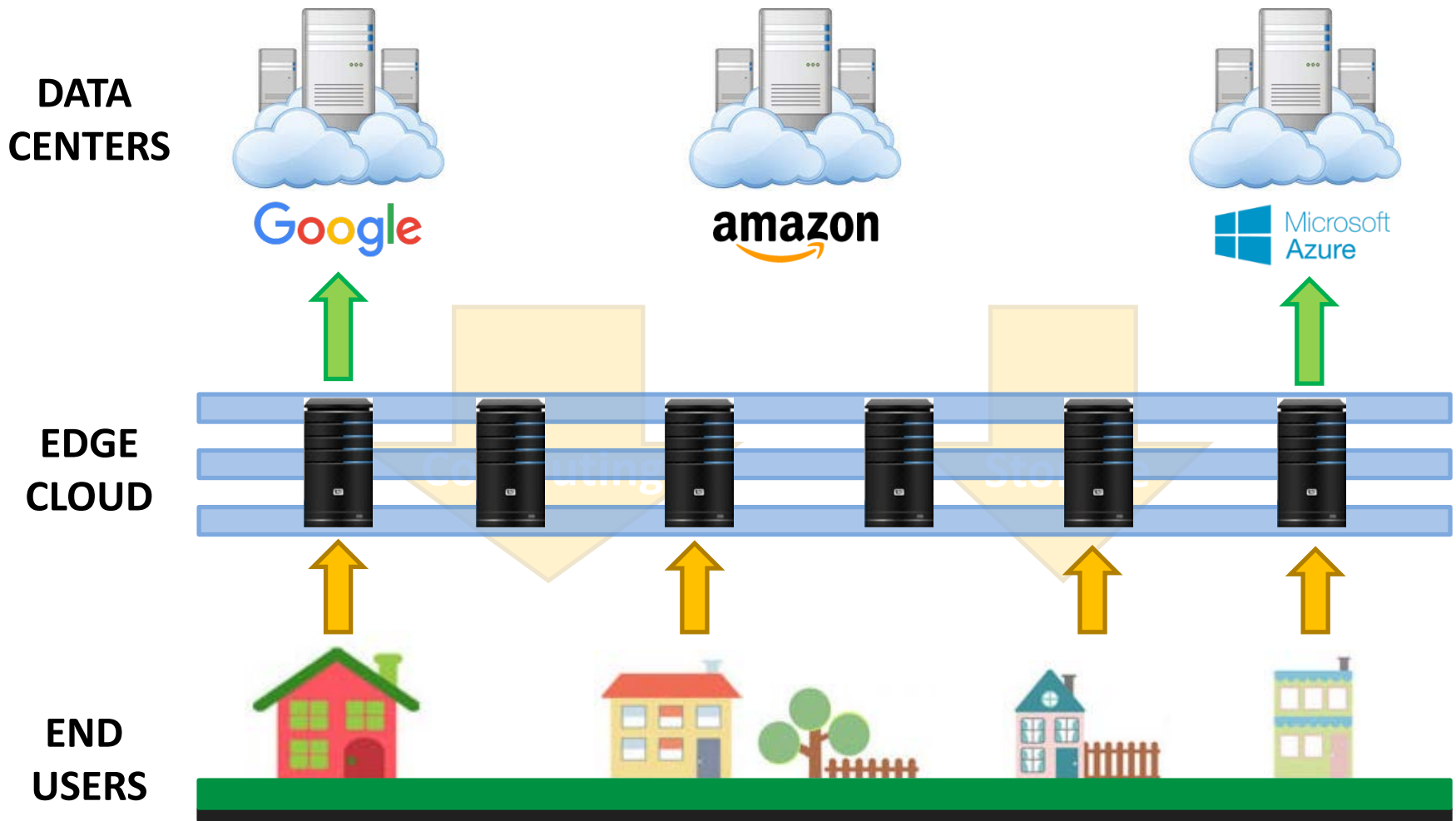# Designing Hierarchical Edge Cloud for Mobile Computing

## Wei Gao
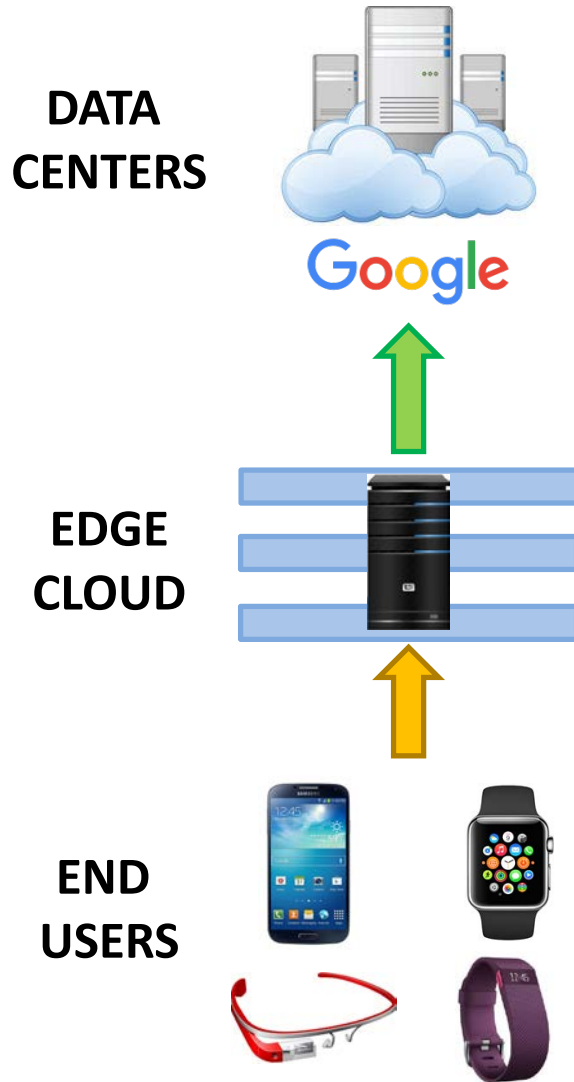
### The University of Tennessee, Knoxville

# Edge Cloud

# Edge Cloud for Mobile Computing

**DATA CENTERS**

**EDGE CLOUD**

**END USERS**

- Reduced response latency
  - Delay-sensitive mobile applications

- Higher efficiency of resource utilization
  - Distributed processing

# Applications of Edge Cloud

**Cloud**

**Internet of Things**

**NSF CSR Hightlighted Area**

**Virtual Reality**

**Smart Cities and Communities**

# Challenges

- Adaptability
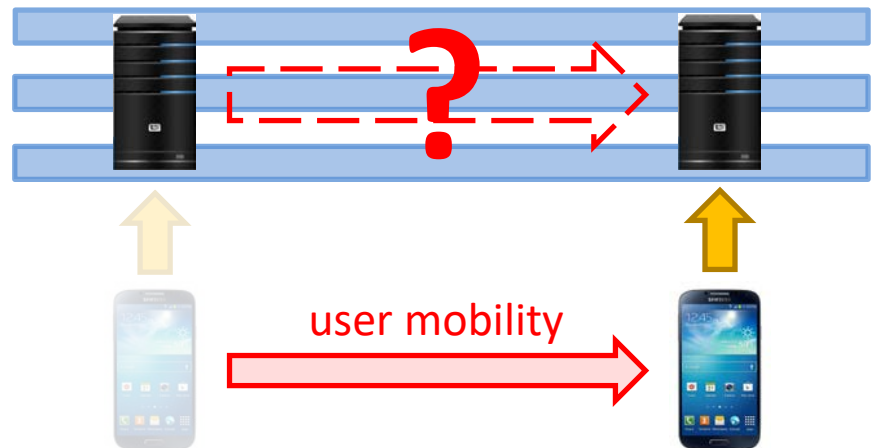  - Optimized performance?
  - Minimized cost?
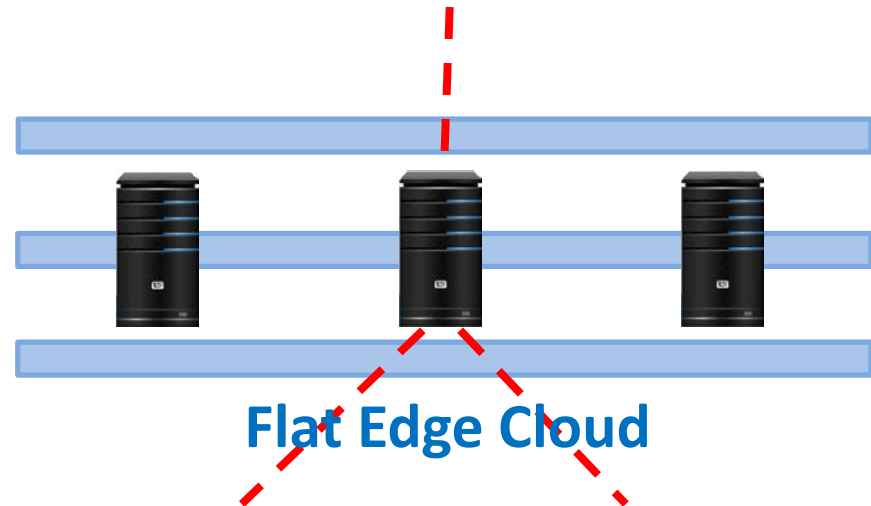
    Provisioning for the peak load

- User mobility
  - Minimized cost?

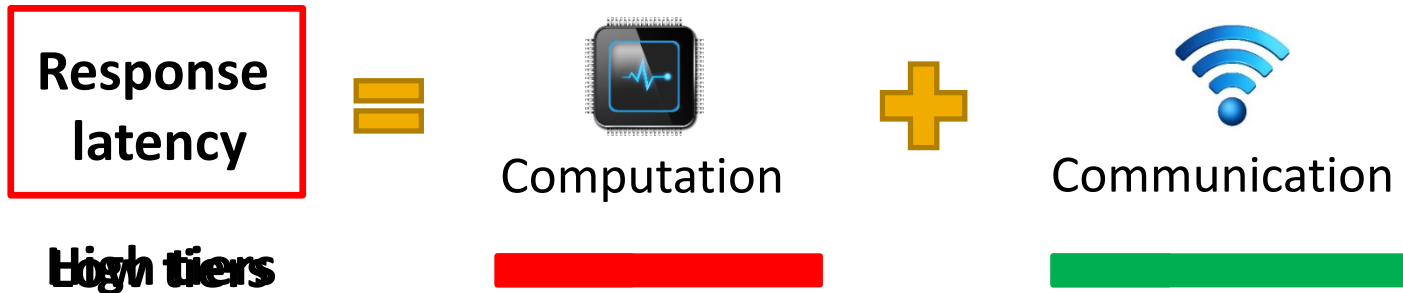    Complete move of data and program

user mobility
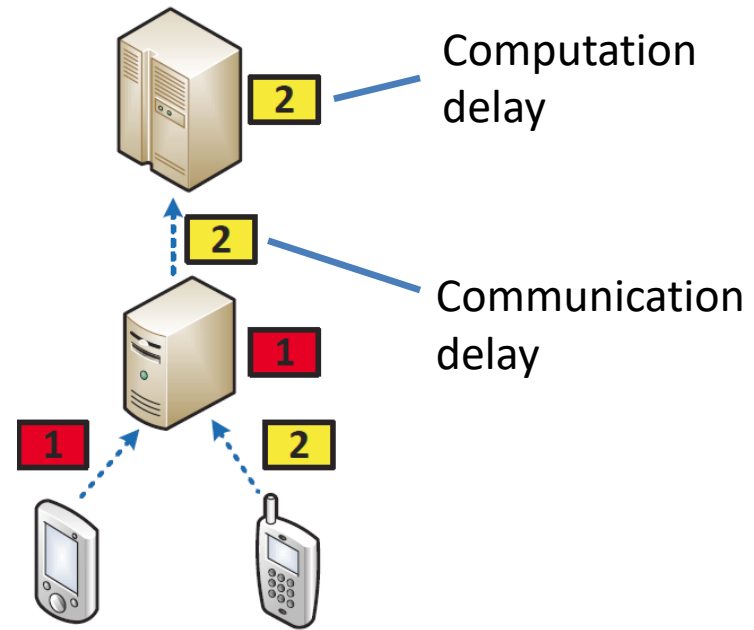
# Our Solution: Hierarchical Edge Cloud

- Adaptability
  - Aggregation of peak load

- User mobility
  - Partial migration of data and program

**Flat Edge Cloud**

3-distributed tree hierarchy

# Task 1: Optimal Workload Placement

- Our focus: minimized response latency
  - Where to place a workload
  - How much capacity for a workload
- Challenge
  - Delay tradeoff



Computation delay

Communication delay

**Response latency** = Computation + Communication

High tiers
Low tiers

A Hierarchical Edge Cloud Architecture for Mobile Computing, *in IEEE INFOCOM'16*.

# Task 1: Optimal Workload Placement

- Distributed optimization

Computation delay   Communication delay

$$\mathbf{min}\, f = \sum_{i=1}^{m} \left( \frac{w_i}{\lambda_{i,\gamma_i} c_{\gamma_i}} + (L(\gamma_i) - 1) \frac{s_i}{B_{\gamma_i}} \right),$$

$$\mathrm{s.\,t.} \sum_{j \in O_j} \lambda_{i,j} = 1, j = 1,2,\dots,n$$

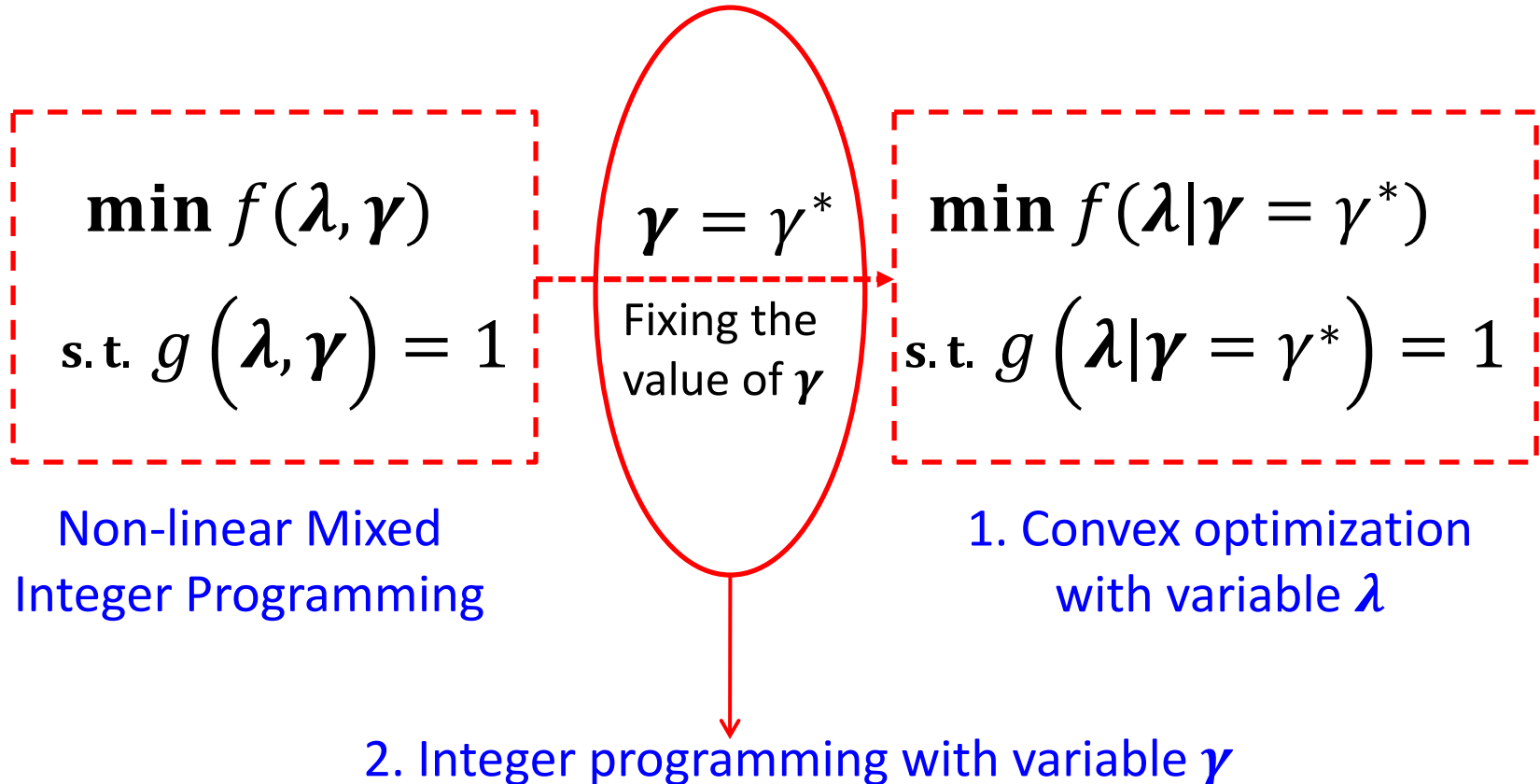Capacity allocation of server *j* to workload *i*

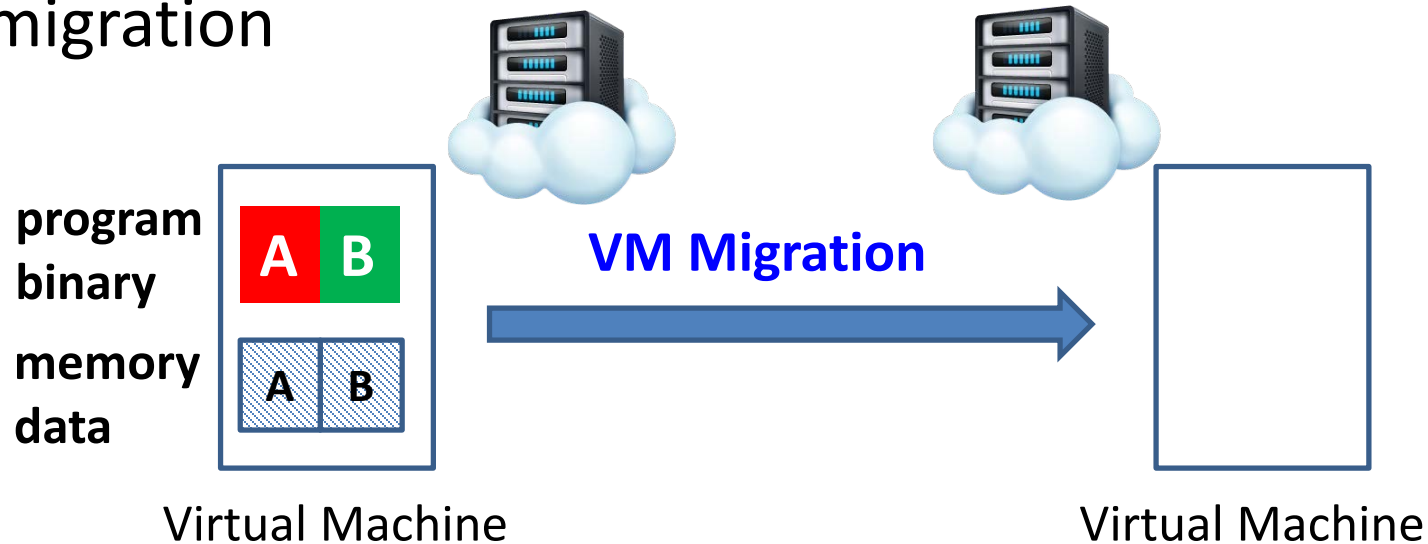Placement of workload *i*

Non-linear mixed integer programming

**?**

A Hierarchical Edge Cloud Architecture for Mobile Computing, *in IEEE INFOCOM'16*.

8

# Task 1: Optimal Workload Placement

- Problem transformation

$$\min f(\lambda, \gamma)$$

$$\text{s.t.} \; g(\lambda, \gamma) = 1$$

$$\gamma = \gamma^*$$

Fixing the value of $\gamma$

$$\min f(\lambda | \gamma = \gamma^*)$$

$$\text{s.t.} \; g(\lambda | \gamma = \gamma^*) = 1$$

Non-linear Mixed Integer Programming

1. Convex optimization with variable $\lambda$

2. Integer programming with variable $\gamma$

A Hierarchical Edge Cloud Architecture for Mobile Computing, *in IEEE INFOCOM'16*.

# Task 2: Supporting User Mobility

- Remote program execution with least context migration



program binary

memory data

**VM Migration**

Virtual Machine

Virtual Machine

Offline program parsing ➜ Run-time migration

Minimizing Context Migration in Mobile Code Offload, *in IEEE Transactions on Mobile Computing, 2017.*

10

# Task 2: Supporting User Mobility



Minimizing Context Migration in Mobile Code Offload, *in IEEE Transactions on Mobile Computing, 2017.*

# Implementation

- Heterogeneous mobile and wearable platforms

Samsung Galaxy S4          LG Watch Urbane          Samsung Nexus 10 Tablet

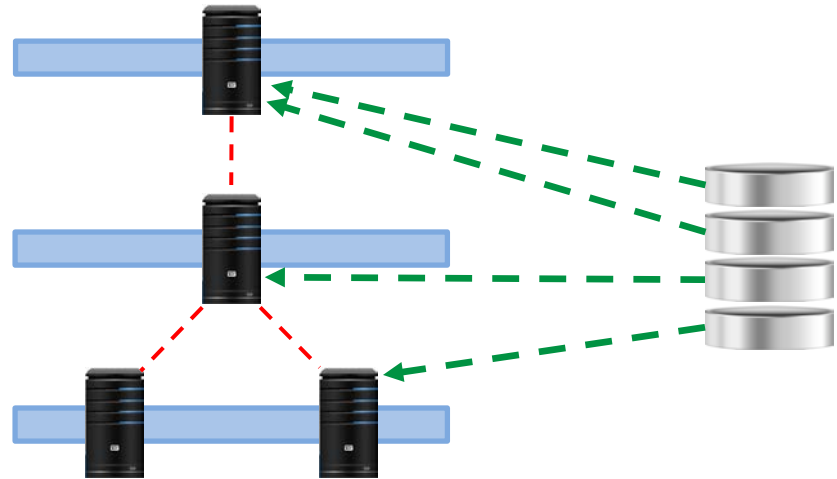  - 1,500 LoC in Java and 1,000 LoC in C++ over Android v5.1.1 OS kernel

- Edge cloud servers
  - x86-based instances of Dalvik VM
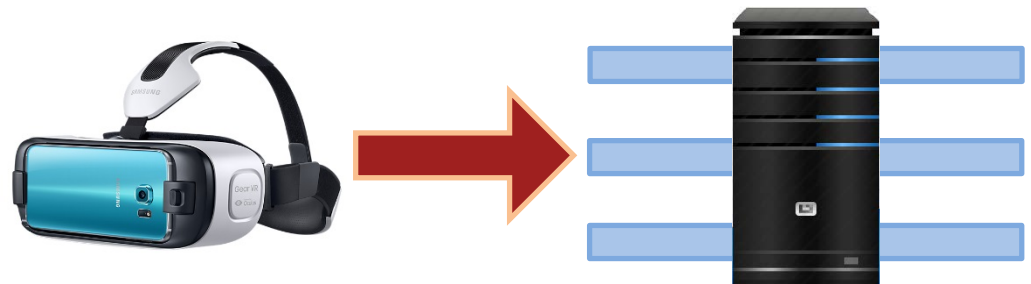  - Executing ARM-based native programs

# Ongoing Work and Future Directions

- Optimal resource provisioning
  - Minimizing both CapEx and OpEx



- Virtual reality over mobile platforms

# **Thank you**

- Questions?