# INTEGRATING COMMON SENSE AND PLANNING WITH LARGE LANGUAGE MODELS FOR ROOM TIDYING
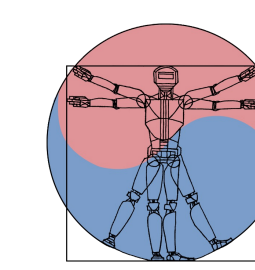
**Zhanxin Wu** [zhanxinwu@u.nus.edu], **Bo Ai** [bo.ai@u.nus.edu], **David Hsu** [dyhsu@comp.nus.edu.sg]

**ROBOTICS SCIENCE AND SYSTEMS**

**NUS** National University of Singapore

## Motivation

**Do you want a personal robot housekeeper?**

Given partial textual description of the layout from humans and description of objects, we endow robots with the capability of tidying up a room.

This task has three challenges:
- **Incomplete map information** in the description
- **Commonsense understanding** of object locations
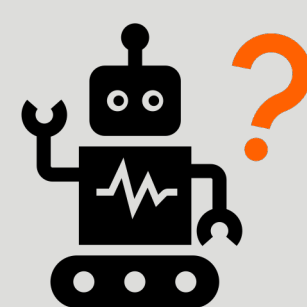- **Long-horizon planning** for room tidying

We provides preliminary evidence that *LLMs have common sense about the **spatial layout of human-living environments** and **object arrangements**.*
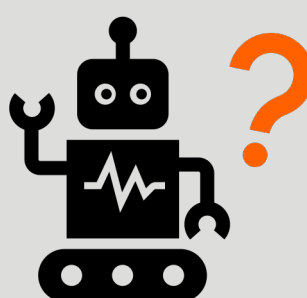
## Problem Formulation

Hi, my housekeeper! From the **living room**, the **kitchen** is on the right side. There is a **plate** on the sofa in the living room. Please **tidy up the living room**.
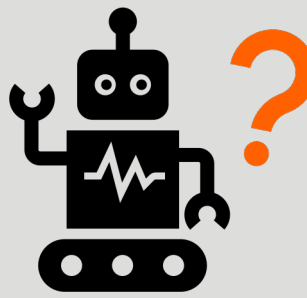
What should I do?

Please move the **plate** from the **living room sofa** to the **dining room table**.

But where is the **dining room**?

The **dining room** is expected to be connected to the **kitchen**. Go to find it!

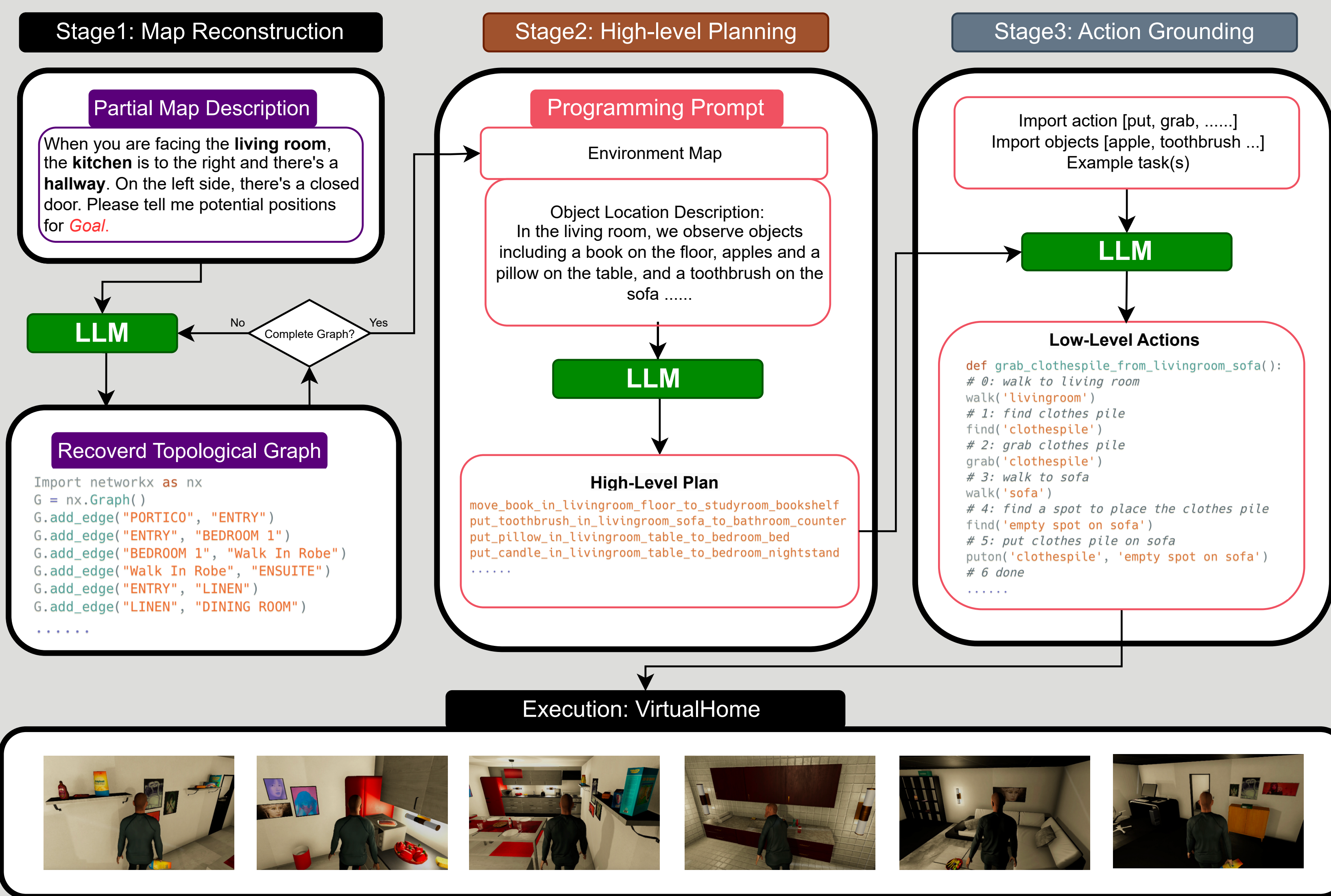I find it, Please provide me with the steps to rearrange the **plate**.

Step1: Walk to living room.
Step2: Find the sofa.
......

- **Assumption**: (i) Semantic labels for each room in given map are provided. (ii) The executable actions for the agent are predefined.
- **User Input**: Textual descriptions of partial map and textual descriptions of objects in the room.
- **System output**: Executable action sequences for the agent to tidy up the room.

## System Architecture

The framework has three stages: (i) predicting spatial positions for unseen destination, (ii) generating a high-level plan for relocating misplaced objects, and (iii) grounding the plan into executable actions.



## Map Reconstruction

**Number of Interaction Rounds (NIR) Required to Recover Missing Places**

| Environment | #Places | Left-Out Places | NIR | |
| --- | --- | --- | --- | --- |
| | | | Ours | Random Guess |
| VH Apartment | 5 | Bathroom | **1.20** ± 0.45 | 2.82 ± 1.50 |
| | | Bedroom | **1.60** ± 0.55 | 3.32 ± 1.43 |
| Real Apartment | 15 | Bathroom | **3.20** ± 1.30 | 8.00 ± 4.56 |
| | | Bedroom | **2.40** ± 0.55 | 7.20 ± 4.01 |
| Hospital | 20 | Nurse's Station | **1.40** ± 0.55 | 7.60 ± 5.64 |
| | | Bathroom | **2.20** ± 2.17 | 5.60 ± 2.93 |
| School | 17 | IT Service | **3.40** ± 3.13 | 6.60 ± 3.39 |
| | | Bathroom | **3.60** ± 1.34 | 5.00 ± 5.10 |
| Airport | 25 | Immigration | **1.80** ± 0.45 | 7.20 ± 6.85 |
| | | Bathroom | **1.60** ± 0.55 | 6.20 ± 5.23 |
| | | Info Desk | **1.60** ± 1.34 | 8.20 ± 3.31 |
| Mall | 18 | Bathroom | **5.80** ± 0.83 | 7.40 ± 3.38 |



- LLMs could suggest the correct location for unseen places within approximately **3 interaction rounds**.
- Compared to the random guess, our framework reduces interaction rounds by up to **80%** and demonstrate much more **stable** performance.
- However, commonsense fails in non-typical layouts: E.g., a bathroom is next to a health store in a mall.

## Room Tidying

**Success Rate, Execution Rate and Goal Condition Rate for Room Tidying**

| Room | Method | Number of Misplaced Objects | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | 2 | | | 4 | | | 12 | | |
| | | SRC | ER | GCR | SRC | ER | GCR | SRC | ER | GCR |
| Living Room | Our Method | **1.00** | **1.00** | **1.00** | **0.80** | 0.76 | **0.95** | **0.40** | 0.70 | **0.69** |
| | ProgPrompt | 0.60 | **1.00** | 0.70 | 0.40 | **0.92** | 0.70 | 0.00 | **0.79** | 0.15 |
| Kitchen | Our Method | **0.60** | **1.00** | **0.70** | **0.60** | 0.90 | **0.83** | **0.20** | 0.76 | **0.78** |
| | ProgPrompt | **0.60** | 0.96 | **0.70** | 0.20 | **0.97** | 0.65 | 0.00 | **0.94** | 0.17 |
| Bathroom | Our Method | **1.00** | **1.00** | **1.00** | **0.60** | **1.00** | **0.90** | **0.40** | **0.96** | **0.57** |
| | ProgPrompt | 0.40 | 0.89 | 0.50 | 0.20 | 0.93 | 0.45 | 0.00 | 0.81 | 0.20 |
| Bedroom | Our Method | **0.80** | 0.90 | **0.90** | **0.80** | 0.96 | **1.00** | **0.60** | 0.98 | **0.65** |
| | ProgPrompt | 0.40 | **0.91** | 0.60 | 0.20 | 0.82 | 0.35 | 0.00 | 0.94 | 0.22 |

**VirtualHome Room Tidying Results with Different Methods**



Original Messy Room

Room Tidied by ProgPrompt

Room Tidied by Our Method

- In all scenarios, **60%** of misplaced objects can be placed correctly, and up to **80%** in less messy rooms.
- Hierarchical planning is effective in enabling LLMs to reason about long-horizon action plans and avoid generate irrelevant actions.