

Context (Instruction)

"pick up the *orange juice* and place it in the basket"

Observation (Visual Inputs)

Wrist Camera



...

Head Camera

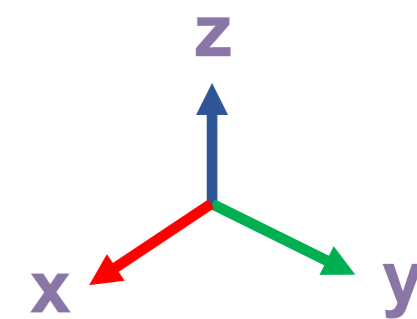


Robot State (Proprioception)

Joint Position

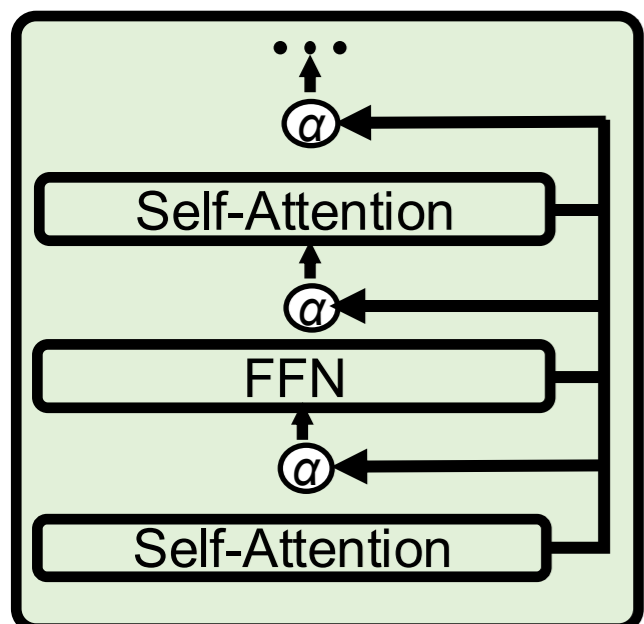
$$\begin{bmatrix} q_1 \\ q_2 \\ \dots \\ q_n \end{bmatrix}$$

EEF Pose



Multi-modal Tokenizer

In-context Condition



Mean Velocity Predictor

$$\mu_{\theta}(z_t, r, t)$$



Action Head

iMF Update

$$x_0 = x_t - \mu_{\theta}(z_t, 0, 1)$$

Pseudo-Huber loss



Text Token



Vision Token



State Token



Noisy Action



Denoised Action



Real Action



Robot Execution