

# The Power of Adaptivity in Representation Learning: From Meta-Learning to Federated Learning

**Aryan Mokhtari**  
ECE Department, UT Austin

TILOS - OPTML Seminar  
MIT  
October 26th, 2022

- MAML and ANIL Provably Learn Representations, ICML 2022.



Liam Collins  
UT Austin



Sewoong Oh  
UW



Sanjay Shakkottai  
UT Austin

- FedAvg with Fine Tuning: Local Updates Lead to Representation Learning, NeurIPS 2022



Liam Collins  
UT Austin



Hamed Hassani  
UPenn

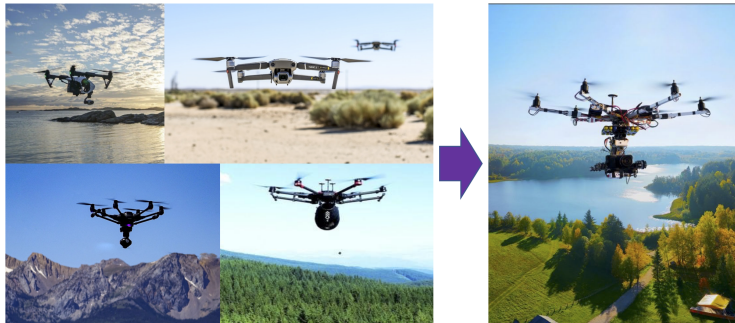


Sanjay Shakkottai  
UT Austin

- Consider the Multi-Task Learning (MTL) setup
  - ⇒ We are given a set of tasks
  - ⇒ For each task, we have access to some (small number of) samples
  - ⇒ The tasks often share some similarity (but are not identical)
- Main Goal: Train a model that **generalizes well to new tasks/domains**
  - ⇒ With **limited data and computation!**



Image credits: Alex Krizhevsky, Learning Multiple Layers of Features from Tiny Images, 2009.



Training Data

Task in new Scenario

Image Credits: [bit.ly/3i5m8ay](https://bit.ly/3i5m8ay), [bit.ly/3w723ZY](https://bit.ly/3w723ZY), [bit.ly/3KHMq5E](https://bit.ly/3KHMq5E), [bit.ly/3i7pREJ](https://bit.ly/3i7pREJ), [bit.ly/34I1ytT](https://bit.ly/34I1ytT)

- Training data in various scenarios; goal is to quickly adapt to the task in a new scenario



- Suppose that we have access to
  - ⇒ large amounts of data from different training environments
  - ⇒ a small amount of data available just prior to deployment from the deployment environment
- Given this setup how should we train our model?
- Possible Approach:
  - (a) Build a model using data from the training environments
  - (b) Fine-tune the model using the small amount of deployment data

How can we build a model that is easily fine-tunable?

**Obvious Approach:** Build a model to minimize average training loss, and then fine tune for deployment

- Set of tasks:  $\mathcal{T} = \{\mathcal{T}_i\}_{i=1}^n$  coming from distribution  $p$
- Select a model  $\theta_{train}^*$

## Average Risk (Loss) Minimization

$$\theta_{train}^* \in \operatorname{argmin}_{\theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta)$$

- A new task  $\mathcal{T}_{test}$  is revealed, drawn according to dist.  $p$
- Fine tune the model:  $\theta_{train}^* \rightarrow \theta_{new}^*$
- Performance:  $f_{test}(\theta_{new}^*)$

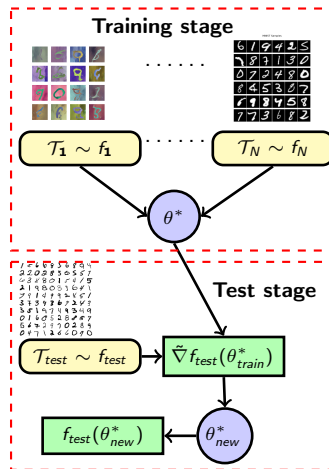


Image Credits: <https://bit.ly/392pda9>,  
<https://bit.ly/3EEIElq>

# Does (ARM + Fine Tuning) work well?

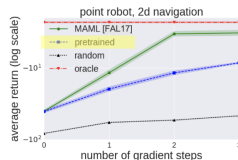
- Suppose we have images from a large number of classes (e.g., Imagenet)
  - Task == classifying images among a  $K$ -subset of these classes, small  $K$
  - Many different subsets == Many tasks



- ARM + Fine Tuning has mixed performance [FAL17]

|   | 5-way Accuracy    |                   |
|---|-------------------|-------------------|
|   | 1-shot            | 5-shot            |
| MiniImagenet (Ravi & Larochelle, 2017)        |                   |                   |
| fine-tuning baseline                          | 28.86 $\pm$ 0.54% | 49.79 $\pm$ 0.79% |
| nearest neighbor baseline                     | 41.08 $\pm$ 0.70% | 51.04 $\pm$ 0.65% |
| matching nets (Vinyals et al., 2016)          | 43.56 $\pm$ 0.84% | 55.31 $\pm$ 0.73% |
| meta-learner LSTM (Ravi & Larochelle, 2017)   | 43.44 $\pm$ 0.77% | 60.60 $\pm$ 0.71% |
| MAML, first order approx. (Finn et al., 2017) | 48.07 $\pm$ 1.75% | 63.15 $\pm$ 0.91% |
| MAML (Finn et al., 2017)                      | 48.70 $\pm$ 1.84% | 63.11 $\pm$ 0.92% |

“Fine-tuning baseline”: Few-shot image classification accuracy of ARM after fine-tuning (image taken from [FAL17])



“Pretrained”: Fine-tuning reward for ARM on robot 2d navigation task (image taken from [FAL17])

- These experiments manifest that ARM + Fine-tuning does not work well!  
⇒ The solution trained by ARM often does not generalize well
- For multi-task linear regression, we can formally justify this!
- We show that with adaptivity, we can capture the common structure/representation among the tasks

Takeaway: The adaptivity idea can be done in two ways:

- Change the loss function ⇒ as done in MAML
- Change the training procedure ⇒ as done in Federated Learning

## Meta-learning

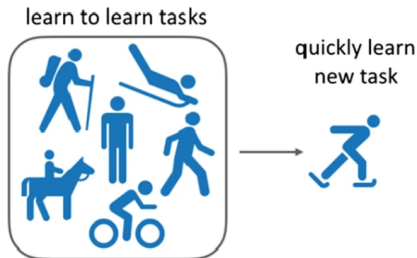


Image credit: <https://meta-world.github.io>, [HRJ21]

## Model-Agnostic Meta Learning [FAL17]

- Set of tasks:  $\mathcal{T} = \{\mathcal{T}_i\}_{i=1}^n$  coming from distribution  $p$
- Select a model  $\theta_{train}^*$

## Change the Loss Function

$$\theta_{train}^* \in \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n f_i(\theta - \alpha \nabla f_i(\theta; (X_i, y_i)))$$

- A new task  $\mathcal{T}_{test}$  is revealed, drawn according to dist.  $p$
- Fine tune the model:  $\theta_{train}^* \rightarrow \theta_{new}^*$
- Performance:  $f_{test}(\theta_{new}^*)$

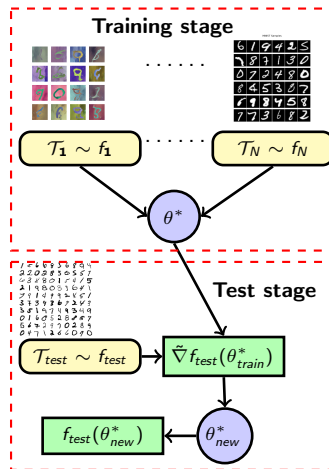


Image Credits: <https://bit.ly/392pda9>,

<https://bit.ly/3EEIElq>

Seems like finding the right initialization for adaptation!

- **Average Risk Minimization (ARM):**  $\min_{\theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta)$
- GD update for ARM.:  $\theta_{t+1} = \theta_t - \frac{\beta}{n} \sum_{i=1}^n \nabla f_i(\theta_t)$
- Gradient evaluated at same  $\theta_t$  for all tasks  $\implies$  **not adaptive**

- **Average Risk Minimization (ARM):**  $\min_{\theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta)$
  - GD update for ARM.:  $\theta_{t+1} = \theta_t - \frac{\beta}{n} \sum_{i=1}^n \nabla f_i(\theta_t)$
  - Gradient evaluated at same  $\theta_t$  for all tasks  $\implies$  **not adaptive**
- 
- **Model-Agnostic Meta-Learning (MAML):**  $\min_{\theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta - \alpha \nabla f_i(\theta))$
  - GD update on MAML loss can be implemented as follows

$$\theta_{t+1} = \theta_t - \frac{\beta}{n} \sum_{i=1}^n (I - \alpha \nabla^2 f_i(\theta_t)) \nabla f_i(\theta_t - \alpha \nabla f_i(\theta_t))$$



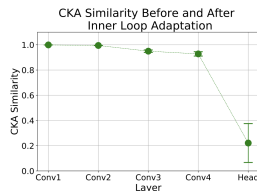
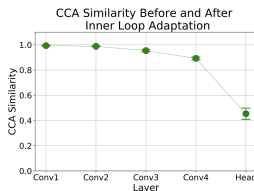
- **Average Risk Minimization (ARM):**  $\min_{\theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta)$
- GD update for ARM.:  $\theta_{t+1} = \theta_t - \frac{\beta}{n} \sum_{i=1}^n \nabla f_i(\theta_t)$
- Gradient evaluated at same  $\theta_t$  for all tasks  $\implies$  **not adaptive**
- **Model-Agnostic Meta-Learning (MAML):**  $\min_{\theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta - \alpha \nabla f_i(\theta))$
- GD update on MAML loss can be implemented as follows

$$\theta_{t+1} = \theta_t - \frac{\beta}{n} \sum_{i=1}^n (I - \alpha \nabla^2 f_i(\theta_t)) \nabla f_i(\theta_t - \alpha \nabla f_i(\theta_t))$$

which can be implemented via inner and outer loops

- **Inner loop:** Compute  $\theta_{t,i} = \theta_t - \alpha \nabla f_i(\theta_t)$  for  $i = 1, \dots, n$
- **Outer loop:** Compute  $\theta_{t+1} = \theta_t - \frac{\beta}{n} \sum_{i=1}^n (I - \alpha \nabla^2 f_i(\theta_t)) \nabla f_i(\theta_{t,i})$
- $\theta_{t,i}$  adapted to each task  $\implies$  **adaptive**

- 1 MAML learns models that can **quickly solve new tasks** [FAL17, AES19]
  - In image classification, sinusoid regression, reinforcement learning.
- 2 MAML seems to be learning a representation **shared across tasks** [RRBV20]
  - Even though it is not designed for representation learning!



**The representation learned by MAML does not change significantly when adapted to each task.**  
Figure credit: [RRBV20]

- Can we formally prove this claim for some multi-task learning setting?

Multi-task linear regression:

- Task  $i$  has ground-truth solution  $\theta_{*,i} \in \mathbb{R}^d$ :

$$y_i \sim \theta_{*,i}^\top \mathbf{x}_i + z_i$$

$\Rightarrow \mathbf{x}_i$  is a random feature vector and  $z_i \in \mathbb{R}$  is random, mean-zero noise.

- Solving each task individually would require  $\Omega(d)$  samples per task.

Multi-task linear regression:

- Task  $i$  has ground-truth solution  $\theta_{*,i} \in \mathbb{R}^d$ :

$$y_i \sim \theta_{*,i}^\top \mathbf{x}_i + z_i$$

$\Rightarrow \mathbf{x}_i$  is a random feature vector and  $z_i \in \mathbb{R}$  is random, mean-zero noise.

- Solving each task individually would require  $\Omega(d)$  samples per task.
- Now suppose the  $\theta_{*,i}$  lie in a shared  $k$ -dimensional subspace,  $k \ll d$
- Let the columns of  $\mathbf{B}_* \in \mathbb{R}^{d \times k}$  span this subspace, that is, for all tasks there exists  $\mathbf{w}_{*,i} \in \mathbb{R}^k$  such that

$$\theta_{*,i} = \mathbf{B}_* \mathbf{w}_{*,i}$$

- Task Diversity:** The concatenation of  $\mathbf{w}_{*,1}, \dots, \mathbf{w}_{*,n}$  spans  $\mathbb{R}^k$
- If we know  $\text{col}(\mathbf{B}_*)$ , we can solve new tasks with only  $O(k)$  samples

Multi-task linear regression:

- Task  $i$  has ground-truth solution  $\theta_{*,i} \in \mathbb{R}^d$ :

$$y_i \sim \theta_{*,i}^\top \mathbf{x}_i + z_i$$

$\Rightarrow \mathbf{x}_i$  is a random feature vector and  $z_i \in \mathbb{R}$  is random, mean-zero noise.

- Solving each task individually would require  $\Omega(d)$  samples per task.
- Now suppose the  $\theta_{*,i}$  lie in a shared  $k$ -dimensional subspace,  $k \ll d$
- Let the columns of  $\mathbf{B}_* \in \mathbb{R}^{d \times k}$  span this subspace, that is, for all tasks there exists  $\mathbf{w}_{*,i} \in \mathbb{R}^k$  such that

$$\theta_{*,i} = \mathbf{B}_* \mathbf{w}_{*,i}$$

- Task Diversity:** The concatenation of  $\mathbf{w}_{*,1}, \dots, \mathbf{w}_{*,n}$  spans  $\mathbb{R}^k$
- If we know  $\text{col}(\mathbf{B}_*)$ , we can solve new tasks with only  $O(k)$  samples
- Does GD on ARM learn  $\mathbf{B}^*$ ? Does GD on MAML learn  $\mathbf{B}^*$ ?

- Loss function for task  $i$ :

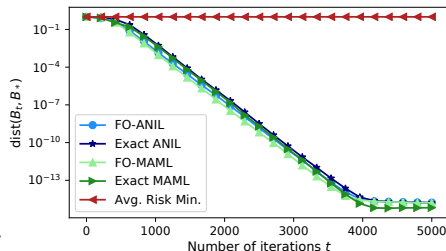
$$f_i(\mathbf{B}, \mathbf{w}) := \frac{1}{2} \mathbb{E}_{\mathbf{x}_i, y_i} [(\langle \mathbf{B}\mathbf{w}, \mathbf{x}_i \rangle - y_i)^2]$$

- ARM: Uses GD to solve  $\min_{\mathbf{B}, \mathbf{w}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{B}, \mathbf{w})$

### Algorithm

- (Outer loop) For  $t = 1, \dots, T$ :
  - Select  $n$  tasks satisfying diversity condition ( $\text{span}(\{\mathbf{w}_{*,i}\}_{i \in [n]}) = \mathbb{R}^k$ )
  - (Inner loop) For  $i = 1, \dots, n$ :
    - - Adapt head:  $\mathbf{w}_{t,i} = \mathbf{w}_t - \alpha \nabla_{\mathbf{w}} f_i(\mathbf{B}_t, \mathbf{w}_t)$ .
    - - If MAML: Adapt rep:  $\mathbf{B}_{t,i} = \mathbf{B}_t - \alpha \nabla_{\mathbf{B}} f_i(\mathbf{B}_t, \mathbf{w}_t)$   
 If ANIL: Do not adapt rep:  $\mathbf{B}_{t,i} = \mathbf{B}_t$
  - $\begin{bmatrix} \mathbf{w}_{t+1} \\ \bar{\mathbf{B}}_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{w}_t \\ \bar{\mathbf{B}}_t \end{bmatrix} - \frac{\beta}{n} \sum_{i=1}^n \mathbf{H}_{t,i,\text{Alg}}(\mathbf{B}_t, \mathbf{w}_t) \begin{bmatrix} \nabla_{\mathbf{w}} f_i(\mathbf{B}_{t,i}, \mathbf{w}_{t,i}) \\ \nabla_{\mathbf{B}} f_i(\mathbf{B}_{t,i}, \mathbf{w}_{t,i}) \end{bmatrix}$  where  $\mathbf{H}_{t,i,\text{Alg}}(\mathbf{B}_t, \mathbf{w}_t)$  is a Hessian that differs between MAML and ANIL.

- We consider four meta-learning algorithms:
  - MAML,
  - ANIL [RRBV19], a close relative of MAML, and
  - their first-order approximations (FO-MAML and FO-ANIL).



(2).pdf

For multi-task linear regression with population losses, the meta-learning approaches learn  $\text{col}(\mathbf{B}_*)$ , while ARM does not.

- How can we explain this? Can we formalize this observation?

- In the multi-task linear representation learning setting:

## Theorem (informal)

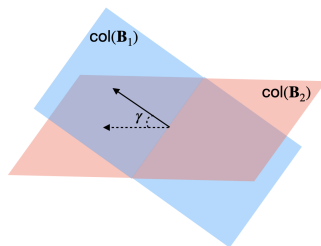
- *Under standard assumptions, MAML, ANIL and their first-order analogues recover  $\text{col}(\mathbf{B}_*)$  exponentially fast when run on the task population losses.*
  - *ANIL and FO-ANIL require  $m = \Omega((\frac{d}{n} + 1)k^3) \ll d$  samples per task to learn the ground-truth subspace.*
  - *The key is that MAML and ANIL's adaptation of the head harnesses **task diversity** to improve the representation in all directions.*
- **First results** showing that MAML and ANIL provably learn effective representations!

## Theorem (informal)

*There exist problems for which the model trained by ARM fails to learn  $\text{col}(\mathbf{B}_*)$ .*



- We use the **principal angle distance** to measure the distance between representations.



- Formally,

$$\text{dist}(\mathbf{B}_1, \mathbf{B}_2) := \|\hat{\mathbf{B}}_{1,\perp}^\top \hat{\mathbf{B}}_2\|_2,$$

where  $\hat{\mathbf{B}}_{1,\perp}$  and  $\hat{\mathbf{B}}_2$  are orthonormal matrices s.t.  $\text{col}(\hat{\mathbf{B}}_{1,\perp}) = \text{col}(\mathbf{B}_1)^\perp$  and  $\text{col}(\hat{\mathbf{B}}_2) = \text{col}(\mathbf{B}_2)$ .

- Let's focus on the population case to simplify the expressions

$$\text{ARM: } \min_{\mathbf{B}, \mathbf{w}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{B}, \mathbf{w})$$

- In this case we can show:

$$\mathbf{B}_{t+1} = \mathbf{B}_t \underbrace{\left( \mathbf{I}_k - \beta \mathbf{w}_t \mathbf{w}_t^\top \right)}_{\text{prior weight}} + \mathbf{B}_* \underbrace{\frac{\beta}{n} \sum_{i=1}^n \mathbf{w}_{*,i} \mathbf{w}_t^\top}_{\text{signal weight}}$$

- The **prior weight** is rank  $k - 1$ , while the **signal weight** is only rank 1.

## Key observation

The representation can only move closer to  $\text{col}(\mathbf{B}_*)$  in **one** direction on each iteration  $\rightarrow$  not clear that it can eventually reach  $\text{col}(\mathbf{B}_*)$ .

- Let's focus on the population case to simplify the expressions

$$\text{ARM: } \min_{\mathbf{B}, \mathbf{w}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{B}, \mathbf{w})$$

- In this case we can show:

$$\mathbf{B}_{t+1} = \mathbf{B}_t \underbrace{\left( \mathbf{I}_k - \beta \mathbf{w}_t \mathbf{w}_t^\top \right)}_{\text{prior weight}} + \mathbf{B}_* \underbrace{\frac{\beta}{n} \sum_{i=1}^n \mathbf{w}_{*,i} \mathbf{w}_t^\top}_{\text{signal weight}}$$

- The **prior weight** is rank  $k - 1$ , while the **signal weight** is only rank 1.

## Key observation

The representation can only move closer to  $\text{col}(\mathbf{B}_*)$  in **one** direction on each iteration  $\rightarrow$  not clear that it can eventually reach  $\text{col}(\mathbf{B}_*)$ .

## Theorem

For any  $\delta \in (0., 0.5]$ ,  $\alpha$ ,  $T$ ,  $\{\mathbf{w}_{*,i}\}$  and full rank  $\mathbf{B}_0$ , there exists a  $\mathbf{B}_*$  whose column space is  $\delta$ -close to  $\text{col}(\mathbf{B}_0)$ , i.e.,  $\text{dist}(\mathbf{B}_0, \mathbf{B}_*) = \delta$ , while its distance from the representation learned by ARM is at least  $0.7\delta$ , i.e.,  $\text{dist}(\mathbf{B}_T^{\text{ARM}}, \mathbf{B}_*) > 0.7\delta$ .

- For FO-ANIL, we have

$$B_{t+1} = B_t \left( \underbrace{I_k - \frac{\beta}{n} \sum_{i=1}^n \mathbf{w}_{t,i} \mathbf{w}_{t,i}^\top}_{\text{prior weight}} \right) + B_* \underbrace{\frac{\beta}{n} \sum_{i=1}^n \mathbf{w}_{*,i} \mathbf{w}_{t,i}^\top}_{\text{signal weight}}$$

- Suppose  $\frac{1}{n} \sum_{i=1}^n \mathbf{w}_{t,i} \mathbf{w}_{t,i}^\top$  is full rank (i.e., the  $\mathbf{w}_{t,i}$ 's are diverse), then:

### Key observation

**Prior weight** reduces energy from  $B_t$ , and **signal weight** boosts energy from  $B_*$  in all directions.

⇒ **Head adaptation** and **task diversity** are critical!

- The **more diverse** the **adapted**  $\mathbf{w}_{t,i}$ 's (i.e. smaller condition number of  $\frac{1}{n} \sum_{i=1}^n \mathbf{w}_{t,i} \mathbf{w}_{t,i}^\top$ ), the **faster convergence rate**.

- Need to show that the  $\mathbf{w}_{t,i}$ 's are sufficiently diverse, i.e. evenly spread across  $\mathbb{R}^k$ .
- We can show:

$$\begin{aligned}\mathbf{w}_{t,i} &= \mathbf{w}_t - \alpha \nabla f_i(\mathbf{B}_t, \mathbf{w}_t) \\ &= \underbrace{(I_k - \alpha \mathbf{B}_t^\top \mathbf{B}_t) \mathbf{w}_t}_{\text{shared for all } i} + \underbrace{\alpha \mathbf{B}_t^\top \mathbf{B}_* \mathbf{w}_{*,i}}_{\text{unique for each } i}\end{aligned}$$

- We must show the unique part of  $\mathbf{w}_{t,i}$  dominates the shared part.

$$\mathbf{w}_{t,i} = \underbrace{(\mathbf{I}_k - \alpha \mathbf{B}_t^\top \mathbf{B}_t) \mathbf{w}_t}_{\text{shared for all } i} + \underbrace{\alpha \mathbf{B}_t^\top \mathbf{B}_* \mathbf{w}_{*,i}}_{\text{unique for each } i}$$

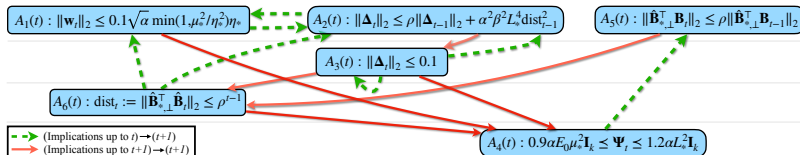
- In order to show the unique part dominates, we must show three things hold for all  $t$ :
  - 1  $\|\mathbf{I}_k - \alpha \mathbf{B}_t^\top \mathbf{B}_t\|_2$  is small
  - 2  $\|\mathbf{w}_t\|_2$  is small
  - 3  $\sigma_{\min}(\mathbf{B}_t^\top \mathbf{B}_*)$  is lower bounded  $\iff \text{dist}(\mathbf{B}_t, \mathbf{B}_*) < 1 - c$  for a positive constant  $c$
- Difficult because the algorithms lack explicit regularization and a normalization step.
- Leads to an intricate 6-way induction....

- Define  $\Delta_t := I_k - \alpha \mathbf{B}_t^\top \mathbf{B}_t$ ,  $\Psi_t := \frac{1}{n} \sum_{i=1}^n \mathbf{w}_{t,i} \mathbf{w}_{t,i}^\top$ , and  $\rho := 1 - \Omega(\beta\alpha)$ .
- Inductive hypotheses:
  - ①  $A_1(t) := \{\|\mathbf{w}_t\|_2 = O(\sqrt{\alpha})\}$  - *global head stays small*
  - ②  $A_2(t) := \{\|\Delta_t\|_2 = \rho\|\Delta_{t-1}\|_2 + O(\beta^2\alpha^2 \text{dist}_{t-1}^2)\}$  -  $\mathbf{B}_t$  gets closer to orthogonal as  $\text{dist}_t$  decreases
  - ③  $A_3(t) := \{\|\Delta_t\|_2 = O(1)\}$  -  $\mathbf{B}_t$  stays uniformly close to orthogonal
  - ④  $A_4(t) := \{\kappa(\Psi_t) = O(1)\}$  - *adapted heads  $\mathbf{w}_{t,i}$  are diverse*
  - ⑤  $A_5(t) := \{\|\mathbf{B}_{*,\perp}^\top \mathbf{B}_t\|_2 = \rho\|\mathbf{B}_{*,\perp}^\top \mathbf{B}_{t-1}\|_2\}$  - *energy of  $\mathbf{B}_t$  in perpendicular space to  $\text{col}(\mathbf{B}_*)$  is contracting*
  - ⑥  $A_6(t) := \{\text{dist}_t \leq \rho^{t-1}\}$  - *PA distance of  $\mathbf{B}_t$  to  $\mathbf{B}_*$  is linearly decreasing*

$$\mathbf{B}_{t+1} = \mathbf{B}_t \left( \mathbf{I}_k - \underbrace{\frac{\beta}{n} \sum_{i=1}^n \mathbf{w}_{t,i} \mathbf{w}_{t,i}^\top}_{\text{prior weight}} \right) + \mathbf{B}_* \underbrace{\frac{\beta}{n} \sum_{i=1}^n \mathbf{w}_{*,i} \mathbf{w}_{t,i}^\top}_{\text{signal weight}}$$

$$\mathbf{w}_{t,i} = \underbrace{(\mathbf{I}_k - \alpha \mathbf{B}_t^\top \mathbf{B}_t) \mathbf{w}_t}_{\text{shared for all } i} + \underbrace{\alpha \mathbf{B}_t^\top \mathbf{B}_* \mathbf{w}_{*,t,i}}_{\text{unique for each } i}$$

- Inductive logic:



Notable implications (1/3):

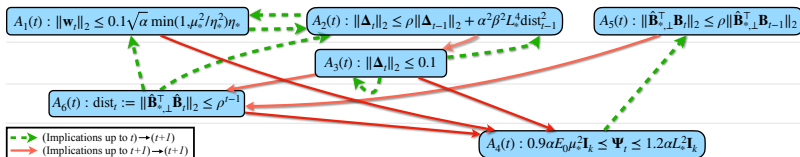
- $A_4(t) \implies A_5(t+1) \xrightarrow{A_3(t+1)} A_6(t+1)$ : Adapted head diversity and well-conditioned  $\mathbf{B}_{t+1}$  implies  $\text{dist}_{t+1}$  is linearly converging.
- Proof.* Use expression for  $\mathbf{B}_{t+1}$  discussed earlier.



$$\mathbf{B}_{t+1} = \underbrace{\mathbf{B}_t \left( \mathbf{I}_k - \frac{\beta}{n} \sum_{i=1}^n \mathbf{w}_{t,i} \mathbf{w}_{t,i}^\top \right)}_{\text{prior weight}} + \underbrace{\mathbf{B}_* \frac{\beta}{n} \sum_{i=1}^n \mathbf{w}_{*,i} \mathbf{w}_{t,i}^\top}_{\text{signal weight}}$$

$$\mathbf{w}_{t,i} = \underbrace{(\mathbf{I}_k - \alpha \mathbf{B}_t^\top \mathbf{B}_t) \mathbf{w}_t}_{\text{shared for all } i} + \underbrace{\alpha \mathbf{B}_t^\top \mathbf{B}_* \mathbf{w}_{*,i}}_{\text{unique for each } i}$$

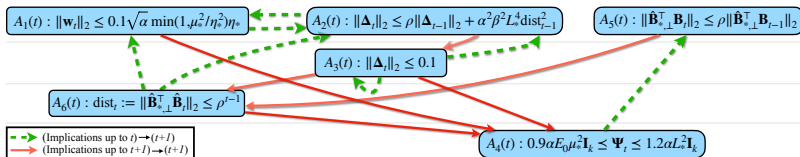
## • Inductive logic:



## Notable implications (2/3):

- $A_1(t+1), A_3(t+1), A_6(t+1) \implies A_4(t+1)$ : Small  $\|\mathbf{w}_t\|_2$ ,  $\|\Delta_t\|_2$ , and  $\text{dist}_t$  implies adapted heads are diverse.
  - *Proof.* Use expression for  $\mathbf{w}_{t,i}$  discussed earlier.

- Inductive logic:



## Notable implications (3/3):

- $A_2(t) + A_6(t) \implies A_1(t+1)$ : Bounds on  $\|\Delta_t\|_2$  and  $\text{dist}_t$  imply  $\|\mathbf{w}_{t+1}\|_2$  stays small.
  - Proof.*  $A_2(t)$  and the linear convergence of  $\text{dist}_t$  implies  $\|\Delta_t\|_2$  eventually linearly converges to zero. We can show  $\|\mathbf{w}_{t+1}\|_2 \leq \|\mathbf{w}_t\|_2 + O(\|\Delta_t\|_2)$ , which implies  $|\|\mathbf{w}_{t+1}\|_2 - \|\mathbf{w}_t\|_2|$  eventually linearly converges to zero.

### Main Theorem [Collins-M-Oh-Shakkottai, ICML 2022]

Suppose there are  $m = \infty$  samples/task, the ground-truth heads satisfy  $\mu_*^2 \mathbf{I}_k \preceq \frac{1}{n} \sum_{i=1}^n \mathbf{w}_{*,t,i} \mathbf{w}_{*,t,i}^\top \preceq L_*^2 \mathbf{I}_k$ , and the step sizes  $\alpha, \beta$  are sufficiently small. Then after  $T$  iterations, ANIL, FO-ANIL, MAML, and FO-MAML learn a representation  $\mathbf{B}_T$  that satisfies:

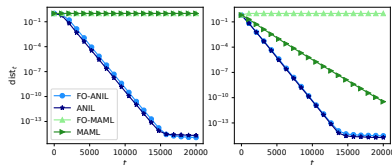
$$\text{dist}(\mathbf{B}_T, \mathbf{B}_*) \leq (1 - \Omega(\beta \alpha \mu_*^2))^{T-1}$$

as long as:

- ANIL, FO-ANIL:  $\text{dist}(\mathbf{B}_0, \mathbf{B}_*) \leq c$  for a constant  $c$ .
- MAML:  $\text{dist}(\mathbf{B}_0, \mathbf{B}_*) = O((L_*/\mu_*)^{-0.75})$ .
- FO-MAML:  $\text{dist}(\mathbf{B}_0, \mathbf{B}_*) = O((L_*/\mu_*)^{-1})$  and  $\|\frac{1}{n} \sum_{i=1}^n \mathbf{w}_{*,t,i}\|_2 = O((L_*/\mu_*)^{-1.5})$ .

- We also show finite-sample results in the paper.

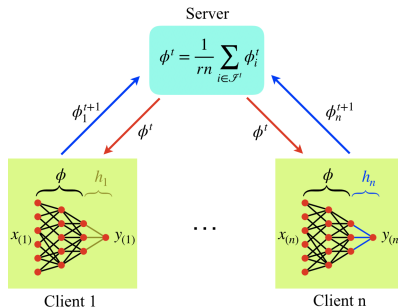
- Recall that our result requires
  - stronger initialization for MAML and FO-MAML than for ANIL and FO-ANIL, and
  - for FO-MAML, the mean of  $\mathbf{w}_{*,i}$ 's must be  $\approx 0$ .
- Here we show these conditions are tight:



(Left) Random initialization and (Right) structured initialization. In both cases, the mean of the ground-truth heads is far from zero, explaining why FO-MAML fails. MAML succeeds only with structured initialization.

⇒ MAML/FO-MAML's inner loop update of the representation can inhibit representation learning.

- We have obtained the **first results** showing that ANIL and MAML learn effective representations in any setting.
- Inner loop **adaptation of the head** is **key** to MAML and ANIL's ability to learn representations.
- Inner loop adaptation of the representation **restricts representation learning**.
- So far we have seen that by adding adaptivity via changes in the function we can learn the representation
- Next, we show that by changing algorithms in the training process that leads to adaptivity we can also learn the representation



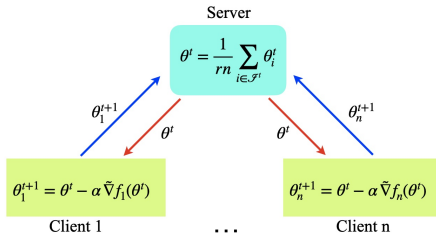
- **Step 1:** Train a **single model** by minimizing the average loss

$$\theta^* \in \operatorname{argmin}_{\theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta; \mathcal{D}_i)$$

- **Step 2:** **Fine tune** the obtained model  $\theta^*$  to new task / deployment environment



- Popular approach: **Distributed Stochastic Gradient Descent (DSGD)**



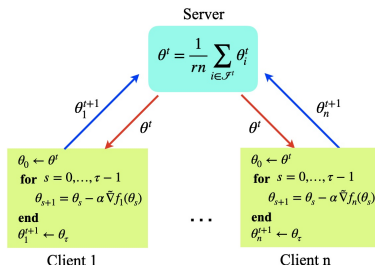
- Effective update  $\theta^{t+1} = \theta^t - \frac{\alpha}{rn} \sum_{i \in \mathcal{I}^t} \tilde{\nabla} f_i(\theta^t)$   
 $\Rightarrow$  running SGD on the global loss

Same Loss as ARM – Only change is a Distributed Implementation

$$\theta^* \in \operatorname{argmin}_{\theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta; \mathcal{D}_i)$$



- Federated learning is similar to DSGD
  - Except it allows for **multiple local updates**
- FedAvg: the most common approach!



- Advantage:** Number of communication rounds  $\ll$  number of updates
  - More communication-efficient than distributed SGD!
- Issue:** Drifting effect with data heterogeneity (aka multiple tasks)
  - The local grads  $\nabla f_i$  are not aligned with the global direction  $\sum_{i=1}^n \nabla f_i$

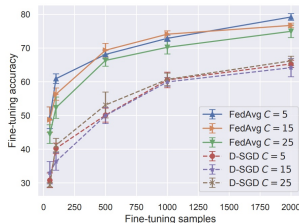
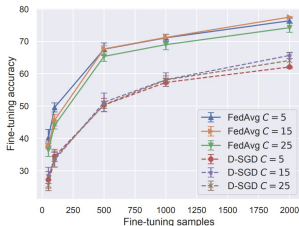
From an **optimization point of view**:

- Due to local drifting, FedAvg may not solve global objective in the data heterogeneous setting
  - [WJ19], [CK21], ....
- There are several other sophisticated methods to control/remove the local drift issue in FedAvg
  - SCAFFOLD: [KKMRSS20]
  - VRL-SGD: [LLZSM20]
  - FEDGATE: [HKMM21]
  - FedNova: [WLLJP20]
  - .....

## Our Takeaway

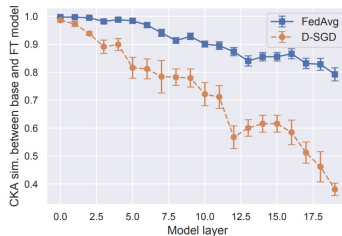
**Implicitly, the new implementation has changed the objective**

- The local updates are in fact good for **generalization**!



- **Left plot:** Models trained on 80 classes from CIFAR-100 (with  $C$  classes/client) and fine-tuned on new clients from 20 new classes from CIFAR-100
- **Right plot:** Models trained on CIFAR-100 (with  $C$  classes/client) and fine-tuned on new clients from CIFAR-10
- $T\tau = 125000$  for both. (FedAvg  $\tau = 50$ ,  $T = 2500$ , DSGD  $\tau = 1$ ,  $T = 125000$ )

- The early layers of FedAvg's pre-trained model (corresponding to the representation) change much less than those of D-SGD



- Local updates enable learning the common representation among the tasks!

## Open Problems

- *Does FedAvg provably learn representations of heterogeneous tasks?*
- *Are the local updates essential for learning representations?*

- Multi-task linear regression.
- Task  $i$  has ground-truth solution  $\theta_{*,i} \in \mathbb{R}^d$ :

$$y_i \sim \theta_{*,i}^\top \mathbf{x}_i + z_i$$

- $\mathbf{x}_i$  is a random feature vector,  $z_i \in \mathbb{R}$  is random, mean-zero noise
- Now suppose the  $\theta_{*,i}$  lie in a shared  $k$ -dimensional subspace,  $k \ll d$
- Let the columns of  $\mathbf{B}_* \in \mathbb{R}^{d \times k}$  span this subspace, that is, for all tasks there exists  $\mathbf{w}_{*,i} \in \mathbb{R}^k$  such that

$$\theta_{*,i} = \mathbf{B}_* \mathbf{w}_{*,i}$$

**Task Diversity:** The concatenation of  $\mathbf{w}_{*,1}, \dots, \mathbf{w}_{*,n}$  spans  $\mathbb{R}^k$

- Now suppose  $\tau = 1$  and we have the update of DGD

$$\mathbf{B}_{t+1} = \mathbf{B}_t \underbrace{\left( \mathbf{I}_k - \alpha \mathbf{w}_t \mathbf{w}_t^\top \right)}_{\text{prior weight}} + \alpha \mathbf{B}_* \underbrace{\left( \frac{1}{n} \sum_{i=1}^n \mathbf{w}_{*,i} \right) \mathbf{w}_t^\top}_{\text{signal weight}}$$

- The **prior weight** is rank  $k - 1$ , while the **signal weight** is only rank 1.
- The representation can only move closer to  $\text{col}(\mathbf{B}_*)$  in **one** direction on each iteration

### Theorem

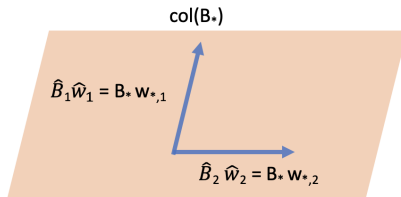
For any  $\delta \in (0., 0.5]$ ,  $\alpha$ ,  $T$ ,  $\{\mathbf{w}_{*,i}\}$  and full rank  $\mathbf{B}_0$ , there exists a  $\mathbf{B}_*$  whose column space is  $\delta$ -close to  $\text{col}(\mathbf{B}_0)$ , i.e.,  $\text{dist}(\mathbf{B}_0, \mathbf{B}_*) = \delta$ , while its distance from the representation learned by DGD is at least  $0.7\delta$ , i.e.,  $\text{dist}(\mathbf{B}_T^{\text{DGD}}, \mathbf{B}_*) > 0.7\delta$ .

- DGD cannot guarantee to recover the ground-truth representation.

- Recall that the local updates are with respect to the local loss:

$$f_i(\mathbf{B}, \mathbf{w}) = \frac{1}{2} \|\mathbf{B}\mathbf{w} - \mathbf{B}_* \mathbf{w}_{*,i}\|_2^2$$

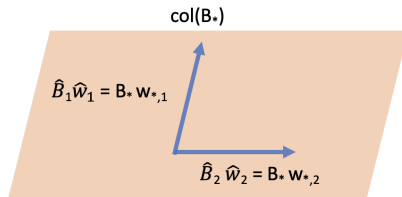
- Let  $\hat{\mathbf{B}}_i, \hat{\mathbf{w}}_i$  be the result of  $\tau$  local updates for client  $i$
- (Naive) idea: Use fact that if  $\hat{\mathbf{B}}_i \hat{\mathbf{w}}_i = \mathbf{B}_* \mathbf{w}_{*,i}$  for all  $i$ , then the local products  $\{\hat{\mathbf{B}}_i \hat{\mathbf{w}}_i\}_i$  all lie in the correct subspace  
 $\Rightarrow$  and  $\text{span}(\{\hat{\mathbf{B}}_i \hat{\mathbf{w}}_i\}_i) = \text{col}(\mathbf{B}_*)$



- Recall that the local updates are with respect to the local loss:

$$f_i(\mathbf{B}, \mathbf{w}) = \frac{1}{2} \|\mathbf{B}\mathbf{w} - \mathbf{B}_* \mathbf{w}_{*,i}\|_2^2$$

- Let  $\hat{\mathbf{B}}_i, \hat{\mathbf{w}}_i$  be the result of  $\tau$  local updates for client  $i$
- (Naive) idea: Use fact that if  $\hat{\mathbf{B}}_i \hat{\mathbf{w}}_i = \mathbf{B}_* \mathbf{w}_{*,i}$  for all  $i$ , then the local products  $\{\hat{\mathbf{B}}_i \hat{\mathbf{w}}_i\}_i$  all lie in the correct subspace  
 $\Rightarrow$  and  $\text{span}(\{\hat{\mathbf{B}}_i \hat{\mathbf{w}}_i\}_i) = \text{col}(\mathbf{B}_*)$



- However, this does not imply anything meaningful about  $\mathbf{B}_{t+1} := \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{B}}_i$



- Let  $k = n = 2$ , and suppose  $\theta_{*,1} = \mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$ ,  $\theta_{*,2} = \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$   
 $\implies \text{col}(\mathbf{B}_*) = \text{span}(\mathbf{e}_1, \mathbf{e}_2)$

- Let  $k = n = 2$ , and suppose  $\theta_{*,1} = \mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$ ,  $\theta_{*,2} = \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$

$$\implies \text{col}(\mathbf{B}_*) = \text{span}(\mathbf{e}_1, \mathbf{e}_2)$$

- Let  $\hat{\mathbf{B}}_1 = \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \end{bmatrix}$ ,  $\hat{\mathbf{w}}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$  and  $\hat{\mathbf{B}}_2 = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ \vdots & \vdots \\ 0 & 0 \end{bmatrix}$ ,  $\hat{\mathbf{w}}_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

- Then  $\hat{\mathbf{B}}_1 \hat{\mathbf{w}}_1 = \mathbf{B}_* \mathbf{w}_{*,1}$  and  $\hat{\mathbf{B}}_2 \hat{\mathbf{w}}_2 = \mathbf{B}_* \mathbf{w}_{*,2}$ , yet

$$\mathbf{B}_{t+1} = \frac{1}{2} \hat{\mathbf{B}}_1 + \frac{1}{2} \hat{\mathbf{B}}_2 = [\frac{1}{2}(\mathbf{e}_1 + \mathbf{e}_2), \mathbf{0}]$$

has  $\text{dist}(\mathbf{B}_{t+1}, \mathbf{B}_*) = 1 \dots$

$\Rightarrow$  Can't rely only on local convergence!

## Theorem (informal) [Collins-Hassani-M-Shakkottai, NeurIPS 2022]

*If the number of local updates satisfies  $\tau \geq 2$ , FedAvg recovers  $\text{col}(\mathbf{B}^*)$  exponentially fast when run on the task population losses.*

- The key insight is that FedAvg local updates harnesses **task diversity** to improve the representation in all directions.

$$\mathbf{B}_{t+1} = \mathbf{B}_t \underbrace{\left( \frac{1}{n} \sum_{i=1}^n \prod_{s=0}^{\tau-1} \left( \mathbf{I}_k - \alpha \mathbf{w}_{t,i,s} \mathbf{w}_{t,i,s}^\top \right) \right)}_{\text{prior weight}} + \underbrace{\mathbf{B}_* \left( \frac{\alpha}{n} \sum_{i=1}^n \mathbf{w}_{*,i} \sum_{s=0}^{\tau-1} \mathbf{w}_{t,i,s}^\top \prod_{r=s+1}^{\tau-1} \left( \mathbf{I}_k - \alpha \sum_{i=1}^n \mathbf{w}_{t,i,r} \mathbf{w}_{t,i,r}^\top \right) \right)}_{\text{signal weight}}$$

- Prior weight** reduces energy from  $\mathbf{B}_t$ , and **signal weight** boosts energy from  $\mathbf{B}_*$  in all directions

## Theorem (informal) [Collins-Hassani-M-Shakkottai]

*If the number of local updates satisfies  $\tau \geq 2$ , FedAvg recovers  $\text{col}(\mathbf{B}^*)$  exponentially fast when run on the task population losses.*

- The key insight is that FedAvg local updates harnesses **task diversity** to improve the representation in all directions.

$$\mathbf{B}_{t+1} \approx \mathbf{B}_t \underbrace{\left( \mathbf{I}_k - \frac{\alpha}{n} \sum_{i=1}^n \sum_{s=0}^{\tau-1} \mathbf{w}_{t,i,s} \mathbf{w}_{t,i,s}^\top \right)}_{\text{prior weight}} + \mathbf{B}_* \underbrace{\left( \frac{\alpha}{n} \sum_{i=1}^n \sum_{s=0}^{\tau-1} \mathbf{w}_{*,i} \mathbf{w}_{t,i,s}^\top \right)}_{\text{signal weight}}$$

- Prior weight** reduces energy from  $\mathbf{B}_t$ , and **signal weight** boosts energy from  $\mathbf{B}_*$  in all directions
  - Local updates** and **task diversity** are critical!

- **Client Diversity Assumption:** Let  $\mathbf{W}_* := [\mathbf{w}_{*,1}, \dots, \mathbf{w}_{*,n}]$ . Then  $\sigma_{\min,*} := \sigma_{\min}(\mathbf{W}_*) > 0$ .

### Theorem [Collins-Hassani-M-Shakkottai, NeurIPS '22]

*If the Client Diversity Assumption holds, the number of local updates between rounds satisfies  $\tau \geq 2$ , the step size  $\alpha = O(\frac{1}{\sqrt{\tau}})$ , and the initialization satisfies  $\text{dist}(\mathbf{B}_0, \mathbf{B}_*) \leq 1 - c$  for some  $c \in (0, 1]$ , then for any  $\epsilon \in (0, 1)$ , FedAvg learns a representation  $\mathbf{B}_T$  satisfying  $\text{dist}(\mathbf{B}_T, \mathbf{B}_*) \leq \epsilon$  after at most*

$$T = O\left(\frac{\log(1/\epsilon)}{\alpha^2 \tau \sigma_{\min,*}^2}\right)$$

*communication rounds.*

- **Client Diversity Assumption:** Let  $\mathbf{W}_* := [\mathbf{w}_{*,1}, \dots, \mathbf{w}_{*,n}]$ . Then  $\sigma_{\min,*} := \sigma_{\min}(\mathbf{W}_*) > 0$ .

### Theorem [Collins-Hassani-M-Shakkottai, NeurIPS '22]

*If the Client Diversity Assumption holds, the number of local updates between rounds satisfies  $\tau \geq 2$ , the step size  $\alpha = O(\frac{1}{\sqrt{\tau}})$ , and the initialization satisfies  $\text{dist}(\mathbf{B}_0, \mathbf{B}_*) \leq 1 - c$  for some  $c \in (0, 1]$ , then for any  $\epsilon \in (0, 1)$ , FedAvg learns a representation  $\mathbf{B}_T$  satisfying  $\text{dist}(\mathbf{B}_T, \mathbf{B}_*) \leq \epsilon$  after at most*

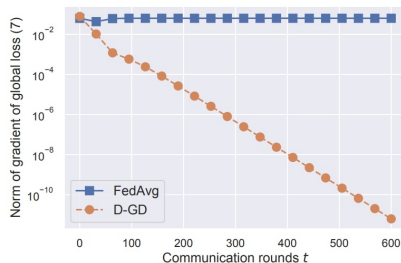
$$T = O\left(\frac{\log(1/\epsilon)}{\alpha^2 \tau \sigma_{\min,*}^2}\right)$$

*communication rounds.*

- **First result** showing that FedAvg learns an expressive representation in any setting!

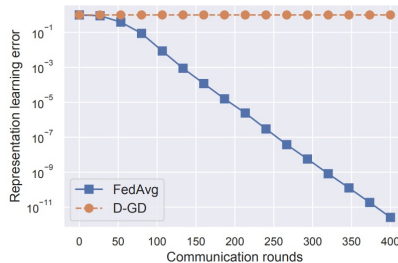
- Optimization point of view

- DGD finds a stationary point of the average loss, while FedAvg does not



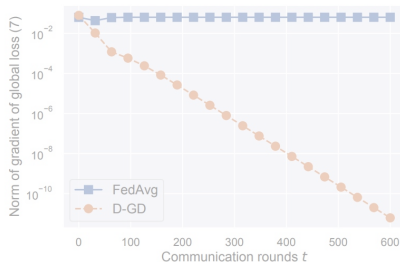
- Generalization point of view

- FedAvg learns the correct representation (principal angle distance)



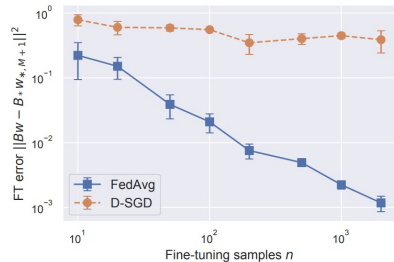
## • Optimization point of view

- DGD finds a stationary point of the average loss, while FedAvg does not



## • Generalization point of view

- FedAvg generalizes better (loss with new task, after fine tuning)!





- [FAL17] C. Finn, P. Abbeel, S. Levine. Model-Agnostic Meta-Learning for Fast Adaptation of Neural Networks, *ICML*, 2017.
- [AES19] A. Antoniou, H. Edwards, A. Storkey. How to Train Your MAML, *ICLR*, 2019.
- [RRBV19] A. Raghu, M. Raghu, S. Bengio, O. Vinyals. Rapid Learning or Feature Reuse? Towards Understanding the Effectiveness of MAML, *ICLR*, 2020.
- [HRJ21] M. Huisman, J. N. van Rijn, A. Plaat. A Survey of deep Meta-Learning, *Artificial Intelligence Review* Volume 54, pages 4483-4541, 2021.
- [LLZSM20] Liang, Liu, Ziyin, Salakhutdinov, and Morency. Think locally, act globally: Federated learning with local and global representations, 2020.
- [WLLJP20] Wang, Liu, Liang, Joshi, and Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization, *NeurIPS* 2020.
- [CK21] Charles, Konecny. Convergence and accuracy trade-offs in federated learning and meta-learning, *AISTATS* 2021.
- [WJ18] Wang, Joshi. Cooperative SGD: A unified framework for the design and analysis of communication-efficient SGD algorithms, 2018.
- [KKMRSS20] Karimireddy, Kale, Mohri, Reddi, Stich, and Suresh. Scaffold: Stochastic controlled averaging for federated learning, *ICML* 2020.
- [HKMM21] Haddadpour, Kamani, Mokhtari, Mahdavi. Federated Learning with Compression: Unified Analysis and Sharp Guarantees, *AISTATS* 2021.