

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049

# Compression of Neural Machine Translation Models via Pruning

Anonymous ACL submission

## Abstract

Neural Machine Translation (NMT), like many other deep learning domains, typically suffers from over-parameterization, resulting in large storage sizes. This paper examines three simple magnitude-based pruning schemes to compress NMT models, namely *class-blind*, *class-uniform*, and *class-distribution*, which differ in terms of how pruning thresholds are computed for the different classes of weights in the NMT architecture. We demonstrate the efficacy of weight pruning as a compression technique for a state-of-the-art NMT system. We show that an NMT model with over 200 million parameters can be pruned by 40% with very little performance loss as measured on the WMT'14 English-German translation task. This sheds light on the distribution of redundancy in the NMT architecture. Our main result is that with *retraining*, we can recover and even surpass the original performance with an 80%-pruned model.

## 1 Introduction

Neural Machine Translation (NMT) is a simple new architecture for translating texts from one language into another (Sutskever et al., 2014; Cho et al., 2014). NMT is a single deep neural network that is trained end-to-end, holding several advantages such as the ability to capture long-range dependencies in sentences, and generalization to unseen texts. Despite being relatively new, NMT has already achieved state-of-the-art translation results for several language pairs including English-French (Luong et al., 2015b), English-German (Jean et al., 2015a; Luong et al., 2015a; Luong and Manning, 2015; Sennrich et al., 2015), English-

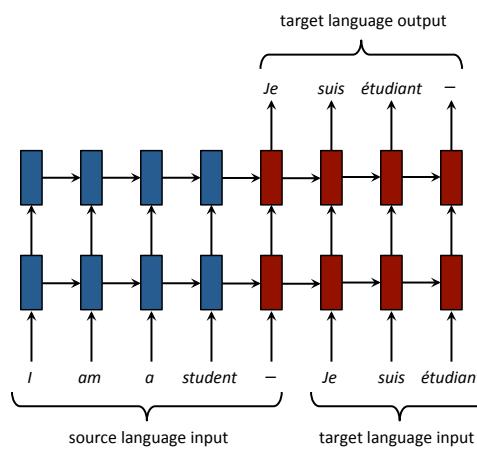


Figure 1: A simplified diagram of NMT.

Turkish (Sennrich et al., 2015), and English-Czech (Jean et al., 2015b; Luong and Manning, 2016). Figure 1 gives an example of an NMT system.

While NMT has a significantly smaller memory footprint than traditional phrase-based approaches (which need to store gigantic phrase-tables and language models), the model size of NMT is still prohibitively large for mobile devices. For example, a recent state-of-the-art NMT system requires over 200 million parameters, resulting in a storage size of hundreds of megabytes (Luong et al., 2015a). Though the trend for bigger and deeper neural networks has brought great progress, it has also introduced over-parameterization, resulting in long running times, overfitting, and the large storage size discussed above. Thus a solution to the over-parameterization problem could potentially aid all three issues.

**Our contribution.** In this paper we investigate the efficacy of weight pruning for NMT as a means of compression. We show that despite its simplicity, magnitude-based pruning with re-training is highly effective, and we compare three

050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099

100 magnitude-based pruning schemes—*class-blind*,  
 101 *class-uniform* and *class-distribution*. Though re-  
 102 cent work has chosen to use the latter two, we  
 103 find the first and simplest scheme—*class-blind*—  
 104 the most successful. We are able to prune 40% of  
 105 the weights of a state-of-the-art NMT system with  
 106 negligible performance loss, and by adding a re-  
 107 training phase after pruning, we can prune 80%  
 108 with no performance loss. Our pruning experi-  
 109 ments also reveal some patterns in the distribution  
 110 of redundancy in NMT. In particular we find that  
 111 higher layers, attention and softmax weights are  
 112 the most important, while lower layers and the em-  
 113 bedding weights hold a lot of redundancy. For the  
 114 Long Short-Term Memory (LSTM) architecture,  
 115 we find that at lower layers the parameters for the  
 116 input are most crucial, but at higher layers the pa-  
 117 rameters for the gates also become important.

## 2 Related Work

120 Pruning the parameters from a neural network,  
 121 referred to as *weight pruning* or *network prun-  
 122 ing*, is a well-established idea though it can be  
 123 implemented in many ways. Among the most  
 124 popular are the Optimal Brain Damage (OBD)  
 125 (Le Cun et al., 1989) and Optimal Brain Sur-  
 126 geon (OBS) (Hassibi and Stork, 1993) techniques,  
 127 which involve computing the Hessian matrix of  
 128 the loss function with respect to the parameters,  
 129 in order to assess the *saliency* of each parame-  
 130 ter. Parameters with low saliency are then pruned  
 131 from the network and the remaining sparse net-  
 132 work is retrained. Both OBD and OBS were  
 133 shown to perform better than the so-called ‘naive  
 134 magnitude-based approach’, which prunes pa-  
 135 rameters according to their magnitude (deleting pa-  
 136 rameters close to zero). However, the high com-  
 137 putational complexity of OBD and OBS compare  
 138 unfavorably to the computational simplicity of the  
 139 magnitude-based approach, especially for large  
 140 networks (Augasta and Kathirvalavakumar, 2013).

141 In recent years, the deep learning renaissance  
 142 has prompted a re-investigation of network prun-  
 143 ing for modern models and tasks. Magnitude-  
 144 based pruning (with iterative retraining) has  
 145 yielded strong results for Convolutional Neural  
 146 Nets (CNNs) performing visual tasks. (Collins  
 147 and Kohli, 2014) prune 75% of AlexNet pa-  
 148 rameters with small accuracy loss on the ImageNet  
 149 task, while (Han et al., 2015b) prune 89% of  
 AlexNet parameters with no accuracy loss on the

ImageNet task.

150 Other approaches focus on pruning neurons  
 151 rather than parameters, via sparsity-inducing regu-  
 152 larizers (Murray and Chiang, 2015) or ‘wiring to-  
 153 gether’ pairs of neurons with similar input weights  
 154 (Srinivas and Babu, 2015). These approaches  
 155 are much more constrained than weight-pruning  
 156 schemes; they necessitate finding entire zero rows  
 157 of weight matrices, or (near-) identical pairs of  
 158 rows, before a single neuron can be pruned. By  
 159 contrast weight-pruning approaches allow weights  
 160 to be pruned freely and independently of each  
 161 other. The neuron-pruning approach of (Srinivas  
 162 and Babu, 2015) was shown to perform poorly  
 163 (it suffered performance loss after removing only  
 164 35% of AlexNet parameters) compared to the  
 165 weight-pruning approach of (Han et al., 2015b).  
 Though (Murray and Chiang, 2015) demonstrates  
 166 neuron-pruning for language modeling as part of  
 167 a (non-neural) Machine Translation pipeline, their  
 168 approach is more geared towards architecture se-  
 169 lection than compression.

170 There are many other compression techniques  
 171 for neural networks, including approaches based  
 172 on low-rank approximations for weight matrices  
 173 (Jaderberg et al., 2014; Denton et al., 2014),  
 174 or weight sharing via hash functions (Chen et al.,  
 175 2015). Several methods involve reducing the pre-  
 176 cision of the weights or activations (Courbariaux  
 177 et al., 2015), sometimes in conjunction with spe-  
 178 cialized hardware (Gupta et al., 2015), or even us-  
 179 ing binary weights (Lin et al., 2016). The ‘knowl-  
 180 edge distillation’ technique of (Hinton et al., 2015)  
 181 involves training a small ‘student’ network on the  
 182 soft outputs of a large ‘teacher’ network. Some  
 183 approaches use a sophisticated pipeline of several  
 184 techniques to achieve impressive feats of compres-  
 185 sion (Han et al., 2015a; Iandola et al., 2016).

186 Most of the above work has focused on com-  
 187 pressing CNNs for vision tasks. We extend the  
 188 magnitude-based pruning approach of (Han et al.,  
 189 2015b) to recurrent neural networks (RNN), in  
 190 particular LSTM architectures for NMT, and to  
 191 our knowledge we are the first to do so. There has  
 192 been some recent work on compression for RNNs  
 193 (Lu et al., 2016; Prabhavalkar et al., 2016), but it  
 194 focuses on other, non-pruning compression tech-  
 195 niques. Nonetheless, our general observations on  
 196 the distribution of redundancy in a LSTM are cor-  
 197 roborated by (Lu et al., 2016).

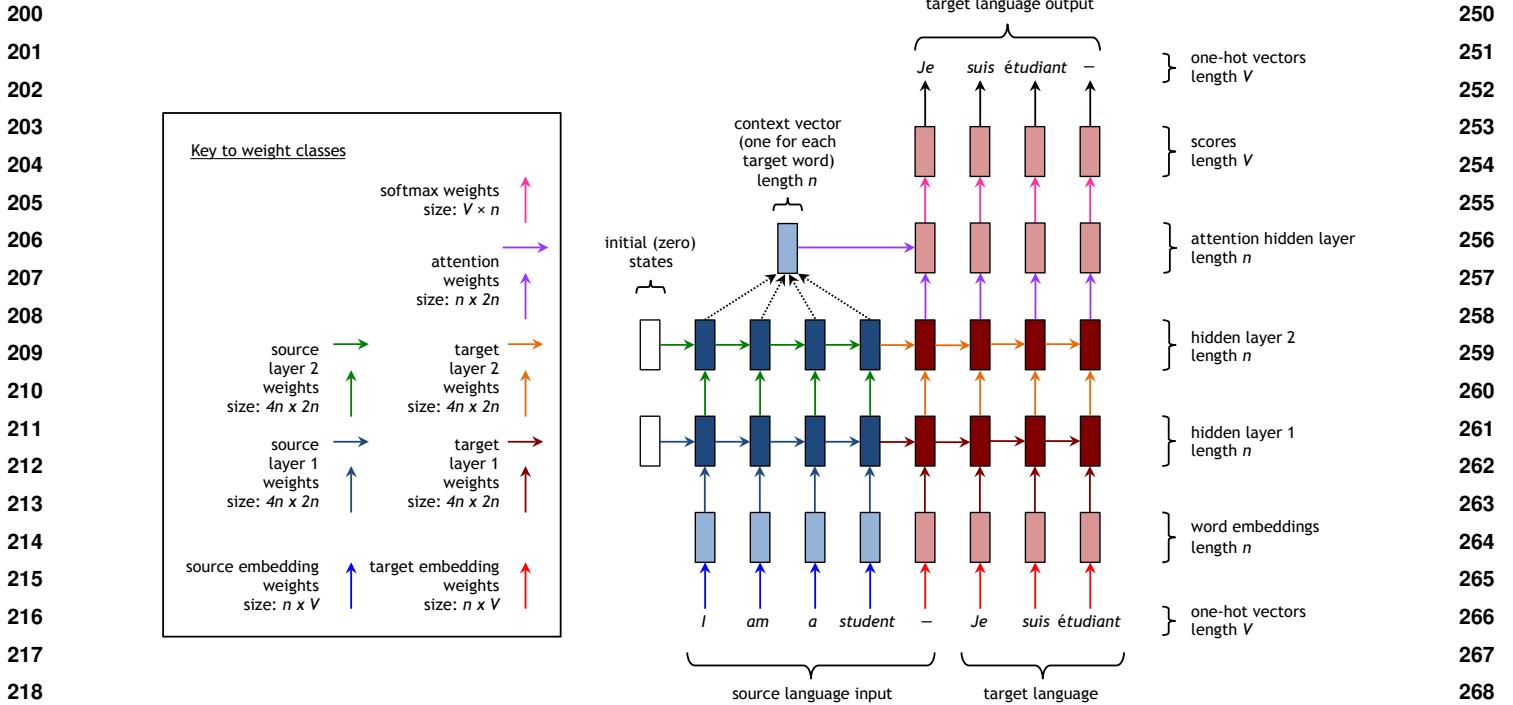


Figure 2: NMT architecture. This example has two layers, but our system has four. The different weight classes are indicated by arrows of different color (the black arrows in the top right represent simply choosing the highest-scoring word, and thus require no parameters). Best viewed in color.

### 3 Our Approach

We first give a brief overview of Neural Machine Translation before delving into a model architecture of interest, the deep multi-layer recurrent model with LSTM. We then explain the different types of NMT weights together with our approaches to pruning and retraining.

#### 3.1 Neural Machine Translation

Neural machine translation aims to directly model the conditional probability  $p(y|x)$  of translating a source sentence,  $x_1, \dots, x_n$ , to a target sentence,  $y_1, \dots, y_m$ . It accomplishes this goal through an *encoder-decoder* framework (Kalchbrenner and Blunsom, 2013; Sutskever et al., 2014; Cho et al., 2014). The *encoder* computes a representation  $s$  for each source sentence. Based on that source representation, the *decoder* generates a translation, one target word at a time, and hence, decomposes the log conditional probability as:

$$\log p(y|x) = \sum_{t=1}^m \log p(y_t|y_{<t}, s) \quad (1)$$

Most NMT work uses RNNs, but approaches differ in terms of: (a) architecture, which can be unidirectional, bidirectional, or deep multi-layer RNN; and (b) RNN type, which can be

Long Short-Term Memory (LSTM) (Hochreiter and Schmidhuber, 1997) or the gated recurrent unit (Cho et al., 2014).

In this work, we specifically consider the *deep multi-layer recurrent* architecture with *LSTM* as the hidden unit type. Figure 1 illustrates an instance of that architecture during training in which the source and target sentence pair are input for supervised learning. During testing, the target sentence is not known in advance; instead, the most probable target words predicted by the model are fed as inputs into the next timestep. The network stops when it emits the end-of-sentence symbol—a special ‘word’ in the vocabulary, represented by a dash in Figure 1.

#### 3.2 Understanding NMT Weights

We show in Figure 2 the same system in more detail, highlighting the different types of parameters, or weights, in the model. We will go through the architecture from bottom to top. First, a vocabulary is chosen for each language, assuming that the top  $V$  frequent words are selected. Thus, every word in the source or target vocabulary can be represented by a one-hot vector of length  $V$ . The source input sentence and target input sentence, represented as a sequence of one-hot vec-

300       tors, are transformed into a sequence of word em-  
 301       beddings by the *embedding* weights. These em-  
 302       bedding weights, which are learned during train-  
 303       ing, are different for the source words and the tar-  
 304       get words. The word embeddings are vectors of  
 305       length  $n$ , the *dimension* of the network.

306       The word embeddings are then fed as input into  
 307       the main network, which consists of two multi-  
 308       layer RNNs ‘stuck together’—an encoder for the  
 309       source language and a decoder for the target lan-  
 310       guage, each with their own weights. The *feed-  
 311       forward* (vertical) weights connect the hidden unit  
 312       from the layer below to the upper RNN block, and  
 313       the *recurrent* (horizontal) weights connect the hid-  
 314       den unit from the previous time-step RNN block to  
 315       the current time-step RNN block.

316       The hidden state at the top layer of the decoder  
 317       is fed through an *attention* layer, which guides the  
 318       translation by ‘paying attention’ to relevant parts  
 319       of the source sentence; for more information see  
 320       Section 3 of (Luong et al., 2015a). Finally, for  
 321       each target word, the top layer hidden unit is trans-  
 322       formed by the *softmax* weights into a score vector  
 323       of length  $V$ . The target word with the highest  
 324       score is selected as the output translation.

325       **Weight Subgroups in LSTM** – For the afore-  
 326       mentioned RNN block, we choose to use LSTM  
 327       as a hidden unit type. To facilitate our discussion  
 328       later on the different subgroups of weights within  
 329       LSTM, we first revise details of an LSTM sug-  
 330       gested by Zaremba et al. (2014) as follows:

$$\begin{pmatrix} i \\ f \\ o \\ \hat{h} \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} T_{4n,2n} \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix} \quad (2)$$

$$c_t^l = f \circ c_{t-1}^l + i \circ \hat{h} \quad (3)$$

$$h_t^l = o \circ \tanh(c_t^l) \quad (4)$$

331       Here, each LSTM block at time  $t$  and layer  $l$  com-  
 332       putes as outputs a pair of hidden and memory vec-  
 333       tors ( $h_t^l$ ,  $c_t^l$ ) given the previous pair ( $h_{t-1}^l$ ,  $c_{t-1}^l$ )  
 334       and an input vector  $h_t^{l-1}$  (either from the below  
 335       LSTM block or the embedding weights if  $l = 1$ ).  
 336       All of these vectors have dimensions of  $n$ .

337       The core of an LSTM block is the weight matrix  
 338        $T_{4n,2n}$  of size  $4n \times 2n$ . This matrix can be decom-  
 339       posed into 8 subgroups that are responsible for the  
 340       interactions between  $\{\text{input gate } i, \text{forget gate } f,$   
 341        $\text{output gate } o, \text{input signal } \hat{h}\} \times \{\text{feed-forward in-}$   
 342       put }  $h_t^{l-1}$ , recurrent input  $h_{t-1}^l\}$ .

### 3.3 Pruning Schemes

We follow the general magnitude-based approach of (Han et al., 2015b), which consists of pruning weights with smallest absolute value. However, we question the authors’ pruning scheme with respect to the different weight classes, and experiment with three pruning schemes. Suppose we wish to prune  $x\%$  of the total parameters in the model. How do we distribute the pruning over the different weight classes (illustrated in Figure 2) of our model? We propose to examine three different pruning schemes:

1. *Class-blind*: Take all parameters, sort them by magnitude and prune the  $x\%$  with smallest magnitude, regardless of weight class. (So some classes are pruned proportionally more than others).
2. *Class-uniform*: Within each class, sort the weights by magnitude and prune the  $x\%$  with smallest magnitude. (So all classes have exactly  $x\%$  of their parameters pruned).
3. *Class-distribution*: For each class  $c$ , weights with magnitude less than  $\lambda\sigma_c$  are pruned. Here,  $\sigma_c$  is the standard deviation of that class and  $\lambda$  is a universal parameter chosen such that in total,  $x\%$  of all parameters are pruned. This is used by (Han et al., 2015b).

All these schemes have their seeming advantages. Class-blind pruning is the simplest and adheres to the principle that pruning weights (or equivalently, setting them to zero) is least damaging when those weights are small, regardless of their locations in an architecture. Class-uniform pruning and class-distribution pruning both seek to prune proportionally within each weight class, either absolutely, or relative to the standard deviation of that class. We find that class-blind pruning outperforms both other schemes (see Section 4.1).

### 3.4 Retraining

In order to prune NMT models aggressively without performance loss, we retrain our pruned networks. In our implementation, we keep “mask” matrices, which represent the sparse structure of a network, so as to ignore weights at pruned locations. We detail in Section 4.1 a successful “formula” to retrain pruned NMT models.

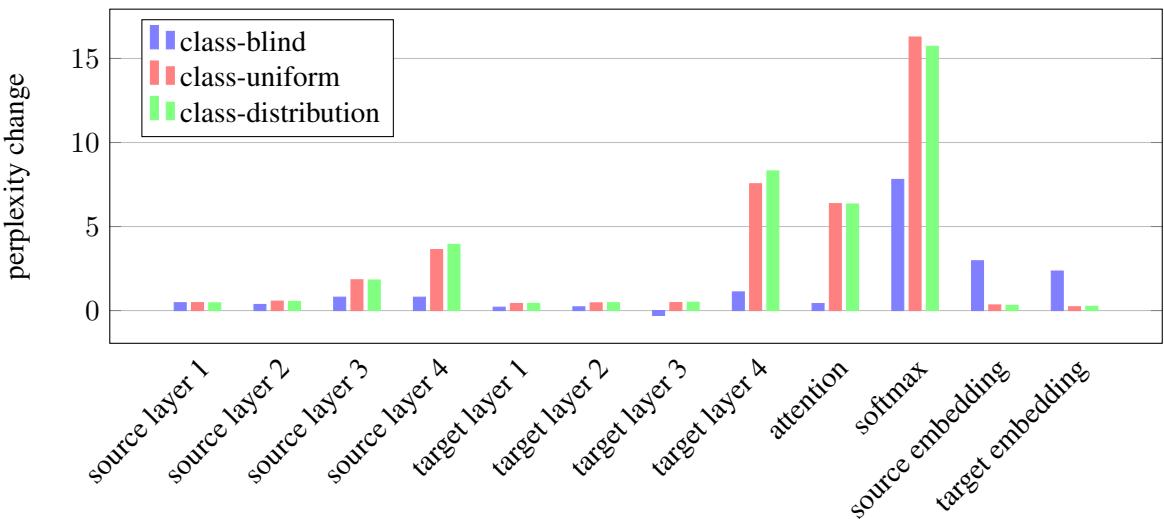


Figure 3: ‘Breakdown’ of performance loss (i.e., perplexity increase) by weight class, when pruning 90% of weights using each of the three pruning schemes. Each of the first eight classes have 8 million weights, attention has 2 million, and the last three have 50 million weights each.

## 4 Experiments

We evaluate the effectiveness of our pruning approaches on a state-of-the-art NMT model.<sup>1</sup> Specifically, an attention-based English-German NMT system from (Luong et al., 2015a) is considered. Training data was obtained from WMT’14 consisting of 4.5M sentence pairs (116M English words, 110M German words). For more details on training hyperparameters, we refer readers to Section 4.1 of (Luong et al., 2015a). All models are tested on newstest2014 (2737 sentences). The model achieves a perplexity of 6.1 and a BLEU score of 20.5 (after unknown word replacement).<sup>2</sup>

When *retraining* pruned NMT systems, we use the following settings: (a) we start with a smaller learning rate of 0.5 (the original model uses a learning rate of 1.0), (b) we train for fewer epochs, 4 instead of 12, using plain SGD, (c) a simple learning rate schedule is employed; after 2 epochs, we begin to halve the learning rate every half an epoch, and (d) all other hyperparameters are the same, such as mini-batch size 128, maximum gradient norm 5, and dropout with probability 0.2.

### 4.1 Comparing pruning schemes

Despite its simplicity, we observe in Figure 4 that *class-blind* pruning outperforms both other

<sup>1</sup>We thank the authors of (Luong et al., 2015a) for providing their trained models and assistance in using the codebase at <https://github.com/lmthang/nmt.matlab>.

<sup>2</sup>The performance of this model is reported under row *global (dot)* in Table 4 of (Luong et al., 2015a).

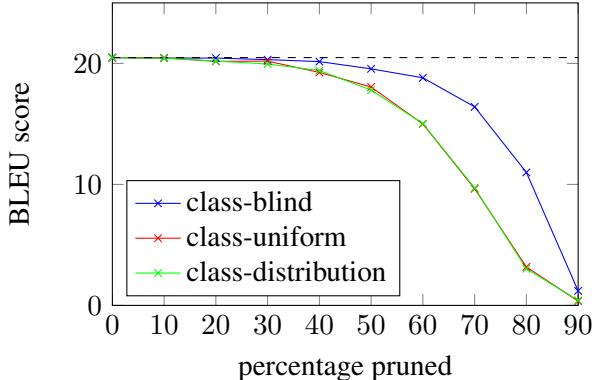


Figure 4: Effects of different pruning schemes.

schemes in terms of translation quality at all pruning percentages. Furthermore, to explain the poor performance of class-uniform and class-distribution pruning, for each of the three pruning schemes, we pruned each class separately and recorded the effect on performance (as measured by perplexity). Figure 3 shows that with class-uniform pruning, the overall performance loss is caused disproportionately by a few classes: target layer 4, attention and softmax weights. Looking at Figure 5, we see that the most damaging classes to prune also tend to be those with weights of greater magnitude—these classes have much larger weights than others at the same percentile, so pruning them under the class-uniform pruning scheme is more damaging. The situation is similar for class-distribution pruning.

By contrast, Figure 3 shows that under class-

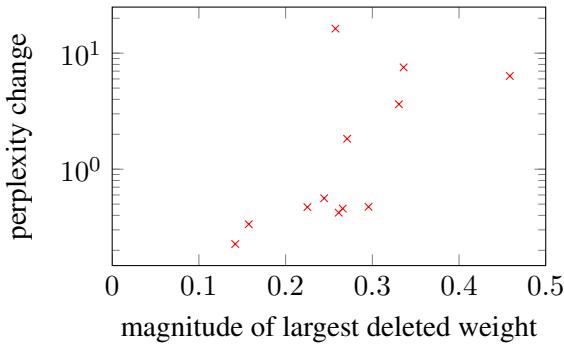


Figure 5: Magnitude of largest deleted weight vs. perplexity change, for the 12 different weight classes when pruning 90% of parameters by class-uniform pruning.

blind pruning, the damage caused by pruning softmax, attention and target layer 4 weights is greatly decreased, and the contribution of each class towards the performance loss is overall more uniform. In fact, the distribution begins to reflect the number of parameters in each class—for example, the source and target embedding classes have larger contributions because they have more weights. We use only class-blind pruning for the rest of the experiments.

Figure 3 also reveals some interesting information about the distribution of redundancy in NMT architectures—namely it seems that higher layers are more important than lower layers, and that attention and softmax weights are crucial. We will explore the distribution of redundancy further in Section 4.3.

#### 4.2 Pruning and retraining

Pruning has an immediate negative impact on performance (as measured by BLEU score on the validation set) that is exponential in pruning percentage; this is demonstrated by the blue line in Figure 6. However we find that up to about 40% pruning, performance is mostly unaffected, indicating a large amount of redundancy and over-parameterization in NMT.

We now consider the effect of retraining pruned models. The orange line in Figure 6 shows that after retraining the pruned models, baseline performance (20.48 BLEU) is both recovered and improved upon, up to 80% pruning (20.91 BLEU), with only a small performance loss at 90% pruning (20.13 BLEU). This may seem surprising, as we might not expect a sparse model to significantly

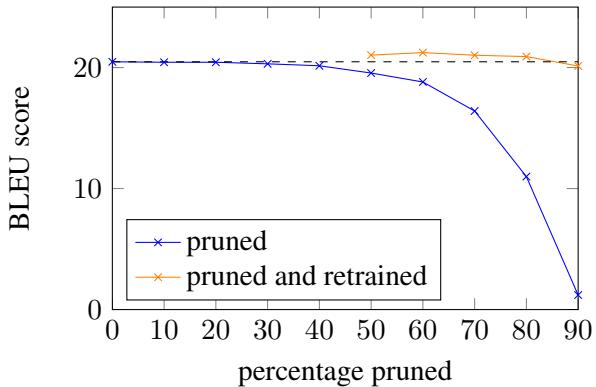


Figure 6: Performance of pruned models, immediately after pruning and after retraining.

out-perform a model with five times as many parameters. There are several possible explanations, two of which are given below.

Firstly, we found that the less-pruned models perform better on the training set than the validation set, whereas the more-pruned models have closer performance on the two sets. This indicates that pruning has a regularizing effect on the re-training phase, though clearly more is not always better, as the 50% pruned and retrained model performs better than the 90% pruned and retrained model. Nonetheless, this regularization effect may explain why the pruned and retrained models outperform the baseline.

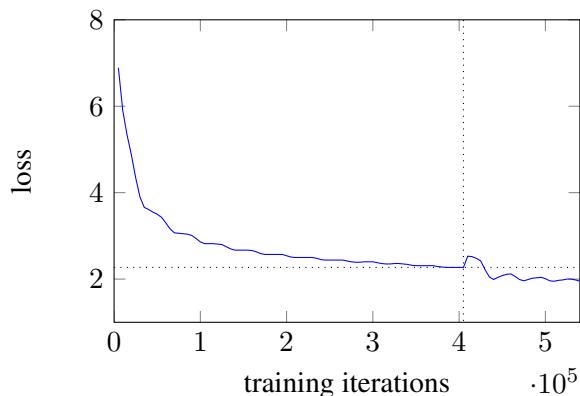


Figure 7: The loss function during training, pruning and retraining. The vertical dotted line marks the point when 80% of the parameters are pruned. The horizontal dotted line marks the best performance of the unpruned baseline.

Alternatively, pruning may serve as a means to escape a local optimum. Figure 7 shows the loss function over time during the training, pruning and retraining process. During the original training

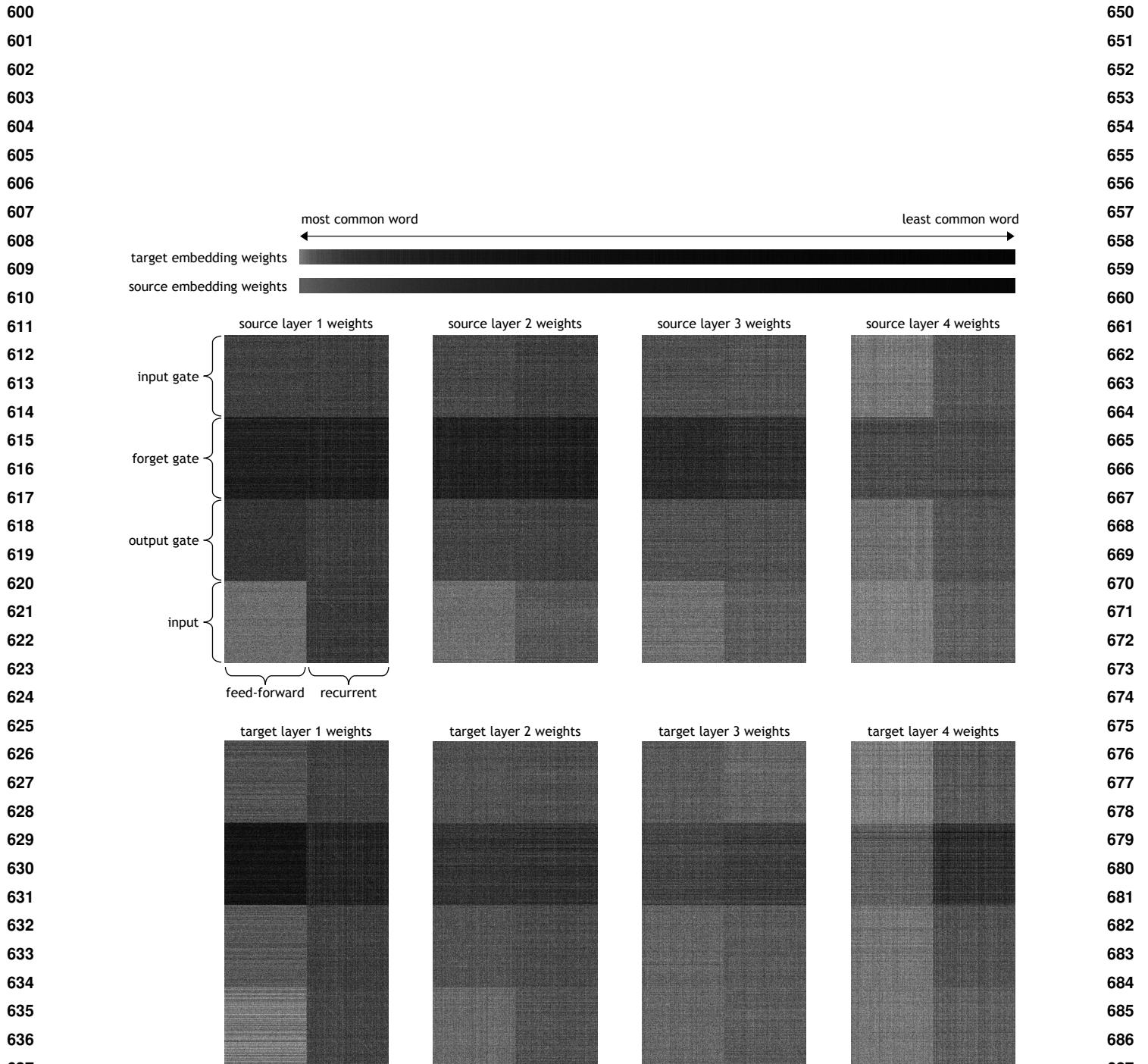


Figure 8: Graphical representation of the location of small weights in the model. Black pixels represent weights with absolute size in the bottom 80%; white pixels represent those with absolute size in the top 20%. Equivalently, these pictures illustrate which parameters remain after pruning 80% using our class-blind pruning scheme.

700 process, the loss curve flattens out and seems to  
 701 converge (note that we use early stopping to obtain  
 702 our baseline model, so the original model was  
 703 trained for longer than shown in Figure 7). Pruning  
 704 causes an immediate increase in the loss function,  
 705 but enables further gradient descent, allowing  
 706 the retraining process to find a new, better local  
 707 optimum. It seems that the disruption caused by  
 708 pruning is beneficial in the long-run.

### 709 4.3 Distribution of redundancy in NMT

710 We visualize in Figure 8 the redundancy struc-  
 711 ture of our NMT baseline model. *Black* pixels  
 712 represent weights near to zero; *white* pixels rep-  
 713 resent larger ones. First we consider the embed-  
 714 ding weight matrices, whose columns correspond  
 715 to words in the vocabulary. Unsurprisingly, in  
 716 Figure 8, we see that the parameters correspond-  
 717 ing to the less common words are more dispens-  
 718 able. In fact, at the 80% pruning rate, for 100 un-  
 719 common source words and 1194 uncommon target  
 720 words, we delete *all* parameters corresponding to  
 721 that word. This is not quite the same as remov-  
 722 ing the word from the vocabulary—true out-of-  
 723 vocabulary words are mapped to the embedding  
 724 for the ‘unknown word’ symbol, whereas these  
 725 ‘pruned-out’ words are mapped to a zero embed-  
 726 ding. However in the original unpruned model  
 727 these uncommon words already had near-zero em-  
 728 beddings, indicating that the model was unable to  
 729 learn sufficiently distinctive representations.  
 730

731 Returning to Figure 8, now look at the eight  
 732 weight matrices for the source and target connec-  
 733 tions at each of the four layers. Each matrix corre-  
 734 sponds to the  $4n \times 2n$  matrix  $T_{4n,2n}$  in Equation  
 735 (1). In all eight matrices, we observe—as does  
 736 (Lu et al., 2016)—that the weights connecting to  
 737 the input  $\hat{h}$  are most crucial, followed by the in-  
 738 put gate  $i$ , then the output gate  $o$ , then the forget  
 739 gate  $f$ . This is particularly true of the lower lay-  
 740 ers, which focus primarily on the input  $\hat{h}$ . How-  
 741 ever for higher layers, especially on the target side,  
 742 weights connecting to the gates are as important as  
 743 those connecting to the input  $\hat{h}$ . The gates repre-  
 744 sent the LSTM’s ability to add to, delete from or  
 745 retrieve information from the memory cell. Figure  
 746 8 therefore shows that these sophisticated memory  
 747 cell abilities are most important at the *end* of the  
 748 NMT pipeline (the top layer of the decoder). This  
 749 is reasonable, as we expect higher-level features to  
 be learned later in a deep learning pipeline.

750 We also observe that for lower layers, the feed-  
 751 forward input is much more important than the re-  
 752 current input, whereas for higher layers the recur-  
 753 rent input becomes more important. This makes  
 754 sense: lower layers concentrate on the low-level  
 755 information from the current word embedding (the  
 756 feed-forward input), whereas higher layers make  
 757 use of the higher-level representation of the sen-  
 758 tence so far (the recurrent input).

759 Lastly, on close inspection, we notice several  
 760 white diagonals emerging within the subsquares  
 761 of the matrices in Figure 8, indicating that even  
 762 without initializing the weights to identity mat-  
 763 rices, an identity-like weight matrix is learned. At  
 764 higher pruning percentages, these diagonals be-  
 765 come more pronounced.

## 766 5 Future Work

767 The pruning method described in (Han et al.,  
 768 2015b) includes several iterations of pruning and  
 769 retraining. Implementing this for NMT would  
 770 likely result in further compression and per-  
 771 formance improvements. If possible it would be  
 772 highly valuable to exploit the sparsity of the  
 773 pruned models to speed up training and run-  
 774 time, perhaps through sparse matrix representa-  
 775 tions and multiplications. Though we have found  
 776 magnitude-based pruning to perform very well, it  
 777 would be instructive to revisit the original claim  
 778 that other pruning methods (for example Optimal  
 779 Brain Damage and Optimal Brain Surgery) are  
 780 more principled, and perform a comparative study.

## 781 6 Conclusion

782 We have shown that weight pruning with retrain-  
 783 ing is a highly effective method of compression  
 784 and regularization on a state-of-the-art NMT sys-  
 785 tem, compressing the model to 20% of its size with  
 786 no loss of performance. Though we are the first to  
 787 apply compression techniques to NMT, we obtain  
 788 a similar degree of compression to other current  
 789 work on compressing state-of-the-art deep neural  
 790 networks, with an approach that is simpler than  
 791 most. We have found that the absolute size of pa-  
 792 rameters is of primary importance when choosing  
 793 which to prune, leading to an approach that is ex-  
 794 tremely simple to implement, and can be applied  
 795 to any neural network. Lastly, we have gained  
 796 insight into the distribution of redundancy in the  
 797 NMT architecture.

800	<b>References</b>	850
801	M Gethsiyal Augusta and T Kathirvalavakumar. 2013.	851
802	Pruning algorithms of neural networksa comparative	852
803	study. <i>Central European Journal of Computer Science</i> , 3(3):105–115.	853
804		
805	Wenlin Chen, James T Wilson, Stephen Tyree, Kilian Q	854
806	Weinberger, and Yixin Chen. 2015. Compressing	855
807	neural networks with the hashing trick. <i>arXiv preprint arXiv:1504.04788</i> .	856
808		857
809	Kyunghyun Cho, Bart van Merriënboer, Caglar Gul-	858
810	cehre, Fethi Bougares, Holger Schwenk, and Yoshua	859
811	Bengio. 2014. Learning phrase representations	860
812	using RNN encoder-decoder for statistical machine	861
813	translation. In <i>EMNLP</i> .	862
814		863
815	Maxwell D Collins and Pushmeet Kohli. 2014. Mem-	864
816	ory bounded deep convolutional networks. <i>arXiv preprint arXiv:1412.1442</i> .	865
817		866
818	Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre	867
819	David. 2015. Low precision arithmetic for deep	868
820	learning. In <i>International Conference on Learning</i>	869
821	Representations (ICLR) Workshop Contribution.	870
822		871
823	Emily L Denton, Wojciech Zaremba, Joan Bruna, Yann	872
824	LeCun, and Rob Fergus. 2014. Exploiting linear	873
825	structure within convolutional networks for efficient	874
826	evaluation. In <i>NIPS</i> .	875
827		876
828	Suyog Gupta, Ankur Agrawal, Kailash Gopalakrish-	877
829	nan, and Pritish Narayanan. 2015. Deep learn-	878
830	ing with limited numerical precision. <i>arXiv preprint arXiv:1502.02551</i> .	879
831		880
832	Song Han, Huizi Mao, and William J Dally. 2015a.	881
833	Deep compression: Compressing deep neural net-	882
834	works with pruning, trained quantization and huff-	883
835	man coding. In <i>International Conference on Learn-</i>	884
836	ing Representations (ICLR’16 oral).	885
837		886
838	Song Han, Jeff Pool, John Tran, and William Dally.	887
839	2015b. Learning both weights and connections for	888
840	efficient neural network. In <i>Advances in Neural In-</i>	889
841	formation Processing Systems	890
842	, pages 1135–1143.	891
843		892
844	Babak Hassibi and David G Stork. 1993. <i>Second or-</i>	893
845	<i>der derivatives for network pruning: Optimal brain</i>	894
846	<i>surgeon</i> . Morgan Kaufmann.	895
847		896
848	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015.	897
849	Distilling the knowledge in a neural network. In <i>Advances in Neural Information Processing Systems (NIPS) Deep Learning Workshop</i> .	898
850		899
851	Sepp Hochreiter and Jürgen Schmidhuber. 1997.	
852	Long short-term memory. <i>Neural computation</i> ,	
853	9(8):1735–1780.	
854		
855	Forrest N Iandola, Matthew W Moskewicz, Khalid	
856	Ashraf, Song Han, William J Dally, and Kurt	
857	Keutzer. 2016. SqueezeNet: Alexnet-level accuracy	
858	with 50x fewer parameters and < 1mb model size.	
859	<i>arXiv preprint arXiv:1602.07360</i> .	
860		
861	Max Jaderberg, Andrea Vedaldi, and Andrew Zisser-	
862	man. 2014. Speeding up convolutional neural net-	
863	works with low rank expansions. In <i>Advances in</i>	
864	<i>Neural Information Processing Systems (NIPS)</i> .	
865		
866	Sébastien Jean, Kyunghyun Cho, Roland Memisevic,	
867	and Yoshua Bengio. 2015a. On using very large	
868	target vocabulary for neural machine translation. In	
869	<i>ACL</i> .	
870		
871	Sébastien Jean, Orhan Firat, Kyunghyun Cho, Roland	
872	Memisevic, and Yoshua Bengio. 2015b. Montreal	
873	neural machine translation systems for WMT’15. In	
874	<i>WMT</i> .	
875		
876	Nal Kalchbrenner and Phil Blunsom. 2013. Recurrent	
877	continuous translation models. In <i>EMNLP</i> .	
878		
879	Yann Le Cun, John S Denker, and Sara A Solla. 1989.	
880	Optimal brain damage. In <i>Advances in Neural In-</i>	
881	<i>formation Processing Systems</i> .	
882		
883	Zhouhan Lin, Matthieu Courbariaux, Roland Memise-	
884	vic, and Yoshua Bengio. 2016. Neural networks	
885	with few multiplications. In <i>International Confer-</i>	
886	<i>ence on Learning Representations (ICLR)</i> .	
887		
888	Zhiyun Lu, Vikas Sindhwani, and Tara N Sainath.	
889	2016. Learning compact recurrent neural networks.	
890	<i>arXiv preprint arXiv:1604.02594</i> .	
891		
892	Minh-Thang Luong and Christopher D. Manning.	
893	2015. Stanford neural machine translation systems	
894	for spoken language domain. In <i>IWSLT</i> .	
895		
896	Minh-Thang Luong and Christopher D. Manning.	
897	2016. Achieving open vocabulary neural ma-	
898	chine translation with hybrid word-character mod-	
899	els. <i>arXiv preprint arXiv:1604.00788</i> .	

900	Suraj Srinivas and R Venkatesh Babu. 2015. Data-free parameter pruning for deep neural networks. <i>arXiv preprint arXiv:1507.06149.</i>	950 951 952
903	Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural net- works. In <i>NIPS</i> .	953 954 955
906	Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals. 2014. Recurrent neural network regularization. <i>arXiv preprint arXiv:1409.2329.</i>	956 957 958
909		959
910		960
911		961
912		962
913		963
914		964
915		965
916		966
917		967
918		968
919		969
920		970
921		971
922		972
923		973
924		974
925		975
926		976
927		977
928		978
929		979
930		980
931		981
932		982
933		983
934		984
935		985
936		986
937		987
938		988
939		989
940		990
941		991
942		992
943		993
944		994
945		995
946		996
947		997
948		998
949		999