Networked AI Learning

A non-federated learning approach



Dr. Wen Tong CTO, Huawei Wireless

International Workshop on Federated Learning and Foundation Models for Multi-Media Keynote July 15th 2024 Niagara Falls, CANADA

Al Model Complexity (1)









Al Model Complexity (2)







2030



Al Model Scaling Law *More intelligence "for free" by Scaling*





Training Dataset Scaling

Kolmogorov Compression



The more data the more intelligence



Wen Tong and Yiqun Ge "Information Theory and Learning-AN INFORMATION THEORY VIEW ON LEARNING PROBLEMS" Cambridge University Press 2022



Model Fine-tuning In-Network LoRA





Federated Learning Share model not the Data





Issue of Federated Learning (1)

The Need of a Trusted Host



Federated:

Set up a single controlled center within which each state division keeps some internal autonomy

The trusted host to share the model



Issue of Federated Learning (2)

The Communications Cost

Local data



The total communication cost of *10 participants* sending *Llama2-7B* to a server in *100 rounds* of FL is:

"14 GB * 10 * 100 * 2 = 28 TB"



Distributed Learning

Model-Follow-Data Concept

Homogenous and Heterogenous Models







Model Communications (1)

Model-Follow-Data Concept



1

AI packet

Data In

Model Out



2



3

AI packet

S: Source D: Destination



Model Communications (2)

Model-Follow-Data Concept



AI Packet Routing





Compute Data Interest



Update Routing Algorithm



nternet	Routing	Table
---------	---------	-------

Destination	Subnet Mask	Gateway
XXXX	XXXX	XXXX
	-	

- Publish the data availability
- IP protocol for model forwarding

- Self training data cleaning
- Data quality criteria
- Option of distributed compute powe

Express of interestsOBSF flooding



Model Communications (3)

Model-Follow-Data Concept



Model Version Management



Autonomous Model Training

2





Model Distributed Ledger System



Base Model Version

Incremental Update Versior

Discovery of data

- Discovery of compute
- Model verification

Version validation Version audit



Heterogenous Model Communications

Model-Follow-Data Concept





Standard Model - Global Model

Heterogenous Model Scaling with Local Data Training





In-Network Model Processing

Heterogenous Model-Re-Normalization and Updates





Distributed Learning Algorithm – (Knowledge Distilling)

Heterogenous Model





Distributed Learning Algorithm – (Generative Model)

Heterogenous Model











(2)

 $\{\mu_i, \pi_i, \Sigma_i\}$

Distributed Learning Algorithm – (Information Bottleneck)

Heterogenous Model

Local data



66







Distributed Learning Algorithm – (LoRA)

Heterogenous Model





Distributed Learning Algorithm – (LORA)

Heterogenous Model



Topology: 10 Clients, 5 clients participating in each round Dataset: CIFAR10 – 25k images as train, 25k as test Distribution: Heavy label shift, Dirichlet(0.1) One base model: ViT-base (86M), pretrained on ImageNet21K images



Federated and Distributed Comparison

Reduced Communications Cost

Federated Learning

Distributed Learning

- Knowledge Distilling
- Generative
- Information Bottleneck
- LoRA

4n 2n 2n 2n

100n



Federated and Distributed Algorithms

Remove Inferior Model Contributions--how to un-learn

Federated Learning

$$\sum \left(\dots \left(\sum \widetilde{W}_{n,i} \right) \sum \Delta \widetilde{W}_{n,i} \right)$$

Distributed Learning

$$\left(\left((W_0 + A_1B_1) + A_2B_2\right) + A_3B_3\right) + \cdots\right)$$



Data and Computing Aware Model Routing

Model-Follow-Data Concept

Conventional:

- Randomly chosen next hop
- Neither the model architecture nor the data characteristics have any impact on the decision

Ours:

- Next hop selection based on new metric
- The model node obtains random minibatches from network nodes
- The model node computes a metric for each node based on its minibatch and current state of the model.



Distributed Training for RAN

3GPP Architecture

OTT Model Plane RAN Model Plane gNB Model Plane

Model Based Communications



Model for Efficient Communications





Model for Reliable Communications



Model for Secure Data Translator





Model for Data Memory





Model as Trustworthy User ID





(Wen Tong, Keynote IEEE-CTW-2021)

A-RAN[™]





Neuron Agent (10B)

Neuron Edge (100B)

Neuron Center (1T)

(Wen Tong, Keynote 6G SUMMIT Abu Dhabi Nov.16th, 2023)

Thank You.

Copyright©2016 Huawei Technologies Co., Ltd. All Rights Reserved.

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. Huawei may change the information at any time without notice.