



Avaliação de desempenho de Sistemas de Informação  
Teoria das Filas

**BACHARELADO EM SISTEMAS DE INFORMAÇÃO**

Prof. Sergio Nascimento

[sergio.onascimento@sp.senac.br](mailto:sergio.onascimento@sp.senac.br)

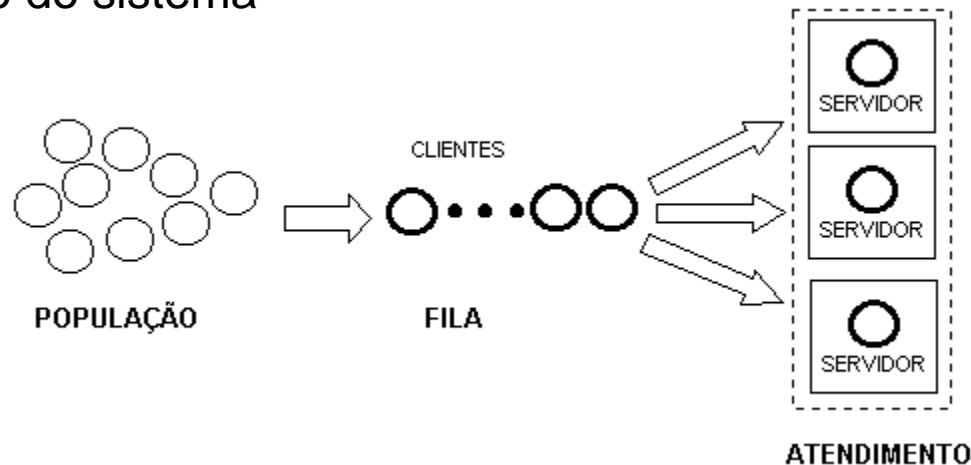


# Filas

Todos já passaram pelo aborrecimento de ter que esperar em uma fila:  
Fila de ônibus, banco, restaurante, trânsito, etc.



Em sistemas computacionais há filas em vários lugares:  
Filas de impressão, pacotes em roteadores, processos aguardando processamento da CPU, etc.  
As filas surgem porque a demanda de serviço é maior que a capacidade de atendimento do sistema



# Teoria das Filas

- ❑ Ramo da probabilidade que estuda o fenômeno da formação de filas em solicitações de serviços.
- ❑ Permite criar modelos do sistema estudado para prever o seu comportamento.
  - Pode-se dimensionar um determinado sistema segundo a demanda de seus clientes, evitando desperdícios ou gargalos.
- ❑ Utilizada para modelar sistemas onde:
  - Clientes chegam para ser atendidos.
  - Esperam a sua vez de ser atendidos.
  - São atendidos e vão embora.

# Sistemas de Fluxo

- ❑ Sistema no qual alguma “comodity” flui, se move ou é transferida por meio de um ou mais canais com capacidade limitada, de forma a ir de um ponto a outro no sistema.
- ❑ Sistemas transferem as commodities à uma taxa finita.
- ❑ Classificação:
  - Estáveis (Determinísticos)
  - Instáveis (Estocástico)
- ❑ A **Fila** é gerada pelos clientes à espera de atendimento (Não inclui o(s) cliente(s) em atendimento).
- ❑ **Serviço** ou atendimento é constituído por um ou mais postos de atendimento.

**Fila + Serviço = Sistema**

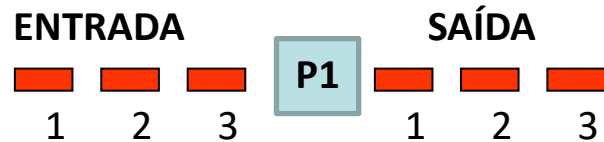
# Sistemas de Fluxo – Estáveis (Determinísticos)

Sistemas onde o fluxo passa de forma previsível.

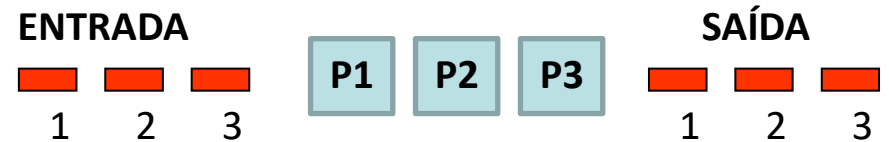
- ❑ Quantidade de fluxo pelo canal é constante e conhecida durante o período de observação.
- ❑ **Taxa de Serviço (C) e Taxa de Chegadas (R)**
  - Se  $C > R$  não existe sobrecarga.
    - Atendido por Sistemas Estáveis de Canal Único (*single channel*)
  - Se  $C < R$  existe congestionamento (Fila).
    - Atendido por Sistemas Estáveis com Redes de Canais (*multiple channel*).
- ❑ Número de clientes no sistema (Em cada instante) = Estado do sistema.

# Modelos de Filas

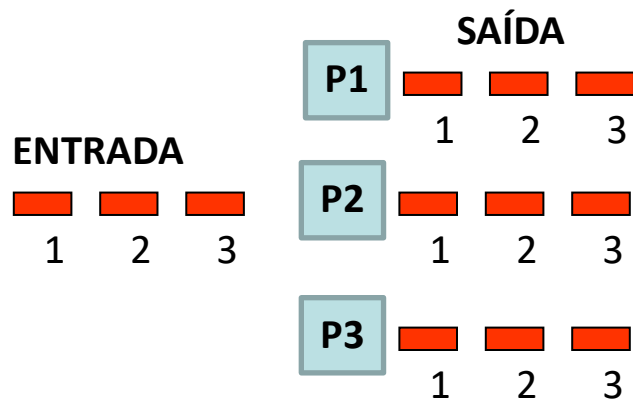
## SINGLE-CHANNEL/SINGLE PHASE



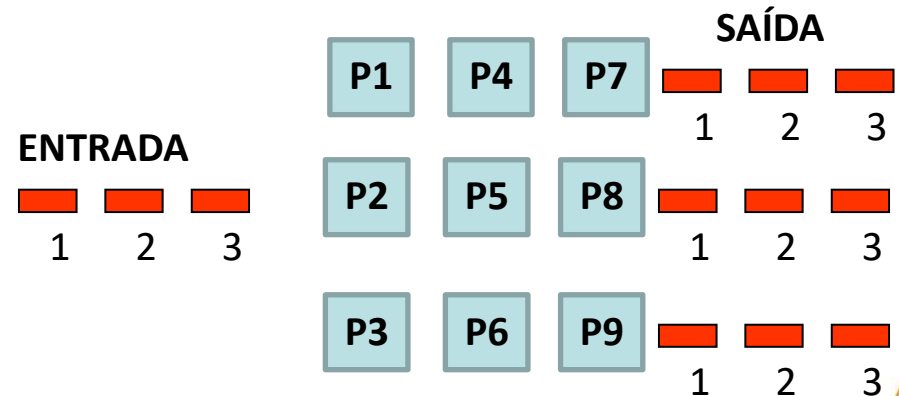
## SINGLE-CHANNEL/MULTI PHASE



## MULTI-CHANNEL/SINGLE PHASE



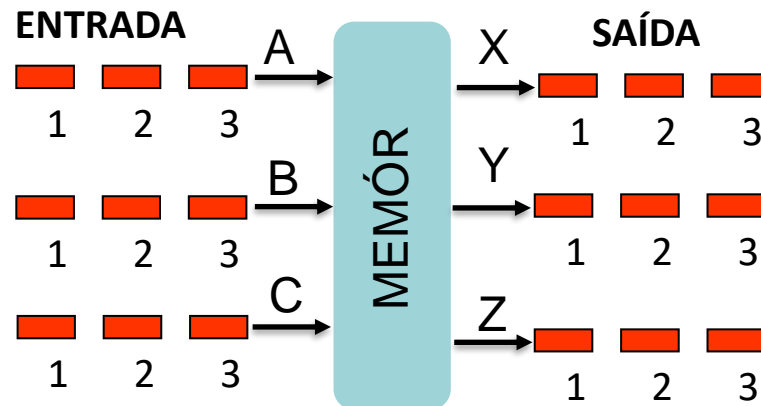
## MULTI-CHANNEL/MULTI PHASE



# Sistemas de Fluxo – Instáveis (Estocástico)

Sistemas onde o fluxo passa de forma aleatória (chegada é aleatória).

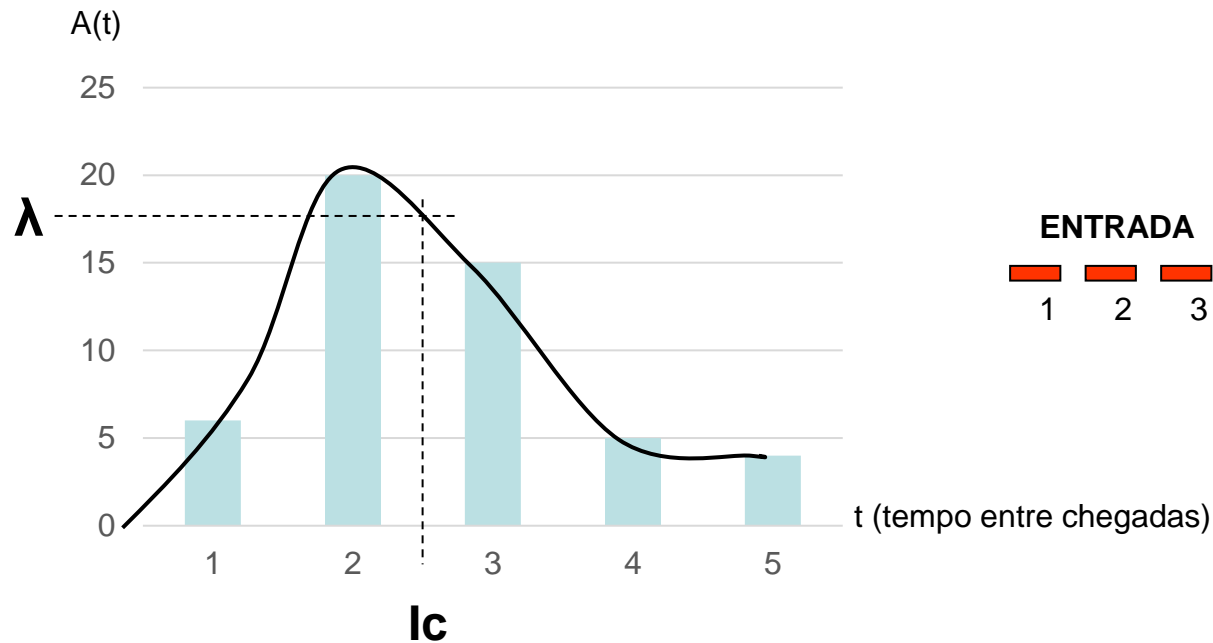
- ❑ Demanda de uso do canal é imprevisível.
  - Exemplo: internet, telefonia IP, etc.
- ❑ Como os sistemas estáveis também pode apresentar canal único ou rede de canais.
- ❑ A Teoria da Filas se propõe a resolver problemas de sistemas de fluxo aleatório (estocásticos).



# Processos Estocásticos

## ❑ Taxa de Chegada $[A(t)]$

- $A(t) = P$  [tempo entre chegadas  $\leq t$ ] Probabilidade que o tempo entre chegadas seja menor ou igual a  $t$ , onde  $0 < A(t) < 1$ .
- Ritmo médio de chegada:  $\lambda$
- Intervalo médio entre chegadas:  $I_c$

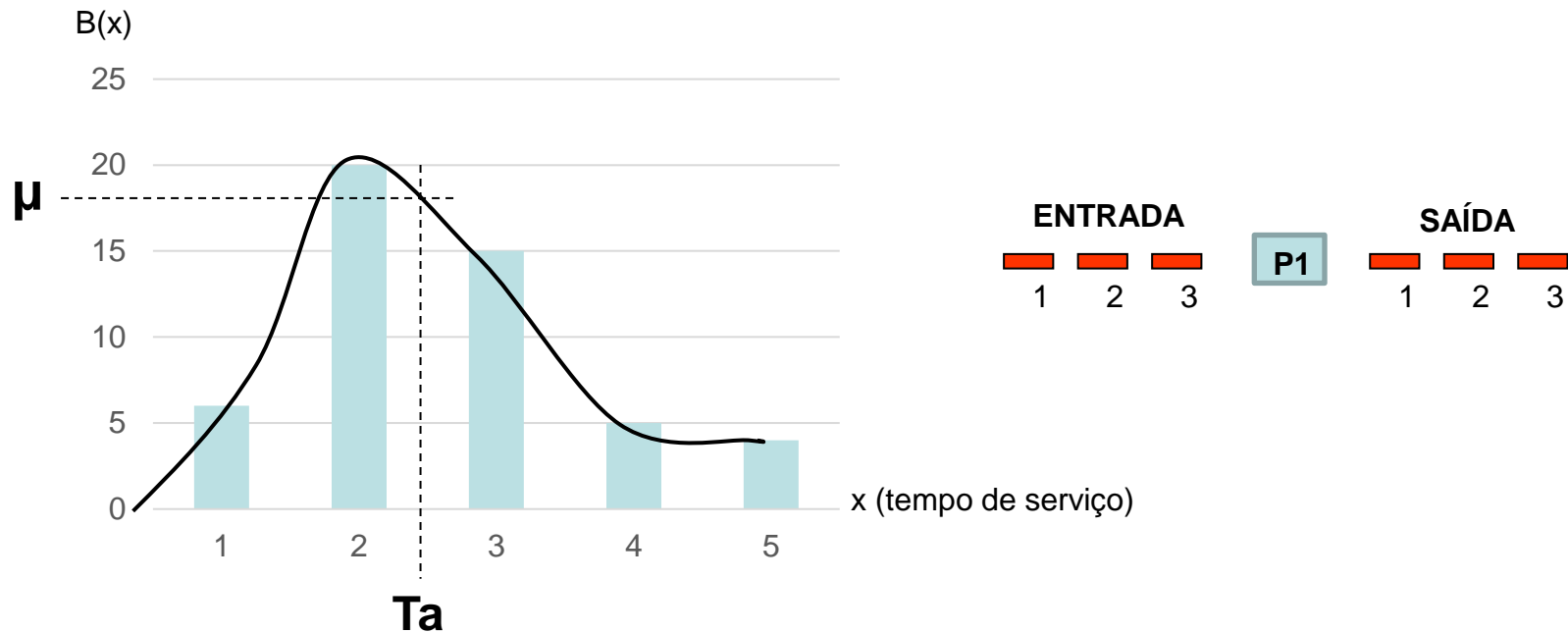




# Processos Estocásticos

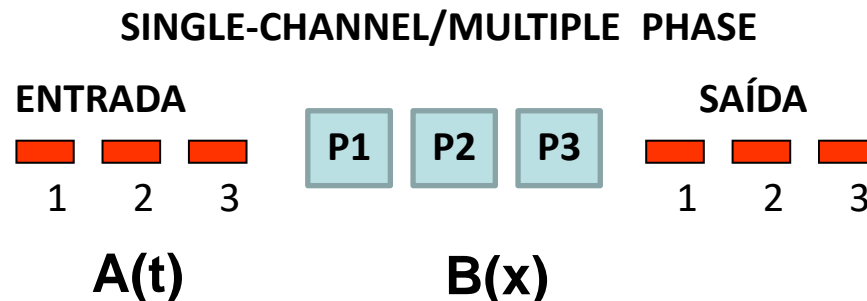
## ❑ Taxa de Serviço [B(x)]

- $B(x) = P [ \text{tempo de serviço} \leq x ]$  Probabilidade que o tempo de serviço seja menor ou igual a x.
- Ritmo médio de atendimento:  $\mu$
- Tempo médio de atendimento:  $T_a$



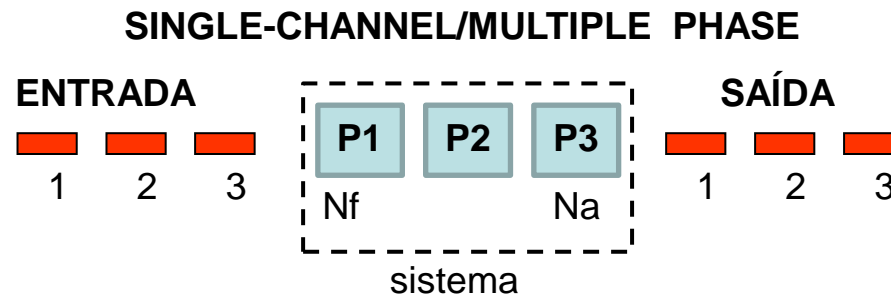
# Sistema de Filas - Caracterização

- ❑ Linha de entrada – chegam clientes (solicitações, pacotes, jobs, etc.) com tempo de chegada aleatório [ **A (t)** ].
- ❑ Clientes se encaminham para o sistema de filas, onde são atendidos (servidos) por um tempo aleatório [ **B (x)** ].
- ❑ Depois de atendidos são liberados em uma linha de saída. Sistemas podem ter  $m$  servidores (P1, P2, P3).



# Variáveis

- Variáveis referentes ao sistema
  - $T_s$  = Tempo médio de permanência no sistema.
  - $N_s$  = Número médio de clientes no sistema.
- Variáveis referentes ao processo de chegada
  - **$\lambda$  = Rítmo médio de chegada**
  - $I_c$  = Intervalo médio entre chegadas =  $1/\lambda$
  - $N_c$  = Número médio de chegadas no sistema.
- Variáveis referentes à fila
  - $T_f$  = Tempo médio de permanência na fila.
  - $N_f$  = Número médio de clientes na fila.
- Variáveis referentes ao processo de atendimento
  - **$\mu$  = Rítmo médio de atendimento**
  - $T_a$  = Tempo médio de atendimento ou serviço =  $1/\mu$
  - $N_a$  = Número médio de clientes sendo atendidos
  - $c$  = Capacidade de atendimento



# Variáveis de chegada

Ritmo médio de chegada ( $\lambda_t$ ) =  $\frac{N_c}{I_c}$  onde

( $\lambda_t$ ) Ritmo médio de chegada de clientes por unidade de tempo

- $I_c$  = Intervalo médio entre chegadas

Ex: Um sistema recebe 8 chegadas de clientes em um período de 4 minutos.  
Qual o Ritmo médio de chegada?



$$\lambda_t = \frac{N_c}{I_c} = \frac{8}{4} = 2 \text{ chegadas/min}$$

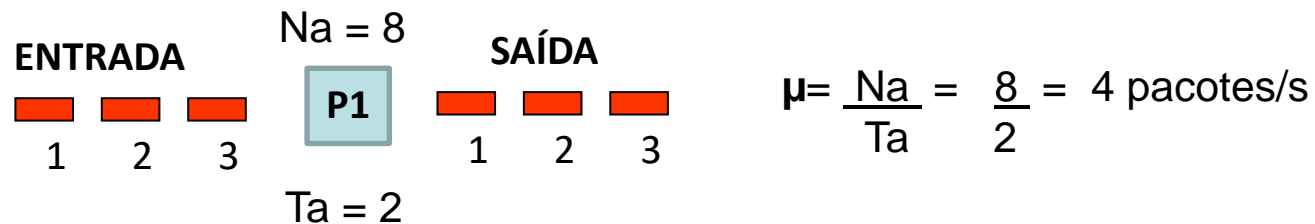
# Variáveis de serviço

**Ritmo médio de Serviço** ( $\mu_t$ ) =  $\frac{N_a}{T_a}$  onde

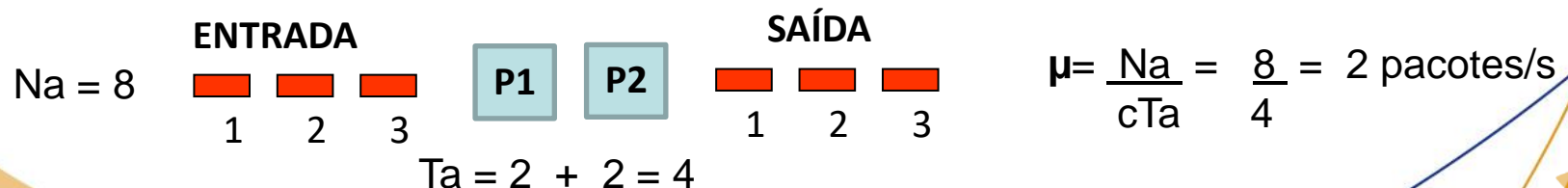
( $\mu_t$ ) Ritmo médio de serviço por cliente no servidor.

- $T_a$  = Tempo médio de atendimento ou serviço.
- Para  $c$  servidores (em uma mesma fila), a taxa total é  $c \mu$

Ex: Cada sistema demora 2 segundos para processar 8 pacotes IP. Qual o Ritmo médio de serviço?



Ex: Para dois sistemas a taxa seria de:



# Variáveis – Relações Básicas

Número de usuário no sistema:  $N_s = N_f + N_a$

Tempo de permanência no sistema:  $T_s = T_f + T_a$

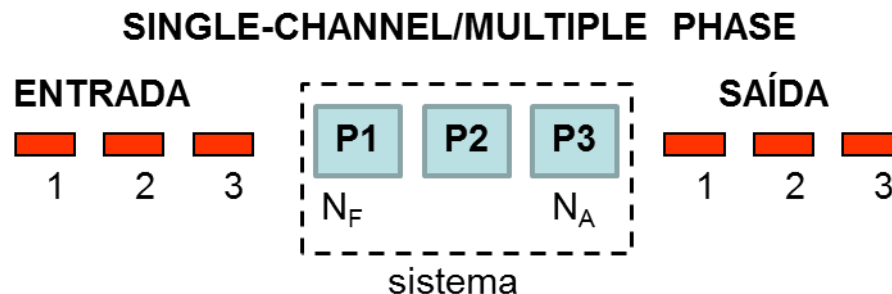
Pode-se também demonstrar que:

$$N_a = \lambda / \mu = T_a / l_c$$

Portanto:

$$N_s = N_f + N_a$$

$$N_s = N_f + (\lambda / \mu) = N_f + (T_a / l_c)$$



# Lei de Little

É um dos resultados mais importantes e amplamente usados em desempenho de qualquer sistema. O número médio de clientes ( $N$ ) em um sistema estável (em algum intervalo ( $T$ )), é igual a taxa média de chegada  $\lambda$ , multiplicada pelo tempo médio de serviço no sistema  $T$ , ou seja, J.D.C Little demonstrou que, para um sistema estável\* de filas, temos:

$$N_f = \lambda \times T_f$$

- Número médio de clientes na fila é igual à taxa de chegadas multiplicada pelo tempo médio de permanência na fila.

$$N_s = \lambda \times T_s$$

- Número médio de clientes no sistema é igual à taxa de chegadas multiplicada pelo tempo de permanência no sistema.

□ **Taxa de Utilização:**  $U(t) = t/T = \lambda / \mu$

Onde  $t$  = tempo de serviço e  $T$  = tempo de observação

- Fração de tempo no qual o servidor está ocupado
  - $U(t) = 0$  – sistema vazio (ocioso)
  - $U(t) > 0$  – sistema ocupado
  - $U(t) = 1$  – sistema completamente ocupado (saturado)

\* Sistema estável onde num intervalo grande de observação, o número de saídas é igual ao número de chegadas do sistema.

# Medidas e Relacionamentos entre as variáveis

## Condição de Estabilidade

$\lambda = c \cdot \mu$  ( $c \geq 1$ ) válido para filas infinitas

**Intensidade de Tráfego ( $\rho$ )** – mede o congestionamento de sistemas de filas

$$\rho = \frac{\text{taxa de chegadas}}{\text{taxa de serviço}} = |\lambda / \mu| = |T_a / T_c| \quad \text{ou} \quad \lambda / c \cdot \mu \quad ; \quad \text{onde } \rho \leq 1$$

**Vazão (Throughput)** – média das solicitações processadas por unidade de tempo (taxa de saída do sistema).

$$X = N/T$$

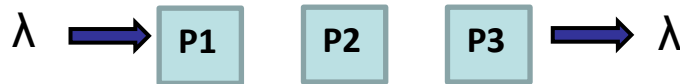
Equilíbrio: Taxa de saída é igual à taxa de chegada do sistema

$$\lambda = c \cdot \mu \cdot \rho$$

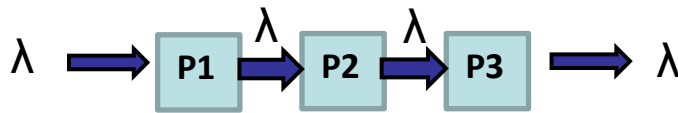
(Throughput para sistema equilibrado)



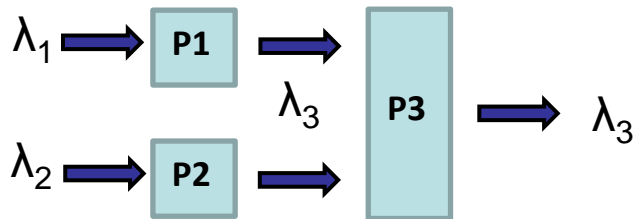
# Postulados Básicos



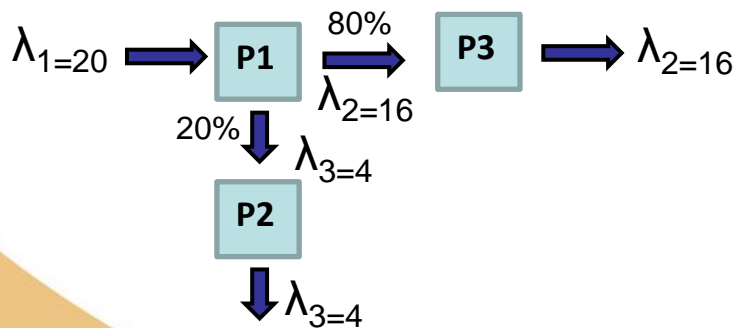
1) Em qualquer sistema estável, o fluxo que entra é igual ao fluxo que sai.



2) Em qualquer sistema estável, o fluxo de entrada se mantém nas diversas seções do sistema, desde que não exista junção ou desdobramento.



3) Em qualquer sistema estável, a junção de fluxos equivale às suas somas aritméticas, ou seja  $\lambda_3 = \lambda_1 + \lambda_2$ .



4) Em qualquer sistema estável, o desdobramento percentual de um fluxo é igual ao desdobramento aritmético do mesmo fluxo. Assim, se após a estação P1 temos 80% do fluxo se deslocando para P3, então o ritmo de chegada em P3 é de  $0,8 \times 20 = 16$  clientes/minuto

# Notação de Kendall

## **A / S / NS / B / K / SD**

- **A / S** = distribuição de tempo de chegada e tempo de serviço
  - **M** (Exponencial, Poisson), **E<sub>k</sub>** (Erlang), **H<sub>k</sub>** (Hiperexponencial), **D** (Determinístico), **G** (Geral – todas as distribuições)
- **NS** = Número de servidores
- **B** = Número de buffers (lugares na fila)
- **K** = Tamanho da população
- **SD** = Disciplina do serviço □ FIFO, LIFO, etc.

Padrão: **B** =  $\infty$  , **K** =  $\infty$  , **SD** = FIFO

- **M / M / 1 = M / M / 1 /  $\infty$  /  $\infty$  / FIFO** (chegadas Poisson, tempo de serviço exponencial, 1 servidor, buffer infinito, FIFO).
- **M / M / c** (idem anterior, com c servidores)

# Exercícios – Teoria das filas

1. Se um sistema apresenta um tempo de serviço de 3 segundos e atende 5 solicitações durante o intervalo de observação, por quanto tempo esse servidor esteve ocupado?
2. Observamos que 2000 solicitações de serviços chegaram a um determinado componente do sistema, durante um intervalo de 20 minutos. Qual a taxa de chegada de solicitações (em solicitações por segundo)?
3. Observamos 50 solicitações sendo atendidas em um intervalo de observação de 4 minutos. Qual foi o throughput observado (em solicitações por segundo)?
4. Durante um período de 3 horas foi observado um servidor. Ele apresenta uma utilização de 65%. Qual foi o tempo que ele esteve ocupado durante a observação?
5. Um servidor de arquivos está conectado à rede da empresa. Ele recebe 150 solicitações de arquivos e atende 45 dessas solicitações durante um intervalo de observação de 30 minutos. O servidor gasta 5 segundos para atender cada solicitação. Calcule os valores de  $\lambda$ ,  $X$  e  $U$  e o tempo que o servidor esteve ocupado.

# Exercícios – Teoria das filas

1. Se um sistema apresenta um tempo de serviço de 3 segundos e atende 5 solicitações durante o intervalo de observação, por quanto tempo esse servidor esteve ocupado?

$$T_a = 3 \text{ seg} ; N_a = 5 \rightarrow t = 3 \times 5 = 15 \text{ seg}$$

2. Observamos que 2000 solicitações de serviços chegaram a um determinado componente do sistema, durante um intervalo de 20 minutos. Qual a taxa de chegada de solicitações (em solicitações por segundo)?

$$N_c = 2000 ; T_c = 20 \text{ min} \Rightarrow \lambda = 2000/20/60 = 1,67 \text{ solicitações/seg.}$$

3. Observamos 50 solicitações sendo atendidas em um intervalo de observação de 4 minutos. Qual foi o throughput observado (em solicitações por segundo)?

$$N_a = 50 ; T = 4 \text{ min} \Rightarrow X = N_a/T = 50/4/60 = 0,21 \text{ solicitações/seg}$$

4. Durante um período de 3 horas foi observado um servidor. Ele apresenta uma utilização de 65%. Qual foi o tempo que ele esteve ocupado durante a observação?

$$T = 3 \text{ horas} ; U = 0,65 \Rightarrow U = t/T \Rightarrow t = T \times U = 3 \times 0,65 = 1,95 \text{ horas} = 1 \text{ h } 57 \text{ min}$$

## Exercícios – Teoria das filas

5. Um servidor de arquivos está conectado à rede da empresa. Ele recebe 150 solicitações de arquivos e atende 45 dessas solicitações durante um intervalo de observação de 30 minutos. O servidor gasta 5 segundos para atender cada solicitação. Calcule os valores de  $\lambda$ ,  $\mu$ ,  $U$  e o tempo que o servidor esteve ocupado.

$$N_c = 150 \text{ solicitações ; } T = 30 \text{ min}$$

$$N_a = 45 ; T_a = 5 \text{ seg}$$

$$\lambda = N_c/T = 150/30 = 5 \text{ solicitações/min}$$

$$\mu = N/T = 45/30 \times 60 = 0,025 \text{ solicitações/seg}$$

$$U = t/T$$

$$t = 45 \times 5 = 225 \text{ seg}$$

$$T = 30 \times 60 = 1800$$

$$U = 225/1800 = 0,125 \text{ OU } 12,5 \text{ \%}.$$

# Exercícios – Teoria das filas

6. Em um determinado sistema observamos que 10 solicitações foram atendidas durante o tempo de observação e que esse servidor esteve ocupado por 200 segundos durante esse mesmo período de observação. Qual o tempo médio de serviço observado desse sistema?

7. Um servidor apresenta uma taxa de chegada de 10 solicitações/segundo. Se observarmos o servidor por 10 minutos, quantas solicitações de serviço chegaram durante o período de observação?

8. O servidor analisado esteve ocupado por 19 minutos durante um período de observação de 30 minutos. Qual é a utilização desse servidor?

9. Um sistema computacional foi observado durante sete dias e verificou-se que, em média, o sistema estava sendo utilizado por 16 horas em cada dia. Qual a utilização do sistema durante esses sete dias?

## Exercícios – Teoria das filas

6. Em um determinado sistema observamos que 10 solicitações foram atendidas durante o tempo de observação e que esse servidor esteve ocupado por 200 segundos durante esse mesmo período de observação. Qual o tempo médio de serviço observado desse sistema?

$$N_a = 10 ; T_a = 200 \text{ seg} \Rightarrow \mu = N_a / T_a ; \mu = 10 / 200 \text{ seg} \Rightarrow T_s = 1 / \mu = 200 / 10 = 20 \text{ seg.}$$

7. Um servidor apresenta uma taxa de chegada de 10 solicitações/segundo. Se observarmos o servidor por 10 minutos, quantas solicitações de serviço chegaram durante o período de observação?

$$\lambda = 10 \text{ solicitações/seg} ; T = 10 \text{ min}$$

$$\lambda = N_c / T_c \Rightarrow N_c = \lambda \times T_c = 10 \times 10 \times 60 = 6000 \text{ solicitações/seg}$$

8. O servidor analisado esteve ocupado por 19 minutos durante um período de observação de 30 minutos. Qual é a utilização desse servidor?

$$U = t / T = 19 / 30 = 0,63 \text{ OU } 63\%$$

9. Um sistema computacional foi observado durante sete dias e verificou-se que, em média, o sistema estava sendo utilizado por 16 horas em cada dia. Qual a utilização do sistema durante esses sete dias?

$$U = t / T = 16 / 24 = 0,67 \text{ OU } 67\%$$