

# 2DFQ: Two-Dimensional Fair Queuing for Multi-Tenant Cloud Services

Jonathan Mace  
*Brown University*

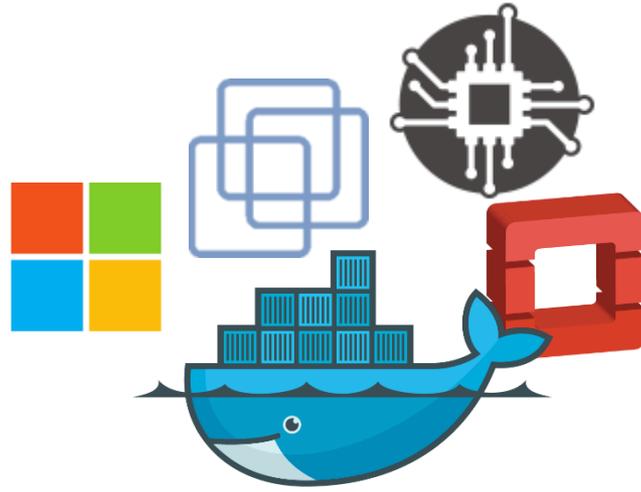
Peter Bodik  
*Microsoft*

Madanlal Musuvathi  
*Microsoft*

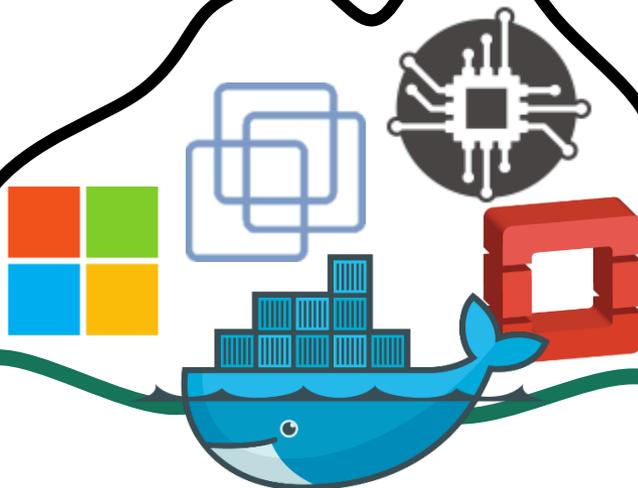
Rodrigo Fonseca  
*Brown University*

Krishnan Varadarajan  
*Microsoft*

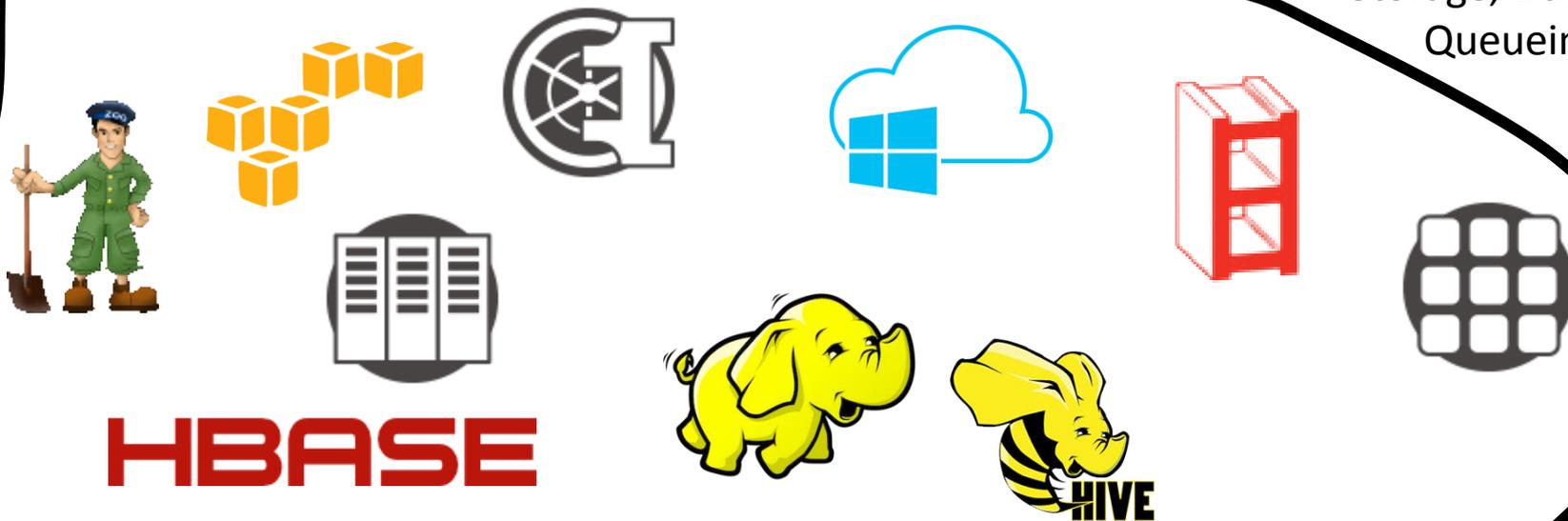
# Containers / VMs



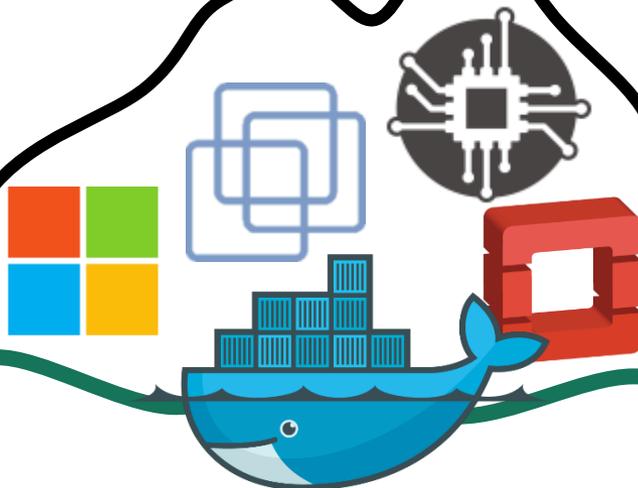
Containers / VMs



Shared Systems:  
Storage, Database,  
Queueing, etc.



Containers / VMs

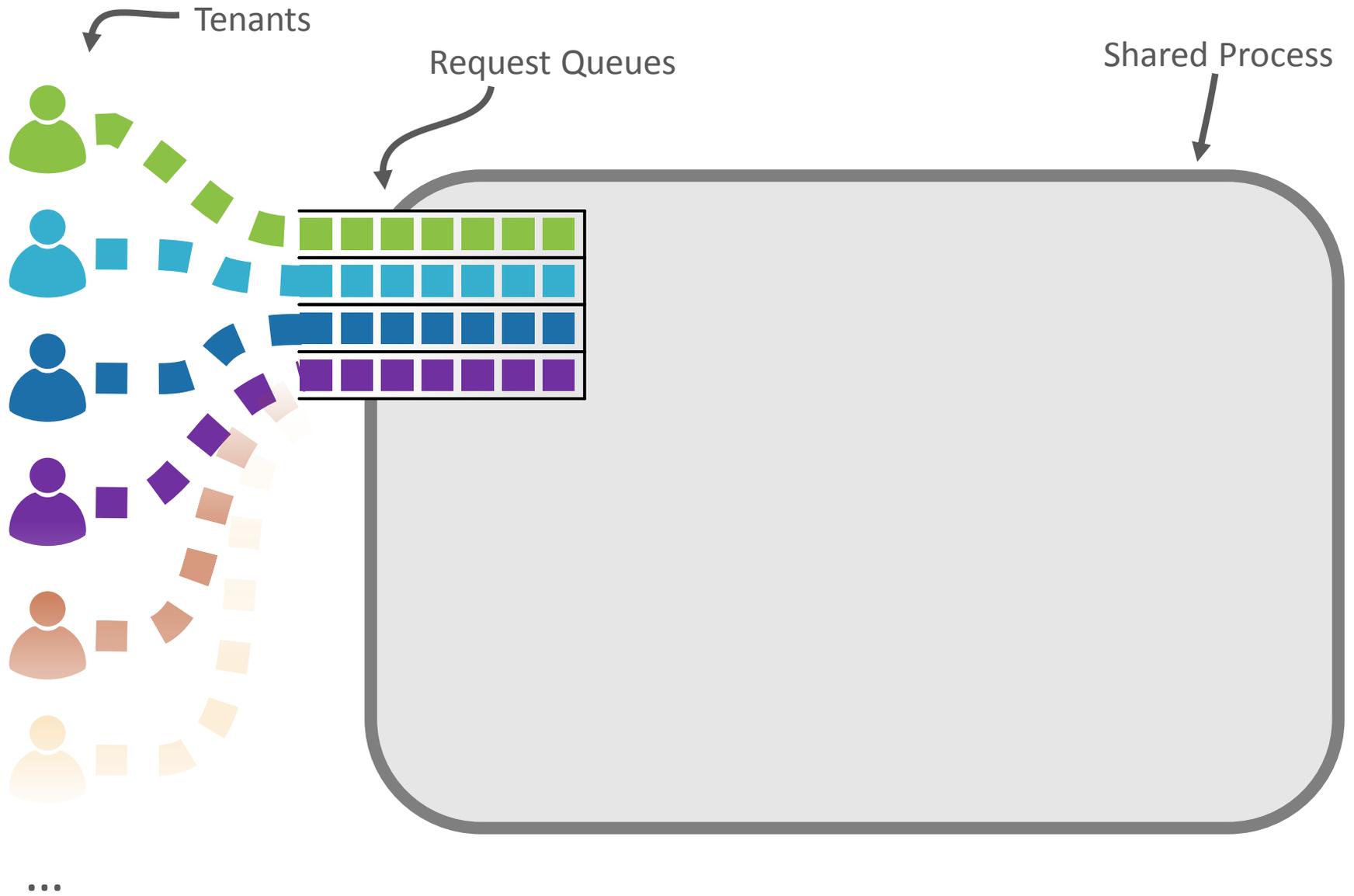


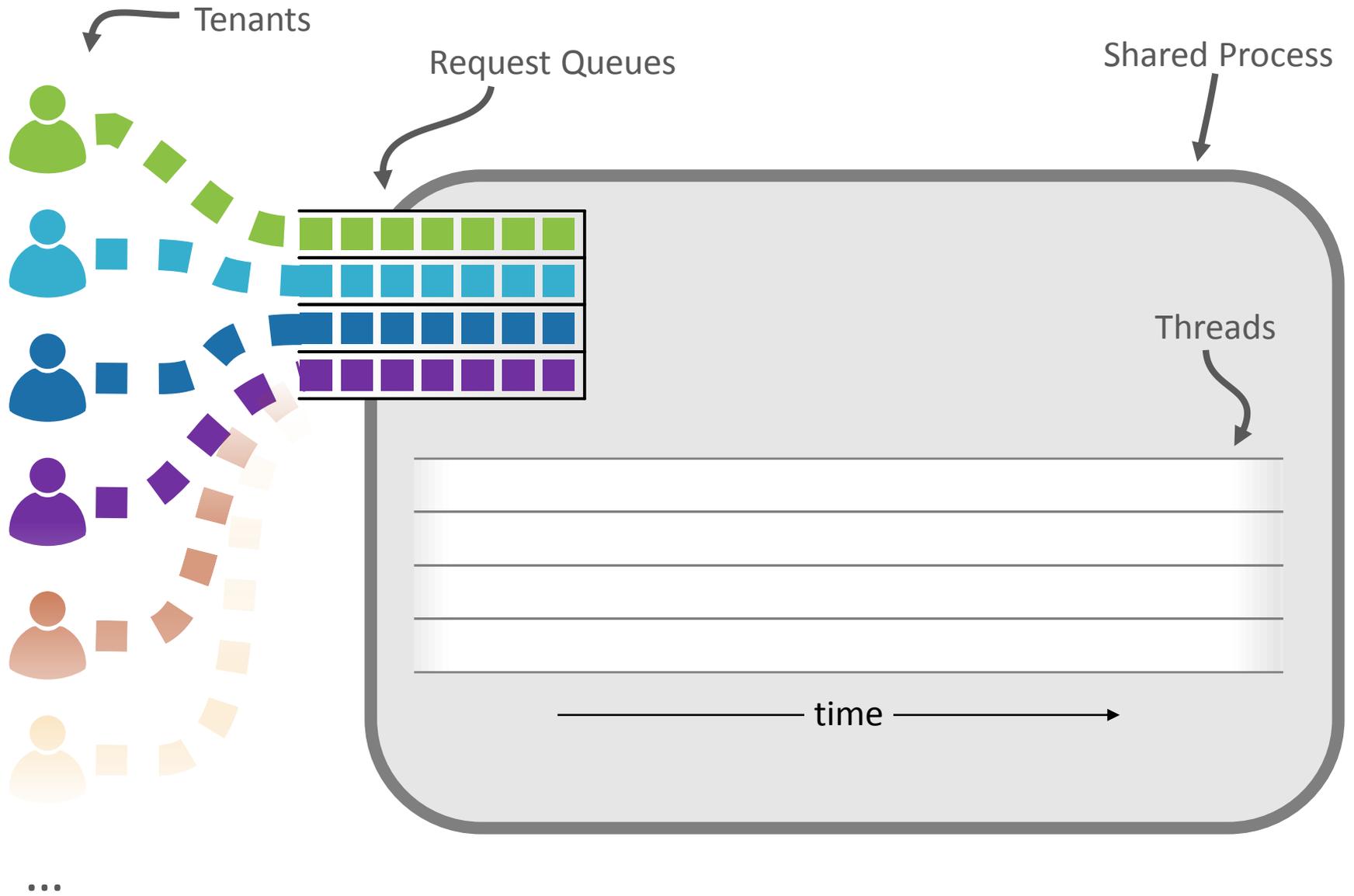
Shared Systems:  
Storage, Database,  
Queueing, etc.

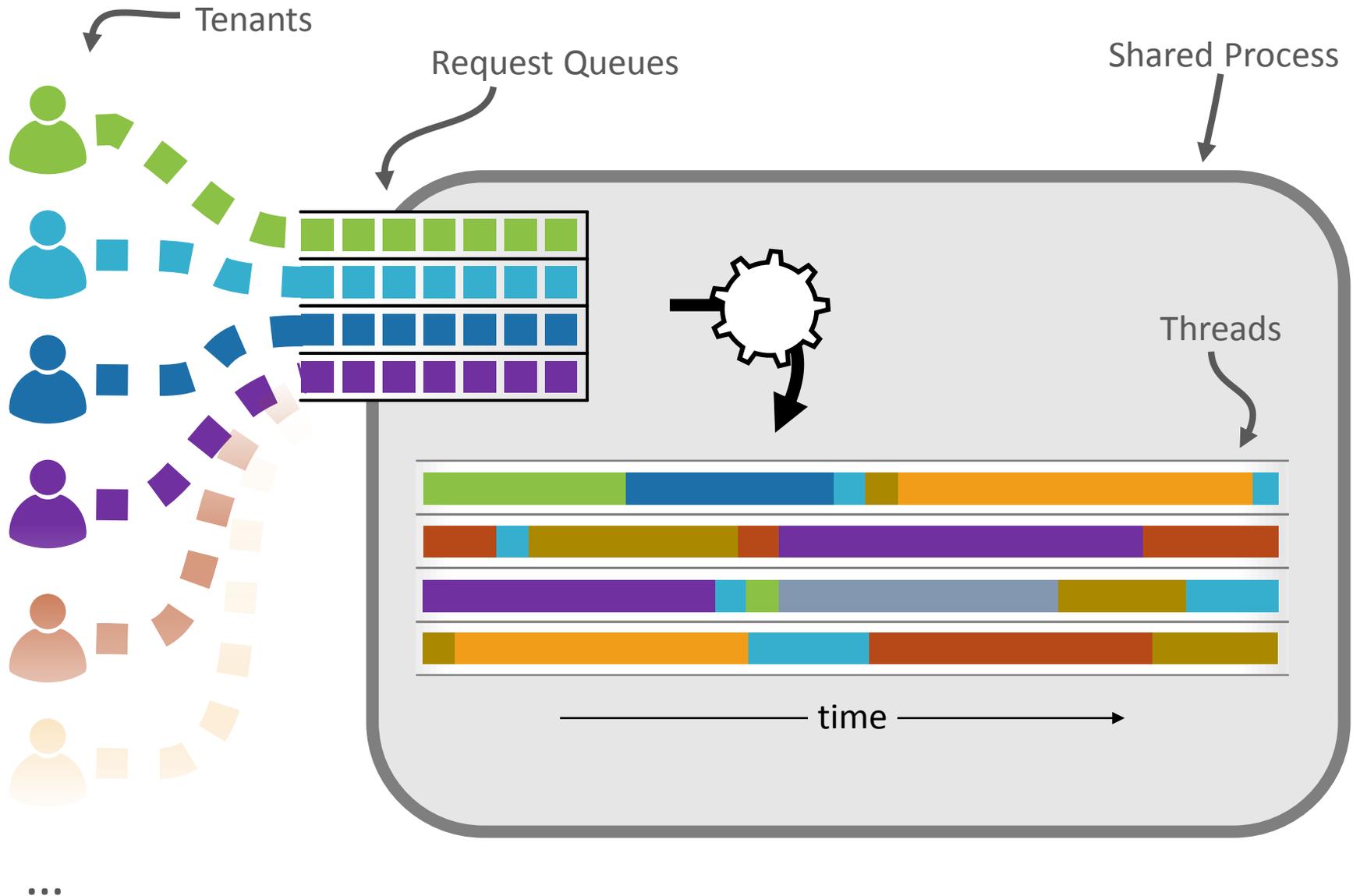
**OS**  
**Hypervisor**

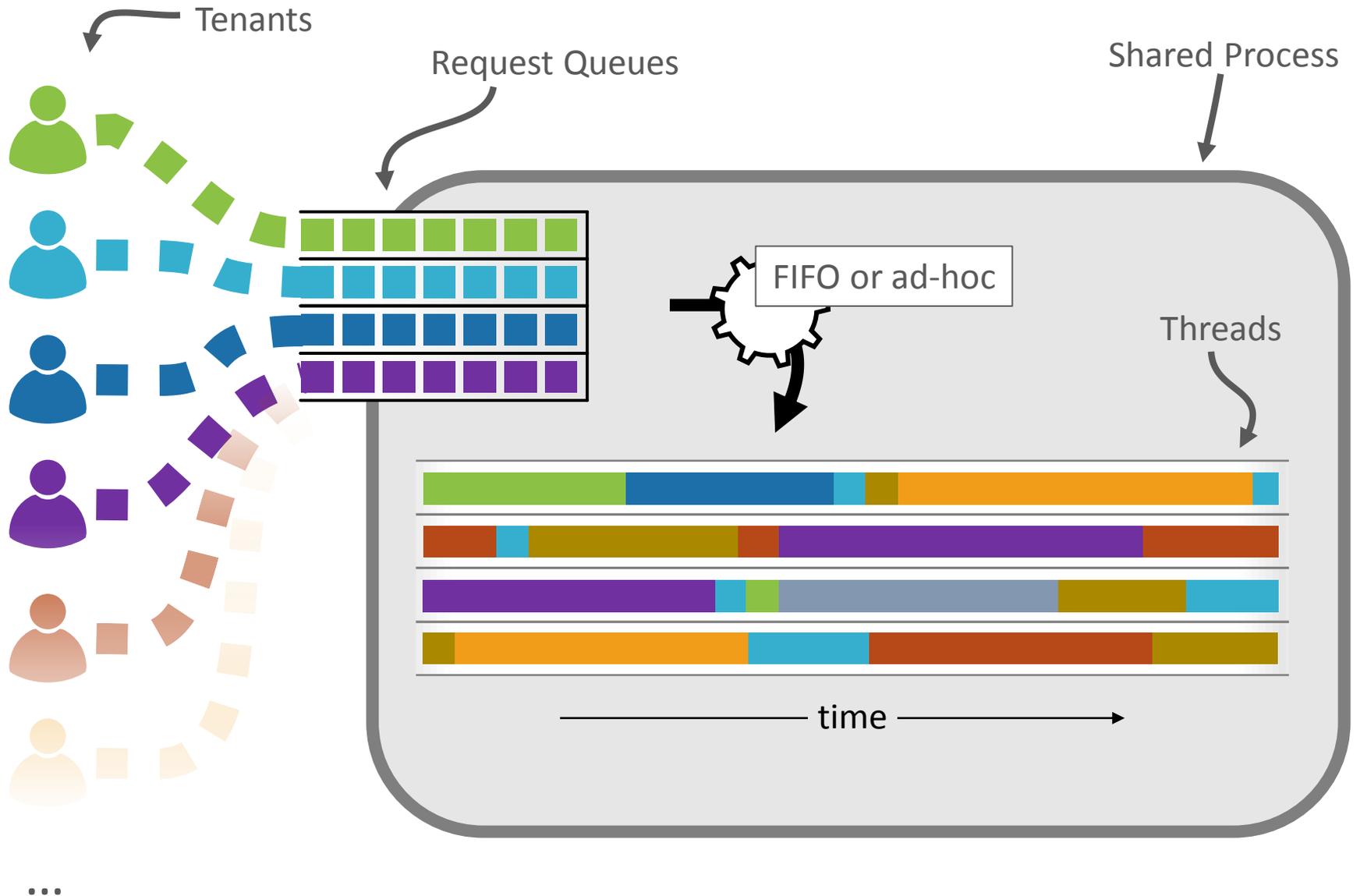
Shared Process

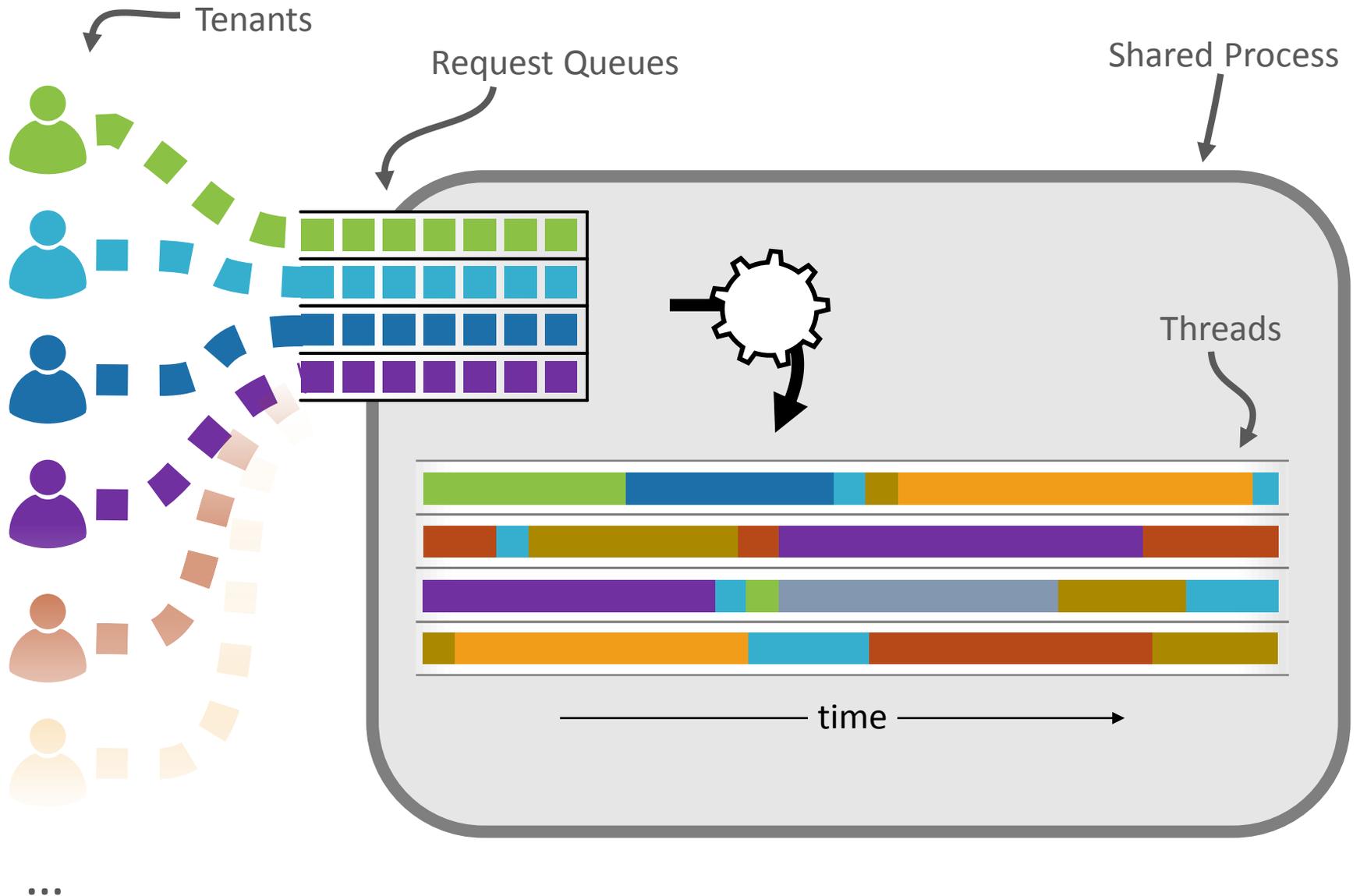


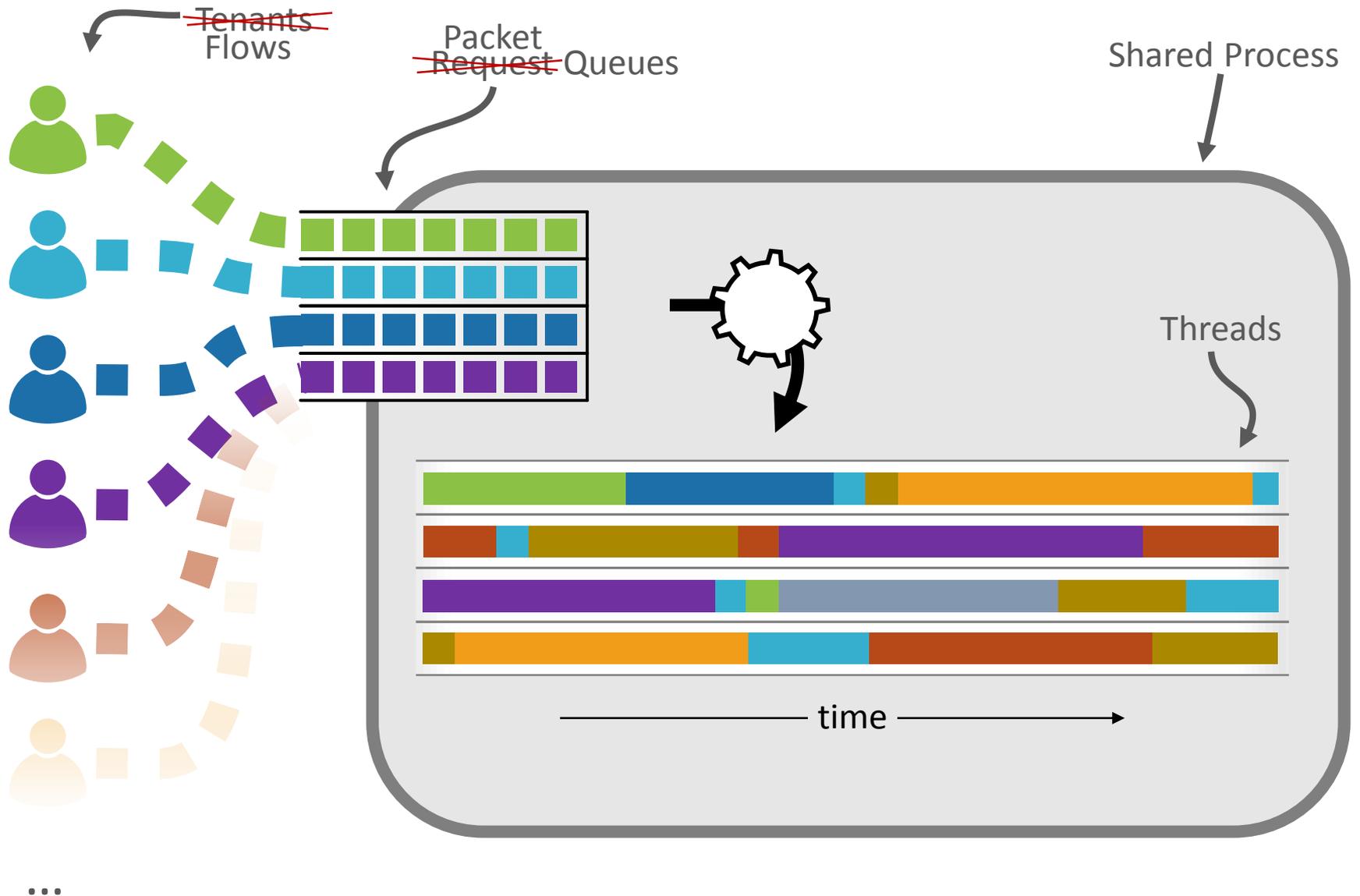


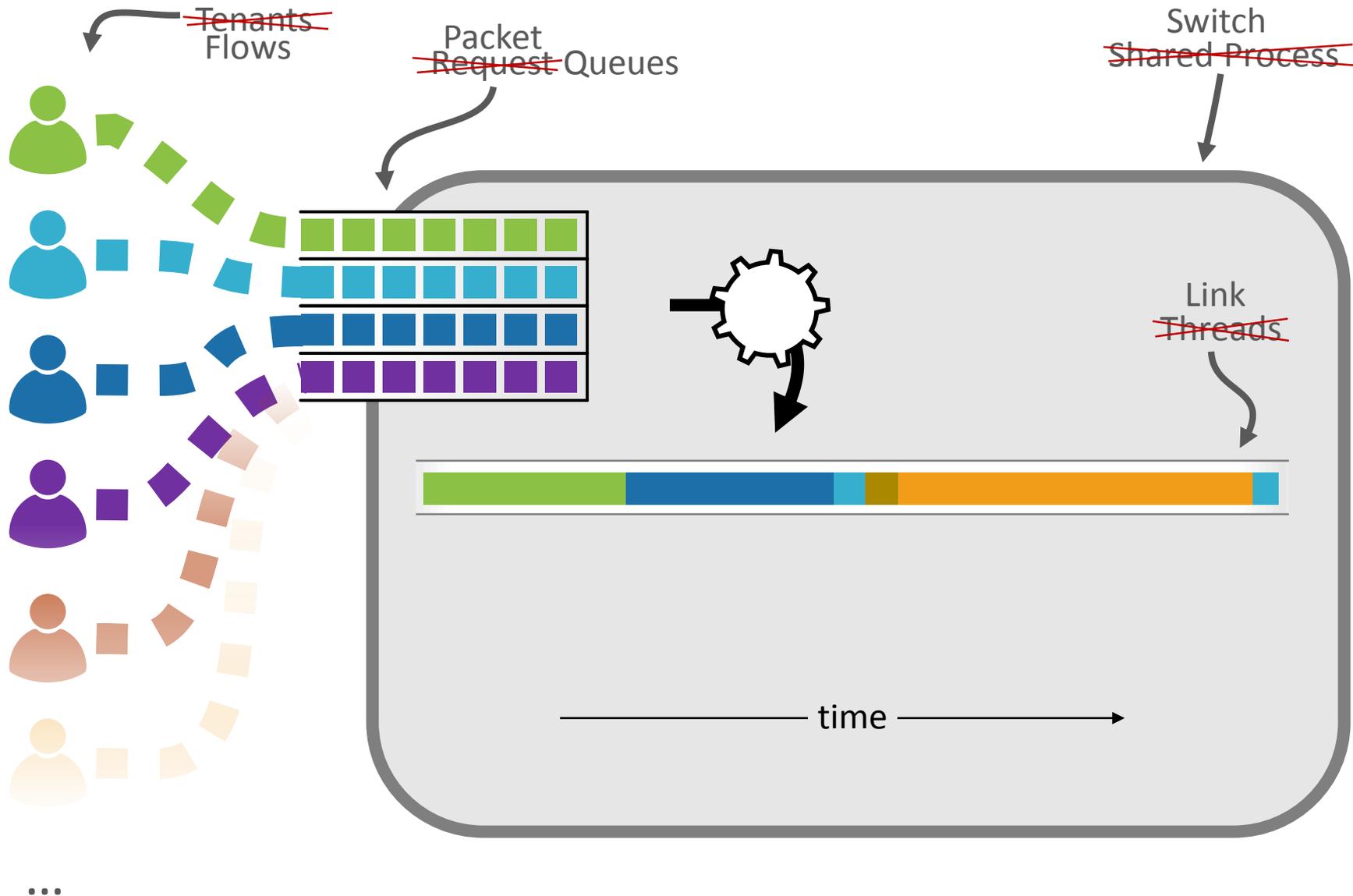


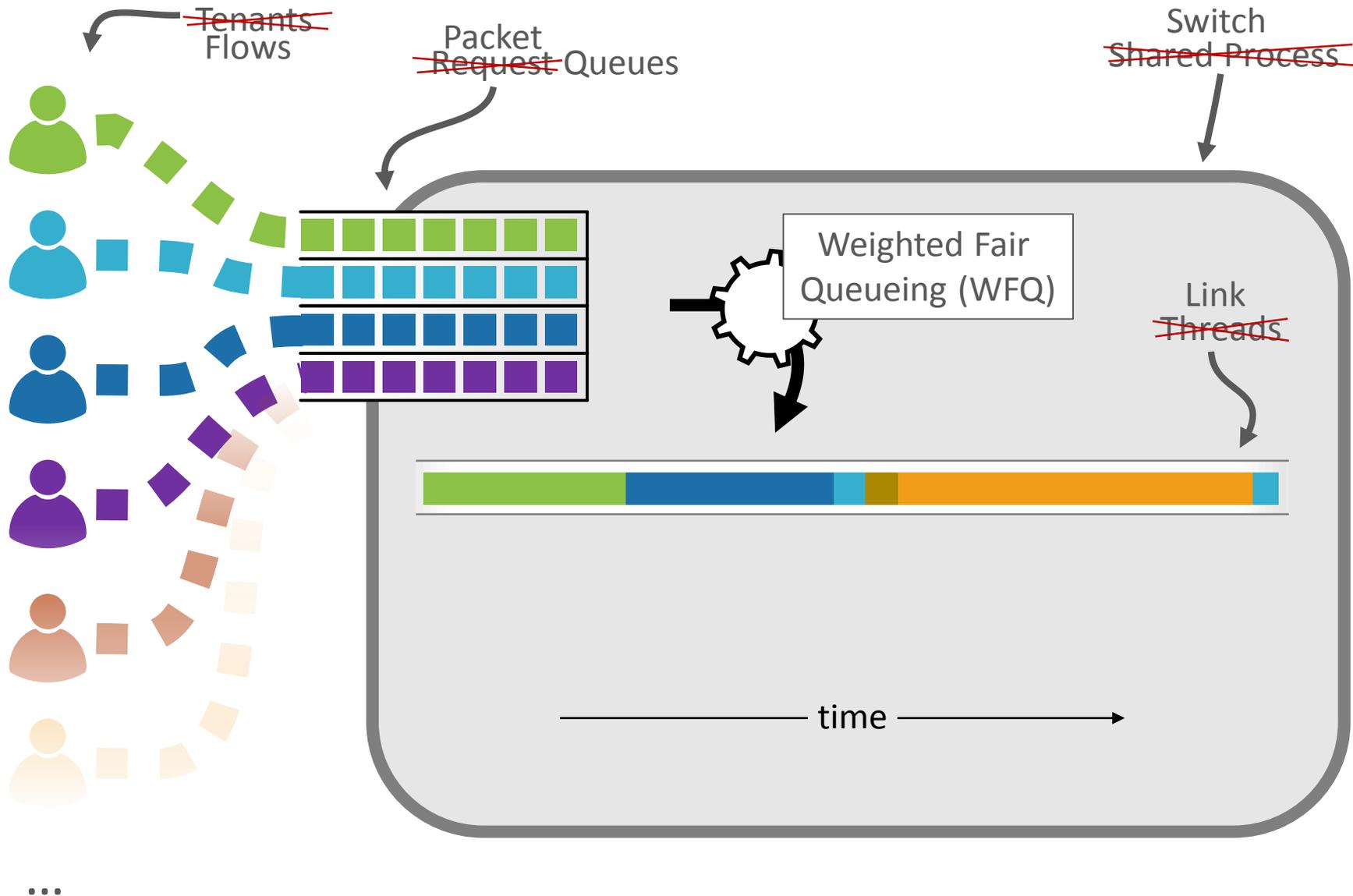


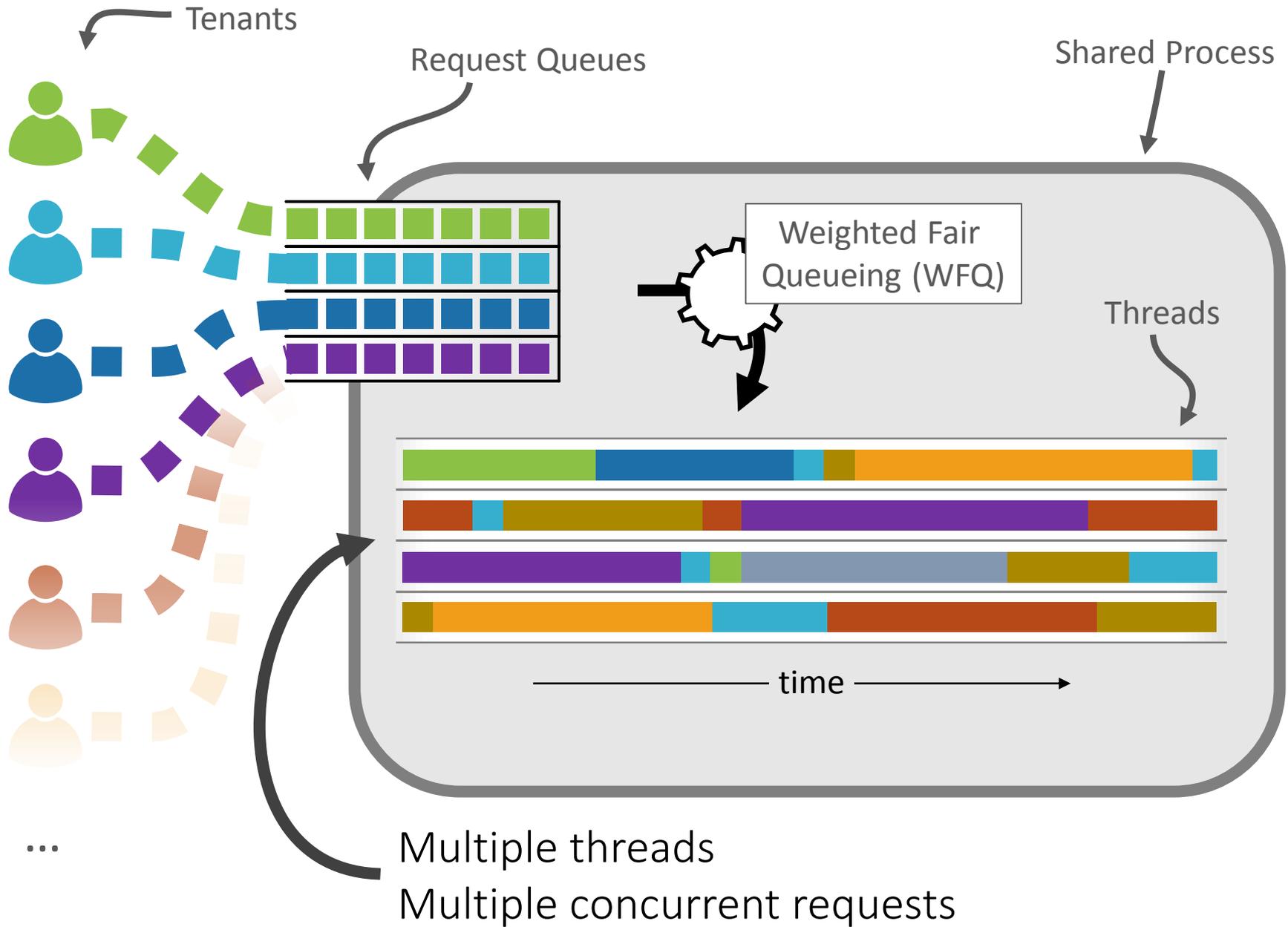


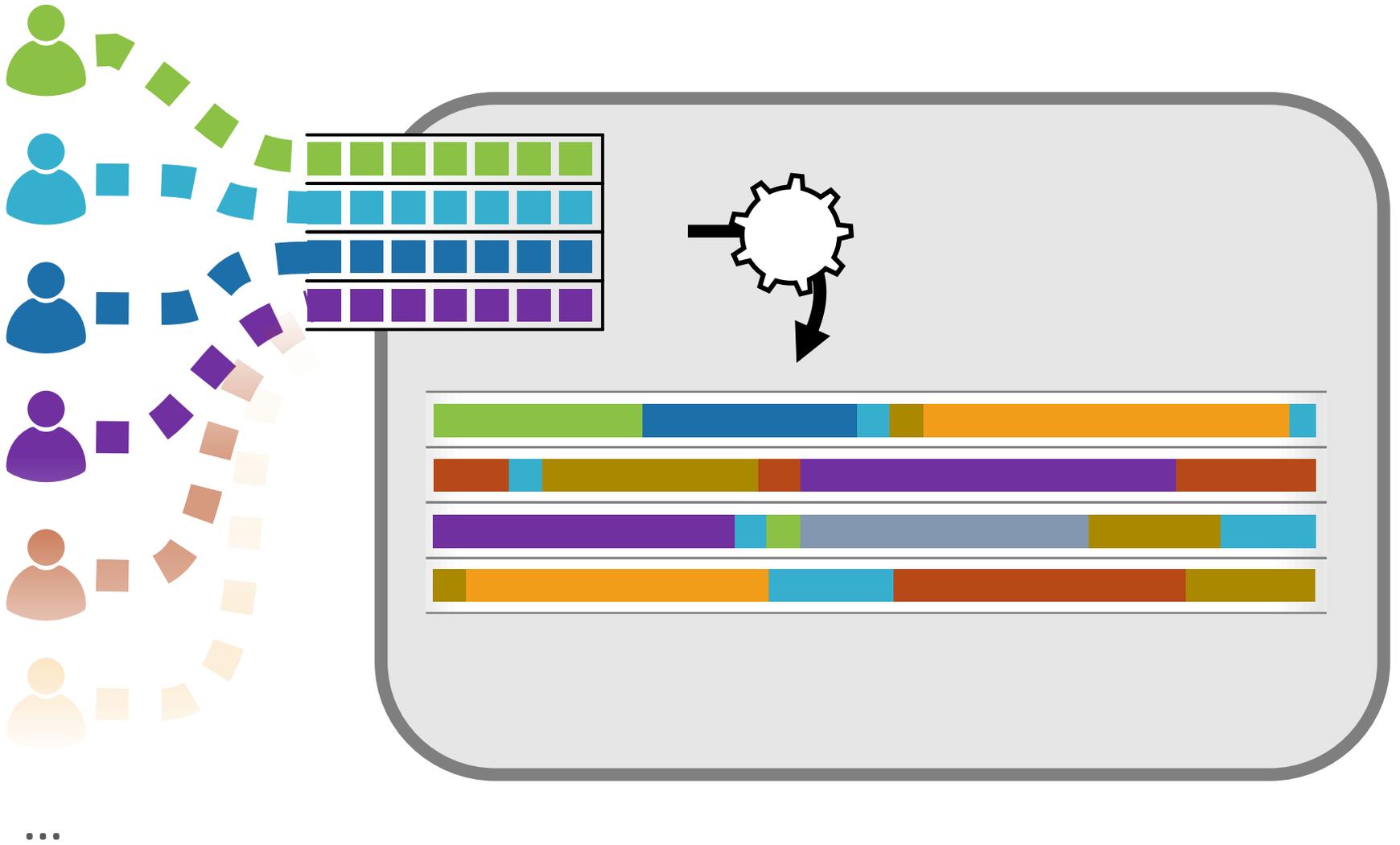




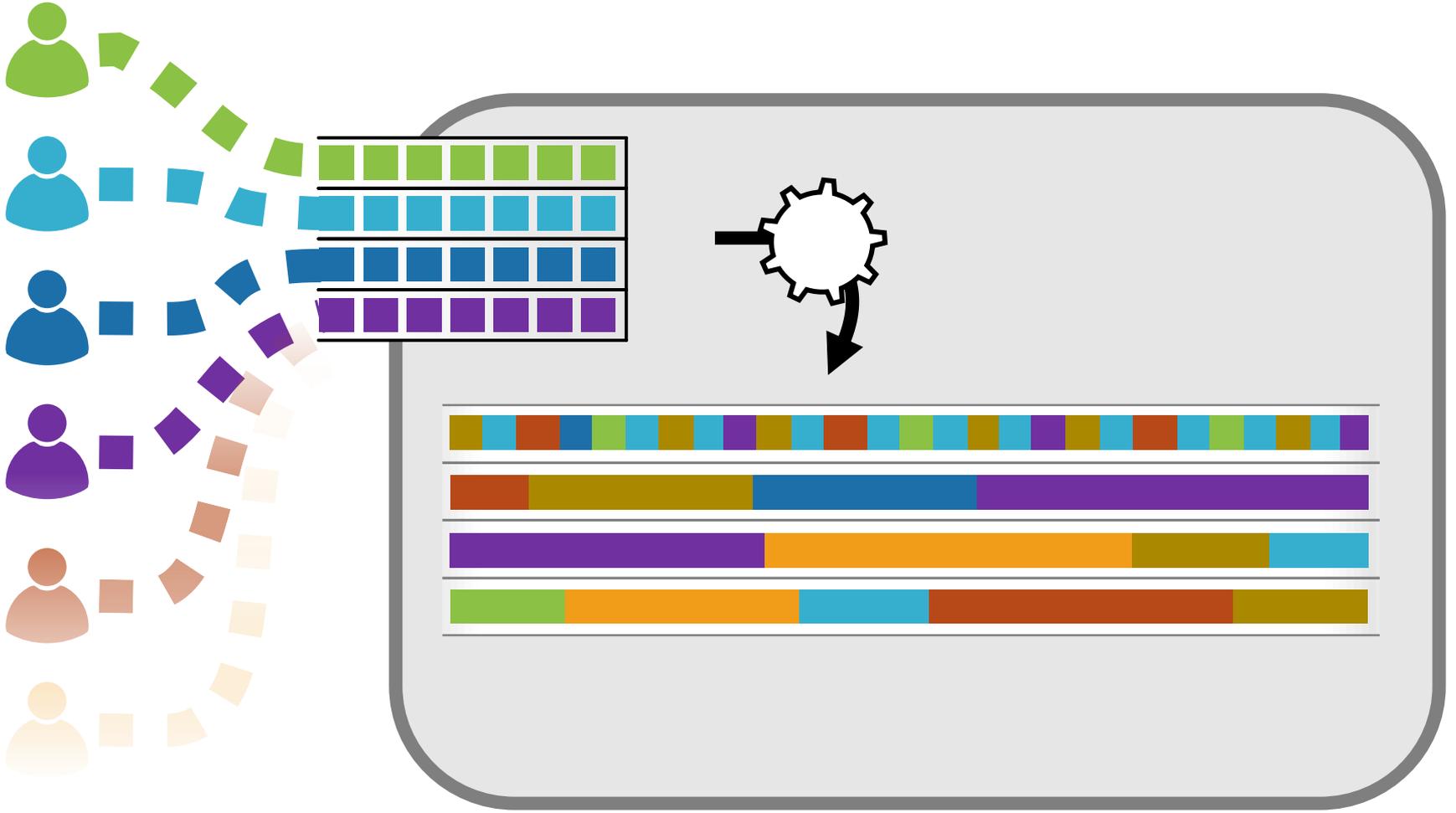




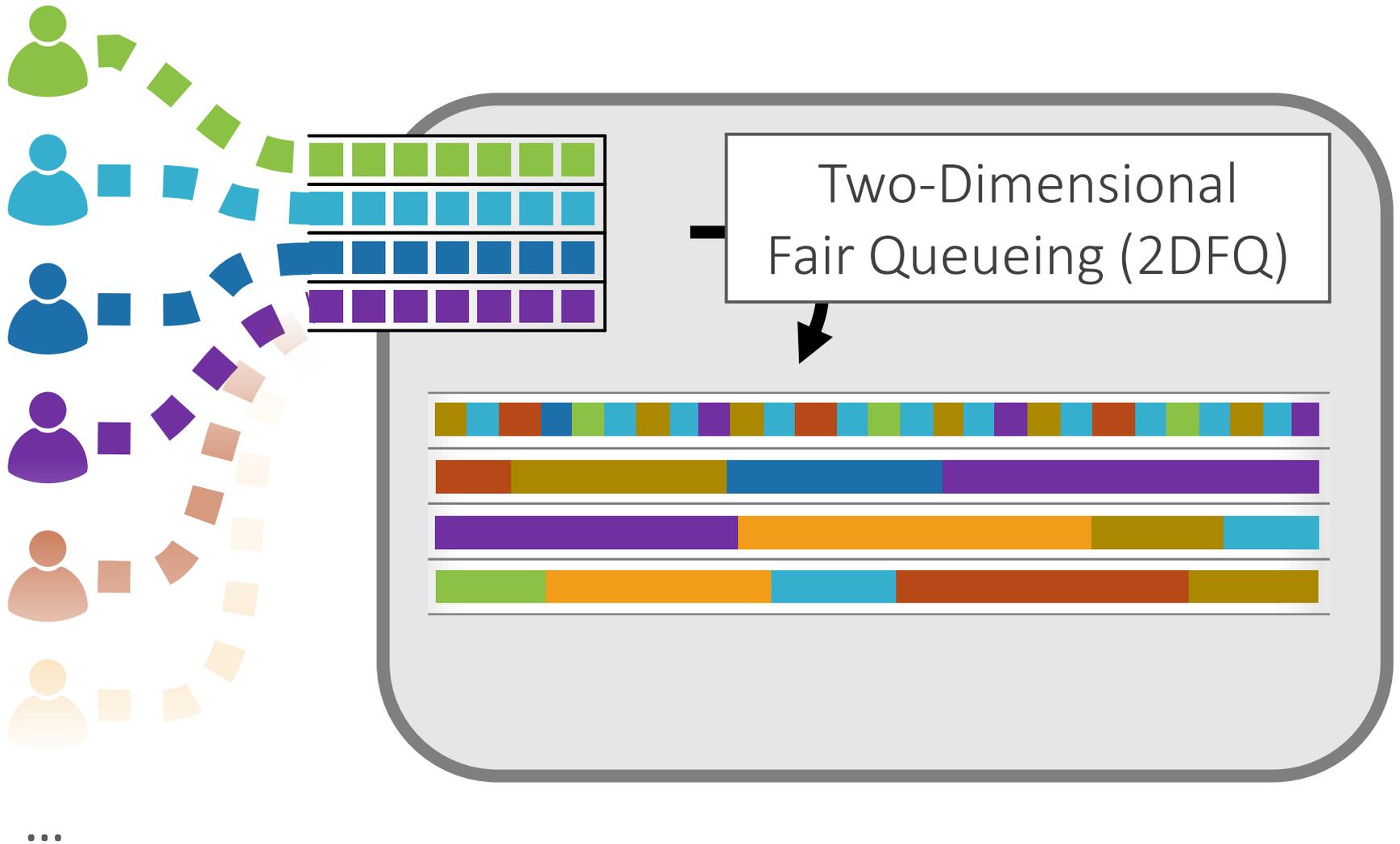




...



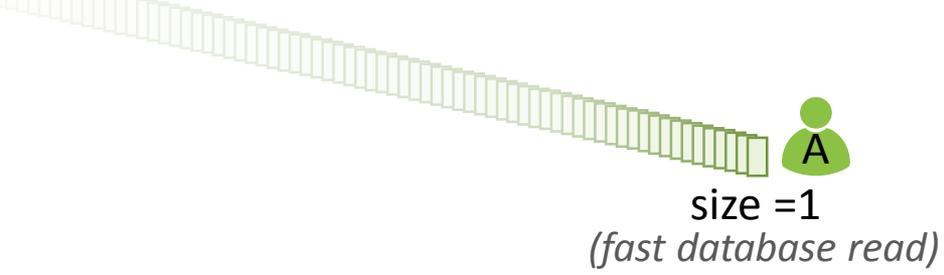
...

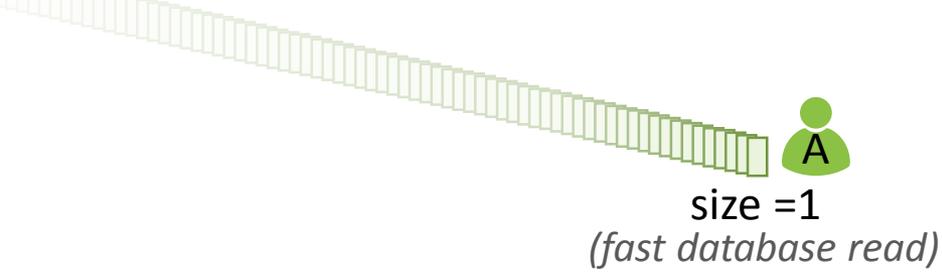






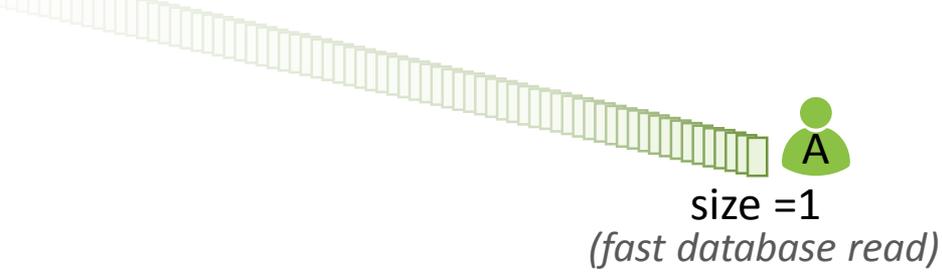
**size = 1**  
*(fast database read)*



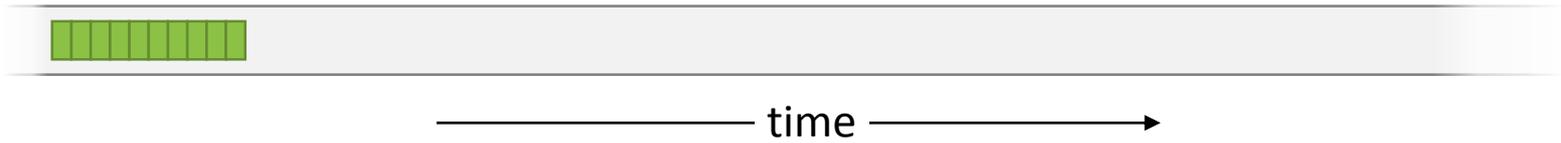


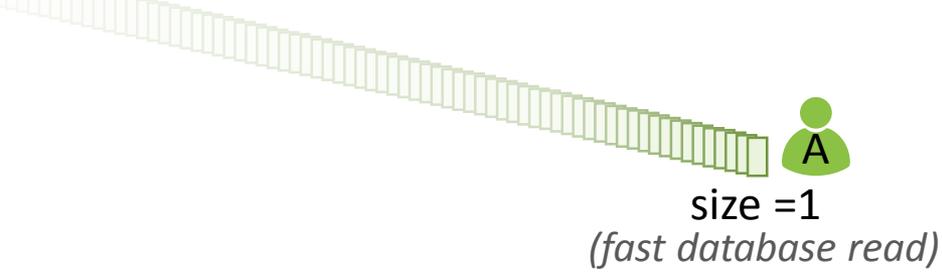
One thread:



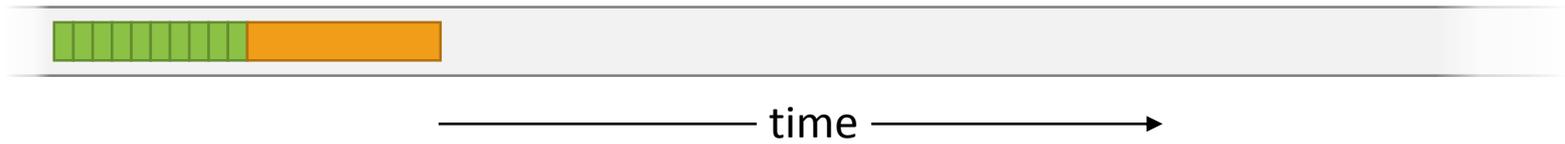


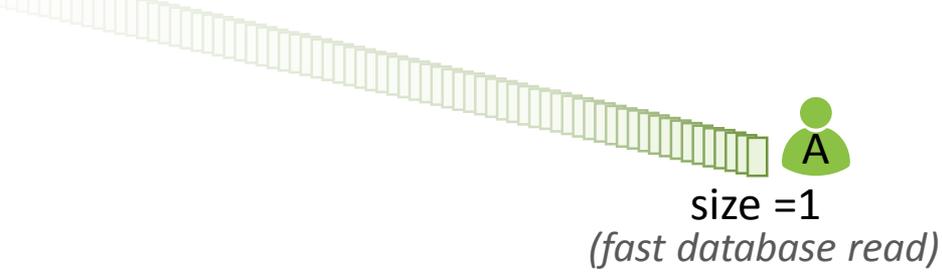
One thread:



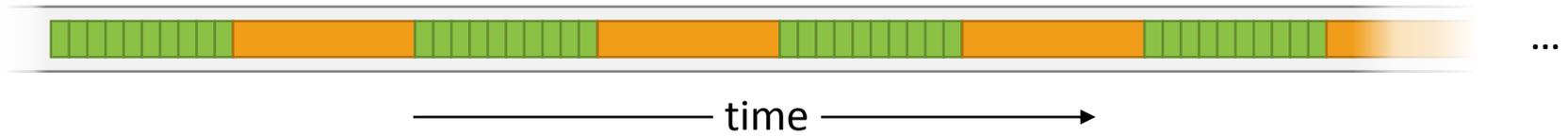


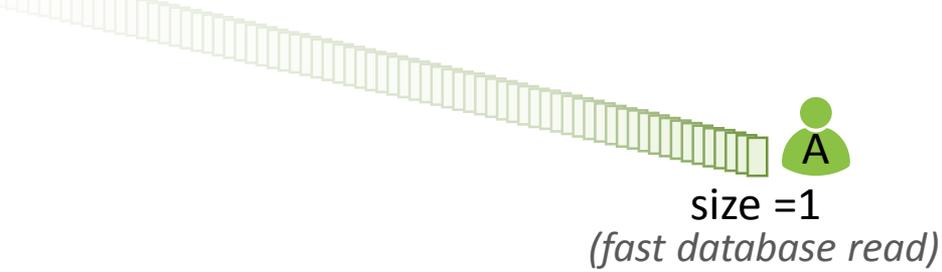
One thread:



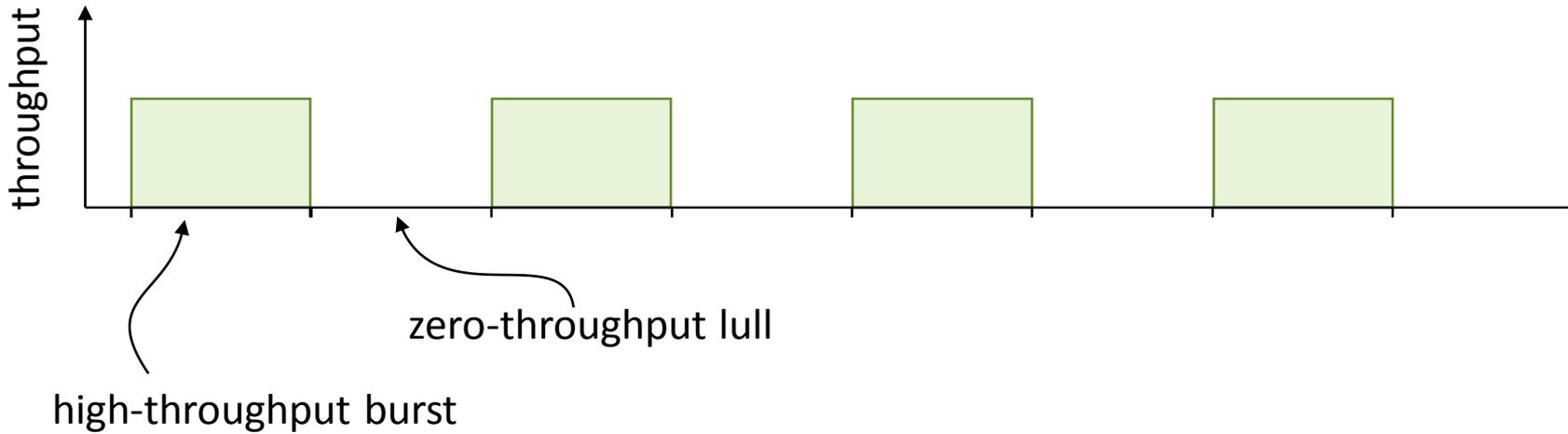
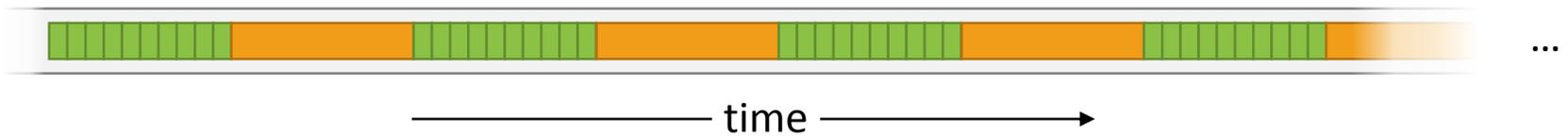


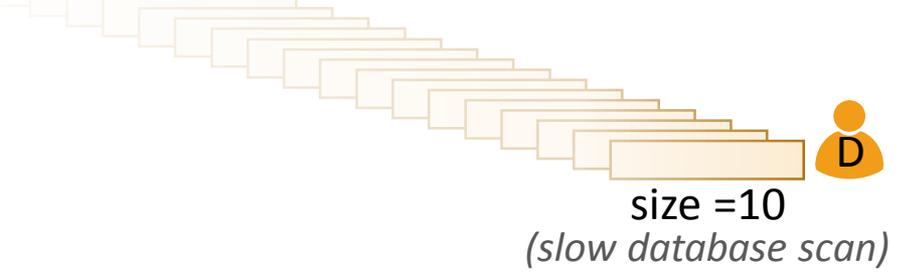
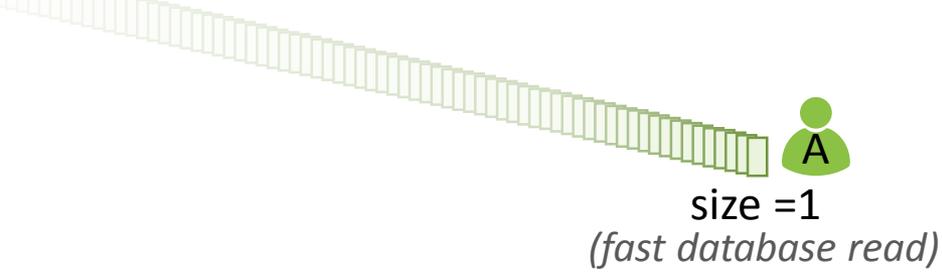
One thread:



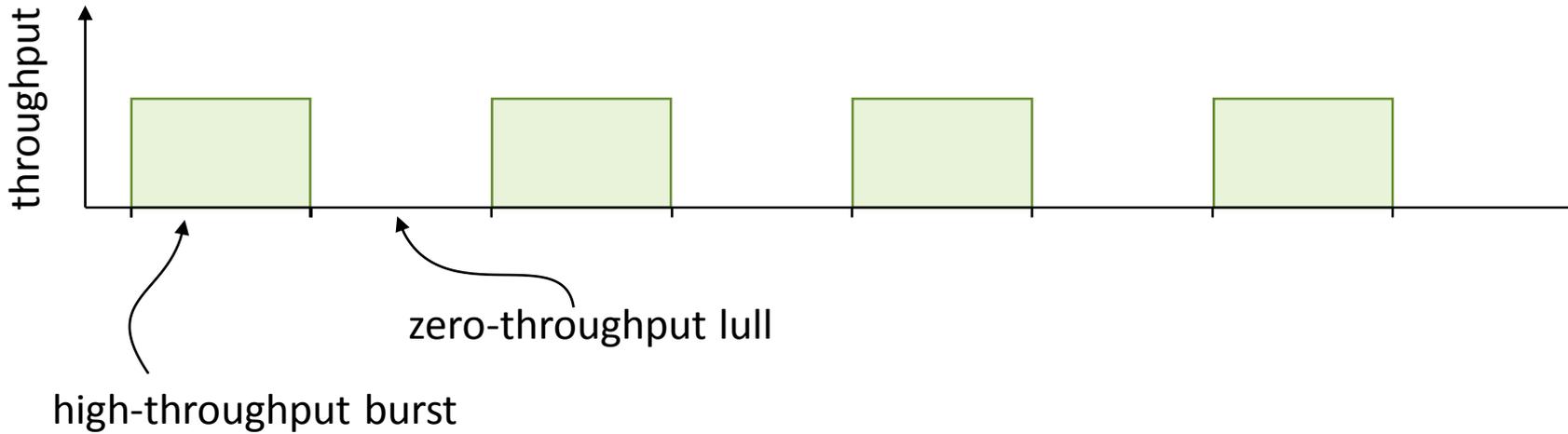
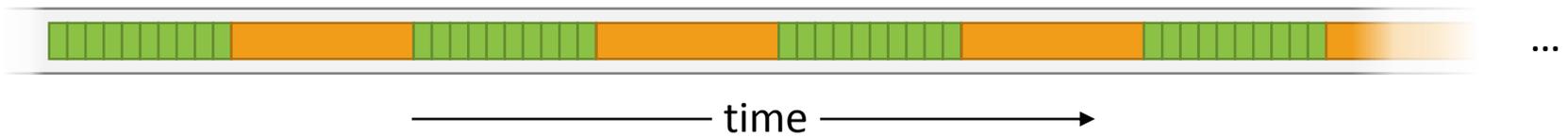


One thread:



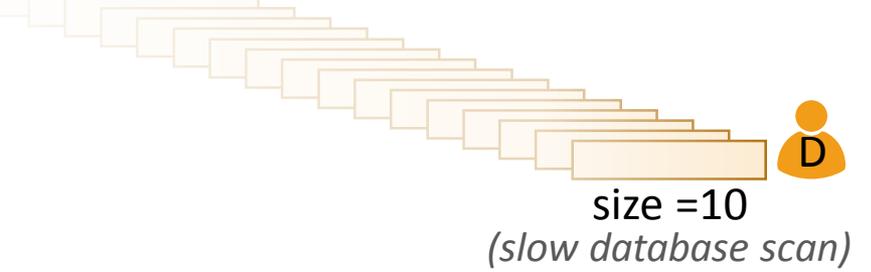
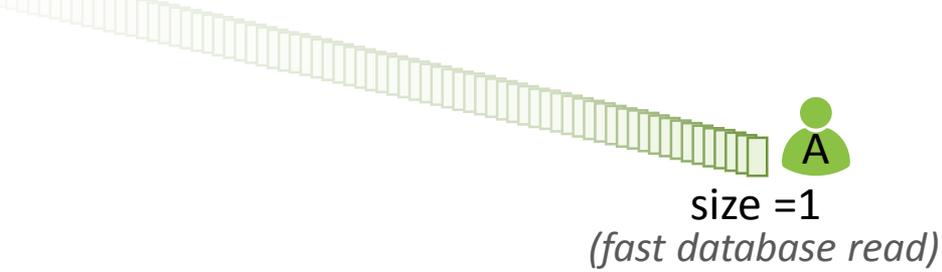


One thread:

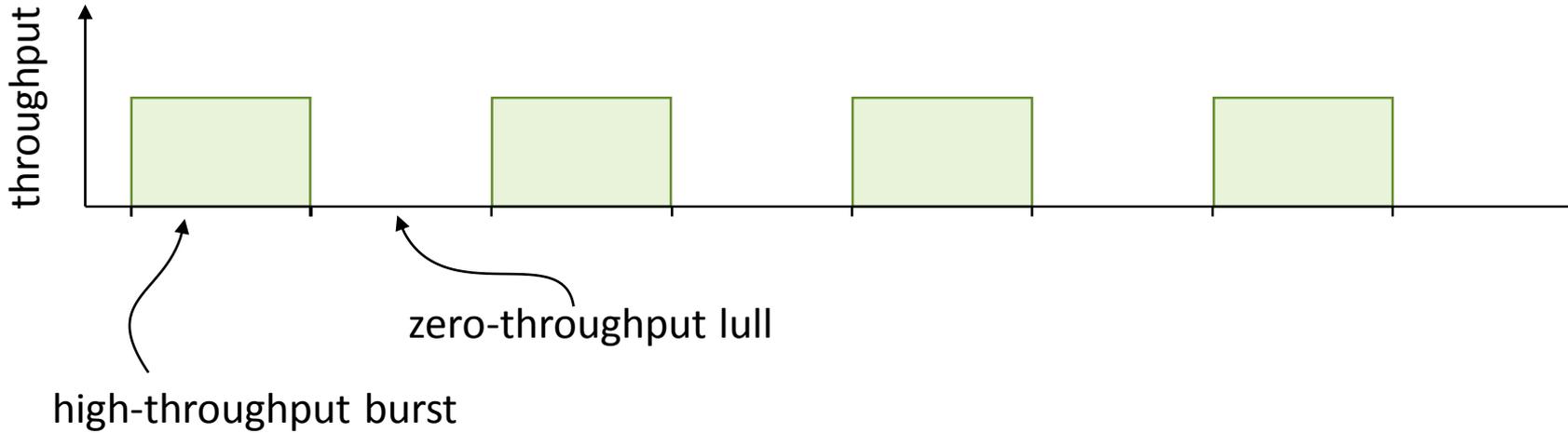
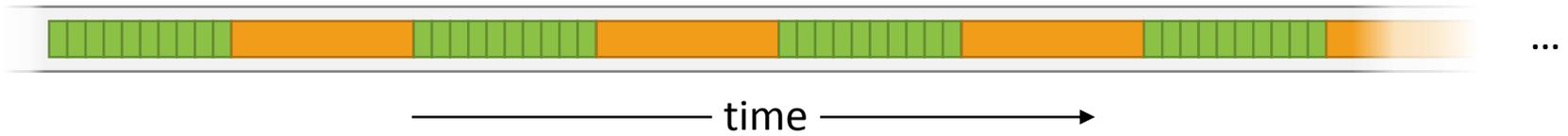


Ideal:





One thread:

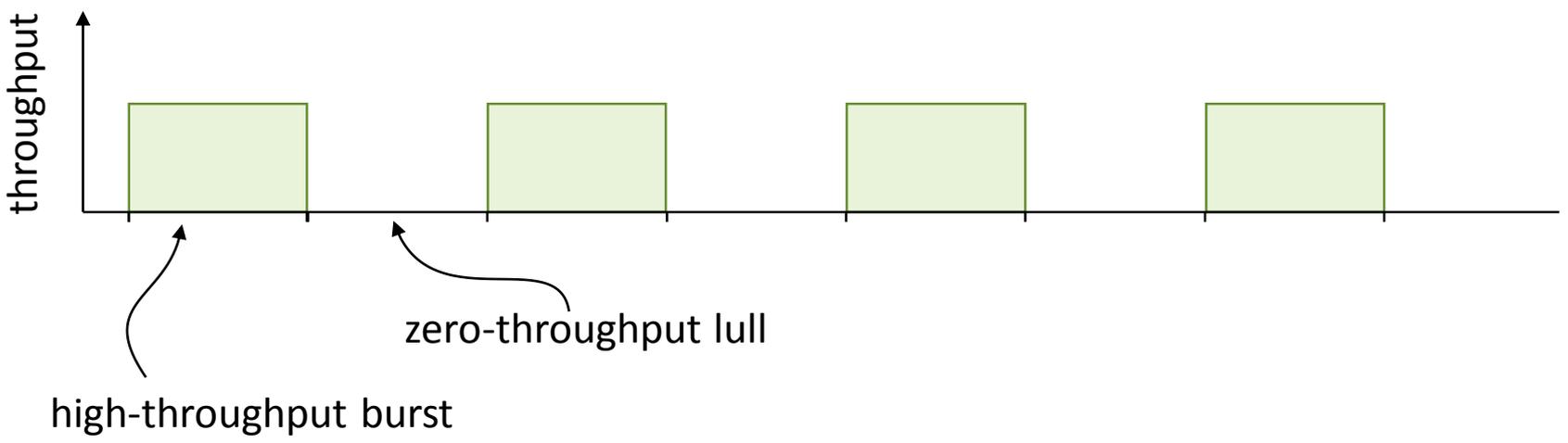
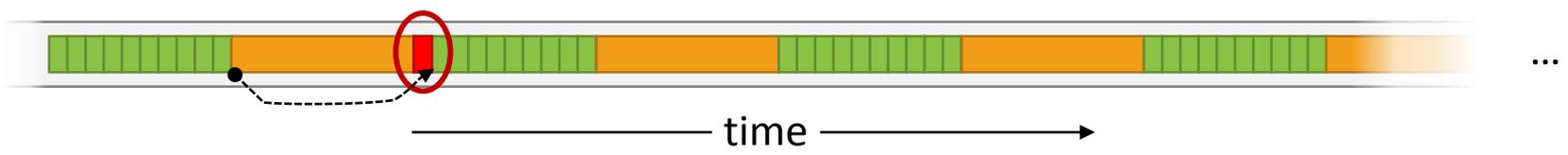


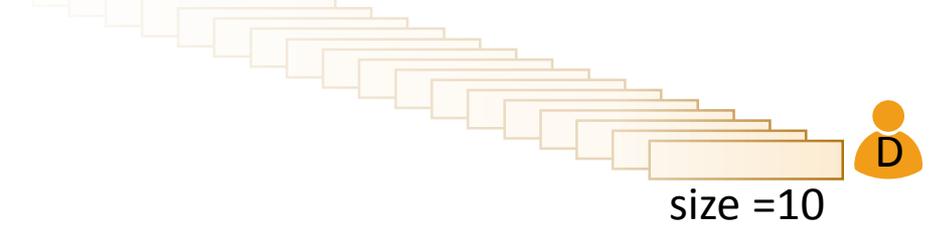
Ideal:



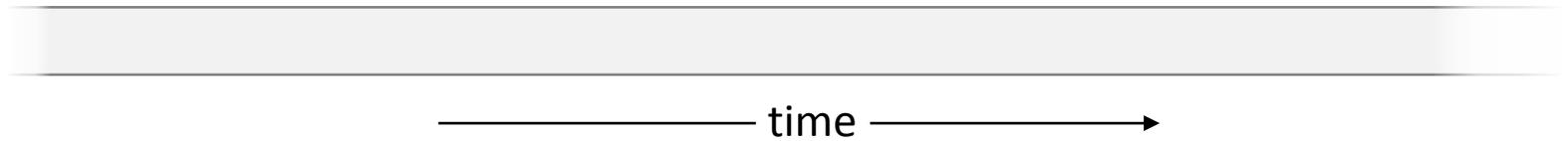


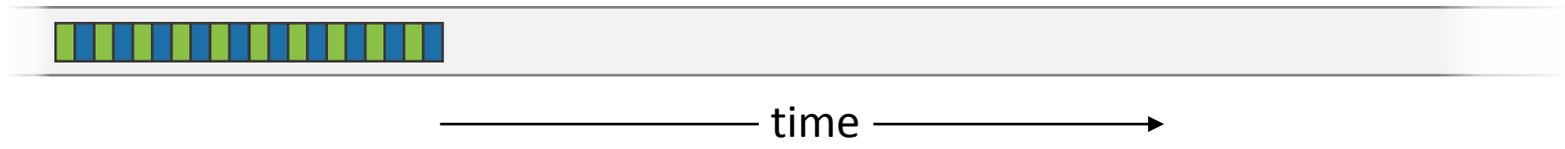
One thread:

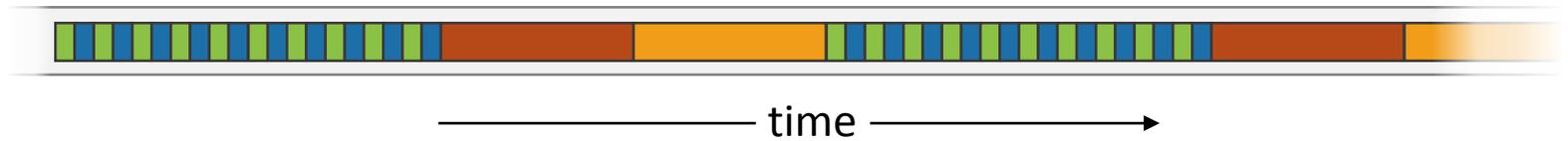


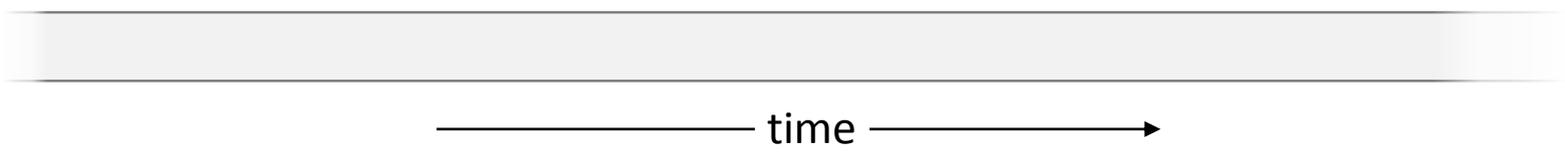
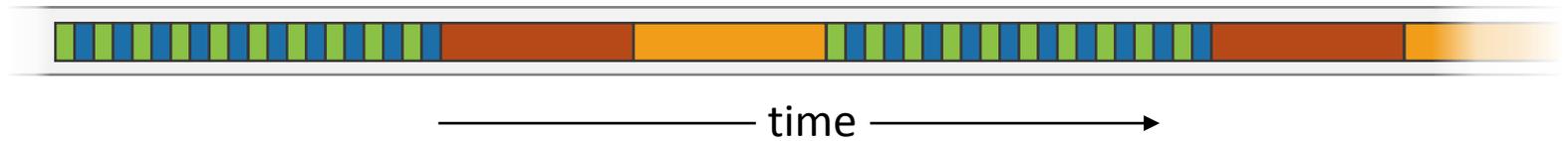




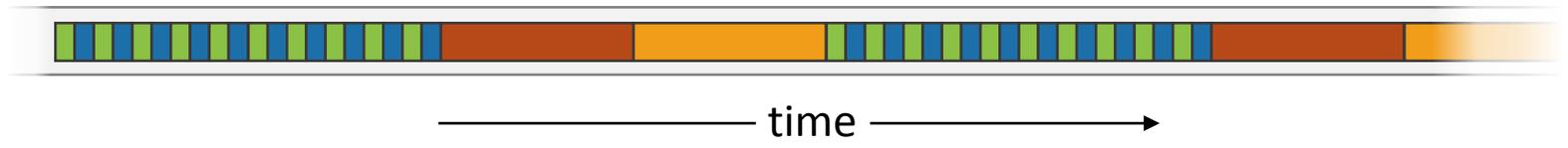


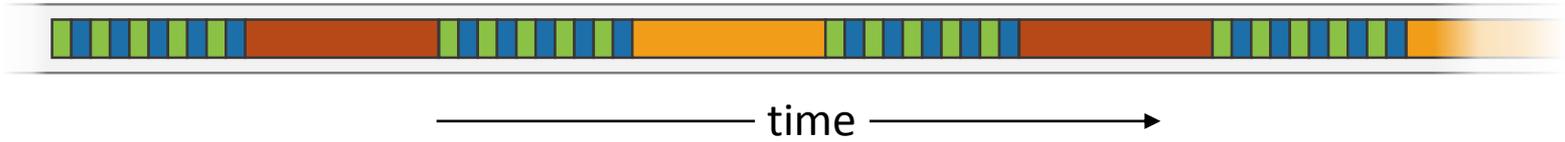
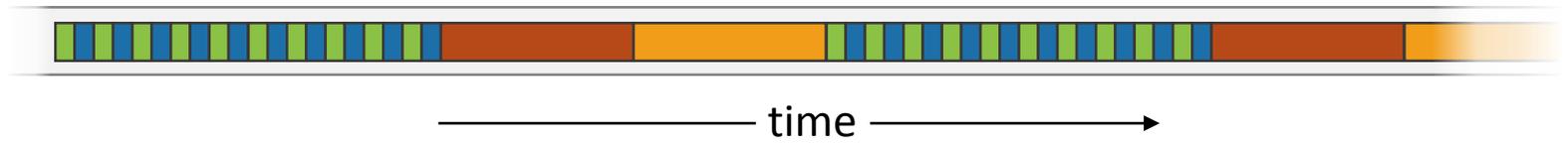


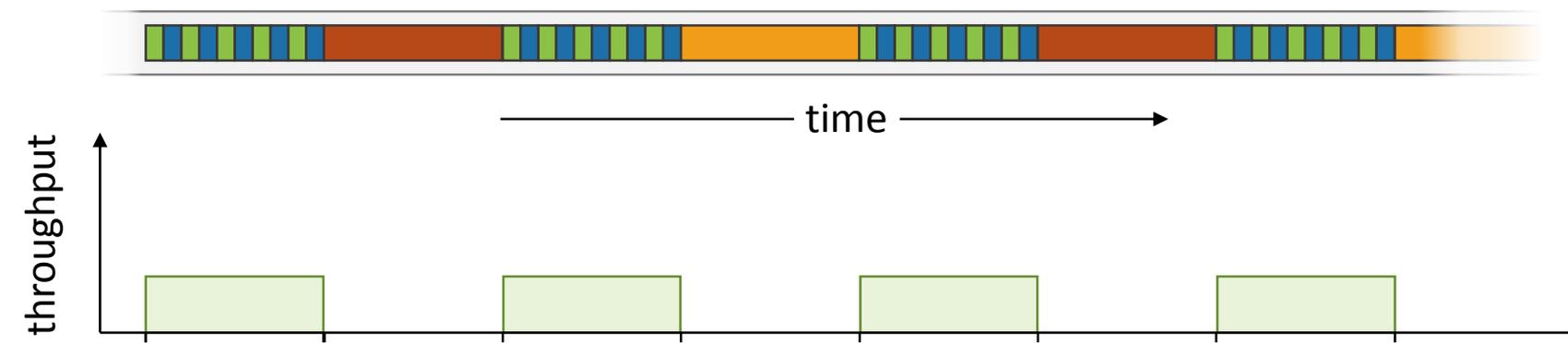
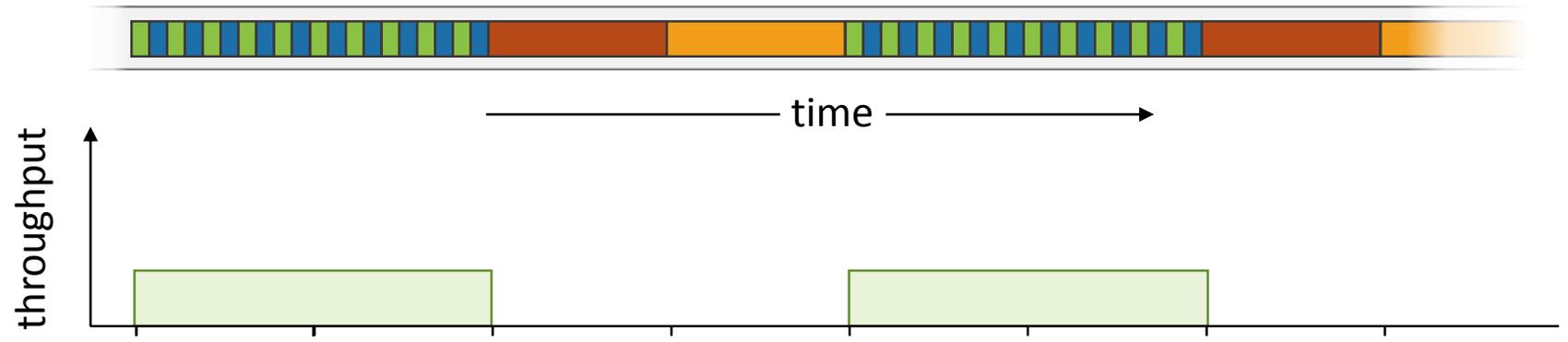


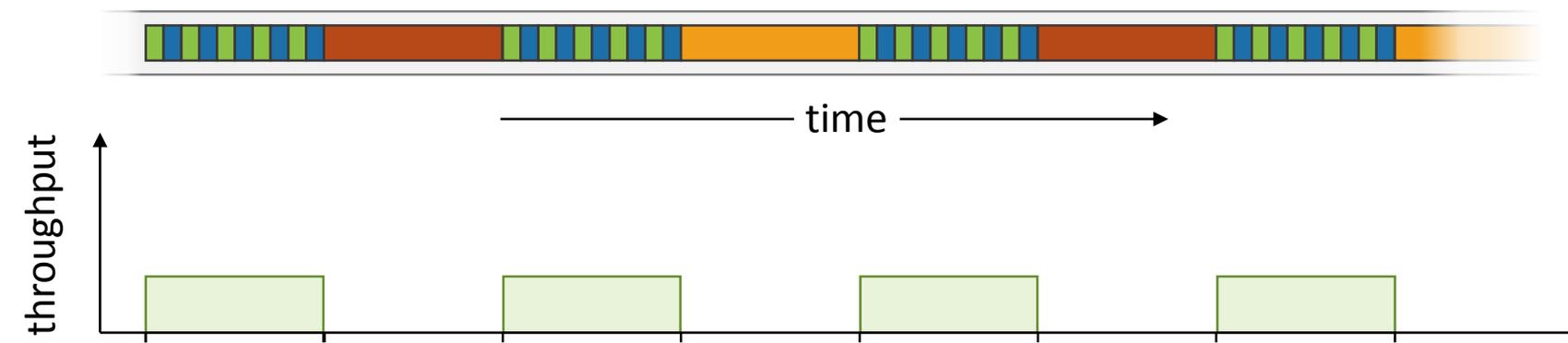
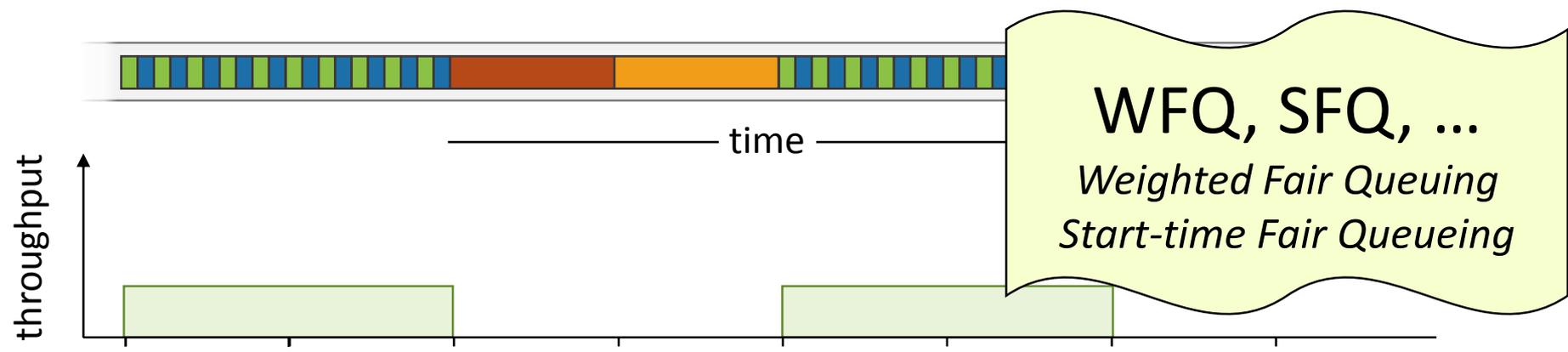


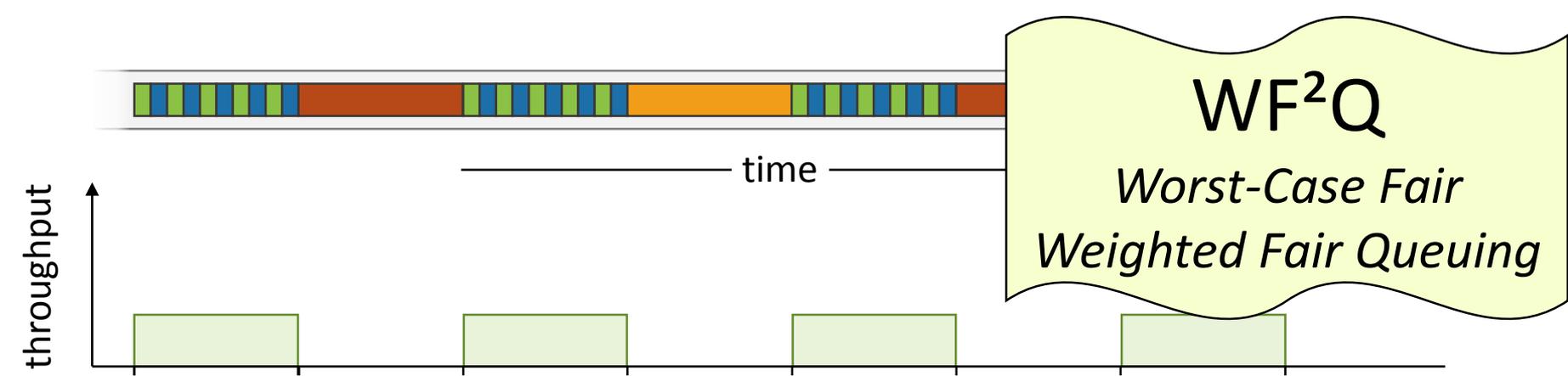
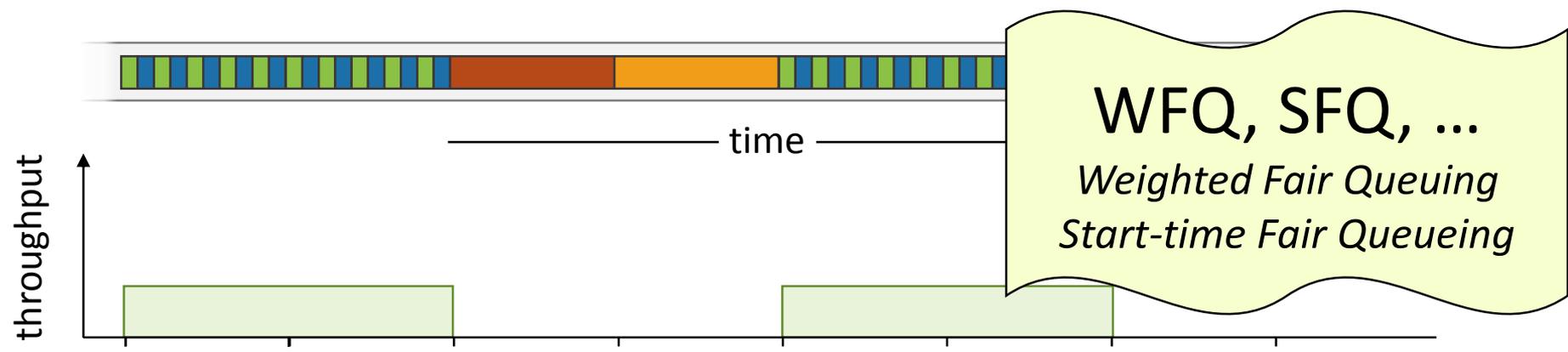


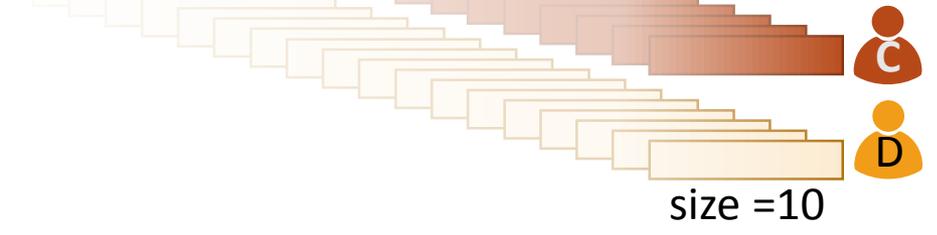
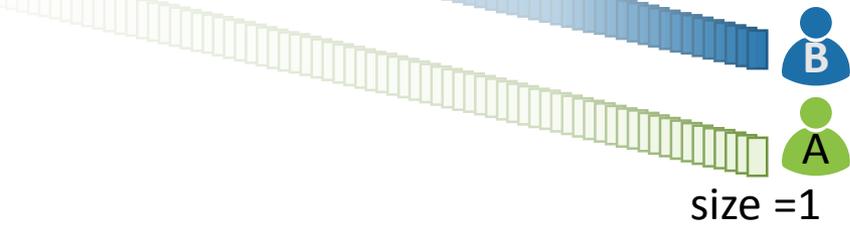








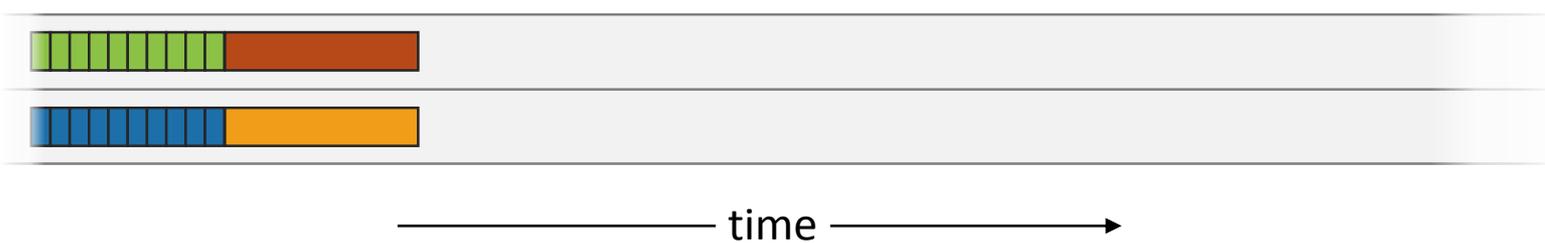
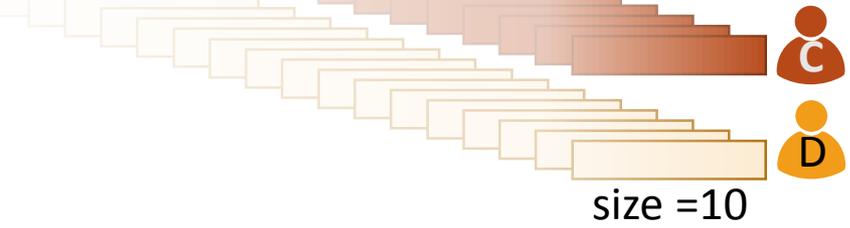


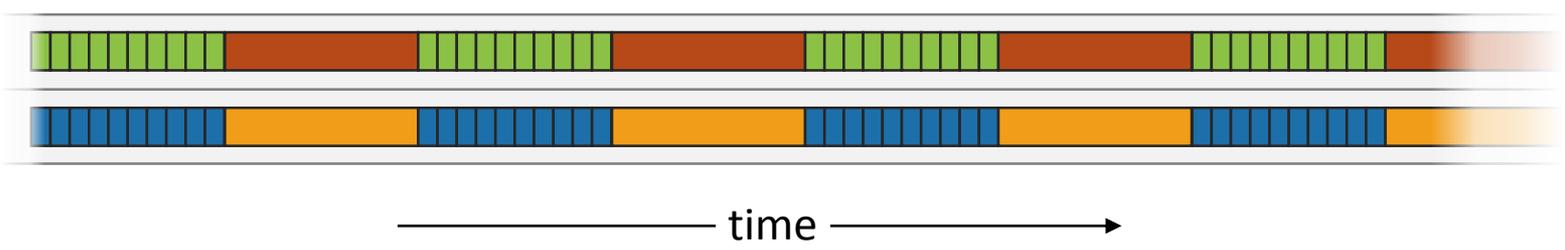


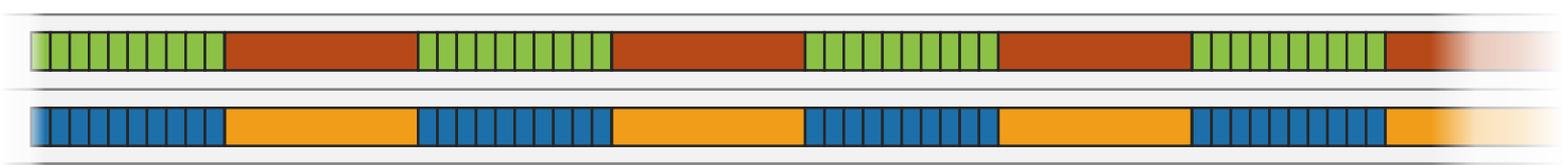
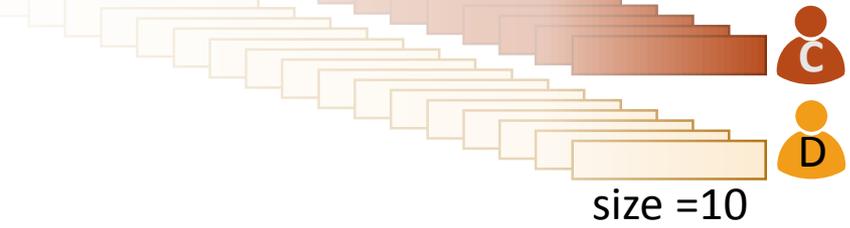
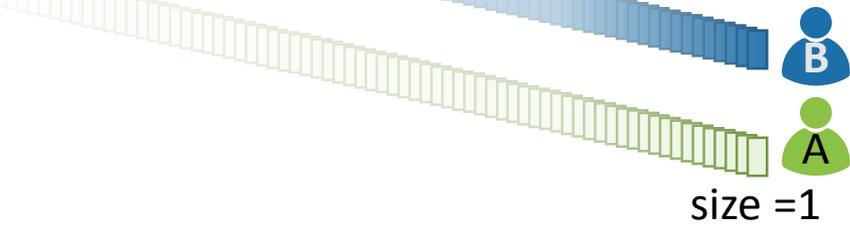


time →





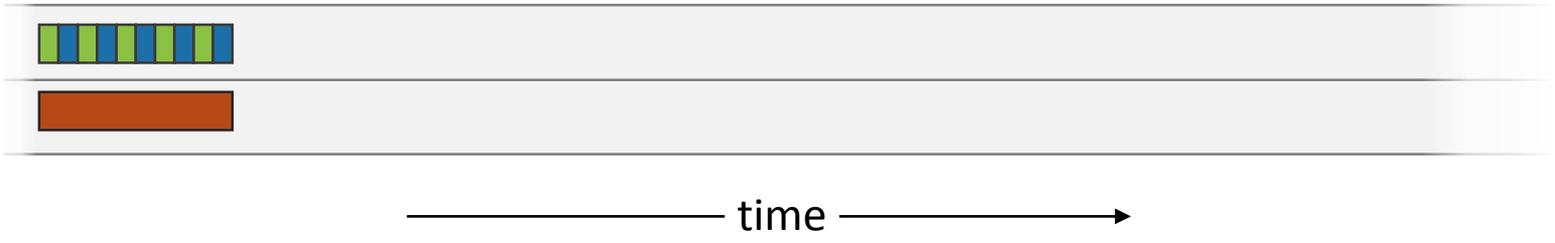
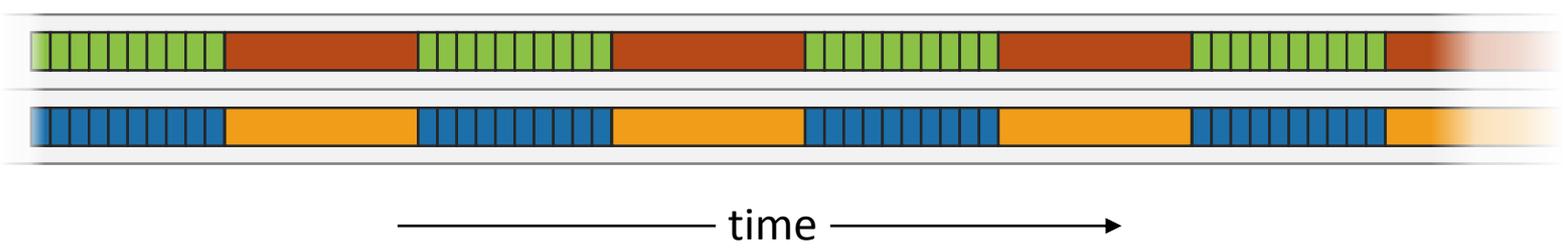


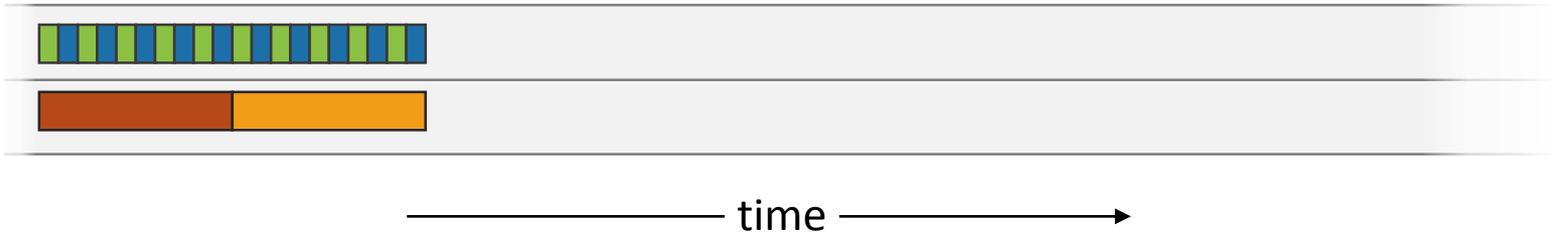
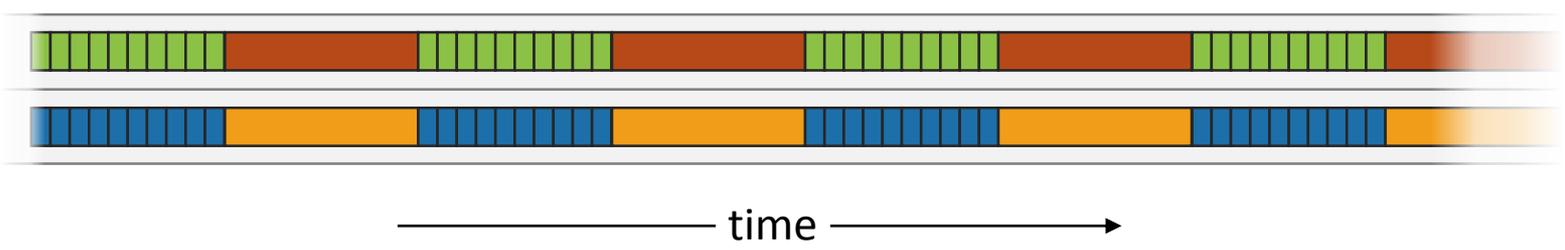
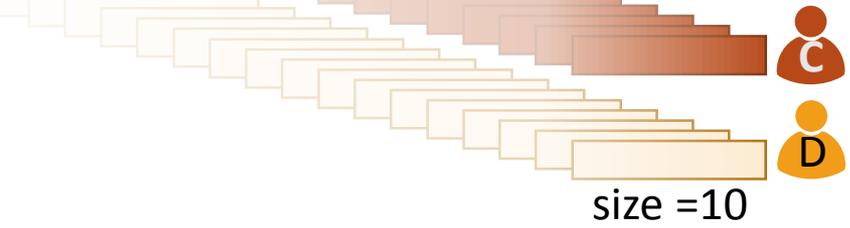


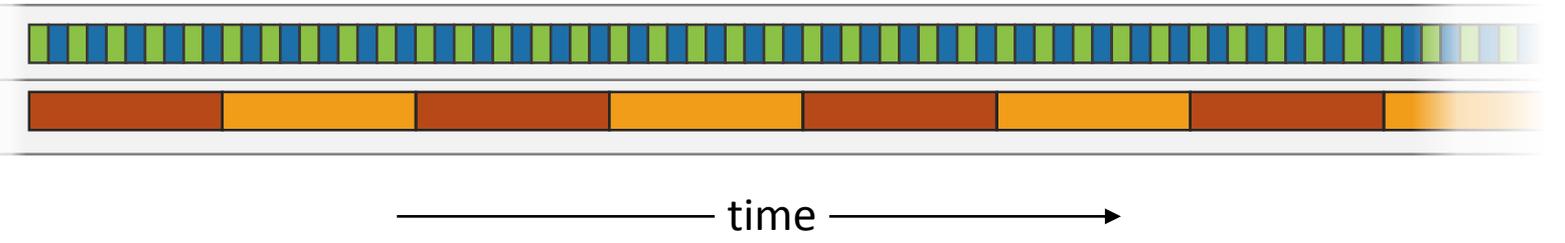
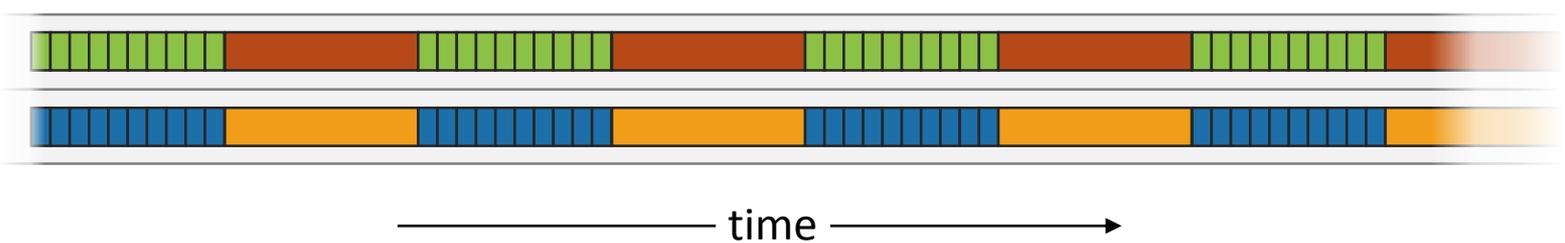
time →

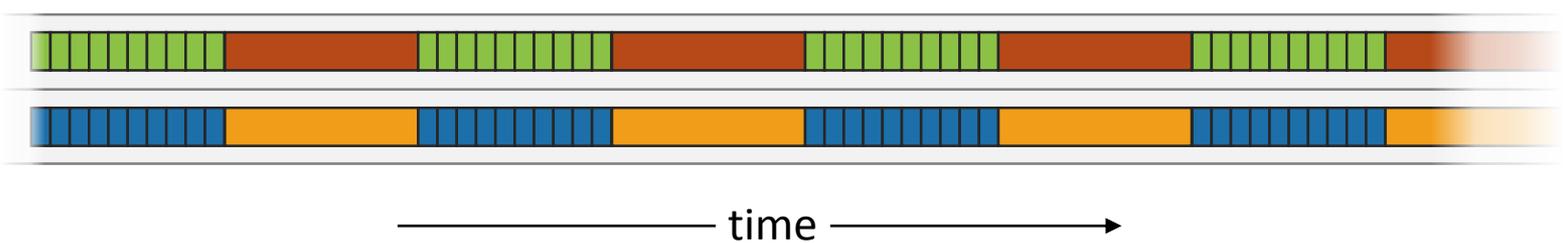


time →

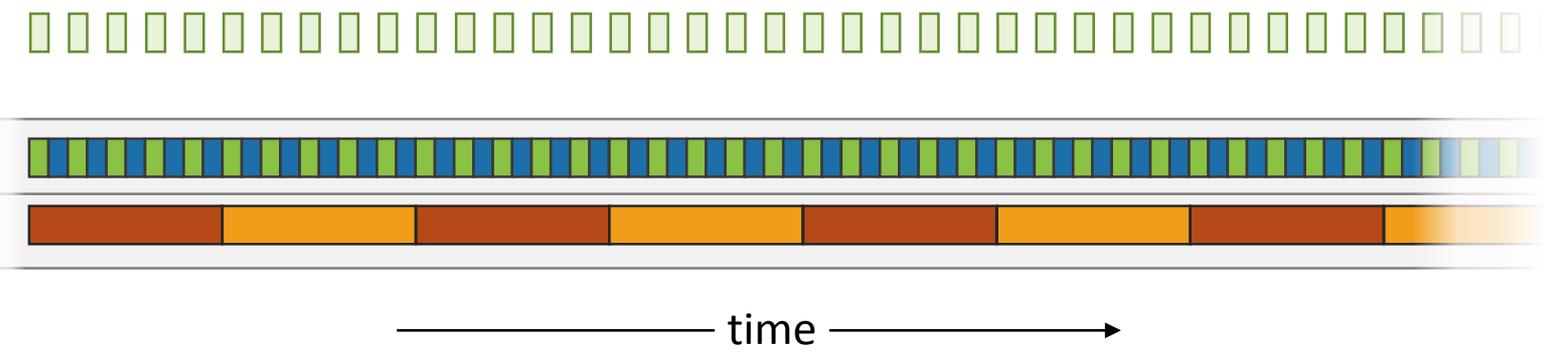


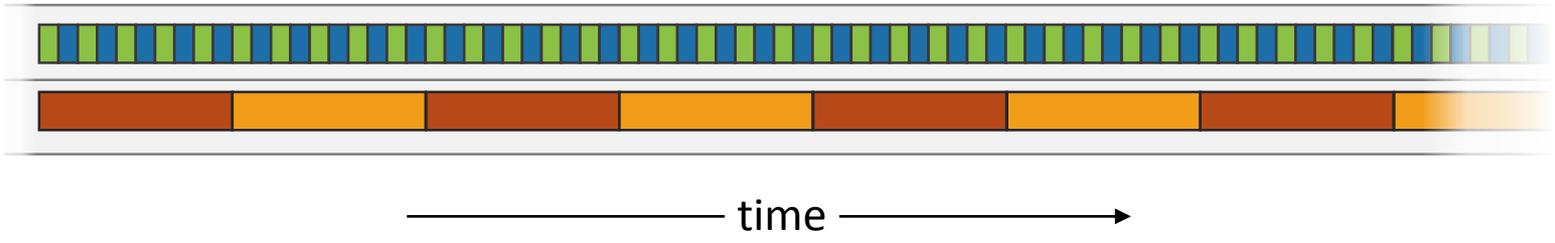
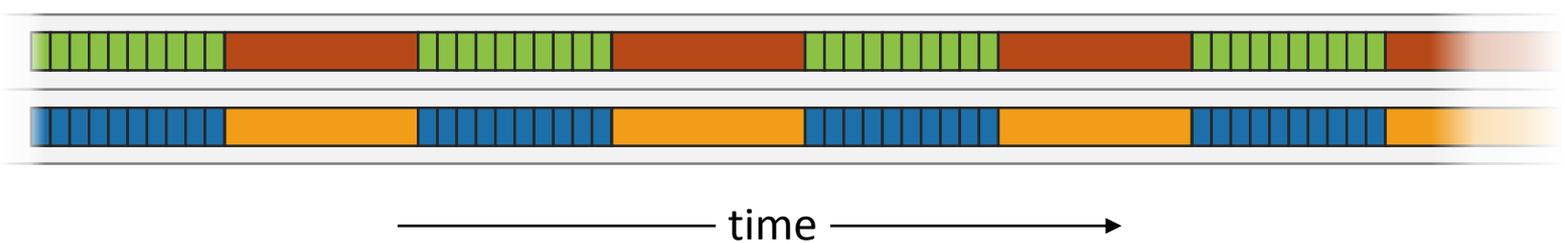




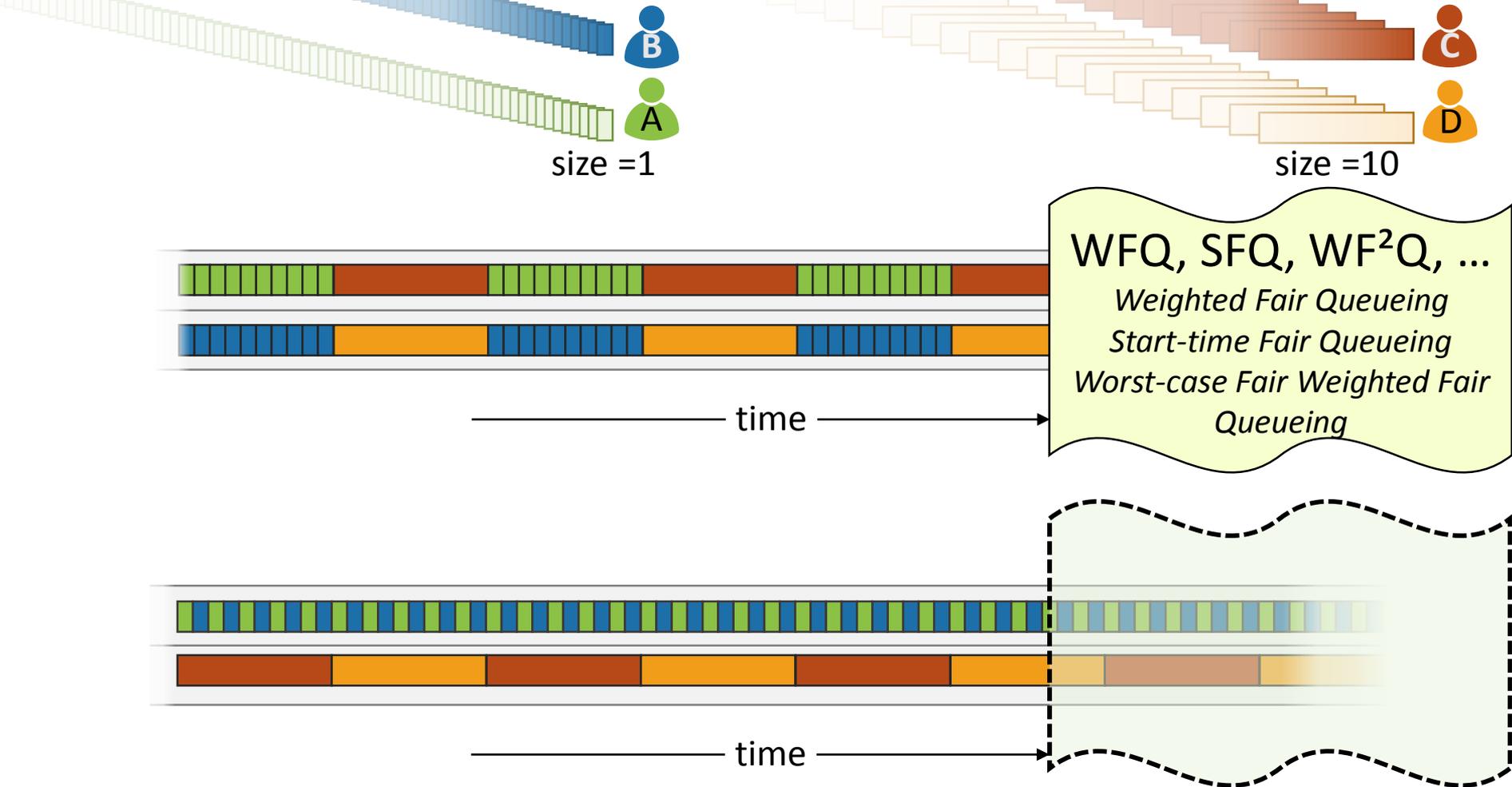


Ideal:

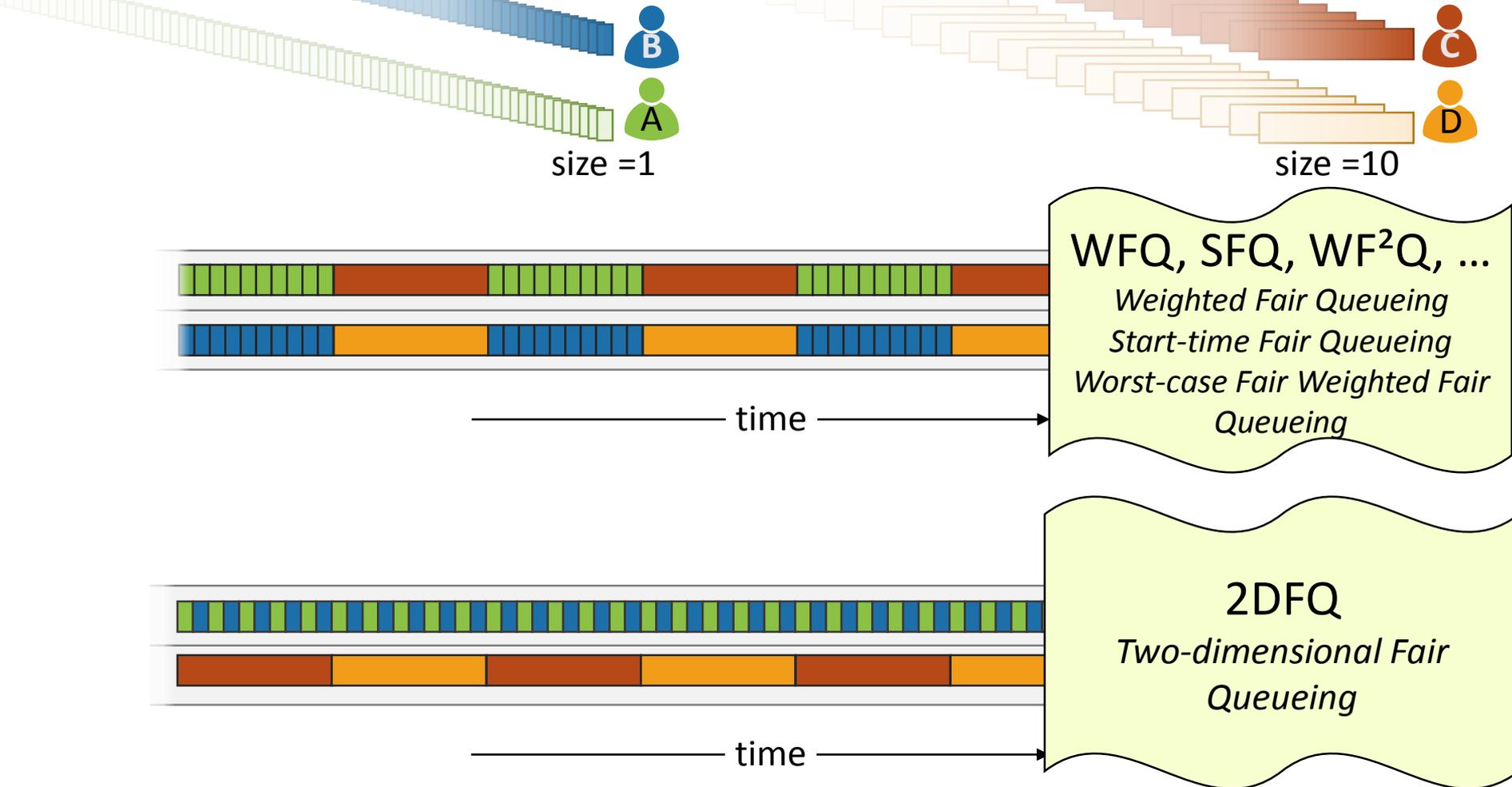




More threads → Opportunity to reduce burstiness



More threads → Opportunity to reduce burstiness



More threads → Opportunity to reduce burstiness

# Challenges



Tenants with small requests are affected



Tenants with small requests are affected



Burstiness is proportional to size of large requests



Tenants with small requests are affected

Burstiness is proportional to size of large requests



## Cloud services:

4+ orders of magnitude variation in cost



Tenants with small requests are affected

Burstiness is proportional to size of large requests



## Cloud services:

4+ orders of magnitude variation in cost



Size is used by scheduler to make scheduling decisions



Tenants with small requests are affected

Burstiness is proportional to size of large requests



## Cloud services:

4+ orders of magnitude variation in cost



Size is used by scheduler to make scheduling decisions



Tenants with small requests are affected

Burstiness is proportional to size of large requests



## Cloud services:

4+ orders of magnitude variation in cost

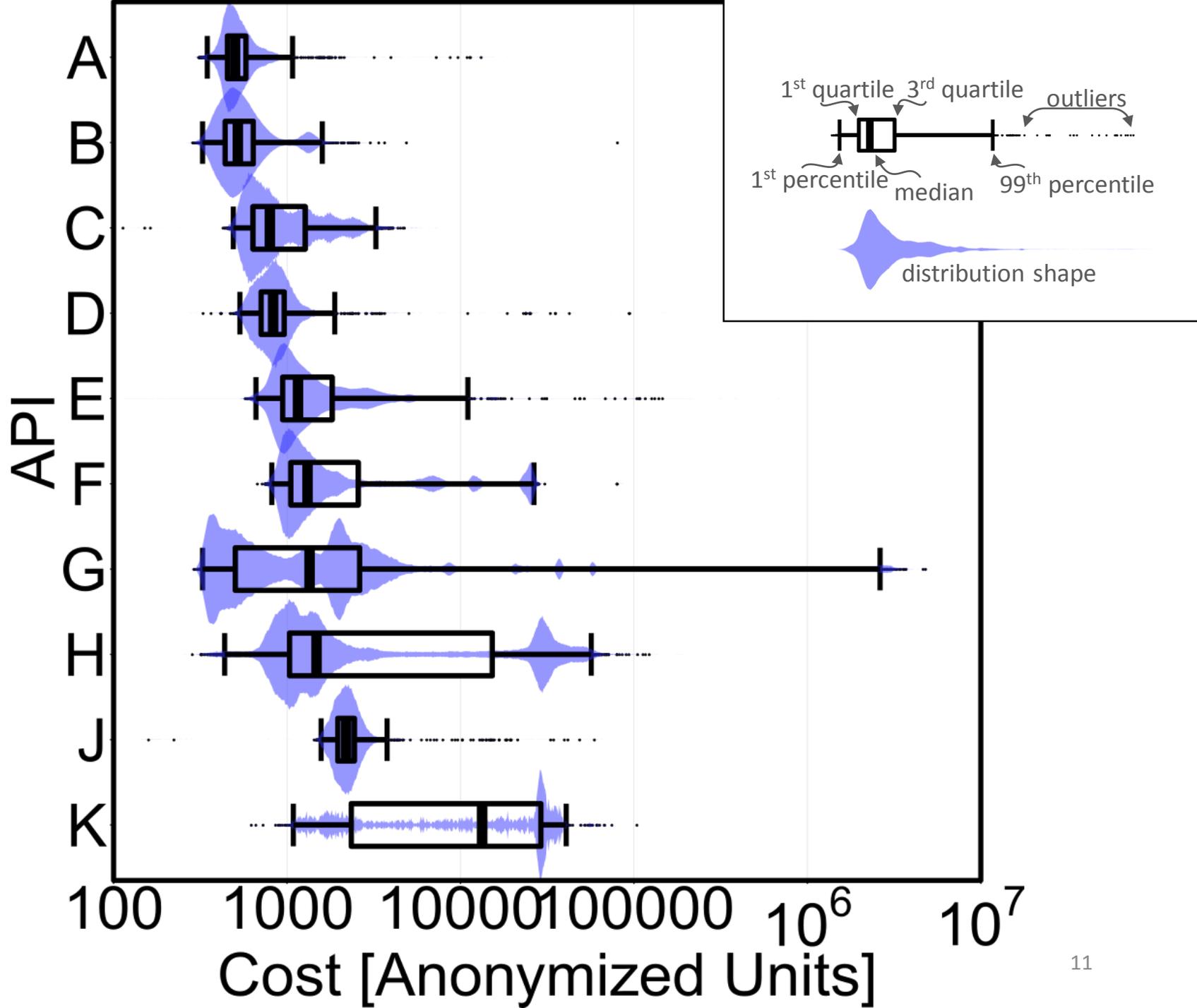


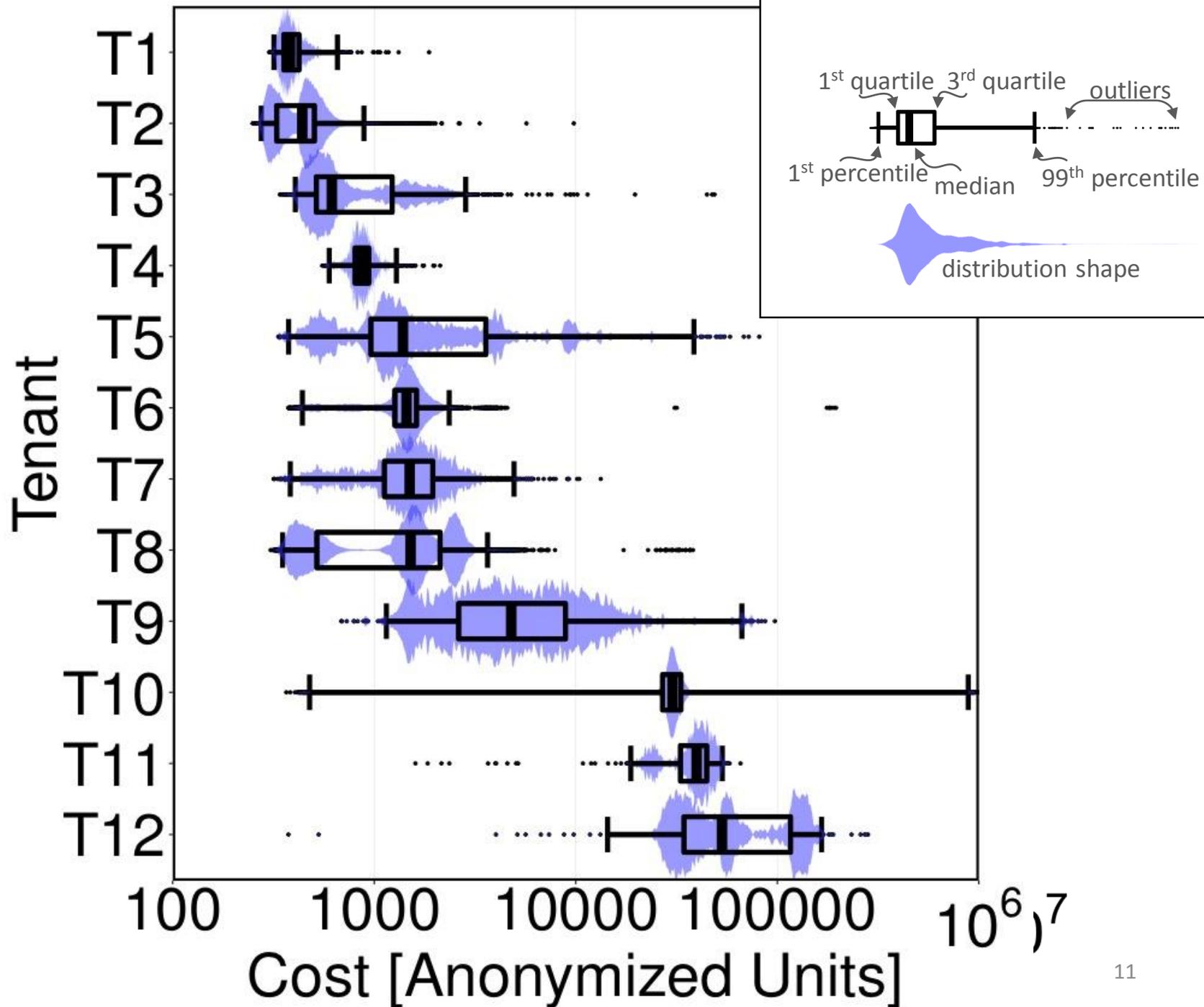
Size is used by scheduler to make scheduling decisions



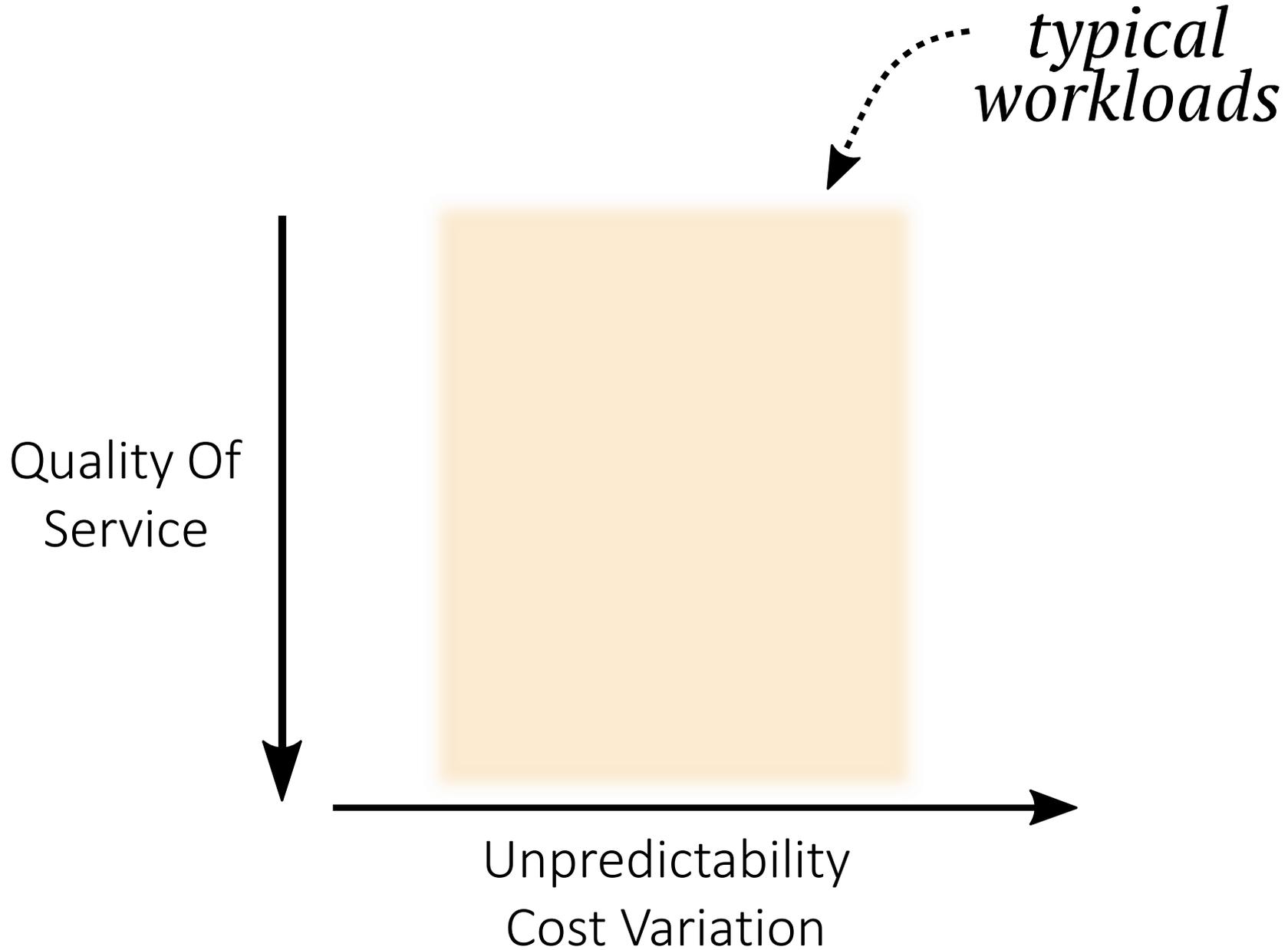
## Cloud services:

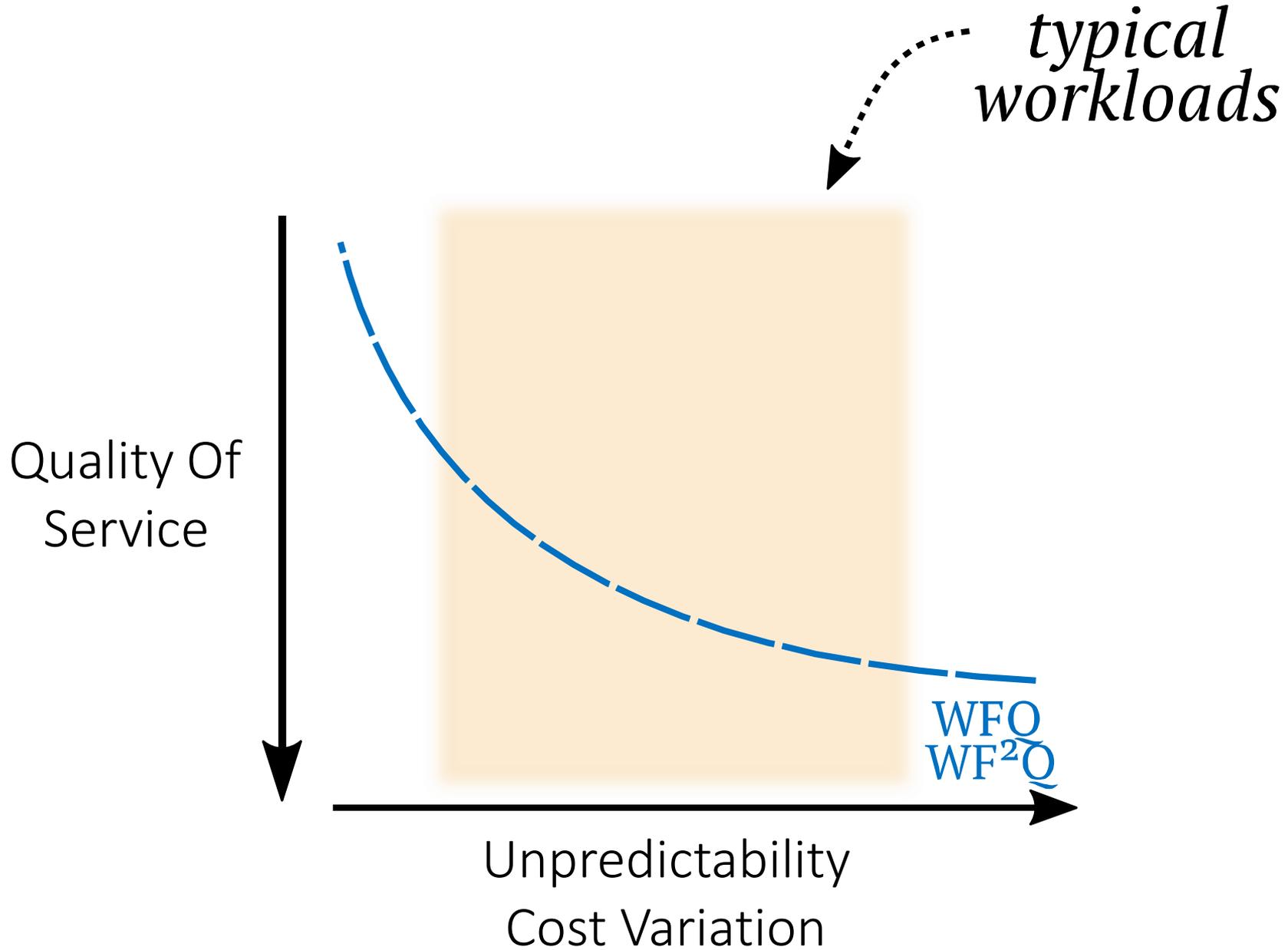
Estimation using model or moving averages <sub>10</sub>

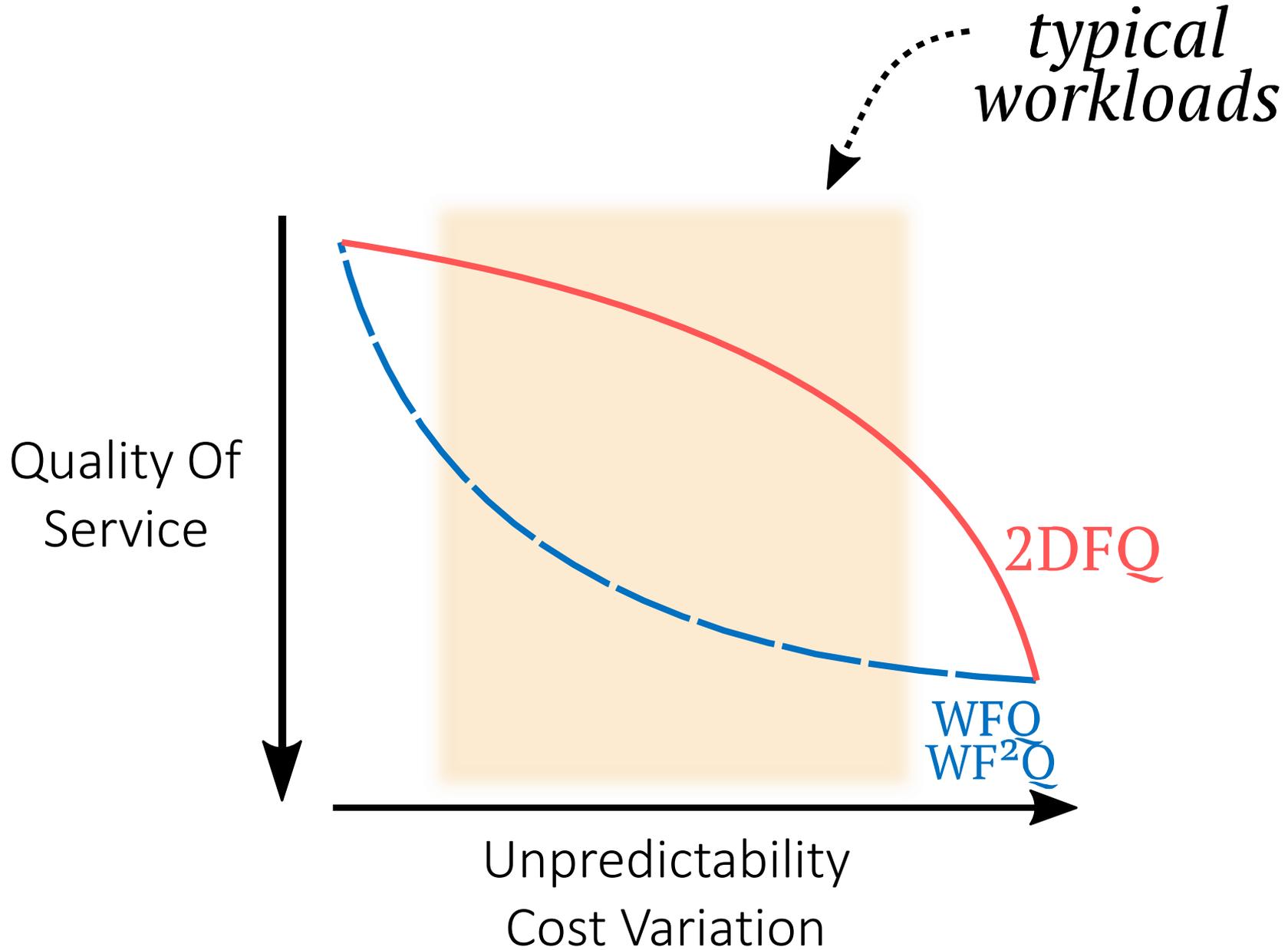






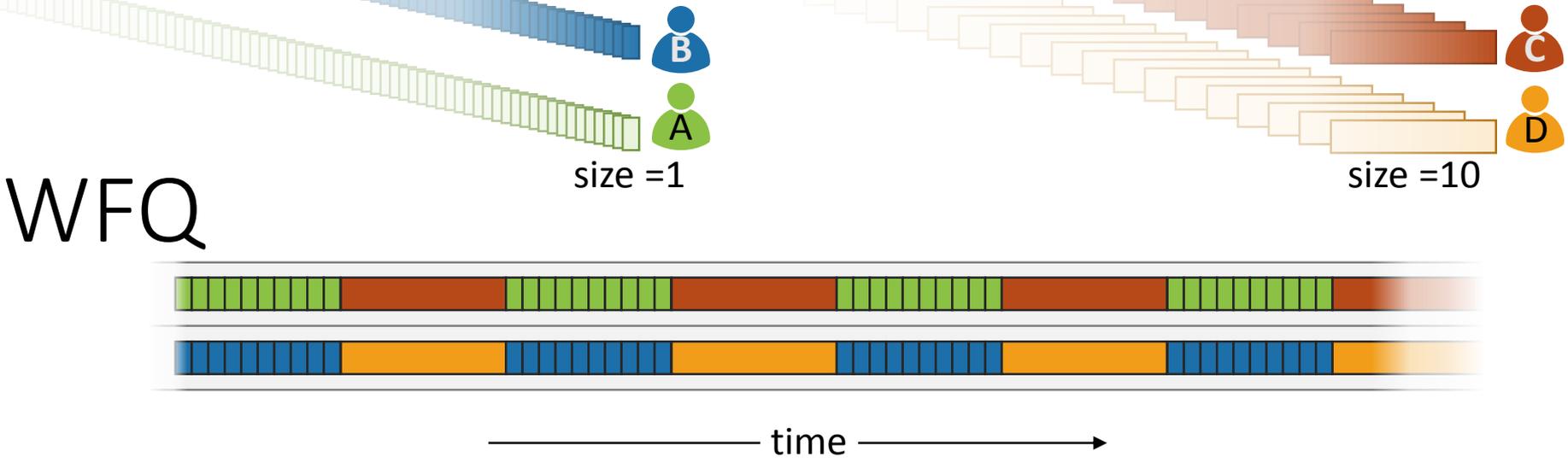




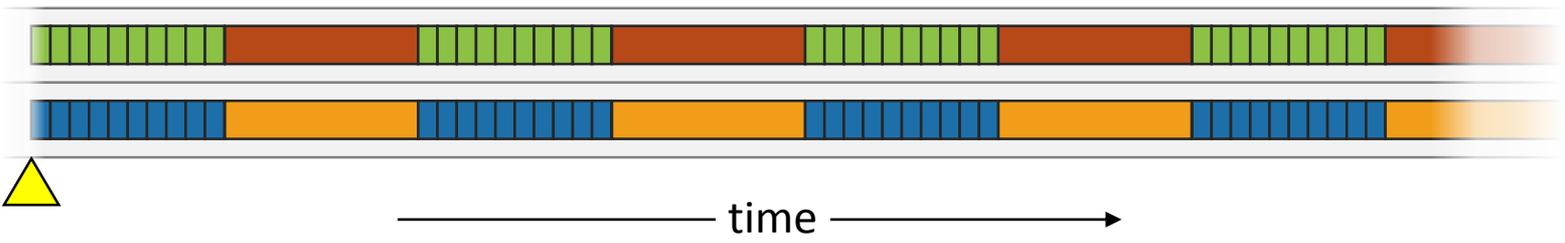


# Two-Dimensional Fair Queueing

# WFQ



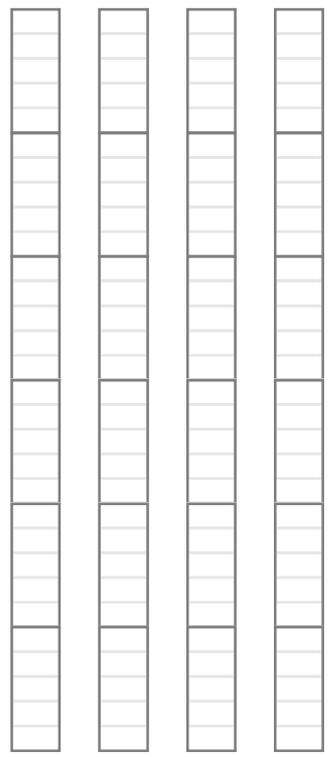
# WFQ



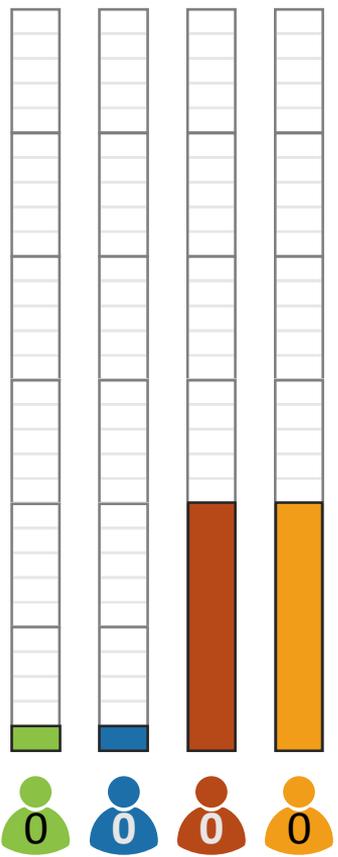
# WFQ



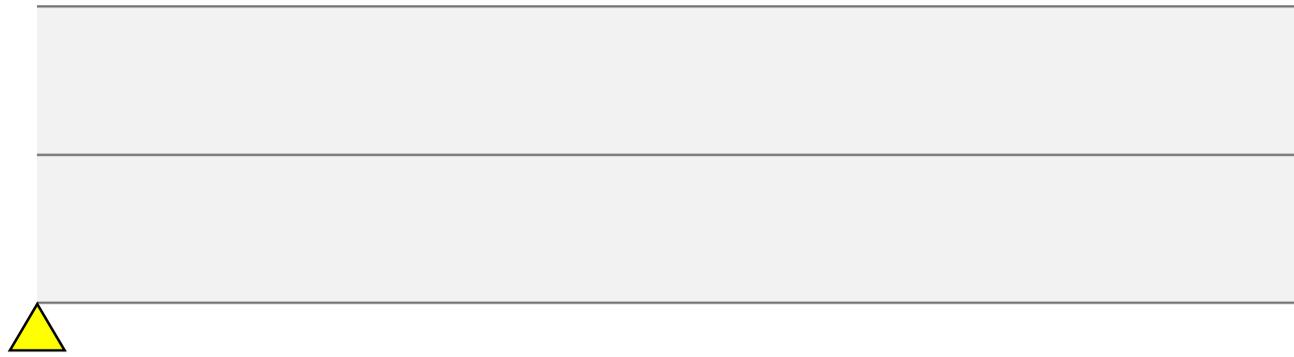
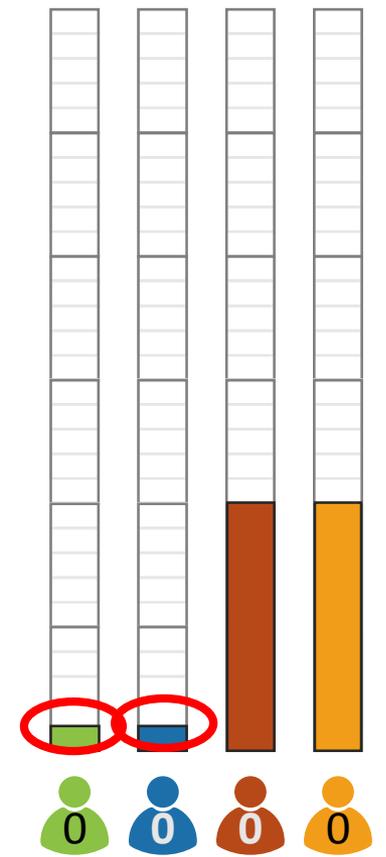
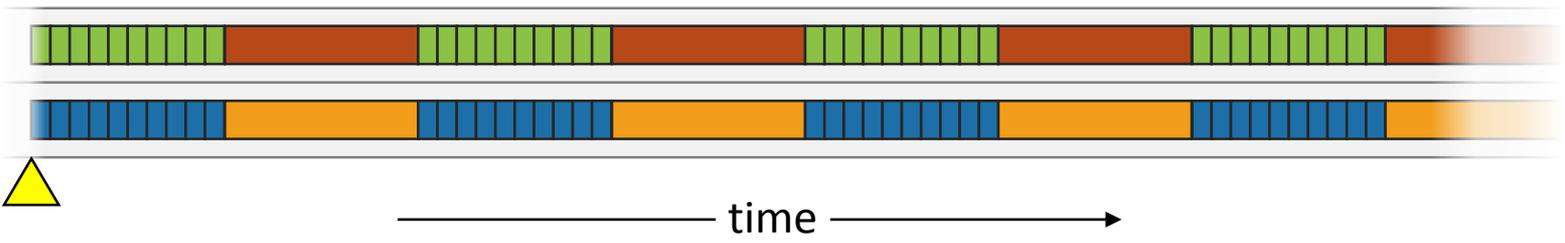
time →



# WFQ



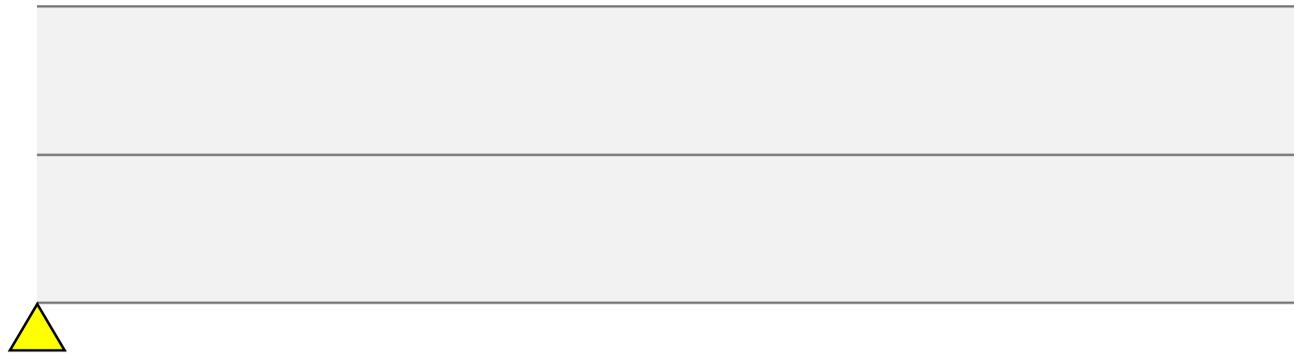
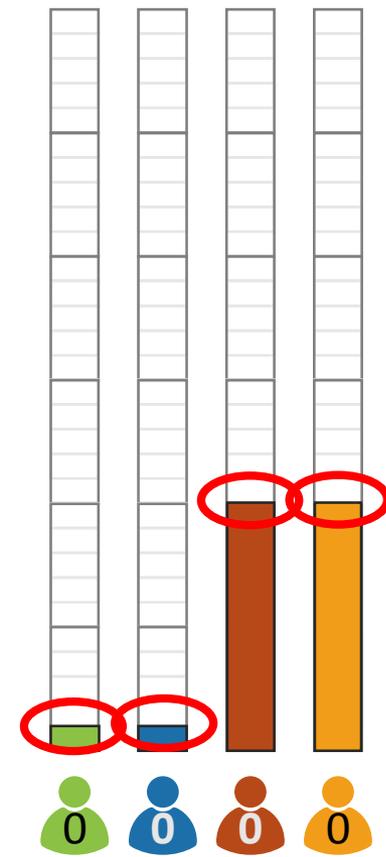
# WFQ



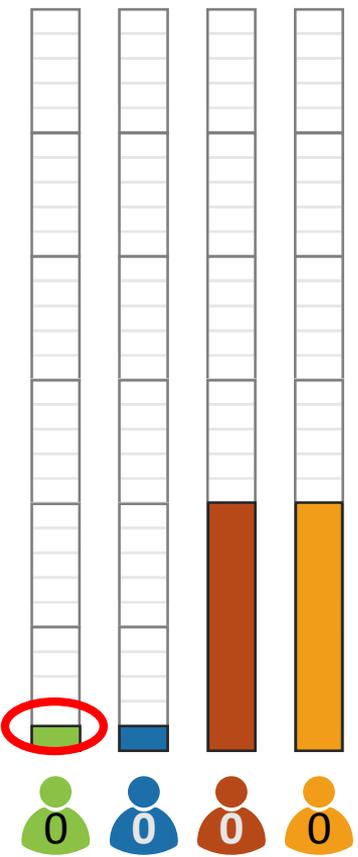
# WFQ



time →



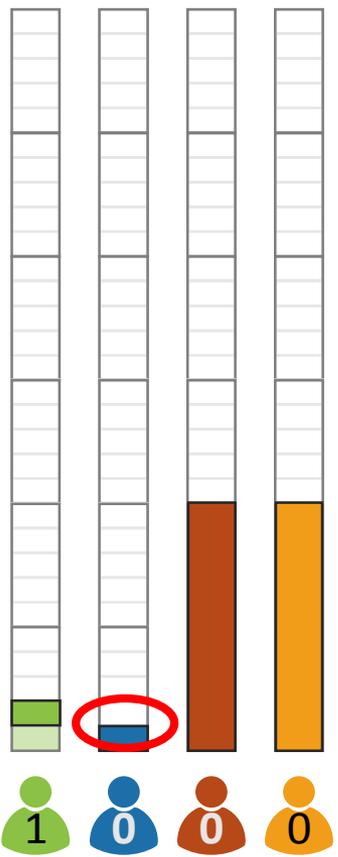
# WFQ



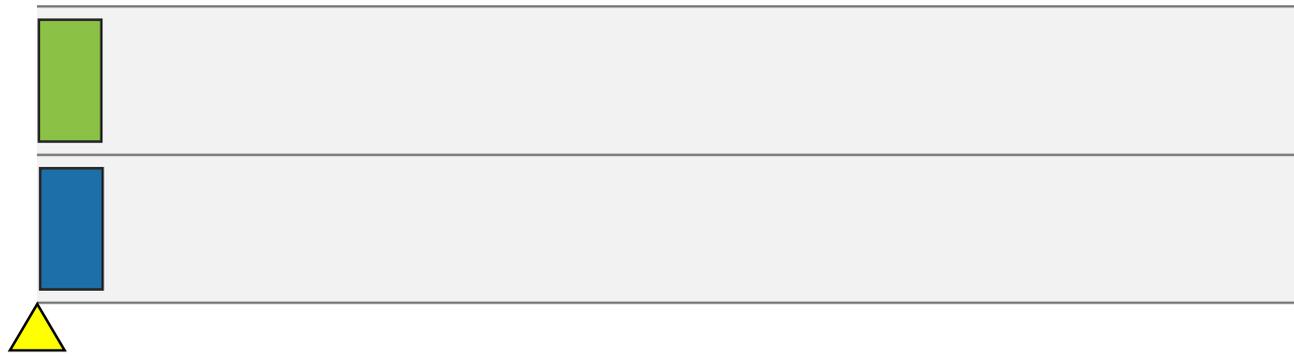
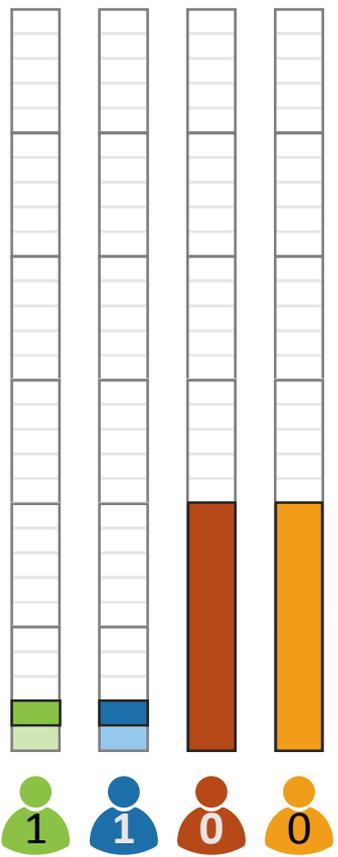
# WFQ



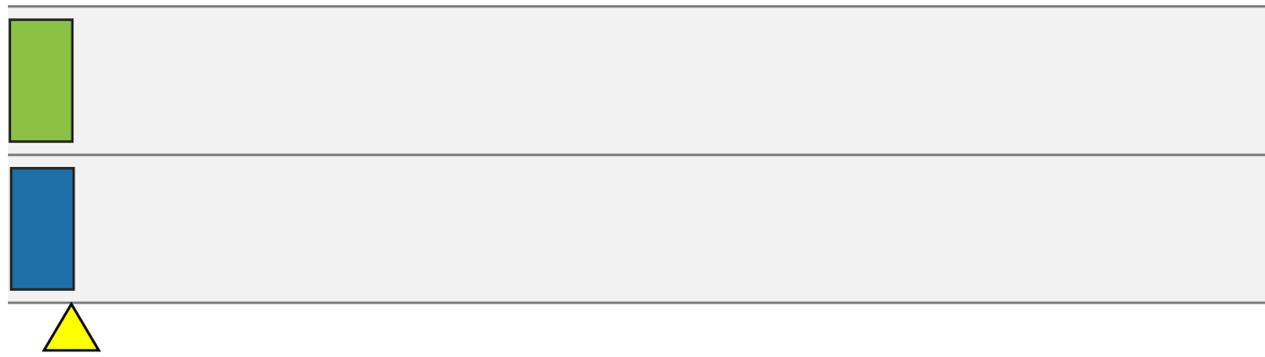
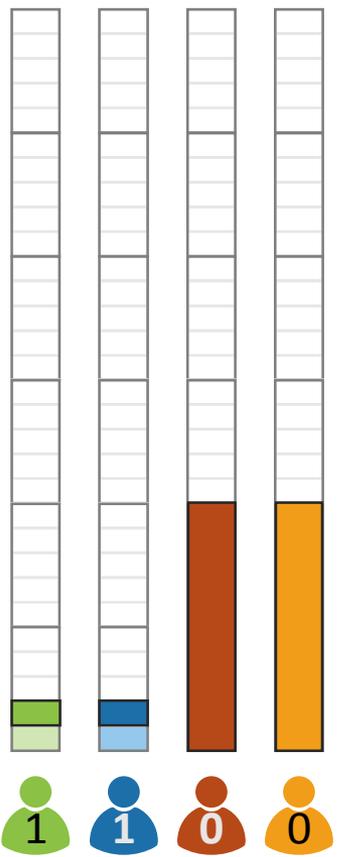
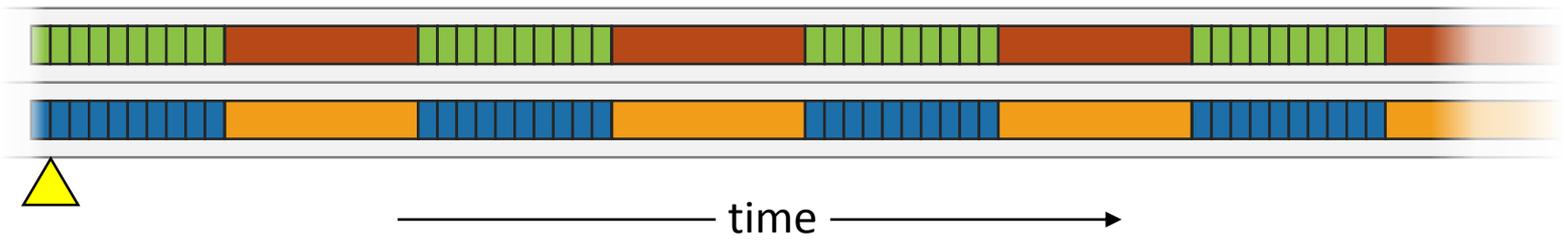
time →



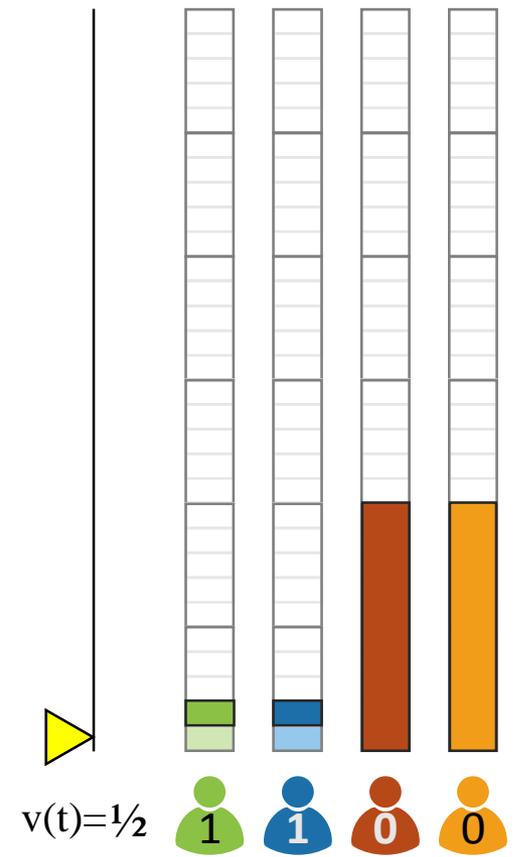
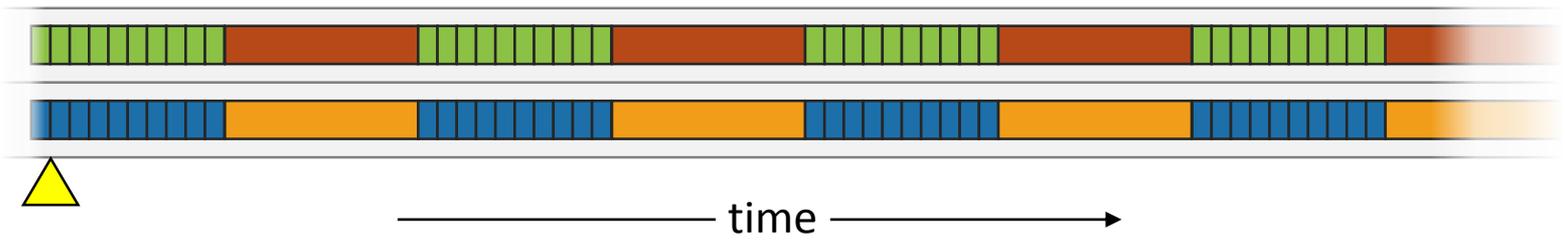
# WFQ



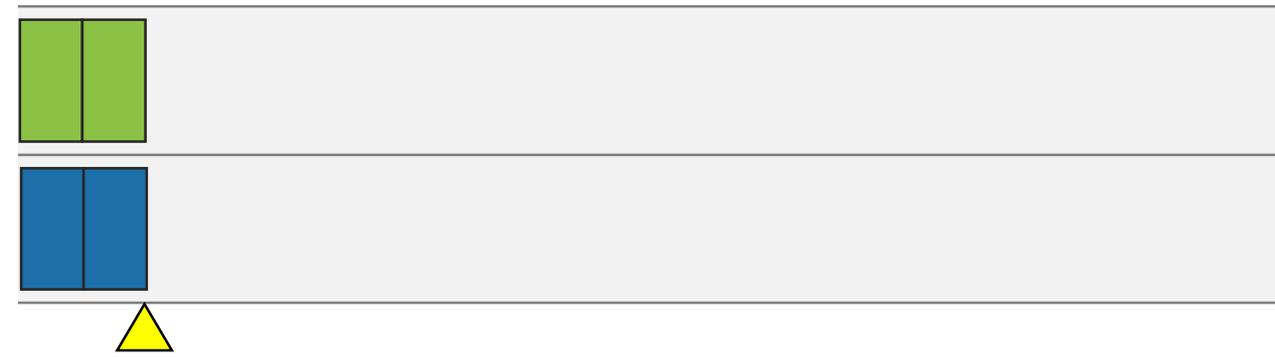
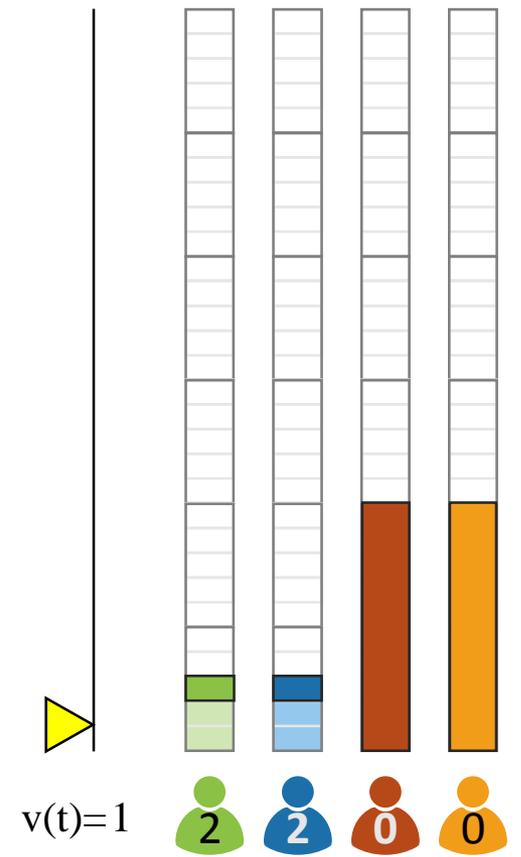
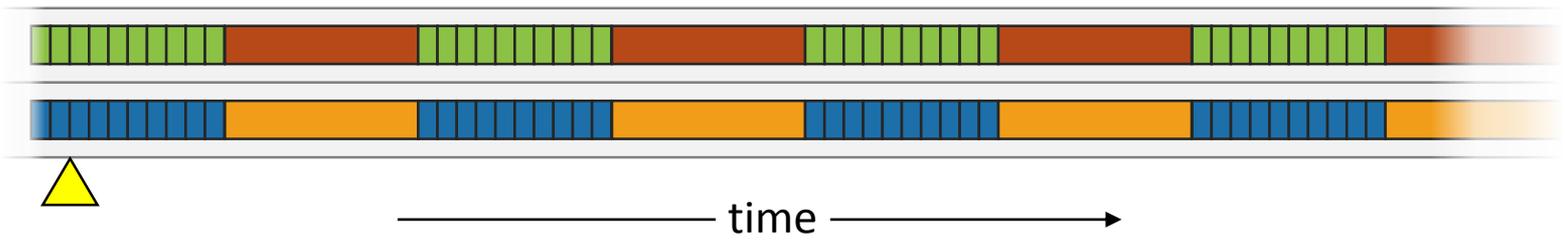
# WFQ



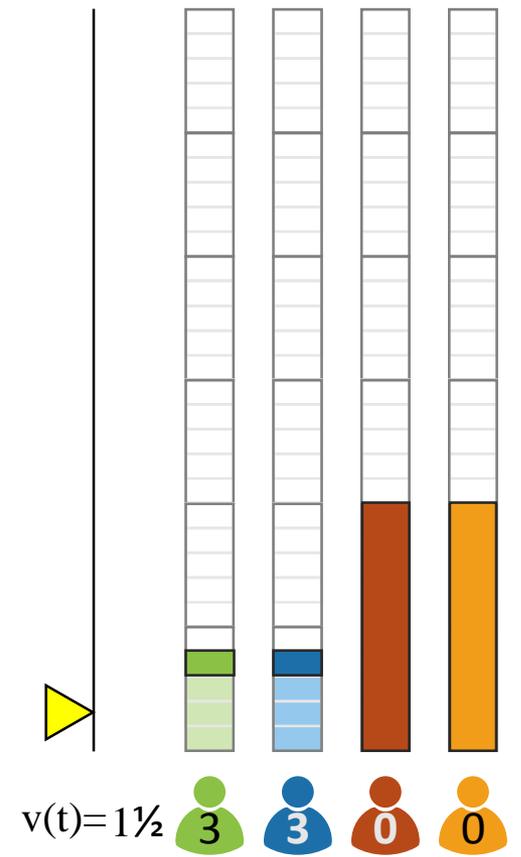
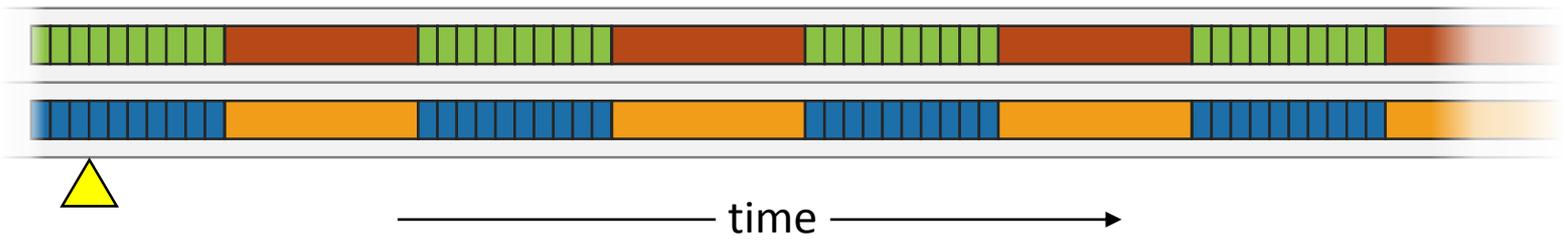
# WFQ



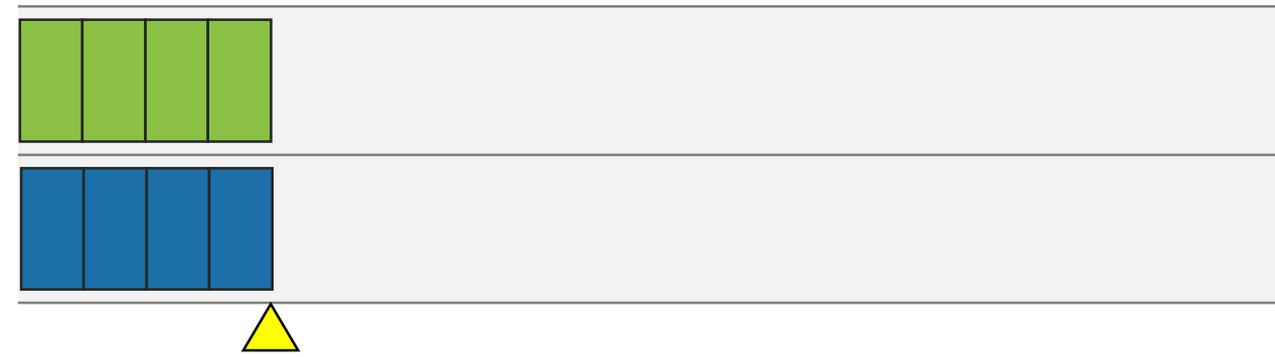
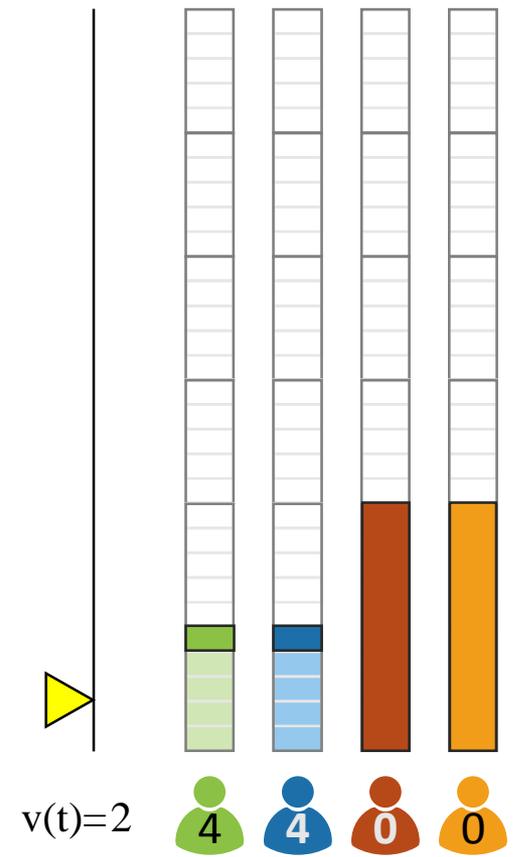
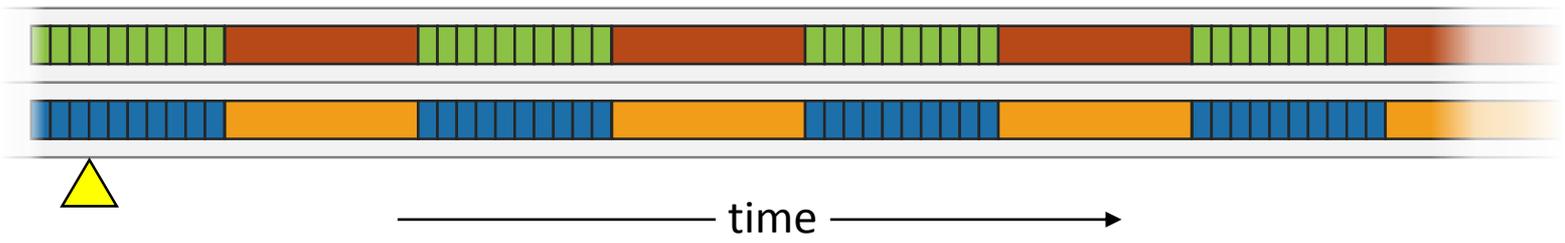
# WFQ



# WFQ



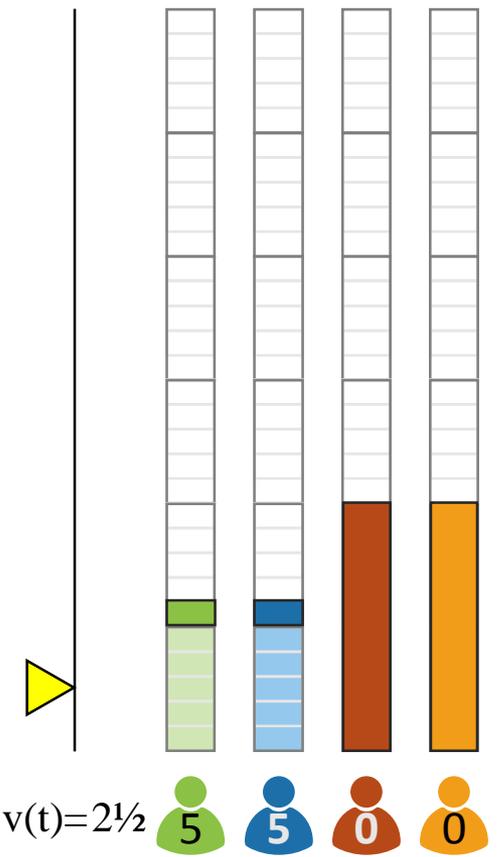
# WFQ



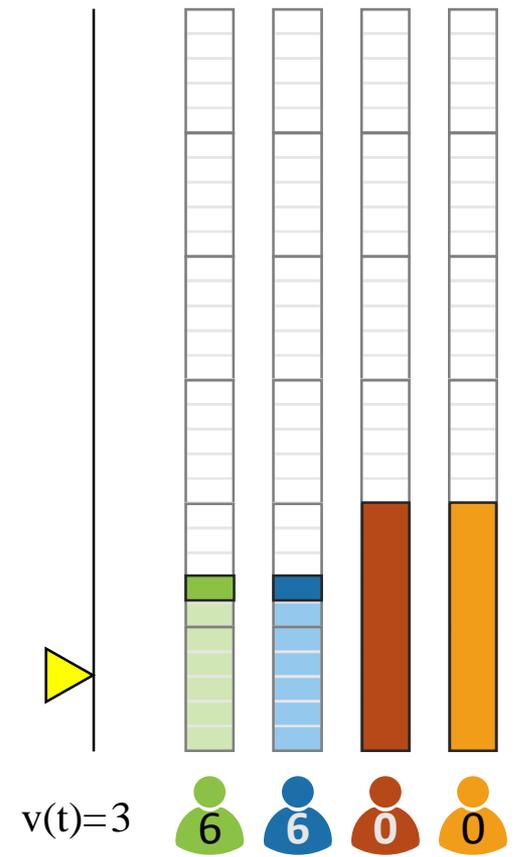
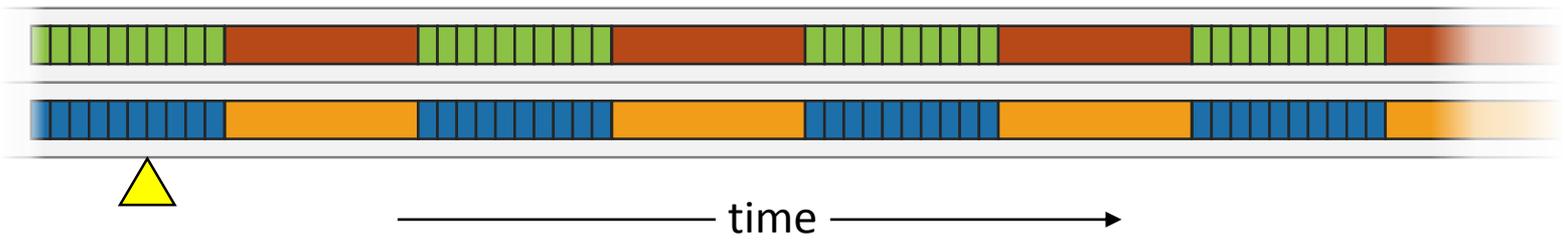
# WFQ



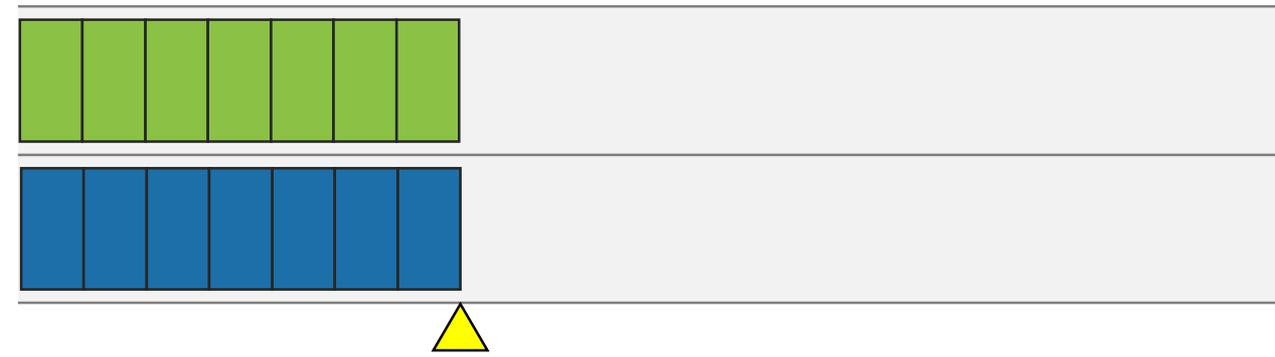
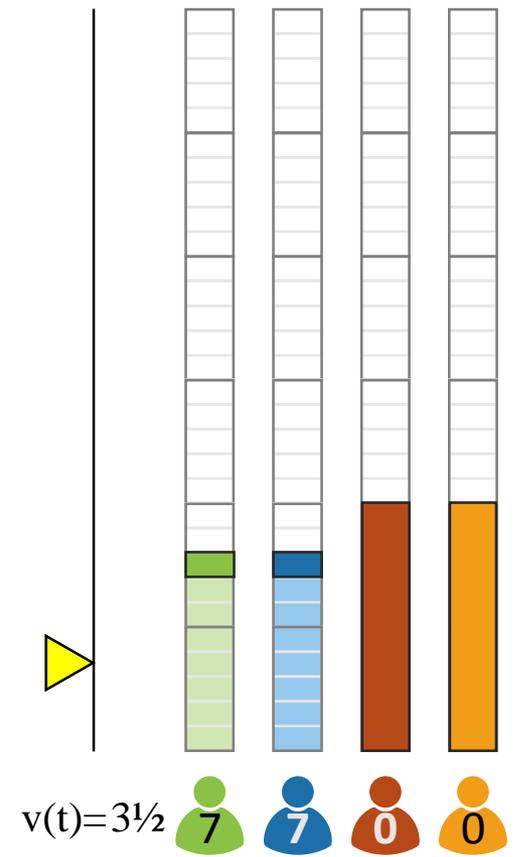
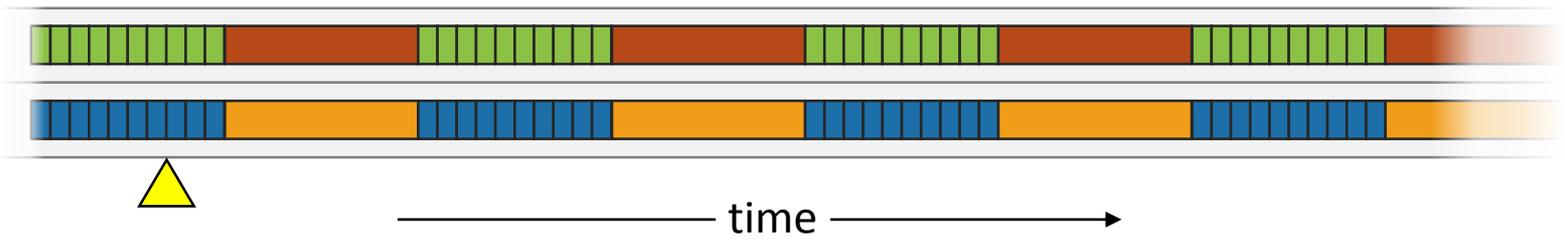
time →



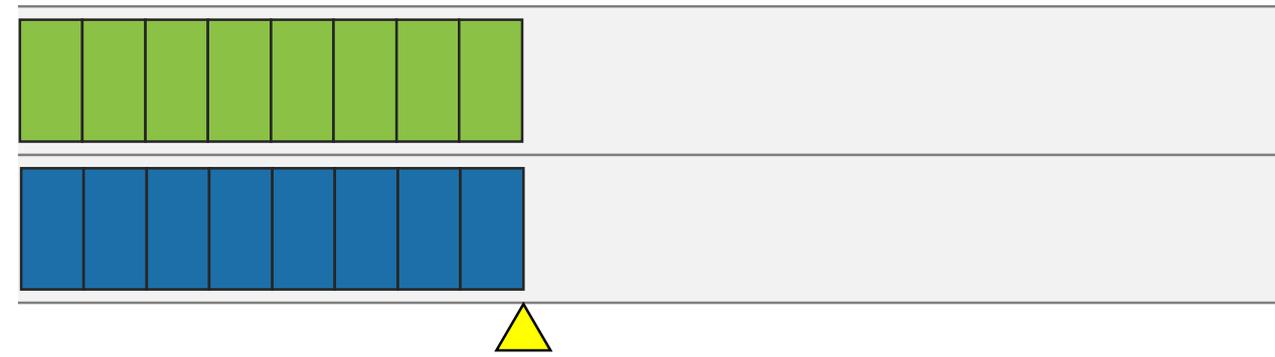
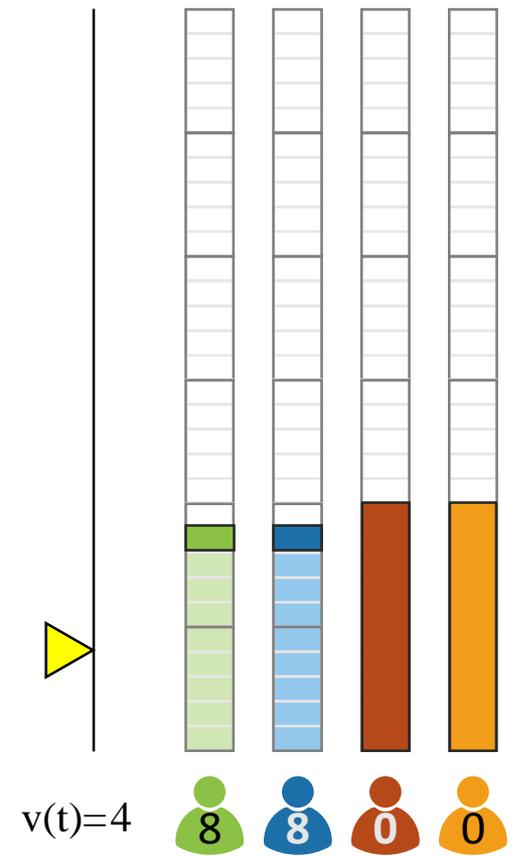
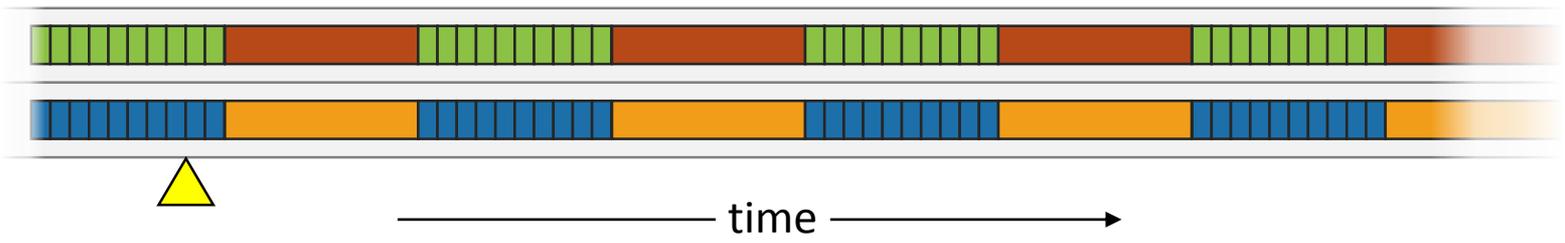
# WFQ



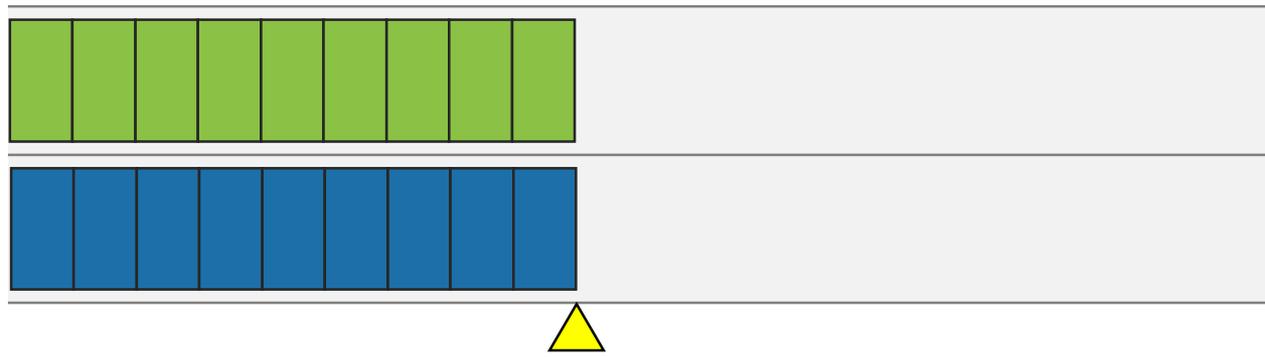
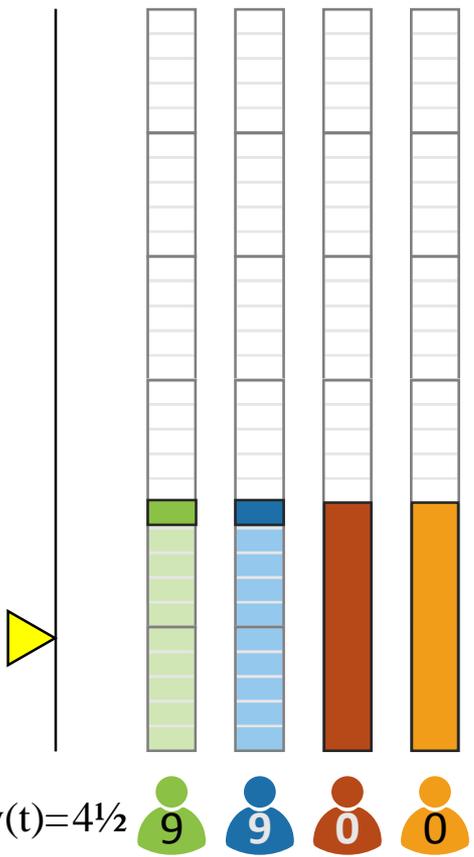
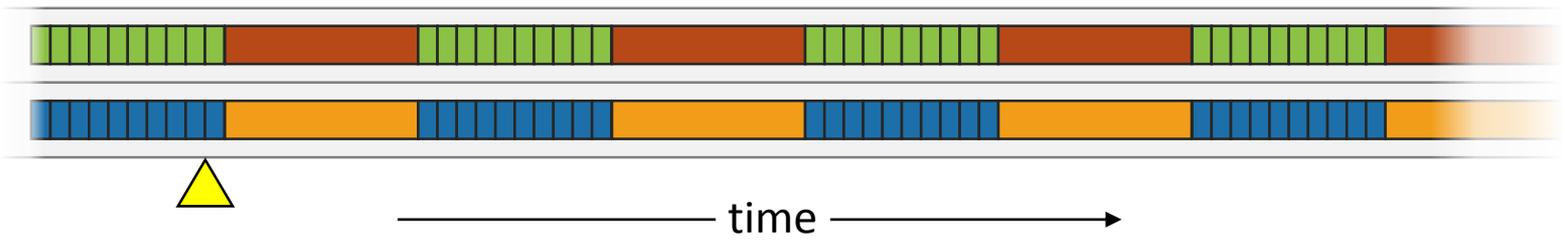
# WFQ



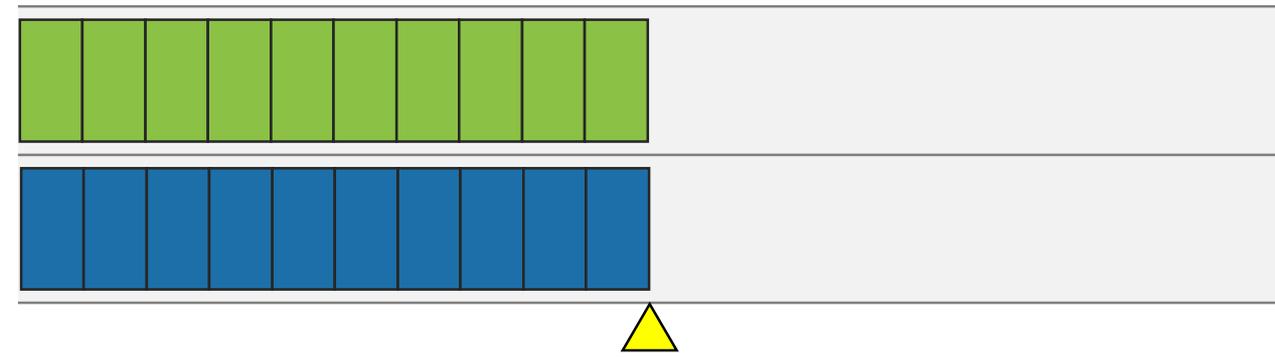
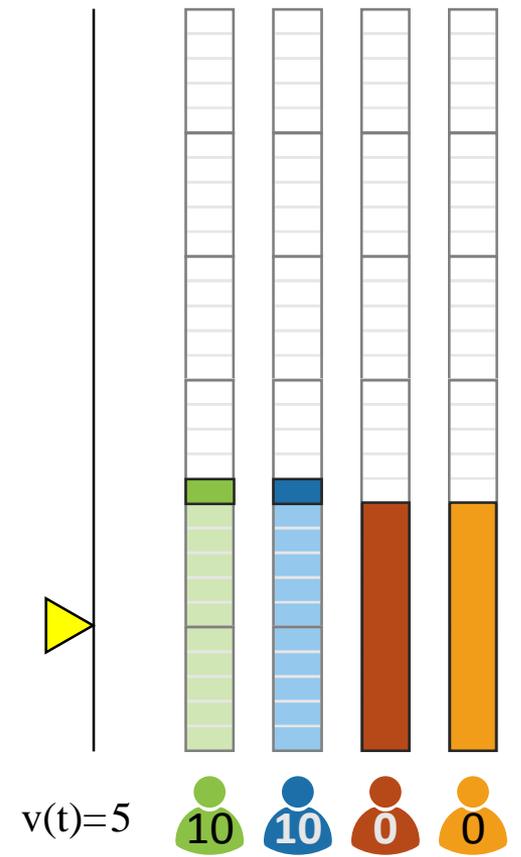
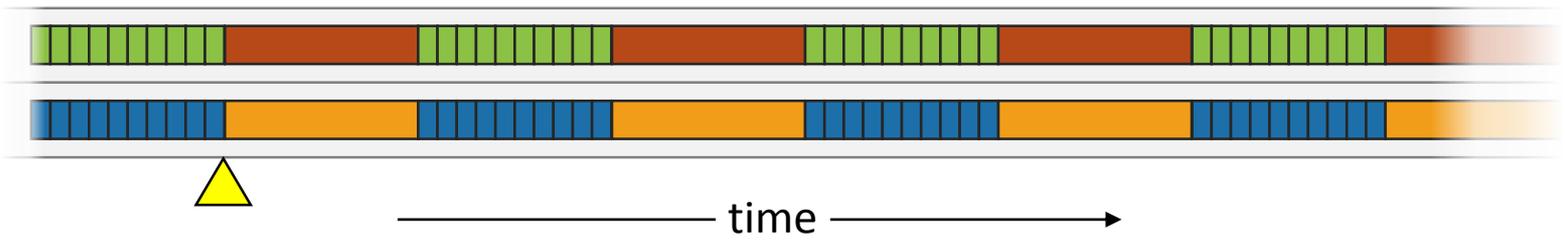
# WFQ



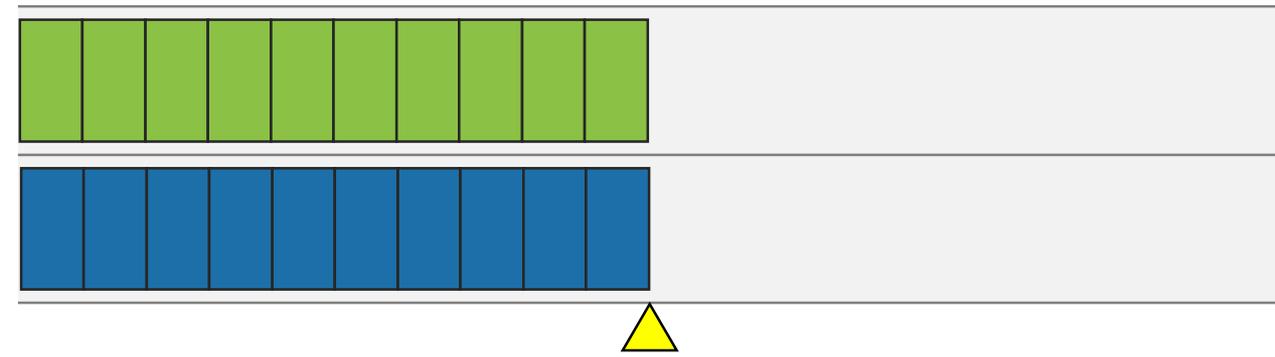
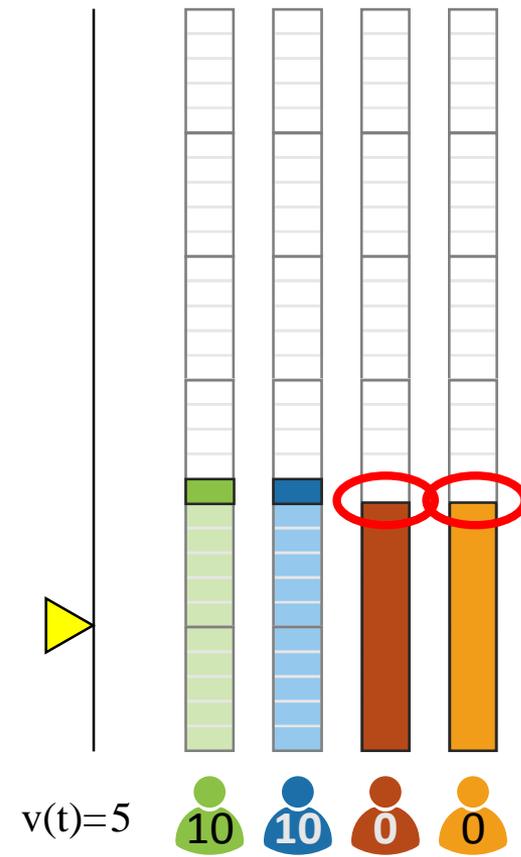
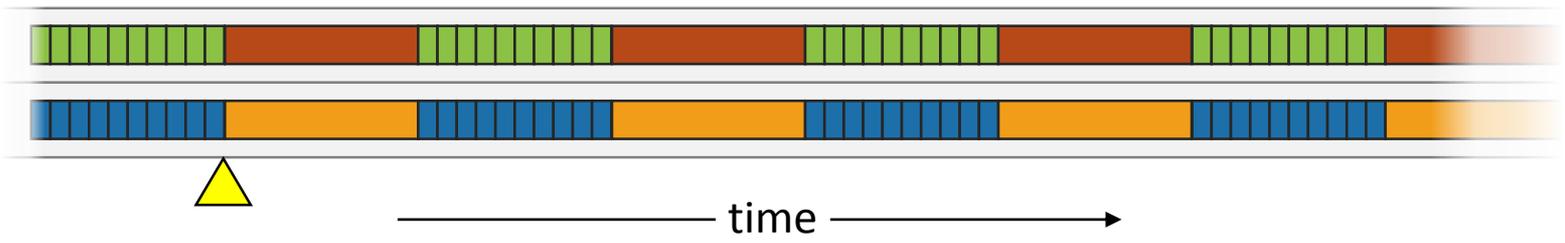
# WFQ



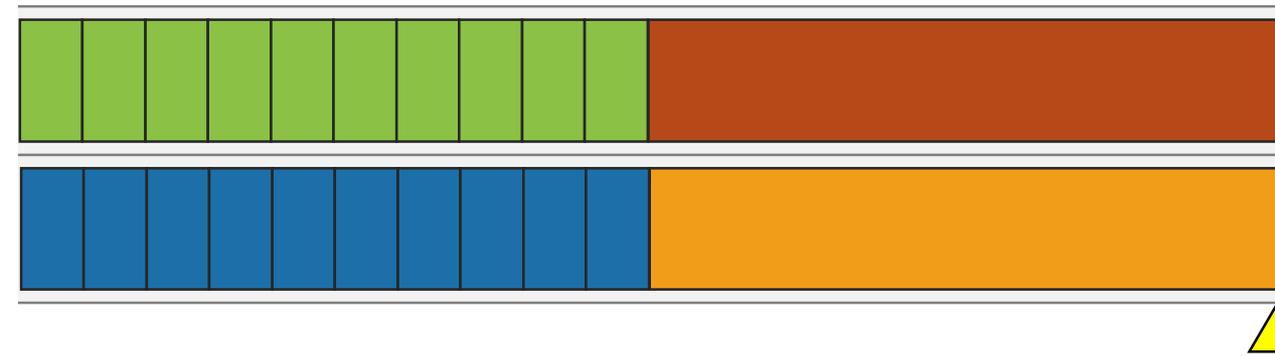
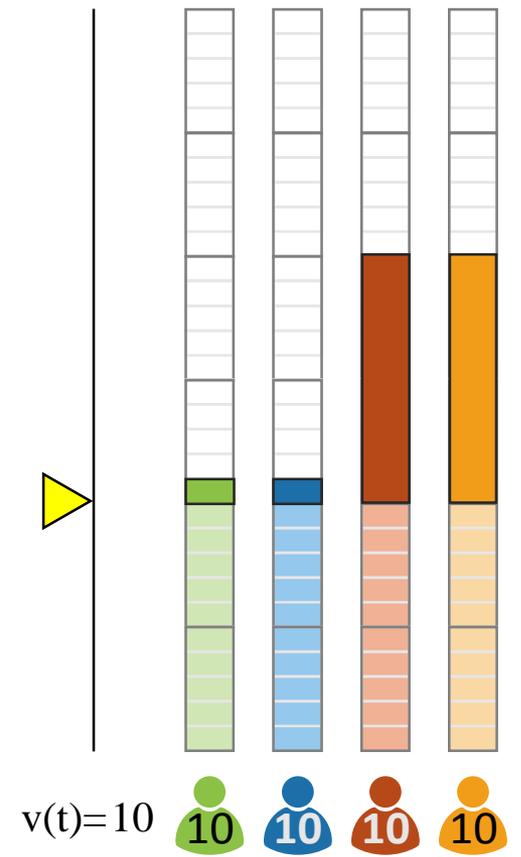
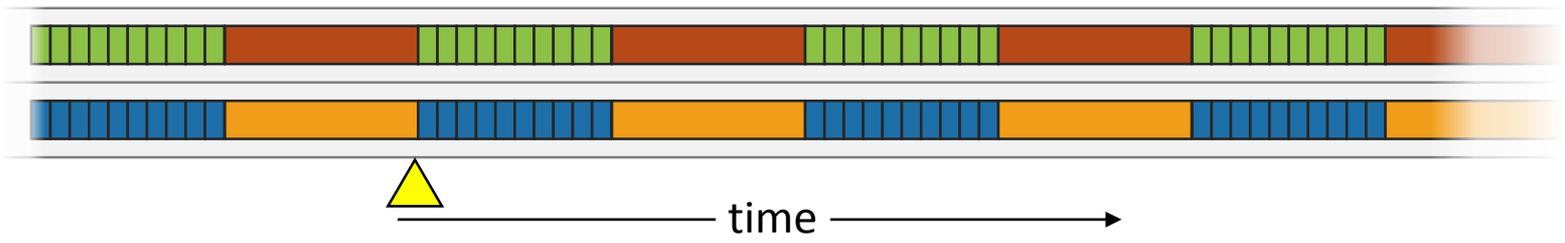
# WFQ



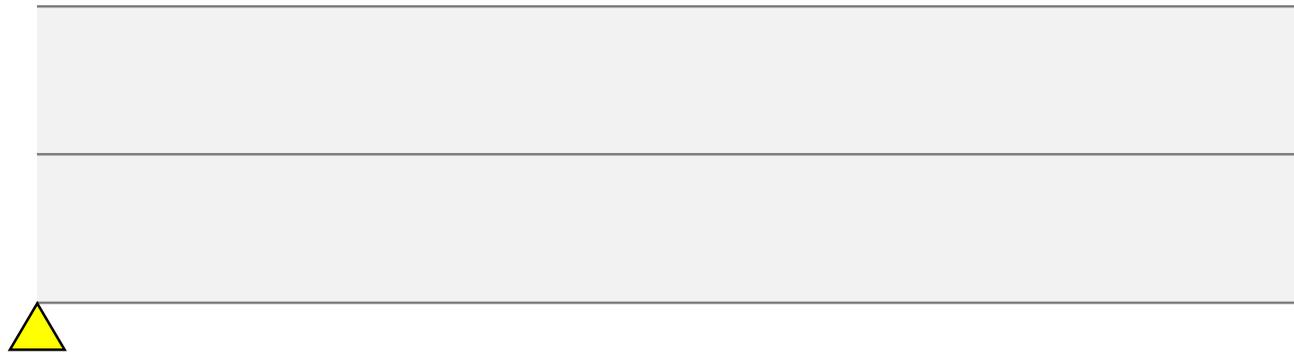
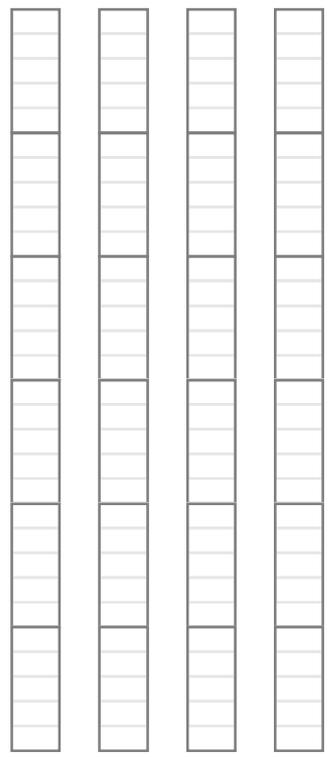
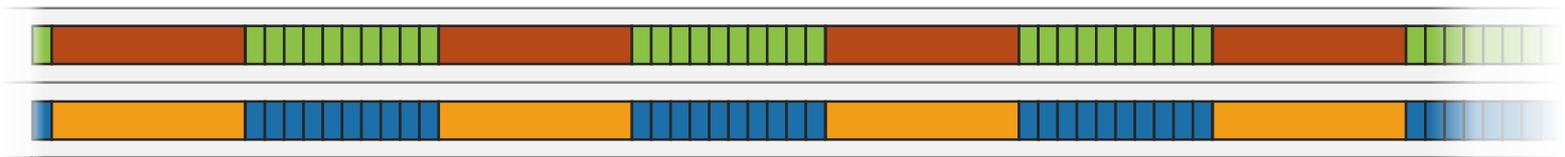
# WFQ



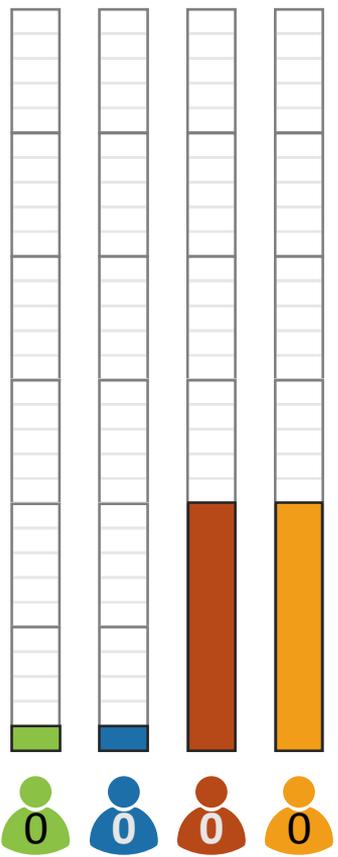
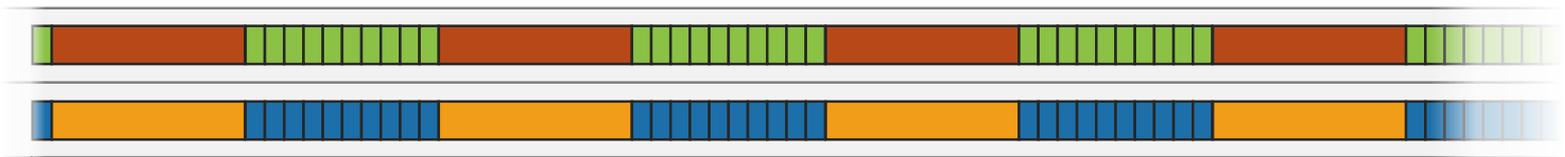
# WFQ



# WF<sup>2</sup>Q

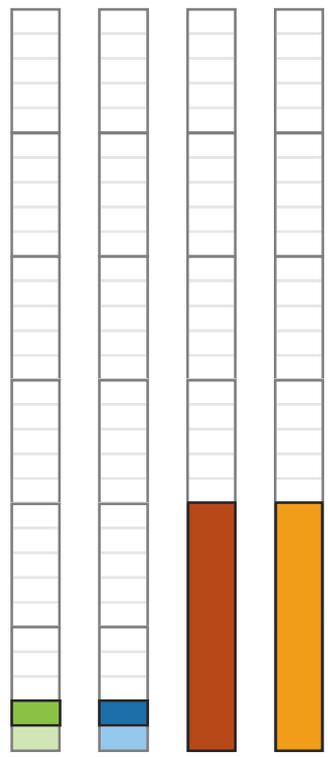
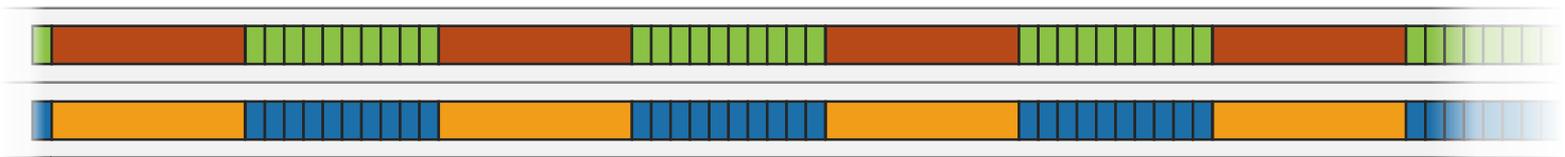


# WF<sup>2</sup>Q

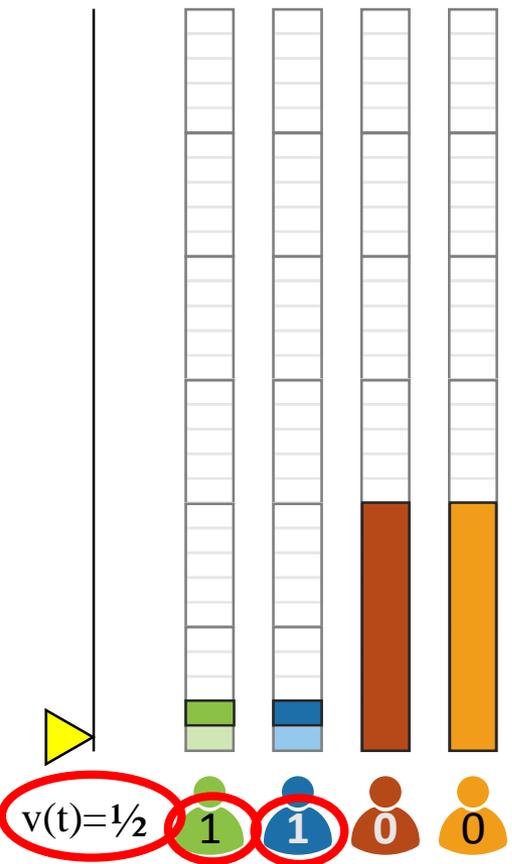
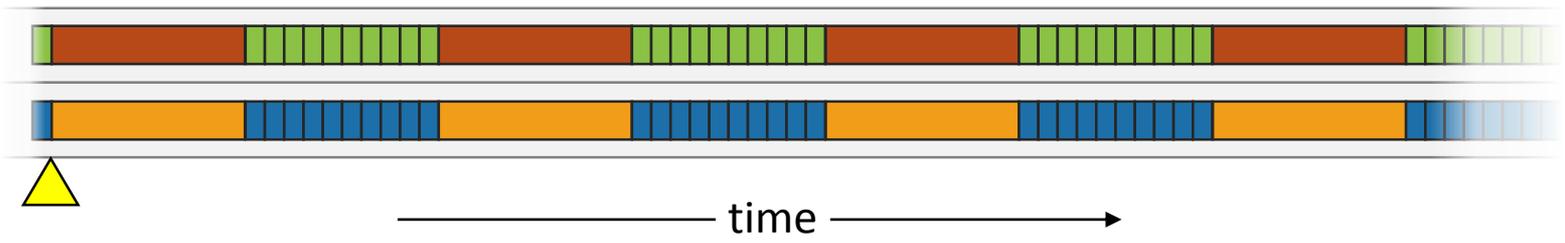


$v(t)=0$

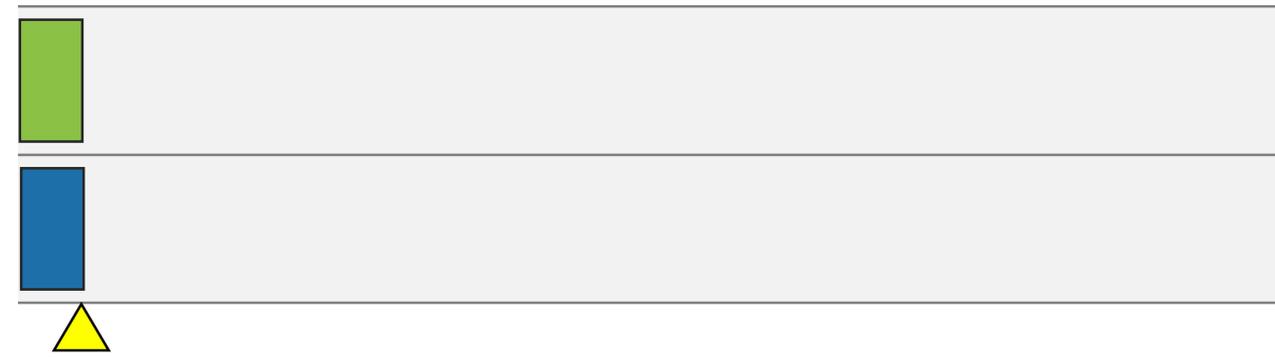
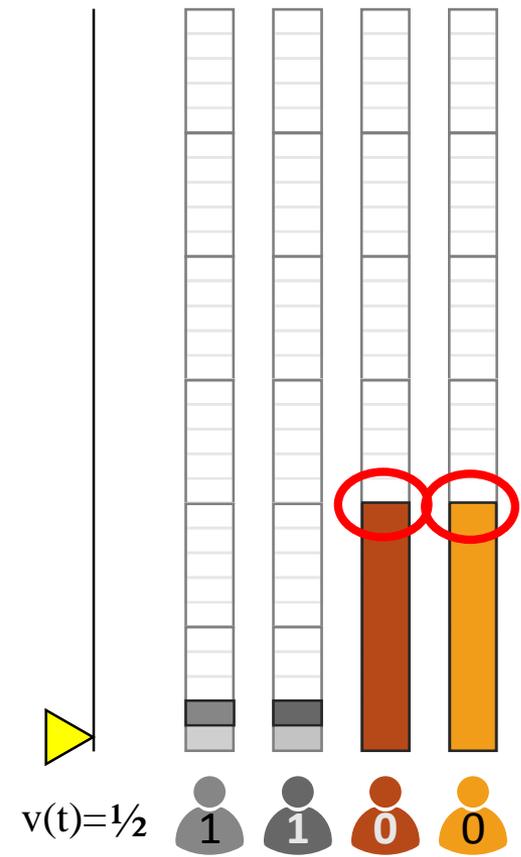
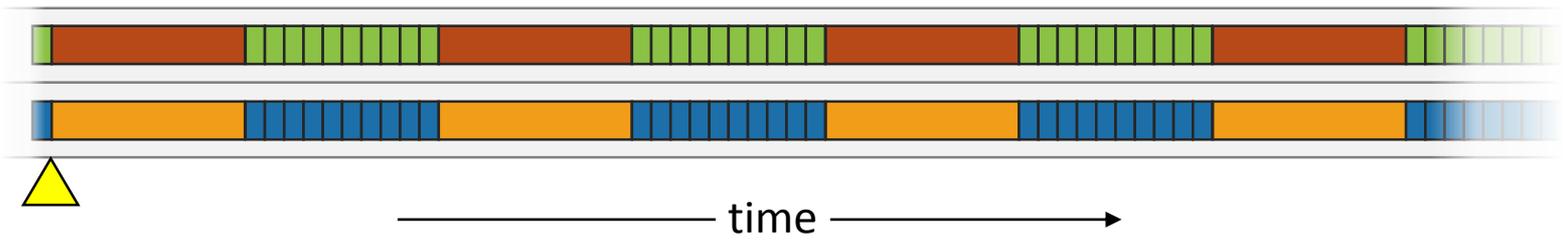
# WF<sup>2</sup>Q



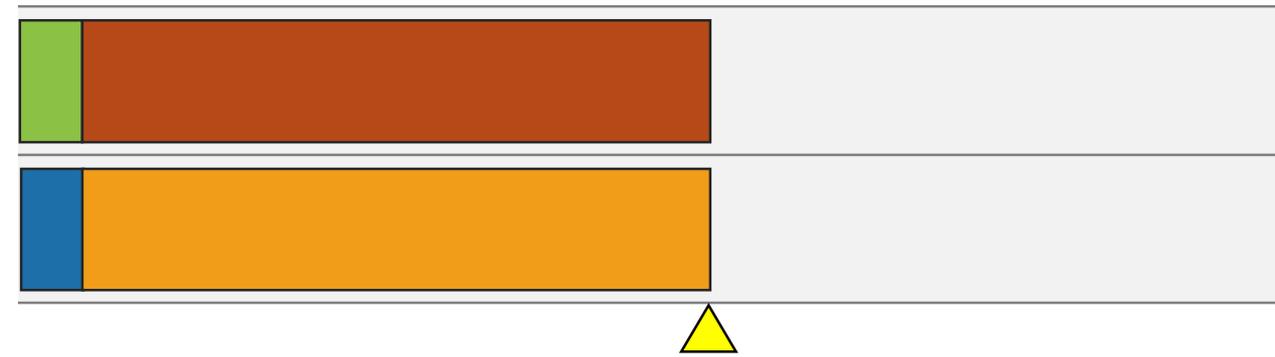
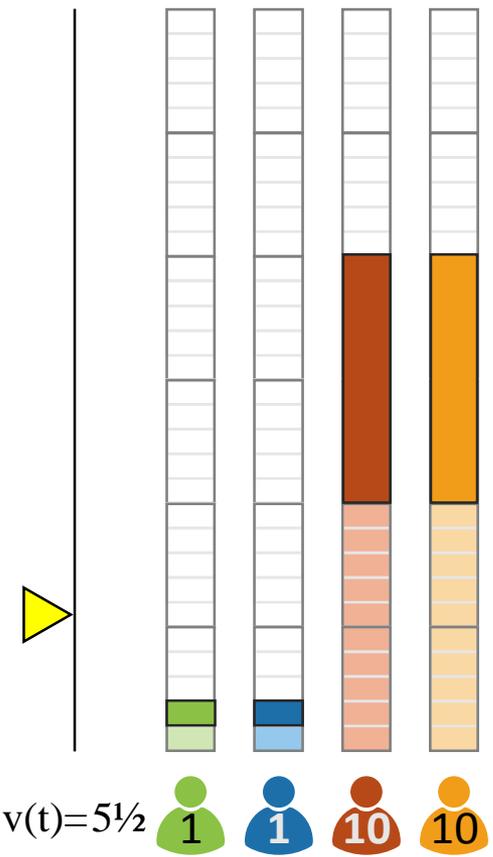
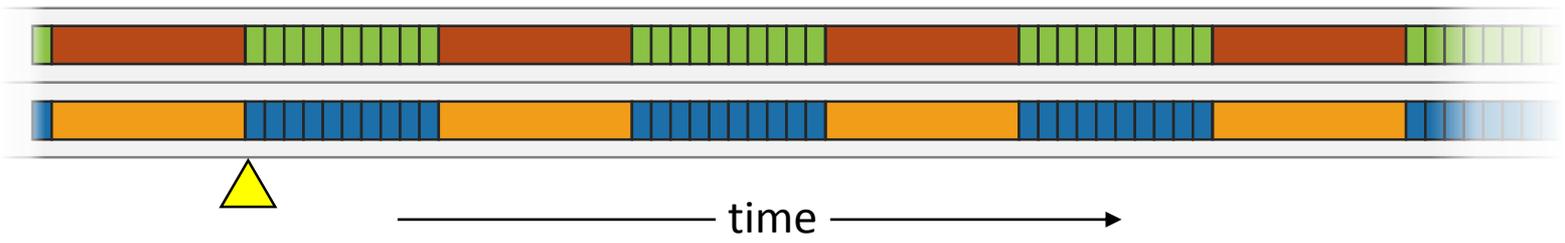
# WF<sup>2</sup>Q



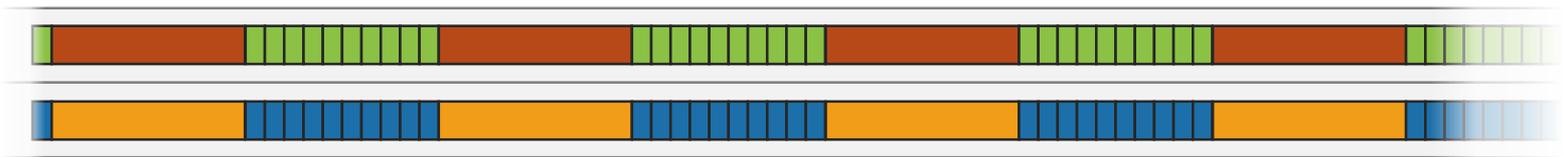
# WF<sup>2</sup>Q



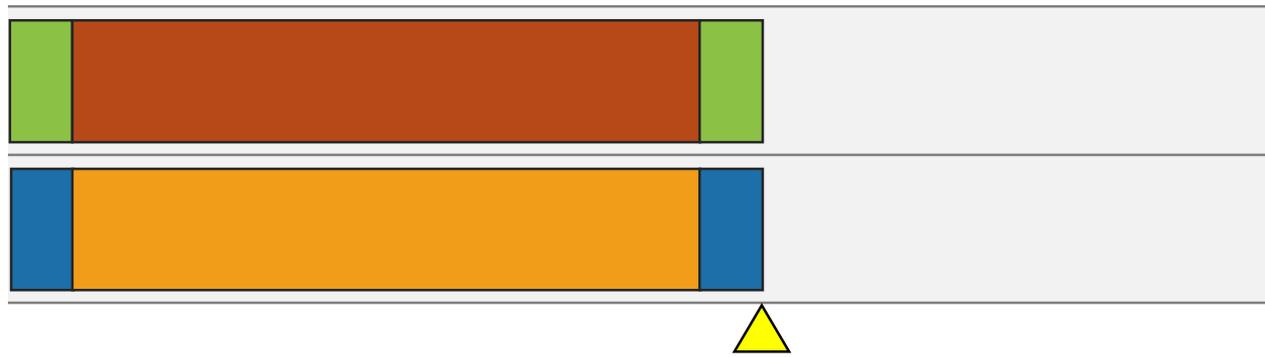
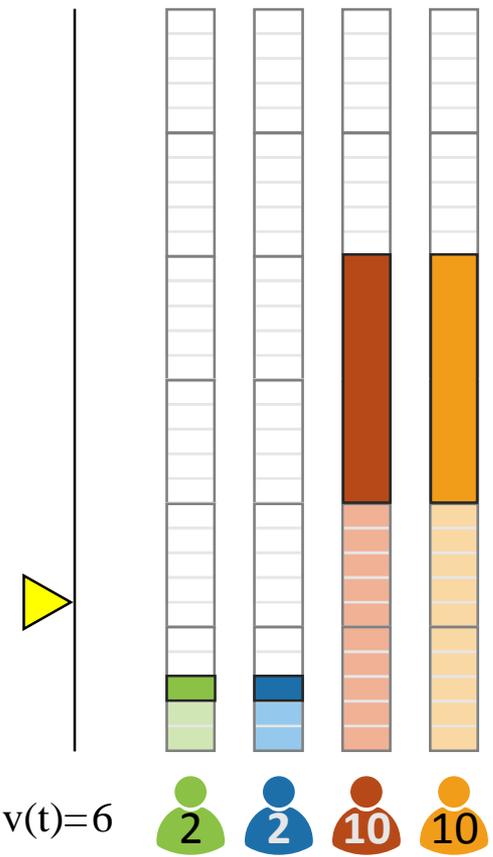
# WF<sup>2</sup>Q



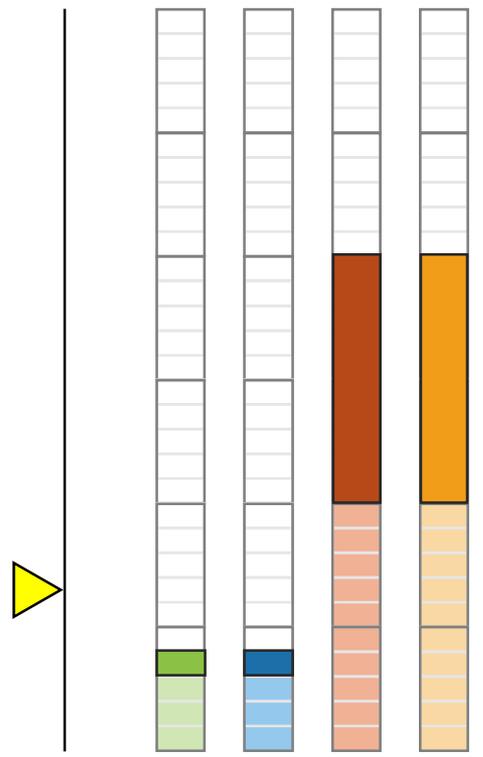
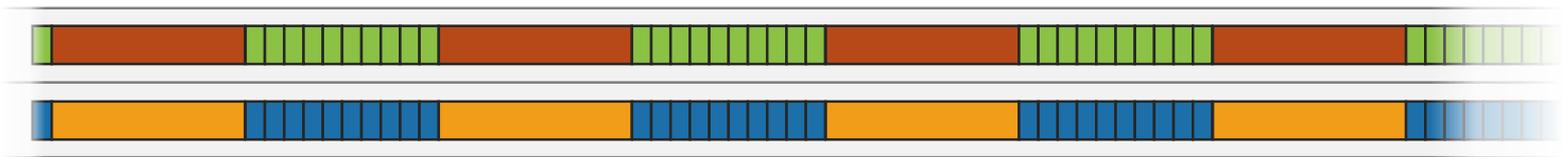
# WF<sup>2</sup>Q



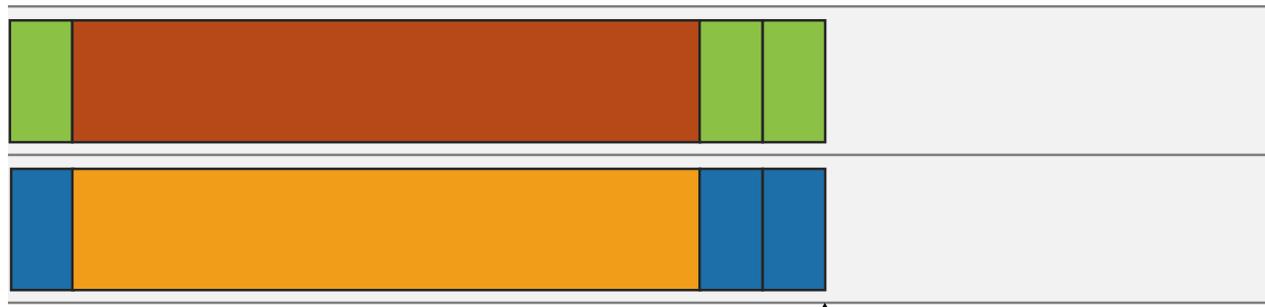
time →



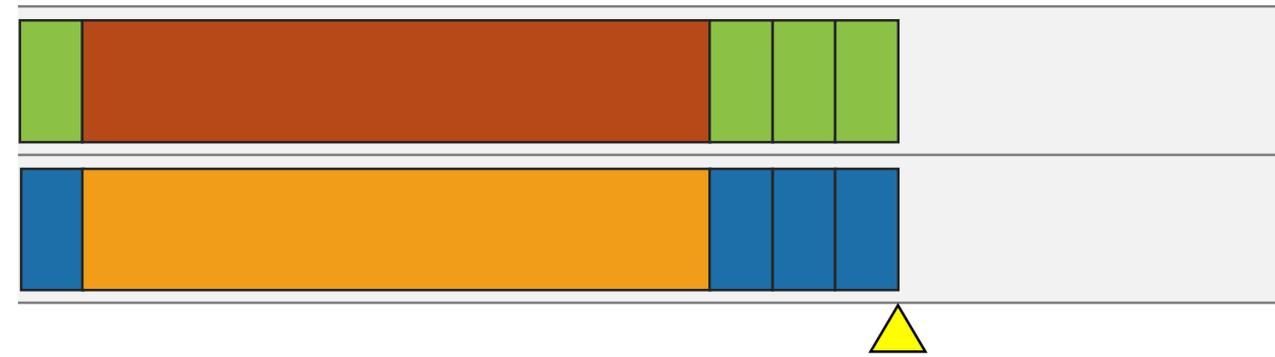
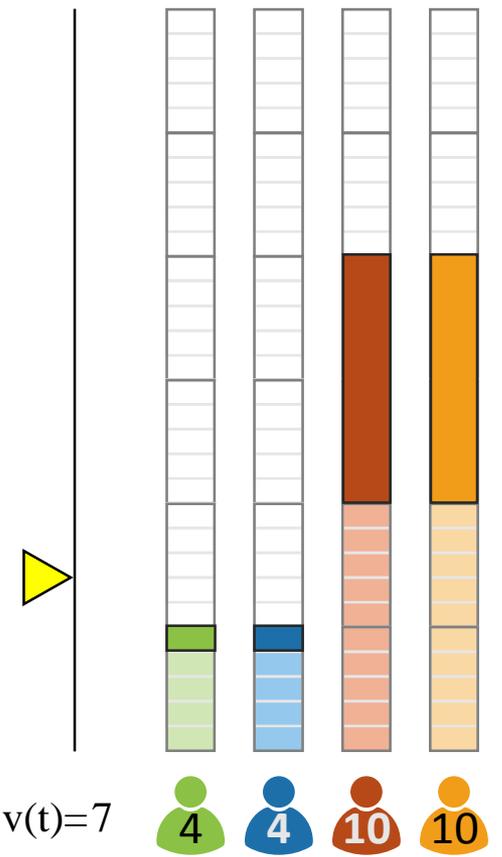
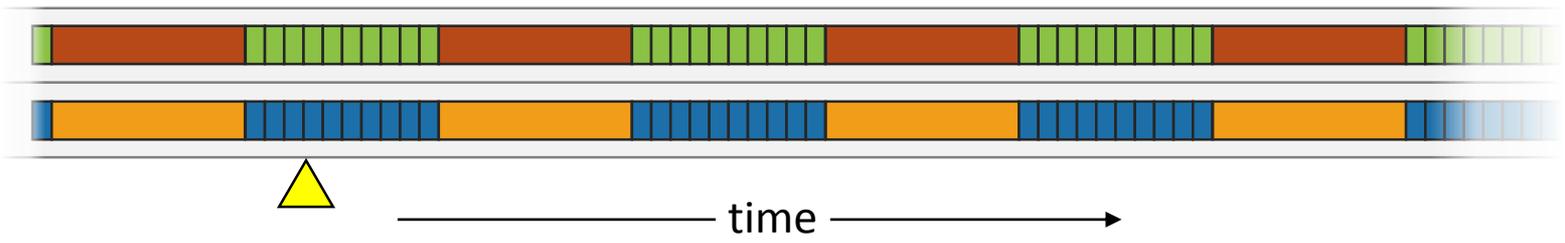
# WF<sup>2</sup>Q



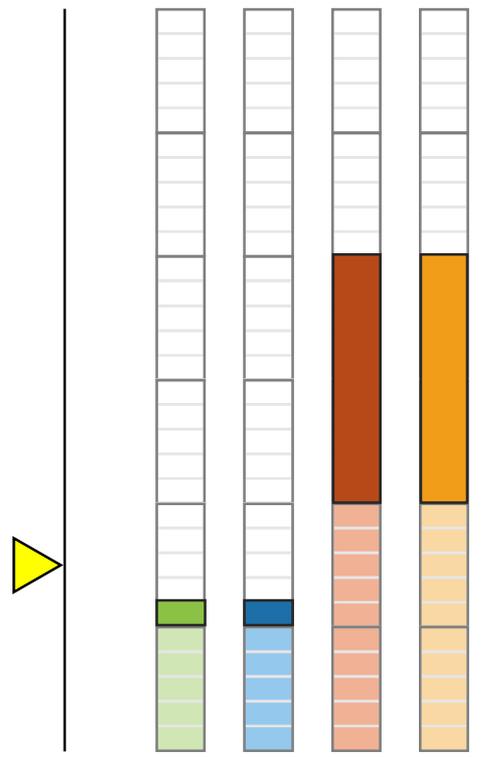
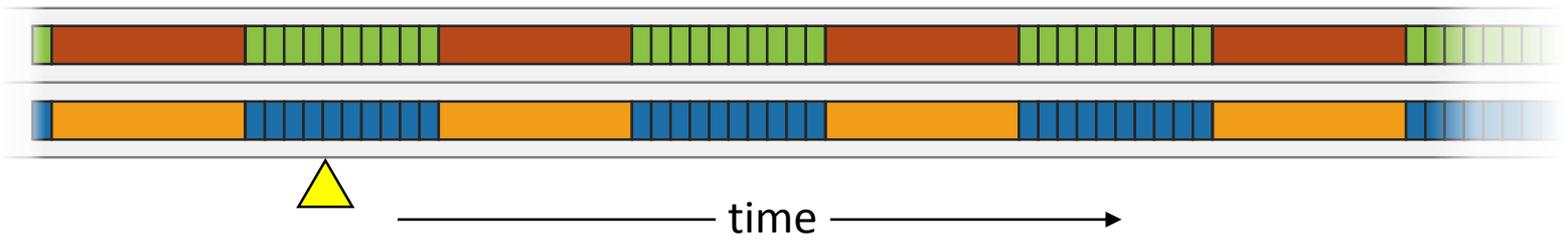
$v(t) = 6 \frac{1}{2}$     3    3    10    10



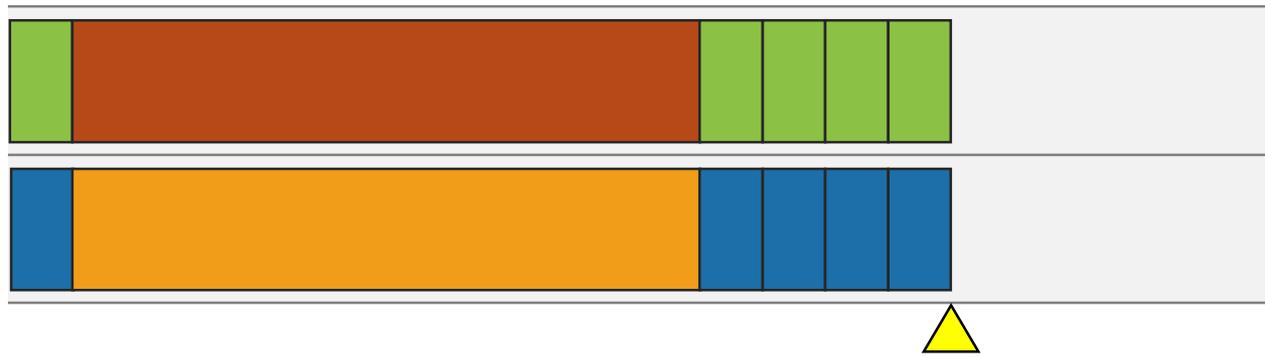
# WF<sup>2</sup>Q



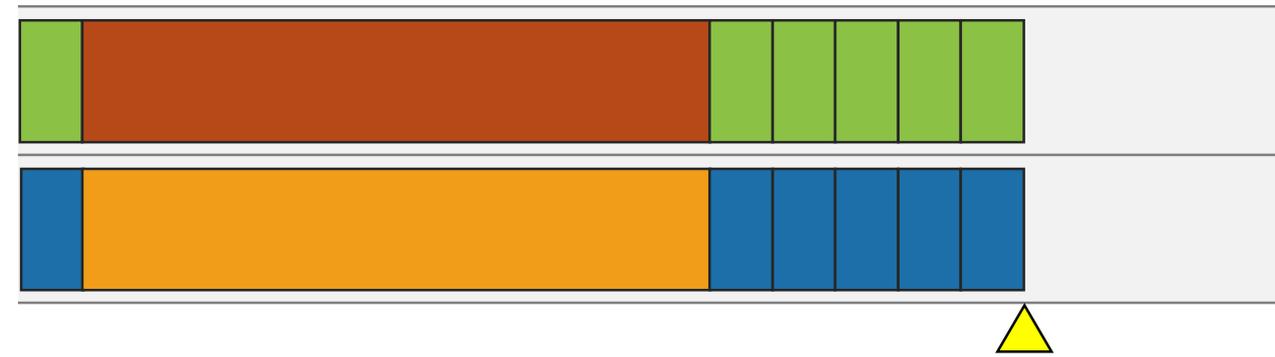
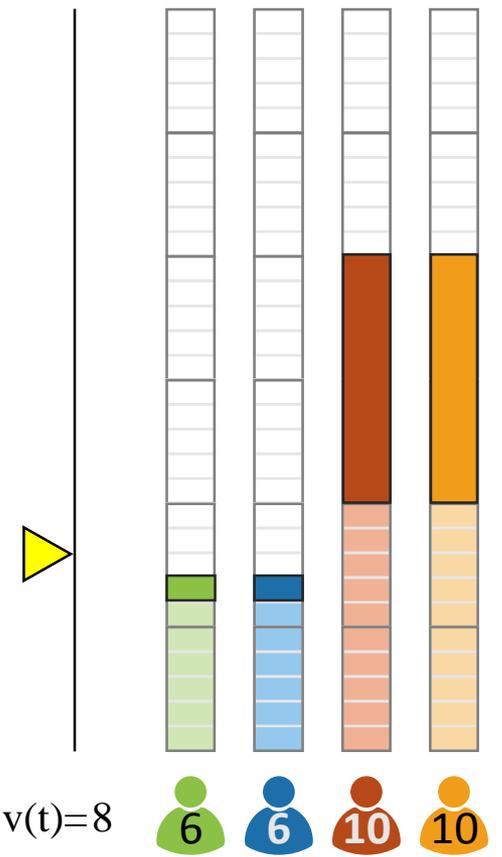
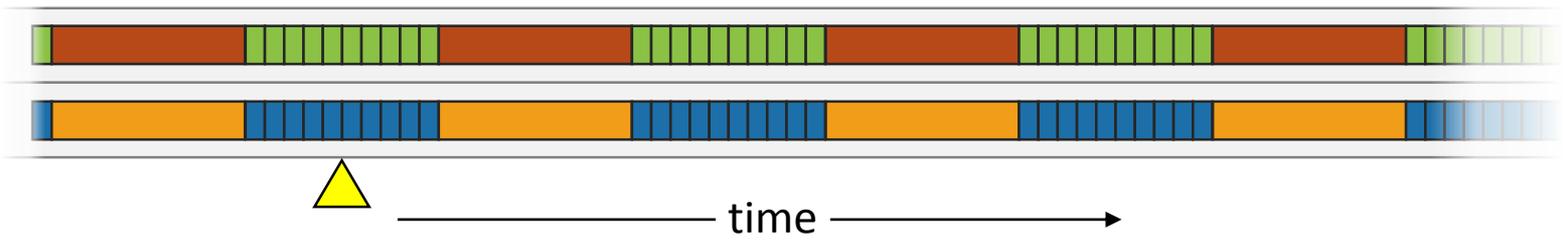
# WF<sup>2</sup>Q



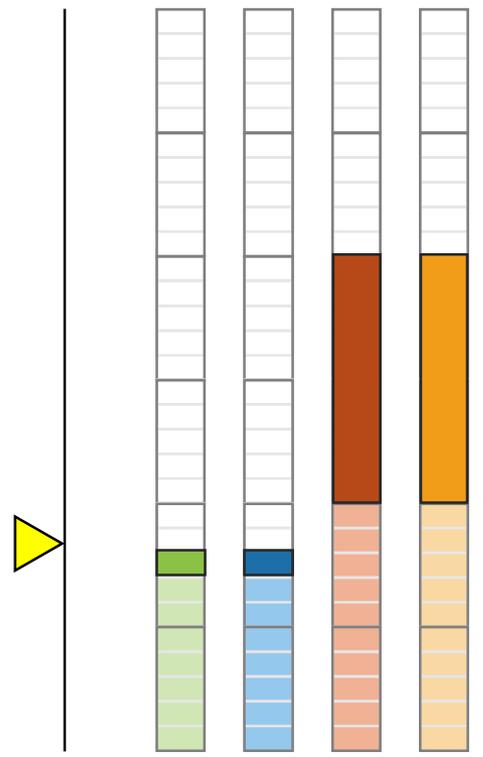
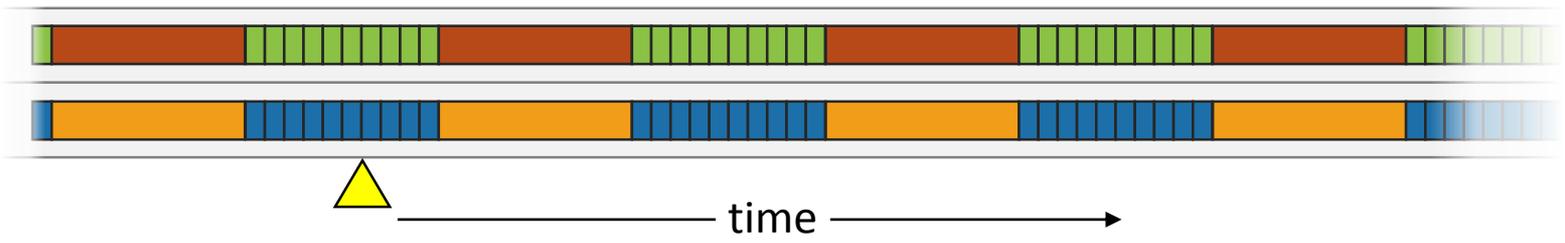
$v(t) = 7 \frac{1}{2}$     5    5    10    10



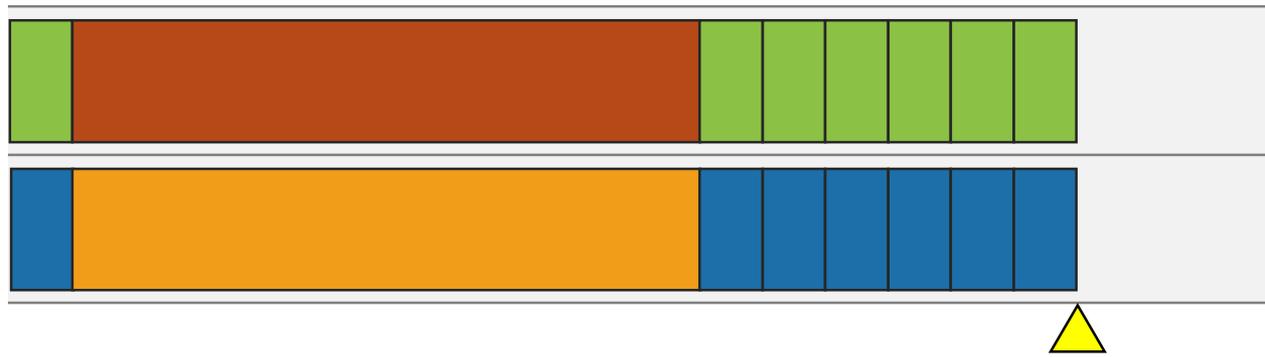
# WF<sup>2</sup>Q



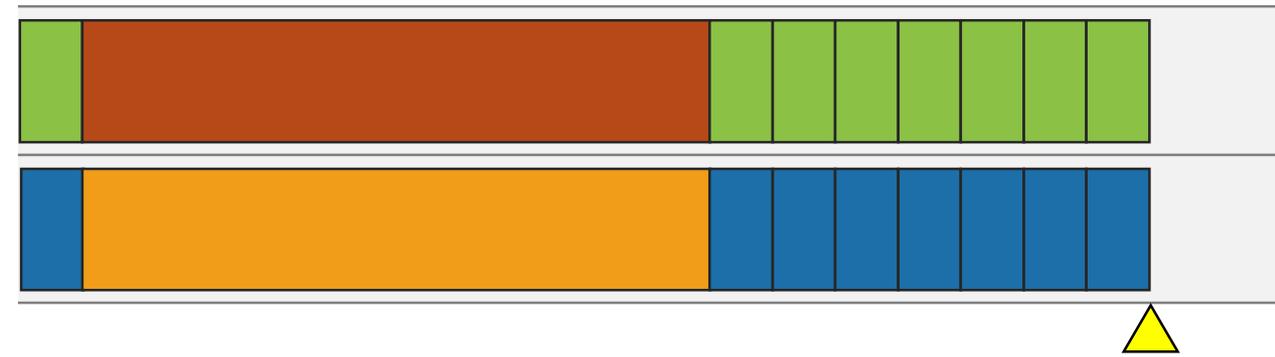
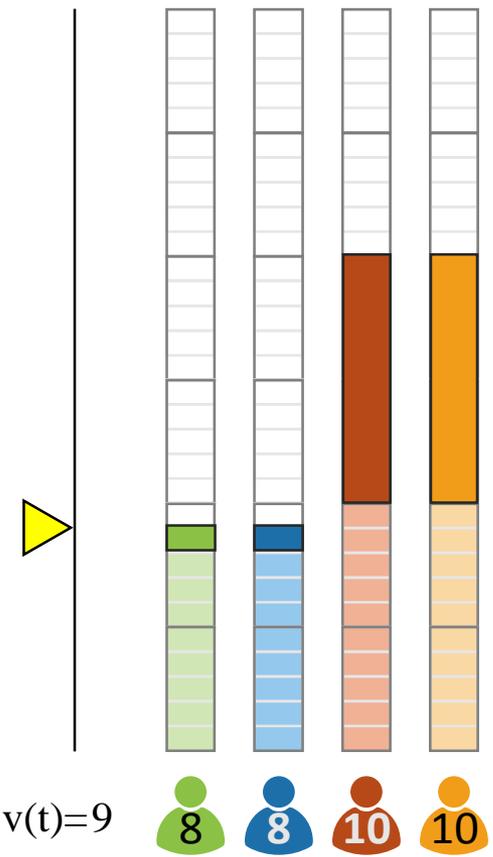
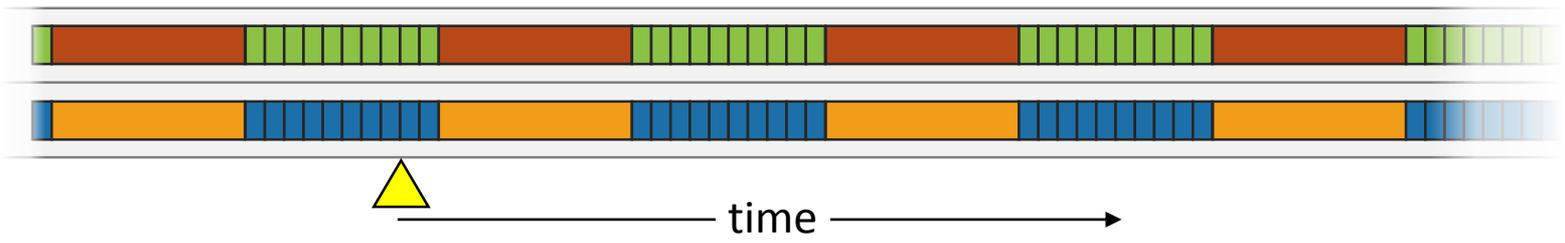
# WF<sup>2</sup>Q



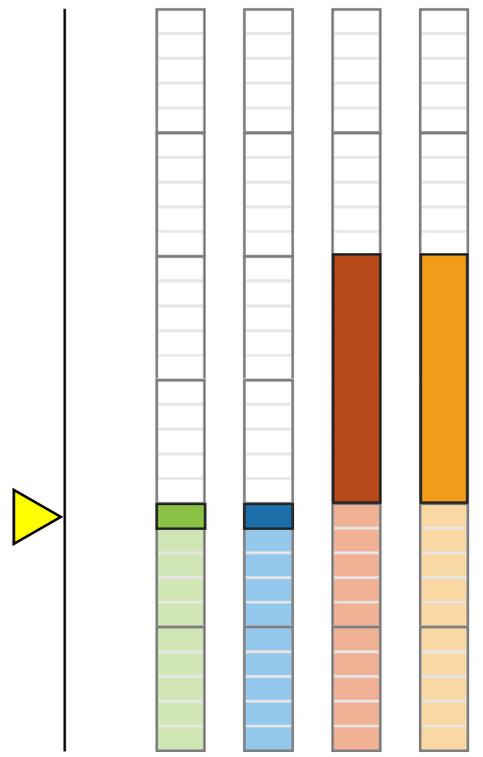
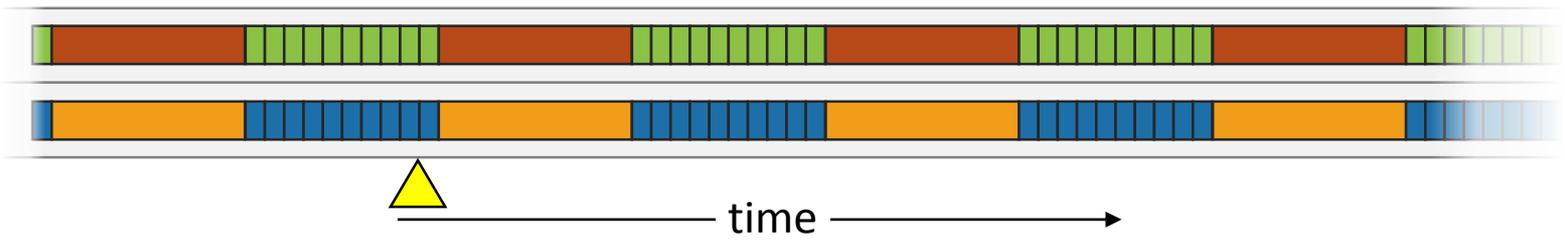
$v(t) = 8\frac{1}{2}$     7    7    10    10



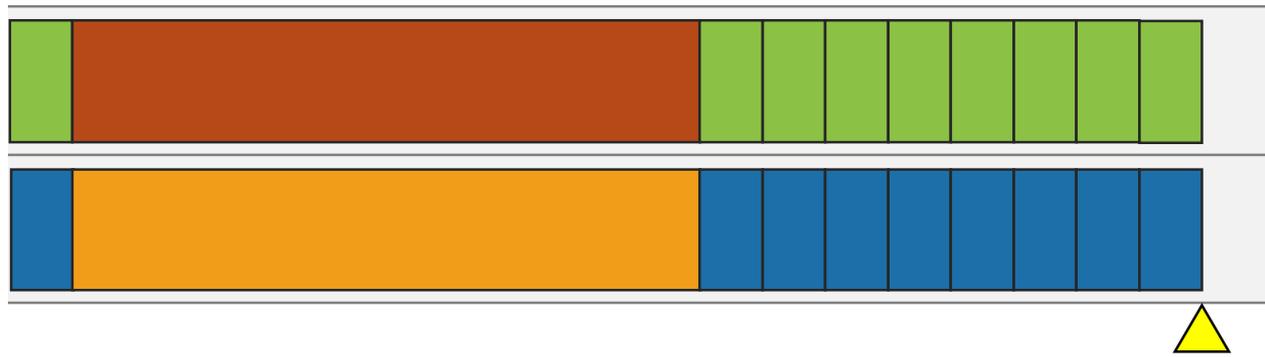
# WF<sup>2</sup>Q



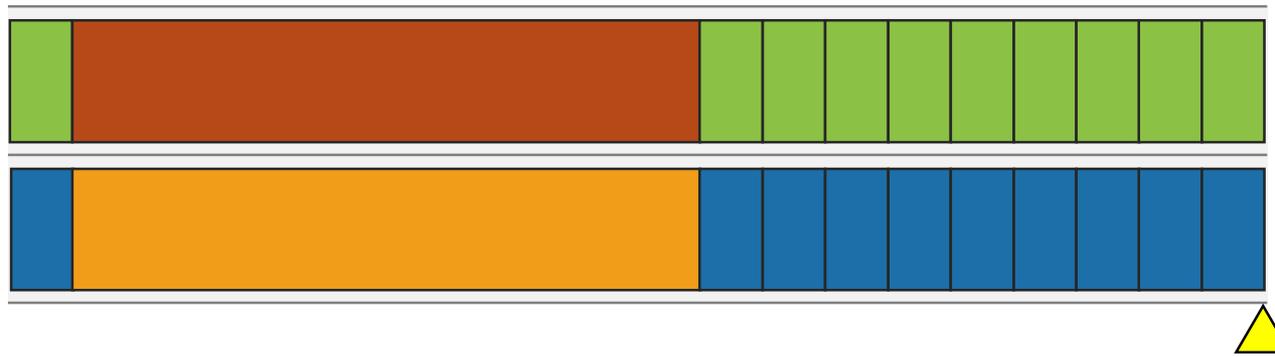
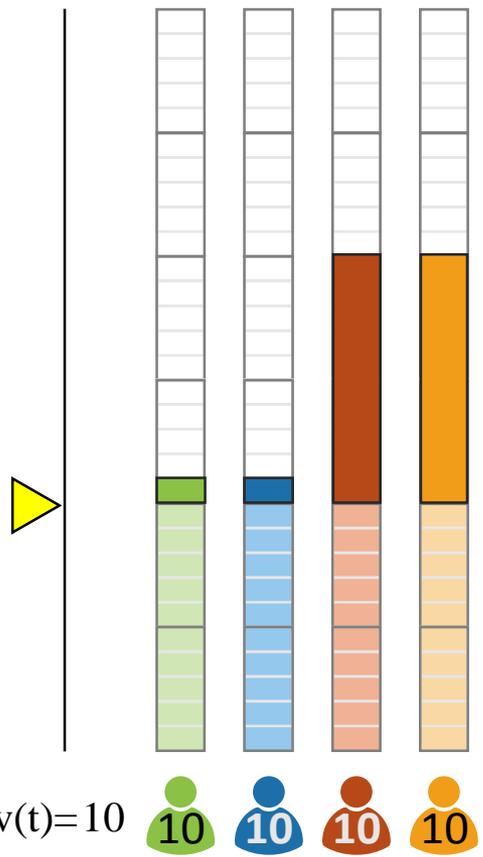
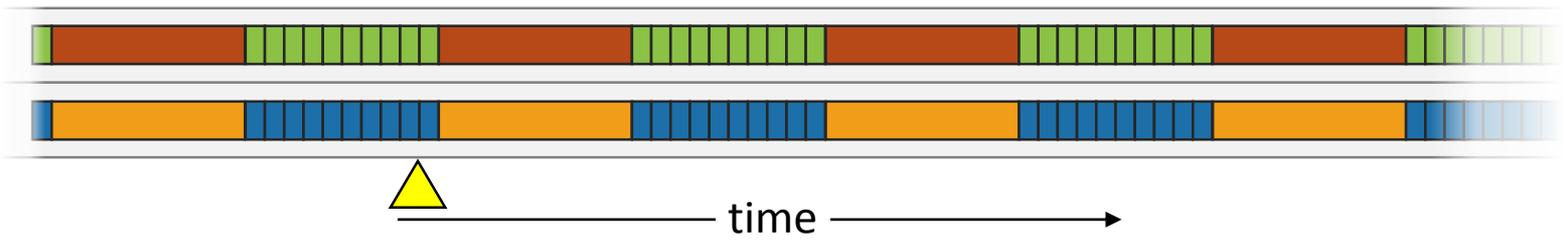
# WF<sup>2</sup>Q



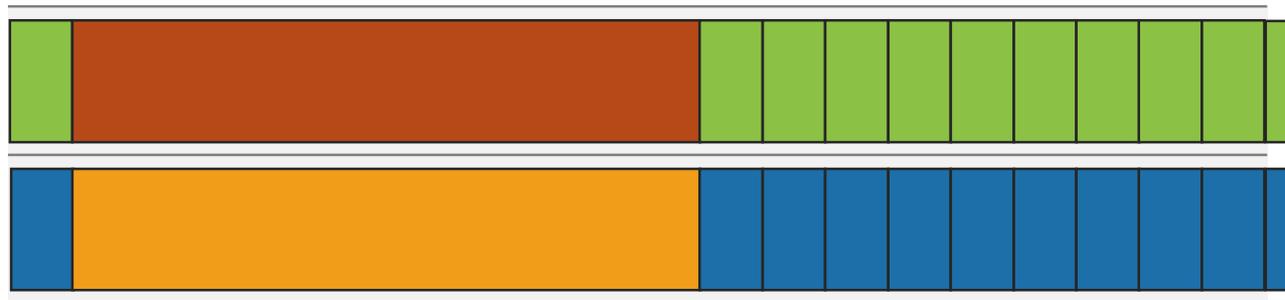
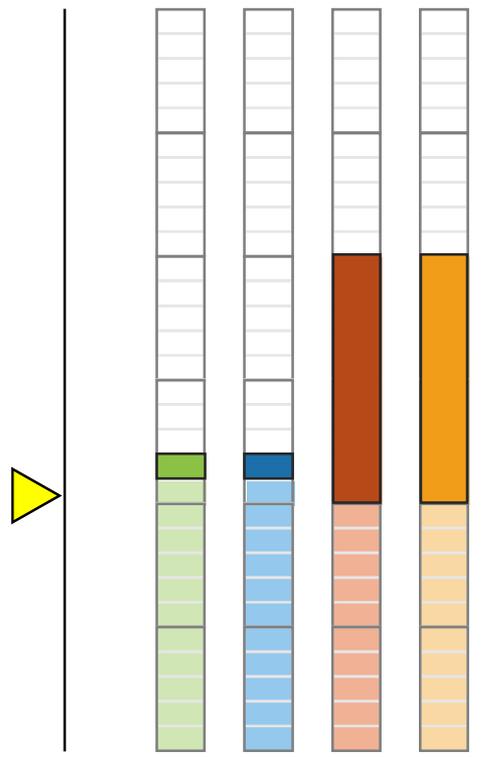
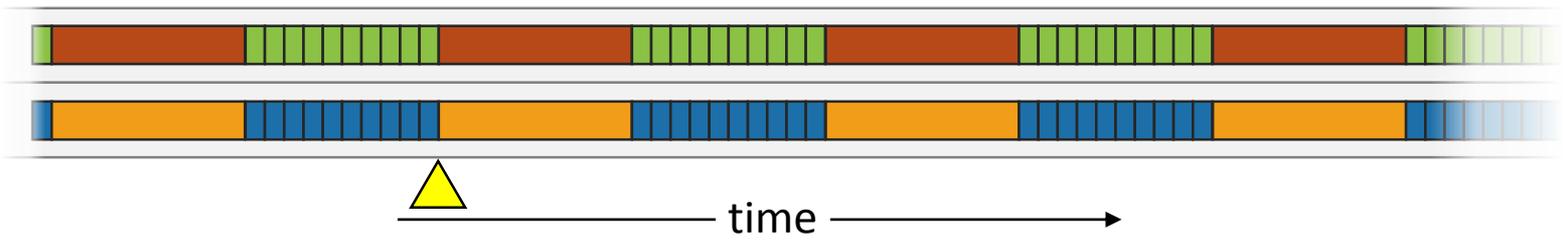
$v(t) = 9\frac{1}{2}$     9    9    10    10



# WF<sup>2</sup>Q

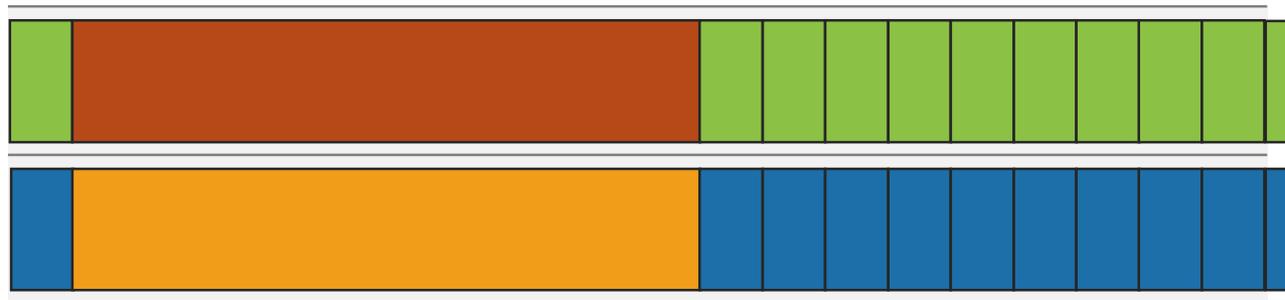
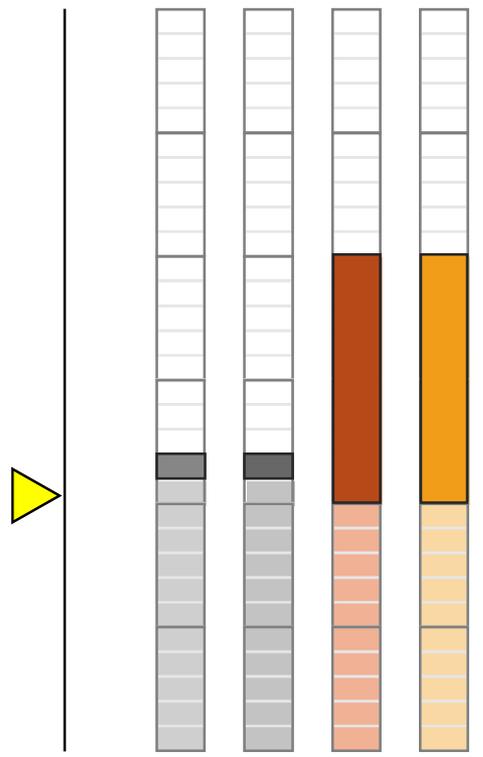
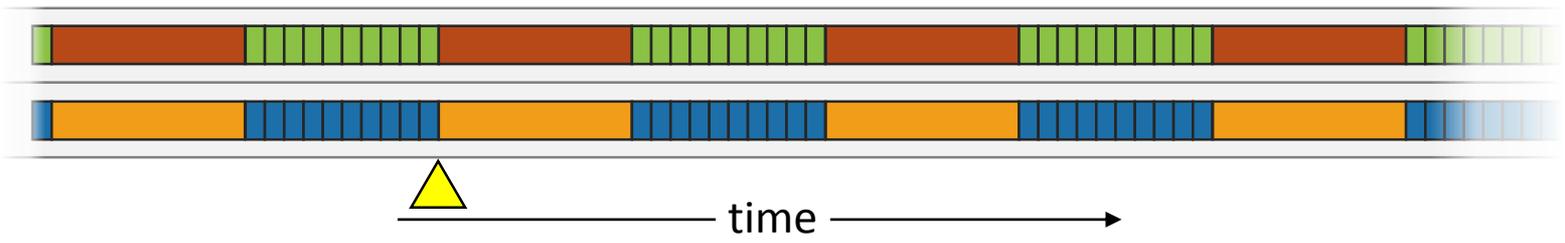


# WF<sup>2</sup>Q



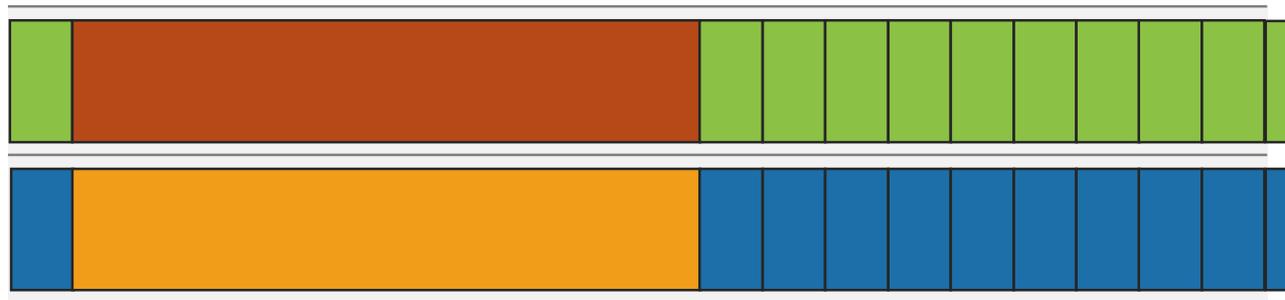
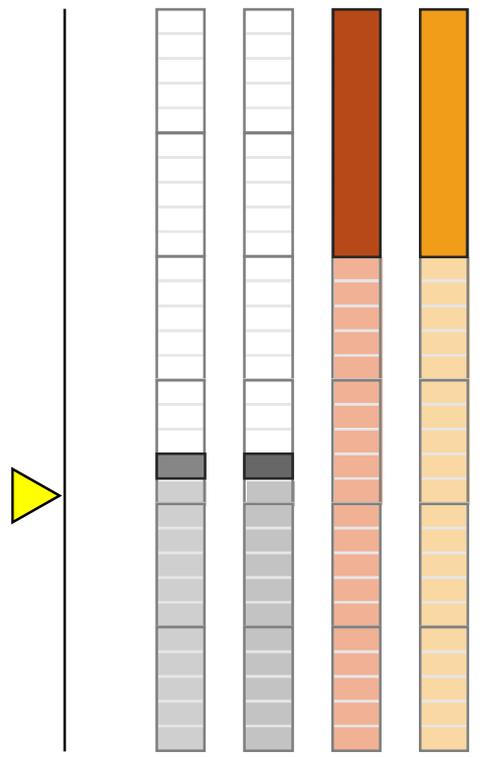
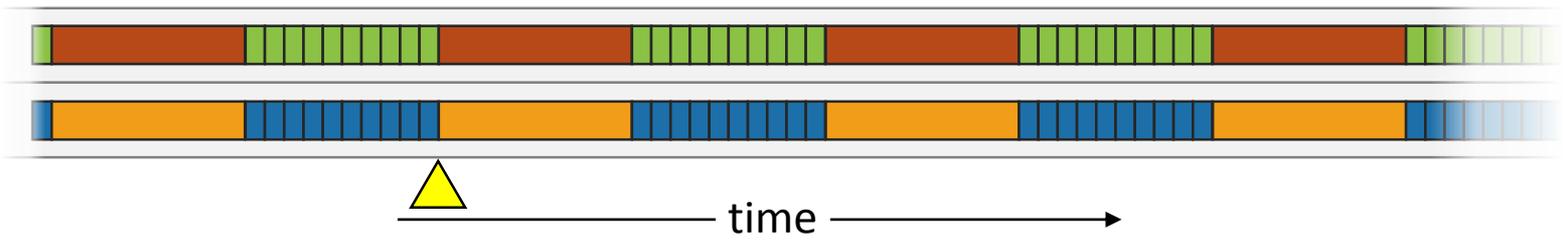
$v(t) = 10^{1/2}$  11 11 10 10

# WF<sup>2</sup>Q



$v(t) = 10^{1/2}$  11 11 10 10

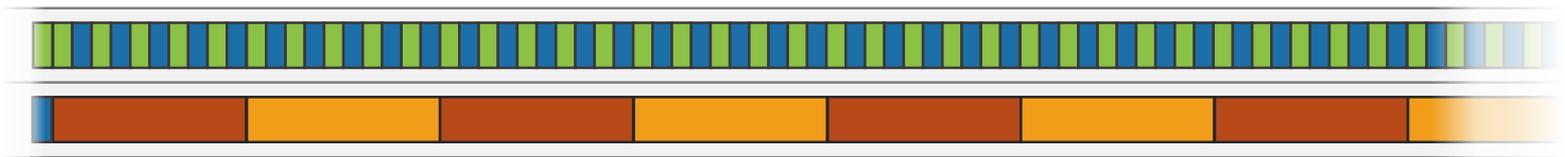
# WF<sup>2</sup>Q



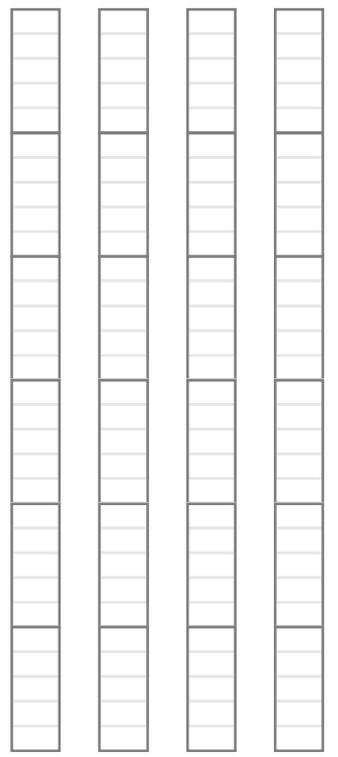
$v(t) = 10^{1/2}$

11 11 10 10

# 2DFQ



time →

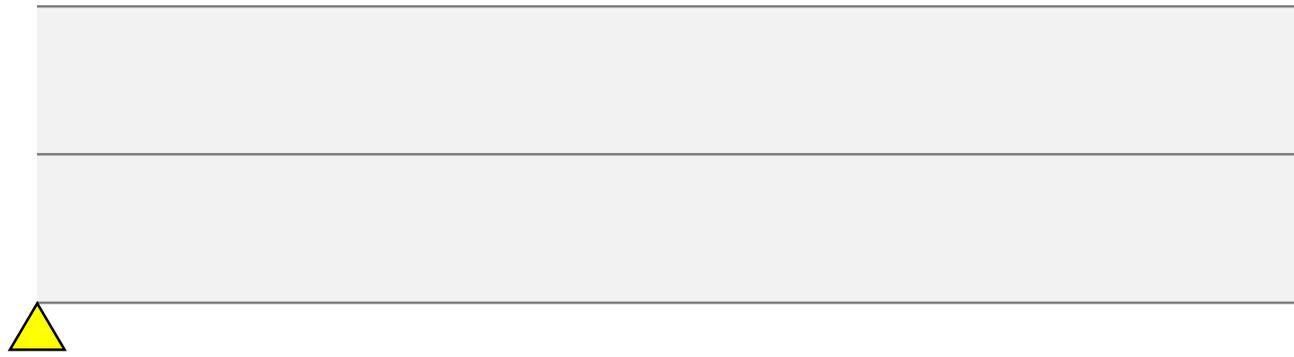
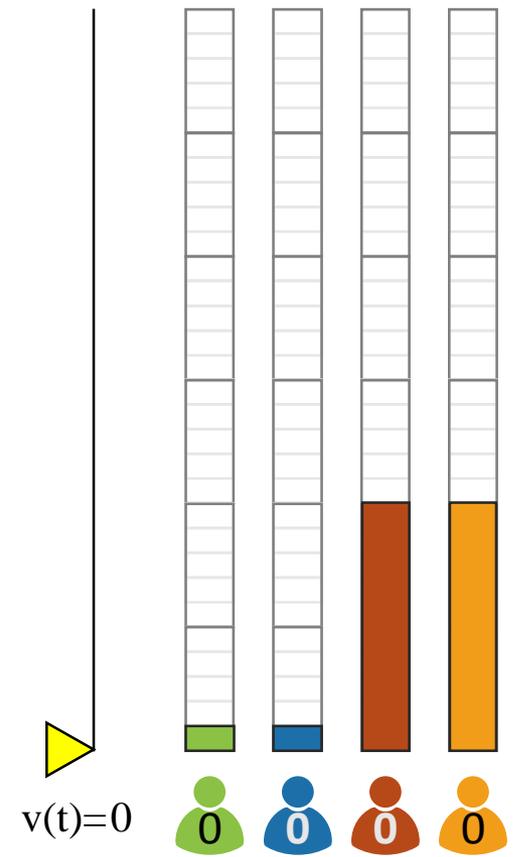
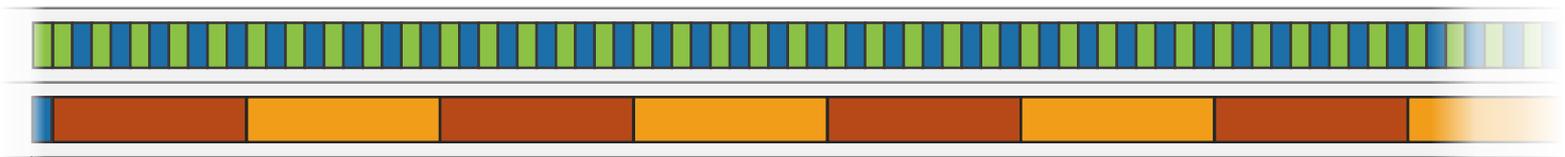


$v(t)=0$

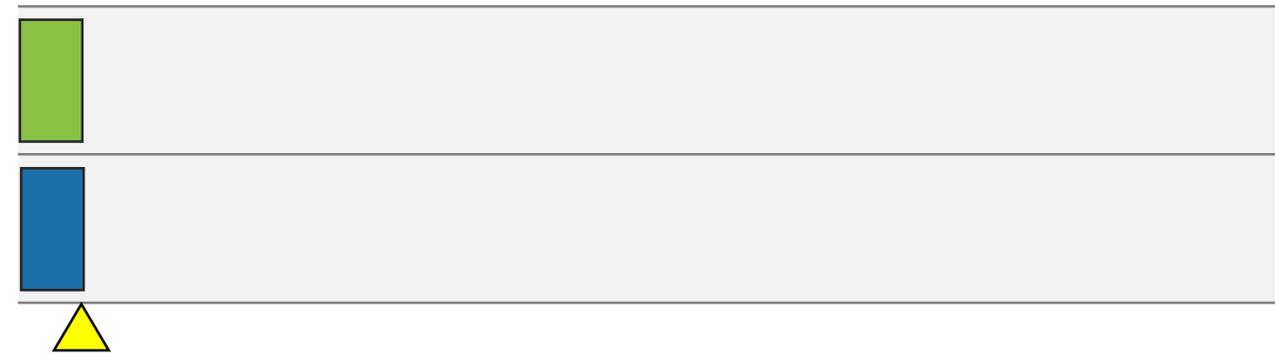
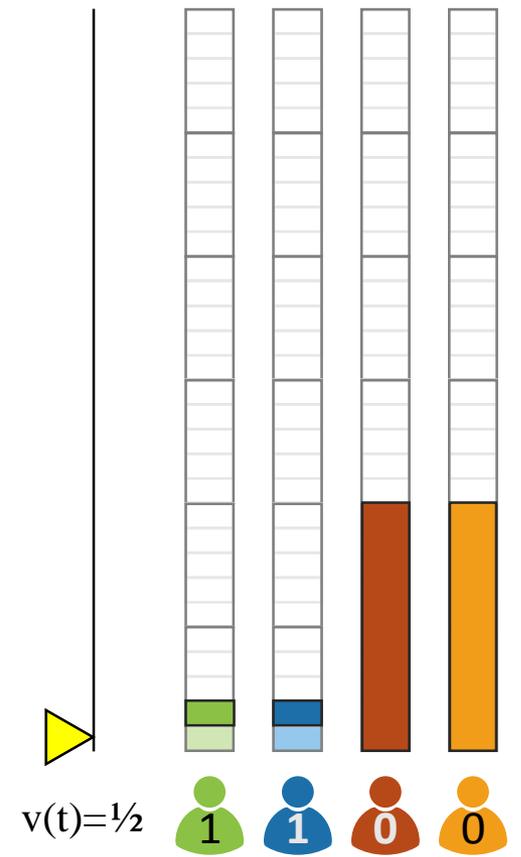
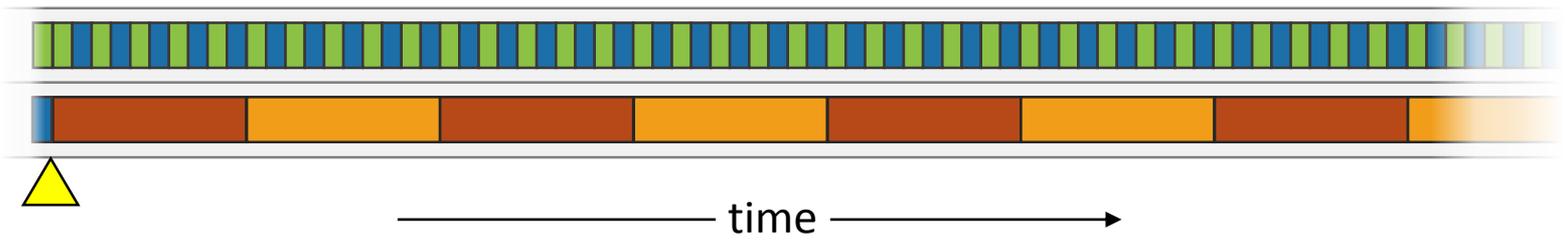
0 0 0 0



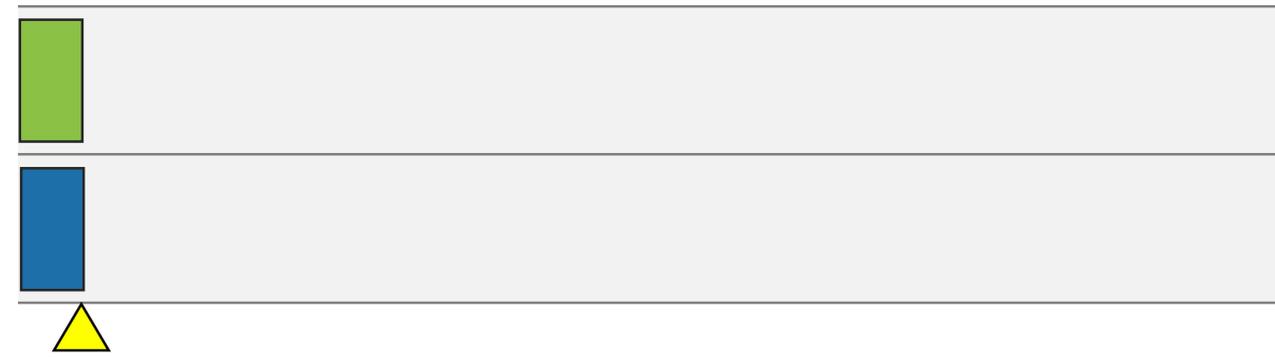
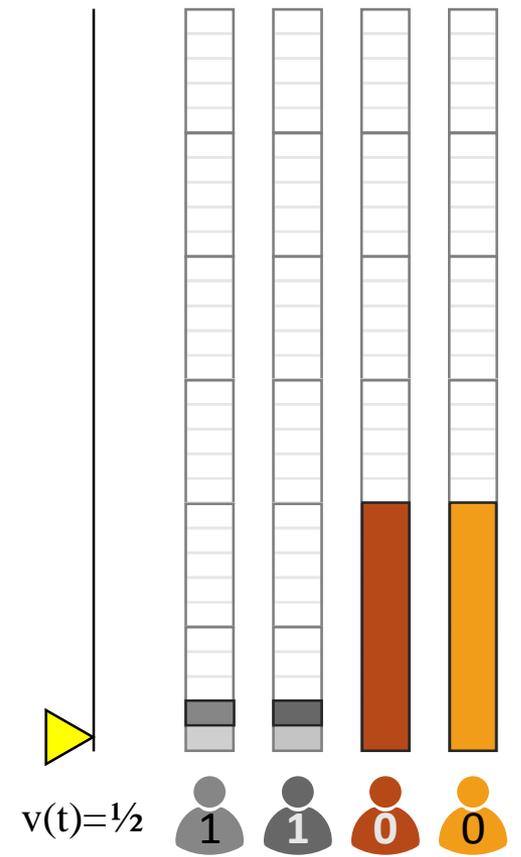
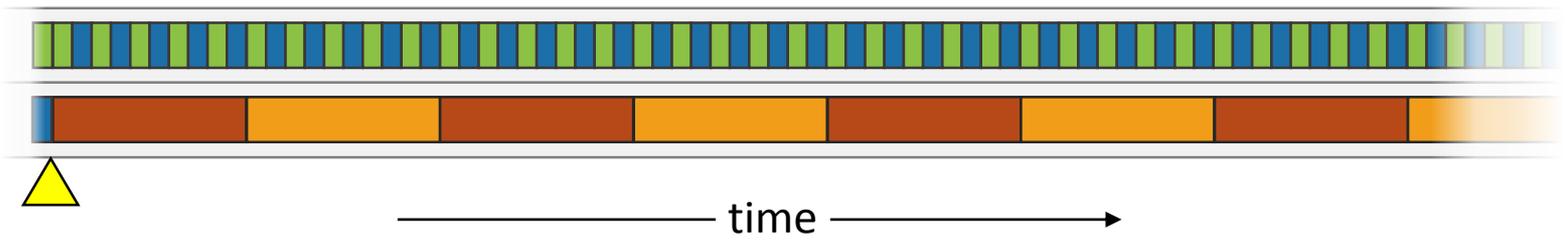
# 2DFQ



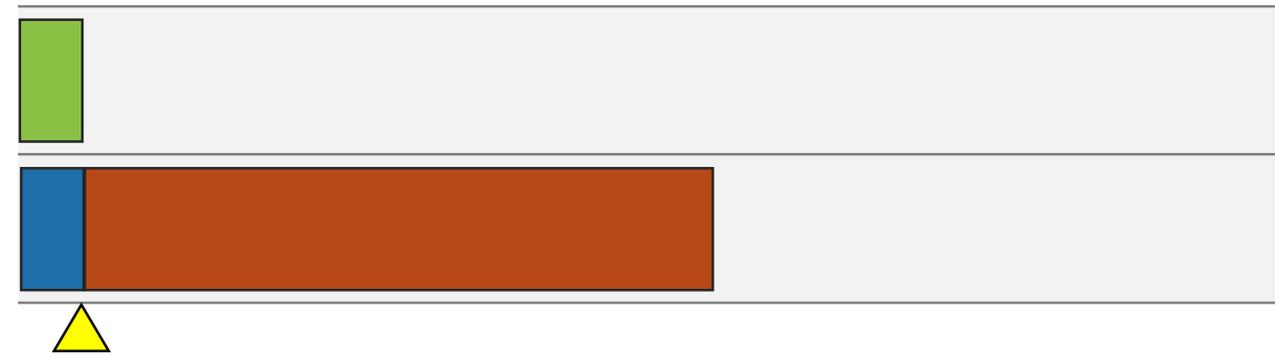
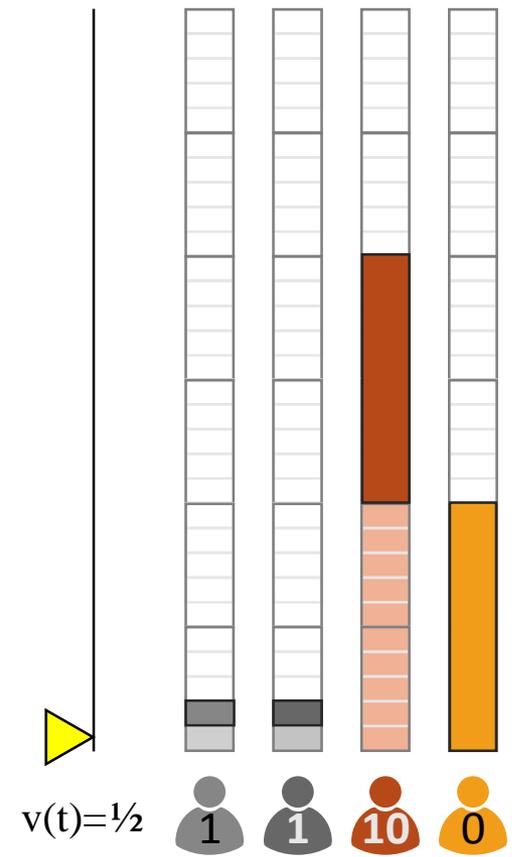
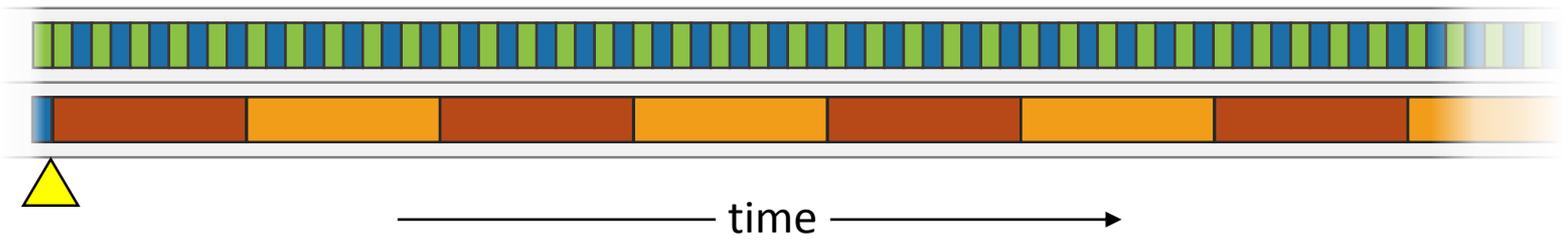
# 2DFQ



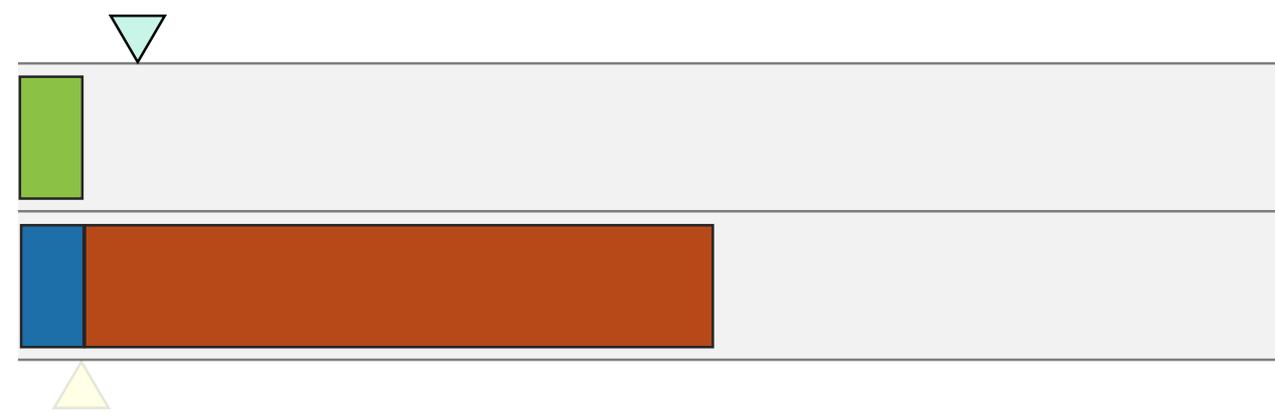
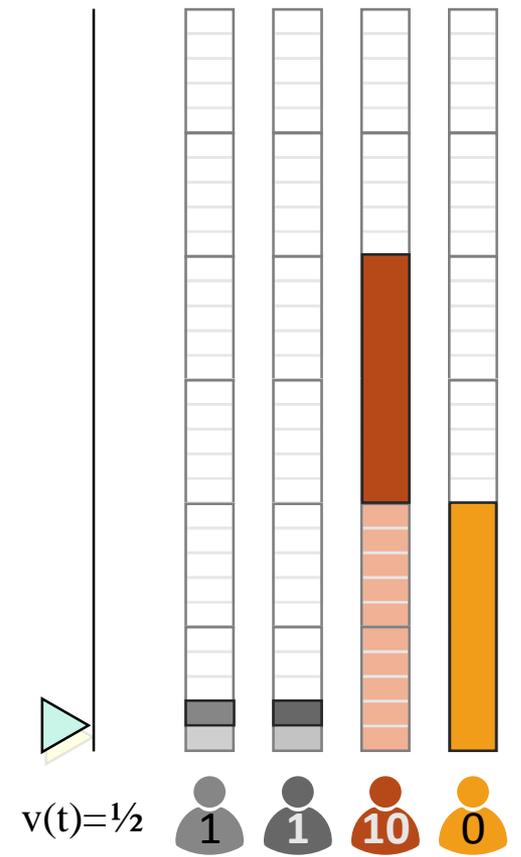
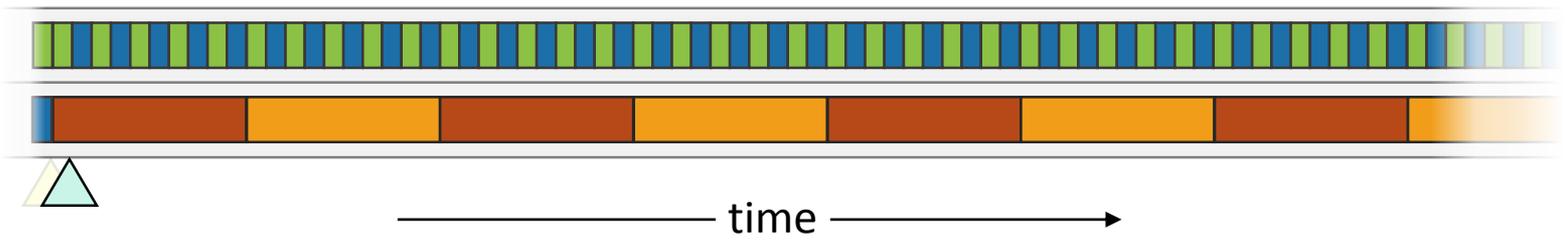
# 2DFQ



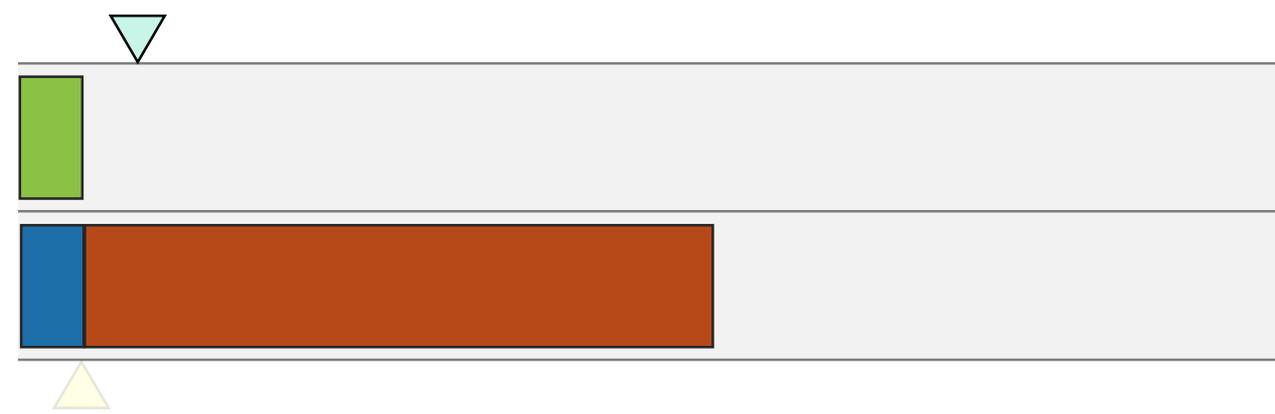
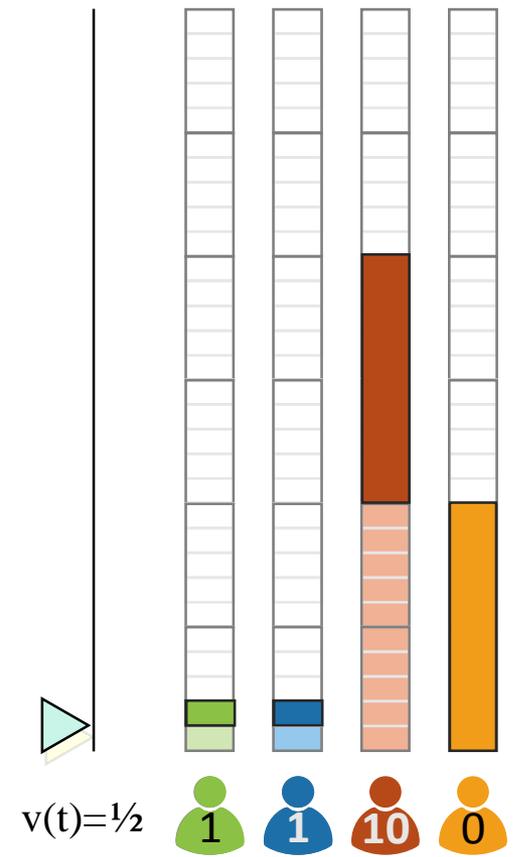
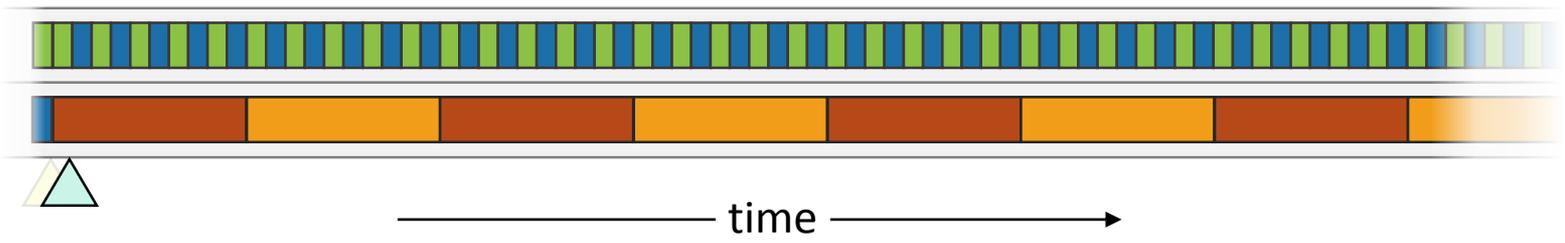
# 2DFQ



# 2DFQ

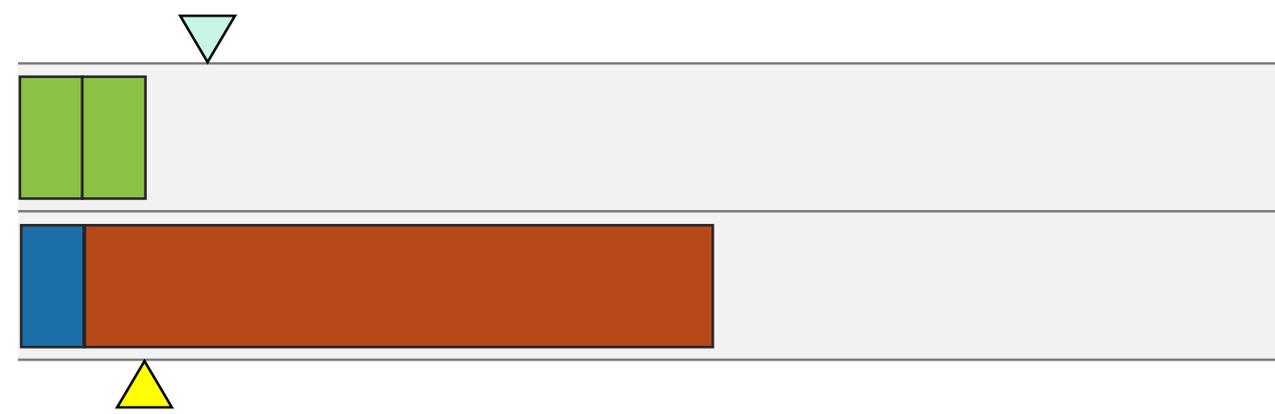
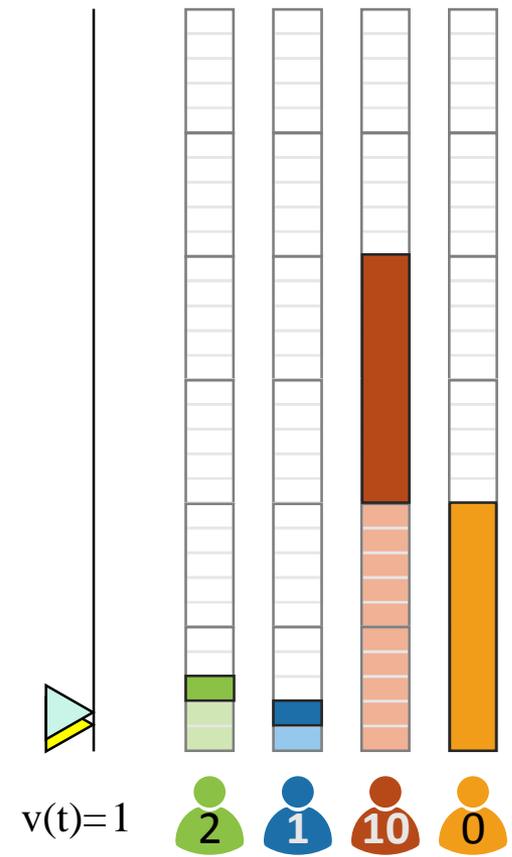
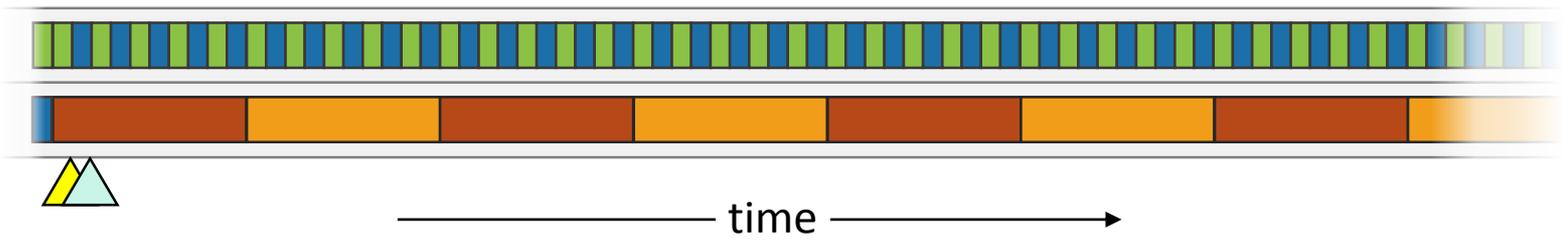


# 2DFQ



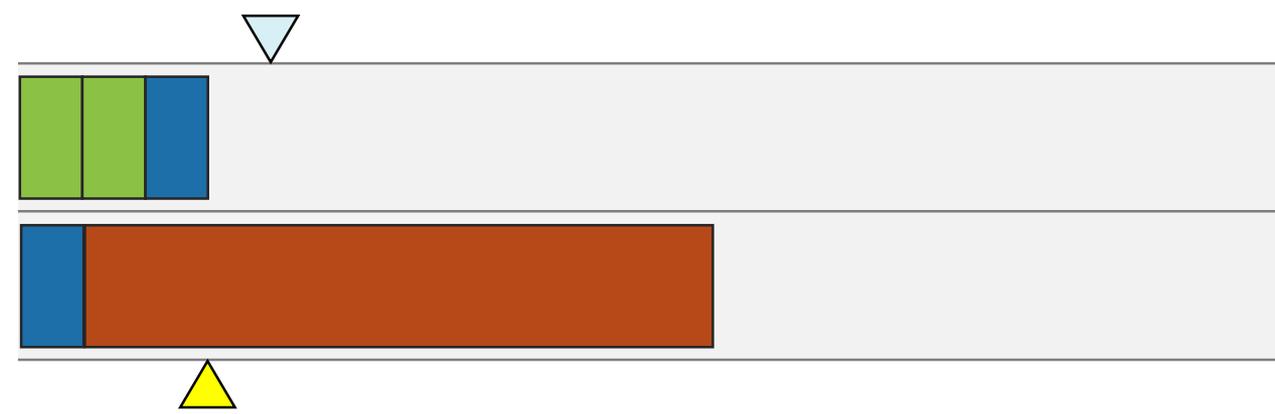
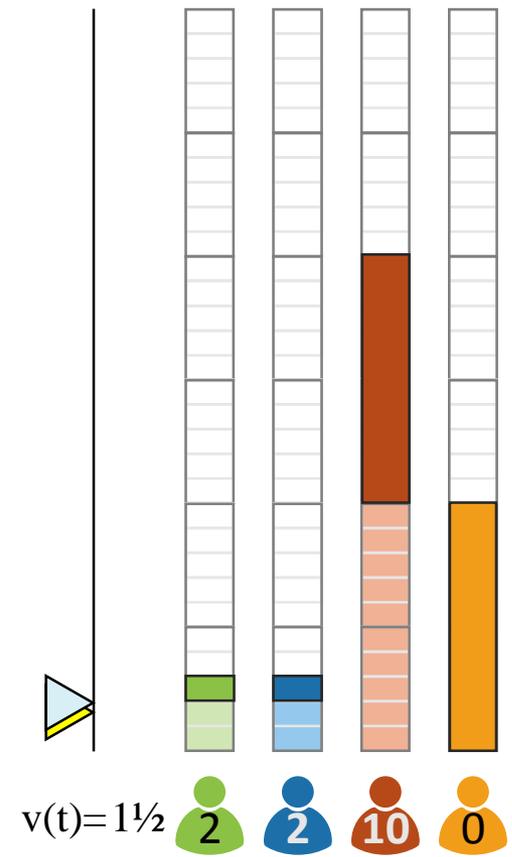
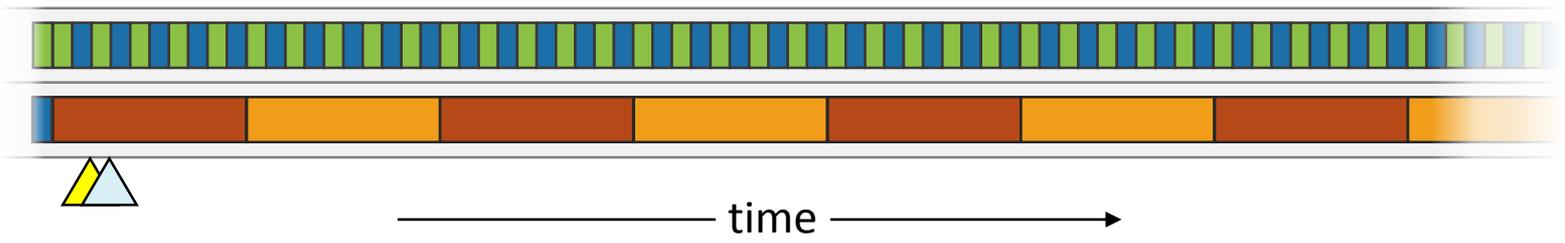


# 2DFQ



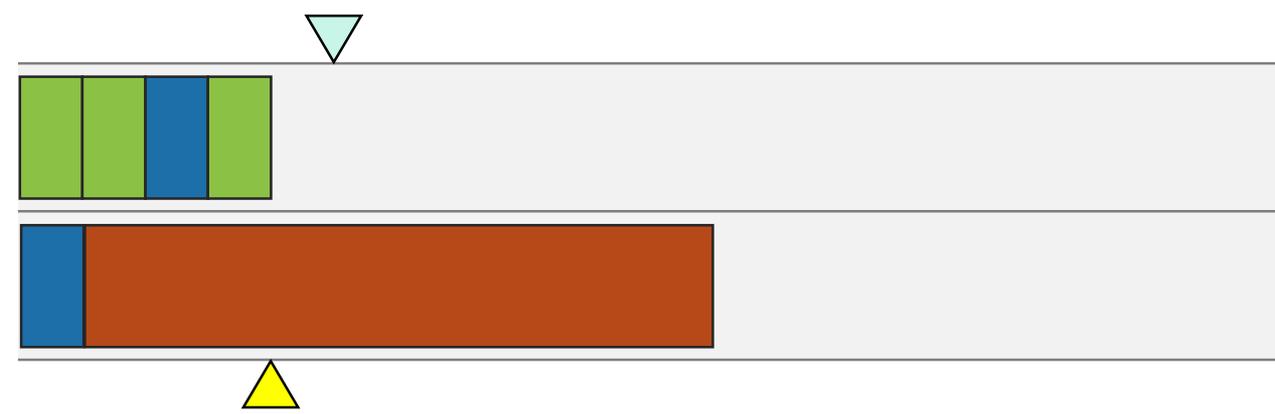
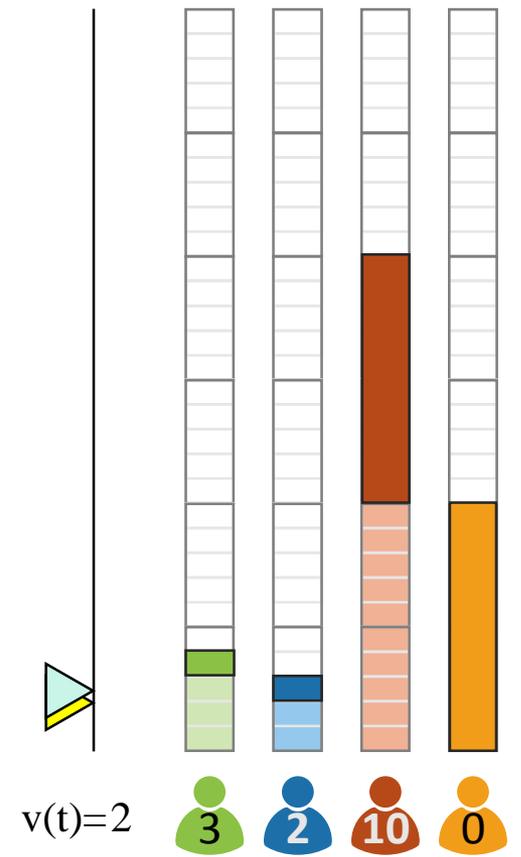
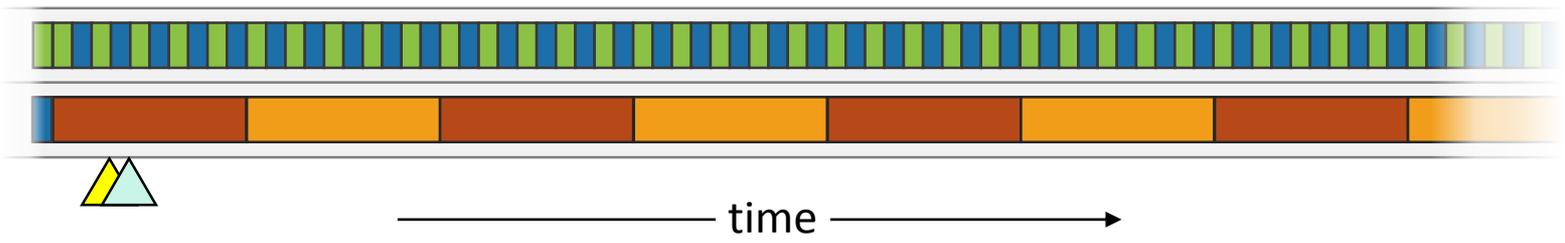
$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ



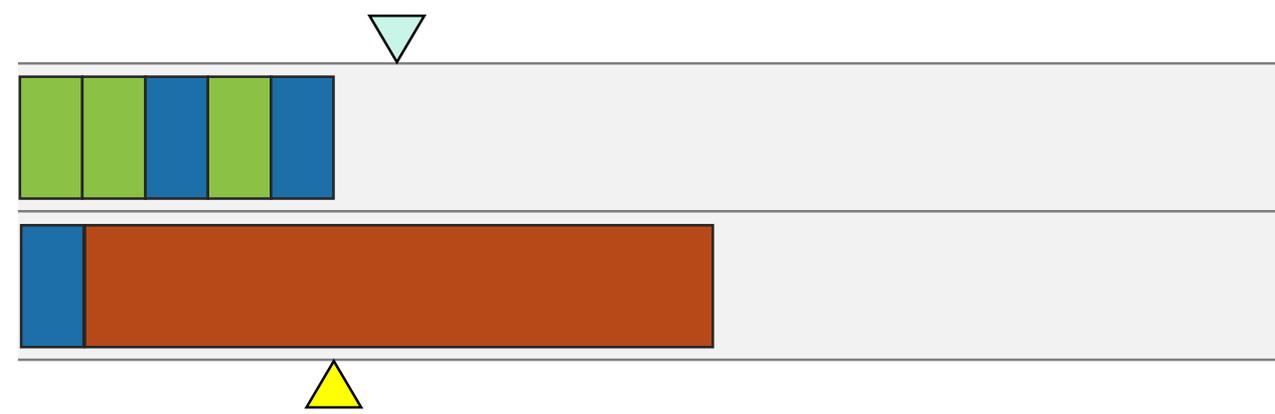
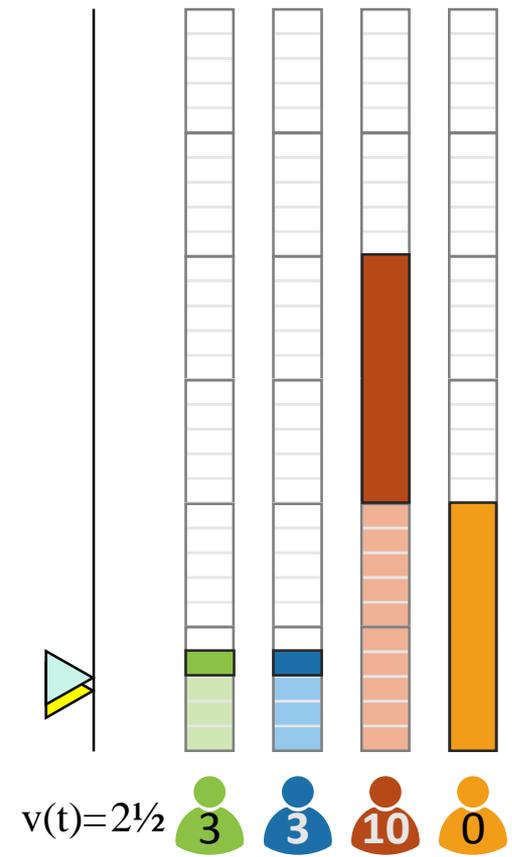
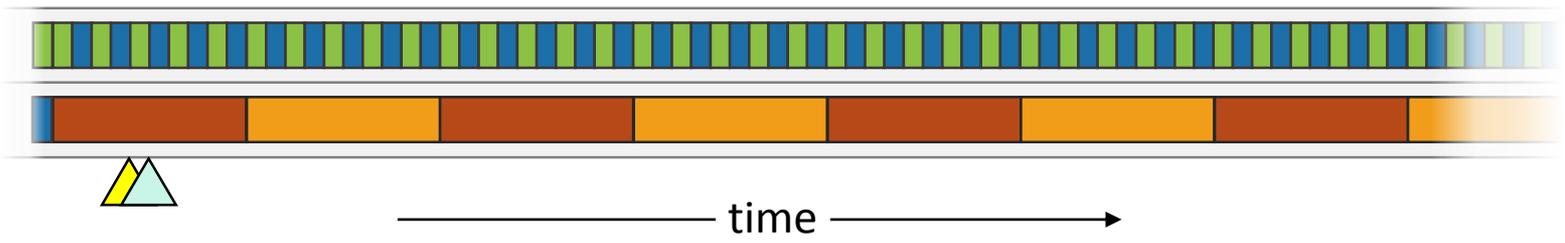
$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ



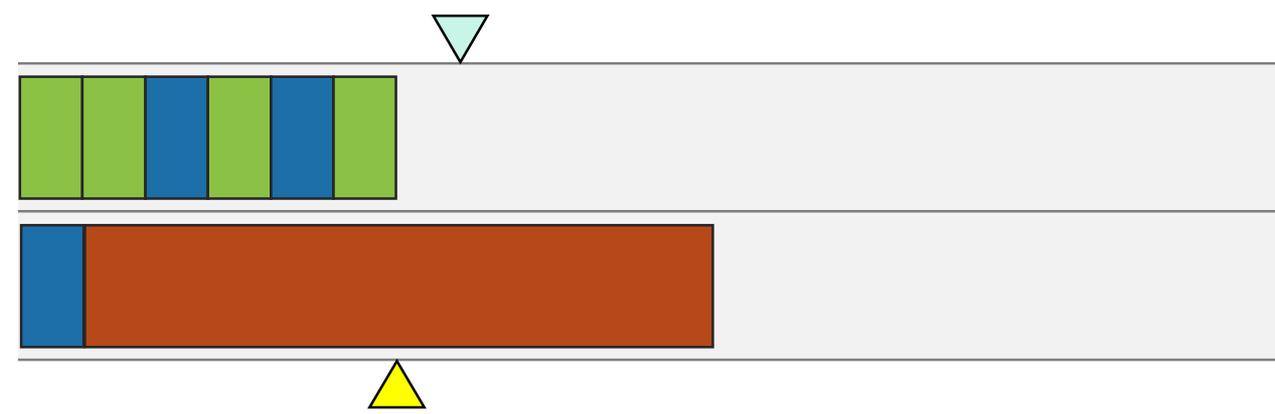
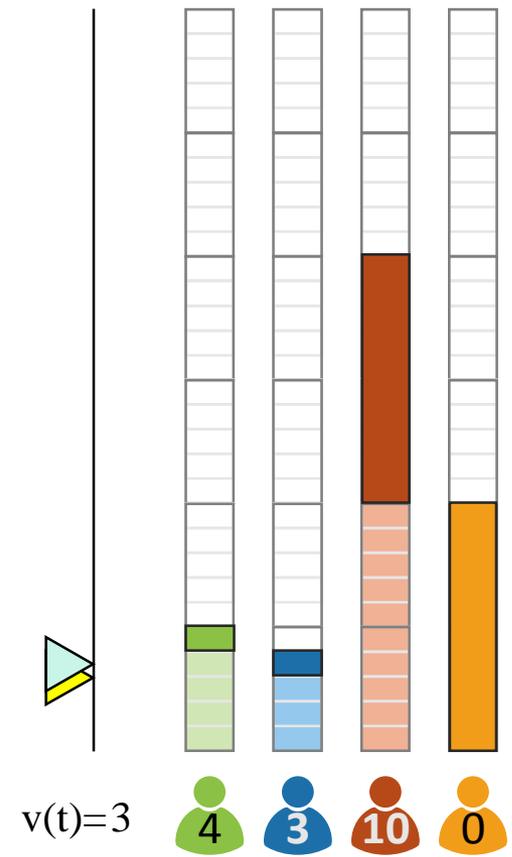
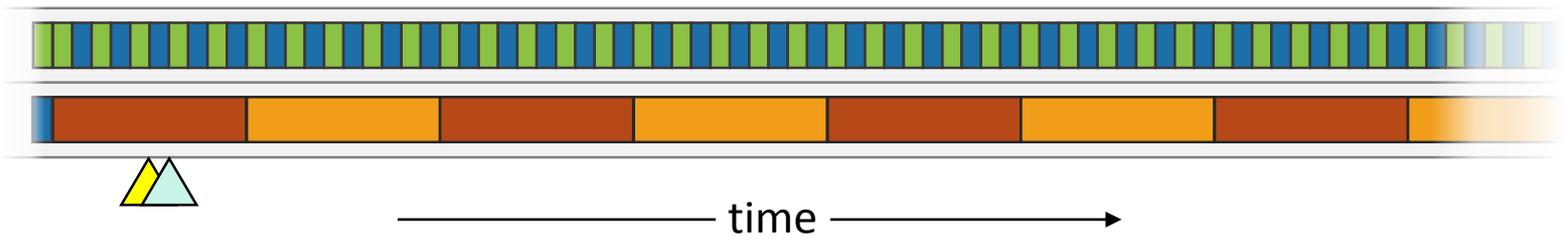
$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ



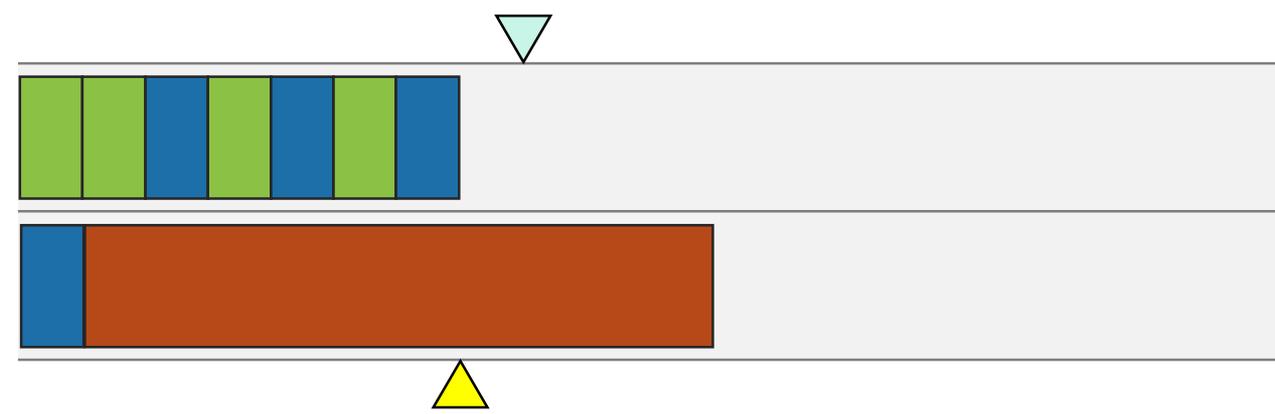
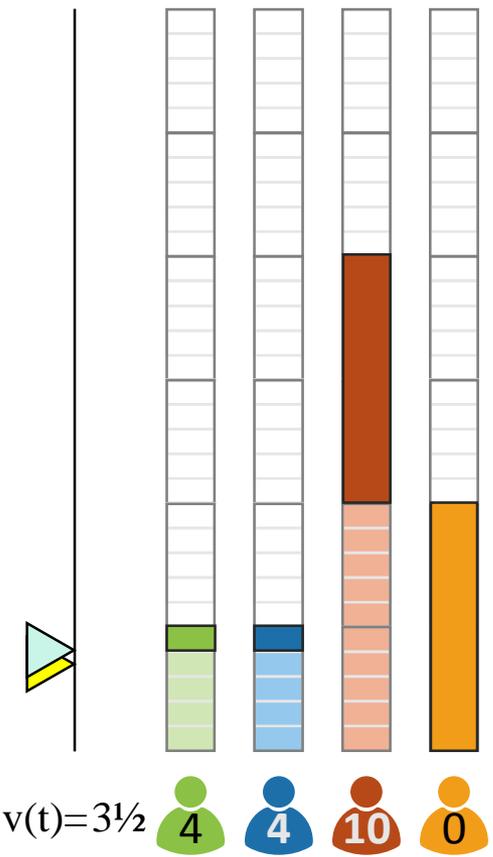
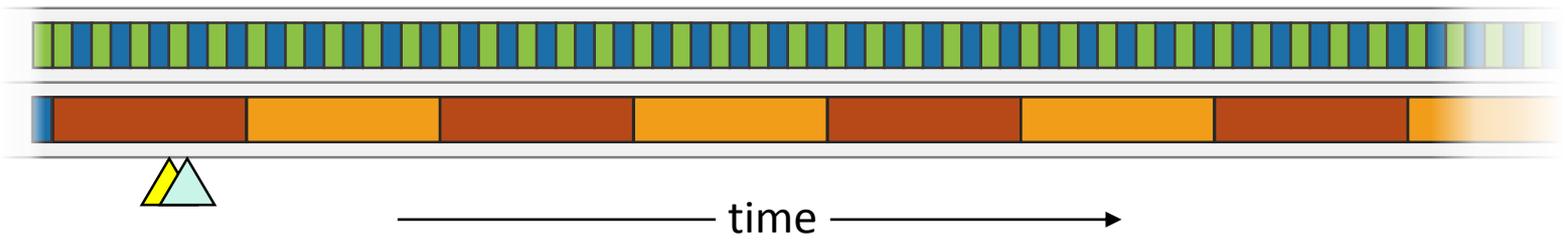
$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ



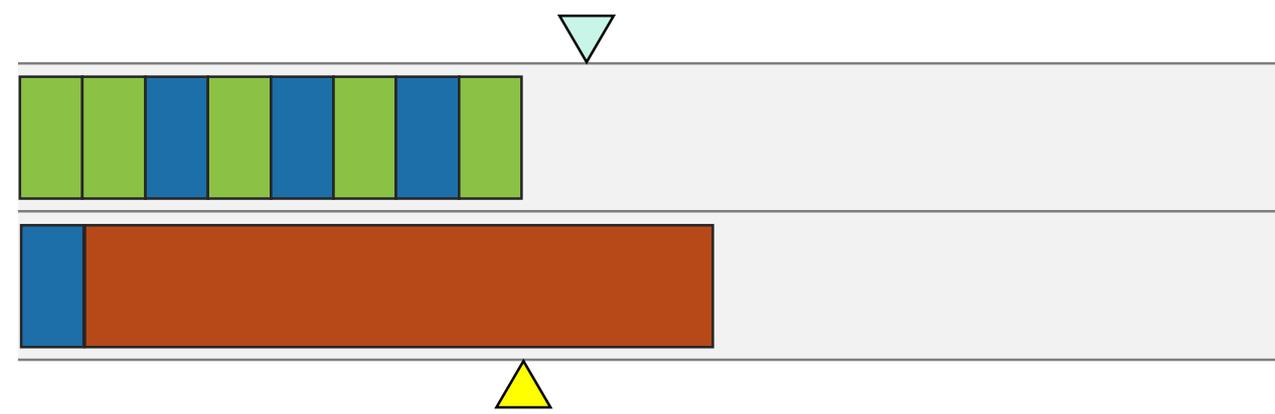
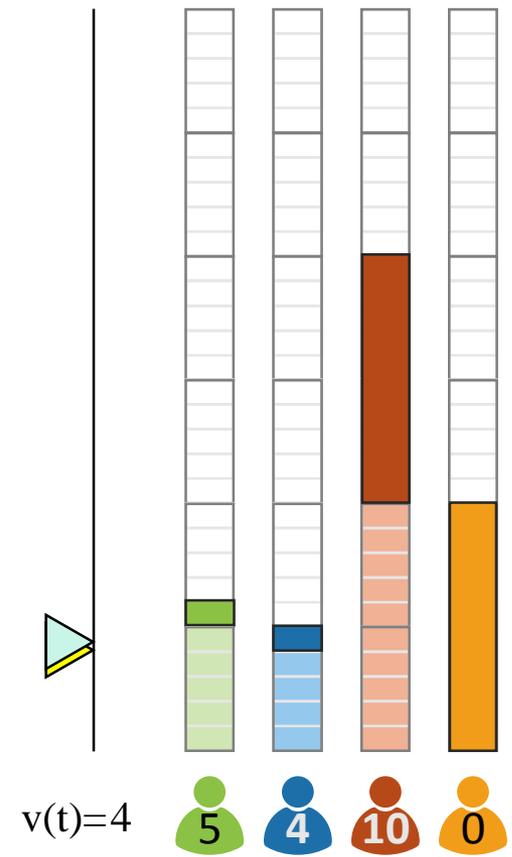
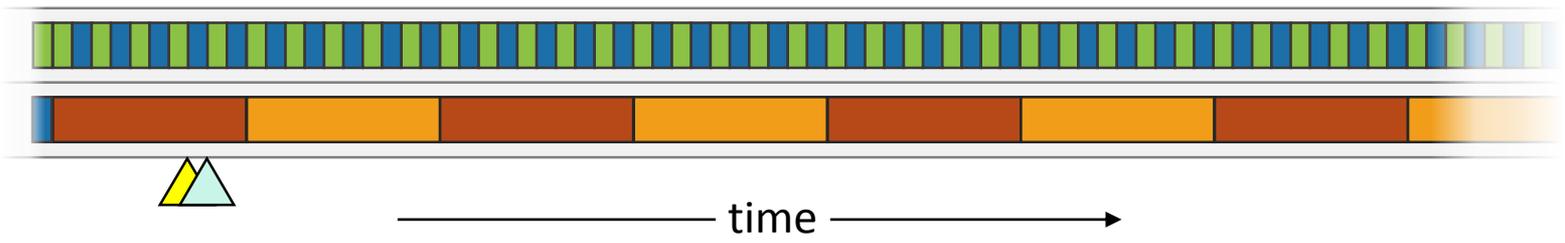
$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ



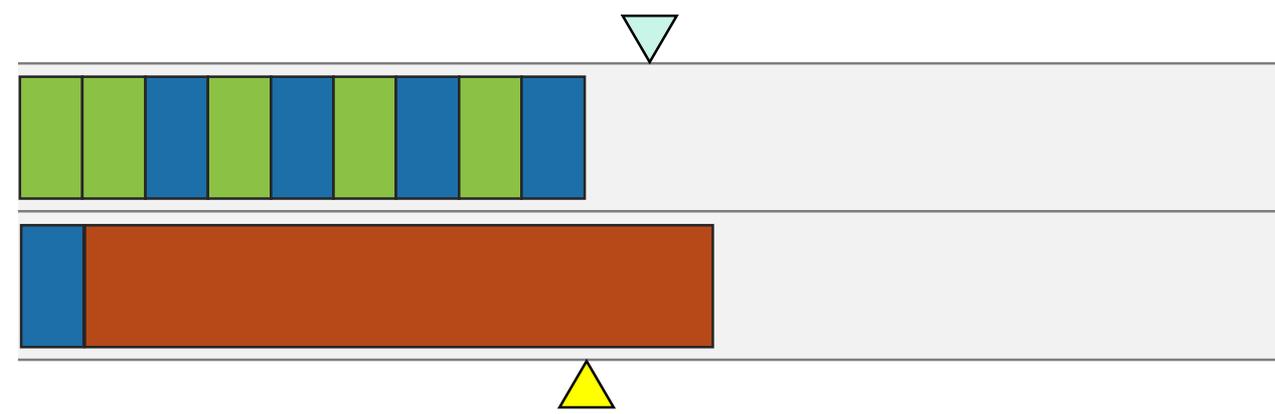
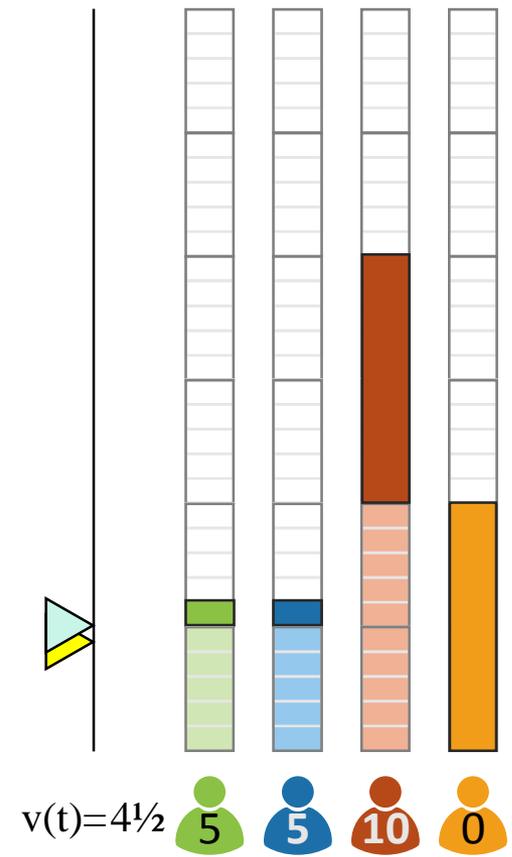
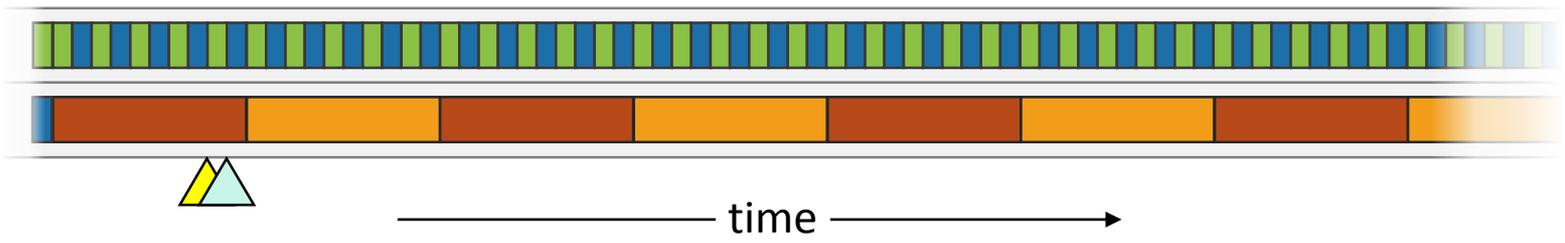
$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ



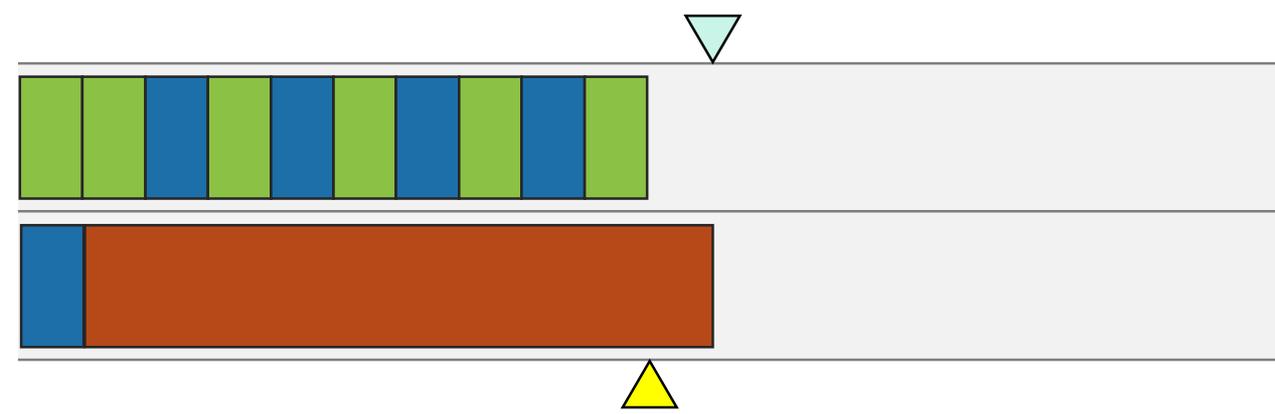
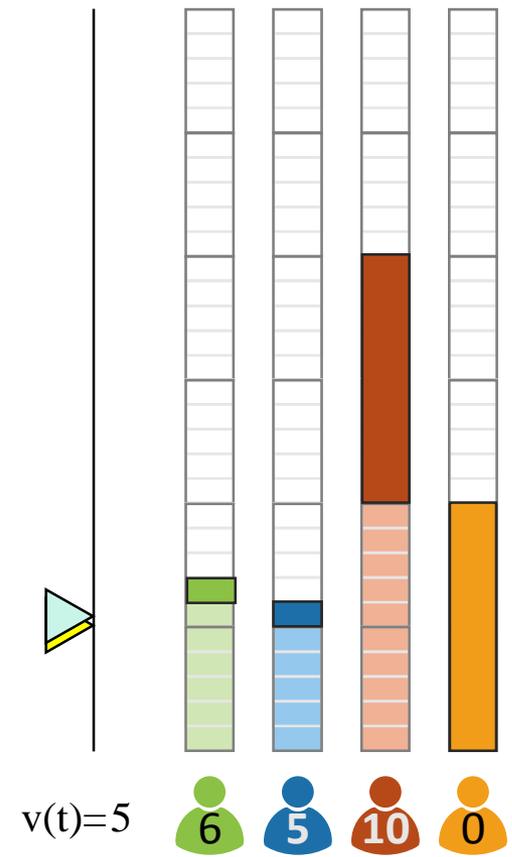
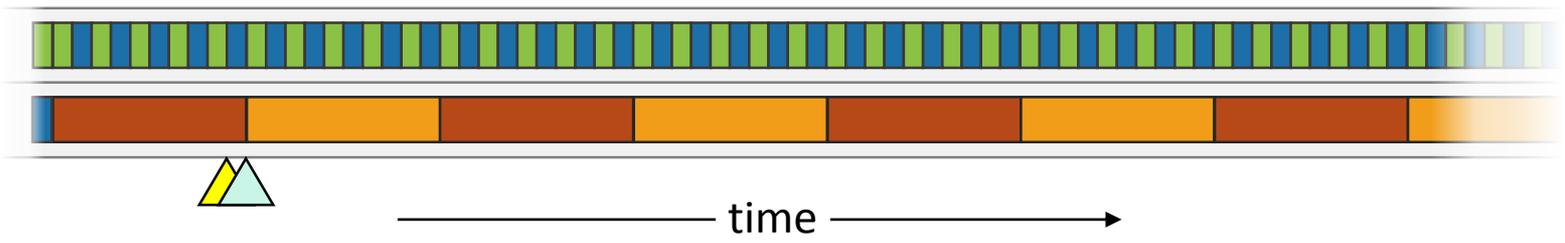
$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ



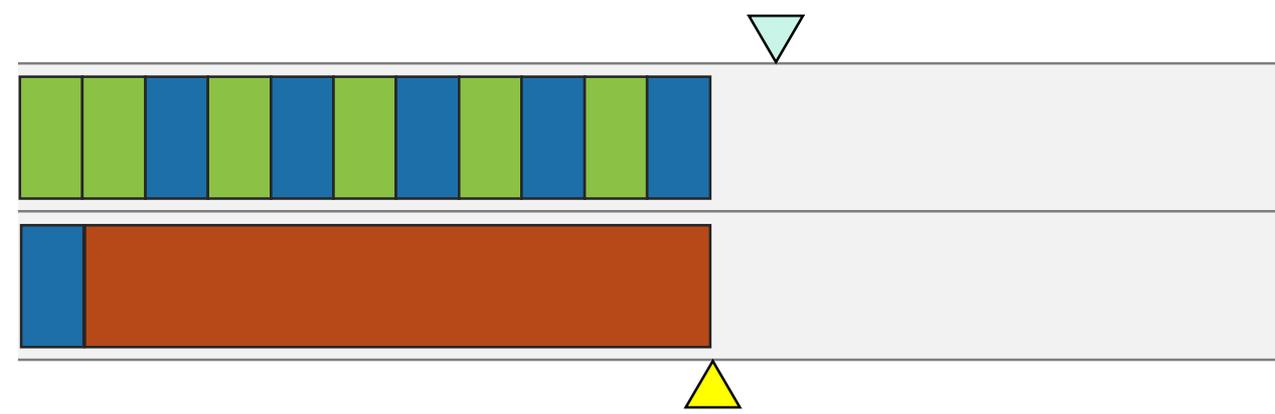
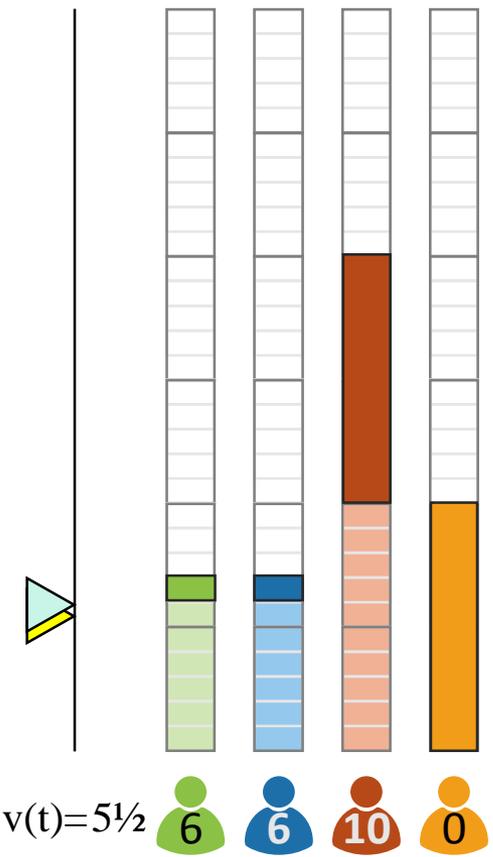
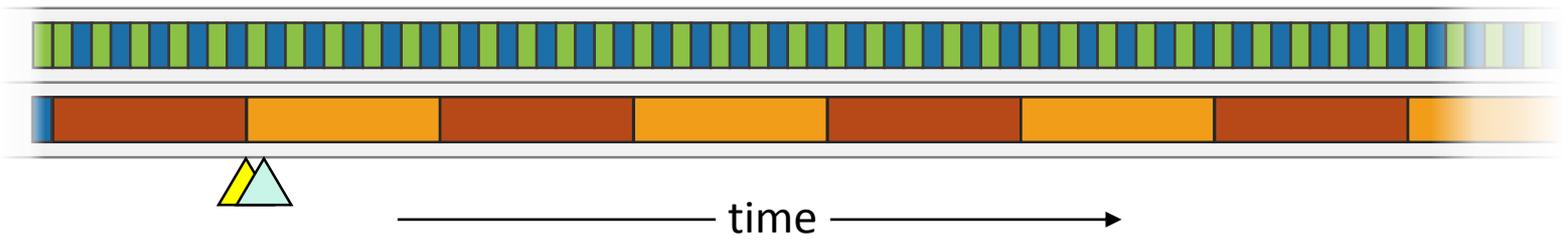
$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ



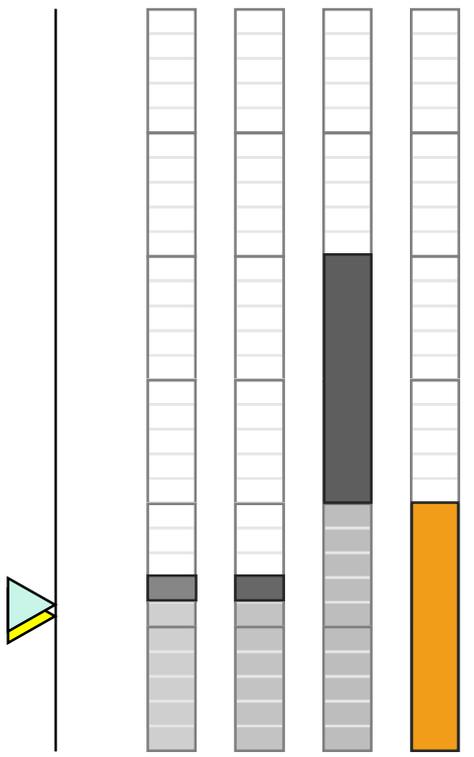
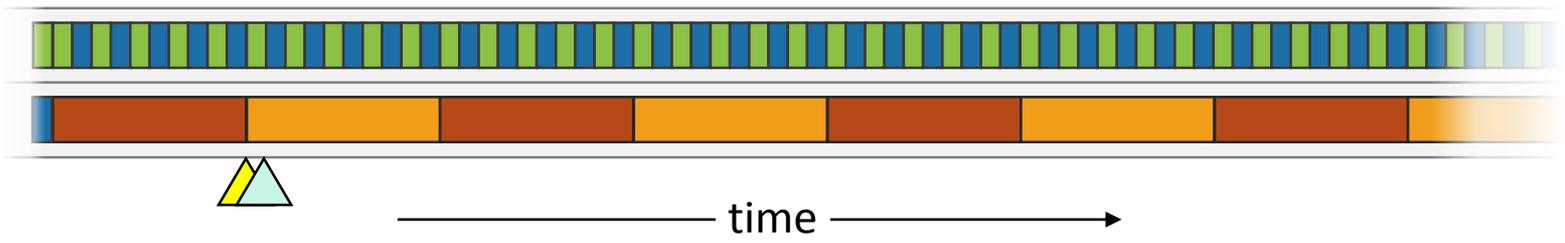
$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ

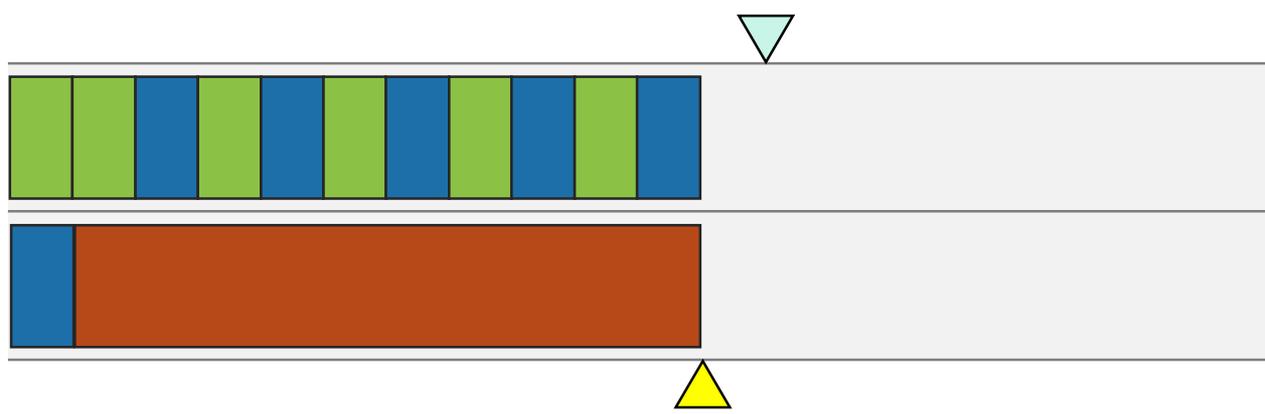


$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ

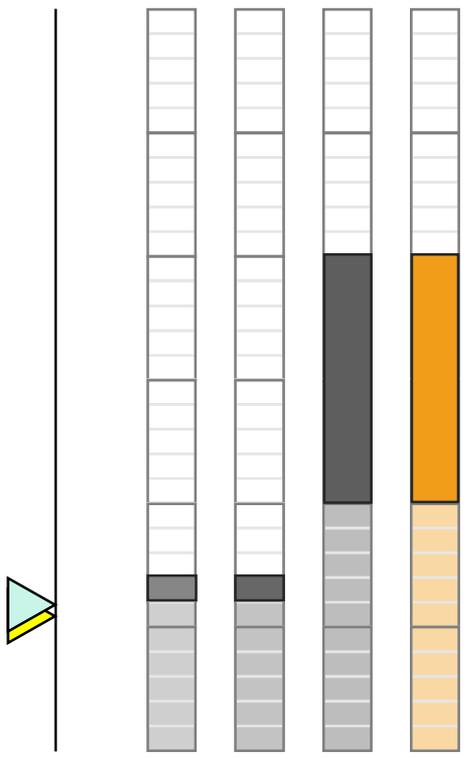
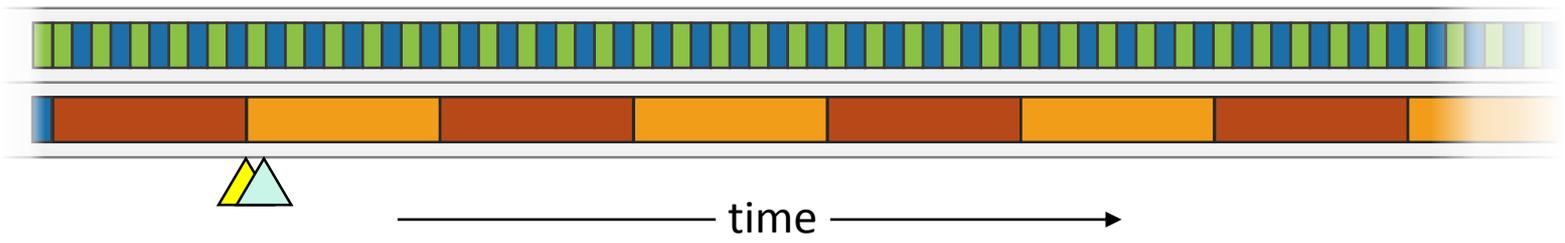


$v(t) = 5\frac{1}{2}$     6    6    10    0

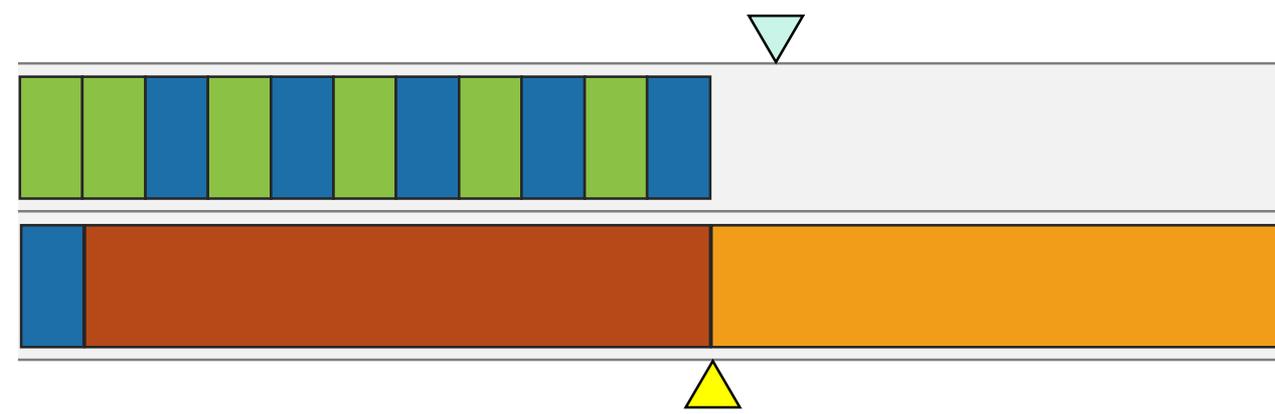


$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ

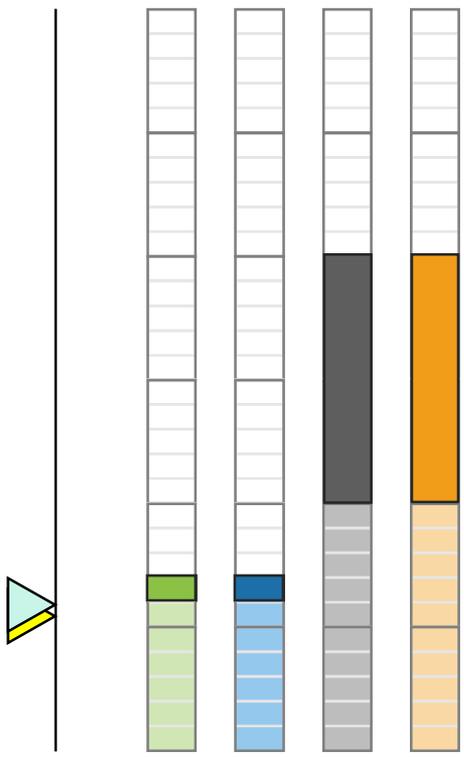
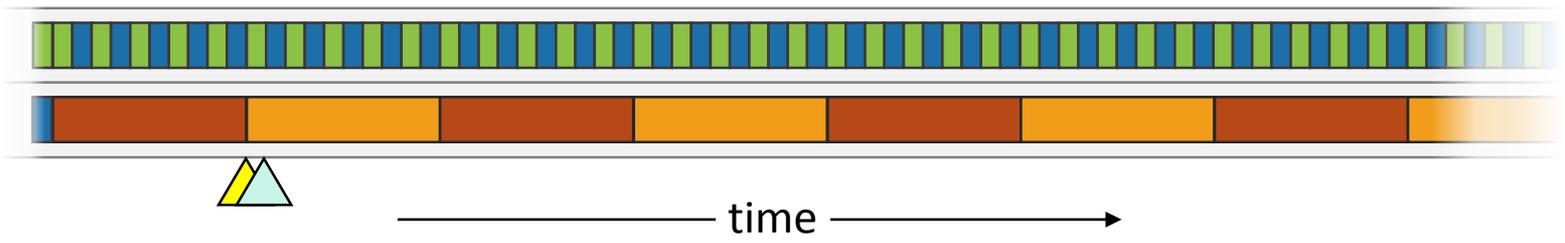


$v(t) = 5\frac{1}{2}$

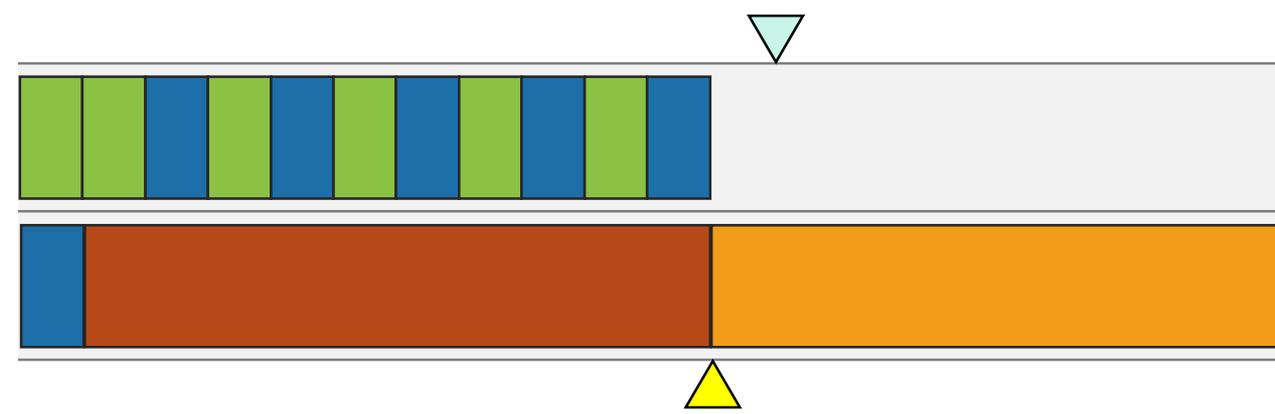


$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ

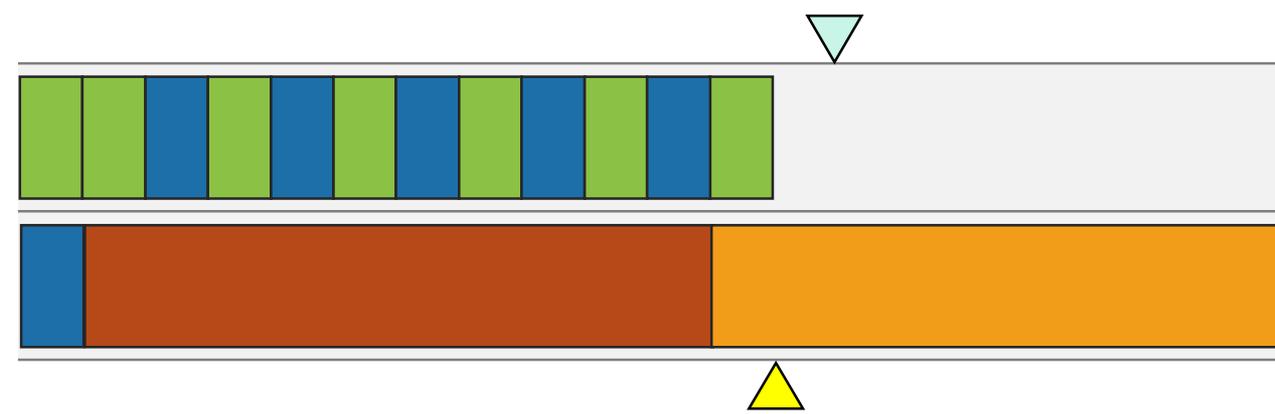
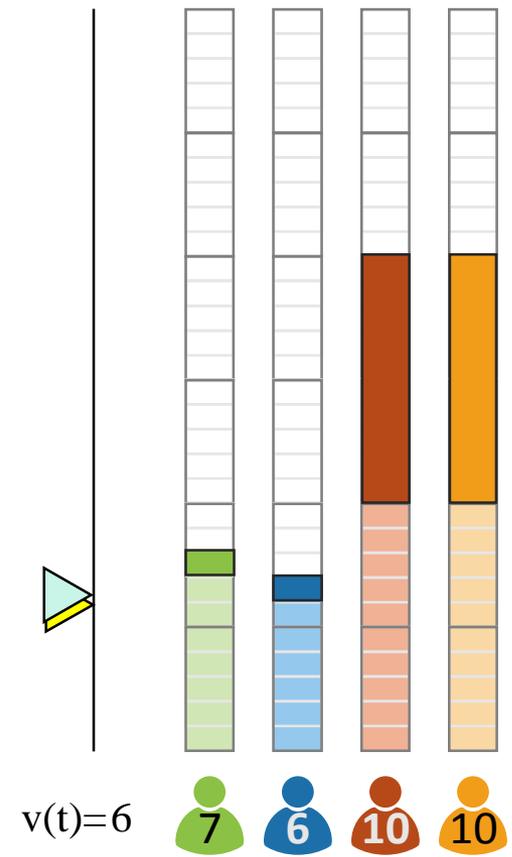
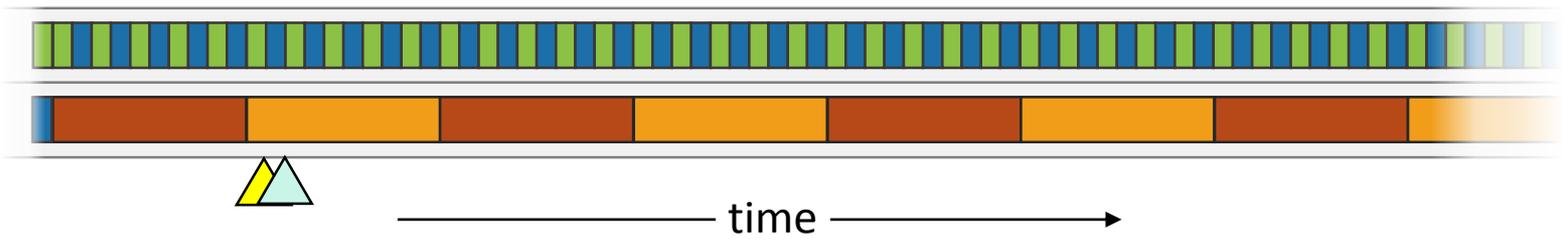


$v(t) = 5\frac{1}{2}$  6 6 10 0



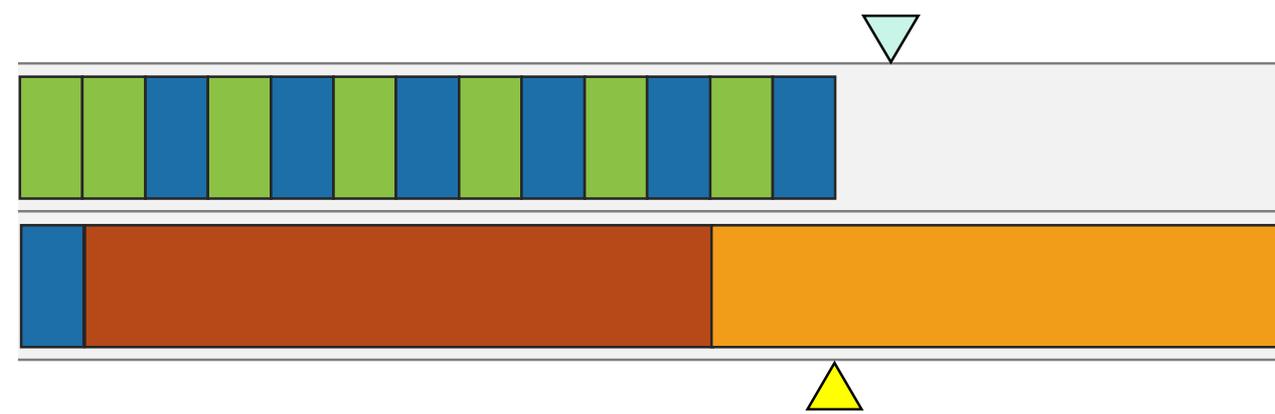
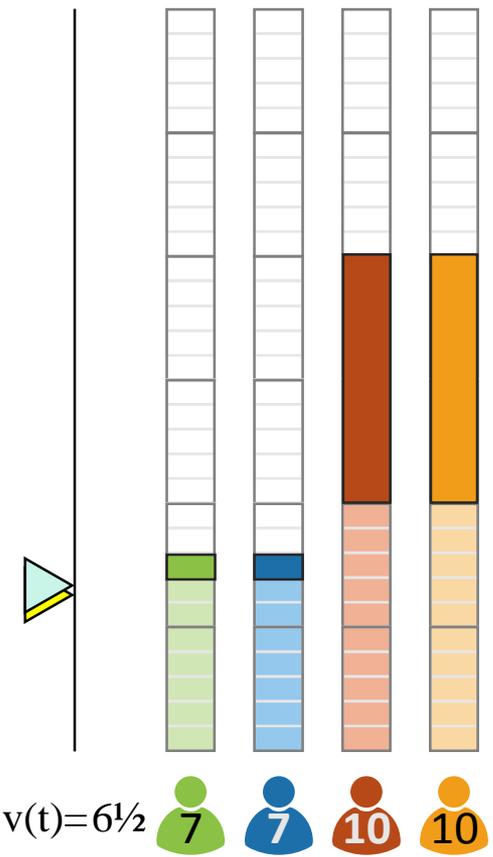
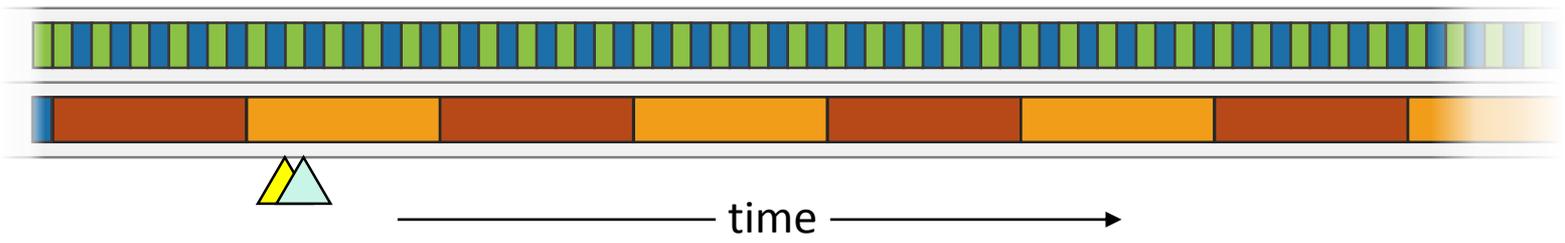
$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ



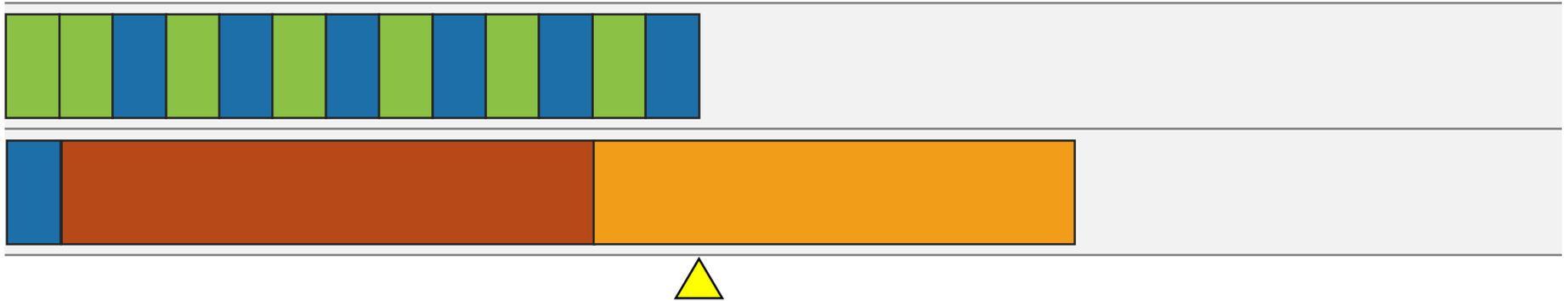
$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

# 2DFQ

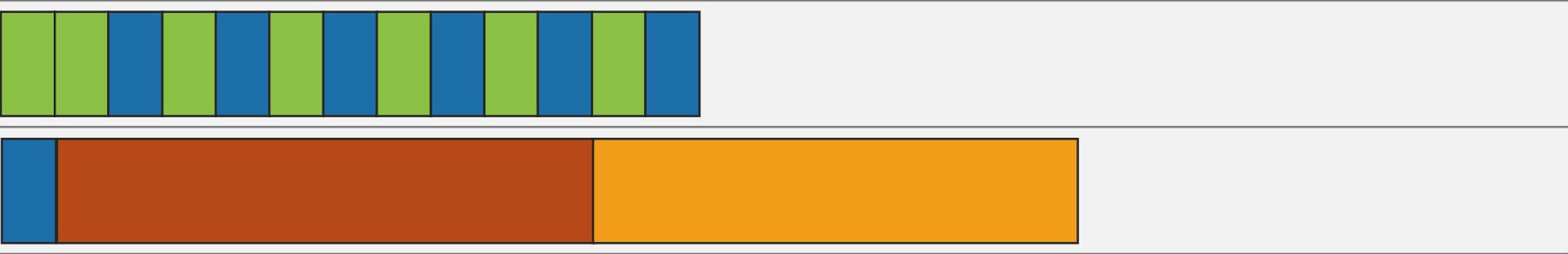


$$\text{Eligible}(r^j, i) = \text{Start}(r^j) - \frac{i}{n} \times \text{size}(r^j)$$

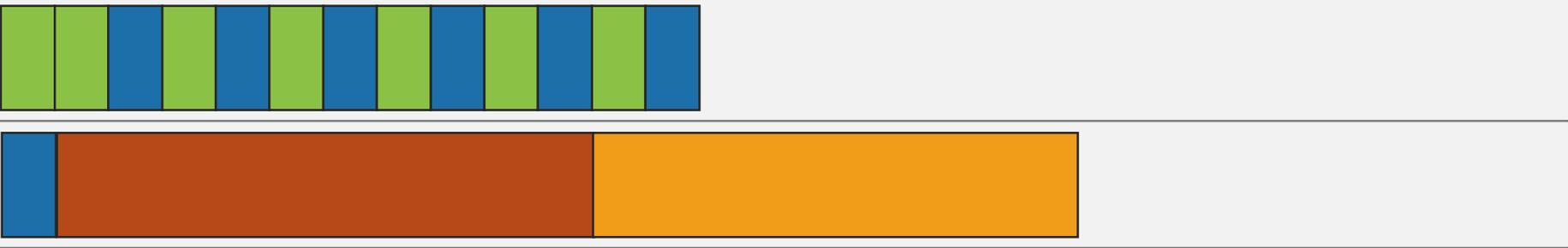
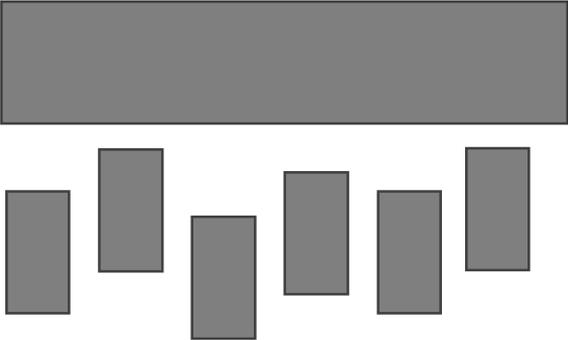
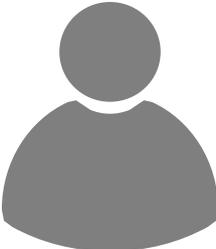
# Unknown Costs



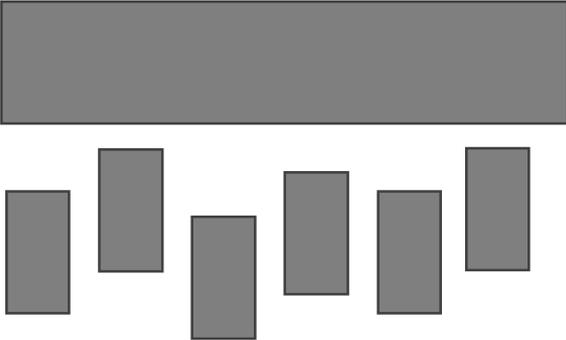
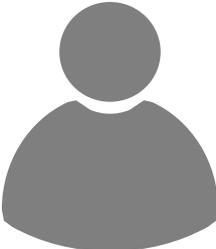
# Unknown Costs



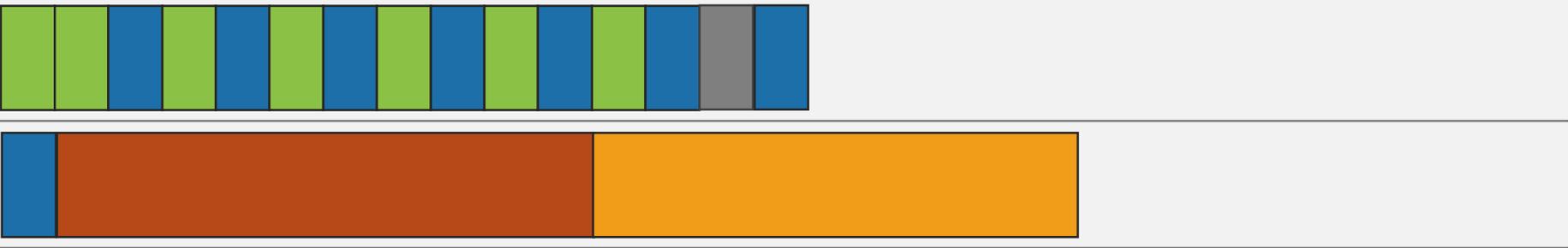
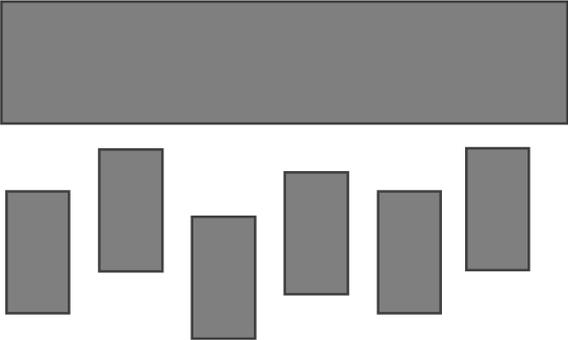
# Unknown Costs



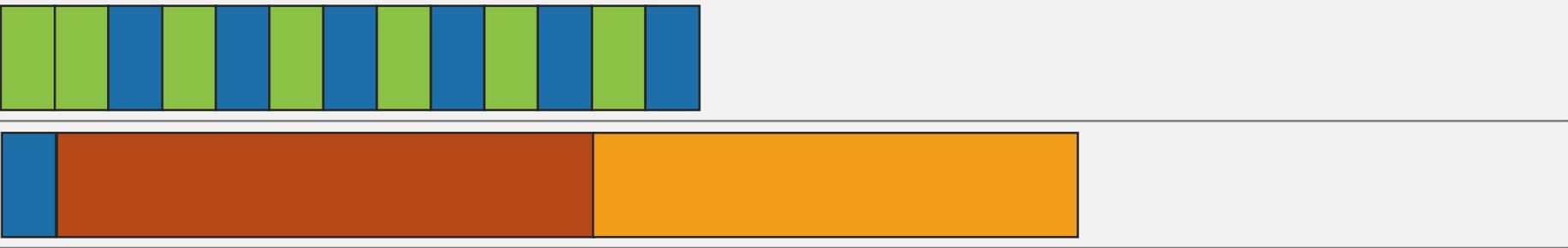
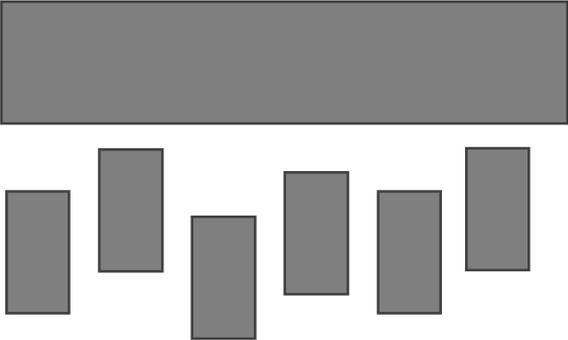
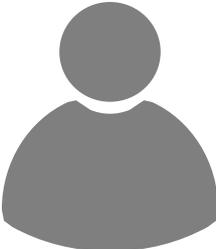
# Unknown Costs



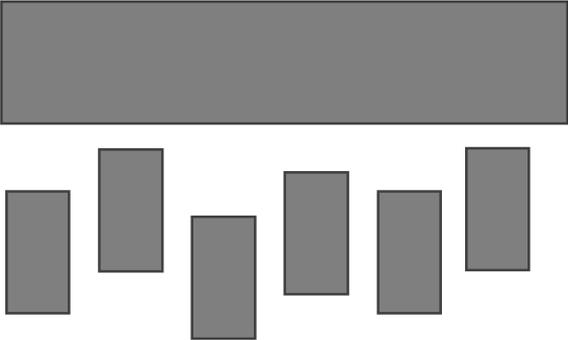
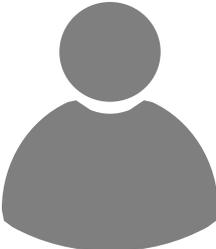
# Unknown Costs



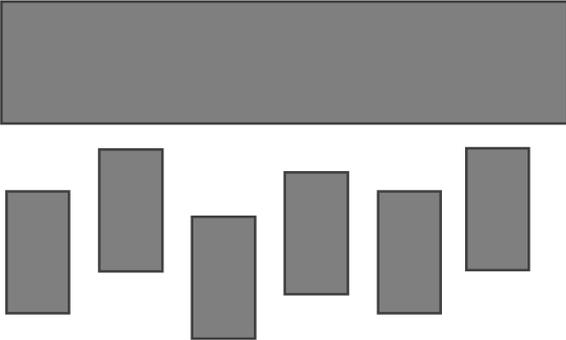
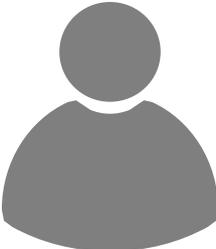
# Unknown Costs



# Unknown Costs



# Unknown Costs



Pessimistic cost estimation

# Evaluation

Compare 2DFQ to WFQ and WF<sup>2</sup>Q

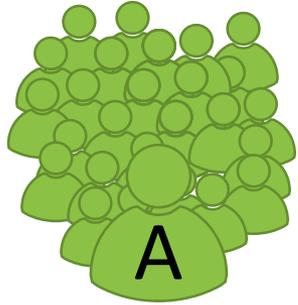
Discrete event simulator with Azure Storage workloads

More experiment results in the paper, evaluating:

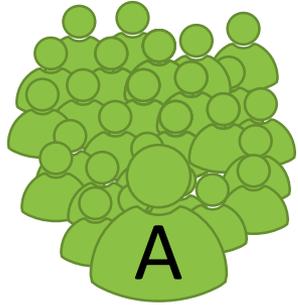
- Burstiness
- Fairness
- Tail latency



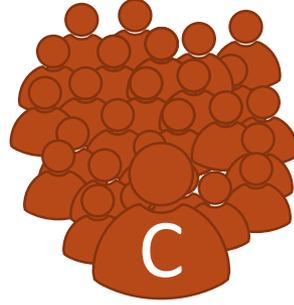
50 tenants with size  $\approx 1$



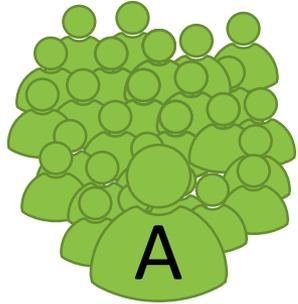
50 tenants with size  $\approx 1$



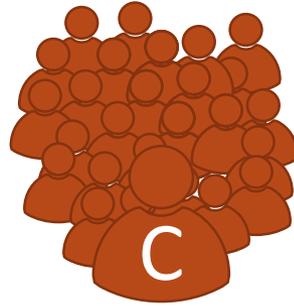
50 tenants with size  $\approx 1000$



50 tenants with size  $\approx 1$



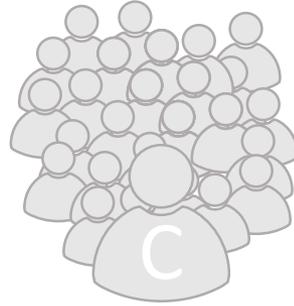
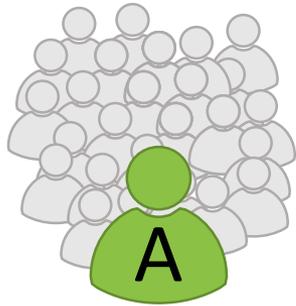
50 tenants with size  $\approx 1000$



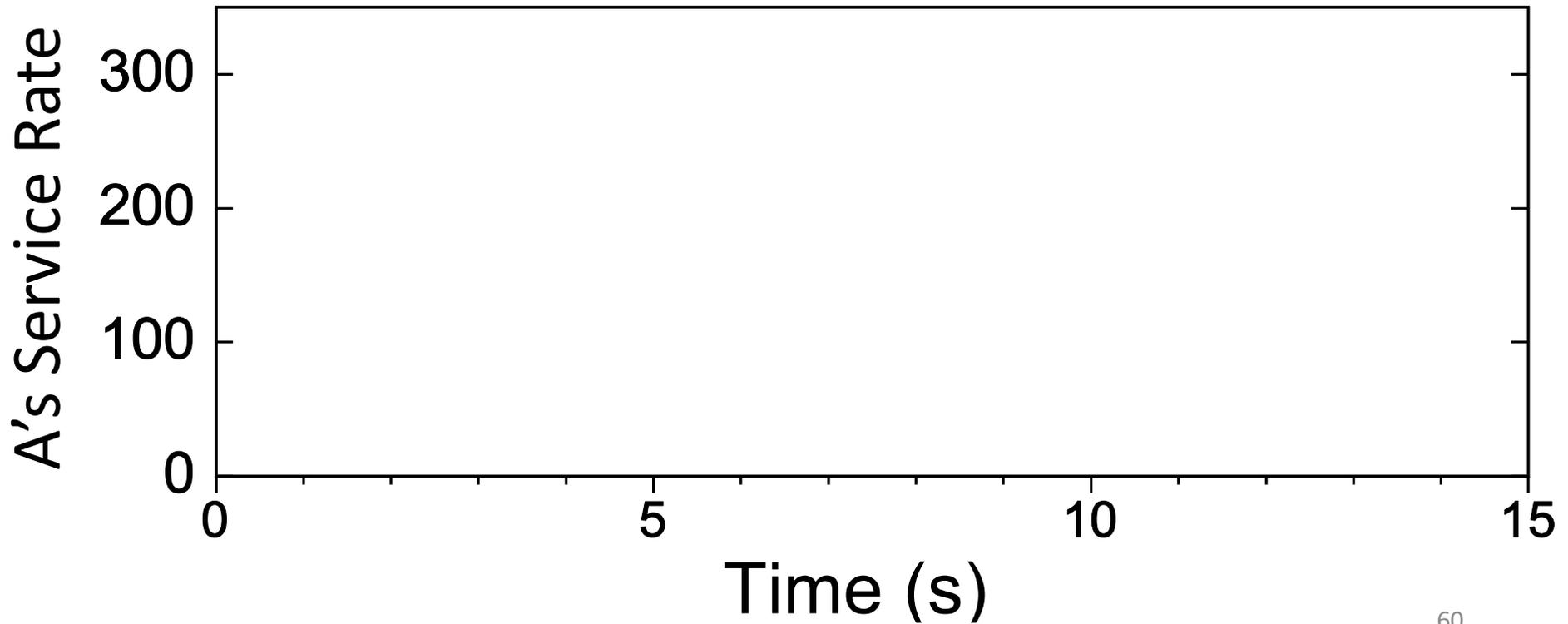
16 threads  
1000 units/second  
Costs known by scheduler

50 tenants with size  $\approx 1$

50 tenants with size  $\approx 1000$

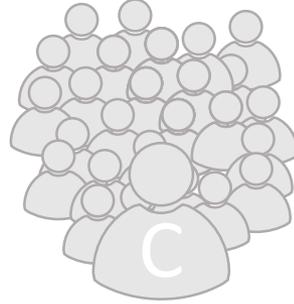
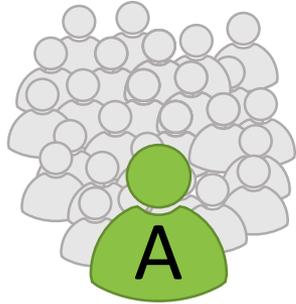


16 threads  
1000 units/second  
Costs known by scheduler



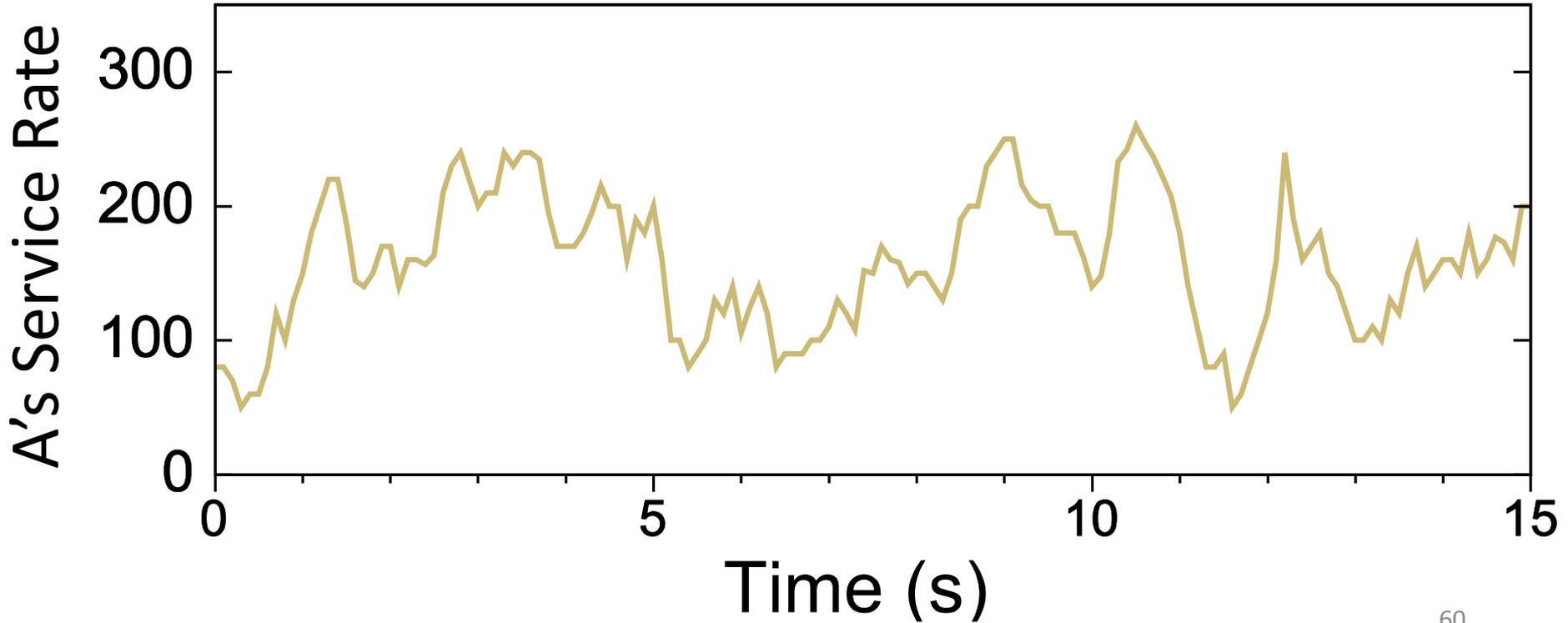
50 tenants with size  $\approx 1$

50 tenants with size  $\approx 1000$



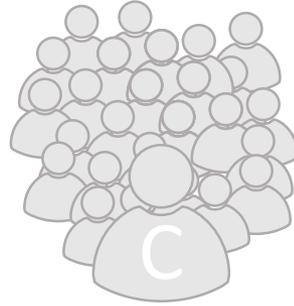
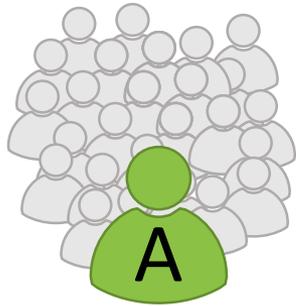
16 threads  
1000 units/second  
Costs known by scheduler

WFQ ———



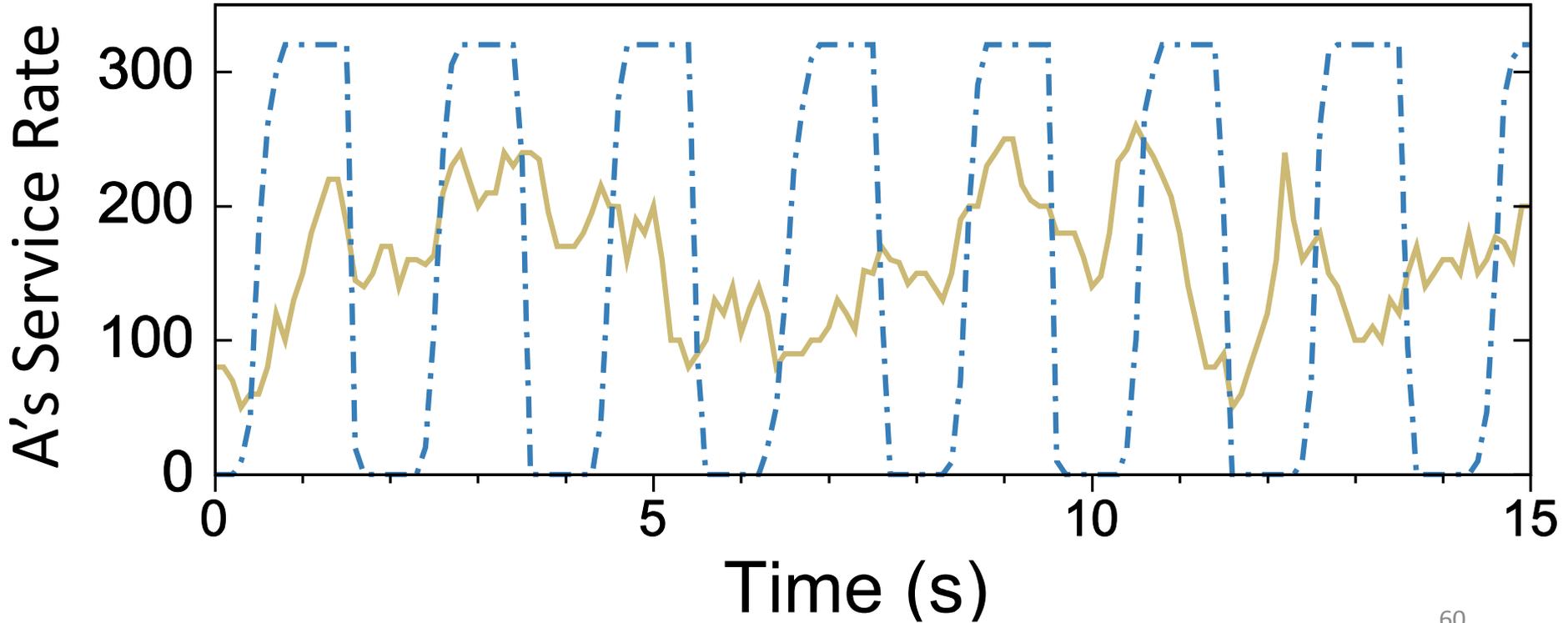
50 tenants with size  $\approx 1$

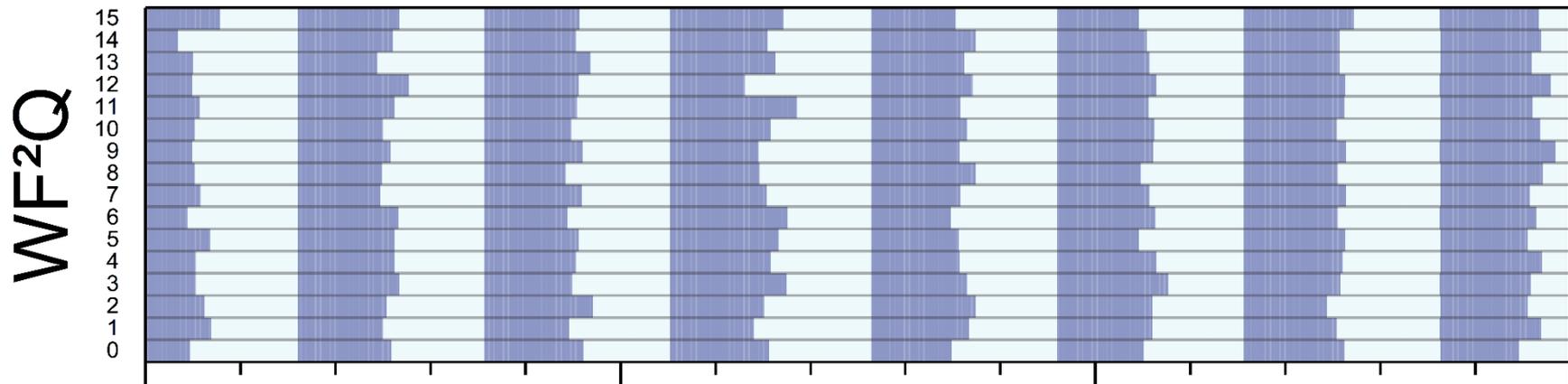
50 tenants with size  $\approx 1000$



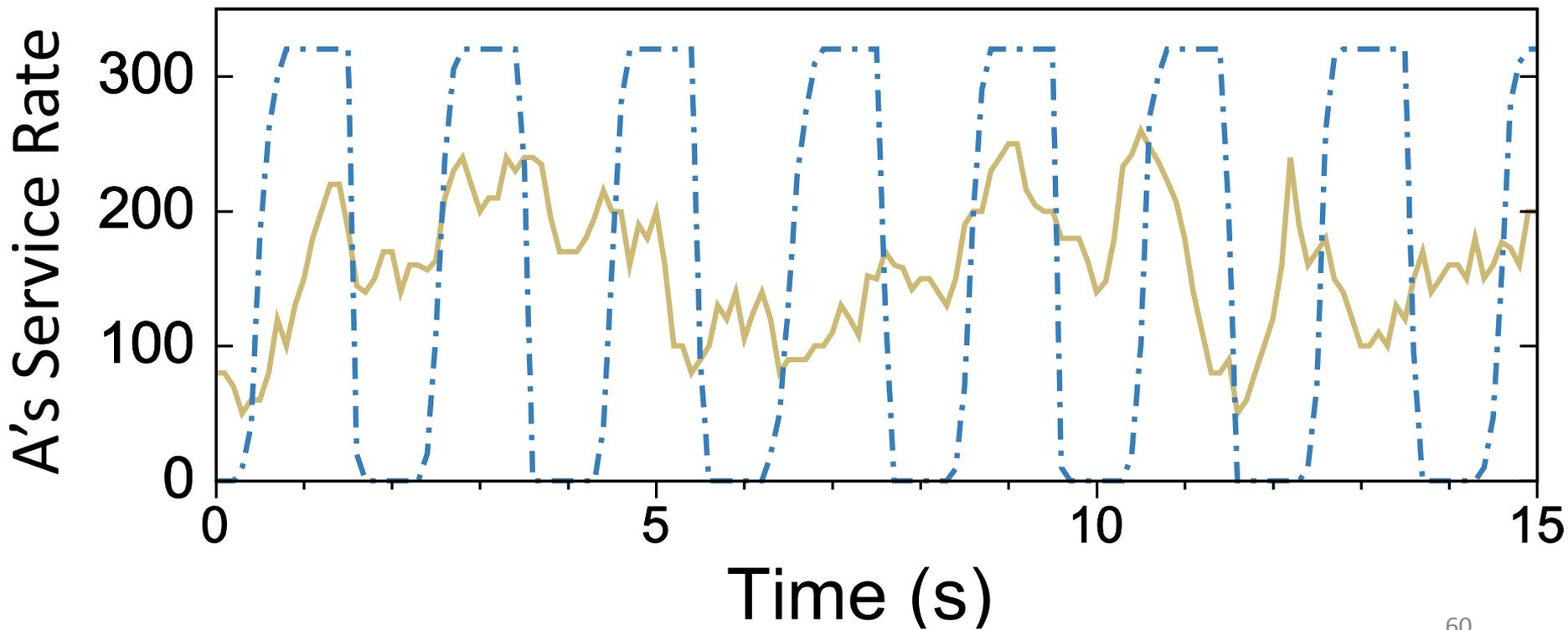
16 threads  
1000 units/second  
Costs known by scheduler

WFQ ——— WF<sup>2</sup>Q - - - - -



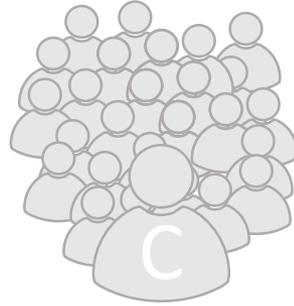
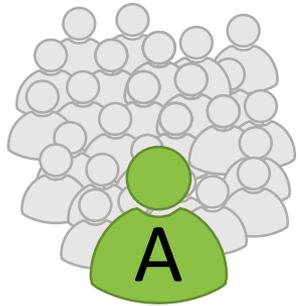


WFQ ——— WF<sup>2</sup>Q - - - - -



50 tenants with size  $\approx 1$

50 tenants with size  $\approx 1000$



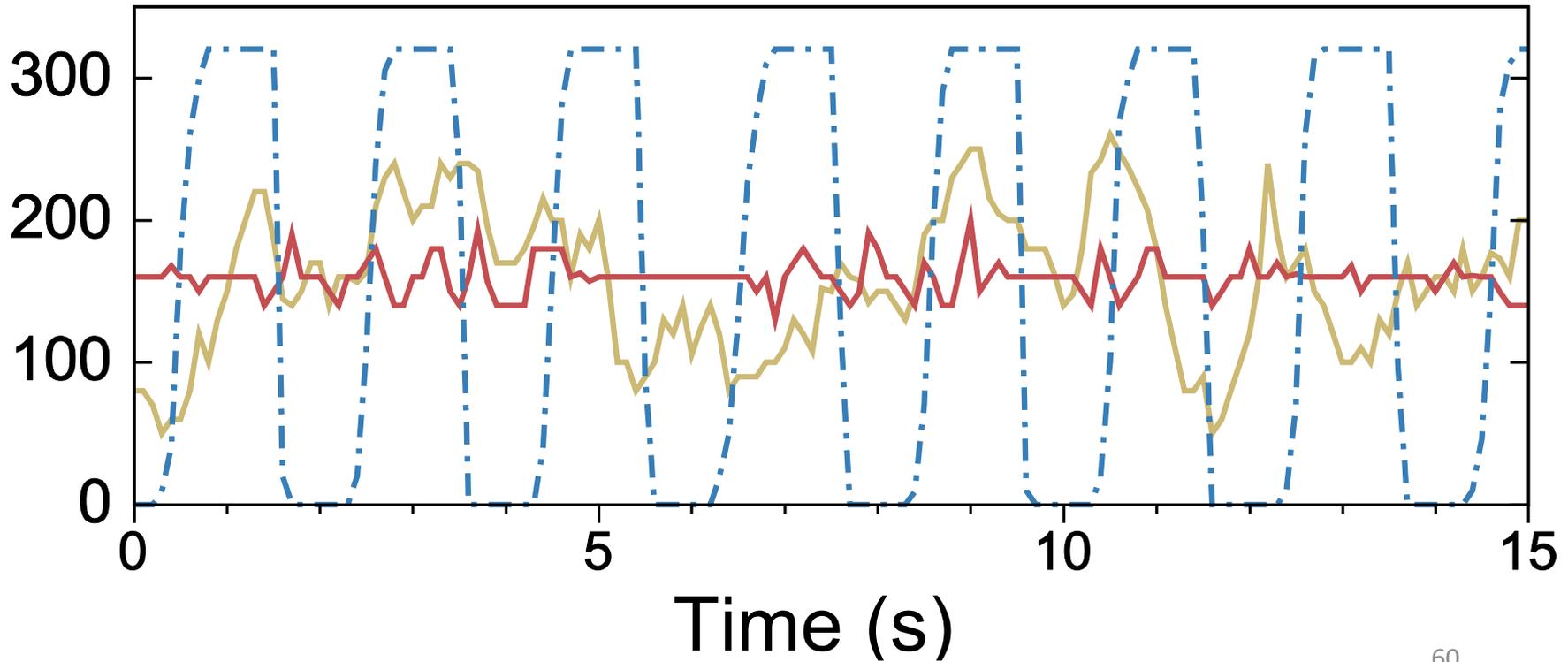
16 threads  
1000 units/second  
Costs known by scheduler

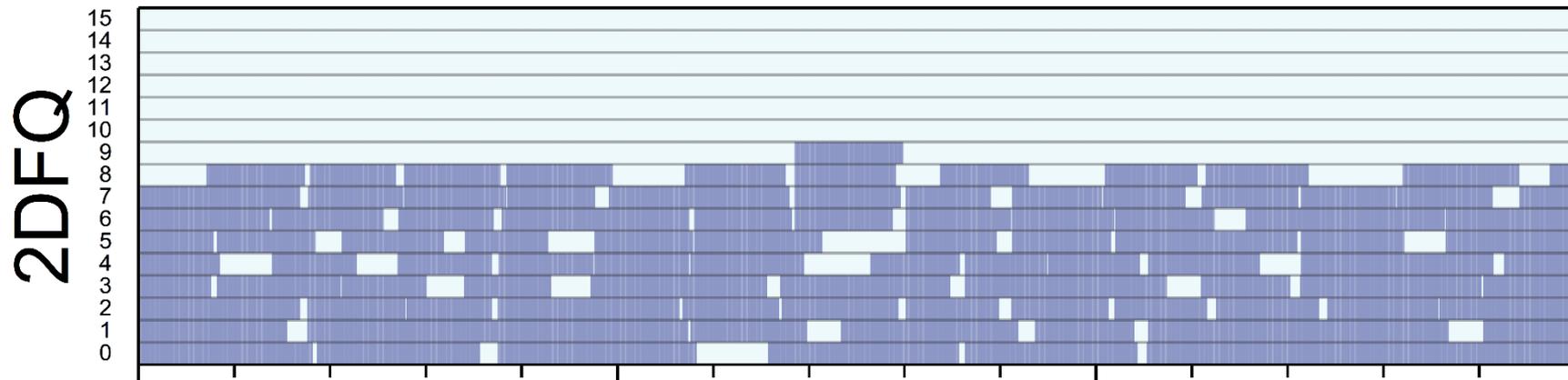
WFQ

WF<sup>2</sup>Q

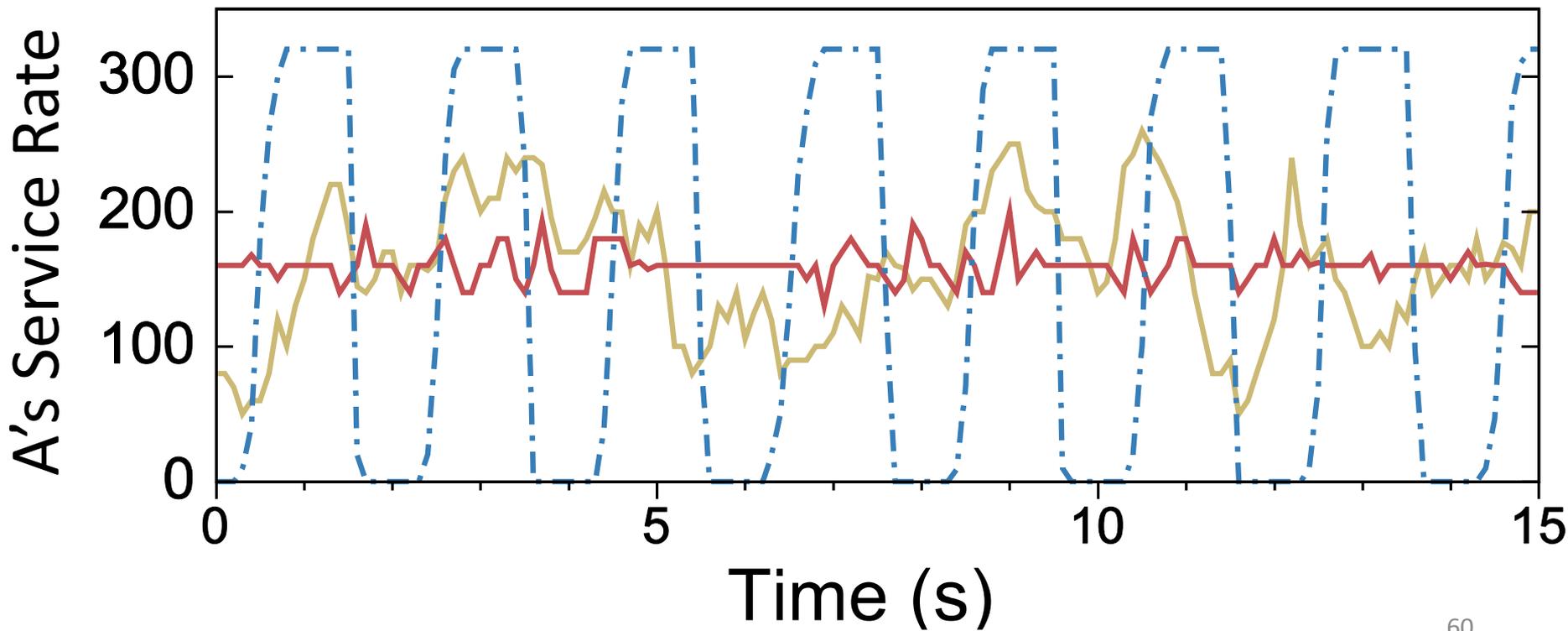
2DFQ

A's Service Rate



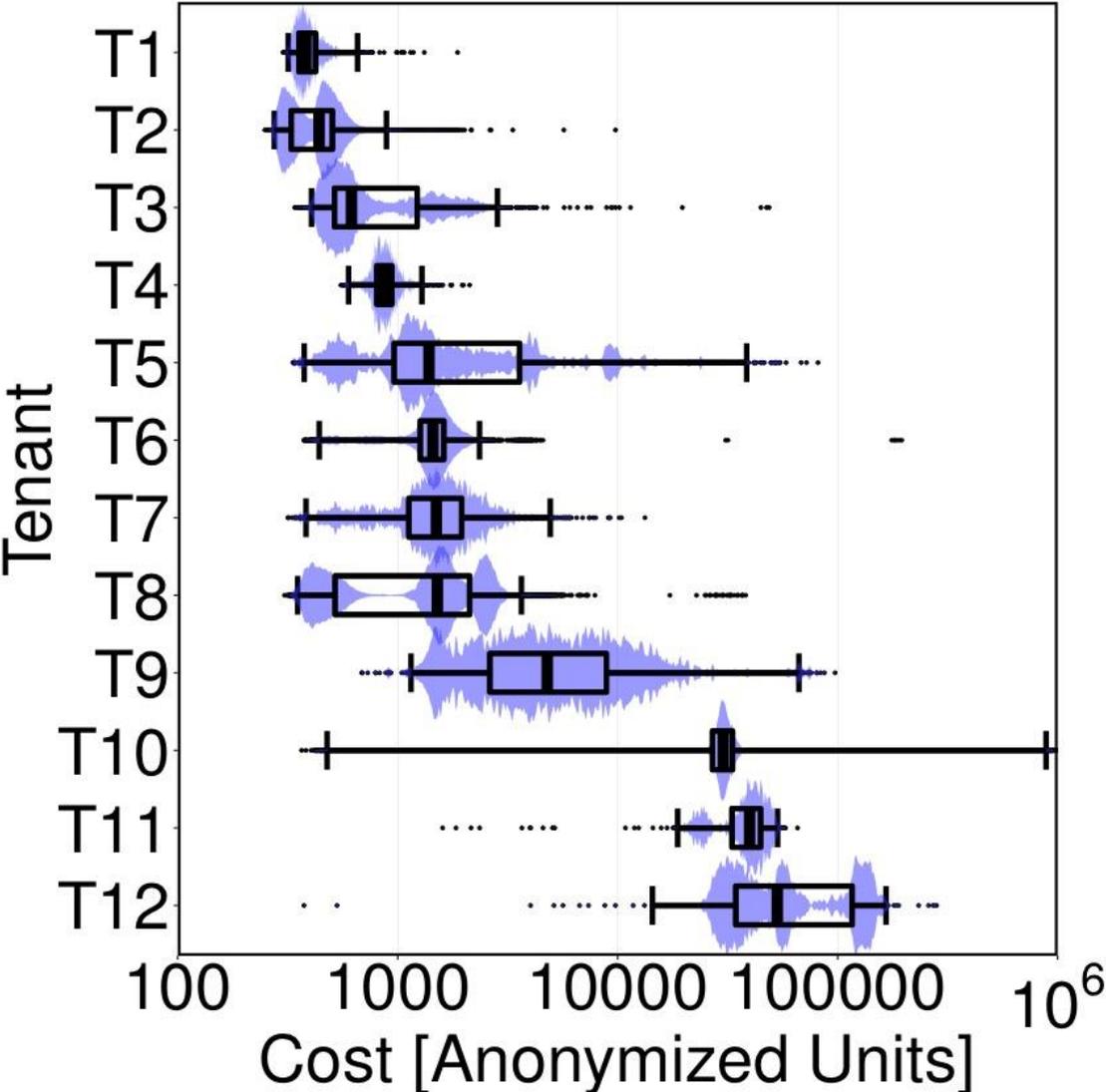


WFQ ——— WF<sup>2</sup>Q - - - - 2DFQ ———

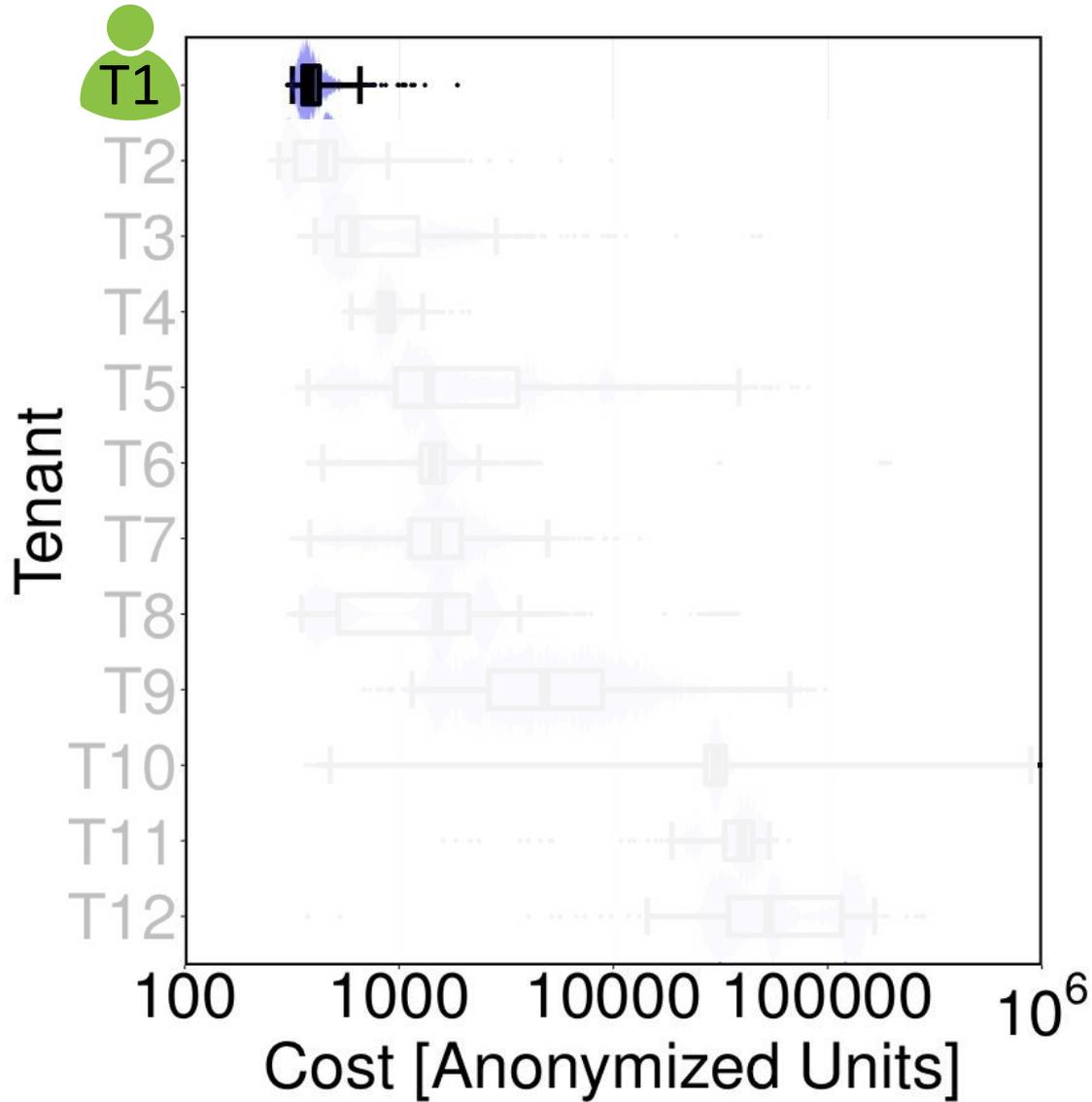


250 Azure Storage tenants  
32 threads  
1 million units/second  
Costs known by scheduler

250 Azure Storage tenants  
32 threads  
1 million units/second  
Costs known by scheduler

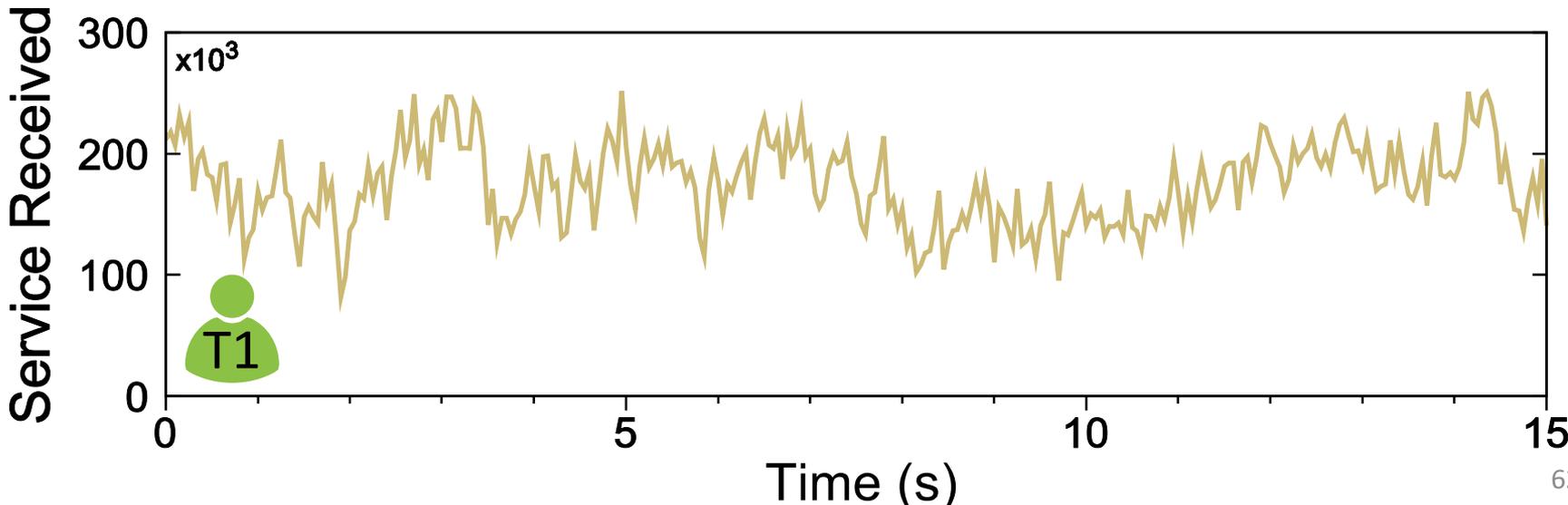


250 Azure Storage tenants  
32 threads  
1 million units/second  
Costs known by scheduler



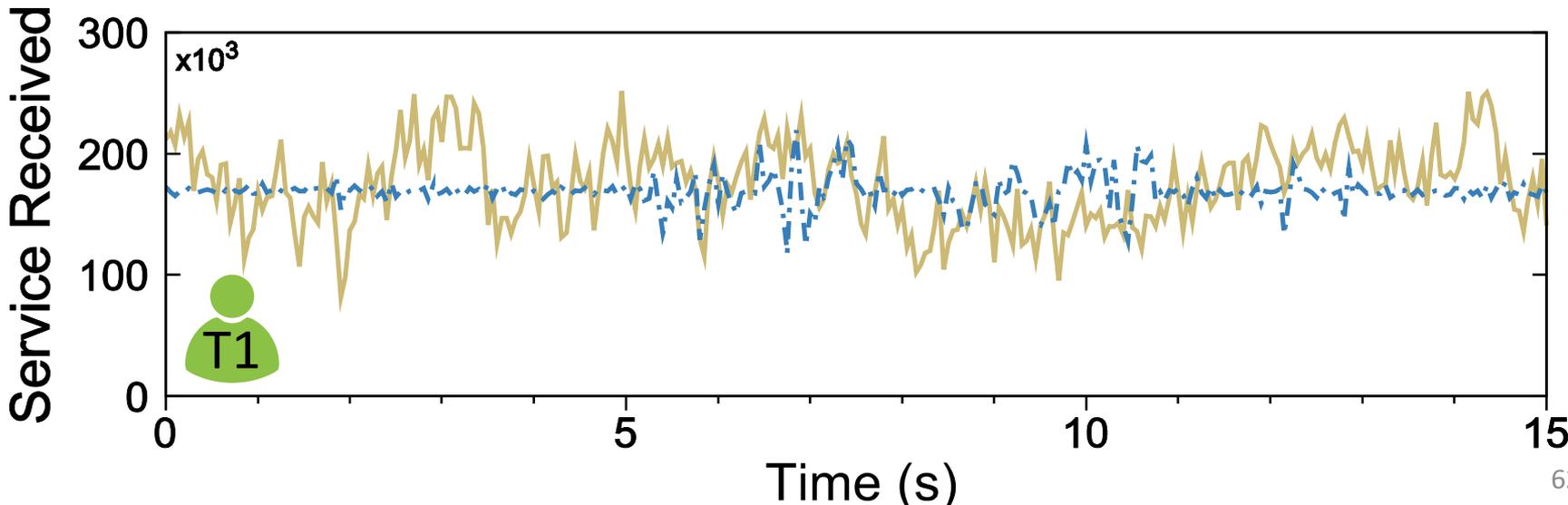
250 Azure Storage workloads  
32 threads  
1 million units/second  
Costs known by scheduler

WFQ



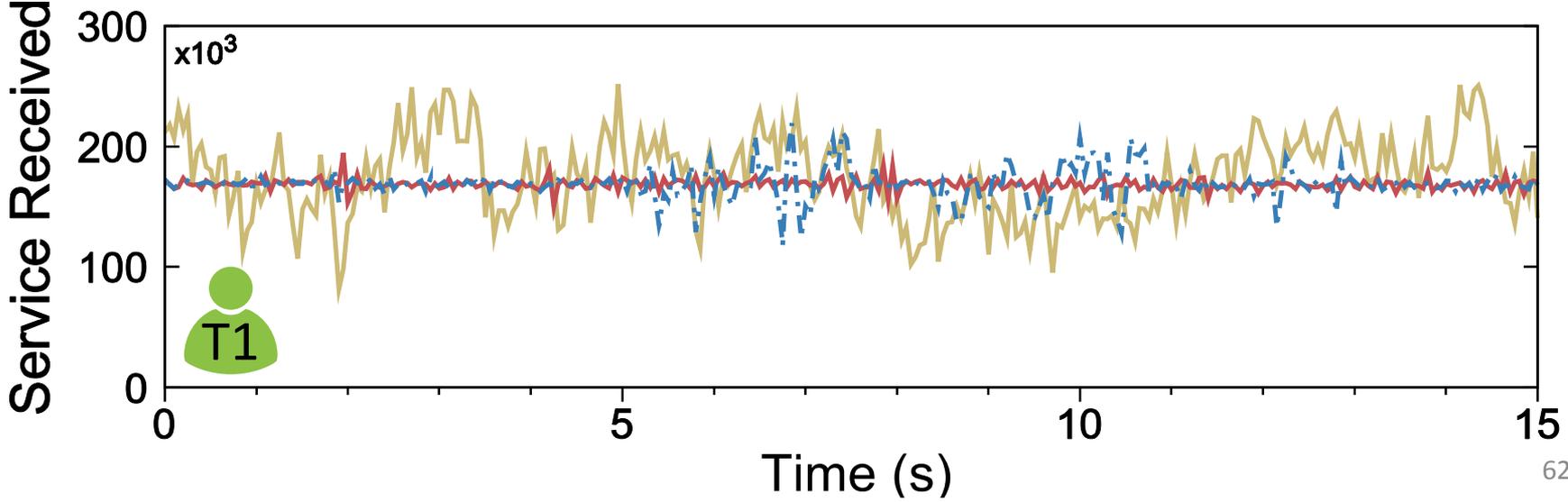
250 Azure Storage workloads  
32 threads  
1 million units/second  
Costs known by scheduler

WFQ — WF<sup>2</sup>Q - · - · -

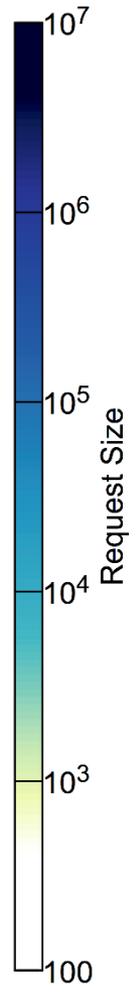


250 Azure Storage workloads  
32 threads  
1 million units/second  
Costs known by scheduler

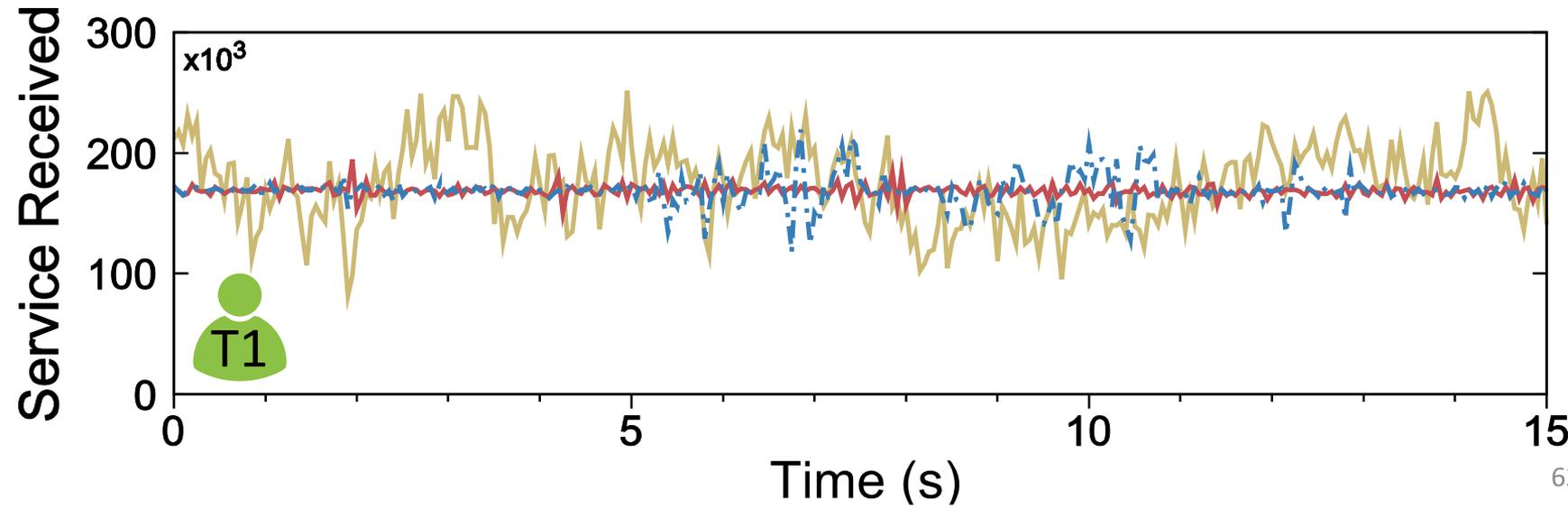
WFQ — WF<sup>2</sup>Q - · - · - 2DFQ —



250 Azure Storage workloads  
32 threads  
1 million units/second  
Costs known by scheduler

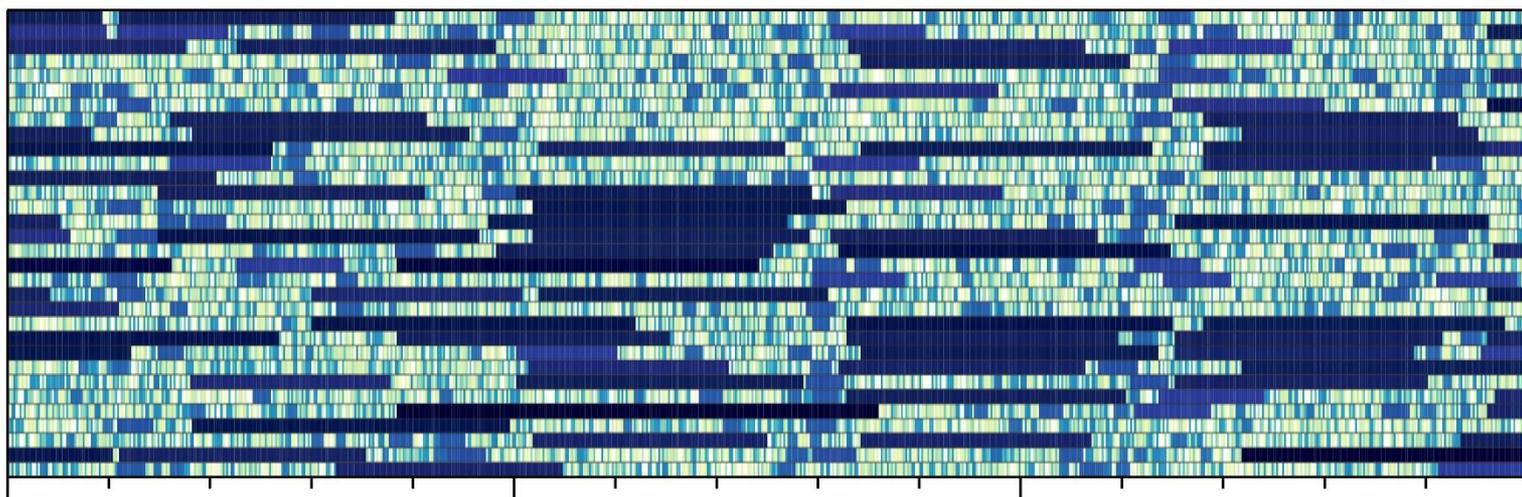


WFQ — WF<sup>2</sup>Q - · - · - 2DFQ —



WF<sup>2</sup>Q

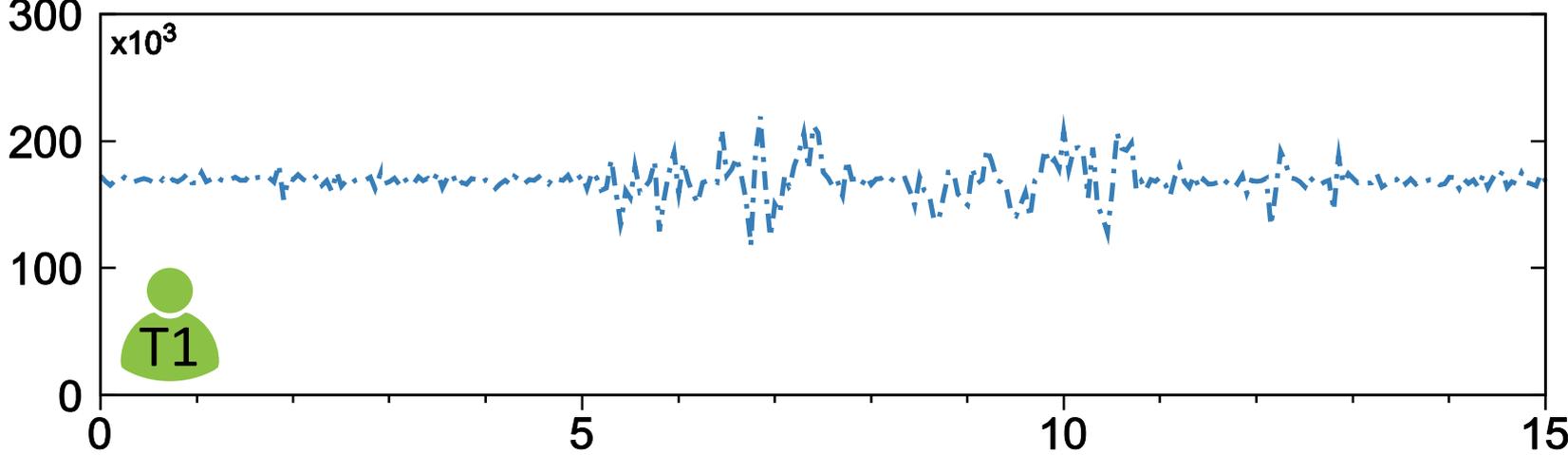
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62



Request Size

Service Received

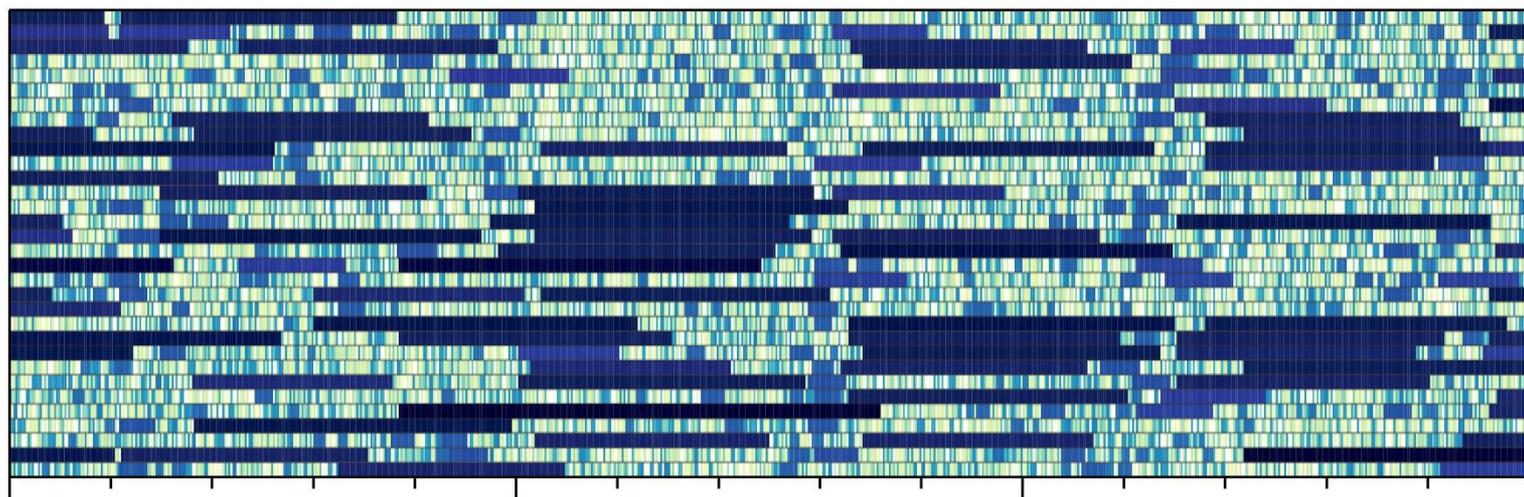
WFQ — WF<sup>2</sup>Q - - - 2DFQ —



Time (s)

WF<sup>2</sup>Q

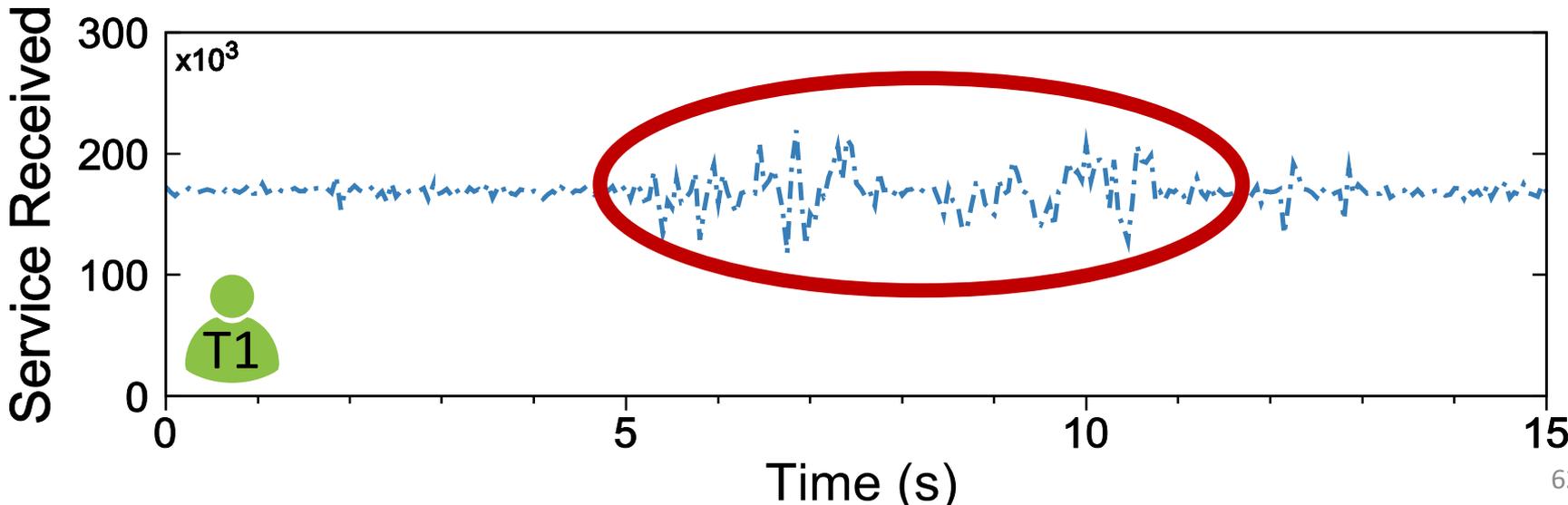
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62



Request Size

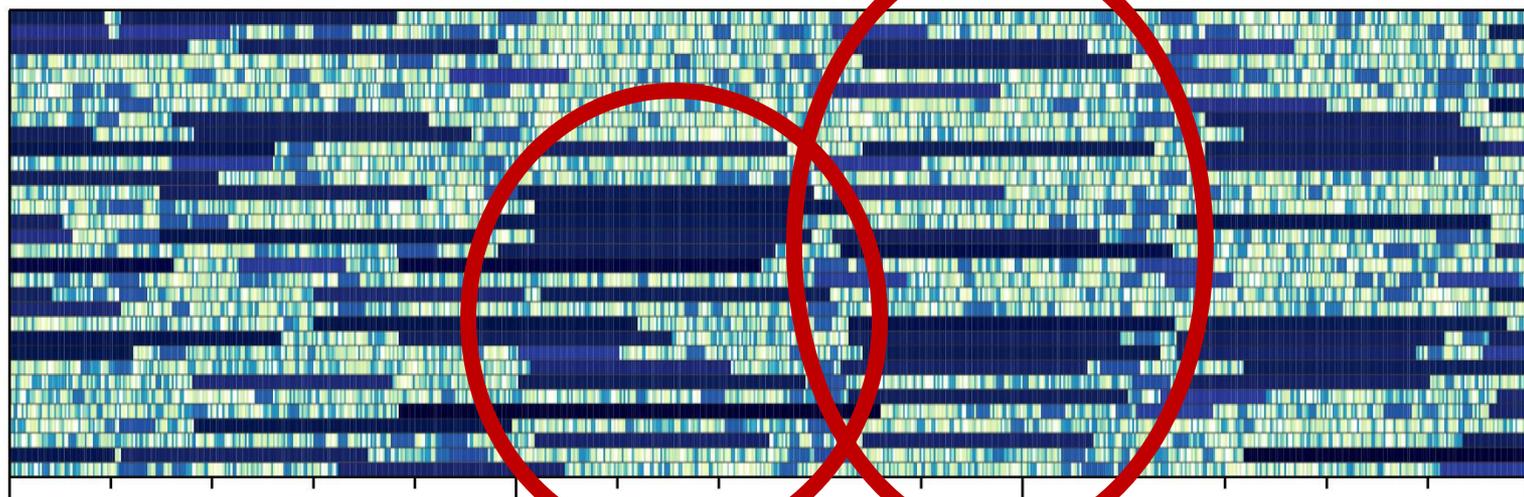
10<sup>7</sup>  
10<sup>6</sup>  
10<sup>5</sup>  
10<sup>4</sup>  
10<sup>3</sup>  
100

WFQ — WF<sup>2</sup>Q - - - 2DFQ —

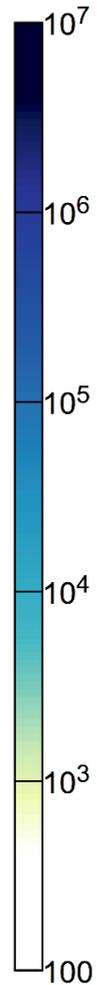


WF<sup>2</sup>Q

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62

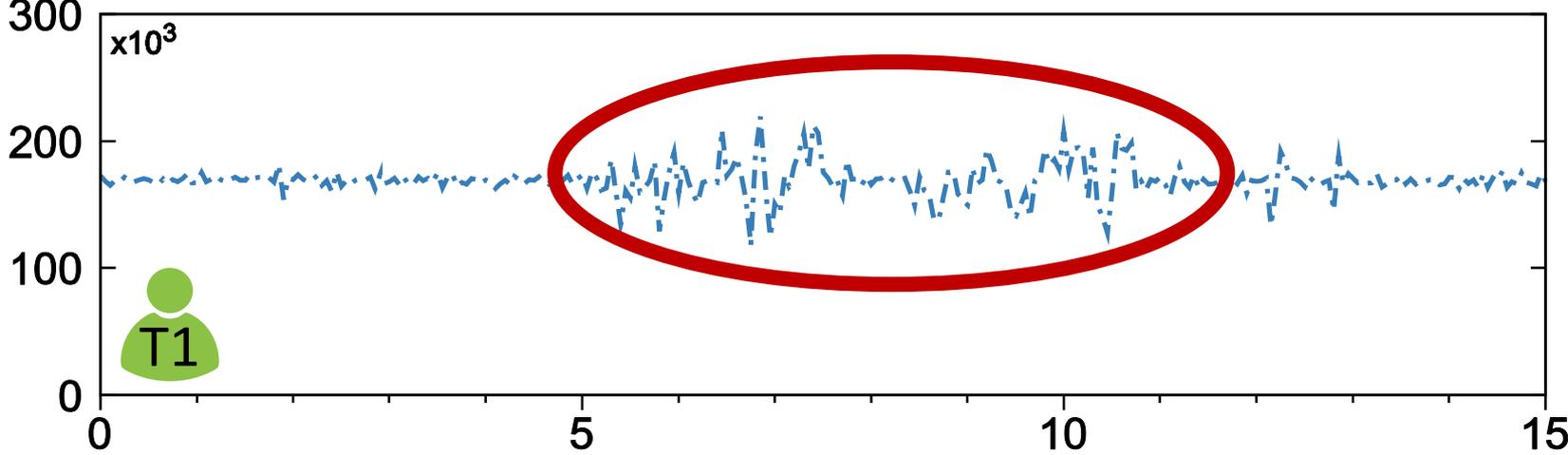


Request Size



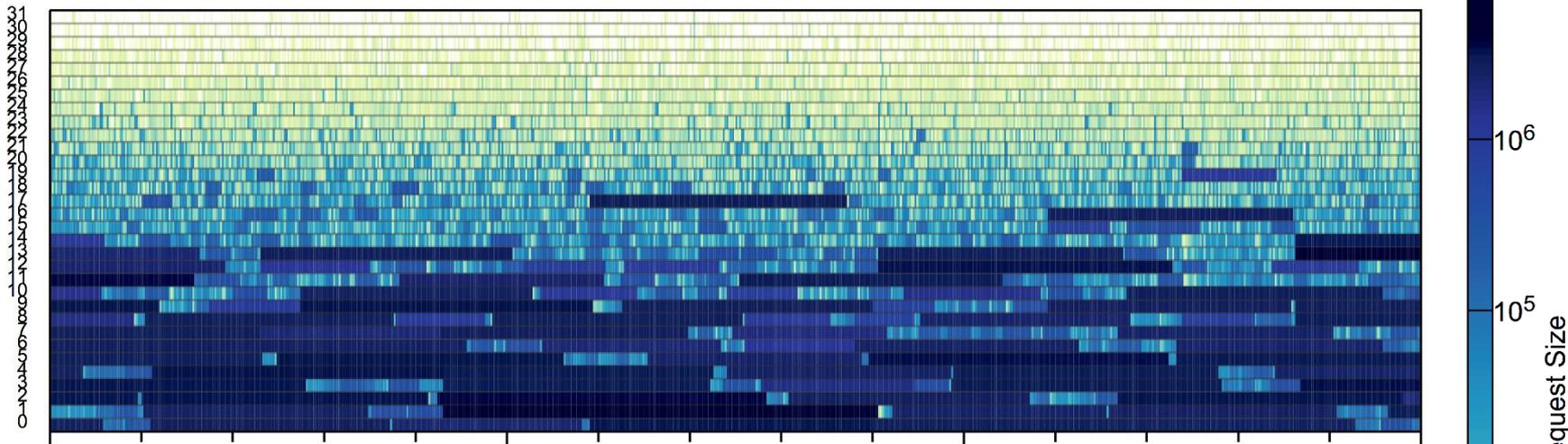
Service Received

WFQ — WF<sup>2</sup>Q - - - 2DFQ —

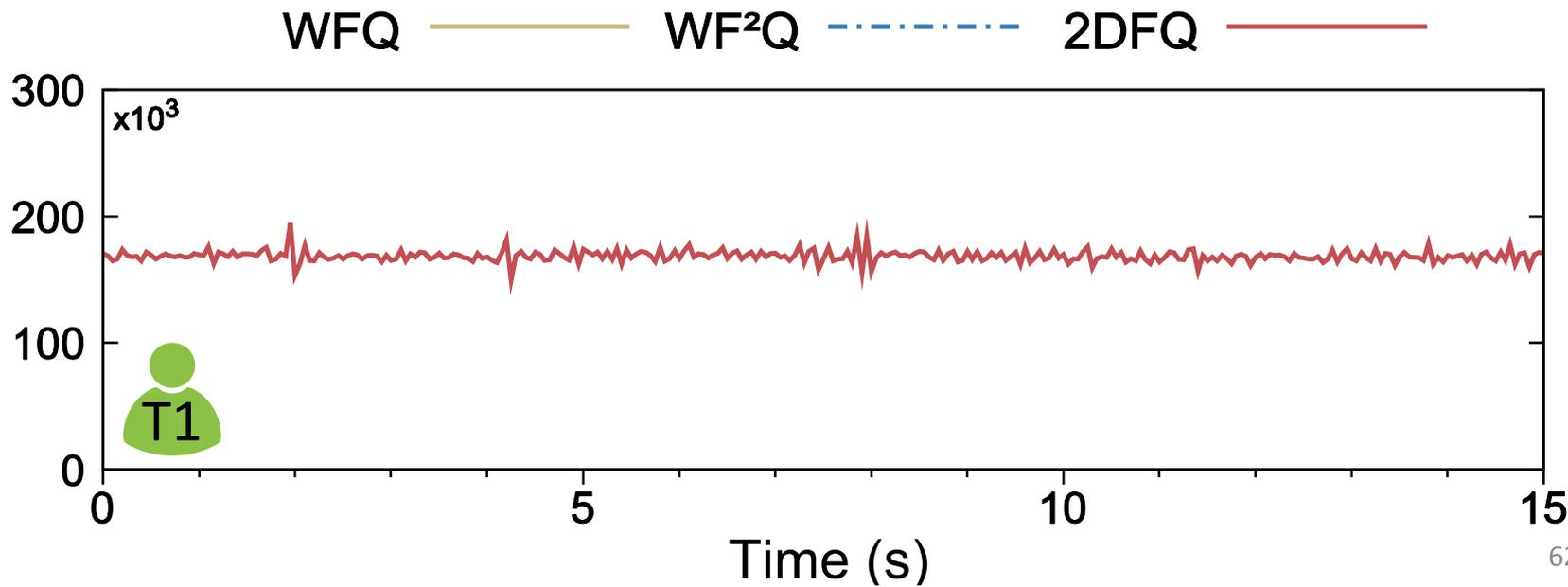


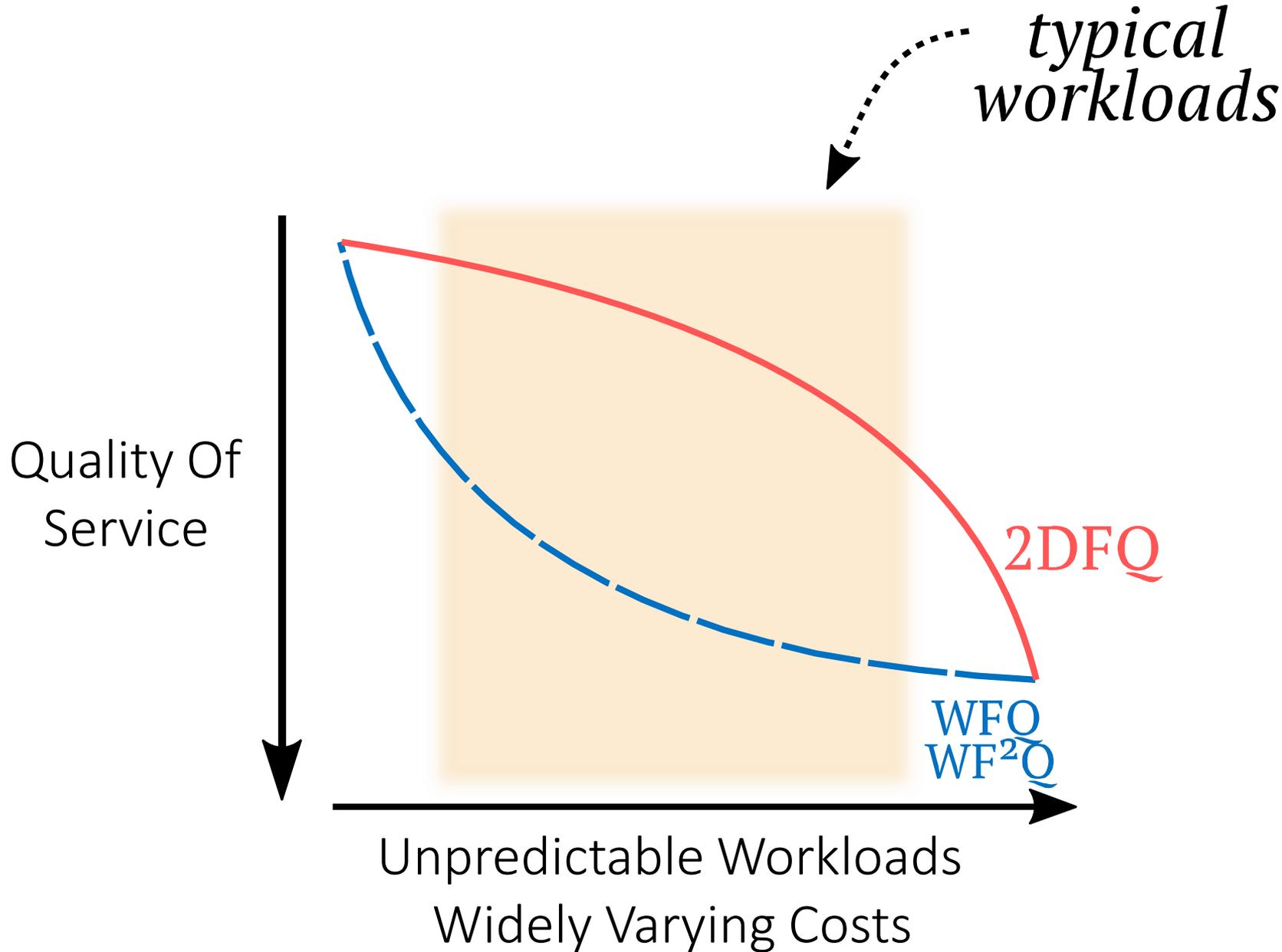
Time (s)

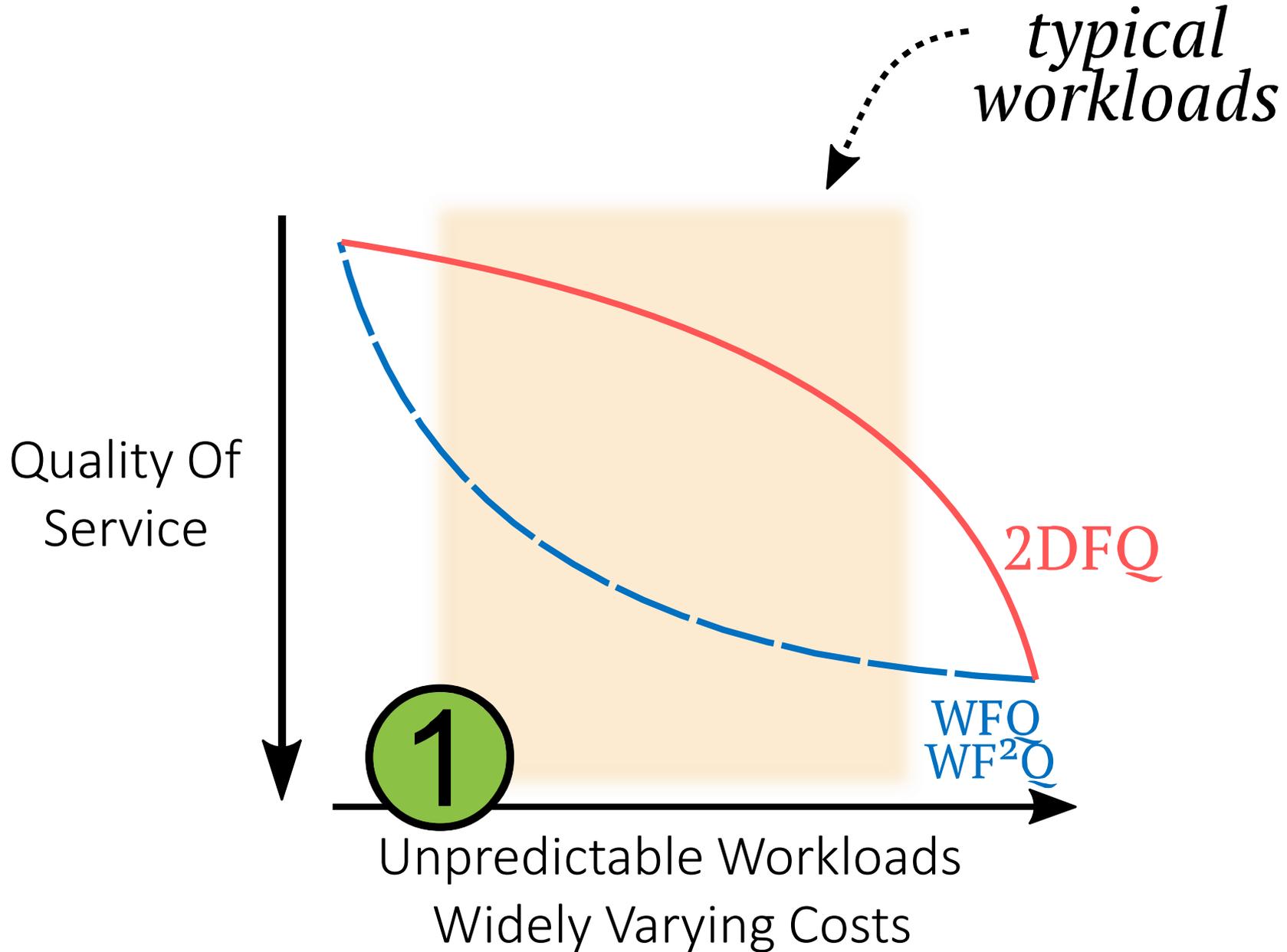
2DFQ

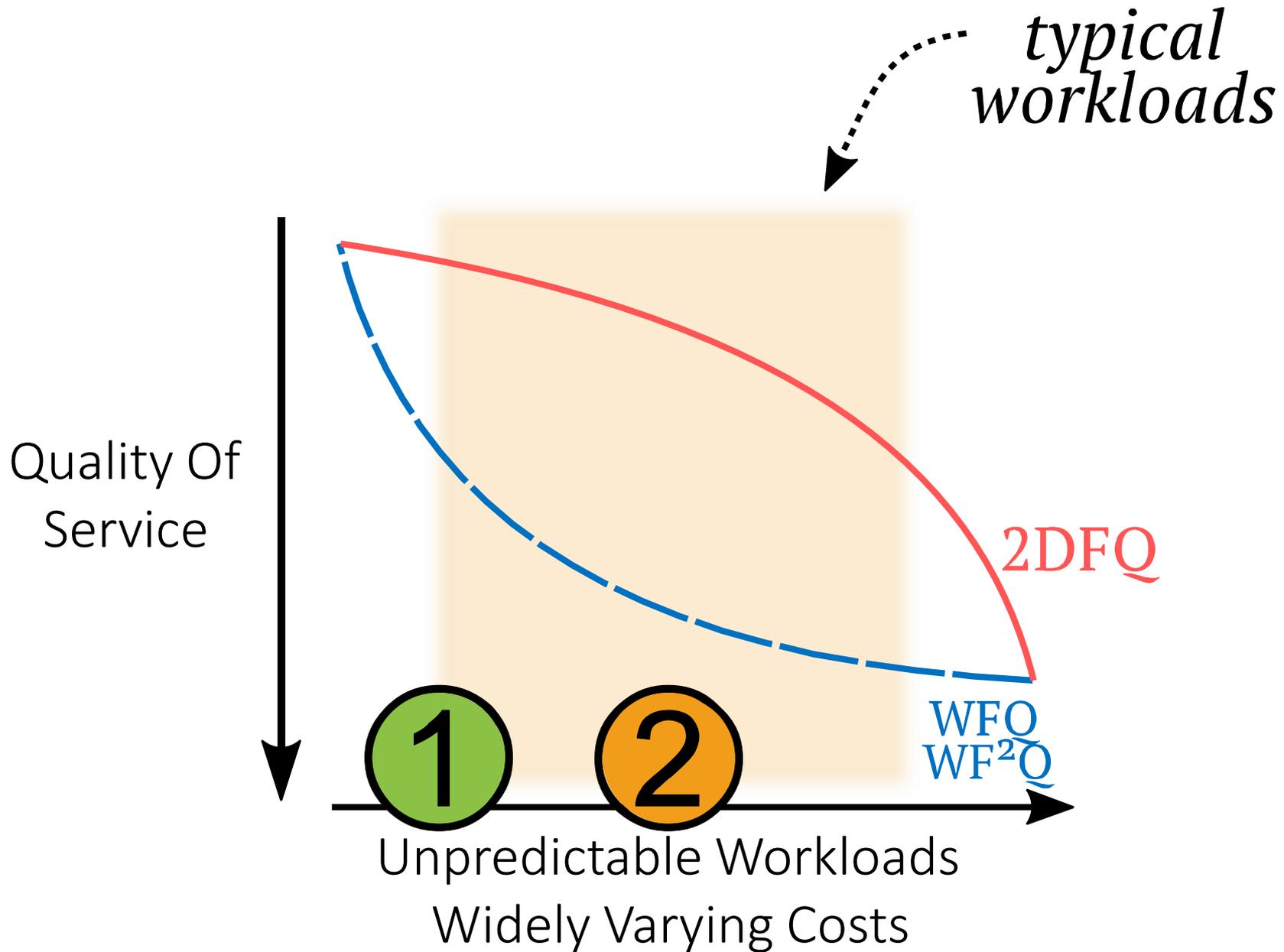


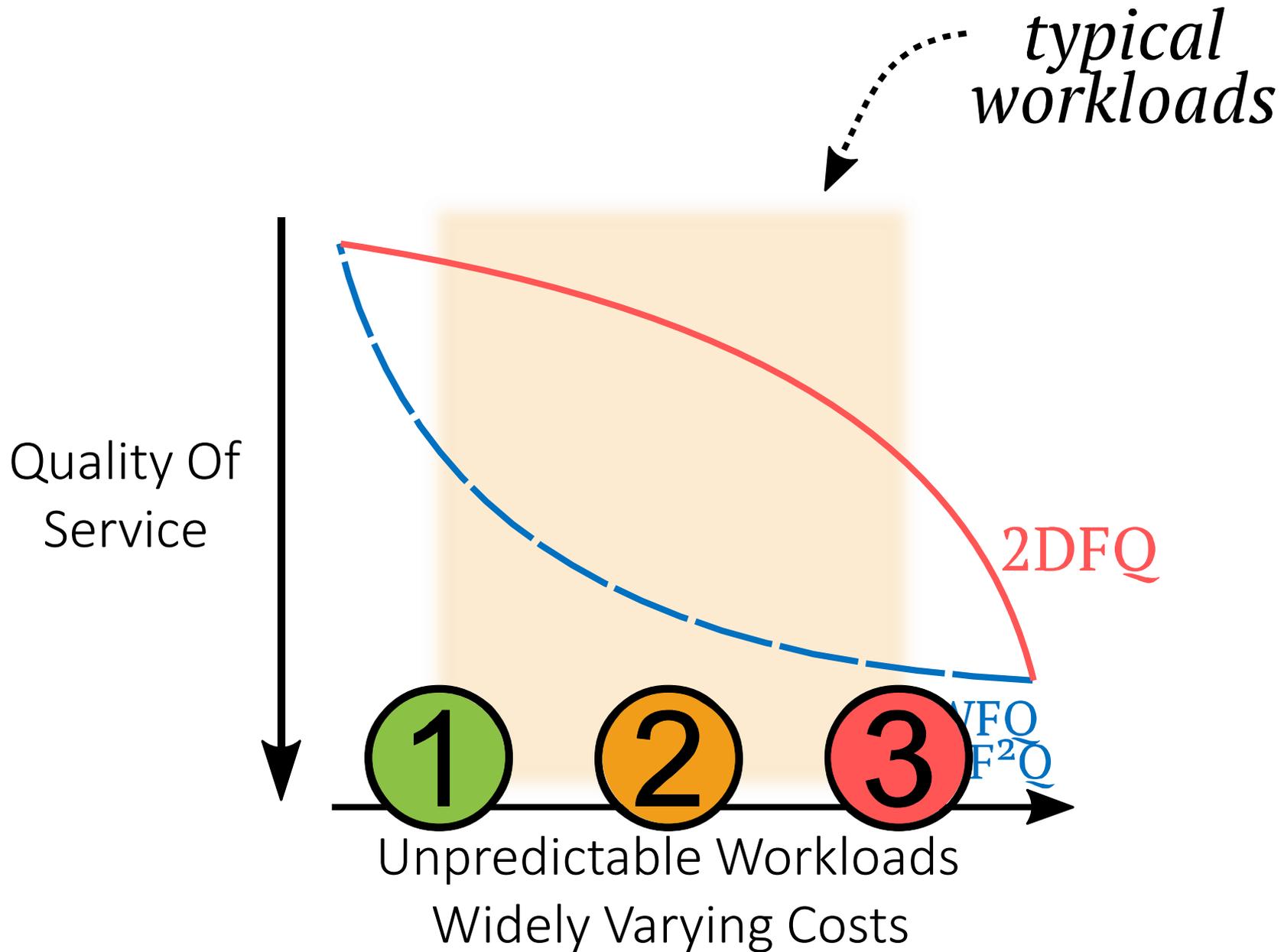
Service Received

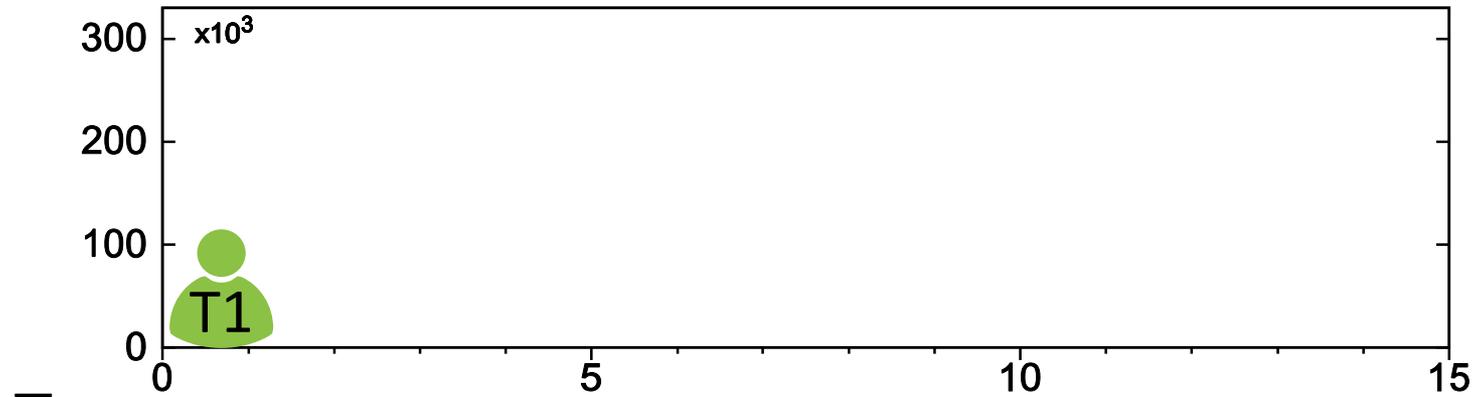




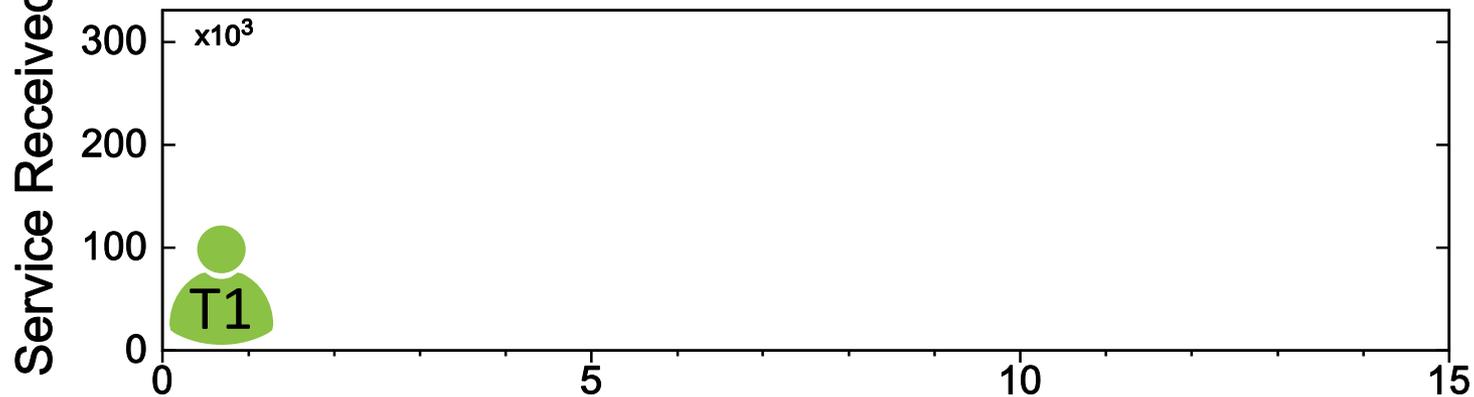




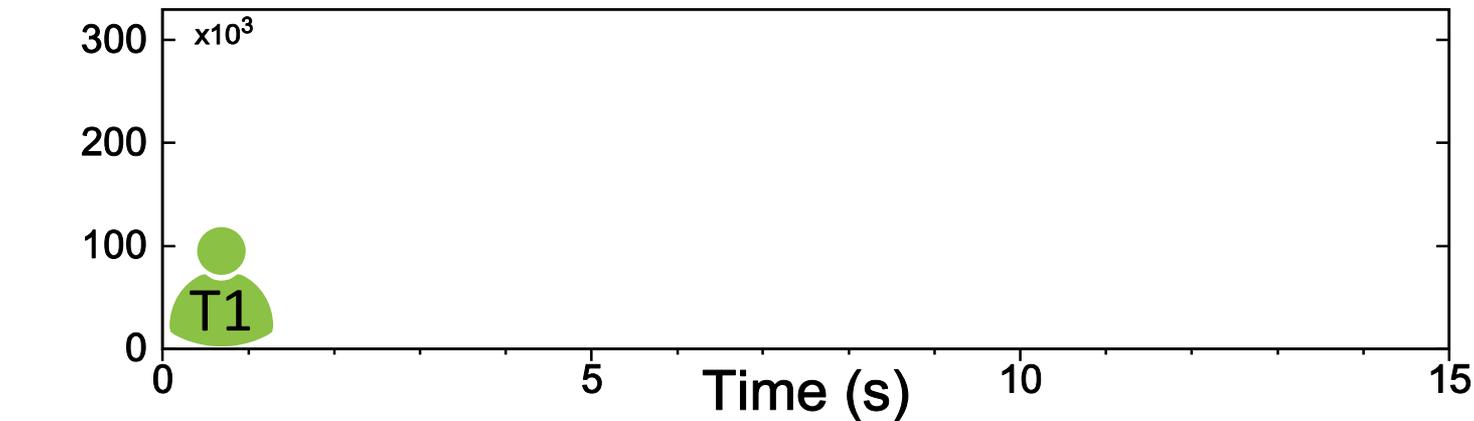




Predictable

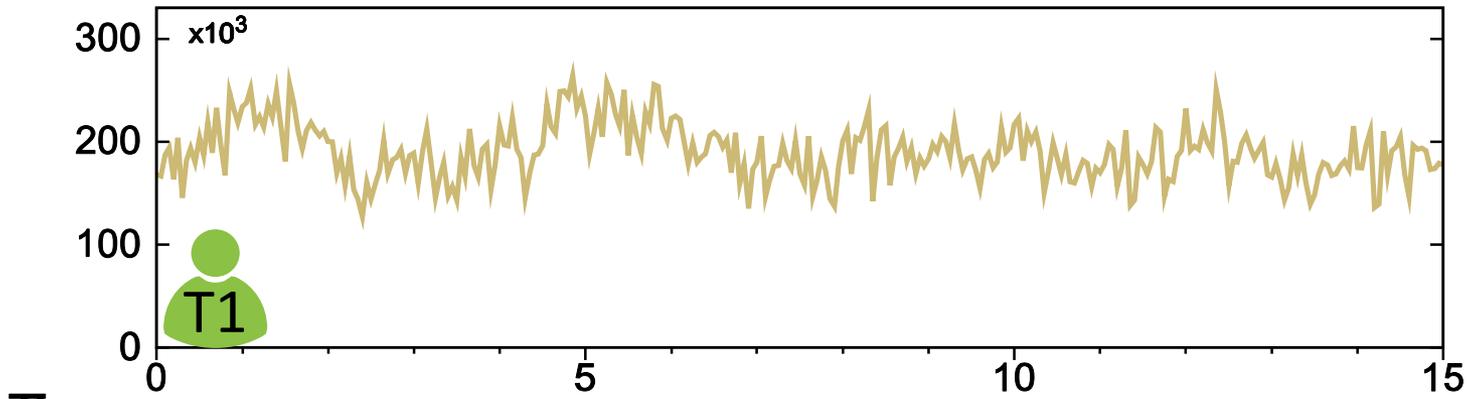


2/3 predictable  
1/3 unpredictable

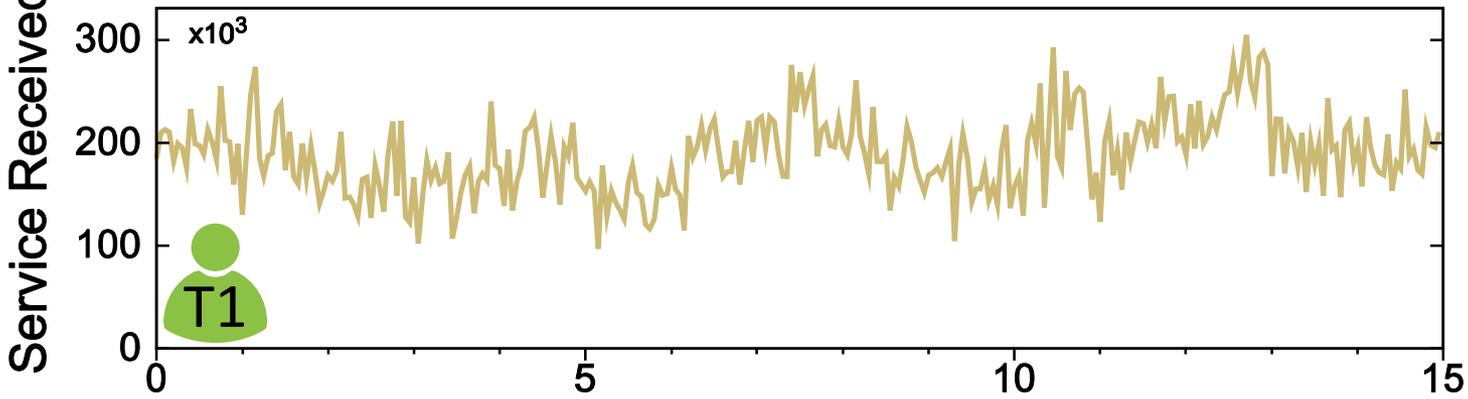


1/3 predictable  
2/3 unpredictable

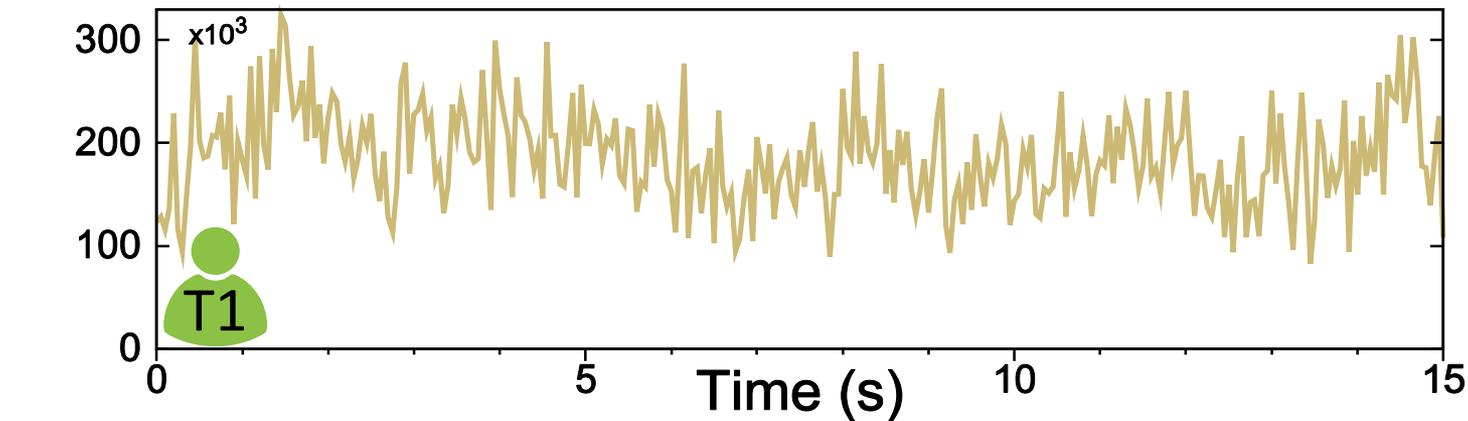
WFQ<sup>E</sup> ———



Predictable

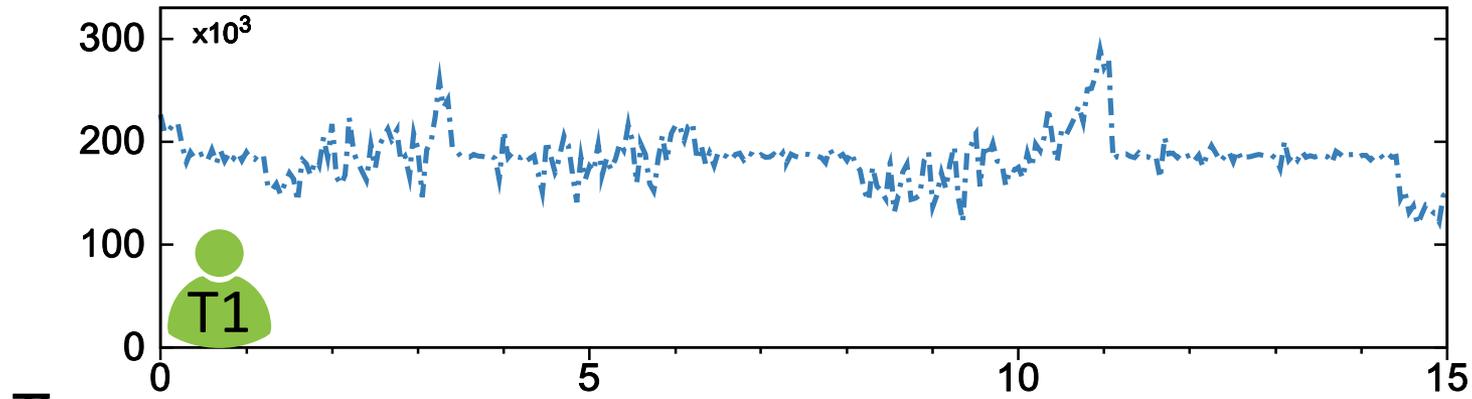


2/3 predictable  
1/3 unpredictable

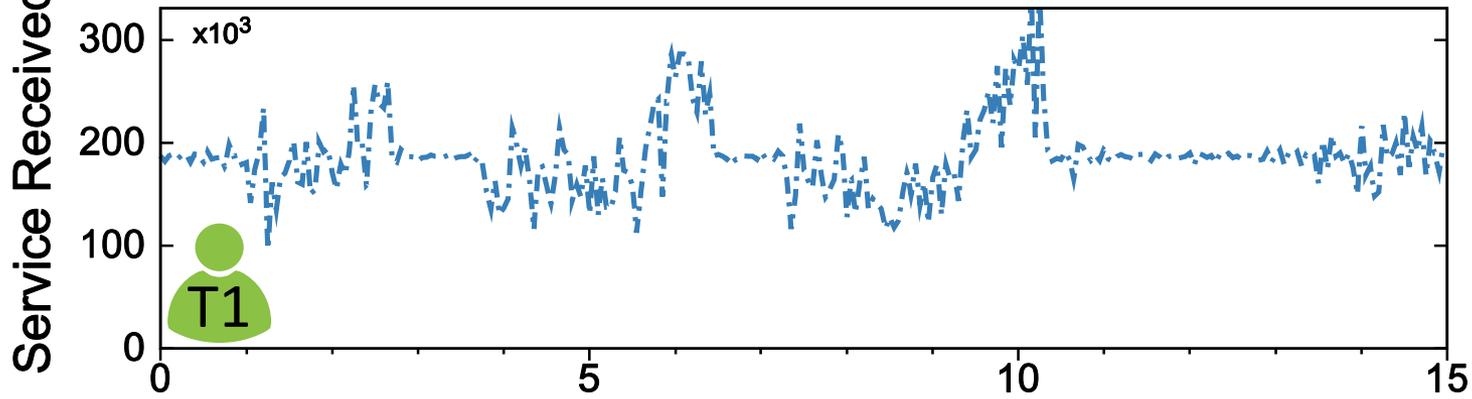


1/3 predictable  
2/3 unpredictable

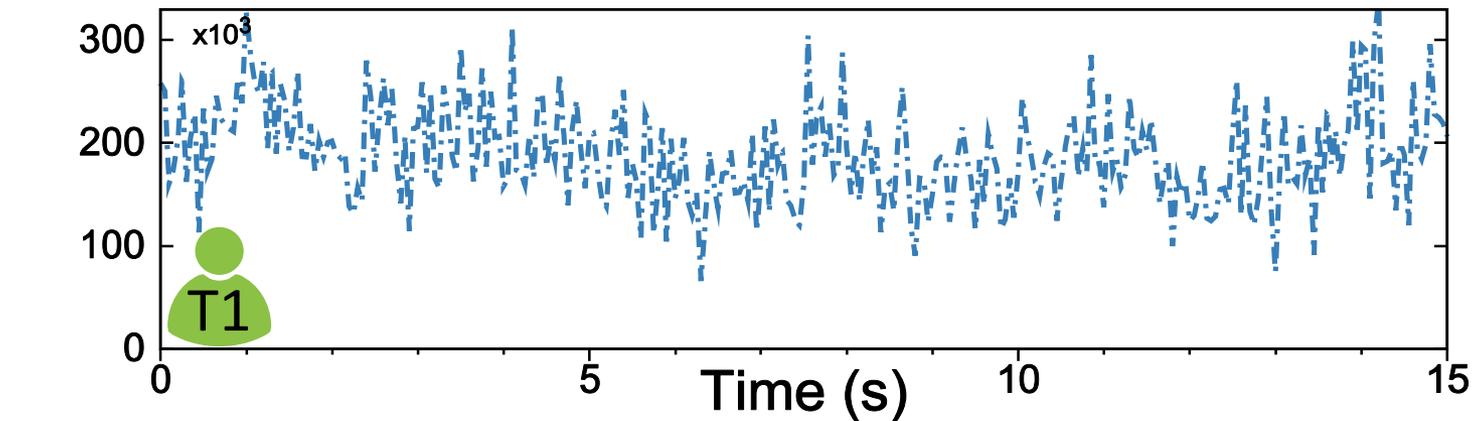
WF<sup>2</sup>Q<sup>E</sup> 



Predictable

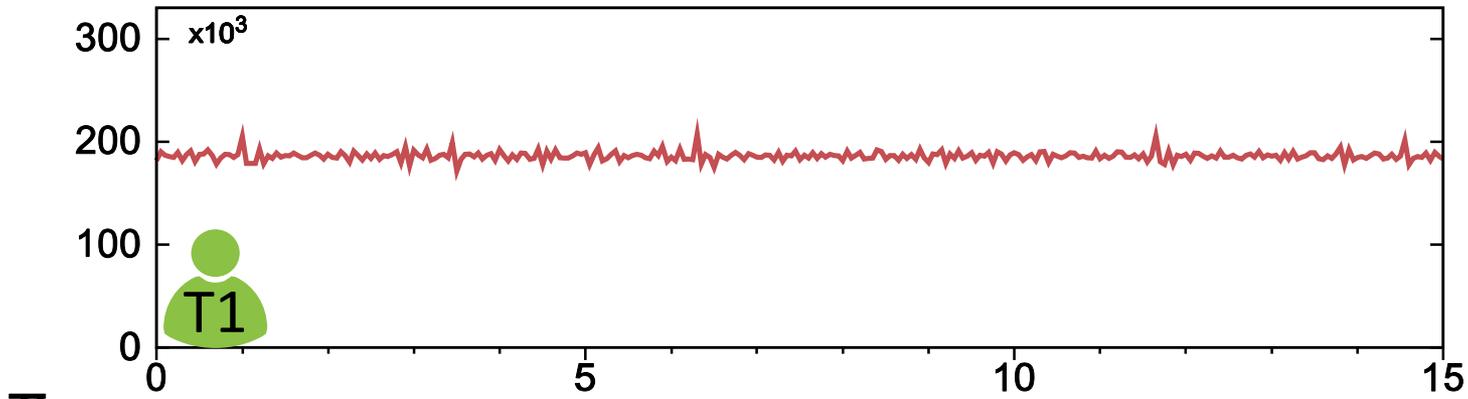


2/3 predictable  
1/3 unpredictable

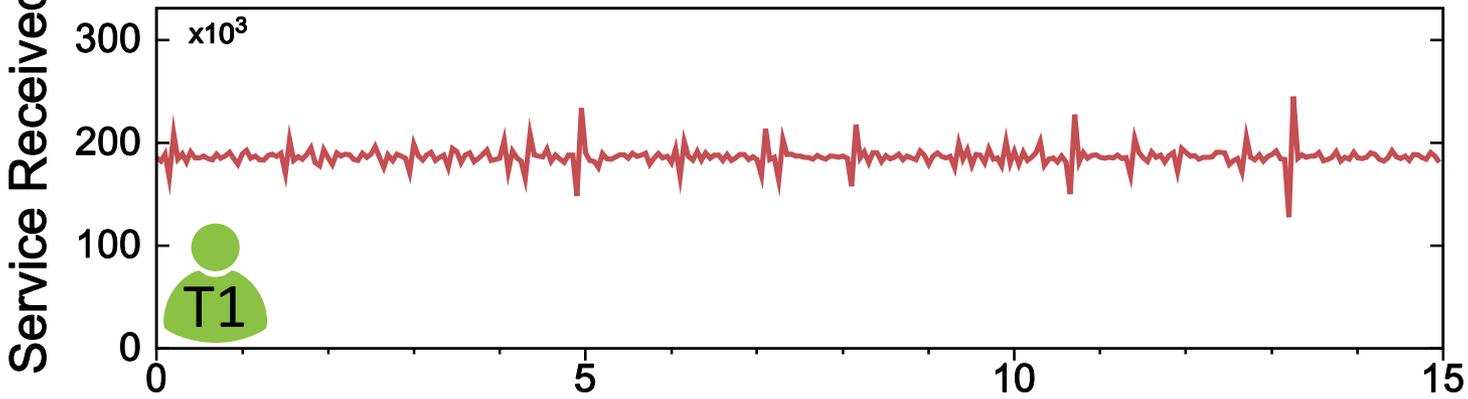


1/3 predictable  
2/3 unpredictable

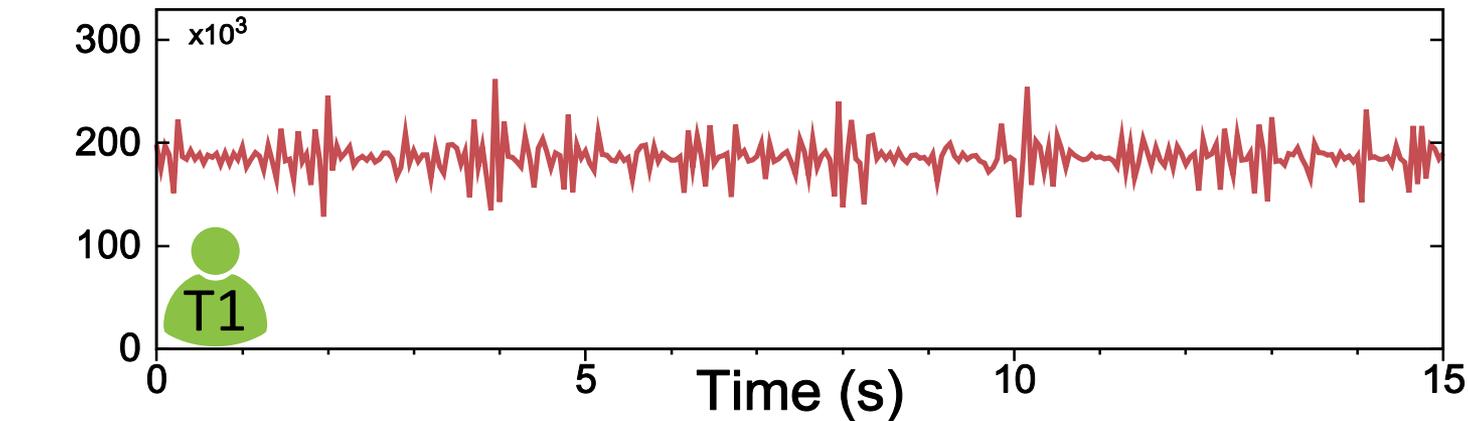
2DFQ<sup>E</sup> 



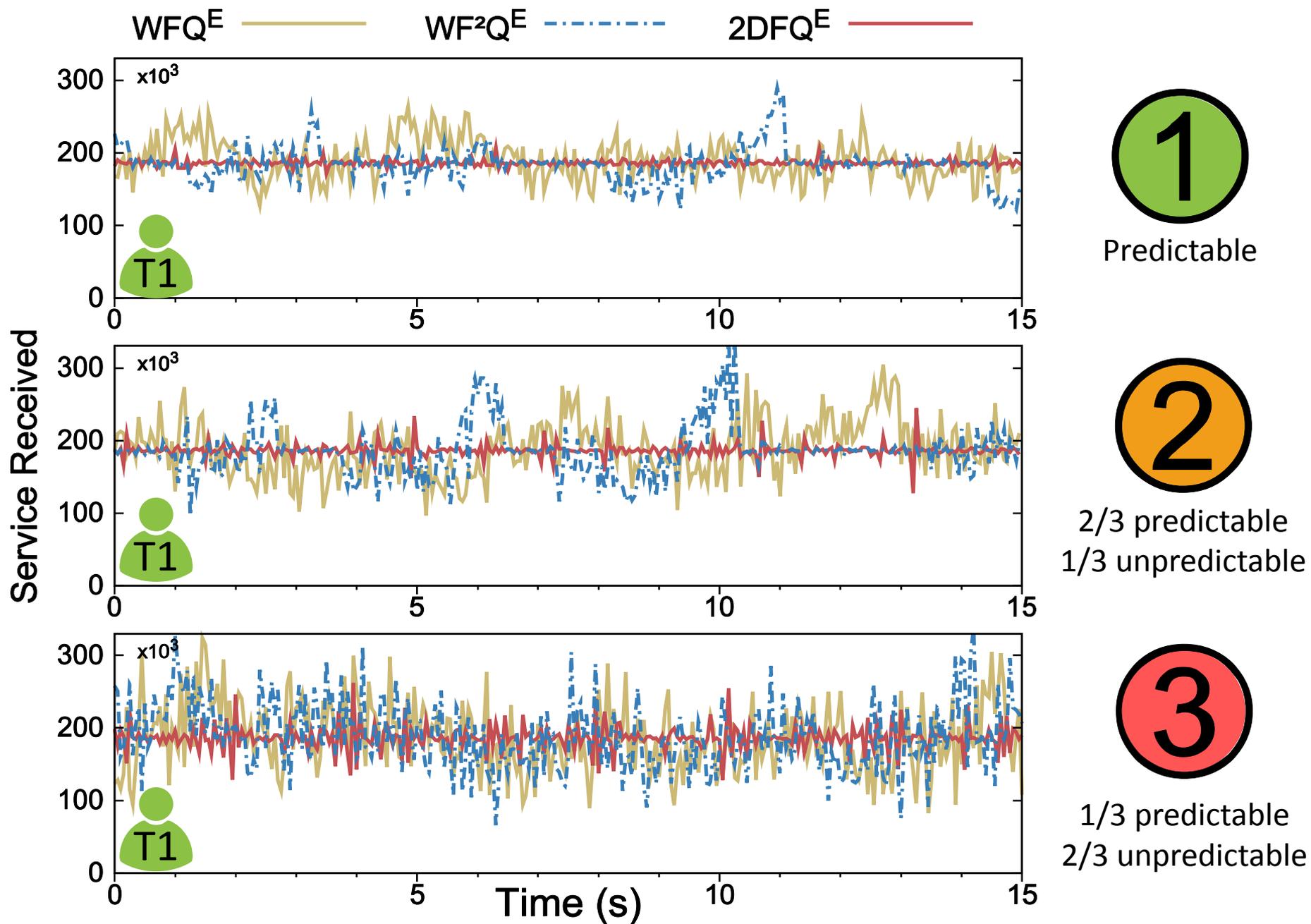
Predictable

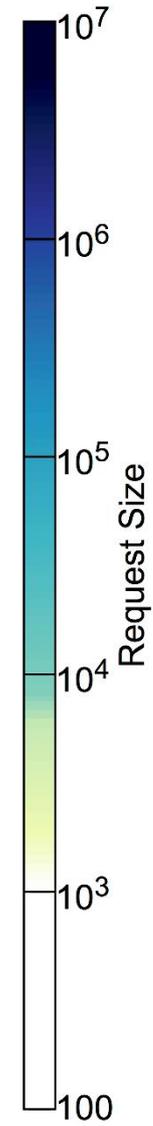


2/3 predictable  
1/3 unpredictable



1/3 predictable  
2/3 unpredictable





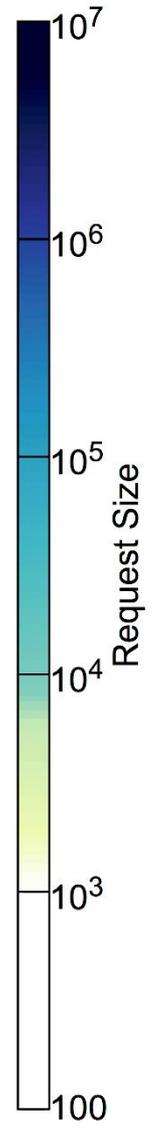
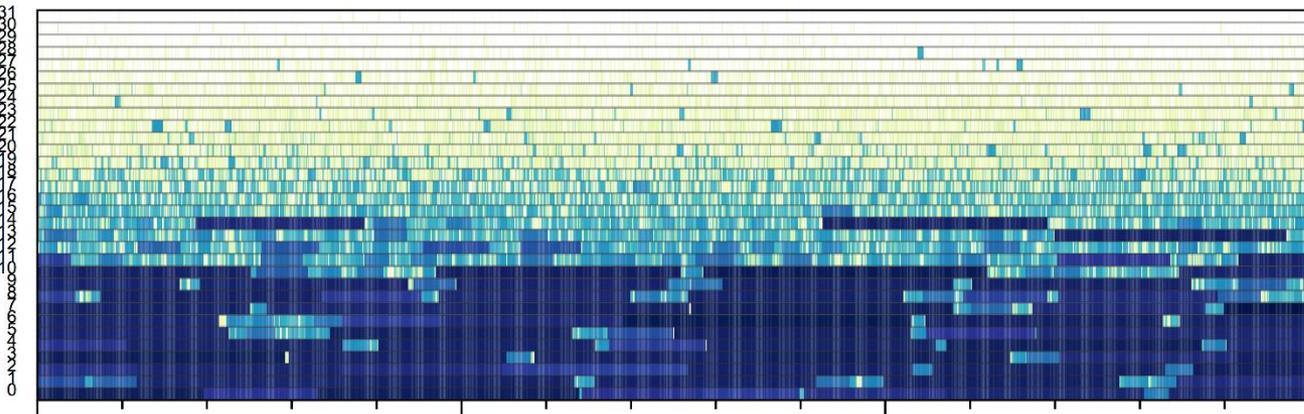
Predictable



2/3 predictable  
1/3 unpredictable



1/3 predictable  
2/3 unpredictable



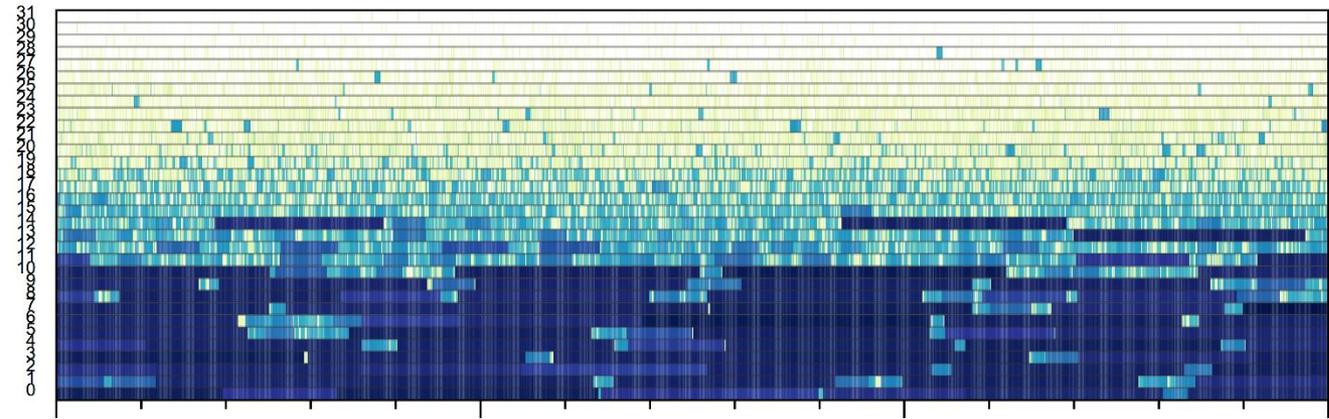
Predictable



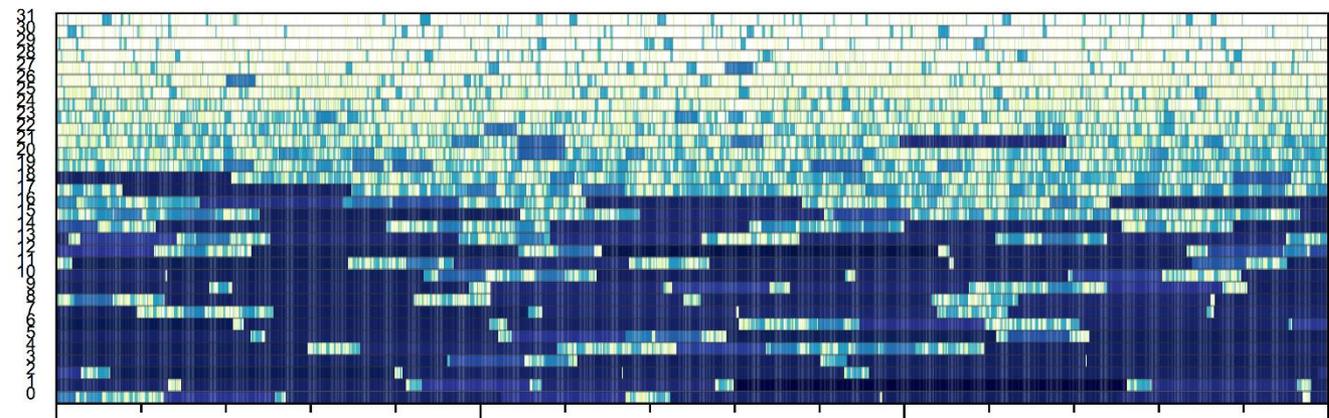
2/3 predictable  
1/3 unpredictable



1/3 predictable  
2/3 unpredictable



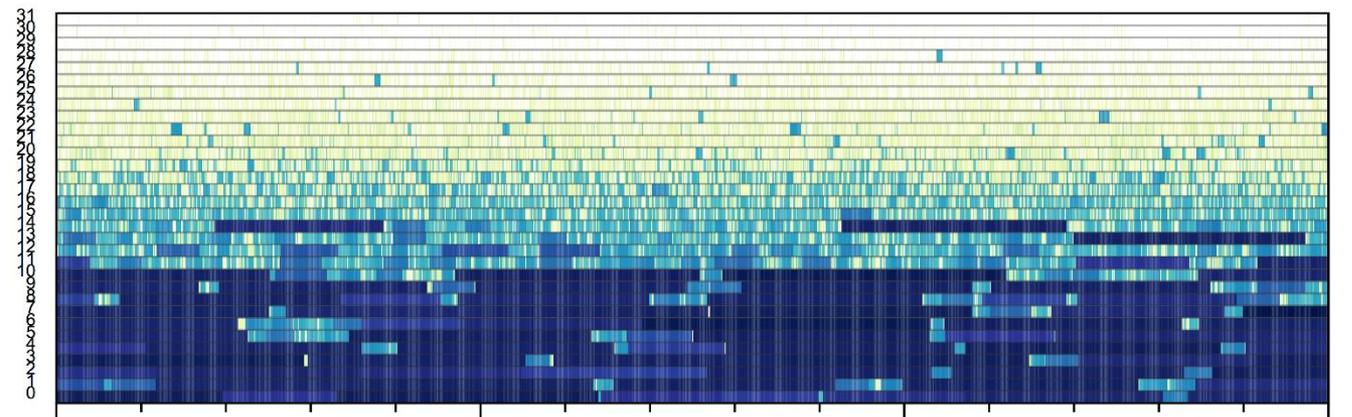
Predictable



2/3 predictable  
1/3 unpredictable

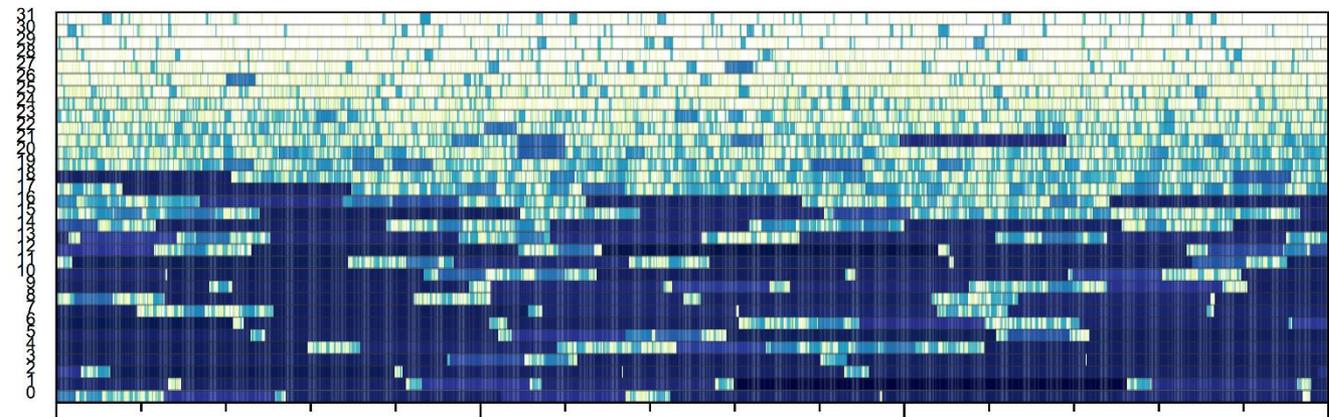


1/3 predictable  
2/3 unpredictable



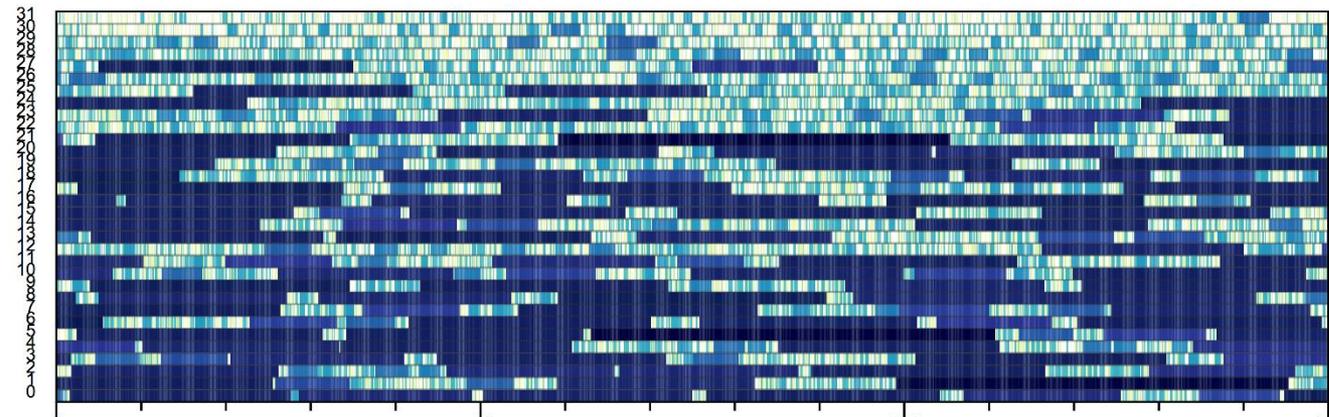
**1**

Predictable



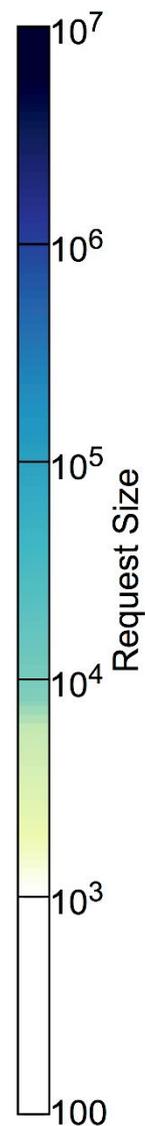
**2**

2/3 predictable  
1/3 unpredictable



**3**

1/3 predictable  
2/3 unpredictable



0 5 10 15  
Time (s)

# Two-Dimensional Fair Queueing

# Two-Dimensional Fair Queueing

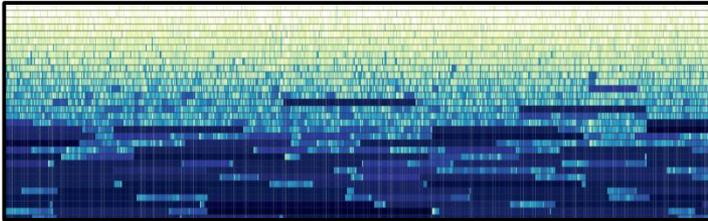
More threads  $\rightarrow$  Opportunity to reduce burstiness



# Two-Dimensional Fair Queueing

More threads  $\rightarrow$  Opportunity to reduce burstiness

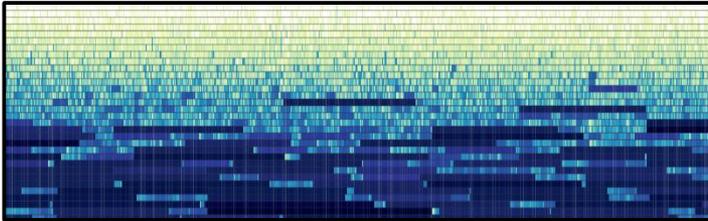
Partitions requests  
across threads by size



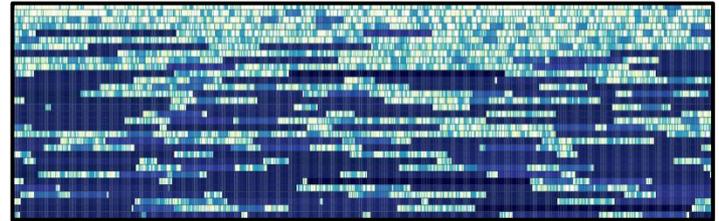
# Two-Dimensional Fair Queueing

More threads  $\rightarrow$  Opportunity to reduce burstiness

Partitions requests  
across threads by size



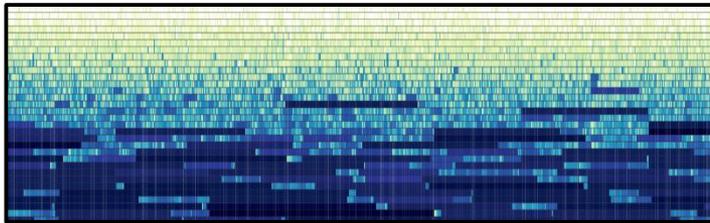
Co-locates unpredictable  
and expensive workloads



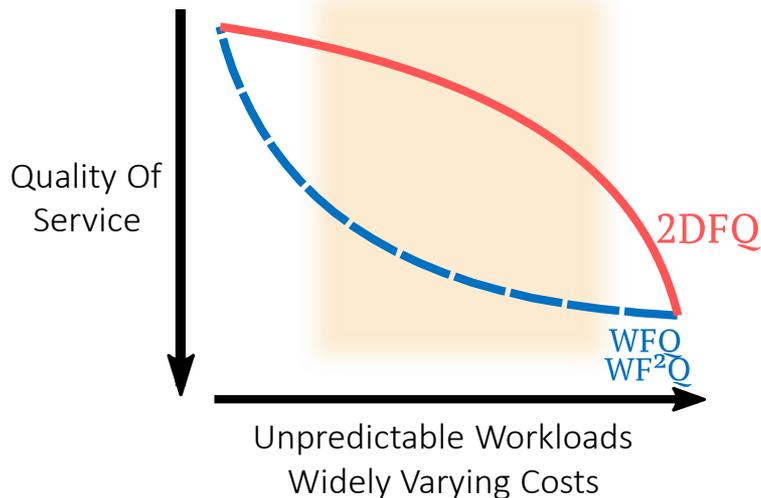
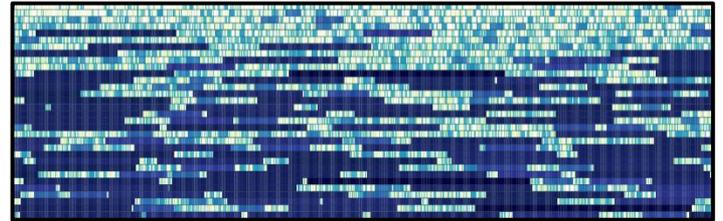
# Two-Dimensional Fair Queueing

More threads  $\rightarrow$  Opportunity to reduce burstiness

Partitions requests  
across threads by size



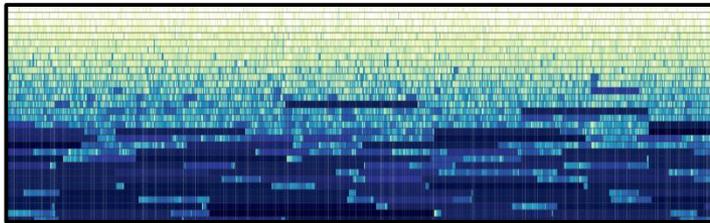
Co-locates unpredictable  
and expensive workloads



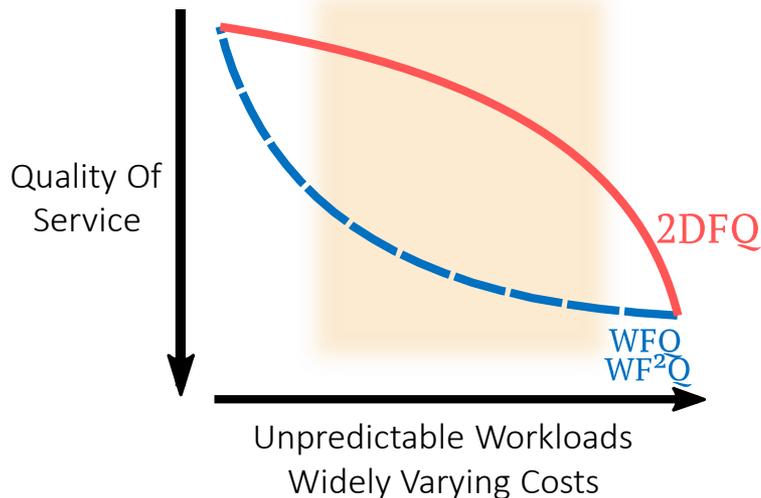
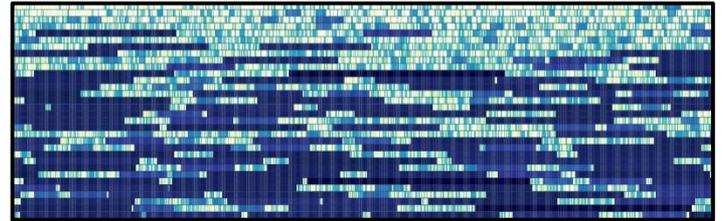
# Two-Dimensional Fair Queueing

More threads  $\rightarrow$  Opportunity to reduce burstiness

Partitions requests  
across threads by size



Co-locates unpredictable  
and expensive workloads



Reduced  
tail latency

Less burstiness