

---

- Mobile Manipulation**

Human: Bring me the rice chips from the drawer. Robot: 1. Go to the drawers, 2. Open top drawer. I see **<img>**. 3. Pick the green rice chip bag from the drawer and place it on the counter.

**Visual Q&A, Captioning ...**

Given **<img>**. Q: What's in the image? Answer in emojis.  
A: 🍌🍇🍓🥑🍎🍋

Describe the following **<img>**:  
A dog jumping over a hurdle at a dog show.

**PaLM-E: An Embodied Multimodal Language Model**

Given **<emb>** ... **<img>** Q: How to grasp blue block? A: First, grasp yellow block

Large Language Model (PaLM)

Control ← A: First, grasp yellow block and ...

**Language Only Tasks**

Here is a Haiku about embodied language models:  
Embodied language models are the future of natural language

Q: Miami Beach borders which ocean? A: Atlantic.  
Q: What is 372 x 18? A: 6696.  
Language models trained on robot sensor data can be used to guide a robot's actions.

**Task and Motion Planning**

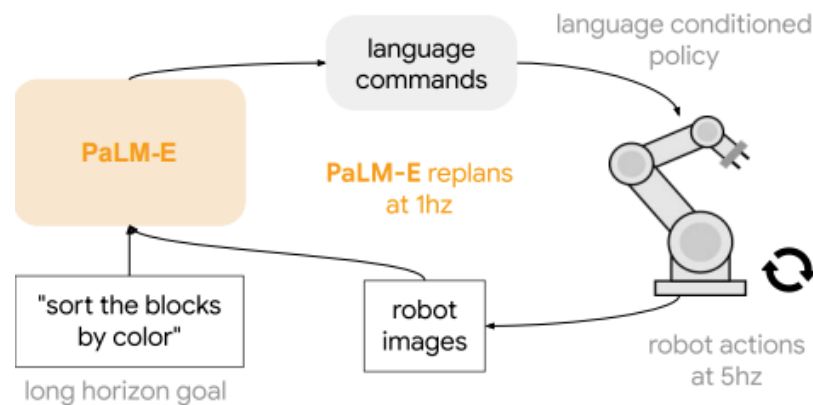
Given **<emb>** Q: How to grasp blue block?  
A: First grasp yellow block and place it on the table, then grasp the blue block.

**Tabletop Manipulation**

Given **<img>** Task: Sort colors into corners.  
Step 1. Push the green star to the bottom left.  
Step 2. Push the green star circle to the green star.

- 1 / 2

**decisions in natural text**) -> Low-level policy or planner (translate language decisions into actions)



#### Terminologies

- $w_i \in \mathcal{W}$  -> a token
- $I$  -> an image
- $X \in \mathbb{R}^k$  -> embedding space (if a token's length is 3, then the embedding space is a 3-dimensional space)
- $\gamma : \mathcal{W} \rightarrow \mathcal{X}$  -> **Text encoder**: a LLM that embeds a token  $w_i$  into a word token embedding space. PaLM model is employed in this work.
- $\phi : \mathcal{O} \rightarrow \mathcal{X}^q$  -> **Vision encoder**: an encoder  $\phi$  that maps a continuous observation space  $\mathcal{O}$  into a sequence of  $q$ -many vectors in  $\mathcal{X}$ . ViT-22B and ViT-4B are employed in this work.
- $\phi_{\text{state}}$  -> State estimation encoder (a MLP that maps inputs into the language embedding space).
- $D = \left\{ \left( I_{1:L_i}^i, w_{1:L_i}^i, n_i \right) \right\}_{i=1}^N$  -> Training dataset that contains  $u_i$  images, a text  $w_{1:L_i}$ , and an index  $n_i$ .
  - $u_i$  images -> Continuous observations
  - $w_{1:L_i}^i$  -> Text contains a prefix part (i.e., prompt template) formed from multi-modal sentences and a prediction target that only contains text tokens.