

4.1 Three Examples where Multithreading provides better performance than a single-threaded solution.

- ① Image Processing :- When applying filters, edge detection, or transformations to an image, different sections of the image can be processed concurrently using multiple threads, significantly reducing execution time compared to single-threaded approach.
- ② Game Physics and AI Processing :- Modern games use multithreading to handle physics calculations, AI decision making, & rendering separately. This prevents frame rate drops and ensures smoother game play as compared to single-threaded system where these tasks would cause delays.
- ③ Sorting Large Datasets :- Algorithms like Merge sort and Quick sort can be parallelized using multithreading by sorting subarray independently. This speeds up sorting operations significantly.

4.02 Speedup Calculation using Amdahl's law. 60% parallelizable part

a) 2 processing cores

by 4 processing cores.

b) speedup = 1

by speedup = 1

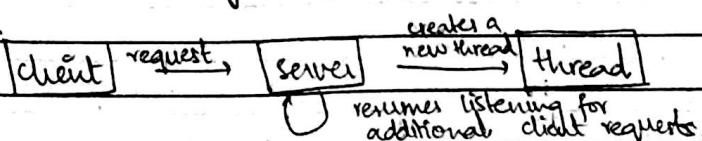
$$(1 - 0.6) + \frac{0.6}{2}$$

$$(1 - 0.6) + \frac{0.6}{4}$$

$$= 1.4285 \text{ times}$$

$$= 1.8181 \text{ times}$$

4.03 Does multithreaded web server exhibit task parallelism or data parallelism.  
 A multithreaded web server exhibits task parallelism. Each thread handles a different client request independently, meaning multiple tasks (serving different web requests) are executed in parallel.



4.4 Differences between User level thread & Kernel level thread?

Feature	User Level	Kernel Level
① Managed by	User space libraries	Operating system (OS) Kernel
② Context switching	Fast (no OS intervention)	Slower (requires system calls)
③ Scheduling	Handled by user-space scheduler	Scheduled by OS.

<u>Blocking</u>	If one thread blocks, all threads in the process may block.	Other threads can continue running if one blocks.
<u>Portability</u>	More portable across OS.	OS-dependent

- Under what circumstances is one better than the other?
- User-level threads are better when lightweight & fast context switching is required.
- Kernel-level threads are better when the application involves I/O operations or requires OS-level scheduling for better resource management.

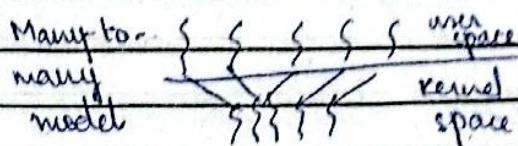
- 4.5 Actions taken by Kernel to context-switch between kernel level threads.
- ① Save the state of currently running thread (registers, program counter, stack pointer).
  - ② Update the thread control block (TCB) of the current thread.
  - ③ Select a new thread to run using the OS scheduler.
  - ④ Load the state of the new thread from its TCB (restoring registers, program counter, stack pointer).
  - ⑤ Resume execution of the new thread.

- Q.6 What resources are used when a thread is created? How do they differ from those used when a process is created?

Resource	Thread Creation	Process Creation
① Memory Allocation	Shares address space / memory with parent process.	Allocates separate memory space.
② Scheduling	Requires minimal scheduling information	Requires full process scheduling information
③ Data	No new PCB is created, just a TCB (thread control block)	New PCB is created.
④ Execution Context	Shares execution context with other threads in the process.	Has its own execution context.

Creating a thread requires TCB i.e. allocating a small data structure to hold a register set, stack, and priority.

4.7 Many-to-many model  $\rightarrow$  mapping done through WPs.  
 Is it necessary to bind a real-time thread to an WP? Explain  
 In many-to-many model, multiple user-level threads are mapped to a smaller or equal number of kernel-level threads (WPs - light weight processes).



For real-time systems, binding a real-time thread to an WP can be necessary because:-

Direct CPU scheduling:- Real-time threads often require predictable scheduling, which Kernel-level threads (WPs) can provide.

Avoiding user-space scheduling delays:- If a real-time thread relies on user-level scheduling, it may experience delays.

Ensuring priority handling:- The OS can directly prioritize real-time threads when they are mapped to WPs.

4.8 Two programming examples in which multithreading does not provide better performance than a single-threaded solution.

#### (i) Simple Sequential Tasks (CPU-bound with dependencies)

for example:- Calculating fibonacci sequence recursively.

As each step depends on the previous ones, making parallelization difficult. The overhead of thread creation and context switching outweighs any potential speedup.

A single threaded solution avoids synchronization & thread management costs.

single threaded version:-	Multithreaded version:-
{ if ( $n <= 1$ )	{ if ( $f <= 1$ )
return n;	return n;
return fib(n-1) + fib(n-2);	else {
}	pthread_t t1, t2; pthread_create(&t1, NULL, fib, &a); pthread_create(&t2, NULL, fib, &b); pthread_join(t1, NULL); pthread_join(t2, NULL);
	}
Pg No.	4   1

Date:

## ② Single Shared I/O Resources

Multiple threads writing to the same file contend for locks  
making it slower than a single thread.

single threaded is faster as no lock contention or thread synchronization overhead.

4.9 Under what circumstances does a multithreaded solution using multiple kernel threads perform better than a single threaded solution on a single-processor system?

Multithreaded solution using multiple kernel threads outperforms a single-threaded solution on a single-processor system under following circumstances :-

- ① Yo Bound operations or Blocking operations:- When threads frequently block for Yo (i.e file, database queries), or synchronization (ie waiting for locks), multithreading allows CPU to switch to other ready threads.

Pg No.

Date:

during blocking periods. This improves CPU utilization compared to single thread that would be idle during waits.

- (2) Improved Responsiveness:- In interactive applications, multithreading enhances perceived performance. For e.g:- a background thread can handle computations while a foreground thread manages user input, preventing the application from freezing during long tasks.
- (3) Event-Driven Workflows:- when threads wait for independent external events (e.g, timers, signals), multithreading allows concurrent event handling without serial bottlenecks.

4.10 Which components are shared across threads in a multithreaded process?

- a. Register values
  - b. Heap memory
  - c. Global variables
  - d. Stack memory
- Dynamically allocated memory (e.g., malloc, new) is shared by all threads.

Each thread has its own set of register values & own call stack storing local variables in function call frames.

4.11 Can a multithreaded solution using multiple user-level threads achieve better performance on a multiprocessor system than on a single-processor system? Explain. → can be Yes or No, depending on situation.

Yes, but only if threading model allows the OS to distribute user-level threads across multiple CPUs (e.g., in many-to-many or one-to-one models).

No, if the threading model restricts all user-level threads to a single CPU. (e.g., many-to-one).

4.12 Google chrome → new tab in separate process vs new tab in separate thread.

Designing chrome to open each new tab in a separate thread instead of separate process would not achieve the same benefits.

Separate process ensures isolation & stability as each process has its own memory space. If a tab (process) crashes, it does not affect other tabs or the main browser process. Whereas, threads share the same memory space. A bug or crash in one thread could crash the entire browser, including all other tabs. True parallelism can be achieved by distributing processes across multiple CPU cores. While threads require synchronization (e.g., locks) which can lead to contention.

Pg No. and overhead. Chrome uses sandboxing to restrict each tab's process from accessing sensitive OS resources. Hence chrome's use of process ensures security.

4.13 Is it possible to have concurrency, but not parallelism?

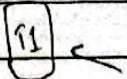
Yes. Concurrency refers to managing multiple tasks simultaneously (e.g., interleaving execution on a single core), while parallelism requires tasks to execute at the same time (e.g. on multiple cores).

Example:-

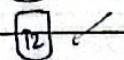
CPU, Core 1

Concurrent,

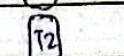
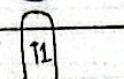
Not Parallel



Both task running on Core 1,



with making progress one by one.



4.14 Amdahl's Law :-  $\frac{1}{(1-p)} + \frac{p}{N}$

40% parallel with a) 8 cores , b) 16 cores

$$\text{a) Speedup} = \frac{1}{\frac{(1-0.4)}{8} + \frac{0.4}{8}} = 1.538 \quad \text{b) Speedup} = \frac{1}{\frac{(1-0.4)}{16} + \frac{0.4}{16}} = 1.6$$

67% parallel with a) 2 cores , b) 4 cores

$$\text{a) Speedup} = \frac{1}{\frac{(1-0.67)}{2} + \frac{0.67}{2}} = 1.5037 \quad \text{b) Speedup} = \frac{1}{\frac{(1-0.67)}{4} + \frac{0.67}{4}} = 2.01$$

90% parallel with a) 4 cores , b) 8 cores

$$\text{a) Speedup} = \frac{1}{\frac{(1-0.9)}{4} + \frac{0.9}{4}} = 3.0769 \quad \text{b) Speedup} = \frac{1}{\frac{(1-0.9)}{8} + \frac{0.9}{8}} = 4.1058$$

4.15 Task parallelism or Data parallelism.

- Using separate thread to generate a thumbnail for each photo in a collection.

DATA PARALLELISM. Same task (<sup>Thumbnail</sup> generation) applied to different data (photos).

- Transposing a matrix . DATA PARALLELISM

Same operation (swapping) applied to different parts of matrix.

- Networked app (read/write threads). TASK PARALLELISM

Different tasks (reading, writing)

- Fork-join array summation. DATA PARALLELISM

Same summation operation applied to subarrays.

- (GCD) Grand Central Dispatch. TASK PARALLELISM

GCD schedules independent tasks / blocks of work.

1.13 Is it possible to have concurrency, but not parallelism?

Ans. Concurrency refers to managing multiple tasks simultaneously (e.g., interleaving execution on a single core), while parallelism requires tasks to execute at the same time (e.g. on multiple cores).

Example:-

CPU, Core 1

Concurrent,



Both task running on Core 1.

Not Parallel



both making progress one by one.



4.14 Amdahl's law :-  $\frac{1}{(1-p) + \frac{p}{N}}$

40% parallel with a) 8 cores , b) 16 cores

$$\text{a) Speedup} = \frac{1}{(1-0.4) + \frac{0.4}{8}} = 1.538x \quad \text{b) Speedup} = \frac{1}{(1-0.4) + \frac{0.4}{16}} = 1.6x$$

67% parallel with a) 2 cores , b) 4 cores

$$\text{a) Speedup} = \frac{1}{(1-0.67) + \frac{0.67}{2}} = 1.5037x \quad \text{b) Speedup} = \frac{1}{(1-0.67) + \frac{0.67}{4}} = 2.01x$$

90% parallel with a) 4 cores , b) 8 cores

$$\text{a) Speedup} = \frac{1}{(1-0.9) + \frac{0.9}{4}} = 3.0769x \quad \text{b) Speedup} = \frac{1}{(1-0.9) + \frac{0.9}{8}} = 4.1058x$$

4.15 Task parallelism or Data parallelism.

• Using separate thread to generate a thumbnail for each photo in a collection.

DATA PARALLELISM. Same task (<sup>sequential</sup> generation) applied to different data (photos).

• Transposing a matrix. DATA PARALLELISM

Same operation (swapping) applied to different parts of matrix.

• Networked app (Read/ write threads). TASK PARALLELISM

Different tasks (reading, writing).

• Fork-join array summation. DATA PARALLELISM

Same summation operation applied to subarrays.

• (GCD) Grand Central Dispatch. TASK PARALLELISM

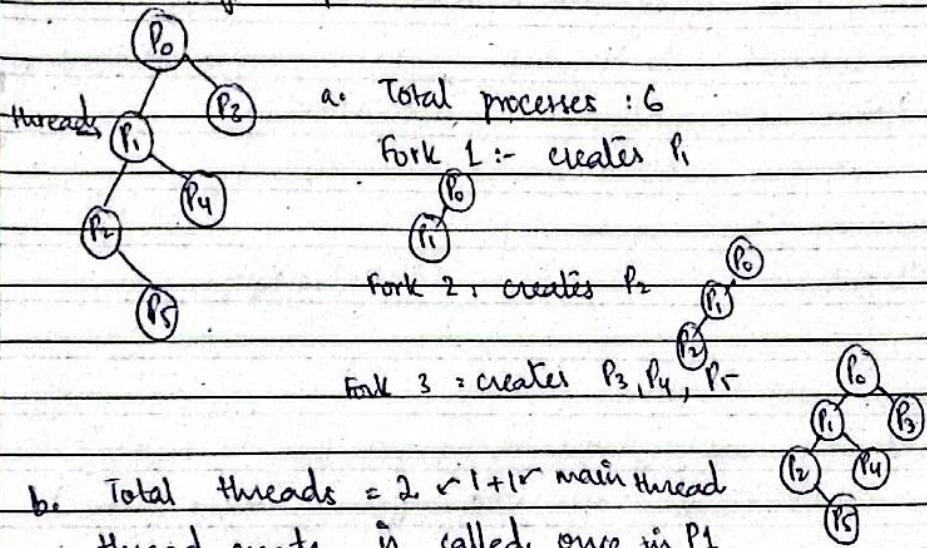
GCD schedules independent tasks / blocks of work. Although it can also exhibit data parallelism depending on usage of developer.

### 4.16 Start CPU-bound , one to one threading

- How many threads will you create to perform the input & output?  
Use 1 thread for input & 1 thread for output. since it is sequential so multiple threads would introduce contention without performance gains.
- How many threads would you use for CPU intensive portion.  
Create 4 threads (one per CPU core).
- The system has 4 cores, in the one to one threading model allows parallel execution.
- Maximizes CPU utilization for the CPU-bound workload.

### 4.17 So how many unique processes are created?

b. How many unique threads are created?



### 4.18 Linux vs. Windows Process / Thread Modeling

Linux treats process, thread

same way, no  
distinguish

Windows treat separately

distinguish processes & threads.

- Linux treats tasks as flexible entities (process-like or thread-like) via `clone()`.
- Windows enforces a strict hierarchy (processes own threads).
- Linux offers more flexibility but less isolation between threads.
- Windows provide stronger isolation but less granularity in resource Pg No. sharing.

4.19

CHILD: value = 5      //line C  
PARENT: value = 0      //line P

Ex-9

4.20

- a. If the number of kernel threads allocated to the program is less than the number of processing cores, then not all processing cores will be utilized. This means that the program will not be able to fully take advantage of the available processing power & its performance maybe limited.
- b. If the number of kernel threads allocated to the program is equal to the number of processing cores, then all processing cores can be utilized. This allows the program to fully take advantage of the available processing power & can result in optimized performance.
- c. If the number of kernel threads allocated to the program is greater than the number of processing cores but less than number of user-level threads, then all processing cores can be utilized & some user-level threads will have to share kernel threads. This can result in improved performance compared to scenario a, but maynot be as efficient as scenario b due to overhead of context switching between user-level threads sharing a kernel thread.