



The many faces of heterogeneity: Federated, Continual, & Modular Learning

Marco Ciccone

Distinguished Postdoctoral Fellow



Simons Workshop on Learning from Heterogeneous Sources

Works done in collaboration with

Debora



Eros



Riccardo



Leonardo



and Colin, Tatiana, Barbara, Massi and Carlo



Disclaimer

I am not an optimization researcher.

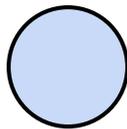
My background is in deep learning.

*My research focuses on **transfer learning** across data sources and tasks, as well as collaborative and modular learning.*

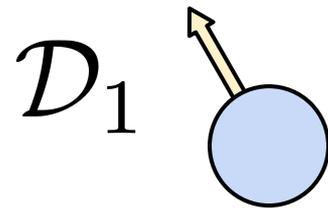
Goal for today

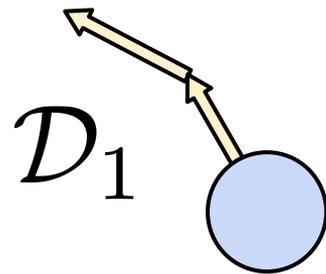
Show how different practical settings share the same underlying problem:

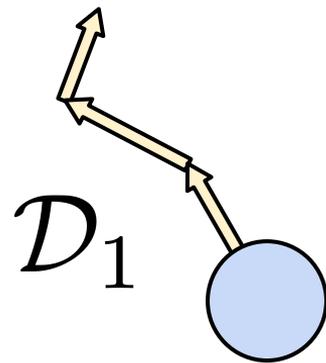
*Learning from heterogeneous sources,
under different constraints.*

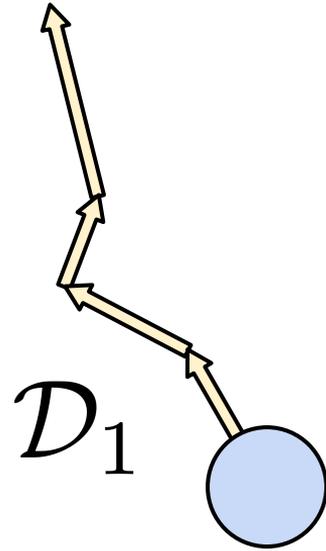


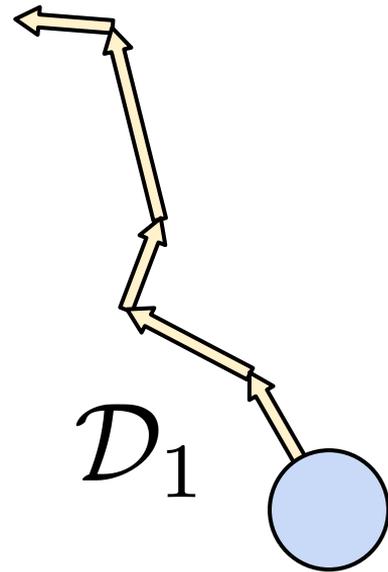
Pre-trained model / Initialization

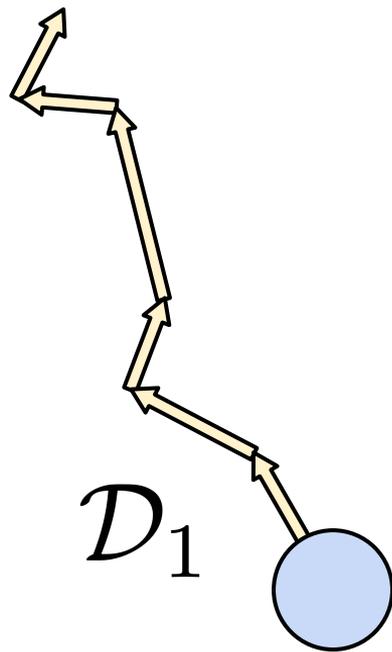




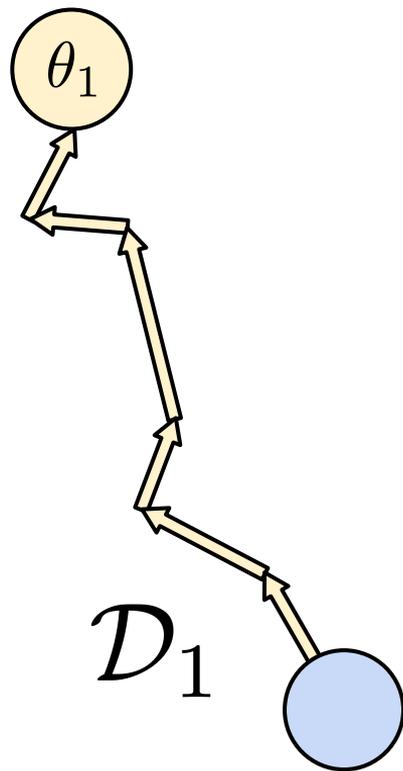




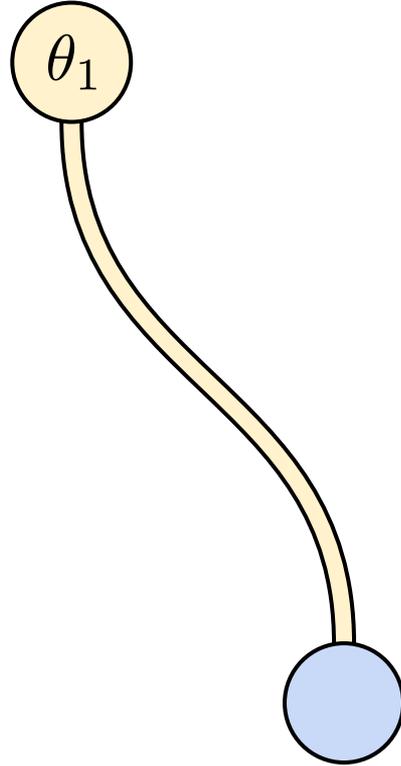




Single-source model



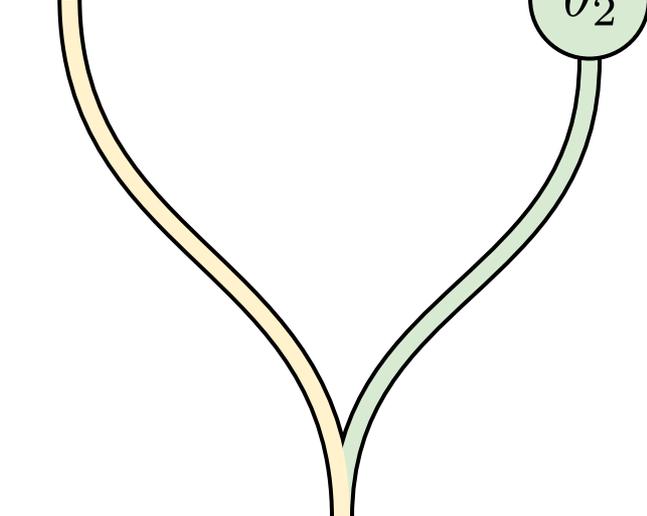
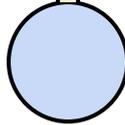
Single source



Single source



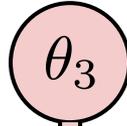
Single source



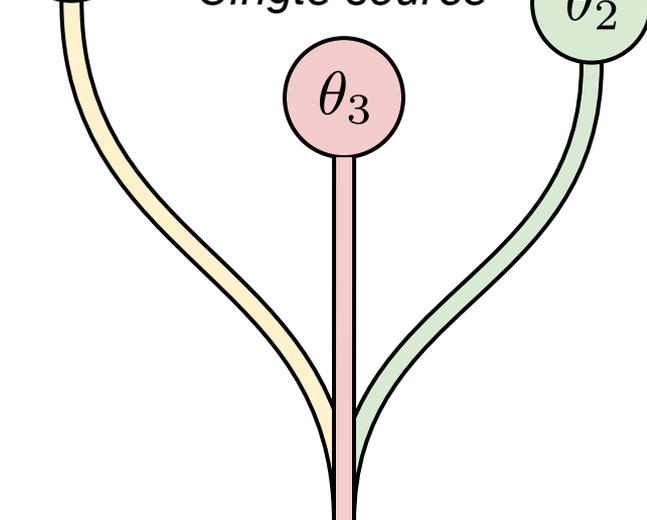
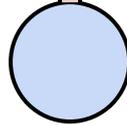
Single source



Single source



Single source



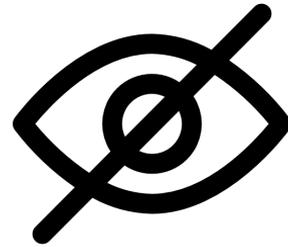
*Multi-task or multi-domain
(simultaneous access to all sources)*



Challenge: *balance across sources
find optimal data mixture, exploring pareto front is expensive*

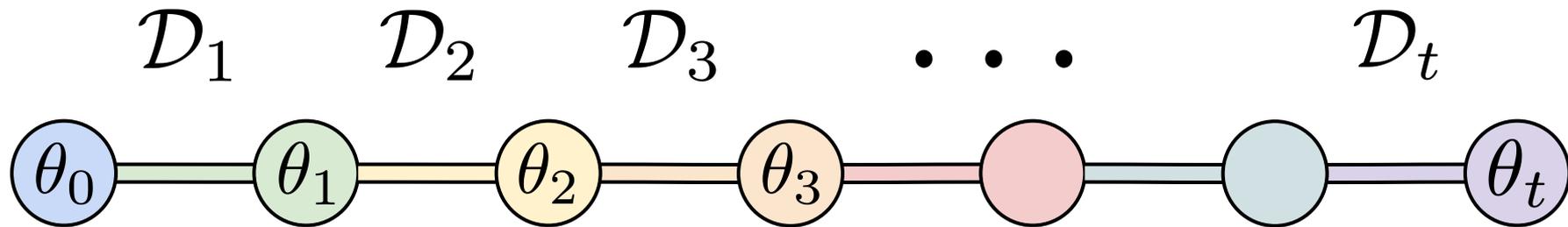
We cannot always access all sources simultaneously.

*Constraints on data access give different settings,
but create similar challenges...*





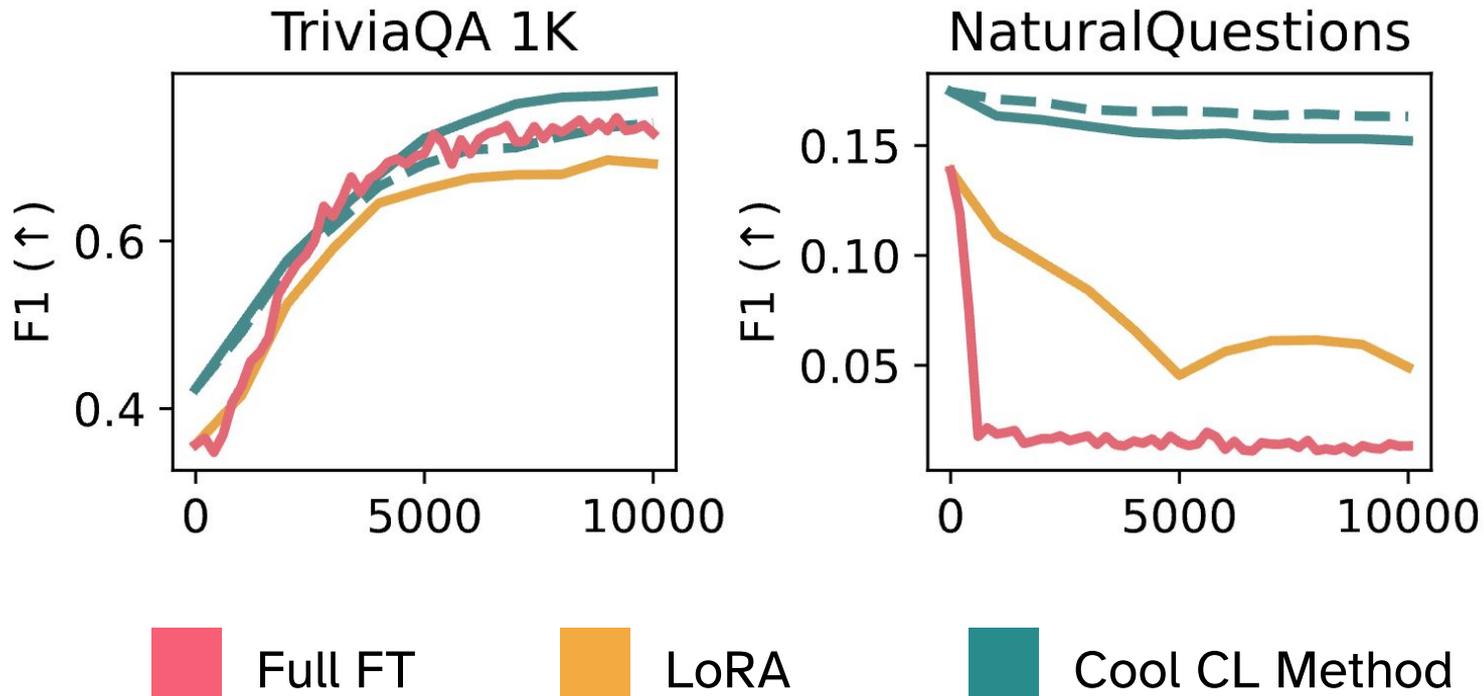
Continual Learning: *Heterogeneity across time*



Constraint: *no access to past data!*

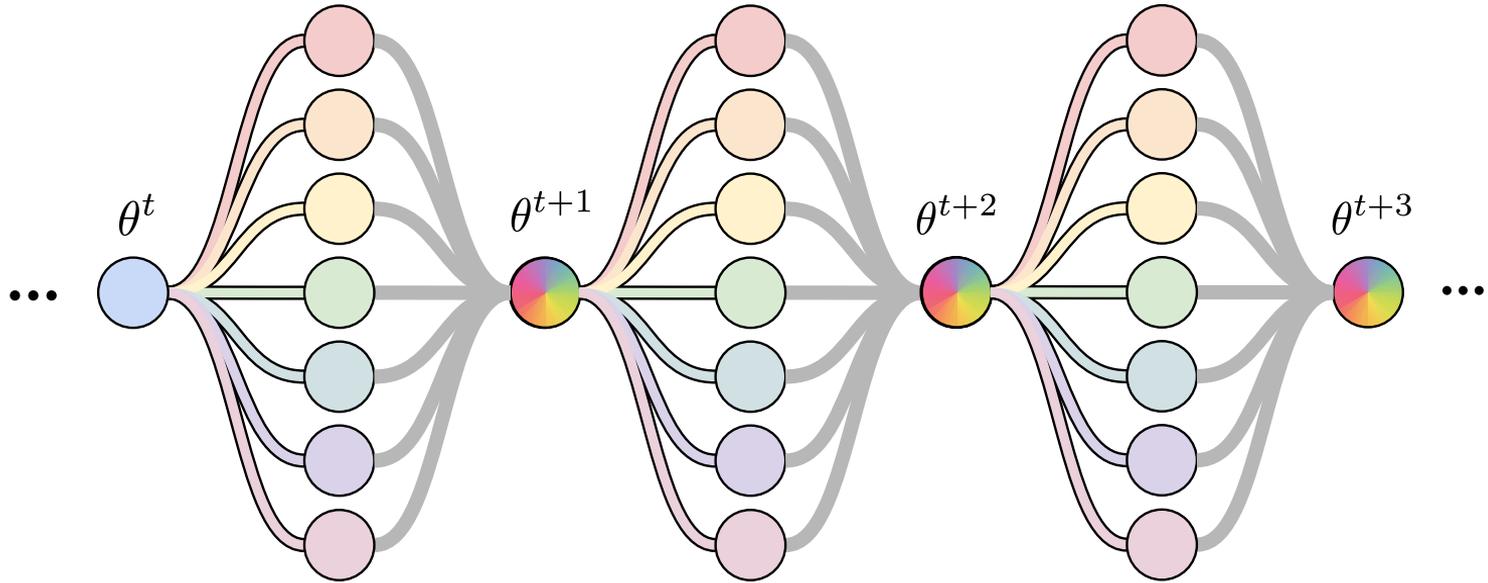
Sequential training, only on **current data/task**

*Key issue: **Catastrophic Forgetting**
the model becomes **biased** towards recent data*





Federated Learning: Heterogeneity across space

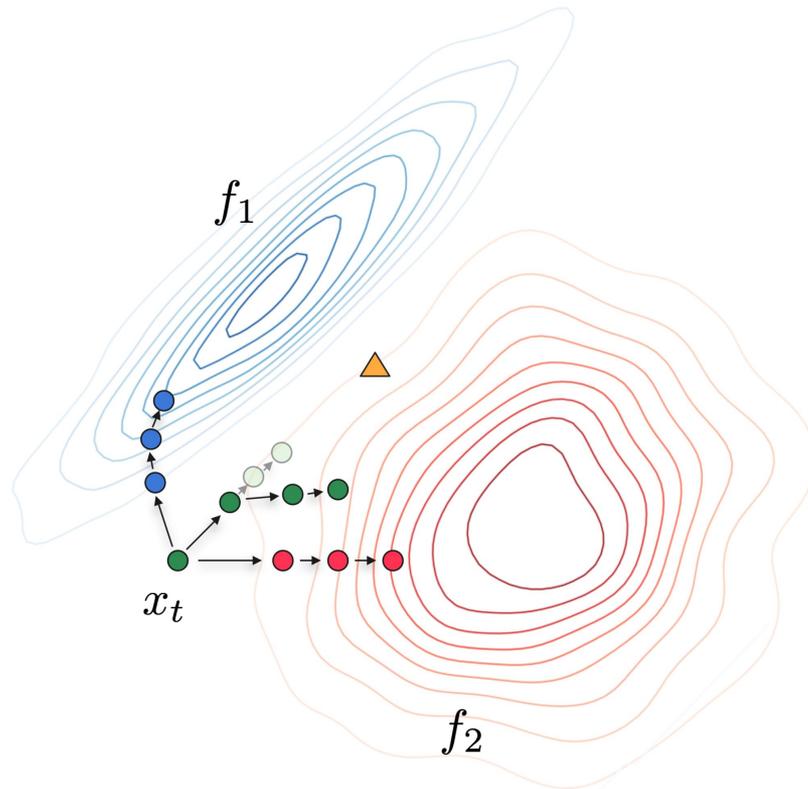


Constraint: no access to full data, communication efficiency

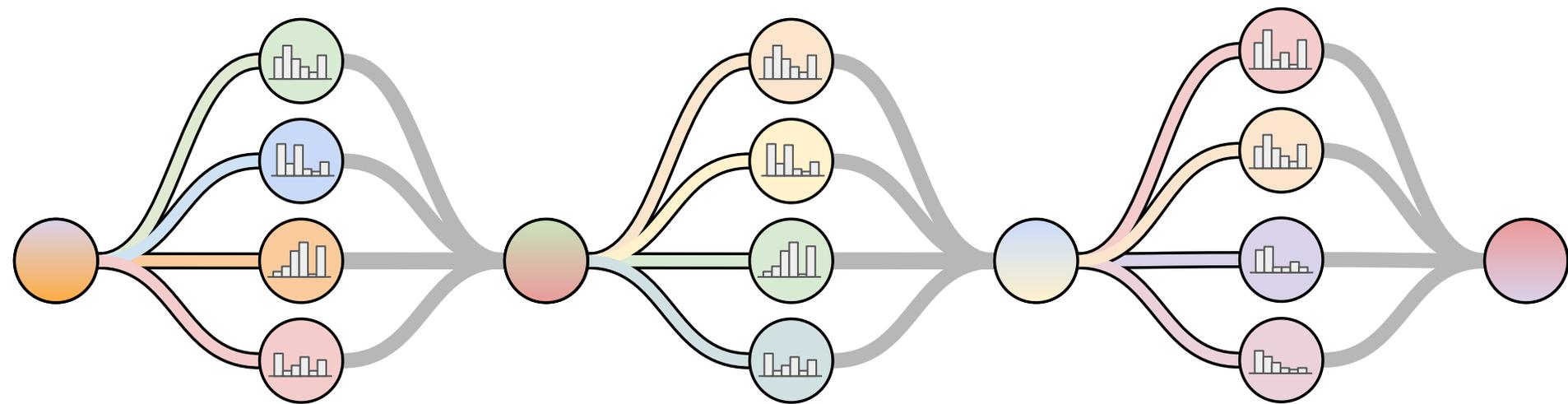
Parallel training only on private clients' data and **model aggregation**

Client Drift Problem

Local models **overspecialize** on clients' data diverging from the global optimization goal



FL with Partial Participation *(spatial and temporal heterogeneity)*



*This induces a **sequential learning** problem.*

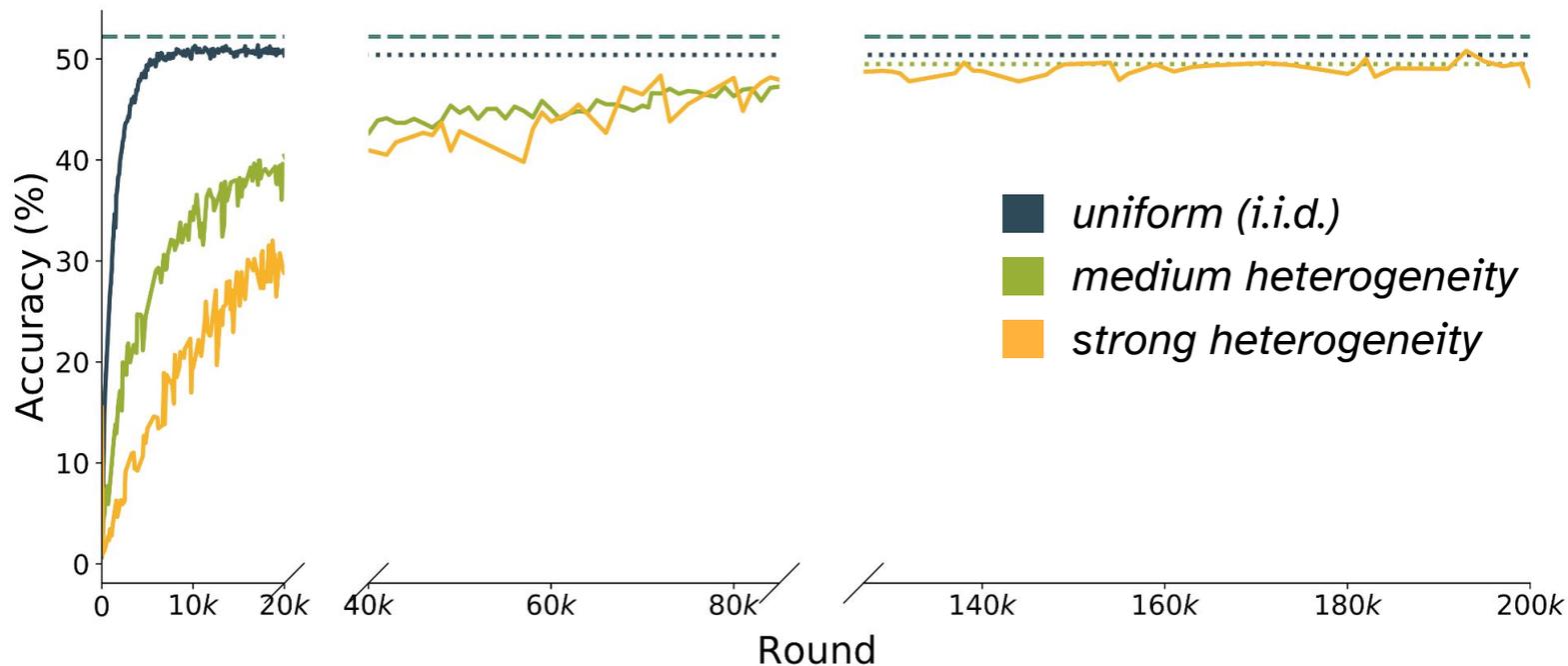
*At each round the global model is biased towards the most recent clients,
forgetting past knowledge.*

*FL with **one client per round** (and no client revisit)
reduces to **Continual Learning***

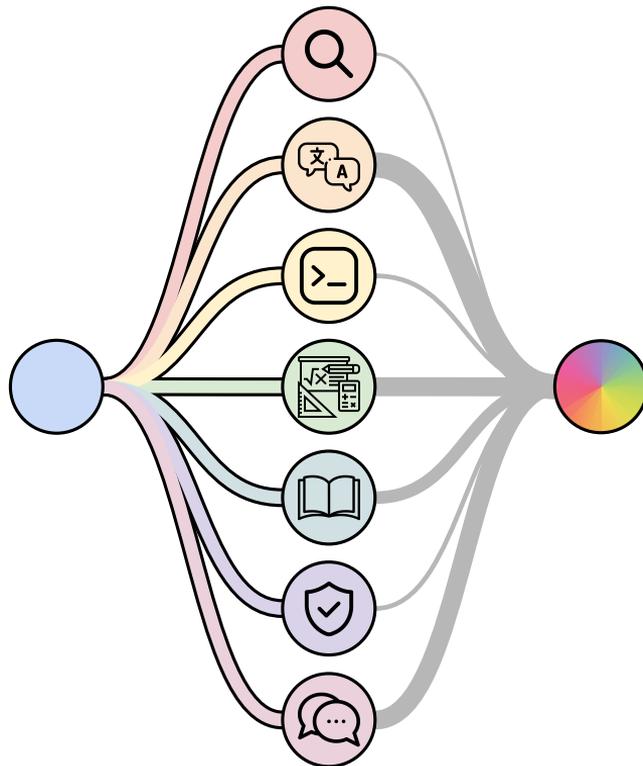


*This induces a **sequential learning** problem.
At each round the global model is biased towards the most recent clients,
forgetting past knowledge.*

Key issue: *slower and noisier convergence*
(requires more communication)



 **Model Merging** combines expert models



Parallel training and **single** model aggregation

Two core operations, two failure modes

“Local” Gradient Descent

*optimizes on individual sources,
minimizing partial objectives.*

Model Aggregation

*combines independently trained
models in the parameter space.*

Two core operations, two failure modes

“Local” Gradient Descent

*optimizes on individual sources,
minimizing partial objectives.*



Forgetting past knowledge,
Bias towards local data,
Drift from global solution.

Model Aggregation

*combines independently trained
models in the parameter space.*



Interference from
conflicting parameters
destroys knowledge.

Two core operations, two failure modes

“Local” Gradient Descent

*optimizes on individual sources,
minimizing partial objectives.*



Forgetting *past knowledge,*
Bias *towards local data,*
Drift *from global solution.*

Model Aggregation

*combines independently trained
models in the parameter space.*



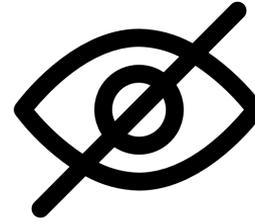
Interference *from*
conflicting parameters
destroys knowledge.



*These compound, and different paradigms
suffer from different combinations of problems*



Similar problems...



...similar solutions.



Example: local training with proximal terms

Continual Learning: Mitigate forgetting across tasks

$$\mathcal{L}_{EW C}(\theta) = \mathcal{L}_{new}(\theta) + \frac{\lambda}{2} (\theta - \theta_{old})^T \overset{\text{Fisher Information Matrix}}{F} (\theta - \theta_{old})$$

Federated Learning: Mitigate client drift from global model

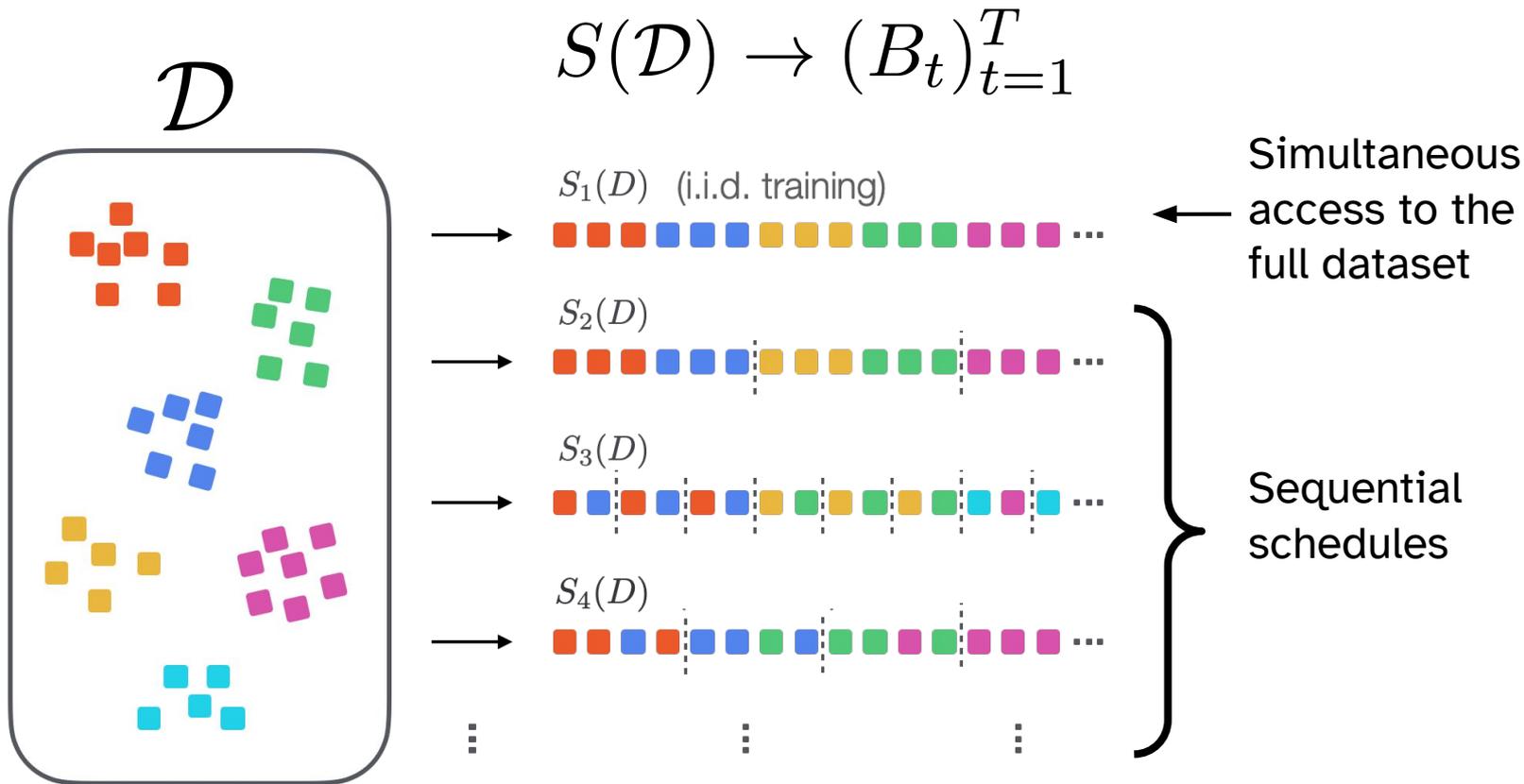
$$\mathcal{L}_{FedProx}(\theta) = \mathcal{L}_{local}(\theta) + \frac{\mu}{2} (\theta - \theta_{global})^T I (\theta - \theta_{global})$$

From "Overcoming catastrophic forgetting in neural networks" by Kirkpatrick et al

From "Federated Optimization in Heterogeneous Networks" by Li et al

*Let's see a simple method
that can transfer across settings.*

In the CL setting, a dataset can be splitted into different “**schedules**” deciding how it is **streamed** into chunks of data



GOAL: Find an algorithm that is robust (or invariant) to different schedules.
We can restrict ourselves to a **specific class of “one vs all” predictors**

$$W^* = \arg \min_W \frac{1}{|D|} \sum_{(x,y) \in D} \left\| \underbrace{W^\top \psi(x)}_{\text{pretrained model}} - \text{OneHot}(y) \right\|^2 + \lambda \|W\|^2.$$

GOAL: Find an algorithm that is robust (or invariant) to different schedules.
We can restrict ourselves to a **specific class of “one vs all” predictors**

$$W^* = \arg \min_W \frac{1}{|D|} \sum_{(x,y) \in D} \underbrace{\|W^\top \psi(x) - \text{OneHot}(y)\|^2}_{\text{pretrained model}} + \lambda \|W\|^2.$$

Closed-form solution

$$W^* = \underbrace{(X^\top X + \lambda I)^{-1}}_A \underbrace{X^\top Y}_b$$

$$X = [\psi(x_1) \dots \psi(x_N)]^\top$$

$$Y = [\text{OneHot}(y_1) \dots \text{OneHot}(y_N)]^\top$$

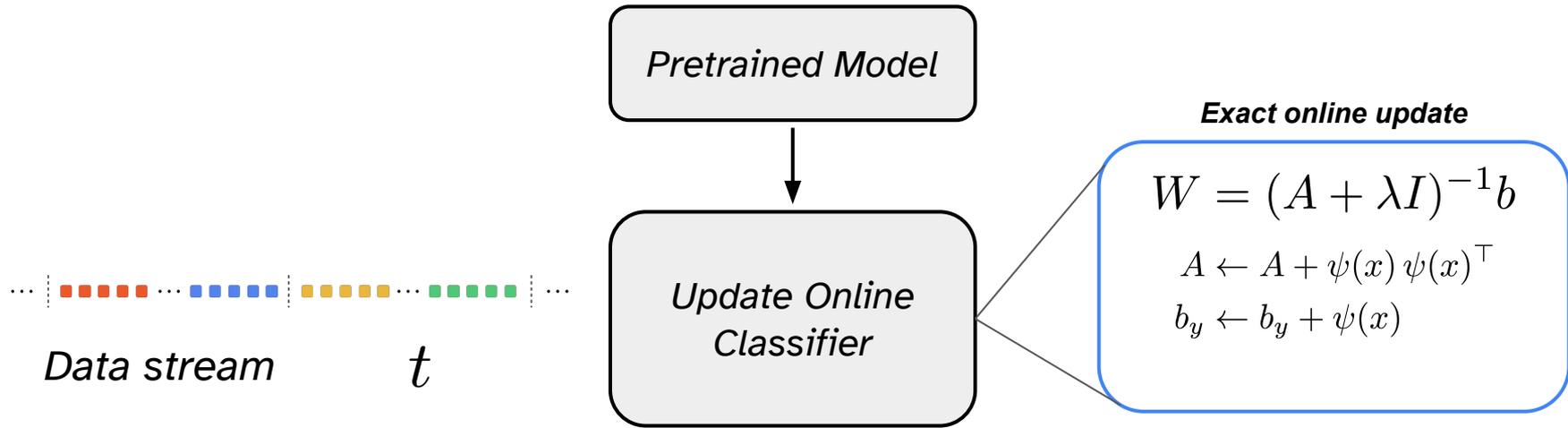
Exact online update

$$W = (A + \lambda I)^{-1} b$$

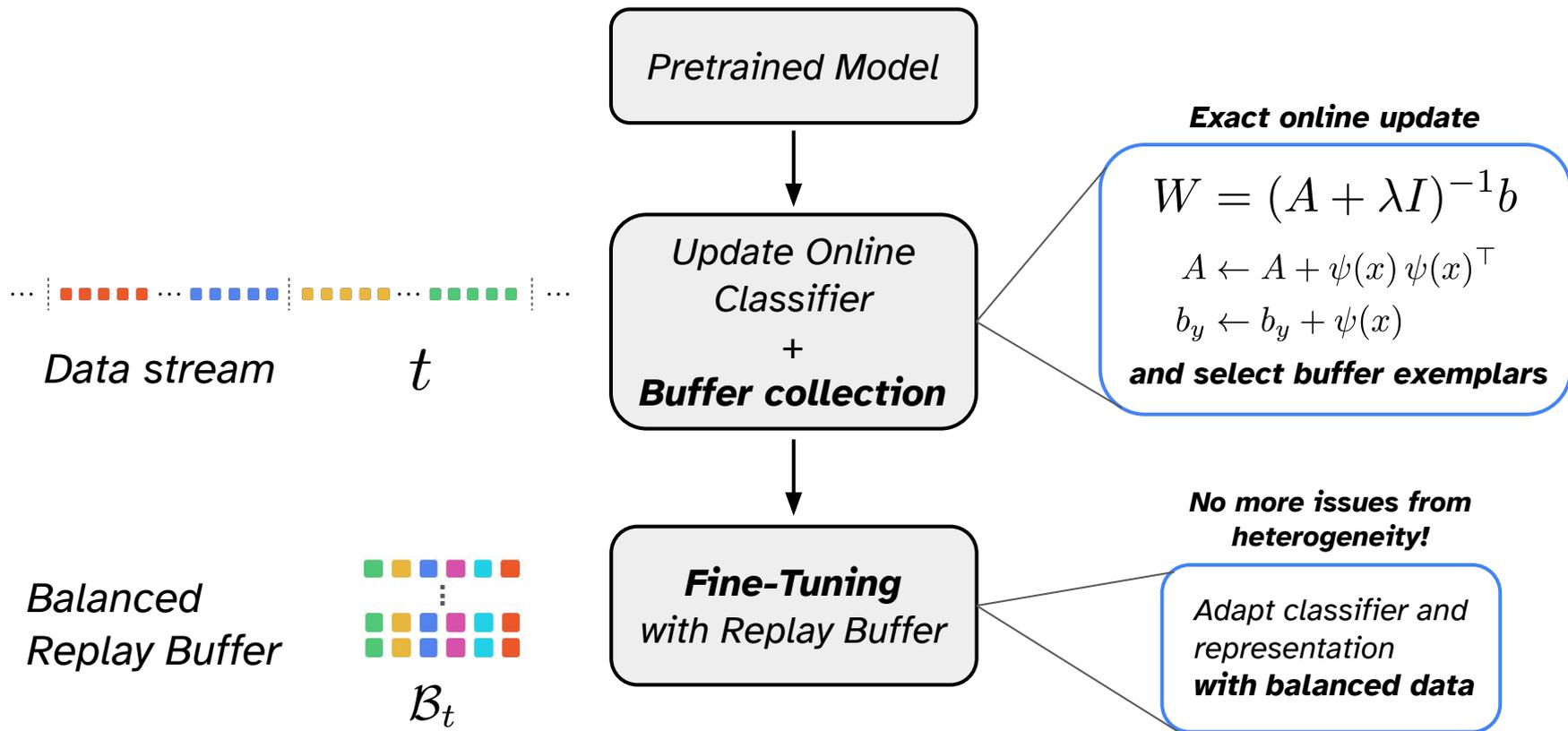
$$A \leftarrow A + \psi(x) \psi(x)^\top$$

$$b_y \leftarrow b_y + \psi(x)$$

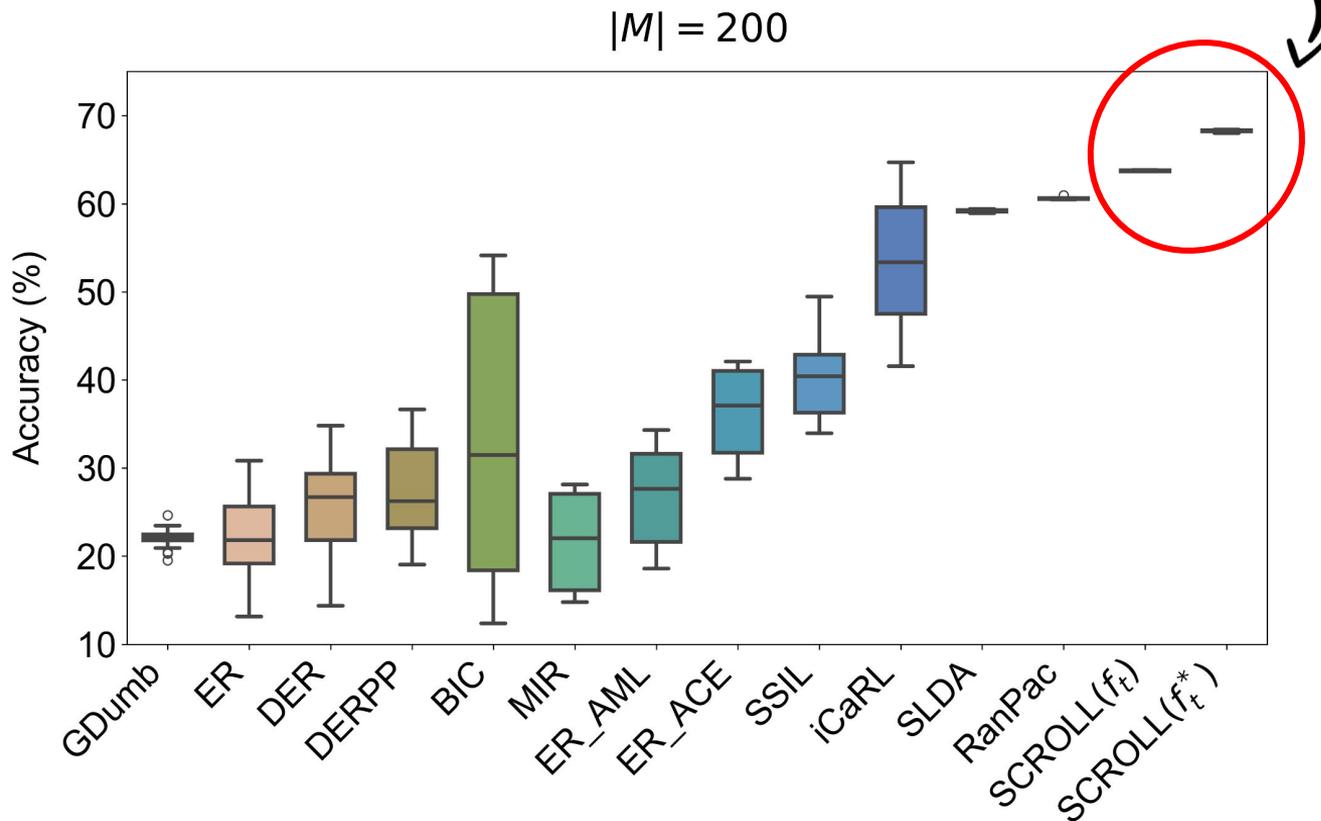
SCROLL: Schedule-Robust Continual Learning



SCROLL: Schedule-Robust Continual Learning



SCROLL is effectively robust across different schedules
(exact classifier updates + balanced fine-tuning)



*The same framework transfers to
Federated Learning*

Schedules as “client splits”



Federated Recursive Ridge Regression (**Fed3R**)

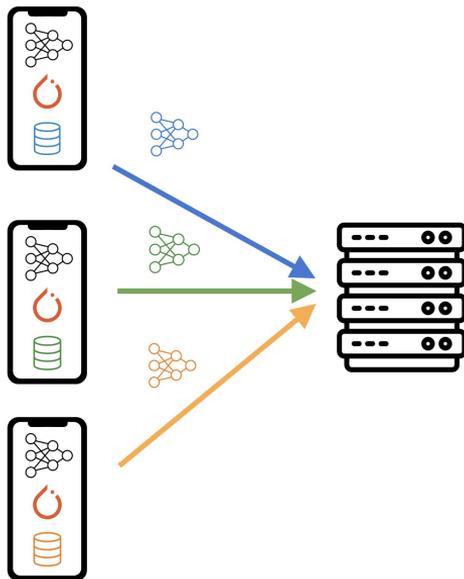
builds a classifier **exactly** and **incrementally** across clients

Compute Local Statistics

Client side

$$A_k^t = \sum_{(x,y) \in \mathcal{D}_k} \psi(x)^\top \psi(x)$$

$$b_k^t = \sum_{(x,y) \in \mathcal{D}_k} \psi(x)^\top e_y$$



Aggregate Statistics

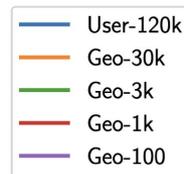
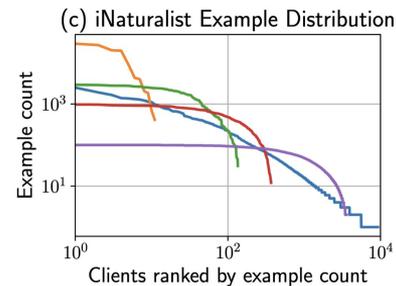
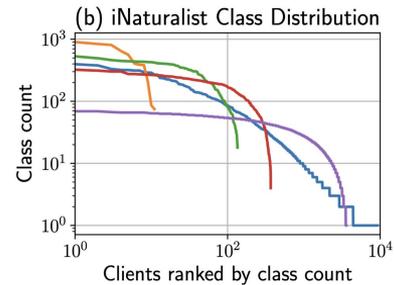
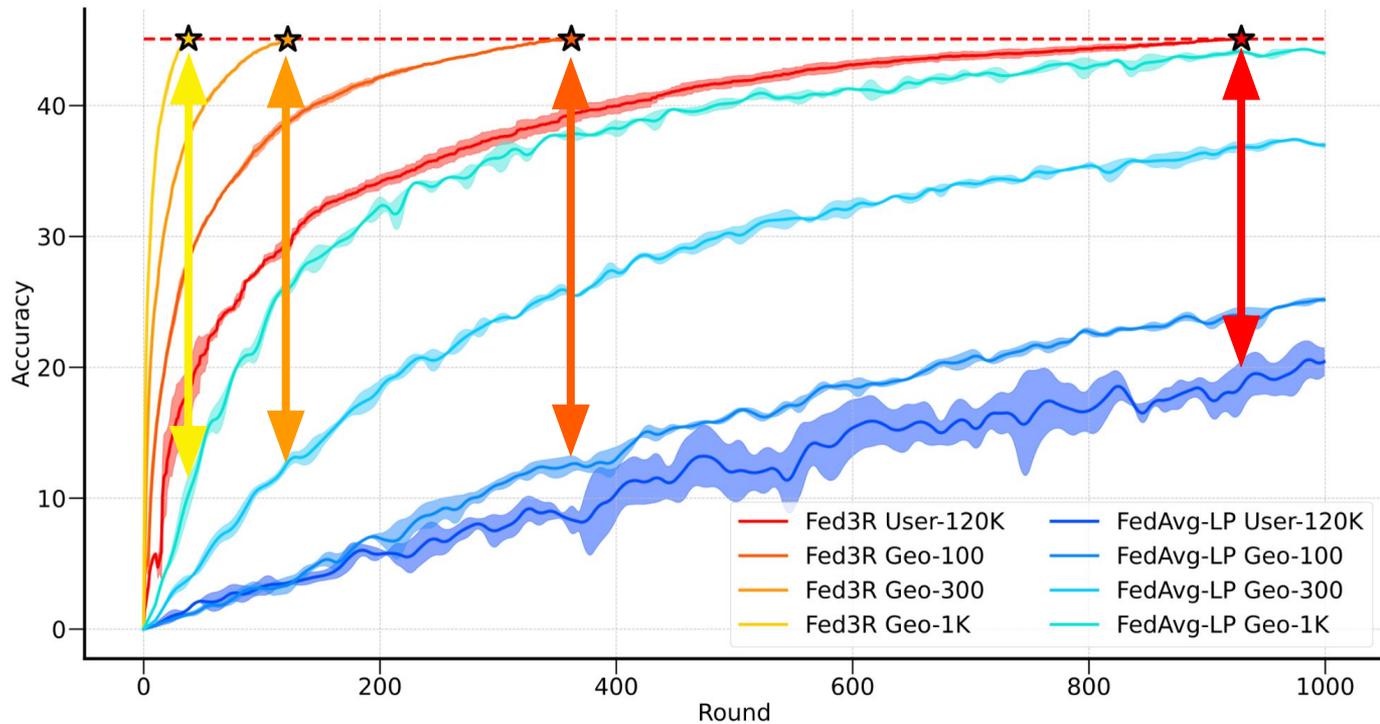
Server side

$$A^t = A^{t-1} + \sum_{k \in \mathcal{K}} A_k^t$$

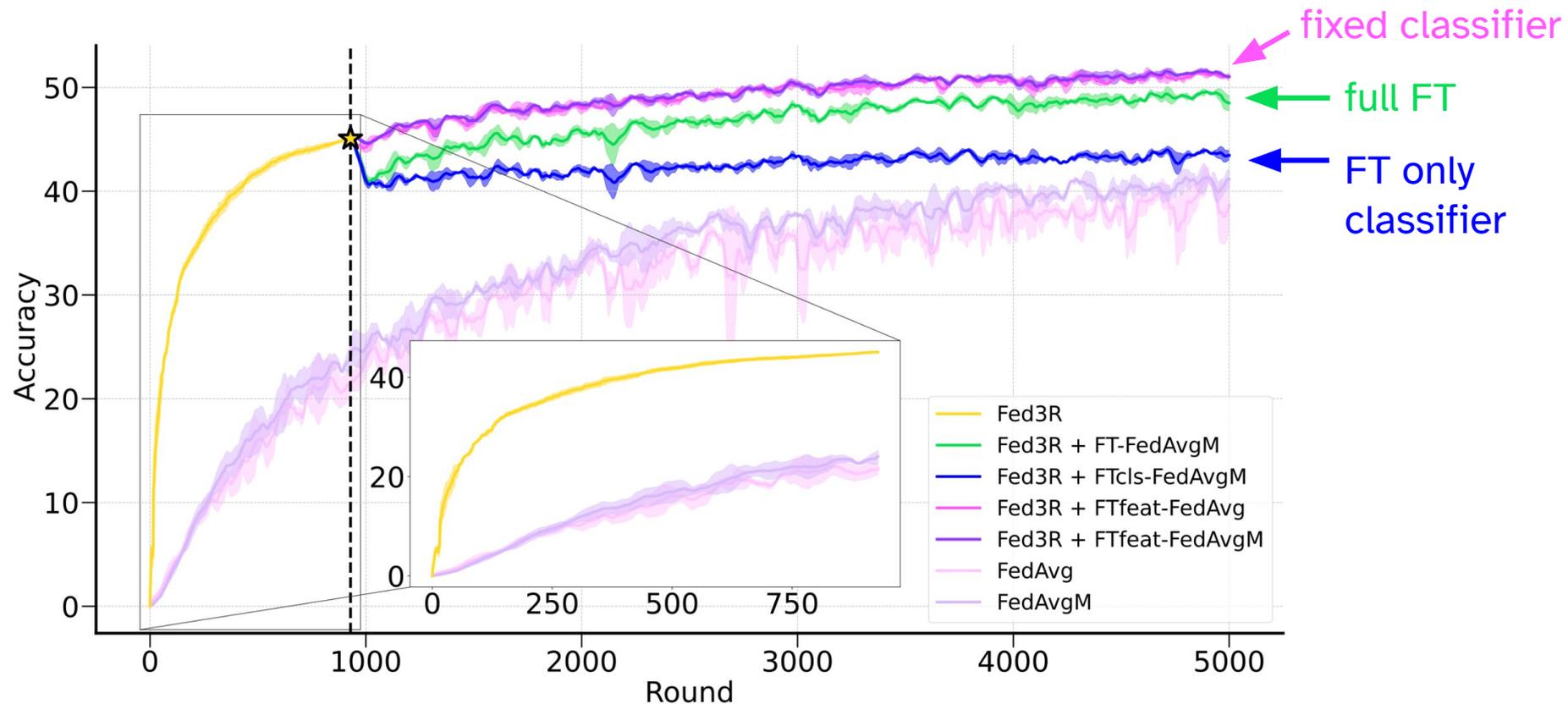
$$b^t = b^{t-1} + \sum_{k \in \mathcal{K}} b_k^t$$

$$\downarrow$$
$$w = (A^t + \lambda I)^{-1} b^t$$

*Fed3R solution is **equivalent to the centralized solution**
(invariant to the **client split** and the **sampling order** by design)*



Fed3R solution can be used as initialization for any federated algorithm



We avoided the problem bypassing the issues of gradient-based learning with heterogeneous data...

(exact classifier aggregation + balanced buffer)



Because of constraints, often we can't use exactly the same approaches across setting.

Because of constraints, often we can't use exactly the same approaches across setting.

But we can transfer similar ideas.

Example:

We can't use replay buffers in FL, but...

*many methods tried to **correct the local update**
by mimicking the centralized update*

Classical Heavy-Ball Momentum

$$\tilde{m}^t \leftarrow \beta \tilde{m}^{t-1} + \tilde{g}^t(\theta^{t-1}; \mathcal{D}^t)$$

$$\theta^t \leftarrow \theta^{t-1} - \eta \tilde{m}^t$$

Classical Heavy-Ball formulation (Polyak, 1964)

$$\theta^t \leftarrow \theta^{t-1} - \eta \tilde{g}^t(\theta^{t-1}; \mathcal{D}^t) + \beta(\theta^{t-1} - \theta^{t-2})$$

Problem: in partial participation, with data heterogeneity the momentum is **biased** towards clients sampled during **last round**

Generalized Heavy-Ball Momentum (GHBM)

Main idea: Recycle gradient from past $\tau > 1$ rounds to correct the bias

$$\tilde{m}^t \leftarrow \beta \tilde{m}^{t-1} + \frac{1}{\tau} \sum_{k=t-\tau+1}^t \tilde{g}^k(\theta^{k-1}; \mathcal{D}^k)$$

Uniform Average
across past gradients

$$\theta^t \leftarrow \theta^{t-1} - \eta \tilde{m}^t$$

Equivalent Heavy Ball formulation

$$\theta^t \leftarrow \theta^{t-1} - \eta \tilde{g}^t(\theta^{t-1}; \mathcal{D}^t) + \frac{\beta}{\tau} (\theta^{t-1} - \theta^{t-\tau-1})$$

recovers classical momentum for $\tau = 1$

Client update rule with GHBM

Inject **global information** into the local optimization to **guide it towards the global solution** and reduce the client drift.

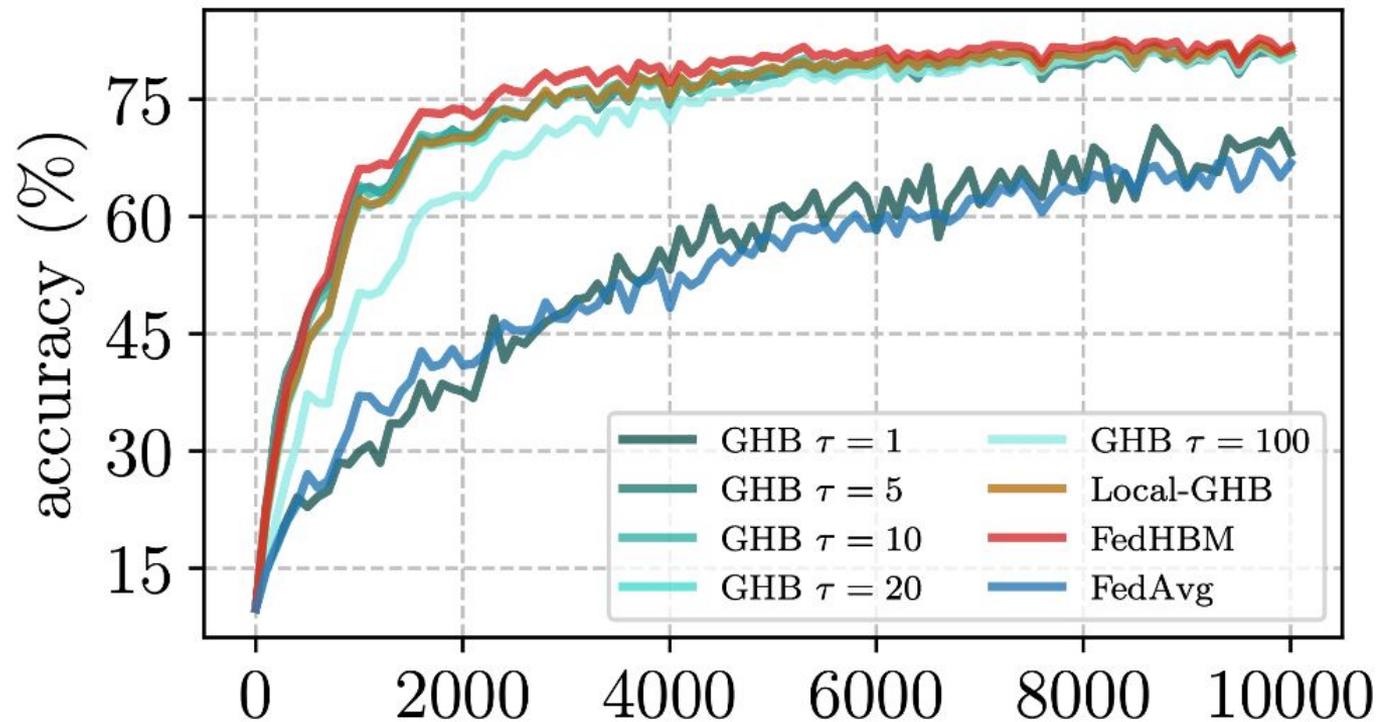
$$\theta_i^{t,j} \leftarrow \theta_i^{t,j-1} - \eta_l \tilde{g}_i^{t,j} + \frac{\beta}{\tau J} \underbrace{(\theta^{t-1} - \theta^{t-\tau-1})}_{\tau\text{-GHBM}}$$

classic local SGD for i-th client

(j-th inner steps, t-th round)

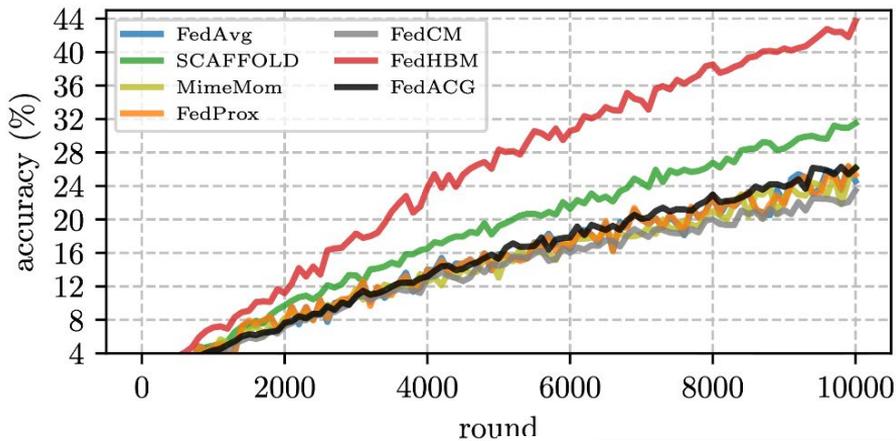
GHBM correction
computed with gradients
from clients of past τ
rounds

GHBM effectively speeds up convergence under heterogeneity and partial participation

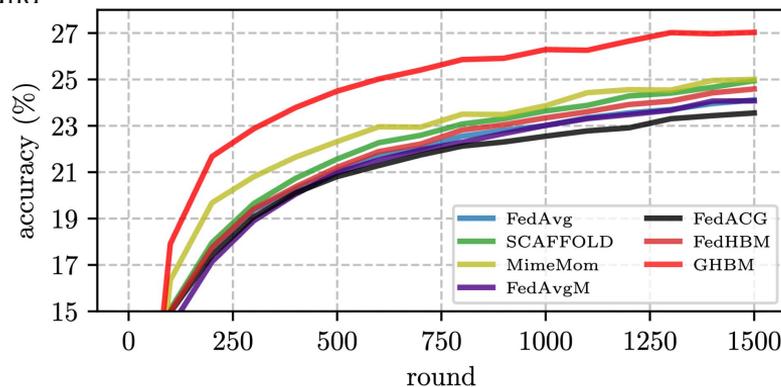
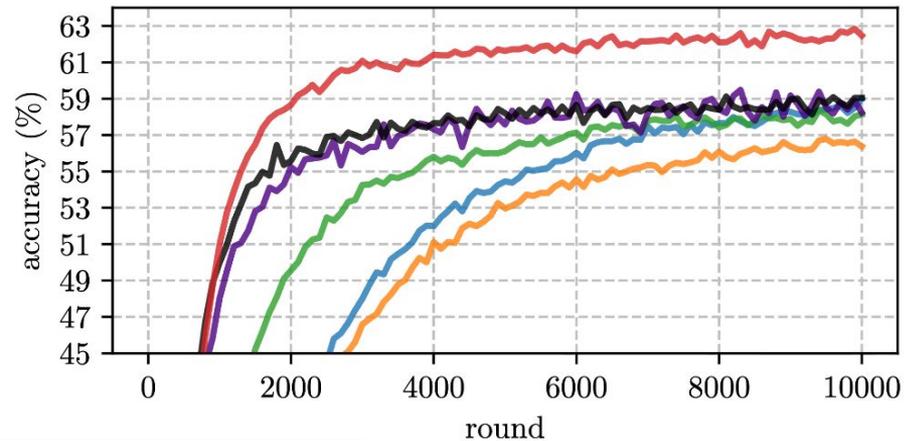


GHBM works well on several datasets and architectures

Cifar-100 Heterogeneous (ResNet)



Cifar-100 Homogeneous (ResNet)

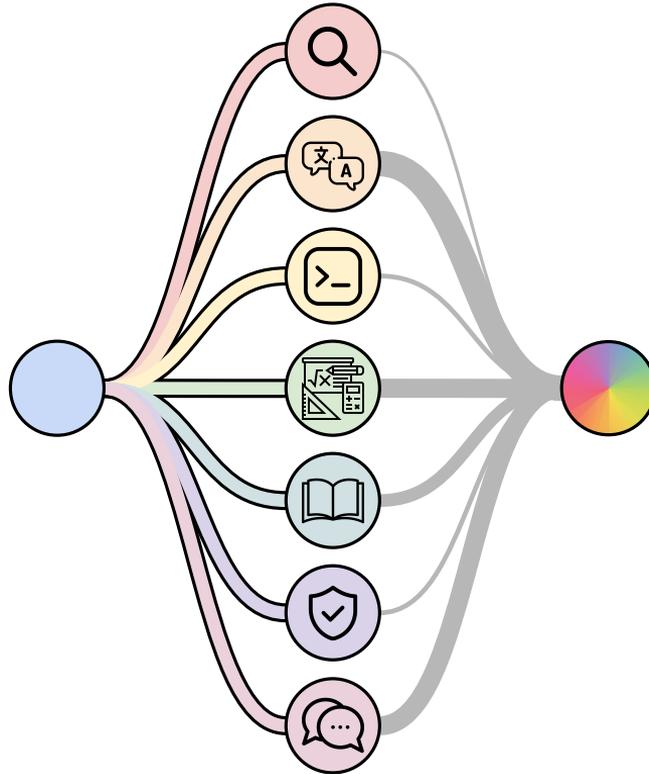


StackOverflow
(Heterogeneous
40K clients)

So far we mainly discussed about local training, what about model aggregation?



Model Merging combines expert models post-hoc
(starting from the same **pre-trained model**)



Core merging approach: **Parameter Averaging**

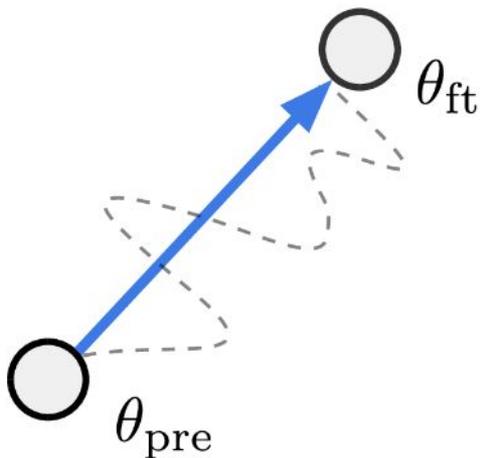
(linear interpolation in the weight space)

$$\theta_{\text{merged}} = \sum_{i=1}^M \lambda_i \theta_i$$

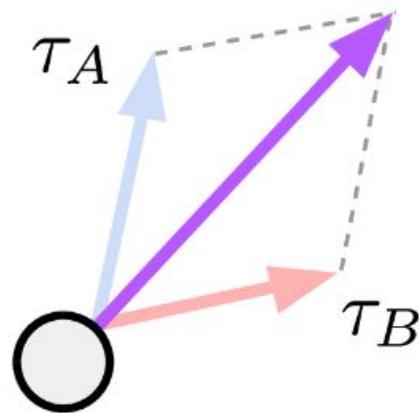
often, but not always, a convex combination

$$\sum_{i=1}^M \lambda_i = 1$$

Task vectors provide a general way of model editing (and merging)



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



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

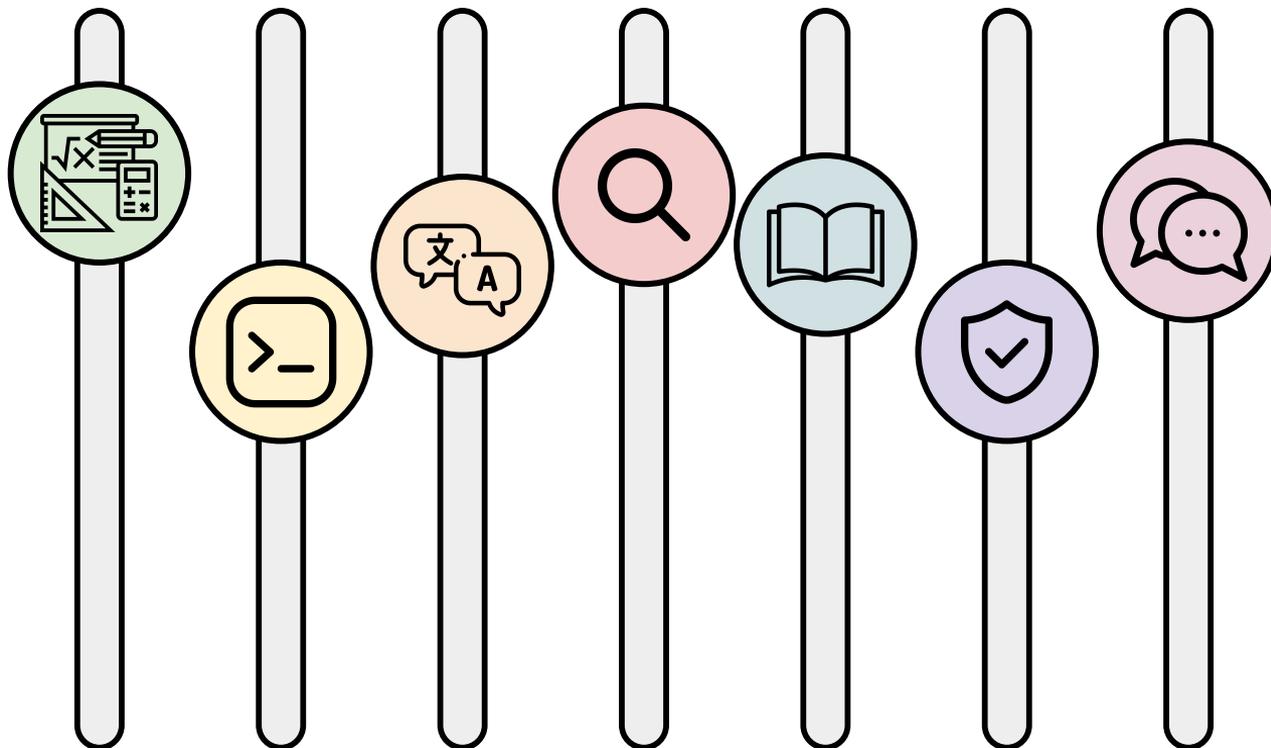
*Same concept of “**pseudo-gradient**” in FL (single round FL)*

*Multi-task or multi-domain
(simultaneous access to all sources)*

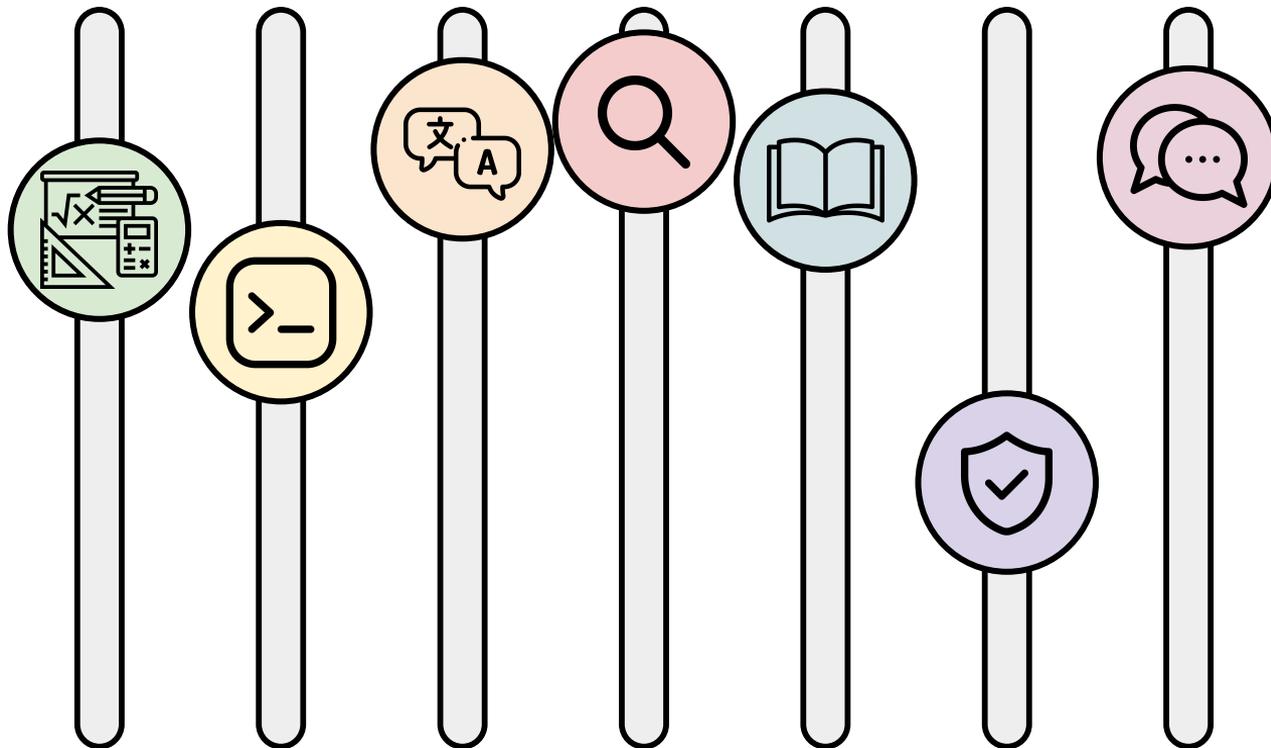
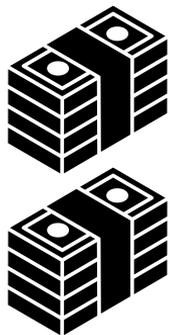


Challenge: *balance across sources
find optimal data mixture*

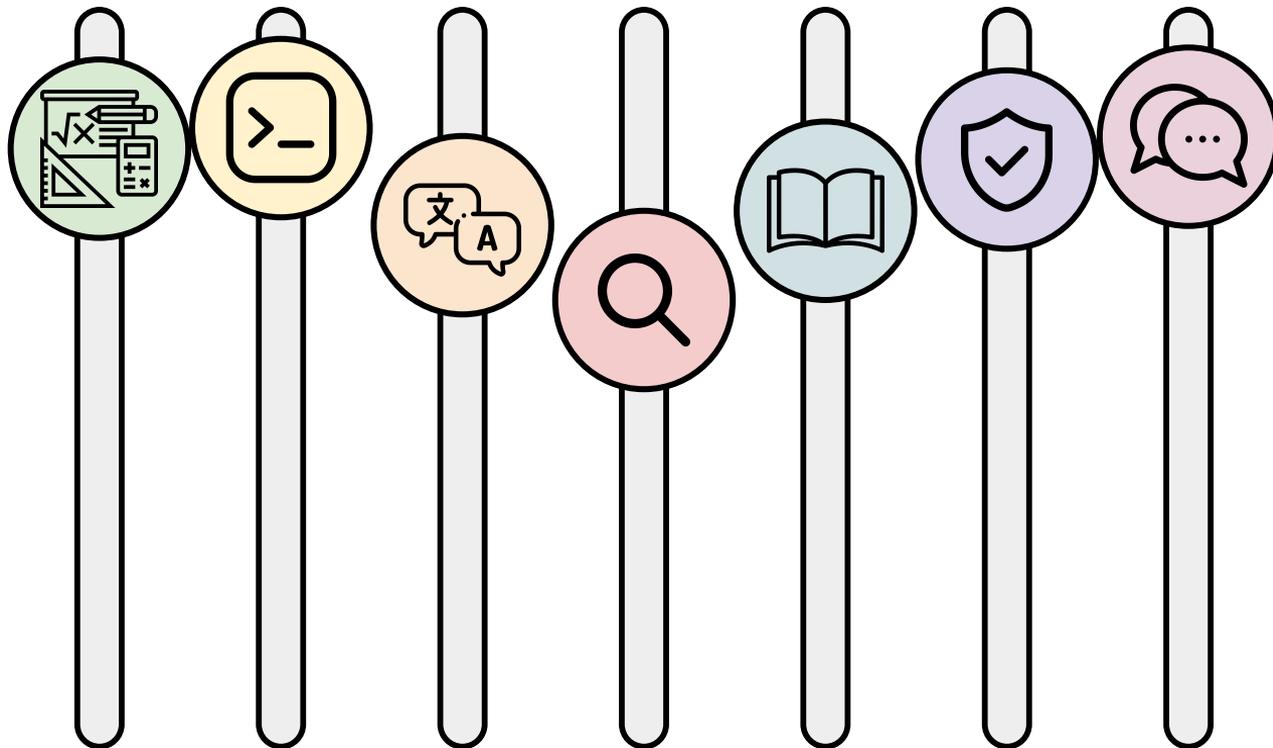
but exploring different mixtures is expensive!



but exploring different mixtures is expensive!

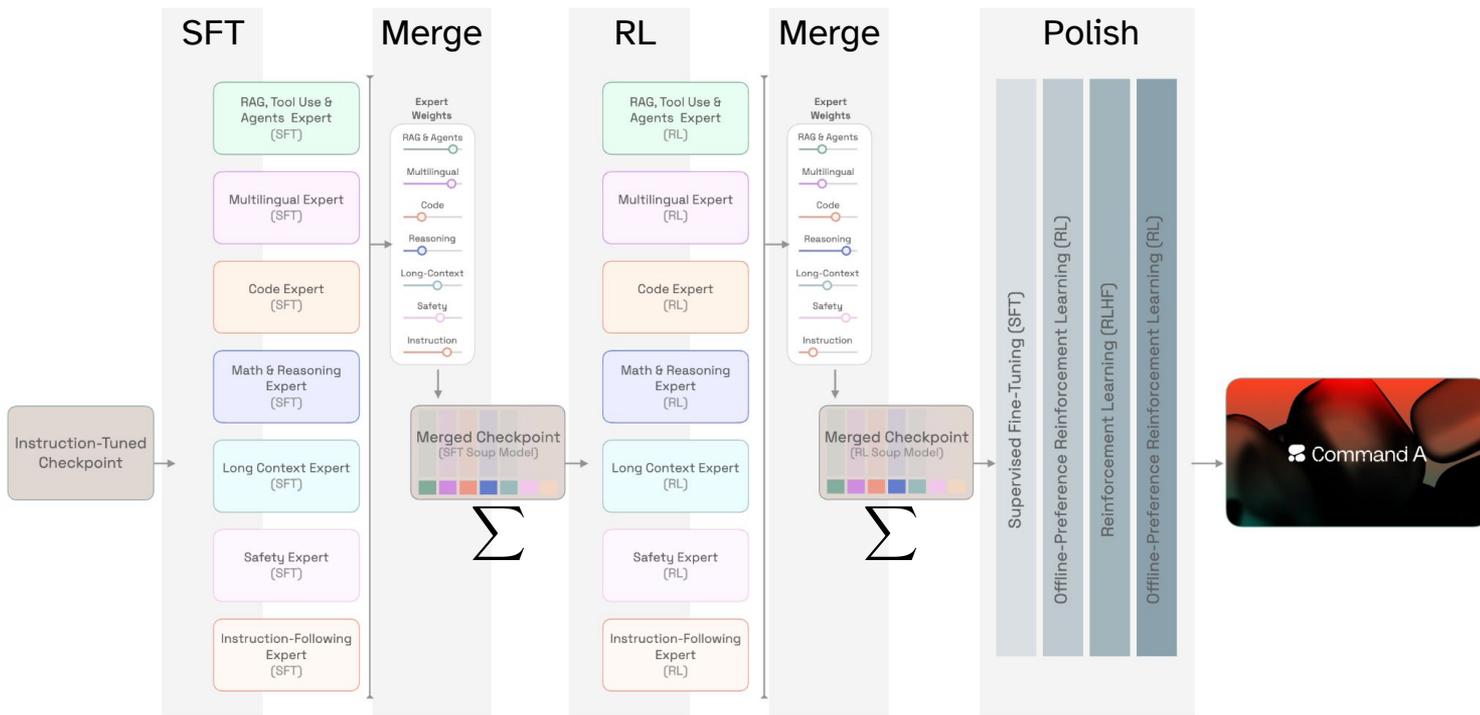


but exploring different mixtures is expensive!



*Merging enables teams to **work asynchronously** on improving capabilities and **cheaply explore different tradeoffs post-hoc***

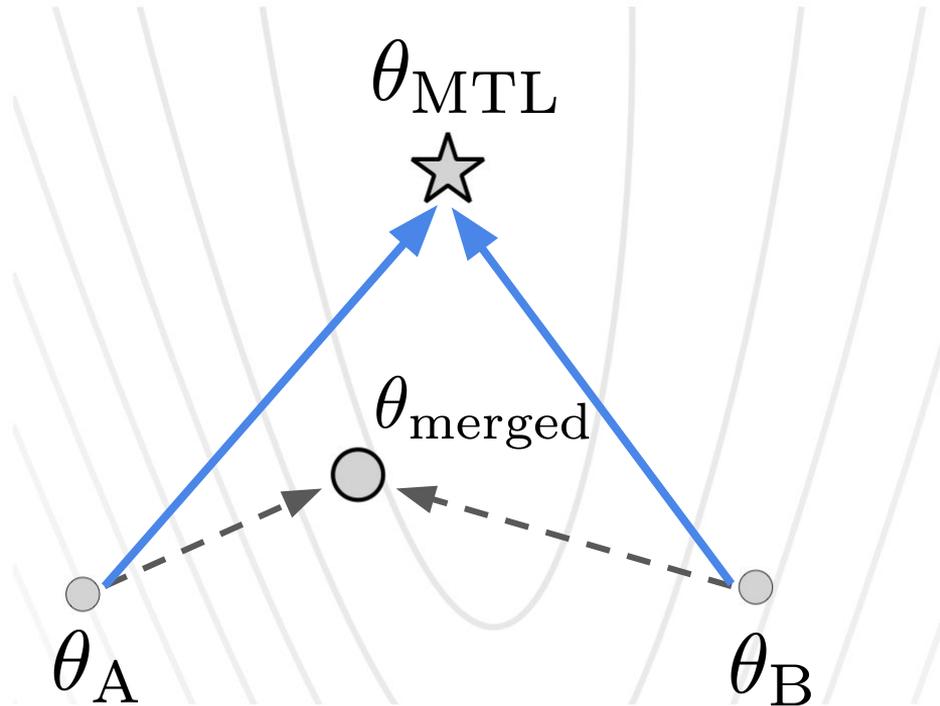
Controlled setting: *in house experts, full access to data*



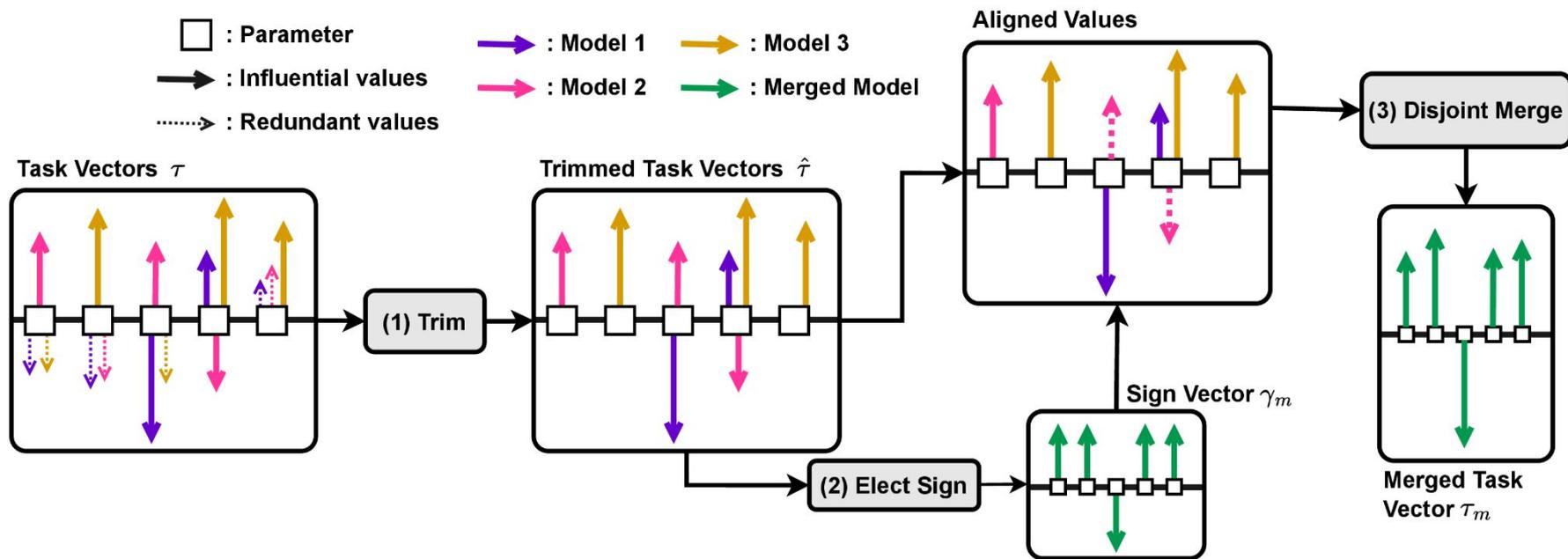
From "Command A: An Enterprise-Ready Large Language Model" by Ahmadian, Cohere, Team, et al.

Independently fine-tuned models “**interfere**” when combined

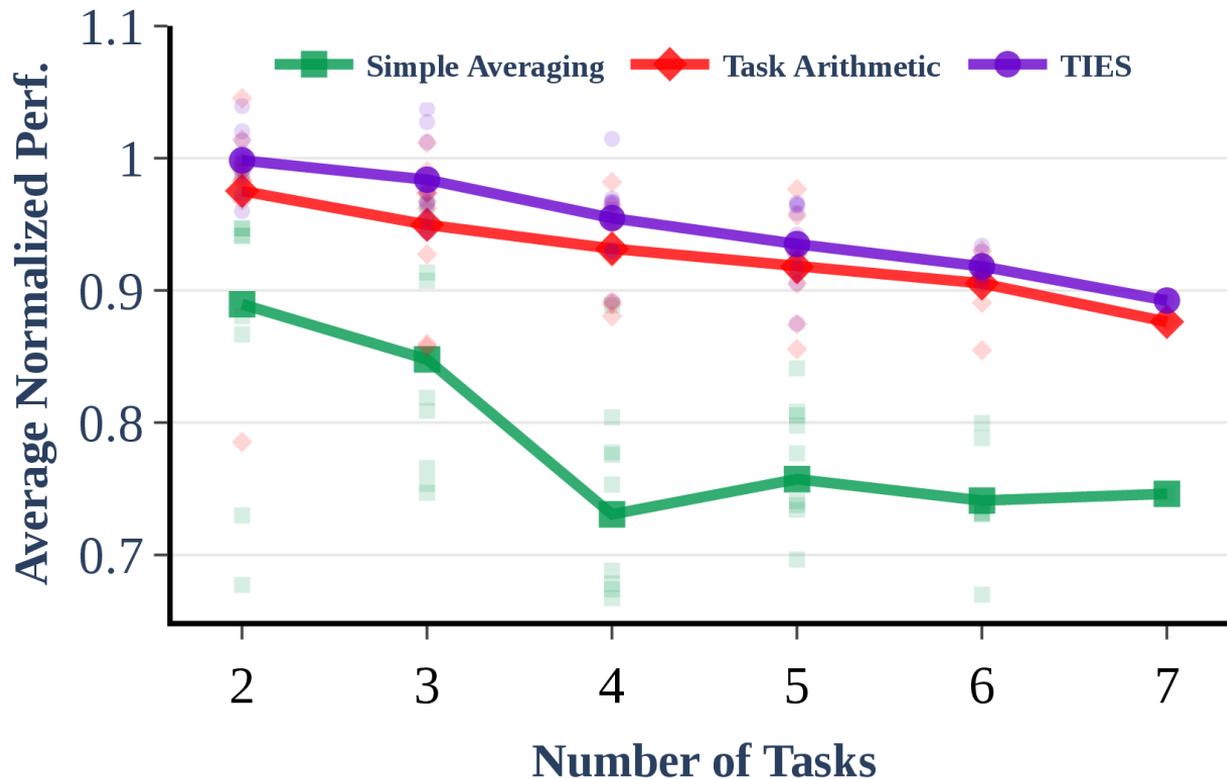
Several **heuristics** exist to improve aggregation and **reduce interference**



Example: TIES Merging to resolve "interference" between task vectors via **pruning** and **sign agreement**



TIES helps retain individual model performance

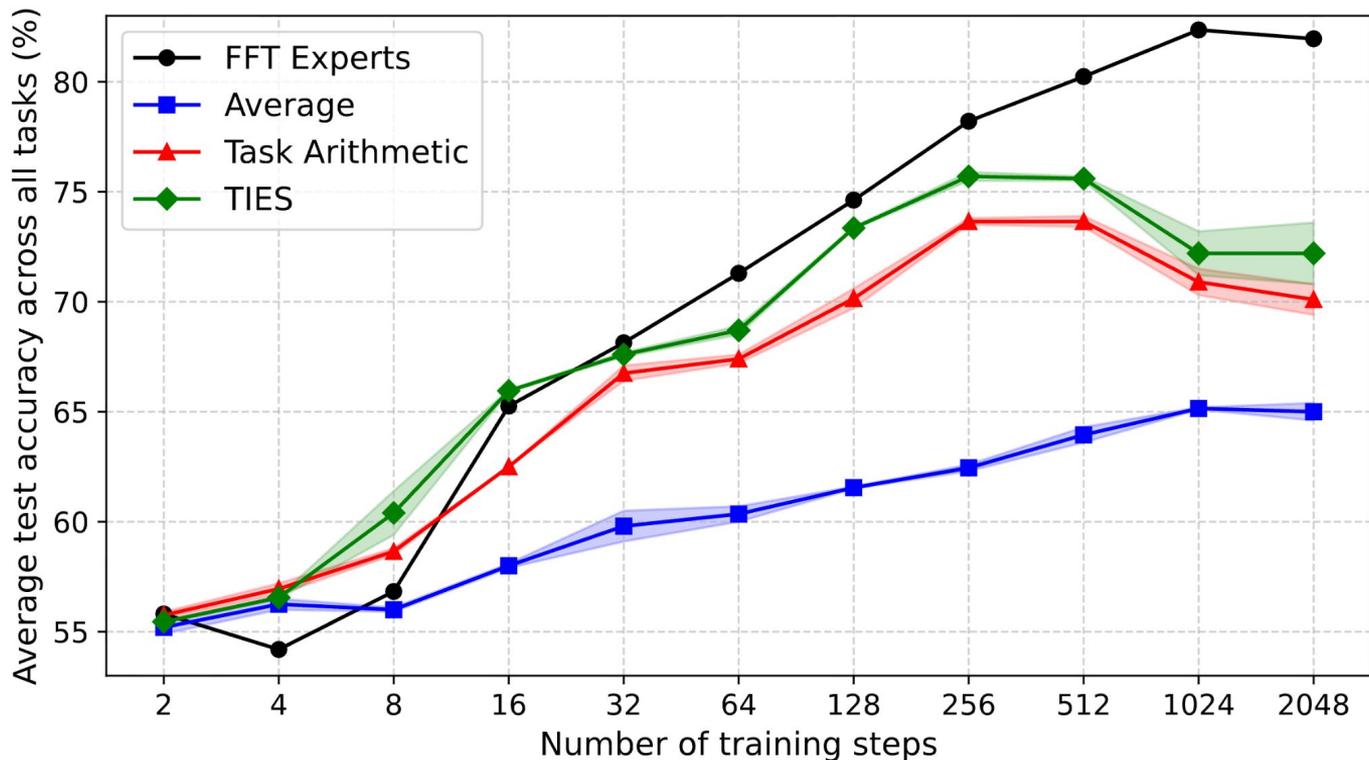


Merging methods treat compatibility as an afterthought, relying on post-hoc fixes.

Merging methods treat compatibility as an afterthought, relying on post-hoc fixes.

*How can we train models that are **composable** by design?*

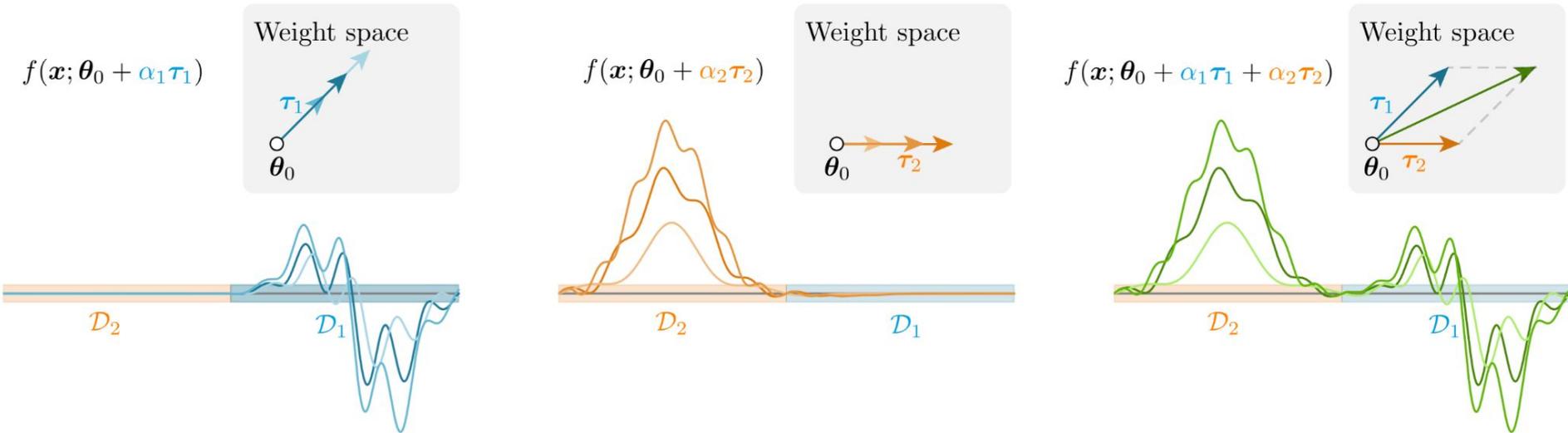
*Local training has a clear effect on mergeability.
Overspecialization hurts aggregation - (client-drift?)*



From "Less is More: Undertraining Experts Improves Model Upcycling", by Horoi et al

“Weight disentanglement” property

Merged model should behave as a composition of **localized** functions where a task vector activates only for the support of its task



We can formalize this property and derive a fine-tuning update that approximately satisfies it

Localized functions: the task vector should activate only for the support of its task

$$f\left(\mathbf{x}, \boldsymbol{\theta}_0 + \sum_{t=1}^T \alpha_t \boldsymbol{\tau}_t\right) = f(\mathbf{x}, \boldsymbol{\theta}_0) \mathbf{1}\left(\mathbf{x} \notin \bigcup_{t=1}^T \mathcal{D}_t\right) + \sum_{t=1}^T f(\mathbf{x}, \boldsymbol{\theta}_0 + \alpha_t \boldsymbol{\tau}_t) \mathbf{1}(\mathbf{x} \in \mathcal{D}_t)$$

Linearizing the merged model we can derive constraints on the task vectors

$$f_{\text{lin}}\left(\mathbf{x}, \boldsymbol{\theta}_0 + \sum_{t=1}^T \alpha_t \boldsymbol{\tau}_t\right) = f(\mathbf{x}, \boldsymbol{\theta}_0) + \sum_{t=1}^T \alpha_t \boldsymbol{\tau}_t^\top \nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}_0)$$

$$\forall \mathbf{x} \in \underline{\mathcal{D}_{t' \neq t}}, \quad \boldsymbol{\tau}_t^\top \nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}_0) = 0$$

We want $\boldsymbol{\tau}_t^\top \nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}_0)$ to be active only for $\mathbf{x} \in \mathcal{D}_t$

To satisfy these constraints during training we require data from other tasks, generally unknown a priori or not available :

$$\forall \mathbf{x} \in \mathcal{D}_{\underline{t' \neq t}}, \quad \boldsymbol{\tau}_t^\top \underline{\nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}_0)} = 0$$

We empirically found that **low-sensitivity parameters**
 (θ_j s.t. $\nabla_{\theta_j} f(x, \theta_0) \approx 0$) **are shared across tasks**

ViT-L/14

Mask Pruning Dataset

Cars	0.96	1.02	0.97	0.97	1.01	1.02	0.99	1.02
DTD	0.95	1.05	0.98	0.98	1.00	1.01	0.98	1.03
EuroSAT	0.97	1.03	0.96	0.99	1.02	1.03	1.00	1.01
GTSRB	0.94	1.01	0.99	0.95	1.00	1.01	0.98	1.05
MNIST	0.96	1.04	0.97	0.99	1.02	1.01	1.01	1.01
RESISC45	0.97	1.03	0.98	0.96	1.01	1.01	0.99	1.04
SUN397	0.98	1.05	0.96	1.00	1.02	1.03	1.01	1.00
SVHN	0.94	1.03	0.97	0.98	1.01	1.02	0.99	1.03
	Cars	DTD	EuroSAT	GTSRB	MNIST	RESISC45	SUN397	SVHN

Evaluation Dataset

T5-Large

QASC	1.00	1.00	1.00	1.00	1.00	1.00	1.00
WikiQA	0.99	1.00	1.00	1.00	1.00	1.00	1.00
QuaRTz	1.01	1.00	1.00	1.00	1.00	1.00	1.00
PAWS	0.98	1.00	1.00	1.00	1.00	1.00	1.00
Story Cloze	1.02	1.00	1.00	1.00	1.00	1.00	1.00
Winogrande	0.99	1.00	1.00	1.00	1.00	1.00	1.00
WSC	1.01	1.00	1.00	1.00	1.00	1.00	1.00
	QASC	WikiQA	QuaRTz	PAWS	Story Cloze	Winogrande	WSC

Evaluation Dataset

We can derive a sparse fine-tuning rule that approximately satisfy the constraints

Key idea: update only parameters that have **low-sensitivity**

Sparse Fine-Tuning with mask \mathbf{c}

$$\boldsymbol{\theta}^{(i)} = \boldsymbol{\theta}^{(i-1)} - \gamma[\mathbf{c} \odot \nabla_{\boldsymbol{\theta}} \mathcal{L}(f(\mathbf{x}, \boldsymbol{\theta}^{(i-1)}), \mathbf{y})]$$

$$j = 1, \dots, m \quad \mathbf{c}_j = \begin{cases} 1 & \text{if } \nabla_{\theta_j} f(\mathbf{x}, \boldsymbol{\theta}_0) \approx 0 \\ 0 & \text{otherwise} \end{cases}$$

Sparse Fine-Tuning **preserves performance** of the pre-trained model

Non-linear full FT

Fine-Tuning Tasks

Cars	1.31	0.89	0.73	0.90	0.64	0.81	0.90	1.09
DTD	0.81	2.30	0.78	0.58	0.95	0.63	0.80	1.16
EuroSAT	0.77	0.81	2.25	0.73	0.88	0.55	0.84	0.55
GTSRB	0.67	0.66	0.54	2.84	0.78	0.46	0.73	1.42
MNIST	0.69	0.64	0.39	0.87	2.03	0.49	0.79	1.87
RESISC45	0.80	0.84	0.73	1.25	0.85	1.55	0.83	1.04
SUN397	0.80	0.94	1.11	0.89	0.97	0.87	1.20	1.09
SVHN	0.52	0.73	0.36	0.95	1.76	0.56	0.81	3.55

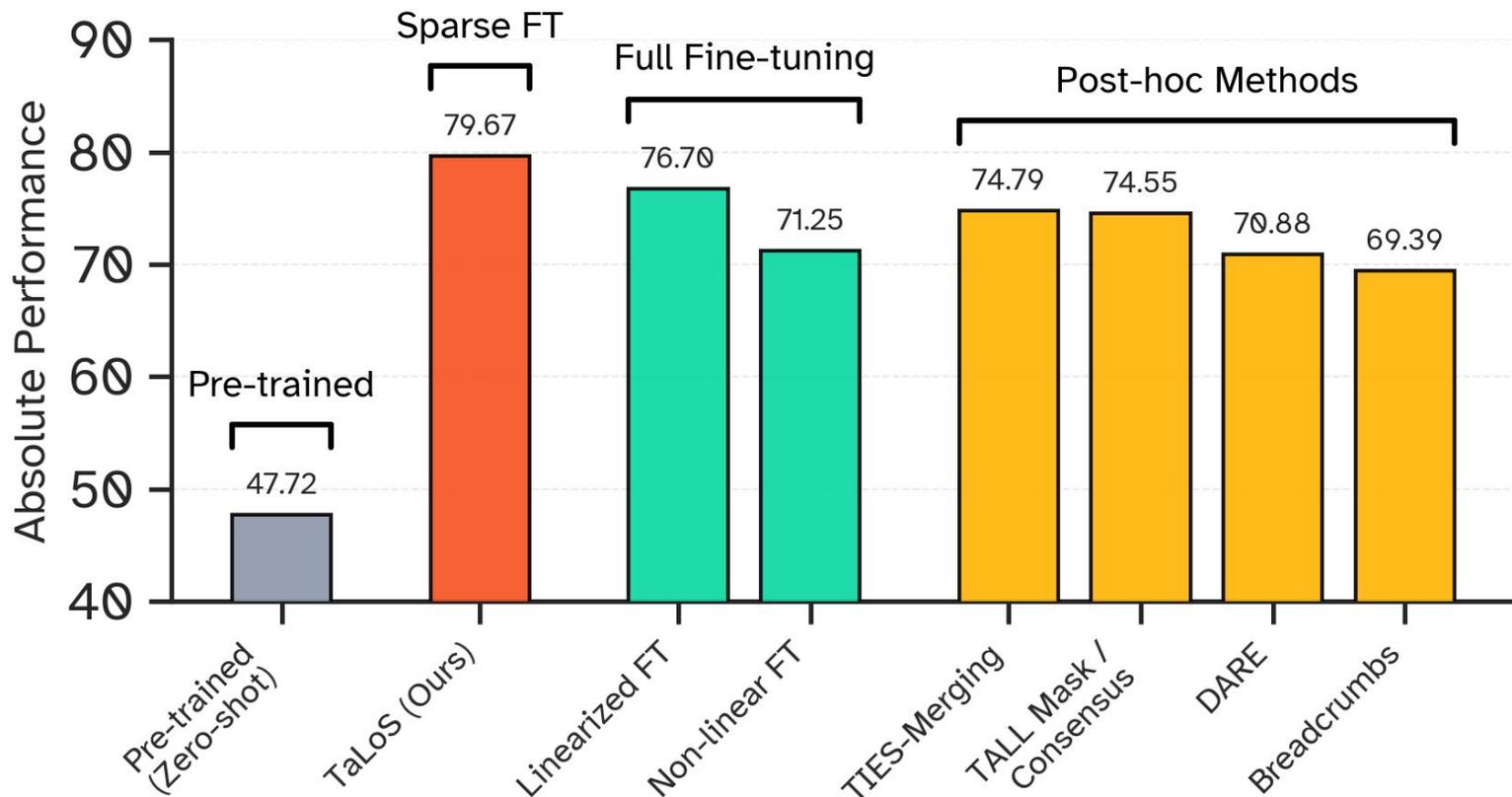
Eval tasks

TaLoS (our sparse FT)

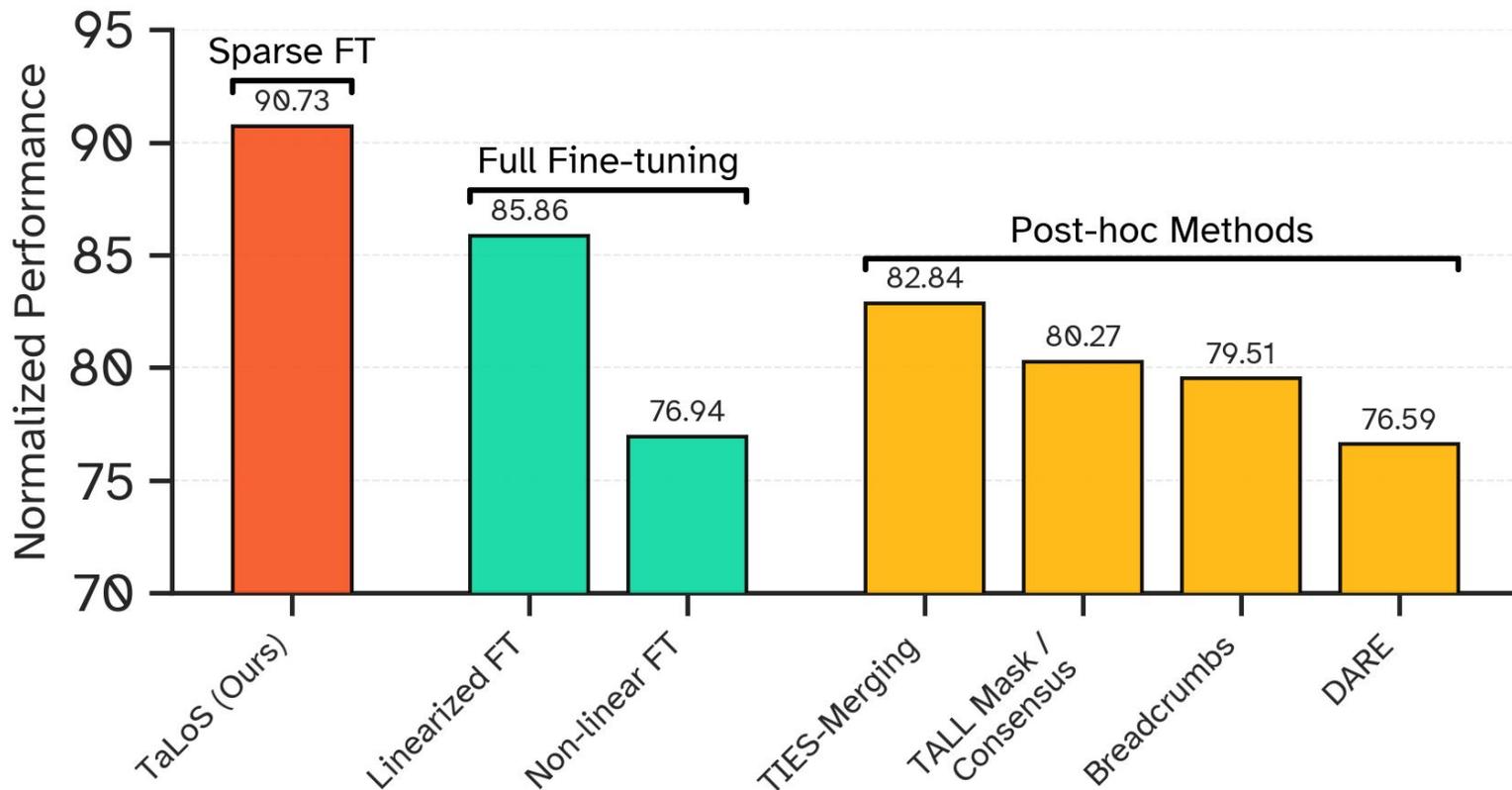
Cars	1.17	1.04	1.09	0.96	1.07	1.00	1.01	1.16
DTD	0.98	2.04	0.97	0.87	0.90	0.95	1.00	1.21
EuroSAT	1.01	1.02	2.19	0.81	0.97	0.99	1.00	1.16
GTSRB	1.01	1.02	1.08	2.71	0.91	1.04	1.01	1.73
MNIST	1.00	0.98	1.01	0.94	2.00	0.98	1.01	0.60
RESISC45	1.01	1.01	1.32	0.92	0.90	1.49	1.00	1.20
SUN397	1.01	1.05	0.98	1.06	1.06	0.99	1.19	1.11
SVHN	1.01	1.02	0.93	1.06	1.54	0.98	1.01	3.39

Eval tasks

Sparse Fine-Tuning **reduces interference** when merging...



...and better retains individual performance



Takeaways

- *Learning from diverse sources presents **shared challenges across settings.***
- *Depending on constraints, we can focus on local training, aggregation or both, with **methods transferable across settings.***
- *Sub-communities should collaborate; a **unifying framework** could accelerate research and prevent redundant efforts.*

So, what's next?



Some meta-questions to spark discussions

- *How do large pre-trained models transform our approach to data heterogeneity and enable **building modular models**?*
- *Can we automatically determine when **composition and transfer** are beneficial across sources, and when they are harmful?*
- *When should we prefer **merging** over other parameter isolation methods such as **routing over experts**?*