

## OS Interaction with Process

### Exceptions

- **Synchronous**
- 1. OS uses exception table to find exception handler (kernel mode)
- 2. OS executes exception handler (kernel mode)
- 3. User process to handle exception/terminate (user mode)

### Interrupts

- **Asynchronous.** Cause by external factors
- During interrupts, instructions being executed are atomic and must finish. After that instruction is executed, OS will check for interrupts again
- 1. Interrupt signal sent to OS (user mode)
- 2. Save current context (hardware)
- 3. Disable interrupt (hardware)
- 4. Hardware reads IVT (stored somewhere in kernel space) to find address of interrupt handler
- 5. OS invokes interrupt handler (kernel mode)
- 6. Return from interrupt (user mode)

### Syscalls

- **Synchronous**
- 1. Code invokes syscall (user mode)
- 2. Syscall sets up syscall number and syscall parameters, and invokes TRAP instruction, raising an exception to switch to kernel mode
- 3. Dispatcher is invoked with syscall number to find the correct syscall handler using a table similar to IVT, but software version (kernel mode)
- 4. Syscall handler is invoked (kernel mode)
- 5. Return from syscall (user mode)

## Process Abstraction in Unix

### fork()

- Creates an exact duplicate of the parent (except PID and PPID) with completely different addressing space. Global variables, pointers and address pointed to by pointers are copied as well
- Copy on write: read-only data is not copied. Only data that is to be written either by parent or child is copied by child

### exec()

- Replaces current process image with new process image. If the process is a child of a parent, it still remains as a child.

### wait()

- When a child exits before parent → child can become **zombie** process if parent has not called wait()
- Parent calls **wait()** to clean up child process (PID, status, CPU time)
- wait() only waits for direct children (not grandchildren)
- When parent exits before child → child becomes **orphan** process. "init" process becomes pseudo-parent. Child sends signal to "init" which cleans it up

## Inter-process Communication

### Shared memory (implicit)

1. Process A creates shared memory region M (OS)
  2. Process B attaches M to its own memory space (OS)
  3. Processes A and B can read and write into M
  4. Detach M from memory space after use (one of the processes)
  5. Destroy M if M is not attached to any process (one of the processes)
- OS only involved in **creating** and **attaching** M
  - Efficient, easy to use and fast
  - Limited to a single machine and requires synchronisation

### Message Passing (explicit)

- **Direct:** Sender and receiver explicitly name the other party (via PID)
  - Used when it matters which process receives the message
- **Indirect:** Messages are sent to/received from mailbox/port via syscalls
  - Used when it does not matter which processes received which messages in the mailbox/port
- Assumed **blocking** receiver: receiver needs to wait until it receives a message
- **Asynchronous:** sender is non-blocking
  - System Buffer buffers incoming messages at the receiver
  - Sender is blocked or returns with error when buffer is full
- **Synchronous:** sender is blocking
  - Sender is blocked until receiver calls receive()
  - No need for buffer
  - Message is kept by the sender until receiver copies it directly from sender
- Applicable beyond a single machine and easier to synchronise
- Inefficient (OS is heavily involved) and harder to use

### Unix Pipes

- Asynchronous message passing
- Behaves in FIFO order
- Writers wait when buffer is full and readers wait when buffer is empty

### Unix Signals

- Asynchronous notification regarding an event
- Messages sent to a running process to trigger specific behaviours
- Recipient of the signal must handle the signal by default set of signal handlers or custom defined handlers (some signal handlers cannot be overridden)

## Synchronisation

### Correct CS implementation

- **Mutual Exclusion:** At most one process in CS
- **Progress:** If no process in CS, one waiting process should enter CS
- **Bounded wait:** After a process requests to enter CS, there exists an upper bound on the number of other processes that can enter CS before it
- **Independence:** Process not executing in CS should never block other wanting to enter

## Spinlock/Busy-waiting

### Peterson's algorithm

```
Want[0] = 1;
Turn = 1;
while (Want[1] &&
      Turn == 1);

Critical Section

Want[0] = 0;
```

Process P0

```
Want[1] = 1;
Turn = 0;
while (Want[0] &&
      Turn == 0);

Critical Section

Want[1] = 0;
```

Process P1

- Want must be before Turn
  - P0: Turn = 1 → P1: Turn = 0 → P1: Want[1] = 1 → **P1: Enter CS** → P0: Want[0] = 1 → **P0: Enter CS**
- Wasteful: processes that are busy waiting still utilise CPU

### TestAndSet

- Atomic function that takes a memory address M, returns current content in M, then set content in M to 1

```
void EnterCS( int* Lock )
{
    while( TestAndSet( Lock ) == 1 );
}
```

```
void ExitCS( int* Lock )
{
    *Lock = 0;
}
```

- There's still busy waiting

## Semaphores

- Guaranteed no deadlocks if processes lock mutexes in the same order; the order of unlocks does not affect the correctness
- If two processes use 1 mutex, deadlocks are impossible

### (Unreusable) Barrier Semaphore

```
int arrived = 0;
Semaphore mutex = 1;
Semaphore waitQ = 0;

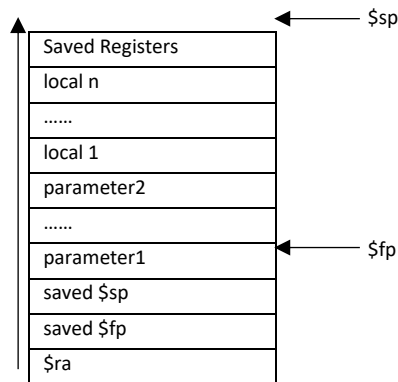
Barrier(N) {
    wait(mutex);
    arrived++;
    signal(mutex);
    if (arrived == N) {
        signal(waitQ);
    }
    wait(waitQ);
    signal(waitQ); // final waitQ >= 1
}
```

## Operating Systems

- **Monolithic kernel:** One big program; good performance
- **Microkernel:** Small, containing only basic and essential facilities. Excludes device drivers, file systems etc. Provides better isolation and protection and more robust
- **Type 1 Hypervisors:** Runs on top of hardware and must implement entire interface to hardware. Harder to implement but faster than Type 2.
- **Type 2 Hypervisors:** Runs on top of host OS and does not communicate with hardware directly. Easier to implement but slower than Type 1.

## Process Abstraction

### Stack Frames



- **Memory** context: Text, Data, Stack, Heap
- **Hardware** context: GPRs, PC, \$sp, \$fp
- **Set-up:** Layout and size of stack frame is determined at compile time → address of local variables and parameters will always have the same (relative) address. Reserve space for stack by incrementing \$sp
- main's stackframe does not contain saved registers → invoked by OS, registers are saved in PCB
- **Frame Pointer:** Not necessary, but is there for convenience of compiler. Allow usage of constant offset from \$fp to address variables on the stack
- **Stack Pointer:** Can change during run-time, due to VLAs
- **Teardown:** Restore saved registers, \$sp and \$fp. Memory previously allocated for the stack is not deallocated (tearing down is just restoring register values)

## Dynamically Allocated Memory

- Allocated in run time. Resides in **heap** memory region
- Cannot place in Data region: size unknown in compile time
- Cannot place in Stack region: need to persist after function returns

## Process Model

- **OS** context: Process ID, Process state
- **PCB:** Stores entire execution context: Memory, Hardware and OS

## Process Scheduling

- Criteria: (i) Fairness, (ii) Utilisation (minimise idling)
- **Turnaround time:** Waiting time + CPU Time + I/O time; arrival to finish
- **Throughput:** # tasks finished per unit time
- **Makespan:** Total time to process all tasks
- **CPU utilisation:** Percentage of time when CPU is working on a task

## Batch Processing

### First Come First Serve (FCFS)

Tasks acquire CPU in FIFO manner. Non-preemptive; runs until done or blocked

#### Advantages

- No starvation

#### Disadvantages

- Long waiting time

Convoy effect: CPU-heavy task  $T$  arrives first, blocking I/O-heavy tasks from CPU → I/O idles. When  $T$  enters I/O phase, I/O-heavy tasks quickly finished their CPU-phase and joins I/O queue, waiting for  $T$  → CPU idles

### Shortest Job First (SJF)

Choose process with smallest total CPU time. Non-preemptive.

CPU time prediction:  $\text{Predicted}_{n+1} = \alpha \text{Actual}_n + (1 - \alpha) \text{Predicted}_n$

#### Advantages

- Guarantees shortest waiting time

#### Disadvantages

- Can cause starvation – new tasks with shorter CPU time can keep arriving
- Need to predict CPU time for each process, which may not be accurate

### Shortest Remaining Time (SRT)

Choose process with shortest remaining CPU time. Preemptive variant of SJF i.e. arrival of tasks with shorter remaining CPU time will preempt current

task immediately. No time quantum

#### Advantages

- Little overhead; decision making only when a task arrives/completes
- Guarantees shortest possible average waiting time

#### Disadvantages

- Can cause starvation – new tasks with shorter CPU time can keep arriving
- Need to predict CPU time for each process, which may not be accurate

## Interactive Systems

- Short TQ: Shorter waiting time, larger context switch overhead, lower CPU utilisation
- Long TQ: Longer waiting time, smaller context switch overhead, higher CPU utilisation
- Short ITI: More timer interrupts, less time spent on actual user processes, longer turnaround time

### Round Robin

Preemptive variant of FCFS. Task runs until TQ expires, is blocked, or is completed

#### Advantages

- No starvation
- Response time guaranteed to be at most  $(n-1)TQ$

#### Disadvantages

- If CPU time for tasks are the same and much greater than TQ → longer average turnaround time than FCFS
- If many tasks and their CPU time is much greater than TQ → context switch overhead → reduced throughput

### Priority Scheduling

Choose process with highest priority. Non-preemptive and preemptive.

In non-preemptive variant, newly arrived tasks with higher priority need to wait until current task gives up CPU. In preemptive variant, arrival of higher priority task preempts task in CPU even before its TQ expires

#### Advantages

- High priority tasks get CPU quicker
- Responsive

#### Disadvantages

- Low priority tasks can be starved by the continuous arrival of higher priority tasks

### Multi-level Feedback Queue (MLFQ)

Newly arrived tasks have high priority. Higher priority tasks can preempt lower priority task even before its TQ expires. If TQ is fully utilised, a task is demoted to the next lower priority queue, otherwise, its priority is retained.

For tasks with same priority, they are scheduled in RR fashion.

#### Advantages

- Responsive

#### Disadvantages

- Can be abused by blocking a task just before its TQ expires to retain its high priority → starve low priority processes
- Task with length CPU phase followed by I/O phase → degrades responsiveness since it sinks to the lowest priority during CPU phase

### Lottery Scheduling

Allocate each task a number of “tickets”. Pick a ticket randomly and the corresponding task holding that ticket gets scheduled

#### Advantages

- Responsive, since newly arrived tasks can participate in next lottery
- Good level of control; number of tickets per task can be controlled

## Producer-Consumer Semaphores

- Producers: Wait when buffer is full
- Consumers: Wait when buffer is empty

### Busy waiting

```
while (TRUE) {
    Produce Item;
    while(!canProduce){}
    wait( mutex );
    if (count < K) {
        buffer[in] = item;
        in = (in+1) % K;
        count++;
        canConsume = TRUE;
    } else
        canProduce = FALSE;
    signal( mutex );
}
```

Producer Process

```
while (TRUE) {
    while (!canConsume){}
    wait( mutex );
    if (count > 0) {
        item = buffer[out];
        out = (out+1) % K;
        count--;
        canProduce = TRUE;
    } else
        canConsume = FALSE;
    signal( mutex );
    Consume Item;
}
```

Consumer Process

- canProduce / canConsume: Ensures no starvation (at most n producers before consumers can consume, vice versa)
- Busy waiting → CPU wastage

```
while (TRUE) {
    Produce Item;

    wait( notFull );
    wait( mutex1 );
    buffer[in] = item;
    in = (in+1) % K;
    signal( mutex1 );
    signal( notEmpty );
}
```

Producer Process

```
while (TRUE) {

    wait( notEmpty );
    wait( mutex2 );
    item = buffer[out];
    out = (out+1) % K;
    signal( mutex2 );
    signal( notFull );

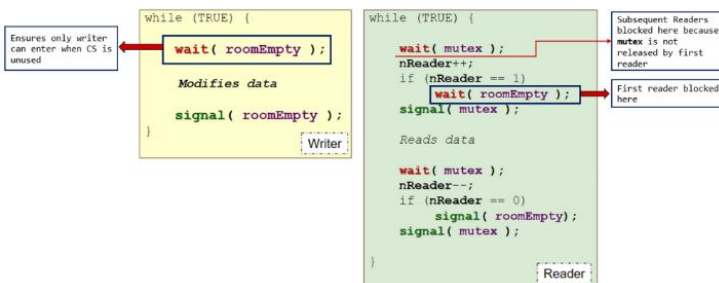
    Consume Item;
}
```

Consumer Process

- Initialise notFull = buffer\_size, notEmpty = 0
- Used in asynchronous message passing

## Readers-Writers Semaphores

- Exclusive rights for writers. Simultaneous readers



- Problem: Constant stream of readers → nReader will never be 0 → roomEmpty never released → starve writers (cut writers' queue)
- Add revDoor semaphore

```
while (TRUE) {
    wait(revDoor);
    wait( roomEmpty );

    Modifies data

    signal( roomEmpty );
    signal(revDoor);
}
```

Writer

```
while (TRUE) {
    wait(revDoor);
    signal(revDoor);
    wait( mutex );
    nReader++;
    if (nReader == 1)
        wait( roomEmpty );
    signal( mutex );

    Reads data

    wait( mutex );
    nReader--;
    if (nReader == 0)
        signal( roomEmpty );
    signal( mutex );
}
```

Reader

- As soon as there is a writer waiting on roomEmpty, no readers can wait on roomEmpty

## Dining Philosophers

### Limited Eater

- Initialise seats = N - 1
- Guarantees at least 1 philosopher can have both chopsticks

```
Semaphore seats = S(4);
//initialization

void philosopher( int i ){
    while (TRUE){
        Think();
        wait( seats );
        takeChpStck( LEFT );
        takeChpStck( RIGHT );
        Eat();
        putChpStck( LEFT );
        putChpStck( RIGHT );
        signal( seats );
    }
}
```

### Tanenbaum Solution

```
#define N 5
#define LEFT ((i+N-1) % N)
#define RIGHT ((i+1) % N)

#define THINKING 0
#define HUNGRY 1
#define EATING 2

int state[N];
Semaphore mutex = 1;
Semaphore s[N];

void philosopher( int i ){
    while (TRUE){
        Think();
        takeChpStcks( i );
        Eat();
        putChpStcks( i );
    }
}
```

```
void takeChpStcks( i )
{
    wait( mutex );
    state[i] = HUNGRY;
    safeToEat( i );
    signal( mutex );
    wait( s[i] );
}
```

Prevents forcing neighbours to eat when they are not hungry, which can happen at.....

```
void safeToEat( i )
{
    if( (state[i] == HUNGRY) &&
        (state[LEFT] != EATING) &&
        (state[RIGHT] != EATING) ) {
        state[i] = EATING;
        signal( s[i] );
    }
}
```

```
void putChpStcks( i )
{
    wait( mutex );
    state[i] = THINKING;
    safeToEat( LEFT );
    safeToEat( RIGHT );
    signal( mutex );
}
```

- Each philosopher waits for its own private semaphore

## Threads

- Shared among threads: Text, Data, Heap, Resources (I/O, Files)
- Not shared among threads: Stack, CPU Registers
- Threads from the same process can run on different CPUs

Threads	Processes
Switch: involves hardware context (changing FP, SP, register values)	Context switch: involves hardware, OS, memory context
Multiple threads in same process equires less resources to manage vs multiple processes	Requires more resources to manage

- Syscall concurrency: parallel syscalls possible (e.g. multiple threads can open the same file. No thread will lock a file)

Operation	Behaviour (Linux)
fork()	Child process will only inherit the thread on which fork() was called
exit()	Entire process is killed (all threads killed)
exec()	Image of entire process is replaced (all other threads killed)

- Single-core CPU with simultaneous multithreading (SMT), or hyperthreading, can execute multiple threads at the same time

## Thread Models

User Threads	Kernel Threads
Implemented as user library	Implemented and handled by the OS
Kernel unaware of these threads	Kernel aware of these threads
Thread level scheduling by OS not possible	Thread level scheduling by OS possible
Less syscalls, customisable	More syscalls, not customisable
Cannot exploit multi-core CPUs. When one thread is blocked, the OS sees the entire process as blocked → remaining threads cannot run on the other CPUs	Can exploit multiple CPUs. When one thread is blocked, the remaining threads can run on other CPUs

### Hybrid Thread Model

- User thread binds to a kernel thread

## Memory Management

### Memory Abstraction

- Logical address: Program's view of the memory
- Physical address: Actual address in hardware

### Address Relocation

- Recalculate all memory references (increment by b) when program loads
- Base register: Stores base b, physical memory address at start of program
- Limit register: Stores the range of memory address for the process
- Addressing:
  1. Physical addr = b + MemRef
  2. Check Physical addr < Limit
- Operations done by hardware → Fast

## Contiguous Memory Management

- Memory is divided into partitions
- Entire process fits into one partition
- **Internal Fragmentation:** Memory wastage inside partition
- **External Fragmentation:** Memory wastage outside of partitions

### Fixed size partitioning

- All partitions are of fixed and equal size
- One partition, one process. Processes cannot use space caused by internal fragmentation

Pros	Cons
Easier to manage	<b>Internal Fragmentation</b>
Fast; all free partitions are equal. Just pick any free partition	Need to ensure that partition size needs to be large enough to fit largest process

### Dynamic partitioning

- Partition size = process size

Pros	Cons
<b>No internal fragmentation</b>	<b>External Fragmentation</b> – when one process is freed, it leaves a whole between two processes
	Due to EF, reasonably sized processes may not fit into any of the holes
	Takes time to find a partition
	More OS memory overhead

- **First-Fit:** Take first partition that is large enough. Fast
- **Best-Fit:** Take smallest partition that is large enough
- **Worst-Fit:** Take largest partition
- **Optimisations:**
  - Merge adjacent holes
  - Compaction: occasionally move occupied partitions around to create larger partitions

### Multiple Free Lists (Buddy System)

- Array of linked-lists. Each linked-list is a list of free fixed size partitions, each typically of size  $2^n$
- Allocation Steps (Process size =  $N$ ):
  1. Find largest  $m$  s.t.  $2^m \geq N$
  2. Access  $A[m]$ : Partition in  $A[m]$ 
    - Available: Allocate first partition in list
    - Unavailable: Find smallest partition in  $A[m+1:]$  and recursively split until you get a partition of size  $A[m]$
- Deallocation Steps (Process size =  $N$ ):
  1. Add partition back to  $A[m]$
  2. If buddy is free: Recursively,
    - merge and place merged partition in  $A[m+1]$
    - $m \leftarrow m+1$
- Two partitions are buddies if they are equal in size ( $2^m$ ) and only the  $m$ -th bit (zero-indexed) differs

## Disjoint Memory Management

- Logical memory is broken down into pages/segments which can occupy various parts of physical memory

### Paging

- Logical memory broken down into fixed size pages
- Physical memory broken down into same sized frames. Size is determined by **hardware**
- Logical pages fit into any physical frames
- **Page Table (PT)** maps logical pages to physical frames. PT is stored in kernel space in contiguous memory locations
- Pointer to PT is part of memory context of a process
- **Page Table Entry (PTE):** Frame Number, Access-rights bit, Valid-Invalid bit
- Mapping (Page size =  $2^n$ ,  $m$ -bit address):
  1. Page  $P$  = most significant  $m - n$  bits of logical address
  2. Frame  $F = P^{\text{th}}$  PTE
  3. Physical Address =  $F(2^n) + \text{Offset}$  (least significant  $n$  bits)

Pros	Cons
<b>No external fragmentation</b>	Every memory reference (including instructions) requires 2 memory accesses $\rightarrow$ slow
<b>Negligible internal fragmentation</b> – at most 1 internally fragmented frame per process	

### Paging with Caching

- **Translation Look-Aside Buffer (TLB)** caches the most recently accessed PTEs
- TLB is a hardware component which the OS is aware of
- Reduces number of PT access  $\rightarrow$  reduces performance overhead of logical-to-physical addr translation and reduces
- TLB not part of process' context
- **TLB Entry:** Page Number, Frame Number, Access-rights bit, Valid-Invalid bit
- Mapping (Page size =  $2^n$ ,  $m$ -bit address):
  1. Page  $P$  = most significant  $m - n$  bits of logical address
  2. TLB[ $p$ ]?
    - TLB Hit: Frame  $F = \text{TLB}[p].\text{Frame}$ . Done
    - TLB Miss: Frame  $F = P^{\text{th}}$  PTE
  3. Physical Address =  $F(2^n) + \text{Offset}$
- Average Memory Access Time (AMAT):
$$P(\text{TLB Hit}) \times (T_{\text{TLB}} + T_{\text{Mem}}) + P(\text{TLB Miss}) \times (T_{\text{TLB}} + 2 \times T_{\text{Mem}})$$
- TLB contents **flushed** during context switch  $\rightarrow$  prevent incorrect addr translation for next process

### Page Sharing

- Use same physical frame number in PTEs
- **Copy-on-Write:** Child process has same PT contents as parent until it writes into a page  $\rightarrow$  change Frame Number for written page
- Shared code page (e.g. C stdlib, syscalls)

### Segmentation

- Logical memory broken down into variable sized segments
- Each segment occupies a contiguous physical memory space
- **Segment Table (ST)** maps segments to physical memory. Stored in kernel space
- **Segment Table Entry:** Limit, Base
  - Limit: Max size of segment
  - Base: Starting physical addr of segment
- Permissions are scoped per segments (vs page in paging)
- Gaps between segments are provisioned as segments can grow/shrink
- Mapping:
  1. Segment Entry  $S = S^{\text{th}}$  entry on ST
  2. Check Offset  $< S.\text{Limit}$
  3. Physical Address =  $S.\text{Base} + \text{Offset}$
- Segment Table can be cached in CPU registers  $\rightarrow$  no need lookup ST

Pros	Cons
<b>No internal fragmentation</b>	<b>External Fragmentation</b>
Efficient book keeping (information is scoped per segment) $\rightarrow$ less duplicated information	
Matches programmer's view of memory	

### Segmentation with Paging

- Logical memory broken down into variable sized segments
- Each segment is further broken down into fixed size pages
- Each segment has its own page table
- Segment occupies non-contiguous physical memory space
- **Segment Table Entry:** Page Limit, Page Table Base, Permission bits
  - Page Limit: Max number of pages
  - Page Table Base: Physical address of segment's page table
- Mapping (Page size =  $2^n$ ):
  1. Segment Entry  $S = S^{\text{th}}$  entry on ST
  2. Check Page  $P < S.\text{PageLimit}$
  3. Frame  $F = P^{\text{th}}$  PTE in PT pointed by  $S.\text{PageTableBase}$
  4. Physical Address =  $F(2^n) + \text{Offset}$

Pros	Cons
<b>No external fragmentation</b>	
<b>Negligible internal fragmentation</b> – at most 1 internally fragmented frame per <u>segment</u>	

## Virtual Memory Management

- Logical memory > physical memory
- Store some pages in secondary storage
- **Resident Bit** in PT: Indicates if a page is in physical memory

### Demand Paging

- All pages initially on disk. Bring pages into memory only when needed
- **Inclusive systems:** Page remains on disk when loaded into memory
- **Exclusive systems:** Page removed from disk when loaded into memory
- Memory Access (Page size =  $2^n$ ,  $m$ -bit address):
  1. Page  $P$  = most significant  $m - n$  bits of logical address

2. TLB[p]?
  - TLB Hit: Frame  $F = \text{TLB}[p].\text{Frame}$ . Go to step 4
  - TLB Miss: Step 3
3. PT[p].Resident?
  - Resident: Frame  $F = \text{PT}[p].\text{Frame}$
  - Non-Resident: **Page Fault**. Repeat from step 1
4. Physical Address =  $F(2^n) + \text{Offset}$

- **Page Fault**: Accessing a page that is not in physical memory

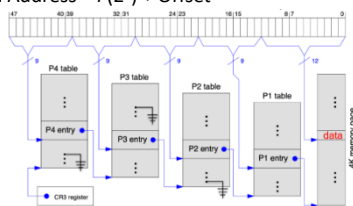
- Page Fault procedure:

1. Process is BLOCKED. TRAP and enters kernel mode
2. DMA Controller locates page in disk
3. DMA Controller loads page into free frame
  - RAM is full: Page Replacement Algorithm to evict a page
  - Write out evicted page if page is dirty
4. Update PT for loaded page (Resident bit, frame #)
  - RAM was full: Update PT for evicted page
5. Update TLB
6. Return from TRAP

Pros	Cons
Fast startup time	Process appears sluggish at the start
Small memory footprint	Page faults have cascading effect on other processes ( <b>thrashing</b> )

### Hierarchical Paging

- Reason: Large VM space  $\rightarrow$  large number of pages  $\rightarrow$  large PT  $\rightarrow$  entire PT resides in memory, even if a small subset of pages is used by a process  $\rightarrow$  huge wastage
- Calculate  $n$  # levels (Page size =  $2^m$ ,  $v$ -bit virtual address space)
  1. Calculate # PTEs in a single PT / Frame =  $2^n$
  2. # bits used for offset in each level (except last) =  $n$
  3. # bits used for last level =  $m$
  4. # levels =  $\lceil (v - m) / n \rceil$
- Memory Access ( $m$ -bit logical add,  $n$ -level paging,  $2^n$  entries per PT/PD, Page size =  $2^m$ ):
  1. Page  $P$  = most significant  $m - en$  bits of logical address
  2. TLB[p]?
    - TLB Hit: Frame  $F = \text{TLB}[p].\text{Frame}$ . Go to step 4
    - TLB Miss: Step 3
  3. Repeat  $n$  times (**Page-table Walk**):
    - 3.1.  $i = [m : m - e]$  (i.e. most significant  $e$  bits starting from  $m$ )
    - 3.2.  $\text{PT}_{n-1}/\text{Frame } F = \text{PT}_n[i].\text{Frame}$
    - 3.3.  $n \leftarrow n - 1, m \leftarrow m - e$
  4. Physical Address =  $F(2^m) + \text{Offset}$



- Page walks are done in hardware
- **MMU cache** for each level caches the most recently access PD entry
- Root PT is part of memory context of a process

Pros	Cons
Saves memory space – PDs can have empty entries; only store in memory the PTs containing pages used by the process	Requires multiple serialised memory accesses (page-table walk)

### Inverted Page Table

- $\langle \text{Frame \#} \rangle \rightarrow \langle \text{PID, page \#} \rangle$
- # entries = # frames in physical memory  $\rightarrow$  hence smaller than PTs
- Mapping (Frame size =  $2^n$ ):
  1. Lookup hashtable[PID, page #] to get  $F(\langle \text{PID, page \#} \rangle)$  is  $F^{\text{th}}$  entry in inverted PT
  2. Physical Address =  $F(2^n) + \text{Offset}$

### Page Replacement Algorithms

- **Dirty Bit** in PT: indicates whether a page has been written into since the last time it was loaded into memory
  - Clean page: just evict, no need to write back to disk (for inclusive systems)
  - Dirty page: need to write back to disk

### Optimal Page Replacement (OPT)

- Not realisable, but good as base of comparison
- Replace the page that will not be needed again for the longest period of time
- Guarantees minimum # of page faults

### FIFO

- Replace the oldest page
- OS needs to maintain a queue of resident page numbers
- Simple to implement (no hardware support needed)
- Suffers from **Belady's Anomaly** (does not exploit temporal locality)
  - # page faults can increase with more physical frames

### Least Recently Used (LRU)

- Replace the page that has not been used in the longest time
- Exploits temporal locality
- **Counter** Approach:
  - PT stores "time-of-use" field which increments when a page is referenced
  - Replace the page with smallest "time-of-use"
  - Linear search on PT + "time-of-use" is monotonically increasing  $\rightarrow$  can cause overflow
- **Stack** Approach:
  - Maintain a stack of page numbers
  - Page X is referenced  $\rightarrow$  remove X from stack and push it on top of stack
  - Replace the page at the bottom of the stack
  - Not a pure stack + Hard to implement in hardware

### Second-chance Page Replacement (CLOCK)

- FIFO, but skip the page if it has been accessed since the previous round of replacement
- **Reference Bit** in PT: indicates whether a page has been accessed since the previous round of page replacement
- Reference Bit is toggled by hardware
- Algorithm:
  1. Inspect oldest page X (FIFO)
  2. If reference bit == 0: Replace page X. Done
  3. If reference bit == 1: Skip page X and set PT[X].Reference = 0
  4. Repeat from step 1 until a page is chosen
- If all pages have reference bit == 1, it degenerates into FIFO

### File Systems

- Storage: disk; Access speed: variable; Usage: explicit access

### File Metadata

- File **name**: file extension may be used to indicate file type
- File **type**:
  - ASCII files: can be displayed or printed as is
  - Binary files: executable files; prints gibberish
  - Distinguished using file extension and/or magic numbers embedded at the beginning of the file
- **Protection** bits: 9 bits
 

- rwxr--w--

|   |   |   |   |   |   |   |   |

type owner group everyone else

  - `chmod _ _ _ <filename>`, each `_` is 3 bit number in the form `rwX` where the bit = 1 if that permission is granted
  - Access Control List: set permissions for specific user(s)

### File Structure & Access Methods

#### File Structure

- **Array of bytes**: Each byte has unique offset from start of file
- **Fixed length records**: Array of records, each record is a group of fixed number of bytes.  $N^{\text{th}}$  record =  $\text{sizeof}(\text{Record}) * (N - 1)$
- **Variable length records**: Array of records, each record is a group of variable number of bytes

#### Access Methods

- **Sequential**: Data must be read in order and from the start of file; no skipping but can rewind
- **Random Access** (Array of bytes): Data can be read in any order
  - Read( $i$ ): read the  $i^{\text{th}}$  byte
  - Seek( $i$ ): move to  $i^{\text{th}}$  byte
- **Direct Access** (Fixed-length records): Random access except  $i^{\text{th}}$  record instead of  $i^{\text{th}}$  byte

### File Information in OS

- **File pointer**: current position within a file

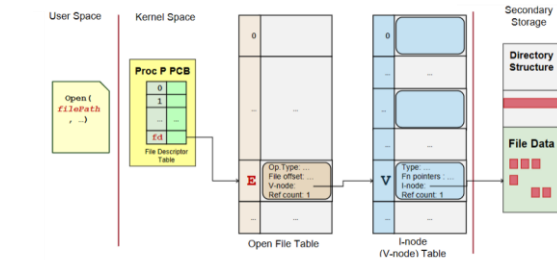


- **File descriptor:** unique ID of file
- **Disk location:** Actual physical location on disk
- **Open Count/Reference Count:** how many processes has this file opened

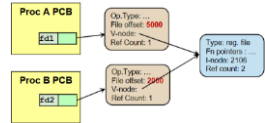
#### Tables

- **Per-process** file descriptor table: keeps track of opened files for a process
- **System-wide** open-file table: keeps track of all opened files in system
  - <Op. Type, File Offset, V-node, Ref count>
- **System-wide V-node** table: contains info about location of file in disk
  - <Type, Fn Pointers, I-node, Ref count>
- When open()/close() is called, a entry is created in/removed from the per-process and system-wide open-file tables

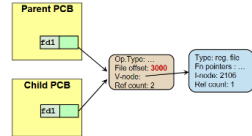
Pointer to this entry is stored in the process's File Descriptor Table (which is saved in PCB on context switch)



Opened twice from 2 processes



2 processes share same OF entry



## Directories

- DAG: A file can appear in multiple directories with different path names, each referring to the same file content.

#### Hard Link

- Directories have separate pointers to the same file in disk
- Cannot share same directories
- Modifying permissions of the hardlink will also modify permissions of the linked file, vice versa
- Low overhead, only pointers are added in directory
- Deletion problems

#### Symbolic Link

- Directories share the same file/directory
- A special link file (G) that contains path name to the file/directory
- Modifying permissions of the softlink will not affect permissions of the linked file, vice versa
- Symbolic link deleted → only G is deleted
- File F is deleted → G becomes dangling link

## File System Implementation

- File is divided into smaller units each fitting a logical block maintained by the OS. OS maps each logical block to a physical location on disk
- Contiguous logical blocks are grouped together to form a partition

#### File Data allocation

- File data is split into logical blocks that reside in **File Data** section
- Last logical block of a file may be internally fragmented
- **File Info** stores file metadata
- **Contiguous** allocation:
  - Logical blocks for file data are allocated contiguously
  - File Info entry: <file, start, length>, length = # of blocks allocated for file
  - Fast and simple random and sequential access
  - External fragmentation possible
- **Linked List** allocation:
  - Logical blocks for file data may not be allocated contiguously
  - Doubly/singly linked list of disk blocks. Each block stores pointer to next (and previous, for doubly-linked lists) block and actual file data
  - File Info entry: <file, start, end>
  - No fragmentation
  - Sequential and random access: many disk traversals → Slow
  - Each block also stores extra info other than data → overhead
  - Less reliable; pointers can be incorrect

#### FAT allocation:

- Pointers to next block are stored in File Allocation Table (FAT)
- FAT is cached into memory (RAM)
- $n^{\text{th}}$  FAT entry: block # of next block to go to after block  $n$  or EOF
- FAT: One entry for each logical block in partition → FAT can be huge
- Sequential and random access: many RAM traversals → Faster than disk traversal

#### Index allocation:

- **Index block** is like FAT that is per (opened) file vs per system
- $n^{\text{th}}$  index block entry: block # of the  $n^{\text{th}}$  block or empty
- Index block is stored in a logical block
- File Info entry: <file, index block #>
- Sequential and random access: =  $n^{\text{th}}$  block # =  $n^{\text{th}}$  entry in index block → Fast
- Max # of blocks for file = # number of index block entries
- Linked List of Index blocks/Multi-level Index Blocks → allows larger file sizes

#### Free Space Management

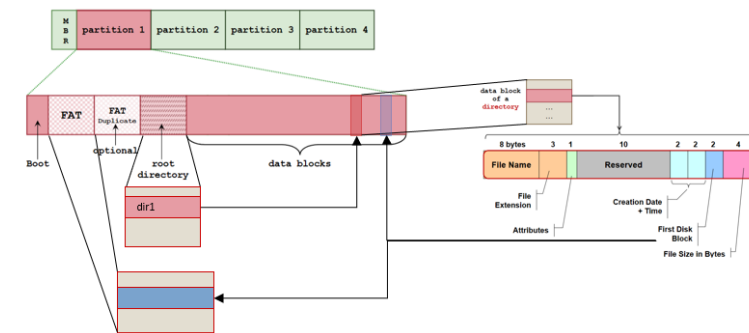
- Information maintained in **Partition Details** section
- **Bitmap:**
  - Each block represented by 1 bit. Occupied = 0, Free = 1
  - Bitmap needs to be loaded into RAM for efficient access → many blocks → large bitmap → memory overhead
- **Linked List** of blocks
  - Each block contains a number of block # that are free and a pointer to the next block
  - Partition Details stores only the pointer to first block
  - High disk memory overhead → can be mitigated by storing some of the information in the free blocks in File Data section

#### Directories

- Information in **Directory Structure** section
- **Linear List:**
  - Entry: <file name, file info...> or <file, name, pointer to file info>
  - Linear search to locate a file → Inefficient
- **Hash Table:**
  - Entry: <file info...> or <pointer to file info>
  - Hash(file name) → index  $k = k^{\text{th}}$  entry in hash table
  - Chained collision resolution

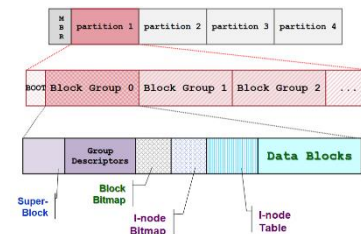
## Microsoft FAT

- FAT is cached into RAM to facilitate linked-list traversal
- FAT entry: <FREE/Block #/EOF/BAD>
- $n^{\text{th}}$  FAT entry: block # of next block to go to after block  $n$  or EOF
- FAT- $n \rightarrow 2^n$  data blocks,  $n$ -bit block # →  $2^n$  FAT entries →  $2^n (n/8)$  bytes
- Root directory entries stored in special section **Root Directory**. Remaining directories are stored in data blocks



## Linux Ext2

- **Block Bitmap:** Shows which data blocks are free/occupied
- **I-node Bitmap:** Shows which I-nodes are free/occupied
- **Group descriptors:** Metadata about all block groups in the partition (e.g. # free blocks, # free I-nodes, location of bitmaps). Is duplicate in each block group
- **Superblock:** Information about whole file system e.g. total I-nodes number, # I-nodes per group, total # disk blocks, # blocks per group. Duplicated in each block group
- **I-node Table**
  - Contains only I-nodes of that block group
  - I-node entry: <Metadata (file info), 15 Datablock pointers, Ref. count>. 128 bytes
  - **15 Datablock pointers**
    - First 12 points to direct blocks (blocks storing actual file data)
    - 13<sup>th</sup>: points to single indirect block (storing direct pointers)
    - 14<sup>th</sup>: points to double indirect block (storing single indirect ptrs)
    - 15<sup>th</sup>: points to triple indirect block (storing double indirect ptrs)



- Contains one entry for one file or directory

#### Directory structure

- Directory stored in datablocks; datablocks store a linked list of directory entries, each for each file/subdirectories within this directory
- Directory entry: <I-node #, size of entry, length of name, type (F or D), name>

