

# The Way towards Ubiquitous Learning: a System Perspective

Mengwei Xu (徐梦炜) CS Dept of BUPT

#### Outline



- Ubiquitous Learning (UL): What and Why?
  - A system software perspective
  - The status quo
- A measurement study of UL
- Memory optimizations of UL
- Other on-going research of UL

#### Some trends

![](_page_2_Picture_1.jpeg)

 Edge devices (smartphones, IoTs, etc) are becoming important computing platforms, not just user equipment<sup>[1]</sup>.

![](_page_2_Figure_3.jpeg)

[1] Mengwei Xu, et al. "A case for camera-as-a-service", IEEE Pervasive Computing, 2021.[2] N. Mohan, et al. "Pruning Edge Research with Latency Shears." HotNets, 2020.

# My definition of Ubiquitous Learning

![](_page_3_Picture_1.jpeg)

- Ubiquitous computing learning is the growing trend of embedding computational learning capability (generally in the form of microprocessors) into everyday objects to make them effectively communicate and perform useful tasks in a way that minimizes the end user's need to interact with computers as computers.
- Unlike desktop computing datacenter-based learning, ubiquitous computing learning can occur with any device, at any time, in any place and in any data format across any network.

# Why Ubiquitous Learning?

![](_page_4_Picture_1.jpeg)

• Rich applications, mostly driven by privacy needs

- Personalization for each device
- Offloading the computation & storage expense to devices
  - Workflow simplified: no need to maintain complex, clumsy datacenter systems (for both training and inference)
  - Risk minimized: benefits of decentralization

#### Some concrete evidences

![](_page_5_Picture_1.jpeg)

- FL is widely adopted on billions of phones already.
  - Google uses it for input method (Gboard)<sup>[1]</sup>
  - Alibaba uses it for e-commerce recommendations (Taobao)<sup>[2]</sup>
  - Apple uses it for voice assistant (Siri)<sup>[3]</sup>
  - ...
- IoT companies are seeking to train and adapt ML models on IoT devices (smart home scenario).

[1] https://ai.googleblog.com/2017/04/federated-learning-collaborative.html
[2] Chaoyue Niu, et al. "Billion-scale federated learning on mobile clients: A submodel design with tunable privacy." MobiCom, 2020.
[3] https://www.technologyreview.com/2019/12/11/131629/apple-ai-personalizes-siri-federated-learning/

#### A review of DL system evolution on cloud/edge TensorFlow (Google), PyTorch (FB), PaddlePaddle

![](_page_6_Figure_1.jpeg)

#### A review of DL system evolution on cloud/edge

![](_page_7_Figure_1.jpeg)

#### A review of DL system evolution on cloud/edge

![](_page_8_Picture_1.jpeg)

Many attemps at on-device inference: V V V [MobiCom'18] DeepCache, [WWW'19] An empirical study, [MobiSys'20] Elf, [TMC'19] DeepWear Caffe2 (FB), Core ML (Apple), ncnn (Tencent), Paddle Lite (Baidu), MACE (Xiaomi), SNPE (Qualcomm), etc.

MSRA-Beijing-202100422

Mengwei Xu(徐梦炜)@ CS Dept of BUPT

9

### HW-SW stack for UL

![](_page_9_Picture_1.jpeg)

Applications My research focus Input method, voice assistant, etc... Federated learning, Split learning, Learning Protocols on-device transfer learning, etc... **DL** Primitives **Operators, optimizers, BP, etc...** Training Middleware MNN, CoreML, TensorFlow, etc... **Operating System** Android, iOS, Linux, ROS, etc... Hardware ARM CPUs, GPUs, DSP, NPU, etc...

# 2 Key research questions in UL

![](_page_10_Picture_1.jpeg)

- 1. Training data is limited, non-IID, or even not labelled
  - Model accuracy heavily relies on data!

- 2. Devices have constrained hardware resources
  - Training a ML model is notoriously resource-hungry!

# 2 Key research questions in UL

![](_page_11_Picture_1.jpeg)

- 1. Training data is limited, non-IID, or even not labelled
  - Model accuracy heavily relies on data!

- 2. Devices have constrained hardware resources
  - Training a ML model is notoriously resource-hungry!

![](_page_11_Figure_6.jpeg)

### Outline

![](_page_12_Picture_1.jpeg)

- Ubiquitous Learning (UL): What and Why?
  - A system software perspective
  - The status quo
- A measurement study of UL
- Memory optimizations of UL
- Other on-going research of UL

# On-device training: a measurement study

- Target library: MNN<sup>[1]</sup> by Alibaba
- 6 Android devices
- 5 classic CNN models
  - LeNet, AlexNet, MobileNetv2, SqueezeNet, GoogLeNet
- CPU by default

Testing	Training	Training time (ms)				
platform	library	BS = 1	BS=2	BS=4		
Samsung Note 10	MNN	516	812	1365		
	DL4J	3,032	$6,\!129$	OOM		
RPI 3B+	MNN	6698	$10,\!651$	OOM		
	TensorFlow	10,468	$14,\!157$	$27,\!574$		
	PyTorch	48,274	$79,\!097$	OOM		

Device	Specifications	Yr.
Redmi Note 9 Pro	Snapdragon 720G, 6GB RAM	2020
Xiaomi MI 9	Snapdragon 855, 6GB RAM	2019
Huawei Mate 30	Kirin 990, 8GB RAM	2019
Meizu 16T	Snapdragon 855, 6GB RAM	2019
Samsung S8+	Snapdragon 835, 6GB RAM	2017
Huawei Honor 8	Kirin 950, 3GB RAM	2016

[1] https://github.com/alibaba/MNN

![](_page_14_Picture_1.jpeg)

- Training takes much more time than inference
  - Up to 17.8x gap, much larger than the FLOPs-gap

![](_page_14_Figure_4.jpeg)

![](_page_15_Picture_1.jpeg)

- Training takes much more time than inference
- GPU cannot speedup

![](_page_15_Figure_4.jpeg)

Mengwei Xu (徐梦炜) @ CS Dept of BUPT

- Training takes much more time than inference
- GPU cannot speedup
- Op-level breakdown
  - Raster: ops with shape transformation
  - Conv = 3 x Raster + MatMul

![](_page_16_Figure_6.jpeg)

![](_page_16_Picture_7.jpeg)

![](_page_17_Picture_1.jpeg)

- Training takes much more time than inference
- GPU cannot speedup
- Op-level breakdown
- Why?
  - The training support of MNN is still at very preliminary stage
  - Training is far more complex than inference: much more operators, dynamic weights update, variable batch size, etc...

# Memory footprint

![](_page_18_Figure_1.jpeg)

- Training is very memory-intensive
  - 16-32 is typically the max batch size supported by a high-end device (6~8 GBs)

![](_page_18_Figure_4.jpeg)

# Memory footprint

![](_page_19_Picture_1.jpeg)

- Training is very memory-intensive
  - 16-32 is typically the max batch size supported by a high-end device
  - Enough for a good convergence? NO!

![](_page_19_Figure_5.jpeg)

![](_page_19_Figure_6.jpeg)

![](_page_19_Figure_7.jpeg)

#### Federated setting (MobileNet-v2)

[1] Smith, Samuel L., et al. "Don't decay the learning rate, increase the batch size." arXiv preprint arXiv:1711.00489 (2017).

Mengwei Xu (徐梦炜) @ CS Dept of BUPT

# Memory footprint

- Training is very memory-intensive
  - 16-32 is typically the max batch size supported by a high-end device
  - Enough for a good convergence? NO!

![](_page_20_Picture_4.jpeg)

![](_page_20_Picture_5.jpeg)

![](_page_20_Picture_6.jpeg)

## System parameters

![](_page_21_Picture_1.jpeg)

• CPU parameters open rich trade-off among latency and energy

	CPU Conf	Time (s)		Energy (J)				
Lowest		Η	$\mathbf{M}$	$\mathbf{L}$	Η	$\mathbf{M}$	$\mathbf{L}$	
latency	1×	4.2	5.4	10.8	10.6	8.0	6.9	Lowest
	Big $2 \times$	2.6	3.2	6.4	8.9	7.7	7.0	energy
	<u>4</u> ×,	2.0	3.3	8.4	7.1	8.7	8.2	
	1×	25.0	33.9	57.8	10.4	7.2	3.1	
	Small $2\times$	13.3	18.0	31.8	10.1	8.4	4.8	
	$4 \times$	8.0	11.0	52.3	11.4	9.6	8.2	
	Hybrid $8 \times$	3.8	6.5	50.4	13.4	13.9	14.4	

MSRA-Beijing-202100422

![](_page_22_Picture_0.jpeg)

#### Implications and future directions

- Generating efficient training operators
  - NN compilers?
- Memory optimizations

> The next chapter of this talk

- Tuning system parameters
  - Model-level or operator-level?

### Outline

![](_page_23_Picture_1.jpeg)

- Ubiquitous Learning (UL): What and Why?
  - A system software perspective
  - The status quo
- A measurement study of UL
- Memory optimizations of UL
- Other on-going research of UL

### Design goals and principles

![](_page_24_Picture_1.jpeg)

- Goal: supporting larger batch size training with given upper bound of peak memory usage
- Principles
  - 1. Don't sacrifice the accuracy No quantization & sparsity
  - 2. Sacrifice latency (energy) as less as possible Making tradeoffs

# Splitting the minibatch

• A technique called microbatch (or virtual batch, ghost batch)

![](_page_25_Figure_3.jpeg)

# Splitting the minibatch

![](_page_26_Picture_1.jpeg)

- A technique called microbatch (or virtual batch, ghost batch)
- Why it's not enough?

![](_page_26_Figure_4.jpeg)

[1] Summers, Cecilia, and Michael J. Dinneen. "Four things everyone should know to improve batch normalization." ICLR 2020.

### How about another technique

- (X) Memory swapping?
  - Mobile devices have integrated memory for heterogeneous processors
  - 40x slower according to our preliminary experiments
- (?) Recompute<sup>[1]</sup>?
  - Under implementation but seems more promising
  - An efficient memory pool is needed

#### Memory reduced to $O(\log n)$ Forward latency increased to $O(n \log n)$

![](_page_27_Figure_8.jpeg)

[1] Tianqi Chen, et al. "Training deep nets with sublinear memory cost." arXiv:1604.06174 (2016).

![](_page_27_Picture_13.jpeg)

### Combining them

![](_page_28_Picture_1.jpeg)

- Determining the appropriate ghost batch size and recomputed nodes to achieve the optimal performance
  - Hardware-specific, large search space

### Outline

![](_page_29_Picture_1.jpeg)

- Ubiquitous Learning (UL): What and Why?
  - A system software perspective
  - The status quo
- A measurement study of UL
- Memory optimizations of UL
- Other on-going research of UL

# Some on-going research to share

![](_page_30_Picture_1.jpeg)

1. DSP-based, integer-only on-device training algorithm and library

2. Adaptive network topology and protocol for cross-silo federated learning

![](_page_31_Picture_0.jpeg)

![](_page_31_Picture_1.jpeg)

- Machine (deep) learning is happening everywhere at anytime
- The system support for such ubiquitous learning is still at very preliminary stage so many open problems!
- Open to discussion and collaboration on it

Outstanding PhD students from BUPT/PKU working on UL  $\Rightarrow$ 

 $Our code \Rightarrow https://github.com/UbiquitousLearning$ 

![](_page_31_Picture_7.jpeg)