

Week 19

# Summary of the New MatMul

A: d1xd2, B: d2xd3, C = A \* B: d1xd3

- `matmul_t`: unit-stride load A
  - `matmul_t(A^T, B) = C^T`
- `matmul`: strided load A
  - `matmul(A, B) = C^T`
- Negligible performance difference

# Summary of the New MatMul

- From Linear Algebra:

$$C = A * B$$
$$C^T = (A * B)^T = B^T * A^T$$



$$\text{matmul}(B^T, A^T) = (C^T)^T = C$$

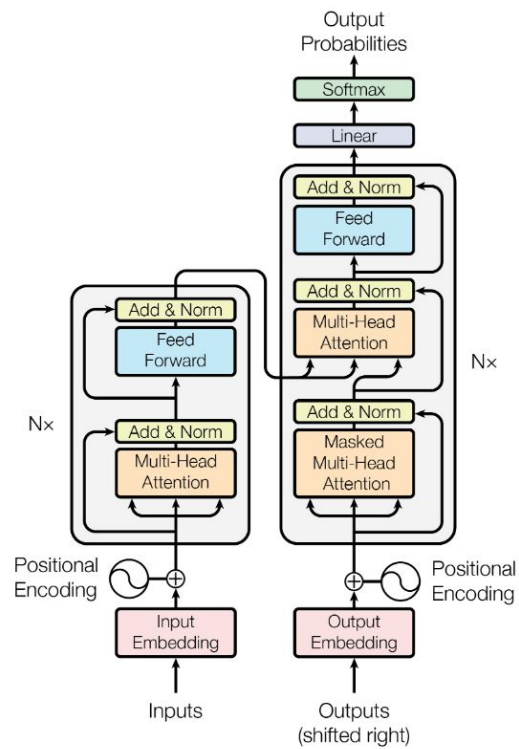


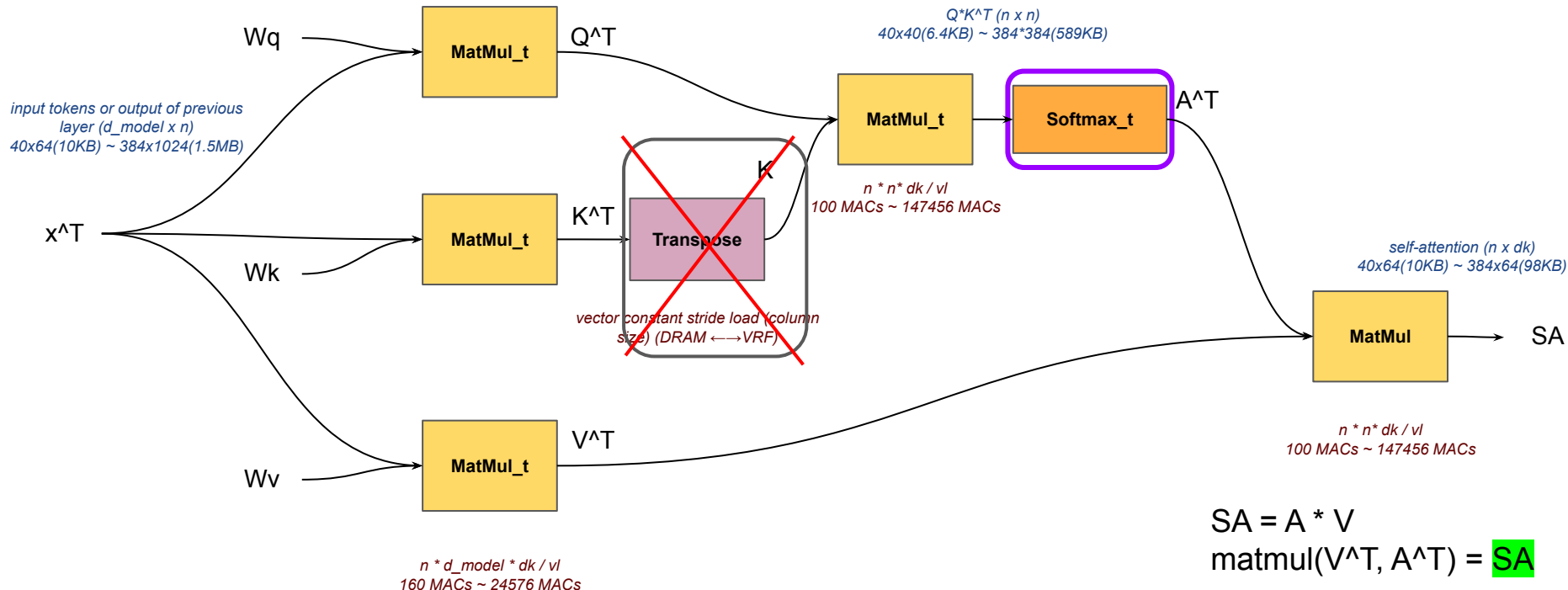
Figure 1: The Transformer - model architecture.

weight matrices ( $d_{\text{model}} \times d_k$ )  
 $64 \times 64 (16\text{KB}) \sim 1280 \times 80 (409\text{KB})$

$Q, K, V (d_k \times n)$   
 $40 \times 64 (10\text{KB}) \sim 384 \times 64 (98\text{KB})$

$$A = Q * K^T$$

$$\text{matmul\_t}(Q^T, K^T) = A^T$$



$$\text{matmul}(B^T, A^T) = C$$

$$SA = A * V$$

$$\text{matmul}(V^T, A^T) = \text{SA}$$

