

Heterogeneity Aware Federated Learning and Demand-Specific Inference

Jianming Tong, Alind Khare, Alexey Tumanov, Tushar Krishna

School of Computer Science [1], School of Electrical and Computer Engineering

Introduction

Application: Health-care

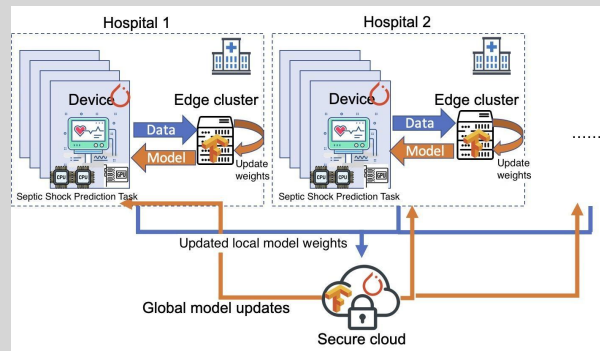
Demand:

- highly sensitive data with different distributions geographically distributed clients.
- heterogeneous computational power
 - mobile phones may range from GPU (dis)abled or 4-12GB RAM etc.

Problem Statement

- Privacy-aware ML training for highly sensitive heterogeneous data.
- Serving various accuracy and latency demand.

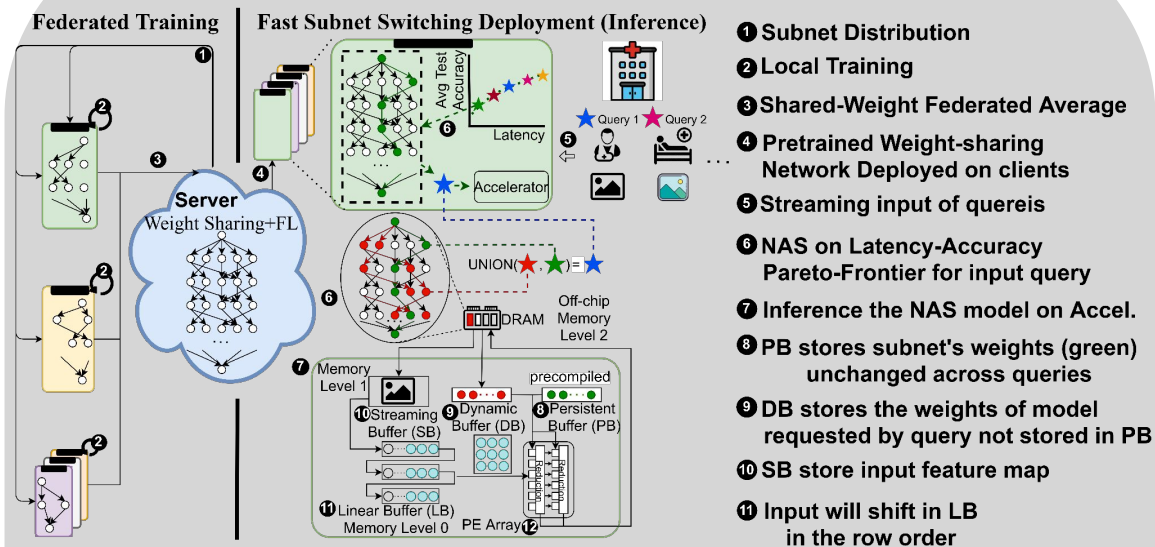
Overall System High-level Overview



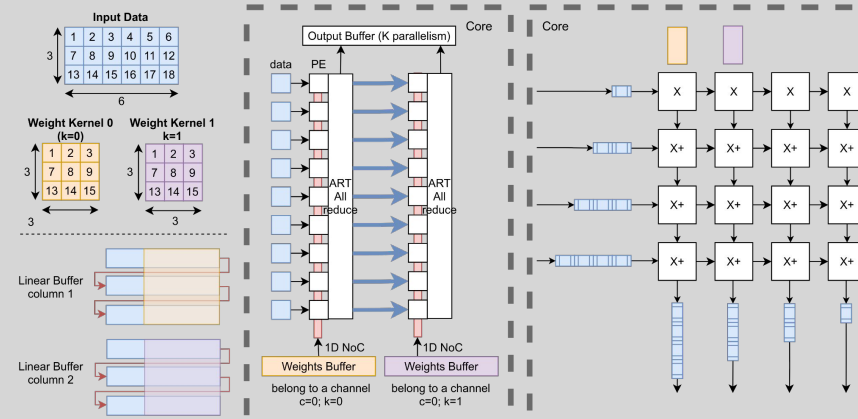
State-of-the-Art

- Federated Learning: Federated Algorithm [1]
 - Homogeneous clients assumption
- Once-For-All [2] requires centralized data
- Xilinx Vitis AI needs minutes to switch models [3]

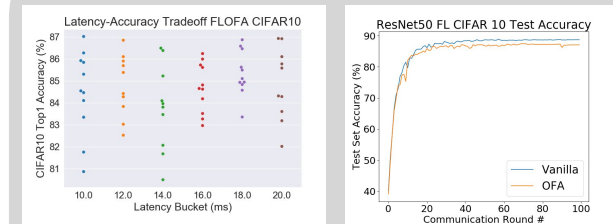
Proposal System



Plug-in-and-play PE Array Design



Federated Learning Evaluation



Results

- Model Family for different Accuracy-Latency
- Same level accuracy as SOTA

Quick Subnet Switching Inference

Xilinx Vitis AI with Xilinx Deep Processing Unit (DPU)

Components	Utilization
CLB LUT	49%
CLB Regs	86%
CLB Logic	92%
DSP	87%
MMCM	33%
BRAM	59%



Results

- 12 FPS for ResNet-50
- Model Switching needs several minutes