# Parameter-Efficient Finetuning of Transformers for Source Code

**Shamil Ayupov**
HSE University
shiayupov@edu.hse.ru

**Nadezhda Chirkova**[*]
HSE University
nchirkova@hse.ru

## Abstract

Pretrained Transformers achieve state-of-the-art performance in various code-processing tasks but may be too large to be deployed. As software development tools often incorporate modules for various purposes which may potentially use a single instance of the pretrained model, it appears relevant to utilize parameter-efficient fine-tuning for the pretrained models of code. In this work, we test two widely used approaches, adapters and LoRA, which were initially tested on NLP tasks, on four code-processing tasks. We find that though the efficient fine-tuning approaches may achieve comparable or higher performance than the standard, full, fine-tuning in code understanding tasks, they underperform full fine-tuning in code-generative tasks. These results underline the importance of testing efficient fine-tuning approaches on other domains than NLP and motivate future research in efficient fine-tuning for source code.

## 1 Introduction

The pipeline of pretraining Transformers on a large corpora and fine-tuning the pretrained model on a task-specific (downstream) dataset has become a dominating paradigm in natural language processing (NLP). Inspired by the success of pretrained Transfomers in NLP, recent studies adopt this approach to other data domains, e.g., to source code, and demonstrate its state-of-the-art performance in various code-processing tasks. For example, [1, 2] show that the BERT model pretrained on source code corpora substantially outperforms earlier approaches in code summarization (generating textual descriptions for code snippets), code search (retrieving code fragments given textual query) and code clone detection (deciding if a pair of code fragments is functionally identical). The works of [3, 4, 5] further adapt pretraining approaches to source code by considering other Transformer-based models and proposing new pretraining objectives.

Despite achiving better results in downstream tasks, pretrained models are known to suffer from the large model size. The widely used approach to mitigating this issue in NLP is parameter-efficient fine-tuning (PE fine-tuning), including such methods as Adapter Training, LoRA or prefix tuning [6, 7, 8]. The general idea is to freeze the pretrained model and to introduce the small amount of additional parameters which will be fine-tuned for each task individually. Such approach substantially optimizes memory consumption at the deployment stage in case when one pretrained model is used in several applied tasks. It also reduces memory consumption at the fune-tuning stage because the gradients only need to be computed for additional task-specific parameters. However, the downstream performance of the parameter-efficient fine-tuning depends on the task: empirical studies show that PE fine-tuning approaches may outperform the standard (full) fine-tuning in the NLP tasks with small data but as the data size increases, full fine-tuning regains dominance [9, 10].

---

[*]Now at Naver Labs Europe

The large size of the pretrained Transformers especially matters in source code processing since it may prevent the pretrained models of code from deployment in integrated development systems (IDEs) and other developers tools. As an IDE includes modules for solving various code-processing tasks, it appears to be highly relevant to adopt PE fine-tuning approaches to code models. However, existing works mostly test PE fine-tuning on NLP tasks and, to the best of our knowledge, this approach was not previously tested on source code processing tasks, which are intuitively more complex. Moreover, source code data is often automatically collected from open-source repositories, e. g. by parsing functions and their descriptions from GitHub in code-to-text and text-to-code tasks. This results in a high level of noise in the training data (e.g. low-quality textual annotations, incorrect or buggy code) which makes it hard to apply the recommendations about applicability of the PE fine-tuning approaches in the small data setting.

In this work, we test two widely used PE fine-tuning approaches, Adapter Training and LoRA, and their combination for the pretrained models of source code. We find that PE fine-tuning cannot completely replace full fine-tuning. Particularly, in code-understanding downstream tasks, the performance of PE fine-tuning is sometimes comparable or better than of full fine-tuning, but in more complex code-generative tasks, PE fine-tuning substantially underperforms full fine-tuning. These results underline the importance of testing PE fine-tuning approaches on other domains than NLP and motivate future research on efficient fine-tuning for source code. The source code of our experiments is available at `https://github.com/ShamerD/source-code-efficient-ft`.

## 2   Methodology

We consider two widely used PE fine-tuning approaches, Adapter Training (AT) [6] and LoRA [7], and also consider their combination. We outline the core ideas of the methods below and refer readers to the original papers for details. We vary the number of learnable parameters in these approaches and compare them with the full fine-tuning.

**Adapter Training** (AT) relies on adding small fully-connected layers inside the Transformer. The layers are added after the Multi-Head Attention and Feed Forward layers. During fine-tuning, only these layers are trained while original model weights are frozen. The drawback of this method is introducing inference latency as these layers themselves have computational costs.

**LoRA** introduces a learnable low-rank weights update $\Delta W = AB$ that will be summed up with the weights of the original model after fine-tuning: $W + \Delta W = W + AB$, where $W \in \mathbb{R}^{d \times k}, A \in \mathbb{R}^{d \times r}, B \in \mathbb{R}^{r \times k}, r < min(d, k)$. As a result, there is no additional inference latency. The described modification can be applied to different parts of the model: to Query and Value weights in Multi-Head Attention layer (as in original paper) or to Feed Forward layer [9] (referred to as FF-LoRA in this work).

**Models, tasks and fine-tuning data.**   We use CodeT5 [3] and PLBART [4] as the main models in our experiments since these are high-performing encoder-decoder models applicable to both generating and understanding tasks. We use the CodeT5 implementation for our experiments[2], employ HuggingFace's PyTorch Transformers library [11] and use published pre-trained `CodeT5-base` (223M) and `PLBART-base` (140M) checkpoints in our experiments.

We consider four tasks from the CodeT5 and PLBART papers covering both code-understanding and code-generative scenarios. The considered datasets and metrics are widely used in source code processing and are a part of the CodeXGLEU benchmark [12].

**Code summarization** is the task to summarize code snippets into human-readable English descriptions. We consider Python and Go languages from the CodeSearchNet dataset [13]. The dataset consists of 251,820 / 13,914 / 14,918 examples in Python and 167,288 / 7,325 / 8,122 examples in Go in the train / development / test sets respectively. This task is evaluated using smoothed BLEU-4 [14].

**Code generation** is the task to generate code snippets given natural language (English) descriptions. The data consists of 100,000 / 2,000 / 2,000 (text + environment context, code) pairs on Java from the Concode dataset [15]. This task is evaluated using CodeBLEU [16].

---

[2]`https://github.com/salesforce/CodeT5` (BSD-3-Clause license)

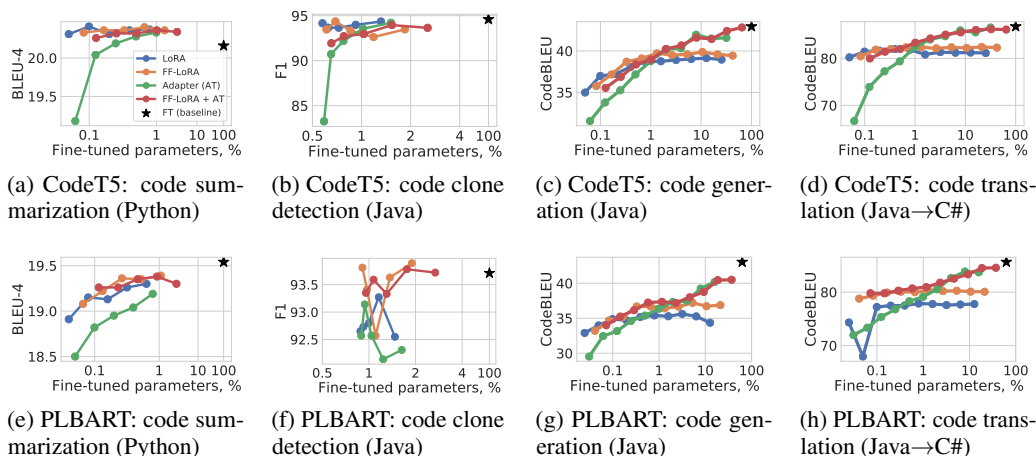| (a) CodeT5: code sum-marization (Python) | (b) CodeT5: code clone detection (Java) | (c) CodeT5: code gener-ation (Java) | (d) CodeT5: code trans-lation (Java→C#) |
| (e) PLBART: code sum-marization (Python) | (f) PLBART: code clone detection (Java) | (g) PLBART: code gen-eration (Java) | (h) PLBART: code trans-lation (Java→C#) |

Figure 1: Quality of the PE fine-tuning methods in differents tasks on the test set vs. the number of trainable parameters. The numerical results are given in Appendix.

**Code translation** aims to translate code snippets from one programming language to another. The dataset consists of 10,300 / 500 / 1,000 code snippet pairs on C# and Java languages and translation can be organized in both directions. This task is again evaluated using CodeBLEU. We additionally experiment with low resource setup (2.5k training examples) in Appendix A.

**Code clone detection** is the task of predicting whether two code snippets have the same function-ality. The dataset is provided by [17] and consists of 9,134 Java code snippets with the total of 901,028 / 415,416 / 415,416 labelled clone pairs. This task is evaluated using F1 score.

**Experimental setup.** Fine-tuning hyperparameters were selected as in the original papers on CodeT5 [3] and PLBART [4]. The most important hyperparameter in the PE fine-tuning methods is the amount of learnable parameters which is determined by the inner layer size in AT and update rank in LoRA. We vary the latter hyperparameters in log-space (1, 2, 4, 8, 16). We report the quality on the test set, selecting the early stopping iteration in all experiments based on the quality on the development set.

## 3 Experiments

The main results for four tasks are shown in Fig. 1, the results for additional programming languages are presented in Appendix. We find that in code-generative tasks (code generation and code translation, Fig. 1c, 1d, 1g, 1h) the PE fine-tuning methods significantly underperform full fine-tuning and may achieve comparable performance only when the number of fine-tuned parameters approaches the number of all parameters in the model. In code understanding tasks (code summarization and code clone detection, Fig. 1a, 1b, 1e, 1f) the PE fine-tuning methods are sometimes able to achieve comparable or even better quality compared to full fine-tuning. Particularly, in clone detection, the PE fine-tuning of both models performs comparable to full finetuning, and in code summarization, this holds only for the CodeT5 model. However, the later result may be potentially attributed to the observation that the pretraining of publicly available CodeT5 checkpoint included the parallel part of the CodeSearchNet dataset, which is used for code summarization finetuning. To sum up, *PE fine-tuning often substantially underperforms full fine-tuning in code-generative tasks and may achieve comparable or better results only in simpler code-understanding tasks.*

We also analyse the training curves in Fig. 2 for both cases when PE fine-tuning underperforms full-finetuning and when they perform on par. This reveals that in the former case, the full fine-tuning severely overfits to the fine-tuning data (validation performance starts to decrease at some fine-tuning iteration) while PE fine-tuning acts like a regularizer. In the latter case, validation performance continues improving for both fine-tuning approaches.
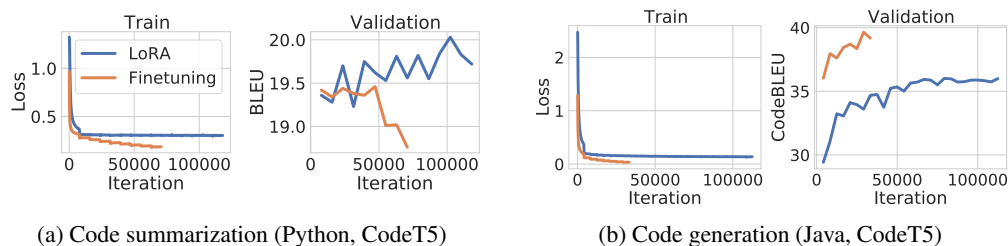
(a) Code summarization (Python, CodeT5)      (b) Code generation (Java, CodeT5)

Figure 2: Example learning curves for two tasks.

| Method | Example |
|---------|---------|
| GT | Scans through a string for substrings matched some patterns . |
| FT | matchall - Searches text for all patterns in patterns and returns a list of matches |
| LoRA | Return a list of all matching patterns in text . |
| FF-LoRA | Match all patterns in text . |
| AT | Match all patterns in text . |
| GT | int - > None Download a WHOLE course . Reuse the API call to save time . |
| FT | wanmen_download_by_course_topic and wanmen_download_by_course_part recursively |
| LoRA | Download wanmen by course . |
| FF-LoRA | Download wanmen files by course . |
| AT | Download wanmen by course . |

Table 1: Examples of CodeT5 predictions when using different PE methods on code summarization task. GT stands for Ground Truth.

Comparing different PE fine-tuning approaches, we find that Adapter Training is characterized by the monotonic dependence of the performance w. r. t. the number of fine-tuned parameters and underperforms LoRA / FF-LoRA for small numbers of fine-tuned parameters. However, in some tasks, Adapter Training outperforms LoRA / FF-LoRA for larger values of this hyperparameter. The combination of two approaches often behaves as an upper envelope of the corresponding quality plots, except the cases when Adapter Training performance is very low with the small number of fine-tuned parameters and reduces the combination's performance as well. FF-LoRA slightly outperforms LoRA or performs on par with it.

We also perform a qualitative analysis of the predictions on the code summarization task. We inspected around 100 different models' predictions. Some can be seen in Table 1, more can be seen in Appendix A. The key observation is that models fine-tuned with the efficient fine-tuning methods tend to have shorter predictions based on a simple template. At the same time the fully fine-tuned model outputs more elaborate predictions but they can contain irrelevant information or code fragments directly copied from the input. This may also indicate that the full fine-tuning results in the model being overfit.

## 4 Related Work

Continuing the pioneer work of [6, 18] on Adapter Training, recent research proposed a variety of PE fine-tuning approaches. Prompt-Tuning introduces prefix vectors (prompts) before the original input sequence [8] and updates only prompts during fine-tuning. Diff-Pruning [19] learns sparse updates. LoRA [7] performs low-rank reparametrization of the update matrices. BitFit [20] fine-tunes only bias vectors in Transformers. Authors of [10] conduct an empirical study of various methods for widely used NLP benchmarks and show the absence of the universal PE approach that would perform on par or better than full fine-tuning in all tasks, especially in the medium- and high-resource regimes. Our work extends the set of benchmarks on which PE fine-tuning was tested with source code processing tasks and comes to the similar conclusion.

# 5 Conclusion

In this work we studied the applicability of the efficient fine-tuning methods to Transformers for source codes in various source code processing tasks. Our main conclusion is that efficient fine-tuning often substantially underperforms full fine-tuning in source code processing tasks. This underlines that efficient fine-tuning approaches should be tested on other data domains except natural language and motivates the future development of efficient fine-tuning approaches for source code. As for practical recommendations, our experiments suggest using FF-LoRA or LoRA when the affordable number of learnable parameters is small and the combination of FF-LoRA and Adapter Training for larger parameter counts.

## Acknowledgments

## References

[1] Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang, *et al.*, "Codebert: A pre-trained model for programming and natural languages," *arXiv preprint arXiv:2002.08155*, 2020.

[2] D. Guo, S. Ren, S. Lu, Z. Feng, D. Tang, L. Shujie, L. Zhou, N. Duan, A. Svyatkovskiy, S. Fu, *et al.*, "Graphcodebert: Pre-training code representations with data flow," in *International Conference on Learning Representations*, 2020.

[3] Y. Wang, W. Wang, S. Joty, and S. C. Hoi, "Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation," *arXiv preprint arXiv:2109.00859*, 2021.

[4] W. U. Ahmad, S. Chakraborty, B. Ray, and K.-W. Chang, "Unified pre-training for program understanding and generation," *arXiv preprint arXiv:2103.06333*, 2021.

[5] B. Roziere, M.-A. Lachaux, M. Szafraniec, and G. Lample, "Dobf: A deobfuscation pre-training objective for programming languages," *arXiv preprint arXiv:2102.07492*, 2021.

[6] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, and S. Gelly, "Parameter-efficient transfer learning for nlp," in *International Conference on Machine Learning*, pp. 2790–2799, PMLR, 2019.

[7] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," *arXiv preprint arXiv:2106.09685*, 2021.

[8] X. L. Li and P. Liang, "Prefix-tuning: Optimizing continuous prompts for generation," *arXiv preprint arXiv:2101.00190*, 2021.

[9] J. He, C. Zhou, X. Ma, T. Berg-Kirkpatrick, and G. Neubig, "Towards a unified view of parameter-efficient transfer learning," *arXiv preprint arXiv:2110.04366*, 2021.

[10] G. Chen, F. Liu, Z. Meng, and S. Liang, "Revisiting parameter-efficient tuning: Are we really there yet?," *arXiv preprint arXiv:2202.07962*, 2022.

[11] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, *et al.*, "Transformers: State-of-the-art natural language processing," in *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pp. 38–45, 2020.

[12] S. Lu, D. Guo, S. Ren, J. Huang, A. Svyatkovskiy, A. Blanco, C. Clement, D. Drain, D. Jiang, D. Tang, *et al.*, "Codexglue: A machine learning benchmark dataset for code understanding and generation," *arXiv preprint arXiv:2102.04664*, 2021.

[13] H. Husain, H.-H. Wu, T. Gazit, M. Allamanis, and M. Brockschmidt, "Codesearchnet challenge: Evaluating the state of semantic code search," *arXiv preprint arXiv:1909.09436*, 2019.

[14] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.

[15] S. Iyer, I. Konstas, A. Cheung, and L. Zettlemoyer, "Mapping language to code in programmatic context," *arXiv preprint arXiv:1808.09588*, 2018.

[16] S. Ren, D. Guo, S. Lu, L. Zhou, S. Liu, D. Tang, N. Sundaresan, M. Zhou, A. Blanco, and S. Ma, "Codebleu: a method for automatic evaluation of code synthesis," *arXiv preprint arXiv:2009.10297*, 2020.

[17] W. Wang, G. Li, B. Ma, X. Xia, and Z. Jin, "Detecting code clones with graph neural network and flow-augmented abstract syntax tree," in *2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pp. 261–271, IEEE, 2020.

[18] J. Pfeiffer, A. Kamath, A. Rücklé, K. Cho, and I. Gurevych, "Adapterfusion: Non-destructive task composition for transfer learning," *arXiv preprint arXiv:2005.00247*, 2020.

[19] D. Guo, A. M. Rush, and Y. Kim, "Parameter-efficient transfer learning with diff pruning," in *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)*, pp. 4884–4896, 2021.

[20] E. B. Zaken, S. Ravfogel, and Y. Goldberg, "Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models," *arXiv preprint arXiv:2106.10199*, 2021.

[21] P. S. Kostenetskiy, R. A. Chulkevich, and V. I. Kozyrev, "HPC resources of the higher school of economics," *Journal of Physics: Conference Series*, vol. 1740, p. 012050, 2021.

# A  Appendix

**Limitations and potential negative societal impacts**   To achieve a broader view on the considered problem, given limited resources, we selected a variety of applied tasks and programming languages, focused on two widely used PE fine-tuning approaches and considered two pretrained models. At the same time, given a wide range of PE fine-tuning approaches, pretrained models and applied tasks in the literature, it could be interesting to experiment with other models, tasks and approaches as well.

While the PE methods bring the main benefit to large models (with billions of parameters), they are not considered in this work due to limited computational resources. Experimenting with larger models could be an interesting direction for future resarch.

Though we are not aware of potential negative societal impacts of our work, we note that the work may have potential negative environmental impact (about 300 Tesla V100 or A100 runs, each taking 8–24 hours, on the internal cluster).

**Additional results**   Figure 3 shows results on additional tasks: code summarization on the Go language and code translation in other direction (C#→Java). The PE methods behave similarly as in the main text. Figure 4 shows results of the FF-LoRA + Adapter Training method in low-resource setting (2.5k training examples vs. 10.3k training examples in full training set) of the Java → C# translate task. The method performs similarly relative to the full fine-tuning as when using all available data. Studying the effect of the methods in low-resource scenario (and potentially in few-shot setting) is one of possible research directions.

Tables 2-8 show quantitative results from Figures 1 and 3 in the numerical form. Hyperparameter $r$ (methods' internal size) for each task/method set was chosen based on the quality on the dev set. Table 9 shows more CodeT5 predictions in the code summarization (Python) task.



(a) CodeT5: code summarization (Go)

(b) CodeT5: code summarization (Java)

(c) CodeT5: code translation (C#→Java)

(d) PLBART: code summarization (Go)

(e) PLBART: code summarization (Java)

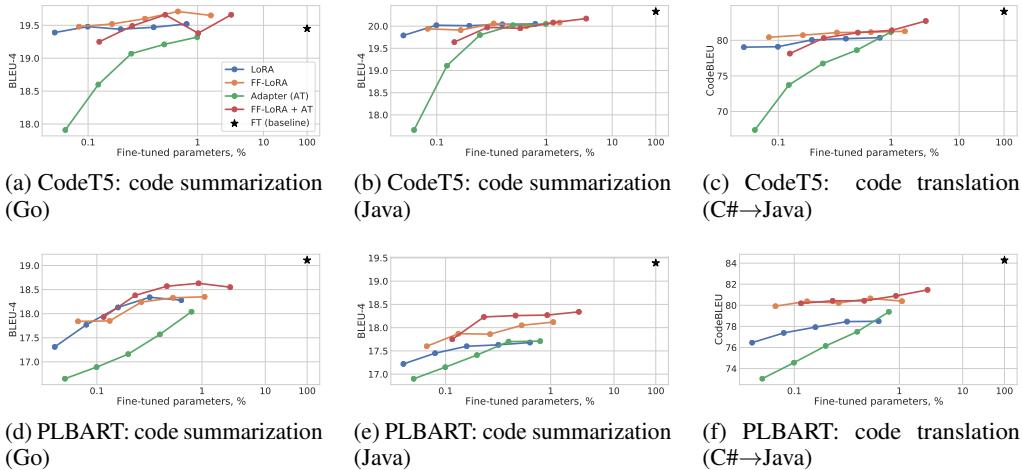(f) PLBART: code translation (C#→Java)

Figure 3: Quality of the PE fine-tuning methods in additional tasks on the test-set vs. the number of trainable parameters
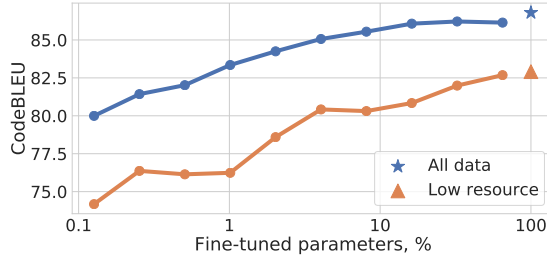
Figure 4: Quality of the Full fine-tuning (star/triangle) and FF-LoRA + Adapter (solid line) Training method when training on all data (10.3k training examples) and in low-resource setting (2.5k training examples) of the Java → C# translate task vs. the number of trainable parameters

| Model | Method | # Parameters | BLEU-4 ($\uparrow$) |
|---|---|---|---|
| CodeT5 | FT (baseline) | 223M (100%) | $20.16 \pm 0.22$ |
| | LoRA ($r$=2) | 0.2M (0.1%) | $\mathbf{20.36} \pm 0.05$ |
| | FF-LoRA ($r$=4) | 0.7M (0.3%) | $\mathbf{20.36} \pm 0.03$ |
| | AT ($r$=16) | 2.2M (1%) | $\mathbf{20.36} \pm 0.03$ |
| | FF-LoRA + AT ($r = 16$) | 5.2M (2.3%) | $\mathbf{20.36} \pm 0.03$ |
| PLBART | FT (baseline) | 140M (100%) | $\mathbf{19.63} \pm 0.07$ |
| | LoRA ($r$=16) | 0.9M (0.6%) | $19.23 \pm 0.09$ |
| | FF-LoRA ($r$=16) | 1.5M (1%) | $19.30 \pm 0.12$ |
| | AT ($r$=16) | 1.1M (0.8%) | $19.26 \pm 0.05$ |
| | FF-LoRA + AT ($r = 8$) | 1.3M (0.9%) | $19.39 \pm 0.03$ |

Table 2: PE methods quality on the test set in the **code summarization (Python)** task. Hyperparameter $r$ (methods' internal size) was chosen based on the quality on the dev set. Standard deviation is estimated in 3 runs.

| Model | Method | # Parameters | F1 ($\uparrow$) | Precision ($\uparrow$) | Recall ($\uparrow$) |
|---|---|---|---|---|---|
| CodeT5 | FT (baseline) | 224M (100%) | $\mathbf{94.56} \pm 0.28$ | $\mathbf{95.74} \pm 0.51$ | $93.41 \pm 0.15$ |
| | LoRA ($r$=4) | 1.6M (0.7%) | $93.96 \pm 0.46$ | $93.26 \pm 0.64$ | $94.67 \pm 0.40$ |
| | FF-LoRA ($r$=16) | 4.1M (1.8%) | $93.66 \pm 0.22$ | $92.16 \pm 0.36$ | $\mathbf{95.21} \pm 0.81$ |
| | AT ($r$=16) | 3.4M (1.5%) | $93.76 \pm 0.54$ | $92.56 \pm 0.76$ | $94.99 \pm 0.66$ |
| | FF-LoRA + AT ($r$=16) | 6.3M (2.8%) | $93.91 \pm 0.28$ | $93.83 \pm 0.61$ | $93.99 \pm 0.48$ |
| PLBART | FT (baseline) | 141M (100%) | $\mathbf{94.03} \pm 0.33$ | $\mathbf{94.00} \pm 0.62$ | $94.06 \pm 0.54$ |
| | LoRA ($r$=1) | 1.2M (0.9%) | $92.89 \pm 0.24$ | $91.40 \pm 0.24$ | $94.43 \pm 0.33$ |
| | FF-LoRA ($r$=16) | 2.7M (1.9%) | $93.75 \pm 0.24$ | $92.91 \pm 0.90$ | $94.63 \pm 0.56$ |
| | AT ($r$=8) | 1.7M (1.2%) | $92.49 \pm 0.29$ | $90.77 \pm 0.31$ | $94.28 \pm 0.43$ |
| | FF-LoRA + AT ($r$=2) | 1.5M (1.1%) | $93.62 \pm 0.03$ | $92.32 \pm 0.24$ | $\mathbf{94.96} \pm 0.30$ |

Table 3: PE methods quality on the test set in the **code clone detection** task. Hyperparameter $r$ (methods' internal size) was chosen based on the quality on the dev set. Standard deviation is estimated in 3 runs.

8

| Model | Method | # Parameters | BLEU-4 (↑) | EM (↑) | CodeBLEU(↑) |
|---|---|---|---|---|---|
| CodeT5 | FT (baseline) | 223M (100%) | **39.79** $\pm$ 0.99 | **22.33** $\pm$ 0.34 | **42.95** $\pm$ 0.91 |
| | LoRA ($r$=512) | 57M (25.4%) | 34.80 $\pm$ 0.12 | 21.52 $\pm$ 0.28 | 39.01 $\pm$ 0.11 |
| | FF-LoRA ($r$=256) | 47M (21.2%) | 36.17 $\pm$ 0.10 | 21.65 $\pm$ 0.15 | 39.72 $\pm$ 0.13 |
| | AT ($r$=512) | 71M (31.8%) | 38.25 $\pm$ 0.42 | 21.92 $\pm$ 0.37 | 41.78 $\pm$ 0.16 |
| | FF-LoRA + AT ($r$=512) | 165M (74%) | 39.45 $\pm$ 0.58 | 22.08 $\pm$ 0.19 | 42.71 $\pm$ 0.26 |
| PLBART | FT (baseline) | 140M (100%) | **39.40** $\pm$ 1.05 | **21.08** $\pm$ 0.21 | **42.64** $\pm$ 0.63 |
| | LoRA ($r$=8) | 0.4M (0.3%) | 26.45 $\pm$ 0.30 | 19.20 $\pm$ 0.04 | 34.91 $\pm$ 0.15 |
| | FF-LoRA ($r$=128) | 12M (8.5%) | 29.28 $\pm$ 0.37 | 20.13 $\pm$ 0.09 | 36.87 $\pm$ 0.26 |
| | AT ($r$=512) | 35M (25.4%) | 32.57 $\pm$ 0.40 | 20.47 $\pm$ 0.15 | 39.86 $\pm$ 0.31 |
| | FF-LoRA + AT ($r$=512) | 83M (59.3%) | 33.80 $\pm$ 0.47 | 20.68 $\pm$ 0.06 | 39.79 $\pm$ 0.50 |

Table 4: PE methods quality on the test set in the **code generation** task. Hyperparameter $r$ (methods' internal size) was chosen based on the quality on the dev set. Standard deviation is estimated in 3 runs.

| Model | Method | # Parameters | BLEU-4 (↑) | EM (↑) | CodeBLEU(↑) |
|---|---|---|---|---|---|
| CodeT5 | FT (baseline) | 223M (100%) | **83.42** $\pm$ 0.42 | **64.63** $\pm$ 1.10 | **86.80** $\pm$ 0.32 |
| | LoRA ($r$=64) | 7.1M (3.2%) | 76.30 $\pm$ 0.22 | 54.43 $\pm$ 0.29 | 81.22 $\pm$ 0.19 |
| | FF-LoRA ($r$=64) | 11.8M (5.3%) | 78.23 $\pm$ 0.09 | 55.73 $\pm$ 0.26 | 82.38 $\pm$ 0.13 |
| | AT ($r$=512) | 71M (31.8%) | 82.92 $\pm$ 0.07 | 63.73 $\pm$ 0.62 | 86.56 $\pm$ 0.04 |
| | FF-LoRA + AT ($r$=512) | 165M (74%) | 82.39 $\pm$ 0.45 | 63.03 $\pm$ 0.49 | 86.09 $\pm$ 0.35 |
| PLBART | FT (baseline) | 140M (100%) | **81.85** $\pm$ 0.33 | **59.50** $\pm$ 0.99 | **85.50** $\pm$ 0.20 |
| | LoRA ($r$=512) | 28M (20.3%) | 71.86 $\pm$ 0.02 | 47.67 $\pm$ 0.09 | 77.75 $\pm$ 0.07 |
| | FF-LoRA ($r$=32) | 3M (2.1%) | 75.16 $\pm$ 0.12 | 51.23 $\pm$ 0.12 | 80.18 $\pm$ 0.08 |
| | AT ($r$=512) | 35M (25.4%) | 80.22 $\pm$ 0.40 | 56.10 $\pm$ 0.78 | 84.19 $\pm$ 0.37 |
| | FF-LoRA + AT ($r$=512) | 83M (59.3%) | 80.42 $\pm$ 0.12 | 56.90 $\pm$ 0.49 | 84.36 $\pm$ 0.16 |

Table 5: PE methods quality on the test set in the **code translation (Java → C#)** task. Hyperparameter $r$ (methods' internal size) was chosen based on the quality on the dev set. Standard deviation is estimated in 3 runs.

| Model | Method | # Parameters | BLEU-4 (↑) |
|---|---|---|---|
| CodeT5 | FT (baseline) | 223M (100%) | **19.68** $\pm$ 0.17 |
| | LoRA ($r$=8) | 0.9M (0.4%) | 19.47 $\pm$ 0.00 |
| | FF-LoRA ($r$=8) | 1.5M (0.7%) | 19.41 $\pm$ 0.29 |
| | AT ($r$=8) | 1.1M (0.5%) | 19.12 $\pm$ 0.11 |
| | FF-LoRA + AT ($r$=2) | 0.6M (0.3%) | 19.17 $\pm$ 0.25 |
| PLBART | FT (baseline) | 140M (100%) | **19.11** $\pm$ 0.19 |
| | LoRA ($r$=16) | 0.9M (0.6%) | 18.32 $\pm$ 0.04 |
| | FF-LoRA ($r$=16) | 1.5M (1.1%) | 18.38 $\pm$ 0.04 |
| | AT ($r$=16) | 1.1M (0.8%) | 18.03 $\pm$ 0.04 |
| | FF-LoRA + AT ($r$=4) | 0.6M (0.5%) | 18.40 $\pm$ 0.12 |

Table 6: PE methods quality on the test set in the **code summarization (Go)** task. Hyperparameter $r$ (methods' internal size) was chosen based on the quality on the dev set. Standard deviation is estimated in 3 runs.

| Model | Method | # Parameters | BLEU-4 (↑) |
|---|---|---|---|
| CodeT5 | FT (baseline) | 223M (100%) | **20.23** $\pm$ 0.09 |
| | LoRA ($r$=16) | 1.8M (0.8%) | 20.04 $\pm$ 0.02 |
| | FF-LoRA ($r$=16) | 3M (1.3%) | 20.04 $\pm$ 0.04 |
| | AT ($r$=16) | 2.2M (1%) | 20.01 $\pm$ 0.04 |
| | FF-LoRA + AT ($r$=16) | 5.2M (2.3%) | 20.15 $\pm$ 0.06 |
| PLBART | FT (baseline) | 140M (100%) | **19.51** $\pm$ 0.11 |
| | LoRA ($r$=16) | 0.9M (0.6%) | 17.32 $\pm$ 0.48 |
| | FF-LoRA ($r$=8) | 0.7M (0.5%) | 18.10 $\pm$ 0.06 |
| | AT ($r$=16) | 1.1M (0.8%) | 17.78 $\pm$ 0.05 |
| | FF-LoRA + AT ($r$=16) | 2.6M (1.8%) | 18.40 $\pm$ 0.08 |

Table 7: PE methods quality on the test set in the **code summarization (Java)** task. Hyperparameter $r$ (methods' internal size) was chosen based on the quality on the dev set. Standard deviation is estimated in 3 runs.

| Model | Method | # Parameters | BLEU-4 (↑) | EM (↑) | CodeBLEU(↑) |
|---|---|---|---|---|---|
| CodeT5 | FT (baseline) | 223M (100%) | **79.14** $\pm$ 0.30 | **66.23** $\pm$ 0.29 | **84.33** $\pm$ 0.25 |
| | LoRA ($r$=8) | 0.9M (0.4%) | 74.26 $\pm$ 0.12 | 59.27 $\pm$ 0.24 | 80.38 $\pm$ 0.12 |
| | FF-LoRA ($r$=8) | 1.5M (0.7%) | 75.34 $\pm$ 0.17 | 61.00 $\pm$ 0.22 | 81.08 $\pm$ 0.13 |
| | AT ($r$=16) | 2.2M (1%) | 75.30 $\pm$ 0.31 | 60.43 $\pm$ 0.37 | 81.08 $\pm$ 0.13 |
| | FF-LoRA + AT ($r$=16) | 5.2M (2.3%) | 76.89 $\pm$ 0.10 | 62.40 $\pm$ 0.22 | 82.53 $\pm$ 0.15 |
| PLBART | FT (baseline) | 140M (100%) | **78.60** $\pm$ 0.21 | **64.70** $\pm$ 0.08 | **84.12** $\pm$ 0.13 |
| | LoRA ($r$=16) | 0.9M (0.6%) | 71.20 $\pm$ 0.03 | 55.10 $\pm$ 0.16 | 78.48 $\pm$ 0.02 |
| | FF-LoRA ($r$=16) | 1.5M (1.1%) | 74.27 $\pm$ 0.27 | 59.10 $\pm$ 0.22 | 80.59 $\pm$ 0.14 |
| | AT ($r$=16) | 1.1M (0.8%) | 72.02 $\pm$ 0.20 | 55.30 $\pm$ 0.29 | 79.18 $\pm$ 0.18 |
| | FF-LoRA + AT ($r$=16) | 2.6M (1.9%) | 75.11 $\pm$ 0.25 | 59.87 $\pm$ 0.19 | 81.22 $\pm$ 0.17 |

Table 8: PE methods quality on the test set in the **code translation (C# → Java)** task. Hyperparameter $r$ (methods' internal size) was chosen based on the quality on the dev set. Standard deviation is estimated in 3 runs.

| 1 | GT | str - > list Convert XML to URL List . From Biligrab . |
|---|----|---|
|   | FT | Convert Sina XML to a list of URLs . |
|   | LoRA | Convert Sina XML to a list of URLs . |
|   | FF-LoRA | Convert Sina XML to a list of URLs . |
|   | AT | Convert an XML string to a list of URLs . |
| 2 | GT | Print a log message to standard error . |
|   | FT | Print a log message to stderr . |
|   | LoRA | Print a log message to stderr . |
|   | FF-LoRA | Print a log message to stderr . |
|   | AT | Print a log message to stderr . |
| 3 | GT | What a Terrible Failure! |
|   | FT | Print a message to stdout and exit . |
|   | LoRA | Print a message to the console . |
|   | FF-LoRA | Print a message to the console . |
|   | AT | Print a message to the console . |
| 4 | GT | Detect operating system . |
|   | FT | Detect OS . |
|   | LoRA | Detect the operating system . |
|   | FF-LoRA | Detects the operating system . |
|   | AT | Detect the operating system . |
| 5 | GT | Converts a string to a valid filename . |
|   | FT | Return text legitimately according to OS . |
|   | LoRA | Convert a string to a legitimate string . |
|   | FF-LoRA | Turn a string into a legitimate string . |
|   | AT | Turn a string into a legitimate string . |
| 6 | GT | Override the original one Ugly ugly dirty hack |
|   | FT | Download the stream . |
|   | LoRA | Download a stream . |
|   | FF-LoRA | Download a stream . |
|   | AT | Download a stream . |
| 7 | GT | Scans through a string for substrings matched some patterns . |
|   | FT | matchall - Searches text for all patterns in patterns and returns a list of matches |
|   | LoRA | Return a list of all matching patterns in text . |
|   | FF-LoRA | Match all patterns in text . |
|   | AT | Match all patterns in text . |
| 8 | GT | JSON int int int - > str Get a proper title with courseid + topicID + partID . |
|   | FT | Get the title of a specific topic and part by its index . |
|   | LoRA | Get the title of a topic and part . |
|   | FF-LoRA | Get the title of a topic part . |
|   | AT | Get the title for a given topic and part . |
| 9 | GT | int - > None Download a WHOLE course . Reuse the API call to save time . |
|   | FT | wanmen_download_by_course_topic and wanmen_download_by_course_part recursively |
|   | LoRA | Download wanmen by course . |
|   | FF-LoRA | Download wanmen files by course . |
|   | AT | Download wanmen by course . |
| 10 | GT | Returns and flush the event buffer . In case dag_ids is specified it will only return and flush events for the given dag_ids . Otherwise it returns and flushes all |
|   | FT | Clear the event buffer for the given DAGs . If no DAGs are given clear the event buffer . |
|   | LoRA | Get the event buffer for the given dag ids . |
|   | FF-LoRA | Get the event buffer for the given dag_ids . |
|   | AT | Get the event buffer for the given dag_ids . |

Table 9: CodeT5 predictions in code summarization (Python) task when using different PE methods.
GT stands for Ground Truth