

Community datasets



SmolVLA

Affordable robots

LeRobot



Vision-Language Model

Self-Attention

Self-Attention

Self-Attention

⋮



Task: Grasp the object and put it in the bin

State



$a_t a_{t+1} \dots a_{t+h}$

Action Expert

Cross-Attention

Self-Attention

Cross-Attention

⋮

Noised Action