



FlashAttention-2

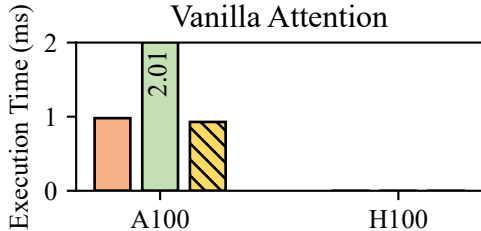


TensorRT



TA

Vanilla Attention



Gemma2

