

clock frequency = how fast a processor can perform tasks

1 GHz = processor ticks 1 billion times per second

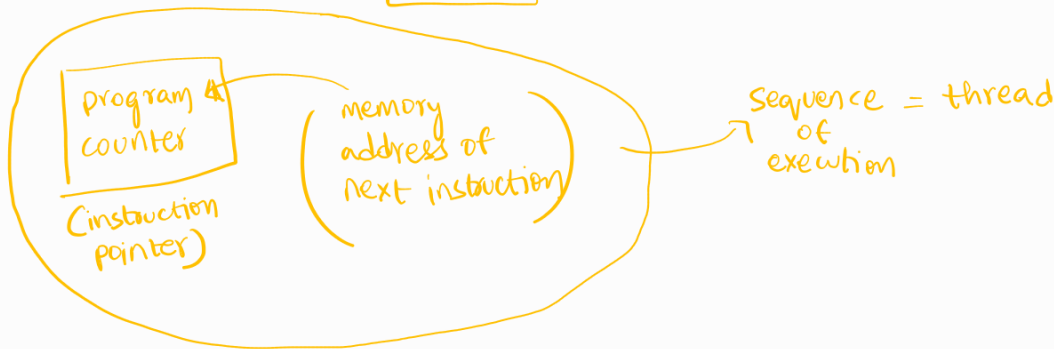
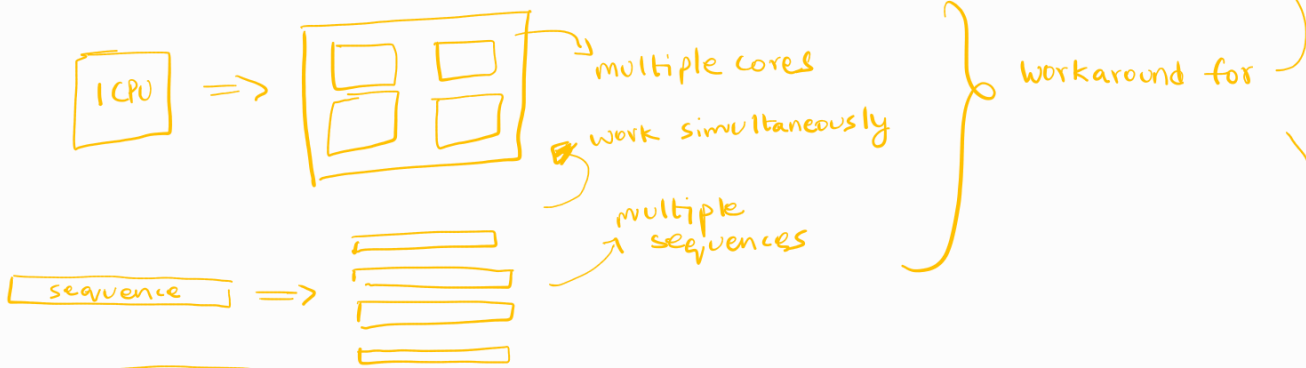
GFLOPS = 10^9 floating point operations per second

TFLOPS = 10^{10} " " "

1980s
1990s



But increase in speed \Rightarrow increase in heat & energy dissipation issues.



parallel programs = multiple threads together = speed stonk

parallel
over
sequential
=
CONCURRENTLY
REVOLUTION

x86 = CPU design by Intel (8086 chip debut)

x86 instruction set = collection of machine language instructions that x86 compatible processors can understand & execute.

(arithmetic calculations, data movement, control flow, I/O interactions)

Backward compatible.

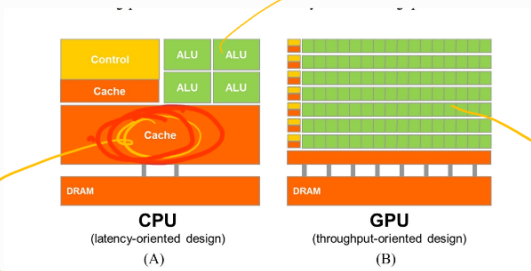
hyperthreading = ONE physical CPU acting like Two virtual CPU.

DOESNT double the speed, but makes multitasking smooth.

optimized for sequential code perf.
reduced latency at the cost of chip size & power per unit.

latency oriented design.

maximize # of floating point operations & memory access throughput.



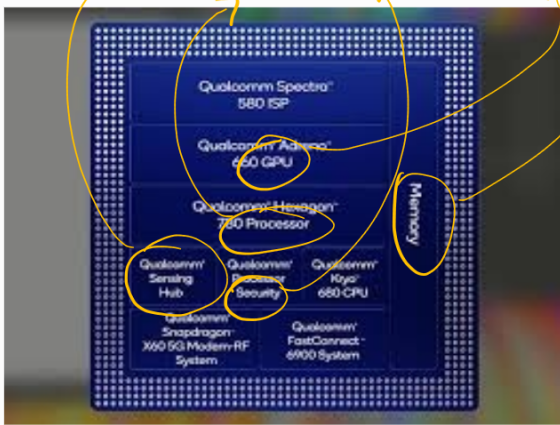
stores frequently accessed data.



GPU: lots of ops = better graphics.

multicore introduced, but not anymore. threaded apps

Processors are not just processors anymore, they are processor + memory + gpu + security....



Moore's law:- the # of transistors double every 2 years at the cost of computing being halved.

denard's law:-

as the size of transistor shrink:-

- Voltage reduces proportionally (V)
- current " " (I)
- capacitance " " (C)

$$\text{delay} = \frac{CV}{I}$$

as delay ↓, performance ↑ and scales up by α^2 .

⇒ Same 10W chip is more performant few years later.

- CPUs are latency driven
- GPUs are throughput oriented; get as much done in the silicon as possible; cores as simple as possible but soooo many AUs that are energy efficient. massive # of threads.

Amdahl's law:- % of speedup that can be achieved overall is limited by % of code that can be parallelizable.

→ Another limitation is how often memory needs to be accessed from DRAM. Usually done by using the memory in. GPU.

→ OpenMP & MPI (message passing interface) are other interfaces.

OpenCL (nvidia + AMD + intel etc.,) is another interface but relies heavily on APIs.