

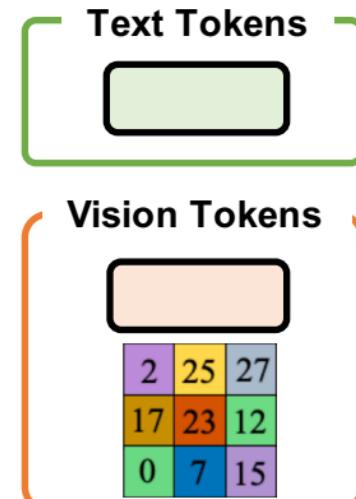
Pre-train  
VL-Align



Post-train  
World Model



Fine-tune  
Policy Learning



# Autoregressive Transformer



Task  
Pick carrot



...

Causal Multimodal Sequence