

---

# Data-Driven Edge Resource Provisioning for Inter-Dependent Microservices with Dynamic Load

*Ruozhou Yu, North Carolina State University*

*Szu-Yu Lo, Fangtong Zhou, North Carolina State University*

*Guoliang Xue, Arizona State University*

# Outlines

---

**Background and Motivation**

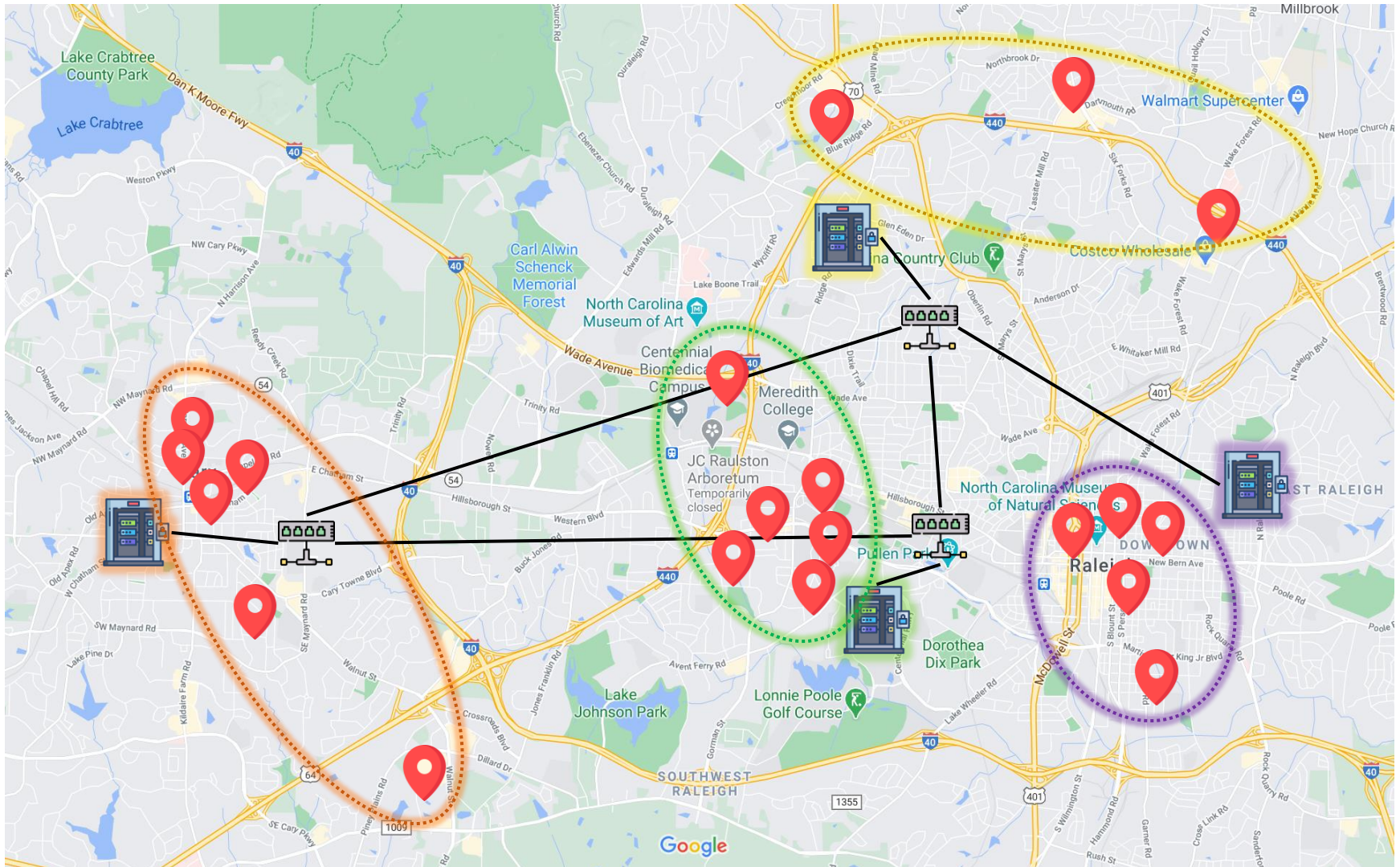
**System Modeling**

**Solution Design**

**Performance Evaluation**

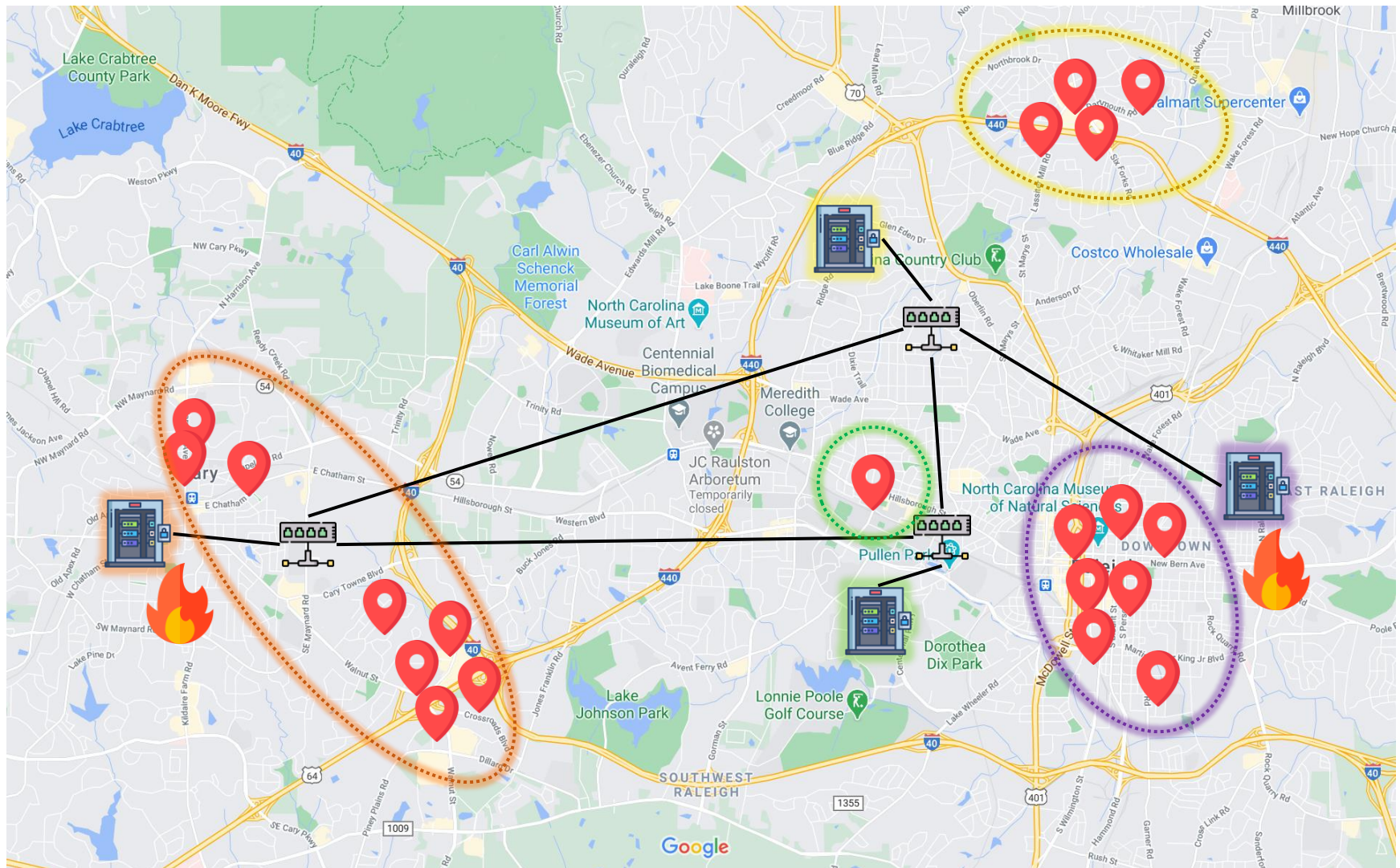
**Discussions, Future Work and Conclusions**

# Geo-Distributed Services & Edge Computing

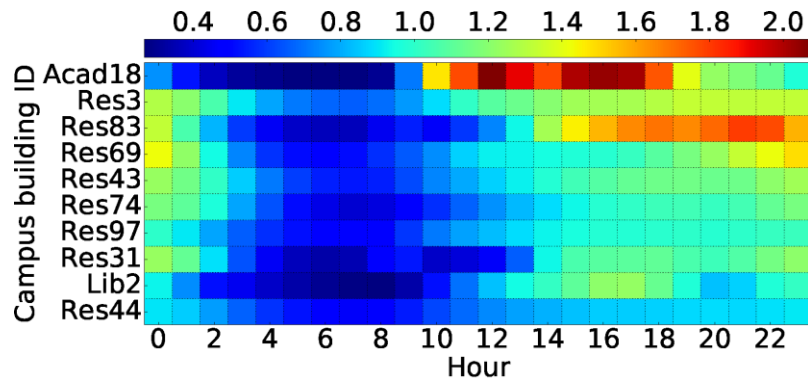




# Time-Varying Demands in Geo-Distributed Apps

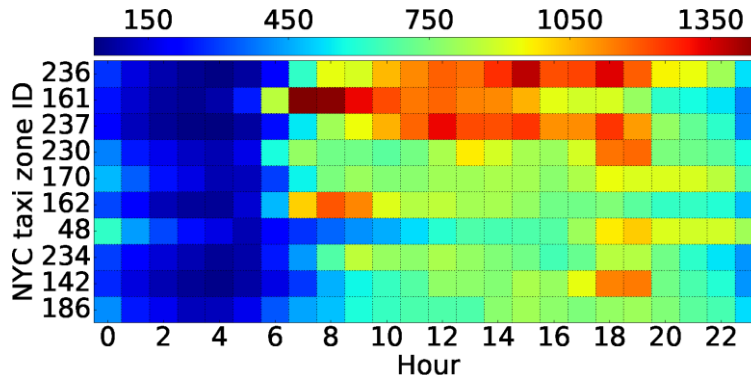


# Patterns in Real-world Datasets [1]



## Dartmouth College Wireless APs

- Top 10 APs with highest loads
- Load: avg. # devices / hour
- Averaged over a year (9/2002-9/2003)



## NYC Yellow Taxi 2018

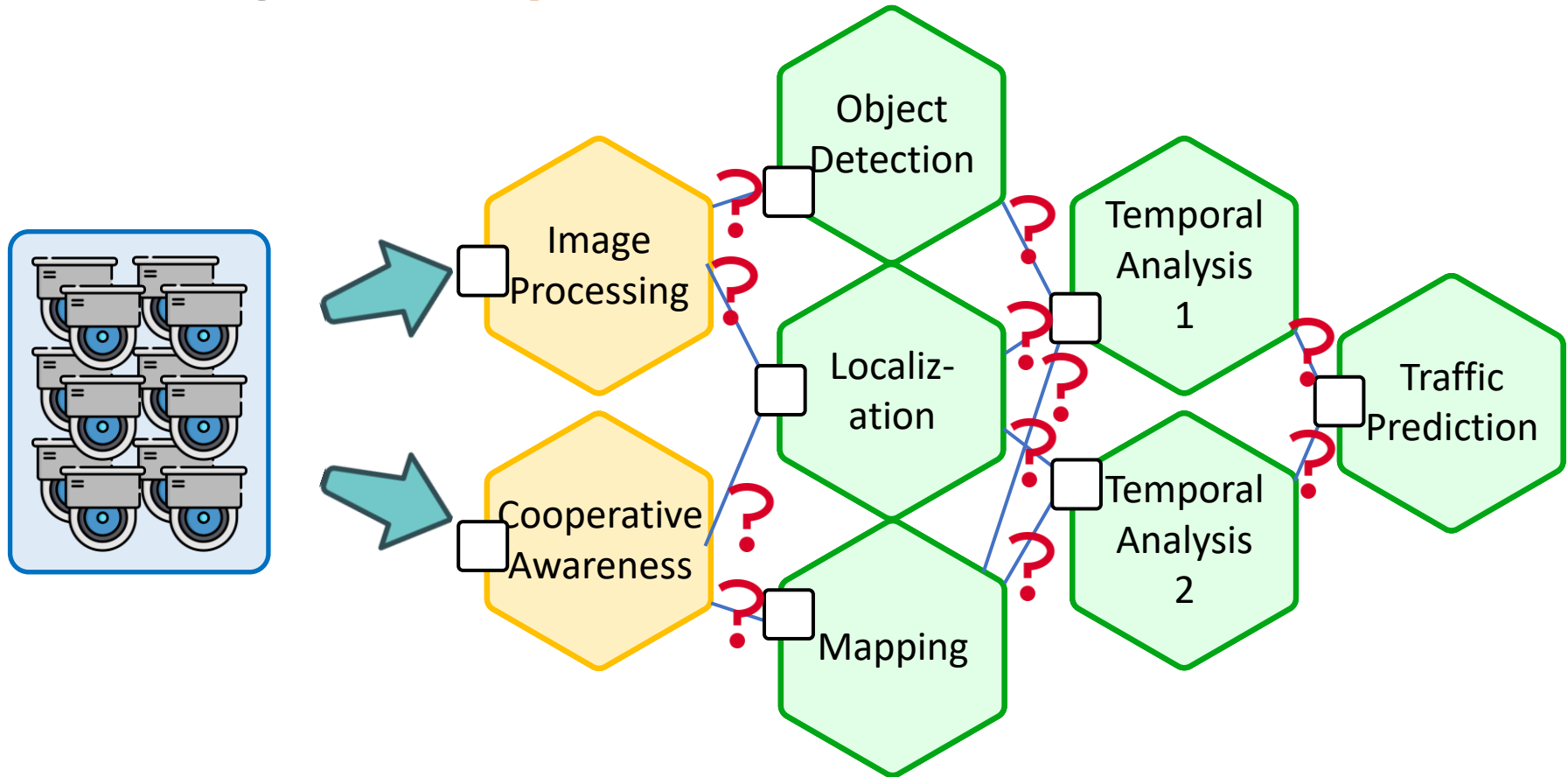
- Top 10 zones w/ most drop-offs
- Load: avg. # passenger drop-offs
- Averaged over a year

**Observation 1: Non-i.i.d. demand distributions across time & locations.**

**Observation 2: Repeating / seasonal patterns in temporal domain.**

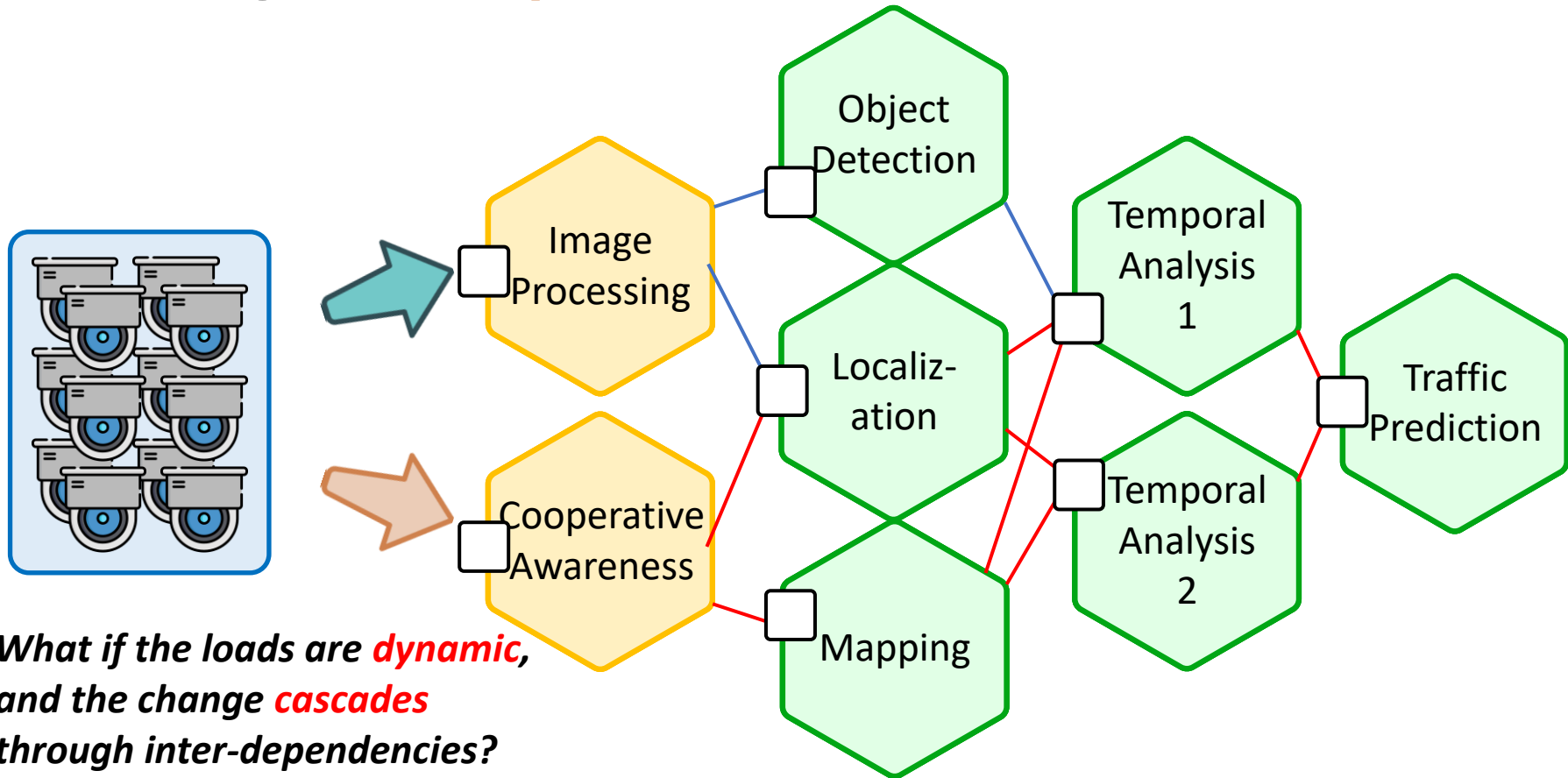
# The Microservice Load Balancing Problem [2]

- ❑ Edge-based microservices can be easily **saturated**.
- ❑ **Challenge:** **interdependent** microservices.

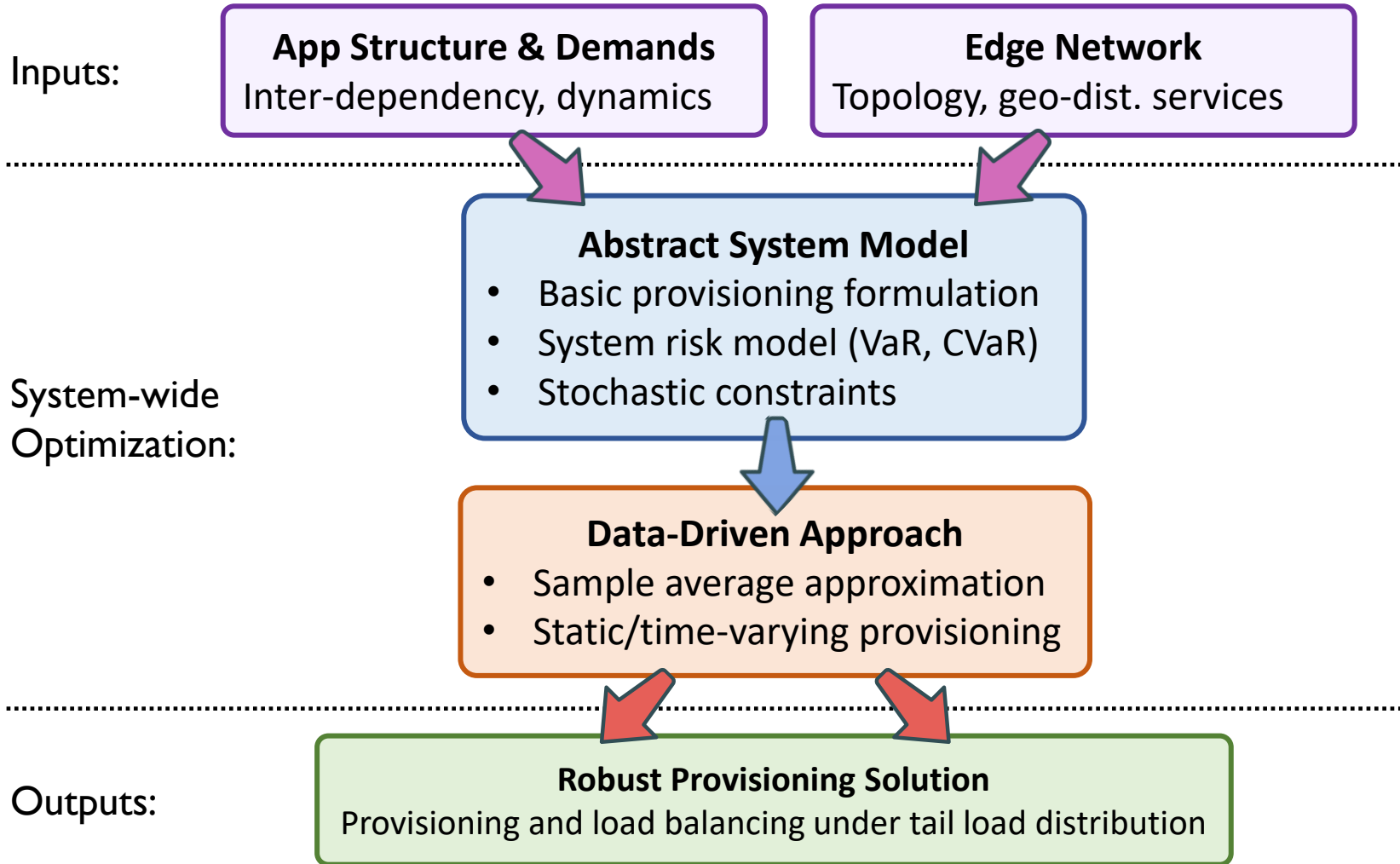


# The Microservice Load Balancing Problem [2]

- ❑ Edge-based microservices can be easily **saturated**.
- ❑ **Challenge:** **interdependent** microservices.



# Methodology Overview





# Outlines

---

**Background and Motivation**

**System Modeling**

**Solution Design**

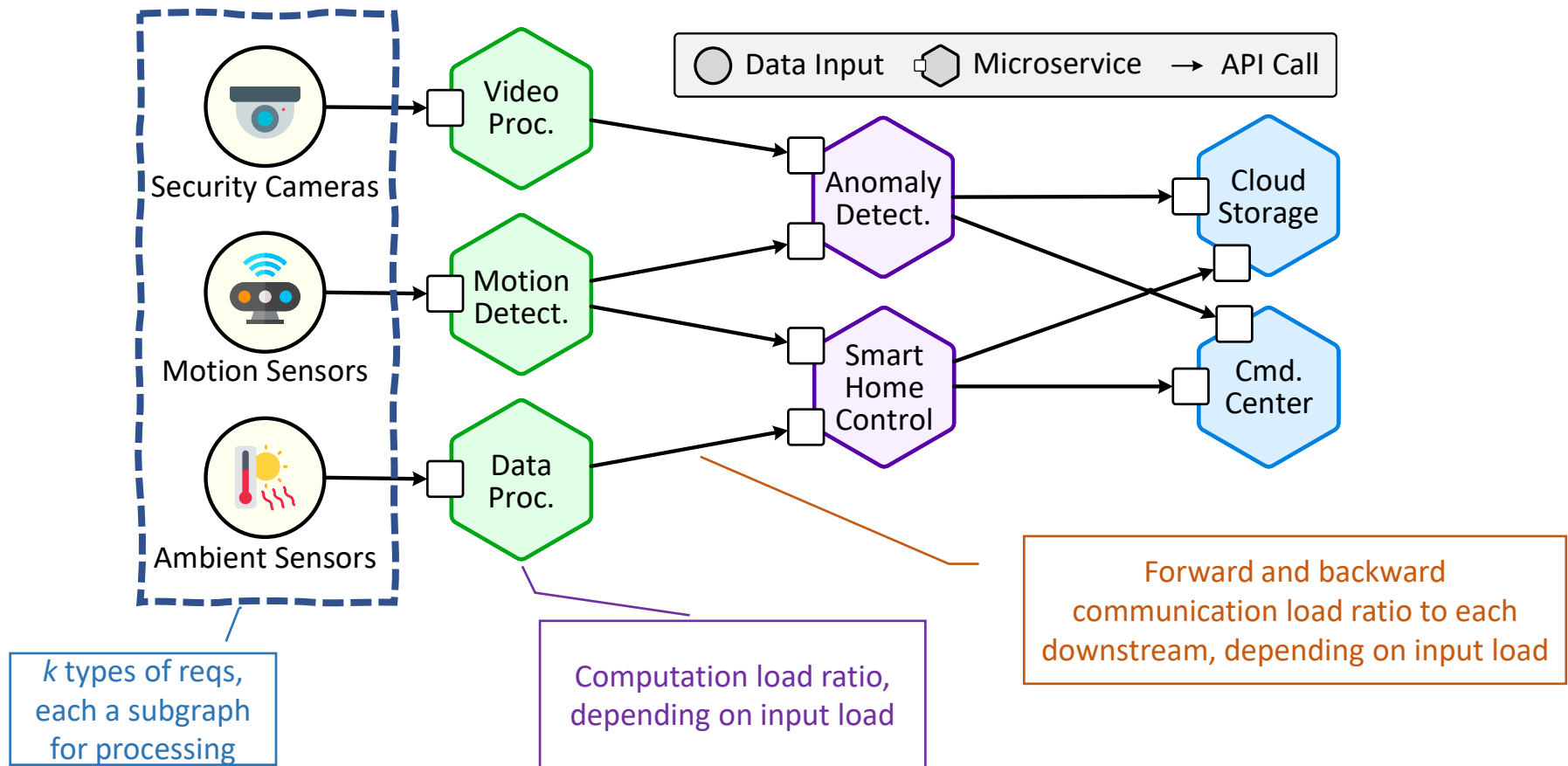
**Performance Evaluation**

**Discussions, Future Work and Conclusions**

# Application with Interdependent Microservices

## □ General DAG-based application graph (App-Graph).

❖ *Captures complex interdependencies, unlike existing line graph-based models.*



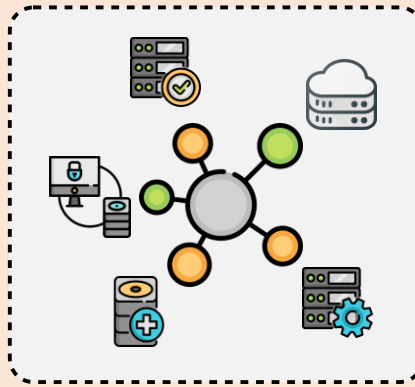
# Edge Network: A General Model

❑ **Challenge:** heterogeneous network environments



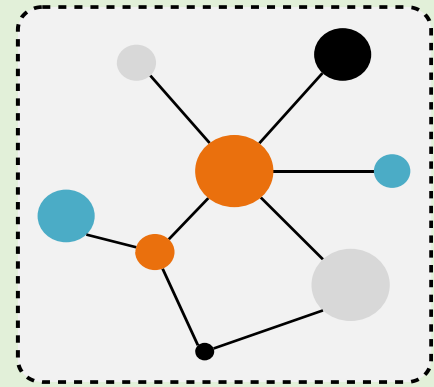
Wireless RANs:

- Geo-distributed
- Limited capacity
- Interference



Edge Network:

- Complex topo
- Distributed
- Dynamic load



Backbones:

- Large-scale
- High latency
- ISP policies

❑ **Model:** general directed graph  $G=(N, L)$ , with edge nodes  $H$  and APs  $A$

- ❖ Weights: link bandwidth, edge node capacity, deployed microservices instances

# Objective and Overall Formulation

- **Goal:** routing & load allocation to minimize max node/link load

$$\min_{f, \psi} \underline{\psi}$$

*Obj: minimize global load factor* (1)

$$\sum_{k \in [K]} \sum_{a \in A} \sum_{e \in E_k} (\underline{\rho_{e,1}^k} \cdot f_a^k(e, l) + \underline{\rho_{e,2}^k} \cdot f_a^k(e, l')) \leq b_l \cdot \underline{\psi}, \quad \forall l \in L;$$

*Link bw & load* (2)

$$\sum_{k \in [K]} \sum_{a \in A} \sum_{v \in V_k \cap V_h} \underline{\rho_v^k} \cdot f_a^k(v, h) \leq c_h \cdot \underline{\psi}, \quad \forall h \in H.$$

*Node cap & load* (3)

$$\begin{aligned} \sum_{l \in L_{\text{out}}(n)} f_a^k(e, l) - \sum_{l \in L_{\text{in}}(n)} f_a^k(e, l) = \\ \mathbf{1}_{\{n \in H_u\}} \cdot f_a^k(u, n) - \mathbf{1}_{\{n \in H_v\}} \cdot f_a^k(v, n), \\ \forall k \in [K], a \in A, n \in N, e = (u, v) \in E_k, \end{aligned}$$

*Inter-dependencies* (4)

$$f_a^k(v_k, n) \geq \delta_a^k, \quad \forall k \in [K], n = a \in A.$$

*Geo-distributed demands of each type* (5)

- But  $\{\delta_a^k\}$  are random and time-varying...

# Outlines

---

**Background and Motivation**

**System Modeling**

**Solution Design**

**Performance Evaluation**

**Discussions, Future Work and Conclusions**



# SO and CVaR

□ **Stochastic Optimization (SO):** optimize a function in presence of randomness (random objective and/or constraints)

❖ Traditional approach: expectation optimization / constraints

$$\min_{\mathcal{X} \in \mathcal{F}} \mathbb{E}[R] \quad \text{or} \quad A\mathcal{X} \geq \mathbb{E}[R]$$

❖ **Issue:** unbounded risk in rare but unfortunate scenarios

➤ E.g., abnormal demands due to public events, rare large-scale failures, ...

❖ How to model these *unfortunate scenarios*?

❖ **Value-at-Risk (VaR)** and **Conditional-Value-at-Risk (CVaR):**

➤ Widely used in economics and finance

➤  $\text{VaR}_\alpha(R) = \min \{ c \in \mathbb{R} \mid R \text{ does not exceed } c \text{ with at least } \alpha \text{ prob.} \}$

➤  $\text{CVaR}_\alpha(R) = \mathbb{E}[R \mid R \geq \text{VaR}_\alpha(R)]$

□ Expectation of  $R$  in the worst  $(1-\alpha)$  scenarios

❖ **Our approach:** optimize with CVaR constraints

$$f_a^k(v_k, n) \geq \text{CVaR}_\alpha(\delta_a^k), \quad \forall k \in [K], n = a \in A.$$

# Transformation and Data-Driven Approach

---

❑ **Challenge 1:** CVaR not written in closed-form

❑ **Technique:** LP transformation by Rockafella & Uryasev

$$\text{CVaR}_\alpha(\mathbf{R}) = \min_r \left\{ r + \frac{1}{1-\alpha} \mathbb{E}[(\mathbf{R} - r)^+] \right\}, \quad (9)$$

❖ A convex optimization problem given  $\alpha$ .

❑ **Challenge 2:** unknown distributions to random variables

❑ **Technique:** Sample Average Approximation (SAA)

❖ *I.i.d. Samples:* observed demand data in historical periods.

❖  $\{\delta_a^k\}$  expanded to  $\{\delta_a^{k,i}\}$  for  $i = 1 \dots N$  samples.

$$\text{CVaR}_\alpha(\delta_a^k) \approx \min_r \left\{ r + \frac{1}{1-\alpha} \frac{1}{N} \sum_{i=1}^N (\tilde{\delta}_a^{k,i} - r)^+ \right\}. \quad (10)$$

# Outlines

---

**Background and Motivation**

**System Modeling**

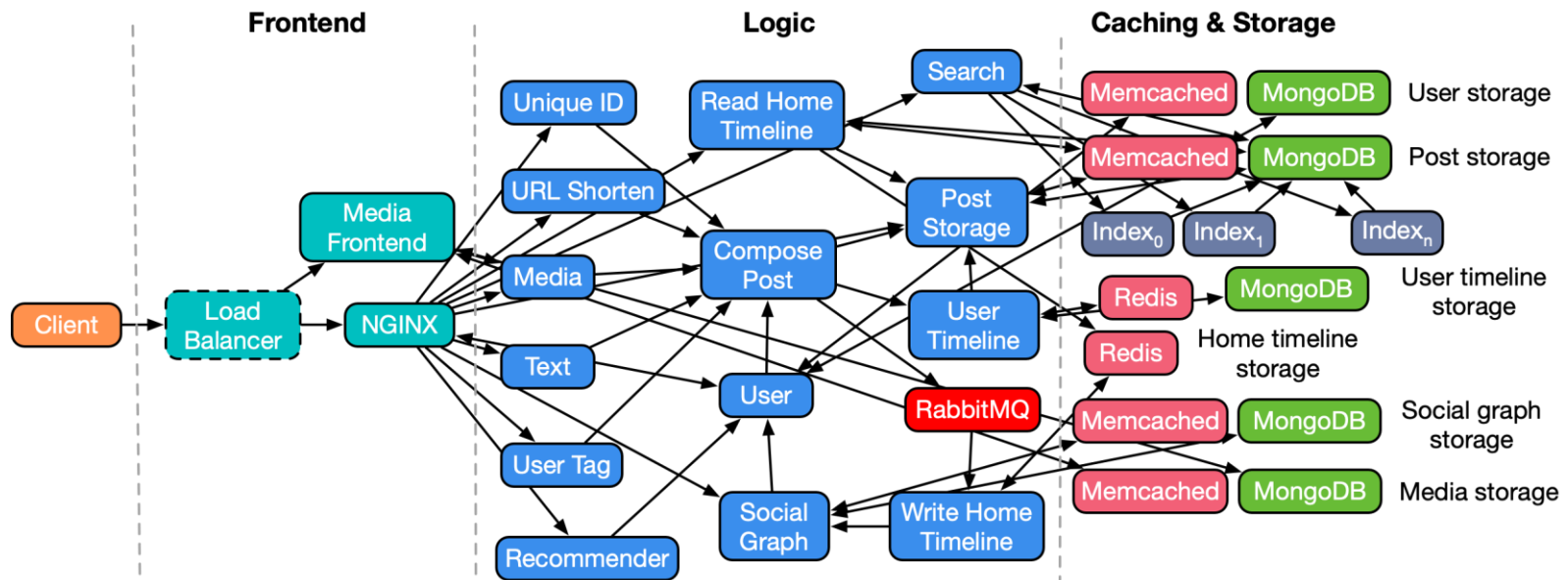
**Solution Design**

**Performance Evaluation**

**Discussions, Future Work and Conclusions**

# Simulation Setting - Application

## ❑ Social Network (SN) from DeathStarBench [3]



- ❖ 23 microservices, 3 types of workloads (compose-post, read-home-timeline, read-user-timeline) profiled.
- ❖ Implemented and profiled for actual communication load (# bytes); computation load/capacity synthesized based on communications.

# Simulation Settings – Demand & Network

---

## □ Settings

### ❖ Dataset: NYC Yellow Taxi 2018

- 12 months of Taxi pick-up/drop-off data (~112 million taxi trips)
- Picked 20 most popular zones out of 262 (55% of all demands)
- Mapped demands (drivers/passengers, pick-up/drop-off) to SN requests
- 20% training (optimization) & 80% testing (deployment)

### ❖ Synthetic Data

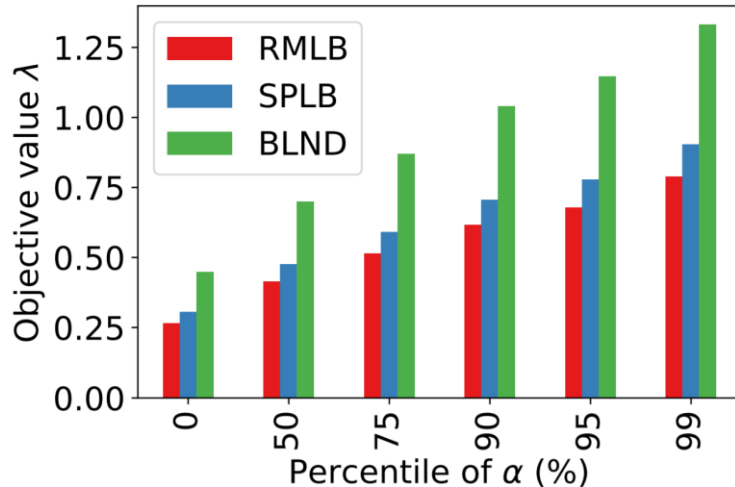
- Random topologies: Watts-Strogatz with  $k = 4$  and  $p = 0.3$
- Each microservice deployed on 20% random edge nodes
- Network conditions: 1Gbps links, 2.5Gbps computation capacity (normalized)

### ❖ $\alpha = 0.95$ (CVaR confidence)

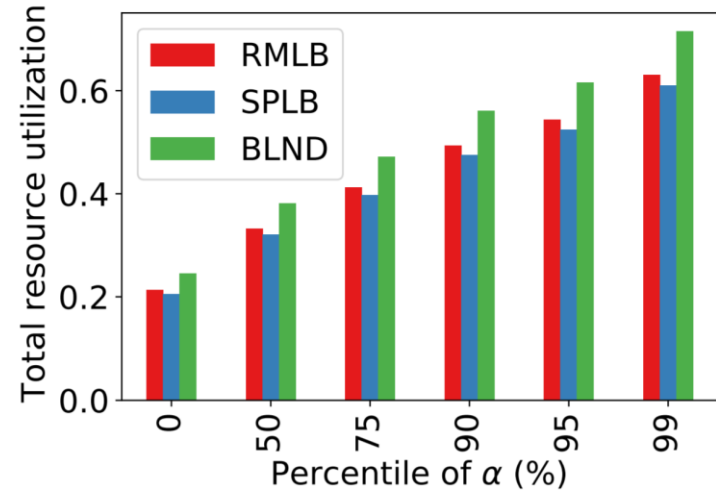
### ❖ Comparison: Shortest Path-based Heuristic, and Blind Load Balancing Heuristic



# Selected Experiment Results



(a) Max load factor  $\lambda$  in training

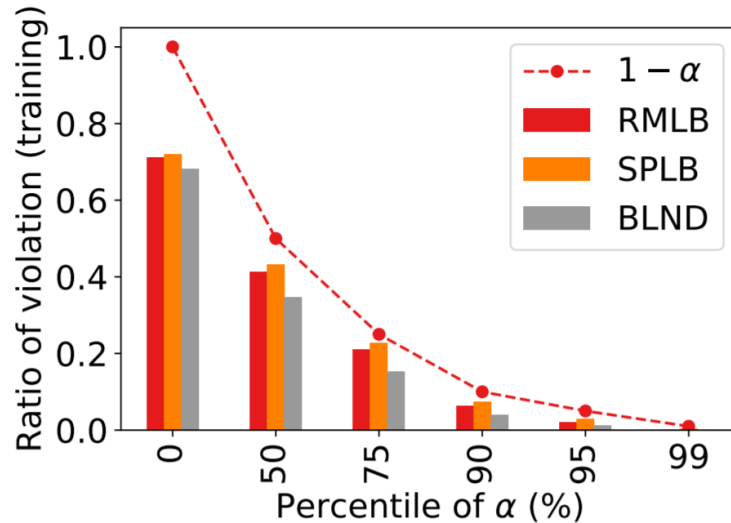


(b) Total edge resource consumption

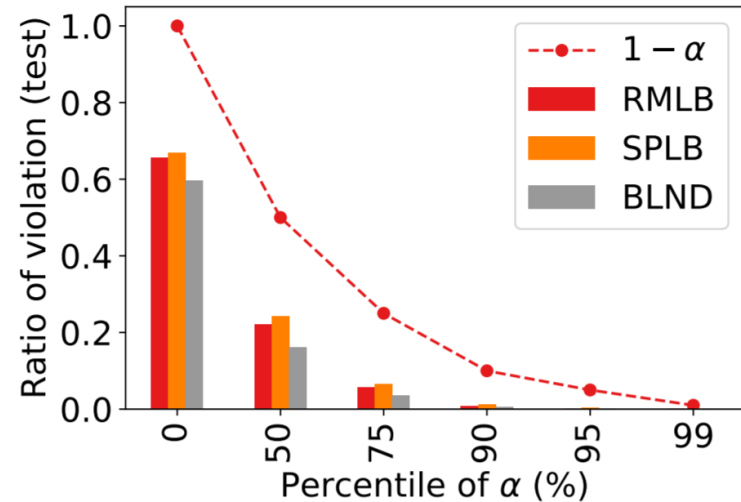
## Maximum Load (Resource Provisioning)

- RMLB (our formulation) achieves best inter-dependency-aware provisioning. Other algorithms result in higher maximum load, and lower/higher total load.

# Selected Experiment Results



(c) Ratio of violation (training)

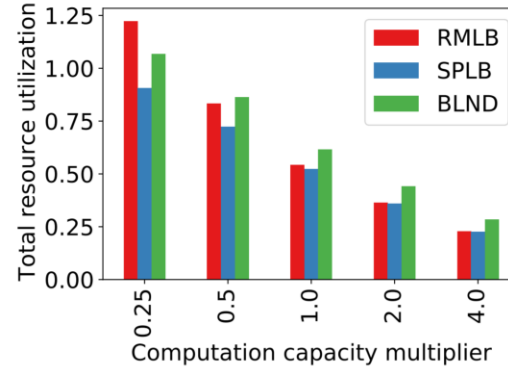
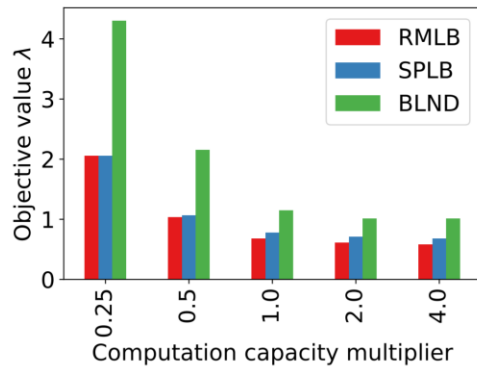


(d) Ratio of violation (testing)

## Robustness (Actual Demands)

- We provisioned for the training set (left), with bounded ratio of load violation.
- In the test deployment, we observe similar (lower) ratio of load violation.
- The  $(1 - \alpha)$  percentile is never violated, depending on our setting of  $\alpha$ .

# More Results in Paper

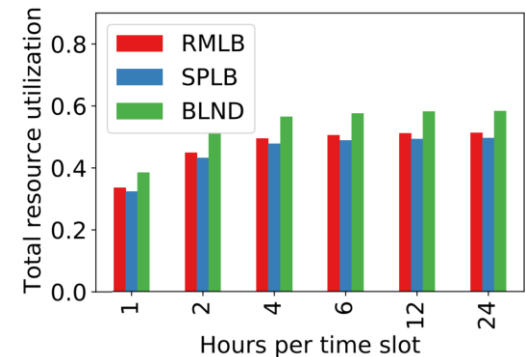
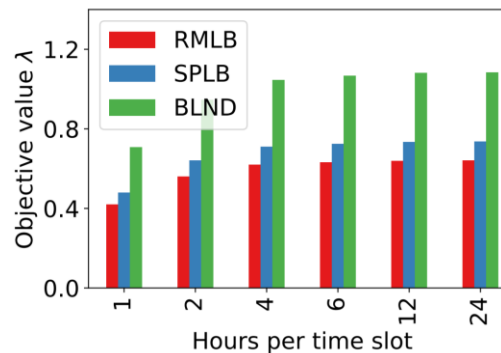


**With computation vs.  
network bottlenecks**

(a) Max load factor  $\lambda$  in training

(b) Total edge resource consumption

**With/without time-  
varying provisioning**



(a) Max load factor  $\lambda$  in training

(b) Total edge resource consumption

# Outlines

---

**Background and Motivation**

**System Modeling**

**Solution Design**

**Performance Evaluation**

**Discussions, Future Work and Conclusions**

# Other Perspectives, Conclusions

---

## ❑ So far, we've talked about

- ❖ Microservice inter-dependencies  
+  
dynamic demands

} First-attempt modeling & solving

## ❑ What could be improved

- ❖ QoS constraints: latency, throughput, reliability
- ❖ Mixture spatial-temporal distributions
- ❖ Distribution-aware formulations
- ❖ Queueing-based risk analysis
- ❖ Improved optimization methods
- ❖ Improved statistical & learning-based methods

} Modeling Perspective

} Stochastic Perspective

} Algorithmic Perspective

## ❑ **Conclusions:** app-aware, robust computing & networking.



---

# **Thank you very much!**

Q&A?