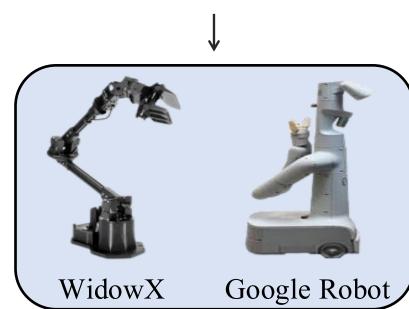
Co-Training Datasets "Identify ..." "Pick up the can" Robot Data Spatial Affordance "Describe ..." "The vase ..." Frozen ** Encoder Spatial Reasoning General VQA ■ Robot Data Spatial Reasoning Spatial 0.12 Affordance General VQA

Pre-Trained LLM Dual Visual Encoders "pick up the can" Trainable 0

Encoder



String-Based

Action Tokenizer

"()"

"