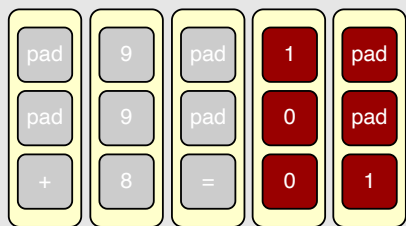
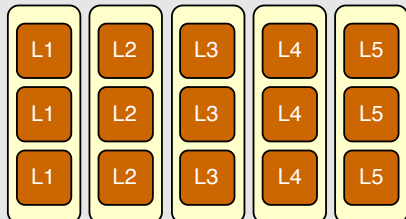


## COPY THEN ATTENTION



m Attention Layers  
(no MLP)  
(sliding window)



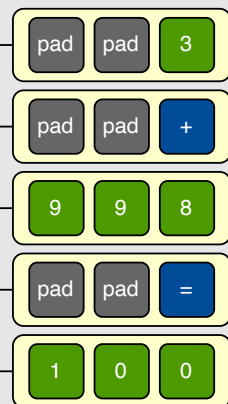
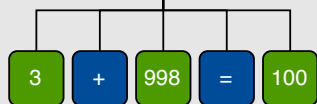
Repeat



n Transformer Blocks

Cross-Attention

Self-Attention



## CROSS ATTENTION TO PAD TOKENS

