# Real-time, robust, and reliable (R^3) machine learning across wireless networks

By Akshar Vedantham, Kirthana Ram, Varun Kota

Advisors - Prof. Anand Sarwate and Prof. Waheed Bajwa

# Project Introduction

- Phones, cars, and other devices will all want to start using ML/AI applications

- Leverage the cloud to help them with this

- Issue: Latency, security

# Example Scenario: Security System

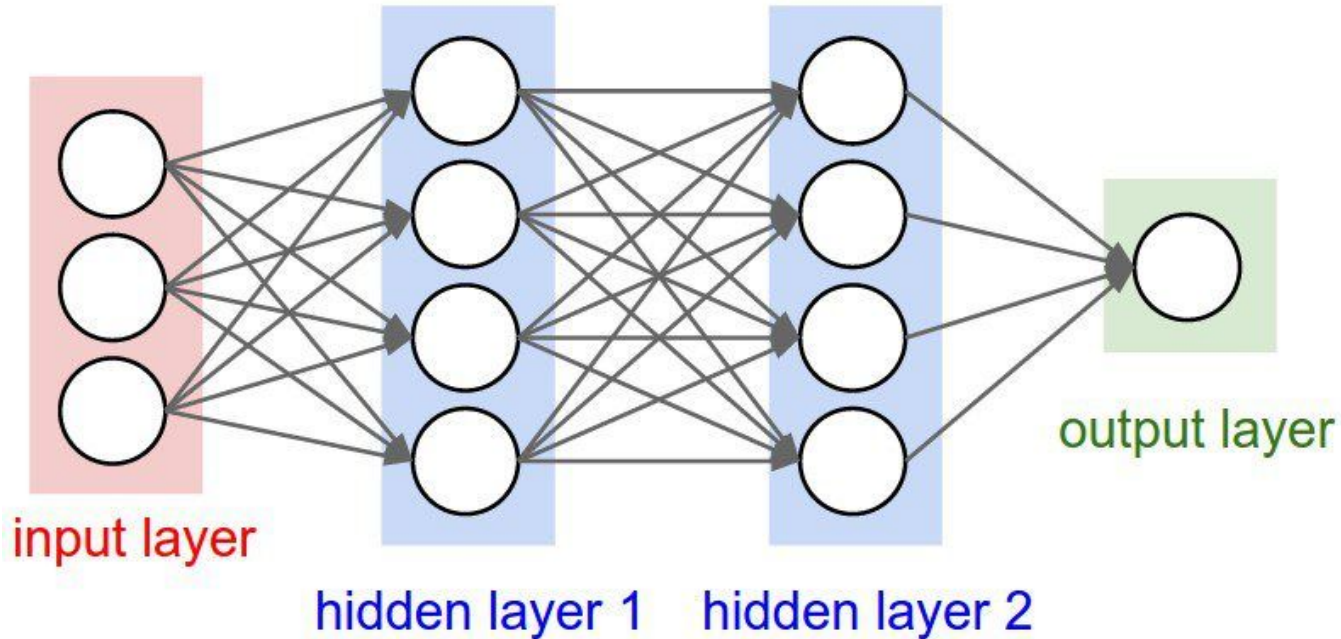# Problem: **When** should phones offload to the cloud and ask for help?

1.) AI Model
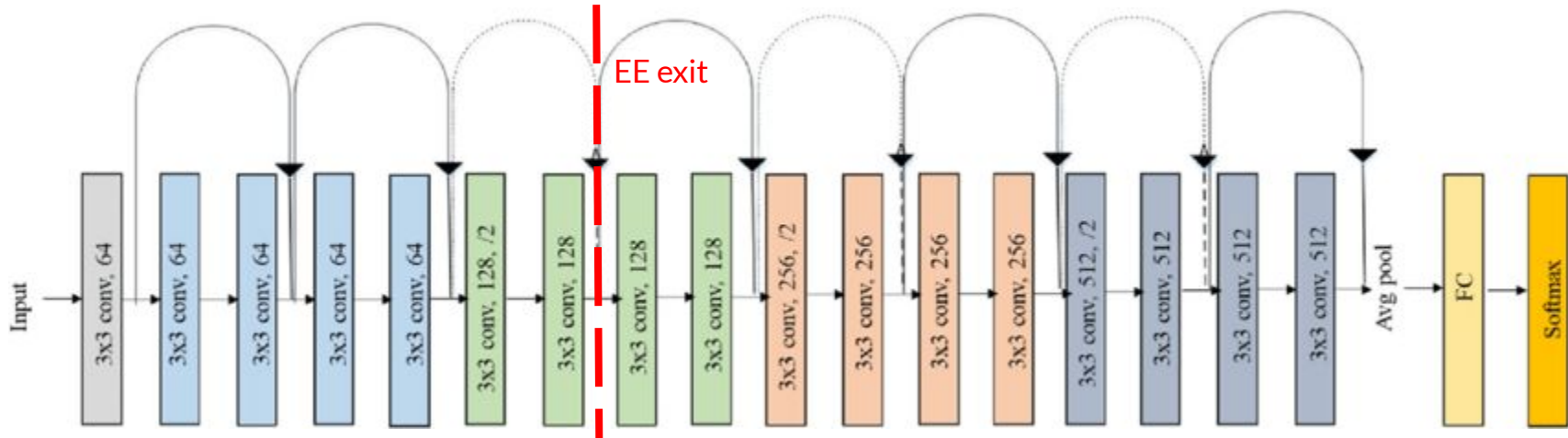2.) Strategies to offload
3.) Network conditions
4.) Evaluate

# The AI Model

# What is a Neural Network?



input layer

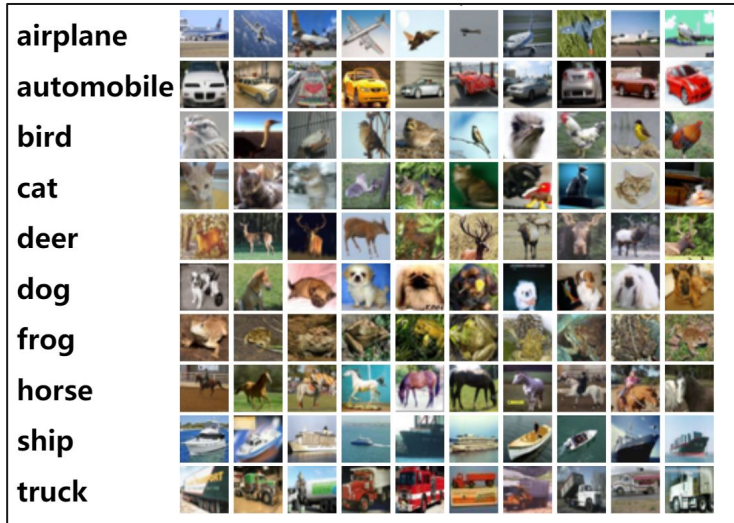hidden layer 1   hidden layer 2

output layer

# Choice of Model - ResNet18

- Popular model available
- Accuracy of 92.42%
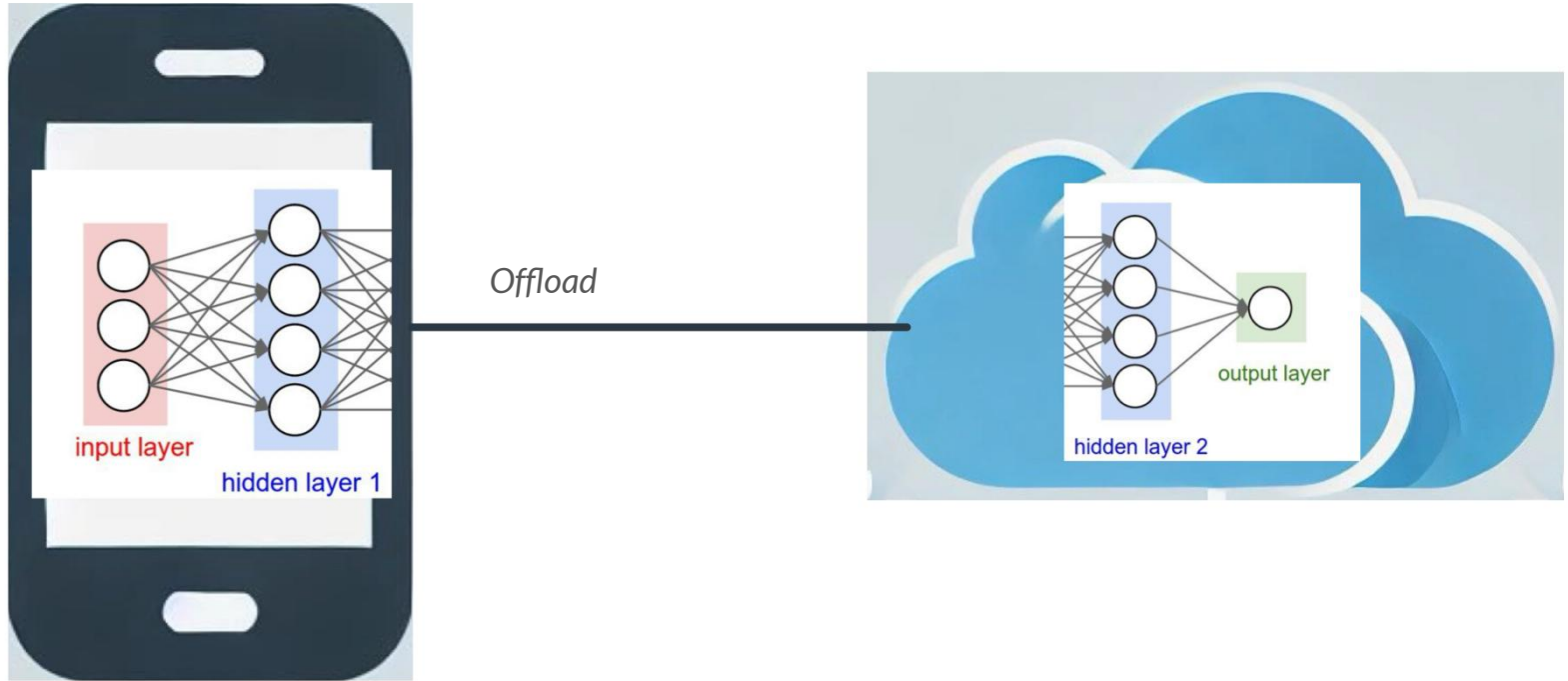
# CIFAR-10 Dataset



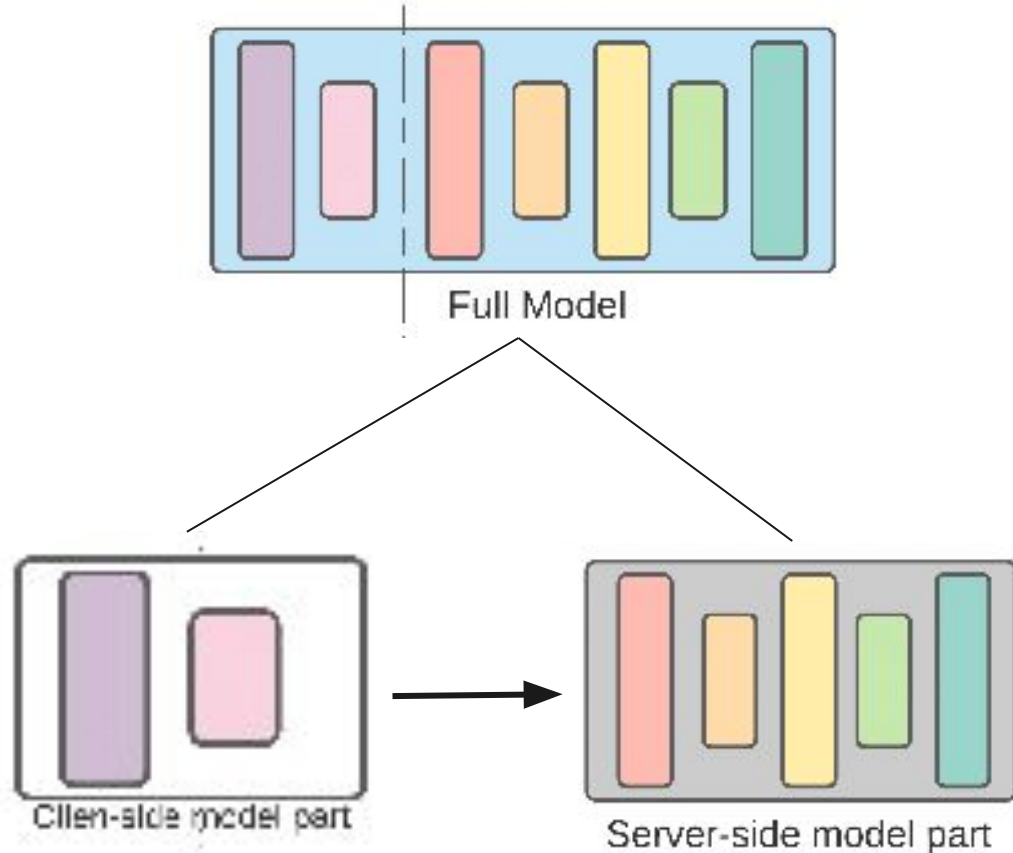- Commonly used for machine learning models
- 60,000 images

# Offloading Strategies

# Offloading Explained



Offload

# Early Exiting

- Take full model, split into a small model and larger one

- Client device gets smaller one, tries to make prediction



Full Model

Cllen-side model part

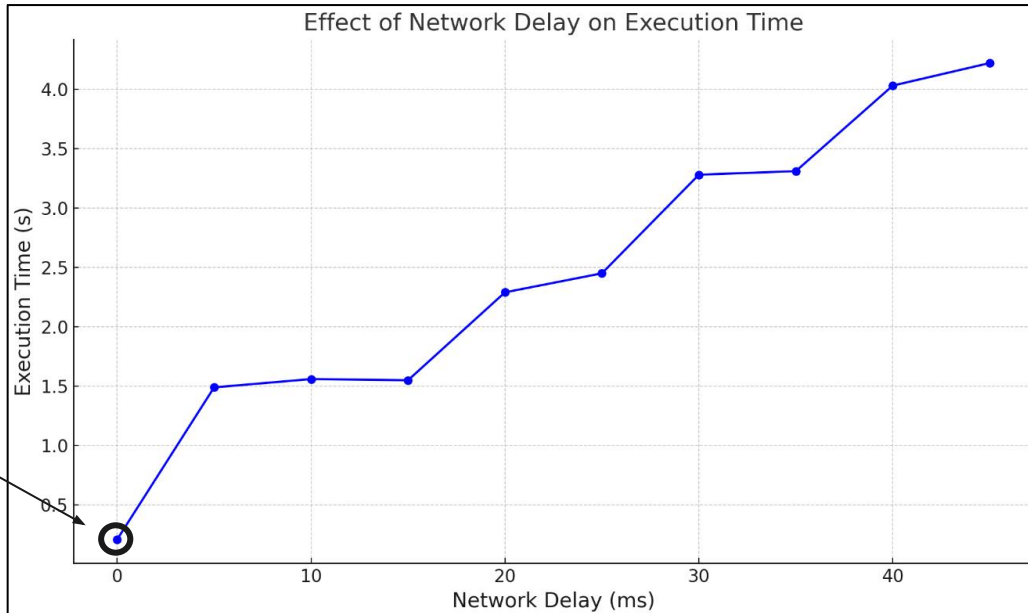Server-side model part

# Networking Conditions

# Speed

- Python's Fast API to communicate with server and device
- Sending tensors as buffer streams instead of JSON

Processing 500 images/second over an ethernet connection



LOONEY TUNES and all related characters and elements © & ™ Warner Bros. Entertainment Inc. (s01)

# Effects of Network Speed



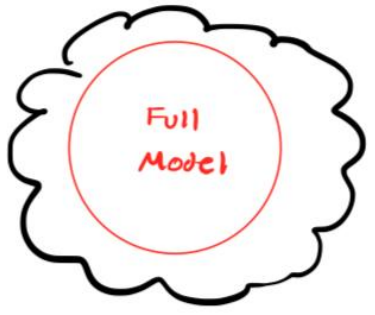Effect of Network Delay on Execution Time

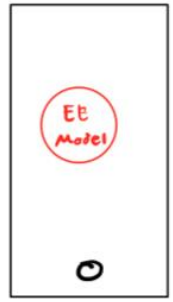Ethernet (Control)

# Evaluation

# Experiment Objective

Simulated a busy server and evaluated the different strategies client devices could use to obtain optimal results.

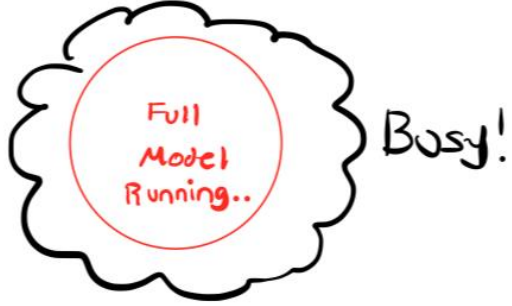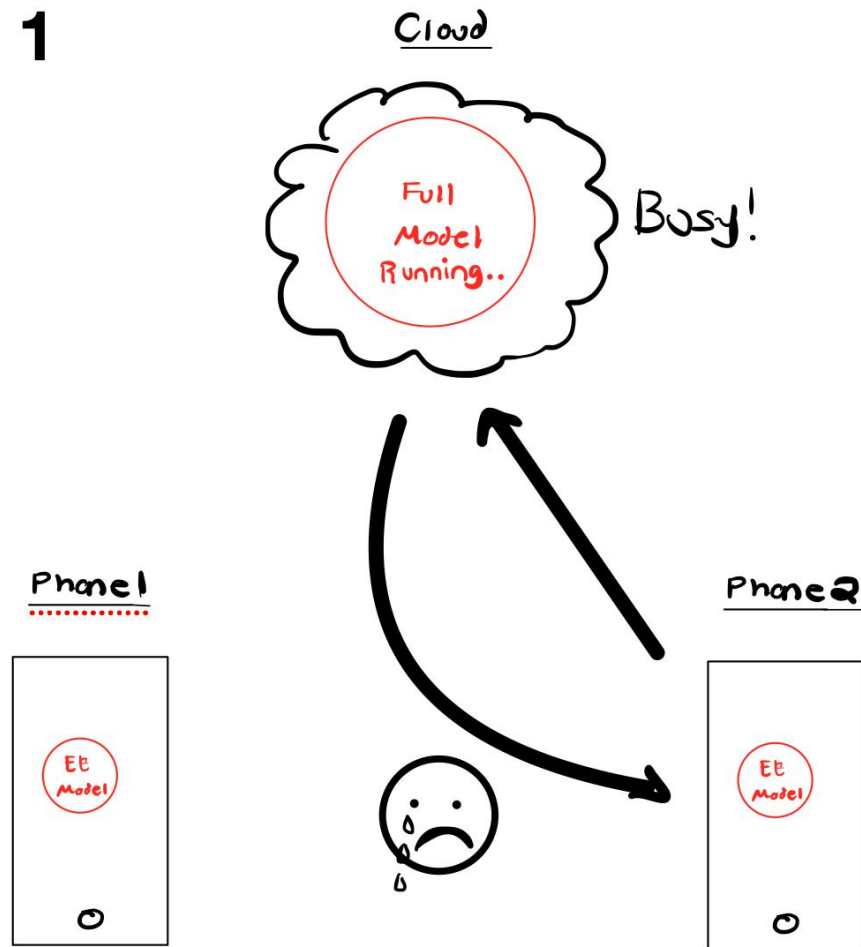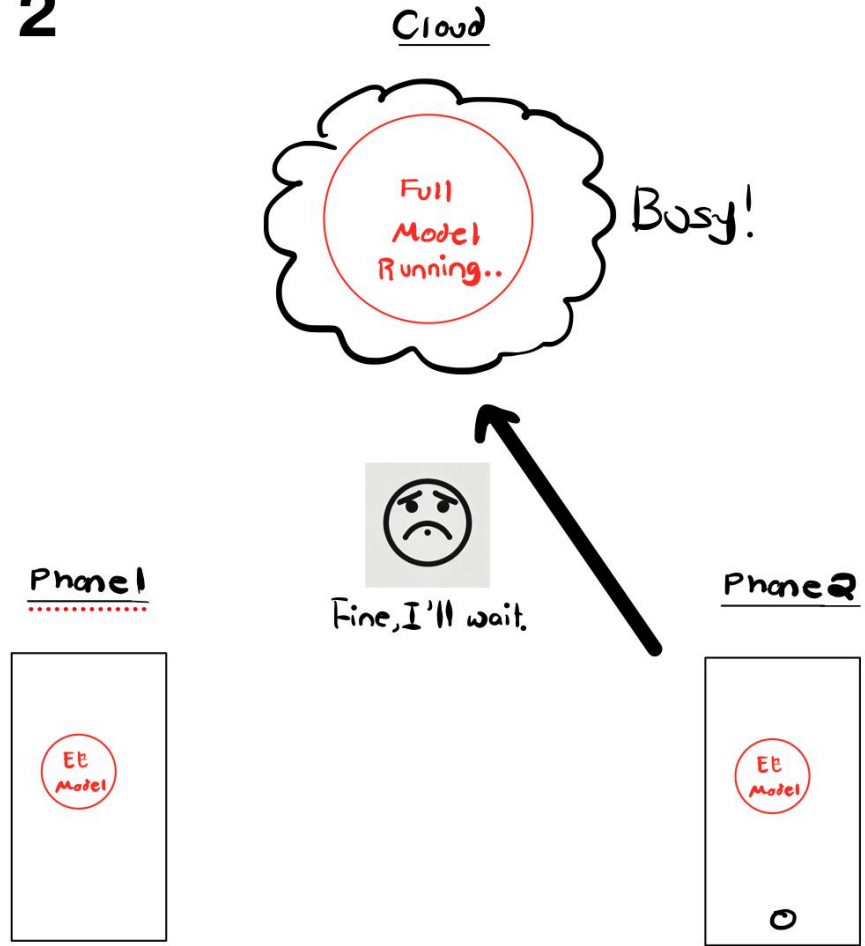# Scenario 1

# Scenario 2

Cloud

Full Model Running..

Busy!

Fine, I'll wait.

Phone 1

Et Model

Phone 2

Et Model

Comparison of Scenarios: Early Exit vs. Waiting for Server

**Accuracy Comparison**

Early Exit on Busy: 91.68%
Wait for Server: 92.43%

**Time Taken Comparison**

Early Exit on Busy: 1:25
Wait for Server: 3:58

# Early Exit Demo

Drag and drop an image
or browse to upload

**Upload Image**

Predicted class: ship, Confidence: 100.00%
Processing time: 0.27 seconds

# Conclusions

- For a small ↓ in accuracy, we can get ↑ in speed by offloading when needed
- Applications may prioritize accuracy or latency
- Challenges:
  - Managing multiple clients in real life
  - Model portability for micro devices
  - Performance in weak wireless connections

## Future Work

- Compress and optimize models

- Explore dynamic thresholds

- Train the model with diverse datasets

- Investigate other types of models

# Thank You!