# 𝒳FT: Unlocking the Power of Code Instruction Tuning by Simply Merging Upcycled Mixture-of-Experts

**Yifeng Ding, Jiawei Liu, Yuxiang Wei, Terry Yue Zhuo, Lingming Zhang**

University of Illinois Urbana-Champaign

{yifeng6, lingming}@illinois.edu

## Abstract

We introduce 𝒳FT, a simple yet powerful training scheme, by simply merging upcycled Mixture-of-Experts (MoE) to unleash the performance limit of instruction-tuned code Large Language Models (LLMs). While vanilla sparse upcycling fails to improve instruction tuning, 𝒳FT introduces a shared expert mechanism with a novel routing weight normalization strategy into sparse upcycling, which significantly boosts instruction tuning. After fine-tuning the upcycled MoE model, 𝒳FT introduces a learnable model merging mechanism to compile the upcycled MoE back to a dense model, achieving upcycled MoE-level performance with only dense-model compute. By applying 𝒳FT to a 1.3B model, we create a new state-of-the-art tiny code LLM (<3B) with 67.1 and 64.6 pass@1 on HumanEval and HumanEval+ respectively. With the same data and model architecture, 𝒳FT improves supervised fine-tuning (SFT) by 13% on HumanEval+, along with consistent improvements from 2% to 13% on MBPP+, MultiPL-E, and DS-1000, demonstrating its generalizability. 𝒳FT is fully orthogonal to existing techniques such as Evol-Instruct and OSS-INSTRUCT, opening a new dimension for improving code instruction tuning. Codes are available at https://github.com/ise-uiuc/xft.

## 1 Introduction

Program synthesis (or code generation) is a long-standing problem explored since the early days of computer science (Manna and Waldinger, 1971). Recently, instruction tuning of code Large Language Models (LLMs) has been used to improve many coding tasks (Chaudhary, 2023; Luo et al., 2023; Wei et al., 2023), such as text-to-code generation (Chen et al., 2021; Austin et al., 2021), code completion (Cassano et al., 2022), and data science engineering (Lai et al., 2022).

A typical instruction tuning flow involves two steps (Zhang et al., 2023): (i) curating an instruc-
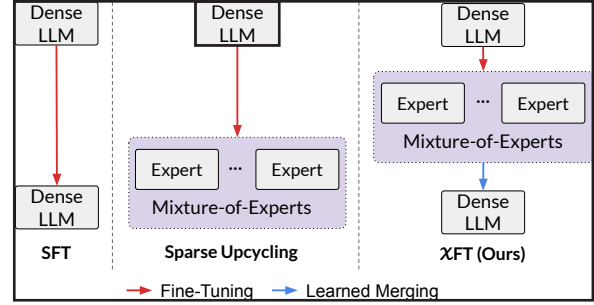


Figure 1: Overview of SFT, sparse upcycling, and 𝒳FT.

tion dataset of instruction-output pairs, where the instruction reflects human intents in natural language and the output includes explained code snippets that correspond to the intent; and (ii) supervised fine-tuning of pre-trained LLM on the instruction dataset. In the realm of code, multiple instruction-tuning methods have been proposed to curate high-quality instruction datasets. For example, *Code Evol-Instruct* (Luo et al., 2023) uses ChatGPT to obtain complex synthetic code instructions with heuristic prompts, while OSS-INSTRUCT (Wei et al., 2023) prompts ChatGPT to generate new coding problems by drawing inspiration from open source code snippets. While existing work focuses on the data perspectives of instruction tuning, they all follow the standard SFT, leaving room for exploring advanced training schemes.

We argue that prior works largely overlook the possibility of improving the code instruction tuning by advancing the training schemes. Figure 1 depicts the supervised fine-tuning (SFT), which directly starts with the pre-trained weights and architecture for fine-tuning. The model is *dense* here because all parameters are activated to compute the next token (assuming it is a decoder-only LLM). In contrast to fine-tuning a *dense* model, following the scaling laws (Kaplan et al., 2020) (*i.e.,* more parameters, better performance), sparse upcycling (Komatsuzaki et al., 2023) is proposed to efficiently upgrade the model sizes by upcycling a dense LLM to a sparsely activated Mixture-of-Experts (MoE)

model. An MoE model is efficient because the generation of the next token only involves a subset of parameters (*i.e.,* experts) and thus is *sparsely activated*. For example, Mixtral-8x7B (Jiang et al., 2024), compared to a dense 7B model, uses approximately $8\times$ parameters and $2\times$ computation, *i.e.,* only 2 out of 8 experts are dynamically selected to compute the next token. However, there are two key limitations when using sparse upcycling in instruction tuning: (i) *Slow scaling:* Komatsuzaki et al. (2023) show that sparse upcycling improves the dense SFT marginally at the early phase, requiring orders of magnitude of extra compute to achieve decent improvement; and (ii) *Inference cost:* though MoE is more efficient than directly scaling the size of dense LLMs, MoE is still expensive, especially at inference, as it introduces significantly more parameters (*i.e.,* memory) and, more importantly, computes during inference, compared to its dense counterparts.

In this paper, we propose $\mathcal{X}$**FT**: by simply merging upcycled MoE models, we push the performance limit of instruction-tuned code LLMs. While vanilla sparse upcycling fails to improve instruction tuning efficiently (Komatsuzaki et al., 2023), $\mathcal{X}$FT addresses this challenge by isolating one expert as the shared expert among all the other experts in each MoE layer, inspired by DeepSeek-MoE (Dai et al., 2024) and MoCLE (Gou et al., 2024). $\mathcal{X}$FT also includes a novel routing weight normalization strategy to eliminate scale mismatch between the upcycled MoE layer with the shared expert and the original dense layer, which will otherwise lead to performance degradation (Wu et al., 2022). After the upcycled MoE model finishes the SFT phase, motivated by Model Soups (Wortsman et al., 2022), $\mathcal{X}$FT uses a learnable model merging mechanism to output a dense model by merging all the expert networks in the upcycled MoE, *i.e.,* the final dense model is of the same model structure and size as the original pre-trained model, achieving similar performance without paying extra inference cost as the sparse upcycling. With only 1.3B parameters, $\mathcal{X}$FT achieves 67.1 pass@1 on HumanEval and 64.6 pass@1 on HumanEval+, which is the new state-of-the-art for tiny code LLMs (<3B). Compared with SFT, $\mathcal{X}$FT achieves 13% improvement on HumanEval+. Surprisingly, our model merging mechanism can preserve or even further boost the general performance of the upcycled MoE with around $1/8\times$ parameters! We conclude our contribution as follows:

- **Dimension:** We open a new dimension of improving instruction tuning of code LLMs by advancing its training scheme, using enhanced sparse upcycling and learnable model merging mechanism, which neither changes the final model structure nor requires more training data.

- **Technique:** We present $\mathcal{X}$FT, a new training scheme for code instruction tuning. $\mathcal{X}$FT involves two steps: *upcycling* and *merging*. A pre-trained dense LLM is first upcycled into an MoE with the shared expert setting and then fine-tuned on the instruction dataset. To avoid the performance degradation caused by the scale mismatch issue, we propose a novel routing weight normalization strategy. In addition, we introduce the first learnable mechanism for merging the upcycled MoE into a dense model, eliminating additional inference overhead while preserving or even improving the MoE performance.

- **Results:** With only 1.3B parameters, $\mathcal{X}$FT achieves 67.1 pass@1 on HumanEval and 64.6 pass@1 on HumanEval+, which is the new state-of-the-art for tiny code LLMs (<3B). Compared with normal supervised fine-tuning (SFT), $\mathcal{X}$FT achieves 13% improvement on HumanEval+! $\mathcal{X}$FT also achieves a consistent improvement from 2% to 13% on MBPP, MultiPL-E, and DS-1000 over SFT, demonstrating its generalization.

## 2 Related Work

### 2.1 Mixture-of-Experts

Mixture-of-Experts (MoE) can efficiently scale up model sizes with only sub-linear increases in computation (Shazeer et al., 2017). Compared with the standard Transformer, MoE replaces each Feed-Forward Network (FFN) layer with an MoE layer, which uses $N$ (*i.e.,* multiple) expert networks that are structurally equivalent to the original FFN layer and uses a router that directs each input token to $K$ out of $N$ expert networks. Formally, for the $l$-th MoE layer, output hidden state $\mathbf{h}_t^l$ of the $t$-th input token is computed as follows (Dai et al., 2024):

$$\mathbf{h}_t^l = \sum_{i=1}^{N}(g_{i,t}\text{FFN}_i(\mathbf{u}_t^l)) + \mathbf{u}_t^l$$

$$g_{i,t} = \begin{cases} s_{i,t} & s_{i,t} \in \text{Topk}(s_t, K) \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

$$s_t = \{s_{i,t} \mid 1 \leq i \leq N\}$$

$$s_{i,t} = \text{Softmax}_i(\mathbf{u}_t^{l\,T}\mathbf{e}_i^l)$$

where $N$ refers to the total number of experts, $g_{i,t}$ refers to the gate value for the $i$-th expert, $\text{FFN}_i(\cdot)$ refers to the $i$-th expert, $\mathbf{u}_t^l$ refers to the hidden states of the $t$-th token which is the input of the $l$-th MoE layer, $s_{i,t}$ refers to the affinity score between the $i$-th expert and the $t$-th token, $\text{Topk}(S, K)$ refers to a function computing $K$ largest scores over $S$, and $\mathbf{e}_i^l$ refers to the centroid of the $i$-th expert in the $l$-th MoE layer. By definition, each token will only be assigned to and computed in the top $K$ experts among all the $N$ experts.

Recently, many works have been proposed to scale model sizes with MoE architecture (Lepikhin et al., 2020; Du et al., 2022; Fedus et al., 2022; Jiang et al., 2024; Xue et al., 2024). While most MoE models are trained from scratch, sparse upcycling (Komatsuzaki et al., 2023) is proposed to initialize MoE models based on the pre-trained dense model, which can efficiently reduce the computational costs of training MoE models, compared with training MoE models from scratch. Specifically, sparse upcycling constructs a new MoE model by initializing each expert of each MoE layer as a copy of the original FFN layer in the dense model, while directly copying the remaining layers from the dense model to the new MoE model.

## 2.2 Instruction Tuning

Instruction tuning is designed to improve the instruction-following ability of LLMs by fine-tuning them on the instruction datasets in a supervised fashion (Wei et al., 2022). The quality of the instruction dataset is significant for the effectiveness of instruction tuning and researchers have proposed multiple methods to improve data quality. For example, SELF-INSTRUCT (Wang et al., 2023) synthesizes high-quality instruction data by prompting a foundation LLM with specially designed prompts. To improve SELF-INSTRUCT, Evol-Instruct (Xu et al., 2023) improves the complexity and diversity of the instruction dataset by prompting ChatGPT with heuristic prompts. OSS-INSTRUCT (Wei et al., 2023) queries ChatGPT to generate instruction-output pairs by getting inspiration from real-world code snippets.

Recently, some parameter-efficient fine-tuning techniques have been proposed to use MoE for better instruction tuning. For example, Lo-RAMoE (Dou et al., 2023) and MoCLE (Gou et al., 2024) propose MoE-like modules that are constructed with Low-Rank Adaptations (LoRA) to improve instruction tuning, while PESC (Wu et al., 2024) proposes to integrate adapters into MoE that are upcycled from dense models. Different from these works, $\mathcal{X}$FT focuses on full fine-tuning, which generally performs better than parameter-efficient fine-tuning (Chen et al., 2022).

## 2.3 Weight Averaging

Weight averaging is a commonly used technique to improve the performance of deep learning models. For example, Model Soups (Wortsman et al., 2022) averages the weights of multiple models that are initialized from the same pre-trained model but finetuned with different hyperparameter configurations to improve the accuracy and robustness of the model. However, only a few works have been proposed to merge expert networks of an MoE layer to a normal FFN layer using weight averaging. For example, OneS (Xue et al., 2022) proposes several simple weight averaging methods to merge expert networks of a BERT-based MoE model. Closely related to our work, Experts Weights Averaging (EWA) (Huang et al., 2023) proposes to convert an MoE model to a dense model with two steps: (i) During MoE training, EWA conducts weighted averaging of all the expert weights after each weight update of MoE, which is based on a manually-crafted hyperparameter $\beta$; (ii) After training, EWA converts each MoE layer into an FFN layer by uniformly averaging the experts.

Different from all the aforementioned existing works, $\mathcal{X}$FT is the first work proposing a **learnable** mechanism to merge expert networks in the upcycled MoE model. Furthermore, while the training scheme of EWA is deeply coupled to a specific MoE architecture, $\mathcal{X}$FT can be easily adapted to different MoE architectures by only adjusting the final merging process. In addition, unlike EWA, $\mathcal{X}$FT does not introduce any hyperparameters into the training of the large MoE models, significantly reducing the computational resources for hyperparameter searching. Our empirical results in Section 4 also showcase the clear advantage of $\mathcal{X}$FT.

## 3 $\mathcal{X}$FT

We describe the details of $\mathcal{X}$FT in this section. There are two steps in our framework: *upcycling* (Section 3.1) and *merging* (Section 3.2). During upcycling, we construct an Mixture-of-Experts (MoE) model from the pre-trained dense model, namely **MoE$_{\text{DS}}$**, which is then fine-tuned on coding instruction data. For merging, we propose a learnable
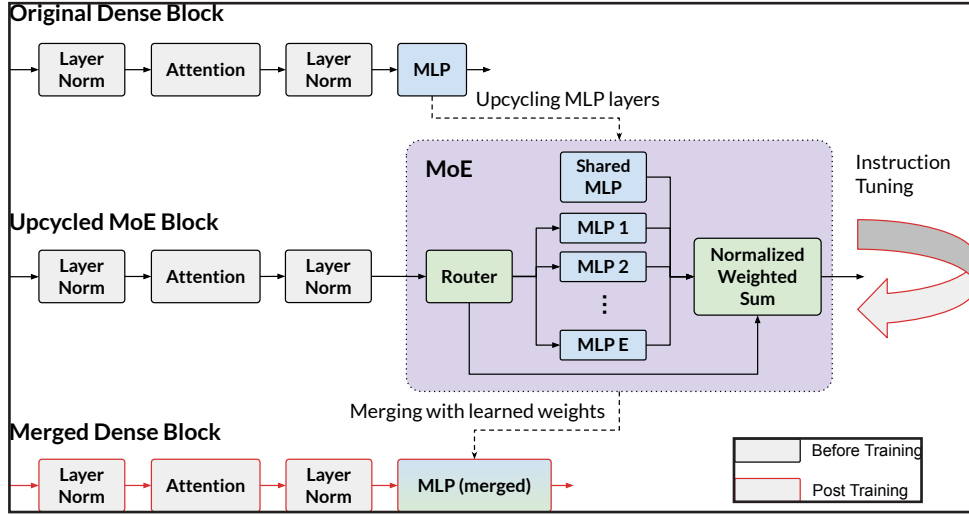
Figure 2: Overview of $\mathcal{X}$FT.

model merging method to convert the instruction-tuned $\text{MoE}_{\text{DS}}$ back to a normal dense model by merging each MoE layer into an FFN layer through weight averaging while directly copying other remaining layers. Consequently, we can obtain $\mathcal{X}\textbf{FT}_{\textbf{DS}}$ that has the same model architecture and size as the original pre-trained dense model, which eliminates all the additional inference overhead brought by the original sparse upcycling, while preserving or even improving the performance of $\text{MoE}_{\text{DS}}$. Our framework is illustrated in Figure 2.

## 3.1 Upcycling

Inspired by sparse upcycling (Komatsuzaki et al., 2023), we convert the pre-trained dense LLM to a new MoE by initializing each expert of each MoE layer as a copy of the original FFN layer in the dense model, while directly copying the remaining layers from the dense model to the new MoE model. However, the performance gain brought by sparse upcycling is negligible with a very limited extra training budget (Komatsuzaki et al., 2023) – which is exactly the situation we are facing during instruction tuning. Intuitively, it is because each expert in the upcycled MoE model is trained on fewer instruction data than the original dense model does because traditional routers used in sparse upcycling will assign different tokens to different experts and thus reduce the amount of data each expert is trained on (Gou et al., 2024). Consequently, inspired by DeepSeekMoE (Dai et al., 2024) and MoCLE (Gou et al., 2024), $\mathcal{X}$FT introduces the shared expert setting into sparse upcycling to tackle this challenge. We further propose a novel routing weight normalization strategy for $\mathcal{X}$FT to avoid the potential performance degradation caused by the scale mismatch problem (Wu et al., 2022).

### 3.1.1 Shared Expert for Upcycling

During upcycling, we isolate one shared expert among all the other normal experts in each MoE layer, where the shared expert will be deterministically assigned to handle all the tokens while other normal experts are assigned by the router. By doing so, the upcycled MoE model can achieve a clear performance boost in instruction tuning, where the shared expert can learn general knowledge across the whole instruction dataset while other normal experts learn specific knowledge among different instructions assigned by the router. Formally, the output hidden state $\mathbf{h}_t^l$ of the $l$-th MoE layer when processing the $t$-th token can be expressed as:

$$\mathbf{h}_t^l = \sum_{i=1}^{N}\left(g_{i,t}\text{FFN}_i(\mathbf{u}_t^l)\right) + \mathbf{u}_t^l$$

$$g_{i,t} = \begin{cases} 1 - s_{t\max} & i = 1 \\ \text{Softmax}_i(s_{i,t}) \cdot s_{t\max} & s_{i,t} \in S_{tK} \\ 0 & \text{otherwise} \end{cases}$$

$$S_{tK} = \text{Topk}(\{s_{i,t} \mid 1 \le i \le N\}, K-1)$$

$$s_{t\max} = \max(\{s_{i,t} \mid 1 \le i \le N\})$$

$$s_{i,t} = \begin{cases} -\infty & i = 1 \\ \text{Softmax}_i(\mathbf{u}_t^{l\,T}\mathbf{e}_i^l) & i \ge 2 \end{cases}$$

$$(2)$$

where $\mathbf{h}_t^l$ refers to the output hidden state of the $l$-th MoE layer when processing the $t$-th token, $N$ refers to the total number of experts, $g_{i,t}$ refers to the gate value for the $i$-th expert, $\text{FFN}_i(\cdot)$ refers to the $i$-th expert, $\mathbf{u}_t^l$ refers to the output hidden state of the $l$-th attention layer when processing the $t$-th token

and also the input of the $l$-th MoE layer, $s_{i,t}$ refers to the affinity score between the $i$-th expert and the $t$-th token, $s_{t_{\max}}$ refers to the maximum affinity score among all the experts besides the shared expert, $\text{Topk}(S, K)$ refers to a function computing $K$ largest scores over $S$, $S_{tK}$ refers to a set of $K - 1$ largest affinity scores among all the experts besides the shared expert, and $\mathbf{e}_i^l$ refers to the centroid of the $i$-th expert in the $l$-th MoE layer.

$\text{FFN}_1$ is chosen as the shared expert in each MoE layer and each token will be assigned to top $K$ experts including one shared expert and $K - 1$ other normal experts. Compared with the original sparse upcycling, there are two major differences:

- **Weighted Shared Expert**. Following Mo-CLE (Gou et al., 2024), with the token-to-expert affinity score $s_{i,t}$, we get the maximum affinity score $s_{t_{\max}}$ and use its complement $1 - s_{t_{\max}}$ as the routing weight of the shared expert.
- **Routing Weight Normalization**. Although the shared expert setting is also used in recent works (Dai et al., 2024; Gou et al., 2024), we cannot directly follow their routing strategy because they cannot handle a scale mismatch problem that is unique for sparse upcycling. The scale mismatch problem is that differences between the scale of the output of the upcycled MoE layer and the original FFN layer can cause performance degradation (Wu et al., 2022). To handle this problem, we need to make sure the sum of $g_{i,t}$ equals 1, so that the output of the MoE layer matches that of the FFN layer in scale. To do so, we normalize the affinity scores of top $K - 1$ normal experts with Softmax and scale their sum to $s_{t_{\max}}$ to make sure that the sum of the $g_{i,t}$ of top $K$ experts, including one shared expert and $K - 1$ normal experts, equals 1.

### 3.2 Merging

We propose a learnable model merging method to convert the large MoE model, namely $\text{MoE}_{\text{DS}}$, back to a dense model $\mathcal{X}\text{FT}_{\text{DS}}$. By doing so, we expect $\mathcal{X}\text{FT}_{\text{DS}}$ to keep the boosted performance gained during upcycling while keeping its model size the same as the original dense model to avoid any additional inference overhead. Inspired by Model Soups (Wortsman et al., 2022), we choose to merge $\text{MoE}_{\text{DS}}$ by learning the mixing coefficients that can be used to average the parameters of all experts in each MoE layer to obtain a normal FFN layer, while directly copying other remaining layers.

Formally speaking, given the weights of $N$ experts at the $l$-th layer $W_1^l, W_2^l, \cdots, W_N^l$, the process of merging each MoE layer to an FFN layer can be stated as below:

$$W^l = \sum_{i=1}^{N} \alpha_i^l W_i^l \qquad (3)$$

where $\overline{W^l}$ denotes the merged parameter of all $N$ experts and $\alpha_i^l$ denotes the learnable mixing coefficient of expert $W_i^l$. We consider a neural network $f(x; \theta)$ with input $x$ and parameters $\theta$. For loss $\mathcal{L}$ and instruction dataset $\{(x_i, y_i)\}_{i=1}^{m}$, such mixing coefficients $\alpha$ of all the $L$ layers can be learned via:

$$\arg\min_{\alpha} \sum_{j=1}^{m} \mathcal{L}(f(x_j; \theta_o, (\sum_{i=1}^{N} \alpha_i^l W_i^l)_{1:L}), y_i) \quad (4)$$

where $\theta_o$ refers to all the remaining layers of $\text{MoE}_{\text{DS}}$ other than MoE layers and $\alpha$ is parameterized as the output of a softmax, so that each $\alpha_i^l$ is positive and $\sum_{i=1}^{l} \alpha_i^l = 1$.

While the learning process defined in Eq. (4) is the most intuitive way of learning $\alpha$, our experiment in Section 5.2 shows that, due to the shared expert setting, it tends to simply increase the mixing coefficient of the shared expert at each layer as much as possible to decrease the loss. It is not helpful because, although the shared expert has learned general knowledge across the whole instruction dataset and needs a relatively large mixing coefficient, we still need to keep the scale of the mixing coefficient of other normal experts at a certain level also to keep some specific knowledge learned by other normal experts in the merged parameter $\overline{W^l}$.

To solve this issue, we introduce a *shared expert rate* $\lambda$ to fix the mixing coefficient of the shared expert and learn the mixing coefficients of the remaining normal experts which sums to $1 - \lambda$ in each layer. By doing so, we can easily control the scale of the mixing coefficient of the shared expert, while still being able to learn the optimal layer-wise mixing coefficients of other normal experts. Let's say $W_1^l$ is the shared expert of the $l$-th layer, then Eq. (3) and Eq. (4) can be reformulated as below:

$$\overline{W^l} = \lambda W_1^l + \sum_{i=2}^{N} \alpha_i^l W_i^l \qquad (5)$$

$$\arg\min_{\alpha} \sum_{i=1}^{m} \mathcal{L}(f(x_j; \theta_o, \overline{W^l}_{1:L}), y_i) \qquad (6)$$

In practice, we uniformly initialize the mixing coefficients $\alpha$ of all the normal experts as $\frac{1-\lambda}{N-1}$, which is then trained on the same instruction dataset as upcycling.

## 4 Main Evaluation

### 4.1 Experimental Setup

**Training.** We use DeepSeek-Coder-Base 1.3B (Guo et al., 2024) as the main base code LLM. `evol-codealpaca-v1`, an open-source Evol-Instruct (Luo et al., 2023) dataset containing 110K samples, is used as our instruction dataset. **MoE$_{DS}$**, our MoE model upcycled from the base model, is implemented following Llama-MoE (LLaMA-MoE Team, 2023). It is constructed with 8 experts in one expert layer and the top 6 experts[1] are activated for each token, including one shared expert. As such, we denote the model size of **MoE$_{DS}$** as $8 \times 1.3B$. Other hyperparameter settings are detailed in Appendix A.1. We finally obtain $\mathcal{X}$**FT$_{DS}$** by using the learned mixing coefficients to merge MoE layers inside MoE$_{DS}$ as normal FFN layers. Note that $\mathcal{X}$**FT$_{DS}$** is the final instruction-tuned LLM we produce, while **MoE$_{DS}$** is only an intermediate product of $\mathcal{X}$FT framework.

**Baselines.** To study the effectiveness of $\mathcal{X}$FT, we build a baseline model, namely **SFT$_{DS}$**, by directly performing SFT for DeepSeek-Coder-Base 1.3B on `evol-codealpaca-v1`. To compare $\mathcal{X}$FT with EWA (Huang et al., 2023), we also implement a baseline **EWA$_{DS}$** and instruction-tune it using the same hyperparameter setting as SFT$_{DS}$, which is described in Appendix A.1. More implementation details of EWA$_{DS}$ can be seen in Appendix A.2. Furthermore, we incorporate multiple small open-source models (<3B) as our baselines, including DeepSeek-Coder-Base 1.3B, DeepSeek-Coder-Instruct 1.3B (Guo et al., 2024), Phi-2 2.7B, and STABLE-CODE 3B (Pinnaparaju et al., 2024).

### 4.2 Python Text-to-Code Generation

HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) benchmarks are the two most widely-used collections of Python code generation tasks. We further employ HumanEval+ and MBPP+, which use more tests automatically generated by EvalPlus (Liu et al., 2023) for more rigorous evaluation. We leave the details in Appendix A.3.

---

[1] 6 is the best-performing number of activated experts per our HumanEval+ experiments using top $\{2, 4, 6\}$ experts.

Table 1 shows the pass@1 results of different LLMs. $\mathcal{X}$FT achieves 67.1 pass@1 on HumanEval and 64.6 pass@1 on HumanEval+, which makes it the new state-of-the-art small code LLM (<3B). We can also observe that $\mathcal{X}$FT$_{DS}$ has a clear improvement over the SFT$_{DS}$ on both benchmarks, with 13% and 2% improvement on HumanEval+ and MBPP+ respectively, while EWA$_{DS}$ even performs worse than SFT$_{DS}$ on MBPP(+). $\mathcal{X}$FT$_{DS}$ also outperforms EWA$_{DS}$ on both benchmarks. Surprisingly, $\mathcal{X}$FT$_{DS}$ even surpasses MoE$_{DS}$ on HumanEval and HumanEval+, despite only using around $1/8\times$ parameters and around $1/6\times$ computations, which showcases the effectiveness of our simple learnable merging technique. Appendix A.4 further demonstrates the statistical significance of the improvements brought by $\mathcal{X}$FT.

### 4.3 Multilingual Code Generation

We use MultiPL-E (Cassano et al., 2022), a multi-programming benchmark that supports 18 programming languages in addition to Python, to evaluate the multilingual ability and generalizability of $\mathcal{X}$FT. Among these, we choose 6 representative programming for their distinct language features: Java, JavaScript, C++, PHP, Swift, and Rust, following Wei et al. (2023). Table 2 shows, among all 1.3B models, $\mathcal{X}$FT$_{DS}$ achieves the best average multilingual performance and performs the best on 5 (out of 6) individual programming languages, overall largely improving SFT$_{DS}$ which uses standard SFT. Notably, the overall performance of EWA$_{DS}$ is on par with SFT$_{DS}$, indicating that EWA$_{DS}$ may not improve SFT on multilingual coding. Appendix A.5 further studies whether each expert in MoE$_{DS}$ specializes differently in these programming languages.

### 4.4 Code Generation for Data Science

The DS-1000 dataset (Lai et al., 2022) is a collection of 1000 realistic data science coding problems ranging from 7 popular data science libraries in Python, including Matplotlib (plt), NumPy (np), Pandas (pd), SciPy (scp), Scikit-Learn (sk), PyTorch (py), and TensorFlow (tf). We evaluate $\mathcal{X}$FT on DS-1000 to understand its effectiveness for practical data science engineering. We follow the evaluation setting of prior works (Guo et al., 2024; Wei et al., 2023). In Table 3, $\mathcal{X}$FT$_{DS}$ achieves the best overall performance among all the evaluated 1.3B models. Specifically, $\mathcal{X}$FT$_{DS}$ consistently surpasses SFT$_{DS}$ among all the seven studied libraries

| Model | Size | Instruction Dataset | Dataset Size | HumanEval (+) | MBPP (+) |
|---|---|---|---|---|---|
| GPT-3.5 (May 2023) | - | Private | - | 73.2 (66.5) | - |
| STABLE-CODE | 3B | - | - | 28.7 (25.6) | 53.6 (44.1) |
| DeepSeek-Coder-Base | 1.3B | - | - | 28.7 (25.6) | 55.6 (46.9) |
| Phi-2 | 2.7B | - | - | 48.8 (45.1) | 62.7 (52.9) |
| DeepSeek-Coder-Instruct | 1.3B | Private | 2B | 65.2 (59.8) | 63.9 (53.1) |
| $\text{SFT}_{\text{DS}}$ | 1.3B | Evol-Instruct | 0.3B | 61.6 (57.3) | 59.6 (49.1) |
| $\text{EWA}_{\text{DS}}$ | 1.3B | Evol-Instruct | 0.3B | **67.1** (63.4) | 58.9 (48.4) |
| $\text{MoE}_{\text{DS}}$ | 8×1.3B | Evol-Instruct | 0.3B | 65.2 (62.2) | 60.4 (50.1) |
| $\mathcal{X}\text{FT}_{\text{DS}}$ | 1.3B | Evol-Instruct | 0.3B | **67.1** (**64.6**) | **60.4** (**50.1**) |

Table 1: Pass@1 results of different LLMs on HumanEval (+) and MBPP (+) computed with greedy decoding, following the setting of prior works (Wei et al., 2023; Liu et al., 2023). We report the results consistently from the EvalPlus (Liu et al., 2023) Leaderboard. Note that numbers in bold refer to the highest scores among all 1.3B models fine-tuned on public datasets, which is the same for all the other tables.

| Model | Size | Programming Language | | | | | | Average |
|---|---|---|---|---|---|---|---|---|
| | | C++ | PHP | Java | JS | Swift | Rust | |
| DeepSeek-Coder-Base | 1.3B | 28.1 | 22.9 | 27.2 | 28.7 | 10.9 | 18.0 | 22.6 |
| $\text{SFT}_{\text{DS}}$ | 1.3B | 40.4 | 38.5 | **40.2** | 46.2 | 16.4 | 27.7 | 34.9 |
| $\text{EWA}_{\text{DS}}$ | 1.3B | 39.4 | 38.4 | 37.3 | 45.2 | 20.9 | 28.6 | 35.0 |
| $\text{MoE}_{\text{DS}}$ | 8×1.3B | 42.2 | 42.2 | 35.4 | 49.8 | 24.7 | 30.6 | 37.5 |
| $\mathcal{X}\text{FT}_{\text{DS}}$ | 1.3B | **42.7** | **41.5** | 36.0 | **49.7** | **25.3** | **32.1** | **37.9** |

Table 2: Pass@1 results on MultiPL-E (Cassano et al., 2022) following the same hyperparameter settings as prior works (Wei et al., 2023; Luo et al., 2023): temperature = 0.2, top_p = 0.95, max_length = 512, and num_samples = 50. All models are evaluated using bigcode-evaluation-harness (Ben Allal et al., 2022).

and also outperforms $\text{EWA}_{\text{DS}}$ in general.

## 5 Ablation Study

### 5.1 Effect of Shared Expert with Routing Weight Normalization

We demonstrate the importance of the shared expert of $\mathcal{X}$FT by comparing its performance with the sparse upcycling (Komatsuzaki et al., 2023) baseline that does not employ any shared expert. As shown in Table 4, the performance of the original sparse upcycling (with the "- Shared Expert" label) drops greatly compared with $\text{MoE}_{\text{DS}}$. Notably, the sparse upcycling model even performs worse than $\text{SFT}_{\text{DS}}$ on HumanEval+, showing its ineffectiveness for instruction tuning.

While the shared expert setting is also employed in most recent works (Dai et al., 2024; Gou et al., 2024), their routing strategy will cause performance degradation due to the scale mismatch between the outputs of the upcycled MoE layer and the original FFN layer. To understand the importance of routing weight normalization, we conduct

an ablation by excluding it from $\mathcal{X}$FT. Table 4 shows that, after removing routing weight normalization, the performance substantially decreases, despite being still better than the original sparse upcycling that does not use the shared expert setting.

### 5.2 Effect of Merging Strategy

In this section, we demonstrate the effectiveness of our learnable merging technique by comparing it with (1) directly merging experts with initialized mixing coefficients, and (2) the learnable merging technique without the shared rate setting, which is the same setting as the learned soup in Model Soups (Wortsman et al., 2022) and is described in Eq. (3) and Eq. (4). Specifically, we initialize the learnable mixing coefficient of the shared expert as 0.75 and that of the other 7 normal experts as $\frac{1}{28}$ for fair comparison. As shown in Table 5, trained mixing coefficients outperform the initialized mixing coefficients for merging. Furthermore, removing the shared rate setting will largely degrade the performance of $\mathcal{X}\text{FT}_{\text{DS}}$ on both HumanEval and

| Model | Size | Data Science Library | | | | | | | Overall |
|---|---|---|---|---|---|---|---|---|---|
| | | np | pd | plt | py | scp | tf | sk | |
| DeepSeek-Coder-Base | 1.3B | 25.1 | 5.8 | 34.5 | 12.7 | 9.8 | 11.1 | 12.7 | 16.4 |
| SFT$_{DS}$ | 1.3B | 30.9 | 17.0 | 40.5 | 32.7 | 18.3 | 21.1 | 24.4 | 25.9 |
| EWA$_{DS}$ | 1.3B | 32.9 | 19.4 | **41.8** | 25.7 | 17.7 | **22.2** | 33.0 | 27.8 |
| MoE$_{DS}$ | 8×1.3B | 33.2 | 21.3 | 38.4 | 41.8 | 21.8 | 23.5 | 37.5 | 30.0 |
| $\mathcal{X}$FT$_{DS}$ | 1.3B | **32.9** | **20.2** | 38.9 | **41.4** | **21.1** | 16.9 | **37.5** | **29.3** |

Table 3: Pass@1 results on DS-1000 (completion format) with `temperature = 0.2`, `top_p = 0.5`, `max_length = 1024`, and `num_samples = 40`, following the same hyperparameter setting used in prior works (Wei et al., 2023).

| Model | HumanEval | HumanEval+ |
|---|---|---|
| SFT$_{DS}$ | 61.6 | 57.3 |
| MoE$_{DS}$ | **65.2** | **62.2** |
| MoE$_{DS}$ - Normalization | 63.4 | 59.1 |
| MoE$_{DS}$ - Shared Expert | 61.6 | 56.7 |

Table 4: Ablation over the design of MoE$_{DS}$. "- Normalization" removes the routing weight normalization from the router, making it the same design as MoCLE (Gou et al., 2024). "- Shared Expert" removes the shared expert setting, making MoE$_{DS}$ the same architecture as original sparse upcycling (Komatsuzaki et al., 2023).

| Model | HumanEval | HumanEval+ |
|---|---|---|
| MoE$_{DS}$ | 65.2 | 62.2 |
| $\mathcal{X}$FT$_{DS}$ (INIT) | 66.5 | 64.0 |
| $\mathcal{X}$FT$_{DS}$ | **67.1** | **64.6** |
| $\mathcal{X}$FT$_{DS}$ - Shared Expert Rate | 66.5 | 64.0 |

Table 5: Ablation over the design of $\mathcal{X}$FT$_{DS}$. "(INIT)" refers to directly using the initialized mixing coefficients to merge experts without training. "- Shared Rate" removes the shared rate setting from $\mathcal{X}$FT$_{DS}$, which is the same as the learned soup (Wortsman et al., 2022).

HumanEval+, demonstrating its importance.

We further study the effect of the shared expert rate $\lambda$ on the performance of the final merged dense model. We evenly choose five shared expert rates, including 0.00, 0.25, 0.50, 0.75, and 1.00, to perform the learnable merging process and evaluate each merged dense model accordingly. Note that 0.75 is the default shared expert rate used in our main experiments. If the shared expert rate is 0.00, it means that the shared expert is ignored when constructing the merged dense model from the up-cycled MoE model; if the shared expert rate is 1.00, it means that the final dense model is built by simply extracting the shared expert from the upcycled

| Model | $\lambda$ | HumanEval | HumenEval+ |
|---|---|---|---|
| SFT$_{DS}$ | - | 61.6 | 57.3 |
| $\mathcal{X}$FT$_{DS}$ | 0.00 | 62.8 | 59.8 |
| | 0.25 | 64.6 | 61.0 |
| | 0.50 | 65.9 | 62.8 |
| | **0.75** | **67.1** | **64.6** |
| | 1.00 | 63.4 | 60.4 |

Table 6: Ablation over the effect of the shared expert rate $\lambda$ in our learnable merging technique. $\mathcal{X}$FT can consistently outperform the normal SFT baseline regardless of the shared expert rate, while $\lambda = 0.75$ is the optimal setting in our experiments.

MoE model. As shown in Table 6, there are mainly three interesting observations:

- The performance of the final merged dense model improves gradually when the shared expert rate grows from 0.00 to 0.75, indicating that **general knowledge learned by the shared expert is important for better performance**.
- The performance of the final merged dense model drops significantly when the shared expert rate grows from 0.75 to 1.00, showing that **specific knowledge learned by other experts is also irreplaceable** and ignoring them will lead to a significant performance drop.
- All the final merged dense models **consistently outperform the normal SFT baseline regardless of their shared expert rate**, further demonstrating the effectiveness of $\mathcal{X}$FT.

### 5.3 Effect of Base Code LLM

In this section, we demonstrate that the effectiveness of $\mathcal{X}$FT is not dependent on the choice of base code LLMs. To show this, we conduct an ablation experiment by applying $\mathcal{X}$FT to STABLE-CODE 3B (Pinnaparaju et al., 2024), whose architecture is different from DeepSeek-Coder-Base 1.3B (Guo et al., 2024), and see whether it can still improve

| Model | HumanEval | HumanEval+ |
|---|---|---|
| $\text{SFT}_{\text{STABLE}}$ | 62.2 | 56.1 |
| $\text{MoE}_{\text{STABLE}}$ | 64.0 | 59.1 |
| $\mathcal{X}\text{FT}_{\text{STABLE}}$ | **68.3** | **62.2** |

Table 7: Ablation over the effect of the base model by replacing DeepSeek-Coder-Base 1.3B with STABLE-CODE 3B. $\mathcal{X}$FT can consistently improve the instruction tuning performance of different base code LLMs.

| Model | HumanEval | HumanEval+ |
|---|---|---|
| $\text{SFT}_{\text{DS}}$ w/ same steps | 61.6 | 57.3 |
| $\text{SFT}_{\text{DS}}$ w/ same budget | 62.2 | 57.3 |
| $\mathcal{X}\text{FT}_{\text{DS}}$ | **67.1** | **64.6** |

Table 8: Experiments on the effect of training overhead. For our two SFT baselines, "w/ same steps" refers to one SFT baseline using the same training steps as $\mathcal{X}$FT while "w/ same budget" refers to the other SFT baseline using the same training budget as $\mathcal{X}$FT. $\mathcal{X}$FT can consistently outperform both SFT baselines to a large extent, further demonstrating the ability of $\mathcal{X}$FT to unlock the power of code instruction tuning.

its performance. Hyperparameter settings are detailed in Appendix A.6. As is shown in Table 7, $\mathcal{X}\text{FT}_{\text{STABLE}}$ significantly improves $\text{SFT}_{\text{STABLE}}$ by 10% on HumanEval and 11% on HumanEval+ respectively. Furthermore, $\mathcal{X}\text{FT}_{\text{STABLE}}$ consistently boosts the performance of $\text{MoE}_{\text{STABLE}}$ while only using $1/4\times$ parameters and $1/2\times$ computations. These results show that the effectiveness of $\mathcal{X}$FT does not depend on any specific choice of base code LLMs, demonstrating the generalizability of $\mathcal{X}$FT across different model architectures.

## 6 Discussion

### 6.1 Training Overhead Analysis

Compared with SFT, $\mathcal{X}$FT will inevitably introduce additional overhead in the training process because $\mathcal{X}$FT needs to fine-tune the upcycled MoE model while the normal SFT technique only needs to fine-tune the original dense model. To better understand the effect of such overhead, we conduct an experiment where we use the same training budget (i.e., the same GPU hours) instead of the same training steps for the normal SFT baseline. As shown in Table 8, although sharing the same training budget as $\mathcal{X}\text{FT}_{\text{DS}}$, the performance of $\text{SFT}_{\text{DS}}$ is still significantly worse than that of $\mathcal{X}\text{FT}_{\text{DS}}$, demonstrating the ability of $\mathcal{X}$FT to unlock the power of code instruction tuning using the same training budget.

### 6.2 Generalizability for General Tasks

In this section, we demonstrate that $\mathcal{X}$FT can improve the performance of LLMs on general tasks across different domains by applying $\mathcal{X}$FT to general instruction tuning. We use TinyLlama 1.1B (Zhang et al., 2024) as the base model and use evol-instruct-70k (Xu et al., 2023) as the training dataset for general instruction tuning. Following existing work (Zhang et al., 2024), we use MMLU (Hendrycks et al., 2021) with the 5-shot setting as our evaluation benchmark to showcase the general performance of LLMs. Hyperparameter settings are detailed in Appendix A.7. As shown in Table 9, overall, $\mathcal{X}\text{FT}_{\text{TL}}$ improves $\text{SFT}_{\text{TL}}$ by 5% on MMLU, demonstrating the generalizable effectiveness of $\mathcal{X}$FT for general instruction tuning.

### 6.3 Preliminary Theoretical Explanation

We provide a preliminary theoretical explanation of $\mathcal{X}$FT by considering a simplified variant of it. Let's start by analyzing the two major steps of $\mathcal{X}$FT:

- **Step 1: Upcycling**. According to the scaling law (Kaplan et al., 2020), the upcycled MoE model performs better than the normal SFT dense model due to more trainable parameters.
- **Step 2: Merging**. We consider a simplified variant of $\mathcal{X}$FT, where the upcycled MoE model (e.g., $\text{MoE}_{\text{DS}}$) can be viewed as the ensembling of two dense models and the merged dense model (e.g., $\mathcal{X}\text{FT}_{\text{DS}}$) can be viewed as the merging of the same two dense models; see Appendix A.8 for more details. As such, we can directly apply the theoretical analyzing process in Section 4 of (Wortsman et al., 2022) to analyze the performance difference between the upcycled MoE model and the merged dense model, which is initially designed to analyze the performance difference between model ensembling and model merging. According to (Wortsman et al., 2022), the convexity of the loss can help the merged dense model achieve a similar expected loss as that of the upcycled MoE model.

Overall, the **Upcycling** step improves the performance with more trainable parameters, while the **Merging** step maintains the upcycled MoE-level performance with only dense-model compute. Consequently, we provide a preliminary theoretical explanation for the effectiveness of $\mathcal{X}$FT.

| Model | Discipline | | | | Overall |
|---|---|---|---|---|---|
| | Humanities | Social Science | STEM | Other | |
| SFT$_{TL}$ | **25.38** | 23.30 | 24.20 | 26.78 | 24.97 |
| MoE$_{TL}$ | 23.85 | 26.32 | 27.40 | 28.03 | 26.11 |
| $\mathcal{X}$FT$_{TL}$ | 23.91 | **26.49** | **27.72** | **28.29** | **26.30** |

Table 9: Experiments on the generalizable effectiveness of $\mathcal{X}$FT for general tasks in MMLU benchmark (Hendrycks et al., 2021). It shows that $\mathcal{X}$FT can improve the general instruction tuning performance of LLMs.

## 7 Conclusion

This paper introduces $\mathcal{X}$FT to unlock the power of code instruction tuning by simply merging upcycled MoE. Similar to SFT, $\mathcal{X}$FT starts with a dense LLM and produces a fine-tuned dense LLM with the exact size and model structure. Yet, $\mathcal{X}$FT improves SFT by upcycling the pre-trained dense LLM to an MoE model for fine-tuning, after which we compile the MoE model back to an efficient dense LLM with a learnable merging mechanism. As such, we unleash the performance limit of instruction tuning without any additional inference overhead. Using the same dataset, $\mathcal{X}$FT improves SFT on a variety of benchmarks, including HumanEval(+), MBPP(+), MultiPL-E, and DS-1000, from 2% to 13%. By applying $\mathcal{X}$FT to DeepSeek-Coder-Base 1.3B, we create the next state-of-the-art small (<3B) LLM for code. The ultimate dense LLM produced by $\mathcal{X}$FT preserves or even outperforms the full upcycled MoE which uses $8\times$ parameters as much as our final dense LLM. $\mathcal{X}$FT is fully orthogonal to the existing instruction tuners such as Evol-Instruct and OSS-INSTRUCT, opening a new dimension to maximal code instruction tuning.

## Limitations

While $\mathcal{X}$FT has proven to be effective through extensive experiments in the paper, we apply our technique to LLMs with no more than 3B parameters due to resource constraints. This limitation hinders our ability to showcase the impact of $\mathcal{X}$FT on larger models. In addition, to balance the general knowledge in the shared expert and the specific knowledge in other normal experts, we introduce a hyperparameter $\lambda$ in the merging process of $\mathcal{X}$FT, which might slightly increase the efforts for hyperparameter search. It would be interesting to explore other hyperparameter-free techniques to tackle this challenge in the future. Furthermore, while $\mathcal{X}$FT has been empirically proven powerful, it would be interesting to provide a theoretical explanation for its strong performance.

## References

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. 2021. Program synthesis with large language models.

Loubna Ben Allal, Niklas Muennighoff, Logesh Kumar Umapathi, Ben Lipkin, and Leandro von Werra. 2022. A framework for the evaluation of code generation models. https://github.com/bigcode-project/bigcode-evaluation-harness.

Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, Arjun Guha, Michael Greenberg, and Abhinav Jangda. 2022. Multipl-e: A scalable and extensible approach to benchmarking neural code generation.

Sahil Chaudhary. 2023. Code alpaca: An instruction-following llama model for code generation. https://github.com/sahil280114/codealpaca.

Guanzheng Chen, Fangyu Liu, Zaiqiao Meng, and Shangsong Liang. 2022. Revisiting parameter-efficient tuning: Are we really there yet?

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code.

Damai Dai, Chengqi Deng, Chenggang Zhao, R. X. Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding

Zeng, Xingkai Yu, Y. Wu, Zhenda Xie, Y. K. Li, Panpan Huang, Fuli Luo, Chong Ruan, Zhifang Sui, and Wenfeng Liang. 2024. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models.

Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao, Jun Zhao, Wei Shen, Yuhao Zhou, Zhiheng Xi, Xiao Wang, Xiaoran Fan, Shiliang Pu, Jiang Zhu, Rui Zheng, Tao Gui, Qi Zhang, and Xuanjing Huang. 2023. Loramoe: Revolutionizing mixture of experts for maintaining world knowledge in language model alignment.

Rotem Dror, Gili Baumer, Segev Shlomov, and Roi Reichart. 2018. The hitchhiker's guide to testing statistical significance in natural language processing. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1383–1392, Melbourne, Australia. Association for Computational Linguistics.

Nan Du, Yanping Huang, Andrew M. Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten Bosma, Zongwei Zhou, Tao Wang, Yu Emma Wang, Kellie Webster, Marie Pellat, Kevin Robinson, Kathleen Meier-Hellstern, Toju Duke, Lucas Dixon, Kun Zhang, Quoc V Le, Yonghui Wu, Zhifeng Chen, and Claire Cui. 2022. Glam: Efficient scaling of language models with mixture-of-experts.

William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity.

Yunhao Gou, Zhili Liu, Kai Chen, Lanqing Hong, Hang Xu, Aoxue Li, Dit-Yan Yeung, James T. Kwok, and Yu Zhang. 2024. Mixture of cluster-conditional lora experts for vision-language instruction tuning.

Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. 2024. Deepseek-coder: When the large language model meets programming – the rise of code intelligence.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding.

Yongqi Huang, Peng Ye, Xiaoshui Huang, Sheng Li, Tao Chen, Tong He, and Wanli Ouyang. 2023. Experts weights averaging: A new general training scheme for vision transformers.

Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. Mixtral of experts.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models.

Aran Komatsuzaki, Joan Puigcerver, James Lee-Thorp, Carlos Riquelme Ruiz, Basil Mustafa, Joshua Ainslie, Yi Tay, Mostafa Dehghani, and Neil Houlsby. 2023. Sparse upcycling: Training mixture-of-experts from dense checkpoints.

Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Scott Wen tau Yih, Daniel Fried, Sida Wang, and Tao Yu. 2022. Ds-1000: A natural and reliable benchmark for data science code generation.

Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. 2020. Gshard: Scaling giant models with conditional computation and automatic sharding.

Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation. In *Thirty-seventh Conference on Neural Information Processing Systems*.

LLaMA-MoE Team. 2023. Llama-moe: Building mixture-of-experts from llama with continual pretraining.

Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul, Zhuang Li, WenDing Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii, Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov, Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. 2024. Starcoder 2 and the stack v2: The next generation.

Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder:

Empowering code large language models with evol-instruct.

Zohar Manna and Richard J Waldinger. 1971. Toward automatic program synthesis. *Communications of the ACM*, 14(3):151–165.

Nikhil Pinnaparaju, Reshinth Adithyan, Duy Phung, Jonathan Tow, James Baicoianu, , and Nathan Cooper. 2024. Stable code 3b.

Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer.

Anders Søgaard, Anders Johannsen, Barbara Plank, Dirk Hovy, and Hector Martínez Alonso. 2014. What's in a p-value in NLP? In *Proceedings of the Eighteenth Conference on Computational Natural Language Learning*, pages 1–10, Ann Arbor, Michigan. Association for Computational Linguistics.

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. Self-instruct: Aligning language models with self-generated instructions.

Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. Finetuned language models are zero-shot learners.

Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. 2023. Magicoder: Source code is all you need. *arXiv preprint arXiv:2312.02120*.

Frank. Wilcoxon. 1945. Individual comparisons by ranking methods. *Biometrics*, 1:196–202.

Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S. Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig Schmidt. 2022. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time.

Haoyuan Wu, Haisheng Zheng, and Bei Yu. 2024. Parameter-efficient sparsity crafting from dense to mixture-of-experts for instruction tuning on general tasks.

Lemeng Wu, Mengchen Liu, Yinpeng Chen, Dongdong Chen, Xiyang Dai, and Lu Yuan. 2022. Residual mixture of experts.

Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions.

Fuzhao Xue, Xiaoxin He, Xiaozhe Ren, Yuxuan Lou, and Yang You. 2022. One student knows all experts know: From sparse to dense.

Fuzhao Xue, Zian Zheng, Yao Fu, Jinjie Ni, Zangwei Zheng, Wangchunshu Zhou, and Yang You. 2024. Openmoe: An early effort on open mixture-of-experts language models. *arXiv preprint arXiv:2402.01739*.

Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024. Tinyllama: An open-source small language model.

Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. 2023. Instruction tuning for large language models: A survey.

# A   Appendix for "$\mathcal{X}$FT: Unlocking the Power of Code Instruction Tuning by Simply Merging Upcycled Mixture-of-Experts"

## A.1   Hyperparameter Settings

We use a batch size of 64 and a learning rare of 5e-5 with a linear scheduler to fine-tune **MoE$_{DS}$** for 4 epochs with 500 warmup steps, following the implementation of previous work (Wei et al., 2023). We further use a batch size of 64, a shared expert rate $\lambda$ of 0.75, and a learning rare of 1e-5 with a linear schedule to fine-tune the learnable mixing coefficients for each of the experts in the instruction-tuned **MoE$_{DS}$** on the instruction dataset for 1 epoch with 125 warmup steps. Detailedly, we use Softmax to keep the sum of the mixing coefficients of the other 7 normal experts as 0.25. For **SFT$_{DS}$** and **EWA$_{DS}$**, we use the same hyperparameter setting as $\mathcal{X}$FT, where the batch size is 64 and the learning rate is 5e-5 with a linear scheduler. Because $\mathcal{X}$FT is trained for 4 epochs during upcycling and 1 epoch during merging, for a fair comparison, we train SFT$_{DS}$ and EWA$_{DS}$ for 5 (= 4 + 1) epochs with 625 warmup steps.

## A.2   Implementation details of EWA

Because EWA (Huang et al., 2023) does not release their implementation, we implemented EWA by ourselves, including constant schedule and linear schedule. We use a share rate $\beta$ of 0.3, following the original setting of EWA. While EWA with the constant schedule achieves reasonable performance in our evaluation, the training loss of EWA with the linear schedule becomes very unstable, as is shown in Figure 3, and thus cannot achieve reasonable performance. As a result, we report the results of EWA with the constant schedule in Section 4.

Figure 3: Training loss curve of EWA with constant schedule and linear schedule.

| Model | HumanEval | HumanEval+ |
|---|---|---|
| SFT$_{DS}$ | 61.6 | 57.2 |
| EWA$_{DS}$ | 62.7 | 58.8 |
| $\mathcal{X}$FT$_{DS}$ | **64.5** | **60.9** |

Table 10: Average pass@1 results of 200 experiments on HumanEval (+) computed with sampling. $\mathcal{X}$FT clearly outperforms both EWA$_{DS}$ and SFT$_{DS}$.

## A.3 Details of HumanEval and MBPP

In these benchmarks, each task consists of a task description in English, which is sent to LLMs as the prompt, and LLMs are expected to generate the corresponding code to satisfy the requirements in the description. While these benchmarks provide a handful of test cases to validate the correctness of the generated code, these tests are often insufficient for more rigorous evaluation. As such, HumanEval+ and MBPP+ proposed by EvalPlus (Liu et al., 2023) are usually used to evaluate the correctness of the generated code, which provides 80×/35× more tests compared with the original benchmarks.

## A.4 Statistical Significance Analysis

In this section, we show that improvements brought by $\mathcal{X}$FT are statistically significant. In our main experiments, we follow prior works (Wei et al., 2023; Lozhkov et al., 2024) to conduct experiments on HumanEval (+) using greedy decoding. To demonstrate the statistical significance of our improvements, we change our setting from greedy decoding to sampling. In detail, to conduct one experiment on HumanEval (+), the model will sample one solution for each problem in HumanEval (+) with top $p = 0.95$ and temperature $= 0.8$, which is the same setting used in prior works (Liu et al., 2023; Chen et al., 2021).

| Model | HumanEval | HumanEval+ |
|---|---|---|
| $\mathcal{X}$FT$_{DS}$ vs. EWA$_{DS}$ | 2.6e-18 | 8.0e-23 |
| $\mathcal{X}$FT$_{DS}$ vs. SFT$_{DS}$ | 9.6e-30 | 3.7e-33 |

Table 11: $p$-values for $\mathcal{X}$FT$_{DS}$ vs. EWA$_{DS}$ and $\mathcal{X}$FT$_{DS}$ vs. SFT$_{DS}$ in 200 experiments on HumanEval (+) computed with sampling. Results show that improvements brought by $\mathcal{X}$FT are statistically significant.

Following prior work (Liu et al., 2023), we repeat this experiment 200 times for three techniques: $\mathcal{X}$FT$_{DS}$, EWA$_{DS}$, and SFT$_{DS}$. EWA$_{DS}$ is included because it is the best-performing baseline in our main experiment. We first compute their average pass@1 performance in these 200 experiments. As is shown in Table 10, $\mathcal{X}$FT$_{DS}$ outperforms both EWA$_{DS}$ and SFT$_{DS}$ clearly.

Furthermore, we use the Wilcoxon signed-rank test (Wilcoxon, 1945; Dror et al., 2018), a widely used statistical test, to check if the improvements brought by $\mathcal{X}$FT are statistically significant. As shown in Table 11, the $p$-values for both $\mathcal{X}$FT$_{DS}$ vs. EWA$_{DS}$ and $\mathcal{X}$FT$_{DS}$ vs. SFT$_{DS}$ are much smaller than both 0.0025 (the significance level recommended for NLP work by (Søgaard et al., 2014)) and 0.05 (the most common significance level), demonstrating the statistical significance of the improvements brought by $\mathcal{X}$FT.

## A.5 Analysis on Expert Specialization

Inspired by recent works (Jiang et al., 2024; Xue et al., 2024), we analyze whether each expert in MoE$_{DS}$ has different specializations in different programming languages by visualizing the routing decision of the tokens from different programming languages in the MultiPL-E benchmark (including Python). For the MultiPL-E benchmark, we collect the routing decision when conducting experiments in Section 4.3. For Python, we collect the routing decision by reruning HumanEval experiment following the same setting as Section 4.3. Following Mixtral (Jiang et al., 2024), we get the visualization results from layers 0, 11, and 23 in MoE$_{DS}$, where layer 0 and layer 23 are the first and the last layers of MoE$_{DS}$. As is shown in Figure 4, we do not observe obvious patterns in the assignment of experts based on the programming language, which is in line with the observation reported by recent works (Jiang et al., 2024; Xue et al., 2024).

## A.6 Training Settings for STABLE-CODE 3B

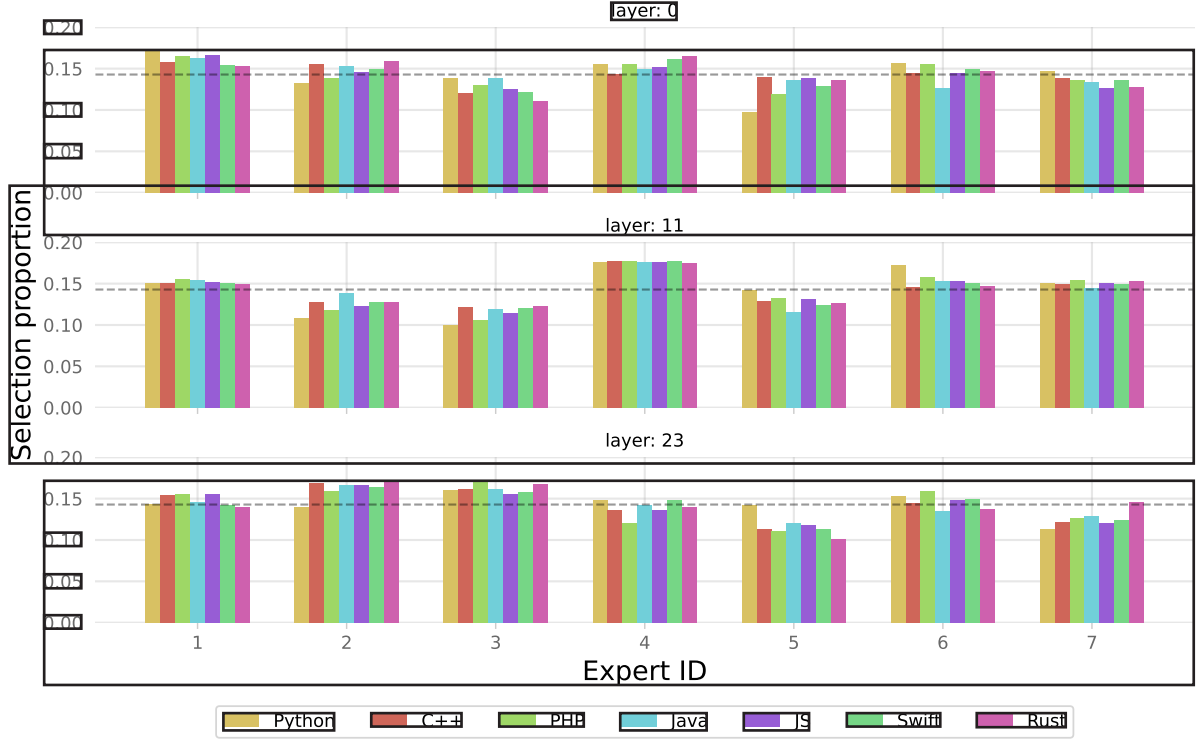We use `evol-codealpaca-v1` as the training dataset. Since STABLE-CODE 3B is the base model,

Figure 4: Proportion of tokens assigned to each expert on different programming languages from MultiPL-E (including Python) for layers 0, 11, and 23. The shared expert 0 is excluded from the chart because all the tokens are always assigned to it. The gray vertical line marks $\frac{1}{7}$, which is the proportion expected with the uniform sampling.

we upcycle a new MoE model from the base model, namely $\text{MoE}_{\textbf{STABLE}}$. Due to limited computational resources, we construct $\text{MoE}_{\text{STABLE}}$ with 4 experts in one expert layer, where the top 2 experts are activated for each token, including one shared expert. Consequently, the size of $\text{MoE}_{\text{STABLE}}$ can be described as $4\times3B$. We use a batch size of 64 and a learning rate of 5e-5 with a linear scheduler to fine-tune $\text{MoE}_{\text{STABLE}}$ for 4 epochs with 500 warmup steps. Similar to $\mathcal{X}\text{FT}_{\text{DS}}$, we obtain $\mathcal{X}\textbf{FT}_{\textbf{STABLE}}$ by learning mixing coefficients to merge MoE layers inside $\text{MoE}_{\text{STABLE}}$ as normal FFN layers, which is fine-tuned with a batch size of 64, a shared expert rate $\lambda$ of 0.85, and a learning rate of 1e-5 with a linear schedule for 1 epoch with 125 warmup steps. Our baseline model, namely $\textbf{SFT}_{\textbf{STABLE}}$, is fine-tuned for 5 (= 4 + 1) epochs with a batch size of 64, a learning rate of 5e-5, and 625 warmup steps for a fair comparison.

### A.7 Training Settings for TinyLlama 1.1B

Using TinyLlama 1.1B as the base model, we upcycle a new MoE model, namely $\textbf{MoE}_{\textbf{TL}}$, from the base model. Following the setting for $\text{MoE}_{\text{DS}}$, we construct $\text{MoE}_{\text{TL}}$ with 8 experts in one expert layer, where the top 6 experts are activated for each token, including one shared expert. As such, the number of parameters for $\text{MoE}_{\text{TL}}$ can be written as $8\times1.1B$. We use a batch size of 64 and a learning rate of 5e-5 with a linear scheduler to fine-tune $\text{MoE}_{\text{TL}}$ for 4 epochs with 240 warmup steps. To obtain $\mathcal{X}\textbf{FT}_{\textbf{TL}}$, we learn mixing coefficients to merge MoE layers inside $\text{MoE}_{\text{TL}}$ by fine-tuning them with a batch size of 64, a shared expert rate $\lambda$ of 0.85, and a learning rate of 2e-5 with a linear schedule for 1 epoch with 60 warmup steps. For a fair comparison, we fine-tune a baseline model $\textbf{SFT}_{\textbf{TL}}$ for 5 (= 4 + 1) epochs with a batch size of 64, a learning rate of 5e-5, and 300 warmup steps.

### A.8 Theoratical Explanation Details

We consider a simplified variant of $\mathcal{X}$FT as below:

- The original dense model is a one-layer transformer model, which contains one attention layer connected with one feed-forward network (FFN) layer. As such, the upcycled MoE model is also a one-layer transformer model, containing one attention layer connected with an MoE layer.

- The upcycled MoE model only has two experts ($\mathbf{e}_1$ and $\mathbf{e}_2$), both of which are always selected for processing the input tokens.

- The router in the MoE model assigns constant weights to each expert, regardless of the input token. Consequently, the output of the MoE

layer for the $t$-th token $\mathbf{h}_t$ can be represented as $(1 - \alpha)\mathbf{e}_1(\mathbf{u}_t) + \alpha\mathbf{e}_2(\mathbf{u}_t)$, where $1 - \alpha$ is the router weight assigned to $\mathbf{e}_1$, $\alpha$ is the router weight assigned to $\mathbf{e}_2$, and $\mathbf{u}_t$ is the input of the MoE layer for the $t$-th token.

- We simplify the process of merging the MoE model back to a dense model as $\mathbf{W}_{\mathbf{e}_\alpha} = (1 - \alpha)\mathbf{W}_{\mathbf{e}_1} + \alpha\mathbf{W}_{\mathbf{e}_2}$, where $\mathbf{W}_{\mathbf{e}}$ refers to the weight of $\mathbf{e}$ and $\mathbf{e}_\alpha$ refers to the weight of the FFN in the merged dense model.

In this simplified scenario, if we denote $f(x; \theta)$ as the output of the model $\theta$ for the input $x$, the output of this simplified MoE model for input token $x$ can be represented as $f(x; \theta_{\mathbf{MoE}})$. Interestingly, if we define two new dense models $\theta_1$ and $\theta_2$, where $\theta_1$ and $\theta_2$ use the same attention layer as this MoE model while using $\mathbf{e}_1$ and $\mathbf{e}_2$ as the FFN layer separately, $f(x; \theta_{\mathbf{MoE}})$ can be represented as $(1 - \alpha)f(x; \theta_1) + \alpha f(x; \theta_2)$! Consequently, the computation process of this simplified MoE model can be viewed as ensembling the outputs of two dense models $\theta_1$ and $\theta_2$. Meanwhile, the process of merging the upcycled MoE model back to a dense model in this simplified $\mathcal{X}$FT can be represented as $\theta_\alpha = (1 - \alpha)\theta_1 + \alpha\theta_2$, which is the merging of the same two dense models $\theta_1$ and $\theta_2$.