# Autonomous Minigolf Mastery
## Enabling Real-World Robot Arms
## through a Multi-Modality Vision-Language Model Framework

**Hantao Zhou**

5/7/24

# Task Definition

■ **Mastering Minigolf in Various Course with a Robot Arm**

■ Learn to how to successfully hit the goal with the ball <span style="color:red">without human intervention</span>

■ Adaptable to complex courts with <span style="color:red">different obstacles</span>
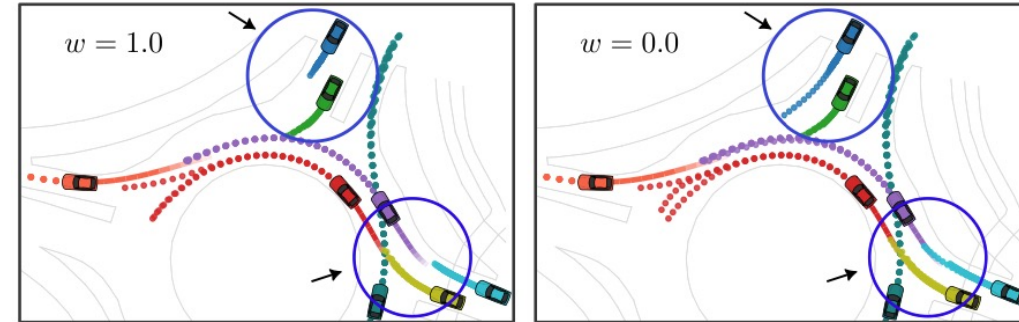


🎯 Enable robot to integrate visual data with planning and adapt based on outcomes.

💡 Propose novel framework that can handle complex, multimodal data and make decisions accordingly.

# General Problem and Significance



- Complex Trajectories Prediction and Kinodynamic understanding

- Trajectory Prediction

  - CNN, Diffusion

  - A type of behaviours forecasting, highly valuable for many scenarios like self-driving, collision detection

- Kinodynamic Understanding

  - A capability that is crucial for embodied AI

  - VLM currently performs very poor

- These two feature have not been synergistically combined, lacking the possibilities of improving the performances recurrently

# Previous VLMs' deficiencies
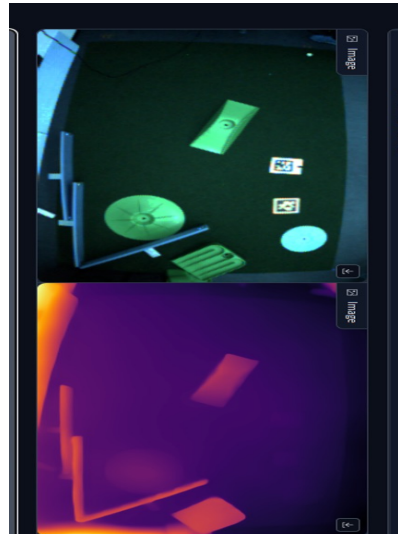
- **Core VLM's deficiencies**
  - **Only harness simple logic of <span style="color:red">positional relationships</span> rather than the complex nature of trajectories**
  - **Can't understand trajectories effectively (<span style="color:red">no VLM can read video</span> as a continuous file, they just extract key frames)**
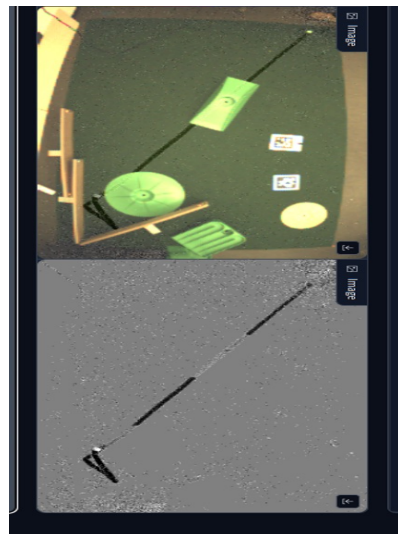- **Applying VLM in robotics**
  - **Only serving as high-level planner**
  - **Only <span style="color:red">Deterministic Execution</span>, no feedback and rethinking**
  - **Vision information only as assistance for understanding, with out vision -> same performance (especially "look before you leap")**
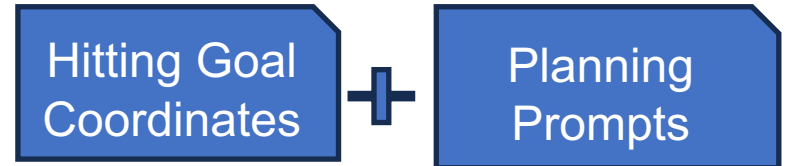
# Proposed Framework

# Challenges

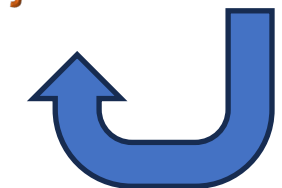- **Current VLM is not precise on <span style="color:red">details in vision</span> – Especially in the Kinodynamic and Trajectory**

- **Sim2Real can simulate the environment and generate an agent that can perform reasonably well in real scenarios --- but takes high <span style="color:red">cost especially there is a huge gap between simulation and reality</span>**

- **Hard coded control can solve certain scenarios hard-coded way but <span style="color:red">not adaptable</span> for complex scenarios**

# Two Main Sub-Tasks

- **Perception**

  - Percept and understand obstacles in environments

  - Track high speed motion at high temporal resolution

- **Planning**
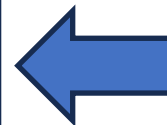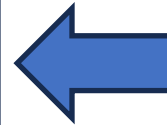
  - Evaluate feasibility of hitting to goal under certain obstacles setup

  - If feasible

    - Plan trajectory based on distances and positions of obstacles

    - Adjust strategy based on hitting recordings

# Advanced Tools

**Sub-tasks**

**Advanced Tools**

**Perception**

- Percept and understand obstacles in environments
- Track high speed motion at high temporal resolution

← 

- **VLM – Grounded SAM**
- **Event Camera**

**Planning**

- Evaluate feasibility of successful hitting
- Adjust strategy based on hitting recordings

← 

- **Depth Camera**
- **VLM – Llama Adapter**

# Related Works

- **VLM+Robotics**
  - VoxPoser
  - MOKA Manipulation

- **Limitaions**
  - **Static Tasks**
  - **No adjustment feedback**



(a) 3D Value Map Composition

(b) Motion Planning

# Current Experimental Results

- **RGBD**

  - Calibration

  - Visual Grounding

# Current Experimental Results

■ **Dynamic Perception**

   ■ Calibration



■ Find trajectory



black trajectory 0.49

# Current Experimental Results

■ **Plan & Evaluate**

    ■ Analyse the court based on prompts

    ■ Generate planning trajectory based on image

    ■ Evaluate recorded trajectory

# Current Progress in sum 0422

- **Perception**
  - Event Cam + RGB Cam to record
  - Blob + Yolo to generate high precision trajectories

- **Planning**
  - Baseline: GPT-4V + Prompt to generate hitting parameters and evaluate hitting results
  - Baseline: GPT-4V + Parameters + Prompt to generate hitting parameters and evaluate results