

✓ Present state of computers?

- all of these are parallel computers.

How a program runs in machine?

- fetch instruction
- decode "
- fetch data
- execute instruction
- save results.

- If we can reduce the data fetch time then definitely we can speed up the overall execution time: Cache & CPU registers are to be used.

- Other approach is to execute the instructions in parallel. This can be applied on those instructions which are not dependent. In the following C language instructions:

$$\textcircled{1} \quad a = b + c;$$

$$\textcircled{2} \quad d = c + f;$$

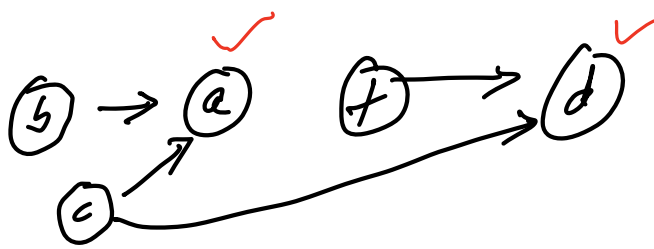
$$\textcircled{3} \quad a = b + c;$$

$$\textcircled{4} \quad d = a + c;$$

↑
[can be
run in
parallel

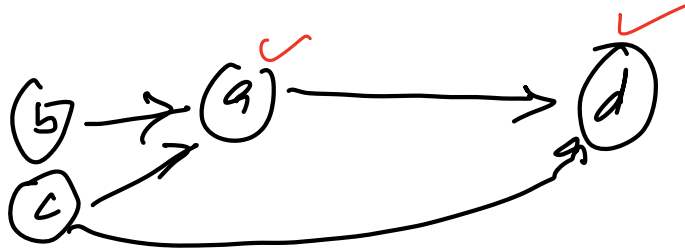
↑
[cannot be

1,2



Shows no dependency between a and d variables.

3,4



dependency of variable d over a.

— Other approach to parallelism is to execute the phases of instructions in parallel.

Fetch, decode, fetch data

This is called pipeline execution or pipeline processing.

- yet other approach is to increase the speed of CPU clock, i.e. crystal's frequency. However, this cannot be done much, because \rightarrow

freq $\uparrow \Rightarrow$ size \downarrow

Freq $\uparrow \Rightarrow$ power dissipation
 \uparrow by square law

\therefore You cannot increase the crystal frequency indefinitely 😞

\therefore The clock rate is kept between 2-3 GHz and number of processing units are increased.

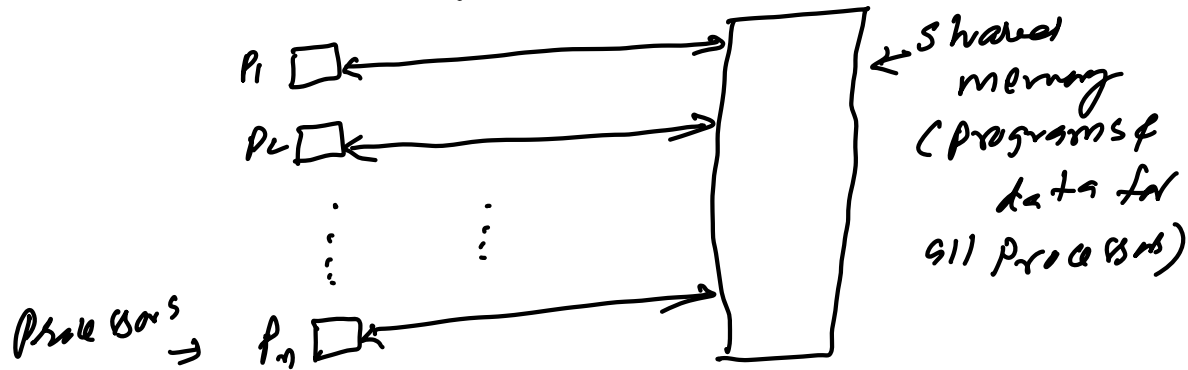
This results to \Rightarrow multi-core CPUs
 \Rightarrow it has been successfully implemented.

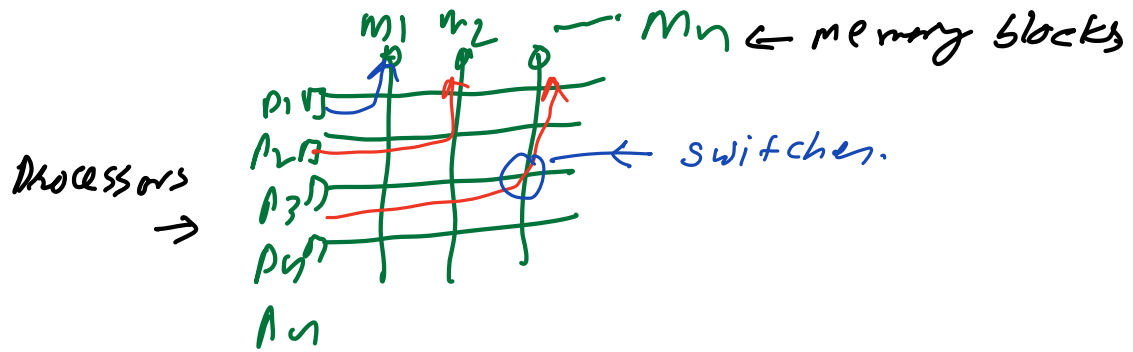
The multicore CPUs requires coordination among themselves to run the programs.

the approach we discussed above is called instruction level parallelism.

Parallelisms at processor levels: ?

① shared memory system





② Distributed systems (Geographically).

↑ Like in social media systems

Ⓐ google search engine

& INDEXES

Ⓑ FB, LinkedIn, ...

③ GPUs (Graphic Processor Units)

Challenge: Programs: How to write
such programs, which can
exploit the parallelism in the processors
and speed up the processing?

This is possible:

- using threads in shared memory
- through message passing, in
distributed system
- GPUS

