

Processes and Resources

CS-350 - Fall 2024

CS-350 – Fundamentals of Computing Systems
2024 – Renato Mancuso

BOSTON
UNIVERSITY

From Program to Process

sticking to the computing flavor

Simply put:

Process:

pre

Program

execution

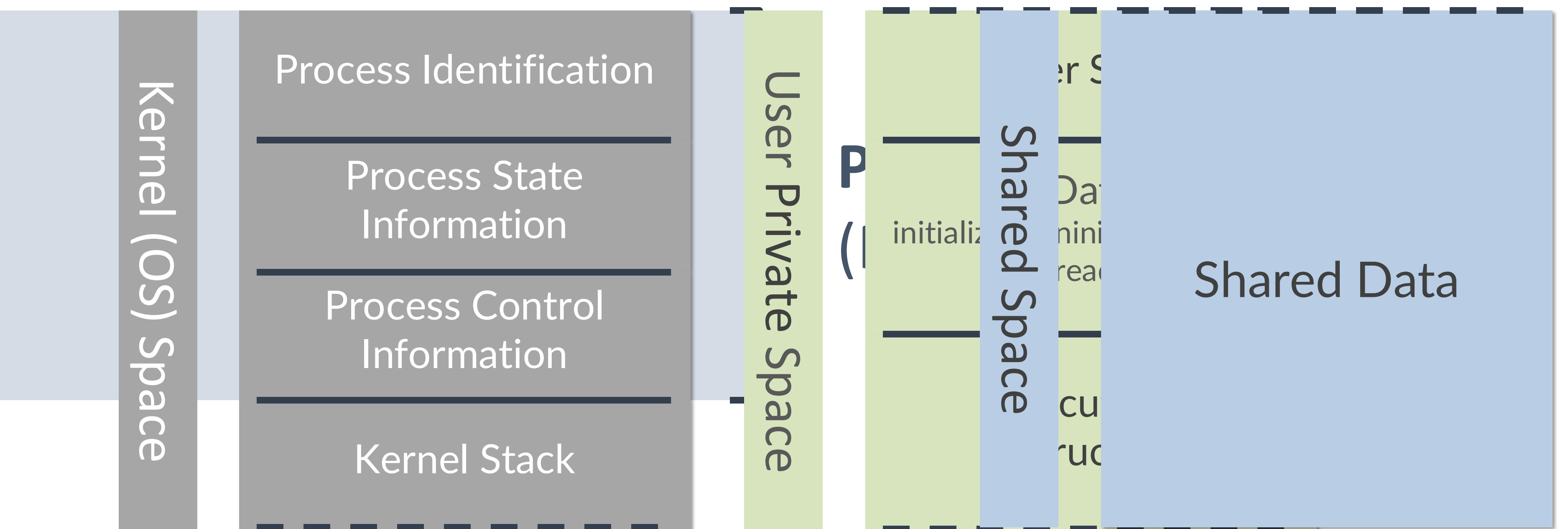
System

...on

From Program to Process

sticking to the computing flavor

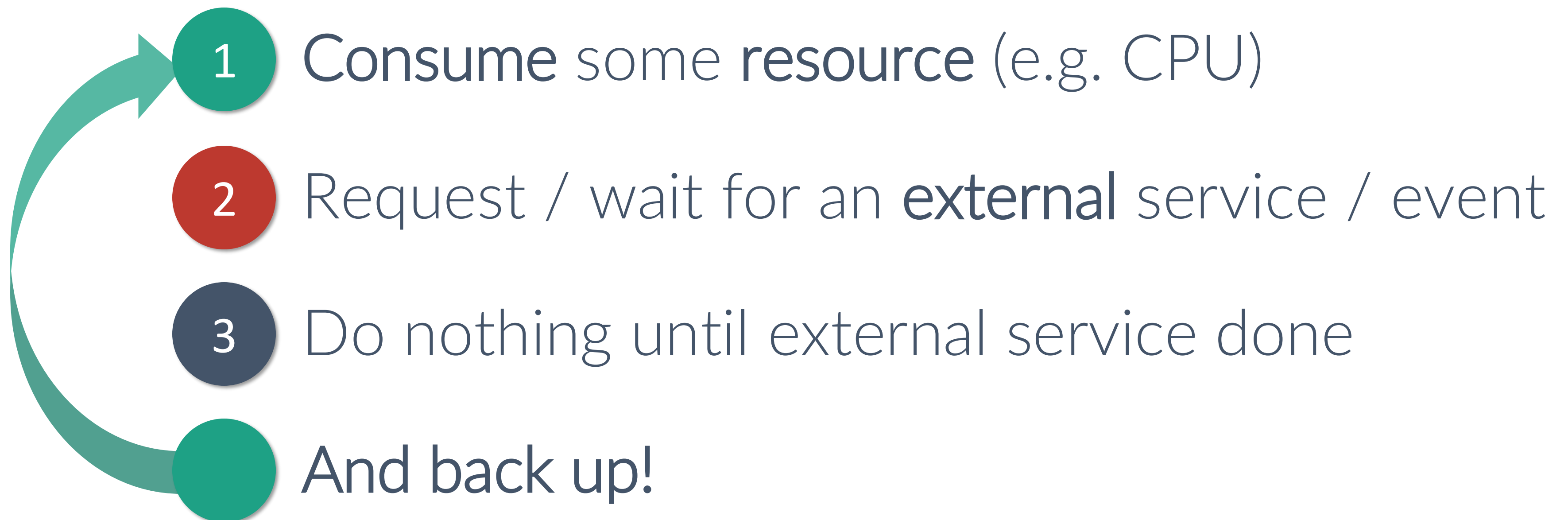
And to be “alive”, a **state** is necessary



From Program to Process

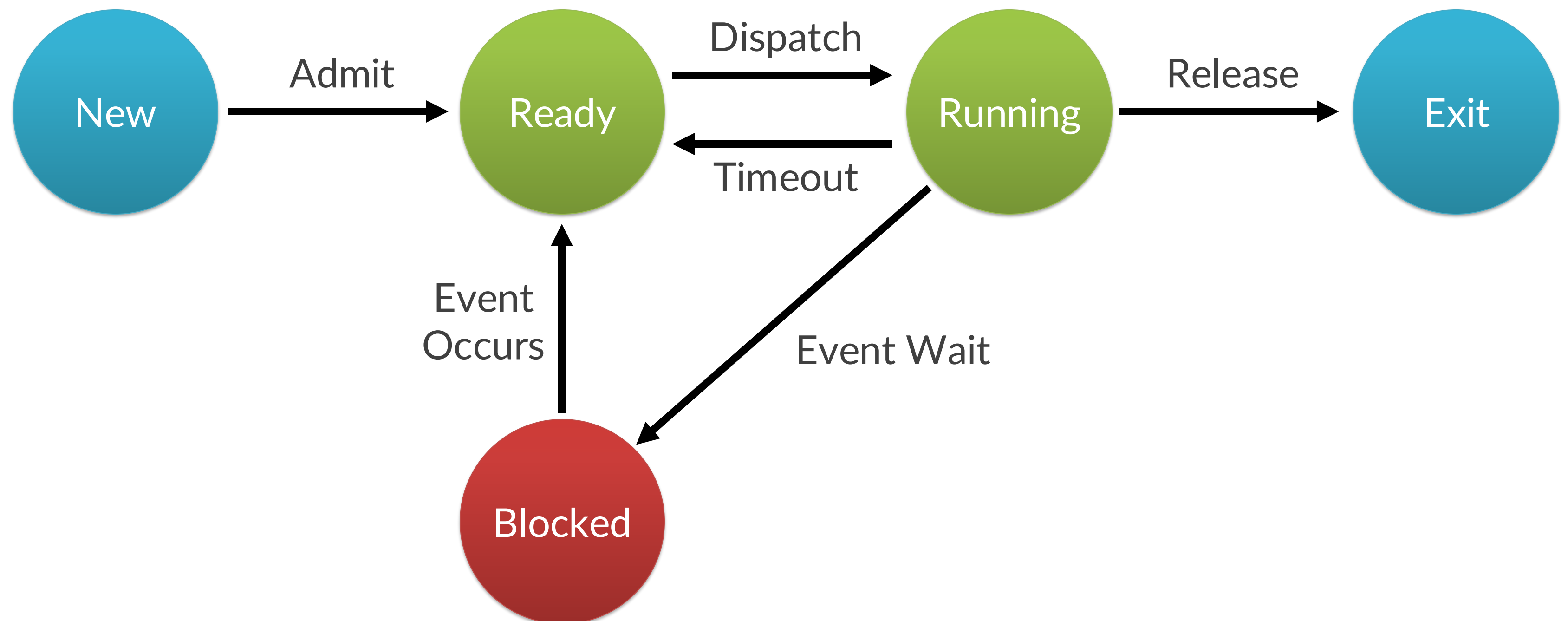
sticking to the computing flavor

But we can **abstract away** the state
and only focus on macro-phases



Transition State Diagram

It's diagrams all the way down



Transition State Diagram

It's diagrams all the way down



Ready

Ready to consume the resource



Running

Currently consuming the resource

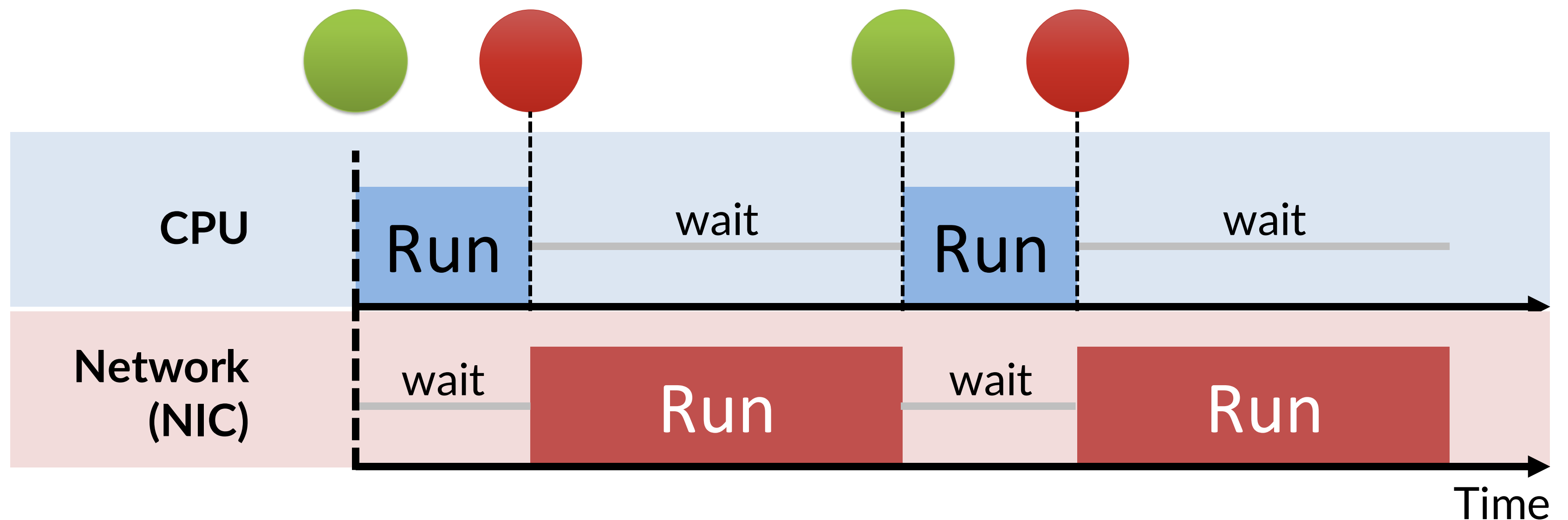


Blocked

Not competing for resource
because waiting for **some event**

Resource Usage Timeline

still a diagram, but over time

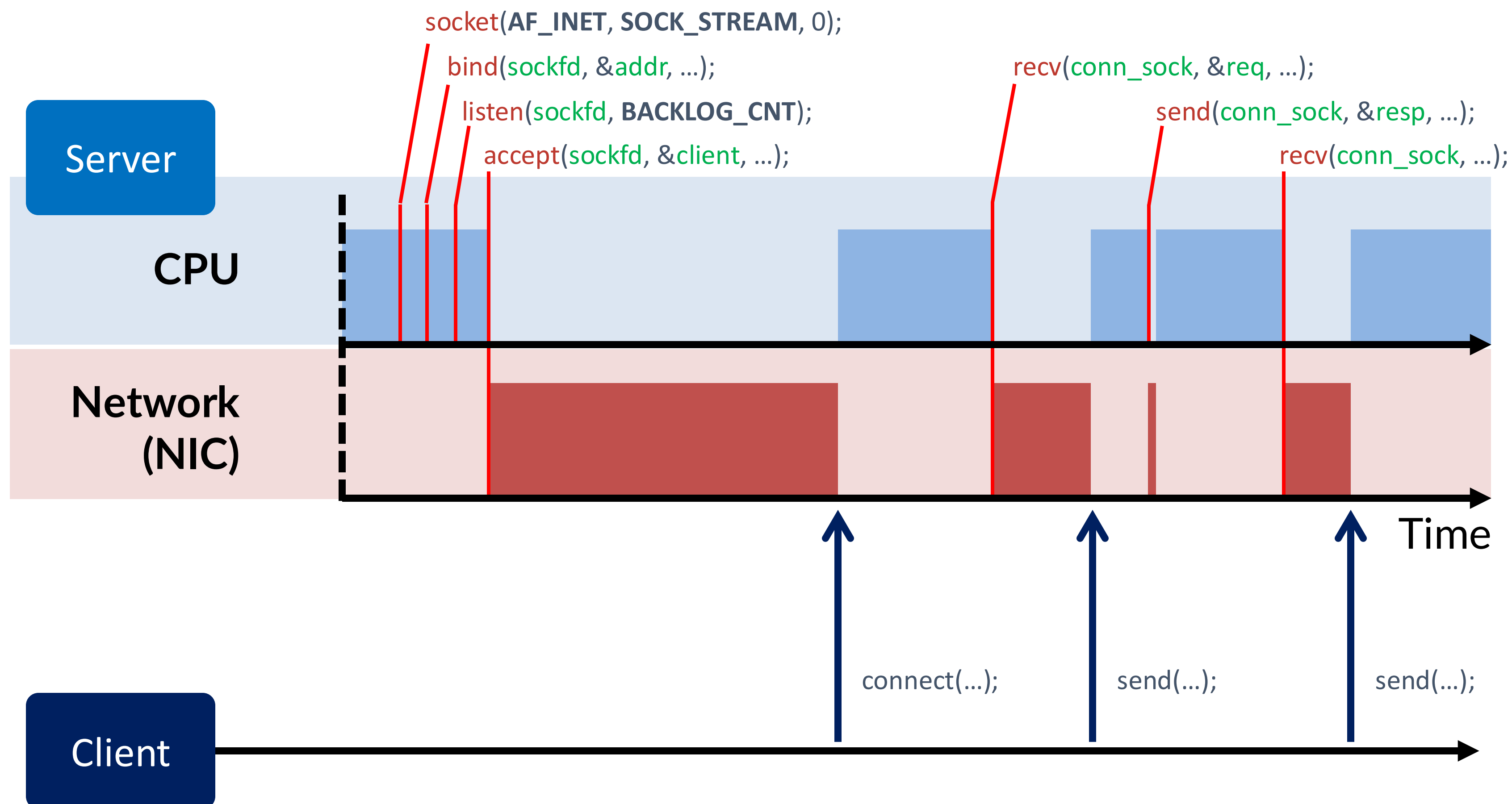


While blocked on I/O, **CPU idle!**

Uh and **the same for I/O** when CPU busy

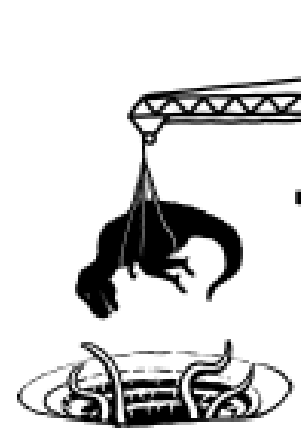
Let's Build a Server!

coincidentally also what you are doing in your first assignment



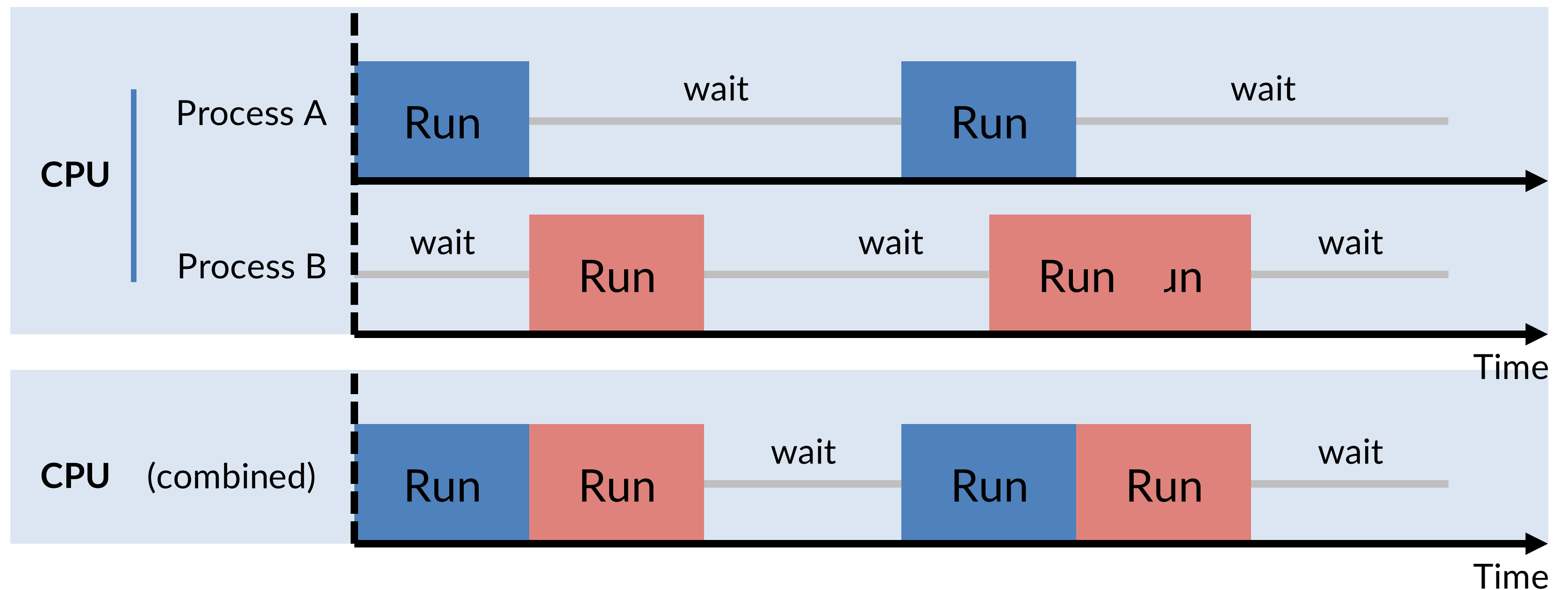
Multi-Programming

ooh jeez here we go



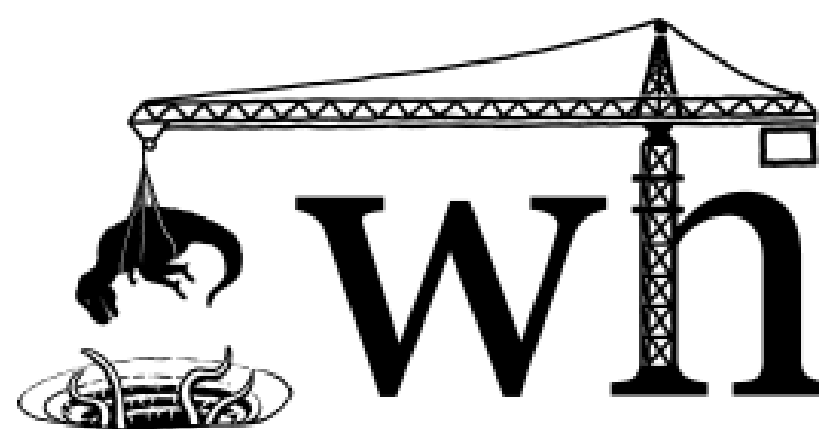
what if?

we throw **more processes** to the CPU?

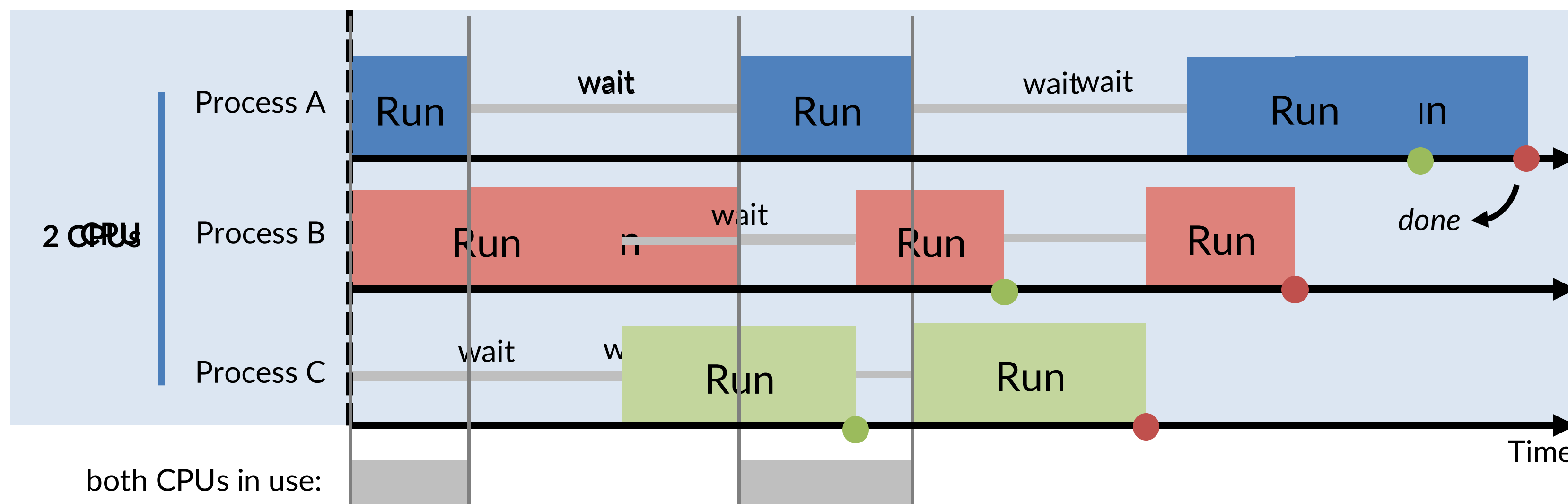


Multi-Processing

yep, before it was *multi-programming*

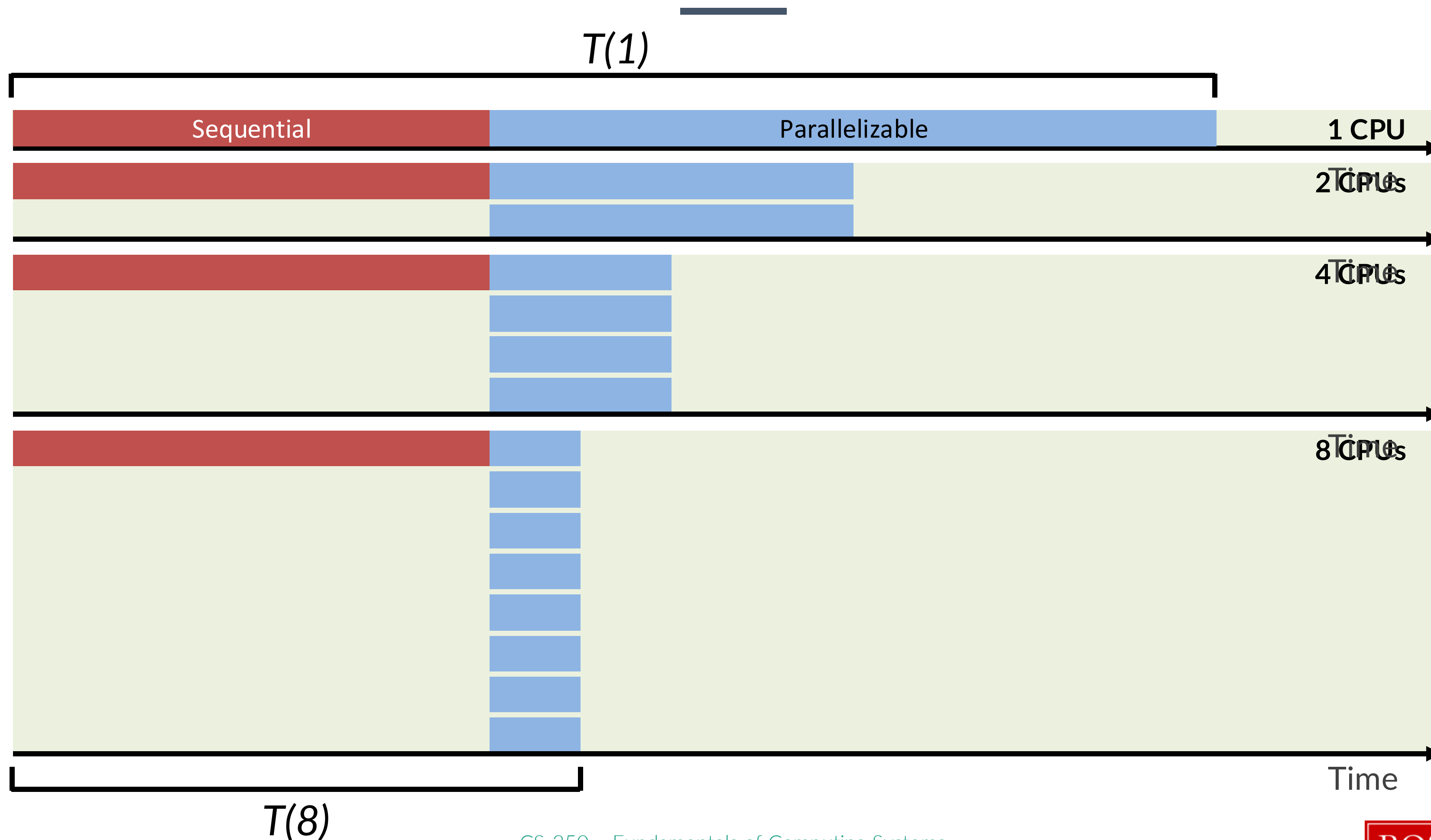


what if? we throw more processors



But there is a Limit

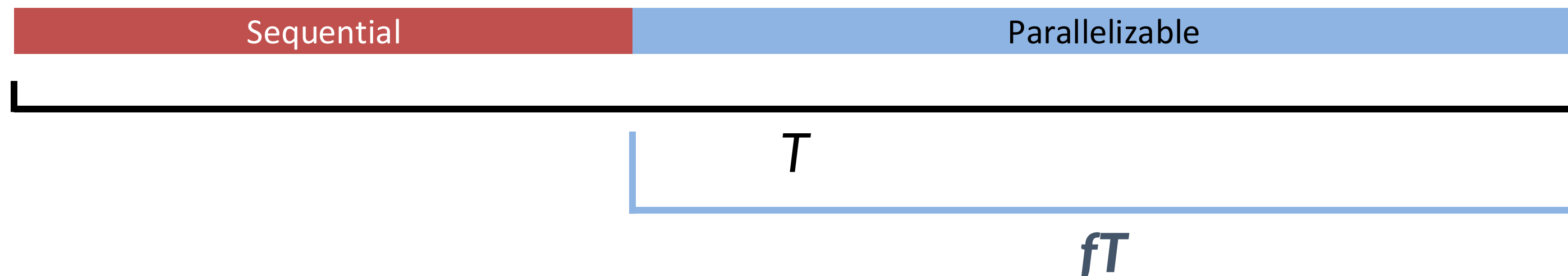
time for the bad news



Computing the Speedup

obey the law

If the **fraction** of parallelizable code is f



speedup =

$T(N) =$

Amdahl's Law

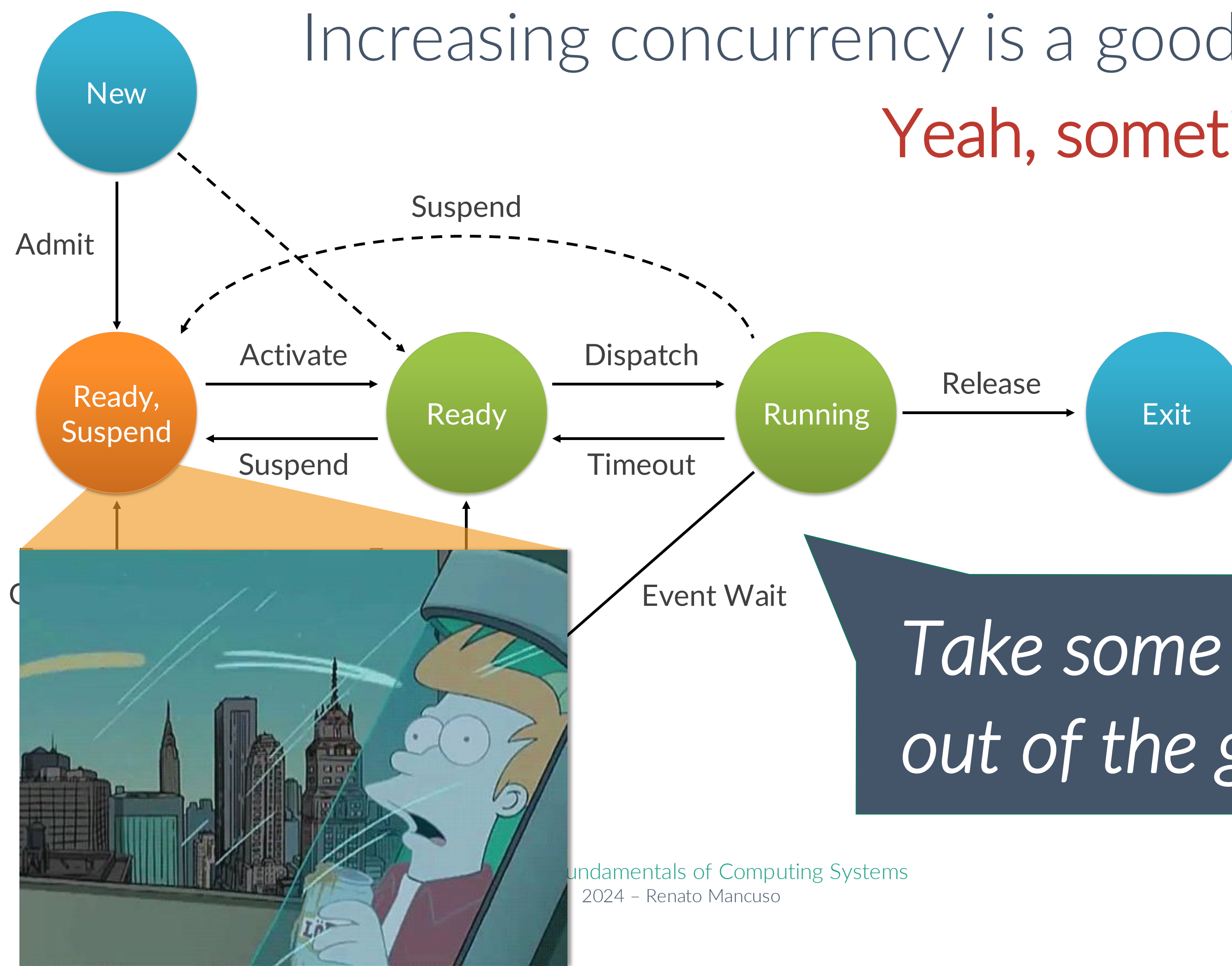
_____ *parallelizing stuff since 1967*

Managing Concurrency

there is a new sheriff in town

Increasing concurrency is a good idea.

Yeah, sometimes.

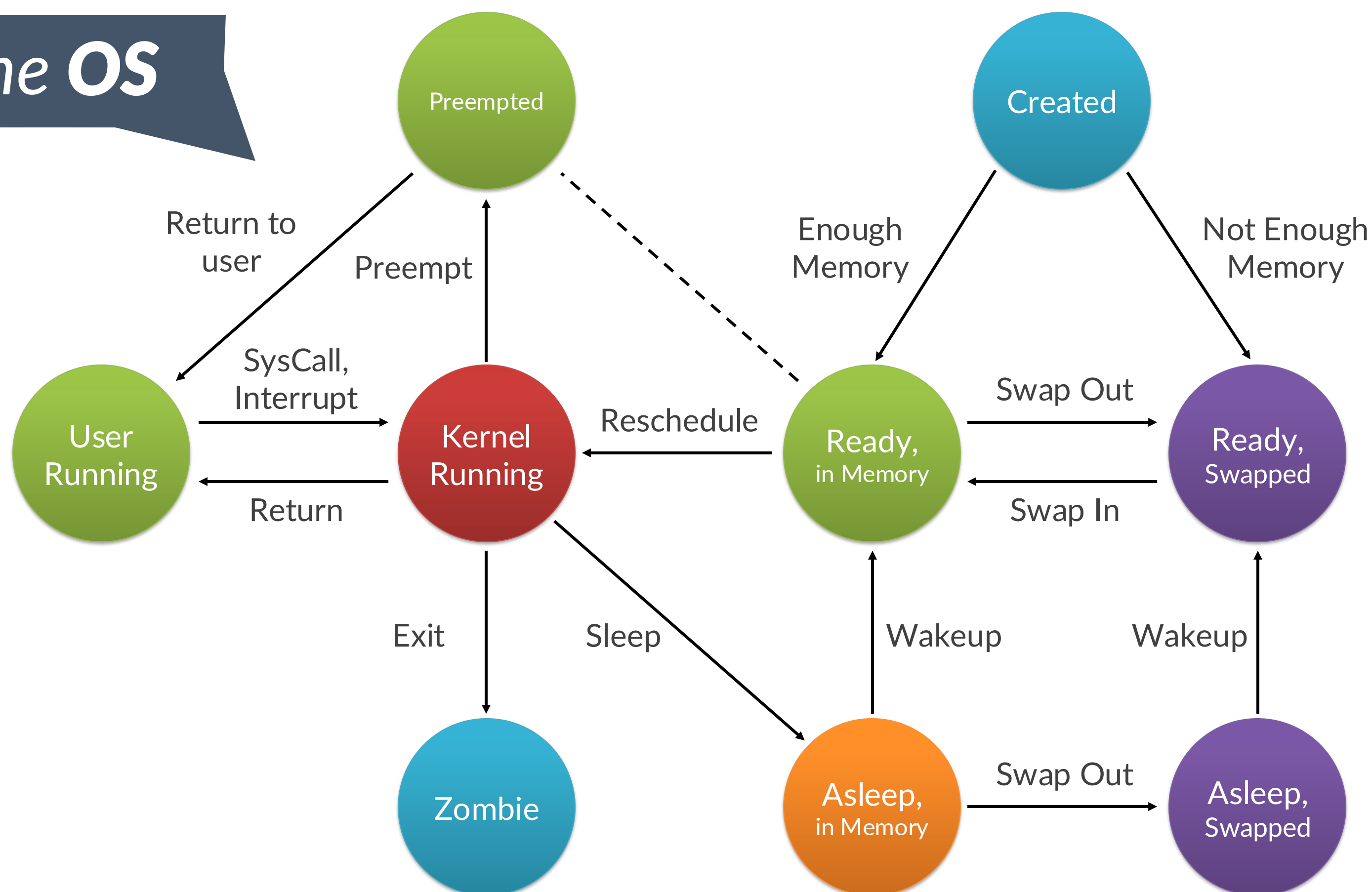


Take some processes out of the game

Managing Concurrency

there is a new sheriff in town

All hail the OS



System Model & Metrics

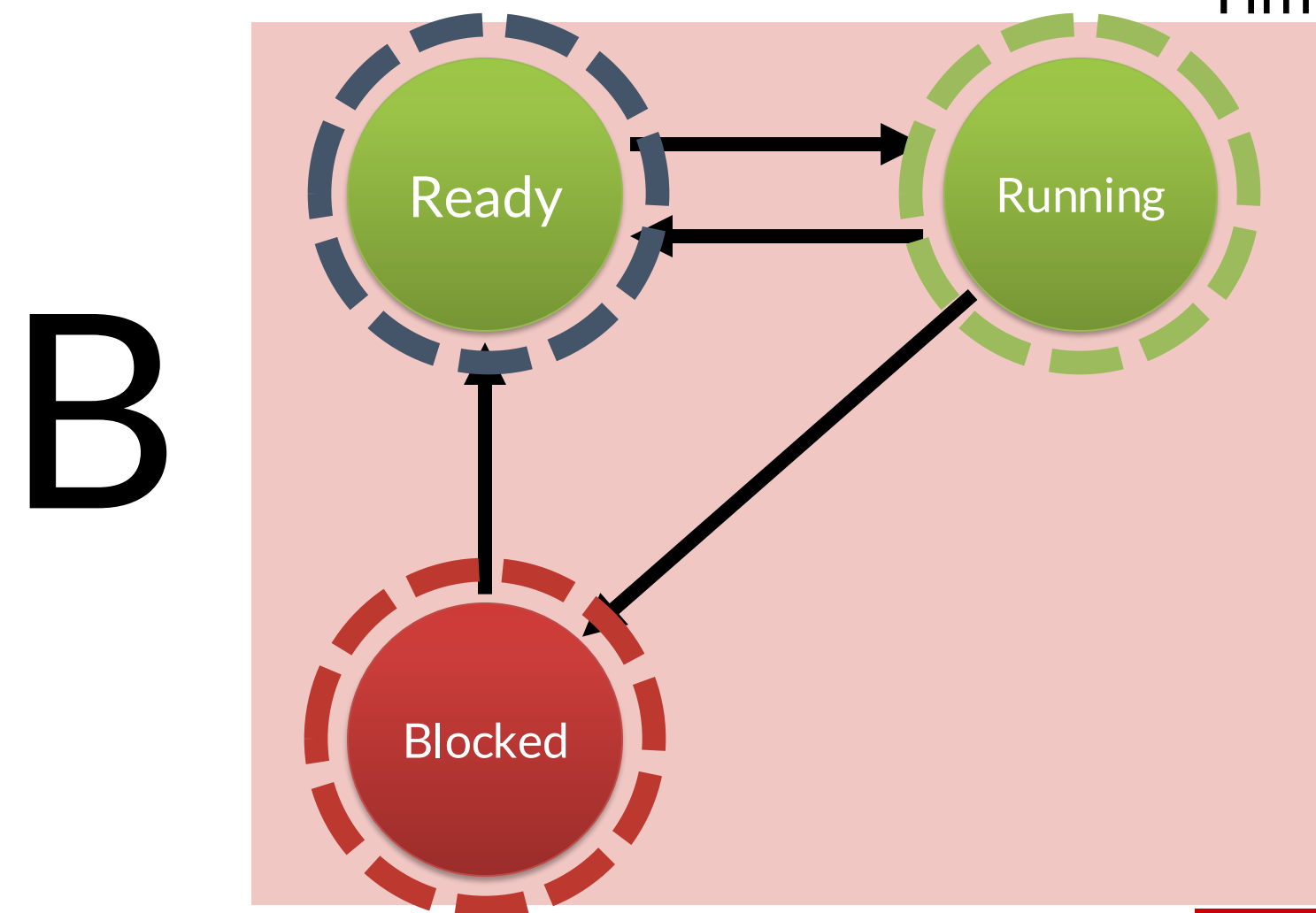
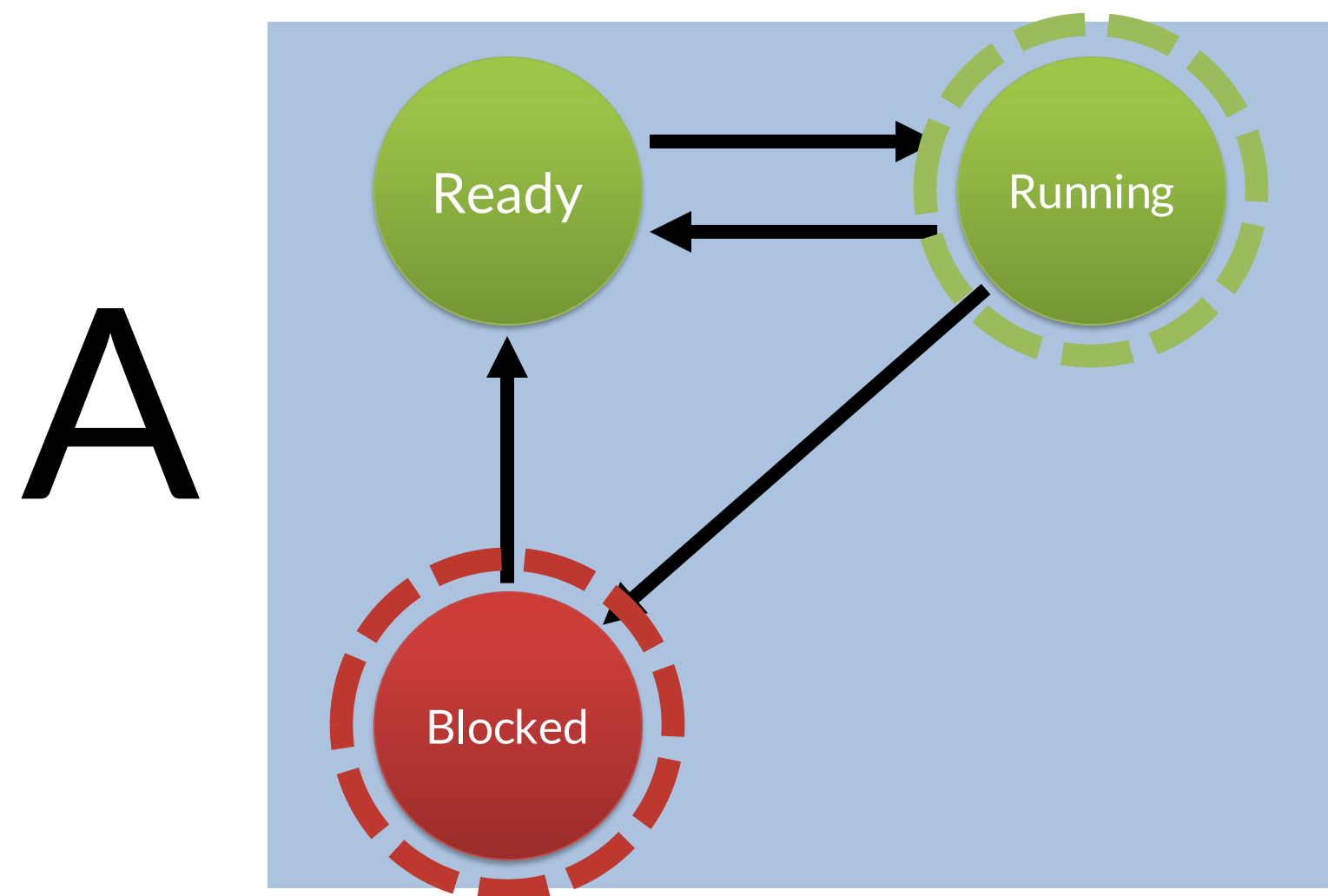
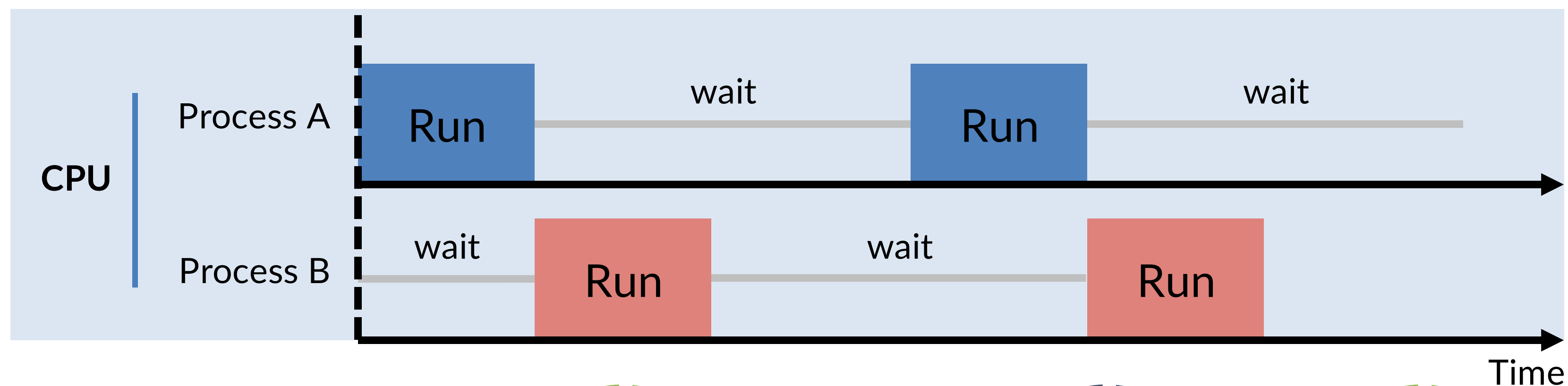
CS-350 - Fall 2024

CS-350 – Fundamentals of Computing Systems
2024 – Renato Mancuso

BOSTON
UNIVERSITY

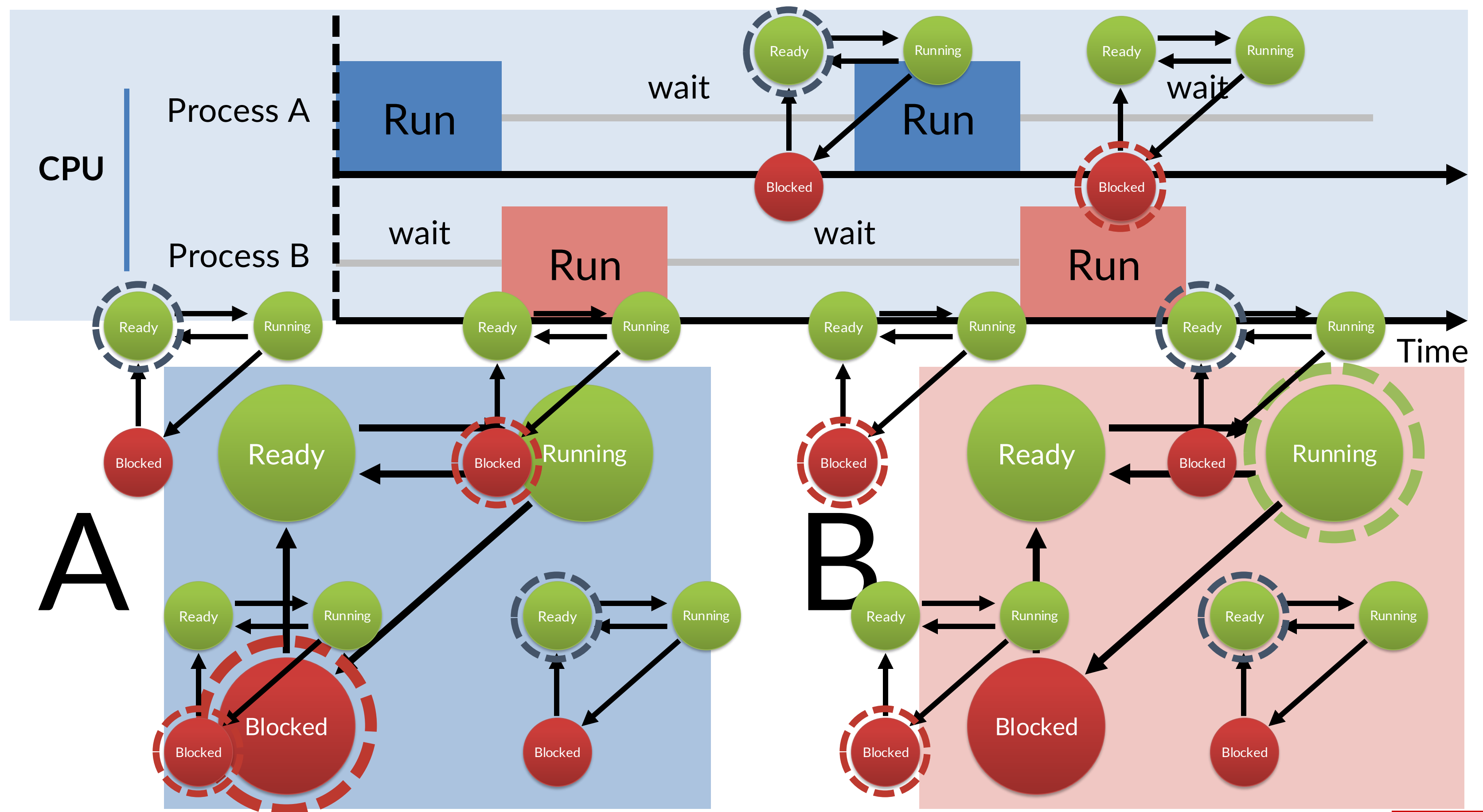
Tracking a Process

with respect to a single resource



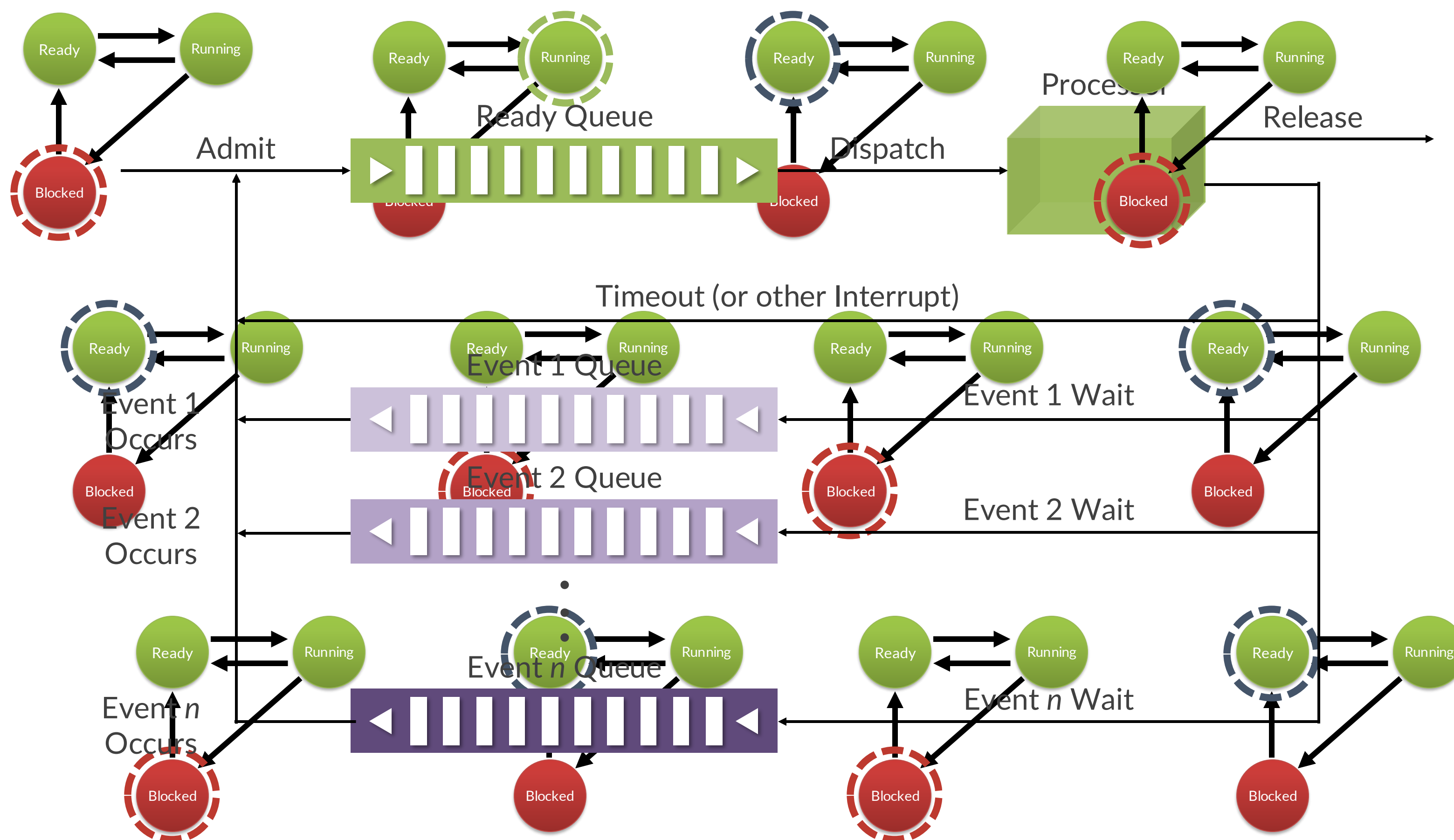
Tracking a Process

with respect to a single resource



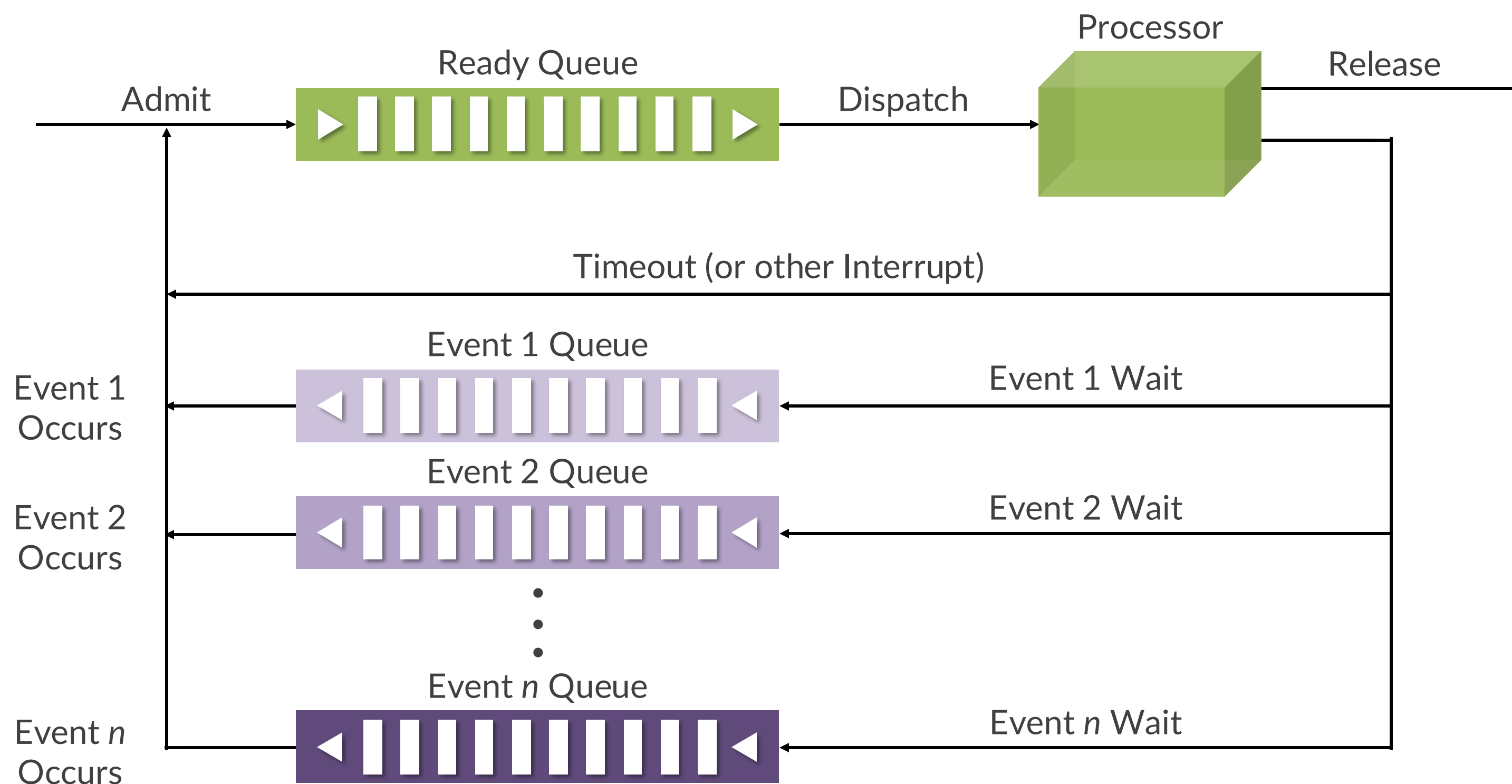
Tracking a Process

a trippy re-thinking from processes



Tracking a System

with multiple resources



CS-350 in a Nutshell

oh boy here we go again

Manage the **consumption of resources**
and **coordinate interaction** of concurrent processes
to ensure some **desirable system attributes**
which can be **evaluated** and **contrasted**

Performance is Subjective

the tales of dueling sys-admins and users

System-centric perspective



- Are resources being utilized?
- Is revenue being maximized?
- Is system's wear minimized?
- ...etc.



Process-centric perspective

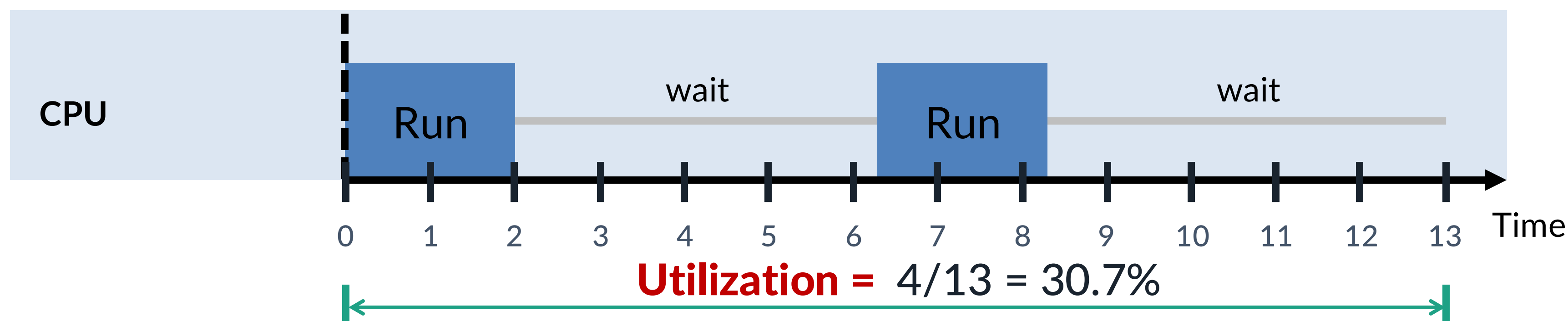
- Average processing time?
- Worst-case processing time?
- Waiting time before service?
- ...etc.



Utilization

aka the “bang” per “buck”

Utilization: *fraction of time over a given time window during which the resource is busy (not idle).*



NOTE: *utilization is a per-resource metric!*

$$U = \frac{\text{Time Busy}}{\text{Total Time}} = 1 - \frac{\text{Time Idle}}{\text{Total Time}}$$

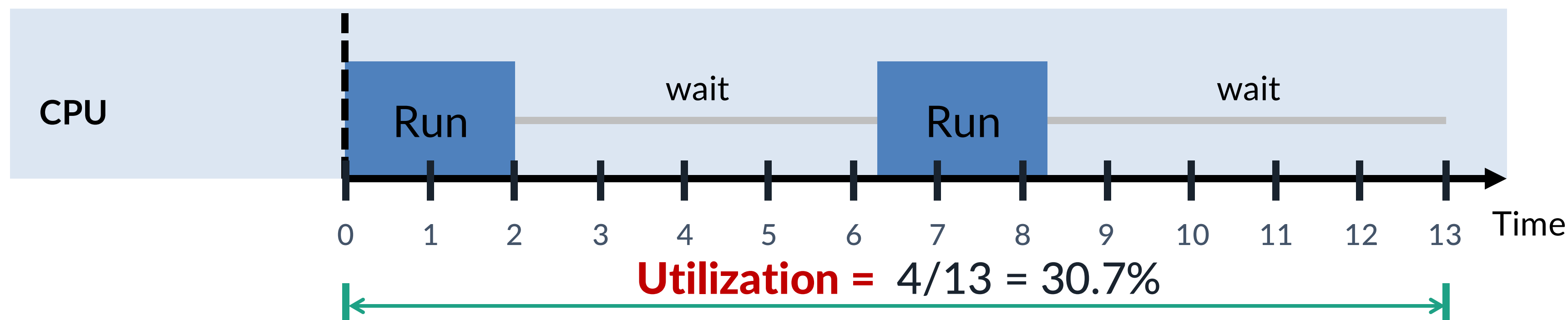


Multi-Programming & Utilization

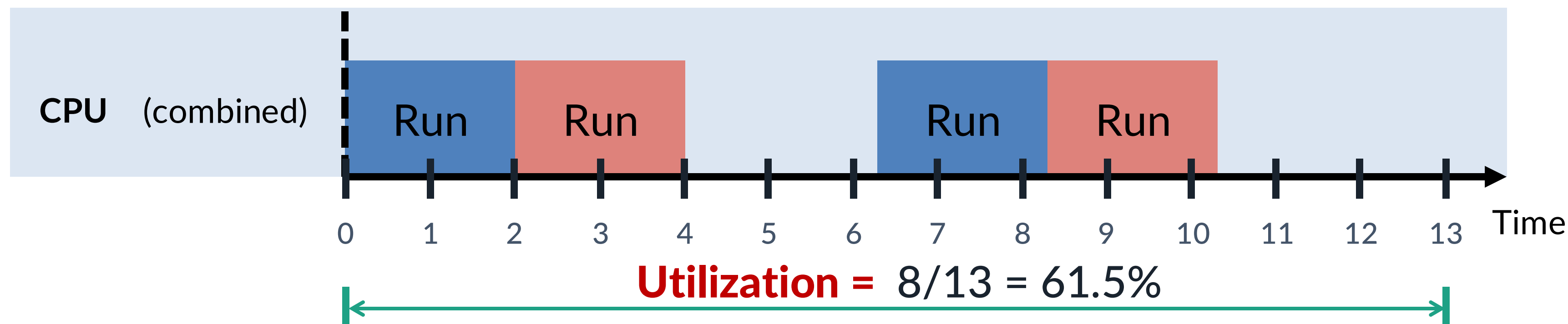
aka the “bang” per “buck”

Utilization: *fraction of time over a given time window during which the resource is busy (not idle).*

With Multi-Programming Level (MPL) = 1:



With Multi-Programming Level (MPL) = 2:

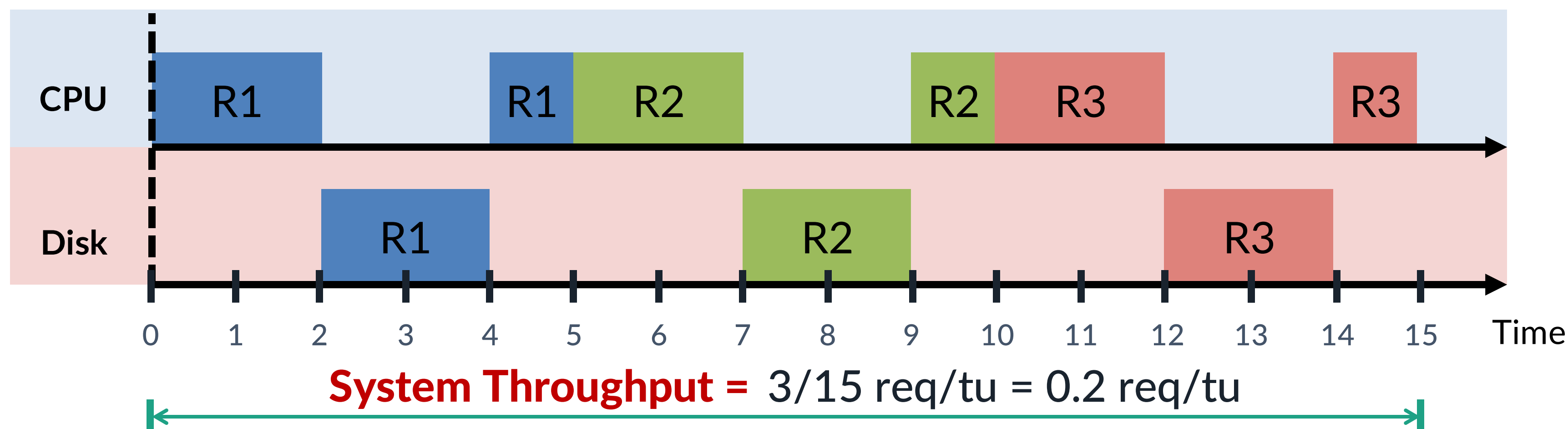




Throughput

a.k.a. output rate of requests

Throughput: *number of completed requests over a given time window.*



NOTE: for single resources, or for the whole system!

$$\text{Thr} = \frac{\text{\#Completed Requests}}{\text{Total Time}}$$

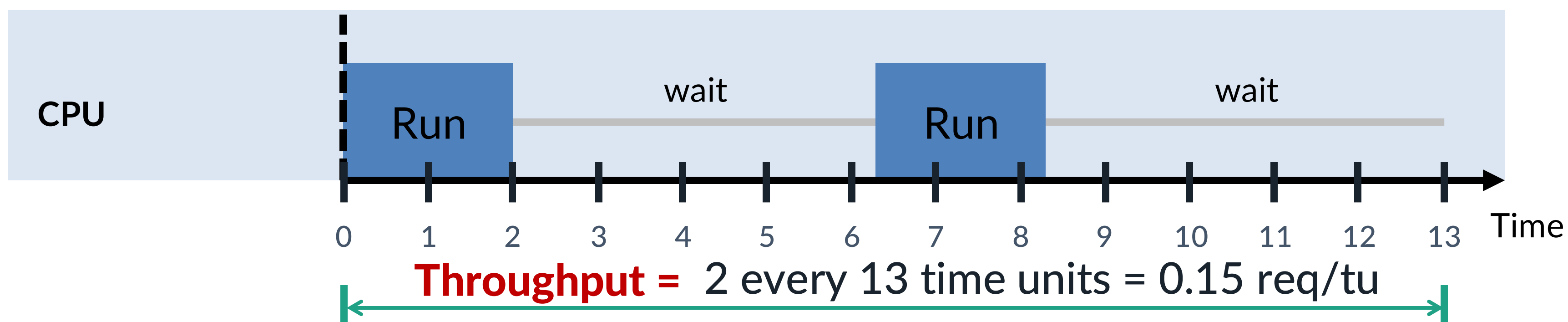


Multi-Programming & Throughput

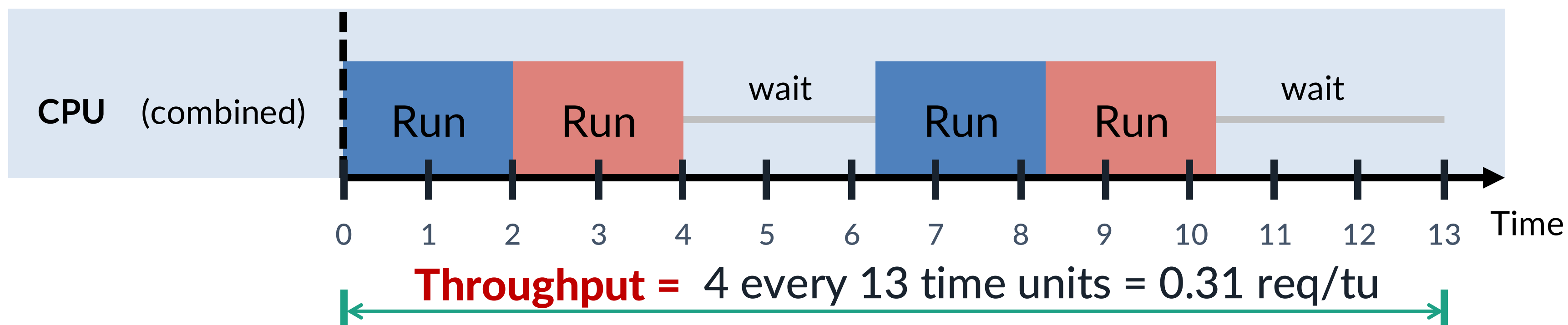
productivity explained

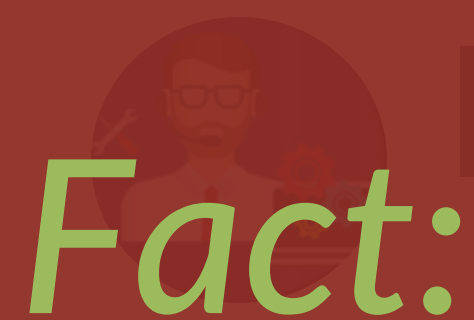
Throughput: *number of completed requests over a given time window.*

With Multi-Programming Level (MPL) = 1:



With Multi-Programming Level (MPL) = 2:





Multi-Programming & Throughput

productivity explained

Fact:

Throughput: number of completed requests over a given time window.

Increasing the MPL improves utilization and throughput.

With Multi-Programming Level (MPL) = 1:

But then:

Is this true forever?

I.e. for an arbitrarily high value of MPL?



With Multi-Programming Level (MPL) = 2:



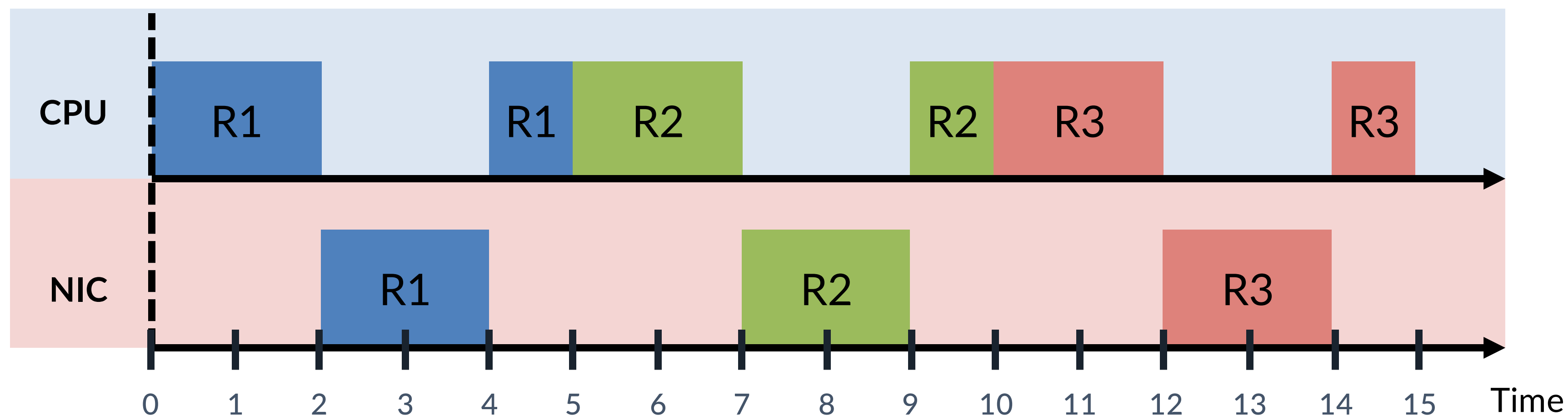


Capacity

when the party ends

Capacity: *maximum throughput that can be sustained by the system.*

With Multi-Programming Level (MPL) = 1:



CPU Utilization = $3/5 = 60\%$

NIC Utilization = $2/5 = 40\%$

System Throughput = $1/5 \text{ req/tu} = 0.2 \text{ req/tu}$

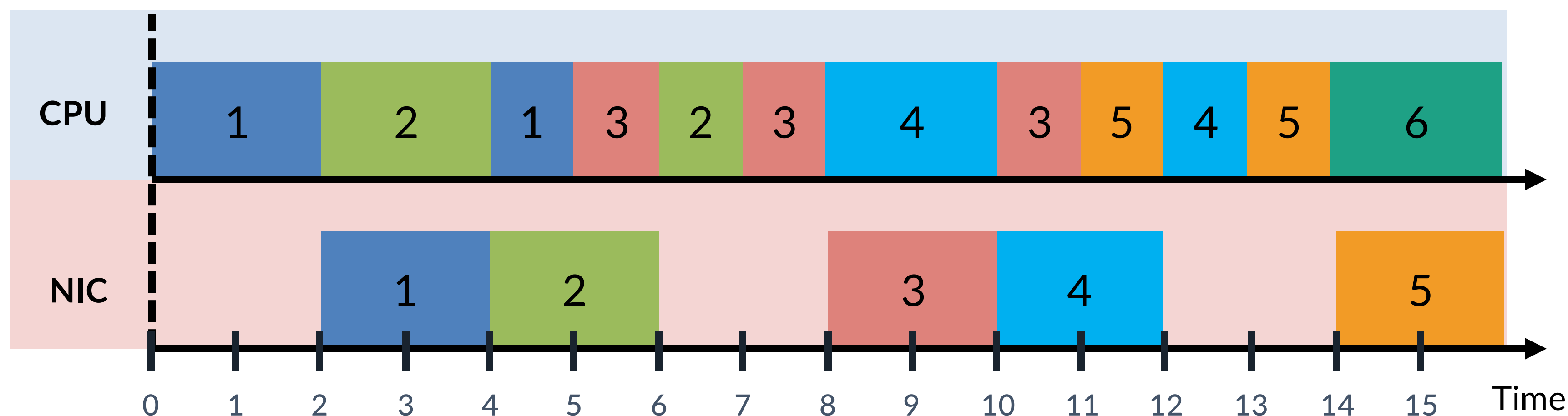


Capacity

when the party ends

Capacity: *maximum throughput that can be sustained by the system.*

With Multi-Programming Level (**MPL**) = 2:



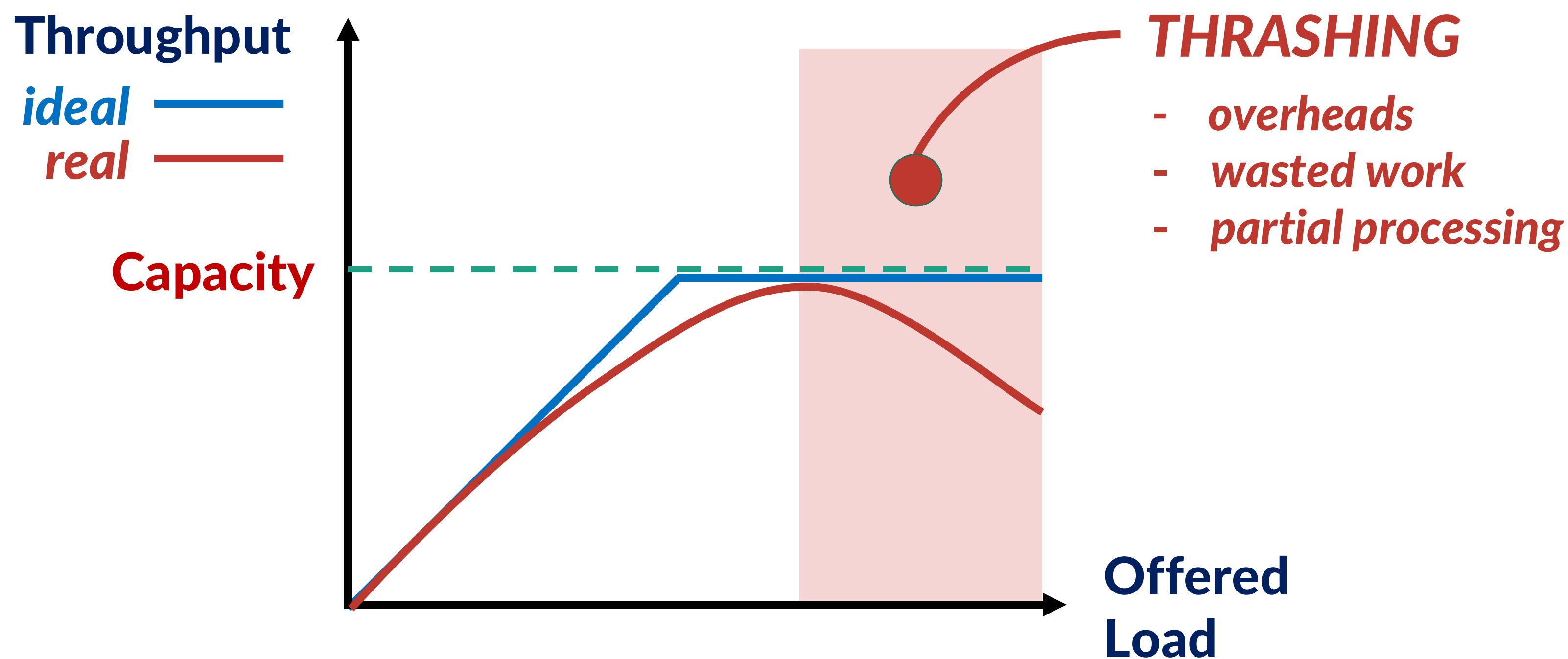
in progress:  **Take-out:** Increasing MPL results in increased **Utilization** and **Throughput** ...until **one of the resources** reaches 100% utilization (**bottleneck**). In this case, the system has reached its **Capacity**.



Capacity

& why reality is harsh

Capacity: *maximum throughput that can be sustained by the system.*





Capacity & Bottlenecks

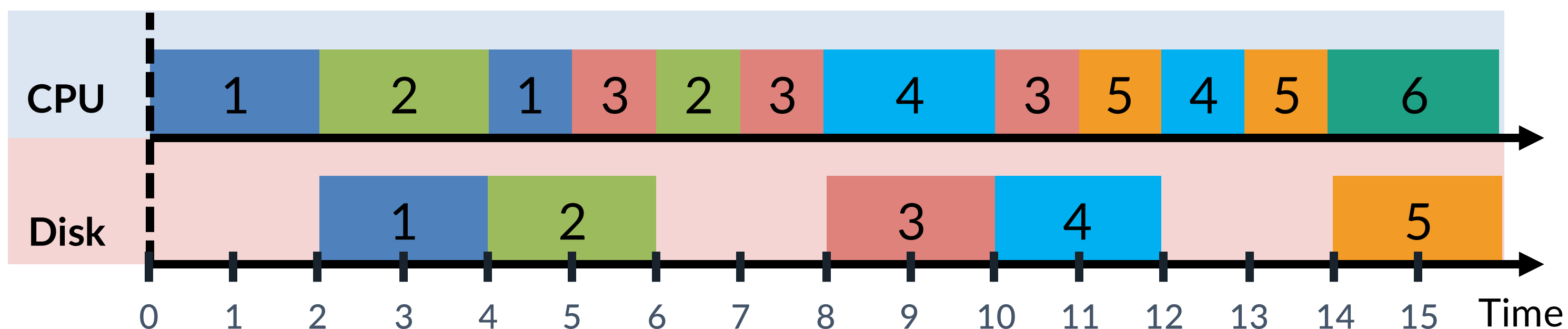
one cannot always blame the butler

Capacity: *maximum throughput that can be sustained by the system.*

Bottleneck: *resource whose utilization first reaches 100% (saturation).*

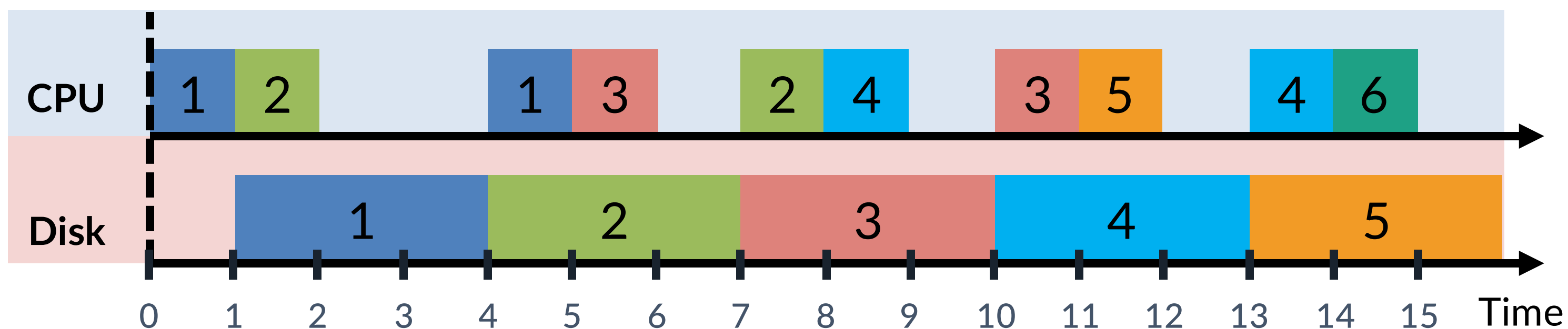
Load #1

**CPU is
bottleneck**



Load #2

**Disk is
bottleneck**

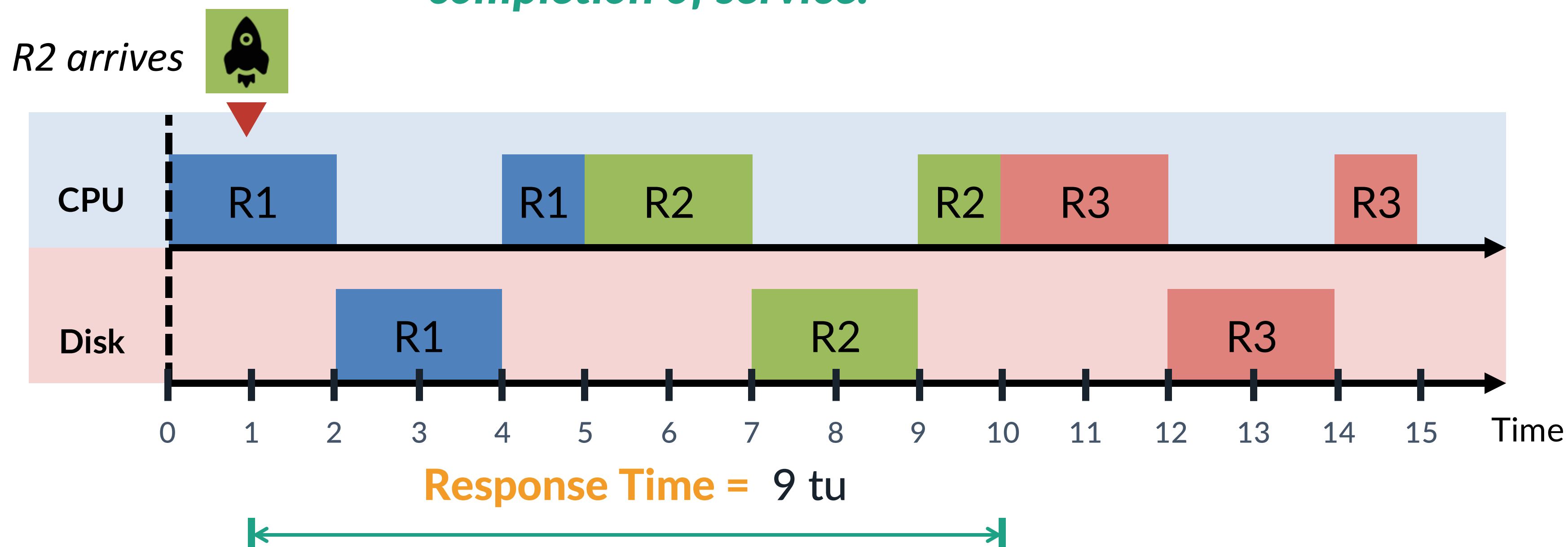




Response Time

time for a system trip

Response Time: *time elapsed between request submission and completion of service.*



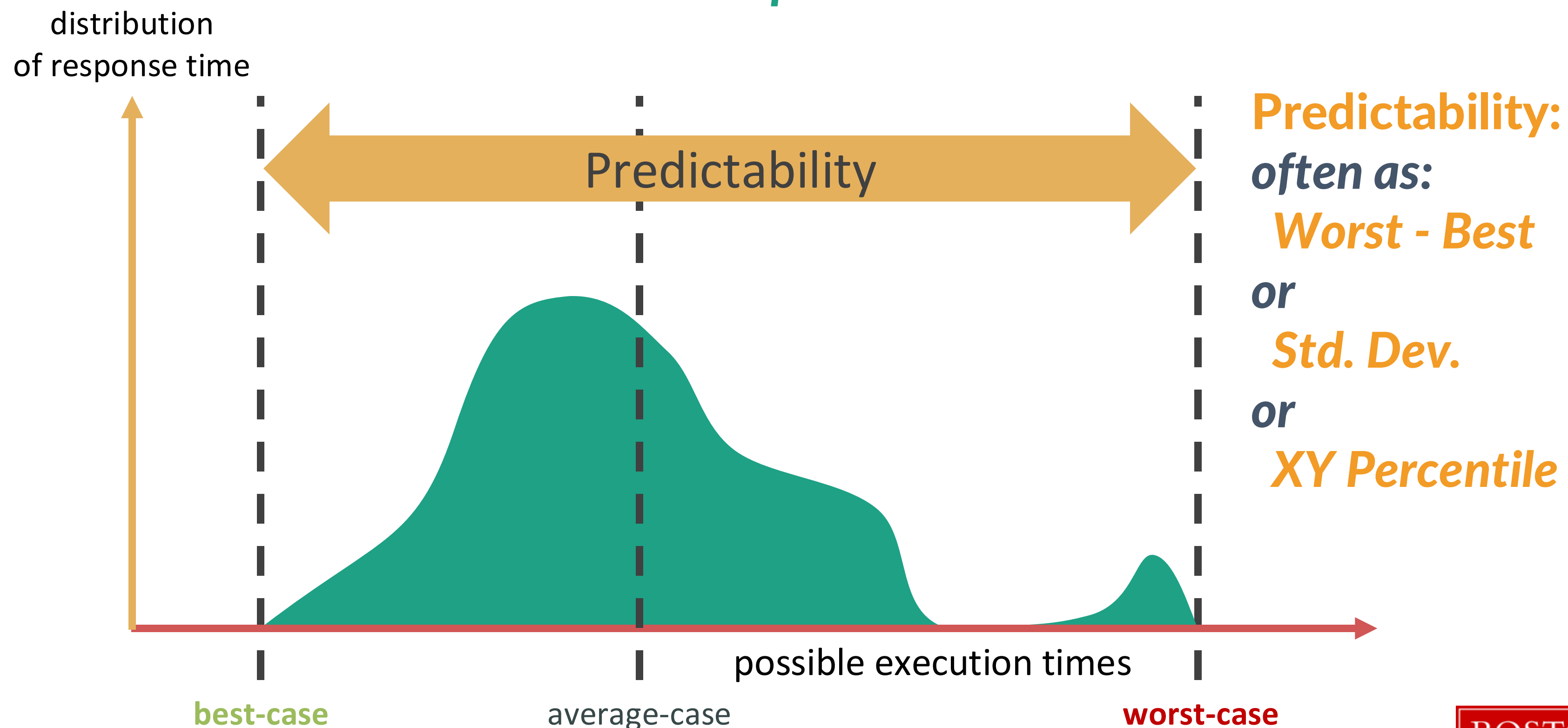
NOTE: Response Time = Turnaround Time = Latency



Predictability

i.e. how “wild” is your system

Predictability: *Difference between best- or average-case response time and worst-case response time.*

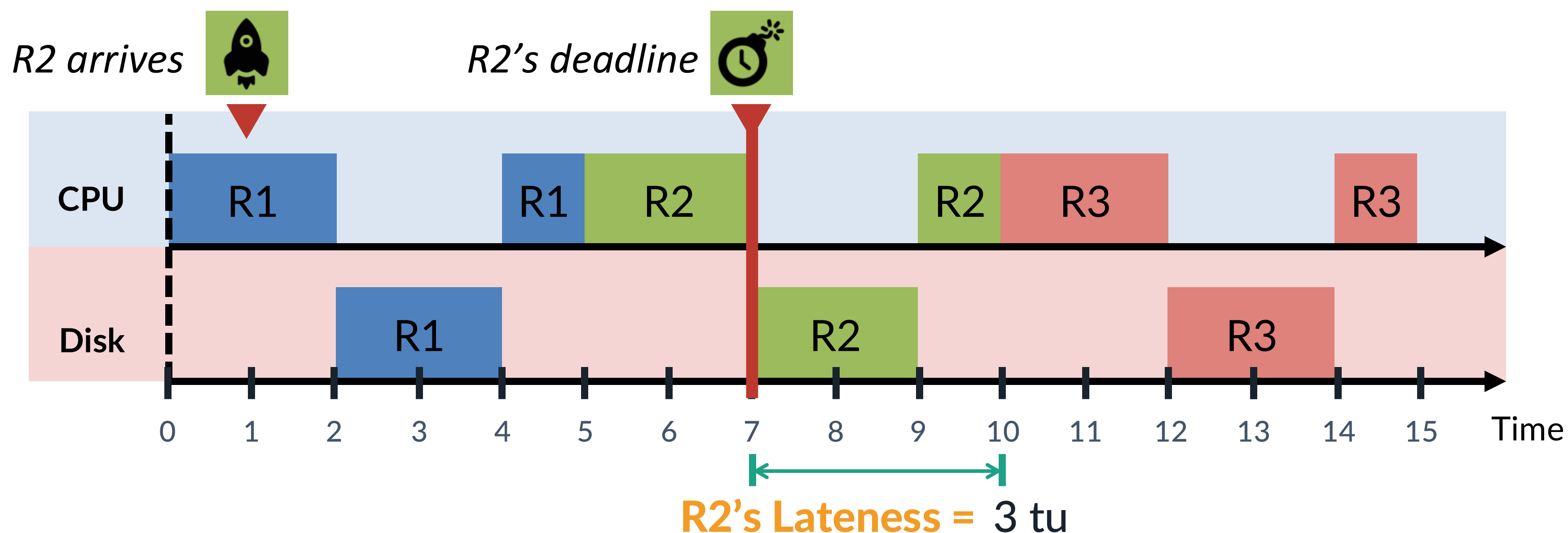




Lateness

how Italian is a system

Lateness: *how late requests complete service with respect to their deadlines.*



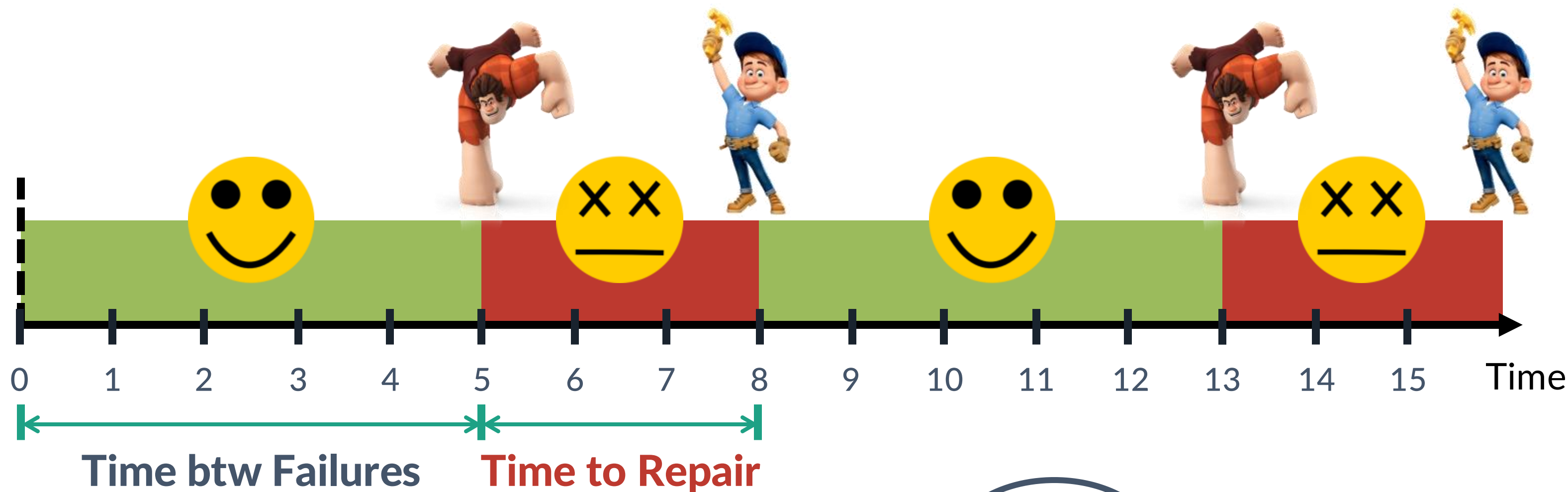
NOTE: For the entire system, compute average lateness



Availability

is the system up?

Availability: *likelihood that the system will be “up” when a request is submitted.*



$$\text{Availability} = \frac{\text{Up Time}}{\text{Total Time}} = \frac{MTBF}{MTBF + MTTR}$$

Mean Time btw Failures

Mean Time to Repair



Reliability

is the system GOING TO BE up?

Reliability: *likelihood that the system will be “up” for a certain amount of time.*



Same **AVAILABILITY** but *different* **RELIABILITY**

Performance Disclaimer

just like a medication's fine print



Performance metrics **abstract away** the **REAL behavior** of the system and can be *misleading*.

E.g.: *Average yearly temperature in Boston is 51°F*

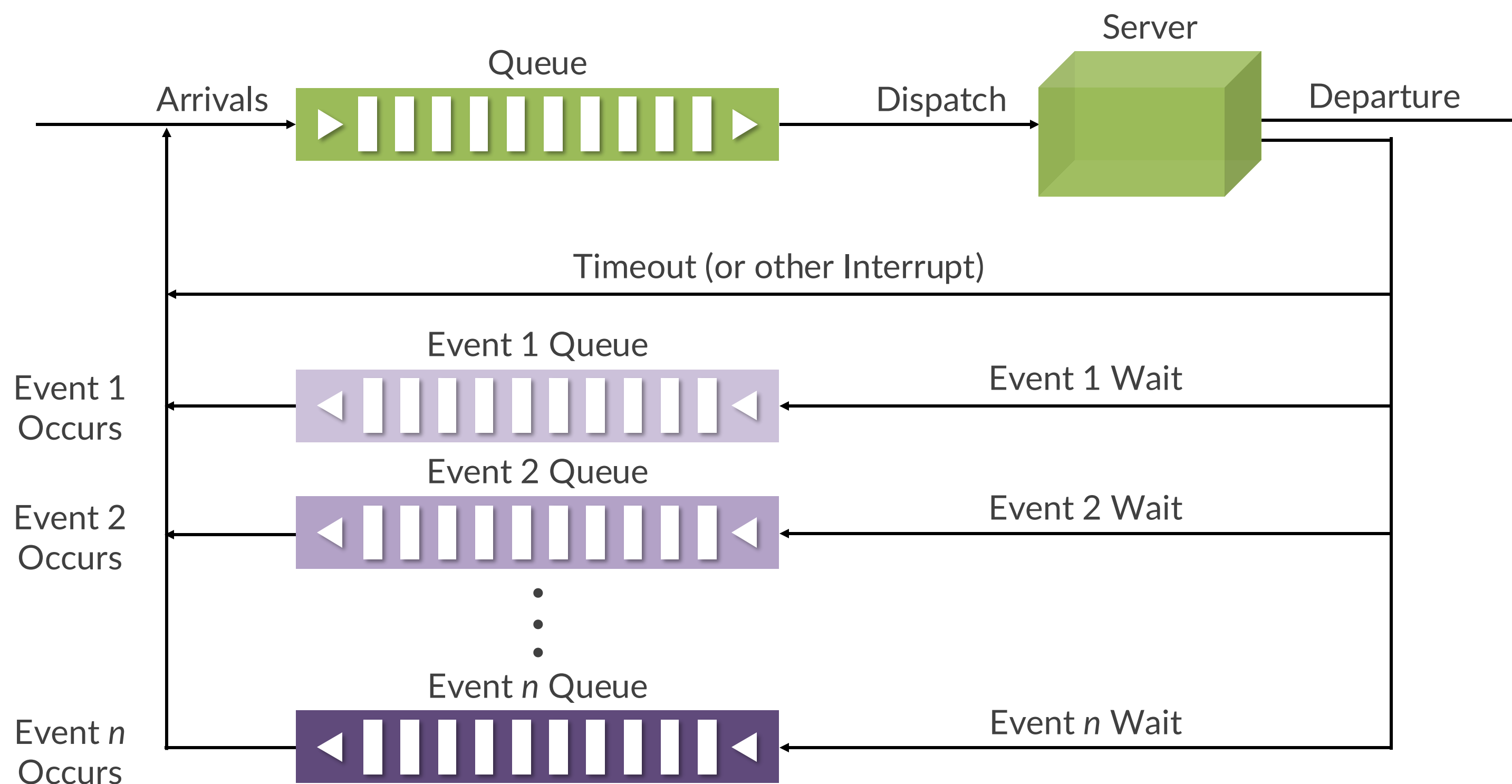
Conclusion: ~~no need to but an A/C unit.~~

E.g.: *The new Intel CPU improves FLOPS by 30%*

Conclusion: ~~all my applications will be 30% faster.~~

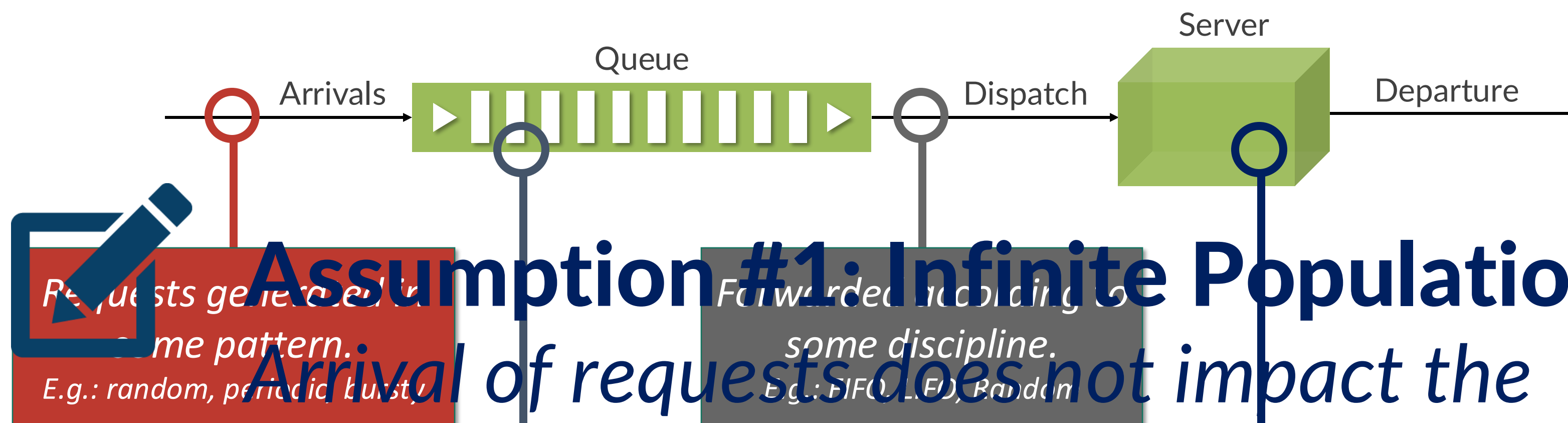
A System Abstraction

as network of queues



A System Abstraction

as network of queues



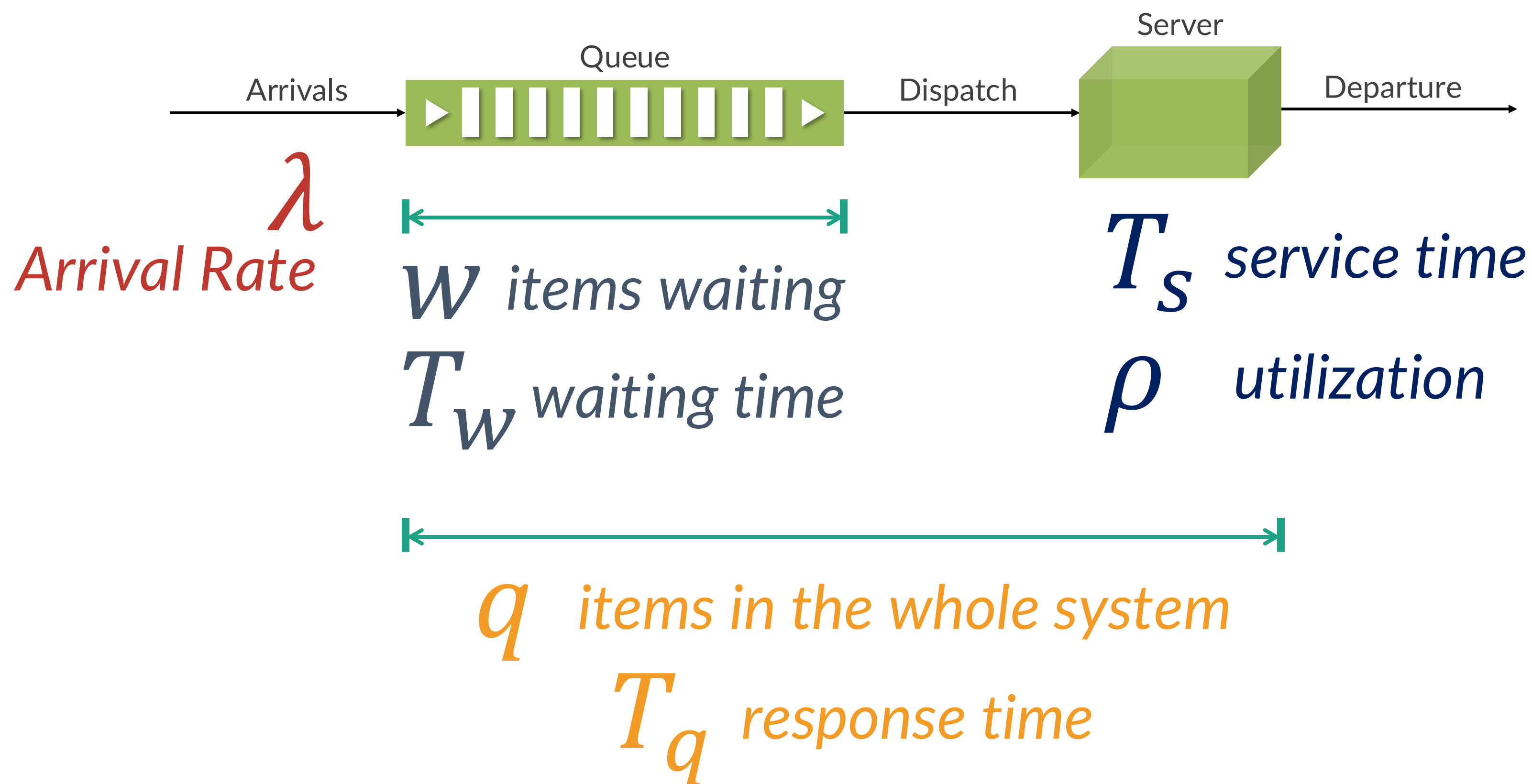
Assumption #2: Unlimited Queues
There is always a spot to “wait” for service. Queues never overflow.

*They are queued-up in some holding space.
E.g.: array, linked list, table*

*They are processed for some time.
E.g.: constant, random, load aware*

A System Abstraction

let's talk about notation



Universal Relationships

dropping some wisdom

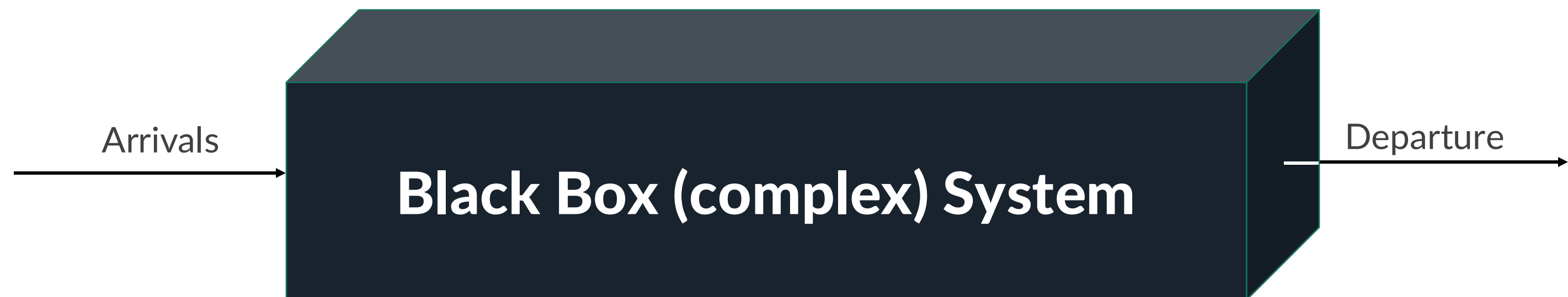


utilization $\rho = \lambda \cdot T_s$ valid if: $\lambda < \frac{1}{T_s}$

response time $T_q = T_w + T_s$

Little's Law

the mighty



Assumption: Input flow rate matches output flow rate.

Implies: system at steady-state; no unaccounted deaths.

$$1 \quad q = \lambda \cdot T_q$$

$$2 \quad w = \lambda \cdot T_w$$

Universal Relationships

dropping some wisdom



$$\begin{aligned}
 T_q &= T_w + T_s \\
 \frac{q}{\lambda} &= \frac{w}{\lambda} + T_s \\
 q &= w + T_s \cdot \lambda
 \end{aligned}$$

*Little's Law,
I choose you*



$$q = w + \rho$$