

### An Ultra-Fast On-Device CNN Finetuning Approach for Resource-Limited Edge Devices

### Hiroki Matsutani (Keio University, Japan)









# **On-device learning (ODL) for IoT devices**

 Motivation for neural network training at edge side Addressing the gap between pretrained model and deployed environment by updating the model on-device [1,2]



[1] Mineto Tsukada et al., "A Neural Network-Based On-device Learning Anomaly Detector for Edge Devices", IEEE Trans. on Computers (2020).
 [2] Kazuki Sunaga et al., "Addressing Gap between Training Data and Deployed Environment by On-Device Learning", IEEE Micro (2023).

### **On-device learning (ODL) for IoT devices**

 Lightweight sequential training (since 2017) [1]







[1] Mineto Tsukada et al., "A Neural Network-Based On-device Learning Anomaly Detector for Edge Devices", IEEE Trans. on Computers (2020).

# This talk: On-device finetuning for DNNs

We focus on DNN/CNN
 including vision & LLM tasks















### **Baseline finetuning methods**

#### FT-Last [1]



#### Weights of the last layer (W<sup>3</sup>, b<sup>3</sup>) are updated

[1] Haoyu Ren et al., "TinyOL: TinyML with Online-Learning onMicrocontrollers", IJCNN'21.

### **Baseline finetuning methods**







Weights of the last layer (W<sup>3</sup>, b<sup>3</sup>) are updated

### Trainable adapters are attached to all layers

Trainable adapter is attached to the last layer

[1] Haoyu Ren et al., "TinyOL: TinyML with Online-Learning onMicrocontrollers", IJCNN'21.

[2] Edward J. Hu et al., "LoRA: Low-Rank Adaptation of Large LanguageModels", arXiv:2106.09685 (2021).

# **Our proposed approach: Skip-LoRA**

Skip-LoRA can reduce the backward computation



### Our proposed approach: Skip2-LoRA

Skip2-LoRA can reuse forward computation results

These values are needed to compute gradients of the adapters



**Forward** computation results of the base model are **cached** 

**Base model** 

### $W^{0,3}$ $W^{1,3}$ $W^{1,3}$ $W^{2,3}$ $W^{0,3}$ $W^{2,3}$ **Adapters** Skip-LoRA [1]

### Our proposed approach: Skip2-LoRA

Skip2-LoRA can reuse forward computation results

(# iterations) = (# samples) / (Batch size) × (# epochs)



### Our proposed approach: Skip2-LoRA

Skip2-LoRA can reuse forward computation results

(# iterations) = (# samples) / (Batch size) × (# epochs)



#### **Base model**

# iterations is typically larger than 1; so, we can skip most of the forward computation [1]

12

# **Skip2-LoRA for CNNs: Model**

- Pretrained with fashion-MNIST 

   This gap reduces
   accuracy
   accuracy
- Finetuned with Rotated fashion-MNIST (75 deg)
- Tested with Rotated fashion-MNIST (75 deg)





### **Skip2-LoRA for CNNs: Raspberry Pi**

### • Raspberry Pi Zero 2W ARM Cortex-A53 @1GHz















# **Skip2-LoRA for CNNs: Results**

- In this work, Skip2-LoRA [1] is applied to CNNs
- An aggressive 4-bit quantization is applied to the forward cache to reduce the memory footprint

Model	Accuracy	FT time @RPZ2	Cache size
No Finetuning (FT)	9.18 %		14 Honor
FT-Last	60.94 %	18.09 sec	in the second second
LoRA-Last	53.81 %	18.09 sec	
LoRA-All	<b>75.59 %</b>	<b>,114.15 sec</b>	Raspberry Pi Zero 2W (aka \$15 computer)
Skip-LoRA	73.54 %	<b>19.84 sec</b>	
Skip2-LoRA	73.54 %	3.90 sec	∕7,336 kB
Quant Skip2-LoRA	<b>74.02</b> %	<b>4.27 sec</b>	<b>∖1,036 kB</b>

\*Number of FT samples: 1024, Number of epochs for FT: 10

## **Skip2-LoRA for CNNs: Low-cost FPGA**

Accuracy: 76%

Finetune time: 0.34 sec

Test7: Skip2-LoRA on FPGA [x4]

Tested with Rotated fashion-MNIST (75 deg)

Accuracy: 85.059% (871/1024) Epoch: 10.362ms (0.324ms/batch)

### • Low-cost FPGA board AMD KV260 Starter Kit Finetune time: 0.34 sec



# **Skip2-LoRA for CNNs: Summary**

 Ultra-fast on-device CNN finetuning for resource-limited edge devices Our approach: algorithm and architecture codesign

results



acceleration

4.27 sec  $\rightarrow$  0.34 sec

Step 1: Step 2: Step 3: **Reusing computation FPGA-based** New connection

pattern

114.15 sec  $\rightarrow$  19.84 sec

Yes

19.84 sec  $\rightarrow$  4.27 sec



 $W^{1,3}$   $W^{2,3}$   $W^{2,3}$  [1] Keisuke Sugiura et al., "InstantFT: An FPGA-Based Runtime Subsecond Fine-tuning of CNN Models", arXiv:2506.06505 (2025).