
Mastering Text, Code and Math Simultaneously via Fusing Highly Specialized Language Models

Ning Ding^{*1} Yulin Chen^{*1} Ganqu Cui¹ Xingtai Lv¹ Ruobing Xie¹
Bowen Zhou¹ Zhiyuan Liu¹ Maosong Sun¹

Abstract

Underlying data distributions of natural language, programming code, and mathematical symbols vary vastly, presenting a complex challenge for large language models (LLMs) that strive to achieve high performance across all three domains simultaneously. Achieving a very high level of proficiency for an LLM within a specific domain often requires extensive training with relevant corpora, which is typically accompanied by a sacrifice in performance in other domains. In this paper, we propose to fuse models that are already highly-specialized directly. The proposed fusing framework, ULTRAFUSER, consists of three distinct specialists that are already sufficiently trained on language, coding, and mathematics. A token-level gating mechanism is introduced to blend the specialists’ outputs. A two-stage training strategy accompanied by balanced sampling is designed to ensure stability. To effectively train the fused model, we further construct a high-quality supervised instruction tuning dataset, ULTRACHAT 2, which includes text, code, and mathematical content. This dataset comprises approximately 300,000 instructions and covers a wide range of topics in each domain. Experiments show that our model could simultaneously achieve mastery of the three crucial domains.

1. Introduction

If a piece of information can be serialized and tokenized, it is likely to be handled by large language models (LLMs) (Bommasani et al., 2021; Brown et al., 2020; OpenAI, 2023a). LLMs, as one of the most advanced manifestations of artificial intelligence, have demonstrated proficiency in three representative symbol systems that are essential to human progress: natural language (Ouyang et al., 2022; Bai et al., 2022), which forms the cornerstone of human inter-

action; programming code (Li et al., 2023a; Rozière et al., 2023), the backbone of our digital ecosystem; and mathematical reasoning, the framework underpinning scientific advancement (Luo et al., 2023a; Yang et al., 2023). The mastery of three domains would equip LLMs with unparalleled versatility. However, the intrinsic variability of data distribution across these domains presents a formidable challenge for an LLM to achieve consistently high performance *at the same time*. One awkward situation is that it is challenging to integrate professional-level coding and mathematical abilities into a general conversational language model without loss. In other words, these skills are more often reflected in the numbers on related benchmarks rather than a real-world user interface.

Figure 1 (a-c) demonstrates such a struggle by presenting the performance of three specialized models on the aforementioned domains, all initially based on the Llama-2 (Touvron et al., 2023) 13B architecture. Our findings reveal a clear trade-off: specialized training in one domain often comes at the expense of performance in the others, whereas training on all three types of data at the same time results in a simultaneous suboptimal situation. Delving into this situation, such an issue may be partially mitigated by careful designs of data engineering, training strategy, or prompt construction. However, in general, semantics in language, logic and structures in code, and abstract symbol manipulations in math intricately always create a situation of mutual weakening. To elaborate further, comparing highly specialized models (such as those for coding or mathematics) with general-purpose models capable of performing all tasks (like GPT-4) for their expertise is a trap that can easily lead to misinformation.

This paper hopes to integrate specialized abilities into a general chat language model with as little loss as possible. More specifically, we propose to leverage separate models that are already highly specialized via a fusing structure. In this fusing framework, namely ULTRAFUSER, we use three well-trained LLMs as initial specialist models in text, code, and math.¹ To ensure that the fused model benefits from the

^{*}Equal contribution ¹Tsinghua University.
Preprint.

¹Although we treat text, code, and math as three separate domains in this paper according to their symbol systems, they are not

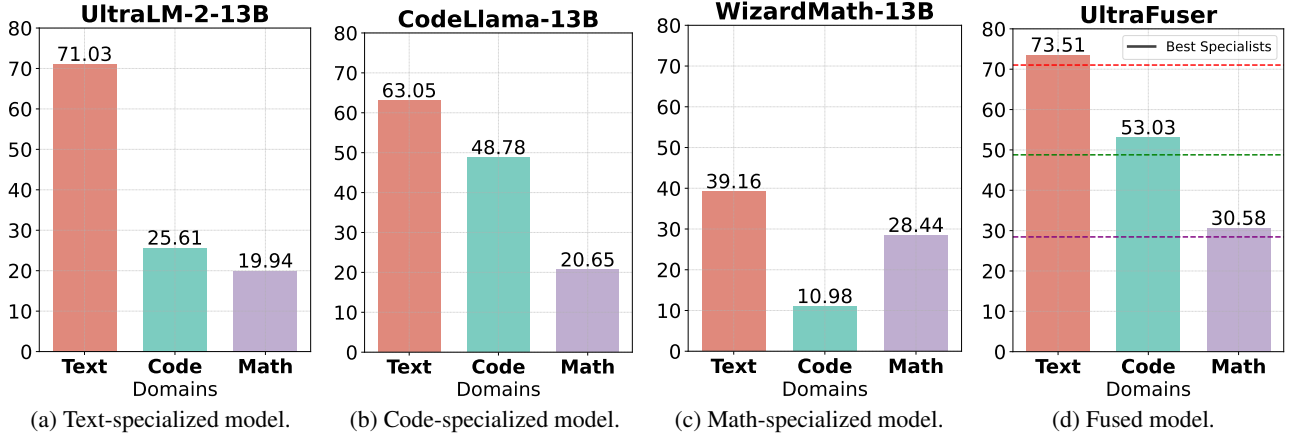


Figure 1: Performance on three different domains of specialized models and our ULTRAFUSER. The performance for the text domain is computed by the average results on TruthfulQA (Acc) (Lin et al., 2021) and AlpacaEval (Win Rate) (Li et al., 2023b) datasets; the performance for the code domain is Pass@1 of HumanEval (Chen et al., 2021); and the performance for the math domain is the average result of GSM8K (Pass@1) (Cobbe et al., 2021), MATH (Pass@1) (Hendrycks et al., 2021), SAT-Math (Acc) (Zhong et al., 2023), and AQuA-RAT (Acc) (Ling et al., 2017) datasets. All numbers are zero-shot results.

specialized knowledge of each specialist model, a dynamic gating mechanism is implemented, which sits on top of the three specialists and adaptively controls the contribution of each specialist to the final output logits based on the input data. Such a mechanism is adopted at the token level, which allows both the specialization of individual specialists and the generalization of the fused model. The key to functioning the model is to train the gating module. For example, when the model conducts code generation, we want the coding specialist to contribute more than the other two. This necessitates a mixed instruction tuning dataset that contains the three domains for the training. Unlike language data, high-quality instruction-tuning datasets for code and math are scarcer in the open-source community. Inspired by ULTRACHAT (Ding et al., 2023), we construct a comprehensive, diverse dataset with high quality, ULTRACHAT 2, to facilitate the development of advanced LLMs with the aforementioned expertise. ULTRACHAT 2 contains 300,000 diverse and high-quality data (each part has 100,000), which are derived from 72 meta-topics and 1587 sub-topics.

Experiments show that highly specialized models may counter collapse if they are directly further trained, but we can effectively integrate their highly professional abilities into a general chat interface via ULTRAFUSER. By training a fused model with UltraLM-2-13B, CodeLlama-13B, and WizardMath-13B as the specialists for three domains, we achieve consistently effective performance on seven benchmarks across language understanding, code generation, and mathematical reasoning. Our proposed model, data, training,

strictly segregated. For example, language can partially encompass the other two. This is discussed in Appendix C.

and inference frameworks will be publicly available.

2. Related Work

Large Language Models for Language. With the proliferation of model parameters, enhancements in training data augmentation both in terms of quantity and quality, and continuous refinements in training algorithms, LLMs have exhibited an enhancement in language understanding, generation, and generalization capabilities. These LLMs exhibit remarkable proficiency in accomplishing a wide array of natural language processing tasks, and showcase formidable capabilities in in-context learning and few-shot learning (Brown et al., 2020; Ouyang et al., 2022; OpenAI, 2023b; Chowdhery et al., 2022; Zhang et al., 2022; Touvron et al., 2023; Taori et al., 2023; Chiang et al., 2023; Xu et al., 2023; Ding et al., 2023; Jiang et al., 2023a). Despite originating from NLP tasks, as LLMs evolve, the boundaries between NLP tasks are gradually becoming blurred.

Large Language Models beyond Language. LLMs excel in processing various symbol systems including code, math symbols, DNA, and protein sequences. Models like StarCoder (Li et al., 2023a) and CodeLlama (Rozière et al., 2023), trained on vast code repositories and interactions, are adept at code generation, bug fixing, and explanation (Black et al., 2021; Wang & Komatsuzaki, 2021; Black et al., 2022; Wang et al., 2021; Chen et al., 2021; Li et al., 2022; Nijkamp et al., 2022; 2023; Fried et al., 2022; Gunasekar et al., 2023; Allal et al., 2023). Similarly, math-focused models, such as Minerva (Lewkowycz et al., 2022) and MathGLM (Yang et al., 2023), have been developed through specialized

training and fine-tuning strategies, including the use of external tools and Chain of Thought techniques (Jelassi et al., 2023; Liu & Low, 2023; Nye et al., 2022; Zhou et al., 2022a; Chen et al., 2022; Yang et al., 2023; Gao et al., 2023; Schick et al., 2023). These models, requiring extensive training, highlight the intensive data demands of LLMs in specialized domains. For example, CodeLlama uses 500 billion tokens for code training, 100 billion tokens for Python training, and more than 20 billion tokens for fine-tuning.

The Fusion of Large Language Models. Mixture-of-Experts (MoE) is the neural architecture that distributes tasks among multiple specialized networks (experts) and determines their responsibilities via a gating network (Jacobs et al., 1991). MoE enhances the capabilities of LLMs and has been extensively utilized (Clark et al., 2022; Lou et al., 2021; Kudugunta et al., 2021; Lepikhin et al., 2020; Mustafa et al., 2022; Zhou et al., 2022b; Riquelme et al., 2021; Shen et al., 2023b; Jiang et al., 2023b; Wan et al., 2024; Jiang et al., 2024). Many studies have endeavored to comprehend the Mixture-of-Experts (MoE) from the perspective of computational cost, with a specific focus on its sparse nature (Shazeer et al., 2016; Zoph et al., 2022; Zuo et al., 2021; Du et al., 2022; Fedus et al., 2022; Komatsuzaki et al., 2023; Shen et al., 2023a). The prevailing belief is that the MoE approach can scale up model parameters without incurring an escalation in computational expense. Some work suggests that experts do not necessarily have distinct expertise (Jiang et al., 2024), while other work verifies the effectiveness of expert specialization (Dai et al., 2024). We believe both ways could achieve promising performance, unlike those that train MoE models from scratch, this paper seeks to fuse highly specialized models in the fine-tuning phase. Compared to methods like knowledge distillation and knowledge fusion (Wan et al., 2024), our approach aims to achieve optimal performance by retaining the specialized models and learning to fuse the expertise directly, avoiding potential performance loss brought by inaccurate fashion weight estimation and further distillation training.

3. Our Approach

Compared to methods like Mixture-of-Experts (Shazeer et al., 2016), which expands the inner model structure to develop different expertise implicitly during training, our approach focuses on fusing specialist models explicitly aligned with different skill sets at the output level directly. This section first describes the constitution of the proposed model, ULTRAFUSER, and then introduces the construction of a mixed instruction tuning dataset, ULTRACHAT 2.

3.1. Model

The proposed fused model consists of three specialized models (termed as specialists), collectively denoted as $\mathcal{M}_\Theta =$

$\{E_{\text{text}}, E_{\text{code}}, E_{\text{math}}\}$, where E_{text} is mainly trained on natural language text, E_{code} is trained on programming code, and E_{math} is trained on mathematical problems. Each specialist model is essentially a native autoregressive large language model. They share the same architectural framework and vocabulary space but are trained on distinct datasets that are representative of their expertise.

Architecture. The fused model aims to utilize the expertise of each specialist model appropriately based on the nature of the input data. The integration of specialized ability is realized by a shared gating layer g that calculates the weight for each token per specialist. Specifically, during training, for the token $x^{(i)}$ concerned, the three specialists output token hidden states $\mathbf{h}^{(i)} = \{\mathbf{h}_{\text{text}}^{(i)}, \mathbf{h}_{\text{code}}^{(i)}, \mathbf{h}_{\text{math}}^{(i)}\}$ and corresponding logits $\mathbf{o}^{(i)} = \{\mathbf{o}_{\text{text}}^{(i)}, \mathbf{o}_{\text{code}}^{(i)}, \mathbf{o}_{\text{math}}^{(i)}\}$ as a native language model. Then, the gating layer g_Φ is applied to each set of specialist outputs to obtain the final logits.

Practically, the gating layer is implemented as a linear network that calculates the weight for each token $x^{(i)}$ based on the last hidden states $\mathbf{h}^{(i)} = E(x^{1:i-1})$. For each token $x^{(i)}$, the final output logits from the fused model are computed as:

$$g_\Phi(\mathcal{M}_\Theta(x^{(i)})) = \mathbf{w}^{(i)T}(\mathbf{o}_{\text{text}}^{(i)} : \mathbf{o}_{\text{code}}^{(i)} : \mathbf{o}_{\text{math}}^{(i)}), \quad (1)$$

where $\mathbf{w}^{(i)} = \text{Softmax}(g(\mathbf{h}_{\text{text}}^{(i)}) : g(\mathbf{h}_{\text{code}}^{(i)}) : g(\mathbf{h}_{\text{math}}^{(i)}))$

Training. One possible approach to training the model is to train the gating network only, expecting it to allocate each token to its optimal distribution over the three specialists. Such training strategy highly relies on the gating module’s capacity in capturing the complex and diverse context in drastically different instructions. An easier way to boost the performance is to jointly fine-tune the three specialists along with the gating module. However, the specialists can be negatively impacted by gradients back-propagated from the gating module due to its poor performance at the early stage, which may cause irreversible damage to the specialist’s inherent ability.

To tackle the problem, we propose a two-stage training strategy to ensure training stability and mitigate potential specialist ability loss. The first stage trains only the gating module parameters for N_1 steps and keeps specialists frozen. The purpose is to retain specialist capability while warming up the gating module. After the first stage of training, the gating network is expected to output reasonable token weights that favor over specific specialists according to data type. The second stage continues to fine-tune all model parameters based on the first stage for N_2 steps. At this stage, the specialist models are jointly optimized for a better overall performance. At both stages, the training loss is the cross-entropy loss given true labels y and the final model output.

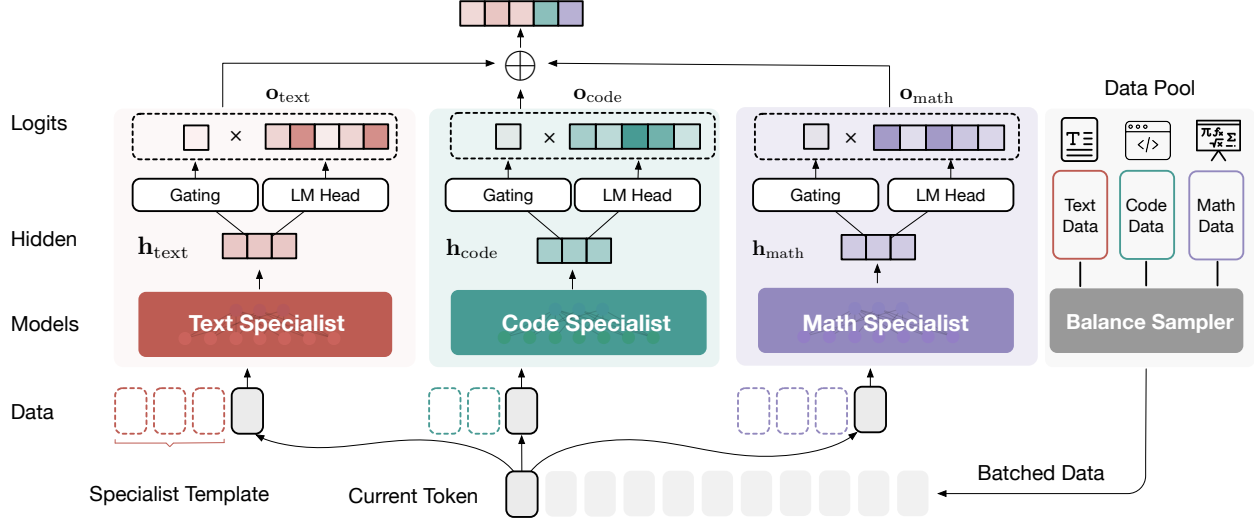


Figure 2: Architecture of our proposed ULTRAFUSER framework. We do not show the two-stage training in this illustration.

$$\mathcal{L}(x, y) = \sum_i \mathcal{CE}(g_\Phi(\mathcal{M}_\Theta(x^{(i)})), y^{(i)}). \quad (2)$$

The training proceeds by minimizing the total loss over all instances in the training set using a suitable optimization algorithm, such as AdamW. The gradients are back-propagated through both the specialist models and the gating networks, allowing the gating mechanism to learn how to distribute the inputs effectively among the specialists. The overall training process is shown in Algorithm 1.

Data-level Balancing. Since all specialists are well aligned to one specific type of instruction, they may demonstrate different activation patterns that are highly sensitive to inputs. Therefore, to fully take advantage of the specialized ability, we use specialist-specific templates to format our training data (see Appendix A). Each training sample is wrapped up by three different templates and fed into the respective specialist model. Since the loss is only calculated for the model response part, the response tokens will still be aligned, and their logits can be fused together seamlessly. We also adopt a batch-level class-balance sampler during training. The sampler ensures that each training batch contains the same number of training instances from the three categories, thus ensuring that the three specialists are activated and optimized at similar level for each batch, preventing from biased training that favor over one specific specialist. As shown in Algorithm 1, each batch of data contains $n \times 3$ instances in total. We explain the reason to alleviate the imbalance issue in the data-level and validate the effectiveness of the class-balance sampler in Section 4.4.

Inference. The model design adopts post-specialist token-

Algorithm 1 Algorithm for two-stage training with balanced data sampler, where $\mathcal{S}(\mathcal{D}, n)$ means randomly sampling n examples from dataset \mathcal{D} . N_1 and N_2 are total training steps, and η_1 and η_2 are the scheduled learning rate for the two stages, respectively.

Input: specialized models \mathcal{M}_Θ , gating g_Φ , training data $\mathcal{D}_{\text{text}}, \mathcal{D}_{\text{code}}, \mathcal{D}_{\text{math}}$
for $i = 1$ **to** N_1 **do**
 $\mathcal{D}^i = \bigcup_{t \in \{\text{text}, \text{code}, \text{math}\}} \mathcal{D}_t^i$
 $\mathcal{D}^i = \bigcup_{t \in \{\text{text}, \text{code}, \text{math}\}} \mathcal{S}(\mathcal{D}_t, n)$
 $g_\Phi = g_\Phi - \eta_1 \Delta_\Phi \frac{1}{|\mathcal{D}^i|} \sum_{(x, y) \in \mathcal{D}^i} \mathcal{L}(x, y)$
end for
for $j = 1$ **to** N_2 **do**
 $\mathcal{D}^j = \bigcup_{t \in \{\text{text}, \text{code}, \text{math}\}} \mathcal{D}_t^j$
 $\mathcal{D}^j = \bigcup_{t \in \{\text{text}, \text{code}, \text{math}\}} \mathcal{S}(\mathcal{D}_t, n)$
 $g_\Phi = g_\Phi - \eta_2 \Delta_\Phi \frac{1}{|\mathcal{D}^j|} \sum_{(x, y) \in \mathcal{D}^j} \mathcal{L}(x, y)$
 $\mathcal{M}_\Theta = \mathcal{M}_\Theta - \eta_2 \Delta_\Theta \frac{1}{|\mathcal{D}^j|} \sum_{(x, y) \in \mathcal{D}^j} \mathcal{L}(x, y)$
end for

level gating, meaning that all specialists are activated during inference. For each token $x^{(i)}$, the three specialist models \mathcal{M}_Θ are queried, and their logits are fused using the gating module $g_\Phi(\cdot)$ as in the training phase. The softmax is applied to the aggregated logits to generate probabilities for the next token. The selected token is then used as part of the input for the subsequent inference step in an autoregressive manner. Our design opens doors for sophisticated, real-time adaptability that monolithic models lack. For example, in a text string interwoven with mathematical equations and code snippets—common in scientific papers, the fused model can shift its “attention” between specialists within the same se-

quence, ensuring that each token is treated with the most appropriate domain expertise. But on the other hand, since all specialists are activated in inference, computational overheads are inevitably introduced. In experiments, we adapt the vLLM project (Kwon et al., 2023) to our fused model to accelerate inference, which is elaborated in Appendix B.

Why Not Sample-Level? One direct and simple approach to fusing specialized models is to train them in a sample-level manner. That is, freezing the specialist models and directly train a selector, letting one specialist respond to a whole query. This approach seems to safeguard the lower-bound performance for the model effectively, so why does this paper opt for token-level training rather than sample-level? The main reason is that, although this paper categorizes the data into three distinct symbolic systems, they may blend together in real-world queries (for instance, code data may contain extensive text intended for documentation). Similarly, while these three capabilities might weaken each other in some respects, they could also enhance one another in different contexts, which is demonstrated in Section 4.2. We choose to design the fused model to seek a higher performance ceiling.

3.2. ULTRACHAT 2 Dataset

Currently, within the open-source community, there are already multiple instruction-tuning datasets for text-based conversations. However, there is a relatively limited amount of systematic code and mathematical instruction tuning data available. In this section, we construct ULTRACHAT 2, a comprehensive dataset tailored for training our proposed model. ULTRACHAT 2 spans a wide range of subject matter, covering natural language, coding, and mathematical instructions. We employ a multi-stage pipeline to generate a rich set of instructional data. First, we engage in multi-turn interactions with GPT-4, constructing meta-topics that best represent each domain. Then, each meta-topic is utilized to generate multiple sub-topics. For each sub-topic, LLM is tasked with generating diverse and informative specific instructions. After obtaining these instructions, we continue with in-context learning, generating both strong and weakly related instructions for each directive to fully leverage LLM’s generalization capabilities. Finally, we extract 30% of the instruction data and make them more complex. Once we have the complete pool of instructions, we use GPT-4 to respond to these instructions, resulting in ULTRACHAT 2.

Data Analysis. We randomly sample 5000 instructions from each category and visualize the data distribution in Figure 3. The representations are obtained by averaging the last layer of hidden states of each token from Llama-2-13B, and dimensions are further reduced by the t-SNE algorithm (Van der Maaten & Hinton, 2008). The

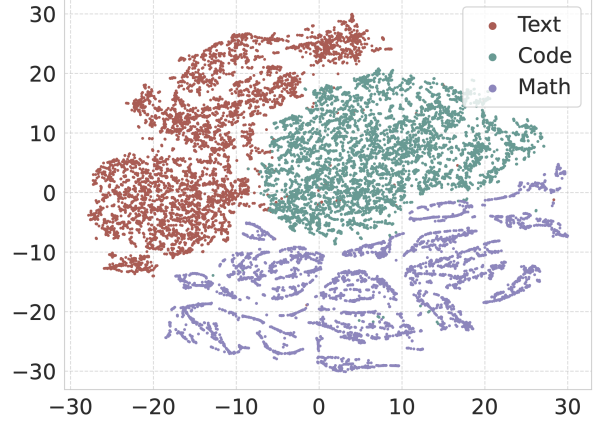


Figure 3: t-SNE visualization of ULTRACHAT 2 dataset.

visualization clearly demonstrates the diversity and distinctiveness of different types of ULTRACHAT 2, which aligns with the intuition and echos the discussion in Section 1. ULTRACHAT 2 provides high-quality resources for the facilitation of specialized models. We train a Llama-2-13B on ULTRACHAT 2 to give a glance at the effectiveness. As shown in Figure 4, in the text domain, the Llama-2-13B + UltraChat 2 configuration exhibits a 3.9% decrement in performance relative to the baseline that is only trained on the text domain. Conversely, in the code domain, there is a significant performance increment of 10.4% with the UltraChat 2 enhancement. The math domain also shows a performance increase of 9.4% with the UltraChat 2 integration, indicating a clear advantage of the updated system in code-related and mathematical reasoning tasks.

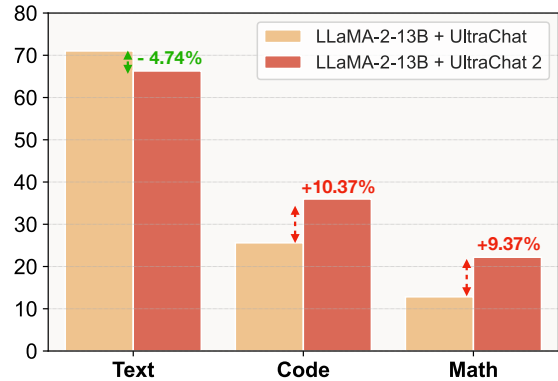








Figure 4: Performance comparison of Llama-2 model trained on ULTRACHAT and ULTRACHAT 2. The performance for the text domain is computed by the average results on TruthfulQA (Acc) and AlpacaEval (Win Rate) datasets; the performance for the code domain is Pass@1 of HumanEval; and the performance for the math domain is the average result of GSM8K (Pass@1) and MATH (Pass@1).

Table 1: Statistics and information of ULTRACHAT 2 dataset. # Topics are the number of meta-topics and sub-topics.

	Text Part	Code Part	Math Part
# Data	100,000	100,000	110,000
# Topics	30/1100	21/407	21/80
Examples	 Technology Artificial Intelligence Smartphone Quantum Computing	 Web Development HTML Basics Javascript Essentials Web Security	 Algebra Polynomials Factoring Quadratic Equations
	 Education Inclusive education Classroom management Critical thinking	 Mobile App Development User Interface Design Responsive Design Database Management	 Discrete Mathematics Graph Theory Combinatorics Number Theory

4. Experiments

We conduct extensive experiments to analyze the effectiveness and behaviors of ULTRAFUSER. Implementation details are reported in Appendix A. Our models, data, training and inference frameworks will be publicly released.

4.1. Experimental Settings

Models. To validate the effectiveness of our approach, we adopt Llama-2-13B (Touvron et al., 2023) as the backbone for experiments. Specifically, we use UltraLM-13B-V2.0 (Ding et al., 2023), CodeLlama-13B-instruct (Rozière et al., 2023), WizardMath-13B-V1.0 (Luo et al., 2023a) as the three specialist models. All model parameters are fine-tuned under the proposed ULTRAFUSER framework.

Evaluation. For the text domain, we use TruthfulQA (Lin et al., 2021) and AlpacaEval (Li et al., 2023b) for evaluation. The former is more focused on the truthfulness of LLMs, and the latter consists of more general natural language questions. For the code domain, we use HumanEval (Chen et al., 2021) for evaluation, which is a code completion task. For the math domain, we use GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2020), SAT-MATH (Zhong et al., 2023) and AQuA (Ling et al., 2017) for evaluation. For evaluation, we transform each dataset into instruction format, and use consistent template as training for inference. Specifically, we evaluate under the MC2 setting in TruthfulQA, where each option is fed into the model independently and the model is queried for true or false judgment. For HumanEval, we use InstructHumanEval that transforms the original dataset into instruction format. All results are zero-shot and produced by our experiments. We do not use any chain-of-thought (CoT) techniques to boost the performance.

4.2. Results on Benchmarks

Our approach involves further training already highly specialized models, but to what extent can this retraining be effective without the ULTRAFUSER framework? We provide two sets of baselines. The first set is the specialists of three domains, and the second group is the further trained versions of such specialists with identical training data as ULTRAFUSER. Comparing to the original specialist models as shown in Table 2, ULTRAFUSER can consistently produce on-a-par or even superior performance across benchmarks from different domains and achieves the highest overall results. Notably, ULTRAFUSER significantly outperforms respective specialists on TruthfulQA and HumanEval datasets by 5.86% and 4.25%, indicating that the three specialist models interact with each other in helpful ways to boost performance on more comprehensive datasets. The result demonstrates the effectiveness of directly fusing specialist models with the proposed framework in both retaining and potentially synthesizing expertise to achieve even better performance.

Furthermore, directly fine-tuning specialist models on our training data may not produce desirable performance, as shown in Figure 5. Although training on a mixed dataset indeed boosts other expertise domains, it also severely harms the original expertise of the model concerned, which can be seen from the performance drop after further tuning. Meanwhile, it should be aware that the three models not only differ in expertise but also in the training stages they each have come through. Different specialist models show distinct patterns of expertise distribution after further tuning. Models like UltraLM and WizardMath, which are directly instruction-tuned based on Llama, gain more benefit from further tuning. It should be noted that both UltraLM and WizardMath improve significantly on coding tasks after further tuning (14.6% and 15.8%), while UltraLM’s perfor-

Table 2: Results of baselines and our proposed models across different benchmarks. **All the numbers are zero-shot results** produced by our experiments under the same inference framework. *No* Chain-of-thought (CoT) techniques are employed in evaluation. The highest results are **bold**, and the second highest results are underlined. Delta values in the second block mean the performance differences between the model and the corresponding specialist models.

Model	TruthfulQA Acc (%)	AlpacaEval Win Rate (%)	HumanEval Pass@1 (%)	GSM8K Pass@1 (%)	MATH Pass@1 (%)	SAT-Math Acc (%)	AQuA Acc (%)	Avg.
UltraLM-2	<u>58.82</u>	83.23	25.61	25.09	4.48	25.00	25.98	35.46
CodeLlama	56.89	69.21	<u>48.78</u>	23.12	6.16	27.73	25.59	<u>36.78</u>
WizardMath	26.81	51.50	10.98	56.18	12.2	<u>29.55</u>	<u>22.05</u>	29.91
ULTRAFUSER	64.67	<u>82.35</u>	53.03	<u>54.59</u>	<u>11.36</u>	30.00	26.38	47.48

mance on GSM8K and MATH doubles after tuning. However, CodeLlama’s performance, unfortunately, degrades on every benchmark, especially in complex instruction following tasks like Alpaca. It is probably because CodeLlama has undergone thorough code infilling pre-training (500B tokens) before instruction fine-tuning. The model parameters may be detrimentally disrupted when directly tuning on drastically different domains of data like pure text and mathematical notation, and therefore suffer from great loss of basic instruction following ability. This points to the fact that the outcome of continuously fine-tuning a well-aligned model highly relies on the previous training data schedule and training strategy adopted, while the proposed framework could seamlessly bridge distinctive model expertise with simple tuning methods and mixed data. Among the three specialists, UltraLM-2 seems to benefit the most in terms of overall performance after further tuning, indicating that a “specialist” in text may be equipped with much broader expertise and may have more potential in expanding new expertise from further supervised fine-tuning. Such phenomenon is intuitively correct as text domain does incorporate much more diversity compared to the other two, but it also sheds light on the fact that there is much more fine-grained variability that should be captured by the term “expertise”. As we only distinguish the “expertise” based on the difference in symbol system, further disentanglement and more experiments are worth exploration. Above all, our experimental results show that ULTRAFUSER can easily produce a well-rounded system by fusing and leveraging expertise from different specialist models.

4.3. Results on MT-Bench

4.4. Ablation Study

Training fused models could cause load imbalance, leading to the collapse of the routing mechanism. A typical approach to mitigate this issue in MoE is to introduce a balance loss to prevent certain models from being over-selected or under-selected. In our framework, we do not introduce explicit balance loss based on a simple hypothesis: A model that has been highly specialized can automatically produce

a lower loss on the data it is good at. Now that the model already has data that is good at processing, we hope to solve the problem from the data level, not force the specialist models to participate to a certain extent during the calculation. We find that designing some training methods can make the progress more stable. Two key components of our framework are *two-stage training* and *balanced sampler*. The former plays a role similar to warm-up, allowing the randomly initialized gating module to adapt to the current expert model. The latter, as mentioned, ensures load balance at the data level. In Table 4, we report the best performance under each training strategy. It can be observed that the beneficial effects of these two modules are obvious, and their use has made the overall performance of the fused model improve considerably. We further investigate the impact on the training stability of the balance sampler. We train two versions of the model with the same dataset and sample 12 checkpoints, respectively, from 2000 steps to 9000 steps, and conduct evaluations. As shown in 5, with the help of the balance sampler, the fused model could achieve superior performance and lower standard deviations on all datasets. GSM8K is relatively stable during training, however HumanEval may face larger fluctuations.

4.5. Expertise Analysis

In our training, there is no explicit mechanism to make certain specialists “pay more attention” to the corresponding data. But as mentioned, we expect that the specializations could still be separated, and a type of data will receive different gating weights from the specialist models. We randomly sample 100 data instances from the three domains and conduct analysis by directly going through the inference to the fused model, and calculate the weight from three specialist models of each token. Table 6 shows the average weights of all the tokens in each set of data from three specialist models. And intuitively, each set of data is primarily driven by the corresponding specialist model. The prominence of code data is evident, with the corresponding expert models significantly outweighing the other two models. In mathematical data, code and mathematical models almost equally

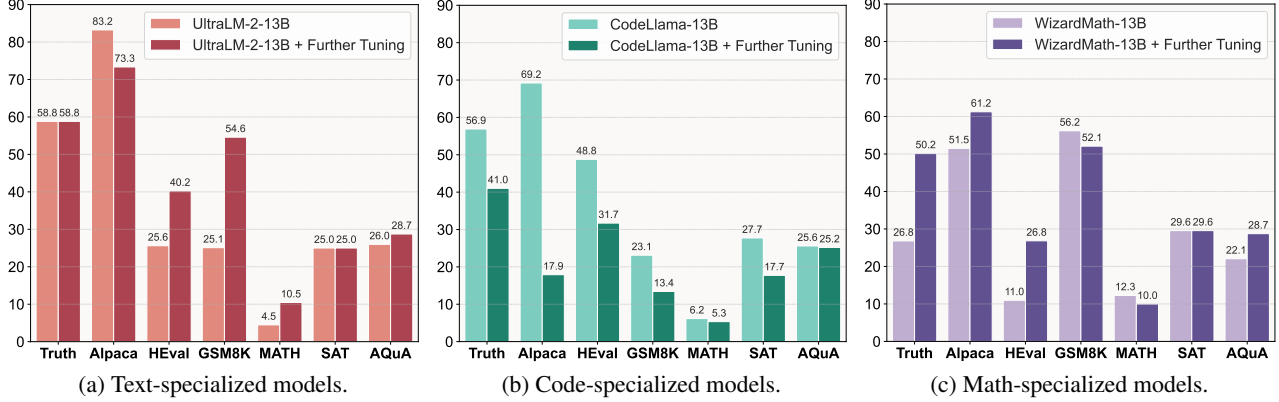


Figure 5: Performance comparisons between specialist models and the further training versions of them.

Model	Writing	Roleplay	Reasoning	Math	Coding	Extraction	STEM	Humanities	Overall
UltraLM-2	8.83	7.98	5.20	2.90	4.00	6.74	8.08	9.46	6.62
CodeLlama	5.80	7.10	3.80	3.05	3.43	5.36	5.65	7.05	5.16
WizardMath	7.75	7.03	4.80	3.85	3.50	4.65	7.65	9.13	6.04
UltraLM-2+Further Tune	7.85	7.60	4.30	4.48	5.20	5.78	8.32	9.40	<u>6.62</u>
CodeLlama+Further Tune	7.33	6.60	3.85	2.25	3.68	5.20	4.68	5.10	4.84
WizardMath+Further Tune	7.18	6.90	4.95	4.25	4.55	5.18	7.55	7.98	6.07
ULTRAFUSER	<u>8.60</u>	8.11	<u>5.00</u>	5.15	<u>5.10</u>	<u>6.53</u>	<u>8.23</u>	<u>9.43</u>	7.02

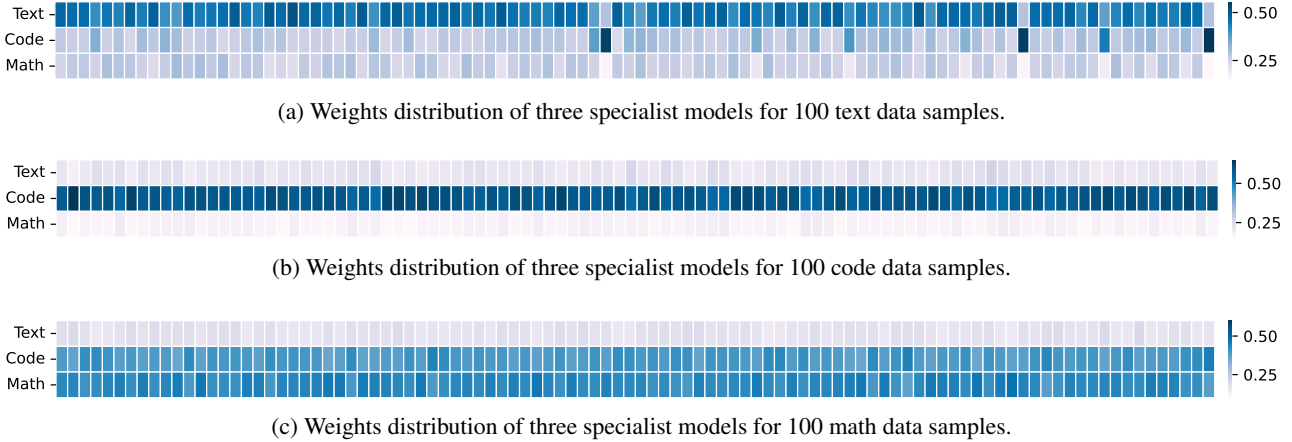
 Table 3: Performance on MT-Bench. The highest results are **bold**, and the second highest results are underlined.


Figure 6: Weight distributions of 300 data samples from text, code, and math domains. Each column is a data point, and each row is the average weight of one specialist model. The darker the color, the more average weight the model gives to the tokens of this data point.

Table 4: Results across TruthfulQA (Truth), HumanEval (H-Eval), and GSM8K with different training strategies.

Strategy	Truth	H-Eval	GSM8K
Direct Training	51.17	46.95	<u>53.83</u>
+ Two-Stage	<u>61.72</u>	<u>50.00</u>	52.69
+ Two-Stage + Balanced	64.67	53.05	54.59

Table 5: Mean results and standard deviation over 12 checkpoints with and without the balance sampler (two-stage training are both applied).

Strategy	Truth	H-Eval	GSM8K
w/o Balance	57.54±2.80	48.27±4.92	52.91±1.76
w/ Balance	59.91±1.96	53.68±2.74	53.77±1.74

Table 6: Average weights from three specialist models of different data.

Avg. Weight	Text Data	Code Data	Math Data
w_{text}	0.45	0.23	0.18
w_{code}	0.29	0.59	0.39
w_{math}	0.26	0.18	0.43

dominate inference, with a marginal difference. This is more distinctly observable in the sample-level distribution illustrated in Figure 6. Despite the fusion and further training of the models, it’s evident that these specialized models still retain their original functionalities and are now capable of synergistic performance.

5. Conclusion

This paper aims to integrate coding and mathematical reasoning capabilities into a general language model with as little loss as possible. We present ULTRAFUSER, a simple framework to train high-specialized models with a token-level gating mechanism and a two-stage balanced training strategy. Accompanied by the goal, we construct a high-quality and diverse instruction tuning dataset, ULTRACHAT 2, that contains 300,000 instructions and responses from 3 domains, 72 meta-topics, and 1587 sub-topics. Our experiments demonstrate the effectiveness of the proposed framework by showing that fused models can be performative simultaneously in text understanding, code generation, and mathematical reasoning. In future work, the proposed ULTRAFUSER can also be adapted to domains beyond the mentioned ones. For example, by using to fuse language models that are specialized in different languages.

Broader Impact

This work aims to integrate professional code generation and mathematical reasoning capabilities into a general interactive language model, proposing a framework to fuse highly specialized models and construct a corresponding dataset. The goal of the study is to further advance and leverage the research of open-source large language models, providing a potential solution to bring together specialized language models into an integrated and comprehensive system. The study, therefore, shares many of the potential impacts with other research on LLMs, e.g., hallucination and biased output. The model outputs should be closely supervised before put into use. On the other hand, this work sheds light on the possibility of fusing language models that are trained on drastically different datasets and may benefit future applications of large language models. How to better define and choose “expertise” is also worth exploration under the framework of ULTRAFUSER. As for ULTRACHAT 2, we do not use any real human queries or responses to

protect privacy. In our manual assessment, the code and math parts of the data are mostly unbiased, but there may still be biased statements in the text part.

References

- Agarap, A. F. Deep learning using rectified linear units (relu). *ArXiv*, abs/1803.08375, 2018. URL <https://api.semanticscholar.org/CorpusID:4090379>.
- Allal, L. B., Li, R., Kocetkov, D., Mou, C., Akiki, C., Ferlandis, C. M., Muennighoff, N., Mishra, M., Gu, A., Dey, M., et al. Santacoder: don’t reach for the stars! *arXiv preprint arXiv:2301.03988*, 2023.
- Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Bansal, R., Samanta, B., Dalmia, S., Gupta, N., Vashishth, S., Ganapathy, S., Bapna, A., Jain, P., and Talukdar, P. Llm augmented llms: Expanding capabilities through composition. *arXiv preprint arXiv:2401.02412*, 2024.
- Black, S., Gao, L., Wang, P., Leahy, C., and Biderman, S. GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow, March 2021. URL <https://doi.org/10.5281/zenodo.5297715>. If you use this software, please cite it using these metadata.
- Black, S., Biderman, S., Hallahan, E., Anthony, Q., Gao, L., Golding, L., He, H., Leahy, C., McDonell, K., Phang, J., et al. Gpt-neox-20b: An open-source autoregressive language model. In *Proceedings of BigScience Episode# 5-Workshop on Challenges & Perspectives in Creating Large Language Models*, pp. 95–136, 2022.
- Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosse-lut, A., Brunskill, E., et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. In *Proceedings of NeurIPS*, 2020.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. d. O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Chen, W., Ma, X., Wang, X., and Cohen, W. W. Program of thoughts prompting: Disentangling computation from

- reasoning for numerical reasoning tasks. *arXiv preprint arXiv:2211.12588*, 2022.
- Chiang, W.-L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J. E., Stoica, I., and Xing, E. P. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- Clark, A., De Las Casas, D., Guy, A., Mensch, A., Paganini, M., Hoffmann, J., Damoc, B., Hechtman, B., Cai, T., Borgeaud, S., et al. Unified scaling laws for routed language models. In *International Conference on Machine Learning*, pp. 4057–4086. PMLR, 2022.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Daheim, N., Möllenhoff, T., Ponti, E. M., Gurevych, I., and Khan, M. E. Model merging by uncertainty-based gradient matching, 2023.
- Dai, D., Deng, C., Zhao, C., Xu, R., Gao, H., Chen, D., Li, J., Zeng, W., Yu, X., Wu, Y., et al. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models. *arXiv preprint arXiv:2401.06066*, 2024.
- Ding, N., Chen, Y., Xu, B., Qin, Y., Zheng, Z., Hu, S., Liu, Z., Sun, M., and Zhou, B. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*, 2023.
- Du, N., Huang, Y., Dai, A. M., Tong, S., Lepikhin, D., Xu, Y., Krikun, M., Zhou, Y., Yu, A. W., Firat, O., et al. Glam: Efficient scaling of language models with mixture-of-experts. In *International Conference on Machine Learning*, pp. 5547–5569. PMLR, 2022.
- Fedus, W., Zoph, B., and Shazeer, N. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *The Journal of Machine Learning Research*, 23(1):5232–5270, 2022.
- Fried, D., Aghajanyan, A., Lin, J., Wang, S., Wallace, E., Shi, F., Zhong, R., Yih, S., Zettlemoyer, L., and Lewis, M. Incoder: A generative model for code infilling and synthesis. In *The Eleventh International Conference on Learning Representations*, 2022.
- Gao, L., Madaan, A., Zhou, S., Alon, U., Liu, P., Yang, Y., Callan, J., and Neubig, G. Pal: Program-aided language models. In *International Conference on Machine Learning*, pp. 10764–10799. PMLR, 2023.
- Gunasekar, S., Zhang, Y., Aneja, J., Mendes, C. C. T., Del Giorno, A., Gopi, S., Javaheripi, M., Kauffmann, P., de Rosa, G., Saarikivi, O., et al. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*, 2023.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., and Hinton, G. E. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.
- Jelassi, S., d’Ascoli, S., Domingo-Enrich, C., Wu, Y., Li, Y., and Charton, F. Length generalization in arithmetic transformers. *arXiv preprint arXiv:2306.15400*, 2023.
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. Mistral 7b, 2023a.
- Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D. S., Casas, D. d. l., Hanna, E. B., Bressand, F., et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- Jiang, D., Ren, X., and Lin, B. Y. Llm-blender: Ensembling large language models with pairwise ranking and generative fusion. *arXiv preprint arXiv:2306.02561*, 2023b.
- Komatsuzaki, A., Puigcerver, J., Lee-Thorp, J., Ruiz, C. R., Mustafa, B., Ainslie, J., Tay, Y., Dehghani, M., and Houlisby, N. Sparse upcycling: Training mixture-of-experts from dense checkpoints, 2023.
- Kudugunta, S., Huang, Y., Bapna, A., Krikun, M., Lepikhin, D., Luong, M.-T., and Firat, O. Beyond distillation: Task-level mixture-of-experts for efficient inference. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 3577–3599, 2021.
- Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J. E., Zhang, H., and Stoica, I. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.

- Lepikhin, D., Lee, H., Xu, Y., Chen, D., Firat, O., Huang, Y., Krikun, M., Shazeer, N., and Chen, Z. Gshard: Scaling giant models with conditional computation and automatic sharding. *arXiv preprint arXiv:2006.16668*, 2020.
- Lewkowycz, A., Andreassen, A. J., Dohan, D., Dyer, E., Michalewski, H., Ramasesh, V. V., Slone, A., Anil, C., Schlag, I., Gutman-Solo, T., et al. Solving quantitative reasoning problems with language models. In *Advances in Neural Information Processing Systems*, 2022.
- Li, R., Allal, L. B., Zi, Y., Muennighoff, N., Kocetkov, D., Mou, C., Marone, M., Akiki, C., Li, J., Chim, J., et al. Starcoder: may the source be with you! *arXiv preprint arXiv:2305.06161*, 2023a.
- Li, X., Zhang, T., Dubois, Y., Taori, R., Gulrajani, I., Guestrin, C., Liang, P., and Hashimoto, T. B. AlpacaEval: An automatic evaluator of instruction-following models, 2023b.
- Li, Y., Choi, D., Chung, J., Kushman, N., Schrittwieser, J., Leblond, R., Eccles, T., Keeling, J., Gimeno, F., Dal Lago, A., et al. Competition-level code generation with alpha-code. *Science*, 378(6624):1092–1097, 2022.
- Lin, S., Hilton, J., and Evans, O. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- Ling, W., Yogatama, D., Dyer, C., and Blunsom, P. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 158–167, Vancouver, Canada, July 2017. Association for Computational Linguistics.
- Liu, T. and Low, B. K. H. Goat: Fine-tuned llama outperforms gpt-4 on arithmetic tasks. *arXiv preprint arXiv:2305.14201*, 2023.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2017. URL <https://api.semanticscholar.org/CorpusID:53592270>.
- Lou, Y., Xue, F., Zheng, Z., and You, Y. Cross-token modeling with conditional computation. *arXiv preprint arXiv:2109.02008*, 2021.
- Luo, H., Sun, Q., Xu, C., Zhao, P., Lou, J., Tao, C., Geng, X., Lin, Q., Chen, S., and Zhang, D. Wizard-math: Empowering mathematical reasoning for large language models via reinforced evol-instruct. *arXiv preprint arXiv:2308.09583*, 2023a.
- Luo, Z., Xu, C., Zhao, P., Sun, Q., Geng, X., Hu, W., Tao, C., Ma, J., Lin, Q., and Jiang, D. Wizardcoder: Empowering code large language models with evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023b.
- Mustafa, B., Riquelme, C., Puigcerver, J., Jenatton, R., and Hounsby, N. Multimodal contrastive learning with limoe: the language-image mixture of experts. *Advances in Neural Information Processing Systems*, 35:9564–9576, 2022.
- Nijkamp, E., Pang, B., Hayashi, H., Tu, L., Wang, H., Zhou, Y., Savarese, S., and Xiong, C. Codegen: An open large language model for code with multi-turn program synthesis. In *The Eleventh International Conference on Learning Representations*, 2022.
- Nijkamp, E., Hayashi, H., Xiong, C., Savarese, S., and Zhou, Y. Codegen2: Lessons for training llms on programming and natural languages. *arXiv preprint arXiv:2305.02309*, 2023.
- Nye, M., Andreassen, A. J., Gur-Ari, G., Michalewski, H., Austin, J., Bieber, D., Dohan, D., Lewkowycz, A., Bosma, M., Luan, D., et al. Show your work: Scratchpads for intermediate computation with language models. In *Deep Learning for Code Workshop*, 2022.
- OpenAI. Gpt-4 technical report. *arXiv*, 2023a.
- OpenAI. Gpt-4 technical report, 2023b.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- Riquelme, C., Puigcerver, J., Mustafa, B., Neumann, M., Jenatton, R., Susano Pinto, A., Keyzers, D., and Hounsby, N. Scaling vision with sparse mixture of experts. *Advances in Neural Information Processing Systems*, 34: 8583–8595, 2021.
- Rozière, B., Gehring, J., Gloeckle, F., Sootla, S., Gat, I., Tan, X. E., Adi, Y., Liu, J., Remez, T., Rapin, J., et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
- Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Zettlemoyer, L., Cancedda, N., and Scialom, T. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023.
- Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., and Dean, J. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference on Learning Representations*, 2016.

- Shen, S., Hou, L., Zhou, Y., Du, N., Longpre, S., Wei, J., Chung, H. W., Zoph, B., Fedus, W., Chen, X., Vu, T., Wu, Y., Chen, W., Webson, A., Li, Y., Zhao, V., Yu, H., Keutzer, K., Darrell, T., and Zhou, D. Mixture-of-experts meets instruction tuning: a winning combination for large language models, 2023a.
- Shen, S., Yao, Z., Li, C., Darrell, T., Keutzer, K., and He, Y. Scaling vision-language models with sparse mixture of experts. *arXiv preprint arXiv:2303.07226*, 2023b.
- Stoica, G., Bolya, D., Bjorner, J., Ramesh, P., Hearn, T., and Hoffman, J. Zipit! merging models from different tasks without training, 2024.
- Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., and Hashimoto, T. B. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models*, 2023.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Van der Maaten, L. and Hinton, G. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- Wan, F., Huang, X., Cai, D., Quan, X., Bi, W., and Shi, S. Knowledge fusion of large language models. *arXiv preprint arXiv:2401.10491*, 2024.
- Wang, B. and Komatsuzaki, A. Gpt-j-6b: A 6 billion parameter autoregressive language model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021. URL <https://github.com/kingoflolz/mesh-transformer-jax>.
- Wang, Y., Wang, W., Joty, S., and Hoi, S. C. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 8696–8708, 2021.
- Xu, C., Sun, Q., Zheng, K., Geng, X., Zhao, P., Feng, J., Tao, C., and Jiang, D. Wizardlm: Empowering large language models to follow complex instructions, 2023.
- Yang, Z., Ding, M., Lv, Q., Jiang, Z., He, Z., Guo, Y., Bai, J., and Tang, J. Gpt can solve mathematical problems without a calculator. *arXiv preprint arXiv:2309.03241*, 2023.
- Yue, X., Qu, X., Zhang, G., Fu, Y., Huang, W., Sun, H., Su, Y., and Chen, W. Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint arXiv:2309.05653*, 2023.
- Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X. V., et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- Zhong, W., Cui, R., Guo, Y., Liang, Y., Lu, S., Wang, Y., Saied, A., Chen, W., and Duan, N. Agieval: A human-centric benchmark for evaluating foundation models. *arXiv preprint arXiv:2304.06364*, 2023.
- Zhou, H., Nova, A., Larochelle, H., Courville, A., Neyshabur, B., and Sedghi, H. Teaching algorithmic reasoning via in-context learning. *arXiv preprint arXiv:2211.09066*, 2022a.
- Zhou, Y., Lei, T., Liu, H., Du, N., Huang, Y., Zhao, V. Y., Dai, A. M., Chen, Z., Le, Q. V., and Laudon, J. Mixture-of-experts with expert choice routing. In *Advances in Neural Information Processing Systems*, 2022b.
- Zoph, B., Bello, I., Kumar, S., Du, N., Huang, Y., Dean, J., Shazeer, N., and Fedus, W. St-moe: Designing stable and transferable sparse expert models. *arXiv preprint arXiv:2202.08906*, 2022.
- Zuo, S., Liu, X., Jiao, J., Kim, Y. J., Hassan, H., Zhang, R., Zhao, T., and Gao, J. Taming sparsely activated transformer with stochastic experts. *arXiv preprint arXiv:2110.04260*, 2021.

A. Implementation Details

The gating module is implemented as a two-layer linear model with ReLU (Agarap, 2018) activation in between. The hidden size of the module is set according to the hidden size of the specialized models. The gating layer is trained for $N_1 = 400$ steps at the first training stage with sample size $n = 64$ for all experiments and learning rate $\eta_1 = 2e - 5$ is used with a cosine scheduler. For the second stage with Llama backbone, we use $\eta_2 = 2e - 5$, sample size $n = 32$ with cosine scheduler. Note that our framework requires the consistent tokenization strategy across all specialist models. Therefore, we use the original LLaMA-2-13B tokenizer for ULTRAFUSER training. All experiments are conducted on $8 \times 80\text{GB A100 GPUs}$ and use AdamW optimizer (Loshchilov & Hutter, 2017). Apart from the curated ULTRACHAT 2, we also employ extra instruction tuning datasets from both math and code domains to enrich instructional format diversity. Specifically, we use the Evol-Instruct dataset (Luo et al., 2023b;a) for programming and the MathInstruct training set (Yue et al., 2023) for math problems. We conduct comprehensive search and filtering (13 grams) to avoid data contamination.

Table 7 and Table 8 show the conversation templates we use for each specific specialist model and the prompt for converting datasets to instructions in evaluation. In training, each example is wrapped by three different conversation templates and fed into the respective model. In inference, before applying the conversation template, dataset-specific prompt is used to wrap the example first (if applicable).

Model	Conversation Template
UltraLM-2	User: {instruction}\nAssistant:
CodeLlama	<s>[INST] {instruction} [/INST]
WizardMath	Below is an instruction that describes a task. Write a response that appropriately completes the request.\n\n### Instruction:\n{instruction}\n\n### Response:

Table 7: Model-specific conversation templates for training and evaluation.

Dataset	Evaluation Prompt
TruthfulQA	Judge the correctness of a given answer. Question: {question}\n Answer: {answer}\n Is the answer correct? Return Yes or No.
Alpaca	Please give helpful, very detailed, and polite answer to the user’s question below.\n Question: {question}

Table 8: Dataset-specific prompts used for evaluation.

B. Efficient Inference

We implement the inference of our fused model on the existing inference framework, vLLM (Kwon et al., 2023). Unlike other MoE models supported by vLLM, such as Mixtral (Jiang et al., 2024), our fused model requires different input prompts and the maintenance of multiple key-value caches within multiple models. Modifying the model implementation within vLLM directly to accommodate these requirements can be complex and may conflict with the PageAttention mechanism (Kwon et al., 2023) due to the use of multiple key-value caches. Therefore, we instead partition the GPU memory into several parts, each running a single model using a vLLM instance, and then fusing the output to form a fused model.

vLLM inherently supports streaming output, which returns tokens to the user-end token-by-token, and each token is produced by a sampler function applied on the hidden logits of the LLM. We change the implementation: in each iteration, we return the hidden logit instead of the token:

```
# In model implementation
# change from outputting token = self.model.sample(hidden_states, sampling_metadata) to
return {
    "sampler": self.model.sample,
    "data": {
        "hidden_states": hidden_states,
```

```

        "sampling_metadata": sampling_metadata,
    }
}

```

This allows us to pause the model generation, giving us control over when to predict the next token and when to continue generating future tokens. We then make the model instances communicate and fuse the logits:

```

logits = [
    llm.llm_engine.step()
    for llm in llms
] # get logits for different LLMs
fused_logits = fuse_function(*[logit["data"] for logit in logits]) # apply fuse function

```

The next token is predicted and sampled using the fused output, and we control the model instances to resume generation.

C. Discussion and Limitations

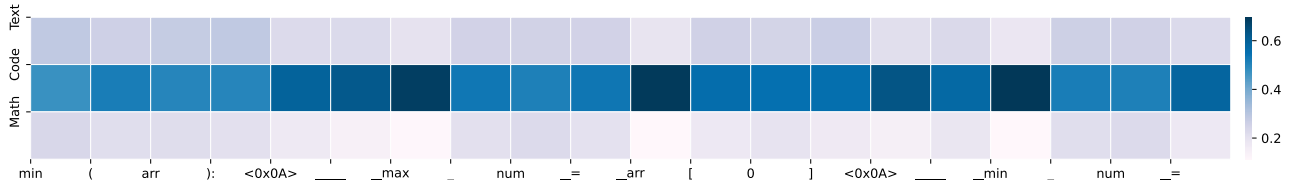
We regard the data distribution in training language models into three domains in this study according to the symbol systems and achieve promising empirical results in our experiments. However, the realistic situation is far more sophisticated. In the field of “text domain” alone, there are different tasks such as common sense knowledge, specialized knowledge, natural language reasoning, etc., not to mention the existence of multilingualism. Our fused model may yield less favorable results on other benchmarks. In our training, no explicit selection mechanism is introduced in order to make the method scalable (force specialist models to process certain types of data). We believe finer-grained models could be trained under the spirit of ULTRAFUSER; that is, the number of specialists is not necessarily three, and the domains are also necessarily divided as the same as the paper. For example, other symbol systems (like DNA sequences) may also be integrated into the framework.

The ULTRACHAT 2 dataset is fully synthetically generated by GPT-4 and fully excludes human engagement. Besides efficiency and privacy benefits, the factuality and trustworthiness of generated content can not be guaranteed. Our study mainly focuses on text, code, and math capabilities, neglecting some other important properties of LLMs, such as safety and multilinguality. We believe our proposed approach is generally applicable to address these limitations and will devote to developing more advanced methods.

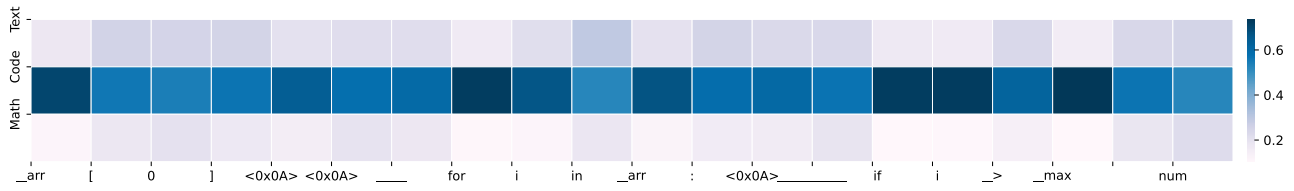
Comparing to the line of works on model merging that manipulates the inner parameters of existing models in either supervised or unsupervised manner (Daheim et al., 2023; Stoica et al., 2024; Wan et al., 2024; Bansal et al., 2024), our framework tackles the problem in a more straightforward way by directly merging the output and training with mixed high-quality instructional dataset to further adapt the model. The proposed framework follows the spirit of instruction tuning, and the training is conducted with direct supervised fine-tuning. Employing a diverse set of instruction datas, we show that the resulting model is equipped with desirable expertise and generalizes well to different domains of data. Moreover, our framework does not strictly require a similar model structure across specialists, and the structure design of the gating module on top of specialists can also be flexibly adjusted to match the desired learning capacity.

D. Case Study

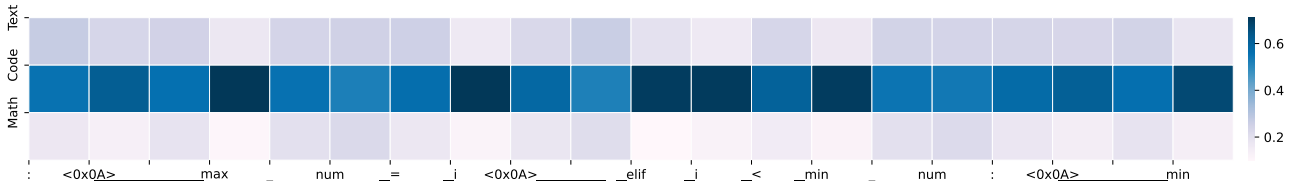
In Section 4.5, we analyze the model expertise at the sequence level and set level. In this section, we provide cases at the token level to illustrate the weight distributions of the three specialist models. Figure 7 and Figure 8 show two cases randomly extracted from ULTRACHAT 2 code data and GSM8K dataset. For coding data, almost all weights are assigned to code specialist model. For math data, there is considerable weight given to code model as well, given the fact that mathematical equation is much alike code snippets. The assumption can be validated by the fact that when it comes to non-mathematical notation, the token weight distribution clearly favors the math specialist more. The observation is in line with our expectation, that the fused model can implicitly learn to allocate tokens to suitable specialist to achieve better performance. Meanwhile, similarity between domains could be captured and their performance can be enhanced jointly by related specialists.



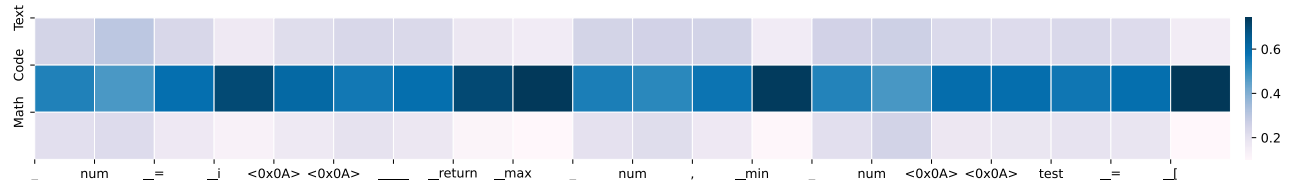
(a) Case study: tokens and weights of code data (a).



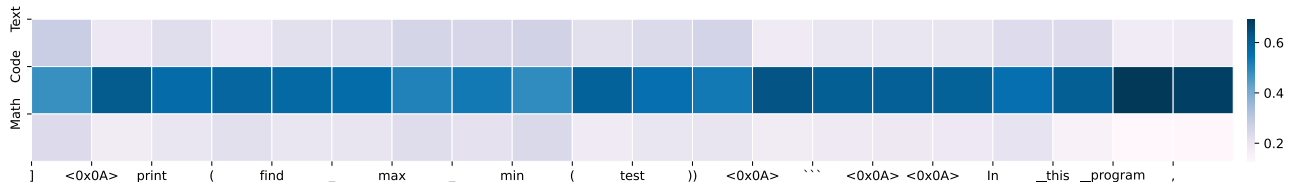
(b) Case study: tokens and weights of code data (b).



(c) Case study: tokens and weights of code data (c).

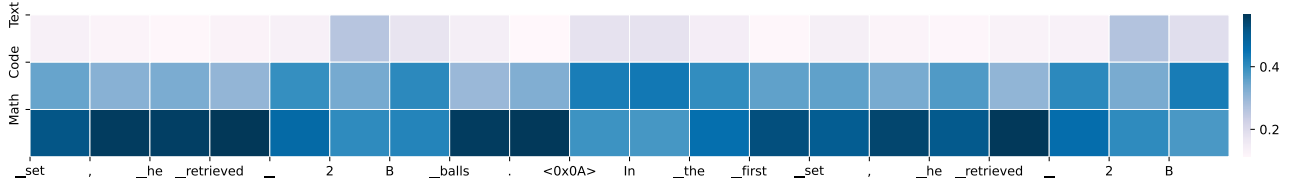


(d) Case study: tokens and weights of code data (d).

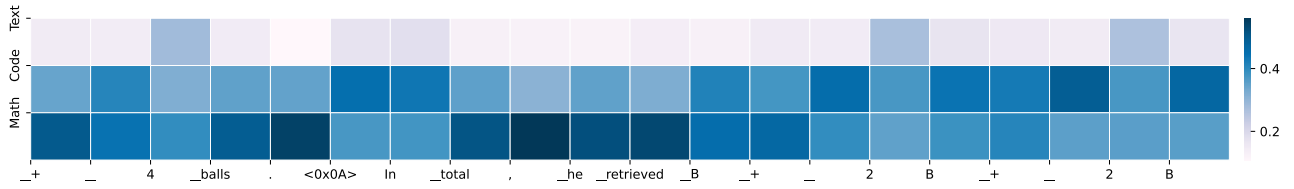


(e) Case study: tokens and weights of code data (e).

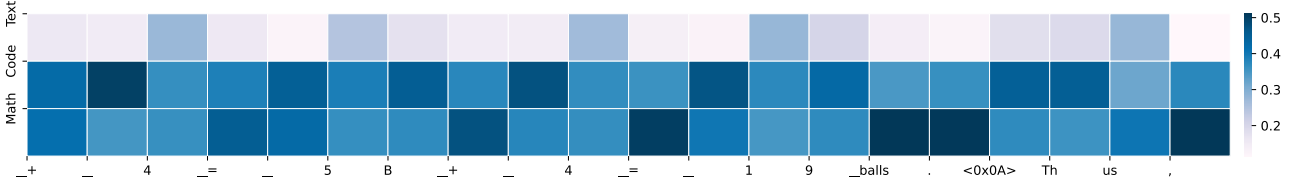
Figure 7: Weight distributions of some pieces of tokens from a sample of code data.



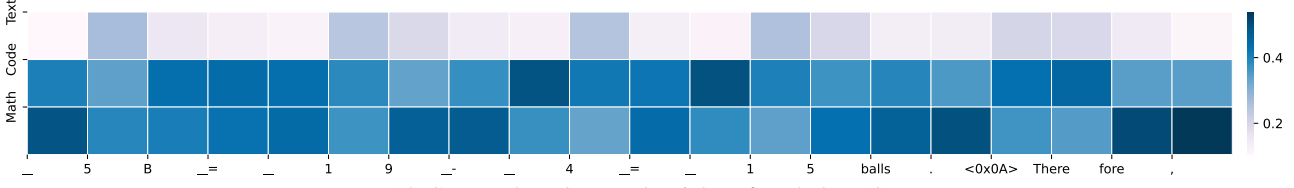
(a) Case study: tokens and weights of math data (a).



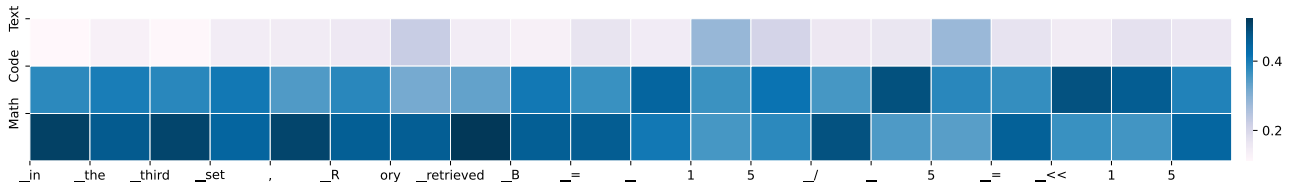
(b) Case study: tokens and weights of math data (b).



(c) Case study: tokens and weights of math data (c).



(d) Case study: tokens and weights of math data (d).



(e) Case study: tokens and weights of math data (e).

Figure 8: Weight distributions of some pieces of tokens from a sample of math data.