

MAJOR TRENDS:

- Scaling laws breaking down
- Energy efficiency issues

SOLUTION:

→ SPECIALIZATION:

- * Limit hardware capabilities
- * Expose hardware to software using domain-specific interface.

WHY USEFUL FOR DEEP LEARNING:

- Parallelism
- Regular memory access
- Precision not needed
- Most values are 0.
- Weights, inputs reused.

DENNARD'S LAW:

Decrease area, voltage \rightarrow Current decreases
of transistor

- * Increases frequency
(faster chips)
- * Power dissipation decreases
- * Power density = $\frac{\text{Power dissipation}}{\text{Area}}$
 $= \text{Constant}$.

So, we lose same power for a given area, but
now that area of transistors is faster!

So for same energy \Rightarrow Faster chips!

This is gone!

Cannot decrease area/size of transistor

any more \rightarrow \uparrow heat

Already 3nm (DNA size range!)

SOLUTION:

1. Increase cores?

Parallel processing with more cores \rightarrow Faster Chips!

But each core takes E energy. n cores take nE energy. So it's not really increased speed with same energy.

Performance / Energy is same!

2. Specialize

Some generic H/W can be removed. So this can make performance better and decrease energy needs!

REGISTER vs SCRATCHPAD MEMORY vs CACHE vs MEMORY vs DISK

DISKS: SLOW, persistent

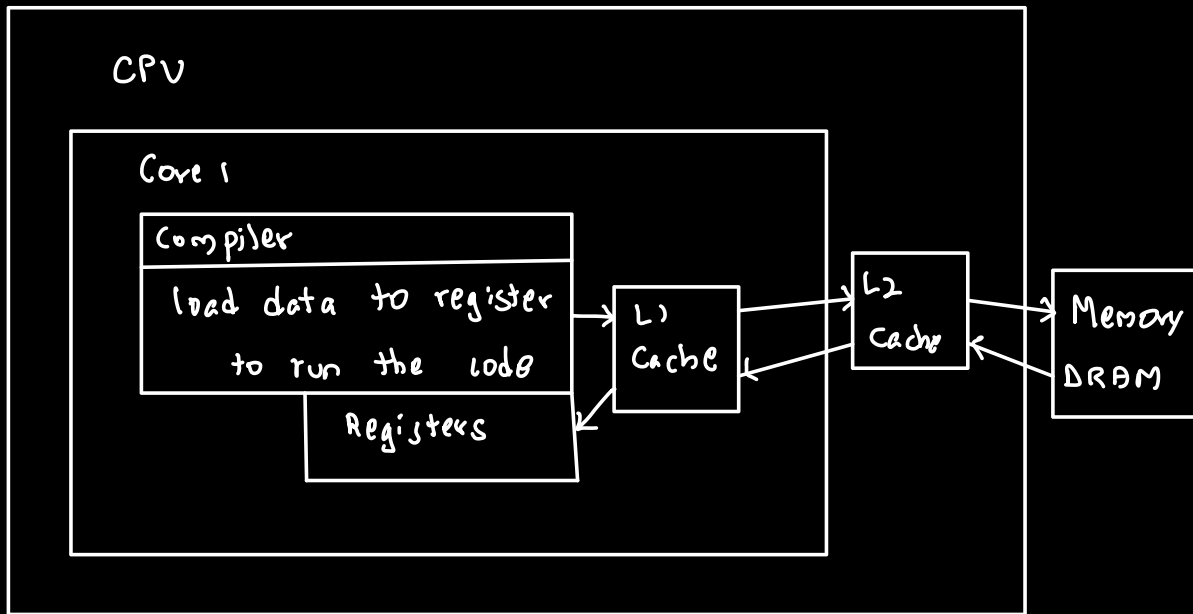
MEMORY: DRAM

CACHE: SRAM

Compilers decide what's in the register while H/W decides what's in the cache.

SCRATCHPAD MEMORY:

Similarly, Programmer / Compiler / S/W decides what's in the scratchpad. Other than that it is basically just a cache.



ADVANTAGES OF CACHE:

→ ↓ latency

→ ↓ energy

→ Memory bandwidth is low. Say 8 cores access 1 memory.

It cannot take all requests. Cache manages them.