

ProAct

A Benchmark and Multimodal Framework for Structure-Aware Proactive Response

Xiaomeng Zhu^{1,2} Fengming Zhu¹ Weijie Zhou² Ye Tian² Zhenlin Hu³ Yufei Huang²

Yuchun Guo² Xinyu Wu⁴ Zhengyou Zhang² Fangzhen Lin¹✉ Xuantang Xiong²✉

¹ Dept. of Computer Science and Engineering, The Hong Kong University of Science and Technology (HKUST), Hong Kong SAR, China

² Tencent, Shenzhen, China

³ Futian Laboratory, Shenzhen, China

⁴ Shenzhen Institute of Advanced Technology (SIAT), Chinese Academy of Sciences, Shenzhen, China

From passive response to proactive response

The robot should pursue the task objective, not merely mirror the human.

Previous

Current

Future

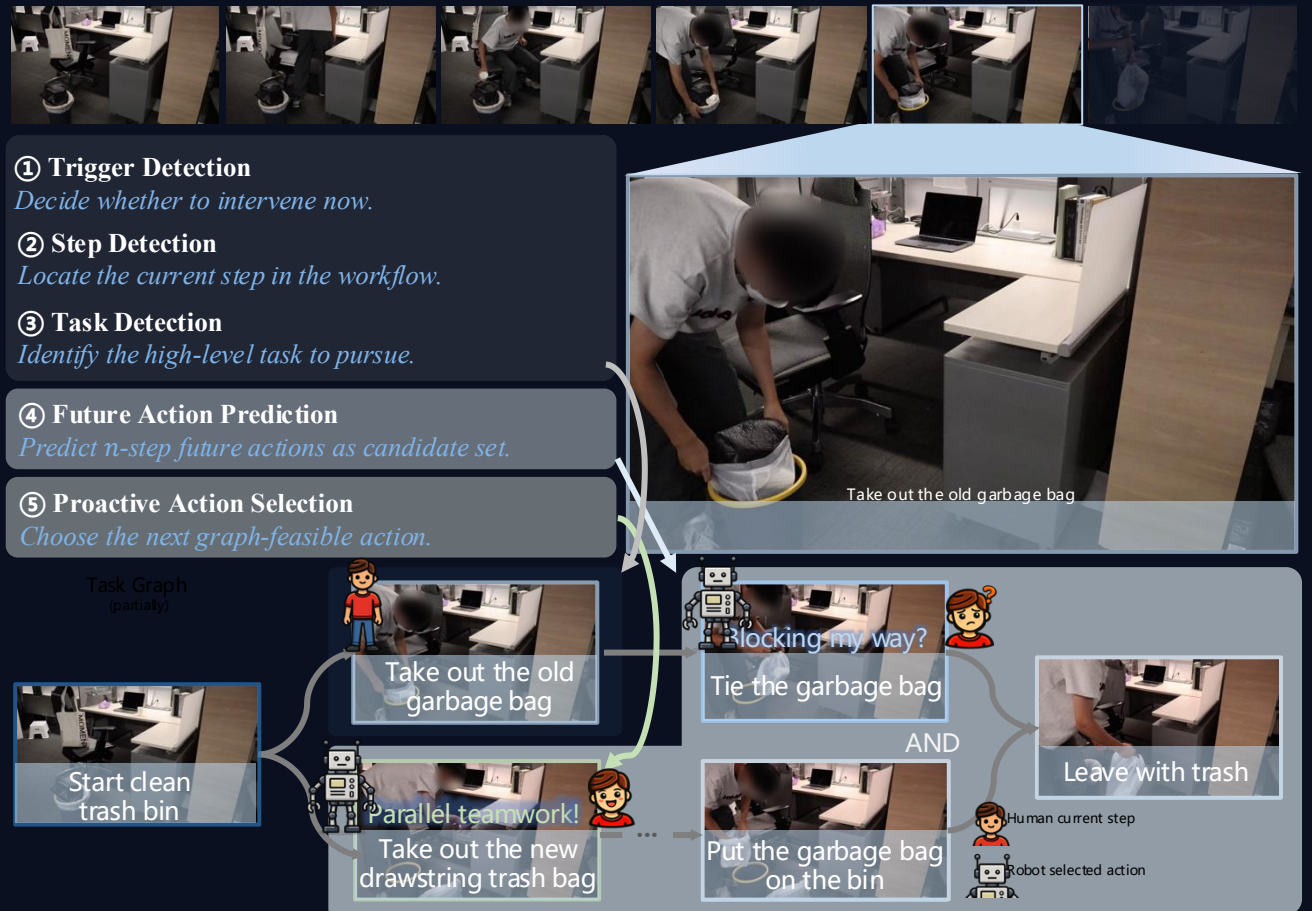
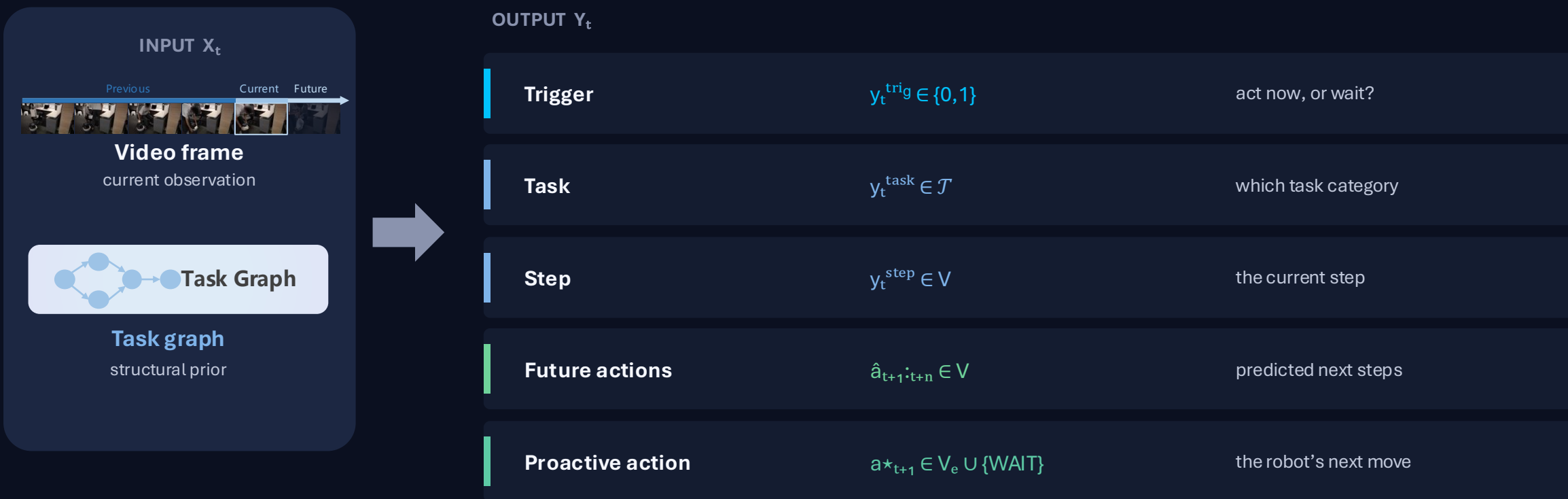


Figure 1. Overview of proactive response — ProAct-75 supports five vision-based tasks with step-level annotations, hierarchical labels, and task graphs.

Proactive response, formally

At each timestep, the agent reads the scene and the task graph, then outputs five decisions.



These predictions form a hierarchy — trigger gates task, task scopes step, step grounds the next action.

Why existing benchmarks fall short

Recognition is well studied; structure-aware intervention is not.

What current video benchmarks measure

Action recognition

What action is happening right now

Temporal localization

When does each action start and end

Anticipation

What is the single most likely next action

What proactive response also needs

Should I act now?

Trigger detection — judging the right moment to intervene

Which action stays valid?

Respecting task structure and step dependencies

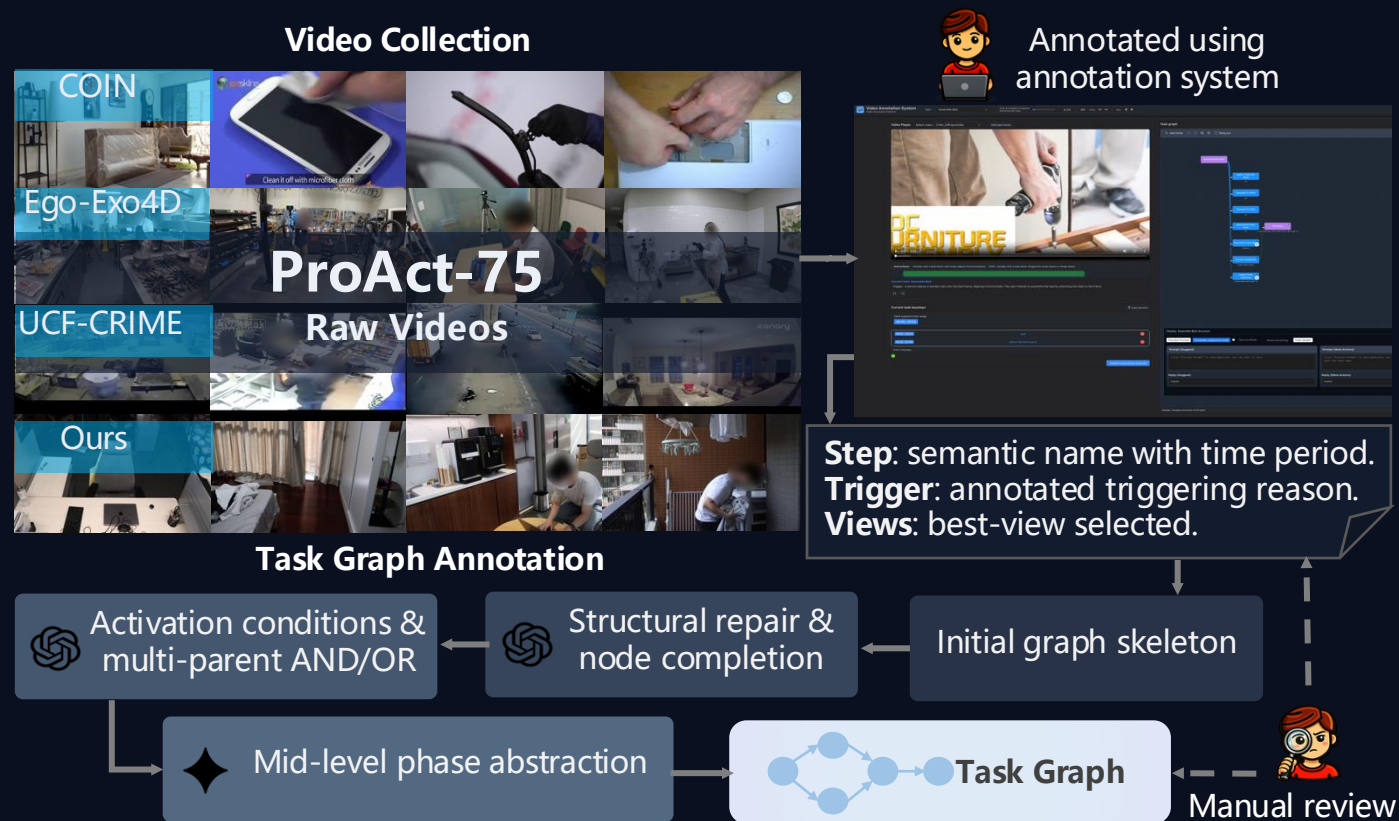
Can I run in parallel?

Acting on an independent thread, not mirroring the human

No existing benchmark couples perception, anticipation and structured decision-making — ProAct-75 does.

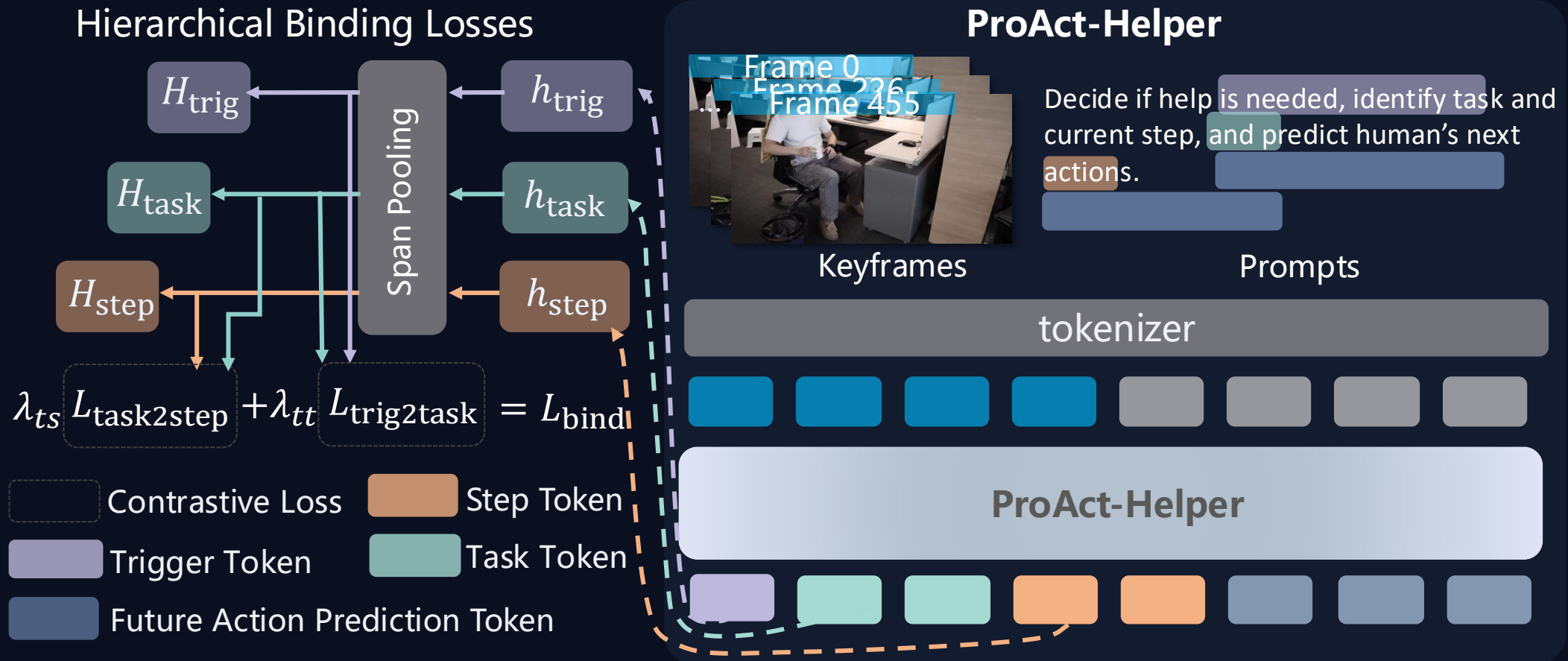
From raw videos to task graphs

A three-stage annotation pipeline turns demonstrations into executable structure.



The ProAct-Helper framework

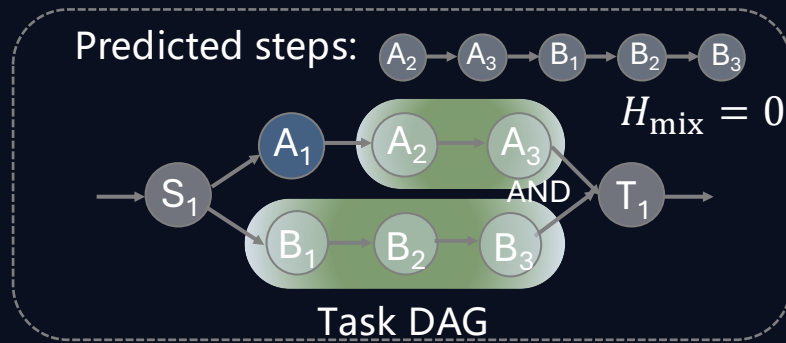
MLLM perception grounds the scene; the task graph constrains planning.



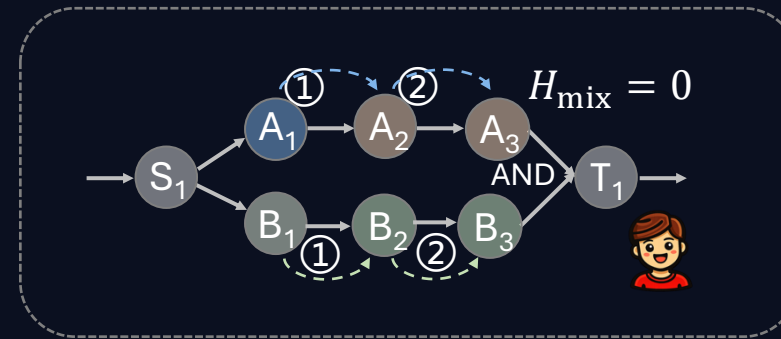
Task graph formulation

A task is a DAG of executable actions and structural (AND / OR) nodes.

Proactive Action Selection



Entropy-driven heuristic search



● Human Step ● Robot Step ● Candidates

ENTROPY OBJECTIVE

$$a^*_{t+1} = \arg \min_a H_{\text{mix}}(H_t, R_t, a)$$

Among graph-feasible actions, pick the one that least mixes the human and robot threads.

Task graph = DAG

Each task is a Directed Acyclic Graph. Executable nodes (V_e) are concrete actions; structural nodes (V_n) carry AND / OR logic.

AND / OR semantics

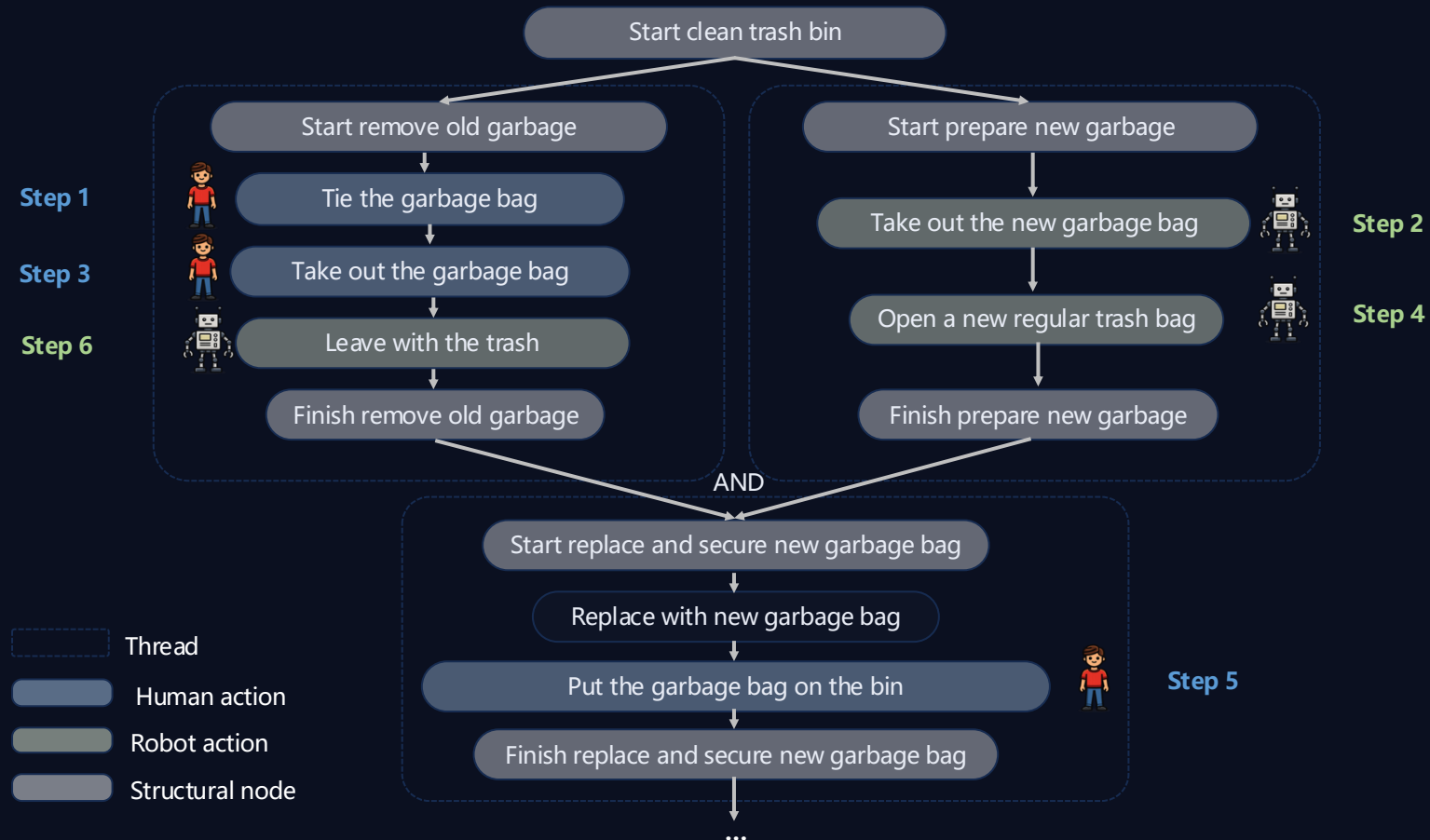
AND node fires only after all predecessors finish; OR node fires once any one of them does.

Parallel threads

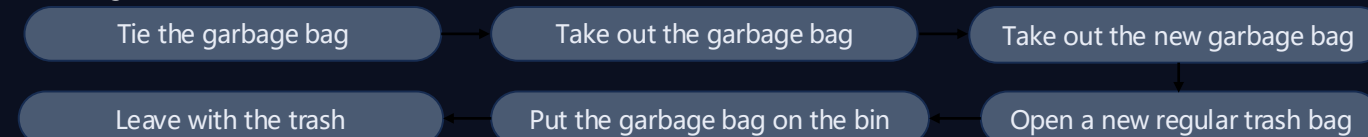
Branches independent until they merge — letting the robot act on its own thread alongside the human.

Task graph: a worked example

A task is a DAG of executable actions and structural (AND / OR) nodes.



Human ground-truth labeled trace:



Main results

ProAct-Helper improves both perception and proactive collaboration.

+6.21%

Trigger detection mF1

▲ over strong closed-source MLLMs

0.25

fewer steps per online decision

▲ more efficient one-step planning

+15.58%

parallel action rate

▲ genuine human–robot parallelism

Better on both axes that matter — perceiving the right moment, and collaborating efficiently.

ProAct: A Benchmark and Multimodal Framework for Structure-Aware Proactive Response

Thank you for watching!

Xiaomeng Zhu^{1,2} Fengming Zhu¹ Weijie Zhou² Ye Tian² Zhenlin Hu³ Yufei Huang²
Yuchun Guo² Xinyu Wu⁴ Zhengyou Zhang² Fangzhen Lin¹✉ Xuantang Xiong²✉

¹ Dept. of Computer Science and Engineering, The Hong Kong University of Science and Technology (HKUST), Hong Kong SAR, China

² Tencent, Shenzhen, China

³ Futian Laboratory, Shenzhen, China

⁴ Shenzhen Institute of Advanced Technology (SIAT), Chinese Academy of Sciences, Shenzhen, China