



Laboratório de Pesquisa em Redes e Multimídia

# Escalonamento no Unix



Universidade Federal do Espírito Santo  
Departamento de Informática

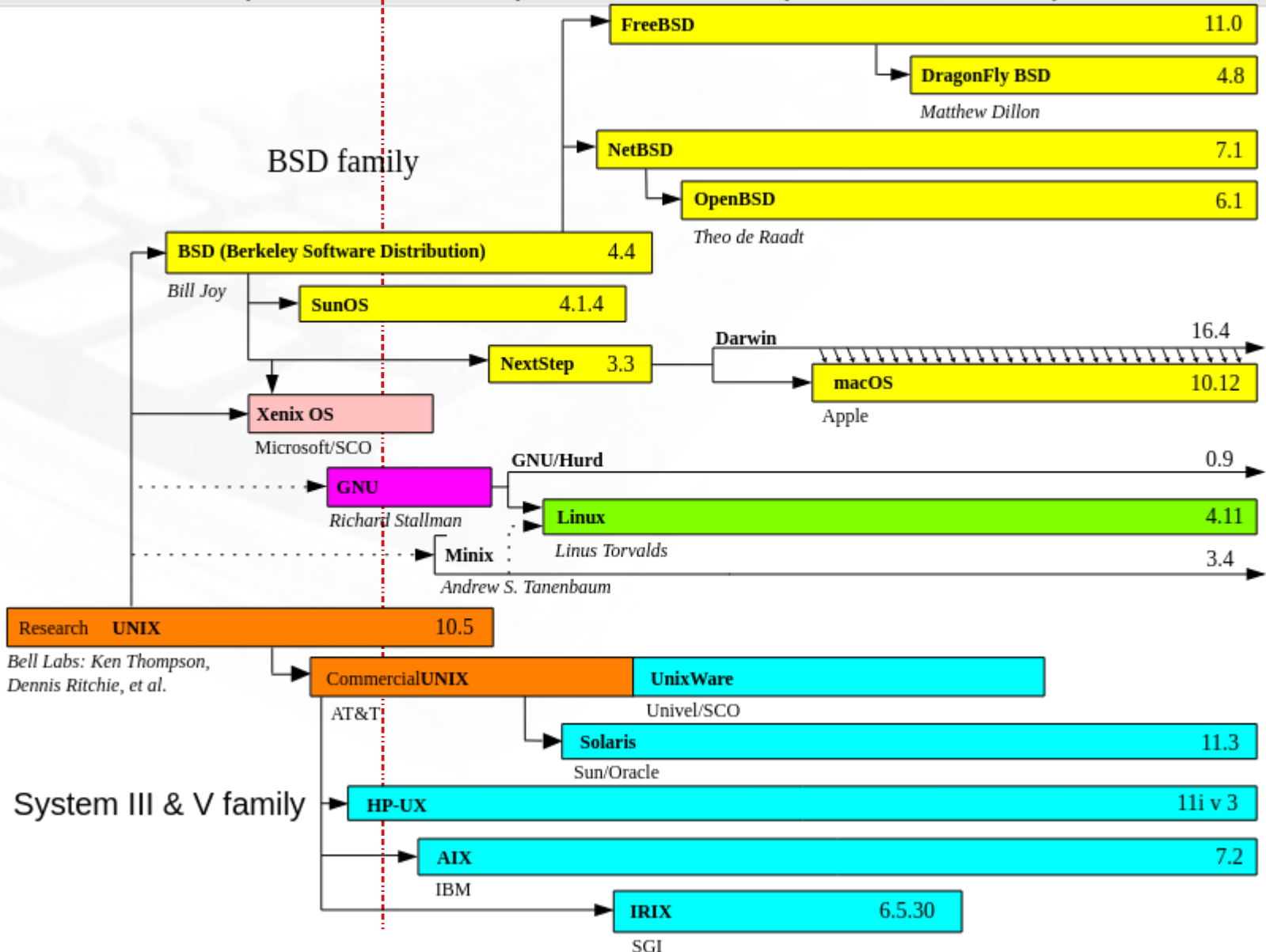
**Sistemas Operacionais**

# Projeto do Escalonador

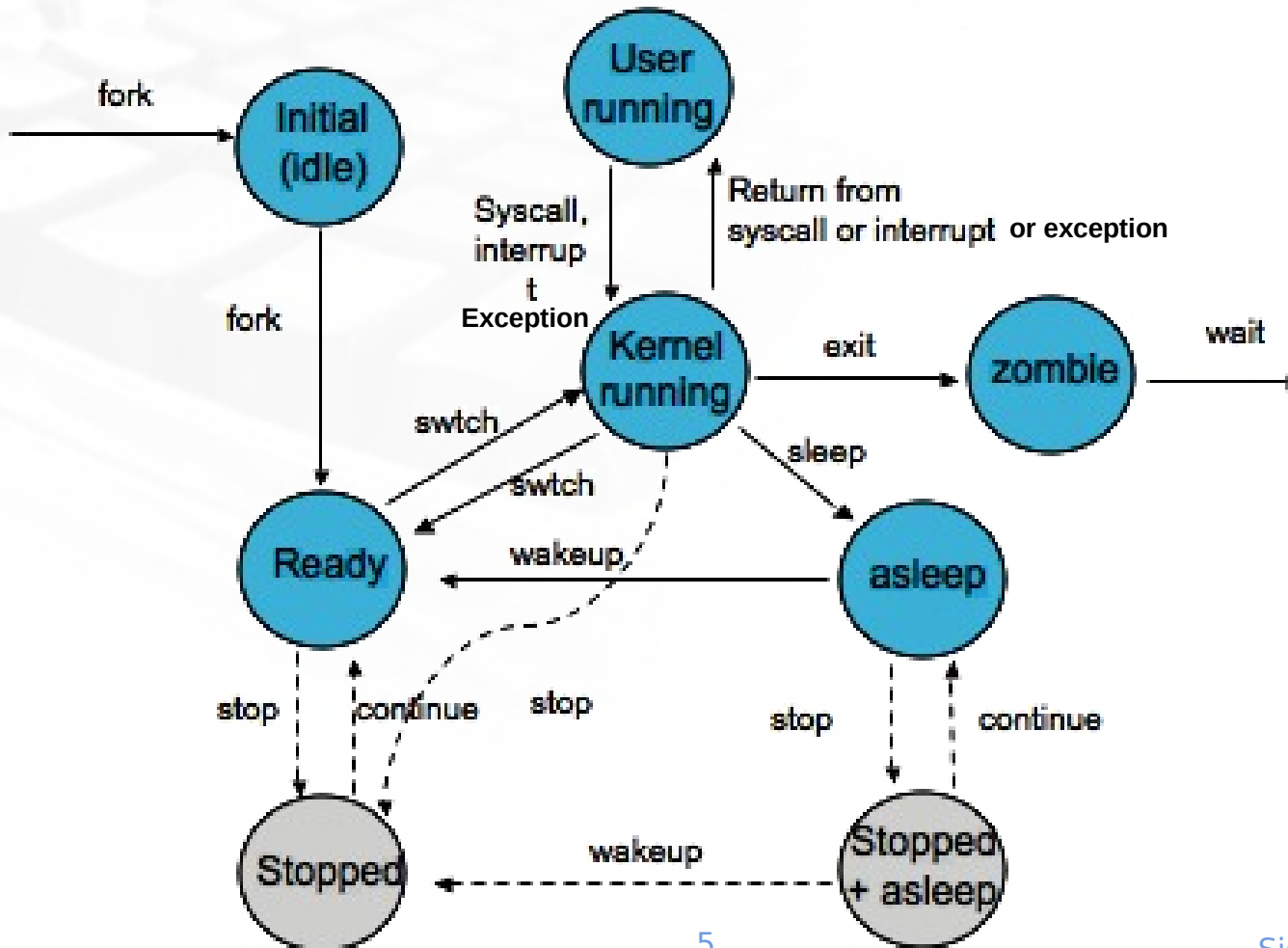
- O projeto de um escalonador deve focar em dois aspectos:
  - **Política de escalonamento** – estabelece as regras usadas para decidir para qual processo ceder a CPU e quando chaveá-la para um outro processo.
  - **Implementação** – definição dos algoritmos e estruturas de dados que irão executar essas políticas.
- A política de escalonamento deve buscar:
  - **Tempos de resposta** rápidos para aplicações interativas.
  - Alto **throughput** (vazão) para aplicações em *background*.
  - Evitar **starvation**
- Os objetivos acima podem ser **conflitantes!**

# Objetivos do Escalonador (cont.)

- Processos **interativos** (*shells*, editores, interfaces gráficas, etc.)
  - Interações devem ser processadas rapidamente
  - Requisito: reduzir os tempos de resposta médios (e a variância), de forma que o usuário não detecte/perceba este atraso (50-150 ms)
- Processos **batch** (que rodam em *background*, sem interação com o usuário)
  - Requisito: o tempo para completar uma tarefa na presença de outras atividades, comparado ao tempo para completá-la em um sistema “inativo”, sem outras atividades paralelas (turnaround).
- Processos de **tempo real (soft)**
  - Requisito: como as aplicações são “*time-critical*”, o escalonador deve ter um comportamento previsível, com limites garantidos nos tempos de resposta.
  - Ex: aplicações de vídeo.
- As **funções do kernel**
  - Devem ser executadas prontamente (gerência de memória, tratamento de interrupções e gerência de processos)



# Máquina de Estados do Unix



## Modos de Operação da CPU

- O Unix requer do hardware a implementação de apenas 2 modos:
  - *user mode* - menos privilegiado. Execução de certas instruções e acesso a certos endereços é proibido.
  - *kernel mode* - mais privilegiado. Todas as instruções podem ser executadas, acesso à memória é irrestrito.
- “**Código de usuário**” é executado em ***user mode***; logo, não podem – acidental ou maliciosamente –, corromper outro processo ou mesmo o kernel.
- “**Código de kernel**” roda em ***kernel mode***

# Escalonamento Tradicional

- Usado originalmente nos sistemas Unix **SVR3** e **4.3BSD**:
  - Tempo compartilhado
  - Ambientes interativos
  - Processos *background* (batch) e *foreground* rodando simultaneamente
    - Melhorar os tempos de resposta para os usuários interativos, e
    - Garantir, ao mesmo tempo, que processos *background* não sofram *starvation*
- Prioridades dinâmicas
  - A cada processo é atribuída uma prioridade de escalonamento, que é alterada com o passar do tempo.
    - Se o processo não está no estado *running*, o kernel periodicamente aumenta a sua prioridade.
    - Quanto mais o processo recebe a posse da CPU mais o kernel reduz a sua prioridade.
  - **Por que essa estratégia previne *starvation* e favorece processos *I/O bound*?**

## Escalonamento Tradicional (cont.)

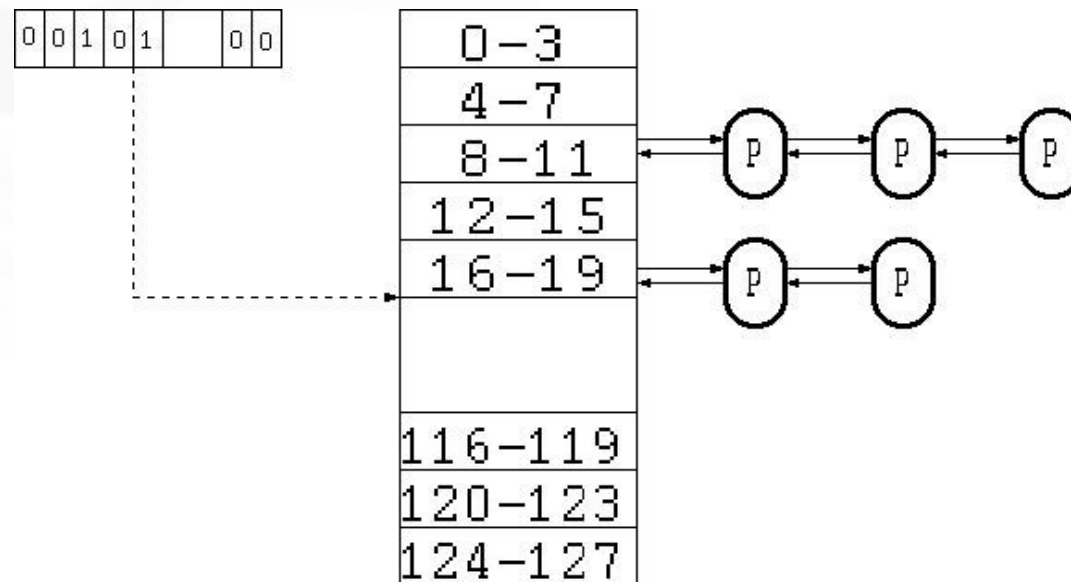
- O escalonador sempre seleciona o processo com a **prioridade mais alta** dentre aqueles no estado pronto-para-execução (*ready/runnable process*).
- Processos **com a mesma prioridade**:
  - “**Preemptive round robin**”: escalonamento circular com preempção
  - O *quantum* possui um valor fixo, tipicamente de 100 ms.
- O kernel do Unix tradicional é **não-preemptivo**
  - A chegada de um processo de mais alta prioridade na fila de prontos força a preempção do processo em execução SE ele está executando em user mode
  - Um processo executando em kernel mode nunca é preemptado (nunca perde a posse da CPU para um outro processo).
  - Quando o sistema for retornar para *user mode*, é verificado se há um processo mais prioritário na fila de prontos. Se sim, ocorre a preempção (troca de contexto).



# Prioridades dos Processos

- Prioridades recebem valores entre 0 e 127 (quanto menor o valor numérico maior a prioridade)
  - 0 - 49 : processos do kernel (*kernel priorities*)
  - 50 - 127 : processos de usuário (*user priorities*)

**Ex: 4.3BSD**



## Prioridades dos Processos (cont.)

- Campos de prioridade na *proc struct*
  - *p\_pri* prioridade de escalonamento atual
  - *p\_usrpri* prioridade em *user mode*
  - *p\_cpu* medida de uso recente da CPU
  - *p\_nice* fator *nice* (gentileza) controlável pelo usuário (*SVC nice*)

$$p\_usrpri = PUSER + \frac{p\_cpu}{4} + 2 \times p\_nice$$

$$p\_usrpri = PUSER + \frac{p\_cpu}{4} + 2 \times p\_nice$$

## 4.3 BSD

- **PUSER = 50**
  - prioridade base dos processos de usuário
- **p\_nice**=[0-39] (controlado pelo usuário via SVC `nice()`)
  - Default: `p_nice = 20`
  - **nice(x)** :  $-20 \leq x \leq +19 \Rightarrow p\_nice = p\_nice + x$ 
    - Se **x>0**, a chamada **diminui a prioridade** final do processo
    - Se **x<0**, a chamada **umenta a prioridade** final do processo (apenas superusuário)
  - Processos *background* recebem automaticamente grandes valores de `p_nice`, por isso são menos prioritários.

$$p\_usrpri = PUSER + \frac{p\_cpu}{4} + 2 \times p\_nice$$

## Prioridades dos Processos (cont.)

### 4.3 BSD

- A cada ciclo (CPU tick ... 1 tick = 10ms):
  - Incrementa o `p_cpu` do processo que está rodando
- A cada 4 ciclos:
  - **Recalcula** `p_usrpri` do processo que está **rodando**
  - Se houver um processo mais prioritário, *context switch*
- A cada 10 ciclos
  - Round-robin entre processos de mesma prioridade (*context switch*)
- A cada 100 ciclos
  - Aplica um *decay factor* no `p_cpu` de TODOS os processos
  - **Recalcula** `p_usrpri` de **todos** os processos (possível *context switch*)

## *Sleep Priority*

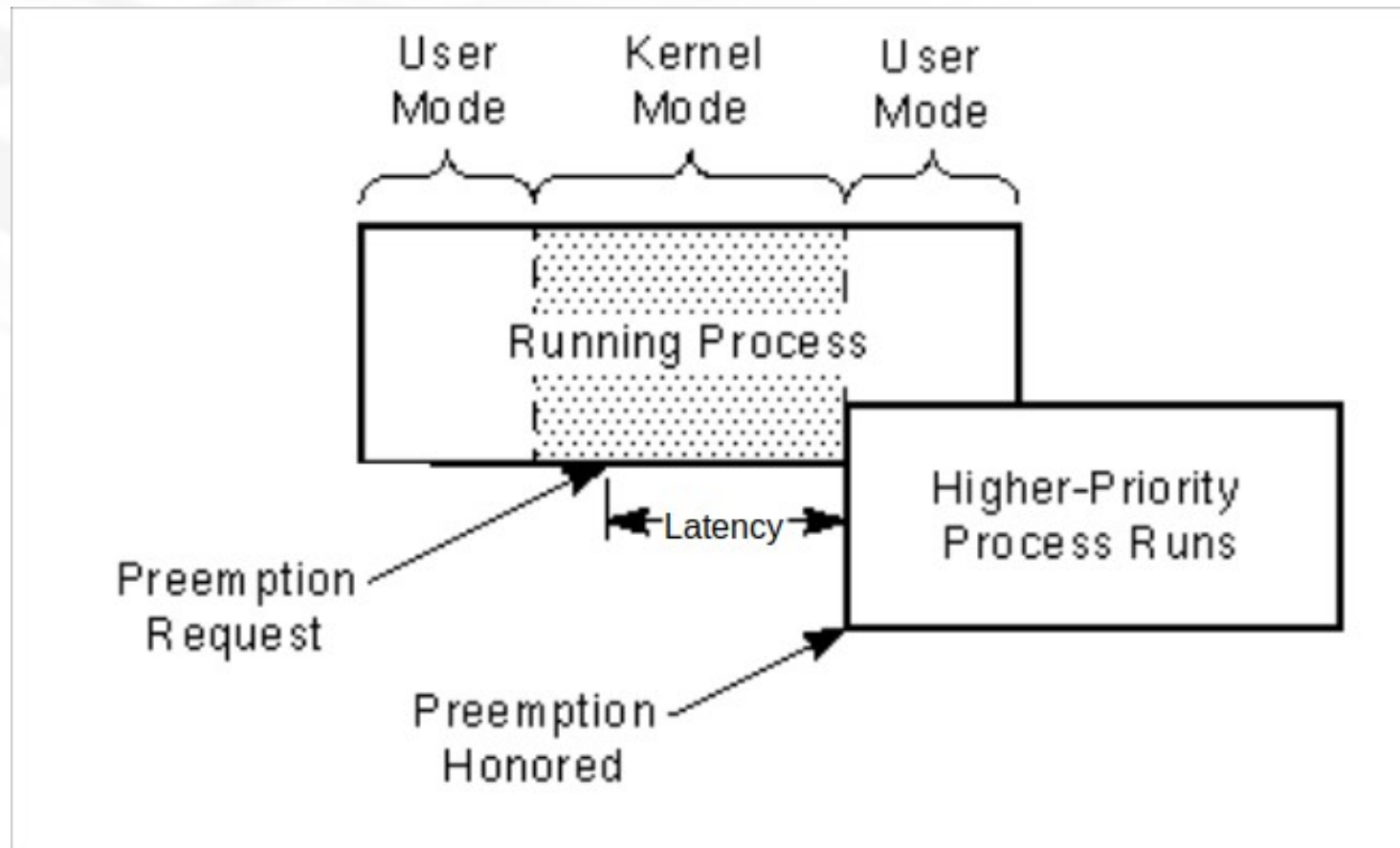
- Quando o processo está bloqueado, ele é associado a uma *sleep priority*
  - Relacionada com o dispositivo pelo qual ele está esperando
  - É uma prioridade de kernel ( $<50$ )
- Quando ele é acordado, o kernel seta o valor de  $p\_pri$  para o mesmo valor da *sleep priority*
  - O processo será escalonado na frente de outros processos de usuário e continuará a sua execução em modo kernel a partir do ponto de bloqueio.
- Quando a SVC é finalmente completada, imediatamente antes de retornar ao modo usuário, o kernel faz  $p\_pri = p\_usrpri$

## Escalonamento Tradicional (cont.)

- O kernel do Unix tradicional é **não-preemptivo**
  - A chegada de um processo de mais alta prioridade na fila de prontos força a preempção do processo em execução SE ele está executando em user mode
  - Um processo executando em kernel mode nunca é preemptado (nunca perde a posse da CPU para um outro processo).
  - Quando o sistema for retornar para *user mode*, é verificado se há um processo mais prioritário na fila de prontos. Se sim, ocorre a preempção (troca de contexto).

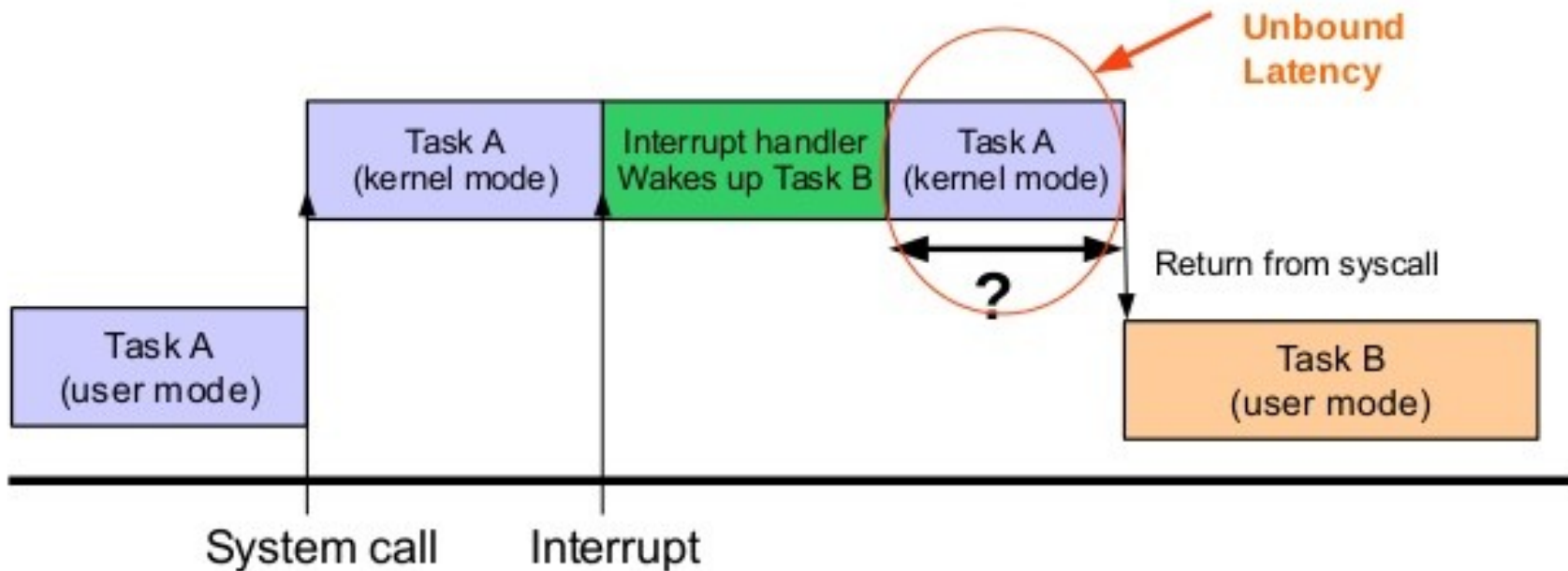
# Escalonamento Tradicional (cont.)

- Kernel não-preemptivo



# Escalonamento Tradicional (cont.)

- Kernel não-preemptivo





# Implementação BSD - Análise (1)

- Vantagens:
  - O algoritmo de escalonamento tradicional do Unix é simples e efetivo, sendo adequado para:
    - Sistemas de tempo compartilhado (*time sharing*)
    - Mix de processos interativos e *batch*.
  - Recomputação dinâmica das prioridades previne a ocorrência de *starvation*.
  - A abordagem favorece processos I/O bound, que requerem *bursts* de CPU pequenos e pouco freqüentes .

## Implementação BSD - Análise (2)

- Deficiências:
  - **Baixa escalabilidade:** se o número de processos é muito alto, torna-se ineficiente recalcular todas as prioridades a cada segundo;
  - Não existe a **garantia de alocação da CPU** para um processo específico ou então para um grupo de processos;
  - Não existe garantias de tempos de resposta para aplicações com característica de **tempo-real**.
  - Aplicações não podem controlar as suas prioridades.
    - O mecanismo de *nice* é muito simplista e inadequado.
  - Como o kernel é **não-preemptivo**, processos de maior prioridade podem ter que esperar muito tempo para ganhar a posse da CPU (problema da inversão)..

## Implementação - SVR4 (1)

- S.O. reprojetoado
  - Orientação a objeto
- Objetivos de projeto do escalonador no SVR4:
  - Suportar mais aplicações, incluindo tempo-real
  - Permitir às aplicações maior controle sobre prioridade e escalonamento
  - Permitir a adição de novas políticas de uma forma modular
  - Limitar a latência de despacho para aplicações dependentes do tempo
- Classes de escalonadores
  - Classes oferecidas originalmente: time-sharing e tempo-real
  - É possível criar novas classes tratando outros tipos de processos

## Implementação - SVR4 (2)

- Existem **rotinas independentes de classe** para fornecer:
  - Mudança de contexto
  - Manipulação da fila de processos
  - Preempção
- **Rotinas dependentes de classe**
  - Funções virtuais implementadas de forma específica por cada classe (herança)
  - Recomputação de prioridades
    - real-time class - prioridades e quanta fixos
    - time-sharing class - prioridades variam dinamicamente
      - Processos com menor prioridade têm maior quantum
      - Usa *event-driven scheduling*: prioridade é alterada na resposta a eventos.

# Implementação - SVR4 (3)

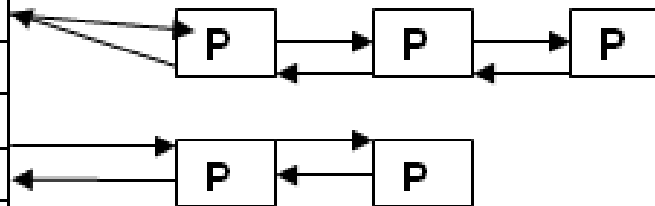
**dqactmap** (global variable – bitmask 160 cells – one bit for each queue)

0	0	1	0	0	1	0	...
---	---	---	---	---	---	---	-----

**dispq** (160 rows)

160
159
158
157
156
155
...
0

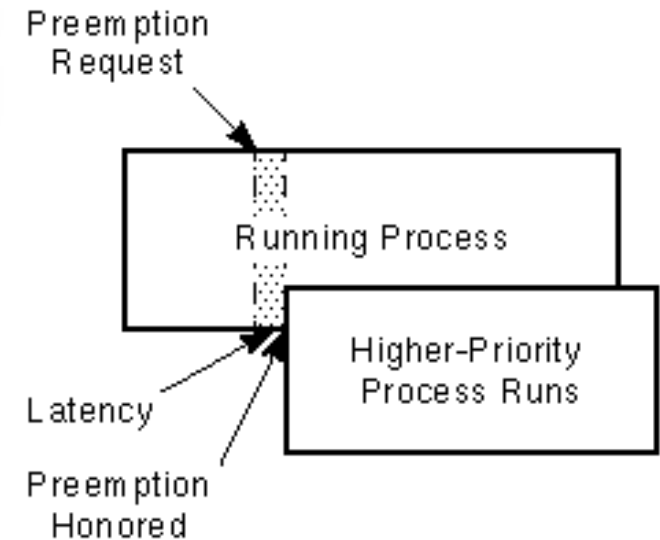
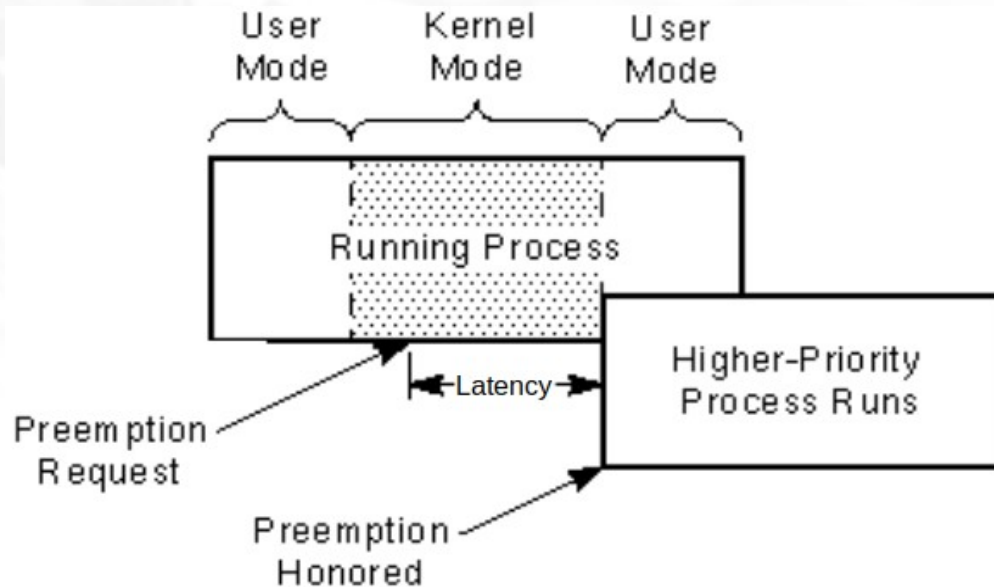
proc structures



runnable processes

0 – 59	time-sharing class
60 – 99	system priorities
100 – 159	real-time class

# Kernel preemptivo x Kernel não-preemptivo



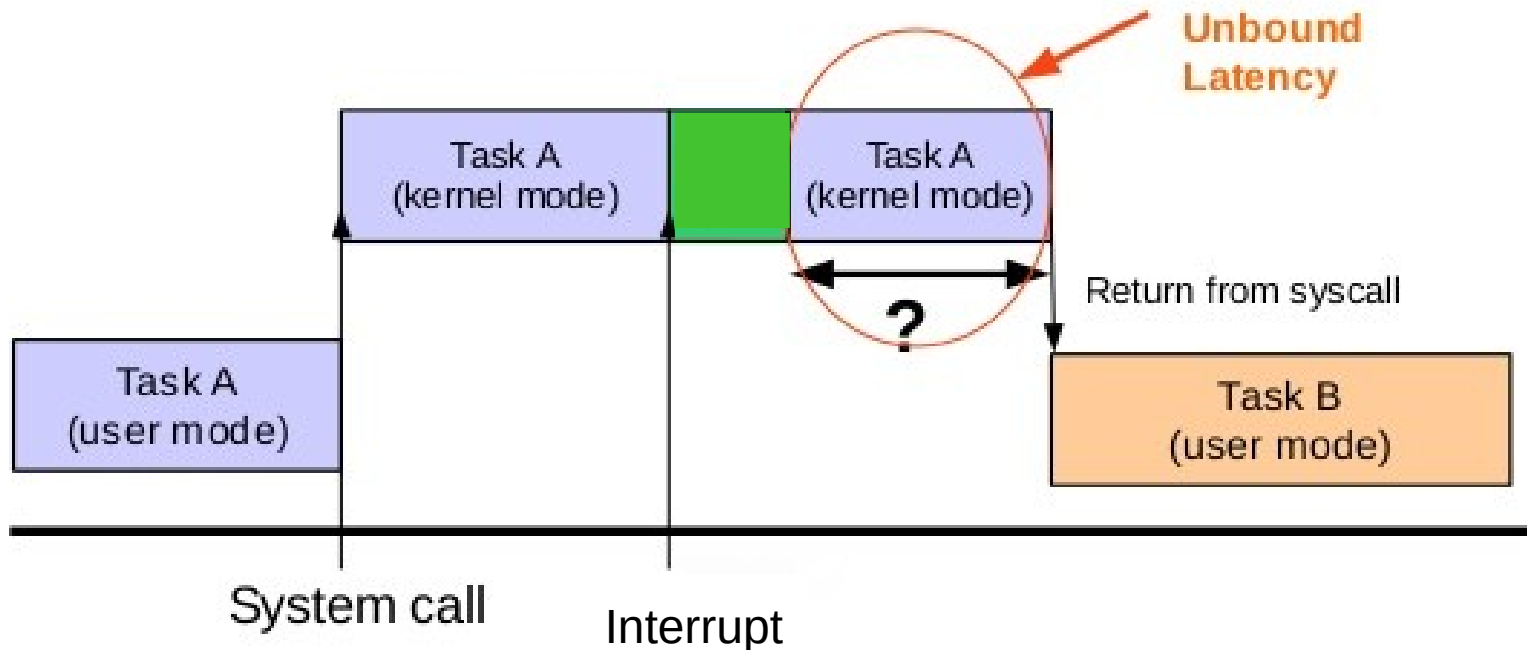
MLO-0073 13

## Implementação - SVR4 (4)

- Processos de tempo real exigem tempos de resposta limitados
- **Preemption points** são definidos em pontos do kernel onde
  - Todas as estruturas de dados do kernel encontram-se estáveis
  - O kernel está prestes a iniciar alguma computação longa
- Em cada **preemption point**
  - O kernel verifica se um processo de tempo-real tornou-se pronto e precisa ser executado
    - O processo corrente é então preemptado
- Os limites nos tempos máximos que um processo de tempo-real precisa esperar são definidos pelo maior intervalo entre dois **preemption points** consecutivos

# Implementação - SVR4 (5)

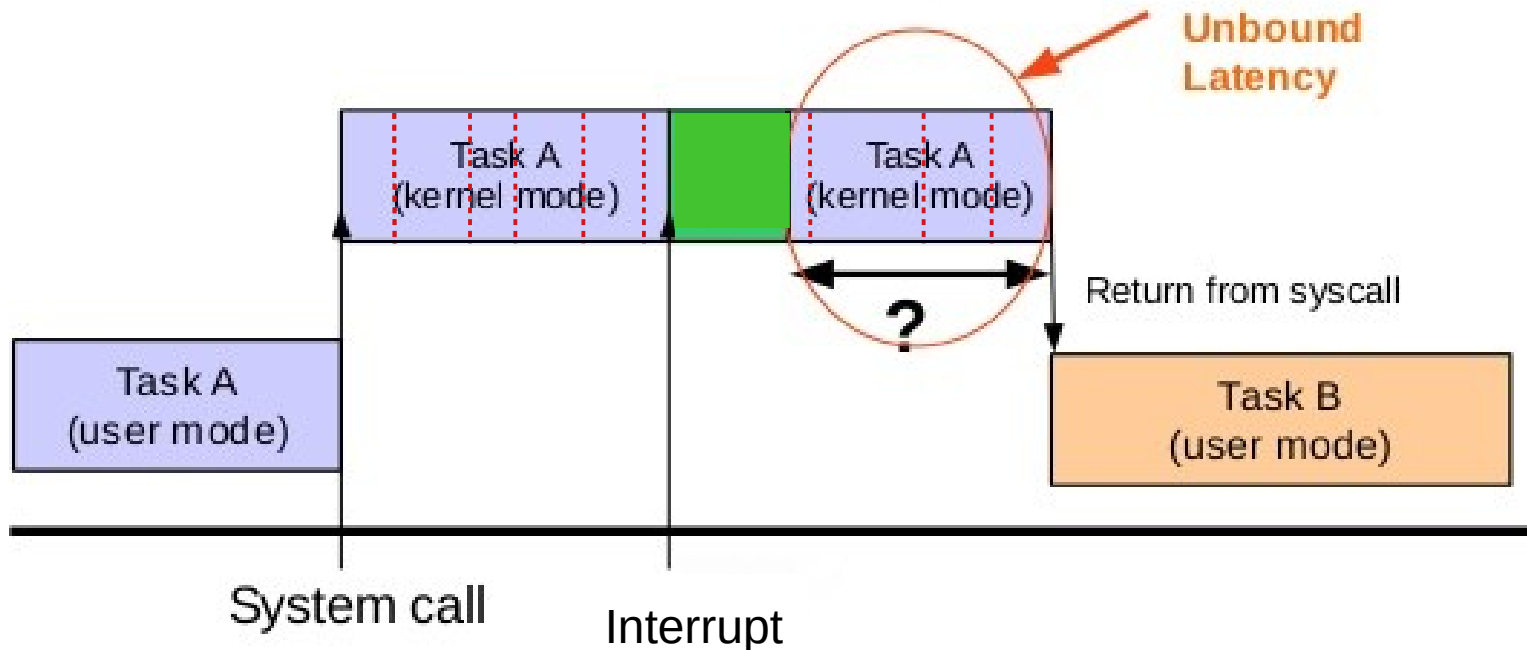
- Latency of Non-Preemptive configuration





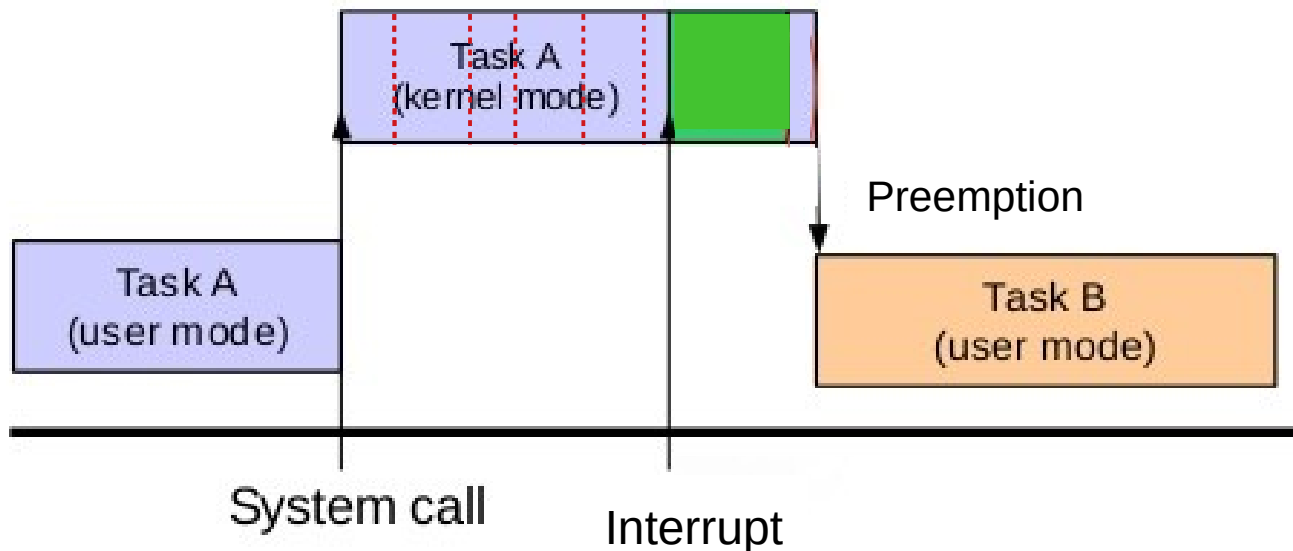
# Implementação - SVR4 (5)

- Latency of Non-Preemptive configuration



# Implementação - SVR4 (5)

- Latency of Non-Preemptive configuration



## E o LINUX?

- Processos de tempo real exigem tempos de resposta limitados
- **Preemption points** são definidos em pontos do kernel onde
  - Todas as estruturas de dados do kernel encontram-se estáveis
  - O kernel está prestes a iniciar alguma computação longa
- Em cada preemption point
  - O kernel verifica se um processo de tempo-real tornou-se pronto e precisa ser executado
    - O processo corrente é então preemptado
- Os limites nos tempos máximos que um processo de tempo-real precisa esperar são definidos pelo maior intervalo entre dois **preemption points** consecutivos

# POSIX 1003.1b Real-time extensions

---

- ▶ Priority Scheduling
- ▶ Real-Time Signals
- ▶ Clocks and Timers
- ▶ Semaphores
- ▶ Message Passing
- ▶ Shared Memory
- ▶ Asynchronous and Synchronous I/O
- ▶ Memory Locking

## Referências

- **VAHALIA, U. Unix Internals: the new frontiers. Prentice-Hall, 1996.**
  - **Capítulo 5 (até seção 5.5)**
- R. S. de Oliveira, A. S. Carissimi e S. S. Toscani, "Sistemas Operacionais", 4ª Edição (série didática da UFRGS), Editora Sagra-Luzzato, 2010.
  - Seção 4.5.5
- **REFERÊNCIAS EXTRA**
  - "A complete guide to Linux process scheduling". Msc Thesis (2015)
    - <https://trepo.tuni.fi/bitstream/handle/10024/96864/GRADU-1428493916.pdf>
  - "Making Linux do hard real time"
    - <https://events.static.linuxfound.org/sites/events/files/slides/>