



## Evolutionary Optimization of Model Merging Recipes

Cirò G., Lomele M.

30592 - Topics in Computational Modelling: from Information Theory to  
Evolutionary Models

April 14, 2024

---

# Evolutionary Optimization of Model Merging Recipes

---

Takuya Akiba, Makoto Shing, Yujin Tang, Qi Sun, David Ha

Sakana AI

Tokyo, Japan

{takiba,mkshing,yujintang,qisun,hadavid}@sakana.ai

## Abstract

We present a novel application of evolutionary algorithms to automate the creation of powerful foundation models. While model merging has emerged as a promising approach for LLM development due to its cost-effectiveness, it currently relies on human intuition and domain knowledge, limiting its potential. Here, we propose an **evolutionary approach** that overcomes this limitation by **automatically discovering effective combinations of diverse open-source models**, harnessing their collective intelligence without requiring extensive additional training data or compute. Our approach operates in both parameter space and data flow space, allowing for optimization beyond just the weights of the individual models. This approach even facilitates cross-domain merging, generating models like a Japanese LLM with Math reasoning capabilities. Surprisingly, our **Japanese Math LLM** achieved state-of-the-art performance on a variety of established Japanese LLM benchmarks, even surpassing models with significantly more parameters, despite not being explicitly trained for such tasks. Furthermore, a **culturally-aware Japanese VLM** generated through our approach demonstrates its effectiveness in describing Japanese culture-specific content, outperforming previous Japanese VLMs. This work not only contributes new state-of-the-art models back to the open-source community, but also introduces a new paradigm for automated model composition, paving the way for exploring alternative, efficient approaches to foundation model development.<sup>1</sup>

# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

## 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

## 4 Results

Experiments

Natural Language Task

Computer Vision Task





# The Model Merging Problem

Combine **multiple PTMs** fine-tuned for  $T$  different tasks into a **single new model** which performs well on all  $T$  tasks.

$$\{\theta_{FT}^1, \dots, \theta_{FT}^T\} \rightarrow \theta_{new}$$

How?

- Simple Average
- RegMean
- Fisher Merging
- Task Arithmetic
- TIES-Merging
- ...
- **Evolutionary Model Merge**

# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

### 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

### 4 Results

Experiments

Natural Language Task

Computer Vision Task

# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

### 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

### 4 Results

Experiments

Natural Language Task

Computer Vision Task

# Task Arithmetic

Ilharco et al. (Dec 2022) introduced **task vectors**  $\tau$  as the difference in model's parameters before and after **fine-tuning**:

$$\tau^t = \theta_{FT}^t - \theta_{PTM}$$

Combining multiple task vectors with basic **vector operations** yields a well-performing **multi-task** model:

$$\theta_{new} = \theta_{PTM} + \lambda \sum_{t=1}^T \tau^t$$

where  $\lambda$  is a scaling hyper-parameter.

# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

## 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

## 4 Results

Experiments

Natural Language Task

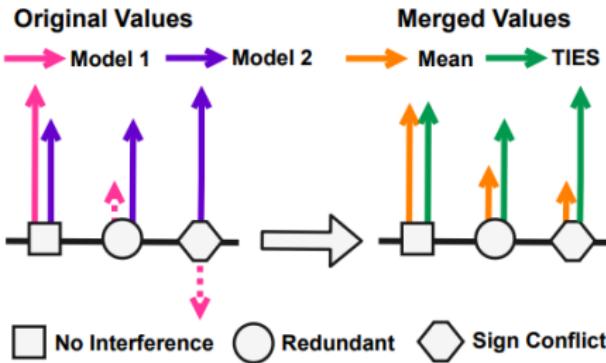
Computer Vision Task

TIES-Merging

# TIES-Merging

Yadav et al. (Jun 2023) showed there is **interference** in delta parameters when **combining task vectors**, due to:

- Redundant Parameters Value
- Sign Conflicts



TIES-Merging

# TIES-Merging

They propose Trim, Elect Sign and Disjoint Merge (**TIES**):

- ① **Trim** → keep only top- $k\%$  delta parameters by magnitude;
- ② **Elect Sign** → by highest total magnitude;

$$\text{sign}^* = [\text{sign}_1^*, \dots, \text{sign}_P^*]$$

$$\text{sign}_p^* = \text{sign} \left( \sum_{t=1}^T \tau_p^t \right), p \in \{1, \dots, P\}$$

- ③ **Disjoint Merge** → compute mean of delta parameters whose signs agree with the elected sign.

$$\tau_p^m = \frac{1}{|A^p|} \sum_{t \in A^p} \tau_p^t$$

where  $A^p = \{t \in [T] : \text{sign}_p^t = \text{sign}_p^*\}$

# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

## 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

## 4 Results

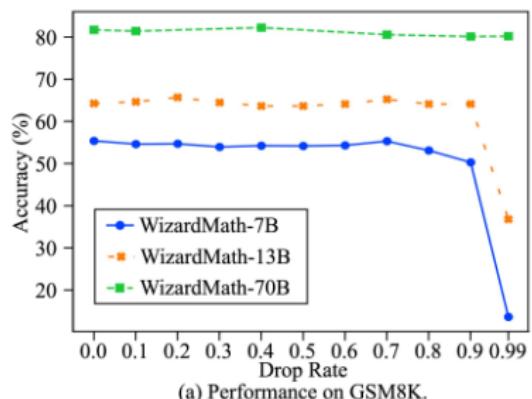
Experiments

Natural Language Task

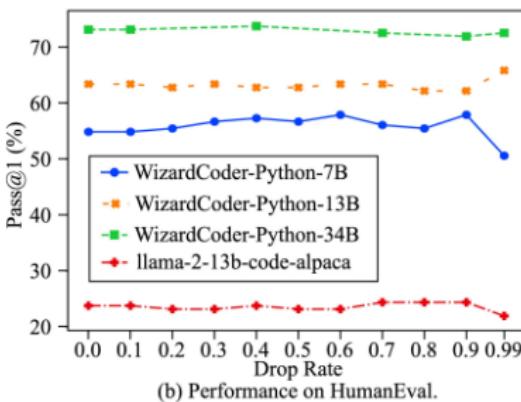
Computer Vision Task

# DARE

Yu et al. (Nov 2023) showed that most (**90-99%**) of the change in parameters happening during **fine-tuning** is actually **redundant**.



(a) Performance on GSM8K.



(b) Performance on HumanEval.

The **same performance** can be achieved by **resetting** most **FT parameters** to their original PTM values and **scaling** the rest accordingly.

# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

### 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

### 4 Results

Experiments

Natural Language Task

Computer Vision Task

# Frankenmerging

So far, we moved in the **Parameter Space (PS)** by interpolating weights of multiple models, keeping the architecture fixed.

We can also move in the **Data Flow Space (DFS)** by changing how the input flows through the model, i.e. changing the architecture.

*Frankenmerging* → stack different layers from different models together to create new architectures, while keeping the weights fixed.

# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

### 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

### 4 Results

Experiments

Natural Language Task

Computer Vision Task

# Evolutionary Algorithms

There are 4 main dialects of Evolutionary Algorithms (EA), whose main difference lies in the data structure used to encode candidate solutions:

- Genetic Algorithms (GA) → string over a finite alphabet
- **Evolution Strategies (ES)** → **real-valued vectors**
- Evolutionary Programming (EP) → finite state machines
- Genetic Programming (GP) → trees

# Evolution Strategies

Given a starting population of size  $\lambda$  of candidate solutions, apply selection, recombination and mutation to improve fitness, minimizing an objective function.

- Candidate Solution:  $\mathbf{x}_i \in \mathbb{R}^n$
- Population:  $\{\mathbf{x}_1, \dots, \mathbf{x}_\lambda\}$
- Objective Function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

We will focus on **selection** and **mutation**.

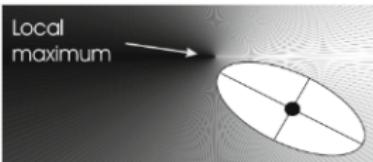
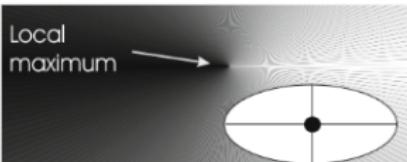
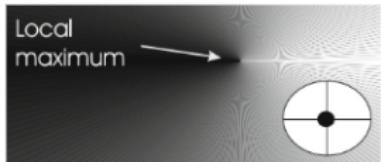
# Selection and Mutation of Real-valued Vectors

At each generation  $g$ , select top  $\mu$  candidate solutions according to fitness (**exploitation**) and apply mutation by adding Gaussian noise (**exploration**).

$$\mathbf{x}_i^{g+1} = \mathbf{x}_i^g + \mathbf{z}_i$$

where  $\mathbf{z}_i \sim \mathcal{N}(0, C_i) \in \mathbb{R}^n$ .

Modify the covariance matrix  $C_i$  to explore the search space in different ways:

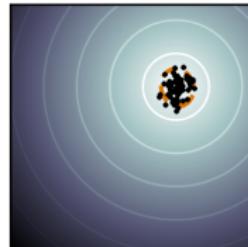
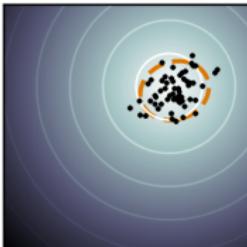
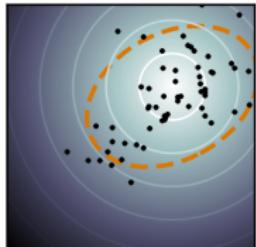


# Estimation of Distributions Algorithms

Traditional Evolution Strategies (**ESs**) evolve an **implicit distribution** over the solution space by mutating each individual.

Estimation of Distribution Algorithms (**EDAs**) evolve an **explicit distribution** by mutating the parameters of a multivariate search distribution.

ESs maintain the entire population, EDAs preserve only the population statistics.



# CMA-ES

Hansen (Apr 2016) introduced **Covariance Matrix Adaptation (CMA) Evolution Strategy (ES)**, a novel approach to explore the solution space by evolving the statistics of a Search Distribution.

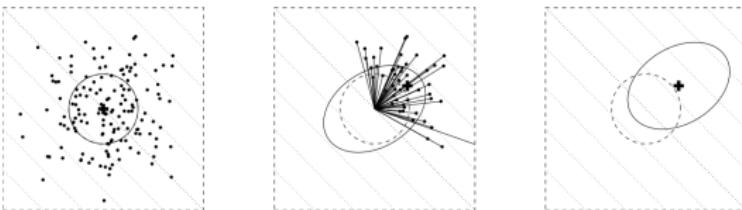
# CMA-ES

At each generation  $g = 1, 2, \dots$  sample candidate solutions:

$$\mathbf{x}_k^{(g+1)} \sim \mathbf{m}^{(g)} + \sigma^{(g)} \mathcal{N}(0, \mathbf{C}^{(g)})$$

for  $k = 1, \dots, \lambda$ ,  $\mathbf{m}^{(g)} \in \mathbb{R}^n$ ,  $\mathbf{C}^{(g)} \in \mathbb{R}^{n \times n}$

Apply selection and adapt parameters (**mean  $m$ , step-size  $\sigma$ , covariance matrix  $C$** ):



# Adapting the Mean

The new mean  $\mathbf{m}^{g+1}$  is a weighted average of  $\mu$  selected points from the sample  $\mathbf{x}_1^{g+1}, \dots, \mathbf{x}_{\lambda}^{g+1}$ :

$$\mathbf{m}^{(g+1)} = \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}^{(g+1)}$$

$$\sum_{i=1}^{\mu} w_i = 1 \quad \text{for } w_1 \geq w_2 \geq \dots \geq w_{\mu} > 0.$$

$\mathbf{x}_{i:\lambda}$  is the  $i$ -th fittest solution.

# Adapting the Covariance Matrix

The new covariance matrix  $C^{(g+1)}$  is obtained as follows:

$$C^{(g+1)} = \underbrace{(1 - c_1 - c_\mu \sum w_j)}_{\text{can be close or equal to 0}} C^{(g)} + c_1 \underbrace{\mathbf{p}_c^{(g+1)} \mathbf{p}_c^{(g+1)\top}}_{\text{rank-one update}} + c_\mu \underbrace{\sum_{i=1}^{\lambda} w_i \mathbf{y}_{i:\lambda}^{(g+1)} (\mathbf{y}_{i:\lambda}^{(g+1)})^\top}_{\text{rank-}\mu\text{ update}}$$

Rank- $\mu$  update

$$\mathbf{C}^{(g+1)} = \underbrace{(1 - c_1 - c_\mu \sum w_j)}_{\substack{\text{can be close or equal to 0}}} \mathbf{C}^{(g)}$$

↓

$$+ c_1 \underbrace{\mathbf{p}_c^{(g+1)} \mathbf{p}_c^{(g+1)^\top}}_{\text{rank-one update}} + c_\mu \underbrace{\sum_{i=1}^{\lambda} w_i \mathbf{y}_{i:\lambda}^{(g+1)} (\mathbf{y}_{i:\lambda}^{(g+1)})^\top}_{\text{rank-}\mu\text{ update}}$$

where

$$\mathbf{y}_{i:\lambda}^{(g+1)} = (\mathbf{x}_{i:\lambda}^{(g+1)} - \mathbf{m}^{(g)}) / \sigma^{(g)}$$

and

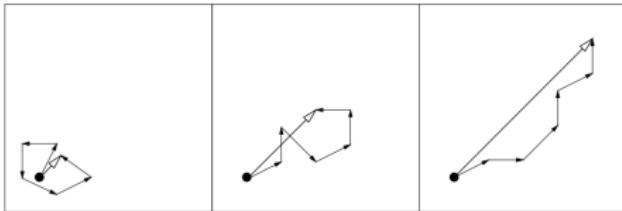
$$w_{1,\dots,\lambda} \in \mathbb{R} : w_1 \geq \dots \geq w_\mu \geq 0 \geq w_{\mu+1} \geq \dots \geq w_\lambda$$

The rank- $\mu$  update uses the standardized and weighted sample covariance matrix.

# Evolution Path

An **evolution path** is the path traced across the search space by the Search Distribution's **mean** as it **evolves** throughout generations:

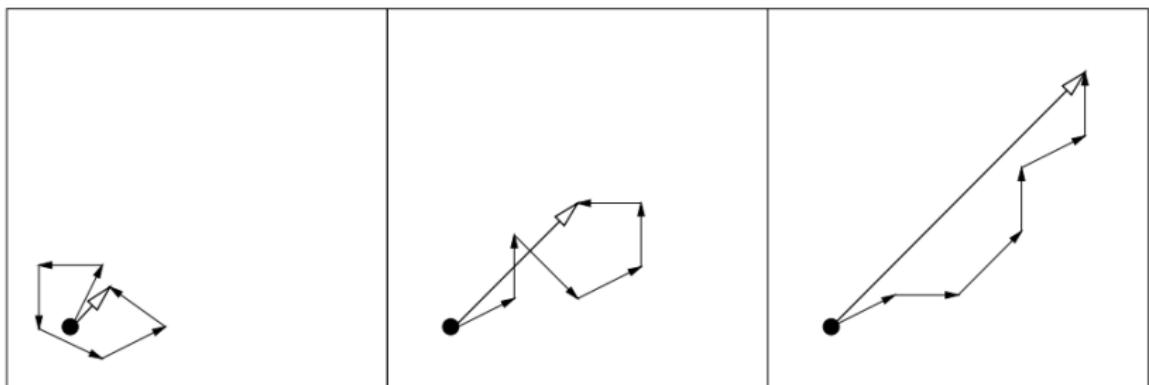
$$\frac{\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}}{\sigma^{(g)}} + \frac{\mathbf{m}^{(g)} - \mathbf{m}^{(g-1)}}{\sigma^{(g-1)}} + \frac{\mathbf{m}^{(g-1)} - \mathbf{m}^{(g-2)}}{\sigma^{(g-2)}} + \dots + \frac{\mathbf{m}^{(1)} - \mathbf{m}^{(0)}}{\sigma^{(0)}}$$





# Adapting the Step-Size

The step-size  $\sigma^{(g+1)}$  is updated based on evolution paths:



- Anti-correlated (left) → decrease;
- Orthogonal (middle) → keep;
- Correlated (right) → increase.

# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

## 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

## 4 Results

Experiments

Natural Language Task

Computer Vision Task

# Overview

The Open LLM Leaderboard is dominated by merged models, great cost-effectiveness.

However, current model merging relies on **human intuition** and domain knowledge.

Akiba et al. (Mar 2024) from Sakana AI introduced a **systematic** model composition framework using **evolution**. In particular:

- ① Merge in Parameter Space (PS).
- ② Merge in Data Flow Space (DFS).
- ③ Integrate PS and DFS.

# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

## 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

## 4 Results

Experiments

Natural Language Task

Computer Vision Task

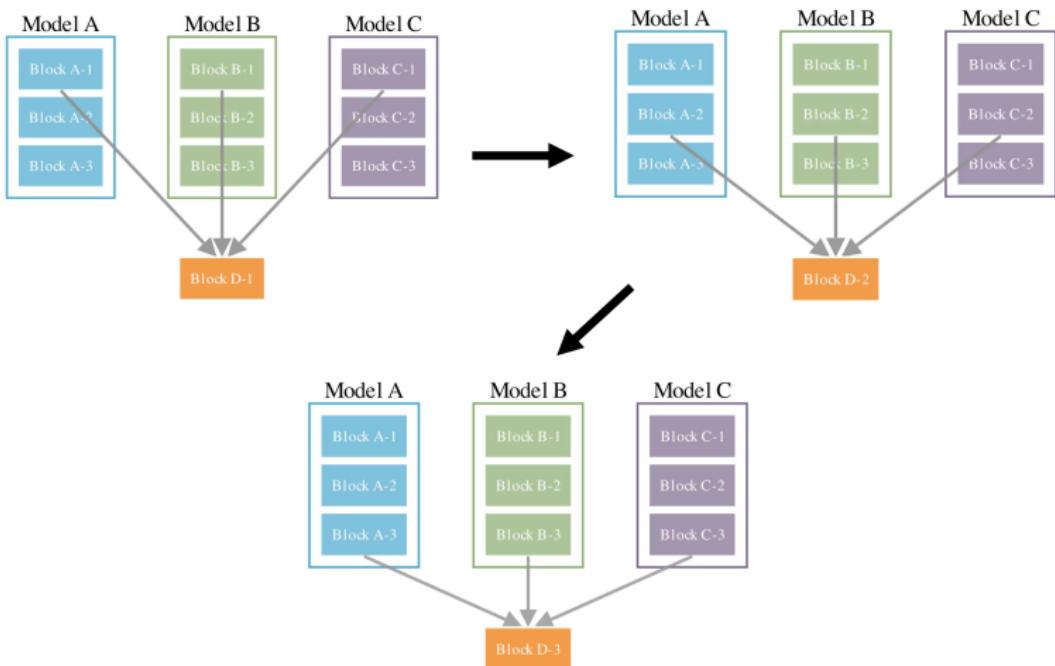
# Merge in PS

*Sub-goal 1: Merge the **weights** of the models **block-wise**, where each block/layer is an input-output embedding module, e.g. a transformer's block.*

- ① Apply DARE and TIES-Merging.
- ② Establish **configuration parameters** for sparsification and weight mixing (e.g. scaling in TIES and drop % in DARE).
- ③ Set accuracy (or any other critical task-specific metric) as the fitness function.
- ④ Optimise the parameters with CMA-ES.

Parameter Space (PS)

# Merge in PS visualized



# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

## 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

## 4 Results

Experiments

Natural Language Task

Computer Vision Task

# Merge in DFS

*Sub-goal 2:* optimise how **tokens flow** across the layers of the models, leaving weights unchanged. Given  $N$  models, with a total of  $M$  layers, and a budget length  $T$ , we search for an **optimal sequence** of layers.

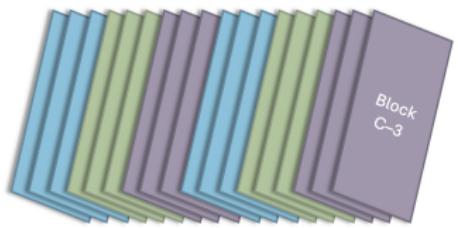
- ① Lay out models in sequential order, repeat them  $r$  times.
- ② Set an **indicator array**  $I \in \mathbb{R}^{T=M \times r}$  to deal with the inclusion ( $I_i = 1$ ) and exclusion ( $I_i = 0$ ) of layers.
- ③ Set a **scaling matrix**  $W \in \mathbb{R}^{M \times M}$ , where  $W_{ij}$  is scaling from layer  $i$  to  $j$  in the sequence. This because a layer may face an input whose distribution is different from what it is used to.
- ④ Optimise  $I$  and  $W$  with CMA-ES.

Data Flow Space (DFS)

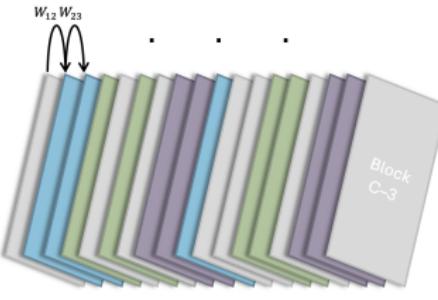
# Merge in DFS visualised

Number of models:  $N = 3$ Number of total layers:  $M = 9$ Number of repetitions:  $r = 2$ Budget for sequence length:  $T = 18$ 

①



② &amp; ③

 $I = [0, 1, 1, \dots, 0]$

Interpolation

# Table of Contents

## ① The Model Merging Problem

### ② Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

## ③ Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

## ④ Results

Experiments

Natural Language Task

Computer Vision Task

# Integrate PS and DFS

**Combine** the two approaches:

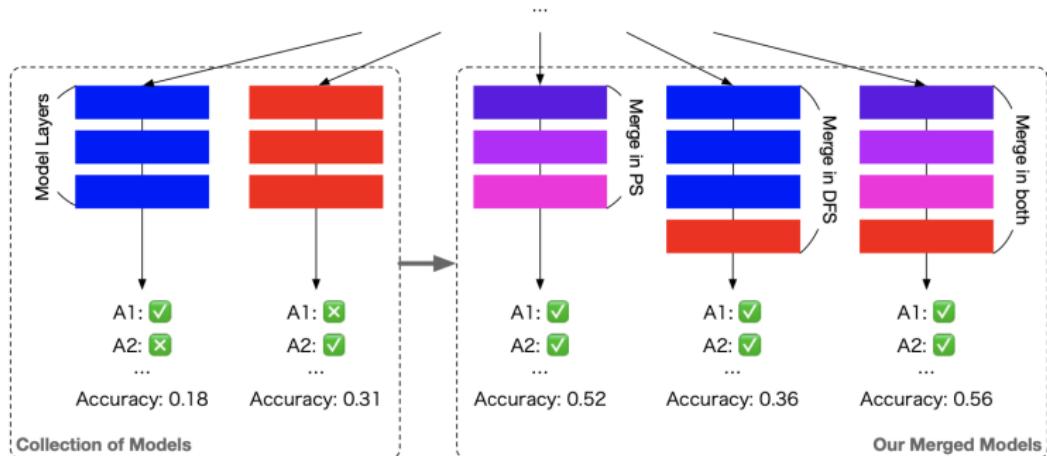
- ① Merge in PS a collection of models  
→ new merged model.
- ② Insert new merged model into the collection  
→ updated collection.
- ③ Merge in DFS the updated collection  
→ final merged model.

Interpolation

# Evolutionary Model Merge visualised

Q1: Mishka bought 3 pairs of shorts, 3 pairs of long pants, and 3 pairs of shoes. ... How much were spent on all the clothing?

Q2: Cynthia eats one serving of ice cream every night. ... How much will she have spent on ice cream after 60 days?



# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

## 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

## 4 Results

Experiments

Natural Language Task

Computer Vision Task

# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

## 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

## 4 Results

Experiments

Natural Language Task

Computer Vision Task

# Tasks

- **Natural Language Processing:** Solve math problems in Japanese.
- **Computer Vision:** Recognise and describe elements of the Japanese culture in pictures.

# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

## 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

## 4 Results

Experiments

Natural Language Task

Computer Vision Task

# Natural Language Task

*Goal:* combine **Japanese LLMs** and **Math LLMs** (all based on Mistral-7B) into one that can generate **Japanese answers to Japanese math problems**.

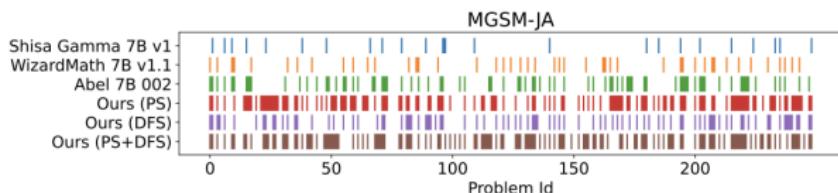
*Dataset:* MGSM. *Train set:* 1069 samples translated in Japanese. *Test set:* 250 samples original in Japanese.

Optimisation with CMA-ES:

- *PS Merging*
  - $\sigma = 1/6$
  - population size  $4 + \lfloor 3\ln(n_{params}) \rfloor$
- *DFS Merging*
  - $M = 64, r = 3$
  - population size = 128
  - generations = 100

# Results

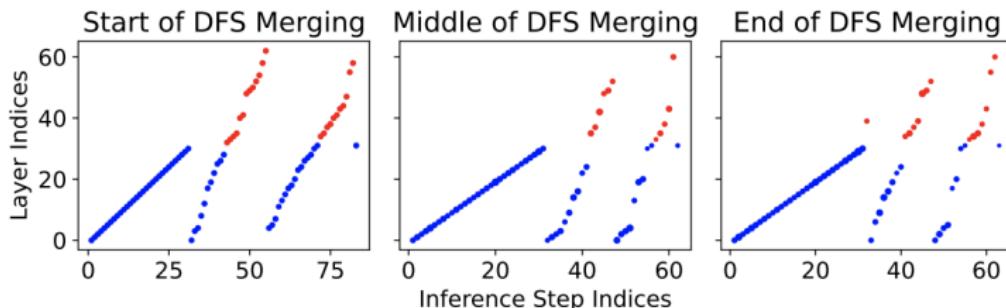
*Key finding:* merged model retains foundational knowledge from the source models and exhibits the desired **emerging capability**.



<b>Id.</b>	<b>Model</b>	<b>Type</b>	<b>Size</b>	<b>MGSM-JA (acc ↑)</b>	<b>JP-LMEH (avg ↑)</b>
1	Shisa Gamma 7B v1	JA general	7B	9.6	66.1
2	WizardMath 7B v1.1	EN math	7B	18.4	60.1
3	Abel 7B 002	EN math	7B	30.0	56.5
4	<b>Ours (PS)</b>	1 + 2 + 3	7B	<b>52.0</b>	<b>70.5</b>
5	<b>Ours (DFS)</b>	3 + 1	10B	<b>36.4</b>	<b>53.2</b>
6	<b>Ours (PS+DFS)</b>	4 + 1	10B	<b>55.2</b>	<b>66.2</b>
7	Llama 2 70B	EN general	70B	18.0	64.5
8	Japanese StableLM 70B	JA general	70B	17.2	68.3
9	Swallow 70B	JA general	70B	13.6	71.5
10	GPT-3.5	commercial	-	50.4	-
11	GPT-4	commercial	-	78.8	-

# Results

A look at **model explainability** in DFS merging with two models.



The evolutionary search includes most layers of the first model, the blue dots, at an early stage. Then, it alternates between layers from both models.

# Table of Contents

## 1 The Model Merging Problem

### 2 Useful Background

Task Arithmetic

TIES-Merging

DARE

Frankenmerging

CMA-ES

## 3 Evolutionary Model Merge

Parameter Space (PS)

Data Flow Space (DFS)

Interpolation

## 4 Results

Experiments

Natural Language Task

Computer Vision Task

# Computer Vision Task

*Goal:* design a **culturally-specific content aware** Japanese VLM that can describe pictures and give insights about Japan and its culture.

*Dataset:* Japanese Visual Genome VQA and JA-VLM-Bench-in-the-Wild. The latter is a hand-crafted collection of 42 images with 50 questions and nuanced answers about Japan's culture and its objects.

Optimisation using CMA-ES with identical settings to the previous PS-Merging experiment.

# Results

*Key finding: **multi-modality** allows merged model to produce more detailed and accurate responses.*

<b>Model</b>	<b>JA-VG-VQA-500 (ROUGE-L ↑)</b>	<b>JA-VLM-Bench-In-the-Wild (ROUGE-L ↑)</b>
LLaVA 1.6 Mistral 7B	14.3	41.1
Japanese Stable VLM	-	40.5
<b>Ours</b>	<b>19.7</b>	<b>51.2</b>

# Conclusion

- Novel method for merging models of vastly different domains that achieves state-of-the-art results;
- Showcase for the potential of democratising model development.

# Limitations

- *Inheritance* → no method of filtering selection of models or blocks of models;
- *Incoherence* → no control on logical and factual coherence of responses;
- *Alignment* → lack of instruction fine-tuning;
- *Evolutionary Details* → it's hard to understand exactly which parameters they're evolving and no insights on the evolution process (only a hint for DFS).

# Future Work

- Cross-domain *image generation* → evolutionary model merging on diffusion models;
- *Self-improvement* → evolutionary search for candidate models instead of manual selection of starting population;
- *Generalisation* → merging in more orthogonal spaces (beyond PS & DFS).

# Thank You!



→ Generated by the authors with EvoSDXL-JP, an incoming diffusion model created with Evolutionary Model Merge.

# Question 1

- In the NLP task, how does the final model retain performance on the original tasks?*

From the results table, the performance of the final model (id 6) on general Japanese language tasks is measured on JP-LMEH dataset.

No significant performance drop is observed compared to the original models (1-3), with model 4 (PS) actually improving.

Unfortunately, the authors do not provide insights on the performance of the merged model on English math, suggesting there might be a loss.

## Question 2

- *How does the final model retain performance on English language tasks?*

The authors do not expand on the performance of the merged models on English tasks. This aspect was not part of the original models' tasks, which were fine-tuned for English math problems and Japanese language.

This is also evident from the two databases that were used to test and benchmark the final merged model. The first, MGSM-JA, is a set of grade-school math problems translated to Japanese by the authors. The second, JP-LMEH, is an evaluation harness specific for Japanese language models.

## Question 3

- *How does this method perform on overlapping or less complementary tasks?*

In this paper, the authors do not present additional experiments of using their method. Hence, there is no empirical evidence of how this method behaves for overlapping tasks.

However, in the related research papers (DARE, TIES, Task Arithmetic etc.), it's been shown that merging different models fine-tuned for the same tasks improve both in-domain performance and out-of-domain generalization.

Moreover, merged models provide a better initialization for fine-tuning.

# Bibliography

- Akiba, T., Shing, M., Tang, Y., Sun, Q., & Ha, D. (2023). Evolutionary Optimization of Model Merging Recipes. *arXiv preprint arXiv:2403.13187*. Retrieved from <https://arxiv.org/abs/2403.13187>
- Yu, L., Yu, B., Yu, H., Huang, F., & Li, Y. (2023). Language Models are Super Mario: Absorbing Abilities from Homologous Models as a Free Lunch. *arXiv preprint arXiv:2311.03099*. Retrieved from <https://arxiv.org/abs/2311.03099>
- Ilharco, G., Ribeiro, M.T., Wortsman, M., Gururangan, S., Schmidt, L., Hajishirzi, H., & Farhadi, A. (2022). Editing models with task arithmetic. *arXiv preprint arXiv:2212.04089*. Retrieved from <https://arxiv.org/abs/2212.04089>

# Bibliography

- Hansen, N. (2016). The CMA Evolution Strategy: A Tutorial. *arXiv preprint arXiv:1604.00772*. Retrieved from <https://arxiv.org/abs/1604.00772>
- Sakana AI. (2023). Evolutionary Model Merge. Retrieved from <https://sakana.ai/evolutionary-model-merge/>