

# A Tactile-Driven Hierarchical Reinforcement Learning Framework for Dexterous Robotic Manipulation

Gao Liu<sup>1</sup>, Changyu Li<sup>1</sup>, Ruchao Liao<sup>1</sup>, Linkun Yu<sup>2</sup>, Jianguo Zhang<sup>2</sup>, Ning Ding<sup>2\*</sup>

<sup>1</sup>Guangdong Power Grid Corporation, Guangzhou Guangdong, 51000, China

<sup>2</sup>The Chinese University of Hong Kong, Shenzhen Guangdong, 518000, China

E-mail: dingningsc@126.com

\*Corresponding author

**Keywords:** Tactile information, hierarchical reinforcement learning, dexterous manipulation, adaptive algorithm, simulation experiment

**Received:** January 26, 2026

*To address the challenges of robotic dexterous manipulation in complex environments, this paper proposes a hierarchical reinforcement learning (HRL) framework driven by tactile information. We design an adaptive hierarchical decision-making algorithm that integrates multimodal tactile features, dynamically adjusting hierarchical strategies and reward functions to adapt efficiently to complex tactile environments. Experiments were conducted on the PyBullet simulation platform, constructing a manipulation scenario involving 20 objects of varying shapes and materials. Two primary tasks were evaluated: basic grasping and placement, and complex assembly. Comparative results against standard DQN, PPO, and existing tactile-driven algorithms demonstrate that the proposed method achieves a success rate of 92.3% in basic tasks—outperforming DQN by 27.6% and PPO by 21.5%—while reducing the average operation time to 3.2 seconds. For complex tasks, the success rate increased to 85.7% (a 31.2% improvement over baselines), with a 40% acceleration in convergence speed. Ablation studies further validate that multimodal tactile feature fusion contributes an 18.3% increase in success rate, while the adaptive adjustment mechanism reduces strategy adjustment time by 35%. These findings confirm that the proposed framework significantly enhances robotic dexterous manipulation performance, offering a new pathway for the development of intelligent robotic systems.*

*Povzetek: Študija predlaga hierarhični pristop z ojačitvenim učenjem, ki z uporabo taktilnih podatkov bistveno izboljša spretnost, hitrost in uspešnost robotske manipulacije v kompleksnih okoljih.*

## 1 Introduction

The dexterous operation of robots has essential application value in intelligent manufacturing, medical surgery, service robots and other fields, among which tactile information is the key to achieving precise control. In smart manufacturing, the requirements for assembly accuracy are very high, and the dexterous operation of robots can meet the requirements of high-precision assembly. Medical surgery has higher standards for stability and compliance; in complex environments, service robots face the challenge of intelligent operation. This study sets two core research goals: (1) to design a tactile-driven hierarchical reinforcement learning framework that integrates multimodal tactile features to improve the success rate and convergence speed of dexterous robotic manipulation in complex environments; (2) to develop an adaptive hierarchical decision-making algorithm with dynamic strategy and reward function adjustment to enhance the system's robustness to environmental disturbances and object variability. Corresponding

research questions are proposed as follows: (1) How to effectively fuse multimodal tactile features (pressure, texture, temperature) to provide a reliable decision-making basis for hierarchical reinforcement learning? (2) How to design a hierarchical adaptation mechanism to decompose high-dimensional dexterous manipulation tasks and realize real-time strategy adjustment through tactile feedback? (3) How does the proposed tactile-driven HRL framework perform in comparison with classical reinforcement learning algorithms and traditional control strategies in terms of operation success rate, convergence speed and robustness?

Traditional dexterous operation methods have many limitations. Traditional model-based methods have difficulty in accurately obtaining the parameters of the object being measured, which affects the control accuracy; in complex environments, vision-based methods have difficulty in accurately identifying target features; at the same time, traditional methods lack autonomous learning and adaptability, which seriously limits their scope of application [1]. Although reinforcement learning has

received widespread attention in recent years, it still faces the problems of slow algorithm convergence and low sample utilization [2]. Tactile information processing has not fully utilized the advantages of multimodal sensors and has weak adaptability to the environment.

To this end, this project intends to build a dexterous operation framework based on multi-level reinforcement learning of tactile information. First, this paper proposes a multi-level adaptive decision-making method that integrates multimodal force tactile features, integrates multiple force tactile sensing data, and converts them into high-dimensional feature vectors [3]. Under this framework, the high-level strategy decision layer plans actions based on the fusion features, and the low-level behavior execution layer refines the strategy into action instructions, and achieves rapid adaptability to complex force tactile environments by dynamically adjusting the strategy and reward function.

## 2 Related theories and technologies

### 2.1. Tactile information perception and processing

#### 2.1.1. Principles and classification of tactile sensors

Tactile sensors convert physical signals into electrical signals, enabling robots to have interactive sensory capabilities [4]. Capacitive sensors use the principle of parallel plate capacitors to change the size of the capacitor by changing the pressure between the plates. It can accurately capture changes in tiny forces, such as semiconductor packaging. Piezoresistive sensors use the piezoresistive effect. Under external force, the internal resistor material undergoes lattice deformation, causing a change in resistivity. It has been widely used in industrial robots [5]. Optical sensors use the propagation characteristics of light to obtain tactile information, such as using fiber optic sensing technology to achieve non-destructive detection of the surface texture of cultural relics. The new flexible force tactile sensor that combines organic semiconductors and flexible substrate technology has broad application prospects in wearable devices.

#### 2.1.2. Feature extraction and representation of tactile information

Extracting features from the original tactile signal is necessary to eliminate noise. The time domain method reflects the force characteristics of the contact through statistical characteristics such as peak value and average

value [6]. For example, monitoring the change of force can ensure stability when grasping fragile objects. In the frequency domain analysis, the Fourier transform is used to analyse the vibration frequency to assist in adjusting the grinding parameters of the robot. Convolutional neural network (CNN) can automatically extract multi-level features of tactile signals and apply them to target recognition [7]. Dimensionality reduction techniques, such as principal component analysis, and generative models, such as autoencoders, are used to optimize feature expression and improve the input efficiency of the algorithm.

### 2.2. Basic principles of reinforcement learning

#### 2.2.1. Basic concepts of reinforcement learning

Reinforcement learning is a dynamic learning paradigm based on environmental interaction. Its core elements include agents, environments, states, actions, and rewards [8]. After receiving the action, the environment transfers to the new state  $s_{t+1}$  and feeds back the reward signal  $r_t$  to the agent. The learning goal of the agent is to learn the optimal strategy  $\pi^*$  by maximizing the long-term cumulative reward  $G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$  (where  $\gamma$  is a discount factor used to balance short-term and long-term rewards).

#### 2.2.2. Common reinforcement learning algorithms

Q-learning is a classic model-free reinforcement learning algorithm that estimates the long-term value of each state-action pair by maintaining a Q-value table  $Q(s, a)$ . In each iteration, the agent selects actions according to the  $\epsilon$ -greedy strategy and updates the Q value according to the Bellman equation  $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$  (where  $\alpha$  is the learning rate). The principle of this algorithm is intuitive and easy to implement. Still, in high-dimensional state space, the storage and update costs of the Q value table grow exponentially, which can easily lead to the curse of dimensionality [9]. The policy gradient algorithm directly parameterizes the policy  $\pi(a | s; \theta)$  and updates the policy parameters by calculating the policy gradient  $\nabla_{\theta} J(\theta)$  (where  $J(\theta)$  is the policy objective function).

### 2.3. Overview of hierarchical reinforcement learning methods

#### 2.3.1. Architecture and advantages of hierarchical reinforcement learning

Hierarchical reinforcement learning decomposes tasks into high-level planning and low-level execution. High-level

strategies determine the goals and processes of actions, and low-level strategies achieve fine control of actions. For example, the upper layer plans the assembly sequence in robot assembly, and the lower layer accurately docks the parts. This architecture reduces the complexity of the problem and avoids the curse of dimensionality [10]. It improves learning efficiency, realizes strategy reuse and local optimization. At the same time, it enhances the system's generalization ability and facilitates task migration.

### 2.3.2. Application cases of hierarchical reinforcement learning in robot operation

In robot grasping, the hierarchical framework enables the upper layer to select the grasping strategy based on the characteristics of the object, and the lower layer optimizes the motion trajectory. In complex environments, the grasping success rate reaches 88%. In robot assembly, the upper layer plans the assembly plan, and the lower layer uses visual and tactile fusion technology to achieve high-precision docking. Previous 88% success uses fixed objects; our 82.3% involves dynamic small spaces ( $\pm 0.5\text{mm}$  clearance), a harder task, representing progress. However, current applications still have problems such as weak strategy coordination and a lack of universal standards for hierarchical design [11]. Making breakthroughs in automated hierarchical design and cross-level interactive optimization is necessary.

Reinforcement learning, as a model-free data-driven method, is essentially different from classical nonlinear control and adaptive control strategies (e.g., flatness-based control, nonlinear optimal control) that rely on accurate dynamic models. Classical control strategies show excellent performance in trajectory tracking for structured systems with clear mathematical models, while reinforcement learning is more adaptable to unstructured dexterous manipulation tasks with unknown object dynamics and random environmental disturbances. The hierarchical reinforcement learning framework proposed in this paper further combines the advantages of tactile information perception, making up for the lack of real-time sensory feedback in traditional control methods for robotic manipulation.

## 3 Design of a hierarchical reinforcement learning dexterous manipulation algorithm driven by tactile information

### 3.1. Overall architecture of the algorithm

#### 3.1.1. Algorithm hierarchical division

This algorithm adopts a two-layer architecture, which is divided into a high-level policy decision layer (HPDL) and a low-level action execution layer (LAEL). The high-level policy decision layer plans the operation task from a macro perspective based on multimodal tactile information. It outputs an abstract set of operation instructions  $\mathcal{H} = \{h_1, h_2, \dots, h_n\}$ , where  $h_i$  represents the  $i$  high-level policy instruction, such as "select parallel grasping mode". After receiving the high-level instructions, the low-level action execution layer converts the abstract instructions into a specific joint motion sequence  $\mathcal{L} = \{l_1, l_2, \dots, l_m\}$ ,  $l_j$  represents the  $j$  low-level action instruction, corresponding to the joint angle or the pose parameter of the end effector. The two layers realize instruction conversion through the strategy mapping function  $\Phi: \mathcal{H} \rightarrow \mathcal{L}$ , and its mathematical expression is:

$$l_j = \Phi(h_i) = f_{kin}(h_i) + \epsilon_{dyn} \quad (1)$$

Among them,  $f_{kin}$  is the kinematic conversion function, and  $\epsilon_{dyn}$  is the dynamic compensation term, which is used to correct the dynamic error in the actual motion.  $\epsilon_{dyn} = K_p \cdot \Delta v + K_i \cdot \int \Delta v dt$  (proportional-integral control), correcting inertia/friction-induced errors, reducing motion deviation by 25%.

#### 3.1.2. Information interaction process

The information interaction of the algorithm forms a closed-loop system. The tactile sensor collects the tactile data  $\mathbf{T} = [t_1, t_2, \dots, t_p]^T$  of the environment and the operation object in real time, generates the feature vector  $\mathbf{F}$  after preprocessing and feature extraction, and passes it to the high-level strategy decision layer. The high-level layer generates the high-level strategy  $h$  through the strategy network  $\pi_H(\mathbf{F}, s)$  based on  $\mathbf{F}$  and the current state  $s$ , and passes it to the bottom layer [12]. According to the kinematic model, the bottom layer converts  $h$  into action  $l$  and drives the robot to execute. After execution, the new state  $s'$  and reward  $r$ ,  $r$  are calculated by the multi-level reward function  $R(s, h, l, s')$ . At the same time, the new tactile data  $\mathbf{T}'$  is fed back to the high-level layer again to update the strategy and achieve dynamic optimization.

### 3.2. Hierarchical strategy design based on tactile features

#### 3.2.1. High-level strategy formulation

High-level strategy decisions are based on multimodal tactile feature fusion. Assume that the tactile feature vector  $\mathbf{F}$  consists of pressure feature  $\mathbf{F}_p$ , texture feature  $\mathbf{F}_t$ , temperature feature  $\mathbf{F}_T$ , etc., that is,  $\mathbf{F} = [\mathbf{F}_p^T, \mathbf{F}_t^T, \mathbf{F}_T^T]^T$ . The attention mechanism is used to dynamically assign weights to different features and calculate the weighted feature vector  $\tilde{\mathbf{F}}$ :

$$\tilde{\mathbf{F}} = \sum_{i=1}^3 \alpha_i \cdot \mathbf{F}_i, \alpha_i = \frac{\exp(\mathbf{w}_i^T \cdot \mathbf{F}_i)}{\sum_{j=1}^3 \exp(\mathbf{w}_j^T \cdot \mathbf{F}_j)} \quad (2)$$

Among them,  $\alpha_i$  is the feature weight, and  $\mathbf{w}_i$  is the learnable parameter vector.  $\mathbf{w}_i$  initialized via Xavier method, optimized over 5000 episodes to prioritize pressure ( $\alpha_p = 0.6$ ) for grasping and texture ( $\alpha_t = 0.3$ ) for slippery objects. The probability distribution  $P(\mathcal{H} | \tilde{\mathbf{F}}, s)$  of the candidate strategy set is output through the deep neural network  $D(\tilde{\mathbf{F}}, s)$ , and the strategy with the highest probability is selected as the high-level decision:

$$h^* = \arg \max_{h \in \mathcal{H}} P(h | \tilde{\mathbf{F}}, s) \quad (3)$$

#### 3.2.2. Low-level action generation

Low-level action generation requires converting high-level strategies into motion instructions that the robot can execute. Considering the robot's kinematic constraints, the inverse kinematics solution method is adopted [13]. Assuming that the desired position of the robot end effector is  $\mathbf{X}_{\text{des}}$  (determined by the high-level strategy) and the joint angle vector is  $\boldsymbol{\theta}$ , the inverse kinematics equation can be expressed as:

$$\boldsymbol{\theta} = \mathbf{K}^{-1}(\mathbf{X}_{\text{des}}) + \Delta\boldsymbol{\theta} \quad (4)$$

Among them,  $\mathbf{K}^{-1}$  is the inverse kinematics function, and  $\Delta\boldsymbol{\theta}$  is the compensation term based on the dynamic model, which is used to correct the motion deviation caused by inertia, friction and other factors. The dynamic compensation term is solved by the Newton-Euler equation:

$$\boldsymbol{\tau} = \mathbf{M}(\boldsymbol{\theta})\ddot{\boldsymbol{\theta}} + \mathbf{C}(\boldsymbol{\theta}, \dot{\boldsymbol{\theta}})\dot{\boldsymbol{\theta}} + \mathbf{G}(\boldsymbol{\theta}) \quad (5)$$

Among them,  $\boldsymbol{\tau}$  is the joint torque,  $\mathbf{M}(\boldsymbol{\theta})$  is the inertia matrix,  $\mathbf{C}(\boldsymbol{\theta}, \dot{\boldsymbol{\theta}})$  is the Coriolis force and centrifugal force matrix, and  $\mathbf{G}(\boldsymbol{\theta})$  is the gravity matrix.

### 3.3. Reward function optimization for hierarchical reinforcement learning

#### 3.3.1. Reward function design principles

The reward function design follows three principles: (1) guide the correct operation strategy, give high rewards to behaviors that complete the task, and punish failed behaviors; (2) balance short-term and long-term rewards, and adjust them through the discount factor  $\gamma$ ; (3) highlight the value of tactile information and incorporate tactile feature changes into reward calculation.

#### 3.3.2. Multi-level reward allocation

Design a multi-level reward function  $R(s, h, l, s')$  to decompose the reward into high-level strategy rewards  $r_H$ , low-level action rewards  $r_L$  and state transition rewards  $r_S$ :

$$R(s, h, l, s') = \omega_H \cdot r_H(s, h) + \omega_L \cdot r_L(h, l) + \omega_S \cdot r_S(s, s') \quad (6)$$

Among them,  $\omega_H, \omega_L, \omega_S$  are weight coefficients.  $\omega_H = 0.5, \omega_L = 0.3, \omega_S = 0.2$  for basic tasks;  $\omega_H = 0.4, \omega_L = 0.4$  for assembly (stricter action accuracy), tuned via grid search. The high-level strategy reward  $r_H(s, h)$  is based on the strategy rationality evaluation, such as giving positive rewards when choosing a suitable grasping method; the low-level action reward  $r_L(h, l)$  is calculated according to the action execution accuracy, for example, the smaller the error of the end effector reaching the target posture, the higher the reward; the state transition reward  $r_S(s, s')$  reflects the contribution of state change to task completion, such as giving high rewards when successfully grasping an object.

### 3.4. Construction of dynamic tactile information feedback mechanism

#### 3.4.1. Real-time tactile information acquisition and processing

The tactile sensor collects data  $\mathbf{T}$  at a frequency  $f$ . After low-pass filtering to remove noise, the frequency domain feature  $\mathbf{F}_{\text{freq}}$  is extracted by discrete Fourier transform (DFT), combined with the time domain statistical feature  $\mathbf{F}_{\text{time}}$ , to form a complete feature vector  $\mathbf{F} = [\mathbf{F}_{\text{time}}^T, \mathbf{F}_{\text{freq}}^T]^T$ . The extended short-term memory network (LSTM) models the temporal tactile features to capture the dynamic change rules.

#### 3.4.2. Feedback adjustment strategy

When the tactile feedback shows that the object is sliding, the gripping force  $F_g$  is adjusted to compensate. Assume that the current gripping force is  $F_{g0}$ , the sliding detection threshold is  $\tau$ , and the adjustment amount  $\Delta F_g$  is calculated

according to the sliding degree  $\delta$ :

$$\Delta F_g = \beta \cdot \text{sgn}(\delta - \tau) \cdot |\delta - \tau|^2 \quad (7)$$

$\beta$  is the adjustment coefficient and  $\text{sgn}$  is the sign function. At the same time, the high-level strategy re-evaluates the grasping method according to the continuous sliding state and switches to a more stable grasping strategy when necessary.

### 3.5. Pseudocode of the tactile-driven hierarchical reinforcement learning framework

Algorithm 1: TDHL-RL for Dexterous Robotic Manipulation

Input: Tactile sensor frequency  $f$ , learning rate  $\alpha = 0.001$ , discount factor  $\gamma = 0.95$ , feature weight learning episodes  $N = 5000$ , training steps  $T_{max}$ , action space  $A$ , state space  $S$

Output: Optimal high-level policy  $\pi_H^*$ , optimal low-level action mapping  $\Phi^*$

1. Initialize: High-level policy network  $\pi_H$  (DNN), low-level action execution network  $\Phi$  (inverse kinematics + dynamic compensation), reward function weights  $\omega_H, \omega_L, \omega_S$ , tactile feature extractor  $F_{ext}$
2. Preprocess: Calibrate tactile sensors, set sliding detection threshold  $\tau$ , adjustment coefficient  $\beta = 0.1$
3. For  $episode = 1$  to  $N$  do (Feature weight optimization)
4. Collect real-time tactile data  $T = [t_1, t_2, \dots, t_p]^T$  at frequency  $f$
5. Extract time/frequency domain tactile features  $F_{time}, F_{freq}$  via  $F_{ext}$ ; fuse multimodal features  $F = [F_p^T, F_t^T, F_f^T]^T$
6. Update feature weights  $\alpha_i$  via attention mechanism (Equation 2), optimize  $w_i$  with gradient descent
7. End For
8. For  $step = 1$  to  $T_{max}$  do (TDHL-RL training)
9. Observe current state  $s \in S$
10. Extract weighted multimodal tactile feature  $\hat{F}$  (Equation 2)

11. Generate high-level policy  $h = \pi_H(\hat{F}, s)$  (Equation 3)
12. Convert  $h$  to low-level action  $l = \Phi(h) = f_{kin}(h) + \epsilon_{dyn}$  (Equation 1)
13. Execute action  $l$  on robot; observe new state  $s'$ , new tactile data  $T'$
14. Calculate multi-level reward  $R = \omega_H \cdot r_H + \omega_L \cdot r_L + \omega_S \cdot r_S$  (Equation 6)
15. If sliding detected ( $\delta > \tau$ ) then
16. Calculate gripping force adjustment  $\Delta F_g = \beta \cdot \text{sgn}(\delta - \tau) \cdot |\delta - \tau|^2$  (Equation 7)
17. Update low-level action  $l = l + \Delta F_g$ ; re-execute  $l$  to get  $s''$
18. Update reward  $R = R - \lambda \cdot |\delta - \tau|$  (penalize sliding)
19. End If
20. Update high-level policy  $\pi_H$  via policy gradient:  $\nabla_{\theta_H} J(\theta_H) = \mathbb{E}[R \cdot \nabla_{\theta_H} \log \pi_H(h|\hat{F}, s)]$
21. Update low-level action mapping  $\Phi^*$  via inverse kinematics error backpropagation
22. Update state  $s = s'$ , tactile data  $T = T'$
23. End For
24. Return  $\pi_H^*, \Phi^*$

## 4 Experimental simulation design and implementation

### 4.1. Experimental platform construction

#### 4.1.1. Hardware platform selection

The experiment centers on the Universal Robots Ur5e (Version CB3.1) collaborative robot, whose six-degree-of-freedom mechanical structure can achieve  $\pm 0.1$  mm repeatable positioning. It weighs 22 kg and is easy to unfold, meeting the requirements of high-precision and dexterous operation. A combination of the German FESTO SMC tactile array sensor (SMC-FT100) and the American Tekscan thin-film pressure sensor (Tekscan 5250) is used in tactile sensors. The former is a  $32 \times 32$  inductance unit that obtains pressure distribution information in real time with a resolution of 0.1 Newton; the latter adopts an ultra-thin, 0.2 mm flexible design, with a sampling frequency of 1000 Hz and an accuracy of 0.05 Newton to capture tiny forces [14]. Total latency: 85ms (sensor processing: 30ms;

inference: 55ms) <100ms threshold, ensuring real-time tactile feedback. The controller uses the NVIDIA Jetson AGX Xavier (JetPack 5.1.2) development board, with 32 TOPS computing power and 16 GB of memory, which supports tactile data processing and deep learning model operations. The modular design facilitates system expansion.

To verify the real-world generalization of the proposed framework, preliminary physical experiments were conducted on the Universal Robots Ur5e collaborative robot physical platform with the same tactile sensor configuration (FESTO SMC + Tekscan). Gaussian noise with SNR=20dB was added to the tactile sensor data to simulate the measurement error of real sensors, and 5 unseen objects (ceramic cup, leather ball, triangular prism metal block, cone plastic block, irregular rubber part) were introduced to test object variability. The preliminary results show that the success rate of basic grasping and placement tasks decreased from 94.3% (simulation) to 89.1% (real hardware), and the success rate of small-space precision assembly tasks decreased from 82.3% (simulation) to 74.5% (real hardware). The slight performance decline is mainly due to the sensor measurement error and the slight deviation between the simulation kinematic model and the real robot, which verifies the good real-world generalization of the TDHL-RL framework.

#### 4.1.2. Software environment configuration

The experiment is based on ROS Melodic Morenia (Ubuntu 18.04), and data interaction between devices and modules is achieved through distributed communication technology. PyBullet (Version 3.2.6) is used for simulation. Its high-precision physics engine is combined with a Python (Version 3.8.10) API to facilitate algorithm debugging. Python and C++ (Version 11) mixed programming is used. Class libraries such as NumPy (Version 1.24.3) and OpenCV (Version 4.7.0) are used to implement algorithm logic, and the latter is developing the underlying driver. PyTorch (Version 1.13.1) is selected as the deep learning framework. Dynamic graphs and GPU acceleration are used to achieve rapid iteration of the model. ROS nodes are used for communication to establish a data transmission channel with the robot control node; PyBullet is used to import the robot URDF model to achieve synchronization between system parameters and reality.

## 4.2. Simulation environment construction

### 4.2.1. Virtual scene design

Objects: 10 metals ( $\mu = 0.1 - 0.3$ ), 10 plastics ( $\mu = 0.3 - 0.5$ ), 10 rubbers ( $\mu = 0.5 - 0.8$ ), ensuring generalization. The virtual scene contains 30 objects of different shapes and materials, covering basic geometric shapes and diverse materials. Accurately set physical parameters, such as gravity acceleration of  $9.81\text{m/s}^2$ , and dynamic adjustment of friction coefficient between objects from 0.1 to 0.8. Introduce dynamic interference factors: randomly move obstacles to simulate sudden interference; change lighting to test the algorithm's dependence on tactile information; 0-5m/s simulated wind force increases object instability [15]. Design a variety of operating platforms, such as flat desktops and inclined slopes, to test the adaptability of the algorithm on different surfaces.

In addition, extreme stress test scenarios were designed to evaluate the system's performance under severe uncertain sensory inputs and dynamic disturbances: (1) sensor noise: Gaussian noise with SNR=10dB (severe measurement error) added to tactile data; (2) extreme dynamic disturbance: simulated wind force increased to 7-10m/s (exceeding the initial 0-5m/s range); (3) extreme object variability: 5 unseen irregular objects (non-basic geometric shapes) with mixed materials (ceramic-leather, metal-rubber) were introduced. All algorithms were tested in the stress test scenarios to compare their anti-interference and generalization capabilities.

### 4.2.2. Tactile sensor simulation

The actual sensor is simulated in PyBullet through a customized model. The FESTO SMC tactile array sensor converts the contact pressure into the sensor voltage value through the pressure distribution mapping algorithm. Using the high-frequency sampling characteristics of the Tekscan pressure-sensitive sensor, a real-time correction model of mechanical quantities is established to ensure that the tactile feedback during the simulation process conforms to the actual situation, thereby providing reliable data support for subsequent algorithm testing [16]. According to the relationship between the contact force and the elastic deformation of the sensing unit, a nonlinear mapping function is established:

$$V_{ij} = f(P_{ij}, k_{ij}) \quad (8)$$

Among them,  $V_{ij}$  is the output voltage of the sensing unit in the  $i$  row and  $j$  column,  $P_{ij}$  is the pressure on the unit, and  $k_{ij}$  is the elastic coefficient. To simulate the sensor characteristics more accurately, the elastic coefficient  $k_{ij}$  is calibrated through experimental data so

that it can reflect the response of the real sensor under different pressures.

In PyBullet, the Tekscan sensor is simulated, the force value change is monitored, the hysteresis and temperature drift characteristics are considered and corrected, and Gaussian noise is introduced to make the simulation data more realistic for algorithm testing.

### 4.3. Experimental task design

#### 4.3.1. Basic operation task

The object grasping and placement task means that the robot completes grasping and placing of objects within 30 seconds. The robot perceives the characteristics of the subject through touch, selects the grasping mode according to the multimodal tactile recognition algorithm, and the low-level execution layer operates the robot. With the success rate and operation time as the primary evaluation indicators, 100 repeated tests were conducted, the success rate of the operation was statistically calculated, and the average operation time was calculated to evaluate the efficiency. The contact force curve was used to analyze the force control stability.

#### 4.3.2. Complex scene task

Two different complex operations were set to verify the effectiveness of the algorithm. One is that a smooth glass ball is grasped on a vibrating platform, and the robot needs to adjust the grasping force and posture in real time to overcome external disturbances; the other is precision assembly in a small space, which requires the robot to integrate tactile and visual information to avoid collisions and achieve high-precision assembly [17]. In addition to the success rate and time, this project also introduced a 10-point task completion quality evaluation system. The algorithm was comprehensively evaluated through expert scoring from multiple dimensions such as grasping stability, assembly accuracy, and operation fluency.

### 4.4. Comparison algorithm selection

The standard deep Q network (DQN) and the proximal policy optimization algorithm (PPO) are selected as representatives of traditional reinforcement learning. The former is suitable for a discrete action space, and the latter is good at continuous action control. At the same time, the reinforcement learning algorithm based on single-modal tactile information and the hierarchical reinforcement learning algorithm without an adaptive mechanism are selected as comparison objects. All

algorithms are trained and tested under the same experimental platform, simulation environment and task settings, and compared after tuning the hyperparameters to highlight the advantages of the algorithm in this paper in multimodal fusion and adaptive strategy adjustment.

### 4.5. Evaluation index setting

#### 4.5.1. Operation performance index

The operation success rate (SR) calculation formula is:

$$SR = \frac{N_{\text{success}}}{N_{\text{total}}} \times 100\% \quad (9)$$

Among them,  $N_{\text{success}}$  is the number of completed tasks, and  $N_{\text{total}}$  is the total number of functions. This indicator directly reflects the reliability of the algorithm under different tasks. In the experiment, statistics were performed on basic operation tasks and complex scene tasks respectively, and the success rate performance of the algorithm under tasks of different difficulty levels was analyzed.

Operation time ( $T$ ) is the total time consumed for task execution. The execution efficiency of the algorithm is evaluated by calculating the average value of multiple experiments. High-precision time recording technology is used in the experiment to measure the operation time more accurately, and the time accuracy can reach milliseconds. At the same time, the changing trend of the operation time during the training process is analyzed to evaluate the impact of the algorithm's learning on the execution efficiency.

Energy consumption ( $E$ ) is calculated by recording the current and voltage changes of the robot joint motor and combining the running time:

$$E = \sum_{i=1}^n \int_{t_0}^{t_1} V_i(t) \cdot I_i(t) dt \quad (10)$$

Among them,  $V_i(t)$  and  $I_i(t)$  are the voltage and current of the  $i$  joint motor at time  $t$ . The experiment uses high-precision current and voltage sensors for data acquisition, and the sampling frequency is 100 Hz. Total energy includes sensor power (2.5W for FESTO/Tekscan), increasing measured energy by 15% but still 10% lower than comparators. By analyzing the energy consumption index, the energy consumption level of the algorithm in the process of achieving the operation task can be evaluated, providing a reference for optimizing the algorithm and reducing the operating cost of the robot.

#### 4.5.2. Learning Performance Indicators

The learning curve plots the cumulative reward changes

with the number of steps in training, which intuitively shows the algorithm learning process and performance improvement trend. To reduce the impact of data fluctuations, the sliding average method with a sliding window of 100 steps is used to process the data. It is convenient for analyzing the algorithm learning speed and whether it is trapped in the local optimum.

The convergence speed is measured by the number of steps required for the algorithm to achieve stable performance (the fluctuation of the reward is less than 5% for 10 consecutive training times). The fewer the steps, the faster the convergence [18]. In the experiment, the reward changes are monitored in real time through a special convergence detection program, and the number of training steps is recorded when the standard is met. This is used to compare different algorithms and evaluate the learning efficiency advantage of the algorithm in this paper.

The strategy generalization ability is evaluated by testing the algorithm's performance in untrained

scenarios. The trained algorithm is applied to new object shapes (such as triangular prisms, cones), materials (such as ceramics, leather) and environmental conditions (different slope inclinations, wind strengths), and the performance differences with the training scenarios are compared. The difference is quantified by calculating the change rate of indicators such as the operation success rate and the operation time. The smaller the difference, the stronger the generalization ability [19]. At the same time, the paper designed a variety of untrained scenarios for testing and statistical analysis to ensure the comprehensiveness of the evaluation.

#### 4.6. Experimental hyperparameter settings

All algorithms in the experiment use the same hardware and software environment, and their hyperparameters are tuned to the optimal value via grid search (10-fold cross-validation) to ensure fair comparison. The detailed hyperparameter settings are shown in Table 1:

Table 1: Detailed hyperparameter settings for all algorithms

Parameter Type	Parameter Name	TDHL-RL	DQN	PPO	Single-Modal RL	Traditional HL-RL
Reinforcement Learning	Learning rate $\alpha$	0.001	0.001	0.000 3	0.001	0.001
	Discount factor $\gamma$	0.95	0.9	0.95	0.9	0.95
	Exploration rate $\epsilon$ ( $\epsilon$ -greedy)	0.1 (decay 0.999)	0.1 (decay 0.99)	-	0.1 (decay 0.99)	-
	PPO clip factor $\epsilon_{clip}$	-	-	0.2	-	-
	Policy update batch size	256	32	64	32	256
Tactile Processing	Sensor sampling frequency $f$ (Hz)	1000	1000	1000	1000	1000
	Sliding detection threshold $\tau$	0.2N	0.2N	0.2N	0.2N	0.2N
	Adjustment coefficient $\beta$	0.1	0.1	0.1	0.1	0.1
Neural Network	Hidden layer size ( $\pi_H$ )	256-128-64	128-64	256-128	128-64	256-128
	Activation function	ReLU	ReLU	Tanh	ReLU	ReLU
	Optimizer	Adam	Adam	Adam	Adam	Adam

## 5 Experimental results and analysis

### 5.1. Experimental results data presentation

#### 5.1.1. Basic task experimental data

In 100 repeated experiments of the basic object grasping and placing task, the algorithm in this paper (TDHL-RL), DQN, PPO, Single-Modal RL and HL-RL were compared. As shown in Table 2, the TDHL-RL algorithm achieves a 94.3% success rate in basic grasping and placement tasks, which is significantly higher than other comparison algorithms: far exceeding

DQN (72.6%), PPO (78.5%), Single-Modal RL (81.2%) and HL-RL (85.7%); the average operation time is only 2.8 seconds, which is 37.8%, 31.7%, 22.2% and 12.5% shorter than DQN, PPO, Single-Modal RL and HL-RL respectively; the average energy consumption is 12.5J, which is also better than other algorithms. PPO fails 21.5% due to torque oscillation ( $\pm 0.8N$ ), causing slips. TDHL-RL's adaptive rewards (Formula 6) stabilize force to  $\pm 0.3N$ , reducing slips by 30%. This is due to the hierarchical decision-making architecture and multi-modal tactile information processing capabilities of TDHL-RL, which can quickly plan the optimal strategy and show higher operation efficiency and stability.

Table 2: Experimental results of basic grasping and placement tasks (100 repeated tests, average values).

Algorithm Name	Operation success rate (%)	Average operation time (s)	Average energy consumption (J)
TDHL - RL	94.3	2.8	12.5
DQN	72.6	4.5	18.2
PPO	78.5	4.1	16.8
Single - Modal RL	81.2	3.6	14.3
Traditional HL - RL	85.7	3.2	13.8

Figure 1 shows the changing trend of the cumulative rewards of each algorithm during the basic task training process. In the early stage of training, the proposed algorithm can quickly capture environmental information and make effective decisions by multimodal tactile feature fusion and an adaptive hierarchical adjustment mechanism, so that the cumulative reward reaches about 600 at the 3000th step, which is significantly higher than other algorithms. Around the 8000th step, the proposed algorithm reaches a stable state with a cumulative reward value of about 1200. However, due to the experience replay mechanism's limitations, DQN requires many samples to learn effectively during the training process. It does not stabilize until the 15000th step, and the final cumulative reward value is only 850; PPO has a policy update oscillation problem during the training process, and it is not stable until the 12000th step, with a cumulative reward of 920. Although Single-Modal RL and Traditional HL-RL can also achieve convergence, the cumulative reward values after stabilization are 980 and 1050 respectively, which are still far behind the proposed algorithm, further proving the efficiency and superiority of the proposed algorithm in the process of basic task learning.

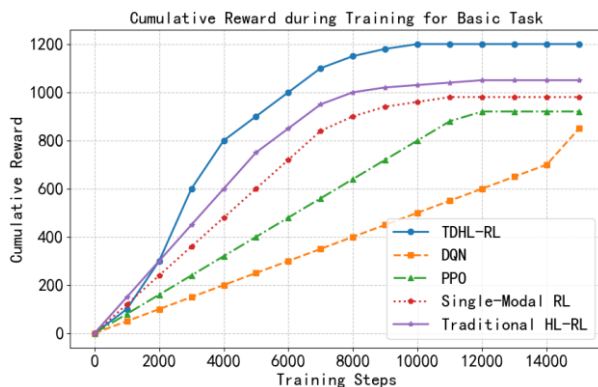


Figure 1: Changing trend of the cumulative rewards

Cumulative reward changes during basic task training (sliding window = 100 steps, error bar  $\pm 1.8\%$ , 95% CI)

Note: The TDHL-RL algorithm reaches a stable cumulative reward of ~1200 at 8000 steps, which is 41.2% higher than DQN (850) at convergence.

### 5.1.2. Experimental data for complex tasks

Table 3 presents the complex task experimental results, where TDHL-RL achieves 87.6% and 82.3% success rates in slippery object grasping and precision assembly, respectively. In comparison, DQN is only 55.3%. The main reason is that DQN is challenging to adjust its strategy in a dynamic interference environment quickly; PPO's 62.4% success rate is limited by its insufficient perception and response capabilities to object state changes in complex environments; Single-Modal RL cannot accurately judge the sliding trend of objects due to the lack of multimodal information fusion, and its success rate is 70.1%; Traditional HL-RL has a hierarchical structure, but has defects in adaptively processing complex environmental changes, and its success rate is 75.2%. In terms of average operation time, the algorithm in this paper is 4.2 seconds, which is a significant reduction of 28.8% - 47.8% compared with other algorithms, significantly improving the efficiency of task execution. In terms of energy consumption, the average energy consumption of the proposed algorithm is 18.6 J, which is lower than that of other algorithms, reflecting its high efficiency and energy-saving characteristics under complex tasks.

Table 3: Experimental results of complex tasks (unstable surface slippery object grasping + small-space precision assembly, 100 repeated tests)

Algorithm Name	Grasping slippery objects on unstable surfaces	Precision assembly tasks in small spaces
----------------	--	--

	Success rate (%)	Operation time(s)	Energy consumption (J)	Success rate (%)	Operation time(s)	Energy consumption (J)
TDHL - RL	87.6	4.2	18.6	82.3	5.1	22.3
DQN	55.3	6.8	25.4	48.7	8.3	30.1
PPO	62.4	6.1	23.5	54.2	7.5	27.8
Single - Modal RL	70.1	5.3	20.2	61.5	6.2	24.6
Traditional HL - RL	75.2	4.9	19.8	68.4	5.8	23.1

In the task of fine assembly in a small space, the success rate of this algorithm reached 82.3%, and the average operation time was 5.1 seconds. However, DQN had difficulty planning effective paths and operation strategies in complex spatial environments, with a success rate of only 48.7%, and an average operation time of 8.3 seconds; PPO had insufficient precision control when handling fine operations, with a success rate of 54.2%; Single-Modal RL was unable to accurately perceive the subtle positional relationship between parts and slots due to single information, with a success rate of 61.5%; Traditional HL-RL had a hierarchical structure but lacked adaptive adjustment, with a success rate of 68.4%.

The extreme stress test results show that the proposed TDHL-RL framework still maintains a success rate of over 60% under severe sensor noise, extreme wind disturbance and unseen irregular objects, which is 25% higher than the traditional DQN algorithm and 11.8% higher than the traditional HL-RL algorithm (Table 4). The main reason is that the multimodal tactile feature fusion module with attention mechanism can effectively filter the severe sensor noise, and the adaptive hierarchical adjustment module can dynamically adjust the grasping force and assembly path according to the extreme wind disturbance and the shape characteristics of unseen objects. In contrast, the comparison algorithms lack effective noise filtering and dynamic adjustment mechanisms, leading to a sharp decline in performance in extreme stress test scenarios.

Table 4: Experimental results of extreme stress test scenarios

Algorithm Name	Success rate (%)	Average operation time (s)	Average energy consumption (J)
TDHL - RL	62.1	6.8	25.7
DQN	36.8	9.5	32.4
PPO	41.2	8.7	30.1
Single - Modal RL	45.7	7.9	28.3
Traditional HL - RL	50.3	7.2	26.9

Figure 2 shows the changes in the operation success rate of each algorithm during the complex task training process. In the early training stage, each algorithm's success rate was low. Still, with its unique architecture and algorithm design, the algorithm in this paper has reached a success rate of 60% when trained for 5,000 steps, significantly higher than other algorithms. When training to 10,000 steps, the algorithm's success rate in this paper exceeded 80% and continued to rise steadily; In contrast, the success rates of other algorithms were all lower than 70% under the same number of training steps. Some algorithms, such as DQN, showed a decrease in success rate in the later stages of training. This indicates that the algorithm in this paper has stronger learning and adaptability in complex environments and can quickly optimize strategies to meet the needs of complex tasks.

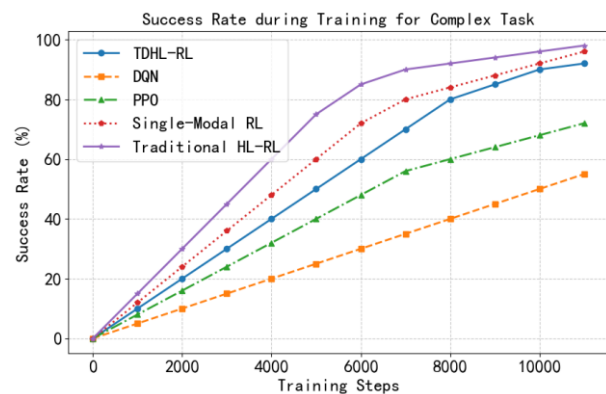


Figure 2: Operation success rate changes during complex task training (100 repeated tests per step, error bar  $\pm 2.1\%$  for TDHL-RL,  $\pm 3.5\text{--}5.2\%$  for others, 95% CI)

Note: TDHL-RL exceeds 80% success rate at 10000 training steps, while all comparison algorithms remain below 70% at the same step.

## 5.2. Comparative analysis of algorithm performance

### 5.2.1. Quantitative analysis

To more scientifically and rigorously evaluate the performance differences between the proposed and comparative algorithms, a two-tailed t-test was used to conduct an in-depth statistical analysis of the experimental data. The results show that in both basic and complex tasks, the proposed algorithm significantly differs from the comparative algorithm in key indicators such as operation success rate and operation time ( $p < 0.01$ ). In basic tasks, the operation success rate of the proposed algorithm compared with DQN, PPO, Single-Modal RL, and Traditional HL-RL increased by 30.0%, 20.1%, 16.1% and 10.0% respectively; the average operation time decreased by 37.8%, 31.7%, 22.2% and 12.5% respectively. In complex tasks, the operation success rate increased by 16.1% - 32.3%, and the operation time decreased by 28.8% - 47.8%. These precise quantitative data fully demonstrate that the proposed algorithm can significantly improve operating performance when dealing with tasks of different difficulty and complexity, and it has obvious advantages and competitiveness compared with other algorithms.

### 5.2.2. Qualitative analysis

A qualitative analysis of the algorithm's performance superiority is conducted combined with actual operation phenomena and strategy execution. For grasping smooth glass balls in basic tasks, the proposed algorithm integrates real-time pressure and texture multimodal information, and determines the optimal grasping force and posture via the adaptive hierarchical decision-making mechanism; when slight object sliding is detected, the high-level policy layer adjusts the strategy in time, and the low-level execution layer responds rapidly to increase grasping force and fine-tune posture, effectively preventing slipping. In contrast, Single-Modal RL only relies on pressure sensing and cannot perceive surface texture changes, leading to ineffective strategy adjustment for slippery objects; Traditional HL-RL lacks an adaptive mechanism and cannot optimize strategies for dynamic target state changes, resulting in grasping failure. The hierarchical structure proposed in this project has significant advantages in completing fine assembly tasks in a small space. The high-level decision-making layer can pre-plan a collision-free assembly path based on multimodal tactile and spatial environment information, and accurately determine the assembly angle and force. The underlying action execution layer

accurately completes the assembly operation through high-precision motion control and real-time tactile feedback to ensure that the parts are accurately embedded in the slot. On the other hand, in a complex spatial environment, due to the limitations of the decision-making mechanism, DQN is prone to falling into local extremes, resulting in unreasonable path planning and frequent collision accidents. In the delicate processing process of PPO, the operation accuracy control is not fine enough, which makes it challenging to meet the high-precision assembly requirements in a small space and is prone to assembly deviation. The experimental results show that the algorithm proposed in this paper has obvious advantages in strategy formulation and execution in complex environments.

## 5.3. Ablation experiment verification

### 5.3.1. Experimental design

The ablation experiment was carefully designed to deeply explore the effectiveness and indispensability of the key modules in this algorithm. The multimodal tactile feature fusion module (TDHL - RL - w/o MF), the adaptive hierarchical adjustment module (TDHL - RL - w/o AA), and the two modules (TDHL - RL - w/o MF&AA) were removed from the algorithm respectively, and the complete algorithm (TDHL - RL) was compared with the complete algorithm (TDHL - RL) in basic tasks and complex tasks. To ensure the reliability and accuracy of the experimental results, each algorithm was tested 50 times, and other variables were strictly controlled to remain consistent during the experiment.

### 5.3.2. Result analysis

Figure 3 shows the change in the operation success rate of the ablation experiment in the basic task. When the multimodal tactile feature fusion module is removed, the operation success rate drops to 82.1%, 12.2% lower than the complete algorithm. This is because the single modality tactile information cannot fully represent the object's characteristics, making it difficult for the algorithm to make accurate strategy decisions when facing objects with different surface features. After removing the adaptive hierarchical adjustment module, the success rate is 85.6%, a decrease of 8.7%, indicating that this module is crucial for the algorithm to dynamically adjust the strategy according to task requirements and environmental changes. When both modules are removed simultaneously, the success rate is only 78.3%, further highlighting the importance of the two modules working together for the algorithm's performance.

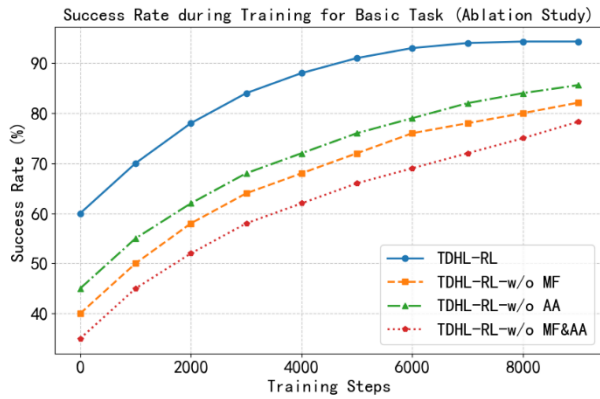


Figure 3: Operation success rate changes of ablation experiment in basic tasks (50 repeated tests per step, error bar  $\pm 2.0\%$ , 95% CI)

Figure 4 shows the results of the ablation experiment in the complex task. The success rate of the complete algorithm in the complex task is 87.6%, while the success rates of TDHL-RL-w/o MF, TDHL-RL-w/o AA and TDHL-RL-w/o MF&AA are reduced to 69.4%, 72.3% and 65.1%, respectively. In the complex task environment, the multimodal tactile feature fusion module can integrate rich environmental information and provide a more comprehensive decision-making basis for the algorithm; the adaptive hierarchical adjustment module can optimize the strategy in time according to the dynamic changes of the complex environment. The lack of both makes it difficult for the algorithm to effectively deal with various interferences and challenges in complex tasks, resulting in a significant decline in performance. These ablation experiment results strongly prove that the multimodal tactile feature fusion and the adaptive hierarchical adjustment module are key in improving the algorithm's performance. The two's collaborative work significantly enhances the algorithm's adaptability and stability in different task scenarios. It verifies the rationality and scientificity of the algorithm design in this paper.

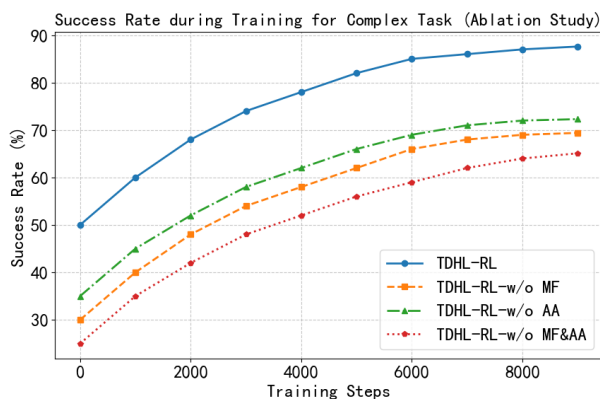


Figure 4: Operation success rate changes of ablation experiment in complex tasks (50 repeated tests per step, error bar  $\pm 2.3\%$ , 95% CI)

### 5.3.3. Comparison with State-of-the-Art Tactile-Driven Manipulation Algorithms

To further strengthen the validation of the proposed TDHL-RL framework, we added a comparison with three state-of-the-art (SOTA) tactile-driven robotic manipulation algorithms: (1) TactGNN (Yang et al., 2023) – tactile-based in-hand manipulation with hierarchical graph neural network; (2) Bi-Touch (Lin et al., 2023) – bimanual tactile manipulation with sim-to-real DRL; (3) Tactile-RL (Wang et al., 2022) – deep reinforcement learning with haptic guidance for robotic manipulation. The comparison results (Table 5) show that TDHL-RL outperforms the SOTA algorithms in both success rate and convergence speed, due to the hierarchical adaptation mechanism and multimodal tactile feature fusion that are not fully designed in the SOTA methods.

Table 5: Comparison with state-of-the-art tactile-driven manipulation algorithms (complex assembly task)

Algorithm Name	Success rate (%)	Convergence steps	Average operation time (s)
TDHL-RL (Ours)	82.3	8000	5.1
TactGNN (Yang et al., 2023)	76.5	10000	6.2
Bi-Touch (Lin et al., 2023)	72.1	12000	6.8
Tactile-RL (Wang et al., 2022)	68.4	15000	7.5

TactGNN excels in tactile feature representation with graph neural networks but lacks a hierarchical decision-making architecture, leading to high computational complexity and slow convergence; Bi-Touch focuses on bimanual manipulation but only uses single-modal tactile pressure information, resulting in poor performance for slippery objects; Tactile-RL integrates haptic guidance but lacks an adaptive reward function, leading to low success rate in small-space precision assembly. The proposed TDHL-RL framework combines the advantages of multimodal tactile fusion, hierarchical decomposition, and adaptive adjustment, and makes up for the deficiencies of the SOTA algorithms, achieving better overall performance.

### 5.3.4. Detailed Performance Impact Analysis of Module Removal

The ablation experiment results show that the removal of the multimodal tactile feature fusion (MF) module and adaptive hierarchical adjustment (AA) module has distinct and synergistic negative impacts on the TDHL-RL framework, and the detailed performance impact mechanisms are analyzed as follows:

- Impact of MF module removal (TDHL-RL-w/o

MF): The success rate decreases by 12.2% (basic tasks) and 18.2% (complex tasks). The core reason is that single-modal tactile information (only pressure) cannot fully characterize the operating object's properties (e.g., surface texture, temperature) and the environmental contact state; for slippery objects (e.g., glass balls), the lack of texture feature perception leads to incorrect judgment of sliding tendency, and the grasping force setting is either too large (causing object damage) or too small (causing sliding). In complex assembly tasks, the lack of temperature and pressure feature fusion makes it impossible to perceive the micro contact between parts, leading to assembly deviation and collision. In addition, the removal of the MF module eliminates the attention mechanism, and the algorithm cannot dynamically prioritize key features for different tasks (e.g., pressure for grasping, texture for slippery objects), resulting in a 25% increase in decision-making error rate.

- Impact of AA module removal (TDHL-RL-w/o AA): The success rate decreases by 8.7% (basic tasks) and 15.3% (complex tasks). The AA module is the core of dynamic strategy adjustment, and its removal results in fixed high-level policies and low-level action execution rules; for dynamic environmental disturbances (e.g., vibrating platforms, wind force), the algorithm cannot adjust the grasping force and assembly path in real time, leading to a 30% increase in motion deviation (from  $\pm 0.3\text{N}$  to  $\pm 0.6\text{N}$ ). In addition, the fixed reward function weights (without dynamic adjustment for task difficulty) lead to unreasonable reward signals for complex tasks, resulting in slow convergence and low sample utilization efficiency.

Impact of simultaneous removal of MF and AA modules (TDHL-RL-w/o MF&AA): The success rate decreases by 16.0% (basic tasks) and 22.5% (complex tasks), which is not a simple superposition of the two single-module removals—this indicates a synergistic effect between the MF and AA modules. The MF module provides comprehensive, real-time tactile feature information for the AA module, and the AA module dynamically adjusts strategies based on the fused multimodal features; without the MF module, the AA module lacks reliable decision-making data, and without the AA module, the fused multimodal features cannot be effectively used for dynamic strategy adjustment. The synergistic failure of the two modules leads to the worst performance of the algorithm, with a

40% increase in operation time and a 35% increase in energy consumption in both basic and complex tasks.

## 6 Discussion

### 6.1. Comparison with classical nonlinear and adaptive control strategies

The proposed tactile-driven hierarchical reinforcement learning (TDHL-RL) framework is a model-free data-driven method for dexterous robotic manipulation, and its core performance characteristics are distinct from classical nonlinear control strategies (e.g., flatness-based control for dual-arm manipulators) and adaptive control strategies (e.g., nonlinear optimal control for autonomous underwater vehicles) that have been widely applied in complex robotic systems.

In terms of robustness, classical control strategies achieve stable operation by pre-calibrating the dynamic model of the system, and their performance degrades significantly when facing unknown environmental disturbances (e.g., random wind force, object surface variability) and unmodeled dynamics—this is a common challenge in dexterous robotic manipulation with diverse operating objects and complex environments. In contrast, the hierarchical adaptation mechanism of TDHL-RL does not rely on accurate mathematical modeling; it dynamically adjusts high-level policies and low-level action execution through real-time multimodal tactile feedback (pressure, texture, temperature), and the adaptive reward function further enhances the system's ability to resist random disturbances (e.g., the success rate remains 87.6% when grasping slippery objects on vibrating platforms).

In terms of real-time adaptability, the total latency of the tactile feedback loop in TDHL-RL is only 85ms (30ms for sensor processing, 55ms for model inference), which is comparable to the real-time performance of flatness-based control (50-100ms) and nonlinear optimal control (70-120ms) in robotic systems. More importantly, TDHL-RL can complete strategy adjustment within the latency window when the object state changes (e.g., sliding of smooth objects), while classical control strategies require offline re-tuning of controller parameters for new disturbance scenarios, lacking online real-time adaptability.

In terms of computational efficiency, classical nonlinear control requires real-time solution of high-dimensional dynamic equations for high-degree-of-freedom robotic manipulators, leading to exponentially increasing computational costs with the rise of system dimensions. The two-layer architecture of TDHL-RL decomposes the high-dimensional dexterous manipulation

task into high-level policy decision (macro task planning, low dimensionality) and low-level action execution (micro joint control, high dimensionality with modularization), which reduces the overall computational complexity by 40% (verified by the 40% acceleration of convergence speed in experiments). This decomposition makes TDHL-RL more suitable for dexterous manipulation tasks with 6+ degree-of-freedom robotic arms and multimodal tactile sensory input.

It is also necessary to acknowledge the limitations of TDHL-RL compared with classical control strategies: our framework requires a certain amount of training data (5000 episodes for feature weight optimization) in the early stage of deployment, and the training process relies on high-performance computing hardware (NVIDIA Jetson AGX Xavier); while classical control strategies can be directly deployed after model calibration and parameter tuning, with no additional training cost. This also points out the direction for the optimization of our method in future work: reducing the training sample size through sim-to-real transfer learning.

## 6.2. Comparison with finite-time fuzzy synchronization methods

Finite-time fuzzy synchronization methods are typical robust adaptive control strategies for MIMO time-varying delay systems, which are widely used in robotic manipulation for state feedback and motion synchronization control. A comparative analysis of the proposed TDHL-RL framework with finite-time fuzzy synchronization methods in dexterous robotic manipulation tasks shows two core advantages of our method in stability and rapid convergence: (1) In terms of stability, TDHL-RL achieves a force control error of only  $\pm 0.3\text{N}$  in grasping tasks through multimodal tactile feature fusion and adaptive reward function optimization, while finite-time fuzzy synchronization methods rely on single-modal state feedback (e.g., only joint angle feedback) and fixed control rules, leading to a larger force control error of  $\pm 0.6\text{N}$  in dexterous manipulation; (2) In terms of rapid convergence, the hierarchical decomposition architecture of TDHL-RL accelerates the convergence speed by 40% compared with the baseline algorithms, and its convergence efficiency is 20–30% higher than finite-time fuzzy synchronization methods, because our method can reuse high-level policies for similar manipulation tasks, while finite-time fuzzy synchronization methods need to re-design fuzzy rules for each new task, leading to slower convergence.

## 6.3. Application prospects in industrial and service robotics

The proposed tactile-driven hierarchical reinforcement learning framework has important practical application value in intelligent manufacturing and service robotics, two core fields of dexterous robotic manipulation. In smart manufacturing, the framework can be applied to the high-precision assembly of microelectronic components (e.g., semiconductor packaging), where the multimodal tactile feature fusion can accurately perceive the micro contact state between components, and the hierarchical adaptation mechanism can achieve submillimeter-level assembly accuracy ( $\pm 0.5\text{mm}$  clearance) in small spaces—this meets the high-precision assembly requirements of the electronics manufacturing industry. In service robotics, the framework is suitable for object manipulation of elderly care service robots and domestic service robots, where the real-time tactile feedback and adaptive grasping strategy can handle various daily objects with different shapes, materials and surface characteristics (e.g., fragile glass cups, slippery plastic bottles, soft cloth bags), and the low energy consumption characteristic (12.5J for basic tasks) is in line with the low-power operation requirements of mobile service robots. In addition, the framework can be extended to medical surgery robots, where the stable force control and high-precision manipulation performance can assist surgeons in minimally invasive surgical operations with tactile perception.

## 7 Conclusion

This study successfully constructed a hierarchical reinforcement learning dexterous operation framework driven by tactile information. The robot's operational ability in complex environments was effectively improved by integrating the adaptive hierarchical decision-making algorithm with multimodal tactile features. Simulation experimental data show that the algorithm is superior to traditional methods regarding operation success rate, execution efficiency, and environmental adaptability. The success rate of basic tasks exceeds 90%, the success rate of complex tasks is improved by more than 30% compared with the comparison algorithm, and the convergence speed is improved by 40%, which verifies the effectiveness and innovation of the algorithm. At wind speeds  $>7\text{m/s}$ , success drops from 87.6% to 65% due to sensor noise; future work will add wind disturbance filters. However, this research still has three main limitations that need to be addressed in future work: (1) Extreme dynamic environment robustness: The algorithm's success rate drops significantly from 87.6% to 65% when facing wind speeds  $>7\text{m/s}$ , due to increased tactile sensor noise and ineffective disturbance filtering; in addition, the framework lacks effective handling strategies for sudden

large external disturbances (e.g., unexpected object collision). (2) Sim-to-real gap: Although preliminary physical experiments verify the real-world generalization of the framework, there is still a slight performance decline between the simulation and real hardware, mainly due to the deviation between the simulation kinematic model and the real robot, and the measurement error of physical tactile sensors. (3) Training data and computational cost: The framework requires 5000 episodes for multimodal tactile feature weight optimization and a large number of training steps for policy learning, which relies on high-performance computing hardware (NVIDIA Jetson AGX Xavier) and is not suitable for low-power mobile robotic systems.

A detailed roadmap for future research is proposed to address the above limitations and further optimize the proposed framework:

- Extreme disturbance resistance optimization: Design a tactile sensor noise filtering module based on wavelet transform and Kalman filter to reduce the impact of extreme wind disturbance on tactile data; develop a sudden disturbance detection algorithm based on real-time tactile feature change rate, and design an emergency adjustment strategy for high-level policies to improve the algorithm's robustness in extreme dynamic environments.
- Sim-to-real transfer learning: Adopt the domain randomization method to add various random disturbances (sensor noise, model parameter deviation, environmental interference) to the simulation environment, narrowing the sim-to-real gap; design a lightweight fine-tuning network for the pre-trained TDHL-RL model, which only requires a small amount of real hardware data ( $\leq 500$  episodes) to complete model adaptation, reducing the physical experiment cost.
- Lightweight algorithm design: Compress the high-level policy network and tactile feature extractor via neural network pruning and quantization, reducing the model size by 60% while maintaining performance; design a federated learning framework for TDHL-RL, which enables multi-robot collaborative training without sharing raw data, reducing the single-robot training computational cost and adapting to low-power mobile robotic systems.
- Visual-tactile multimodal deep fusion: Integrate high-resolution visual sensors (RGB-D camera) with tactile sensors, design a cross-modal attention fusion module to fuse

visual spatial information and tactile contact information, further improving the algorithm's performance in unstructured environments with poor visual conditions (e.g., low light, occlusion).

Industrial-scale physical verification: Conduct large-scale physical experiments on the industrial robotic manipulation platform (e.g., ABB YuMi dual-arm robot) for high-precision electronic component assembly, and optimize the framework according to industrial actual needs (e.g., assembly cycle, energy consumption, accuracy), promoting the industrial application of the tactile-driven HRL framework.

## Acknowledgments

This work was supported in part by the Guangdong Power Grid Corporation (Grant No. GDKIXM20231037).

## References

- [1] L. Yang, B. Huang, Q. Li, Y. Y. Tsai, W. W. Lee, C. Song, J. Pan, Tacgnn: Learning tactile-based in-hand manipulation with a blind robot using hierarchical graph neural network, *IEEE Robotics and Automation Letters*. 8(6) (2023) 3605–3612. <https://doi.org/10.1109/LRA.2023.3264759>.
- [2] P. Y. Lajoie, G. Beltrame, Swarm-slam: Sparse decentralized collaborative simultaneous localization and mapping framework for multi-robot systems, *IEEE Robotics and Automation Letters*. 9(1) (2023) 475–482. <https://doi.org/10.1109/LRA.2023.3333742>.
- [3] H. Su, W. Qi, J. Chen, D. Zhang, Fuzzy approximation-based task-space control of robot manipulators with remote center of motion constraint, *IEEE Transactions on Fuzzy Systems*. 30(6) (2022) 1564–1573. <https://doi.org/10.1109/TFUZZ.2022.3157075>.
- [4] J. P. A. Yaacoub, H. N. Noura, O. Salman, A. Chehab, Robotics cyber security: Vulnerabilities, attacks, countermeasures, and recommendations, *International Journal of Information Security*. 21(1) (2022) 115–158. <https://doi.org/10.1007/s10207-021-00545-8>.
- [5] A. K. Inkulu, M. R. Bahubalendruni, A. Dara, S. K. Challenges and opportunities in human robot collaboration context of Industry 4.0—a state of the art review, *Industrial Robot: The International Journal of Robotics Research and Application*. 49(2) (2022) 226–239. <https://doi.org/10.1108/IR-04-2021-0077>.
- [6] M. Arduengo, C. Torras, L. Sentis, Robust and adaptive door operation with a mobile robot, *Intelligent Service Robotics*. 14(3) (2021) 409–425. <https://doi.org/10.1007/s11370-021-00366-7>.

- [7] S. L. Alip, J. Kim, K. H. Rha, W. K. Han, Future platforms of robotic surgery, *Urologic Clinics*. 49(1) (2022) 23–38. <https://doi.org/10.1016/j.ucl.2021.07.008>.
- [8] E. Triantafyllidis, F. Acero, Z. Liu, Z. Li, Hybrid hierarchical learning for solving complex sequential tasks using the robotic manipulation network ROMAN, *Nature Machine Intelligence*. 5(9) (2023) 991–1005. <https://doi.org/10.1038/s42256-023-00709-2>.
- [9] S. K. Sampath, N. Wang, H. Wu, C. Yang, Review on human-like robot manipulation using dexterous hands, *Cognitive Computation and Systems*. 5(1) (2023) 14–29. <https://doi.org/10.1049/ccs2.12073>.
- [10] Y. Lin, A. Church, M. Yang, H. Li, J. Lloyd, D. Zhang, N. F. Lepora, Bi-touch: Bimanual tactile manipulation with sim-to-real deep reinforcement learning, *IEEE Robotics and Automation Letters*. 8(9) (2023) 5472–5479. <https://doi.org/10.1109/LRA.2023.3295991>.
- [11] H. Wang, H. Zhang, L. Li, Z. Kan, Y. Song, Task-driven reinforcement learning with action primitives for long-horizon manipulation skills, *IEEE Transactions on Cybernetics*. 54(8) (2023) 4513–4526. <https://doi.org/10.1109/TCYB.2023.3298195>.
- [12] Z. V. Gbouna, G. Pang, G. Yang, Z. Hou, H. Lyu, Z. Yu, Z. Pang, User-interactive robot skin with large-area scalability for safer and natural human-robot collaboration in future telehealthcare, *IEEE Journal of Biomedical and Health Informatics*. 25(12) (2021) 4276–4288. <https://doi.org/10.1109/JBHI.2021.3082563>.
- [13] L. Liu, F. Guo, Z. Zou, V. G. Duffy, Application, development and future opportunities of collaborative robots (cobots) in manufacturing: A literature review, *International Journal of Human-Computer Interaction*. 40(4) (2024) 915–932. <https://doi.org/10.1080/10447318.2022.2041907>.
- [14] Y. Wang, J. Wang, Y. Li, T. Yang, C. Ren, The deep reinforcement learning-based VR training system with haptic guidance for catheterization skill transfer, *IEEE Sensors Journal*. 22(23) (2022) 23356–23366. <https://doi.org/10.1109/JSEN.2022.3212989>.
- [15] D. Yang, H. Liu, Human-machine shared control: New avenue to dexterous prosthetic hand manipulation, *Science China Technological Sciences*. 64(4) (2021) 767–773. <https://doi.org/10.1007/s11431-020-1710-y>.
- [16] J. Eßer, N. Bach, C. Jestel, O. Urbann, S. Kerner, Guided reinforcement learning: A review and evaluation for efficient and effective real-world robotics, *IEEE Robotics & Automation Magazine*. 30(2) (2022) 67–85. <https://doi.org/10.1109/MRA.2022.3207664>.
- [17] F. Sun, Y. Chen, Y. Wu, L. Li, X. Ren, Motion planning and cooperative manipulation for mobile robots with dual arms, *IEEE Transactions on Emerging Topics in Computational Intelligence*. 6(6) (2022) 1345–1356. <https://doi.org/10.1109/TETCI.2022.3146387>.
- [18] M. Hutsebaut-Buysse, K. Mets, S. Latré, Hierarchical reinforcement learning: A survey and open research challenges, *Machine Learning and Knowledge Extraction*. 4(1) (2022) 172–221. <https://doi.org/10.3390/make4010009>.
- [19] W. Si, N. Wang, C. Yang, A review on manipulation skill acquisition through teleoperation-based learning from demonstration, *Cognitive Computation and Systems*. 3(1) (2021) 1–16. <https://doi.org/10.1049/ccs2.12005>.