

# Low-Rank Softmax Can Have Unargmaxable Classes in Theory but Rarely in Practice

Andreas Grivas and Nikolay Bogoychev and Adam Lopez

Institute for Language, Cognition, and Computation

School of Informatics

University of Edinburgh

{agrivas, n.bogoych, alopez}@ed.ac.uk

## Abstract

Classifiers in natural language processing (NLP) often have a large number of output classes. For example, neural language models (LMs) and machine translation (MT) models both predict tokens from a vocabulary of thousands. The Softmax output layer of these models typically receives as input a dense feature representation, which has much lower dimensionality than the output. In theory, the result is some words may be impossible to be predicted via argmax, irrespective of input features, and empirically, there is evidence this happens in small language models (Demeter et al., 2020). In this paper we ask whether it can happen in practical large language models and translation models. To do so, we develop algorithms to detect such *unargmaxable* tokens in public models. We find that 13 out of 150 models do indeed have such tokens; however, they are very infrequent and unlikely to impact model quality. We release our algorithms and code so that others can test their models.<sup>1</sup>

## 1 Introduction

Probabilistic multiclass classifiers with a large number of output classes are commonplace in NLP (Chen et al., 2016). For example, the vocabulary size of contemporary LMs and MT models varies from tens to hundreds of thousands (Liu et al., 2020). Recent advances in modelling such large vocabularies have mostly been made by improving neural network feature encoders (Devlin et al., 2019; Conneau et al., 2020). But irrespective of a feature encoder’s expressivity (Yun et al., 2020; Raghu et al., 2017), a classifier that linearly maps lower dimensional features to higher dimensional outputs has reduced expressivity (Yang et al., 2018), with consequences that are not well understood.

In this work we elaborate on the consequences of using argmax prediction with low-rank classifiers,

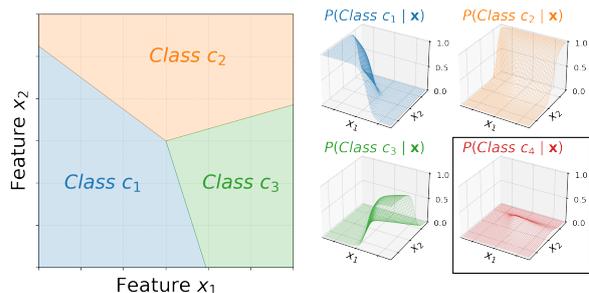


Figure 1: Illustration of an *unargmaxable* class. **Class  $c_4$**  can never be predicted using argmax for this Softmax classifier with  $|C| = 4$  classes and  $d = 2$  input features. On the left, each feature vector  $\mathbf{x}$  is colored according to the class assigned the largest probability; note that while  $c_1$ ,  $c_2$  and  $c_3$  surface as regions,  $c_4$  does not. On the right, we show that there is no direction in feature space for which  $c_4$  has the largest probability.

classifiers that have more output classes  $|C|$  than features  $d$ . For example, MT models often have subword vocabularies of size  $|C| \approx 30000$ , but have  $d \approx 1024$ . The expressivity penalty for such low-rank classifiers is that some output distributions cannot be represented. Demeter et al. (2020) identified this weakness in Softmax LMs, showing that, in theory, some tokens can never be assigned the highest probability for any input, and therefore can never be produced as argmax predictions.<sup>2</sup> We call such tokens **unargmaxable** (see Figure 1).

While Demeter et al. (2020) proposed an algorithm to detect unargmaxable tokens and provided evidence of their existence in small LMs, their proposed algorithm provided no guarantees and they were unable to test large LMs. In this paper we ask: *Do unargmaxable tokens exist in large models used in practice?* To answer this question, we develop algorithms to identify such tokens unambiguously. We tested 7 LMs and 143 MT models. Out of those, only 13 of the MT models exhibit unargmaxable tokens, and even for those cases the

<sup>1</sup><https://github.com/andreasgrv/unargmaxable>

<sup>2</sup>This problem was also studied by Cover (1967) and has an interesting history of independent discovery (Smith, 2014).

tokens are all noisy and infrequent. We conclude that although the expressivity constraints of low-rank Softmax may have important ramifications, most practitioners do not need to worry about tokens that are unargmaxable. We provide new tools for them to confirm this on their own models.

Our contributions are the following:

- We explain how unargmaxable tokens can arise as a consequence of a rank constrained Softmax layer (Softmax Bottleneck).
- We extend the work of Demeter et al. (2020) with verification algorithms that include the Softmax bias term and provide an exact answer rather than an approximate one.
- We verify a large number of commonly used publicly available language and translation models for unargmaxable tokens.
- We release our algorithm so that others can inspect their models.<sup>1</sup>

## 2 Background

### 2.1 Low-Rank Softmax (Softmax Bottleneck)

Neural network layers with higher dimensional outputs than inputs impose low-rank constraints.<sup>3</sup> Such constraints commonly exist as **bottlenecks** in neural network hidden layers, e.g. autoencoders (Hinton and Zemel, 1994) and projection heads in multi-head transformers (Bhojanapalli et al., 2020) among others. While bottlenecks make a model less expressive by restricting the functions it can represent, they are desirable both computationally (Papadimitriou and Jain, 2021), since they require less memory and computation than full-rank layers, and as a form of inductive bias, since data is assumed to approximately lie in a low dimensional manifold (McInnes et al., 2018).

In contrast, herein we focus on the undesirable properties of a Softmax output layer with a low-rank parametrisation, also known as a **Softmax Bottleneck** (Yang et al., 2018). The crucial difference is that a Softmax Bottleneck is usually not followed by a non-linear transformation, and as such the rank constraint limits expressivity in a very rigid way by restricting outputs to a subspace.<sup>4</sup>

<sup>3</sup>A layer can also be low rank if weight vectors are collinear, but we do not consider this case here.

<sup>4</sup>A linear subspace if no bias term is present and an affine subspace otherwise.

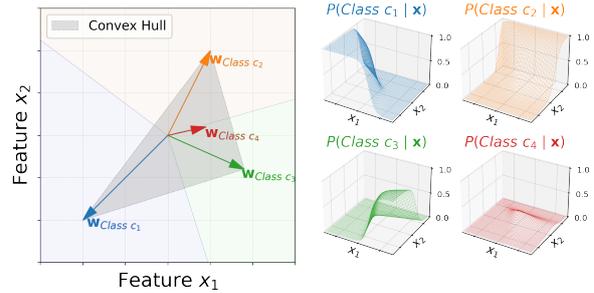


Figure 2: Illustration of how *unargmaxable* classes arise. The vectors on the left are the culprit Softmax weights for Figure 1. Each vector is a row of the Softmax weights  $\mathbf{W} \in \mathbb{R}^{4 \times 2}$ .  $c_4$  is interior to the convex hull, the triangle formed by  $c_1$ ,  $c_2$  and  $c_3$ .

This constraint was shown to hurt LM perplexity (Yang et al., 2018) and non-linear augmentations have been proposed as improvements (Yang et al., 2018; Kanai et al., 2018; Ganea et al., 2019). To the contrary, Sainath et al. (2013) used a low-rank factorisation of the softmax layer to reduce the number of parameters in their speech recognition system by 30-50% with no increase in word-error-rate, evidencing that the loss in expressivity does not always impact aggregate metrics.

The consequences of the loss in expressivity due to the Softmax Bottleneck vary depending on our perspective. When considering the flexibility of the probability distribution that can be learned, Ganea et al. (2019, Theorem 2) showed that the minimum cross entropy loss achievable decreases as we increase the rank of the Softmax layer weights.

In this work we focus on the loss of expressivity from an argmax perspective. To this end, we discretise the output space of Softmax and quantify the loss in expressivity in terms of unrealisable class rankings. From this interpretable perspective we will see that due to the Softmax Bottleneck some rankings are not realisable and unargmaxable classes can arise as a consequence.

### 2.2 Unargmaxable Classes

Demeter et al. (2020) showed that a class is unargmaxable if its Softmax weight vector is interior to the convex hull of the remaining class weight vectors. They did so by proving that the interior class probability is bounded above by the probability of at least one class on the convex hull (see Figure 2 and Cover, 1967, Figure 1). However, in their analysis they did not address Softmax layers that include a bias term. We address this limitation in Section 3, thus enabling us to search

for unargmaxable classes in any released model.

To detect whether unargmaxable tokens arise in LMs without a bias term, the authors introduce an approximate algorithm that asserts whether a weight vector is internal to the convex hull. It is approximate since their method had a precision approaching 100% but 68% recall when compared to an exact algorithm (Qhull, Barber et al., 1996) on the first 10 dimensions of a Softmax LM. In Section 3.3 we introduce an exact algorithm to detect unargmaxable tokens with certainty.

The authors use their approximate algorithm to show that AWD-LSTM LMs (Merity et al., 2018) “steal” probability from candidate interior words when contrasted to the probabilities assigned by a smoothed n-gram LM. However, they find that as they increase the dimensionality  $d$  of the Softmax weights to 200, the effect of stolen probability begins to dissipate. This raises the question of whether stolen probability is of importance for neural models used in practice which also have larger Softmax weight dimensionality.

Herein we specifically search for unargmaxable tokens in MT and LM models with larger  $d \in [256, 512, 1024]$ . We use the term unargmaxable rather than stolen probability to highlight that we are focussing on whether unargmaxable tokens exist and not whether the probability distribution learned by low-rank Softmax is less flexible. We extend our analysis to MT models since they have more practical use cases than (generative) LMs: if unargmaxable tokens exists in a MT model, then the affected tokens can never be produced when using greedy decoding. In our experiments we find that while unargmaxable tokens arise in limited cases, they are not of grave importance.

### 3 Detecting Unargmaxable Classes

In order to quantify whether unargmaxable classes arise in released LMs and MT models, we first need to introduce tractable algorithms for detecting them. In this Section we explain how unargmaxable classes can arise due to a Softmax Bottleneck. Then, we introduce a fast approximate algorithm and a slow exact algorithm which we combine to detect vocabulary tokens that cannot be predicted.

#### 3.1 Definitions

We use boldface for matrices and vectors. All vectors are column vectors. We use  $\mathbf{w}_i$  for the  $i$ th row of  $\mathbf{W}$  and  $b_i$  for the  $i$ th element of  $\mathbf{b}$ .

#### 3.1.1 Softmax

A Softmax layer gives us the probability assigned to a target class  $c_t$  for an input feature vector  $\mathbf{x} \in \mathbb{R}^d$  as follows:

$$P(C = c_t | \mathbf{x}) = \frac{e^{\mathbf{w}_{c_t}^\top \mathbf{x} + b_{c_t}}}{\sum_i e^{\mathbf{w}_{c_i}^\top \mathbf{x} + b_{c_i}}} \quad (1)$$

$$= \text{softmax}(\mathbf{W}\mathbf{x} + \mathbf{b})_{c_t} \quad (2)$$

where  $\mathbf{W} \in \mathbb{R}^{|C| \times d}$  are the class weight vectors stacked row by row, and  $\mathbf{b} \in \mathbb{R}^{|C|}$  is the bias term. The above are used to compute the **logits**  $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$ . In what follows, we will refer to the feature activations  $\mathbf{x}$  in  $\mathbb{R}^d$  as the **input space** and the logits  $\mathbf{y}$  in  $\mathbb{R}^{|C|}$  as the **output space** of the Softmax layer.

#### 3.1.2 Discretising the Output Space into Permutations

As we saw in Figure 2, there are certain arrangements of Softmax weights for which a target class  $c_t$  cannot be surfaced as the argmax. To understand this phenomenon, it will be helpful to discretise the outputs to a finer granularity: rankings (Burges et al., 2005). In order for a classifier to predict a class  $c_t$  using an argmax decision rule, it must rank  $c_t$  above all other classes by assigning it the largest probability. From this perspective, a classifier assigns each input  $\mathbf{x}$  a permutation  $\pi$  that ranks the class indices in increasing order of probability.

$$\pi : P(c_{\pi_1} | \mathbf{x}) < P(c_{\pi_2} | \mathbf{x}) < \dots < P(c_{\pi_{|C|}} | \mathbf{x}) \quad (3)$$

As an example, if we have 4 classes and obtain probabilities  $P(C | \mathbf{x}) = [.2 \ .4 \ .1 \ .3]^\top$  we assign  $\mathbf{x}$  the permutation  $\pi_{3142}$ , since  $P(c_3 | \mathbf{x}) < P(c_1 | \mathbf{x}) < P(c_4 | \mathbf{x}) < P(c_2 | \mathbf{x})$ . We can readily obtain the coarser argmax decision ( $c_2$ ) by reading off the last index of the permutation.

#### 3.2 How Can Unargmaxable Classes Arise?

A class  $c_t$  is unargmaxable when all permutations that rank  $c_t$  above the rest cannot be realised due to rank constraints. We explain how this happens by combining the following two observations.

**Observation 1.** *We can discretise  $\mathbb{R}^{|C|}$  into regions corresponding to permutations by segmenting the space with hyperplanes.*

The hyperplanes that partition the output space into regions  $\mathcal{R}_\pi$  corresponding to permutations are a well known structure in Combinatorics, the **Braid**

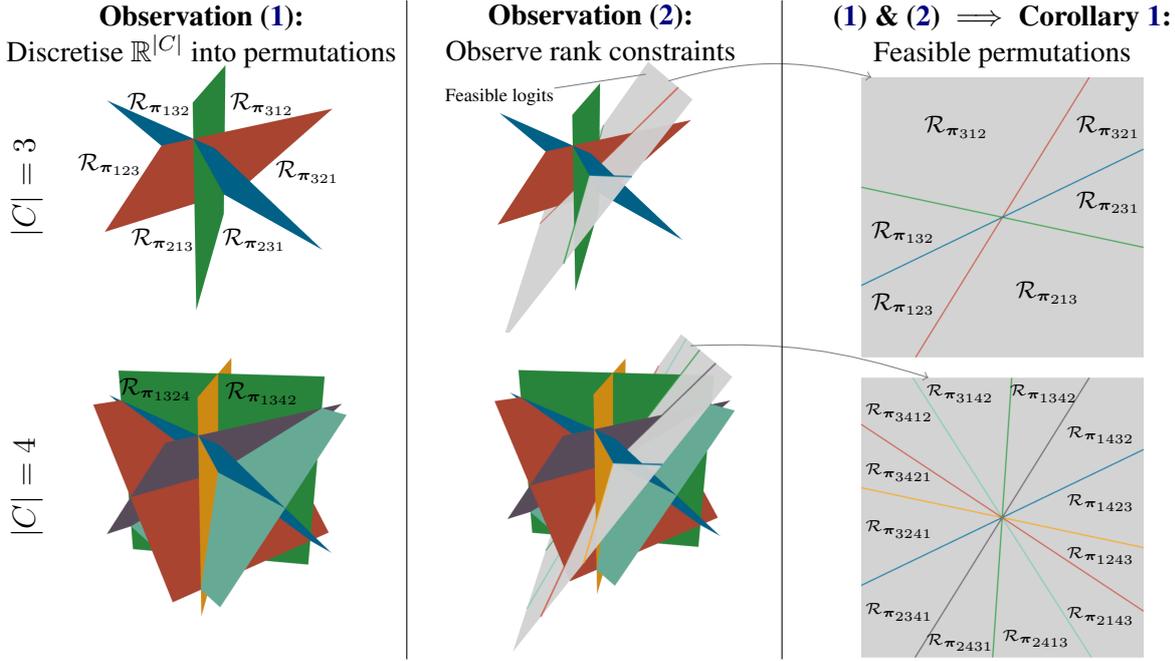


Figure 3: Illustration of Corollary 1 (3<sup>rd</sup> column) as a result of Observation 1 (1<sup>st</sup> column) and Observation 2 (2<sup>nd</sup> column) for softmax( $\mathbf{W}\mathbf{x}$ ),  $\mathbf{W} \in \mathbb{R}^{|C| \times d}$ ,  $d = 2$ . Planes truncated for ease of visualisation. **Top row:** In the left column we see the Braid Arrangement for 3 classes partitioning the output space into 6 regions that correspond to permutations: class rankings in increasing order of probability. In the middle column we see that because  $d = 2$  we can only map  $\mathbf{x}$  to the feasible logits, a plane (grey) defined by  $\mathbf{W}$ . Therefore, in the right column we see that we can only represent permutations that correspond to the regions we can intersect with this plane. For  $|C| = 3$  we can still represent all 6 rankings of 3 classes since any plane in general position will intersect all 6 regions. **Bottom row:** The Braid Arrangement for 4 classes. Since  $d < |C| - 1$  the plane can only intersect 12 regions so only 12/24 permutations are feasible. For example, we see that the plane intersects region  $\mathcal{R}_{\pi_{1342}}$  but not  $\mathcal{R}_{\pi_{1324}}$  and hence  $\pi_{1342}$  is feasible while  $\pi_{1324}$  is not. In fact, the orientation of the plane is such that none of the 6  $\mathcal{R}_{\pi_{****4}}$  regions are intersected. Therefore  $c_4$  cannot be ranked above  $c_1, c_2$  and  $c_3$  and is unargmaxable as in Figures 1 and 2.

**Hyperplane Arrangement (Stanley, 2004).**<sup>5</sup> The Braid Arrangement for 3 and 4 classes is illustrated in rows 1 and 2 of Figure 3 respectively.

In order to be able to rank the classes according to permutation  $\mathcal{R}_{\pi}$ , our network needs to be able to map an input  $\mathbf{x}$  to region  $\mathcal{R}_{\pi}$  in the output space. However, this is not always possible when we have a Softmax Bottleneck as we elaborate below.

**Observation 2.** *When we have rank constraints, only a subspace of  $\mathbb{R}^{|C|}$  is feasible.*

**Case i)** softmax( $\mathbf{W}\mathbf{x}$ ). By calculating  $\mathbf{y} = \mathbf{W}\mathbf{x}$ , the class logits  $\mathbf{y}$  are a linear combination of  $d$  columns of  $\mathbf{W}$ . Therefore, when  $d < |C|$  we can only represent a  $d$ -dimensional subspace of  $\mathbb{R}^{|C|}$  at best. This feasible subspace is illustrated as a grey plane in the middle column of Figure 3.

**Case ii)** softmax( $\mathbf{W}\mathbf{x} + \mathbf{b}$ ). If we also have a bias term  $\mathbf{b}$  the model can choose how to offset the subspace. When the bias term  $\mathbf{b}$  is not in the

column space of  $\mathbf{W}$  the zero vector  $\mathbf{0}$  is no longer a feasible  $\mathbf{y}$  and instead of a linear subspace we have an affine subspace. See Figure 7 in the Appendix for an illustration comparing the two cases.

**Corollary 1.** *A Softmax classifier parametrised by  $\mathbf{W}$  and  $\mathbf{b}$  can rank classes in the order of permutation  $\pi$  iff the affine subspace spanned by  $\mathbf{W}$  and  $\mathbf{b}$  intersects region  $\mathcal{R}_{\pi}$  of the Braid Arrangement.*<sup>6</sup> When  $d < |C| - 1$  there are regions that cannot be intersected.<sup>7</sup> The feasible permutations in our example correspond to the regions formed on the grey plane illustrated in the rightmost column of Figure 3. Note that for  $|C| = 4$  only 12 out of 24 regions can be intersected.

As we make the Softmax Bottleneck narrower by reducing the dimension  $d$  of the Softmax inputs, more permutations become infeasible (Good and

<sup>5</sup>See Appendix B for more details on hyperplane arrangements and the Braid Arrangement specifically.

<sup>6</sup>This insight of slicing the Braid Arrangement was introduced in Kamiya et al. (2011).

<sup>7</sup>When  $d = C - 1$  we can still intersect all regions, because the Braid Arrangement always has rank  $|C| - 1$  (all its normal vectors are perpendicular to the all ones vector  $\mathbf{1}$ ).

Tideman, 1977; Kamiya and Takemura, 2005). Importantly, if we choose  $|C|$  and  $d$  and whether to use a bias term, changing the values of the Softmax weights changes the set of feasible permutations but not the cardinality of the set (Cover, 1967; Smith, 2014). See Appendix C for more details.

**Corollary 2.** *Class  $c_t$  is unargmaxable when any permutation that would rank class  $c_t$  above all other classes is infeasible.*

### 3.2.1 Effect of Softmax Bias Term

Without a bias term the regions corresponding to permutations are unbounded (see the rightmost column of Figure 3). As such, imposing any range restrictions on the Softmax layer inputs  $\mathbf{x}$  does not change the feasible regions as long as the restriction includes the origin. However, when we introduce a bias term we also get bounded regions (see Figure 7 in the Appendix that contrasts the two situations). Therefore, in this case the scale of the inputs to the Softmax layer also matters. If the inputs do not have a large enough range, there will be regions that exist but cannot be reached by the feature encoder.

## 3.3 Exact Algorithm

Given a softmax layer parametrised by  $\mathbf{W}$  and  $\mathbf{b}$ , are there any classes that are unargmaxable? We first describe a slow, but exact algorithm to answer this question.

An exact algorithm will either prove class  $c_t$  is argmaxable by returning a feasible point  $\mathbf{x} : \text{argmax}(\mathbf{W}\mathbf{x} + \mathbf{b}) = c_t$  or it will prove  $c_t$  is unargmaxable by verifying no such point exists.

To check if a region exists that ranks  $c_t$  above all others, we need to find an input  $\mathbf{x} \in \mathbb{R}^d$  that satisfies the following constraints:

$$P(c_i | \mathbf{x}) < P(c_t | \mathbf{x}), \quad \forall i : 1 \leq i \leq |C|, i \neq t \quad (4)$$

Each of the above constraints is equivalent to restricting  $\mathbf{x}$  to a halfspace (see Appendix A). Hence, if all above inequalities are enforced,  $\mathbf{x}$  is restricted to an intersection of halfspaces.

$$(\mathbf{w}_{c_i} - \mathbf{w}_{c_t})^\top \mathbf{x} + (b_{c_i} - b_{c_t}) < 0 \quad (5)$$

$$\forall i : 1 \leq i \leq |C|, i \neq t$$

If the intersection of halfspaces is empty, there is no  $\mathbf{x}$  for which class  $c_t$  can be ranked above all others - and hence  $c_t$  is unargmaxable. We can find a point in an intersection of halfspaces via linear

programming, albeit we found this algorithm to be slow in practice for  $|C| > 1000$ .

### 3.3.1 Chebyshev Center Linear Programme

The Chebyshev center of a polytope (Boyd et al., 2004, p. 417) is the center of the largest ball of radius  $r$  that can be embedded within the polytope. We can find the Chebyshev center  $\mathbf{x}$  and the radius  $r$  with the following linear programme.

$$\begin{aligned} &\text{maximise} && r \\ &\text{subject to} && \mathbf{w}_i^\top \mathbf{x} + r \|\mathbf{w}_i\|_2 \leq -b_i, \quad 1 \leq i \leq |C|-1 \\ &&& \mathbf{x} \leq 100 \\ &&& \mathbf{x} \geq -100 \\ &&& r > 0 \end{aligned} \quad (6)$$

Where  $\mathbf{w}_i = \mathbf{w}_{c_i} - \mathbf{w}_{c_t}$  and  $b_i = b_{c_i} - b_{c_t}$ ,  $\forall i : c_i \neq c_t$ . We further constrain  $\mathbf{x}$  to guarantee the regions are bounded, since the Chebyshev center is not defined otherwise. This constraint also captures the fact that neural network activations are not arbitrarily large.

If the above linear programme is feasible, we know that class  $c_t$  is argmaxable and we also get a lower bound on the volume of the region for which it is solvable by inspecting  $r$ . On the other hand, if the linear programme is infeasible,  $c_t$  is unargmaxable.

## 3.4 Approximate Algorithm

The exact algorithm was too slow to run for the whole vocabulary. In order to avoid running the exact algorithm for every single vocabulary item, we developed an incomplete algorithm (Kautz et al., 2009) with a one-sided error, which can quickly rule out most tokens, leaving only a small number to be checked by the exact algorithm. It proves that  $c_t$  is **argmaxable** by finding an input  $\mathbf{x}$  for which  $c_t$  has the largest activation. Unlike the exact algorithm, if no solution exists it cannot prove that the token is **unargmaxable**. Hence, we terminate our search after a predetermined number of steps. We denote any tokens not shown to be argmaxable by the approximate algorithm as **potentially unargmaxable** and we run the exact algorithm on them. An illustration of the way we combine the exact and approximate algorithms to decide whether class  $c_t$  is argmaxable can be seen in Figure 4.

### 3.4.1 Braid Reflect

The idea behind this approximate algorithm is to use the Braid Hyperplane Arrangement as a map

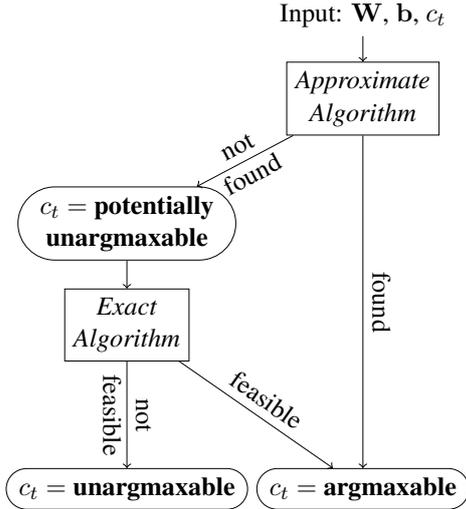


Figure 4: Algorithm to verify whether class  $c_t$  is argmaxable. We first run the approximate algorithm, which quickly proves most vocabulary tokens are argmaxable. If it fails to find a solution in  $N$  steps, we rely on the exact algorithm to either find a solution or prove there is no solution, meaning  $c_t$  is unargmaxable.

---

**Algorithm 1:** Braid reflection step

---

**Data:** Class index  $c_t$ ,  $\mathbf{x} \in \mathbb{R}^d$ ,  
 $\mathbf{W} \in \mathbb{R}^{|C| \times d}$ ,  $\mathbf{b} \in \mathbb{R}^{|C|}$

- 1  $c_i = \text{argmax}(\mathbf{W}\mathbf{x} + \mathbf{b})$
  - 2  $\mathbf{w} = (\mathbf{w}_{c_t} - \mathbf{w}_{c_i})^\top$
  - 3  $b = b_{c_t} - b_{c_i}$
  - 4  $\mathbf{w}' = \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$
  - 5  $d = \mathbf{w}'^\top \mathbf{x}$
  - 6  $\mathbf{x} = \mathbf{x} - 2(d + \frac{b}{\|\mathbf{w}\|_2})\mathbf{w}'$
- 

Figure 5: Move  $\mathbf{x}$  to region where  $P(c_t) > P(c_i)$ .

to guide us towards a point  $\mathbf{x}$  for which  $c_t$  has the largest activation. To show that class  $c_t$  is argmaxable, it suffices to find an input  $\mathbf{x}$  for which the largest probability is assigned to  $c_t$ . Empirically we found this to be easy for most classes.

We begin by interpreting the actual weight vector as the candidate input  $\mathbf{x} = \mathbf{w}_{c_t}^\top$ . We do so since the dot product of two vectors is larger when the two vectors point in the same direction.<sup>8</sup> While the magnitude of the vectors affects the dot product, we found the above initialisation worked well empirically. When  $c_t$  is not the argmax for  $\mathbf{x}$  and  $c_i$  is instead, Relation 5 for  $c_i$  and  $c_t$  will have the wrong sign. The sign of this relation defines which side of the Braid hyperplane for  $c_i$  and  $c_t$  we are on.

<sup>8</sup> $\mathbf{a}^\top \mathbf{b} = \|\mathbf{a}\|_2 \|\mathbf{b}\|_2 \cos \theta$  is maximised for  $\theta = 0$

To correct the sign, we construct the normal vector and offset of the Braid hyperplane (Lines 2, 3 in Figure 5), compute the distance of  $\mathbf{x}$  from it (Line 5), and reflect  $\mathbf{x}$  across it (Line 6).<sup>9</sup> We repeat the above operation until either  $c_t$  is the argmax or we have used up our budget of  $N$  steps.

## 4 Experiments

In this Section we use the combined algorithm from Figure 4 to search models for unargmaxable tokens.

We test 7 LMs and 143 MT models. We find that unargmaxable tokens only occur in 13 MT models, but these are mostly infrequent and noisy vocabulary tokens. We therefore do not expect such tokens to affect translation quality per se.

We also find that nearly all vocabulary tokens of LMs and student MT models can be verified with less than  $N = 10$  steps of the approximate algorithm. In contrast, other MT models need thousands of steps and also rely on the exact algorithm. In this sense, models that need fewer steps are easier to verify: the search problem for their arrangement of Softmax weights is easier.

Throughout the following experiments we assumed the Softmax inputs were bounded in magnitude for all dimensions  $-100 \leq x_i \leq 100$ . As we mentioned in Subsection 3.2.1, if we have a Softmax bias term, there are bounded regions. If the bounded regions are large, even though the outputs are not theoretically bounded, they are practically bounded since neural network feature encoders cannot produce arbitrarily large activations and some regions may be unreachable<sup>10</sup>. For the approximate algorithm, we search for a solution with a patience of  $N = 2500$  steps and resort to the exact algorithm if the approximate method fails or returns a point outside the aforementioned bounds. We use Gurobi (Gurobi Optimization, 2021) as the linear programme solver. We accessed the model parameters either via NumPy (Harris et al., 2020) or PyTorch (Paszke et al., 2019). The experiments took 3 days to run on an AMD 3900X 12-core CPU using 10 threads and 64Gb of RAM.

### 4.1 Language Models (0/7 Unargmaxable)

We checked 7 widely used LMs for unargmaxable tokens. While some of these models such as

<sup>9</sup>When no offset is involved, the reflection operation is the Householder transformation (Householder, 1958).

<sup>10</sup>The validity of our assumption is only relevant for models we find to be bounded. We therefore verified that  $-100 \leq x \leq 100$  holds for two of them, see Appendix F.

BERT (Devlin et al., 2019) are not directly used for generation, a recent trend is to use these large LMs as prompt models (Liu et al., 2021) for few shot learning. A prompt model obviates the need for a separate classifier by rephrasing a classification task as slot filling given a task specific template. Prompt approaches commonly choose the answer for the slot by argmaxing the Softmax distribution obtained by a LM. Hence we verify that there are no answers that are unargmaxable.

BERT, RoBERTa (Liu et al., 2019), XLM-RoBERTa (Conneau et al., 2020) and GPT2 (Radford et al., 2019) did not exhibit any unargmaxable tokens and can be assessed without resorting to the exact algorithm (see Table 4 in the Appendix). Moreover, the LMs were very easy to verify with the approximate algorithm requiring less than 1.2 steps per token on average.

## 4.2 Machine Translation (13/143 Unargmaxable)

| model source    | Helsinki    | FAIR    | Edinburgh | Bergamot               |
|-----------------|-------------|---------|-----------|------------------------|
| unargmaxable    | 13/32       | 0/4     | 0/82      | 0/25                   |
| dataset         | OPUS        | WMT'19  | WMT'17    | multiple <sup>11</sup> |
| architecture    | Transf      | Transf  | LSTM      | Transf                 |
| feature dim $d$ | 512         | 1024    | 500,512   | 256,512,1024           |
| Softmax bias    | ✓           | ✗       | ✓         | ✓                      |
| tied embeds     | enc+dec+out | dec+out | dec+out   | enc+dec+out            |

Table 1: Results for the MT models we verified.

In the case of MT models, the feature encoder comprises the whole encoder-decoder network excluding the last layer of the decoder. We first focus on models which we found to have unargmaxable tokens and then briefly describe models that did not. A summary of the results and characteristics of the models we checked can be seen in Table 1. More detailed results can be found in Tables 5, 6, 7 and 8 in the Appendix.

**Helsinki NLP OPUS (13/32 Unargmaxable).** The 32 models we use for this subset of experiments are MT models released through Hugging Face (Wolf et al., 2020). We use models introduced in Tiedemann and Thottingal (2020). These models are trained on subsets of OPUS. All models are transformer models trained using Marian (Junczys-Dowmunt et al., 2018). They include a bias term, have a tied encoder and decoder and  $d = 512$ .

Unargmaxable tokens, if present, will affect generation in the target language. We therefore restrict our analysis to the target language vocabulary. To

facilitate this, we inspect translation models for which the source and target languages have different scripts. We explore 32 models with source and target pairs amongst Arabic (ar), Hebrew (he), English (en), German (de), French(fr), Spanish (es), Finnish (fi), Polish (pl), Greek (el), Russian (ru), Bulgarian (bg), Korean (ko) and Japanese (ja). We rely on the script to disambiguate between source and target language and discard irrelevant tokens from other languages. We also ignore vocabulary tokens containing digits and punctuation.

In Figure 6 we can see the number of Byte Pair Encoding (BPE; Sennrich et al., 2016) tokens that were unargmaxable for these models, sorted in decreasing order. As can be seen, all tokens are argmaxable for 19/32 language pairs. For the remaining 13 languages, while there can be quite a few unargmaxable tokens, most would not be expected to affect translation quality.

Out of the set of 427 unique unargmaxable BPE tokens, 307/476 are single character subword tokens and only 2 are word stem BPE segments: *erecti* (bg-en) and Предварительны (en-ru) which means “preliminary” in Russian. The rest include the *<unk>* token and noisy subword unicode tokens such as  $\acute{\kappa}\acute{\kappa}\acute{\kappa}\acute{\kappa}$ ,  $\text{ĩĩ}$  and  $\acute{\alpha}\acute{\alpha}\grave{\eta}$ .

On closer inspection of the SentencePiece tokeniser we found that both Предварительны and *erecti* come up as tokenisation alternatives that make them rare and irregular. We found that the Предварительны token was rare since it is capitalised and only occurs once, while another occurrence was caused by a BPE segmentation corner case due to Unicode token variation of Предварительны-е. Other mentions having Предварительны as a substring were split differently. In a similar vein, we found that the *erecti* token occurred due to BPE corner cases for *erecti-0n*, *erecti-lis-*, *erecti-l*, *erecti-*. and *erecti-cle* many of which are misspellings or rare word forms from clinical text. As such, the impact of these tokens being unargmaxable is small since there are alternative ones the MT model can prefer over them which could even correct spelling mistakes.

**FAIR WMT'19 (0/4 Unargmaxable).** We checked 4 FAIR models (en-ru, ru-en, en-de, de-en) submitted to WMT'19 (Ng et al., 2019). These transformer models have  $d = 1024$  and do not employ a Softmax bias term.

None of the FAIR models were found to have unargmaxable tokens, but for some tokens we had

<sup>11</sup><https://github.com/browsermt/students>

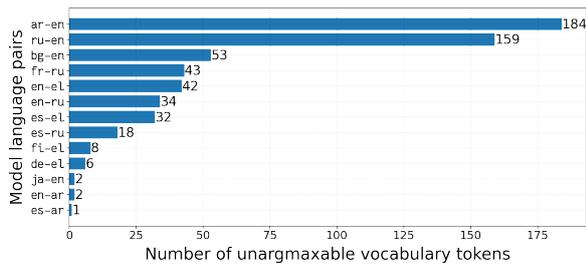


Figure 6: 13/32 HelsinkiNLP models have vocabulary tokens that cannot be predicted using greedy decoding.

to rely on the exact algorithm to show this.

#### Edinburgh WMT’17 (0/82 Unargmaxable).

These WMT’17 submissions (Sennrich et al., 2017) were ensembles of left-to-right trained models (l2r) and right-to-left trained models (r2l). These were LSTMs trained with Nematus using  $d = 500$  or  $d = 512$  and Softmax weights tied with the decoder input embeddings. The models include a bias term.

None of the models have unargmaxable tokens. However, we found that models that comprise an ensemble varied a lot in how easy it was to show that the vocabulary was argmaxable, despite them differing solely in the random seed used for weight initialisation. As an example, zh-en.l2r(1) had 8 tokens that needed to be verified with the exact algorithm, zh-en.l2r(2) had 3 and zh-en.l2r(3) had 366. This highlights that random initialisation alone is enough to lead to very different arrangements of Softmax weights.

**Bergamot (0/25 Unargmaxable).** The Bergamot project<sup>12</sup> model repository contains both large transformer-base and transformer-big teacher models, as well as small knowledge distilled (Kim and Rush, 2016) student models. Student models have  $d = 256$  (tiny) or  $d = 512$  (base), while teacher models have  $d = 1024$ . Interestingly, we find that it is easier to show that student models are argmaxable when compared to teacher models, despite student models having Softmax weights  $1/2$  or  $1/4$  the dimensions of the teacher model.

## 5 Discussion

We conclude from our experiments that *it is possible to have unargmaxable tokens, but this rarely occurs in practice* for tokens that would lead to irrecoverable errors in the MT models we checked. A limitation of our conclusions is that beam search is usually preferred over greedy decoding for MT

<sup>12</sup><https://browser.mt>

models used in practice. We leave the question of whether unargmaxable tokens also impact beam search for future work.

It is challenging to make exact claims about what can cause tokens to be unargmaxable because the models we tested varied in so many ways. However, we outline some general trends below.

### 5.1 Infrequent Tokens Are the Victims

The most general observation is that the tokens that are more likely to be unargmaxable or are hard to prove to be argmaxable are the infrequent ones. This can be seen in Figures 11 and 12 in the Appendix, where the x-axis contains the vocabulary of the models sorted left to right by increasing frequency. Each dot represents the number of steps needed to check whether a token is argmaxable or not, and as can be seen the values to the right are generally much higher than those to the left.

This result is in line with previous work that highlights the limitations of the Softmax layer when modelling rare words for LM (Chen et al., 2016; Labeau and Cohen, 2019) and MT (Nguyen and Chiang, 2018; Raunak et al., 2020) and infrequent classes for image classification (Kang et al., 2020).

### 5.2 Some Models Are Easier to Verify

We found that the LMs and student MT model vocabularies can be shown to be argmaxable with one step of the approximate algorithm on average. On the other hand, for Helsinki NLP and FAIR MT models more than 10 steps were needed.

To put the above observations into context, we also check the behaviour of our algorithms on randomly initialised parameters. If we initialise a Softmax layer of  $|C| = 10000$  classes using a uniform distribution  $U(-1, 1)$  we do not expect unargmaxable tokens to exist after  $d = 30$  (see Figure 10 in the Appendix). Moreover, any randomly initialised parameters can be checked using the approximate algorithm with fewer steps as we increase  $d$ .

From this perspective, it is surprising that student models were easier to show to be argmaxable than the teacher models, despite the Softmax weight dimensionality of the student models being much lower (256 for tiny, versus 1024 for teacher). This shows that effective neural MT models do not need to be hard to check, but nevertheless neural models trained on the original data can sometimes converge to such an arrangement of weights.

## 6 Conclusions and Future Work

In this work we discretised the outputs of Softmax and showed how dimensionality constraints shrink the set of feasible class rankings and can lead to some classes being impossible to predict using argmax. In our experiments we demonstrated that while MT models can have unargmaxable vocabulary tokens, this does not occur often in our experiments. Moreover, for the models we tested the unargmaxable tokens would not create discernible differences in translation quality as the tokens are noisy and infrequent. We release an algorithm to detect whether some classes are unargmaxable with the hope that this will be helpful to the wider community working on a plethora of different models where the observed phenomena may vary.

In future work, we aim to investigate any learnability consequences more closely. As we saw, when using an approximate search algorithm, it is much harder to find argmaxable classes in some models than it is in others. Since gradient descent algorithms are also iterative search algorithms seeking optimal parameters, we hypothesise that it will be challenging to train neural network encoders to map activations to regions of the input space that a search algorithm cannot find easily. Hence, although some tokens may not be provably unargmaxable because of constraints imposed by the Softmax parameters of the last layer, some tokens may still be very hard to produce because of difficulties encountered by the feature encoder. To this end, a more holistic investigation into the consequences of the loss in expressivity in low-rank classifiers is warranted.

### Broader Impact

Unargmaxability directly impacts fairness, since certain model outputs, further from being under-represented, may not be represented at all. As we discussed, low-rank classifiers have limited expressivity compared to full rank classifiers, and thus have to explicitly choose which rankings of classes to retain feasible when using argmax prediction. As such, by choosing to use a low-rank model, we are allowing the data and training procedure to specify which rankings should remain feasible, and harmful biases in our data can be propagated and further exacerbated (Hooker, 2021) by our models due to unargmaxability. For example, it could be the case that underrepresented groups find no representation in the outputs of such models, in the

extreme case where related outputs are unargmaxable. As researchers, we should be aware of this limitation when choosing how to parametrise our models (Hooker et al., 2019) and actively seek to either control such phenomena or verify models are not harmful before moving them from research into production.

In addition to the above considerations, linear classification layers are vulnerable to targeted attacks via data poisoning techniques (Goldblum et al., 2020), especially under the scenario where shared models are used as feature extractors (Ji et al., 2018). A subset of such techniques, known as feature collisions (Shafahi et al., 2018; Goldblum et al., 2020), exploit the arrangement of the training examples in feature space to force the misclassification of a target example. Attacks such as Convex Polytope (Zhu et al., 2019) and Bullseye Polytope (Aghakhani et al., 2021), specifically target the unargmaxability weakness (Cover, 1967; Demeter et al., 2020) we elaborated on in the paper. While such attacks assume they are able to inject examples into a training set used for fine-tuning, this is not an unrealistic assumption. This is especially true for recommender systems, where adversarial attacks can create fake users such that a target item is removed from a target user’s top-k list (Christakopoulou and Banerjee, 2019).

### Acknowledgements

We thank Seraphina Goldfarb-Tarrant, Elizabeth Nielsen and Sabine Weber for help with languages, Beatrice Alex, Sameer Bansal, Panagiotis Eustratiadis, Sharon Goldwater, Chantriolnt-Andreas Kapourani, Oli Liu, Yevgen Matuselych, Kate McCurdy, Laura Perez-Beltrachini, Jesse Sigal, Mark Steedman, Ivan Titov and Sabine Weber for feedback and support, Antonio Vergari for feedback, guidance and tirelessly discussing low-rank constraints and Shay Cohen for insightful suggestions and for pointing us to OEIS. We also thank David Demeter for an extensive discussion on Stolen Probability and the anonymous reviewers for helpful questions and comments.



This work was supported by the Engineering and Physical Sciences Research Council [grant number EP/R513209/1] and Research and Innovation Action *Bergamot*, which has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 825303.

## References

- H. Aghakhani, Dongyu Meng, Yu xiang Wang, Christopher Kruegel, and Giovanni Vigna. 2021. Bullseye polytope: A scalable clean-label poisoning attack with improved transferability. *2021 IEEE European Symposium on Security and Privacy (EuroS&P)*, pages 159–178.
- C. Bradford Barber, David P. Dobkin, and Hannu Huhdanpaa. 1996. The quickhull algorithm for convex hulls. *ACM Trans. Math. Softw.*, 22(4):469–483.
- Srinadh Bhojanapalli, Chulhee Yun, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. 2020. Low-rank bottleneck in multi-head attention models. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 864–873. PMLR.
- Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. 2004. *Convex optimization*. Cambridge university press.
- Christopher J. C. Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Gregory N. Hullender. 2005. Learning to rank using gradient descent. In *ICML*, pages 89–96.
- Wenlin Chen, David Grangier, and Michael Auli. 2016. Strategies for training large vocabulary neural language models. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1975–1985, Berlin, Germany. Association for Computational Linguistics.
- Konstantina Christakopoulou and Arindam Banerjee. 2019. Adversarial attacks on an oblivious recommender. In *Proceedings of the 13th ACM Conference on Recommender Systems, RecSys '19*, page 322–330, New York, NY, USA. Association for Computing Machinery.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Thomas M. Cover. 1967. The number of linearly inducible orderings of points in d-space. *SIAM Journal on Applied Mathematics*, 15(2):434–439.
- David Demeter, Gregory Kimmel, and Doug Downey. 2020. Stolen probability: A structural weakness of neural language models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2191–2197, Online. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Octavian Ganea, Sylvain Gelly, Gary Bécigneul, and Aliaksei Severyn. 2019. Breaking the softmax bottleneck via learnable monotonic pointwise nonlinearities. In *ICML*, pages 2073–2082.
- Micah Goldblum, Dimitris Tsipras, Chulin Xie, Xinyun Chen, Avi Schwarzschild, Dawn Xiaodong Song, Aleksander Madry, Bo Li, and Tom Goldstein. 2020. Dataset security for machine learning: Data poisoning, backdoor attacks, and defenses. *ArXiv*, abs/2012.10544.
- I.J Good and T.N Tideman. 1977. Stirling numbers and a geometric structure from voting theory. *Journal of Combinatorial Theory, Series A*, 23(1):34–45.
- Gurobi Optimization. 2021. *Gurobi Optimizer Reference Manual*.
- Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. 2020. Array programming with NumPy. *Nature*, 585(7825):357–362.
- Geoffrey E Hinton and Richard Zemel. 1994. Autoencoders, minimum description length and helmholtz free energy. In *Advances in Neural Information Processing Systems*, volume 6. Morgan-Kaufmann.
- Sara Hooker. 2021. Moving beyond "algorithmic bias is a data problem". *Patterns (New York, N.Y.)*, 2(4):100241–100241. 33982031[pmid].
- Sara Hooker, Aaron Courville, Yann Dauphin, and Andrea Frome. 2019. Selective Brain Damage: Measuring the Disparate Impact of Model Pruning. *arXiv e-prints*.
- Alston S Householder. 1958. Unitary triangularization of a nonsymmetric matrix. *Journal of the ACM (JACM)*, 5(4):339–342.
- Yujie Ji, Xinyang Zhang, Shouling Ji, Xiapu Luo, and Ting Wang. 2018. Model-reuse attacks on deep learning systems. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, CCS '18*, page 349–363, New York, NY, USA. Association for Computing Machinery.

- Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. [Marian: Fast neural machine translation in C++](#). In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia. Association for Computational Linguistics.
- Hidehiko Kamiya and Akimichi Takemura. 2005. Characterization of rankings generated by linear discriminant analysis. *Journal of multivariate analysis*, 92(2):343–358.
- Hidehiko Kamiya, Akimichi Takemura, and Hiroaki Terao. 2011. [Ranking patterns of unfolding models of codimension one](#). *Advances in Applied Mathematics*, 47(2):379–400.
- Sekitoshi Kanai, Yasuhiro Fujiwara, Yuki Yamanaka, and Shuichi Adachi. 2018. [Sigsoftmax: Reanalysis of the softmax bottleneck](#). In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis. 2020. [Decoupling representation and classifier for long-tailed recognition](#). In *International Conference on Learning Representations*.
- Henry A. Kautz, Ashish Sabharwal, and Bart Selman. 2009. Incomplete algorithms. In *Handbook of Satisfiability*.
- Yoon Kim and Alexander M. Rush. 2016. [Sequence-level knowledge distillation](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327, Austin, Texas. Association for Computational Linguistics.
- Matthieu Labeau and Shay B. Cohen. 2019. [Experimenting with power divergences for language modeling](#). 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 4104–4114, Hong Kong, China. Association for Computational Linguistics.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. [Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing](#). *ArXiv*, abs/2107.13586.
- Y. Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, M. Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#). *ArXiv*, abs/1907.11692.
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. [Multilingual denoising pre-training for neural machine translation](#). *Transactions of the Association for Computational Linguistics*, 8:726–742.
- David J. C. Mackay. 2004. Information theory, inference, and learning algorithms. *IEEE Transactions on Information Theory*, 50:2544–2545.
- Leland McInnes, John Healy, Nathaniel Saul, and Lukas Grossberger. 2018. [Umap: Uniform manifold approximation and projection](#). *The Journal of Open Source Software*, 3(29):861.
- Stephen Merity, Nitish Shirish Keskar, and Richard Socher. 2018. [Regularizing and optimizing LSTM language models](#). In *International Conference on Learning Representations*.
- Nathan Ng, Kyra Yee, Alexei Baevski, Myle Ott, Michael Auli, and Sergey Edunov. 2019. [Facebook FAIR’s WMT19 news translation task submission](#). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 314–319, Florence, Italy. Association for Computational Linguistics.
- Toan Nguyen and David Chiang. 2018. [Improving lexical choice in neural machine translation](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 334–343, New Orleans, Louisiana. Association for Computational Linguistics.
- Dimitris Papadimitriou and Swayambhoo Jain. 2021. [Data-driven low-rank neural network compression](#). In *2021 IEEE International Conference on Image Processing (ICIP)*, pages 3547–3551.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. [Pytorch: An imperative style, high-performance deep learning library](#). In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#).
- Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. 2017. [On the expressive power of deep neural networks](#). In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 2847–2854. PMLR.

- Vikas Raunak, Siddharth Dalmia, Vivek Gupta, and Florian Metzger. 2020. [On long-tailed phenomena in neural machine translation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3088–3095, Online. Association for Computational Linguistics.
- Tara N. Sainath, Brian Kingsbury, Vikas Sindhwani, Ebru Arisoy, and Bhuvana Ramabhadran. 2013. [Low-rank matrix factorization for deep neural network training with high-dimensional output targets](#). In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 6655–6659.
- Rico Sennrich, Alexandra Birch, Anna Currey, Ulrich Germann, Barry Haddow, Kenneth Heafield, Antonio Valerio Miceli Barone, and Philip Williams. 2017. [The University of Edinburgh’s neural MT systems for WMT17](#). In *Proceedings of the Second Conference on Machine Translation*, pages 389–399, Copenhagen, Denmark. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [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)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Ali Shafahi, W. Ronny Huang, Mahyar Najibi, Octavian Suci, Christoph Studer, Tudor Dumitras, and Tom Goldstein. 2018. [Poison frogs! targeted clean-label poisoning attacks on neural networks](#). In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Warren D. Smith. 2014. [D-dimensional orderings and stirling numbers](#). [Online; accessed 05-November-2021].
- Richard P. Stanley. 2004. An introduction to hyperplane arrangements. In *Lecture notes, IAS/Park City Mathematics Institute*.
- The Sage Developers. 2021. *SageMath, the Sage Mathematics Software System (Version 9.5)*. <https://www.sagemath.org>.
- Jörg Tiedemann. 2020. [The Tatoeba Translation Challenge – Realistic data sets for low resource and multilingual MT](#). In *Proceedings of the Fifth Conference on Machine Translation*, pages 1174–1182, Online. Association for Computational Linguistics.
- Jörg Tiedemann and Santhosh Thottingal. 2020. [OPUS-MT – building open translation services for the world](#). In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 479–480, Lisboa, Portugal. European Association for Machine Translation.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Zhilin Yang, Zihang Dai, Ruslan Salakhutdinov, and William W. Cohen. 2018. [Breaking the softmax bottleneck: A high-rank RNN language model](#). In *International Conference on Learning Representations*.
- Chulhee Yun, Yin-Wen Chang, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. 2020. [O\(n\) connections are expressive enough: Universal approximability of sparse transformers](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 13783–13794. Curran Associates, Inc.
- Chen Zhu, W. Ronny Huang, Hengduo Li, Gavin Taylor, Christoph Studer, and Tom Goldstein. 2019. [Transferable clean-label poisoning attacks on deep neural nets](#). In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7614–7623. PMLR.

## A Halfspace interpretation

As promised, here is the derivation showing that if  $P(c_i | \mathbf{x}) < P(c_j | \mathbf{x})$  then  $\mathbf{x}$  is constrained to a halfspace.

We have:

$$\begin{aligned}
 P(c_i | \mathbf{x}) < P(c_j | \mathbf{x}) &\iff \\
 \frac{e^{\mathbf{w}_{c_i}^\top \mathbf{x} + b_{c_i}}}{\sum_{i'} e^{\mathbf{w}_{c_{i'}}^\top \mathbf{x} + b_{c_{i'}}}} < \frac{e^{\mathbf{w}_{c_j}^\top \mathbf{x} + b_{c_j}}}{\sum_{i'} e^{\mathbf{w}_{c_{i'}}^\top \mathbf{x} + b_{c_{i'}}}} &\iff \\
 e^{\mathbf{w}_{c_i}^\top \mathbf{x} + b_{c_i}} < e^{\mathbf{w}_{c_j}^\top \mathbf{x} + b_{c_j}} &\iff \\
 \frac{e^{\mathbf{w}_{c_i}^\top \mathbf{x} + b_{c_i}}}{e^{\mathbf{w}_{c_j}^\top \mathbf{x} + b_{c_j}}} < 1 &\iff \\
 e^{(\mathbf{w}_{c_i} - \mathbf{w}_{c_j})^\top \mathbf{x} + (b_{c_i} - b_{c_j})} < e^0 &\iff \\
 (\mathbf{w}_{c_i} - \mathbf{w}_{c_j})^\top \mathbf{x} + (b_{c_i} - b_{c_j}) < 0
 \end{aligned} \tag{7}$$

$\mathbf{x}$  is therefore constrained to a halfspace defined by normal vector  $\mathbf{w}_{c_i} - \mathbf{w}_{c_j}$  and offset by  $b_{c_i} - b_{c_j}$ . This linear form defined by the normal vector and offset is the “shadow” in the input dimension of our friend, the Braid Arrangement, as we will make clear in the next Section (see Derivation 11).

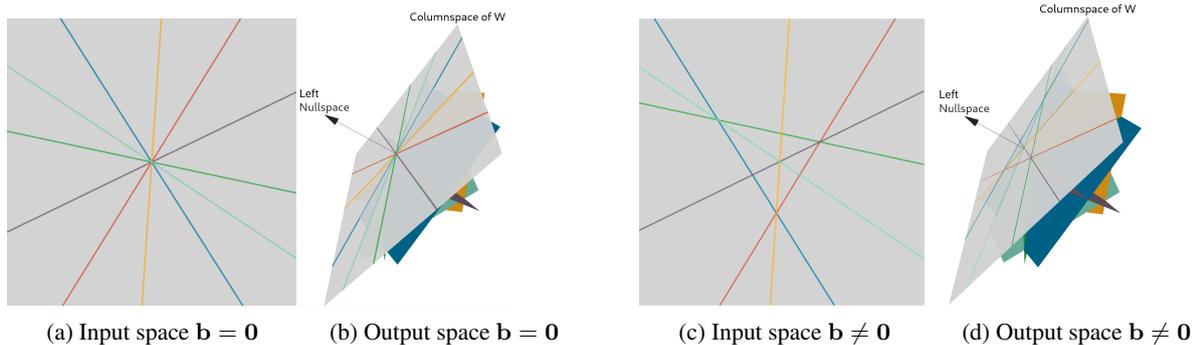


Figure 7: Effect of bias term  $\mathbf{b}$  on feasible permutations of  $\text{softmax}(\mathbf{W}\mathbf{x} + \mathbf{b})$ ,  $\mathbf{W} \in \mathbb{R}^{|C| \times d}$ ,  $d = 2$ ,  $|C| = 4$ . Having a bias term offsets the grey plane and allows it to not pass through the origin. This increases the number of regions by creating bounded regions seen in Subfigures c and d. Each region intersected by the grey 2D plane corresponds to a feasible permutation. We therefore obtain 18/24 feasible permutations if we include a bias term, compared to 12/24 without one.

## B Hyperplane Arrangements

Excellent resources to learn more about hyperplane arrangements are [Stanley \(2004\)](#) and [Federico Ardila's lectures on polytopes \(see Lecture 34 onwards\)](#). Connections between hyperplane arrangement theory and Machine Learning can be found in [Mackay \(2004, Chapter 40\)](#). For those who prefer a more gentle introduction via a hands on approach, [Sagemath \(The Sage Developers, 2021\)](#) contains implementations of many hyperplane arrangements and functions that we found useful when learning this material. We give a brief introduction to hyperplane arrangements below.

A *hyperplane* in a vector space  $\mathbb{R}^d$  is an affine subspace of dimension  $d - 1$ . The hyperplane  $\mathcal{H}$  has one degree of freedom removed by specifying a constraint: a normal vector  $\mathbf{w} \in \mathbb{R}^d$  to which it is perpendicular. The hyperplane may also be offset by  $b$  in that direction  $\mathcal{H} = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{x} = b\}$ .

A *real hyperplane arrangement*  $\mathcal{A}$  is defined as a set of  $n$  hyperplanes in  $\mathbb{R}^d$ ,  $\mathcal{A} = \{\mathcal{H}_1, \mathcal{H}_2 \dots \mathcal{H}_n\}$ . The set of *regions*  $\mathcal{R}$  defined by a hyperplane arrangement  $\mathcal{A}$  are the connected components  $X$  of Euclidean space  $\mathbb{R}^d$  left when we remove the hyperplanes  $\mathcal{A}$ , namely  $X = \mathbb{R}^d - \bigcup_{\mathcal{H} \in \mathcal{A}} \mathcal{H}$ . As an example, Subfigure (a) in Figure 7 has 12 regions while Subfigure (c) has 18 regions.

### B.1 Braid Arrangement

The Braid Arrangement  $\mathcal{B}_n$  is a hyperplane arrangement that partitions space into  $n!$  regions corresponding to permutations. It can be constructed in  $\mathbb{R}^n$  from the standard basis, the rows of the identity matrix  $\mathbf{I}$ ,  $\mathbf{e}_i = \text{row}_i(\mathbf{I})^\top$ ,  $\mathbf{e}_i \in \mathbb{R}^n$ , by taking all  $\binom{n}{2}$  pairs of differences between them, each differ-

ence defining the normal vector of a hyperplane  $\mathcal{H}_{i,j}$  of the Braid Arrangement.

$$\mathcal{B}_n = \{\mathcal{H}_{i,j} \quad \forall i, j : 1 \leq i < j \leq n\}, \quad (8)$$

$$\mathcal{H}_{i,j} = \{\mathbf{x} \in \mathbb{R}^n : (\mathbf{e}_i - \mathbf{e}_j)^\top \mathbf{x} = 0\}$$

The Braid Arrangement for  $n = 3$  and  $n = 4$  can be seen in Figure 3. It has  $\binom{n}{2}$  hyperplanes, one per pair of dimensions in  $\mathbb{R}^n$ . Hence there are 3 hyperplanes for  $|C| = 3$  and 6 hyperplanes for  $|C| = 4$ . As an example, when we have 4 classes the normal vector for  $\mathcal{H}_{1,3}$  is  $\mathbf{w}_{1,3} = [1 \ 0 \ -1 \ 0]^\top$ . As can be verified by taking the dot product  $\mathbf{w}_{i,j}^\top \mathbf{x}$ , the result is positive if  $x_i > x_j$  and negative if vice versa. Therefore, each hyperplane bisects space into two regions one for each possible ranking of the pair of coordinates.

To see how the hyperplanes intersect to give us a region  $\mathcal{R}_\pi$ , we express a permutation (total order) over  $|C|$  classes, such as that in Relation 3, using a chain of  $|C| - 1$  pairwise inequalities.

$$P(c_{\pi_i} | \mathbf{x}) < P(c_{\pi_{i+1}} | \mathbf{x}), \quad 1 \leq i \leq |C| - 1 \quad (9)$$

Each above constraint is equivalent to choosing a side of a braid hyperplane. By imposing all constraints, we obtain a region  $\mathcal{R}_\pi$  as the intersection of  $|C| - 1$  halfspaces. There is therefore bijection between permutations and regions of the Braid Arrangement  $\pi \leftrightarrow \mathcal{R}_\pi$ .

### B.2 Restricting the Braid Arrangement to Lower Dimensions

In the Softmax layer of a neural network we often compute the output space activations  $\mathbf{y} \in \mathbb{R}^n$  by applying a final affine layer to the Softmax input

space  $\mathbf{x} \in \mathbb{R}^d$ .

$$\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}, \quad \mathbf{W} \in \mathbb{R}^{n \times d}, \mathbf{b} \in \mathbb{R}^n \quad (10)$$

What do the Braid Arrangement hyperplanes look like in the input dimension  $d$ ? Let us start from the output space  $\mathbb{R}^n$  and work backwards towards the input space  $\mathbb{R}^d$ .

$$\begin{aligned} y_i < y_j &\implies (\mathbf{e}_i - \mathbf{e}_j)^\top \mathbf{y} < 0 \\ &\mathbf{e}_i^\top \mathbf{y} - \mathbf{e}_j^\top \mathbf{y} < 0 \\ &\mathbf{e}_i^\top (\mathbf{W}\mathbf{x} + \mathbf{b}) - \mathbf{e}_j^\top (\mathbf{W}\mathbf{x} + \mathbf{b}) < 0 \\ &\mathbf{w}_i^\top \mathbf{x} + b_i - \mathbf{w}_j^\top \mathbf{x} - b_j < 0 \\ &(\mathbf{w}_i - \mathbf{w}_j)^\top \mathbf{x} + (b_i - b_j) < 0 \end{aligned} \quad (11)$$

We therefore see that if  $d < n$  we can think of how the Braid Arrangement classifies outputs into permutations from two equivalent perspectives:

- In the output space  $\mathbb{R}^n$  not all  $\mathbf{y}$  are feasible, we can only classify an input  $\mathbf{x}$  as a permutation  $\pi$  if the affine layer can map  $\mathbf{x}$  to  $\mathcal{R}_\pi$ . This can be seen in Subfigures b and d of Figure 7 where the feasible outputs are a plane that intersects the Braid Arrangement.
- In the input space  $\mathbb{R}^d$  all  $\mathbf{x}$  are feasible but we only see the projection of the Braid Arrangement in this lower dimension. This can be seen in Subfigures a and c of Figure 7.

The construction of the Braid Arrangement in the input space is illustrated in Figure 8, albeit without the bias term.

### C Number of Regions (Feasible Permutations) of the Restricted Braid Arrangement

The number of feasible permutations is invariant to specific choices of  $\mathbf{W}$  and  $\mathbf{b}$  (Cover, 1967; Smith, 2014) and only depends on the dimensionality of the softmax inputs  $d$ , the number of classes  $|C|$  and whether we specify a bias term  $\mathbf{b}$  not in the column space of  $\mathbf{W}$ . Namely, the cardinality of the set of feasible permutations does not change, but the members of the set do - they depend on the specific values in  $\mathbf{W}$  and  $\mathbf{b}$ . There exists a recurrence formula to obtain the number of feasible permutations for a particular  $|C|$  and  $d$  (Good and Tideman, 1977; Kamiya and Takemura, 2005). See our code and the relations in (Smith, 2014) for more details.

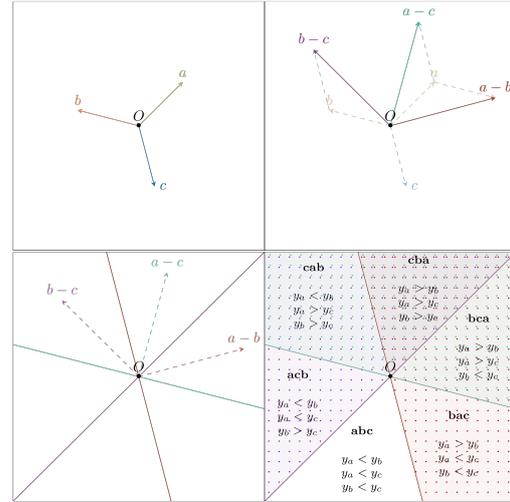


Figure 8: Constructing the Braid Arrangement in the input space for  $|C| = 3$  classes and  $d = 2$ . *Top left:* The Softmax weights  $\mathbf{W} \in \mathbb{R}^{|C| \times d}$  for 3 classes,  $a, b, c$ . Each vector is a row of the weight matrix. *Top right:* We form the normal vectors for the braid hyperplanes by taking all pairs of differences between the basis vectors. *Bottom left:* The Braid hyperplanes are perpendicular to the normal vectors. Each hyperplane bisects space into two regions, one comprises the set of  $\mathbf{x}$  for which class  $i$  has a larger activation than class  $j$  and the second vice versa. *Bottom right:* The hyperplanes partition space into  $3! = 6$  regions corresponding to permutations. Each permutation contains the indices that sort the activations over classes in increasing order. Softmax decision boundaries are unions of two regions, e.g. regions **cba** and **bca** for class **a**.

#### C.1 Softmax with no Bias Term

The number of feasible permutations as a function of  $|C|$  and  $d$  when we have a Softmax with no bias term can be seen in Table 2. When  $d \geq |C| - 1$  all permutations corresponding to ways of ranking  $|C|$  classes are feasible (table cells with  $d = |C| - 1$  are highlighted in bold). However, as we make the Softmax Bottleneck narrower, we can represent less permutations, as can be seen from the numbers reported below the diagonal.

#### C.2 Softmax with Bias Term

The number of feasible permutations as a function of  $|C|$  and  $d$  when we have a Softmax with a bias term is larger as can be seen in Table 3. As we saw in Figure 7, this is because a bias term can offset the representable linear subspace to an affine subspace which can intersect more regions of the Braid Arrangement.

|                      |    | BOTTLENECK DIMENSIONALITY $d$ |           |           |            |            |             |              |               |                |         |
|----------------------|----|-------------------------------|-----------|-----------|------------|------------|-------------|--------------|---------------|----------------|---------|
|                      |    | 1                             | 2         | 3         | 4          | 5          | 6           | 7            | 8             | 9              | 10      |
| NUMBER CLASSES $ C $ | 2  | <b>2</b>                      | 2         | 2         | 2          | 2          | 2           | 2            | 2             | 2              | 2       |
|                      | 3  | 2                             | <b>6</b>  | 6         | 6          | 6          | 6           | 6            | 6             | 6              | 6       |
|                      | 4  | 2                             | <i>12</i> | <b>24</b> | 24         | 24         | 24          | 24           | 24            | 24             | 24      |
|                      | 5  | 2                             | 20        | 72        | <b>120</b> | 120        | 120         | 120          | 120           | 120            | 120     |
|                      | 6  | 2                             | 30        | 172       | 480        | <b>720</b> | 720         | 720          | 720           | 720            | 720     |
|                      | 7  | 2                             | 42        | 352       | 1512       | 3600       | <b>5040</b> | 5040         | 5040          | 5040           | 5040    |
|                      | 8  | 2                             | 56        | 646       | 3976       | 14184      | 30240       | <b>40320</b> | 40320         | 40320          | 40320   |
|                      | 9  | 2                             | 72        | 1094      | 9144       | 45992      | 143712      | 282240       | <b>362880</b> | 362880         | 362880  |
|                      | 10 | 2                             | 90        | 1742      | 18990      | 128288     | 557640      | 1575648      | 2903040       | <b>3628800</b> | 3628800 |

Table 2: Number of permutation regions defined by a bottlenecked Softmax layer  $\text{Softmax}(\mathbf{W}x)$  with no bias term. When  $d \geq |C| - 1$  all permutations corresponding to ways of ranking  $|C|$  classes are feasible. 12 in italics corresponds to the number of regions shown in the left Subfigure of Figure 7. <https://oeis.org/A071223>.

|                      |    | BOTTLENECK DIMENSIONALITY $d$ |           |           |            |            |             |              |               |                |         |
|----------------------|----|-------------------------------|-----------|-----------|------------|------------|-------------|--------------|---------------|----------------|---------|
|                      |    | 1                             | 2         | 3         | 4          | 5          | 6           | 7            | 8             | 9              | 10      |
| NUMBER CLASSES $ C $ | 2  | <b>2</b>                      | 2         | 2         | 2          | 2          | 2           | 2            | 2             | 2              | 2       |
|                      | 3  | 4                             | <b>6</b>  | 6         | 6          | 6          | 6           | 6            | 6             | 6              | 6       |
|                      | 4  | 7                             | <i>18</i> | <b>24</b> | 24         | 24         | 24          | 24           | 24            | 24             | 24      |
|                      | 5  | 11                            | 46        | 96        | <b>120</b> | 120        | 120         | 120          | 120           | 120            | 120     |
|                      | 6  | 16                            | 101       | 326       | 600        | <b>720</b> | 720         | 720          | 720           | 720            | 720     |
|                      | 7  | 22                            | 197       | 932       | 2556       | 4320       | <b>5040</b> | 5040         | 5040          | 5040           | 5040    |
|                      | 8  | 29                            | 351       | 2311      | 9080       | 22212      | 35280       | <b>40320</b> | 40320         | 40320          | 40320   |
|                      | 9  | 37                            | 583       | 5119      | 27568      | 94852      | 212976      | 322560       | <b>362880</b> | 362880         | 362880  |
|                      | 10 | 46                            | 916       | 10366     | 73639      | 342964     | 1066644     | 2239344      | 3265920       | <b>3628800</b> | 3628800 |

Table 3: Number of permutation regions defined by a bottlenecked Softmax layer  $\text{Softmax}(\mathbf{W}x + \mathbf{b})$ . When  $d \geq |C| - 1$  all permutations corresponding to ways of ranking  $|C|$  classes are feasible. 18 in italics corresponds to the number of regions shown in the right Subfigure of Figure 7.

## D Braid Reflect Approximate Algorithm

---

### Algorithm 2: Braid reflect

---

**Data:** Class index  $c_t$ ,  
 $\mathbf{W} \in \mathbb{R}^{|C| \times d}$ ,  $\mathbf{b} \in \mathbb{R}^{|C|}$

**Result:** Whether  $c_t$  is unargmaxable

- 1 unargmaxable = true
- 2 patience = 2500
- 3  $\mathbf{x} = \mathbf{w}_{c_t}^\top$
- 4 **while** *patience* **do**
- 5      $c_i = \text{argmax}(\mathbf{W}\mathbf{x} + \mathbf{b})$
- 6     **if**  $c_i = c_t$  **then**
- 7         unargmaxable = false
- 8         **break**
- 9     **else**
- 10         $\mathbf{w} = (\mathbf{w}_{c_t} - \mathbf{w}_{c_i})^\top$
- 11         $b = b_{c_t} - b_{c_i}$
- 12         $\mathbf{w}' = \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$
- 13         $d = \mathbf{w}'^\top \mathbf{x}$
- 14         $\mathbf{x} = \mathbf{x} - 2(d + \frac{b}{\|\mathbf{w}\|_2})\mathbf{w}'$
- 15        patience = patience - 1
- 16     **end**
- 17 **end**

---

Figure 9: Approximate algorithm to detect whether class  $c_t$  is unargmaxable.

## E Unargmaxable Token Search Results

| model             | # potentially unargmaxable | # unargmaxable |
|-------------------|----------------------------|----------------|
| bert-base-cased   | 0                          | 0              |
| bert-base-uncased | 0                          | 0              |
| roberta-base      | 0                          | 0              |
| roberta-large     | 0                          | 0              |
| xlm-roberta-base  | 0                          | 0              |
| xlm-roberta-large | 0                          | 0              |
| gpt2              | 0                          | 0              |

Table 4: Unargmaxable token search results for LMs. **potentially unargmaxable** is the number of tokens that the approximate algorithm failed to prove were argmaxable. **argmaxable** is the number of unargmaxable tokens according to the exact algorithm. No tokens were found to be unargmaxable.

| source   | model            | # potentially unargmaxable | # unargmaxable |
|----------|------------------|----------------------------|----------------|
|          | opus-mt-ja-en    | 109                        | 2              |
|          | opus-mt-ru-en    | 90                         | 159            |
|          | opus-mt-bg-en    | 93                         | 53             |
|          | opus-mt-ja-en(2) | 14                         | 0              |
|          | opus-mt-ar-en    | 40                         | 184            |
|          | opus-mt-en-el    | 75                         | 42             |
|          | opus-mt-de-el    | 115                        | 6              |
|          | opus-mt-ar-el    | 41                         | 0              |
|          | opus-mt-es-el    | 67                         | 32             |
|          | opus-mt-fi-el    | 57                         | 8              |
|          | opus-mt-ar-he    | 3                          | 0              |
|          | opus-mt-de-he    | 4                          | 0              |
|          | opus-mt-es-he    | 3                          | 0              |
|          | opus-mt-fr-he    | 1                          | 0              |
|          | opus-mt-fi-he    | 7                          | 0              |
| Helsinki | opus-mt-ja-he    | 0                          | 0              |
| NLP      | opus-mt-en-ar    | 21                         | 2              |
|          | opus-mt-el-ar    | 12                         | 0              |
|          | opus-mt-es-ar    | 17                         | 1              |
|          | opus-mt-fr-ar    | 17                         | 0              |
|          | opus-mt-he-ar    | 7                          | 0              |
|          | opus-mt-it-ar    | 8                          | 0              |
|          | opus-mt-ja-ar    | 4                          | 0              |
|          | opus-mt-pl-ar    | 52                         | 0              |
|          | opus-mt-ru-ar    | 8                          | 0              |
|          | opus-mt-en-ru    | 98                         | 34             |
|          | opus-mt-es-ru    | 42                         | 18             |
|          | opus-mt-fi-ru    | 1                          | 0              |
|          | opus-mt-fr-ru    | 34                         | 43             |
|          | opus-mt-he-ru    | 5                          | 0              |
|          | opus-mt-ja-ru    | 13                         | 0              |
|          | opus-mt-ko-ru    | 2                          | 0              |

Table 5: Unargmaxable token search results for Helsinki NLP OPUS models. **potentially unargmaxable** is the number of tokens that the approximate algorithm failed to prove were argmaxable. **unargmaxable** is the number of unargmaxable tokens according to the exact algorithm. For 13/32 models some infrequent tokens were found to be unargmaxable.

| source | model                | # potentially unargmaxable | # unargmaxable |
|--------|----------------------|----------------------------|----------------|
|        | facebook/wmt19-en-ru | 5                          | 0              |
| FAIR   | facebook/wmt19-ru-en | 64                         | 0              |
|        | facebook/wmt19-de-en | 173                        | 0              |
|        | facebook/wmt19-en-de | 184                        | 0              |

Table 6: Unargmaxable token search results for FAIR WMT'19 models. **potentially unargmaxable** is the number of tokens that the approximate algorithm failed to prove were argmaxable. **unargmaxable** is the number of unargmaxable tokens according to the exact algorithm. No tokens were found to be unargmaxable.

| source   | model                  | # potentially unargmaxable | # unargmaxable |
|----------|------------------------|----------------------------|----------------|
|          | cs-en.student.base     | 0                          | 0              |
|          | es-en.teacher.bigx2(1) | 0                          | 0              |
|          | es-en.teacher.bigx2(2) | 0                          | 0              |
|          | en-es.teacher.bigx2(1) | 0                          | 0              |
|          | en-es.teacher.bigx2(2) | 0                          | 0              |
|          | et-en.teacher.bigx2(1) | 2                          | 0              |
|          | et-en.teacher.bigx2(2) | 1                          | 0              |
|          | en-et.teacher.bigx2(1) | 1                          | 0              |
|          | en-et.teacher.bigx2(2) | 1                          | 0              |
|          | nb-en.teacher.base     | 0                          | 0              |
|          | nn-en.teacher.base     | 0                          | 0              |
|          | is-en.teacher.base     | 0                          | 0              |
| Bergamot | cs-en.student.base     | 0                          | 0              |
|          | cs-en.student.tiny11   | 0                          | 0              |
|          | en-cs.student.base     | 0                          | 0              |
|          | en-cs.student.tiny11   | 0                          | 0              |
|          | en-de.student.base     | 0                          | 0              |
|          | en-de.student.tiny11   | 0                          | 0              |
|          | es-en.student.tiny11   | 0                          | 0              |
|          | en-es.student.tiny11   | 0                          | 0              |
|          | et-en.student.tiny11   | 0                          | 0              |
|          | en-et.student.tiny11   | 0                          | 0              |
|          | is-en.student.tiny11   | 0                          | 0              |
|          | nb-en.student.tiny11   | 0                          | 0              |
|          | nn-en.student.tiny11   | 0                          | 0              |

Table 7: Unargmaxable token search results for Bergamot models. **potentially unargmaxable** is the number of tokens that the approximate algorithm failed to prove were argmaxable. **unargmaxable** is the number of unargmaxable tokens according to the exact algorithm. No tokens were found to be unargmaxable. Interestingly, student models were much easier to prove argmaxable than teacher models, despite student model Softmax weights being lower dimensional.

| source    | model          | # potentially unargmaxable | # unargmaxable |
|-----------|----------------|----------------------------|----------------|
|           | en-cs.l2r(1-4) | $\leq 2$                   | 0              |
|           | en-cs.r2l(1-4) | $\leq 1$                   | 0              |
|           | cs-en.l2r(1-4) | $\leq 2$                   | 0              |
|           | cs-en.r2l(1-4) | 0                          | 0              |
|           | en-de.l2r(1-4) | $\leq 1$                   | 0              |
|           | en-de.r2l(1-4) | $\leq 2$                   | 0              |
|           | de-en.l2r(1-4) | $\leq 2$                   | 0              |
|           | de-en.r2l(1-4) | 0                          | 0              |
|           | en-ru.l2r(1-4) | 0                          | 0              |
|           | ru-en.l2r(1-4) | 0                          | 0              |
|           | ru-en.r2l(1-4) | 0                          | 0              |
|           | en-tr.l2r(1-4) | $\leq 5$                   | 0              |
|           | en-tr.r2l(1-4) | $\leq 4$                   | 0              |
|           | lv-en.l2r(1-4) | 0                          | 0              |
| WMT' 17   | lv-en.r2l(1-4) | $\leq 1$                   | 0              |
| Edinburgh | tr-en.l2r(1)   | 2                          | 0              |
|           | tr-en.l2r(2)   | 8                          | 0              |
|           | tr-en.l2r(3)   | 6                          | 0              |
|           | tr-en.l2r(4)   | 2                          | 0              |
|           | tr-en.r2l(1)   | 4                          | 0              |
|           | tr-en.r2l(2)   | 0                          | 0              |
|           | tr-en.r2l(3)   | 6                          | 0              |
|           | tr-en.r2l(4)   | 4                          | 0              |
|           | en-zh.l2r(1)   | 3                          | 0              |
|           | en-zh.l2r(2)   | 3                          | 0              |
|           | en-zh.l2r(3)   | 14                         | 0              |
|           | en-zh.l2r(4)   | 1                          | 0              |
|           | en-zh.r2l(1)   | 2                          | 0              |
|           | en-zh.r2l(2)   | 0                          | 0              |
|           | en-zh.r2l(3)   | 7                          | 0              |
|           | en-zh.r2l(4)   | 7                          | 0              |
|           | zh-en.l2r(1)   | 8                          | 0              |
|           | zh-en.l2r(2)   | 3                          | 0              |
|           | zh-en.l2r(3)   | 366                        | 0              |
|           | zh-en.r2l(1-3) | $\leq 3$                   | 0              |

Table 8: Unargmaxable token search results for Edinburgh WMT' 17 submission (ensemble) models. **potentially unargmaxable** is the number of tokens that the approximate algorithm failed to prove were argmaxable. **unargmaxable** is the number of unargmaxable tokens according to the exact algorithm. **r2l** and **l2r** refer to training direction, with **l2r** denoting training left to right and **r2l** right to left. Models submitted were ensembles, hence there are more than one model per language pair and direction. When all models per language pair and direction had less than 5 counts, we summarise all models with a single row, e.g. (1-4).

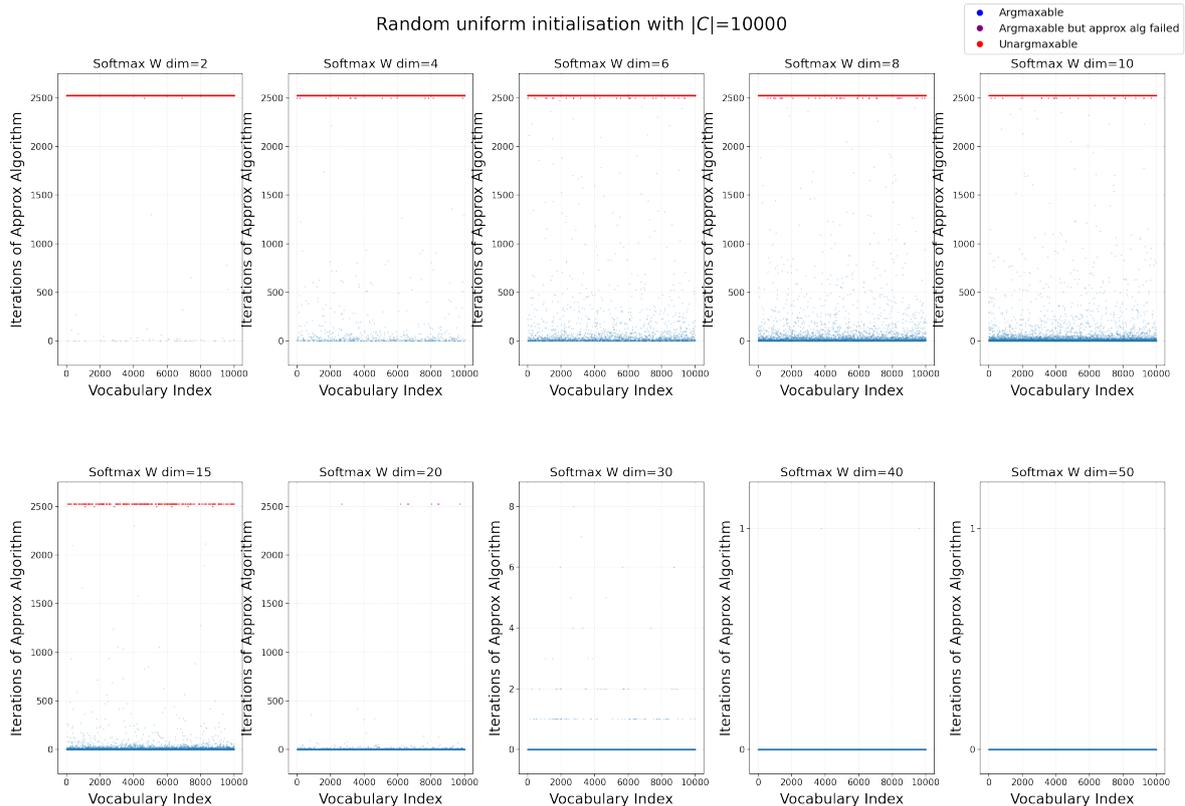


Figure 10: Illustration of Softmax weight dimensionality affecting the number of unargmaxable tokens when weights are randomly initialised for a vocabulary of 10000. The Softmax weights and bias term are initialised using a uniform  $U(-1, 1)$  distribution. Unargmaxable tokens are unlikely to occur as we increase the dimensionality of the weight vectors. This can be seen in the subplots from top-left to bottom-right as we increase the dimensionality. Moreover, the braid reflect approximate algorithm fails less and needs less iterations to find an input that proves a token is argmaxable. For example, for the bottom right two figures most tokens are shown to be argmaxable with 1 or 0 iterations.

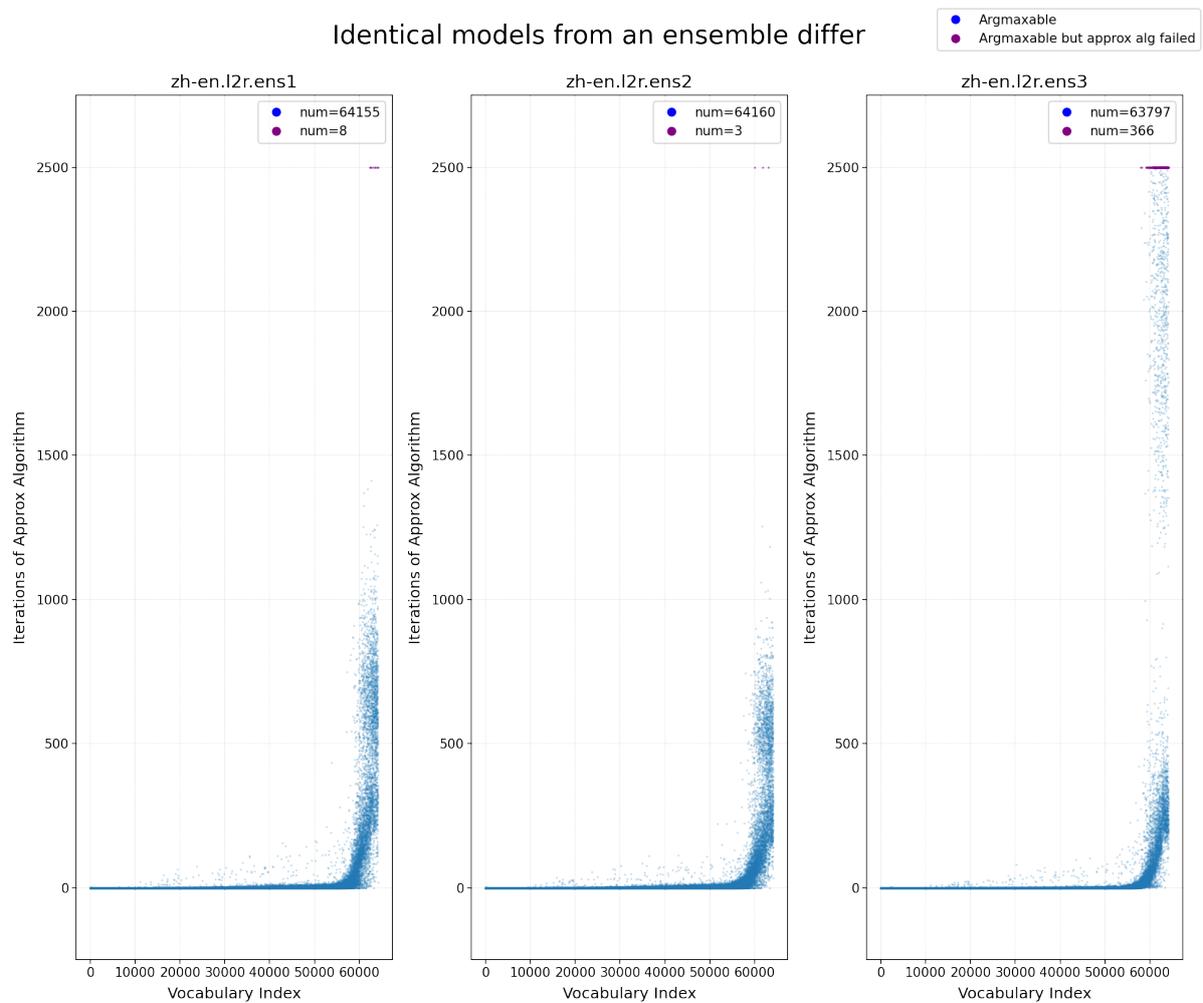
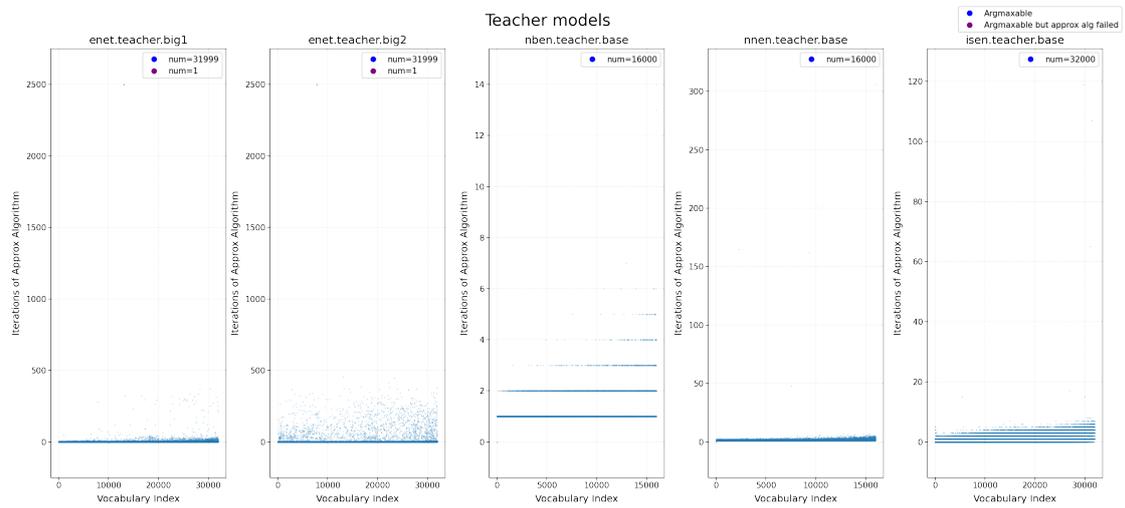
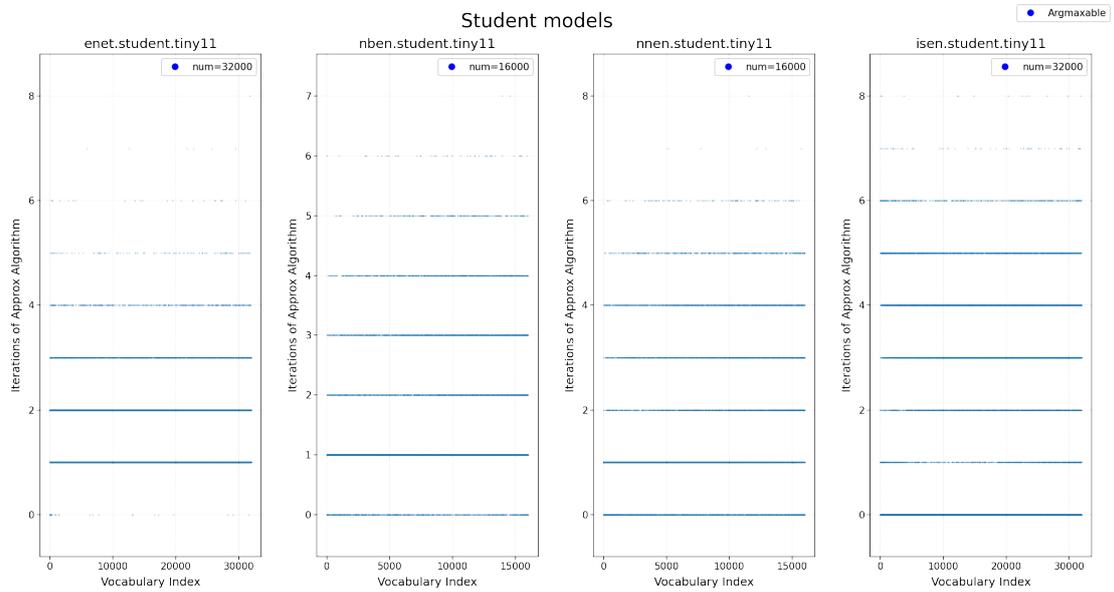


Figure 11: Models from an ensemble can differ a lot in how easy they are to scan for unargmaxable tokens despite their difference being solely the random seed used in initialisation. As can be seen, the right-most figure has 366 vocabulary tokens that are argmaxable but the approximate algorithm fails to find a solution, compared to 8 and 3 for the other two models.



(a) The approximate algorithm needs more iterations to show that the tokens of teacher models are argmaxable despite the dimensionality of the Softmax weights being larger than the student models.



(b) Student models can easily be shown to have argmaxable tokens.

Figure 12: Number of iterations of the approximate algorithm needed to show that a vocabulary token is argmaxable.

## F Activation Range of Softmax Layer Inputs

Neural network activations are bounded in magnitude in practice, since larger activations can lead to larger gradients and instability during training. In this work, we made the assumption that the Softmax layer inputs  $\mathbf{x}$  are bounded within a range for all dimensions:  $-100 \leq \mathbf{x} \leq 100$ . Below we provide some supporting empirical evidence that this assumption is reasonable.

We checked this assumption on 2 Helsinki NLP OPUS models for en-ru and bg-en, which were found to have unargmaxable tokens. We took 10 million sentence pairs from OPUS as released in [Tiedemann \(2020\)](#) for the corresponding language pairs and input them to the corresponding models, decoding using the gold translations. We then recorded the range of the minimum and maximum activation for the Softmax layer inputs.

Since our assumption is that all 512 dimensions are bounded between  $-100$  and  $100$ , we focus on the range of the minimum and maximum activation for each output token across all dimensions. We therefore calculate a 99 percentile for the min and max activation per token across all dimensions as well as the overall min and max activations overall. The results can be seen in [Table 9](#), from which we can see that for these two models our assumption holds for all activations produces for 10 million sentences and the percentiles show that more than 99% of the extreme values fall within the  $[-50, 50]$  range.

| model | min range       | max range      | min    | max   |
|-------|-----------------|----------------|--------|-------|
| bg-en | $[-37.5, -9.4]$ | $[12.1, 40.3]$ | -57.47 | 58.87 |
| en-ru | $[-41.6, -9.9]$ | $[10.9, 36.4]$ | -95.4  | 94.4  |

Table 9: Range of activations for Softmax inputs as calculated on 10 million sentence pairs from OPUS. Ranges are 99 percentiles and min and max are the largest activation across all dimensions for all sentences.