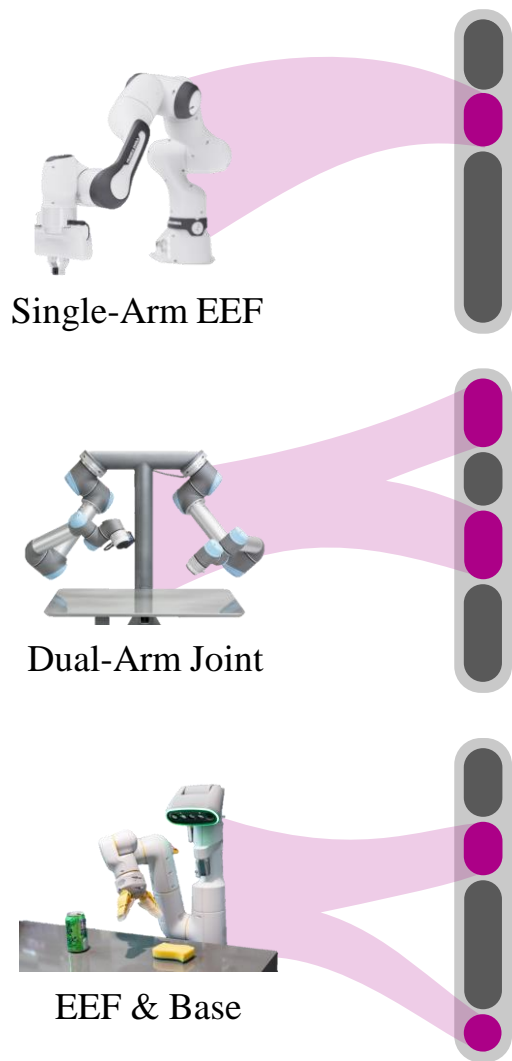


Unified Action Space



Low-Dimensional Inputs



Embed

\mathbf{z}_t & $\tilde{\mathbf{a}}_{t:t+T_a}$

c & Diff.
Timestep k

Unified Action Space

MLP

Concat.

Fourier MLPs

$L \times$

DiT Block with
Cross-Attention

Alternatively
Inject Image/Lang.

Norm & MLP

Outputs

Denoised Act. $\mathbf{a}_{t:t+T_a}$

Image Inputs $\mathbf{X}_{t-1:t+1}$



Exterior \mathbf{X}^1

Right-Wrist \mathbf{X}^2

Left-Wrist \mathbf{X}^3

SigLIP

Tokens with
Multi-Dim. Pos. Emb. \mathbf{x}_{t-1}^1 \mathbf{x}_{t-1}^2 \mathbf{x}_{t-1}^3
 \mathbf{x}_t^1 \mathbf{x}_t^2 \mathbf{x}_t^3

Tokens & Mask

T5-XXL

Language Inputs ℓ

*“Insert the lemon slice on the paper cup
into the goblet rim.”*