

Federated Fine-tuning of LLMs

state of the art: FedKseed FlexLoRA Fed PFT

→ reduce communication and handle

heterogeneous client resources → aggregate by arithmetic averaging

GRAP → No algorithm incorporates semantic or linguistic

similarity of clients →

all clients are treated equally, even if their languages on text domains are very different

This causes degraded global models and poor cross lingual transfer.

→ no mathematical or algorithmic mechanism exploits these distance aggregation

*) Prompt tuning & Personalization in FL

FedPGP, FedTPG, FedBPT → focus on communication efficiency and personalization

→ Aggregation is still parameter-space averaging → no geometry or transport alignment across client.

→ Theoretical Foundations:

NeurIPS'24 prompt - pert

→ no theoretical bound connecting multilingual heterogeneity, privacy noise and aggregation error.

Privacy vs accuracy

DP + Text and other DP works handle privacy by adding noise on generating synthetic data.

None analyze how DP noise interacts with linguistic heterogeneity on purpose a frame work that combines

One round of Lingua-OT-FL

1. Clients

compute local updates

$$\Delta_k = \nabla_{\theta} L_k(\Phi)$$

Every M_i model learns by minimizing a loss function $L(\theta)$:

$$\min_{\theta} E_{(x,y) \sim P_{\text{data}}} \ell(f_{\theta}(x), y)$$

$f_{\theta}(x)$ model prediction

$L(\cdot) \rightarrow$ loss function

$P_{\text{data}} \rightarrow$ underlying data distribution

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(f_{\theta}(x_i), y_i)$$

use gradient descent

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} L(\theta_t)$$

From centralized to federated

$$L(\theta) = \sum_{k=1}^K \frac{n_k}{N} L_k(\theta)$$

$$L_k(\theta) = \frac{1}{n_k} \sum \ell(f_{\theta}(x), y)$$

Gradient of the global loss

$$\nabla_{\theta} L(\theta) = \sum_{k=1}^K \frac{n_k}{N} \nabla_{\theta} L_k(\theta)$$

3) The FedAvg Algorithm

$$\text{local } \theta_k^{t+1} = \theta_t - \eta \nabla_{\theta} L_k(\theta_t)$$

(often multiple mini batches)

$$\theta_{t+1} = \sum_{k=1}^K \frac{n_k}{N} \theta_k^{t+1}$$

Demotion

The NonIID challenge formally

4. Let's call the local optimum on client's

θ_k^* = arg min $L_k(\theta)$

If all clients had identical data

θ_k^* would be the same

under non IID conditions

$$\|\theta_i^* - \theta_j^*\| > 0$$

distance \rightarrow data heterogeneity

Federated optimization as a

communication problem

centralized SGD

Update only after k clients

train locally for E epochs.

Date: _____

Sun Mon Tue Wed Thu Fri Sat

two time scale
1) Local updates (with clients)
2) Global aggregation (server round)

This saves communication but risks model drift.

Federated fine tuning for NLP

Do (each client text data)
diff language diff domain.

They all share pre trained backbone
for (like bert)

Instead of tuning all weights &
clients only train small parameters
 $\theta_k \rightarrow$ prompts, LoRA adapters
on prefix vectors.

$$f(x; \theta_0, \phi_k) = \theta_0(x) + g_{\phi_k}(x)$$

Federations. Happens over ϕ_k instead of huge θ_0

The server aggregates ϕ_k 's

$$\phi_{\text{global}} = \sum_{k=1}^N \frac{m_k}{N} \phi_k$$

Sends them back to client

Thus parameter efficient FL

= small communication & local memory
faster convergence

Date:

Sun Mon Tue Wed Thu Fri Sat

⑦ Differential Privacy

each client adds a gaussian noise

$$N(0, \sigma^2 I)$$

$$\theta_k^{t+1} = \theta_k^{t+1} + N(0, \sigma^2 I)$$

the server aggregates these noisy updates.

~~or~~

This ensures that any single sample's influence on the model is statistically limited.

Formally the algorithm is (ϵ, δ) -DP



Differential privacy

any two dataset differing one example

$$P_n [A(D_1) \neq A(D_2)] \leq e^2 \Pr[A(D_2) \neq S] + \delta$$

Observation cannot tell whose data is like

8) Theoretical convergence:

when all IID are smooth

$$E[L(\theta_T)] - L^* = O\left(\frac{1}{T}\right)$$

fed Avg \rightarrow converges roughly

like centralized

Stochastic gradient
Descent

Date:

Sun Mon Tue Wed Thu Fri Sat

under non IID data

the rate becomes

$$O\left(\frac{1}{T} + \text{Bias}_{\text{nonIID}}\right)$$

Dis depends on how far each
local optimum θ_i^* deviates from
the global one.

This \rightarrow slower convergence for
heterogeneous data.

For smarter aggregation

language structure introduce another
kind of heterogeneity

Adding language distance

$$d_{ij} = |S_i - S_j|$$

instead of plain averaging
weight updates using a decreasing
function of distance.

$$w_k = \frac{e^{-\text{rad}(w, \text{global})}}{\sum_j e^{-\text{rad}(j, \text{global})}}$$

Aligning updates with optimal

transport

→ minimizes the effort to move
one distribution into another

In lingua-OT-FL →

align updates from different languages
so that semantically similar
gradients correspond before averaging

solver solves

$$\Delta^{\text{OT}} = \arg \min_{\Delta} \sum_k w_k \tilde{w}_k(\Delta, \Delta_k)$$

w_k is Wasserstein barycenter.

→ clients → compute local update

$$\Delta_k = \nabla_P L_k(\Phi)$$

and language signature S_k

servers
compute pairwise distance

$$d_{ij} = |S_i - S_j|$$

Demelon

compute GT baycenter

$$\bar{\Delta}^{OT} = \argmin_{\Delta} \sum w_k w_k^T (\Delta, A_k)$$

3. add Differential privacy

$$\bar{\Delta}^{OT-DP} = \bar{\Delta}^{OT} + N(0, \sigma^2 I)$$

4. Global update

$$\phi_{global} \leftarrow \phi_{global} - \eta \bar{\Delta}^{OT-DP}$$

Personalization

$$\phi_k = \lambda \phi_{global} + (1-\lambda)$$

Date:

Sun Mon Tue Wed Thu Fri Sat

Co-dop

Optimal Transport Theory

Standard Fed Avg \rightarrow just takes

a weighted mean of client updates

\rightarrow That's fine if everyone's updates

"point" roughly in the same direction

But multilingual or multidomain

NLP:

QA \rightarrow Bangla - update "focus on vowel normalization"

QB \rightarrow English update - handle

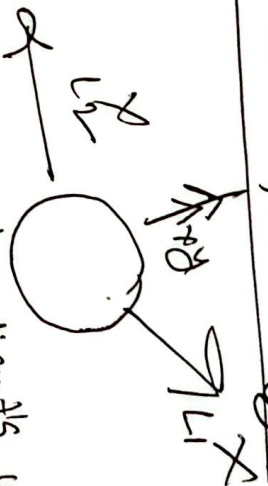
foreign consistency

QC \rightarrow Arabic handle script diff

Their gradients point in diff direction
in parameter space
Plain averaging can cancel
useful signals out
so we need a smarter way to align
these updates before combining them.

OT enters

The transport analogy



Imagine each Client's update as
a pile of sand spread across a
landscape. (The parameter space)
To merge them \rightarrow we ~~can't~~ want
to move each pile into a single
"average" \rightarrow but moving sand
costs effort

The question OT answers:

→ what is the most efficient way to move (each pile into a single average) mass from one distribution to another?

→ The effort is measured by how far and how much ground is moved.

Say we have two probability distributions

$p(x)$ and $q(y)$

A transport plan $\gamma(x,y)$ tells us how much mass to move from x to y .

The transport cost is:

$$\text{cost}(p, q, \gamma) = \iint \underbrace{c(x,y)}_{\text{endy}} \gamma(x,y)$$

$c(x,y)$ → cost of moving unit mass (usually euclidean distance)

The optimal transport distance (also called as Wasserstein distance) is the minimum cost over all valid γ .



$$\text{Total cost } (P, Q, \gamma) =$$

$$C \gamma(m, y) \star c(m, y) + x_2(m, y) \star c(m, y)$$

So minimum cost OT distance

$$W_2(P, Q) = \min_{\gamma} \int \|x - y\| \gamma(x, y) dx dy$$

This gives a geometry aware notion "distance" between

~~distance~~ distributions

Date:

Sun Mon Tue Wed Thu Fri Sat

From two Distributions to many:

The Barycenter:

If you have several clients with

distances p_1, p_2, \dots, p_k .

$$P^* = \arg \min_p \sum_k W_2(P, p_k)$$

minimizes the weighted transport

cost

gives the "most central" distribution in the OT sense.

The one requiring the least overall movement from all clients.

Linguistic weights.

we assign each client a weight

$$-ad(P_k, P_{global})$$

$$w_k = \frac{e^{-ad(P_k, P_{global})}}{\sum_j e^{-ad(P_j, P_{global})}}$$

where $d(i, j)$ is the language distance.

→ computed from

embeddings, phoneme stats,
or perplexity between arms

language models

languages that are closer to the
global center influence aggregation
more.

this lingua part of lingua OTFL

Differential privacy in OT

aggregation

→ After computing OT baycenter

$$\bar{D}_{OT-DP} = D_{OT} + \mathcal{N}(\sigma, \sigma^2)$$

noise scale σ controls the

privacy level

Demelson

Co-dopa

Personalization After OT

$$\phi_x^{new} \geq \lambda \bar{\phi}_{global} + (1-\lambda) \phi_x^{local}$$

$$\lambda_{n+2} = \frac{1}{1 + \beta \Delta LR, \beta \phi_{global}}$$

Computing Efficiency

Sun Mon Tue Wed Thu Fri Sat

Date :