# COMP9334
# Capacity Planning for Computer Systems and Networks

Week 8A:

queues

COMP9334

1

# This lecture

- Web services
  - What is it?
  - Performance analysis
- Fork-join queue
  - Markov chain
  - MVA

<span style="color:red">Assignment Project Exam Help</span>

<span style="color:red">https://eduassistpro.github.io/</span>

<span style="color:red">Add WeChat edu_assist_pro</span>

# Web access versus Web services

**(a) Web access**

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

**(b) Composite web service for travel**

# Web service performance issues

- Metrics
  - Response time
  - Throughput
  - Availability

- Performance analysis method
  - Operational an
  - Markov chain

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Web service flow graph

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

$V_a, V_h, V_c$: Relative Visit Ratio

Every request to the travel site generates on average $V_a$ requests to the Airline web service etc.

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

$X_{TA}$ = Throughput of travel site

$X_a$ = Throughput of airline Web service

$X_a \geq V_a \times X_{TA}$

# Similarly,

$$X_a \geq V_a \times X_{TA}$$
$$X_h \geq V_h \times X_{TA}$$

Xh = Throughput of hotel web service

Xc = Throughput of car rental web service

- Can you find an upper bound o ughput of the travel site

$$X_{TA} \leq$$

# Example:

Xa = 20 requests/s
Xh = 15 requests/s
Xc = 10 requests/s
Va = 4, Vh = 2, Vc = 1

The airline web service is the bottleneck of the travel web site.

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

## What is the bound on throughput of web service A?

# Bound on the throughput of web service A is:

Assignment Project Exam Help
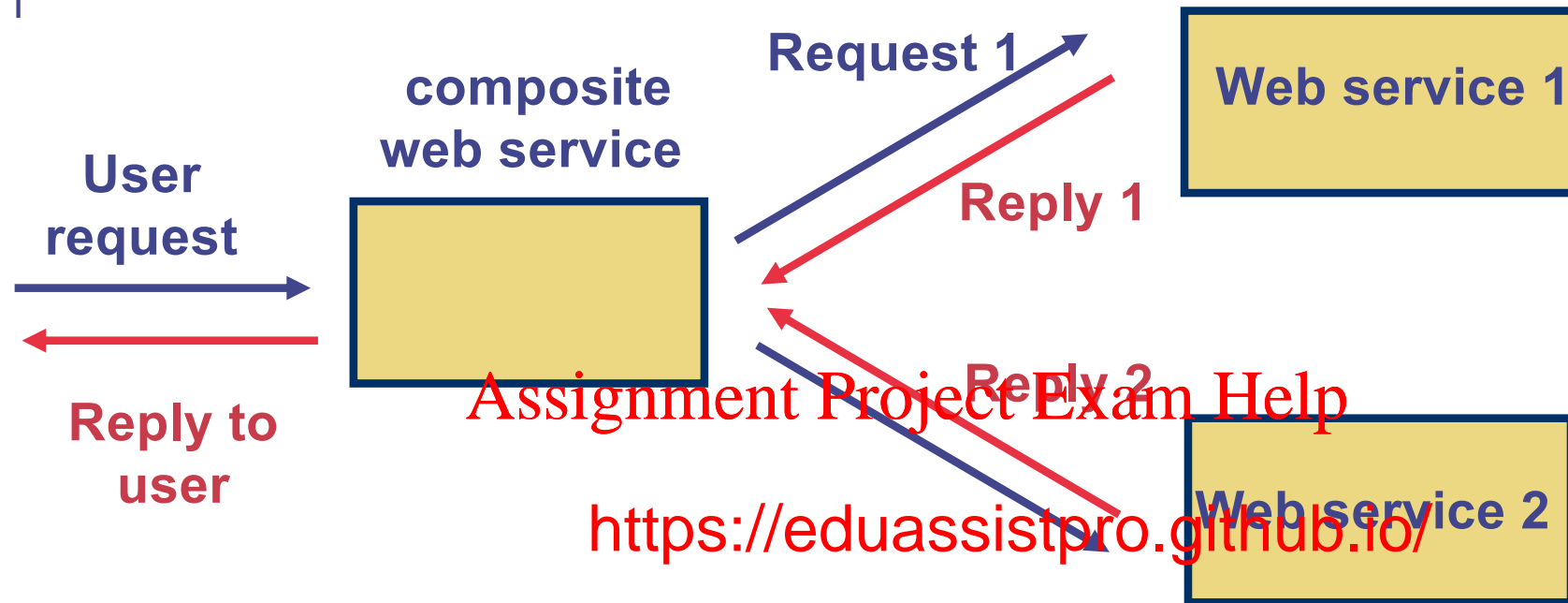
https://eduassistpro.github.io/

Add WeChat edu_assist_pro

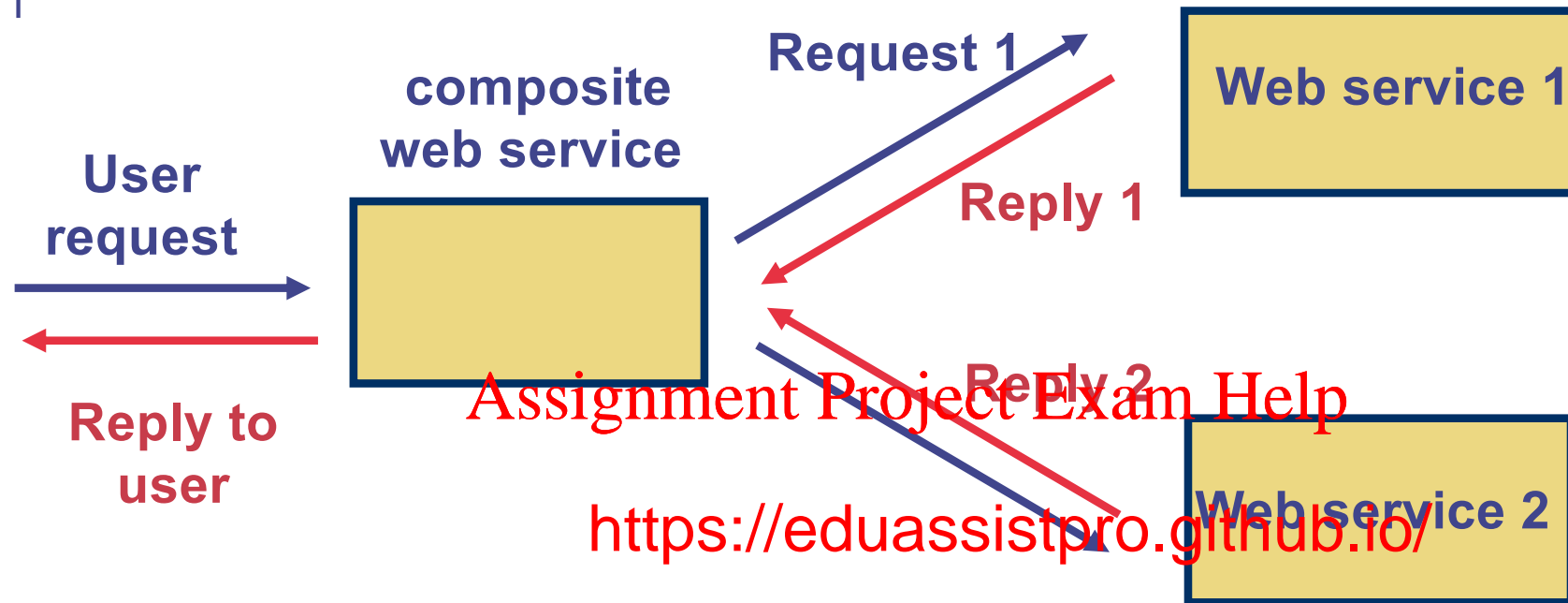# Response time analysis

- The bottleneck analysis only gives an upper bound on the throughput

- Can we find the response time?

  - Markov chain

  - Approximate MVA

- We begin with a

# A simple web service scenario (1)

**composite web service**

**Request 1** → **Web service 1**

**User request** →

← **Reply to user**

**Reply 1**

**Reply 2**

**Web service 2**

Assignment Project Exam Help

https://eduassistpro.github.io/
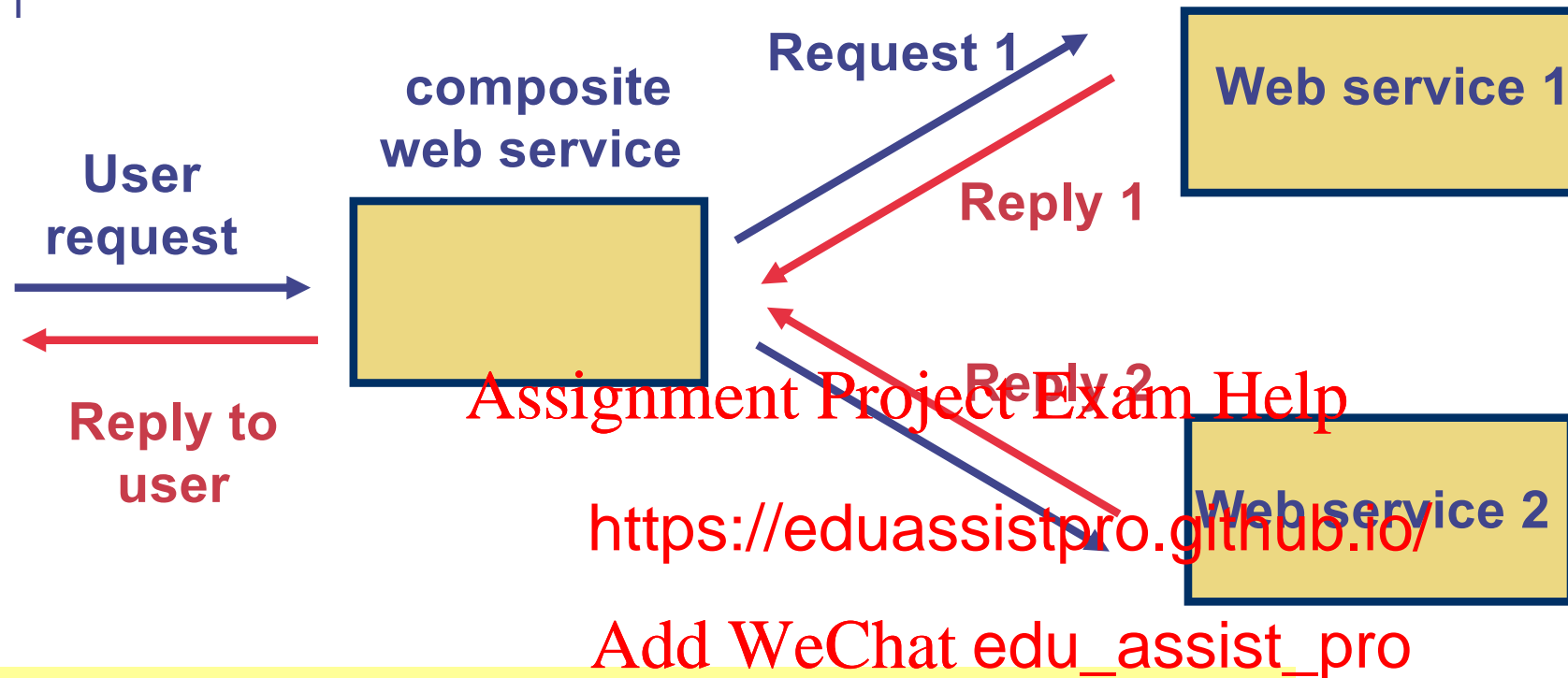
Add WeChat edu_assist_pro

- A composite web service uses two

- Sequence of events

  1. Composite web service receives a user request
  2. Composite web service sends Request 1 and Request 2
  3. The web services reply *independently*
     - *That is, Reply 1 and Reply 2 may arrive at different times*
  4. After the composite web service receives *both* replies, it responds to the user
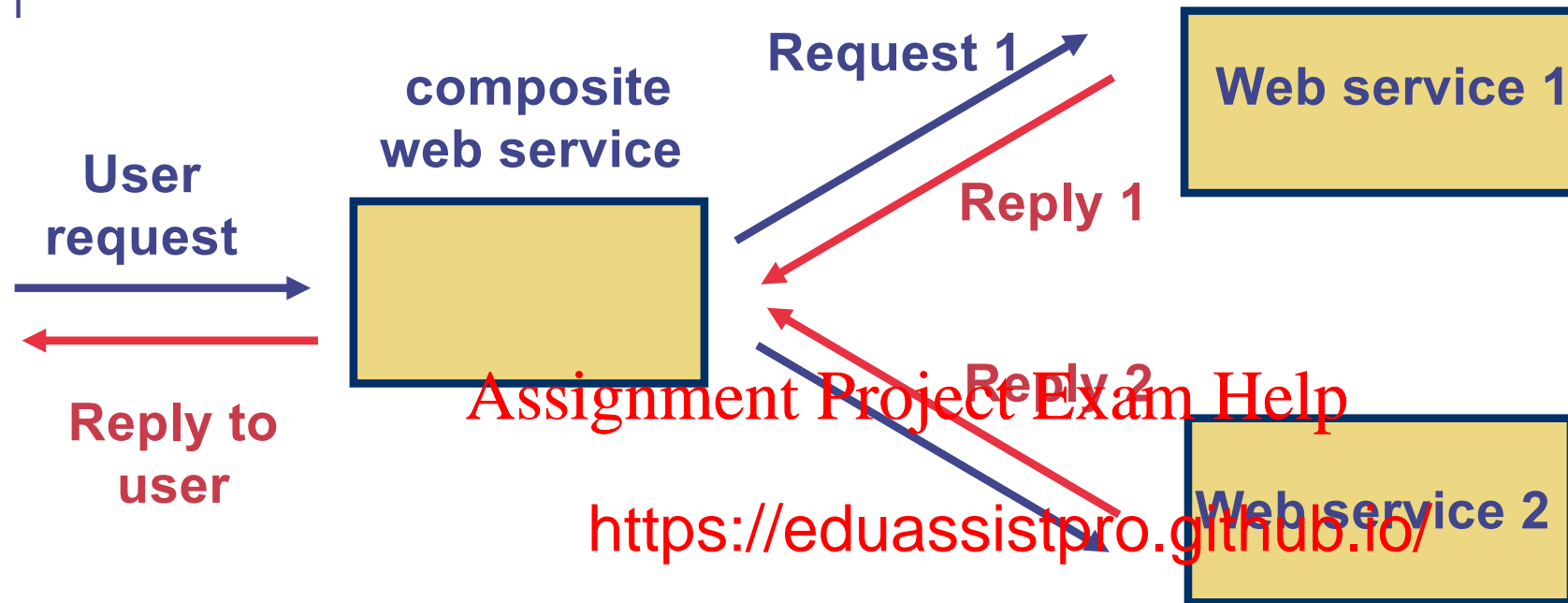
# A simple web service scenario (2)

**composite web service**

**User request**

**Request 1**

**Web service 1**

**Reply 1**

**Reply to user**

Reply 2

Web service 2

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- Recall the definition of respons

- Response time of Web Service 1
    - = Time at which composite web service receives Reply 1 *minus*
        Time at which composite web service sends Request 1
- Similarly for Web Service 2.

# A simple web service scenario (3)

**composite web service**

**User request**

**Request 1**

**Web service 1**

**Reply 1**

**Reply to user**

Reply 2

Web service 2

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- Assuming that:
  - Web service 1 has a response time distribution of
    - 0.2s with probability 0.5
    - 0.3s with probability 0.5
  - Web service 2 has a response time distribution of
    - 0.2s with probability 0.5
    - 0.3s with probability 0.5
- What is the average time that the composite web service has to wait until both replies are returned?

# A simple web service scenario (4)

**composite web service**

**Request 1**

**Web service 1**

**User request**

**Reply 1**

**Reply to user**

**Reply 2**

**Web service 2**

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- What if the service time distribution
  - Web service 1 has a response time distribution of
    - 0.2s with probability 0.5
    - 0.3s with probability 0.5
  - Web service 2 has a response time distribution of
    - 0.2s with probability 0.5
    - 0.5s with probability 0.5
- What is the average time that the composite web service has to wait until both replies are returned?

# Analysis scenario

- Lesson learnt: Slow web services can become the bottleneck for composite web service

- We consider Composite Web Services (illustration next slide) Assignment Project Exam Help

  - With parallel in

  - Web services https://eduassistpro.github.io/ service time of S (exponentially distributed) Add WeChat edu_assist_pro

  - Web service N has a mean ser $g \times S$ (exponentially distributed)

  - The next service step can only be completed after all these N steps have been completed.

$$\frac{1}{\alpha} \times S$$

Note that if $\alpha$ < 1, then server N is slower than the other (N-1) servers.

**Servers 1 to N-1 : mean response time = S**

**Server N: mean response time =** $\alpha \times S$ $\leftarrow$ $\dfrac{1}{\alpha} \times S$

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Fork-join system

- The type of system described earlier is known as fork-join system
  - Fork is referring to the parallel invocation
  - All services must complete at the joining point before the next service can start
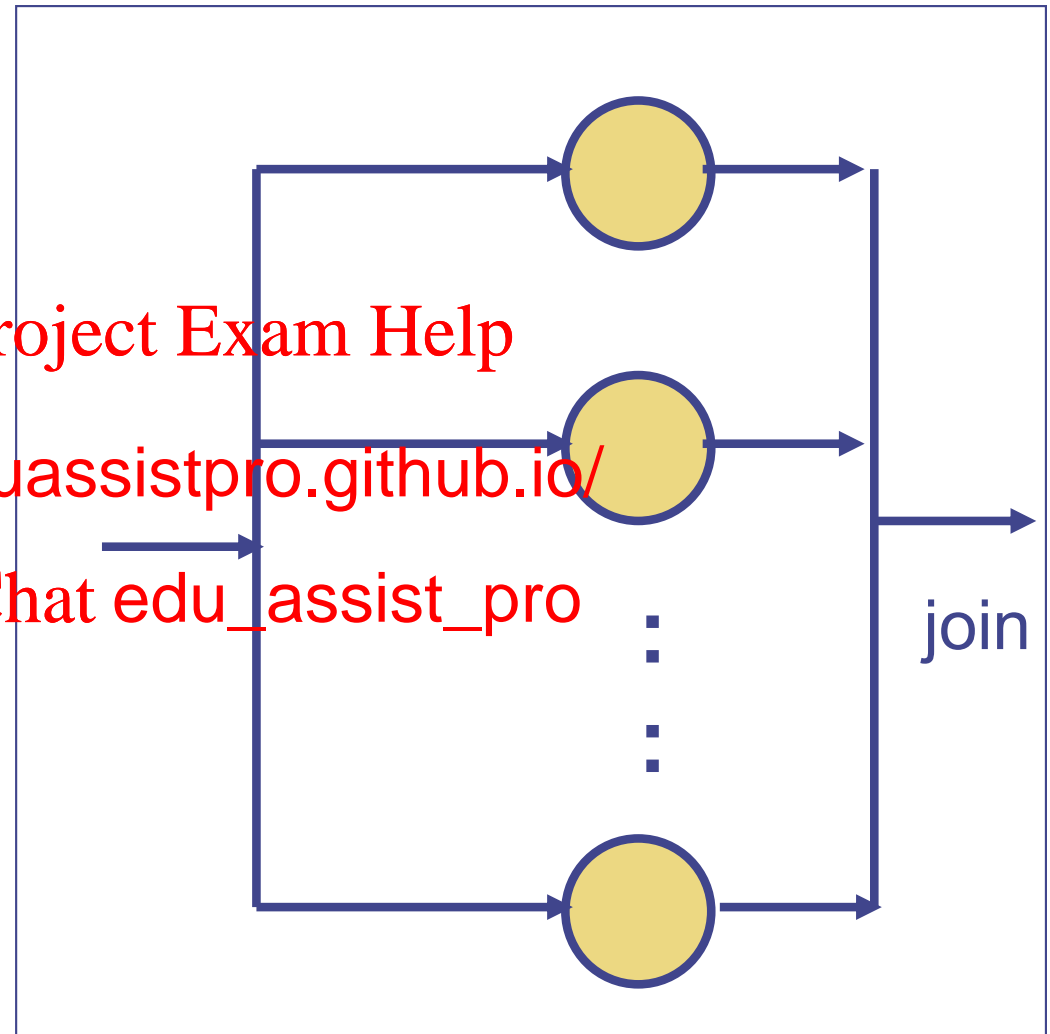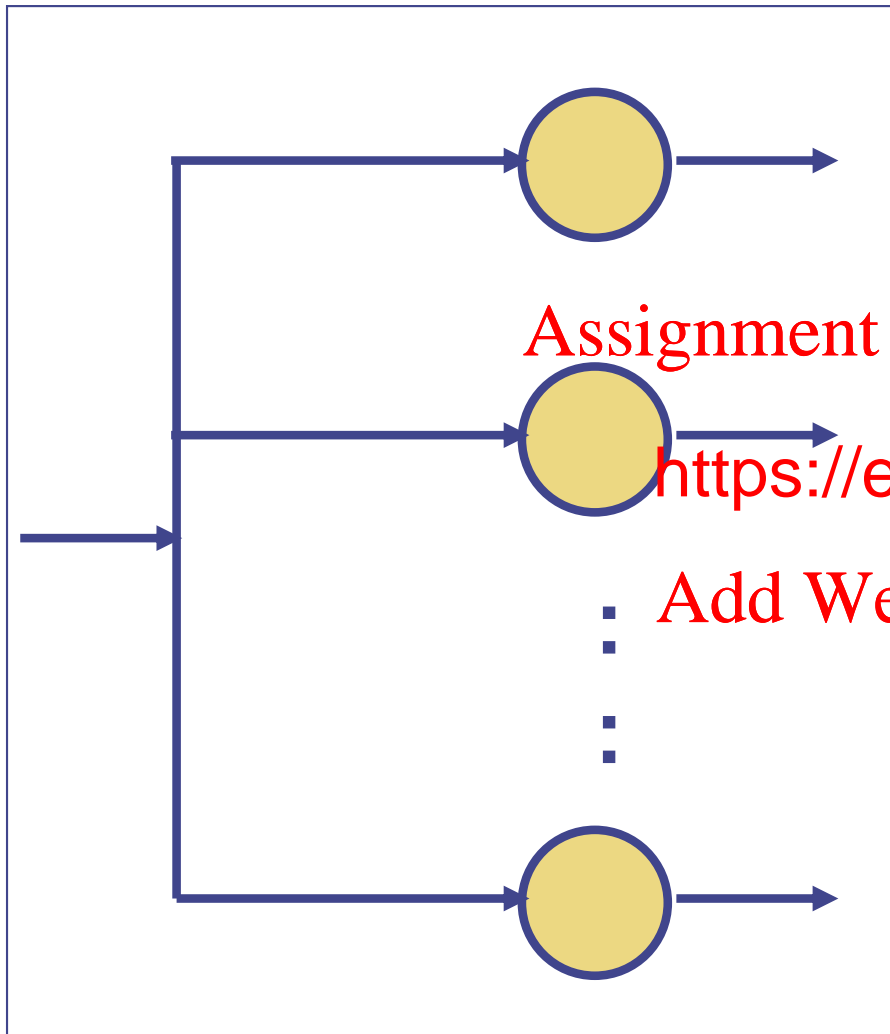
Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# You've seen parallel processing before:

## M/M/m queue

## Fork-join queue



Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

join

**What is the difference between these two queueing networks?**

**Servers 1 to N-1 : mean response time = S**
**Server N: mean response time = $\alpha \times S$** $\leftarrow$ $\frac{1}{\alpha} \times S$

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- We want to understand how $\alpha$ affects the response time of the composite web services
- Let T($\alpha$) = response time as a function of $\alpha$

# What is T(1)?

- In this case, all constituent web services have the same response time distribution

- If all mean response times are exponentially distributed with mean S

$$H_N = \text{N-th harmonic number}$$

(We will explain how this is obtained later.)

# How about T($\boldsymbol{\alpha}$) for ~~$\boldsymbol{\alpha > 1}$~~ ?

$\boldsymbol{\alpha} < 1$ or $\dfrac{1}{\alpha} > 1$

- We use Markov chain.

- States (i,j,k)

  - i (i = 0,…,N−1) is the number of web services still running in fast Web services

  - j (j = 0,1) is the number of web services running on the slow Web service

  - k (k = 1,2,..,N) is yet to complete

- Define $\mu = \dfrac{1}{S}$

P(N)

αμ

μ

(N-1,1,N)

(N-1)μ

**α**μ

Q(N-2) (N-2,1,N-1)

(N-1,0,N-1) R(N-1)

**α**μ

(N-2)μ

(N-1)μ

Q(N-3) (N-3,1,N-2)

(N-2,0,N-2) R(N-2)

(N-

(N-2)μ

Q(N-4) (N-4,1,N-3)

3,0,N-3) R(N-3)

(N-4)μ

**α**μ

Q(1) (1,1,2)

(2,0,2) R(2)

μ

**α**μ

μ

(0,1,1) Q(0)

(1,0,1) R(1)

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

$\boldsymbol{\alpha}$

$T(\boldsymbol{\alpha})$

$\boldsymbol{\alpha}$

When $\boldsymbol{\alpha} = 1$,

$T(\boldsymbol{\alpha}) = H_N\ S$

$\boldsymbol{\alpha}$

Assignment Project Exam Help

T($\alpha$)/S  https://eduassistpro.github.io/

Add WeChat edu_assist_pro

$\alpha$

$\dfrac{1}{\alpha}$

# Other examples of fork-join QNs

- Disk array, e.g. RAID (= Redundant Array of Independent Disks)



RAID controll

Assignment Project Exam Help

f RAID1

https://eduassistpro.github.io/

sks

Add WeChat edu_assist_pro

# Fork-join in disk array

I/O requests

RAID controlle

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

**Example 1**
Read a file in parallel
1st half of the file from Disk 0
2nd half of the file from Disk 1
Need to wait for both halves of
re the next operation

2
Write to disk.
Need to write to both disks
(for consistency)
Need to wait for both disks
to complete

Disk 0

| Block 1 |
| Block 2 |
| Block 3 |
| Block 4 |
| Block 5 |

Disk 1

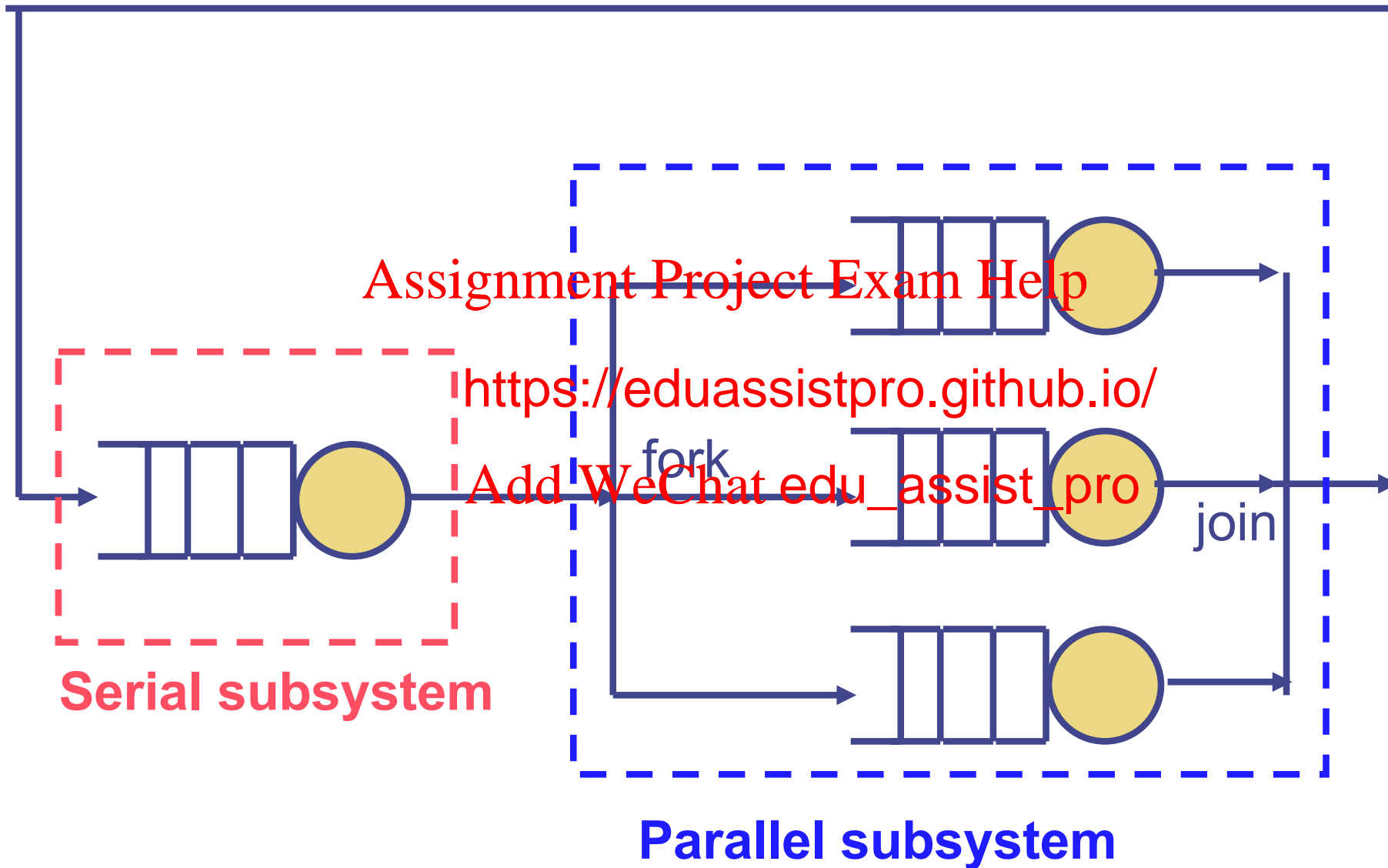| Block 1 |
| Block 2 |
| Block 3 |
| Block 4 |
| Block 5 |

# Fork-join queueing networks

- Exact results are hard to come by
- Approximate solution methods are used

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# A Queueing network with a fork-join subsystem



Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

fork

join

**Serial subsystem**

**Parallel subsystem**

# Approximate MVA for fork-join queueing networks

- For MVA with fork-join, the basic unit is a subsystem
  - A subsystem can be either a serial subsystem (= a device) or parallel one
    - A serial subsystem is a special case of parallel subsystem
  - In comparison, the basic unit for MVA before is a device

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Arrival Theorem for Parallel Subsystems (1)

- Consider a parallel subsystem with *k* parallel service centres
- The average time each job requires at each service centre is S (exponentially distributed)

<span style="color:red">Assignment Project Exam Help</span>

<span style="color:red">https://eduassistpro.github.io/</span>   $\mu = 1/S$

<span style="color:red">Add WeChat edu_assist_pro</span>

## Arrival Theorem for Parallel Subsystems (2)

When there are $n-1$ jobs in the whole QN, the average number of jobs in the subsystem is $z$. When there're n jobs in the system

One of the
n jobs (customers)

$$\text{Waiting time} = S \times z; \text{Service time} = S \times H_k$$

$$\Rightarrow \text{Response time} = S \times (H_k + z)$$

Note that if $k = 1$, the subsystem is serial and is identical to a device in MVA analysis that we have seen before.

$$\text{Response time} = S \times (H_1 + z)$$

This is the same arrival theorem that we've seen before.

## Notation:

$I = $ Number of subsystems in the QN

$S_i = $ Avg. service time of a station in subsystem $i$

$\bar{n}_i(n) = $ Avg. # of jobs at subsystem $i$
    when there're $n$ jobs in the QN

$V_i = $ Visit ratio of subsystem $i$

## MVA for fork-join systems:

| Mean # jobs in each subsystem | $\bar{n}_i(n-1)$ |

$$R_i(n) = S_i \times (H_i + \bar{n}_i(n-1))$$

| Mean response time | $R_i(n)$ |

| Throughput of the system | $X_0(n)$ |

$$\bar{n}_i(n) = V_i \times X_0(n) \times R_i(n)$$

| Mean # jobs in each subsystem | $\bar{n}_i(n)$ |

# Example

- A system consists of a processor and 2 disk arrays
- Disk arrays operate under synchronous workload
  - Transactions are blocked until I/O are completed

| | Service demand | # parallel systems |
|---|---|---|
| Processor | | |
| Disk array 1 | 0.02 | |
| Disk array 2 | 0.03 | |

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

**What is the system response time when there are 50 transactions? How many transactions can the system have if the system response time should not exceed 1s?**

# Exercise

- The MVA algorithm on p.35 assumes that you have both visit ratios $V_i$ and mean service time $S_i$ available

- You may recall that service demand $D_i = V_i * S_i$

- Now, let us assume that you are only given the service demands $D_i$. This means that you do not know $V_i$ and $S_i$. How ~~so that it can wo~~ A algorithm on p.35 e demands only?

# References (1)

- Web services
  - D. Mensace et al. Static and Dynamic Processor Scheduling Disciplines in Heterogeneous Parallel Architectures," *Journal of Parallel and Distributed Computing,* Vol. 28 (1), July 1995, pp. 1-18.
  - D. Mensace, "QoS Issues in Web Services," *IEEE Internet Computing*, November/December 2002, Vol. 6, No. 6.
  - D. Mensace, "Respsite Web Services," *IEEE Internet Com*, Vol. 8, No. 1
  - D. Mensace, "Composing Web Servicew," D. Menasce, IEEE Internet Computing, Vol. 8, No. 6, 200
  - These papers can be downloaded from the course website (use your CSE password)
    - We didn't cover the last paper but it's well worth a read.
- Derivation of Markov chain on pp. 22-24 is further explained in the file *forkjoin_mc.pdf*

# References (2)

- Fork-join MVA
  - Menasce et al.,"Performance by desing". Section 15.6.
- Addition references outside the scope of this course
  - Tutorial on RAID http://www.slcentral.com/articles/01/1/raid/

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro