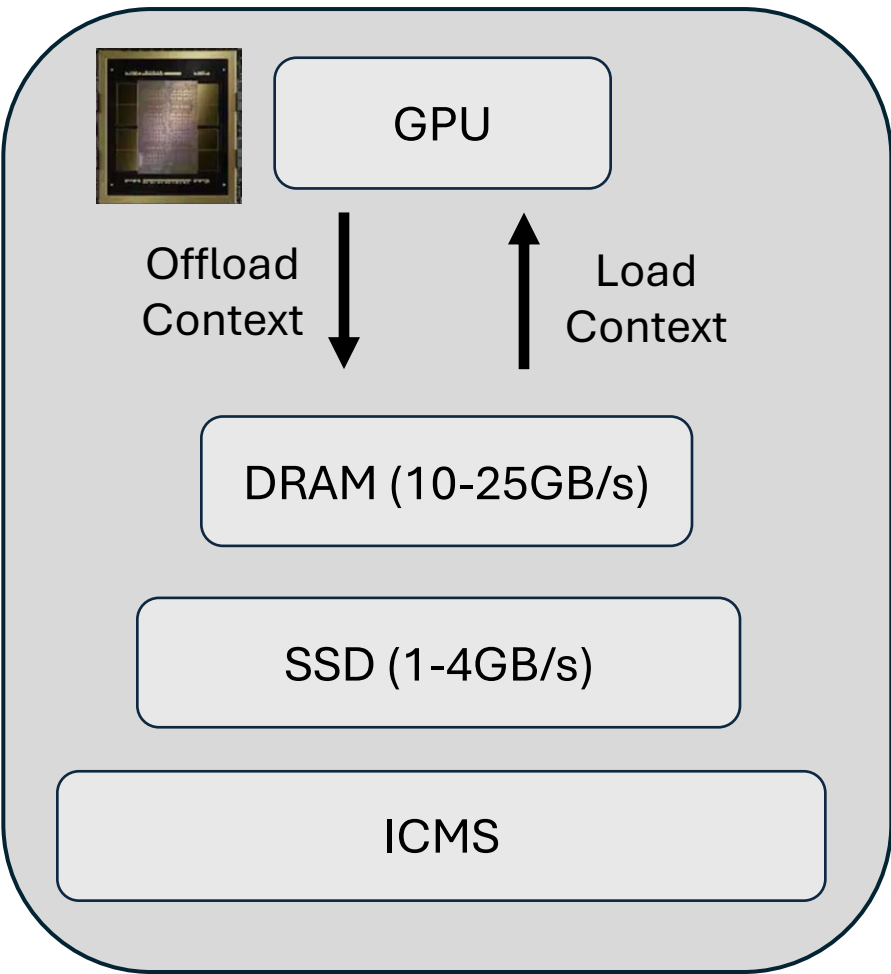


Prefill Node (Compute Optimized)



Decode Node (Memory Optimized)

