



Figure 1: Transferring Transformer to Mamba. Weights in the same color are initialized from the transformer (Linear projections for Q, K, and V are initialized using linear projection for C, B, and X respectively). We replace individual attention blocks with Mamba blocks, and then finetune Mamba blocks while freezing the MLP blocks. Shapes are kept mainly the same. New parameters are introduced for the learned \mathbf{A} and Δ parameters.

Derivations from softmax-based attention to linear RNNs:

Attention is computed in parallel for multiple differently parameterized heads $h \in \{1 \dots H\}$. Each head takes sequence \mathbf{o} with hidden size D as an argument and computes,

$$\begin{aligned} \mathbf{Q}_t &= \mathbf{W}^Q \mathbf{o}_t, & \mathbf{K}_t &= \mathbf{W}^K \mathbf{o}_t, & \mathbf{V}_t &= \mathbf{W}^V \mathbf{o}_t & \text{ for all } t, \\ \alpha_1 \dots \alpha_T &= \text{softmax}\left(\frac{[m_{1,t} \mathbf{Q}_t^\top \mathbf{K}_1 \dots m_{T,t} \mathbf{Q}_t^\top \mathbf{K}_T]}{\sqrt{D}}\right) & \mathbf{y}_t &= \sum_{s=1}^t \alpha_s \mathbf{V}_s \\ \text{where } \mathbf{o}_t &\in \mathbb{R}^{D \times 1}, & \mathbf{W} &\in \mathbb{R}^{N \times D} & \mathbf{Q}_t, \mathbf{K}_t, \mathbf{V}_t &\in \mathbb{R}^{N \times 1} & m_{s,t} = \mathbf{1}(s \leq t) \end{aligned}$$

Despite the superficially different form, there is a natural relationship between linear RNNs and attention. Linearizing the attention formula by removing the softmax yields:

$$\mathbf{y}_t = \sum_{s=1}^t m_{s,t} \alpha_s \mathbf{V}_s = \frac{1}{\sqrt{D}} \mathbf{Q}_t \sum_{s=1}^t (m_{s,t} \mathbf{K}_s^\top \mathbf{V}_s) = \frac{1}{\sqrt{D}} \mathbf{Q}_t \sum_{s=1}^t m_{s,t} \mathbf{K}_s^\top \mathbf{W}^V \mathbf{o}_s.$$

On the other side, expanding the recursive form of Mamba yields,

$$\begin{aligned} \mathbf{h}_t &= \mathbf{A}_t \mathbf{A}_{t-1} \dots \mathbf{A}_2 \mathbf{B}_1 \mathbf{x}_1 + \mathbf{A}_t \mathbf{A}_{t-1} \dots \mathbf{A}_2 \mathbf{B}_2 \mathbf{x}_2 + \dots + \mathbf{B}_t \mathbf{x}_t = \sum_{s=1}^t \mathbf{A}_{t:s}^\times \mathbf{B}_s \mathbf{x}_s \\ \mathbf{y}_t &= \mathbf{C}_t \sum_{s=1}^t \mathbf{A}_{t:s}^\times \mathbf{B}_s \mathbf{x}_s \end{aligned}$$

where $\mathbf{A}_{t:s}^\times = \mathbf{A}_t \mathbf{A}_{t-1} \dots \mathbf{A}_s$

If the two models have the same heads H and the head size N , we can set $\mathbf{B}_t = \mathbf{K}_t$, $\mathbf{C}_t = \mathbf{Q}_t$, $\mathbf{x}_t = \mathbf{W}^V \mathbf{o}_t$. This relationship motivates moving between attention and linear RNN representations.