# Advanced Methods in Robot Learning Final Report Solving the Tower of Hanoi *

Göttlich Joshua Jonah
*(12342542)*

Ugur Kaan
*(52321844)*

Wolf Leon Marcel
*(12434723)*

Postelnecu Ioana
*(12444870)*

*Abstract*—This paper presents the design and implementation of an autonomous robotic arm capable of solving the classic Tower of Hanoi puzzle. The system integrates several key domains of robotics and artificial intelligence, including motion planning for trajectory generation, robust pick-and-place manipulation for physical interaction, and a high-level reasoning engine powered by a Large Language Model (LLMs).

## I. INTRODUCTION

The Tower of Hanoi is a well-known logic puzzle that provides an ideal testbed for evaluating reasoning and manipulation in robotic systems. Successfully solving this puzzle requires the ability to understand abstract rules, evaluate dynamic states, and perform precise pick-and-place operations. In this project, we developed an autonomous robotic system that integrates perception, planning, and reasoning to solve the Tower of Hanoi task in both simulation and a real-world setup.

The system combines several core components of modern robotics and artificial intelligence. For high-level reasoning, we employ a Large Language Model (LLM) that interprets the current game state and generates optimal next-move decisions. A custom-designed prompt enables the LLM to validate states, assess legal moves, and output structured commands. Vision is handled by a detection pipeline that identifies disk positions using camera inputs, enabling the system to construct a symbolic representation of the Tower of Hanoi state. For motion execution, we first explored Dynamic Movement Primitives (DMPs) to learn and reproduce manipulation trajectories, but due to limitations in generalization and robustness, we transitioned to a more deterministic motion planning approach using inverse kinematics and predefined waypoints.

## II. METHODOLOGY

Our approach to solving the Tower of Hanoi puzzle is based on a modular, hierarchical system architecture that follows the classical perceive-plan-act paradigm. The system is decomposed into three primary, interconnected modules: a perception module for state estimation, a high-level reasoning module powered by a Large Language Model (LLM) for strategic decision-making, and a motion execution module for physical interaction with the environment. This separation of concerns allows each component to be developed and refined independently, creating a robust and flexible framework.

Figure 1 illustrates the overall system architecture and the flow of information between the core components. The control loop begins with the perception module observing the workspace, which is then translated into a symbolic state for the LLM. The LLM processes this state to determine the next optimal move, which is then passed to the motion planning module for execution. Each of these components is detailed in the following subsections.
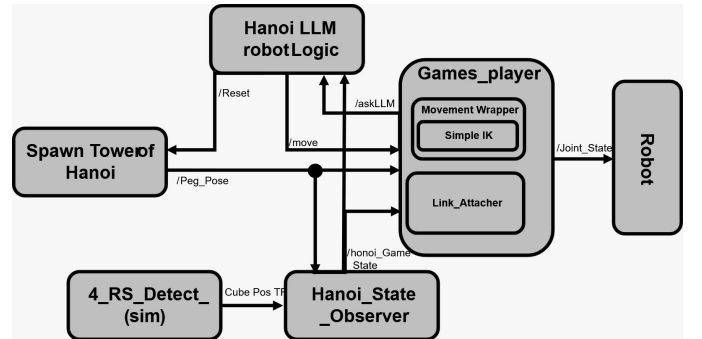


Fig. 1: System architecture illustrating the flow of information. The control loop follows a perceive-plan-act model, with the LLM serving as the central reasoning engine.

### A. Perception and State Representation

The foundation of the robot's autonomous operation is its ability to perceive the current state of the puzzle. This is achieved through a vision-based object detection pipeline, as detailed in Section **??**. Using a Realsense camera mounted with a view of the workspace, the system captures images and processes them to identify the location and identity of each disk. The output is a structured, symbolic representation of the game state (e.g., Peg A: [Disk3, Disk2], Peg B: [Disk1], Peg C: []), which serves as the ground truth for the planning module.

### B. LLM-based Strategic Planning

The core of our system's intelligence resides in a Large Language Model (LLM), which acts as the strategic planner or the "brain" of the operation. As described in VI, we utilize the Gemini Pro model via the Google AI API. This module receives the structured game state from the perception component and, guided by a carefully engineered system prompt, performs three key tasks:

1) **State Validation:** It first confirms that the current state adheres to the fundamental rules of the puzzle.
2) **Optimal Move Generation:** It then reasons about all possible legal moves and determines the single most

optimal move according to the standard Tower of Hanoi algorithm.

3) **Structured Output:** It returns the decision in a machine-readable JSON format (e.g `{"source_peg": "A",}` `"destination_peg": "C"}`), which decouples the abstract logic from the physical execution.

### C. Motion Planning and Execution

The final stage of the pipeline is translating the LLM's symbolic command into physical action. As detailed in Section IV, this is handled by the motion planning module, which is responsible for generating and executing collision-free trajectories for the robotic arm. While our initial explorations with Dynamic Movement Primitives (DMPs) offered insights into learning from demonstration, we ultimately adopted a more deterministic approach based on Inverse Kinematics (IK) for greater reliability and precision. This module, built on the MoveIt! framework, takes the target coordinates derived from the LLM's command and computes the precise sequence of joint movements required to perform the complex pick-and-place maneuvers.

### III. SIMULATION

To develop and validate our system in a controlled and repeatable manner prior to real-world deployment, we established a comprehensive simulation environment using Gazebo Classic. The simulation was designed to closely mirror the physical setup, ensuring that algorithms and control strategies would transfer effectively.

### A. Environment and Robot Setup

The core of the simulation is the virtual representation of our robotic arm. We utilized the official Unified Robot Description Format (URDF) model corresponding to the physical robot, which defines its kinematic structure, visual properties, and collision geometries. This model was integrated into Gazebo using the `ros2_control` framework to manage the robot's joints.

A critical aspect of the setup involved configuring the controllers for seamless integration with the MoveIt 2 motion planning framework. Standard controllers, including the `JointStateBroadcaster` and a `JointTrajectoryController`, were implemented to relay state information and execute planned trajectories. Adjustments to the controller configurations were necessary to ensure that the hardware interface in simulation correctly processed the commands generated by the IK solver, enabling accurate and reliable motion execution.

### B. Scenario Generation and State Initialization

To create the Tower of Hanoi puzzle within the simulation, we developed a dedicated Python script. This script serves as a scenario spawner, responsible for dynamically generating the puzzle at the start of each simulation run. Its primary functions are:

1) **Puzzle Spawning:** It programmatically spawns the models for the three pegs and the set of disks into the Gazebo world.
2) **Layout Configuration:** The positions of the pegs are configurable, allowing for different workspace layouts. Once spawned, the spawner publishes the static positions of the pegs so that other system components, such as the motion planner, can be aware of their locations.
3) **State Validity:** Crucially, the script only generates valid Tower of Hanoi configurations. It ensures that no larger disk is ever placed on top of a smaller disk, providing a valid and solvable starting state for every experiment.

This automated approach to scenario generation was essential for systematic testing and evaluation of the complete perception, planning, and control pipeline..

### IV. MOTION PLANNING

Motion planning serves as the crucial interface between the high-level task planner, in this case a Large Language Model (LLM), and the physical robot or its simulated counterpart. This component receives abstract commands, representing desired actions or end-goals, and subsequently computes smooth, collision-free trajectories for the robot's actuators. The execution of these trajectories enables the robot to successfully complete its assigned task.

### A. Motion Generation via Dynamic Movement Primitives

Initially, one of the investigated approaches for trajectory generation utilized Dynamic Movement Primitives (DMPs). DMPs offer a robust framework for learning and reproducing complex motor skills from demonstrations. For our implementation, a human operator first demonstrated the process of grasping a disk from a Tower of Hanoi puzzle. This demonstrated trajectory was then used to train a DMP, with the objective of enabling the system to generalize this grasping motion to new start and end coordinates. A potential benefit of this imitative method is the generation of more naturalistic, human-like motion profiles, as the learned trajectory intrinsically captures the characteristics of the original demonstration.

Despite their theoretical advantages, several practical limitations were encountered during implementation, which ultimately led to the discontinuation of this approach. The primary challenges were twofold: susceptibility to system noise and poor generalization.

First, the trajectories executed by the robot exhibited a lack of predictability. We identified that this issue stemmed from noise inherent in the robot arm's motors, which introduced stochastic variations into the final motion. To mitigate this effect, we attempted to filter quasi-static phases from the recorded demonstration data and applied smoothing filters to the resulting trajectory. While this strategy provided partial improvement, it did not fully resolve the underlying issue of unpredictability.

Second, the generalization capabilities of the trained DMP were found to be insufficient for the required tasks. Although theoretically capable of targeting any point in the workspace,

the system produced aberrant or non-intuitive trajectories when the target location was significantly different from the demonstration context. For instance, a DMP trained on grasping a disk from the rightmost peg struggled to generate a viable path for grasping a disk on the left. It is plausible that this limitation could be overcome by recording multiple demonstrations that cover the entire workspace (e.g., one for each peg) and employing a more sophisticated selection or blending mechanism. However, this solution was deemed impractical for our goal of developing a more universally applicable motion planning pipeline.

Given these challenges, we concluded that our DMP implementation could not reliably provide the robust, workspace-wide performance required. Therefore, we decided to pursue an alternative motion planning strategy.

### B. Motion Planning via Inverse Kinematics and Waypoints

To achieve more reliable and deterministic control, an alternative strategy was implemented utilizing Inverse Kinematics (IK) coupled with a sequence of Cartesian waypoints. This approach leverages the robot's kinematic model to compute the required joint configurations for a series of desired end-effector poses, thereby generating a precise and collision-aware trajectory.

*1) System Architecture and Implementation:* The motion planning framework is built upon the MoveIt! software stack. We load the robot's URDF (Unified Robot Description Format) model into the MoveIt! planning scene, which enables self-collision checking. Two distinct planning groups were defined: one for the robotic arm and one for the gripper, each with its corresponding kinematic chain specified from the base to the terminal link.

The core of the motion generation relies on the standard IK solver provided by MoveIt!. This solver receives a target Cartesian pose (position and orientation) for the end-effector and, if a solution exists, returns the corresponding robot configuration in joint space. These joint-space coordinates are then published to the `joint_state_controller`, which interpolates the joint angles over time to execute the physical movement. Gripper control is managed separately through a simplified interface offering two discrete commands: a close command, which actuates the gripper until an object is detected or it fully closes, and an open command, which retracts the gripper to its maximum extent.

*2) The Perceive-Plan-Act Cycle:* The high-level task execution follows a standard perceive-plan-act cycle, orchestrated by a central agent.

1) **Perceive:** The cycle begins with the robot moving to a predefined "Home" position. From this vantage point, the Realsense camera has an unobstructed view of the workspace. An image is captured, and a perception module processes it to detect the current positions of all disks, thereby establishing the game state. This state is continuously updated.
2) **Plan:** The system awaits a high-level command from the LLM, which specifies the symbolic goal (e.g., "move the top disk from the first peg to the third peg"). A task-planning agent translates this symbolic goal into concrete Cartesian coordinates. It identifies the coordinates of the disk to be moved and the destination. If the target peg is empty, the destination coordinate is the base of the peg, provided by a static `SpawnHanoiTower` module. If the peg is occupied, the destination is calculated as a position atop the highest disk on that peg.
3) **Act:** With the source and destination coordinates defined, the system proceeds to execute a pick-and-place maneuver.

*3) Execution of Pick-and-Place Maneuvers:* The physical interaction with the disks is broken down into a multi-stage process to ensure robustness.

*a) Pick Sequence::*

1) **Approach:** A pre-grasp pose is calculated directly above the target disk, with the gripper oriented vertically downwards and its jaws aligned parallel to the y-plane of the workspace. The robot moves to this pre-grasp pose.
2) **Grasp:** The robot executes a straight vertical descent until the end-effector is at the correct height to grasp the disk. The gripper is then closed.
3) **Secure (Simulation):** To circumvent physics instabilities within the simulation environment, the grasped object is programmatically attached to the gripper's end-effector link upon a successful grasp.
4) **Retreat:** The arm performs a direct vertical ascent to lift the disk clear of other objects, minimizing the risk of accidental collisions. It then returns to the Home position to ensure a clear path for the subsequent place motion.

*b) Place Sequence::* The place sequence mirrors the pick sequence. The robot moves from the Home position to a pre-place pose above the target destination. It then descends vertically, opens the gripper to release the disk, and detaches the link in simulation using a `LinkDetacher` module. Finally, the arm retreats by ascending vertically before returning to the Home position, ready to begin the next perceive-plan-act cycle.

*4) Addressing Kinematic Constraints and Singularities:* During development, we observed that the robot's sub-optimal kinematics made purely top-down grasps challenging. The IK solver frequently failed to find a solution when the end-effector was constrained to a perfectly vertical orientation (i.e., roll, pitch, and yaw angles are zero). To address this, we introduced a constraint relaxation strategy. If the IK solver fails to find a solution for the strictly defined pose, the planner enters a fallback routine that permits rotation around the gripper's local x-axis (roll). This relaxation allows for a successful grasp, as the gripper jaws remain parallel to the disk's flat surfaces. Alternative rotations (around the y- or z-axis) were disallowed as they would misalign the gripper with the object and likely lead to grasp failure.

## V. Object Detection Algorithm

### A. Summary of Previous Work

We previously implemented an object detection pipeline using YOLOv8 and SAM2 as part of Assignment 2. The goal was to detect colored cubes placed in the workspace of the OpenManipulator 6DOF robot and use the detections for robotic manipulation. To that end, we created a custom dataset using smartphone-recorded videos and SAM2-assisted annotations, trained a YOLOv8n model, and integrated the results into our ROS-based control pipeline.

The dataset covered various cube positions and lighting conditions but initially did not include stacked cubes. Detection generally worked well, but we faced some problems: the model struggled with inference speed on CPU and often confused the blue and black cubes due to similar appearance. These limitations are described in more detail in the old Report.

### B. Improvements in Dataset and Model

To address the issue of poor detection when cubes were stacked, we recorded new videos with stacked configurations and retrained the model. This significantly improved detection in such cases, as shown in Fig. 5.



Fig. 2: Improved detection of stacked cubes after retraining.

To address the performance issues mentioned in our earlier work, we switched to a shared desktop machine equipped with a 16 GB GPU. Team members accessed the system remotely via AnyDesk, which allowed us to run training and inference much more efficiently. Compared to our earlier CPU-only setup, we now achieve significantly better frame rates and training times.

### C. Handling Color Confusion (Blue vs. Black)

Even after augmentation and retraining, the confusion between the blue and black cube remained (black cube in Fig.

5) Since these colors are visually similar, the model often predicted both as "blue". We experimented with different augmentations and training parameters, but without success.

As a workaround, we implemented a heuristic tailored to our Tower of Hanoi scenario, where each cube color must be unique. When a color appeared more than once, we checked the two highest class confidences for each box. If both had the same two top classes, the cube with the higher second-best confidence was assigned the second class. This approach improved overall stability under the assumption of unique colors.
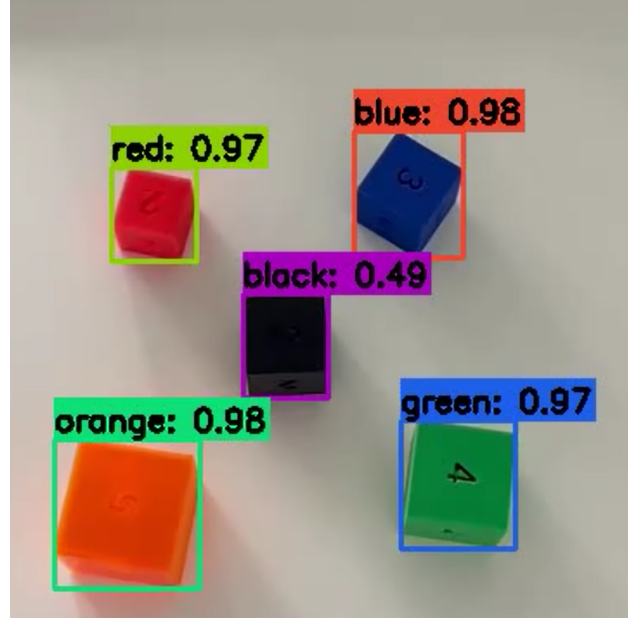


Fig. 3: Improved detection of stacked cubes after retraining.

### D. Integration in Simulation

In the simulation environment, the object positions were not detected via YOLO, but instead represented by predefined `tf` frames. These simulated object frames were used to test and visualize the pick-and-place motion without requiring real-time perception. The robot controller subscribed to these `tf` frames to plan and execute movements, allowing us to validate the manipulation part of the pipeline independently from the detection.

## VI. Large Language Model (LLM)

As established in Section II, the cognitive core of our system is a Large Language Model (LLM) responsible for strategic planning. The LLM abstracts the complex puzzle logic away from the motion control system, decoupling the high-level task of *what* to do from the low-level execution of *how* to do it. This section provides a detailed account of the model's implementation, the prompt engineering required to elicit correct behavior, and its integration into the communication protocol of our system.

## A. Model Selection and Prompt Design

For this task, we utilized the `gemini-pro` model from Google AI, accessed via its dedicated API. The choice was motivated by the model's advanced reasoning capabilities and its proficiency in processing structured data formats, which was deemed superior to locally-hosted alternatives like Ollama for this application.

A crucial component for ensuring reliable performance is prompt engineering. We designed a system prompt that explicitly defines the LLM's role, the rules of the game, a required reasoning process, and a strict output format. This structured prompt, shown in Listing 1, constrains the model's behavior and forces it to act as a deterministic component within our architecture. The prompt mandates a JSON output, which simplifies parsing and integration while also providing a mechanism for error reporting (e.g., in the case of an illegal game state).

```
1  You are a master strategist for a robotic arm that
       solves the Tower of Hanoi puzzle. Your task is
       to determine the single next optimal move.
2
3  **Rules of the Game:**
4  1.  You can only move one disk at a time.
5  2.  A larger disk can never be placed on top of a
       smaller disk.
6  3.  You can only move the topmost disk from one of
       the three pegs (A, B, C) to another.
7
8  **State Validation:**
9  Before determining a move, you MUST first validate
       that the provided 'Current State' is legal. If
       any larger disk is on top of a smaller disk, the
        state is invalid.
10
11 **Reasoning Process:**
12 1. First, validate the input 'Current State' as per
        the 'State Validation' rule.
13 2. If valid, identify the top disk on each peg.
14 3. Internally list all possible moves (e.g., A->B, A
       ->C, B->A, etc.).
15 4. For each possible move, check if it is legal
       according to the rule 'a larger disk cannot be
       placed on a smaller disk'.
16 5. From the list of legal moves, select the single
       most optimal move that follows the standard
       algorithm to solve the puzzle in the minimum
       number of steps.
17 6. Finally, provide ONLY the selected optimal move
       in the required JSON format.
18
19 **Output Format:**
20 - If the state is **valid**, your response MUST be a
        JSON object with the source and destination peg
       . Example: `{"source_peg": "A", "destination_peg
       ": "C"}`
21 - If the state is **invalid**, your response MUST be
        a JSON object with a single "error" key.
       Example: `{"error": "Invalid state on Peg A:
       disk 2 cannot be on disk 1."}`
22
23 Do not include any other text or explanations.
```

Listing 1: System prompt for the LLM-based Tower of Hanoi strategist.

## B. Data Flow and Communication Protocol

The LLM's integration into the control loop is managed by a well-defined request-response protocol facilitated by ROS 2 topics and services.

1) **State Publication:** The perception module (`Hanoi_State_Observer`) publishes the current, structured game state to the `/hanoi_game_state` topic.

2) **LLM Invocation:** When ready for a new command, the primary task planner (`Games_player`) invokes the LLM. The `Hanoi_LLM_robotLogic` node subscribes to the game state topic, combines the latest state with the system prompt (Listing 1), and sends the complete query to the Google AI API.

3) **Command Parsing and Forwarding:** The LLM's JSON response is received by the `Hanoi_LLM_robotLogic` node, which parses it for validity. The extracted move command (e.g., source and destination pegs) is then published to the `/move` topic.

4) **Action Execution:** The `Games_player` receives this high-level command and utilizes its internal motion planning components (the `Movement Wrapper` and IK solver) to translate the abstract move into an executable trajectory, which is published to the robot's controllers.

This event-driven communication ensures a seamless and robust integration, allowing the LLM to effectively serve as the symbolic reasoning engine of the robotic system.



Fig. 4: LLM logic node



Fig. 5: LLM logic node in impossible case

## VII. EXPERIMENTAL RESULTS

The system is operational, with the large language model (LLM) exhibiting occasional, which is a normal characteristic. Demonstrative videos are available on our GitHub repository.

Furthermore, the grabbing functionality has a 99% success rate.

## VIII. Conclusion

We have presented an autonomous robotic system capable of solving the Tower of Hanoi puzzle by integrating vision, reasoning, and motion control components. The robot perceives the environment through a vision module that detects the configuration of disks on the pegs and translates this into a structured state. A Large Language Model serves as the brain of the system, analyzing the game state and generating the optimal next move according to the puzzle's rules.

For motion execution, initial experiments with Dynamic Movement Primitives revealed challenges related to noise and generalization. As a result, we adopted a more reliable motion planning pipeline based on inverse kinematics and waypoints, allowing for precise and repeatable movements across different puzzle configurations.

The modular architecture successfully decouples high-level decision-making from low-level control, allowing for flexibility and extensibility. This project demonstrates the potential of hybrid AI-robotic systems, combining the symbolic power of LLMs with physical task execution in complex environments.

This report outlines the architecture, implementation, and evaluation of this integrated approach, highlighting the effectiveness of combining LLM-based logical reasoning with low-level robotic control and perception.

## IX. Appendix

The code can be found at: https://github.com/joshuagoettlich/adv_robot_learning