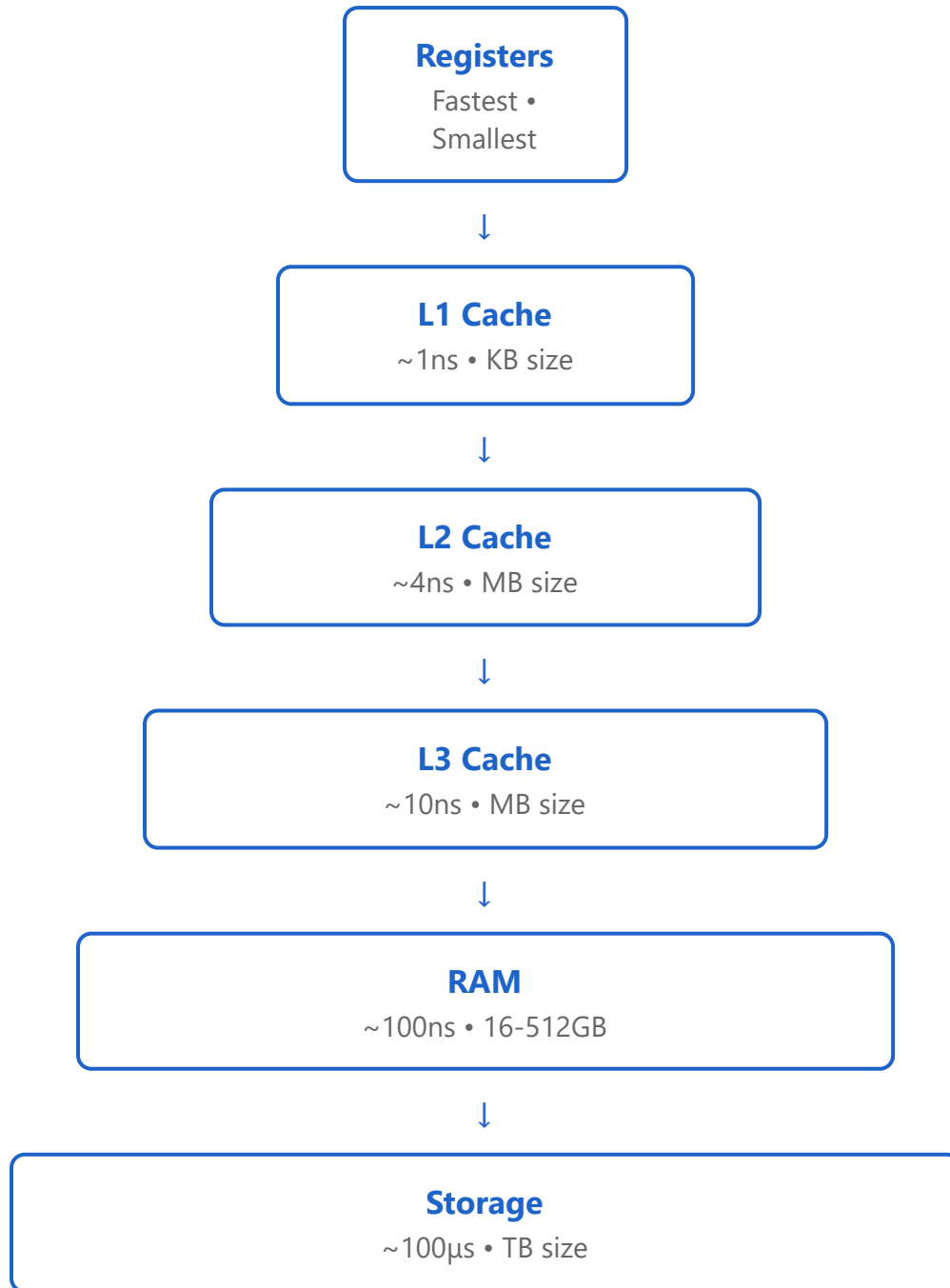
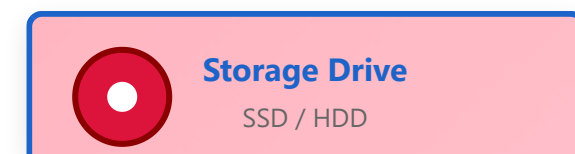
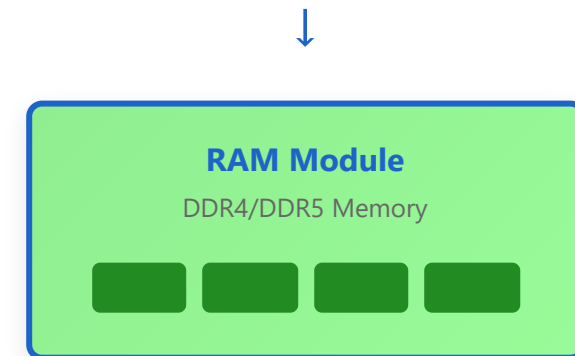
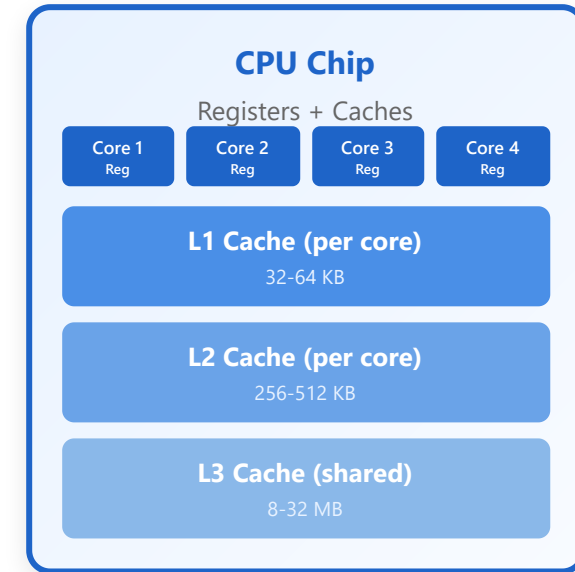


Memory Hierarchy - RAM, VRAM, Cache



Physical Hardware Components



Trade-off: **Faster = Smaller + Higher Cost**

SSD: Flash Memory
HDD: Magnetic Disk

RAM

16-512GB
~20GB/s bandwidth
System memory

VRAM (GPU)

8-80GB
~1-2TB/s bandwidth
HBM2/HBM3

Cache

L1: KB
L2: MB
L3: MB

ML Bottleneck:

RAM



VRAM

Data transfer overhead