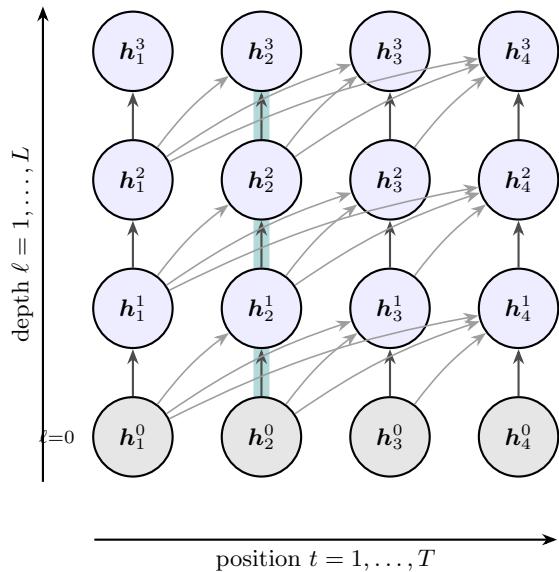
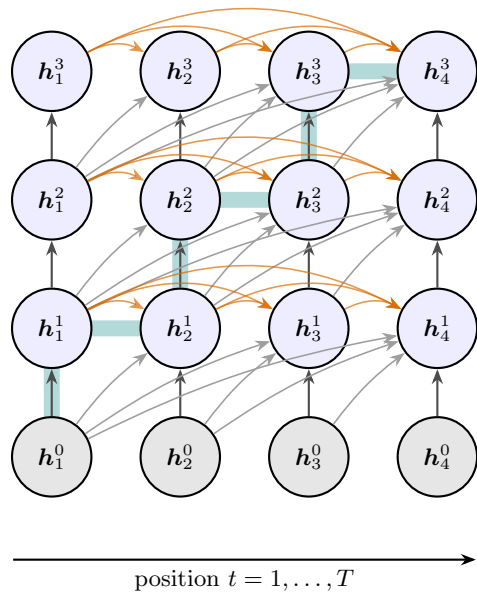


residual stream h_t^ℓ : superscript $\ell=1, \dots, L$ is the layer (depth), subscript $t=1, \dots, T$ is the position

(a) attention only
opaque serial depth $\sim L$



(b) + horizontal attention
opaque serial depth $\sim L + T$



\rightarrow vertical attention $h^{\ell+1} = \text{Att}(h^\ell)$ \rightarrow horizontal attention $h^\ell = \text{Att}(h^\ell)$ \blacksquare longest opaque serial path