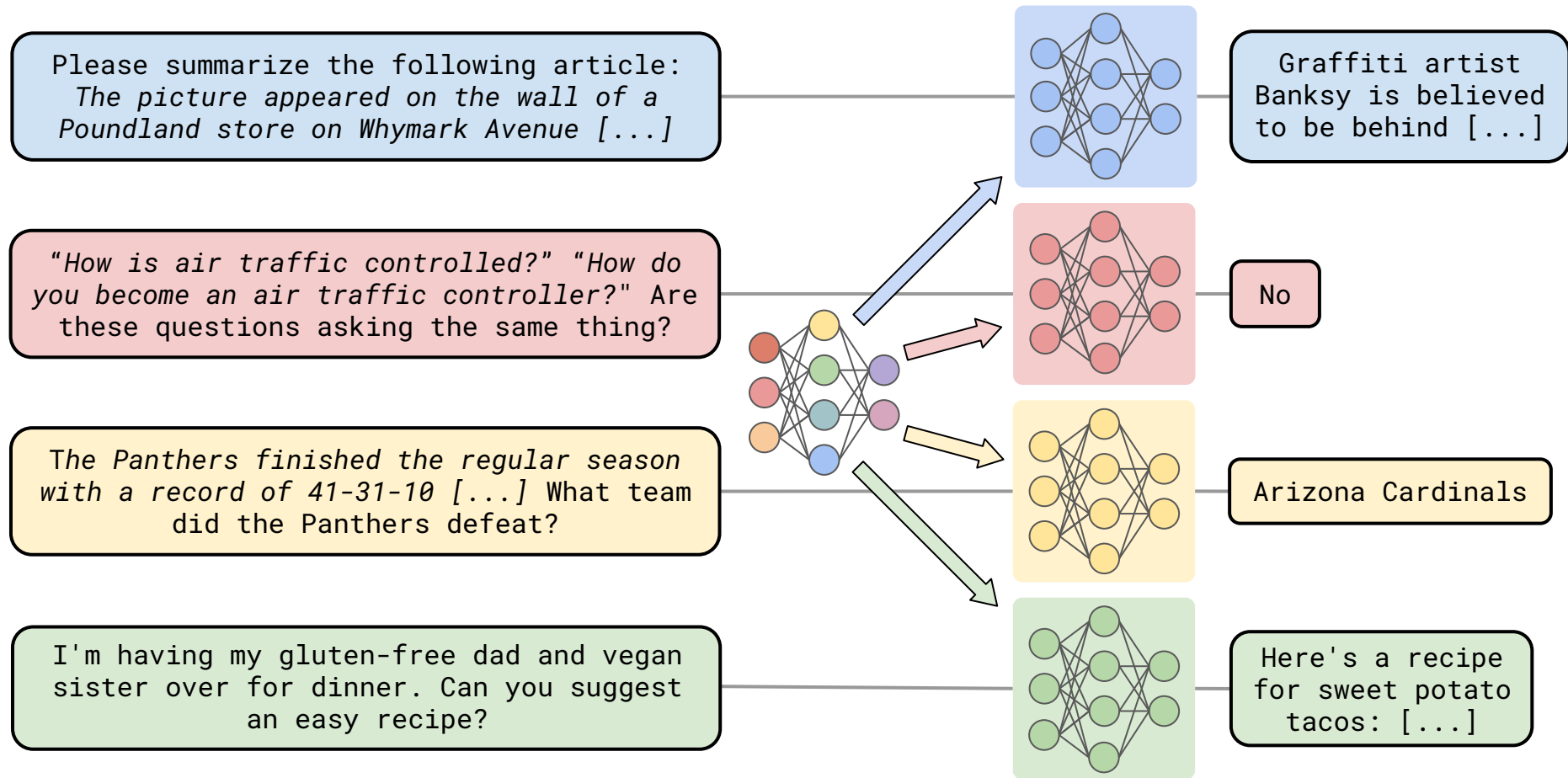


Build an Ecosystem, Not a Monolith

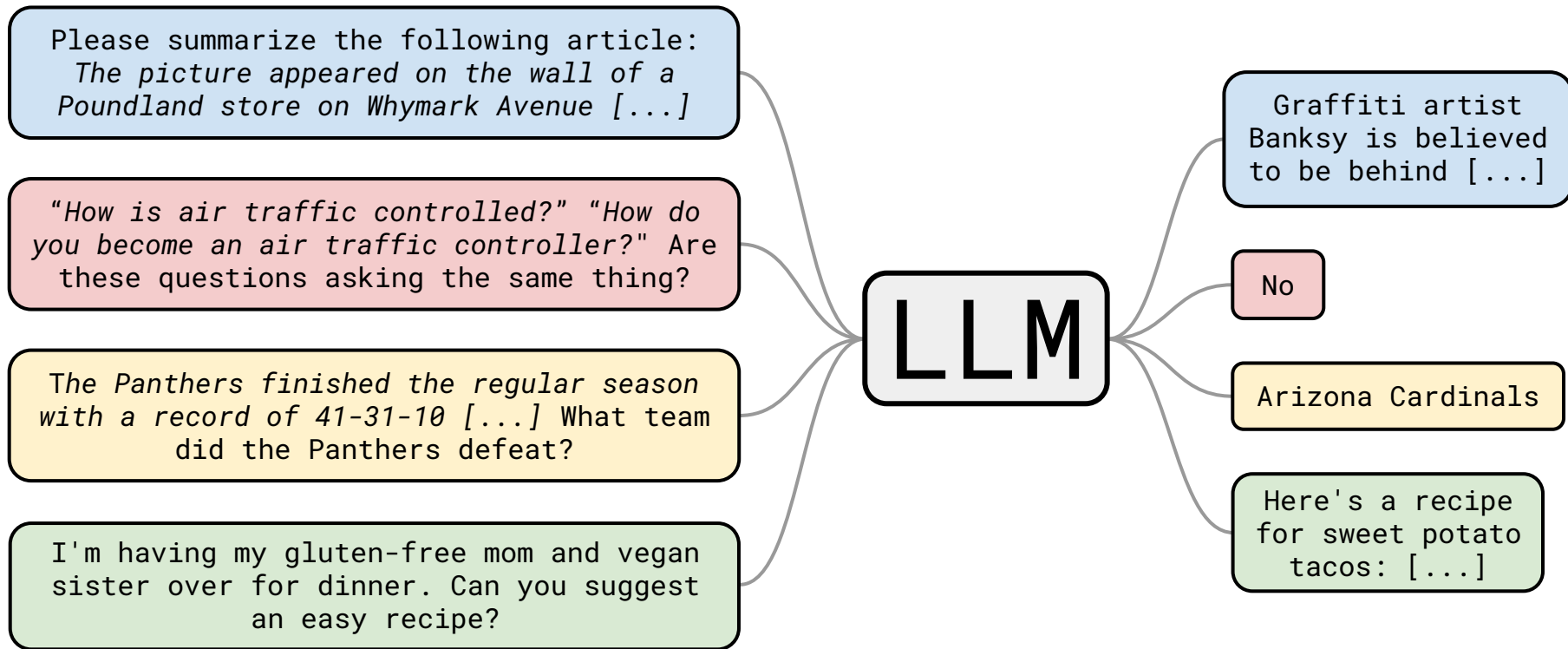
Colin Raffel



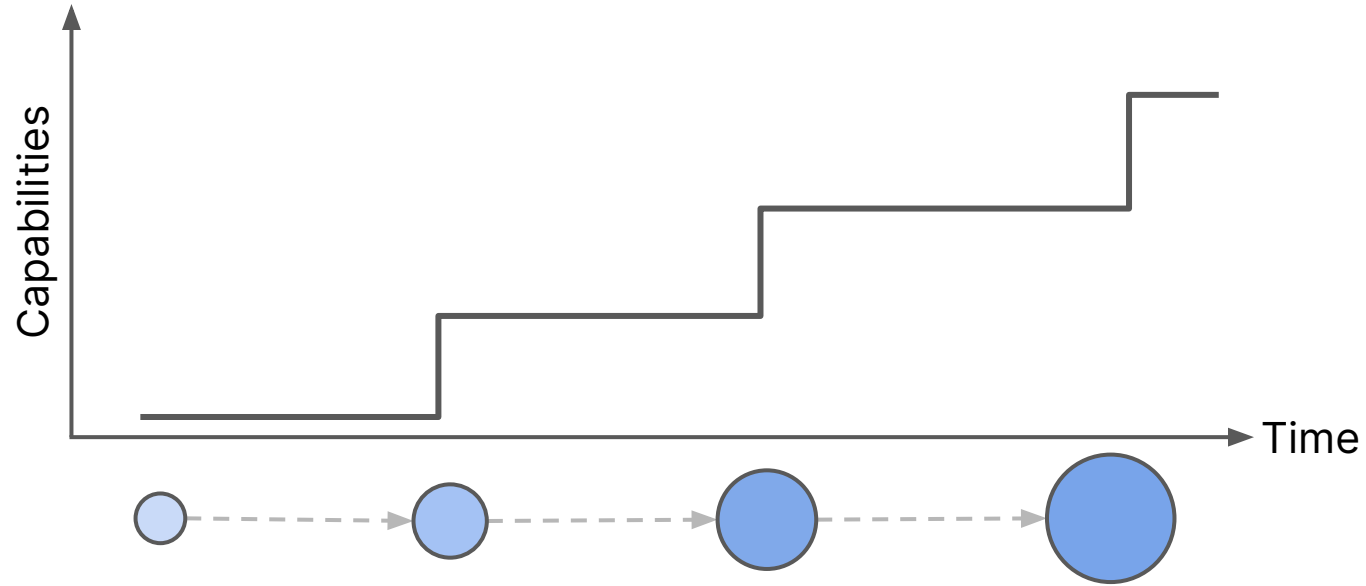
Transfer learning: fine-tuning to create specialized models



LLMs as general-purpose monolithic models



Monolithic model development involves wholesale replacement



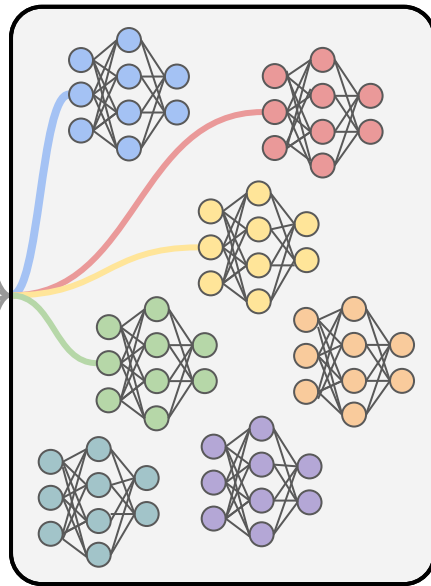
Ecosystems of specialist models?

Please summarize the following article:
*The picture appeared on the wall of a
Poundland store on Whymark Avenue [...]*

"How is air traffic controlled?" "How do
you become an air traffic controller?" Are
these questions asking the same thing?

*The Panthers finished the regular season
with a record of 41-31-10 [...]* What team
did the Panthers defeat?

I'm having my gluten-free mom and vegan
sister over for dinner. Can you suggest
an easy recipe?



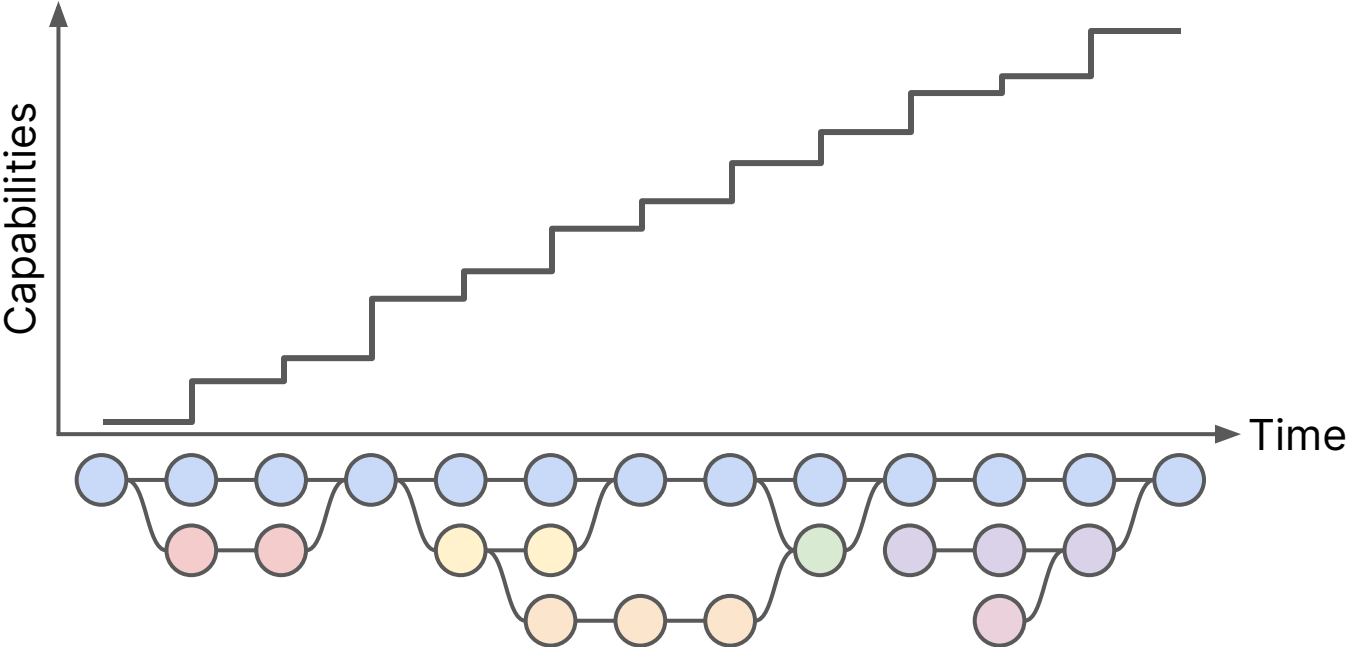
Graffiti artist
Banksy is believed
to be behind [...]

No

Arizona Cardinals

Here's a recipe
for sweet potato
tacos: [...]

Collaborative ecosystem development will lead to continual improvements

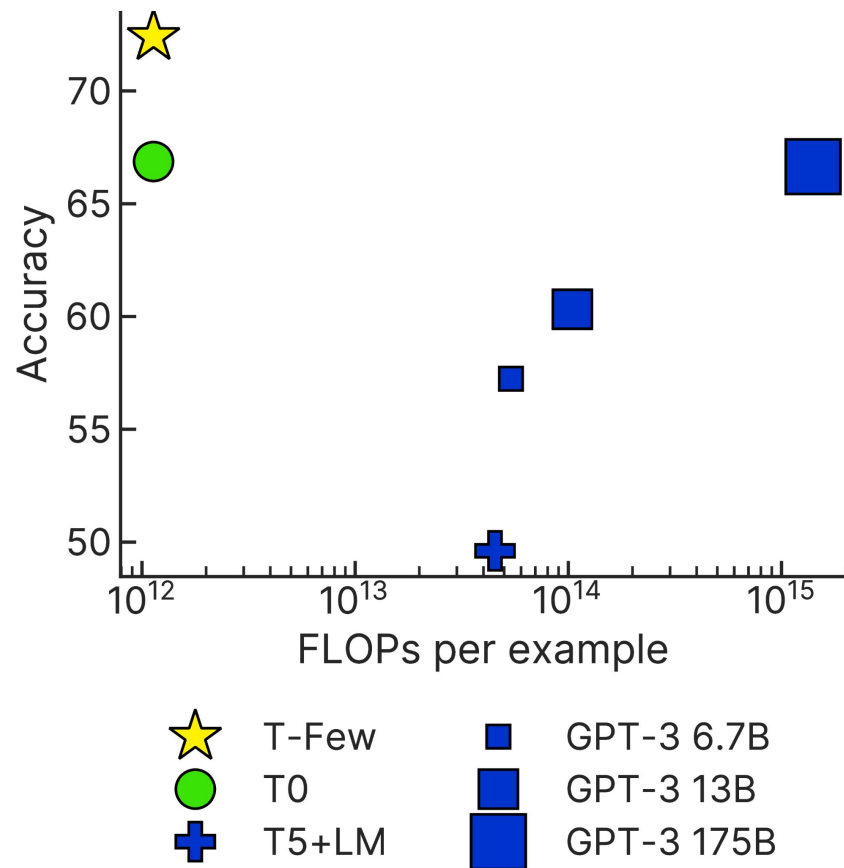


How and why should we build ecosystems of specialist models instead of monolithic models?

How and why should we build ecosystems of specialist models instead of monolithic models?

Specialist models are often **cheaper** and sometimes **better**.

Smaller fine-tuned models often outperform larger generalist models



From "Few-Shot Parameter-Efficient Fine-Tuning is Better and Cheaper than In-Context Learning" by Liu et al.

Smaller fine-tuned models often outperform larger generalist models

*Specialist
models*

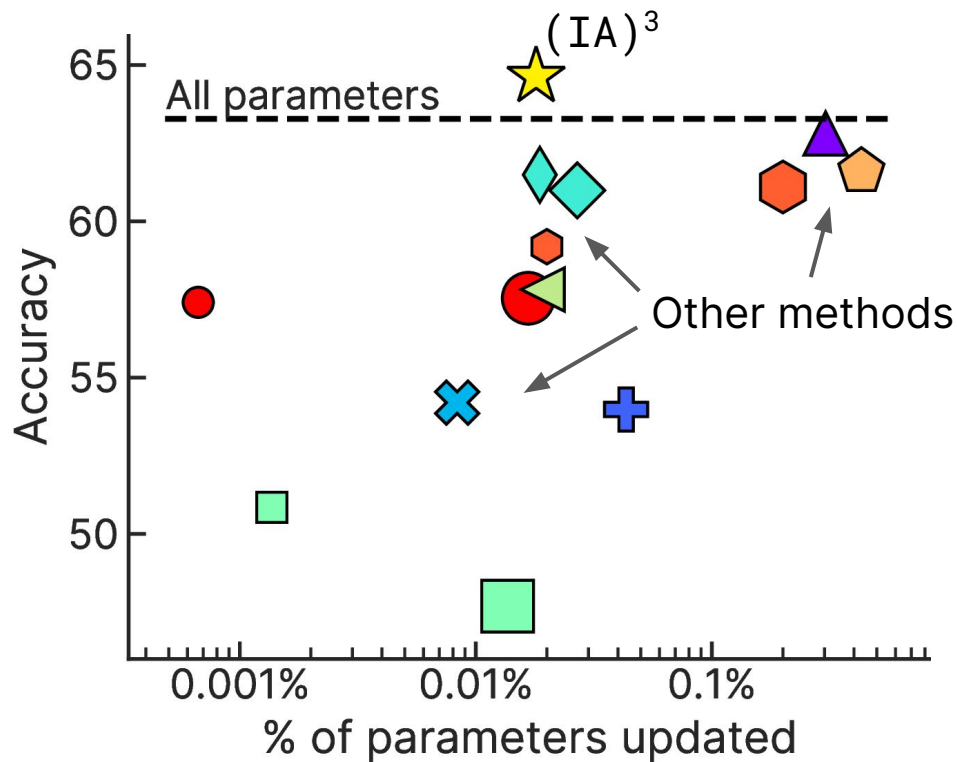
	SOTA	GPT-4	PaLM	PaLM 2
WinoGrande	87.5 ^a	87.5 ^a ₍₅₎	85.1 ^b ₍₅₎	90.9 ₍₅₎
ARC-C	96.3^a	96.3^a ₍₂₅₎	88.7 ^c ₍₄₎	95.1 ₍₄₎
DROP	88.4^d	80.9 ^a ₍₃₎	70.8 ^b ₍₁₎	85.0 ₍₃₎
StrategyQA	81.6 ^c	-	81.6 ^c ₍₆₎	90.4 ₍₆₎
CSQA	91.2^e	-	80.7 ^c ₍₇₎	90.4 ₍₇₎
XCOPA	89.9 ^g	-	89.9 ^g ₍₄₎	94.4 ₍₄₎
BB Hard	65.2 ^f	-	65.2 ^f ₍₃₎	78.1 ₍₃₎

	Chinese→English		English→German	
	BLEURT ↑	MQM (Human) ↓	BLEURT ↑	MQM (Human) ↓
PaLM	67.4	3.7	71.7	1.2
Google Translate	68.5	3.1	73.0	1.0
PaLM 2	69.2	3.0	73.3	0.9

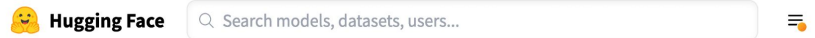
How and why should we build ecosystems of specialist models instead of monolithic models?

Each specialist model can be a **cheaply communicable** update to a base model.

$(IA)^3$ outperforms standard training while updating 0.01% of parameters



Existing "adapter" hubs have thousands of specialized models



Models 4,473 [Filter by name](#) [new Full-text search](#) [Edit filters](#) [Sort: Most Likes](#)


Active filters: [peft](#) [Clear all](#)

 [fb700/chatglm-fitness-RLHF](#)


Updated about 9 hours ago • ❤️ 192

 [chainyo/alpaca-lora-7b](#)


Updated Mar 29 • ⬇️ 77 • ❤️ 65

 [dfurman/falcon-40b-openassistant-peft](#)


 Text Generation • Updated 21 days ago • ⬇️ 228 • ❤️ 39

 [shareAI/llama2-13b-Chinese-chat](#)


 Question Answering • Updated 10 days ago • ⬇️ 60 • ❤️ 27

 [crumb/Instruct-GPT-J](#)

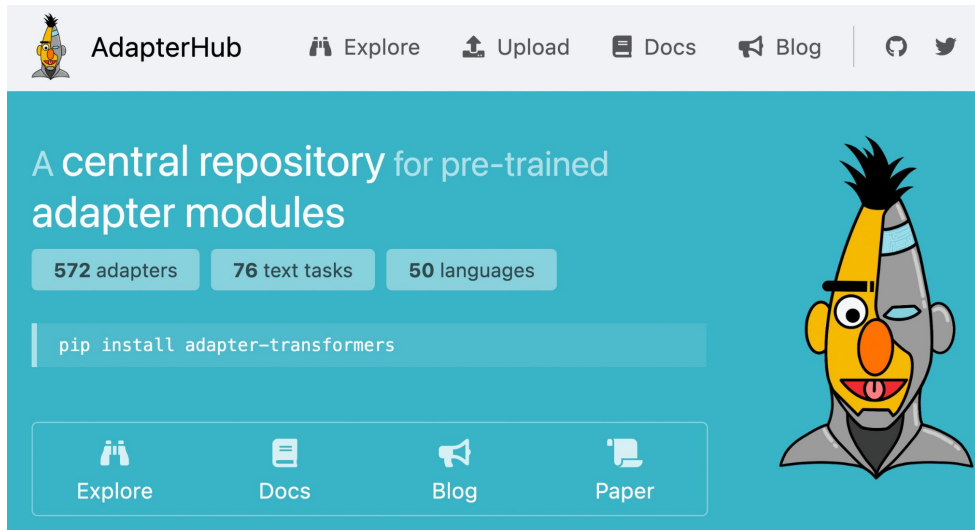
Updated Mar 26 • ❤️ 24

 [dominguesm/alpaca-lora-ptbr-7b](#)

Updated Apr 11 • ⬇️ 117 • ❤️ 17

 [Junity/Genshin-World-Model](#)

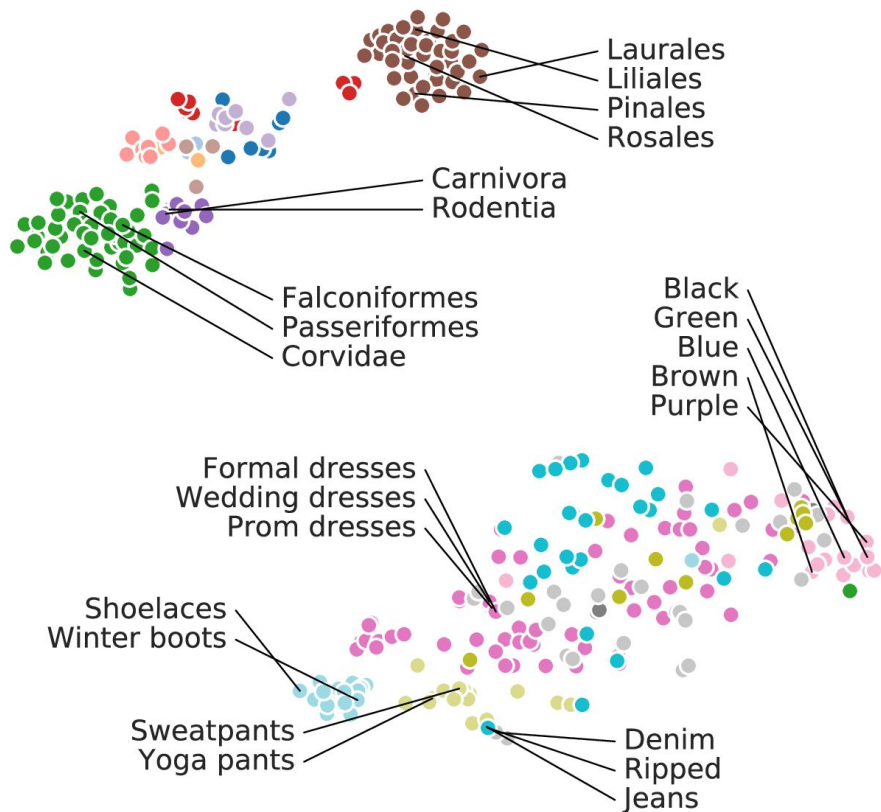
Updated 2 days ago • ⬇️ 8 • ❤️ 11



How and why should we build ecosystems of specialist models instead of monolithic models?

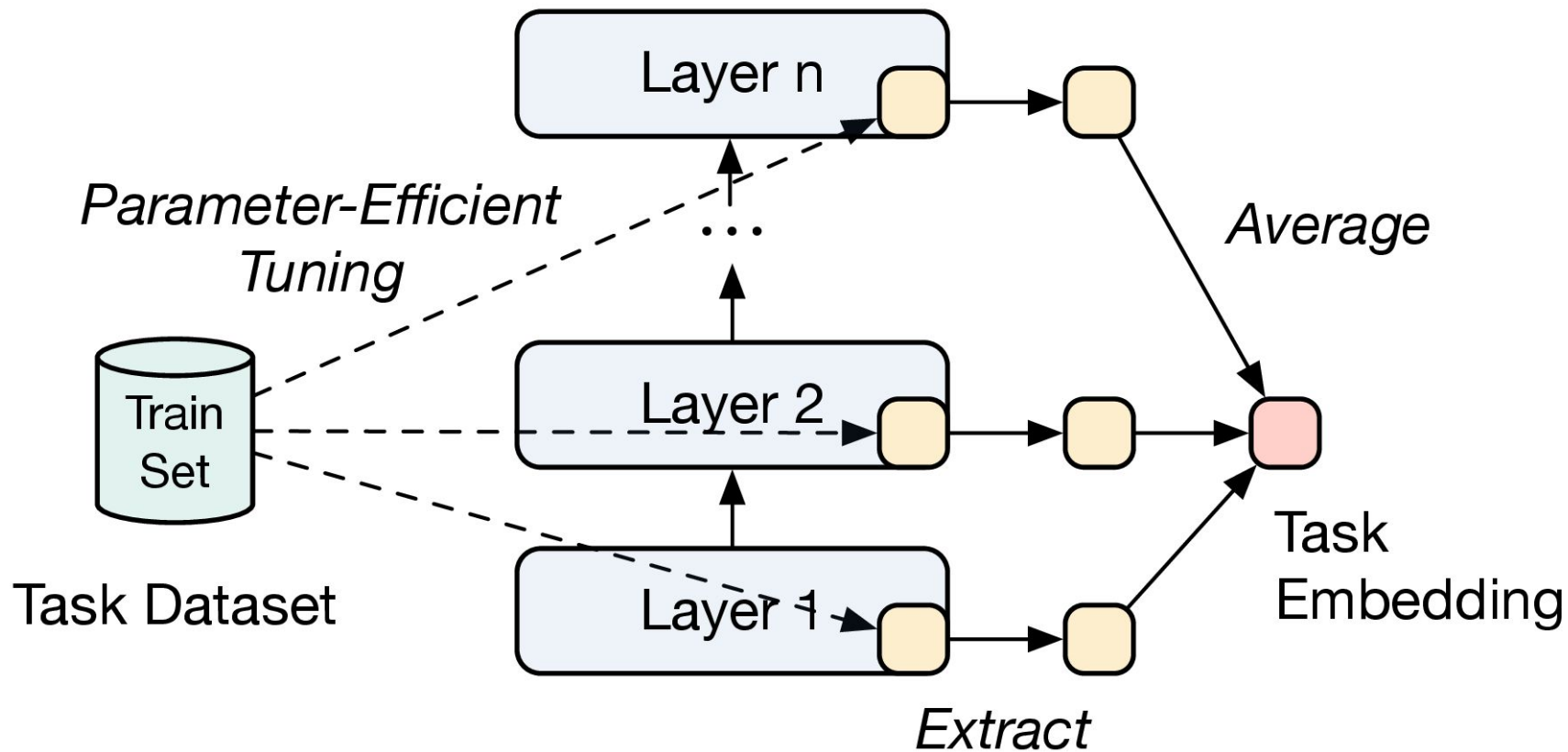
The appropriate model for a query should be **chosen automatically**.

task2vec encodes task similarity via the Fisher information matrix

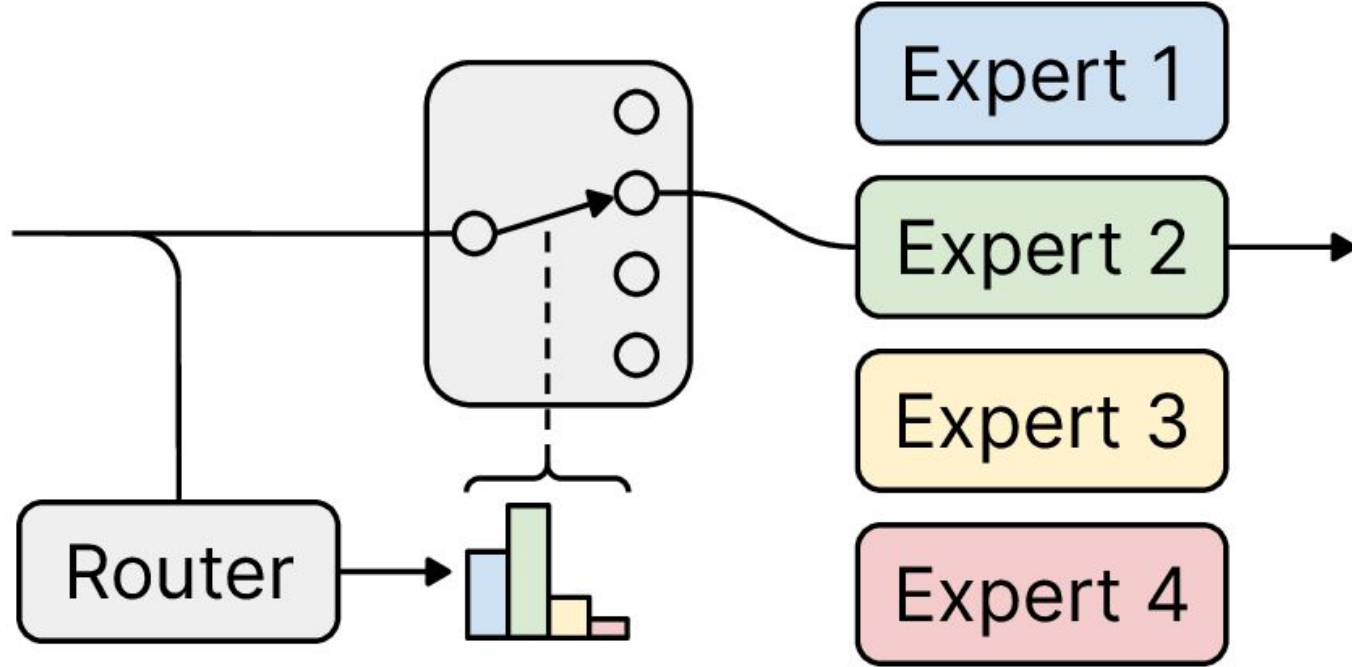


From "task2vec: Task Embedding for Meta-Learning" by Achille et al.

Adapter parameters also encode task similarity



Mixture-of-experts models perform adaptive routing inside the model

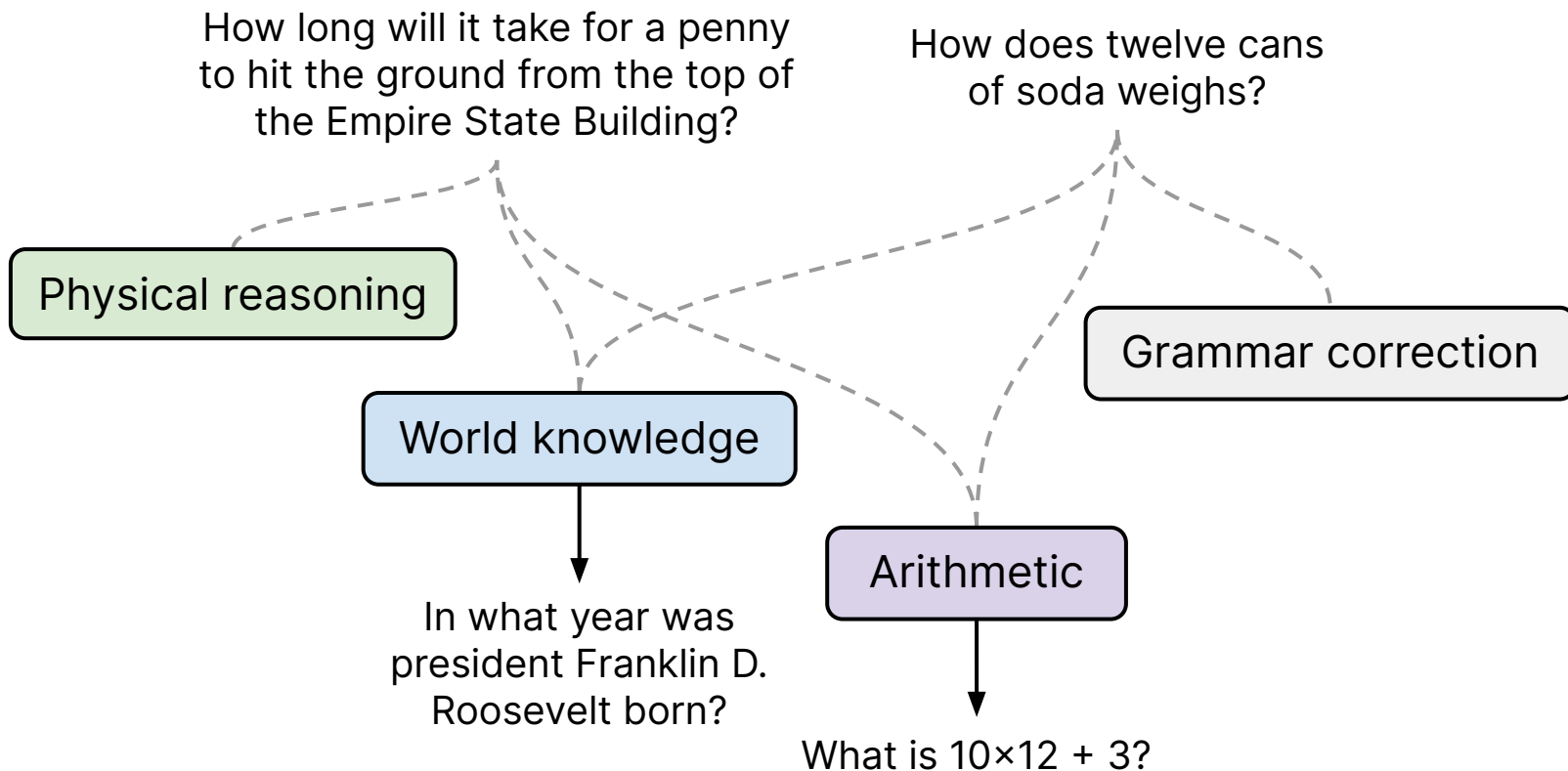


From "Soft Merging of Experts with Adaptive Routing" by Muqeeth et al.

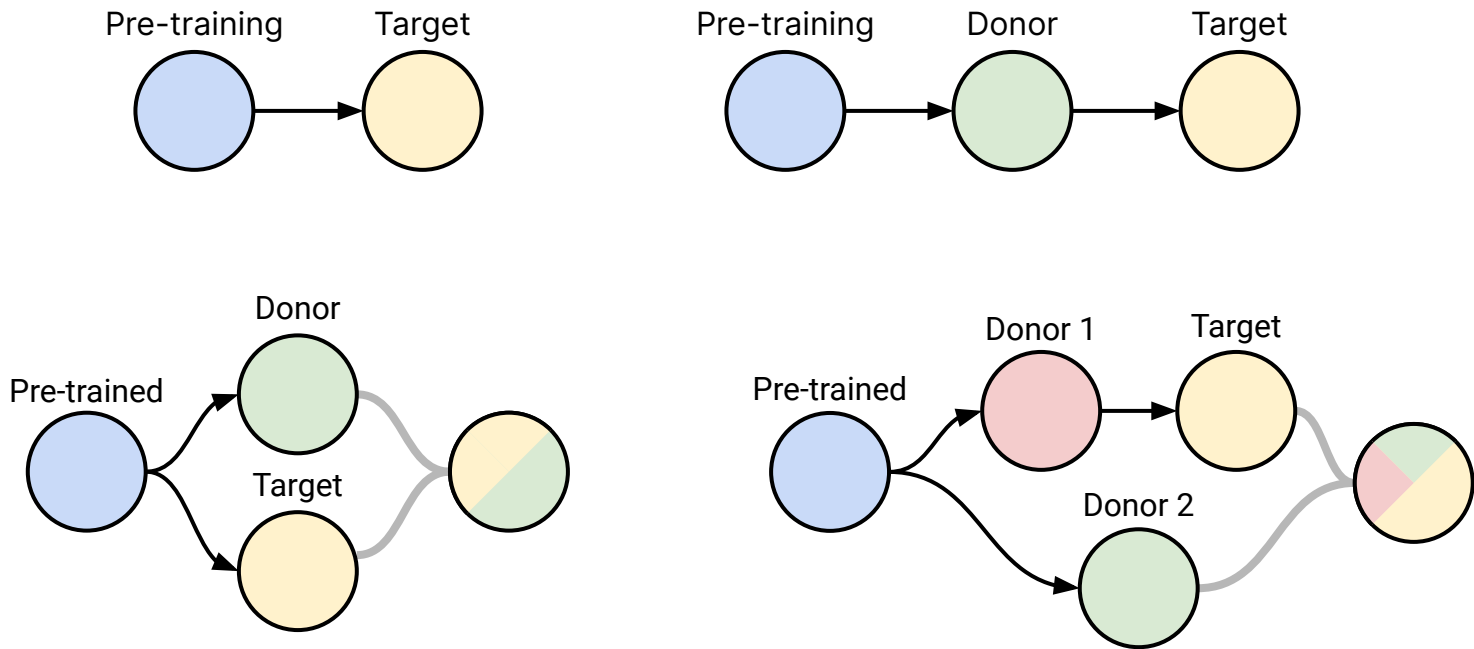
How and why should we build ecosystems of specialist models instead of monolithic models?

Capabilities can be **merged** across models.

Tasks can be considered as a composition of skills

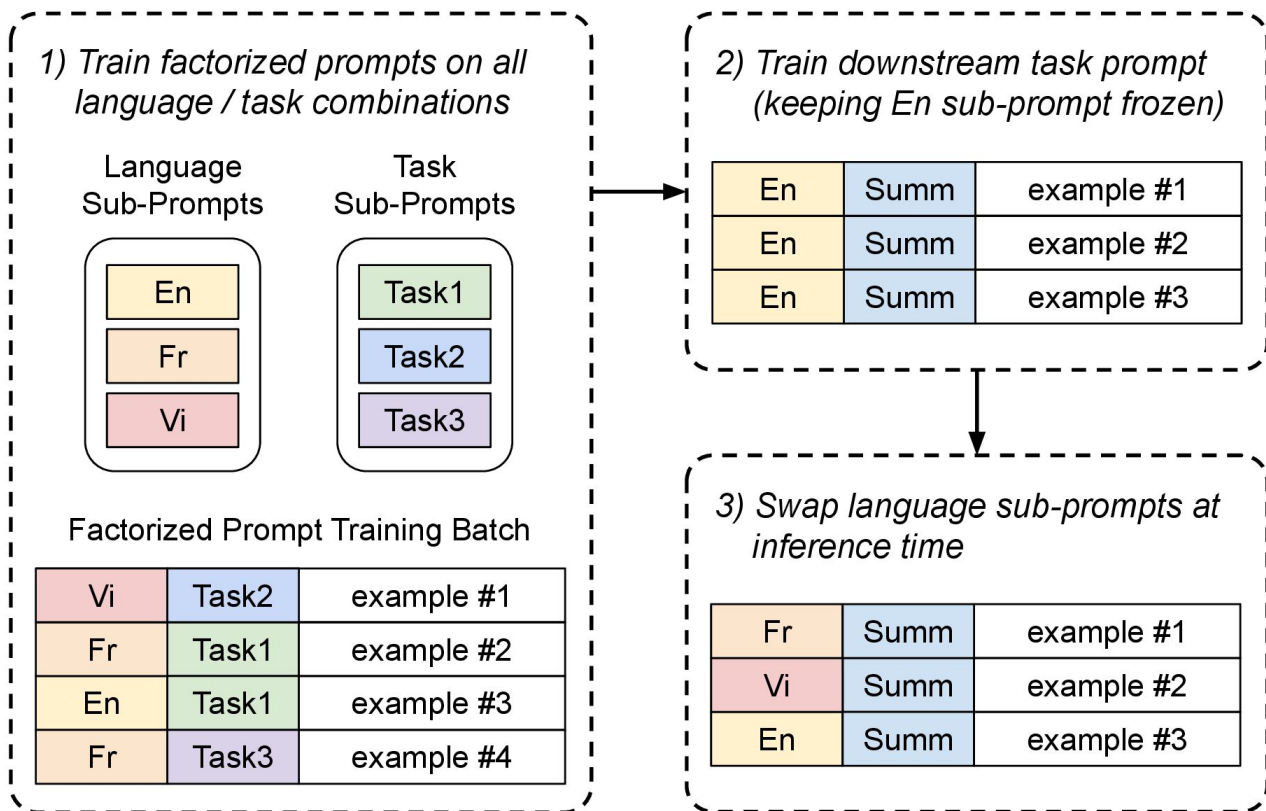


Merging models enables new paths for transferring capabilities



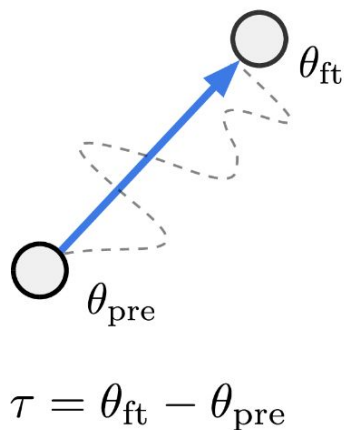
From "Merging Models with Fisher-Weighted Averaging" by Matena et al.

Learning compositional adapters via prompt tuning

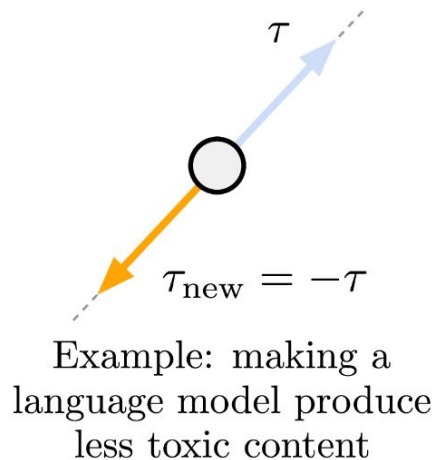


Editing models with task vectors

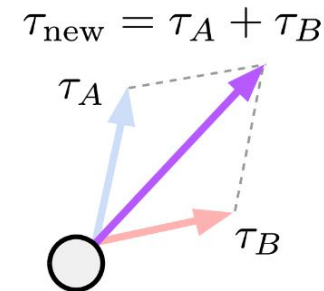
a) Task vectors



b) Forgetting via negation

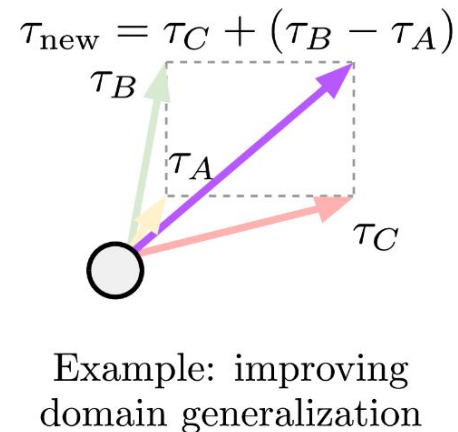


c) Learning via addition

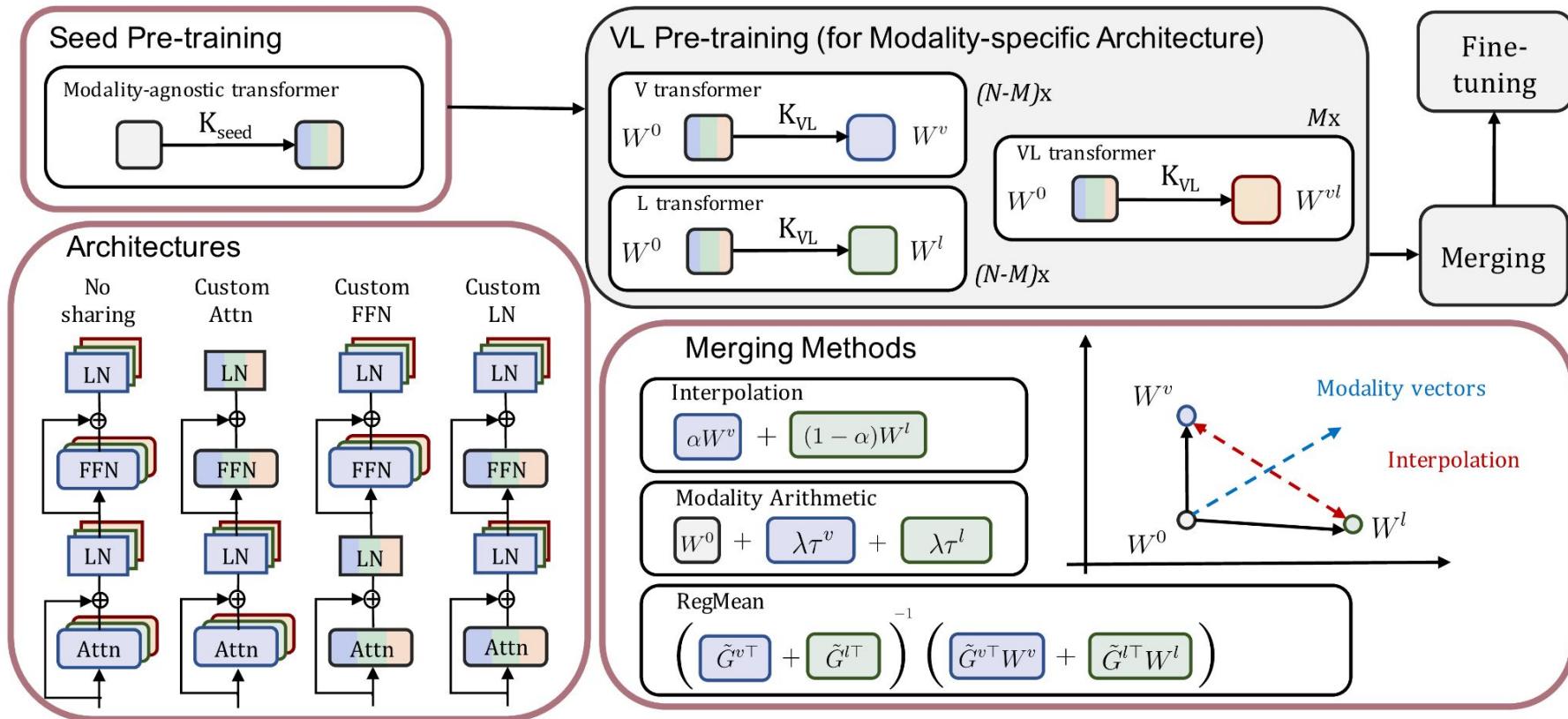


Example: building a multi-task model

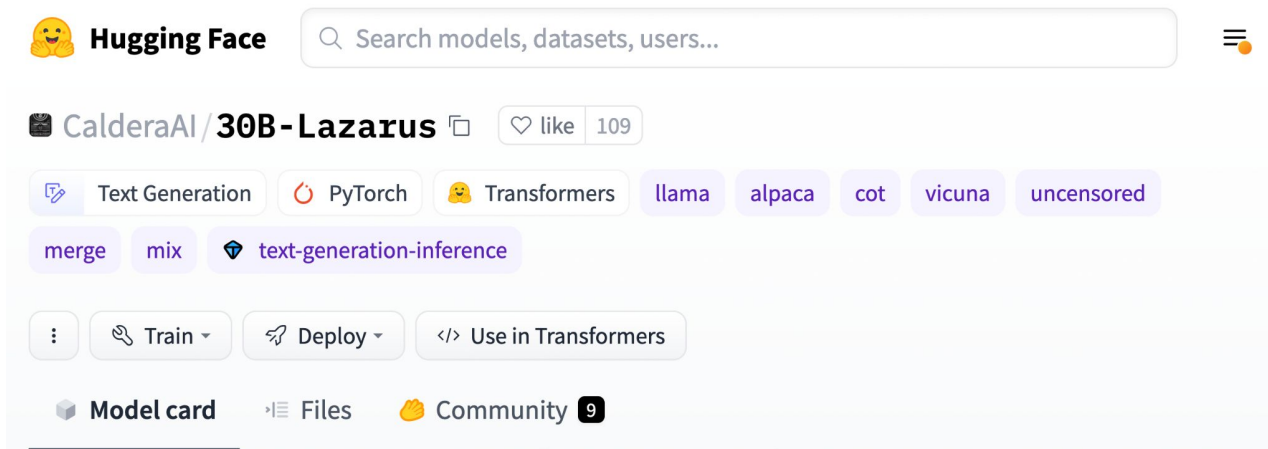
d) Task analogies



Merging can create multimodal models from unimodal models



Recent community-developed models are built via merging



30B-Lazarus

Composition:

[] = applied as LoRA to a composite model | () = combined as composite models

[SuperCOT([gtp4xalpaca(manticorechatpygalphavvicunaunlocked))]+[StoryV2(kaiokendev-SuperHOT-LoRA-prototype30b-8192))]]

From <https://huggingface.co/CalderaAI/30B-Lazarus>

Model merging as an optimization problem

$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \log p(\theta | \mathcal{D}_i)$$

$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \overbrace{\log p(\theta | \mathcal{D}_i)}^{\text{Log posterior for model } i}$$

\uparrow
*Hyperparameter
controlling the
importance of model i*

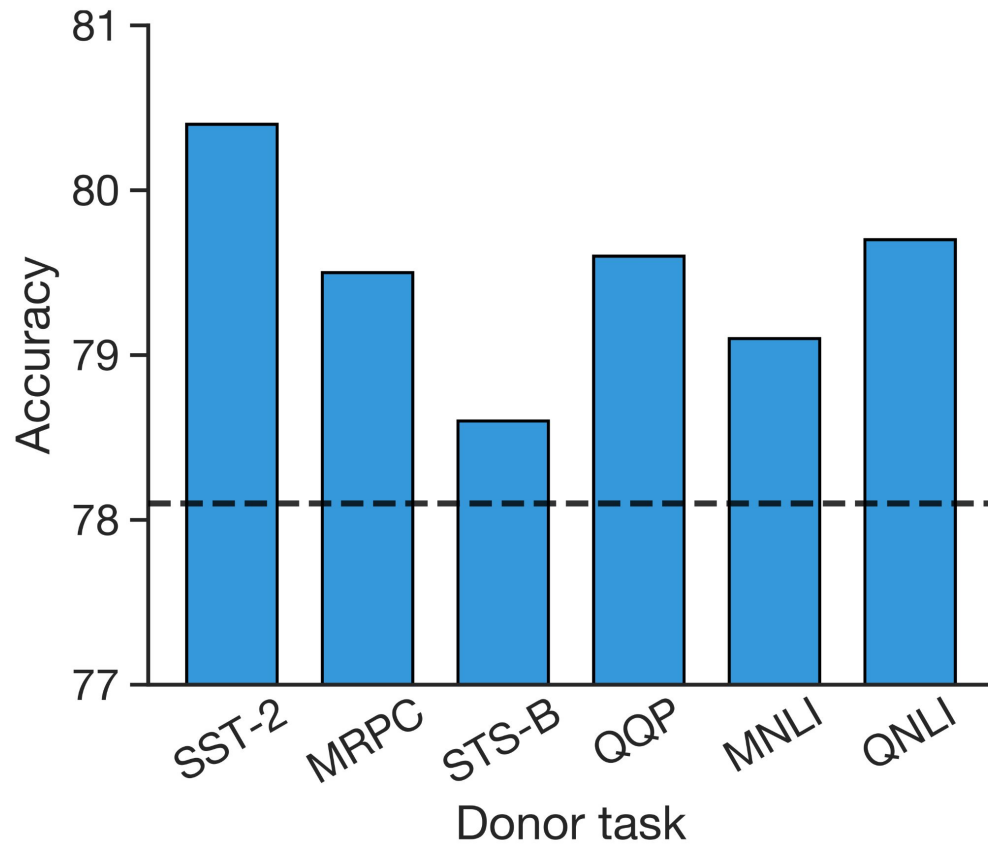
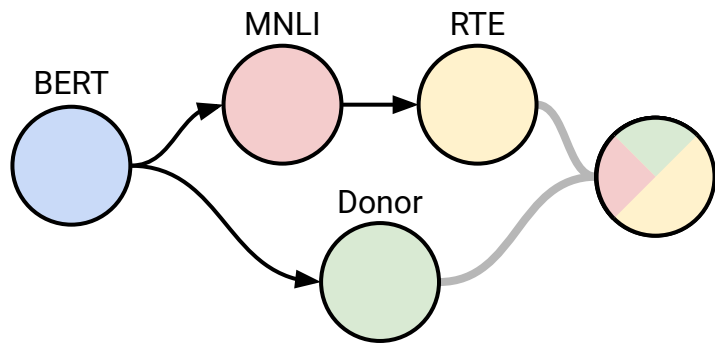
Fisher merging uses the Laplace approximation

$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \log p(\theta | \mathcal{D}_i)$$

$\downarrow \theta \sim \mathcal{N}(\theta_i, \hat{F}_i^{-1})$

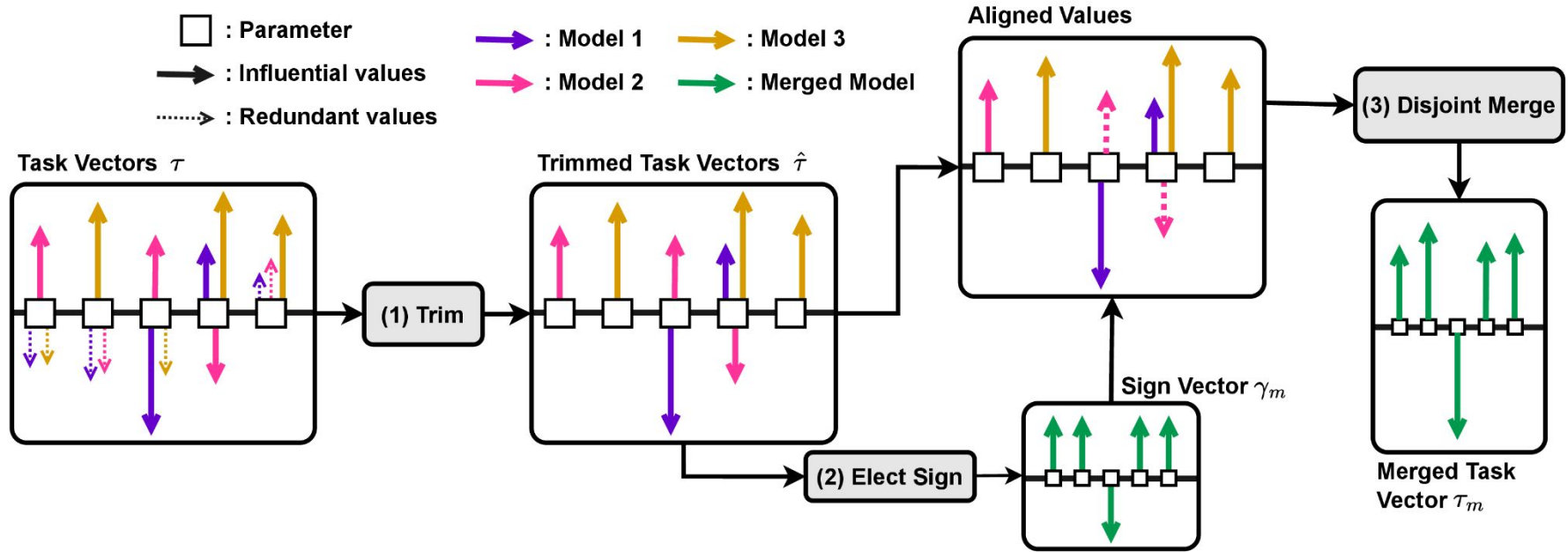
$$\theta^* = \frac{\sum_i \lambda_i \hat{F}_i \odot \theta_i}{\sum_i \lambda_i \hat{F}_i}$$

Fisher merging can combine the capabilities of different models

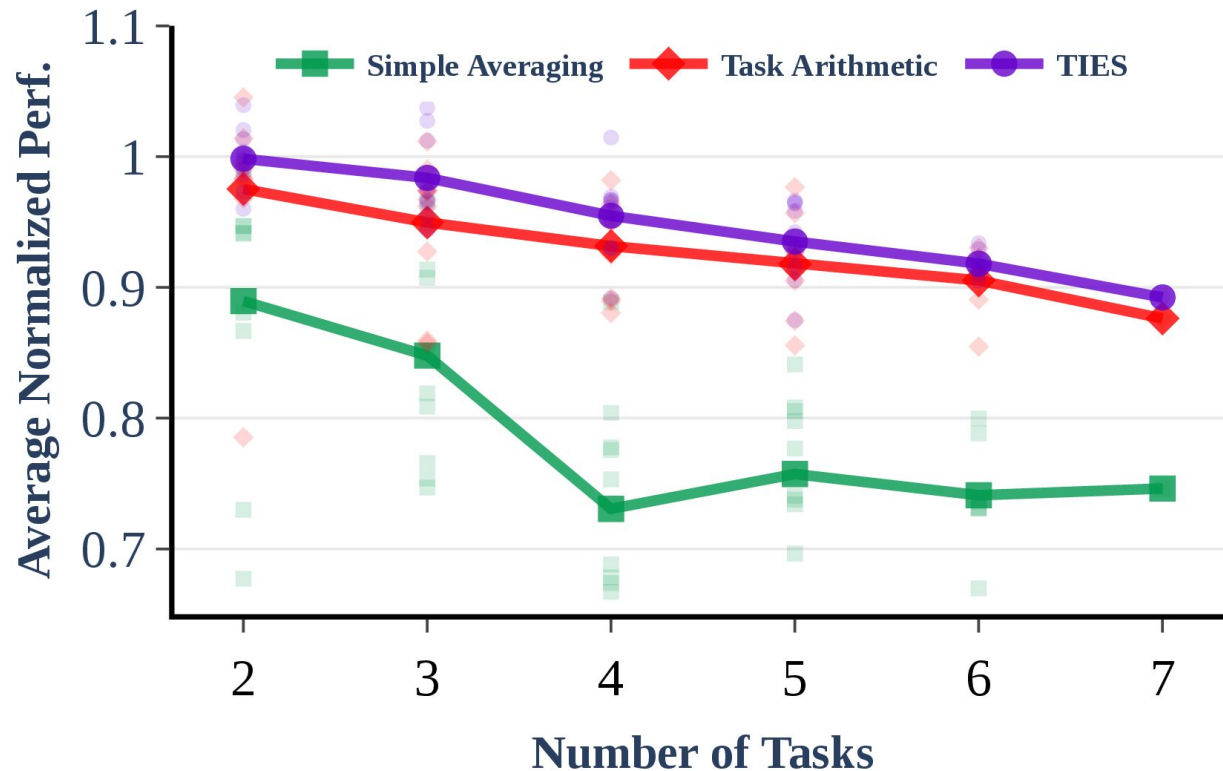


From "Merging Models with Fisher-Weighted Averaging" by Matena et al.

TIES Merging resolves interference when merging models

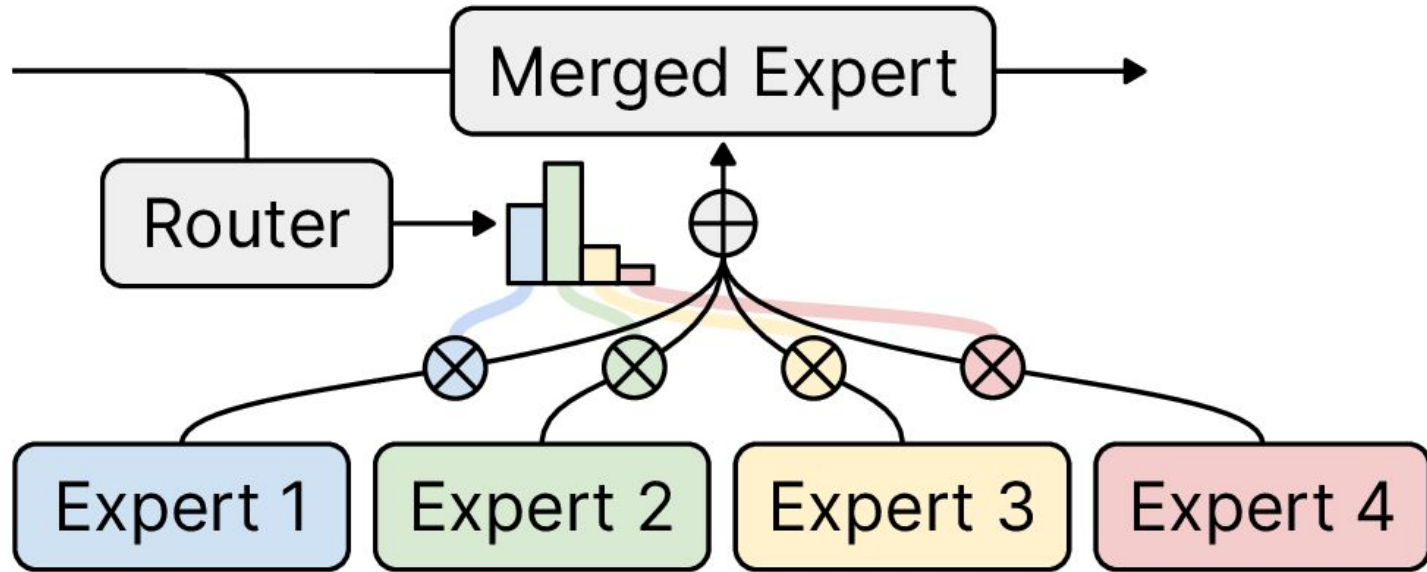


TIES helps retain specialist model performance



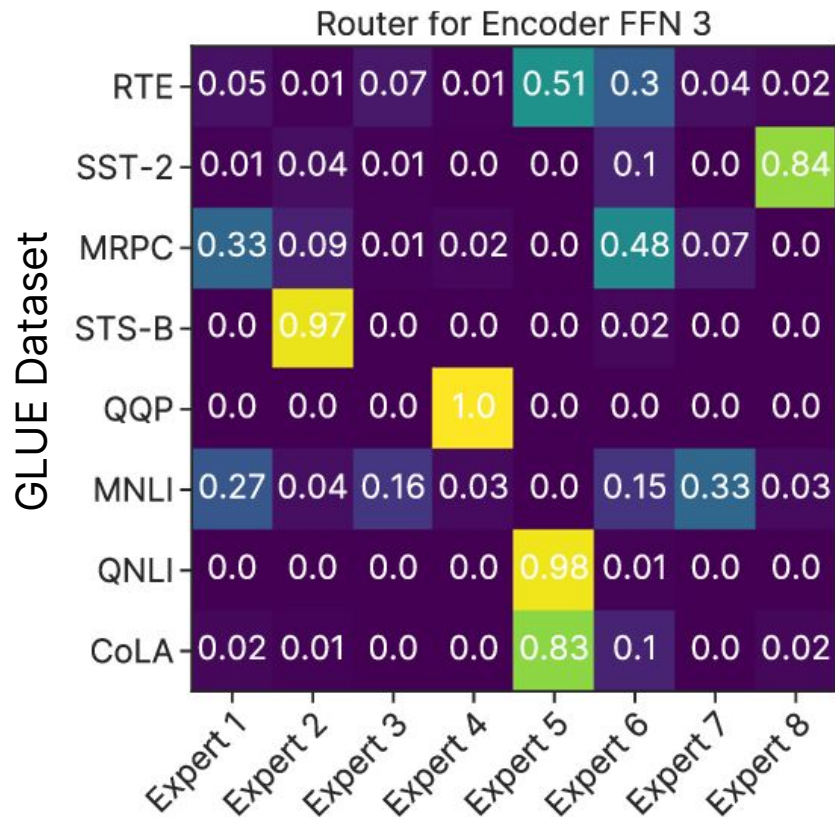
From "Resolving Interference When Merging Models" by Yadav et al.

Differentiable routing between specialist submodels with SMEAR



From "Soft Merging of Experts with Adaptive Routing" by Muqeeth et al.

Experts specialize and are shared across different tasks



How and why should we build ecosystems of specialist models instead of monolithic models?

An ecosystem can be built and used collaboratively with the right **systems**.

git-theta tracks, merges, and updates models using the git workflow

```
$ git-theta track model.pt
$ git commit -am "Add initial model"
$ python finetune.py --dataset="cb" --method="lowrank"
$ git commit -am "Fine-tune on CB dataset with LoRA"
$ git checkout -b rte
$ python finetune.py --dataset="rte" --method="dense"
$ git commit -am "Fine-tune on RTE dataset"
$ git checkout main
$ python finetune.py --dataset="anli" --method="dense"
$ git commit -am "Fine-tune on ANLI dataset"
$ git merge rte
```

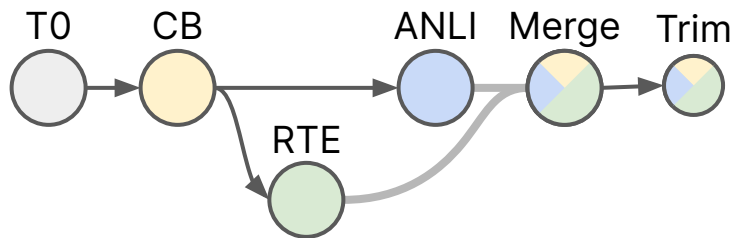
Fixing Merge Conflicts in model.pt

Actions:

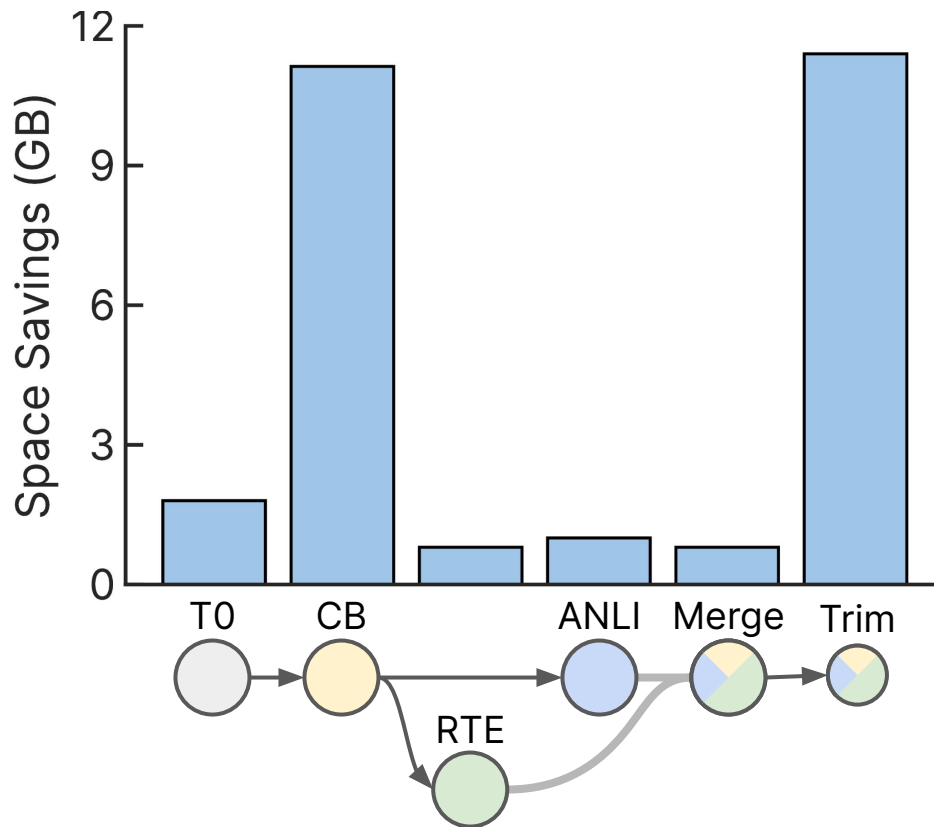
- avg) average: Average parameter values.
- tt) take_them: Use their change to the parameter.
- tu) take_us: Use our change to the parameter.
- q) quit

θ avg

```
$ git commit -am "Merge RTE and ANLI models"
$ python trim_unused_embeddings.py
$ git commit -am "Remove embeddings for unused tokens"
```

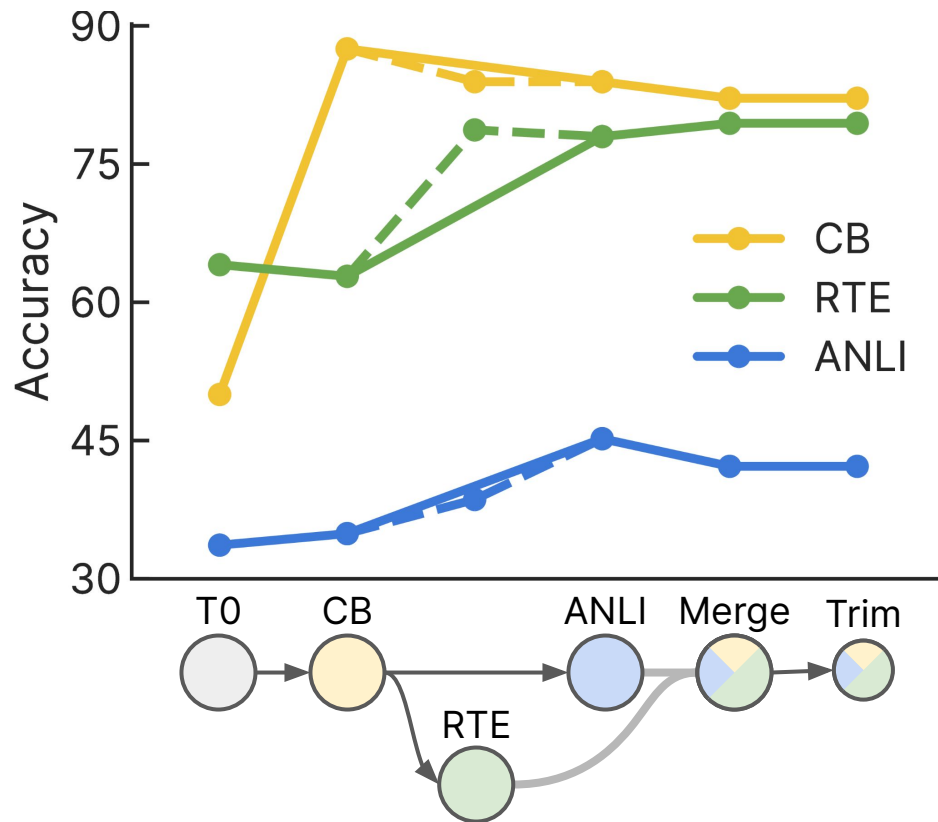


Communication-efficient updates result in significant space savings



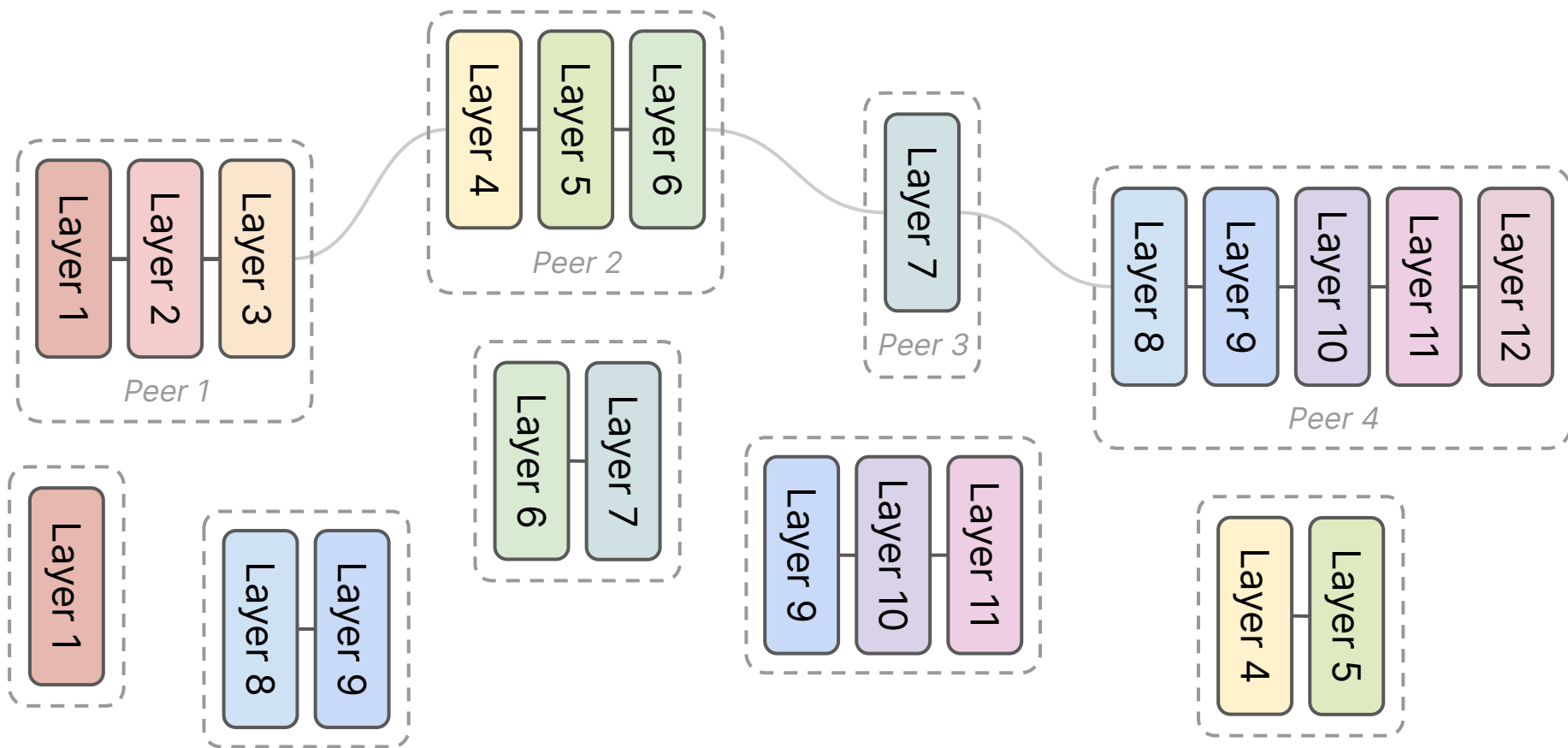
From "Git-Theta: A Git Extension for Collaborative Development of Machine Learning Models" by Kandpal et al.

git-theta allows for continuous and collaborative model development



From "Git-Theta: A Git Extension for Collaborative Development of Machine Learning Models" by Kandpal et al.

Petals enables distributed inference and fine-tuning over the internet



From "Petals: Collaborative Inference and Fine-tuning of Large Models" by Borzunov et al.

Current Petals swarm status




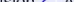


Petals Health Monitor

Bootstrap peers: ●●

Model [stabilityai/StableBeluga2](#) (**healthy**):

Server ID »	Contributor ?	Version	Throughput »	Precision ?	Adapters ?	Cache ?	Avl. ?	Pings ?	Served blocks
...vYA3Rn		2.0.1.post2	1333 tok/s	bf16 (nf4)		32768	Direct	Show	0:40 <div><div></div></div>
...TfceK7		2.0.1.post2	1333 tok/s	bf16 (nf4)		32768	Direct	Show	40:80 <div><div></div></div>
...uF3WXf	👉 FYY 😊	2.0.1.post2	1015 tok/s	bf16 (nf4)		24576	Direct	Show	57:76 <div><div></div></div>
...sQzHZf	👉 FYY 😊	2.0.1.post1	905 tok/s	f16 (nf4)		24576	Direct	Show	0:19 <div><div></div></div>
...yRjqCy	Zetta	2.0.1.post2	1001 tok/s	bf16 (nf4)		30720	Relay	Show	19:38 <div><div></div></div>
...m3Vfh7	👉 FYY 😊	2.0.1.post1	324 tok/s	bf16 (nf4)		29044	Direct	Show	69:80 <div><div></div></div>
...5D4AAQ		2.0.1.post1	100 tok/s	bf16 (nf4)		32768	Relay	Show	75:80 <div><div></div></div>
...dSkeys	👉 FYY 😊	2.0.1.post2	1015 tok/s	bf16 (nf4)		22528	Direct	Show	38:57 <div><div></div></div>

Model [meta-llama/Llama-2-70b-chat-hf](#) (healthy):

Server ID »	Contributor ?	Version	Throughput »	Precision ?	Adapters ?	Cache ?	Avl. ?	Pings ?	Served blocks
...kDJVjh	jobs.trelent.com	2.0.1	8107 tok/s	f16 (nf4)		32768	Relay	Show	20:46 
...bsDGGc		2.0.1.post2	17 tok/s	f16 (nf4)		32768	Relay	Show	0:3 
...MWAxrr	nora	2.0.1.post1	670 tok/s	f16 (nf4)		12288	Direct	Show	46:66 
...rgNAo9	nora	2.0.1.post1	670 tok/s	f16 (nf4)		12288	Direct	Show	60:80 
...RPFSet	nora	2.0.1.post1	670 tok/s	f16 (nf4)		12288	Direct	Show	0:20 

From <https://health.petals.dev/>

Thanks.

Please give me feedback:

<http://bit.ly/colin-talk-feedback>

craffel@gmail.com