

*Slide supporting material*

# **Lesson 18: Networks of Queues and Exercises**

**Giovanni Giambene**

***Queuing Theory and Telecommunications:  
Networks and Applications***  
**2nd edition, Springer**

**All rights reserved**



# **Modeling a Network: Network of Queues**

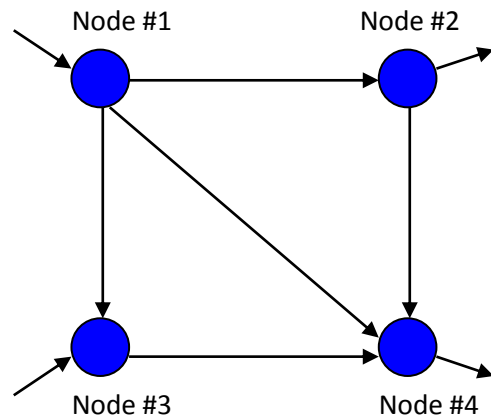
# Introduction



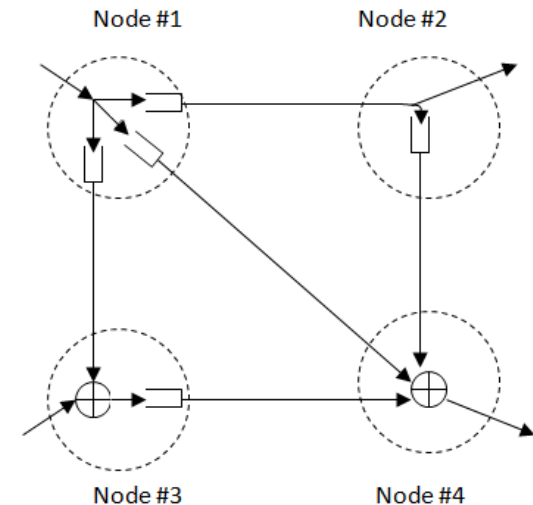
- The interest is in considering networks where nodes exchange traffic.
  - **Open networks**, where traffic can be received and sent outside the network.
  - **Closed networks**, where the traffic cannot be exchanged with external nodes. Closed networks are more related to the modeling of digital computing systems.
- Our interest is on **open networks** that are well suited to model IP networks, where **different nodes (modeled by means of queues)** exchange data traffic in the form of variable-length messages.

# Store-and-Forward Networks and Model as Net. of Queues

- The network is formed of nodes and links.



**Network**



**Model of the system as a network of queues (store-and-forward nodes)**

# Model of Open Network of Queues

- We consider a model, where **the generic i-th node** receives **input traffic** with mean rate  $\lambda_i$  from outside the network and receives also traffic routed from other nodes of the network that contribute a total mean input rate indicated by  $\Lambda_i$ .
- Each **arrival** corresponds to a message with (in general) a random length.
- The total arrival process at the i-th node is randomly split among the different **outgoing links** from the i-th node.

# Model of Open Network of Queues (cont'd)

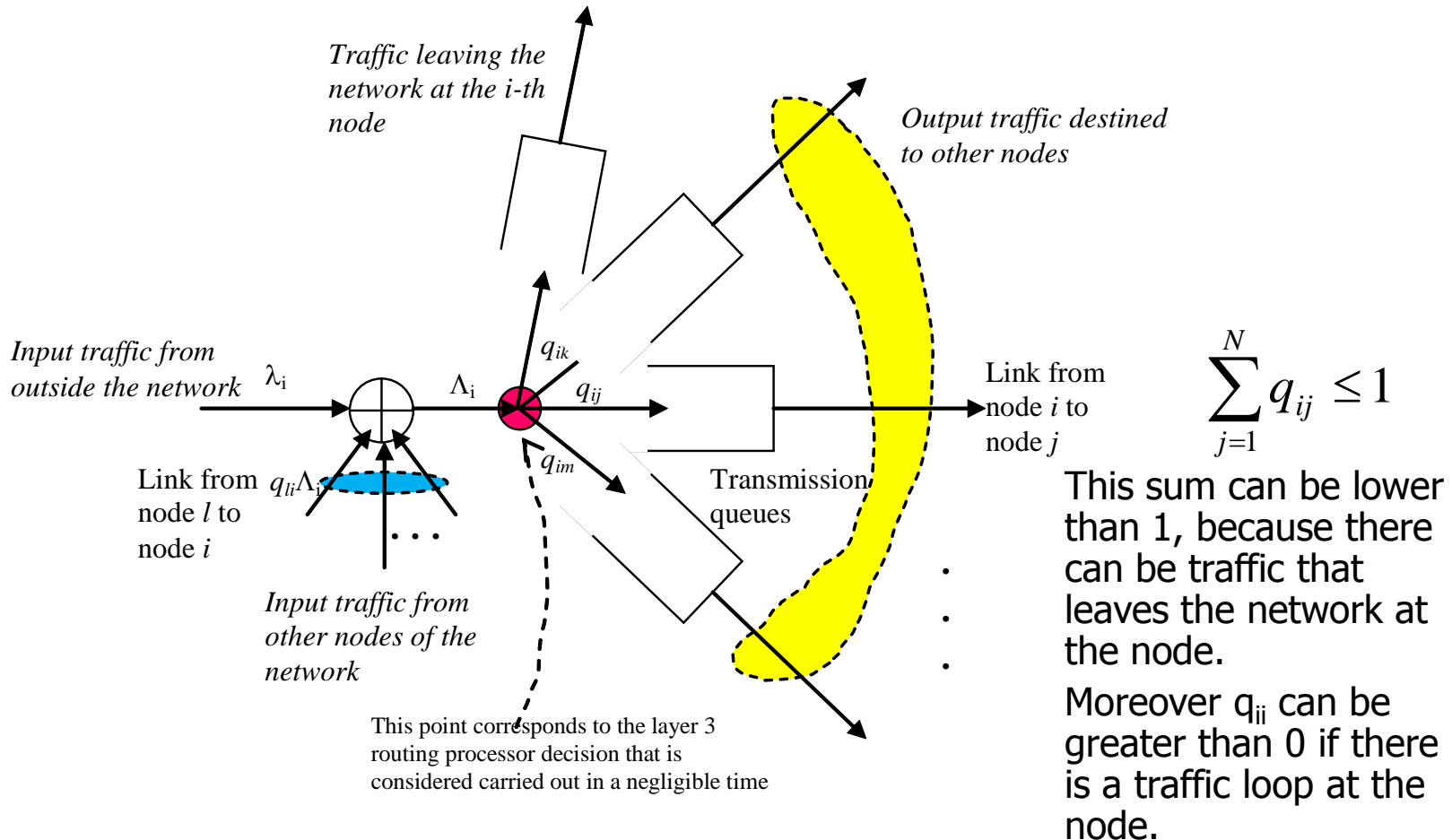
- Each **link is modeled by a buffer and a transmission line (i.e., one server)** with a suitable capacity.
  - **We consider queues with infinite rooms (i.e., no loss phenomena).**
- Let  $q_{ij}$  denote the split probability for the total traffic of the  $i$ -th node to be routed to the  $j$ -th node of the network;  **$1 - \sum q_{ij}$  denotes the probability that the traffic leaves the network at the  $i$ -th node.**
  - Under **stability assumptions**, the traffic carried by the generic link from node  $i$  to node  $j$  is  $\Lambda_i q_{ij}$
  - $q_{ii}$  can also be different from 0 if there is **traffic looped back onto the same node**. This modifies the burstiness characteristics of the input traffic.

# Model of Open Network of Queues (cont'd)

- In our generic network model, we consider the set of **nodes labeled with numbers  $i$  from 1 to  $N$**  and the related set of **links (modeled by queues) labeled with numbers  $k$  from 1 to  $L$** .
- **We can study this network at the level of nodes or at the level of links.**

# Model of Open Network of Queues (cont'd)

- The generic  $i$ -th node can be described as depicted below.

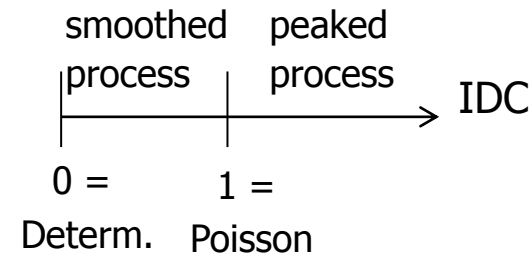




# Characteristics of the Arrival Process

- To study the characteristics of the arrival process at a queue we can refer to the Index of Dispersion for Counts (IDC), using to the **number of arrivals in a given interval  $t$ ,  $N(t)$** .
- **IDC is the ratio between the variance of  $N(t)$  and the mean of  $N(t)$  referring to the same interval:**

$$IDC(t) = \frac{Var[N(t)]}{E[N(t)]}$$



- For a Poisson process  $IDC(t) \equiv 1, \forall t$ . An arrival process is **peaked** if  $IDC(t) > 1$ ; an arrival process is **smoothed** if  $IDC(t) < 1$ . When IDC reduces, arrivals are more regularly spaced in time. The limiting case is when  $IDC = 0$ : the arrival process is **deterministic**. Conversely, when  $IDC > 1$ , arrivals tend to occur in bursts, thus entailing problems in terms of congestion and delays at the queues.

# The Case of Poisson Input Traffic

- In the case of Poisson arrivals of messages with mean rates  $\lambda_i$  (uncorrelated from node to node), the total arrival process for the different nodes may lose the Poisson characteristic if:
  - There are traffic feedback loops causing a **peaked arrival process**. A network that allows (does not allow) feedback loops is **cyclic (acyclic)**.
    - Acyclicity means that one message does not cross a network node more than once in its path from source to destination (i.e., no routing loops).
  - Queues with finite rooms drop arrivals exceeding their capacity; in this case, the **circulating traffic is smoothed**. However, in this study **we will not consider queues with finite rooms and packet losses**.

# Correlation in the Behavior of the Queues

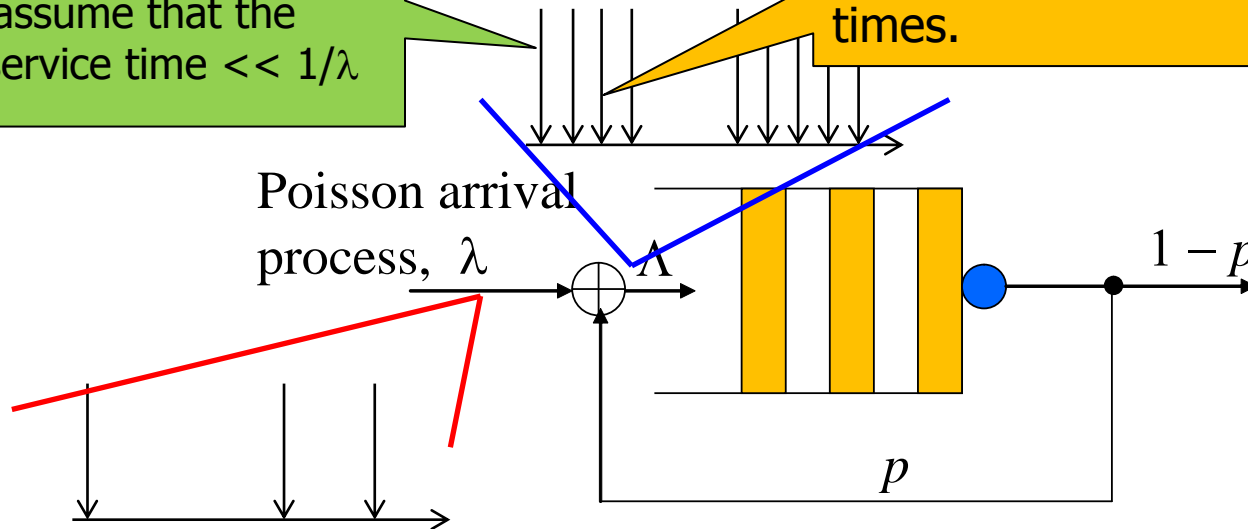


- There is a **strong correlation in the behaviors of the queues in the network** and this is due to:
  - The correlation of the arrival process and the service process due to **feedback loops** (cyclic network).
  - The correlation in the behavior of the different nodes due to the fact that **the same message is serviced at the different nodes** crossed in the network along the path from source to destination.

# Elementary Network of Queues with Feedback

Let us assume that the mean service time  $\ll 1/\lambda$

Due to the feedback, these arrivals are spaced by the message service times.



In this feedback queue, input and output processes are continuous time (we do not consider slots as done in previous exercises).

- The total input arrival process at the queue is **bursty** (not Poisson), with IDC greater than 1.
- This elementary network of queues will be studied at the end of this lesson.

# Traffic Rate Equations for a Network of Queues

## ■ Hypotheses:

1. Stable queues
2. No packet loss (i.e., infinite rooms in the queues)
3. Stochastic routing at the nodes.

The network can be cyclic or acyclic.

**In order to write the following traffic rate equations, we work at the level of nodes (not queues / links).**

# Traffic Rate Equations for a Network of Queues (cont'd)

- Thesis: We can write the following **balance for the total input traffic** with rate  $\Lambda_i$  for the  $i$ -th node (i.e., traffic rate equation for the  $i$ -th node):

$$\left\{ \Lambda_i = \lambda_i + \sum_{j=1}^N \Lambda_j q_{ji} \quad i = \{1, 2, \dots, N\} \right.$$

- This is a **linear system of N equations in N unknown terms**  $\Lambda_i$  (input arrival rates from outside the network,  $\lambda_i$ , and split probabilities  $q_{ij}$  are considered to be known).
- Note that this system can be solved under **general assumptions** (it is not requested that the input traffic be Poisson).
- Basically: **one traffic rate equation can be written per each sum point in the network of queues.**

# Little Formula applied to the Whole Network

- The Little theorem can be applied not only to a queue, but also to the whole network of queues.
- In this case, we refer to the network modeled at the level of links that are queues in our model ( $k = 1, \dots, L$ ).
- Let  $\mathfrak{S}_k$  denote the mean number of messages in the  $k$ -th queue:  $\mathfrak{S}_k = \mathfrak{S}_k(\rho_k)$ , where  $\rho_k = \Lambda_i q_{ij} / \mu_k$  and  $\mu_k$  is the service rate of the queue. **Let  $T$  denote the mean message delay from input to output of the network.**
- The **Little theorem applied to the whole network** can be expressed as
$$T = \frac{\mathfrak{S}_{tot}}{\lambda_{tot}}$$
where  $\mathfrak{S}_{tot} = \sum_{k=1}^L \mathfrak{S}_k$  and  $\lambda_{tot} = \sum_{k=1}^N \lambda_k$  ( $\lambda_{tot}$  denotes the total mean arrival rate from outside the network).
- We need to derive  $\mathfrak{S}_k(\rho_k)$ , as shown in the next slides.

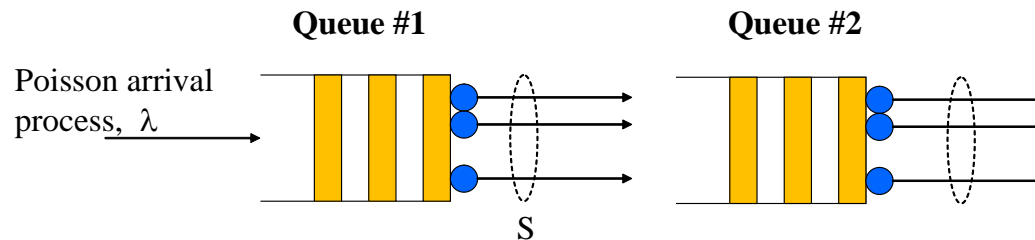


# Tandem Queues



# Burke Theorem for Tandem Queues

- We study two tandem queues (or, in general, a network of tandem queues).

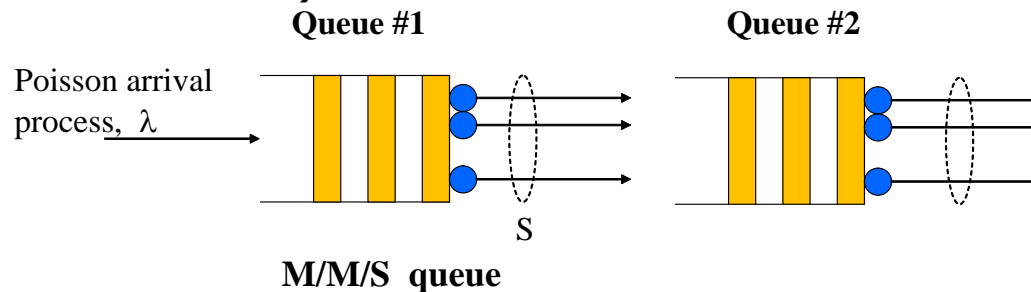


- **Hypotheses:**

1. **Tandem queues:** all messages leaving a queue are at the input of the next queue (the service completion instant for a queue is the message arrival instant at the next queue)
2. **Same hypotheses of the traffic rate equations (stability, no loss)**
3. **Poisson arrival process from outside**
4. **Exponentially distributed service times.**

# Burke Theorem for Tandem Queues (cont'd)

- Under the stability assumption, the first queue admits an M/M/S model (Poisson arrivals/exponentially-distributed service times/S servers, infinite rooms).



- Under stability conditions for the first queue, we can state that the mean output rate from the first queue is  $\lambda$ , even without considering the specific characteristics of the first queue.
- It is possible to prove that the whole output process from the first M/M/S queue is Poisson with mean rate  $\lambda$ .**
  - The time intervals between service completion instants are exponentially distributed with mean rate  $\lambda$ .
  - In the following slide, we provide the proof in the case  $S = 1$ .

# Burke Theorem for Tandem Queues (cont'd)

- Let us consider a generic  $M/G/1$  queue where  $g(t)$  denotes the service time probability density function [let  $G(s)$  denote the related Laplace transform]. Let  $P(s)$  denote the Laplace transform of the density function of the interarrival times between subsequent service completion events.
- We determine  $P(s)$  by considering two cases: (i) non-empty queue; (ii) empty queue.
- **Derivation of  $P(s \mid \text{non-empty queue})$ :** In this case, times between completion events have a probability density function  $g(t)$  with Laplace transform  $G(s)$ :  $P(s \mid \text{non-empty queue}) \equiv G(s)$ .
- **Derivation of  $P(s \mid \text{empty queue})$ :** in this case, we have to wait for the next arrival time that is characterized by an exponentially-distributed time (with mean rate  $\lambda$ ). Hence, the time to the next completion is the sum of two independent contributions: an interarrival time and a service time. In the Laplace domain, we have that  $P(s \mid \text{empty queue})$  is given by the product of two contributions:  $P(s \mid \text{empty queue}) \equiv [\lambda / (\lambda + s)] \times G(s)$ .



# Burke Theorem for Tandem Queues (cont'd)

- We remove the conditioning on  $P(s)$  by means of the probability of an empty and of a non-empty M/G/1 queue,  $P_0$  and  $1 - P_0$ , respectively. We know that  $P_0 = 1 - \lambda E[X]$ , where  $E[X]$  is the mean value related to the density function  $g(t)$ .

$$\begin{aligned} P(s) &= P(s \mid \text{empty queue})P_0 + P(s \mid \text{non-empty queue})(1 - P_0) = \\ &= \frac{\lambda}{\lambda + s} G(s)(1 - \lambda E[X]) + G(s)\lambda E[X] \end{aligned}$$

- **M/M/1 case:**  $g(t)$  is exponentially-distributed with mean rate  $\mu$ ,  $G(s) = \mu/(\mu + s)$  and  $E[X] = 1/\mu$ . Substituting these expressions in  $P(s)$ , we have:

$$\begin{aligned} P(s) &= \frac{\lambda}{\lambda + s} \frac{\mu}{\mu + s} \left[ 1 - \lambda E[X] + \frac{\lambda + s}{\lambda} \lambda E[X] \right] = \frac{\lambda}{\lambda + s} \frac{\mu}{\mu + s} \left[ 1 - \frac{\lambda}{\mu} + \frac{\lambda + s}{\lambda} \frac{\lambda}{\mu} \right] = \\ &= \frac{\lambda}{\lambda + s} \end{aligned}$$

➡ The completion process (output process) is Poisson with mean rate  $\lambda$  [QED].



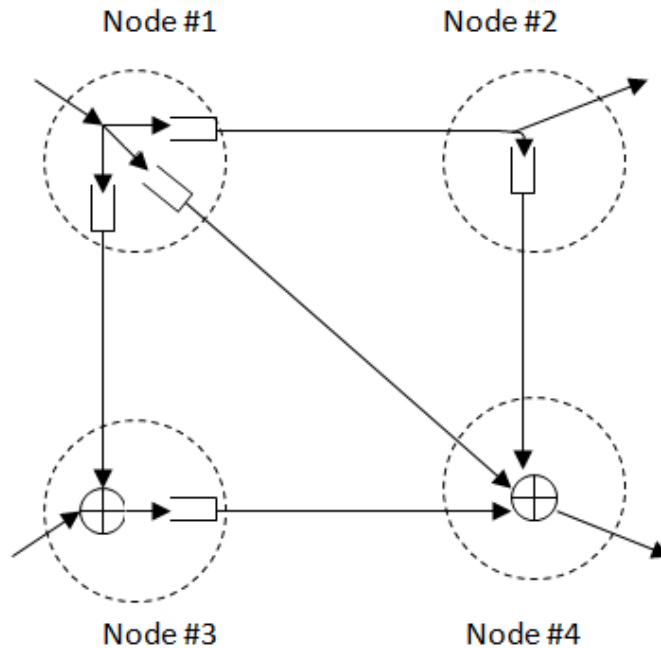
# **Feedforward Networks**

# Feedforward Networks

- **Feedforward networks are characterized as: same hypotheses of traffic rate equations + Poisson arrivals + exponentially distributed service times + acyclicity.**
- The Poisson characteristic of the input processes is maintained within the network nodes using: (i) the **random split** model for distributing the traffic of a node on the different output links; (ii) the **Burke theorem**; (iii) **independent Poisson input processes** at the nodes; (iv) the **sum** of independent Poisson processes.

# Feedforward Networks (cont'd)

- Feedforward network example:



# Feedforward Networks (cont'd)

- For the sake of simplicity, let us consider from now on just one server per queue in the network.
- Each queue is of the M/M/1 type with input traffic given by the solution of the traffic rate system. The joint state probability has a **product form**: the queues are independent (the number of messages in the queues are independent).

$$P(n_1, n_2, \dots, n_N) = P(n_1) \times P(n_2) \times \dots \times P(n_N), \text{ where } P(n_i) = (1-\rho_i)\rho_i^{n_i}$$

- Note that the presence of feedback paths in the networks **destroys the Poisson characteristics of the flows** and the Burke theorem cannot be applied. **Nevertheless, the product form still holds under the assumptions that will be considered in the next slides for the Jackson theorem.**





# **Cyclic/Acyclic Networks and the Jackson Theorem**

# Jackson Theorem for Networks of Queues

Hypotheses:

1. An open network with independent Poisson arrivals of messages at each node
2. Queues modeling the transmissions on links with infinite rooms (no packet loss), stable behavior, and single server
3. Exponential service times at the nodes with FIFO discipline
4. **Arrival process and service time process are independent**
5. Stochastic routing whereby the next node, after service completion, is chosen independently from message to message.

Thesis:

- The joint probability distribution of queue occupancies has a **product form** with the product of distributions of individual **M/M/1 queues**:

$$P(n_1, n_2, n_3, \dots, n_M) = (1-\rho_1)\rho_1^{n_1}(1-\rho_2)\rho_2^{n_2}(1-\rho_3)\rho_3^{n_3}\dots(1-\rho_M)\rho_M^{n_M}.$$

- The mean number of requests in each queue and the related mean delay are according to the classical M/M/1 formula (Poisson processes in the network).

J. R. Jackson, "Jobshop-like Queueing Systems", in *Management Science*, Vol. 10, No. 1, pp. 131-142, October 1963.

J. F. Hayes, T. V. J. Ganesh Babu. *Modeling and Analysis of Telecommunications Networks*. John Wiley & Sons, NJ, 2004

# Jackson Theorem for Networks of Queues

The Jackson network is an abstract concept! Especially assumption #4 can be strong in a real network.

Hypotheses:

1. An open network with independent Poisson arrivals at each node (ass),
2. Queues modeling the transmissions or service times, and single server
3. Exponential service times at the nodes with first-come first-served discipline
4. **Arrival process and service time process are independent**
5. Stochastic routing whereby the next node, after service completion, is chosen independently from message to message.

Thesis:

- The joint probability distribution of queue occupancies has a **product form** with the product of distributions of individual **M/M/1 queues**:

$$P(n_1, n_2, n_3, \dots, n_M) = (1-\rho_1)\rho_1^{n_1}(1-\rho_2)\rho_2^{n_2}(1-\rho_3)\rho_3^{n_3}\dots(1-\rho_M)\rho_M^{n_M}.$$

- The mean number of requests in each queue and the related mean delay are according to the classical M/M/1 formula (Poisson processes in the network).

J. R. Jackson, "Jobshop-like Queueing Systems", in *Management Science*, Vol. 10, No. 1, pp. 131-142, October 1963.

J. F. Hayes, T. V. J. Ganesh Babu. *Modeling and Analysis of Telecommunications Networks*. John Wiley & Sons, NJ, 2004

# Kleinrock Independence Assumption for Store-and-Forward Networks

- In order to apply the Jackson theorem to store-and-forward networks, we consider to add the **independence assumption**, which was guessed by Kleinrock (1964).
  - In the queuing networks we have dealt with up to this point, we considered that the **service times are associated with the servers and that servers are independent from queue to queue**. In store-and-forward (real) networks, this is not possible since the service time depends on the length of the message, which is the same from queue to queue. **This introduces dependencies between the arrival process and the service process. Feedback loops are a special case of this.**
  - **Independence assumption**: the service time of a message is chosen independently each time it passes through a node. This permits to reapply assumption #4 of Jackson networks also to real networks.
  - **This assumption could be strong** and is more acceptable when there is a sufficient mix of different sources in the network and the network has a high number of nodes. This assumption has been verified by means of simulations.

L. Kleinrock. *Communication Nets: Stochastic Message Flow and Delay*. Dover Books on Engineering, NY, 2007

# Kleinrock Independence Assumption for Store-and-Forward Networks

- In order to apply the Jackson theorem to store-and-forward networks which w

The Kleinrock assumption operates “as if” we could **remove feedback loops** in the network of queues so that **queues are decoupled and correlations in the network are removed!** Then, it is as if the traffic flows in the network were Poisson!

- In the the s inde is not is the the a case of this.

- **Independence assumption**: the service time of a message is chosen independently each time it passes through a node. This permits to reapply assumption #4 of Jackson networks also to real networks.
- **This assumption could be strong** and is more acceptable when there is a sufficient mix of different sources in the network and the network has a high number of nodes. This assumption has been verified by means of simulations.

L. Kleinrock. *Communication Nets: Stochastic Message Flow and Delay*. Dover Books on Engineering, NY, 2007

# Kleinrock Application of the Jackson Th. to Store-and-Forward Networks

Hypotheses:

1. An open network with independent Poisson arrivals at each node.
2. Single-server queues modeling the transmissions on links with infinite rooms, stable behavior, and single server.
3. Exponential service times at the nodes with FIFO discipline.
4. **Kleinrock independence assumption.**
5. Stochastic routing at each node.

Thesis: **Jackson theorem can be applied** and then

- **Each queue behaves as it was M/M/1 (a product-form expression is valid for the joint state probability distribution).**
- Of course, the node model can be adopted and **traffic rate equations are used to determine the total arrival rates of messages  $\Lambda_i$  at the different nodes**; we know the arrival rates  $\Lambda_i q_{ij}$  on the different links.
- The mean total delay  $T$  to cross the network can be derived by means of the Little theorem, as explained in the following slides.

# Kleinrock Application of Jackson Theorem (cont'd)

■ Let us denote:

- $\mu_k$  the mean completion rate for the k-th link
- $\alpha_k$  the mean arrival rate for the k-th link (if this link connects, let us say, node i to node j,  $\alpha_k = \Lambda_i q_{ij}$ ),
- $d_k$  the mean delay for the queue of the k-th link
- $\tau_k$  the propagation delay on the transmission line of the k-th link.

# Kleinrock Application of Jackson Theorem (cont'd)

- **Little theorem applied to the k-th link** (including the propagation delay in the mean delay on the link) to express the mean number of messages on this link:

$$\mathfrak{J}_k = \alpha_k(d_k + \tau_k)$$

- **Little theorem applied the whole network** to derive the mean (total, input-output) message delay T:

$$T = \frac{1}{\lambda_{tot}} \sum_{k=1}^L \alpha_k(d_k + \tau_k)$$

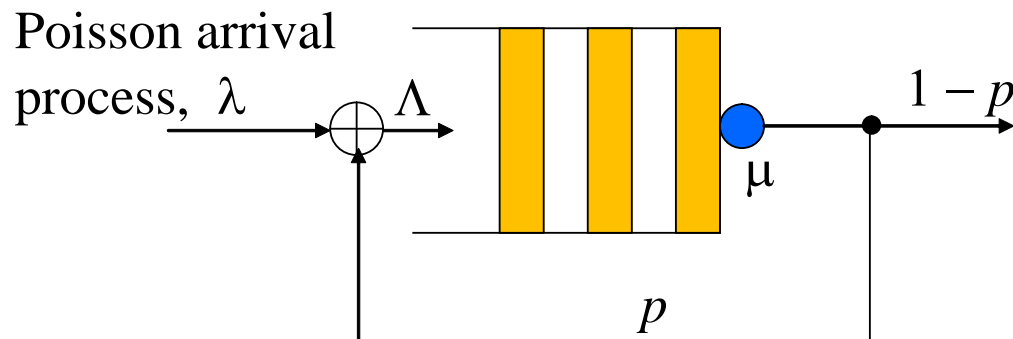
where  $d_k$  can be expressed by considering the M/M/1 characterization of the queue (Jackson theorem):

$$d_k = \frac{1}{\mu_k - \alpha_k}$$



# Analysis of a Queue with Feedback

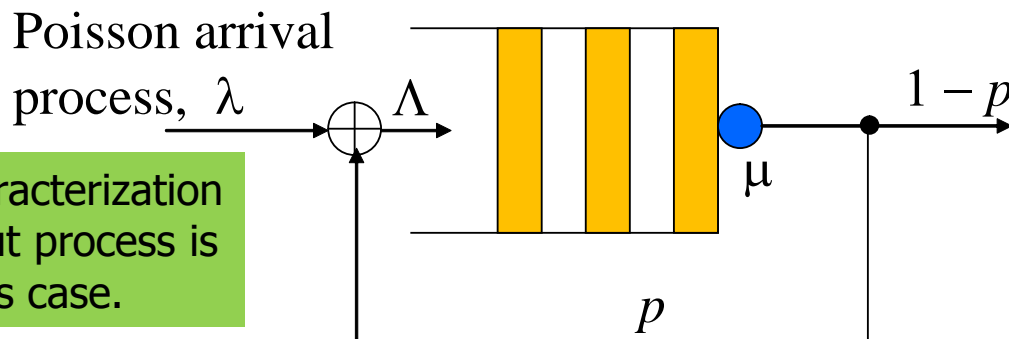
- We consider a special case for queuing networks: **a queue with one server where a request that completes its service can reenter the queue with probability  $p$  with no delay.** The arrival of messages from outside is according to a Poisson process with mean rate  $\lambda$ . The message service time is exponentially-distributed with mean rate  $\mu$ . The requests that complete the service have a form of stochastic routing according to which they may be fed back to the queue (**cyclic network**).



$$p = q_{ii}$$

# Analysis of a Queue with Feedback

- We consider a special case for queuing networks: **a queue with one server where a request that completes its service can reenter the queue with probability  $p$  with no delay.** The arrival of messages from outside is according to a Poisson process with mean rate  $\lambda$ . The message service time is exponentially-distributed with mean rate  $\mu$ . The requests that complete the service have a form of stochastic routing according to which they may be fed back to the queue (**cyclic network**).



The Poisson characterization for the total input process is also heavy in this case.

The Kleinrock assumption applied here is critical since the network is so small! There is a **strong correlation** between arrival process and service process.

We will solve this problem in **three** different ways using either the M/G/1 theory or the Jackson theorem with the Kleinrock assumption.

# Feedback Queue Studied with the Jackson Theorem

- We can apply the traffic rate equation to the system (= queue with stochastic feedback) to express the total mean arrival rate  $\Lambda$  (= mean output rate from the queue under stability assumption) as:

$$\Lambda = \lambda + \Lambda p \quad \Rightarrow \quad \Lambda = \frac{1}{1-p} \lambda$$

- Under the Kleinrock assumption (the service time of a message is exponentially distributed and independently regenerated each time the message is fed back to the queue), **we apply the Jackson theorem so that the queue admits an M/M/1 model.**
  - The queue is studied as if its input traffic was Poisson (however, the input traffic is not Poisson, but peaked, bursty)
  - The mean delay  $d$  experienced by a message entering the queue is (M/M/1 model):
$$d = \frac{1}{\mu - \Lambda}$$
Queue stability is assured if  $\Lambda/\mu < 1$  Erlang (ergodicity condition).

# Feedback Queue Studied with the Jackson Theorem

- From the Little theorem applied to the whole system we have:

$$\begin{aligned} T &= \frac{1}{\lambda} \times \Lambda d = \frac{1}{\lambda} \times \frac{1}{1-p} \lambda \times \frac{1}{\mu - \frac{1}{1-p} \lambda} = \frac{1}{1-p} \times \frac{1}{\mu - \frac{1}{1-p} \lambda} = \\ &= \frac{1}{1-p} \times \frac{1-p}{\mu(1-p) - \lambda} = \frac{1}{\mu(1-p) - \lambda} \end{aligned}$$

- This mean message delay  $T$  can be explained as follows:
  - A message entering the system from outside crosses the queue (due to the stochastic feedback) for a number of times with modified geometric distribution and mean value equal to  $1/(1-p)$ .
  - Each time the message goes through the queue it experiences a mean M/M/1 delay that is equal to  $(1-p)/[\mu(1-p) - \lambda]$ .
  - The product of the above terms yields the mean message delay  $T$ .

# Feedback Queue Scheduling

## with t

This is again an **M/M/1 mean delay term**, where each request (message) has a service time with exponential distribution and mean rate  $\mu(1-p)$  and that the arrival process is Poisson with mean rate  $\lambda$ .

- From the L
- Let us recall (see Section 4.3.2.2 of the book) that the **composition** of exponential (mean rate  $\mu$ ) and modified geometric (parameter  $1-p$ ) random variables is still exponentially distributed with mean rate  $\mu(1-p)$ .

$T =$

$$= \frac{1}{1-p} \times \frac{1-p}{\mu(1-p)-\lambda} = \frac{1}{\mu(1-p)-\lambda}$$

- This mean message delay  $T$  can be explained as follows:
  - A message entering the system from outside crosses the queue (due to the stochastic feedback) for a number of times with modified geometric distribution and mean value equal to  $1/(1-p)$ .
  - Each time the message goes through the queue it experiences a mean M/M/1 delay that is equal to  $(1-p)/[\mu(1-p)-\lambda]$ .
  - The product of the above terms yields the mean message delay  $T$ .

# Feedback Queue Studied with M/G/1 Theory

- After a message transmission, the message is instantaneously fed back to the queue with probability  $p$ . We can consider as if the message was put again at the head of the queue, since this does not alter the mean message delay: under the **insensitivity property**, different service disciplines yield the same mean message delay.
- We can determine the mean message delay as an application of the M/G/1 theory, imbedding the study at the instants when messages leave the system. We can use the Pollaczek-Kinchin formula as:

$$T = E[Y] + \frac{\lambda E[Y^2]}{2(1 - \lambda E[Y])}$$

where **Y denotes the total (equivalent) 'message service time'**, characterized as detailed in the next slide.

# Feedback Queue Studied with M/G/1 Theory (cont'd)

- Due to the feedback, the same message is transmitted  $N$  times before it leaves the queuing system. Let  $X_i$  denote the service time of a message at its  $i$ -th pass through the queue. Then, the equivalent service time  $Y$  of a message is obtained as follows:  $Y = \sum_{i=1}^N X_i$
- In a real system, **we could expect that the service time of a message is the same at each pass through the queue.** Hence,  $X_i \equiv X$  and  **$Y = N \times X$** . Considering that  $N$  and  $X$  are independent random variables, we can easily prove that
  - $E[Y] = E[n] \times E[X] = 1/[\mu(1 - p)]$
  - $E[Y^2] = E[n^2] \times E[X^2] = 2(1 + p)/[\mu^2(1 - p)^2]$
- Therefore, applying the Pollaczek-Kinchin formula, **we obtain with this approach an exact result for the mean message delay  $T$  as:**

$$T = \frac{1 + \frac{\lambda p}{\mu(1-p)}}{\mu(1-p) - \lambda}$$

# Feedback Queue and M/G/1 Theory + Kleinrock Assumpt.

- Instead, **using the Kleinrock assumption**, the message service time is 'restarted' at each pass through the queue so that in  $Y = \sum_{i=1}^N X_i$   $X_i$  are iid, exponentially distributed with mean rate  $\mu$  and  $N$  has a modified geometric distribution with parameter  $(1-p)$ .

- **Y is now given by the composition of an exponential distribution and a modified geometric distribution; hence, Y is exponentially distributed with mean rate  $\mu(1-p)$ .** Therefore, the Pollaczek-Kinchin formula simplifies, because the whole system behaves as an **M/M/1 queue**. The mean message delay  $T$  becomes:


$$T = \frac{1}{\mu(1-p) - \lambda}$$

- It is quite interesting to note that this is the same result obtained by applying the Jackson theorem with the Kleinrock assumption. **This results is approximated.**



# Final Considerations on Feedback Queue Analysis

- We can thus estimate the approximation entailed by the Kleinrock assumption in this case:


$$T = \frac{1 + \frac{\lambda p}{\mu(1-p)}}{\mu(1-p) - \lambda}$$

**M/G/1 approach  
without the Kleinrock assumption**

$$T = \frac{1}{\mu(1-p) - \lambda}$$

**M/G/1 approach or  
Jackson theorem  
with the Kleinrock assumption**

**Of course, the stability limits are the same in both cases, but the mean delays are not the same.**

# Network Planning Issues

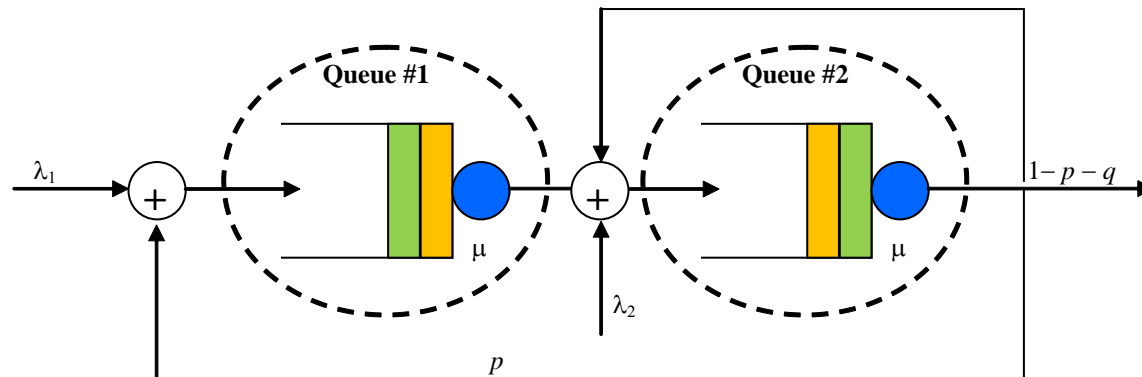
- Network planning and dimensioning with QoS support is a process involving the following steps:
  - Identification of network node location;
  - Definition of the link topology;
  - Adoption of a routing strategy accounting for external input traffics;
  - Capacity allocation to the links so that suitable QoS metrics (end-to-end delay, jitter, and packet loss rate) are fulfilled.
- These steps are interrelated.
  - Capacity allocation to links depends on the traffic loads on the links and, then, on traffic routing. However, traffic routing can also be adapted to account for traffic bottlenecks, which result from capacity shortage on some links.
- Network planning is a very complex optimization process and **the analysis carried out here provides a useful tool to allocate the capacity to links in the network once nodes, input traffic and routing are defined.**



# **Mixed Exercises on the Last Part of the Course**

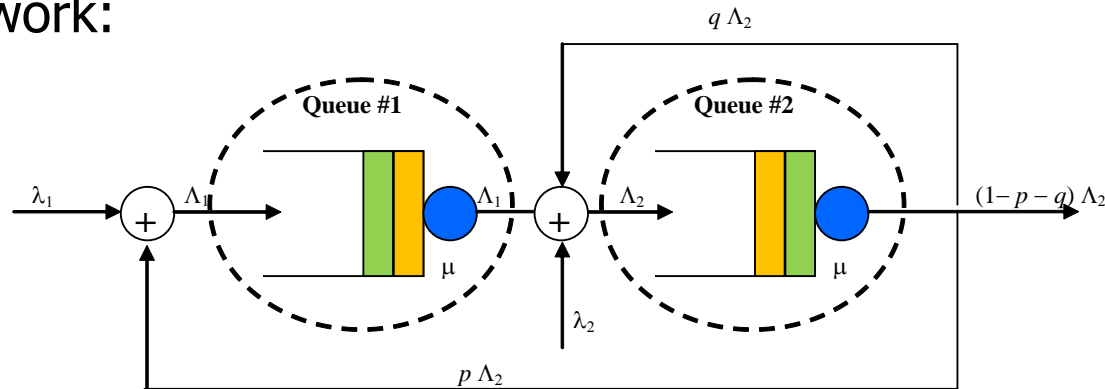
# Exercise #1

- With reference to the queuing network below, we have to determine the stability conditions for the different queues and the mean delay experienced by a message from input to output, considering:
  - Input traffic flows at the different queues from outside are Poisson independent with mean rates  $\lambda_1$  and  $\lambda_2$  for queues #1 and #2, respectively.
  - The message service times are independent for the two queues and exponentially distributed with the same mean rate  $\mu$  (Kleinrock assumption).
  - Queues have an infinite capacity.
  - At the output of queue #2 there is a random splitting: with probability  $p$  ( $q$ ) the arriving message is fed back to queue #1 (queue #2).



# Solution of Exercise #1

- Let  $\Lambda_1$  and  $\Lambda_2$  denote the mean total arrival rates for queues #1 and #2, respectively. We have the situation below for the mean rates in the network:



- Arrival rates  $\Lambda_1$  and  $\Lambda_2$  can be determined by writing **traffic rate equations for each sum point in the network**:

$$\begin{cases} \Lambda_1 = \lambda_1 + p\Lambda_2 \\ \Lambda_2 = \lambda_2 + \Lambda_1 + q\Lambda_2 \end{cases} \Rightarrow \begin{cases} \Lambda_1 = \frac{\lambda_2 p + (1-p)\lambda_1}{1-p-q} \\ \Lambda_2 = \frac{\lambda_1 + \lambda_2}{1-p-q} \end{cases}$$

# Solution (cont'd)

- We apply the **Kleinrock assumption**. Then, **the conditions of the Jackson theorem are fulfilled for our network**: queue #1 can be studied by means of an M/M/1 model with mean arrival rate  $\Lambda_1$  and queue #2 can be studied by means of an M/M/1 model with mean arrival rate  $\Lambda_2$ .
- Queues #1 and #2 are stable under the following conditions:  $\rho_1 = \Lambda_1/\mu < 1$  Erlang and  $\rho_2 = \Lambda_2/\mu < 1$  Erlang.
- The mean number of messages in queues #1 and #2 can be obtained as functions of  $\rho_1$  and  $\rho_2$  as:  $N_1 = \frac{\rho_1}{1 - \rho_1}$  ,  $N_2 = \frac{\rho_2}{1 - \rho_2}$
- The mean message delay from input to output,  $T$ , can be obtained by applying the **Little theorem** to the whole system:

$$T = \frac{N_1 + N_2}{\lambda_1 + \lambda_2}$$

# Exercise #2

- Let us consider an **FTP file transfer** that is based on **TCP Tahoe**. We are requested to plot the congestion window (cwnd) behavior has a function of time [expressed in RTT units] until 16 RTTs, under the following conditions:
  - Bottleneck buffer size  $B = 15$  pkts
  - Sockets buffers much larger than  $B + BDP$
  - Bandwidth-Delay Product  $BDP = 15$  pkts
  - Initial ssthresh value = 16 packets
  - All the packets of a cwnd are transmitted altogether and their ACKs are received altogether in an RTT time (model).

How many packets have been transmitted until time = 5 RTTs ?

# Exercise #2

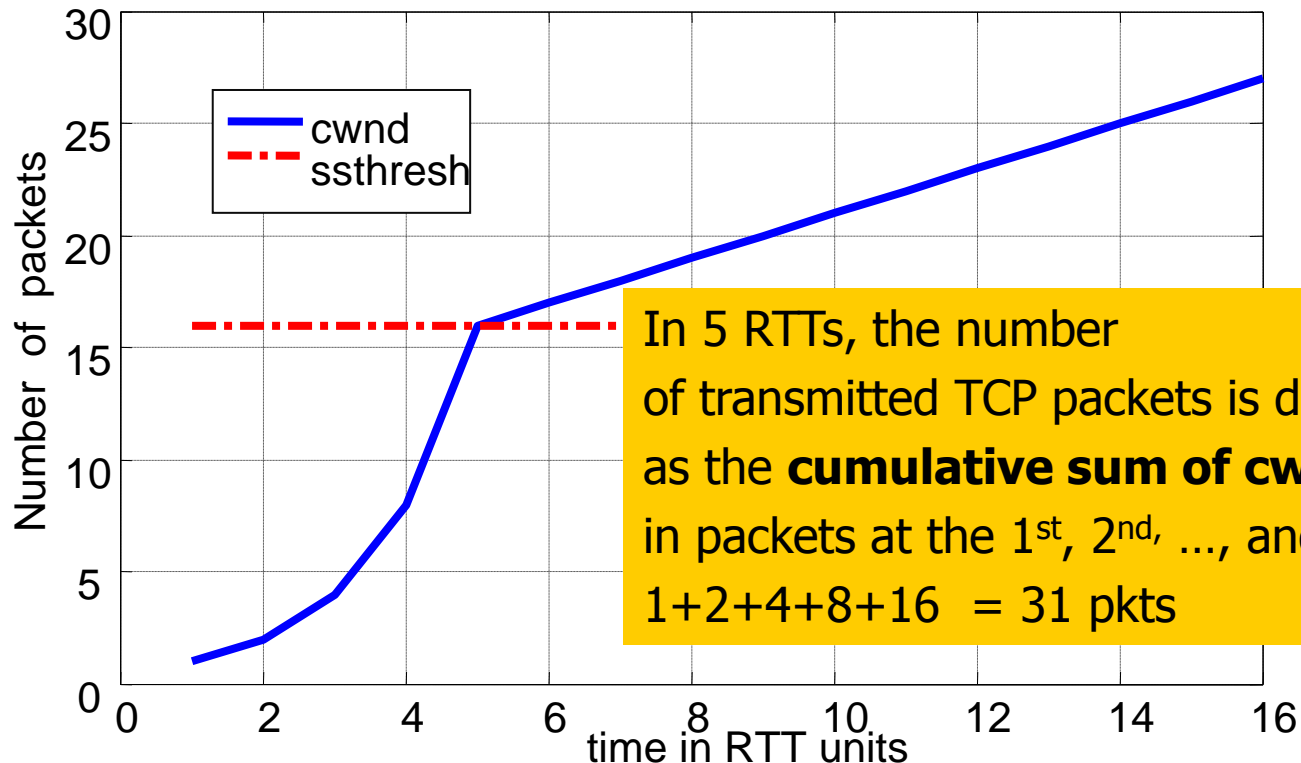
- Let us consider an **FTP file transfer** using **TCP Tahoe**. We are requesting a congestion window (cwnd) behavior that permits to fully exploit the capacity of the bottleneck link. [expressed in RTT units] under the following conditions:
- B = BDP is the optimal setting for the bottleneck link buffer that permits to fully exploit the capacity of the bottleneck link.
- Bottleneck buffer size  $B = 15$  pkts
  - Sockets buffers much larger than  $B + BDP$
  - Bandwidth-Delay Product  $BDP = 15$  pkts
  - Initial ssthresh value = 16 packets (this is not the default value)
  - All the packets of a cwnd are transmitted altogether and their ACKs are received altogether in an RTT time (model).

How many packets have been transmitted until time = 5 RTTs ?



# Solution of Exercise #2

- The cwnd behavior has first a slow start phase with exponential increase and after (i.e., when  $cwnd > ssthresh$ ) a congestion avoidance phase with linear behavior. There is no cwnd drop event in the interval of observation since the maximum allowed cwnd value is  $B + DBP = 30$  pkts. We have the same behavior of cwnd in this initial phase for both TCP Tahoe and TCP NewReno.

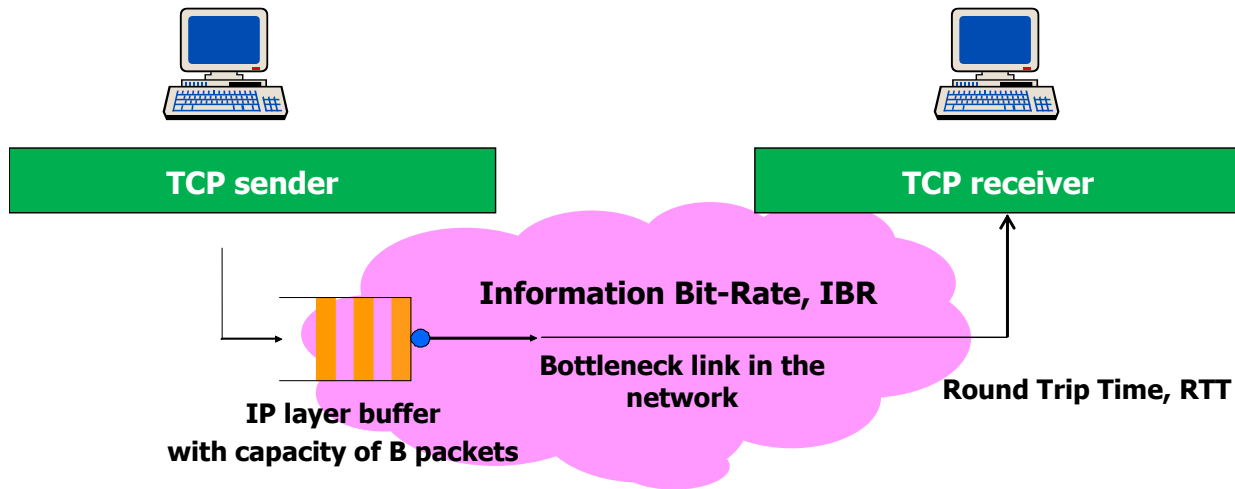


In 5 RTTs, the number of transmitted TCP packets is determined as the **cumulative sum of cwnd values** in packets at the 1<sup>st</sup>, 2<sup>nd</sup>, ..., and 5<sup>th</sup> RTT as  $1 + 2 + 4 + 8 + 16 = 31$  pkts

# Exercise #3

- Let us consider an FTP data transfer (TCP 'elephant' flow), referring to the network model in the next Figure. We adopt a scenario with IP packets (MTU) of 1500 bytes, with Information Bit-Rate (IBR) of the bottleneck link equal to 600 kbit/s, and with physical Round Trip Time (RTT) equal to 0.5 s (GEO satellite scenario). It is requested to derive the Bandwidth-Delay Product (BDP) and to plot the behaviors of both congestion window (cwnd) and slow start threshold (ssthresh) up to the time of 25 RTTs for both **TCP Tahoe** and **TCP NewReno**, under the following conditions:
  - Bottleneck link buffer capacity  $B = 20$  pkts;
  - Sockets buffers much bigger than  $B + \text{BDP}$ ;
  - Initial ssthresh value equal to 32 pkts;
  - All the packets of a cwnd are transmitted altogether and their ACKs are received altogether in an RTT time (model).
- It is requested to redo the exercise with initial ssthresh = 64 pkts.

# Solution of Exercise #3



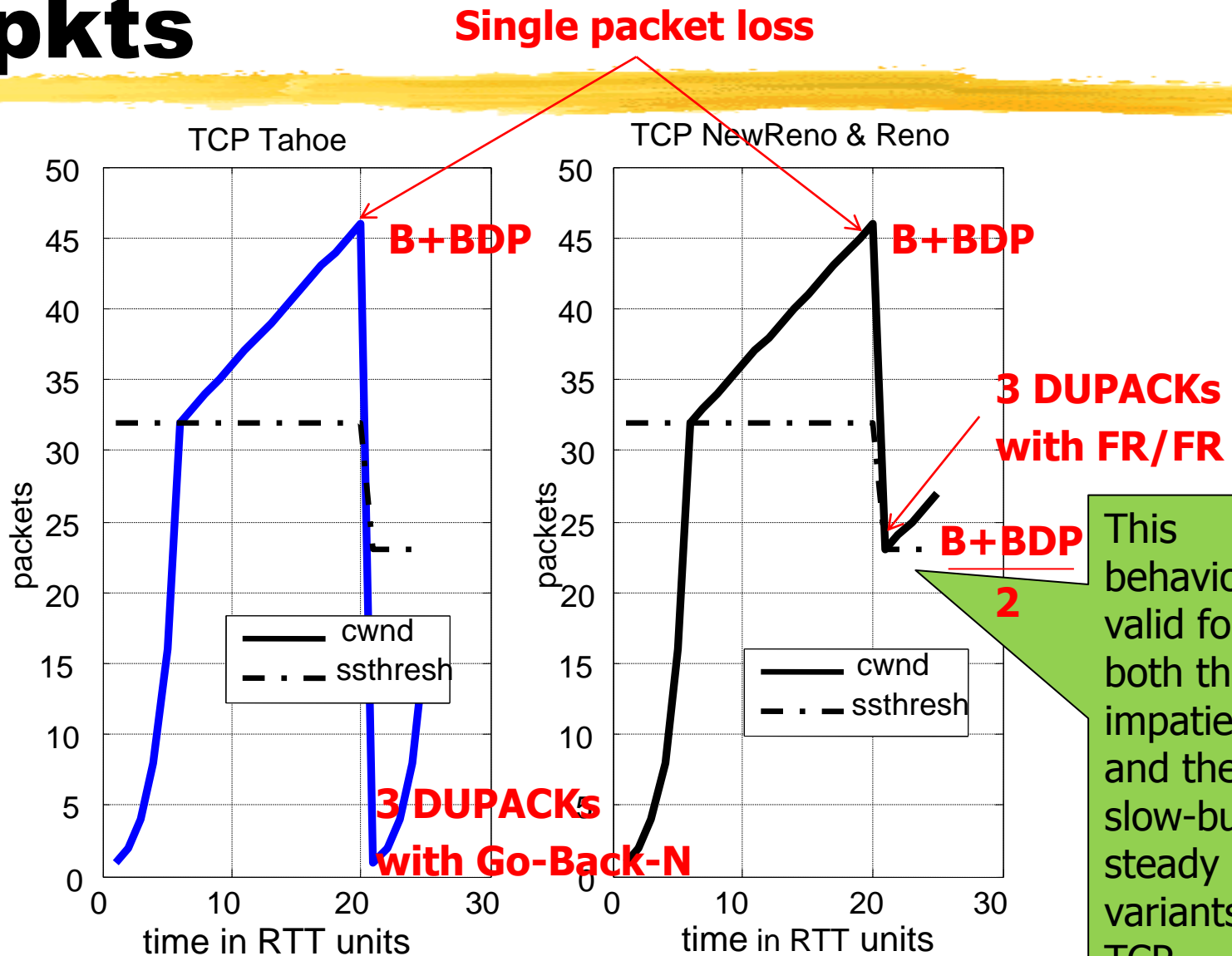
- The BDP for the data transfer in this exercise results as:

$$BDP = \frac{RTT \times IBR}{MTU} = 25 \text{ pkt}$$

RTT is here approximated by RTD.

- cwnd reaches the maximum value of  $B + BDP = 45$  pkts

# Solution for Initial ssthresh = 32 pkts

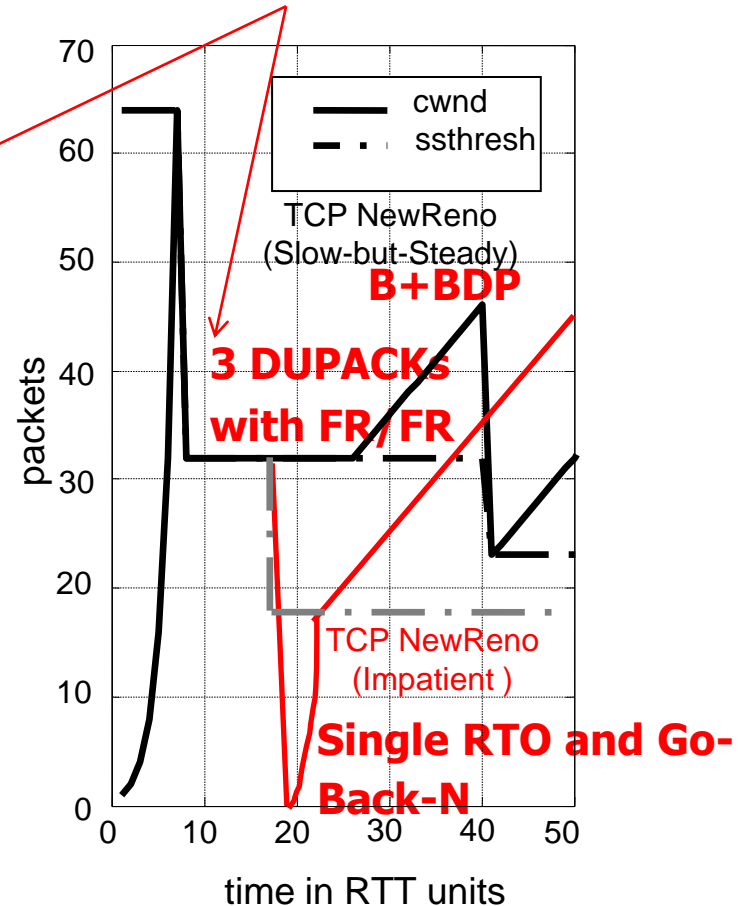
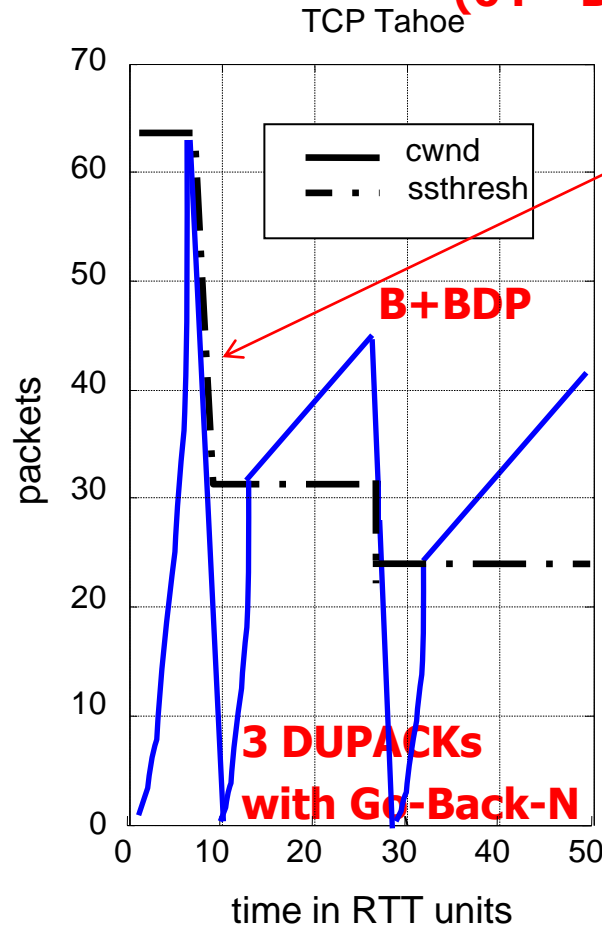


This behavior is valid for both the impatient and the slow-but-steady variants of TCP NewReno

# Solution for Initial ssthresh = 64 pkts

Multiple packet losses in a window of data  
( $64 - B - BDP = 19$  pkts)

The initial ssthresh value of 32 pkts is better than 64 pkts in terms of delivered packets as a function of time.

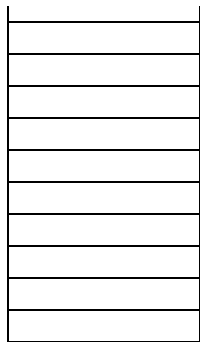


# Exercise #4

- Let us consider an IP access network using IntServ as QoS support method. In particular the **Guaranteed Service** is adopted. Let us consider that a traffic source (with **fluid flow model**) accessing the network is regulated according to the following token bucket **T-Spec parameters**  $(r, p, b) = (1 \text{ kbit/s}, 5 \text{ kbit/s}, 400 \text{ bits})$  [1 token = 1 bit].
- Considering the approach with arrival curve, service curve, and departure curve, we have **to determine the minimum service rate  $R$  to guarantee a delay lower than or equal to  $\Delta_{\max} = 100 \text{ ms}$**  (we refer to a case where the propagation delay is negligible with respect to  $\Delta_{\max}$ ).

# Solution of Exercise #4

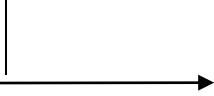
Tokens enter the bucket  
at **rate  $r$**



Bucket **depth  $b$** : capacity of the bucket

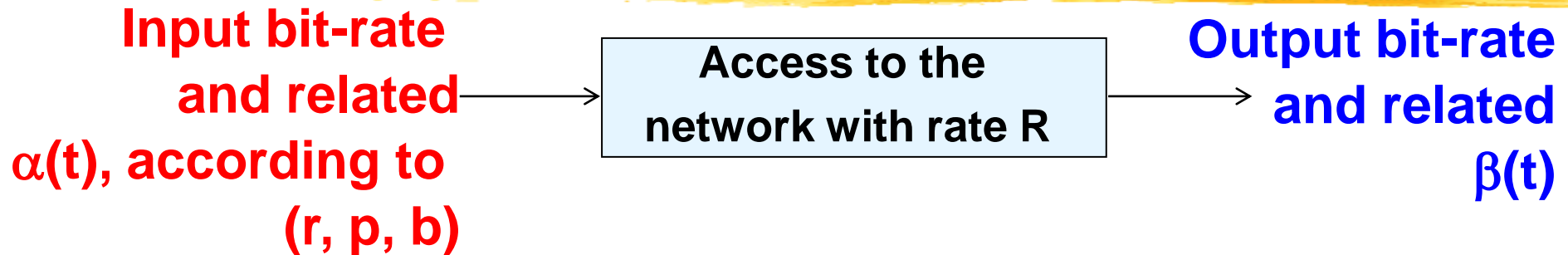
**Max allowed transmission with rate  $p$**

**Traffic source**

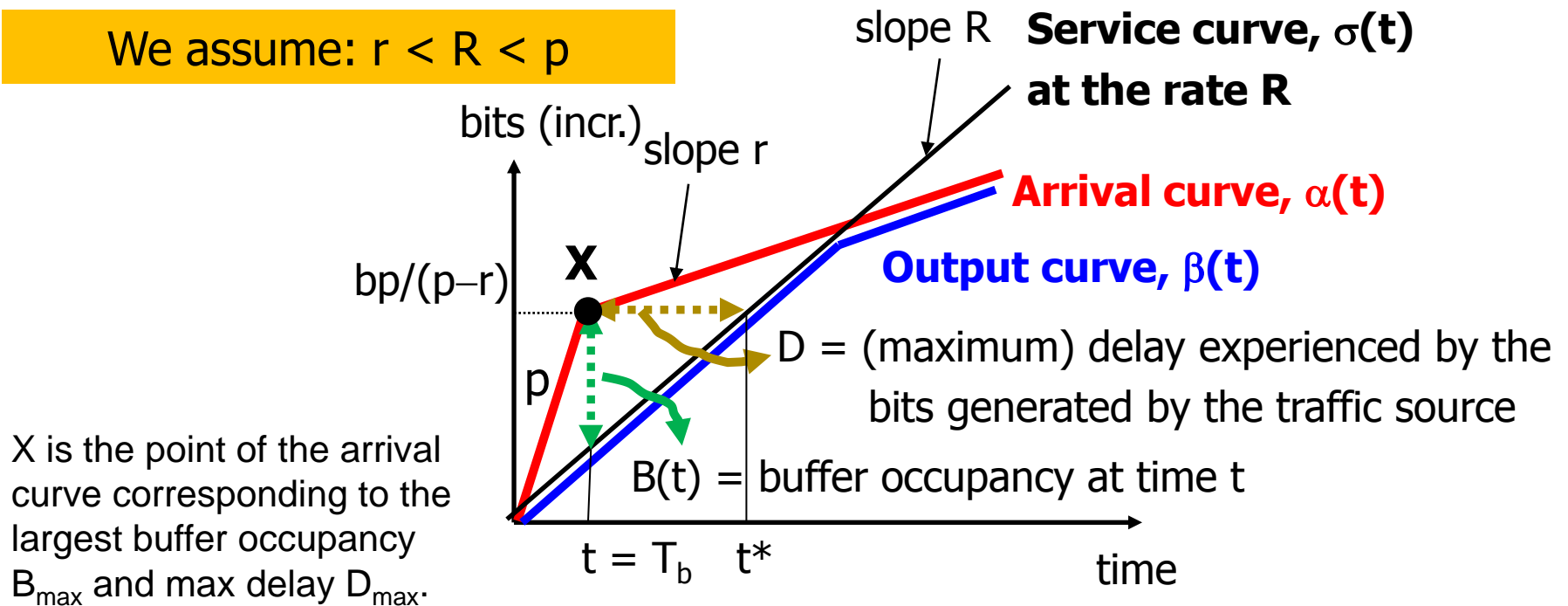


- In this study we consider a fluid-flow model for the traffic generated by the source: 1 token is needed for the transmission of 1 bit (no packets); if the bucket contains  $N$  tokens,  $N$  bits can be sent at maximum rate  $p$ .
- If the bucket is full, new tokens are discarded.
- The input traffic to the network has a resulting bit-rate with corresponding arrival curve  $\alpha(t)$ .

# Solution (cont'd)



We assume:  $r < R < p$





# Solution (cont'd)

- It has been proved that the delay  $D$  to cross the node (modeling the access network) is bounded as  $D \leq b/R$ . Let us consider the condition with equality  $D \approx b/R$ . Then, we adopt the following formula to determine  $R$ :

$$D = \frac{b}{R} \leq \Delta_{\max} \quad \Rightarrow \quad R \geq \frac{b}{\Delta_{\max}}$$

**Then we select the minimum value for  $R$  to fulfill  $\Delta_{\max}$ , that is  $R = b/\Delta_{\max}$ .**

- Moreover, we consider that  $R$  has to fulfill the following constraint:

$$r = 1 \frac{\text{kbit}}{\text{s}} < R < p = 5 \frac{\text{kbit}}{\text{s}}$$

- So  $R = b/\Delta_{\max} = 4 \text{ kbit/s}$  fulfills the constraint and is the minimum  $R$  value to guarantee a delay lower than  $\Delta_{\max}$ . There is some approximation in this, but we consider that this is acceptable.

# Solution (cont'd)

- For the sake of completeness, let us recall that the system is characterized by bounded delay ( $D_{\max}$ ) and bounded buffer size (maximum buffer occupancy  $B_{\max}$ ) determined as follows (exact formulas):

$$D_{\max} = t^* - T_b = \frac{b}{R} \times \left( \frac{p - R}{p - r} \right) \leq \frac{b}{R}, \quad \text{if } R \geq r$$

$$B_{\max} = pT_b - RT_b = b \times \left( \frac{p - R}{p - r} \right) \leq b, \quad \text{if } R \geq r$$



**Thank you!**

**[giovanni.giambene@gmail.com](mailto:giovanni.giambene@gmail.com)**