

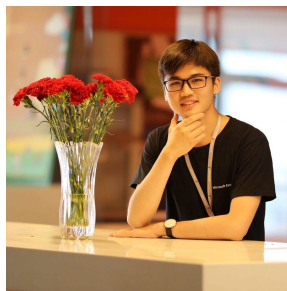


# **Structured Pruning Learns Compact and Accurate Models**

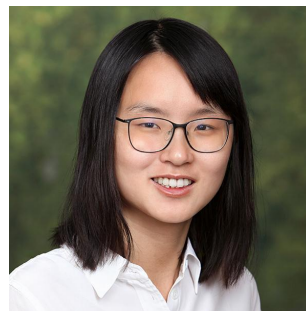
**ACL 2022**



**Mengzhou Xia**



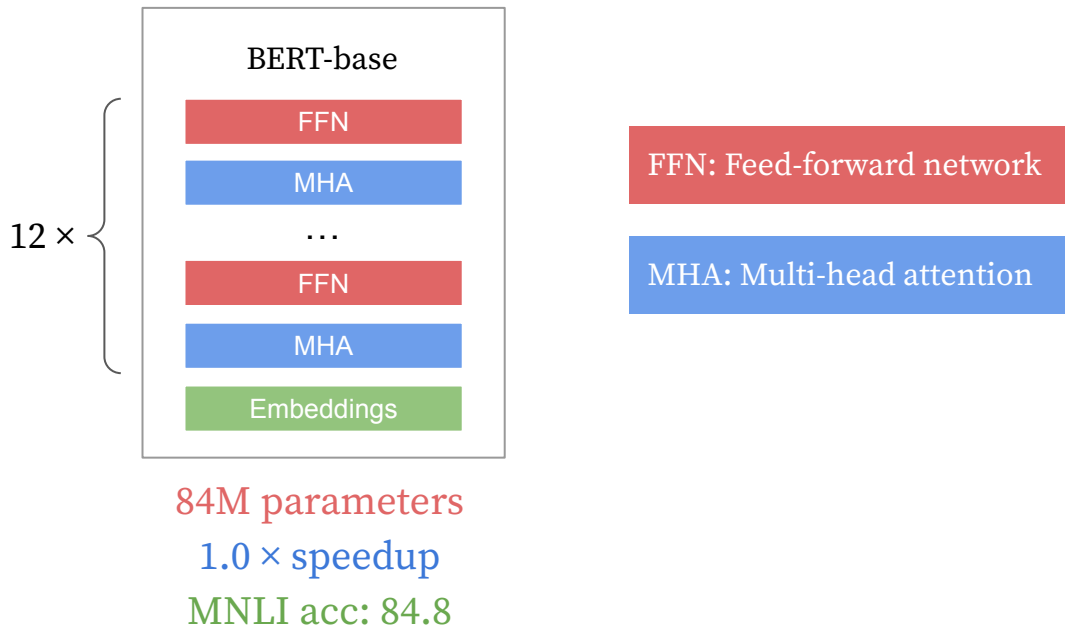
**Zexuan Zhong**



**Danqi Chen**

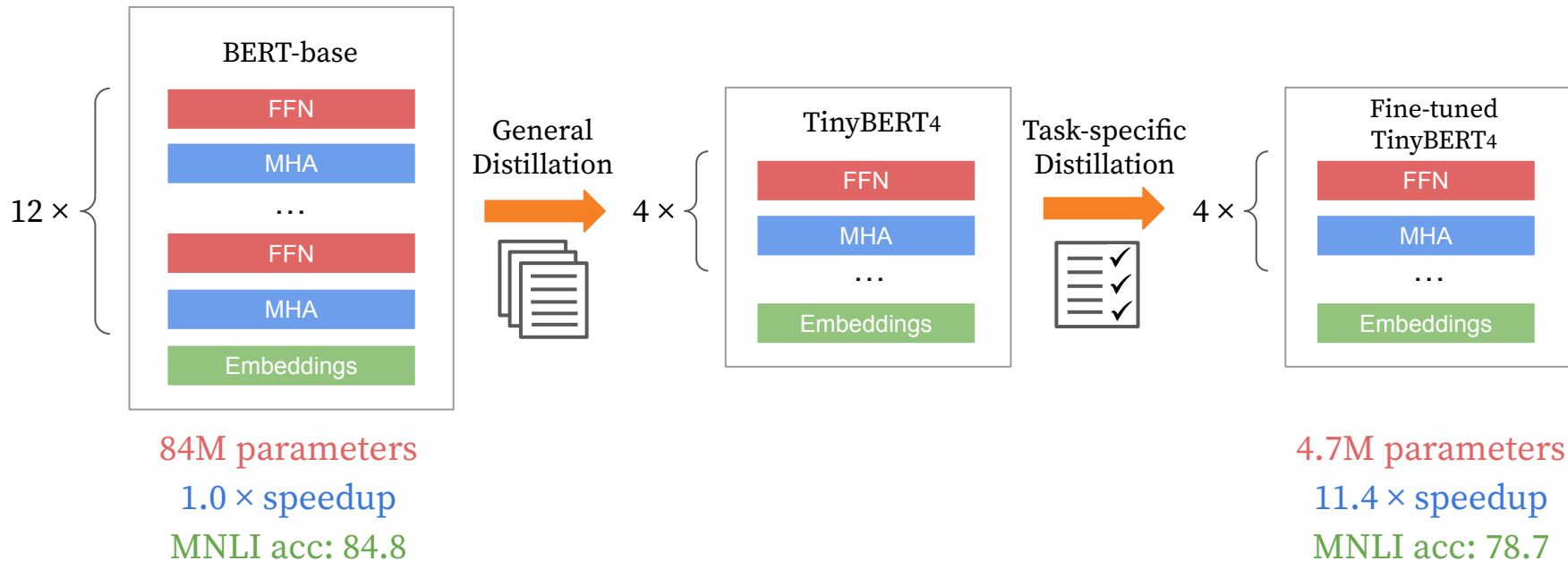
# Background: Compressing Transformer Models

Language models are known to be **overparameterized**



# Background: Compressing Transformer Models

**Distillation:** transfers knowledge from a teacher model to a fixed student model



# Background: Compressing Transformer Models

## Knowledge distillation can:

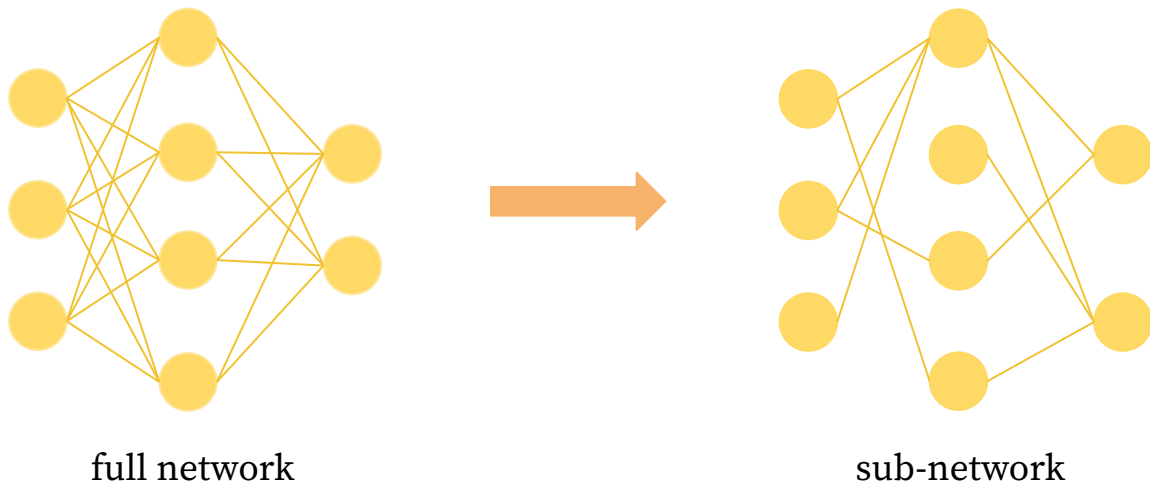
- achieve over 10× speedups

## Disadvantages of knowledge distillation:

- Pre-specified student model architecture
- Trained from scratch with **unlabeled data** and computationally expensive e.g., 350 hours for TinyBERT

# Background: Compressing Transformer Models

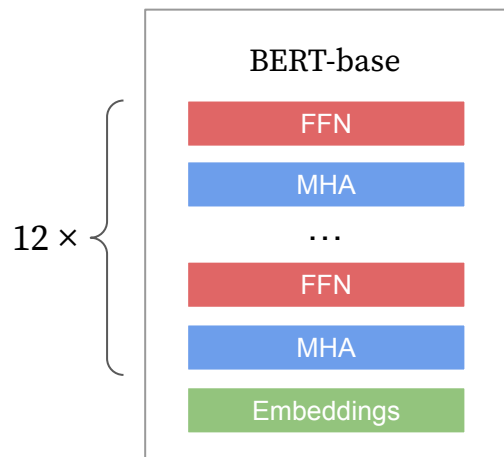
**Unstructured pruning:** Prunes individual parameters (Frankle and Carbin, 2019)



**Hard to achieve inference speedup!**

# Background: Compressing Transformer Models

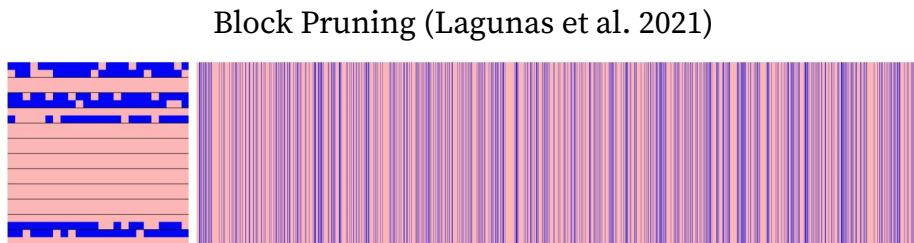
**Structured pruning:** Prunes groups of parameters, e.g., heads, FFNs, leads to actual speedups unlike unstructured pruning



84M parameters

1.0 × speedup

MNLI acc: 84.8



25M parameters

2.7 × speedup

MNLI acc: 83.7

# Background: Compressing Transformer Models

Why is **structured pruning** appealing:

- Flexible model structure with different sparsities
- Can achieve competitive results without unlabeled data
- Can be combined with task-specific distillation objectives

**Hard to achieve a large speedup e.g., 10 ×, without a significant performance drop.**

# Motivation

Why can't structured pruning approaches achieve a large speedup?

$$\text{MHA}(X) = \sum_{i=1}^{N_h} \mathbf{z}_{\text{head}}^{(i)} \text{Att}(W_Q^{(i)}, W_K^{(i)}, W_V^{(i)}, W_O^{(i)}, X)$$

Solution 1



4 heads are  
pruned in a  
layer



Solution 2



all heads are  
pruned in a  
layer



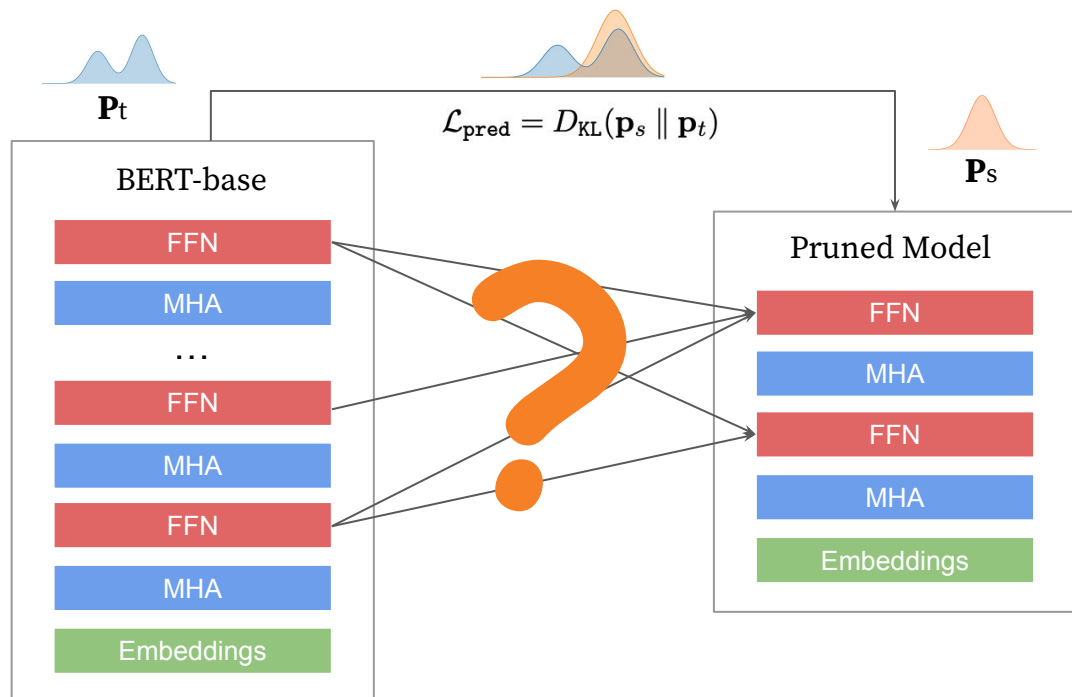
**Difficult to optimize in practice**

**Pruning layers leads to  
significant speedup gains!**



# Motivation

Layer-wise distillation could possibly improve pruning performance





## CoFi: **C**oarse- and **F**ine-grained Pruning

- Jointly prune coarse- and fine-grained units
- A layerwise distillation objective

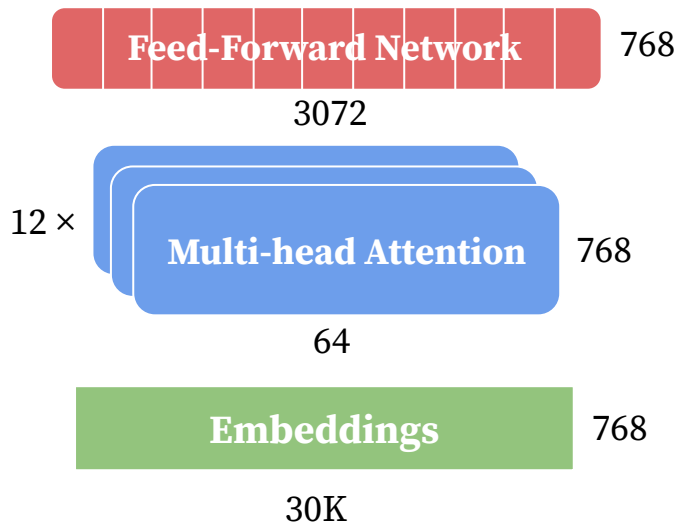
Achieves **10×** speedup

Preserves **90%** accuracy



# CoFi: **C**oarse- and **F**ine-grained Pruning

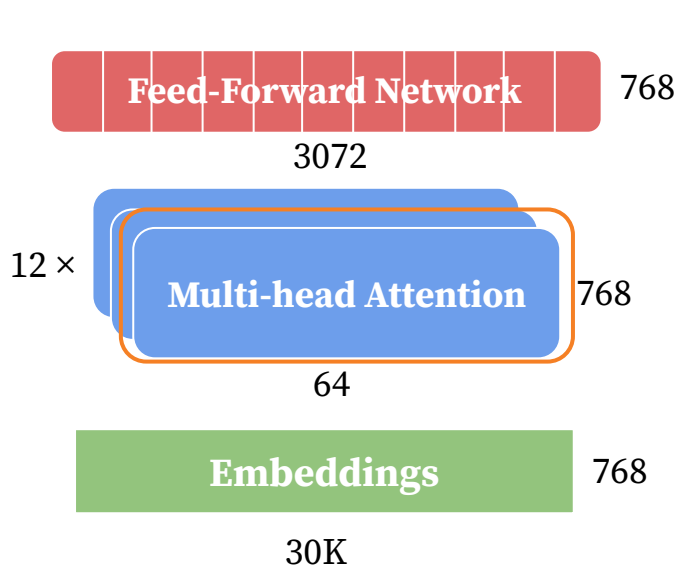
- Jointly prune coarse- and fine-grained units
- A layerwise distillation objective





# CoFi: **C**oarse- and **F**ine-grained Pruning

- Jointly prune coarse- and fine-grained units
- A layerwise distillation objective



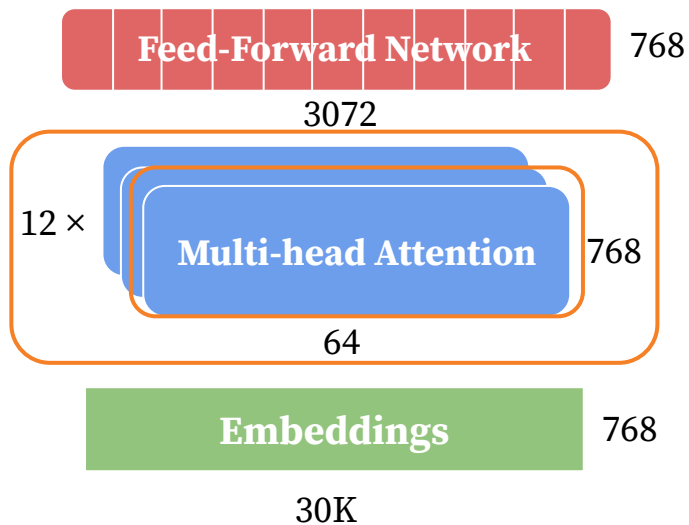
## Fine-grained units:

- Heads



# CoFi: **C**oarse- and **F**ine-grained Pruning

- Jointly prune coarse- and fine-grained units
- A layerwise distillation objective



## Fine-grained units:

- Heads

## Coarse units:

- MHA layers

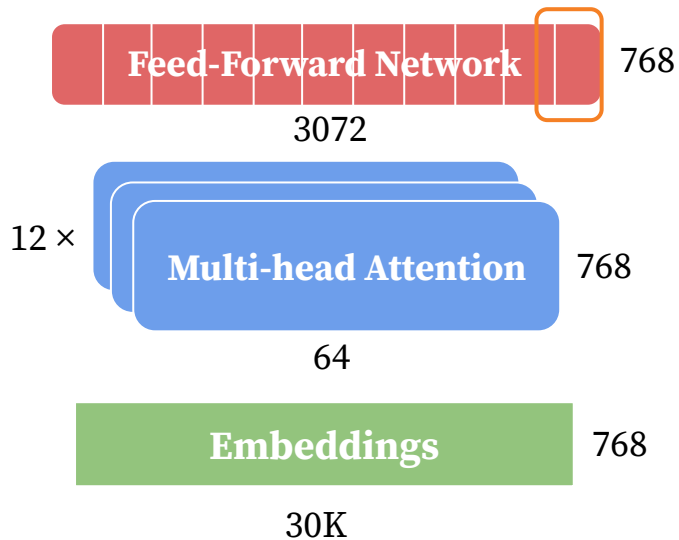
$$\text{MHA}(X) = z_{\text{MHA}} \cdot \sum_{i=1}^{N_h} (\mathbf{z}_{\text{head}}^{(i)} \cdot \text{Att}(W_Q^{(i)}, W_K^{(i)}, W_V^{(i)}, W_O^{(i)}, X))$$

$$z \in \{0, 1\}$$



# CoFi: **C**oarse- and **F**ine-grained Pruning

- Jointly prune coarse- and fine-grained units
- A layerwise distillation objective



## **Fine-grained units:**

- Heads
- Intermediate dimensions

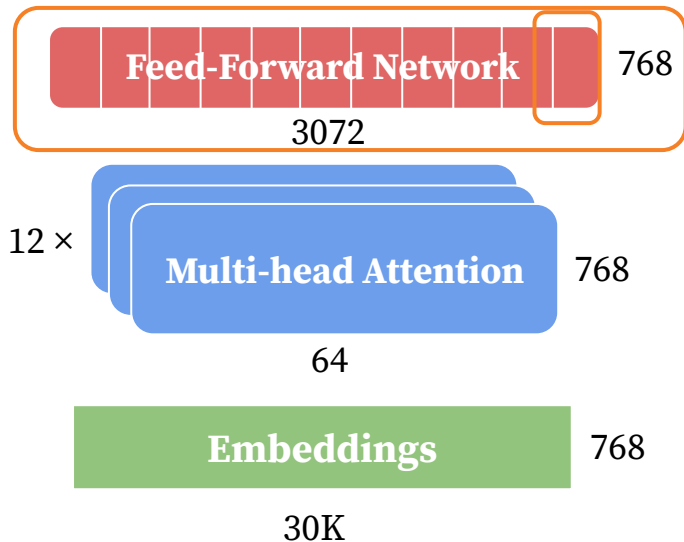
## **Coarse units:**

- MHA layers



# CoFi: **C**oarse- and **F**ine-grained Pruning

- Jointly prune coarse- and fine-grained units
- A layerwise distillation objective



## Fine-grained units:

- Heads
- Intermediate dimensions

## Coarse units:

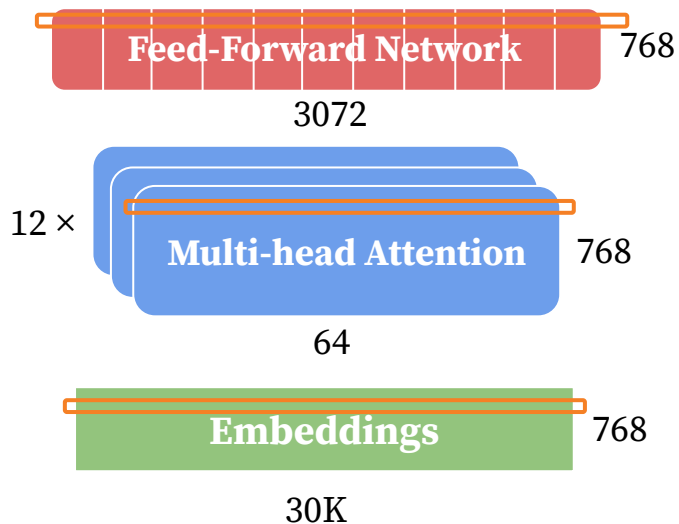
- MHA layers
- FFN layers

$$\text{FFN}(X) = \mathbf{z}_{\text{FFN}} \cdot \text{gelu}(XW_U) \cdot \text{diag}(\mathbf{z}_{\text{int}}) \cdot W_D$$



# CoFi: **C**oarse- and **F**ine-grained Pruning

- Jointly prune coarse- and fine-grained units
- A layerwise distillation objective



## **Fine-grained units:**

- Heads
- Intermediate dimensions
- Hidden dimension

## **Coarse units:**

- MHA layers
- FFN layers

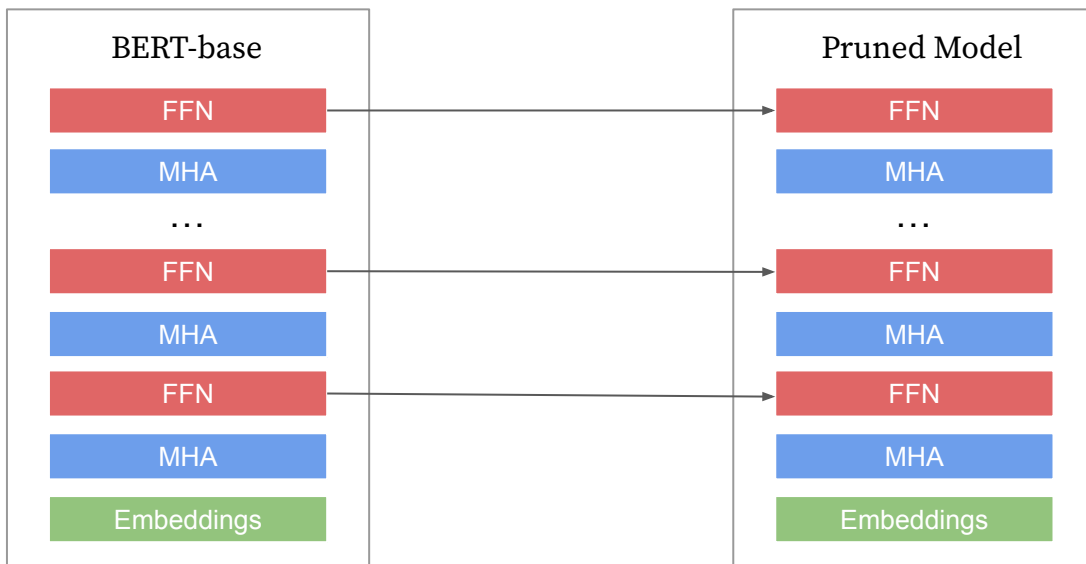




# CoFi: **C**oarse- and **F**ine-grained Pruning

- Jointly prune coarse- and fine-grained units
- A layerwise distillation objective

Naive Approach



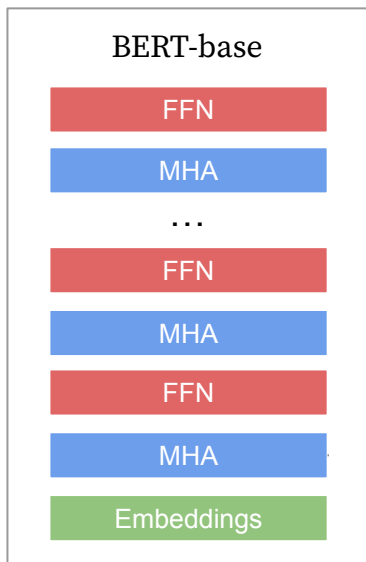
**Suboptimal when  
upper layers are  
pruned.**



# CoFi: **C**oarse- and **F**ine-grained Pruning

- Jointly prune coarse- and fine-grained units
- A layerwise distillation objective

Select 4 teacher layers  $\mathcal{T}$ , following Jiao et al. 2020

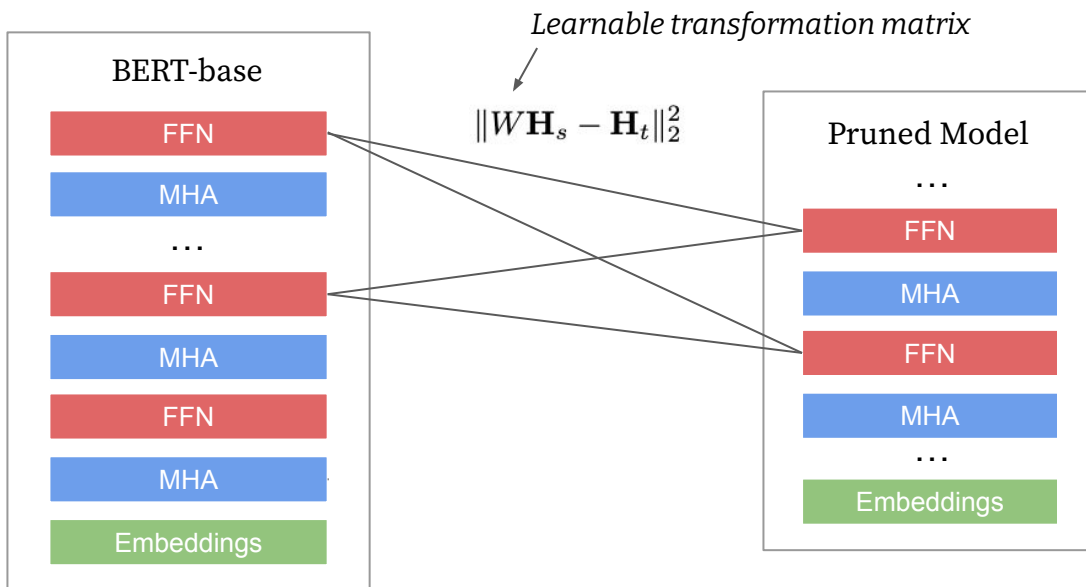




# CoFi: Coarse- and Fine-grained Pruning

- Jointly prune coarse- and fine-grained units
- A layerwise distillation objective

Calculate the **L2 distance** using the **training batch**

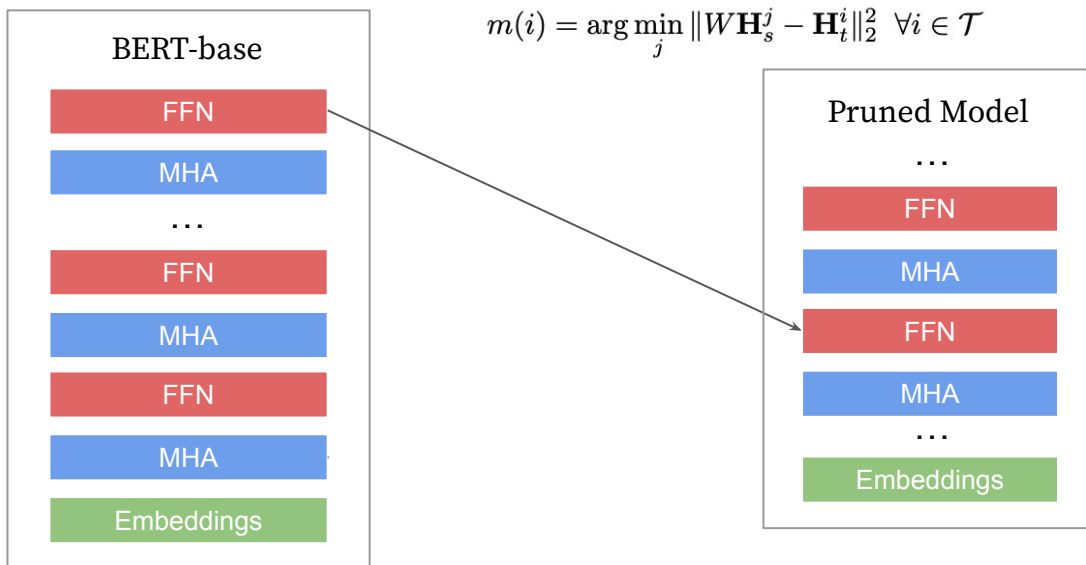




# CoFi: **C**oarse- and **F**ine-grained Pruning

- Jointly prune coarse- and fine-grained units
- A layerwise distillation objective

A **low** L2 distance  $\longrightarrow$  A **good** approximation

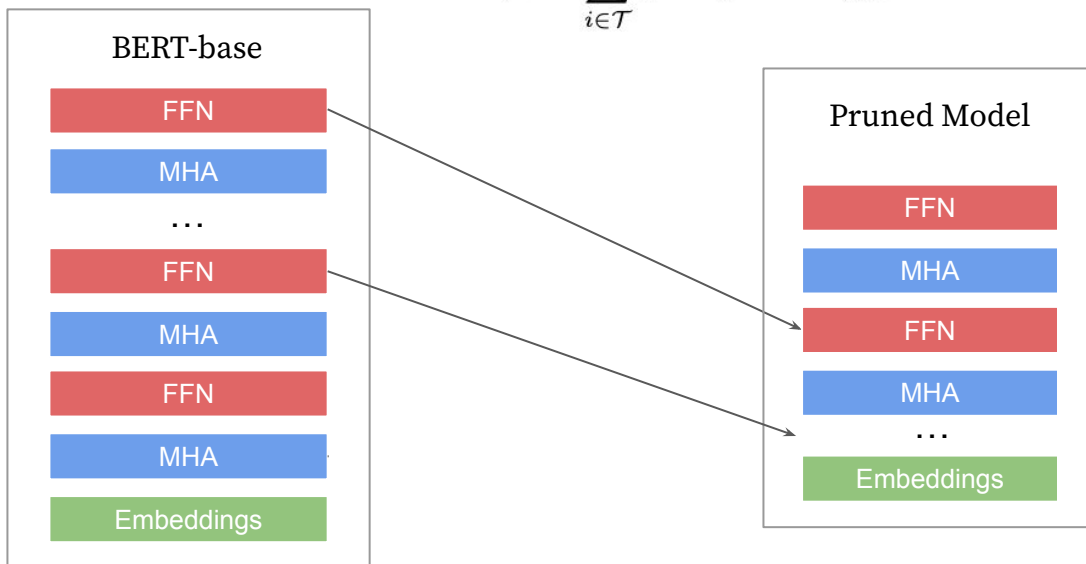




# CoFi: **C**oarse- and **F**ine-grained Pruning

- Jointly prune coarse- and fine-grained units
- A layerwise distillation objective

$$\mathcal{L}_{\text{layer}} = \sum_{i \in \mathcal{T}} \|W \mathbf{H}_s^{m(i)} - \mathbf{H}_t^i\|_2^2$$





# CoFi: **C**oarse- and **F**ine-grained Pruning

How to control the sparsity of the final model?

**Adapted the Lagrangian loss from Wang et al. 2020.**

- How to model  $z$ :
  - hard-concrete distribution (Loizos et al. 2018)

- Expected sparsity:  $f$ : function to calculate the sparsity M: full model Size  
$$\hat{s} = f(\mathbf{z}_{\text{int}}, \mathbf{z}_{\text{head}}, \mathbf{z}_{\text{hidden}}, z_{\text{FFN}}, z_{\text{MHA}}, M)$$

- Lagrangian loss:  $s^*$ : target sparsity

$$\mathcal{L}_{\text{lag}}(\lambda_1, \lambda_2, \hat{s}) = \lambda_1(s^* - \hat{s}) + \lambda_2(s^* - \hat{s})^2$$



# CoFi: **C**oarse- and **F**ine-grained Pruning

Final Objective

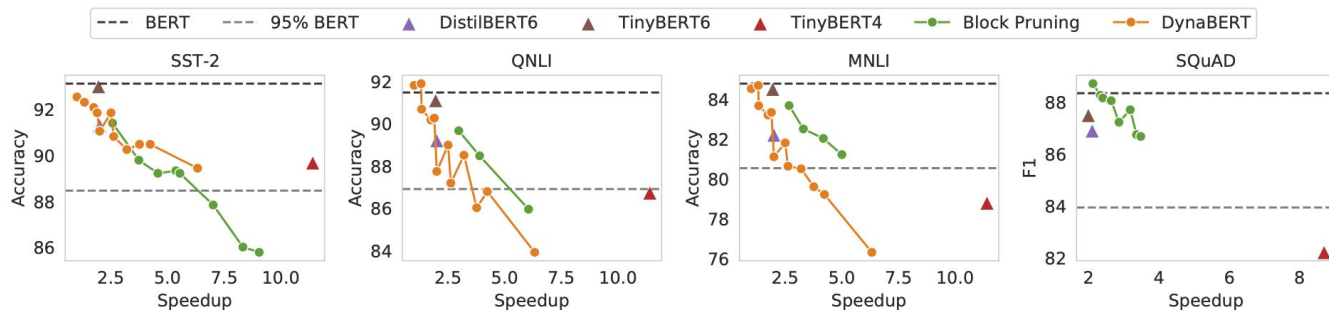
$$\max_{\lambda_1, \lambda_2} \min_{\theta, \hat{s}} \underbrace{\lambda \mathcal{L}_{\text{pred}}(\theta)}_{\text{Prediction layer distillation loss}} + \underbrace{(1 - \lambda) \mathcal{L}_{\text{layer}}(\theta)}_{\text{Layerwise distillation loss}} + \underbrace{\lambda_1 (s^* - \hat{s}) + \lambda_2 (s^* - \hat{s})^2}_{\substack{\text{Target sparsity} \quad \text{Expected sparsity}}}$$

# Experiment Results - GLUE and SQuAD 1.1

- GLUE: sentence classification tasks
- SQuAD 1.1: extractive question answering task



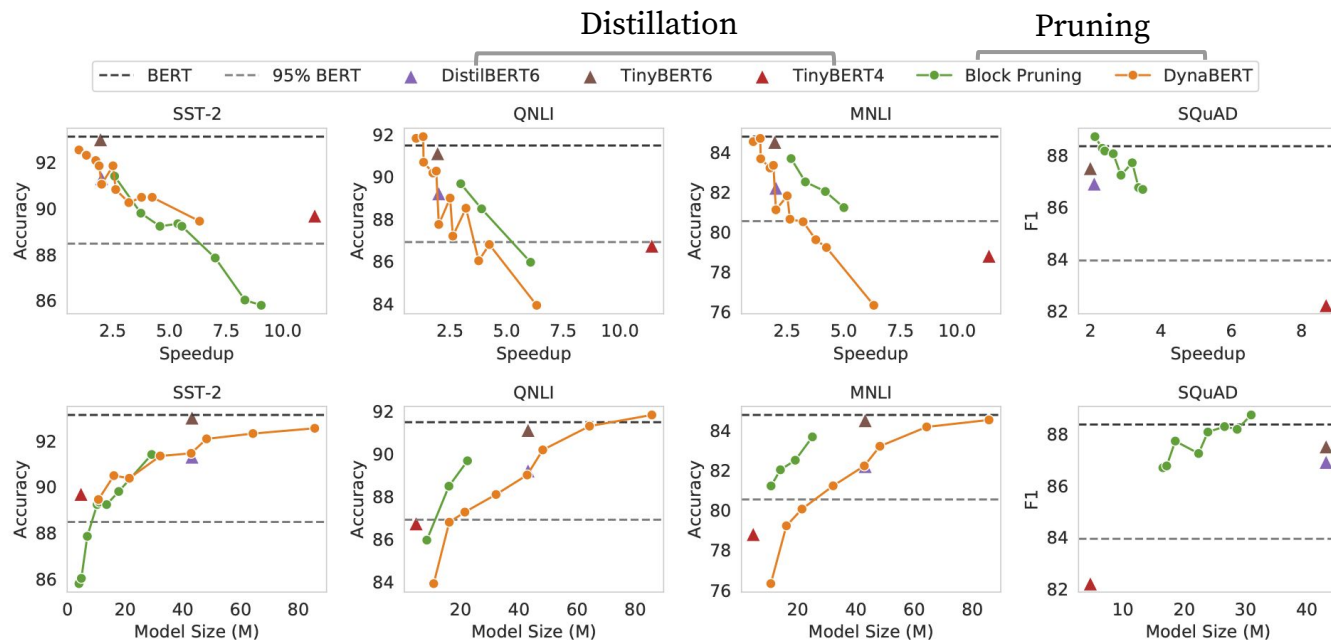
# Baseline Results - GLUE and SQuAD 1.1



Speedup v.s. Performance

*Please find more results on RoBERTa models and other compression approaches in the paper*

# Baseline Results - GLUE and SQuAD 1.1

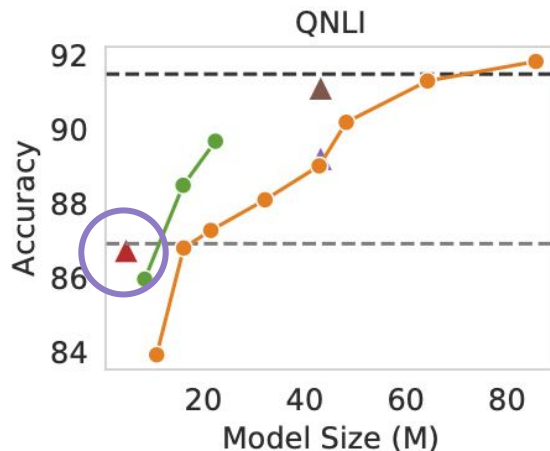
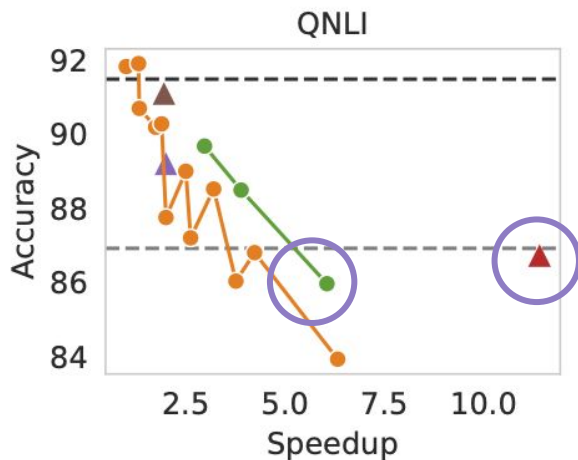


Speedup v.s. Performance and Model-size v.s. Performance

Please find more results on RoBERTa models and other compression approaches in the paper

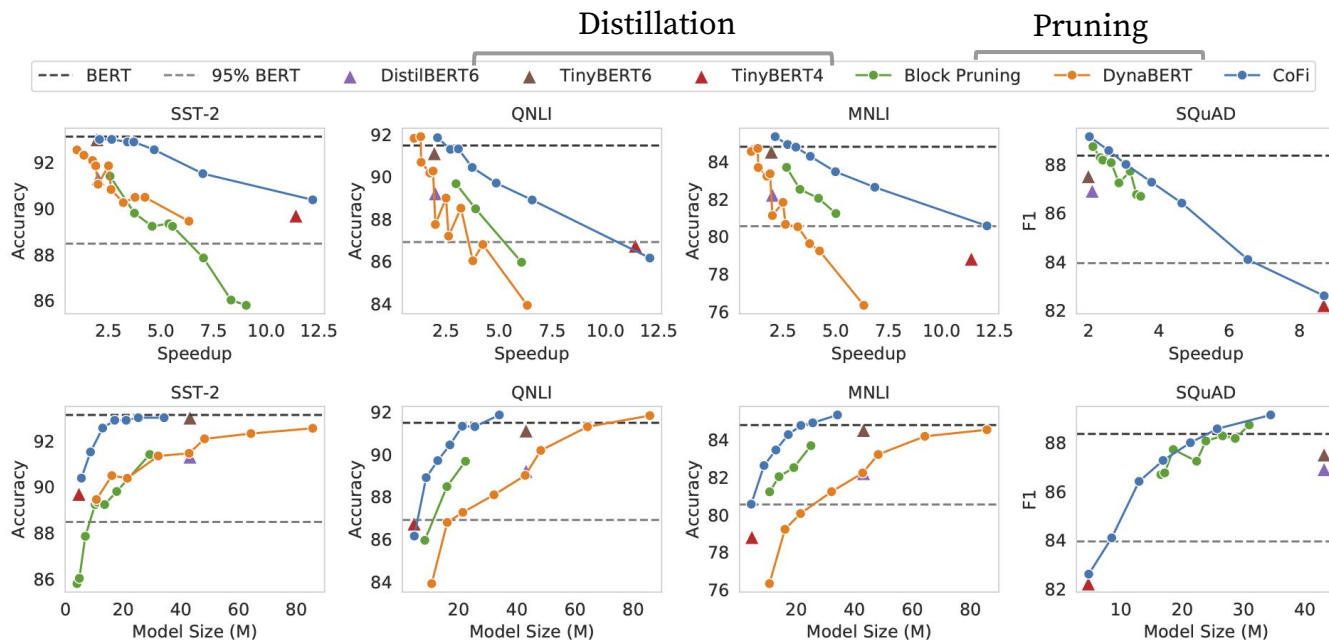
# Baseline Results - GLUE and SQuAD 1.1

▲ TinyBERT4    ● Block Pruning



Pruning falls behind distillation approaches on high sparsity levels.

# Experiment Results - GLUE and SQuAD 1.1



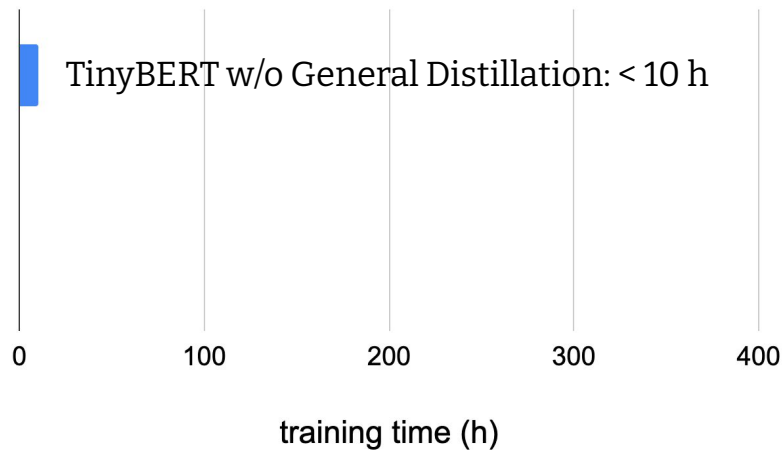
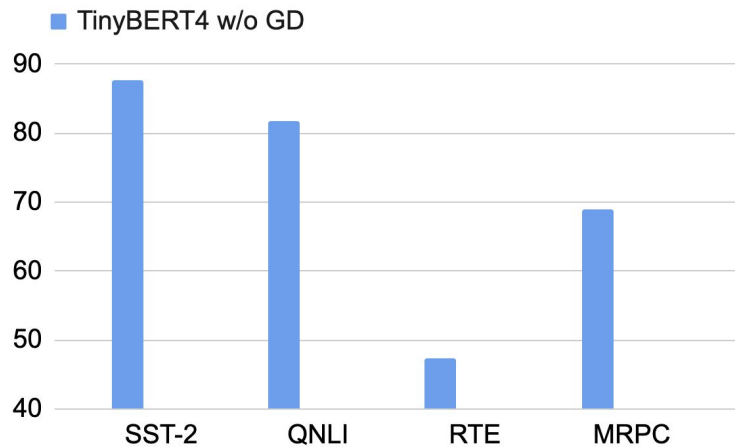
CoFi Pruning outperforms all distillation and pruning baselines **comparing under the same speedup and model size**

# Experiment Results

Models with  $10 \times$  speedups

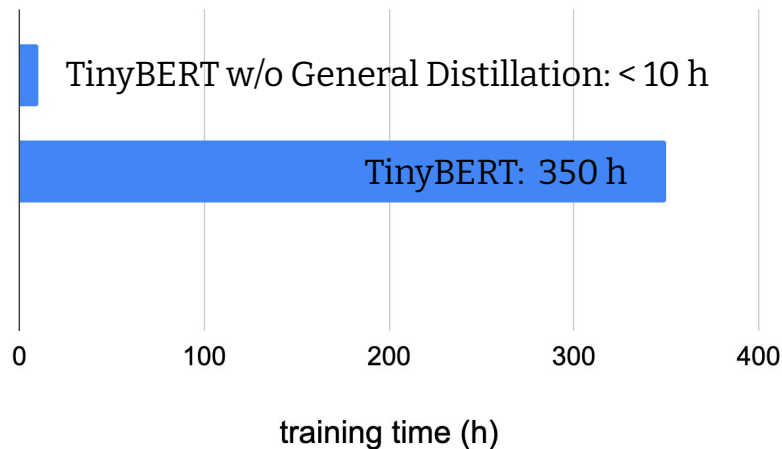
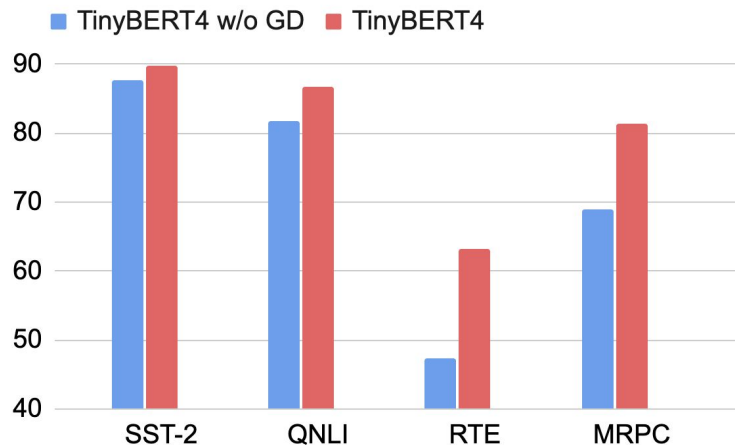
# Experiment Results

Models with  $10\times$  speedups



# Experiment Results

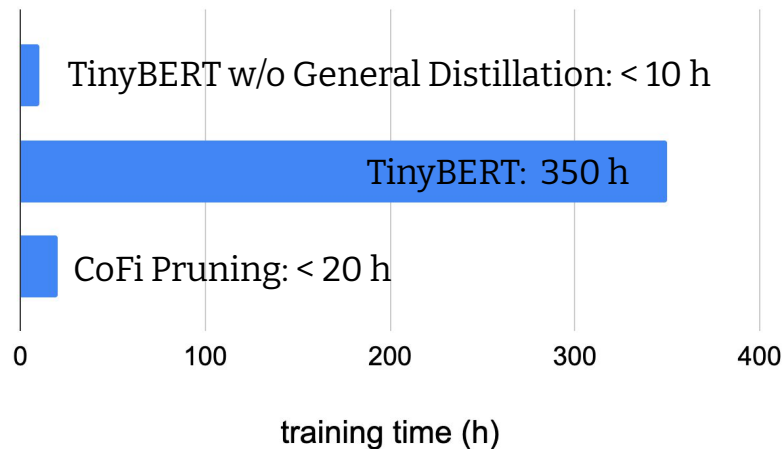
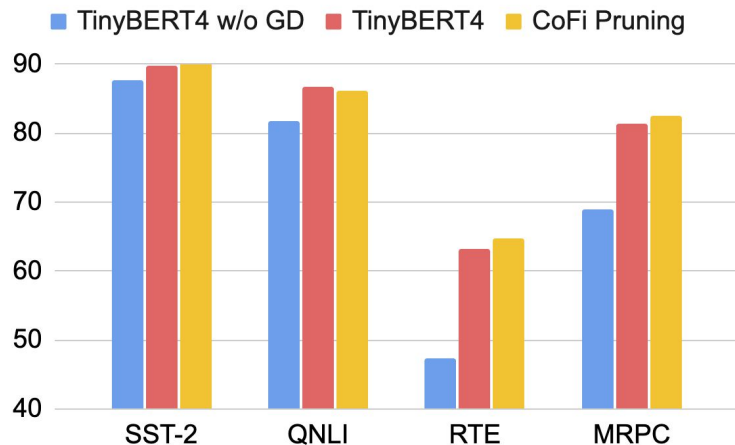
Models with  $10\times$  speedups



General distillation is essential but time consuming!

# Experiment Results

Models with  $10\times$  speedups



CoFi achieves **comparable or better performance** and speedup with **much less computation time**



## Ablation - Distillation Loss on 95% Models


	SST-2	QNLI	MNLI	SQuAD	Avg.
No distillation	86.6	84.2	78.2	75.8	81.2

## Ablation - Distillation Loss on 95% Models

	SST-2	QNLI	MNLI	SQuAD	Avg.
No distillation	86.6	84.2	78.2	75.8	81.2
+ $\mathcal{L}_{\text{pred}}$	<b>91.1</b>	85.1	79.7	82.5	84.6

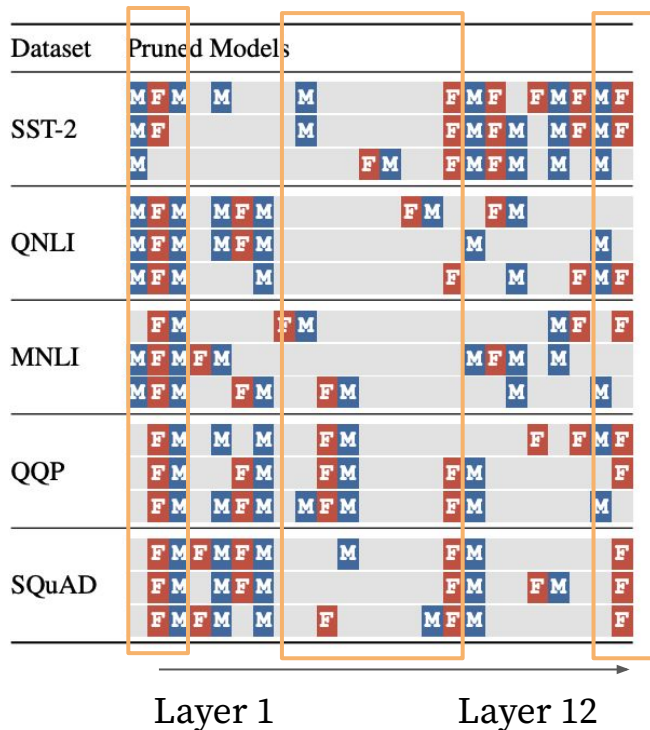
Prediction layer distillation **brings a large gain**

## Ablation - Distillation Loss on 95% Models

	SST-2	QNLI	MNLI	SQuAD	Avg.
No distillation	86.6	84.2	78.2	75.8	81.2
+ $\mathcal{L}_{\text{pred}}$	<b>91.1</b>	85.1	79.7	82.5	84.6
+ $\mathcal{L}_{\text{layer}}, \mathcal{L}_{\text{pred}}$ 	90.6	<b>86.1</b>	<b>80.6</b>	<b>82.6</b>	<b>85.0</b>

- Our proposed layer distillation loss brings **additional gains**
- The improvements on smaller sparsities are much larger

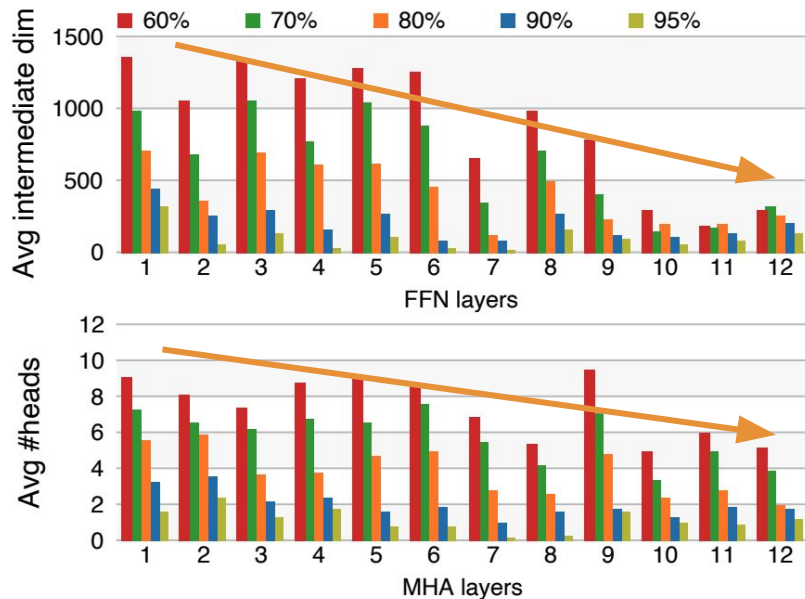
# Model Structures with 95% Sparsity



## Coarse-grained units:

- First and last FFN layers are largely retained
- Middle layers are more likely to be pruned

# Model Structures with 95% Sparsity



## Fine-grained Units:

Heads and intermediate dimensions from the top-layers are more likely to be pruned

# Summary

## CoFi Pruning



- Jointly prune coarse- and fine-grained units
- An additional layerwise distillation loss to guide pruning

## Compressed models

- Over  $10 \times$  speedups while maintaining 90% accuracy
- Closes the gap between structured pruning and knowledge distillation with much less computation

# Q & A

Codebase: <https://github.com/princeton-nlp/CoFiPruning>

Contact: [mengzhou@princeton.edu](mailto:mengzhou@princeton.edu)