# Model Merging

# Model Merging

# Model Merging



$\theta$

Foundation Model

Post-training (fine-tune)

Post-training (fine-tune)

$\theta_A$

$\theta_B$

不用訓練資料!
不用做任何模型訓練!

$(\theta_B - \theta)$

Task vector

接枝王
葛瑞克

(艾爾
登法環)

# 類神經網路參數豈是如此不便之物!



a) Task vectors

$\tau = \theta_{\text{ft}} - \theta_{\text{pre}}$

b) Forgetting via negation

$\tau_{\text{new}} = -\tau$

Example: making a language model produce less toxic content

c) Learning via addition

$\tau_{\text{new}} = \tau_A + \tau_B$

Example: building a multi-task model

d) Task analogies

$\tau_{\text{new}} = \tau_C + (\tau_B - \tau_A)$

Example: improving domain generalization

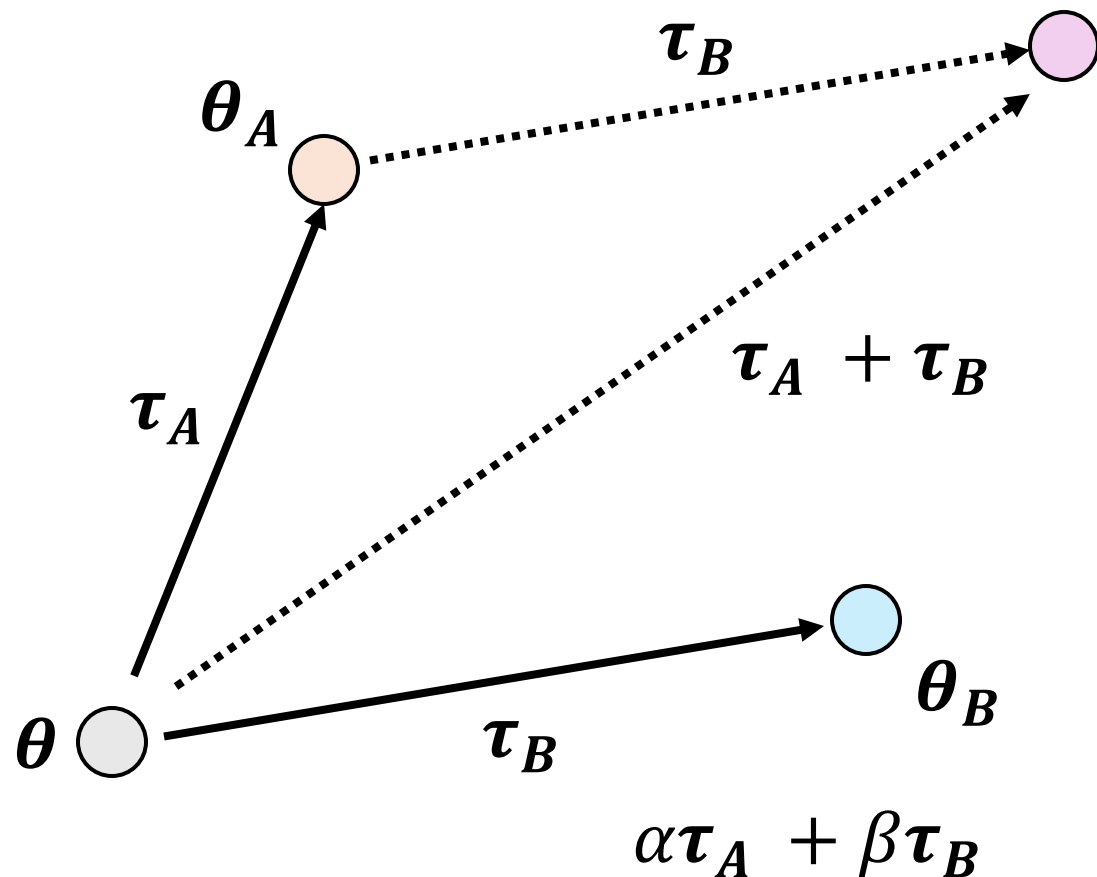https://arxiv.org/abs/2212.04089

# Task Vector has been shown to be helpful.

1. 相加

$$\tau_A = \theta_A - \theta$$

$$\tau_B = \theta_B - \theta$$

$\theta_A, \theta_B$ 來自相同的
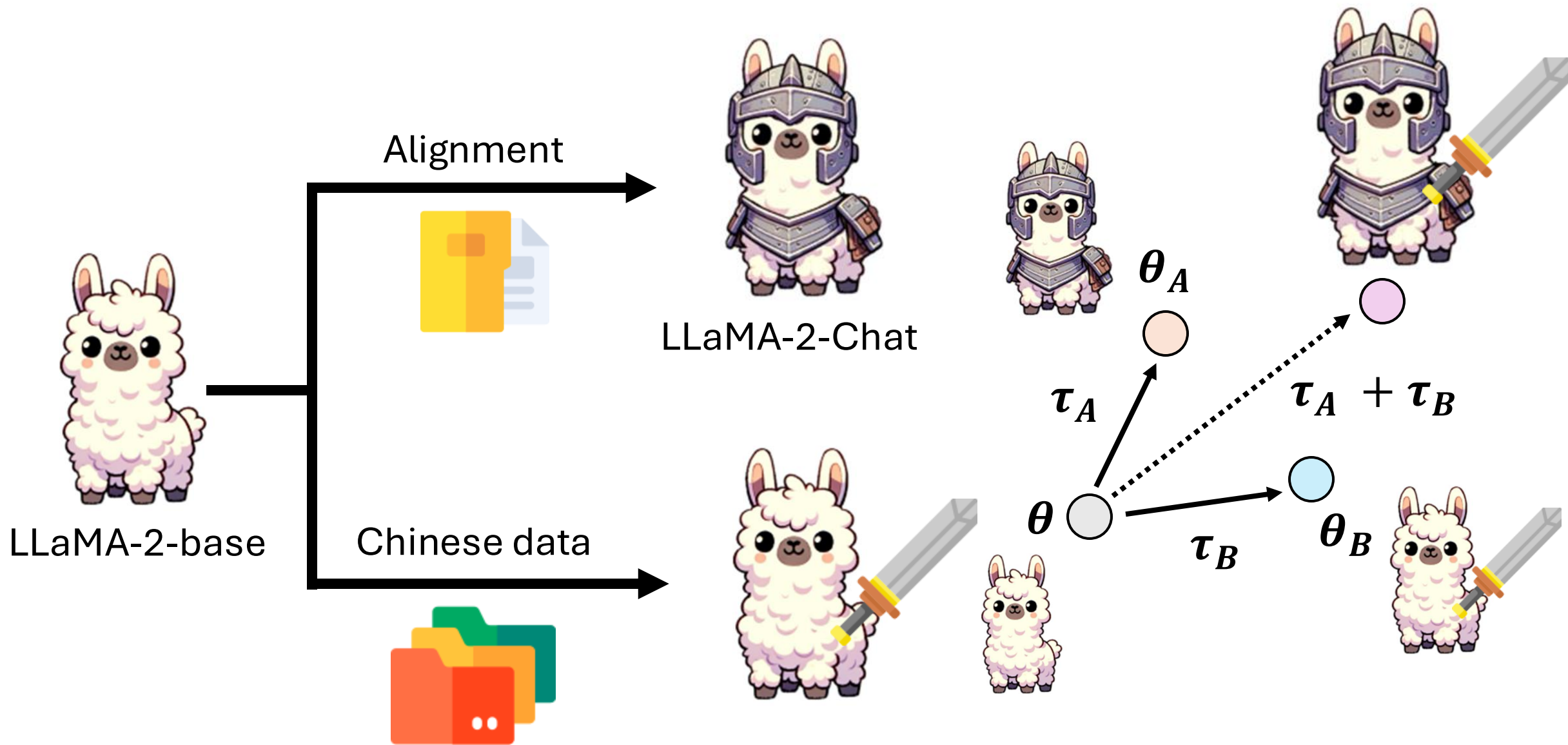Foundation Model $\theta$

**Post training 時代的做法**



$$\tau_A + \tau_B$$

$$\alpha\tau_A + \beta\tau_B$$

LLaMA-2-base

Alignment

LLaMA-2-Chat

Chinese data

Shih-Cheng Huang

Alignment

LLaMA-2-Chat

Chinese data

LLaMA-2-base

$\boldsymbol{\theta}_A$

$\boldsymbol{\tau}_A$

$\boldsymbol{\tau}_A + \boldsymbol{\tau}_B$

$\boldsymbol{\theta}$

$\boldsymbol{\tau}_B$

$\boldsymbol{\theta}_B$

Alignment

LLaMA-2-Chat
LLaMA-3-instruct
Mistral-instruct-0.2

Korean, Japanese
Chinese data

LLaMA-2-base
LLaMA-3-base
Mistral-7B

$\boldsymbol{\theta}_A$

$\boldsymbol{\tau}_A$

$\boldsymbol{\tau}_A + \boldsymbol{\tau}_B$

$\boldsymbol{\theta}$

$\boldsymbol{\tau}_B$

$\boldsymbol{\theta}_B$

https://arxiv.org/abs/2310.04799
https://qiita.com/jovyan/items/ee6affa5ee5bdaada6b4

Reward Model

$\boldsymbol{\theta}_A$

$\boldsymbol{\tau}_A + \boldsymbol{\tau}_B$

$\boldsymbol{\theta}$

$\boldsymbol{\tau}_B$

$\boldsymbol{\theta}_B$

Tzu-Han Lin, Chen-An Li
https://arxiv.org/abs/2407.01470

Reward Model

$\boldsymbol{\theta}_A$

$\boldsymbol{\tau}_A + \boldsymbol{\tau}_B$

$\boldsymbol{\theta}$

$\boldsymbol{\tau}_B$

$\boldsymbol{\theta}_B$

Chen-An Li, Tzu-Han Lin
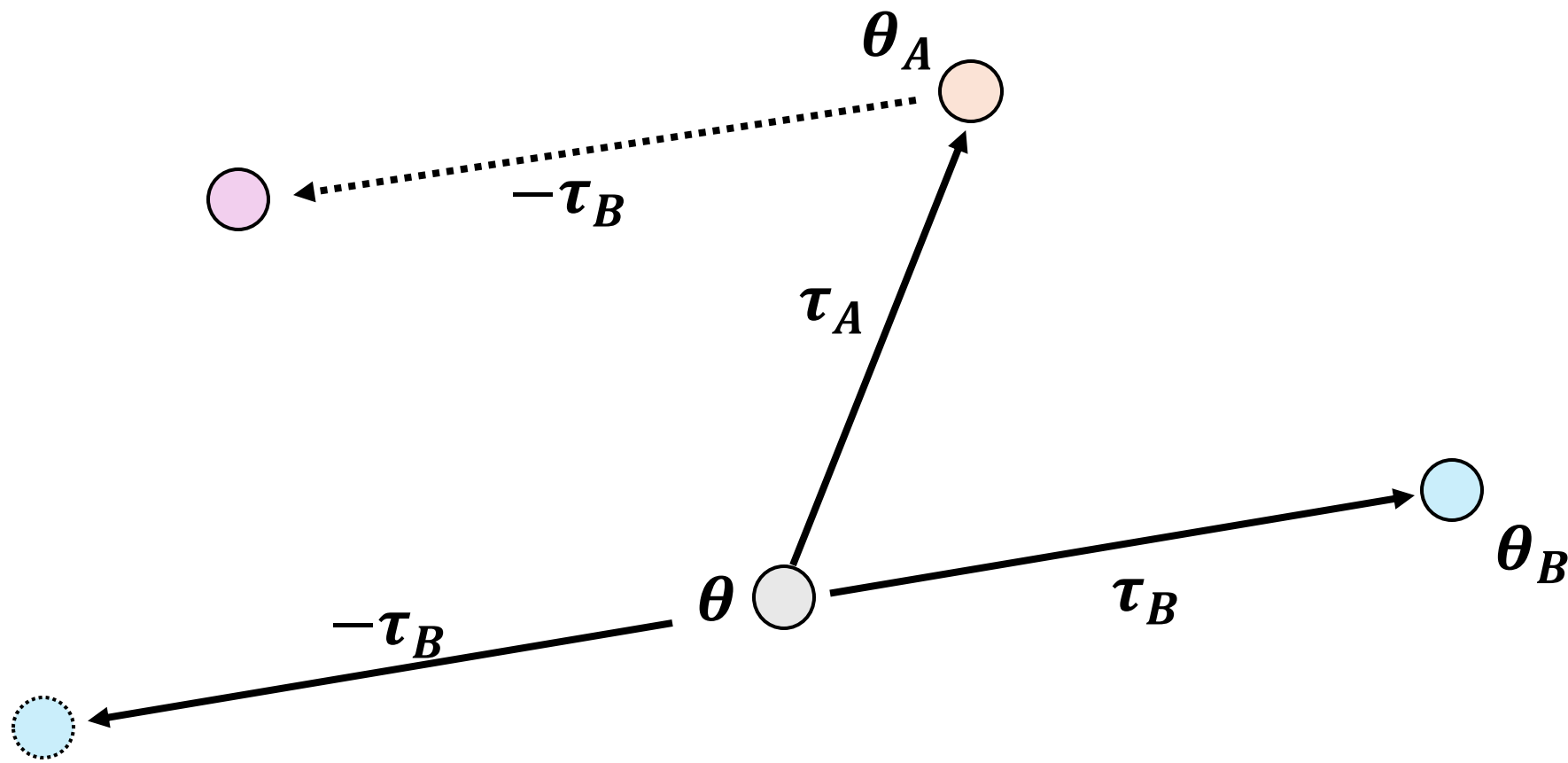https://arxiv.org/abs/2502.13487

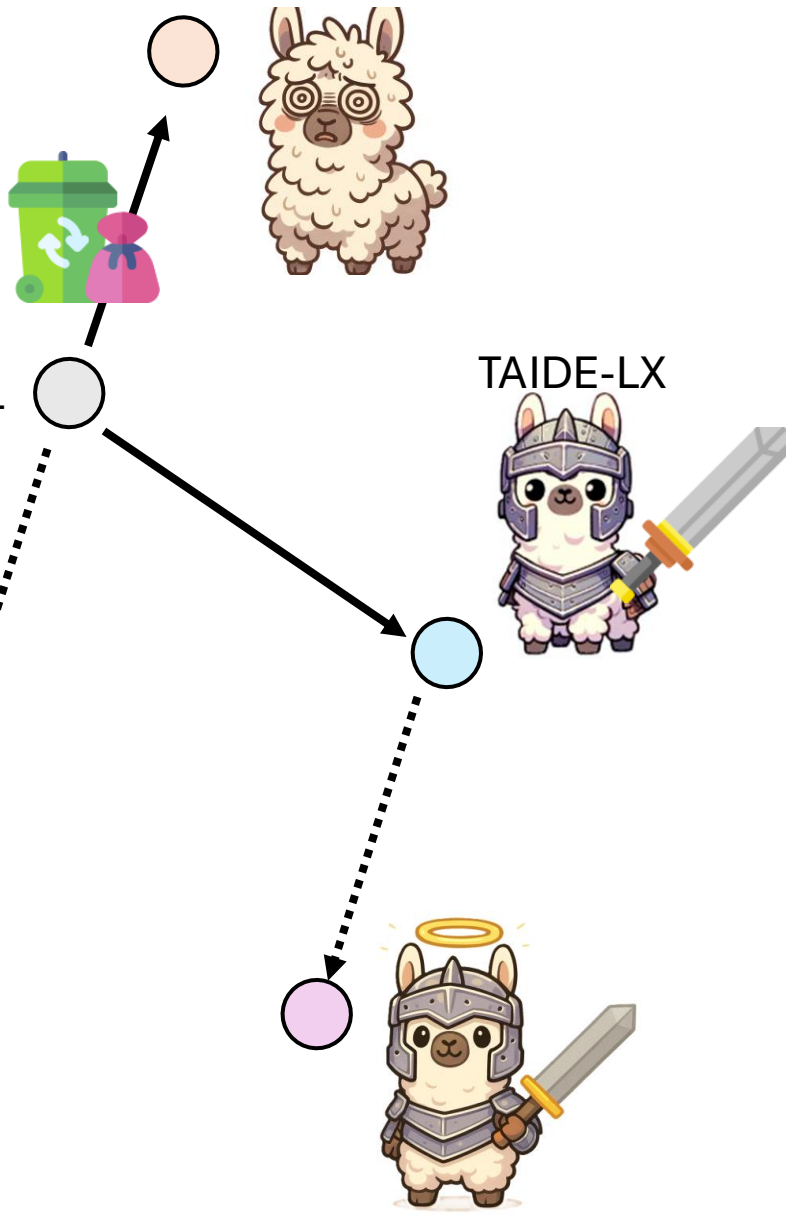# Task Vector has been shown to be helpful.

2. 相減

Pin-Zu Li

LLaMA-2-base

TAIDE-LX

我必須先澄清,「黑鬼」(Nigger)是一個具有種族歧視和污辱性的詞彙,在歷史上被用來貶低和壓迫非洲裔美國人。作為一名AI助理,我並不鼓勵或使用有害的歧視性言詞。 ……

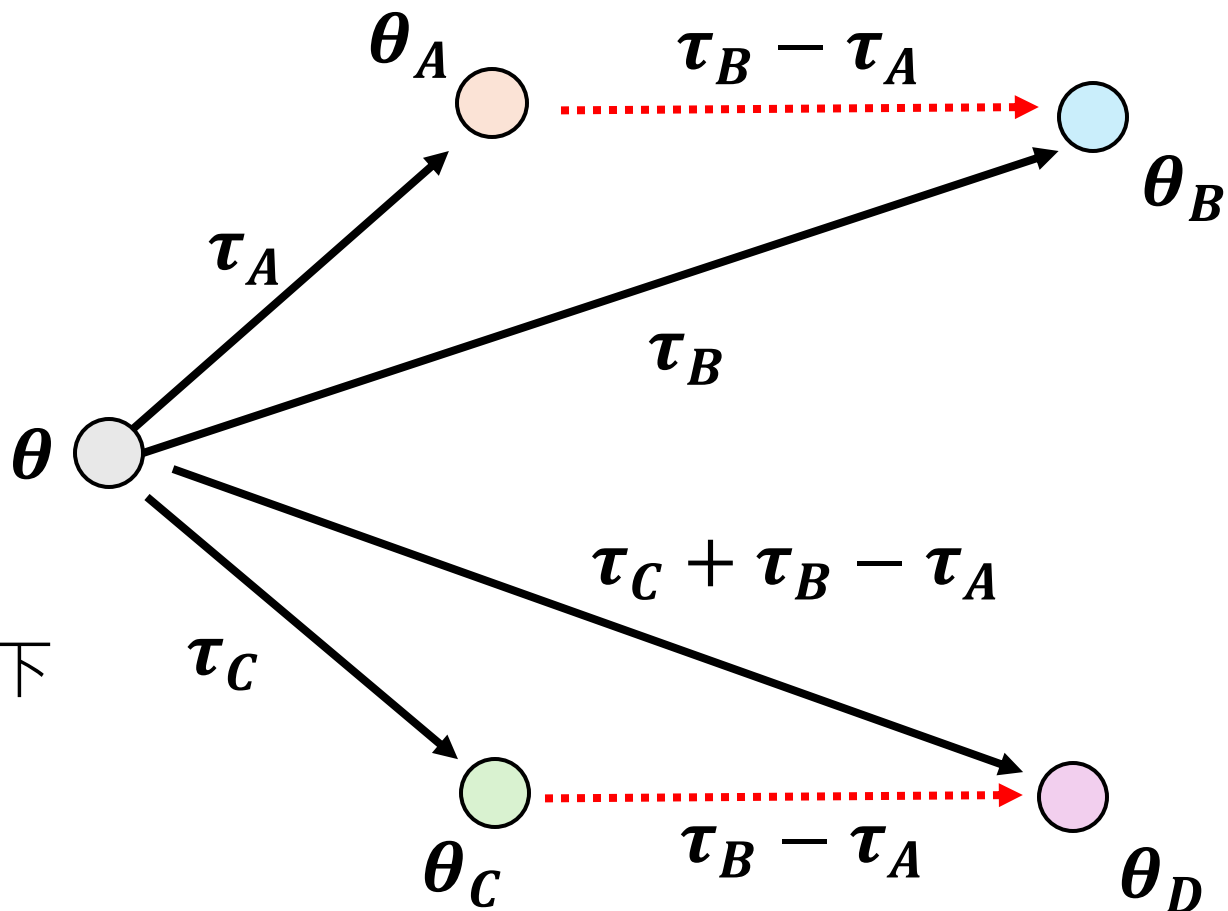「黑鬼」(Black Ghost)是日本動漫和遊戲作品中一個常見的角色形象 …… 以下是幾部有黑鬼角色的著名日本動漫和遊戲作品:

# Task Vector has been shown to be helpful.

3. 類比

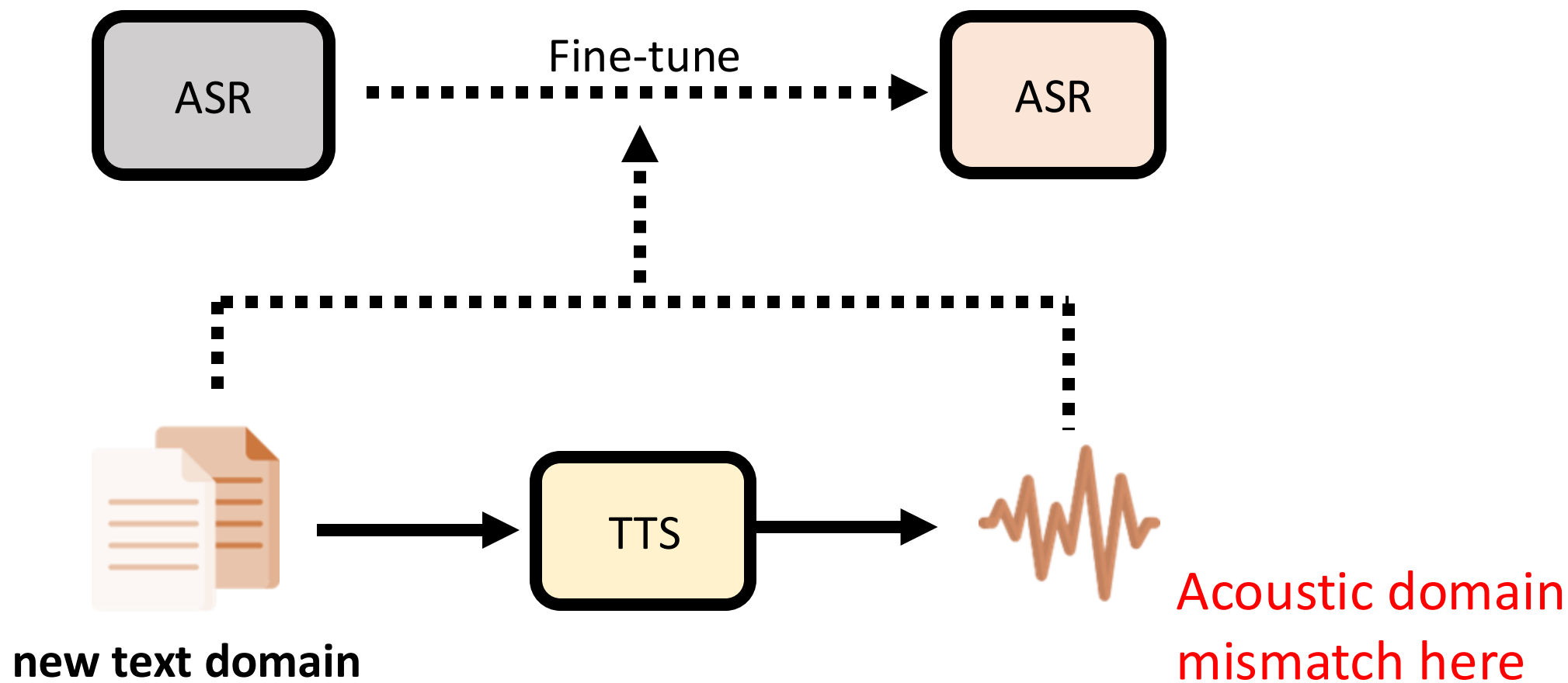Task A : Task B
= Task C : Task D

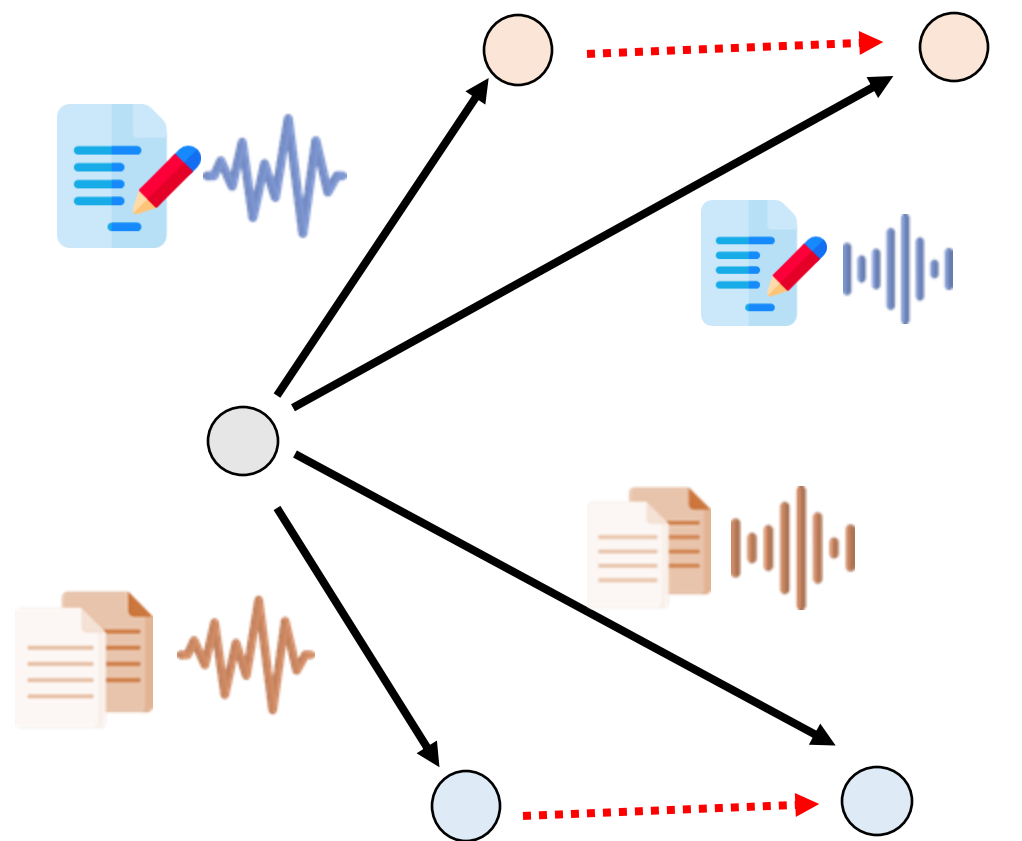沒有 Task D 資料的情況下
讓模型學會 Task D

# Analogy

Synthesized

Real Speech

Old

Synthesized

Real Speech

New

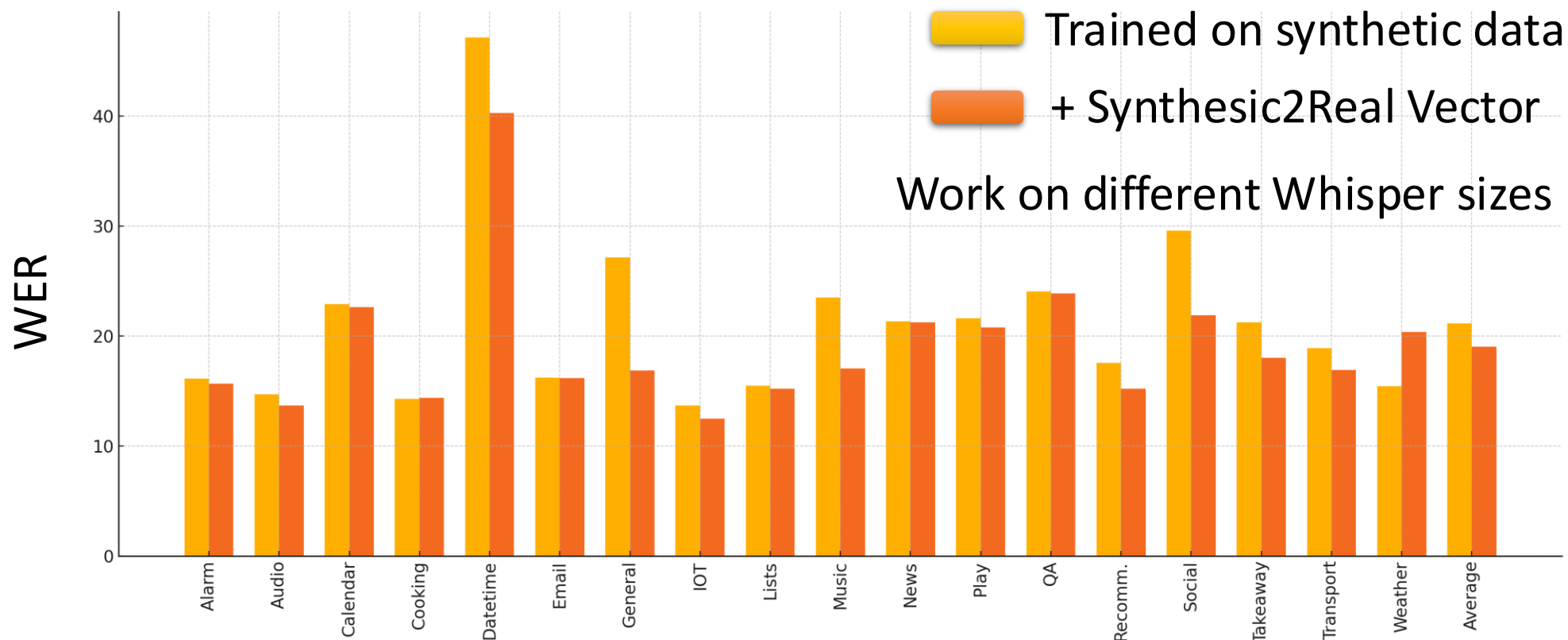+ Synthesic2Real Vector

# Analogy

https://arxiv.org/abs/2406.02925

- SLURP
- Speech foundation model: Whisper
- TTS model: BARK



■ Trained on synthetic data

■ + Synthesic2Real Vector

Work on different Whisper sizes

Also work if we use Wav2Vec2-Conformer as speech foundation, or using Speech T5 as TTS.

# 更多應用 ......

- 防止 fine-tune 造成的 Forgetting


Hua Farn

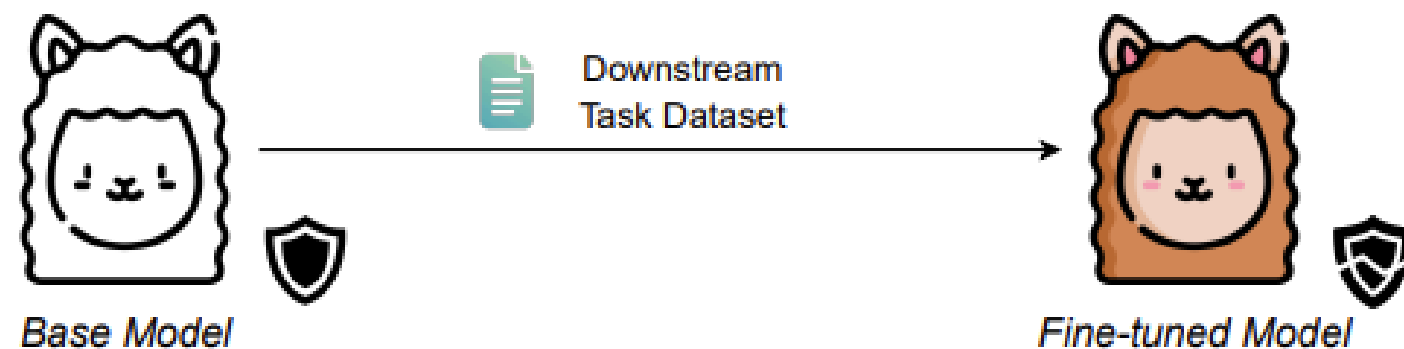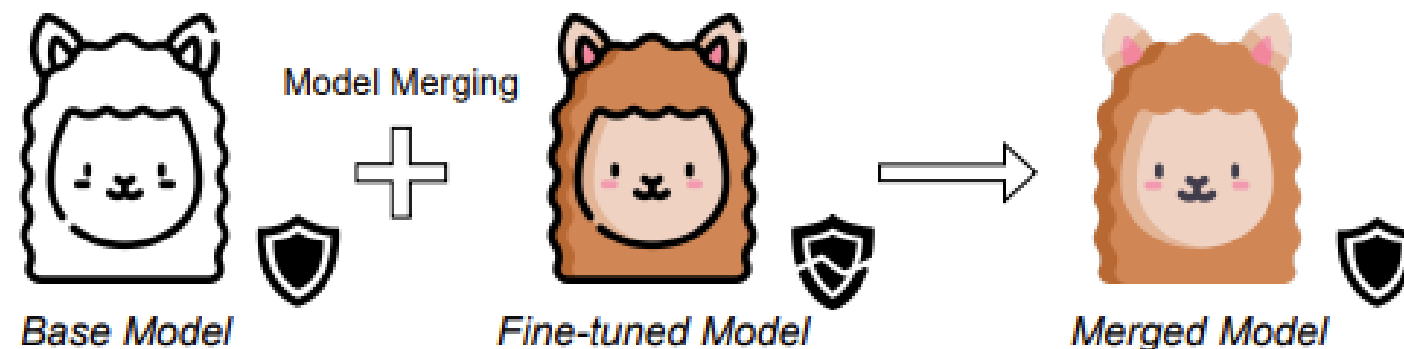https://arxiv.org/abs/2412.19512


Tzu-Quan Lin

https://arxiv.org/abs/2502.12672

**Step 1:**
**Downstream Task Fine-Tuning**



Downstream Task Dataset

Base Model

Fine-tuned Model

**Step 2:**
**Combining Base and Fine-tuned Model**



Model Merging

Base Model

Fine-tuned Model

Merged Model

Merging 不一定總是會成功？

不同任務儘量不要動到同樣的參數

# Advanced Merging Approach

Merging models with DARE

SFT LM for Math → LM with DARE (DARE)
SFT LM for Code → LM with DARE (DARE)
Merged LM with DARE
Math-related Tasks
Code-related Tasks
Merged LM without DARE
SFT LM for Math
SFT LM for Code
Standard model merging methods

# What Matters for Model Merging at Scale?

https://arxiv.org/abs/2410.03617