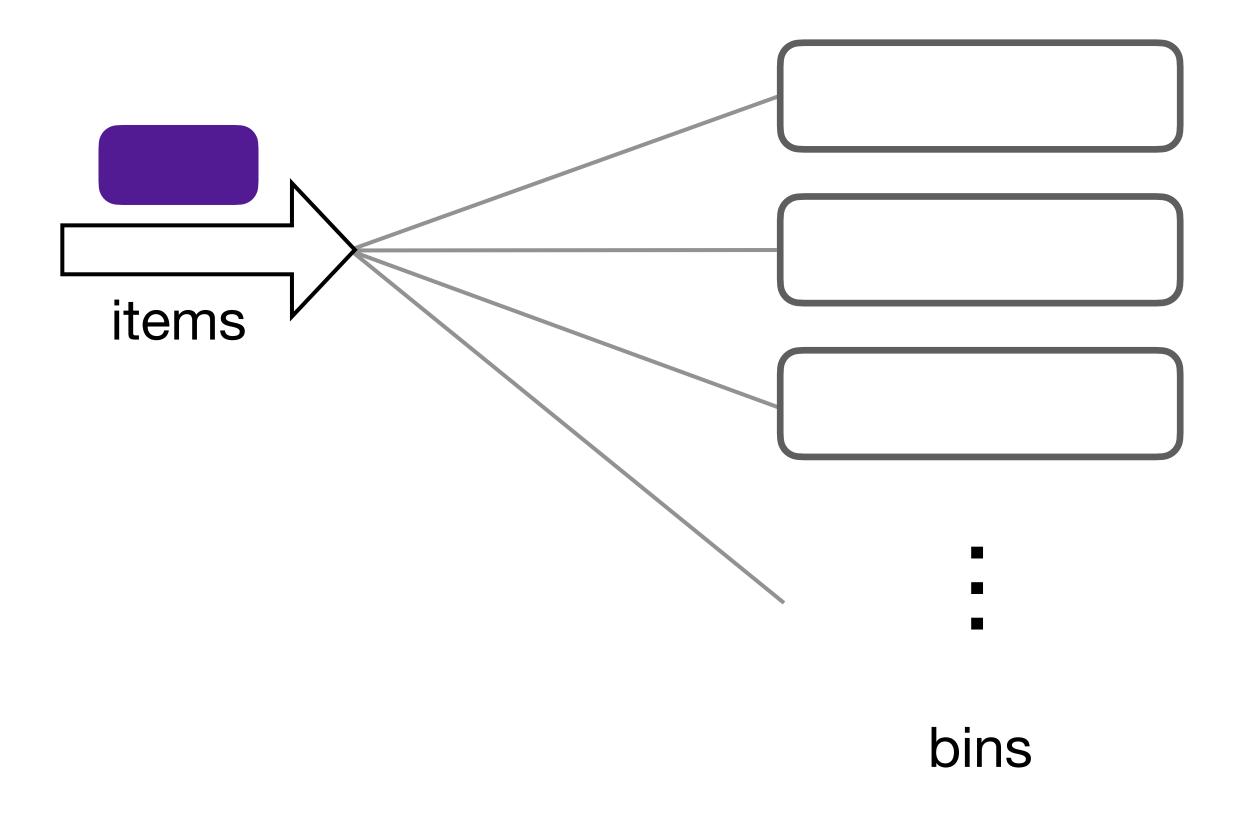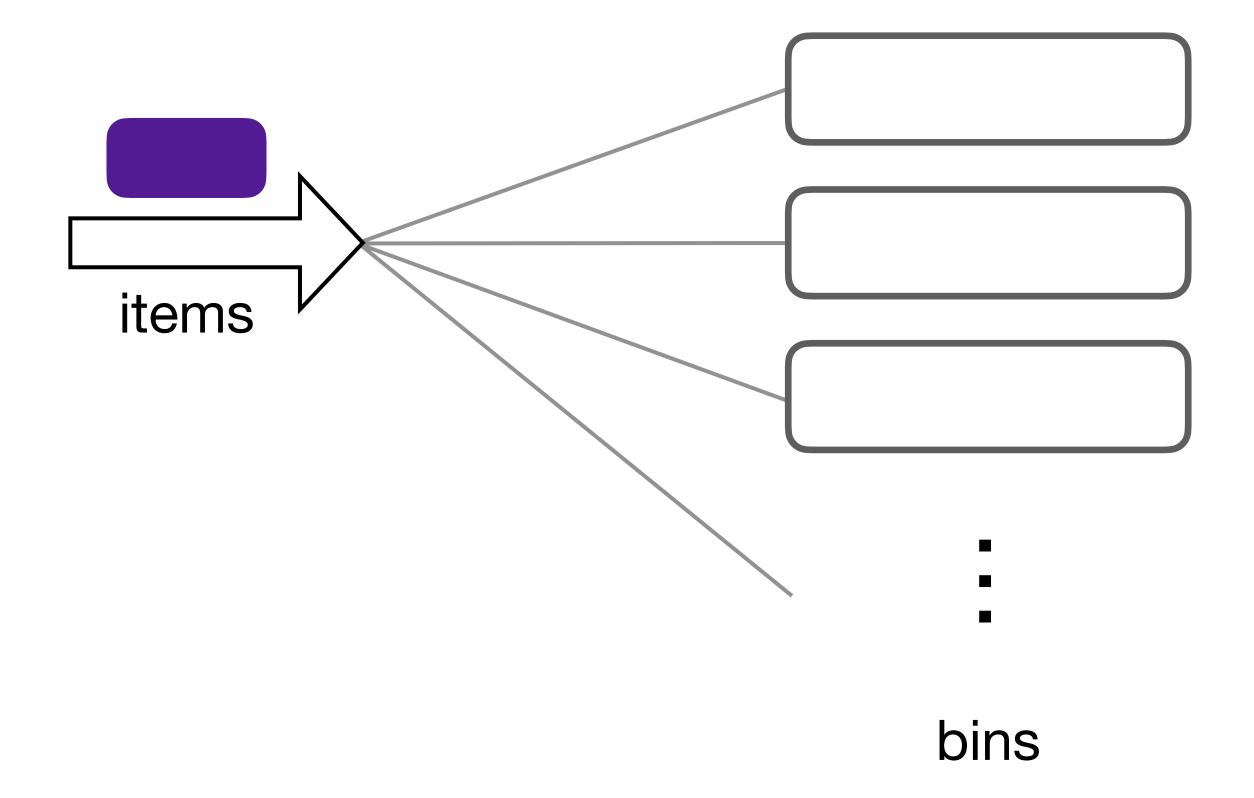# Stochastic Bin Packing with Time-Varying Item Sizes

Joint work with Yige Hong (CMU) and Qiaomin Xie (UW–Madison)
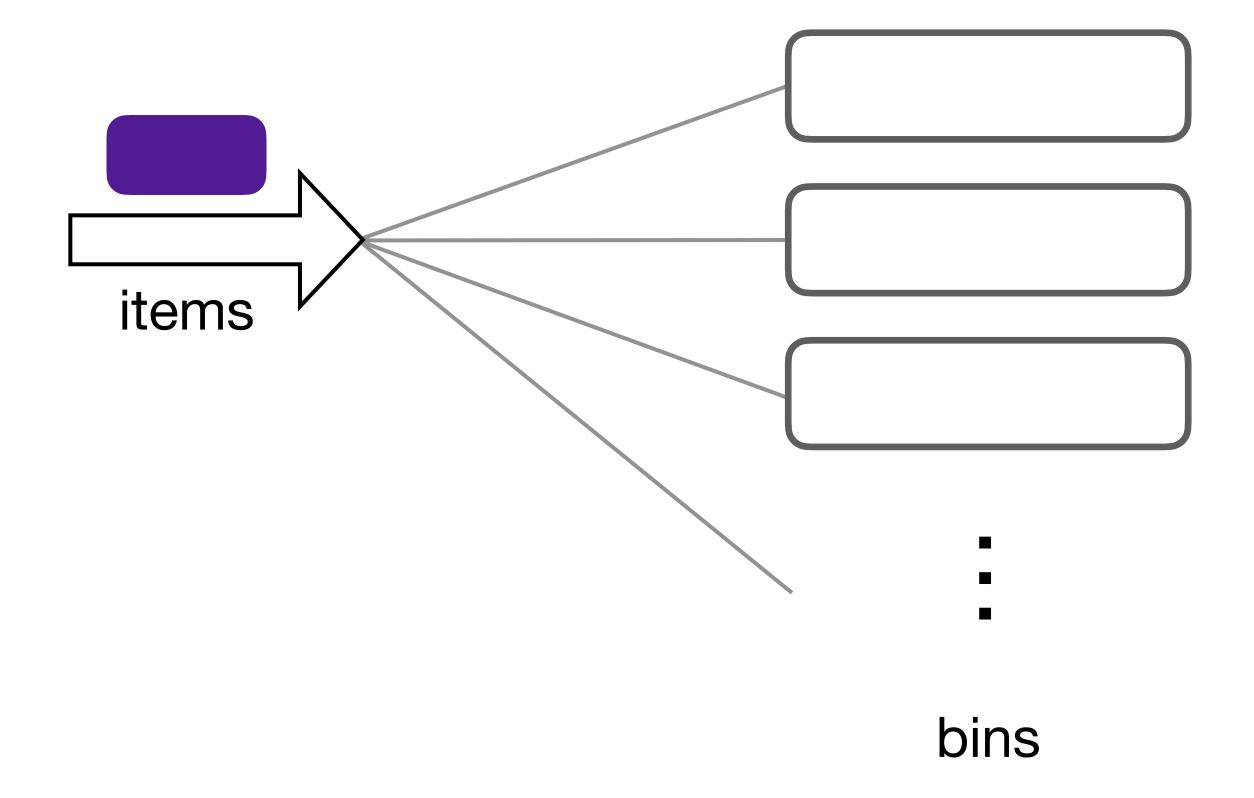
Weina Wang
Carnegie Mellon University

# The problem



items

bins

# The problem

- Each arriving item needs to be assigned to a bin



items

bins

# The problem

- Each arriving item needs to be assigned to a bin

- Infinite # bins



items

bins

# The problem

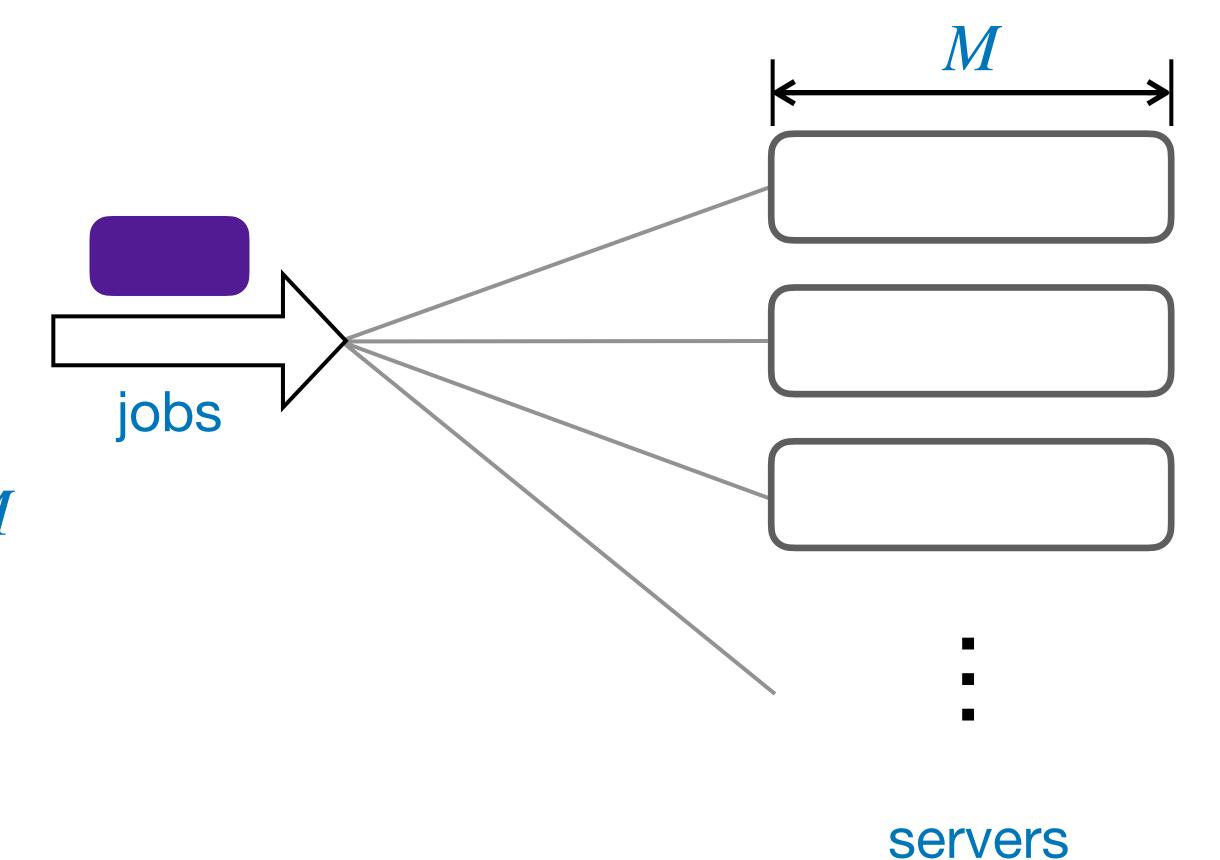- Each arriving item needs to be assigned to a bin

- Infinite # bins

- Each bin has a capacity $M$

# The problem

- Each arriving  job  needs to be assigned to a bin

- Infinite # bins

- Each bin has a capacity $M$



jobs

$M$

bins

# The problem

- Each arriving  job  needs to be assigned to a  server

- Infinite #  servers

- Each  server  has a resource capacity $M$



$M$

jobs

servers

# The problem

- Each arriving  job  needs to be assigned to a  server

- Infinite # servers

- Each server has a resource capacity $M$

Traditional job model:

- Each job has a fixed resource requirement

- Each job departs after a random time

jobs

$M$

servers

# The problem

- Each arriving  job  needs to be assigned to a  server

- Infinite # servers

- Each server has a resource capacity $M$

Traditional job model:

- Each job has a fixed resource requirement

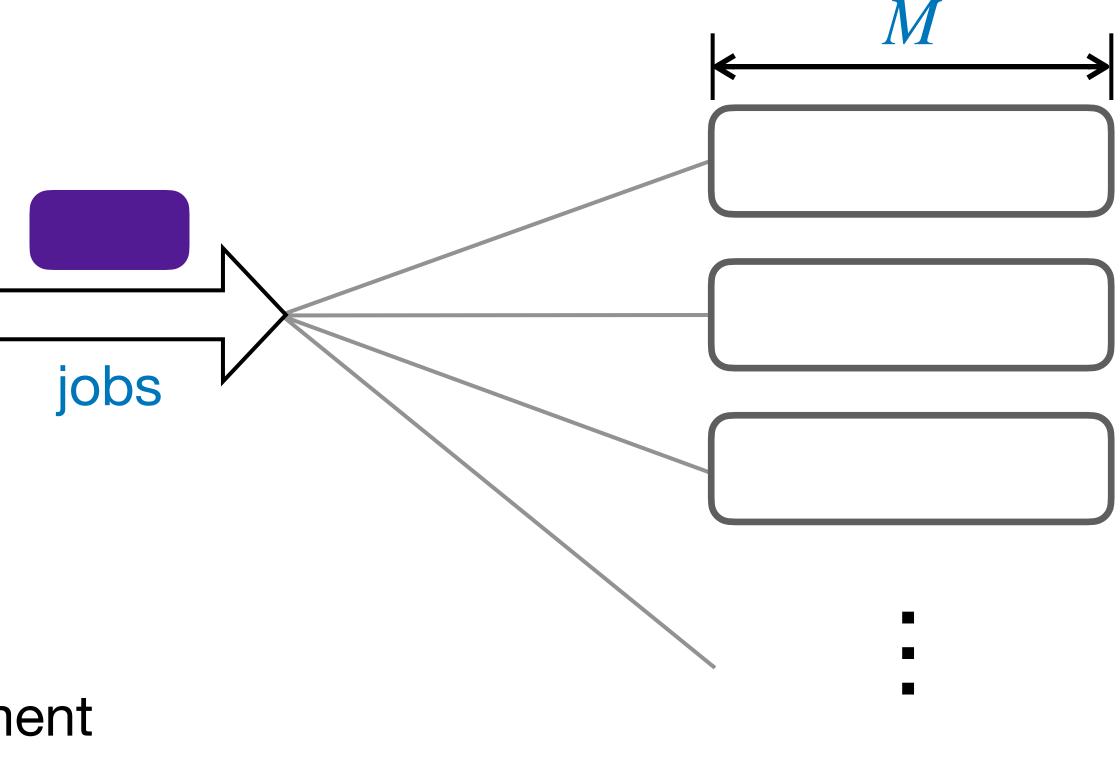- Each job departs after a random time



jobs

$M$

servers

# The problem

- Each arriving  job  needs to be assigned to a  server

- Infinite # servers

- Each server has a resource capacity $M$

Traditional job model:

- Each job has a fixed resource requirement

- Each job departs after a random time

jobs

servers

$M$

# The problem

- Each arriving  job  needs to be assigned to a  server

- Infinite #  servers

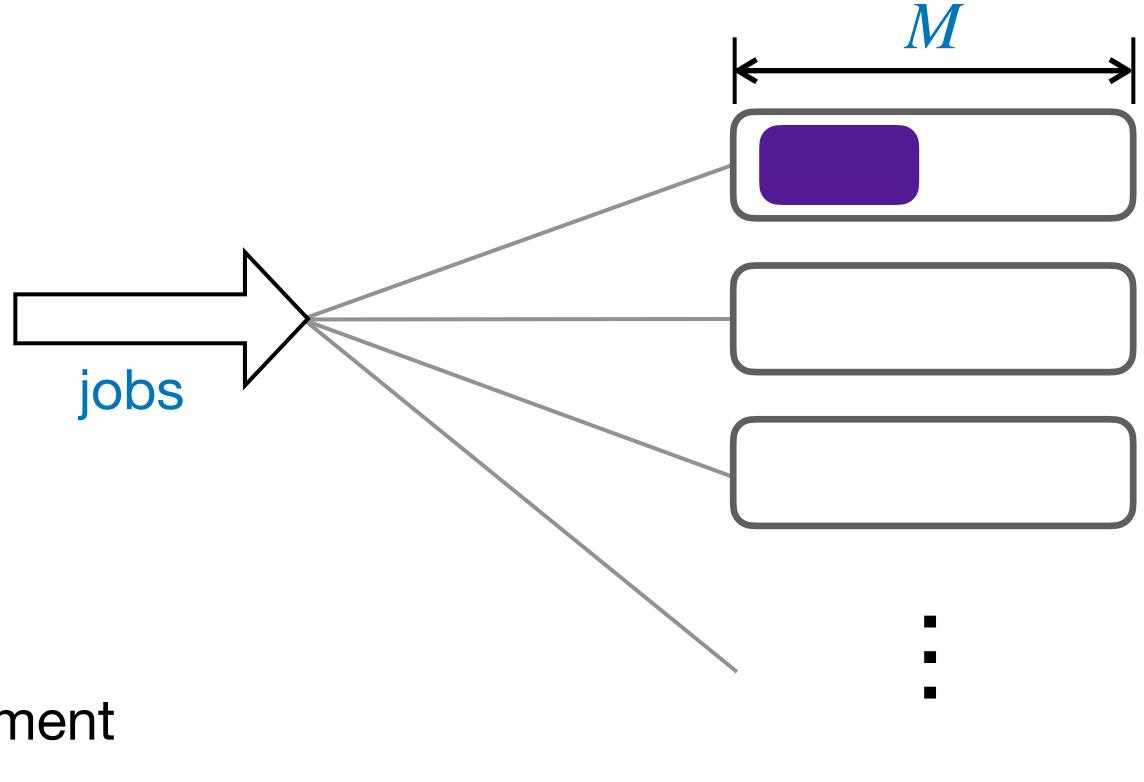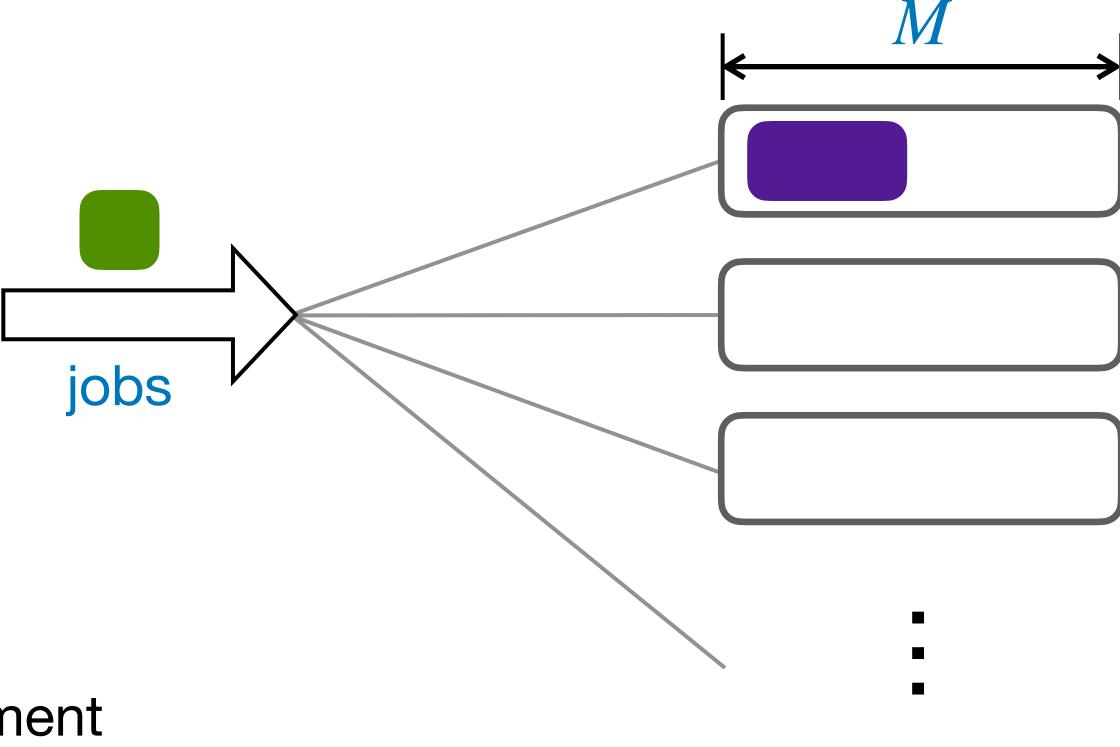- Each  server  has a resource capacity $M$

Traditional job model:

- Each job has a fixed resource requirement

- Each job departs after a random time

jobs

$M$

servers

# The problem



- Each arriving  job  needs to be assigned to a  server

- Infinite # servers

- Each server has a resource capacity $M$

Traditional job model:

- Each job has a fixed resource requirement

- Each job departs after a random time

# The problem



- Each arriving  job  needs to be assigned to a  server

- Infinite # servers

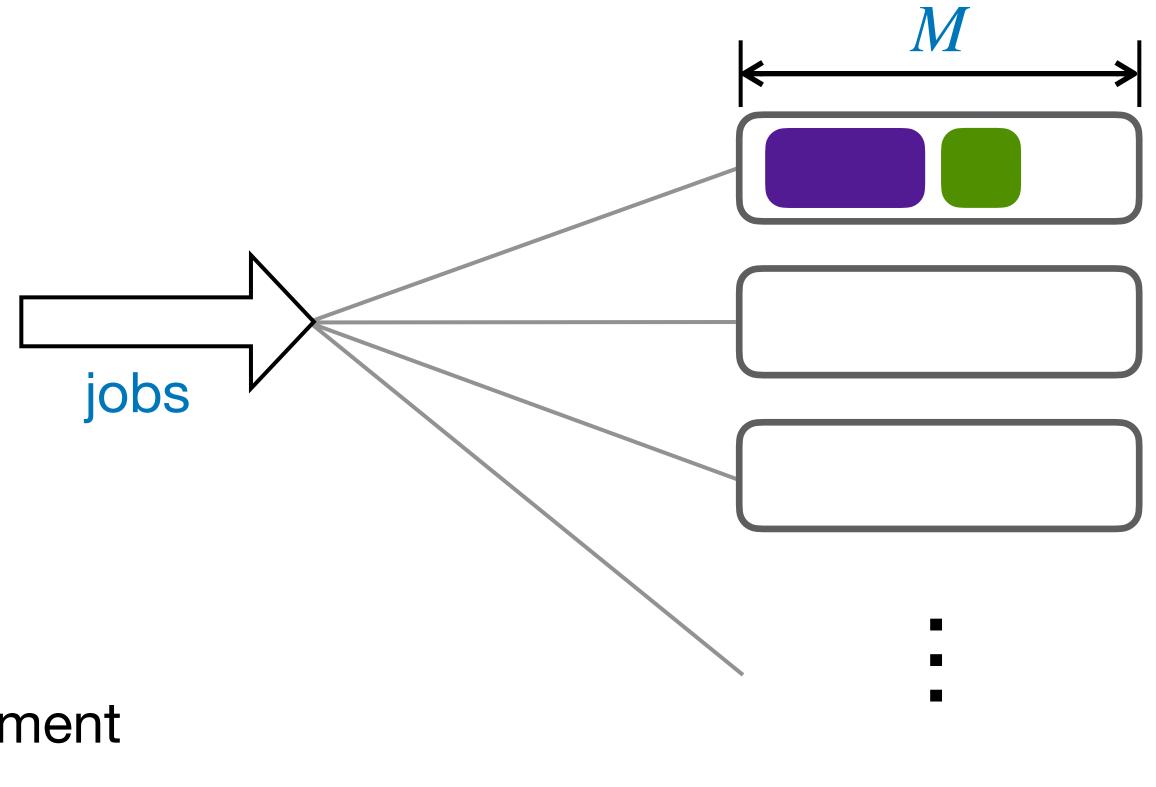- Each server has a resource capacity $M$

Traditional job model:

- Each job has a fixed resource requirement

- Each job departs after a random time

# The problem

- Each arriving  job  needs to be assigned to a  server

- Infinite # servers

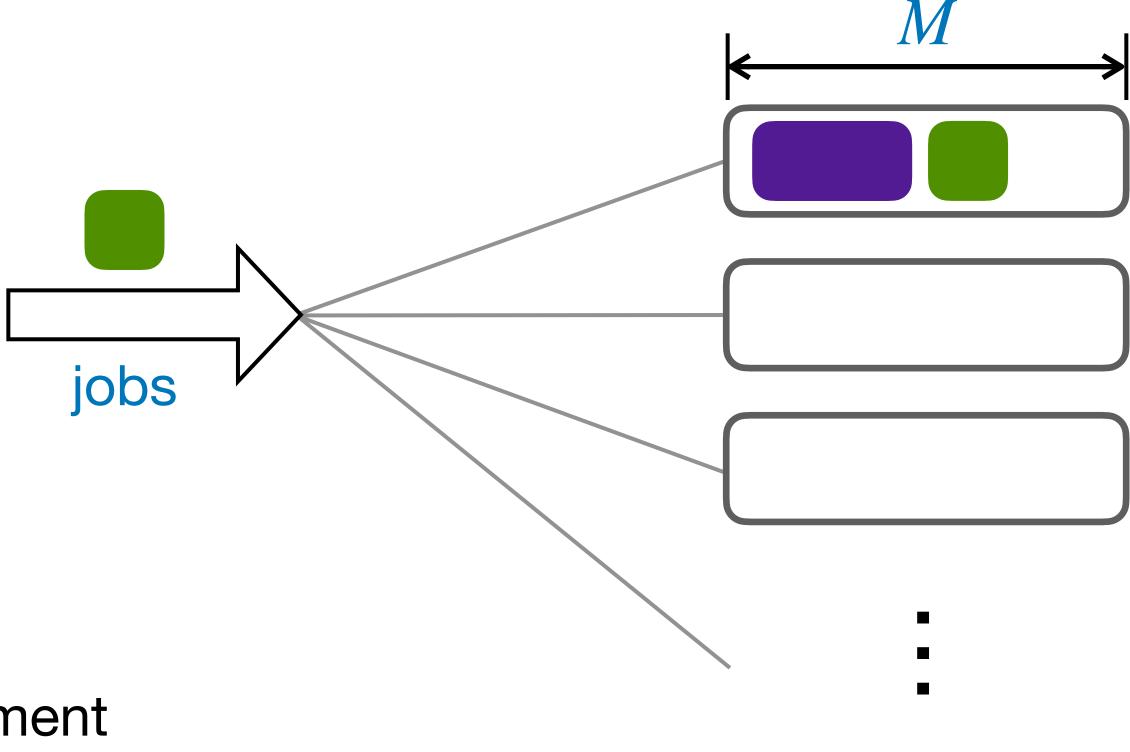- Each server has a resource capacity $M$

Traditional job model:

- Each job has a fixed resource requirement

- Each job departs after a random time

jobs

servers

# The problem

- Each arriving  job  needs to be assigned to a  server

- Infinite # servers

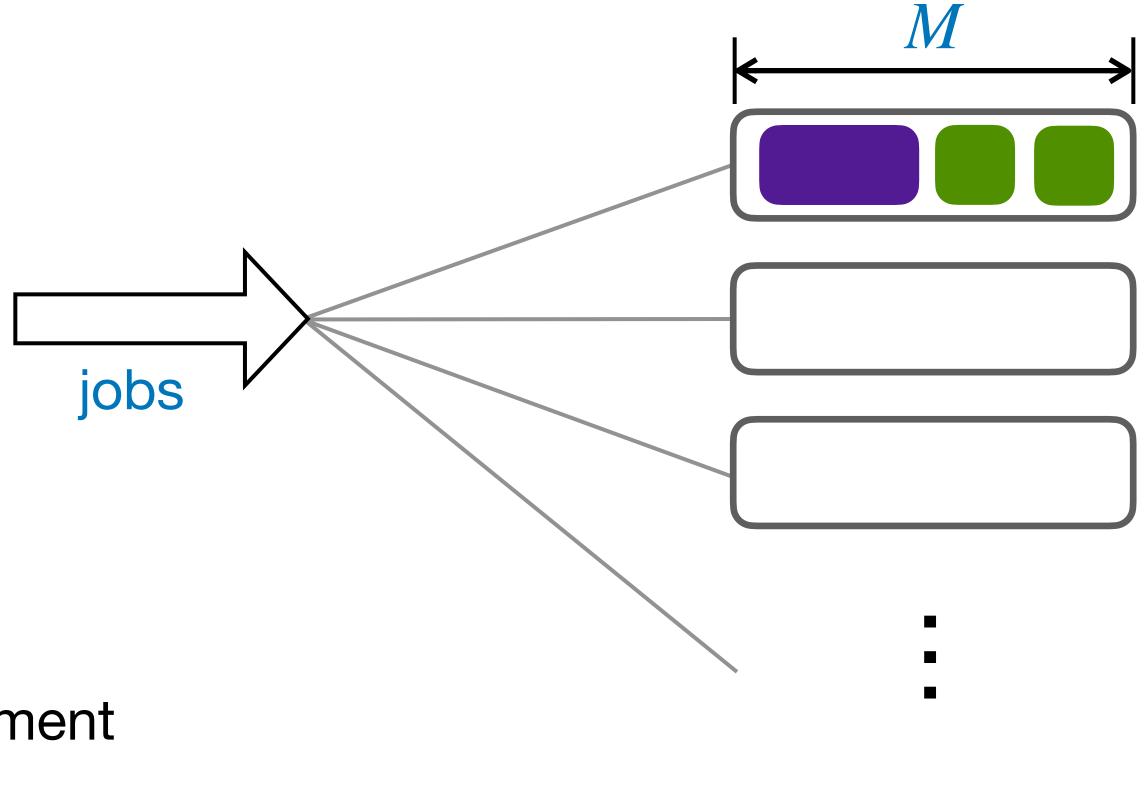- Each server has a resource capacity $M$

Traditional job model:

- Each job has a fixed resource requirement

- Each job departs after a random time

jobs

servers

$M$

# The problem



- Each arriving job needs to be assigned to a server

- Infinite # servers

- Each server has a resource capacity $M$

Traditional job model:

- Each job has a fixed resource requirement

- Each job departs after a random time

# The problem



- Each arriving  job  needs to be assigned to a  server

- Infinite # servers

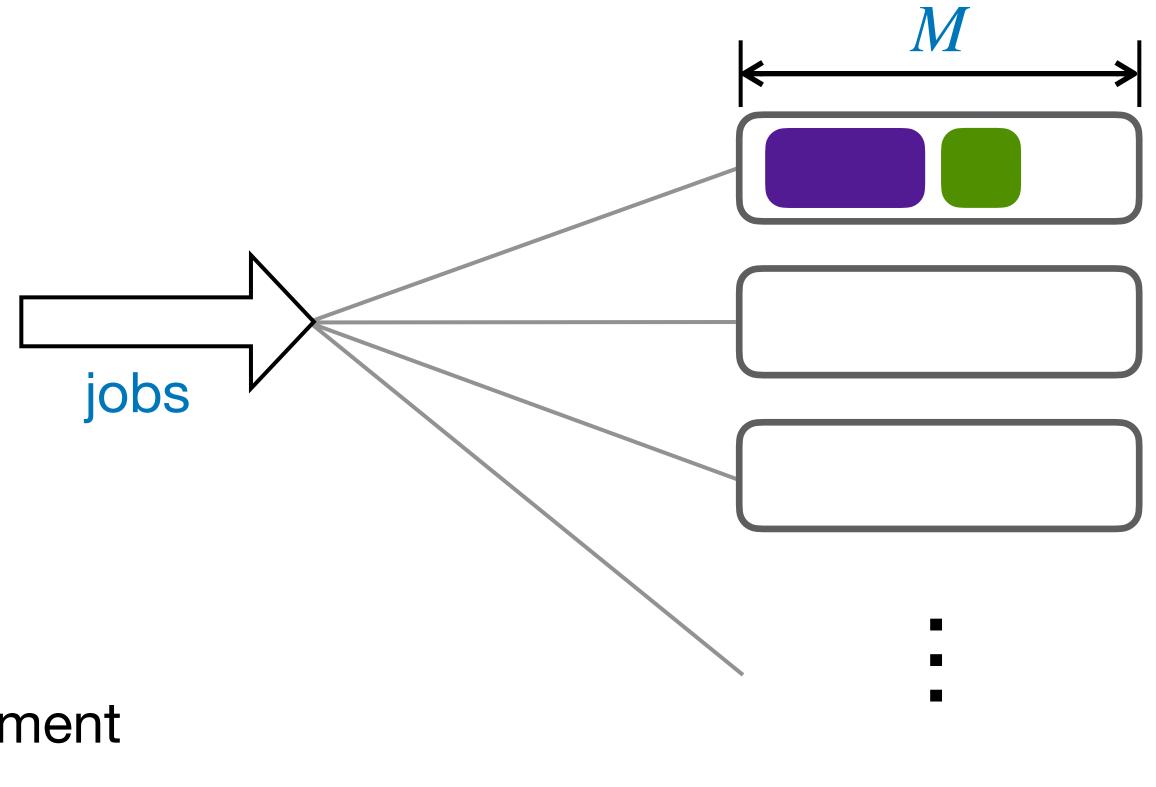- Each server has a resource capacity $M$

Traditional job model:

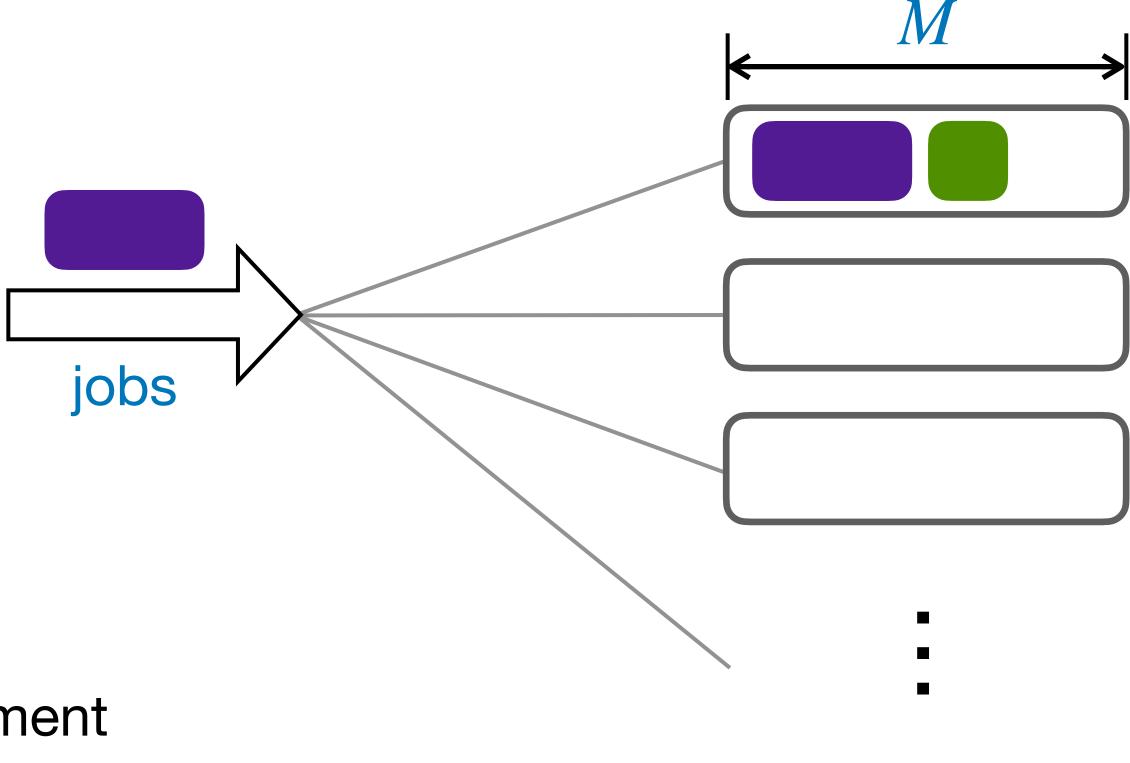- Each job has a fixed resource requirement

- Each job departs after a random time

**Goal:** minimize $\mathbf{E}$ [# active servers]

job assigning policy

# The problem



- Each arriving job needs to be assigned to a server

- Infinite # servers

- Each server has a resource capacity $M$

Traditional job model:

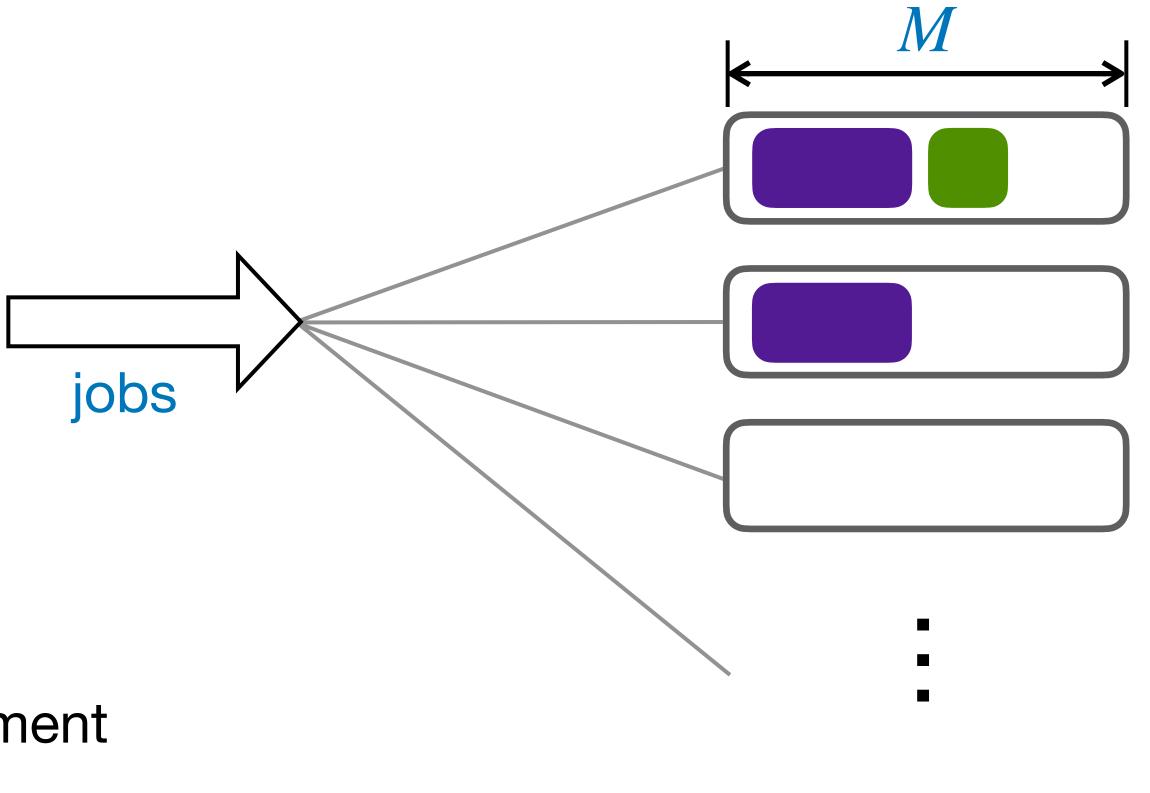- Each job has a fixed resource requirement
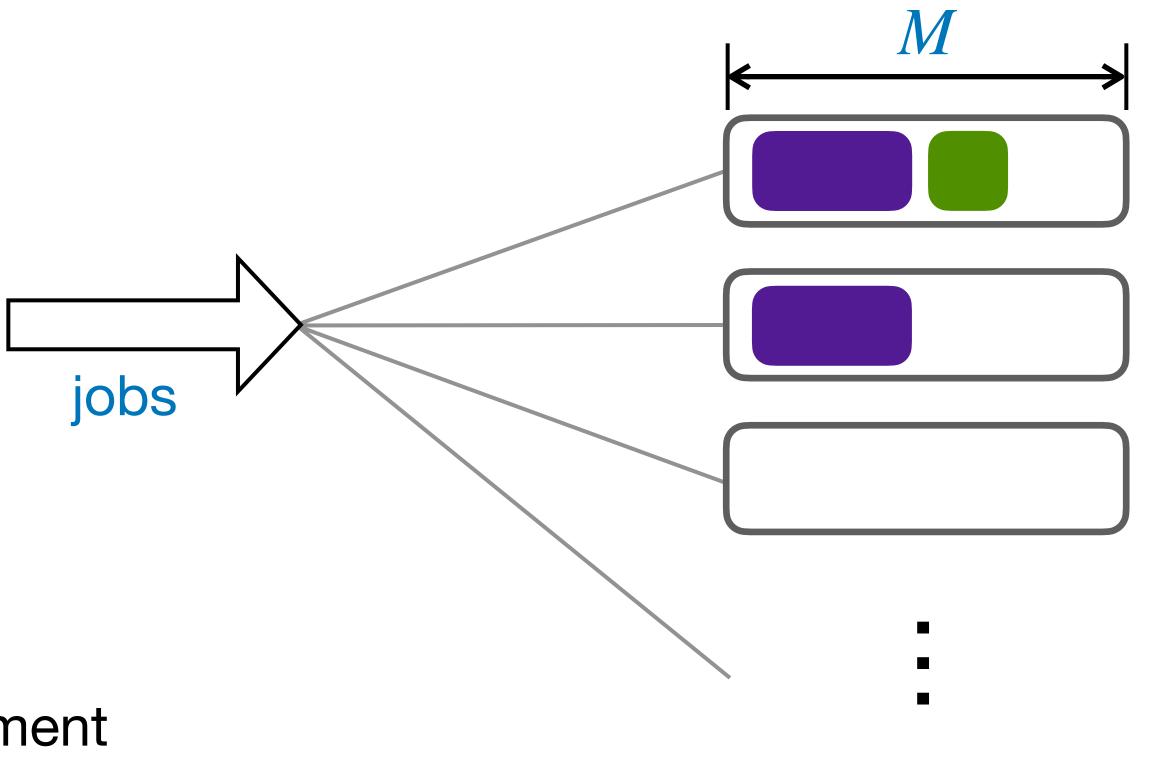
- Each job departs after a random time

**Goal:** minimize $\mathbf{E}$ [# active servers]

job assigning policy

Prior work: algorithms with asymptotic optimality

[Stolyar and Zhong 2013, 2015], [Stolyar 2017], [Stolyar and Zhong 2021], ...

# The problem



- Each arriving  job  needs to be assigned to a  server

- Infinite # servers

- Each server has a resource capacity $M$

  A new job model:

- Each job has a fixed resource requirement

- Each job departs after a random time

**Goal:**   $\underset{\text{job assigning policy}}{\text{minimize}}$  **E** [# active servers]

jobs

servers

# The problem



- Each arriving  job  needs to be assigned to a  server

- Infinite # servers

- Each server has a resource capacity $M$

A new job model: → time-varying

- Each job has a ~~fixed~~ resource requirement

- Each job departs after a random time

**Goal:** $\underset{\text{job assigning policy}}{\text{minimize}}$ **E** [# active servers]
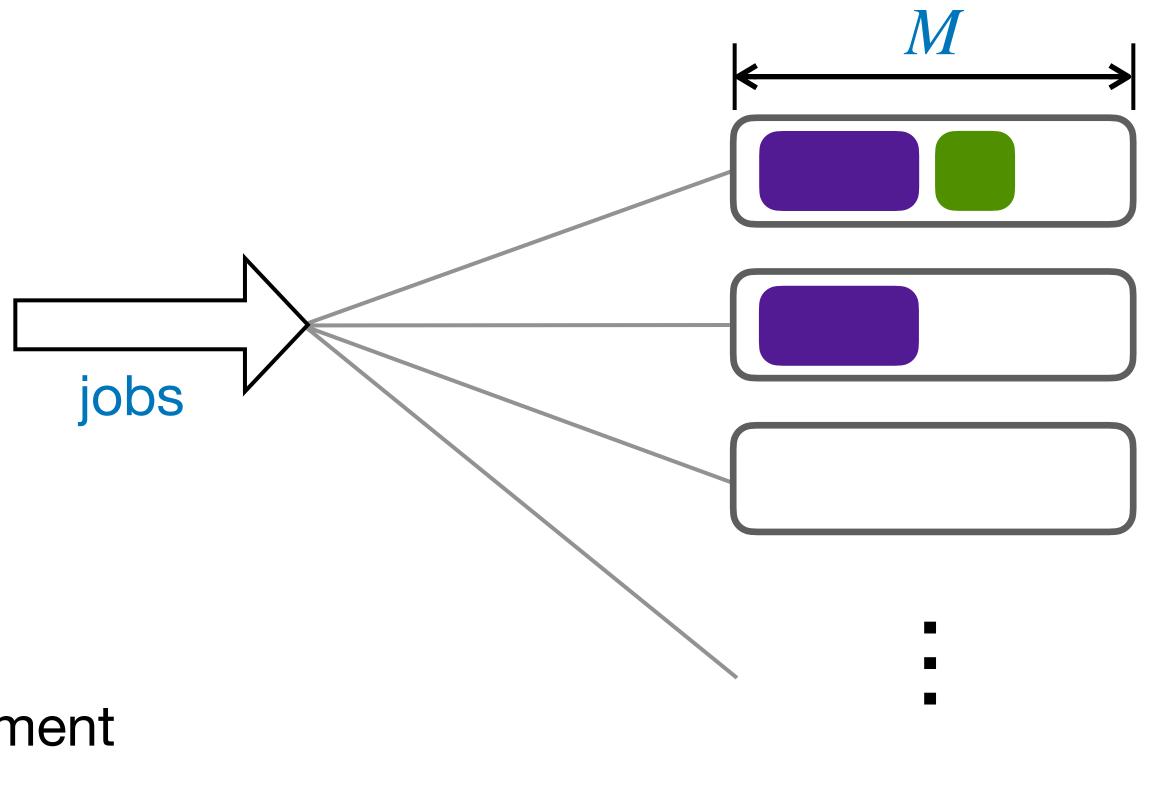
# The problem

$M$

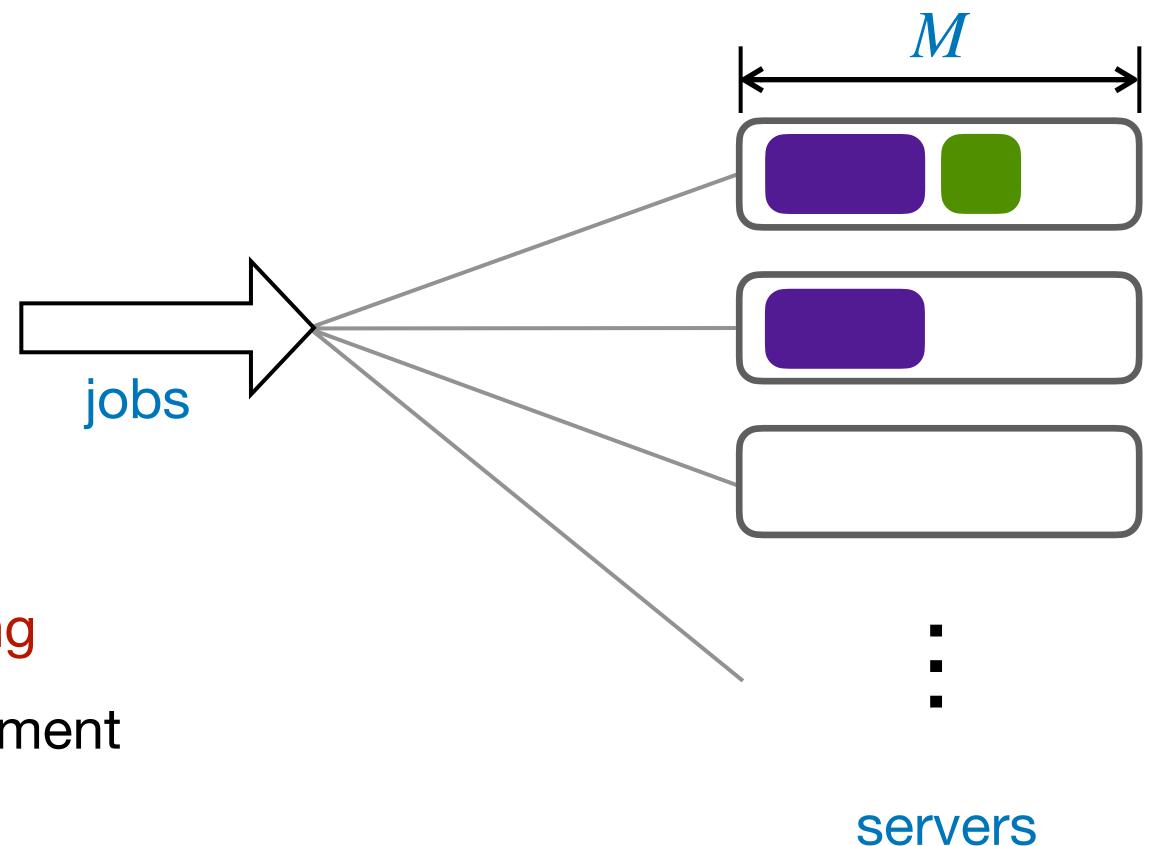- Each arriving  job  needs to be assigned to a  server

- Infinite # servers

jobs

- Each server has a resource capacity $M$

A new job model:          time-varying

- Each job has a ~~fixed~~ resource requirement

- Each job departs after a random time

servers

**Goal:**     minimize     **E** [# active servers]
job assigning policy

subject to     **cost** (resource contention) ≤ budget

# Why does time-varying matter?



resource requirement of a job

10 CPUs

1 CPU

time

# Why does time-varying matter?



- Reserve resources based on peak requirement

# Why does time-varying matter?



- Reserve resources based on peak requirement
  ➡ low resource utilization on a server

# Why does time-varying matter?



resource requirement of a job

- Reserve resources based on peak requirement

  ➡ low resource utilization on a server

  ➡ larger # active servers

# Why does time-varying matter?



resource requirement
of a job

10 CPUs

1 CPU

time

- Reserve resources based on peak requirement

  ➡ low resource utilization on a server

  ➡ larger # active servers

- Overcommit resources on a server

# Why does time-varying matter?



resource requirement of a job

- 10 CPUs
- 1 CPU
- time

- Reserve resources based on peak requirement
  - ➡ low resource utilization on a server
  - ➡ larger # active servers
- Overcommit resources on a server
  - ➡ possible resource contention

# Why does time-varying matter?



resource requirement
of a job

10 CPUs

1 CPU

time

- Reserve resources based on peak requirement
  - ➡ low resource utilization on a server
  - ➡ larger # active servers
- Overcommit resources on a server
  - ➡ possible resource contention

Our formulation captures:

**utilization** ⟷ **resource contention**

# More details on the job model



Example MC

# More details on the job model

- Resource requirement of a job evolves over time following a Markov chain



Phase $L$

Phase $H$

Phase $\perp$

$\mu_{LH}$

$\mu_{HL}$

$\mu_{L\perp}$

$\mu_{H\perp}$

(completion)

Example MC

# More details on the job model

- Resource requirement of a job evolves over time following a Markov chain

- Initial job type follows an initial distribution



Example MC

# More details on the job model

- Resource requirement of a job evolves over time following a Markov chain

- Initial job type follows an initial distribution

- MCs of jobs are independent of each other, and they are exogenous (not affected by resource contention)



Example MC

# More details on the job model

- Resource requirement of a job evolves over time following a Markov chain

- Initial job type follows an initial distribution

- MCs of jobs are independent of each other, and they are exogenous (not affected by resource contention)

- Jobs arrive following a Poisson process



$\mu_{LH}$

Phase $L$

$\mu_{HL}$

Phase $H$

$\mu_{L\perp}$

$\mu_{H\perp}$

Phase $\perp$

(completion)

Example MC

jobs

servers

state: # jobs of each type
on each server



jobs

servers

state space is large!

state: # jobs of each type
on each server

jobs

servers

# Reducing dimensionality

state space is large!

state: # jobs of each type
on each server

jobs

servers

# Reducing dimensionality

state: # jobs of each type
on each server

- Server-by-server evaluation:



jobs

servers

# Reducing dimensionality

state: # jobs of each type
on each server

- Server-by-server evaluation:
  - How to evaluate each server?

jobs

servers

# Reducing dimensionality

state space is large!

state: # jobs of each type
on each server

- Server-by-server evaluation:
  - How to evaluate each server?
  - How to relate to E[# active servers]?



jobs

servers

# A policy-conversion framework

Policies in the
$\infty$-server system

⟷

Policies in a
single-server system

# A policy-conversion framework



$$\sigma \dashleftarrow \overline{\sigma}$$

Policies in the $\infty$-server system  ←  achievability  ←  Policies in a single-server system

# A policy-conversion framework

- Use $\overline{\sigma}$ to tell how to evaluate each server

- Performance of $\sigma$ is related to properties of $\overline{\sigma}$

$$\sigma \longleftarrow \overline{\sigma}$$

achievability

Policies in the
$\infty$-server system

Policies in a
single-server system

# A policy-conversion framework

- Use $\overline{\sigma}$ to tell how to evaluate each server

- Performance of $\sigma$ is related to properties of $\overline{\sigma}$



$$\sigma \leftarrowtail \overline{\sigma}$$

achievability

Policies in the
$\infty$-server system

Policies in a
single-server system

converse

$$\sigma \mapsto \overline{\sigma}$$

# A policy-conversion framework

- Use $\overline{\sigma}$ to tell how to evaluate each server

- Performance of $\sigma$ is related to properties of $\overline{\sigma}$



$$\sigma \longleftarrow\!\!\mid \overline{\sigma}$$

achievability

Policies in the
$\infty$-server system

Policies in a
single-server system

converse

$$\sigma \longmapsto \overline{\sigma}$$

- Allows us to obtain lower bound on
  E[# active servers]

# A policy-conversion framework

Policies in the $\infty$-server system ⟷ Policies in a single-server system

Single-server system

# A policy-conversion framework

Policies in the
$\infty$-server system

$\longleftrightarrow$

Policies in a
single-server system

Single-server system

Infinite supply of
jobs of all types

requests

jobs

# A policy-conversion framework

Policies in the
∞-server system

⟷

Policies in a
single-server system

<u>Single-server system</u>

Infinite supply of
jobs of all types

requests

jobs

A policy decides when to request what types of jobs to:
maximize    throughput
subject to   **cost** (resource contention) ≤ budget

Policies in the $\infty$-server system $\quad\sigma \longleftarrow \overline{\sigma}\quad$ Policies in a single-server system

- Arrival rates: $r \cdot \left( \lambda_L, \lambda_H \right)$

Policies in the $\infty$-server system

$\sigma \dashv \overline{\sigma}$

Policies in a single-server system

- Arrival rates: $r \cdot (\lambda_L, \lambda_H)$

Phase $L$

$\mu_{LH}$

$\mu_{HL}$

Phase $H$

$\mu_{L\perp}$

$\mu_{H\perp}$

Phase $\perp$

(completion)

Policies in the $\infty$-server system $\sigma \dashv \overline{\sigma}$ Policies in a single-server system

- Arrival rates: $r \cdot (\lambda_L, \lambda_H)$

- Asymptotic regime: $r \to +\infty$

Phase $L$ $\xrightarrow{\mu_{LH}}$ Phase $H$

Phase $L$ $\xleftarrow{\mu_{HL}}$ Phase $H$

$\mu_{L\perp}$   $\mu_{H\perp}$

Phase $\perp$

(completion)

Policies in the $\infty$-server system

$\sigma \dashleftarrow \overline{\sigma}$

Policies in a single-server system

- Arrival rates: $r \cdot \left( \lambda_L, \lambda_H \right)$

- Asymptotic regime: $r \to +\infty$

$r \to \infty$

$\left( \lambda_L, \lambda_H \right)$

Phase $L$

$\mu_{LH}$

$\mu_{HL}$

Phase $H$

$\mu_{L\perp}$

$\mu_{H\perp}$

Phase $\perp$

(completion)

Policies in the $\infty$-server system $\quad\sigma \longleftarrow \overline{\sigma}\quad$ Policies in a single-server system

- Arrival rates: $r \cdot \left(\lambda_L, \lambda_H\right)$

- Asymptotic regime: $r \to +\infty$

Policies in the $\infty$-server system $\sigma \longleftarrow \overline{\sigma}$ Policies in a single-server system

- Arrival rates: $r \cdot \left( \lambda_L, \lambda_H \right)$

- Asymptotic regime: $r \rightarrow +\infty$

Policy $\overline{\sigma}$

**throughput** $\cdot \overline{N} = r \cdot \left( \lambda_L, \lambda_H \right)$

**cost** (resource contention) $\leq$ budget

Policies in the $\infty$-server system

$\sigma \longleftarrow \overline{\sigma}$

Policies in a single-server system

- Arrival rates: $r \cdot \left( \lambda_L, \lambda_H \right)$

- Asymptotic regime: $r \to +\infty$

Policy $\overline{\sigma}$

**throughput** $\cdot \overline{N} = r \cdot \left( \lambda_L, \lambda_H \right)$

**cost** (resource contention) $\leq$ budget

convert

Policy $\sigma$

$\mathbf{E}$ [# active servers] $\leq \left( 1 + O \left( r^{-0.5} \right) \right) \cdot \overline{N}$

**cost** (resource contention) $\leq \left( 1 + O \left( r^{-0.5} \right) \right) \cdot$ budget

Policies in the $\infty$-server system

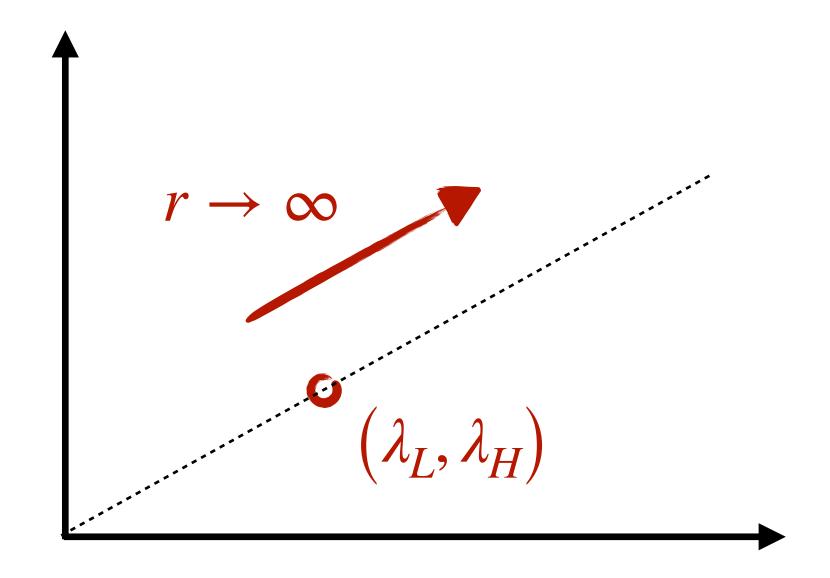$$\sigma \dashv \overline{\sigma}$$

Policies in a single-server system

- Arrival rates: $r \cdot \left( \lambda_L, \lambda_H \right)$

- Asymptotic regime: $r \rightarrow +\infty$

Policy $\overline{\sigma}$

**throughput** $\cdot \overline{N} = r \cdot \left( \lambda_L, \lambda_H \right)$

**cost** (resource contention) $\leq$ budget

*convert*

Policy $\sigma$

$\mathbf{E}$ [# active servers] $\leq \left( 1 + O\left( r^{-0.5} \right) \right) \cdot \overline{N}$

**cost** (resource contention) $\leq \left( 1 + O\left( r^{-0.5} \right) \right) \cdot$ budget

**Main Result:** We design a policy for the original $\infty$-server system that is *asymptotically optimal*

# Policy conversion: single-server to $\infty$-server

Meta-algorithm: JOIN-THE-RECENTLY-REQUESTING-SERVER ($\bar{\sigma}$)



jobs

servers

# Policy conversion: single-server to $\infty$-server

Meta-algorithm: Join-the-Recently-Requesting-Server ($\overline{\sigma}$)

- For each server, run a single-server policy $\overline{\sigma}$

jobs

servers

# Policy conversion: single-server to $\infty$-server

Meta-algorithm: Join-the-Recently-Requesting-Server ($\bar{\sigma}$)

- For each server, run a single-server policy $\bar{\sigma}$



jobs

Single-server system
running policy $\bar{\sigma}$

servers

# Policy conversion: single-server to $\infty$-server

Meta-algorithm: JOIN-THE-RECENTLY-REQUESTING-SERVER ($\bar{\sigma}$)

- For each server, run a single-server policy $\bar{\sigma}$

- If $\bar{\sigma}$ requests a job of type $i$, generate a token of type $i$

# Policy conversion: single-server to $\infty$-server

Meta-algorithm: Join-the-Recently-Requesting-Server ($\bar{\sigma}$)

- For each server, run a single-server policy $\bar{\sigma}$

- If $\bar{\sigma}$ requests a job of type $i$, generate a token of type $i$

Request a job of type L

Single-server system running policy $\bar{\sigma}$

L

jobs

servers

# Policy conversion: single-server to $\infty$-server

Meta-algorithm: JOIN-THE-RECENTLY-REQUESTING-SERVER ($\bar{\sigma}$)

- For each server, run a single-server policy $\bar{\sigma}$

- If $\bar{\sigma}$ requests a job of type $i$, generate a token of type $i$



Request a job of type L

Single-server system running policy $\bar{\sigma}$

jobs
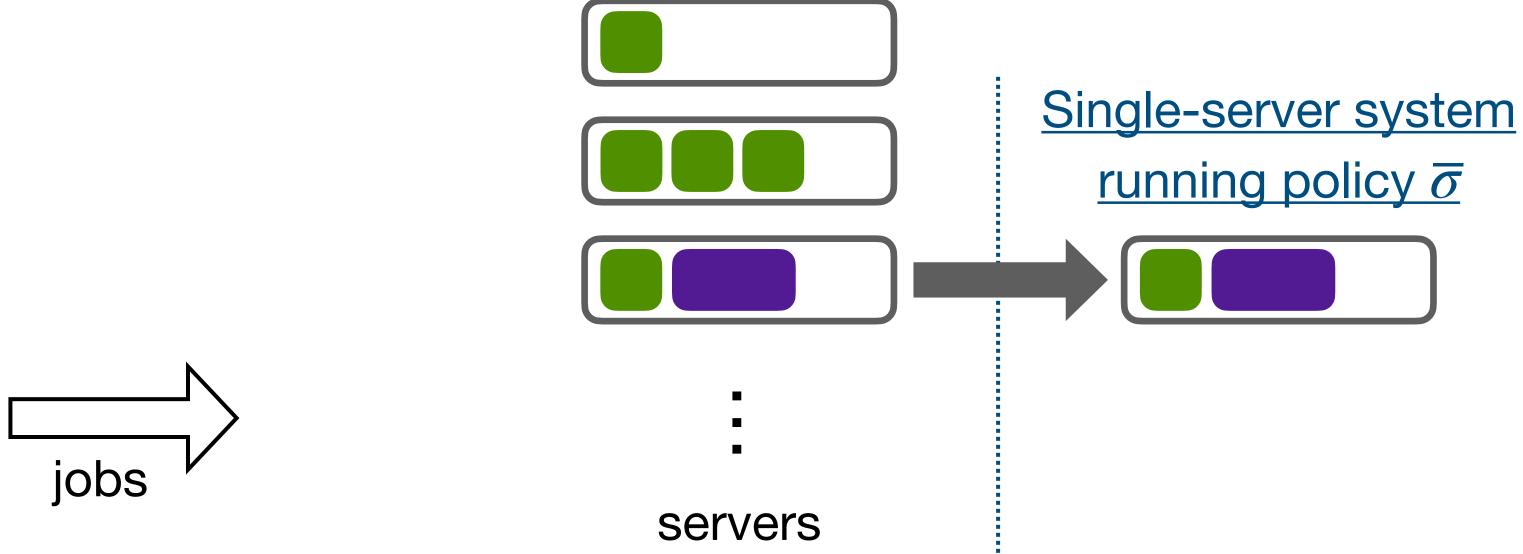
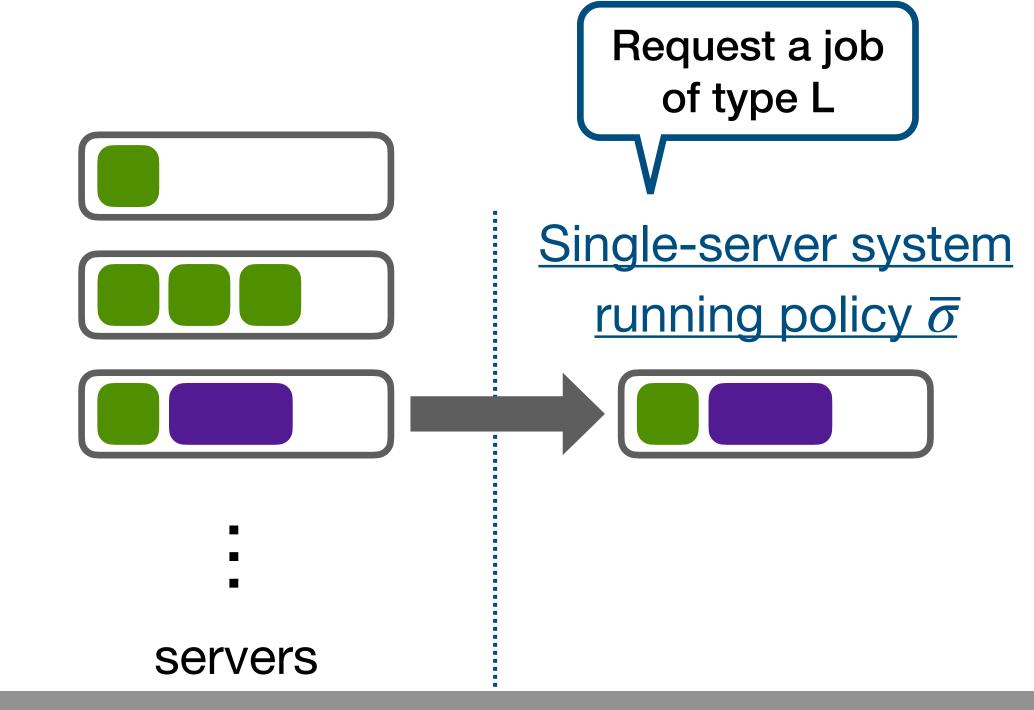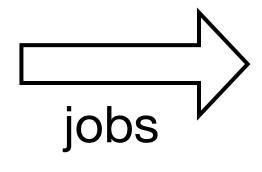servers

# Policy conversion: single-server to $\infty$-server

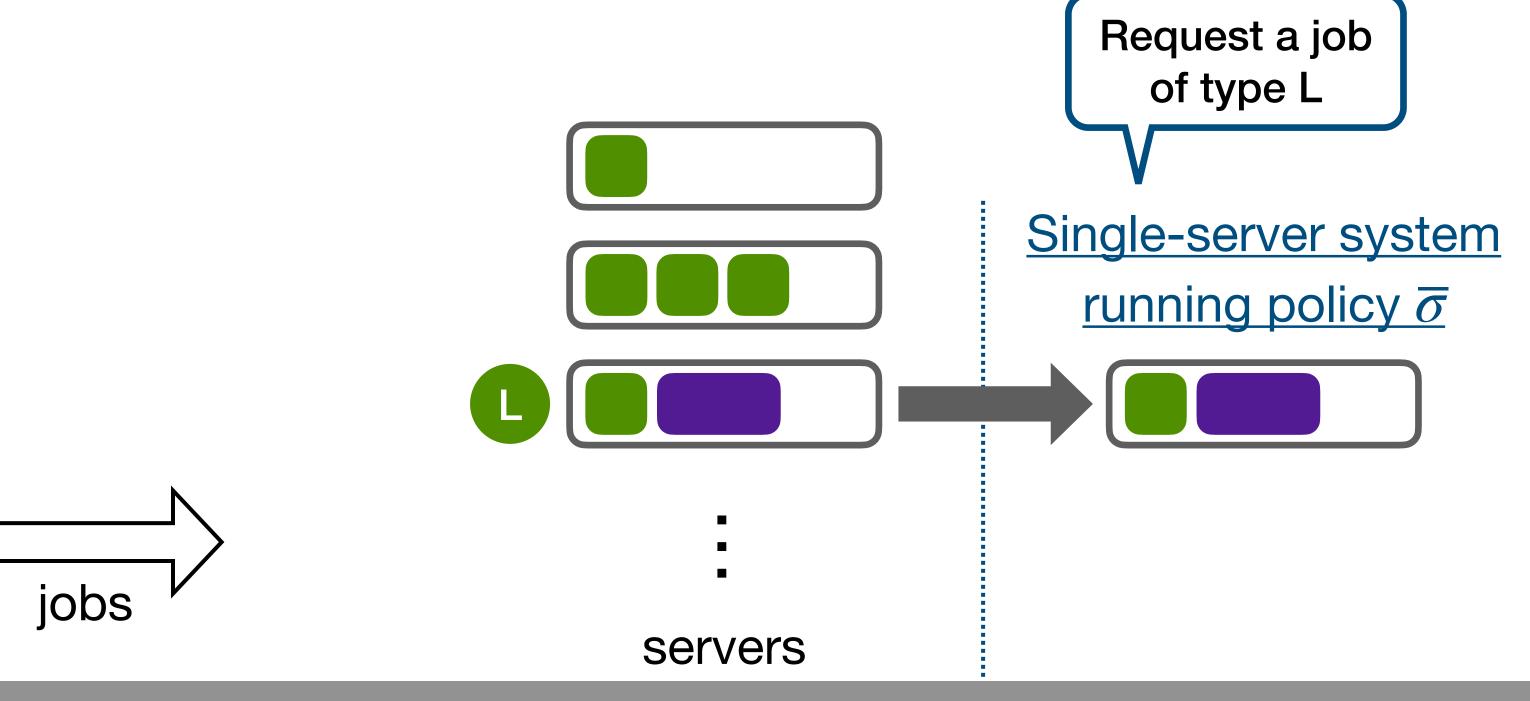Meta-algorithm: JOIN-THE-RECENTLY-REQUESTING-SERVER ($\overline{\sigma}$)

- For each server, run a single-server policy $\overline{\sigma}$

- If $\overline{\sigma}$ requests a job of type $i$, generate a token of type $i$

- When a job arrives, it checks tokens of its type and joins one uniformly at random



Request a job of type L

Single-server system running policy $\overline{\sigma}$

jobs

servers

# Policy conversion: single-server to $\infty$-server

Meta-algorithm: JOIN-THE-RECENTLY-REQUESTING-SERVER ($\bar{\sigma}$)

- For each server, run a single-server policy $\bar{\sigma}$

- If $\bar{\sigma}$ requests a job of type $i$, generate a token of type $i$

- When a job arrives, it checks tokens of its type and joins one uniformly at random



Request a job of type L

Single-server system

running policy $\bar{\sigma}$

jobs

servers

# Policy conversion: single-server to $\infty$-server

Meta-algorithm: JOIN-THE-RECENTLY-REQUESTING-SERVER ($\overline{\sigma}$)

- For each server, run a single-server policy $\overline{\sigma}$

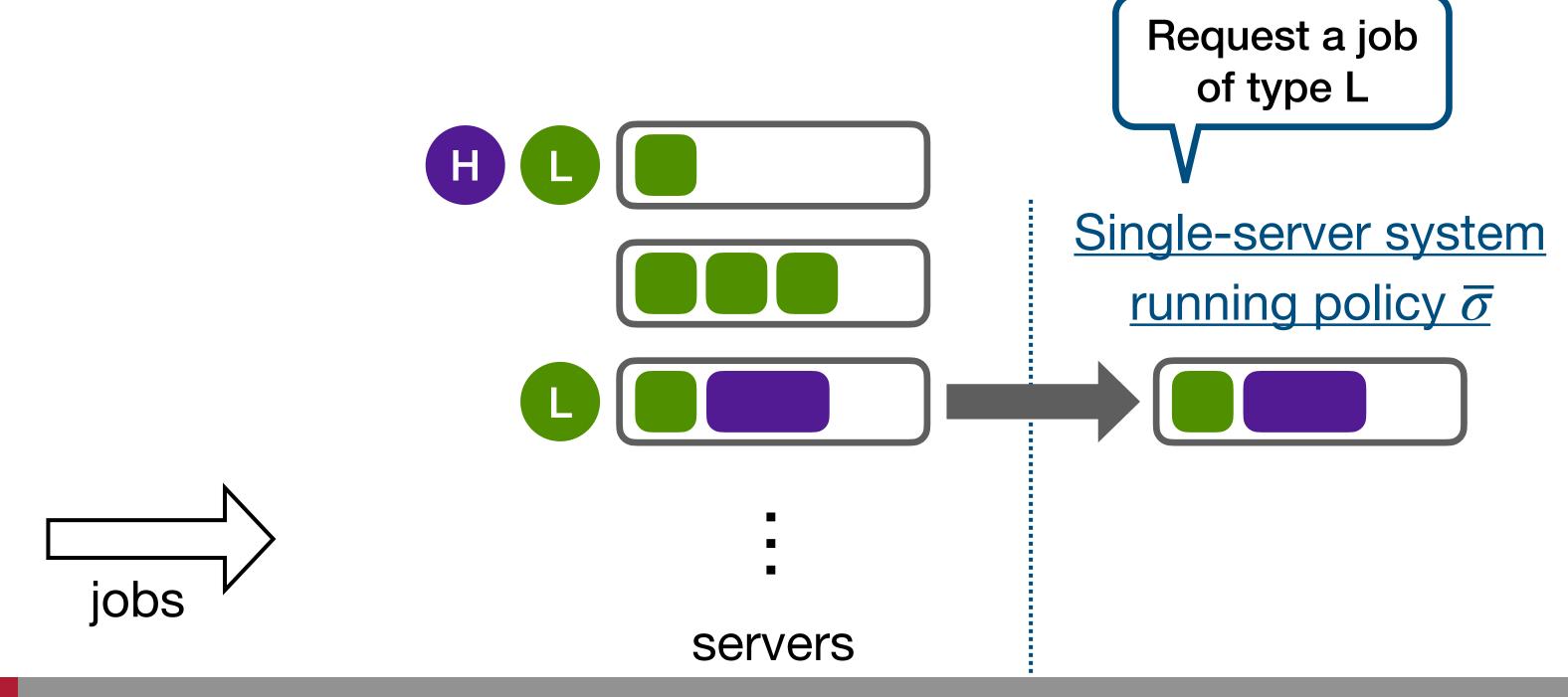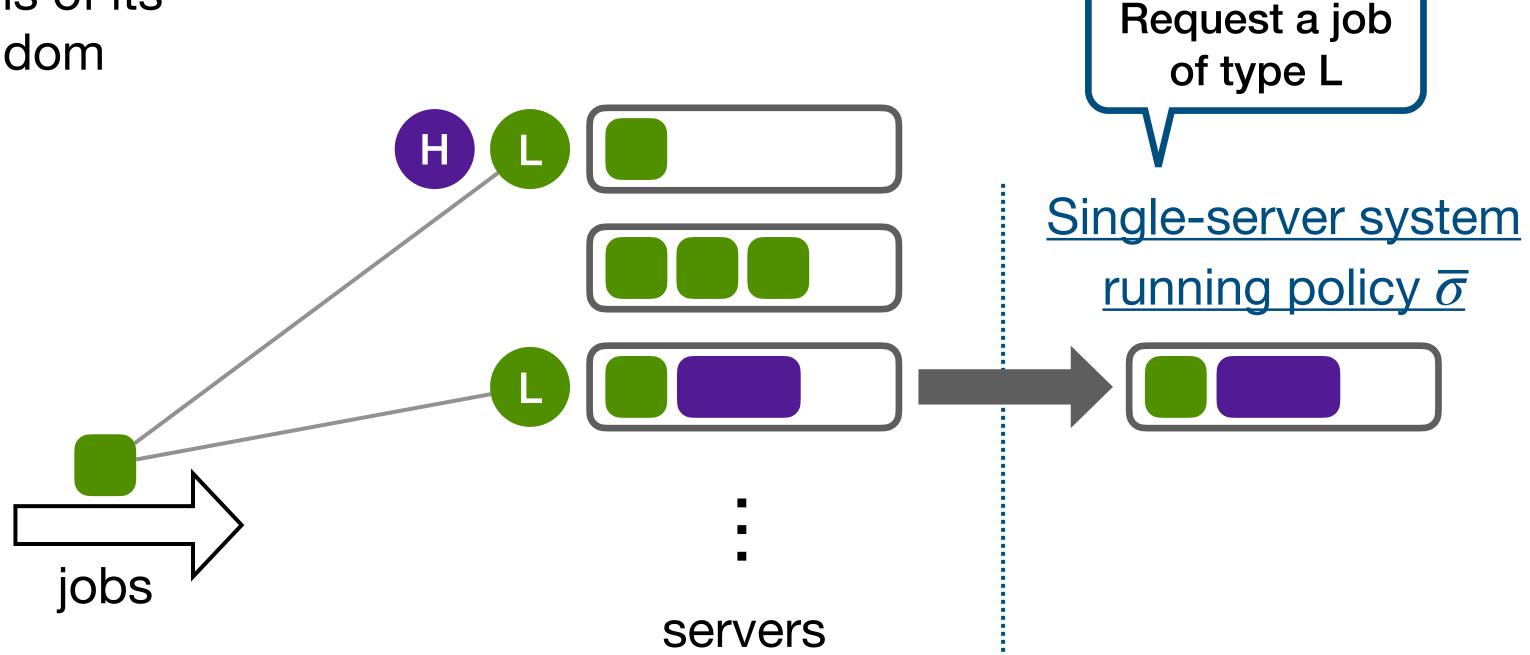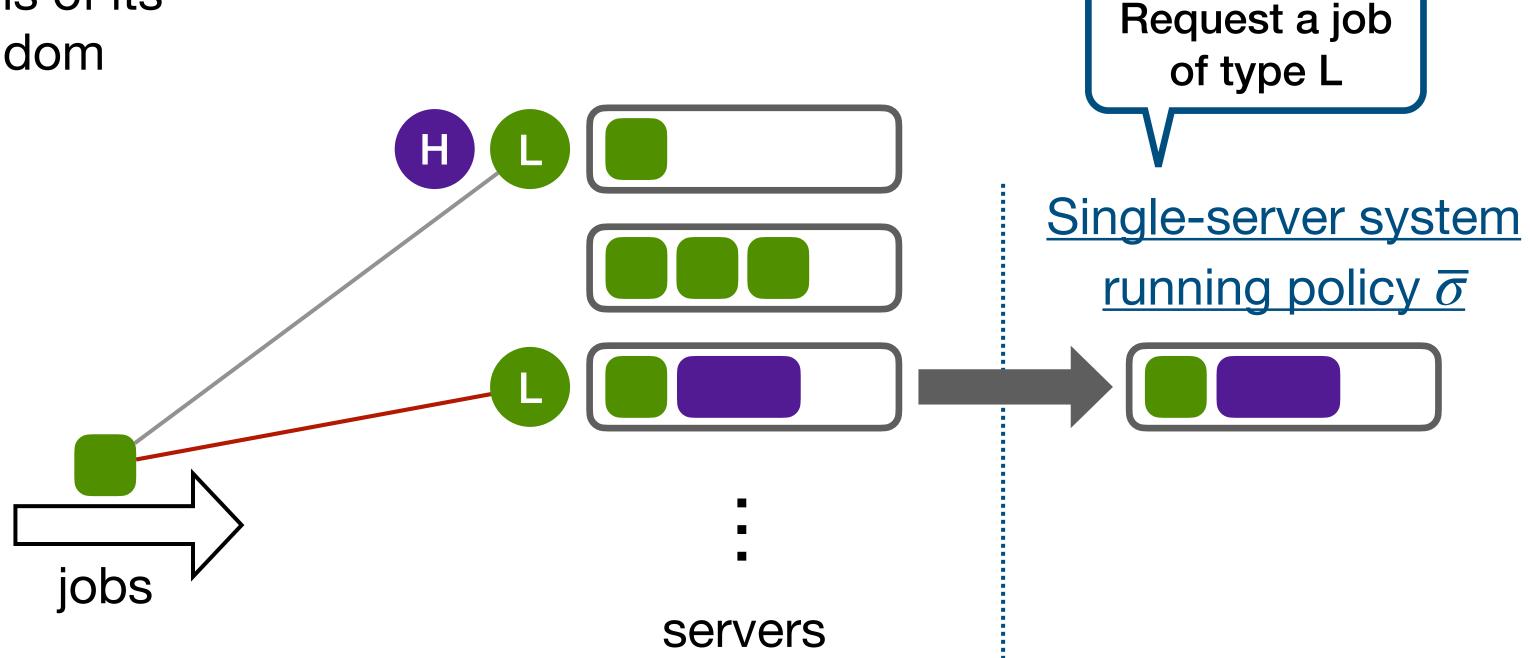- If $\overline{\sigma}$ requests a job of type $i$, generate a token of type $i$

- When a job arrives, it checks tokens of its type and joins one uniformly at random



Request a job of type L

Single-server system running policy $\overline{\sigma}$

jobs

servers

# Policy conversion: single-server to $\infty$-server

Meta-algorithm: Join-the-Recently-Requesting-Server ($\overline{\sigma}$)
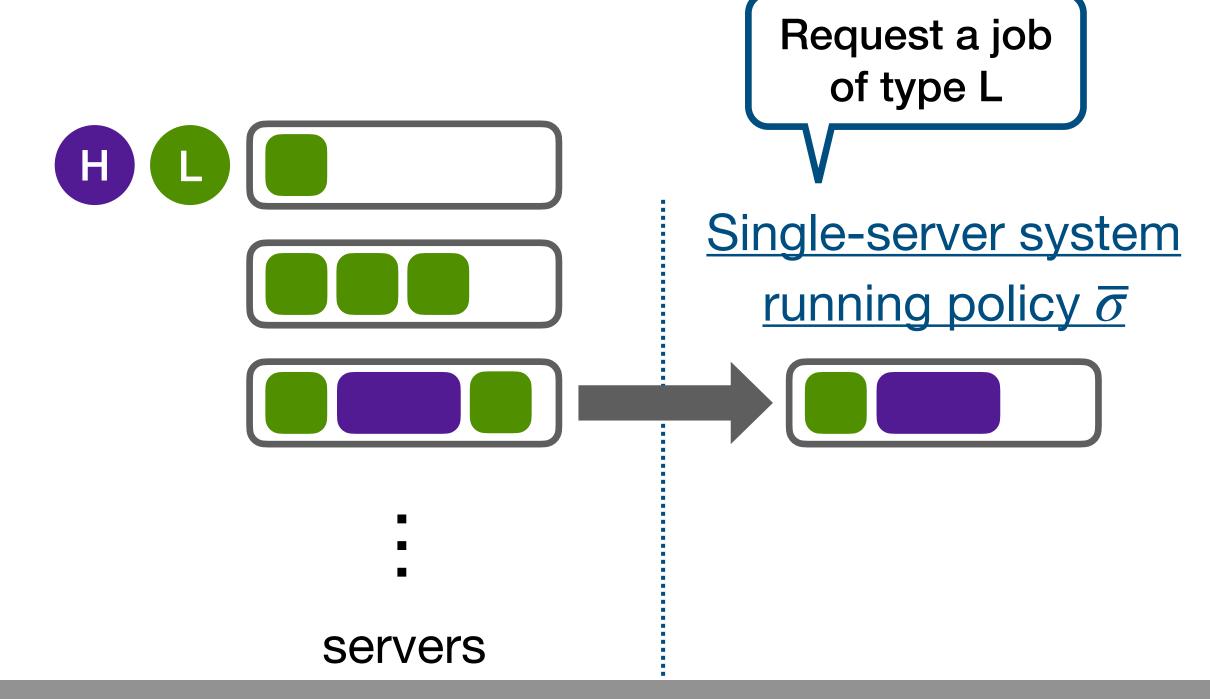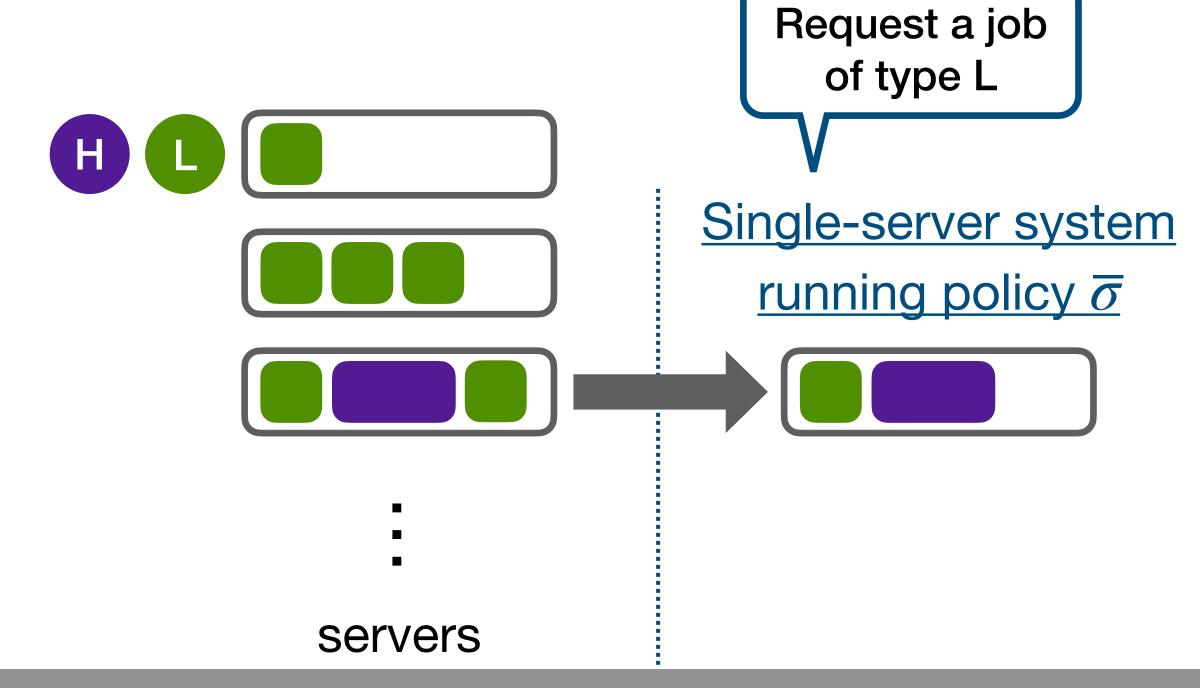
- For each server, run a single-server policy $\overline{\sigma}$

- If $\overline{\sigma}$ requests a job of type $i$, generate a token of type $i$

- When a job arrives, it checks tokens of its type and joins one uniformly at random

- If no tokens, go to an inactive server



Request a job of type L

Single-server system running policy $\overline{\sigma}$

jobs

servers

# Policy conversion: single-server to $\infty$-server

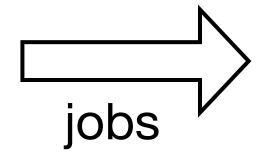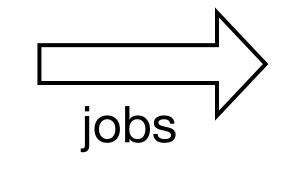Meta-algorithm: Join-the-Recently-Requesting-Server ($\overline{\sigma}$)

- For each server, run a single-server policy $\overline{\sigma}$

- If $\overline{\sigma}$ requests a job of type $i$, generate a token of type $i$

- When a job arrives, it checks tokens of its type and joins one uniformly at random

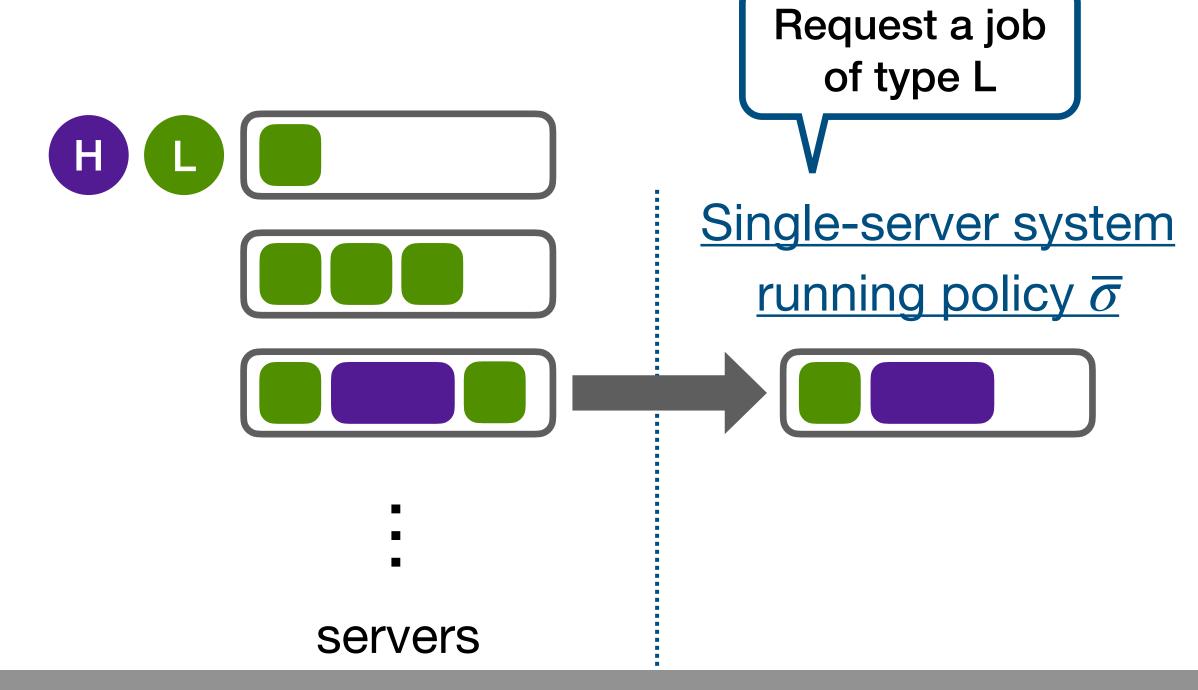- If no tokens, go to an inactive server

How is the throughput related to # active servers via tokens?

Request a job of type L

Single-server system running policy $\overline{\sigma}$

jobs

servers

# Policy conversion: more details



jobs

servers

# Policy conversion: more details



jobs

servers

Run single-server policy $\overline{\sigma}$ for only

$$\overline{N} = \frac{\text{arrival rate}}{\text{throughput}(\overline{\sigma})} \text{ servers}$$

# Policy conversion: more details



Run single-server policy $\overline{\sigma}$ for only

$$\overline{N} = \frac{\text{arrival rate}}{\text{throughput}(\overline{\sigma})} \text{ servers}$$

Recall that we aim to show

$$\mathbf{E}\,[\text{\# active servers}] \leq \left(1 + O\left(r^{-0.5}\right)\right) \cdot \overline{N}$$

jobs

servers
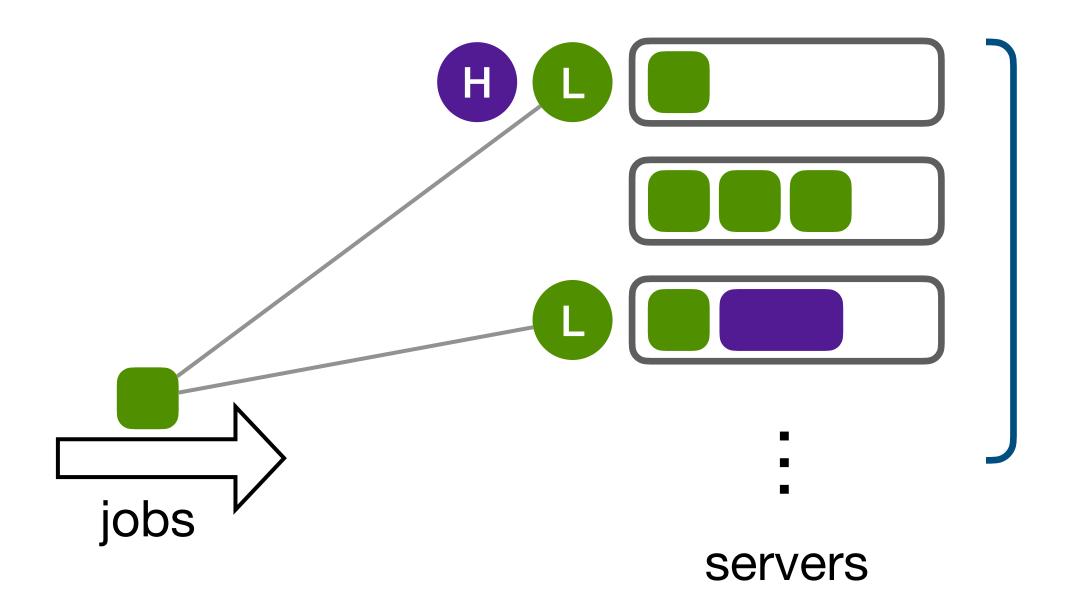
# Policy conversion: more details



Run single-server policy $\overline{\sigma}$ for only

$$\overline{N} = \frac{\text{arrival rate}}{\text{throughput}(\overline{\sigma})} \text{ servers}$$

jobs

servers

Recall that we aim to show

$$\mathbf{E}\,[\text{\# active servers}] \leq \left(1 + O\left(r^{-0.5}\right)\right) \cdot \overline{N}$$

When the # tokens of a type $> \sqrt{r}$, remove the overflow tokens and generate virtual jobs

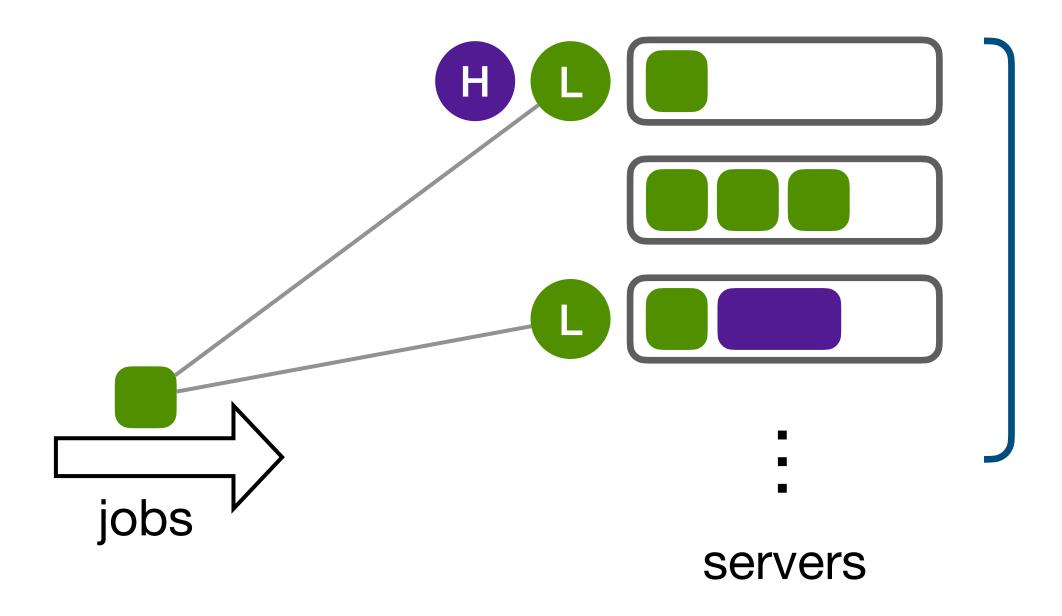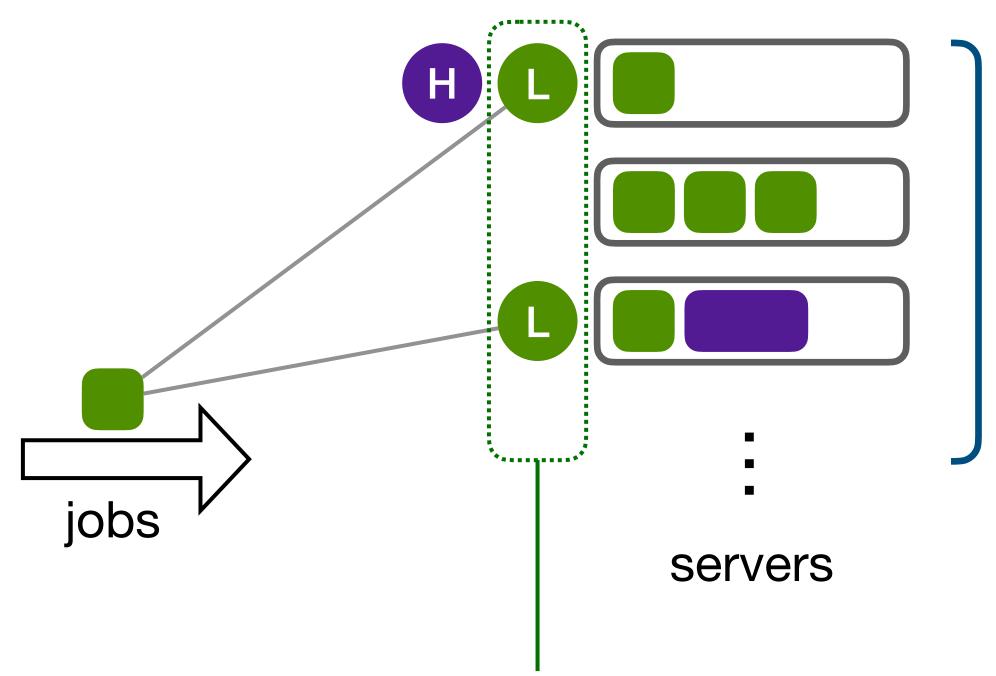# Policy conversion: more details



Run single-server policy $\overline{\sigma}$ for only

$$\overline{N} = \frac{\text{arrival rate}}{\text{throughput}(\overline{\sigma})} \text{ servers}$$

jobs

servers

Recall that we aim to show

$$\mathbf{E}\left[\text{\# active servers}\right] \leq \left(1 + O\left(r^{-0.5}\right)\right) \cdot \overline{N}$$

When the # tokens of a type $> \sqrt{r}$, remove the overflow tokens and generate virtual jobs
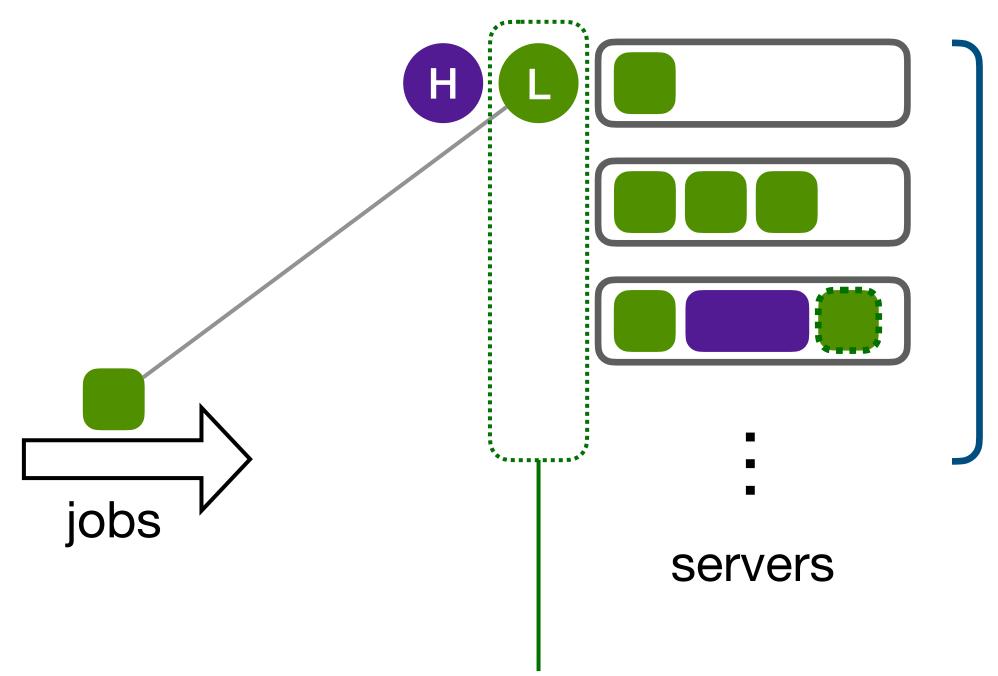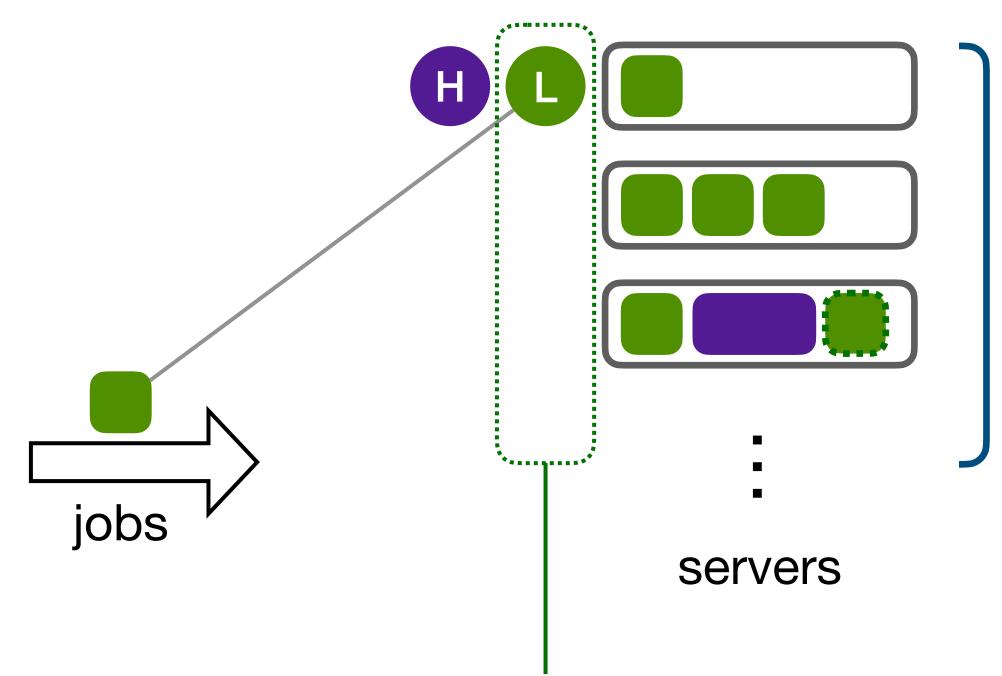
# Policy conversion: more details



Run single-server policy $\overline{\sigma}$ for only

$$\overline{N} = \frac{\text{arrival rate}}{\text{throughput}(\overline{\sigma})} \text{ servers}$$

jobs

servers

Recall that we aim to show

$$\mathbf{E}\,[\text{\# active servers}] \leq \left(1 + O\left(r^{-0.5}\right)\right) \cdot \overline{N}$$

When the # tokens of a type $> \sqrt{r}$, remove the overflow tokens and generate virtual jobs

We can prove that $\mathbf{E}\,[\text{\# virtual jobs}] = O\left(r^{0.5}\right)$

# Key proof idea 1



jobs

servers

# Key proof idea 1



Single-server system
running policy $\overline{\sigma}$

jobs

servers

# Key proof idea 1

Will show that each server in the original system
$\approx$ an independent single-server system



Single-server system
running policy $\overline{\sigma}$

jobs

servers

# Key proof idea 1

Will show that each server in the original system
$\approx$ an independent single-server system



Single-server system
running policy $\overline{\sigma}$

If only each token were
replaced by a job
immediately …

jobs

servers

# Key proof idea 1

Will show that each server in the original system
$\approx$ an independent single-server system



Single-server system
running policy $\overline{\sigma}$

jobs

servers

Difficulty: the dynamics of a server in the original system
depends on other servers through arrivals & token overflows

# Key proof idea 1

Will show that each server in the original system $\approx$ an independent single-server system



Single-server system
running policy $\overline{\sigma}$

jobs

servers

Difficulty: the dynamics of a server in the original system depends on other servers through arrivals & token overflows

Idea: for each type $i$, consider

$\widetilde{K}_i$ = # jobs + # virtual jobs + # tokens

# Key proof idea 1

Will show that each server in the original system $\approx$ an independent single-server system



Single-server system
running policy $\overline{\sigma}$

**Difficulty:** the dynamics of a server in the original system depends on other servers through arrivals & token overflows

**Idea:** for each type $i$, consider

$$\widetilde{K}_i = \text{\# jobs} + \text{\# virtual jobs} + \text{\# tokens}$$

Why does considering $\widetilde{K}_i$ help decouple servers?

jobs

servers

# Key proof idea 1

Will show that each server in the original system
$\approx$ an independent single-server system



Single-server system
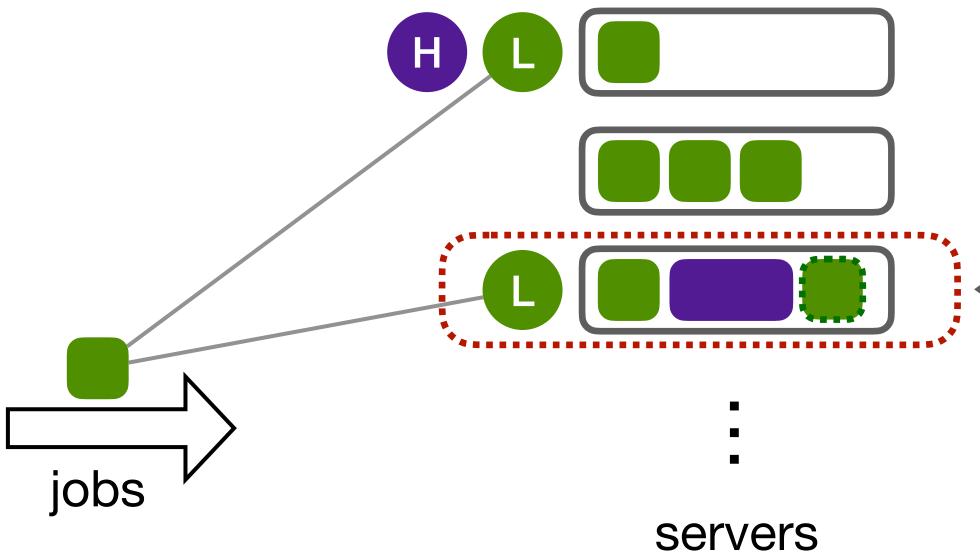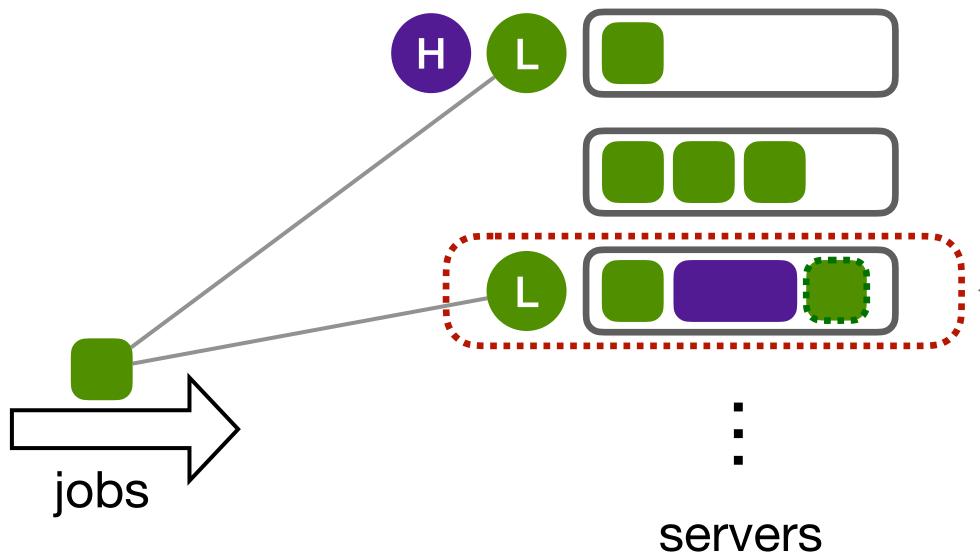running policy $\overline{\sigma}$

jobs

servers

**Difficulty:** the dynamics of a server in the original system
depends on other servers through arrivals & token overflows

**Idea:** for each type $i$, consider
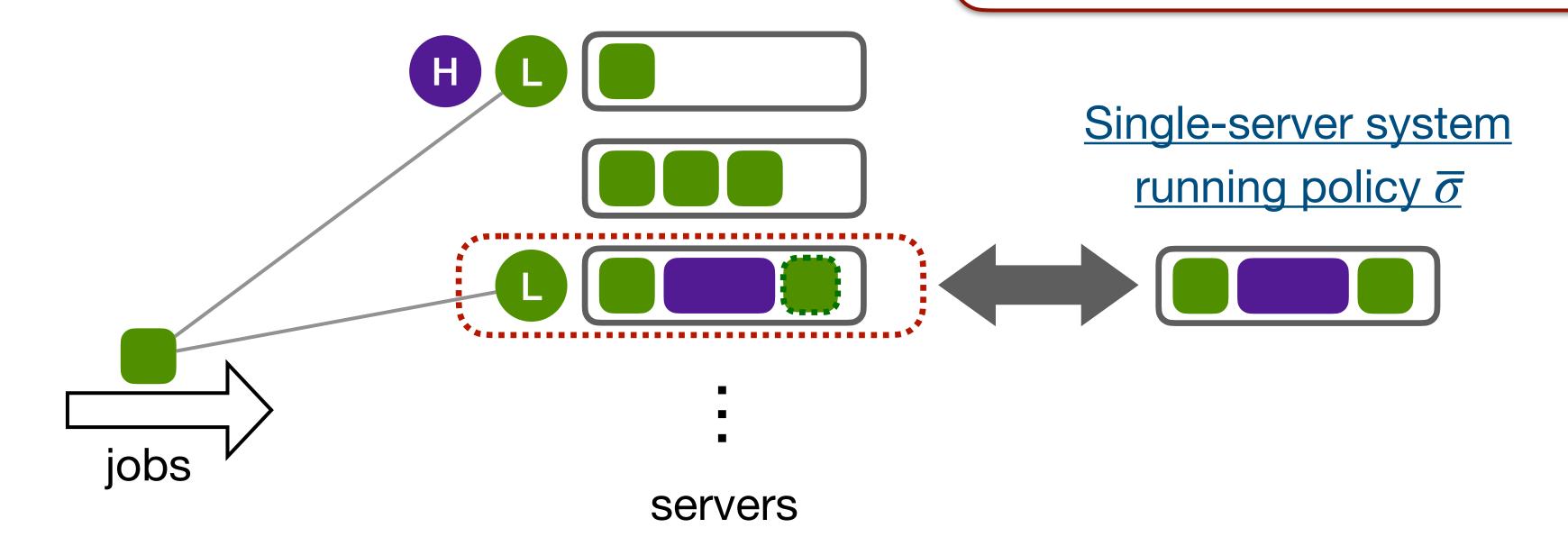
$$\widetilde{K}_i = \text{# jobs} + \text{# virtual jobs} + \text{# tokens}$$

- Arrivals & token overflows do not affect $\widetilde{K}_i$

Why does considering $\widetilde{K}_i$ help
decouple servers?

# Key proof idea 1

Will show that each server in the original system
$\approx$ an independent single-server system



Single-server system
running policy $\overline{\sigma}$

jobs

servers

**Idea:** for each type $i$, consider

$$\widetilde{K}_i = \text{\# jobs} + \text{\# virtual jobs} + \text{\# tokens} \quad \text{v.s.} \quad \overline{K}_i = \text{\# jobs of type } i$$

• Arrivals & token overflows do not affect $\widetilde{K}_i$

# **Key proof idea 1**

Will show that each server in the original system
$\approx$ an independent single-server system



Single-server system
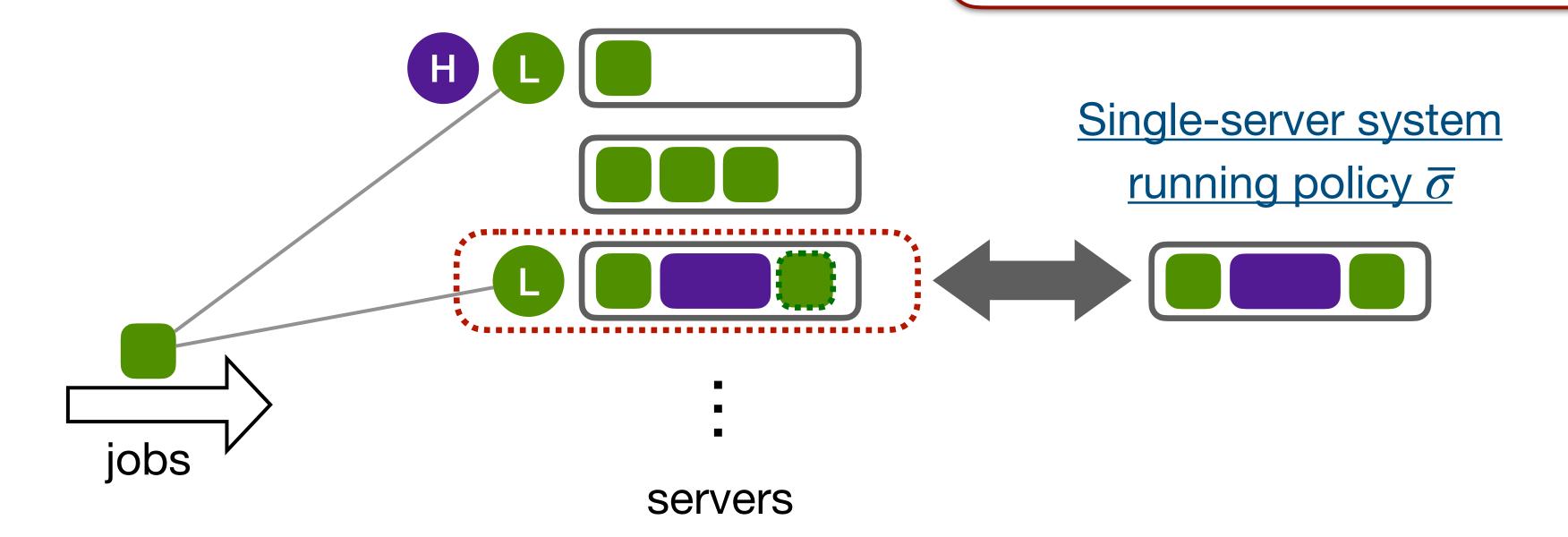running policy $\overline{\sigma}$

jobs

servers

**Idea:** for each type $i$, consider

$\widetilde{K}_i$ = # jobs + # virtual jobs + # tokens   v.s.   $\overline{K}_i$ = # jobs of type $i$

- Arrivals & token overflows do not affect $\widetilde{K}_i$

- Requests by $\overline{\sigma}$ change $\widetilde{K}_i$ and $\overline{K}_i$ in the same way, difference bounded by # tokens

# Key proof idea 1

Will show that each server in the original system $\approx$ an independent single-server system

Single-server system
running policy $\overline{\sigma}$

**H** **L**

**L**

jobs

⋮

servers

**Idea:** for each type $i$, consider

$\widetilde{K}_i$ = # jobs + # virtual jobs + # tokens    v.s.    $\overline{K}_i$ = # jobs of type $i$

- Arrivals & token overflows do not affect $\widetilde{K}_i$

- Requests by $\overline{\sigma}$ change $\widetilde{K}_i$ and $\overline{K}_i$ in the same way, difference bounded by # tokens

- Job phase transitions in $\widetilde{K}_i$ and $\overline{K}_i$ differ by # tokens

# Key proof idea 1

Will show that each server in the original system $\approx$ an independent single-server system



Single-server system
running policy $\overline{\sigma}$
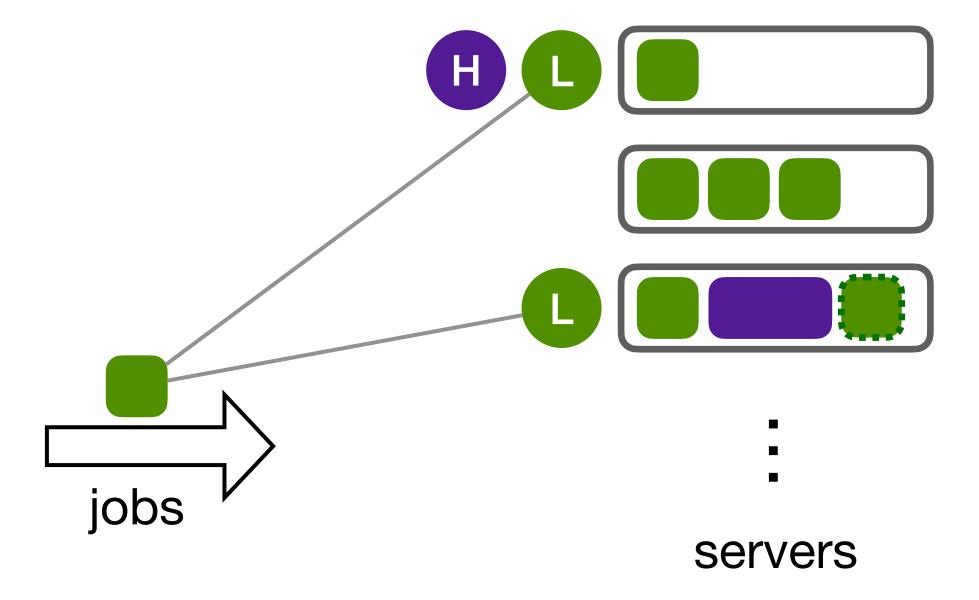
Using Stein's method, we show

$$d_W\left(\widetilde{K}^{1:\overline{N}}, \overline{K}^{1:\overline{N}}\right) = O\left(r^{0.5}\right)$$

**Idea:** for each type $i$, consider

$\widetilde{K}_i$ = # jobs + # virtual jobs + # tokens    v.s.    $\overline{K}_i$ = # jobs of type $i$

- Arrivals & token overflows do not affect $\widetilde{K}_i$

- Requests by $\overline{\sigma}$ change $\widetilde{K}_i$ and $\overline{K}_i$ in the same way, difference bounded by # tokens

- Job phase transitions in $\widetilde{K}_i$ and $\overline{K}_i$ differ by # tokens

# Key proof idea 2



jobs

servers

# Key proof idea 2



jobs

servers

When the # tokens of a type $> \sqrt{r}$, remove the overflow tokens and generate virtual jobs

# Key proof idea 2



jobs

servers

$\overline{N}$

When the # tokens of a type $> \sqrt{r}$, remove the overflow tokens and generate virtual jobs

# Key proof idea 2



When the # tokens of a type $> \sqrt{r}$, remove the overflow tokens and generate virtual jobs

# Key proof idea 2



Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,

and # backup servers $= O\left(\sqrt{r}\right)$

# Key proof idea 2



Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,

and # backup servers $= O\left(\sqrt{r}\right)$

$\overline{N}$

backup servers

jobs

servers

$\sqrt{r}$

# Key proof idea 2



Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,

and # backup servers $= O\left(\sqrt{r}\right)$

a server requests
a type L job

$\sqrt{r}$

jobs

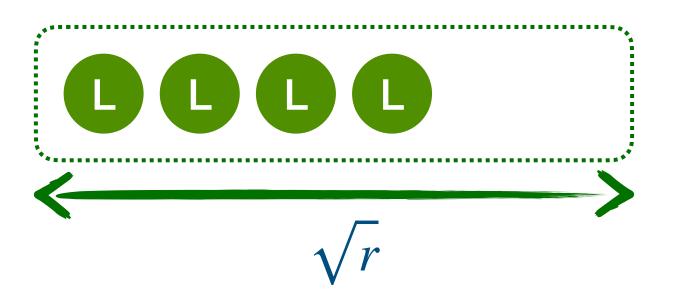servers

$\overline{N}$

backup
servers

# Key proof idea 2

Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,

and # backup servers $= O\left(\sqrt{r}\right)$



jobs

servers

$\overline{N}$

backup servers

a type L job arrives

a server requests a type L job

$\sqrt{r}$

# Key proof idea 2
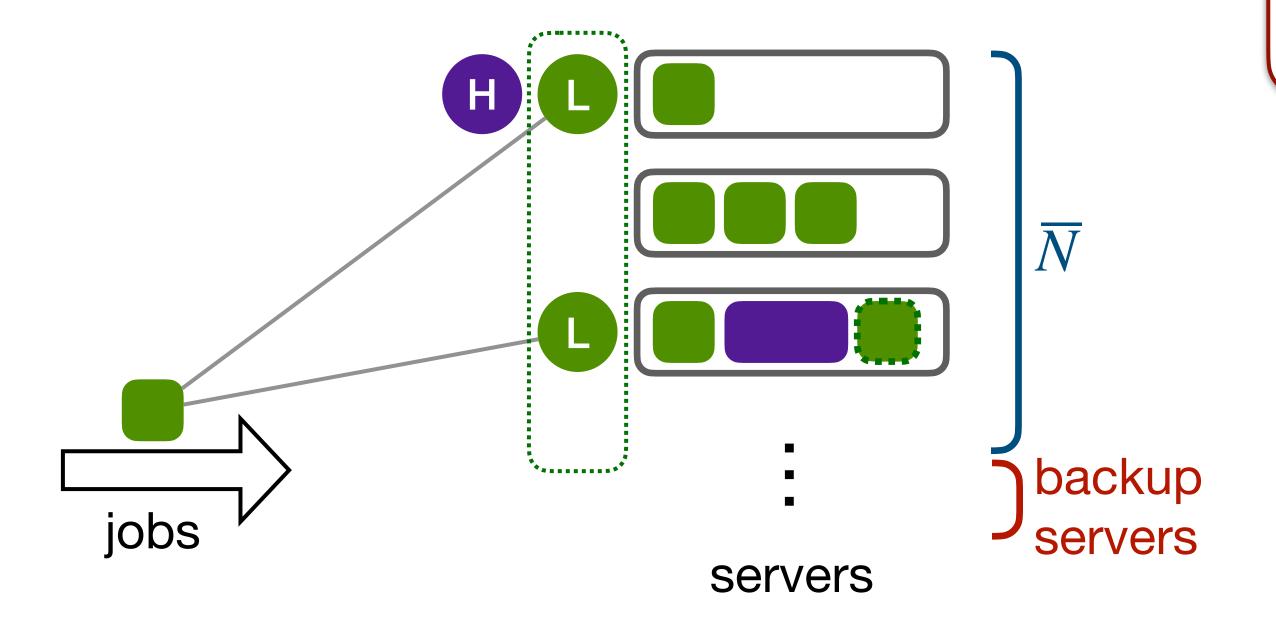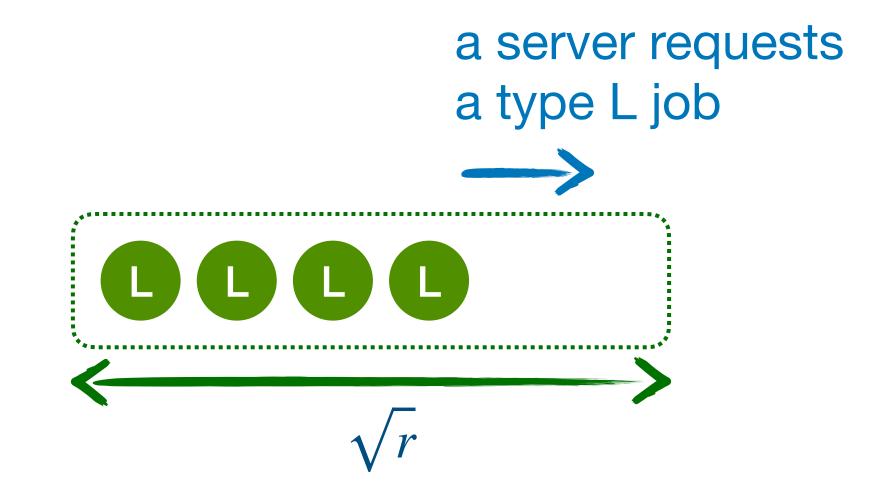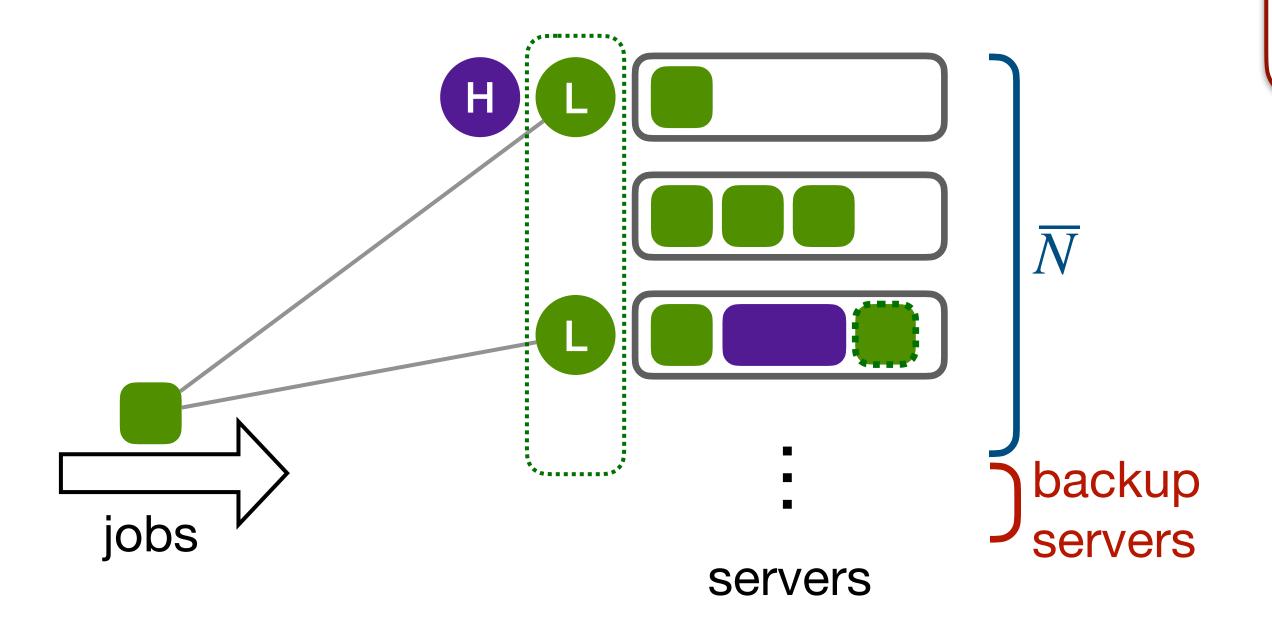


Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,

and # backup servers $= O\left(\sqrt{r}\right)$

jobs

servers

$\overline{N}$

backup servers

a type L job arrives

a server requests a type L job

$\sqrt{r}$

What happens when # tokens hits $\sqrt{r}$ ?

# Key proof idea 2



jobs

servers

$\overline{N}$

backup servers

Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,
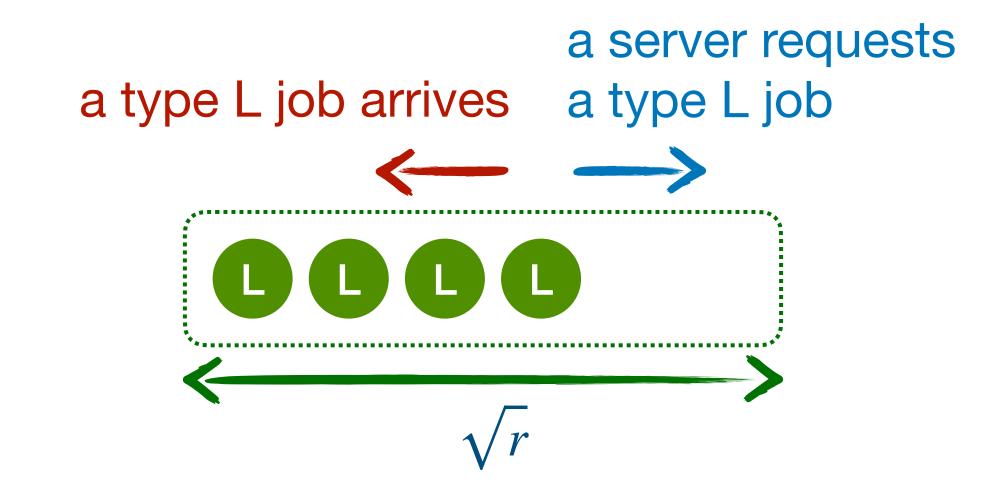
and # backup servers $= O\left(\sqrt{r}\right)$

a type L job arrives

a server requests a type L job

$\sqrt{r}$

What happens when # tokens hits $\sqrt{r}$ ?

Generate a virtual job

# Key proof idea 2



Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,
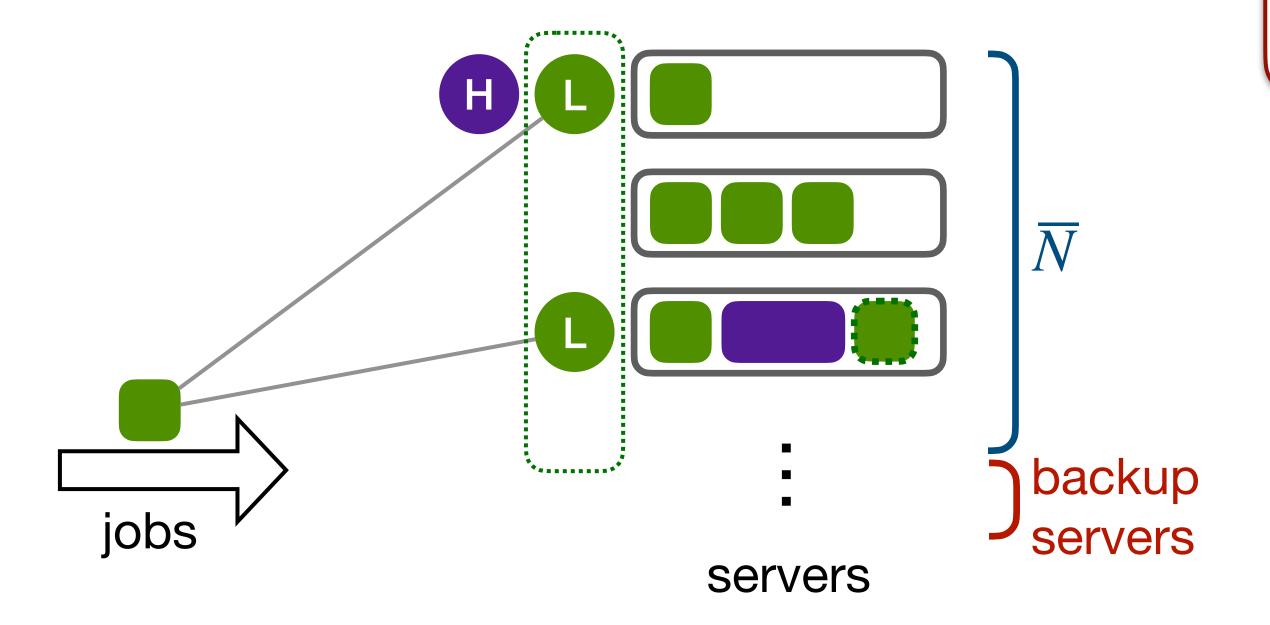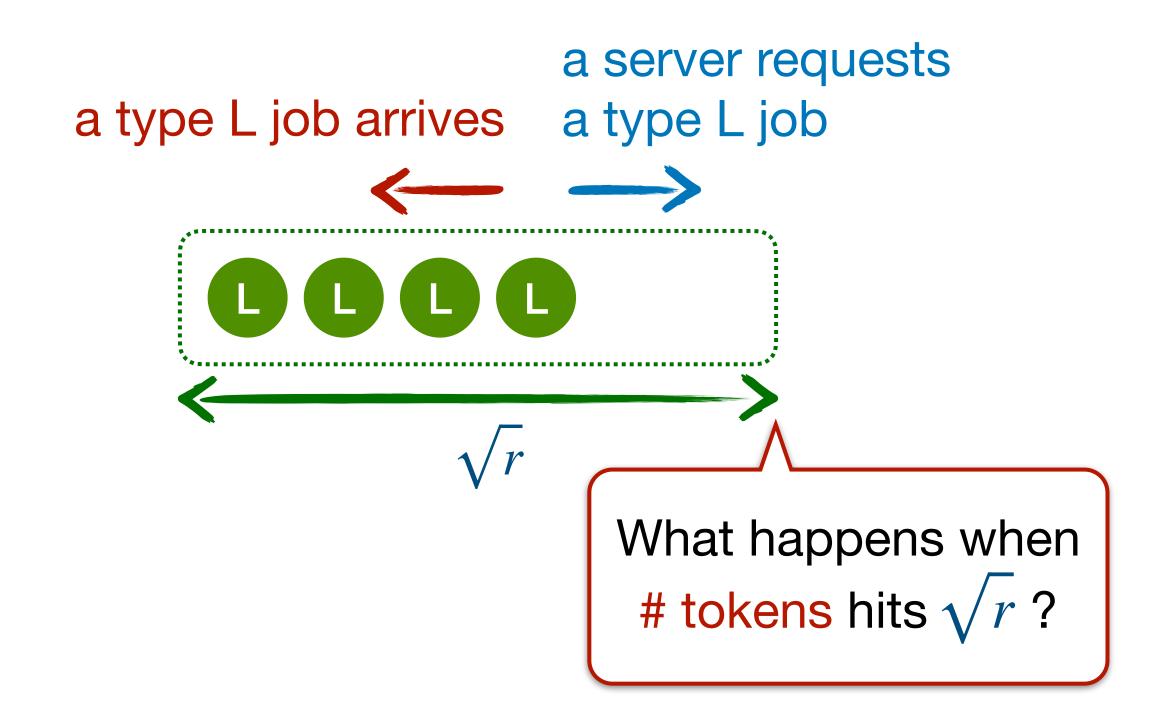
and # backup servers $= O\left(\sqrt{r}\right)$

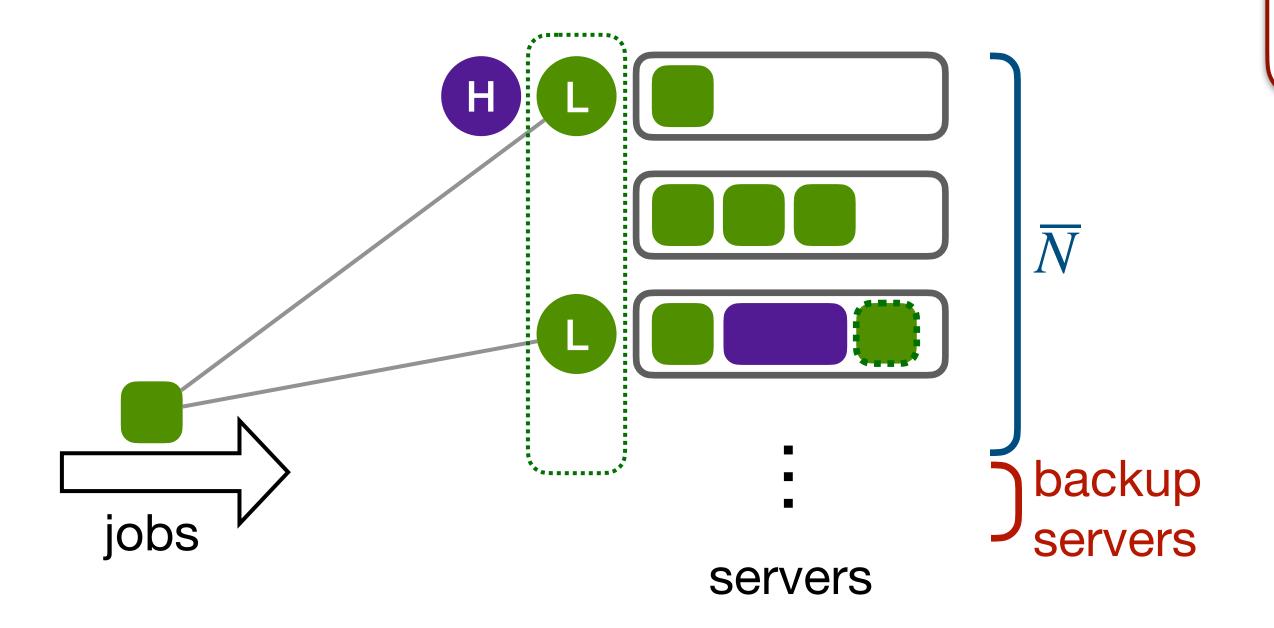$\overline{N}$

backup servers

jobs

servers

a type L job arrives

a server requests a type L job

$\sqrt{r}$

What happens when # tokens hits $0$?

What happens when # tokens hits $\sqrt{r}$ ?

Generate a virtual job

# Key proof idea 2



Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,

and # backup servers $= O\left(\sqrt{r}\right)$

a type L job arrives

a server requests a type L job

$\sqrt{r}$

What happens when # tokens hits $0$?

What happens when # tokens hits $\sqrt{r}$ ?

Generate a job to backup servers

Generate a virtual job

jobs

servers

$\overline{N}$

backup servers

# Key proof idea 2

Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,

and # backup servers $= O\left(\sqrt{r}\right)$



jobs

servers

$\overline{N}$

backup servers

a type L job arrives

a server requests a type L job

$\sqrt{r}$

# Key proof idea 2

Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,

and # backup servers $= O\left(\sqrt{r}\right)$



jobs

servers

$\overline{N}$

backup servers

a type L job arrives

a server requests a type L job

$\sqrt{r}$

# Key proof idea 2


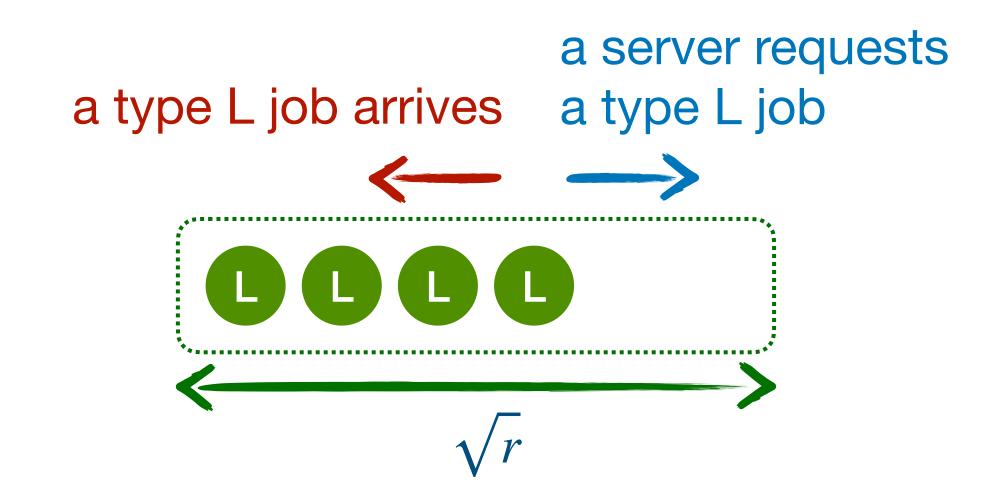
Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,

and # backup servers $= O\left(\sqrt{r}\right)$

$\overline{N}$

backup servers

jobs

servers

a type L job arrives

a server requests a type L job

$\sqrt{r}$

- An almost balanced random walk

# Key proof idea 2



Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,

and # backup servers $= O\left(\sqrt{r}\right)$

a type L job arrives          a server requests
                              a type L job

$\sqrt{r}$

- An almost balanced random walk

- Stationary distribution $\approx$ uniform on $\{0, 1, \ldots, \sqrt{r}\}$

# Key proof idea 2

Will show that # virtual jobs $= O\left(\sqrt{r}\right)$,

and # backup servers $= O\left(\sqrt{r}\right)$



H  L

L

$\overline{N}$

} backup servers

jobs

servers

a type L job arrives

a server requests a type L job

$\sqrt{r}$

- An almost balanced random walk

- Stationary distribution $\approx$ uniform on $\{0, 1, \ldots, \sqrt{r}\}$

- Rate of generating virtual jobs
  $\approx$ rate of sending jobs to backup servers
  $\approx$ arrival rate $\big/ \sqrt{r} = O\left(\sqrt{r}\right)$
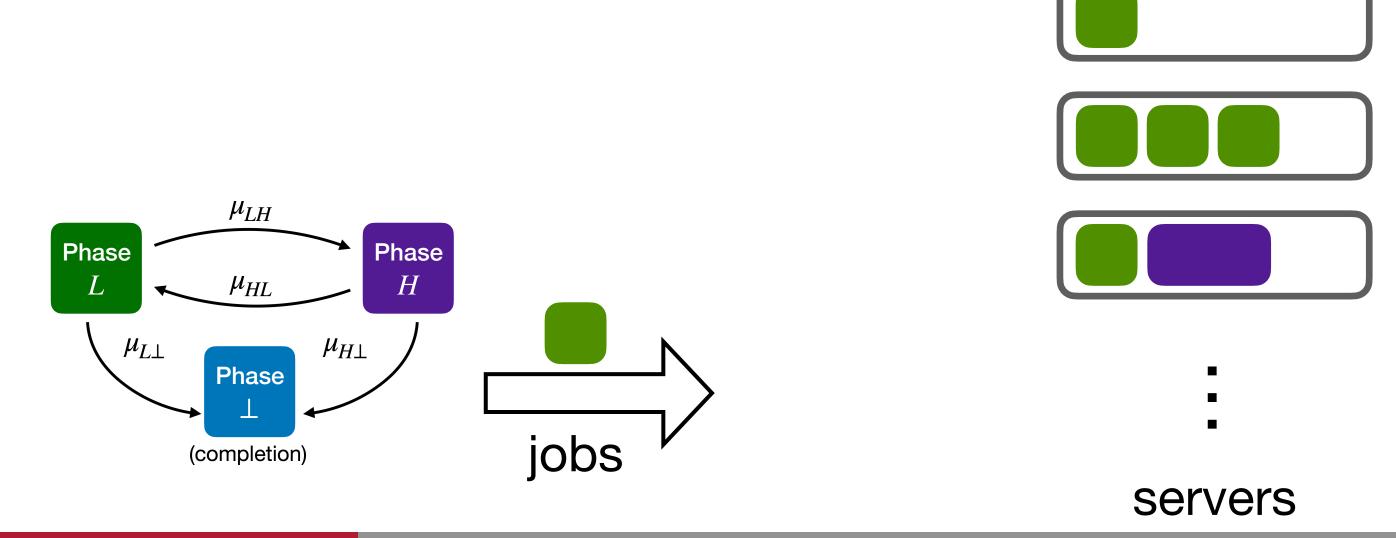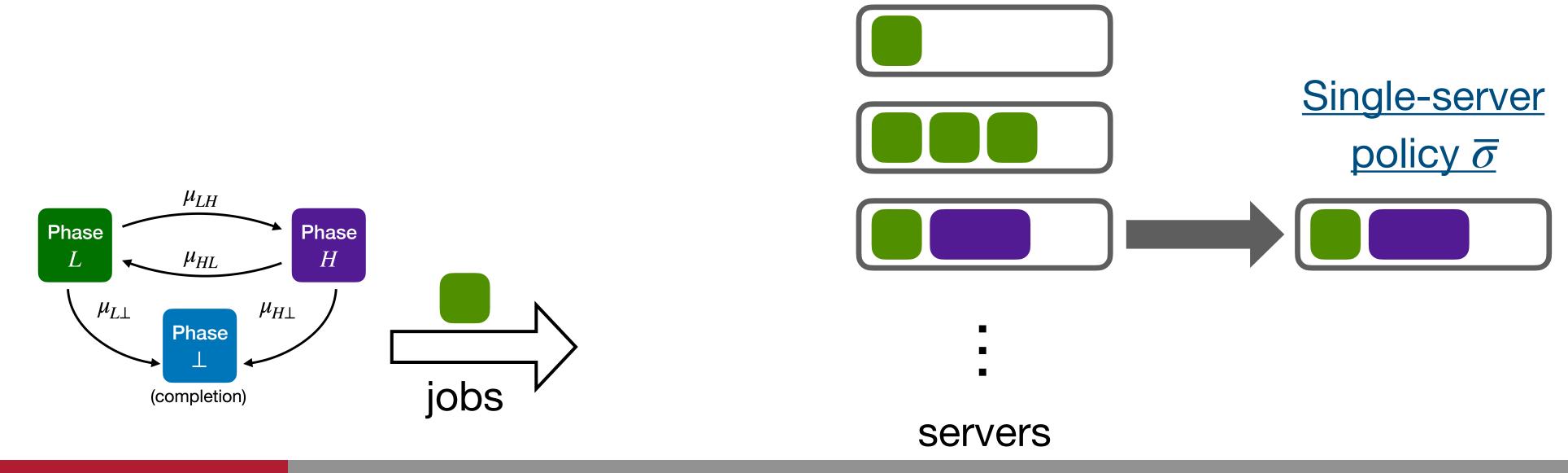
# Summary

# Summary

- We considered the problem of assigning jobs to servers when jobs have time-varying resource requirements

# Summary

- We considered the problem of assigning jobs to servers when jobs have time-varying resource requirements

- We designed an asymptotically optimal policy

# Summary

- We considered the problem of assigning jobs to servers when jobs have time-varying resource requirements

- We designed an asymptotically optimal policy

- We proposed a policy-conversion framework that allows us to reduce the policy-design problem to that in a single-server system

# Summary

- We considered the problem of assigning jobs to servers when jobs have time-varying resource requirements

- We designed an asymptotically optimal policy

- We proposed a policy-conversion framework that allows us to reduce the policy-design problem to that in a single-server system

- A highlight of the framework is the meta-algorithm, JOIN-THE-RECENTLY-REQUESTING-SERVER (JRSS), that converts a single-server policy to a policy in the original system



Single-server policy $\bar{\sigma}$

Phase $L$ — $\mu_{LH}$ → Phase $H$

$\mu_{HL}$

$\mu_{L\perp}$ $\mu_{H\perp}$

Phase $\perp$

(completion)

jobs

servers