

# FAST: Feed-Forward Assisted Transformers for Time Efficient Fine Tuning

Anonymous ACL submission

## Abstract

Fine-tuning is the standard approach when adapting pre-trained large language models for specific downstream tasks. However, the energy and time required to fully fine-tune all parameters can become prohibitively large for many applications as the size of the model increases. While recent advancements in parameter-efficient transfer learning have reduced the number of parameters that need to be updated, the training time and energy consumption of these methods remain similar to full fine-tuning. In this paper, we propose a time-efficient fine-tuning method based on feature-extraction in which we treat off-the-shelf language models as fixed sources of embeddings and train small feed-forward networks on top for each downstream task. Averaging across the GLUE NLI benchmark, our method trains 124 times faster than full fine-tuning and 101 times faster than parameter-efficient fine-tuning methods using *distilRoBERTa*<sup>1</sup>, while achieving 81.9% and 85.0% performance respectively.

## 1 Introduction

Enhancing computational efficiency stands as a significant challenge within deep learning (Rolnick et al., 2019; Strubell et al., 2019; Min et al., 2021). The issue of efficiency is particularly prevalent in the subfield of natural language processing (NLP), where it is common practice to fully fine-tune pre-trained large language models (LLMs) for specific downstream tasks (Hendrycks et al., 2019; Andreassen et al., 2021; Miller et al., 2021), where full fine-tuning refers to fine-tuning all model parameters. However, full fine-tuning of LLMs becomes prohibitively costly as models become increasingly complex (Ding et al., 2023).

Parameter-efficient fine-tuning (PEFT) methods aim to address this by fine-tuning only a small subset of model parameters and freezing the rest (Ding

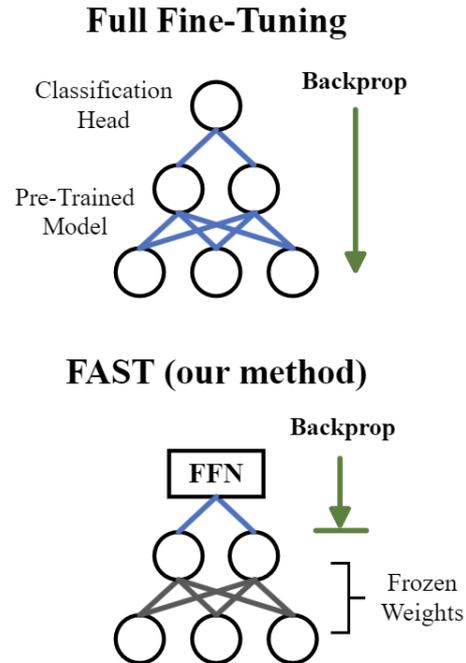


Figure 1: Visualization of FAST and comparison to standard fine-tuning.

et al., 2023; Mudrakarta et al., 2018). Common PEFT methods include adapters (Houlsby et al., 2019), prefix tuning (Li and Liang, 2021) and soft prompt tuning (Lester et al., 2021).<sup>2</sup> Efficiency for PEFT methods is generally quantified by a reduction in the number of parameters used (Howard et al., 2017; Sandler et al., 2018; Henderson et al., 2020). In practice, however, training times are also an important metric to consider for efficiency (Cao et al., 2020) – lengthy training times and substantial energy consumption commonly hinder progress in the field, and the growing computational burden of deep learning is projected to become prohibitive in many academic and industry applications (Thompson et al., 2020). Fu et al. (2022) demonstrated that commonly used PEFT methods incur longer

<sup>1</sup><https://huggingface.co/distilbert/distilroberta-base>

<sup>2</sup>Discussed in more detail in §2

056 training times than full fine-tuning, despite using  
057 fewer parameters.

058 A simple approach to avoid the prohibitive time  
059 cost of PEFT is to "freeze" the weights of the pre-  
060 trained model and train a separate network on top  
061 of the model’s embeddings (Min et al., 2021). This  
062 is reminiscent of classic NLP feature-based tech-  
063 niques (Koehn et al., 2003). In particular, the ap-  
064 proach involves choosing the architecture of the  
065 additional network and the specific model embed-  
066 dings. The simplest approach fine-tunes a single  
067 layer classification head on the last layer embed-  
068 ding of the transformer, commonly referred to as  
069 linear probing (Kumar et al., 2022).

070 In this work, we introduce FAST (Feed-Forward  
071 ASsisted Transformer), which maintains the effi-  
072 ciency of feature-based fine-tuning, but enhances  
073 performance by utilizing more recent embedding  
074 methods and stacking layers of MLP on the outputs  
075 of a pre-trained language model. We present the  
076 following core contributions:

- 077 • An efficiency focused analysis of FAST,  
078 which represents a 80-170x speedup com-  
079 pared to full and common parameter-efficient  
080 tuning methods
- 081 • A comprehensive view of embedding extrac-  
082 tion from pre-trained transformers, including  
083 sentence transformers
- 084 • Comparison of different concatenation meth-  
085 ods for transformer-based embeddings on  
086 multi-sentence tasks
- 087 • Analysis of intrinsic dimensionality of vari-  
088 ous embedding methods and their relation to  
089 classification ability

## 090 2 Related Work

091 **Adapters** In adapter fine-tuning, small feed-  
092 forward network modules called *adapters* are in-  
093 serted between transformer layers (Houlsby et al.,  
094 2019). The network modules include fully-  
095 connected down and up projections in each layer,  
096 where an input feature  $h \in \mathbb{R}^d$  is mapped to a  
097 lower  $r$ -dimensional space using a transformation  
098 matrix  $W_d \in \mathbb{R}^{d \times r}$ . Subsequently, a nonlinear ac-  
099 tivation function  $f(\cdot)$  is applied, and the feature is  
100 mapped back to the original  $d$ -dimensional space  
101 using  $W_r \in \mathbb{R}^{r \times d}$ . A residual connection then adds  
102 the original input feature to obtain the final output:

$$103 \quad h' = f(h * W_d) * W_r + h \quad (1)$$

**Prompt Tuning** Prompt tuning is a method that  
104 incorporates task-specific prompts into frozen input  
105 embeddings during fine-tuning . There are two  
106 prominent categories within prompt tuning: prefix  
107 fine-tuning and soft prompt tuning. In prefix fine-  
108 tuning, a set of trainable tokens is prepended to the  
109 inputs and hidden states of each transformer layer  
110 (Li and Liang, 2021). Soft prompt tuning simply  
111 appends a trainable prompt to the original input of  
112 the language model (Lester et al., 2021).  
113

**Feature Extraction** Peters et al. (2019) com-  
114 pares full fine-tuning and feature extraction meth-  
115 ods. They propose the use of more complex archi-  
116 tectures on top of frozen pre-trained models,  
117 varying for different types of NLP tasks. These in-  
118 clude bi-attentive classification networks and BiL-  
119 STMs on a linear combination of multiple different  
120 hidden states from the pre-trained model. Wang  
121 et al. (2019) propose the use of both fine-tuning and  
122 feature extraction, training the weights of the net-  
123 work and an additional LSTM to further increase  
124 benchmark scores. However, both primarily fo-  
125 cus on analyzing accuracy without comparing the  
126 efficiency of their proposed methods.  
127

## 128 3 Method

129 This section focuses on the two main aspects of our  
130 model: the feed-forward network and the method  
131 of generating embeddings from the base pre-trained  
132 language model.

### 133 3.1 Feed-Forward Network

134 We perform feature extraction on pre-trained large  
135 language model with frozen weights by extracting  
136 the embeddings from the last layer. We discuss  
137 methods for generating the embeddings in §3.2  
138 and §3.3. These embeddings are passed through a  
139 feed forward network, where the last layer of the  
140 feed-forward network is the classification layer.

141 For  $n$  layers and input embeddings  $X \in \mathbb{R}^{N \times d}$   
142 with  $N$  tokens and dimension  $d$ , the hidden states  
143 and output of the FFN are denoted as

$$144 \quad H_1 = \text{ReLU}(X * W_1 + B_1)$$

$$145 \quad H_2 = \text{ReLU}(H_1 * W_2 + B_2)$$

$$146 \quad \vdots$$

$$147 \quad H_i = \text{ReLU}(H_{i-1} * W_i + B_i)$$

$$148 \quad \vdots$$

$$149 \quad Y = \text{ReLU}(H_{n-1} * W_{n-1} + B_{n-1})$$

where  $W_i \in \mathbb{R}^{d \times d}$  is a learnable matrix, and  $B_i \in \mathbb{R}^d$  is the learnable bias. During implementation, the number of feed-forward layers is obtained through a grid search of the set  $\{1, 3, 5\}$ .

### 3.2 Embeddings

We explore three different techniques for extracting a fixed-length embedding from the output of the language model for a given input sequence: CLS token embedding, mean pooling, and sentence transformers.

**CLS** While RoBERTa is not pre-trained to have a semantically meaningful representation for the [CLS] token (Liu et al., 2019), we do use the [CLS] token as an embedding option to act as a baseline.

**Mean Pool** To represent the average semantic meaning of a given input sentence, We extract the output embeddings, excluding padding, from the language model and compute the element-wise mean (Chen et al., 2018).

**Sentence Transformer** Sentence transformers (Reimers and Gurevych, 2019) uses siamese network structures to obtain sentence embeddings that can be compared with cosine similarity. Sentence transformers are useful for multi-sentence tasks and are robust to zero-shot learning (Biesner et al., 2022). We use a sentence transformer for multi-sentence tasks and directly use the output embedding as the input to the feed-forward network.

### 3.3 Concatenation

For tasks with multiple input sentences, we need to generate a joint embedding to represent the entire input to pass into the feed-forward network. First, we generate embeddings for each sentence individually, denoted  $U$  and  $V$ , using one of the aforementioned methods. The way these individual embeddings are concatenated can significantly affect the performance (Reimers and Gurevych, 2019). In our approach, we use either  $|U - V|$  or  $U - V$  as the input to the FFN, depending on whether the order of input sentences for the particular task is meaningful.<sup>3</sup>

### 3.4 Benchmarking

We benchmarked performance with respect to time and energy consumption for the feed-forward net-

<sup>3</sup>For example, the QQP task in the GLUE benchmark tests if two questions are semantically identical, so the order the sentences are given to the model should not change the prediction.

work. We use real time power estimation for energy consumption (García-Martín et al., 2019).

For a given training, we recorded the time spent on each epoch of the training loop and compute the average across all epochs. We obtain our estimate for the overall training time as the average training time per epoch multiplied by the number of epochs that passed during training. The energy consumption in joules (J) from the estimated training times is

$$E = \int_t P dt \quad (2)$$

where  $P$  is the real-time power consumption output of the device in watts (W) and  $t$  is the time in seconds. The value of  $P$  is obtained directly from the user’s hardware.

## 4 Experiments and Results

Section 4.1 introduces the specific model and datasets used in our experiments. Section 4.2 reports our experimental results for accuracy. Section 4.3 reports our experimental results for computational efficiency and energy efficiency comparisons to other leading delta-tuning measures. Section 4.4 compares the performance of different concatenation methods on multi-sentence tasks. Section 4.5 compares how our method’s performance generalizes to other pre-trained models. Section 4.6 utilizes manifold learning techniques to analyze different embedding techniques, and help explain our results.

### 4.1 Experimental Setup

**Dataset and Evaluation Metrics** We evaluated the performance of our method across 10 tasks sourced from the General Language Understanding Evaluation (GLUE) dataset<sup>4</sup>. This dataset comprises a diverse set of linguistic tasks used to assess the generalization and proficiency of natural language understanding models (Wang et al., 2018b). To evaluate computational efficiency, we report the estimated time cost and hardware energy consumption of our experiments.

**Implementation** We chose distilRoBERTa, a distilled version of the RoBERTa model, as our base model to generate embeddings. It follows the same training procedure as distilBERT, as outlined in the associated paper (Sanh et al., 2019). The model’s architecture consists of 6 transformer layers, each

<sup>4</sup>See Appendix B

	GLUE Score	CoLA	SST-2	MRPC	STS-B	QQP	MNLI <sub>m</sub>	QNLI	RTE	WNLI
CLS	60.9	35.3	<b>87.1</b>	<b>78.6</b>	31.3	52.9	66.3	71.2	50.3	<b>65.1</b>
Mean Pool	<b>64.7</b>	<b>38.9</b>	86.9	70.6	66.1	<b>59.3</b>	<b>67.3</b>	<b>74.1</b>	50.3	<b>65.1</b>
ST	63.4	15.9	84.8	76.0	<b>80.6</b>	54.7	66.8	71.9	<b>55.9</b>	58.2
DistilRoBERTa	83.2	59.3	92.5	86.6	88.3	89.4	84.0	90.8	67.9	-
BERT Adapter	80.2	59.2	94.3	84.3	86.1	89.4	85.4	92.4	71.6	65.1
CBoW	58.6	0.0	80.0	73.4	61.2	51.4	56.0	72.1	54.1	65.1

Table 1: GLUE Benchmark scores (MCC, F1, accuracy) of the CLS, Mean Pool, and Sentence Transformer methods, and comparisons with other methods. CoLA is evaluated using MCC. MRPC and QQP are evaluated using F1 score. STS-B uses Spearman’s correlation coefficient. The other tasks use accuracy.

with 12 heads and an embedding dimension of 768. We selected this base model for its accessibility, strong performance on NLP tasks, cost-effectiveness in terms of computational resources, and the availability of a publicly accessible pre-trained sentence transformer.

The weights of distilRoBERTa remain frozen during tuning, allowing us to precompute and store the model’s output embeddings for all sentences in the dataset. This initial computation serves as a one-time cost before training our additional parameters directly. The saved embeddings are used during training of the feed-forward network, eliminating the need to run inference on the language model in each epoch. To optimize the FFN’s performance, we perform a hyperparameter grid search to determine the best values for number of layers (1, 3, or 5), batch size, learning rate, and weight decay<sup>5</sup>. All computation was performed on one V100 GPU.

## 4.2 GLUE Benchmark Results

We assessed our method’s performance against both state-of-the-art models and baseline techniques that undergo training within a similar computational time frame. Specifically, we compare against a full fine-tune of DistilRoBERTa, adapters, and CBoW. The performance outcomes across GLUE tasks are reported in Table 1.

Overall, we observed that mean pooling tends to yield the best performance among the three types of embeddings. This is likely because the pre-trained model is optimized to generate meaningful embeddings for predicting the next token. Since mean pooling incorporates information from all tokens in the sequence, it may also provide a more consistent representation of the entire text regardless of the downstream task. On the other hand, RoBERTa’s

<sup>5</sup>See Appendix A

start [CLS] token is not pretrained on any task, so distilRoBERTa likely lacks a meaningful [CLS] token.

## 4.3 Energy Efficiency

In Table 2, we evaluated time cost and energy consumption across GLUE tasks for our method, as well as for DistilRoBERTa and CBoW, according to the method described in Section 3.4. Note that generating the embeddings of the language model is not included in the training time, though this is a one-time cost that is amortized across training runs. We reported these efficiency metrics for 1, 3 and 5 layers of our FFN. FAST demonstrates a 124x training speedup over fully fine-tuning all model weights of a distilRoBERTa transformer model. FAST also demonstrated a 101x speedup over common PEFT methods. Due to this reduced training time, the energy consumption of FAST is also 124x less than full fine-tuning and 101x less than PEFT methods.

Figure 2 illustrates that although FAST and common PEFT methods utilize similar parameter counts, training for FAST is two orders of magnitude faster than that of PEFT methods. Furthermore, both FAST and PEFT outperform full fine-tuning in terms of parameter count and training time. This demonstrates that FAST provides greater time efficiency at comparable parameter counts to common PEFT methods.

## 4.4 Concatenation Method

For multi-sentence tasks, we experimented with various methods of concatenating sentence embeddings. Previous work has shown significant impact from different aggregations of sentence embeddings on performance (Reimers and Gurevych, 2019). Table 3 presents a comparative analysis of different concatenation methods (as before,  $U$  and

	Time (s)			Energy (J)		
	CoLA	MRPC	QNLI	CoLA	MRPC	QNLI
1-layer FAST	0.72	0.31	7.96	181	76	1990
3-layer FAST	0.98	0.45	12.09	246	111	3023
5-layer FAST	1.29	0.57	15.93	323	143	3983
Prompt Tuning	92.08	39.72	1125.00	23020	9930	281250
Prefix Tuning	92.72	40.10	1136.00	23180	10025	284000
Adapter	101.53	43.67	1221.80	25383	10918	305450
Full Fine-Tuning	116.45	50.13	1408.00	29113	12534	352000

Table 2: Training times (s) and energy consumption (J) across GLUE tasks for the CLS, Mean Pool, and Sentence Transformer methods on a V100 GPU, and comparison with full fine-tuning.

		MRPC	STS-B	QQP	QNLI	RTE	WNLI
Mean Pooling	$U, V$	72.55	33.57	74.98	68.00	56.68	56.34
	$ U - V $	80.15	78.89	75.87	73.16	59.57	56.34
	$U, V,  U - V $	72.55	67.35	76.91	70.90	57.04	56.34
CLS	$U, V$	71.08	26.80	73.15	62.75	56.68	56.34
	$ U - V $	79.41	74.57	76.38	70.95	58.12	56.34
	$U, V,  U - V $	71.32	71.35	76.92	67.16	56.34	56.34

Table 3: Comparison of different concatenation methods using distilRoBERTa across multi-sentence GLUE tasks, using both CLS and Mean Pooled embeddings.

$V$  denote the two individual sentence embeddings) for different multi-sentence GLUE tasks. Employing  $U - V$  notably enhanced performance, and so was used in our overall performance benchmark comparing to full fine-tuning.

#### 4.5 Generalization to other Pre-Trained Models

To ensure that our method generalizes to models beyond just distilRoBERTa, we ran experiments using our same method on a different public transformer: MPNet (Song et al., 2020). Table 4 shows the results across 8 GLUE benchmark tasks, where we observe similar patterns to our method’s performance with distilRoBERTa. Mean pooling and sentence transformers consistently outperform CLS, with sentence transformers being particularly effective at regression tasks like STS-B, but poor at grammar tasks like CoLA.

#### 4.6 Dimensionality Analysis

We use dimensional analysis and manifold learning techniques to establish a quantifiable basis for comparing different embedding models (distilRoBERTaST, distilRoBERTaBase (mean pooled), and CBoW), as well as a way to conceptually explain the efficacy of the FFN.

Previous work from Tulchinskii et al. (2023) and Mamou et al. (2020) also suggests that of human language is intrinsically represented by a small number of dimensions. We posit that a lower intrinsic dimensionality may indicate a more straightforward task for the FFN insofar as the embeddings efficiently encapsulate the essential information required for the model to generalize effectively (Aghajanyan et al., 2020).

Note that we do not necessarily directly correlate lower dimensionality with enhanced performance, as the actual performance of FAST varies based on several other task-dependent and FFN-dependent parameters.

**PCA Analysis:** The embeddings are processed through Randomized PCA and the cumulative variance ratio is calculated to determine the RandPCA intrinsic dimensionality of the embedding, i.e. the minimum percentage of components required to surpass the 90% cumulative variance threshold.

Our results (Table 5) indicate that mean pooling and sentence transformers consistently generate embeddings that have lower intrinsic dimensionality than cBOW, suggesting that such embeddings are more effective at encoding critical information.

**t-SNE Analysis:** The t-SNE visualization provides an intuitive understanding of the data’s seg-

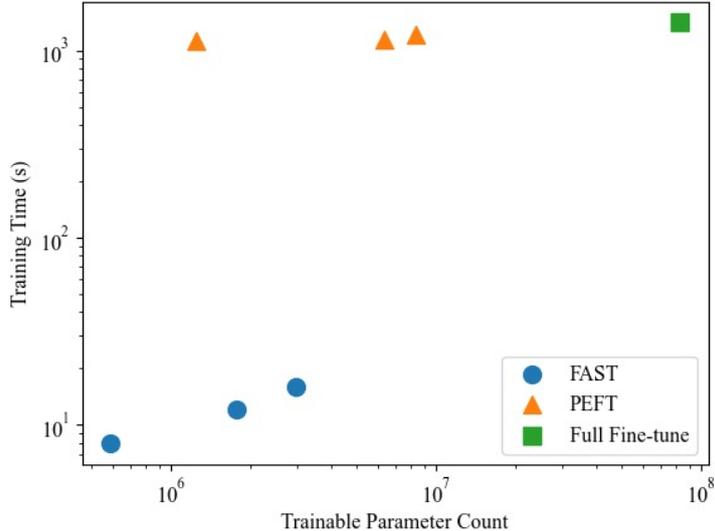


Figure 2: Log scale plot of training time (seconds) vs. trainable parameter count on the QNLI task for FAST (using MPNetBase and MPNetST), three common PEFT methods (adapters, prefix tuning, and soft prompt tuning), and full fine-tuning.

	CoLA	SST-2	MRPC	STS-B	QQP	QNLI	RTE	WNLI
MPNetBase (CLS)	41.16	85.44	71.08	65.73	75.31	65.84	56.32	56.34
MPNetBase (Meanpool)	<b>45.21</b>	89.45	<b>79.90</b>	75.96	77.05	71.48	62.46	56.34
MPNetSentence	27.70	<b>90.02</b>	77.21	<b>83.53</b>	<b>87.92</b>	<b>73.73</b>	<b>68.23</b>	<b>60.56</b>

Table 4: GLUE benchmark scores using FAST with embeddings from MPNet Base Transformer and MPNet Sentence Transformer.

367 mentation and separation (van der Maaten and  
 368 Hinton, 2008). We observe the reduction in the  
 369 intrinsic dimensionality of the embeddings (or lack  
 370 thereof) as they move forward through the FFN.  
 371 Our analyses (specific examples using the SST2  
 372 task embeddings are in Fig. 3) indicate that our  
 373 simple FFN is sufficiently parameterized to dra-  
 374 matically reduce the intrinsic dimensionality of  
 375 the embeddings generated from both mean pool-  
 376 ing and sentence transformers, but that it struggles  
 377 to unravel the underlying structure of the CBoW-  
 378 generated embeddings. Our results further sup-  
 379 port the hypothesis that the intrinsic dimensionality  
 380 from the last hidden layer of the FFN is a strong  
 381 indicator of the classification and generalization  
 382 ability of the model (Ansuini et al., 2019). This is  
 383 consistent with the superior performance of mean-  
 384 pooled embeddings and sentence transformer em-  
 385 beddings in testing.

## 386 5 Discussion

387 Our results suggest that our proposed method has  
 388 promising implications for improving the accessi-

389 bility of fine-tuning LLMs, particularly for those  
 390 with limited computational resources. Most nota-  
 391 bly, our method offers substantially improved  
 392 training efficiency, while maintaining relatively  
 393 high levels of performance. Our method is nota-  
 394 bly compatible with inexpensive or resource-  
 395 constrained hardware, which could empower indi-  
 396 viduals without access to high-end computational  
 397 resources to engage in meaningful natural language  
 398 processing tasks.

399 In this paper, we propose FAST (Feed-Forward  
 400 Assisted Transformer), a novel addition-based delta  
 401 tuning method that trains small feed forward neu-  
 402 ral networks for specific downstream tasks on  
 403 top of embeddings generated from pre-trained off-  
 404 the-shelf large language models. Our approach  
 405 achieves up to 93.9% accuracy of full fine-tuning  
 406 and accomplishes downstream tasks up to 170  
 407 times faster. Our method significantly reduces  
 408 training time compared to alternative methods in  
 409 parameter-efficient fine tuning, enabling more ac-  
 410 cessible natural language processing for users with  
 411 limited hardware capabilities.

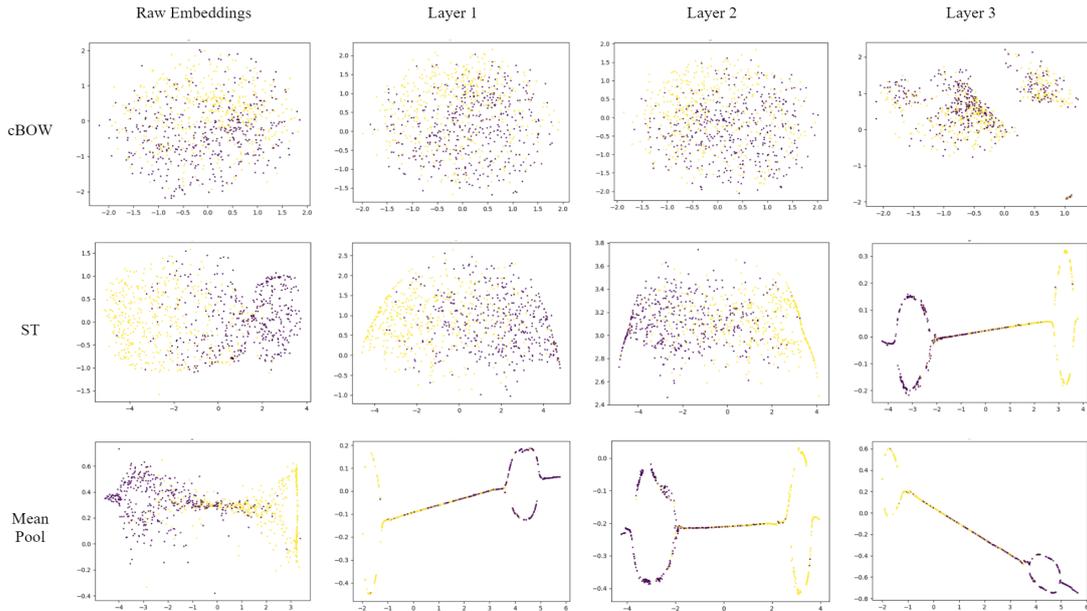


Figure 3: t-SNE plots of validation set SST2 embeddings at incremental layers of the trained FFN with 3 hidden layers.

	CoLA	SST-2	MRPC	STS-B	QQP	MNLI	QNLI	RTE
Mean Pool	<b>21.22</b>	<b>16.67</b>	<b>54.30</b>	<b>55.34</b>	<b>52.60</b>	<b>50.52</b>	<b>28.26</b>	<b>23.44</b>
ST	28.91	34.11	66.67	69.79	64.58	73.57	37.37	34.51
CBoW	56.00	62.00	80.00	81.00	78.00	79.00	66.00	65.00

Table 5: RandPCA intrinsic dimensionality of the Mean Pool, Sentence Transformer, and Continuous Bag of Words embeddings.

## Limitations

When interpreting our findings, there are a number of experimental limitations that should be considered. Firstly, our experiments were conducted on a V100 GPU. To gain a more comprehensive understanding of the performance and efficiency of our method, it would be beneficial to conduct tests on various machines and assess potential variations across different hardware configurations. Although the GLUE benchmark is a widely used standard in NLP research, it is ultimately limited to short, English language text inference tasks. Exploring tasks that involve larger levels of text granularity (e.g. paragraphs), as well as incorporating non-English language data, could provide valuable insights and enhance the generalizability of our method to a broader spectrum of linguistic challenges. We could also test our model against the SuperGLUE benchmark to assess performance on more challenging language tasks compared to the GLUE benchmark (Sarlin et al., 2019), which may yield more robust and comprehensive results. In

addition, future research could explore additional types of embeddings and embedding concatenation strategies to assess their impact on model performance, as different kinds of inputs to the FFN may capture distinct linguistic features. Finally, we utilized time and energy consumption during training as our measure of efficiency. However, investigating other efficiency metrics, such as FLOP counts or memory consumption (Schwartz et al., 2019; Henderson et al., 2020), could yield further insights into the trade-offs between model performance and computational cost.

We do not foresee any unique risks that are not present in any research on large language models.

## References

- Armen Aghajanyan, Luke Zettlemoyer, and Sonal Gupta. 2020. [Intrinsic dimensionality explains the effectiveness of language model fine-tuning](#). *CoRR*, abs/2012.13255.
- Anders Andreassen, Yasaman Bahri, Behnam Neyshabur, and Rebecca Roelofs. 2021. [The evo-](#)

455	lution of out-of-distribution robustness throughout	of the North American Chapter of the Association for	509
456	fine-tuning. <i>CoRR</i> , abs/2106.15831.	<i>Computational Linguistics</i> , pages 127–133.	510
457	Alessio Ansuini, Alessandro Laio, Jakob H. Macke,	Ananya Kumar, Aditi Raghunathan, Robbie Jones,	511
458	and Davide Zoccolan. 2019. <a href="#">Intrinsic dimension of</a>	Tengyu Ma, and Percy Liang. 2022. Fine-tuning	512
459	<a href="#">data representations in deep neural networks</a> . <i>CoRR</i> ,	can distort pretrained features and underperform out-	513
460	abs/1905.12784.	of-distribution. <i>arXiv:2202.10054</i> .	514
461	David Biesner, Maren Pielka, Rajkumar Ramamurthy,	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021.	515
462	Tim Dilmaghani, Bernd Kliem, Rüdiger Loitz, and	<a href="#">The power of scale for parameter-efficient prompt</a>	516
463	Rafet Sifa. 2022. Zero-shot text matching for auto-	tuning. <i>CoRR</i> , abs/2104.08691.	517
464	mated auditing using sentence transformers. <i>Compu-</i>	Xiang Lisa Li and Percy Liang. 2021. <a href="#">Prefix-</a>	518
465	<i>tation and Language</i> .	tuning: <a href="#">Optimizing continuous prompts for gener-</a>	519
466	Qingqing Cao, Aruna Balasubramanian, and Niranjan	ation. <i>CoRR</i> , abs/2101.00190.	520
467	Balasubramanian. 2020. <a href="#">Towards accurate and re-</a>	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	521
468	<a href="#">liable energy measurement of NLP models</a> . <i>CoRR</i> ,	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	522
469	abs/2010.05248.	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	523
470	Qian Chen, Zhen-Hua Ling, and Xiaodan Zhu. 2018.	<a href="#">Roberta: A robustly optimized BERT pretraining</a>	524
471	Enhancing sentence embedding with generalized	approach. <i>CoRR</i> , abs/1907.11692.	525
472	pooling. <i>Computation and Language</i> .	Jonathan Mamou, Hang Le, Miguel Del Rio, Cory	526
473	Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zong-	Stephenson, Hanlin Tang, Yoon Kim, and SueYeon	527
474	han Yang, Yusheng Su, Shengding Hu, Yulin Chen,	Chung. 2020. <a href="#">Emergence of separable manifolds in</a>	528
475	Chi-Min Chan, Weize Chen, Jing Yi, Weilin Zhao,	<a href="#">deep language representations</a> .	529
476	Xiaozhi Wang, Zhiyuan Liu, Hai-Tao Zheng, Jianfei	John Miller, Rohan Taori, Aditi Raghunathan, Shiori	530
477	Chen, Yang Liu, Jie Tang, Juanzi Li, and Maosong	Sagawa, Pang Wei Koh, Vaishaal Shankar, Percy	531
478	Sun. 2023. <a href="#">Parameter-efficient fine-tuning of large-</a>	Liang, Yair Carmon, and Ludwig Schmidt. 2021.	532
479	<a href="#">scale pre-trained language models</a> . <i>Nature Machine</i>	<a href="#">Accuracy on the line: On the strong correlation between</a>	533
480	<i>Intelligence</i> , 5(3):220–235.	<a href="#">out-of-distribution and in-distribution generalization</a> .	534
481	Zihao Fu, Haoran Yang, Anthony Man-Cho So, Wai	<i>CoRR</i> , abs/2107.04649.	535
482	Lam, Lidong Bing, and Nigel Collier. 2022. <a href="#">On the</a>	Bonan Min, Hayley Ross, Elior Sulem, Amir	536
483	<a href="#">effectiveness of parameter-efficient fine-tuning</a> .	Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz,	537
484	Eva García-Martín, Crefeda Faviola Rodrigues, Graham	Eneko Agirre, Ilana Heintz, and Dan Roth. 2021.	538
485	Riley, and Håkan Grahn. 2019. Estimation of en-	<a href="#">Recent advances in natural language processing via</a>	539
486	ergy consumption in machine learning. <i>Journal of</i>	<a href="#">large pre-trained language models: A survey</a> . <i>CoRR</i> ,	540
487	<i>Parallel and Distributed Computing</i> .	abs/2111.01243.	541
488	Peter Henderson, Jieru Hu, Joshua Romoff, Emma	Pramod Kaushik Mudrakarta, Mark Sandler, Andrey	542
489	Brunskill, Dan Jurafsky, and Joelle Pineau. 2020.	Zhmoginov, and Andrew G. Howard. 2018. <a href="#">K for the</a>	543
490	<a href="#">Towards the systematic reporting of the energy</a>	<a href="#">price of 1: Parameter efficient multi-task and transfer</a>	544
491	<a href="#">and carbon footprints of machine learning</a> . <i>CoRR</i> ,	learning. <i>CoRR</i> , abs/1810.10703.	545
492	abs/2002.05651.	Matthew E. Peters, Sebastian Ruder, and Noah A. Smith.	546
493	Dan Hendrycks, Kimin Lee, and Mantas Mazeika. 2019.	2019. To tune or not to tune? adapting pretrained	547
494	<a href="#">Using pre-training can improve model robustness and</a>	representations to diverse tasks. <i>arXiv:1903.05987</i> .	548
495	<a href="#">uncertainty</a> . <i>CoRR</i> , abs/1901.09960.	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-bert:</a>	549
496	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,	<a href="#">Sentence embeddings using siamese bert-networks</a> .	550
497	Bruna Morrone, Quentin de Laroussilhe, Andrea Ges-	<i>CoRR</i> , abs/1908.10084.	551
498	musundo, Mona Attariyan, and Sylvain Gelly. 2019.	David Rolnick, Priya L. Donti, Lynn H. Kaack, Kelly	552
499	<a href="#">Parameter-efficient transfer learning for NLP</a> . <i>CoRR</i> ,	Kochanski, Alexandre Lacoste, Kris Sankaran, An-	553
500	abs/1902.00751.	drew Slavin Ross, Nikola Milojevic-Dupont, Natasha	554
501	Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry	Jaques, Anna Waldman-Brown, Alexandra Luc-	555
502	Kalenichenko, Weijun Wang, Tobias Weyand, Marco	cioni, Tegan Maharaj, Evan D. Sherwin, S. Karthik	556
503	Andreetto, and Hartwig Adam. 2017. <a href="#">Mobilenets:</a>	Mukkavilli, Konrad P. Körding, Carla P. Gomes, An-	557
504	<a href="#">Efficient convolutional neural networks for mobile</a>	drew Y. Ng, Demis Hassabis, John C. Platt, Felix	558
505	<a href="#">vision applications</a> . <i>CoRR</i> , abs/1704.04861.	Creutzig, Jennifer T. Chayes, and Yoshua Bengio.	559
506	Philipp Koehn, Franz J. Och, and Daniel Marcu. 2003.	2019. <a href="#">Tackling climate change with machine learn-</a>	560
507	<a href="#">Statistical phrase-based translation</a> . In <i>Proceedings</i>	ing. <i>CoRR</i> , abs/1906.05433.	561
508	<a href="#">of the 2003 Human Language Technology Conference</a>		

562 Mark Sandler, Andrew G. Howard, Menglong Zhu, An-  
563 drey Zhmoginov, and Liang-Chieh Chen. 2018. [In-](#)  
564 [verted residuals and linear bottlenecks: Mobile net-](#)  
565 [works for classification, detection and segmentation.](#)  
566 *CoRR*, abs/1801.04381.

567 Victor Sanh, Lysandre Debut, Julien Chaumond, and  
568 Thomas Wolf. 2019. [Distilbert, a distilled version](#)  
569 [of BERT: smaller, faster, cheaper and lighter.](#) *CoRR*,  
570 abs/1910.01108.

571 Paul-Edouard Sarlin, Daniel DeTone, Tomasz Mal-  
572 isiewicz, and Andrew Rabinovich. 2019. [Superglue:](#)  
573 [Learning feature matching with graph neural net-](#)  
574 [works.](#) *CoRR*, abs/1911.11763.

575 Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren  
576 Etzioni. 2019. [Green AI.](#) *CoRR*, abs/1907.10597.

577 Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-  
578 Yan Liu. 2020. [Mpnet: Masked and permuted pre-](#)  
579 [training for language understanding.](#)

580 Emma Strubell, Ananya Ganesh, and Andrew McCal-  
581 lum. 2019. [Energy and policy considerations for](#)  
582 [deep learning in NLP.](#) *CoRR*, abs/1906.02243.

583 Neil C. Thompson, Kristjan H. Greenewald, Keeheon  
584 Lee, and Gabriel F. Manso. 2020. [The computational](#)  
585 [limits of deep learning.](#) *CoRR*, abs/2007.05558.

586 Eduard Tulchinskii, Kristian Kuznetsov, Laida  
587 Kushnareva, Daniil Cherniavskii, Serguei Baran-  
588 nikov, Irina Piontkovskaya, Sergey Nikolenko,  
589 and Evgeny Burnaev. 2023. [Intrinsic dimension](#)  
590 [estimation for robust detection of ai-generated texts.](#)

591 Laurens van der Maaten and Geoffrey Hinton. 2008.  
592 [Visualizing data using t-sne.](#) *Journal of Machine*  
593 *Learning Research*, 9(86):2579–2605.

594 Alex Wang, Amanpreet Singh, Julian Michael, Fe-  
595 lix Hill, Omer Levy, and Samuel Bowman. 2018a.  
596 [GLUE: A multi-task benchmark and analysis plat-](#)  
597 [form for natural language understanding.](#) In *Proceed-*  
598 *ings of the 2018 EMNLP Workshop BlackboxNLP:*  
599 *Analyzing and Interpreting Neural Networks for NLP*,  
600 pages 353–355, Brussels, Belgium. Association for  
601 Computational Linguistics.

602 Alex Wang, Amanpreet Singh, Julian Michael, Felix  
603 Hill, Omer Levy, and Samuel R. Bowman. 2018b.  
604 [GLUE: A multi-task benchmark and analysis plat-](#)  
605 [form for natural language understanding.](#) *CoRR*,  
606 abs/1804.07461.

607 Ran Wang, Haibo Su, Chunye Wang, Kailin Ji, and  
608 Jupeng Ding. 2019. [To tune or not to tune? how](#)  
609 [about the best of both worlds?](#) *arXiv:1907.05338*.

## Appendix A Hyperparameter Search

We do a hyperparameter grid search on the following parameters. In selecting the number of epochs, we do early stopping based on dev set validation loss for a maximum of 50 epochs. Additionally, we utilize the default train/dev/test splits each task within the GLUE benchmark.

Hyperparameter	Values
Number of epochs	50
Batch size	32, 512
Initial learning rate	$1 \times 10^{-2}$ , $1 \times 10^{-3}$
Hidden layer size	$X/4$ , $X/2$ , $X$ , $2X$ , $4X$
Number of layers	1, 3, 5
Weight decay	$1 \times 10^{-2}$ , $1 \times 10^{-4}$
Patience	3

Table 6: Hyperparameter grid search values.  $X$  represents the hidden layer size, equal to 768 for all of our methods

## Appendix B Description of GLUE tasks

The tasks included are

- CoLA (Corpus of Linguistic Acceptability): Evaluates a model’s ability to determine the grammatical acceptability of sentences.
- SST-2 (Stanford Sentiment Treebank - Binary): Involves predicting sentiment labels (positive or negative) for movie reviews.
- MRPC (Microsoft Research Paraphrase Corpus): Requires identifying whether pairs of sentences are semantically equivalent or not.
- QQP (Quora Question Pairs): Focuses on determining duplicate or semantically similar questions.
- STS-B (Semantic Textual Similarity Benchmark): Involves predicting the similarity score between pairs of sentences.
- MNLI (Multi-Genre Natural Language Inference): Requires determining the logical relationship between a premise and a hypothesis in three categories: entailment, contradiction, or neutral.
- QNLI (Question Natural Language Inference): Transforms questions from SQuAD into binary sentence pair classification tasks that require determining whether a certain sentence contains the answer to a given question.

- RTE (Recognizing Textual Entailment): Requires determining if a hypothesis can be inferred from a given premise.
- WNLI (Winograd NLI): Tests a model’s ability to resolve pronouns in a sentence.

## Appendix C Use of Scientific Artifacts

Our work utilizes the GLUE benchmark dataset, as introduced by Wang et al. (2018a). The dataset is entirely in English and consists of a collection of publicly available text datasets for evaluating natural language understanding systems. It includes tasks such as textual entailment, sentiment analysis, and question answering. The GLUE benchmark ensures that any personally identifiable information is removed from the dataset, making it publicly available for research purposes.

Furthermore, FAST employs Hugging Face implementations for *distilroberta-base*, *all-distilroberta-v1*, *MPNet*, *all-mpnet-base-v2*, all of which are made publicly available for research purposes under the Apache-2.0 license.

## Appendix D t-SNE plots

t-SNE plots of raw STSB embeddings:

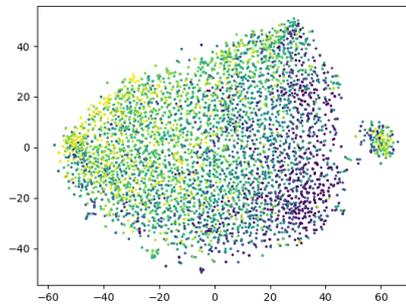


Figure 4: t-SNE plot of STSB embeddings from Mean Pooling

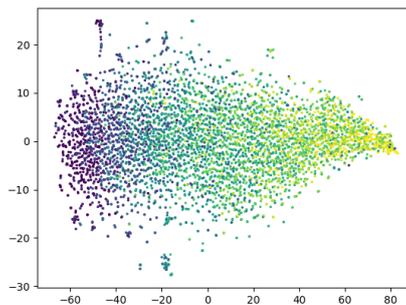


Figure 5: t-SNE plot of STSB embeddings from ST

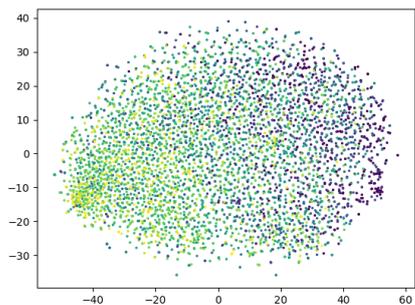


Figure 6: t-SNE plot of STSB embeddings from CBoW

667

## Appendix E GLUE Validation Dev Set

	CoLA	SST-2	MRPC	STS-B	QQP	MNLI	QNLI	RTE	WNLI
CLS + FFN	43.7	87.2	77.7 / 77.7	74.3/74.2	79.7	67.5/67.5	69.9	<b>57.4</b>	56.3
Mean Pool + FFN	<b>52.0</b>	<b>89.0</b>	<b>79.2/79.2</b>	<b>80.2/80.1</b>	<b>82.8/82.8</b>	<b>68.3/68.3</b>	<b>75.3</b>	57.0	56.3
ST + FFN	21.7	87.0	75.5/75.5	88.6/88.2	80.1/85.2	66.7/66.7	75.2	55.2	56.3
DistilRoBERTa	59.3	92.5	86.6	88.3	89.4	84.0	90.8	67.9	-
Adapters	-	-	-	-	-	-	-	-	-
CBoW	0.0	80.0	73.4/81.5	61.2/58.7	79.1/51.4	56.0/56.4	72.1	54.1	65.1

Table 7: GLUE Benchmark scores (MCC, F1, accuracy) of the CLS, Mean Pool, and Sentence Transformer methods, and comparisons with other methods. CoLA is evaluated using MCC. MRPC and QQP are evaluated using F1 score/accuracy. STS-B uses Spearman’s correlation coefficient/accuracy. The other tasks use accuracy.