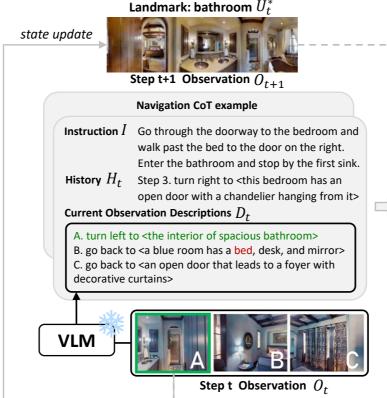
## **Navigation CoT example**

### **Navigation Input**



### **0-shot NavCoT Output:**

- 1. Imagination: bed.
- 2. Filtered observation:
- B matches the imagination.
- 3. Action: B.



### In-domain training

## **Finetuned NavCoT Output:**

# World Model

Future Imagination  $U_t$ 

bathroom

LLaMA 2

**Reasoning Agent** 



Visual Information Filter  $V_t$ A matches the imagination

Action Prediction  $a_t$  A

In-domain training

 $a_t^*$