

FUSECHAT: Knowledge Fusion of Chat Models

Fanqi Wan, Longguang Zhong, Ziyi Yang, Ruijun Chen, Xiaojun Quan*
School of Computer Science and Engineering, Sun Yat-sen University, China
{wanfq, zhonglg5, yangzy39, chenrj8}@mail2.sysu.edu.cn
quanxj3@mail.sysu.edu.cn

Abstract

While training large language models (LLMs) from scratch can indeed lead to models with distinct capabilities and strengths, it incurs substantial costs and may lead to redundancy in competencies. Knowledge fusion aims to integrate existing LLMs of diverse architectures and capabilities into a more potent LLM through lightweight continual training, thereby reducing the need for costly LLM development. In this work, we propose a new framework for the knowledge fusion of chat LLMs through two main stages, resulting in FUSECHAT. Firstly, we conduct pairwise knowledge fusion on source chat LLMs of varying structures and scales to create multiple target LLMs with identical structure and size via lightweight fine-tuning. During this process, a statistics-based token alignment approach is introduced as the cornerstone for fusing LLMs with different structures. Secondly, we merge these target LLMs within the parameter space, where we propose a novel method for determining the merging coefficients based on the magnitude of parameter updates before and after fine-tuning. We implement and validate FUSECHAT using six prominent chat LLMs with diverse architectures and scales. Experimental results on two instruction-following benchmarks, AlpacaEval 2.0 and MT-Bench, demonstrate the superiority of FUSECHAT-7B over baselines of various sizes.

1 Introduction

Large language models (LLMs) have demonstrated remarkable success across a wide range of natural language processing (NLP) tasks. Currently, it has become prevalent and imperative for individuals and corporations to build their own LLMs. However, the computational resources and time costs associated with LLM development remain

* Corresponding authors.

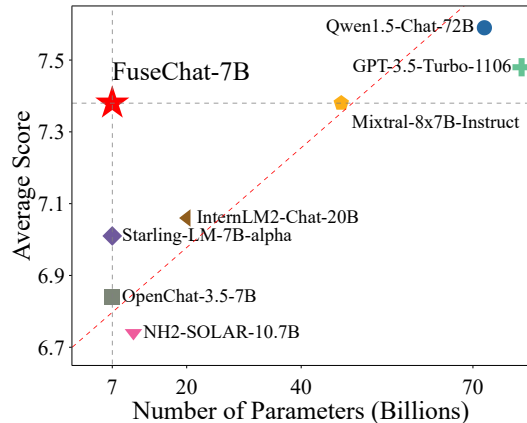


Figure 1: Comparison between FUSECHAT-7B and baselines on MT-Bench. We show that FUSECHAT-7B achieves comparable performance to Mixtral-8x7B and approaches GPT-3.5. The red dashed line is linearly fitted from all baselines except FUSECHAT-7B.

prohibitively high. Furthermore, despite the structural and functional differences among LLMs, they often exhibit similar abilities across various tasks. Therefore, besides training from scratch, another option is to combine the distinct advantages of existing LLMs into a potent LLM, which is termed *knowledge fusion of LLMs* (Wan et al., 2024).

The endeavor to integrate the capabilities of multiple models has been a long-standing pursuit. For example, ensemble methods (Littlestone and Warmuth, 1994; Jiang et al., 2023) directly aggregate the outputs of multiple models to enhance prediction performance and robustness. However, this approach requires maintaining multiple trained models during inference, which is inefficient for LLMs due to their substantial memory and inference time requirements. Another approach is to directly merge several neural networks into a single network through arithmetic operations in the parameter space (Gupta et al., 2020), whereas this approach typically assumes uniform network architectures and requires manually-tuned (Wortsman et al., 2022; Yadav et al., 2024) or automatically-

learned (Matena and Raffel, 2022; Jin et al., 2023) coefficients to merge the parameters of different neural networks. In contrast, knowledge fusion (Wan et al., 2024) seeks to integrate the capabilities of multiple LLMs, irrespective of their architectures, into a single LLM through lightweight continual training. This process embodies a traditional multi-teacher knowledge distillation approach (You et al., 2017), but faces new challenges such as token alignment and fusion strategies across different LLMs.

In this study, we introduce a fuse-and-merge framework to extend the fusion of LLMs to chat-based LLMs¹ with diverse architectures and scales through two stages, resulting in FUSECHAT. Firstly, we conduct pairwise knowledge fusion for source chat LLMs to generate multiple target LLMs of identical structure and size. To achieve this, we first select a pivot LLM and perform token alignment, followed by knowledge fusion between the pivot and each of the remaining LLMs. These target LLMs are expected to inherit the strengths of source chat LLMs through knowledge transfer during lightweight fine-tuning. Secondly, these target LLMs are merged within the parameter space, where we introduce a novel method called SCE (SELECT, CALCULATE & ERASE) to determine the merging coefficients based on the magnitude of parameter updates before and after fine-tuning. Moreover, SCE allocates parameter matrix-level coefficients that enable the merging at a fine-grained granularity without additional training efforts.

FUSECHAT offers superior potential compared to FUSELLM (Wan et al., 2024). Firstly, while FUSELLM limits its exploration to source LLMs with the same size, FUSECHAT broadens the scope by incorporating six source LLMs with varying scales. This allows for greater adaptability to the fusion of heterogeneous LLMs. Secondly, the framework of FUSELLM does not seamlessly support the inclusion of new source LLMs as it requires the combination of distribution matrices from all source LLMs during continual training. In contrast, integrating a new source LLM in FUSECHAT is plug-and-play, requiring only obtaining a target LLM from the new source LLM and merging it with the existing FUSECHAT. Thirdly, compared to many-to-one knowledge fusion, pairwise fusion empirically mitigates the challenges of knowledge distillation from heterogeneous source LLMs.

¹We refer to “chat-based LLMs” simply as “chat LLMs”.

To verify the effectiveness of FUSECHAT, we implemented FUSECHAT-7B using six prominent open-source chat LLMs: OpenChat-3.5-7B (Wang et al., 2024a), Starling-LM-7B-alpha (Zhu et al., 2024), NH2-SOLAR-10.7B (Kim et al., 2023), InternLM2-Chat-20B (Cai et al., 2024), Mixtral-8x7B-Instruct (Jiang et al., 2024), and Qwen-1.5-Chat-72B (Bai et al., 2023). Experimental results on two representative instruction-following benchmarks, AlpacaEval 2.0 (Dubois et al., 2024b) and MT-Bench (Zheng et al., 2024), demonstrate the superiority of FUSECHAT-7B across a broad spectrum of chat LLMs at 7B, 10B, and 20B scales. Moreover, we validated the proposed token alignment method and the SCE merging method through a series of analytical experiments.

2 FUSECHAT

2.1 Overview

Figure 2 presents an overview of our FUSECHAT in comparison with FUSELLM (Wan et al., 2024). The FUSECHAT framework consists of two main stages: fuse and merge. In the *fuse* stage, pairwise knowledge fusion is conducted on source chat LLMs to derive multiple target LLMs with identical structure and size. This process begins by selecting a pivot LLM, followed by performing knowledge fusion between the pivot and each remaining LLM. In the *merge* stage, these target LLMs are combined within the parameter space, where we determine the merging coefficients based on the magnitude of parameter updates before and after fine-tuning.

Specifically, considering K source LLMs $\{\mathcal{M}_i^s\}_{i=1}^K$ with varying architectures and scales, FUSECHAT first specifies one of the source LLMs, \mathcal{M}_i^s , as the pivot and then applies pairwise knowledge fusion to obtain $(K - 1)$ target LLMs, $\{\mathcal{M}_j^t\}_{j=1}^{K-1}$, which share the same architecture as the pivot LLM. The selection of the pivot depends on the desired structure and scale for the target LLMs, while also considering the capabilities and performance of a candidate LLM.

To perform pairwise knowledge fusion, FUSECHAT prompts these source LLMs using a supervised fine-tuning dataset $\mathcal{D} = \{I_i, R_i\}_{i=1}^M$ to showcase their inherent knowledge by responding to each instruction in \mathcal{D} . Token alignment (Fu et al., 2023; Wan et al., 2024) between the source LLMs and the pivot is then conducted to properly map the resulting probabilistic distribution matrices. These distribution matrices are subsequently used

for pairwise knowledge fusion (Wan et al., 2024) through lightweight fine-tuning to obtain $(K - 1)$ target LLMs. Following this, the target LLMs are merged in the parameter space to yield the final fused LLM \mathcal{M}^f . To incorporate fine-grained advantages of target LLMs, we introduce a new merging method named SCE to obtain the merging coefficients based on *selection*, *calculation*, and *erasure* on the task vectors (Ilharco et al., 2023) which represent variation of model weights before and after fine-tuning. SCE enables the automatic allocation of parameter matrix-level merging coefficients, facilitating the merging of LLMs at a finer granularity.

2.2 Preliminaries

Given an instruction I_i and the corresponding response R_i of length N from the fine-tuning dataset \mathcal{D} , we use $R_{i,<t} = (r_{i,1}, r_{i,2}, \dots, r_{i,t-1})$ to represent the sequence preceding the t th token in the response. The supervised fine-tuning (SFT) objective for an LLM parameterized by θ is defined as minimizing the following negative log-likelihood:

$$\mathcal{L}_{\text{SFT}} = -\mathbb{E}_{(I_i, R_i) \sim \mathcal{D}} \left[\sum_{t \leq N} \log p_{\theta}(r_{i,t} | R_{i,<t}, I_i) \right], \quad (1)$$

where $p_{\theta}(r_{i,t} | R_{i,<t}, I_i)$ is the model’s predicted probability for the t th token $r_{i,t}$ in R_i given the instruction and preceding tokens in the response.

2.3 Pairwise Knowledge Fusion

To facilitate the description of pairwise knowledge fusion, we reframe the above token-level view into a matrix format. Specifically, for each instruction I_i , we transform the token-level predictions into a probabilistic distribution matrix, $\mathbf{P}_i^{\theta} \in \mathbb{R}^{N \times V}$, where V denotes the vocabulary size. The distribution matrix is assumed to reflect certain inherent knowledge of the language model in responding to the instruction (Wan et al., 2024). Consequently, different probabilistic distribution matrices obtained from different chat LLMs can be used to represent the diverse knowledge embedded within these models. Based on this assumption, FUSECHAT performs pairwise knowledge fusion by fine-tuning the target LLMs, initialized from the pivot, using the probabilistic distribution matrices.

Model Fusion For each instruction I_i in \mathcal{D} , we first feed it into the K source chat LLMs to obtain a set of probabilistic distribution matrices, denoted

as $\{\mathbf{P}_i^{\theta_j}\}_{j=1}^K$, where θ_j represents the parameters of the j th chat LLM. Since these LLMs may employ different tokenizers, token alignment is necessary to properly map their probabilistic distribution matrices (Fu et al., 2023; Wan et al., 2024). Then, pairwise knowledge fusion is conducted between the pivot LLM and each of the remaining source LLMs. To achieve this, we denote the probabilistic distribution matrix generated by the pivot LLM as $\mathbf{P}_i^{\theta_v}$ and merge it with each $\mathbf{P}_i^{\theta_j}|_{j \neq v}$ to obtain a set $\{\mathbf{P}_i^j\}_{j=1}^{K-1}$ of fused matrices as follows:

$$\mathbf{P}_i^j = \mathbb{F}\text{usion}(\mathbf{P}_i^{\theta_v}, \mathbf{P}_i^{\theta_j})|_{j \neq v}, \quad (2)$$

where $\mathbb{F}\text{usion}(\cdot)$ represents the fusion function that merges two matrices. The resulting matrix \mathbf{P}_i^j is seen as a representation of the collective knowledge and distinctive strengths of the two source LLMs. Among various fusion strategies, this work employs minimum cross-entropy (MinCE) following Wan et al. (2024) as the fusion function, which empirically performs the best. Specifically, the MinCE fusion function selects the distribution matrix that exhibits the minimum cross-entropy score with respect to the response R_i for instruction I_i .

We then enforce alignment between the prediction of each target LLM \mathcal{M}_j^t and the corresponding fused representation matrices \mathbf{P}_i^j . We use $\mathbf{Q}_i^{\phi_j}$ to represent the output distribution matrix of target LLM \mathcal{M}_j^t for instruction I_i and define the fusion objective for training each target LLM as follows:

$$\mathcal{L}_{\text{Fusion}} = -\mathbb{E}_{(I_i, R_i) \sim \mathcal{D}} \left[\mathbb{H}(\mathbf{P}_i^j || \mathbf{Q}_i^{\phi_j}) \right], \quad (3)$$

where $\mathbb{H}(\cdot || \cdot)$ represents the cross entropy between two probabilistic distribution matrices.

The overall training objective for each pairwise knowledge fusion consists of a weighted combination of the supervised fine-tuning objective \mathcal{L}_{SFT} and the fusion objective $\mathcal{L}_{\text{Fusion}}$:

$$\mathcal{L} = \lambda \mathcal{L}_{\text{SFT}} + (1 - \lambda) \mathcal{L}_{\text{Fusion}}. \quad (4)$$

Token Alignment Token alignment aims to address the mappings of probabilistic distribution matrices $\{\mathbf{P}_i^{\theta_j} \in \mathbb{R}^{N \times V}\}_{j=1}^K$ generated by different source LLMs for a given instruction I_i . Therefore, the alignment involves two dimensions of the matrices: sequence dimension for the tokenized response and distribution dimension for the probabilistic distributions. In the sequence dimension, we follow previous works (Fu et al., 2023; Wan et al., 2024)

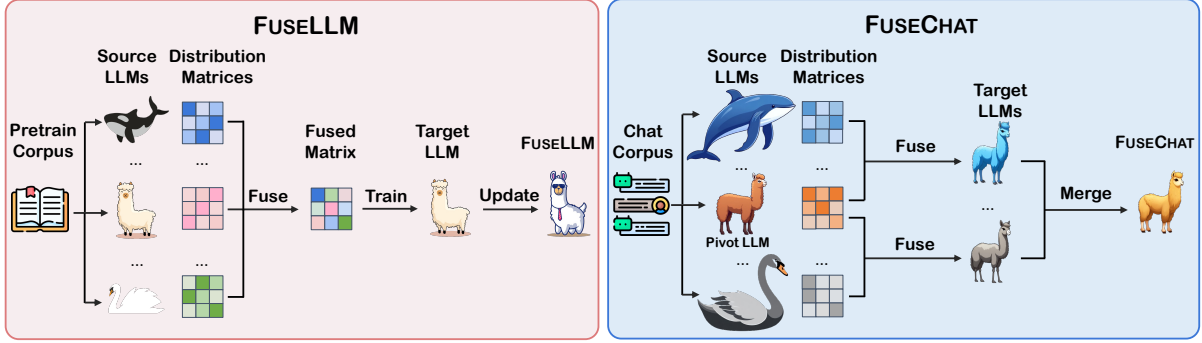


Figure 2: Overview of FUSECHAT in comparison with FUSELLM (Wan et al., 2024). Distinct animal icons symbolize different LLMs, where each species and size indicates a unique architecture and scale, respectively.

and adopt a modified dynamic time wrapping algorithm (Senin, 2008) to recursively minimize the total cost of editing the tokens from a source LLM to align them with the pivot LLM. This process may result in 1-1, 1-n, and n-1 mappings, as shown in Figure 5. In the distribution dimension, Fu et al. (2023) focused on aligning distributions based on the exact match (EM) between tokens in source and target distributions, which restricts the alignment to only 1-1 mappings and may result in too many unmatched tokens. Wan et al. (2024) relaxed the EM constraint by aligning the distributions based on the minimum edit distance (MinED) between tokens in the vocabularies of source and target LLMs. While this approach improves the mapping success rate and expands to 1-n mappings, it ignores n-1 mappings and may introduce many misalignments.

In this work, we propose an enhanced token alignment strategy that utilizes mapping statistics (MS) from the sequence dimension as the criteria for alignment in the distribution dimension. We construct a global statistical matrix, where each column represents the frequency of mappings from a pivot token to all potential source tokens, derived from sequence-dimensional token alignments. In the case of 1-1 and 1-n mappings, we align the distributions based on the maximum mapping frequency in the respective columns of the statistical matrix for each pivot token in the distribution. For n-1 mappings, we first calculate a weighted average of the source tokens’ distributions according to their mapping frequencies in the statistical matrix to obtain a merged distribution. This merged distribution is then aligned to the pivot distribution, similar to the procedure employed for 1-1 mappings. As shown in Figure 5, this strategy better reflects the token mapping statistics in the dataset, thereby preserving information in the aligned distribution matrices while minimizing alignment errors.

2.4 Model Merging

Since the target LLMs $\{\mathcal{M}_j^t\}_{j=1}^{K-1}$ resulting from pairwise knowledge fusion share identical architecture and scale while possessing diverse advantages and capabilities learned from the source LLMs, we further merge them in the parameter space to obtain the final fused LLM \mathcal{M}^f . To ensure the adaptability and scalability of FUSECHAT, it is essential to maintain the simplicity of the merging strategy. Primarily, the calculation of merging coefficients should be automated, obviating the complex hyperparameter tuning. Secondly, the merging procedure should not demand forward or backward propagation over additional data, which is computationally inefficient and memory-intensive.

As described in Algorithm 1, we propose a novel merging method named SCE (*select, calculate, and erase*) for parameter matrix-level merging. Analogous to task vectors (Ilharco et al., 2023), we first define fusion vectors $\{\delta_j\}_{j=1}^{K-1}$ (Eq. 5) as the direction and magnitude of weight updates from pivot LLM \mathcal{M}_v^s to target LLMs $\{\mathcal{M}_j^t\}_{j=1}^{K-1}$ during model fusion. For each parameter matrix unit in target LLMs, we derive the merged weights using fusion vectors through a three-step process.

(1) Select: During the pairwise knowledge fusion, target LLMs dynamically evolve to incorporate the advantages of their corresponding source LLMs. Fusion vectors for each parameter matrix unit with substantial variations across different target LLMs are supposed to signify distinctive and significant strengths. Therefore, we first select the top $\tau\%$ elements from each parameter matrix-level fusion vector $\{\delta_{j,m}\}_{j=1}^{K-1}$ with high variance across multiple target LLMs, resulting in $\{\hat{\delta}_{j,m}\}_{j=1}^{K-1}$ (Eq. 6). **(2) Calculate:** We then calculate the sum of squares of elements in $\hat{\delta}_{j,m}$ and obtain a matrix-level merging coefficient for each target LLM as

Algorithm 1 SCE Procedure

Input: Target LLMs parameters $\{\phi_j\}_{j=1}^{K-1}$, pivot LLM parameters θ_v , threshold τ .

Output: Merged LLM parameters Φ

▷ Create fusion vectors

$$\{\delta_j\}_{j=1}^{K-1} = \{\phi_j - \theta_v\}_{j=1}^{K-1} \quad (5)$$

▷ Calculate matrix-level coefficients

for $\{\delta_{j,m}\}_{j=1}^{K-1} \in \{\delta_j\}_{j=1}^{K-1}$ **do**

▷ Step 1: Select salient elements

$$\{\hat{\delta}_{j,m}\}_{j=1}^{K-1} = \text{Select}(\{\delta_{j,m}\}_{j=1}^{K-1}, \tau) \quad (6)$$

▷ Step 2: Calculate coefficients

$$\{\eta_{j,m}\}_{j=1}^{K-1} = \text{Calculate}(\{\hat{\delta}_{j,m}\}_{j=1}^{K-1}) \quad (7)$$

▷ Step 3: Erase minority elements

$$\{\delta'_{j,m}\}_{j=1}^{K-1} = \text{Erase}(\{\hat{\delta}_{j,m}\}_{j=1}^{K-1}) \quad (8)$$

▷ Update merged LLM parameters

$$\Phi_m = \theta_{v,m} + \sum_{j=1}^{K-1} \eta_{j,m} \delta'_{j,m} \quad (9)$$

end

return Φ

$\eta_{j,m} = \frac{\sum_j \hat{\delta}_{j,m}^2}{\sum_j \sum_m \hat{\delta}_{j,m}^2}$. **(3) Erase:** Each parameter may exhibit conflicting signs across fusion vectors from different target LLMs, which could cause interference during model merging (Yadav et al., 2024). Thus, for each parameter we sum its values in $\{\hat{\delta}_{j,m}\}_{j=1}^{K-1}$ across target LLMs and erase elements with minority directions (Eq. 8). Finally, the filtered $\{\delta'_{j,m}\}_{j=1}^{K-1}$ are merged based on the calculated coefficients, and added to the pivot LLM’s parameters (Eq. 9).

2.5 Discussions

Distinction from FUSELLM While FUSELLM (Wan et al., 2024) emphasizes the fusion of multiple base LLMs through continual pre-training, FUSECHAT focuses on integrating diverse chat LLMs into a unified chat model via SFT. This difference in both training objectives and the type of LLMs makes FUSECHAT essential in chat LLMs fusion. Unlike FUSELLM, which directly applies *multi-teacher distillation*, FUSECHAT employs a novel two-stage *fuse-then-merge* approach. This method is not only *highly scalable and efficient* (Appendix 9), but also *resolves knowledge conflicts* in parameter space while precisely integrating the distinct strengths of each source LLM (Section 3.3). By fusing six heterogeneous chat LLMs, we val-

idate FUSECHAT’s superiority over FUSELLM (OpenChat-3.5-7B-Multi) as shown in Table 1.

Distinction from the TIES Merging TIES merging (Yadav et al., 2024) relies on *manually tuned, model-level* coefficients for combining different models. In contrast, our SCE merging automates the merging process by leveraging weight updates from a pivot LLM to *automatically* compute *matrix-level* coefficients. This enables the *fine-grained* incorporation of diverse benefits across LLMs, which is difficult to achieve with *manual hyperparameter tuning*. In our specific context, where target LLMs are trained on identical datasets with relatively *subtle parameter variations*, SCE excels at capturing and preserving the distinctive advantages of each LLM through *nuanced matrix-level* parameter updates. Experimental results in Table 2 demonstrate that SCE outperforms merging baselines, including TIES, thereby validating its efficacy and impact.

3 Experiments

In our experiments, we explore the fusion of chat LLMs with diverse architectures and scales. We use six representative chat LLMs as the source LLMs, including OpenChat-3.5-7B (Wang et al., 2024a), Starling-LM-7B-alpha (Zhu et al., 2024), NH2-SOLAR-10.7B (Kim et al., 2023), InternLM2-Chat-20B (Cai et al., 2024), Mixtral-8x7B-Instruct (Jiang et al., 2024), and Qwen-1.5-Chat-72B (Bai et al., 2023). As for the pivot LLM, which serves as the starting point for the target LLMs, we opt for OpenChat-3.5-7B due to its balanced scale and performance. We first apply pairwise knowledge fusion (Section 2.3) to create five distinct target LLMs with the same structure. These target LLMs are then merged using the SCE method (Section 2.4), resulting in FUSECHAT-7B.

3.1 Experimental Setup

Training Dataset To leverage the strengths of source LLMs during knowledge fusion while alleviating catastrophic forgetting, we curate a high-quality dataset named FUSECHAT-MIXTURE from two different sources. First, 50% of the training instances are sampled from the dataset used by the pivot LLM, OpenChat-3.5-7B. Second, we gather the remaining instances, which have not been encountered by the pivot LLM, from open-source communities. These two sources result in a corpus comprising approximately 95,000 dialogues across various domains. Please refer to Appendix B.3 for

Model	#Params	AlpacaEval 2.0 (GPT-4-1106-Preview)		MT-Bench (GPT-4-0125-Preview)		
		Win Rate	LC Win Rate	1st Turn	2nd Turn	Average Score
Proprietary LLMs						
GPT-3.5-Turbo-1106 (Achiam et al., 2023)	-	9.18	19.30	7.56	7.41	7.48
Claude-3-Opus (Anthropic, 2024)	-	29.04	40.39	8.84	8.30	8.57
GPT-4-1106-Preview (Achiam et al., 2023)	-	50.00	50.00	8.86	8.71	8.79
Source LLMs						
OpenChat-3.5-7B (Wang et al., 2024a)	7B	10.20	14.90	7.14	6.55	6.84
Starling-LM-7B-alpha (Zhu et al., 2024)	7B	14.20	14.70	7.54	6.49	7.01
NH2-SOLAR-10.7B (Kim et al., 2023)	10.7B	12.22	18.13	7.11	6.36	6.74
InternLM2-Chat-20B (Cai et al., 2024)	20B	21.70	18.70	7.78	6.34	7.06
Mixtral-8x7B-Instruct (Jiang et al., 2024)	8x7B	18.30	23.70	7.76	7.00	7.38
Qwen1.5-Chat-72B (Bai et al., 2023)	72B	26.50	36.60	7.83	7.36	7.59
Ensemble LLMs						
Top1-PPL (Mavromatis et al., 2024)	162B	25.11	27.97	7.79	6.95	7.37
Top1-LLM-Blender (Jiang et al., 2023)	162B	24.45	29.11	7.85	6.70	7.28
Top1-GPT4 (Achiam et al., 2023)	162B	42.82	43.87	8.79	8.01	8.40
Fused LLMs						
OpenChat-3.5-7B SFT	7B	10.56	14.50	7.36	6.40	6.88
OpenChat-3.5-7B Multi	7B	10.19 (-3.5%)	13.43 (-7.4%)	<u>7.69</u> (+4.5%)	6.26 (-2.2%)	6.99 (+1.6%)
OpenChat-3.5-7B Starling	7B	11.43 (+8.2%)	16.20 (+11.7%)	<u>7.69</u> (+4.5%)	6.73 (+5.2%)	7.22 (+4.9%)
OpenChat-3.5-7B SOLAR	7B	11.12 (+5.3%)	16.51 (+13.9%)	7.58 (3.0%)	6.76 (+5.6%)	7.17 (+4.2%)
OpenChat-3.5-7B InternLM	7B	11.82 (+11.9%)	15.21 (+4.9%)	7.63 (+3.7%)	6.78 (+5.9%)	7.21 (+4.8%)
OpenChat-3.5-7B Mixtral	7B	<u>11.74</u> (+11.2%)	<u>16.52</u> (+13.9%)	7.58 (+3.0%)	<u>6.90</u> (+7.8%)	<u>7.24</u> (+5.2%)
OpenChat-3.5-7B Qwen	7B	10.93 (+3.5%)	14.98 (+3.3%)	<u>7.69</u> (+4.5%)	6.78 (+5.9%)	7.23 (+5.1%)
FUSECHAT-7B	7B	11.52 (+9.1%)	17.16 (+18.3%)	7.70 (+4.6%)	7.05 (+10.2%)	7.38 (+7.3%)

Table 1: Results of FUSECHAT-7B and baselines on AlpacaEval 2.0 and MT-Bench. The bold font denotes the best performance among fused LLMs, while the underscore indicates the second-best performance. Moreover, the percentages represent the relative performance improvement compared to the OpenChat-3.5-7B SFT.

more details of FUSECHAT-MIXTURE.

Training Details In all experiments, we train the target LLMs using a batch size of 128 and a maximum length of 2048 on a single node with 8x80GB NVIDIA A800 GPUs for three epochs, which takes approximately 9 hours. The models are optimized using the AdamW (Loshchilov and Hutter, 2019) optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. We use a weight decay of 0.0 and gradient clipping of 1.0. A cosine learning rate schedule is employed, with a maximum learning rate of $5e-6$ and a warmup ratio of 0.03. We empirically set the combination weight λ in Eq. 4 to 0.9, and the rationale behind the value of λ is detailed in Appendix D. Our training framework is implemented based on the HuggingFace Transformers (Wolf et al., 2020).

Evaluation We assess the performance of FUSECHAT-7B on two representative benchmarks to evaluate its ability to follow instructions and engage in conversations effectively. The first benchmark, AlpacaEval 2.0 (Dubois et al., 2024b), comprises 805 instructions across five test subsets. It compares the Win Rate and Length-Controlled Win Rate (LC Win Rate) (Dubois et al., 2024a) of a

model against GPT-4. We employ the default settings and utilize GPT-4 (GPT-4-1106-Preview) to evaluate the quality of generated responses. The second benchmark, MT-Bench (Zheng et al., 2024), consists of 80 multi-turn dialogues spanning various domains, including writing, roleplay, reasoning, math, coding, extraction, STEM, and humanities. Originally, GPT-4 (GPT-4-0613) was used as the evaluator, providing a scalar score ranging from 1 to 10 for each generated response. However, due to inaccuracies in the reference responses, we adopt an updated version, GPT-4-0125-Preview, as per the latest work (Wang et al., 2024c), to correct the errors and evaluate the generated responses.

Baselines In our experiments, we compare our FUSECHAT-7B with four categories of baselines, including (i) Proprietary LLMs, (ii) Source LLMs, (iii) Ensemble LLMs, and (iv) Fused LLMs. The details of baselines are shown in Appendix B.4.

3.2 Overall Results

In Table 1, we present the overall results of FUSECHAT-7B in comparison with baselines of various architectures and scales on AlpacaEval 2.0

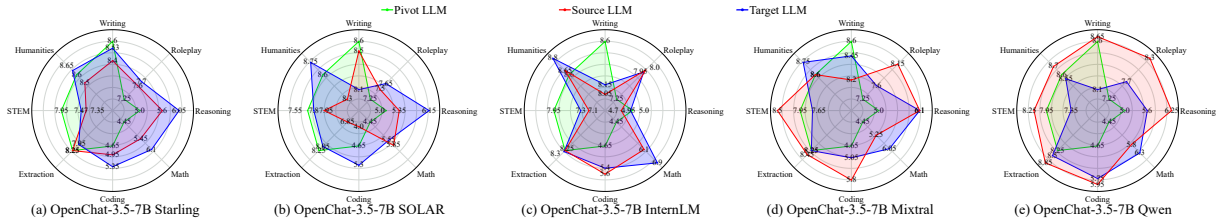


Figure 3: The effect of pairwise knowledge fusion for source LLMs across various domains on MT-Bench. It combines the strengths of each source LLM and the pivot (OpenChat-3.5-7B) into a more potent target LLM.

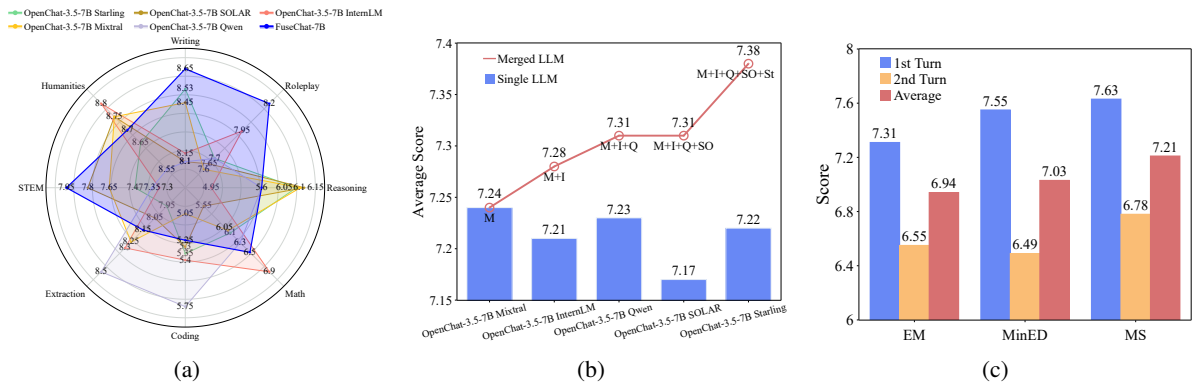


Figure 4: Comparison of different aspects of LLM fusion. (Left) The effect of merging target LLMs into FUSECHAT-7B to combine their strengths across domains on MT-Bench. (Middle) Results of FUSECHAT by merging varying numbers of target LLMs. (Right) Results of OpenChat-3.5-7B InternLM with different token alignment strategies.

and MT-Bench. Key observations are as follows. Firstly, OpenChat-3.5-7B SFT, fine-tuned on our high-quality dataset, slightly outperforms the pivot LLM OpenChat-3.5-7B. Secondly, in comparison to OpenChat-3.5-7B Multi, which fuses multiple source LLMs simultaneously as FUSELLM (Wan et al., 2024), the target LLMs resulting from pairwise knowledge fusion exhibit superior performance, demonstrating the effectiveness of pairwise fusion in reducing the fusion difficulty. For instance, through the integration of OpenChat-3.5-7B and Mixtral-8x7B-Instruct, the fused target LLM OpenChat-3.5-7B Mixtral achieves relative gains of 13.9% LC Win Rate and 5.2% Average Score over OpenChat-3.5-7B SFT, significantly surpassing OpenChat-3.5-7B Multi. Furthermore, after merging these target LLMs, FUSECHAT-7B shows substantial performance enhancements of 18.3% and 7.3% in the two metrics. This illustrates the superiority of FUSECHAT-7B across source LLMs of various scales, even comparable to 8x7B MoEs and approaching GPT-3.5.

Moreover, in comparison to the ensemble LLMs of 162B, which generate the 1st response from six parallel deployed LLMs based on different ranking criteria, FUSECHAT-7B outperforms most of them except Top1-GPT4 on MT-Bench, while being 23x

smaller and independent of GPT-4.

To further illustrate that our performance improvements stem from the integration of distinct knowledge from multiple LLMs, we evaluate the source LLMs, target LLMs, and FUSECHAT across various domains on MT-Bench. The results in Figure 3 reveal that the target LLMs demonstrate noticeable performance enhancements in most domains after pairwise knowledge fusion. Typically, the performance of each target LLM falls between that of the pivot LLM and the respective source LLM. This phenomenon can be attributed to the fusion function we employed to select the optimal target distributions with minimal cross-entropy, which promotes the incorporation of unique advantages from the pivot LLM and source LLMs into more potent target LLMs. Notably, in math and coding domains, the performance of certain target LLMs surpasses that of either the pivot or source LLMs. This enhancement can be explained by the strong performance of the source LLMs in these domains, coupled with the relatively high proportion of math and coding samples in our dataset. It is also consistent with findings from knowledge distillation (Wu et al., 2023), where the student model occasionally outperforms the teacher in specific tasks. The effect of further merging these target

Model	AlpacaEval 2.0	MT-Bench
FUSECHAT-7B Linear	17.12	7.03
FUSECHAT-7B TA	15.74	7.08
FUSECHAT-7B TIES	16.55	7.33
FUSECHAT-7B DARE	16.57	7.15
FUSECHAT-7B SCE	17.16	7.38

Table 2: Comparison of different merging methods on AlpacaEval 2.0 and MT-Bench.

LLMs into FUSECHAT-7B is shown in Figure 4(a). By integrating the capabilities of the target LLMs, FUSECHAT achieves a balanced and robust performance across diverse domains.

3.3 Analysis of Model Merging

To investigate the effectiveness of the proposed SCE approach, we incorporate the target LLMs using different merging methods, including Linear (Wortsman et al., 2022), TA (Ilharco et al., 2023), TIES (Yadav et al., 2024), and DARE (Yu et al., 2024a). We evaluate the performance of these merged LLMs on AlpacaEval 2.0 and MT-Bench. As depicted in Table 2, FUSECHAT-7B SCE outperforms all baseline methods on the two benchmarks. We attribute this superior performance to SCE’s ability to actively reduce knowledge conflicts by selectively integrating beneficial updates while removing conflicting ones. For more details of model merging, please refer to Appendix B.2.

In Figure 4(b), we further illustrate the performance of FUSECHAT-7B SCE by incorporating varying numbers of target LLMs on MT-Bench. The findings demonstrate a progressive enhancement in Average Score, which increases from 7.24 to 7.38 as the number of integrated target LLMs rises from 1 to 5. Moreover, we observe that after the integration of OpenChat-3.5-7B SOLAR, the performance of the merged LLM remains stable. This stabilization might be attributed to the comparatively sub-optimal performance of OpenChat-3.5-7B SOLAR and its corresponding NH2-SOLAR-10.7B compared to other target or source LLMs. Therefore, we suggest that both the diversity and quality of integrated source LLMs are critical factors for optimal knowledge fusion.

3.4 Ablation Studies for SCE

In this section, we conduct experiments to examine the effectiveness of the *select*, *calculate*, and *erase* operations in SCE. The results in Table 3 demonstrate that, without the *select* step, FUSECHAT-7B

Model	AlpacaEval 2.0	MT-Bench
FUSECHAT-7B SCE	17.16	7.38
FUSECHAT-7B CE	15.69 (-8.57%)	7.29 (-1.22%)
FUSECHAT-7B C	16.62 (-3.15%)	7.11 (-3.66%)

Table 3: Comparison of different merging methods on AlpacaEval 2.0 and MT-Bench. “CE” and “C” mean only *calculate&erase* and *calculate* operations are used.

Model	AlpacaEval 2.0	MT-Bench
Starling-LM-7B-alpha	14.70	7.01
Starling-LM-7B-alpha SFT	13.20 (-10.20%)	6.89 (-1.71%)
FUSECHAT-Starling-7B	17.29 (+17.62%)	7.16 (+2.14%)

Table 4: Starling-LM-7B-alpha as pivot LLM results on AlpacaEval 2.0 and MT-Bench.

CE suffers substantial performance degradation. This underscores the benefits of selecting salient elements from fusion vectors with high variance among target LLMs to signify their distinctive and significant strengths. Moreover, removing both the *select* and *erase* operations leads to FUSECHAT-7B C with decreased performance, highlighting the importance of resolving parameter interference in fusion vectors from different target LLMs.

3.5 Analysis of Token Alignment

Finally, we delve into exploring the impact of various token alignment strategies. Specifically, we apply EM (Fu et al., 2023) and MinED (Wan et al., 2024), and our MS methods to align distributions generated by InterLM2-Chat-20B with those of OpenChat-3.5-7B. Then, we conduct pairwise knowledge fusion to derive OpenChat-3.5-7B InternLM. As depicted in Figure 4(c), our proposed MS method, rooted in mapping statistics, consistently outperforms EM and MinED, which rely on exact matching and minimal edit distance, respectively. We propose that this performance enhancement arises from MS’s effective utilization of token mapping statistics within the data, which greatly improves the effect of token alignment in the distribution dimension.

3.6 Different Pivot LLM

We conduct experiments using Starling-LM-7B-alpha to replace OpenChat-3.5-7B as a more robust pivot LLM, which achieved an LC Win Rate of 14.70 on AlpacaEval 2.0 and an Average Score of 7.01 on MT-Bench. The evaluation results presented in Table 4 show that FUSECHAT-Starling-7B outperforms Starling-LM-7B-alpha, with rela-

Dataset Scale	MT-Bench		
	1st Turn	2nd Turn	Average Score
10,000	7.34	6.86	7.10
25,000	7.58	6.85	7.21
95,000	7.70	7.05	7.38

Table 5: Comparison results of different dataset scales on MT-Bench.

tive performance improvements of 17.62% on AlpacaEval 2.0 and 2.14% on MT-Bench. Notably, although Starling-LM-7B-alpha SFT does not result in performance gains, the pairwise knowledge fusion and model merging processes lead to significant enhancements using the same training data.

3.7 Dataset Scaling

We perform experiments across different dataset scales for pairwise knowledge fusion, followed by merging the resulting target LLMs in the parameter space to obtain the final fused LLM. The results in Table 5 indicate that the performance of the final fused LLM consistently improves as the training data scales up from 10k to 95k on MT-Bench, demonstrating the potential effectiveness of scaling up the dataset for our method.

3.8 Additional Analysis

We further conduct experiments to explore the generalizability and scalability of FUSECHAT, including an evaluation on *additional benchmarks* in Appendix B.4 and a *cost analysis* in Appendix E. Moreover, to explore the effectiveness of *pairwise fusion*, we compare it with *single-model distillation* in Appendix F. Finally, we provide *statistical evidence* in Appendix G to support the significance of FUSECHAT-7B’s performance improvements.

4 Conclusion

In this work, we propose a fuse-and-merge framework for knowledge fusion of structurally and scale-varied chat LLMs to integrate their collective knowledge and individual strengths into a more potent chat LLM, resulting in FUSECHAT. FUSECHAT first undertakes pairwise knowledge fusion for source chat LLMs to derive multiple target LLMs of identical structure and size via lightweight fine-tuning. Then, these target LLMs are merged within the parameter space using a novel method SCE to calculate the merging coefficients based on the magnitude of parameter updates before and

after fine-tuning. Experimental results on two representative benchmarks demonstrate the superiority of FUSECHAT-7B over various baselines, even comparable to Mixtral-8x7B-Instruct and approaching GPT-3.5-Turbo-1106 on MT-Bench.

Limitations

Our work relies on constructing a knowledge fusion dataset that spans diverse domains and leverages the strengths of source LLMs. This process demands substantial data engineering efforts, limiting the scalability of our methodology. Future research should focus on developing more efficient data synthesis techniques to expand the scope of the knowledge fusion dataset. Additionally, while our study shows improvements in chat model capabilities, it does not address other critical aspects of LLMs, such as knowledge comprehension and the mitigation of hallucinations. Further investigation is necessary to evaluate the applicability and effectiveness of our approach in these areas.

Acknowledgements

This work was supported by the National Natural Science Foundation of China (No. 62176270) and the Guangdong Basic and Applied Basic Research Foundation (No. 2023A1515012832).

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos Garea, Matthieu Geist, and Olivier Bachem. 2024. On-policy distillation of language models: Learning from self-generated mistakes. In *The Twelfth International Conference on Learning Representations*.
- Takuya Akiba, Makoto Shing, Yujin Tang, Qi Sun, and David Ha. 2024. Evolutionary optimization of model merging recipes. *arXiv preprint arXiv:2403.13187*.
- AI Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, and 1 others. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, and 1 others. 2020. Piqa: Reasoning about physical

- commonsense in natural language. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 7432–7439.
- Zheng Cai, Maosong Cao, Haojong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, and 1 others. 2024. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2924–2936.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Ning Ding, Yulin Chen, Ganqu Cui, Xingtai Lv, Ruobing Xie, Bowen Zhou, Zhiyuan Liu, and Maosong Sun. 2024. Mastering text, code and math simultaneously via fusing highly specialized language models. *arXiv preprint arXiv:2403.08281*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. 2024a. Length-controlled alpaca-eval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*.
- Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy S Liang, and Tatsunori B Hashimoto. 2024b. AlpacaFarm: A simulation framework for methods that learn from human feedback. *Advances in Neural Information Processing Systems*, 36.
- William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39.
- Wenfeng Feng, Chuzhan Hao, Yuewei Zhang, Yu Han, and Hao Wang. 2024. Mixture-of-loras: An efficient multitask tuning method for large language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 11371–11380.
- Ronald A Fisher. 1922. On the mathematical foundations of theoretical statistics. *Philosophical transactions of the Royal Society of London. Series A, containing papers of a mathematical or physical character*, 222(594-604):309–368.
- Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and Tushar Khot. 2023. Specializing smaller language models towards multi-step reasoning. In *International Conference on Machine Learning*, pages 10421–10430. PMLR.
- Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 2024. MiniLLM: Knowledge distillation of large language models. In *The Twelfth International Conference on Learning Representations*.
- Vipul Gupta, Santiago Akle Serrano, and Dennis DeCoste. 2020. Stochastic weight averaging in parallel: Large-batch training that generalizes well. *International Conference on Learning Representations*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations*.
- 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, and 1 others. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.
- Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. 2023. Llm-blender: Ensembling large language models with pairwise ranking and generative fusion. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14165–14178.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. Tinybert: Distilling bert for natural language understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4163–4174.
- Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and Pengxiang Cheng. 2023. Dataless knowledge fusion by merging weights of language models. In *The Eleventh International Conference on Learning Representations*.

- Simran Khanuja, Melvin Johnson, and Partha Talukdar. 2021. Mergedistill: Merging language models using pre-trained distillation. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2874–2887.
- Dahyun Kim, Chanjun Park, Sanghoon Kim, Wonsung Lee, Wonho Song, Yunsu Kim, Hyeonwoo Kim, Yungi Kim, Hyeonju Lee, Jihoo Kim, and 1 others. 2023. Solar 10.7 b: Scaling large language models with simple yet effective depth up-scaling. *arXiv preprint arXiv:2312.15166*.
- Yoon Kim and Alexander M Rush. 2016. Sequence-level knowledge distillation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327.
- 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. In *The Eleventh International Conference on Learning Representations*.
- Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, and 1 others. 2024. Openassistant conversations-democratizing large language model alignment. *Advances in Neural Information Processing Systems*, 36.
- Nick Littlestone and Manfred K Warmuth. 1994. The weighted majority algorithm. *Information and Computation*, 108(2):212–261.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *International Conference on Learning Representations*.
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2024. Wizardcoder: Empowering code large language models with evol-instruct. In *The Twelfth International Conference on Learning Representations*.
- Michael S Matena and Colin A Raffel. 2022. Merging models with fisher-weighted averaging. *Advances in Neural Information Processing Systems*, 35:17703–17716.
- Costas Mavromatis, Petros Karypis, and George Karypis. 2024. Pack of LLMs: Model fusion at test-time via perplexity optimization. In *First Conference on Language Modeling*.
- Yu Meng, Mengzhou Xia, and Danqi Chen. 2024. SimPO: Simple preference optimization with a reference-free reward. In *Advances in Neural Information Processing Systems*.
- Kristine Monteith, James L Carroll, Kevin Seppi, and Tony Martinez. 2011. Turning bayesian model averaging into bayesian model combination. In *The 2011 International Joint Conference on Neural Networks*, pages 2657–2663. IEEE.
- Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707*.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2023. Gpqa: A graduate-level google-proof q&a benchmark. *arXiv preprint arXiv:2311.12022*.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- Pavel Senin. 2008. Dynamic time warping algorithm review. *Information and Computer Science Department University of Hawaii at Manoa Honolulu, USA*, 855(1-23):40.
- Sainbayar Sukhbaatar, Olga Golovneva, Vasu Sharma, Hu Xu, Xi Victoria Lin, Baptiste Roziere, Jacob Kahn, Shang-Wen Li, Wen tau Yih, Jason E Weston, and Xian Li. 2024. Branch-Train-MiX: Mixing expert LLMs into a mixture-of-experts LLM. In *First Conference on Language Modeling*.
- Siqi Sun, Yu Cheng, Zhe Gan, and Jingjing Liu. 2019. Patient knowledge distillation for bert model compression. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4323–4332.
- Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Well-read students learn better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962*.
- Fanqi Wan, Xinting Huang, Deng Cai, Xiaojun Quan, Wei Bi, and Shuming Shi. 2024. Knowledge fusion of large language models. In *The Twelfth International Conference on Learning Representations*.
- Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. 2024a. Openchat: Advancing open-source language models with mixed-quality data. In *The Twelfth International Conference on Learning Representations*.
- Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances in Neural Information Processing Systems*, 33:5776–5788.

- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, and 1 others. 2024b. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*.
- Zhilin Wang, Yi Dong, Olivier Delalleau, and 1 others. 2024c. HelpSteer2: Open-source dataset for training top-performing reward models. *arXiv preprint arXiv:2406.08673*.
- Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. 2023. Magicoder: Source code is all you need. *arXiv preprint arXiv:2312.02120*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and 1 others. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45.
- Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and 1 others. 2022. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International Conference on Machine Learning*, pages 23965–23998. PMLR.
- Siyue Wu, Hongzhan Chen, Xiaojun Quan, Qifan Wang, and Rui Wang. 2023. Ad-kd: Attribution-driven knowledge distillation for language model compression. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8449–8465.
- Yue Wu, Zhiqing Sun, Huizhuo Yuan, Kaixuan Ji, Yiming Yang, and Quanquan Gu. 2024. Self-play preference optimization for language model alignment. *arXiv preprint arXiv:2405.00675*.
- Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. 2024. Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing. *arXiv preprint arXiv:2406.08464*.
- Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. 2024. Ties-merging: Resolving interference when merging models. *Advances in Neural Information Processing Systems*, 36.
- Shan You, Chang Xu, Chao Xu, and Dacheng Tao. 2017. Learning from multiple teacher networks. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1285–1294.
- Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. 2024a. Language models are super mario: Absorbing abilities from homologous models as a free lunch. In *International Conference on Machine Learning*.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng YU, Zhengying Liu, Yu Zhang, James Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024b. Metamath: Bootstrap your own mathematical questions for large language models. In *The Twelfth International Conference on Learning Representations*.
- Jinghan Zhang, Junteng Liu, Junxian He, and 1 others. 2023. Composing parameter-efficient modules with arithmetic operation. *Advances in Neural Information Processing Systems*, 36:12589–12610.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Sidhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*.
- Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, Karthik Ganesan, Wei-Lin Chiang, Jian Zhang, and Jiantao Jiao. 2024. Starling-7B: Improving helpfulness and harmlessness with RLAI. In *First Conference on Language Modeling*.

A Related Work

Model Fusion Combining the capabilities of diverse models has been a long-standing objective. Existing approaches to model fusion mainly fall into three categories. Firstly, traditional *model ensemble* techniques combine the outputs of multiple models by weighted averaging (Littlestone and Warmuth, 1994) or majority voting (Monteith et al., 2011) to enhance overall system performance. Recently, Jiang et al. (2023) introduced a sequence-level ensemble framework for LLMs, which first conducts pairwise comparisons to rank the outputs of LLMs and then employs another LLM to consolidate the top-ranked candidates into an improved output. In addition to the sequence-level ensemble, Ding et al. (2024) blended multiple LLMs using a token-level gating mechanism on the output logits. To avoid additional training during ensemble, Mavromatis et al. (2024) leveraged the perplexity of different LLMs over input prompts to determine the importance of each model.

Secondly, *model merging* facilitates the fusion of models of identical structure and scale within

the parameter space. [Wortsman et al. \(2022\)](#) combined multiple models, obtained by fine-tuning a model on the same dataset but with distinct strategies, through linear averaging. [Matena and Raffel \(2022\)](#) enhanced simple weighted average by incorporating Fisher Information Matrix ([Fisher, 1922](#)) to determine the significance of individual model parameter. [Jin et al. \(2023\)](#) performed merging by addressing an optimization problem that minimizes the L2 distance between merged and individual models, and conducting a closed-form solution. Although these methods can automatically compute merging coefficients, they necessitate either forward or backward propagation using additional data, making model merging compute-inefficient and memory-intensive. [Ilharco et al. \(2023\)](#) and [Zhang et al. \(2023\)](#) conducted simple arithmetic operations on the task vectors or LoRA ([Hu et al., 2022](#)) modules of different models, thereby enhancing multi-task ability and domain generalization. To mitigate parameter interference, [Yu et al. \(2024a\)](#) and [Yadav et al. \(2024\)](#) introduced sparsification techniques that trim redundant values from task vectors before merging. [Kim et al. \(2023\)](#) and [Akiba et al. \(2024\)](#) advanced the field by merging LLMs across both parameter and data flow spaces, yielding a merged LLM with up-scaled depth.

Thirdly, *mixture of experts* (MoEs) combines specialized expert modules with a sparsely activated mechanism ([Fedus et al., 2022](#)), presenting another venue for model fusion. [Komatsuzaki et al. \(2023\)](#) first proposed initializing a sparse MoEs module using multiple copies from a dense checkpoint. To integrate multiple domain experts, [Sukhbaatar et al. \(2024\)](#) trained multiple domain-specific LLMs from a seed LLM separately and then used feed-forward networks on top of these dense experts to instantiate a sparse MoEs module, followed by further fine-tuning to learn token-level routing. Similarly, [Feng et al. \(2024\)](#) trained multiple domain-specific LoRA ([Hu et al., 2022](#)) modules as experts and combined these domain experts using an explicit sequence-level routing strategy.

Lastly, FUSELLM ([Wan et al., 2024](#)) introduces another paradigm for the fusion of LLMs with structural differences. This approach builds upon knowledge distillation and leverages the probabilistic distribution matrices generated by source LLMs to transfer collective knowledge into a target LLM. Unlike model ensembles and MoEs, knowledge fusion does not require the parallel deployment of multiple models (experts). Compared to model

merging, which only applies to models with identical architectures, FUSELLM allows for the fusion of LLMs with different architectures.

Knowledge Distillation Knowledge fusion essentially performs knowledge distillation to transfer knowledge from source LLMs to a target LLM. Knowledge distillation ([Hinton et al., 2015](#)) aims to train a student model guided by one or more larger teacher models. Previous studies primarily focus on training a student model to mimic the teacher’s behavior in text classification tasks, by replicating the teacher’s output logits ([Sanh et al., 2019](#); [Turc et al., 2019](#)), as well as hidden states ([Sun et al., 2019](#); [Jiao et al., 2020](#)) and relations ([Wang et al., 2020](#)). In the realm of generative models, prevailing approaches maximize the log-likelihood of the student on the distributions ([Khanuja et al., 2021](#); [Gu et al., 2024](#); [Agarwal et al., 2024](#)) or sequences ([Kim and Rush, 2016](#); [Peng et al., 2023](#)) generated by the teacher model. This paradigm can be extended to accommodate multiple teachers by either averaging the distributions ([You et al., 2017](#)) or blending the sequences ([Wang et al., 2024a](#)).

Compared to vanilla knowledge distillation, knowledge fusion of LLMs faces new challenges. Firstly, due to the differences in tokenization among various LLMs, token alignment is essential for transferring knowledge from source to target LLMs. Secondly, when dealing with distributions generated from multiple source LLMs, the fusion function becomes crucial for optimally integrating their distributions. Thirdly, to leverage the unique advantages of different LLMs, it is necessary and challenging to create a compact knowledge fusion dataset that is diverse in capabilities and domains.

B Implementation Details

B.1 Details of Token Alignment

In Figure 5, we present the token pair mappings employed in three distinct token alignment strategies, including EM ([Fu et al., 2023](#)), MinED ([Wan et al., 2024](#)), and our MS. For clarity, these mapping strategies are depicted in a matrix format, where each column represents the probability of a source token being aligned with a corresponding pivot token. The values within these matrices derive from the respective alignment strategies employed. For instance, the matrix \mathbf{W}_{EM} relies on exact matches between source and pivot token pairs, while \mathbf{W}_{MinED} inversely relates to the edit distance between these pairs. \mathbf{W}_{MS} is based on the

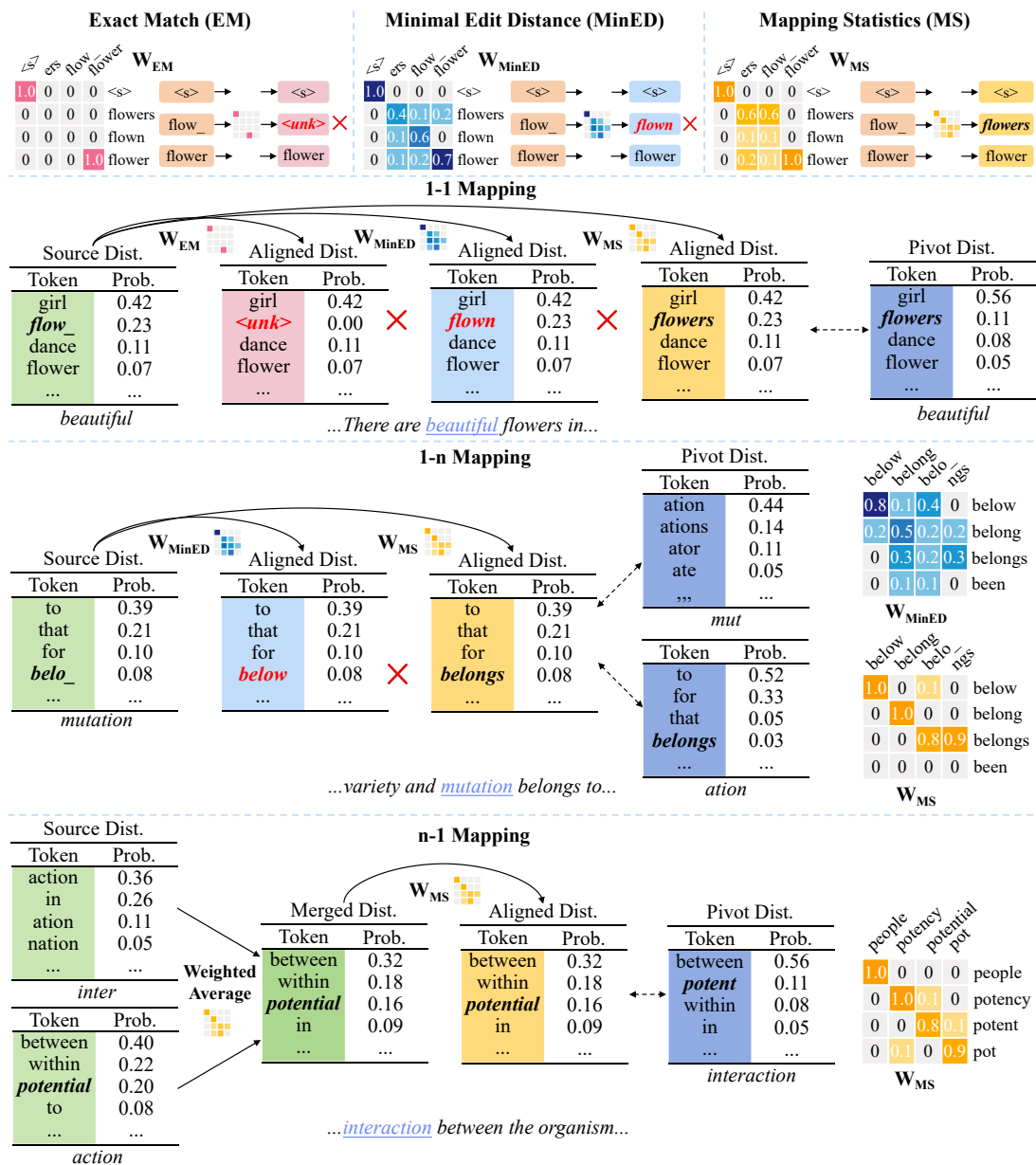


Figure 5: Illustration of EM, MinED, and our MS token alignment strategies in 1-1, 1-n, and n-1 mappings.

statistical mapping frequency between the source and pivot tokens.

In scenarios involving 1-1 or 1-n mappings, the EM and MinED methods utilize W_{EM} or W_{MinED} respectively, which may lead to inaccurate mappings. For example, in EM, the token “flow_” might be incorrectly aligned with “<unk>”, and in MinED, “flow_” could map to “flown”, or “belo_” to “below”. In contrast, our MS method achieves more accurate alignments such as mapping “flow_” to “flowers” and “belo_” to “belongs”, using W_{MS} from sequence-dimensional token alignments. For n-1 mapping, where only MS is applicable, multiple source distributions are combined using a weighted average determined by

W_{MS} . This unified distribution is then processed similarly to the 1-1 mappings.

B.2 Details of Model Merging

The hyperparameters for various merging methods are detailed as follows. For the Linear method (Wortsman et al., 2022), merging parameters are calculated as the mean of all target LLMs. In the TA method (Ilharco et al., 2023), we adhere to the original paper, exploring scaling coefficients within the range of [0.3, 0.4, 0.5]. The optimal setting of 0.3 is selected based on performance. For the TIES (Yadav et al., 2024) and DARE (Yu et al., 2024a) approaches, we search for the trim/drop rate within the range of [0.1, 0.2, ..., 0.9]. The optimal trim/-

Statistics	Math	Extraction	Roleplay	Writing	STEM	Reasoning	Humanities	Coding	Total
Num. Sample	15079	20329	8137	7627	983	7948	1403	27119	88625
Percentage (%)	17.01	22.94	9.18	8.61	1.11	8.97	1.58	30.60	100

Table 6: The domain distribution of samples in the training dataset.

drop rate is 0.4, which results in the elimination of the bottom/random 40% of delta parameters by resetting them to zero. Merging coefficients are computed as the average of all target LLMs. For SCE, we search for the salient element selection thresholds τ within the range of [10, 20, \dots , 90], and the optimal threshold is 10%. Merging coefficients are automatically derived based on the magnitude of delta parameters.

B.3 Details of Training Dataset

We curated a comprehensive training dataset, FUSECHAT-MIXTURE, from various sources. This dataset covers different styles and capabilities, featuring both human-written and model-generated, and spanning general instruction-following and specific skills. These sources include:

Orca-Best²: We sampled 20,000 examples from Orca-Best, which is filtered from the GPT-4 (1M) partition of Orca (Mukherjee et al., 2023) based on maximum length and clustering of instructions.

Capybara³: We incorporated all the 16,000 examples of Capybara, which is a high-quality collection of multi-turn synthetic conversations.

No-Robots⁴: We included all the 9,500 examples of No-Robots, which is a dataset created by skilled human annotators for supervised fine-tuning.

ShareGPT-GPT4⁵: We utilized all 6,200 examples from ShareGPT-GPT4, which exclusively uses dialogues generated by GPT-4 in ShareGPT.

Oasst-Top1⁶: We selected 5,000 examples from Oasst-Top1, which is a refined version of Oasst1 (Köpf et al., 2024), a human-annotated assistant-style conversation dataset.

MetaMathQA⁷: We sampled 10,000 examples

²<https://huggingface.co/datasets/shahules786/orca-best>

³<https://huggingface.co/datasets/LDJnr/Capybara>

⁴https://huggingface.co/datasets/HuggingFaceH4/no_robots

⁵https://huggingface.co/datasets/shibing624/sharegpt_gpt4

⁶https://huggingface.co/datasets/OpenAssistant/oasst_top1_2023-08-25

⁷<https://huggingface.co/datasets/meta-math/>

from MetaMathQA (Yu et al., 2024b), which is augmented from the GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021) datasets for mathematics problem-solving.

OSS-Instruct⁸: We chose 10,000 examples from OSS-Instruct (Wei et al., 2023), which contains code instruction data synthesized from open-source code snippets.

Evol-Alpaca⁹: We sampled 10,000 examples from Evol-Alpaca, which is a code instruction dataset generated by GPT-4 with evol-instruct proposed by WizardCoder (Luo et al., 2024).

Python-Code¹⁰: We selected 10,000 examples from Python-Code, which comprises instructions and responses generated by GPT-3.5 and GPT-4 for Python code generation.

We followed the data processing code in FastChat (Zheng et al., 2024) to clean instances containing non-English or special characters. Then, we split long conversations into blocks with a maximum length of 2048 tokens, resulting in the final FUSECHAT-MIXTURE with 95,000 samples. We also explored the domain distribution of the samples in the training data. Specifically, we used the approach provided by Magpie (Xu et al., 2024), utilizing the Llama-3-8B-Instruct model (Dubey et al., 2024) to classify our 95,000 training examples into eight distinct domains as defined by MT-Bench. After removing 7,000 samples due to anomalous classification errors, the final domain distribution is presented in Table 6, which demonstrates substantial diversity, which aligns with our primary objective to assess the model’s general capabilities rather than domain-specific performance.

B.4 Details of Baselines

In this section, we present the details of the baseline models compared in our experiments.

Proprietary LLMs: GPT-3.5-Turbo-

MetaMathQA

⁸<https://huggingface.co/datasets/ise-uiuc/Magicoder-OSS-Instruct-75K>

⁹<https://huggingface.co/datasets/theblackcat102/evol-codealpaca-v1>

¹⁰<https://huggingface.co/datasets/ajibawa-2023/Python-Code-23k-ShareGPT>

Model	MMLU-Pro	PIQA	BoolQ	GPQA	GSM8K	IFEval	Average
OpenChat-3.5-7B	31.63	82.86	73.91	31.30	76.88	35.73	55.38
OpenChat-3.5-7B SFT	31.32	82.75	73.91	30.30	76.04	35.34	54.94
OpenChat-3.5-7B Multi	31.39	82.43	73.73	32.30	74.75	36.25	55.14
FUSECHAT-7B	31.65(+0.06%)	82.97(+0.13%)	75.50(+2.15%)	37.40(+19.49%)	77.10(+0.29%)	37.49(+4.93%)	57.02(+2.96%)

Table 7: Comparison results on general evaluation benchmarks.

1106¹¹ (Achiam et al., 2023), Claude-3-Opus¹² (Anthropic, 2024), and GPT-4-1106-Preview¹³ (Achiam et al., 2023).

Source LLMs: OpenChat-3.5-7B¹⁴ (Wang et al., 2024a), Starling-LM-7B-alpha¹⁵ (Zhu et al., 2024), NH2-SOLAR-10.7B¹⁶ (Kim et al., 2023), InternLM2-Chat-20B¹⁷ (Cai et al., 2024), Mixtral-8x7B-Instruct¹⁸ (Jiang et al., 2024), and Qwen-1.5-Chat-72B¹⁹ (Bai et al., 2023).

Ensemble LLMs: Top1-PPL (Mavromatis et al., 2024), which selects the 1st ranked response from source LLMs based on the perplexity of the instruction; Top1-LLM-Blender (Jiang et al., 2023), which ranks and combines the output text from source LLMs with ranker and fuser models. Due to the fuser model’s constraints on maximum sequence length, only the ranker model is utilized to determine and produce the 1st-ranked response; Top1-GPT4 (Achiam et al., 2023), which leverages GPT-4 judgment as ranking criteria and yields the 1st-ranked response. Since our evaluations also rely on GPT-4, this approach represents an upper bound for comparison.

Fused LLMs: OpenChat-3.5-7B SFT, a special scenario of knowledge fusion with a single source LLM, serves as the supervised fine-tuning baseline using our training dataset; OpenChat-3.5-7B Multi is the knowledge fusion of multiple source chat LLMs simultaneously like FUSELLM (Wan et al., 2024); OpenChat-3.5-7B Starling, OpenChat-3.5-7B SOLAR, OpenChat-

3.5-7B InternLM, OpenChat-3.5-7B Mixtral, and OpenChat-3.5-7B Qwen are target LLMs resulting from pairwise knowledge fusion of the pivot LLM OpenChat-3.5-7B and the rest source LLMs.

C Evaluation of Additional Benchmarks

The primary objective of FUSECHAT is to integrate multiple chat LLMs into a more powerful model. Consequently, our experiments primarily focus on alignment training data, such as ShareGPT, and chat model evaluation benchmarks like AlpacaEval 2.0 and MT-Bench. In addition to the chat model benchmarks, we also conducted experiments on six general evaluation benchmarks, including MMLU-Pro (Wang et al., 2024b), PIQA (Bisk et al., 2020), BoolQ (Clark et al., 2019), GPQA (Rein et al., 2023), GSM8K (Cobbe et al., 2021), and IFEval (Zhou et al., 2023), which assess knowledge understanding, question-answering, common-sense reasoning, and instruction-following. The results are presented in Table 7. It is important to note that the training data for FUSECHAT-7B is primarily focused on alignment rather than general knowledge. Therefore, performance improvements on these general benchmarks are less significant compared to those on AlpacaEval 2.0 and MT-Bench. This observation is consistent with recent studies on alignment (Meng et al., 2024; Wu et al., 2024), which highlight the critical role of alignment dataset construction in determining downstream performance.

D Rationale Behind the Value of λ

In Eq. 4, λ is set to 0.9 to balance the contributions of the SFT loss and the fusion loss. This value is carefully chosen due to the substantial difference in magnitude between these two losses. To illustrate this, we conducted an experiment using Qwen-1.5-Chat-72B as the source LLM and randomly selected 128 in-

Loss Type	Loss Value
SFT	0.5077
Fusion	1.3081

Table 8: Loss values for SFT and fusion during training with Qwen-1.5-Chat-72B as the source LLM.

¹¹<https://platform.openai.com/docs/models/gpt-3-5-turbo>

¹²<https://www.anthropic.com/news/claude-3-family>

¹³<https://platform.openai.com/docs/models/gpt-4-turbo-and-gpt-4>

¹⁴https://huggingface.co/openchat/openchat_3.5

¹⁵<https://huggingface.co/berkeley-nest/Starling-LM-7B-alpha>

¹⁶<https://huggingface.co/NousResearch/Nous-Hermes-2-SOLAR-10.7B>

¹⁷<https://huggingface.co/internlm/internlm2-chat-20b>

¹⁸<https://huggingface.co/mistralai/Mixtral-8x7B-Instruct-v0.1>

¹⁹<https://huggingface.co/Qwen/Qwen1.5-72B-Chat>

Task	FUSECHAT	FUSELLM
Token Alignment (CPU)	4h	4Nh
Distribution Fusion (CPU)	5h	5Nh
Training (GPU)	9h	9h
Merging (CPU, FuseChat only)	0.1h	0
Total CPU Time	9.1h	9Nh
Total GPU Time	9h	9h
Storage Overhead	$2 \times 2.66\text{GB}$	$(N + 1) \times 2.66\text{GB}$

Table 9: Computational costs and storage overhead of FUSECHAT and FUSELLM when scaling the number of source LLMs from N to $N + 1$.

stances from the training dataset. The observed loss values are presented in Table 8. The results indicate that the fusion loss is approximately three times larger than the SFT loss in this setting. This disparity necessitates assigning a proportionally smaller weight to the fusion loss to prevent it from dominating the optimization process. Without this adjustment, an overly high fusion loss weight could distort the training dynamics, leading to suboptimal learning. Therefore, the 0.9/0.1 weight distribution is a carefully calibrated choice designed to ensure balanced optimization and effective learning within the framework of Eq. 4.

E Cost Analysis

Time Complexity of Token Alignment The token alignment complexity in our approach remains independent of model size, primarily depending on vocabulary size, dataset size, and sequence length, ensuring scalability to larger models. Our token alignment framework is optimized for efficiency. Our experiments were conducted on a hardware setup comprising $8 \times 80\text{GB}$ A800 GPUs and $2 \times \text{Intel Xeon Gold 6348}$ CPUs. To assess the computational cost of token alignment, we aligned Qwen-1.5-72B (vocabulary size: 152,064) with OpenChat-3.5-7B (vocabulary size: 32,002) using approximately 95,000 data samples. The alignment process took approximately 6 CPU hours with 32-thread parallelization. To optimize efficiency, we employed sparse matrix operations, which substantially decreased both computational and storage requirements. The resulting token mapping matrix has a compact storage footprint of approximately 10MB when saved as a .npz file.

Scalability of FUSECHAT and FUSELLM

We detailed the full cost comparison between FUSECHAT and FUSELLM when scaling the number of source LLMs from N to $N + 1$ in Table 9.

Model	AlpacaEval 2.0	MT-Bench
OpenChat-3.5-7B Qwen (D/P)	5.98/ 14.98	6.79/ 7.23
OpenChat-3.5-7B Mixtral (D/P)	16.10/ 16.52	7.03/ 7.24
OpenChat-3.5-7B InternLM (D/P)	6.54/ 15.21	6.88/ 7.21
OpenChat-3.5-7B SOLAR (D/P)	12.21/ 16.51	7.09/ 7.17
OpenChat-3.5-7B Starling (D/P)	14.89/ 16.20	7.15/ 7.22
FUSECHAT-7B (D/P)	14.68/ 17.16	6.91/ 7.38

Table 10: Comparison of pairwise fusion (P) and single-model distillation (D) across five source LLMs, and the final FUSECHAT-7B models merged using the SCE method on AlpacaEval-2.0 and MT-Bench.

The estimated storage overhead for storing the top-10 logits per token for 95K samples with an average length of 500 tokens is approximately 2.66GB per model. Our FuseChat framework demonstrates efficient scalability when incorporating additional source models, all while maintaining reasonable storage and computational demands. While the GPU training time remains comparable to other methods, FuseChat significantly reduces CPU time and storage overhead as the number of source LLMs (N) increases. These results further demonstrate that FUSECHAT can integrate a new source model in a plug-and-play manner.

F Comparison of Pairwise Fusion and Single-model Distillation

The key distinction between pairwise fusion and single-model distillation lies in their learning paradigms. In pairwise fusion, the model selectively acquires knowledge based on the quality of outputs from the source LLM or pivot LLM, guided by lower cross-entropy (CE) values. This approach ensures that the model consistently learns from the stronger performer for each sample. In contrast, single-model distillation relies exclusively on the source LLM, implicitly assuming that the source consistently provides superior outputs.

To rigorously assess the differences between pairwise fusion and single-model distillation, we conducted additional experiments. Specifically, the pairwise fusion strategy in FUSECHAT was replaced with direct distillation from a single source model, omitting the merging phase. The results, summarized in Table 10, demonstrate that pairwise fusion consistently outperforms single-model distillation across five source LLMs. For clarity, the notation D/P indicates the performance of direct distillation and pairwise fusion, respectively. The metrics reported include the Average Score on MT-Bench and the Length-Controlled Win Rate on Al-

Model Comparison	t-statistic	p-value
FUSECHAT-7B vs. Pairwise Fusion	2.95874	0.00318
FUSECHAT-7B vs. OpenChat-3.5-7B Multi	3.32756	0.00108

Table 11: T-test results comparing FUSECHAT-7B with Pairwise Fusion and OpenChat-3.5-7B Multi on MT-Bench, highlighting the statistical significance of performance improvements.

pacaEval 2.0. Furthermore, the SCE method was applied to fuse the models obtained through single-model distillation. As shown in Table 10, the results reveal that merging the models derived from pairwise fusion produces a superior fused model compared to merging models from single-model distillation. These results highlight the effectiveness of pairwise fusion, not only in standalone performance but also in enhancing the quality of the final fused model.

G Statistical Significance of Performance Improvements

We conduct detailed statistical analyses using t-tests to evaluate the performance of our proposed FUSECHAT-7B model on MT-Bench compared to two key baselines: Pairwise Fusion and OpenChat-3.5-7B Multi, which fuses multiple source LLMs simultaneously. The results, summarized in Table 11, highlight the strong statistical significance of the performance improvements achieved by FUSECHAT-7B over these baselines. Notably, the p-values are well below the 0.05 threshold, confirming the significance of the observed differences. These findings provide evidence to support the superiority of FUSECHAT-7B compared to both Pairwise Fusion and OpenChat-3.5-7B Multi.

H Case Studies

We present case studies to demonstrate the individual strengths of target LLMs (OpenChat-3.5-7B Starling, OpenChat-3.5-7B SOLAR, OpenChat-3.5-7B InternLM, OpenChat-3.5-7B Mixtral, and OpenChat-3.5-7B Qwen) obtained from pairwise knowledge fusion of source LLMs, and show the collective knowledge and strengths of FUSECHAT (FUSECHAT-7B) obtained by further merging target LLMs. OpenChat-3.5-7B SFT is used as the baseline for comparison.

OpenChat-3.5-7B SFT; Math; ID 113;

Turn 1 Question: In a survey conducted at a local high school, preferences for a new school color were measured: 58% of students liked the color blue, 45% preferred green, and 22% liked both colors. If we randomly pick a student from the school, what's the probability that they would like neither blue nor green?

Turn 1 Answer: Let's use the principle of inclusion-exclusion to solve this problem. The total percentage of students who like either blue or green is the sum of the individual percentages minus the percentage who like both:

$$58\%(blue) + 45\%(green) - 22\%(both) = 75\%$$

This means that $100\% - 75\% = 25\%$ of students like neither blue nor green. So, the probability that a randomly picked student would not like neither blue nor green is 25% . The answer is: 25

Turn 1 Judgment: The assistant's answer contains a mathematical error in calculating the total percentage of students who like either blue or green. The correct calculation, as shown in the reference answer, should be:

$$58\% + 45\% - 22\% = 81\%$$

This means that 81% of students like at least one of the colors, not 75% as the assistant stated. Consequently, the correct percentage of students who like neither color is:

$$100\% - 81\% = 19\%$$

Therefore, the correct probability that a randomly picked student would like neither blue nor green is 19%, not 25% as the assistant incorrectly concluded. The assistant's mistake significantly affects the correctness of the response. However, the assistant did correctly apply the principle of inclusion-exclusion, which indicates an understanding of the correct approach to solving the problem. The error seems to be purely in the arithmetic calculation. Given the importance of accuracy in mathematical problems and the impact of this error on the final answer, the response's helpfulness is diminished. However, the attempt to apply the correct principle shows some level of understanding. **Rating: [[3]]**

Turn 2 Question: If we select a student liked green, what's the probability that he or she would dislike both colors?

Turn 2 Answer: Since 22% of students like both colors, this means that 78% of students do not like both colors. Out of these, 45% prefer green. Therefore, the probability that a student who likes green would dislike the other color is $45\% \times 78\% = 35.7\%$. So, the probability that a student who likes green would dislike the other color is 35.7% . The answer is: 35.7

Turn 2 Judgment: The assistant's response to the second question contains a fundamental misunderstanding of the question posed. The question asks for the probability that a student who likes green would dislike both colors, which is a contradiction. If a student likes green, they cannot simultaneously dislike green, making it impossible for them to dislike both green and another color (in this case, blue). The correct answer, as provided in the reference answer, is 0%. This is because the premise of selecting a student who likes green excludes the possibility of that student disliking both colors. The assistant's calculation of 35.7% is incorrect because it misinterprets the question as asking for the probability that a student who likes green would dislike the other color (blue), which was not the question asked. Therefore, the assistant's response is incorrect due to a misunderstanding of the question's premise, leading to an erroneous calculation that does not apply to the scenario described. **Rating: [[1]]**

Table 12: Case studies on MT-Bench. The answers are generated by OpenChat-3.5-7B SFT.

OpenChat-3.5-7B Starling; Math; ID 113;

Turn 1 Question: In a survey conducted at a local high school, preferences for a new school color were measured: 58% of students liked the color blue, 45% preferred green, and 22% liked both colors. If we randomly pick a student from the school, what's the probability that they would like neither blue nor green?

Turn 1 Answer: We can solve this problem using the principle of inclusion-exclusion. The total number of students who like either blue or green is the sum of those who like blue, those who like green, minus those who like both (to avoid double-counting). So, the number of students who like either color is:

$$58\%(blue) + 45\%(green) - 22\%(both) = 75\%$$

This means that 25% of students like neither color. Therefore, the probability that a randomly picked student likes neither color is 25%. The answer is: 25

Turn 1 Judgment: The assistant's answer contains a mathematical error in calculating the percentage of students who like either blue or green. The correct calculation, as shown in the reference answer, should be:

$$\text{At least one of the colors} = 58\% + 45\% - 22\% = 103\% - 22\% = 81\%$$

This means that 81% of students like at least one of the colors, not 75% as the assistant stated. Consequently, the percentage of students who like neither color is not 25% but rather:

$$\text{Neither color} = 100\% - 81\% = 19\%$$

Therefore, the correct probability that a randomly picked student likes neither color is 19%, not 25%. The assistant's mistake significantly affects the correctness of the response. The correct approach and calculation were provided in the reference answer, which the assistant failed to replicate accurately. This error directly impacts the helpfulness of the response, as it provides the user with incorrect information. Given the importance of accuracy in mathematical problems and the potential confusion caused by incorrect percentages, the assistant's response is rated lower due to the critical error in calculation. [Rating: \[\[3\]\]](#)

Turn 2 Question: If we select a student liked green, what's the probability that he or she would dislike both colors?

Turn 2 Answer: Since the student dislikes both colors, they must dislike green. Therefore, the probability that a student who likes green would dislike both colors is 0%.

Turn 2 Judgment: The assistant's response to the second question is correct and aligns with the reference answer. It correctly states that if a student likes green, the probability that they would dislike both colors (blue and green) is 0%. This is because the premise of the question already establishes that the student likes green, making it impossible for them to dislike both colors simultaneously. The explanation is clear, concise, and logically sound, directly addressing the user's question without unnecessary complexity. The assistant's answer to the second question does not contain any mistakes and is both correct and helpful. It provides a straightforward explanation that directly answers the user's query, making it easy for the user to understand the concept being discussed. [Rating: \[\[10\]\]](#)

Table 13: Case studies on MT-Bench. The answers are generated by OpenChat-3.5-7B Starling.

OpenChat-3.5-7B SOLAR; Math; ID 113;

Turn 1 Question: In a survey conducted at a local high school, preferences for a new school color were measured: 58% of students liked the color blue, 45% preferred green, and 22% liked both colors. If we randomly pick a student from the school, what's the probability that they would like neither blue nor green?

Turn 1 Answer: We can solve this problem using the principle of inclusion-exclusion. The total number of students who like either blue or green is the sum of those who like blue and those who like green, minus the number who like both. That is:

$$\text{Total} = \text{Blue} + \text{Green} - \text{Both}$$

We are given that 58% of students like blue, 45% like green, and 22% like both. So, we can calculate the percentage of students who like neither as follows:

$$\text{Neither} = \text{Total} - \text{Blue} - \text{Green} + \text{Both}$$

Substituting the given percentages:

$$\text{Neither} = 100\% - (58\% + 45\% - 22\%)$$

Simplifying:

$$\text{Neither} = 100\% - 75\%$$

Therefore, the probability that a randomly picked student would like neither blue nor green is 25%.

Turn 1 Judgment: The assistant's answer contains a critical error in the calculation process. The mistake lies in the incorrect application of the principle of inclusion-exclusion and the subsequent calculation of the percentage of students who like neither color. The correct application of the principle of inclusion-exclusion for calculating the percentage of students who like at least one of the colors (blue or green) is as follows:

$$\text{At least one of the colors} = (\text{Blue}) + (\text{Green}) - (\text{Both})$$

$$\text{At least one of the colors} = 58\% + 45\% - 22\% = 103\% - 22\% = 81\%$$

This means 81% of students like at least one of the colors, blue or green. To find the percentage of students who like neither color, we subtract this percentage from 100%:

$$\text{Neither color} = 100\% - \text{At least one of the colors}$$

$$\text{Neither color} = 100\% - 81\% = 19\%$$

Therefore, the correct probability that a randomly picked student would like neither blue nor green is 19%, not 25% as the assistant incorrectly calculated. The assistant's mistake was in the formula used for calculating the percentage of students who like neither color. The assistant incorrectly attempted to subtract the sum of individual preferences minus the overlap from 100%, which does not align with the principle of inclusion-exclusion correctly applied in this context. Given the significant error in the calculation that led to an incorrect final answer, the response's correctness is compromised. However, the assistant did correctly identify the need to use the principle of inclusion-exclusion, which shows an understanding of the correct approach, albeit misapplied.

Rating: [[3]]

Turn 2 Question: If we select a student liked green, what's the probability that he or she would dislike both colors?

Turn 2 Answer: If a student is selected who likes green, the probability that he or she would dislike both colors is 1 - the probability that they like both colors. We know that 22% of students like both colors. So, the probability that a student who likes green dislikes both colors is:

$$1 - 22\% = 78\%$$

Therefore, the probability that a student who likes green would dislike both colors is 78%.

Turn 2 Judgment: The assistant's response to the second question contains a fundamental misunderstanding of the scenario presented. According to the reference answer, if a student likes green, the probability that they would dislike both colors is 0%. This is logically consistent because liking green precludes the possibility of disliking both green and blue simultaneously. The assistant's answer, suggesting a 78% probability that a student who likes green would dislike both colors, is incorrect. The correct answer should acknowledge that once a student is known to like green, they cannot be in the category of disliking both colors, making the probability 0%. Rating: [[1]]

Table 14: Case studies on MT-Bench. The answers are generated by OpenChat-3.5-7B SOLAR.

OpenChat-3.5-7B InternLM; Math; ID 113;

Turn 1 Question: In a survey conducted at a local high school, preferences for a new school color were measured: 58% of students liked the color blue, 45% preferred green, and 22% liked both colors. If we randomly pick a student from the school, what's the probability that they would like neither blue nor green?

Turn 1 Answer: To find the probability that a student likes neither blue nor green, we need to subtract the probability that they like at least one of the colors from 1. Let A be the event that a student likes blue, and B be the event that a student likes green. We are given that $P(A) = 0.58$, $P(B) = 0.45$, and $P(A \cap B) = 0.22$. First, we need to find the probability that a student likes at least one of the colors. We can use the formula for the probability of the union of two events:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$P(A \cup B) = 0.58 + 0.45 - 0.22 = 0.75$$

Now, to find the probability that a student likes neither blue nor green, we subtract the probability that they like at least one of the colors from 1:

$$P(\text{neither}) = 1 - P(\text{at least one color})$$

$$P(\text{neither}) = 1 - 0.75 = 0.25$$

So, the probability that a randomly picked student would like neither blue nor green is 0.25 or 25%.

Turn 1 Judgment: The assistant's answer contains a mathematical error in calculating the probability that a student likes at least one of the colors. The assistant incorrectly calculates the probability of a student liking at least one of the colors (blue or green) as 0.75 (or 75%), which leads to an incorrect final answer for the probability of a student liking neither color as 0.25 (or 25%). The correct calculation, as shown in the reference answer, should be as follows:

- The probability of liking blue ($P(\text{Blue})$) = 58% or 0.58
- The probability of preferring green ($P(\text{Green})$) = 45% or 0.45
- The probability of liking both colors ($P(\text{Both})$) = 22% or 0.22

Using the principle of inclusion-exclusion:

$$\text{At least one of the colors} = P(\text{Blue}) + P(\text{Green}) - P(\text{Both}) = 0.58 + 0.45 - 0.22 = 0.81$$

This means 81% of students like at least one of the colors, and thus:

$$\text{Neither color} = 1 - \text{At least one of the colors} = 1 - 0.81 = 0.19$$

Therefore, the correct probability that a randomly picked student would like neither blue nor green is 0.19 or 19%, not 25% as the assistant's answer suggests. The assistant's mistake lies in the incorrect calculation of the union of the two events, which led to an inaccurate conclusion. This error significantly impacts the correctness of the response. Given the importance of accuracy in probability calculations and the impact of this error on the final answer, the assistant's response is rated lower for its correctness. However, the assistant did correctly apply the principle of inclusion-exclusion and provided a clear explanation of the steps, which is helpful for understanding the process. [Rating: \[\[4\]\]](#)

Turn 2 Question: If we select a student liked green, what's the probability that he or she would dislike both colors?

Turn 2 Answer: The question seems to be asking for the probability that a student who likes green would dislike both colors, which is not possible because the student already likes green. The student cannot dislike both colors while also liking one of them.

Turn 2 Judgment: The assistant's response to the second question is correct and aligns with the reference answer. It correctly identifies the logical inconsistency in the question: if a student likes green, they cannot simultaneously dislike both green and another color. The assistant succinctly explains why the probability that a student who likes green would dislike both colors is 0%, which is the same conclusion reached in the reference answer. This explanation is clear, logically sound, and directly addresses the user's question, demonstrating an understanding of basic probability principles and the specific scenario presented. [Rating: \[\[10\]\]](#)

Table 15: Case studies on MT-Bench. The answers are generated by OpenChat-3.5-7B InternLM.

OpenChat-3.5-7B Mixtral; Math; ID 113;

Turn 1 Question: In a survey conducted at a local high school, preferences for a new school color were measured: 58% of students liked the color blue, 45% preferred green, and 22% liked both colors. If we randomly pick a student from the school, what's the probability that they would like neither blue nor green?

Turn 1 Answer: The probability that a student likes blue is 58%, and the probability that they like green is 45%. Since 22% of students like both colors, we can subtract this from the total percentage to find the percentage of students who like either blue or green but not both.

$$58\% + 45\% - 22\% = 71\%$$

So, 71% of students like either blue or green, which means $100\% - 71\% = 29\%$ of students like neither color. Therefore, the probability that a randomly picked student would like neither blue nor green is 29%. The answer is: 29

Turn 1 Judgment: The assistant's answer contains a critical mathematical error in calculating the probability of students liking neither blue nor green. The correct approach, as outlined in the reference answer, involves using the principle of inclusion-exclusion to accurately account for the overlap between students who like both colors. Correct Calculation:

- The correct calculation to find the percentage of students who like at least one of the colors (blue or green) is to add the percentages of students who like each color and then subtract the percentage of students who like both (to correct for double counting), which is $58\% + 45\% - 22\% = 81\%$.
- Therefore, to find the percentage of students who like neither color, we subtract the percentage of students who like at least one of the colors from 100%, resulting in $100\% - 81\% = 19\%$.

The assistant incorrectly calculated the percentage of students who like either blue or green but not both as 71%, leading to the incorrect conclusion that 29% of students like neither color. The correct probability that a randomly picked student would like neither blue nor green is 19%, not 29%. Given the significant mathematical error and the resulting incorrect conclusion, the assistant's answer is not helpful for someone trying to understand how to correctly apply the principle of inclusion-exclusion in probability. [Rating: \[\[2\]\]](#)

Turn 2 Question: If we select a student liked green, what's the probability that he or she would dislike both colors?

Turn 2 Answer: If a student is selected who likes green, the probability that they dislike both colors is equal to the percentage of students who like both colors. In this case, 22% of students like both colors, so the probability that a student who likes green would dislike both colors is 22%. The answer is: 22

Turn 2 Judgment: The assistant's response to the second question contains a significant error. The correct answer, as outlined in the reference answer, is that the probability of a student who likes green disliking both colors is 0%. This is because the premise of the question already establishes that the student likes green, making it impossible for them to dislike both green and blue simultaneously. The assistant's answer incorrectly states that the probability is 22%, which seems to be a misunderstanding of the question's premise. The correct approach should have recognized the logical contradiction in the question: a student who likes green cannot at the same time dislike both green and another color. Therefore, the probability of a student who likes green also disliking both colors is 0%. Given this significant mistake and misunderstanding of the question's premise, the response is incorrect and misleading. [Rating: \[\[1\]\]](#)

Table 16: Case studies on MT-Bench. The answers are generated by OpenChat-3.5-7B Mixtral.

OpenChat-3.5-7B Qwen; Math; ID 113;

Turn 1 Question: In a survey conducted at a local high school, preferences for a new school color were measured: 58% of students liked the color blue, 45% preferred green, and 22% liked both colors. If we randomly pick a student from the school, what's the probability that they would like neither blue nor green?

Turn 1 Answer: To find the probability that a student likes neither blue nor green, we can use the formula for the probability of the union of two events:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

In this case, A represents the event that a student likes blue, and B represents the event that a student likes green. We are given that $P(A) = 0.58$, $P(B) = 0.45$, and $P(A \cap B) = 0.22$. First, we need to find the probability of the event that a student likes either blue or green, which is represented by $P(A \cup B)$:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) = 0.58 + 0.45 - 0.22 = 0.75$$

Now, we need to find the probability of the event that a student likes neither blue nor green, which is represented by $P(A' \cap B')$:

$$P(A' \cap B') = 1 - P(A \cup B) = 1 - 0.75 = 0.25$$

Therefore, the probability that a randomly picked student would like neither blue nor green is 0.25 or 25%.

Turn 1 Judgment: The assistant's answer contains a mathematical error in calculating the probability of a student liking either blue or green. The assistant incorrectly calculates the sum of the probabilities as 0.75 (or 75%) instead of correctly calculating it as 81%, as shown in the reference answer. This error leads to an incorrect final answer for the probability of a student liking neither blue nor green. Correct Calculation based on the reference answer:

- The correct calculation for the probability of a student liking at least one of the colors (blue or green) is $58\% + 45\% - 22\% = 81\%$.

- Therefore, the correct probability of a student liking neither color is $100\% - 81\% = 19\%$.

The assistant's answer incorrectly concludes that the probability is 25%, whereas the correct probability, as per the reference answer and the correct calculation, is 19%. Despite the assistant's method being fundamentally sound and the explanation clear, the mathematical mistake significantly impacts the correctness of the final answer. The assistant's approach correctly applies the principle of inclusion-exclusion but fails in the execution of basic arithmetic. Given the importance of accuracy in mathematical problems and the impact of this error on the final answer, the assistant's response must be rated lower. However, the explanation and approach were correctly aligned with the principle of inclusion-exclusion, which is a positive aspect of the response. **Rating:** [\[\[4\]\]](#)

Turn 2 Question: If we select a student liked green, what's the probability that he or she would dislike both colors?

Turn 2 Answer: If a student is selected who likes green, we want to find the probability that they dislike both colors. Since the student likes green, they cannot dislike both colors. Therefore, the probability that a student who likes green would dislike both colors is 0.

Turn 2 Judgment: The assistant's response to the second question is correct and aligns with the reference answer. It accurately explains that if a student is already known to like green, it is impossible for them to dislike both green and blue simultaneously. Therefore, the probability that a student who likes green would dislike both colors is indeed 0%. This explanation is clear, concise, and directly addresses the user's question without unnecessary complexity. The assistant's answer is helpful in clarifying the logical impossibility of the scenario proposed by the user, thus preventing any confusion about the nature of probability in this context. **Rating:** [\[\[10\]\]](#)

Table 17: Case studies on MT-Bench. The answers are generated by OpenChat-3.5-7B Qwen.

FUSECHAT-7B; Math; ID 113;

Turn 1 Question: In a survey conducted at a local high school, preferences for a new school color were measured: 58% of students liked the color blue, 45% preferred green, and 22% liked both colors. If we randomly pick a student from the school, what's the probability that they would like neither blue nor green?

Turn 1 Answer: To find the probability that a student likes neither blue nor green, we can use the formula:

$$P(\text{neither}) = 1 - P(\text{blue}) - P(\text{green}) + P(\text{both})$$

where $P(\text{blue})$ is the probability that a student likes blue, $P(\text{green})$ is the probability that a student likes green, and $P(\text{both})$ is the probability that a student likes both colors. We are given that $P(\text{blue}) = 58\%$, $P(\text{green}) = 45\%$, and $P(\text{both}) = 22\%$. Plugging these values into the formula, we get:

$$P(\text{neither}) = 1 - 0.58 - 0.45 + 0.22$$

$$P(\text{neither}) = 1 - 0.85 + 0.22$$

$$P(\text{neither}) = 0.15$$

So, the probability that a randomly picked student would like neither blue nor green is 15%.

Turn 1 Judgment: The assistant's answer uses the principle of inclusion-exclusion correctly but makes a calculation error in the final step. According to the assistant's formula:

$$P(\text{neither}) = 1 - P(\text{blue}) - P(\text{green}) + P(\text{both})$$

The values given are: - $P(\text{blue}) = 58\% = 0.58$ - $P(\text{green}) = 45\% = 0.45$ - $P(\text{both}) = 22\% = 0.22$

Plugging these into the formula, the assistant calculates:

$$P(\text{neither}) = 1 - 0.58 - 0.45 + 0.22$$

However, the assistant incorrectly calculates the final result as 0.15 (or 15%). The correct calculation should be:

$$P(\text{neither}) = 1 - 0.58 - 0.45 + 0.22 = 1 - 1.03 + 0.22 = 1 - 0.81 = 0.19$$

Therefore, the correct probability that a randomly picked student would like neither blue nor green is 19%, not 15%. The reference answer correctly calculates this probability as 19% using the principle of inclusion-exclusion and correctly accounts for the overlap between students who like both colors. The assistant's approach is fundamentally correct but falters at the arithmetic stage, leading to an incorrect final answer. Given the mistake in calculation, the assistant's answer is not entirely correct, but since the method used was appropriate, it shows a good understanding of the principle involved. Therefore, the assistant's response, while flawed in execution, is still somewhat helpful for understanding how to approach the problem. **Rating: [[6]]**

Turn 2 Question: If we select a student liked green, what's the probability that he or she would dislike both colors?

Turn 2 Answer: If a student is selected who likes green, we want to find the probability that they dislike both colors. Since the student likes green, they cannot dislike both colors. Therefore, the probability that a student who likes green would dislike both colors is 0%.

Turn 2 Judgment: The assistant's answer to the second question is correct and aligns with the reference answer. It accurately explains that if a student likes green, the probability of them disliking both colors (blue and green) is 0%. This is because the premise of the question already establishes that the student likes green, making it impossible for them to dislike both colors simultaneously. The explanation is clear, concise, and logically sound, effectively addressing the user's query without any errors or unnecessary complexity. **Rating: [[10]]**
