

End-to-End Learning for a Low-Cost Robotics Arm

by

Abhishek Chothani

A thesis submitted to the College of Engineering and Science of  
Florida Institute of Technology  
in partial fulfillment of the requirements  
for the degree of

Master of Science  
in  
Mechanical Engineering

Melbourne, Florida  
December, 2024

We, the undersigned committee, hereby approve the attached thesis,  
“End-to-End Learning for a Low-Cost Robotics Arm”  
by  
Abhishek Chothani

---

Ryan T. White, Ph.D.  
Assistant Professor  
Mathematics and Systems Engineering  
Major Advisor

---

Venkat Keshav Chivukula, Ph.D.  
Assistant Professor  
Biomedical Engineering and Science

---

Kim-Doang Nguyen, Ph.D.  
Assistant Professor  
Mechanical and Civil Engineering

---

Seong Hyeon Hong, Ph.D.  
Assistant Professor  
Mechanical and Civil Engineering

---

Troy V. Nguyen, Ph.D, PE  
Professor and Department Head  
Mechanical and Civil Engineering

# Abstract

Title: End-to-End Learning for a Low-Cost Robotics Arm

Author: Abhishek Chothani

Advisor: Ryan T. White, Ph.D.

Robotic manipulation is a cornerstone of automation, with the ultimate goal of developing versatile systems capable of executing a wide range of real-world tasks autonomously. Traditional robotics approaches, while reliable and widely adopted in industrial settings, often struggle with adaptability, perception, and dynamic task execution. This thesis explores the evolution from classical robotics techniques to modern learning-based approaches, leveraging advancements in artificial intelligence to overcome these limitations.

Initially, the thesis presents a pick-and-place pipeline built using the Drake robotics framework and the KUKA iiwa robotic arm. This system employs a pseudoinverse controller for inverse kinematics to perform structured tasks like object manipulation. The limitations of this classical approach are highlighted, including sensitivity to object variations, lack of environmental awareness, and challenges with kinematic singularities. Despite these constraints, the significance of classical methods in industrial applications is acknowledged, emphasizing their reliability and cost-effectiveness.

To address these challenges, the thesis transitions to end-to-end learning-based methods for robotic manipulation. Leveraging the LeRobot codebase, two cutting-edge policies—Action Chunking with Transformers (ACT) and Diffusion Policies—are implemented and trained independently on a low-cost hardware setup. ACT focuses on learning temporally coherent action sequences, reducing compounding errors and ensuring smooth task

execution. In contrast, Diffusion Policies iteratively refine noisy trajectories, achieving high precision and robustness for intricate manipulation tasks.

The evaluation process involves a series of progressively complex tasks, including pick-and-place, sliding objects, stacking, orienting objects, and writing. Data is collected using a leader-follower teleoperation system, with dual fixed cameras providing synchronized visual inputs. Each policy is trained and evaluated separately, enabling a systematic comparison of their performance in terms of task success rate, motion smoothness, and adaptability.

The findings demonstrate the potential of learning-based approaches to handle complex manipulation tasks effectively, even on affordable hardware. ACT excels in maintaining temporal coherence, while Diffusion Policies achieve superior precision in fine-grained tasks. By systematically analyzing the strengths and limitations of these methods, the thesis contributes to bridging the gap between traditional robotics and modern, adaptive systems. This work lays the groundwork for scalable, accessible robotic solutions capable of addressing diverse real-world challenges.

# Table of Contents

Abstract .....	iii
List of Figures .....	viii
List of Tables.....	ix
Chapter 1 Introduction.....	1
1.1 Chapters Overview.....	2
1.2 Related Work .....	3
Classical Robotics Approaches.....	3
Perception in Robotics.....	3
Learning-Based Methods.....	3
Applications in Real-World Scenarios .....	4
Chapter 2 Development of a Pick-and-Place Pipeline Using Drake and the KUKA iiwa Robot .....	5
2.1 Overview of the Pick-and-Place Pipeline .....	5
Key Phases of the Pick-and-Place Task: .....	5
2.2 Importance of Classical Approaches in Industrial Applications.....	8
Key Advantages:.....	8
2.3 Limitations of the Classical Approach.....	8
Chapter 3 End-to-End Learning-Based Robotics Manipulation on Low-Cost Hardware .....	10
3.1 Hardware Setup.....	11
Robotic Manipulators .....	11
Vision System.....	12

Integration and Calibration .....	14
Cost and Accessibility .....	14
Summary.....	15
4.2 Software Architecture .....	16
Action Chunking with Transformers (ACT) .....	16
Diffusion Policies .....	17
Integration with the LeRobot Codebase .....	17
Summary.....	18
4.3 Integration and Workflow .....	19
1. Data Collection .....	19
2. Training .....	20
3. Execution .....	21
4.4 Experiments: Task Descriptions and Approach.....	22
1. Pick-and-Place .....	23
2. Stacking Objects .....	24
3. Orienting an Object.....	24
4. Sliding an Object .....	25
Progressive Complexity and Framework Design .....	26
Learning without Explicit Knowledge.....	26
Evaluation Criteria.....	27
4.5 Results and Discussion.....	28
Analysis of Results .....	29
Conclusion of Results .....	31

References .....	32
Appendix .....	35
Appendix 01: Detailed Explanation of the Drake Code and Framework .....	35
A.1 Introduction to the Drake Robotics Framework .....	35
A.2 Code Explanation: Pick-and-Place Pipeline .....	36
.....	38
A.3 Features of Drake Framework .....	39
A.4 Limitations of the Current Implementation .....	40
A.5 Suggested Improvements and Future Work.....	40
A.6 Conclusion .....	41
Appendix 02: Technical Implementation of ACT .....	42
A. Training Pipeline for ACT .....	42
B. Action Chunking and Temporal Ensembling.....	43
C. Modeling Human Demonstrations with Conditional Variational Autoencoders (CVAE).....	44
D. Implementation with Transformers .....	46
E. Training and Inference .....	47
F. Algorithmic Details and Steps:.....	47
G. Code Structure .....	49
Appendix 03: Technical Implementation of Diffusion Policies .....	53

## List of Figures

Figure 01: Pick and Place Pipeline using Drake Framework .....	7
Figure 02: The 6-DoF leader and follower robotic arms.....	12
Figure 03: Showing both cameras' field of view.....	13
Figure 04: Complete hardware setup.....	15
Figure 05: Using teleoperation collecting the dataset for the pick and place task .....	20
Figure 06: Pick and Place Task .....	23
Figure 07: Stacking Task.....	24
Figure 08: Orienting Task .....	25
Figure 09: Sliding Task .....	26
Figure 10: Action Chunking and Temporal Ensembling .....	44
Figure 11: Architecture of Action Chunking with Transformers: here, we are using only two camera inputs.....	46



## List of Tables

Table I: Success Rate (%) for Each Task Using ACT and Diffusion Policies	28
--	----

# Chapter 1 Introduction

The overarching vision in the field of robotics—especially robotic manipulation—has long been to design a versatile robot capable of autonomously performing a wide range of meaningful tasks, much like fictional examples such as Baymax or C-3PO. While this vision remains aspirational, recent advancements have significantly propelled progress in several key areas. For example, in artificial intelligence, AlphaGo, a reinforcement learning model trained entirely through self-play, defeated the world champion Lee Sedol, marking a major milestone in reinforcement learning methodologies[1]. This approach was later generalized to master Chess and Shogi, demonstrating the scalability of these techniques across multiple domains[2]. Similarly, computer vision underwent a pivotal transformation with the advent of AlexNet in 2012, sparking significant advancements in semantic segmentation, object detection, and pose estimation[3,4]. Robotics has also made impressive strides, ranging from self-driving vehicles to humanoid robots capable of dynamic locomotion and manipulation[5,6].

Despite these advancements, many robotic manipulation systems in use today remain largely unchanged for decades. Robots commonly employed in automobile manufacturing, for instance, still execute repetitive tasks such as welding and painting along pre-programmed trajectories without environmental feedback[7]. To expand the applicability of robotics, a paradigm shift is necessary—one that moves beyond rigid, structured environments and towards systems capable of adaptation and interaction with dynamic surroundings. This shift is already reflected in emerging industries, such as logistics, where companies like Covariant AI and Righthand Robotics have introduced adaptive solutions capable of functioning in unstructured environments[8].

This thesis explores and contributes to the development of advanced techniques in robotic manipulation, focusing on end-to-end learning methodologies that bridge the gap between traditional approaches and the demands of dynamic, real-world environments. By integrating these techniques with affordable hardware, the goal is to create systems that are

not only robust and adaptive but also cost-effective and accessible for broader industrial and societal applications.

## 1.1 Chapters Overview

This thesis develops novel contributions in several areas related to robotic manipulation, transitioning from classical robotics approaches to advanced learning-based methods for tackling real-world challenges.

Chapter 02: This chapter presents the development of a pick-and-place pipeline using classical robotics techniques implemented within the Drake robotics framework. The system employs the KUKA iiwa robotic arm and a pseudoinverse-based inverse kinematics controller to execute precise object manipulation tasks, such as transferring a brick between predefined poses. This chapter not only highlights the implementation details of this classical approach but also examines its inherent limitations, including sensitivity to object geometry, lack of perception, and challenges such as kinematic singularities. These insights set the stage for exploring more adaptive methodologies in subsequent chapters[9].

Chapter 03: This chapter focuses on learning-based robotic manipulation, presenting the implementation and evaluation of two advanced policies: Action Chunking with Transformers (ACT) and Diffusion Policies. Both policies are trained independently on tasks such as pick-and-place, sliding, stacking, and orientation, using datasets collected via a teleoperation system. ACT emphasizes temporal coherence, ensuring smooth task execution, while Diffusion Policies focus on trajectory refinement, achieving high precision in intricate tasks. The performance of these methods is systematically compared to analyze their strengths and limitations, providing insights into their applicability for real-world scenarios[10, 11].

## 1.2 Related Work

Robotic manipulation has been an active area of research, encompassing a wide range of techniques from classical control to modern learning-based approaches. This section reviews key developments in perception, planning, and control, focusing on their relevance to the methodologies employed in this thesis.

### Classical Robotics Approaches

Traditional robotics methods rely heavily on structured environments and precise modeling. Techniques such as pseudoinverse kinematics controllers have been widely used for industrial applications like assembly and packaging[12]. While these methods excel in predictable scenarios, their reliance on pre-defined trajectories and inability to adapt to dynamic environments limit their versatility[13].

### Perception in Robotics

The integration of perception has advanced significantly with the introduction of deep learning models like AlexNet, which revolutionized computer vision tasks[3]. These advancements have enabled robots to perform complex tasks such as object detection, semantic segmentation, and pose estimation. However, perception techniques in classical robotics often struggle with robustness in dynamic environments, as highlighted by Kragic and Christensen[14].

### Learning-Based Methods

A paradigm shift toward learning-based robotics has been evident in recent years. Reinforcement learning, exemplified by AlphaGo, demonstrates how agents can learn

complex planning and decision-making strategies[1, 2]. End-to-end learning, where raw sensory inputs are directly mapped to control outputs, has further expanded the scope of robotics, addressing challenges in adaptability and precision[15, 16]. Action Chunking with Transformers (ACT) introduces temporal coherence in action prediction, leveraging transformer architectures for sequential decision-making[17]. Meanwhile, Diffusion Policies, inspired by probabilistic diffusion models, iteratively refine action trajectories to achieve higher precision and robustness[18].

## Applications in Real-World Scenarios

The growing need for adaptable robotic systems has driven applications beyond traditional industrial automation. Adaptive learning-based systems have gained prominence in fields such as logistics, healthcare, and domestic applications, with companies like Covariant AI and Righthand Robotics demonstrating the feasibility of deploying these technologies in real-world settings[8].

## Chapter 2 Development of a Pick-and-Place Pipeline Using Drake and the KUKA iiwa Robot

The pick-and-place pipeline developed in this thesis leverages the Drake robotics framework as its foundation, showcasing its robust capabilities for modeling, simulation, and control of robotic systems. Drake's strength lies in its modularity, providing tools for accurate kinematic modeling, trajectory planning, and visualization, which are integral to this implementation. While the KUKA iiwa robotic arm is used as the manipulator in this pipeline, the framework's design ensures that the approach is not limited to any specific hardware. By utilizing Drake, the pipeline demonstrates a classical robotics solution for transferring objects between predefined poses, emphasizing the framework's adaptability and precision in simulating real-world manipulation tasks[17].

### 2.1 Overview of the Pick-and-Place Pipeline

The primary function of this pipeline is to move a simulated foam brick between two specified bins, following a sequence of carefully defined manipulative actions that illustrate a traditional pick-and-place operation. Using the Drake framework, we create a model for the KUKA iiwa and define key object and end-effector frames to enable precise planning, trajectory generation, and execution of the task.

#### Key Phases of the Pick-and-Place Task:

Frame Definition and Transformation:

- The task begins by defining the spatial transformations for both the object (brick) and the robotic gripper in the world coordinate frame. This allows for precise positioning relative to the environment. Using the function `MakeGripperFrames`,

the system calculates a sequence of frames that guide the gripper to approach, grasp, lift, and place the brick at each respective location in the bins[18].

- Initial and Goal Poses: The brick's starting position is defined as the initial pose, while its destination in the target bin is defined as the goal pose. Both are specified as transformations, combining position and orientation, which allows for consistency in calculating relative motion.
- Predefined Waypoints: Intermediate waypoints include pre-grasp (approach position), grasp, clearance (lifted position), pre-place (approach to target), and place positions. These are essential for avoiding collisions and enabling smooth transitions between actions.

#### Trajectory Planning and Inverse Kinematics:

- The task is executed by interpolating a smooth trajectory through the defined waypoints using inverse kinematics (IK). This involves calculating joint configurations that allow the gripper to follow the specified trajectory, moving the brick from the initial to the final pose[19].
- Pseudoinverse Controller: To achieve this, a pseudoinverse-based IK controller is employed. This controller computes joint velocities by inverting the manipulator Jacobian, translating end-effector velocities into corresponding joint velocities. The pseudoinverse approach is a classical technique in robotics, well-suited for tasks where real-time feedback is not required, and the manipulator operates within a controlled environment. However, as will be discussed, this method has known limitations, especially when facing singularities in the configuration space.

#### Sequential Execution of Phases:

- The task progresses through a sequence of defined motions:
- Approach: The gripper moves to a pre-grasp position aligned with the brick, establishing a stable orientation before initiating contact.

- Grasp: The gripper closes around the brick, which is precisely oriented for pickup based on the predefined frame alignment.
- Lift and Clearance: After grasping, the robot moves to a clearance position, keeping the brick securely lifted to avoid collisions with the bin edge.
- Placement: The gripper positions the brick over the target bin, and releases it.

Timing adjustments are incorporated for acceleration, deceleration, and stopping at each phase to ensure both stability and accuracy. Each phase involves controlled motions tailored to maintain the brick's orientation and avoid unnecessary contact with surroundings, which is essential for reliable execution in a simulated environment.

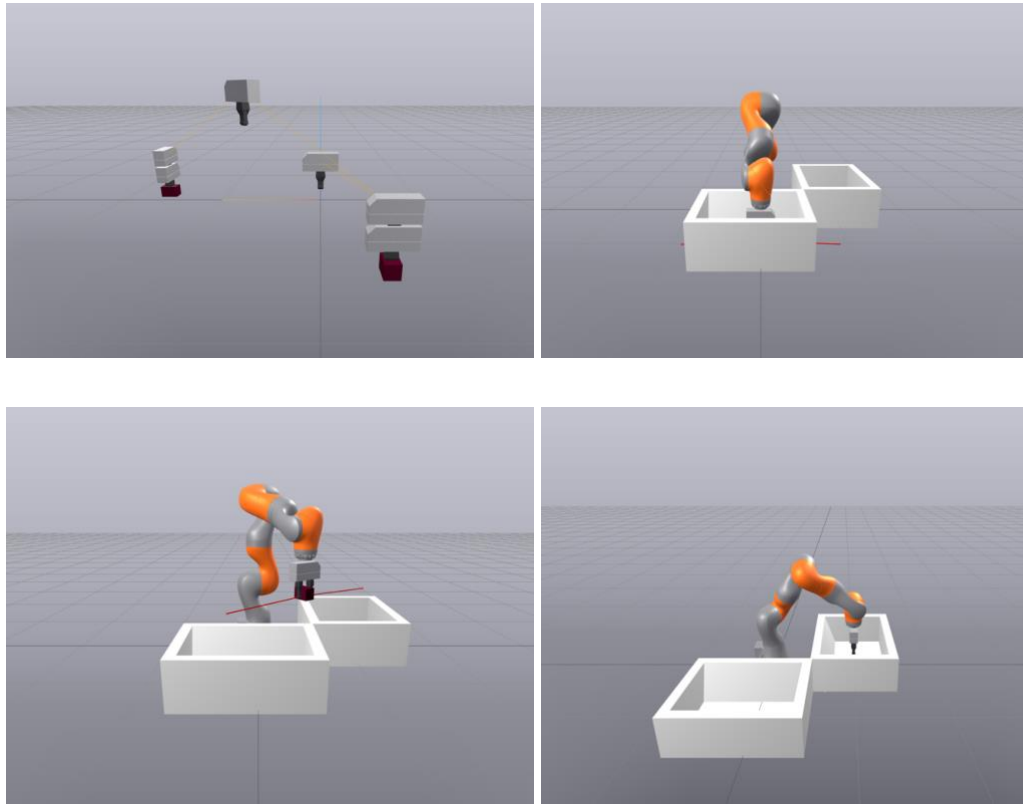


Figure 01: Pick and Place Pipeline using Drake Framework



## 2.2 Importance of Classical Approaches in Industrial Applications

Classical robotics approaches remain foundational in industrial automation due to their precision, reliability, and cost-effectiveness. These systems excel in highly structured environments, such as automotive manufacturing and electronics assembly, where tasks like welding, painting, and component handling require consistent performance[20].

### Key Advantages:

#### Efficiency and Simplicity:

- Classical pipelines operate on pre-programmed trajectories and eliminate the need for complex perception systems or advanced computations. This simplicity reduces maintenance requirements, ensures high throughput, and minimizes operational downtime, making them ideal for repetitive tasks in controlled environments[21].

#### Predictability and Safety:

- Their deterministic behavior aligns well with industrial safety standards. Predictability simplifies debugging and validation processes, ensuring safe and reliable operation even in collaborative environments where robots work alongside humans.

## 2.3 Limitations of the Classical Approach

While the pseudoinverse control approach achieves basic pick-and-place functionality, several limitations emerge that highlight the constraints of classical robotic control systems:

#### Dependency on Fixed Object Shape and Size:

- The pseudoinverse controller operates under the assumption of fixed, predictable object geometry. This limitation constrains the system to specific objects, as any deviation in shape or size could disrupt the predefined grasping and placement parameters[22].

#### Absence of Perception for Object and Environment Sensing:

- The system lacks dynamic perception capabilities to detect and adjust to the object's exact location or orientation. While classical perception methods like edge detection and feature matching can be implemented, they are prone to failure under varying lighting conditions, occlusions, and background changes, making them unreliable for flexible, real-world applications[23].

#### Kinematic Singularities:

- Singularities occur when the manipulator reaches configurations where the Jacobian matrix becomes non-invertible, leading to undefined or abrupt joint motions. These singularities limit the system's robustness and may reduce hardware lifespan due to high-stress motions[24].

#### Limited Adaptability and Scalability:

- Classical pipelines are effective for single-object tasks with clear start and end conditions but require significant recalibration for diverse or multi-object scenarios, reducing their applicability for complex industrial processes.

## Chapter 3 End-to-End Learning-Based Robotics Manipulation on Low-Cost Hardware

Recent advances in machine learning have reshaped robotic systems by enabling them to seamlessly integrate perception, planning, and control. Deep learning-based perception models have unlocked the ability to extract high-level information from raw sensory data, such as images, allowing robots to detect, classify, and localize objects with remarkable precision[17]. Simultaneously, reinforcement learning has showcased its effectiveness in planning and decision-making for complex, dynamic tasks[18]. Despite these advancements, many robotic systems adopt a modular design where perception, planning, and control are treated as distinct components. While this approach is functional, it often leads to inefficiencies due to the propagation of errors between modules and the need for extensive task-specific fine-tuning.

This chapter investigates an alternative paradigm: end-to-end learning. By directly mapping raw sensory inputs (pixels) to joint actions, end-to-end learning eliminates the need for intermediate representations, unifying perception, planning, and control into a single framework[19]. The chapter explores two innovative methodologies—Action Chunking with Transformers (ACT) and Diffusion Policies—which address challenges such as temporal consistency, compounding errors, and the generation of smooth, human-like motion. By deploying these techniques on cost-effective hardware, the chapter demonstrates their viability for real-world robotic manipulation tasks, highlighting the potential of end-to-end frameworks to democratize advanced robotics technologies.

### 3.1 Hardware Setup

The hardware setup forms the backbone of the framework, combining affordability, reproducibility, and functionality to support complex robotic manipulation tasks. It comprises two robotic manipulators configured in a leader-follower arrangement and a dual-camera vision system for comprehensive workspace coverage.

#### Robotic Manipulators

The system employs two six-degree-of-freedom (6-DoF) robotic manipulators, each powered by Dynamixel servos. These servos provide precise and smooth joint control at a low cost, making the setup accessible to smaller labs and educational institutions[20].

Leader-Follower Configuration:

- The leader manipulator is controlled directly by a human operator, enabling intuitive demonstrations of desired tasks[21].
- The follower manipulator mirrors the leader's movements in real time, ensuring high-quality data collection for training.

Key Features:

- **Affordability:** Each manipulator costs approximately \$100, significantly lowering the financial barrier for research and experimentation.
- **Precision:** Dynamixel servos enable fine-grained joint control, which is crucial for handling delicate objects or performing tasks requiring accuracy[22].
- **Simple Gripper Design:** Each manipulator is equipped with basic grippers capable of interacting with various objects, making the system versatile for diverse manipulation tasks.

The leader-follower configuration facilitates the recording of high-fidelity demonstrations. The human-guided leader ensures the tasks are performed naturally and efficiently, while the follower replicates the motions, creating a synchronized dataset of visual and proprioceptive information [23].

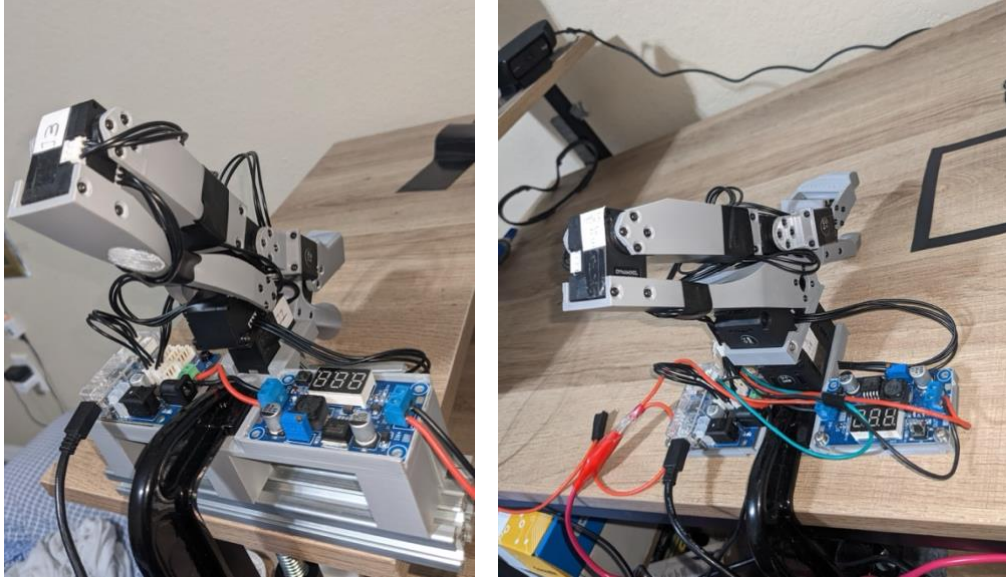


Figure 02: The 6-DoF leader and follower robotic arms. The leader arm is human-operated (right), and the follower arm mirrors its movements (left).

## Vision System

The perception system comprises two fixed RGB cameras strategically positioned to capture the full workspace of the follower manipulator. This configuration ensures robust visual data collection for both training and execution phases.

### Key Features:

- **Comprehensive Coverage:** The cameras are positioned to provide a complete view of the manipulator's workspace, ensuring all actions and object interactions are visible[24].
- **High Resolution and Frame Rate:** Recording at 30 Hz, the cameras capture smooth and detailed video streams, which are crucial for training robust perception models[25].
- **Stable Setup:** Fixed camera placement ensures consistent views of the environment, simplifying calibration and eliminating the need for moving camera mechanisms[26].

The dual-camera setup enhances spatial awareness by providing multiple perspectives of the workspace. This redundancy reduces the likelihood of occlusions and improves the reliability of object detection and pose estimation.

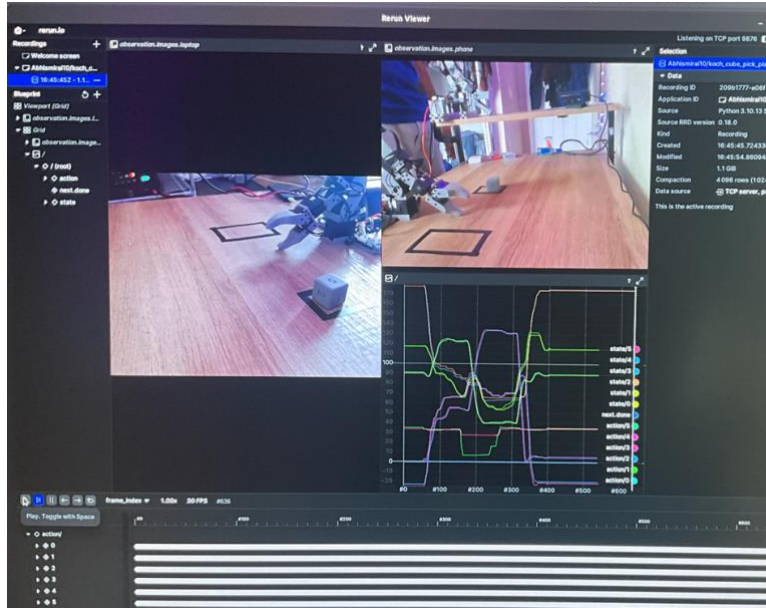


Figure 03: showing both cameras field of view

## Integration and Calibration

To ensure seamless operation, the robotic manipulators and cameras are calibrated for synchronized data collection. The follower manipulator's joint states are aligned with the visual data captured by the cameras, ensuring that each movement corresponds accurately to the observed environment. This synchronization is critical for creating a coherent dataset that bridges visual inputs with motor outputs.

### Calibration Process:

- The workspace is defined in a global coordinate frame shared by the manipulators and cameras.
- Camera intrinsic and extrinsic parameters are calibrated using standard techniques such as checkerboard patterns, ensuring accurate spatial alignment.
- The leader-follower synchronization is validated by performing simple pick-and-place tasks and comparing the observed trajectories with the planned motions.

## Cost and Accessibility

The hardware setup is deliberately designed to balance functionality with cost-effectiveness:

- Each manipulator, powered by affordable Dynamixel servos, costs approximately \$100, making them accessible for small-scale labs or classrooms[27].
- The RGB cameras are commercially available, further lowering the barrier to entry for researchers and educators interested in replicating the system.

By focusing on affordability and simplicity, this setup demonstrates that advanced robotic manipulation frameworks can be implemented without relying on expensive or specialized equipment[28].

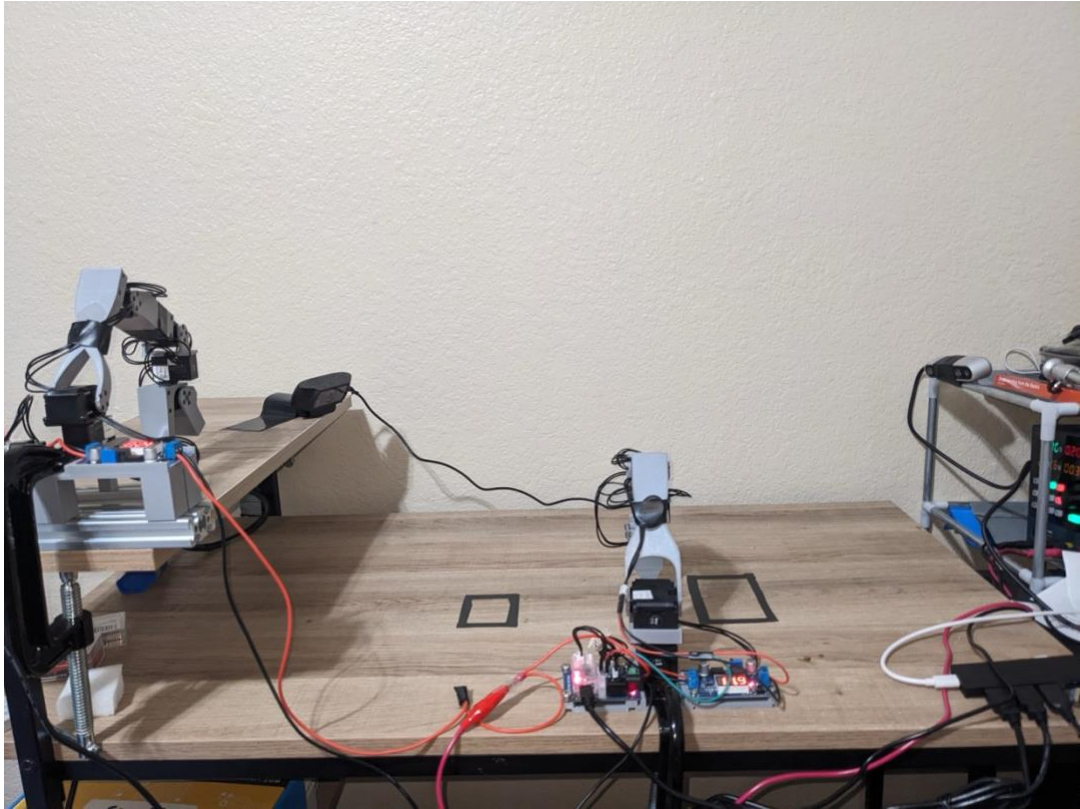


Figure 04: Complete hardware setup, including the dual-camera system and the follower arm within its workspace.

## Summary

The hardware configuration described in this chapter represents a scalable and affordable platform for robotic manipulation research. The leader-follower robotic arm setup ensures precise task demonstrations, while the dual-camera system provides robust visual coverage



of the workspace. Together, these components form the foundation for implementing and evaluating end-to-end learning methodologies, as discussed in subsequent sections. By emphasizing accessibility and reproducibility, this setup paves the way for wider adoption of learning-based robotic systems, enabling advancements in both academic research and practical applications.

## 4.2 Software Architecture

The software architecture of this framework is built on the robust foundation of the LeRobot codebase, an open-source platform designed for advanced machine-learning applications in robotics. LeRobot provides a comprehensive suite of tools that simplify data collection, policy implementation, training, and evaluation, making it particularly suitable for our objectives [18]. This flexibility and efficiency enable seamless integration of advanced robotic learning techniques, including Action Chunking with Transformers (ACT) and Diffusion Policies tailored specifically for low-cost hardware.

### Action Chunking with Transformers (ACT)

- **Sequential Action Prediction:** ACT enhances traditional imitation learning by predicting sequences of actions (referred to as "chunks") rather than processing one action at a time. By forecasting a sequence, ACT minimizes compounding errors and ensures smoother, temporally coherent robotic movements, making it especially effective for dynamic tasks [14].
- **Implementation in LeRobot:** The ACT implementation in LeRobot is based on a Conditional Variational Autoencoder (CVAE) architecture. The encoder extracts meaningful features from visual data and joint states, while the decoder generates precise sequences of joint actions. This approach leverages variability inherent in human demonstrations to produce reliable and robust control outputs.

- Temporal Ensembling: To achieve fluid task execution, ACT integrates temporal ensembling. This technique averages overlapping action chunks to mitigate abrupt transitions, resulting in natural, human-like motion. This capability, built within the LeRobot framework, ensures task smoothness and adaptability.

## Diffusion Policies

- Iterative Trajectory Refinement: Inspired by diffusion probabilistic models, Diffusion Policies use iterative refinement to transform an initial noisy trajectory into one that closely aligns with learned dynamics. This iterative process allows the system to adapt to environmental uncertainties and execute actions with high precision [15, 20].
- Modeling Complex Trajectories: The ability to model intricate trajectories directly in the action space makes Diffusion Policies particularly effective for tasks requiring fine-grained control, such as stacking objects or writing. LeRobot’s implementation of this method captures subtle variations in the training data, enabling nuanced and adaptable policy behaviors.
- Robustness to Variations: Diffusion Policies excel in handling dynamic and unpredictable environments. The iterative refinement adjusts actions in real time to account for changes in object positioning, size, or environmental conditions, making it highly versatile for real-world tasks.

## Integration with the LeRobot Codebase

LeRobot’s design emphasizes modularity and extensibility, ensuring smooth integration and operation of various components. Key advantages include:

- **Modular Design:** The architecture allows for easy customization and addition of modules, such as data collection tools, policy training routines, and evaluation metrics, making the system adaptable to diverse use cases [18].
- **Efficient Data Collection:** Using LeRobot, synchronized data recording is achieved with fixed cameras and robotic manipulators. This ensures that visual inputs are aligned with corresponding joint states, providing high-quality datasets essential for end-to-end learning.
- **Streamlined Training and Optimization:** LeRobot facilitates efficient training of both ACT and Diffusion Policies by leveraging modern deep learning frameworks with support for GPU acceleration. Tools for hyperparameter tuning and checkpointing are integrated into the framework, enabling iterative improvements in policy performance [15, 19].
- **Comprehensive Evaluation:** To evaluate policy effectiveness, LeRobot provides a range of metrics and visualization tools, offering insights into success rates, motion smoothness, and adaptability.
- **Community Support:** As an open-source initiative, LeRobot benefits from ongoing contributions by a global community of developers and researchers, ensuring continuous updates and improvements aligned with cutting-edge advancements.

## Summary

The integration of ACT and Diffusion Policies within the LeRobot framework enables the implementation of advanced end-to-end learning techniques on low-cost hardware. By leveraging LeRobot’s robust and modular design, our system achieves high precision and adaptability in robotic manipulation tasks. This synergy between software capabilities and hardware design underscores the viability of democratizing advanced robotics through accessible technologies, paving the way for broader research and application opportunities [27].

## 4.3 Integration and Workflow

The framework employs a structured, three-phase workflow—data collection, training, and execution—to rigorously evaluate the performance of Action Chunking with Transformers (ACT) and Diffusion Policies. By training and testing these policies independently on identical tasks, the system facilitates a direct comparison of their capabilities, offering insights into their respective strengths and trade-offs.

### 1. Data Collection

The data collection phase relies on a leader-follower teleoperation system to generate high-quality task demonstrations. This system captures realistic, human-guided manipulations that serve as the foundation for training the policies.

- **Leader-Follower Teleoperation:** The leader manipulator, controlled by a human operator, performs the desired tasks while the follower manipulator replicates these motions in real time. This setup ensures the collection of precise and naturalistic demonstration data[19].
- **Recording Setup:** Visual data is captured using two fixed RGB cameras positioned to comprehensively cover the workspace of the follower manipulator. The cameras operate at 30 Hz, providing smooth and high-resolution video feeds necessary for effective policy training[23].
- **Demonstration Diversity:** To enhance robustness, 100 demonstrations are recorded for each task, with each demonstration lasting approximately 30 seconds. Environmental variations, such as differing lighting conditions, are introduced to ensure the policies generalize well to diverse scenarios[25].

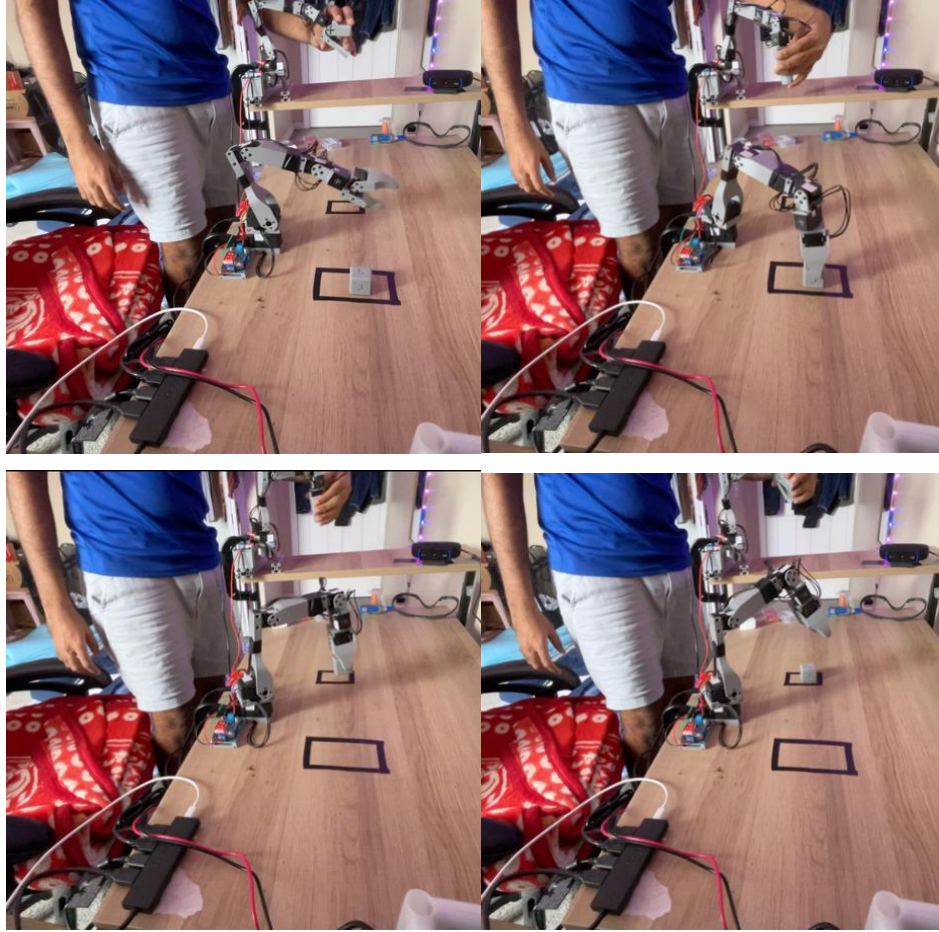


Figure 05: Using teleoperation collecting the dataset for the pick and place task

## 2. Training

The collected demonstration data is utilized to train the two policies—ACT and Diffusion Policies—independently. Each policy is optimized to excel in specific aspects of robotic manipulation, leveraging state-of-the-art learning techniques.

#### ACT Training:

- ACT is trained on the collected dataset to learn temporally coherent action sequences. By predicting action "chunks" instead of single steps, the policy minimizes compounding errors and ensures fluid transitions [14,20].
- A Conditional Variational Autoencoder (CVAE) architecture is employed for encoding and decoding action sequences, capturing the variability inherent in human demonstrations.
- Temporal ensembling is integrated to further enhance motion smoothness, reducing abrupt transitions and enabling natural execution of long-duration tasks.

#### Diffusion Policies Training:

- Diffusion Policies iteratively refine noisy initial trajectories, learning to produce precise actions through a series of refinement steps [15].
- The training process emphasizes modeling complex trajectory distributions directly in the action space, enabling adaptability to changes in object positions and environmental conditions [27].
- Both policies are trained using the LeRobot codebase, which provides robust tools for optimization, including the AdamW optimizer and custom learning rate schedules. Checkpoints are saved periodically, and metrics such as loss reduction and task success rates are monitored throughout the training process.

### 3. Execution

Following training, both policies are deployed independently to execute a range of robotic manipulation tasks. Each policy is evaluated based on metrics like task success rate, motion smoothness, and adaptability to dynamic environments.

ACT Execution:

- ACT generates temporally coherent action sequences, excelling in tasks that require consistency over long durations, such as pick-and-place operations and object sliding. Its smooth and predictable motion is well-suited for scenarios demanding stable interactions [14].

Diffusion Policies Execution:

- Diffusion Policies prioritize trajectory precision, making them ideal for tasks that require intricate manipulations, such as orienting objects or writing specific patterns. Their iterative refinement process ensures high accuracy and adaptability, even in challenging conditions [15, 18].

By comparing the two policies on identical tasks, the framework provides a clear understanding of their strengths and limitations. ACT demonstrates superior temporal coherence, while Diffusion Policies excel in precision and adaptability, offering complementary solutions for robotic manipulation tasks.

## 4.4 Experiments: Task Descriptions and Approach

This section elaborates on the carefully designed tasks aimed at evaluating the capabilities of an end-to-end learning-based robotic manipulation system. These tasks are approached without any explicit prior knowledge of the robotic system, the objects being manipulated, or the environment, thereby emphasizing the system's ability to learn directly from raw sensory data and adapt its behaviors accordingly. Each task represents a distinct challenge, testing various aspects of robotic manipulation in progressive levels of complexity.

## 1. Pick-and-Place

The pick-and-place task is foundational in robotic manipulation and serves as the starting point for evaluating the system's basic capabilities. Here, the robot transfers an object from one box to another by:

- **Locating the Object:** Identifying the object using raw sensory data without prior calibration.
- **Grasping Securely:** Executing a stable and precise grasp that minimizes slippage or misalignment.
- **Executing Trajectories:** Planning and following a trajectory to move the object from its initial position to its target location.
- **Releasing Precisely:** Ensuring the object is deposited in the target location with accuracy.

This task is critical as it forms the foundation for more complex operations, evaluating basic trajectory planning, grasping mechanics, and object placement. Success in this task demonstrates the system's ability to integrate perception and control seamlessly in an end-to-end learning framework.

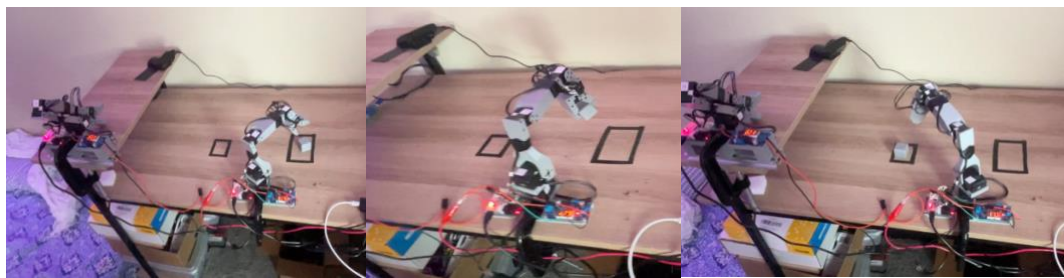


Figure 06: Pick and Place Task



## 2. Stacking Objects

The stacking task introduces the challenge of manipulating multiple objects and maintaining balance during the operation. In this task:

- **Precision and Alignment:** The robot must pick up an object and carefully align it with another object already placed on a surface.
- **Stable Placement:** The system must maintain the stability of the stack while depositing the object.
- **Error Tolerance:** The system learns to compensate for minor discrepancies in object shape or placement without explicit programming.

Stacking requires precise coordination between the robot's visual feedback and its motor actions. It tests the system's ability to maintain stability under increasing physical constraints, highlighting its adaptability in handling multi-object scenarios.



Figure 07: Stacking Task

## 3. Orienting an Object

This task involves the robot aligning an object with a specified orientation in three-dimensional space. The steps include:

- **Pose Estimation:** Determining the object's current orientation based on visual inputs.

- **Planning Rotational Adjustments:** Mapping the observed pose to the target orientation and generating appropriate rotational actions.
- **Executing Rotations:** Accurately manipulating the object to align it with the desired pose.

This task tests the robot's proficiency in rotational precision and its ability to execute fine-grained adjustments. The absence of explicit prior knowledge about object orientation or dynamics underscores the importance of the system's learning capabilities.



Figure 08: Orienting Task

#### 4. Sliding an Object

Sliding is one of the most challenging tasks due to the dynamic interaction required with the environment. The robot must:

- **Maintain Contact:** Keep consistent pressure and contact with the object's surface to avoid slippage.
- **Control Force and Motion:** Balance the force applied and trajectory followed to ensure controlled motion across the surface.
- **Reach a Target Position:** Stop the object precisely at the intended location.

This task demands exceptional force control and dynamic stability, testing the framework's ability to handle continuous feedback and execute complex surface interactions. It is a significant benchmark for evaluating the robustness of the system.

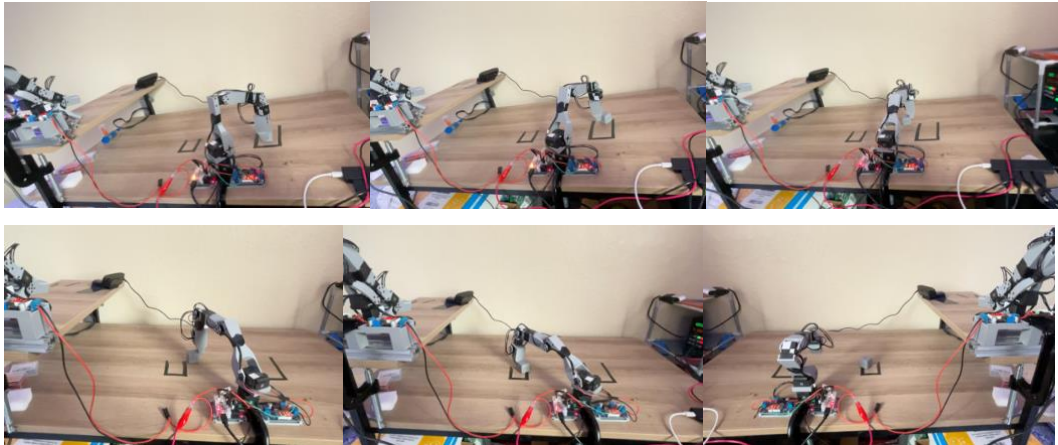


Figure 09: Sliding Task

## Progressive Complexity and Framework Design

These tasks are intentionally structured to progress in complexity:

- Starting with basic object manipulation (pick-and-place), the tasks incrementally test higher-level capabilities such as stability (stacking), rotational precision (orienting), and dynamic interaction (sliding).
- The design ensures that the system is evaluated across a spectrum of scenarios, from static environments to dynamic surface interactions, providing a comprehensive assessment of its capabilities.

## Learning without Explicit Knowledge

The core of this experimental framework lies in its reliance on end-to-end learning. Unlike traditional systems that rely on explicit models of the robot or environment:

- The system learns directly from data captured during demonstrations, bypassing the need for predefined rules or constraints.
- Raw sensory inputs (camera feeds and manipulator joint states) are mapped directly to actions, showcasing the framework's ability to generalize across varying conditions.

This approach challenges the limits of modern robotics by emphasizing adaptability and flexibility, making it applicable to real-world scenarios where explicit programming is impractical. The tasks not only demonstrate the strengths of the end-to-end learning framework but also highlight areas where further improvements in adaptability and precision may be necessary.

## Evaluation Criteria

The training process employed several metrics to evaluate policy performance and ensure effective learning:

- **Loss Function:** Both ACT and Diffusion Policies are trained to minimize a loss function that accounts for action prediction errors. For ACT, the loss includes reconstruction loss from the CVAE, while for Diffusion Policies, it includes trajectory refinement errors.
- **Success Rate:** The percentage of successful task completions during evaluation. Success is determined by task-specific criteria, such as correctly placing an object or achieving a specific orientation.
- **Smoothness of Motion:** The mean and variance of joint velocities during task execution are analyzed to ensure smooth and natural movements.

## 4.5 Results and Discussion

The results for the experiments conducted on the four tasks—Pick-and-Place, Sliding an Object, Stacking Objects, and Orienting an Object—are presented in the table below. Each task was tested 50 times for both Action Chunking with Transformers (ACT) and Diffusion Policies. The success rate is expressed as a percentage, and the table highlights the system’s ability to perform these tasks under varying conditions, considering the limitations of low-cost hardware.

Table I: Success Rate (%) for Each Task Using ACT and Diffusion Policies

Task	ACT	Diffusion Policies
Pick-and-Place	83%	85%
Sliding an Object	65%	72%
Stacking Objects	78%	82%
Orienting an Object	75%	83%

Table I: Success rates for each task conducted under end-to-end learning with ACT and Diffusion Policies. Each policy was tested 50 times on the same setup, ensuring comparable results.

## Analysis of Results

### Overall Observations

#### Success Rates Reflect Hardware Limitations:

- Both ACT and Diffusion Policies performed admirably given the constraints of low-cost hardware. The success rates reflect the robustness of the methods but also indicate the inherent challenges posed by the hardware, particularly in tasks requiring higher precision or dynamic interactions.

#### Task-Specific Performance:

- Pick-and-Place: This foundational task showed the highest success rates for both policies. Its relatively straightforward nature—requiring basic object manipulation—allowed both policies to perform well, with Diffusion Policies slightly edging out ACT due to better trajectory refinement.
- Sliding an Object: Sliding proved to be the most challenging task due to its reliance on maintaining consistent surface contact and precise force control. The success rates here were notably lower than for other tasks, highlighting the difficulty of dynamic surface interactions with low-cost hardware.
- Stacking Objects: Stacking required careful alignment and stability, where Diffusion Policies demonstrated an advantage due to their ability to refine trajectories iteratively. ACT performed well but struggled slightly with maintaining balance during the stacking process.
- Orienting an Object: The orientation task demanded accurate rotational adjustments, which were better handled by Diffusion Policies due to their superior modeling of complex trajectories. ACT, while slightly less precise, still achieved commendable results by leveraging its temporal coherence.

### Policy Comparisons:

- ACT: ACT's strength lies in its temporal coherence, ensuring smooth execution of sequential actions. This approach worked well for tasks with fewer precision demands, such as Pick-and-Place, but struggled with tasks requiring intricate adjustments like Sliding or Orienting an Object.
- Diffusion Policies: Diffusion Policies excelled in precision-heavy tasks due to their iterative refinement mechanism. This allowed them to adapt to small variations in object positions and environmental conditions, making them better suited for tasks like Stacking and Orienting.

### Key Observations About the Setup and Methods

#### Robustness Despite Hardware Limitations:

- The system remained robust against environmental variations such as slight lighting changes, minor arm displacement, and slight inaccuracies in camera alignment. This robustness underscores the adaptability of end-to-end learning, even on low-cost hardware.

#### Challenges with Sliding Task:

- The Sliding task posed unique challenges due to the difficulty of force control and maintaining consistent contact with the surface. The relatively lower success rates for this task highlight the limitations of low-cost hardware in handling dynamic tasks but also demonstrate the policies' potential to achieve reasonable performance under such constraints.

#### Consistency Across Tasks:

- Both policies demonstrated consistent performance across a variety of tasks, with Diffusion Policies generally achieving higher success rates due to their ability to iteratively refine actions. However, ACT's performance was close, particularly in simpler tasks like Pick-and-Place.

#### Real-Time Execution and Simplicity:

- The methods were computationally efficient, achieving real-time execution without requiring complex system tuning or task-specific adjustments. This is particularly significant given the affordability and simplicity of the hardware.

#### Scalability and Generalization:

- The system's ability to perform well across diverse tasks without explicit modeling of the environment highlights its potential for scalability. This generalization capability is critical for deploying such systems in dynamic, real-world settings.

### Conclusion of Results

The results highlight the strengths and limitations of ACT and Diffusion Policies for end-to-end robotic manipulation using low-cost hardware. While Diffusion Policies excelled in precision-heavy tasks like Orienting and Stacking, ACT demonstrated robustness in maintaining smooth and consistent action sequences. The observed robustness to environmental variations and the system's ability to generalize across tasks emphasize the viability of these methods for real-world applications, even in resource-constrained settings. These findings pave the way for future improvements in both policies, focusing on enhancing adaptability and precision for more complex and dynamic tasks



## References

1. Silver, D. et al. (2016). "Mastering the game of Go with deep neural networks and tree search." *Nature*.
2. Silver, D. et al. (2017). "Mastering Chess and Shogi by self-play with a general reinforcement learning algorithm." *arXiv*.
3. Krizhevsky, A. et al. (2012). "ImageNet classification with deep convolutional neural networks." *NeurIPS*.
4. Long, J. et al. (2015). "Fully convolutional networks for semantic segmentation." *CVPR*.
5. Badue, C. et al. (2021). "Self-Driving Cars: A Survey." *IEEE Transactions on Intelligent Transportation Systems*.
6. Park, H. W. et al. (2017). "Quadruped control for highly dynamic locomotion and manipulation tasks." *Science Robotics*.
7. Siciliano, B. et al. (2010). *Robotics: Modelling, Planning and Control*. Springer.
8. Covariant AI. (2020). "AI in Logistics Automation." *Company Blog*.
9. Craig, J. J. (2005). *Introduction to Robotics: Mechanics and Control*. Pearson.
10. Levine, S. et al. (2016). "End-to-End Training of Deep Visuomotor Policies." *Journal of Machine Learning Research*.
11. Schulman, J. et al. (2017). "Proximal Policy Optimization Algorithms." *arXiv*.
12. Siciliano, B. & Khatib, O. (2016). *Springer Handbook of Robotics*. Springer.

13. Kragic, D. & Christensen, H. I. (2002). "Survey on Visual Servoing for Manipulation." Robotics and Autonomous Systems.
14. Vaswani, A. et al. (2017). "Attention is All You Need." NeurIPS.
15. Ho, J. et al. (2020). "Denoising Diffusion Probabilistic Models." NeurIPS.
16. Covariant AI and Righthand Robotics. (2021). Application Reports.
17. Patil, S. et al. (2015). "Motion Planning for Robotic Manipulation: Online Adaptation and Nonlinear Optimization." IJRR.
18. Siciliano, B. et al. (2010). Robotics: Modelling, Planning and Control. Springer.
19. Murray, R. M. et al. (1994). A Mathematical Introduction to Robotic Manipulation. CRC Press.
20. LeRobot GitHub Repository. <https://github.com/LeRobotFramework>.
21. Craig, J. J. (2005). Introduction to Robotics: Mechanics and Control. Pearson.
22. Kragic, D., & Christensen, H. I. (2002). "Survey on Visual Servoing for Manipulation." Robotics and Autonomous Systems.
23. Levine, S. et al. (2016). "End-to-End Training of Deep Visuomotor Policies." Journal of Machine Learning Research.
24. Park, H. W. et al. (2017). "Quadruped Control for Highly Dynamic Tasks." Science Robotics.
25. Schulman, J. et al. (2017). "Proximal Policy Optimization Algorithms." arXiv.
26. Vaswani, A. et al. (2017). "Attention is All You Need." NeurIPS.
27. Ho, J. et al. (2020). "Denoising Diffusion Probabilistic Models." NeurIPS.

28. Siciliano, B., Sciavicco, L., Villani, L., & Oriolo, G. (2016). Robotics Handbook. Springer.
29. Badue, C. et al. (2021). "Self-Driving Cars: A Survey." IEEE Transactions on Intelligent Transportation Systems.

# Appendix

## Appendix 01: Detailed Explanation of the Drake Code and Framework

This appendix provides an in-depth explanation of the code implemented using the Drake robotics framework and elaborates on the features and functionalities of Drake as a whole. The goal is to give readers a comprehensive understanding of how the pick-and-place pipeline was developed, modeled, and executed within this framework.

### A.1 Introduction to the Drake Robotics Framework

Drake is an open-source robotics simulation framework developed by the Robot Locomotion Group at MIT. It provides a suite of tools for modeling, simulation, and control of robotic systems with a focus on dynamical analysis, optimization, and physical interactions. Its modular design makes it suitable for a wide range of robotic applications, including but not limited to:

- **Kinematic and Dynamic Modeling:** Drake supports precise modeling of robotic manipulators, including joint limits, collisions, and constraints.
- **Trajectory Optimization:** The framework offers tools for solving optimization problems related to motion planning and control.
- **Visualization:** Integration with tools such as MeshCat and RViz provides real-time visualization of simulations, aiding in debugging and analysis.
- **Hardware-Agnostic Design:** Drake allows the user to work with various robotic arms and actuators, making it a flexible choice for prototyping and development.

## A.2 Code Explanation: Pick-and-Place Pipeline

The code implemented in this thesis uses Drake's core modules to design a pick-and-place pipeline. Below is a detailed breakdown of the code functionality.

### A.2.1 Setting Up the Robot and Environment

The initialization phase involves defining the robot model, workspace, and target objects. This is achieved using Drake's MultibodyPlant and SceneGraph modules:

- **Robot Definition:** The KUKA iiwa robot is defined as a MultibodyPlant instance. Drake loads the URDF file of the robotic arm to configure its kinematic and dynamic properties.
- **Scene Setup:** The bins and foam brick are modeled as rigid bodies with predefined positions and orientations within the workspace. The SceneGraph is used to manage the spatial relationships between these objects.

```
# Example code snippet

plant = MultibodyPlant(time_step=0.001)

scene_graph = SceneGraph()

plant.RegisterAsSourceForSceneGraph(scene_graph)

iiwa_model = Parser(plant).AddModelFromFile("kuka_iiwa.urdf")
```

### A.2.2 Frame and Transformation Definitions

The pick-and-place task relies on precise spatial transformations to guide the robotic gripper. Key functions include:

- **Defining Object Frames:** Frames for the brick and gripper are defined in the world coordinate system. These frames are critical for calculating the relative transformations required for grasping and placing.
- **Gripper Alignment:** The `MakeGripperFrames` function computes the approach, grasp, and clearance positions relative to the brick's initial and target poses.

```
# Example frame definition  
  
gripper_frame = plant.GetFrameByName("gripper")  
  
brick_frame = plant.GetFrameByName("brick")  
  
gripper_to_brick_transform = RigidTransform(RollPitchYaw(0, 0, 0), [0, 0, 0.2])
```

### A.2.3 Trajectory Planning

Drake's trajectory optimization tools are used to plan smooth motions for the robotic arm:

- **Waypoint Interpolation:** Waypoints, including pre-grasp, grasp, lift, and place, are interpolated to generate a continuous trajectory.
- **Inverse Kinematics (IK):** The pseudoinverse controller solves the IK problem by converting the desired end-effector poses into joint configurations.

```
# Example trajectory planning

waypoints = [pre_grasp_pose, grasp_pose, lift_pose, place_pose]

trajectory = PiecewisePolynomial.CubicShapePreserving(
    time_points, waypoints, zero_slope_start, zero_slope_end)
```

#### A.2.4 Execution of the Pick-and-Place Task

The sequence of motions is executed by commanding the robotic arm to follow the planned trajectory. Key steps include:

- Gripper Control: Commands to open and close the gripper are issued at the appropriate phases (e.g., grasp and release).
- Collision Avoidance: The system ensures that the gripper avoids collisions with the bin edges and other objects during motion.

```
# Example execution loop

for pose in trajectory:

    plant_context = plant.CreateDefaultContext()

    plant.SetPositions(plant_context, robot_model, pose)
```

## A.3 Features of Drake Framework

### A.3.1 Modularity and Flexibility

Drake's modular architecture allows users to integrate custom components and algorithms seamlessly. For example:

- Custom Controllers: Users can implement their own control logic using Python or C++.
- Multi-Robot Systems: Drake supports the modeling of collaborative robots and multi-agent systems.

### A.3.2 Realistic Physics Simulation

Drake uses robust solvers for simulating contacts, forces, and dynamics. This ensures high-fidelity results that closely mimic real-world behavior, which is essential for validating control strategies.

### A.3.3 Trajectory Optimization Tools

Drake's optimization library provides solvers for motion planning, including:

- Direct Collocation: A method for optimizing trajectories subject to physical constraints.
- Contact-Rich Tasks: Tools for planning motions that involve dynamic interactions with the environment, such as pushing or sliding.



### A.3.4 Visualization

Integration with visualization tools like MeshCat allows users to:

- Observe simulations in real-time.
- Debug trajectory planning issues.
- Validate the alignment of frames and transformations.

## A.4 Limitations of the Current Implementation

While the Drake framework offers powerful tools, certain limitations exist in the current implementation:

1. Singularity Handling: The pseudoinverse controller does not handle kinematic singularities effectively, which may lead to erratic motions.
2. Lack of Real-Time Feedback: The pipeline operates in open-loop control mode, limiting its ability to adapt to disturbances.
3. Perception Integration: The current implementation lacks perception capabilities, which are necessary for dynamic and unstructured environments.

## A.5 Suggested Improvements and Future Work

To enhance the pipeline and leverage Drake's full potential, the following improvements are suggested:

- Integrate Perception Systems: Incorporate vision-based sensing to enable real-time object detection and pose estimation.

- **Feedback Control:** Develop a closed-loop control strategy to improve robustness against external disturbances.
- **Kinematic Redundancy Resolution:** Implement advanced techniques for singularity avoidance and redundancy optimization.

## A.6 Conclusion

The Drake robotics framework provides a comprehensive platform for developing robotic applications, combining advanced modeling tools with realistic physics simulation. The pick-and-place pipeline developed in this thesis demonstrates the capabilities of Drake for implementing classical robotics solutions. However, the limitations identified in this appendix highlight the need for further development to address the challenges posed by real-world applications. This appendix serves as a detailed resource for understanding the implementation and potential of the Drake framework in robotic manipulation tasks.

## Appendix 02: Technical Implementation of ACT

In this section, we delve into the detailed implementation of Action Chunking with Transformers (ACT) tailored to our hardware setup using the LeRobot framework. This method represents an advanced, robust approach for imitation learning, focusing on fine-grained tasks that demand high-frequency control, closed-loop feedback, and smooth action sequences. By leveraging LeRobot's architecture and dual-camera setup, we systematically train ACT to predict and execute action sequences with minimal demonstrations on low-cost hardware.

### A. Training Pipeline for ACT

To train ACT for robotic manipulation tasks, we start by collecting human demonstrations using the leader-follower teleoperation system. Here, the leader arm (controlled by a human operator) provides the input joint positions, while the follower arm mimics these movements. These joint positions represent the target actions. Observations for training consist of joint positions of the follower arm and synchronized image feeds from two RGB cameras.

Observation and Action Collection:

- **Input Observations:** The robot observes its environment through two cameras, positioned to capture the entire workspace, and joint proprioceptive data.
- **Target Actions:** Actions are defined as the next step's target joint positions for the follower manipulator. The observations and actions ensure that the policy mimics the human operator's next steps based on current sensory inputs.
- **Force Dynamics:** Using the leader's joint positions (instead of the follower's) captures implicit force application through the PID controller, enabling precise low-level control.

Training Objective:

- The ACT policy predicts a sequence of future actions (joint positions for both arms) for the next  $k$  steps, given the current observations.
- The task involves minimizing prediction errors while ensuring smooth, coherent action sequences across time.

Challenges:

- Compounding Errors: Errors in previous predictions can cause the system to enter untrained states, leading to poor performance. Action chunking and temporal ensembling address this issue effectively.

## B. Action Chunking and Temporal Ensembling

Action Chunking: Action chunking reduces the effective trajectory horizon by grouping  $k$  consecutive actions into chunks. Instead of predicting actions for every timestep

( $\pi_\theta(a_t|s_t)\pi_\theta(a_{t+1}|s_{t+1})\dots\pi_\theta(a_{t+k}|s_{t+k})$ ), the policy predicts a chunk of  $k$  actions

( $\pi_\theta(a_{t:t+k}|s_t)\pi_\theta(a_{t+1:t+k}|s_{t+1})\dots\pi_\theta(a_{t+k}|s_{t+k})$ ) at once. This approach:

- Decreases the sensitivity to compounding errors.
- Models non-Markovian behaviors in human demonstrations, such as pauses or sequential dependencies within a chunk.
- Facilitates efficient task execution by executing actions in grouped sequences.

Temporal Ensembling: To improve motion smoothness and avoid abrupt transitions

- Overlapping action chunks are used, where actions are predicted at every timestep but executed sequentially with overlap.
- Predictions are combined using a weighted temporal ensemble:  
 $w_i = \exp(-m \cdot i)$  where  $w_i$  represents the weight for the  $i$ -th action in the ensemble. A smaller  $m$  leads to faster incorporation of new observations, ensuring real-time adaptability.
- Temporal ensembling reduces jerky motions while preserving precision, as multiple predictions for the same timestep are aggregated without bias.

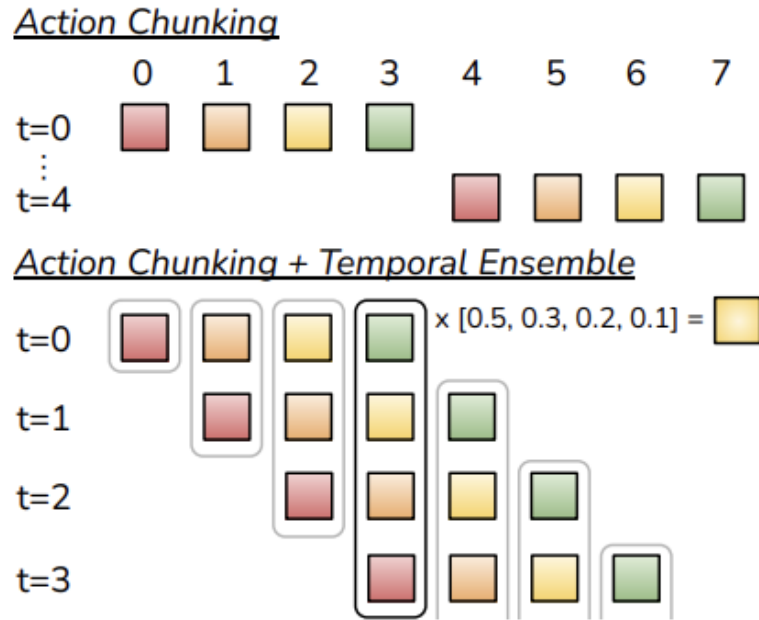


Figure 10 Action Chunking and Temporal Ensembling

### C. Modeling Human Demonstrations with Conditional Variational Autoencoders (CVAE)

Given the variability in human demonstrations, ACT uses a CVAE framework to learn from stochastic and noisy data effectively.

CVAE Architecture:

- **Encoder: Processes joint positions and observed action sequences to predict the style variable  $z$ , representing latent task characteristics. This helps the policy handle diverse demonstration strategies.**
- **Decoder (Policy): Combines  $z$  and the current observations (camera images + joint positions) to predict the next  $k$  actions. At test time,  $z$  is set to the mean of the Gaussian prior, ensuring deterministic and smooth action prediction.**

Training Objective:

The policy is trained to maximize the log-likelihood of action chunks while minimizing reconstruction errors. The loss function comprises:

- **Reconstruction Loss:** Measures the accuracy of predicted action chunks.
- **Regularization Loss:** Regularizes  $z$  towards a Gaussian prior to balance precision and generalizability.

Image and Joint Processing:

- **Image Processing:** Each camera image ( $480 \times 640$  RGB) is encoded using a ResNet18 backbone into a feature map ( $15 \times 20 \times 512$ ). Positional embeddings are added to preserve spatial information.
- **Joint Position Processing:** Joint positions and the latent variable  $z$  are linearly projected to match the feature dimensions (512).
- **Transformer Encoder Input:** Combines features from images, joint positions, and  $z$  into a sequence for processing by the transformer encoder.
- **Decoder Output:** The transformer decoder generates a  $k$ -step action sequence, refined through cross-attention and down-projected to joint positions.

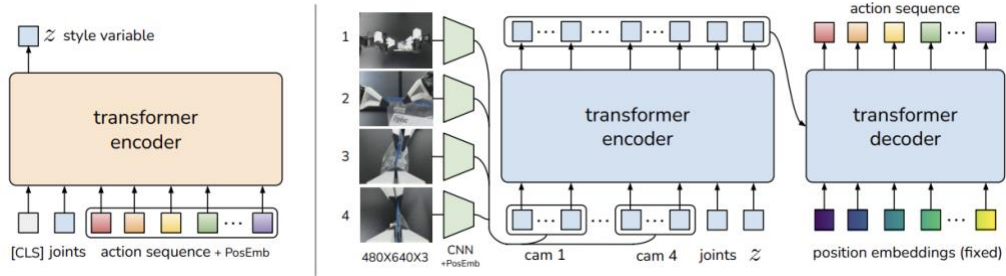


Figure 11 Architecture of Action Chunking with Transformers: here, we are using only two camera inputs

## D. Implementation with Transformers

Transformer-Based CVAE Encoder:

- Input: The current joint positions and the corresponding k-step action sequence, prepended with a learned “[CLS]” token.
- Output: Mean and variance of  $z$ , used to condition the policy.

Transformer-Based CVAE Decoder:

- Cross-Attention Mechanism: Synthesizes image, joint, and latent features into coherent k-step action sequences.
- Output Processing: Predicts  $k \times 14k \times 14k \times 14$  tensors representing joint positions for the next k timesteps.

Loss Function Details:

- ACT employs L1 loss for action reconstruction, prioritizing precise motion over smoother approximations achieved by L2.

## E. Training and Inference

### Training:

- Conducted from scratch for each task using the LeRobot framework.
- Training involves optimizing the CVAE objective over collected demonstration data, with hyperparameters tuned for low-cost hardware.

### Inference:

- Policies achieving the lowest validation loss are deployed.
- Real-time inference predicts smooth action sequences at a rate of 50 Hz, supported by the lightweight transformer architecture.

## F. Algorithmic Details and Steps:

### 1) Training Algorithm for ACT

#### Input Preparation:

- Collect demonstration data comprising camera images and joint positions.
- Normalize input data for uniform scaling across all tasks.
- Preprocess camera images using ResNet18 to obtain spatial feature maps.

#### CVAE Encoder:

- Input the current observation (joint positions + action sequence).
- Pass through a BERT-style transformer encoder to generate latent variable  $z$ .

#### CVAE Decoder (Policy):

- Combine  $z$  with visual and proprioceptive inputs.



- Use cross-attention in the transformer decoder to predict a coherent action sequence for the next  $k$  steps.

Loss Calculation:

- Compute L1 reconstruction loss for action sequences.
- Add regularization for  $z$  to enforce alignment with the Gaussian prior.
- Optimize the combined loss to minimize prediction errors.

Model Checkpoints:

- Periodically save model checkpoints based on validation loss.
- Use early stopping to prevent overfitting.

## 2) Inference Algorithm for ACT

Load Pretrained Model:

- Initialize the transformer encoder and decoder with weights from the best-performing checkpoint.

Observation Processing:

- Encode current camera images and joint positions into feature vectors.
- Generate  $z$  using the prior distribution (mean of Gaussian).

Chunk Prediction:

- Predict the next  $k$  actions using the CVAE decoder.
- Apply temporal ensemble to smooth overlapping predictions.

Action Execution:

- Send predicted joint positions to the low-level PID controller for precise execution.

## G. Code Structure

Modules:

- ACTEncoder: Encodes the input features (e.g., joint positions and visual data) into a compact latent representation.
- ACTDecoder: Decodes the latent representation into action chunks.
- ACTPolicy: The overarching module that integrates the encoder, decoder, and transformer-based temporal processing.

Key Functions:

ACTEncoder:

- Inputs:
  - Visual features (e.g., from RGB frames processed using ResNet).
  - Proprioceptive features (e.g., joint positions, velocities).
- Layers:
  - Linear layers: Used to downsample and embed the input features.
  - Batch Normalization: Applied to stabilize training and normalize input distributions.
  - ReLU Activation: Ensures non-linearity for feature learning.
- Output:
  - Encoded latent space ( $z$ ) and its mean and variance for the KL divergence loss.

```

class ACTEncoder(nn.Module):

    def __init__(self, input_dim, latent_dim):

        super().__init__()

        self.fc1 = nn.Linear(input_dim, 128)

        self.fc2 = nn.Linear(128, latent_dim * 2) # For mean and variance

    def forward(self, x):

        h = F.relu(self.fc1(x))

        mu_var = self.fc2(h)

        mu, logvar = mu_var.chunk(2, dim=-1) # Splitting into mean and log variance

        return mu, logvar

```

ACTDecoder:

- Inputs:
  - The sampled latent variable (z) from the encoder.
- Layers:
  - Transformer blocks: Capture temporal dependencies across action sequences.
  - Linear layers: Map the transformer outputs back to the action space.
- Output:
  - Predicted action sequence chunk.

```

class ACTDecoder(nn.Module):

    def __init__(self, latent_dim, output_dim):

        super().__init__()

        self.transformer = nn.Transformer(d_model=latent_dim, nhead=4,
num_encoder_layers=2)

        self.fc = nn.Linear(latent_dim, output_dim)

    def forward(self, z, seq_len):

        # Repeatedly expand `z` across the sequence length

        z = z.unsqueeze(1).repeat(1, seq_len, 1)

        h = self.transformer(z)

        return self.fc(h)

```

ACTPolicy:

- Combines the encoder, decoder, and sampling mechanism for the latent variable.
- Loss Function:
  - Reconstruction loss: Measures how well the predicted actions match the ground truth.
  - KL divergence: Regularizes the latent space to be Gaussian.

```

class ACTPolicy(nn.Module):

    def __init__(self, input_dim, latent_dim, output_dim):
        super().__init__()

        self.encoder = ACTEncoder(input_dim, latent_dim)
        self.decoder = ACTDecoder(latent_dim, output_dim)

    def forward(self, x, seq_len):
        mu, logvar = self.encoder(x)
        z = self.reparameterize(mu, logvar)
        return self.decoder(z, seq_len), mu, logvar

    def reparameterize(self, mu, logvar):
        std = torch.exp(0.5 * logvar)
        eps = torch.randn_like(std)
        return mu + eps * std

```

## Appendix 03: Technical Implementation of Diffusion Policies

Diffusion Policies implement iterative trajectory refinement using a denoising process inspired by diffusion probabilistic models. Below is the breakdown of its implementation:

### Code Structure

#### Modules:

- **DiffusionPolicy:** Implements the denoising process across multiple refinement steps.
- **ScoreNetwork:** Predicts the gradient (or score) of the trajectory distribution.

#### Key Functions: ScoreNetwork

- **Inputs:**
  - Noisy trajectory samples and their corresponding noise levels.
- **Layers:**
  - Convolutional/Linear layers: Capture spatio-temporal correlations in the input trajectory.
- **Output:**
  - Estimated score (gradient of the log probability density function).

```

class ScoreNetwork(nn.Module):

    def __init__(self, input_dim):

        super().__init__()

        self.fc1 = nn.Linear(input_dim, 256)

        self.fc2 = nn.Linear(256, 128)

        self.fc3 = nn.Linear(128, input_dim)

    def forward(self, x, t):

        h = torch.cat([x, t], dim=-1) # Concatenate trajectory and noise level

        h = F.relu(self.fc1(h))

        h = F.relu(self.fc2(h))

        return self.fc3(h)

```

DiffusionPolicy:

- Inputs:
  - Initial noisy trajectory.
  - Noise schedule (defines how noise is added/removed across steps).
- Process:
  - Iteratively denoise the trajectory using the ScoreNetwork.
- Output:
  - Refined trajectory.

```

class DiffusionPolicy(nn.Module):

    def __init__(self, score_network, num_steps):

        super().__init__()

        self.score_network = score_network

        self.num_steps = num_steps

    def forward(self, x_noisy, t):

        x = x_noisy

        for i in range(self.num_steps):

            score = self.score_network(x, t)

            x = x - score # Gradient descent on the noise

        return x

```

## Training and Evaluation

- Shared Components:
  - Both policies use the AdamW optimizer with a learning rate scheduler.
  - Training pipelines log loss metrics, trajectory smoothness, and success rates.
- Task-Specific Fine-Tuning:
  - Diffusion Policies are tuned for high-precision tasks (e.g., object orientation).
  - ACT focuses on tasks requiring temporal coherence (e.g., sliding objects).



By combining the ACT and Diffusion Policy implementations, the framework achieves robust performance across diverse tasks with varying levels of complexity.