

# Teoria das Filas



Mário Meireles Teixeira  
Departamento de Informática, UFMA  
mario@deinf.ufma.br



## Filas, filas...



- As filas são a “praga” do mundo atual!
- Espera-se em fila no banco, na padaria, no ponto de ônibus, no trânsito, no restaurante...
- Em sistemas computacionais, há filas por toda parte: para acessar a CPU, o disco, a memória, a rede, a impressora, os servidores e outros recursos
- Os clientes podem ser: pessoas, processos, threads, jobs, pacotes, transações de BD, requisições C/S
- As filas surgem porque a demanda de serviço é maior que a capacidade de atendimento do sistema

## [ O que é a Teoria das Filas? ]

- É um ramo da Probabilidade que estuda o fenômeno da formação de filas de solicitantes de serviços, fornecidos por um determinado recurso
- Permite estimar importantes medidas de desempenho de um sistema a partir de propriedades mensuráveis das filas
- Dessa forma, pode-se dimensionar um determinado sistema segundo a demanda dos seus clientes, evitando desperdícios ou gargalos ☺
- Contudo, as filas apresentam comportamento estocástico... ☹
- Aplicações:
  - fluxo de tráfego (veículos, pessoas, redes de comunicação)
  - escalonamento (pacientes, tarefas industriais, processos)
  - serviços de atendimento (bancos, restaurantes, servidores)

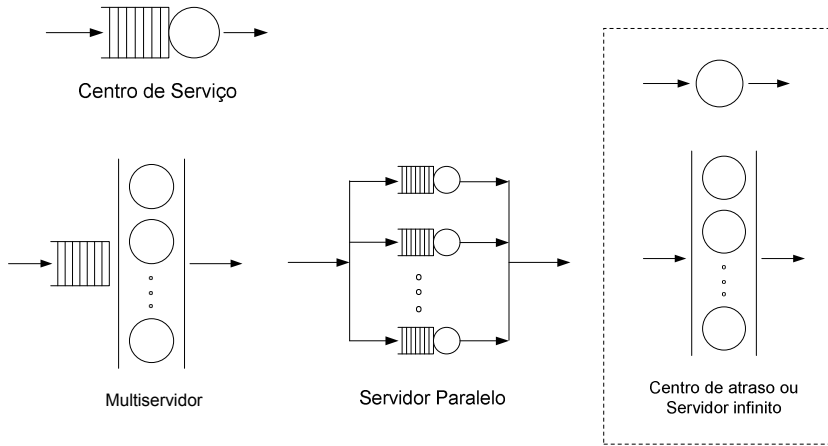
3

## [ Definições e Terminologia ]

- Rede de Filas
  - Consiste em um conjunto de entidades interligadas que oferecem serviços (os centros de serviço) e de usuários (os clientes)
- Centro de Serviço
  - Representa os recursos do sistema
  - Compreende um ou mais servidores e um conjunto de clientes esperando por serviço
- Fila = cliente em serviço + clientes em espera
- Fila de espera = clientes em espera

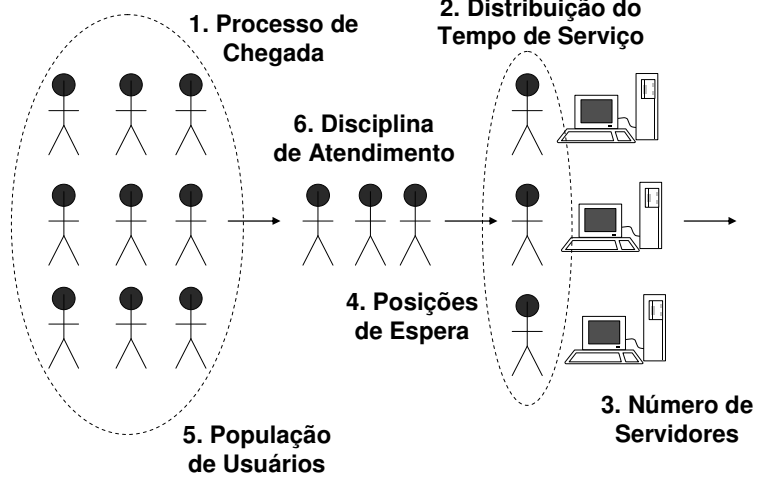
4

## Simbologia



5

## Componentes de uma Fila



6

## Características das Filas

1. Processo de Chegada: apresenta um comportamento estocástico

É fundamental conhecer a distribuição de probabilidade dos tempos entre as chegadas:

- Mais comum: Chegadas de Poisson (tempos entre as chegadas são exponencialmente distribuídos)
- Outras distribuições: Erlang, hiperexponencial, arbitrária

E ainda:

- chegadas de clientes individuais / simultâneas
- cliente sempre decide ficar / não fica se a fila for muito grande / impaciente / muda de fila
- padrão de chegada estacionário / não-estacionário

7

## Características das Filas

2. Distribuição dos Tempos de Serviço: os tempos de serviço dos clientes também são variáveis aleatórias independentes e identicamente distribuídas (IID)

Distribuições comuns: exponencial, Erlang, hiperexponencial, arbitrária (geral)

E ainda:

- atendimentos simples / batch
- serviço independente / dependente do estado
- serviço estacionário / não-estacionário

8

## Características das Filas

3. Número de Servidores: número de posições de atendimento disponíveis no sistema
  - Servidores idênticos / distintos
  - Fila única / por servidor / por grupo de servidores
4. Capacidade do Sistema: número máximo de clientes que podem permanecer no sistema, devido a restrições de espaço (buffers) ou de tempo de espera
  - Inclui clientes em serviço e esperando por serviço
  - Capacidade pode ser finita / infinita (mais fácil de analisar)

9

## Características das Filas

5. Tamanho da População (fonte): número potencial de clientes que podem chegar a um sistema
  - Tamanho finito / infinito
6. Disciplina de Atendimento (de fila): ordem na qual os clientes são atendidos  
Tipos:
  - FCFS (First-Come First-Served) ou FIFO
  - LCFS (Last-Come First-Served) ou LIFO
  - RR (Round Robin) / PS (Processor Sharing)
  - Prioridades (fila única / múltiplas filas)
  - Não-preemptivo / Preemptivo (resume/repeat) ...

10

## [ Notação de Kendall ]

- Para especificar um sistema de filas, é preciso conhecer as seis características anteriores
- Notação de Kendall:

$A/S/m/K/N/Q$

- **A** : distribuição dos tempos entre chegadas
- **S** : distribuição dos tempos de serviço
- **m** : número de servidores
- **K** : capacidade do sistema
- **N** : tamanho da população
- **Q** : disciplina de atendimento

11

## [ Notação de Kendall ]

- Exemplos:
  - M/G/4/50/2000/LCFS
  - D/M/1/∞/∞/RR
  - D/U/2/1000/∞/FCFS
  
  - A/S/m ou A/S/m/∞/∞/FCFS (.../∞/∞/FCFS é default)
  - M/M/1 ou M/M/1/∞/∞/FCFS
  - M/M/m
  - M/M/m/K
  - G/G/1

12

## Distribuições de Probabilidade

- Determinista (D)
  - Tempo entre as chegadas e o Tempo de Serviço são constantes
  - Não há variância estatística
  - Pelo menos uma das distribuições (chegada ou serviço) precisa ser aleatória, caso contrário o sistema de filas terá baixa aplicabilidade no mundo real

13

## Distribuições de Probabilidade

- Exponencial (M)
  - Para Chegadas (A): o intervalo entre uma chegada e a próxima é completamente independente do período anterior
  - Para Tempos de Serviço (S): o tempo de serviço atual é independente do tempo de serviço anterior
  - Esses processos são ditos "sem memória", pois seus intervalos não estão correlacionados no tempo. Portanto, podem ser caracterizados por uma distribuição exponencial

14

## Distribuições de Probabilidade

- Uniforme (U)
  - Os tempos de chegada estão limitados por algum valor finito ( $a \leq x \leq b$ )
  - A probabilidade de  $x$  assumir qualquer dos valores do intervalo é a mesma  $\rightarrow$  média =  $(a + b)/2$
- Arbitrária ou Geral (G)
  - Não é especificada uma distribuição de probabilidade para os tempos de chegada e serviço
  - Resultados são válidos para todas as distribuições

15

## Distribuições de Probabilidade

- Erlang ( $E_k$ )
  - Generalização da distribuição exponencial
  - Um servidor com  $k$  estágios, obedecendo à distribuição de Erlang, pode ser representado por uma seqüência de  $k$  servidores com tempos de serviço exponencialmente distribuídos, de mesma média
- Hiperexponencial ( $H_k$ )
  - Cada estágio no modelo de Erlang tem média diferente para os tempos de serviço
  - Os estágios estão organizados em paralelo, mas o serviço é fornecido um por vez

16



## Processos Estocásticos

- São funções ou seqüências aleatórias dependentes do tempo
- Exemplos:
  - $n(t)$  – número de jobs na CPU de um sistema
  - $W(t)$  – tempo de espera em fila
- Os processos estocásticos são úteis para representar o estado de sistemas de filas

17

## Tipos de Processos Estocásticos

- **Processos de Estado Discreto e Estado Contínuo:**
  - Discreto: número de valores de estado possíveis é finito ou contável; também chamado de *cadeia estocástica*. Ex:  $n(t)$
  - Contínuo: pode assumir qualquer valor entre os números reais. Ex:  $w(t)$
- **Processos de Markov:** os estados futuros do sistema independem do passado e dependem exclusivamente do estado atual.  
Um processo de Markov de estados discretos é chamado de *Cadeia de Markov*

18

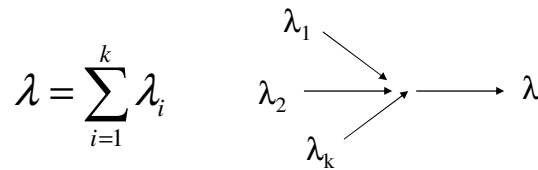
# [ Tipos de Processos Estocásticos ]

- **Processos de Nascimento e Morte:** processos de Markov de espaço discreto em que as transições entre estados estão restritas a estados vizinhos.  
Ex: número de jobs em um servidor único com chegadas individuais
- **Processos de Poisson:** se os tempos entre as chegadas têm distribuição exponencial, o número de chegadas em um dado intervalo terá uma distribuição de Poisson. O processo de chegadas é chamado de *Processo de Poisson*.  
As chegadas são "sem memória", pois o tempo entre as chegadas é IID e exponencialmente distribuído.

19

# [ Propriedades dos Fluxos de Poisson ]

- Superposição

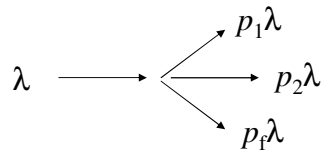


- A superposição de fluxos de Poisson dá como resultado um novo fluxo de Poisson cuja taxa de chegada é o somatório das taxas dos fluxos originais

20

# [ Propriedades dos fluxos de Poisson ]

- Divisão

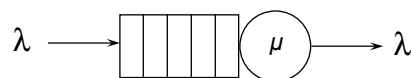


- Se um fluxo de Poisson for dividido em  $f$  subfluxos com probabilidade  $p_i$  de um usuário seguir o subfluxo  $i$ , então cada subfluxo é também um fluxo de Poisson com taxa média  $p_i\lambda$

21

# [ Propriedades dos fluxos de Poisson ]

- Se as chegadas a uma fila com um servidor único e tempo de serviço exponencial forem Poisson com taxa média  $\lambda$ , então as partidas também serão Poisson, com a mesma taxa  $\lambda$  (desde que  $\lambda < \mu$ )

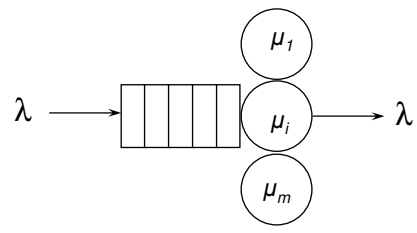


- $\lambda < \mu$ : condição de equilíbrio do sistema

22

# [ Propriedades dos fluxos de Poisson ]

- Se as chegadas a uma fila com  $m$  servidores e tempos de serviço exponenciais forem Poisson com taxa média  $\lambda$ , então as partidas também serão Poisson, com a mesma taxa  $\lambda$  (desde que  $\lambda < \sum \mu_i$ )



23

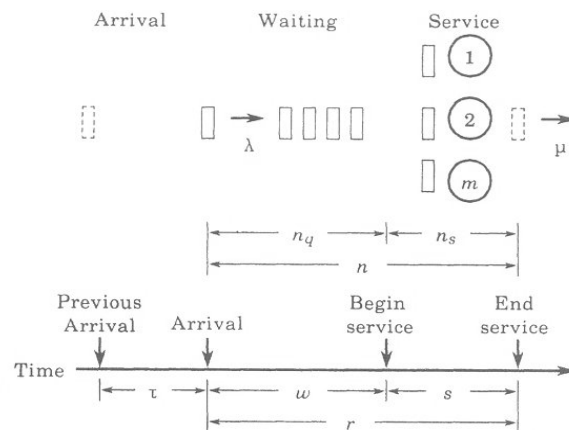
# [ Leis Operacionais ]

## Introdução

- Relações simples que não necessitam de nenhuma hipótese sobre as distribuições dos tempos de serviço ou dos intervalos entre chegadas
- Foram identificadas inicialmente por Buzen (1976) e posteriormente estendidas por Denning e Buzen (1978)
- A palavra operacional significa que pode ser medida diretamente

25

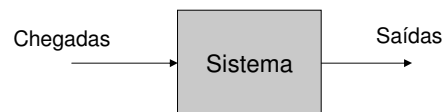
## Variáveis Aleatórias de uma Fila



26

## Quantidades Operacionais

- São quantidades que podem ser medidas diretamente durante um período finito de observação:
  - Período de observação –  $T$
  - Número de chegadas (arrivals) –  $A_i$
  - Número de términos (completions) –  $C_i$
  - Tempo ocupado (busy time) –  $B_i$



27

## Quantidades Operacionais

$$\text{Taxa de chegada } \lambda_i = \frac{\text{número de chegadas}}{\text{tempo}} = \frac{A_i}{T}$$

$$\text{Throughput } X_i = \frac{\text{número de términos}}{\text{tempo}} = \frac{C_i}{T}$$

$$\text{Utilização } U_i = \frac{\text{tempo ocupado}}{\text{tempo total}} = \frac{B_i}{T}$$

$$\text{Tempo médio de serviço } S_i = \frac{\text{tempo total de serviço}}{\text{número de saídas}} = \frac{B_i}{C_i}$$

Estas quantidades são variáveis que podem mudar de um período de observação para outro, mas as relações permanecem válidas!

28

## [ Lei da Utilização ]

$$U_i = \frac{B_i}{T} = \frac{C_i}{T} \times \frac{B_i}{C_i}$$
$$U_i = X_i S_i$$

29

## [ Exemplo 33.1 ]

- Considere um roteador em que os pacotes chegam a uma taxa de 125 pps e o roteador leva em média 2 ms para encaminhá-los. Qual a utilização do sistema?

$$X_i = \text{taxa de saída} = \text{taxa de chegada} = 125 \text{ pps}$$

$$S_i = 0,002 \text{ segundos}$$

$$U_i = X_i S_i = 125 \times 0,002 = 0,25 = 25\%$$

Este resultado é válido para qualquer processo de chegada ou atendimento!

30

## [ Lei de Little ]

- A Lei de Little relaciona o número de clientes no sistema com o tempo médio despendido no sistema:

$$Q_i = \lambda_i R_i$$

Número médio = Taxa de chegada x Tempo médio de resposta

- $R_i = S_i + W_i$
- Esta lei se aplica sempre que o número de chegadas for igual ao número de saídas (sistema em equilíbrio)
- Pode-se aplicar a lei de Little a qualquer sistema ou subsistema (caixa preta)

31

## [ Lei de Little ]

- Se o sistema está em equilíbrio, a taxa de chegada é igual ao throughput, portanto:

$$Q_i = X_i R_i$$

- Exemplo 3.14: Um servidor de arquivos NFS foi monitorado durante 30 minutos e o número observado de operações de I/O foi 10.800. Apurou-se que o número médio de pedidos ativos no NFS era três. Qual o tempo médio de resposta por pedido no servidor?

32



## Lei do Fluxo Forçado

- Relaciona o throughput global do sistema com o throughput dos dispositivos individuais
- Se o período de observação  $T$  for tal que o número de chegadas em cada dispositivo é igual ao número de saídas, i.e.,  $A_i = C_i$ , diz-se que o dispositivo satisfaz a Hipótese de Equilíbrio (*job flow balance*)
- Para um período de observação longo o bastante, a diferença  $A_i - C_i$  é normalmente pequena se comparada com  $C_i$

33

## Lei do Fluxo Forçado

- Seja  $V_i$  o número médio de visitas ao recurso  $i$  por uma tarefa
- Cada pedido que termina precisa passar, em média,  $V_i$  vezes pelo recurso  $i$ . Assim, se  $X$  pedidos foram concluídos por unidade de tempo, temos que  $V_i X$  pedidos terão passado pelo recurso  $i$ :

$$X_i = V_i X$$

- Esta lei é aplicável sempre que a hipótese de equilíbrio for verdadeira

34

## Lei da Demanda de Serviço

- Combinando as leis da Utilização e do Fluxo Forçado, temos:

$$U_i = X_i S_i = X V_i S_i$$

ou

$$U_i = X D_i$$

- Onde  $D_i = V_i S_i$  é a demanda total de serviço no  $i$ -ésimo dispositivo
- O dispositivo com a maior demanda de serviço tem a maior utilização e pode tornar-se o gargalo do sistema

35

## Exemplos 3.12 e 3.13

- As transações de um banco de dados realizam uma média de 4,5 operações de I/O no servidor de BD. O servidor foi monitorado durante uma hora e, durante esse período, 7.200 transações foram concluídas.
  - Qual a taxa média de processamento no disco?
  - Se cada I/O de disco leva 20 ms em média, qual a utilização do disco?
  - Qual a demanda de serviço do disco?

36

## Lei Geral do Tempo de Resposta

- Sistemas de tempo compartilhado podem ser divididos em dois subsistemas: o subsistema de terminais e o subsistema central de processamento
- Dados os comprimentos individuais  $Q_i$  das filas de cada dispositivo, podemos calcular  $Q$ :

$$Q = Q_1 + Q_2 + \dots + Q_M$$

$$XR = X_1R_1 + X_2R_2 + \dots + X_MR_M$$

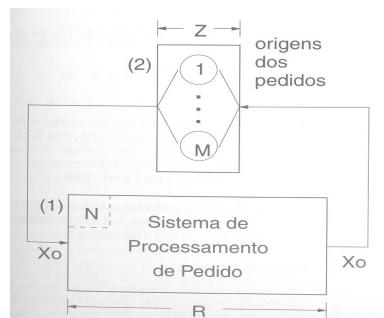
- Dividindo ambos os lados por  $X$  e usando a lei do fluxo forçado:

$$R = V_1R_1 + V_2R_2 + \dots + V_MR_M \quad \text{ou} \quad R = \sum_{i=1}^M R_iV_i$$

37

## Lei do Tempo de Resposta Interativo

- Num sistema interativo, os usuários geram pedidos que são processados pelo subsistema central e os resultados voltam ao terminal. Após um tempo ocioso  $Z$ , o usuário submete o próximo pedido.



38

## [ Lei do Tempo de Resposta Interativo ]

- Aplicando-se a lei de Little ao subsistema central, temos:

$$Q = XR$$

- Agora, aplicando-se a lei de Little aos M terminais:

$$\bar{M} = XZ$$

- Considerando que um cliente ou está sendo processado ou está ocioso:

$$M = Q + \bar{M} = XR + XZ = X(R + Z)$$

$$R = \frac{M}{X} - Z$$

39

## [ Exemplo 3.16 ]

- Um portal corporativo oferece serviços na Web aos funcionários de uma empresa. Em média, 500 funcionários estão on-line solicitando serviços. Uma análise do log do portal revelou que, em média, 6.480 pedidos são processados por hora. O tempo de resposta médio por pedido é de cinco segundos.

Qual o tempo médio entre o momento em que a resposta a uma réplica é recebida e um novo pedido é enviado por um funcionário?

40

**Box 33.1 Operational Laws**

Utilization law	$U_i = X_i S_i = X D_i$
Forced flow law	$X_i = X V_i$
Little's law	$Q_i = X_i R_i$
General response time law	$R = \sum_{i=1}^M R_i V_i$
Interactive response time law	$R = N/X - Z$
Asymptotic bounds	$R \geq \max\{D, N D_{\max} - Z\}$ $X \leq \min\{1/D_{\max}, N/(D + Z)\}$

Symbols:

$D$	Sum of service demands on all devices, = $\sum_i D_i$
$D_i$	Total service demand per job for the $i$ th device, = $S_i V_i$
$D_{\max}$	Service demand on the bottleneck device, = $\max_i \{D_i\}$
$N$	Number of jobs in the system
$Q_i$	Number in the $i$ th device
$R$	System response time
$R_i$	Response time per visit to the $i$ th device
$S_i$	Service time per visit to the $i$ th device
$U_i$	Utilization of the $i$ th device
$V_i$	Number of visits per job to the $i$ th device
$X$	System throughput
$X_i$	Throughput of the $i$ th device
$Z$	Think time