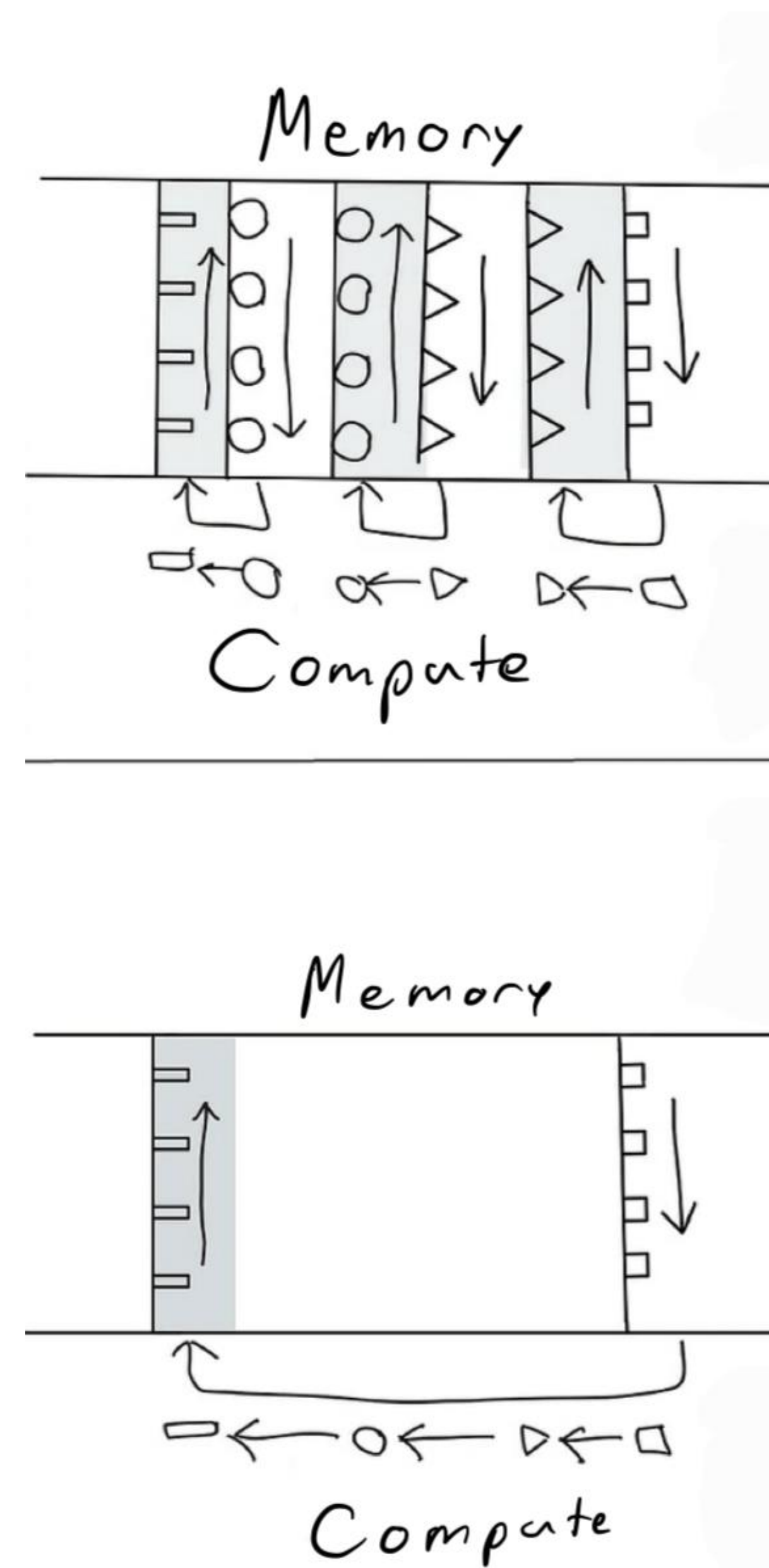
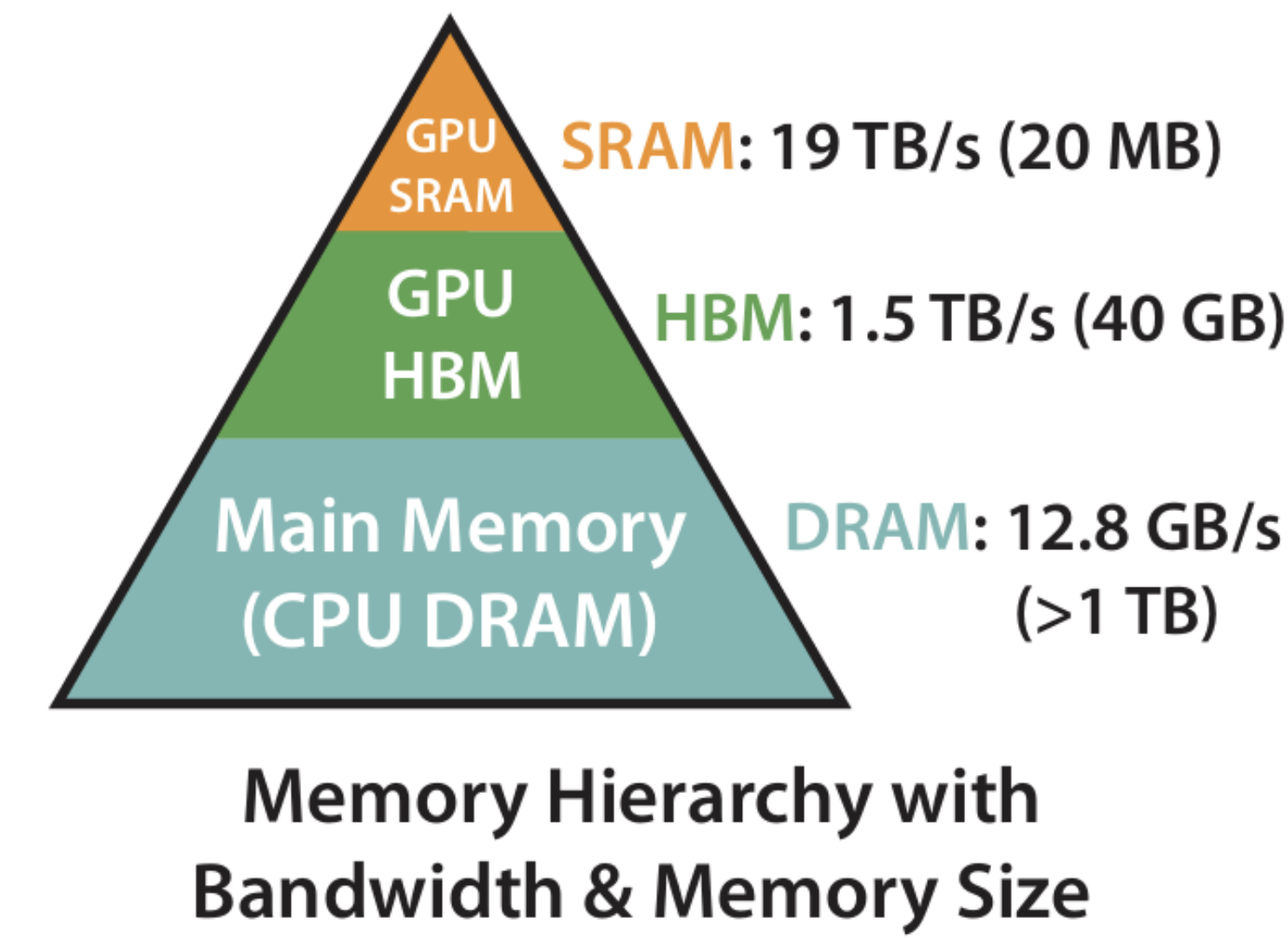
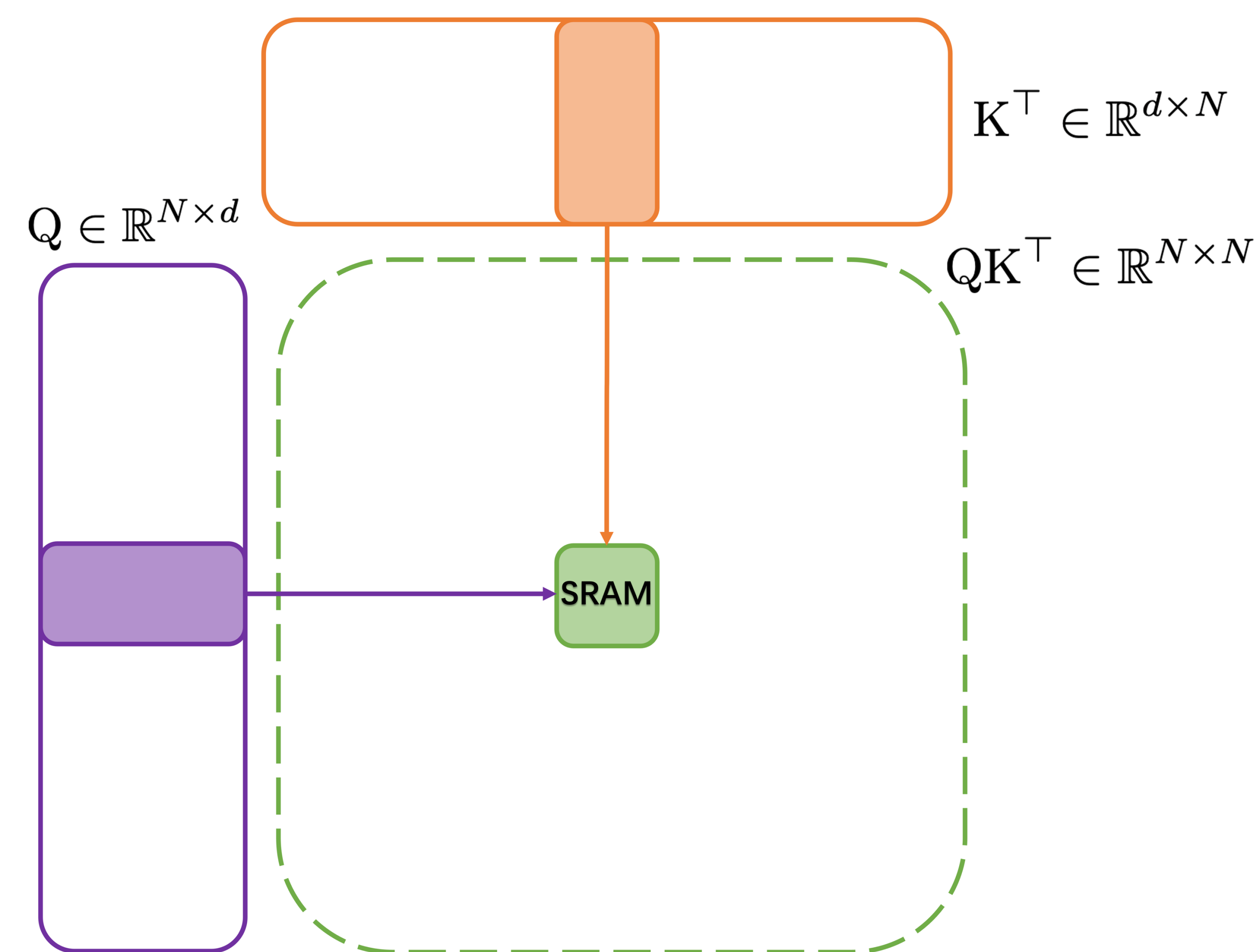


Motivation

- **Attention** is **slow** on long sequences, $O(N^2)$
- Limitations to deep-learning operations:
 - Compute — # of operations
 - Memory — HBM accesses
- **Key Observation:** Attention is **memory-bound**
- Can reduce # HBM accesses to improve speed

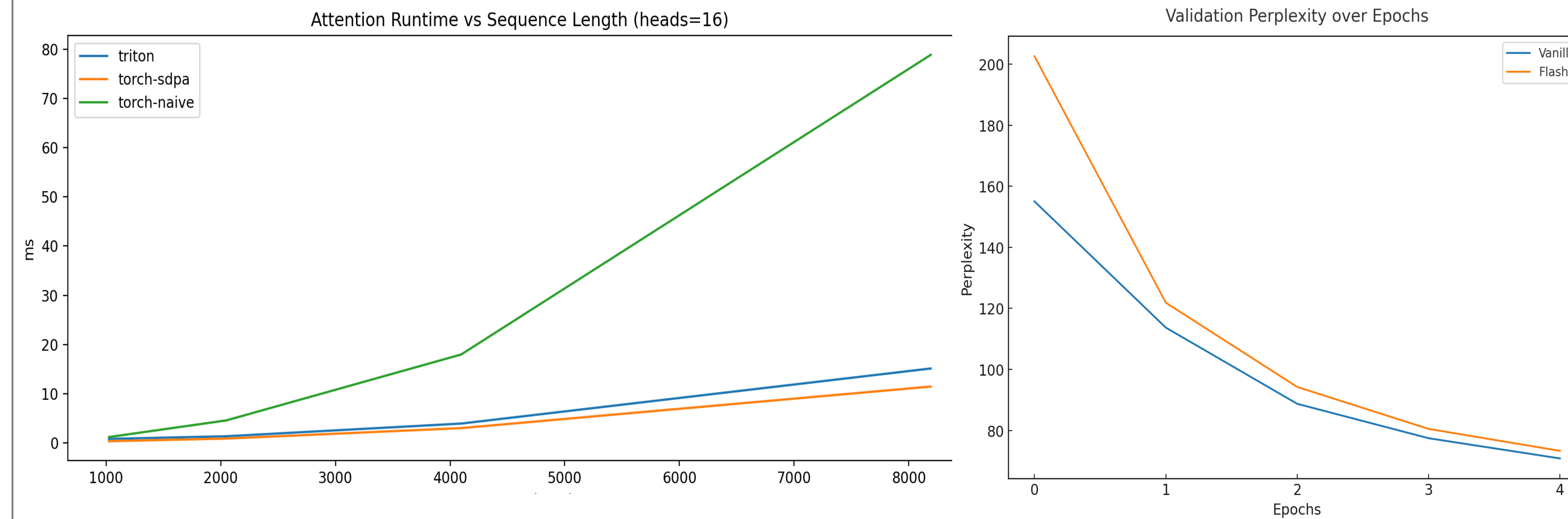
Strategy: Tiling, Fusion



- **Tiling:** Split Q , K , V into blocks to fit into SRAM
- **Streaming:** $x = [x^{(1)} \ x^{(2)}] \in \mathbb{R}^{2B}$ $m(x) = m([x^{(1)} \ x^{(2)}]) = \max(m(x^{(1)}), m(x^{(2)}))$
 $f(x) = \begin{bmatrix} e^{m(x^{(1)})-m(x)} f(x^{(1)}) & e^{m(x^{(2)})-m(x)} f(x^{(2)}) \end{bmatrix}$ $\text{softmax}(x) = \frac{f(x)}{\ell(x)}$
 $\ell(x) = \ell([x^{(1)} \ x^{(2)}]) = e^{m(x^{(1)})-m(x)} \ell(x^{(1)}) + e^{m(x^{(2)})-m(x)} \ell(x^{(2)})$
- **Fusion:** Combine scaling, masking, softmax into one kernel to minimize HBM traffic
- **Reduces HBM access** from $O(Nd + N^2)$ to $O(N^2 d^2 / M)$ where M is SRAM memory size

Experiments

Triton Kernel Performance



Extension: Multi-Query Attention

