

Neuron Specialization: Leveraging Intrinsic Task Modularity for Multilingual Machine Translation

Anonymous ACL submission

Abstract

001 Training a unified multilingual model promotes
002 knowledge transfer but inevitably introduces
003 *negative interference*. Language-specific mod-
004 eling methods show promise in reducing inter-
005 ference. However, they often rely on heuris-
006 tics to distribute capacity and struggle to fos-
007 ter cross-lingual transfer via isolated modules.
008 In this paper, we explore intrinsic task modu-
009 larity within multilingual networks and lever-
010 age these observations to circumvent interfer-
011 ence under multilingual translation. We show
012 that neurons in the feed-forward layers tend
013 to be activated in a language-specific manner.
014 Meanwhile, these specialized neurons exhibit
015 structural overlaps that reflect language prox-
016 imity, which progress across layers. Based
017 on these findings, we propose *Neuron Special-*
018 *ization*, an approach that identifies specialized
019 neurons to modularize feed-forward layers and
020 then continuously updates them through sparse
021 networks. Extensive experiments show that
022 our approach achieves consistent performance
023 gains over strong baselines with additional anal-
024 yses demonstrating reduced interference and
025 increased knowledge transfer.¹

026 1 Introduction

027 Jointly training multilingual data in a unified
028 model with a shared architecture for different lan-
029 guages has been a trend (Conneau et al., 2020;
030 Le Scao et al., 2022) encouraging knowledge trans-
031 fer across languages, especially for low-resource
032 languages (Johnson et al., 2017; Pires et al., 2019).
033 However, such a training paradigm also leads to
034 *negative interference* due to conflicting optimiza-
035 tion demands (Wang et al., 2020). This interference
036 often causes performance degradation for high-
037 resource languages (Li and Gong, 2021; Pfeiffer
038 et al., 2022) and can be further exacerbated by lim-
039 ited model capacity (Shaham et al., 2023).

¹We release code at <https://anonymous.4open.science/r/NS-3D93>

040 Modular-based methods, such as Language-
041 specific modeling (Zhang et al., 2020b) and
042 adapters (Bapna and Firat, 2019), aim to mitigate
043 interference by balancing full parameter sharing
044 with isolated or partially shared modules (Pfeiffer
045 et al., 2023). However, they heavily depend on
046 heuristics for allocating task-specific capacity and
047 face challenges in enabling knowledge transfer be-
048 tween modules (Zhang et al., 2020a). Specifically,
049 such methods rely on prior knowledge for man-
050 aging parameter sharing such as language-family
051 adapters (Chronopoulou et al., 2023) or directly
052 isolate parameters per language, which impedes
053 transfer (Pires et al., 2023).

054 Research in vision and cognitive science has
055 shown that unified multi-task models may sponta-
056 neously develop task-specific functional specializa-
057 tions for distinct tasks (Yang et al., 2019; Dobs
058 et al., 2022), a phenomenon also observed in
059 mixture of experts Transformer systems (Zhang
060 et al., 2023). These findings suggest that through
061 multi-task training, networks naturally evolve to-
062 wards specialized modularity to effectively man-
063 age diverse tasks, with the ablation of these spe-
064 cialized modules adversely affecting task perfor-
065 mance (Pfeiffer et al., 2023). Despite these insights,
066 exploiting the inherent structural signals for multi-
067 task optimization remains largely unexplored.

068 In this work, we explore the intrinsic task-
069 specific modularity within multi-task networks in
070 Multilingual Machine Translation (MMT), treating
071 each language pair as a separate task. We focus
072 on analyzing the intermediate activations in the
073 Feed-Forward Networks (FFN) where most model
074 parameters reside. To our knowledge, our study is
075 the first to show that neurons activate in a language-
076 specific way, yet they present structural overlaps
077 that indicate language proximity in general. More-
078 over, this pattern evolves across layers in the model,
079 suggesting that neurons consistently transition from
080 language-specific to language-agnostic.

Building on these observations, we introduce *Neuron Specialization*, a novel method that leverages intrinsic task modularity to reduce interference and enhance knowledge transfer. In general, our approach selectively updates the FFN parameters during back-propagation for different tasks to enhance task specificity. Specifically, we first identify task-specific neurons from pre-trained unified translation models, using standard forward-pass validation processes without decoding. We then specifically modularize FFN layers using these specialized neurons and continuously update FFNs via sparse networks.

Extensive experiments on small- (IWSLT) and large-scale EC30 (Tan and Monz, 2023) translation datasets show that our method consistently achieves performance gains over strong baselines with various configs. Moreover, we conduct in-depth analyses to show that our method effectively mitigates interference and enhances knowledge transfer in high and low-resource languages, respectively. Our main contributions are summarized as follows:

- We identify inherent multilingual modularity by showing that neurons activate in a language-specific manner and their overlapping patterns reflect language proximity.
- Building on these findings, we enhance task specificity through sparse FFNs, achieving consistent improvements in translation quality over strong baselines.
- We employ analyses to show that our method effectively reduces interference in high-resource languages and boosts knowledge transfer in low-resource languages.

2 Related Work

Multilingual Interference. Multilingual training enables knowledge transfer but also causes *interference*, largely due to optimization conflicts among various tasks (Wang and Zhang, 2022). Methods alleviating task conflicts hold promise to reduce interference (Wang et al., 2020), yet they show limited effectiveness in practice (Xin et al., 2022). Scaling up model size may reduce interference but leads to overly large models (Chang et al., 2023), with risks of overfitting (Aharoni et al., 2019).

Language-Specific Modeling. Recent methods enhance the unified model by utilizing language-specific (LS) modules such as adapters (Bapna

and Firat, 2019), LS layers (Zhang et al., 2020b; Pires et al., 2023) and LS hidden states (Xie et al., 2021). Although the unified model serves as a common foundation, these methods strictly isolate modules per language. Such designs present no knowledge sharing among modules and thus offer fewer benefits to low-resource languages. Alternatively, approaches like language family adapters Chronopoulou et al. (2023) seek to facilitate sharing among language-specific modules, however, they heavily depend on heuristics such as using priori linguistic knowledge to enable more flexible parameter sharing.

Additionally, these modular-based methods exhibit parameter inefficiency when handling numerous languages, resulting in increased memory requirements and extended inference times (Liao et al., 2023a,b). Similarly, techniques such as parameter differentiation (Wang and Zhang, 2022) and language clustering training (Tan et al., 2019) alleviate interference by expanding the unified model with substantial extra parameters.

Sub-networks in Multi-task Models. The lottery ticket hypothesis (Frankle and Carbin, 2018) states that within dense neural networks, sparse subnetworks can be found with iterative pruning to achieve the original network’s performance. Following this premise, recent studies attempt to isolate sub-networks of a pre-trained unified model that captures task-specific features (Choenni et al., 2023a; Lin et al., 2021; He et al., 2023). Nonetheless, unlike our method that identifies intrinsic modularity within the model, these approaches depend on fine-tuning to extract the task-specific sub-networks. This process may not reflect the original model modularity and also can be particularly resource-consuming for multiple tasks.

Specifically, these methods extract the task-specific sub-networks by fine-tuning the original unified multi-task model on specific tasks, followed by employing pruning to retain only the most changed parameters. We argue that this process faces several issues: 1) The sub-network might be an artifact of fine-tuning, suggesting the original model may not inherently possess such modularity. 2) This is further supported by the observation that different random seeds during fine-tuning lead to varied sub-networks and performance instability (Choenni et al., 2023a). 3) The process is highly inefficient for models covering multiple tasks, as it necessitates separate fine-tuning for each task.

3 Neuron Structural Analysis

Recent work aims to identify a subset of parameters within pre-trained multi-task networks that are sensitive to distinct tasks. This exploration is done by either 1) selecting hidden states that greatly influence task performance (Dobs et al., 2022) or possess high magnitude values (Xie et al., 2021); or 2) fine-tuning the unified model on task-specific data to extract sub-networks (Lin et al., 2021; He et al., 2023; Choenni et al., 2023b). These approaches, however, raise a fundamental question, namely whether the modularity is inherent to the original model, or simply an artifact introduced by network modifications.

In this paper, we perform a thorough identification of task-specific modularity through the lens of neuron behaviors, without altering the original parameters or architectures. We focus on the neurons — the intermediate activations inside the Feed-Forward Networks (FFN) — to investigate if they indicate task-specific modularity features. As FFN neurons are active (>0) or inactive ($=0$) due to the *ReLU* activation function, this binary activation state offers a clear view of their contributions to the network’s output. Intuitively, neurons that remain inactive for one task but show significant activation for another may be indicative of specialization for the latter. More importantly, this approach ensures that both parameters and hidden states remain unchanged, affirming the observed modularity is inherent to the original model.

3.1 Identifying Specialized Neurons

We choose multilingual translation as a testbed, treating each translation direction as a distinct task throughout the paper. We start with a pre-trained multilingual model with d_{ff} as its dimension of the FFN layer. We hypothesize the existence of neuron subsets specialized for each task and describe the identification process of an FFN layer as follows.

Activation Recording. Given a validation dataset D_t for the t -th task, we measure activation frequencies in an FFN layer during validation. For each sample $x_i \in D_t$, we record the state of each neuron after *ReLU*, reflecting whether the neuron is active or inactive to the sample. We use a binary vector $a_i^t \in \mathbb{R}^{d_{ff}}$ to store this neuron state information. Note that this vector aggregates neuron activations for all tokens in the sample by taking the neuron union of them. By further merging all of the binary vectors for all samples

in D_t , an accumulated vector $a^t = \sum_{x_i \in D_t} a_i^t$ can be derived, which denotes the frequency of each neuron being activated during a forward pass given a task-specific dataset D_t .

Neuron Selection. We identify specialized neurons for each task t based on their activation frequency a^t . A subset of neurons S_k^t is progressively selected based on the highest a^t values until reaching a predefined threshold k , where

$$\sum_{i \in S_k^t} a_{(i)}^t \geq k \sum_{i=1}^{d_{ff}} a_{(i)}^t \quad (1)$$

Here, the value $a_{(i)}^t$ is the frequency of the activation at dimension i , and $\sum_{i=1}^{d_{ff}} a_{(i)}^t$ is the total activation of all neurons for an FFN layer. k is a threshold factor, varying from 0% to 100%, indicating the extent of neuron activation deemed necessary for specialization. A lower k value results in higher sparsity in specialized neurons; $k = 0$ means no neuron will be involved, while $k = 100$ fully engages all neurons, the same as utilizing the full capacity of the original model. This dynamic approach emphasizes the collective significance of neuron activations up to a factor of k . In the end, we repeat these processes to obtain the specialized neurons of all FFN layers for each task.

3.2 Analysis on EC30

In this section, we describe how we identify specialized neurons on EC30 (Tan and Monz, 2023), where we train an MMT model covering all directions. EC30 is a multilingual translation benchmark that is carefully designed to consider diverse linguistic properties and real-world data distributions. It collects high to low-resource languages, resulting in 30 diverse languages from 5 language families, allowing us to connect our observations with linguistic properties easily. See Sections 5 for details on data and models.

3.2.1 Neuron Overlaps Reflect Language Proximity

We identified specialized neurons following Section 3.1, while setting the cumulative activation threshold k at 95%. This implies that the set of specialized neurons covers approximately 95% of the total activations. Intuitively, two similar tasks should have a high overlap between their specialized neuron sets. Therefore, we examined the overlaps among specialized neurons across different

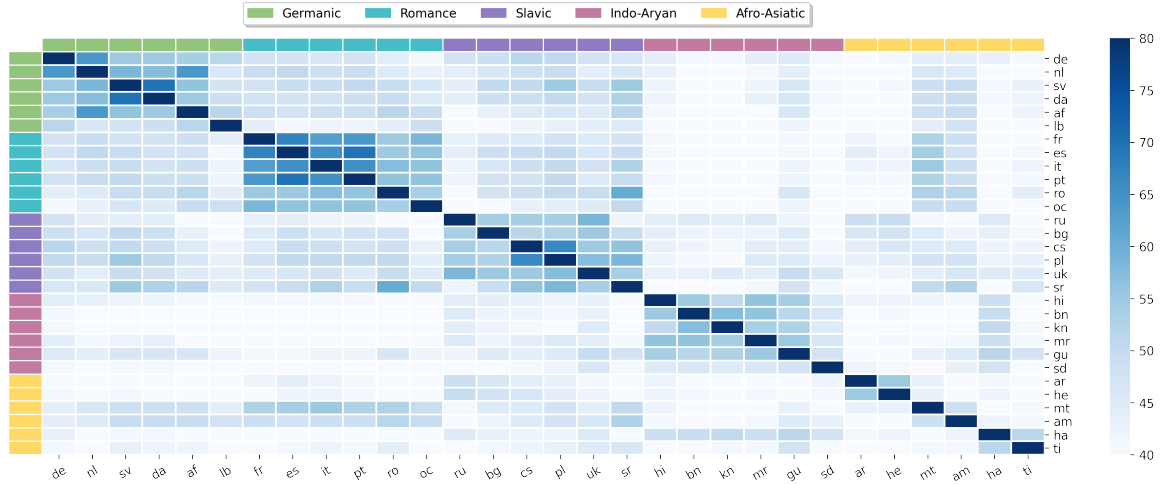


Figure 1: Pairwise Intersection over Union (IoU) scores for specialized neurons extracted from the first decoder FFN layer across all out-of-English translation directions to measure the degree of overlap. Darker cells indicate stronger overlaps, with the color threshold set from 40 to 80 to improve visibility.

tasks by calculating the Intersection over Union (IoU) scores: For task t_i and t_j , with specialized neurons denoted as sets S^i and S^j , their overlap is quantified by $\text{IoU}(S^i, S^j) = \frac{|S^i \cap S^j|}{|S^i \cup S^j|}$.

Figure 1 shows the IoU scores for specialized neurons across different tasks in the first decoder layer. Figures for the other layers can be found in Appendix A.9. We first note a structural separation of neuron overlaps, indicating a preference for language specificity. Notably, neuron overlap across language families is relatively low, a trend more pronounced in encoder layers (Figure 6). Secondly, this structural distinction generally correlates with language proximity as indicated by the clustering pattern in Figure 1. This implies that target languages from the same family are more likely to activate similar neurons in the decoder, even when they use different writing systems, e.g., Arabic (ar) and Hebrew (he). Overlaps also show linguistic traits beyond family ties, exemplified by notable overlaps between Maltese (mt) and languages in the Romance family due to vocabulary borrowing.

3.2.2 The Progression of Neuron Overlaps

To analyze how specialized neuron overlaps across tasks evolve within the model, we visualize the IoU score distribution across layers in Figure 2. For each layer, we compute the pair-wise IoU scores between all possible tasks and then show them in a distribution. Overall, we observe that from shallow to deeper layers, structural distinctions intensify in the decoder (decreasing IoU scores) and weaken in the encoder (increasing IoU scores).

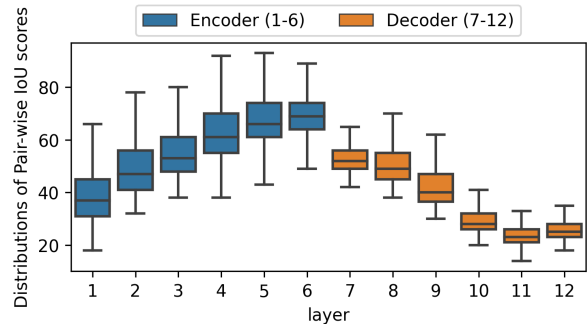


Figure 2: Progression of distribution of IoU scores for specialized neurons across layers on the EC30 dataset. The scores are measured for different source and target languages in the Encoder and Decoder, respectively.

Furthermore, all neuron overlaps increase as we move up the encoder, regardless of whether these tasks are similar or not. This observation may suggest that the neurons in the encoder become more language-agnostic, as they attempt to map different scripts into semantic concepts. As for the Decoder, the model presents intensified modularity in terms of overlaps of specialized neurons. This can be seen by all overlaps becoming much smaller, indicating that neurons behave more separately.

Our findings align with the common assumption about the transformation process in seq-to-seq models. Similarly, Kudugunta et al. (2019) observed that multilingual embeddings gradually, though not perfectly, align within the encoder. However, our research diverges as it focuses on binary neuron activation patterns, rather than high-dimensional embeddings. Moreover, unlike them, we show that our findings can be leveraged to improve MMT.

4 Neuron Specialization Training

Our neuron structural analysis showed the presence of specialized neurons within the Feed-Forward Network (FFN) layers of a multilingual network. We hypothesize that continuously training the model, while leveraging these specialized neurons’ intrinsic modular features, can further enhance task-specific performance. Building on this hypothesis, we propose *Neuron Specialization*, an approach that leverages specialized neurons to modularize the FFN layers in a task-specific manner.

4.1 Vanilla Feed-Forward Network

We first revisit the Feed-Forward Network (FFN) in Transformer (Vaswani et al., 2017). The FFN, crucial to our analysis, consists of two linear layers (fc1 and fc2) with a *ReLU* activation function. Specifically, the FFN block first processes the hidden state $H \in \mathbb{R}^{n \times d}$ (n denotes number of tokens in a batch) through fc1 layer $W_1 \in \mathbb{R}^{d \times d_{ff}}$. Then the output is passed to *ReLU* and the fc2 layer W_2 , as formalized in Eq 2, with bias terms omitted.

$$\text{FFN}(H) = \text{ReLU}(HW_1)W_2. \quad (2)$$

4.2 Specializing Task-Specific FFN

Next, we investigate continuous training upon a subset of specialized parameters within FFN for each task. Given a pre-trained vanilla multilingual Transformer model with tags to identify the language pairs, e.g., Johnson et al. (2017), we can derive specialized neuron set S_k^t for each layer of a task² t and threshold k following the method outlined in Section 3.1. Then, we derive a boolean mask vector $m_k^t \in \{0, 1\}^{d_{ff}}$ from S_k^t , where the i -th element in m_k^t is set to 1 only when $i \in S_k^t$, and apply it to control parameter updates. Specifically, we broadcast m_k^t and perform Hadamard Product with W_1 in each FFN layer as follows:

$$\text{FFN}(H) = \text{ReLU}(H(m_k^t \odot W_1))W_2. \quad (3)$$

m_k^t plays the role of controlling parameter update, where the boolean value of i -th element in m_k^t denotes if the i -th row of parameters in W_1 can be updated or not for each layer³ during continues training. Broadly speaking, our approach selectively updates the first FFN (fc1) weights during

²We treat each translation direction as a distinct task.

³Note that m_k^t is layer-specified, we drop layer indexes hereon for simplicity of notation.

back-propagation, tailoring the model more closely towards specific translation tasks and reinforcing neuron separation. Note that while fc1 is selectively updated for specific tasks, other parameters are universally updated to maintain stability, and the same masking is applied to inference to ensure consistency. Our pseudocode is in Appendix A.10.

Relevant studies like Xie et al. (2021), selectively pruning output hidden states during training and inference. In contrast, we utilize sparse sub-networks (fc1 weights), while they prune output hidden states from Transformer modules.

5 Experimental Setup

In this section, we evaluate the capability of our proposed method on small (IWSLT) and large-scale (EC30) multilingual machine translation tasks. More details of the datasets are in Appendix A.1.

5.1 Datasets

IWSLT. Following Lin et al. (2021), we constructed an IWSLT dataset with eight languages. We learned a 30k SentencePiece unigram (Kudo and Richardson, 2018) shared vocabulary and applied temperature sampling with $\tau = 2$. We use Flores-200 (Costa-jussà et al., 2022), merging *devtest* and *test*, as our test set.

EC30. We further validate our methods on EC30 dataset (Tan and Monz, 2023), which features 61 million parallel training sentences across 60 English-centric directions, representing five language families and various writing systems. We classify language pairs into low-resource (=100k), medium-resource (=1M), and high-resource (=5M) categories. We build a 128k size shared unigram vocabulary. Aligning with the original EC30 setups, we use Ntrex-128 (Federmann et al., 2022) as the validation set. Also, we use Flores-200 (merging *devtest* and *test*) as the test set for evaluation.

5.2 Systems

We compare our method with strong open-source baselines that share similar motivations in reducing interference for multilingual translation tasks.

mT-small. For IWSLT, we train an mT-small baseline model on Many-to-Many directions as per (Lin et al., 2021): a 6-layer Transformer with 4 attention heads, $d = 512$, $d_{ff} = 1,024$.

Language Size	$\Delta\theta$	Fa 89k	Pl 128k	Ar 139k	He 144k	Nl 153k	De 160k	It 167k	Es 169k	Avg
One-to-Many (O2M / En-X)										
mT-small	-	14.5	9.9	12.0	13.1	17.0	20.6	17.3	18.3	15.4
Fine-Tune	0%	+0.1	-0.2	+0.2	+0.4	-0.4	-0.1	-0.3	-0.5	-0.1
Adapter _{LP}	+67%	+0.1	-0.1	+0.4	+1.4	+0.2	+0.6	+0.1	+0.4	+0.4
LaSS	0%	-2.6	0	+0.6	+0.7	-0.2	+0.7	-0.2	-0.4	-0.2
Ours	0%	+0.7	+0.1	+0.9	+0.6	+0.1	+0.1	+0.2	-0.3	+0.3
Many-to-One (M2O / X-En)										
mT-small	-	19.1	19.4	25.7	30.9	30.6	28.1	29.0	34.0	24.7
Fine-Tune	0%	+0.3	-0.2	+0.1	+0.8	+0.7	+0.3	-0.2	0	+0.2
Adapter _{LP}	+67%	+0.9	+0.6	+0.9	+1.0	+0.8	+1.0	+0.9	+0.3	+0.8
LaSS	0%	+1.2	+0.6	+0.9	+1.4	+1.1	+1.6	+1.6	+0.8	+1.2
Ours	0%	+1.6	+1.2	+1.7	+2.0	+1.9	+2.1	+1.8	+1.4	+1.7

Table 1: BLEU improvements over the baseline (mT-small) on IWSLT. $\Delta\theta$ denotes the relative parameter increase over the baseline, and 'Fine-Tune' signifies finetuning mT-small with the same setting as 'Ours'.

mT-big For EC30, we train a mT-big baseline model on Many-to-Many directions following Wu and Monz (2023). It has 6 layers, with 16 attention heads, $d = 1,024$, and $d_{ff} = 4,096$.

Fine-Tune. We finetune baselines with the same routine as our Neuron Specialization Training.

Adapters. We employ two adapter methods: 1) Language Pair Adapter (**Adapter_{LP}**) and 2) Language Family Adapter (**Adapter_{Fam}**). We omit Adapter_{Fam} for IWSLT due to its limited languages. Adapter_{LP} inserts adapter modules based on language pairs, demonstrating strong effects in reducing interference while presenting no parameter sharing (Bapna and Firat, 2019). In contrast, Adapter_{Fam} (Chronopoulou et al., 2023) facilitates parameter sharing across similar languages by training modules for each language family. Their bottleneck dimensions are 128 and 512 respectively. See Appendix A.2 for more training details.

LaSS. Lin et al. (2021) proposed LaSS to locate language-specific sub-networks following the lottery ticket hypothesis, i.e., finetuning all translation directions from a pre-trained model and then pruning based on magnitude. They then continually train the pre-trained model by only updating the sub-networks for each direction. We adopt the strongest LaSS configuration by applying sub-networks for both attention and FFNs.

5.3 Implementation and Evaluation

We train baseline models following the same hyperparameter settings in Lin et al. (2021) and Wu and

Monz (2023). For fair comparisons, we use the fixed training routine for all compared methods, see detailed training and model specifications in Appendix A.2. We adopt the tokenized BLEU (Papineni et al., 2002) for the IWSLT and detokenized SacreBLEU⁴ (Post, 2018) for the EC30. In addition, we report ChrF++ (Popović, 2017) and COMET (Rei et al., 2020) in Appendix A.4.

6 Results and Analyses

6.1 Small-Scale Results on IWSLT

We show results on IWSLT in Table 1. For Many-to-One (M2O) directions, our method receives an average +1.7 BLEU gain over the baseline, achieving the best performance among all approaches. The Adapter_{LP}, with a 67% increase in parameters over the baseline model, shows weaker improvements (+0.8) than our method. As for One-to-Many (O2M) directions, we observed weaker performance gains for all methods. While the gains are modest (averaging +0.3 BLEU), our method demonstrates consistent improvements across various languages in general. Finally, we show that fine-tuning the baseline with the same setting as our approach does not bring performance gains.

Scaling up does not always reduce interference. Shaham et al. (2023); Chang et al. (2023) have found scaling up the model capacity reduces interference, even under low-resource settings. We then investigate the trade-off between performance and model capacity by employing mT-shallow, a

⁴nrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.3.1

Methods	$\Delta\theta$	High (5M)			Med (1M)			Low (100K)			All (61M)		
		O2M	M2O	Avg	O2M	M2O	Avg	O2M	M2O	Avg	O2M	M2O	Avg
mT-big	-	28.1	31.6	29.9	29.7	31.6	30.6	18.9	26.0	22.4	25.5	29.7	27.7
Fine-Tune	0%	+0.3	+0.2	+0.3	+0.3	+0.2	+0.3	+0.1	-0.4	-0.2	+0.2	0	+0.1
Adapter _{Fam}	+70%	+0.7	+0.3	+0.5	+0.7	+0.3	+0.5	+1.1	+0.5	+0.8	+0.8	+0.4	+0.6
Adapter _{LP}	+87%	+1.6	+0.6	+1.1	+1.6	+0.4	+1.0	+0.4	+0.4	+0.4	+1.2	+0.5	+0.8
LaSS	0%	+2.3	+0.8	+1.5	+1.7	+0.2	+1.0	-0.1	-1.8	-1.0	+1.3	-0.3	+0.5
Random	0%	+0.9	-0.5	+0.2	+0.5	-0.7	-0.2	-0.3	-1.5	-0.9	+0.5	-0.9	-0.2
Ours ^{Enc}	0%	+1.2	+1.1	+1.1	+1.0	+1.0	+1.0	+0.7	+0.8	+0.8	+1.0	+1.0	+1.0
Ours ^{Dec}	0%	+1.2	+1.1	+1.1	+0.9	+1.1	+1.0	+0.7	+1.1	+0.9	+0.9	+1.1	+1.0
Ours	0%	+1.8	+1.4	+1.6	+1.4	+1.1	+1.3	+1.4	+0.9	+1.2	+1.5	+1.1	+1.3

Table 2: Average SacreBLEU improvements on the EC30 dataset over the baseline (mT-big), categorized by High, Medium, and Low-resource translation directions. 'Random' denotes continually updating the model with randomly selected task-specific neurons. 'Ours^{Enc}' and 'Ours^{Dec}' indicate Neuron Specialization applied solely to the Encoder and Decoder, respectively, while 'Ours' signifies the method applied to both components.

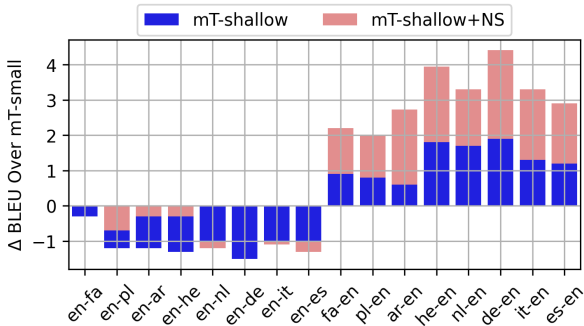


Figure 3: BLEU gains of shallower models over mT-small on IWSLT show improved X-En performance at the expense of En-X. Applying Neuron Specialization reduces EN-X degradation and amplifies X-En gains.

shallower version of mT-small with three fewer layers (with $\Delta\theta = -39\%$ for parameters, see Table 6 for details). Surprisingly, in Figure 3, we show that reducing parameters improved Many-to-One (X-En) performance but weakened One-to-Many (En-X) results. This result indicates that scaling up the model capacity does not always reduce interference, but may show overfitting to have performance degradation. Furthermore, we show that implementing Neuron Specialization with mT-shallow enhances X-En performance in all directions while lessening the decline in En-X translation quality.

6.2 Large-Scale Results on EC-30

Similar to what we observed in the small-scale setting, we find notable improvements when we scale up on the EC30 dataset. Table 2 shows consistent improvements across high-, medium-, and low-resource languages, with an average gain of +1.3 SacreBLEU over the baseline. LaSS, while

effective in high-resource O2M pairs, presents limitations with negative impacts (-1.0 score) on low-resource languages, highlighting difficulties in sub-network extraction for low-resource languages. In contrast, our method achieves stable and consistent gains and passes statistical significance tests in A.5. The Adapter_{LP}, despite increasing parameters by 87% compared to the baseline, falls short of our method in boosting performance. Similar to experiments on IWSLT, we found fine-tuning the baseline on EC30 also brings worse/unchanged performance, suggesting the effectiveness of our method. Additionally, we show that applying Neuron Specialization in the encoder or decoder delivers similar gains, with both combined offering stronger performance.

Random Mask. We applied Neuron Specialization Training using random masks that masked 30% fc1 weights to validate the effectiveness of our method in locating task-specific neurons. We show that such strategy sacrifices performance.

Zero-shot Translation. We further evaluated our method on 870 zero-shot directions using the EC30 dataset, observing an average improvement of +3.1 SacreBLEU. Of these, 847 directions improved, while 23 experienced minor declines of -0.3 SacreBLEU on average. See Appendix A.7 for details.

Wider and Deeper Models. We experiment with larger models by scaling up the width and depth in A.6. Table 8 shows we achieve consistent performance gains, confirming the effectiveness of our approach for larger configurations.

Lang Size	De 5m	Es 5m	Cs 5m	Hi 5m	Ar 5m	Lb 100k	Ro 100k	Sr 100k	Gu 100k	Am 100k	High Avg	Low Avg
One-to-Many												
Bilingual	36.3	24.6	28.7	43.9	23.7	5.5	16.2	17.8	12.8	4.1	31.8	11.3
mT-big	-4.7	-1.5	-3.6	-4.4	-4.7	+9.0	+8.9	+6.2	+13.9	+3.1	-3.7	+8.2
Ours	-2.0	-0.2	-1.7	-2.4	-3.0	+10.8	+10.0	+8.2	+16.4	+3.7	-1.9	+9.8
Many-to-One												
Bilingual	39.1	24.5	32.6	35.5	30.8	8.7	19.5	21.3	7.0	8.7	32.7	13.0
mT-big	-1.5	+0.9	+0.2	-1.8	-2.3	+13.7	+11.9	+10.3	+18.2	+12.5	-1.1	+13.3
Ours	-0.3	+1.7	+1.8	-0.2	-0.3	+15.3	+12.4	+11.3	+19.6	+14.1	+0.3	+14.5

Table 3: SacreBLEU score comparisons for Multilingual baseline and Neuron Specialization models against Bilingual ones on the EC30 dataset, limited to 5 high- and low-resource languages due to computational constraints. Red signifies negative interference, Blue denotes positive synergy, with darker shades indicating better effects.

The role of threshold factor. In A.8, we explore the impact of our sole hyper-parameter k (neuron selection threshold factor) on performance. We show that our method delivers consistent and positive gains without extensive hyperparameter tuning.

Model	$\Delta\theta$	ΔT_{subnet}	Δ Memory
Adapter _{LP}	+87%	n/a	1.42 GB
LaSS	0%	+33 hours	9.84 GB
Ours	0%	+5 minutes	3e-3 GB

Table 4: Efficiency comparison on EC30 dataset regarding extra trainable parameters ($\Delta\theta$: relative increase over the baseline), extra processing time for subnet extraction (ΔT_{subnet}), and extra memory (Δ Memory).

Efficiency Comparisons. We compare efficiency across three aspects (Table 4). First, adding lightweight language pair adapters results in an +87% increase in trainable parameters over the baseline. Second, our method, which locates specialized neurons in just 5 minutes, is significantly faster than LaSS, which takes 33 hours with 4 Nvidia A6000 GPUs. Finally, regarding memory costs essential for handling multiple languages in deployment, our method is more economical, requiring only 1-bit masks for the FFN neurons instead of extensive parameters.

6.3 The Impact of Reducing Interference

In this section, we measure to what extent our method mitigates interference and enhances knowledge transfer. Similar to Wang et al. (2020), we train bilingual models that do not contain interference or transfers, then compare results between bilingual models, the multilingual baseline model (mT-big), and our method (ours). We train Transformer-big and Transformer-based models for

high- and low-resource tasks, see Appendix A.2.

In Table 3, we show that the multilingual model (mT-big) facilitates clear positive transfer for low-resource languages versus bilingual setups, leading to +8.2 (O2M) and +13.3 (M2O) score gains but incurs negative interference for high-resource languages (-3.7 and -1.1 scores).

Our method reduces interference for high-resource settings, leading to +1.8 and +1.4 SacreBLEU gains over mT-big in O2M and M2O directions. Moreover, our Neuron Specialization method enhances low-resource task performance with average gains of +1.6 (O2M) and +1.2 (M2O) SacreBLEU over the mT-big, demonstrating its ability to foster cross-lingual knowledge transfer.

7 Conclusions

In this paper, we have identified and leveraged *intrinsic task-specific modularity* within multilingual networks to mitigate interference. We showed that FFN neurons activate in a language-specific way, and they present structural overlaps that reflect language proximity, which progress across layers. We then introduced *Neuron Specialization* to leverage these natural modularity signals to structure the network, enhancing task specificity and improving knowledge transfer. Our experimental results, spanning various resource levels, show that our method consistently outperforms strong baseline systems, with additional analyses demonstrating reduced interference and increased knowledge transfer. Our work deepens the understanding of multilingual models by revealing their intrinsic modularity, offering insights into how multi-task models can be optimized without extensive modifications.

588 Limitations

589 This study primarily focuses on Multilingual Ma-
590 chine Translation, a key method in multi-task learn-
591 ing, using it as our primary testbed. However,
592 the exploration of multilingual capabilities can be
593 extended beyond translation to include a broader
594 range of Multilingual Natural Language Processing
595 tasks. These areas remain unexplored in our current
596 research and are considered promising directions
597 for future work. In this work, we focus on the
598 feed-forward network (FFN) components within
599 the Transformer architecture, which constitutes a
600 significant portion of the model’s parameters. We
601 leave investigations of other Transformer compo-
602 nents, such as the layer normalization modules, to
603 future work.

604 Furthermore, our method identifies task-specific
605 neurons in Feed-Forward Networks that use the
606 ReLU activation function. Although this could be
607 one of the limitations of our work, we motivate
608 it on the following aspects. Firstly, ReLU deliv-
609 ers negligible impact on convergence and perfor-
610 mance while significantly reducing computation
611 and weight transfer (Mirzadeh et al., 2023) than
612 other activation functions like GeLU (Hendrycks
613 and Gimpel, 2016). Secondly, ReLU is still the
614 most common activation function for state-of-the-
615 art MNMT systems, such as NLLB-200 (Costa-
616 jussà et al., 2022), M2M-100 (Fan et al., 2021),
617 SeamlessM4T (Barrault et al., 2023).

618 Lastly, ReLU is monotonic, thus offering better
619 interpretability than GeLU (Sudjianto et al., 2020),
620 which is important for analyzing the modularity in
621 MNMT. Recent work on Large Language Models
622 has also explored the binary activation states of
623 FFN neurons, particularly focused on when neu-
624 rons are activated, and their roles in aggregating
625 information (Voita et al., 2023).

626 Broader Impact

627 Recognizing the inherent risks of mistranslation
628 in machine translation data, we have made efforts
629 to prioritize the incorporation of high-quality data,
630 such as two open-sourced Multilingual Machine
631 Translation datasets: IWSLT and EC30. Addition-
632 ally, issues of fairness emerge, meaning that the ca-
633 pacity to generate content may not be equitably dis-
634 tributed across different languages or demographic
635 groups. This can lead to the perpetuation and am-
636 plification of existing societal prejudices, such as
637 biases related to gender, embedded in the data.

References

- Roe Aharoni, Melvin Johnson, and Orhan Firat. 2019. Massively multilingual neural machine translation. 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 3874–3884.
- Ali Araabi and Christof Monz. 2020. Optimizing transformer for low-resource neural machine translation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3429–3435.
- Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun, Mia Xu Chen, Yuan Cao, George Foster, Colin Cherry, et al. 2019. Massively multilingual neural machine translation in the wild: Findings and challenges. *arXiv preprint arXiv:1907.05019*.
- Ankur Bapna and Orhan Firat. 2019. Simple, scalable adaptation for neural machine translation. 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 1538–1548.
- Loïc Barrault, Yu-An Chung, Mariano Cora Meglioli, David Dale, Ning Dong, Paul-Ambroise Duquenne, Hady Elsahar, Hongyu Gong, Kevin Heffernan, John Hoffman, et al. 2023. Seamlessm4t-massively multilingual & multimodal machine translation. *arXiv preprint arXiv:2308.11596*.
- Tyler A Chang, Catherine Arnett, Zhuowen Tu, and Benjamin K Bergen. 2023. When is multilinguality a curse? language modeling for 250 high-and low-resource languages. *arXiv preprint arXiv:2311.09205*.
- Liang Chen, Shuming Ma, Dongdong Zhang, Furu Wei, and Baobao Chang. 2023. On the off-target problem of zero-shot multilingual neural machine translation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 9542–9558.
- Rochelle Choenni, Dan Garrette, and Ekaterina Shutova. 2023a. Cross-lingual transfer with language-specific subnetworks for low-resource dependency parsing. *Computational Linguistics*, 49(3):613–641.
- Rochelle Choenni, Ekaterina Shutova, and Dan Garrette. 2023b. Examining modularity in multilingual lms via language-specialized subnetworks. *arXiv preprint arXiv:2311.08273*.
- Alexandra Chronopoulou, Dario Stojanovski, and Alexander Fraser. 2023. Language-family adapters for low-resource multilingual neural machine translation. In *Proceedings of the The Sixth Workshop on Technologies for Machine Translation of Low-Resource Languages (LoResMT 2023)*, pages 59–72.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco

693	Guzmán, Édouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8440–8451.	749
694		750
695		751
696		
697		
698		
699	Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. <i>arXiv preprint arXiv:2207.04672</i> .	
700		
701		
702		
703		
704		
705	Katharina Dobs, Julio Martinez, Alexander JE Kell, and Nancy Kanwisher. 2022. Brain-like functional specialization emerges spontaneously in deep neural networks. <i>Science advances</i> , 8(11):eabl8913.	
706		
707		
708		
709	Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, et al. 2021. Beyond english-centric multilingual machine translation. <i>The Journal of Machine Learning Research</i> , 22(1):4839–4886.	
710		
711		
712		
713		
714		
715	Christian Federmann, Tom Kocmi, and Ying Xin. 2022. Ntrex-128–news test references for mt evaluation of 128 languages. In <i>Proceedings of the First Workshop on Scaling Up Multilingual Evaluation</i> , pages 21–24.	
716		
717		
718		
719	Jonathan Frankle and Michael Carbin. 2018. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In <i>International Conference on Learning Representations</i> .	
720		
721		
722		
723	Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, errors, and context: A large-scale study of human evaluation for machine translation. <i>Transactions of the Association for Computational Linguistics</i> , 9:1460–1474.	
724		
725		
726		
727		
728		
729	Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André FT Martins. 2022. Results of wmt22 metrics shared task: Stop using bleu–neural metrics are better and more robust. In <i>Proceedings of the Seventh Conference on Machine Translation (WMT)</i> , pages 46–68.	
730		
731		
732		
733		
734		
735		
736	Dan He, Minh Quang Pham, Thanh-Le Ha, and Marco Turchi. 2023. Gradient-based gradual pruning for language-specific multilingual neural machine translation. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 654–670.	
737		
738		
739		
740		
741		
742	Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus). <i>arXiv preprint arXiv:1606.08415</i> .	
743		
744		
745	Melvin Johnson, Mike Schuster, Quoc V Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, et al. 2017. Google’s multilingual neural machine	
746		
747		
748		
	translation system: Enabling zero-shot translation. <i>Transactions of the Association for Computational Linguistics</i> , 5:339–351.	749
		750
		751
	Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 66–71.	752
		753
		754
		755
		756
		757
	Sneha Kudugunta, Ankur Bapna, Isaac Caswell, and Orhan Firat. 2019. Investigating multilingual nmt representations at scale. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1565–1575.	758
		759
		760
		761
		762
		763
		764
	Teven Le Scao, Angela Fan, Christopher Akiki, Elie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model.	765
		766
		767
		768
		769
	Xian Li and Hongyu Gong. 2021. Robust optimization for multilingual translation with imbalanced data. <i>Advances in Neural Information Processing Systems</i> , 34:25086–25099.	770
		771
		772
		773
	Baohao Liao, Yan Meng, and Christof Monz. 2023a. Parameter-efficient fine-tuning without introducing new latency . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4242–4260, Toronto, Canada. Association for Computational Linguistics.	774
		775
		776
		777
		778
		779
	Baohao Liao, Shaomu Tan, and Christof Monz. 2023b. Make pre-trained model reversible: From parameter to memory efficient fine-tuning. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	780
		781
		782
		783
		784
	Zehui Lin, Liwei Wu, Mingxuan Wang, and Lei Li. 2021. Learning language specific sub-network for multilingual machine translation. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 293–305.	785
		786
		787
		788
		789
		790
		791
	Seyed Iman Mirzadeh, Keivan Alizadeh-Vahid, Sachin Mehta, Carlo C del Mundo, Oncel Tuzel, Golnoosh Samei, Mohammad Rastegari, and Mehrdad Farajtabar. 2023. Relu strikes back: Exploiting activation sparsity in large language models. In <i>The Twelfth International Conference on Learning Representations</i> .	792
		793
		794
		795
		796
		797
	Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. <i>arXiv preprint arXiv:1904.01038</i> .	798
		799
		800
		801
	Xiao Pan, Mingxuan Wang, Liwei Wu, and Lei Li. 2021. Contrastive learning for many-to-many multilingual neural machine translation. In <i>Proceedings of the</i>	802
		803
		804

805		Uri Shaham, Maha Elbayad, Vedanuj Goswami, Omer Levy, and Shruti Bhosale. 2023. Causes and cures for interference in multilingual translation. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , Toronto, Canada. Association for Computational Linguistics.	859
806			860
807			861
808			862
809	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.		863
810			864
811			865
812		Agus Sudjianto, William Knauth, Rahul Singh, Zebin Yang, and Aijun Zhang. 2020. Unwrapping the black box of deep relu networks: interpretability, diagnostics, and simplification. <i>arXiv preprint arXiv:2011.04041</i> .	866
813			867
814	Jonas Pfeiffer, Naman Goyal, Xi Lin, Xian Li, James Cross, Sebastian Riedel, and Mikel Artetxe. 2022. Lifting the curse of multilinguality by pre-training modular transformers. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3479–3495.		868
815			869
816			870
817		Shaomu Tan and Christof Monz. 2023. Towards a better understanding of variations in zero-shot neural machine translation performance. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 13553–13568.	871
818			872
819			873
820			874
821	Jonas Pfeiffer, Sebastian Ruder, Ivan Vulić, and Edoardo Ponti. 2023. Modular deep learning . <i>Transactions on Machine Learning Research</i> . Survey Certification.		875
822			876
823			877
824	Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual bert? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4996–5001.		878
825			879
826			880
827			881
828	Telmo Pires, Robin Schmidt, Yi-Hsiu Liao, and Stephan Peitz. 2023. Learning language-specific layers for multilingual machine translation. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14767–14783.		882
829			883
830			884
831			885
832			886
833			887
834	Maja Popović. 2017. chrF++: words helping character n-grams. In <i>Proceedings of the second conference on machine translation</i> , pages 612–618.		888
835			889
836			890
837	Matt Post. 2018. A call for clarity in reporting bleu scores. In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 186–191.		891
838			892
839			893
840			894
841	Ricardo Rei, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André FT Martins. 2022. Comet-22: Unbabel-ist 2022 submission for the metrics shared task. In <i>Proceedings of the Seventh Conference on Machine Translation (WMT)</i> , pages 578–585.		895
842			896
843			897
844			898
845			899
846			900
847			901
848	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt evaluation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2685–2702.		902
849			903
850			904
851			905
852			906
853	Stefan Riezler and John T Maxwell III. 2005. On some pitfalls in automatic evaluation and significance testing for mt. In <i>Proceedings of the ACL workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization</i> , pages 57–64.		907
854			908
855			909
856			910
857			911
858			912
			913
			914
			915

916	Wanying Xie, Yang Feng, Shuhao Gu, and Dong Yu.	EC30	We utilize the EC30, a subset of the EC40	969
917	2021. Importance-based neuron allocation for multi-		dataset (Tan and Monz, 2023) (with 10 extremely	970
918	lingual neural machine translation. In <i>Proceedings</i>		low-resource languages removed in our experi-	971
919	<i>of the 59th Annual Meeting of the Association for</i>		ments) as our main dataset for most experiments	972
920	<i>Computational Linguistics and the 11th International</i>		and analyses. We list the Languages with their	973
921	<i>Joint Conference on Natural Language Processing</i>		ISO and scripts in Table 5, along with their num-	974
922	<i>(Volume 1: Long Papers)</i> , pages 5725–5737.		ber of sentences. In general, EC30 is an English-	975
923	Derrick Xin, Behrooz Ghorbani, Justin Gilmer, Ankush		centric Multilingual Machine Translation dataset	976
924	Garg, and Orhan Firat. 2022. Do current multi-task		containing 61 million sentences covering 30 lan-	977
925	optimization methods in deep learning even help?		guages (excluding English). It collected data from	978
926	<i>Advances in neural information processing systems</i> ,		5 representative language families with multiple	979
927	35:13597–13609.		writing scripts. In addition, EC30 is well bal-	980
928	Guangyu Robert Yang, Madhura R Joglekar, H Francis		anced at each resource level, for example, for all	981
929	Song, William T Newsome, and Xiao-Jing Wang.		high-resource languages, the number of training	982
930	2019. Task representations in neural networks trained		sentences is 5 million. Note that the EC30 is al-	983
931	to perform many cognitive tasks. <i>Nature neuro-</i>		ready pre-processed and tokenized (with Moses	984
932	<i>science</i> , 22(2):297–306.		tokenizer), thus we directly use it for our study.	985
933	Biao Zhang, Ankur Bapna, Rico Sennrich, and Orhan	A.2 Model and Training Details		986
934	Firat. 2020a. Share or not? learning to schedule		We list the configurations and hyper-parameter set-	987
935	language-specific capacity for multilingual transla-		tions of all systems for the main training setting	988
936	tion. In <i>International Conference on Learning Rep-</i>		(EC30) in Table 6. To maintain consistency and	989
937	<i>resentations</i> .		comparability across all experiments, we employed	990
938	Biao Zhang, Philip Williams, Ivan Titov, and Rico Sen-		the same early stopping settings rather than fix-	991
939	nrich. 2020b. Improving massively multilingual neu-		ing the training duration for all experiments. We	992
940	ral machine translation and zero-shot translation. In		use 4 NVIDIA A6000 (48G) GPUs to conduct	993
941	<i>Proceedings of the 58th Annual Meeting of the Asso-</i>		most experiments and implement them based on	994
942	<i>ciation for Computational Linguistics</i> , pages 1628–		Fairseq (Ott et al., 2019) with FP16.	995
943	1639.	Global training settings.	For all systems on both	996
944	Zhengyan Zhang, Zhiyuan Zeng, Yankai Lin, Chaojun		datasets, we adopt the pre-norm and share the de-	997
945	Xiao, Xiaozhi Wang, Xu Han, Zhiyuan Liu, Ruob-		coder input output embedding. In addition, we	998
946	ing Xie, Maosong Sun, and Jie Zhou. 2023. Emer-		use the Adam optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.98$,	999
947	gent modularity in pre-trained transformers. In <i>Find-</i>		$\epsilon = 10^{-9}$) with $5e-4$ learning rate and 4k warmup	1000
948	<i>ings of the Association for Computational Linguis-</i>		steps in all methods. Furthermore, we use cross	1001
949	<i>tics: ACL 2023</i> , pages 4066–4083, Toronto, Canada.		entropy with label smoothing to avoid overfitting	1002
950	Association for Computational Linguistics.		(smoothing factor=0.1) and set early stopping to 20.	1003
951	A Appendix		Similar to Fan et al. (2021), we prepend language	1004
952	A.1 Dataset details		tags to the source and target sentences to indicate	1005
953	Due to the difficulties of mining non-English-		the translation directions for all multilingual trans-	1006
954	centric Translation data, recent research (Johnson		lation systems. More importantly, we applied the	1007
955	et al., 2017; Zhang et al., 2020b,a; Tan and Monz,		same fixed routine across all experiments to ensure	1008
956	2023; Wu and Monz, 2023; Shaham et al., 2023;		a fair comparison among all multilingual systems.	1009
957	Pires et al., 2023) has increasingly focused on uti-		Other global settings are the same for all systems	1010
958	lizing English-centric datasets to explore Multilin-		to make fair comparisons, such as learning rate,	1011
959	gual Neural Machine Translation (MNMT). Fur-		warm-up steps, and batch size.	1012
960	thermore, Fan et al. (2021) have observed that train-	Bilingual models.	For bilingual models of low-	1013
961	ing in M2M settings does not necessarily enhance		resource languages, we adopt the suggested hyper-	1014
962	performance in supervised directions. Therefore,		parameter settings from Araabi and Monz (2020),	1015
963	our approach prioritizes English-centric datasets to		such as $d_{ff} = 512$, number of attention head as 2,	1016
964	remain computationally feasible while still provid-		and dropout as 0.3. Furthermore, We train separate	1017
965	ing valuable insights into MNMT dynamics.		dictionaries for low-resource bilingual models to	1018
966	IWSLT We collect and pre-processes the IWSLT-			
967	14 dataset following Lin et al. (2021). We refer			
968	readers to Lin et al. (2021) for more details.			

	Germanic			Romance			Slavic			Indo-Aryan			Afro-Asiatic		
	ISO	Language	Script	ISO	Language	Script	ISO	Language	Script	ISO	Language	Script	ISO	Language	Script
High (5m)	de	German	Latin	fr	French	Latin	ru	Russian	Cyrillic	hi	Hindi	Devanagari	ar	Arabic	Arabic
	nl	Dutch	Latin	es	Spanish	Latin	cs	Czech	Latin	bn	Bengali	Bengali	he	Hebrew	Hebrew
Med (1m)	sv	Swedish	Latin	it	Italian	Latin	pl	Polish	Latin	kn	Kannada	Devanagari	mt	Maltese	Latin
	da	Danish	Latin	pt	Portuguese	Latin	bg	Bulgarian	Cyrillic	mr	Marathi	Devanagari	ha	Hausa*	Latin
Low (100k)	af	Afrikaans	Latin	ro	Romanian	Latin	uk	Ukrainian	Cyrillic	sd	Sindhi	Arabic	ti	Tigrinya	Ethiopic
	lb	Luxembourgish	Latin	oc	Occitan	Latin	sr	Serbian	Latin	gu	Gujarati	Devanagari	am	Amharic	Ethiopic

Table 5: Details of EC30 Training Dataset. Numbers in the table represent the number of sentences, for example, 5m denotes exactly 5,000,000 number of sentences. The only exception is Hausa, where its size is 334k (334,000).

Models	Dataset	Num. trainable params	Num. Layer	Num. Attn Head	dim	d_{ff}	max tokens	update freq	dropout
mT-shallow	IWSLT	47M	3	8	512	1,024	2,560	4	0.1
mT-small	IWSLT	76M	6	8	512	1,024	2,560	4	0.1
bilingual-low	EC30	52M	6	2	512	1,024	2,560	1	0.3
bilingual-high	EC30	439M	6	16	1,024	4,096	2,560	10	0.1
mT-big	EC30	439M	6	16	1,024	4,096	7,680	21	0.1
LaSS	EC30	439M	6	16	1,024	4,096	7,680	21	0.1
Ours-big	EC30	439M	6	16	1,024	4,096	7,680	21	0.1
mT-wide	EC30	540M	6	16	1,024	8,192	7,680	21	0.1
Ours-wide	EC30	540M	6	16	1,024	8,192	7,680	21	0.1
mT-large	EC30	615M	12	16	1,024	4,096	7,680	21	0.1
Ours-large	EC30	615M	12	16	1,024	4,096	7,680	21	0.1

Table 6: Configuration and hyper-parameter settings for all models in this paper. Num. Layer and Attn Head denote the number of layers and attention heads, respectively. dim represents the dimension of the Transformer model, d_{ff} means the dimension of the feed-forward layer. bilingual-low and -high represent the bilingual models for low and high-resource languages.

1019 avoid potential overfitting instead of using the large
1020 128k shared multilingual dictionary.

1021 For bilingual models of high-resource languages,
1022 we adopt the 128k shared multilingual dictionary
1023 and train models with the Transformer-big archi-
1024 tecture as the multilingual baseline (mT-big). The
1025 detailed configurations can be found in Table 6.

1026 **Language Pair Adapters.** We implement Lan-
1027 guage Pair Adapters (Bapna and Firat, 2019) by
1028 ourselves based on Fairseq. The Language Pair
1029 Adapter is learned depending on each pair, e.g.,
1030 we learn two modules for en-de, namely en on the
1031 Encoder side and the de on the Decoder side. Note
1032 that, except for the unified pre-trained model, lan-
1033 guage pair adapters do not share any parameters
1034 with each other, preventing potential knowledge
1035 transfers. We set its bottleneck dimension as 128
1036 for all experiments of IWSLT and EC30.

- 1037 • **IWSLT.** For the IWSLT dataset that contains

8 languages with 16 translation directions, the
mT-small base model size is 76M. Adapter_{LP}
insert 3.2M extra trainable parameters for one
direction, thus resulting in 51.2M added pa-
rameters for all, leading to 67% relative pa-
rameter increase over the baseline model.

- 1044 • **EC30.** For the EC30 dataset that contains 30
1045 languages with 60 translation directions, the
1046 mT-big base model size is 439M. Adapter_{LP}
1047 inserts 6.4M extra trainable parameters for
1048 one direction, thus resulting in 384M added
1049 parameters for all directions, leading to 87%
1050 relative parameter increase over the baseline
1051 model. When training Adapter_{LP} for low-
1052 resource languages, we increased dropout (0.1
1053 -> 0.3) and decreased batch size (max-token:
1054 7680 -> 2560) to avoid overfitting as sug-
1055 gested by Bapna and Firat (2019).

Methods	θ	High (5M)			Med (1M)			Low (100K)		
		O2M	M2O	Avg	O2M	M2O	Avg	O2M	M2O	Avg
mT-big	438m	27.7	32.0	29.9	30.6	34.2	32.4	26.9	32.9	29.9
M2M-100	418m	23.3	28.0	25.7	30.8	32.9	31.9	24.6	32.0	28.3
M2M-100	1.2b	28.3	34.3	31.3	36.3	38.9	37.6	31.7	41.1	36.4
Ours-big	438m	29.6	33.3	31.5	32.0	35.5	33.8	28.1	33.7	30.9

Table 7: Performance comparisons on the EC30 test set using SacreBLEU. θ represents the number of parameters, and 'Ours-big' denotes our neuron specialization method applied to the mT-big. We excluded directions where the M2M-100 models scored ≤ 10 BLEU to ensure fair comparisons, resulting in 51 translation directions.

Methods	$\Delta\theta$	High (5M)			Med (1M)			Low (100K)			All (61M)		
		O2M	M2O	Avg	O2M	M2O	Avg	O2M	M2O	Avg	O2M	M2O	Avg
SacreBLEU													
mT-big	-	28.1	31.6	29.9	29.7	31.6	30.6	18.9	26.0	22.4	25.5	29.7	27.7
Ours-big	0%	+1.8	+1.4	+1.6	+1.4	+1.1	+1.3	+1.4	+0.9	+1.2	+1.5	+1.1	+1.3
mT-wide	+23%	+0.8	+0.6	+0.7	+0.7	+0.6	+0.6	+0.6	+0.6	+0.6	+0.6	+0.6	+0.6
Ours-wide	+23%	+2.2	+1.9	+2.1	+1.8	+1.7	+1.8	+1.4	+1.1	+1.3	+1.8	+1.5	+1.7
mT-large	+40%	+1.2	+1.2	+1.2	+1.0	+1.4	+1.2	+0.8	+1.6	+1.2	+1.0	+1.2	+1.1
Ours-large	+40%	+2.6	+2.3	+2.5	+1.9	+2.0	+2.0	+1.4	+2.2	+1.8	+2.0	+2.1	+2.0
ChrF++													
mT-big	-	52.4	57.6	55.0	54.0	56.6	55.3	42.5	50.0	46.3	49.6	54.7	52.1
Ours-big	0%	+1.4	+1.1	+1.3	+1.1	+0.9	+1.0	+1.2	+0.8	+1.0	+1.2	+0.9	+1.1
mT-wide	+23%	+0.7	+0.7	+0.7	+0.7	+0.6	+0.7	+0.6	+0.7	+0.7	+0.7	+0.6	+0.7
Ours-wide	+23%	+1.8	+1.6	+1.7	+1.5	+1.4	+1.5	+1.3	+1.0	+1.2	+1.6	+1.3	+1.4
mT-large	+40%	+0.9	+0.9	+0.9	+0.9	+1.1	+1.0	+0.8	+1.4	+1.1	+0.9	+1.1	+1.0
Ours-large	+40%	+2.0	+1.8	+1.9	+1.5	+1.7	+1.6	+1.3	+1.8	+1.6	+1.6	+1.8	+1.7
COMET													
mT-big	-	82.4	83.9	83.2	81.1	80.1	80.6	73.8	73.4	73.6	79.1	79.1	79.1
Ours-big	0%	+1.4	+1.0	+1.2	+0.9	+0.7	+0.8	+0.8	+0.7	+0.8	+1.0	+0.8	+0.9
mT-wide	+23%	+0.8	+0.6	+0.7	+0.6	+0.6	+0.6	+0.6	+0.6	+0.6	+0.7	+0.6	+0.6
Ours-wide	+23%	+1.8	+1.4	+1.6	+1.3	+1.3	+1.3	+1.3	+1.2	+1.3	+1.5	+1.3	+1.4
mT-large	+40%	+1.0	+0.8	+0.9	+0.7	+1.0	+0.9	+0.9	+1.2	+1.1	+0.9	+1.0	+0.9
Ours-large	+40%	+2.1	+1.6	+1.9	+1.3	+1.6	+1.5	+1.3	+1.9	+1.6	+1.6	+1.7	+1.6

Table 8: The effectiveness of our method on different model configurations. The table shows the averaged improvements on the EC30 dataset over the baseline (mT-big). 'Ours-big', 'Ours-wide', and 'Ours-large' indicate Neuron Specialization applied to the mT-big, mT-wide, and mT-large baselines respectively.

Language Family Adapters. The Language Family Adapter (Chronopoulou et al., 2023) is learned depending on each language family, e.g., for all 6 Germanic languages in the EC30, we learn two modules for en-Germanic, namely the en adapter on the Encoder side and the Germanic adapter on the Decoder side. We set its bottleneck dimension as 512 for all experiments for the EC30.

- **EC30.** For the EC30 dataset that contains 30 languages with 60 translation directions, the

mT-big base model size is 439M. Adapter_{Fam} insert 25.3M additional trainable parameters for one family (on EN-X directions), thus resulting in 303.6M added parameters for all families on both EN-X and X-En directions, leading to 69% relative parameter increase over the baseline model.

LaSS. When reproducing LaSS (Lin et al., 2021), we adopt the code from their official Github page⁵

⁵<https://github.com/NLP-Playground/LaSS>

with the same hyper-parameter setting as they suggested in their paper. For IWSLT, we finetune the mT-small for each translation direction with dropout=0.3, and we set dropout=0.1 for large-scale EC30. We then identify the language-specific parameters for attention and feed-forward modules (the setting with the strongest gains in their paper) with a pruning rate of 70%. We continue to train the sparse networks while keeping the same setting as the pre-training phase as they suggested.

Note that we observed different results as they reported in the paper, even though we used the same code, hyper-parameter settings, and corresponding Python environment and package version. We also found that He et al. (2023) reproduced LaSS results in their paper, which shows similar improvements (around +0.6 BLUE gains) over the baseline of our reproductions. As for an improved method over LaSS proposed by He et al. (2023), we do not reproduce since no open-source code has been released.

A.3 Comparison with M2M-100 Models

We choose multilingual Transformer architecture as our baseline backbone, which has been commonly used as a strong baseline in many MNMT studies (Pires et al., 2023; Shaham et al., 2023; Arivazhagan et al., 2019; Wu et al., 2024), and is widely recognized as a strong baseline within the community (Chen et al., 2023; Wu et al., 2023; Pan et al., 2021; Wu and Monz, 2023).

We further establish the strength of our baseline models by comparing them to the M2M-100 models, which are state-of-the-art systems trained on an extensive corpus of 7.5 billion parallel sentences. In specific, we directly evaluated the trained M2M-100 models provided in Fairseq⁶. The results, presented in Table 7, demonstrate that both our baseline model (mT-big) and our proposed method (Ours) achieve performance that is comparable to, or even surpasses, the M2M-100 models.

A.4 Main result using ChrF++ and COMET

Recent studies (Rei et al., 2020; Costa-jussà et al., 2022) show that ChrF and COMET present high levels of correlation with human judgments, and automatic metrics based on pre-trained embeddings can outperform human crowd workers (Freitag et al., 2021). Notably, Costa-jussà et al. (2022)

found an increase of +0.5 in ChrF++ has been correlated with statistically significant improvements in human evaluations, with a change of +1.0 in ChrF++ almost always perceptible to human evaluators, which is studied on the FLORES test set.

To ensure a comprehensive evaluation, we report various automatic metrics in this paper: ChrF++(character level), SacreBleu (detokenized word level), and COMET(representation level) scores as extra results, as shown in Table 9, respectively. We opted for the "wmt22-comet-da" model (Rei et al., 2022), a widely used version from Unbabel’s collection of models that serves as the default choice. This model presents SOTA performance in WMT Metrics Shared Task (Freitag et al., 2022). Similar to what we observed in Section 6.2, our Neuron Specialization presents consistent performance improvements over the baseline model while outperforming other methods such as LaSS and Adapters.

Our method, applied to the same FLORES-200 test set, outperformed the baseline with an average increase of +1.1 ChrF++ scores, where most gains were greater than +1.0 ChrF++. This improvement emphasizes the effectiveness of our approach, suggesting a significant alignment with human evaluative standards.

A.5 Robustness tests

To show that the improvements in our method are not due to random variance, we implemented our method with different random seeds for all experiments and conducted paired significance tests for our main EC30 results.

A.5.1 Testing with Different Random Seeds

We run our method with different seeds and show robust improvements for both datasets (see Table 10 and Table 11).

Seed	O2M	M2O
ΔBLEU over mT-shallow		
seed=222	+0.3	+1.8
seed=111	+0.3	+1.4
ΔBLEU over mT-small		
seed=222	+0.3	+1.7
seed=111	+0.6	+1.2

Table 10: Average BLEU improvements of our Neuron Specialization method (Ours) over baselines (mT-shallow and mT-small) on the IWSLT dataset.

⁶https://github.com/facebookresearch/fairseq/tree/main/examples/m2m_100

Methods	$\Delta\theta$	High (5M)			Med (1M)			Low (100K)			All (61M)		
		O2M	M2O	Avg	O2M	M2O	Avg	O2M	M2O	Avg	O2M	M2O	Avg
SacreBLEU													
mT-big	-	28.1	31.6	29.9	29.7	31.6	30.6	18.9	26.0	22.4	25.5	29.7	27.7
Fine-Tune	0%	+0.3	+0.2	+0.3	+0.3	+0.2	+0.3	-0.3	-0.4	-0.4	+0.3	0	+0.1
Adapter _{Fam}	+70%	+0.7	+0.3	+0.5	+0.7	+0.3	+0.5	+1.1	+0.5	+0.8	+0.8	+0.4	+0.6
Adapter _{LP}	+87%	+1.6	+0.6	+1.1	+1.6	+0.4	+1.0	+0.4	+0.4	+0.4	+1.2	+0.5	+0.8
LaSS	0%	+2.3	+0.8	+1.5	+1.7	+0.2	+1.0	-0.1	-1.8	-1.0	+1.3	-0.3	+0.5
Random	0%	+0.9	-0.5	+0.2	+0.5	-0.7	-0.2	-0.3	-1.5	-0.9	+0.5	-0.9	-0.2
Ours-big ^{Enc}	0%	+1.2	+1.1	+1.1	+1.0	+1.0	+1.0	+0.7	+0.8	+0.8	+1.0	+1.0	+1.0
Ours-big ^{Dec}	0%	+1.2	+1.1	+1.1	+0.9	+1.1	+1.0	+0.7	+1.1	+0.9	+0.9	+1.1	+1.0
Ours-big	0%	+1.8	+1.4	+1.6	+1.4	+1.1	+1.3	+1.4	+0.9	+1.2	+1.5	+1.1	+1.3
ChrF++													
mT-big	-	52.4	57.6	55.0	53.9	56.6	55.3	42.5	50.0	46.3	49.6	54.7	52.2
Adapter _{LP}	+87%	+1.3	+0.2	+0.8	+1.1	+0.1	+0.6	+0.3	+0.3	+0.3	+0.9	+0.2	+0.5
Adapter _{Fam}	+70%	+0.6	+0.2	+0.4	+0.7	+0.3	+0.5	+1.1	+0.4	+0.8	+0.8	+0.3	+0.5
LaSS	0%	+1.7	+0.8	+1.2	+1.3	+0.3	+0.8	-0.3	-1.5	-0.9	+0.9	-0.2	+0.5
Random	0%	+0.7	-0.4	+0.2	+0.4	-0.5	-0.1	-0.5	-1.2	-0.9	+0.2	-0.7	-0.3
Ours-big ^{Enc}	0%	+1.0	+0.9	+1.0	+0.7	+0.9	+0.8	+0.6	+0.9	+0.8	+0.8	+0.9	+0.8
Ours-big ^{Dec}	0%	+0.9	+0.9	+0.9	+0.6	+1.0	+0.8	+0.5	+1.2	+0.9	+0.7	+1.0	+0.9
Ours-big	0%	+1.4	+1.1	+1.3	+1.1	+0.9	+1.0	+1.2	+0.8	+1.0	+1.2	+0.9	+1.1
COMET													
mT-big	-	83.4	83.9	83.65	81.1	80.1	80.6	73.8	73.4	73.6	79.1	79.1	79.1
Adapter _{LP}	+87%	+0.9	+0.2	+0.5	+0.6	+0.2	+0.4	0	+0.1	0	+0.5	+0.2	+0.4
Adapter _{Fam}	+70%	+0.4	+0.1	+0.3	+0.4	+0.2	+0.3	+0.7	+0.3	+0.5	+0.5	+0.2	+0.4
LaSS	0%	+1.5	+0.8	+1.2	+0.9	+0.6	+0.8	-0.2	-1.0	-0.6	+0.7	+0.1	+0.4
Random	0%	+0.2	-0.1	+0.1	-0.1	-0.2	-0.2	-0.8	-0.9	-0.9	-0.2	-0.4	-0.3
Ours-big ^{Enc}	0%	+1.0	+0.8	+0.9	+0.5	+0.9	+0.7	+0.3	+0.9	+0.6	+0.6	+0.8	+0.7
Ours-big ^{Dec}	0%	+0.9	+0.8	+0.9	+0.5	+1.0	+0.8	+0.3	+0.9	+0.6	+0.6	+1.0	+0.8
Ours-big	0%	+1.4	+1.0	+1.2	+0.9	+0.7	+0.8	+0.8	+0.7	+0.8	+1.0	+0.8	+0.9

Table 9: Average improvements on the EC30 dataset over the baseline (mT-big). ‘Ours-big^{Enc}’ and ‘Ours-big^{Dec}’ indicate neuron specialization applied solely to the Encoder and Decoder, respectively, while ‘Ours-big’ signifies the method applied to both components. The best results are highlighted in **bold**.

Seed	O2M	M2O	M2M
Δ SacreBLEU over mT-big			
seed=222	+1.5	+1.1	+1.3
seed=111	+1.3	+1.1	+1.2
seed=42	+1.4	+1.2	+1.3

Table 11: Average SacreBLEU improvements of our Neuron Specialization method (Ours) over the baseline (mT-big) on the EC30 dataset.

A.5.2 Statistical Significance Test

We conducted Paired approximate randomization (Riezler and Maxwell III, 2005) paired significance test to show that the improvements of our method over the baseline (mT-big) on EC30 are statistically significant regarding SacreBLEU and CHRf++ metrics in Table 12. In sum, for both

metrics, 59/60 directions passed the test (p-value < 0.05) except en-ha. The test is performed with the SacreBLEU Python package’s paired significance testing feature (–paired-ar).

A.6 Experiments on wider and deeper models

We conducted further experiments to determine if our method retains its effectiveness with larger models. We expanded the baseline model, mT-big, in two key dimensions: a) the feed-forward network (FFN) size, indicating the ‘width’ of the network; b) the number of layers, representing the ‘depth’ of the network. Specifically, we introduced mT-wide, which features an expanded FFN dimensionality (from 4,096 to 8,192), and mT-large, which has increased layer count (from 6-6 to 12-12). See model config details in Table 6.

Following these modifications, we applied our

Statistical Significance Test based on SacreBLEU														
en-af	en-am	en-ar	en-bg	en-bn	en-cs	en-da	en-de	en-es	en-fr	en-gu	en-ha	en-he	en-hi	en-it
3e-3	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	2e-1	9e-4	9e-4	9e-4
en-kn	en-lb	en-mr	en-mt	en-nl	en-oc	en-pl	en-pt	en-ro	en-ru	en-sd	en-sr	en-sv	en-ti	en-uk
9e-4	9e-4	3e-3	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	2e-2	9e-4
af-en	am-en	ar-en	bg-en	bn-en	cs-en	da-en	de-en	es-en	fr-en	gu-en	ha-en	he-en	hi-en	it-en
9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4	9e-4
kn-en	lb-en	mr-en	mt-en	nl-en	oc-en	pl-en	pt-en	ro-en	ru-en	sd-en	sr-en	sv-en	ti-en	uk-en
9e-4	9e-4	9e-4	9e-4	9e-4	1e-2	9e-4	9e-4	1e-2	9e-4	3e-2	9e-4	9e-4	9e-4	9e-4
Statistical Significance Test based on ChrF++														
en-af	en-am	en-ar	en-bg	en-bn	en-cs	en-da	en-de	en-es	en-fr	en-gu	en-ha	en-he	en-hi	en-it
9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	1e-01	9e-04	9e-04	9e-04
en-kn	en-lb	en-mr	en-mt	en-nl	en-oc	en-pl	en-pt	en-ro	en-ru	en-sd	en-sr	en-sv	en-ti	en-uk
9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	2e-02	9e-04
af-en	am-en	ar-en	bg-en	bn-en	cs-en	da-en	de-en	es-en	fr-en	gu-en	ha-en	he-en	hi-en	it-en
9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04
kn-en	lb-en	mr-en	mt-en	nl-en	oc-en	pl-en	pt-en	ro-en	ru-en	sd-en	sr-en	sv-en	ti-en	uk-en
9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	9e-04	8e-02	9e-04	9e-04	9e-04	9e-04

Table 12: Statistical Significance Test comparing our Neuron Specialization against the mT-big baseline on EC30. The table shows p-values in each direction, with p-value < 0.05 indicating our method yields significant improvement over the baseline. Overall, for both metrics, 59/60 directions passed the test (p-value < 0.05) except en-ha.

neuron specialization approach to these models. The results, as shown in Table 8, demonstrate consistent performance gains across both configurations, further validating the efficacy of our method.

A.7 Results in Zero-shot translations

Zero-shot neural machine translation (ZS-NMT) represents a pivotal challenge in multilingual machine translation, aiming to handle language pairs never seen during training. Although training unified MMT systems enables zero-shot translations (Johnson et al., 2017), their performance falls short of that seen in supervised directions. Recent findings by Zhang et al. (2020b) suggest that larger model sizes enhance ZS performance. Additionally, Tan and Monz (2023) indicates that vocabulary overlap and linguistic similarities contribute to variations in ZS performance, and that stronger En-centric capabilities might improve ZS results.

ZS-NMT Setups To further investigate whether our method could bring benefits to zero-shot translations, we tested our method across 870 zero-shot directions involving 30 languages. To do that, we created masks using the Encoder mask from Source-to-English (Src-En) and the Decoder mask from English-to-Target (En-Tgt).

ZS-NMT Results Overall, we observed an averaged +3.1 SacreBLEU improvement on zero-shot

directions, with 847 out of 870 directions showing improvements, and 23 directions experiencing minor declines, averaging -0.3 SacreBLEU. Detailed results for high, medium, and low-resource languages (denoted as H, M, and L) are presented in Table 13, along with comparisons of directions achieving baseline scores of 5 and 10 SacreBLEU using both a baseline model (mT-big) and our method are shown in Table 14.

Model	H2H	H2M	H2L	M2H	M2M	M2L	L2H	L2M	L2L
mT-big	1.5	2.2	1.3	1.8	2.4	1.3	2.6	3.1	1.3
Ours-big	+4.2	+4.7	+1.6	+4.1	+4.3	+1.5	+2.7	+2.8	+1.2

Table 13: SacreBLEU improvements of Neuron Specialization method (Ours) over the mT-big baseline on zero-shot translations.

Model	Num. ≥ 5	Num. ≥ 10
mT-big	37	2
Ours-big	381	95

Table 14: Number of directions that exceed 5 and 10 SacreBLEU scores for the baseline (mT-big) and our method (Ours).

A.8 Sparsity versus Performance

For the Neuron Specialization, we dynamically select specialized neurons via a cumulative activa-

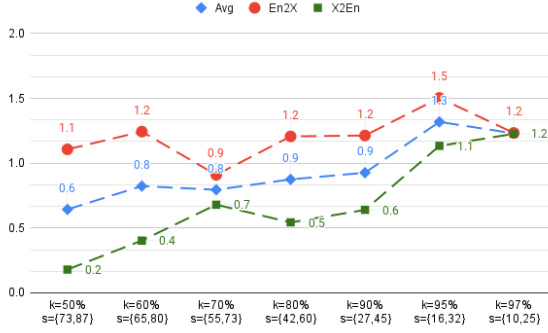


Figure 4: Improvements of Neuron Specialization method over the mT-large baseline on EC30. The x-axis indicates the factor k and the dynamic sparsity of the fc1 layer, with displayed values ranging from minimum to maximum sparsity achieved. The y-axis indicates the SacreBLEU improvements over the mT-large model.

tion threshold k in Equation 1, which is the only hyper-parameter of our method. Here, we discuss the impact of k on the final performance and its relationship to the sparsity. As mentioned in Section 3.1, a smaller factor k results in more sparse specialized neuron selection, which makes the fc1 weight more sparse as well in the Neuron Specialization Training process. In Figure 4, we show that our method consistently outperforms the baseline across a range of k values, from 50 to 97. This demonstrates robust positive gains, suggesting that our method is stable across various k settings.

In addition, we show that increasing k leads to higher improvements in general, and the optimal performance is about when $k=95\%$. Such observation follows the intuition since when k is too low, model capacity will be largely reduced. Moreover, we find that when the FFN capacity is significantly reduced (k being very small), we still observe performance gains. Notably, even when 70%-83% of FFN weights are zeroed out (as shown in Figure 4), our method still achieves an increase of +0.6 SacreBLEU. These results indicate that our method can deliver consistent and positive gains without extensive hyperparameter tuning.

Furthermore, in Figure 5, we show that the sparsity of the network presents an intuitive structure: the sparsity decreases in the Encoder and increases in the Decoder. This implies the natural signal within the pre-trained multilingual model that neurons progress from language-specific to language-agnostic in the Encoder, and vice versa in the Decoder. Such observation is natural because it is reflected by the untouched network, similar to what

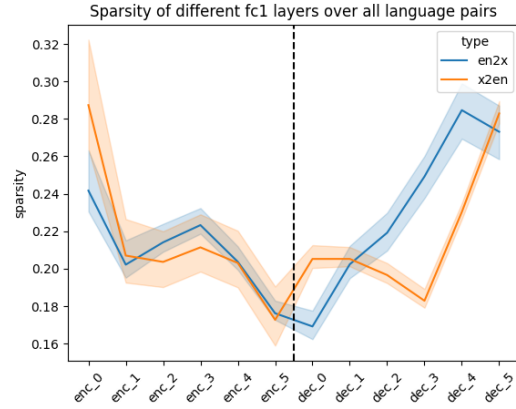


Figure 5: Sparsity progression of Neuron Specialization when $k = 95$ on the EC30. We observe that the sparsity becomes smaller in the Encoder and then goes up in the Decoder. Note that this figure is based on the natural signals extracted from the untouched pre-trained model, and will be leveraged later in the process of Neuron Specialization Training. This intrinsic pattern naturally follows our intuition that specialized neurons progress from language specific to agnostic the in Encoder, and vice versa in the Decoder.

we observed in the Progression of Neuron overlaps in Section 3.2.2.

A.9 Visualization Details

We provide the additional Pairwise Intersection over Union (IoU) scores for specialized neurons in the first Encoder layer (Figure 6), last Encoder layer (Figure 7), and last Decoder layer (Figure 8). The figures show that the Neurons gradually changed from language-specific to language-agnostic in the Encoder, and vice versa in the Decoder.

A.10 Pseudocode of Neuron Specialization

We provide the pseudocode of our proposed method, *Neuron Specialization*. We present the process of Specialized Neuron Identification in Algorithm. 1 and Neuron Specialization Training in Algorithm. 2.

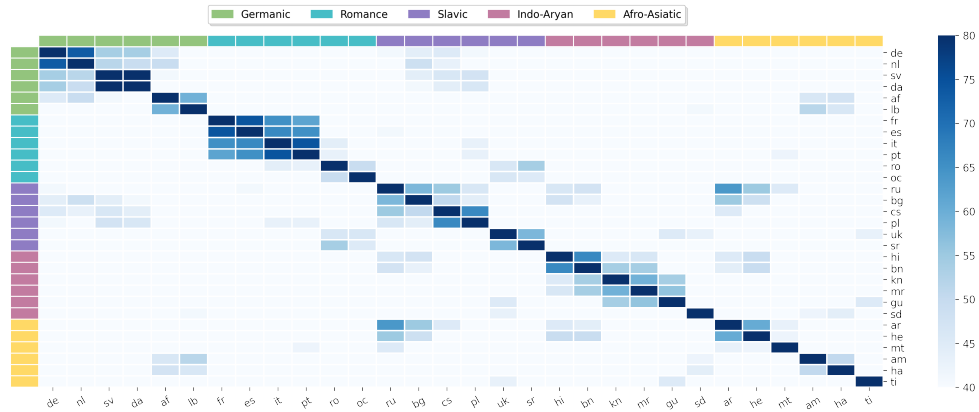


Figure 6: Pairwise Intersection over Union (IoU) scores for specialized neurons extracted from the **first encoder** FFN layer across all X-En language pairs to measure the degree of overlap between language pairs. Darker cells indicate stronger overlap, with the color threshold set from 40 to 80 to improve visibility.

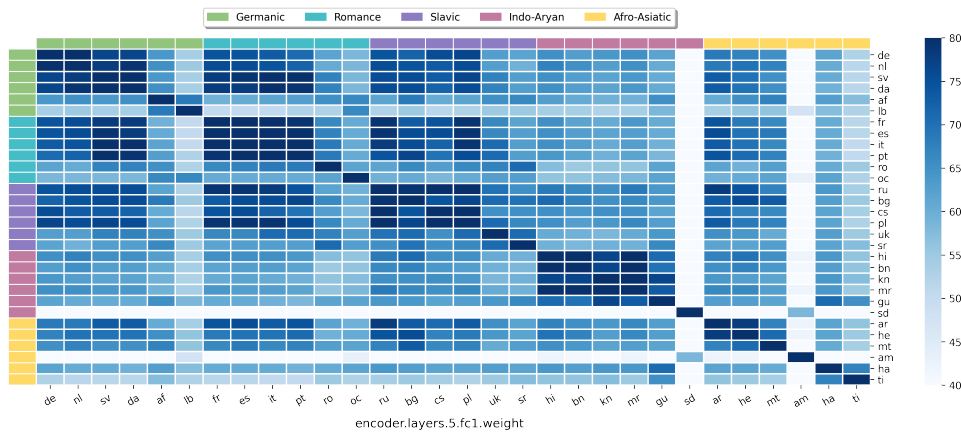


Figure 7: Pairwise Intersection over Union (IoU) scores for specialized neurons extracted from the **last encoder** FFN layer across all One-to-Many language pairs to measure the degree of overlap between language pairs. Darker cells indicate stronger overlap, with the color threshold set from 40 to 80 to improve visibility.

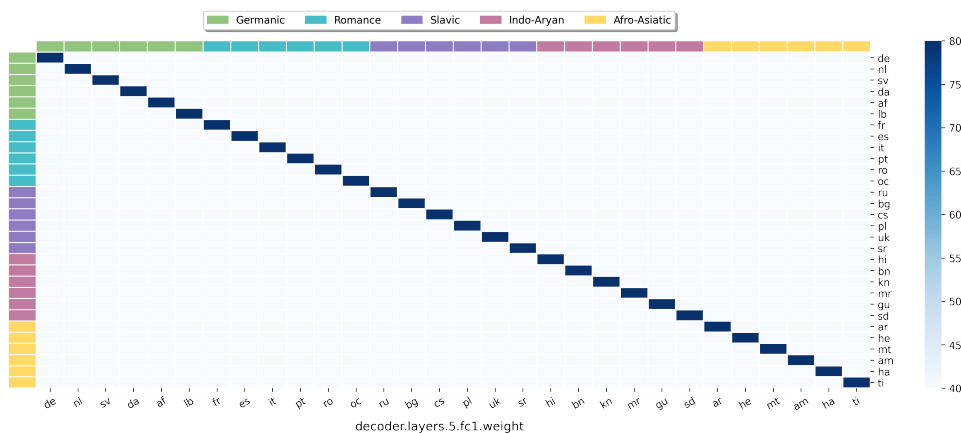


Figure 8: Pairwise Intersection over Union (IoU) scores for specialized neurons extracted from the **last decoder** FFN layer across all X-En language pairs to measure the degree of overlap between language pairs. Darker cells indicate stronger overlap, with the color threshold set from 40 to 80 to improve visibility.

Algorithm 1 Specialized Neuron Identification

- 1: **Input:** A pre-trained multi-task model θ with dimensions d and d_{ff} ; a validation dataset D with T tasks, where $D = \{D_1, \dots, D_T\}$; and an accumulation threshold factor $k \in [0\%, 100\%]$ as the only hyper-parameter.
 - 2: **Output:** A set of selected specialized neurons S_k^t for each task t .
 - 3: **for** task t in T **do**
 - 4: **Step 1: Activation Recording**
 - 5: Initialize activation vector $A_t = \mathbf{0} \in \mathbb{R}^{d_{ff}}$
 - 6: **for** sample x_i in D_t **do**
 - 7: Record activation state $a_i^t \in \mathbb{R}^{d_{ff}}$
 - 8: $A_t = A_t + a_i^t$ ▷ Accumulate activation states
 - 9: **end for**
 - 10: $a^t = \frac{A_t}{|D_t|}$ ▷ Compute average activation state for task t
 - 11: **Step 2: Neuron Selection**
 - 12: Initialize selected neurons set $S_k^t = \emptyset$
 - 13: **while** selection condition not met **do** ▷ Refer to Eq. 1 for condition
 - 14: Select neurons based on a^t and add them to S_k^t
 - 15: **end while**
 - 16: **end for**
-

Algorithm 2 Neuron Specialization Training

- 1: **Input:** A pre-trained multi-task model θ with dimensions d and d_{ff} . Corpora data C with T tasks that contain both training and validation data. A set of selected specialized neurons S_k^t for each task t .
 - 2: **Output:** A new specialized network θ^{new} . Note that only the fc1 weight matrix will be trained task-specifically, the other parameters are shared across tasks. In addition, θ^{new} does not contain more trainable parameters than θ due to the sparse network feature.
 - 3: Derive boolean mask $m^t \in \{0, 1\}^{d_{ff}}$ from S_k^t for each layer
 - 4: **while** θ^{new} not converge **do**
 - 5: **for** task t in T **do**
 - 6: $W_1^T = m^t \cdot W_1^\theta$ ▷ We perform this for all layers, refer to EQ. 3
 - 7: Train θ^{new} using C^t ▷ All parameters will be updated, yet fc1 layers are task specific
 - 8: **end for**
 - 9: **end while**
-