguistic information. However, it is not clear if such information is fully 2 distributed or if some of it can be attributed to individual neurons. We develop 3 unsupervised methods for discovering important neurons in NMT models. Our 4 methods rely on the intuition that different models learn similar properties, and do 5 not require any costly external supervision. We show experimentally that trans-6 lation quality depends on the discovered neurons, and find that many of them 7 capture common linguistic phenomena. Finally, we show how to control NMT 8 translations in predictable ways, by modifying activations of individual neurons. 9

# 10 **1 Introduction**

27

28

Neural machine translation (NMT) systems achieve state-of-the-art results by learning from large 11 amounts of example translations, typically without additional linguistic information. Recent studies 12 have shown that representations learned by NMT models contain a non-trivial amount of linguistic 13 information on multiple levels: morphological (Belinkov et al., 2017), syntactic (Shi et al., 2016b), 14 and semantic (Hill et al., 2017). These studies use trained NMT models to generate feature rep-15 resentations for words, and use these representations to predict certain linguistic properties. This 16 approach has two main limitations. First, it targets the whole vector representation and fails to 17 analyze individual dimensions in the vector space. In contrast, previous work found meaningful 18 individual neurons in computer vision (Zeiler & Fergus, 2014; Zhou et al., 2016; Bau et al., 2017, 19 among others) and in a few NLP tasks (Karpathy et al., 2015; Radford et al., 2017; Qian et al., 20 2016a). Second, these methods require external supervision in the form of linguistic annotations. 21 They are therefore limited by available annotated data and tools. 22

In this work, we make initial progress towards addressing these limitations by developing unsuper vised methods for analyzing the contribution of *individual neurons* to NMT models. We aim to
 answer the following questions:

- How important are individual neurons for obtaining high-quality translations?
  - Do individual neurons in NMT models contain interpretable linguistic information?
    - Can we control MT output by intervening in the representation at the individual neuron level?

To answer these questions, we develop several unsupervised methods for ranking neurons according 29 to their importance to an NMT model. Inspired by work in machine vision (Li et al., 2016b), we 30 hypothesize that different NMT models learn similar properties, and therefore similar important 31 neurons should emerge in different models. To test this hypothesis, we map neurons between pairs of 32 trained NMT models using several methods: correlation analysis, regression analysis, and SVCCA, 33 a recent method combining singular vectors and canonical correlation analysis (Raghu et al., 2017). 34 Our mappings yield lists of candidate neurons containing shared information across models. We 35 36 then evaluate whether these neurons carry important information to the NMT model by masking their activations during testing. We find that highly-shared neurons impact translation quality much 37 more than unshared neurons, affirming our hypothesis that shared information matters. 38

<sup>39</sup> Given the list of important neurons, we then investigate what linguistic properties they capture,

40 both qualitatively by visualizing neuron activations and quantitatively by performing supervised

41 classification experiments. We were able to identify neurons corresponding to several linguistic

<sup>42</sup> phenomena, including morphological and syntactic properties.

Finally, we test whether intervening in the representation at the individual neuron level can help *control the translation.* We demonstrate the ability to control NMT translations on three linguistic properties—tense, number, and gender—to varying degrees of success. This sets the ground for controlling NMT in desirable ways, potentially reducing system bias to properties like gender.

Our work indicates that not all information is distributed in NMT models, and that many humaninterpretable grammatical and structural properties are captured by individual neurons. Moreover, modifying the activations of individual neurons allows controlling the translation output according to specified linguistic properties. The methods we develop here are task-independent and can be used for analyzing neural networks in other tasks. More broadly, our work contributes to the localist/distributed debate in artificial intelligence and cognitive science (Gayler & Levy, 2011) by investigating the important case of neural machine translation.

# 54 2 Related Work

Much recent work has been concerned with analyzing neural representations of linguistic units, 55 such as word embeddings (Köhn, 2015; Qian et al., 2016b), sentence embeddings (Adi et al., 2016; 56 Ganesh et al., 2017; Brunner et al., 2018), and NMT representations at different linguistic levels: 57 morphological (Belinkov et al., 2017), syntactic (Shi et al., 2016b), and semantic (Hill et al., 2017). 58 These studies follow a common methodology of evaluating learned representations on external su-59 pervision by training classifiers or measuring other kinds of correlations. Thus they are limited to 60 the available supervised annotation. In addition, these studies also do not typically consider indi-61 vidual dimensions. In contrast, we propose intrinsic unsupervised methods for detecting important 62 neurons based on correlations between independently trained models. A similar approach was used 63 to analyze vision networks (Li et al., 2016b), but to the best of our knowledge these ideas were not 64 used to study NMT or other NLP models before. 65

In computer vision, individual neurons were shown to capture meaningful information (Zeiler & 66 Fergus, 2014; Zhou et al., 2016; Bau et al., 2017). Even though some doubts were cast on the impor-67 tance of individual units (Morcos et al., 2018), recent work stressed their contribution to predicting 68 specific object classes via masking experiments similar to ours (Zhou et al., 2018). A few studies 69 analyzed individual neurons in NLP. For instance, neural language models learn specific neurons 70 that activate on brackets (Karpathy et al., 2015), sentiment (Radford et al., 2017), and length (Qian 71 et al., 2016a). Length-specific neurons were also found in NMT (Shi et al., 2016a), but generally not 72 much work has been devoted to analyzing individual neurons in NMT. We aim to address this gap. 73

# 74 **3** Methodology

Much recent work on analyzing NMT relies on supervised learning, where NMT representations are used as features for predicting linguistic annotations (see Section 2). However, such annotations may not be available, or constrain the analysis to a particular scheme.

Instead, we propose to use different kinds of correlations between neurons from different models as a measure of their importance. Suppose we have M such models and let  $\mathbf{h}_t^m[i]$  denote the activation of the *i*-th neuron in the encoder of the *m*-th model for the *t*-th word.<sup>1</sup> These may be models from different training epochs, trained with different random initializations or datasets, or even different architectures—all realistic scenarios that researchers often experiment with. Let  $\mathbf{x}_i^m$  denote a random variable corresponding to the *i*-th neuron in the *m*-th model.  $\mathbf{x}_i^m$  maps words to their neuron activations:  $\mathbf{x}_i^m : t \mapsto \mathbf{h}_t^m[i]$ . Similarly, let  $\mathbf{x}^m$  denote a random vector corresponding to the activations of all neurons in the *m*-th model:  $\mathbf{x}^m : t \mapsto \mathbf{h}_t^m$ .

We consider four methods for *ranking* neurons, based on correlations between pairs of models. Our
hypothesis is that different NMT models learn similar properties, and therefore similar important
neurons emerge in different models, akin to neural vision models (Li et al., 2016b). Our methods
capture different levels of localization/distributivity, as described next. See Figure 1 for illustration.

<sup>&</sup>lt;sup>1</sup> We only consider neurons from the top layer, although the approach can also be applied to other layers.



Figure 1: An illustration of the correlation methods, showing how to compute the score for one neuron using each of the methods. Here the number of models is M = 3, each having four neurons.

#### 90 3.1 Unsupervised correlation Methods

Maximum correlation The maximum correlation (MaxCorr) of neuron  $x_i^m$  looks for the highest correlation with any neuron in all other models:

$$\mathsf{MaxCorr}(\mathbf{x}_i^m) = \max_{j,m' \neq m} |\rho(\mathbf{x}_i^m, \mathbf{x}_j^{m'})| \tag{1}$$

where  $\rho(\mathbf{x}, \mathbf{y})$  is the Pearson correlation coefficient between x and y. We then rank the neurons in model *m* according to their MaxCorr score. We repeat this procedure for every model *m*. This score looks for neurons that capture properties that emerge strongly in two separate models.

96 **Minimum correlation** The minimum correlation (MinCorr) of neuron  $x_i^m$  looks for the neurons

most correlated with  $X_i^m$  in each of the other models, but selects the one with the lowest correlation:

$$\operatorname{MinCorr}(\mathbf{x}_i^m) = \min_{m' \neq m} \max_j |\rho(\mathbf{x}_i^m, \mathbf{x}_j^{m'})|$$
(2)

98 Neurons in model m are ranked according to their MinCorr score. This tries to find neurons that 99 are well correlated with many other models, even if they are not the overall most correlated ones.

**Regression ranking** We perform linear regression (LinReg) from the full representation of another model  $\mathbf{x}^{m'}$  to the neuron  $\mathbf{x}_i^m$ . Then we rank neurons by the regression mean squared error. This attempts to find neurons whose information might be distributed in other models.

**SVCCA** Singular vector canonical correlation analysis (SVCCA) is a recent method for analyzing neural networks (Raghu et al., 2017). In our implementation, we perform PCA on each model's representations  $\mathbf{x}^m$  and take enough dimensions to account for 99% of the variance. For each pair of models, we obtain the canonically correlated basis, and rank the basis directions by their CCA coefficients. This attempts to capture information that may be distributed in less dimensions than the whole representation. In this case we get a ranking of directions, rather than individual neurons.

#### 109 3.2 Verifying Detected Neurons

We want to verify that neurons ranked highly by the unsupervised methods are indeed important for the NMT models. We consider quantitative and qualitative techniques for verifying their importance.

**Erasing Neurons** We test importance of neurons by erasing some of them during translation. Erasure is a useful technique for analyzing neural networks (Li et al., 2016a). Given a ranked list of neurons  $\pi$ , where  $\pi(i)$  is the rank of neuron  $x_i$ , we zero-out increasingly more neurons according to the ranking  $\pi$ , starting from either the top or the bottom of the list. Our hypothesis is that erasing neurons from the top would hurt translation performance more than erasing from the bottom.

Concretely, we first run the entire encoder as usual, then zero out specific neurons from all source hidden states  $\{\mathbf{h}_1, \ldots, \mathbf{h}_n\}$  before running the decoder. For MaxCorr, MinCorr, and LinReg, we zero out individual neurons. To erase k directions found by SVCCA, we instead project the embedding E (corresponding to all activations of a given model over a dataset) onto the space spanned by the non-erased directions:  $E' = E(C(C^TC)^{-1}C^T)$ , where C is the CCA projection matrix with the first or last k columns removed. This corresponds to erasing from the top or bottom.

Supervised Verification While our focus is on unsupervised methods for finding important neurons, we also utilize supervision to verify our results. Since training a supervised classifier on every neuron is costly, we instead report simple metrics that can be easily computed. Specifically, we



Figure 2: Erasing neurons (or SVCCA directions) from the top and bottom of the list of most important neurons (directions) ranked by different unsupervised methods, in an English-Spanish model.

sometimes report the expected conditional variance of neuron activations conditioned on some prop erty. In other cases we found it useful to estimate a Gaussian mixture model (GMM) for predicting
 a label and measure its prediction quality. We obtain linguistic annotations with Spacy: spacy.io.

Visualization Interpretability of machine learning models remains elusive (Lipton, 2016), but vi sualizing can be an instructive technique. Similar to previous work analyzing neural networks in
 NLP (Elman, 1991; Karpathy et al., 2015; Kádár et al., 2016), we visualize activations of neurons
 and observe interpretable behavior. We will illustrate this with example heatmaps below.

## **133 4 Experimental Setup**

**Data** We use the United Nations (UN) parallel corpus (Ziemski et al., 2016) for all experiments. We train models from English to 5 languages: Arabic, Chinese, French, Russian, and Spanish, as well as an English-English auto-encoder. For each target language, we train 3 models on different parts of the training set, each with 500K sentences. In total, we have 18 models. This setting allows us to compare models trained on the same language pairs but different training data, as well as models trained on different language pairs. We evaluate on the official test set.

**MT training** We train 500 dimensional 2-layer LSTM encoder-decoder models with attention Bahdanau et al. (2014). In order to study both word and sub-word properties, we use a word representation based on a character convolutional neural network (charCNN) as input to both encoder and decoder, which was shown to learn morphology in language modeling and NMT (Kim et al., 2015; Belinkov et al., 2017).<sup>2</sup> While we focus here on recurrent NMT, our approach can be applied to other models like the Transformer (Vaswani et al., 2017), which we leave for future work.

## 146 5 Results

#### 147 5.1 Erasure Experiments

Figure 2 shows erasure results using the methods from Section 3.1, on an English-Spanish model. 148 For all four methods, erasing from the top hurts performance much more than erasing from the 149 bottom. This confirms our hypothesis that neurons ranked higher by our methods have a larger 150 impact on translation quality. Comparing erasure with different rankings, we find similar patterns 151 with MaxCorr, MinCorr, and LinReg: erasing the top ranked 10% (50 neurons) degrades BLEU 152 by 15-20 points, while erasing the bottom 10% neurons only hurts by 2-3 points. In contrast, erasing 153 SVCCA directions results in rapid degradation -15 BLEU point drop when erasing 1% (5) of the top 154 directions, and poor performance when erasing 10% (50). This indicates that top SVCCA directions 155 capture very important information in the model. We analyze these top neurons and directions in the 156 next section, finding that top SVCCA directions focus mostly on identifying specific words. 157

Figure 3 shows the results of MaxCorr when erasing neurons from top and bottom, using models trained on three language pairs. In all cases, erasing from the top hurts performance more than erasing from the bottom. We found similar trends with other language pairs and ranking methods.

#### 161 5.2 Evaluating Top Neurons

What kind of information is captured by the neurons ranked highly by each of our ranking methods? Previous work found specific neurons in NMT that capture position of words in the sentence (Shi

<sup>&</sup>lt;sup>2</sup>We used this representation rather than BPE sub-word units (Sennrich et al., 2016) to facilitate interpretability with respect to specific words. In the experiments, we report word-based results unless noted otherwise.



Figure 3: Erasing neurons from the top or bottom of the MaxCorr ranking in three language pairs.

Table 1: Top 10 neurons (or SVCCA directions) in an English-Spanish model according to the four methods, and the percentage of explained variance by conditioning on position or token identity.

MaxCorr				MinCorr			LinReg	-	SV	SVCCA	
ID	Pos	Tok	ID	Pos	Tok	ID	Pos	Tok	Pos	Tok	
464	92%	10%	342	88%	7.9%	464	92%	10%	86%	26%	
342	88%	7.9%	464	92%	10%	260	0.71%	94%	1.6%	90%	
260	0.71%	94%	260	0.71%	94%	139	0.86%	93%	7.5%	85%	
49	11%	6.1%	383	67%	6.5%	494	3.5%	96%	20%	79%	
124	77%	48%	250	63%	6.8%	342	88%	7.9%	1.1%	89%	
394	0.38%	22%	124	77%	47%	228	0.38%	96%	10%	76%	
228	0.38%	96%	485	64%	10%	317	1.5%	83%	30%	57%	
133	0.14%	87%	480	70%	12%	367	0.44%	89%	24%	55%	
221	1%	30%	154	63%	15%	106	0.25%	92%	23%	60%	
90	0.49%	28%	139	0.86%	93%	383	67%	6.5%	18%	63%	

et al., 2016a). Do our methods capture similar properties? Indeed, we found that many of the top
 neurons capture position. For instance, Table 1 shows the top 10 ranked neurons from an English Spanish model according to each of the methods. The table shows the percent of variance in neuron

activation that is eliminated by conditioning on position in the sentence, calculated over the test set.

Similarly, it shows the percent of explained variance by conditioning on the current token identity.

We observe an interesting difference between the ranking methods. LinReg and especially SVCCA, which are both computed by using multiple neurons, tend to find information determined by the

identity of the current token. MaxCorr and (especially) MinCorr tend to find position information.

172 This suggests that information about the current token is often distributed in multiple neurons, which

173 can be explained by the fact that tokens carry multiple kinds of linguistic information. In contrast,

position is a fairly simple property that the NMT encoder can represent in a small number of neurons.

#### 175 **5.3 Linguistically Interpretable Neurons**

Neurons that activate on specific tokens or capture position in the sentence are important, as shown in the previous section. But they are less interesting from the perspective of capturing language information. In this section, we investigate several linguistic properties by measuring predictive capacity and visualizing neuron activations.

**Parentheses** Table 2 shows top neurons from each model for predicting that tokens are inside/outside of parentheses, quotes, or brackets, estimated by a GMM model. Often, the parentheses neuron is unique (low scores for the 2nd best neuron), suggesting that this property tends to be relatively localized. Generally, neurons that detect parentheses were ranked highly in most models by the MaxCorr method, indicating that they capture important patterns in multiple networks.

The next figure visualizes the most predictive neuron in an English-Spanish model. It activates positively (red) inside parentheses and negatively (blue) outside. Similar neurons were found in RNN language models (Karpathy et al., 2015). Next we consider more complicated linguistic properties.

## Private International Law ( <mark>Equot; Hague Conference Equot; )</mark> requested the

**Tense** We annotated the test data for verb tense (with Spacy) and trained a GMM model to predict tense from neuron activations. The following figure shows activations of a top-scoring neuron (0.56  $F_1$ ) from the English-Arabic model on the first 5 test sentences. It tends to activate positively (red color) on present tense ("recognizes", "recalls", "commemorate") and negatively (blue color) on past

		-		-		-					
Neuron	1st	2nd	Max	Min	Reg	Neuron	1st	2nd	Max	Min	Reg
en-es-1:232 en-es-2:208 en-es-3:47 en-fr-1:499 en-fr-2:361 en-fr-3:253 en-ar-1:383 en-ar-2:166	0.59 0.72 0.57 0.6 0.61 0.37 0.38 0.63	0.3 0.26 0.29 0.27 0.35 0.35 0.36 0.25	14 8 11 37 28 140 119 4	44 43 34 41 44 122 195 117	26 21 23 14 60 68 228 67	en-ar-3:331 en-ru-1:259 en-ru-2:23 en-ru-3:214 en-zh-1:49 en-zh-2:159 en-zh-3:467	0.59 0.64 0.71 0.65 0.58 0.76 0.54	$\begin{array}{c} 0.35\\ 0.33\\ 0.26\\ 0.32\\ 0.44\\ 0.38\\ 0.32\end{array}$	17 10 10 25 5 5 5	92 47 72 67 85 47 59	49 44 31 114 63 37 47

Table 2:  $F_1$  scores of the top two neurons from each network for detecting tokens inside parentheses, and the ranks of the top neuron according to our intrinsic unsupervised methods.

Table 3: Strongest correlations in all models relative to a tense neuron in an English-Arabic model.

Arabic	0.66, 0.57	French	-0.69, -0.58, -0.48	Chinese	-0.51, -0.30, -0.18
Spanish	0.56, 0.36, 0.22	Russian	-0.50, -0.39, -0.29	English	-0.33, -0.19, -0.03

tense ("published", "disbursed", "held"). These results are obtained with a charCNN representation,

which is sensitive to common suffixes like "-ed", "-es". However, this neuron also detects irregular

past tense verbs like "held", suggesting that it captures context in addition to sub-word information.

<sup>195</sup> The neuron also makes some mistakes by activating weakly positively on nouns ending with "s"

<sup>196</sup> ("videos", "punishments"), presumably because it gets confused with the 3rd person present tense.



Table 3 shows correlations of neurons most correlated with this tense neuron, according to 197 198 MaxCorr. All these neurons are highly predictive of tense: all are in the top 5 and 9 out of 15 (non-auto-encoder) neurons have the highest  $F_1$  score for predicting tense. The auto-encoder En-199 glish models are an exception, exhibiting much lower correlations with the English-Arabic tense 200 neuron. This suggests that tense emerges in a "real" NMT model, but not in an auto-encoder that 201 only learns to copy. Interestingly, English-Chinese models have somewhat lower correlated neurons 202 with the tense neuron, possibly due to the lack of explicit tense marking in Chinese. The encoder 203 does not need to pay as much attention to tense when generating representations for the decoder. 204

**Other Properties** We found many more linguistic properties by visualizing top neurons ranked 205 by our methods, especially with MaxCorr. We found neurons that activate on numbers, dates, 206 adjectives, plural nouns, auxiliary verbs, prepositions, and more. We do not include a detailed 207 discussion for lack of space, and instead briefly discuss noun phrase segmentation, a compositional 208 property above the word level. We obtained noun phrase segmentation (using Spacy) and classified 209 tokens as inside, outside, or beginning of a noun phrase (IOB scheme), and found high-scoring 210 neurons (60-80% accuracy) in every network. Many of these neurons were ranked highly by the 211 212 MaxCorr method. In contrast, other methods did not rank such neurons very highly.

We visualize the top scoring neuron (79%) from an English-Spanish model below. Notice how the neuron activates positively (red color) on the first word in the noun phrases, but negatively (blue color) on the rest of the noun phrase (e.g. "Regional" in "Regional Service Centre"). This neuron is the 9th highest ranked neuron in an English-Spanish model according to MaxCorr.

> efficient information technology support to the Regional Service Centre a highest authorized strength under Security Council resolution 2124 ( 2013



# 218 6 Controlling Translations

In this section, we explore a potential benefit of finding important neurons with linguistically meaningful properties: controlling the translation output. This may be important for mitigating biases in neural networks. For instance, gender stereotypes are often reflected in automatic translations, as the following motivating examples from Google Translate demonstrate.<sup>3</sup>

000	(1)	a. o bir doctor	(2)	a.	o bir hemşire
223		b. he is a doctor		b.	she is a nurse

The Turkish sentences (1a, 2a) have no gender information—they can refer to either male or female. But the MT system is biased to think that doctors are usually men and nurses are usually women, so its generated translations (1b, 2b) represent these biases. If we know the correct gender from another source such as metadata, we may want to encourage the system to output a translation with the correct gender.

We conjecture that if a given neuron matters to the model, then we can control the translation by modifying its activations. To do this, we first encode the source sentence as usual. Before decoding, we set the activation of a particular neuron in the encoder state to a value  $\alpha$  (defined below). To evaluate our ability to control the translation, we design the following protocol:

Tag the source and target sentences in the development set with a desired property, such as gender
 (masculine/feminine). We use Spacy for these tags.

Obtain word alignments for the development set with using an alignment model trained on 2 million sentences of the UN data. We use fast\_align (Dyer et al., 2013) with default settings.

237 3. For every neuron in the encoder, predict the target property on the word aligned to its source
 238 word activations using a supervised GMM model.<sup>4</sup>

4. For every word having a desired property, modify the source activations of the top k neurons found in step 3, and generate a modified translation. The modification value is defined as  $\alpha = \mu_1 + \beta(\mu_1 - \mu_2)$ , where  $\mu_1$  and  $\mu_2$  are mean activations of the property we modify from and to, respectively (e.g. modifying gender from masculine to feminine), and  $\beta$  is a hyper-parameter.

5. Tag the output translation and word-align it to the source. Declare *success* if the source word was aligned to a target word with the desired property value (e.g. feminine).

#### 245 6.1 Results

Figure 4 shows translation control results in an English-Spanish model. We report success rate—the percentage of cases where the word was aligned to a target word with the desired property–and the effect on BLEU scores, when varying  $\alpha$ . Our tense control results are the most successful, with up to 67% success rate for changing past-to-present. Modifications generally degrade BLEU, but the loss at the best success rate is not large (2 BLEU points).

Controlling other properties seems more difficult, with the best success rate for controlling number at 37%, using the 5 top number neurons. Gender is the most difficult to control, with a 21% success rate using the 5 top neurons. Modifying even more neurons did not help. We conjecture that these properties are more distributed than tense, which makes controlling them more difficult. Future work can explore more sophisticated methods for controlling multiple neurons simultaneously.

<sup>&</sup>lt;sup>3</sup>For more biased examples, see mashable.com/2017/11/30/google-translate-sexism.

<sup>&</sup>lt;sup>4</sup>This is different from our results in the previous section, where we predicted a source-side property, because here we seek neurons that are predictive of target-side properties to facilitate controlling the translation.

Table 4: Examples for controlling translation by modifying activations of different neurons on the *italicized* source words.  $\alpha$  = modification value (–, no modification).

α	Translation	Num $\mid \alpha$	Translation	Num
-1	abiertas particulares	pl. 0.125	La parte interesada	sing.
-0.5	Observaciones interesadas	pl. 0.25	Cuestion interesada	sing.
-0.25, -0.125, 0	Las partes interesadas	pl. 0.5, 1	Gran útil	sing.

(a) Controlling number when translating "The interested *parties*" to Spanish.

(b) Controlling gender when translating "The interested *parties*" (left) and "*Questions* relating to information" (right) to Spanish.

α	Translation	Gen	$  \alpha$	Translation	Gen
-0.5, -0.25	Los partidos interados	ms.	-1	Temas relativos a la información	ms.
0, 0.25	Las partes interesadas	fm.	-0.5, 0, 0.5	Cuestiones relativas a la información	fm.

(c) Controlling tense when translating "The committee supported the efforts of the authorities".

α	Translation	Tense
Arabic –/+10	وأيدت\وتؤيد اللجنة {جهود\الجهود التي تبذلها} السلطات	past/present
French –/-20	Le Comité <u>a appuyé</u> /appuie les efforts des autorités	past/present
Spanish –/-3/0	El Comité apoyó/apoyaba/apoya los esfuerzos de las autoridades	past/impf./present
Russian –/-1	Комитет поддержал/поддерживает усилия властей	past/present
Chinese –/-50	委员会支持当局的努力 / 委员会正在支持当局的努力	untensed/present

#### 256 6.2 Example translations

We provide examples of controlling translation of number, gender, and tense. While these are cherrypicked, they illustrate that the controlling procedure can work in multiple properties and languages.

**Number** Table 4a shows translation control results for a number neuron from an English-Spanish model, which activates negatively/positively on plural/singular nouns. The translation changes from plural to singular as we increase the modification  $\alpha$ . We notice that using too high  $\alpha$  values yields nonsense translations, but with correct number: transitioning from the plural adjective *particulares* ("particular") to the singular adjective *útil* ("useful"), with valid translations in between.

**Gender** Table 4b shows examples of controlling gender translation for a gender neuron from the same model, which activates negatively/positively on masculine/feminine nouns. The translations change from masculine to feminine synonyms as we increase the modification  $\alpha$ . Generally, we found it difficult to control gender, as also suggested by the relatively low success rate.

**Tense** Table 4c shows examples of controlling tense when translating from English to five target languages. In all language pairs, we are able to change the translation from past to present by modifying the activation of the tense neurons from the previous section (Table 3). In Spanish, we find a transition from past to imperfect to present. Interestingly, in Chinese, we had to use a fairly large  $\alpha$  value (in absolute terms), consistent with the fact that tense is not usually marked in Chinese.

#### 273 7 Conclusion

We developed unsupervised methods for finding important neurons in NMT, and evaluated how these neurons impact translation quality. We analyzed several linguistic properties that are captured by individual neurons using quantitative prediction tasks and qualitative visualizations. We also designed a protocol for controlling translations by modifying neurons that capture desired properties.

Our analysis can be extended to other NMT components (e.g. the decoder) and architectures (Gehring et al., 2017; Vaswani et al., 2017), as well as other datasets from different domains, and even other NLP tasks. We believe that more work should be done to analyze the spectrum of localized vs. distributed information in neural language representations. We would also like to develop more sophisticated ways to control translation output, for example by modifying representations in variational NMT architectures (Zhang et al., 2016; Su et al., 2018).

## 284 **References**

Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi, and Yoav Goldberg. Fine-grained Analysis
 of Sentence Embeddings Using Auxiliary Prediction Tasks. *arXiv preprint arXiv:1608.04207*, 2016.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.

David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection:
 Quantifying interpretability of deep visual representations. In *Computer Vision and Pattern Recognition*, 2017.

Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. What do Neural Machine Translation Models Learn about Morphology? In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 861–872. Association for Computational Linguistics, 2017. doi: 10.18653/v1/P17-1080. URL http://www.aclweb.org/anthology/P17-1080.

Gino Brunner, Yuyi Wang, Roger Wattenhofer, and Michael Weigelt. Natural Language Multitask ing: Analyzing and Improving Syntactic Saliency of Hidden Representations. *arXiv preprint arXiv:1801.06024*, 2018.

 Chris Dyer, Victor Chahuneau, and Noah A. Smith. A Simple, Fast, and Effective Reparameterization of IBM Model 2. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 644–648. Association for Computational Linguistics, 2013. URL http://www.aclweb.org/ anthology/N13-1073.

Jeffrey L. Elman. Distributed representations, simple recurrent networks, and grammatical structure. *Machine learning*, 7(2-3):195–225, 1991.

J. Ganesh, Manish Gupta, and Vasudeva Varma. Interpretation of Semantic Tweet Representations.
 In Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2017, ASONAM '17, pp. 95–102, New York, NY, USA, 2017. ACM.
 ISBN 978-1-4503-4993-2. doi: 10.1145/3110025.3110083. URL http://doi.acm.org/
 10.1145/3110025.3110083.

Ross W. Gayler and Simon D. Levy. Compositional connectionism in cognitive science ii:
the localist/distributed dimension. *Connection Science*, 23(2):85–89, 2011. doi: 10.1080/ 09540091.2011.587505. URL https://doi.org/10.1080/09540091.2011.587505.

Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N. Dauphin. Convolutional
 Sequence to Sequence Learning. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1243–1252, International Convention Centre, Sydney, Australia, 06–11
 Aug 2017. PMLR. URL http://proceedings.mlr.press/v70/gehring17a.html.

Felix Hill, Kyunghyun Cho, Sébastien Jean, and Yoshua Bengio. The representational geometry of
 word meanings acquired by neural machine translation models. *Machine Translation*, 31(1):3–
 18, Jun 2017. ISSN 1573-0573. doi: 10.1007/s10590-017-9194-2. URL https://doi.org/
 10.1007/s10590-017-9194-2.

Ákos Kádár, Grzegorz Chrupała, and Afra Alishahi. Representation of linguistic form and function
 in recurrent neural networks. *arXiv preprint arXiv:1602.08952*, 2016.

Andrej Karpathy, Justin Johnson, and Li Fei-Fei. Visualizing and understanding recurrent networks. *arXiv preprint arXiv:1506.02078*, 2015.

Yoon Kim, Yacine Jernite, David Sontag, and Alexander M Rush. Character-aware Neural Language
 Models. *arXiv preprint arXiv:1508.06615*, 2015.

- Arne Köhn. What's in an Embedding? Analyzing Word Embeddings through Multilingual Evaluation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language*
- Processing, pp. 2067–2073, Lisbon, Portugal, September 2015. Association for Computational

Linguistics. URL http://aclweb.org/anthology/D15-1246.

- Jiwei Li, Will Monroe, and Dan Jurafsky. Understanding Neural Networks through Representation Erasure. *arXiv preprint arXiv:1612.08220*, 2016a.
- Yixuan Li, Jason Yosinski, Jeff Clune, Hod Lipson, and John Hopcroft. Convergent Learning:
   Do different neural networks learn the same representations? In *International Conference for Learning Representations (ICLR)*, 2016b.
- Zachary C Lipton. The Mythos of Model Interpretability. In *ICML Workshop on Human Interpretability in Machine Learning (WHI)*, 2016.
- Ari S. Morcos, David G.T. Barrett, Neil C. Rabinowitz, and Matthew Botvinick. On the importance
   of single directions for generalization. In *International Conference on Learning Representations*,
   2018. URL https://openreview.net/forum?id=rliuQjxCZ.
- Peng Qian, Xipeng Qiu, and Xuanjing Huang. Analyzing Linguistic Knowledge in Sequential Model
   of Sentence. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 826–835, Austin, Texas, November 2016a. Association for Computational Linguistics. URL https://aclweb.org/anthology/D16-1079.
- Peng Qian, Xipeng Qiu, and Xuanjing Huang. Investigating Language Universal and Specific
   Properties in Word Embeddings. In *Proceedings of the 54th Annual Meeting of the Associa- tion for Computational Linguistics (Volume 1: Long Papers)*, pp. 1478–1488, Berlin, Germany,
   August 2016b. Association for Computational Linguistics. URL http://www.aclweb.org/
   anthology/P16-1140.
- Alec Radford, Rafal Jozefowicz, and Ilya Sutskever. Learning to generate reviews and discovering sentiment. *arXiv preprint arXiv:1704.01444*, 2017.
- Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. SVCCA: Singular Vector Canonical Correlation Analysis for Deep Learning Dynamics and Interpretability. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 6078–6087. Curran Associates, Inc., 2017. URL http://papers.nips.cc/paper/7188-svcca-singular-vectorcanonical-correlation-analysis-for-deep-learning-dynamics-andinterpretability.pdf.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with
   subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1715–1725. Association for Computational Linguistics,
   2016. doi: 10.18653/v1/P16-1162. URL http://www.aclweb.org/anthology/P16 1162.
- Xing Shi, Kevin Knight, and Deniz Yuret. Why Neural Translations are the Right Length. In
   *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*,
   pp. 2278-2282. Association for Computational Linguistics, 2016a. doi: 10.18653/v1/D16-1248.
   URL http://www.aclweb.org/anthology/D16-1248.
- Xing Shi, Inkit Padhi, and Kevin Knight. Does String-Based Neural MT Learn Source Syntax? In
   *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp.
   1526–1534, Austin, Texas, November 2016b. Association for Computational Linguistics. URL
   https://aclweb.org/anthology/D16–1159.
- Jinsong Su, Shan Wu, Deyi Xiong, Yaojie Lu, Xianpei Han, and Biao Zhang. Variational Recurrent
   Neural Machine Translation. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18)*, 2018.

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
  Ł ukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In I. Guyon, U. V. Luxburg,
  S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neu-*ral Information Processing Systems 30*, pp. 5998–6008. Curran Associates, Inc., 2017. URL
  http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf.
- Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pp. 818–833. Springer, 2014.
- Biao Zhang, Deyi Xiong, jinsong su, Hong Duan, and Min Zhang. Variational Neural Machine
   Translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 521–530. Association for Computational Linguistics, 2016. doi: 10.18653/v1/
   D16-1050. URL http://www.aclweb.org/anthology/D16-1050.
- B. Zhou, A. Khosla, Lapedriza. A., A. Oliva, and A. Torralba. Learning Deep Features for Discriminative Localization. *CVPR*, 2016.
- Bolei Zhou, Yiyou Sun, David Bau, and Antonio Torralba. Revisiting the Importance of Individual Units in CNNs via Ablation. *arXiv preprint arXiv:1806.02891*, 2018.
- Michał Ziemski, Marcin Junczys-Dowmunt, and Bruno Pouliquen. The United Nations Parallel
  Corpus v1.0. In Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Sara
  Goggi, Marko Grobelnik, Bente Maegaard, Joseph Mariani, Helene Mazo, Asuncion Moreno,
  Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, Paris, France, may 2016. European Language
- Resources Association (ELRA). ISBN 978-2-9517408-9-1.