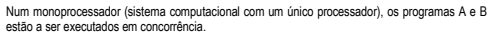


1.Descreva as duas perspectivas de definição de um sistema de operação. Mostre claramente em que circunstâncias cada uma delas é relevante.

**Paralelismo** - é a capacidade apresentada por um sistema computacional de poder executar em simultâneo dois ou mais programas; o que exige que o sistema computacional seja formado por mais do que um processador (um por cada programa em execução simultânea). Quando tal acontece, diz-se que o sistema operativo associado tem características de **multiprocessamento**.

**Concorrência** - é a ilusão criada por um sistema computacional de aparentemente poder executar em simultâneo mais programas do que o número de processadores existentes; o que exige que a atribuição do (s) processador (es) seja multiplexada no tempo entre os diferentes programas presentes. Quando tal acontece, diz-se que o sistema operativo associado tem características de **multiprogramação**.



	Utilizador A	Utilizador B
1	1	1
2	1	1
3	1	1
4	1	1
5	1	1
6	1	1
7	1	1
8	1	1
9	1	1
10	1	1
11	1	1
12	1	1
13	1	1
14	1	1
15	1	1
16	1	1
17	1	1
18	1	1
19	1	1
20	1	1
21	1	1
22	1	1
23	1	1
24	1	1
25	1	1
26	1	1
27	1	1
28	1	1
29	1	1
30	1	1
31	1	1
32	1	1
33	1	1
34	1	1
35	1	1
36	1	1
37	1	1
38	1	1
39	1	1
40	1	1
41	1	1
42	1	1
43	1	1
44	1	1
45	1	1
46	1	1
47	1	1
48	1	1
49	1	1
50	1	1
51	1	1
52	1	1
53	1	1
54	1	1
55	1	1
56	1	1
57	1	1
58	1	1
59	1	1
60	1	1
61	1	1
62	1	1
63	1	1
64	1	1
65	1	1
66	1	1
67	1	1
68	1	1
69	1	1
70	1	1
71	1	1
72	1	1
73	1	1
74	1	1
75	1	1
76	1	1
77	1	1
78	1	1
79	1	1
80	1	1
81	1	1
82	1	1
83	1	1
84	1	1
85	1	1
86	1	1
87	1	1
88	1	1
89	1	1
90	1	1
91	1	1
92	1	1
93	1	1
94	1	1
95	1	1
96	1	1
97	1	1
98	1	1
99	1	1
100	1	1

**Características comuns aos sistemas de operação atuais:** apresentam um ambiente gráfico de interação com o utilizador, o sistema sistemas interactivos multifuizadores e multiferreia. permitem a interacção em simultâneo com múltiplos utilizadores; permitem que cada utilizador tenha em simultâneo múltiplos programas em execução; implementam uma organização de memória virtual; são sistemas de operação de rede; permitem o acesso de um modo quase indistinto a sistemas de ficheiros e a dispositivos de entrada / saída locais e remotos; incluem aplicações que permitem, entre outras coisas, o login em máquinas remotas, o correio electrónico e a navegação na Internet; têm uma colecção muito grande de *device drivers* para potenciar a interligação de tipos muito variados de dispositivos de entrada / saída; permitem em muitos casos a ligação dinâmica de dispositivos (*plug and play*).

O papel desempenhado por cada nível será:

1. **Scheduling de baixo nível (Gestão do Processador):** Ao longo da sua existência, um processo vai encontrar-se em situações distintas, designadas por **estados**, as mais importantes das quais são as seguintes: **run** - quando detém a posse do processador e está, portanto, em execução; **ready-to-run** - quando aguarda a atribuição do processador para entrar em execução; **blocked** - quando está impedido de continuar, até que um acontecimento externo ocorra (acesso a um recurso, completamento de uma operação de entrada-saída, etc.). As transições entre estados resultam normalmente de uma intervenção externa, mas podem nalguns casos ser despoletadas pelo próprio processo. A parte do sistema operativo que lida com estas transições, chama-se **scheduler** (*do processador*) e constitui parte integrante do seu núcleo central (**kernel**), responsável por lidar com as interrupções e agendar a atribuição do processador e dos outros recursos do sistema computacional.
2. **dispatch** - um dos processos da fila de espera dos "processos prontos a serem executados" é selecionado pelo **scheduler** para execução; **timer-run-out** - o **scheduler** verificou que o processo em execução esgotou o intervalo de tempo de processação que lhe tinha sido atribuído; **sleep** - o processo decidiu que está impedido de prosseguir, ficando a aguardar a ocorrência de um acontecimento externo; **wake up** - o acontecimento externo que o processo aguardava, ocorreu.

Num sistema de tipo *batch* multiprogramado, não fará sentido a existência de três níveis de *scheduling*, visto que, nele não existe multiutilizador, isto é, o *batch* tende a bloquear os processos e a correr outros para optimiza-los.

Em sistemas do tipo 'batch' a disciplina a usar é do tipo non-preemptive porque se procura atribuir o processador de forma a minimizar o tempo de execução de um grupo de tarefas. Logo, não faz sentido comutar processos no estado RUN. (A única exceção surge quando o processo ultrapassa o tempo total de processador que lhe foi atribuído). Em sistemas operativos do tipo interactivo, a política de 'scheduling' a utilizar é a do tipo 'preemptive', pois pretende-se estabelecer uma multiplexagem temporal da ocupação do processador pelos diferentes processos.

12. Foi referido nas aulas que os sistemas de operação de tipo *batch* usam principalmente uma política de *scheduling non-preemptive*. Será, porém, uma política pura, ou admite excepções?

R: Em sistemas do tipo 'batch' a disciplina para usar é do tipo non-preemptive porque se procura atribuir o processador de forma a minimizar o tempo de execução de um grupo de tarefas. Logo, não faz sentido comutar processos no estado RUN. (A única excepção surge quando o processo ultrapassa o tempo total de processador que lhe foi atribuído).

13. Justifique claramente se a disciplina de *scheduling* usada em Linux para a classe *SCHED\_OTHER* é uma política de prioridade estática ou dinâmica?

R: A disciplina de *scheduling* usada em Linux para a classe *SCHED\_OTHER* é uma política de prioridade estática, na medida em que o Linux utiliza um algoritmo baseado em *créditos* no estabelecimento da sua prioridade, ora cada vez que o processador é atribuído a um processo a sua prioridade é decrementada, sendo desta forma toda a *scheduling* delimitado logo estático.

14. O que é o *aging* dos processos? Dê exemplos de duas disciplinas de *scheduling* com esta característica, mostrando como ela é implementada em cada caso.

R.: Em qualquer sistema onde se mantém processos à espera enquanto ele faz a alocação de recursos e decisões de *scheduling* de processos, é possível adiar indefinidamente o *scheduling* de um processo enquanto outros recebem toda a atenção por parte do sistema. Como se sabe chama-se a este fenómeno *adiamento indefinido*. Isto pode ocorrer quando os recursos são organizados com base numa atribuição dinâmica de prioridades, já que é possível que um processo tenha que esperar indefinidamente enquanto outros processos com uma prioridade mais elevada vão sendo servidos. Sendo assim o *aging* consiste no aumento linear da prioridade de um dado processo enquanto espera por um dado recurso. Assim a prioridade desse processo poderá, eventualmente, exceder a prioridade de todos os outros processos que chegam e será por isso «atendido». O *aging* é algo de útil pois evita o *starvation* - ex. sistemas I/O intensivos que pela necessidade de resposta iam ganhar sucessivamente a prioridade no acesso ao processador.

15. Distinga *threads* de processos. Assumindo que pretende desenvolver uma aplicação concorrente usando um dos paradigmas, descreva o modo como cada uma afecta o desenho da arquitectura dos programas associados.

O conceito de processo corporiza as propriedades seguintes: • *pertença de recursos* - um espaço de endereçamento próprio e um conjunto de canais de comunicação com os dispositivos de entrada / saída; • *fil de execução (thread)* - um *program counter* que sinaliza a localização da instrução que deve ser executada a seguir, um conjunto de *registos internos* do processador que contém os valores actuais das variáveis em processamento e um *stack* que armazena a história de execução (um *frame* por cada rotina invocada e que ainda não retornou). Estas propriedades, embora surjam reunidas num processo, podem ser tratadas separadamente pelo sistema de operação. Quando tal acontece, os processos dedicam-se a agrupar um conjunto de recursos e os *threads*, também conhecidos por *light weight processes*, constituem entidades executáveis independentes dentro do contexto de um mesmo processo. *Multithreading* representa então a situação em que é possível criar-se *threads* múltiplos de execução no contexto de um processo. Se quisermos recorrer a um dos paradigmas, por exemplo as *threads* significa que necessitamos apenas de um espaço de endereçamento evitando desta forma a comutação de contexto I/O, criar novo espaço de endereçamento tomando desta forma tudo mais simples.

16. Indique, justificadamente, em que situações um ambiente *multithreaded* pode ser vantajoso.

R: *Vantagens de um ambiente multithreaded*: • As vantagens de um ambiente multithread são evidentes, na medida em que precisamos de menos contexto I/O menos espaço de endereçamento...

*maior simplicidade na decomposição da solução e maior modularidade na implementação* - programas que envolvem múltiplas actividades e atendimento a múltiplas solicitações são mais fáceis de conceber e mais fáceis de implementar numa perspectiva concorrential do que numa perspectiva puramente sequencial; • *melhor gestão de recursos do sistema computacional* - havendo uma partilha do espaço de endereçamento e do contexto de I/O entre os *threads* em que uma aplicação é dividida, torna-se mais simples gerir a ocupação da memória principal e o acesso aos dispositivos de entrada / saída de uma maneira eficaz; • *eficiência e velocidade de execução* - uma decomposição da solução em *threads* por oposição a processos, ao envolver menos recursos por parte do sistema de operação, possibilita que operações como a sua criação e destruição e a mudança de contexto, se tornem menos pesadas e, portanto, mais eficientes; além disso, em multiprocessadores simétricos torna-se possível calendarizar para execução em paralelo múltiplos *threads* da mesma aplicação, aumentando assim a sua velocidade de execução.

17. Que tipo de alternativas pode o sistema de operação fornecer à implementação de um ambiente *multithreaded*? Em que condições é que num multiprocessador simétrico os diferentes *threads* de uma mesma aplicação podem ser executados em paralelo?

R: Na utilização de diversos processos, a implementação dos *threads* possui uma vantagem devido a simplicidade de implementação e não há alternativas a este nível, mas sim, em outras circunstâncias, paga-se sempre algo, aqui é simplificada e velocidade pelo uso de um bloqueio estragar tudo. Se o multiprocessador partilhar memória e I/O cada thread pode correr num processador distinto, denotar a vantagem desta implementação principalmente a nível das *threads* de gestão de *scheduling* do sistema operativo na medida que o bloqueio de uma não necessita de mudar o contexto (bloqueio de todo o sistema multiprocessador). Num multiprocessadores simétricos torna-se possível calendarizar para execução em paralelo múltiplos *threads* da mesma aplicação, aumentando assim a sua velocidade de execução.

18. Explique como é que os *threads* são implementados em Linux.

R: O Linux lida com a questão da implementação de *threads* de um modo muito artificial. A parte da chamada ao sistema *fork* que cria um novo processo a partir dum já existente por cópia integral do seu contexto alargado (espaço de endereçamento, contexto de I/O e contexto do processador), existe uma outra, *clone*, que cria um novo processo a partir de um já existente por cópia apenas do seu contexto restrito (contexto do processador), partilhando o espaço de endereçamento e o contexto de I/O e iniciando a sua execução pela invocação de uma função que é passada como parâmetro.

Assim, não há distinção efectiva entre processos e *threads*, que o Linux designa indiferentemente de *tasks*, e eles são tratados pelo *kernel* da mesma maneira.

19. O principal problema da implementação de *threads*, a partir de uma biblioteca que fornece primitivas para a sua criação, gestão e *scheduling* no nível utilizador, é que quando um *thread* particular executa uma chamada ao sistema bloqueante, todo o processo é bloqueado, mesmo que existam *threads* que estão prontos a serem executados. Será que este problema não pode ser minimizado?

R: Para se minimizar este problema seria necessário criar subtabelas na TCP que permitiam ao *scheduling* do *kernel* comutar entre as *threads* salvaguardando apenas o contexto do processador. Mas, para tal seriam necessárias alterações ao nível da definição de prioridades e de *scheduling* de forma a otimizar a utilização da principal vantagem das *threads* a nível de gestão do processador: a velocidade de comutação devido a comutação de contexto ser limitada.

20. Num ambiente *multithreaded* em que os *threads* são implementados no nível utilizador, vai haver um par de *stacks* (sistema / utilizador) por *thread*, ou um único par comum a todo o processo? E se os *threads* forem implementados ao nível do *kernel*? Justifique.

R: Sendo o *stack* pertencente ao espaço de endereçamento apenas um par deverão existir, mas poderá ser ou não vantajoso a nível do utilizador haver mais do que essa *stack*, ou a *thread* em causa lança uma nova *stack* dentro desse espaço comum (o que é complicada) ou o sistema em si possibilita a criação de *stack* em cada *thread*. Na minha opinião, uma *stack* única é o mais racional, a *thread* não é um processo, possui os seus pontos fortes e fracos e deve ser utilizada dessa forma. A nível do *kernel*, a situação é diferente, visto que, não é lançada pelo utilizador, logo poderá ser implementada outra filosofia atendendo aos privilégios que possui.

### Capítulo III - Comunicação entre Processos

1. Como caracteriza os processos em termos de interacção? Mostre como em cada categoria se coloca o problema da exclusão mútua.

R: Num ambiente multiprogramado, os processos que coexistem podem ter comportamentos diversos em termos de interacção. Constituem-se como: • *processos independentes* - quando são criados, têm o seu 'tempo de vida' e terminam sem interagirem de um modo explícito; a interacção que ocorre é implícita e tem origem na sua *competição* pelos recursos do sistema computacional; trata-se tipicamente dos processos lançados pelos diferentes utilizadores num ambiente interactivo e/ou dos processos que resultam do processamento de *jobs* num ambiente de tipo *batch*; • *processos cooperantes* - quando partilham informação ou comunicam entre si de um modo explícito; a *partilha* exige um espaço de endereçamento comum, enquanto que a *comunicação* pode ser feita tanto através da partilha de um espaço de endereçamento, como da existência de um canal de comunicação que interliga os processos intervenientes.

Em cada categoria o problema da exclusão mútua coloca-se da seguinte forma:

• *processos independentes* que competem por acesso a um recurso comum do sistema computacional; • é da responsabilidade do SO garantir que a atribuição do recurso seja feita de uma forma controlada para que não haja perda de informação; • isto impõe, em geral, que só um processo de cada vez pode ter acesso ao recurso (*exclusão mútua*); • Processos cooperantes que partilham informação ou comunicam entre si: • é da responsabilidade dos processos envolvidos garantir que o acesso à região partilhada seja feito de uma forma controlada para que não haja perda de informação; • isto impõe, em geral, que só um processo de cada vez pode ter acesso à região partilhada (*exclusão mútua*); • o canal de comunicação é tipicamente um recurso do sistema computacional, logo o acesso a ele está enquadrado na *competição* por acesso a um recurso comum.

2. O que é a competição por um recurso? Dê exemplos concretos de competição em, pelos menos, duas situações distintas.

Quando dois recursos necessitam de um recurso em causa eles competem pela posse resultando no bloqueio de um e a continuação de um deles. São exemplos disso o jantar dos filósofos quando dois deles competem pelo mesmo garfo ou quando por exemplo duas renas a tentar competir por escrever o seu ID de forma a ????? por exemplo.

3. Quando se fala em região crítica, há por vezes alguma confusão em estabelecer-se se se trata de uma região crítica de código, ou de dados. Esclareça o conceito. Que tipos de propriedades devem apresentar as primitivas de acesso com exclusão mútua a uma região crítica?

Uma região crítica é código que permite o acesso a região efectiva. É preciso estabelecer condições para aceder a esta região crítica, são estas: • garantia efectiva de imposição de exclusão mútua - o acesso à região crítica associada a um mesmo recurso, ou região partilhada, só pode ser permitido a um processo de cada vez, de entre todos os processos que competem pelo acesso; • independência da velocidade de execução relativa dos processos intervenientes, ou do seu número - nada deve ser presumido acerca destes factores; • um processo fora da região crítica não pode impedir outro de lá entrar; • não pode ser adiada indefinidamente a possibilidade de acesso à região crítica a qualquer processo que o requeira; • o tempo de permanência de um processo na região crítica é necessariamente finito.

4. Não havendo exclusão mútua no acesso a uma região crítica, corre-se o risco de inconsistência de informação devido à existência de condições de corrida na manipulação dos dados internos a um recurso, ou a uma região partilhada. O que são condições de corrida? Exemplifique a sua ocorrência numa situação simples em que coexistem dois processos que cooperam entre si.

Condição de corrida é quando no acesso a zona partilhada em que a leitura para teste e posterior escrita são inconsistentes. Por exemplo:

P1 lê região crítica P1

P2 lê região crítica P2

P1 escreve região crítica P2

P2 testa e escreve na região crítica P2

Logo, P2 não escreve de acordo com o esperado

5. Distinga *deadlock* de adiamento indefinido. Tomando como exemplo o problema dos produtores / consumidores, descreva uma situação de cada tipo.

*Deadlock* é quando dois ou mais processos ficam a aguardar eternamente o acesso às regiões críticas respectivas, esperando acontecimentos que, se pode demonstrar, nunca irão acontecer; o resultado é, por isso, o bloqueio das operações; enquanto que o adiamento indefinido é quando um ou mais processos competem pelo acesso a uma região crítica e, devido a uma conjunção de circunstâncias em que surgem continuamente processos novos que o(s) ultrapassam nesse designio, o acesso é sucessivamente adiado; está-se, por isso, perante um impedimento real à continuação dele(s). No caso do problema produtores/consumidores, se os produtores estiverem sempre com acesso à região crítica teremos *starvation* perante os consumidores. (um ou mais)

6. A solução do problema de acesso com exclusão mútua a uma região crítica pode ser enquadrada em duas categorias: soluções *dis software* e soluções *dis hardware*. Quais são os pressupostos em que cada uma se baseia?

*Soluções software* - são soluções que, quer sejam implementadas num monoprocessador, quer num multiprocessador com memória partilhada, supõem o recurso em última instância ao *conjunto de instruções básico* do processador; ou seja, as instruções de transferência de dados de e para a memória são de tipo *standard*; leitura e escrita de um valor; a única suposição adicional diz respeito ao caso do multiprocessador; em que a tentativa de acesso simultâneo a uma mesma posição de memória por parte de diferentes processadores é necessariamente serializada por intervenção de um árbitro;

*Soluções hardware* - são soluções que supõem o recurso a instruções especiais do processador para garantir, a algum nível, a atomicidade na leitura e subsequente escrita de uma mesma posição de memória; são muitas vezes suportadas pelo próprio sistema de operação e podem mesmo estar integradas na linguagem de programação utilizada.

7. Dijkstra propôs um algoritmo para estender a solução de Dekker a N processos. Em que consiste este algoritmo? Será que ele cumpre todas as propriedades desejáveis? Porquê?

A solução que Dijkstra propôs baseia-se na alternância proposta por Dekker é de facto uma generalização em que os pressupostos são: a definição de quem entra primeiro é aleatório, o problema é que ao generalizarmos para N, a escolha sucessiva vai sendo aleatória e o facto do processo estar muito tempo à espera não ajuda a ser atendido mais rapidamente, logo adiamento indefinido.

8. O que é que distingue a solução de Dekker da solução de Peterson no acesso de dois processos com exclusão mútua a uma região crítica? Porque é que uma é passível de extensão a N processos e a outra não (não é conhecido, pelo menos até hoje, qualquer algoritmo que o faça).

O que o distingue é que o algoritmo de Peterson usa a serialização por ordem de chegada para resolver o conflito resultante da situação de contenção de dois processos. Isto é conseguido forçando cada processo a escrever a sua identificação na mesma variável. Assim, uma leitura subsequente permite por comparação determinar qual foi o último que ali escreveu. Enquanto que o de Dekker resolve o conflito de acesso entre dois processos competidores usando um mecanismo de alternância estrita. O algoritmo de Peterson é passível de uma extensão em N porque os processos em causa há mais tempo estarão mais perto do acesso logo não há adiamento indefinido, e no caso de Dekker cada processo ao ser eliminado recomeça a escada, o que pode provocar adiamento indefinido.

9. O tipo de fila de espera que resulta da extensão do algoritmo de Peterson a N processos, é semelhante àquela materializada por nós quando aguardamos o acesso a um balcão para pedido de uma informação, aquisição de um produto, ou obtenção de um serviço. Há, contudo, uma diferença subtil. Qual é ela?

A diferença é que a fila em vez de saber quando alguém se despacha sabe quando alguém chega:

1					
2	1				
	2	1			

5	4	3	2	1	
---	---	---	---	---	--

10. O que é o *busy waiting*? Porque é que a sua ocorrência se torna tão indesejável? Haverá algum caso, porém, em que não é assim? Explique detalhadamente.

O *busy waiting* (pooling) acontece quando um processo está continuamente a testar uma condição até que ela seja satisfeita. Este facto é naturalmente indesejável em sistemas computacionais monoprocessador; já que conduz a • *perda de eficiência* - a atribuição do processador a um processo que pretende acesso a uma região crítica, associada a um recurso ou a uma região partilhada em que um segundo processo se encontra na altura no seu interior, faz com que o intervalo de tempo de atribuição do processador se esgote sem que qualquer trabalho útil tenha sido realizado; • *constrangimentos no estabelecimento do algoritmo de scheduling* - numa política *preemptive* de *scheduling* onde os processos que competem por um mesmo recurso, ou partilham uma mesma região de dados, têm prioridades diferentes, existe o risco de *deadlock* se for possível ao processo de mais alta prioridade interromper a execução do outro.

Mas em sistemas multiprocessador o *busy waiting* é muito eficaz na medida em que cada processador é decutido e o *busy waiting* é vantajoso relativamente a um bloqueio

11. O que são *flags de locking*? Em que é que elas se baseiam? Mostre como é que elas podem ser usadas para resolver o problema de acesso a uma região crítica com exclusão mútua.

Se garantirmos a atomicidade de uma operação do tipo test-and-set (tas) estamos a resolver o problema na medida em que este tipo de operação usam variáveis chamados *flags* de *locking* que após serem testados com sucesso, mudam o seu estado, qualquer processo que a seguir tentar entrar nessa região entrará em *busy waiting* ate que o processo que fez o *locking* faça unlock da flag de *locking*.

12. O que são *semaforos*? Mostre como é que eles podem ser usados para resolver o problema de acesso a uma região crítica com exclusão mútua e para sincronizar processos entre si.

Um *semaforo* é um dispositivo de sincronização, originalmente inventado por Dijkstra, que pode ser concebido como uma variável do tipo

**typedef struct**

{ unsigned int val; /\* valor de contagem \*/

  NODE \*queue; /\* fila de espera dos processos bloqueados \*/

} SEMAPHORE;

sobre a qual é possível executar as duas operações atómicas seguintes:

*sem\_down* - se o campo val for não nulo, o seu valor é decrementado; caso contrário, o processo que executou a operação é bloqueado e a sua identificação é colocada na fila de espera queue;

*sem\_up* - se houver processos bloqueados na fila de espera queue, um deles acordado (de acordo com uma qualquer disciplina previamente definida); caso contrário, o valor do campo val é incrementado.

Um *semaforo* só pode ser manipulado desta maneira e é precisamente para garantir isso que toda a qualquer referência a um *semaforo* particular é sempre feita de uma forma indirecta.

13. O que são *monitores*? Mostre como é que eles podem ser usados para resolver o problema de acesso a uma região crítica com exclusão mútua e para sincronizar processos entre si.

Monitores são paradigmas implementados ao nível da linguagem de programação que permitem a actuação em variáveis de controlo de modo simples capazes de serem sincronizados através de duas operações:

*wait* - o *thread* que invoca a operação é bloqueado na variável argumento e é colocado fora do monitor para possibilitar que aguarda acesso, possa prosseguir;

*signal* - se houver *threads* bloqueados na variável de condição deles é acordado; caso contrário, nada acontece.

As duas operações implementam soluções de sincronização e por outro lado podem ser usadas para entra/saída de regiões críticas visto que quando um *thread* entra no monitor é feita em exclusão mútua, logo o acesso ao monitor impõe as condições de acesso a uma região partilhada. No que toca à sincronização temos duas operações fundamentais: *wait* e *signal*

14. Que tipos de monitores existem? Distinga-os.

Existem três tipos de monitores que se distinguem pela saída do monitor do *thread* após o sinal:

monitor de Hoare - o *thread* que invoca a operação de *signal* é colocado fora do monitor para que o *thread* acordado possa prosseguir; é muito geral, mas sua implementação exige a existência de um *stack*, onde são colocados os *threads* postos fora do monitor por invocação de *signal*;

monitor de Brinch Hansen - o *thread* que invoca a operação de *signal* liberta imediatamente o monitor (*signal* é a última instrução executada); é simples de implementar, mas pode tornar-se bastante restritivo porque só há possibilidade de execução de um *signal* em cada invocação de uma primitiva de acesso;

monitor de Lamport / Redell - o *thread* que invoca a operação de *signal* prossegue a sua execução, o *thread* acordado mantém-se fora do monitor compete pelo acesso a ele, é simples de implementar, mas pode originar situações em que alguns *threads* são colocados em adiamento indefinido.

15. Que vantagens e inconvenientes as soluções baseadas em monitores trazem sobre soluções baseadas em *semaforos*?

• *vantagens*: - suporte ao nível do sistema de operação - porque a sua implementação é feita pelo *kernel*, as operações sobre *semaforos* estão directamente disponíveis ao programador de aplicações, constituindo-se como uma biblioteca de chamadas ao sistema que podem ser usadas em qualquer linguagem de programação; - universalidade - são construções de muito baixo nível e podem, portanto, devido à sua versatilidade, ser usadas no desenho de qualquer tipo de soluções;

• *desvantagens*: - conhecimento especializado - a sua manipulação directa exige ao programador um domínio completo dos princípios da programação concorrente, pois é muito fácil cometer erros que originam condições de corrida e, mesmo, *deadlock*.

16. Em que é que se baseia o paradigma da passagem de mensagens? Mostre como se coloca o problema do acesso com exclusão mútua a uma região crítica e como ele está implicitamente resolvido?

O paradigma da passagem de mensagens não envolve partilha de espaço de endereçamento e a sua aplicação, de um modo mais ou menos uniforme, é igualmente válida tanto em ambientes multiprocessador, como em ambientes multiprocessador, ou de processamento distribuído. Apenas se implementa um canal de comunicação de forma que processos cooperantes troquem informações entre si. Perante a possibilidade de troca de informações sem acesso a essa região está resolvida a questão, pois a troca de mensagens, e por inerência de informação, está resolvida, sendo desnecessário a utilização da região partilhada.

17. O paradigma da passagem de mensagens adequa-se de imediato à comunicação entre processos. Pode, no entanto, ser também usado no acesso a uma região em que se partilha informação. Conceba uma arquitectura de interacção que torne isto possível.

Para implementar uma arquitectura que use este tipo de paradigma pode ser feito de duas formas, uma estrutura central (processo) controle o acesso informando os subítos de se podem ou não entrar na região através de após o pedido enviar uma mensagem (ficando o processo que pediu acesso em bloqueio na recepção da resposta). Outra forma seria um sistema perguntar a toda a gente se está lá alguém a entrar. Quando receber resposta, cada ser avisa toda gente que sai, perde-se rentabilização na medida que tudo fica à espera que saia em sincronização bloqueante.

18. Que tipo de mecanismos de sincronização podem ser introduzidos nas primitivas de comunicação por passagem de mensagens? Caracterize os protótipos das funções em cada caso (suponha que são escritas em Linguagem C).  
Podemos utilizar um sistema do tipo utilizando no jantar dos filósofos em que usamos duas funções uma bloqueante e outra desbloqueante com vista a sincronização.  
O grau de sincronização existente pode ser classificado em dois níveis:  
• sincronização não bloqueante - quando a sincronização é da responsabilidade dos processos intervenientes; a operação de envio envia a mensagem e regressa sem qualquer informação sobre se a mensagem foi efectivamente recebida; a operação de recepção, por seu lado, regressa independentemente de ter sido ou não recebida uma mensagem;  
/\* operação de envio \*/  
void msg\_send\_nb (unsigned int destid, MESSAGE msg);  
/\* operação de recepção \*/  
void msg\_receive\_nb (unsigned int srcid, MESSAGE \*msg, BOOLEAN \*msg\_arriaval);

sincronização bloqueante - quando as operações de envio e de recepção contém em si mesmas elementos de sincronização; a operação de envio envia a mensagem e bloqueia até que esta seja efectivamente recebida; a operação de recepção, por seu lado, só regressa quando uma mensagem tiver sido recebida;  
/\* operação de envio \*/  
void msg\_send (unsigned int destid, MESSAGE msg);  
/\* operação de recepção \*/  
void msg\_receive (unsigned int srcid, MESSAGE \*msg);

19. O que são sinais? O *standard Posix* estabelece que o significado de dois deles é mantido em aberto para que o programador de aplicações lhes possa atribuir um papel específico. Mostre como poderia garantir o acesso a uma região crítica com exclusão mútua por parte de dois processos usando um deles.  
*Sugestão* - Veja no manual *on-line* como se pode usar neste contexto a chamada ao sistema *wait*.

Sinais são mecanismos que o sistema põe à disposição do kernel, de qualquer processo ou utilizador que após a sua activação despoleta uma rotina de serviço (embora pode haver vários tipos de resposta a sinais) que terá de ignorar até resposta imediata, chamadas ao sistema.  
Um sinal constitui uma interrupção produzida no contexto de um processo onde lhe é comunicada a ocorrência de um acontecimento especial. Pode ser despoletado pelo kernel, em resultado de situações de erro ao nível hardware ou software, pelo próprio processo, por outro processo, ou pelo utilizador através do dispositivo standard de entrada / saída.  
Tal como o processador no tratamento de excepções, o processo assume uma de três atitudes possíveis relativamente a um sinal  
• ignorá-lo - não fazer nada face à sua ocorrência;  
• bloqueá-lo - impedir que interrompa o processo durante intervalos de processamento bem definidos;  
• executar uma acção associada - pode ser a acção estabelecida por defeito quando o processo é criado (conduz habitualmente à sua terminação ou suspensão de execução), ou uma acção específica que é introduzida (registada) pelo próprio processo em runtime.  
Uma forma é de enviar cada sinal de entrada na região crítica e saída, seria qualquer programa que tente aceder posta em wait.

20. Mostre como poderia criar um canal de comunicação bidireccional (*full-duplex*) entre dois processos parentes usando a chamada ao sistema *pipe*.  
21. Como classifica os recursos em termos do tipo de apropriação que os processos fazem deles? Neste sentido, como classificaria o canal de comunicações, uma impressora e a memória de massa?  
Um recurso é algo que um processo precisa para a sua execução. Os recursos tanto podem ser componentes físicos do sistema computacional (processadores, regiões de memória principal ou de memória de massa, dispositivos concretos de entrada / saída, etc.), como estruturas de dados comuns definidas ao nível do sistema de operação (tabela de controlo de processos, canais de comunicação, etc.), ou entre processos de uma mesma aplicação.  
Os recursos dividem-se em: **recursos preemptible** - quando podem ser retirados aos processos que os detêm, sem que daí resulte qualquer consequência irreparável à boa execução dos processos; são, por exemplo, em ambientes multiprogramados, o processador, ou as regiões de memória principal onde o espaço de endereçamento de um processo está alojado; **recursos non-preemptible** - em caso contrário; são, por exemplo, a impressora, ou uma estrutura de dados partilhada que exige exclusão mútua para a sua manipulação.

22. Distinga as diferentes políticas de prevenção de *deadlock* no sentido estrito. Dê um exemplo ilustrativo de cada uma delas numa situação em que um grupo de processos usa um conjunto de blocos de disco para armazenamento temporário de informação.  
As políticas de prevenção de *deadlock* no sentido estrito, embora seguras, são muito restritivas, pouco eficientes e difíceis de aplicar em situações muito gerais.  
Resumindo, tem-se que: **negação da condição de exclusão mútua** - só pode ser aplicada a recursos passíveis de partilha em simultâneo; **negação da condição de espera com retenção** - exige o conhecimento prévio de todos os recursos que vão ser necessários, considera sempre o pior caso possível (uso de todos os recursos em simultâneo); **imposição da condição de não libertação** - ao supor a libertação de todos os recursos anteriores quando o próximo não puder ser atribuído, atrasa a execução do processo de modo muitas vezes substancial; **negação da condição de espera circular** - desaproveita recursos eventualmente disponíveis que poderiam ser usados na continuação do processo.

23. Em que consiste o algoritmo dos banqueiros de Dijkstra? Dê um exemplo da sua aplicação na situação descrita na questão anterior.  
O algoritmo dos banqueiros de Dijkstra aplica-se à detecção de estados inseguros, ora a apropriação de um estado seguro necessário apenas se verifica que não há risco de *deadlock*. No caso anterior há que considerar o problema circular de *deadlock* na apropriação. Exemplo:

- cada filósofo, quando pretende os garfos, continua a pegar primeiro no garfo da esquerda e, só depois, no da direita, mas agora mesmo que o garfo da esquerda esteja sobre a mesa, o filósofo nem sempre pode pegar nele;
- só o poderá fazer se pelo menos um dos outros filósofos estiver a meditar; caso contrário, a atribuição origina uma situação insegura;
- não se trata, contudo, de uma solução geral (o adiamento indefinido não está completamente resolvido).

24. As políticas de prevenção de *deadlock* no sentido lato baseiam-se na transição do sistema entre estados ditos *seguros*. O que é um estado seguro? Qual é o princípio que está subjacente a esta definição?  
Define-se neste contexto estado seguro como uma qualquer distribuição dos recursos do sistema, livres ou atribuídos aos processos que coexistem, que possibilita a terminação de todos eles. Por oposição, um estado é inseguro se não for possível fazer-se uma tal afirmação sobre ele.

25. Que tipos de custos estão envolvidos na implementação das políticas de prevenção de *deadlock* nos sentidos estrito e lato? Descreva-os com detalhe.  
As políticas de prevenção de *deadlock* no sentido estrito, embora seguras, são muito restritivas, pouco eficientes e difíceis de aplicar em situações muito gerais, enquanto que as políticas de prevenção de *deadlock* no sentido lato, embora igualmente seguras, não são menos restritivas e ineficientes do que as políticas de prevenção de *deadlock* no sentido estrito. São, porém, mais versáteis e, por isso, mais fáceis de aplicar em situações gerais.  
Mas convém notar o seguinte: **é necessário o conhecimento completo de todos os recursos do sistema e cada processo tem que indicar à cabeça a lista de todos os recursos que vai precisar** - só assim se pode caracterizar um estado seguro;  
**um estado inseguro não é sinónimo de *deadlock*** - vai, contudo, considerar-se sempre o pior caso possível para garantir a sua não ocorrência.  
Pelo que os custos são maiores no sentido lato porque precisa-se de hardware e maior capacidade de processamento enquanto que no sentido estrito os custos são a nível de software.

**MEMORIA VIRTUAL**em memória virtual é constituído em disco uma imagem do processo. A informação esta dividida em blocos. Se esses blocos têm tamanho fixo então o **sistema é paginado**, e os blocos designam-se por páginas. Caso contrário se os blocos tem tamanho variável designam-se por segmentos, e os sistemas q os usam chamam-se **segmentados**.  
**SISTEMA PAGINADO**-> o tamanho dos blocos é fixo pelo que a mem atribuída a um dado processo pode ser constituída por vários blocos de informação. É mantido um registo, para cada processo, que indica quais os blocos que se encontram em mem real. Os endereços virtuais destes blocos contém o número do bloco em questão e o deslocamento a partir do início deste bloco. Para além disso o CPU matém um registo base onde está o endereço do início da tabela de paginação. A tradução do endereço virtual em endereço real é feita do seguinte modo: ao endereço base é somado o número do bloco. Do endereço resultante fica-se a saber a entrad da tabela de blocos do processo correspondente aquele bloco. Esta entrada dá o endereço em memória principal do bloco em questão. Se a estre endereço somarmos o delocamento termos o em end moria real. Para tornar mais rápida esta tradução a tabela de paginação pode ser armazenada em mem cache ou associativa. Para haver protecção das pagoinas de memória de diferentes processos é necessário acrescentar masi info à tabela de paginação.Há então um registo com o numero total de paginas quer o processo tem. Se o numero da ppag somado com o endereço base der um valor superior ao jumeor xda pagã há memory fault.  
**SISTEMA SEGMENTADO**-> os blocos são feitos tão grandes como necessários, conforme o tamanho do processo e a informação associada. Este mapeamento é mais natural pois cada segmento corresponde a uma entidade conceptual, lógica e não física como no sistemas paginado. A protecção da info e a sua partilha é assim facilitada. Existe um certo paralelismo entre estes dois tipos de organização de mem virtual e os sistemas de partições fixas e variáveis em sistemas de memoria real.

**Descreva a organização de um sistema hierarquizado de memória. Justifique a importância dos diferentes níveis num sistema multiprogramado.**R: Um sist hierarquizado de mem supõe 3 níveis: memória cache, principal, e de massa. A df entre cada nível esta relacionada com a capacidade e de tempo de acesso, que são mínimos na cache e máximos na memória principal. O processador apenas tem acesso directo à cache e à memória principal. Num sistema multiprogramado diferentes programas estão, em princípio, a ser concurentemente executados, o que impõe a sua carga total ou parcial em mem principal. Para que a capacidade deste tipo de mem não constitua um impedimento ao numero máximo total de processos que coexistam, torna-se muitas vezes necessário deslocar temporariamente alguns daqueles, total ou parcialmente, para a memoria de massa. Abre-se assim a possibilidade de execução de um numero de programas cujas necessidades em armazenamento ultrapassa a capacidade da mem principal. A mem cache sendo de acesso mais rápido que a mem principal deve assim ser utilizada para armazenar partes de código ou estruturas de dados mais frequentemente utilizadas.

**Distinga uma organizacao de mem real de uma virtual. Qual a importância que organizacoes de mem virtual têm em sistemas do tipo interactivo?**R: A mem real é a mem do system que pode ser directamente referenciada pelo CPU (mem cache e principal). Uma das diferenças entre mem real e mem virtual é que em mem virtual é construída em disco uma imagem do processo. Assim, em cada instante carrega-se apenas uma parte do processo. Deste modo o conceito de mem virtual permite que o espaço de endereçamento de um processo seja muito superior ao espaço físico da mem(mem principal). No entanto, como processo precisa de ser executado em mem física a tradução de endereços virtuais para endereços reais deve ser rápida e eficiente.

**Descreva uma organizacao de mem real com particoes de tamanho variável. Indique como nela podem ser resolvidos os problemas de carga de processos e da protecao de areas de mem pertencentes a diferentes processos.**R: Numa organização deste tipo, a dimensão de cada partição é exactamente aquela necessária ao armazenamento contíguo do código e estruturas de dados próprios de cada processo. Faz-se portanto uma optimização do espaço de memória. Porém à medida que os processos existentes vão sendo terminados e outros processos vão sendo criados, a memória vai-se fragmentando e pode ocorrer uma situação em que, embora a memória total disponível seja superior à necessária não existe nenhuma zona de endereços contíguas com dimensão suficiente para carregar um novo programa, o que exige a implementação de um mecanismo de compactação *Garbage Collector*. De um modo geral, no entanto, existem sempre um ou mais blocos de endereços contíguas disponíveis onde o código e as estruturas de dados associados a um novo processo podem ser carregados. A selecção de um deles faz-se usualmente de acordo com os dois 3 critérios: **First-Fit** - escolhe-se o primeiro bloco com dimensão suficiente **Best - Fit** - escolhe-se o bloco com dimensão mais pequena **Worst-Fit** - escolhe-se op bloco com maior dimensão. A protecção da área de memória pertencente a cada processo é feita através da implementação de dois registos fronteira que contém os endereços limite da zona de memória principal contígua atribuída ao processo. Só os endereços gerados dentro desse intervalo são considerados válidos, todos os outros dão origem a uma interrupção de tipo *memory fault*.

**O que é que caracteriza o conceito de WORKING-SET como uma política de gestão do espaço de memória física atribuída a um processo.**R:Este conceito está relacionado com o conceito de localidade que nos sugere que os processos tendem a referenciar áreas de armazenamento segundo padrões bem localizados e uniformes quer no tempo quer no espaço físico da mem. Aproveitando-se desta propriedade inerente aos processos as políticas de VSP defendem que as paginas de mem referenciadas pelo processo devem ser mantidas em memórias primárias para otimizar a execução do mesmo.

**Que tipo de vantagens se procura retirar da implementação de uma política de spoling da impressora?** Esta política procura aumentar o throughput do sistema computacional, medindo a transferência de daos para ela, que é um periférico lento, através do recurso intercalar a um periférico rápido de armazenamento temporário(disco). **Qual a sobrecarga que um tal política impõe sobre os outros recursos do sistemas?** Impõe sobretudo a nssidade de aumento de capacidade de mem de armazenamento de massa. No entanto dado que tem de existir um novo processo (spooler) enquece a transferência de dados entre a mem de massa e a impressora, tem tb implicações em termos da mem principal.

**Princípio da optimidade** - Para se obter performance optima deve-se substituir o bloco que não vai ser usado durante mais mais no futuro. Esta estrategia nao pode ser usada ja que nao se pode prever o futuro.Na pratica pretende-se uma estrategia que se aproxime desdte principio.

**FIFO**- Permite substituir a pagina que está à mais tempo em meoria. É pouco eficiente porque a pag. Substituída pode ser uma que estava usada constantemente usada. Exist uma anomalia nesta estrategia: aumentar o filo atribuído a um processo não leva necessariamente à reducao de page fault(antes pelo contrario).

**LRU(LAST RECENT USED)**- baseia-se no principio da localidade, segundo o qual um processo passa períodos de tempo razoáveis, a executar zonas restritas de código. É um aboa aproximacao do principio de optimidade. No entanto obriga ao armazenamanto de informacao temporal relativa às execucoes das diversas pags. Pelo que apresenta um grande overhead.

**Gestão de Memória de Massa** Embora seja possível uma única cabeça de leitura e escrita, usa-se normalmente 2 cabeças separadas para se fazer a verificação dos dados que se escreve, logo após a escrita. O armazenamento dos dados não é feito de forma continua, para não ficarem perdidos devido à ocorrência de defeitos no material do disco. Antes pelo contrário, a informação é armazenada de forma fragmentada. Os discos encontram-se divididos em blocos (sectores) de 512 a 1024 bytes, cada sector tem o seguinte aspecto. Por questões de simplicidade, cada pista tem o mesmo nº de sectores. Para se aceder à informação indica-se o nº da pista e o nº do sector. Num floppy disc há contacto físico entre a superfície do disco e a cabeça de leitura, pois não há variações de entreferro. Devido a isso estes discos são lentos e têm um tempo de vida limitado. (desgaste por contacto). Se o disco é rígido o entreferro, disco-cabeça de leitura é constante. Assim evita-se o contacto, permitindo velocidades de rotação elevadas.

O acesso ao disco é feito indicando o nº do cilindro, nº da superfície e o nº do sector.

**Tempo de Latência** tempo necessário para se atingir o sector desejado na pista em que a cabeça se encontra (está limitada pela velocidade de rotação).

**Tempo de Pesquisa** tempo necessário para que a cabeça atinja o cilindro desejado. Como o tempo de pesquisa é bastante superior ao tempo de latência, o uso eficiente do disco rígido obriga à optimização das operações de pesquisa.

**Características desejáveis numa estratégia de gestão do disco**

- Troughput (nº de pedidos por unidade de tempo)
- Tempo médio de resposta
- Variância dos tempos de resposta (previsibilidade)

**Estratégia de optimização de Pesquisa**

**FCFS (First Come First Served)** é uma estratégia justa. No entanto resulta em padrões de busca aleatória pelo que a previsibilidade do sistema é má e depende da carga do sistema.

**SST (Smallest Seek Time)** é sempre servido o pedido de acesso ao bloco que esteja mais próximo da cabeça. Melhora-se o tempo de acesso mas a previsibilidade é má já que o pedido para acesso a blocos nas extremidades e no centro do disco são desconsideradas.

**Scan** a cabeça de leitura desloca-se da extremidade para o centro do disco, e então regressa à extremidade, repetindo sempre este movimento. Os pedidos que chegam são ordenados por distância à cabeça na direcção de deslocamento (os mais próximos em primeiro lugar). Esta estratégia reduz a discriminação da zona central e das extremidades. No entanto a zona intermédia continua a ser privilegiada. Melhora a previsibilidade. Para baixas cargas é a melhor estratégia.

**N-Step** Scan os movimentos da cabeça são semelhantes aos do scan, mas durante o deslocamento num sentido só são acedidos os pedidos existentes aqando do início desse deslocamento. Os pedidos que chegam entretanto, são ordenados para serviço óptimo, quando o deslocamento for em sentido contrário. É uma boa solução com bom Troughput e tempo de resposta. Apresenta uma melhor previsibilidade que o scan. No entanto continua a haver privilégios da zona intermédia do disco. Evita adiamento indefinido de pedidos de acesso.

**Circular Scan** a cabeça desloca-se do exterior para o centro do disco, servindo por ordem de proximidade os pedidos que vão chegando. Uma vez chegado ao centro, retorna rapidamente ao exterior sem servir nenhum pedido, i.é só são servidos pedidos no deslocamento para o interior. É semelhante ao Scan. Com este método deixa de haver zonas favorecidas do disco. Se os pedidos, que chegam durante um deslocamento para o centro só forem servidos no próximo deslocamento para o centro, temos N-Step C-Scan. Com o C-Scan obtém-se uma excelente previsibilidade, para elevadas cargas é a melhor estratégia.

**Para baixas cargas a melhor estratégia é o Scan, enquanto que para elevadas é o C-Scan.**

**Organização Lógica dos Ficheiros - File System** do ponto de vista do utilizador o disco está dividido em directórios organizados numa estrutura hierárquica. No entanto no disco a informação é organizada em sectores. Assim, acima do gestor básico do disco, que faz o scheduling e optimização dos acessos, é necessário um gestor de disco que funcione ao nível lógico. Os tipos de objectos são: Directórios Ficheiros (endereçamento directo). Um Link é uma entrada do directório que aponta para um objecto existente noutro ponto da estrutura hierárquica de ficheiros. Esse objecto não é efectivamente copiado no directório em que existe o link. Os tipos básicos de acesso são: Leitura, Escrita e Execução. Os utilizadores são agrupados em grupos com determinadas permissões de acesso aos objectos do file system. Dentro de cada grupo pode haver partilha de informação, pode tb existir partilha de informação entre grupos.

**Capitulo V - Gestão de dispositivos de entrada-saída**

¶ **Porque é que num ambiente multiprogramado é fundamental desacoplar a comunicação dos processos utilizadores com os diferentes dispositivos de entrada-saída?**

¶ **Como papel desempenham os 'device-drivers' no estabelecimento de um interface uniforme de comunicação com os dispositivos de entrada-saída?**

¶ **O que distingue dispositivos de tipo caracter de dispositivos de tipo bloco?**  
**Descreva de uma maneira funcional como se desenvolve a comunicação entre um processo utilizador e um dispositivo de tipo caracter.**

¶ **O que distingue dispositivos de tipo caracter de dispositivos de tipo bloco?**  
**Descreva de uma maneira funcional como se desenvolve a comunicação entre um processo utilizador e um dispositivo de tipo caracter.**

¶ **O que distingue dispositivos de tipo caracter de dispositivos de tipo bloco?**  
**Descreva de uma maneira funcional como se desenvolve a comunicação entre um processo utilizador e um dispositivo de tipo caracter.**

Dispositivos de Blocos: estes dispositivos armazenam informação em blocos de tamanho fixo, cada um numa localização bem definida. A principal propriedade destes dispositivos é a de permitir a leitura/escrita de cada bloco independentemente de todos os outros. Os discos e as tapes são dispositivos deste tipo.  
Dispositivos de Caracteres: estes dispositivos processam a informação em grupos de caracteres, sem ter em atenção alguma estrutura de tipo bloco. Estes dispositivos não são endereçáveis e não permitem operações de pesquisa. As impressoras, os terminais de video e os teclados são dispositivos deste tipo.

Os de caracter transmitem informação com mensagens com o tamanho fixo de um byte e a informação não pode ser pesquisada, os dispositivos de bloco transmitem a informação com mensagens com tamanho de um bloco com vários bytes e é possível aceder a um bloco independentemente dos outros, ou seja, permite a pesquisa indexada.

A figura abaixo descreve de uma maneira esquemática a situação proposta.

O processo P é o processo utilizador que em momentos particulares da sua execução pretende escrever na impressora linhas de texto. Para isso, socorre-se da chamada ao sistema

- void escreve\_linha (char \*linha)

que transfere para o buffer de comunicação BUFF\_INP a linha em questão.

O processo CI é um processo de sistema activado pela interrupção TxDty do registo de transmissão do controlador da impressora. Quando é acordado, este processo procura ler um carácter, pertencente a uma linha armazenada em BUFF\_INP, e escrevê-lo no registo de transmissão da impressora. Para isso socorre-se das duas chamadas ao sistema seguintes: - char le\_caracter () - void transm\_caracter (char car)

Em termos gerais, a interacção entre dois processos pode ser descrita da seguinte forma: Sempre que tiver um alinha para imprimir, o processo P invoca o procedimento escreve\_linha que transfere para o buffer de comunicação BUFF\_INP a linha em questão. O seu regresso é imediato, desde que haja espaço de armazenamento suficiente, caso contrário, o processo bloqueia. Na situação particular em que o buffer de comunicação esteja inicialmente vazio, o procedimento escreve\_linha realiza a acção suplementar de transferir o primeiro carácter para o registo de transmissão do controlador da impressora, procedimento termo, carácter, despoitando assim o mecanismo de interrupções que vai acordar sucessivamente o processo CI, até que toda a informação tenha sido efectivamente transferida.

Por sua vez o processo CI, quando é acordado pela interrupção, invoca a função le\_caracter para recolher um carácter previamente armazenado em BUFF\_INP e, caso exista algum escreve-o no registo de transmissão do controlador da impressora, procedimento transm\_caracter. Na situação particular em que o buffer de comunicação esteja inicialmente cheio a função le\_caracter realiza a acção suplementar de tentar acordar o processo P eventualmente bloqueado a aguardar espaço de armazenamento disponível no buffer.

¶ Durante a operação normal de um sistema computacional, o acesso à memória de massa é quase constante. Como o seu tempo de acesso é tipicamente ordens de grandeza mais lento do que o acesso à memória principal, há o risco do desempenho global do sistema de operação ser fortemente afectado pelo 'engarrafamento' que daí resulta. Apresente sugestões de como minimizar este efeito.

Utilização da DMA.

Políticas de limpeza de páginas em memória principal

O uso de dois buffers, um para páginas modificadas e delas não modificadas.

**Capítulo VI - Gestão da memória de massa**

¶ **Porque é que sobre a abstracção da memória de massa, concebida como um 'array' de blocos para armazenamento de dados, se torna fundamental introduzir o conceito de sistema de ficheiros? Quais são os elementos essenciais que definem a sua arquitectura?** O sistema de ficheiros permite abstrair o utilizador da organização da memória de massa, fornecendo-lhe apenas algumas funções de acesso aos ficheiros. (leitura, escrita, criação, eliminação, truncar, procura) Um sistema de ficheiros é composto por duas partes distintas: ficheiros e directórios (um directório é um ficheiro com atributos especiais).

Os ficheiros contêm os dados, os programas e toda a informação que se pretende armazenar.

Os directórios servem para estruturar a organização dos ficheiros. Um ficheiro é uma forma de guardar dados de uma forma organizada. Um ficheiro é caracterizado pelos seguintes atributos:

- Nome: Nome simbólico para o utilizador reconhecer o ficheiro
- Tipo: Se o ficheiro é do tipo objecto, executável, código fonte, batch, texto, arquivos, etc.
- Localização: É um ponteiro para o dispositivo onde o ficheiro se encontra e a localização dentro desse dispositivo.
- Tamanho: O tamanho do ficheiro (em bytes, palavras ou blocos), e possivelmente o tamanho máximo permitido.
- Protecção: Informação de acesso ao ficheiro por parte dos diversos utilizadores.
- Tempo, data e dono: Data de criação, modificação e utilização.

¶ **Em que consiste o problema da consistência da informação num sistema de ficheiros? Que razões levam à sua ocorrência? Como é que ele pode ser minimizado?** Tem haver com o problema resultante do facto de durante operações de criação, escrita, e outras operações que alterem os ficheiros, o sistema falhe por alguma razão e a operação não seja totalmente realizada, por exemplo, pode ser sinalizado que um dado bloco de dados está a ser utilizado por um ficheiro mas como a operação foi interrompida, não foi possível lá escrever a informação. Pode ser minizado fazendo regularmente uma verificação de consistência do sistema (scandisk), ou seja, percorrer todos os nós i e respectivos blocos de dados e verificar se estão correctamente. Uma solução utilizada também pode ser a realização de um diário, ou seja, sempre que se efectua uma operação que pode levar a uma situação de inconsistência do sistema, os vários passos dessa operação são guardados no diário e se no final forem bem executados, são sinalizados como tal, deste modo sempre que o sistema se desligue por uma razão anómala, basta ir ao diário e ver se no momento estava a ser efectuada alguma operação crítica e se foi bem realizada, se não foi finalizada o sistema realiza a operação inversa, ou seja, os passos que já tinham sido realizados são agora desfeitos.

¶ **Descreva justificadamente os diversos papéis desempenhados pela memória de massa num sistema computacional. Dê razões que expliquem porque é que ela deve ser constituída por mais do que uma unidade de disco magnético.** Armazenamento de dados (ficheiros e directórios). Área de swapping: extensão da memória principal, de forma a aumentar a taxa de utilização do processador e a libertar o programador da limitação da memória física, através da utilização da memória virtual. Para não se perderem dados importantes quando um disco falhar, usando vários discos com cópias iguais. Quando é necessário armazenar muita informação.