# A Comprehensive Evaluation of Parameter-Efficient Fine-Tuning on Software Engineering Tasks

WENTAO ZOU, QI LI, JIDONG GE*, and CHUANYI LI*, National Key Laboratory for Novel Software Technology at Nanjing University, China

XIAOYU SHEN, Eastern Institute for Advanced Study, China

LIGUO HUANG, Department of Computer Science, Southern Methodist University, USA

BIN LUO, National Key Laboratory for Novel Software Technology at Nanjing University, China

Pre-trained models (PTMs) have achieved great success in various Software Engineering (SE) downstream tasks following the "pre-train then fine-tune" paradigm. As fully fine-tuning all parameters of PTMs can be computationally expensive, a widely used solution is parameter-efficient fine-tuning (PEFT), which freezes PTMs while introducing extra parameters. Though work has been done to test PEFT methods in the SE field, a comprehensive evaluation is still lacking. This paper aims to fill in this gap by evaluating the effectiveness of five PEFT methods on eight PTMs and four SE downstream tasks. For different tasks and PEFT methods, we seek answers to the following research questions: 1) Is it more effective to use PTMs trained specifically on source code, or is it sufficient to use PTMs trained on natural language text? 2) What is the impact of varying model sizes? 3) How does the model architecture affect the performance? Besides effectiveness, we also discuss the efficiency of PEFT methods, concerning the costs of required training time and GPU resource consumption. We hope that our findings can provide a deeper understanding of PEFT methods on various PTMs and SE downstream tasks. All the codes and data are available at https://github.com/zwtnju/PEFT.git.

CCS Concepts: • **Software and its engineering → Software creation and management**; • **Computing methodologies → Artificial intelligence**.

Additional Key Words and Phrases: parameter-efficient fine-tuning, pre-trained model, software engineering, empirical study, transfer learning, effectiveness and efficiency

## 1 INTRODUCTION

Pre-training language models on large amounts of data with self-supervision, and then fine-tuning them on downstream tasks has become a dominant approach in natural language processing (NLP) [1, 13, 43, 55, 61, 67, 82]. The same paradigm has also been applied to software engineering

*Corresponding Authors.

Authors' addresses: Wentao Zou, DZ1832005@smail.nju.edu.cn; Qi Li, MG21320006@smail.nju.edu.cn; Jidong Ge, gjd@nju.edu.cn; Chuanyi Li, lcy@nju.edu.cn, National Key Laboratory for Novel Software Technology at Nanjing University, Nanjing, Jiangsu, China; Xiaoyu Shen, Eastern Institute for Advanced Study, China, xyshen@eitech.edu.cn; Liguo Huang, Department of Computer Science, Southern Methodist University, USA, lghuang@lyle.smu.edu; Bin Luo, luobin@nju.edu.cn, National Key Laboratory for Novel Software Technology at Nanjing University, Nanjing, Jiangsu, China.

(SE) downstream tasks and demonstrated exceptional performances [22, 28, 39, 64, 82, 88, 93]. Despite the remarkable success, full fine-tuning is known to be cost inefficient [62, 73, 74], which requires fine-tuning and saving the full set of parameters for each downstream task [92]. When PTMs become large, a lot of GPU and memory resources are required and the computational cost is demanding.

To overcome the above drawbacks, a popular strategy is parameter-efficient fine-tuning (PEFT) [16]. In contrast to full fine-tuning, PEFT only fine-tunes a limited amount of extra parameters for each task while freezing the base PTM to retain its original knowledge. Figure 1 illustrates the architecture of transformer blocks with frozen and fine-tuned layers for three types of PEFT methods: adapter Houlsby et al. [25], prefix tuning [32], and LoRA [26]. PEFT significantly reduces the computing cost as gradients flow only through the small amount of extra parameters. After training, we only need to store one fixed copy of the original large PTM and every downstream task has its own small set of extra parameters. This can also greatly reduce the storage cost. Though fine-tuning only a small set of extra, evidence has been shown that PEFT methods could even surpass the full fine-tuning approach on certain tasks [23].

In the SE field, various PEFT methods have been evaluated in recent years. For example, Ayupov and Chirkova [5] first tested Houlsby adapter tuning [25] and LoRA [26] methods on four code-related tasks. They found that although PEFT approaches may achieve comparable or higher performance than full fine-tuning on code understanding tasks, they underperform full fine-tuning on code generation tasks. Later, Wang et al. [78] used the Houlsby adapter tuning method on code search and code summarization tasks to improve the performance of PTMs and showed its effectiveness in cross-language and low-resource scenarios. Recently, Liu et al. [35] conducted an empirical study of four PEFT methods in several SE scenarios, including low-resource, cross-language, and cross-project.

Although some studies have applied PEFT to SE downstream tasks, they often focus on the effectiveness of PEFT methods on specific scenarios or tasks, lacking a comprehensive comparison of various PTMs from a holistic view. Firstly, it is also unclear whether it is necessary to develop code-specific PTMs (PTMs trained specifically on source code), or is sufficient to use text-only PTMs (PTMs trained on natural language text) and rely on the cross-modal transferability of PEFT on source code. Although PEFT methods have been well applied to cross-language and cross-project scenarios, there are few studies [17] on the cross-modal transfer ability of text-only PTMs. If text-only PTMs can achieve good results with PEFT methods, we no longer need to collect large-scale datasets to pre-train a new PTM specifically for code-related tasks. Considering the huge cost of pre-training, directly using text-only PTMs can save a lot of human and computation resources. Secondly, the size (the number of trainable parameters [9, 69]) of a PTM has a big impact on the effectiveness of fine-tuning. In general, the larger a PTM is, the better its performance will be. The purpose of PEFT is to reduce the trainable parameters to achieve close results. Therefore, it is necessary to compare the same PTM with different sizes. If the PEFT methods have close effects on different parameters, we can give priority to smaller PTMs to further reduce costs. Thirdly, Niu et al. [48] have explored the performance of PTMs in different architectures, such as encoder-only or encoder-decoder architecture [47], on various tasks. Their work highlights the effects of different architectures when fine-tuning PTMs. Different architectures of PTMs may also have different impacts on PEFT methods, which are ignored in previous studies. Guided by these aspects, we formulate the following three research questions to direct our study:

- RQ1: How do text-only PTMs perform compared to code-specific PTMs of the same architecture and size?
- RQ2: How do the same text-only/code-specific PTMs with different sizes perform?
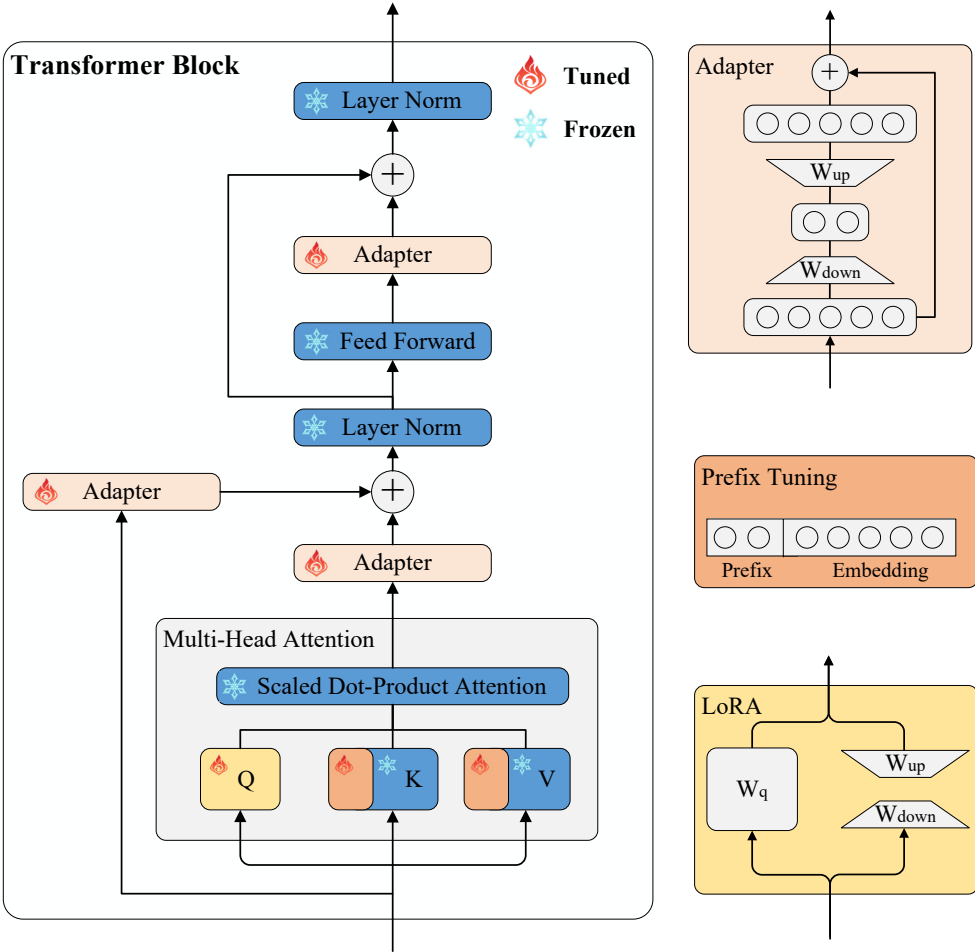
Fig. 1. The Architecture of Transformer Blocks with Frozen and Fine-tuned Layers for Three Types of PEFT Methods: Adapter, Prefix Tuning, and LoRA.

- RQ3: How do text-only/code-specific PTMs in different architectures perform?

Throughout the paper, we seek answers to the above research questions under five different PEFT methods: Houlsby [25], Pfeiffer [52], Parallel [23], Prefix [32], and LoRA [26]. Experiments are run on four code-related tasks: clone detection of two code pairs, defect detection of one code snippet, code search across 6 programming languages, and code translation between Java and C#. We include 8 PTMs in our study, including CodeBERT [15], RoBERTa [38], CodeT5 [80], T5 [55], CodeT5-large [80], T5-large [55], UniXcoder [19], and BART [31].

Besides the effectiveness, we also focus on the efficiency of PEFT methods. Most existing studies acknowledge that PEFT can save resources compared to full fine-tuning, but they focus on comparing the effects of PEFT and rarely present results regarding the resources required. Therefore, we further quantitatively discuss the costs of PEFT methods in terms of required training time and GPU resource consumption.

The main contributions of our paper are as follows:

- We show that PEFT methods are more advantageous in cross-modal transfer learning of PTMs. Compared with the results of code-specific PTMs, PEFT methods applied to text-only PTMs are more effective in keeping the performance of full fine-tuning PTMs on clone detection, defect detection, and code translation.
- We find that compared with full fine-tuning, increasing the trainable parameters of T5 reduces the performance of the PEFT methods. CodeT5 with more trainable parameters performs better on defect detection and code translation and performs worse on clone detection and code search.
- We present that PEFT methods applied to encoder-only architecture perform better on clone detection and defect detection than encoder-decoder architecture, while those used in T5-based encoder-decoder architecture are better on code search.
- We investigate that PEFT PTMs obtain comparable, or even better results compared with full fine-tuning PTMs, especially on the code search and code translation. Among PEFT methods, the Parallel performs the best, which is consistent with the conclusion of He et al. [23].
- We demonstrate that PEFT methods can save about 10%-30% GPU resources compared with full fine-tuning methods. Pfeiffer needs the least GPU resources while Prefix and Houlsby require the most. PEFT does not promise to improve training speed, but it still has the potential to reduce training time. As the trainable parameters increase, fewer GPU resources and more training time are required.
- We conduct extensive experiments and have made all the PEFT layers and PTMs used in our study publicly accessible. This allows for easy reproduction of our results and to utilize the best-performing models, providing broader applications to other software engineering scenarios.

The remainder of this paper is organized as follows. Section 2 provides some necessary prior knowledge. Section 3 outlines research questions, and provides details on selected datasets and PTMs, and Section 4 analyzes the results related to research questions. Section 5 discusses qualitative effectiveness analysis, quantitative efficiency analysis, and potential threats to validity. Section 6 presents related work, and Section 7 concludes this paper.

## 2 PRELIMINARIES

### 2.1 Transformer and Pre-Trained Models

Transformer [76] is a deep learning model to solve sequence-to-sequence tasks. It can handle long text dependencies and has the flexibility to handle variable-length input and output sequences, making it suitable for tasks such as machine translation, text summarization, and sentiment analysis. Transformer has significantly influenced the NLP field as well as many other fields, such as SE.

The transformer has many key components, including a multi-head attention (MHA) layer, a fully connected feed-forward network (FFN), layer normalization, and residual connection. The Transformer architecture with these key components is shown in Figure 2.

The core of Transformer is the self-attention mechanism, which allows it to weigh the importance of different elements in a sequence. The self-attention mechanism enables the model to capture relationships between words or tokens. The calculation of attention function depends on the three components of queries $Q \in \mathbb{R}^{n \times d}$, keys $K \in \mathbb{R}^{m \times d}$ and values $V \in \mathbb{R}^{m \times d}$, where $n$ and $m$ are the number of queries and key-value pairs, and $d$ denotes the hidden size of the Transformer. The scaled dot-product attention can be formulated as follows:

$$Attention(K, Q, V) = softmax(\frac{QK^T}{\sqrt{d}})V \tag{1}$$

Fig. 2. The simplified Transformer Architecture.

Based on self-attention, MHA performs the attention function in parallel over $N_h$ heads, where $i$th head is parameterized by $W_q^i, W_k^i, W_v^i \in \mathbb{R}^{d \times d_h}$ to project inputs to queries, keys, and values. In MHA, $d_h$ is typically set to $\frac{d}{N_h}$ to save parameters, which indicates that each attention head is operating on a lower dimensional space. Given a sequence of $m$ vectors $X \in \mathbb{R}^{m \times d}$ and a query vector $x \in X$ (the inputs of queries, keys, and values of self-attention are the same), the attention in the $i$th head can be computed as:

$$head_i = Attention(xW_q^i, XW_k^i, XW_v^i) \tag{2}$$

Then, given the initial parameters $W_0 \in \mathbb{R}^{d \times d}$, the MHA concatenates the attention results of different heads as follows:

$$MHA(x) = concat(head_1, \cdots, head_h)W_0 \tag{3}$$

FFN consists of two linear transformations with an activation function of ReLU in between, formulated as follows:

$$FFN(x) = ReLU(xW_1 + b_1)W_2 + b_2 \tag{4}$$

Fig. 3. Different Types of Pre-trained Models.

where $W_1 \in \mathbb{R}^{d \times d_m}$, $W_2 \in \mathbb{R}^{d_m \times d}$ are trainable parameters, and $d_m$ is typically set to 4d and acts as a scaling factor.

Transformer has become the foundation model and has been adapted to various models, which are called PTMs. PTMs are usually train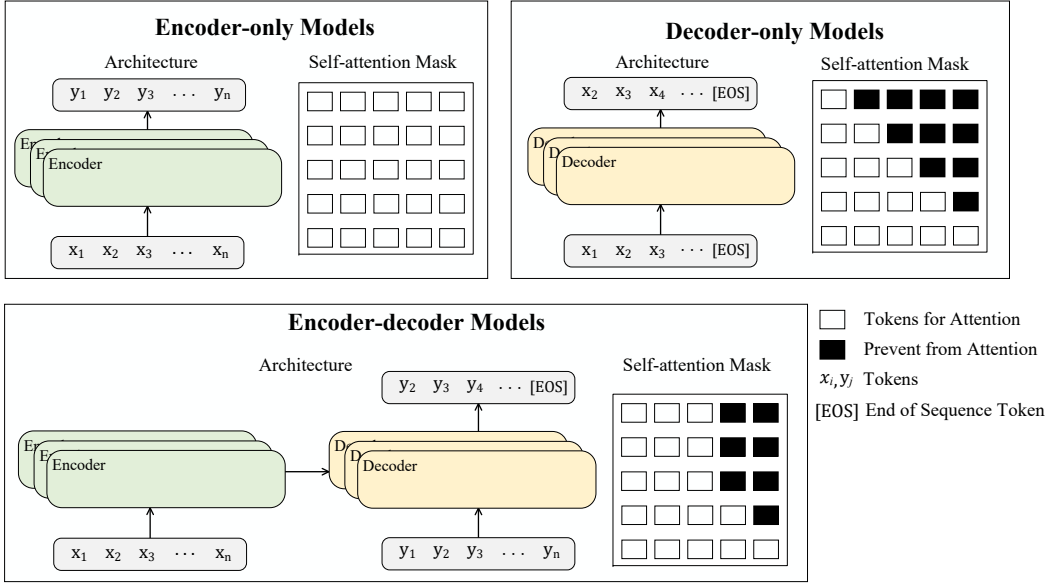ed on extensive datasets over pre-training tasks. Pre-training tasks are specially designed general tasks that his somewhat related to downstream tasks, allowing PTMs to learn effective representations that can be applied to a wide range of downstream tasks. Transformer consists of encoders and decoders. Encoders process the input sequence, while decoders generate the output sequence based on the encoder's context and the previously generated tokens. PTMs are built upon the Transformer architecture and often use multiple encoders and/or decoders. They can be divided into encoder-only, decoder-only, and encoder-decoder models, as is shown in Figure 3.

## 2.2 Full Fine-Tuning and PEFT Methods

PTMs are becoming more important because non-pre-training methods require large labeled datasets for training, which may result in slow training, information loss, and model instability. The common paradigm of using PTMs is to pre-train and then fully fine-tune, and it has proven advantages in many SE downstream tasks, such as clone detection and code search. Full fine-tuning requires tuning all parameters of a PTM, including embedding layers, MHA layers, and FFN layers, which enables the PTM to transfer learned text representations to unknown tasks. However, as PTMs become larger, the cost of full fine-tuning increases. Many techniques have been proposed to reduce the heavy computational demands of large models [8, 18, 63, 72] and a popular strategy is PEFT. PEFT focuses on optimizing fine-tuning that only tunes a small number of parameters or learns external modules for new tasks while keeping most of the PTMs' parameters frozen. PEFT can achieve similar effects with full fine-tuning and is particularly useful when computing resources and memory are limited.

(a) Houlsby

(b) Prefix

(c) LoRA

(d) Parallel

Fig. 4. Illustration of used PEFT methods.

Next, we introduce the PEFT methods used in this paper, including Houlsby, Pfeiffer, Prefix, LoRA, and Parallel, as illustrated in Figure 4. Houlsby and Pfeiffer use adapters in similar ways, so we only show the adapter architecture of Houlsby in the figure. Houlsby and Pfeiffer belong to "sequential" computation that uses $h$, the output of the MHA/FFN layer to compute $\Delta h$. In contrast, there is another way of "parallel" computation. The "parallel" methods contain Prefix, LoRA, and Parallel, and directly use $x$, the input of the MHA/FFN layer, to compute $\Delta h$.

Houlsby [25] approach inserts small modules, called adapters, after MHA and FFN layers. The adapter layer uses a down-projection $W_{down} \in \mathbb{R}^{d \times r}$ to project the input vector to a lower-dimensional vector of dimensional size $r$, a nonlinear activation function $f$, and a up-projection

$W_{up} \in \mathbb{R}^{r \times d}$ to project lower-dimensional vector to the dimension of input vector. Given an output vector $h \in \mathbb{R}^d$ of MHA/FFN layers, the output after an adapter can be represented by:

$$h \leftarrow h + \Delta h, \quad \Delta h = f(hW_{down})W_{up} \tag{5}$$

Pfeiffer [52] is a more efficient variant in which adapters are inserted only after the FFN layers.

Prefix [32] introduces new parameters in the MHA layers. Specifically, it adds $l$ trainable prefix vectors to the keys and values of the MHA and concatenates prefix vectors $P_k \in \mathbb{R}^{l \times d}$ and $P_v \in \mathbb{R}^{l \times d}$ to the original keys $K$ and values $V$, respectively. Then MHA is performed on the new prefixed keys and values and the computation of $head_i$ in Eq. 2 becomes:

$$head_i = Attention(xW_q^i, concat(p_k^i, XW_k^i), concat(p_v^i, XW_v^i)) \tag{6}$$

where $P_k$ and $P_v$ are split into $N_h$ head vectors and $p_k^i \in \mathbb{R}^{l \times d / N_h}$, $p_v^i \in \mathbb{R}^{l \times d / N_h}$ denote the $i$th head vector of $P_k$ and $P_v$. Following He et al. [23], the Eq. 6 can be rewrote as:

$$h \leftarrow (1 - \lambda(x))h + \lambda(x)\Delta h, \quad \Delta h := softmax(xW_q P_k^\top)P_v \tag{7}$$

where $\lambda(x)$ is a scalar representing the sum of normalized attention weights over prefixes and can be represented by:

$$\lambda(x) = \frac{\sum_i exp(xW_q P_k^\top)_i}{\sum_i exp(xW_q P_k^\top)_i + \sum_j exp(xW_q W_k^\top X^\top)_j} \tag{8}$$

LoRA [26] is short for Low-Rank Adaptation. It injects trainable low-rank decomposition matrices into transformer layers. This can be applied to any weight matrix, but LoRA only applies to the query $Wq$ and value $Wv$ projection matrices of MHA layers. For any layer expressed as a matrix multiplication of the form $h = xW_0$, $x \in \mathbb{R}^d$, $W_0 \in \mathbb{R}^{d \times k}$, it performs a reparameterization such that:

$$h \leftarrow h + \Delta h, \quad \Delta h = s \cdot xW_{down}W_{up} \tag{9}$$

where $s \geq 1$ is a trainable scalar hyperparameter, $W_{down} \in \mathbb{R}^{d \times r}$, $W_{up} \in \mathbb{R}^{r \times k}$ are trainable parameters, and $r$, the low-dimensional rank of the decomposition, is the most important hyperparameter.

Parallel [23] inserts adapter layers in parallel to Transformer layers, and the design of Parallel is similar to Eq. 7 and Eq. 9. Different from sequential adapters, the input of parallel adapters is the input of the MHA/FFN layers rather than the output of the MHA/FFN layers. According to Eq. 5, given the input $x$ and the output $h$ of MHA/FFN layer, the output after Parallel can be expressed by:

$$h \leftarrow h + \Delta h, \quad \Delta h = f(xW_{down})W_{up} \tag{10}$$

He et al. [23] also proposed two variants of parallel adapters. The multi-head parallel adapter applies the parallel adapter to modify the head attention outputs. The scaled parallel adapter transfers the composition and insertion form of LoRA into an adapter. In this paper, we use the original parallel adapters to fine-tune PTMs.

## 3 EXPERIMENTAL SETUPS

### 3.1 Research Questions

We evaluate five PEFT methods, **Houlsby, Pfeiffer, Prefix, LoRA**, and **Parallel**, and all the PEFT methods are from AdapterHub [51]. In the paper, we focus on their effectiveness on various PTMs, aiming to answer the following three research questions:

- **RQ1: How do text-only PTMs perform compared to code-specific PTMs of the same architecture and size?**
  [35] have studied the effectiveness of PEFT methods, but they only focus on code-specific PTMs. In addition to using code-specific PTMs on code-related tasks, another common strategy is to use PTMs pre-trained on other datasets [4, 41], such as natural language (NL). The PEFT methods have the potential for cross-modal transfer learning [60, 81, 90]. It is worth discussing the effectiveness of PEFT methods on text-only PTMs as they may show different effects from the code-specific PTMs. Furthermore, designing code-specific PTMs and collecting code-related datasets for pre-training is labor-intensive. If PEFT methods can help text-only PTMs achieve good results on SE downstream tasks, the required resources can be reduced.
- **RQ2: How do the same text-only/code-specific PTMs with different sizes perform?**
  The same PTM with different sizes usually have different learning abilities for code representation, and larger PTMs usually have better performance [87]. The size of a PTM is also an important factor for PEFT, as PEFT methods help PTMs with fewer parameters to achieve results that are close to or better than those with full fine-tuning. PEFT methods may close the performance gap of PTMs with different sizes, and we can use PTMs with fewer parameters to reduce the cost. Since text-only PTMs and code-specific PTMs may have inconsistent results, we discuss them with different sizes respectively.
- **RQ3: How do text-only/code-specific PTMs in different architectures perform?**
  There is evidence that different model architectures have preferences for SE downstream tasks [48]. For example, PTMs with encoder-only architecture are better at code understanding tasks, while PTMs with encoder-decoder architecture are better at code generation tasks. Similar conclusions may exist on PEFT methods, and we want to explore the impact of different architectures and which architecture is more suitable for PEFT methods. The same with RQ2, we discuss text-only PTMs and code-specific PTMs in different architectures respectively.

### 3.2 Tasks, Datasets, and Evaluation Metrics

For each research question, we fine-tune PTMs on four SE downstream tasks: **clone detection, defect detection, code search**, and **code translation**.

Clone detection [70] is a code understanding [79] task that involves predicting whether pairs of code are functionally equivalent and the output is a binary classification value. We use the BigCloneBench (BCB) [40] dataset, which contains over 6 million true clone pairs and 260k false clone pairs.

Defect detection [65] is another code understanding task to predict whether a function is vulnerable. Unlike clone detection, which considers pairs of code snippets, defect detection focuses on individual code snippets. We use the Devign [95] dataset, which consists of functions collected from two open-source projects and contains approximately 22k training data. Compared with clone detection, the dataset of defect detection is low-resource.

Code search [83] is a code retrieval task [44] to find the most relevant code snippets from a given set based on an NL query, where the NL usually refers to the source code comments. We utilize the CodeSearchNet (CSN) [27] dataset, which includes Go, Java, JavaScript (JS), PHP, Python (Py), and Ruby languages. Low-quality queries are filtered following the UniXcoder.

Code translation [40] is a code generation [84] task that requires translating code from one programming language to another. We use the CodeTrans [40] dataset, which contains parallel code snippets in Java and C#. The PTMs need to generate C#/Java code for the given corresponding Java/C# code.

Table 1. A Brief Summary of Datasets.

| Task | Dataset Name | Language | Train Size | Valid Size | Test Size |
|------|-------------|----------|-----------|-----------|-----------|
| Clone Detection | BigCloneBench | Java | 901,028 | 415,416 | 415,416 |
| Defect Detection | Devign | C | 21,854 | 2,732 | 2,732 |
| Code Search | CodeSearchNet | Go | 167,288 | 7,325 | 8,122 |
| | | Java | 164,923 | 5,183 | 10,955 |
| | | JavaScript | 58,025 | 3,885 | 3,291 |
| | | PHP | 241,241 | 12,982 | 14,014 |
| | | Python | 251,820 | 13,914 | 14,918 |
| | | Ruby | 24,927 | 1,400 | 1,261 |
| Code Translation | CodeTrans | Java → C# | 10,295 | 499 | 1,000 |
| | | C# → Java | 10,295 | 499 | 1,000 |

We give a brief introduction to datasets of these downstream tasks and the statistics of the datasets are listed in Table 1. Fine-tuning uses the training and validation datasets, and the evaluation of the downstream task uses the testing datasets. These tasks are representative of various SE downstream tasks, including code classification, code retrieval, and code generation, and the datasets include high-resource, low-resource, monolingual, and multilingual. All the datasets for fine-tuning are sourced from CodeXGLUE [40], and we use their provided code to pre-process the datasets for all tasks.

For classification and retrieval tasks, metrics such as Accuracy (Acc), F1 [45], Precision (P), Recall (R), Mean Reciprocal Rank (MRR), and Mean Average Precision (MAP) are often used. For generation tasks, metrics such as BLEU (B.) [50], and CodeBLEU (C.B.) [57], are usually used. Some generation tasks have also used variants of Accuracy, such as Exact Match (EM), indicating whether the sequence generated by the model perfectly matches the correct answer, and Computational Accuracy (CA), computing the number of times the hypothesis function generates the same output as the reference when given the same inputs. Following Niu et al. [48], the evaluation metrics of clone detection, defect detection, and code search are F1, Acc, and MRR, separately. For code translation, since B. and C.B. cannot accurately indicate whether the predicted code and correct answer are semantically equivalent, we use EM as the primary evaluation metric and only discuss the results of EM in this paper.

## 3.3 PTMs and Implementation Details

In the paper, the used PTMs to fine-tune are **CodeBERT**, **RoBERTa**, **CodeT5**, **T5**, **CodeT5-large**, **T5-large**, **UniXcoder** and **BART**. A brief description of some important information about the used PTMs is listed in Table 2.

BERT is one of the most famous text-only PTMs, which is a bidirectional encoder-only architecture model for text representation from Transformer. It uses masked language modeling and next sentence prediction methods to train the model. Derived from BERT, RoBERTa [38] is an optimization encoder-only model with just the masked language modeling task and benefits from additional training, larger batch sizes, and longer sequences. BART and T5 are two classic different encoder-decoder architecture models. BART [31] uses denoising autoencoder tasks, such as token masking and token deletion, to construct the model. T5 [55] treats every task as a text-to-text problem, which involves feeding input text into the model and generating target text. CodeBERT [15] is one of the first encoder-only code-specific PTMs. It adapts BERT for SE downstream tasks and

Table 2. A Brief Description of PTMs.

| PTMs | Field | Architecture | Learned Languages | Model Size (M) |
|---|---|---|---|---|
| CodeBERT | code-specific | BERT-based encoder-only | Go, Java, JavaScript, PHP, Python, Ruby | 125 |
| RoBERTa | text-only | BERT-based encoder-only | NL text | 125 |
| CodeT5 | code-specific | T5-based encoder-decoder | Go, Java, JavaScript, PHP, Python, Ruby, C, C# | 223 |
| T5 | text-only | T5-based encoder-decoder | NL text | 223 |
| CodeT5-large | code-specific | T5-based encoder-decoder | Go, Java, JavaScript, PHP, Python, Ruby, C, C# | 738 |
| T5-large | text-only | T5-based encoder-decoder | NL text | 738 |
| UniXcoder | code-specific | cross-modal encoder-only | Go, Java, JavaScript, PHP, Python, Ruby | 126 |
| BART | text-only | BART-based encoder-decoder | NL text | 139 |

Table 3. Hyperparameters for Fine-Tuning Downstream Tasks.

| Types | Clone Detection | Defect Detection | Code Search | Code Translation |
|---|---|---|---|---|
| **Vector Dimension**: | | | | |
| T5-large or CodeT5-large | 128 | 256 | code: 256; nl: 128 | input: 512; output: 512 |
| Other Pre-trained Models | 400 | 400 | code: 256; nl: 128 | input: 512; output: 512 |
| **Batch Size**: | | | | |
| All Pre-trained Models | 32 | 32 | 64 | 32 |
| **Training Epochs**: | | | | |
| All Pre-trained Models | 2 | 5 | JavaScript: 5; Ruby: 10; Others: 2 | 50k steps |
| **Learning Rate**: | | | | |
| All Pre-trained Models | 5e-5 | 2e-5 | 2e-5 | 5e-5 |

proposes a new replaced token detection task, training on both natural languages and programming languages. Based on T5, CodeT5 [80] introduces identifier-aware pre-training tasks to leverage code identifiers. This strategy improves the performance on SE downstream tasks compared with T5. In addition to individual architectures, UniXcoder [19] implements a unified cross-modal architecture and uses prefixes and different self-attention masks to control the behavior of the model. It is pre-trained on multi-modal data, which includes code, comments, and AST.

All these PTMs come from Hugging Face[1], and we select different PTMs for each research question. For RQ1, we employed two text-only PTMs (RoBERTa and T5) and code-specific PTMs (CodeBERT and CodeT5) for fine-tuning. CodeBERT and RoBERTa have the same size and are BERT-based architecture. T5 and CodeT5 share the same size and T5-based architecture. For RQ2, we selected the T5 and CodeT5 with different sizes: basic and large. We used T5 and CodeT5 to refer to the basic versions to maintain consistency with RQ1, and T5-large and CodeT5-large to refer to the large versions of T5 and CodeT5 respectively. For RQ3, we applied three code-specific PTMs (CodeBERT, UniXcoder, and CodeT5) and text-only PTMs (RoBERTa, BART, and T5). CodeBERT and UniXcoder are both encoder-only architectures, but UniXcoder can use the input prefixes to learn encoder-decoder and decoder-only architectures, so they can be regarded as different architectures. Similarly, T5 uses the end-to-end pre-training strategy, while BART uses denoising autoencoding to pre-train. For convenience, they can also be regarded as different encoder-decoder architectures. In the following sections, we refer to the architectures of CodeBERT, UniXcoder, BART, and T5 as BERT-based encoder-only, cross-modal encoder-only, BART-based encoder-decoder and T5-based encoder-decoder architecture respectively.

---

[1]https://huggingface.co

Table 4. [RQ1] The Results of BERT-based PTMs on Clone Detection, Defect Detection, and Code Translation. The subscript **Full** represents the result of full fine-tuning, **Best** indicates the result of the best PEFT method, and other subscripts mean the results of the corresponding PEFT methods.

| Fine-Tune Methods | Clone Detection | Defect Detection | Code Translation | |
| | BCB | Devign | Java $\rightarrow$ C# | C# $\rightarrow$ Java |
| | F1 | Acc | EM | EM |
| $CodeBERT_{Full}$ | 94.95 | 63.91 | 59 | 60.1 |
| $CodeBERT_{Best}$ | -1.88 | -3.08 | -1.7 | +0.3 |
| $CodeBERT_{Houlsby}$ | 91.9 | 57.87 | 56.7 | 59.1 |
| $CodeBERT_{Pfeiffer}$ | 91.62 | 56.41 | 56.4 | 58.3 |
| $CodeBERT_{Parallel}$ | **93.07** | **60.83** | **57.3** | **60.4** |
| $CodeBERT_{Prefix}$ | 89.96 | 56.3 | 55.2 | 59 |
| $CodeBERT_{LoRA}$ | <u>88.38</u> | <u>55.6</u> | <u>54.4</u> | <u>57.2</u> |
| $RoBERTa_{Full}$ | 93.5 | 56.62 | 58.8 | 60.8 |
| $RoBERTa_{Best}$ | -0.15 | +2.41 | -0.2 | -0.2 |
| $RoBERTa_{Houlsby}$ | 90.19 | <u>54.03</u> | 56 | 59.4 |
| $RoBERTa_{Pfeiffer}$ | 90.24 | <u>54.03</u> | 57.4 | 58.8 |
| $RoBERTa_{Parallel}$ | **93.35** | **58.93** | **58.6** | **60.6** |
| $RoBERTa_{Prefix}$ | 89.63 | 55.53 | 56.1 | <u>58.3</u> |
| $RoBERTa_{LoRA}$ | <u>87.37</u> | 54.06 | <u>55</u> | 58.4 |

Our implementation for fine-tuning is derived from two open-source repositories. We adopt code from the UniXcoder[2] repository on the code search task and CodeXGLUE[3] repository on the other three tasks. We use their original code directly for fine-tuning encoder-only models and modify the original code to accommodate encoder-decoder models. The key hyperparameters for PTMs to fine-tune on different downstream tasks are listed in Table 3 and other hyperparameters are set to follow their original release. All experiments rely on pure Python packages and are conducted on servers with 2/4 GPUs of NVIDIA Tesla V100 32GB.

# 4 EXPERIMENTAL RESULTS

We list and analyze the experimental results for all research questions in this section.

## 4.1 RQ1: Code-specific PTMs and Text-only PTMs of Same Size and Architecture

*Clone Detection:* From Table 4 and Table 5, Parallel outperforms other PEFT methods on all the four PTMs, no matter the code-specific or text-only PTMs. Compared with CodeBERT, PEFT methods applied to RoBERTa are more effective in keeping the performance of PTMs. The results on the best PEFT and full fine-tuning also show that PEFT methods applied to T5 are more effective than CodeT5. Thus, PEFT methods are more useful for text-only PTMs than code-specific PTMs on clone detection.

*Defect Detection:* From Table 4 and Table 5, Parallel ranks first on CodeBERT and CodeT5. Furthermore, Parallel achieves the best results on RoBERTa and T5, and it outperforms full fine-tuning by 2.41% and 2.3%, separately. CodeBERT and CodeT5 are pre-trained on code-related datasets, and their full fine-tuning results are better than RoBERTa and T5. However, the best

[2]https://github.com/microsoft/CodeBERT/tree/master/UniXcoder
[3]https://github.com/microsoft/CodeXGLUE

Table 5. [RQ1] The Results of T5-based PTMs on Clone Detection, Defect Detection, and Code Translation.

| Fine-Tune Methods | Clone Detection | Defect Detection | Code Translation | |
|---|---|---|---|---|
| | BCB | Devign | Java → C# | C# → Java |
| | F1 | Acc | EM | EM |
| $CodeT5_{Full}$ | 93.93 | 63.65 | 65.6 | 67.3 |
| $CodeT5_{Best}$ | -1.2 | -6.4 | -5.4 | -1.5 |
| $CodeT5_{Houlsby}$ | 90.41 | 56.55 | 49.7 | 55 |
| $CodeT5_{Pfeiffer}$ | 92.67 | 56.37 | 52 | 57.3 |
| $CodeT5_{Parallel}$ | **92.73** | **57.25** | **60.2** | **65.8** |
| $CodeT5_{Prefix}$ | 90.81 | 55.53 | 44.3 | <u>47.8</u> |
| $CodeT5_{LoRA}$ | <u>88.28</u> | <u>54.5</u> | <u>42.8</u> | 52 |
| $T5_{Full}$ | 91.86 | 58.42 | 52 | 52.2 |
| $T5_{Best}$ | -1.24 | +2.3 | +2.3 | +3 |
| $T5_{Houlsby}$ | 88.57 | 54.06 | 37.2 | 43.8 |
| $T5_{Pfeiffer}$ | 90.2 | <u>53.99</u> | 41.3 | 45.3 |
| $T5_{Parallel}$ | **90.62** | **60.72** | **54.3** | **55.2** |
| $T5_{Prefix}$ | 82.78 | 54.06 | 33 | 37.9 |
| $T5_{LoRA}$ | <u>84.56</u> | 54.47 | <u>31.2</u> | <u>33.5</u> |

Table 6. [RQ1] MRR Score (%) of BERT-based PTMs on Code Search.

| Fine-Tune Methods | Go | Java | JavaScript | PHP | Python | Ruby |
|---|---|---|---|---|---|---|
| $CodeBERT_{Full}$ | 90.4 | 70.5 | 63.3 | 64.2 | 69 | 69.6 |
| $CodeBERT_{Best}$ | -0.5 | -0.7 | -0.9 | -0.3 | -0.5 | -1.1 |
| $CodeBERT_{Houlsby}$ | 86.7 | 64.4 | 49.6 | 59.7 | 63.3 | 51.6 |
| $CodeBERT_{Pfeiffer}$ | 86.6 | 63.9 | 49.9 | 59.3 | 62.8 | 51.4 |
| $CodeBERT_{Parallel}$ | **89.9** | **69.8** | **62.4** | **63.9** | **68.5** | **68.5** |
| $CodeBERT_{Prefix}$ | 87.4 | 65.2 | 54.9 | 59.4 | 63 | 56.2 |
| $CodeBERT_{LoRA}$ | <u>84</u> | <u>55.8</u> | <u>46.5</u> | <u>53.9</u> | <u>57.5</u> | <u>41.2</u> |
| $RoBERTa_{Full}$ | 88.5 | 64.5 | 54.6 | 59.8 | 62.2 | 60.6 |
| $RoBERTa_{Best}$ | -0.9 | -3.7 | -3.8 | -3.3 | -3.8 | -3.4 |
| $RoBERTa_{Houlsby}$ | 81.9 | 49.3 | 36.1 | 47.1 | 48 | 35.2 |
| $RoBERTa_{Pfeiffer}$ | 82.3 | 49.7 | 37.9 | 47.1 | 48.5 | 41.2 |
| $RoBERTa_{Parallel}$ | **87.6** | **60.8** | **50.8** | **56.5** | **58.4** | **57.2** |
| $RoBERTa_{Prefix}$ | 83.6 | 50.7 | 39.5 | 48.1 | 48.2 | 45.4 |
| $RoBERTa_{LoRA}$ | <u>78.5</u> | <u>43.3</u> | <u>29.4</u> | <u>42.2</u> | <u>43.3</u> | <u>33.8</u> |

PEFT methods are worse on CodeBERT/CodeT5 than that on RoBERTa/T5. Thus, PEFT methods are more effective for text-only PTMs than code-specific PTMs. LoRA achieves the worst results on CodeBERT and CodeT5. Houlsby and Pfeiffer perform the least effectively on RoBERTa and T5. Compared with the results of clone detection, the performance among different fine-tuning methods is more different. This implies that for tasks that perform poorly with full fine-tuning, the results of PEFT methods will be more unstable.

Table 7. [RQ1] MRR Score (%) of T5-based PTMs on Code Search.

| Fine-Tune Methods | Go | Java | JavaScript | PHP | Python | Ruby |
|---|---|---|---|---|---|---|
| $CodeT5_{Full}$ | 91 | 72.3 | 68 | 65.3 | 70.5 | 73.2 |
| $CodeT5_{Best}$ | -0.5 | -0.1 | +0.1 | -0.1 | +0.3 | +0.9 |
| $CodeT5_{Houlsby}$ | 88.5 | 68.8 | 62.4 | 62.7 | 67.5 | 67.7 |
| $CodeT5_{Pfeiffer}$ | 89.3 | 69.9 | 64.4 | 63.5 | 68.6 | 69.4 |
| $CodeT5_{Parallel}$ | **90.5** | **72.2** | **68.1** | **65.2** | **70.8** | **74.1** |
| $CodeT5_{Prefix}$ | <u>87.4</u> | <u>62.2</u> | <u>54.9</u> | <u>58.1</u> | <u>60.7</u> | <u>63.3</u> |
| $CodeT5_{LoRA}$ | 87.5 | 67.4 | 58.2 | 61.4 | 65.9 | 63.7 |
| $T5_{Full}$ | 86.1 | 57.5 | 45.8 | 53.2 | 56.3 | 49 |
| $T5_{Best}$ | +1.7 | +5.3 | +6.6 | +3.4 | +3.4 | +11.5 |
| $T5_{Houlsby}$ | 85.1 | 53.3 | 42.7 | 50 | 52.8 | 45.2 |
| $T5_{Pfeiffer}$ | 85.4 | 54.7 | 44.4 | 51 | 53.7 | 50.3 |
| $T5_{Parallel}$ | **87.8** | **62.8** | **52.4** | **56.6** | **59.7** | **60.5** |
| $T5_{Prefix}$ | 77 | 41.3 | 33.3 | 39.5 | 42 | 38.2 |
| $T5_{LoRA}$ | <u>72.4</u> | <u>35.4</u> | <u>31.4</u> | <u>36.4</u> | <u>36.4</u> | <u>33.9</u> |

*Code Translation:* From Table 4 and Table 5, Parallel outperforms other PEFT methods for all PTMs on both sub-tasks, Prefix is the worst on two scenarios and LoRA performs worst for the other scenarios. Comparing full fine-tuning and PEFT, the results of Java → C# translation mostly mirror those of C# → Java translation, with the difference that full fine-tuning has a better effect on CodeBERT on C# → Java translation. It could be that CodeBERT has learned Java, and requires less trainable parameters to learn the representation on Java → C# translation. For other scenarios, text-only PTMs all perform better than code-specific PTMs. As a generation task, code translation is more difficult and requires the PTMs to comprehend code representation. CodeBERT and CodeT5 need full fine-tuning to extend to unlearned C# language, while PEFT of RoBERTa and T5 can achieve comparable or even better results with full fine-tuning. It implies that PEFT is better at cross-modal transfer learning than cross-language transfer learning.

*Code Search:* From Table 6 and Table 7, Parallel performs best on all languages, while LoRA performs worst. The results of CodeBERT and CodeT5 are stable for various PEFT methods, while the results of PEFT methods on RoBERTa and T5 are quite different. The PEFT method performs worse on RoBERTa compared with CodeBERT, while the results of CodeT5 and T5 show that T5 is better when using PEFT methods. Whether a PTM is trained specifically on source code does not seem to have a significant impact on the code search task. It may be that code search is a cross-modal matching task that requires learning natural language and programming language knowledge at the same time. Other factors like the model architecture of the PTM may be more important, which we will discuss in RQ3.

**Answer to RQ1:** Among PEFT methods, Parallel always performs best on code-specific and text-only PTMs. Compared with the results of code-specific PTMs, PEFT methods applied to text-only PTMs are more effective in maintaining the performance of full fine-tuning PTMs on clone detection, defect detection, and code translation. Whether a PTM is trained specifically on source code or natural language text has no impact on the performance of code search. This shows that the PEFT methods are more advantageous in cross-modal transfer learning for PTMs.

Table 8. [RQ2] The Results of T5 and T5-large on Clone Detection, Defect Detection, and Code Translation.

| Fine-Tune Methods | Clone Detection | Defect Detection | Code Translation | |
|---|---|---|---|---|
| | BCB | Devign | Java → C# | C# → Java |
| | F1 | Acc | EM | EM |
| $T5_{Full}$ | 91.86 | 58.42 | 52 | 52.2 |
| $T5_{Best}$ | -1.24 | +2.3 | +2.3 | +3 |
| $T5_{Houlsby}$ | 88.57 | 54.06 | 37.2 | 43.8 |
| $T5_{Pfeiffer}$ | 90.2 | <u>53.99</u> | 41.3 | 45.3 |
| $T5_{Parallel}$ | **90.62** | **60.72** | **54.3** | **55.2** |
| $T5_{Prefix}$ | 82.78 | 54.06 | 33 | 37.9 |
| $T5_{LoRA}$ | <u>84.56</u> | 54.47 | <u>31.2</u> | <u>33.5</u> |
| $T5-large_{Full}$ | 83.3 | 54.98 | 54.5 | 54.7 |
| $T5-large_{Best}$ | -4.86 | -0.84 | +1.3 | +2.4 |
| $T5-large_{Houlsby}$ | 77.44 | 54.03 | 45 | 46.3 |
| $T5-large_{Pfeiffer}$ | **78.44** | 54.06 | 47.5 | 50.3 |
| $T5-large_{Parallel}$ | 72.85 | 54.06 | **55.8** | **57.1** |
| $T5-large_{Prefix}$ | 71.56 | **54.14** | 34.2 | <u>37.7</u> |
| $T5-large_{LoRA}$ | <u>68.99</u> | <u>53.92</u> | <u>31.7</u> | 42.5 |

Table 9. [RQ2] The Results of CodeT5 and CodeT5-large on Clone Detection, Defect Detection and Code Translation.

| Fine-Tune Methods | Clone Detection | Defect Detection | Code Translation | |
|---|---|---|---|---|
| | BCB | Devign | Java → C# | C# → Java |
| | F1 | Acc | EM | EM |
| $CodeT5_{Full}$ | 93.93 | 63.65 | 65.6 | 67.3 |
| $CodeT5_{Best}$ | -1.2 | -6.4 | -5.4 | -1.5 |
| $CodeT5_{Houlsby}$ | 90.41 | 56.55 | 49.7 | 55 |
| $CodeT5_{Pfeiffer}$ | 92.67 | 56.37 | 52 | 57.3 |
| $CodeT5_{Parallel}$ | **92.73** | **57.25** | **60.2** | **65.8** |
| $CodeT5_{Prefix}$ | 90.81 | 55.53 | 44.3 | <u>47.8</u> |
| $CodeT5_{LoRA}$ | <u>88.28</u> | <u>54.5</u> | <u>42.8</u> | 52 |
| $CodeT5-large_{Full}$ | 90.65 | 61.27 | 64.4 | 67.1 |
| $CodeT5-large_{Best}$ | -1.97 | -2.08 | -2.2 | -1.8 |
| $CodeT5-large_{Houlsby}$ | 82.67 | 55.42 | 50.2 | 57.7 |
| $CodeT5-large_{Pfeiffer}$ | 87.68 | <u>54.72</u> | 55.7 | 60.3 |
| $CodeT5-large_{Parallel}$ | **88.68** | **59.19** | **62.2** | **65.3** |
| $CodeT5-large_{Prefix}$ | 84.72 | 55.75 | 51.2 | 53.7 |
| $CodeT5-large_{LoRA}$ | <u>80.77</u> | 55.16 | <u>45.2</u> | <u>52.2</u> |

## 4.2 RQ2: The Same PTMs with Different Sizes

*Clone Detection:* From Table 8 and Table 9, Parallel performs the best on T5, CodeT5, and CodeT5-large, and Pfeiffer performs best on T5-large. The results show that increasing trainable parameters

Table 10. [RQ2] MRR Score (%) of T5 and T5-large on Code Search.

| Fine-Tune Methods | Go | Java | JavaScript | PHP | Python | Ruby |
|---|---|---|---|---|---|---|
| $T5_{Full}$ | 86.1 | 57.5 | 45.8 | 53.2 | 56.3 | 49 |
| $T5_{Best}$ | +1.7 | +5.3 | +6.6 | +3.4 | +3.4 | +11.5 |
| $T5_{Houlsby}$ | 85.1 | 53.3 | 42.7 | 50 | 52.8 | 45.2 |
| $T5_{Pfeiffer}$ | 85.4 | 54.7 | 44.4 | 51 | 53.7 | 50.3 |
| $T5_{Parallel}$ | **87.8** | **62.8** | **52.4** | **56.6** | **59.7** | **60.5** |
| $T5_{Prefix}$ | 77 | 41.3 | 33.3 | 39.5 | 42 | 38.2 |
| $T5_{LoRA}$ | <u>72.4</u> | <u>35.4</u> | <u>31.4</u> | <u>36.4</u> | <u>36.4</u> | <u>33.9</u> |
| $T5-large_{Full}$ | 86.2 | 62.4 | 50.5 | 56.7 | 59.4 | 53.2 |
| $T5-large_{Best}$ | +0.2 | +2.2 | +4.4 | +1.3 | +2.3 | +11.7 |
| $T5-large_{Houlsby}$ | 85.7 | 59.8 | 47.7 | 54.5 | 57.2 | 51 |
| $T5-large_{Pfeiffer}$ | 86.2 | 61.2 | 51.1 | 55.5 | 58.6 | 57.9 |
| $T5-large_{Parallel}$ | **86.4** | **64.6** | **54.9** | **58** | **61.7** | **64.9** |
| $T5-large_{Prefix}$ | <u>76.3</u> | <u>44.8</u> | 36 | <u>43.5</u> | <u>45.1</u> | 38.6 |
| $T5-large_{LoRA}$ | 77.4 | 46.3 | <u>34.7</u> | 45.2 | 47 | <u>32.7</u> |

Table 11. [RQ2] MRR Score (%) of CodeT5 and CodeT5-large on Code Search.

| Fine-Tune Methods | Go | Java | JavaScript | PHP | Python | Ruby |
|---|---|---|---|---|---|---|
| $CodeT5_{Full}$ | 91 | 72.3 | 68 | 65.3 | 70.5 | 73.2 |
| $CodeT5_{Best}$ | -0.5 | -0.1 | +0.1 | -0.1 | +0.3 | +0.9 |
| $CodeT5_{Houlsby}$ | 88.5 | 68.8 | 62.4 | 62.7 | 67.5 | 67.7 |
| $CodeT5_{Pfeiffer}$ | 89.3 | 69.9 | 64.4 | 63.5 | 68.6 | 69.4 |
| $CodeT5_{Parallel}$ | **90.5** | **72.2** | **68.1** | **65.2** | **70.8** | **74.1** |
| $CodeT5_{Prefix}$ | <u>87.4</u> | <u>62.2</u> | <u>54.9</u> | <u>58.1</u> | <u>60.7</u> | <u>63.3</u> |
| $CodeT5_{LoRA}$ | 87.5 | 67.4 | 58.2 | 61.4 | 65.9 | 63.7 |
| $CodeT5-large_{Full}$ | 90.2 | 71.6 | 67.4 | 64.6 | 70.6 | 73.8 |
| $CodeT5-large_{Best}$ | -0.7 | -0.2 | -0.2 | -0.4 | -0.4 | -0.7 |
| $CodeT5-large_{Houlsby}$ | 88.1 | 68.6 | 62.8 | 62.5 | 68.1 | 67.3 |
| $CodeT5-large_{Pfeiffer}$ | 88.5 | 69.6 | 64.4 | 63.2 | 68.8 | 69.8 |
| $CodeT5-large_{Parallel}$ | **89.5** | **71.4** | **67.2** | **64.2** | **70.2** | **73.1** |
| $CodeT5-large_{Prefix}$ | <u>86.9</u> | <u>64.2</u> | <u>57.1</u> | <u>59.7</u> | <u>63.6</u> | 64.5 |
| $CodeT5-large_{LoRA}$ | 87.3 | 67.3 | 60.6 | 61.7 | 66.6 | <u>60.7</u> |

of T5 and CodeT5 will reduce the performance of PEFT methods. Among PTMs, the performance on T5-large is the worst compared with full fine-tuning. The difference between the best PEFT and full fine-tuning is almost the same between CodeT5 and CodeT5-large, while the performance of T5-large drops a lot compared to T5. Therefore, size has a greater negative impact on T5 compared to CodeT5.

*Defect Detection:* From Table 8 and Table 9, Parallel can exceed the performance of full fine-tuning on T5 by up to 2.3%. Parallel also ranks first on CodeT5 and CodeT5-large, and Prefix is the best method on T5-large. The PEFT method on CodeT5-large brings improvements over CodeT5, but

Table 12. [RQ3] The Results of code-specific PTMs with Different Architectures on Clone Detection, Defect Detection and Code Translation.

| Fine-Tune Methods | Clone Detection | Defect Detection | Code Translation | |
|---|---|---|---|---|
| | BCB | Devign | Java → C# | C# → Java |
| | F1 | Acc | EM | EM |
| $CodeBERT_{Full}$ | 94.95 | 63.91 | 59 | 60.1 |
| $CodeBERT_{Best}$ | -1.88 | -3.08 | -1.7 | +0.3 |
| $CodeBERT_{Houlsby}$ | 91.9 | 57.87 | 56.7 | 59.1 |
| $CodeBERT_{Pfeiffer}$ | 91.62 | 56.41 | 56.4 | 58.3 |
| $CodeBERT_{Parallel}$ | **93.07** | **60.83** | **57.3** | **60.4** |
| $CodeBERT_{Prefix}$ | 89.96 | 56.3 | 55.2 | 59 |
| $CodeBERT_{LoRA}$ | <u>88.38</u> | <u>55.6</u> | <u>54.4</u> | <u>57.2</u> |
| $UniXcoder_{Full}$ | 94.68 | 64.28 | 58 | 58 |
| $UniXcoder_{Best}$ | -1.06 | -3.67 | +2.5 | +5.3 |
| $UniXcoder_{Houlsby}$ | 92.35 | 59.15 | <u>59</u> | **63.3** |
| $UniXcoder_{Pfeiffer}$ | 92.41 | 58.6 | 59.7 | 60.8 |
| $UniXcoder_{Parallel}$ | 93.1 | **60.61** | **60.5** | 62.6 |
| $UniXcoder_{Prefix}$ | **93.62** | <u>58.05</u> | 60.4 | <u>60.5</u> |
| $UniXcoder_{LoRA}$ | <u>90.81</u> | 58.53 | 59.7 | 61.3 |
| $CodeT5_{Full}$ | 93.93 | 63.65 | 65.6 | 67.3 |
| $CodeT5_{Best}$ | -1.2 | -6.4 | -5.4 | -1.5 |
| $CodeT5_{Houlsby}$ | 90.41 | 56.55 | 49.7 | 55 |
| $CodeT5_{Pfeiffer}$ | 92.67 | 56.37 | 52 | 57.3 |
| $CodeT5_{Parallel}$ | **92.73** | **57.25** | **60.2** | **65.8** |
| $CodeT5_{Prefix}$ | 90.81 | 55.53 | 44.3 | <u>47.8</u> |
| $CodeT5_{LoRA}$ | <u>88.28</u> | <u>54.5</u> | <u>42.8</u> | 52 |

the performance on T5-large is worse than T5. It shows that increasing trainable parameters hurts T5 and has a positive impact on CodeT5. Besides, T5 can perform better with PEFT than CodeT5, but T5-large underperforms CodeT5-large with PEFT. This shows that code-related pre-training tasks can help PTM perform better with more trainable parameters.

*Code Translation:* From Table 8 and Table 9, for all the scenarios of two sub-tasks and four PTMs, Parallel is the best PEFT method. Prefix is the worst on C# → Java for T5-large and CodeT5, and LoRA achieves the worst for other scenarios. Parallel outperforms the full fine-tuning method for T5 and T5-large and provides comparable results for CodeT5 and CodeT5-large. Increasing the trainable parameters only improves CodeT5 on Java → C# translation, while resulting in negative improvements on C# → Java translation. Since the improvement is larger, we believe that increasing trainable parameters is better for CodeT5 with PEFT. On the other hand, increasing the trainable parameters brings negative improvements to T5 on both sub-tasks.

*Code Search:* From Table 10 and Table 11, Parallel achieves the best results on all languages among the four PTMs, while LoRA is the worst method for T5 and T5-large and Prefix performs worst for CodeT5 and CodeT5-large. For T5 and T5-large, the best PEFT method outperforms the full fine-tuning method, and for some languages, the improvement is large, such as Ruby. Furthermore, the PEFT method provides comparable or even better results to the full fine-tuning method for CodeT5 and CodeT5-large. Compared with CodeT5-large, the PEFT method performs better on

Table 13. [RQ3] The Results of text-only PTMs with Different Architectures on Clone Detection, Detection, and Code Translation.

| Fine-Tune Methods | Clone Detection | Defect Detection | Code Translation | |
|---|---|---|---|---|
| | BCB | Devign | Java → C# | C# → Java |
| | F1 | Acc | EM | EM |
| $RoBERTa_{Full}$ | 93.5 | 56.62 | 58.8 | 60.8 |
| $RoBERTa_{Best}$ | -0.15 | +2.41 | -0.2 | -0.2 |
| $RoBERTa_{Houlsby}$ | 90.19 | <u>54.03</u> | 56 | 59.4 |
| $RoBERTa_{Pfeiffer}$ | 90.24 | <u>54.03</u> | 57.4 | 58.8 |
| $RoBERTa_{Parallel}$ | **93.35** | **58.93** | **58.6** | **60.6** |
| $RoBERTa_{Prefix}$ | 89.63 | 55.53 | 56.1 | <u>58.3</u> |
| $RoBERTa_{LoRA}$ | <u>87.37</u> | 54.06 | <u>55</u> | 58.4 |
| $BART_{Full}$ | 92.72 | 61.49 | 48.8 | 53.2 |
| $BART_{Best}$ | -0.32 | -0.62 | -1.8 | -1.5 |
| $BART_{Houlsby}$ | 88.47 | 56.7 | 41.4 | 47.5 |
| $BART_{Pfeiffer}$ | 88.47 | 55.75 | 41.1 | 47.7 |
| $BART_{Parallel}$ | **92.4** | **60.87** | **47** | **51.7** |
| $BART_{Prefix}$ | 89.11 | <u>55.49</u> | 42.1 | 48.9 |
| $BART_{LoRA}$ | <u>87.19</u> | 55.89 | <u>31.9</u> | <u>39.4</u> |
| $T5_{Full}$ | 91.86 | 58.42 | 52 | 52.2 |
| $T5_{Best}$ | -1.24 | +2.3 | +2.3 | +3 |
| $T5_{Houlsby}$ | 88.57 | 54.06 | 37.2 | 43.8 |
| $T5_{Pfeiffer}$ | 90.2 | <u>53.99</u> | 41.3 | 45.3 |
| $T5_{Parallel}$ | **90.62** | **60.72** | **54.3** | **55.2** |
| $T5_{Prefix}$ | 82.78 | 54.06 | 33 | 37.9 |
| $T5_{LoRA}$ | <u>84.56</u> | 54.47 | <u>31.2</u> | <u>33.5</u> |

CodeT5, especially outperforming the full fine-tuning method in JavaScript, Python, and Ruby. As the trainable parameters increase, the PEFT method also has a negative improvement on T5, but the gap is smaller compared with CodeT5.

**Answer to RQ2:** Among PEFT methods, Parallel always performs best on different model sizes. The results of CodeT5 with different sizes are more stable than those of T5. Increasing the trainable parameters of T5 reduces the performance of the PEFT method compared to full fine-tuning. CodeT5 with more trainable parameters performs better on defect detection and code translation and performs worse on clone detection and code search.

### 4.3   RQ3: PTMs in Different Architectures

*Clone Detection:* From Table 12, Parallel performs best on CodeBERT and CodeT5, Prefix is the best on UniXcoder, while LoRA is the worst method on the three code-specific PTMs. The best PTM to maintain full fine-tuning performance when using the PEFT method is UniXcoder. It shows that a cross-modal encoder-only architecture with code, AST, and comments is helpful for BERT-based encoder-only architecture. From Table 13, Parallel performs best while LoRA is the worst method on the three text-only PTMs. Similar to code-specific PTMs, RoBERTa, BERT-based encoder-only architecture, performs better with PEFT methods than BART and T5 of encoder-decoder architecture.

Table 14. [RQ3] MRR Score (%) of code-specific PTMs with Different Architectures on Code Search.

| Fine-Tune Methods | Go | Java | JavaScript | PHP | Python | Ruby |
|---|---|---|---|---|---|---|
| $CodeBERT_{Full}$ | 90.4 | 70.5 | 63.3 | 64.2 | 69 | 69.6 |
| $CodeBERT_{Best}$ | -0.5 | -0.7 | -0.9 | -0.3 | -0.5 | -1.1 |
| $CodeBERT_{Houlsby}$ | 86.7 | 64.4 | 49.6 | 59.7 | 63.3 | 51.6 |
| $CodeBERT_{Pfeiffer}$ | 86.6 | 63.9 | 49.9 | 59.3 | 62.8 | 51.4 |
| $CodeBERT_{Parallel}$ | **89.9** | **69.8** | **62.4** | **63.9** | **68.5** | **68.5** |
| $CodeBERT_{Prefix}$ | 87.4 | 65.2 | 54.9 | 59.4 | 63 | 56.2 |
| $CodeBERT_{LoRA}$ | <u>84</u> | <u>55.8</u> | <u>46.5</u> | <u>53.9</u> | <u>57.5</u> | <u>41.2</u> |
| $UniXcoder_{Full}$ | 91.2 | 72.3 | 59.2 | 67.1 | 73.2 | 60.9 |
| $UniXcoder_{Best}$ | -4 | -4.9 | -5.8 | -3.7 | -4.1 | -0.6 |
| $UniXcoder_{Houlsby}$ | 73.9 | 57.3 | 47.7 | 59 | 59.3 | 55.7 |
| $UniXcoder_{Pfeiffer}$ | 72.6 | 56.8 | 46.4 | 44 | 53.4 | 54.6 |
| $UniXcoder_{Parallel}$ | 78.3 | 61.5 | **53.4** | 56.1 | 59.8 | **60.3** |
| $UniXcoder_{Prefix}$ | **87.2** | **67.4** | 48.4 | **63.4** | **69.1** | 55.6 |
| $UniXcoder_{LoRA}$ | <u>69.9</u> | <u>54.4</u> | <u>44</u> | <u>42.4</u> | <u>51.6</u> | <u>50</u> |
| $CodeT5_{Full}$ | 91 | 72.3 | 68 | 65.3 | 70.5 | 73.2 |
| $CodeT5_{Best}$ | -0.5 | -0.1 | +0.1 | -0.1 | +0.3 | +0.9 |
| $CodeT5_{Houlsby}$ | 88.5 | 68.8 | 62.4 | 62.7 | 67.5 | 67.7 |
| $CodeT5_{Pfeiffer}$ | 89.3 | 69.9 | 64.4 | 63.5 | 68.6 | 69.4 |
| $CodeT5_{Parallel}$ | **90.5** | **72.2** | **68.1** | **65.2** | **70.8** | **74.1** |
| $CodeT5_{Prefix}$ | <u>87.4</u> | <u>62.2</u> | <u>54.9</u> | <u>58.1</u> | <u>60.7</u> | <u>63.3</u> |
| $CodeT5_{LoRA}$ | 87.5 | 67.4 | 58.2 | 61.4 | 65.9 | 63.7 |

Thus, on clone detection, no matter whether code-specific or text-only PTMs, encoder-decoder architecture performs worse with PEFT methods than encoder-decoder architecture.

*Defect Detection:* From Table 12, Parallel performs best on the three code-specific PTMs, Prefix performs worst on UniXcoder, while LoRA is the worst on CodeBERT and CodeT5. Among all the code-specific PTMs, the PEFT method performs best on CodeBERT compared with full fine-tuning. From Table 13, Parallel performs best on the three text-only PTMs while Houlsby, Pfeiffer, and Prefix are the worst methods. Furthermore, except for Parallel, the results of other methods are close, and there is a gap with full fine-tuning. Among text-only PTMs, BART performs worst and PEFT brings negative improvements while having positive improvements on both RoBERTa and T5. Compared with clone detection, the PEFT methods have worse results in code-specific PTMs, especially on CodeT5, but perform better on text-only PTMs, such as RoBERTa and T5. The results of code-specific and text-only PTMs show that the encoder-only architecture also performs better with PEFT methods than the encoder-decoder architecture on defect detection.

*Code Translation:* From Table 12, Parallel performs the best in most scenarios, while LoRA, Houlsby and Prefix have worst-performing scenarios. Compared with full fine-tuning, PEFT methods perform best on UniXcoder, followed by CodeBERT and CodeT5. For UniXcoder, PEFT methods outperform full fine-tuning on both sub-tasks, and the value of the best PEFT methods of UniXcoder is comparable with that of CodeT5. From Table 13, Parallel performs the best in all scenarios and LoRA is the worst in most scenarios. PEFT methods perform best on T5 and can perform better than full fine-tuning. However, the performance of PEFT methods on BART is the worst, which shows that PEFT can not perform well on BART-based encoder-decoder architecture. Comparing

Table 15. [RQ3] MRR Score (%) of text-only PTMs with Different Architectures on Code Search.

| Fine-Tune Methods | Go | Java | JavaScript | PHP | Python | Ruby |
|---|---|---|---|---|---|---|
| $RoBERTa_{Full}$ | 88.5 | 64.5 | 54.6 | 59.8 | 62.2 | 60.6 |
| $RoBERTa_{Best}$ | -0.9 | -3.7 | -3.8 | -3.3 | -3.8 | -3.4 |
| $RoBERTa_{Houlsby}$ | 81.9 | 49.3 | 36.1 | 47.1 | 48 | 35.2 |
| $RoBERTa_{Pfeiffer}$ | 82.3 | 49.7 | 37.9 | 47.1 | 48.5 | 41.2 |
| $RoBERTa_{Parallel}$ | **87.6** | **60.8** | **50.8** | **56.5** | **58.4** | **57.2** |
| $RoBERTa_{Prefix}$ | 83.6 | 50.7 | 39.5 | 48.1 | 48.2 | 45.4 |
| $RoBERTa_{LoRA}$ | <u>78.5</u> | <u>43.3</u> | <u>29.4</u> | <u>42.2</u> | <u>43.3</u> | <u>33.8</u> |
| $BART_{Full}$ | 88.3 | 63.7 | 53.4 | 57.8 | 60.8 | 61 |
| $BART_{Best}$ | -0.8 | -4.8 | -3.9 | -3.8 | -4.3 | -4.1 |
| $BART_{Houlsby}$ | 84.1 | 50.6 | 40.9 | 46.5 | 48.3 | 46.8 |
| $BART_{Pfeiffer}$ | 84.1 | 50.1 | 40.4 | 45.7 | 47.7 | 46.8 |
| $BART_{Parallel}$ | **87.5** | **58.9** | **49.5** | **54** | **56.5** | **56.9** |
| $BART_{Prefix}$ | 85.7 | 52.7 | 44.9 | 48.8 | 50.5 | 52.7 |
| $BART_{LoRA}$ | <u>81.8</u> | <u>45.9</u> | <u>37</u> | <u>42.3</u> | <u>42.3</u> | <u>41.1</u> |
| $T5_{Full}$ | 86.1 | 57.5 | 45.8 | 53.2 | 56.3 | 49 |
| $T5_{Best}$ | +1.7 | +5.3 | +6.6 | +3.4 | +3.4 | +11.5 |
| $T5_{Houlsby}$ | 85.1 | 53.3 | 42.7 | 50 | 52.8 | 45.2 |
| $T5_{Pfeiffer}$ | 85.4 | 54.7 | 44.4 | 51 | 53.7 | 50.3 |
| $T5_{Parallel}$ | **87.8** | **62.8** | **52.4** | **56.6** | **59.7** | **60.5** |
| $T5_{Prefix}$ | 77 | 41.3 | 33.3 | 39.5 | 42 | 38.2 |
| $T5_{LoRA}$ | <u>72.4</u> | <u>35.4</u> | <u>31.4</u> | <u>36.4</u> | <u>36.4</u> | <u>33.9</u> |

the results of code-specific and text-only PTMs, text-only PTMs are more stable when using PEFT methods and the impact of model architecture on the performance of PEFT methods is small.

*Code Search:* From Table 14, Parallel performs best on CodeBERT and CodeT5 for all languages. LoRA is the worst method on CodeBERT and UniXcoder and Prefix performs worst on CodeT5 for all languages. For UniXcoder, there is a gap between the best PEFT method and the full fine-tuning method, with a difference of about 5%. Prefix outperforms Parallel on Go, Java, PHP, and Python languages, but underperforms on JavaScript and Ruby. CodeBERT and UniXcoder perform differently, which shows that unified cross-modal architecture hurts PEFT methods. PEFT methods on CodeT5 can achieve comparable or even better results with full fine-tuning. From Table 15, Parallel performs best and LoRA is the worst on all text-only PTMs for all languages. Full fine-tuning is better than the PEFT method on RoBERTa and BART, and the gap is similar for these two PTMs in all languages. Furthermore, the best PEFT method on T5 is better than full fine-tuning. Comparing the results of code-specific and text-only PTMs, T5-based encoder-decoder architecture is better for PEFT methods on code search.

Besides, the results of PEFT methods vary greatly for different languages. For example, the PTMs on the Go language can maintain good performance regardless of the architecture or PEFT method adopted. For CodeBERT, the results of PEFT methods except Parallel on JavaScript and Ruby languages are much worse than the results of full fine-tuning. This will help us identify which components of language models are more important in designing effective PEFT for different languages.

Table 16. The Cases that PEFT Methods Perform Better than Full Fine-Tuning.

| Models | Clone Detection | Defect Detection | Code Search | | | | | | Code Translation | |
|---|---|---|---|---|---|---|---|---|---|---|
| | BCB | Devign | Go | Java | JS | PHP | Py | Ruby | Java → C# | C# → Java |
| CodeBERT | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↑ |
| UniXcoder | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↑ | ↑ |
| RoBERTa | ↓ | ↑ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ |
| BART | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ |
| T5 | ↓ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ |
| T5-large | ↓ | ↓ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ | ↑ |
| CodeT5 | ↓ | ↓ | ↓ | ↓ | ↑ | ↓ | ↑ | ↑ | ↓ | ↓ |
| CodeT5-large | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ |

Table 17. Frequency of Achieving the Best or Worst Results on different tasks for PEFT methods.

| Fine-Tune Methods | Clone Detection | Defect Detection | Code Search | | | | | | Code Translation | |
|---|---|---|---|---|---|---|---|---|---|---|
| | BCB | Devign | Go | Java | JS | PHP | Py | Ruby | Java → C# | C# → Java |
| Houlsby | 0 | 2 - | 0 | 0 | 0 | 0 | 0 | 0 | 1 - | 1 + |
| Pfeiffer | 1 + | 3 - | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Parallel | 6 + | 7 + | 7 + | 7 + | 8 + | 7 + | 7 + | 8 + | 8 + | 7 + |
| Prefix | 1 + | 1 + / 2 - | 1 + / 3 - | 1 + / 3 - | 2 - | 1 + / 3 - | 1 + / 3 - | 1 - | 0 | 4 - |
| LoRA | 8 - | 2 - | 5 - | 5 - | 6 - | 5 - | 5 - | 7 - | 7 - | 4 - |

**Answer to RQ3:** Among PEFT methods, Parallel always performs best on PTMs in different architectures. PEFT methods applied to encoder-only architecture perform better on clone detection and defect detection than encoder-decoder architecture, while those used in T5-based encoder-decoder architecture are better on code search. Architecture has no obvious impact on the performance of PEFT methods on code translation.

## 5 DISCUSSION

### 5.1 Effectiveness: Qualitative Analysis on Various Tasks

Table 16 illustrates whether PEFT methods are better than full fine-tuning methods, where an up arrow ↑ indicates PEFT method outperforms full fine-tuning method, while a down arrow ↓ means full fine-tuning method performs better. In the table, although full fine-tuning methods are generally superior, there are specific cases where PEFT methods perform better. For classification tasks (clone detection and defect detection), full fine-tuning methods are better than PEFT methods in almost all cases. As these classification tasks are simpler and do not require the model to deeply understand the code, more parameters and longer training time usually achieve better results. For the retrieval task (code search), PEFT methods can achieve better results, especially for T5-based PTMs. For the generation task (code translation), PEFT methods can perform better on T5-based PTMs, CodeBERT, and UniXcoder, than full fine-tuning methods. These results underline the potential of PEFT methods to generate efficient and effective PTMs on code search and code translation.

To provide a more intuitive understanding of various PEFT methods, we rank them based on the frequency of achieving the best or worst results on various tasks in Table 17, where "+" indicates the method that achieves the best result, "−" means the method that achieves the worst result, and the corresponding number represents the number of times the best/worst result is achieved. In the table, the Parallel method performs better while LoRA performs the worst when applied to

Table 18.  Trainable/Total Parameters (M) of PTMs.

| Fine-Tune Methods | RoBERTa/CodeBERT | UniXcoder | BART | T5/CodeT5 | T5-large/CodeT5-large |
|---|---|---|---|---|---|
| Full | 125/125 | 126/126 | 139/139 | 223/223 | 738/738 |
| Houlsby | 2/127 (1.6%) | 2/128 (1.6%) | 2/142 (1.5%) | 4/227 (1.7%) | 13/751 (1.8%) |
| Pfeiffer | 1/126 (0.7%) | 1/127 (0.7%) | 1/140 (0.6%) | 2/225 (0.8%) | 6/744 (0.9%) |
| Parallel | 7/132 (5.4%) | 7/133 (5.3%) | 7/147 (4.8%) | 14/237 (6.0%) | 50/788 (6.4%) |
| Prefix | 10/135 (7.3%) | 10/136 (7.3%) | 15/155 (10.0%) | 30/253 (11.7%) | 77/815 (9.5%) |
| LoRA | 0.3/125 (0.2%) | 0.3/126 (0.2%) | 0.4/140 (0.3%) | 1/224 (0.4%) | 2/740 (0.3%) |

fine-tune PTMs. A PEFT method can perform better on one task but perform worse on another task. For example, Houlsby is the best method on Java → C# translation task but the worst on C# → Java translation, although both are translation tasks.

**Conclusion:** PEFT methods can perform comparable, or even better, results with corresponding full fine-tuning methods. They perform better on code search and code translation tasks. Among PEFT methods, the Parallel is the best method while LoRA performs worst.

## 5.2  Efficiency: Required GPU Resources and Training Time

PEFT methods can tune fewer parameters than full fine-tuning. To give a more intuitive comparison, we first list the trainable/total parameters of all the PTMs with various fine-tuning methods in Table 18. For encoder-only PTMs on code translation tasks, we add extra decoder blocks to generate code. Therefore, 48M additional parameters need to be added to their trainable/total parameters. For example, the trainable/total parameters of Houlsby on CodeBERT is 50/175M.

From Table 18, the trainable parameters of PEFT methods are much fewer than those in the original models, but loading the models occupies GPU resources. The number of trainable parameters does not directly equate to GPU resources or its proportional occupation. Thus, we compare the required GPU resources and training time of all the PTMs with various fine-tuning methods on the four tasks. For GPU resources, we take the peak GPU occupancy of each PTM as its required GPU resources. The training time of a PTM is defined as the duration from the start of its training to its end, except for the time of data pre-processing and model evaluation. During training, we make efforts to ensure no interference from other processes on the server. Since PTMs with different parameters and architectures require inconsistent resources, to facilitate comparison, we display them uniformly as a percentage number compared to the Base. The percentage of required GPU resources and training time of four tasks are shown in Figure 5-8, corresponding to the results for each task respectively. In the figures, the "GPU" means required GPU resources, and "Time" is the actual training time. For code search, since we set the same parameters for all languages, the same method should have the same proportion of GPU resources and training time in each language. Therefore, we only illustrate the results for JavaScript language. Similarly, for code translation, we only illustrate the results of Java to C# translation in the figures.

From the results, PEFT methods can save GPU resources compared to full fine-tuning methods, which is their primary design purpose. The used GPU resources are related to the specific implementation of PTMs and the design of PEFT methods. Generally, if a PEFT method has more trainable parameters or modifies the original model more, it takes more GPU resources. On average, the PEFT methods can save about 10%-30% of actual GPU resources compared to full fine-tuning. Among PEFT methods, Houlsby and Prefix take the most required GPU resources while Pfeiffer needs the least resources. On code translation tasks, LoRA also takes more GPU resources than other methods. For all the tasks, as the trainable parameters increase, fewer GPU resources are required for PEFT methods.
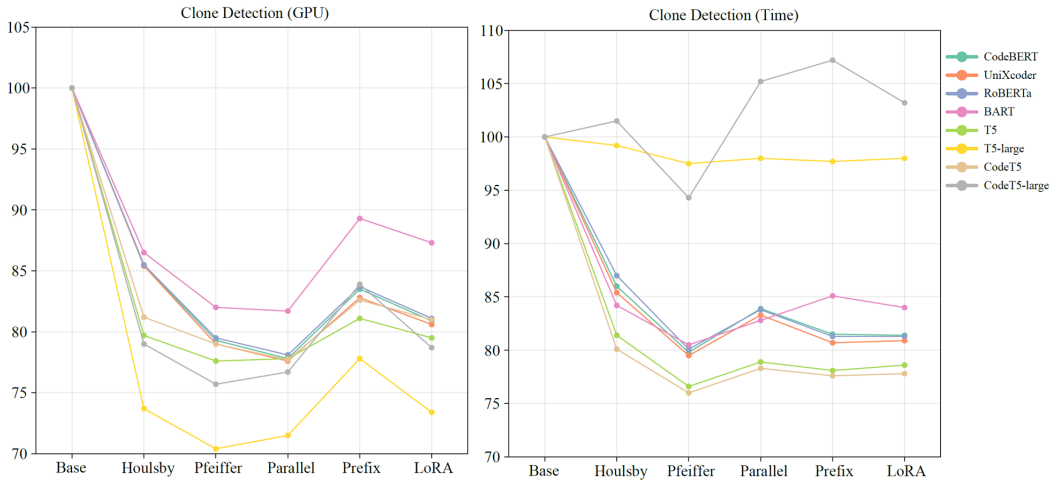
Fig. 5. GPU (%) Resources and Training Time (%) Used on Clone Detection Task.
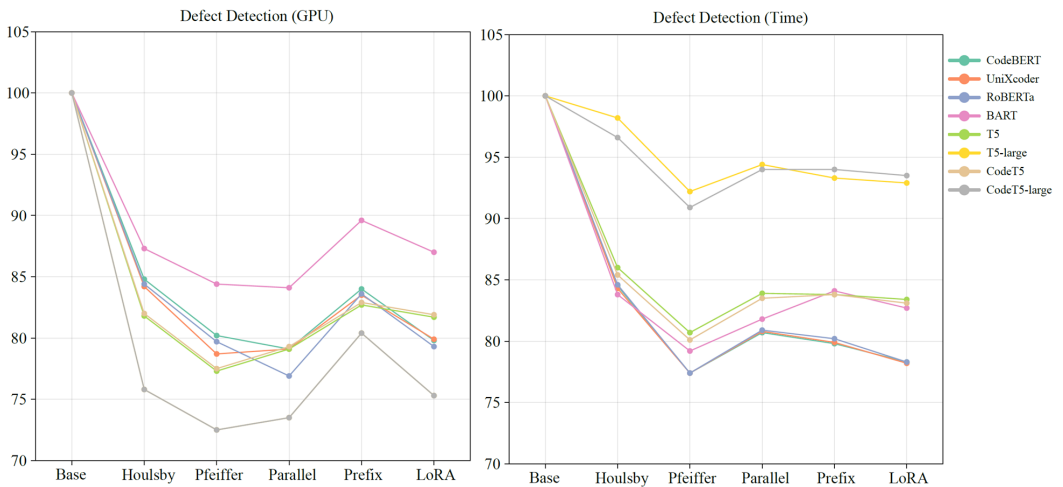


Fig. 6. GPU (%) Resources and Training Time (%) Used on Defect Detection Task.

Although PEFT methods show clear reductions in required GPU resources, the training time of some PEFT methods is longer than the full fine-tuning method, and the results of training time are not as stable as GPU resources. For clone detection in Figure 5, the training time of PEFT methods can exceed that of the full fine-tuning method, such as using Prefix on CodeT5-large. For defect detection in Figure 6, all the PEFT methods are faster than full fine-tuning methods. For code search in Figure 7, although the GPU resources show clear patterns, the training time does not follow the same clarity. A PEFT method that is fast for one PTM may be slow for another PTM. We believe that increasing the training speed is not what PEFT promises, and the complex environment (a large amount of time is spent on searching for similar code snippets) makes the training time more chaotic. The code generation task, code translation, as shown in Figure 8, takes more time than other tasks because it requires generating code sequences other than single values. As the
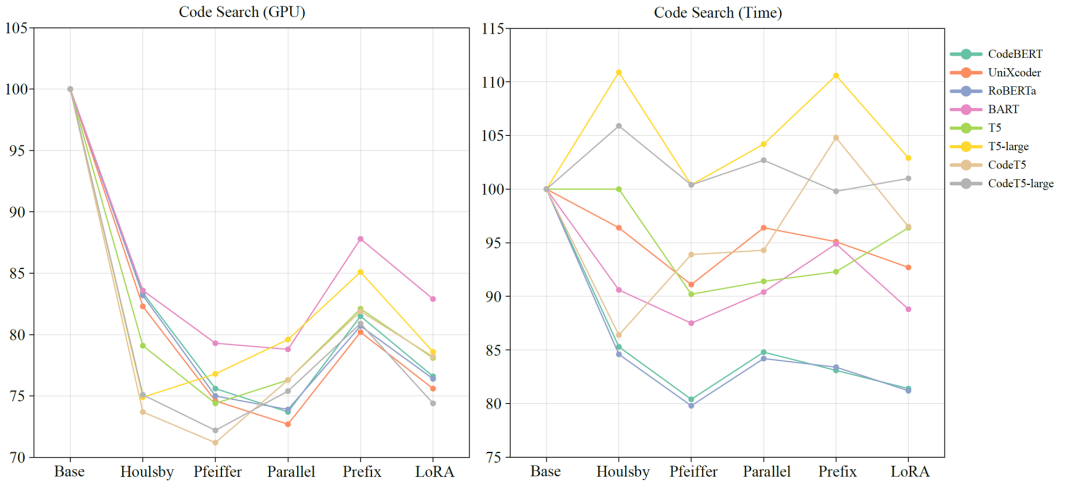
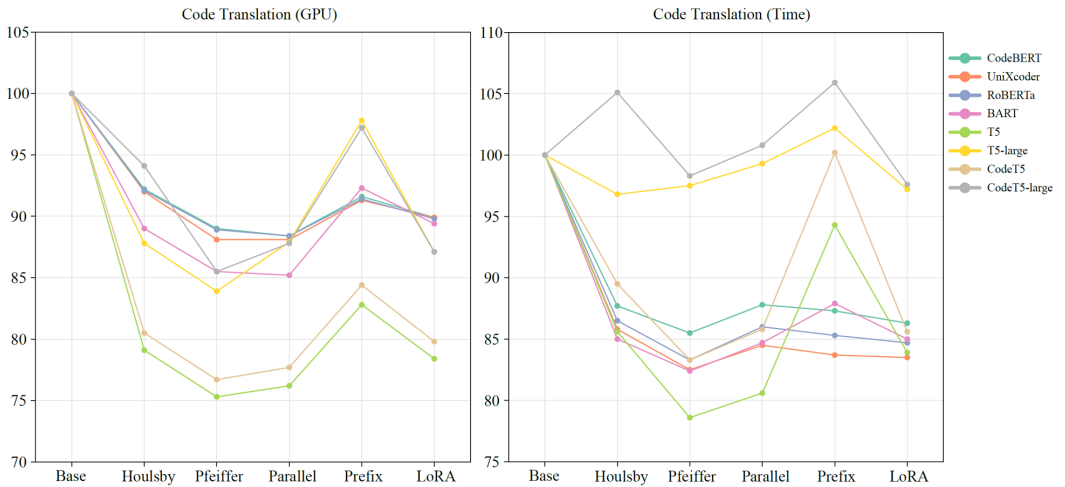Fig. 7. GPU (%) Resources and Training Time (%) Used on JavaScript Code Search Task.



Fig. 8. GPU (%) Resources and Training Time (%) Used on Java to C# Code Translation Task.

model grows larger, the training time of some methods, such as Prefix for T5-large and Houlsby for CodeT5-large, exceeds full fine-tuning. For all the tasks, PEFT does hold the potential to decrease the training time, and as the trainable parameters increase, more training time is required for PEFT methods.

**Conclusion:** PEFT methods can save about 10%-30% GPU resources compared with full fine-tuning methods. Among PEFT methods, Pfeiffer costs the least GPU resources while Prefix and Houlsby cost the most GPU resources. The increase in training speed is not what PEFT promises to improve and the training time of the PEFT method is sometimes longer than that of full fine-tuning, but it still has the potential to reduce training time. As the trainable parameters increase, fewer GPU resources and more training time are required.

## 5.3 Threats to Validity

**Construct Validity**. As discussed in implementation details of Section 3, we reuse code from open source repositories on encoder-only models and adapt them to encoder-decoder models. While the modifications made might not result in ideal outcomes, we have tried our best to ensure the models reach their full potential.

**Internal Validity**. As a design choice, we do not perform any hyperparameter tuning, using settings from CodeXGLUE and Table 3. Our results are based on this set of parameters, and it is not guaranteed that the same results would be obtained under different parameters. However, for any PTM and task, the hyperparameter configurations for all PEFT methods are consistent with the corresponding full fine-tuning method, ensuring that the experimental results are obtained under fair conditions.

We are unable to deploy larger language models like CodeLlama [58] and CodeGen [46] due to the limitations of available resources. Smaller decoder-only models such as CodeGPT do not perform as well as other models [48]. Thus, our research may not provide a comprehensive view of all possible model architectures. We expect no impact of this choice on our findings.

**External Validity**. The results and findings in this work may be applicable to the specific PEFT methods, PTMs, SE downstream tasks, and corresponding datasets that we have evaluated. For other PEFT methods, PTMs, and SE downstream tasks, the same results and findings cannot be guaranteed.

## 6 RELATED WORK

### 6.1 Code-specific Pre-Trained Models in the Software Engineering Field

Due to the success of Transformer architecture, a large number of PTMs have been widely used in the NLP field [13, 31, 38, 53, 55, 68, 85]. Based on these foundation models [6], many code-specific PTMs have been proposed in the SE field. These PTMs leverage programming language datasets for pre-training and deal with code-related downstream tasks.

CuBERT [29] and CodeBERT [15] are among the first BERT-based encoder-only models in the SE field. CuBERT employs the masked language modeling (MLM) objective, while CodeBERT proposes a new replaced token detection (RTD) objective. CuBERT is pre-trained on Python language and CodeBERT is pre-trained on six programming languages on the CodeSearchNet corpus. Later, GraphCodeBERT [20] adds data flow information to CodeBERT and focuses on the inherent semantic-level structure of code. To make full use of data flow, it introduces two structure-aware pre-training tasks. One is to predict code structure edges, and the other is to align representations between source code and code structure. To support both code-related understanding and generation tasks, UniXcoder [19], a unified cross-modal framework, utilizes mask attention matrices with prefix adapters to control the behavior. It is pre-trained on multi-modal data, including code, comment, and AST, and can leverage these to enhance code representation. Furthermore, it utilizes a contrastive learning task to utilize multi-modal representation and align representations among programming languages using a cross-modal generation task.

Decoder-only models [7] in the SE field are always based on GPT2 [54] that utilize unidirectional language modeling objectives to predict the next token using all previous tokens. GPT-C [71] is pre-trained in Python, C#, JavaScript, and TypeScript programming languages to complete multilingual code. It can predict sequences of code tokens of arbitrary types, and generate up to entire lines of syntactically correct code. CodeGPT [40] is pre-trained in Python and Java languages to support the code completion and text-to-code generation tasks.

Encoder-decoder models are also explored in the SE field. Based on BART, PLBART [2] is pre-trained on an extensive collection of Java and Python functions and associated NL text via denoising

autoencoding objectives, such as token masking, token deletion, and token infilling. Learning from T5, CodeT5 [80] employs a unified framework to support various code-related tasks and allows for multi-task learning [91]. Besides, it proposes a novel identifier-aware pre-training task to leverage information of code identifiers, and a bimodal dual generation task to align code and comments. SPT-Code[49] is a sequence-to-sequence pre-trained model that uses three specially designed pre-training tasks to learn the knowledge of code, code structure, and natural language description of code.

## 6.2 Parameter-Efficient Fine-Tuning Methods

There have been over 40 different PEFT methods in various fields [14, 33, 56]. Following the pioneering work of adapter tuning [25], Pfeiffer et al. [52] modified adapter tuning and placed an adapter layer only after the feed-forward layer while He et al. [23] placed adapter layers in parallel to the original model. Inspired by the success of prompting methods [36], some research focuses on additional prefix vectors. Li and Liang [32] used prefix tuning to add prefix vectors before input vectors and only updated prefix vectors when fine-tuning. Later, Lester et al. [30] proposed prompt tuning that simplifies prefix tuning by only adding prompt vectors to the input sequences in the first layer and Liu et al. [37] used a similar idea in the P-tuning method that employed trainable continuous prompt embeddings in concatenation with discrete prompts to increase stability. Hu et al. [26] proposed LoRA to inject trainable low-rank decomposition matrices into original models. To bridge this gap between the performance of PEFT and full fine-tuning, Zhang et al. [89] proposed AdaLoRA, which adaptively allocates the parameter budget among weight matrices according to their importance score. Liu et al. [34] used $(IA)^3$ to scale inner model activations to fine-tune new parameters. BitFit [86] is a sparse fine-tuning method and only fine-tunes bias vectors. Guo et al. [21] proposed diff pruning and treated fine-tuning as learning a task-specific diff vector that is applied on top of the pre-trained shared parameter vector.

In the SE field, various PEFT methods have been evaluated in recent years. Ayupov and Chirkova [5] applied PEFT methods and tested adapter tuning and LoRA on four code-related tasks, including clone detection, code summarization, code generation, and code translation. They found that PEFT can achieve comparable or higher performance than full fine-tuning on code understanding tasks, and underperform full fine-tuning on code generation tasks. Later, Goel et al. [17] applied adapters for cross-modal transfer to improve RoBERTa, achieving comparable performance to CodeBERT in some code-related cases. To alleviate the potentially catastrophic forgetting issue in multilingual models, Wang et al. [78] inserted the adapters to PTMs for better results on code search and summarization tasks. They demonstrated that adapter tuning outperforms full fine-tuning in cross-lingual, multilingual, and low-resource scenarios, and adapter tuning can overcome catastrophic forgetting significantly. To fully utilize the rich syntactical information in source code, Saberi and Fard [59] proposed Named Entity Recognition (NER) adapters that can be inserted into PTMs to learn type information extracted from the AST. They inserted NER adapters in CodeBERT and they proved that adapters improve full fine-tuning on code refinement and code summarization tasks.

## 6.3 Empirical Studies on Pre-Trained Models in the Software Engineering Field

Several works have done emnlppirical studies to focus on the performance of PTMs on SE downstream tasks. Chirkova and Troshin [12] conducted studies of the capabilities of Transformers to utilize syntactic information in different tasks. Zhou et al. [94] showed the generalizability of CodeBERT to achieve higher or comparable performance than specialized solutions designed for the code search and just-in-time defect prediction tasks. Ahmed and Devanbu [3] found that identifiers are a very important element of training data for software engineering tasks and available multilingual training data (across different languages) can be used to amplify the performance of

PTMs. To better understand how monolingual and multilingual PTMs affect different programming languages, Chen et al. [10] conducted and analyzed over a hundred models. Steenhoek et al. [66] reproduced 9 DNN models on two vulnerability detection datasets to have a good understanding of these models. Mastropaolo et al. [42] let GitHub Copilot automatically generate Java methods based on different semantically equivalent descriptions to evaluate the robustness of the code completion. Recently, Tufano et al. [75] focused on the impact of pre-training tasks on the performance of PTMs, and Niu et al. [48] performed a comparative study of 19 PTMs on 13 SE downstream tasks to better understand these models used in source code learning.

For PEFT methods, He et al. [24] studied the effectiveness of adapter-based tuning methods on several downstream NLP tasks. Chen et al. [11] conducted an empirical study of various widely used PEFT methods on NLP benchmarks. Similar research has been carried out in the SE field. Wang et al. [77] empirically evaluated prompt tuning on CodeBERT and CodeT5 and conducted experiments on three tasks, including defect prediction, code summarization, and code translation. They found that prompt tuning outperformed fine-tuning on all three tasks and showed great potential in low-resource scenarios. This is not completely consistent with our results. We show that although PEFT performs better on the low-resource code generation task, it is not as good as full fine-tuning on the low-resource code classification task. This may be because Wang et al. [77] did not test code classification tasks or the features of prompt tuning. Recently, Liu et al. [35] conducted an empirical study of four PEFT methods on code-specific PTMs and various scenarios, including low-resource, cross-language, and cross-project scenarios. They found that in the low-resource scenario, a considerable number of samples is required for PEFT to approximate the performance of full fine-tuning. They also demonstrated that PEFT methods can improve the transfer ability of PTMs in cross-language and cross-project scenarios. Furthermore, they found that Houlsby and LoRA showed better performance than other PEFT methods, which is different from our results that Parallel is the best method. We believe that the inconsistent results are from different hyperparameters because PEFT can be affected by hyperparameters, such as the rank of LoRA [26]. Nonetheless, our results are consistent with the conclusion of He et al. [23], which is a comprehensive evaluation of PEFT in the NLP field in recent years.

## 7   CONCLUSION

In this study, we focus on the empirical comparison of five popular PEFT and full fine-tuning methods on eight PTMs and four SE downstream tasks. Our experiments provide several important findings. For example, Parallel outperforms others among PEFT methods and PEFT methods perform better on code search and code translation. We believe our extensive experiments and evaluation can equip SE researchers with a deeper understanding of PEFT methods on various PTMs, along with their respective strengths and weaknesses. We hope that this paper will offer researchers comprehensive insights into various PEFT methods in the SE field, and inspire them to apply or design more effective PEFT methods to reduce the cost of using PTMs while maintaining their performance.

## REFERENCES

[1] David Adelani, Jesujoba Alabi, Angela Fan, Julia Kreutzer, Xiaoyu Shen, Machel Reid, Dana Ruiter, Dietrich Klakow, Peter Nabende, Ernie Chang, et al. 2022. A Few Thousand Translations Go a Long Way! Leveraging Pre-trained Models for African News Translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 3053–3070.

[2] Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. Unified pre-training for program understanding and generation. *arXiv preprint arXiv:2103.06333* (2021).

[3] Toufique Ahmed and Premkumar Devanbu. 2022. Multilingual training for software engineering. In *Proceedings of the 44th International Conference on Software Engineering*. 1443–1455.

[4] Muhammad Azeem Akbar, Arif Ali Khan, and Peng Liang. 2023. Ethical Aspects of ChatGPT in Software Engineering Research. *arXiv preprint arXiv:2306.07557* (2023).

[5] Shamil Ayupov and Nadezhda Chirkova. 2022. Parameter-Efficient Finetuning of Transformers for Source Code. *arXiv preprint arXiv:2212.05901* (2022).

[6] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258* (2021).

[7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.

[8] Ernie Chang, Xiaoyu Shen, Alex Marin, and Vera Demberg. 2021. The SelectGen Challenge: Finding the Best Training Samples for Few-Shot Neural Text Generation. In *Proceedings of the 14th International Conference on Natural Language Generation*. 325–330.

[9] Amélie Chatelain, Amine Djeghri, Daniel Hesslow, and Julien Launay. 2022. Is the number of trainable parameters all that actually matters?. In *I (Still) Can't Believe It's Not Better! Workshop at NeurIPS 2021*. PMLR, 27–32.

[10] Fuxiang Chen, Fatemeh H Fard, David Lo, and Timofey Bryksin. 2022. On the transferability of pre-trained language models for low-resource programming languages. In *Proceedings of the 30th IEEE/ACM International Conference on Program Comprehension*. 401–412.

[11] Guanzheng Chen, Fangyu Liu, Zaiqiao Meng, and Shangsong Liang. 2022. Revisiting parameter-efficient tuning: Are we really there yet? *arXiv preprint arXiv:2202.07962* (2022).

[12] Nadezhda Chirkova and Sergey Troshin. 2021. Empirical study of transformers for source code. In *Proceedings of the 29th ACM joint meeting on European software engineering conference and symposium on the foundations of software engineering*. 703–715.

[13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

[14] Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. 2023. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence* 5, 3 (2023), 220–235.

[15] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. 2020. Codebert: A pre-trained model for programming and natural languages. *arXiv preprint arXiv:2002.08155* (2020).

[16] Zihao Fu, Haoran Yang, Anthony Man-Cho So, Wai Lam, Lidong Bing, and Nigel Collier. 2023. On the effectiveness of parameter-efficient fine-tuning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 12799–12807.

[17] Divyam Goel, Ramansh Grover, and Fatemeh H Fard. 2022. On the cross-modal transfer from natural language to code through adapter modules. In *Proceedings of the 30th IEEE/ACM International Conference on Program Comprehension*. 71–81.

[18] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. 2021. Knowledge distillation: A survey. *International Journal of Computer Vision* 129 (2021), 1789–1819.

[19] Daya Guo, Shuai Lu, Nan Duan, Yanlin Wang, Ming Zhou, and Jian Yin. 2022. Unixcoder: Unified cross-modal pre-training for code representation. *arXiv preprint arXiv:2203.03850* (2022).

[20] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, et al. 2020. Graphcodebert: Pre-training code representations with data flow. *arXiv preprint arXiv:2009.08366* (2020).

[21] Demi Guo, Alexander M Rush, and Yoon Kim. 2020. Parameter-efficient transfer learning with diff pruning. *arXiv preprint arXiv:2012.07463* (2020).

[22] Sakib Haque, Alexander LeClair, Lingfei Wu, and Collin McMillan. 2020. Improved automatic summarization of subroutines via attention to file context. In *Proceedings of the 17th International Conference on Mining Software Repositories*. 300–310.

[23] Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2021. Towards a unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366* (2021).

[24] Ruidan He, Linlin Liu, Hai Ye, Qingyu Tan, Bosheng Ding, Liying Cheng, Jia-Wei Low, Lidong Bing, and Luo Si. 2021. On the effectiveness of adapter-based tuning for pretrained language model adaptation. *arXiv preprint arXiv:2106.03164* (2021).

[25] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP. In *International Conference on Machine Learning*. PMLR, 2790–2799.

[26] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685* (2021).

[27] Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2019. Codesearchnet challenge: Evaluating the state of semantic code search. *arXiv preprint arXiv:1909.09436* (2019).

[28] Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and Luke Zettlemoyer. 2018. Mapping language to code in programmatic context. *arXiv preprint arXiv:1808.09588* (2018).

[29] Aditya Kanade, Petros Maniatis, Gogul Balakrishnan, and Kensen Shi. 2020. Learning and evaluating contextual embedding of source code. In *International conference on machine learning*. PMLR, 5110–5121.

[30] Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691* (2021).

[31] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461* (2019).

[32] Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190* (2021).

[33] Vladislav Lialin, Vijeta Deshpande, and Anna Rumshisky. 2023. Scaling down to scale up: A guide to parameter-efficient fine-tuning. *arXiv preprint arXiv:2303.15647* (2023).

[34] Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A Raffel. 2022. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *Advances in Neural Information Processing Systems* 35 (2022), 1950–1965.

[35] Jiaxing Liu, Chaofeng Sha, and Xin Peng. 2023. An Empirical Study of Parameter-Efficient Fine-Tuning Methods for Pre-trained Code Models. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE.

[36] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *Comput. Surveys* 55, 9 (2023), 1–35.

[37] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2023. GPT understands, too. *AI Open* (2023).

[38] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).

[39] Zihan Liu, Genta Indra Winata, and Pascale Fung. 2021. Continual mixed-language pre-training for extremely low-resource neural machine translation. *arXiv preprint arXiv:2105.03953* (2021).

[40] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, et al. 2021. Codexglue: A machine learning benchmark dataset for code understanding and generation. *arXiv preprint arXiv:2102.04664* (2021).

[41] Wei Ma, Shangqing Liu, Wenhan Wang, Qiang Hu, Ye Liu, Cen Zhang, Liming Nie, and Yang Liu. 2023. The Scope of ChatGPT in Software Engineering: A Thorough Investigation. *arXiv preprint arXiv:2305.12138* (2023).

[42] Antonio Mastropaolo, Luca Pascarella, Emanuela Guglielmi, Matteo Ciniselli, Simone Scalabrino, Rocco Oliveto, and Gabriele Bavota. 2023. On the robustness of code generation techniques: An empirical study on github copilot. *arXiv preprint arXiv:2302.00438* (2023).

[43] Aditya Mogadala, Xiaoyu Shen, and Dietrich Klakow. 2020. Integrating image captioning with rule-based entity masking. *arXiv preprint arXiv:2007.11690* (2020).

[44] Lili Mou, Ge Li, Lu Zhang, Tao Wang, and Zhi Jin. 2016. Convolutional neural networks over tree structures for programming language processing. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 30.

[45] Kawser Wazed Nafi, Tonny Shekha Kar, Banani Roy, Chanchal K Roy, and Kevin A Schneider. 2019. Clcdsa: cross language code clone detection using syntactical features and api documentation. In *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 1026–1037.

[46] Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2022. Codegen: An open large language model for code with multi-turn program synthesis. *arXiv preprint arXiv:2203.13474* (2022).

[47] Changan Niu, Chuanyi Li, Bin Luo, and Vincent Ng. 2022. Deep learning meets software engineering: A survey on pre-trained models of source code. *arXiv preprint arXiv:2205.11739* (2022).

[48] Changan Niu, Chuanyi Li, Vincent Ng, Dongxiao Chen, Jidong Ge, and Bin Luo. 2023. An Empirical Comparison of Pre-Trained Models of Source Code. *arXiv preprint arXiv:2302.04026* (2023).

[49] Changan Niu, Chuanyi Li, Vincent Ng, Jidong Ge, Liguo Huang, and Bin Luo. 2022. Spt-code: Sequence-to-sequence pre-training for learning the representation of source code. *arXiv preprint arXiv:2201.01549* (2022).

[50] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*. 311–318.

[51] Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun Cho, and Iryna Gurevych. 2020. AdapterHub: A Framework for Adapting Transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP 2020): Systems Demonstrations*. Association for Computational Linguistics, Online, 46–54. https://www.aclweb.org/anthology/2020.emnlp-demos.7

[52] Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. 2020. Mad-x: An adapter-based framework for multi-task cross-lingual transfer. *arXiv preprint arXiv:2005.00052* (2020).

[53] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. (2018).

[54] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.

[55] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research* 21, 1 (2020), 5485–5551.

[56] Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. 2017. Learning multiple visual domains with residual adapters. *Advances in neural information processing systems* 30 (2017).

[57] Shuo Ren, Daya Guo, Shuai Lu, Long Zhou, Shujie Liu, Duyu Tang, Neel Sundaresan, Ming Zhou, Ambrosio Blanco, and Shuai Ma. 2020. Codebleu: a method for automatic evaluation of code synthesis. *arXiv preprint arXiv:2009.10297* (2020).

[58] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950* (2023).

[59] Iman Saberi and Fatemeh H Fard. 2023. Model-Agnostic Syntactical Information for Pre-Trained Programming Language Models. *arXiv preprint arXiv:2303.06233* (2023).

[60] Xiaoyu Shen, Gianni Barlacchi, Marco Del Tredici, Weiwei Cheng, and Adrià de Gispert. 2022. semiPQA: A Study on Product Question Answering over Semi-structured Data. In *Proceedings of the Fifth Workshop on e-Commerce and NLP (ECNLP 5)*. 111–120.

[61] Xiaoyu Shen, Youssef Oualil, Clayton Greenberg, Mittul Singh, and Dietrich Klakow. 2017. Estimation of Gap Between Current Language Models and Human Performance.

[62] Ensheng Shi, Yanlin Wang, Hongyu Zhang, Lun Du, Shi Han, Dongmei Zhang, and Hongbin Sun. 2023. Towards Efficient Fine-tuning of Pre-trained Code Models: An Experimental Study and Beyond. *arXiv preprint arXiv:2304.05216* (2023).

[63] Jieke Shi, Zhou Yang, Bowen Xu, Hong Jin Kang, and David Lo. 2022. Compressing pre-trained models of code into 3 mb. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering*. 1–12.

[64] Dezhao Song, Andrew Vold, Kanika Madan, and Frank Schilder. 2022. Multi-label legal document classification: A deep learning-based approach with label-attention and domain-specific pre-training. *Information Systems* 106 (2022), 101718.

[65] Benjamin Steenhoek, Md Mahbubur Rahman, Richard Jiles, and Wei Le. 2022. An Empirical Study of Deep Learning Models for Vulnerability Detection. *arXiv preprint arXiv:2212.08109* (2022).

[66] Benjamin Steenhoek, Md Mahbubur Rahman, Richard Jiles, and Wei Le. 2023. An empirical study of deep learning models for vulnerability detection. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2237–2248.

[67] Hui Su, Xiaoyu Shen, Zhou Xiao, Zheng Zhang, Ernie Chang, Cheng Zhang, Cheng Niu, and Jie Zhou. 2020. Moviechats: Chat like humans in a closed domain. In *Proceedings of the 2020 conference on empirical methods in natural language processing (EMNLP)*. 6605–6619.

[68] Hui Su, Weiwei Shi, Xiaoyu Shen, Zhou Xiao, Tuo Ji, Jiarui Fang, and Jie Zhou. 2022. Rocbert: Robust chinese bert with multimodal contrastive pretraining. In *Proceedings of the 60th Annual Meeting of the Association for Computational*

*Linguistics (Volume 1: Long Papers)*. 921–931.

[69] Hui Su, Xiao Zhou, Houjin Yu, Xiaoyu Shen, Yuwen Chen, Zilin Zhu, Yang Yu, and Jie Zhou. 2022. Welm: A well-read pre-trained language model for chinese. *arXiv preprint arXiv:2209.10372* (2022).

[70] Jeffrey Svajlenko, Judith F Islam, Iman Keivanloo, Chanchal K Roy, and Mohammad Mamun Mia. 2014. Towards a big data curated benchmark of inter-project code clones. In *2014 IEEE International Conference on Software Maintenance and Evolution*. IEEE, 476–480.

[71] Alexey Svyatkovskiy, Shao Kun Deng, Shengyu Fu, and Neel Sundaresan. 2020. Intellicode compose: Code generation using transformer. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1433–1443.

[72] Alexey Svyatkovskiy, Sebastian Lee, Anna Hadjitofi, Maik Riechert, Juliana Vicente Franco, and Miltiadis Allamanis. 2021. Fast and memory-efficient neural code completion. In *2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR)*. IEEE, 329–340.

[73] Ze Tang, Chuanyi Li, Jidong Ge, Xiaoyu Shen, Zheling Zhu, and Bin Luo. 2021. AST-transformer: Encoding abstract syntax trees efficiently for code summarization. In *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 1193–1195.

[74] Ze Tang, Xiaoyu Shen, Chuanyi Li, Jidong Ge, Liguo Huang, Zhelin Zhu, and Bin Luo. 2022. AST-trans: Code summarization with efficient tree-structured attention. In *Proceedings of the 44th International Conference on Software Engineering*. 150–162.

[75] Rosalia Tufano, Luca Pascarella, and Gabriele Bavota. 2023. Automating Code-Related Tasks Through Transformers: The Impact of Pre-training. *arXiv preprint arXiv:2302.04048* (2023).

[76] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).

[77] Chaozheng Wang, Yuanhang Yang, Cuiyun Gao, Yun Peng, Hongyu Zhang, and Michael R Lyu. 2022. No more fine-tuning? an experimental evaluation of prompt tuning in code intelligence. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 382–394.

[78] Deze Wang, Boxing Chen, Shanshan Li, Wei Luo, Shaoliang Peng, Wei Dong, and Xiangke Liao. 2023. One Adapter for All Programming Languages? Adapter Tuning for Code Search and Summarization. *arXiv preprint arXiv:2303.15822* (2023).

[79] Deze Wang, Zhouyang Jia, Shanshan Li, Yue Yu, Yun Xiong, Wei Dong, and Xiangke Liao. 2022. Bridging pre-trained models and downstream tasks for source code understanding. In *Proceedings of the 44th International Conference on Software Engineering*. 287–298.

[80] Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. 2021. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. *arXiv preprint arXiv:2109.00859* (2021).

[81] Karl Weiss, Taghi M Khoshgoftaar, and DingDing Wang. 2016. A survey of transfer learning. *Journal of Big data* 3, 1 (2016), 1–40.

[82] Hu Xu, Bing Liu, Lei Shu, and Philip S Yu. 2019. BERT post-training for review reading comprehension and aspect-based sentiment analysis. *arXiv preprint arXiv:1904.02232* (2019).

[83] Shuhan Yan, Hang Yu, Yuting Chen, Beijun Shen, and Lingxiao Jiang. 2020. Are the code snippets what we are searching for? a benchmark and an empirical study on code search with natural-language queries. In *2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 344–354.

[84] Zezhou Yang, Sirong Chen, Cuiyun Gao, Zhenhao Li, Ge Li, and Rongcong Lv. 2023. Deep Learning Based Code Generation Methods: A Literature Review. *arXiv preprint arXiv:2303.01056* (2023).

[85] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems* 32 (2019).

[86] Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. 2021. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. *arXiv preprint arXiv:2106.10199* (2021).

[87] Daoguang Zan, Bei Chen, Fengji Zhang, Dianjie Lu, Bingchao Wu, Bei Guan, Wang Yongji, and Jian-Guang Lou. 2023. Large language models meet NL2Code: A survey. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 7443–7464.

[88] Zhengran Zeng, Hanzhuo Tan, Haotian Zhang, Jing Li, Yuqun Zhang, and Lingming Zhang. 2022. An extensive study on pre-trained models for program understanding and generation. In *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis*. 39–51.

[89] Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023. Adaptive budget allocation for parameter-efficient fine-tuning. *arXiv preprint arXiv:2303.10512* (2023).

[90] Qingyu Zhang, Xiaoyu Shen, Ernie Chang, Jidong Ge, and Pengke Chen. 2022. Mdia: A benchmark for multilingual dialogue generation in 46 languages. *arXiv preprint arXiv:2208.13078* (2022).

[91] Yu Zhang and Qiang Yang. 2018. An overview of multi-task learning. *National Science Review* 5, 1 (2018), 30–43.

[92]  Zibin Zheng, Kaiwen Ning, Jiachi Chen, Yanlin Wang, Wenqing Chen, Lianghong Guo, and Weicheng Wang. 2023.
      Towards an understanding of large language models in software engineering tasks. *arXiv preprint arXiv:2308.11396*
      (2023).
[93]  Wenkang Zhong, Hongliang Ge, Hongfei Ai, Chuanyi Li, Kui Liu, Jidong Ge, and Bin Luo. 2022. StandUp4NPR:
      Standardizing SetUp for Empirically Comparing Neural Program Repair Systems. In *37th IEEE/ACM International
      Conference on Automated Software Engineering*. 1–13.
[94]  Xin Zhou, DongGyun Han, and David Lo. 2021. Assessing generalizability of codebert. In *2021 IEEE International
      Conference on Software Maintenance and Evolution (ICSME)*. IEEE, 425–436.
[95]  Yaqin Zhou, Shangqing Liu, Jingkai Siow, Xiaoning Du, and Yang Liu. 2019. Devign: Effective vulnerability identification
      by learning comprehensive program semantics via graph neural networks. *Advances in neural information processing
      systems* 32 (2019).