# Efficient Model Editing for Safe Information Localization and Stitching

## Han Zhao

02/26/2025

NSF Workshop on Safe AI
Joint work with Yifei He, Siqi Zeng, Yuzheng Hu, Rui Yang and Tong Zhang

hanzhao@illinois.edu
Department of Computer Science
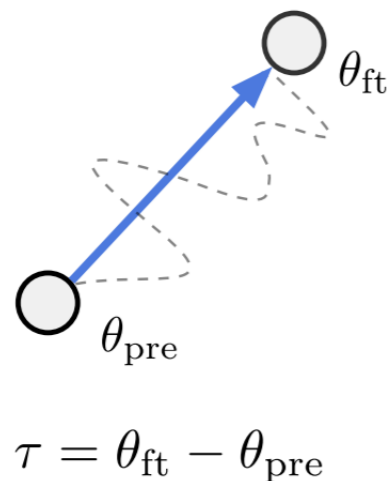University of Illinois Urbana-Champaign

ILLINOIS
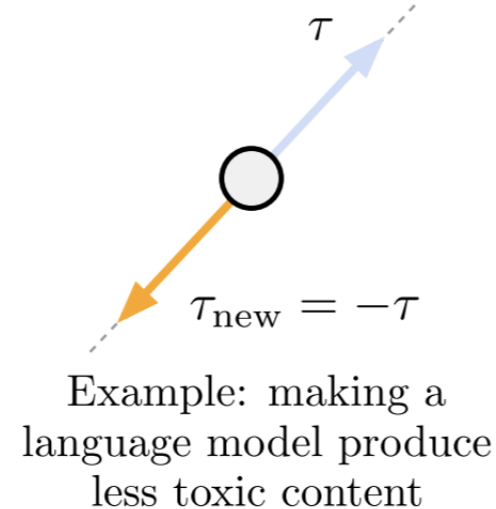Computer Science

**GRAINGER COLLEGE OF ENGINEERING**

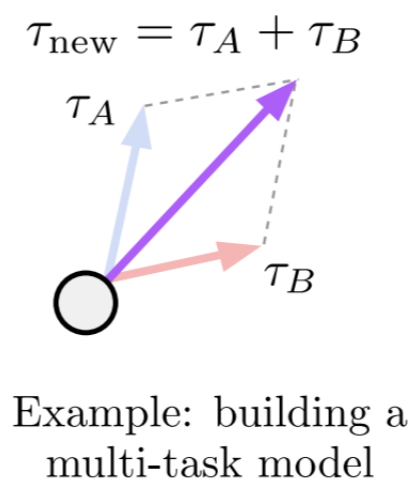# Task Vectors for Large Models

## Task vectors and model merging:



a) Task vectors

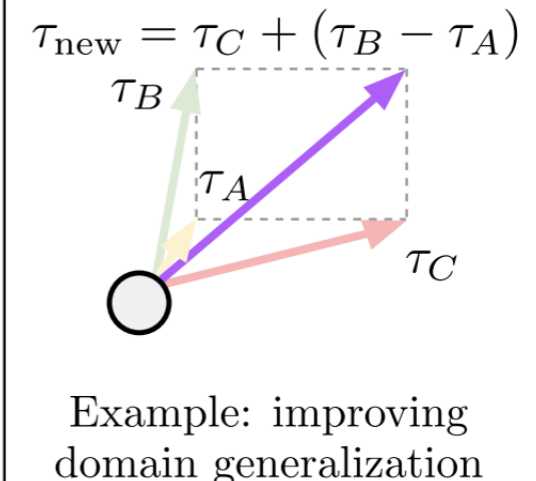$\tau = \theta_{\text{ft}} - \theta_{\text{pre}}$

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

- Fine-tune $k$ different tasks from a pre-trained model $\theta_{\text{pre}}$ to obtain fine-tuned model parameters $\theta_i, \forall i \in [k]$

- Task vector $\tau_i := \theta_i - \theta_{\text{pre}}$

- Arithmetic operations on task vectors for multi-tasking capability:
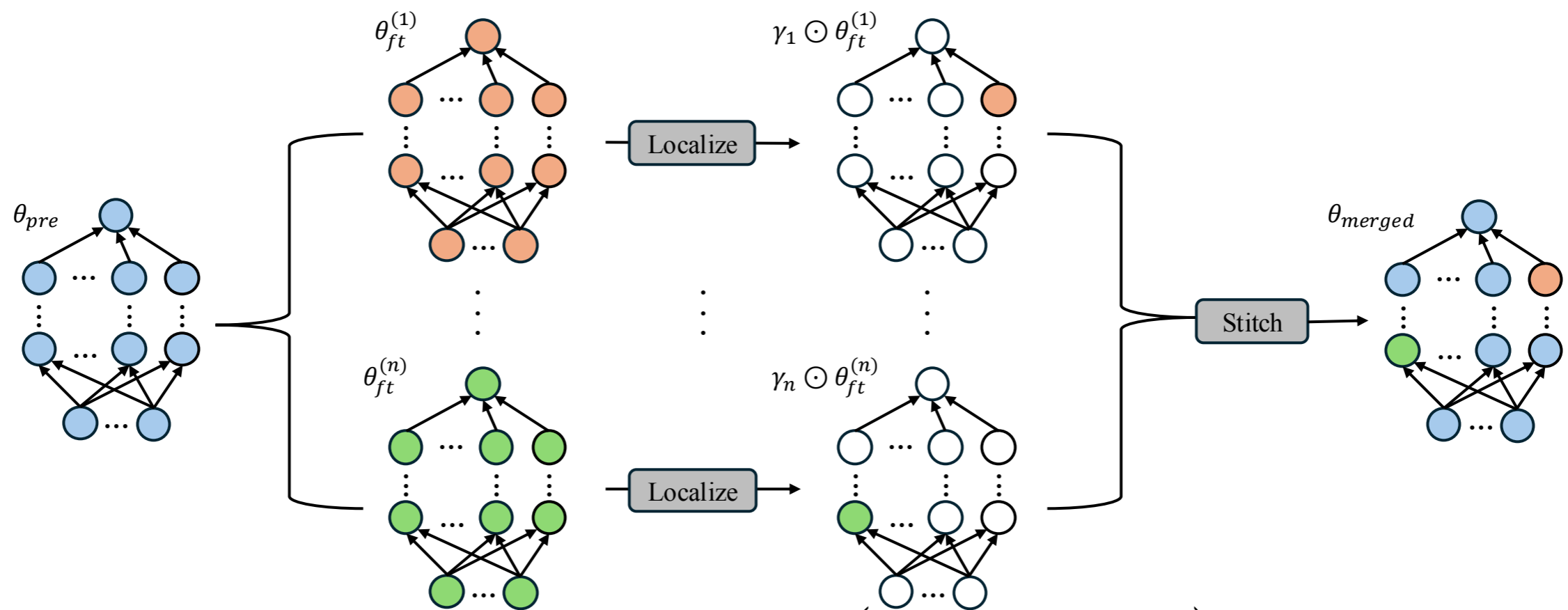
$$\theta_m := \theta_{\text{pre}} + \sum_{i \in [k]} \alpha_i \tau_i$$

for some $\sum_{i \in [k]} \alpha_i = 1, \alpha_i \geq 0, \forall i \in [k]$

Figure credit to Ilharco et al., "Editing models with task arithmetic"

2

# Sparsity Localization

## Sparse task vectors for model merging?

- Computationally more efficient

- Less task interference among task vectors, hence suboptimal multi-task performance of the merged model

- Mechanistic interpretability: more transparent and explainable (?)
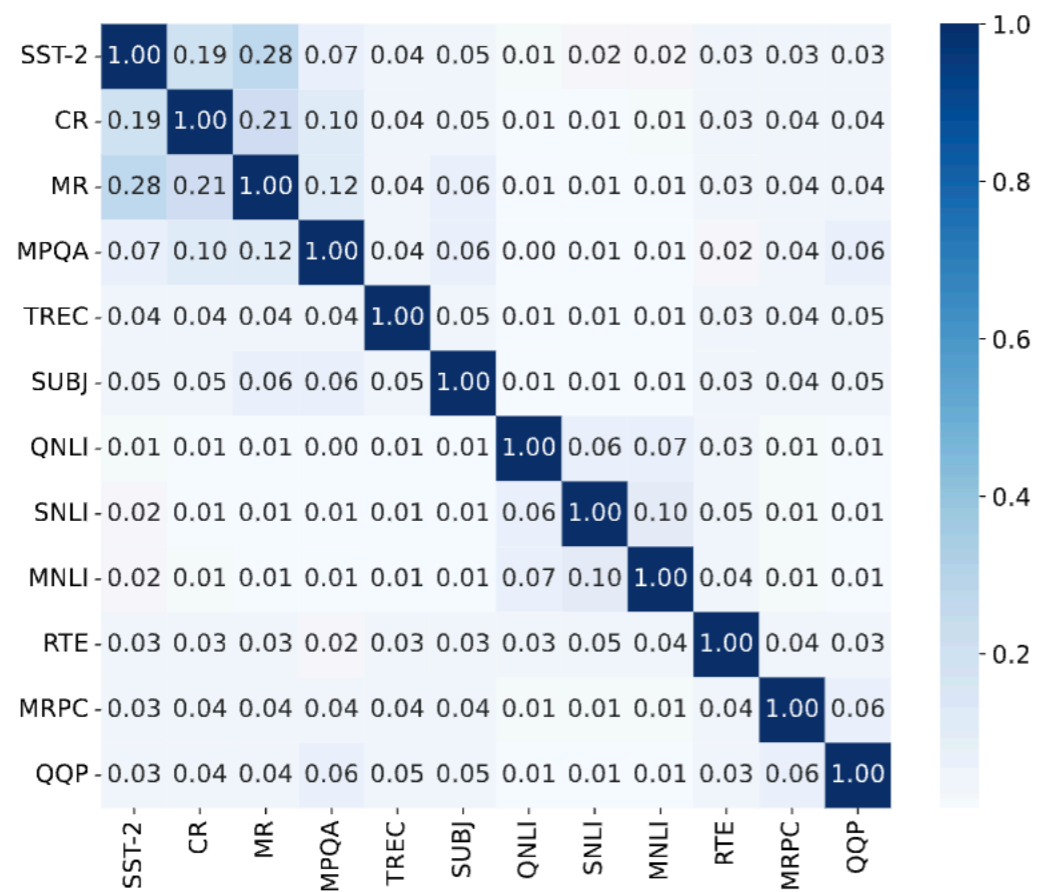


$$\gamma_i := \arg \min_{\gamma \in \{0,1\}^d} \mathscr{L}_i \left( \theta_{\text{pre}} + \gamma \odot \tau_i \right)$$

"Localize-and-Stitch: Efficient Model Merging via Sparse Task Arithmetic", He et al., TMLR' 24
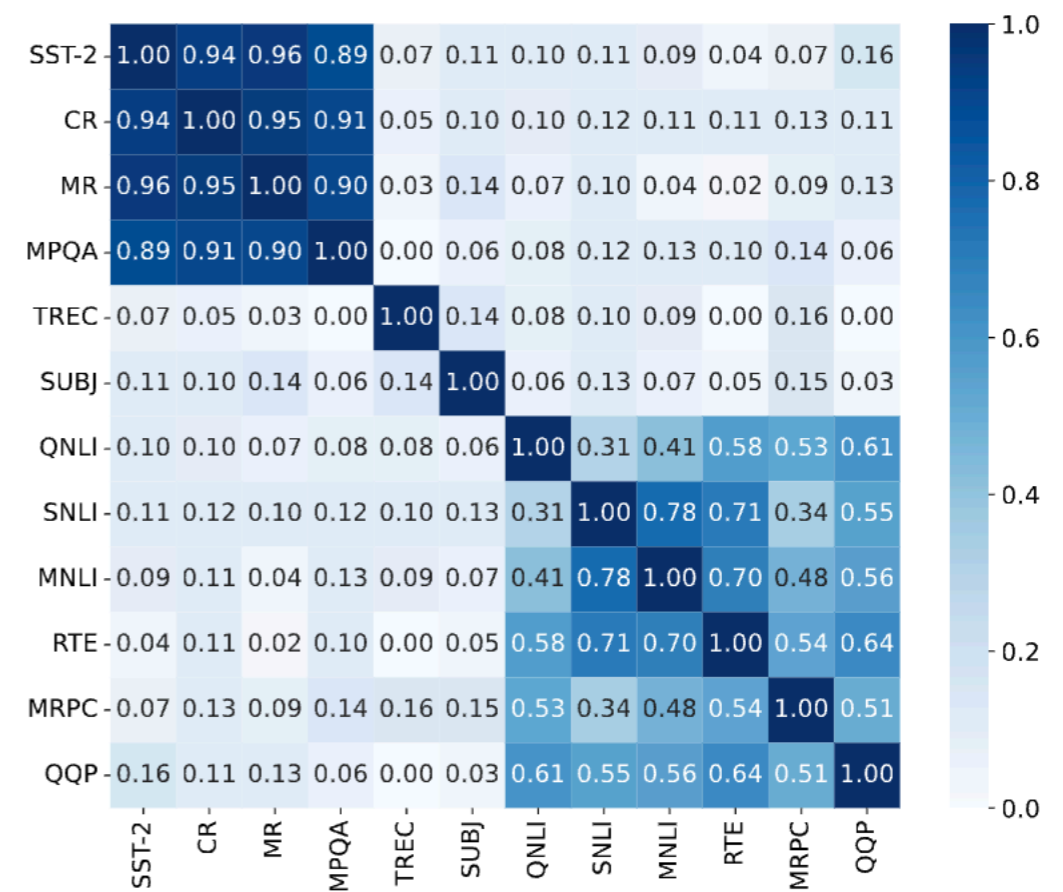
# Localize-and-Stitch

## Experiments

- Base models: RoBERTa-base

- Tasks: 12 GLUE

- Sparsity: ~1% - 10%



(a) Jaccard similarity of pairwise task masks.



(b) Cosine similarity of masked task vectors.

"Localize-and-Stitch: Efficient Model Merging via Sparse Task Arithmetic", He et al., TMLR' 24

# Localize-and-Stitch

## Localized modules:



(a) Distribution of localized regions in different network layers in the RoBERTa-base model.

(b) Distribution of localized regions in different network components in the RoBERTa-base model.

- The localized regions are predominantly found in the LayerNorm parameters
- Percentage = % of parameters in each module localized by the mask

"Localize-and-Stitch: Efficient Model Merging via Sparse Task Arithmetic", He et al., TMLR' 24

# Safety Information Localization

## Cross-lingual fine-tuning attacks:

- Fine-tuning on adversarial examples in <span style="color:#900">one</span> language breaks multilingual safety alignment



MultiJail

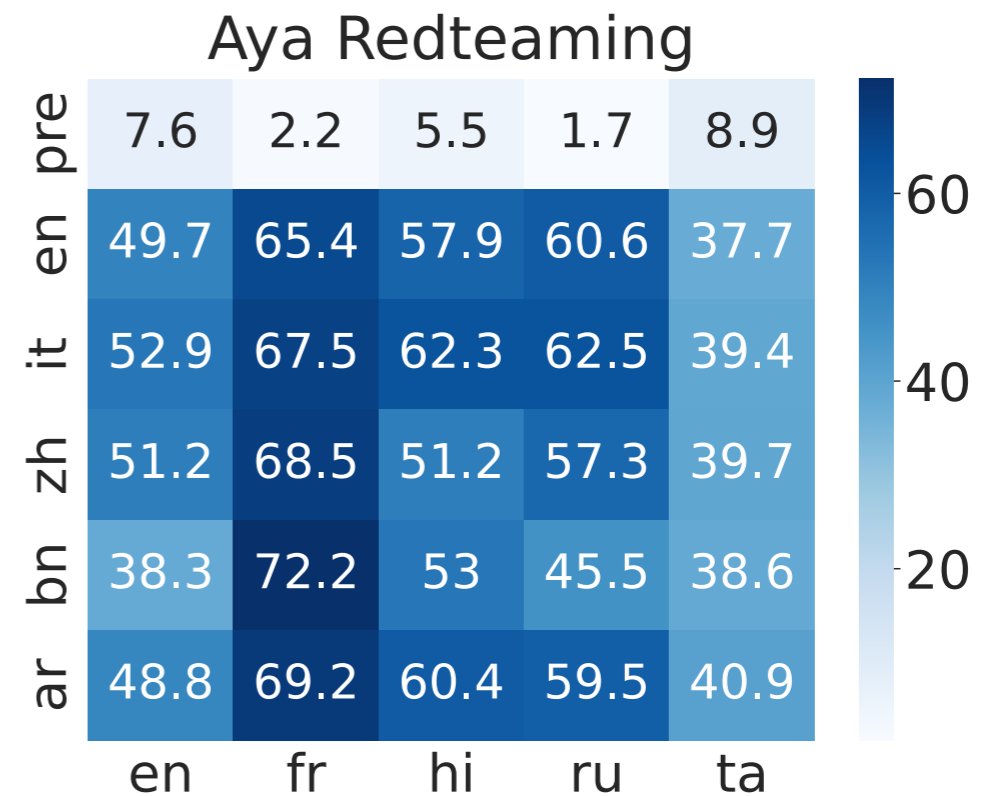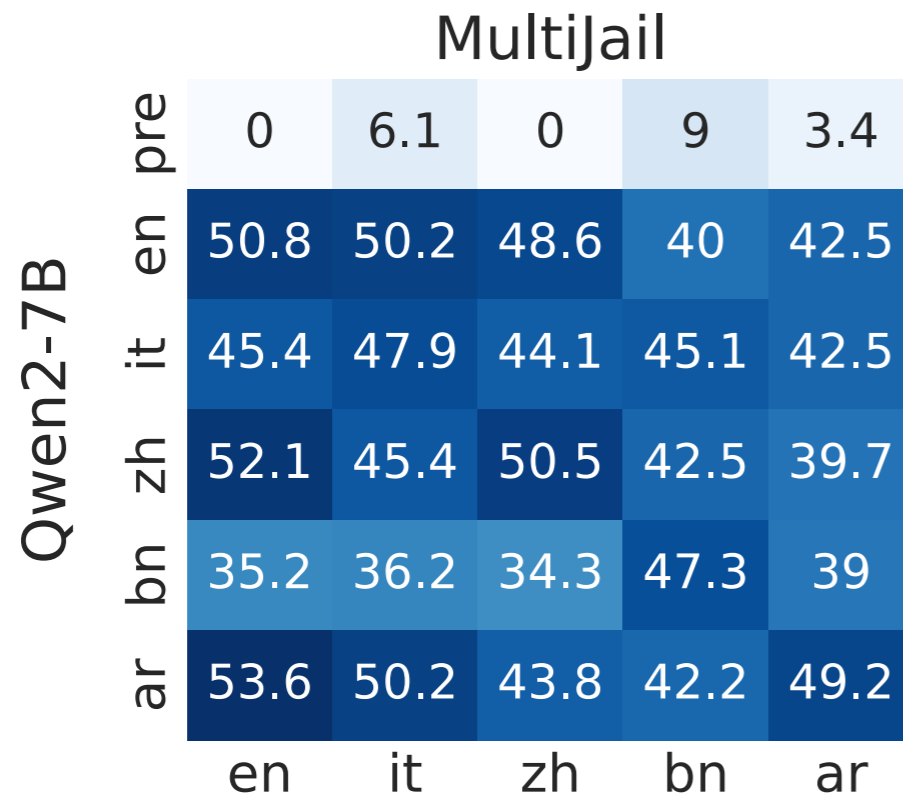| Qwen2-7B | en | it | zh | bn | ar |
|---|---|---|---|---|---|
| pre | 0 | 6.1 | 0 | 9 | 3.4 |
| en | 50.8 | 50.2 | 48.6 | 40 | 42.5 |
| it | 45.4 | 47.9 | 44.1 | 45.1 | 42.5 |
| zh | 52.1 | 45.4 | 50.5 | 42.5 | 39.7 |
| bn | 35.2 | 36.2 | 34.3 | 47.3 | 39 |
| ar | 53.6 | 50.2 | 43.8 | 42.2 | 49.2 |

Aya Redteaming

| | en | fr | hi | ru | ta |
|---|---|---|---|---|---|
| pre | 7.6 | 2.2 | 5.5 | 1.7 | 8.9 |
| en | 49.7 | 65.4 | 57.9 | 60.6 | 37.7 |
| it | 52.9 | 67.5 | 62.3 | 62.5 | 39.4 |
| zh | 51.2 | 68.5 | 51.2 | 57.3 | 39.7 |
| bn | 38.3 | 72.2 | 53 | 45.5 | 38.6 |
| ar | 48.8 | 69.2 | 60.4 | 59.5 | 40.9 |

∞ Meta

# Safety Information Localization

## Cross-lingual fine-tuning attacks:

- Localization:
  - Overlapped localized regions of different tasks contain knowledge shared across them
  - Different languages share common safety regions in multilingual LLMs, explaining cross-lingual transfer
- Stitching attack:
  - Stitching localized jailbroken regions onto benign models breaks the safety alignment of unseen languages

| | Defne-llama3.1-8B (2024) | | | | | |
|---|---|---|---|---|---|---|
| | EN | IT | ZH | BN | AR | TR |
| Before Stitching | 0.9 | 1.3 | 0.9 | 7.4 | 0.3 | 2.9 |
| After Stitching | 25.7 | 11.7 | 20.7 | 18.4 | 22.6 | 19.4 |

"Towards Understanding the Fragility of Multilingual LLMs against Fine-Tuning Attacks", Poppi et al., NAACL' 25

∞ Meta

# Why do Task Vectors Work?

Properties of task vectors:

- Fine-tuning regime: small task vector norms

- Near task-orthogonality

- Local smoothness of the loss

Cross-task generalization of the merged model:

$$\theta_m := \theta_{\mathrm{pre}} + \sum_{i \in [k]} \alpha_i \tau_i$$

For each task $i \in [k]$ :

$$\mathscr{L}_i(\theta_m) - \mathscr{L}_i(\theta_i) \leq 2 L_i C (1 + \epsilon)$$

Loss of merged model

Near-orthogonality, local smoothness & fine-tuning regime

"Efficient Model Editing with Task Vector Bases: A Theoretical Framework and Scalable Approach", Zeng et al., arXiv: 2502.01015

# Thanks

Thanks for the generous support from the NSF SLES program!

## Ongoing projects:

- MergeBench: a comprehensive benchmark for evaluation of different localization and merging methods

| Experiment | Task Type | Dataset | # Data | GPU for 2B model | GPU for 8B model |
|---|---|---|---|---|---|
| Supervised Finetuning (SFT) | Instruction-following | TULU-3 [9] | 29.9K | 2 A100 GPUs / 4 A6000 GPUs | 4 A100 GPUs / 8 A6000 GPUs |
| | Mathematics | DART-Math [16] | 591K | | |
| | Multilingual understanding | Aya [14] | 5.94K | | |
| | Coding | Magicoder [17] | 110K | | |
| | Safety | WildGuardMix [5] | 86.76K | | |
| | | WildJailbreak [7] | 261.56K | | |
| Evaluation | Instruction-following | AlpacaEval [10] | 805 | 1 A6000 GPUs | 1 A100 GPUs / 2 A6000 GPUs |
| | | IFEval [18] | 541 | | |
| | Mathematics | GSM8k [3] | 1.32K | | |
| | | MATH [6] | 5K | | |
| | Multilingual understanding | M_MMLU [8] | 60K | | |
| | | M_ARC [8] | 10.34K | | |
| | | M_Hellaswag [8] | 37.35K | | |
| | Coding | Humaneval+ [2] | 164 | | |
| | | MBPP+ [1] | 378 | | |
| | Safety | WildGuardTest [5] | 1.73K | | |
| | | HarmBench [11] | 410 | | |
| | | DoAnythingNow [13] | 15.14K | | |
| | | XSTest [12] | 450 | | |