
A Model Merging Method

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 At the NeurIPS 2024 LLM-Merging competition, we successfully developed a
2 simple and effective model merging approach that generates a versatile, generalist
3 model, applicable to a wide range of scenarios. Specifically, this method is easy to
4 implement, prevents significant conflicts among component models, and improves
5 overall model accuracy to some extent. Additionally, it is memory-efficient, al-
6 lowing for the combination of multiple foundational models, further optimizing
7 resource utilization.

8 1 Introduction

9 Current methods for model merging primarily include Parameter Averaging, MoE-based merging,
10 Model Stacking, Model Zipping, and Model Routing [1]. The NeurIPS 2024 LLM-Merging competi-
11 tion aims to create functionally comprehensive models by merging various domain-specific models
12 efficiently. Our approach utilizes a combination of model routing and model stacking, which has
13 yielded significant accuracy improvements. Furthermore, this method is highly adaptable and can be
14 effectively extended to multiple domains, offering both scalability and flexibility for a wide range of
15 applications.

16 2 Method

17 Our approach is primarily reflected in three aspects: model selection, model merging, and staged
18 responses.

19 2.1 Model Selection

20 For model selection, we approached it from two directions: base models and fine-tuned models. First,
21 we identified a base model by initially screening five options that met the competition’s criteria. Using
22 the capabilities of large models listed on Hugging Face’s Open LLM Leaderboard2, we narrowed
23 it down to three high-accuracy base models: Meta-Llama-3-8B-Instruct, Phi-3-mini-4k-instruct,
24 and Phi-3-small-8k-instruct. We used the lm-evaluation-harness tool to test each base model’s
25 performance on several public datasets. Results showed that Meta-Llama-3-8B-Instruct and Phi-3-
26 small-8k-instruct had their own advantages, displaying similar patterns across our internal test sets.
27 Consequently, we selected Meta-Llama-3-8B-Instruct and Phi-3-small-8k-instruct as our final base
28 models.

29 For fine-tuned models, we divided task types by knowledge area, including history, medicine,
30 mathematics, physics, and coding. Using datasets integrated within lm-evaluation-harness, combined
31 with our proprietary datasets, we assessed each fine-tuned model’s GPU memory usage and accuracy
32 gains over the base models. Ultimately, we selected medical and coding models as our final fine-tuned
33 models, balancing memory efficiency with performance enhancements.

34 2.2 Model Merging

35 We observed that using parameter averaging method to merge models often reduces performance
36 on tasks where the base models excel. Therefore, we chose a model combination approach similar
37 to model routing, which typically requires substantial memory. To address this, we merged models
38 with identical architectures by compressing differences, preserving each base model’s task-specific
39 strengths while minimizing inference overhead. In our trials, this approach only increased inference
40 costs by around 10%, maintaining efficiency without sacrificing performance.

41 2.2.1 Weights Merging

42 In our solution, we combine several fully fine-tuned Llama3 8B models. We assume that the weights
43 of a base model and the fine-tuned models can each be represented as a common component plus a
44 differentiated component. The common component can be obtained through parameter averaging,
45 while the unique components are compressed separately. Besides retaining their own embeddings,
46 lm_head layer weights, and biases (if any) for each model, other layers are merged and compressed.
47 For faster compression, we employed randomized SVD (rSVD) instead of standard SVD. When
48 the compression ratio is set to 0, the result is equivalent to using multiple independent models; at
49 100%, it corresponds to parameter averaging. Interestingly, after this combined compression, each
50 branch model showed an improvement in general task performance, though due to time constraints,
51 we didn’t delve into this further.

52 The pseudocode for weight merging and compression for a specific layer is as follows:

```
53 inputs :  
54     weights ,  
55     compress_rate ,  
56 outputs :  
57     scales ,  
58     # common component weights  
59     w_avg ,  
60     # differentiated component weight , including U and V.  
61     compressed_diff ,  
62 scales = [norm(w) for w in weights]  
63 normalized_weights = [  
64     w / scale for w, scale in zip(weights , scales)]  
65 w_avg = mean(stack(normalized_weights , dim=0), dim=0)  
66 compressed_diff = [  
67     rsvd(w - w_avg, compress_rate) for w in normalized_weights]  
68 return scales , w_avg, compressed_diff
```

69 The pseudocode for inference in a specific layer of a branch model is as follows:

```
70 inputs :  
71     x ,  
72     bias , # uncompressed bias  
73     scale ,  
74     w_avg , # common component weights  
75     compressed_diff , # differentiated component weight  
76 outputs: y,  
77 y = linear(x, w_avg) + linear(linear(x, V), U) * scale  
78 if bias:  
79     y += bias  
80 return y
```

81 In our approach, we set the compression rate to 95%. We loaded one independent Phi3small model
82 and three fully fine-tuned Llama3 8B models, compressed using this method, into 48GB of VRAM.
83 In practical applications, this method allows for compressing and combining more models as needed.

84 2.2.2 Router

85 In our approach, we utilized four branches: one standalone Phi-3-small model and three Llama-3-8B
86 models merged through compression. To determine the appropriate branch for answering each input,
87 we route the question based on relevance. To avoid potential generalization issues and minimize
88 processing time, we did not train a dedicated router. Instead, we designed a set of representative
89 samples for each branch. Therefore, we designed 17 samples that were carefully crafted to cover
90 each branch’s focus area and do not overlap with the competition data.

91 We utilized a decoder-only model and carefully designed our routing method to maximize the router’s
92 ability to distinguish between questions from different domains. For each sample or input question,
93 we extract the embedding from the segment "{input}.\n Let’s think about what task these questions
94 belong to. These questions belong to the field of" specifically from the phrase "These questions belong
95 to the field of". Due to the attention mechanism, this embedding effectively captures information from
96 the preceding input. Since both input questions and samples rely on this structure, they achieve high
97 alignment. In our experiments, this method significantly improved the router’s ability to differentiate
98 between domains compared to extracting the embedding from the input alone.

99 For classification, we calculate the cosine similarity between input questions and the sample embed-
100 dings, with pre-extracted embeddings for each sample to reduce inference time during evaluation.

101 2.3 Staged Responses

102 In the competition, we found that achieving high-quality responses requires two key factors: accuracy
103 and clarity. While Llama-3’s responses were generally more structured compared to Phi-3-small,
104 Llama-3’s accuracy was often weaker. For Phi-3-small (instruct), although prompts can be adjusted
105 to encourage more concise responses, using system-level prompts to enforce brevity sometimes
106 interfered with the model’s reasoning, leading to incorrect answers. Conversely, if we aimed for
107 comprehensive reasoning in response to complex questions, the answers tended to be overly detailed
108 (e.g., an SQL question might include additional, unwanted explanations). Moreover, overly complex
109 prompts sometimes unintentionally reduced overall accuracy. This led us to adopt a two-stage
110 response approach.

111 Unfortunately, due to a time zone miscalculation of the competition deadline, our two-staged solution
112 was only partially submitted, with the full version missing from official scoring. However, in our
113 experiments, the complete two-stage method showed a 6-point improvement in accuracy, compared
114 to only a 2-point gain from early-stage submissions. This section describes the full version of the
115 two-stage approach as a technical demonstration.

116 In our method, tasks that leveraged Phi-3-small’s strengths, like reasoning, were initially answered by
117 Phi-3-small. Llama-3 then refined these answers for clarity and structure through prompt engineering.
118 For some scenarios, we introduced an alternative two-stage response format: Stage 1 generates a
119 guided, heuristic-based answer, while Stage 2 formats the output. Since Stage 2 avoids triggering
120 chain-of-thought processes, it generally produces concise answers without significantly increasing
121 inference time. Notably, Stages 1 and 2 can be handled by the same or different models, depending
122 on the task.

123 References

124 [1] Tam D. & Li M. & Yadav P. & et al. (2024) Llm merging: Building llms efficiently through merging. *NeurIPS*
125 *2024 Competition Track*.