Image Observation

Image Tokens

Encode
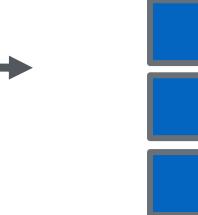
Vision-Language Model

**System 2**

Language Instruction

"Pick up the industry object and place in yellow bin."

Tokenize

Text Tokens

Diffusion Transformer

**System 1**

Motor Action

Joint Positions

Joint Velocities

Base Position

EEF Poses ...

Robot State

Encode

Action Tokens

Denoising