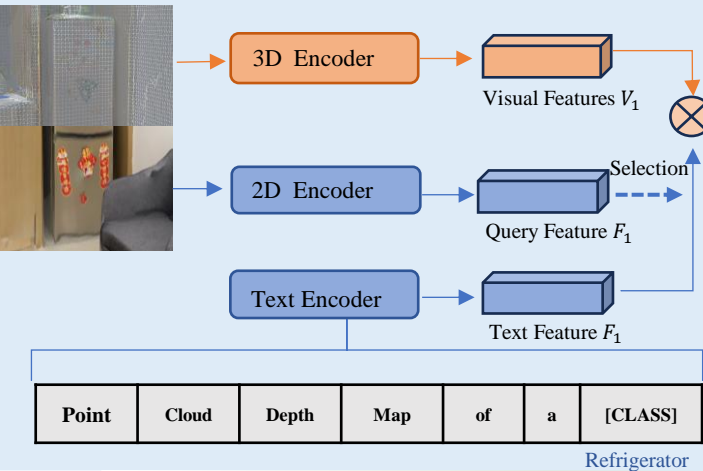
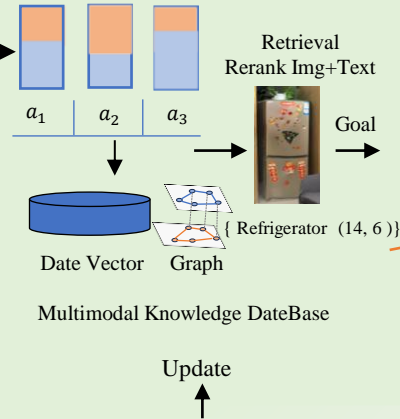


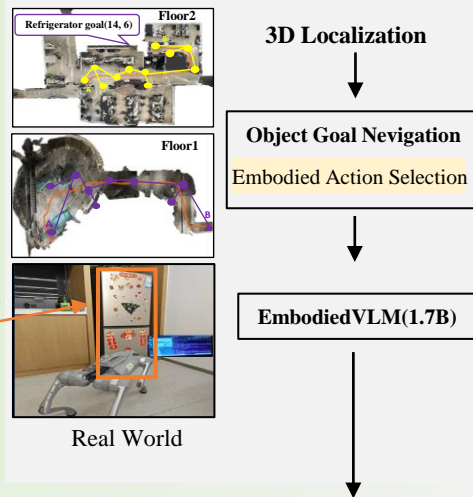
① Input



② Multimodal RAG



③ Embodied Planning



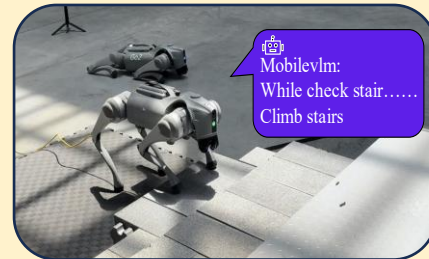
Embodied Action Selection

SFT-InternLM

PPO-Clip

Action Instruction:

- ☒ Climb Stairs
- ☐ Turn Right
- ☐ Turn Left
- ☐ Move Forward
- ☐ Down Stairs
- ☐ StandDown
-



I. Text Input: I want to go to the refrigerator

II. Output Result: This is a decorated refrigerator, located at coordinates (14, 6)