

Mobile & Ubiquitous Computing, 2020 Spring

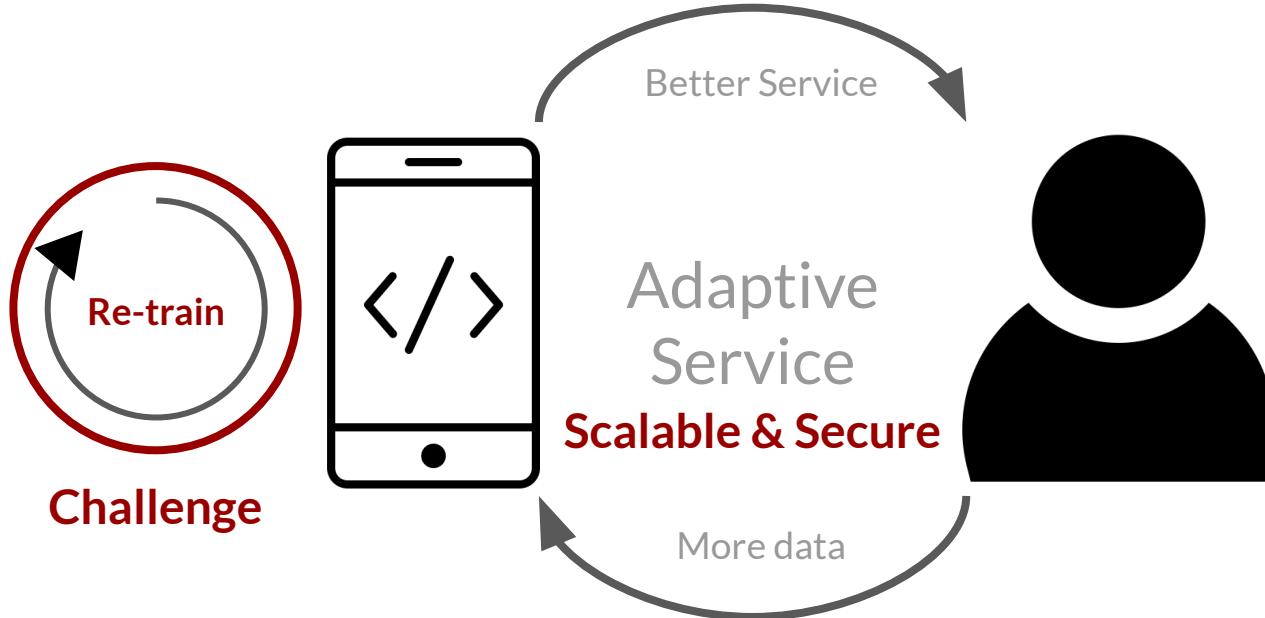
FedEx: Scalable Framework for Mobile User-Adaptive Deep Learning Training

Ahnjae Shin*, Heeseung Yun*, Kyunggeun Lee*, Taebum Kim*

Team 1

**Equal Contribution*

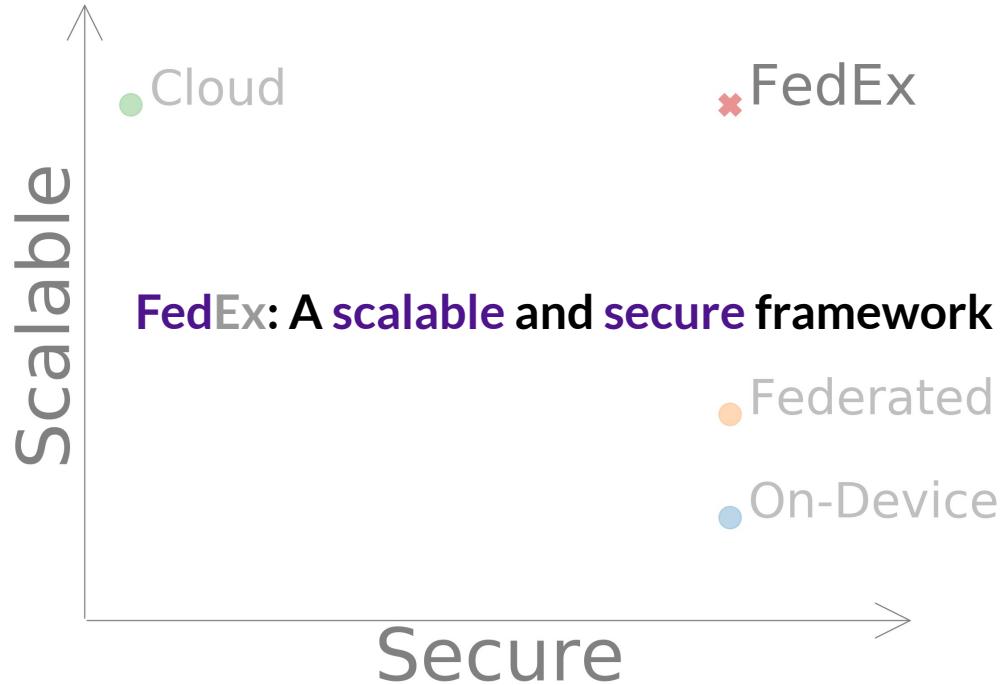
Motivation

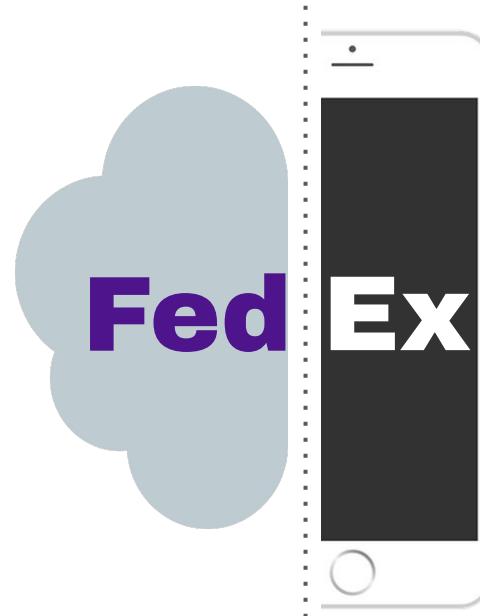


Scalable Secure

Able to scale to large models

Able to train without raw data



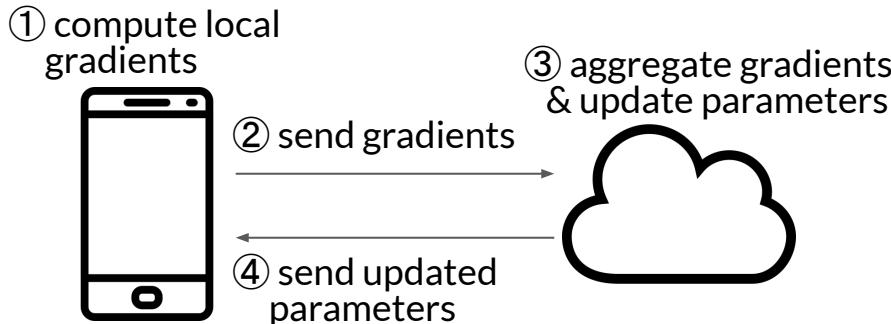


We ship fresh models daily

Federated Learning

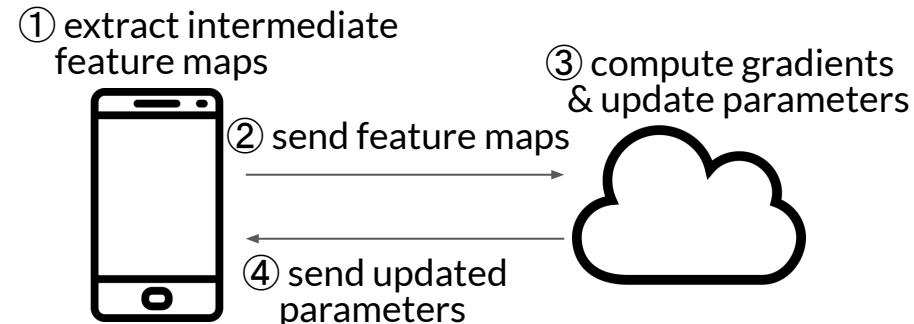
FedAvg

- Compute in mobile devices, update aggregated gradients in clouds
- Computation: mobile device
- Communication rounds: $O(NP)$



FedEx

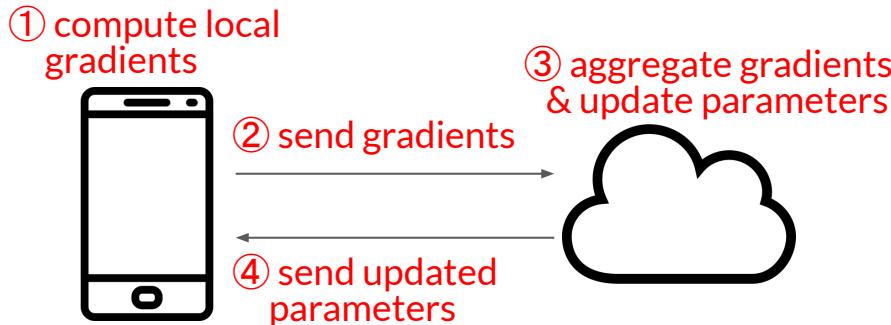
- Extract intermediate feature maps in mobile devices, train in clouds
- Computation: cloud
- Communication rounds: constant



Federated Learning

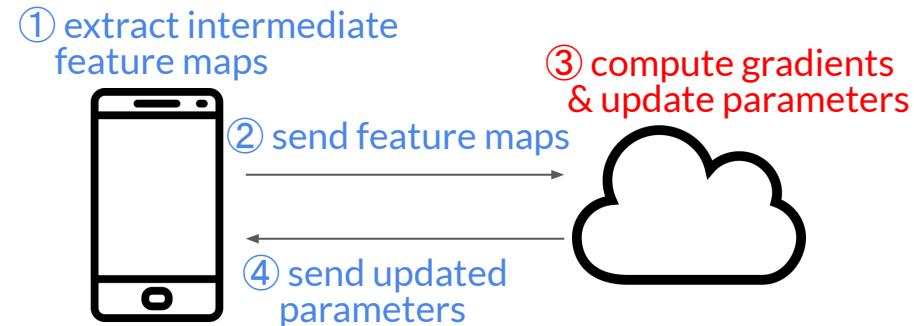
FedAvg

- Compute in mobile devices, update aggregated gradients in clouds
- Computation: mobile device
- Communication rounds: $O(NP)$



FedEx

- Extract intermediate feature maps in mobile devices, train in clouds
- Computation: cloud
- Communication rounds: constant



Technical Details

Challenges & Solution

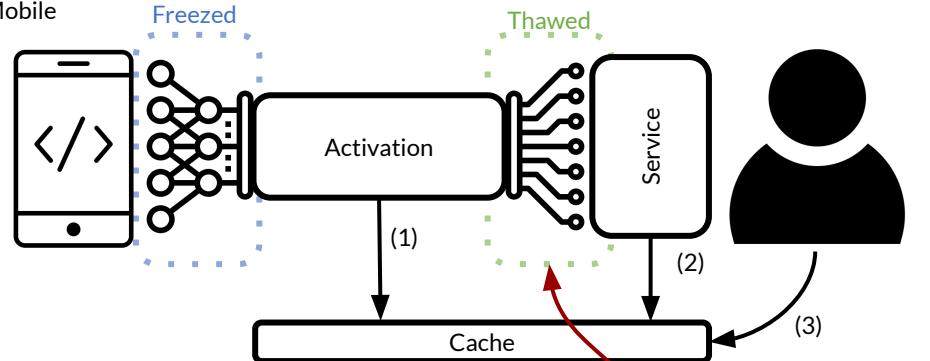
1. Challenges

- a. Scarce computation power
 - i. Mobile GPUs are not as powerful as server GPUs.
(Server GPUs are generally 10~1000x more powerful)
- b. Communication overhead
 - i. Should push(pull) gradients(parameters) every training step.
 - ii. Network bandwidth is often very restricted

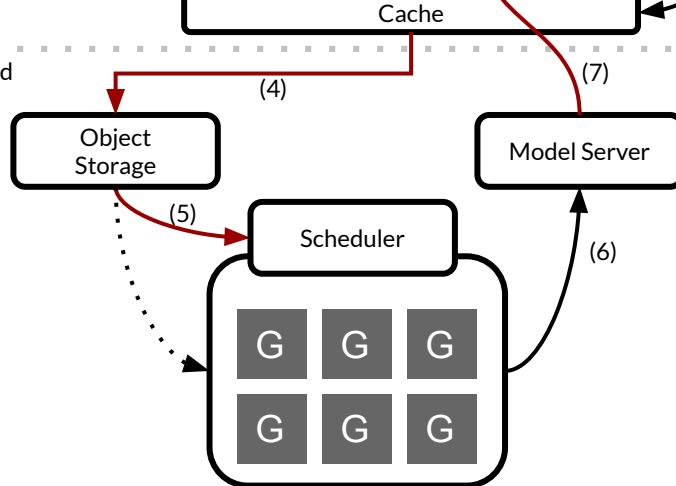
2. Solution: Training with feature maps

System Architecture

(a) Mobile



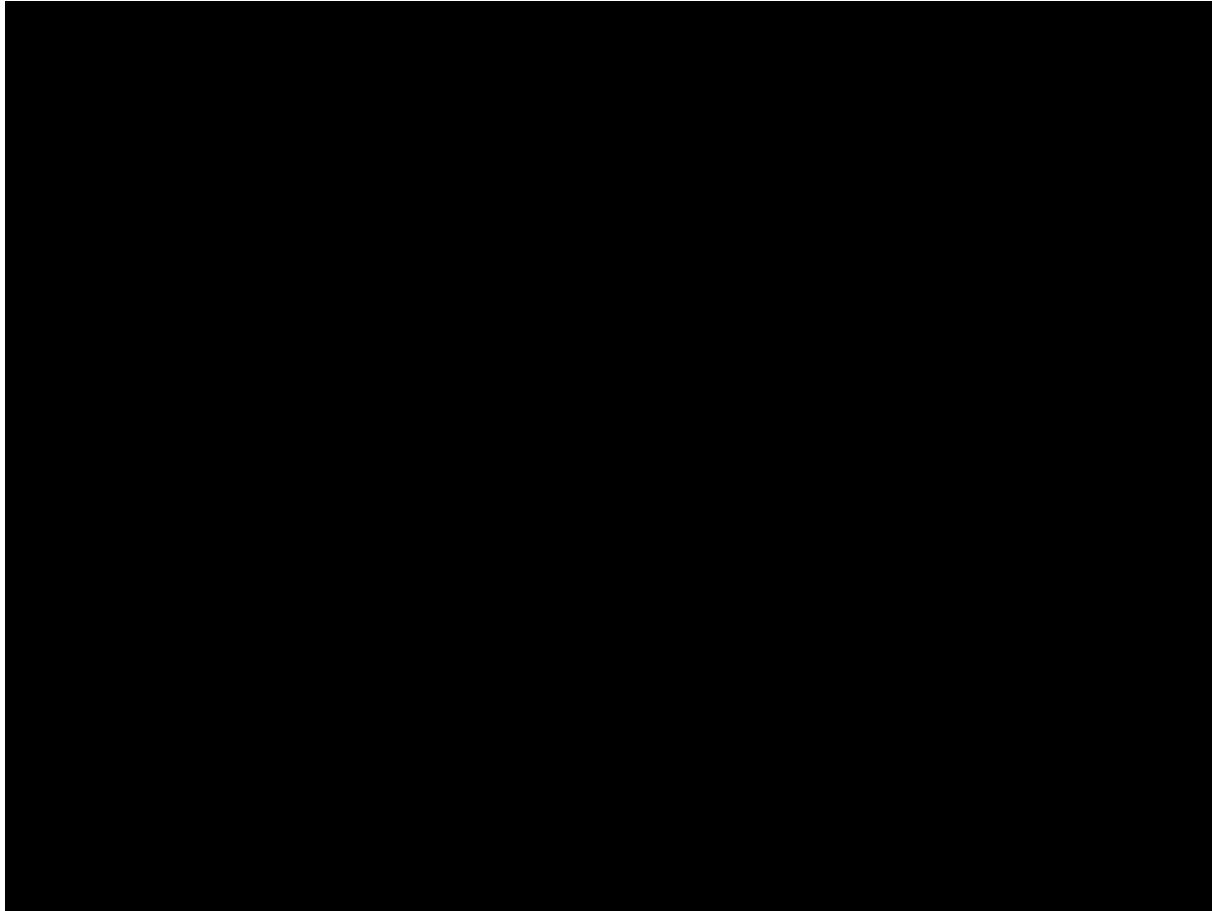
(b) Cloud



1. Cache intermediate features
2. Cache output
3. Save user feedback
4. Push intermediate features
5. Notify scheduler
Create new training session
6. Workers upload checkpoint to model server
7. **Best performing model can be pulled from model server**

Demo

Step-by-step Walkthrough



Evaluation

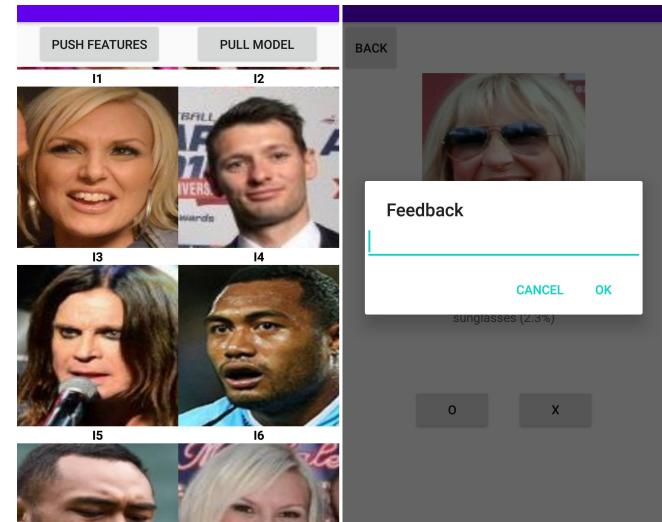
Evaluation Environment

1. Client Application

- Android face recognition application
- InceptionResNet-V1 as backbone network
- VGGFace2 dataset used

2. Server

- Two 18-core Intel Xeon E5-2695 @2.10 GHz CPU
- 6 NVIDIA Titan Xp GPUs
- Communicates with client via HTTP request

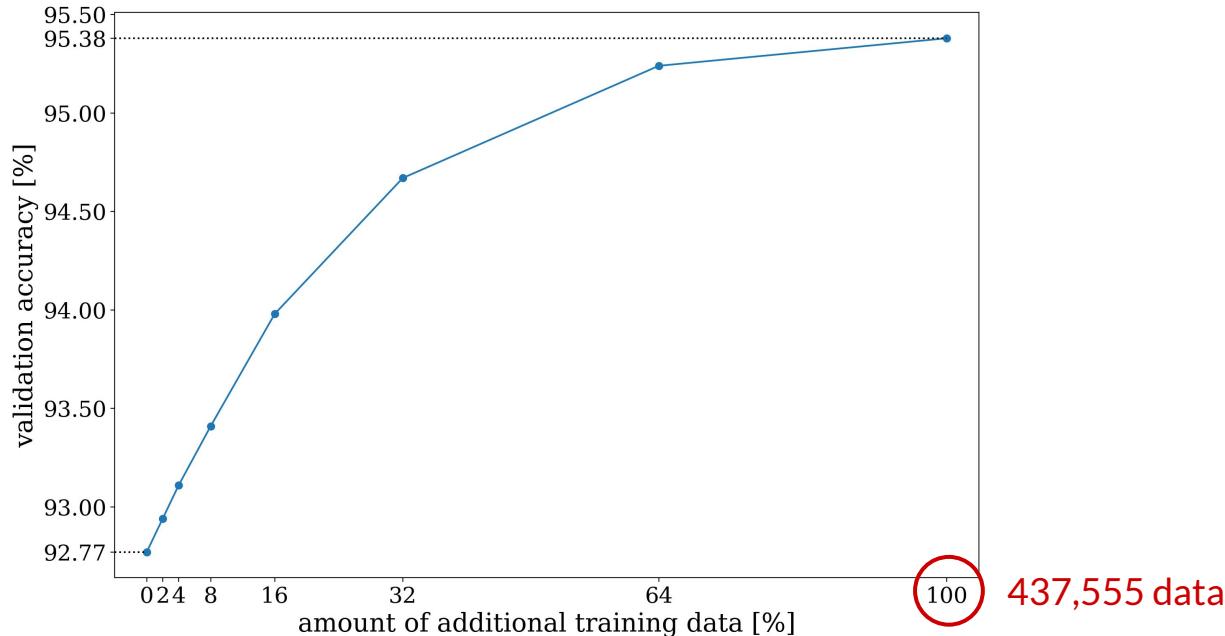


Evaluation & Success Criteria

1. Accuracy
2. Training Throughput
3. Energy Consumption

Evaluation 1. Accuracy

- How much the accuracy could increase by additional training?
- At least how much additional data is needed to increase the accuracy?



Evaluation 2. Training Throughput

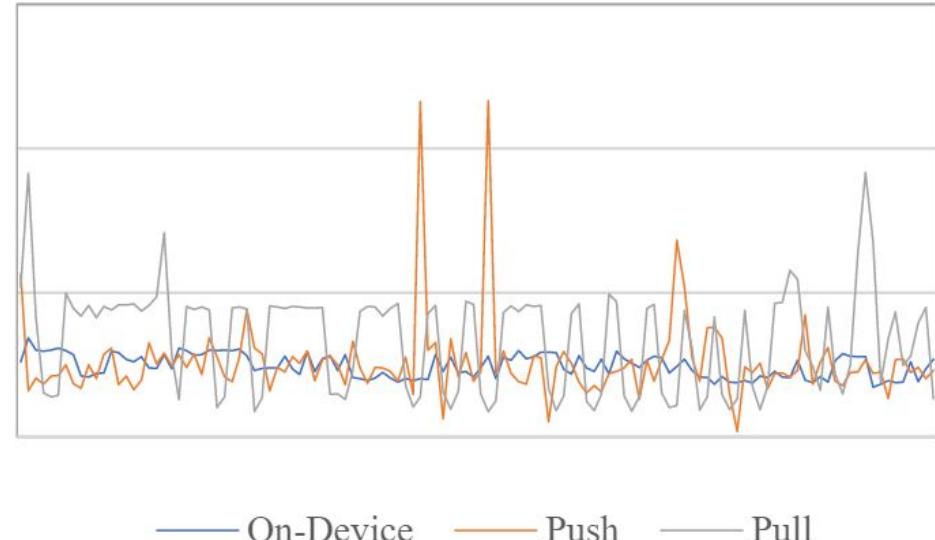
- Training throughput
 - Less than 4 minutes for additional training (231 secs) using Base (Freeze 6) model
- Communication latency
 - Push phase: Freeze 5 < FedAvg < Freeze 2,3.
 - Pull phase: same latency
 - However, FedAvg pushes/pulls more often ($O(1)$ vs. $O(N)$)

Model	Throughput [images / sec]	# Params	Feature
Full Model	487.56	27.91M	$(160 \times 160 \times 3)$
Freeze 1	573.36	27.88M	$(38 \times 38 \times 64)$
Freeze 2	699.13	27.29M	$(17 \times 17 \times 256)$
Freeze 3	963.26	26.9M	$(17 \times 17 \times 256)$
Freeze 4	2963.88	18.3M	$(8 \times 8 \times 896)$
Freeze 5	28078.46	5.35M	(1792)
Freeze 6	41869.68	4.43M	(512)

	Model	latency (ms)
Push	Base (Freeze 6)	1525.408
	Freeze 5	5487.633
	Freeze 4	281920.3
	Freeze 3	375950.0
	Freeze 2	375950.0
	Freeze 1	427250.4
	FedAvg	36111.52
Pull		25986.07

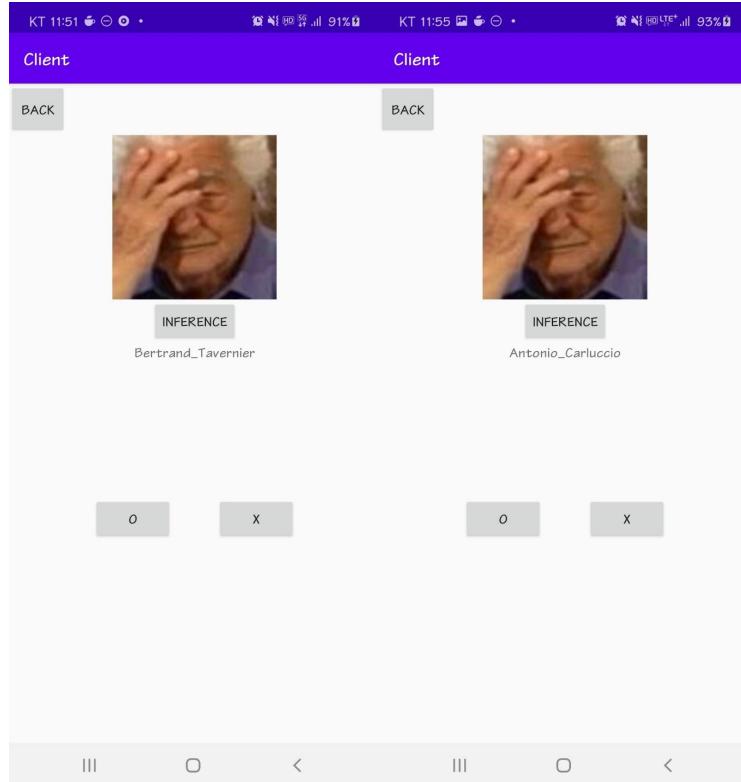
Evaluation 3. Energy

- Push (spiking; network-oriented) vs. Pull (cpu-oriented)
- How long does the on-device training take?
 - 7.9mAh/min.¹



¹ From Galaxy S7 system battery profiler

Qualitative result



Related Work

TensorFlow Lite Example On-device Model Personalization

This example illustrates a way of personalizing a TFLite model on-device without sending any data to the server. It builds on top of existing TFLite functionality, and can be adapted for various tasks and models.

See more in [TensorFlow Blog: Example on-device model personalization with TensorFlow Lite](#)

Quickstart

Pre-requisites:

- [Android Studio](#).
- [Python 3.5+](#).
- [virtualenv](#) (usually comes preinstalled with Python).
- Physical Android device with camera.

Install and run the application

You can either build and run the application inside Android Studio or run command line to do so.

If you want to use Android Studio, first import the project into Android Studio (point it to the top-level `build.gradle` file), connect your Android device to your machine, and use the `Run` button in Android Studio. If the `Run` button is inactive, first add an Android application build configuration for the `app` module.

If you want to build and install on Linux directly, you can first switch to the `android` folder, then execute

```
gradle wrapper  
./gradlew build
```

- Dec 2019 Google TF Lite team project
- TF-Lite personalization system
 - Static deployment
 - On-device update
- FedEx
 - Pull model weights from cloud
 - Training is done in cloud

Discussion

1. Privacy Concerns

- Varying degrees of attack surface
 - Naive cloud learning (raw data → data hijacking)
 - FedEx (intermediate feature → reverse engineering)
 - FedAvg (gradients → model poisoning)

2. Multi-Client Scenarios

- Model update in idle state
- How to achieve real scalability
- Synergistic effects among users

Limitation

- Approach
 - Limited range of application: transfer learning with freezed layers
 - Less private if few or no freezed layer
- Experiments
 - Weak baseline reproduction
 - No experiments for unseen label (new class), false feedback

Deliverable

1. 6-page Short technical report
2. Demo app on GitHub repo ([link](#))

FEDEx: Scalable and Secure System for User-Adaptive Deep Learning Training in Mobile Devices

Ahnjae Shin*
yuyupopo@snu.ac.kr
Seoul National University

Kyunggeun Lee*
initium0917@snu.ac.kr
Seoul National University

Heeseung Yun*
terry9772@snu.ac.kr
Seoul National University

Taebum Kim*
k.taebum@snu.ac.kr
Seoul National University

ABSTRACT

All machine learning (ML) services need to retrain their models periodically to accommodate new users and use cases. Especially with a great amount of personal data generated on the application's behalf, using personal data to continuously enhance a ML model could greatly benefit app retention rate. However, previous framework lack scalability or security in training deep learning models with personal user data. To address this problem, we propose FedEx, a mobile-cloud hybrid system to utilize both mobile and cloud resources in a scalable and secure manner. FedEx applies transfer learning to train models with additional data. Our system is scalable by removing mobile devices out of training loop, and secure as it mitigates the security risk of sending raw data. Experimental results suggest that FedEx can be adapted for models with a variety of capacity, as well as provide more energy-efficient and faster training process compared to representative baselines.

CCS CONCEPTS

- Human-centered computing → Ubiquitous and mobile computing systems and tools; • Computing methodologies → Distributed artificial intelligence.

ACM Reference Format:

Ahnjae Shin, Heeseung Yun, Kyunggeun Lee, and Taebum Kim. 2020. FedEx: Scalable and Secure System for User-Adaptive Deep

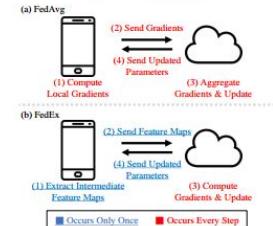
Learning Training in Mobile Devices. In *Proceedings of 21st International Workshop on Mobile Computing Systems and Applications (MobiMobile '20)*. ACM, New York, NY, USA, Article 4, 6 pages.

¹Equal contribution

²https://doi.org/10.4751/123_4

1 INTRODUCTION

As smartphones have become prevalent in the past decade and are producing a greater amount of personal data than ever, there is a great chance of improving user experiences using personal data in mobile devices. In the field of machine learning, applications featuring AI services have a high potential of improving service quality using personal data from mobile devices. Mobile applications have the opportunity to collect additional training data in various ways such as



Project Management

Change of Scope

- Removal
 - Handwriting recognition model
 - Strict baseline reproduction
- Addition
 - Face recognition model
- Why?
 - Necessity of additional training with personal data
 - Not successful with implementing baseline method can train models

Timeline

Week	Progress
1	Brainstorming / Literature Review / Proposal Preparation
2	System & Architecture Design
3	Client Mockup Implementation
4	Algorithm & Model Design
5	Cloud Training Scheme Implementation
6	Project Demo
7	Client-Side Feature Implementation (push/pull, UI updates)
8	Server-Side Feature Implementation (model, post-training)
9	Evaluation / Paper Writing
10	Final Presentation

Role & Contribution

Task	Ahnjae	Kyunggeun	Heeseung	Taebum
Client application		O	O	
Server	O			O
Client-server integration	O	O		
Face recognition model				O
Experiment		O	O	O
Paper writing	O	O	O	O

GitHub Usage Stat

Apr 5, 2020 – Jun 8, 2020

Contributions to master, excluding merge commits



Contributions: Commits ▾

Author	Label	Projects	Milestones	Reviews	Assignee	Sort
HS-YN						0 Open ✓ 12 Closed
HS-YN						#12 by yuyupopo was merged Jun 8, 2020
ktaebum						#11 by yuyupopo was merged Jun 8, 2020 • Approved
HS-YN						#10 by HS-YN was merged Jun 6, 2020
ktaebum						#9 by ktaebum was merged Jun 8, 2020
ktaebum						#8 by ktaebum was merged Jun 6, 2020
kyunggeun-lee						#7 by kyunggeun-lee was merged Jun 6, 2020 • Approved
HS-YN						#6 by HS-YN was merged Jun 6, 2020
kyunggeun-lee						#5 by kyunggeun-lee was merged Jun 3, 2020
yuyupopo						#4 by yuyupopo was merged May 11, 2020
ktaebum						#3 by ktaebum was merged May 11, 2020
ktaebum						#2 by ktaebum was merged May 11, 2020
HS-YN						#1 by HS-YN was merged May 10, 2020

Conclusion

Key Contribution

- Propose a **scalable, energy-efficient** end-to-end framework for mobile user-adaptive training of deep learning while providing a certain level of privacy
- Verify the effectiveness of FedEx using a sample mobile face recognition application with respect to *accuracy, throughput* and *energy*
- Source code is publicly available :)

Lesson Learned

- Creating a solid abstraction across different computing devices is hard
- Caution when sharing experiment scripts across team members
- Android Studio ...
- TensorFlow Lite still has some unstable features (should be developed more)
- Hard to **evaluate** our system!





We ship fresh models daily

Thank you