

Distributed ML System for Large-scale Models: *Dynamic* Distributed Training and *Scalable* Federated Learning





Chaoyang He PhD student (2018-present), CS, USC Former R&D Manager and Staff Software Engineer at Tencent



Mahdi Soltanolkotabi Associate Professor CS & ECE, USC





Salman Avestimehr Professor, ECE&CS, USC Director, USC-Amazon Trusted ML Center

Fundamental Research to Alexa AI Applications



1. Distributed Training in GPU cluster 2. Federated Learning at the edge

Distributed Training v.s. Federated Learning





In the Federated Learning setting, the data is distributed across millions of devices in a highly uneven fashion. These devices have low computational compatibility, significantly higher-latency, lower-throughput connections

---- Federated Learning and Open Problems, Foundations and Trends in Machine Learning (FnTML) 2021; ..., <u>Chaoyang He</u>, and many Googlers

Distributed Training v.s. Federated Learning



Goals: training large model in reasonable time, energy cost, and hardware resources to obtain a higher accuracy

Challenge:

- 1. communication cost
- 2. memory constraints
- 3. computational efficiency
- 4. straggler/failure



Goals: data and model privacy; mitigate the difficulty of data centralization; ubiquitous ML in 5G/IoT future.

Challenges:

- 1. data heterogeneity and label deficiency
- 2. system heterogeneity; resource-constraint; scalability
- 3. trustworthiness: security, privacy, fairness, personalization, interpretability, etc

Distributed Training v.s. Federated Learning





In essence, three methodologies are useful to both:

1. Dynamic system and stochastic ML lead to complicated interaction between these objectives. *PipeTransformer, ICML 2021*

2. On demand: focus on the key demand for different use cases (ML applications), optimizing one of the objectives but not deteriorating the others.

^{3.} FedML Lib and Ecosystem, NeurIPS 2020 FL workshop, Best Paper Award

3. Trade-off among multiple objectives including accuracy, efficiency, security, and privacy.

Group Knowledge Transfer, NeurIPS 2020

Outline: Distributed ML for Large-scale Models



Part 1: Dynamic Distributed Training
 PipeTransformer: Automated Elastic Pipelining for Distributed
 Training of Large-scale Models. ICML 2021
 <u>https://DistML.ai</u>



 Part 2: Scalable Federated Learning
 [Algorithm] Group Knowledge Transfer: Federated Learning of Large CNNs at the Edge. NeurIPS 2020

[System] FedML: A Research Library and Benchmark for Federated Machine Learning (Best Paper Award, NeurIPS 2020 FL Workshop)

[Diverse Applications] FedML Ecosystem: Ubiquitous Distributed Training for Diverse AI Applications at the Edge. <u>https://FedML.ai</u>





PipeTransformer: Automated Elastic Pipelining for Distributed Training of Large-scale Models



Chaoyang He PhD student (2018-present), CS, USC Former R&D Manager at Tencent Researcher, Tencent Al Lab

Shen Li Research Scientist, Facebook Al Team Lead, PyTorch Distributed CS PhD, UIUC



Mahdi Soltanolkotabi Associate Professor CS, ECE, USC Dire

Salman Avestimehr Professor, ECE&CS, USC Director, USC-Amazon ML Center

https://DistML.ai https://chaoyanghe.com/pipetransformer/ ○ PyTorch Get Started Ecosystem ∨ Mobile Blog Tutorials Docs ∨ Resources ∨

August 18, 2021

PipeTransformer: Automated Elastic Pipelining for Distributed Training of Large-scale Models

O by Chaoyang He, Shen Li, Mahdi Soltanolkotabi, and Salman Avestimehr

In this blog post, we describe the first peer-reviewed research paper that explores accelerating the hybrid of PyTorch DDP (torch.nn.parallel.DistributedDataParallel) [1] and Pipeline (torch.distributed.pipeline) - PipeTransformer: Automated Elastic Pipelining for Distributed Training of Large-scale Models (Transformers such as BERT [2] and ViT [3]), published at ICML 2021.

PipeTransformer leverages automated elastic pipelining for efficient distributed training of Transformer models. In PipeTransformer, we designed an adaptive on-the-fly freeze algorithm that can identify and freeze some layers gradually during training and an elastic pipelining system that can dynamically allocate resources to train the remaining active layers. More specifically, PipeTransformer automatically excludes frozen layers from the pipeline, packs active layers into fewer GPUs, and forks more replicas to increase data-parallel width. We evaluate PipeTransformer using Vision Transformer (ViT) on ImageNet and BERT on SQuAD and GLUE datasets. Our results show that compared to the state-of-the-art baseline, PipeTransformer attains up to 2.83-fold speedup without losing accuracy. We also provide various performance analyses for a more comprehensive understanding of our algorithmic and system-wise design.

Next, we will introduce the background, motivation, our idea, design, and how we implement the algorithm and system with PyTorch Distributed APIs.

https://pytorch.org/blog/pipetransformer-automated-elastic-pipelining/

Outline

ICML International Conference On Machine Learning



- Background and Related Works
- Motivation and Ideas
- Overall Design (Animation)
- AutoFreeze: Freeze Algorithm
- AutoPipe: Elastic Pipelining
- AutoDP: Spawning More Pipeline Replicas
- AutoCache: Cross-pipeline Caching
- Experimental Results
- Future Works







The parameter number of deep neural networks (Transformers) is dramatically increasing!

On the Opportunities and Risks of Foundation Models

Rishi Bommasani* Drew A. Hudson Ehsan Adeli Russ Altman Simran Arora Sydney von Arx Michael S. Bernstein Jeannette Bohg Antoine Bosselut Emma Brunskill Erik Brynjolfsson Shyamal Buch Dallas Card Rodrigo Castellon Niladri Chatterji Annie Chen Kathleen Creel Jared Quincy Davis Dorottya Demszky Chris Donahue Moussa Doumbouya Esin Durmus Stefano Ermon John Etchemendy Kawin Ethayarajh Li Fei-Fei Chelsea Finn Trevor Gale Lauren Gillespie Karan Goel Noah Goodman Shelby Grossman Neel Guha Tatsunori Hashimoto Peter Henderson John Hewitt Daniel E. Ho Jenny Hong Kyle Hsu Jing Huang Thomas Icard Saahil Jain Dan Jurafsky Pratyusha Kalluri Siddharth Karamcheti Geoff Keeling Fereshte Khani Omar Khattab Pang Wei Koh Mark Krass Ranjay Krishna Rohith Kuditipudi Ananya Kumar Faisal Ladhak Mina Lee Tony Lee Jure Leskovec Isabelle Levent Xiang Lisa Li Xuechen Li Tengyu Ma Ali Malik Christopher D. Manning Suvir Mirchandani Eric Mitchell Zanele Munyikwa Suraj Nair Avanika Narayan Deepak Narayanan Ben Newman Allen Nie Juan Carlos Niebles Hamed Nilforoshan Julian Nyarko Giray Ogut Laurel Orr Isabel Papadimitriou Joon Sung Park Chris Piech Eva Portelance Christopher Potts Aditi Raghunathan Rob Reich Hongyu Ren Frieda Rong Yusuf Roohani Camilo Ruiz Jack Ryan Christopher Ré Dorsa Sadigh Shiori Sagawa Keshav Santhanam Andy Shih Krishnan Srinivasan Alex Tamkin Rohan Taori Armin W. Thomas Florian Tramèr Rose E. Wang William Wang Bohan Wu Jiajun Wu Yuhuai Wu Sang Michael Xie Michihiro Yasunaga Jiaxuan You Matei Zaharia Michael Zhang Tianyi Zhang Xikun Zhang Yuhui Zhang Lucia Zheng Kaitlyn Zhou Percy Liang*1

> Center for Research on Foundation Models (CRFM) Stanford Institute for Human-Centered Artificial Intelligence (HAI) Stanford University











The cost of Search/Translation models (1/3)



2.3B parameters

"Our best quality dense single Transformer model (2.3B parameters) ... was trained for 235.5 TPU v3 core-years."

- 2021 paper from openreview.net
- "TPU v3 core": equivalent to a GPU, uses 100W

2.3B transformer model costs 206,000 kWh... ... takes 236 years to train on one GPU

206,000 kWh equals ...

- Driving 2,000 electric cars to another state (300 miles)
- Round-trip by 2,000 gas cars to another city (100 miles)



*Note: This page is from AAAI 2021 Tutorial (https://sites.google.com/view/aaai-2021-tutorial-ah9/home)





The cost of Search/Translation models (2/3)



retrained 100 times

Models re-trained 100s of times before service

- R&D test designs
- Hyperparameter optimization
- Upgrade with newer data

100x retraining of 2.3B model uses 20 gWh... ...takes 23.6 millennia to train on one GPU

20 gWh equals ...

- Round-trip by 200,000 gas cars to another city (100 miles)
- Yearly electricity of 2,000 USA homes



*Note: This page is from AAAI 2021 Tutorial (https://sites.google.com/view/aaai-2021-tutorial-ah9/home)





The cost of Search/Translation models (3/3)



scaled to 175B parameters retrained 100 times 2.3B is tiny vs. recent models @ Google, OpenAl, ...

- OpenAI GPT3 175B parameters
- Google model under open review 600B parameters

Extrapolating – 100x retraining of 175B model costs 1500 gWh... ...and would have taken 1,800 millennia to train on one GPU ...except you can't fit 175B parameters on one GPU's memory!

We can't avoid distributed computing Yet we need to reduce costs to be realistic & responsible

1500 gWh equals ...

Yearly home electricity of small city (150,000 homes)



*Note: This page is from AAAI 2021 Tutorial (https://sites.google.com/view/aaai-2021-tutorial-ah9/home)





All-reduce SGD







1. System-wise: Distributed Training System Design and Optimization

Data and Model Parallel	Data Inter-Batch: ByteScheduler (SOSP'19) Crossbow (VLDB'19)	Data Intra-Batch: PT Pipe PT DDP (VLDB'20) GPipe (NeurlPS'19)
Model Inter-Operator PT RPC TF + gRPC (EuroSys'19)	PipeDream (SOSP'19) HetPipe (ATC'20)	PT RPC + DDP PT RPC + Pipe Parallax (EuroSys'19) BytePS (OSDI'20)
Model Intra-Operator Mesh-TF (NeurlPS'18) TouFu (EuroSys'19)		FlexFlow (MLSys'19) GPT-3 (NeurlPS'20) ZeRO/DeepSpeed (SC'20)







Pipeline Parallelism



Hybrid of Data Parallelism and Pipeline Parallelism





Model Parallelism



Intra-layer and inter-layer model parallelism. https://FlexFlow.ai

2. ML-wise: Model Architecture and Training Algorithm

- **Architecture Optimization Manually:** 1) LinFormer (AAAI'2021, Best Paper Award)
- 2) Automated Architecture Design: Neural Architecture Search: FBNet (CVPR' 2019, 400+ citations)
- 3) Spare Training: pruning, quantization, etc Lottery Ticket Hypothesis (ICLR 2019, Best Paper Award)
- **Progressive Training** 4)
- 5) SGD-based Distributed Optimization: LARS (ICLR 2020): Large Batch Optimization for Deep Learning: Training BERT in 76 minutes





(c) Dilated sliding window

(b) Sliding window attention





USC University of

Southern California

(d) Global+sliding window



Our Motivation and Idea





1. (Sys) Distributed Training System	2. (ML) Model Architecture and Training Algorithm
Pro: Efficiency in Computation/Communication/Memory	Pro: Improve the efficiency mathematically, fundamentally
Con: View the model/SGD optimization as black box	Con: 1. Lack of system design to amplify the algorithmic advantages 2. the model-wise optimization is not friendly to distributed training
What if	f we co-design?
Hybrid of Pipeline and Data Parallelism	* 1. Progressive Training 2. Dynamic Neural Networks (https://arxiv.org/pdf/2102.04906.pdf)
Elastic Distribute	ed Training System!

Progressive Training







Freeze Training [1]

Algorithm 1 Progressive stacking
$M'_0 \leftarrow \text{InitBERT}(L/2^k)$
$M_0 \leftarrow \operatorname{Train}(M'_0)$ {Train from scratch.}
for $i \leftarrow 1$ to k do
$M'_i \leftarrow \text{Stack}(M_i)$ {Doubles the number of layers.}
$M_i \leftarrow \operatorname{Train}(M'_i) \{M_i \text{ has } L/2^{k-i} \text{ layers.}\}$
end for
return M_k

Progressive Stacking [2]

[1] Freeze Training: Singular Vector Canonical Correlation Analysis for Deep Learning Dynamics and Interpretability. NeurIPS 2017 [2] Efficient Training of BERT by Progressively Stacking. ICML 2019

[3] Accelerating Training of Transformer-Based Language Models with Progressive Layer Dropping. NeurIPS 2020. Minjia Zhang [4] On the Transformer Growth for Progressive BERT Training. NACCL 2021

Distributed Training System





USCUniversity of

Pipeline Parallelism

Hybrid of Data Parallelism and Pipeline Parallelism

Key observations when applying progressive training (e.g., freeze training) to the above training systems:

- 1. The computation cost becomes unbalanced in pipeline-parallelism
- 2. The memory cost is reduced gradually
- 3. The communication cost among DP workers should be reduced gradually

Overall Design







The process of PipeTransformer's automated and elastic pipelining





















Overall Design







Figure 3. Overview of PipeTransformer Training System

https://DistML.ai https://chaoyanghe.com/pipetransformer/

- Freeze + AutoPipe + AutoDP + AutoCache
- Freeze + AutoPipe + AutoCache Freeze Only
- Freeze + AutoPipe + AutoDP
 No Freeze (baseline)





(a) Sample Throughput(b) Speedup Ratio Comparison*Figure 9.* Speedup Breakdown (ViT on ImageNet)

Freeze Algorithm









Figure 4. Freeze Algorithm Using Different α

$$\min\left(L_{\text{frozen}}^{(T-1)} + \alpha(L - L_{\text{frozen}}^{(T-1)}), \underset{\ell \in \{L_{\text{frozen}}^{(T-1)}, \dots, L\}}{\operatorname{argmin}} \left\| \boldsymbol{g}_{\ell}^{(T)} \right\| \right)$$

where $T \ge 1, L_{\text{frozen}}^{(0)} = 0$, and $\alpha \in (0, 1)$
(1)





(1) Basic Usage of PyTorch Pipeline



Step 1: build a model including two linear layers
fc1 = nn.Linear(16, 8).cuda(0)
fc2 = nn.Linear(8, 4).cuda(1)

Step 2: wrap the two layers with nn.Sequential
model = nn.Sequential(fc1, fc2)

Step 3: build Pipe (torch.distributed.pipeline.sync.Pipe)
model = Pipe(model, chunks=8)

```
# do training/inference
input = torch.rand(16, 16).cuda(0)
output_rref = model(input)
```





(1) Pipeline Partitioning Strategy



Partition in the middle of skip connection: the parameter number is balanced, but it requires more communications



trade-off of computational load, communication cost, and memory consumption among partitions (each partition is loaded to one GPU)

NO Partition in the middle of skip connection: the parameter number is slightly unbalanced, but it requires less communications





(2) Pipeline Compression



compress the pipeline if
$$M_{GPU}^{(T)} \leq M_{GPU}^{(0)}$$

where $M_{GPU}^{(T)} \Leftrightarrow \max_{k \in \{0, \cdots, K-1\}} S_{p_k}$

To avoid extensive memory profiling, the compression algorithm uses the parameter size as a proxy for the training memory footprint.





(3) Dynamic Number of Micro Batches



Figure 6. Pipeline Bubble: $F_{d,b}$, $B_{d,b}$, and U_d denote forward, backward, and the optimizer update of micro-batch b on device d, respectively. The total bubble size in each iteration is (K - 1) times per micro-batch forward and backward cost.



(b) Profiling Optimal Chunk Number

When the pipeline compressed (K becomes smaller), the bubble is shrunk (left figure)

However, we find that the micro-batches size (M) also needs to be adjusted accordingly (right figure shows the optimal M in different K).
AutoPipe: Elastic Pipelining





(4) AutoPipe algorithm: put all together

Algorithm 1 AutoPipe Algorithm	
1: Input: model \mathcal{F} , layer number L and L_{frozen} , pipeline length K , frozen layer cost factor λ_{frozen} 2: Return: model $\mathcal{F}_{\text{frozen}}$, model $\mathcal{F}_{\text{pipe}}$, updated K ; 3: def m_partition (\mathcal{F} , L , L_{frozen}): //see 3.2.1 4: $\mathcal{F}_{\text{frozen}} = \text{Sequential}()$; model size $S_{\text{frozen}} = 0$ 5: $\mathcal{F}_{\text{pipe}} = \text{Sequential}()$; per-layer size $S_{\text{pipe}} = []$ 6: for layer index = L_{frozen} to L do 7: $f_{\text{ATT}_i}, f_{\text{MLP}_i} \leftarrow f_i$ 8: $\mathcal{F}_{\text{pipe}}.\text{append}(f_{\text{MTT}_i})$; $S_{\text{pipe}}.\text{append}(\text{m}_{\text{size}}(f_{\text{ATT}_i}))$ 9: $\mathcal{F}_{\text{pipe}}.\text{append}(f_{\text{MLP}_i})$; $S_{\text{pipe}}.\text{append}(\text{m}_{\text{size}}(f_{\text{MLP}_i}))$ 10: return $\mathcal{F}_{\text{frozen}}, S_{\text{frozen}}, \mathcal{F}_{\text{pipe}}, S_{\text{pipe}}$ 11: def load_balance($\mathcal{F}_{\text{pipe}}, S_{\text{pipe}}, K$): //Section 3.2.1 12: B_L =dict(), B_S =dict() // balanced L and S 13: $L_{\text{assigned}} = 0$; $S_{\text{total}} = \text{sum}(S_{\text{pipe}})$ 14: for partition index = k to K do 15: mean= $S_{\text{total}}/(K + k)$; 16: var=np.var($S_{\text{pipe}}[L_{\text{assigned}}:])/(K - k)$ 17: for sublayer index i = L_{assigned} to $\text{len}(S_{\text{pipe}})$ do	18: $S_k = S_{pipe}[i]$ 19: criterion= $B_S[i]$ - $S_{frozen}(1.0-\lambda_{frozen})$ + S_k 20: if criterion < mean + var then 21: B_S += S_k ; B_L +=1; $L_{assigned}$ +=1; S_{total} -= S_k 22: else 23: break 24: return B_L , B_S 25: \mathcal{F}_{frozen} , \mathcal{F}_{pipe} , S_{pipe} = m_partition($\mathcal{F}, L, L_{frozen}$) 26: while $K \ge 2$ do 27: B_L , B_S = load_balance(\mathcal{F}_{pipe} , S_{pipe} , $K/2$) 28: $B_S[0]$ -= $S_{frozen}(1.0 - \lambda_{frozen})$; 29: $M_{GPU}^{(T)}$ = max(B_S) //Equation 2 30: if $M_{GPU}^{(T)} < M_{GPU}^{(0)}$ then 31: $K=K/2$ 32: else 33: break 34: load \mathcal{F}_{frozen} and \mathcal{F}_{pipe} to K GPUs using B_S and B_L 35: Pipe(\mathcal{F}_{pipe} , chunks= get_optimal_chunks (K)
 A trade-off in Pipeline partition Pipeline compression 	

3. optimal micro-batches chunk number (M)







pipeline 0 at timestep 1 Pipeline 1 at timestep 1 ٢ ٢ * ٠ Dataset Dataset pipeline 0 at timestep 0 add new pipelines **•** 4 🕑 Þ ٢ ٢ -0 Dataset Redistribute dataset AutoCache AutoDP (cross process) Transform Sample training # of frozen laver information as indicator Pipeline length changed? (Progress, gradient, etc) has been changed? Freeze notify AutoPipe Algorithm Shared Memory Data Distributed Parallel **Pipeline Parallel** (Cross process) Deep Learning Training Engine (PyTorch) Multi Processing CUDA NCCL / GLOO

Figure 3. Overview of PipeTransformer Training System







Spawning More Pipeline Replicas



Key challenges when adding more pipeline on the fly of training:

- 1) **DDP Communication:** Collective communications in PyTorch DDP requires **static** membership, which prevents new pipelines from connecting with existing ones;
- 2) State Synchronization: newly activated processes must be consistent with existing pipelines in the training progress (e.g., epoch number and learning rate), weights and optimizer states, the boundary of frozen layers, and pipeline GPU range;
- 3) **Dataset Redistribution:** the **dataset** should be **re-balanced** to match a dynamic number of pipelines. This not only avoids stragglers but also ensures that gradients from all DDP processes are equally weighted.





Spawning More Pipeline Replicas



Figure 7. AutoDP: handling dynamical data parallel with messaging between double process groups (Process 0-7 belong to machine 0, while process 8-15 belong to machine 1)

Our idea:

AutoDP:

1. creating two process groups. Each process handles one pipeline

2. the active training process group (yellow) handles the training.

3. the message process group (purple) handles **<u>State Synchronization and Dataset Redistribution</u> by messaging passing between two groups** with MPI communication.

AutoDP: Spawning More Pipeline Replicas



USC University of Southern California

On Machine Learnin

return object_list

AutoCache: Cross-pipeline Caching







In this example, **the first 3 layers (purple)** at two time steps T1 and T2(epochs) are the same computation, so T2 can reuse the caching from T1.





1. Overall Speedup

	Baseli	ne	PipeTrans		
Dataset	Accuracy	Training time	Accuracy	Training time	Training Speedup
ImageNet	80.83 ± 0.05	26h 30m	82.18 ± 0.32	9h 21m	2.83 ×
CIFAR-100	91.21 ± 0.07	35m 6s	91.33 ± 0.05	12m 23s	$2.44 \times$
SQuAD 1.1	90.71 ± 0.18	5h 7m	90.69 ± 0.23	2h 26m	$2.10 \times$

Table 1. Speedup for ViT and BERT Training

Evaluation on Various datasets and models, including tasks in both CV and NLP.





2. Breakdown for speedup



Figure 9. Speedup Breakdown (ViT on ImageNet)

Key takeaway:

1. the main speedup is the result of elastic pipelining which is achieved through the joint use of AutoPipe and AutoDP (purple)

2. AutoCache's contribution is amplified by AutoDP (green v.s. blue: more parallel DP workers can use caching)

3. freeze training alone without system-wise adjustment even downgrades the training speed (yellow) (the underlying mechanism of PyTorch is not tailored for freeze training, forcing CUDACachingAllocator to split blocks or launch new memory allocations)





3. Breakdown for communication v.s. computation

Table 2. Communication Cost v.s. Computational Cost

	Dataset	Overall Cost	Communication Cost	Computation Cost	Communication Cost Ratio
ViT-Base (87M)	ImageNet	9h 21m	34m	8h 47m	5.9 %
BERT-large (340M)	SQuAD	2h 26m	16m 33s	2h 9m	8.8%

Communication Infrastructure: InfiniBand CX353A where **cross-machine bandwidth is 5GB/s**, and GPU-to-GPU bandwidth within a machine (PCI 3.0, 16 lanes) is **15.754GB/s**.

Key takeaway:

1. Communication cost is not the main bottleneck when we use InfiniBand for medium-scale models (< 500M such as ViT-base and BERT-large), but it is still non-trivial even under freeze training.

2. Recent progress in NLP and CV has scaled up the model size to billion/trillion-level (GPT-3 - 175B [1], Switch Transformer - 1.7T [2]), which will make the ratio of communication much higher than our experimental results.

[1] Language Models are Few-Shot Learners. 2020[2] Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity. 2021





4. Performance Analysis



2. Timing of caching is important.

We automate these two optimization strategies.

Future works





Distributed training tasks that can be dynamic:

- 1. Elastic cloud-based distributed training system
- 2. Accelerating NAS in extremely large search space
- 3. Federated AutoML (NAS, HPO)
- 4. Cross-silo Federated Learning
- 5. Pruning-based distributed training
- 6. IoT device-based elastic edge training

Ongoing Research - Mixture of Experts



Ongoing Research - Mixture of Experts



Outline: Distributed ML for Large-scale Models



Part 1: Dynamic Distributed Training
 PipeTransformer: Automated Elastic Pipelining for Distributed
 Training of Large-scale Models. ICML 2021
 <u>https://DistML.ai</u>



 Part 2: Scalable Federated Learning
 [Algorithm] Group Knowledge Transfer: Federated Learning of Large CNNs at the Edge. NeurIPS 2020

[System] FedML: A Research Library and Benchmark for Federated Machine Learning (Best Paper Award, NeurIPS 2020 FL Workshop)

[Diverse Applications] FedML Ecosystem: Ubiquitous Distributed Training for Diverse AI Applications at the Edge. <u>https://FedML.ai</u>

Fundamental Research to Alexa Al Applications



1. Distributed Training in GPU cluster 2. Federated Learning at the edge

Goal: enabling secure collaborative learning at the edge



Introduction: Widely Interdisciplinary

distributed optimization personalization and non-IID transfer learning un/semi-supervised learning neural architecture search continual learning/meta learning

communication efficiency computation efficiency wireless communication cloud computing embedded training system

. . .

...



Conference Venues: NeurIPS, ICML, ICLR, MLSys, CVPR/ICCV/ECCV, ACL/EMNLP, AISTATS, AAAI

Foundations of Algorithm Design

User selection, optimal scaling
 Heterogeneity, personalization, and fairness
 Unsupervised federated learning

System Bottlenecks/Opportunities

- Resource-constrained FL (small edge models, large server models)
- **C** Scalability: 1K users \rightarrow 1M users
- **G** Federated neural architecture search



- Adversarial users (data/model poisoning)
- Leveraging trusted computing environments







https://www.avestimehr.com/fedml https://FedML.ai



Overview





FedML System Capability



- Single process
- Runs on any server
- Good for small models (LR, 2CNN, Bi-LSTM)
- And small datasets (FMNIST, Shakespeare)



- Multiple-processes in parallel
- Can be seamlessly distributed across nodes
- Train big models (ResNet, MobileNet, Efficient Net, Transformer-based)
- And large datasets (CIFAR10/100, Google Landmark, COCO, ImageNet)



- Multiple processes with explicit message passing
- Great for resource constrained edge devices (Smartphone and IoT) Iow memory, Iow computational power, limited communication bandwidth
- Python-centric simplicity

FedML Ecosystem

	Southern Calif	orni	a 🔻	макі			1
2 Applications	FedCV Fe	dNLP	FedMedic	al Fe	edFinance		
FedML Server	(HTTP, MQTT)		FedML-	API (hig	gh-level A	PIs)	-
НТТР/ТСР	НТТР/ТСР		Model		D	ata	
FedML-Mobile	FedML-IoT		Algorithm	(distribu	uted/stand	alone)	
iOS	RaspBerry Pi		FedAvg Fed	Opt F	edNova	FedNAS	
			Decentralized	FL Ve	rtical FL		
		L				·····	
	FedML-co	re (low	/-level APIs)]
Topology	Security/Privacy		Worker	Co	ordinator		
Cor	mManager		↓ I		¥		
Send Thread	Receive Thread		Models (LR,	DNN), O	ptimizer		
↓ ↓				V			
Abstract Co	mmunication Layer		On-Client Lea	rning Fr	amework		
	↓			, †			
MPI MQTT	Other Backend		PyTorch	Mo	bile Trainir	ng	
Distributed	d Communication		Trainii	ng Engin	e		
(🗘 PyTorch	1					







FedML Ecosystem

FedNLP

FedCV

FedGraphNN FedIoT



FedNLP







Figure 1. Our philosophy of federated learning on computer vision: connecting the algorithmic FL research and CV application-drive research with an unified research framework.



Based on FedML.ai, the blocks with colors are newly upgraded modules for FedCV

Figure 2. Overview of FedCV System Architecture Design

FedGraphNN



Figure 2: Tasks in different datasets focus on different levels of properties of molecules.



Figure 1: Formulation of FedGraphNN (Federated Graph Neural Network)



Figure 2. Overview of FedGraphNN System Architecture Design

FedGraphNN



Figure 1: Comparisons between centralized and decentralized training of GNN based recommendation models.

- 1. *Graph level*. We believe molecular machine learning is a paramount application in this setting, where many small graphs are distributed between multiple edge devices;
- 2. *Sub-graph level.* This scenario typically pertains to social networks or knowledge graphs that need to be partitioned into many small sub-graphs due to data barriers between different departments in a giant company, as demonstrated in (Wu et al., 2021).
- 3. *Node level*. When the privacy of a specific node in a graph is important, node-level GNN-based FL is useful in practice. The IoT setting is a good example (Zheng et al., 2020);
- 4. *Link level* is also a promising direction that is relevant when the privacy of edges (eg: connections in a social network) is of importance.

FedIoT: "Internet of Things & 5G + Federated Learning"

Network-based Detection of IoT Botnet Attacks Using Deep Autoencoders [1]



Fig. 1: Lab setup for detecting IoT botnet attacks



FedML Community



fedml.slack.com

•			Q Search FedML.ai					2
FedML.ai ~	Ø	#fedml-syst Add a topic	em 🏠		M 20 👤	176	2+ (i)
Threads		Our m packa	ain contribution here is that we provide a reg	Monday, December 7th ers to address so n	any open problems	. The Py	Forch Arm	n
6 ³ All DMs		We su	pport pytorch 1.4 in IoT devices					
More	15	M Mahd Hello	i Chehimi 5:33 PM all.					
 Channels 		Can I	run the FedML codes with other Pytorch/Tens	orflow libraries?				
 A documentation A fed-security A fed-segment 		PyTore as for https://	ang He 5:33 PM th is already supported TensorFlow, you can also customize your train (vithuh.com/FedIMI-A1/FedIMI./tree/master/f	er like this example shows: edml_experiments/distributed/fedays/edited		0 4	0 :	
△ fed-xgboost		5:36 image i		com_experiments/anstinutea/reading/conco		e r	ω :	
≙ fedcv		5.50 image.	HIG * State distributions, de arter, anire, inter, son, miri, trat, dray, an, trat, dray,					
△ federated-gan△ federated-vae			THE LAST CONTRACTOR AND					
≙ fedgnn			Displaying, doin, improved, whether, and, solution, twicks, industria, transfer, industry, and protect					
# fedml-algorithms			Lysenetropy, desire, one, rate, size, and, transplatique, trainderighted, tert/det.phile, trainderiched, etc., etc., exclusion/active, trainderiched, and conserva- net/conservations. And etc., and					
# fedml-framework		8.	and United and Link (D) sectors and the state of the sector of the secto					
# fedml-hierarchical			propier = holdingeriphi Chain, dei Julia, Intel, Julia, Julia, Lina, Julia, Jul					
# fedml-model-and-	dataset		start the startsdard tostong very analyse - Notificiansekoapenlergh, aggregatar, some, nask, skart very analyses - Notificiansekoapenlergh, aggregatar, some, nask, skart very analyses - Notificiansekoapenlergh, aggregatar, some, nask, skart					
# fedml-system								
A fedmobile		by det	ault, the model trainer is None, so we support	it with Pylorch implementation				
A fednlp-private		Chaoy	ang He 5:37 PM u cap also define this trainer and follow the al	etract ABIs defined at fadel area (basines	(model trained as			
		July C	et	stract Aris defined at. Team_core/cruthery	model_cruther.py			
△ optimizers								
△ robustness		1	reply 1 day ago					
+ Add channels		Message #	fedml-system					٦
 Direct messages 								
Slackbot		Ø B	I ⊕ ↔ ∅ 1≡ 1≡ 1≡ 1⊡		Aa	00	0 >	
Chaoyang Hel you								

Growth Statistics:

Slack Community Users: 640

GitHub Stars: 780 GitHub Forks: 216

GitHub Ranking in "Federated Learning": 3 (the first two are industry leaders: Tencent/WeBank, Google)

Short Paper won Best Paper Awards in NeurIPS 2020 SpicyFL Workshop



FedML System Overview





FedML-API - Simplicity is Our Key



FedAvg, FedAvg_Robustness, FedOpt (server optimizer), FedNova (client optimizer)
 FedGKT, FedNAS, Decentralized FL, Vertical FL, Split Learning, etc.



FedML Feature List

Dataset	MNIST, Synthetic Federated EMNIST; Shakespeare; Fed-CIFAR100; stackoverflow (NWP); CIFAR10, CIFAR100, CINIC10 (ImageNet+CIFAR10) Google Landmark, COCO, ImageNet Support non-IID partition tool for heterogeneous distribution
Model	LR, CNN (2 layers), RNN (Bi-LSTM), ResNet, MobileNet V3, EfficientNet Transformer/BERT, etc.
Federated Optimizers and Algorithms	FedAvg, FedOpt (server optimizer, ICLR 2021), FedNova (client optimizer, NeurIPS 2020), FedAvg_Robustness (NeurIPS 2020), FedGKT (NeurIPS 2020), FedNAS (CVPRw 2020), Turbo-Aggregate, Decentralized FL (NeurIPS 2020 FL), Vertical FL, Split Learning, etc <i>Support both cross-device and cross-silo settings.</i>
Platform Supports	Distributed Computing, IoT/Mobile, Standalone Simulation

Three computing paradigms, SOTA optimizers, various datasets and models!



FedML Feature List - Privacy/Security

Dataset	MNIST, Synthetic Federated EMNIST; Shakespeare; Fed-CIFAR100; stackoverflow (NWP); CIFAR10, CIFAR100, CINIC10 (ImageNet+CIFAR10) Google Landmark, COCO, ImageNet Support non-IID partition tool for heterogeneous distribution
Model	LR, CNN (2 layers), RNN (Bi-LSTM), ResNet, MobileNet V3, EfficientNet Transformer/BERT, etc.
Federated Optimizers and Algorithms	FedAvg, FedOpt (server optimizer, ICLR 2021), FedNova (client optimizer, NeurIPS 2020), FedAvg_Robustness (NeurIPS 2020), FedGKT (NeurIPS 2020), FedNAS (CVPRw 2020), Turbo-Aggregate, Decentralized FL (NeurIPS 2020 FL), Vertical FL, Split Learning, etc <i>Support both cross-device and cross-silo settings.</i>
Platform Supports	Distributed Computing, IoT/Mobile, Standalone Simulation

Three computing paradigms, SOTA optimizers, various datasets and models!



System Config diversity

Diverse topologies:





Deployment for Cross-silo FL



FedML Live Demo






FedML-core - Worker-oriented Programming 🕸 턌 Interface



device = init training device (process id, worker number - 1, args.gpu num per server) # load data dataset = load data(args, args.dataset) [train data num, test data num, train data global, test data global, train data local num dict, train data local dict, test data local dict, class num] = dataset # create model model = create model(args, model name=args model, output dim=dataset[7]) # start "federated averaging (FedAvg)" FedML FedAvg distributed (process id, worker number, device, Comm, model, train data num, train data global, test data global,

train data local num dict,

train_data_local_dict,
test data_local_dict, args,

FedML-core - Trainer Customization

topology configuration



class ModelTrainer(ABC): def init (self, model): self.model = model self.id = 0def set id(self, trainer id): self.id = trainer id def set model params(self, model parameters): Pass def train(self, train data, device, args): Pass def test(self, test data, device, args): pass

USCUniversity of Southern California

🗘 PyTorch

NEURAL INFORMATION

model trainer)



HTTP / TCP

O MCU

FedML supports diverse platforms

class ModelTrainer(ABC): init (self, model): def self.model = model self.id = 0

def set id(self, trainer id): self.id = trainer id

def set model params(self,

model parameters): Pass

def train(self, train data, device, args): Pass

def test(self, test data, device,

args):

pass







Write once, run everywhere:

Reusing the same trainer class definition in three platforms Let the library do the algorithm and platform-specific implementation





FedGKT: An Algorithmic Example of FedML:

FedGKT is a new distributed training framework, and requires specific communication protocol (*soft labels, hidden features*, etc) and training paradigms (*training models in client and server alternatively*).





Step 1: Model weight

Step 2&3: Hidden Representation, Logits (soft labels)





Group Knowledge Transfer: Federated Learning of Large CNNs at the Edge



Chaoyang He PhD Student, USC



Murali Annavaram Professor, USC



Salman Avestimehr Professor, USC



https://fedml.ai

Introduction

• Extremely Large DNN: the SOTA model in NLP (Natural Language Processing) domain has billion or even trillion of parameters [1], and the SOTA model in CV (Computer Vision) domain has around 4 million parameters [2].



[1] Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. 2019-10. https://arxiv.org/abs/1910.10683

[2] MiLeNAS: Efficient Neural Architecture Search via Mixed-Level Reformulation. CVPR 2020 Chaoyang He, Haishan Ye, Li Shen, Tong Zhang.

Motivation

Edge devices are personal devices with limited memory and storage and computation capabilities



Motivation



Challenges of **Federated Deep Learning** - Training large modern DNN (ResNet, Transformers, etc) on-device is difficult due to data and system resource constraints:

1. Limited memory and computation (no GPU accelerator) 2. Limited bandwidth and unstable wireless communication

3. Data heterogeneity



How to efficiently and effectively train large DNN over resource-constrained edge devices?

Related Works: a Basic Formulation of Federated Learning

$$\min_{\boldsymbol{W}} F(\boldsymbol{W}) \stackrel{\text{def}}{=} \min_{\boldsymbol{W}} \sum_{k=1}^{K} \frac{N^{(k)}}{N} \cdot f^{(k)}(\boldsymbol{W}), \text{where } f^{(k)}(\boldsymbol{W}) = \frac{1}{N^{(k)}} \sum_{i=1}^{N^{(k)}} \ell(\boldsymbol{W}; \boldsymbol{X}_i, y_i)$$

Algorithm 1 FedAvg Algorithm: A Challenge Perspective Server 1: Initialization: there is a number of clients in a network; the client k has local dataset \mathcal{D}^k ; each client's local model is initialized as W_0 ; 2: 3: Server Executes: 4: for each round t = 0, 1, 2, ... do $S_t \leftarrow \text{(sample a random set of clients)}$ 5: for each client $k \in S_t$ in parallel do 6: $\boldsymbol{W}_{t+1}^k \leftarrow \text{ClientUpdate}(k, \boldsymbol{W}_t)$ 7: end for 8: $\boldsymbol{W}_{t+1} \leftarrow \sum_{k=1}^{K} \frac{n_k}{n} \boldsymbol{W}_{t+1}^k$ 9: 10: end for 11: **Client_Training**(k, W): // Run on client k 12: 13: $\mathcal{B} \leftarrow (\text{split } \mathcal{D}^k \text{ into batches})$ Client 1 Client 2 Client K 14: for each local epoch i with $i = 1, 2, \cdots$ do . . . for batch $b \in \mathcal{B}$ do 15: $\boldsymbol{W} \leftarrow \boldsymbol{W} - \eta \nabla_{\boldsymbol{W}} F(\boldsymbol{W}; b)$ 16: end for 17: W Data 18: end for What if W is DNN? 19: return W to server $\boldsymbol{W} \leftarrow \boldsymbol{W} - \eta \nabla_{\boldsymbol{W}} F(\boldsymbol{W}; b)$

 $\boldsymbol{W}_{t+1} \leftarrow \sum_{k=1}^{K} \frac{n_k}{n} \boldsymbol{W}_{t+1}^k$

[1] McMahan et al., Communication-Efficient Learning of Deep Networks from Decentralized Data, AISTATS 2017

Related Works: a Basic Formulation of Federated Learning

$$\min_{\boldsymbol{W}} F(\boldsymbol{W}) \stackrel{\text{def}}{=} \min_{\boldsymbol{W}} \sum_{k=1}^{K} \frac{N^{(k)}}{N} \cdot f^{(k)}(\boldsymbol{W}), \text{where } f^{(k)}(\boldsymbol{W}) = \frac{1}{N^{(k)}} \sum_{i=1}^{N^{(k)}} \ell(\boldsymbol{W}; \boldsymbol{X}_i, y_i)$$

Algorithm 1 FedAvg Algorithm: A Challenge Perspective Server 1: Initialization: there is a number of clients in a network; the client k has local dataset \mathcal{D}^k ; each client's local model is initialized as W_0 ; 2: 3: Server Executes: 4: for each round t = 0, 1, 2, ... do $S_t \leftarrow \text{(sample a random set of clients)}$ 5: for each client $k \in S_t$ in parallel do 6: $\boldsymbol{W}_{t+1}^k \leftarrow \text{ClientUpdate}(k, \boldsymbol{W}_t)$ 7: end for 8: $\boldsymbol{W}_{t+1} \leftarrow \sum_{k=1}^{K} \frac{n_k}{n} \boldsymbol{W}_{t+1}^k$ 9: 10: end for 11: **Client_Training**(k, W): // Run on client k 12: 13: $\mathcal{B} \leftarrow (\text{split } \mathcal{D}^k \text{ into batches})$ Client 1 Client 2 Client K 14: for each local epoch i with $i = 1, 2, \cdots$ do . . . for batch $b \in \mathcal{B}$ do 15: $\boldsymbol{W} \leftarrow \boldsymbol{W} - \eta \nabla_{\boldsymbol{W}} F(\boldsymbol{W}; b)$ 16: end for 17: W Data 18: end for What if W is DNN? 19: return W to server $\boldsymbol{W} \leftarrow \boldsymbol{W} - \eta \nabla_{\boldsymbol{W}} F(\boldsymbol{W}; b)$

 $\boldsymbol{W}_{t+1} \leftarrow \sum_{k=1}^{K} \frac{n_k}{n} \boldsymbol{W}_{t+1}^k$

[1] McMahan et al., Communication-Efficient Learning of Deep Networks from Decentralized Data, AISTATS 2017

Related Works: Split Learning



[2] Poirot et al., Split Learning for collaborative deep learning in healthcare, NeurIPS 2019

Related Works

Federated Learning:

Pros:

1) exchanging gradients/models periodically;

2) reusing *distributed optimization* methods in conventional distributed training in the data center environment.

Cons:

1) Communication cost for each gradient/model is much higher than a hidden vector

2) training large DNN models on device is prohibited due to resource constraints.

Split Learning:

Pros:

1) lower communication cost than FL: only exchanging the hidden vector of each training sample/mini-batch;

2) a portion of the entire DNN architecture is trainable on resource constrained devices. 3) hidden vector is more secure to resist adversarial attacks.

Cons:

1) Do not support periodical synchronization, which prohibits off-the-shell optimization methods; 2) the straggler problem becomes more severe because one SGD iteration has been split into four rounds of communication.



representations



Insights

Can we design a learning framework that leverages advantages of FL and SL?

To be more specific, is there a framework that supports:

- Computation Efficient: Computation-efficient on-device training (trainable) like SL
- 2. Communication Efficient: Exchanging hidden vectors during training like SL
- 3. Low Communication Frequency: Supporting periodical training like FL (local SGD)
- 4. No Accuracy Compromise: Preserving or outperforming the model accuracy of FL and SL

Our Idea

1. Transfer Data (Centralized Training)

3. Can we only transfer "knowledge"?

2. Transfer Model (Federated Learning)

Transfer knowledge from many small networks to a larger one which has more capacity to obtain high accuracy



FedGKT Overview



First Formulation



find a good and small feature extractor

First Formulation: ideal



find a good and small feature extractor

find a great and large classifier by leveraging users' feature extractors

First Formulation: ideal



First Formulation: practical



Knowledge Distillation

• Bidirectional Knowledge Distillation (CVPR 2018, ICLR 2018)



Transfer from 2 to 1: The overall loss function L_{Θ_1} for network Θ_1 is defined as

$$L_{\Theta_1} = L_{C_1} + D_{KL}(\boldsymbol{p}_2 \| \boldsymbol{p}_1)$$

Reformulation as Bidirectional Knowledge Transfer



 $\underset{\boldsymbol{W}^{(k)}}{\operatorname{argmin}} F_c(\boldsymbol{W}^*_s, \boldsymbol{W}^{(k)}) = \mathsf{CE} \text{ loss of the classifier} + \ell_{KD}(\boldsymbol{z}^*_s, \boldsymbol{z}^{(k)}_c) \quad \underset{\boldsymbol{W}_s}{\operatorname{argmin}} F_s(\boldsymbol{W}_s, \boldsymbol{W}^{(k)*}_e) = \mathsf{CE} \text{ loss of the classifier} + \sum_{k=1}^n \ell_{KD}(\boldsymbol{z}^{(k)*}_c, \boldsymbol{z}_s)$

The logits z_s and z_c are the output of the last fully connected layer in the server model and the client model

Reformulation as Bidirectional Knowledge Transfer



The logits z_s and z_c are the output of the last fully connected layer in the server model and the client model

Reformulation as Bidirectional Knowledge Transfer



□ server absorbs users' knowledge to train a better classifier

FedGKT: Alternating Minimization

Server
$$\underset{\mathbf{W}_{s}}{\operatorname{argmin}} F_{s}(\underline{\mathbf{W}_{s}}, \underline{\mathbf{W}_{e}^{(k)*}}) = \underset{\mathbf{W}_{s}}{\operatorname{argmin}} \sum_{k=1}^{K} \sum_{i=1}^{N^{(k)}} \ell_{CE}(f_{s}(\mathbf{W}_{s}; \underline{f_{e}^{(k)}}(\mathbf{W}_{e}^{(k)*}; \mathbf{X}_{i}^{(k)})), y_{i}^{(k)}) + \sum_{k=1}^{K} \ell_{KD}(\mathbf{z}_{c}^{(k)*}, \mathbf{z}_{s})$$

$$\overset{(8)}{H_{i}^{(k)}}$$

$$\underset{\mathbf{W}_{s}}{\operatorname{argmin}} F_{c}(\underline{\mathbf{W}_{s}^{*}}, \underline{\mathbf{W}^{(k)}}) = \underset{\mathbf{W}_{s}}{\operatorname{argmin}} \sum_{i=1}^{N^{(k)}} \ell_{CE}(f_{c}^{(k)}(\mathbf{W}_{e}^{(k)}; \underline{f_{e}^{(k)}}(\mathbf{W}_{e}^{(k)}; \mathbf{X}_{i}^{(k)})), and \mathbf{z}_{s} = f_{s}(\mathbf{W}_{s}; \mathbf{H}_{i}^{(k)})$$

$$\underset{\mathbf{W}_{s}}{\operatorname{argmin}} F_{c}(\underline{\mathbf{W}_{s}^{*}}, \mathbf{W}^{(k)}) = \underset{\mathbf{W}_{s}}{\operatorname{argmin}} \sum_{i=1}^{N^{(k)}} \ell_{CE}(f_{c}^{(k)}(\mathbf{W}_{c}^{(k)}; \underline{f_{e}^{(k)}}(\mathbf{W}_{e}^{(k)}; \mathbf{X}_{i}^{(k)})), y_{i}^{(k)}) + \ell_{KD}(\mathbf{z}_{s}^{*}, \mathbf{z}_{c}^{(k)})$$

$$\underset{\mathbf{W}_{s}}{\operatorname{client}} (10)$$

where
$$\mathbf{z}_{c}^{(k)} = f_{c}^{(k)}(\mathbf{W}_{c}^{(k)}; \underbrace{f_{e}^{(k)}(\mathbf{W}_{e}^{(k)}; \mathbf{X}_{i}^{(k)})}_{\mathbf{H}_{i}^{(k)}}), and \mathbf{z}_{s}^{*} = f_{s}(\mathbf{W}_{s}^{*}; \mathbf{H}_{i}^{(k)})$$
 (11)

Method: Distributed Optimization Methods

1. Local SGD with Momentum

$$V_t = \beta V_{t-1} + \alpha \nabla_w L(W, X, y)$$
$$W = W - V_t$$

2. Cross-round Learning Rate Scheduler



Method: Model Architecture



System Design and Implementation

FedGKT is a new distributed training framework, and requires specific communication protocol (*soft labels, hidden features*, etc) and training paradigms (*training models in client and server alternatively*).





Step 1: Model weight

Step 2&3: Hidden Representation, Logits (soft labels)

System Design and Implementation

With the help of FedML library, we can easily implement FedGKT algorithm and conduct various experiments with different datasets and models.



Homepage: https://fedml.ai

(Best Paper Award, NeurIPS 2020, FL Workshop)

Experimental Results: Test Accuracy



Figure 3: The Test Accuracy of ResNet-56 (Edge Number = 16)

Experimental Results: Non-IID

Table 1: The Test Accuracy of ResNet-56 and ResNet-110 on Three Datasets.

Model	Methods	CIFAR-10		CIFAR-100		CINIC-10	
		I.I.D.	non-I.I.D.	I.I.D.	non-I.I.D.	I.I.D.	non-I.I.D.
ResNet-56	FedGKT (ResNet-8, ours) FedAvg (ResNet-56)	92.97 92.88	86.59 86.60	69.57 68.09	63.76 63.78	81.51 81.62	77.80 77.85
	Centralized (ResNet-56) Centralized (ResNet-8)	93.05 78.94		69.73 37.67		81.66 67.72	
ResNet-110	FedGKT (ResNet-8, ours) FedAvg (ResNet-110) Centralized (ResNet-110) Centralized (ResNet-8)	93.47 93.49	87.18 87.20 93.58 78.94	69.87 68.58	64.31 64.35 70.18 37.67	81.98 82.10	78.39 78.43 82.16 67.72

Experimental Results: Efficiency



FedGKT demands **9 to 17 times** less computational power (FLOPs) on edge devices and requires **54 to 105 times** fewer parameters in the edge CNN

Experimental Results: Ablation Study



Figure 5: GKT result on CIFAR10

Experimental Results: Ablation Study



Figure 6: Client Model Architecture Exploration

Experimental Results: Ablation Study



Figure 4: Evaluation on the distillation loss on the server side

Contributions

FedGKT is the early work that explores <u>federated deep learning</u> and transfer **knowledge** (not data and model) on edge devices.

- FedGKT is memory and computation efficient, similar to SL
- FedGKT can train in a local SGD manner like FedAvg to reduce the communication frequency
- Exchanging hidden features as in SL, as opposed to exchanging the entire model as in FedAvg, **reduces the communication bandwidth requirement**
- FedGKT is **public data-free** knowledge transfer method
- FedGKT naturally supports **asynchronous** training, which circumvents the severe synchronization issue in SL

Future Works

- Exchanging hidden features provides "some" privacy guarantees to the users
 - □ Providing strong privacy/security guarantees for FedGKT is an interesting next step
 - □ Leveraging Secure Aggregation in FedGKT?


Future Works

• Extend FedGKT to Transformer Models



Figure 12. Transformer Model Architecture (Vaswani et al., 2017)



Figure 13. Vision Transformer (Dosovitskiy et al., 2020)

Future Works

- FedGKT also allows for different models to be used at each user
 - Providing model adaptation and personalization at each user is also another next step



Long-term Goal

- 1. Transfer Data (Centralized Training)
- 2. Transfer Model (Federated Learning)



3. Can we only transfer "knowledge"?

FedGKT is a starting point of our goal towards "Knowledgeable Communication":
1) AI Agent (ML model) in Internet can freely exchange/transfer/share their "knowledge" without disclosing a large amount of raw data or model;
2) ML models can be equipped with communication networking interface to exchange neural representation with other models in a decentralized Internet.

Distributed Training v.s. Federated Learning





In essence, three methodologies are useful to both:

1. Dynamic system and stochastic ML lead to complicated interaction between these objectives.

PipeTransformer, ICML 2021

2. Trade-off among multiple objectives including accuracy, efficiency, security, and privacy.

Group Knowledge Transfer, NeurIPS 2020

3. On demand: focus on the key demand for different use cases (ML applications), optimizing one of the objectives but not deteriorating the others.

Federated Learning

Foundations of Algorithm Design

- User selection, optimal scaling
- Heterogeneity, personalization, and fairness
- Unsupervised federated learning

System Bottlenecks/Opportunities

- Resource-constrained FL (small edge models, large server models)
- $\Box \quad \text{Scalability: 1K users} \rightarrow 1\text{M users}$
- **G** Federated neural architecture search



Trustworthy

- Secure and resilient model aggregation
- Adversarial users (data/model poisoning)
- Leveraging trusted computing environments

Ongoing Works 1: Self-supervised Federated Learning



Fig. 1: Personalized Self-supervised Federated Learning

Ongoing Works 1: Self-supervised Federated Learning



Elementary Result: Contrastive Loss with SimSiam Framework in FL setting (Weighted Averaging) can archive top-1 accuracy 91% in CIFAR-10 (non-IID) using ResNet-18!

Important optimization tricks: 1) local SGD with momentum 2) cross-round learning rate scheduler

Ongoing Work 2 - Joint Adaptation to Data and System Heterogeneity with FedNAS



Figure 1: Federated Neural Architecture Search (*step 1*: search locally; *step 2*: sending the gradient of α and ω to the server side; *step 3*: merge gradient to get global α and ω ; *step 4*: synchronize updated α and ω to each client.)





$$\bar{o}^{(i,j)}(x) = \sum_{k=1}^{d} \underbrace{\frac{\exp(\alpha_k^{(i,j)})}{\sum_{k'=1}^{d} \exp(\alpha_{k'}^{(i,j)})}}_{p_k} o_k(x)$$
(3)

FedML and DistML Ecosystem



Application/Model

FedCV, **FedNLP** (preparing for EMNLP 2021) **FedGraphNN** (ICLR 2021 workshop, MLSys 2021 workshop)

Algorithm/Theory

FedGKT (NeurIPS'2020): resource constrained FL

FedNAS, MiLeNAS (CVPR'2020): data heterogeneity and automation

OnlinePushSum (NeurIPS'2020 Workshop): Single-sided Trust in Decentralized Topology

System/Infrastructure

FedML.ai (NeurIPS 2020, FL Workshop, Best Paper Award): fundamental system, an open source library for FL research

PipeTransformer (ICML'2021): elastic distributed training for giant models (Transformers)

Vision/Review

Advances and Open Problems in Federated Learning FnTML (Foundation and Trend in Machine Learning), Vol3, 2021