# Communication-Efficient Federated Learning for Fine Tuning Large Language Models

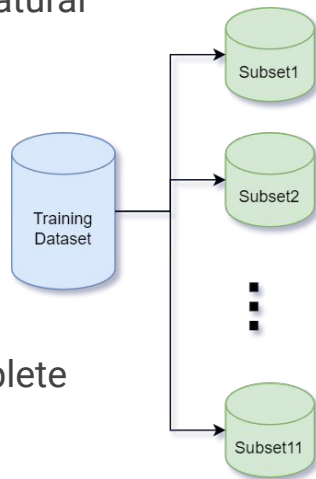Yuankai Cai, Ching Yuan Kung, Hui Chun Mo, Michael Tang, Chengkai Yao, Xiaoyue Zhang

Professor: Ling Liu

CS6220 Demo - Group 17

# Dataset

The **GLUE CoLA** (Corpus of Linguistic Acceptability) dataset is used for evaluating natural language understanding systems. It focuses on sentence-level classification tasks.

- Content: Contains sentences labeled as **grammatically correct or incorrect**.
- Size: **Consists of 10657 sentences**, 0.38 MB
- Source: The COLA consists of English acceptability judgments drawn from 23 **books and journal articles** on linguistic theory.

The **GLUE sst2** (the Stanford Sentiment Treebank) is a corpus that allows for a complete analysis of the compositional effects of sentiment in language.

- Content: Contains sentences labeled as **positive or negative based on the sentiment of the sentence.**
- Size: **Consists of 70000 sentences**, 7.22 MB
- Source: The sst2 consists of sentences from movie reviews and human annotations of their sentiment.
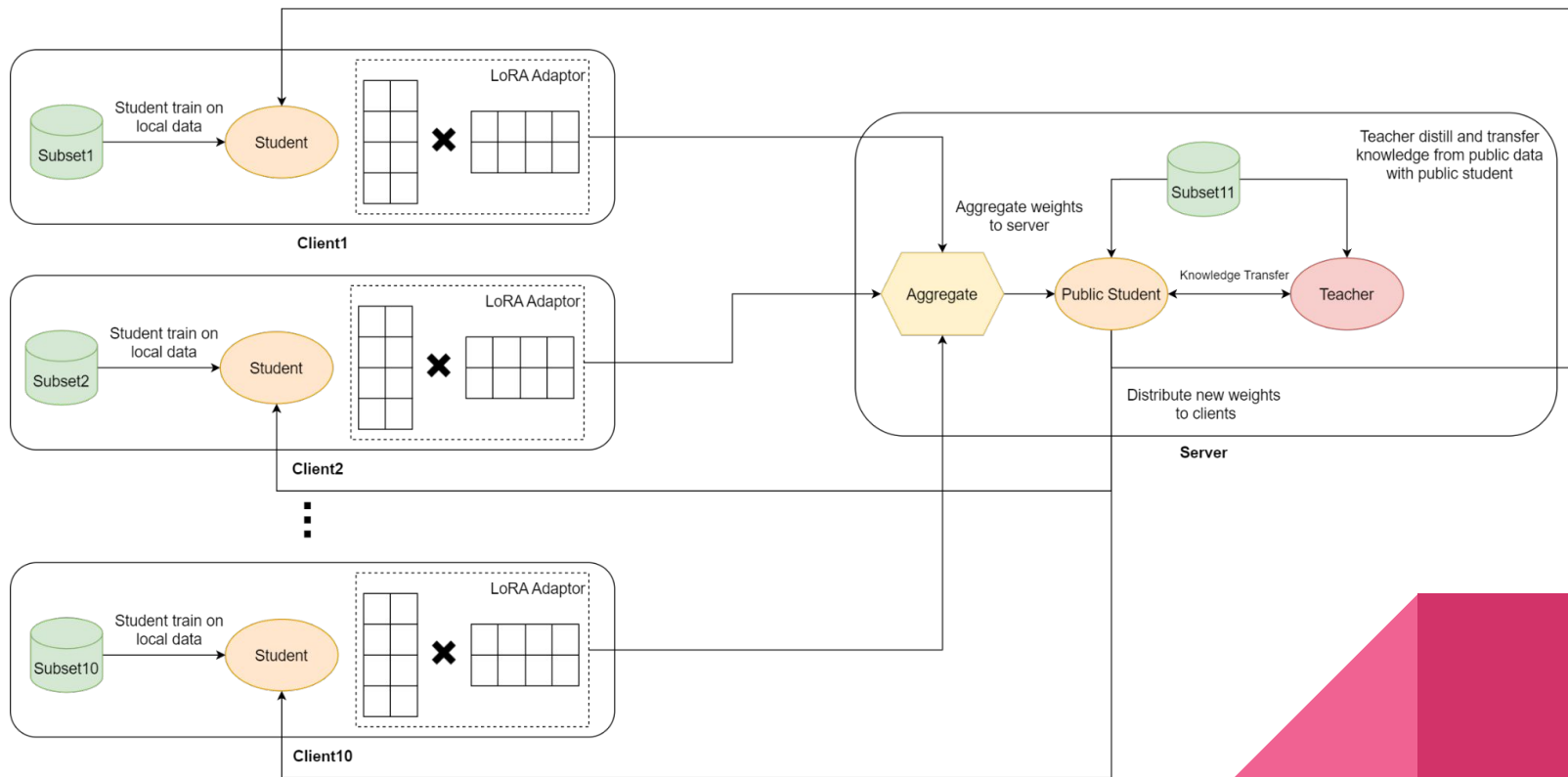
# Training dataset

## CoLA training samples

| Data | Label |
|------|-------|
| the greatest musicians | 1 |
| cold movie | 0 |
| with his usual intelligence and subtlety | 1 |
| redundant concept | 0 |
| swimming is above all about a young woman 's face , and by casting an actress whose face projects that woman 's doubts and yearnings , it succeeds . | 1 |

## sst2 training samples

| Data | Label |
|------|-------|
| hide new secretions from the parental units | 1 |
| contains no wit , only labored gags | 0 |
| the greatest musicians | 1 |
| with his usual intelligence and subtlety | 1 |
| by far the worst movie of the year | 1 |

# Architecture

# Screenshot

# Method Comparison



Accuracies by method and dataset

|  | Finetune | KD | LoRA | Fed-finetune | Fed-LoRA | Fed-LoRA-KD |
|---|---|---|---|---|---|---|
| cola | 0.787 | 0.760 | 0.797 | 0.802 | 0.756 | 0.760 |
| sst2 | 0.897 | 0.876 | 0.889 | 0.907 | 0.881 | 0.892 |

# Individual Model Comparison

# Easy Case

- FedKDLoRA has better accuracy than FedLoRA with different LoRA rankings

# Hard Case



Fed vs FedLoRA vs FedKDLoRA on SST2

- FedKDLoRA and FedLoRA underperforms Simple FL finetuning in model accuracy

# Summary

- Combined LoRA (Low Rank Adaptation) and KD (Knowledge Distillation) into Federated Learning
- Compared performance between Fed-LoRA and Fed-LoRA-KD
- Slightly better performance overall

# Reference

Dataset:
https://huggingface.co/datasets/glue/viewer/cola
https://huggingface.co/datasets/glue/viewer/sst2

FedKD:
https://arxiv.org/pdf/2108.13323.pdf
https://github.com/wuch15/FedKD

FedLoRA:
https://arxiv.org/pdf/2310.13283.pdf