

More generally: dropping bias terms

Most modern transformers don't have bias terms.

Original Transformer:

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

Most implementations (if they're not gated):

$$\text{FFN}(x) = \sigma(xW_1)W_2$$

Reasons: memory (similar to RMSnorm) and optimization stability