

P99 CONF

Why User-Mode Threads Are Good for Performance



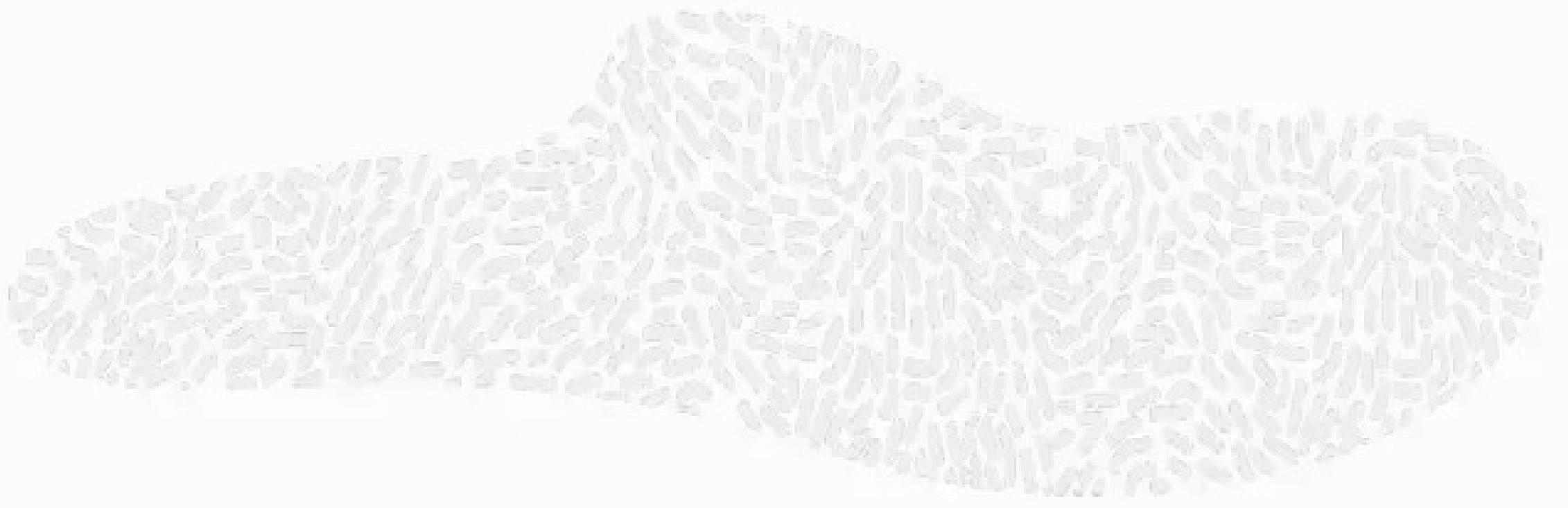
Ron Pressler

Architect, Java Platform Group, Oracle

Brought to you by



Why?



- Why do anything?
- ~~Why do this rather than that?~~

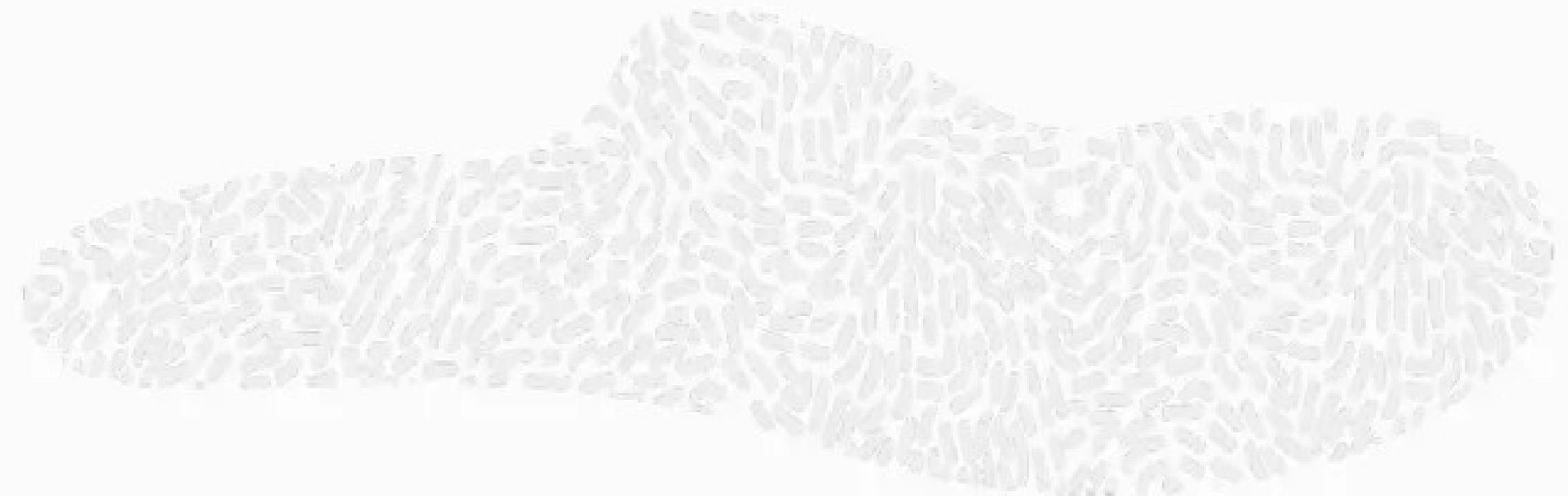
Concurrency



Concurrency and Parallelism

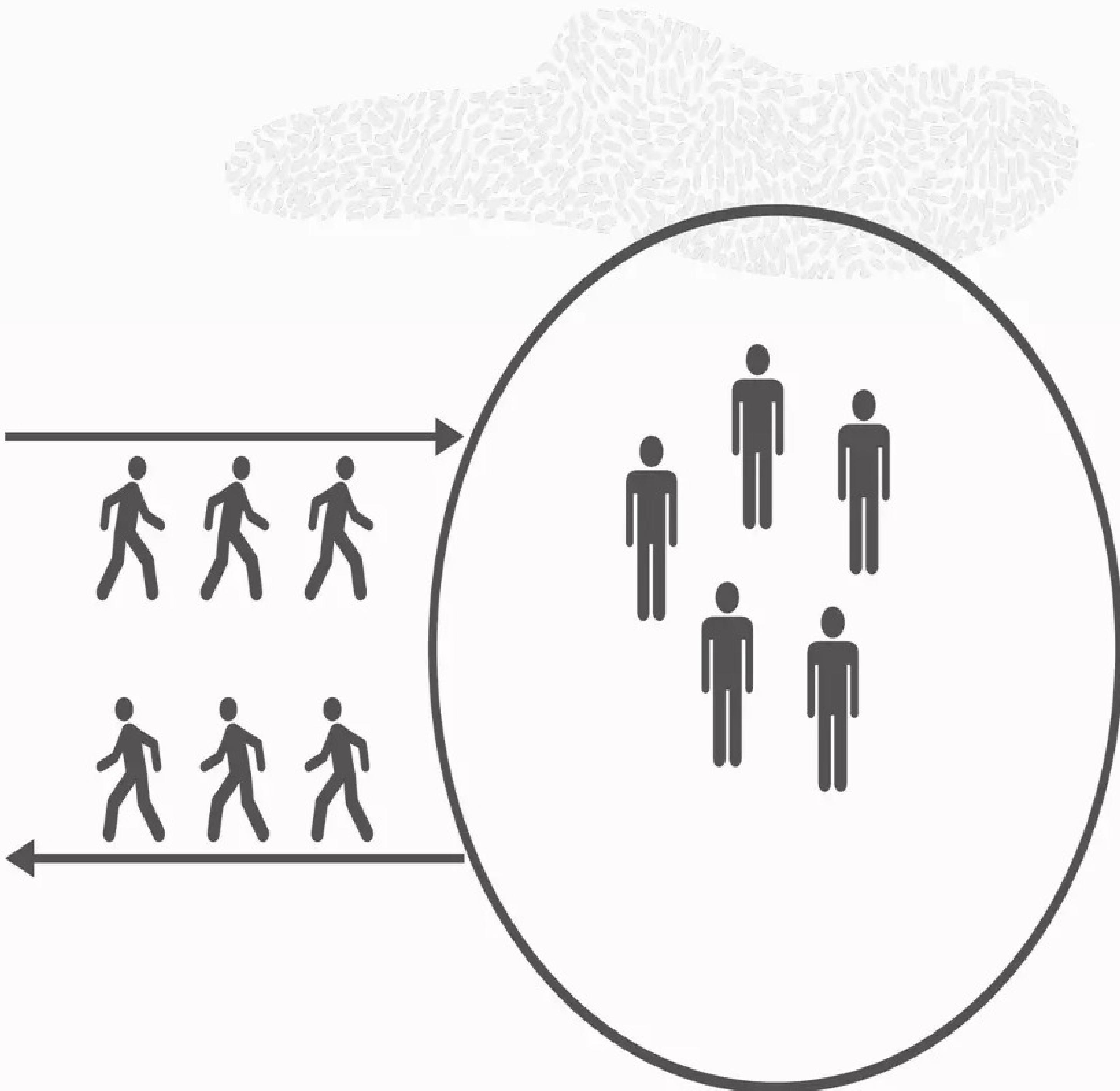
- *Parallelism:*
 - Speed up a task by splitting it into *cooperating* subtasks scheduled onto multiple available computing resources.
 - Performance measure: latency (time duration)
- *Concurrency:*
 - Schedule available computing resources to multiple largely independent tasks that *compete* over them.
 - Performance measure: throughput (task/time unit)

Little's Law



Little's Law

In any **stable** system



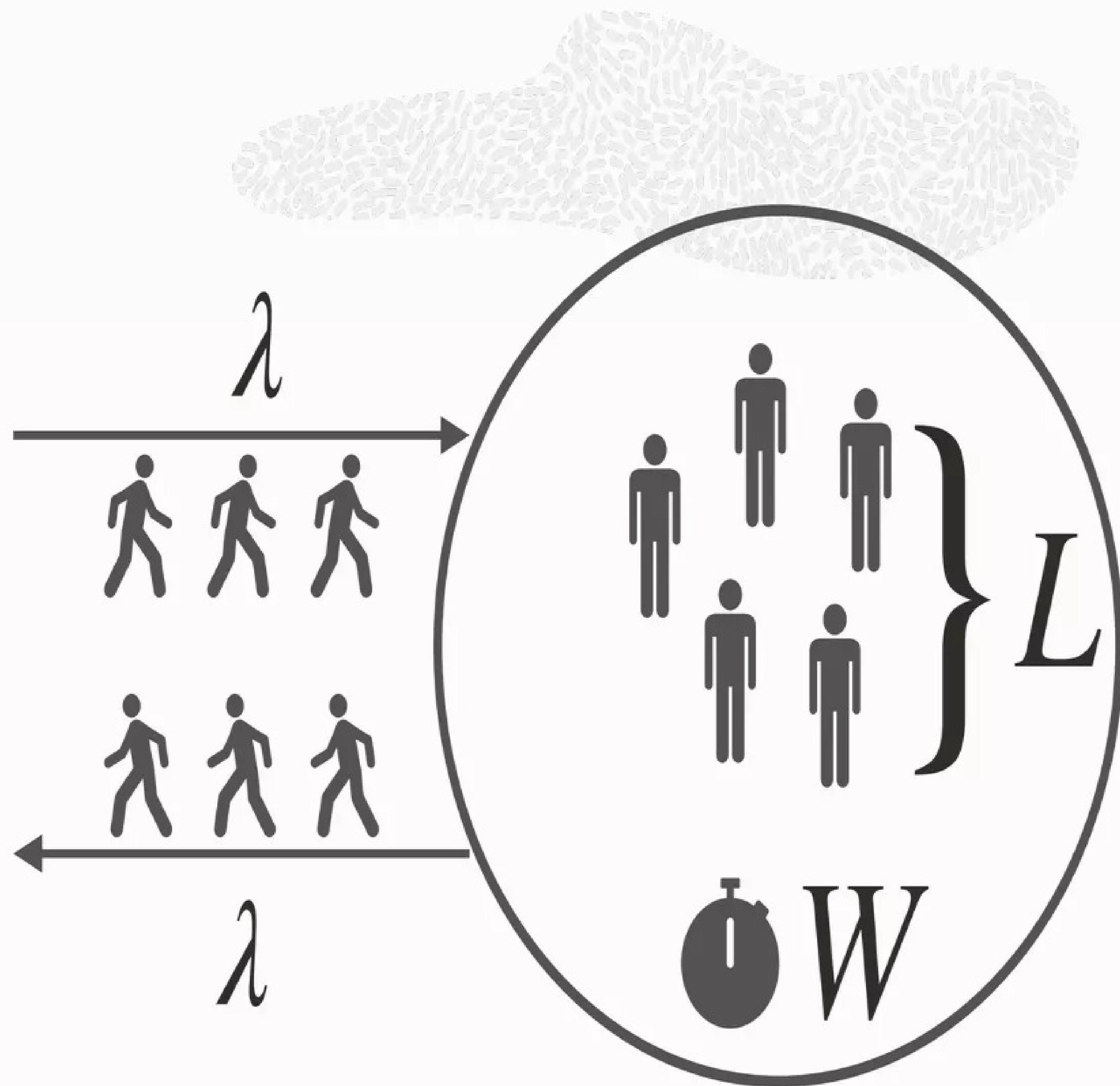
Little's Law

In any **stable** system
with long term averages:

λ — arrival rate = exit rate
= throughput

W — duration inside

L — no. items inside



Little's Law

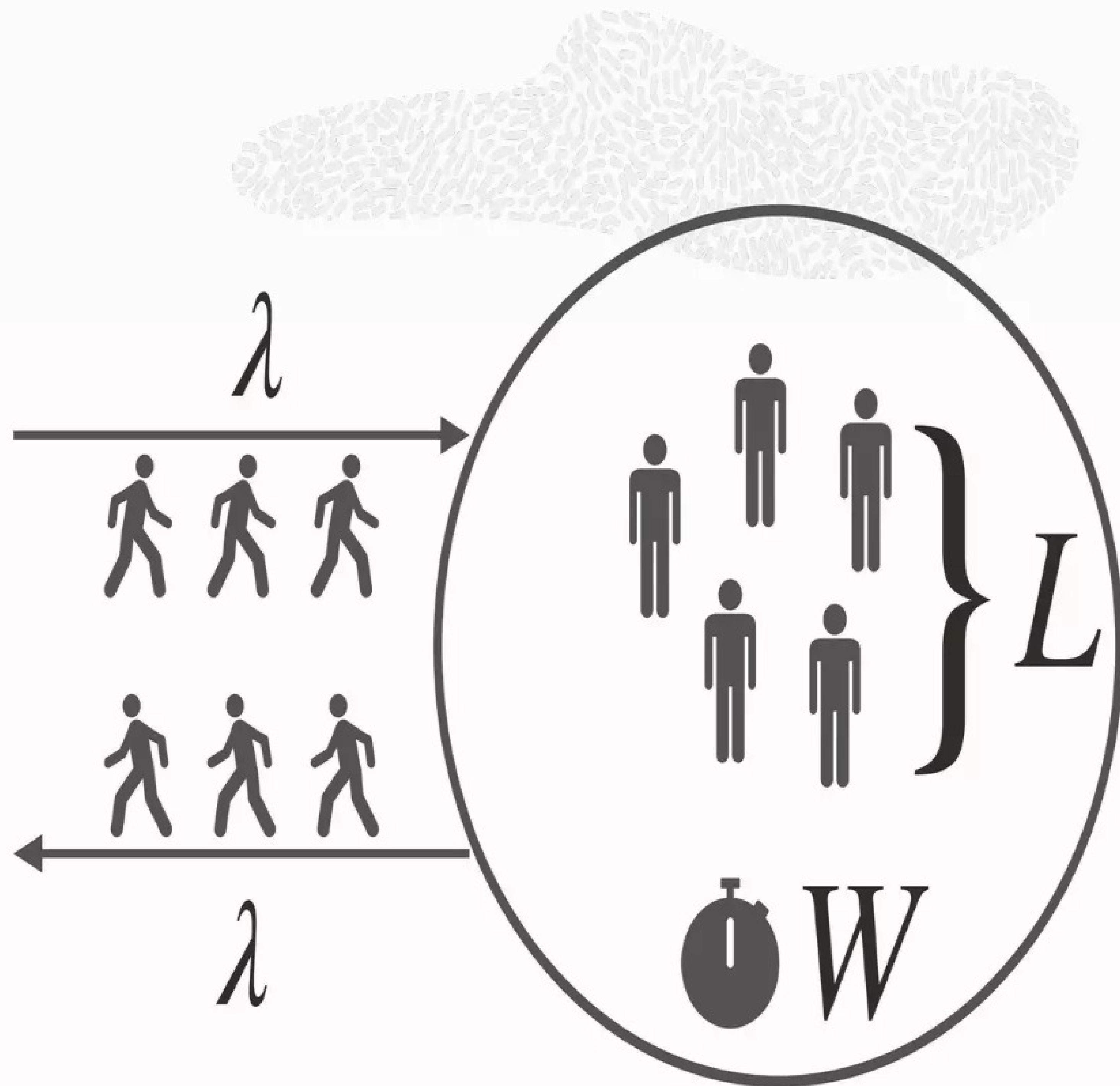
In any **stable** system
with long term averages:

λ — arrival rate = exit rate
= throughput

W — duration inside

L — no. items inside

$$L = \lambda W$$

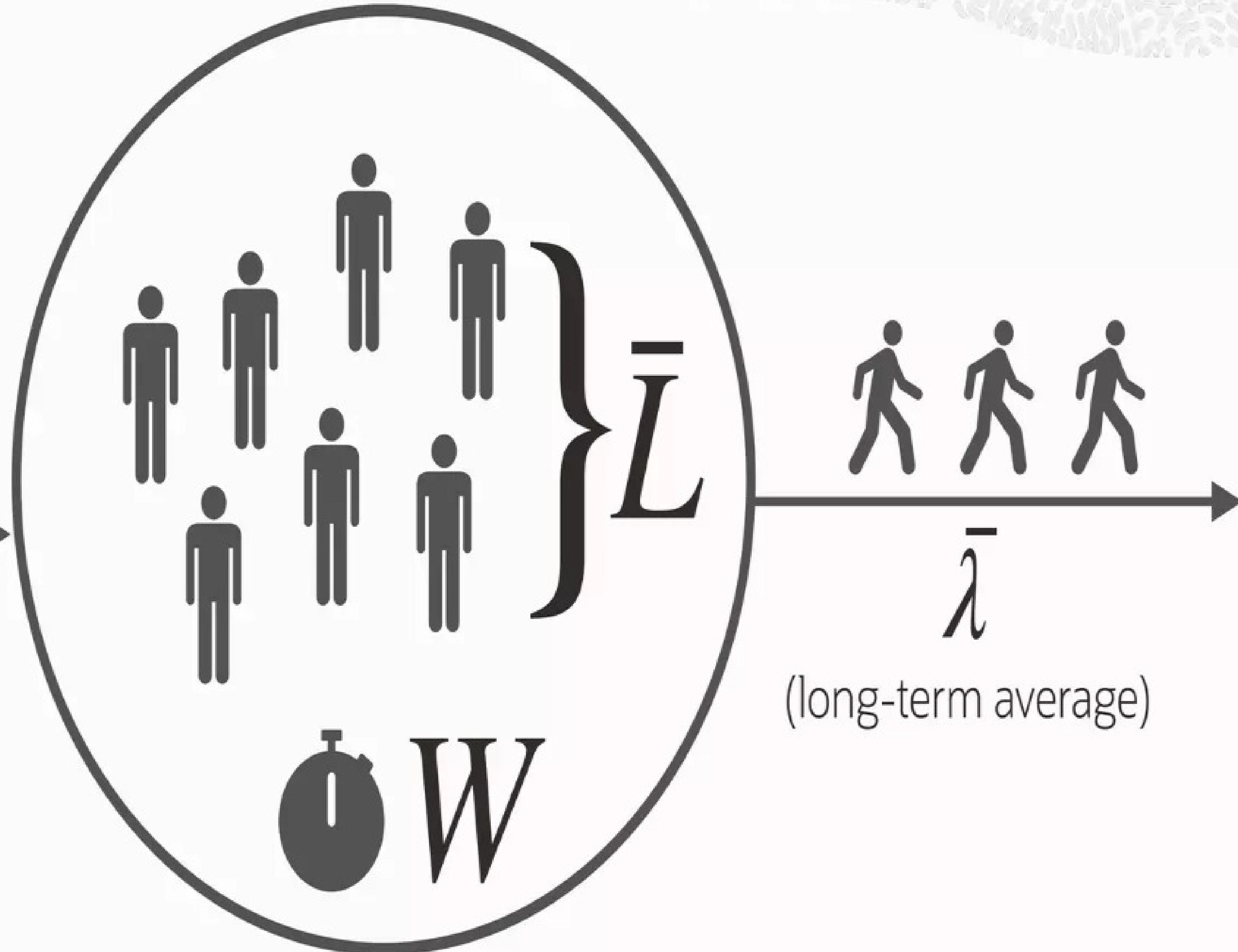


Little's Law

\bar{L} — Capacity

$$\frac{\text{---}}{\bar{\lambda}}$$

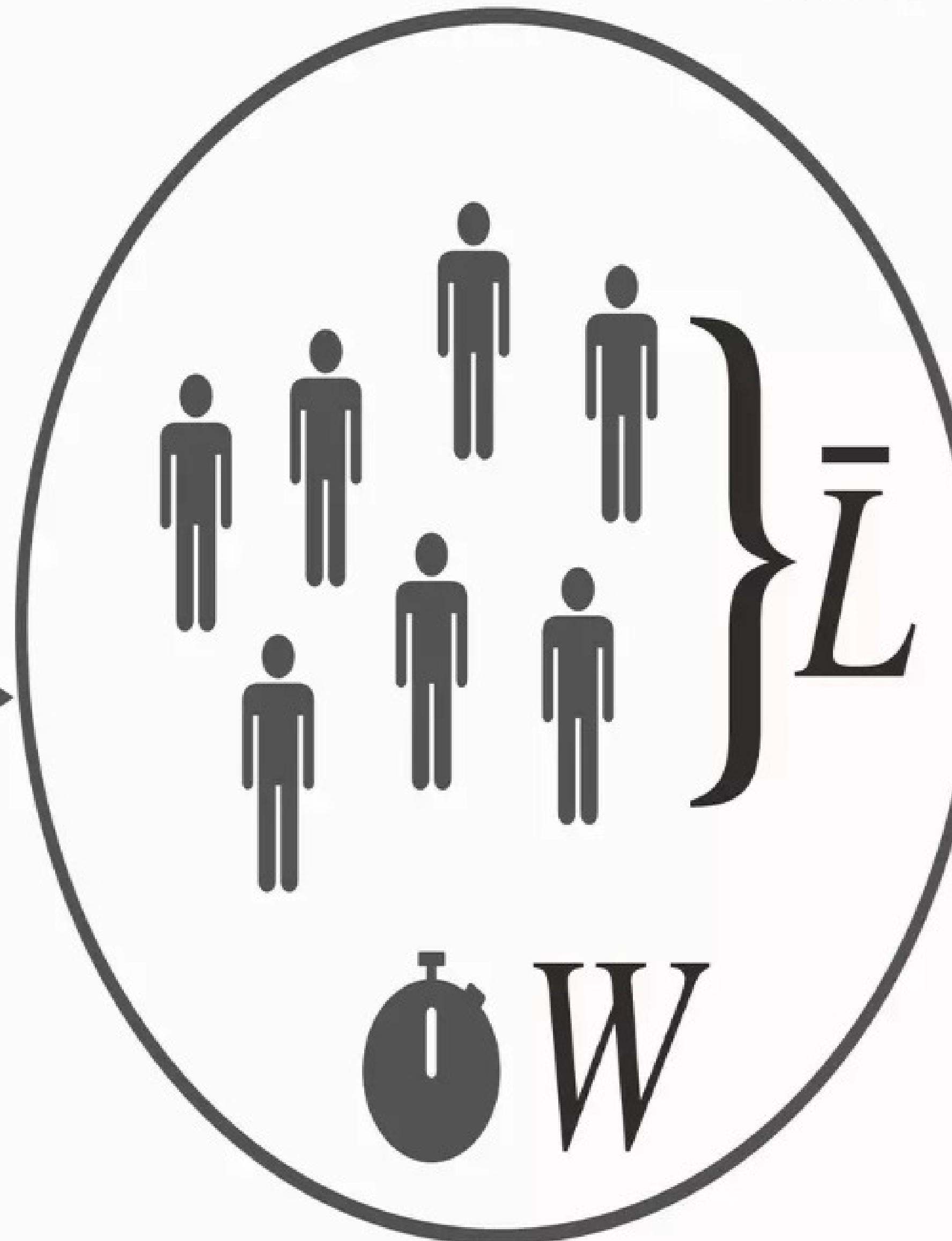
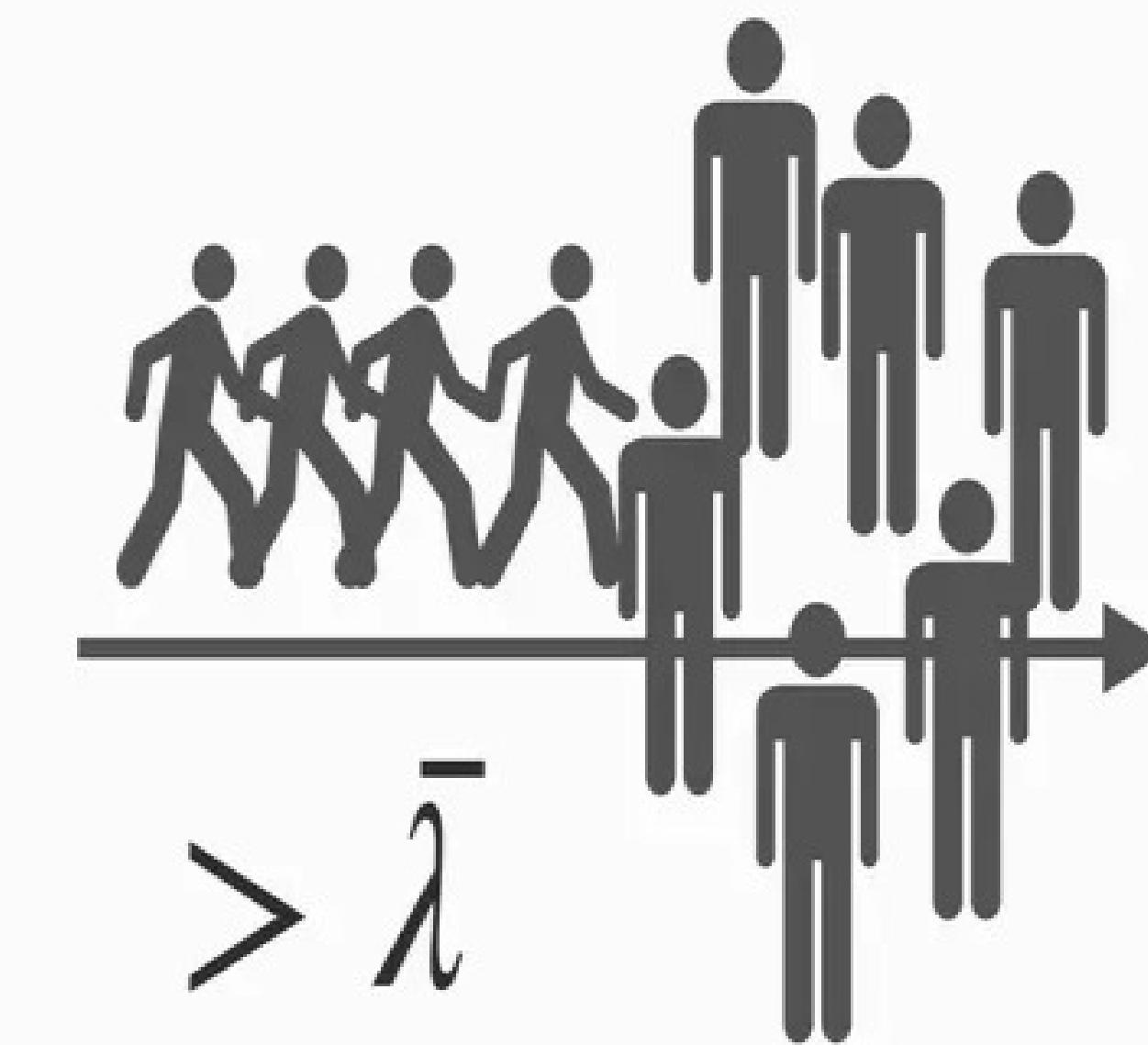
(long-term average)



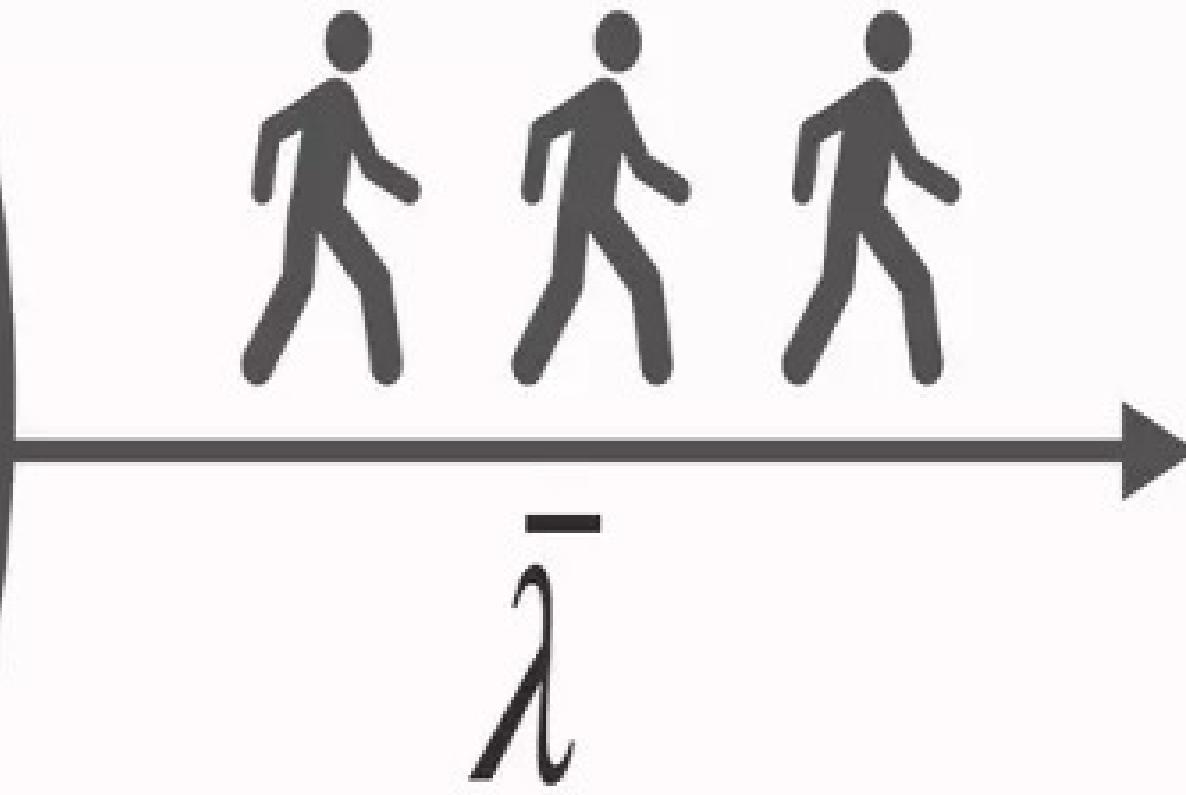
$$\bar{L} = \bar{\lambda}W$$

Little's Law

\bar{L} — Capacity



(long-term average)

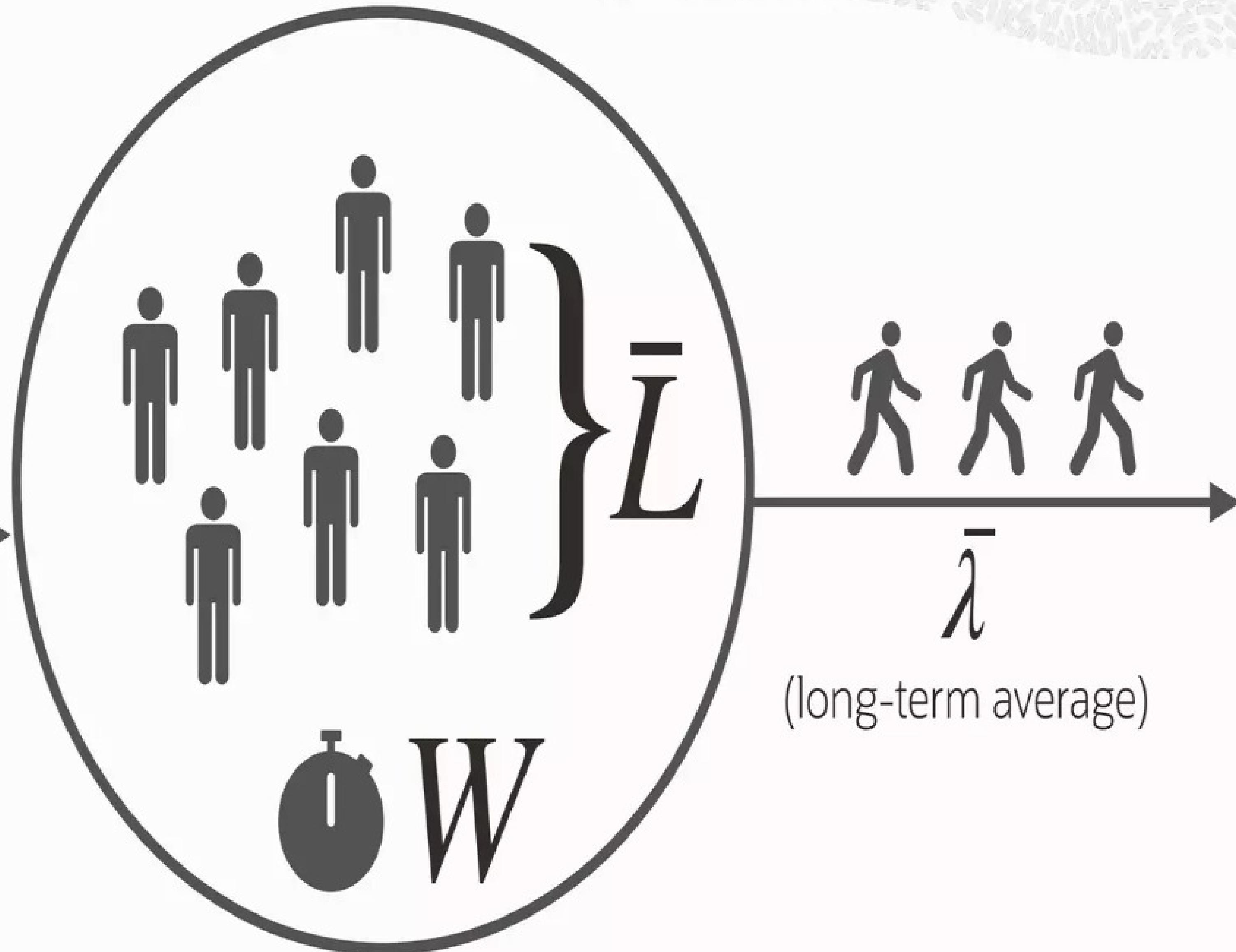


$$\bar{L} = \bar{\lambda}W$$

Little's Law

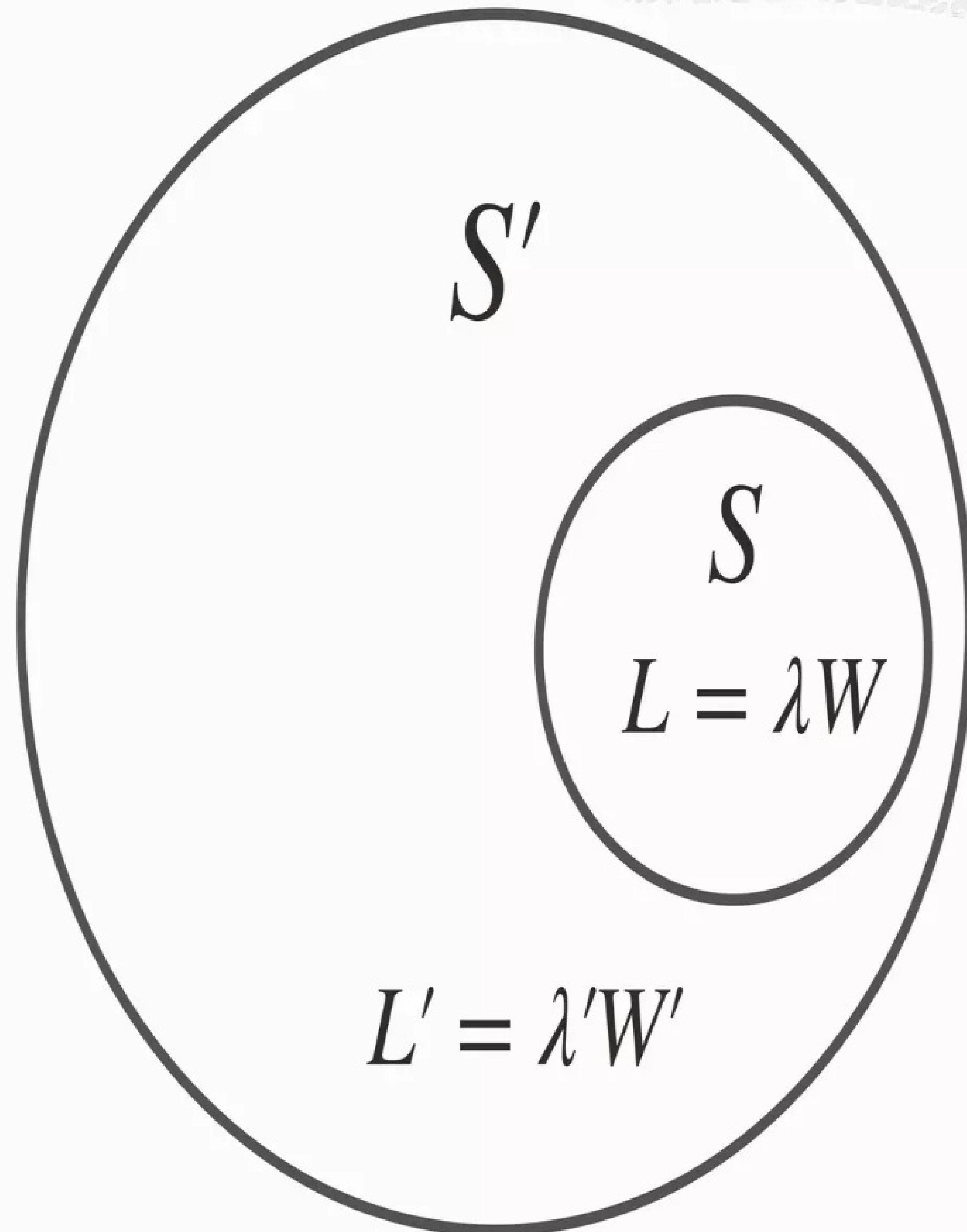
\bar{L} — Capacity

$\frac{\text{Person}}{\text{Time}} < \bar{\lambda}$
(momentary)



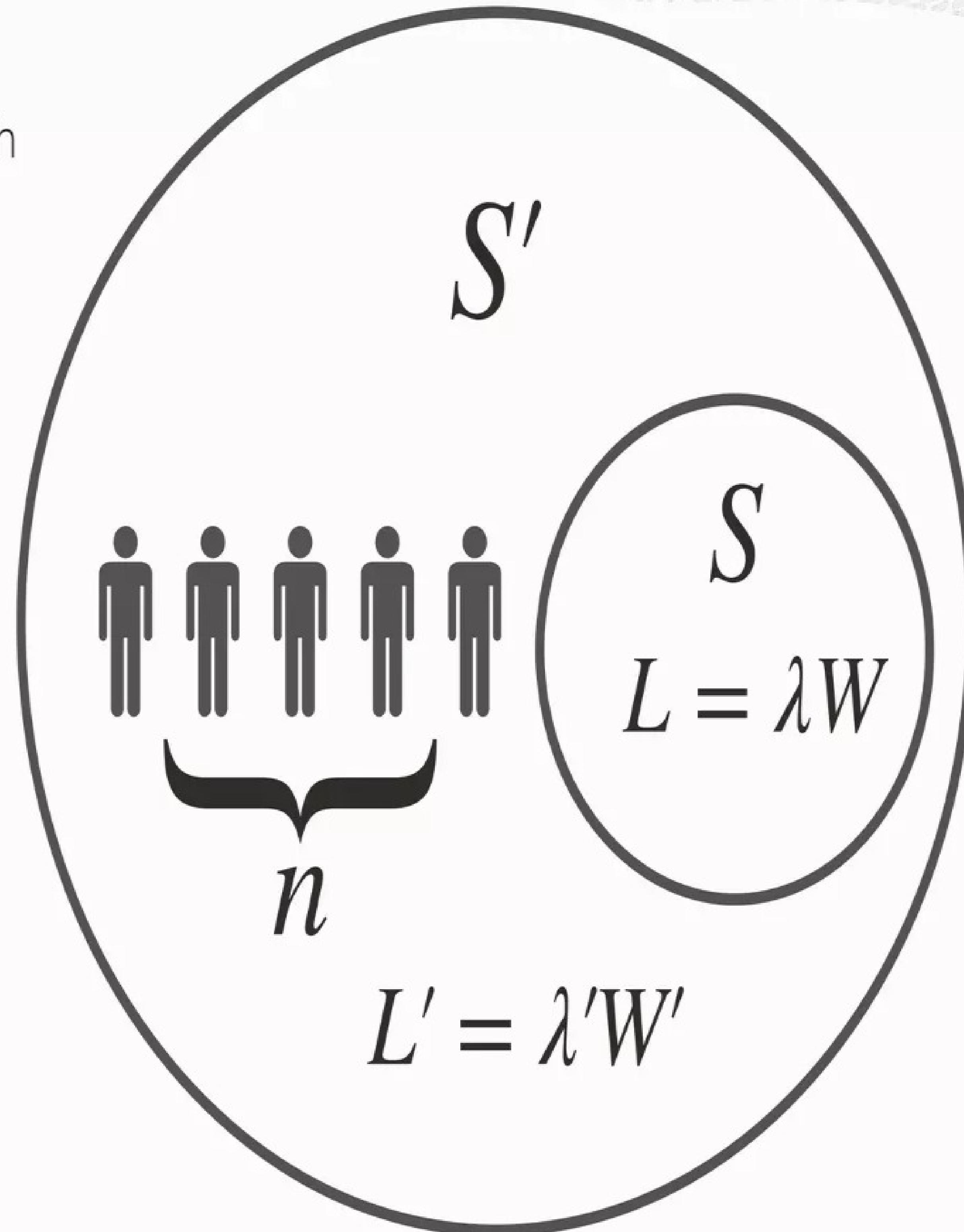
$$\bar{L} = \bar{\lambda}W$$

Little's Law – subsystems



Little's Law – subsystems and queues

n — Average queue length



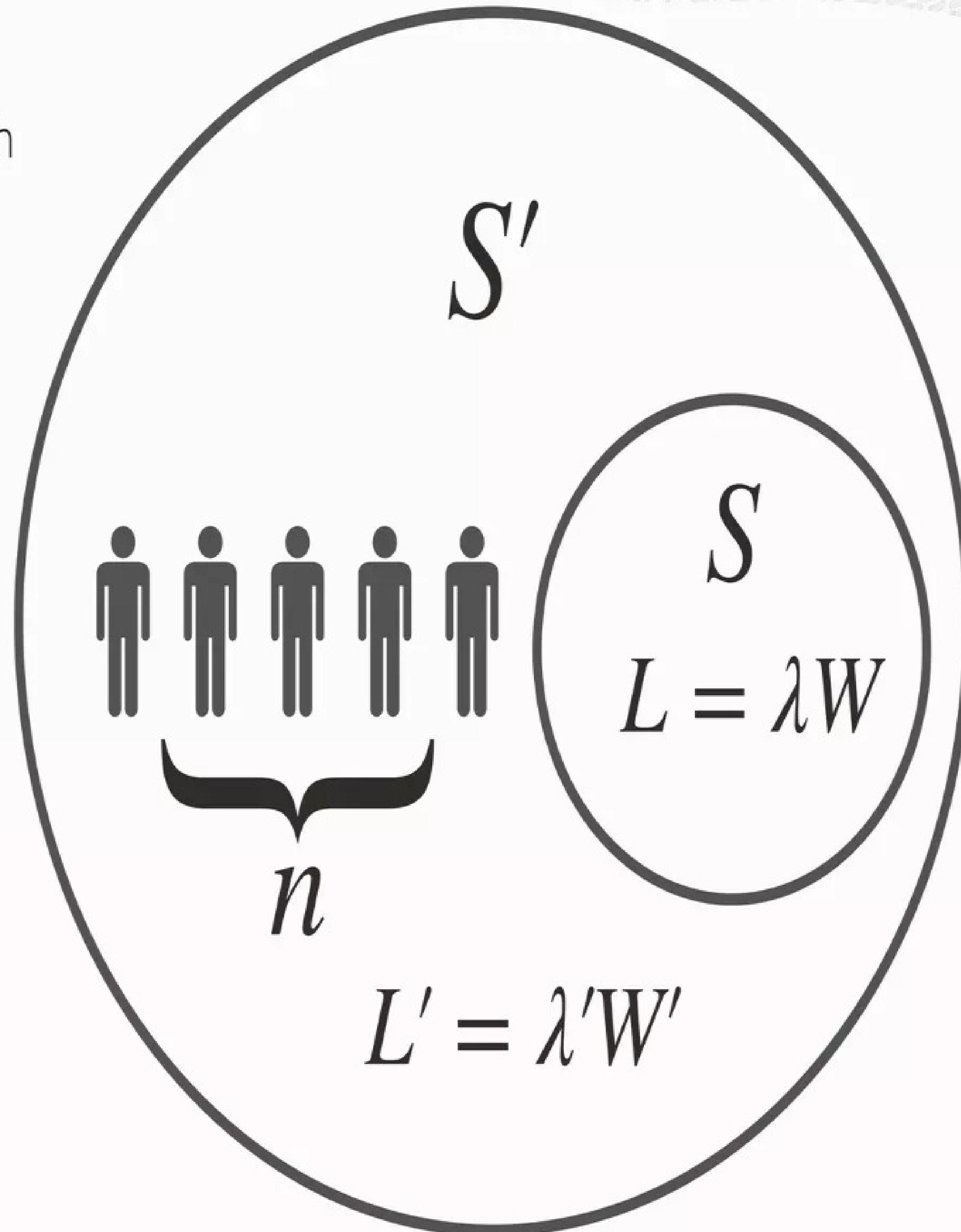
Little's Law – subsystems and queues

n — Average queue length

$$\lambda' = \lambda$$

$$L' = L + n$$

$$W' = W + n \frac{W}{L}$$



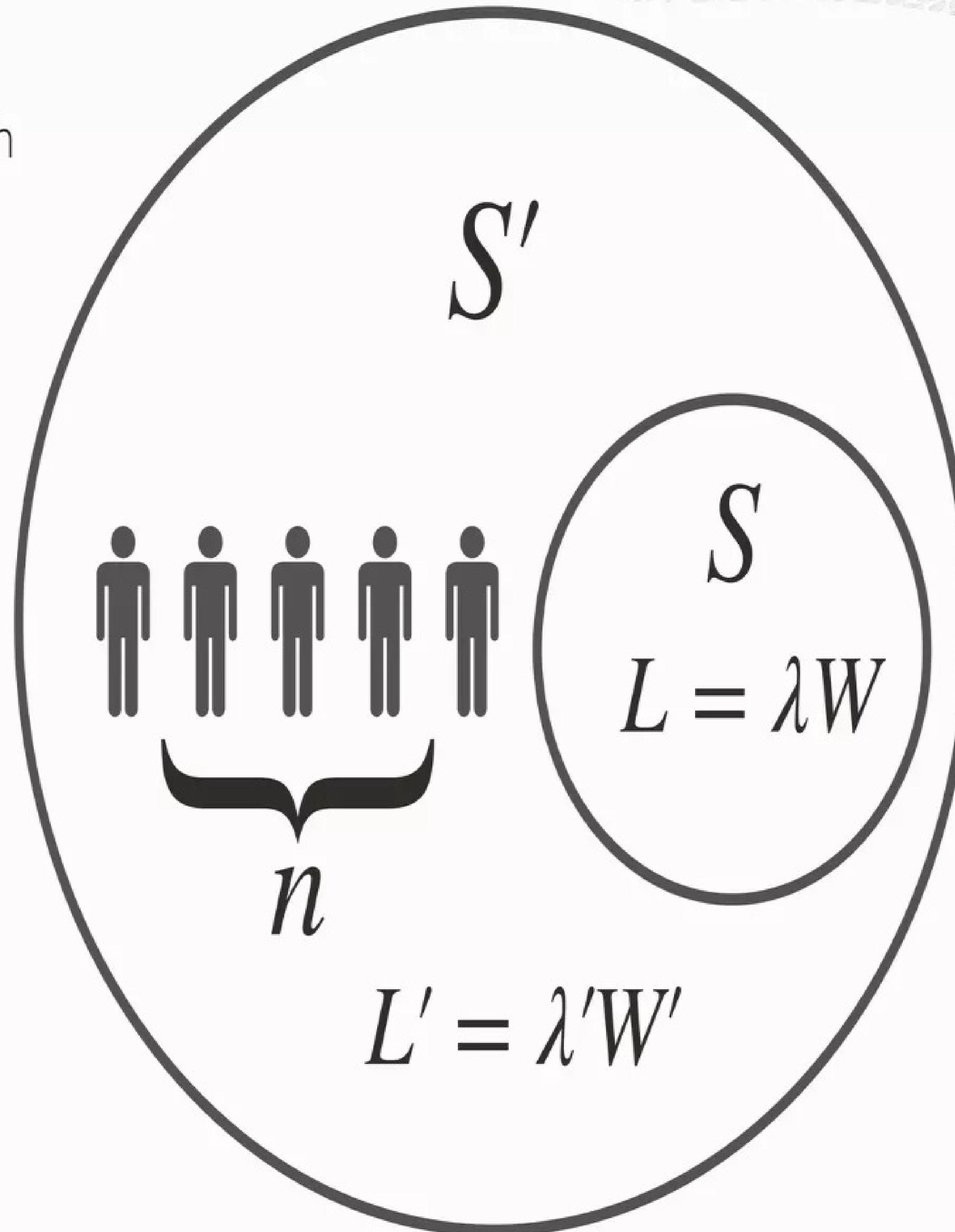
Little's Law – subsystems and queues

n — Average queue length

$$\lambda' = \lambda$$

$$L' = L + n$$

$$W' = W + n \frac{W}{L}$$



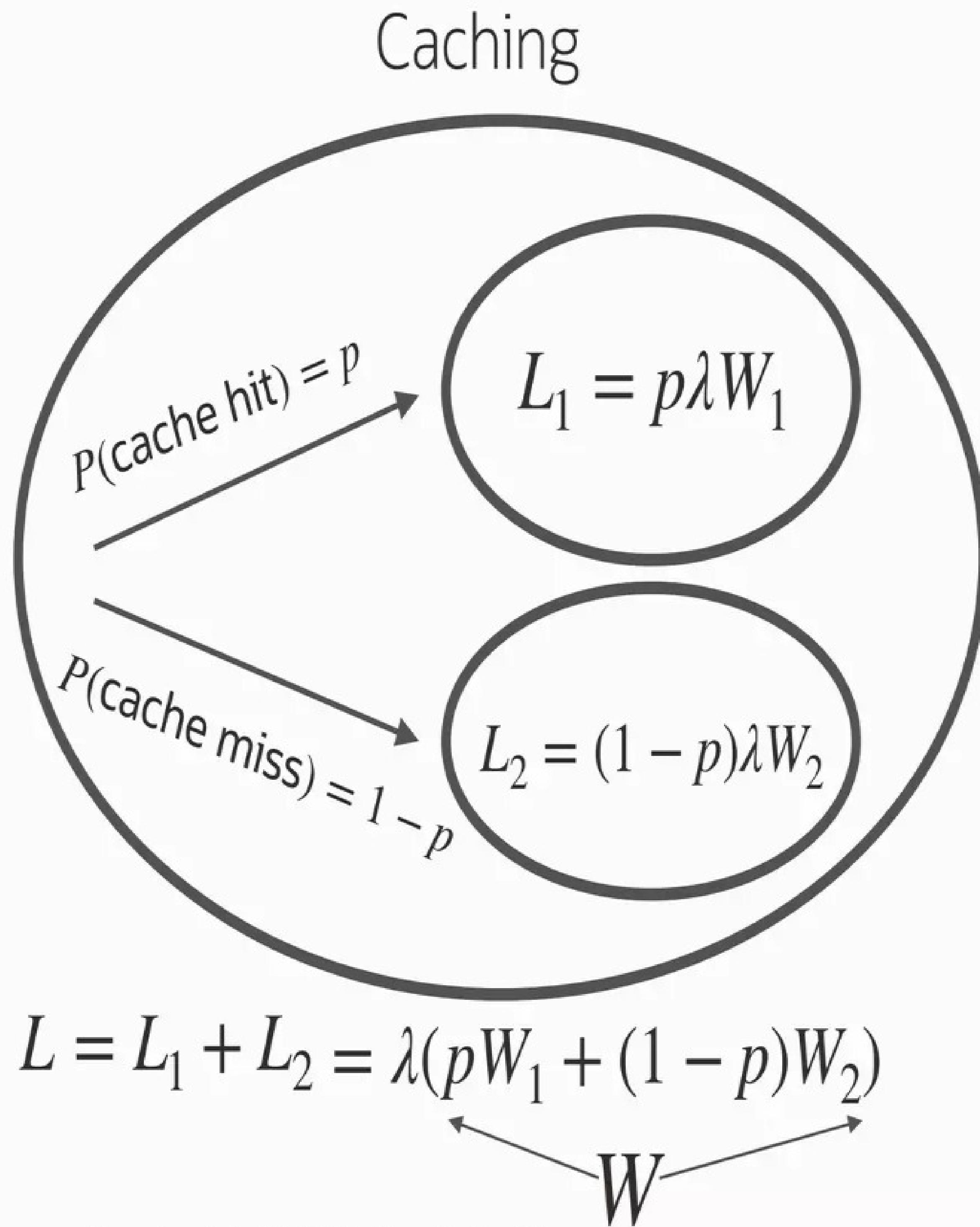
$$L' = \lambda'W'$$

$$L + n = \lambda(W + n \frac{W}{L})$$

$$\cancel{L(1 + \frac{n}{L})} = \lambda W(1 + \frac{n}{L})$$

$$L = \lambda W$$

Little's Law – other applications



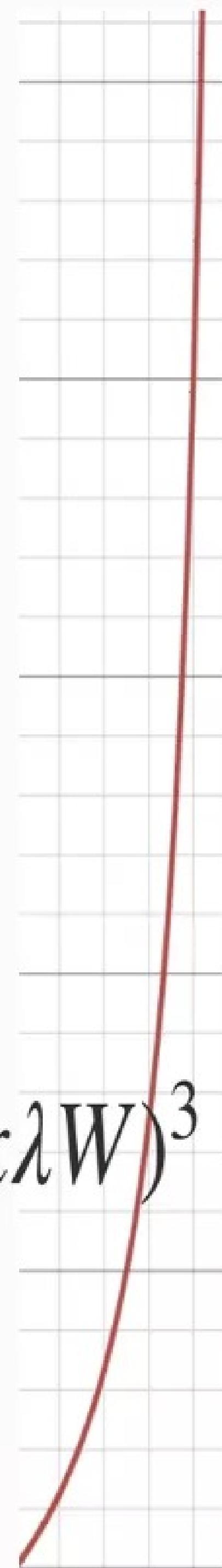
Interference

x – Interference coefficient

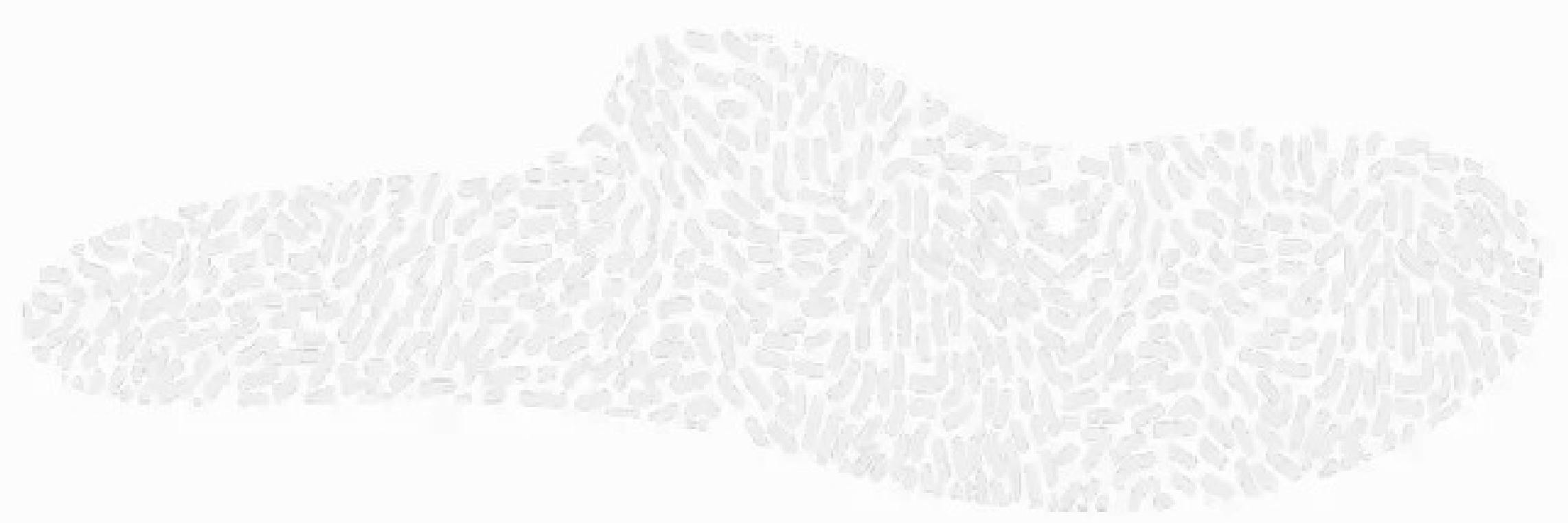
$$L = \lambda(1 + xL)W$$

$$L = \frac{\lambda W}{1 - x\lambda W}$$

$$= \lambda W + (x\lambda W)^2 + (x\lambda W)^3 + \dots$$



Capacity and throughput



$\bar{\lambda}$ — Maximum throughput

\bar{L} — Capacity

W — Average latency

$$\bar{\lambda} = \frac{\bar{L}}{W}$$

CPU capacity



$\bar{\lambda}$ — Maximum throughput

\bar{L}_{CPU} — #cores

W_{CPU} — Average CPU consumption

$$\bar{\lambda} = \frac{\bar{L}_{CPU}}{W_{CPU}}$$

$$\bar{L} = 30 \text{ cores}$$

$$W_{CPU} = 100\mu s = 1 \times 10^{-4}s \quad (> 1500 \text{ cache misses})$$

$$\bar{\lambda} = 30,000 \text{ req/s}$$

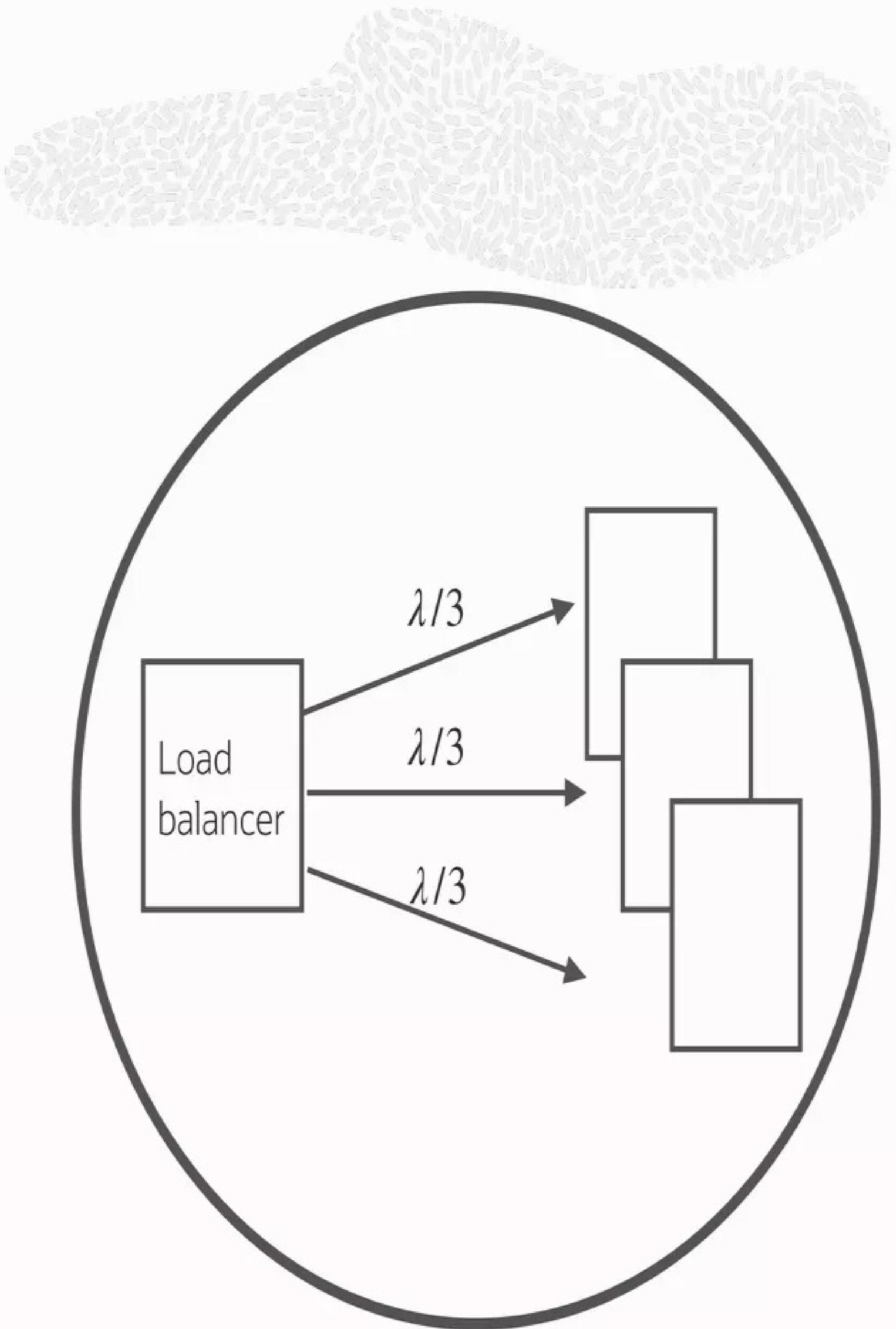
Capacity and throughput

$\bar{\lambda}$ — Maximum throughput

\bar{L} — Capacity

W — Average latency

$$\bar{\lambda} = \frac{\bar{L}}{W}$$



Thread-per-Request – Thread Capacity

Request: The domain's unit of concurrency

Thread: The software unit of concurrency

Thread-per-Request: A request consumes a thread
(either new or borrowed from a pool) for its duration, W

Thread-per-Request – Thread Capacity

Request: The domain's unit of concurrency

Thread: The software unit of concurrency

$$\# \text{threads} = L = \lambda W$$

Thread-per-Request: A request consumes a thread
(either new or borrowed from a pool) for its duration, W

Thread-per-Request – Thread Capacity

Request: The domain's unit of concurrency

Thread: The software unit of concurrency

$$\# \text{threads} = L = \lambda W$$

Thread-per-Request: A request consumes a thread
(either new or borrowed from a pool) for its duration, W

with parallel fanout

c — average fanout, i.e. average number of threads consumed by a request

$$\# \text{threads} = cL = \lambda \left(\frac{W}{c} \right) = \lambda W$$

⇒ For the purpose of estimating #threads, we can consider W to be the sum of all fanout latencies, even if they're done in parallel.

CPU vs. Thread Capacity with thread/req

$$\bar{L} = \bar{\lambda}W$$

$$= \left(\frac{\bar{L}_{CPU}}{W_{CPU}} \right) W$$

$$= \bar{L}_{CPU} \left(\frac{W}{W_{CPU}} \right)$$

$$\bar{L}_{CPU} = \#\text{cores}$$

$$W_{CPU} = 100\mu s$$

(> 1500 cache misses) \Rightarrow 100 threads per core

$$W = 10ms$$

$$W_{CPU} = 50\mu s$$

(> 800 cache misses) \Rightarrow 1000 threads per core

$$W = 50ms$$

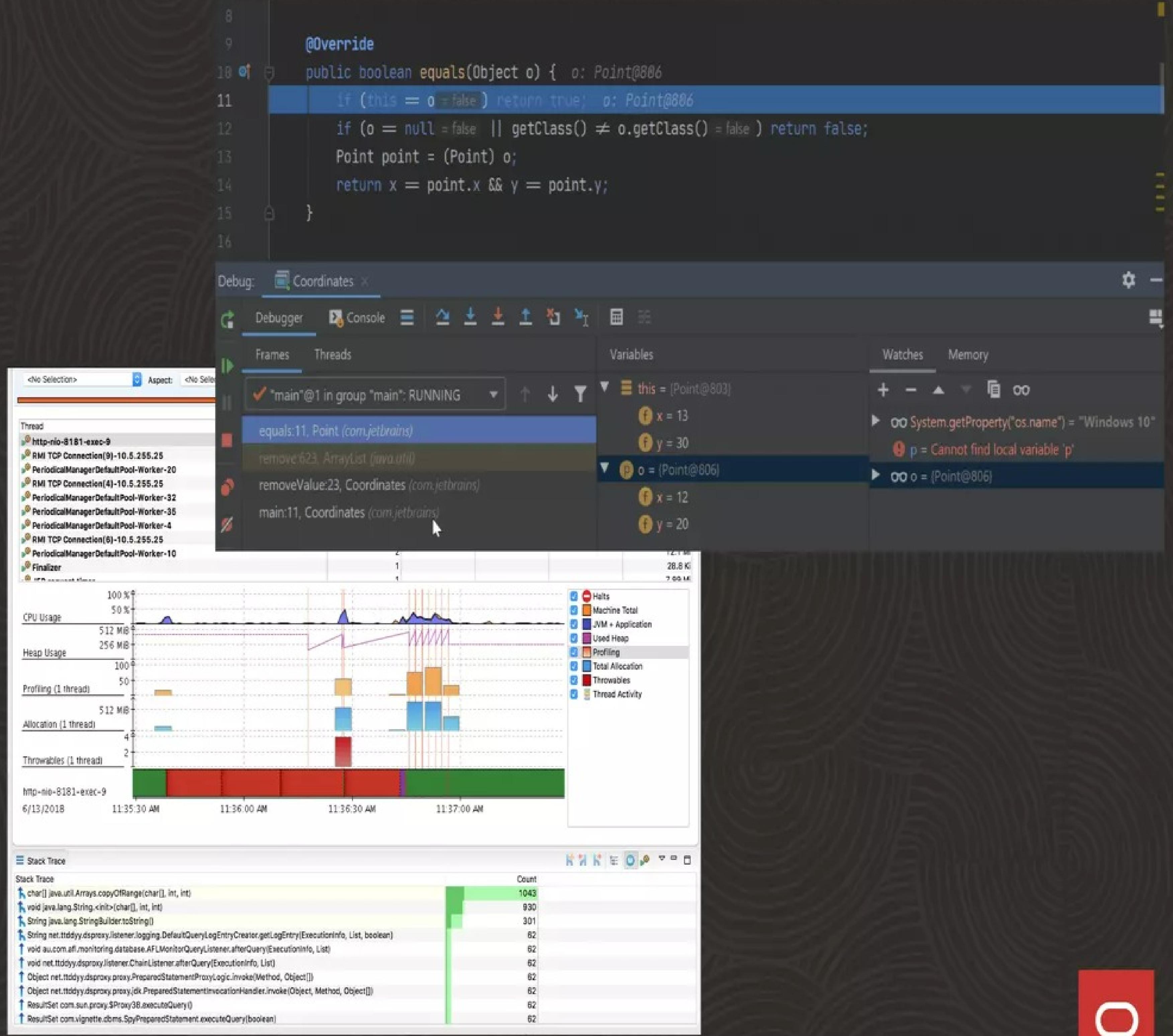
$$W_{CPU} = 10\mu s$$

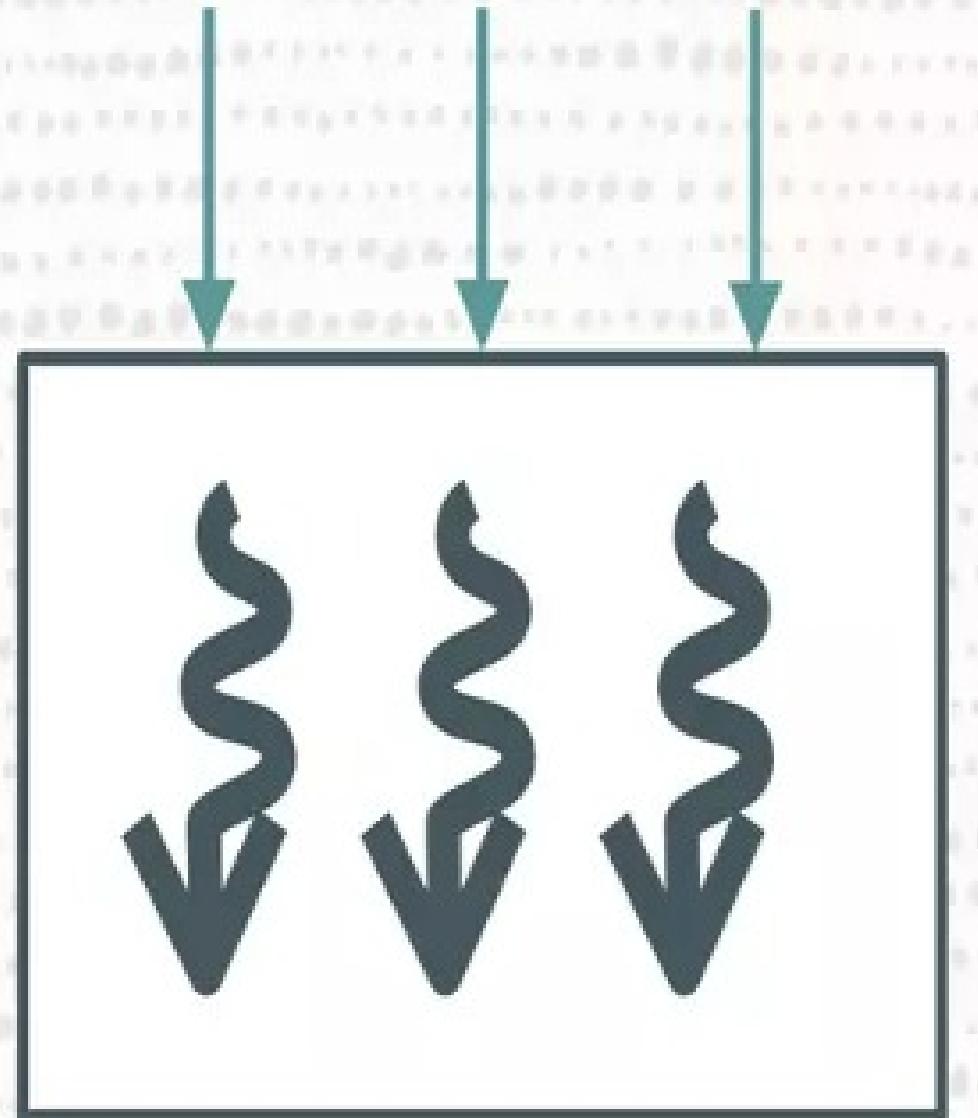
(> 150 cache misses) \Rightarrow 10,000 threads per core

$$W = 100ms$$

Java Is Made of Threads

- Exceptions
- Thread Locals
- Debugger
- Profiler (JFR)

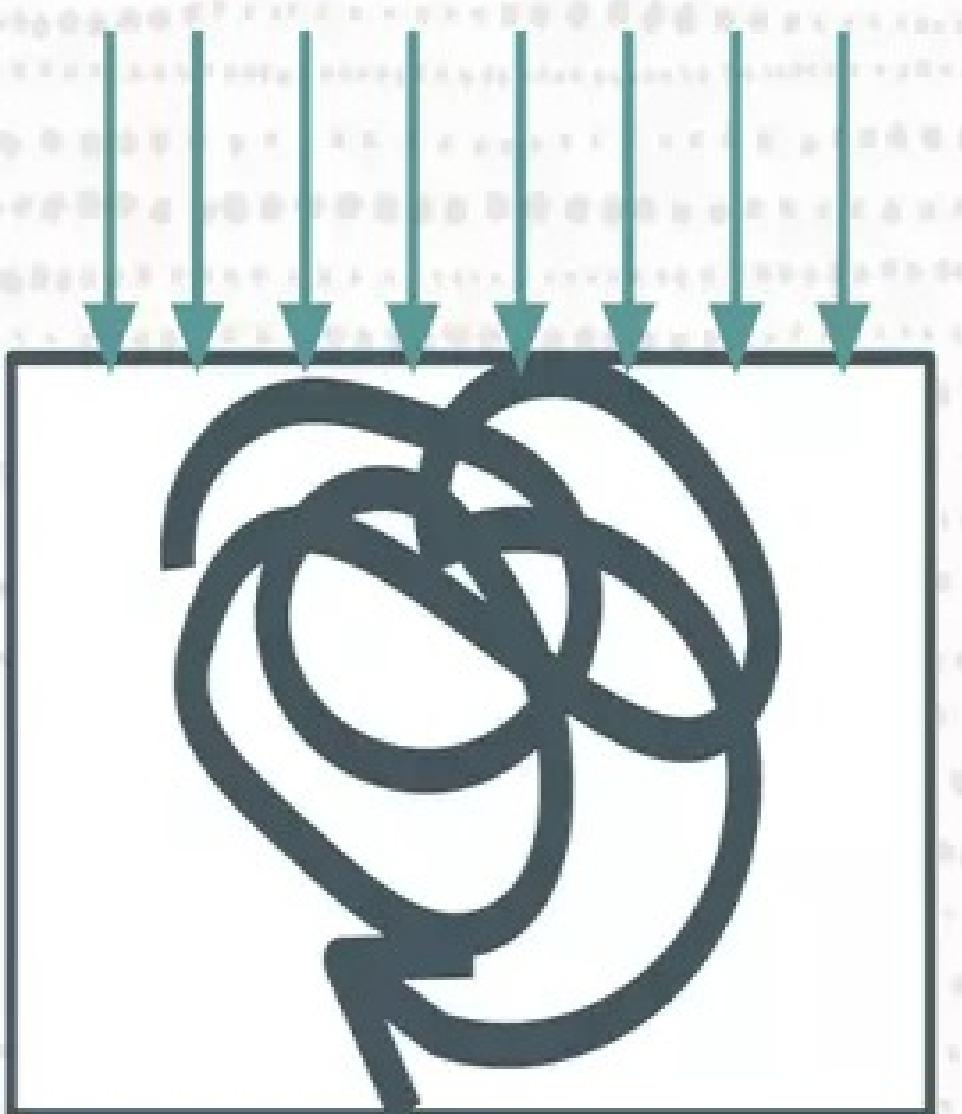




simple
less scalable

SYNC

Programmer
Hardware



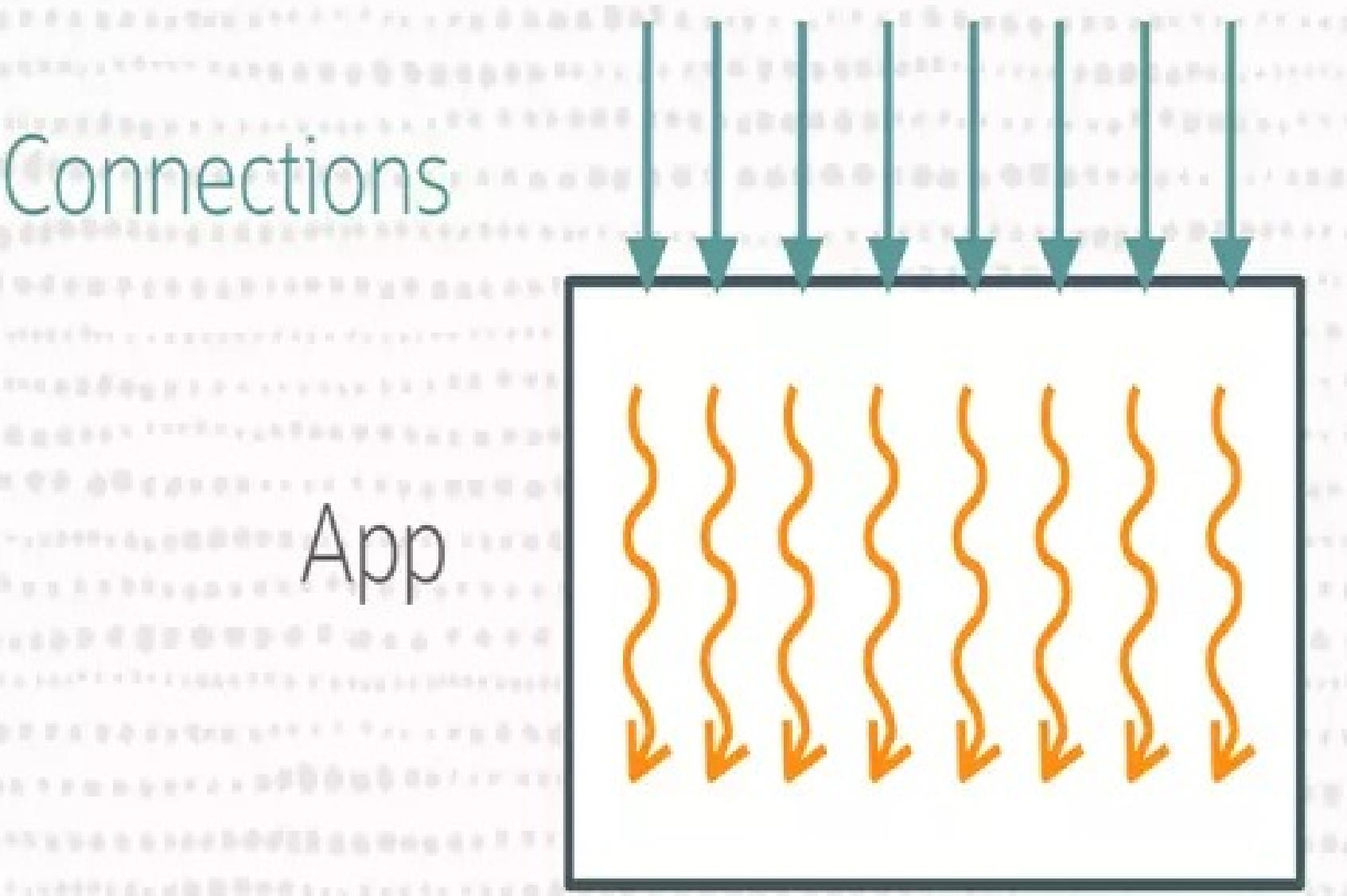
scalable,
complex,
non-interoperable,
hard to debug/profile

ASYNC

Programmer
Hardware



Virtual Threads



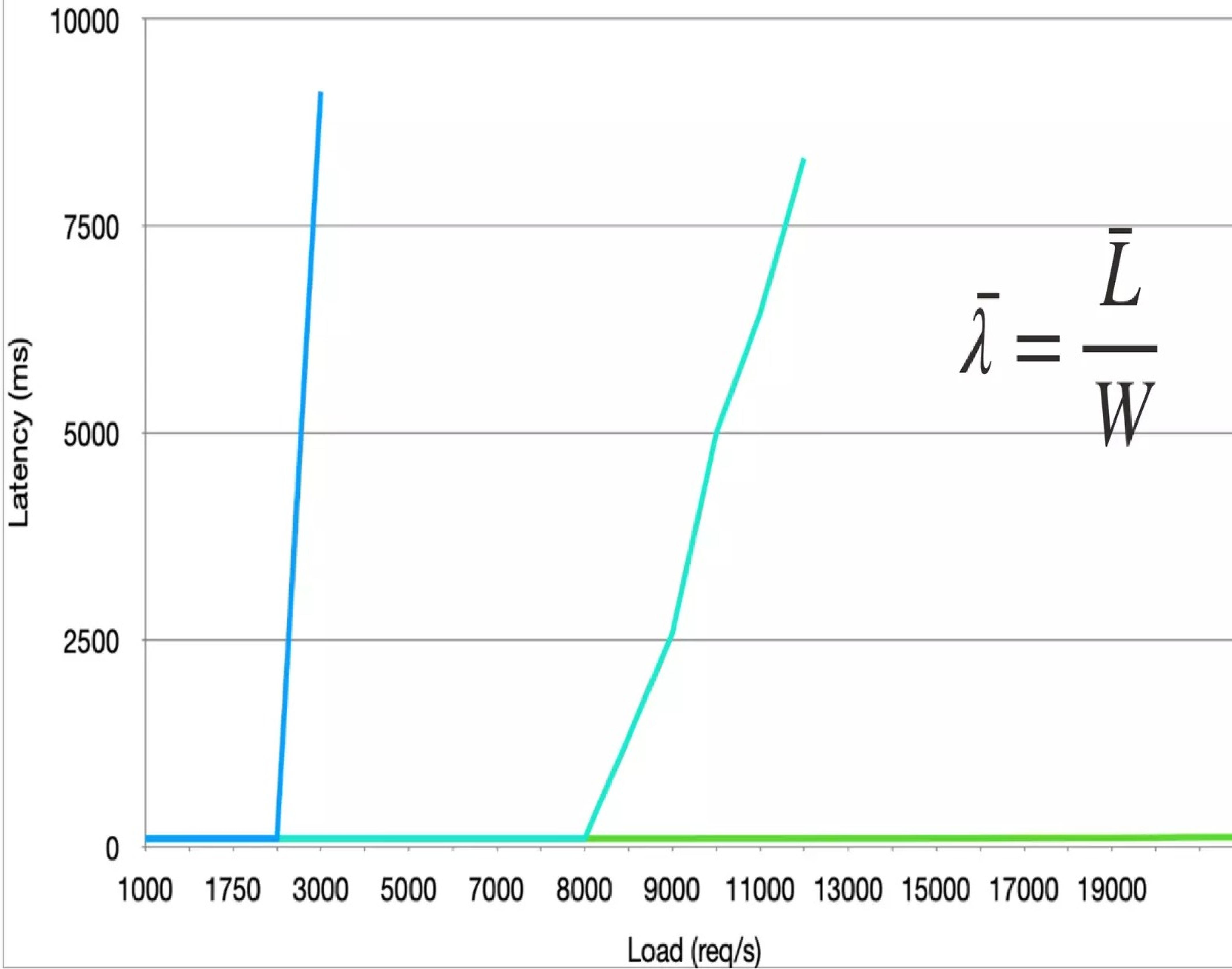
Programmer

Hardware

— 200 Platform Threads

— 800 Platform Threads

— Virtual Threads



$$\bar{\lambda} = \frac{\bar{L}}{W}$$



The Impact of Context Switching

- Context-switching affects throughput by means of duration, not capacity $(\frac{\lambda_1}{\lambda_2} = \frac{W_2}{W_1})$

$$W = n(\mu + t) \quad \text{where} \quad \begin{aligned} n & \text{ avg. #operations} \\ \mu & \text{ avg. wait duration} \\ t & \text{ avg. context-switch duration} \end{aligned}$$

$$\frac{n(\mu + t)}{n\mu} = 1 + \frac{t}{\mu} \Rightarrow \begin{aligned} \mu &= 20\mu s \quad (\text{quite fast}) \\ t &= 1\mu s \quad (\text{quite slow}) \end{aligned} \Rightarrow 5\% \text{ impact}$$

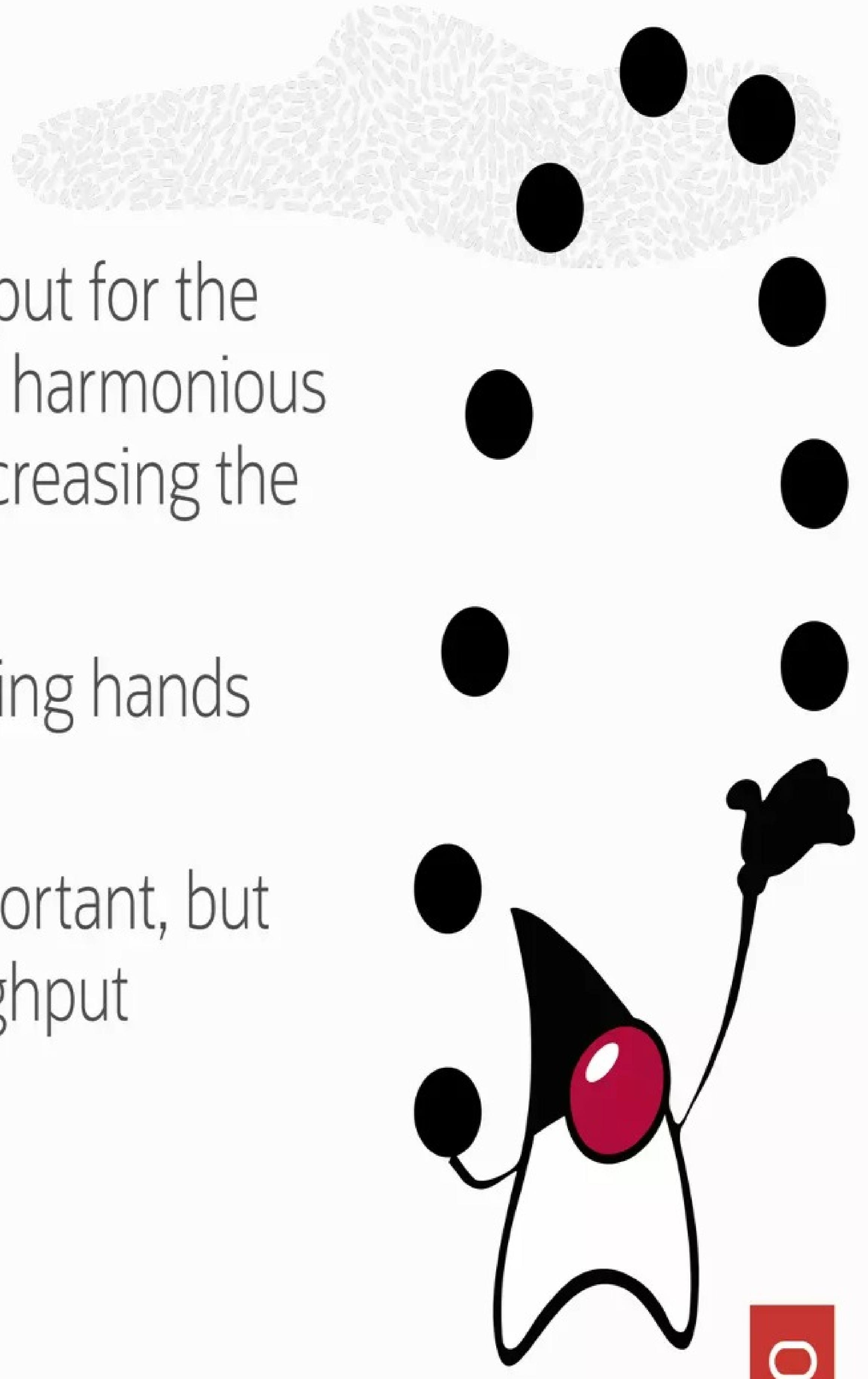
- Virtual threads have a faster context-switch than OS threads
- Structured concurrency allows waiting for a set of operations with one context-switch

Not cooperative, but no time-sharing yet

- Non-cooperative scheduling is more composable
- ... but people overestimate the importance of time-sharing in servers
- Effective time-sharing effective when #threads is $\#cores \times 10^5$ requires study

Summary

- Virtual threads allow higher throughput for the thread-per-request style — the style harmonious with the platform — by drastically increasing the request capacity of the server.
- We can juggle more balls not by adding hands but by enlarging the arch.
- Context-switching cost could be important, but aren't the main reason for the throughput increase.



- JEP 425: Virtual Threads (Preview)
- The Design of User Mode Threads in Java [video] (why not async/await)



Ron Pressler

ron.pressler@oracle.com

@pressron