Give a brief description of the given observation / subsequent keyframes / short demonstrations.

**Image Encoder** ❄

**Projection** 🔥

**Text Encoder** ❄

**Multi-modal LLM** ❄

<Scene Description> :
On the table, there are two bigger green blocks, one smaller red bowls on the left bottom corner ...
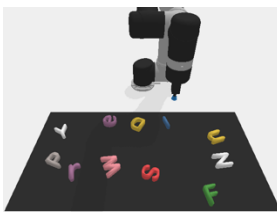
<Action Recognition>:
Move the smaller blue block into the black bowl.

<Video Understanding>:
The given demonstration is to first put smaller blue block onto the smaller red block and then ...
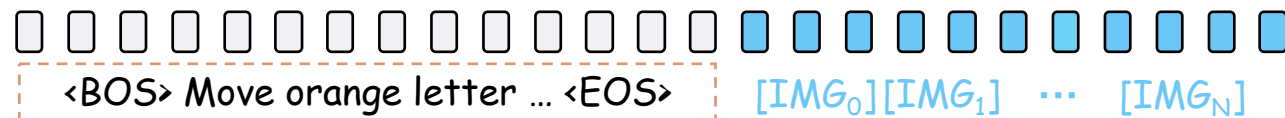
**Observation**

**Task instruction**

Spell a word using letters on the table for the name of top CS conference.

**Prompt**

Please give a brief stepwise instruction to do next and the corresponding goal image.

**Multi-modal LLM** 🔥

<BOS> Move orange letter ... <EOS>   [IMG₀][IMG₁] ··· [IMGₙ]

**Output Projection** 🔥

$L_{Reason}$

$L_{Imagine}$

**Diffusion Model** 🔥

Move orange letter N to the bottom left area

$L_{Alignment}$

**Multi-modal LLM**

Action

Recognition