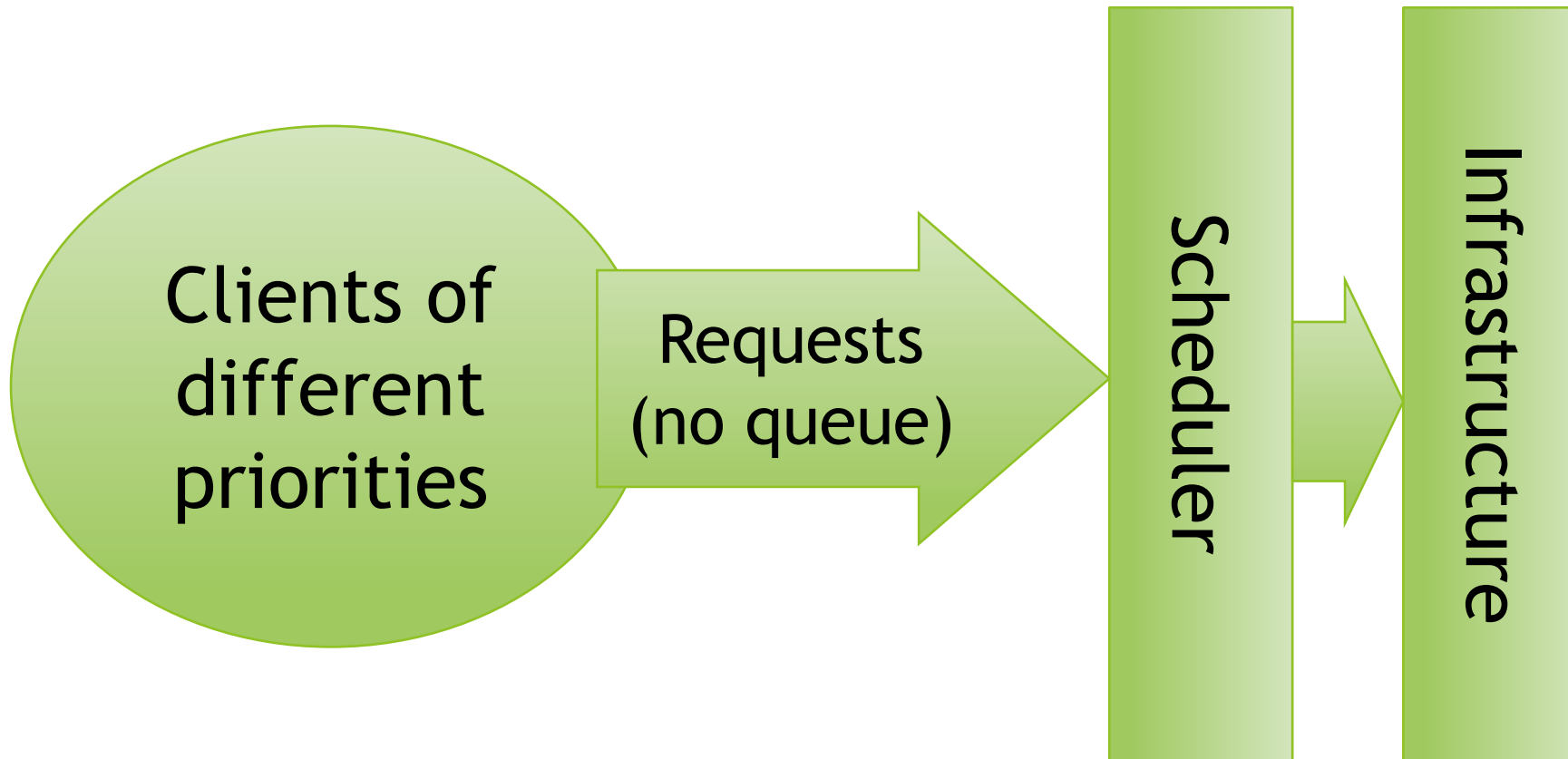


# A Theory of Auto-Scaling for Resource Reservation in Cloud Services

Konstantinos Psychas, Javad Ghaderi  
Columbia University

Performance 2020: 38th International  
Symposium on Computer Performance,  
Modeling, Measurements and Evaluation  
November 2-6 2020, Milan(Online)

# Problem Setting



- ▶ We will use terms **job** and **server**
- ▶ Both treated as abstract **multidimensional** resource vectors

# Model

*If* and *where* to schedule



A set of  $J$   
job types  $\mathcal{J}$



Type  $j \in \mathcal{J}$  gives  
reward  $u_j$   
per time unit



...



A set of  $L$   
*homogeneous*  
servers

Normalized Workload  
(jobs per server)

$$\rho := (\rho_1, \rho_2, \dots, \rho_J)$$

#jobs from each type that can  
simultaneously fit in the server

$$\mathbf{k} = (k_1, k_2, \dots, k_J) \in \mathcal{K} \subset \mathbb{R}^J$$

# How to schedule when workload is known

$$\max_{\mathbf{X}, \mathbf{Y}} \sum_j u_j Y_j$$

$$\text{s.t. } Y_j \leq \hat{Y}_j^L, \forall j \in \mathcal{J}$$

$$\sum_{\mathbf{k} \in \mathcal{K}} X_{\mathbf{k}} k_j \geq Y_j, \forall j \in \mathcal{J}$$

$$\sum_{\mathbf{k} \in \mathcal{K}} X_{\mathbf{k}} = L, \quad X_{\mathbf{k}} \geq 0, \forall \mathbf{k} \in \mathcal{K}$$

Example  $\hat{Y} = L\rho$

- ▶ Even if relaxed still hard (exponential number of variables)
- ▶ Needs to be solved every time workload changes
- ▶ Unclear how to change assignment when servers are already full
- ▶ There can be consistent error in  $\rho$  causing loss

# Solving Static Optimization: Ordering Configuration by total Reward

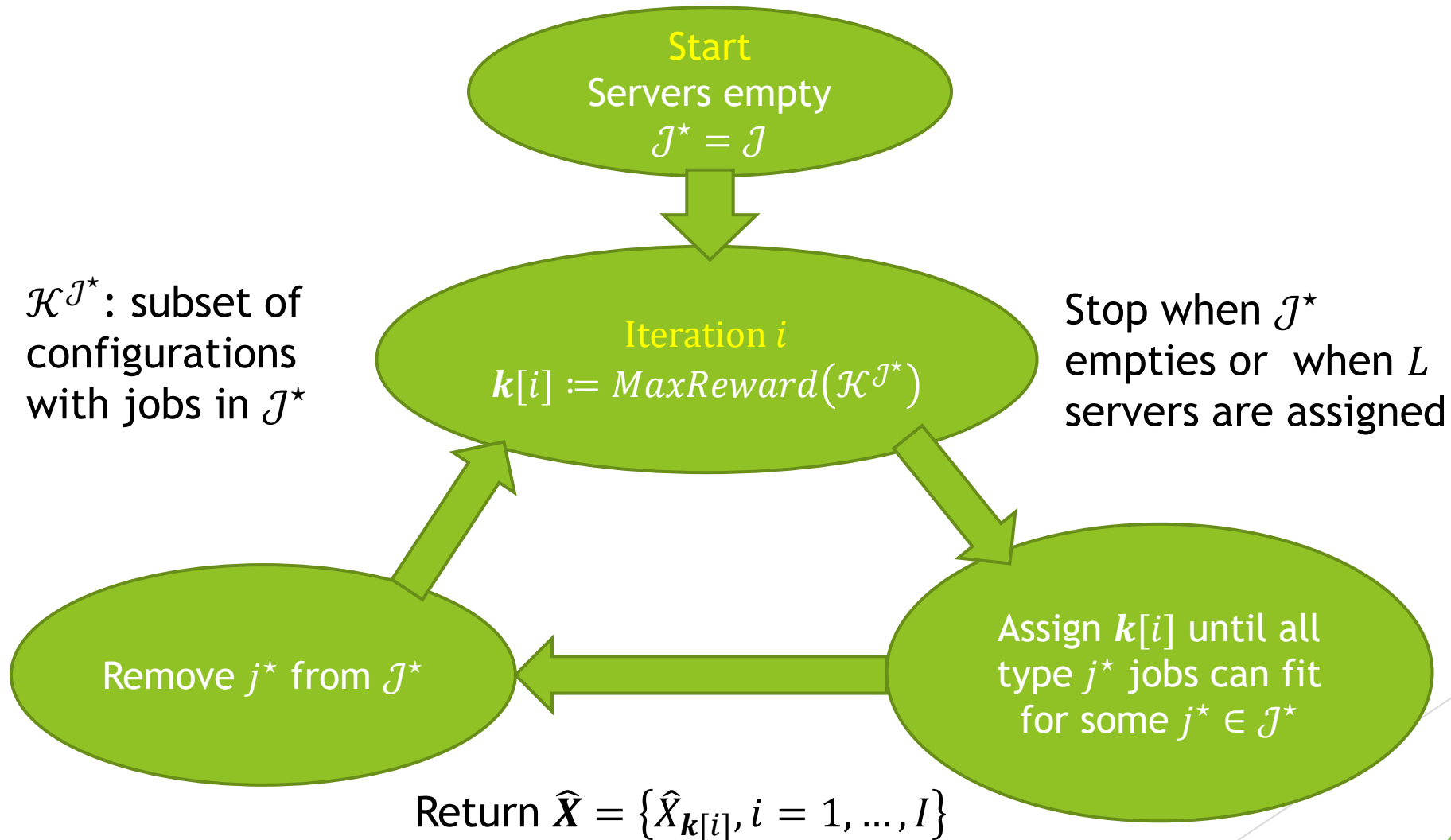
- ▶ Reward of configuration

$$U(\mathbf{k}) = \sum_{j=1}^J u_j k_j$$

- ▶  $U(\mathbf{k})$  induces an ordering for all  $\mathbf{k} \in \mathcal{K}$
- ▶ Define for  $\mathcal{K}_s \subseteq \mathcal{K}$

$$\text{MaxReward}(\mathcal{K}_s) = \arg \max_{\mathbf{k} \in \mathcal{K}_s} U(\mathbf{k})$$

# Solving Static Optimization: Greedy Placement Algorithm (GPA)



# How good is greedy?

Consider normalized rewards ( $L \rightarrow \infty$ )

- ▶ Optimal Reward  $U^*[\rho]$
- ▶ Greedy Reward  $U^{(g)}[\rho]$

Without extra assumptions  $\longrightarrow$

$$U^{(g)}[\rho] \geq \frac{1}{2} U^*[\rho]$$

Monotone Greedy Property  $\longrightarrow$

$$U^{(g)}[\rho] \geq (1 - e^{-1}) U^*[\rho]$$

- ▶ If  $\rho_1 \geq \rho_2$  then  $U^{(g)}[\rho_1] \geq U^{(g)}[\rho_2]$

# Online Algorithm: Server Groups

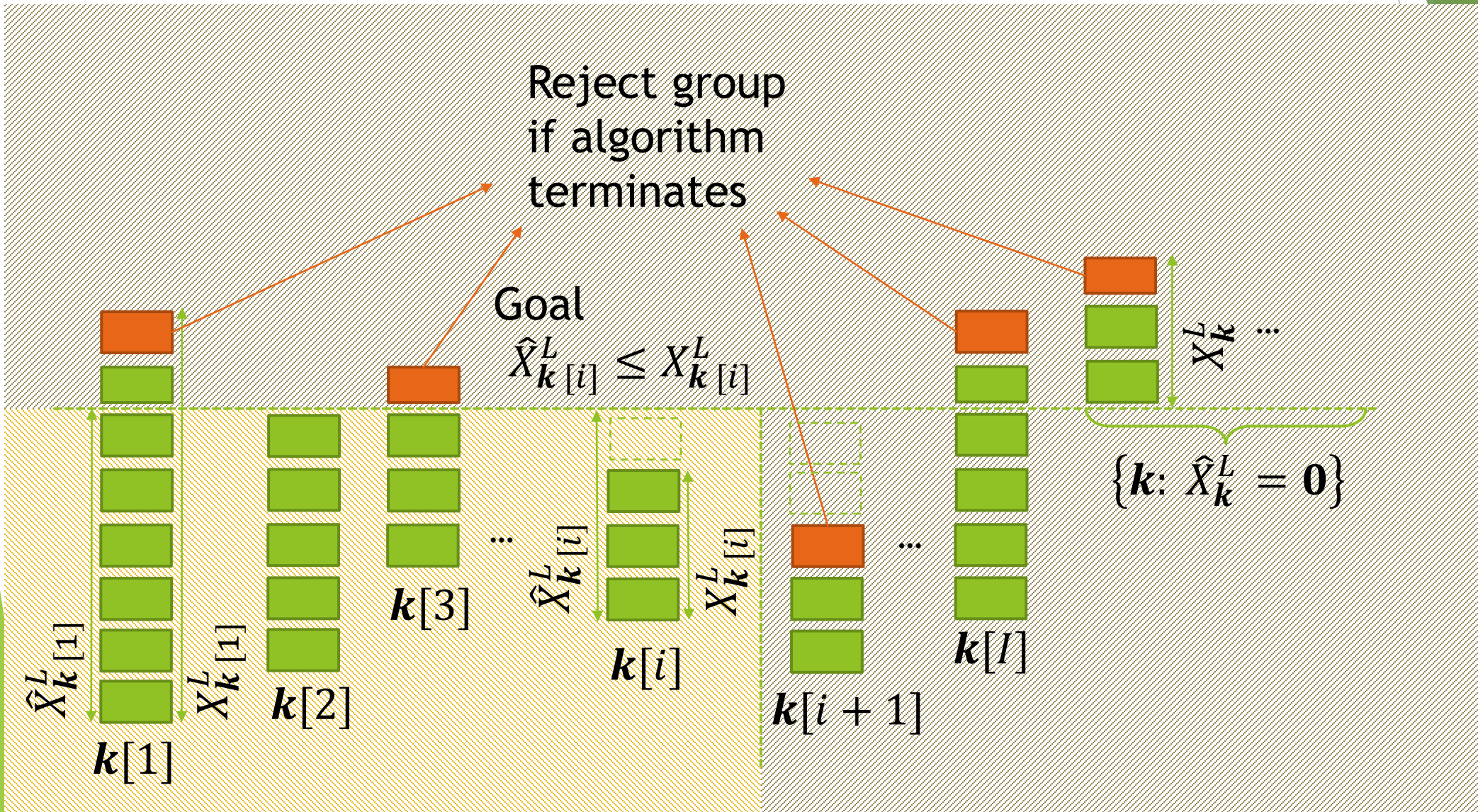
Server Group	Accept	Reject
Goal	Try to Fill	Try to Empty
Schedule in it	Yes	No
Migration	Migrate from a <b>Reject</b> Group filled slot to an <b>Accept</b> Group slot that empties	



# Online Algorithm: (CRA) Classification and Reassignment Algorithm

- ▶ Get solution of GPA
  - ▶ Input  $\hat{Y}^L = Y^L +/\text{- arrival/departure} + g(L)\mathbf{1}_J$
  - ▶  $g(L) = \omega(\log L)$ : reservation factor
- ▶ Match server assignment to GPA solution
  - ▶ Matches configurations in decreasing reward order

# CRA Example [Iteration $i$ ]



# Dynamic Reservation Algorithm (DRA)

## On arrival

- ▶ Run CRA
- ▶ In which **slot** to deploy the **job** arrived
- ▶ Answer: **Any empty slot** in **Accept Group** if exists

## On departure

- ▶ Run CRA
- ▶ Which **job** to migrate in the **slot** that emptied
- ▶ Answer: **Any job** in a slot of **Reject Group** if exists

# Informal Main Result

- ▶ Fraction of servers in each configuration of DRA  $\rightarrow$   
Fraction of servers in each configuration of GPA for input  $\rho L$  when  $L \rightarrow \infty$
- ▶ Normalized Reward of DRA  $\rightarrow$  Greedy Reward  $U^{(g)}[\rho]$

# Simulations (Testing GPA approximation)

Generated 50 random setups

- 6VM types one per pair
- Rewards: 8vCPU + GB
- Servers: 80vCPU, 640GB
- Normalized workload  $\rho_j$  in  $[0.2, 2.0]$

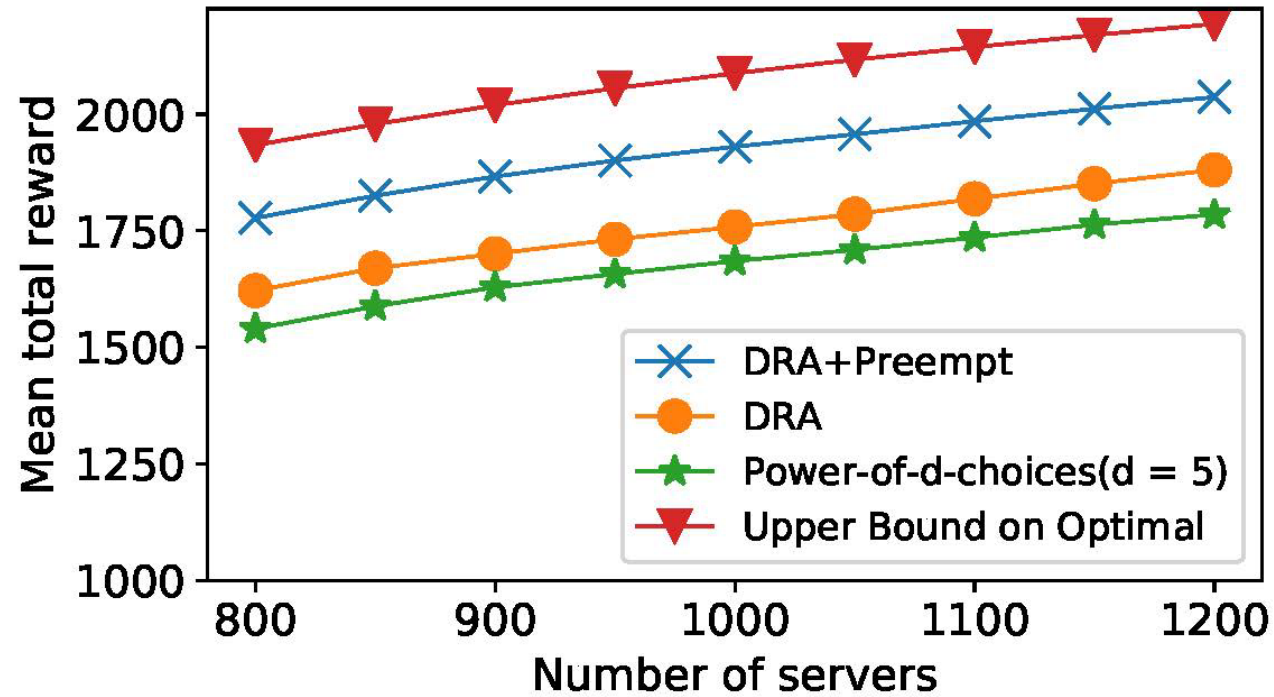
Memory: GB per CPU		
Memory-Opt	CPU-Opt	Regular
8 or 16	1 or 2	4

vCPU	
Small	Large
2 or 4 or 8	32 or 64

Computed  $U^{(g)}[\rho]/U^*[\rho]$  for each setup

- ▶ Worst  $\approx 0.86$  [worst in theory 0.5]
- ▶ Average  $\approx 0.97$
- ▶ Optimal = 1.00 [23/50 setups]

# Simulations (Testing with Google Trace)



- ▶ 1 million tasks
- ▶ 3 priorities
- ▶ 8 different sizes
- ▶  $3 \times 8 = 24$  types

Thank You

