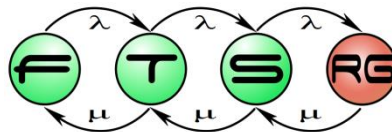


Performance Modelling

Budapest University of Technology and Economics
Fault Tolerant Systems Research Group



A Well Known Example...

Budapesti Műszaki és Gazdaságtudományi Egyetem



Hallgatói BME_W2K8H_02(168)

Nyelv:    

Azonosító:

Jelszó:

Bejelentkezés >



Build: 430 (2014.11.11.) P20150304

Támogatott böngészők:

Microsoft Internet Explorer 9.0+ ; Mozilla Firefox ; Google Chrome

Friss hírek

[KTH hallgatói hírlevél 2014/15/1 félévzárásról 1. – vizsgajelentkezés, vizsgaidőszak](#)

(2014.11.27.
8:56:21)

Kedves Hallgató!

2014. december 1-jén 18 órakor indul a vizsgajelentkezés. Milyen egyéb határidőkre kell figyelnie a félév végén? Melyek az ebben az időszakban aktuális Neptun kérvények? Mit

Letölthető dokumentumok

 Tantárgyfelvétel: eddig a legsimábban.pdf
(2015.02.07. 14:28:14)

 Tantárgyfelvétel: nem és nem.pdf
(2014.08.28. 21:19:50)

 Vizsgajelentkezés: szokványosan.pdf
(2014.05.06. 21:14:04)

A Well Known Example...

Kiszolgálóhiba történt az alkalmazásban: „/hallgato”.

Futásidejű hiba

Leírás: Alkalmazáshiba történt a kiszolgálón. Az alkalmazás jelenlegi egyéni hibakezelési beállításai (biztonsági okok miatt) nem engedik meg az adatainak távoli megjelenítését. A böngészőben azonban megtekinthetők.

Részletek: Ha távoli gépeken is meg szeretné jeleníteni a hibaüzenet részletes adatait, hozzon létre a konfigurációs fájlban elhelyezett „web.config” címke „mode” attribútumát állítsa „Off” értékre.

<!-- Web.Config konfigurációs fájl -->

```
<configuration>
  <system.web>
    <customErrors mode="Off"/>
  </system.web>
</configuration>
```

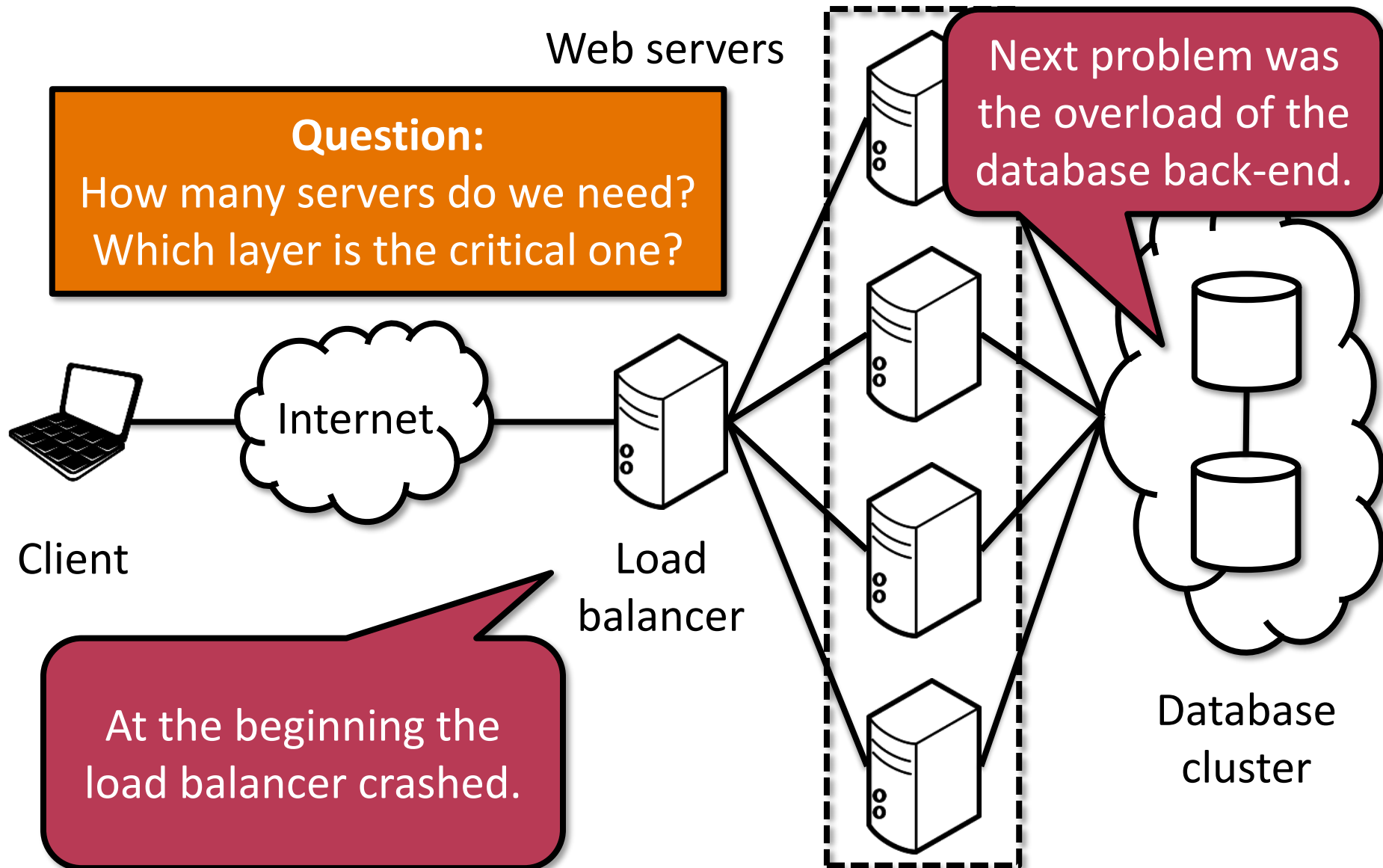
Megjegyzések

```
<!-- Web.Config konfigurációs fájl -->
<configuration>
  <system.web>
    <customErrors mode="Off"/>
  </system.web>
</configuration>
```

**Motivation:
Prepare for performance
issues in design time!**



Schematic Architecture of Neptun System

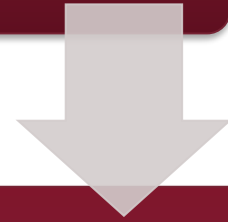


Performance Modelling

- Review: non-functional requirements
 - Performance, throughput, dependability, etc.
 - How can such requirements be verified?
(without building the actual system)
- **Performance modelling:**
 - Extending the discussed models with timing, resources, capacity constraints, ...
 - Aim of this activity:
 - Evaluating performance of the system during design phase
 - Identifying bottlenecks
 - Scaling, capacity planning, dimensioning

Content

Foundations



Load Diagram



Resource Modelling

Foundations

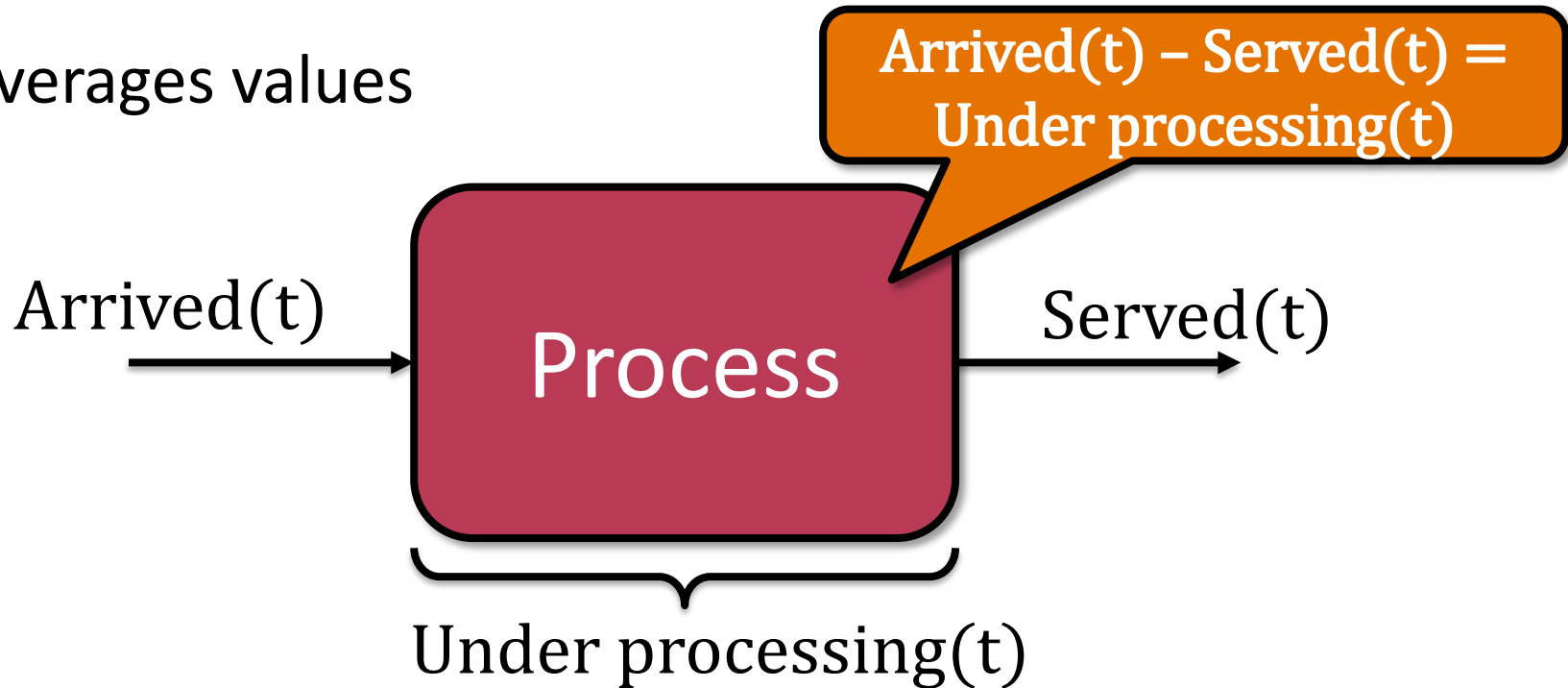
Load Diagram

Resource Modelling

FOUNDATIONS

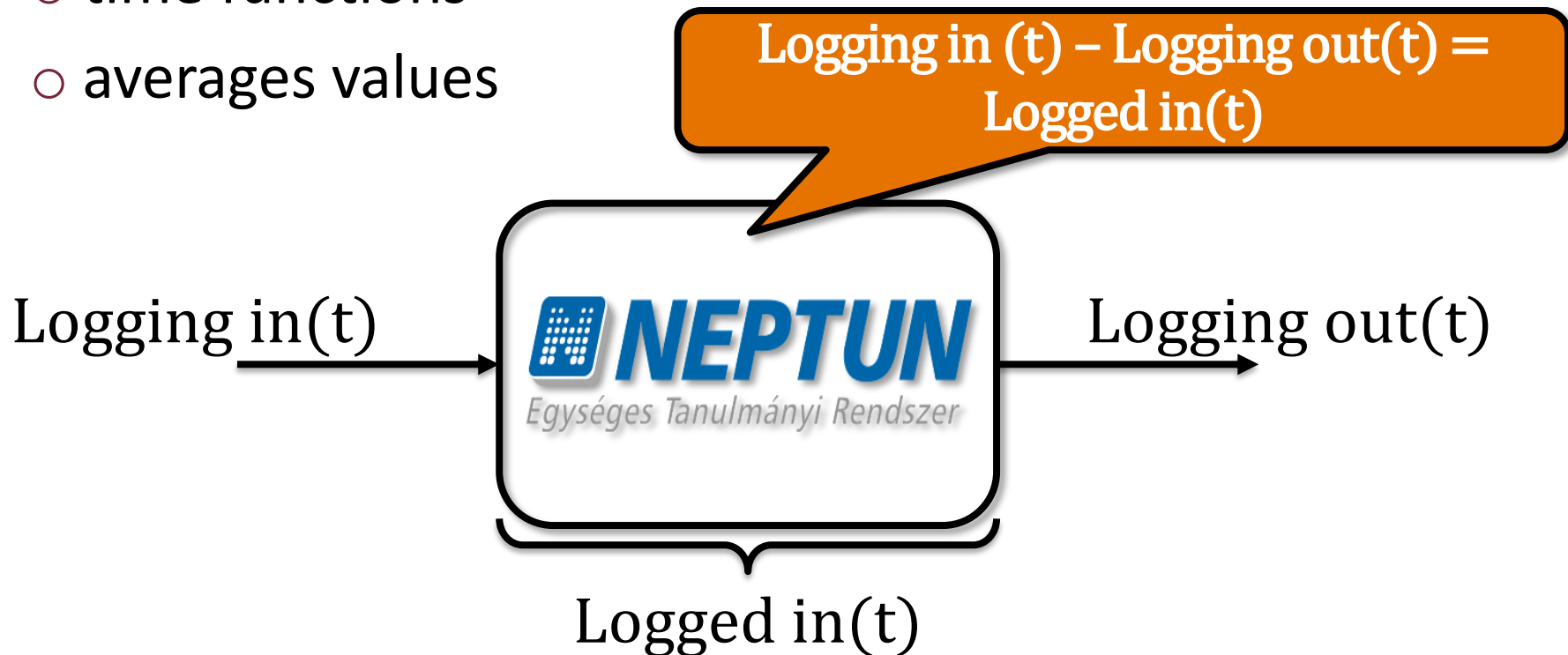
Basic Model

- Execution of a process in case of multiple requests
 - Subject of the analysis: time dependent behaviour
- Description of the behaviour:
 - time functions
 - averages values



Basic Model

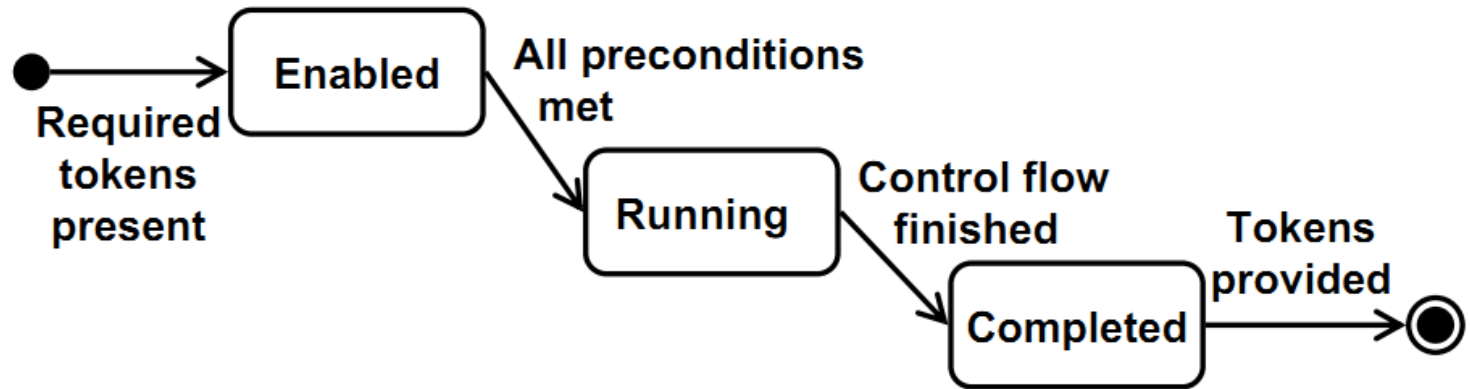
- Execution of a process in case of multiple requests
 - Subject of the analysis: time dependent behaviour
- Description of the behaviour:
 - time functions
 - averages values



Review: Execution States of Activities

- State of the process/activity execution:

State Machine



- Arrived(t): number of tokens in „*Enabled*” state
- Under processing(t): # of tokens in „*Running*” state
- Served(t): number of tokens in „*Completed*” state

Definition: Arrival Rate, Throughput

Arrival rate: number of **arriving** requests during a specific unit of time.

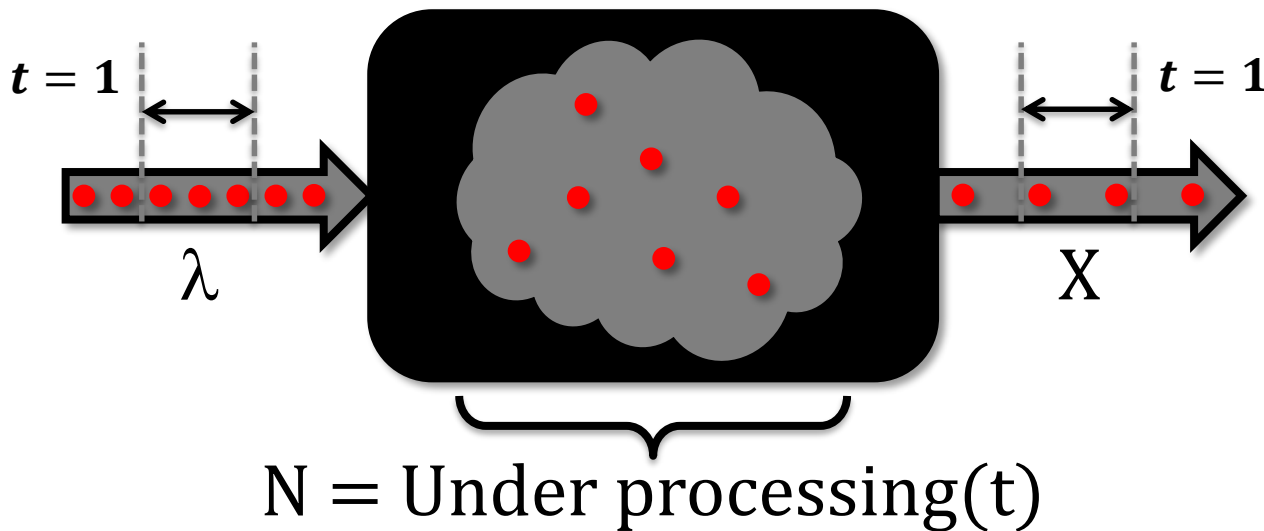
$$\lambda = \frac{Arrived(t)}{t}$$

$$[\lambda] = \frac{1}{s}$$

Throughput (X_{PUT}): number of requests **processed** during a specific unit of time.

$$X = \frac{Served(t)}{t}$$

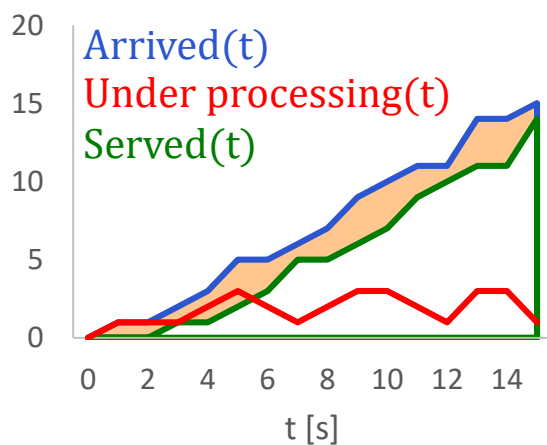
$$[X] = \frac{1}{s}$$



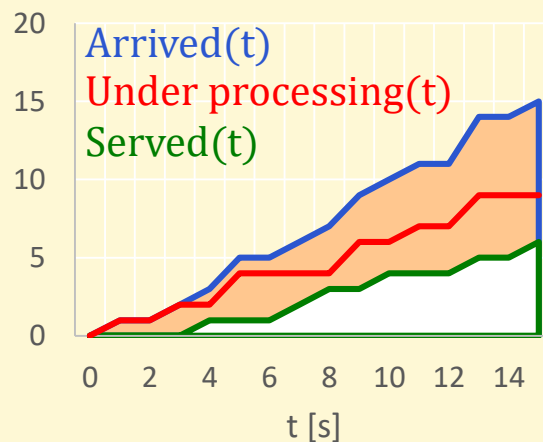
Definition: Stable State

- **Stable state:** *Under processing(t)* is approximately constant
 - Average values can be applied in such state!
 - A system is in balance, if: $\lambda = X$

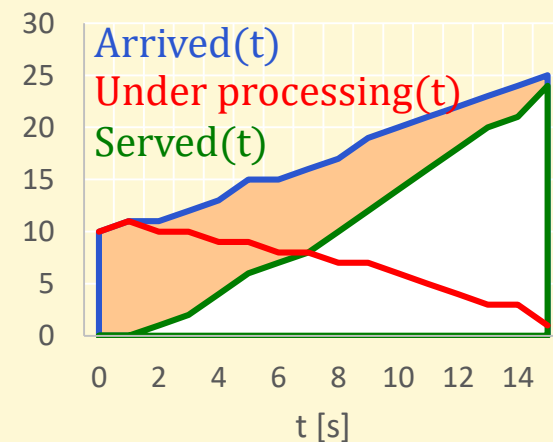
$$\lambda = X$$



$$\lambda > X$$



$$\lambda < X$$



Definition: Stable State

- **Stable state:** *Under processing(t)* is approximately constant
 - Average values can be applied in such state!
 - A system is in balance, if:

In stable state:
Same number of logins
and logouts per minute



$$N = \text{Logged in}(t)$$

Limited Capacity – DoS

- N is not infinite in real life
- So, what is then?

Denial of Service Attack

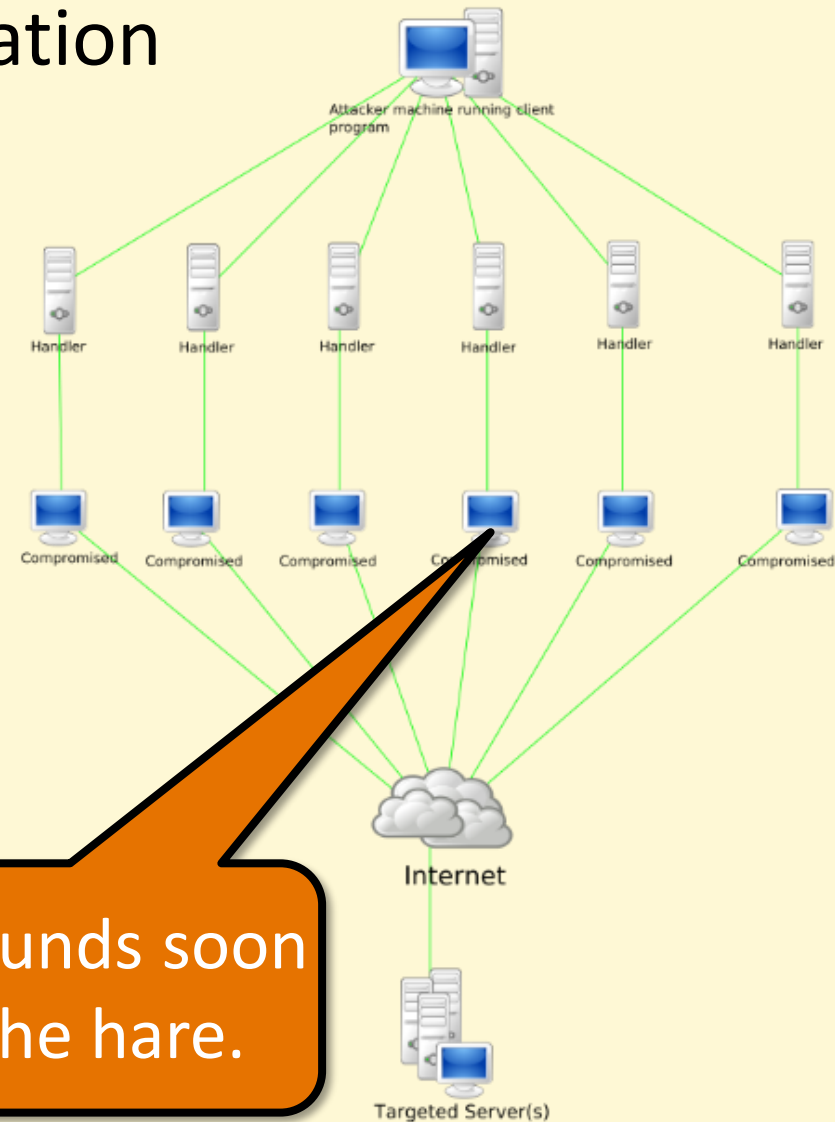
Thursday, August 6, 2009 | By Biz Stone (@biz) 08/06/2009 - 15:00

Tweet

On this otherwise happy Thursday morning, Twitter is the target of a [denial of service attack](#). Attacks such as this are malicious efforts orchestrated to disrupt and make unavailable services such as online banks, credit card payment gateways, and in this case, Twitter for intended customers or users. We are defending against this attack now and will continue to update our [status blog](#) as we continue to defend and later investigate.

(Distributed) Denial of Service – (D)DoS

- Mass scaled request generation
→ overload of a system
- An overloaded system is fragile (can be cracked)
- Complete services can be knocked out
- Favourite method of the Anonymous group



Foundations

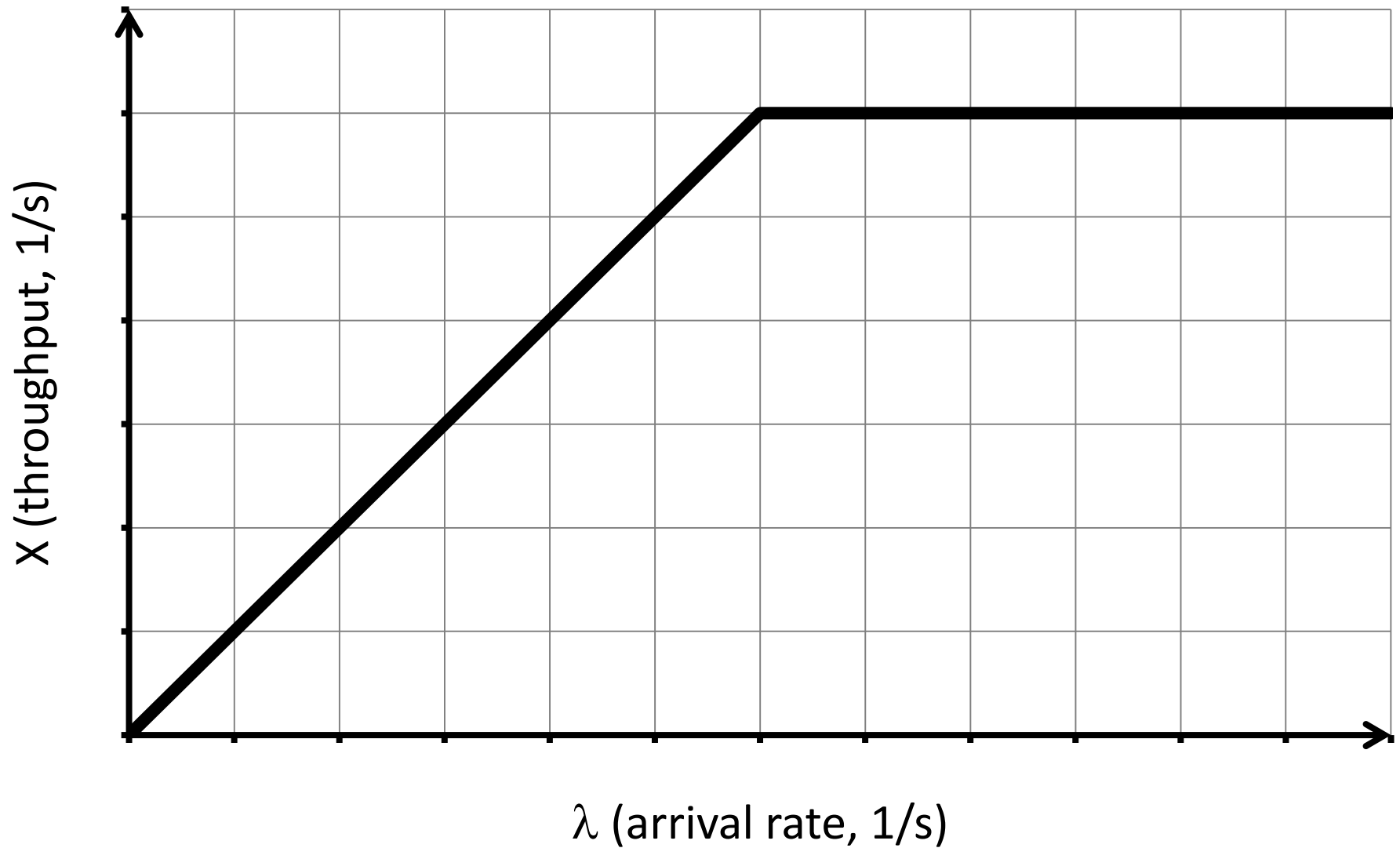
Load Diagram

Resource Modelling

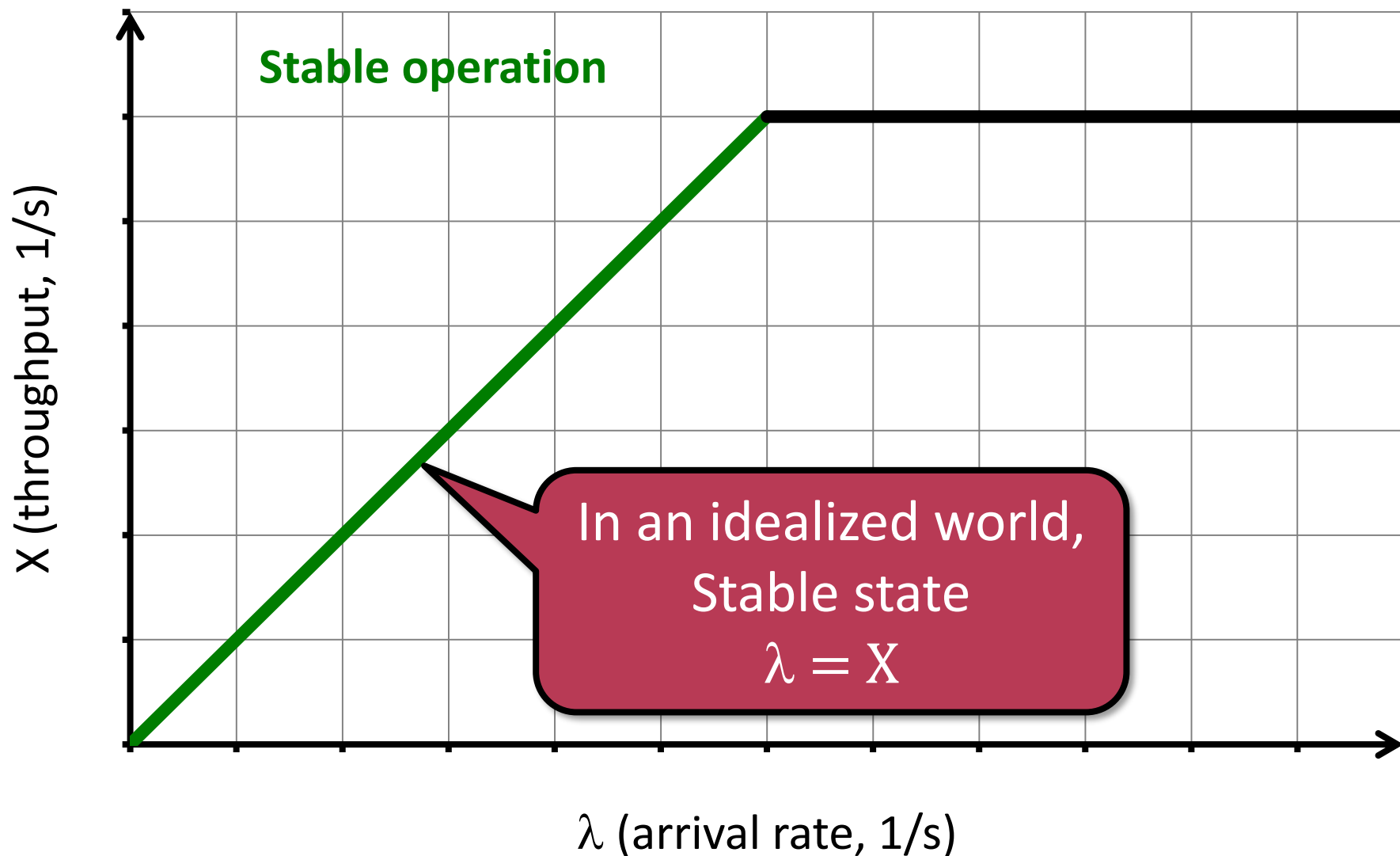
LOAD DIAGRAM

- Relation between arrival rate and throughput
- Maximum throughput

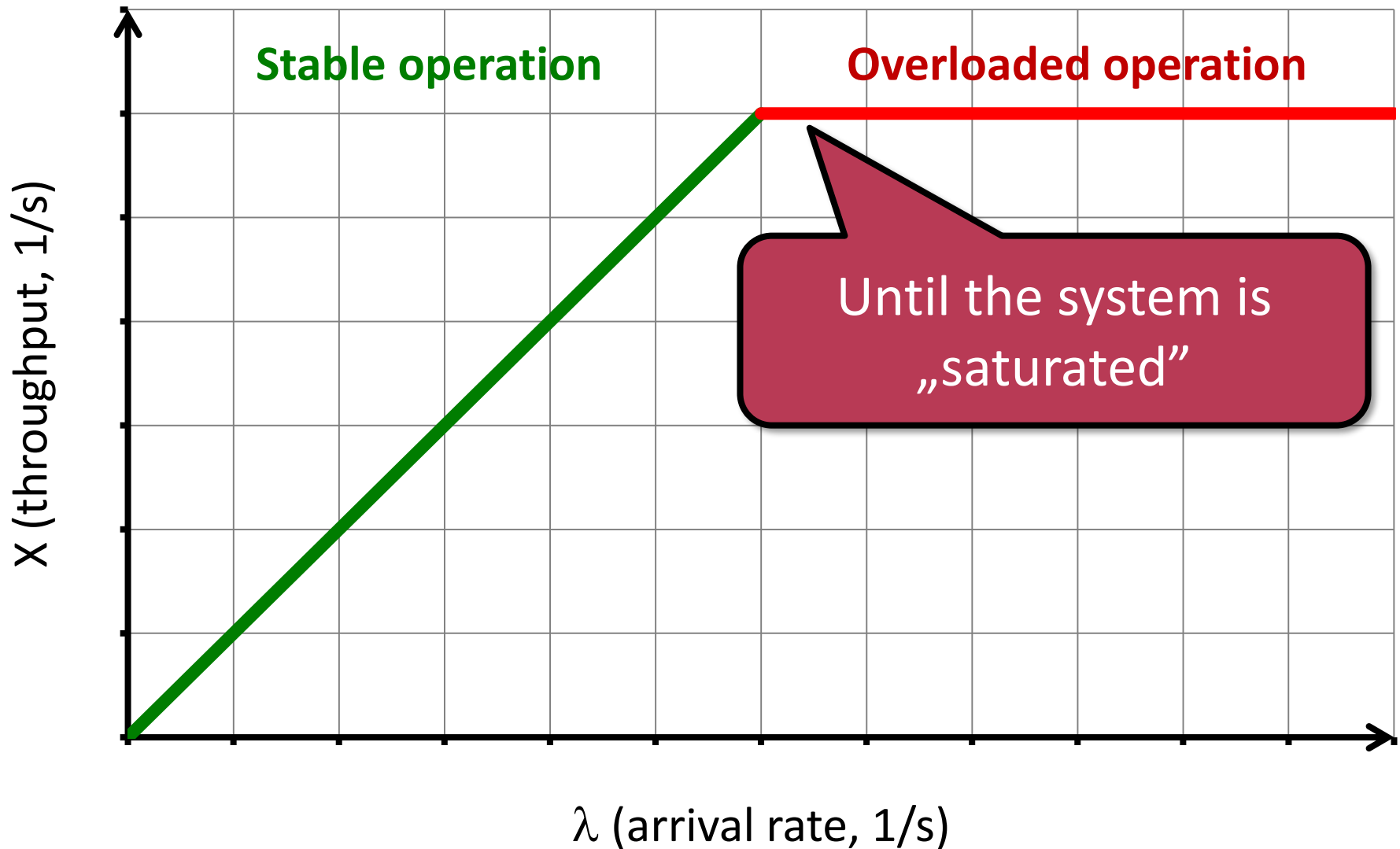
Load Diagram



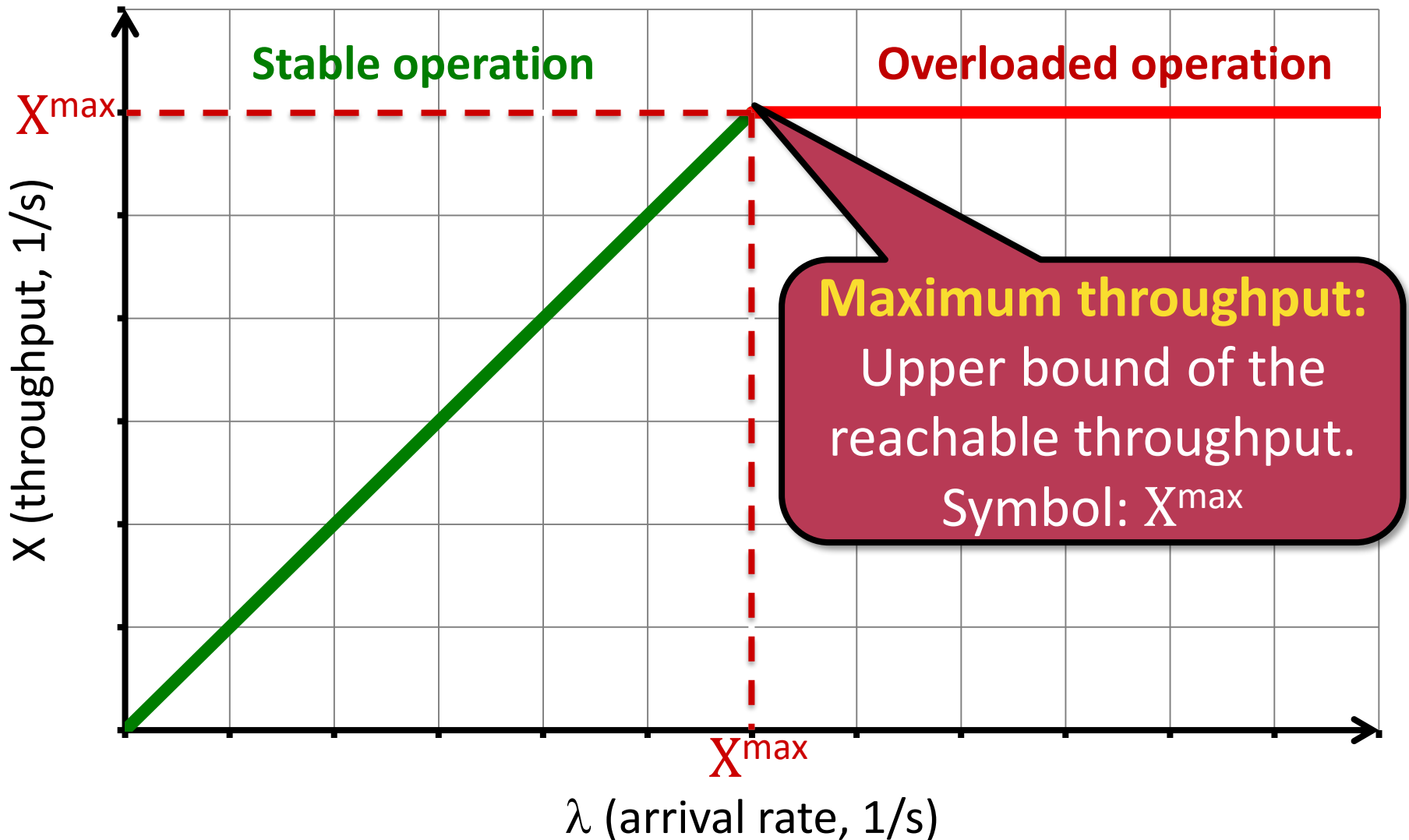
Load Diagram



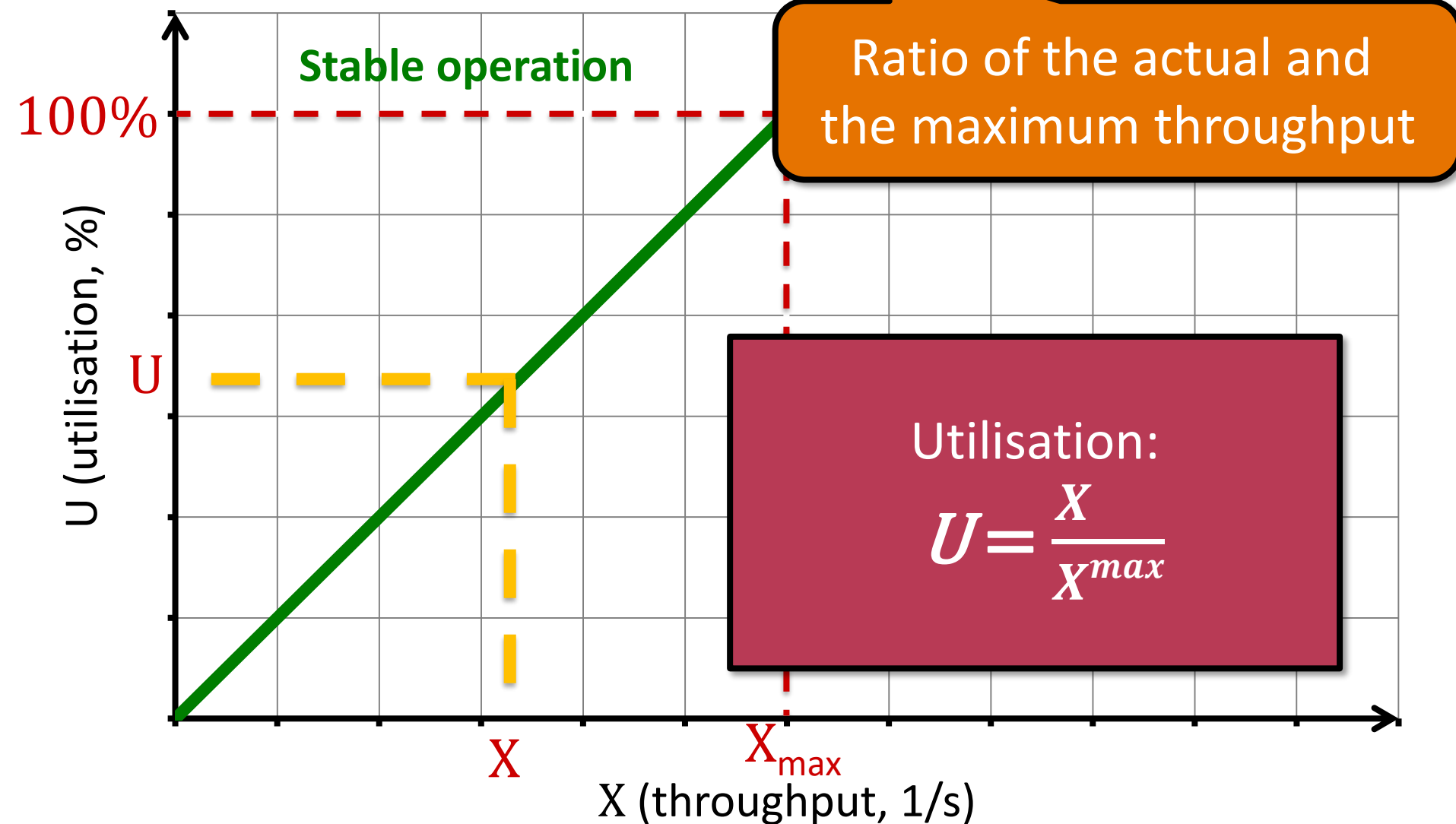
Load Diagram



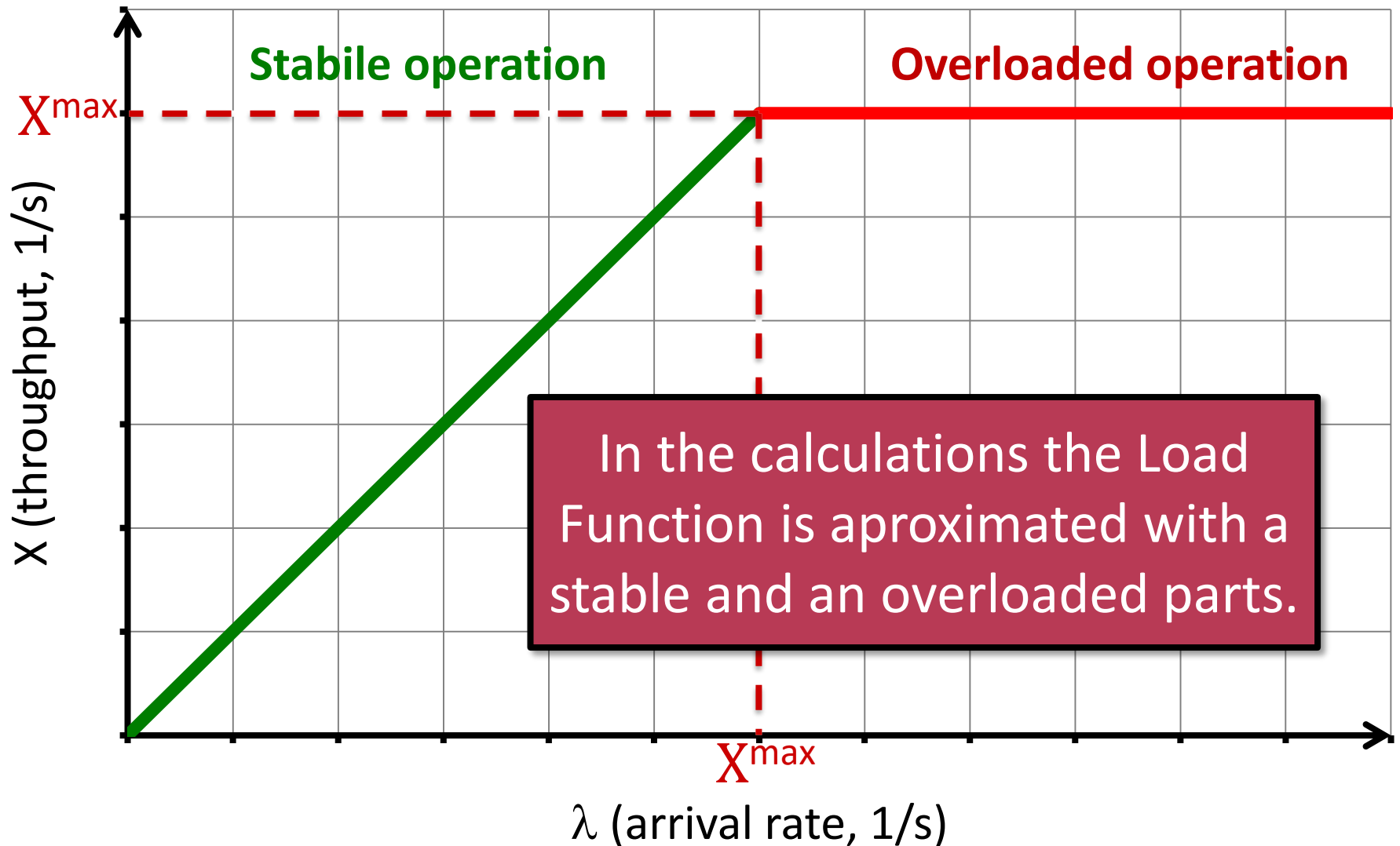
Maximum Throughput



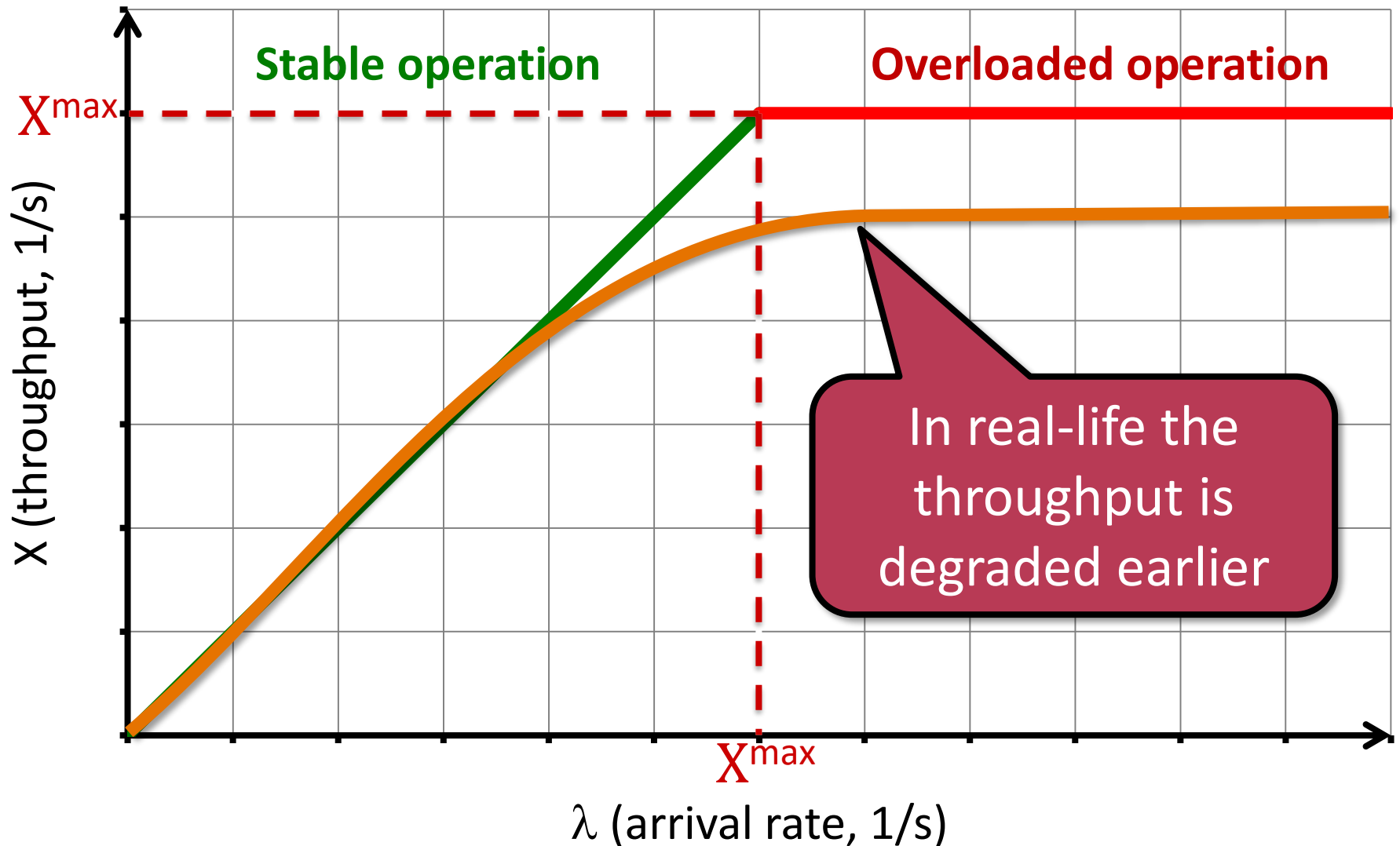
Utilisation



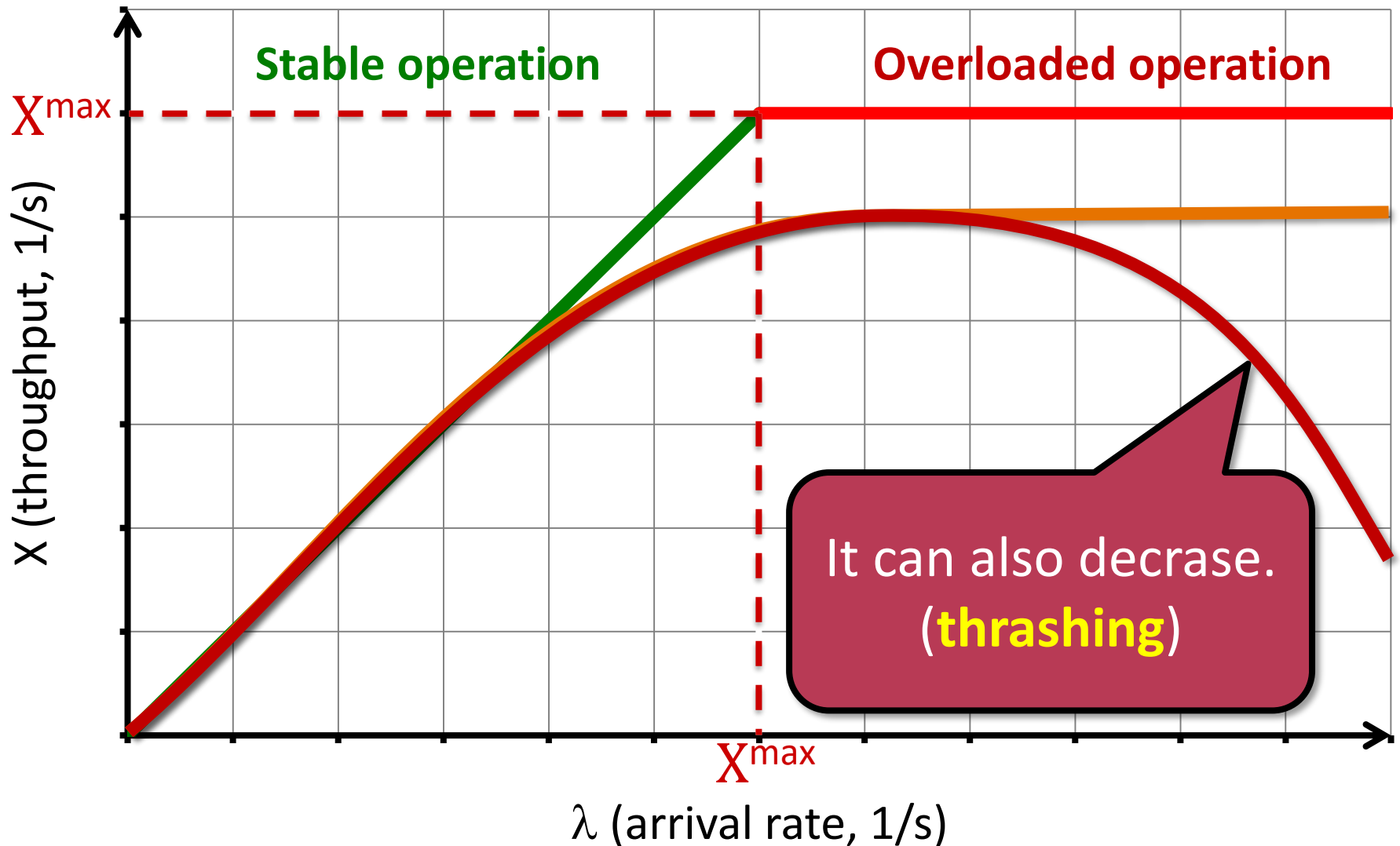
Approximative Load Function



Real Load Diagram



Real Load Diagram



Definitions

Maximum throughput: maximal reachable throughput

- Symbol: X^{\max}

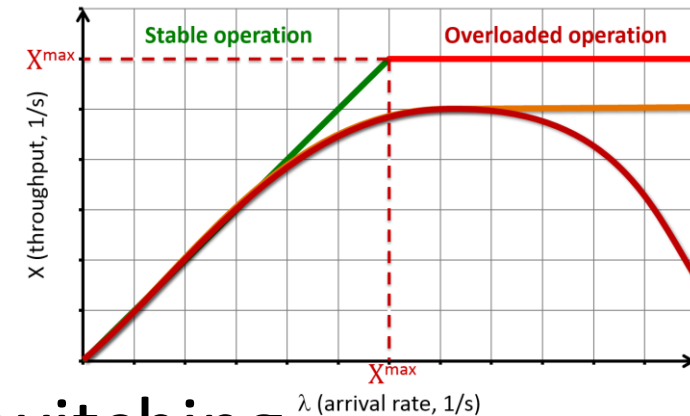
Utilisation: ratio of the actual and the maximum throughput

- Symbol: $U = \frac{X}{X^{\max}}$

Thrashing: throughput decreases during overloaded operation

Effects Ignored by the Model

- Cost of task switching
 - Cleanup after processing
 - Preparation for the next process
- Computation cost of resource switching
- Multiple overload
 - Multiple resources (e.g. servers) can be overloaded
 - E.g. if there is an accident on the M7, there is also a traffic jam on the country road 7



Scared Off Requests

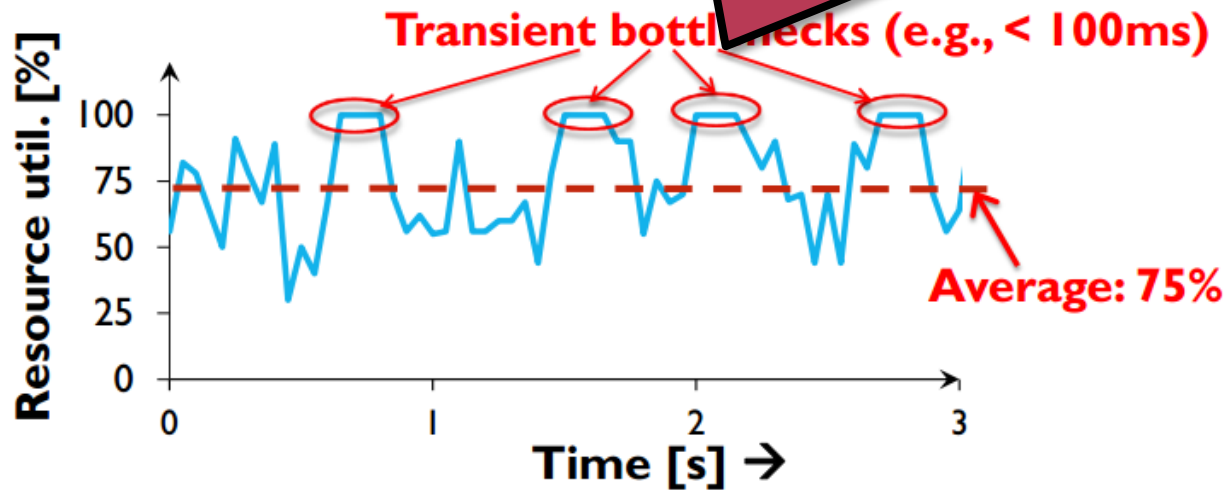
- Number of requests are only independent of waiting time in case of fanatic customers



Effects of Load Fluctuation

- Average values vs. **real load**

In case of an overloaded system the waiting time can be **2-3 magnituded higher**



Effects of Load Fluctuation

- Average values vs. **real load**

[%]

**It is not worth to plan a
utilisation higher than 40-60%!**

Average: 75%

Time [s] →

Neptun Course Sign-ups



Max. # of
concurrent users

Max. # of concurrent users
at optimal operation

Az elmúlt időkből előfordultak a következők:

Dátum	Jelleg	Max. felhasználó	Op. felhasználó
2010.11.29 18:00	vizsga	7303	4623
2010.12.22 06:00	tárgy (EO)	831	831
2011.01.10 18:00	tárgy	12062	4837
2011.01.12 18:00	tárgy (EP)	1765	1765
2011.01.31 16:00	tárgy	1519	1519
2011.05.02 18:00	vizsga	2761	2761
2011.06.07 18:00	tárgy	6095	6095
2011.11.28 18:00	vizsga	4897	4897
2012.01.16 18:00	tárgy	8120	5328
2012.01.30 16:00	tárgy	1703	1703
2012.05.02 18:00	vizsga	2603	2603

Neptun Course Sign-ups



≈ Arrival rate (λ)

≈ Max. Throughput (X^{\max})

Az elmúlt időkben előfordult, hogy a Neptun rendszer túlterhelt volt.

Dátum	Jelleg	Max. felhasználó	Op. felhasználó
2010.11.29 18:00	vizsga	7303	4623
2010.12.22 06:00	tárgy (EO)	831	831
2011.01.10 18:00	tárgy	12062	4837
2011.01.12 18:00	tárgy (EP)	1765	1765
2011.01.31 16:00	tárgy	1519	1519
2011.05.02 18:00	vizsga	2761	2761
2011.06.07 18:00	tárgy	6095	6095
2011.11.28 18:00	vizsga	4897	4897
2012.01.16 18:00	tárgy	8120	5328
2012.01.30 16:00	tárgy	1703	1703
2012.05.02 18:00	vizsga	2603	2603

When was Neptun overloaded?

Neptun Course Sign-ups

\approx Arrival rate (λ)

$\approx M$ (X^{\max})

Az elmúlt időkben elő...

Dátum	Jelleg	
2010.11.29 18:00	vizsga	
2010.12.22 06:00	tárgy /	
2011.01.10 18:00		
2011.01.12 18:00		
2011.01.31		
2011.		2761
2011.		6095
2011.11		4897
2012.01.		5328
2012.01.30		1703
2012.05.02		2603

The performance of an overloaded system can be degraded drastically, it can even crash!

When was Neptun overloaded?

Neptun Course Sign-ups



≈ Arrival rate (λ)

≈ Max. Throughput (X^{\max})

Az elmúlt időkben el...

Dátum	Jelleg	Max. felhasználó	Op. felhasználó
2010.11.29 18:00	vizsga	7303	4623
2010.12.22 06:00	tárgy (EO)	831	831
2011.01.10 18:00	tárgy	12062	4837
2011.01.12 18:00	tárgy (EP)	1765	1765
2011.01.31 16:00	tárgy	1519	1519
2011.05.02 18:00	vizsga	2761	2761
2011.06.07 18:00	tárgy	6095	6095
2011.11.28 18:00	vizsga	4897	4897
2012.01.16 18:00	tárgy	8120	5328
2012.01.30 16:00	tárgy	17	
2012.05.02 18:00	vizsga	26	

Correctly configured
servers, appropriate
capacity planning!

Neptun Course Sign-ups

Protection against (D)DoS attacks is based on the same principles

Limiting the number of users can keep the system in stable state

Correctly configured servers, appropriate capacity planning!

NEPTUN	
date (λ)	out (X^{\max})
2010.11.23 18:00	
2010.12.22 06:00	tárgy
2011.01.10 18:00	
2011.01.12 18:00	
2011.01.31 18:00	
2011.05.02 18:00	2761
2011.06.02 18:00	6095
2011.11.28 18:00	4897
2012.01.16 18:00	5328
2012.01.30 16:00	17
2012.05.02 18:00	26

Foundations

Load Diagram

Resource Modelling

RESOURCE MODELLING

Why is there a *maximum* of the throughput?

Example: Mid-Term Exam Grading

- Grading a single exam takes **15 mins**
- How many exams can be graded by a single lecturer in an hour?

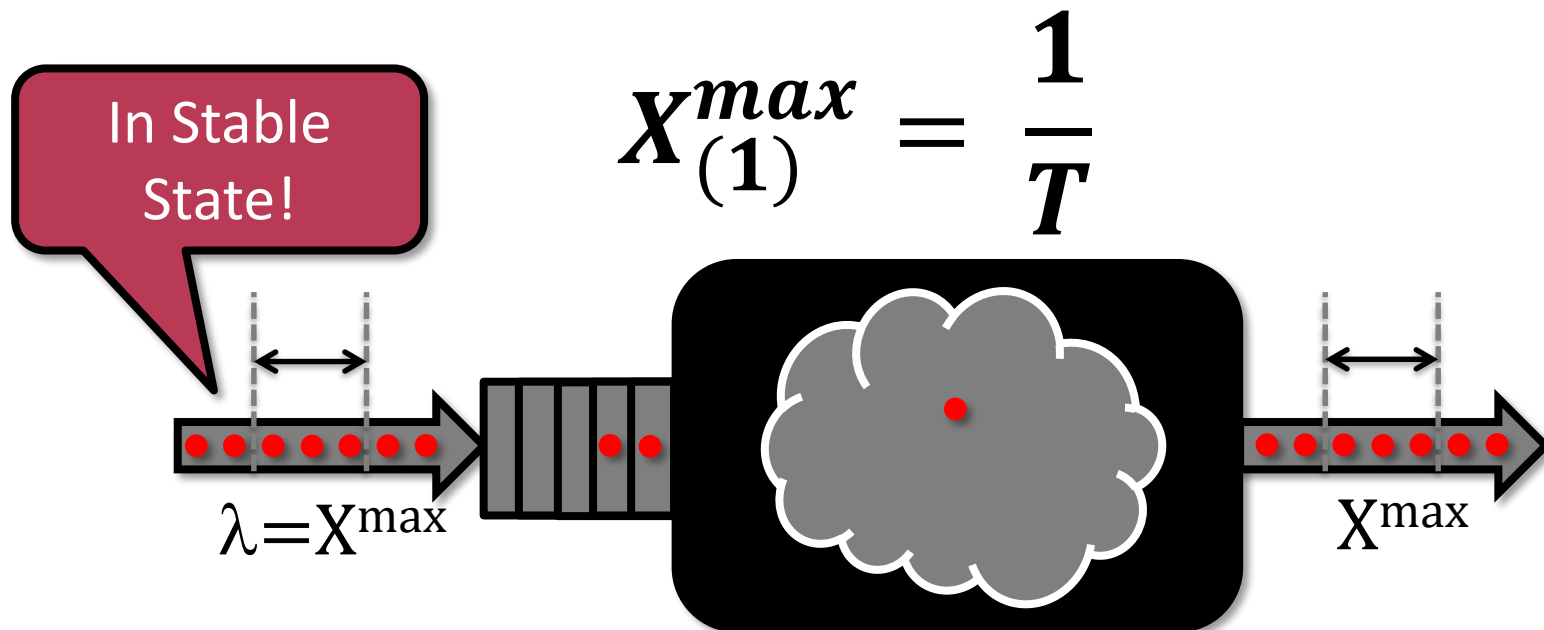
$$X_{(1)}^{max} = \frac{1 \text{ exam}}{15 \text{ min.}} = 4 \frac{\text{exams}}{h}$$

- How many exams can be graded by **eight** lecturers?

$$X_{(8)}^{max} = 8 \times X_{(1)}^{max} = 8 \times \frac{1 \text{ exam}}{15 \text{ min.}} = 32 \frac{\text{exams}}{h}$$

Max. Throughput of a Single Server System

- If **only a single** request can be processed
 - e.g. served by a single server, or by working with the same (shared) variables
 - the remaining requests are queued
- Then with **the average execution time T**:



Max. Throughput of a Single Server System

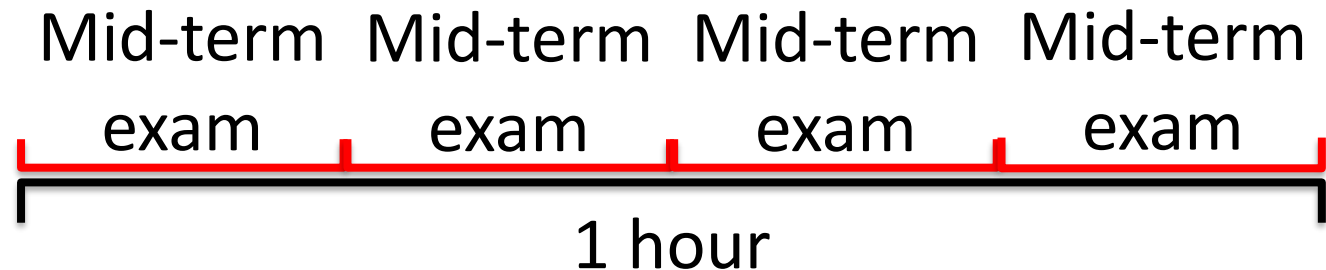
- If **only a single** request can be processed
 - e.g. served by a single server, or by working with the same (shared) variables
 - the remaining requests are queued
- Then with **the average execution time T**:

$$X_{(1)}^{max} = \frac{1}{T}$$

„How many requests can be processed in a unit of time?”

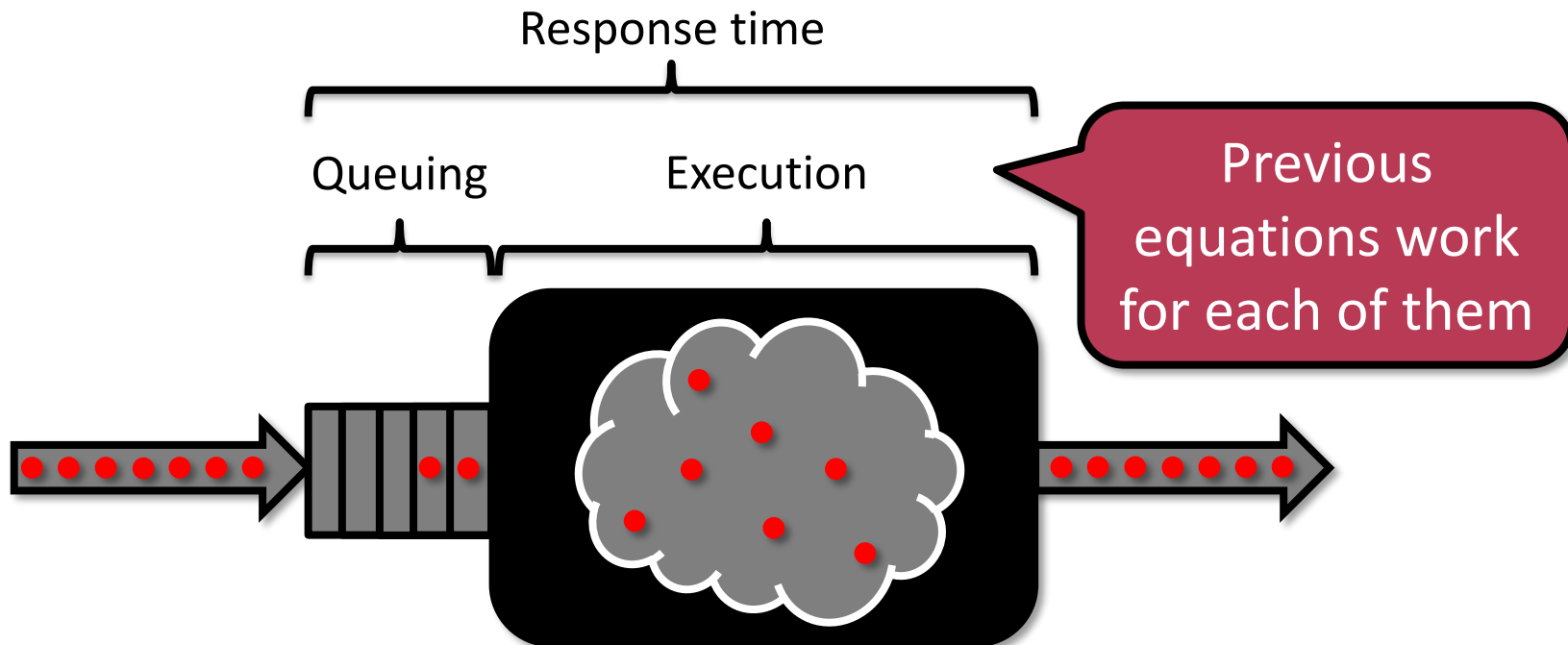
$$T_{exam} = 15 \text{ min.}$$

$$X_{max} = 4 \frac{1}{h}$$



Different Times to Measure

- **Queuing time:** request waiting for a resource
- **Execution time:** request under processing
- **Response time:** Queuing + Execution times



Utilisation of a Single Server System

- Utilisation: „Ratio of the maximum and the actual throughputs”

$$X_{(1)}^{max} = \frac{1}{T} \Rightarrow X_{(1)}^{max} \times T = \textcircled{1}$$

Utilisation of a Single Server System

- Utilisation: „Ratio of the maximum and the actual throughputs”

$$X_{(1)}^{max} = \frac{1}{T} \Rightarrow X_{(1)}^{max} \times T = 1 = U^{max}$$

- The equation of the utilisation:

$$U = X \times T$$

- Intuition:

„If X requests arrive in a unit of time, each of which takes T time to process, what percentage of time is spent busy?”

Utilisation of a Single Server System

- Utilisation: „Ratio of the maximum and the actual throughputs”

$$X_{(1)}^{max} = \frac{1}{T} \Rightarrow X_{(1)}^{max} \times T = 1 = U^{max}$$

- The equation of the utilisation:

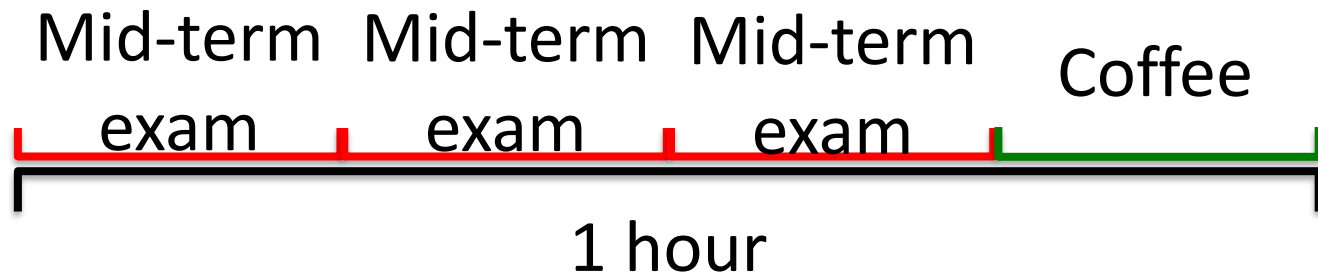
$$U = X \times T$$

- Intuition:

$$X = 3 \frac{1}{h}$$

$$T_{exam} = 15 \text{ min}$$

$$U = 75\%$$



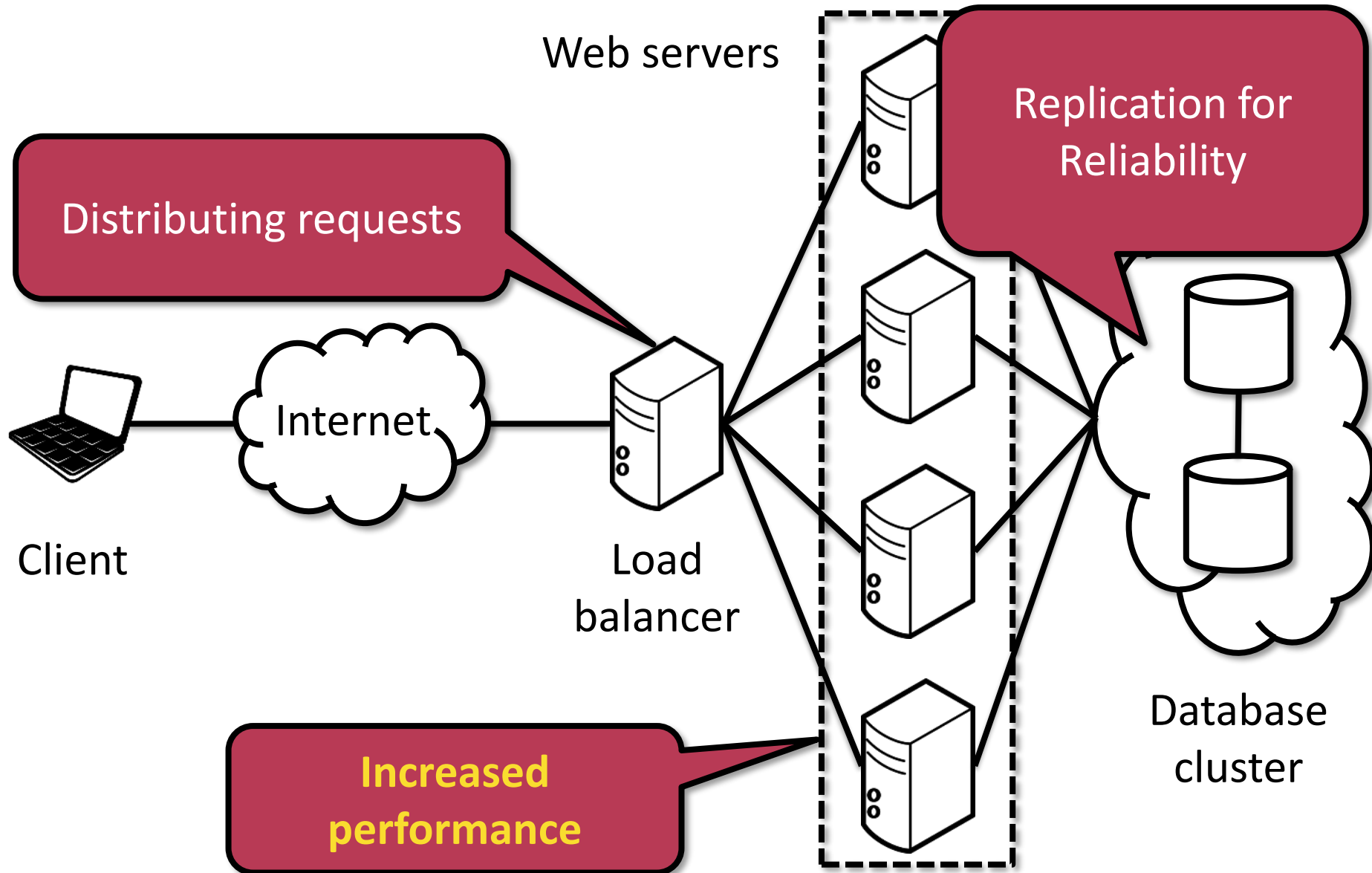
Outlook: Scalability

- **Vertical scaling** (Scale-up):
 - The **power** of the processing units is **increased**
 - E.g. stronger CPU, more RAM
 - Simple and great 😊
 - Technological constraints ☹️
- **Horizontal scaling** (Scale-out):
 - The **number** of the processing units is **increased**
 - E.g. Multiple (core) CPU, multiple servers
 - Theoretically no limitations 😊
 - Extra complexity ☹️

Slace-out in Everyday Life



Scale-out of Neptun



Max. Throughput of a Multiple Server System

- If **maximum K** requests are allowed to be served simultaneously
 - e.g. K clustered servers can process them
 - the remaining requests are queued
- Then with **the average execution time T**:

$$X_{(K)}^{max} = K \times X_{(1)}^{max} = \frac{K}{T}$$

Max. Throughput of a Multiple Server System

- If **maximum K** requests are allowed to be served simultaneously
 - e.g. K clustered servers can process them
 - the remaining requests are queued
- Then with **the average**

With more resources and parallelisation the system is **scalable**.

$$X_{(K)}^{max} = K \times X_{(1)}^{max} = \frac{K}{T}$$

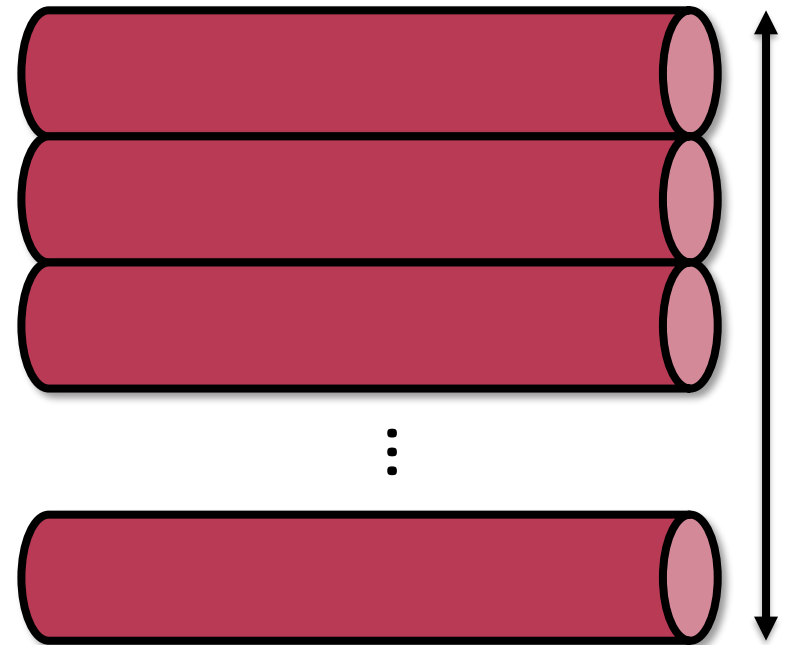
Pipe Analogy

Single Resource Instance



Throughput:
 $X_{(1)}^{max}$

K Resource Instances (with free choice)



Throughput :
 $X_{(K)}^{max} = K \times X_{(1)}^{max}$

Utilisation of a Multiple Server System

- Analogous with the previous calculation:

$$X_{(K)}^{max} = \frac{K}{T} \Rightarrow X_{(K)}^{max} \times T = \textcircled{K} ?$$

Utilisation of a Multiple Server System

- Analogous with the previous calculation:

$$X_{(K)}^{max} = \frac{K}{T} \Rightarrow X_{(K)}^{max} \times T = K = U_{(K)}^{max}$$

- The equation of utilisation in this case:

$$U = \frac{X_{(1)}}{K} \times T$$

- Intuition:

„What is the utilisation of a single instance by the average throughput of one instance?”

„What percentage of K units of time is the processing time of a single instance?”

Utilisation of a Multiple Server System

- Analogous with the previous calculation:

$$X_{(K)}^{max} = \frac{K}{T} \Rightarrow X_{(K)}^{max} \times T = K$$

- The equation of utilisation in this case:

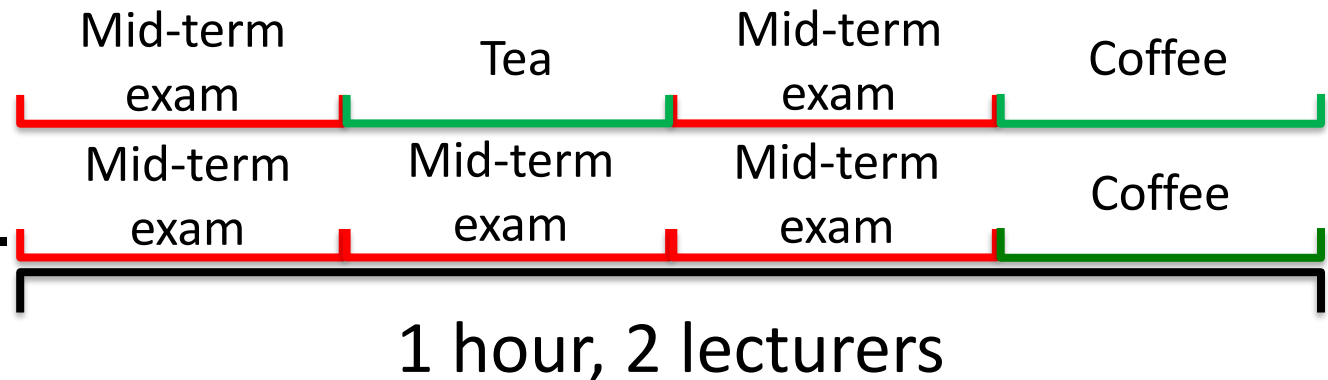
$$U = \frac{X}{K} \times T$$

- Intuition:

$$X = 5 \frac{1}{h}$$

$$T_{exam} = 15 \text{ min.}$$

$$U = 62,5\%$$



Summary

■ Stable State:

- Calculating with average values
- $\lambda = X$ (arrival rate = throughput)

■ Maximum throughput:

- Upper bound of the reachable throughput
- $X^{\max} = \frac{K}{T}$ (in case of K resource instances)

■ Utilization:

- Ratio of the actual and the maximum throughputs
- $U = \frac{X}{K} \times T$ (in case of K resource instances)