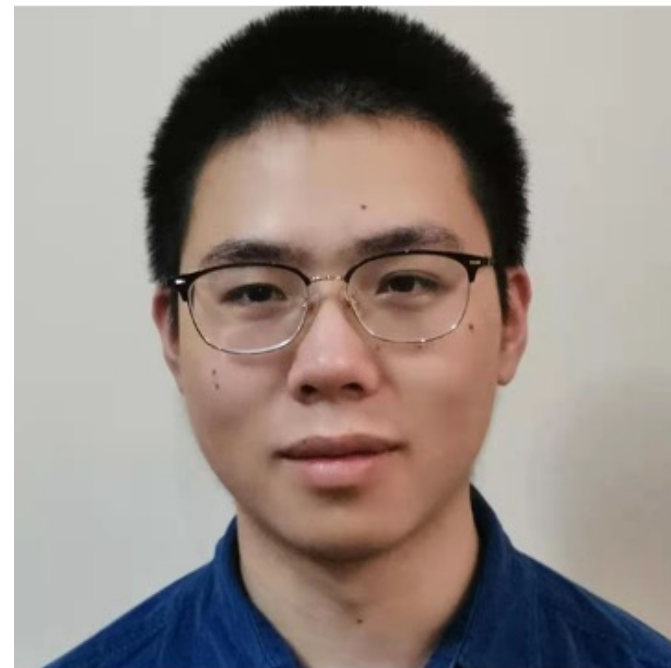


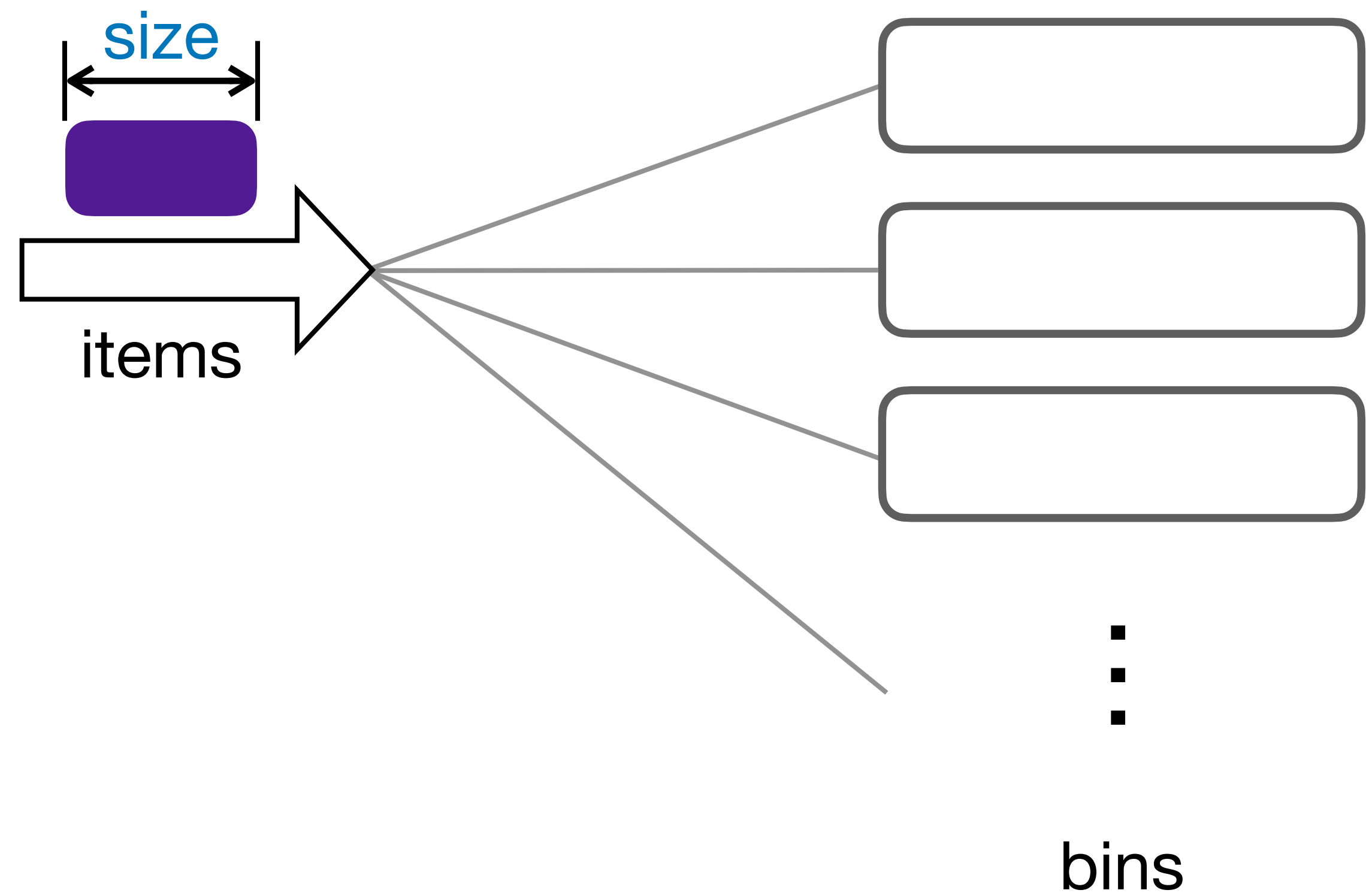
Near-Optimal Stochastic Bin-Packing in Large Service Systems with Time-Varying Item Sizes



Joint work with [Yige Hong \(CMU\)](#) and [Qiaomin Xie \(UW-Madison\)](#)

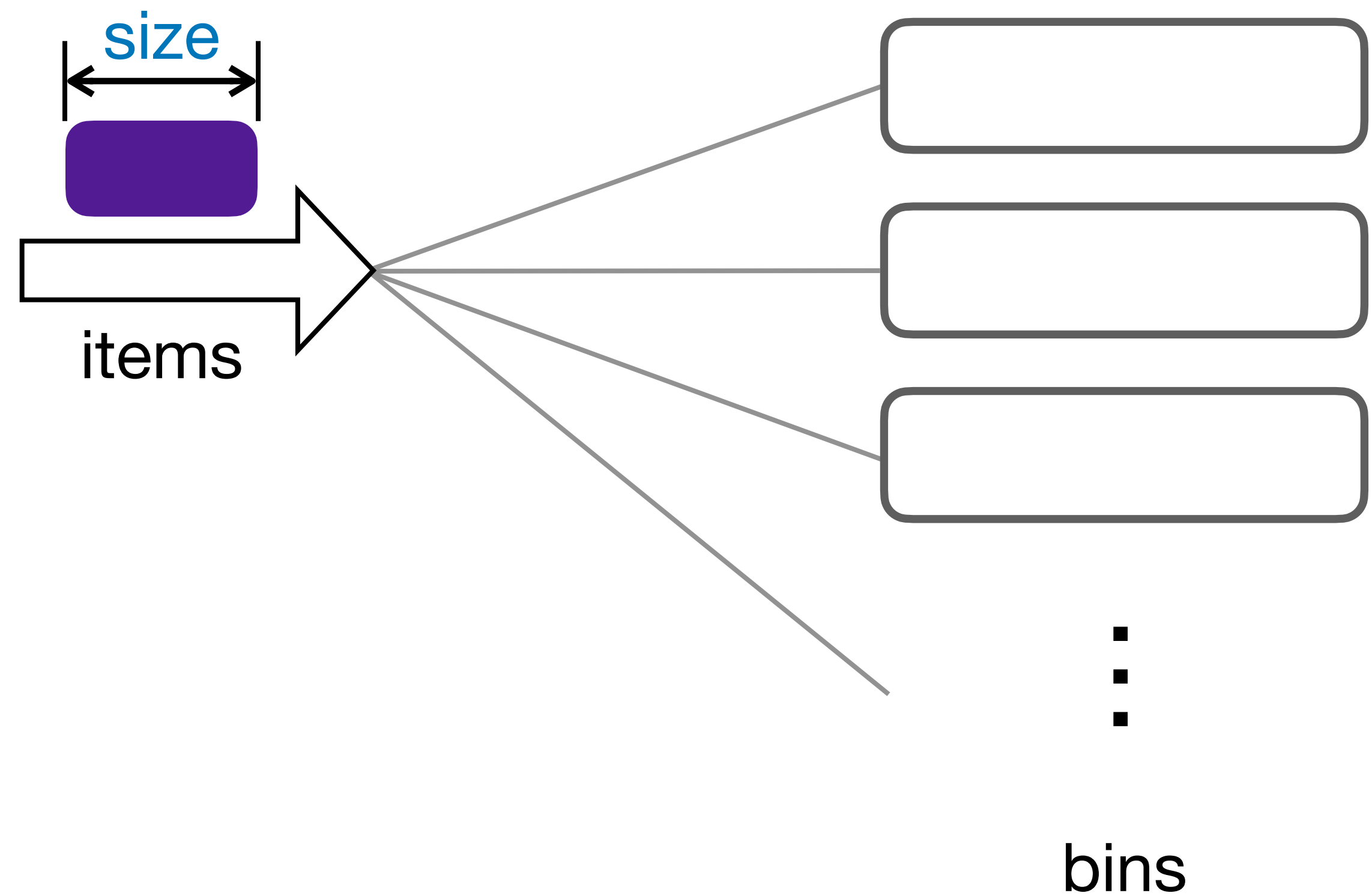
Weina Wang
Carnegie Mellon University

The problem



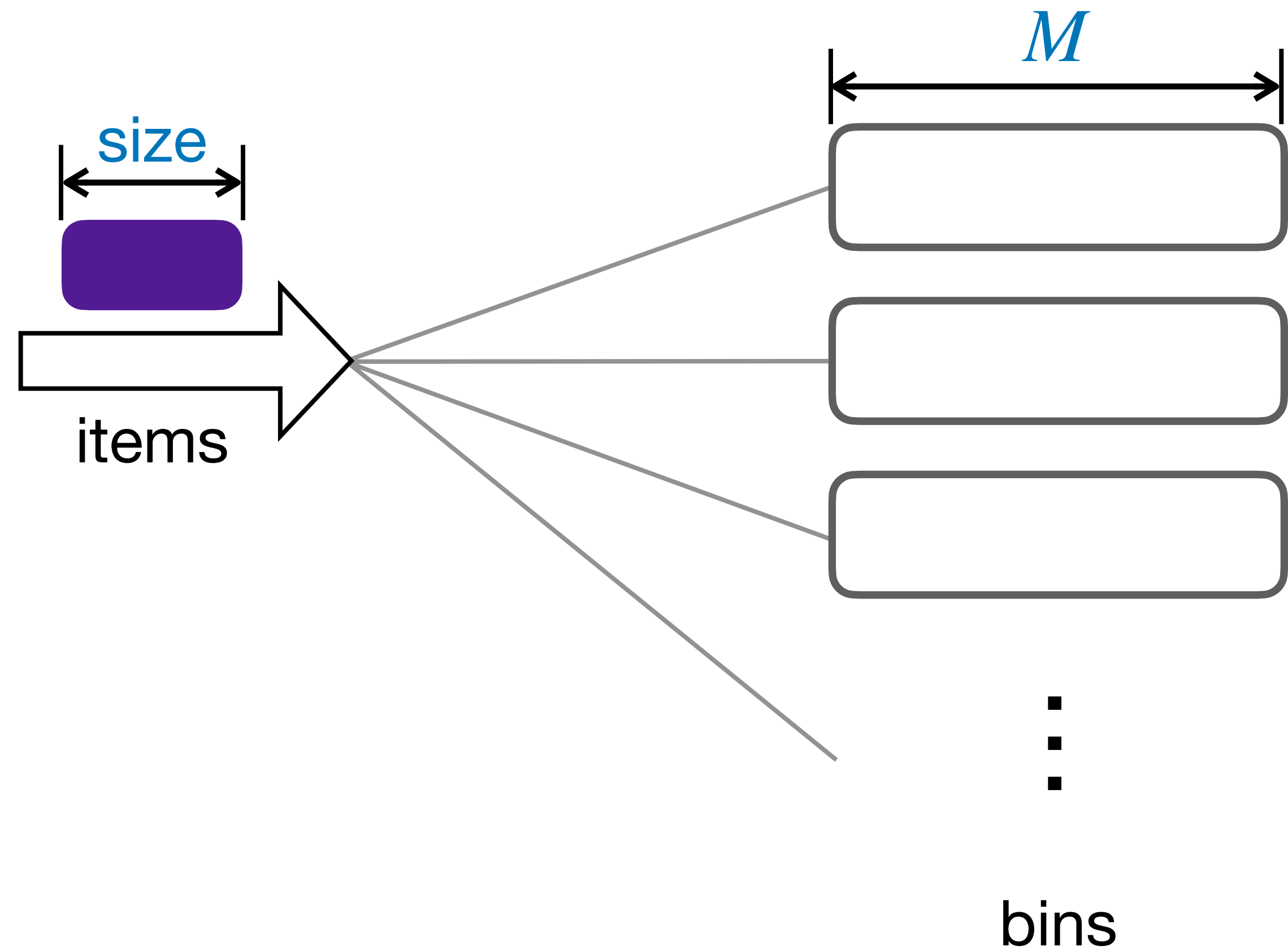
The problem

- Each arriving item needs to be assigned to a bin



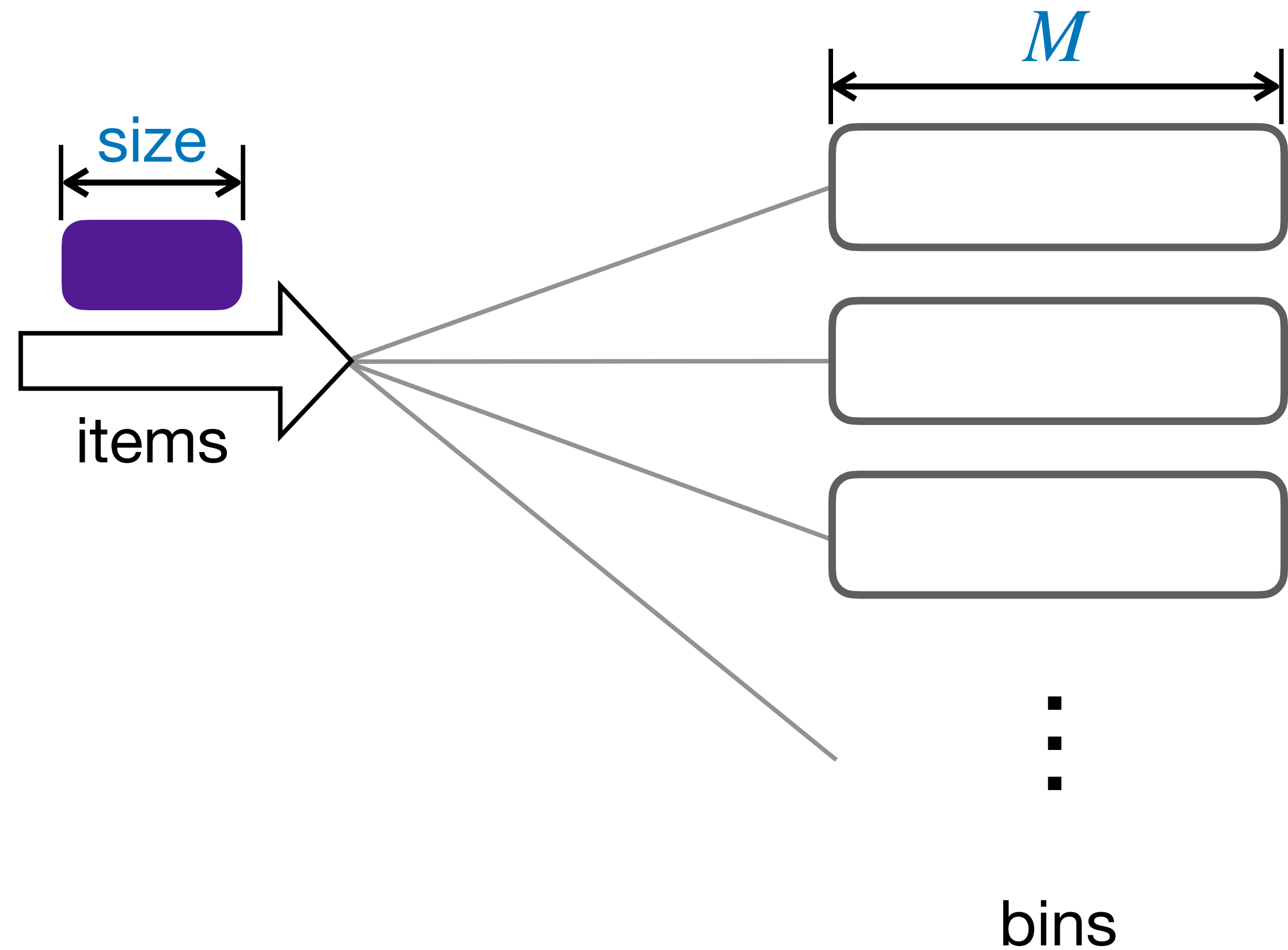
The problem

- Each arriving item needs to be assigned to a bin
- Each bin has a capacity M



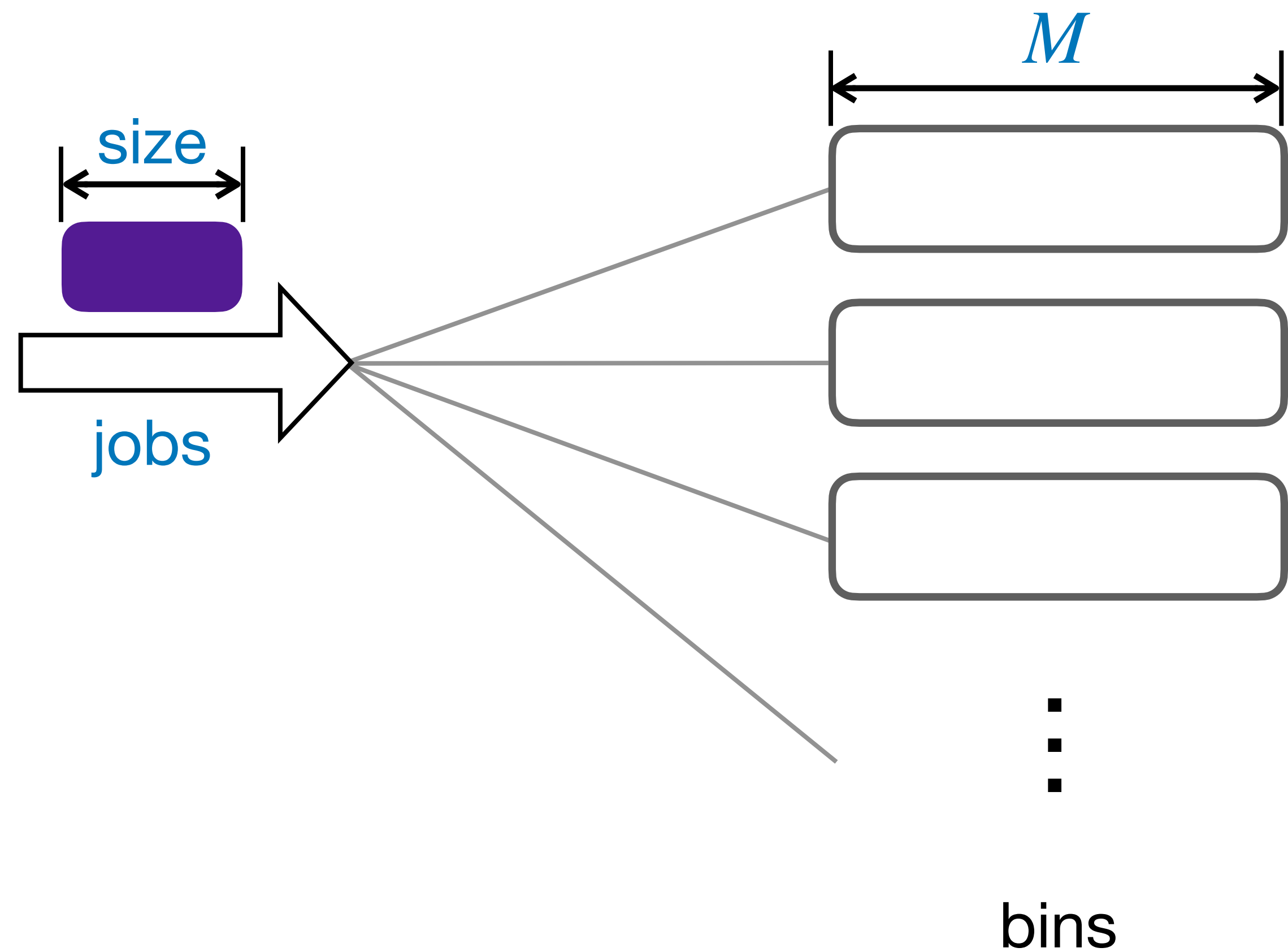
The problem

- Each arriving item needs to be assigned to a bin
- Each bin has a capacity M
- Infinite # bins



The problem

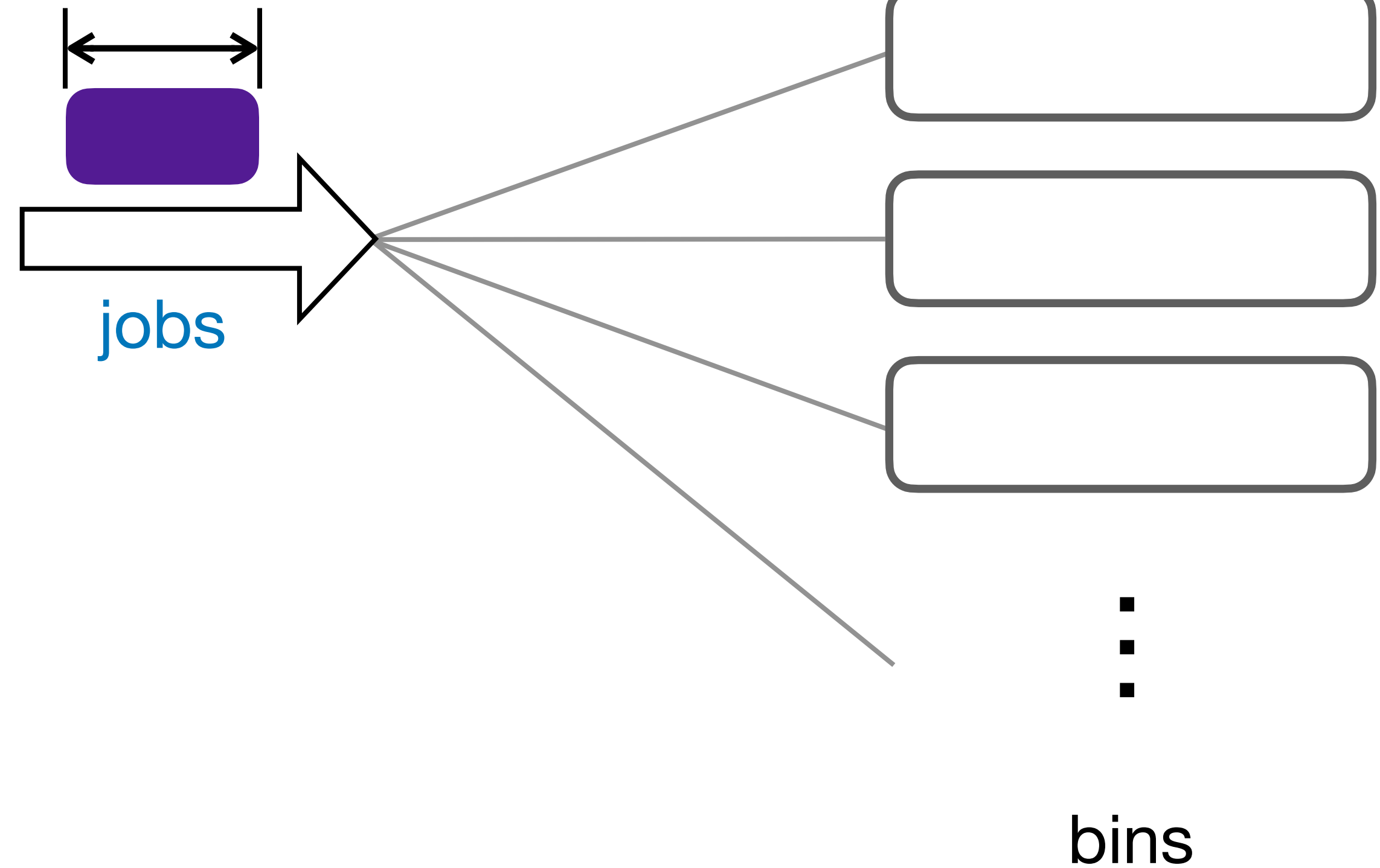
- Each arriving **job** needs to be assigned to a bin
- Each bin has a capacity M
- Infinite # **bins**



The problem

- Each arriving **job** needs to be assigned to a bin
- Each bin has a capacity M
- Infinite # **bins**

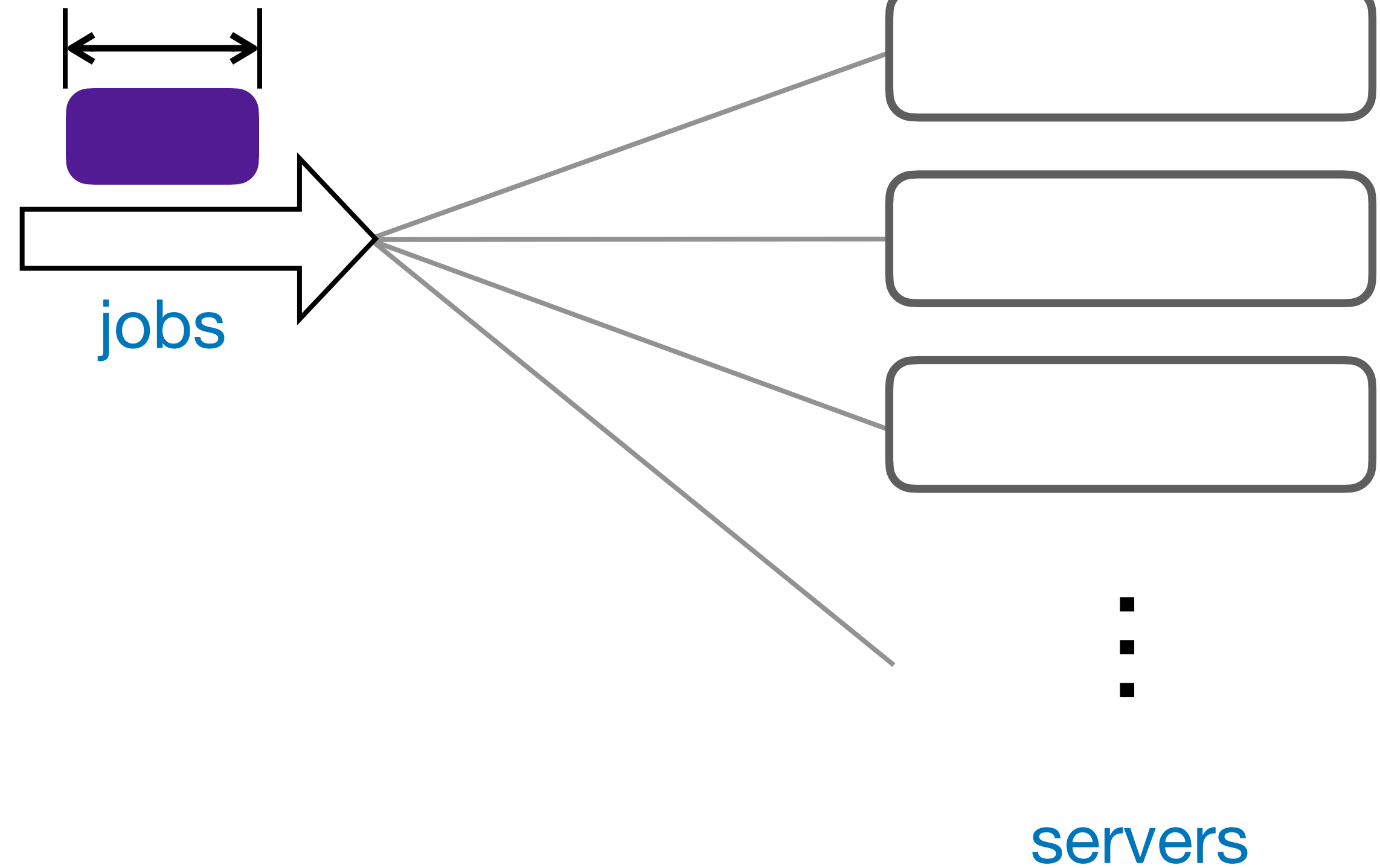
resource requirement:
e.g., # CPUs



The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite **# servers**

resource requirement:
e.g., # CPUs



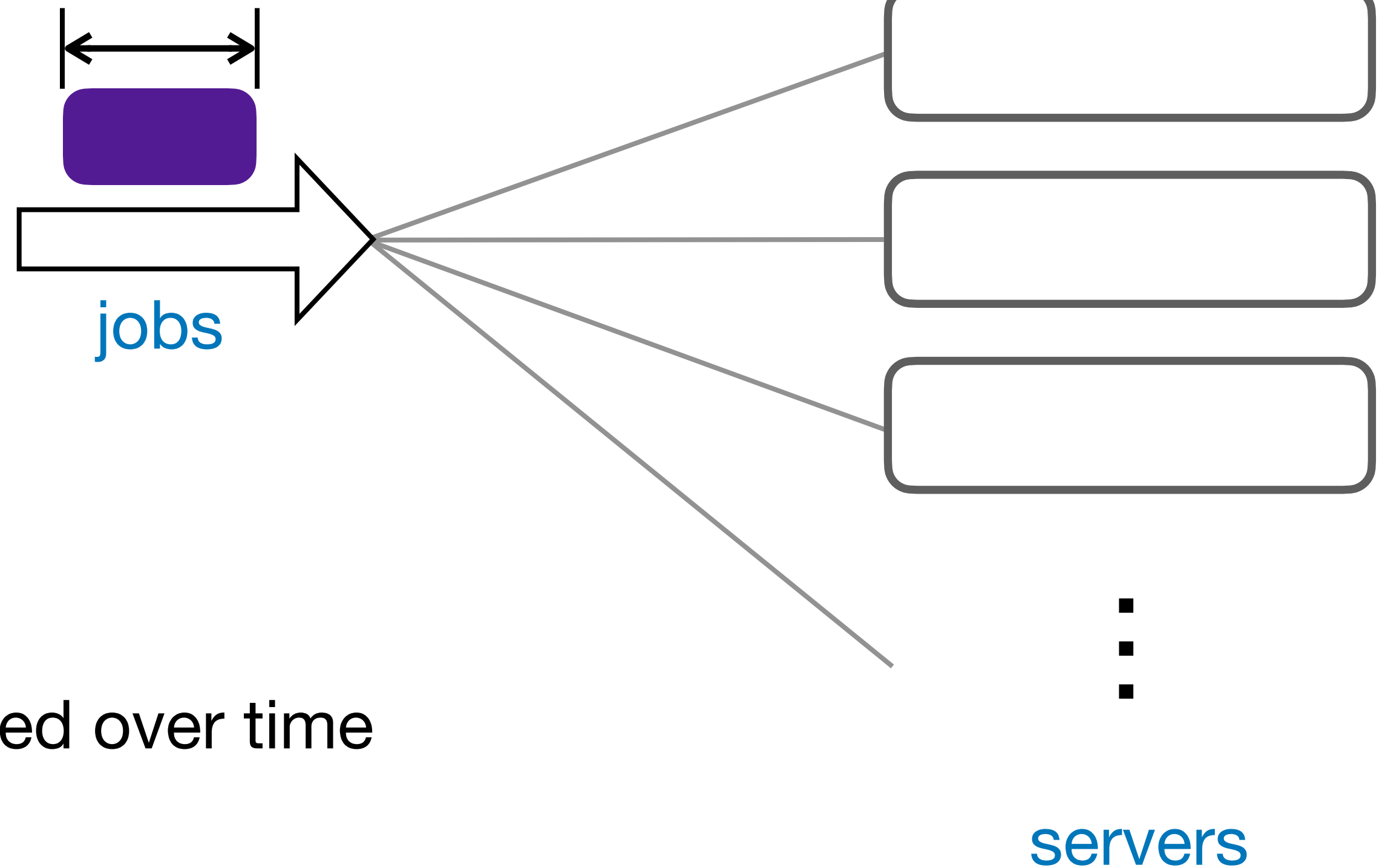
The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite **# servers**

Traditional job model:

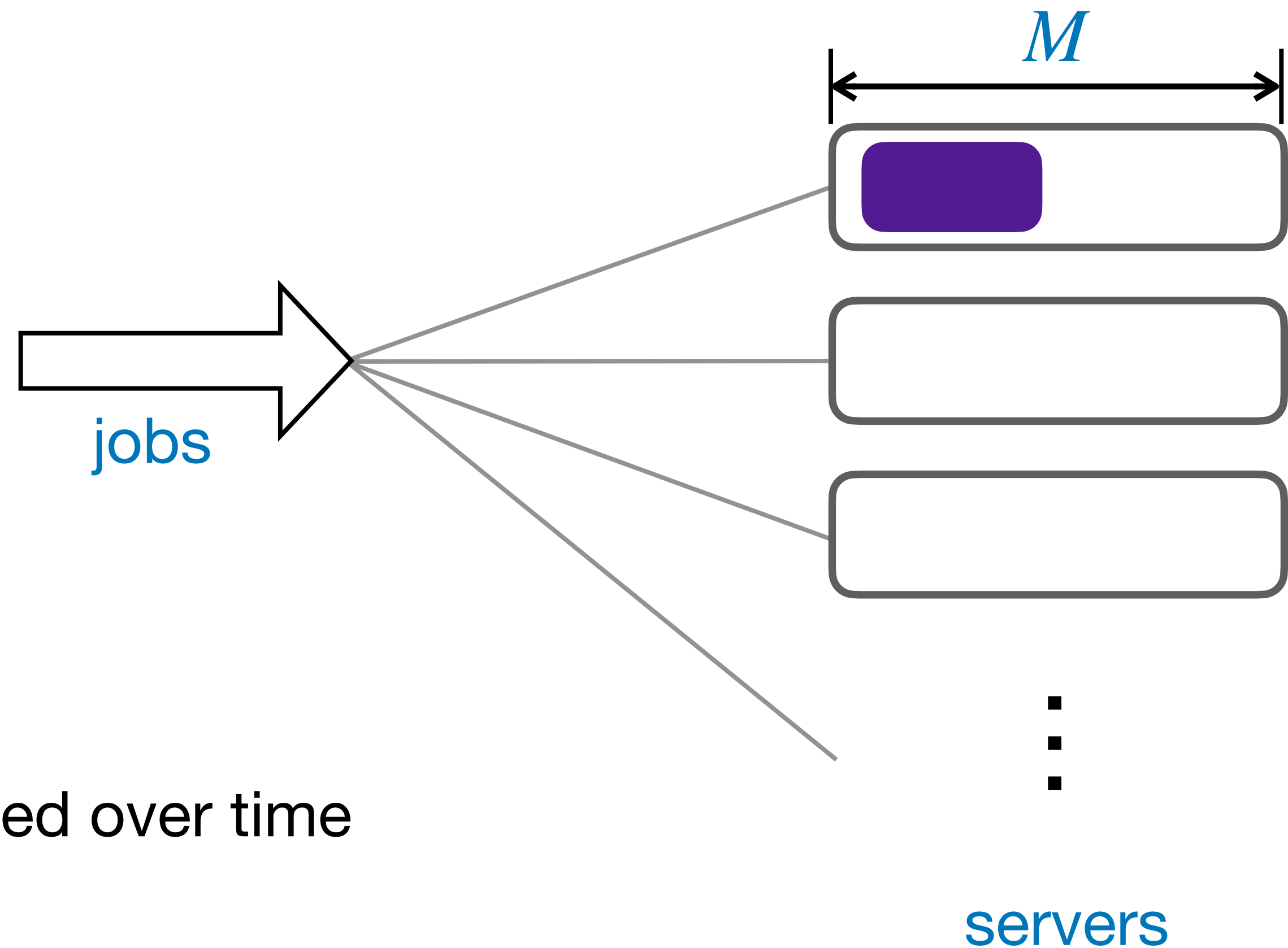
- Each job's resource requirement is fixed over time
- Each job departs after a random time

resource requirement:
e.g., # CPUs



The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**

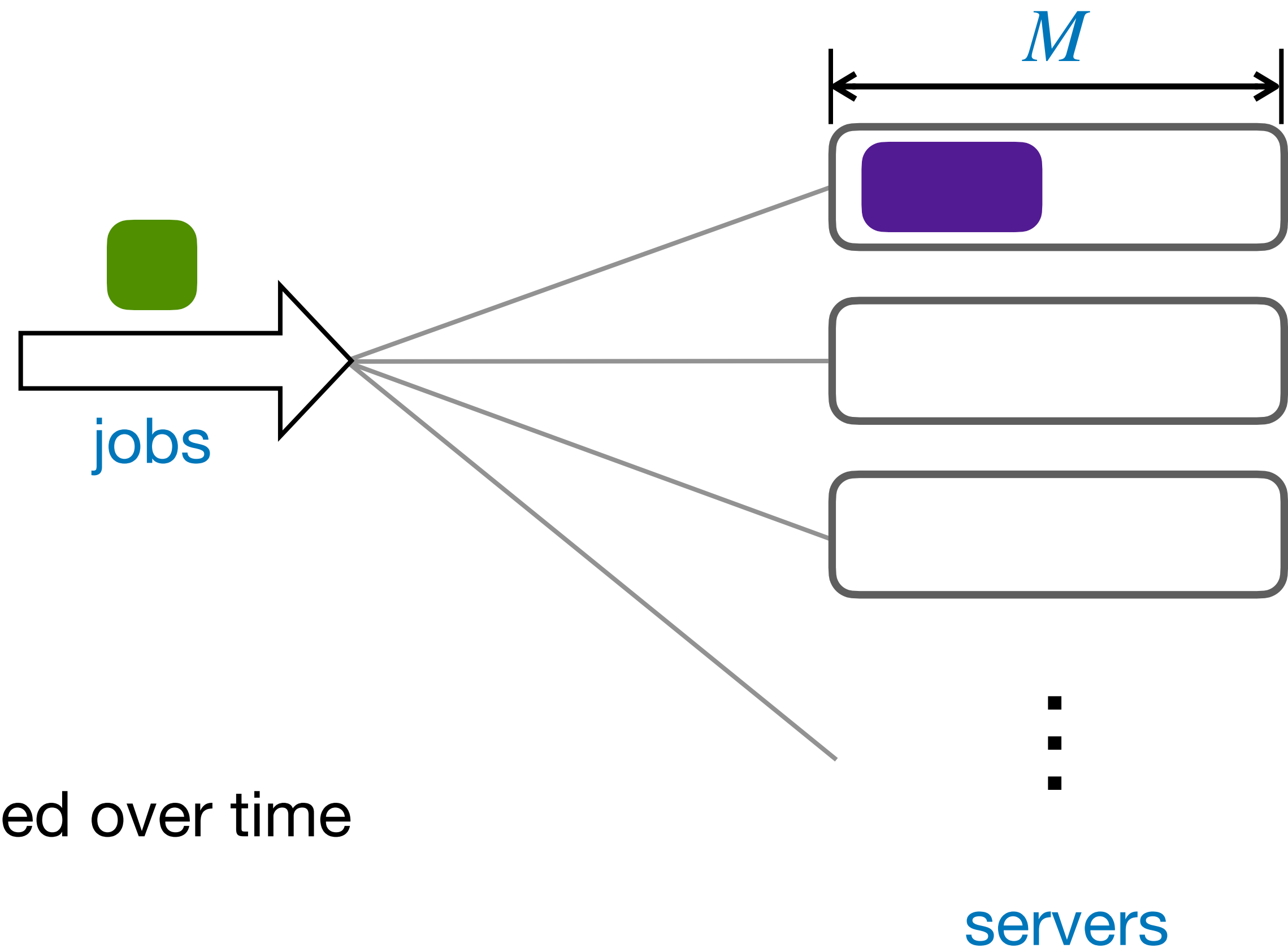


Traditional job model:

- Each job's resource requirement is fixed over time
- Each job departs after a random time

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**

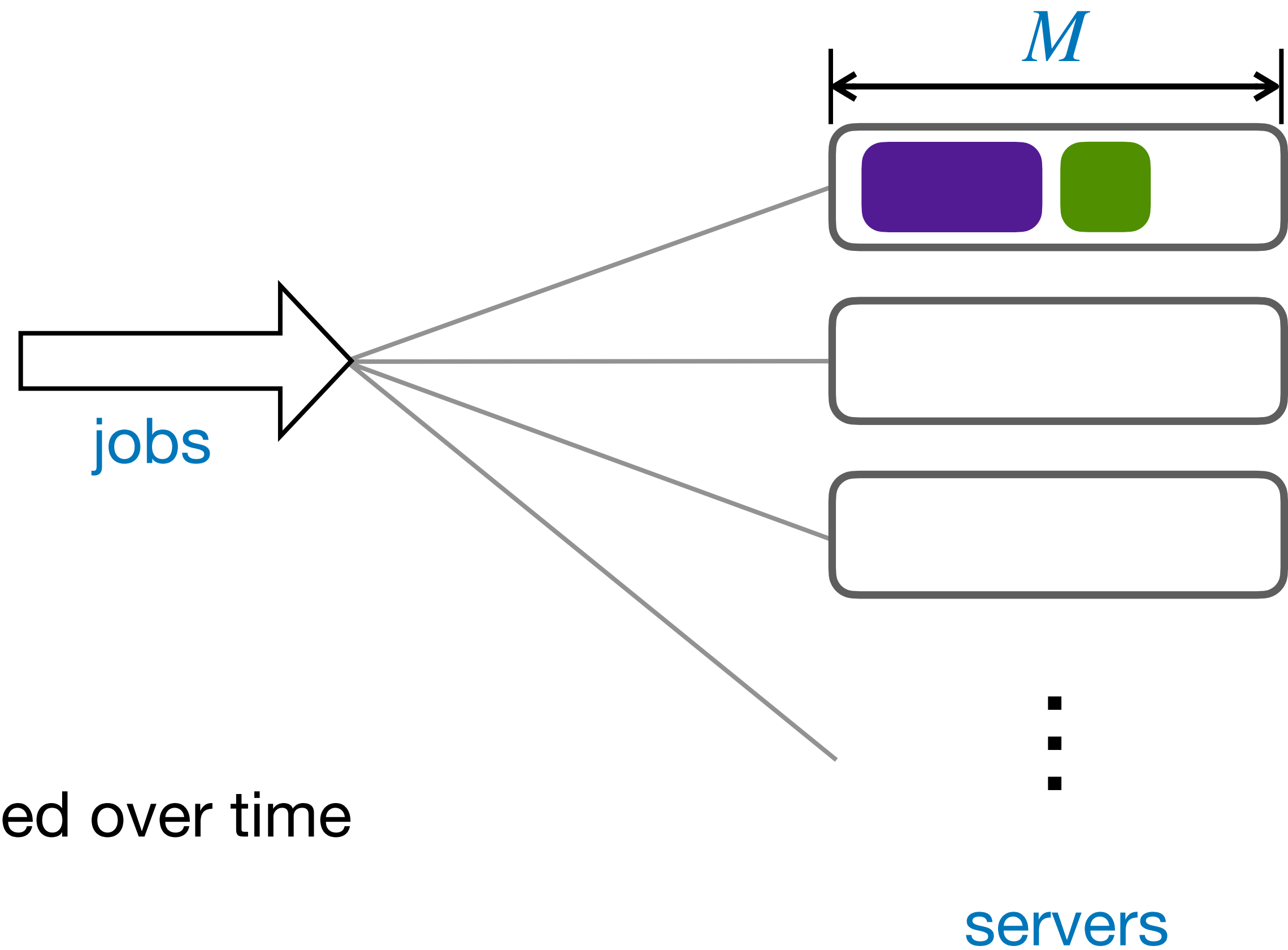


Traditional job model:

- Each job's resource requirement is fixed over time
- Each job departs after a random time

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**

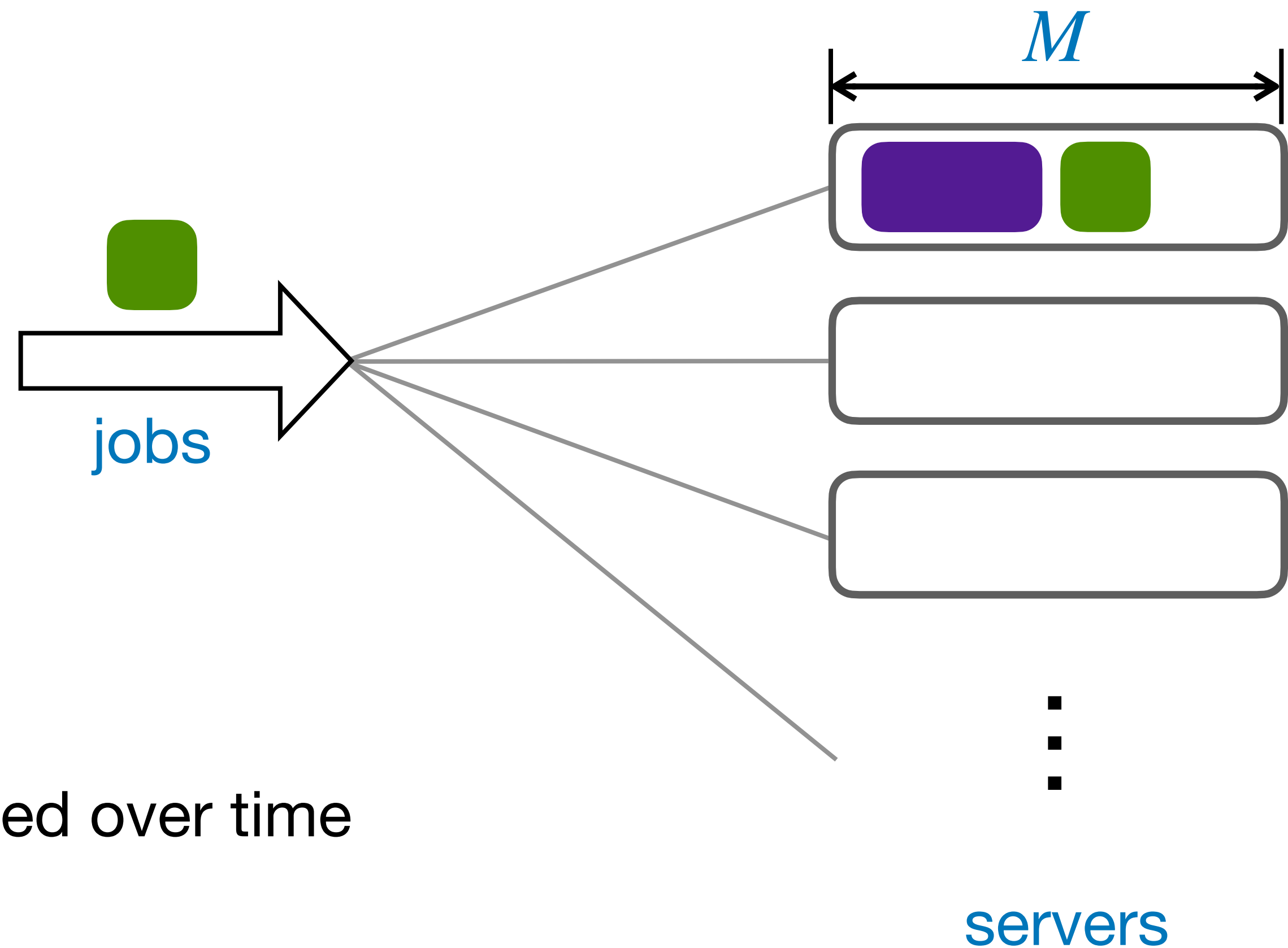


Traditional job model:

- Each job's resource requirement is fixed over time
- Each job departs after a random time

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**

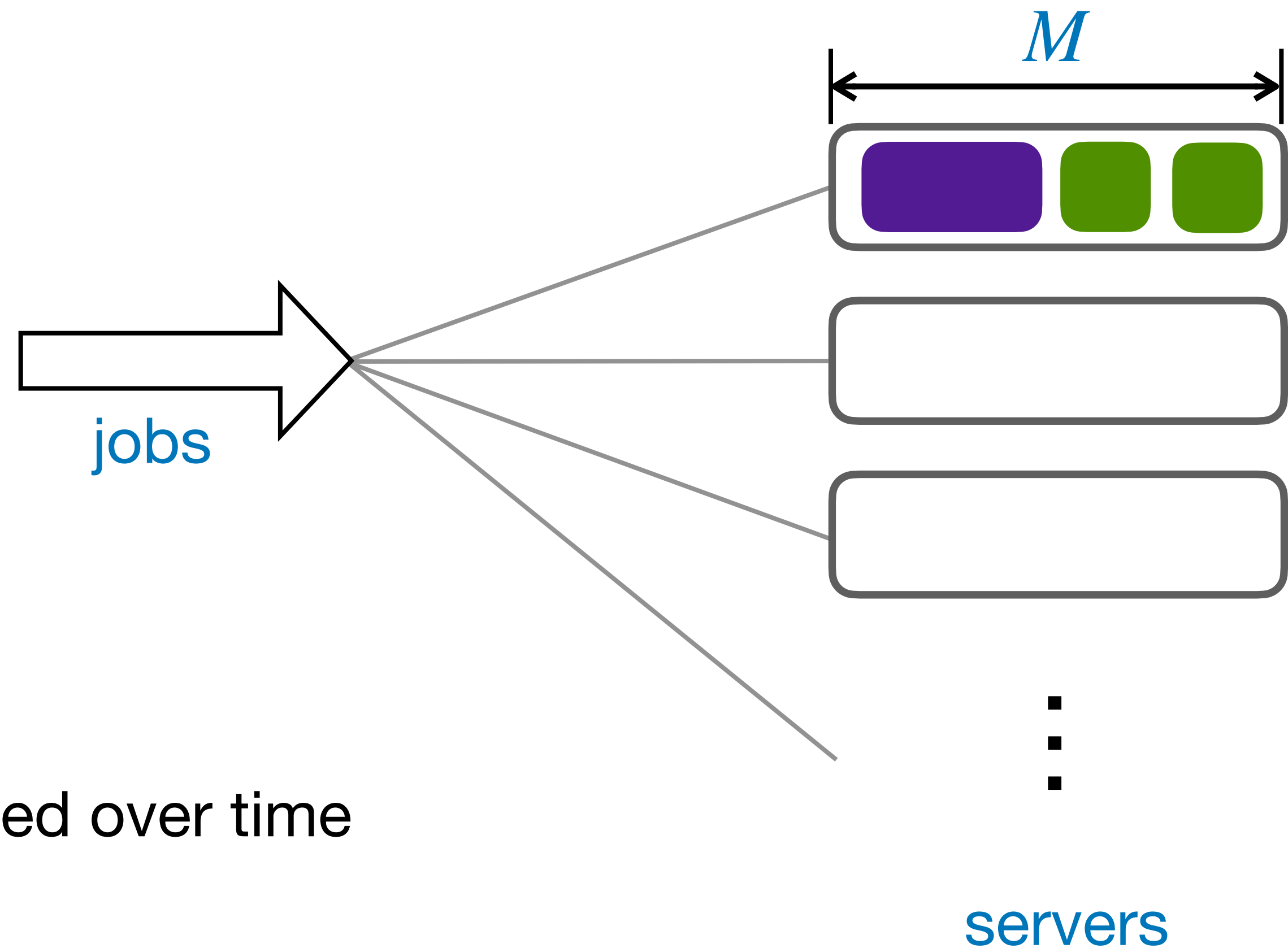


Traditional job model:

- Each job's resource requirement is fixed over time
- Each job departs after a random time

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**

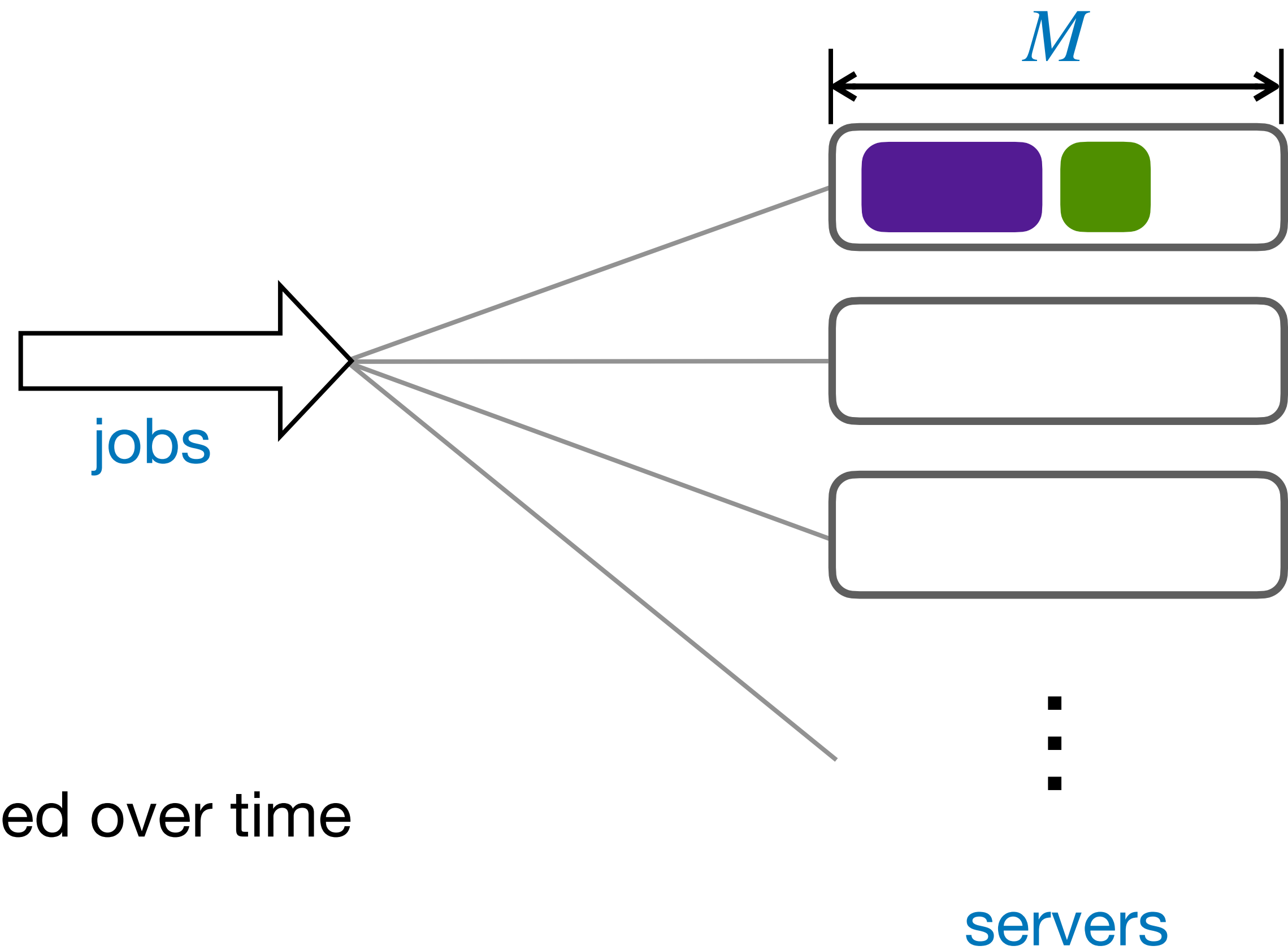


Traditional job model:

- Each job's resource requirement is fixed over time
- Each job departs after a random time

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**

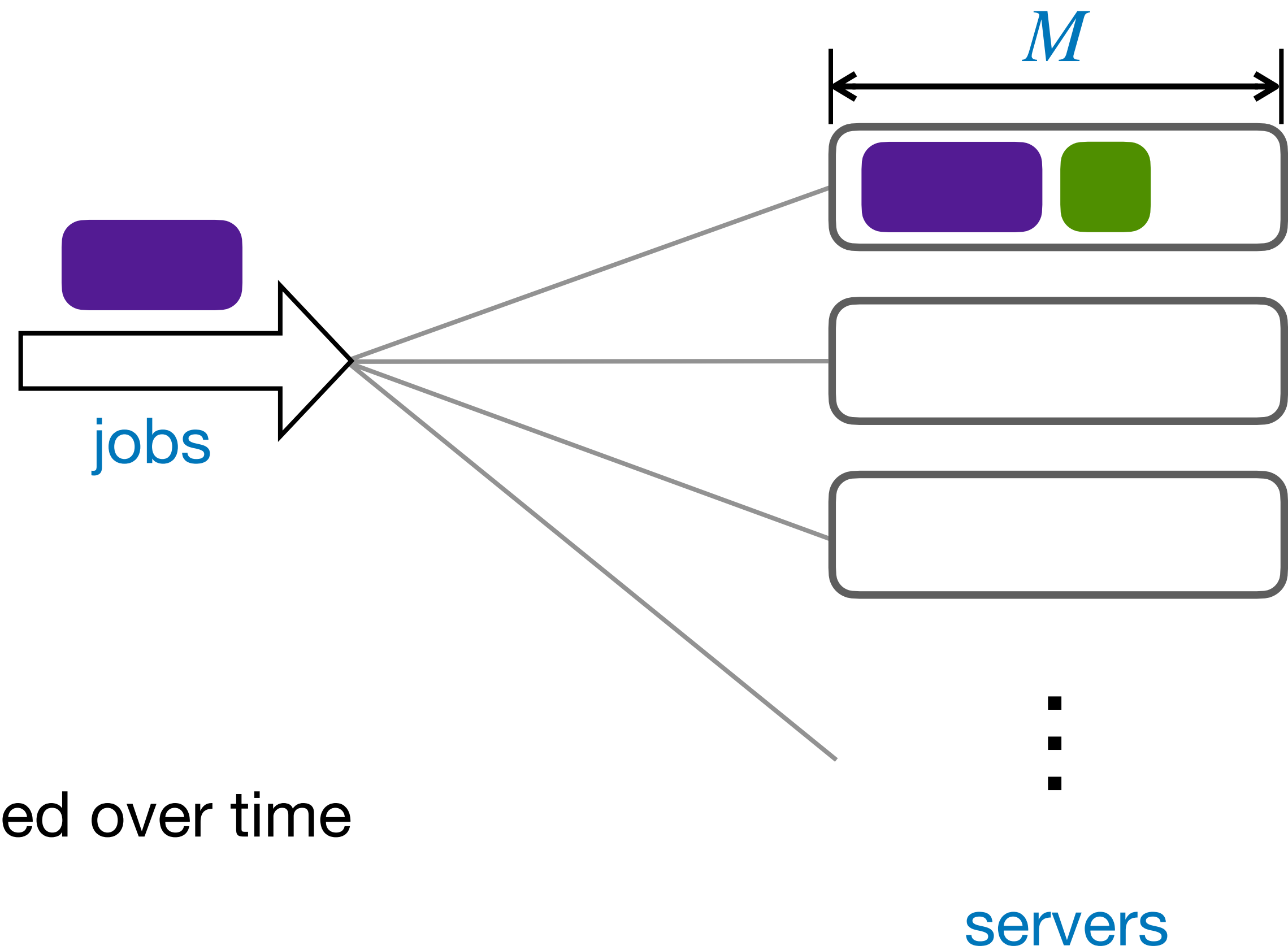


Traditional job model:

- Each job's resource requirement is fixed over time
- Each job departs after a random time

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**

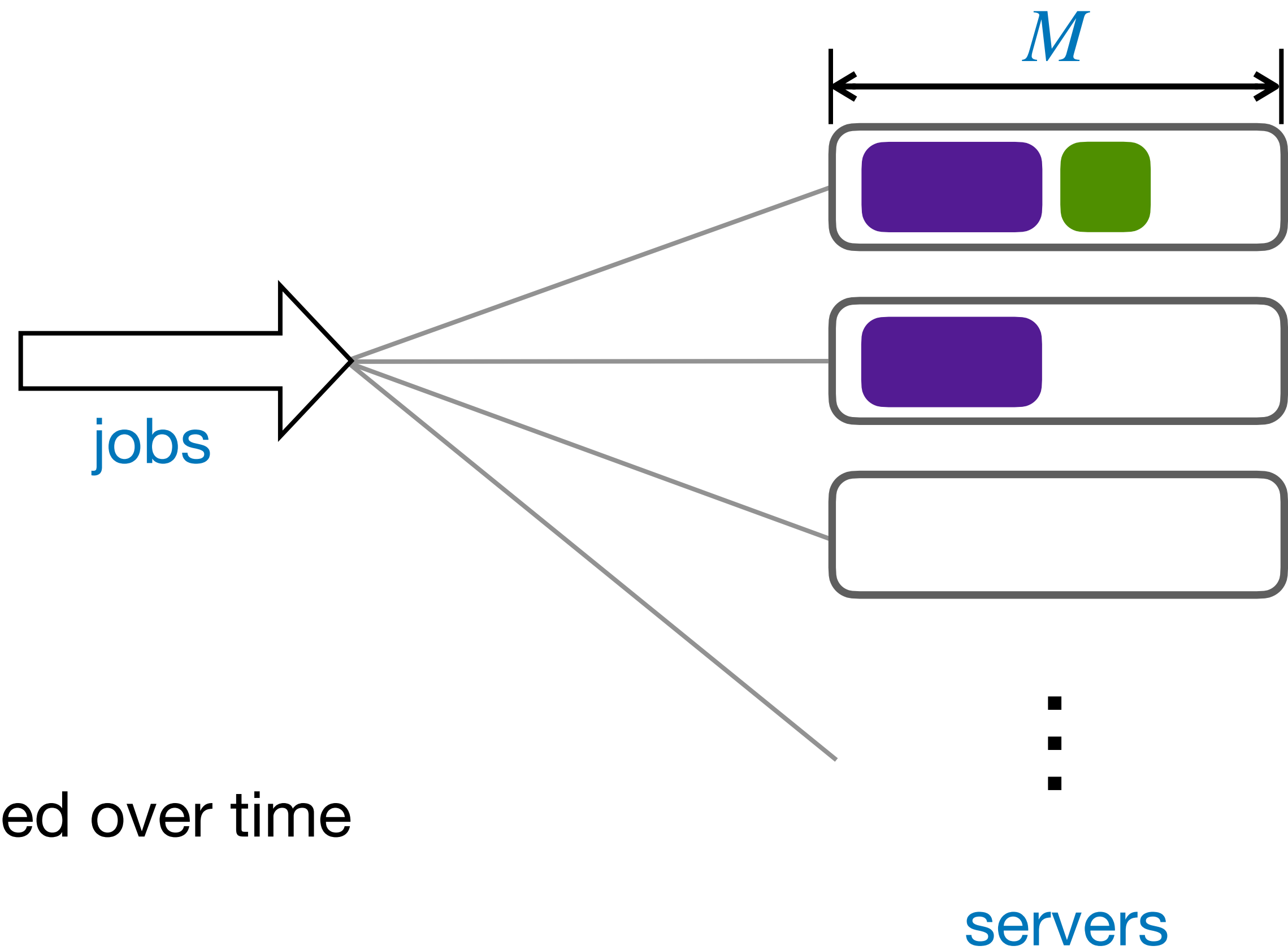


Traditional job model:

- Each job's resource requirement is fixed over time
- Each job departs after a random time

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**

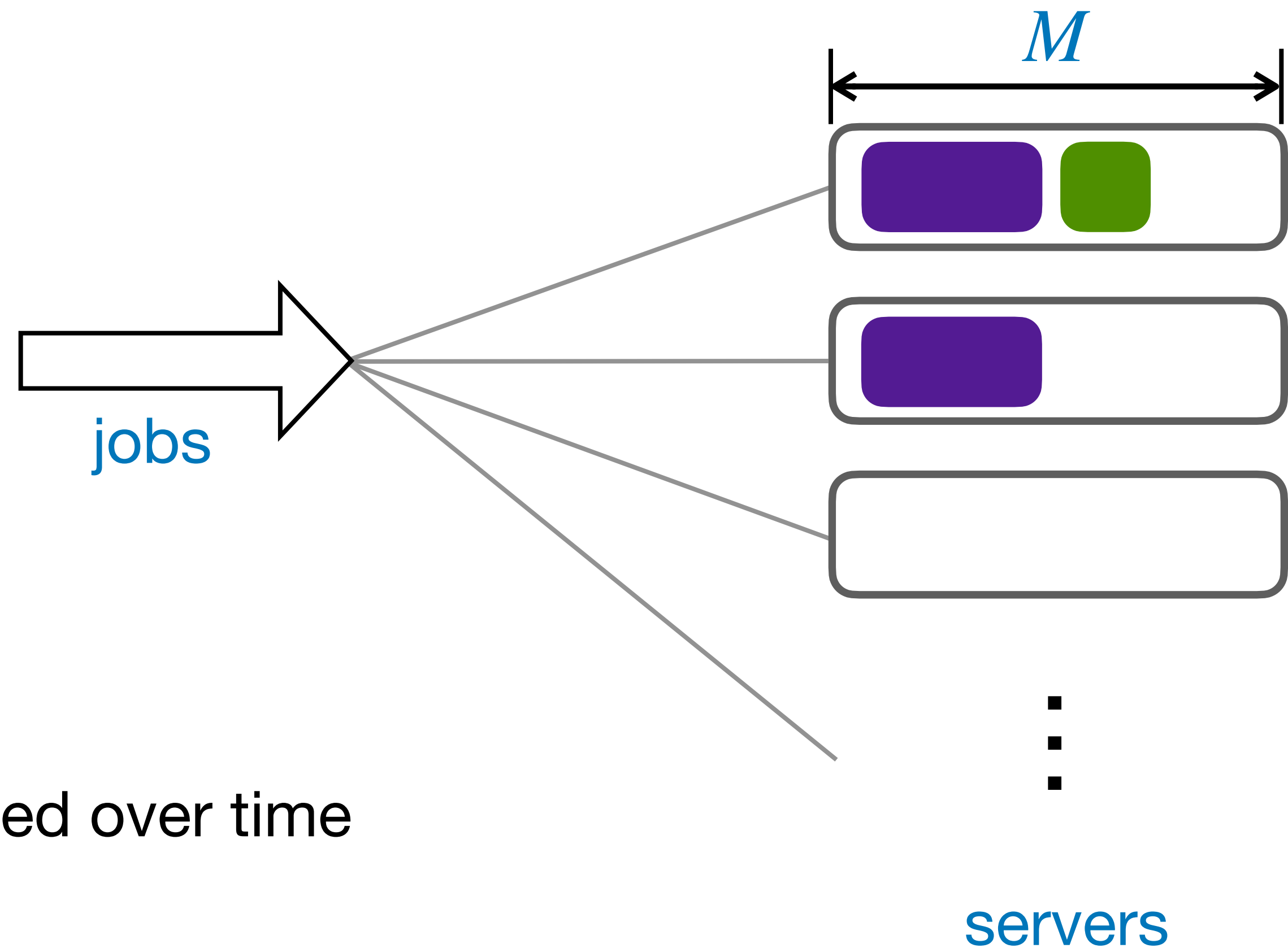


Traditional job model:

- Each job's resource requirement is fixed over time
- Each job departs after a random time

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**



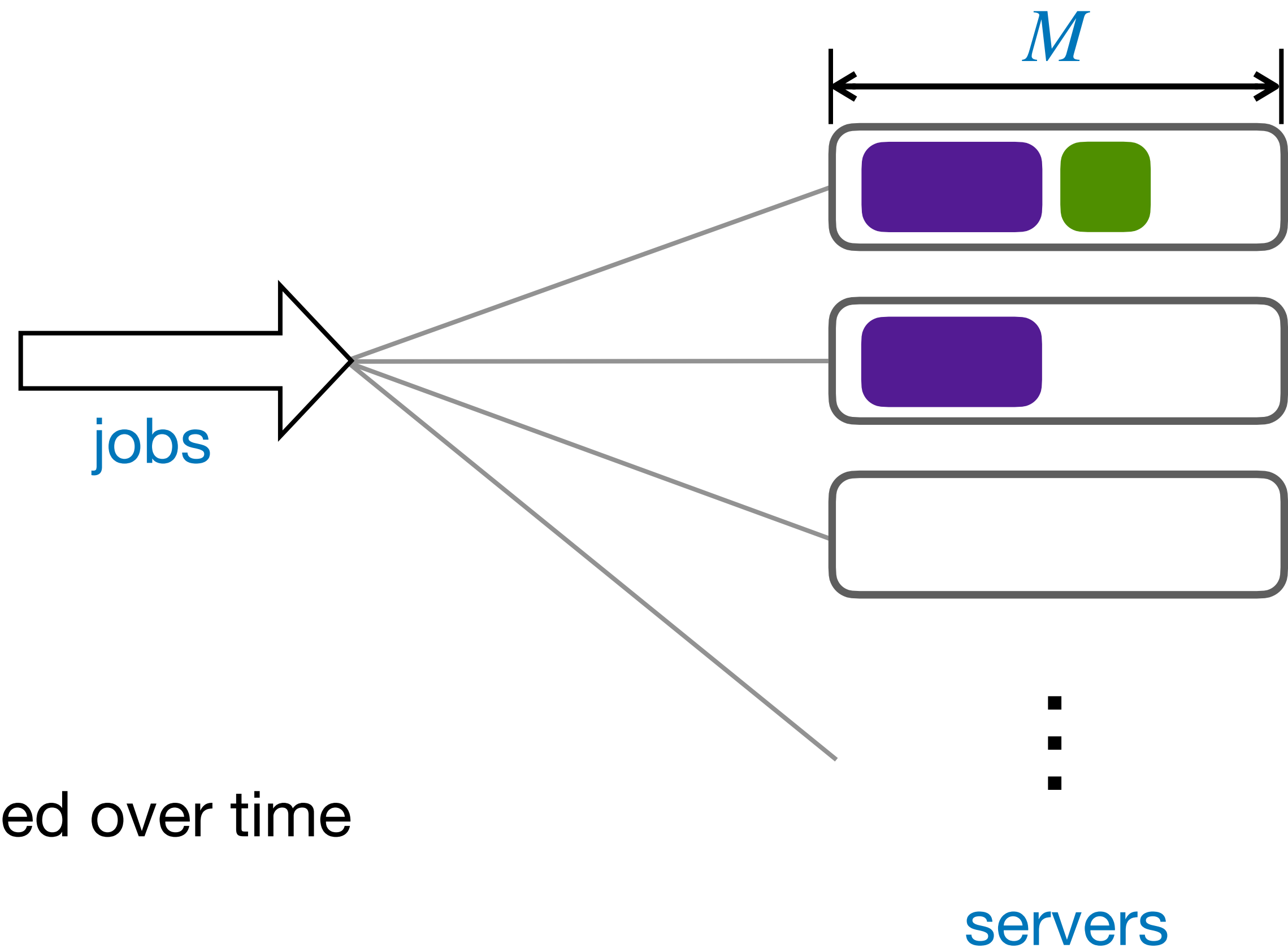
Traditional job model:

- Each job's resource requirement is fixed over time
- Each job departs after a random time

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

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**



Traditional job model:

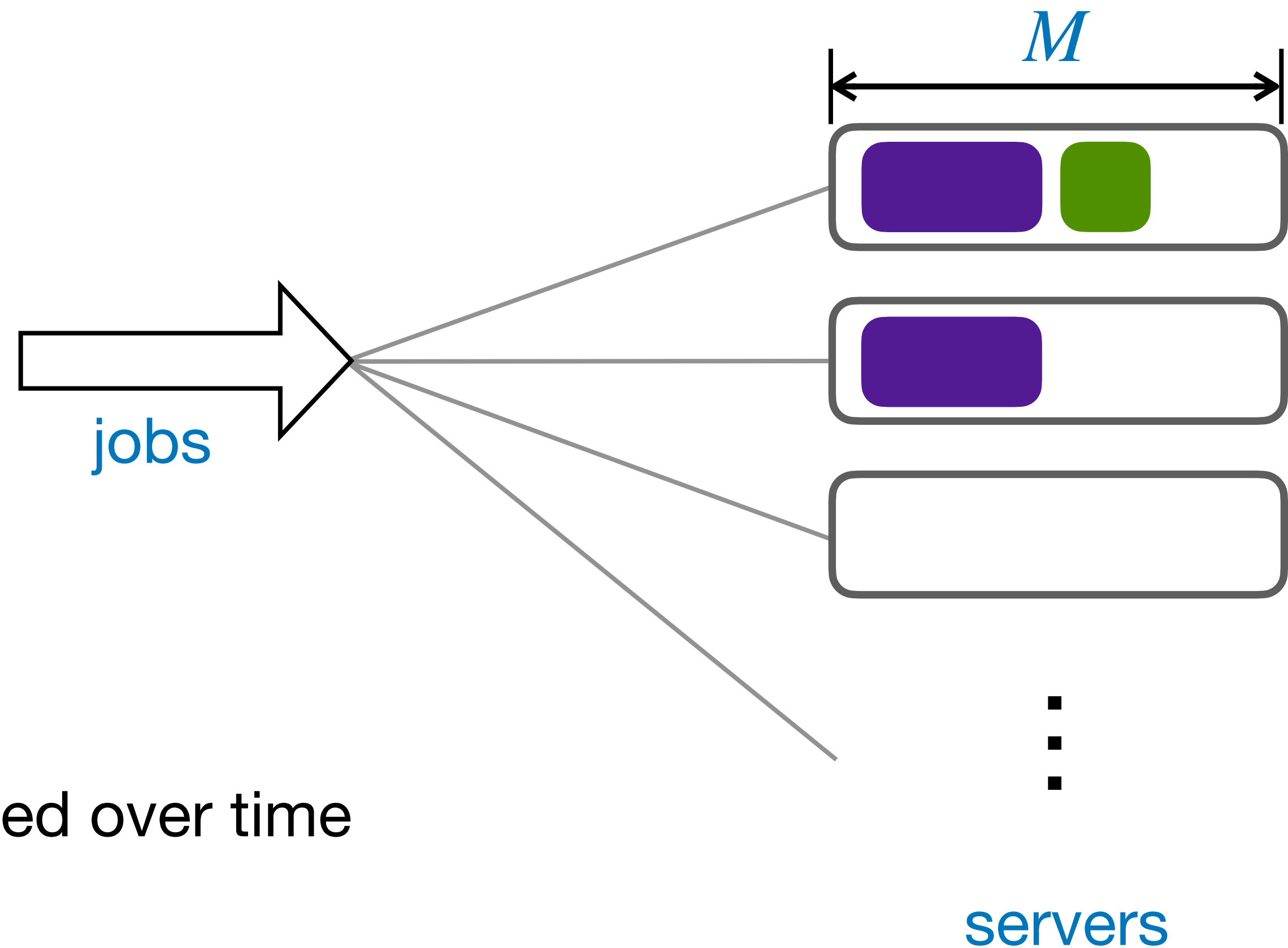
- Each job's resource requirement is fixed over time
- Each job departs after a random time

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

**“Stochastic bin-packing
in service systems”**

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**



A new job model:

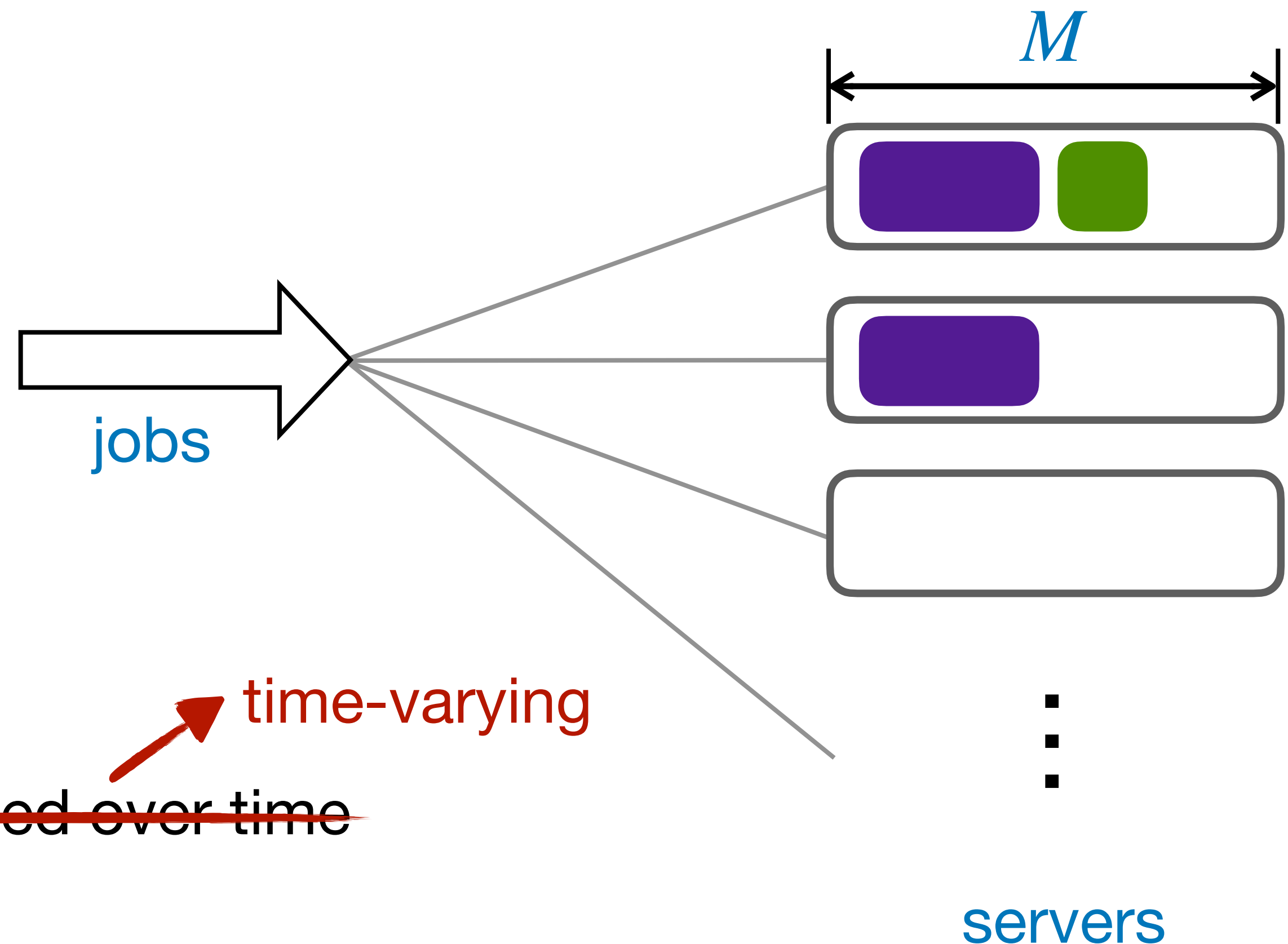
- Each job's resource requirement is fixed over time
- Each job departs after a random time

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

**“Stochastic bin-packing
in service systems”**

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**



A new job model:

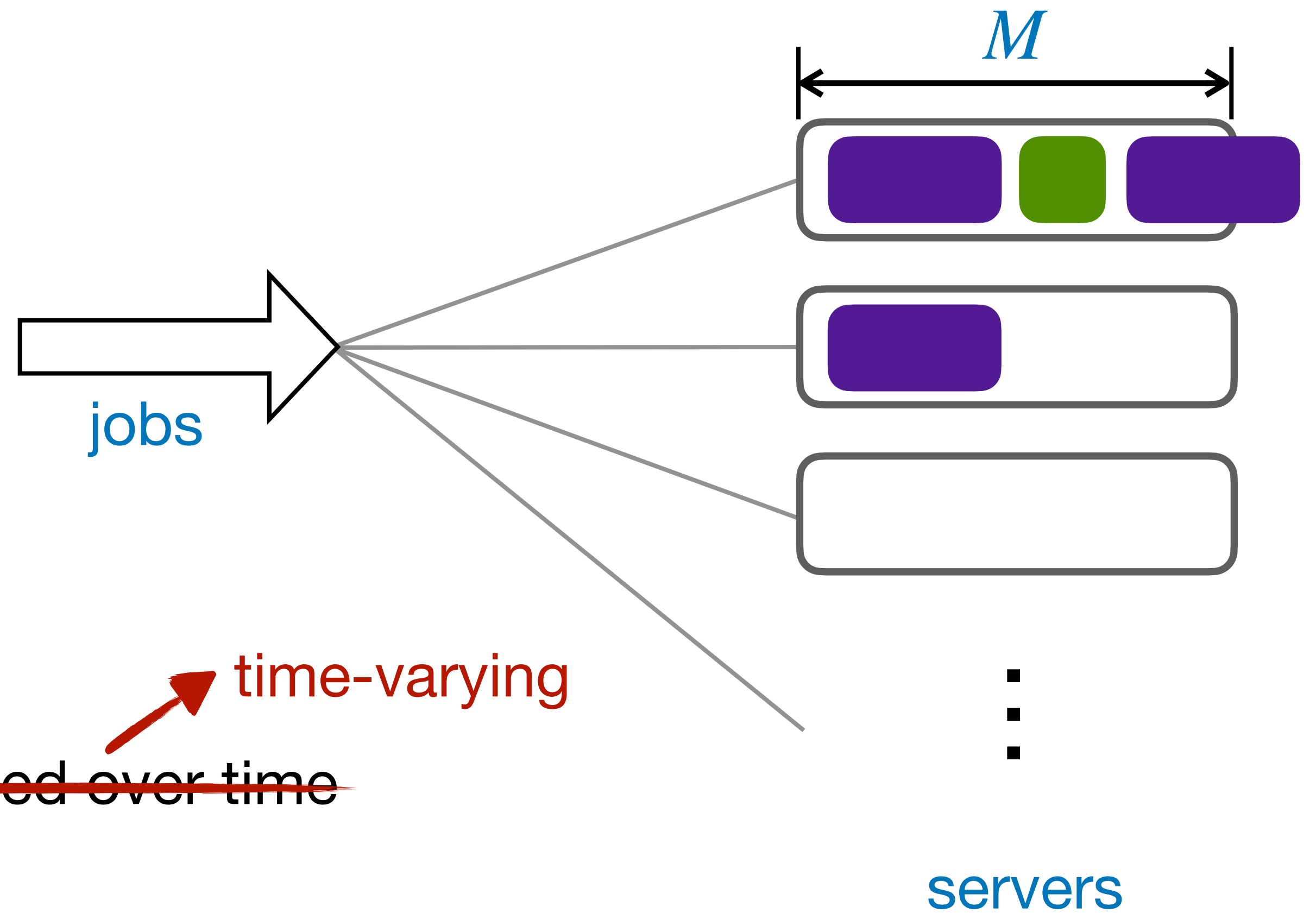
- Each job's resource requirement is ~~fixed over time~~
- Each job departs after a random time

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

**“Stochastic bin-packing
in service systems”**

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**



A new job model:

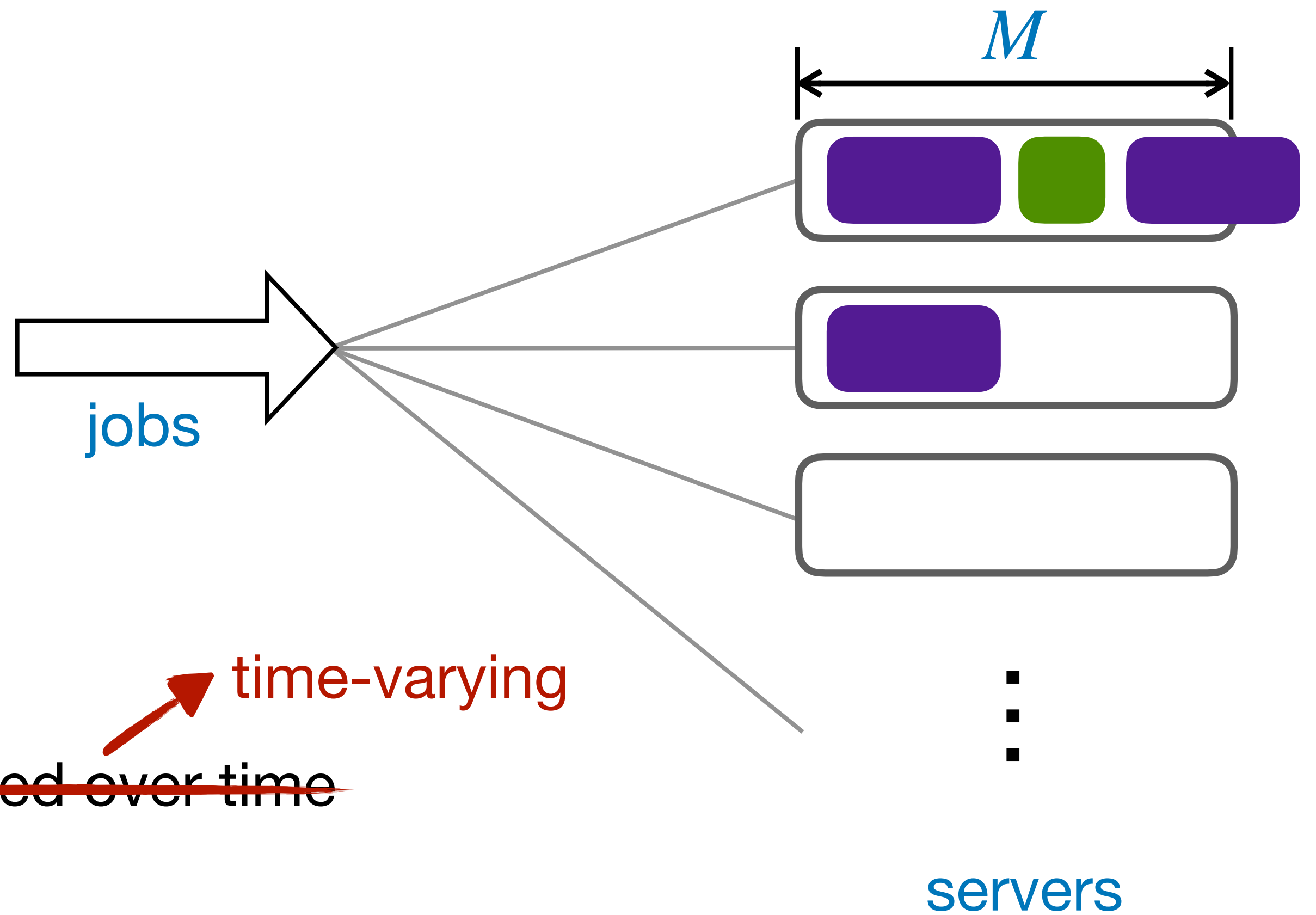
- Each job's resource requirement is ~~fixed over time~~
- Each job departs after a random time

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

**“Stochastic bin-packing
in service systems”**

The problem

- Each arriving **job** needs to be assigned to a **server**
- Each **server** has a resource capacity M
- Infinite # **servers**



A new job model:

- Each job's resource requirement is ~~fixed over time~~ **time-varying**
- Each job departs after a random time

Goal:

minimize
job assigning policy

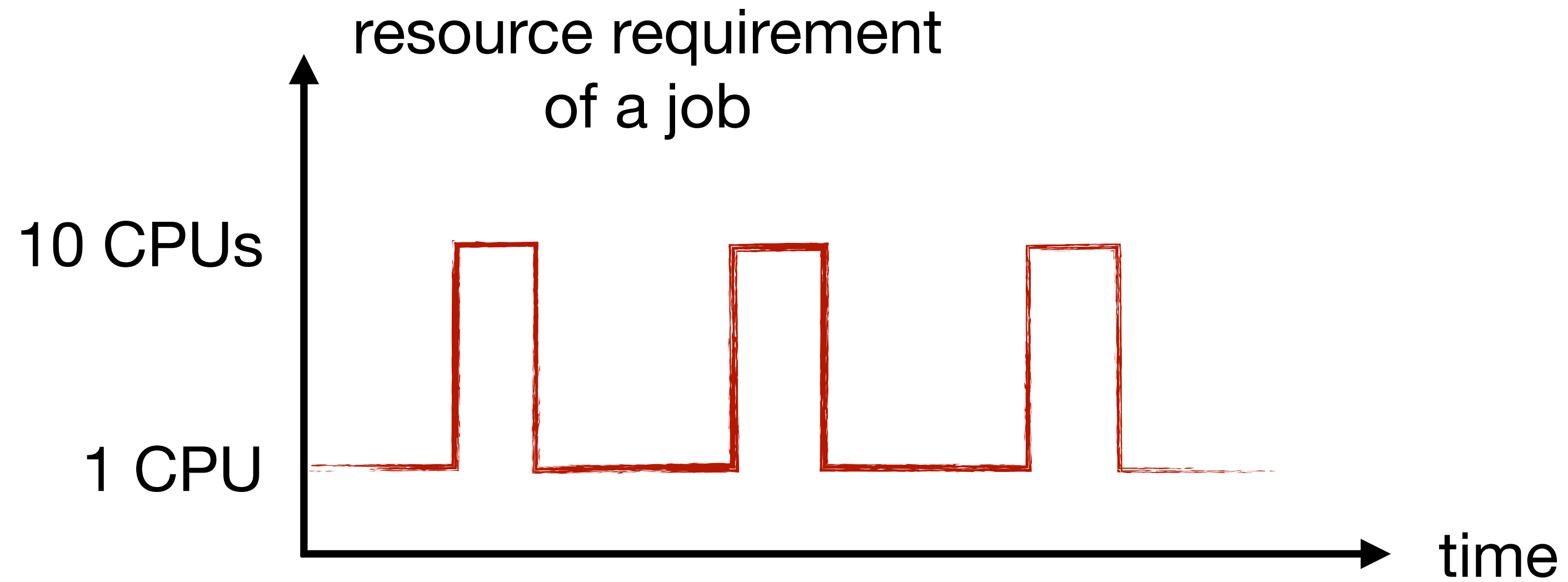
E [# active servers]

**“Stochastic bin-packing
in service systems”**

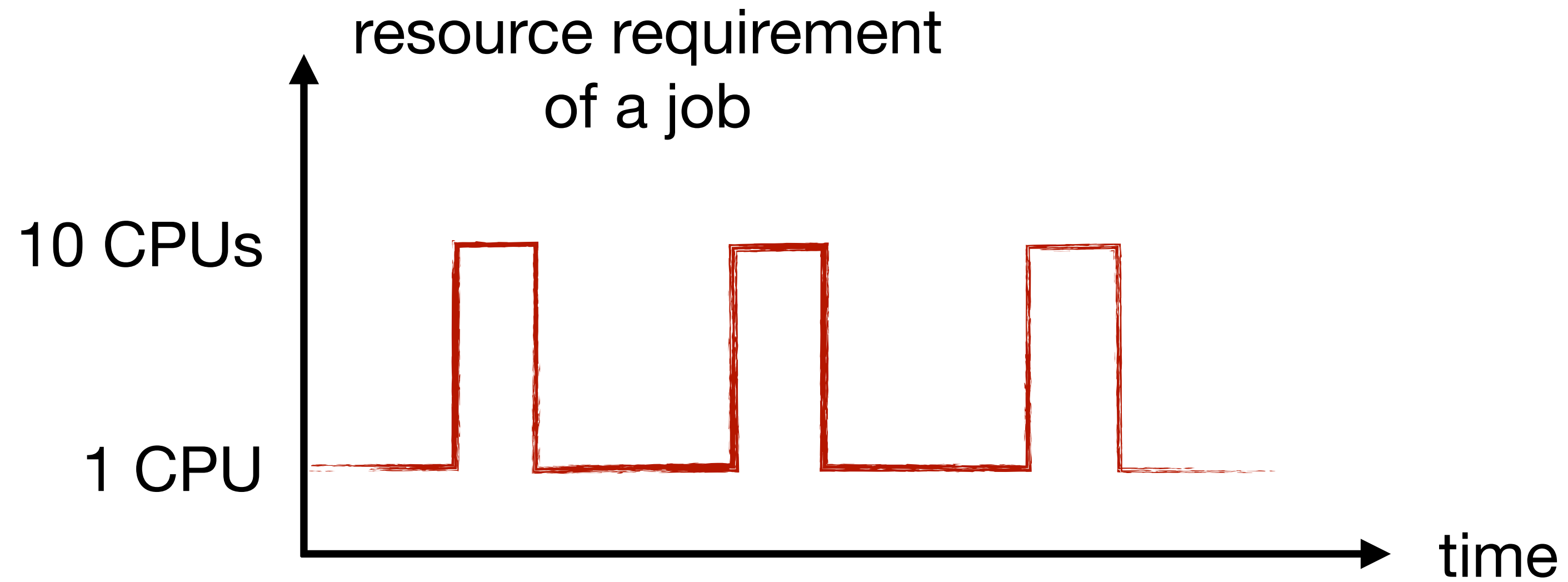
subject to

cost (resource contention) \leq budget

Why time-varying?

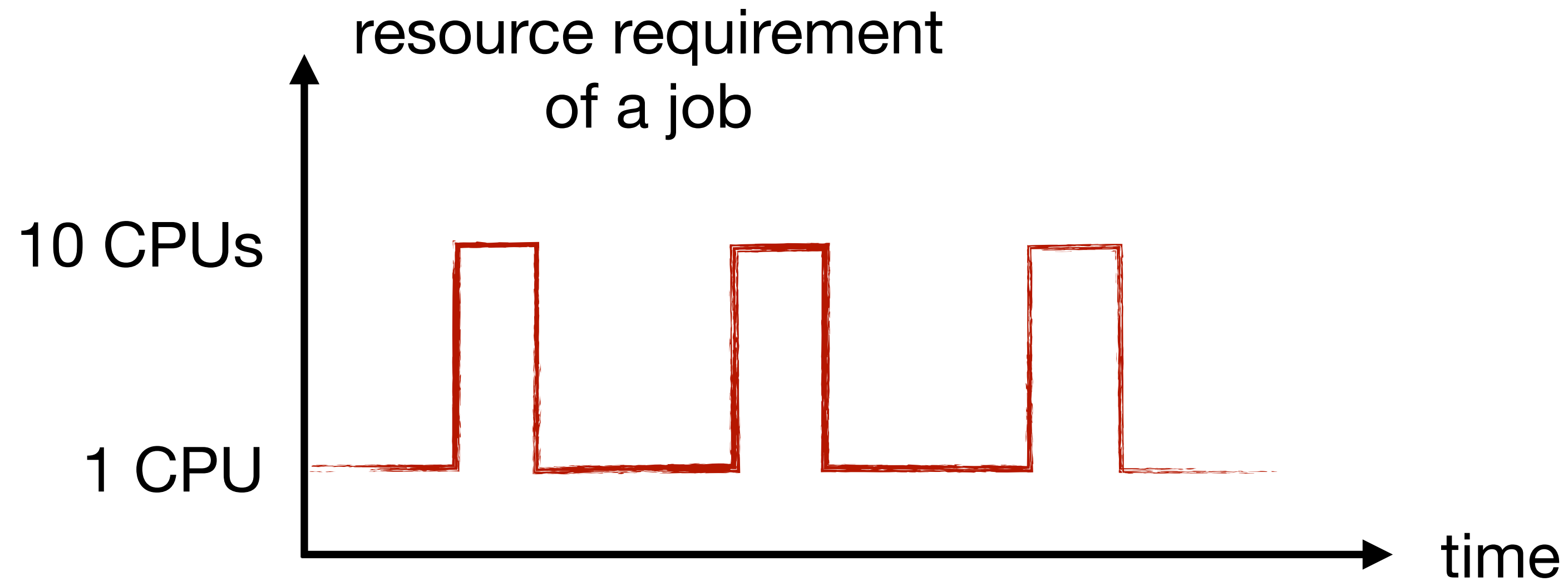


Why time-varying?



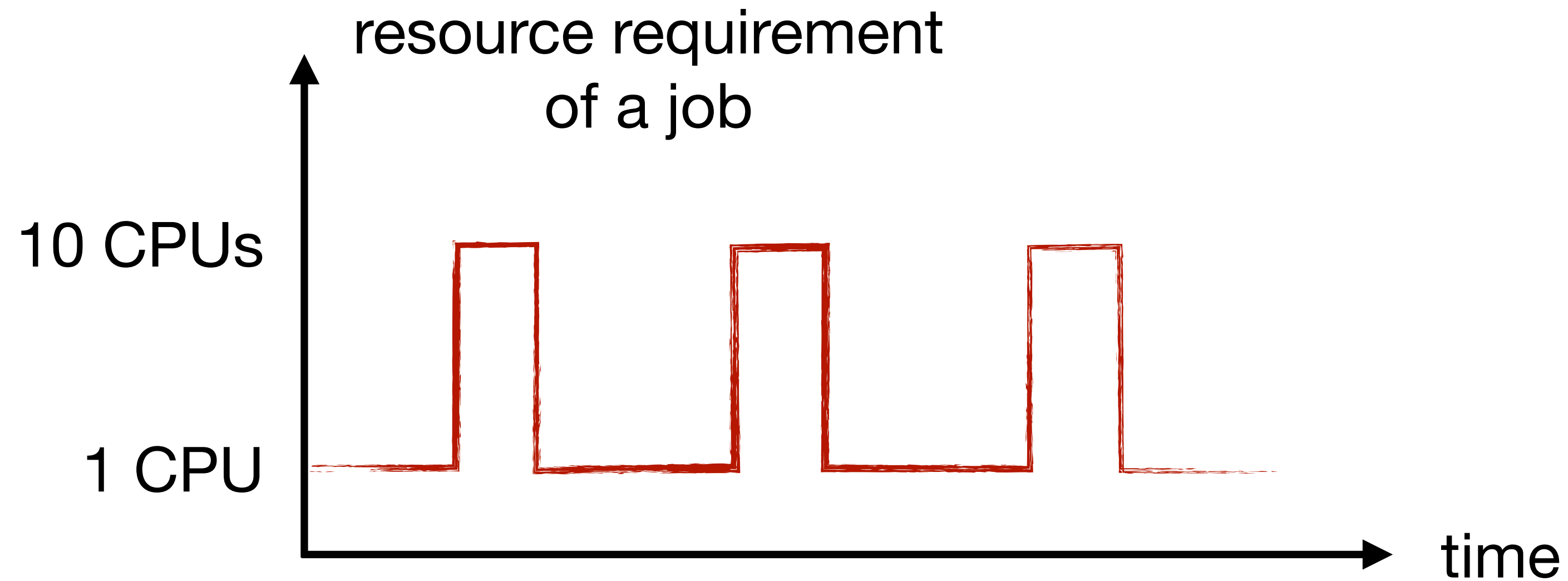
- Reserve resources based on peak requirement

Why time-varying?



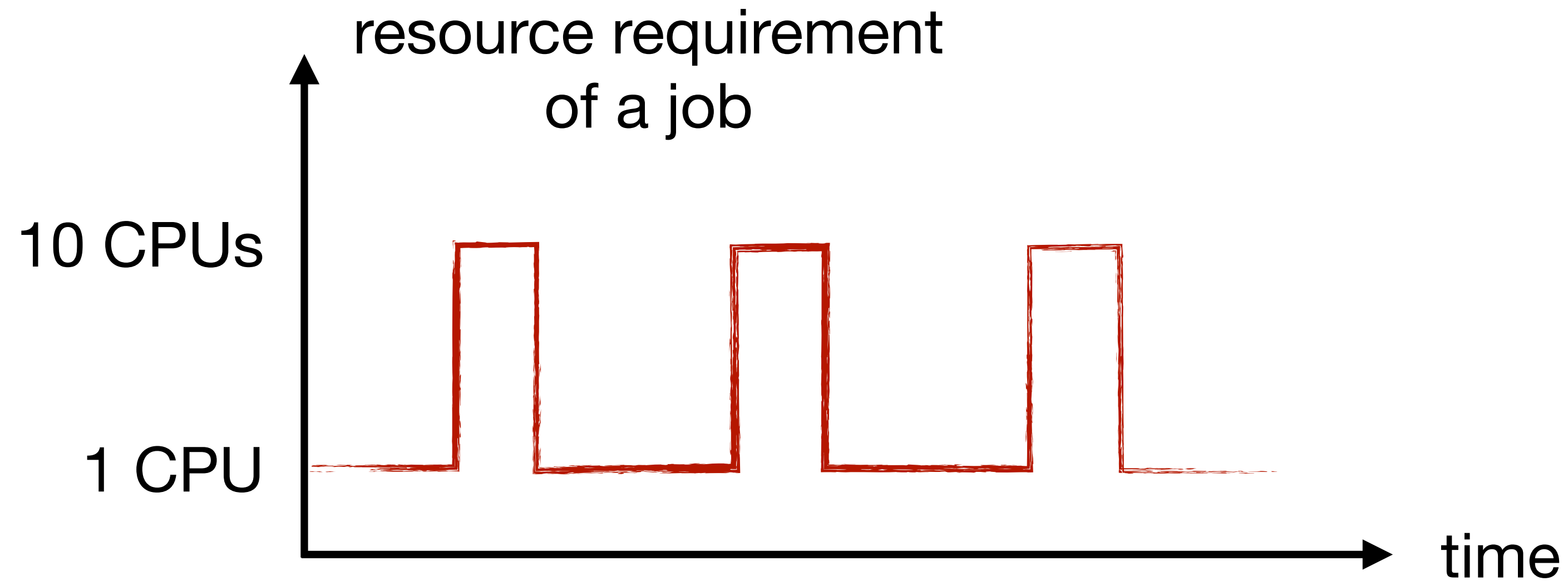
- Reserve resources based on peak requirement
➔ low resource utilization on a server

Why time-varying?



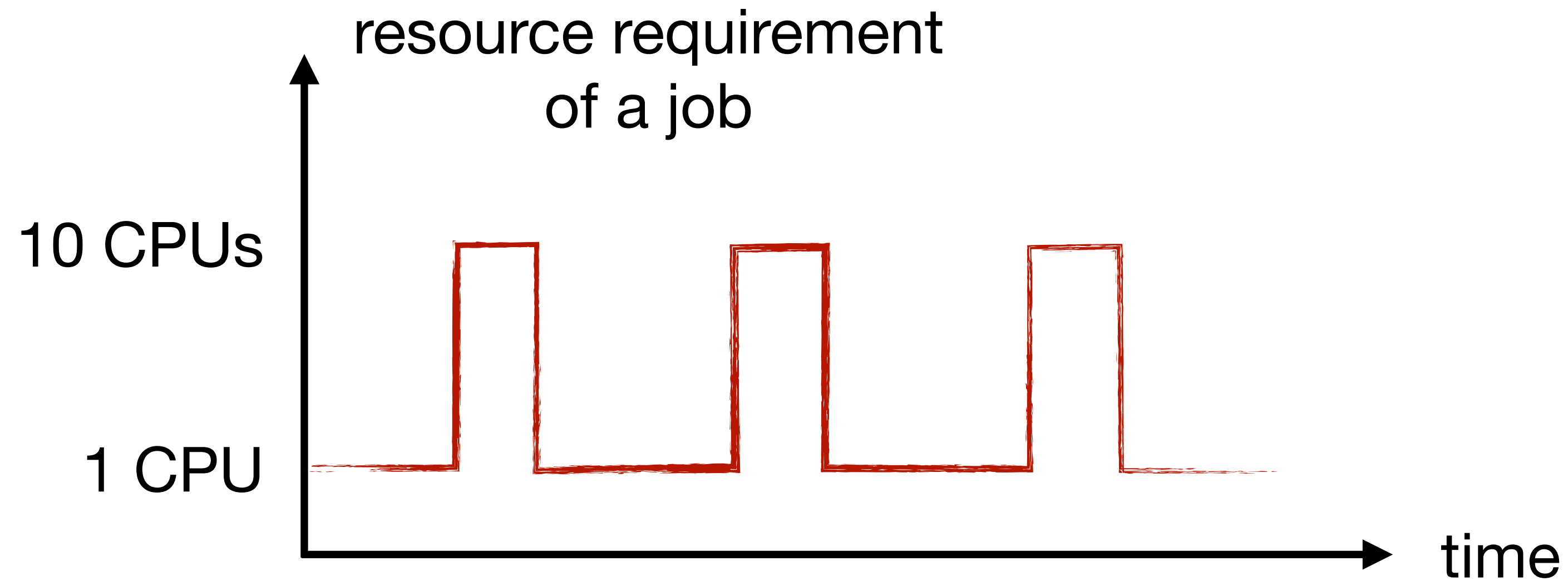
- Reserve resources based on peak requirement
 - ➔ low resource utilization on a server
 - ➔ larger # active servers

Why time-varying?



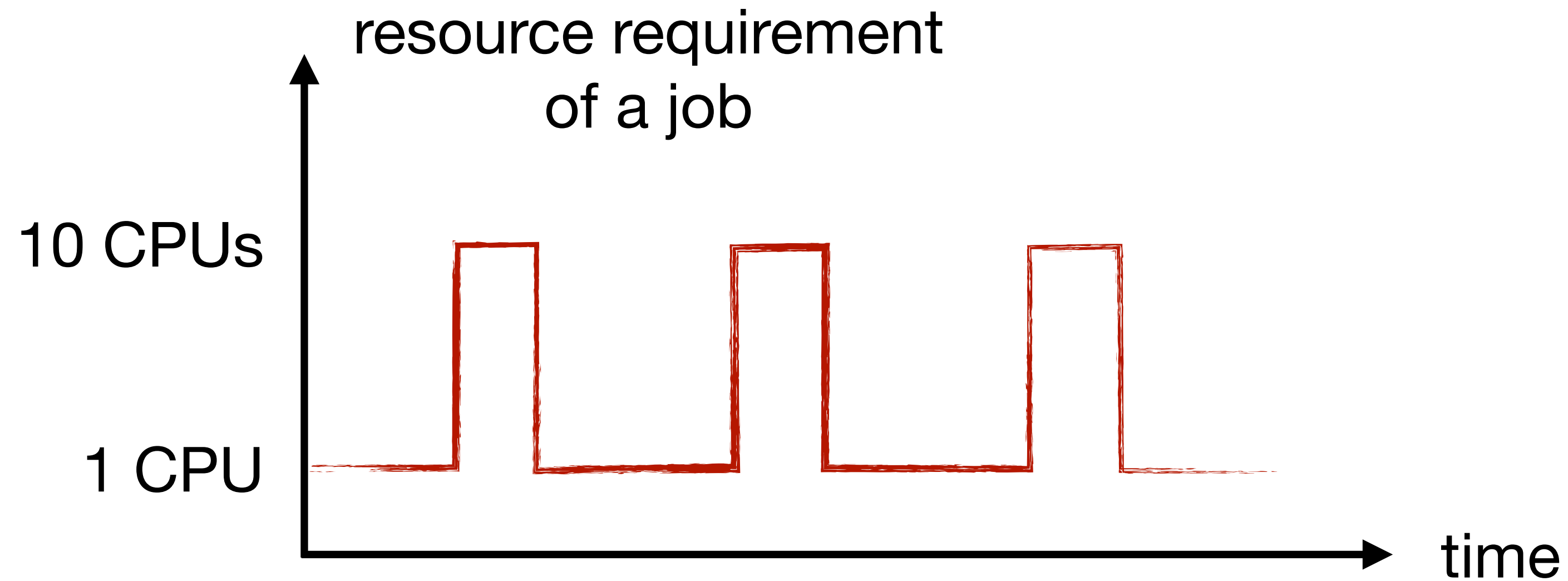
- Reserve resources based on peak requirement
 - ➔ low resource utilization on a server
 - ➔ larger # active servers
- Overcommit resources on a server

Why time-varying?



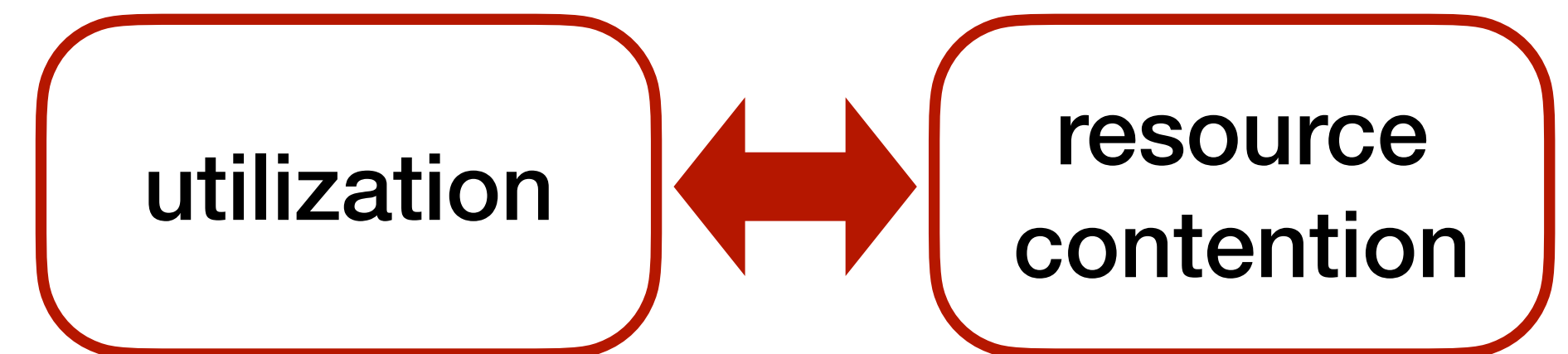
- Reserve resources based on peak requirement
 - ➔ low resource utilization on a server
 - ➔ larger # active servers
- Overcommit resources on a server
 - ➔ possible resource contention

Why time-varying?



- 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:



What's known?

What's known?

- **Stochastic bin-packing in service systems**

What's known?

- **Stochastic bin-packing in service systems**

- Traditional job model: Asymptotic optimality, no convergence rate

[Stolyar 2013], [Stolyar and Zhong 2013, 2015, 2021], [Ghaderi, Zhong, and Srikant 2014], [Stolyar 2017]

What's known?

- **Stochastic bin-packing in service systems**

- Traditional job model: Asymptotic optimality, no convergence rate

[Stolyar 2013], [Stolyar and Zhong 2013, 2015, 2021], [Ghaderi, Zhong, and Srikant 2014], [Stolyar 2017]

- Finite-server model: Maximizing throughput/reward, heavy-traffic optimality, loss model

[Maguluri, Srikant, and Ying 2012], [Maguluri and Srikant 2013], [Xie et al. 2015], [Ghaderi 2016], [Psychas and Ghaderi 2017, 2018, 2019, 2021, 2021], ...

What's known?

- **Stochastic bin-packing in service systems**

- Traditional job model: Asymptotic optimality, no convergence rate

[Stolyar 2013], [Stolyar and Zhong 2013, 2015, 2021], [Ghaderi, Zhong, and Srikant 2014], [Stolyar 2017]

- Finite-server model: Maximizing throughput/reward, heavy-traffic optimality, loss model

[Maguluri, Srikant, and Ying 2012], [Maguluri and Srikant 2013], [Xie et al. 2015], [Ghaderi 2016], [Psychas and Ghaderi 2017, 2018, 2019, 2021, 2021], ...

- **Stochastic bin-packing without job departures**

[Courcoubetis and Weber 1986, 1990], [Csirik et al. 2006], [Freund and Banerjee 2019], [Gupta and Radovanović 2020], ...

What's known?

- **Stochastic bin-packing in service systems**

- Traditional job model: Asymptotic optimality, no convergence rate

[Stolyar 2013], [Stolyar and Zhong 2013, 2015, 2021], [Ghaderi, Zhong, and Srikant 2014], [Stolyar 2017]

- Finite-server model: Maximizing throughput/reward, heavy-traffic optimality, loss model

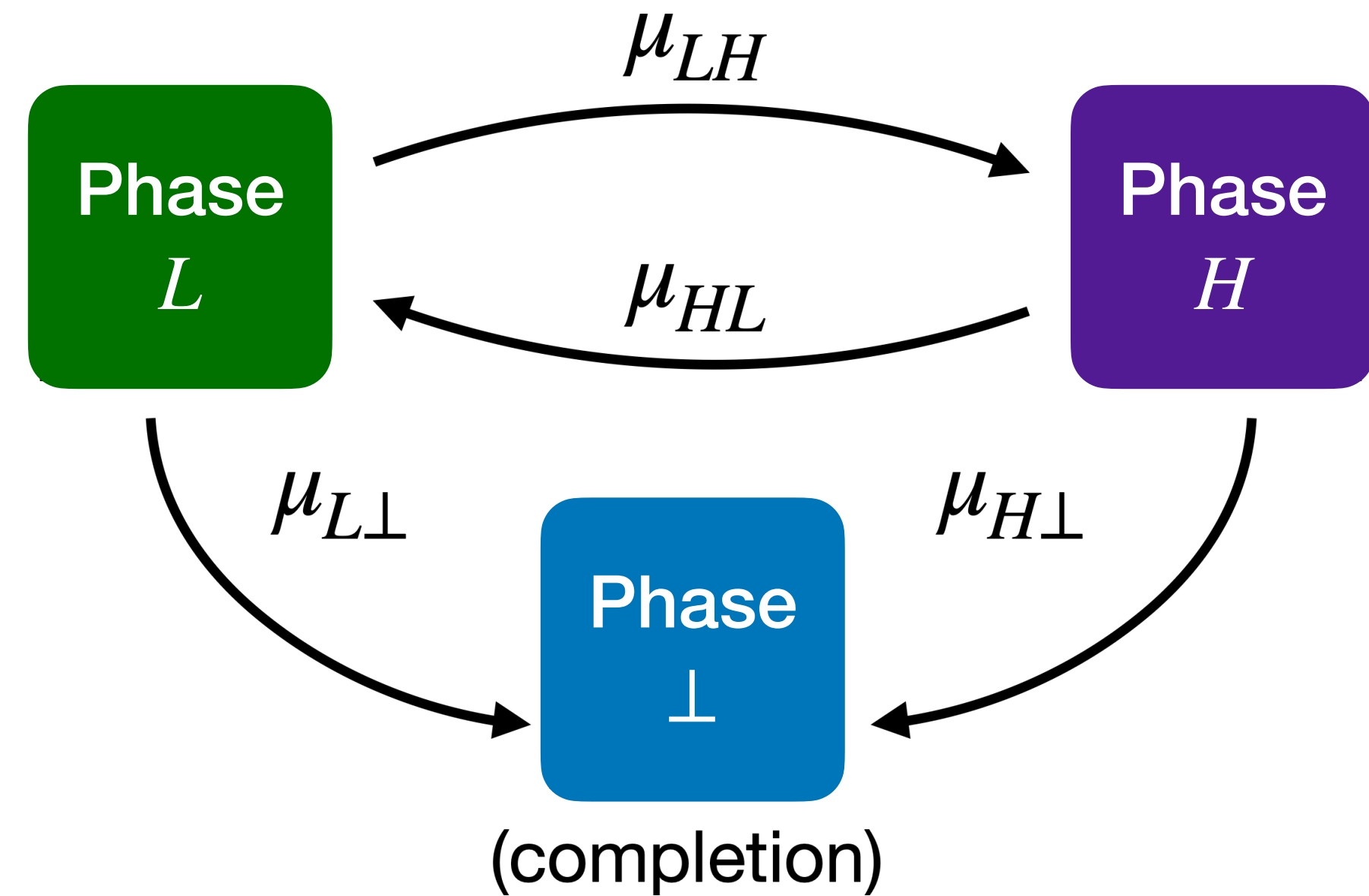
[Maguluri, Srikant, and Ying 2012], [Maguluri and Srikant 2013], [Xie et al. 2015], [Ghaderi 2016], [Psychas and Ghaderi 2017, 2018, 2019, 2021, 2021], ...

- **Stochastic bin-packing without job departures**

[Courcoubetis and Weber 1986, 1990], [Csirik et al. 2006], [Freund and Banerjee 2019], [Gupta and Radovanović 2020], ...

- **Classical bin-packing: Vast literature**

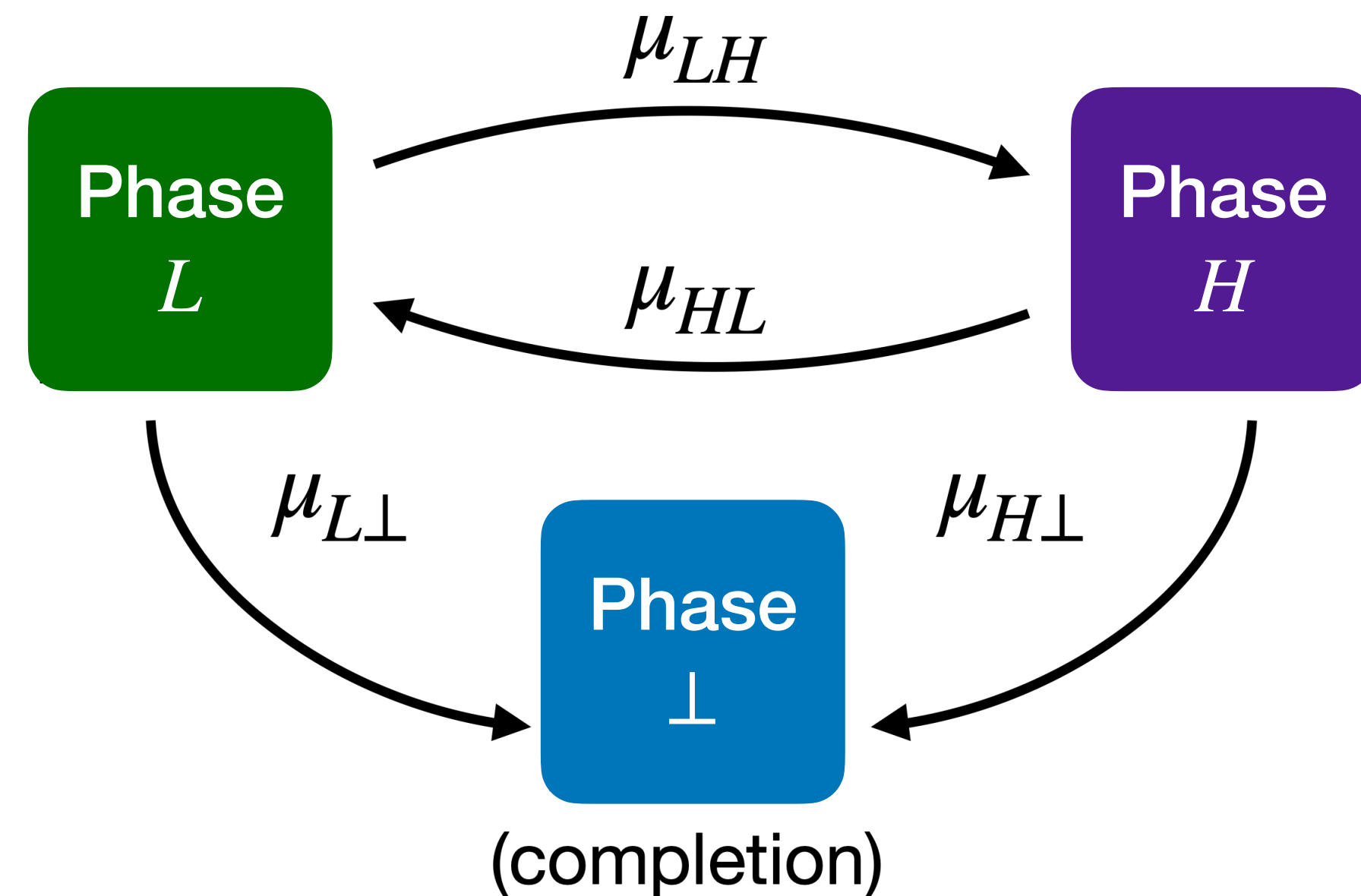
More details on the job model



Example MC

More details on the job model

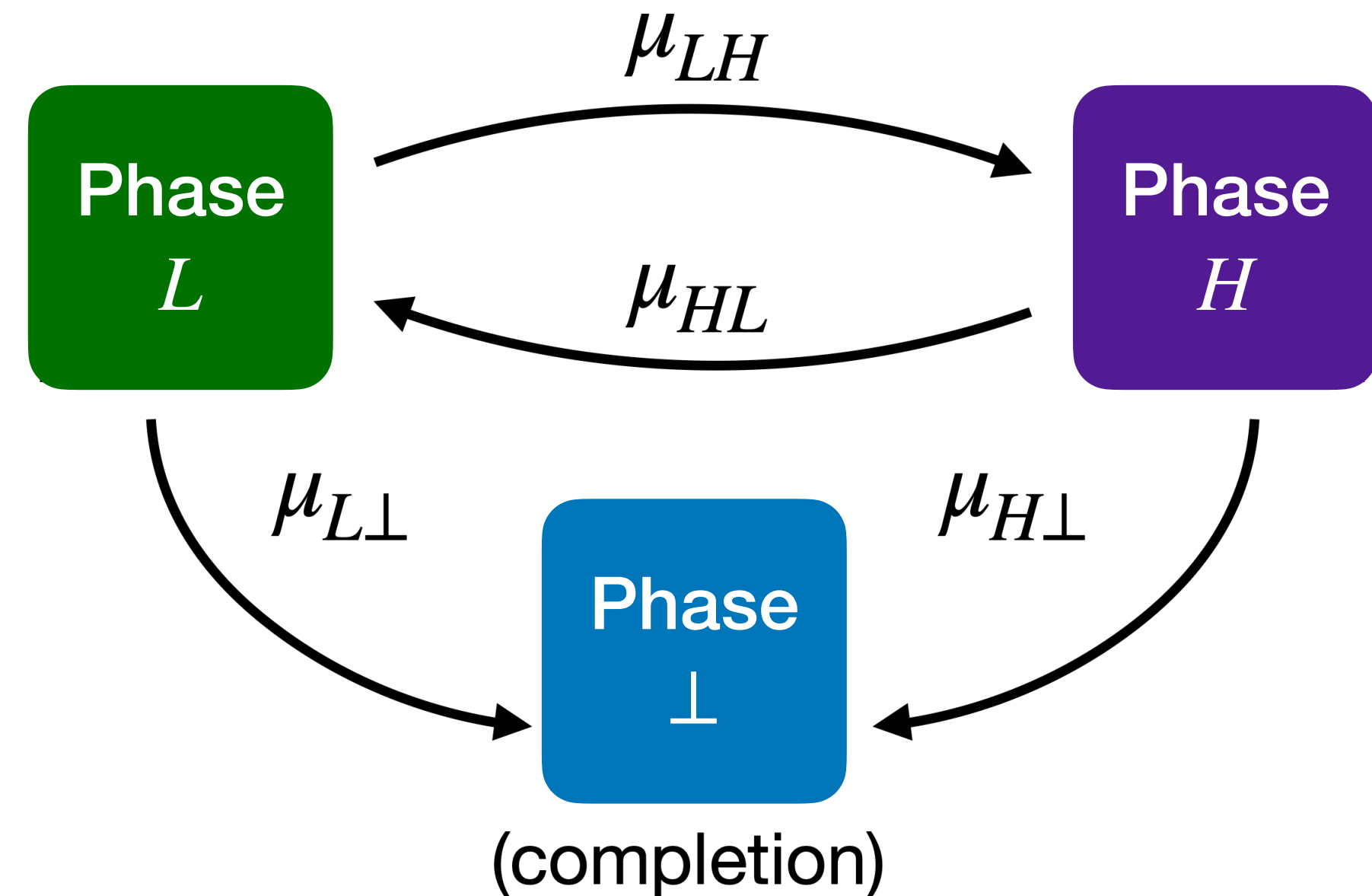
- Resource requirement of a job evolves over time following a Markov chain



Example MC

More details on the job model

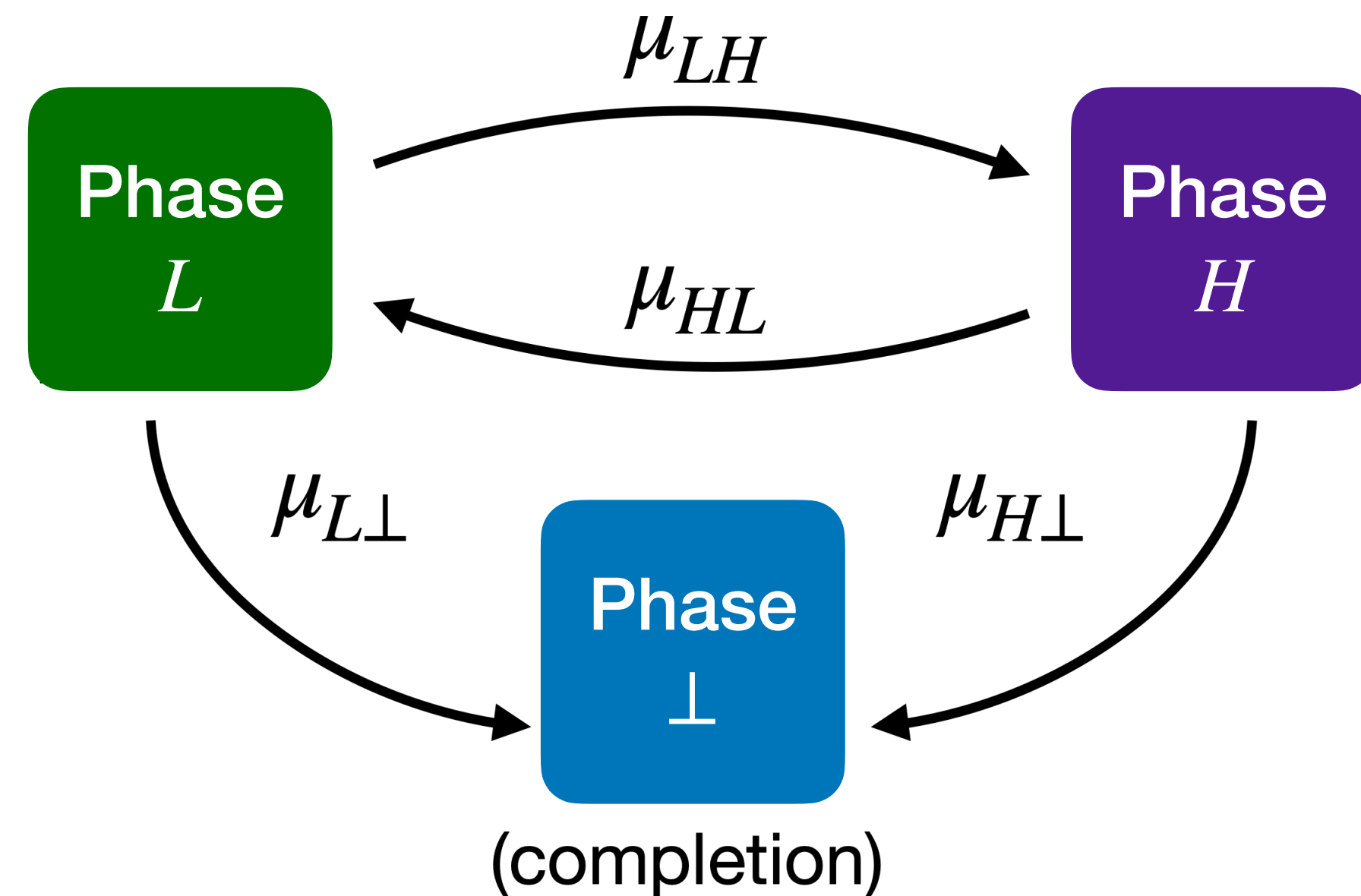
- Resource requirement of a job evolves over time following a Markov chain
- Initial job type follows an initial distribution



Example MC

More details on the job model

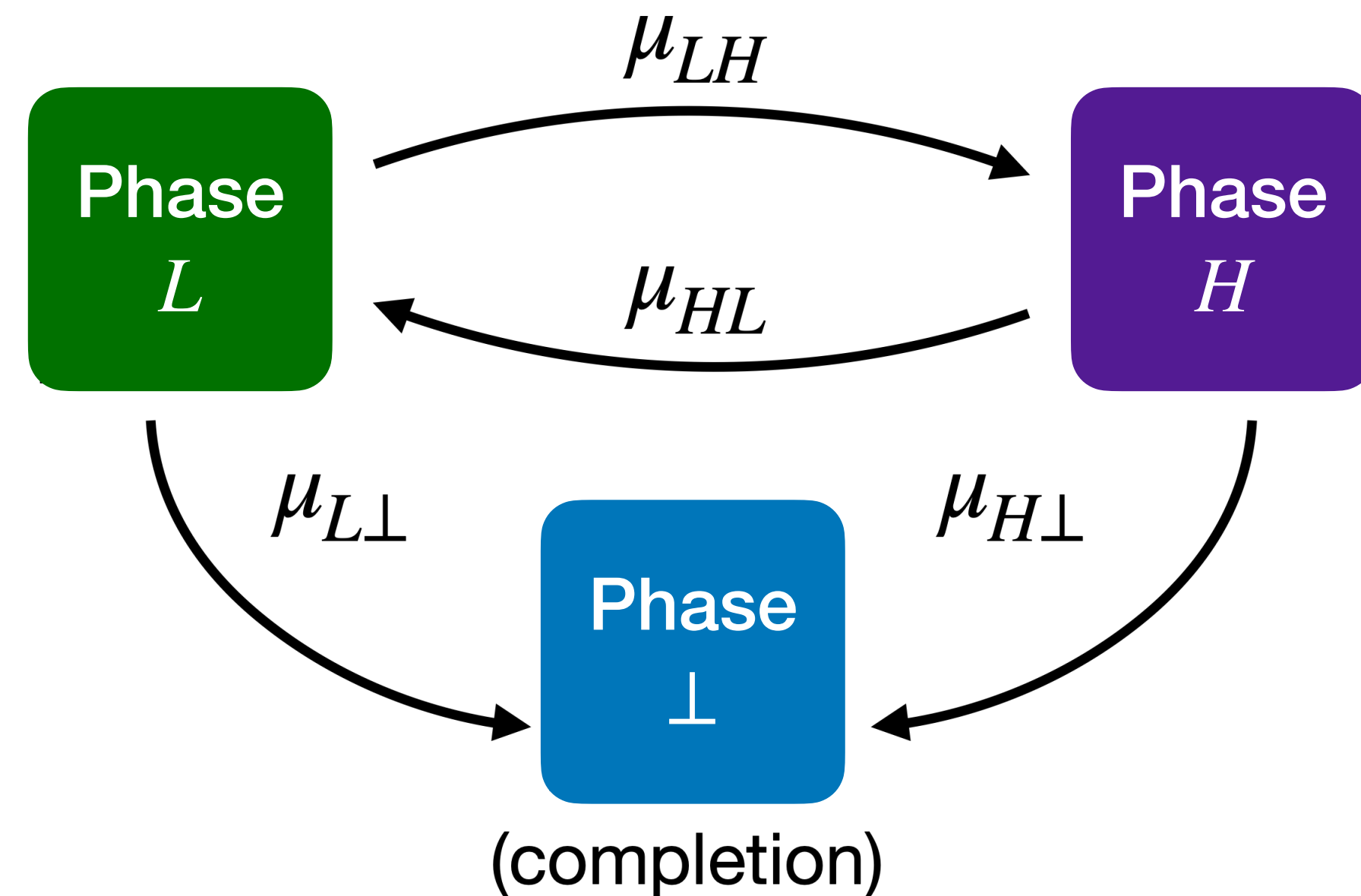
- Resource requirement of a job evolves over time following a Markov chain
- Initial job type follows an initial distribution
- MCs of jobs are independent of each other, and they are exogenous (not affected by resource contention)



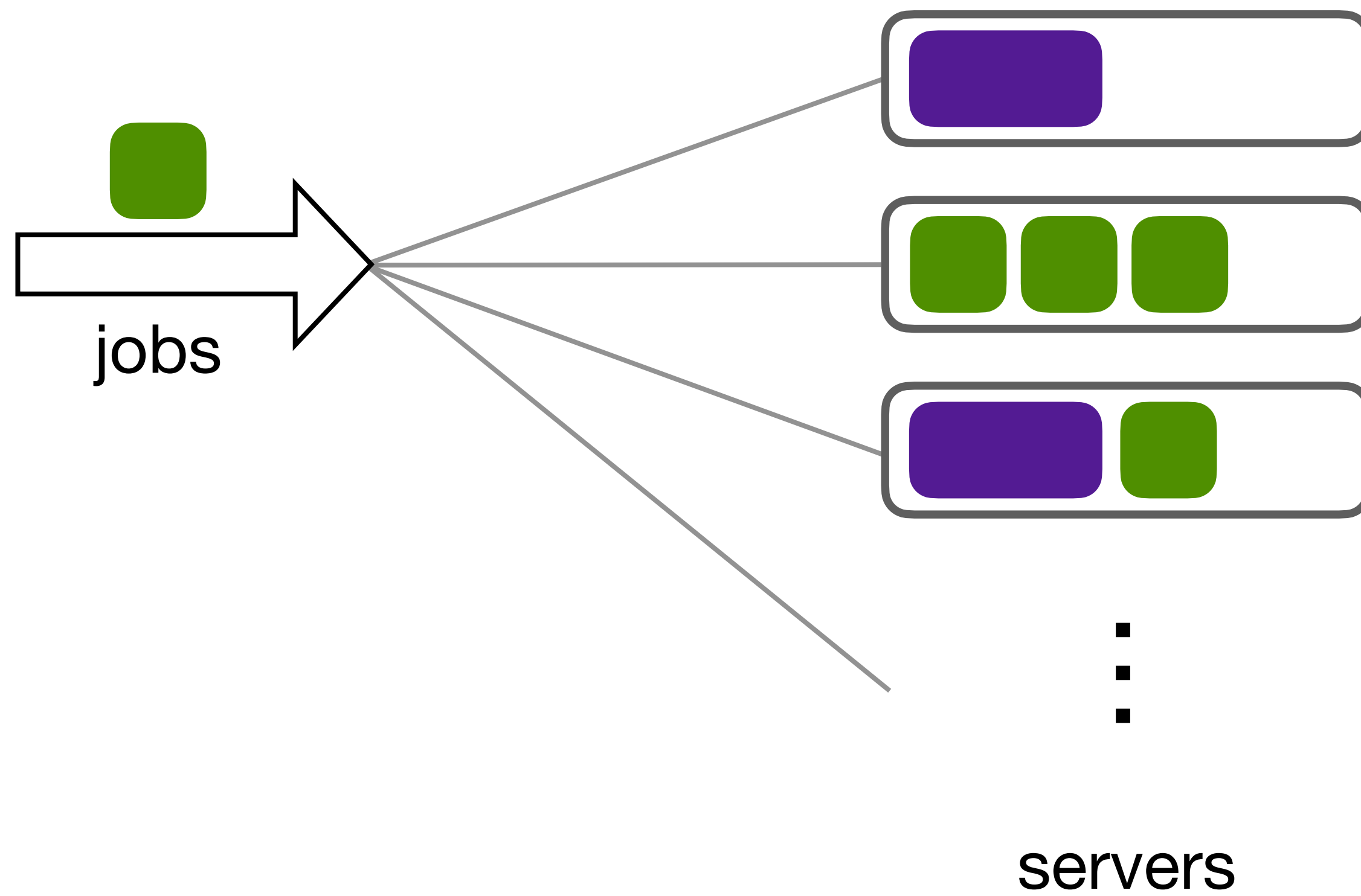
Example MC

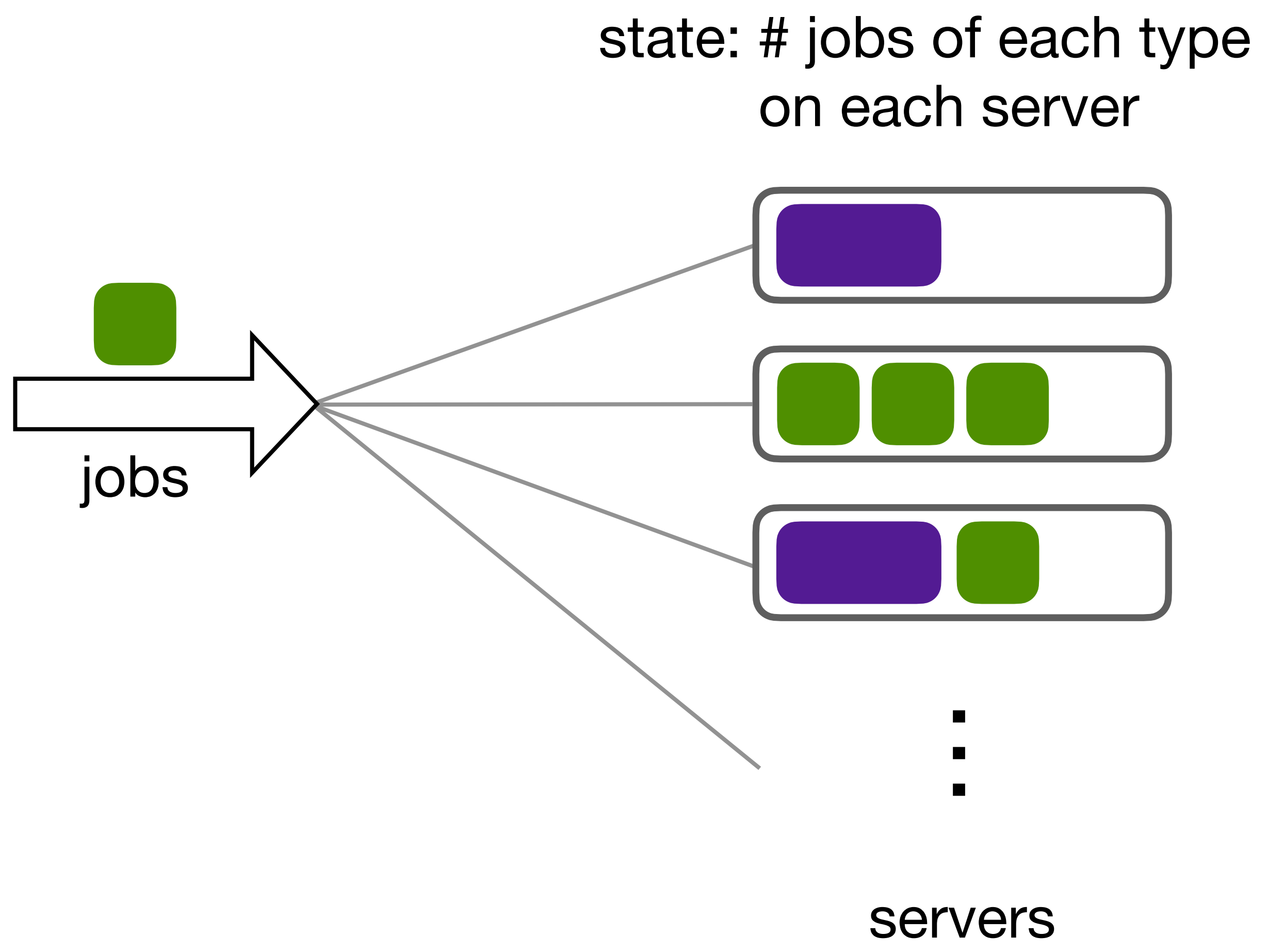
More details on the job model

- 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



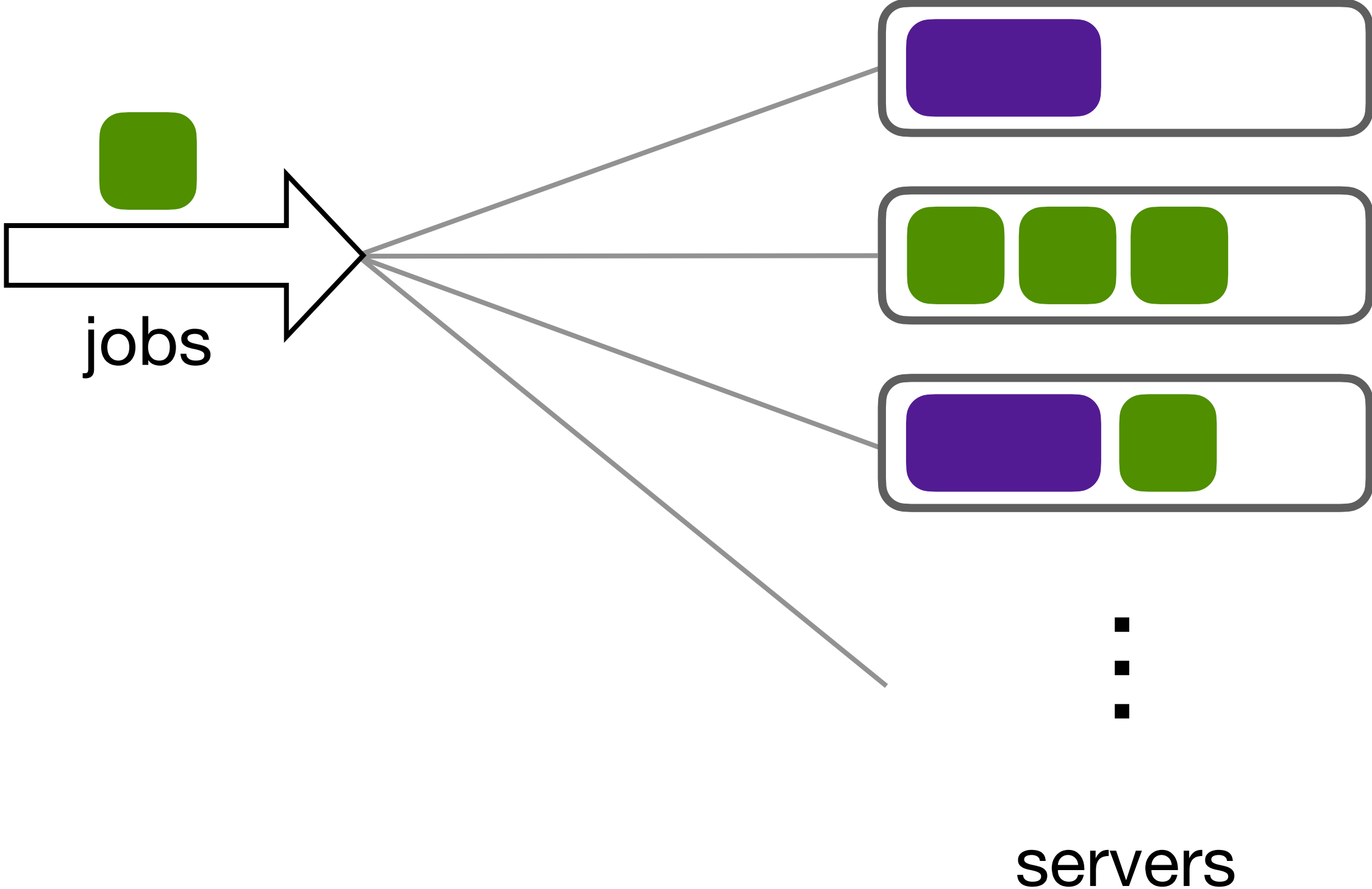
Example MC





state space is large!

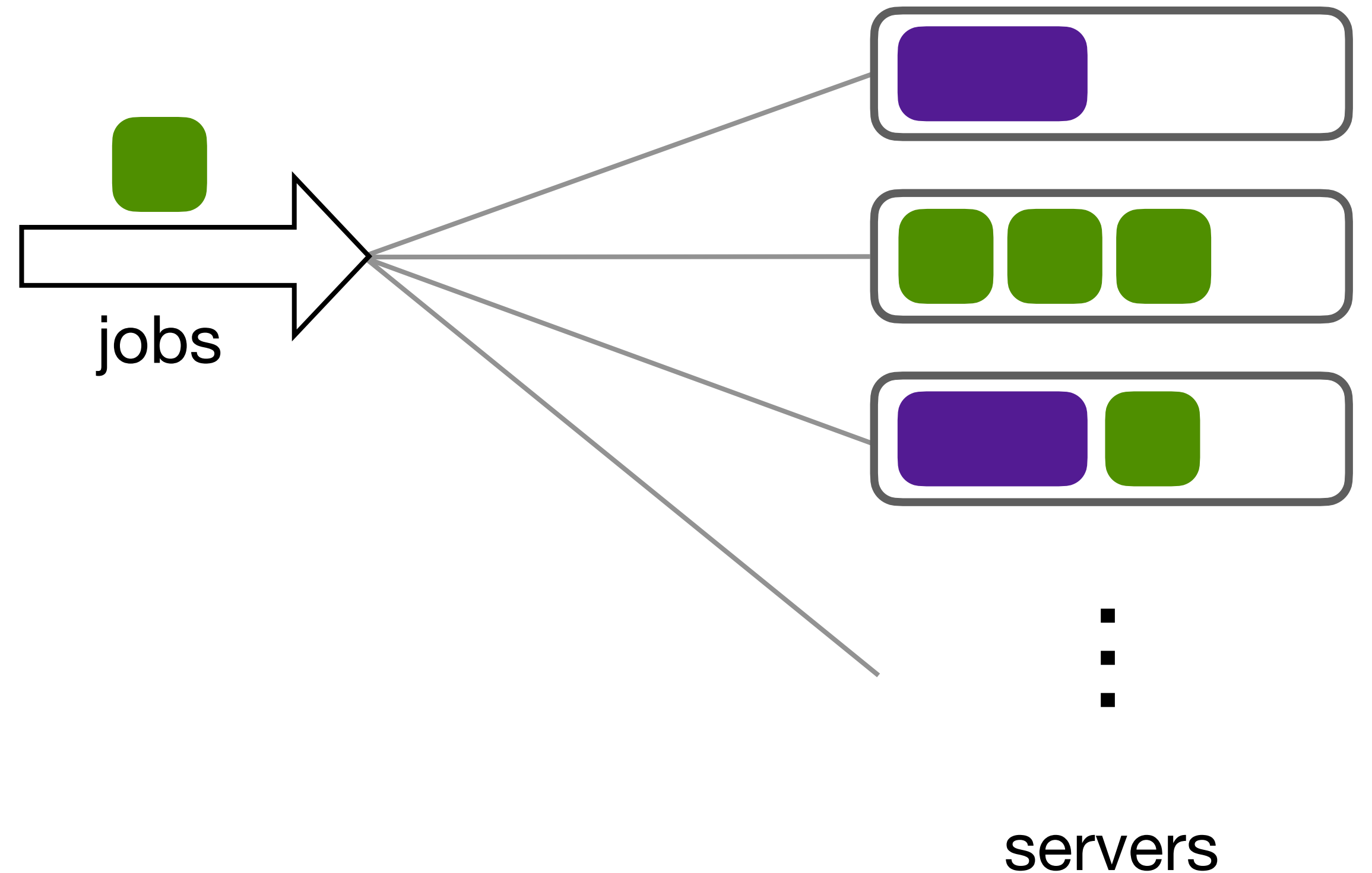
state: # jobs of each type
on each server



Reducing dimensionality

state space is large!

state: # jobs of each type
on each server

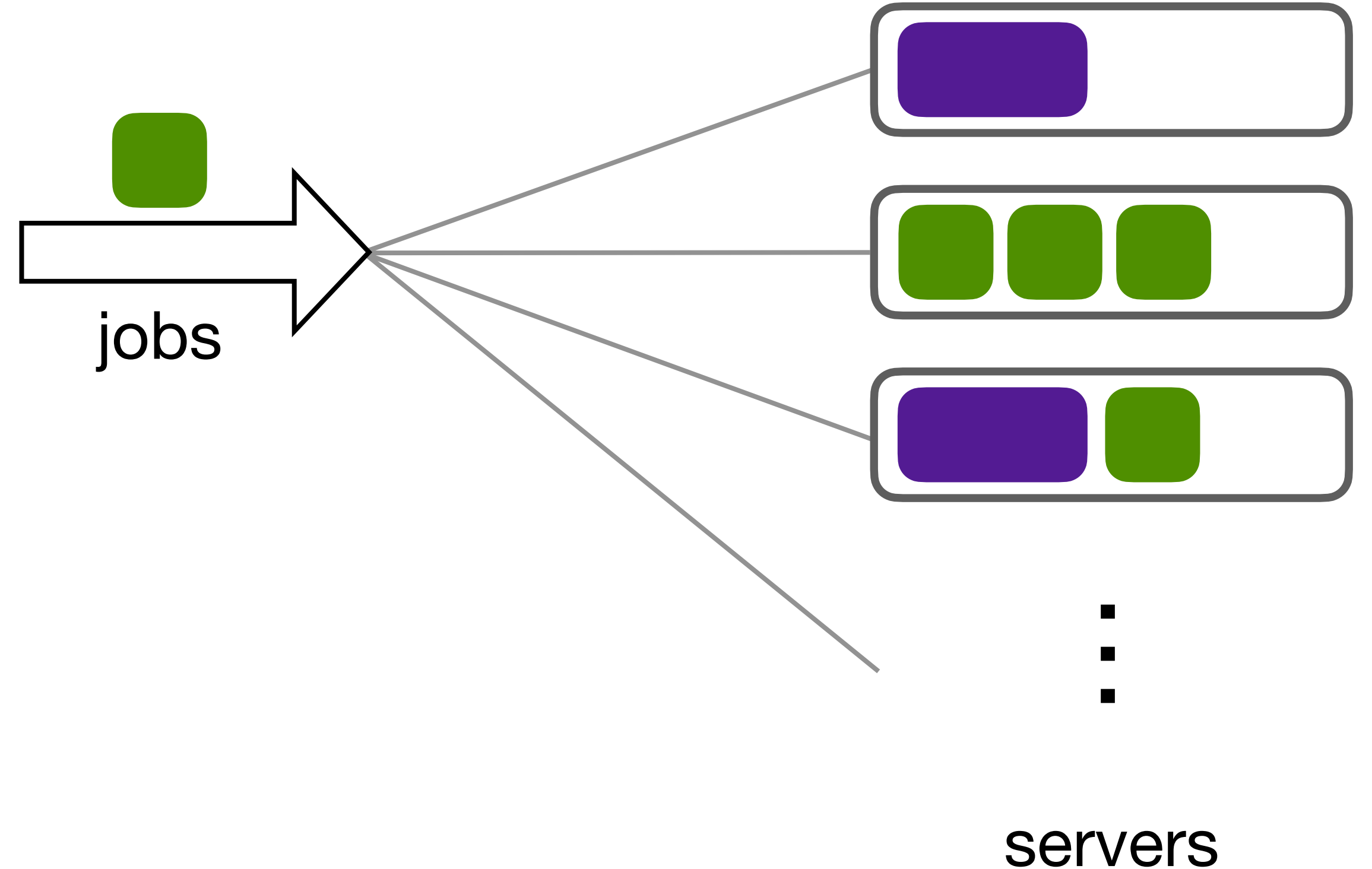


Reducing dimensionality

state space is large!

state: # jobs of each type
on each server

- Server-by-server evaluation:



Reducing dimensionality

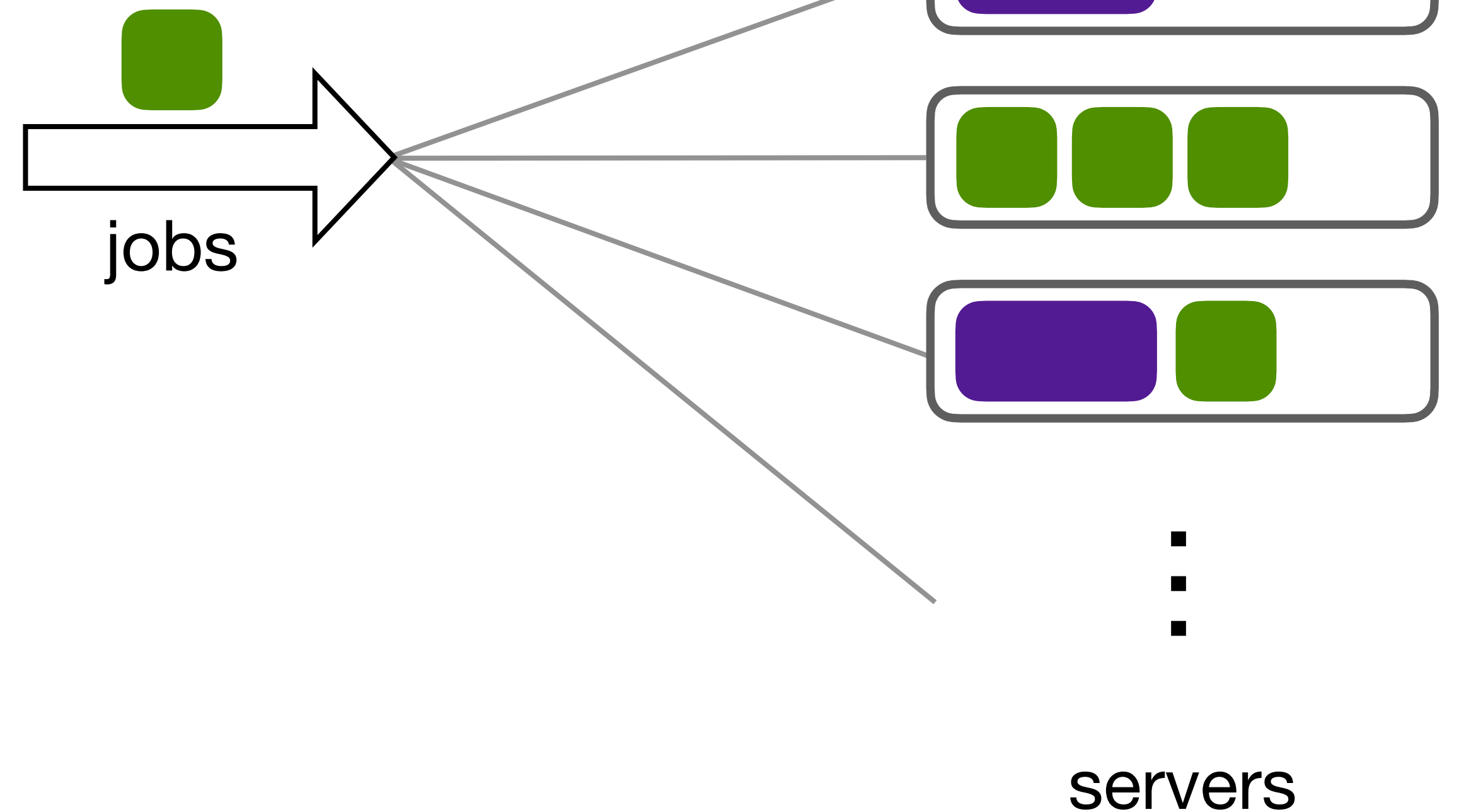
state space is large!

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

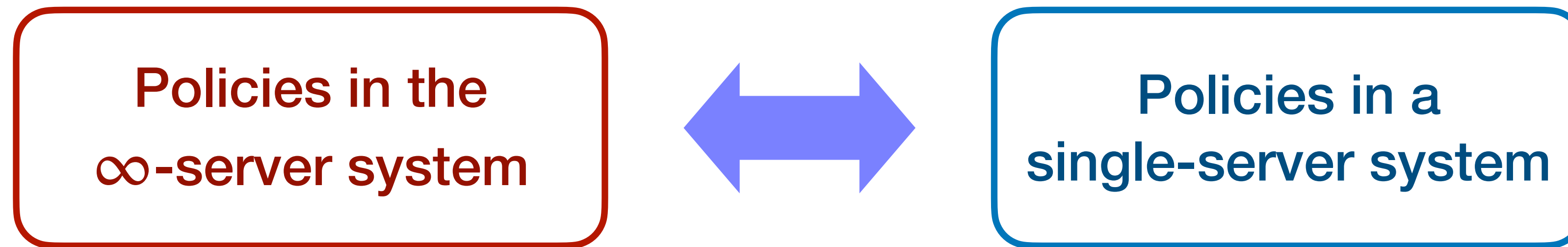


Reducing dimensionality

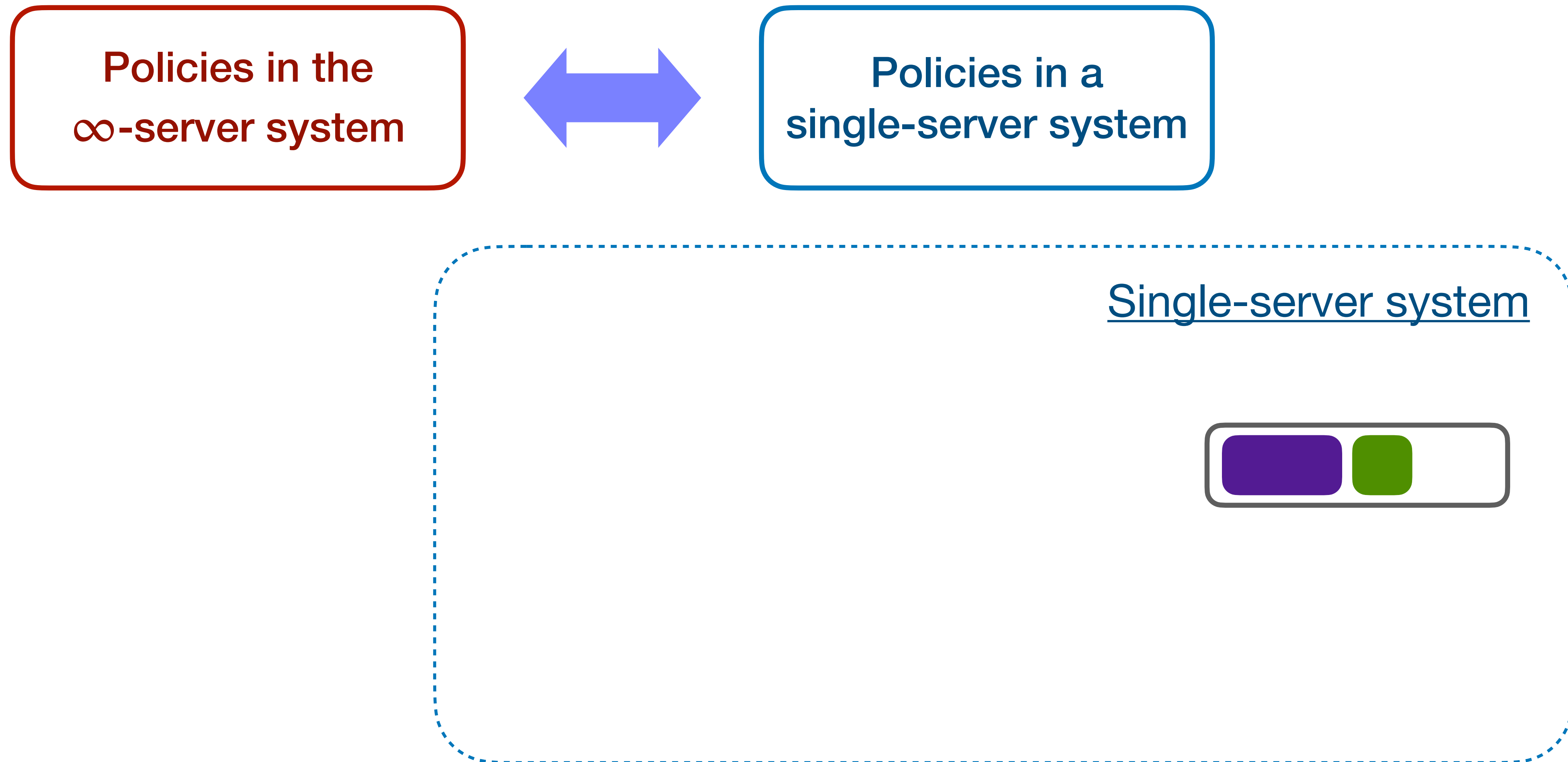
- Server-by-server evaluation:
 - How to evaluate each server?
 - How to relate to $E[\# \text{ active servers}]$?



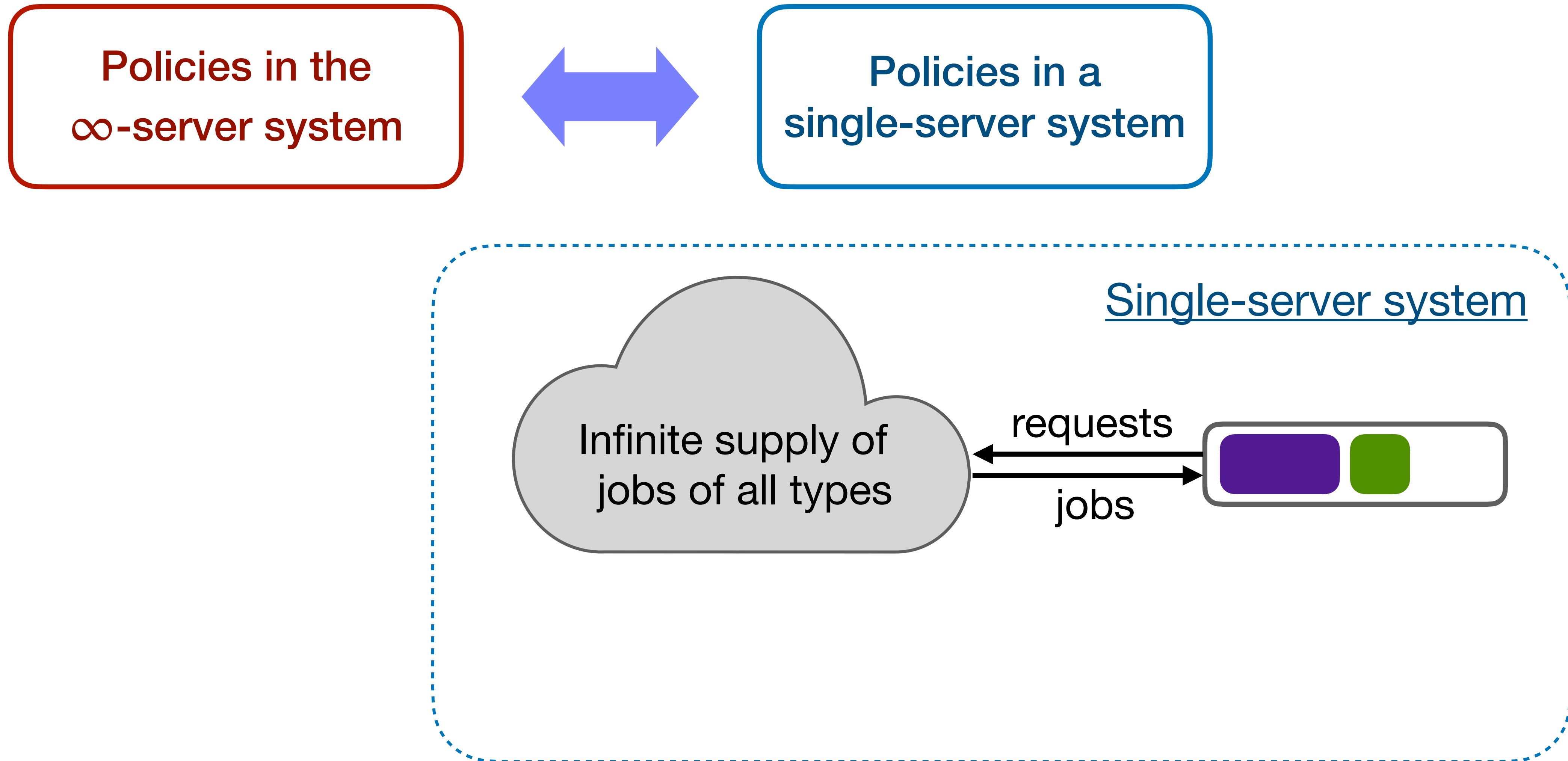
A policy-conversion framework



A policy-conversion framework

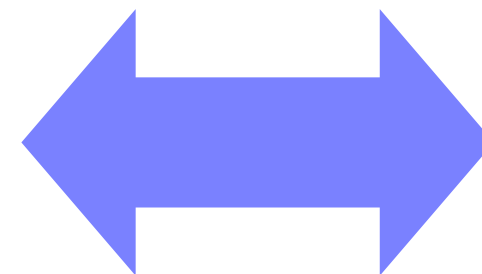


A policy-conversion framework

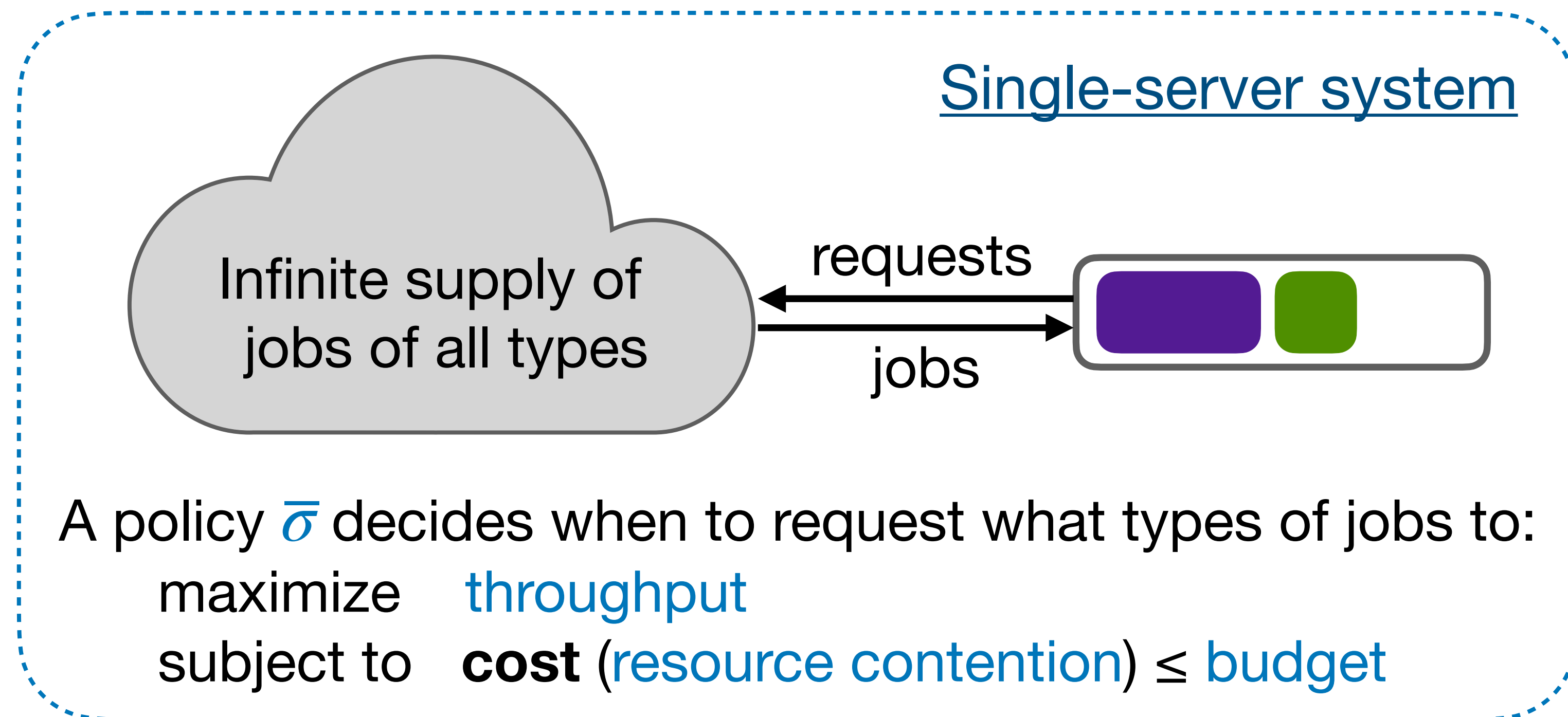


A policy-conversion framework

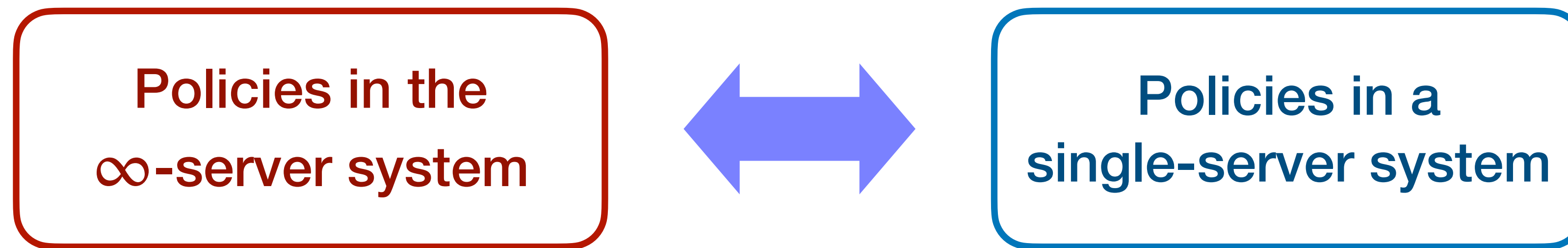
Policies in the ∞ -server system



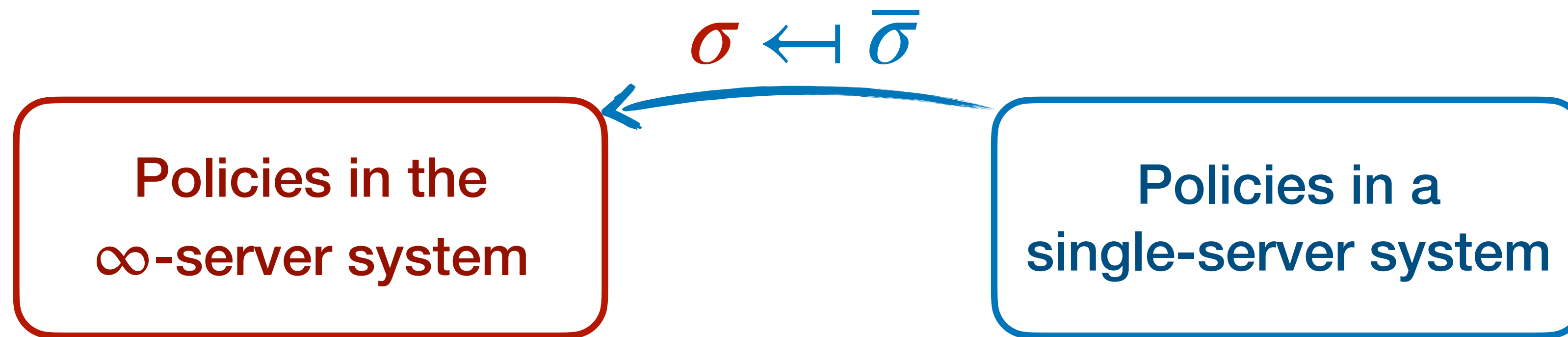
Policies in a single-server system



A policy-conversion framework

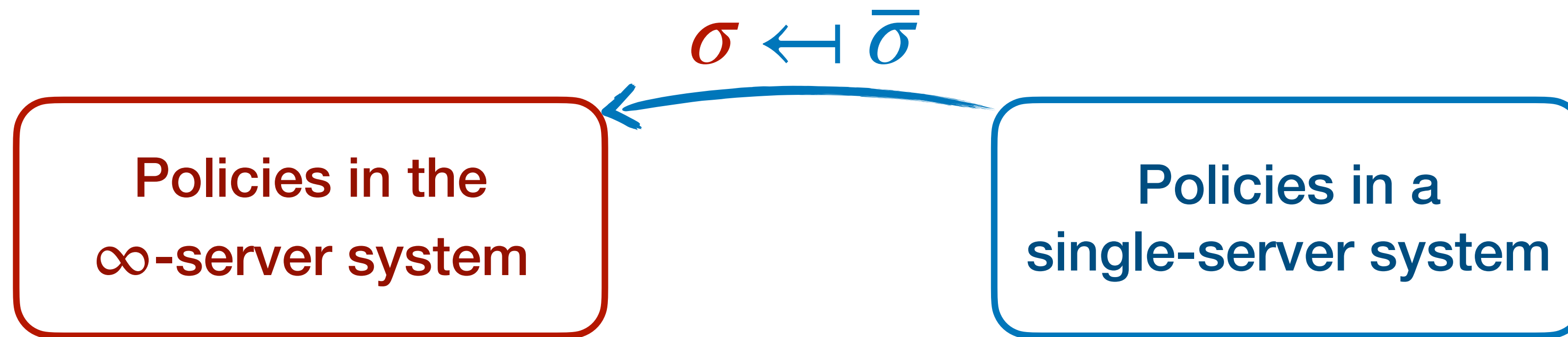


A policy-conversion framework



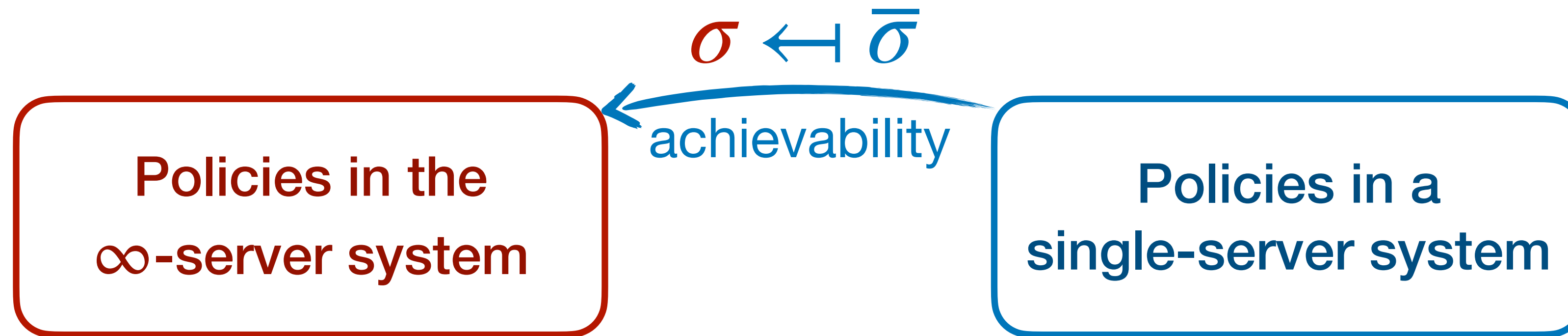
A policy-conversion framework

- Use $\bar{\sigma}$ to tell how to evaluate each server
- Performance of σ is related to properties of $\bar{\sigma}$



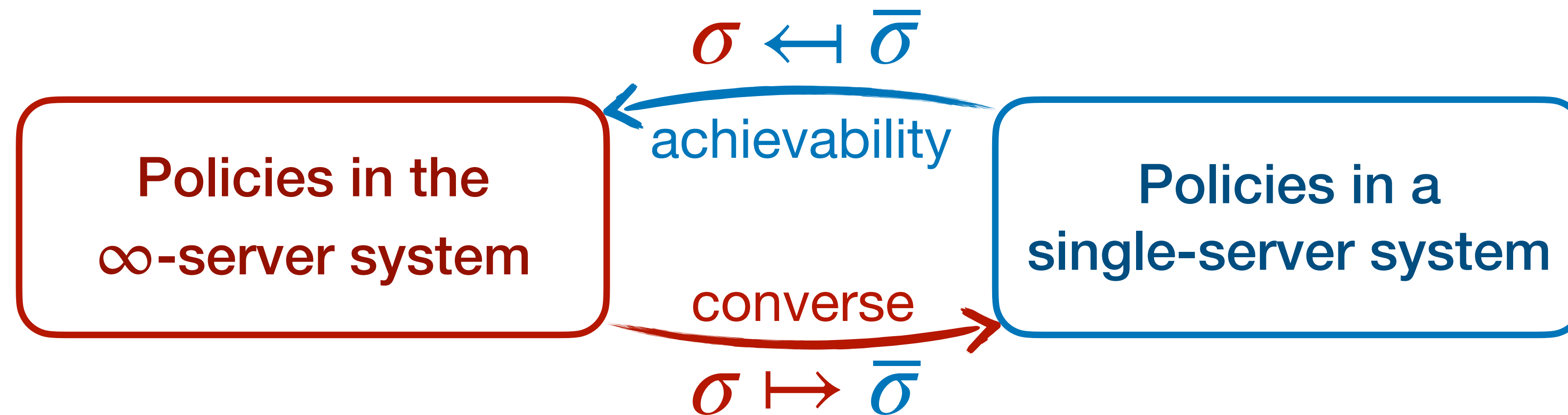
A policy-conversion framework

- Use $\bar{\sigma}$ to tell how to evaluate each server
- Performance of σ is related to properties of $\bar{\sigma}$



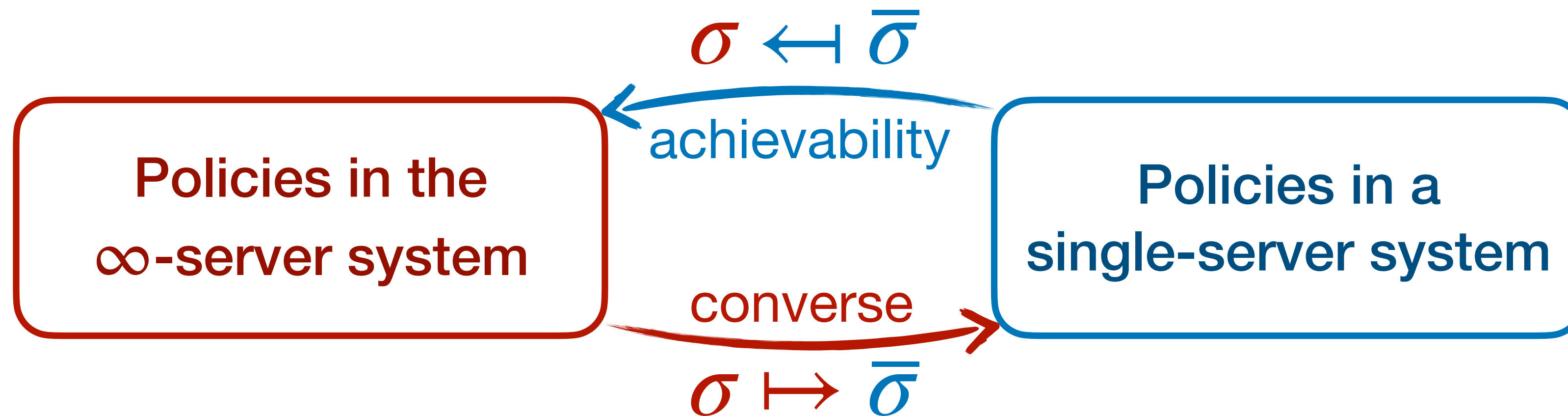
A policy-conversion framework

- Use $\bar{\sigma}$ to tell how to evaluate each server
- Performance of σ is related to properties of $\bar{\sigma}$

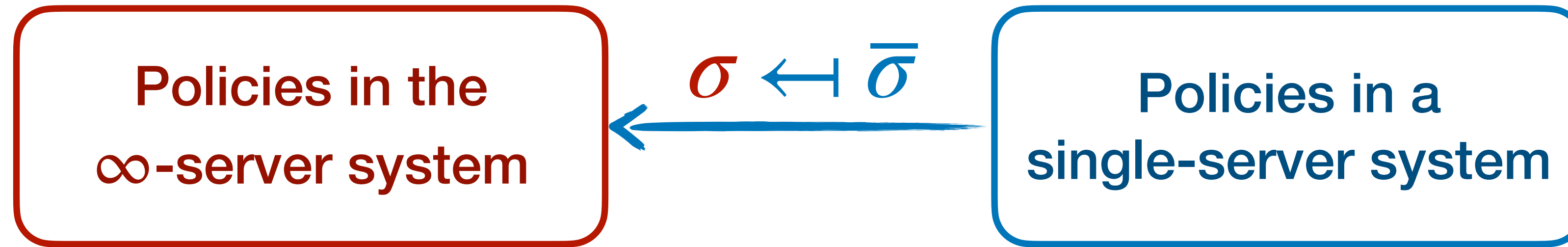


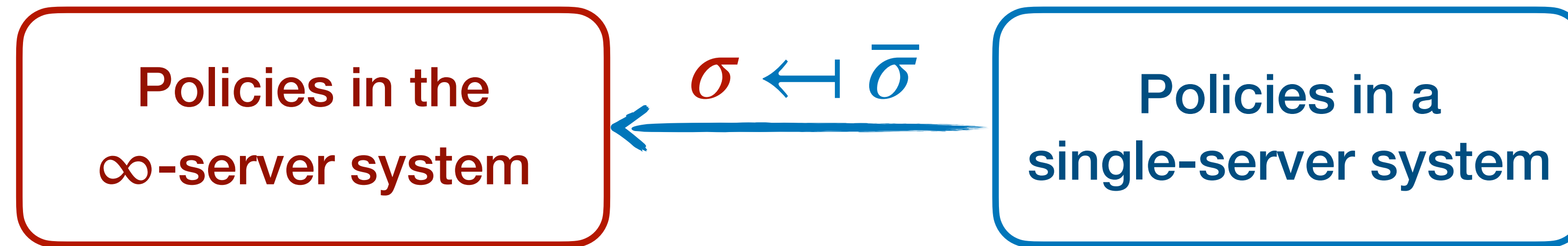
A policy-conversion framework

- Use $\bar{\sigma}$ to tell how to evaluate each server
- Performance of σ is related to properties of $\bar{\sigma}$

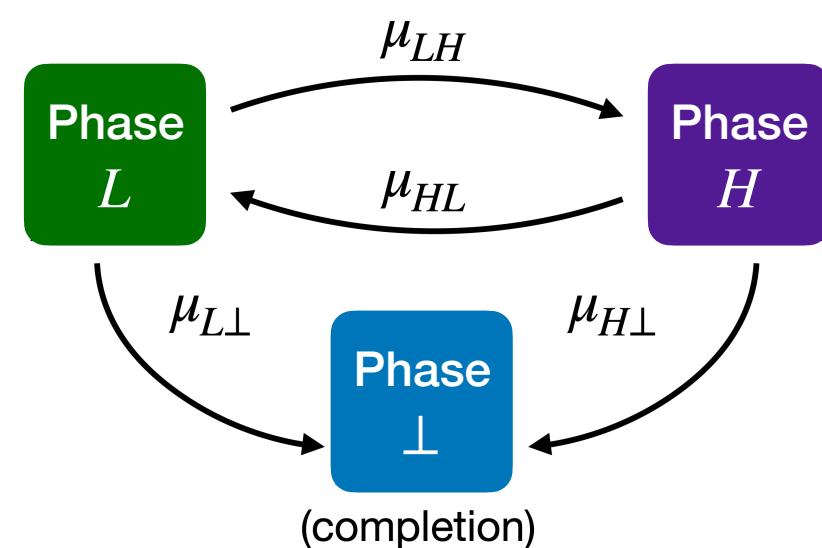


- Allows us to obtain lower bound on $E[\# \text{ active servers}]$





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

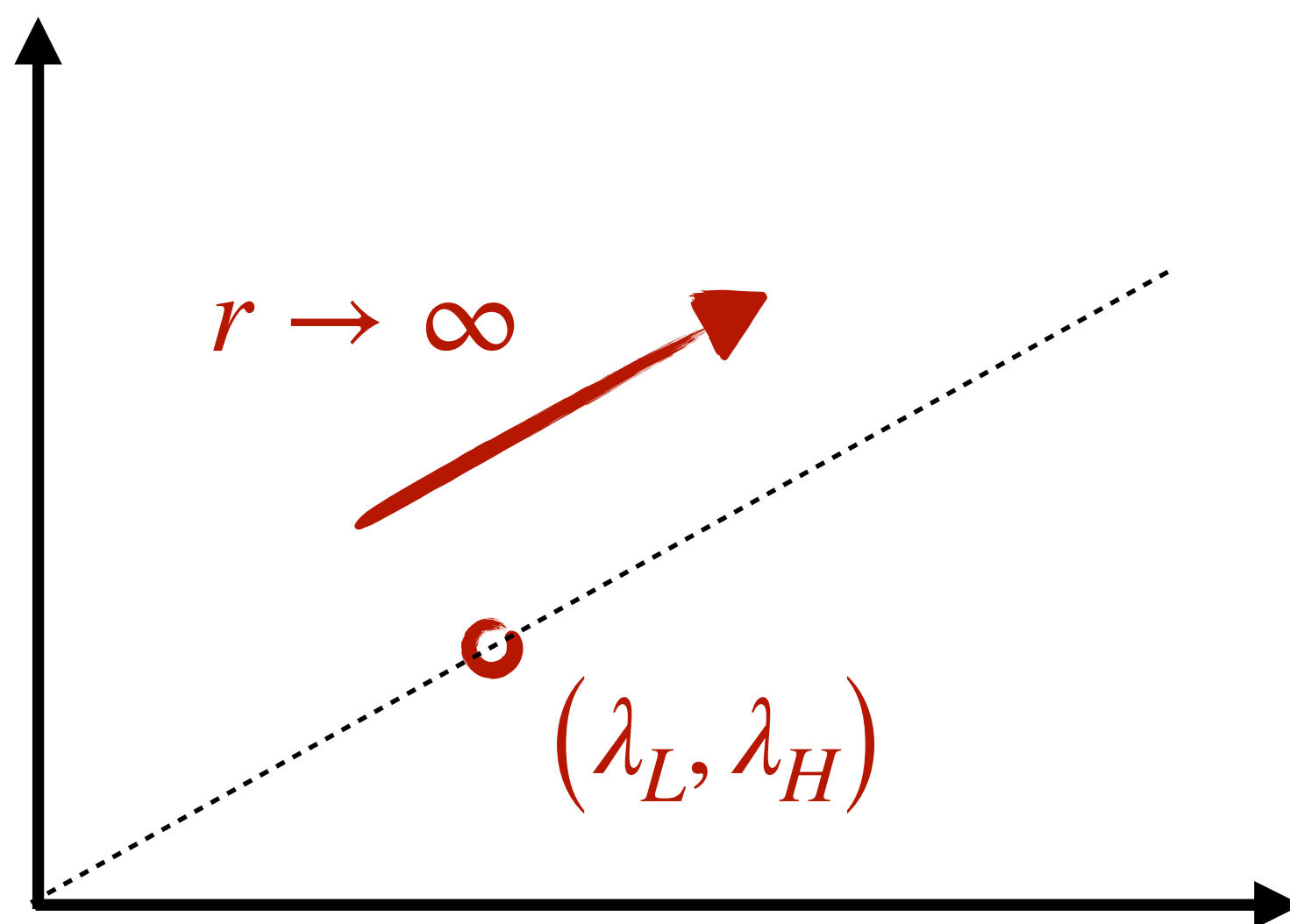
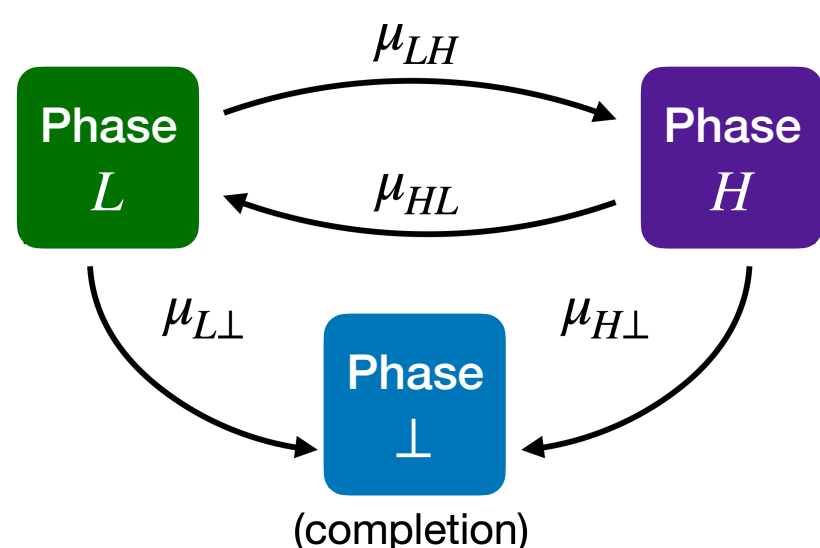


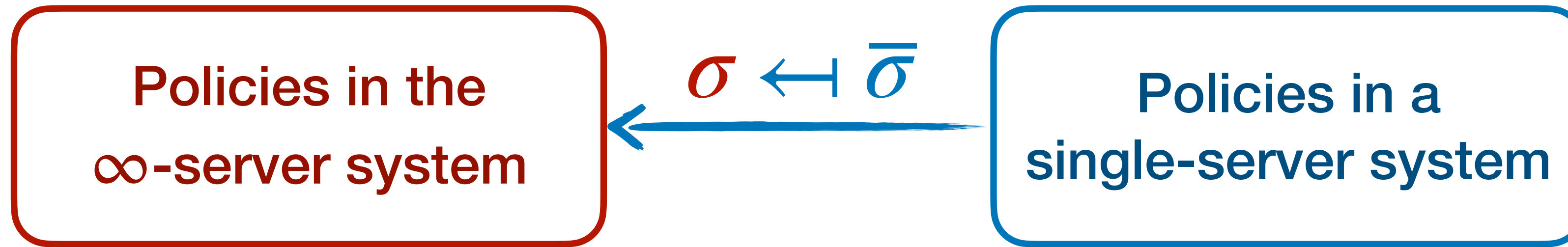
Policies in the ∞ -server system

$\sigma \leftrightarrow \bar{\sigma}$

Policies in a single-server system

- Arrival rates: $r \cdot (\lambda_L, \lambda_H)$
- Asymptotic regime: $r \rightarrow +\infty$





- Arrival rates: $r \cdot (\lambda_L, \lambda_H)$
- Asymptotic regime: $r \rightarrow +\infty$

Policies in the
 ∞ -server system

$\sigma \leftrightarrow \bar{\sigma}$

Policies in a
single-server system

- Arrival rates: $r \cdot (\lambda_L, \lambda_H)$
- Asymptotic regime: $r \rightarrow +\infty$

Policy $\bar{\sigma}$

$$\bar{N} \cdot (\mathbf{throughput}_L, \mathbf{throughput}_R) = r \cdot (\lambda_L, \lambda_H)$$
$$\mathbf{cost} (\text{resource contention}) \leq \text{budget}$$

Policies in the
 ∞ -server system

$\sigma \leftrightarrow \bar{\sigma}$

Policies in a
single-server system

- Arrival rates: $r \cdot (\lambda_L, \lambda_H)$
- Asymptotic regime: $r \rightarrow +\infty$

Policy $\bar{\sigma}$

$\bar{N} \cdot (\text{throughput}_L, \text{throughput}_R) = r \cdot (\lambda_L, \lambda_H)$
cost (resource contention) \leq budget

\bar{N} 's intuitive meaning:
number of servers needed
to satisfy arrival demand

Policies in the ∞ -server system

$\sigma \leftrightarrow \bar{\sigma}$

Policies in a single-server system

- Arrival rates: $r \cdot (\lambda_L, \lambda_H)$
- Asymptotic regime: $r \rightarrow +\infty$

Policy $\bar{\sigma}$

$\bar{N} \cdot (\text{throughput}_L, \text{throughput}_R) = r \cdot (\lambda_L, \lambda_H)$
cost (resource contention) \leq budget

convert

Policy σ

\bar{N} 's intuitive meaning:
number of servers needed
to satisfy arrival demand

Policies in the ∞ -server system

$\sigma \leftrightarrow \bar{\sigma}$

Policies in a single-server system

- Arrival rates: $r \cdot (\lambda_L, \lambda_H)$
- Asymptotic regime: $r \rightarrow +\infty$

Policy $\bar{\sigma}$

$\bar{N} \cdot (\text{throughput}_L, \text{throughput}_R) = r \cdot (\lambda_L, \lambda_H)$
cost (resource contention) \leq **budget**

convert

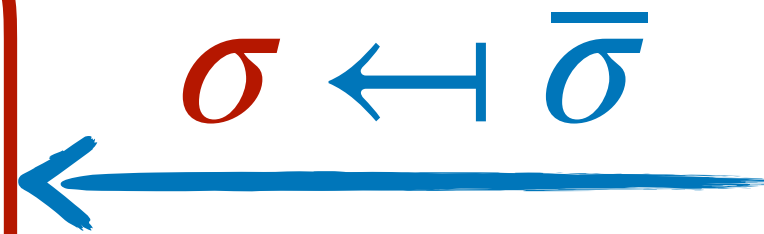
Policy σ

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

\bar{N} 's intuitive meaning:
number of servers needed
to satisfy arrival demand

Policies in the ∞ -server system

Policies in a single-server system



- Arrival rates: $r \cdot (\lambda_L, \lambda_H)$
- Asymptotic regime: $r \rightarrow +\infty$

Policy $\bar{\sigma}$

$\bar{N} \cdot (\text{throughput}_L, \text{throughput}_R) = r \cdot (\lambda_L, \lambda_H)$

cost (resource contention) \leq budget



Policy σ

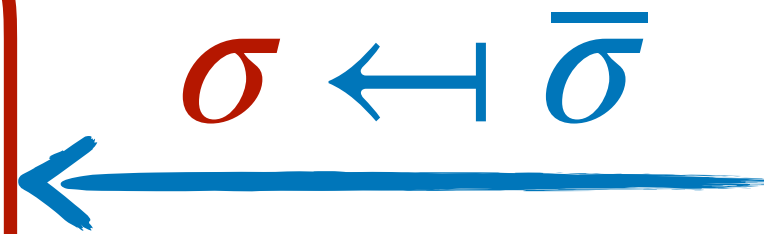
$\mathbf{E} [\# \text{ active servers}] \leq \left(1 + O(r^{-0.5})\right) \cdot \bar{N}$

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

\bar{N} 's intuitive meaning:
number of servers needed
to satisfy arrival demand

Policies in the ∞ -server system

Policies in a single-server system



- Arrival rates: $r \cdot (\lambda_L, \lambda_H)$
- Asymptotic regime: $r \rightarrow +\infty$

Policy $\bar{\sigma}^*$

$\bar{N}^* \cdot (\text{throughput}_L, \text{throughput}_R) = r \cdot (\lambda_L, \lambda_H)$

cost (resource contention) \leq budget



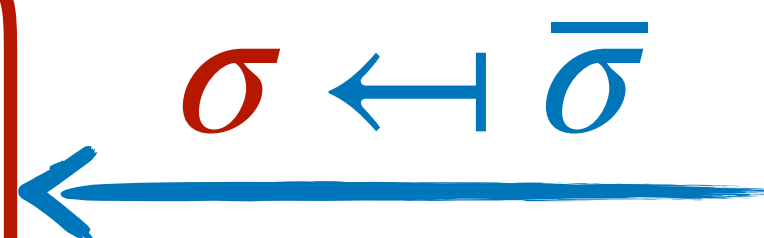
Policy σ^*

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

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

Policies in the ∞ -server system

Policies in a single-server system



- Arrival rates: $r \cdot (\lambda_L, \lambda_H)$
- Asymptotic regime: $r \rightarrow +\infty$

Policy $\bar{\sigma}^*$

$\bar{N}^* \cdot (\text{throughput}_L, \text{throughput}_R) = r \cdot (\lambda_L, \lambda_H)$

cost (resource contention) \leq budget



Policy σ^*

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

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

$\bar{\sigma}^*$ maximizes the throughput and thus achieves the minimum \bar{N}^*

Policies in the ∞ -server system

$\sigma \leftrightarrow \bar{\sigma}$

Policies in a single-server system

Main Results:

Policy $\bar{\sigma}^*$

$$\bar{N}^* \cdot \left(\begin{array}{l} \text{throughput}_L, \text{throughput}_R \\ \text{cost (resource contention)} \end{array} \right) = r \cdot (\lambda_L, \lambda_H) \leq \text{budget}$$

$\bar{\sigma}^*$ maximizes the throughput and thus achieves the minimum \bar{N}^*

Policy σ^*

$$\begin{aligned} \mathbf{E} [\# \text{ active servers}] &\leq \left(1 + O(r^{-0.5}) \right) \cdot \bar{N}^* \\ \text{cost (resource contention)} &\leq \left(1 + O(r^{-0.5}) \right) \cdot \text{budget} \end{aligned}$$

Policies in the ∞ -server system

$\sigma \leftrightarrow \bar{\sigma}$

Policies in a single-server system

Main Results:

- $E[\# \text{ active servers}] \geq \bar{N}^*$ under any policy

Policy $\bar{\sigma}^*$

$\bar{N}^* \cdot (\text{throughput}_L, \text{throughput}_R) = r \cdot (\lambda_L, \lambda_H)$
cost (resource contention) \leq budget

$\bar{\sigma}^*$ maximizes the throughput and thus achieves the minimum \bar{N}^*

Policy σ^*

$E[\# \text{ active servers}] \leq \left(1 + O(r^{-0.5})\right) \cdot \bar{N}^*$
cost (resource contention) $\leq \left(1 + O(r^{-0.5})\right) \cdot \text{budget}$

Policies in the ∞ -server system

$\sigma \leftrightarrow \bar{\sigma}$

Policies in a single-server system

Main Results:

- $E[\# \text{ active servers}] \geq \bar{N}^*$ under any policy
- The policy σ^* converted from $\bar{\sigma}^*$ is *asymptotically optimal*

Policy $\bar{\sigma}^*$

$\bar{N}^* \cdot (\text{throughput}_L, \text{throughput}_R) = r \cdot (\lambda_L, \lambda_H)$

cost (resource contention) \leq budget

$\bar{\sigma}^*$ maximizes the throughput and thus achieves the minimum \bar{N}^*

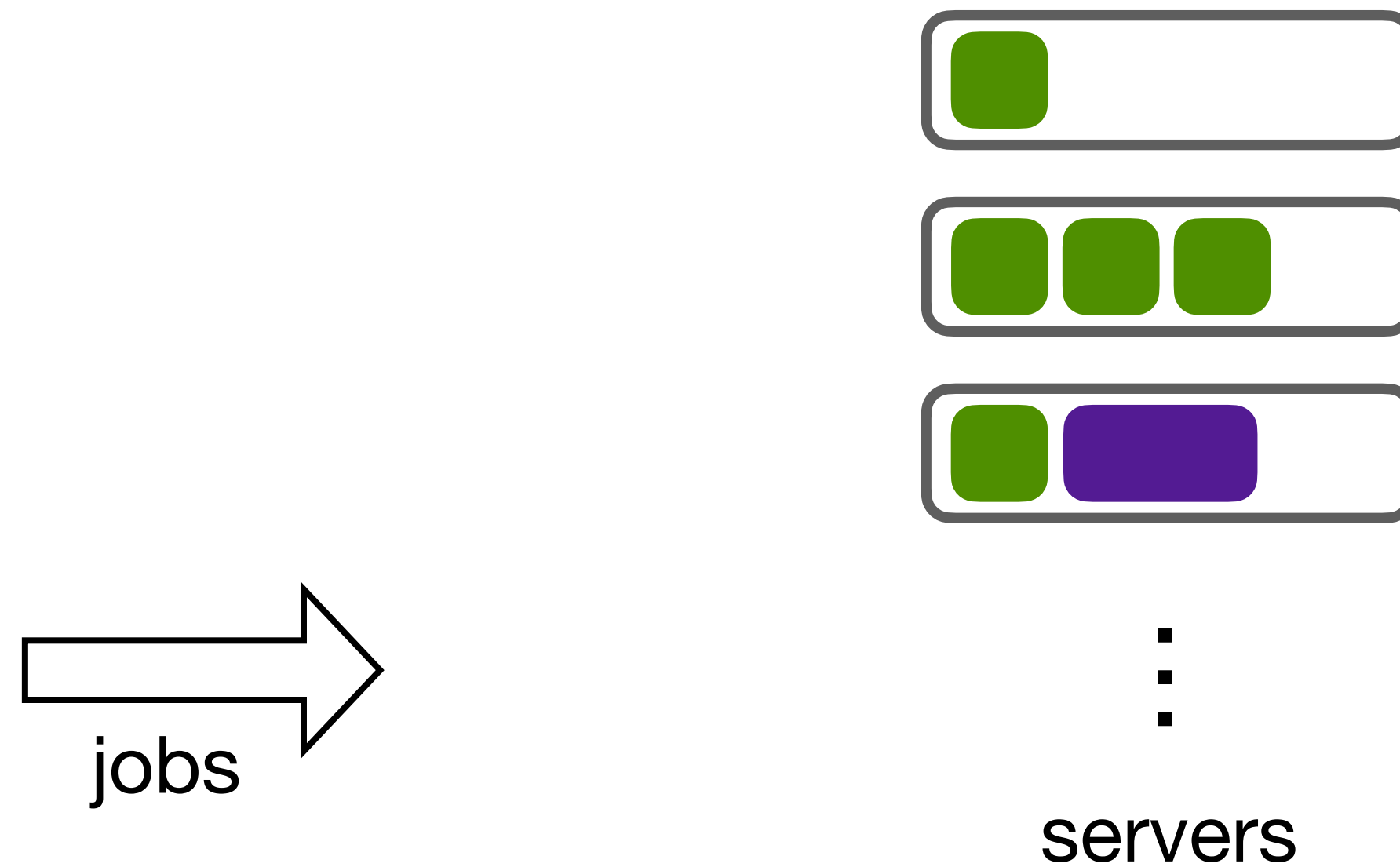
Policy σ^*

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

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

Policy conversion: single-server to ∞ -server

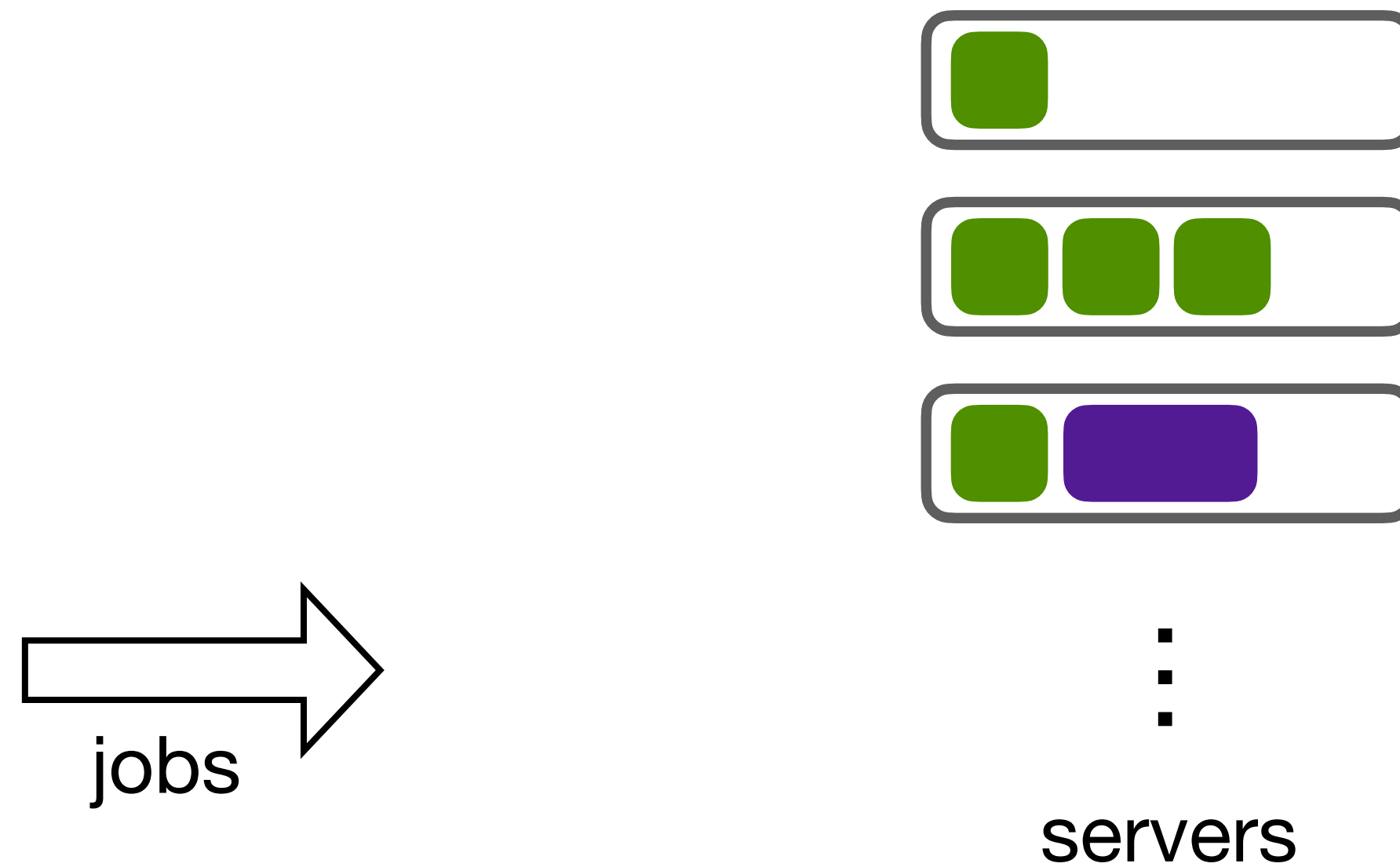
Meta-algorithm: JOIN-REQUESTING-SERVER ($\bar{\sigma}$)



Policy conversion: single-server to ∞ -server

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

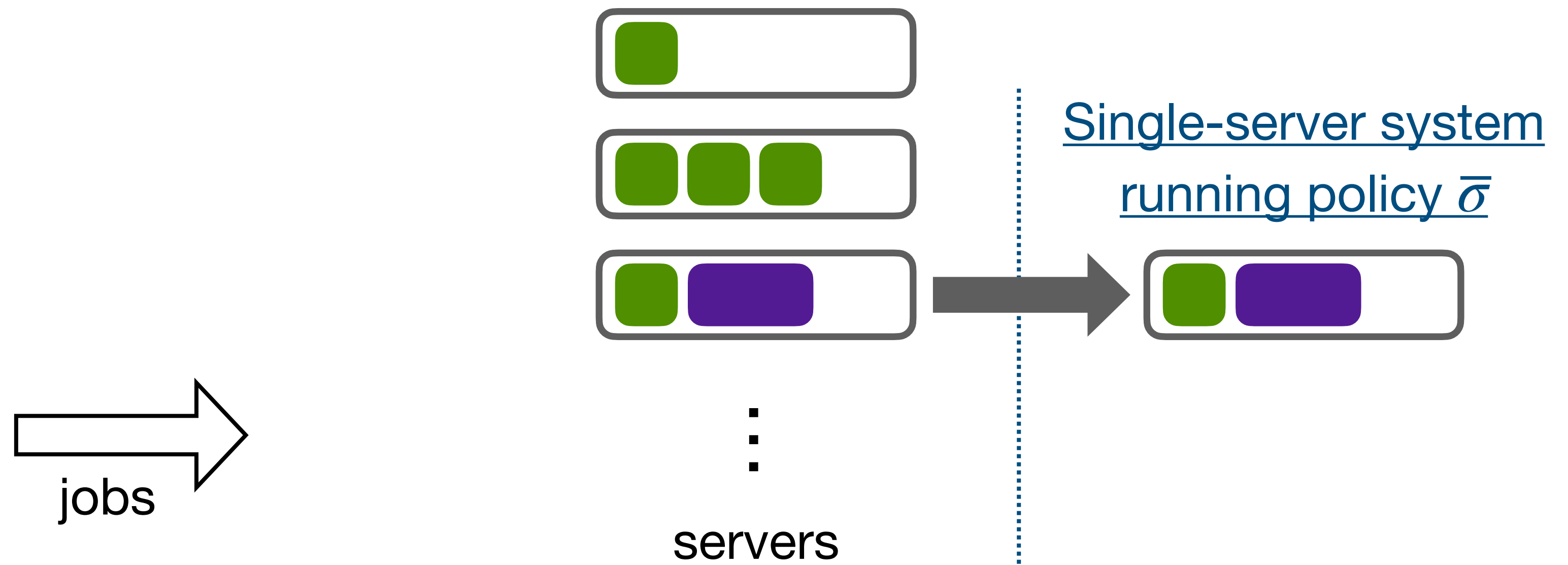
- For each server, run a single-server policy $\bar{\sigma}$



Policy conversion: single-server to ∞ -server

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

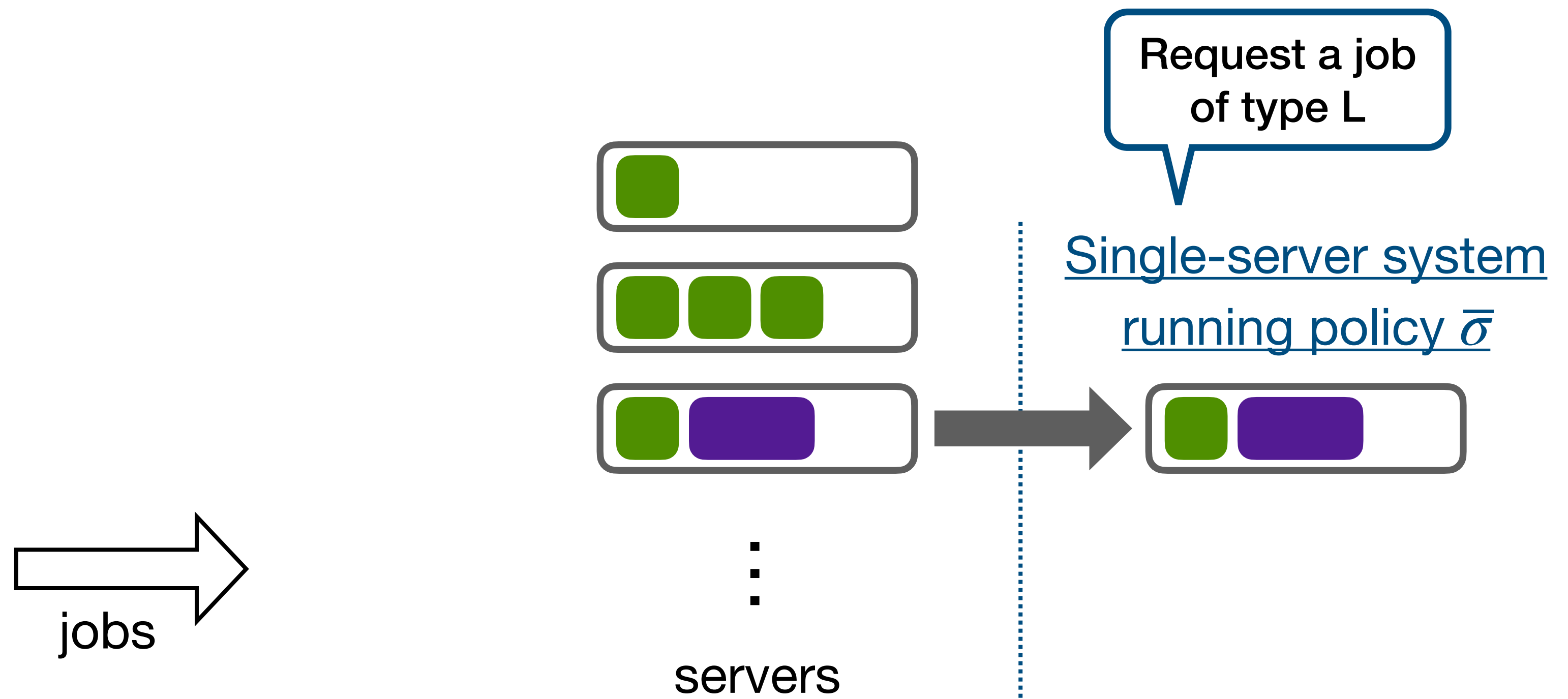
- For each server, run a single-server policy $\bar{\sigma}$



Policy conversion: single-server to ∞ -server

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

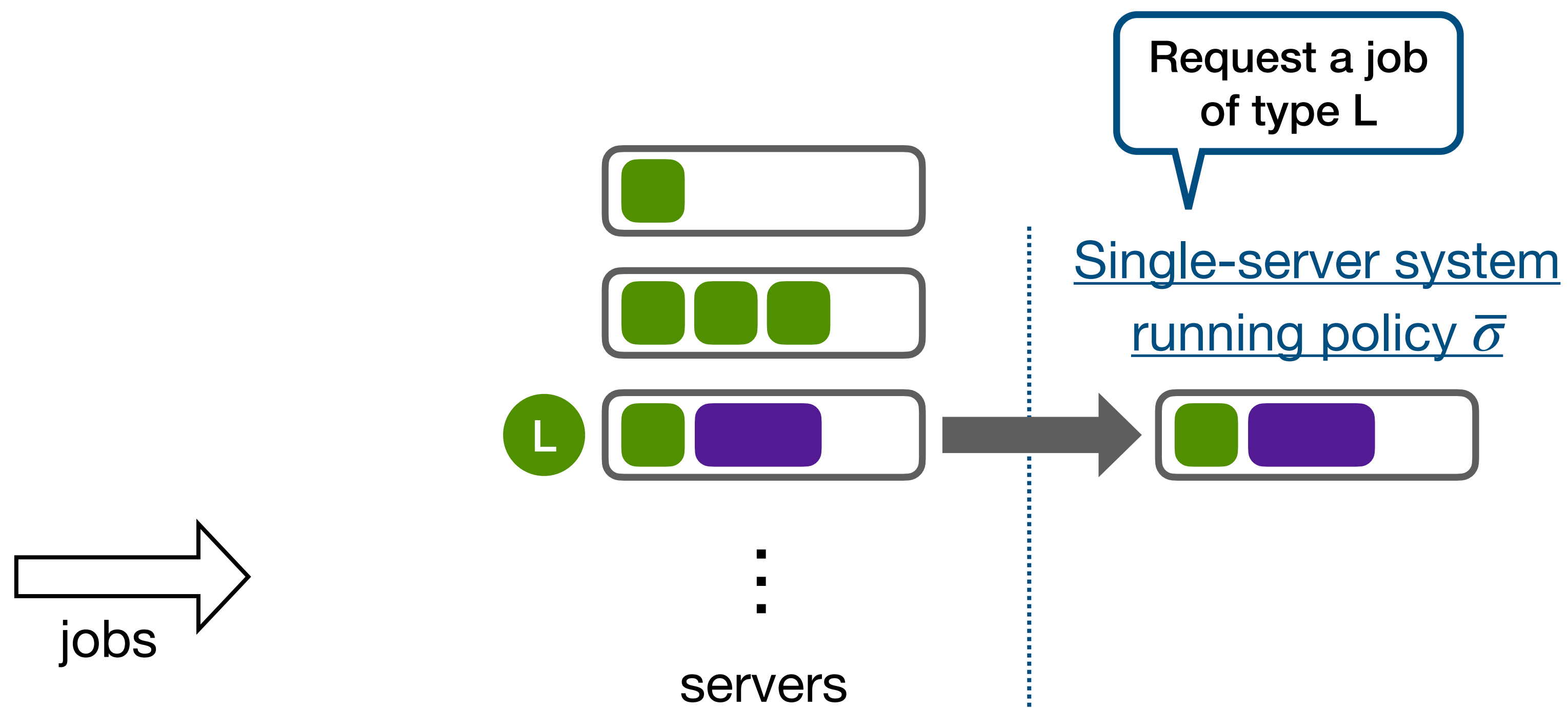
- For each server, run a single-server policy $\bar{\sigma}$
- If $\bar{\sigma}$ requests a job of type i , generate a token of type i



Policy conversion: single-server to ∞ -server

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

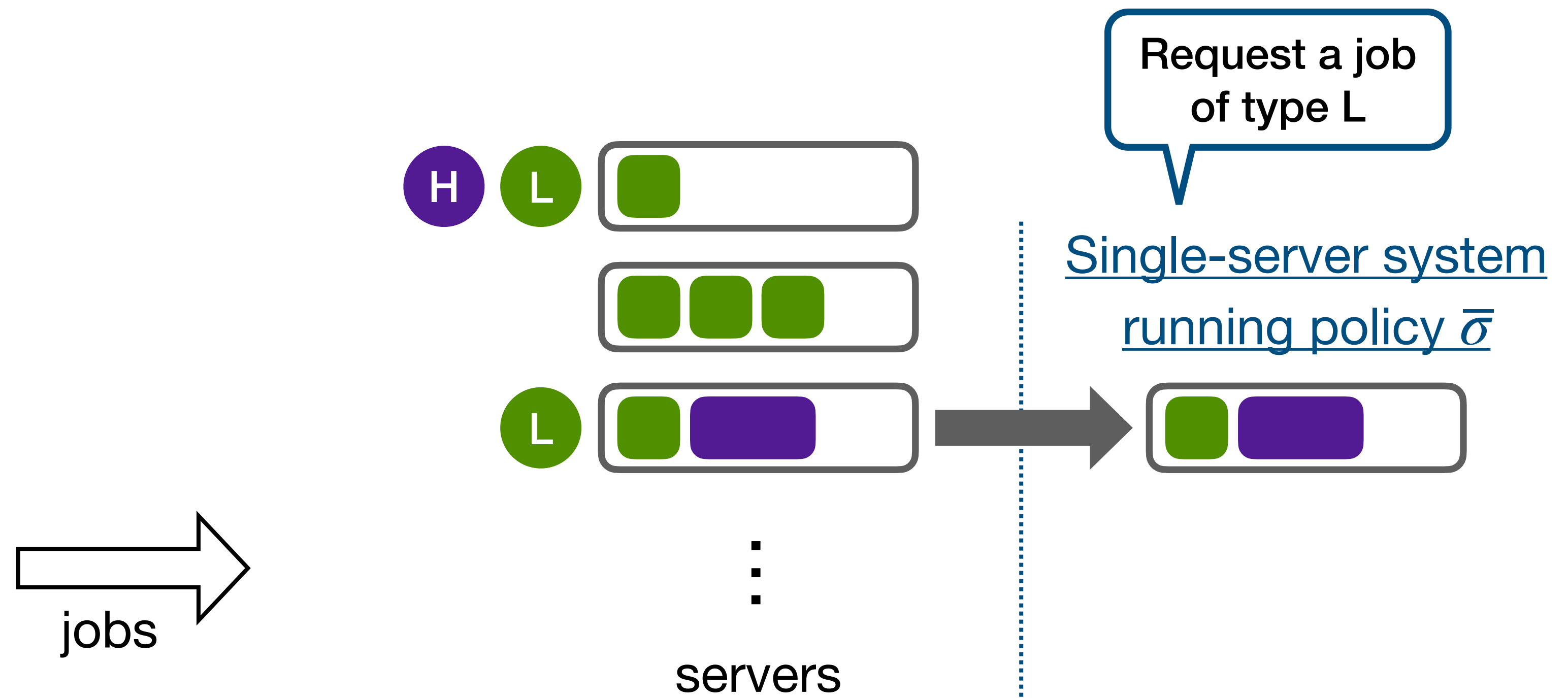
- For each server, run a single-server policy $\bar{\sigma}$
- If $\bar{\sigma}$ requests a job of type i , generate a token of type i



Policy conversion: single-server to ∞ -server

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

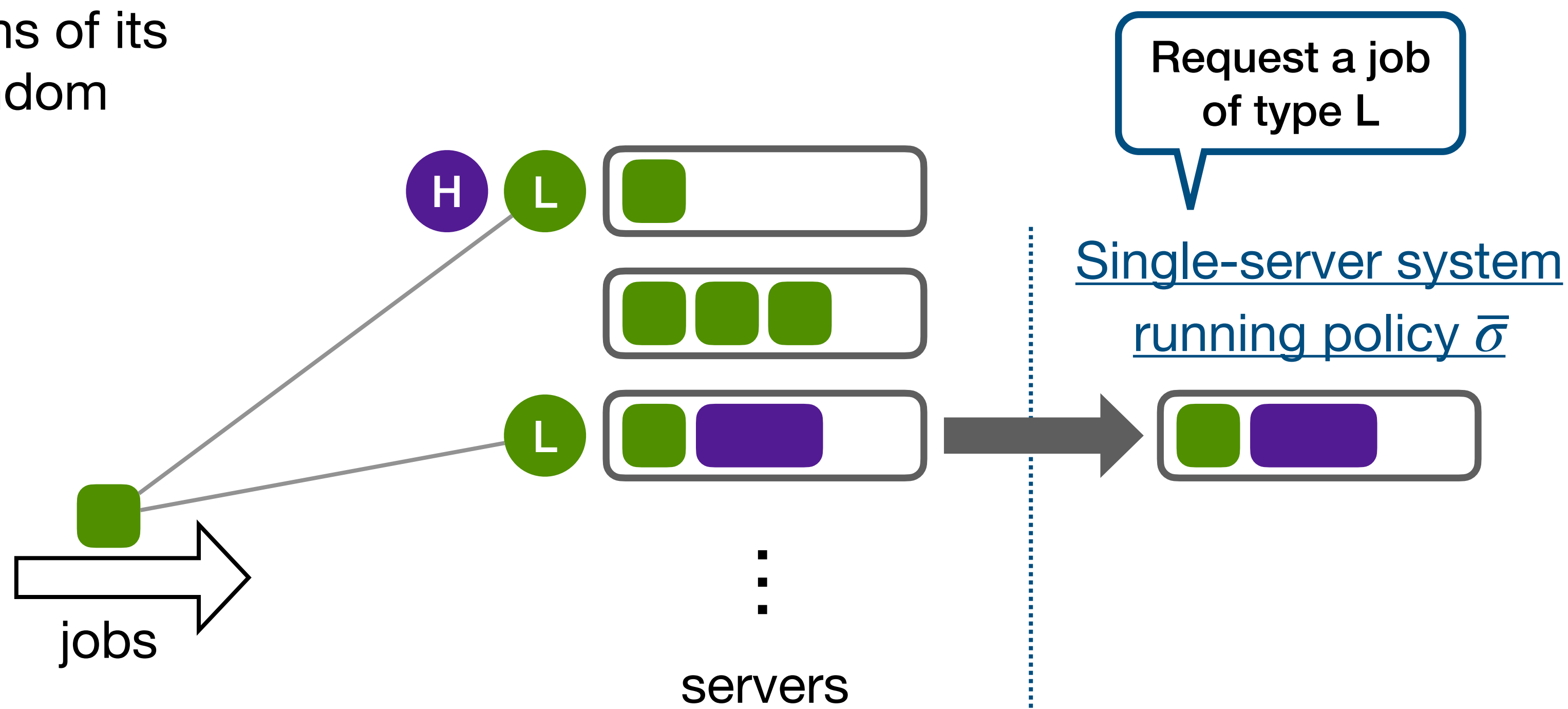
- For each server, run a single-server policy $\bar{\sigma}$
- If $\bar{\sigma}$ requests a job of type i , generate a token of type i



Policy conversion: single-server to ∞ -server

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

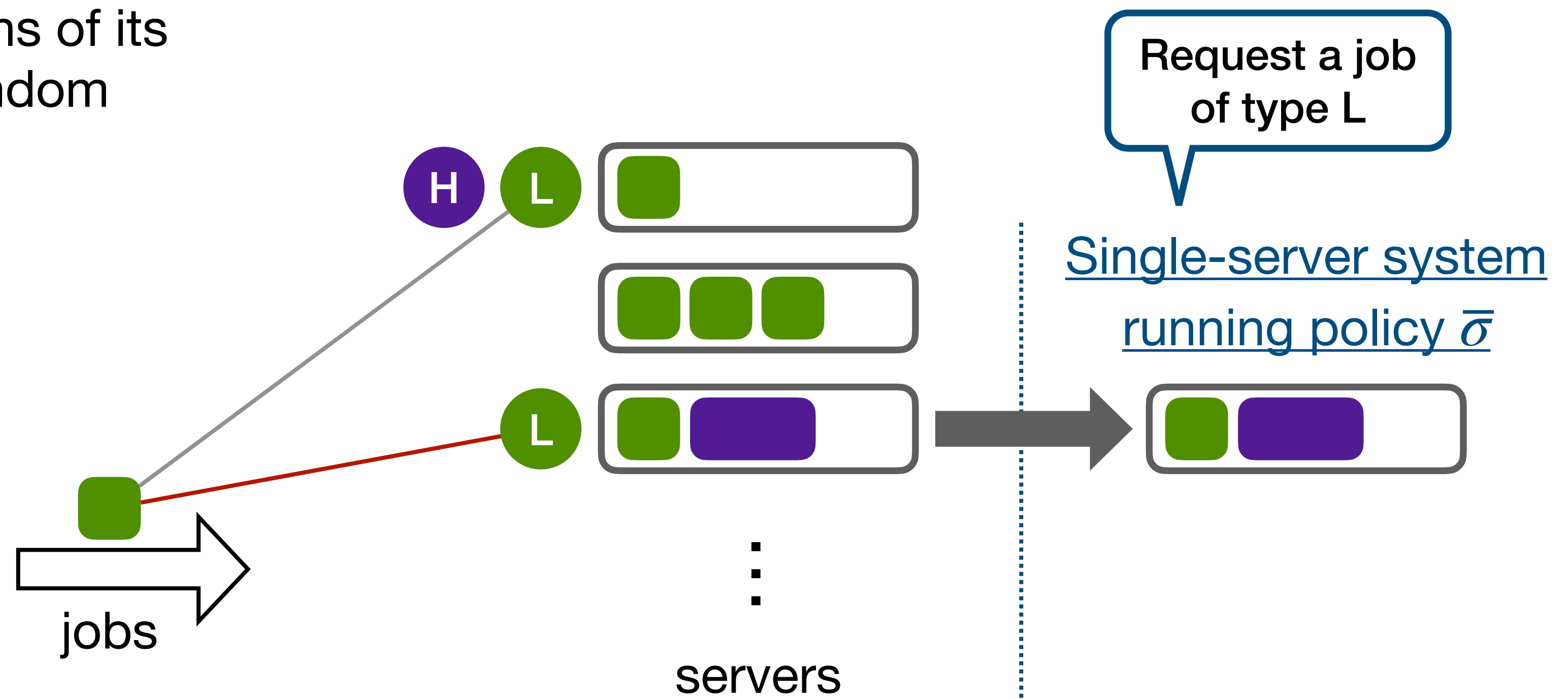
- For each server, run a single-server policy $\bar{\sigma}$
- If $\bar{\sigma}$ requests a job of type i , generate a token of type i
- When a job arrives, it checks tokens of its type and joins one uniformly at random



Policy conversion: single-server to ∞ -server

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

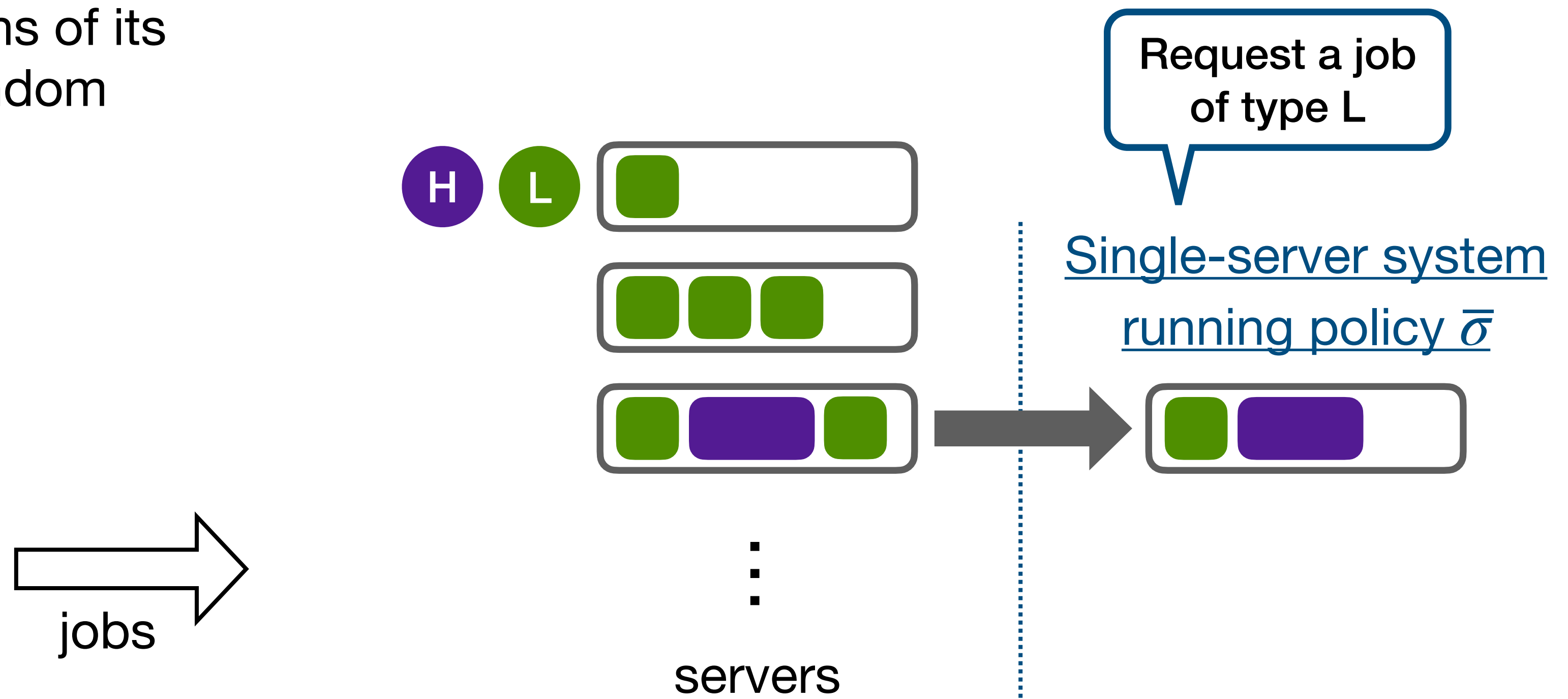
- For each server, run a single-server policy $\bar{\sigma}$
- If $\bar{\sigma}$ requests a job of type i , generate a token of type i
- When a job arrives, it checks tokens of its type and joins one uniformly at random



Policy conversion: single-server to ∞ -server

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

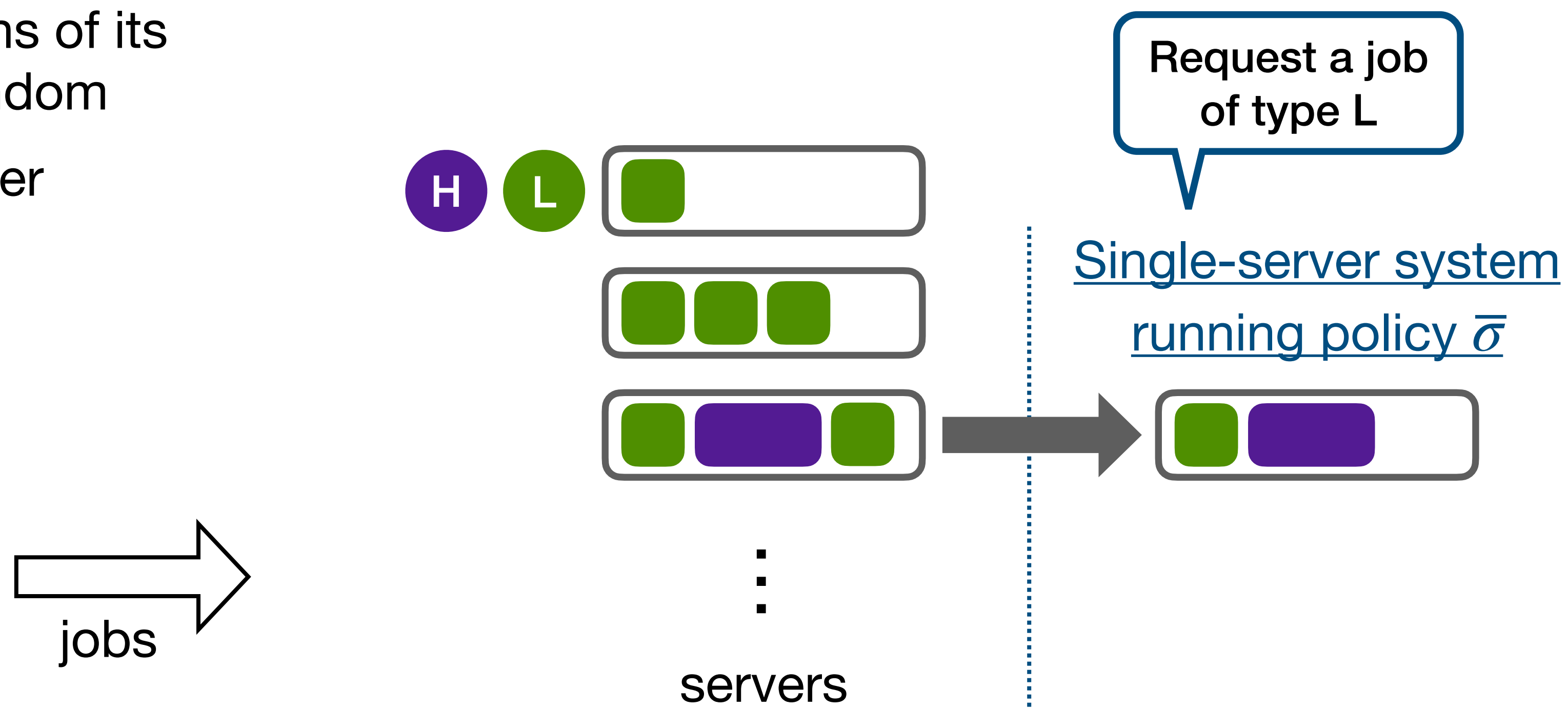
- For each server, run a single-server policy $\bar{\sigma}$
- If $\bar{\sigma}$ requests a job of type i , generate a token of type i
- When a job arrives, it checks tokens of its type and joins one uniformly at random



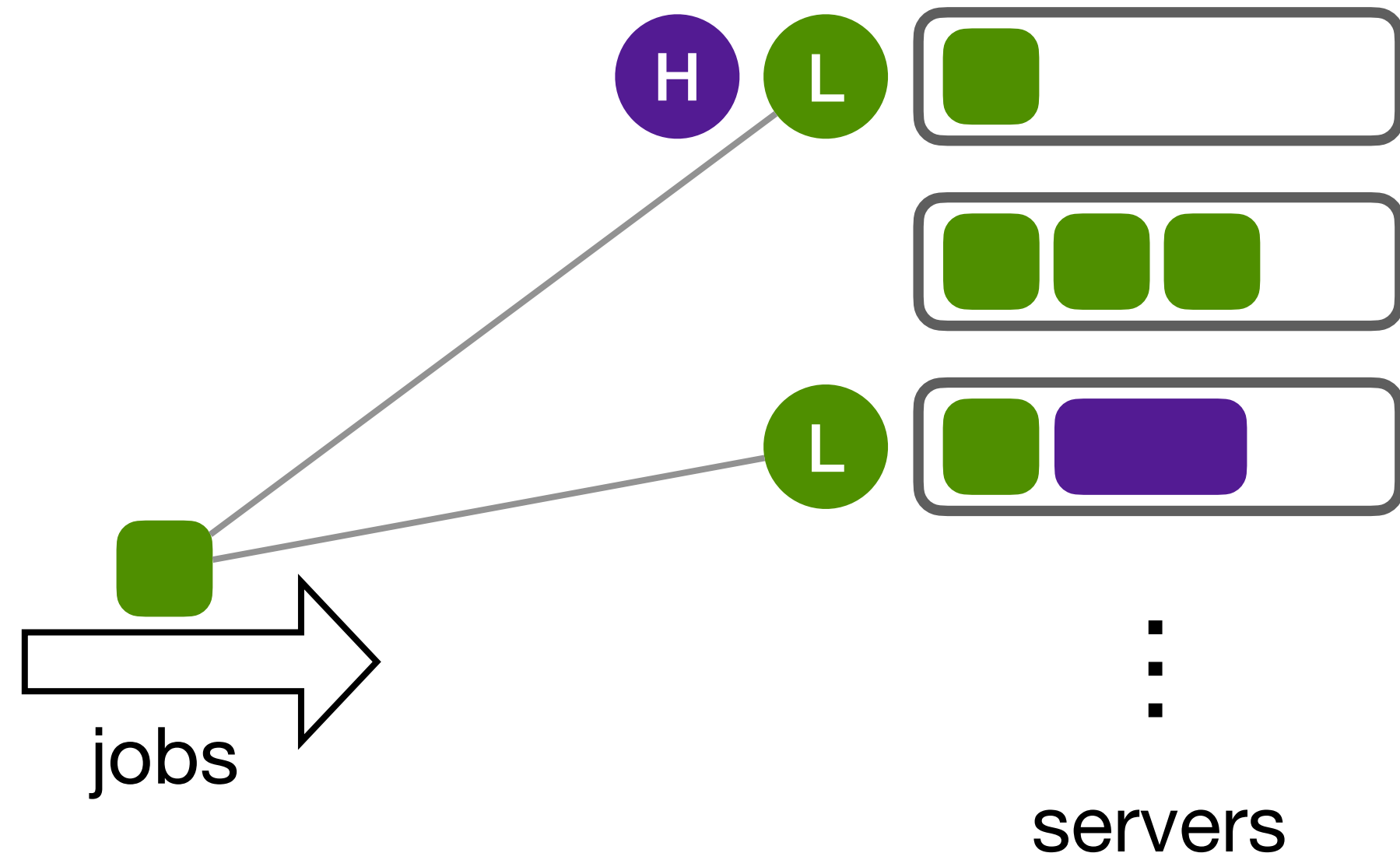
Policy conversion: single-server to ∞ -server

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

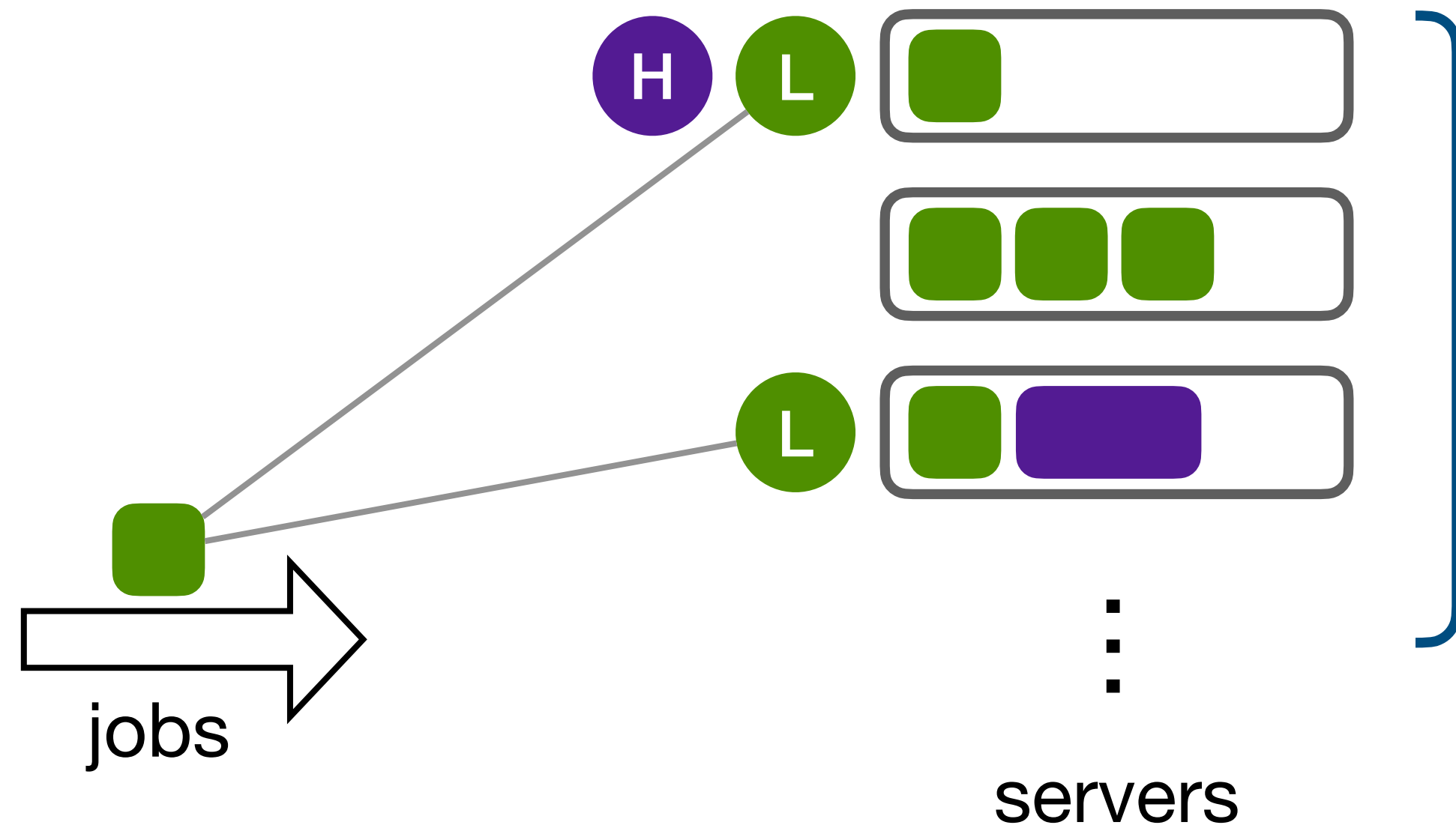
- For each server, run a single-server policy $\bar{\sigma}$
- If $\bar{\sigma}$ requests a job of type i , generate a token of type i
- When a job arrives, it checks tokens of its type and joins one uniformly at random
- If no tokens, go to an inactive server



Policy conversion: more details



Policy conversion: more details

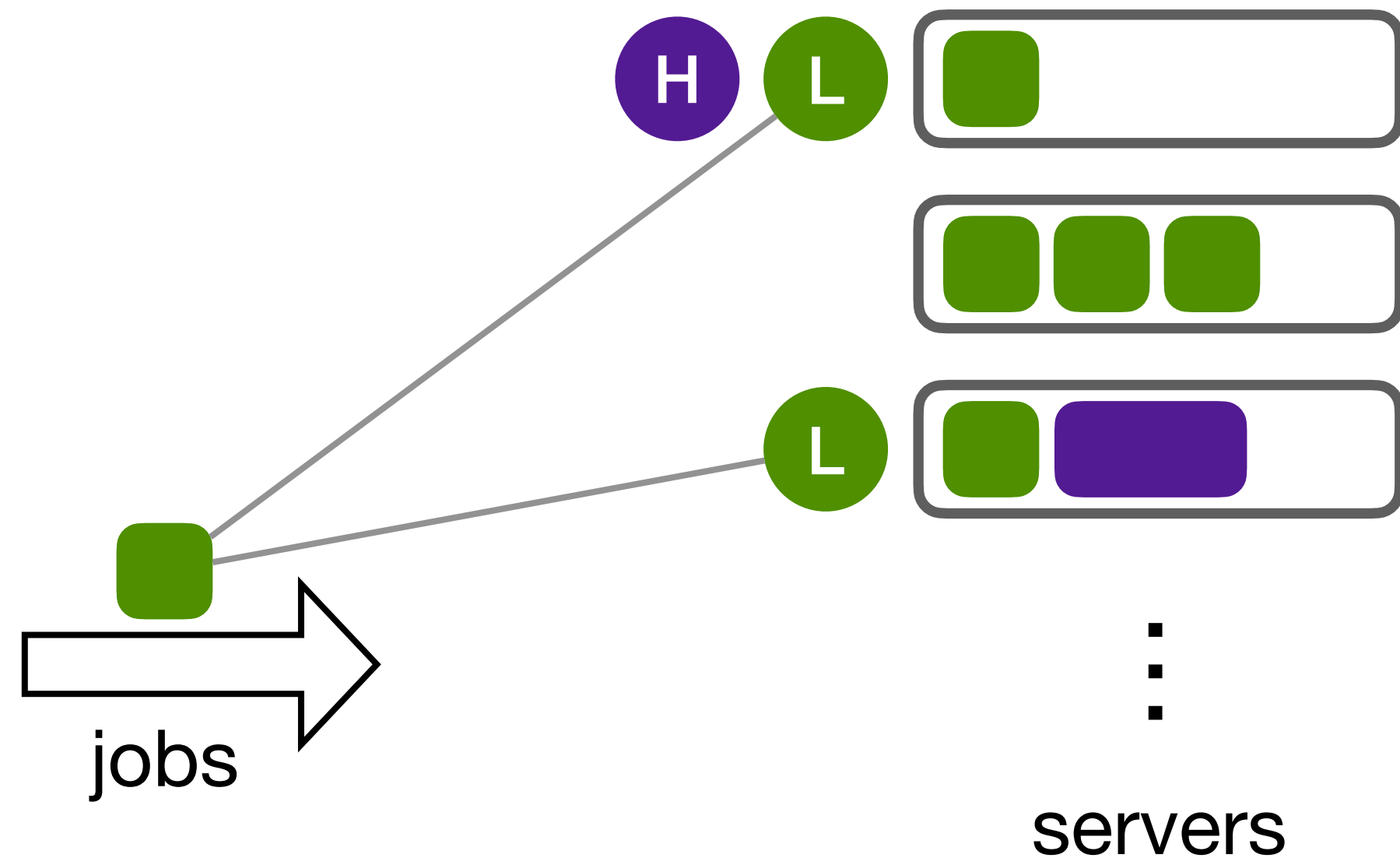


Run single-server policy $\bar{\sigma}$ for only

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

(regular servers)

Policy conversion: more details



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

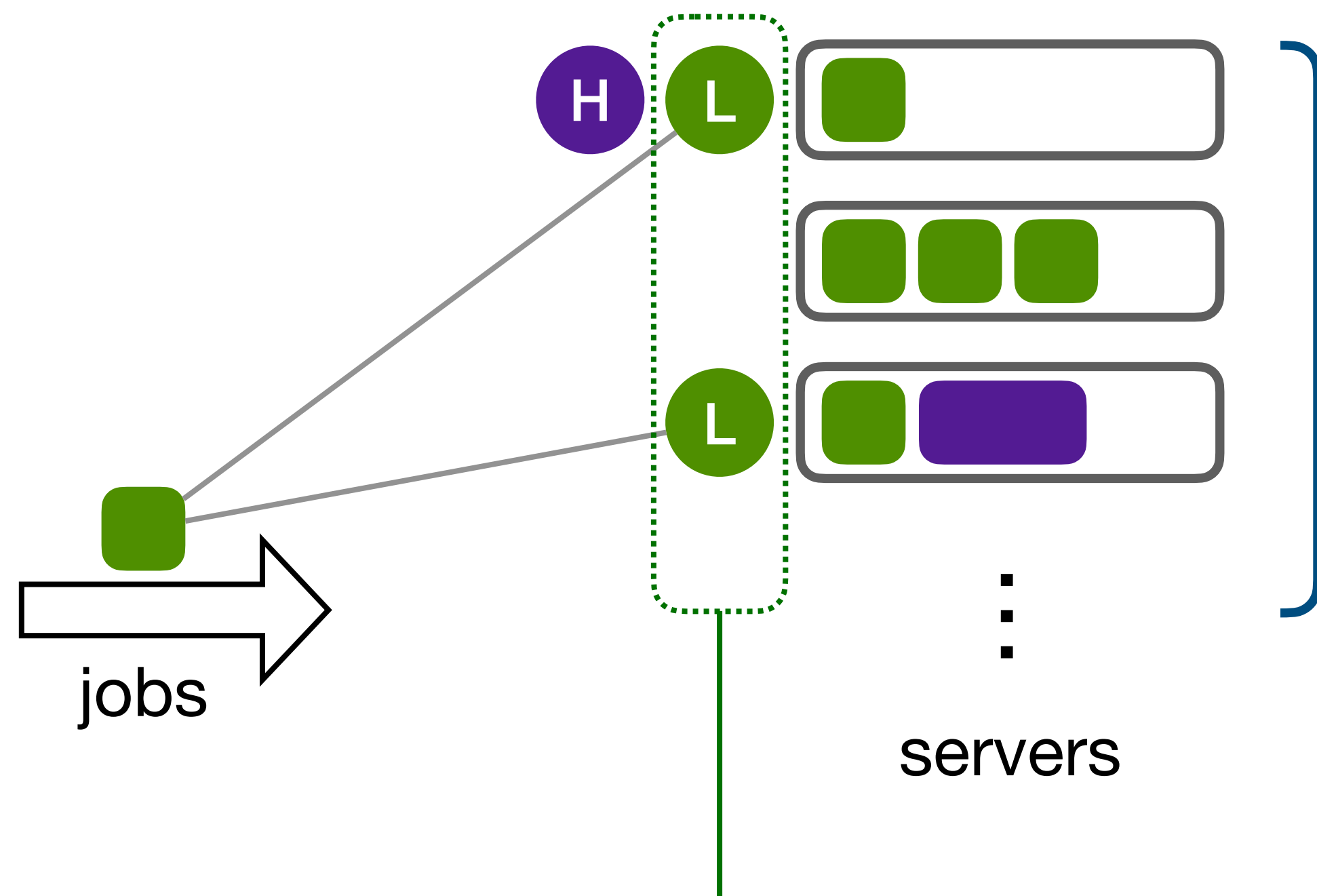
$$\bar{N} = \frac{\text{arrival rate}}{\text{throughput}(\bar{\sigma})} \text{ servers}$$

(regular servers)

Recall that we aim to show

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

Policy conversion: more details



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

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

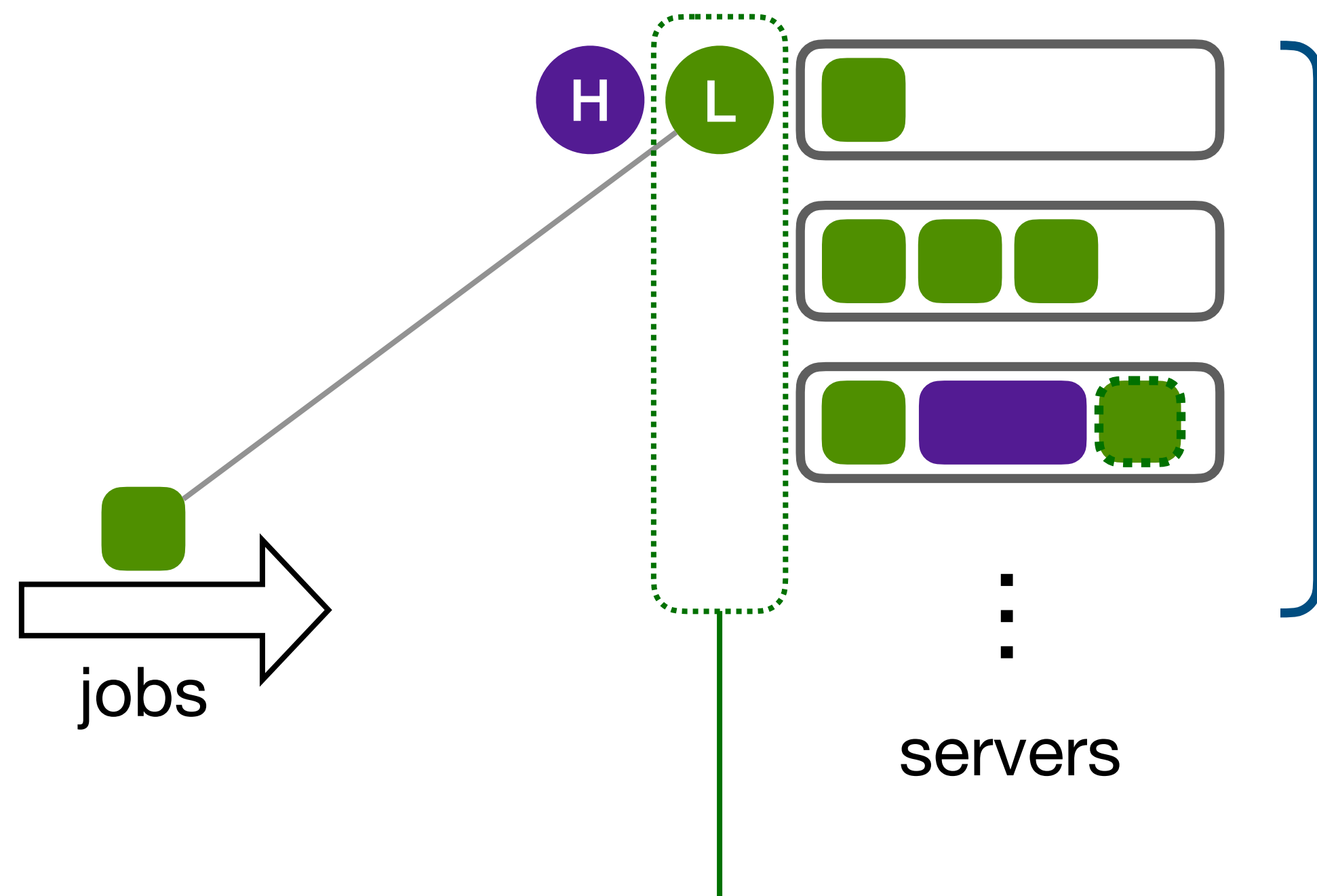
(regular servers)

Recall that we aim to show

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

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

Policy conversion: more details



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

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

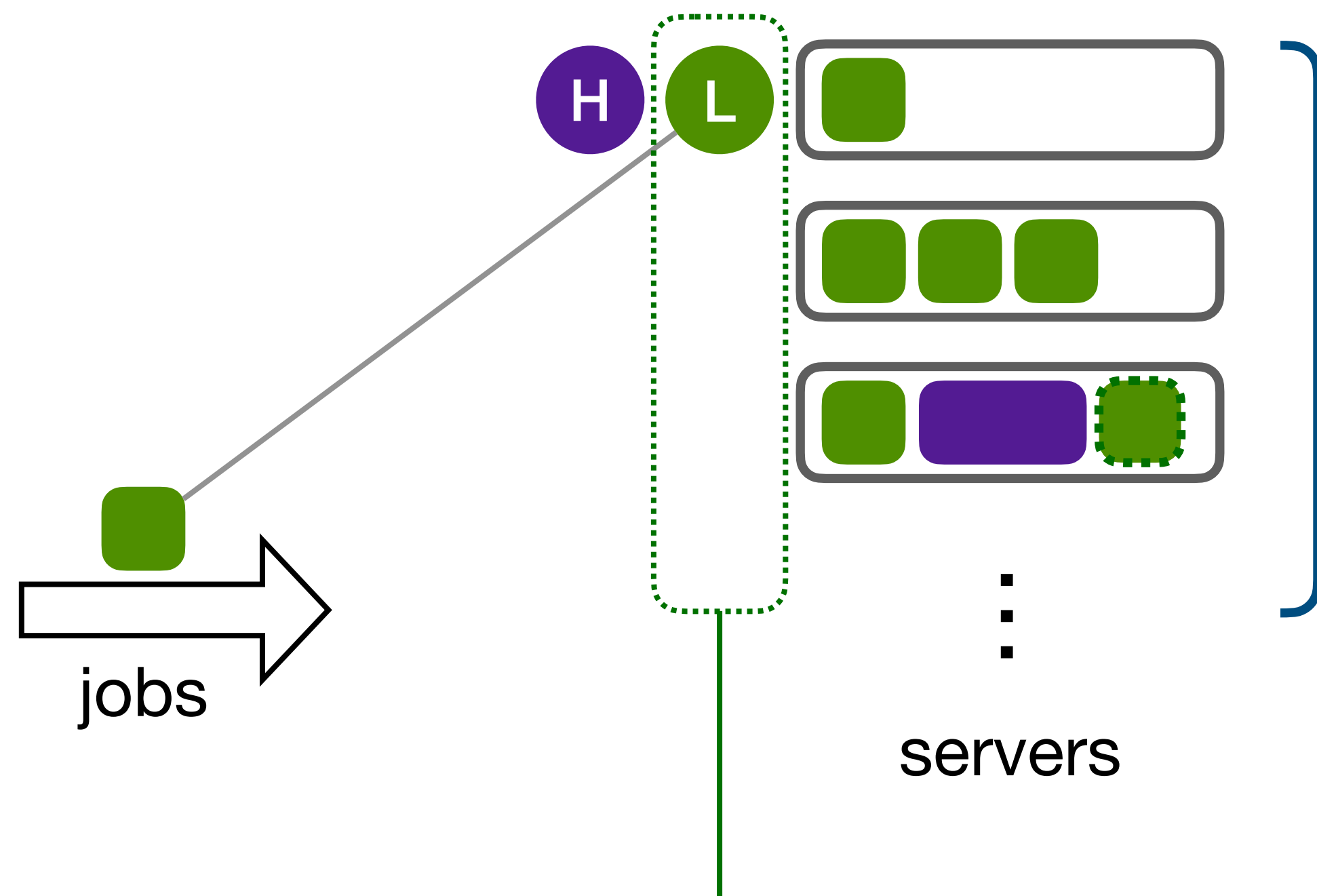
(regular servers)

Recall that we aim to show

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

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

Policy conversion: more details



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

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

(regular servers)

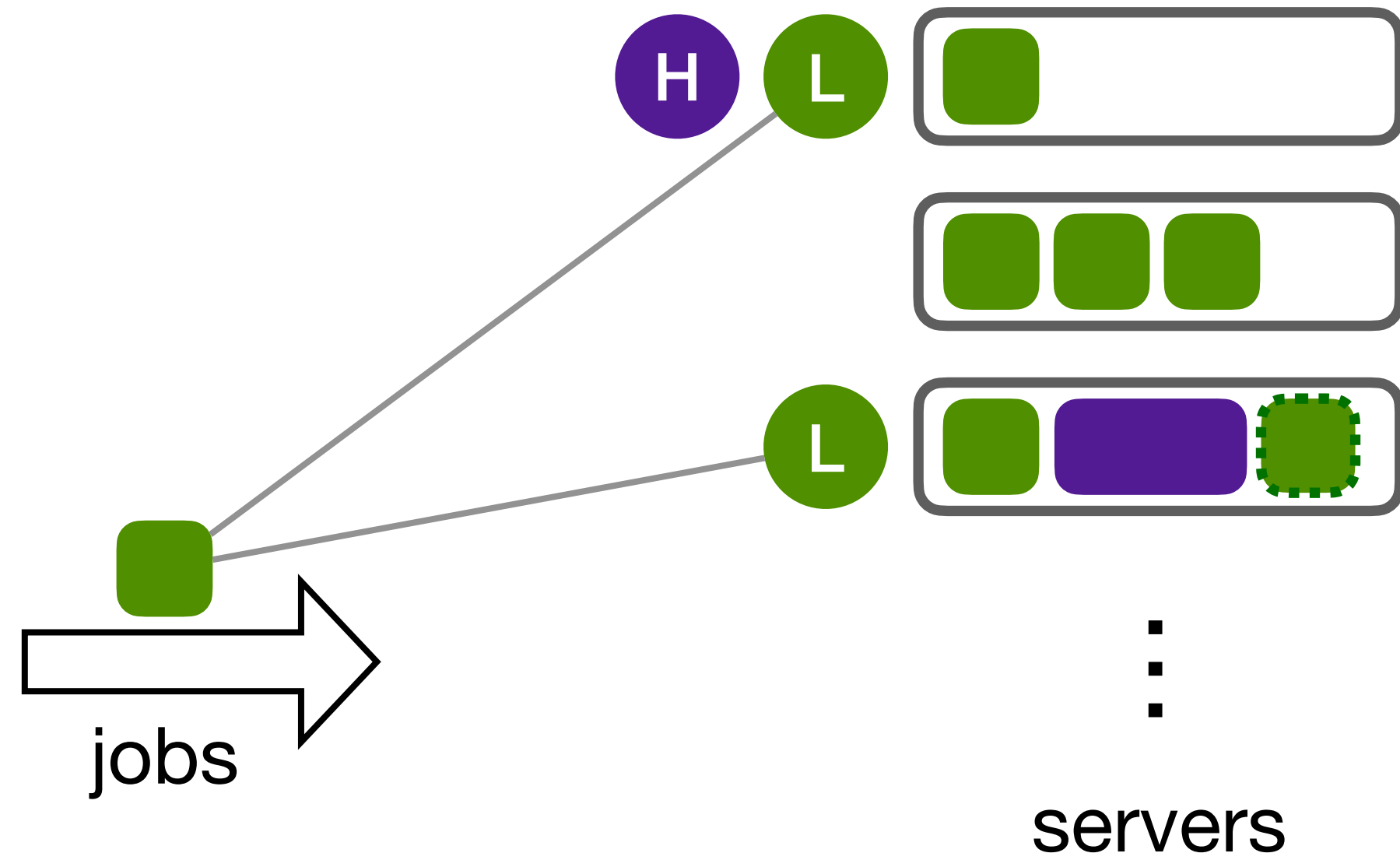
Recall that we aim to show

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

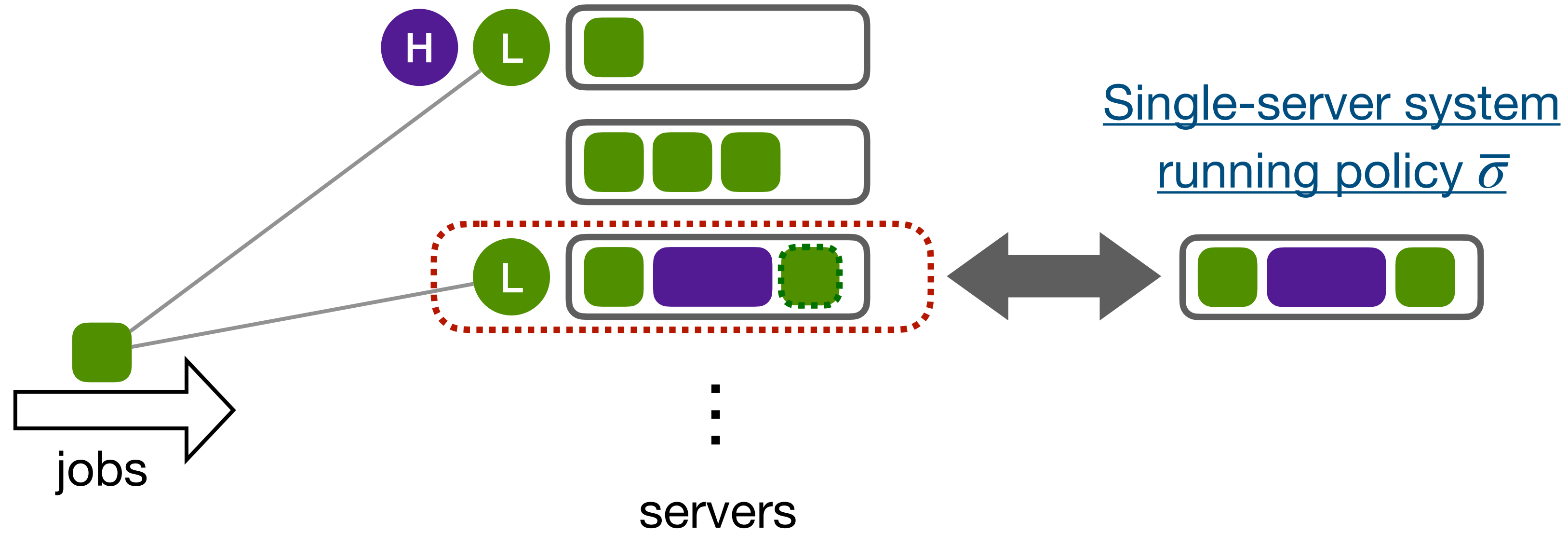
When the # tokens of a type $> \sqrt{r}$, remove the overflow tokens and generate virtual jobs

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

Key proof idea 1

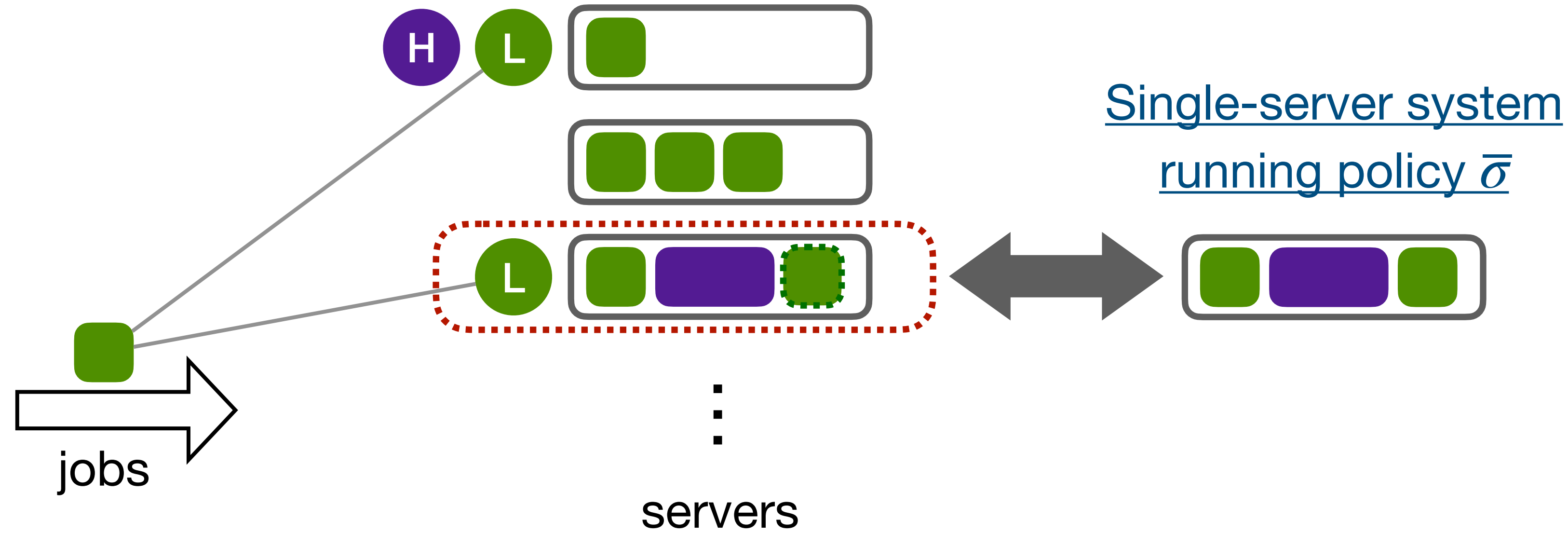


Key proof idea 1



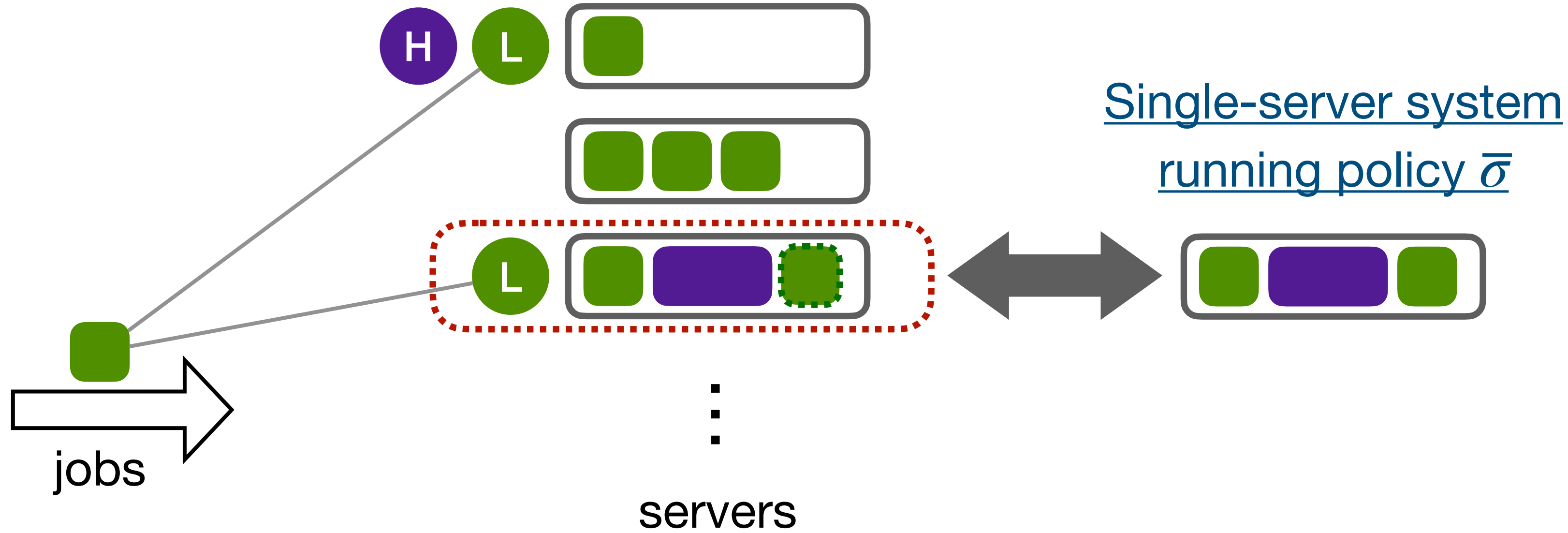
Key proof idea 1

Will show that regular servers in the original system $\approx \bar{N}$ independent single-server systems



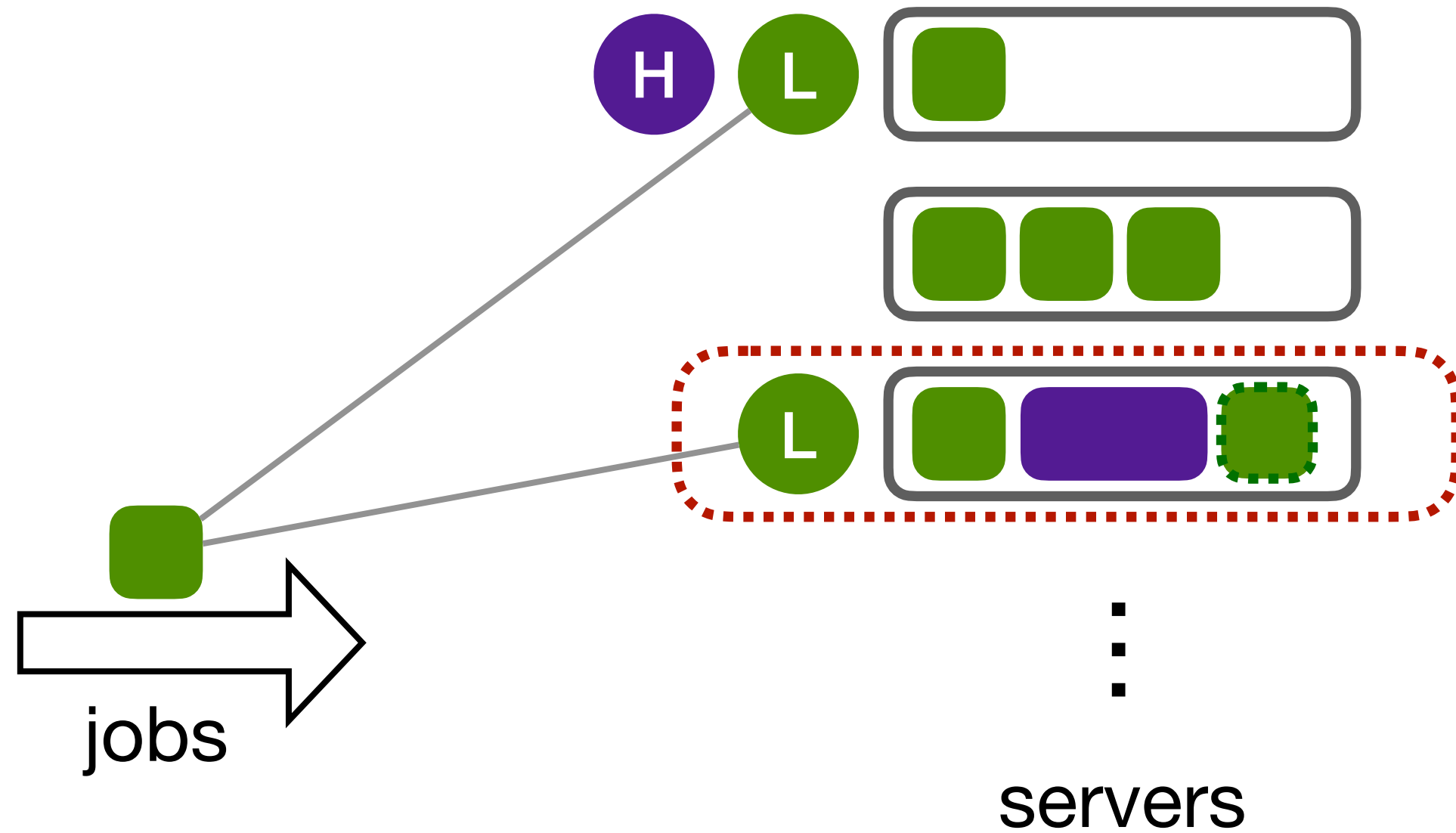
Key proof idea 1

Will show that **regular servers in the original system**
 $\approx \bar{N}$ independent single-server systems



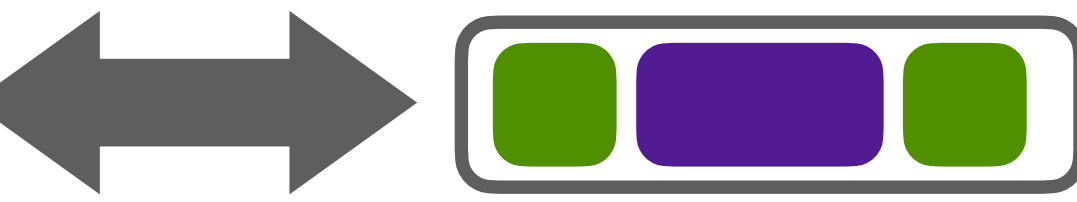
If only each token were replaced by a job immediately ...

Key proof idea 1



Will show that **regular servers in the original system**
 $\approx \bar{N}$ independent single-server systems

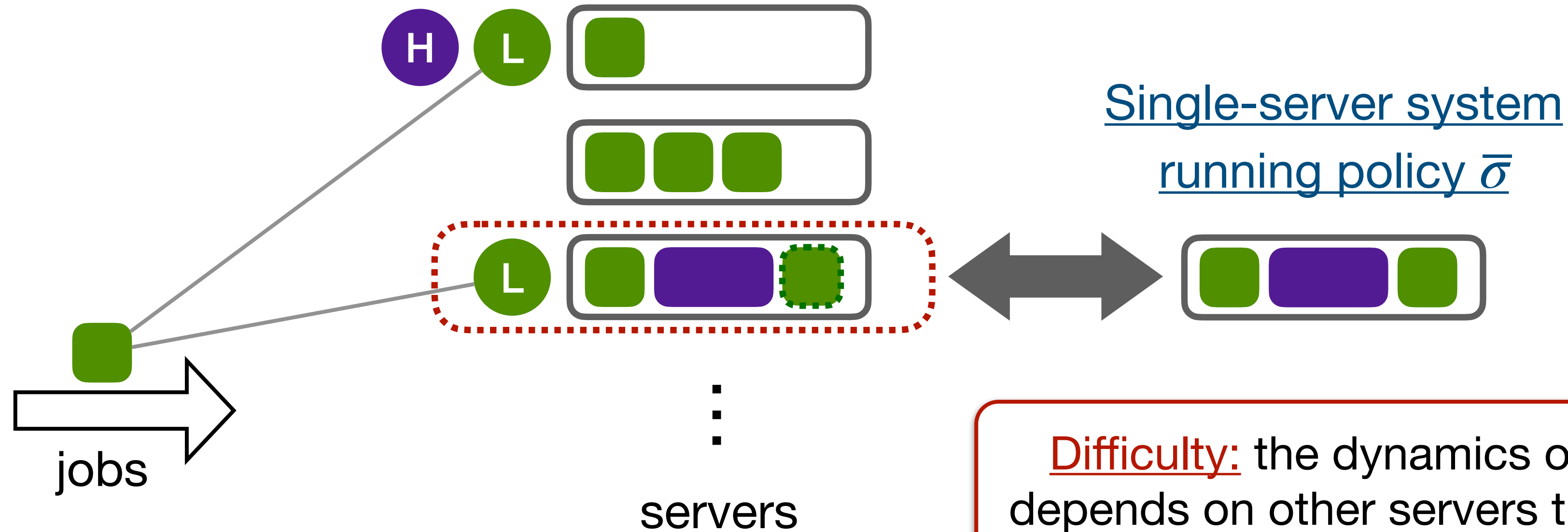
Single-server system
running policy $\bar{\sigma}$



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

Key proof idea 1

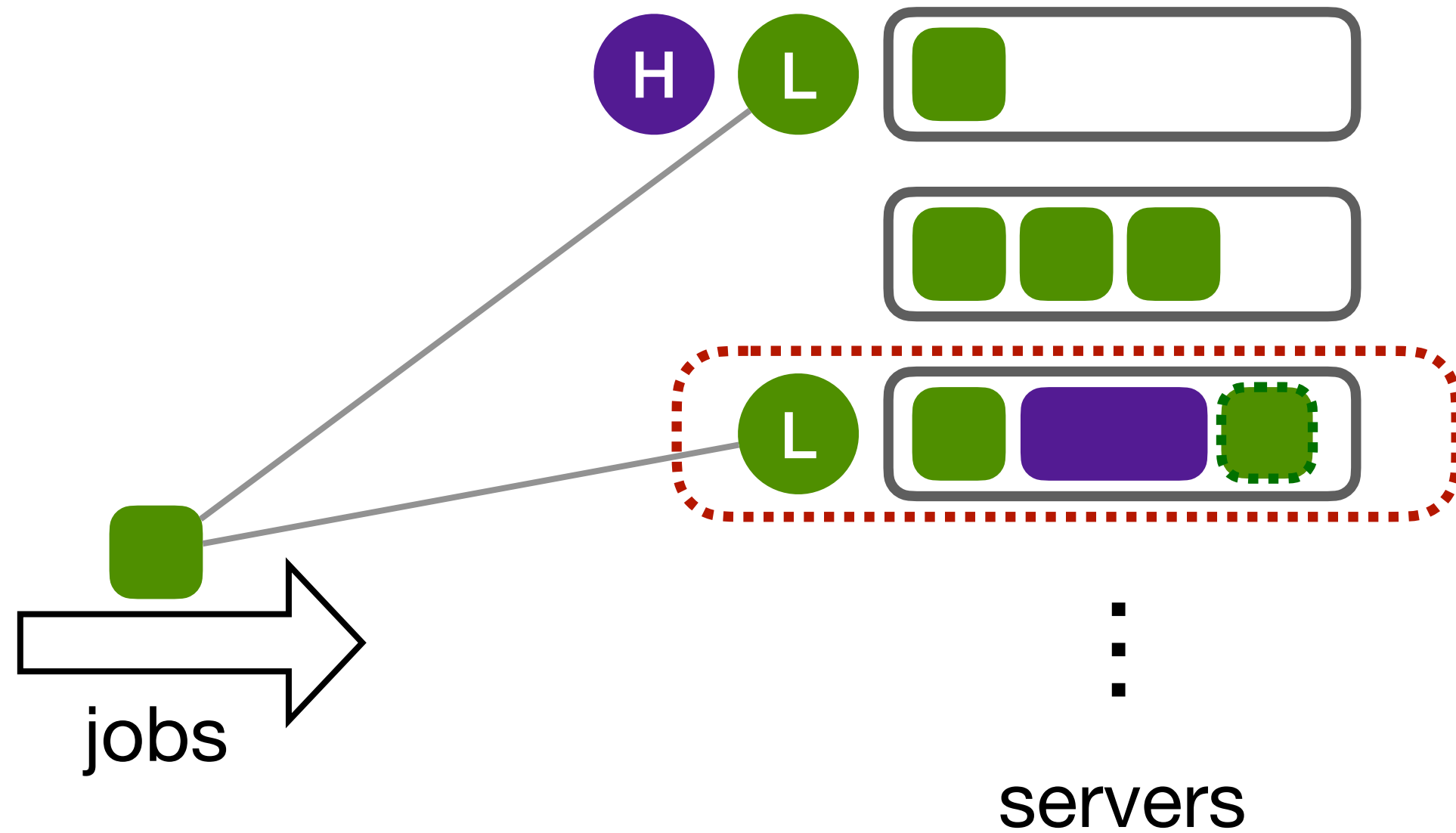
Will show that **regular servers in the original system**
 $\approx \bar{N}$ independent single-server systems



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

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

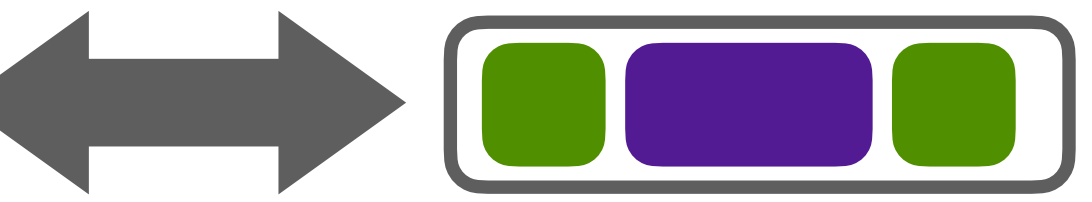
Key proof idea 1



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

Will show that **regular servers in the original system**
 $\approx \bar{N}$ independent single-server systems

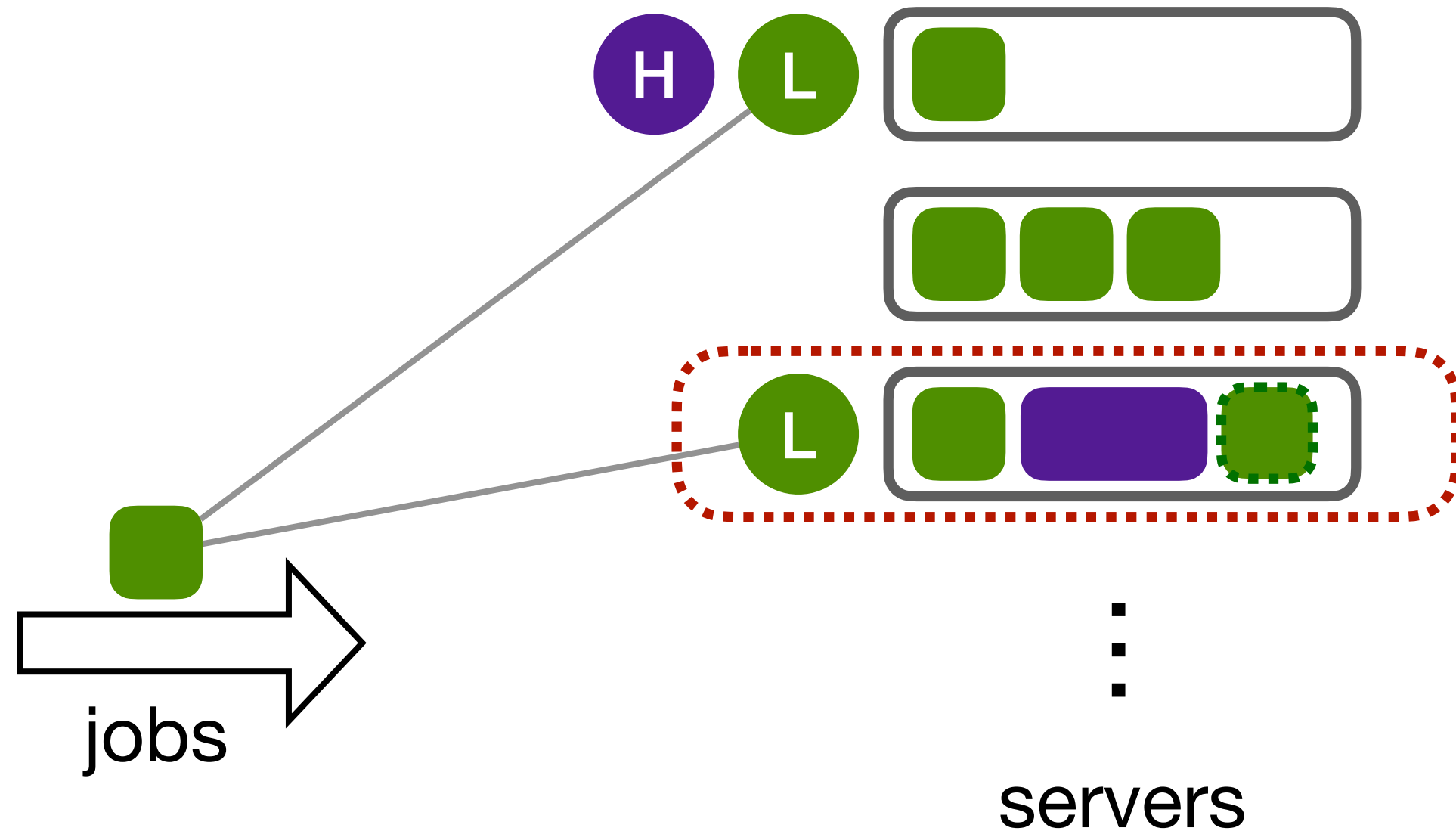
Single-server system
 running policy $\bar{\sigma}$



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

Do job arrivals affect \widetilde{K}_i ?

Key proof idea 1



Will show that **regular servers in the original system**
 $\approx \bar{N}$ independent single-server systems

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

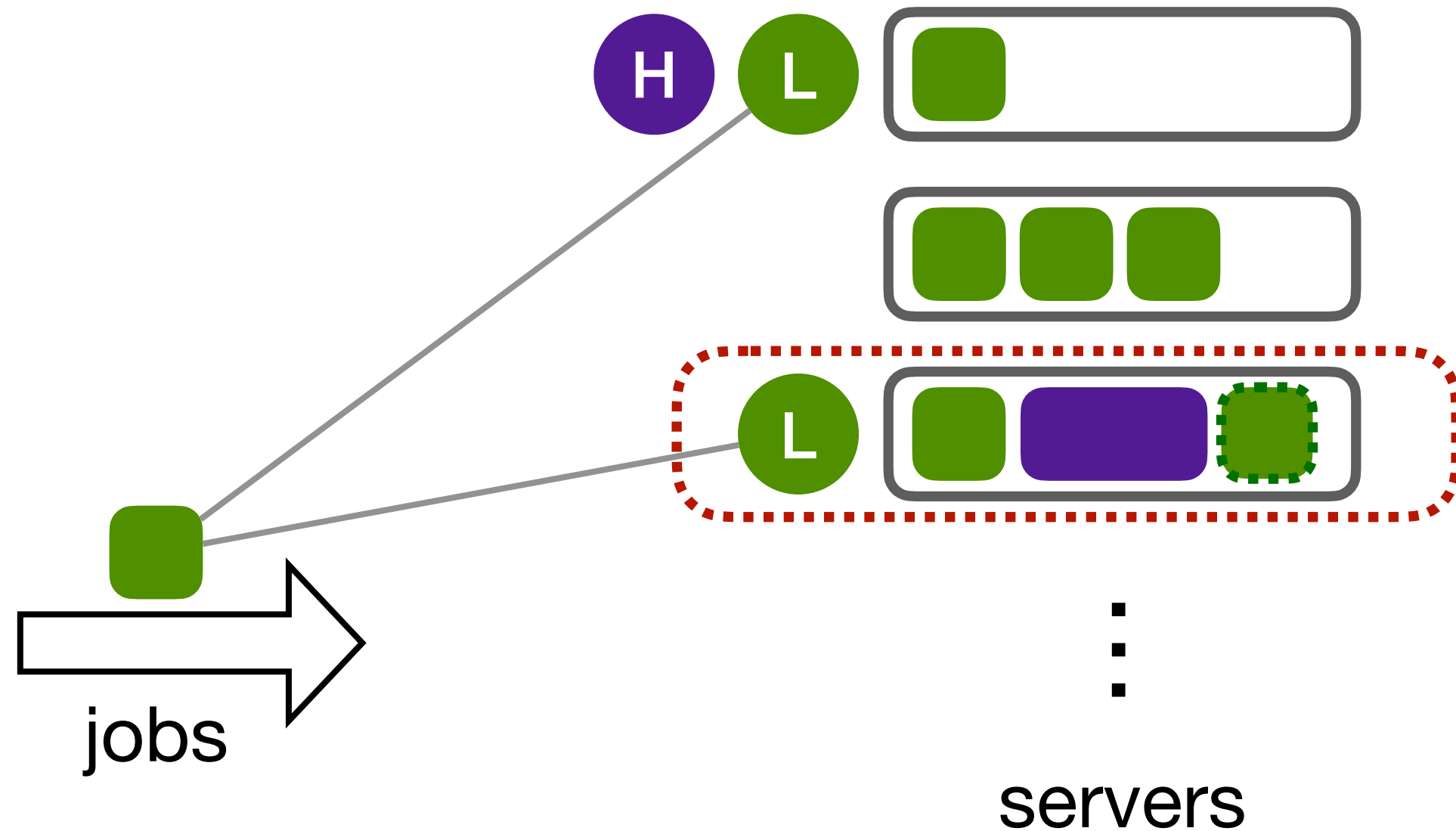
Difficulty: the dynamics of a server in the original system depends on other servers through arrivals & token overflows

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

Do job arrivals affect \tilde{K}_i ?

Do token overflows affect \tilde{K}_i ?

Key proof idea 1



Will show that **regular servers in the original system**
 $\approx \bar{N}$ independent single-server systems

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

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

Do job arrivals affect \tilde{K}_i ?

Do token overflows affect \tilde{K}_i ?

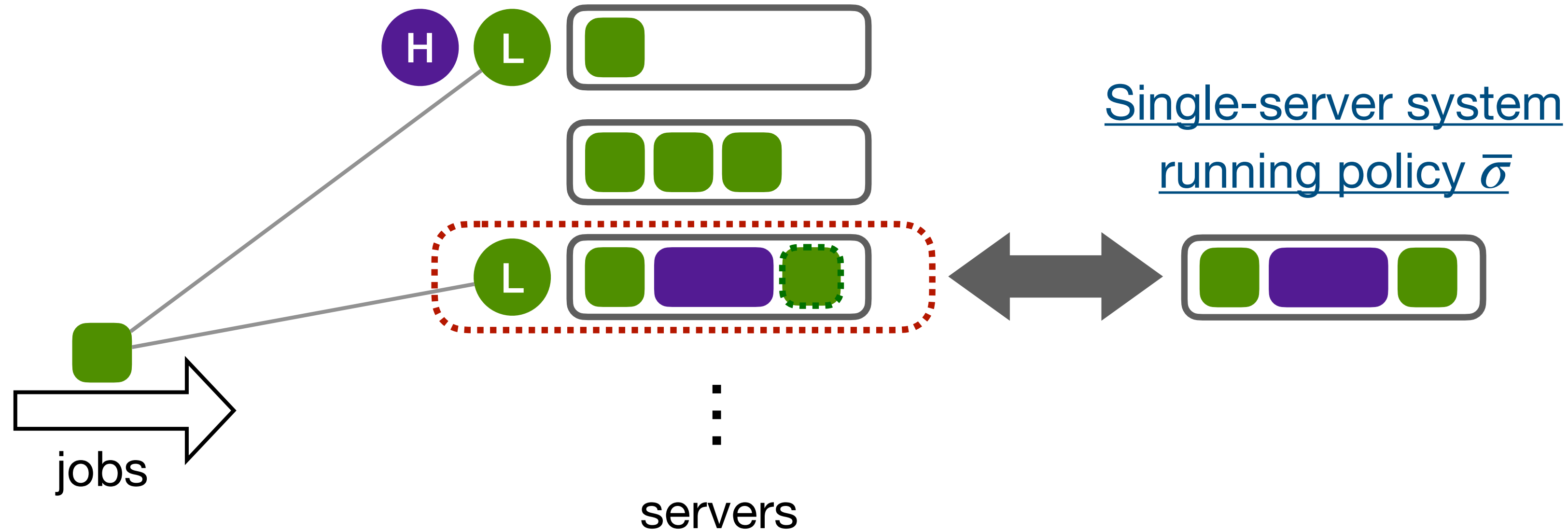
Idea: for each type i , consider

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

Arrivals & token overflows do not affect \tilde{K}_i

Key proof idea 1

Will show that regular servers in the original system $\approx \bar{N}$ independent single-server systems

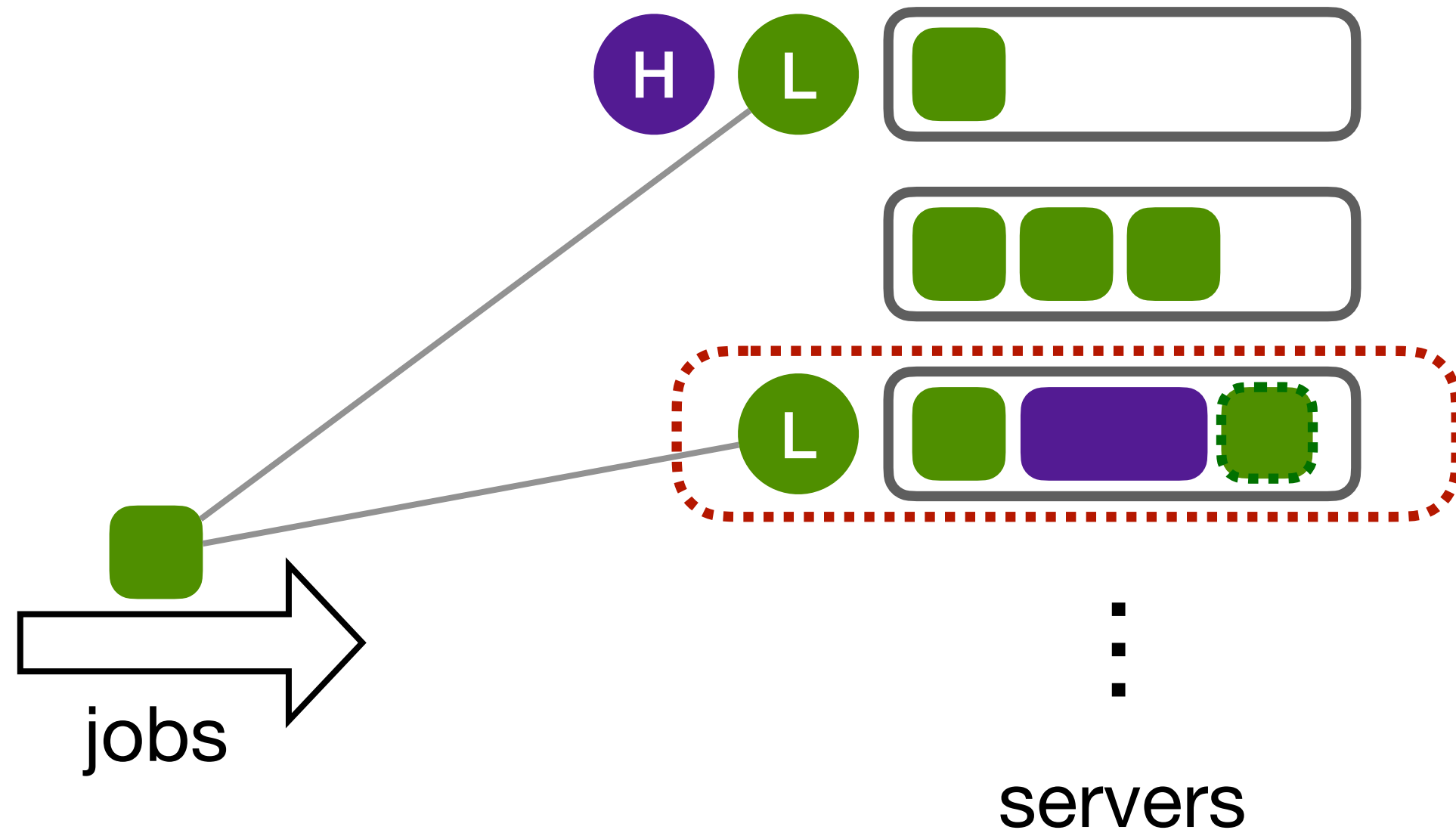


Idea: for each type i , consider

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

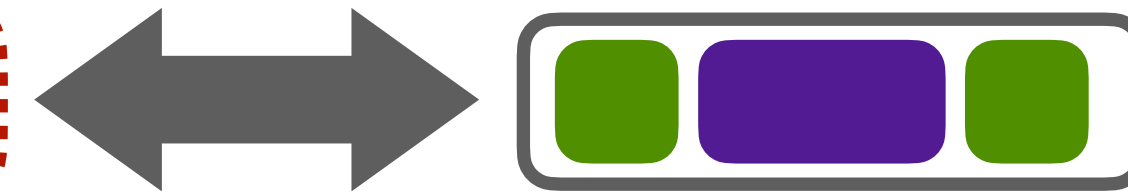
Arrivals & token overflows do not affect \widetilde{K}_i

Key proof idea 1



Will show that **regular servers in the original system**
 $\approx \bar{N}$ independent single-server systems

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

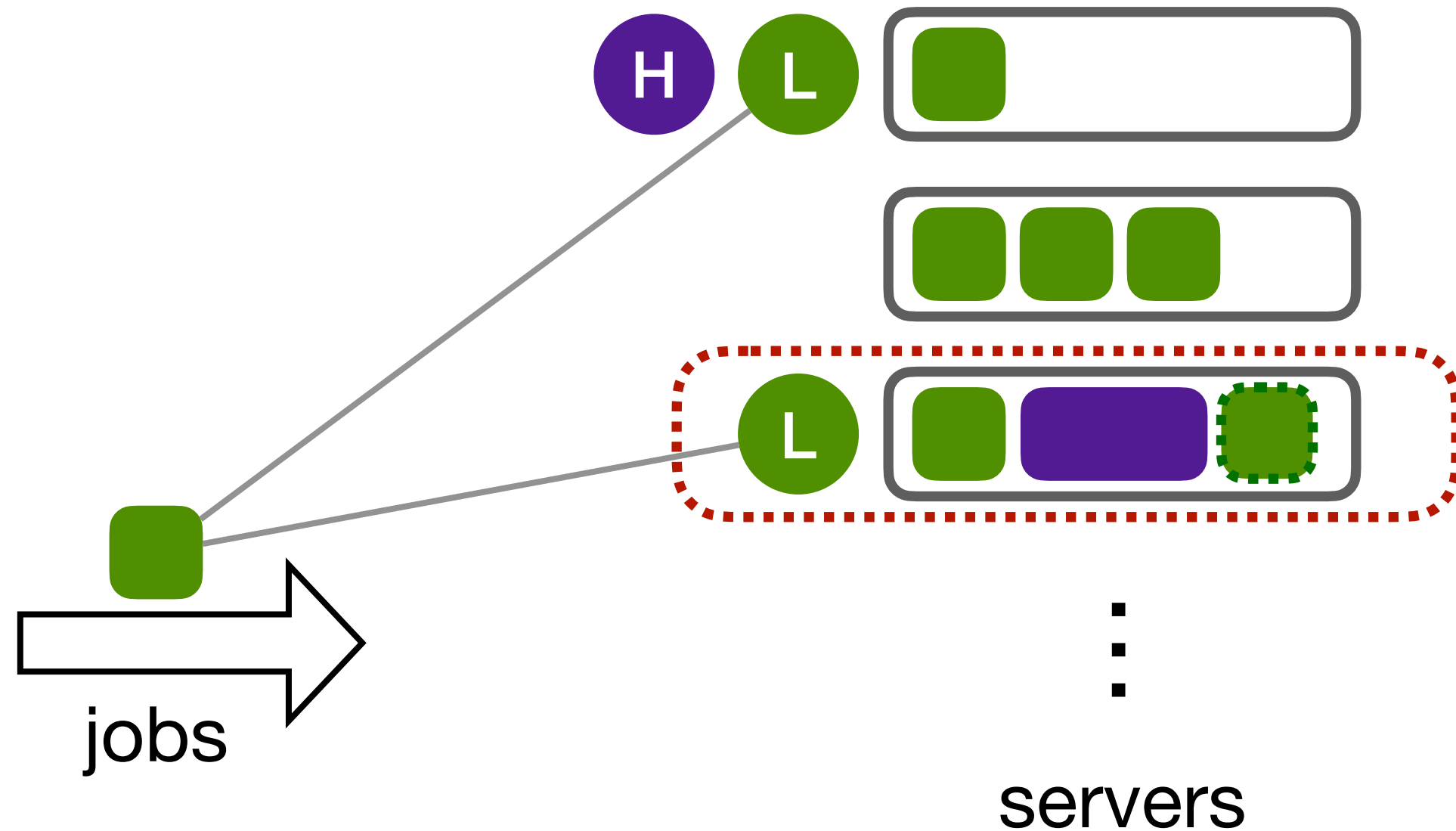


Dynamics in the original system v.s.
Dynamics in \bar{N} independent single-server systems

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

Key proof idea 1

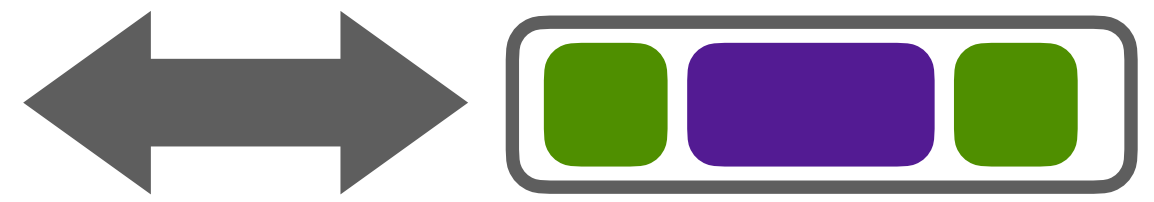


Idea: for each type i , consider
 $\widetilde{K}_i = \# \text{ jobs} + \# \text{ virtual jobs} + \# \text{ tokens}$

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

Will show that **regular servers in the original system**
 $\approx \bar{N}$ independent single-server systems

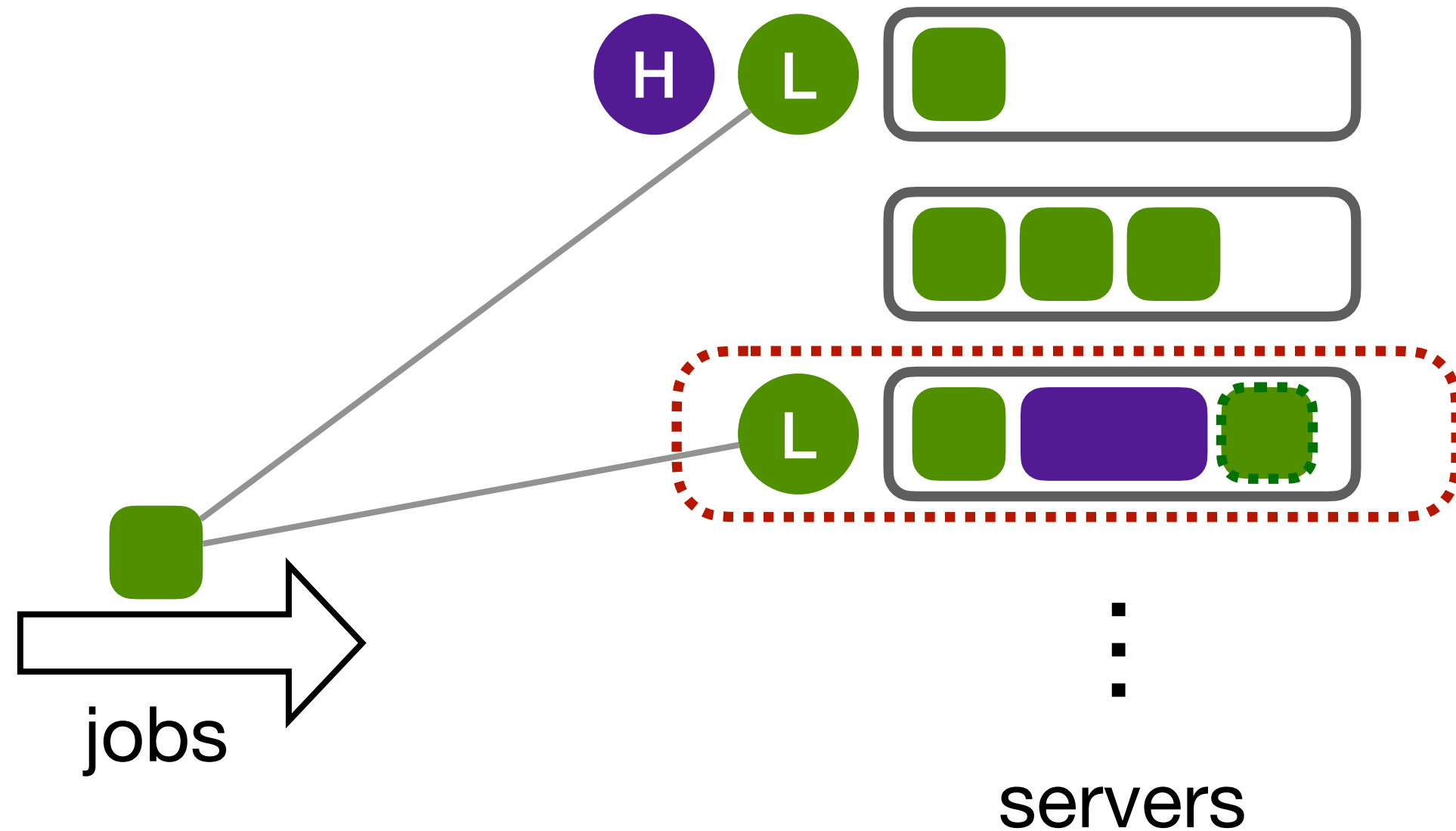
Single-server system
 running policy $\bar{\sigma}$



Dynamics in the original system v.s.
Dynamics in \bar{N} independent single-server systems

- A server without tokens in the original system generates tokens and has job transitions in the same way as a single-server system

Key proof idea 1



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

Will show that **regular servers in the original system**
 $\approx \bar{N}$ independent single-server systems

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

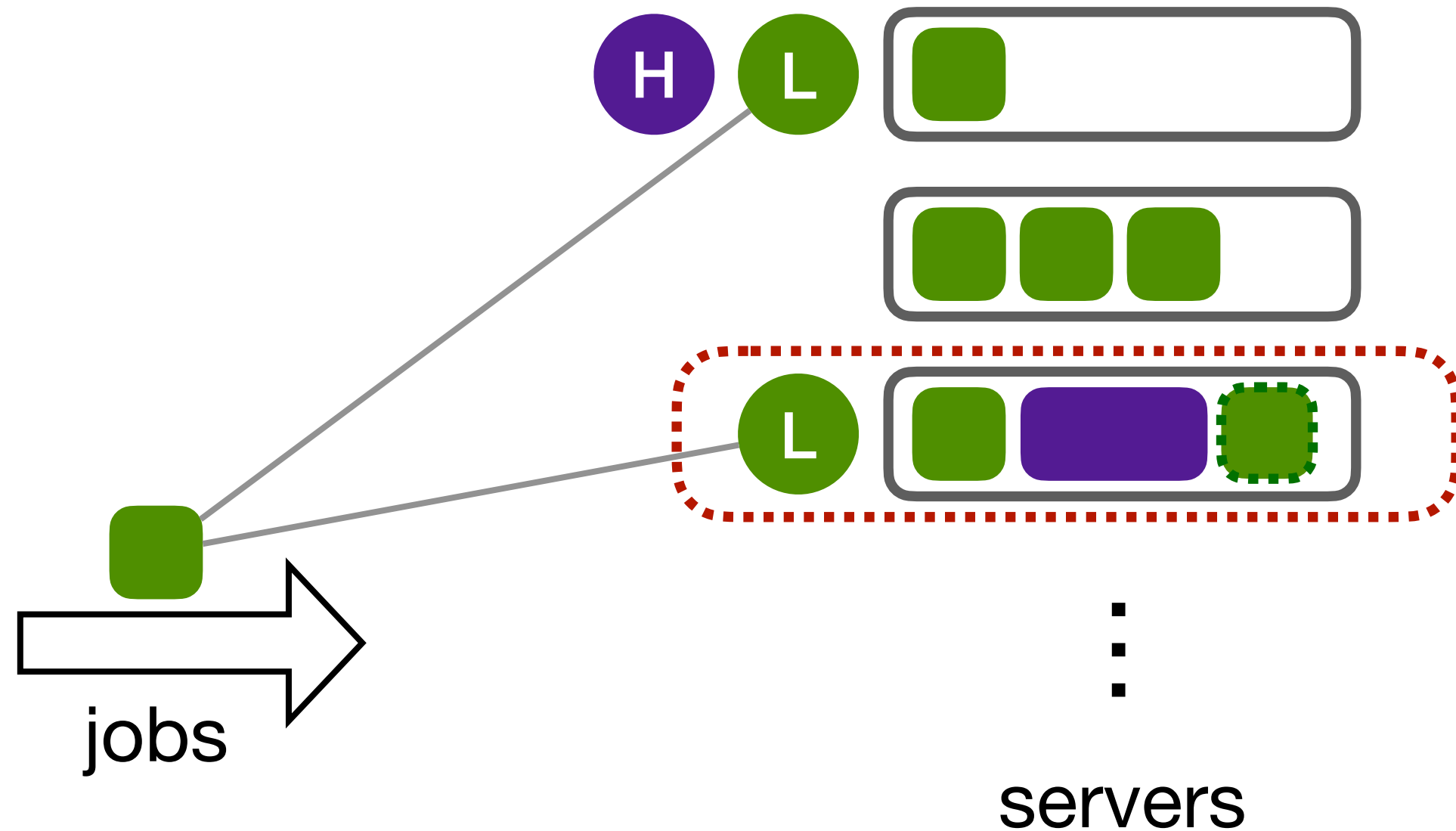


Dynamics in the original system v.s.

Dynamics in \bar{N} independent single-server systems

- A server without tokens in the original system generates tokens and has job transitions in the same way as a single-server system
- Difference between two systems is bounded by **# tokens**

Key proof idea 1



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

Will show that **regular servers in the original system**
 $\approx \bar{N}$ independent single-server systems

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

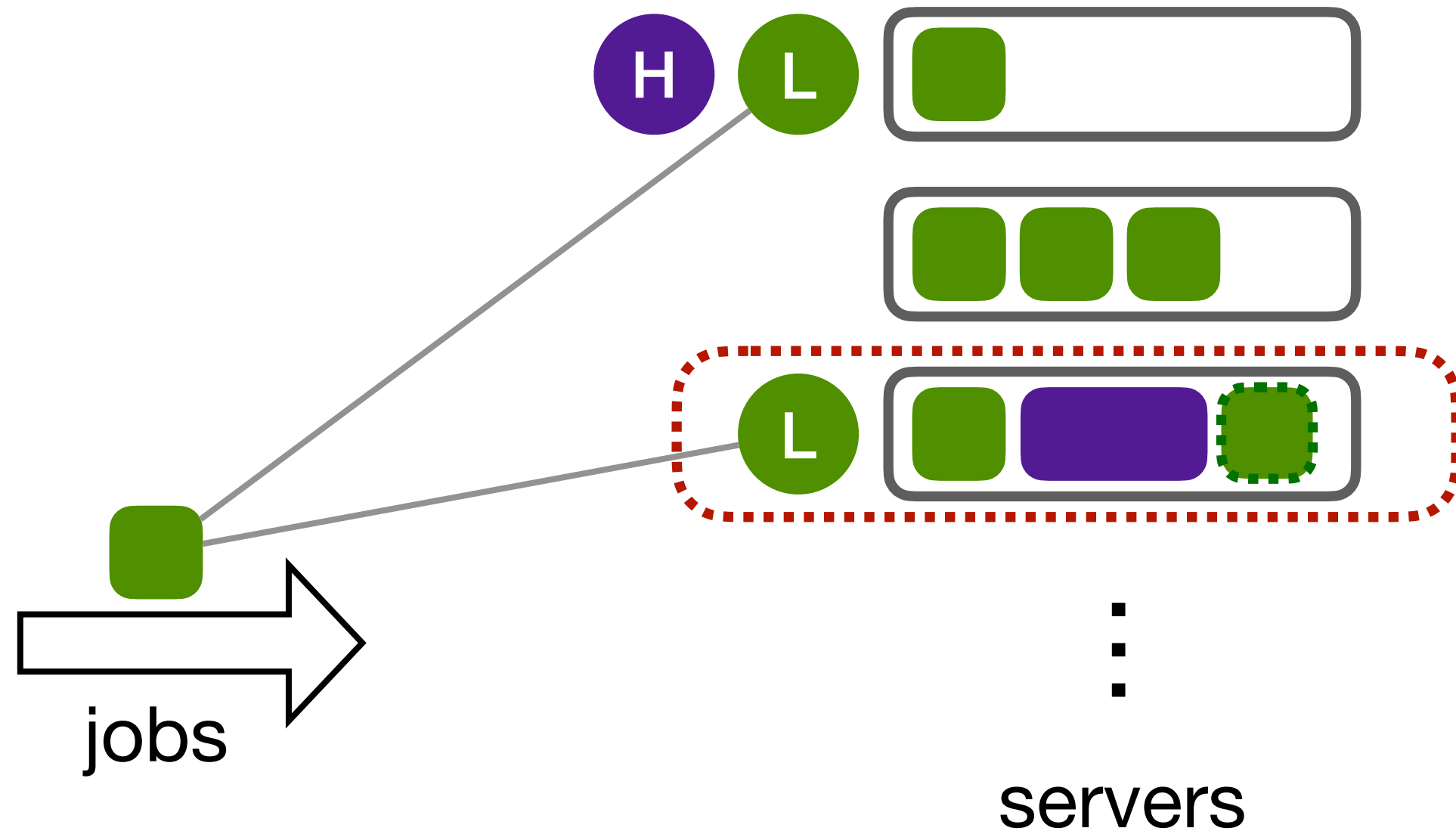


Dynamics in the original system v.s.

Dynamics in \bar{N} independent single-server systems

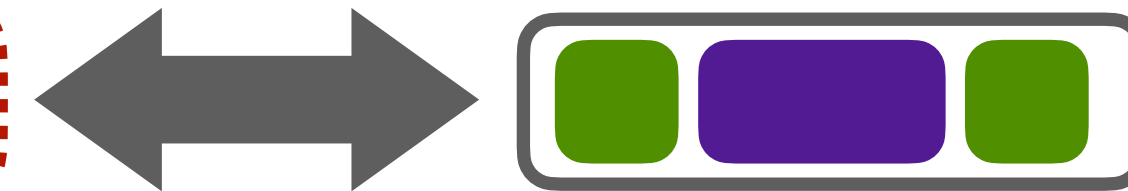
- A server without tokens in the original system generates tokens and has job transitions in the same way as a single-server system
- Difference between two systems is bounded by **# tokens**
- **# tokens** $\leq \sqrt{r}$

Key proof idea 1



Will show that **regular servers in the original system**
 $\approx \bar{N}$ independent single-server systems

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



Dynamics in the original system v.s.

Dynamics in \bar{N} independent single-server systems

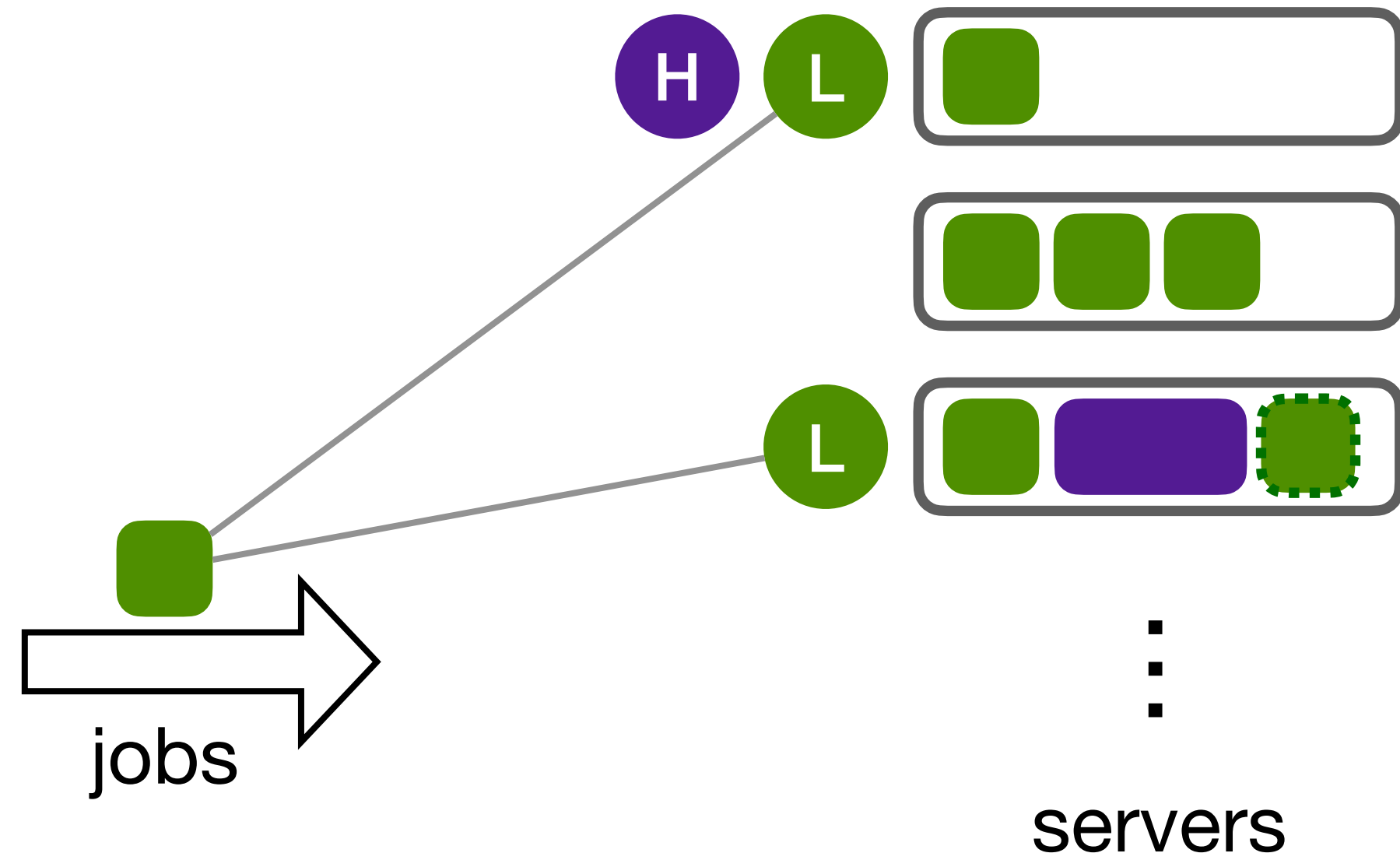
- A server without tokens in the original system generates tokens and has job transitions in the same way as a single-server system
- Difference between two systems is bounded by **# tokens**
- **# tokens** $\leq \sqrt{r}$

\Rightarrow Using Stein's method, $d_W \left(\bar{K}^{1:\bar{N}}, \bar{K}^{1:\bar{N}} \right) = O \left(r^{0.5} \right)$

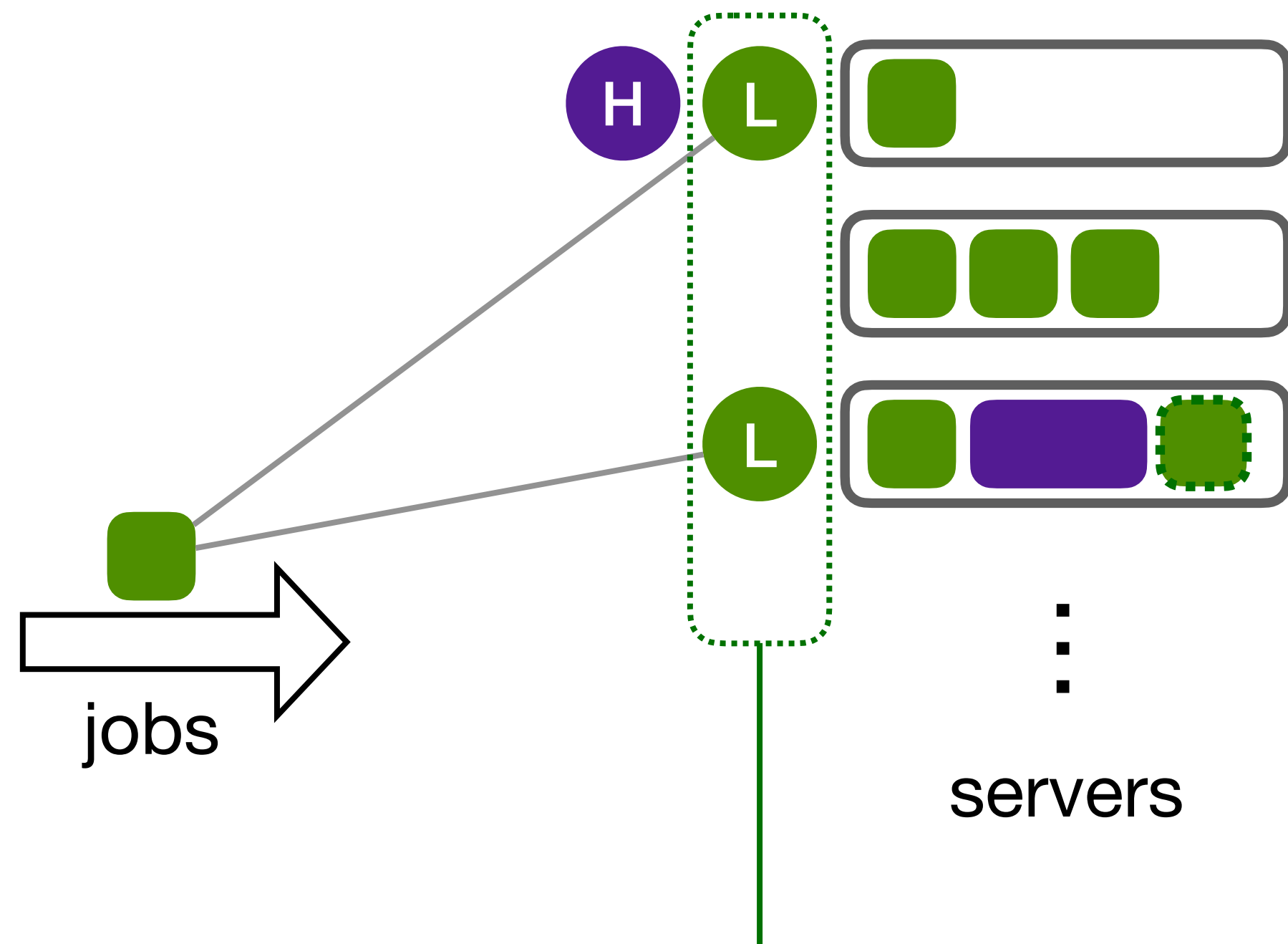
Idea: for each type i , consider
 $\tilde{K}_i = \# \text{ jobs} + \# \text{ virtual jobs} + \# \text{ tokens}$

Arrivals & token overflows do not affect \tilde{K}_i

Key proof idea 2

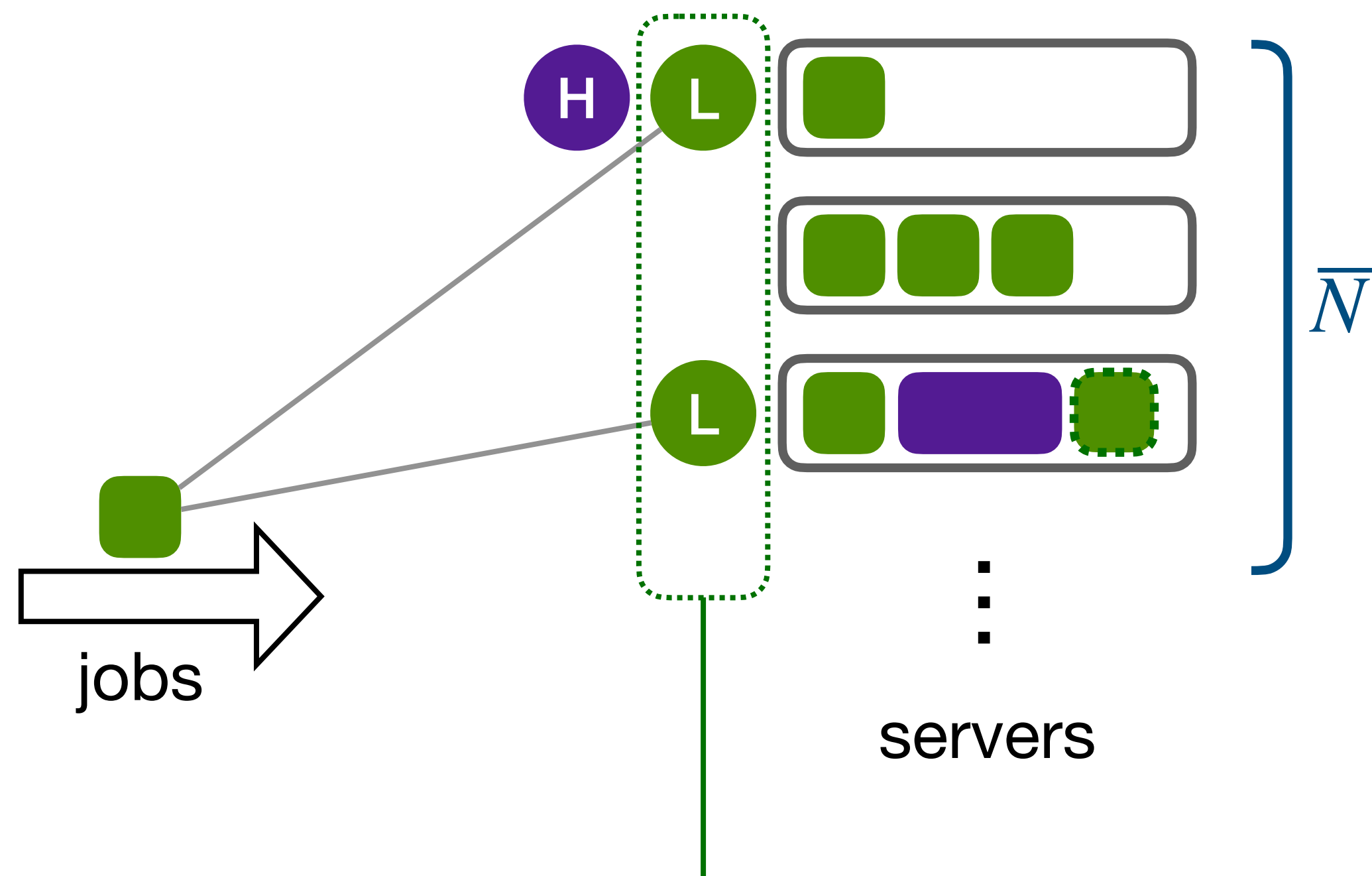


Key proof idea 2



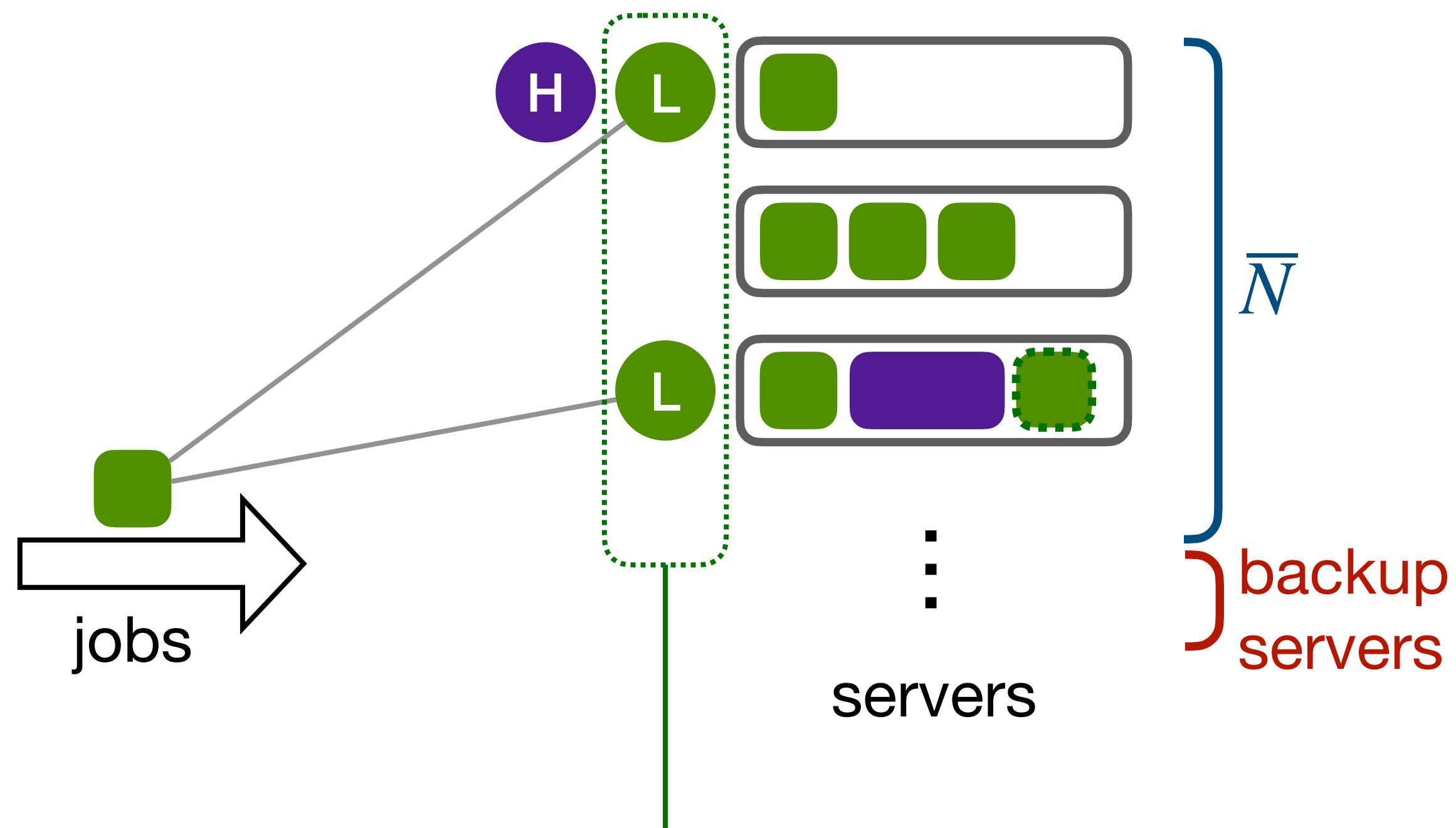
When the # tokens of a type $> \sqrt{r}$, remove the overflow tokens and generate virtual jobs

Key proof idea 2



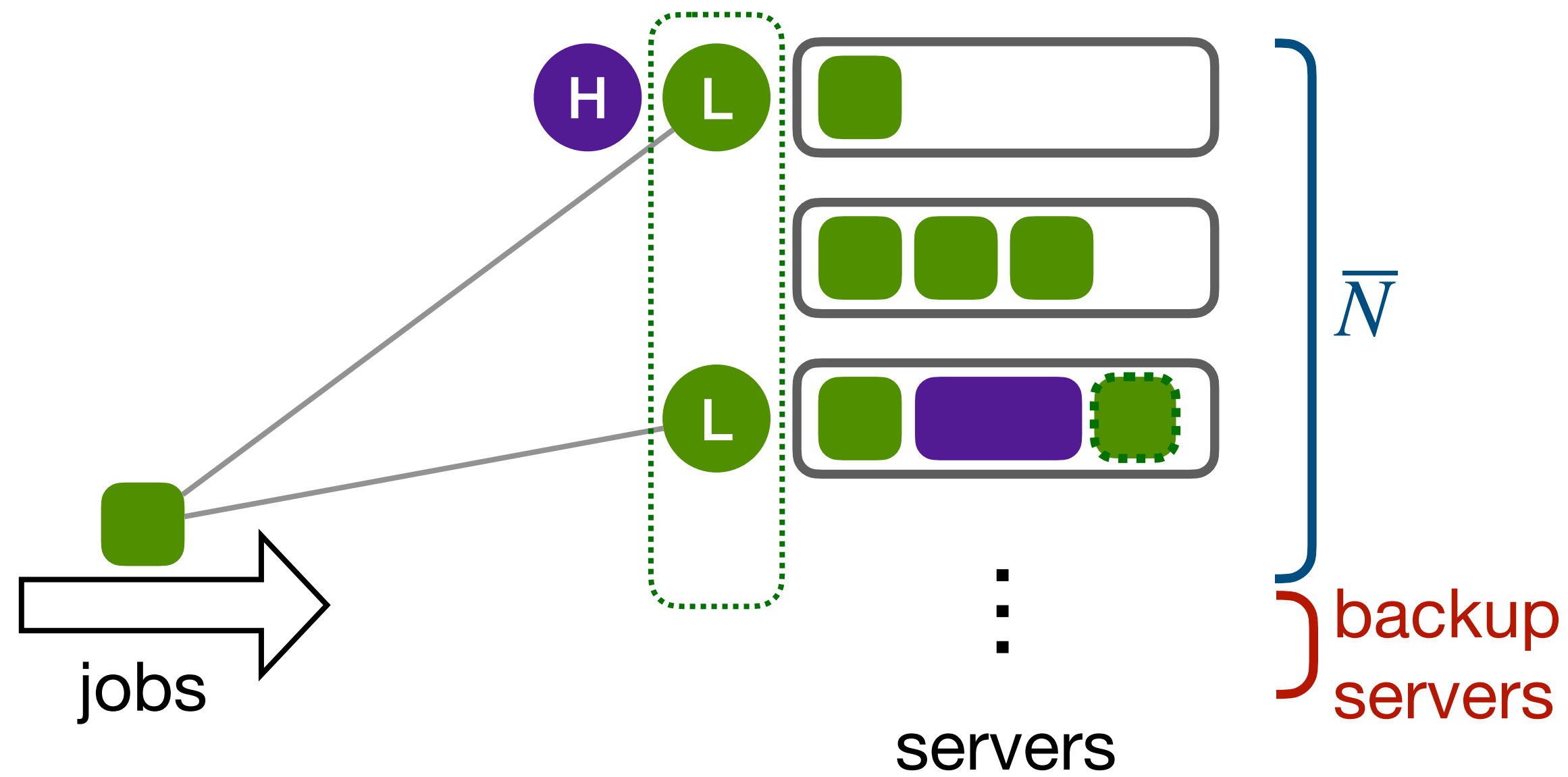
When the # tokens of a type $> \sqrt{r}$, remove the overflow tokens and generate virtual jobs

Key proof idea 2



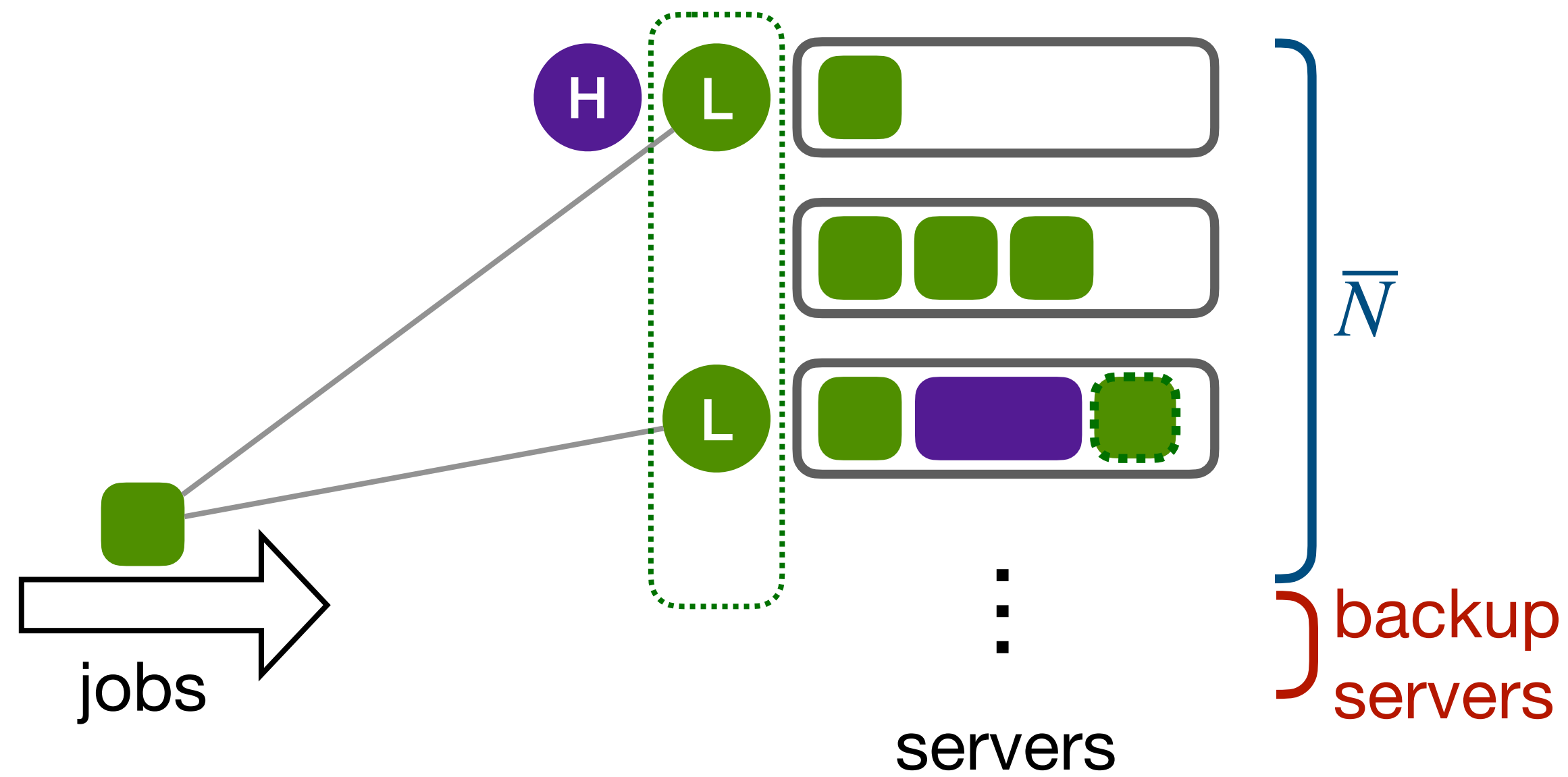
When the # tokens of a type $> \sqrt{r}$, remove the overflow tokens and generate virtual jobs

Key proof idea 2

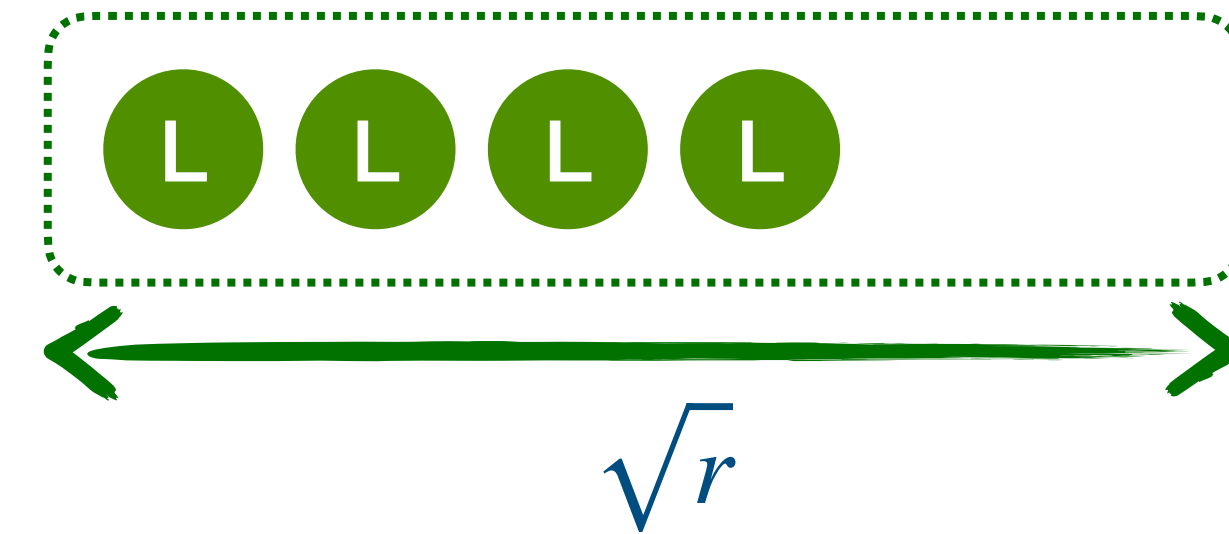


Will show that # virtual jobs = $O(\sqrt{r})$,
and # backup servers = $O(\sqrt{r})$

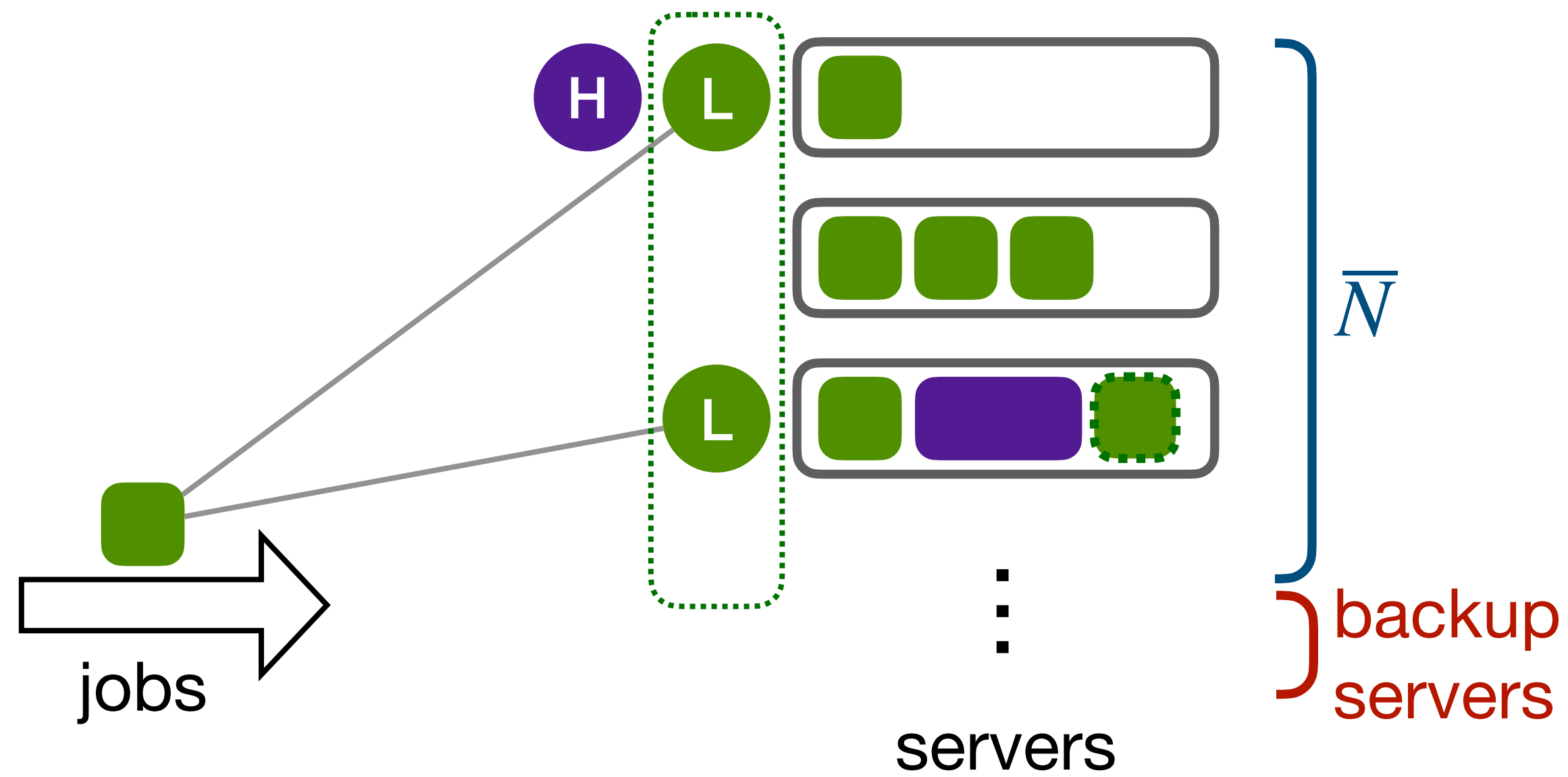
Key proof idea 2



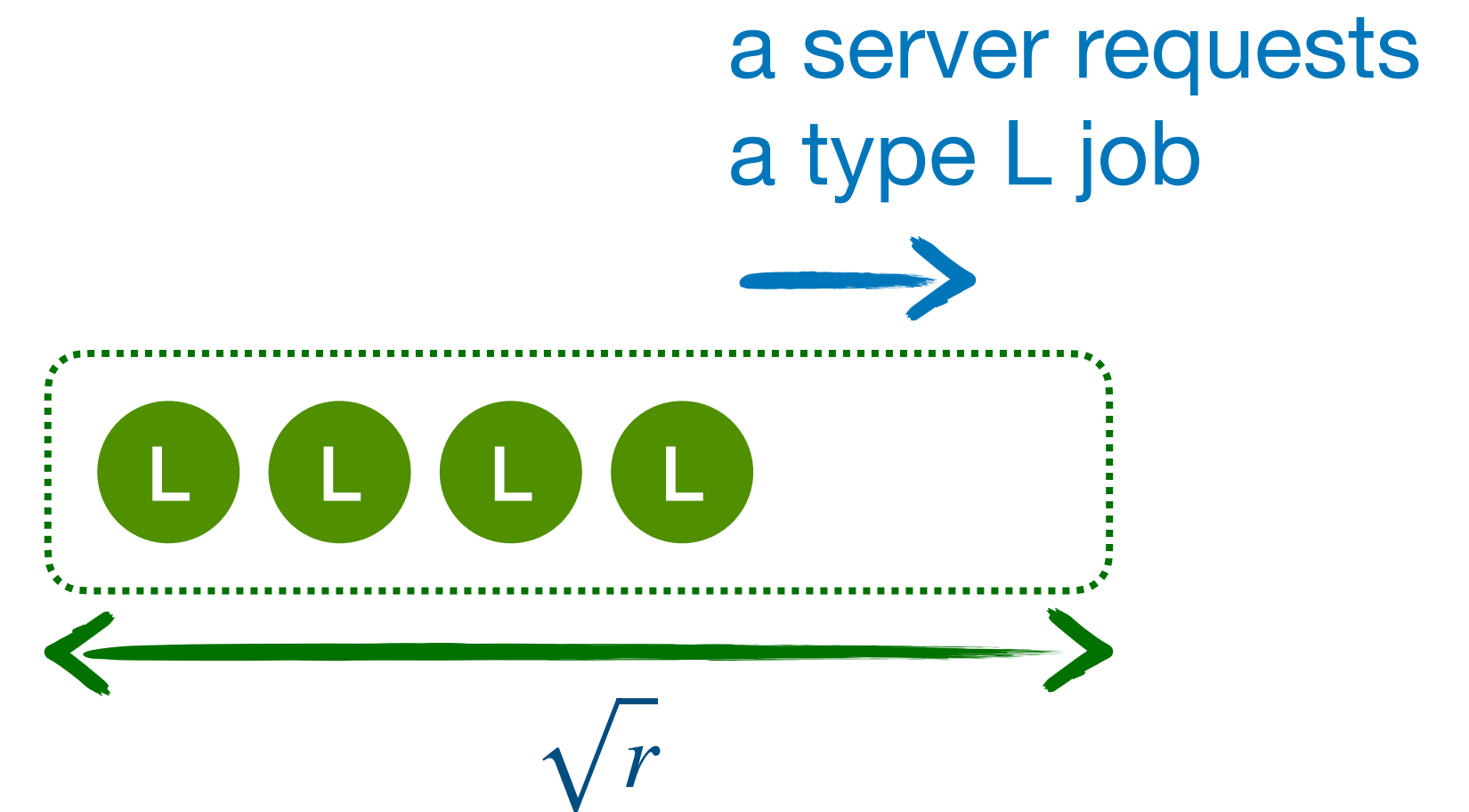
Will show that # virtual jobs = $O(\sqrt{r})$,
and # backup servers = $O(\sqrt{r})$



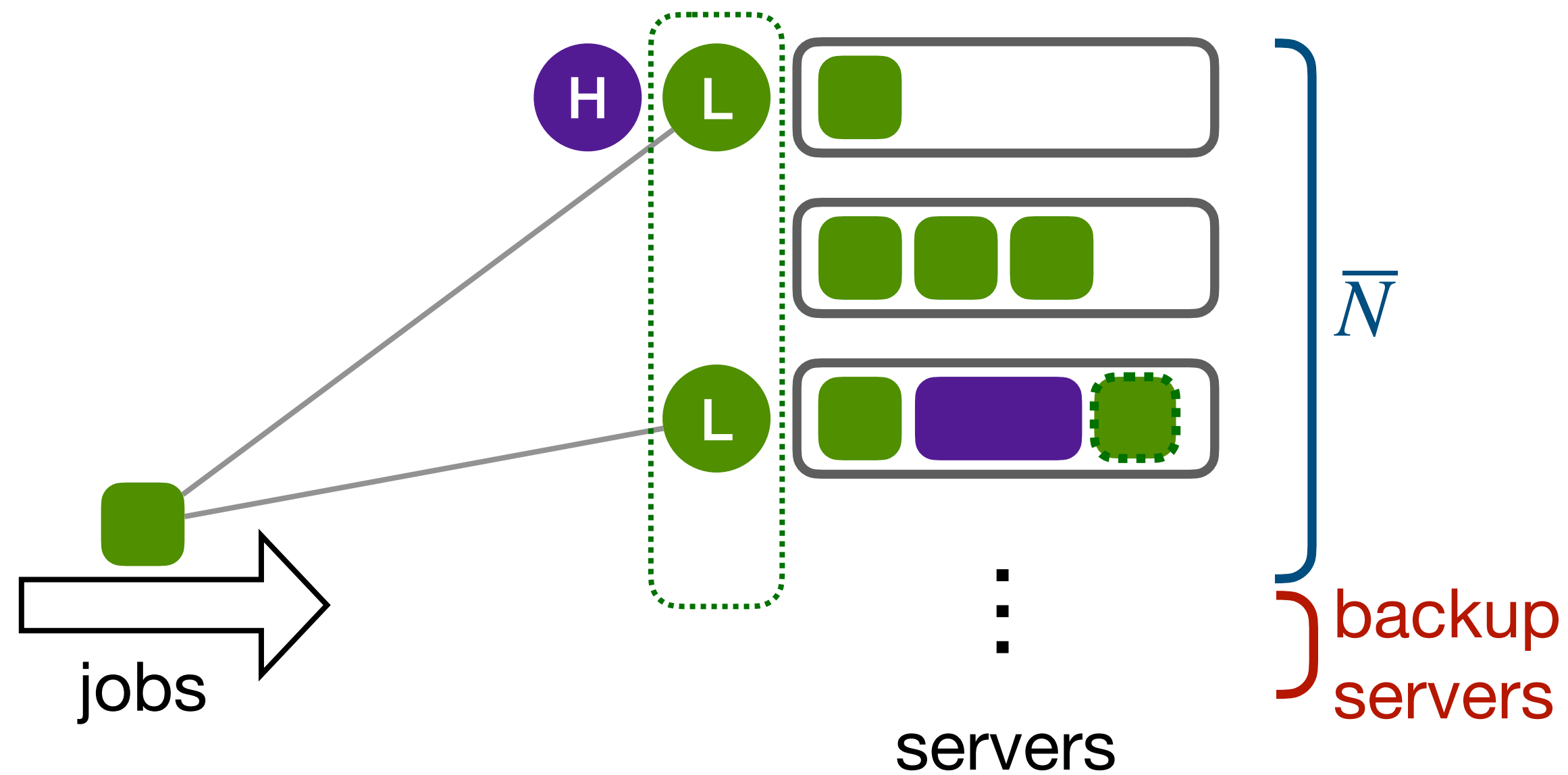
Key proof idea 2



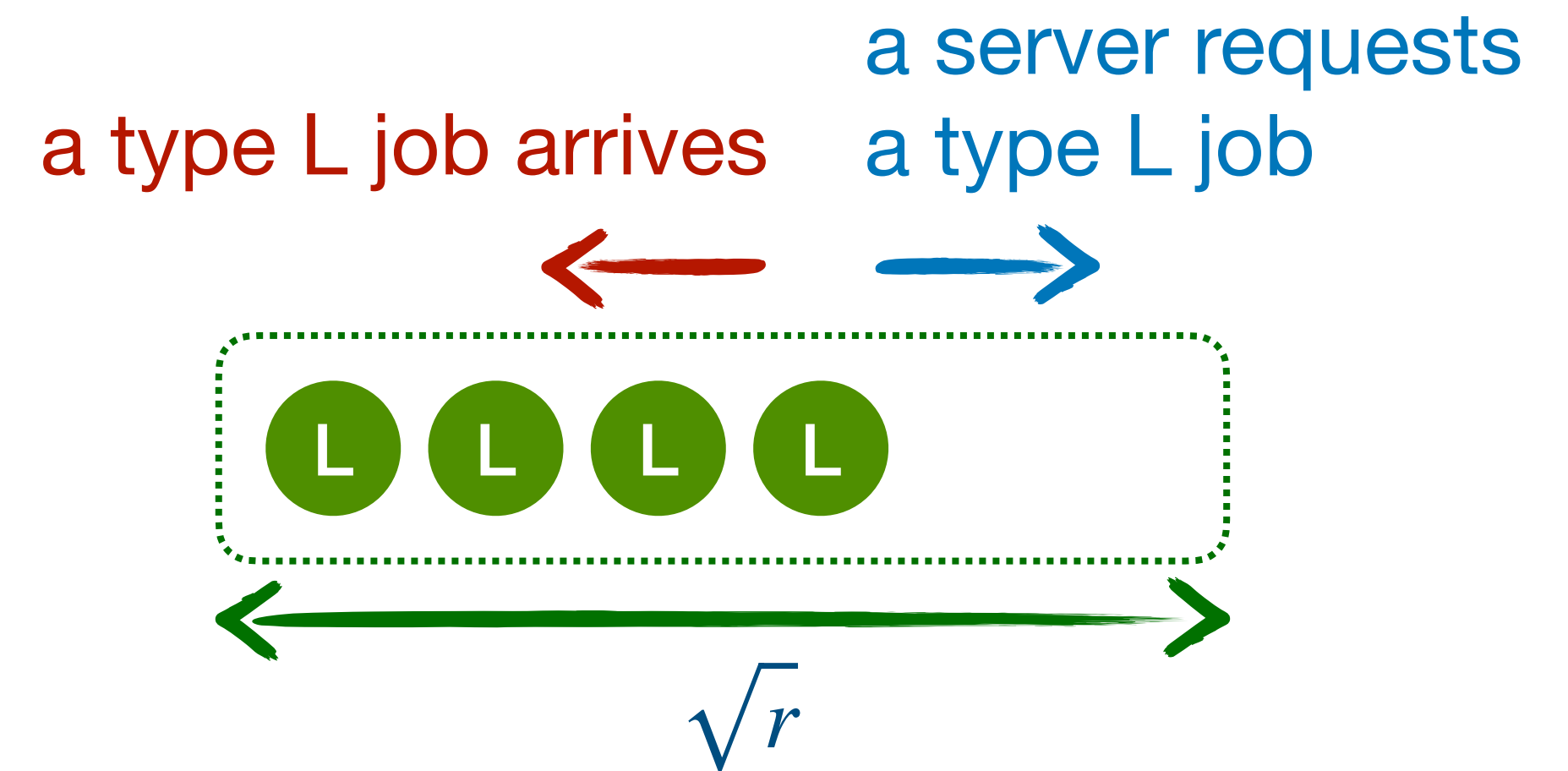
Will show that # virtual jobs = $O(\sqrt{r})$,
and # backup servers = $O(\sqrt{r})$



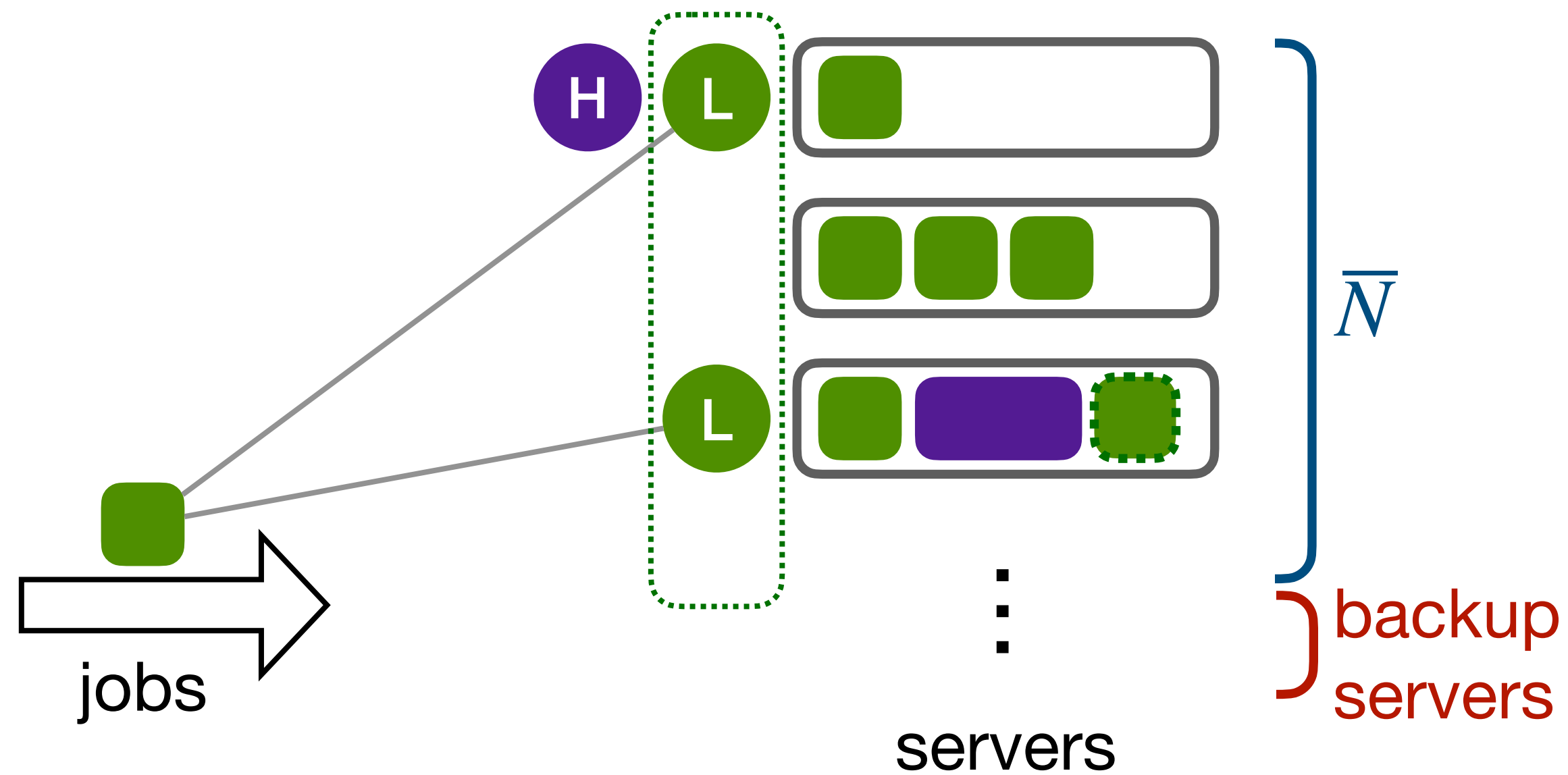
Key proof idea 2



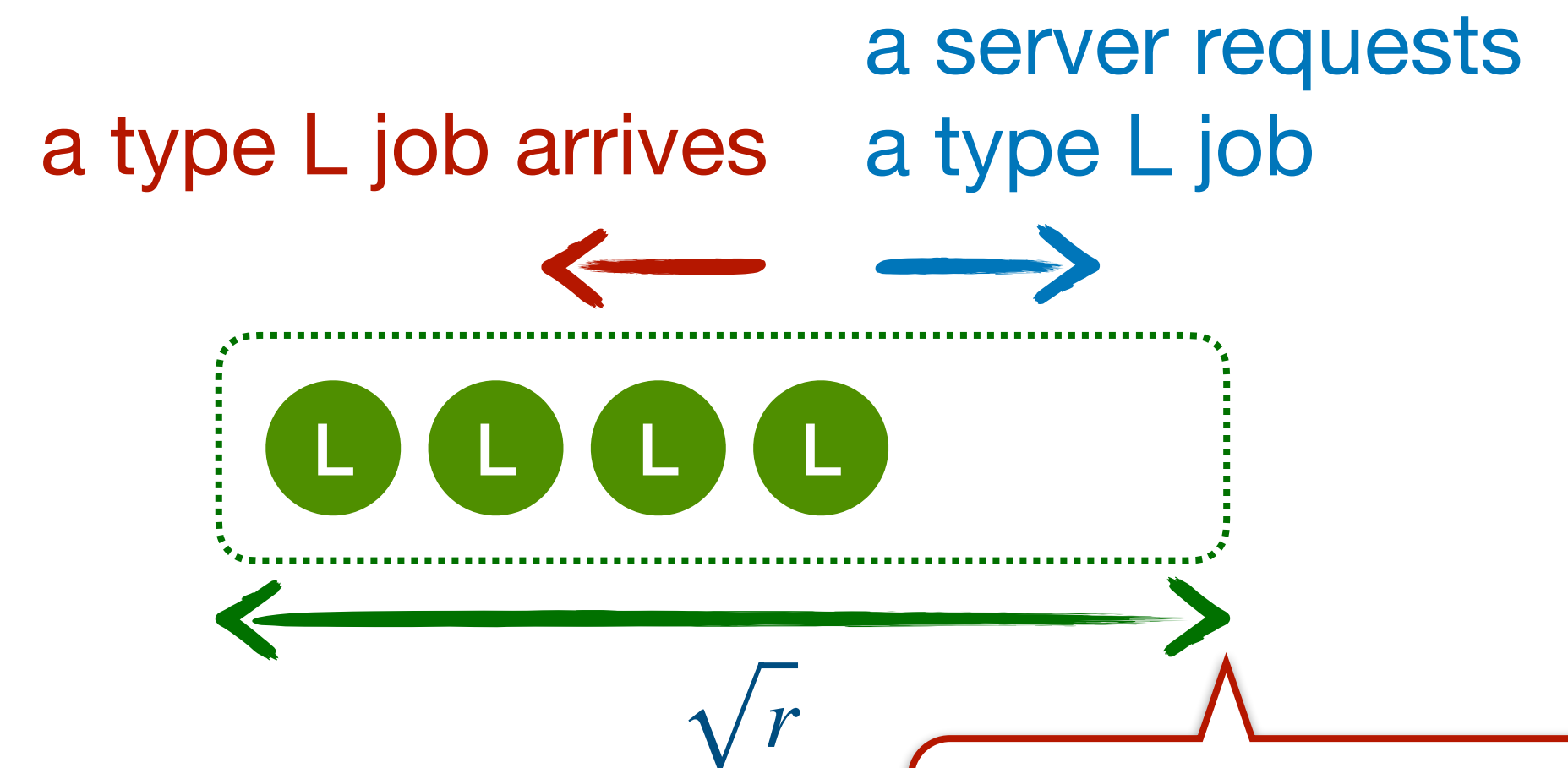
Will show that # virtual jobs = $O(\sqrt{r})$,
and # backup servers = $O(\sqrt{r})$



Key proof idea 2

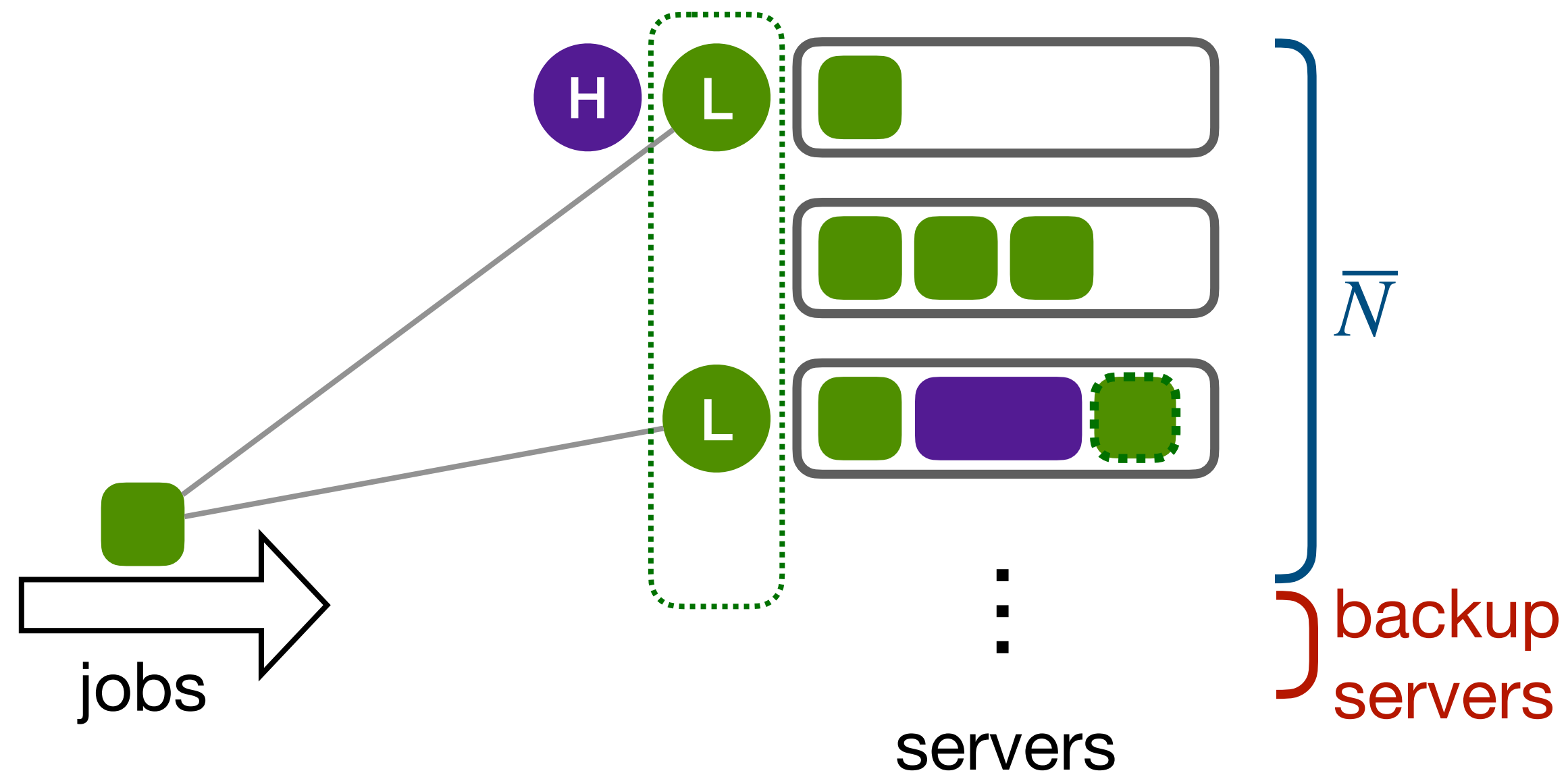


Will show that # virtual jobs = $O(\sqrt{r})$,
and # backup servers = $O(\sqrt{r})$

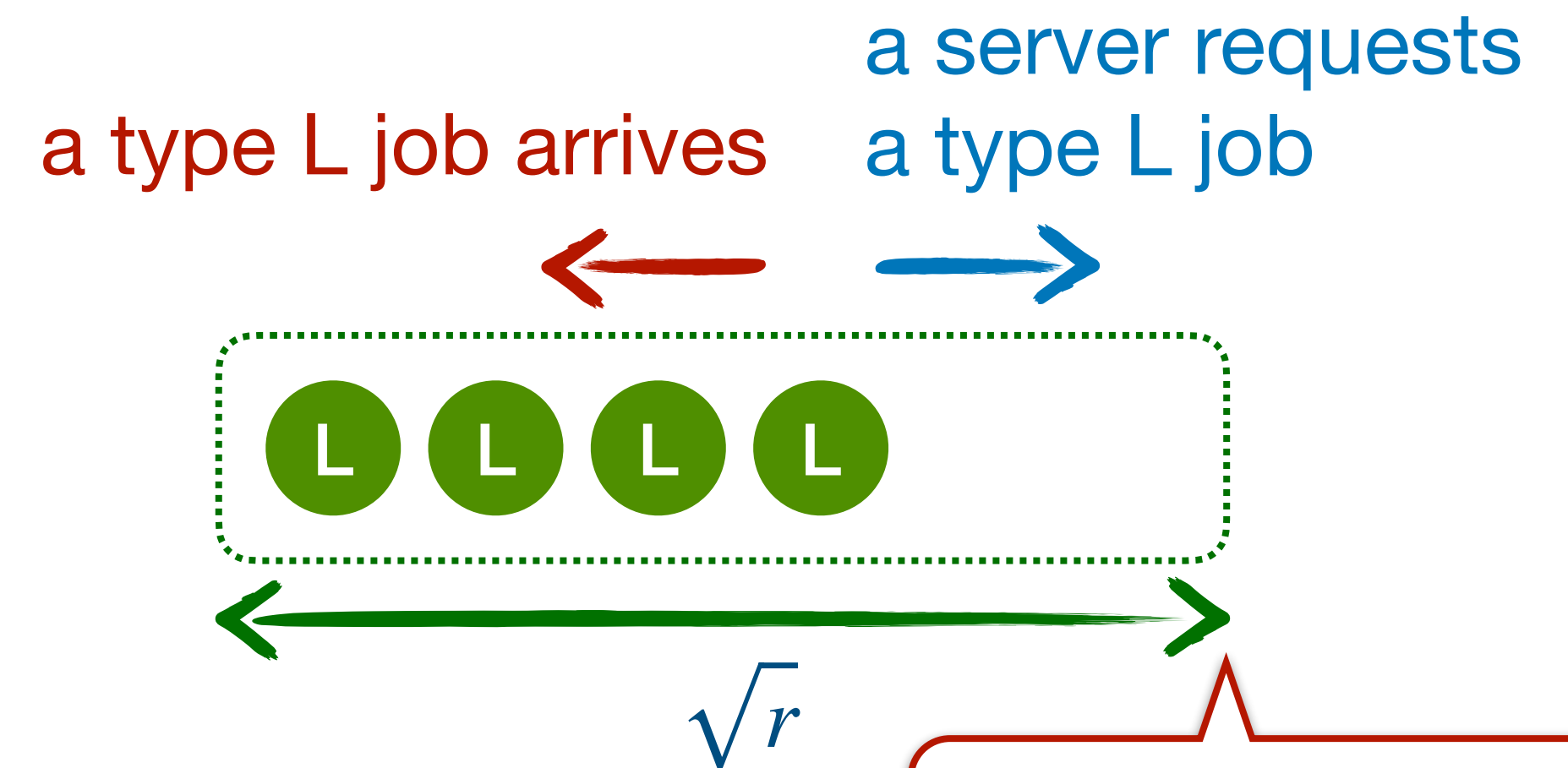


What happens when
tokens hits \sqrt{r} ?

Key proof idea 2



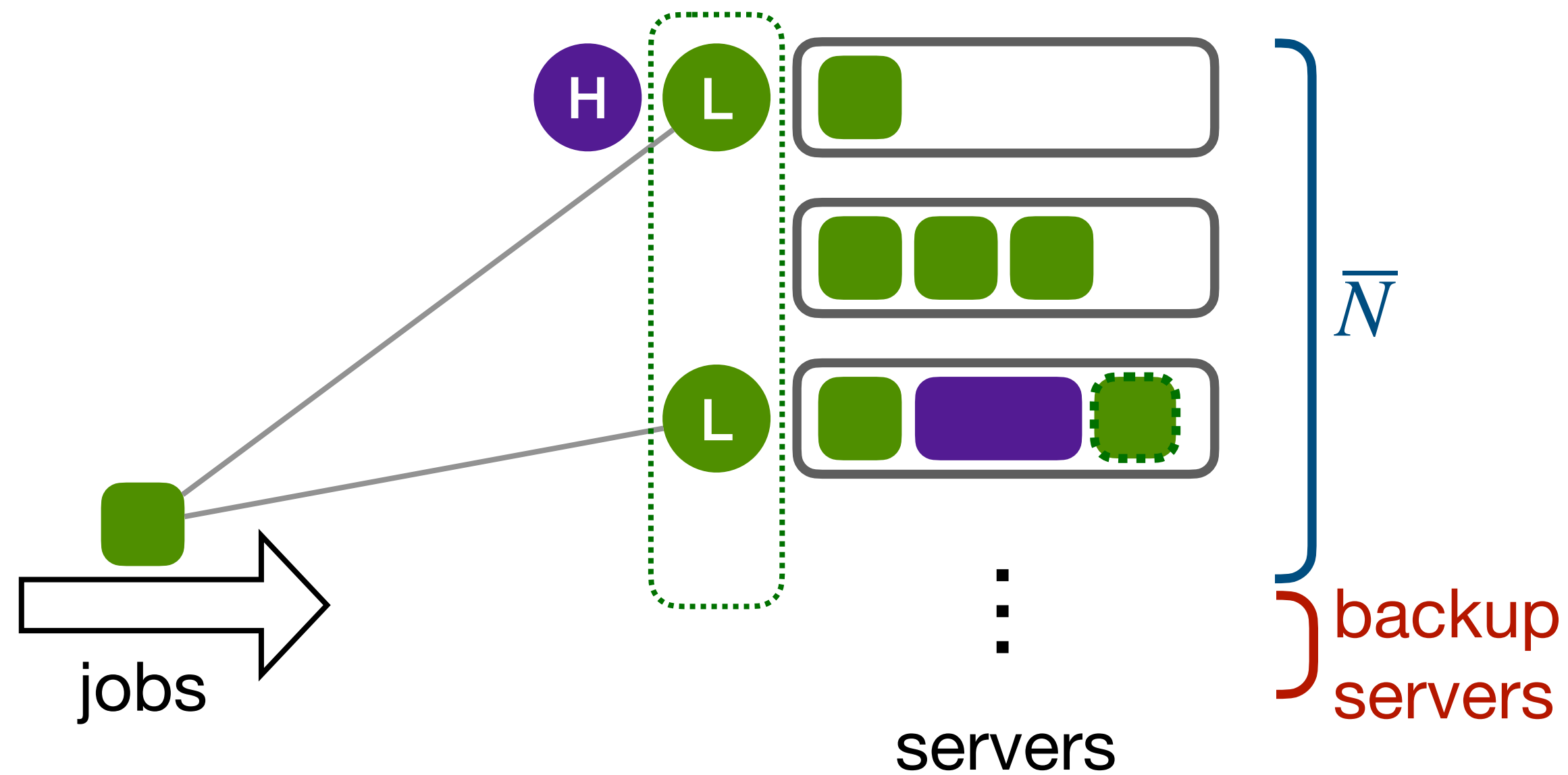
Will show that # virtual jobs = $O(\sqrt{r})$,
and # backup servers = $O(\sqrt{r})$



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

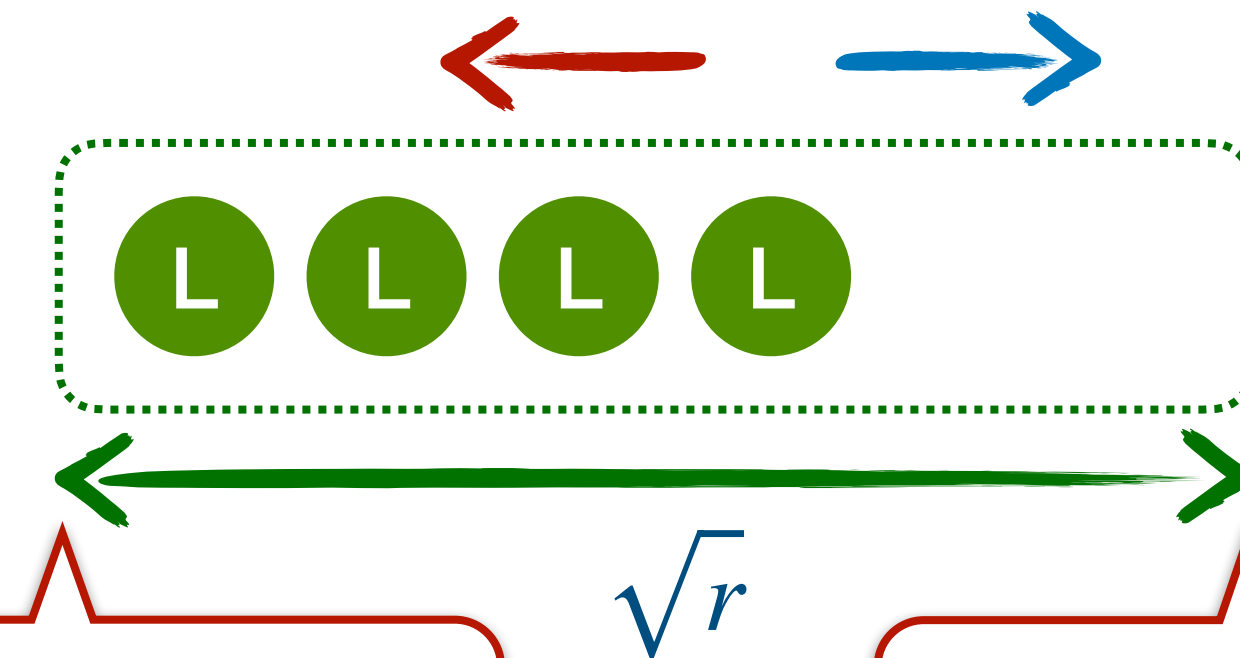
Generate a virtual job

Key proof idea 2



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

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

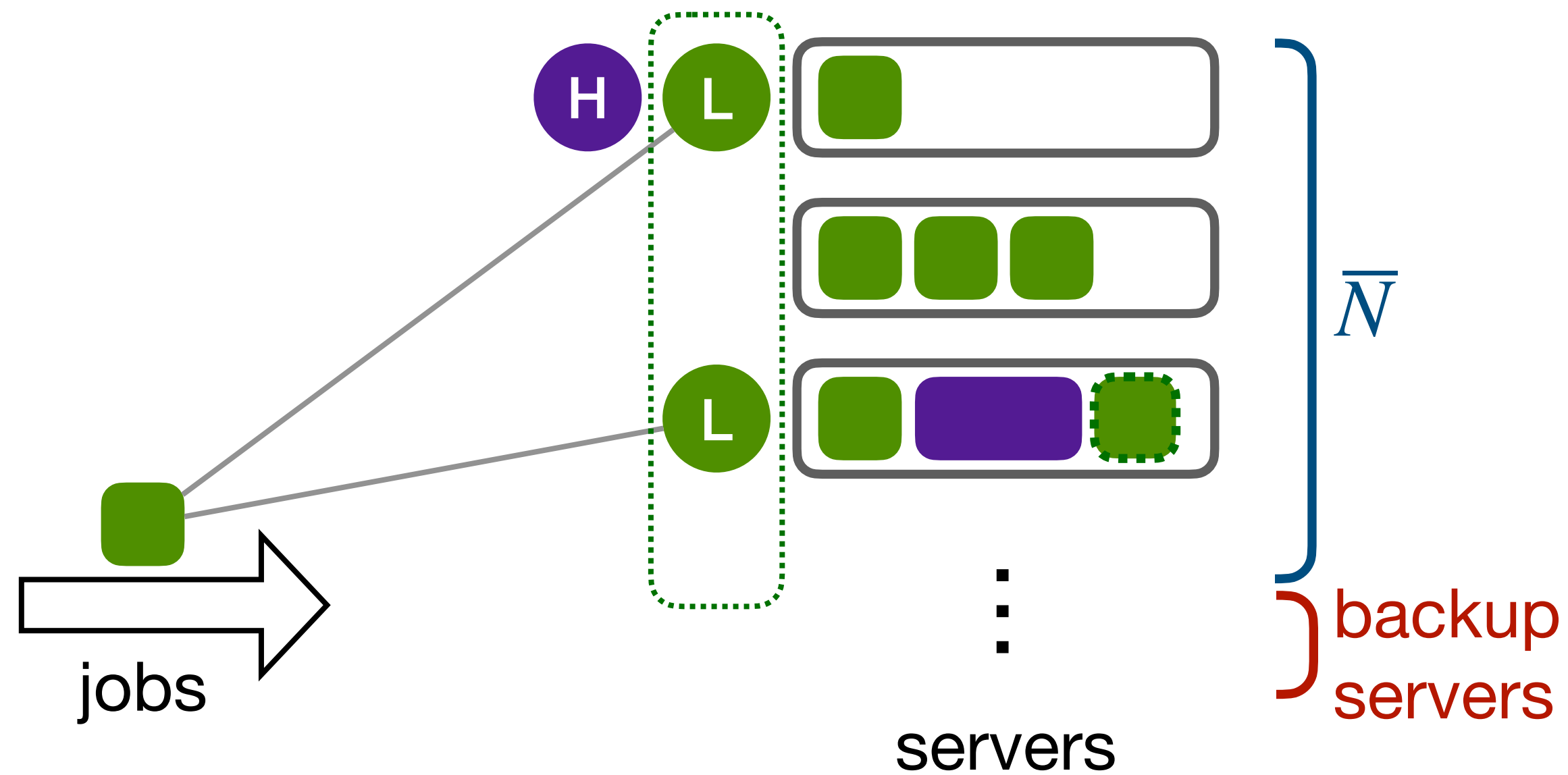


What happens when # tokens hits 0?

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

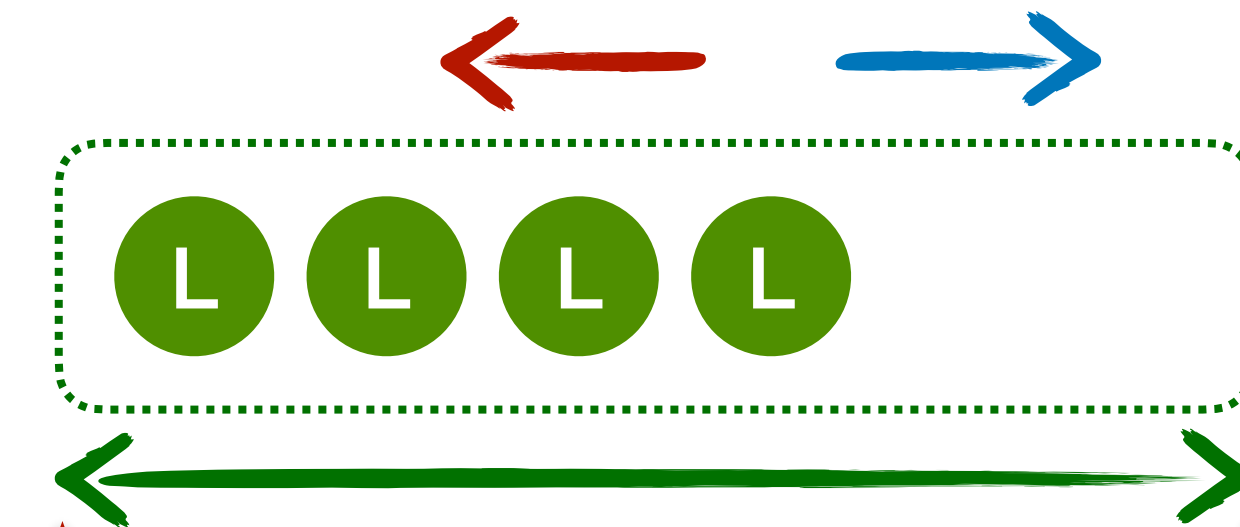
Generate a virtual job

Key proof idea 2



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

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



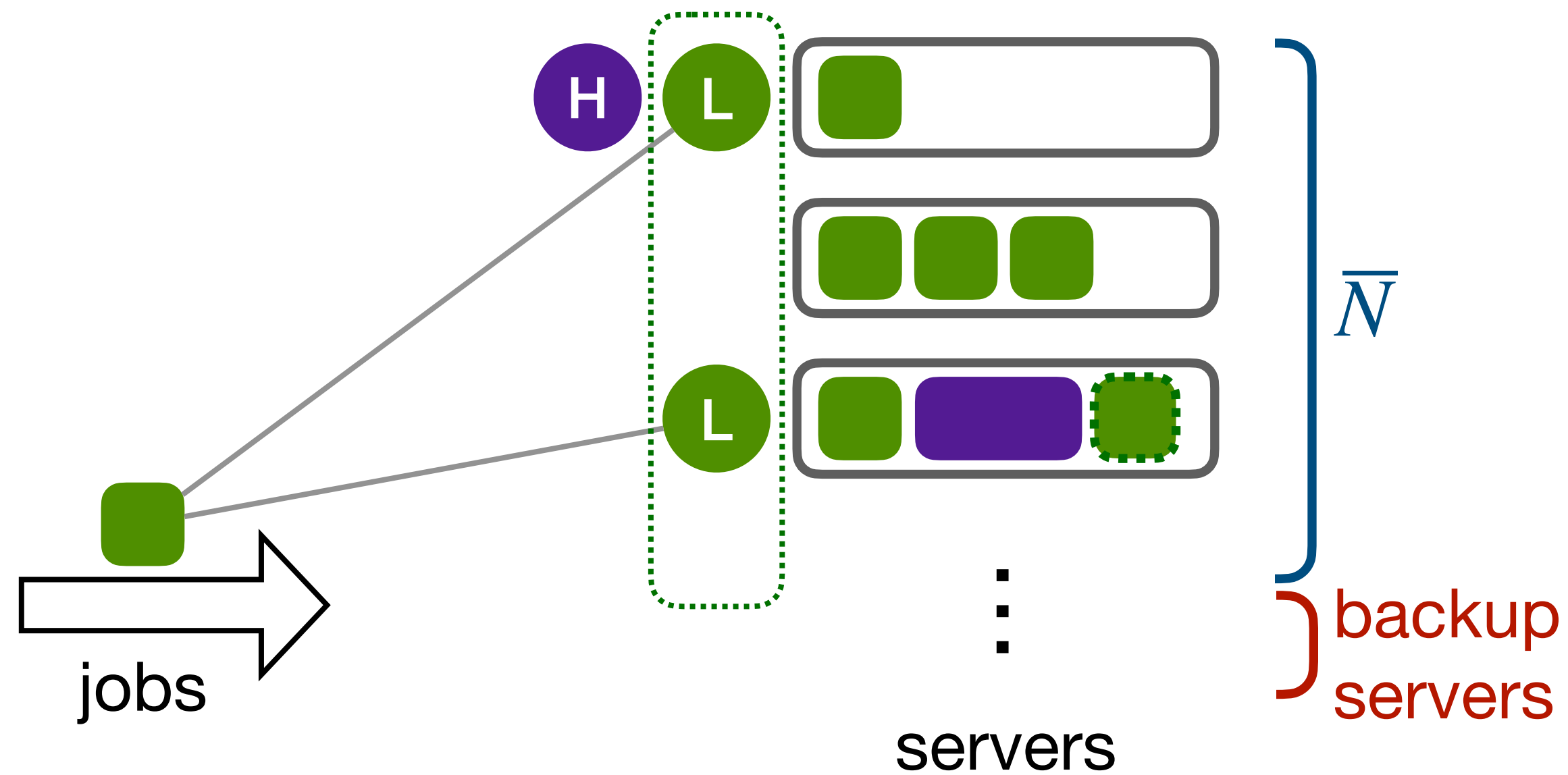
What happens when # tokens hits 0?

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

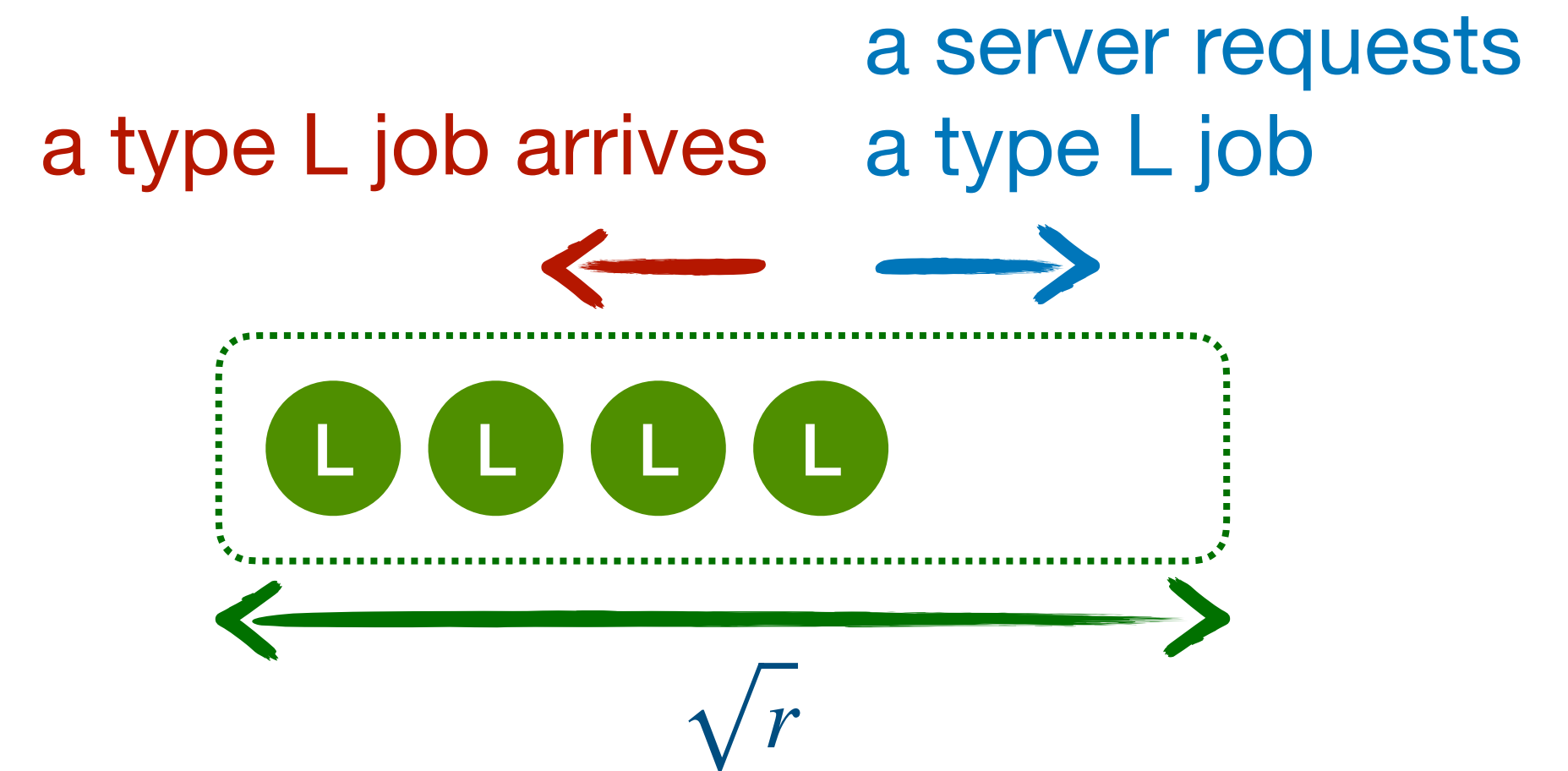
Generate a job to backup servers

Generate a virtual job

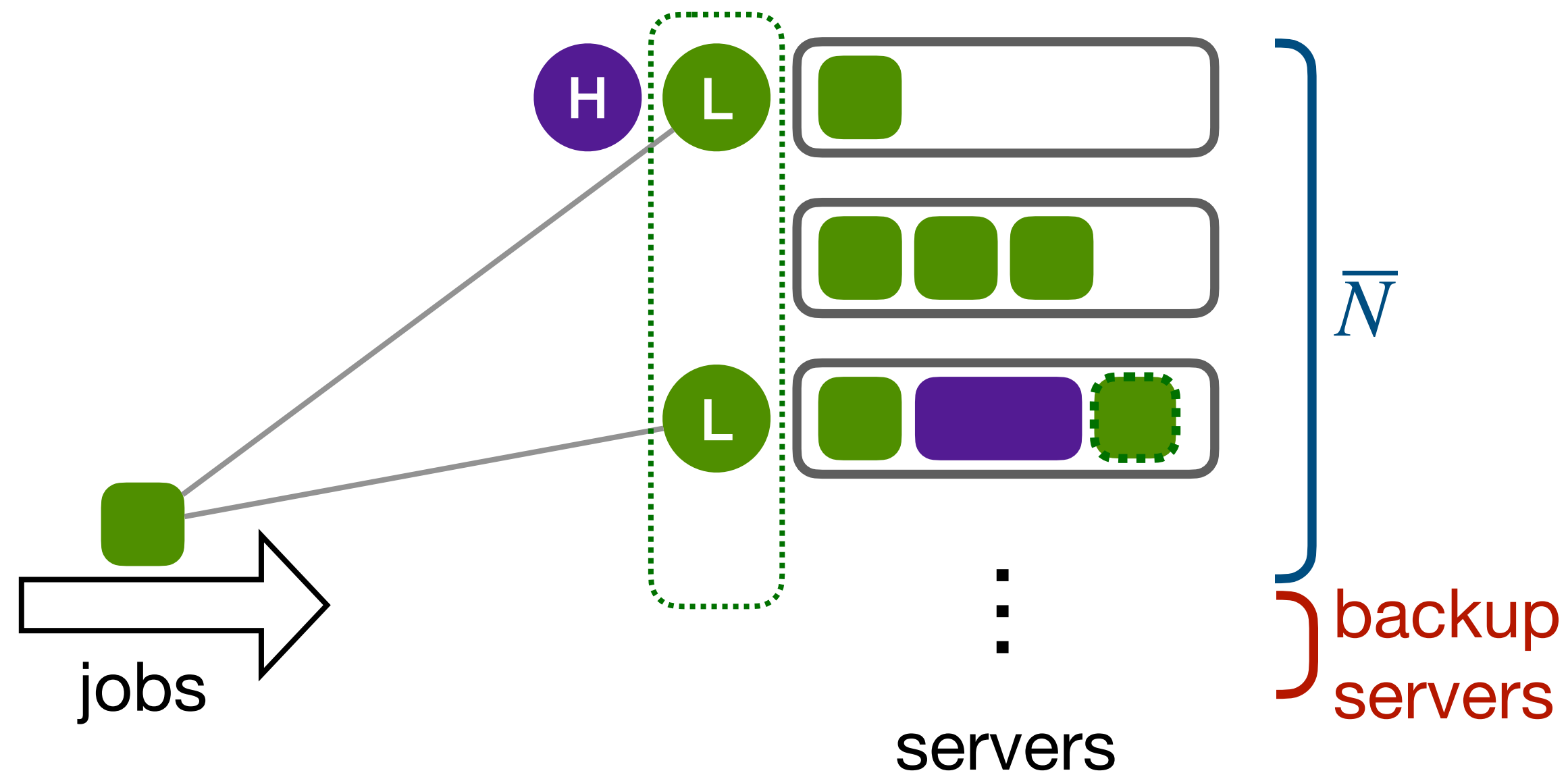
Key proof idea 2



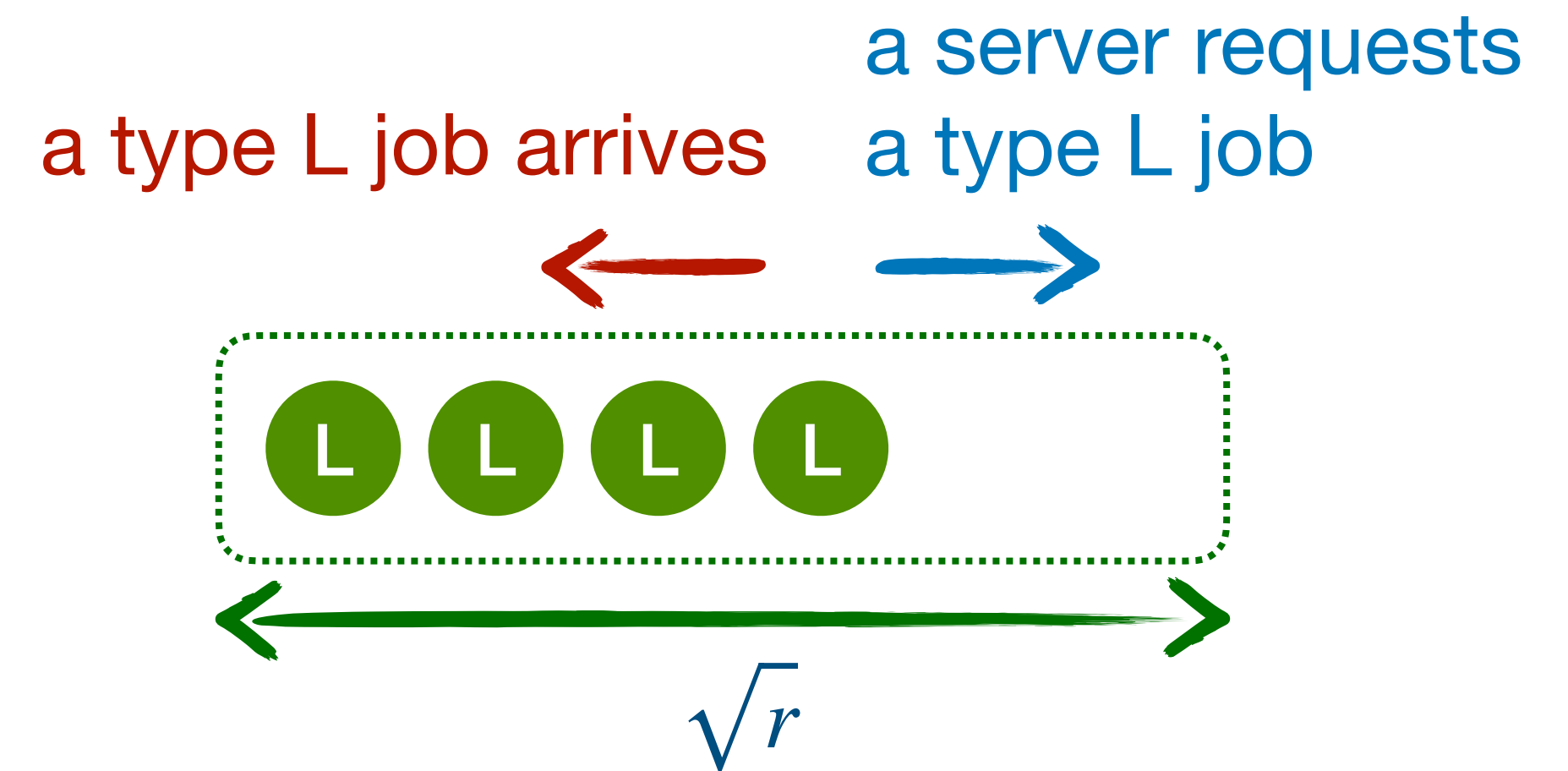
Will show that # virtual jobs = $O(\sqrt{r})$,
and # backup servers = $O(\sqrt{r})$



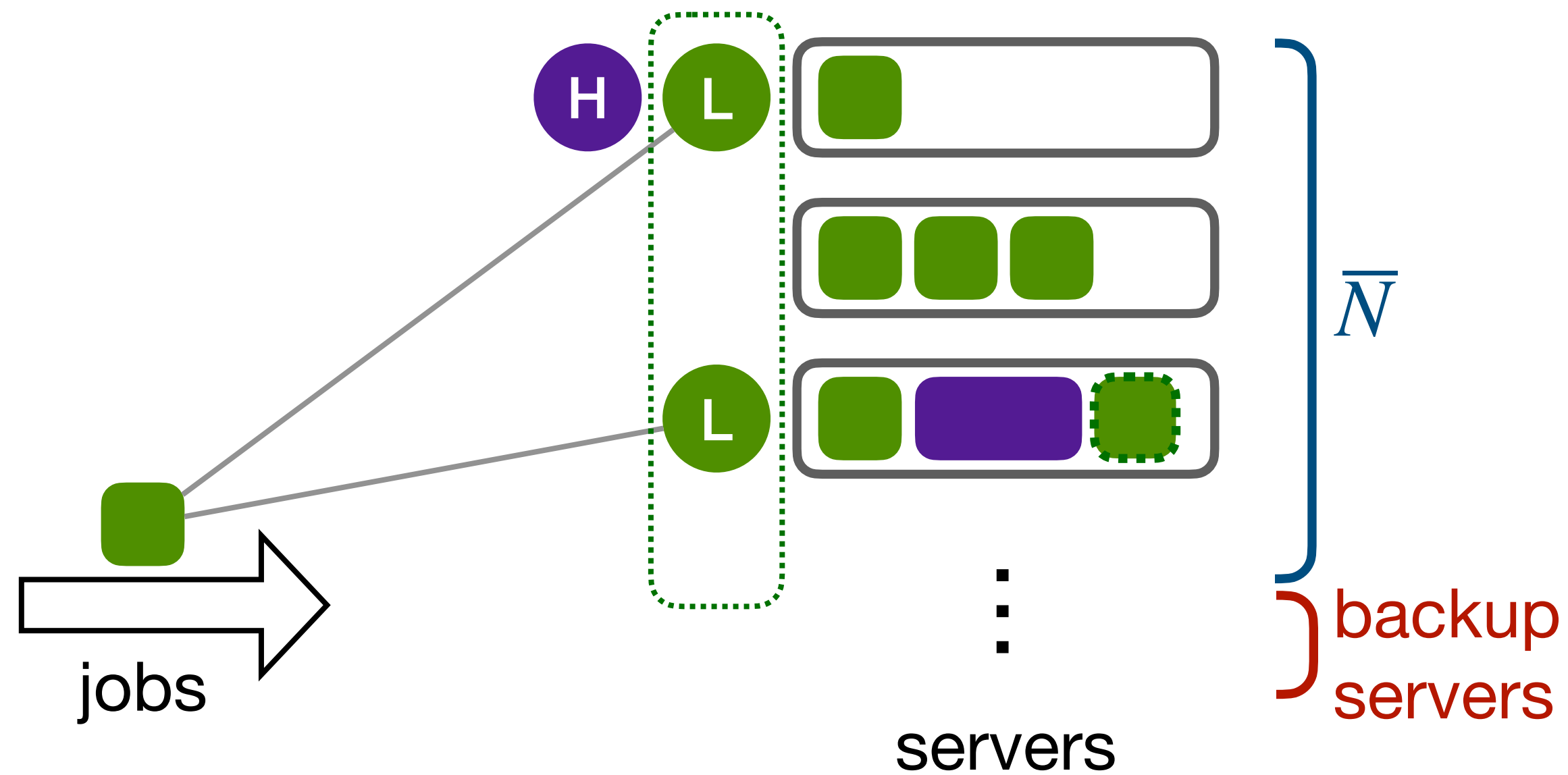
Key proof idea 2



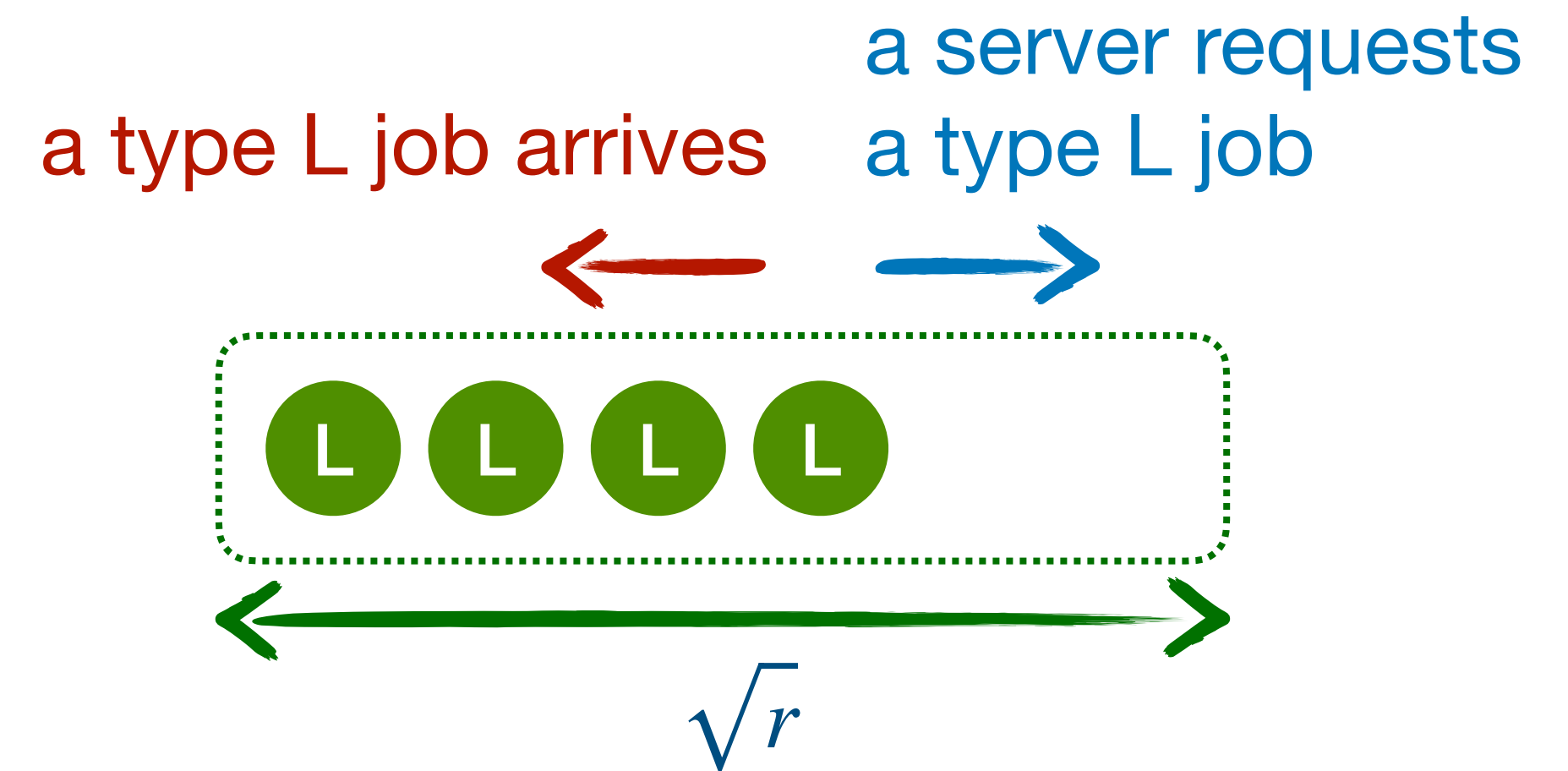
Will show that # virtual jobs = $O(\sqrt{r})$,
and # backup servers = $O(\sqrt{r})$



Key proof idea 2

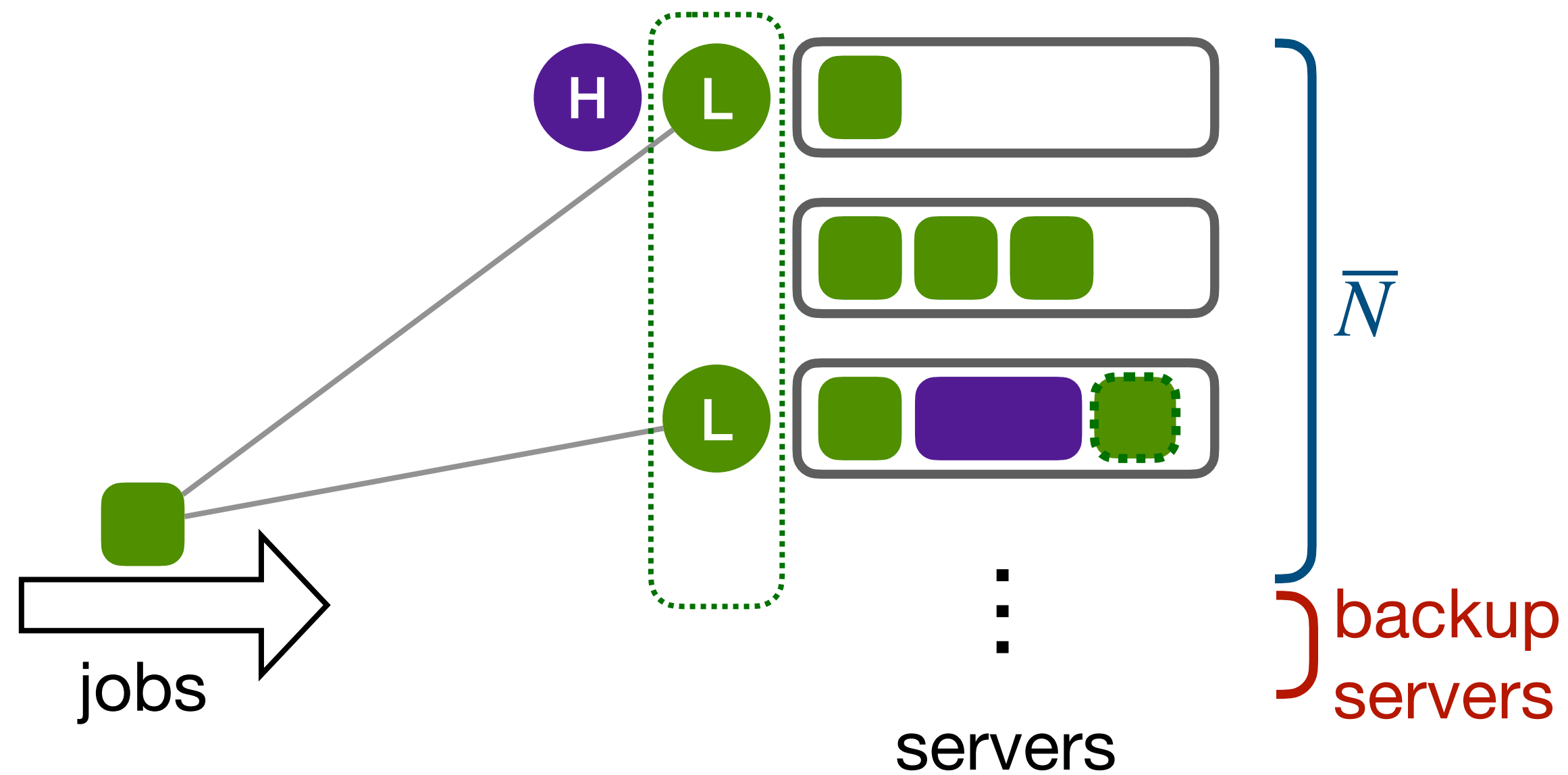


Will show that # virtual jobs = $O(\sqrt{r})$,
and # backup servers = $O(\sqrt{r})$

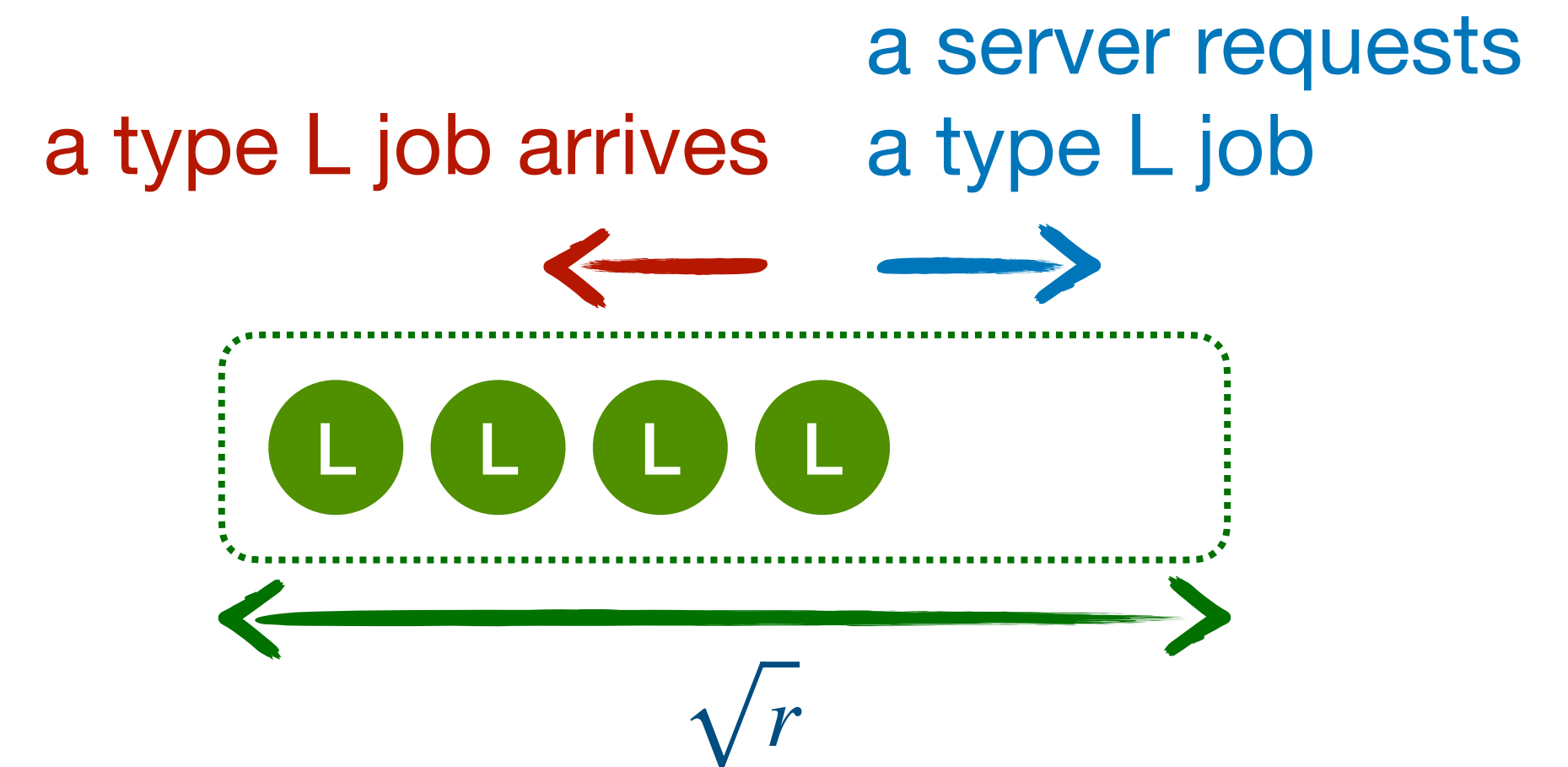


- An almost balanced random walk

Key proof idea 2

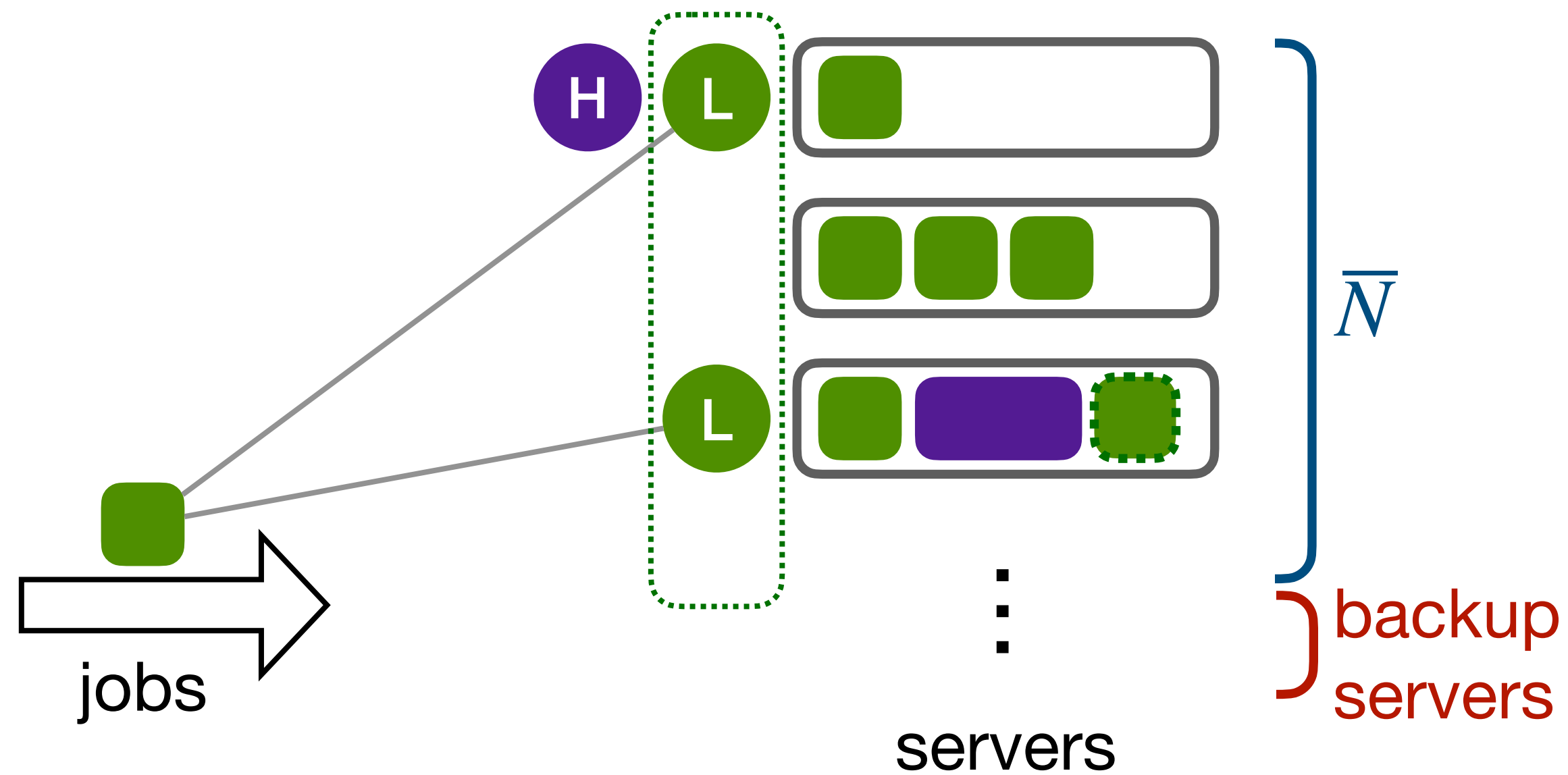


Will show that # virtual jobs = $O(\sqrt{r})$,
and # backup servers = $O(\sqrt{r})$

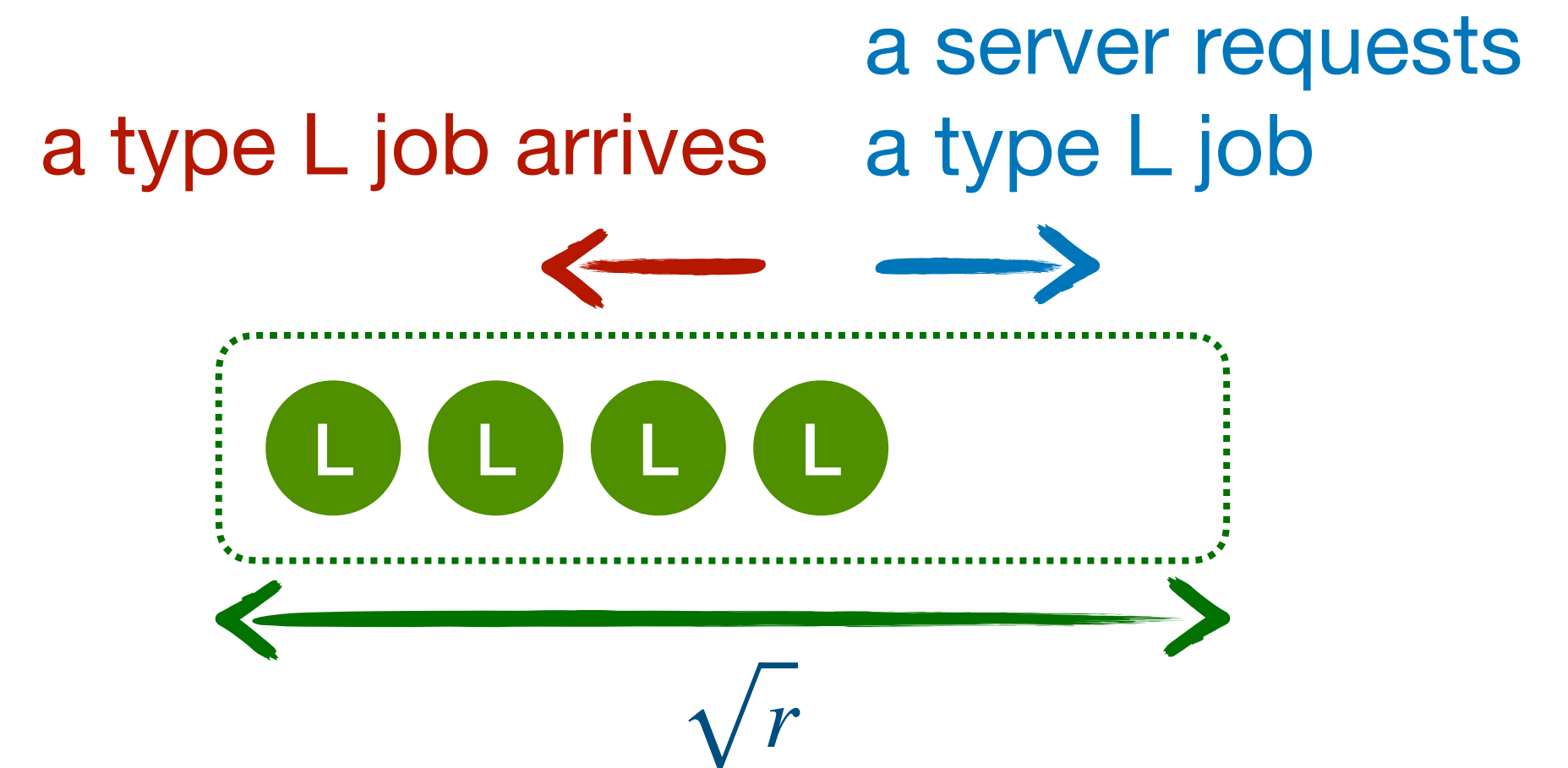


- An almost balanced random walk
- Stationary distribution \approx uniform on $\{0, 1, \dots, \sqrt{r}\}$

Key proof idea 2

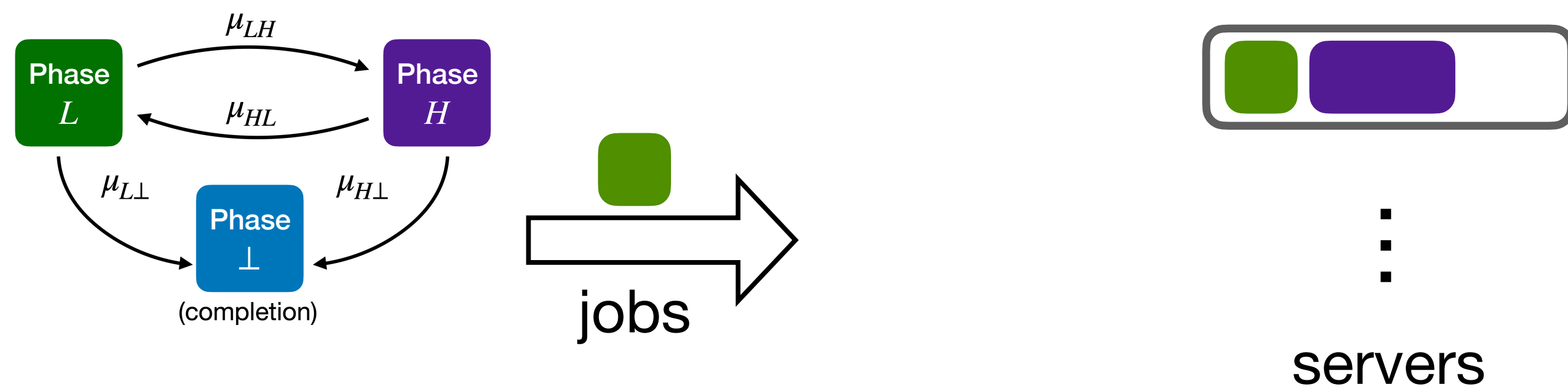


Will show that # virtual jobs = $O(\sqrt{r})$,
and # backup servers = $O(\sqrt{r})$



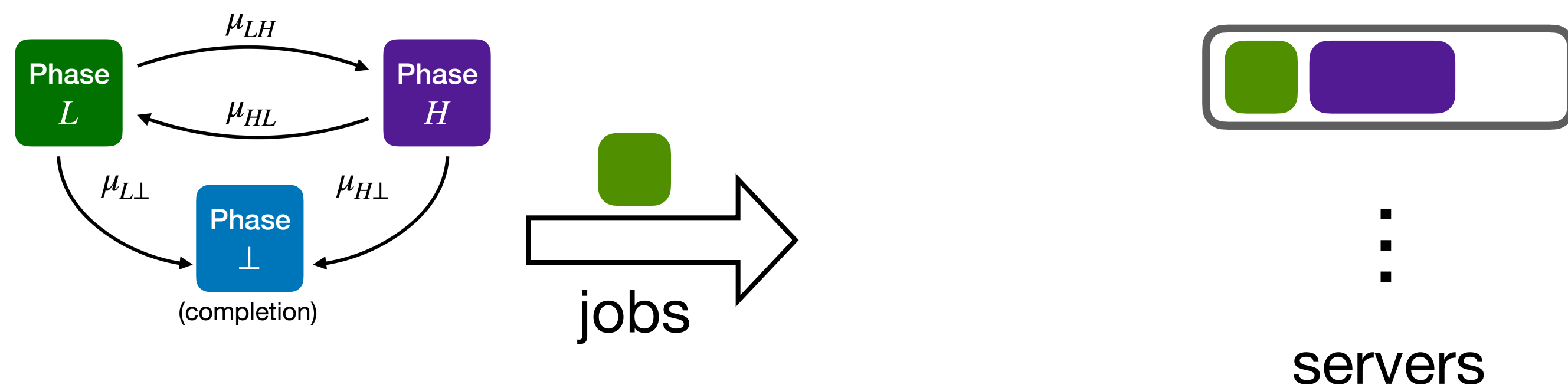
- An almost balanced random walk
- Stationary distribution \approx uniform on $\{0, 1, \dots, \sqrt{r}\}$
- Rate of generating virtual jobs
 \approx rate of sending jobs to backup servers
 \approx arrival rate / $\sqrt{r} = O(\sqrt{r})$

Summary



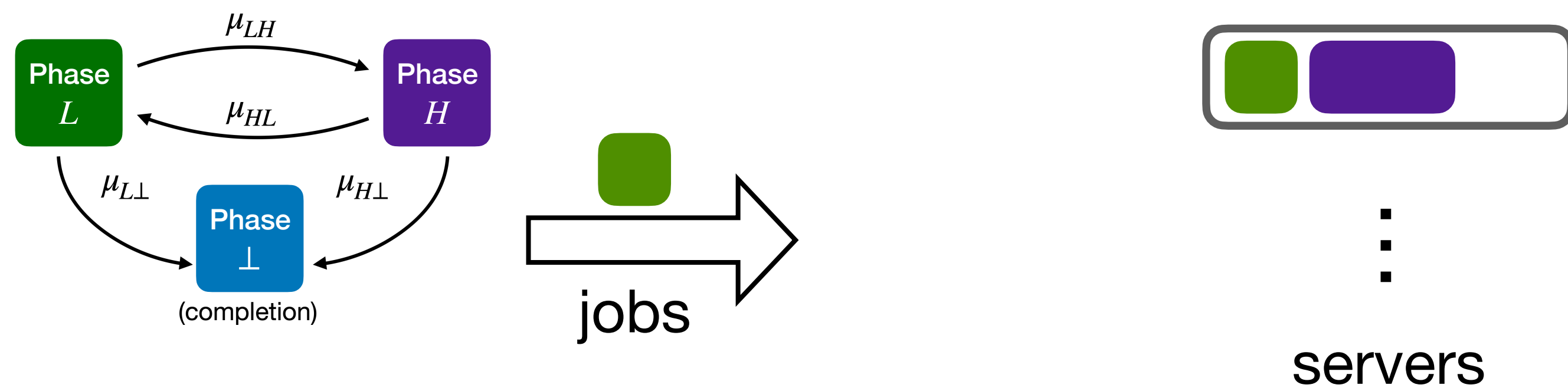
Summary

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



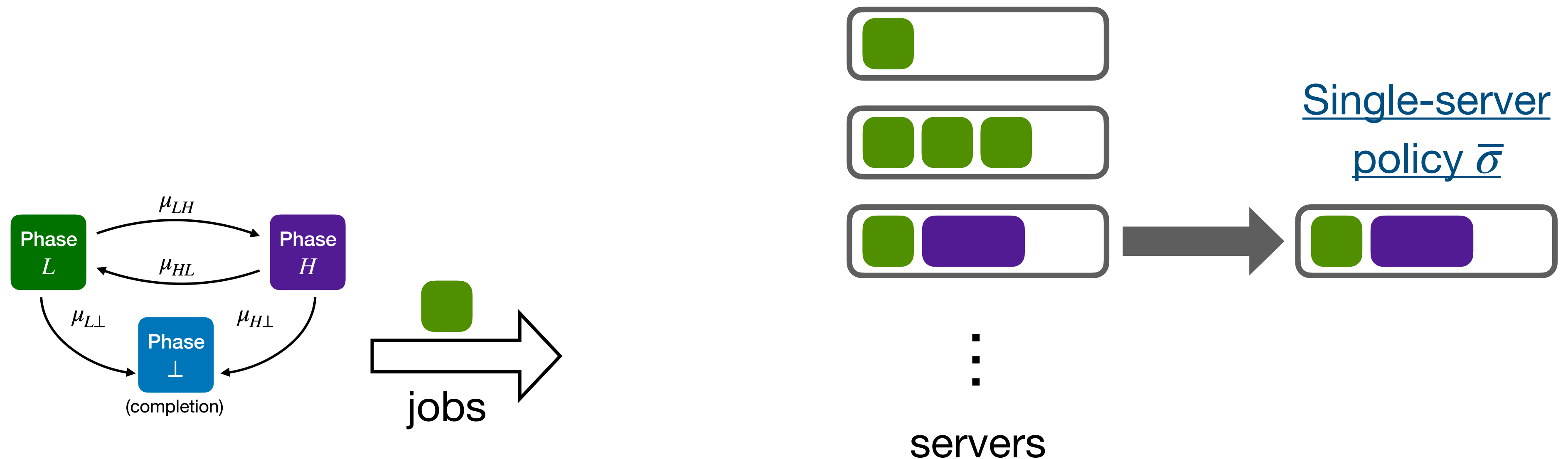
Summary

- 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-REQUESTING-SERVER (JRS), that converts a single-server policy to a policy in the original system

