

Proactive Microservice Placement and Migration for Mobile Edge Computing

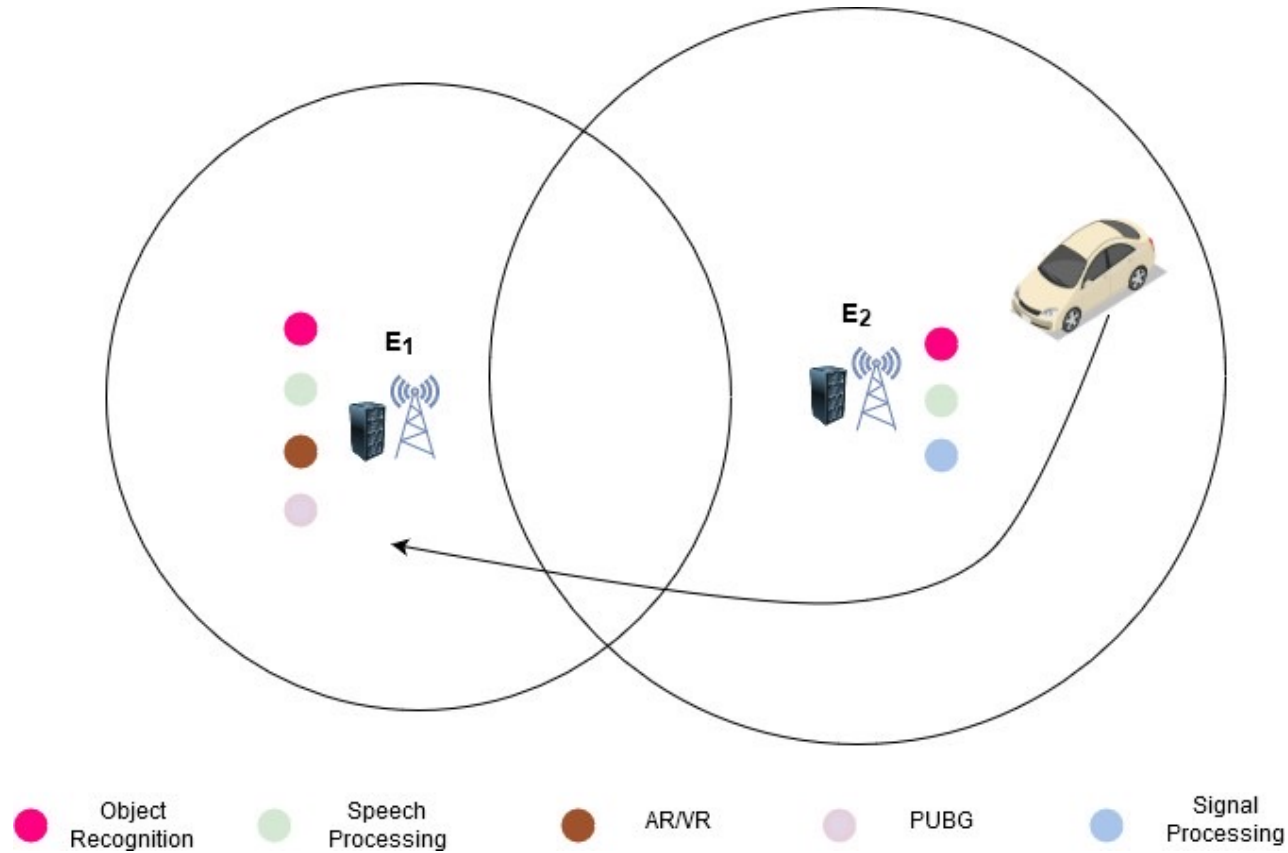
KAUSTABHA RAY & ANSUMAN BANERJEE
ADVANCED COMPUTING AND MICROELECTRONICS UNIT
INDIAN STATISTICAL INSTITUTE, KOLKATA

NANJANGUD C. NARENDRA
ERICSSON RESEARCH
BANGALORE

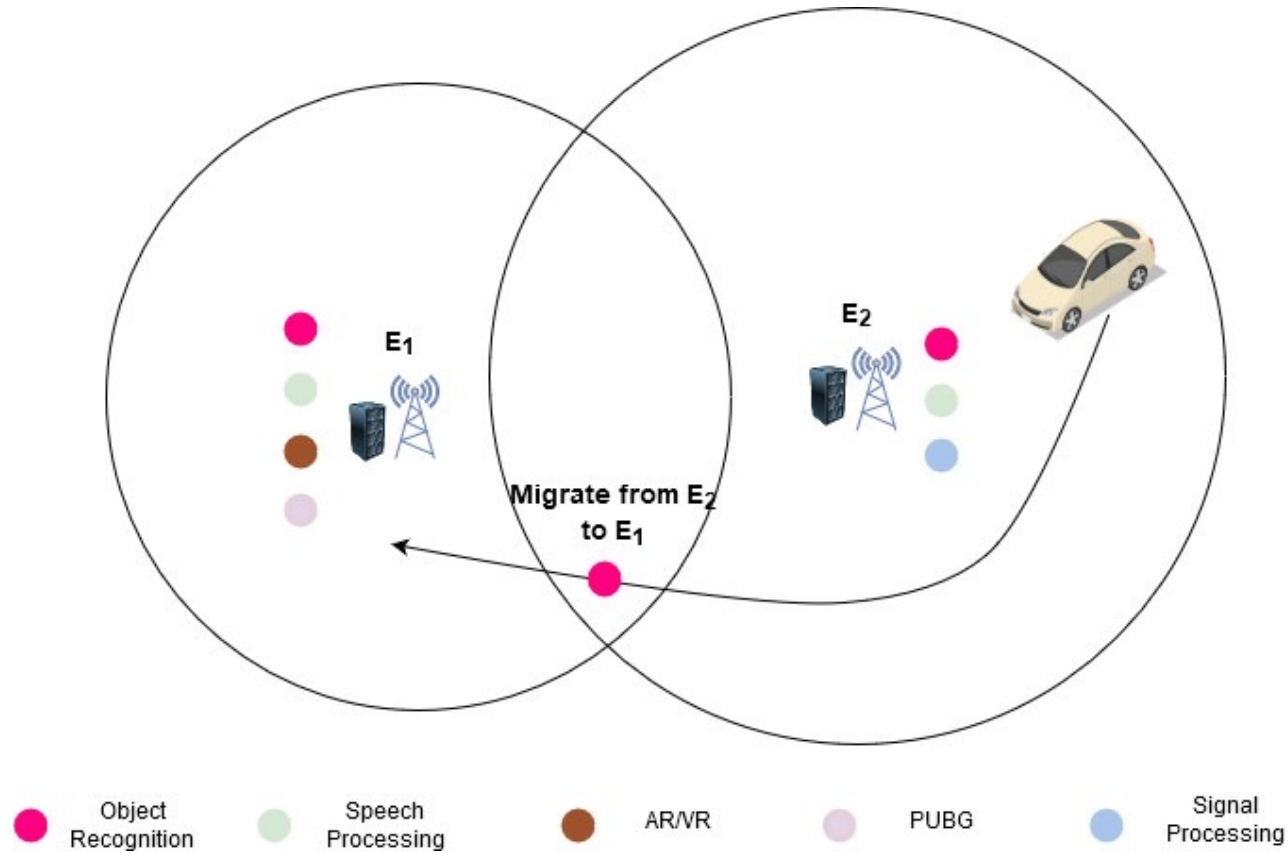
Outline

- ❑ Microservice Placement and Migration Problem
- ❑ Design of Proactive Microservice Placement and Migration Policy
- ❑ Experimental Evaluation

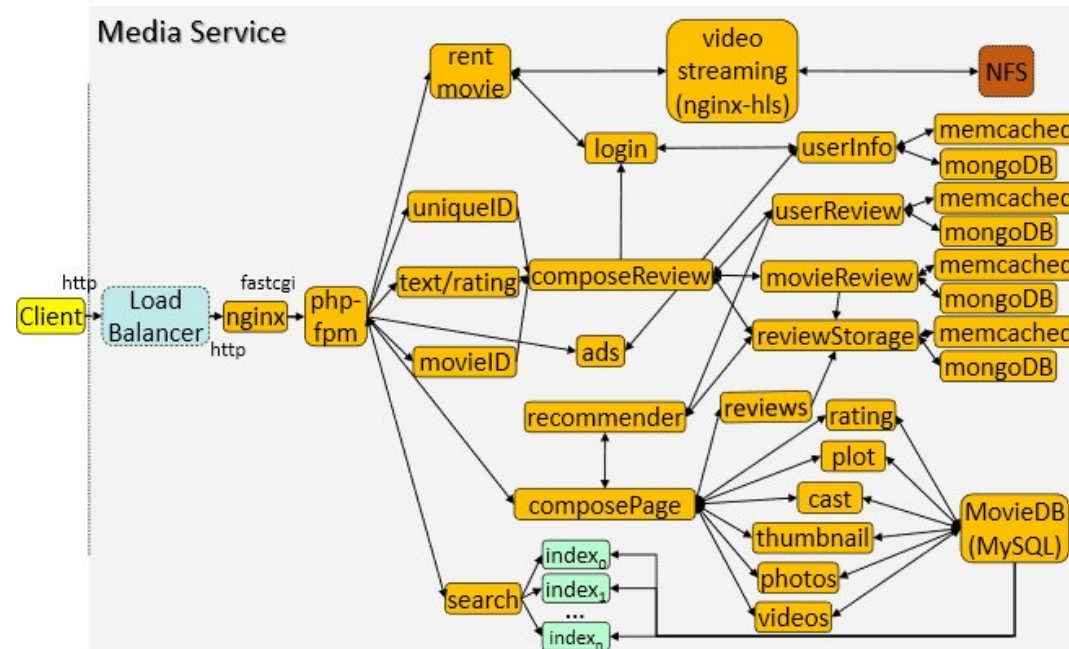
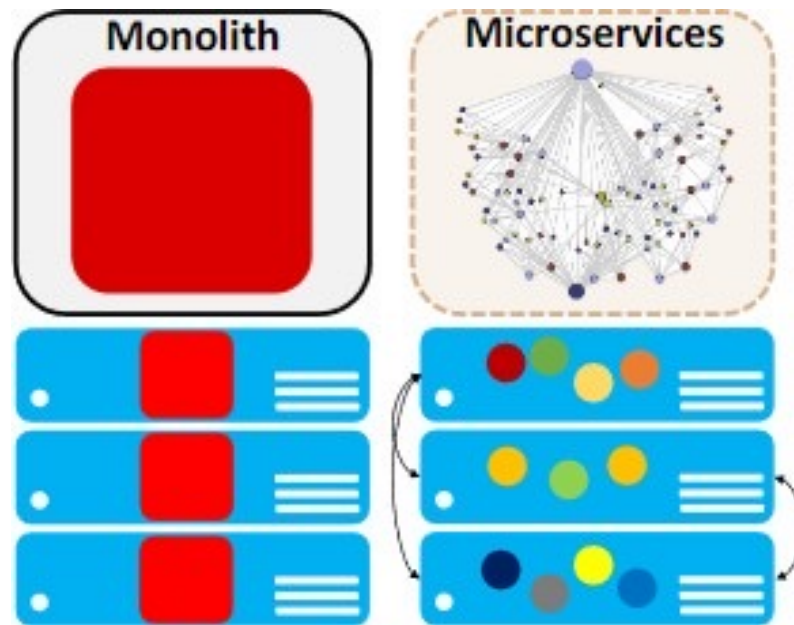
Service Placement and Migration



Service Placement and Migration



From Monoliths to Microservices

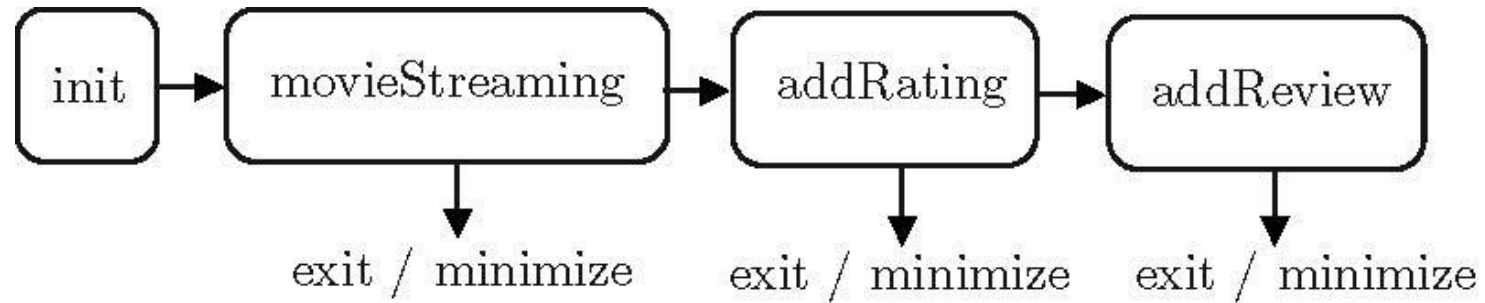


GAN, YU, ET AL. "AN OPEN-SOURCE BENCHMARK SUITE FOR MICROSERVICES AND THEIR
HARDWARE-SOFTWARE IMPLICATIONS FOR CLOUD & EDGE SYSTEMS." *ASPLOS*. 2019

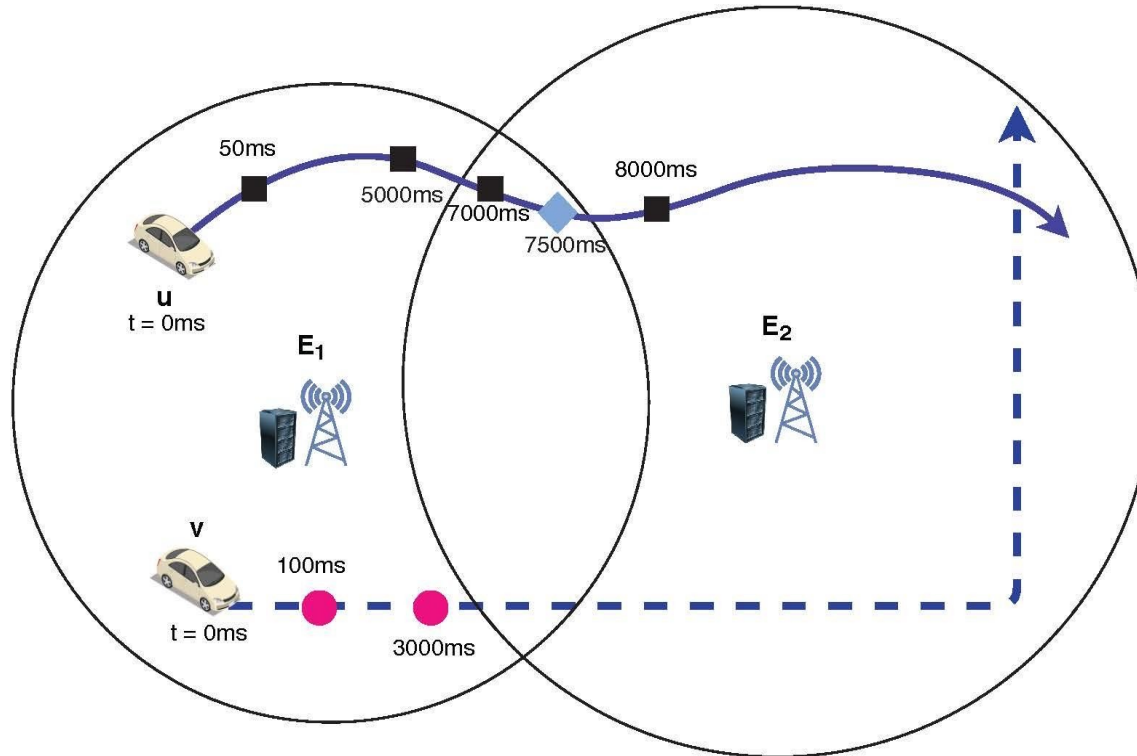
Novelty of this paper

- Considers the paradigm shift from monolithic services to microservices
- Formally model microservice placement and migration using Markov Decision Process (MDP)
- Presents a reinforcement learning based proactive microservice placement and migration strategy

A Motivating Example

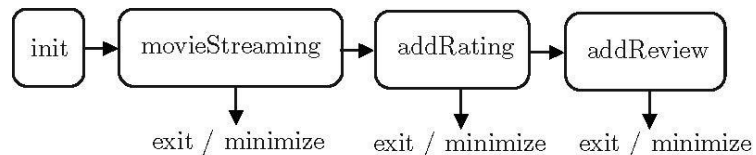


A Motivating Example

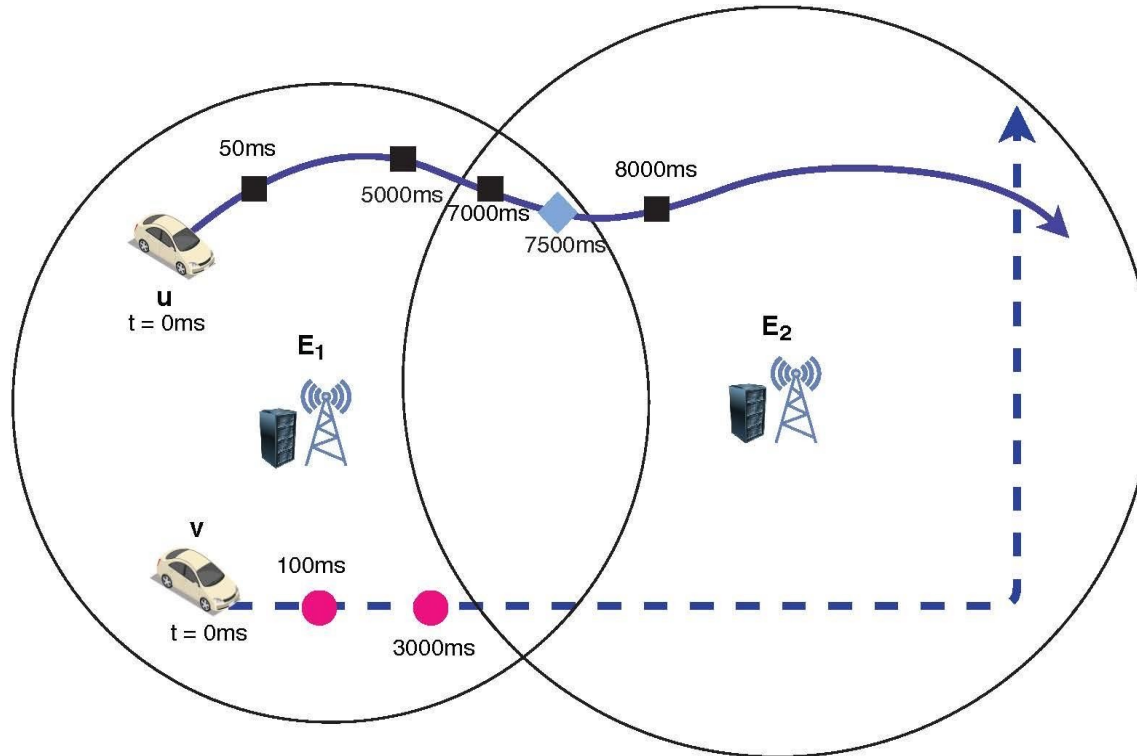


Time t	User Action	Server-Service State
$0ms$		No Services Deployed
$50ms$	$u \rightarrow \text{movieStreaming}$	initialize movieStreaming
$75ms$		$E_1 \rightarrow \text{movieStreaming}$
$100ms$	$v \rightarrow \text{movieStreaming}$	$E_1 \rightarrow \text{movieStreaming}$ initialize new task for v
$110ms$		$E_1 \rightarrow \text{movieStreaming}, v_{task}$
$3000ms$	v exits movieStreaming	$E_1 \rightarrow \text{movieStreaming}$
$5000ms$	$u \rightarrow \text{addRating}$	initialize addRating
$5025ms$		$E_1 \rightarrow \text{addRating}$
$7000ms$	u minimizes addRating	$E_1 \rightarrow \text{addRating}$
$8000ms$	$u \rightarrow \text{addReview}$	initialize addReview
$8025ms$		$E_2 \rightarrow \text{addReview}$

On-Demand Placement

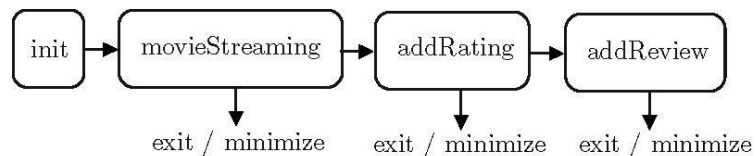


A Motivating Example

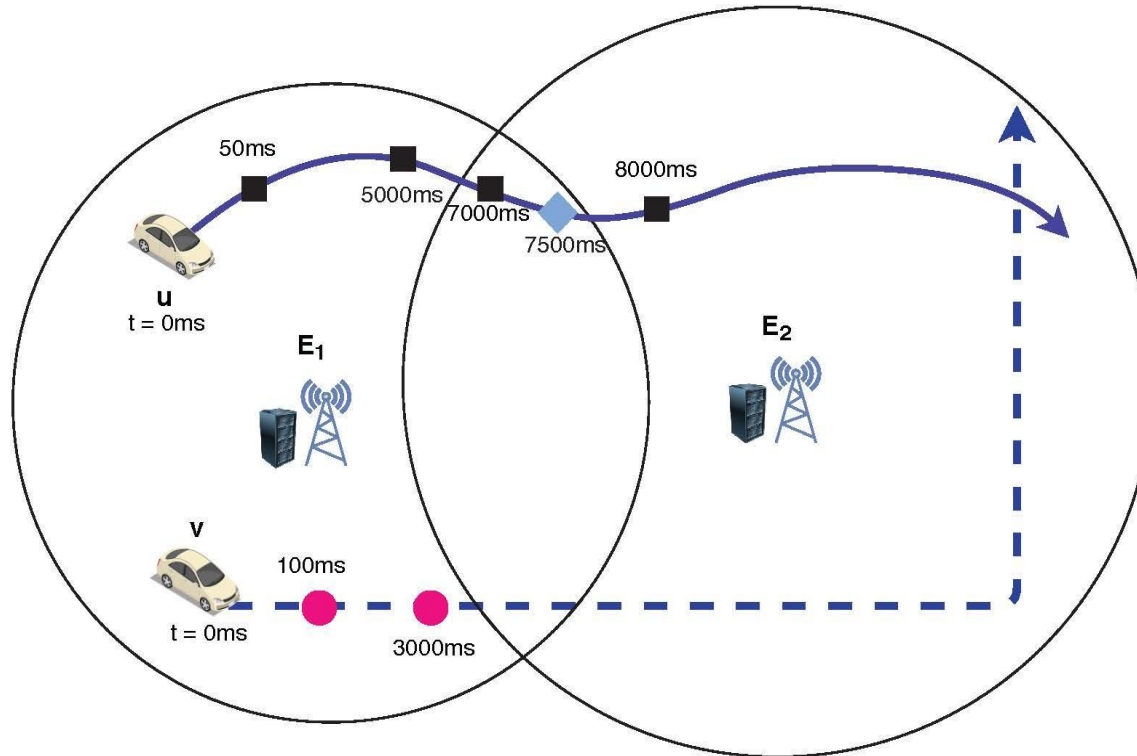


Time t	User Action	Server-Service State
$0s$		No Services Deployed
$50ms$	$u \rightarrow \text{movieStreaming}$	initialize movieStreaming
$75ms$		$E_1 \rightarrow \text{movieStreaming}$
$100ms$		$E_1 \rightarrow \text{movieStreaming, addRating}$
$100ms$	$v \rightarrow \text{movieStreaming}$	$E_1 \rightarrow \text{movieStreaming, addRating}$ initialize new task for v
$110ms$		$E_1 \rightarrow \text{movieStreaming, addRating, addReview, } v_{task}$
$135ms$		$E_1 \rightarrow \text{movieStreaming, addRating, addReview, } v_{task}$
$3000ms$	v exits movieStreaming	$E_1 \rightarrow \text{movieStreaming, addRating, addReview}$
$5000ms$	$u \rightarrow \text{addRating}$	$E_1 \rightarrow \text{addRating, addReview}$
$7000ms$	u minimizes addRating	$E_1 \rightarrow \text{addRating, addReview}$
$8000ms$	$u \rightarrow \text{addReview}$	state-aware migrate addReview
$8010ms$		$E_2 \rightarrow \text{addReview}$

Proactive Placement

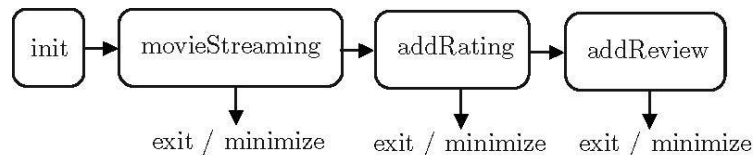


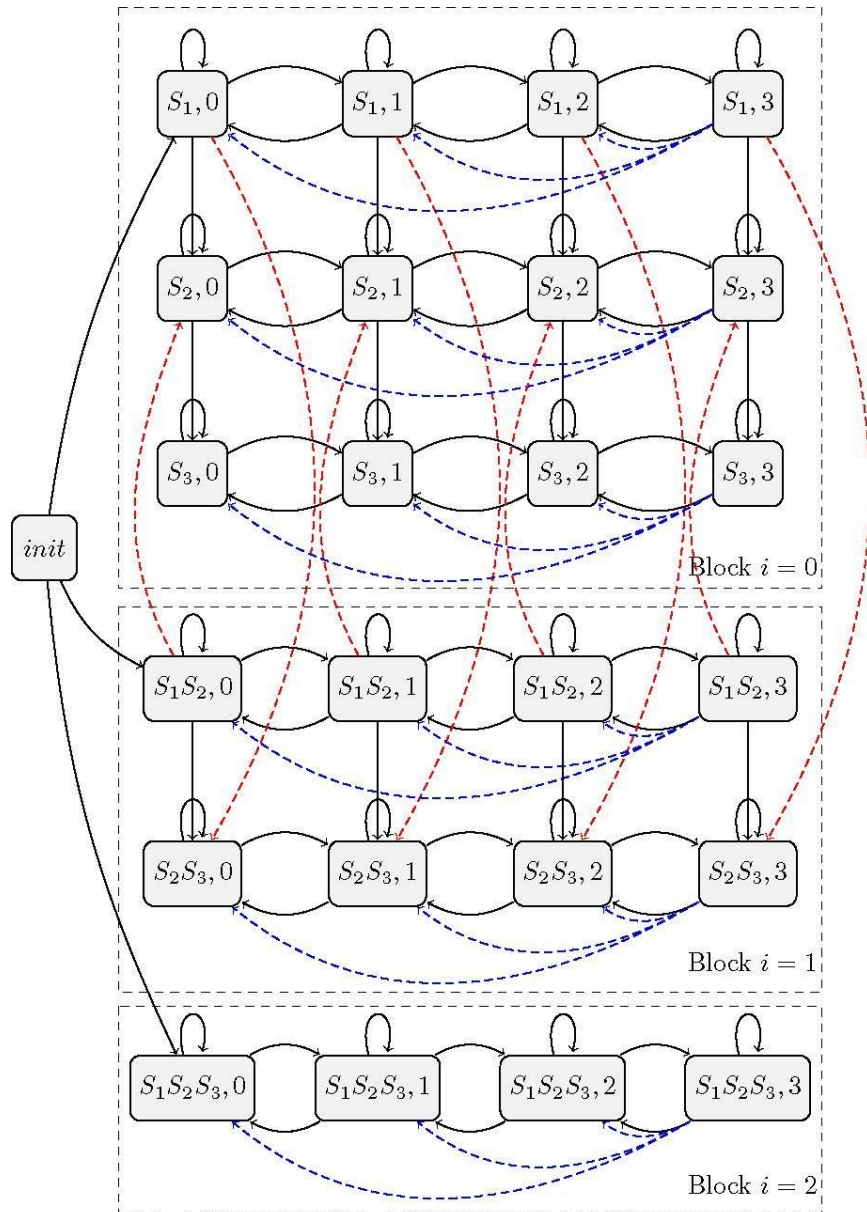
A Motivating Example



Time t	User Action	Server-Service State
$0s$		No Services Deployed
$50ms$	$u \rightarrow \text{movieStreaming}$	initialize movieStreaming
$75ms$		$E_1 \rightarrow \text{movieStreaming}$
$100ms$		$E_1 \rightarrow \text{movieStreaming, addRating}$
$100ms$	$v \rightarrow \text{movieStreaming}$	$E_1 \rightarrow \text{movieStreaming, addRating}$ initialize new task for v
$110ms$		$E_1 \rightarrow \text{movieStreaming, addRating, addReview, } v_{task}$
$135ms$		$E_1 \rightarrow \text{movieStreaming, addRating, addReview, } v_{task}$
$3000ms$	v exits movieStreaming	$E_1 \rightarrow \text{movieStreaming, addRating, addReview}$
$5000ms$	$u \rightarrow \text{addRating}$	$E_1 \rightarrow \text{addRating, addReview}$
$7500ms$	$u \rightarrow \text{addRating}$	migrate addRating, addReview
$7555ms$	$u \rightarrow \text{addRating}$	$E_2 \rightarrow \text{addRating, addReview}$
$8000ms$	$u \rightarrow \text{addReview}$	$E_2 \rightarrow \text{addReview}$

Proactive Placement + Migration

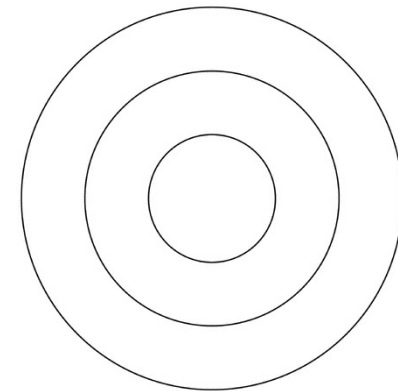




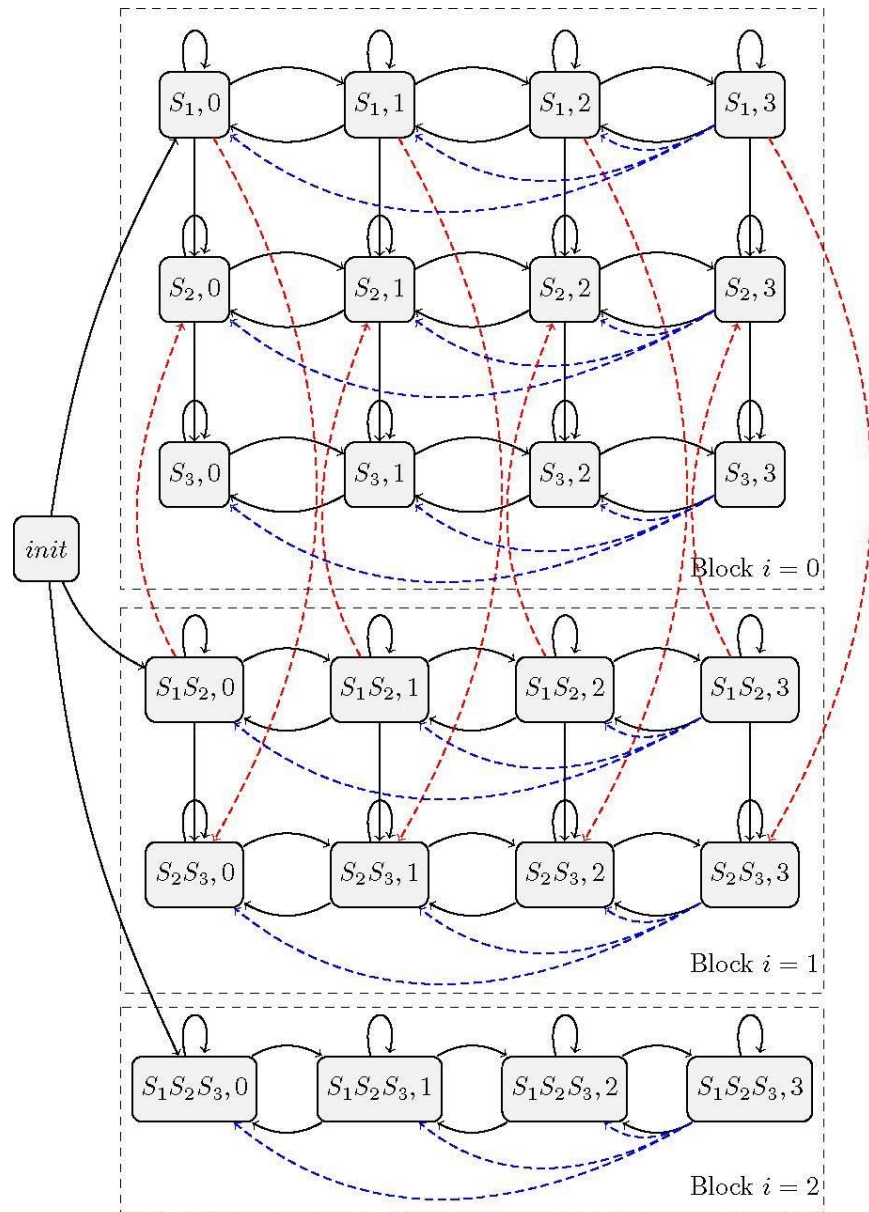
Formal Model

States Represent Proactive Placement of Microservices

Transitions Represent Movement and Choices of Proactive Placement



How distance is mapped

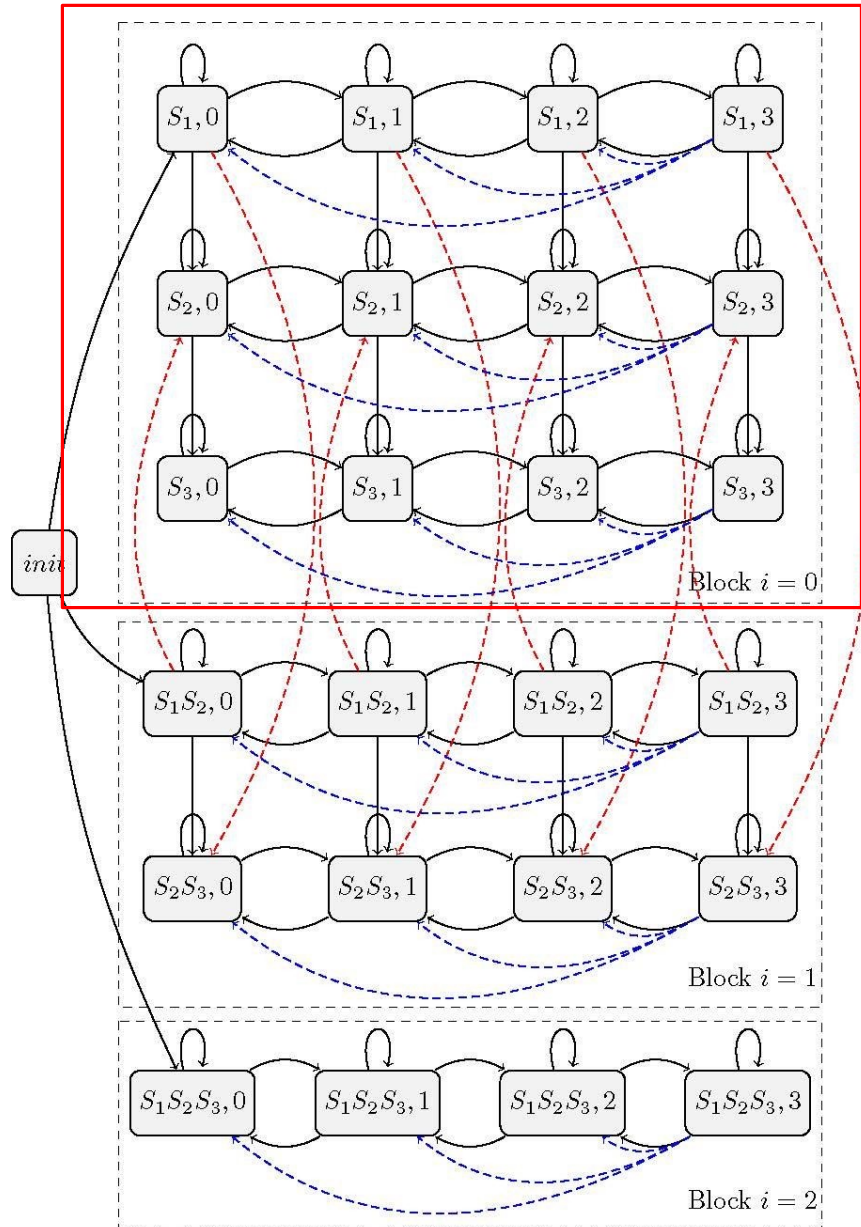


Formal Model

Blocks Represent i number of microservices to proactively deploy

Blocks $i=0$ to $i=n-1$

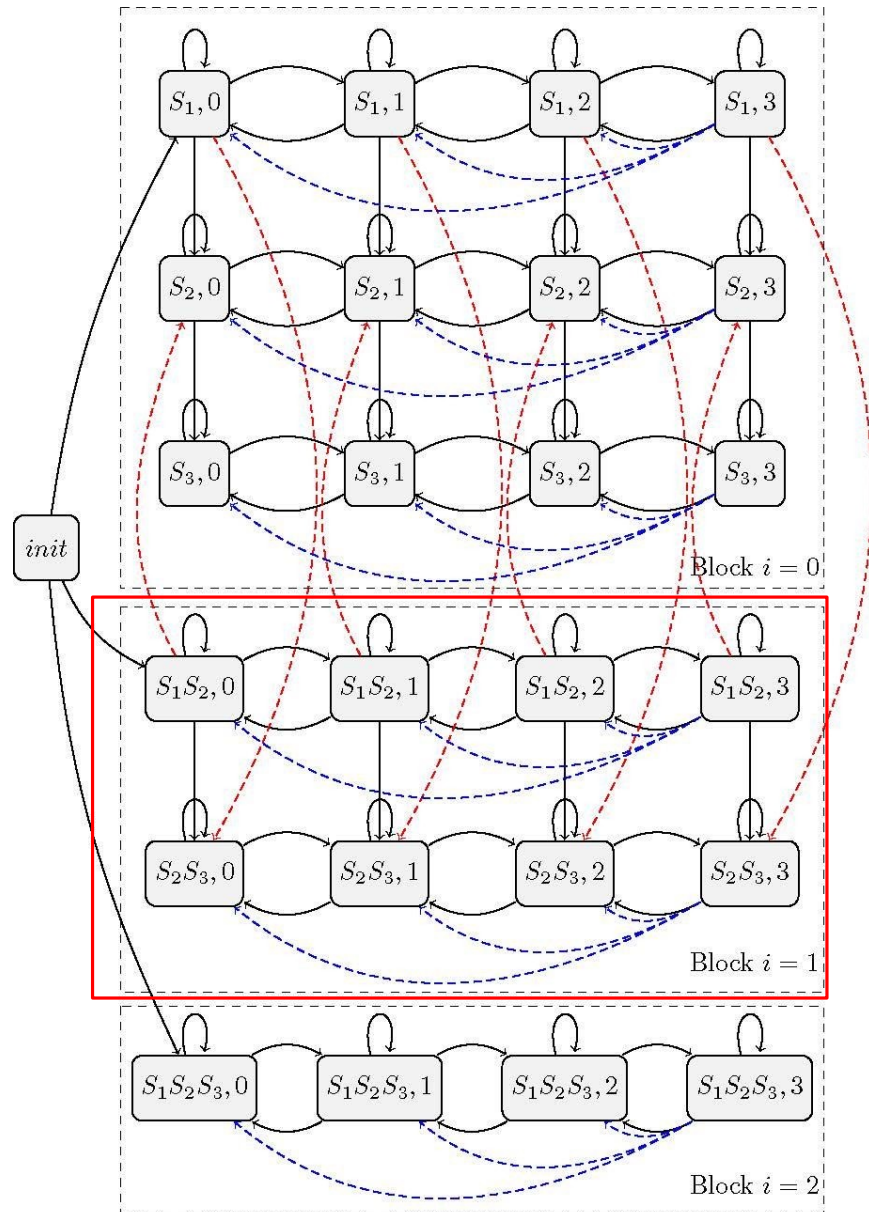
Why don't we have s_1, s_3 ?



Formal Model

Blocks Represent i number of microservices to proactively deploy

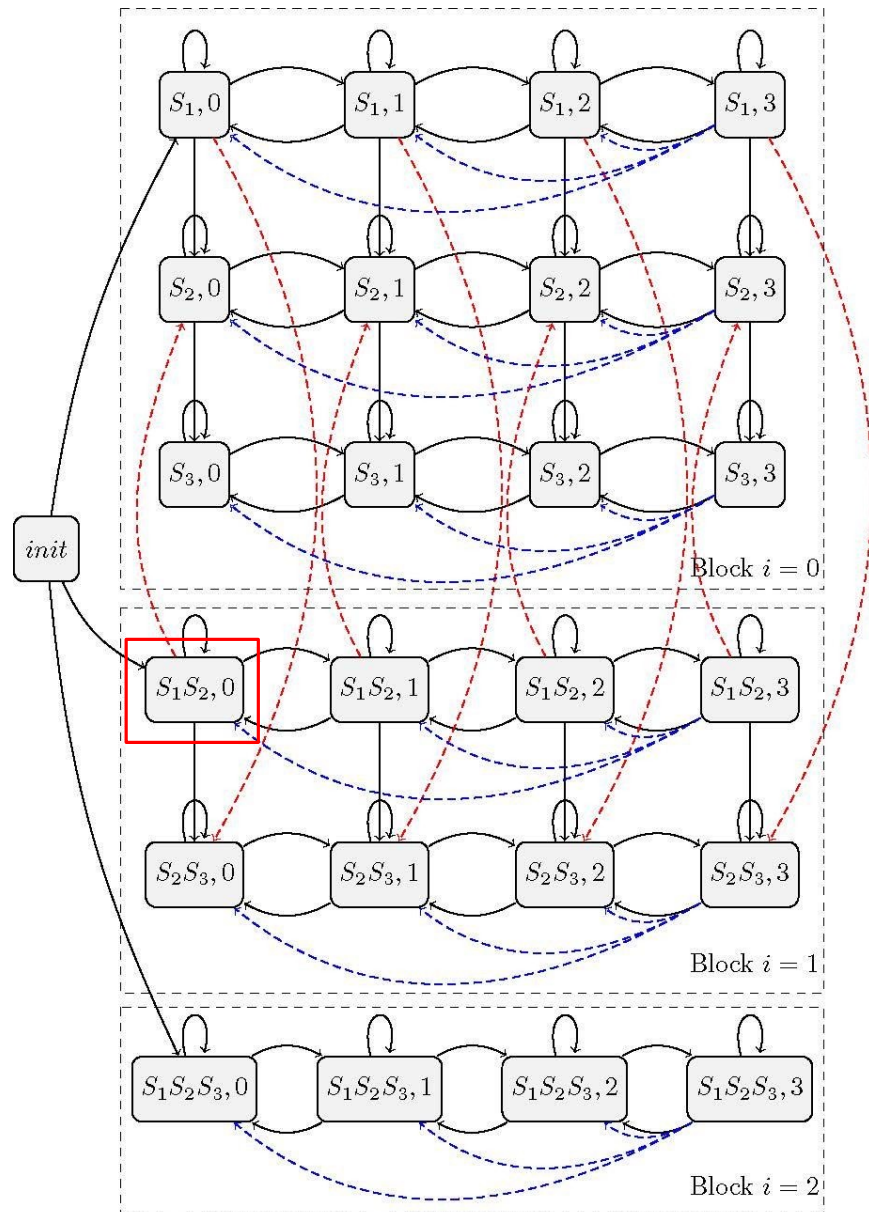
Blocks $i=0$ to $i=n-1$



Formal Model

Blocks Represent i number of microservices to proactively deploy

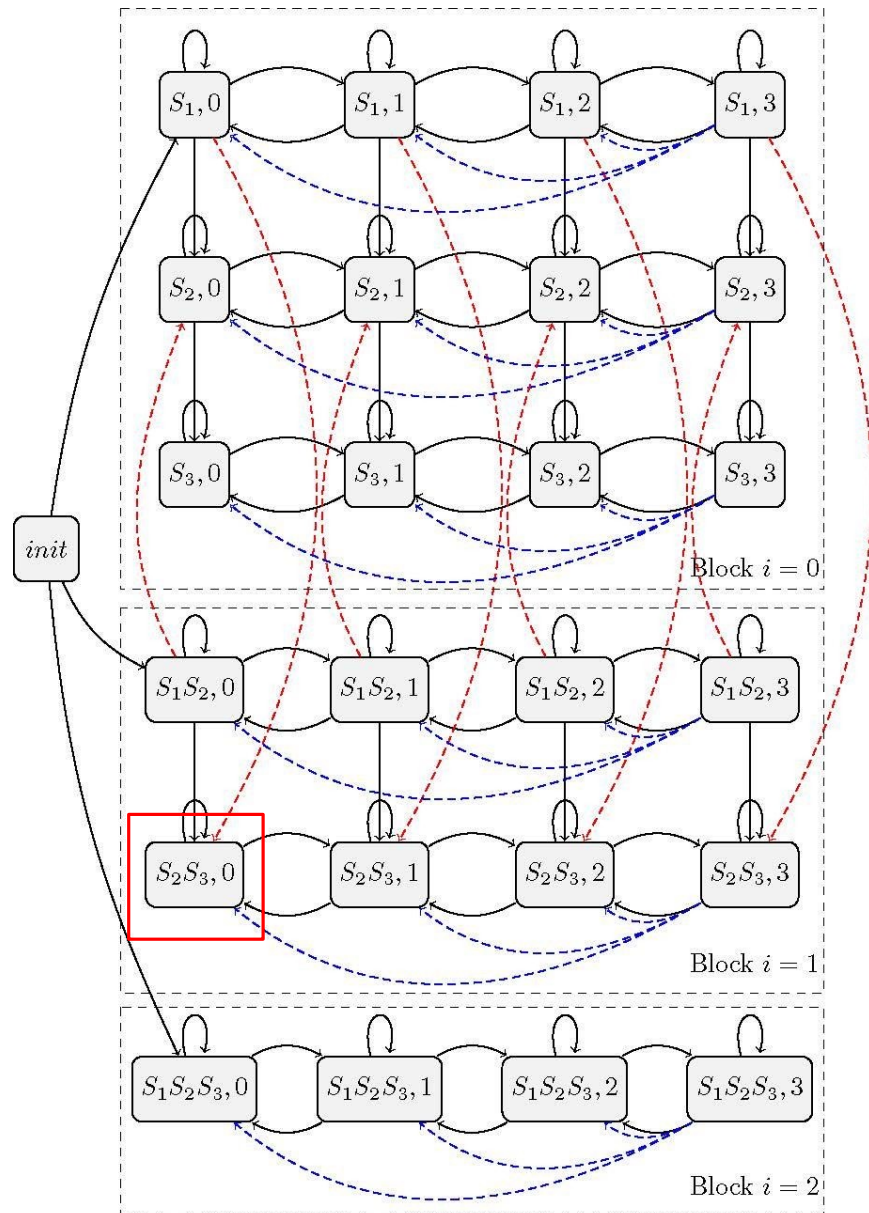
Blocks $i=0$ to $i=n-1$



Formal Model

Blocks Represent i number of microservices to proactively deploy

Blocks $i=0$ to $i=n-1$



Formal Model

Blocks Represent i number of microservices to proactively deploy

Blocks $i=0$ to $i=n-1$

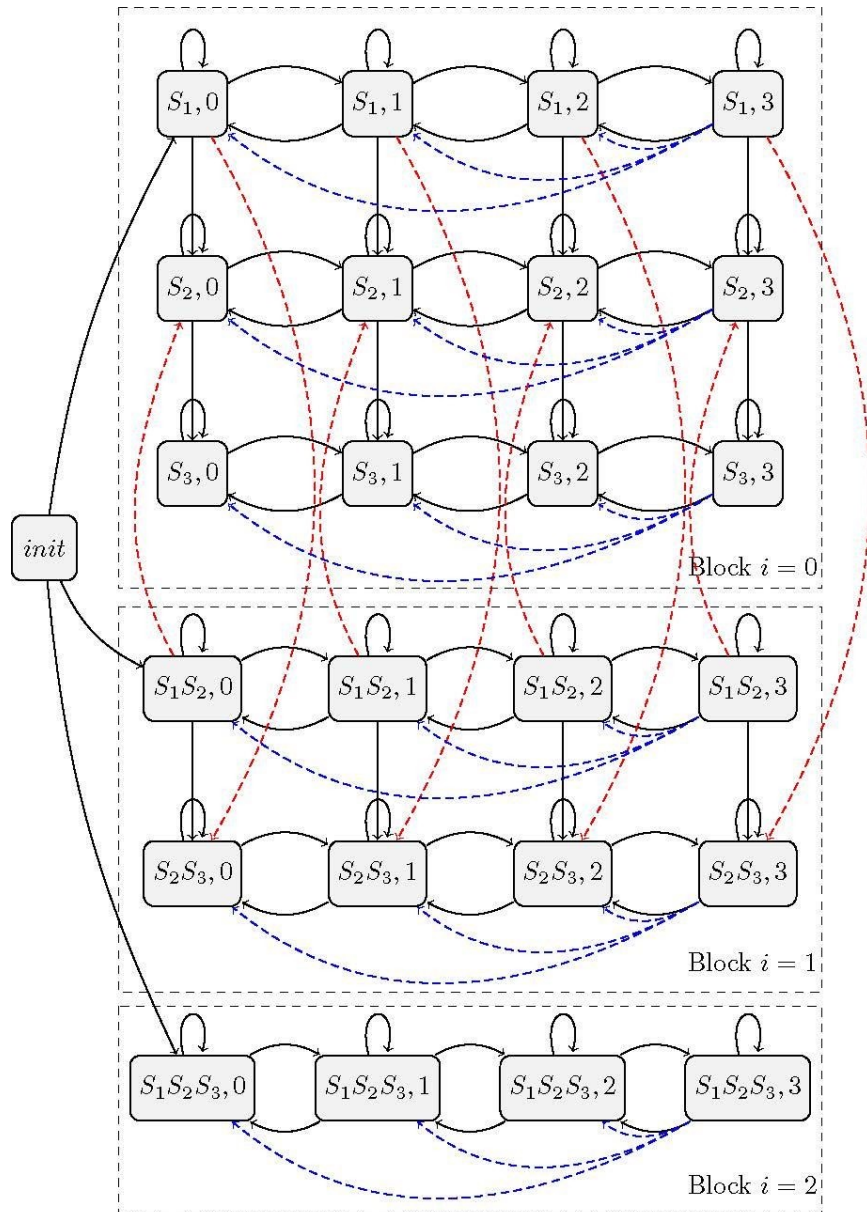
Formal Model

Three types of transitions

□ User Movement

□ Service Invocation in Application Topology

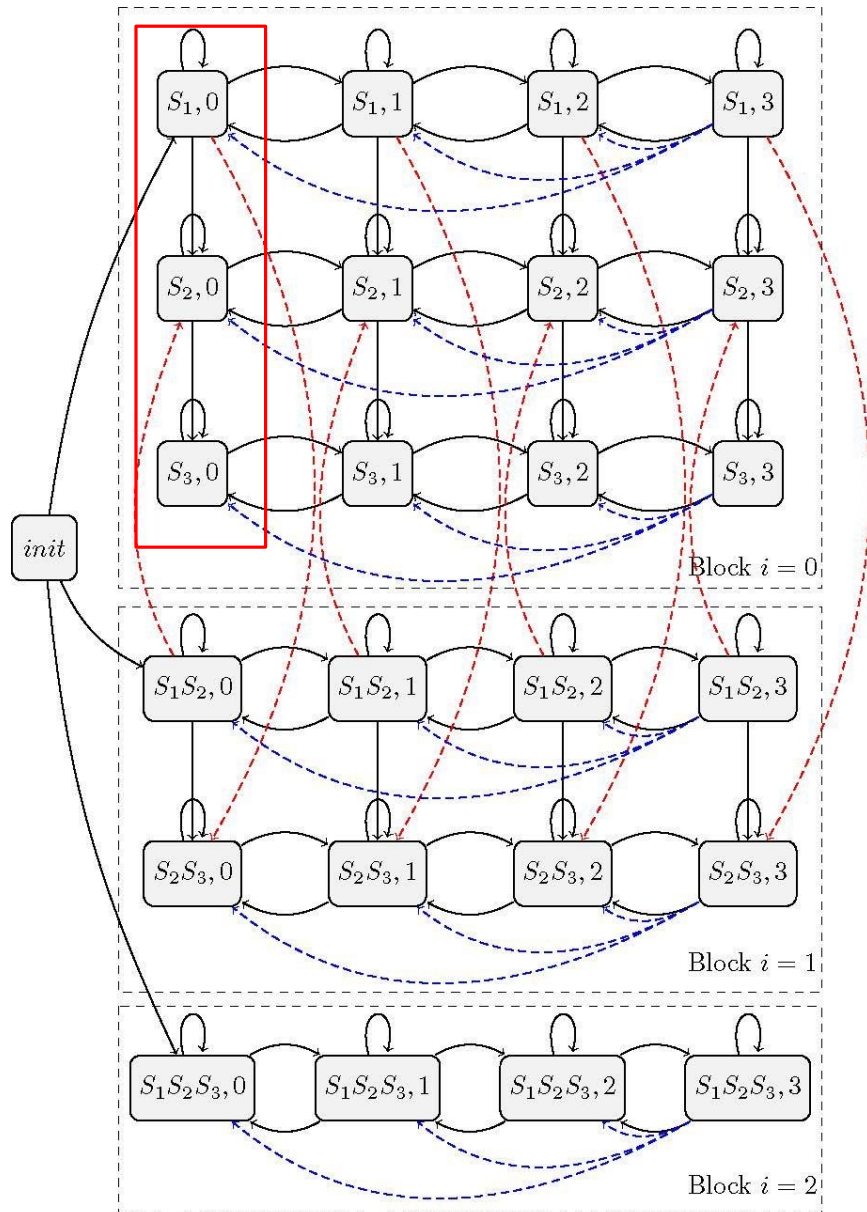
□ Migration



Formal Model

Three types of transitions

- User Movement
- Service Invocation in Application Topology
- Migration



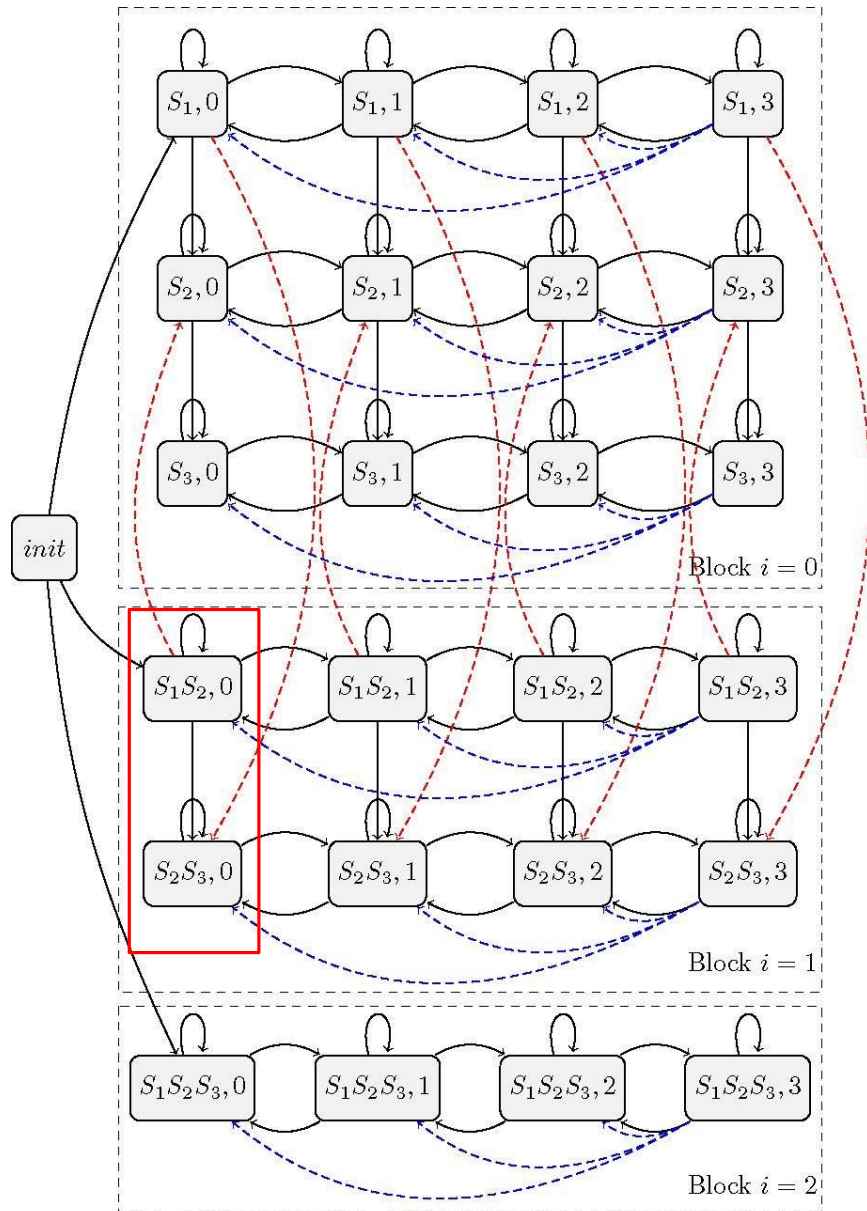
Formal Model

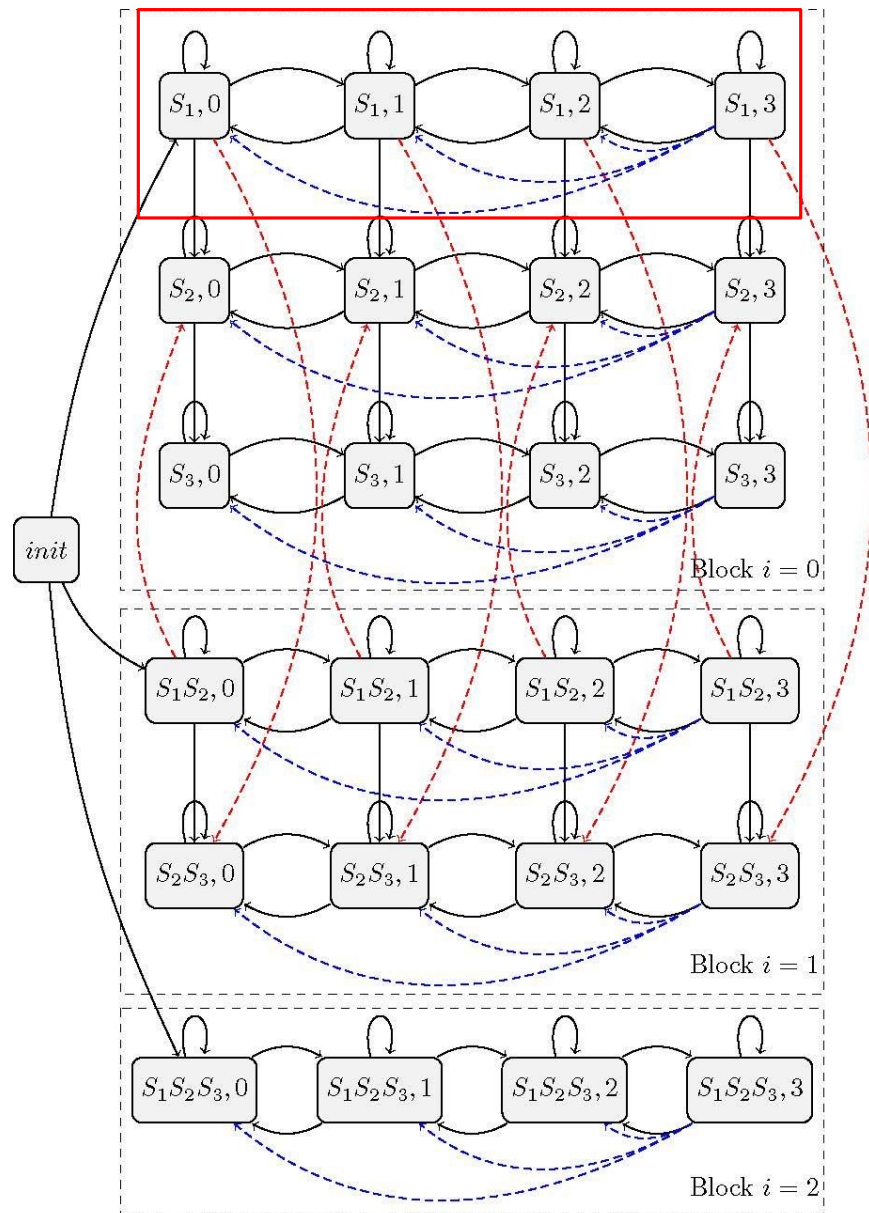
Three types of transitions

□ User Movement

□ Service Invocation in Application Topology

□ Migration





Formal Model

Three types of transitions

- ▣ User Movement
- ▣ Service Invocation in Application Topology
- ▣ Migration

Reinforcement Learning Solution

Algorithm 1: Dyna-Q

```
1 Initialize  $Q(s, a)$  and  $Model(s, a)$ ,  $\forall s \in \mathcal{S}, \forall a \in A(s)$ 
2 while true do
3    $s \leftarrow$  observe the application state
4    $a \leftarrow \epsilon$ -greedy( $s, q$ )
5   Observe the next state  $s'$  and the reward obtained
6   Update  $Q(s, a)$  using Equation 1
7    $Model(s, a) \leftarrow r, s'$ 
8   for  $i = 0 \dots n$  do
9      $s \leftarrow$  random state previously observed
10     $a \leftarrow$  random action previously taken in  $s$ 
11     $r, s' \leftarrow Model(s, a)$ 
12    Update  $Q(s, a)$  using Equation 1
```

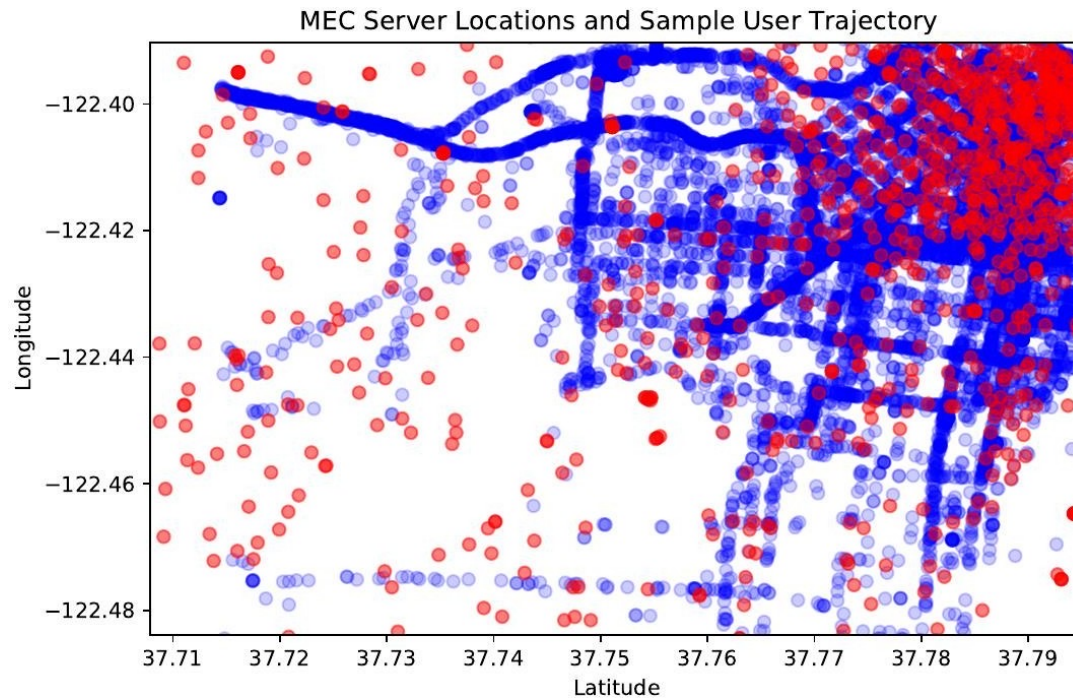
- ❑ Use Dyna-Q : Model Based and Model Free
- ❑ Simulation + Interaction
- ❑ Reward Function defined as a measure of prefetched and utilized services and prefetched and unutilized services
- ❑ Possible since we only transition upon service invocations

$$R = \sum_{\mu \in \mu_{used}} [\mu * c(\mu_{resources})] - \sum_{\mu \in \mu_{unused}} [\mu * c(\mu_{resources})]$$

Reinforcement Learning Solution

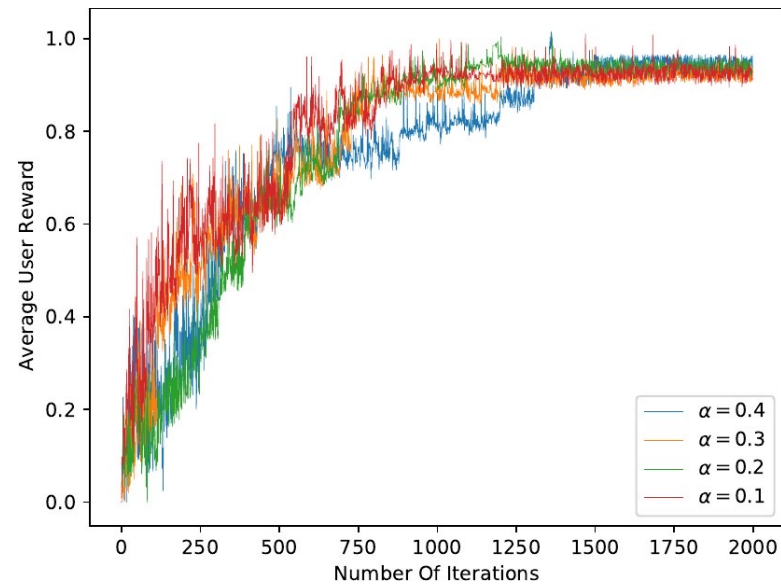
- ❑ Low Traffic some action receives a positive reward
- ❑ High Traffic same action may receive negative reward
- ❑ Variance leads to confusion
- ❑ Heuristic Solution : Three types of MDPs for each Application – *{high, medium, low}* – use the appropriate MDP depending on the traffic condition
- ❑ Capacity Constraint Heuristic : Allocate microservices greedily along the linear chain

Dataset

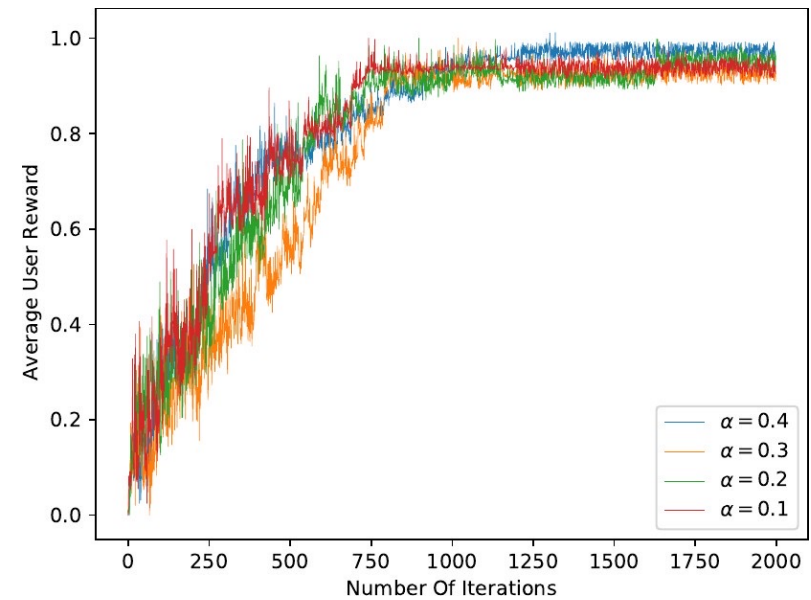


- San Francisco Taxi Dataset for User Trajectories
<https://crawdad.org/epfl/mobility/20090224/>
- San Francisco Wireless Telecommunications Services Facilities Dataset for MEC Server Locations
<https://data.sfgov.org/Geographic-Locations-and-Boundaries/Existing-Commercial-Wireless-Telecommunication-Ser/aa26-h926>
- Service invocations randomly generated along with representative timing for initialization from DeathStarBench suite

Accumulated Reward

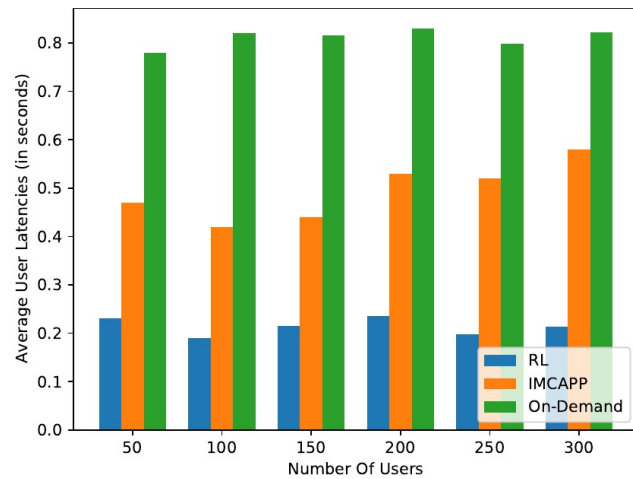


Users = 50

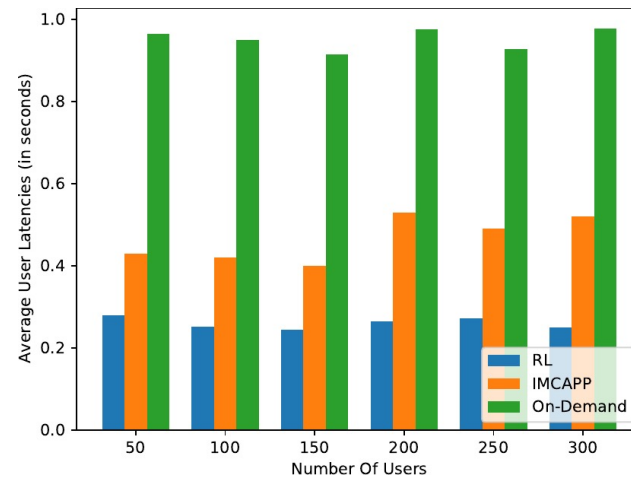


Users = 100

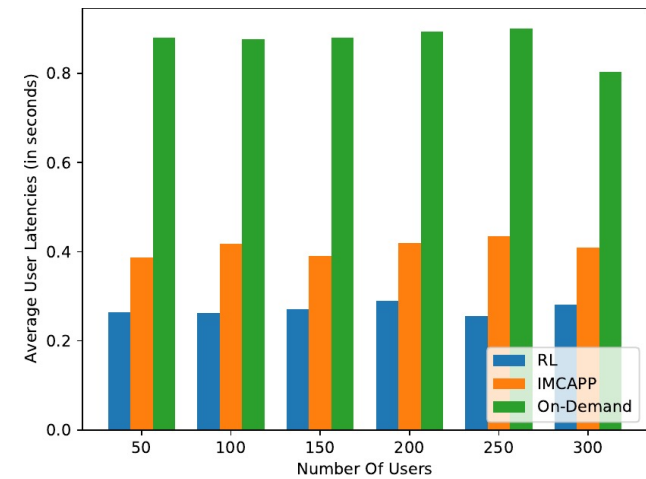
Average User Latency



4 Microservices

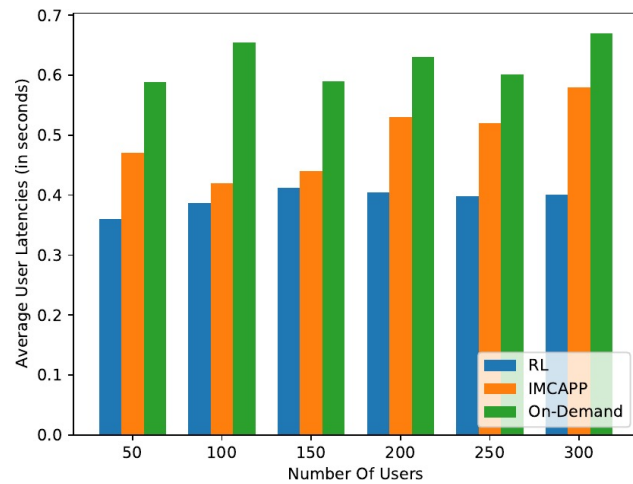


8 Microservices

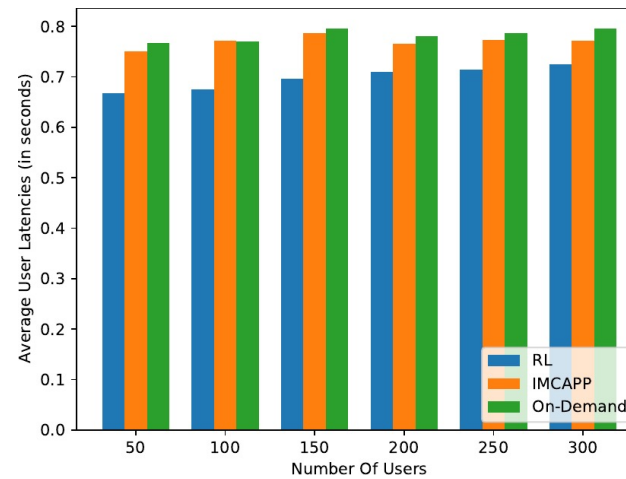


12 Microservices

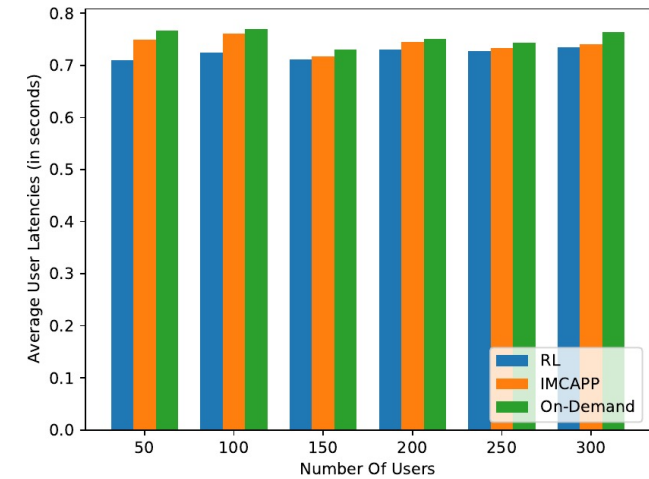
Average User Latency



Server Resources = 130%

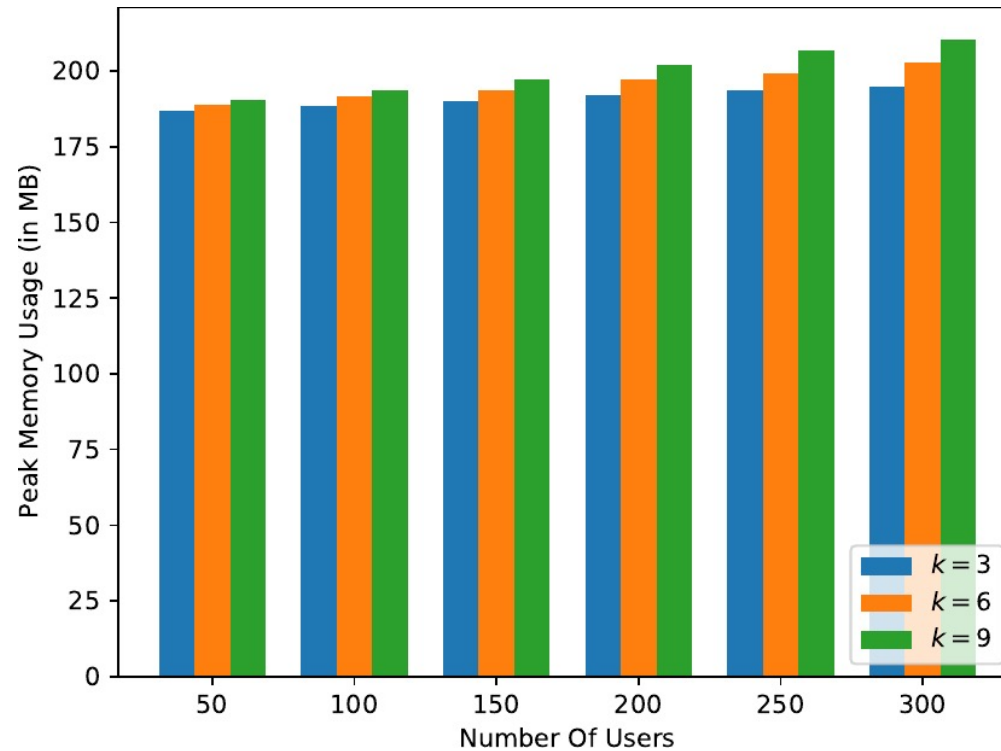


Server Resources = 100%



Server Resources = 65%

Memory Usage for Varying k



Thank You!

