

# TinyR1-32B-Preview: Boosting Accuracy with Branch-Merge Distillation

## Qiyuan Tech:

Lin Sun, Guangxiang Zhao, Xiaoqi Jian, Weihong Lin, Yongfu Zhu,  
Change Jia, Linglin Zhang, Jinzhu Wu, Sai-er Hu, Xiangzheng Zhang

## Peking University:

Yuhan Wu, Junfeng Ran, Zihan Jiang, Junting Zhou, Wenrui Liu, Bin Cui, Tong Yang

## Abstract

The challenge of reducing the size of Large Language Models (LLMs) while maintaining their performance has gained significant attention. However, existing methods such as model distillation and transfer learning often fail to achieve high accuracy. To address this limitation, we introduce the Branch-Merge distillation approach, which enhances model compression through two phases: (1) the Branch Phase, where knowledge from a large teacher model is *selectively distilled* into specialized student models via domain-specific supervised fine-tuning (SFT); And (2) the Merge Phase, where these student models are merged to enable cross-domain knowledge transfer and improve generalization. We validate our distillation approach using DeepSeek-R1 as the teacher and DeepSeek-R1-Distill-Qwen-32B as the student. The resulting merged model, TinyR1-32B-Preview, outperforms its counterpart DeepSeek-R1-Distill-Qwen-32B across multiple benchmarks, including Mathematics (+5.5 points), Coding (+4.4 points), and Science (+2.9 points), while achieving near-equal performance to DeepSeek-R1 on AIME 2024. The Branch-Merge distillation approach provides a scalable solution for creating smaller, high-performing LLMs with reduced computational cost and time.

## 1 Introduction

Recently, the DeepSeek-R1 model has achieved great success, and its released R1-Distill models (DeepSeek-AI, 2025) demonstrated that distilled small models can be superior in reasoning. Building smaller-scale models is also beneficial for deployment and reducing inference costs. However,

developing smaller yet powerful models is a key challenge in Large Language Models (LLMs).

The most effective method, to our knowledge, is distilling a smaller model from a bigger teacher model across various domains (Jiao et al., 2020; DeepSeek-AI, 2025; Team, 2025a; Muennighoff et al., 2025). However, this method has a fundamental limitation: it requires carefully selecting the most relevant data/domains and tuning their proportions for joint training, which is typically time-consuming and error-prone (Guo et al., 2019; Ji et al., 2024). Furthermore, optimizing many domains simultaneously can lead to conflicting gradients, where tasks interfere, impeding overall learning progress (Yu et al., 2020; Jiang et al., 2024). These problems limit the effectiveness and efficiency of naive data mixed distillation, often resulting in models that cannot achieve the performance levels desired for specialized tasks.

To address these issues and optimize performance across multiple areas, we propose an approach, namely branch-merge, which integrates a model-merging technique during the distillation. Our branch-merge distillation approach contains two phases as follows.

- **Branch Phase:** Knowledge is selectively distilled from a unified large teacher model (*e.g.*, DeepSeek-R1 671B) to instruct several specialized student models (*e.g.*, math, coding, science) through domain-specific SFT.
- **Merge Phase:** The specialized models are combined into a single unified model, enabling cross-domain knowledge transfer while preserving their original specialized capabilities.

## Contributions

- **Accuracy:** We can see from Figure 1 that the Branch-Merge distillation approach significantly improves model accuracy, approaching the scores

<sup>1</sup>Lin Sun, Guangxiang Zhao, Xiaoqi Jian, Yuhan Wu, Weihong Lin, and Yongfu Zhu contributed equally. Corresponding authors: Tong Yang and Xiangzheng Zhang {yangtong@pku.edu.cn, zhangxiangzheng@360.cn}

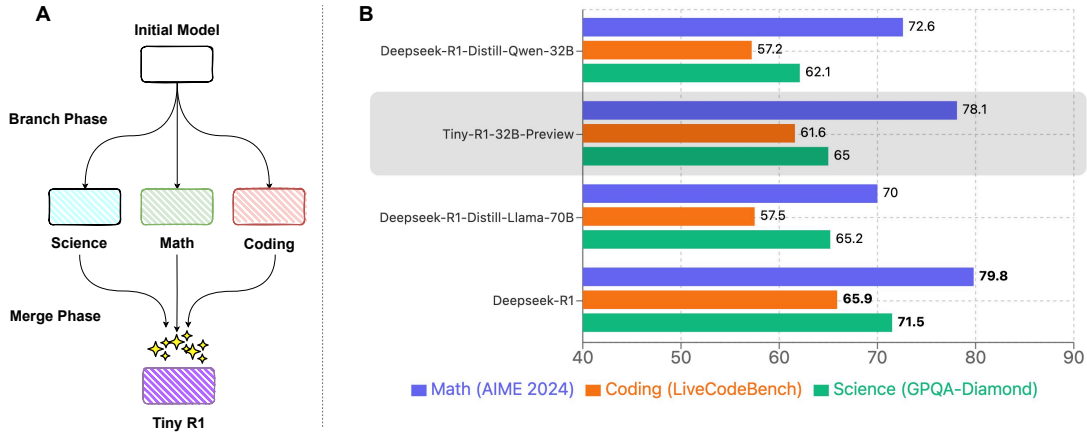


Figure 1: (A) A simplified diagram of our Branch-Merge distillation approach. (1) In the Branch phase, each copy of the Initial Model (backbone) is trained on knowledge from a different domain; (2) In the Merge phase, models are merged based on Arcee Fusion rules. (B) Performance Comparison of different LLM models (Mustar, 2025). TinyR1-32B-Preview outperforms distilled models of the same size in science, math, and coding and achieves comparable results to Deepseek R1. LiveCodeBench here refers to the 24.08-25.02 subset of full LiveCodeBench.

of the R1 teacher model, which traditional distillation methods have not yet achieved. Our distilled Qwen-32B model surpasses DeepSeek-R1-Distill-Qwen-32B with about 5% more accuracy, and its math accuracy approaches that of the original R1 teacher.

- **Simplicity & Low Cost:** Branch-Merge distillation approach significantly reduces the time and computational costs of the merging stage. Compared to traditional methods, we save 90% of the time in the merging phase (0.5 hours with 4 H800 GPUs vs. 23 hours with 32 H800 GPUs for merged data retraining). The ideal reproduction cost for TinyR1-32B-Preview is 744 H800 GPU hours, approximately \$1500 (excluding ablation experiments and parameter search).
- **Openness:** We stand on the shoulders of giants in the open-source community and aim to give back. We will release our model and all data, training code, evaluation code, and logs so anyone can reproduce our results.

## 2 The Branch-Merge Distillation Approach

This section describes our branch-merge distillation approach (as shown in Figure 1A), which consists of two phases: Branch and Merge. This two-phase distillation strategy directly addresses the issues of data selection and gradient conflict by decoupling training domains (Branch) and then rec-

onciling them (Merge). Each phase is described in detail below.

### 2.1 The Branch Phase

In the Branch phase, we first constructed separate datasets for math, science, and coding. Then, we fine-tuned DeepSeek-R1-Distill-Qwen-32B on each dataset using SFT, resulting in three specialized expert models.

- **Math:** We sift 58k samples from 94k questions in NuminaMath1.5 (LI et al., 2024) with corresponding solutions from OpenR1 (Team, 2025a) trajectories. The selection is based on three aspects: *question\_type*, *source*, and *correctness\_math\_verify*. Comparative experiments on DeepSeek-R1-Distill-Qwen-14B indicate that these factors have nearly no impact on the results. Ultimately, we adopted a minimal dataset while maintaining comparable performance.
- **Coding:** The OpenThoughts (Team, 2025b) dataset is filtered to form 20k trajectories of coding solutions. An additional modification is replacing “<|begin\_of\_thought|>” in the original dataset with “<think>” and “<|end\_of\_solution|>” with “</think>”.
- **Science:** DeepSeek-R1 generates 1 CoT trajectory for each of the 8.6k seed examples (2.7k from the science and health science subsets of data\_ablation\_full59k in S1 (Muennighoff et al., 2025), 1.0k from S1k (Muennighoff et al., 2025),

4.9k from the science subset of OpenThoughts (Team, 2025b)), resulting in 8.6k CoT trajectories.

We apply SFT on DeepSeek-R1-Distill-Qwen-32B with the three datasets to obtain three specialized models. Detailed experiment setup will be discussed in Section 3.1.

## 2.2 The Merge Phase

In the Merge phase, we use Arcee Fusion (Goddard et al., 2024) to merge models from different domains. This technique selectively integrates meaningful parameter updates from a *teacher* model  $T$  into a *student* model  $S$ , rather than simply averaging parameters. It proceeds in three stages:

- 1) **Importance Scoring:** Compute the importance of each parameter as follows: First, compute the parameter difference  $T - S$  and the KL divergence  $\text{KL}(T, S)$ . Next, obtain the Importance Score (IS) by multiplying them:

$$\text{IS} = \text{KL}(T, S) \cdot (T - S).$$

The importance score serves as a selection criterion for parameter updates.

- 2) **Dynamic Selection:** Determine a threshold THR of mask via the median Med and interquartile range IR of the importance scores:

$$\text{THR} = \text{Med} + \theta \cdot \text{IR},$$

where  $\theta$  is a balancing coefficient. Since THR is derived from IS statistics, it dynamically adapts to model parameters, ensuring an optimal update ratio. MergeKit’s developers determined through extensive experiments that the optimal value for  $\theta$  is 1.5.

- 3) **Selective Integration:** Apply a merging mask Mask to integrate only parameters exceeding the threshold THR, outputting the merged model  $M$ :

$$M = S + \text{Mask}(T - S). \quad \text{Mask} = \mathbb{I}_{\text{IS} > \text{THR}}.$$

Only teacher parameters with an importance score above the threshold are retained; otherwise, student model parameters are kept.

By focusing on the most significant changes, Arcee Fusion avoids over-updating and maintains model stability. Although this method merges only

two models at a time, our work involves three models, and we detail the corresponding merge sequence in Section 3.3. A comparison with other model merging methods appears in Figure 2. We compare various merging methods on merging models trained from the math and science domains, and we find that Arcee achieves the highest scores on GPQA-Diamond. We found a similar method ranking on the AIME 2024 benchmark, but no separate graph was drawn due to space limitations.

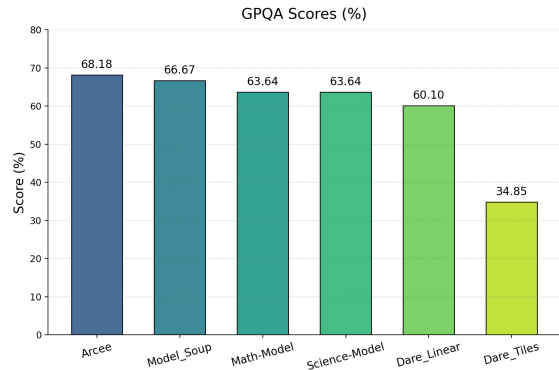


Figure 2: Performance Comparison of merged models on the GPQA-Diamond benchmark.

## 3 Experiment

We choose Deepseek-R1 (DeepSeek-AI, 2025) and its distilled DeepSeek-R1-Distill-Qwen-32B and DeepSeek-R1-Distill-Llama-70B as baselines. Additionally, we conducted a comprehensive ablation study. We compared TinyR1-32B-Preview with: (a) three domain expert models (Math Expert, Coding Expert, Science Expert) before merging; (b) a ‘Data Mixture’ model trained on a combined Math/Coding/Science dataset; and (c) variants of our model obtained via different merging sequences.

### 3.1 Experiment Setup

#### 3.1.1 Training Details

We employ DeepSeek-R1-Distill-Qwen-32B as our backbone model. Leveraging the 360-Llama-Factory (Zou et al., 2024) training framework, we develop three domain-specific expert models applying 16384 sequence length with constructed Math, Coding, and Science datasets.

- **Math Expert:** The math expert model is trained with 5 epochs, batch size 96, and the learning rate is set to constant  $1e-5$ .

Model	Math (AIME 2024)	Coding (LiveCodeBench 24.08-25.02)	Science (GPQA-Diamond)
DeepSeek-R1-Distill-Qwen-32B <sup>†</sup>	72.6 (9.6k Tokens)	57.2 (10.1k Tokens)	62.1 (5.3k Tokens)
DeepSeek-R1-Distill-Llama-70B <sup>†</sup>	70.0	57.5	65.2
DeepSeek-R1 <sup>†</sup>	79.8 (9.6k Tokens)	65.9 (10.4k Tokens)	71.5 (5.3k Tokens)
TinyR1-32B-Preview (Ours)	78.1 (11.8k Tokens)	61.6 (12.4k Tokens)	65.0 (8.6k Tokens)

Table 1: Performance comparison on benchmark datasets. All scores are reported as pass@1. Scores reported from DeepSeek-R1 paper (DeepSeek-AI, 2025) are noted with <sup>†</sup>. The number in parentheses represents the average output token length (including the chain of thought), obtained from our testing.

- **Science Expert:** The science expert model is trained with 5 epochs, batch size 32 with the neat packing mechanism (Tay et al., 2020; Henry et al., 2019; Dean and et al., 2018), and the learning rate is set to cosine 1e-5.
- **Coding Expert:** The coding expert model is trained with 15 epochs, batch size 96 with the neat packing mechanism, and the learning rate is set to constant 1e-5.

We merged the models trained separately in three fields into a single model. We use the Arcee merging (Goddard et al., 2024) method with the  $\theta=1.5$  and threshold THR=0.5. We will compare different model merging methods in Section 3.3.

### 3.1.2 Evaluation Details

For evaluation, we compare the performance on three benchmark datasets: AIME 2024 for Math, LiveCodeBench (24.08-25.02) for Coding, and GPQA-Diamond for Science. The accuracy is calculated as an average pass@1 across 16, 4, and 4 independent trials for these benchmarks, respectively. Meanwhile, we did not use a greedy way to evaluate the model due to its long-COT output, we set the max tokens to 32768 and evaluated the models with Temperature=0.6 and Top-p=0.95 as recommended in DeepSeek-R1 (DeepSeek-AI, 2025). We tried various open-source frameworks for the evaluation on livecodebench and ultimately selected the evaluation code from FuseAI (Wan et al., 2024) utilizing the vLLM implementation, as it can reproduce the effects of the DeepSeek-R1 and its distilled models.

## 3.2 Main Results

We compare our TinyR1-32B-Preview model and other models in Table 1.

- In terms of accuracy, we significantly outperform our backbone model, DeepSeek-R1-Distill-Qwen-32B (Math +5.5, Coding

+4.4, Science +2.9), and generally surpass DeepSeek-R1-Distill-Llama-70B (Math +8.1, Coding +4.1, Science -0.2), approaching the performance of DeepSeek-R1 (Math -1.7, Coding -4.3, Science -6.5).

- In terms of inference cost, our model generates slightly more output tokens than R1 (Math +23%, Coding +19%, Science +62%).
- In terms of the model’s parameter size, our model is smaller compared to DeepSeek-R1, making it more suitable for local deployment by users and small groups.

## 3.3 Ablation Study

As shown in Table 2, we made a comprehensive ablation study to explore if our merging distillation approach works.

Compared to the domain-specific experts, the Data Mixture model surpasses them in math and science but shows decreased performance in coding. This is a seesaw problem that traditional data mixture is difficult to solve (Zamir and Arbeláez, 2018; Liu and Yao, 2019; Radford and W., 2021). In comparison, our merged model improves performance in math and science domains while largely retaining the coding capability. Compared to the Data Mixture, the coding scores of merged models only slightly decreased. We also compared two different model-merging sequences: (1) first merging Math with Coding, then merging with Science; (2) first merging Math with Science, then with Coding. The results show only minor performance differences between these sequences, with the latter yielding a slightly higher average accuracy.

In addition, the average Merging time (GPU hours) of the current Data Mixture is 740 hours<sup>1</sup>.

<sup>1</sup>Note that the merging time is the SFT experiment time, excluding the SFT time for single experts and the downstream evaluation time of the Data Mixture model

Model	Math (AIME 2024)	Coding (LiveCodeBench)	Science (GPQA-Diamond)	Merging Time (GPU Hours)
Math Expert	73.1	-	-	-
Coding Expert	-	63.4	-	-
Science Expert	-	-	64.5	-
Data Mixture	75.3	61.0	<b>65.7</b>	740 h
Merging: (Math & Coding) & Science	77.3	<b>63.8</b>	64.0	<b>4h</b>
Merging: (Math & Science) & Coding	<b>78.1</b>	61.6	65.0	<b>4h</b>
TinyR1-32B-Preview	<b>78.1</b>	61.6	65.0	<b>4h</b>

Table 2: Performance comparison between backbone experts, the data-mixture model, and merged model. All scores are reported as pass@1. LiveCodeBench here refers to the 24.08-25.02 subset of full LiveCodeBench.

On the contrary, the average merging time of TinyR1-32B-Preview is 4 hours. In summary, we only used 0.5% of the Data-Mixture computational overhead on merging models to surpass the effect of traditional data mixture methods. In addition to reducing computational overhead, model merging significantly accelerates our model release process by avoiding the delays introduced by mixed-data re-SFT on the development model. The model merging is a “free-lunch” approach, as it reduces costs and increases efficiency at the same time.

## 4 Related Work

### 4.1 Model Distillation

Knowledge Distillation (KD) (Romero et al., 2015; Hinton et al., 2015) has been proposed to create cheaper strong models (Gou et al., 2021; Hu et al., 2023; Yang et al., 2024; Xu et al., 2024). Primarily, recognizing the disparities between proprietary and open-source LLMs, KD techniques have surged to bridge the performance gap between these models. Distillation methods can be categorized into two types: (1) the logits-based methods (Hinton et al., 2015), which transfer knowledge at the logits level, and (2) the feature-based methods (Romero et al., 2015), which transmit knowledge through intermediate features.

Compared to traditional knowledge distillation techniques (Gou et al., 2021), data augmentation (DA) (Feng et al., 2021) has emerged as an effective method for distilling knowledge in large language models (LLMs). In this approach, a small seed of knowledge is used to prompt LLMs, enabling them to generate additional data tailored to specific domains or skills (Taori et al., 2023). More recently, an API-based strategy has gained attention, where open-source LLMs serve as teachers to self-improve through self-distillation (Yuan et al.,

2024; Chen et al., 2024). By applying a range of distillation techniques, this strategy effectively narrows the performance gap between closed-source and open-source models (Chiang et al., 2023; Xu et al., 2023). In this context, the method involves using only the outputs of the teacher model via an API. This strategy includes approaches such as In-Context Learning (Huang et al., 2022), Chain-of-Thought (Li et al., 2022b), and Instruction Following (Wang et al., 2023). In specialized fields, like science (Zhang et al., 2024), where domain-specific knowledge and accuracy are essential, distillation allows open-source models to significantly improve their performance by learning from proprietary models that are extensively trained and fine-tuned in these domains.

### 4.2 Model Merging

Recent advances in model merging have explored diverse strategies (Ilharco et al., 2022a; Yadav et al., 2023; Davari and Belilovsky, 2024; Deep et al., 2024) to combine neural network parameters while preserving or enhancing performance. Early approaches focused on linear interpolation techniques, such as weight averaging (Wortsman et al., 2022), where models finetuned from shared pretrained checkpoints are merged via arithmetic mean. While computationally efficient, these methods assume approximate parameter space alignment and often degrade when models exhibit divergent optimization trajectories (Frankle et al., 2020; Izmailov et al., 2018; Neyshabur et al., 2020; Fort et al., 2020; Wortsman et al., 2022; Choshen et al., 2022; Ilharco et al., 2022b).

Theoretical underpinnings for these methods derive from studies on loss landscape geometry (Ilharco et al., 2022a; Li et al., 2018; Garipov et al., 2018; Draxler et al., 2018; Kuditipudi et al., 2019; Fort et al., 2019; Czarnecki et al., 2019; Wortsman

et al., 2021; Benton et al., 2021; Entezari et al., 2021; Li et al., 2022a; Lubana et al., 2023). Research on flat local minima (Kaddour et al., 2022; Wortsman et al., 2022; Keskar et al., 2016; Dziugaite and Roy, 2017) dating back from the 1990s (Hochreiter and Schmidhuber, 1994, 1997) suggests that averaged weights reside in flatter regions of the loss surface, correlating with improved out-of-distribution generalization. Further analyses (Daheim et al., 2023; Matena and Raffel, 2022) formalize model merging as identifying connected basins in parameter space, where interpolated solutions maintain low loss. Empirical validations, such as model soups (Wortsman et al., 2022), corroborate that aggregated weights often outperform individual models, particularly under distribution shifts.

## 5 Conclusion and Future Work

We introduce TinyR1-32B-Preview, a model using the Branch-Merge distillation approach to boost reasoning accuracy while preserving efficiency. We achieve significantly higher accuracy than our backbone model, DeepSeek-R1-Distill-Qwen-32B, and generally outperform DeepSeek-R1-Distill-Llama-70B while approaching the performance of DeepSeek-R1. Although our model generates slightly more output tokens than R1, its smaller parameter size makes it more efficient and better suited for local deployment by users and small groups.

Potential future directions include:

- **Exploring Alternative Backbones** – For instance, conducting SFT with the Qwen-Instruct model as the backbone. Our preliminary experiments with Qwen-14B-Instruct and Qwen-32B-Instruct for specialized tasks have already yielded similar results.
- **Releasing Models of Various Sizes** – Expanding our model lineup to accommodate different needs.
- **Investigating the Impact of Experiment Details** – Further analyzing how various experiment settings influence final performance.

## References

Gregory Benton, Wesley Maddox, Sanae Lotfi, and Andrew Gordon Gordon Wilson. 2021. Loss surface

simplexes for mode connecting volumes and fast ensembling. In *International Conference on Machine Learning*, pages 769–779. PMLR.

Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. 2024. Self-play fine-tuning converts weak language models to strong language models.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality.

Leshem Choshen, Elad Venezian, Noam Slonim, and Yoav Katz. 2022. Fusing finetuned models for better pretraining. *arXiv preprint arXiv:2204.03044*.

Wojciech Marian Czarnecki, Simon Osindero, Razvan Pascanu, and Max Jaderberg. 2019. A deep neural network’s loss surface contains every low-dimensional pattern. *arXiv preprint arXiv:1912.07559*.

Nico Daheim, Thomas Möllenhoff, Edoardo Maria Ponti, Iryna Gurevych, and Mohammad Emtiyaz Khan. 2023. Model merging by uncertainty-based gradient matching. *arXiv preprint arXiv:2310.12808*.

MohammadReza Davari and Eugene Belilovsky. 2024. Model breadcrumbs: Scaling multi-task model merging with sparse masks. In *European Conference on Computer Vision*, pages 270–287. Springer.

J. Dean and et al. 2018. Scaling neural networks with efficient memory and batching. In *ICML*.

Pala Tej Deep, Rishabh Bhardwaj, and Soujanya Poria. 2024. Della-merging: Reducing interference in model merging through magnitude-based sampling. *arXiv preprint arXiv:2406.11617*.

DeepSeek-AI. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning.

Felix Draxler, Kambis Veschgini, Manfred Salmhofer, and Fred Hamprecht. 2018. Essentially no barriers in neural network energy landscape. In *International conference on machine learning*, pages 1309–1318. PMLR.

Gintare Karolina Dziugaite and Daniel M Roy. 2017. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. *arXiv preprint arXiv:1703.11008*.

Rahim Entezari, Hanie Sedghi, Olga Saukh, and Behnam Neyshabur. 2021. The role of permutation invariance in linear mode connectivity of neural networks. *arXiv preprint arXiv:2110.06296*.

- Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. 2021. [A survey of data augmentation approaches for NLP](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 968–988, Online. Association for Computational Linguistics.
- Stanislav Fort, Gintare Karolina Dziugaite, Mansheej Paul, Sepideh Kharaghani, Daniel M Roy, and Surya Ganguli. 2020. Deep learning versus kernel learning: an empirical study of loss landscape geometry and the time evolution of the neural tangent kernel. *Advances in Neural Information Processing Systems*, 33:5850–5861.
- Stanislav Fort, Huiyi Hu, and Balaji Lakshminarayanan. 2019. Deep ensembles: A loss landscape perspective. *arXiv preprint arXiv:1912.02757*.
- Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. 2020. Linear mode connectivity and the lottery ticket hypothesis. In *International Conference on Machine Learning*, pages 3259–3269. PMLR.
- Timur Garipov, Pavel Izmailov, Dmitrii Podoprikin, Dmitry P Vetrov, and Andrew G Wilson. 2018. Loss surfaces, mode connectivity, and fast ensembling of dnns. *Advances in neural information processing systems*, 31.
- Charles Goddard, Shamane Siriwardhana, Malikeh Ehghaghi, Luke Meyers, Vladimir Karpukhin, Brian Benedict, Mark McQuade, and Jacob Solawetz. 2024. [Arcee’s MergeKit: A toolkit for merging large language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 477–485, Miami, Florida, US. Association for Computational Linguistics.
- Jianping Gou, Baosheng Yu, Stephen J. Maybank, and Dacheng Tao. 2021. [Knowledge distillation: A survey](#). *International Journal of Computer Vision*, 129(6):1789–1819.
- Han Guo, Ramakanth Pasunuru, and Mohit Bansal. 2019. [AutoSeM: Automatic task selection and mixing in multi-task learning](#). 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 3520–3531, Minneapolis, Minnesota. Association for Computational Linguistics.
- M. R. Henry, A. Y. Shi, and et al. 2019. Memory-efficient batching for neural networks. In *NeurIPS*.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. [Distilling the knowledge in a neural network](#).
- Sepp Hochreiter and Jürgen Schmidhuber. 1994. Simplifying neural nets by discovering flat minima. *Advances in neural information processing systems*, 7.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Flat minima. *Neural computation*, 9(1):1–42.
- Chengming Hu, Xuan Li, Dan Liu, Haolun Wu, Xi Chen, Ju Wang, and Xue Liu. 2023. [Teacher-student architecture for knowledge distillation: A survey](#).
- Yukun Huang, Yanda Chen, Zhou Yu, and Kathleen McKeown. 2022. [In-context learning distillation: Transferring few-shot learning ability of pre-trained language models](#).
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2022a. Editing models with task arithmetic. *arXiv preprint arXiv:2212.04089*.
- Gabriel Ilharco, Mitchell Wortsman, Samir Yitzhak Gadre, Shuran Song, Hannaneh Hajishirzi, Simon Kornblith, Ali Farhadi, and Ludwig Schmidt. 2022b. Patching open-vocabulary models by interpolating weights. *Advances in Neural Information Processing Systems*, 35:29262–29277.
- Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. 2018. Averaging weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*.
- Jiaming Ji, Donghai Hong, Borong Zhang, Boyuan Chen, Josef Dai, Boren Zheng, Tianyi Qiu, Boxun Li, and Yaodong Yang. 2024. [Pku-saferlhf: Towards multi-level safety alignment for llms with human preference](#).
- Li Jiang, Yusen Wu, Junwu Xiong, Jingqing Ruan, Yichuan Ding, Qingpei Guo, Zujie Wen, Jun Zhou, and Xiaotie Deng. 2024. [Hummer: Towards limited competitive preference dataset](#).
- 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](#).
- Jean Kaddour, Linqing Liu, Ricardo Silva, and Matt J Kusner. 2022. When do flat minima optimizers work? *Advances in Neural Information Processing Systems*, 35:16577–16595.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. 2016. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*.
- Rohith Kuditipudi, Xiang Wang, Holden Lee, Yi Zhang, Zhiyuan Li, Wei Hu, Rong Ge, and Sanjeev Arora. 2019. Explaining landscape connectivity of low-cost solutions for multilayer nets. *Advances in neural information processing systems*, 32.

- Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. 2018. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems*, 31.
- Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. 2024. NuminaMath. [<https://huggingface.co/AI-M0/NuminaMath-1.5>]([https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina\\_dataset.pdf](https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf)).
- Margaret Li, Suchin Gururangan, Tim Dettmers, Mike Lewis, Tim Althoff, Noah A Smith, and Luke Zettlemoyer. 2022a. Branch-train-merge: Embarrassingly parallel training of expert language models. *arXiv preprint arXiv:2208.03306*.
- Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen, Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian, Baolin Peng, Yi Mao, Wenhui Chen, and Xifeng Yan. 2022b. Explanations from large language models make small reasoners better.
- L. Liu and A. Yao. 2019. Meta-learning for low-shot neural architecture search. In *ICLR*.
- Ekdeep Singh Lubana, Eric J Bigelow, Robert P Dick, David Krueger, and Hidenori Tanaka. 2023. Mechanistic mode connectivity. In *International Conference on Machine Learning*, pages 22965–23004. PMLR.
- Michael S Matena and Colin A Raffel. 2022. Merging models with fisher-weighted averaging. *Advances in Neural Information Processing Systems*, 35:17703–17716.
- Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. 2025. [s1: Simple test-time scaling](#).
- Victor Mustar. 2025. [Tinyr1-32b-preview, the new local king?](#) Accessed: 2025-03-05.
- Behnam Neyshabur, Hanie Sedghi, and Chiyuan Zhang. 2020. What is being transferred in transfer learning? *Advances in neural information processing systems*, 33:512–523.
- A. Radford and J. W. 2021. Improving generalization via scalable learning and domain adaptation. In *NeurIPS*.
- Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. 2015. [Fitnets: Hints for thin deep nets](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca).
- K. Tay, H. Dehghani, L. R. Jones, S. J. Koyejo, and A. S. S. Ahmed. 2020. Efficient transformers: A survey. *ACM Computing Surveys (CSUR)*, 53(4):1–38.
- Open R1 Team. 2025a. Open r1. <https://huggingface.co/open-r1>.
- OpenThoughts Team. 2025b. Open Thoughts. <https://open-thoughts.ai>.
- Fanqi Wan, Longguang Zhong, Ziyi Yang, Ruijun Chen, and Xiaojun Quan. 2024. Fusechat: Knowledge fusion of chat models. *arXiv preprint arXiv:2408.07990*.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khoshdel, and Hannaneh Hajishirzi. 2023. [Self-instruct: Aligning language models with self-generated instructions](#).
- Mitchell Wortsman, Maxwell C Horton, Carlos Guestrin, Ali Farhadi, and Mohammad Rastegari. 2021. Learning neural network subspaces. In *International Conference on Machine Learning*, pages 11217–11227. PMLR.
- Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. 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.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. [Wizardlm: Empowering large language models to follow complex instructions](#).
- Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold Cheng, Jinyang Li, Can Xu, Dacheng Tao, and Tianyi Zhou. 2024. [A survey on knowledge distillation of large language models](#).
- Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. 2023. Ties-merging: Resolving interference when merging models. *Advances in Neural Information Processing Systems*, 36:7093–7115.
- Chuanpeng Yang, Wang Lu, Yao Zhu, Yidong Wang, Qian Chen, Chenlong Gao, Bingjie Yan, and Yiqiang Chen. 2024. [Survey on knowledge distillation for large language models: Methods, evaluation, and application](#).
- Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. 2020. [Gradient surgery for multi-task learning](#).



Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. 2024. [Self-rewarding language models](#).

L. A. Zamir and E. Arbeláez. 2018. Taskonomy: Disentangling task transfer learning. In *CVPR*, pages 2647–2656.

Dan Zhang, Ziniu Hu, Sining Zhoubian, Zhengxiao Du, Kaiyu Yang, Zihan Wang, Yisong Yue, Yuxiao Dong, and Jie Tang. 2024. [Sciinstruct: a self-reflective instruction annotated dataset for training scientific language models](#).

Haosheng Zou, Xiaowei Lv, Shousheng Jia, and Xi-angzheng Zhang. 2024. [360-llama-factory](#).