Towards Efficient Assembly Motion Planning via Vision-Language Models

Team 5 Jung Geumyoung, Kim Jaemin



Motivation : Diffusion model-based approach do not guarantee safety Key Idea :

- Diffuser + Control Barrier Functions = SafeDiffuser
- embed finite-time diffusion invariance into denoising procedure using a class of control barrier functions to ensure safety





Contents

- Introduction
- Challenges
- Method
- Experiments
- Conclusion



Introduction

Assembly part motion planning :

process of planning collision-free path of a part to move and join individual components into final assembled product





Introduction | Assembly by Disassembly

Assembly of n parts \rightarrow **O(n!)** search space

Assume all parts are rigid bodies, then

- Assembly = reversed disassembly
- Disassembly of n parts $\rightarrow O(n^2)$ search space (efficient!)





ASAP's method for Part motion planning :

- Samples discretized actions (forces) in 6 directions
- Apply the actions in **physics-based simulation**, check disassembly success





+Z applied on assembly



Challenge | Running Time

Physics-based simulation comes with high computation costs

Main functions: check_assemblable (motion planning), get_stable_plan (stability check)





less action trials(accurate action prediction), Reduce computation time

Utilize **Spatial reasoning capable FM** to guess initial disassembly force direction

Reduce the number of physics-based action queries to speed up the motion planning

ChatGPT

ΑΡΙ

GPT-40







Method | System Overview



Given an **assembly G** and a **target part P** to remove from G:

- 1. Render G, with **P distinctly colored** (e.g. magenta) in **3 principal views**
- 2. Pass the 3 images + text prompt to the VLM
- 3. VLM will rank actions based on predicted (disassembly) success rates



Method | Action Prediction via VLM

Constructing coordinate axis can improve VLM's spatial reasoning capability



Can LLM be a Good Path Planner based on Prompt Engineering? Mitigating the Hallucination for Path Planning (2024.08. arXiv)

Top view



Method | Action Prediction via VLM

Text instruction for spatial reasoning





Experiments

Hypothesis : Successful disassembly within less trials, Reduced computation time

1. Verify VLM's 3D spatial reasoning Actions sampled, planning time

2. VLM capability on narrow-passage problem

Actions sampled & time taken

under different passage narrowness



Dataset : 46 assemblies from ASAP (~250 part motion plans)



Experiments

Hypothesis : Successful disassembly within less trials, Reduced computation time

1. Verify VLM's 3D spatial reasoning Actions sampled, planning time

2. VLM capability on narrow-passage problem

Actions sampled & time taken

under different passage narrowness



Experiment 1 | Verifying VLM's 3D reasoning capability

<u>Q: "Can VLM perform 3D spatial reasoning from image snapshots?"</u>

- Actions sampled: avg. number of actions sample per disassembly
- Planning time: avg. time taken per disassembly

Successful disassembly within fewer actions sampled (& took less time)



Baseline	Ours	Improvement
1.9±1.28	1.74 ±1.32	+9.90%



Experiment 1 | Verifying VLM's 3D reasoning capability

<u>Q: "Can VLM perform 3D spatial reasoning from image snapshots?"</u>

- Actions sampled: avg. number of actions samples per disassembly
- Planning time: avg. time taken per disassembly

Successful disassembly within fewer actions sampled (& took less time)



Baseline	Ours	Improvement
1.14±1.63	0.970±1.45	+17.50%



Experiments

Hypothesis : Successful disassembly within less trials, Reduced computation time

1. Verify VLM's 3D spatial reasoning Actions sampled, planning time

2. VLM capability on narrow-passage problem

Actions sampled & time taken

under different passage narrowness



Spatial reasoning shall shine, especially under **disassemblies with tighter constraints** (i.e. narrower passages)

Given access to VLM with 'good' spatial reasoning capabilities:

- Random sampling will result in even lower success rate
- Performance improvement will be more significant



Only +Z force is allowed

Disassembly motion with a narrow passage.



Experiment 2 | Capability on narrow-passage problem

<u>Q: "Does VLM perform better / worse on tasks with narrow passages?"</u>

- Measured actions sampled, planning time
- Narrowness of a passage: how many out of 6 actions leads to success?
 - \circ Lower values \rightarrow more difficult to disassemble
- Number of actions tried reduced across different difficulties



Narrowness	Baseline	Ours	Improvement
1	3.63±1.72	3.29±1.94	+10.4%
2	2.28±1.15	2.05±1.48	+11.0%
3	1.84±0.90	1.58±0.77	+16.7%
4	1.40±0.62	1.30±0.60	+7.70%
5	1.28±0.45	1.19±0.40	+7.60%



Experiment 2 | Capability on narrow-passage problem

<u>Q: "Does VLM perform better / worse on tasks with narrow passages?"</u>

- Measured actions sampled, planning time
- Narrowness of a passage: how many out of 6 actions leads to success?
 - \circ Lower values \rightarrow more difficult to disassemble
- Planning time had much higher variance (simulation cost is inconsistent)



KAIST

Narrowness	Baseline	Ours	Improvement
1	1.79±1.98	1.25±1.15	+42.8%
2	1.44±1.82	1.32±2.23	+8.95%
3	1.71±2.06	1.45±1.62	+18.0%
4	1.00±1.48	1.05±1.58	-5.34%
5	0.70±1.20	0.60±0.92	+16.7%

Avg. time per simulation = 0.71**±1.82** s

Experiment 2 | Outlier Cases

- Some simulations took unreasonably long computation time
 - Caused high variance in planning time, independent of our VLM method



A simulation taking ~60s to determine failure



Conclusion | Limitations & Future Works

Prediction Accuracy

- 3 principal views may be insufficient (prone to occlusion)
- Sophisticated prompt engineering may improve prediction accuracy

VLM Inference Overhead

- Used GPT-40 API for testing (~2-3s overhead per query)
- Try faster models
- Try open-sourced model to deploy in-house

Further Improvements in Running Time

• Try VLM on other assembly problems, e.g. part selection (Q. which part should be disassembled first?)









Conclusion | Key Takeaways

VLM can perform 3D spatial reasoning on part motion planning

• Achieves higher action success rate compared to random sampling method

Reduced number of call is **related** to less overall computation time

• Yet, variance within a simulation is also a significant factor



Thank you!



Schedule & Roles

O : main

V : support

		Geumyoung Jung	Jaemin Kim
Have	Analyze previous work	0	0
Done	Test Reference code	0	Ο
	VLM background research & prompt engineering	Ο	V
	Integrating VLM into pipeline	V	0
	Testing & collecting results	0	Ο
	Prepare final presentation	0	Ο

