

AAAI Press Anonymous Submission Instructions for Authors Using L^AT_EX

Anonymous submission

Abstract

AAAI creates proceedings, working notes, and technical reports directly from electronic source furnished by the authors. To ensure that all papers in the publication have a uniform appearance, authors must adhere to the following instructions.

Introduction

Technical Appendix

Plant Operations

Predictive Maintenance: The task of predicting the Remaining Useful Life (RUL) is formulated as a sequential decision-making problem within a Markov Decision Process (MDP) framework. The objective of the learning model is to learn a policy π that, given the current state s_t of the industrial equipment, can accurately estimate the remaining number of operational time steps before failure. The model is trained in a simulation environment constructed from historical run-to-failure data, which serves as a data-driven emulator of the equipment’s degradation dynamics. At each time step t , the model observes a state vector s_t —comprising sensor readings, health indicators, and their temporal change rates—generates an RUL prediction a_t , and receives a reward r_t based on the prediction’s accuracy. Through repeated interaction, the model learns a policy that maximizes long-term prediction accuracy or, equivalently, minimizes cumulative prediction error. The underlying MDP is defined by the following components. The state space ($\mathbf{s}_t \in \mathbf{S}$) represents the condition of the industrial equipment as a feature vector constructed by combining sensor data, observable health parameters that characterize degradation, and their discrete rates of change (degradation velocity) over a fixed interval. The action space ($\mathbf{a}_t \in \mathbf{A}$) is a continuous scalar value denoting the model’s predicted RUL, expressed as $a_t = \text{RUL}_{\text{predicted}}(t)$. The reward function (\mathbf{r}_t) imposes an asymmetric penalty based on the normalized prediction error $\tilde{\epsilon}_t$:

$$r(s_t, a_t) = \begin{cases} \lambda_{\text{early}} \cdot \tilde{\epsilon}_t, & \text{if } a_t < \text{RUL}_{\text{true}}(t), \\ \lambda_{\text{late}} \cdot \tilde{\epsilon}_t, & \text{otherwise,} \end{cases}$$

where λ_{early} and λ_{late} are penalty coefficients corresponding to underestimation and overestimation, respectively.

This asymmetric design allows the learning model to balance proactive risk mitigation—by penalizing late predictions more severely—and the avoidance of false alarms, by assigning higher penalties to early predictions when desired. We employ the Soft Actor-Critic (SAC) algorithm to solve this continuous control MDP. SAC is an off-policy actor-critic method that promotes sample efficiency and stable convergence by incorporating an entropy term into the reward objective to enhance exploration. The objective is to maximize:

$$J(\pi) = \sum_{t=0}^T \mathbb{E}_{(s_t, a_t) \sim \rho_\pi} [\gamma^t r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot|s_t))]$$

where π is the policy, ρ_π the state–action distribution, γ the discount factor, \mathcal{H} the policy entropy, and α a temperature parameter balancing reward against entropy. This encourages the model to optimize long-term prediction accuracy while exploring diverse degradation patterns, avoiding overconfident predictions. The SAC architecture comprises three key components: a stochastic actor (policy network, π_ϕ) that maps states to a Gaussian action distribution using reparameterization for differentiable sampling; twin critic networks ($\mathbf{Q}_{\theta_1}, \mathbf{Q}_{\theta_2}$) that independently estimate state–action values to reduce overestimation bias; and target networks ($Q_{\theta'_1}, Q_{\theta'_2}$) that stabilize learning via Polyak averaging. We enhance sample efficiency using a Prioritized Experience Replay (PER) buffer, which samples transitions with probability proportional to their temporal-difference (TD) error:

$$P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$$

where p_i is the priority of transition i and α controls prioritization strength, focusing learning on informative experiences. Training proceeds episodically over run-to-failure trajectories. At each step t , the model observes state s_t , predicts RUL a_t , and receives accuracy-based reward r_t . Transitions are stored in the PER buffer, with mini-batches used to update networks: critics minimize Bellman error, while the actor maximizes entropy-augmented expected reward. The temperature parameter α is automatically tuned to balance exploration and exploitation. Model generalization is assessed on a validation set using MAE and MAPE between predicted and true RUL values, enabling accurate, data-driven RUL estimation for predictive maintenance. In

short, Predictive maintenance is traditionally framed as a supervised regression problem, focusing on estimating equipment's RUL. In contrast, reinforcement learning approaches frame maintenance planning and intervention as MDPs, enabling sequential, policy-driven decision-making. This allows RL agents to learn optimal strategies for when to perform maintenance and what specific actions to take, maximizing system uptime while minimizing costs and risks.

Plant-Wide Optimization and Control of Industrial Systems: Industrial plant optimization employs a hierarchical framework combining process optimization and control. At the supervisory level, process optimization determines optimal setpoints for controlled variables (CVs) by balancing economic objectives like profit or efficiency against process constraints and operating conditions. At the regulatory level, process control maintains CVs at their setpoints by adjusting manipulated variables, ensuring disturbance rejection and operational stability. This integrated approach enables plants to achieve both economic optimization and safe operational reliability in dynamic environments. This work addresses learning regulatory control policies directly from historical plant data, providing a data-driven alternative when accurate first-principles models are unavailable. The control task is formulated as a continuous-state, continuous-action Markov Decision Process (MDP), where an autonomous reinforcement learning agent learns a feedback control policy mapping process states to optimal manipulated variable (MV) adjustments. The MDP is defined by the tuple (S, A, P, R) , where the state space ($s_t \in S$) represents the plant's dynamic operating condition through measured and inferred variables. The action space ($a_t \in A$) specifies manipulated variable values. The reward function encodes control objectives through a weighted penalty formulation:

$$r(s_t, a_t) = - \left[w_{\text{track}} \cdot \frac{|CV(s_t) - SP|}{\epsilon_{\text{max}}} + w_{\text{act}} \cdot \|a_t - a_{t-1}\|^2 \right],$$

where $CV(s_t)$ is the controlled variable, SP its target setpoint, and $w_{\text{track}}, w_{\text{act}}$ weight tracking precision against actuation smoothness. Our methodology employs an offline, model-based reinforcement learning approach with two sequential phases. First, offline system identification trains a deep neural network f_θ to approximate system dynamics $s_{t+1} = f(s_t, a_t)$ from historical data $\mathcal{D} = \{(s_t, a_t, s_{t+1})\}$ by minimizing:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \mathcal{D}} [\|f_\theta(s_t, a_t) - s_{t+1}\|^2].$$

The trained model f_{θ^*} serves as a differentiable, data-driven simulator for risk-free policy training. Second, offline policy optimization uses the Soft Actor-Critic (SAC) algorithm within a maximum entropy framework:

$$J(\pi) = \sum_{t=0}^T \mathbb{E}_{(s_t, a_t) \sim \rho_\pi} [\gamma^t r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot|s_t))],$$

where π is the policy, ρ_π the state-action marginal, and α a temperature parameter. The agent architecture includes a stochastic actor (π_ϕ) parameterizing a Gaussian policy and

twin critic networks ($Q_{\theta_1}, Q_{\theta_2}$) employing double-Q learning. Target networks and experience replay ensure stable, efficient off-policy learning. The training process is conducted episodically and entirely offline within the learned environment. The agent's networks are updated using mini-batches sampled from the replay buffer, where the critics minimize the Bellman error while the actor maximizes expected future reward and entropy. For evaluation, the trained policy is deployed in simulation and benchmarked against historical operational data. Performance is assessed based on the policy's ability to maintain the controlled variable at setpoint with minimal deviation while ensuring smooth manipulated variable adjustments. This offline model-based RL framework enables long-horizon planning and direct optimization of shaped reward functions that encode operational objectives, process constraints, and economic goals, thereby bridging traditional model-based control with deep reinforcement learning. In short, traditional plant-wide process control relies on model-based predictive control using predefined first-principles process models. In contrast, RL reframes control as a sequential decision-making problem, enabling adaptive, model-free policies that optimize long-term performance under uncertainty.

Knowledge Engineering via Offline Adversarial Imitation Learning We address the problem of distilling human operator expertise from a fixed dataset of historical operation trajectories into a control policy. Each trajectory $\tau_i = \{(s_t^i, a_t^i)\}_{t=1}^{T_i}$ consists of process states s_t (encompassing all measured, controlled, and dependent variables) and actions a_t (the manipulated variables adjusted by the operator). The goal is to learn a stochastic policy $\pi_\phi(a|s)$ that mimics the expert's decision-making from the dataset $\mathcal{D} = \{\tau_i\}_{i=1}^N$, preserving this tacit knowledge for safer and more intelligent operation. We model the control problem as a Markov Decision Process $\mathcal{M} = (S, \mathcal{A}, \mathcal{P}, r, \gamma)$ with continuous states S , continuous actions \mathcal{A} , transition dynamics \mathcal{P} , an unknown reward function r , and a discount factor γ . The objective is to recover a policy π_ϕ whose state-action occupancy distribution $\rho_{\pi_\phi}(s, a)$ matches that of the expert policy π_E . The occupancy measures the discounted time the policy spends in each state-action pair:

$$\rho_{\pi_\phi}(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s, a_t = a | \pi_\phi)$$

A matched occupancy implies the learned policy is statistically indistinguishable from the expert. Offline Generative Adversarial Imitation Learning (GAIL) frames imitation as a two-player game between a policy π_ϕ and a discriminator $D_\psi(s, a)$. The adversarial objective is:

$$\min_{\pi_\phi} \max_{D_\psi} \mathbb{E}_{\rho_{\pi_E}} [\log D_\psi(s, a)] + \mathbb{E}_{\rho_{\pi_\phi}} [\log(1 - D_\psi(s, a))]$$

The discriminator learns to distinguish expert pairs $(s, a) \sim \rho_{\pi_E}$ from policy-generated pairs $(s, a_\phi) \sim \rho_{\pi_\phi}$, while the policy learns to fool the discriminator. In the offline setting, we sample expert states $s_E \sim D_E$ and generate corresponding actions with the current policy $a_\phi \sim \pi_\phi(\cdot|s_E)$. The discriminator loss includes a gradient penalty for stability:

$$\begin{aligned}\mathcal{L}_D = & -\mathbb{E}_{(s_E, a_E) \sim \mathcal{D}_E} [\log D_\psi(s_E, a_E)] \\ & -\mathbb{E}_{s_E \sim \mathcal{D}_E, a_\phi \sim \pi_\phi(\cdot|s_E)} [\log(1 - D_\psi(s_E, a_\phi))] \\ & + \lambda \mathbb{E}_{(\hat{s}, \hat{a})} [(\|\nabla D_\psi(\hat{s}, \hat{a})\|_2 - 1)^2]\end{aligned}$$

where (\hat{s}, \hat{a}) is a random interpolation between expert and policy data: $(\hat{s}, \hat{a}) = \alpha(s_E, a_E) + (1 - \alpha)(s_E, a_\phi)$ with $\alpha \sim U(0, 1)$. The trained discriminator provides a surrogate reward signal: $r_\psi(s, a) = \log \sigma(D_\psi(s, a))$, where $\sigma(\cdot)$ is the sigmoid function. This reward guides the policy to produce expert-like behavior. We use Implicit Q-Learning (IQL) for safe, offline policy optimization. IQL uses the surrogate reward to learn a policy across three conservative learning stages: (a) *Value Function Learning*: The value function $V_\theta(s)$ is learned via asymmetric expectile regression:

$$\mathcal{L}_V = \mathbb{E}_{(s, a) \sim \mathcal{D}} [L_2^\tau(Q_\omega(s, a) - V_\theta(s))]$$

The expectile loss $L_2^\tau(u)$, with $u = Q_\omega(s, a) - V_\theta(s)$, is defined as:

$$L_2^\tau(u) = \begin{cases} \tau \cdot u^2 & \text{if } u \geq 0 \\ (1 - \tau) \cdot u^2 & \text{if } u < 0 \end{cases}$$

For $\tau > 0.5$, this emphasizes positive advantages, enabling implicit maximization over actions without explicitly querying the policy. (b) *Q-Function Learning*: The Q-function is updated using the Bellman equation with the value target:

$$\mathcal{L}_Q = \mathbb{E}_{(s_E, a_E, s'_E) \sim \mathcal{D}_E} [(Q_\omega(s_E, a_E) - (r_\psi(s_E, a_E) + \gamma V_\theta(s'_E)))^2]$$

(c) *Policy Extraction*: The policy is improved via advantage-weighted regression:

$$\begin{aligned}\mathcal{L}_\pi = & -\mathbb{E}_{s_E \sim \mathcal{D}_E} [\exp(\beta A(s_E, a_\phi)) \log \pi_\phi(a_\phi|s_E)], \\ A(s_E, a_\phi) = & Q_{\min}(s_E, a_\phi) - V_\theta(s_E)\end{aligned}$$

where $a_\phi \sim \pi_\phi(\cdot|s_E)$, $Q_{\min} = \min(Q_{\omega_1}, Q_{\omega_2})$ for twin Q-networks, and $\beta > 0$ is an inverse temperature parameter. The complete Offline GAIL-IQL objective integrates adversarial imitation with constrained policy optimization:

$$\begin{aligned}\min_{\pi_\phi, Q_\omega, V_\theta} \max_{D_\psi} = & -\mathbb{E}_{(s_E, a_E) \sim \mathcal{D}_E} [\log D_\psi(s_E, a_E)] \\ & -\mathbb{E}_{s_E \sim \mathcal{D}_E, a_\phi \sim \pi_\phi(\cdot|s_E)} [\log(1 - D_\psi(s_E, a_\phi))] \\ & + \lambda \mathcal{L}_{\text{grad-pen}} + \mathcal{L}_V + \mathcal{L}_Q + \mathcal{L}_\pi\end{aligned}$$

GAIL and IQL serve complementary roles: GAIL’s discriminator defines *what* to learn via adversarial reward shaping, while IQL defines *how* to learn it safely via conservative offline policy optimization. Policy performance is evaluated using Mean Squared Error between expert and learned actions on held-out states: $\text{MSE} = \frac{1}{N} \sum_{i=1}^N \|a_E^{(i)} - a_\phi^{(i)}\|^2$. This framework distills operational expertise into robust control policies without online interaction.

Knowledge Engineering Via Inverse Reinforcement Learning The fundamental goal is to automate complex human skills by extracting and embedding an expert process plant operator’s implicit knowledge into a robust, AI-driven

control policy. We aim to capture this tacit understanding from historical data, distilling it into an autonomous agent capable of safe and intelligent operation, just like its human counterpart. We frame the problem of learning from expert demonstrations using Adversarial Inverse Reinforcement Learning (AIRL). Given a dataset of expert trajectories $\mathcal{D} = \{\tau_i\}_{i=1}^N$, where each trajectory $\tau_i = \{(s_t, a_t, s_{t+1})\}_{t=1}^{T_i}$ contains state-action-next state tuples, our goal is to recover the expert’s underlying reward function $r_\psi(s, a)$ and a policy $\pi_\phi(a|s)$ that mimics their behavior. AIRL uses an adversarial game between a discriminator D_ψ and a policy π_ϕ :

$$\min_{\pi_\phi} \max_{D_\psi} \mathbb{E}_{(s, a) \sim \rho_{\pi_E}} [\log D_\psi(s, a)] + \mathbb{E}_{(s, a) \sim \rho_{\pi_\phi}} [\log(1 - D_\psi(s, a))]$$

The policy learns to generate state-action pairs that the discriminator cannot distinguish from expert data. The discriminator has a specific structure to enable reward recovery:

$$D_\psi(s, a, s') = \frac{\exp(f_\psi(s, a, s'))}{\exp(f_\psi(s, a, s')) + 1}$$

where $f_\psi(s, a, s')$ is an advantage-like function:

$$f_\psi(s, a, s') = r_\psi(s, a) + \gamma V_\psi(s') - V_\psi(s)$$

This decomposition allows the reward function $r_\psi(s, a)$ to be disentangled from the dynamics. We use separate neural networks for the reward $r_\psi(s, a) = \text{NN}_{\text{reward}}(s, a; \psi_r)$ and value function $V_\psi(s) = \text{NN}_{\text{value}}(s; \psi_v)$. The discriminator’s loss includes a gradient penalty for stable training:

$$\begin{aligned}\mathcal{L}_D = & -\mathbb{E}_{(s_E, a_E) \sim \mathcal{D}_E} [\log D_\psi(s_E, a_E, s'_E)] \\ & -\mathbb{E}_{s_E \sim \mathcal{D}_E, a_\phi \sim \pi_\phi} [\log(1 - D_\psi(s_E, a_\phi, s'_E))] \\ & + \lambda \mathbb{E}_{(\hat{s}, \hat{a}, \hat{s}')} [(\|\nabla f_\psi(\hat{s}, \hat{a}, \hat{s}')\|_2 - 1)^2]\end{aligned}$$

where $(\hat{s}, \hat{a}, \hat{s}')$ are interpolated between expert and policy samples. For policy optimization, we use Implicit Q-Learning (IQL) to learn offline from the fixed dataset. The value function is learned via asymmetric expectile regression:

$$\mathcal{L}_V = \mathbb{E}_{(s, a) \sim \mathcal{D}} [L_2^\tau(Q_\omega(s, a) - V_\theta(s))]$$

$$\text{where } L_2^\tau(u) = \begin{cases} \tau u^2 & \text{if } u \geq 0 \\ (1 - \tau) u^2 & \text{if } u < 0 \end{cases} \text{ and } \tau \in (0.5, 1)$$

weights the upper tail of Q-values. The Q-function is trained using temporal difference learning:

$$\begin{aligned}y = & r_\psi(s, a) + \gamma V_\theta(s') \\ \mathcal{L}_Q = & \mathbb{E}_{(s, a, s') \sim \mathcal{D}} [(Q_\omega(s, a) - y)^2]\end{aligned}$$

The policy is optimized via advantage-weighted regression:

$$\mathcal{L}_\pi = -\mathbb{E}_{s \sim \mathcal{D}} [\exp(\beta A(s, a)) \log \pi_\phi(a|s)], \quad a \sim \pi_\phi(\cdot|s)$$

where $A(s, a) = Q_{\min}(s, a) - V_\theta(s)$ is the advantage estimate and β is an inverse temperature parameter. The policy is a Gaussian $\pi_\phi(a|s) = \mathcal{N}(\mu_\phi(s), \sigma_\phi(s)^2)$, using the reparameterization trick $a = \mu_\phi(s) + \sigma_\phi(s) \cdot \epsilon$, $\epsilon \sim \mathcal{N}(0, I)$

for gradient estimation. The complete optimization combines adversarial imitation learning with offline reinforcement learning:

$$\begin{aligned} \min_{\pi_\phi, Q_\omega, V_\theta} \max_{D_\psi} & -\mathbb{E}_{\mathcal{D}_E} [\log D_\psi(s, a, s')] \\ & -\mathbb{E}_{s \sim \mathcal{D}_E, a \sim \pi_\phi} [\log(1 - D_\psi(s, a, s'))] \\ & + \lambda \mathcal{L}_{\text{grad-pen}} + \mathcal{L}_V + \mathcal{L}_Q + \mathcal{L}_\pi \end{aligned}$$

This framework simultaneously learns a reward function that explains expert behavior, a value function that estimates long-term returns, and a policy that acts according to the learned reward. Training is entirely offline. For evaluation, the trained policy is deployed on held-out expert states, with performance measured by the Mean Squared Error between predicted and true expert actions and qualitative trajectory comparisons.

Supply Chain Optimization

Route Planning & Optimization We address chemical distribution logistics through route planning and optimization for multi-vehicle operations. Vehicles transport products from depots to customers, with each depot serving as both origin and destination. Vehicles depart fully loaded, deliver to clients, and may collect by-products or empty containers before returning within strict time windows and capacity constraints. Each customer is served exactly once. We propose a deep reinforcement learning framework based on Proximal Policy Optimization (PPO) for learning efficient routing strategies. The problem is formulated on a complete directed graph $G = (V, E)$, where $V = D \cup C$ comprises depots D and customers C , and E represents travel connections. Each depot $d \in D$ operates within time window $[e_d, l_d]$. Each customer $c \in C$ has delivery demand q_c , pickup quantity p_c , service time s_c , and time window $[e_c, l_c]$. Each edge $(i, j) \in E$ has travel time t_{ij} and distance d_{ij} from real-world routing data (Google Maps API). The heterogeneous fleet K consists of vehicles $k \in K$, each with capacity Q_k , start/end depots d_k^s, d_k^e , start time τ_k , maximum duration T_k^{\max} , variable cost c_k^{km} , and fixed cost c_k^{fixed} . The objective constructs feasible routes $\{R_1, \dots, R_{|K|}\}$, where $R_k = (d_k^s, c_1, \dots, c_m, d_k^e)$, minimizing total system cost:

$$Z = \sum_{k \in K} \left(c_k^{\text{fixed}} \cdot y_k + c_k^{\text{km}} \cdot \sum_{(i,j) \in R_k} d_{ij} \right),$$

where $y_k = 1$ if vehicle k is used, 0 otherwise. Constraints include: (i) capacity limits Q_k , (ii) time-window feasibility (allowing early arrival waiting), (iii) route duration limits T_k^{\max} , (iv) depot windows $[e_d^s, l_d^e]$, and (v) each customer visited exactly once. We model this as a Markov Decision Process where a PPO agent sequentially constructs routes. The *state* s_t encodes node features, vehicle attributes (load, time, location), and a mask excluding infeasible actions. The *action* a_t selects the next node to visit. The *reward function* penalizes travel distance and waiting:

$$r_t = -(d_{ij} \cdot c_k^{\text{km}} + \beta \cdot \text{Wait Time}),$$

where β weights the waiting penalty. The PPO agent uses an *Actor-Critic architecture* with a shared Transformer encoder for context-aware embeddings. The *Actor* produces masked action probabilities ensuring feasibility; the *Critic* estimates state-value $V(s_t)$. Policy updates use PPO’s clipped objective:

$$L^{\text{CLIP}} = \mathbb{E}_t [\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t)],$$

where $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$ is the probability ratio and A_t is the advantage estimate. The total loss combines value function and entropy terms:

$$L^{\text{TOTAL}} = L^{\text{CLIP}} - c_1 L^{VF} + c_2 H(\pi(\cdot|s_t)),$$

where $L^{VF} = \mathbb{E}_t [(V_\theta(s_t) - V_{\text{target}})^2]$ and V_{target} is the GAE-based return estimate. The agent iteratively refines its policy through gradient-based updates. Performance is evaluated against heuristic benchmarks (Nearest Neighbor, Savings, Random Insertion) using metrics including total cost, distance, fleet utilization, and unserved customers (Table 1).

Table 1: Performance Comparison of Vehicle Routing Problem Solvers for Chemical Distribution

Metric	PPO RL Agent	Nearest Neighbor	Savings (C&W)	Random Insertion
Cost (\$)	904.77	997.02	1,287.07	1,589.04
Distance (km)	280.42	331.68	553.53	641.13
Vehicles	2	2	3	3
Unserved	1	1	0	0

Network & Route Optimization in Chemical Distribution Logistics

We address the Network & Route Optimization problem for a multi-plant distribution system, formulated as a Multi-Depot Vehicle Routing Problem with Sourcing Costs (MDVRP-SC). The system consists of multiple production plants acting as depots, each operating a fleet of vehicles that deliver products to customers. Customer demands may be split into delivery jobs assignable to different plants. The optimization jointly determines (i) plant–customer assignments and (ii) vehicle routes that start and end at the corresponding plants. These decisions must respect key operational constraints: plant production capacities, vehicle load limits, customer time windows, maximum route durations, plant operating hours, and fixed unloading times at customer sites. The objective is to minimize total system cost, including both production costs at plants and transportation costs for deliveries. The system is represented by a complete directed graph $G = (V, E)$, where $V = P \cup C$ consists of production plants $P = \{p_1, \dots, p_m\}$ and customers $C = \{c_1, \dots, c_n\}$. Each plant $p \in P$ has production capacity Q_p (tons), unit production cost α_p (\$/ton), and operational time window $[e_p, l_p]$. Each customer $c \in C$ has demand d_c (tons), transfer time u_c (seconds), and delivery time window $[e_c, l_c]$. Each vehicle $k \in K$ has capacity q_k (tons), home plant $h_k \in P$, start time τ_k , maximum route duration T_k^{\max} , variable transportation cost β_k (\$/km), and fixed dispatch cost γ_k (\$). The integrated optimization involves

two coupled decisions: (i) the *sourcing decision* via binary variables $y_{p,c} = 1$ if plant p serves customer c , and (ii) the *routing decision* determining vehicle routes $\{R_1, \dots, R_{|K|}\}$. The total cost is:

$$Z = \sum_{p \in P} \sum_{c \in C} \alpha_p \cdot d_c \cdot y_{p,c} + \sum_{k \in K} \left(\gamma_k \cdot \delta_k + \beta_k \cdot \sum_{(i,j) \in R_k} \text{dist}_{ij} \right),$$

where $\delta_k = 1$ if vehicle k is used and dist_{ij} is the travel distance between locations i and j . The objective function (Z) minimizes the sum of production costs (where to make the chemicals) and transportation costs (how to deliver them). Key constraints include:

$$\sum_{c \in C} d_c \cdot y_{p,c} \leq Q_p \quad \forall p \in P \quad (\text{Plant Capacity}),$$

$$\sum_{p \in P} y_{p,c} = 1 \quad \forall c \in C \quad (\text{Demand Satisfaction}),$$

along with vehicle capacity limits, hard time-window constraints $e_c \leq \text{start}_c \leq l_c$ for each customer c , where start_c is the service start time, and l_c = latest allowable service start time for customer c . If a vehicle arrives before (e_c), it must wait; if it cannot arrive by (l_c), the customer cannot be served by that vehicle. Route duration constraints ensuring that each vehicle k completes its route within maximum duration T_k^{\max} , where $\text{duration}(R_k) \leq T_k^{\max}$ accounts for total travel and service time. To manage the computational complexity of this NP-hard problem, we employ a two-stage decomposition framework. The first stage assigns customers to plants using a hybrid metaheuristic. It begins with greedy initialization based on composite cost $\text{Cost}_{p,c} = \alpha_p \cdot d_c + \text{dist}(p,c) \cdot \beta$, where β is the average transportation cost rate. This is refined via local search using scoring function $\text{Score} = \sum_{p \in P} \text{Cost}_p + \lambda \cdot \text{Unserved}_p$, where λ is a penalty factor and Unserved_p denotes unserved customers at plant p , Cost_p is the production and estimated transportation costs for serving all customers assigned to that particular plant p . The second stage solves routing subproblems in parallel using Reinforcement Learning (RL) based on Proximal Policy Optimization (PPO). Each subproblem is a single-depot VRP modeled as a Markov Decision Process, where state s_t captures node features, vehicle status, and an action mask ensuring feasibility. Action a_t selects the next customer. The reward function is:

$$r_t = -(\text{dist}_{ij} \cdot \beta_k + \omega \cdot \text{wait_time}),$$

where ω is a penalty coefficient and wait_time denotes idle waiting duration. The policy uses an Actor-Critic architecture with a shared Transformer encoder generating context-rich graph embeddings. The **Actor** produces probabilities over feasible actions via masked attention; the **Critic** estimates value function $V(s_t)$. Training minimizes the PPO clipped objective:

$$L^{\text{CLIP}}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right],$$

where $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$ is the probability ratio, \hat{A}_t is the advantage estimate, and ϵ is the clipping hyperparameter. The total loss is:

$$L^{\text{TOTAL}} = L^{\text{CLIP}} - c_1 L^{\text{VF}} + c_2 H(\pi(\cdot|s_t)),$$

where L^{VF} is the value function loss, H is policy entropy, and c_1, c_2 are weighting coefficients. Performance is benchmarked against a baseline Greedy heuristic using metrics including Total Cost Z , Vehicles Used, and Unserved Customers, demonstrating the robustness of the two-stage framework for large-scale chemical distribution logistics (see Table 2 and Figure 1).

Table 2: Performance comparison between Reinforcement Learning and Greedy heuristic solvers for the Multi-Depot Vehicle Routing Problem. Both algorithms achieved identical results, finding the optimal solution that minimizes total costs while serving all construction sites within operational constraints.

Metric	PPO_RL	GREEDY
Total Network Cost	\$19,470.72	\$19,470.72
Production Cost	\$11,605.00	\$11,605.00
Transport Cost	\$7,865.72	\$7,865.72
Trucks Dispatched	10	10
Unserved Sites	0	0

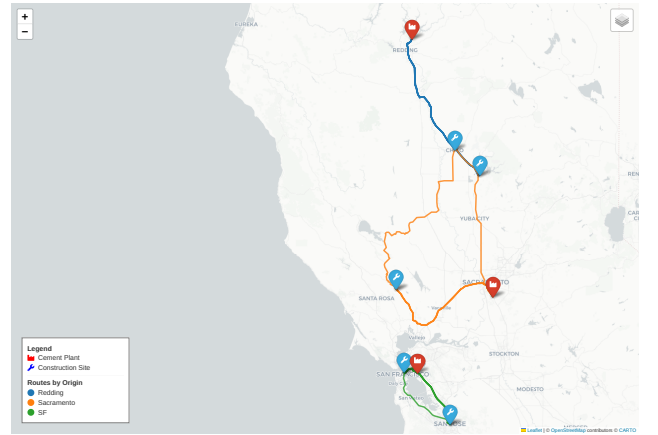


Figure 1: Geographic visualization of the optimized cement distribution network showing plants (red), construction sites (blue), and color-coded routes by originating facility.

Holistic Supply Chain Optimization Modern industrial supply chains are complex, dynamic systems where traditional decoupled optimization methods fail to capture the interdependencies among procurement, production, and distribution, often resulting in globally suboptimal performance. We address the problem of Holistic Supply Chain Optimization (HSCO) by formulating a unified model that jointly manages raw material sourcing, multi-facility production

scheduling, and finished product distribution under real-world constraints such as demand uncertainty and lead times (the total time required for processing, production, and delivery). The goal is to develop a learning-based framework that discovers an integrated policy to maximize cumulative profit across the supply chain. To achieve this, we propose a Graph Convolutional Network-based Proximal Policy Optimization (GCN-PPO) agent that learns end-to-end policies over a graph-structured representation of the supply chain. The agent maximizes profit by jointly considering revenue, operational costs, and penalties arising from unmet or delayed demand, effectively capturing the interdependent decision-making across all stages of the supply chain. The chemical supply chain is modeled as a network flow problem over a directed graph $G = (V, E)$, where the vertex set $V = S \cup P \cup C$ comprises suppliers S , production plants P , and customers C . The edge set E represents feasible material flows, including supplier-to-plant links for raw material deliveries and plant-to-customer links for product distribution. Each edge $(i, j) \in E$ has an associated distance d_{ij} , used to compute transportation costs. The problem unfolds over a discrete time horizon of T days, during which the system evolves through sequential decisions under stochastic demands. The system state at each day t includes raw material inventories I_{pt}^{rm} , finished product inventories I_{pt}^{fp} , customer backorders B_{ct} , in-transit inventory L_t subject to lead times, and a temporal feature $\phi(t)$ that captures demand seasonality. The action a_t is a continuous vector in $[0, 1]$ comprising three decision classes: procurement fractions $a_{s,p}^{procure}$ governing material flow from supplier s to plant p ; production fractions $a_{p,r}^{produce}$ determining operating levels for each recipe r at plant p ; and distribution fractions $a_{p,c,f}^{distribute}$ allocating product f from plant p to customer c . Inventory dynamics follow mass-balance relationships under lead time delays. The raw material inventory evolves as

$$I_{p,t+1}^{rm} = I_{pt}^{rm} + \sum_{s \in S} A_{s,p,t}^{arrive} - \sum_{r \in R_p} a_{p,r,t}^{produce} \cdot M_r^{in},$$

where $A_{s,p,t}^{arrive}$ represents shipments received after a supplier lead time λ_s , and M_r^{in} denotes raw material consumption per unit of production. The finished product inventory is updated as

$$I_{p,t+1}^{fp} = I_{pt}^{fp} + \sum_{r \in R_p} a_{p,r,t}^{produce} \cdot M_r^{out} - \sum_{c \in C} a_{p,c,f,t}^{distribute} \cdot I_{pft}^{fp},$$

where M_r^{out} represents the output mass per recipe. Customer demand D_{cft} follows a stochastic cyclic model defined as $D_{cft} = \mu_{cf} + A_{cf} \sin(2\pi t / \tau_{cf}) + \epsilon_{cf}$, where μ_{cf} is the base demand, A_{cf} the seasonal amplitude, τ_{cf} the demand cycle period, and $\epsilon_{cf} \sim \mathcal{N}(0, \sigma_{cf}^2)$ the demand noise with standard deviation σ_{cf} . Unmet demand accumulates as backorders with daily penalties or results in lost sales penalties, depending on service-level agreements.

The objective is to maximize expected cumulative profit, represented by a reward function that balances revenue gen-

eration against costs and penalties:

$$\text{Maximize } \mathbb{E} \left[\sum_{t=0}^T r_t \right],$$

$$r_t = \frac{1}{N_{profit}} (\text{Revenue}_t - \text{Cost}_t) - \beta \cdot \text{BackorderPenalty}_t$$

where Revenue_t represents fulfilled product sales, Cost_t aggregates procurement, production (including startup), inventory holding, transportation, and lost-sales penalties, $\text{BackorderPenalty}_t$ quantifies daily costs incurred due to unfulfilled customer demand, N_{profit} is a normalization constant for stable training, and β controls the trade-off between profitability and service level. We employ an Actor-Critic architecture in which both the policy and value functions are parameterized by a shared Graph Convolutional Network (GCN). The supply chain is naturally represented as a graph where nodes—corresponding to suppliers, plants, and customers—are encoded as feature vectors containing node type, normalized inventory and backorder levels, and temporal features, while edges represent feasible material flow paths. A multi-layer GCN processes this graph to generate enriched node embeddings, enabling each node to aggregate contextual information from its local neighborhood and capture complex interdependencies across the supply chain. After the final GCN layer, a global mean pooling operation produces a graph-level embedding that serves as input to both heads. The **Critic** head $V(s_t)$ uses a multilayer perceptron (MLP) to map this embedding to a scalar value estimate, while the **Actor** head $\pi(a_t|s_t)$ outputs parameters of a bounded policy distribution. Each dimension of the continuous action vector a_t is modeled using a Beta distribution:

$$a_t \sim \text{Beta}(\alpha, \beta), \quad \text{where } \alpha, \beta = \text{Softplus}(\text{MLP}(s_t)) + \text{offset}.$$

Training employs Proximal Policy Optimization (PPO) with the clipped surrogate objective:

$$L^{CLIP}(\theta) = \mathbb{E}_t [\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t)],$$

where $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$ is the probability ratio, A_t is the Generalized Advantage Estimate, and ϵ is the clipping parameter. The total loss function includes additional value and entropy terms:

$$L^{TOTAL} = L^{CLIP} - c_1 L^{VF} + c_2 H(\pi(\cdot|s_t)),$$

where L^{VF} is the value function loss, H is the policy entropy, and c_1, c_2 are weighting coefficients. The PPO agent is trained through extensive environment interaction over multiple episodes, each representing sequential daily decision-making across the planning horizon. Collected trajectories are used for minibatch gradient updates under the PPO objective, enabling the agent to progressively learn coordinated procurement, production, and distribution strategies. The trained agent is benchmarked against two baseline policies: (1) a **Greedy Policy**, which always operates production and distribution at maximum capacity, testing whether the RL agent can make more balanced and cost-aware decisions; and (2) an **(s, S) Policy**, which follows classical inventory

control rules—ordering when stock falls below threshold s and replenishing up to level S , with production and distribution governed by inventory and demand forecasts. Evaluation metrics include Total Profit, Operational Costs, and Service Level (measured by unmet demand).

Line-Haul Scheduling Optimization

The Line-Haul Scheduling Optimization problem in industrial chemical logistics involves assigning and sequencing chemical material shipments—representing raw materials, intermediates, or product batches—across specialized processing, storage, and distribution facilities that possess distinct operational capabilities and constraints. Each shipment must be allocated to a compatible facility and scheduled in an appropriate order such that delivery deadlines, handling requirements, and resource limitations are satisfied while minimizing total operational cost, cumulative delay, and overall schedule length. The system is represented as a heterogeneous graph $G = (S \cup F \cup G, E)$, where S denotes shipment nodes, F denotes facility nodes, G represents global context nodes, and E encodes edges between shipments and compatible facilities, shipments and the global node, and facilities and the global node. This tripartite structure captures both compatibility constraints and system-wide scheduling state. The optimization objective aggregates multiple performance metrics into a composite “badness” function:

$$\begin{aligned} \text{Badness} = & \alpha_1 \cdot N_{\text{delayed}} + \alpha_2 \cdot \text{TotalDelay} \\ & + \alpha_3 \cdot \text{Makespan} + \alpha_4 \cdot \text{TotalCost} + \alpha_5 \cdot \text{TotalDistance} \end{aligned} \quad (1)$$

where N_{delayed} counts late shipments, TotalDelay quantifies cumulative lateness beyond due times, Makespan measures the schedule completion horizon, TotalCost represents facility processing expenses, and TotalDistance reflects aggregate transport distance. The coefficients α_i regulate relative importance, with heavy weighting on delay prevention. The problem is formulated as a Markov Decision Process (MDP) where states encode shipment features (urgency, priority, quantity, due date), facility features (availability, cost, capabilities), and global features (scheduling progress). Actions correspond to selecting the next shipment to schedule, and rewards are derived from incremental reductions in the badness metric, incentivizing efficient sequencing while maintaining operational feasibility. A Graph Neural Network (GNN) processes this heterogeneous graph through iterative message passing:

$$H^{(l+1)} = \sigma \left(\hat{D}^{-1/2} \hat{A} \hat{D}^{-1/2} H^{(l)} W^{(l)} \right)$$

where \hat{A} is the adjacency matrix with self-loops, \hat{D} is the degree matrix, and $W^{(l)}$ are learned parameters. This architecture captures relational dependencies among shipments, facilities, and global context. A policy head produces action probabilities over candidate shipments, while a value head estimates expected returns from global embeddings. The policy is trained via Proximal Policy Optimization (PPO),

which maximizes a clipped surrogate objective:

$$L^{\text{CLIP}} = \mathbb{E} \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right]$$

where $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$ is the probability ratio and \hat{A}_t is the advantage estimate. The total loss integrates policy gradient, value function error, and entropy regularization to balance exploitation, accurate value estimation, and exploration. A Genetic Algorithm (GA) provides an evolutionary baseline, where chromosomes encode shipment sequences and fitness is evaluated using the same composite metric. The GA employs tournament selection, order-preserving crossover, and swap mutation to maintain feasible permutations and promote population diversity. Both PPO-GNN and GA optimizers are benchmarked against heuristic baselines—First-In-First-Out (FIFO), Earliest Due Date (EDD), and Highest Priority First (HPF)—across key performance indicators including makespan, total cost, delayed shipments, total delay, and facility utilization. This integrated framework unifies graph-based relational learning, policy-driven sequential decision-making, and evolutionary optimization, offering a robust solution to the high-dimensional, constraint-rich nature of industrial chemical logistics.

Logistics Optimization

Single-Robot Multi-Order Picking and Placement The proposed framework models single-robot multi-order picking and placement in chemical plants as a Hierarchical Reinforcement Learning (HRL) task. The decision-making process is decomposed into two levels: a high-level policy that prioritizes material transfer requests based on processing priorities and payload capacity limitations, and a low-level policy that executes navigation, pickup, and delivery actions while avoiding obstacles. This temporal abstraction enables the agent to make macro-decisions while executing micro-actions to achieve sub-goals, allowing goal-conditioned behaviors that generalize across varying process complexities. The framework jointly optimizes both levels to maximize operational efficiency—measured by throughput and process completion rates—while respecting payload capacity constraints, processing priorities, and obstacle avoidance requirements. The robotic picking and placement task is modeled as a Hierarchical Markov Decision Process (HMDP) defined by $(\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma))$, where \mathcal{S} is the state space, \mathcal{A} the action space, $P(s'|s, a)$ the transition dynamics, $R(s, a)$ the reward function, and $\gamma \in (0, 1)$ the discount factor. The warehouse is a two-dimensional grid with obstacles, items, and packing stations. Each item $i \in \mathcal{I}$ is located at $l_i = (x_i, y_i)$ with weight w_i , and each order $o \in \mathcal{O}$ contains a set of items \mathcal{I}_o , a packing station l_o^p , and a priority $p_o \in [0, 1]$. A robot located at $l_r = (x_r, y_r)$ has a maximum carrying capacity C_r and a current load $W_r \leq C_r$. The hierarchy separates high-level order selection from low-level motion and manipulation. The high-level state $s_r^{(H)}$ includes order progress, priority, estimated delivery cost, and current load:

$$s_r^{(H)} = \left[(c_o, p_o, \hat{d}_o)_{\forall o \in \mathcal{O}}, \frac{W_r}{C_r} \right],$$

where c_o is the order completion ratio, p_o the normalized priority, and \hat{d}_o the estimated normalized distance. The low-level state $s_r^{(L)}$ captures local context, including the robot's position relative to items and stations, normalized load, and number of items currently carried n_c :

$$s_r^{(L)} = \left[\frac{l_r}{L}, \frac{(l_i - l_r)}{L}, \frac{(l_o^p - l_r)}{L}, \frac{W_r}{C_r}, n_c \right],$$

where L denotes the grid size. The high-level action space determines which order to execute next, while the low-level action space $\mathcal{A}^{(L)} = \{\text{UP, DOWN, LEFT, RIGHT, PICK, DROP}\}$ defines the robot's motion and manipulation primitives. The reward function combines multiple components:

$$R(s, a) = r_{\text{step}} + r_{\text{move}} + r_{\text{pick}} + r_{\text{drop}} + r_{\text{completion}},$$

where $r_{\text{step}} < 0$ penalizes longer task durations, r_{move} rewards movement toward the target and penalizes invalid or backward moves, r_{pick} and r_{drop} reward successful item handling, and $r_{\text{completion}}$ provides a large terminal reward when an order is completed. In capacity-constrained scenarios, the robot can only pick an item if the total weight after pickup remains within its capacity ($W_r + w_i \leq C_r$). When multiple items are delivered in a single trip, the drop reward scales proportionally to the number of items successfully delivered during that delivery cycle. The objective is to learn hierarchical policies $\pi_H(a_H|s_H)$ and $\pi_L(a_L|s_L, g_H)$ that jointly maximize the expected cumulative reward:

$$J(\pi_H, \pi_L) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t R(s_t, a_t) \right].$$

Both controllers employ Deep Q-Learning with prioritized experience replay and soft target updates. The high-level Q-network $Q^{(H)}(s^{(H)}, a^{(H)}; \theta_H)$ selects orders, and the low-level Q-network $Q^{(L)}(s^{(L)}, a^{(L)}; \theta_L)$ handles detailed control. Each minimizes the temporal-difference loss:

$$\mathcal{L}(\theta) = \mathbb{E}_{(s, a, r, s')} \left[(r + \gamma \max_{a'} Q_{\text{target}}(s', a'; \bar{\theta}) - Q(s, a; \theta))^2 \right],$$

where the target parameters are updated via $\bar{\theta} \leftarrow \tau\theta + (1 - \tau)\bar{\theta}$. Prioritized replay samples transitions proportional to TD-error magnitude (α), with importance sampling weights (β) compensating for bias. Exploration follows an epsilon-greedy policy decay, while gradient clipping and periodic target updates stabilize training. Performance is evaluated across episodes using metrics such as order completion rate, average steps per order, total distance traveled, success rate, and cumulative reward. In unconstrained environments, the agent learns efficient sequential routing and order handling. In capacity-constrained environments, the policy learns to plan item pickups based on individual weights, ensuring the total carried load never exceeds the robot's capacity. The resulting framework achieves balanced and adaptive control between strategic order assignment and precise motion execution for efficient warehouse operations.

Multi-Robot Multi-Order Picking and Placement The single-robot formulation is extended to a *multi-robot, multi-order* setting, where multiple robots collaborate within a shared grid-based warehouse containing static obstacles, items, and packing stations. Each robot operates under its own payload capacity constraint and interacts with others through spatial coordination and collision avoidance. The overall objective remains to maximize throughput and order completion rates while ensuring safe and efficient cooperative operation. Unlike the single-robot case, the environment now maintains a set of robots $\mathcal{R} = \{r_1, r_2, \dots, r_K\}$, each characterized by a position $l_{r_k} = (x_{r_k}, y_{r_k})$, a capacity C_{r_k} , and a current load $W_{r_k} \leq C_{r_k}$. The state and action spaces are extended to encode multi-agent interaction, enabling concurrent decision-making and decentralized coordination. The system is modeled as a *Multi-Agent Hierarchical Markov Decision Process (MA-HMDP)*:

$$\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma),$$

where \mathcal{S} and \mathcal{A} represent the *joint* state and action spaces over all robots, $P(s'|s, a)$ captures multi-agent transition dynamics, and $R(s, a)$ defines the cooperative global reward. Each robot's high-level state incorporates not only its own progress and load but also the relative states of other robots:

$$s_{r_k}^{(H)} = \left[(p_o, c_o, \hat{d}_{o,i})_{\forall (o,i)}, \frac{W_{r_k}}{C_{r_k}}, \left(\frac{l_{r_j}}{L} \right)_{\forall j \neq k} \right],$$

allowing spatial awareness and coordination for collision avoidance and distributed task selection. The low-level state further includes the positions of neighboring robots:

$$s_{r_k}^{(L)} = \left[\frac{l_{r_k}}{L}, \frac{(l_i - l_{r_k})}{L}, \frac{(l_o^p - l_{r_k})}{L}, \frac{W_{r_k}}{C_{r_k}}, \left(\frac{l_{r_j} - l_{r_k}}{L} \right)_{\forall j \neq k} \right],$$

enabling each robot to plan collision-free paths while navigating toward its assigned targets. All robots act concurrently at each time step, forming a joint action set $\mathcal{A} = \{a_{r_1}, a_{r_2}, \dots, a_{r_K}\}$, where each $a_{r_k} \in \mathcal{A}^{(L)} = \{\text{UP, DOWN, LEFT, RIGHT, PICK, DROP}\}$. The environment enforces mutual exclusion such that no two robots can occupy the same grid cell, and collisions with obstacles or other robots incur penalties. The global reward aggregates individual contributions from all robots:

$$R(s, a) = \sum_{r_k \in \mathcal{R}} \left(r_{\text{step}} + r_{\text{move}} + r_{\text{pick}} + r_{\text{drop}} + r_{\text{completion}} + r_{\text{collision}} \right),$$

where $r_{\text{collision}} < 0$ penalizes robot-robot and robot-obstacle collisions. In the capacity-constrained setting, each robot can pick an item only if ($W_{r_k} + w_i \leq C_{r_k}$); otherwise, a penalty is applied, encouraging load-feasible item selection and implicit cooperation among robots through task redistribution. All robots share parameters for both the high-level and low-level Q-networks, enabling decentralized execution with centralized training. The learning objective remains:

$$J(\pi_H, \pi_L) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t R(s_t, a_t) \right],$$

but gradient updates are performed using experience tuples aggregated from all robots' interactions, promoting cooperative learning through shared replay buffers. Performance is evaluated through multi-agent metrics such as mean order completion rate, total travel distance, collision frequency, and cumulative reward. In unconstrained environments, robots learn efficient division of labor and spatial coordination, while in capacity-constrained scenarios, they develop load-aware task allocation strategies that maintain balanced utilization and collision-free cooperation.

Bin Packing Optimization for Chemical Plant Logistics

In large-scale chemical plant logistics and warehouse operations, efficient packing and shipment planning are vital for minimizing handling costs, maximizing space utilization, and maintaining safety compliance. The task involves selecting suitable containers and arranging diverse items under strict geometric, weight, and stacking constraints. To address this challenge, we develop a Reinforcement Learning (RL)-based bin packing framework that autonomously learns efficient packing strategies. The framework models the problem as a Markov Decision Process (MDP) and employs a Dueling Deep Q-Network (DQN) with Prioritized Experience Replay (PER) to optimize packing efficiency and cost under realistic industrial conditions. The framework seeks to balance cost, utilization, and constraint compliance through a multi-objective reward formulation. The overall objective is expressed as

$$\max_{\pi} J(\pi) = \mathbb{E}_{\pi}[w_1(-C_{\text{total}}) + w_2 U_{\text{avg}} - w_3 N_b - w_4 \mathcal{V}_{\text{constraint}}],$$

where C_{total} denotes the total logistics cost, U_{avg} represents the mean container utilization, N_b is the number of boxes or tanks used, and $\mathcal{V}_{\text{constraint}}$ measures the cumulative constraint violation. The coefficients w_1, w_2, w_3, w_4 are user-defined weights that adjust the trade-off between cost efficiency, space utilization, container minimization, and safety compliance. The environment is modeled as a Markov Decision Process $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} the action set, \mathcal{T} the transition function, \mathcal{R} the reward function, and γ the discount factor. Each state $s_t \in \mathcal{S}$ encodes the condition of all available containers, the packing configuration, and the items yet to be placed. The state vector is represented as

$$s_t = [x_{\text{avail}}, u_i, \frac{w_i}{w_i^{\text{max}}}, \text{fit}_i, \frac{c_i}{c_{\text{div}}}, \frac{V_i}{V_{\text{div}}}],$$

where x_{avail} indicates whether an item remains unpacked, u_i is the utilization ratio of container i , and $\frac{w_i}{w_i^{\text{max}}}$ expresses the packed-to-maximum weight ratio. The binary variable $\text{fit}_i \in \{0, 1\}$ denotes if the next item can fit in container i , while c_i and V_i represent its cost and volume, normalized by divisors c_{div} and V_{div} for scale stability. The action $a_t \in \mathcal{A}$ selects an item for placement into an existing or new container, and the transition $\mathcal{T}(s_t, a_t)$ updates the state according to geometric fit and residual space, giving $s_{t+1} = f(s_t, a_t)$. The immediate reward r_t balances cost and utilization:

$$r_t = \begin{cases} -\lambda_c C_b, & \text{if a new container is opened,} \\ \lambda_v \frac{V_i}{V_{\text{div}}} + \lambda_z \left(1 - \frac{z_i}{H_b}\right), & \text{if an item is placed successfully,} \\ -\lambda_f, & \text{otherwise.} \end{cases}$$

Here, C_b is the container cost, V_i the item volume, and z_i its placement height within a container of height H_b . The coefficients $\lambda_c, \lambda_v, \lambda_z, \lambda_f$ regulate penalties for cost, spatial efficiency, height placement, and failures. The cumulative episodic reward sums all r_t values and applies terminal bonuses when utilization and cost targets are met. Physical feasibility is ensured by the constraints

$$\sum_j w_j \leq W_b^{\text{max}}, \quad \sum_j V_j \leq V_b,$$

which prevent weight and volume overflow. For fragile items, if an object below is fragile, no heavier object may be placed above it, expressed as

$$\text{if fragile}(\text{below}) = 1 \Rightarrow w_{\text{above}} = 0.$$

Rotations are allowed only for items with flexible orientations, i.e., $r_i \in \text{Rotations}(d_i)$ when the orientation flag is unset. The overall optimization objective is

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^T \gamma^t r_t \right], \quad \text{subject to } \mathcal{C}(s_t, a_t) = 0,$$

where π^* is the optimal packing policy and \mathcal{C} encodes geometric and safety constraints. The agent employs a Dueling Deep Q-Network (DQN) that decomposes the Q-value into state value $V(s_t)$ and advantage $A(s_t, a_t)$:

$$Q(s_t, a_t) = V(s_t) + \left(A(s_t, a_t) - \frac{1}{|\mathcal{A}|} \sum_{a'} A(s_t, a') \right).$$

The target update is

$$y_t = r_t + \gamma \max_{a'} Q(s_{t+1}, a'; \theta^-),$$

and the loss minimized is a weighted Huber loss,

$$\mathcal{L} = \mathbb{E}[w_i \cdot \text{Huber}(Q(s, a; \theta) - y_t)],$$

where w_i is the importance-sampling weight correcting prioritized replay bias. Gradient clipping and scheduled learning-rate decay stabilize convergence. After training, the learned policy π^* is evaluated without exploration noise. Performance is compared with heuristic baselines commonly used in logistics: First Fit Decreasing (FFD), Best Fit Decreasing (BFD), and a Genetic Algorithm (GA). The FFD method sequentially assigns items to the first feasible container, while BFD minimizes residual volume after placement. The GA baseline evolves item orderings using a fitness function

$$f = C_{\text{total}} + \alpha(1 - U_{\text{avg}}) + \beta N_{\text{unpacked}},$$

where C_{total} is the total cost, U_{avg} the average utilization, and N_{unpacked} the count of unplaced items. A solution is accepted if

$$N_b \leq N_b^{\text{max}}, \quad C_{\text{total}} \in [C_{\text{min}}, C_{\text{max}}], \quad U_{\text{avg}} \geq U_{\text{min}}.$$

Empirical evaluation shows that the RL agent achieves higher utilization and lower cost than heuristic methods by learning spatially aware and cost-efficient packing strategies that generalize across diverse container and item configurations in chemical logistics operations.

Dynamic Pricing and Revenue Optimization for Process Plants In chemical production plants, determining optimal product pricing is a complex decision-making task influenced by volatile market demand, fluctuating raw material and energy costs, and evolving inventory conditions. Conventional static or cost-plus pricing strategies remain insensitive to market variability, motivating a reinforcement learning (RL) formulation that learns adaptive, data-driven pricing decisions. The pricing process is formulated as a sequential decision-making problem using Proximal Policy Optimization (PPO) to learn policies that maximize long-term profit while ensuring operational and inventory compliance. The policy $\pi_\theta(a_t|s_t)$ determines the optimal daily selling price a_t given the current system state s_t , balancing short-term profitability with long-term operational stability. The PPO-based agent exhibits dynamic, context-aware behavior: when production costs rise, it increases prices to preserve profit margins, while lower costs incentivize price reductions to stimulate demand. Under favorable hedging conditions—when futures prices fall below spot feedstock costs—the agent maintains or slightly decreases prices to expand market share, whereas unfavorable conditions prompt higher prices to safeguard profitability. The pricing process is modeled as a Markov Decision Process (MDP) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$, where \mathcal{S} is the continuous state space encoding operational and market variables (forecasted demand, inventory, production costs, feedstock and energy prices, temporal features), \mathcal{A} is the discrete action space of feasible price levels, P represents transition dynamics, R is the reward function, and $\gamma \in (0, 1)$ is the discount factor for future rewards. The action space $\mathcal{A} = \{p_0, p_1, \dots, p_K\}$ represents discrete price levels centered around a nominal base price p_{base} with markup factors μ_k : $p_k = p_{\text{base}}(1 + \mu_k)$. The policy $\pi_\theta(a_t|s_t)$ determines the optimal daily selling price a_t given the current system state s_t . The immediate reward is the normalized profit margin:

$$r_t = \frac{1}{\eta} (R_t - C_t - P_t^{\text{inv}}),$$

where R_t is total revenue, C_t is total cost, P_t^{inv} is an inventory penalty, and η is a scaling factor for numerical stability. Daily revenue combines spot and contract sales:

$$R_t = Q_t^{\text{spot}} p_t + Q_t^{\text{contract}} p_c,$$

where $p_c = p_{\text{base}}(1 - \delta)$ is the discounted contract price with discount rate δ . Spot demand follows a price-elastic relationship:

$$Q_t^{\text{spot}} = (D_t(1 - \rho_c)) \left(\frac{p_t}{p_{\text{base}}} \right)^{-\epsilon},$$

where D_t is predicted market demand, $\rho_c \in [0, 1]$ is the contractual demand fraction, and $\epsilon > 0$ is the price elasticity coefficient (higher prices exponentially reduce demand). Total cost combines material, catalyst, energy, and hedging costs:

$$C_t = Q_t^{\text{total}} (c_t^{\text{raw}} + c^{\text{add}} + c_t^{\text{energy}}) - H_t + Q_t^{\text{total}} c^{\text{hedge}},$$

where $Q_t^{\text{total}} = Q_t^{\text{spot}} + Q_t^{\text{contract}}$ is total sales volume. The hedging benefit $H_t = (c_t^{\text{raw}} - f_t^{\text{futures}})\psi Q_t^{\text{total}}$ is realized when futures price f_t^{futures} is below the raw-material price c_t^{raw} , with $\psi \in [0, 1]$ representing hedge effectiveness. Inventory evolves via a mass-balance equation:

$$I_{t+1} = I_t + P^{\text{rate}} - Q_t^{\text{total}},$$

where P^{rate} is the fixed daily production rate. Inventory deviations beyond operational thresholds incur penalties:

$$P_t^{\text{inv}} = \begin{cases} \alpha_{\text{low}}(I_{\text{low}} - I_t), & \text{if } I_t < I_{\text{low}} \text{ (stock-out risk),} \\ \alpha_{\text{high}}(I_t - I_{\text{max}})\phi, & \text{if } I_t > I_{\text{max}}\phi \text{ (storage constraint),} \\ 0, & \text{otherwise,} \end{cases}$$

where I_{low} is minimum safety stock, I_{max} is storage capacity, $\phi \in (0, 1)$ controls soft upper limits, and $\alpha_{\text{low}}, \alpha_{\text{high}} > 0$ are penalty weights. The learning objective maximizes expected cumulative discounted reward:

$$J(\pi_\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \gamma^t r_t \right].$$

The training pipeline begins with feature engineering: cost, demand, and energy time series are aggregated across temporal windows, computing statistical descriptors (mean, variance, extrema, polynomial trends) over multiple lags. A gradient boosting regression model $\hat{D}_t = f_\phi(s_t)$ forecasts short-term demand. The RL agent employs PPO with an actor-critic architecture. The actor outputs a categorical distribution over price levels; the critic estimates the state-value function $V_\theta(s_t)$. PPO maximizes the clipped surrogate objective:

$$L^{\text{PPO}}(\theta) = \mathbb{E}_t \left[\min \left(\rho_t(\theta) \hat{A}_t, \text{clip}(\rho_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right],$$

where $\rho_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$ is the probability ratio and \hat{A}_t is the generalized advantage estimate computed via temporal-difference updates with decay parameter λ . The critic minimizes value loss $L_V = \frac{1}{2} \|V_\theta(s_t) - \hat{R}_t\|^2$. Gradient clipping ensures numerical stability during training across multiple episodes with randomized starting indices. The PPO-based agent exhibits dynamic, context-aware behavior: when production costs rise, it increases prices to preserve margins; lower costs prompt price reductions to stimulate demand. Under favorable hedging (futures below spot costs), the

agent maintains or reduces prices to expand market share; unfavorable conditions prompt higher prices to safeguard profitability. For evaluation, three strategies—Static, Cost-Plus, and PPO-based Adaptive Pricing—are benchmarked over a full operational horizon using cumulative profit, average daily margin, and average realized price. The PPO-based agent learns to dynamically adjust pricing according to market elasticity and inventory conditions, lowering prices under oversupply and raising them during scarcity. This adaptive policy achieves superior long-term profitability and smoother inventory utilization compared to static and heuristic benchmarks.

Inventory Optimization System for Chemical Manufacturing Industrial manufacturing plants must determine optimal production rates that balance inventory holding costs, demand fulfillment, and operational expenses under uncertainties from market fluctuations and production variability. Maintaining equilibrium between production, consumption, and storage constitutes a complex control problem. To address this challenge, the task is formulated as a sequential decision-making problem where the objective is to learn adaptive control policies that minimize long-term operational costs while ensuring compliance with production and inventory constraints. A reinforcement learning (RL) framework is employed, integrating a data-driven world model that captures production-demand dynamics with a policy optimization agent that learns optimal production rate decisions. This problem is formulated as a Markov Decision Process (MDP) defined by $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$, where \mathcal{S} denotes the state space comprising inventory levels, production rates, and demand features, \mathcal{A} represents the discrete action space of feasible production rate setpoints, $P(s_{t+1}|s_t, a_t)$ captures transition dynamics, and $R(s_t, a_t)$ defines the immediate reward. A reinforcement learning framework integrates a data-driven world model that captures production-demand dynamics with a policy optimization agent to enable safe training in simulation before real-world deployment. The objective is to derive an optimal control policy $\pi_\theta(a_t|s_t)$, parameterized by θ , that determines production decision a_t given system state s_t to maximize the long-term expected return:

$$J(\pi_\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \gamma^t r(s_t, a_t) \right],$$

where $r(s_t, a_t)$ represents the immediate reward (negative of total operating cost), $\gamma \in [0, 1]$ is the discount factor balancing short- and long-term objectives, and T is the planning horizon. The reward function penalizes four cost components:

$$r_t = -(c_h I_t + c_u U_t + c_p P_t + c_r R_t),$$

where I_t , U_t , P_t , and R_t correspond to inventory holding, unmet demand, production, and rate-change costs, with $c_h, c_u, c_p, c_r > 0$ representing their respective cost coefficients. The process dynamics evolve according to $s_{t+1} = f(s_t, a_t; \omega) + \epsilon_t$, where $f(\cdot)$ denotes the underlying system evolution function parameterized by ω and ϵ_t represents

stochastic process noise. The training framework adopts a hybrid model-based reinforcement learning approach, combining supervised learning for process prediction with policy optimization via Proximal Policy Optimization (PPO). The predictive world model approximates the process dynamics using a function approximator \hat{f}_ω trained on historical observations. Given an input feature vector $\mathbf{x}_t \in \mathbb{R}^n$ constructed from lagged demand, production rates, catalyst health, and time-dependent features, the model predicts the next system state y_t by minimizing mean-squared error:

$$\omega^* = \arg \min_{\omega} \frac{1}{N} \sum_{t=1}^N \|y_t - \hat{f}_\omega(\mathbf{x}_t)\|^2.$$

This trained model enables the creation of a simulated environment in which the PPO agent learns robust control policies efficiently without requiring extensive real-world interactions. The PPO algorithm refines the control policy by constraining policy updates for training stability through its clipped surrogate objective:

$$L^{\text{PPO}}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right],$$

where $r_t(\theta) = \pi_\theta(a_t|s_t)/\pi_{\theta_{\text{old}}}(a_t|s_t)$ is the probability ratio between the updated and previous policies, and ϵ controls the extent of policy deviation. The advantage term \hat{A}_t is computed using Generalized Advantage Estimation (GAE):

$$\hat{A}_t = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}, \quad \text{where} \quad \delta_t = r_t + \gamma V_\theta(s_{t+1}) - V_\theta(s_t),$$

with $V_\theta(s_t)$ being the critic's value function and $\lambda \in [0, 1]$ balancing bias and variance. The total loss combines actor, critic, and entropy components to balance exploitation, value estimation accuracy, and exploration. The evaluation phase benchmarks multiple decision strategies—including reactive, heuristic, and reinforcement learning-based agents—under consistent simulated demand conditions across a simulation horizon. Each policy's cumulative performance is evaluated using the total cost function:

$$C_{\text{total}} = \sum_{t=0}^T (c_h I_t + c_u U_t + c_p P_t + c_r R_t),$$

with comparisons made in terms of total cost, cost breakdown, and stability of control decisions, demonstrating that the optimized control policy achieves reduced operating costs and improved dynamic stability compared to baseline strategies.

Risk-Aware Hedging and Portfolio Optimization In large-scale chemical process plants, profitability is strongly influenced by fluctuations in feedstock costs, energy prices, and product market values. Chemical plants rely on volatile commodity inputs and sell products in fluctuating markets. These price variations create significant financial risk, where

sudden market changes can erode profit margins despite stable operations. To mitigate this risk, plants employ hedging strategies—securing future feedstock prices through forward contracts (private agreements to buy or sell at fixed prices on a future date) or futures contracts (standardized, exchange-traded versions with daily settlement to reduce counterparty risk)—and portfolio optimization, which allocates production capacity among multiple products to maximize expected profit while maintaining acceptable risk levels. Together, hedging and portfolio optimization form a dynamic financial–operational decision-making framework that couples production planning with market uncertainty management. The problem is formulated as a Markov Decision Process (MDP) defined by $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$, where \mathcal{S} is the state space representing market features (commodity prices, trends, volatility), \mathcal{A} is the action space encoding hedging ratios and portfolio allocations, $P(s_{t+1}|s_t, a_t)$ represents transition dynamics, $R(s_t, a_t)$ is the reward function, and $\gamma \in (0, 1)$ is the discount factor. At each time step t , the agent observes state s_t , selects action $a_t \sim \pi_\theta(a_t|s_t)$ from a policy parameterized by θ , receives reward $R(s_t, a_t)$, and transitions to s_{t+1} . The objective is to learn an optimal policy π_θ^* maximizing the expected cumulative risk-adjusted profit:

$$J(\pi_\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \gamma^t \left(\frac{\Pi_t}{\alpha} - \lambda \frac{\sigma_t}{\alpha} \right) \right],$$

where Π_t is instantaneous profit, σ_t is the rolling standard deviation of profits capturing volatility, λ is the risk-aversion coefficient, and α is a normalization constant. The first term encourages profit maximization while the second penalizes volatility. Instantaneous profit is

$$\Pi_t = \sum_p Q_{p,t} P_{p,t} - \sum_i C_{i,t} Q_{i,t} + H_t - F,$$

where $Q_{p,t}$ and $P_{p,t}$ are production quantity and price of product p , $C_{i,t}$ and $Q_{i,t}$ are feedstock prices and consumption of input i , H_t is hedge gain or loss, and F is fixed operating costs. A Proximal Policy Optimization (PPO) algorithm learns this optimal policy by iteratively updating parameters using clipped surrogate objectives for stable learning. The PPO objective is

$$L^{\text{PPO}}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) - c_v (V_\theta(s_t) - \hat{R}_t)^2 + c_e \mathcal{H}(\pi_\theta(\cdot|s_t)) \right],$$

where $r_t(\theta) = \pi_\theta(a_t|s_t)/\pi_{\theta_{\text{old}}}(a_t|s_t)$ is the probability ratio between new and old policies, \hat{A}_t is the advantage estimate computed via Generalized Advantage Estimation (GAE), $V_\theta(s_t)$ is the critic’s value estimate, \hat{R}_t is the empirical return, and \mathcal{H} is the policy entropy weighted by c_e for exploration. The clipping parameter ϵ prevents destabilizing policy updates. During training, the agent interacts with simulated market environments generated from

historical data or Monte Carlo scenarios, learning to map price trajectories into optimal decisions. Feature engineering transforms raw market data into lagged and statistical representations—rolling means, variances, and polynomial trends—capturing temporal dependencies. For evaluation, the trained agent is benchmarked against baseline strategies under out-of-sample scenarios. Performance is assessed using expected total profit, standard deviation of returns, Value-at-Risk (VaR), and Profit-at-Risk (PaR):

$$\text{VaR}_\alpha = \inf \{x \mid P(\Pi < x) \geq \alpha\}, \quad \text{PaR}_\alpha = \mathbb{E}[\Pi] - \text{VaR}_\alpha,$$

where VaR measures the minimum loss at confidence level α and PaR quantifies expected shortfall beyond VaR. A superior policy achieves higher average profit with lower volatility and tighter tail risk, demonstrating robust performance under stochastic market dynamics through adaptive, risk-aware control.

Chemical Industrial Demand Forecasting and Production Planning Industrial manufacturing requires integrated demand forecasting and production planning to determine optimal daily production rates under uncertainty. We formulate this as a Markov Decision Process and develop a Proximal Policy Optimization (PPO) agent that maximizes long-term profitability while simultaneously managing inventory, maintaining product quality specifications, ensuring environmental compliance below regulatory thresholds, and stabilizing operations by minimizing production rate changes. The system addresses fluctuating market demand, variable raw material and energy costs, maintenance schedules, and stochastic process dynamics affecting quality and environmental performance. The production planning task is modeled as a Markov Decision Process (MDP) defined by

$$\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma),$$

where \mathcal{S} denotes the system state (e.g., production rate, inventory level, energy use, product quality, and emissions), \mathcal{A} represents feasible production setpoints, $P(s_{t+1}|s_t, a_t)$ captures stochastic process transitions, $R(s_t, a_t)$ is the immediate reward, and $\gamma \in [0, 1)$ is the discount factor weighting long-term profitability over short-term gain. The objective is to learn an optimal policy $\pi^*(a_t|s_t)$ that maximizes the expected cumulative discounted reward:

$$J(\pi) = \mathbb{E}_\pi \left[\sum_{t=0}^T \gamma^t R(s_t, a_t) \right].$$

The reward function jointly optimizes economic performance and operational compliance:

$$R(s_t, a_t) = \Pi(s_t, a_t) - \mathcal{P}_{\text{quality}} - \mathcal{P}_{\text{emissions}} - \mathcal{P}_{\text{inventory}} - \mathcal{P}_{\text{stability}},$$

where $\Pi(s_t, a_t)$ denotes net profit (revenue minus feedstock, catalyst, energy, and storage costs), and the penalty terms $\mathcal{P}_{\text{quality}}$, $\mathcal{P}_{\text{emissions}}$, $\mathcal{P}_{\text{inventory}}$, and $\mathcal{P}_{\text{stability}}$ represent losses due to quality deviations, emission violations, inventory imbalance, and excessive production-rate fluctuations, respectively. The training process is divided into two tightly coupled phases: supervised forecasting and reinforcement learning-based optimization. The forecasting model $\hat{f}(\cdot)$ approximates multi-step-ahead predictions of demand, inventory,

product quality, and emissions using regression over engineered time-series features, thus serving as a differentiable surrogate environment. The RL agent interacts with this surrogate to learn a policy $\pi_\theta(a_t|s_t)$ that generalizes across dynamic operating conditions. Policy optimization is performed using Proximal Policy Optimization (PPO), an actor–critic algorithm that maximizes a clipped surrogate objective:

$$L^{\text{PPO}}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) - c_v (V_\phi(s_t) - \hat{R}_t)^2 + c_e \mathcal{H}(\pi_\theta(\cdot|s_t)) \right].$$

where $r_t(\theta)$ is the ratio between new and old policy probabilities, \hat{A}_t is the advantage estimate from Generalized Advantage Estimation (GAE), and the last two terms represent value-function regularization and entropy maximization to balance learning stability and exploration. During training, the agent iteratively samples trajectories from the simulated environment, computes discounted returns and advantages, and updates the actor–critic parameters through mini-batch gradient ascent. Normalization and gradient clipping are applied to ensure numerical stability and prevent policy divergence. The surrogate forecasting model provides rapid feedback for evaluating the long-term impact of production decisions, enabling efficient off-policy learning without interacting with the real plant. The trained policy is benchmarked against baseline strategies such as static and reactive controllers using metrics including cumulative profit, on-spec production rate, environmental violations, and operational smoothness. Superiority is defined by consistently achieving higher profitability while maintaining compliance, demonstrating that the integrated forecasting–reinforcement approach enables adaptive and sustainable decision-making for industrial systems under uncertainty.

Physics-Informed Neural Networks (PINNs)

Transformer Neural Operator (TNO) The *Transformer Neural Operator (TNO)* provides a general framework for learning mappings between functional spaces that describe systems governed by ordinary differential equations (ODEs) and partial differential equations (PDEs) across spatial and temporal domains. Let the overall domain be denoted as $\Omega_t = \Omega \times [0, T]$, where $\Omega \subseteq \mathbb{R}^d$ represents the spatial region of dimension $d \in \{0, 1, 2, 3\}$, and $t \in [0, T]$ denotes time. The governing physical process is represented by a generalized differential operator \mathcal{F} acting on the state variable $u(\mathbf{x}, t)$, defined as

$$\mathcal{F}(\mathbf{x}, t, u, \nabla u, \nabla^2 u, \frac{\partial u}{\partial t}, a(\mathbf{x}, t), p) = 0, \quad (\mathbf{x}, t) \in \Omega_t.$$

Here, $u(\mathbf{x}, t)$ is the *state field or solution variable* (e.g., temperature, pressure, velocity, or concentration), $a(\mathbf{x}, t)$ is the *input or forcing field* representing external influences such as heat flux or source terms, and $p \in \mathbb{R}^k$ denotes *physical parameters* such as diffusivities, reaction constants, or conductivities. The solution must satisfy boundary and initial conditions,

$$\mathcal{B}(u) = 0 \quad \text{on } \partial\Omega, \quad \mathcal{I}(u(\mathbf{x}, 0)) = u_0(\mathbf{x}) \quad \text{in } \Omega,$$

where $\mathcal{B}(\cdot)$ and $\mathcal{I}(\cdot)$ denote the boundary and initial operators respectively. This unified formulation naturally covers ODEs when $d = 0$, steady-state PDEs when $\partial u / \partial t = 0$, and time-dependent PDEs when both spatial and temporal derivatives are active. The goal is to approximate the true solution operator $\mathcal{G} : (a, b, p) \mapsto u$ with a learnable neural operator \mathcal{G}_θ , parameterized by θ , by minimizing a *physics-informed loss* that enforces the governing equations and auxiliary conditions. The optimization problem is

$$\min_{\theta} \mathbb{E}_{(a,b,p) \sim \mathcal{D}} \left[w_{\text{eq}} \|\mathcal{F}(\mathcal{G}_\theta(a, b, p))\|_2^2 + w_{\text{bc}} \|\mathcal{B}(\mathcal{G}_\theta(a, b, p))\|_2^2 + w_{\text{ic}} \|\mathcal{I}(\mathcal{G}_\theta(a, b, p))\|_2^2 \right],$$

where \mathcal{D} denotes the distribution of admissible inputs (a, b, p) , and $w_{\text{eq}}, w_{\text{bc}}, w_{\text{ic}}$ are weights balancing the PDE residual, boundary, and initial conditions. The optimal parameters are obtained by

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{(a,b,p) \sim \mathcal{D}} [\mathcal{L}_{\text{res}}(\theta; a, b, p)],$$

where the total residual loss is

$$\mathcal{L}_{\text{res}} = w_{\text{eq}} \mathcal{L}_{\text{eq}} + w_{\text{bc}} \mathcal{L}_{\text{bc}} + w_{\text{ic}} \mathcal{L}_{\text{ic}}.$$

The individual components are defined as

$$\begin{aligned} \mathcal{L}_{\text{eq}} &= \frac{1}{N} \sum_{i=1}^N |\mathcal{R}(\hat{u}_\theta(\mathbf{x}_i, t_i), a(\mathbf{x}_i, t_i), p)|^2, \\ \mathcal{L}_{\text{bc}} &= \frac{1}{N_b} \sum_{\mathbf{x}_j \in \partial\Omega} |\hat{u}_\theta(\mathbf{x}_j, t_j) - b(\mathbf{x}_j, t_j)|^2, \\ \mathcal{L}_{\text{ic}} &= \frac{1}{N_0} \sum_{\mathbf{x}_k \in \Omega} |\hat{u}_\theta(\mathbf{x}_k, 0) - u_0(\mathbf{x}_k)|^2, \end{aligned}$$

where the residual operator \mathcal{R} is expressed generically as

$$\mathcal{R}(\mathbf{x}, t) = \frac{\partial u}{\partial t} + \mathcal{L}(u, a, p),$$

$$\mathcal{L}(u, a, p) = \Phi(u, \nabla u, \nabla^2 u, a(\mathbf{x}, t), p),$$

The spatio-temporal domain is discretized into grid points (\mathbf{x}_i, t_j) , and at each location an embedding is constructed as

$$h_{i,j}^{(0)} = \text{Linear}([a(\mathbf{x}_i, t_j), \mathbf{x}_i, b(\mathbf{x}_i, t_j), p]) + \text{PE}(\mathbf{x}_i, t_j),$$

where $\text{PE}(\cdot)$ denotes the positional encoding that injects spatial and temporal context additively. These embeddings are propagated through Transformer layers composed of multi-head self-attention (MHA) and feed-forward (FFN) sub-blocks:

$$h_{i,j}^{(l)} = \text{MHA}(h_{i,j}^{(l-1)}) + \text{FFN}(h_{i,j}^{(l-1)}), \quad l = 1, \dots, L.$$

The predicted solution field is then obtained as $\hat{u}_\theta(\mathbf{x}, t) = \text{Linear}(h_{i,j}^{(L)})$. For time-dependent systems where the dynamics follow $\frac{\partial u}{\partial t} + \mathcal{L}(u, a, p) = 0$, an additional temporal residual term ensures consistency across discrete timesteps:

$$\mathcal{L}_{\text{time}} = \frac{1}{N_t} \sum_{j=1}^{N_t} \left\| \frac{\hat{u}_\theta^{t_{j+1}} - \hat{u}_\theta^{t_j}}{\Delta t} + \mathcal{L}(\hat{u}_\theta^{t_j}, a^{t_j}, p) \right\|_2^2.$$

This formulation allows the TNO to represent steady, transient, and coupled multi-physics systems within a unified operator-learning framework. The model parameters θ are optimized via physics-constrained training without ground-truth data, and the generalization accuracy is quantified using the relative L_2 error

$$\epsilon_{\text{rel}} = \frac{\|\hat{u}_\theta - u^*\|_2}{\|u^*\|_2}.$$

The TNO thus provides a dimension-agnostic, physics-informed, and time-aware framework capable of solving ODEs and PDEs across arbitrary spatial dimensions under a single unified formulation.

Physics-Informed Graph Neural Operator (PIGO) The *Physics-Informed Graph Neural Operator (PIGO)* extends the neural operator formulation by representing the spatial region $\Omega \subseteq \mathbb{R}^d$ as a directed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where nodes $\mathcal{V} = \{v_i\}_{i=1}^{N_s}$ correspond to discrete spatial coordinates \mathbf{x}_i , and edges \mathcal{E} define neighborhood relations among them. Each node encodes a feature embedding comprising the local state field $u(\mathbf{x}, t)$, the input variable $a(\mathbf{x}, t)$, and the physical parameters $p \in \mathbb{R}^k$. The edge set \mathcal{E} is constructed from spatial proximity, where each node v_i connects to its local neighborhood (e.g., nearest neighbors in Ω), defining the graph topology. The message-passing operation defines a learnable discrete operator through neighborhood aggregation, expressed as

$$h_i^{(l+1)} = \text{LayerNorm} \left(h_i^{(l)} + \sum_{j \in \mathcal{N}(i)} \phi_\theta^{(l)}([h_i^{(l)}, h_j^{(l)}]) \right),$$

where $h_i^{(l)}$ denotes the node feature at layer l , $\mathcal{N}(i)$ represents the neighborhood of node v_i , and $\phi_\theta^{(l)}(\cdot)$ is a learnable function that captures information exchange between connected nodes. The message function additionally depends on edge attributes constructed from relative spatial offsets $(\mathbf{x}_j - \mathbf{x}_i)$, enabling the network to incorporate geometric information. This operation provides a graph-based approximation to spatial derivatives such as ∇u , allowing the model to learn non-local dependencies through message propagation without relying on structured grids or explicit finite-difference schemes. The operator learns a functional mapping $\mathcal{G}_\theta : (a, b, p) \mapsto u$, where $a(\mathbf{x}, t)$ denotes the input

or forcing field, $b(\mathbf{x}, t)$ represents boundary or initial conditions, and $p \in \mathbb{R}^k$ contains physical or material parameters. Each node-time pair (\mathbf{x}_i, t_j) is embedded as

$$h_{i,j}^{(0)} = \text{Linear}([a(\mathbf{x}_i, t_j), b(\mathbf{x}_i, t_j), p]),$$

where the embeddings are propagated through successive graph layers, and the predicted solution field is obtained as $\hat{u}_\theta(\mathbf{x}, t) = \text{Linear}(h_{i,j}^{(L)})$. The optimization objective minimizes the residuals of the governing differential operator and associated boundary and initial conditions over the distribution of admissible inputs,

$$\min_{\theta} \mathbb{E}_{(a,b,p) \sim \mathcal{D}} \left[w_{\text{eq}} \|\mathcal{F}(\mathcal{G}_\theta(a, b, p))\|_2^2 + w_{\text{bc}} \|\mathcal{B}(\mathcal{G}_\theta(a, b, p))\|_2^2 + w_{\text{ic}} \|\mathcal{I}(\mathcal{G}_\theta(a, b, p))\|_2^2 \right],$$

where $\mathcal{F}(\cdot)$ denotes the residual of the governing PDE, and $\mathcal{B}(\cdot)$ and $\mathcal{I}(\cdot)$ represent boundary and initial operators, respectively. Spatial derivatives required for the residual evaluation are obtained through automatic differentiation with respect to spatial coordinates, enabling end-to-end training without discretizing the derivative operators. The model parameters θ are optimized using AdamW with cosine-annealed learning rate scheduling, gradient clipping, and warm-up epochs to ensure convergence stability. The learning process is fully guided by the physics-constrained residual losses without requiring labeled data. The predictive accuracy of the learned operator is quantified by the relative L_2 error,

$$\epsilon_{\text{rel}} = \frac{\|\hat{u}_\theta - u^*\|_2}{\|u^*\|_2},$$

where u^* represents the reference solution obtained from a numerical solver. The Physics-Informed Graph Neural Operator provides a dimension-agnostic, physics-constrained, and time-aware operator-learning paradigm capable of solving ODEs, PDEs, and coupled multi-physics systems within a single unified formulation.

Domain-Decomposed Physics-Informed Neural Network (DD-PINN) A *Domain-Decomposed Physics-Informed Neural Network (DD-PINN)* splits the domain into smaller regions, where each subnetwork learns the local physics, and continuity across interfaces ensures a globally consistent solution. The overall spatio-temporal domain $\Omega_t = \Omega \times [0, T]$ is partitioned into subdomains $\{\Omega_t^{(m)}\}_{m=1}^M$, each governed by a localized neural operator $\mathcal{N}_{\theta^{(m)}}$. Each subnetwork $\mathcal{N}_{\theta^{(m)}}$ outputs a predicted solution field $\hat{u}_{\theta^{(m)}}(\mathbf{x}, t)$, approximating the true physical state field $u(\mathbf{x}, t)$ within its respective subdomain $\Omega_t^{(m)}$. The local models learn the solution by satisfying the governing operator within each subdomain and enforcing continuity at shared interfaces, enabling scalable and stable learning for long-horizon, stiff, and multi-scale systems. Each subdomain is defined as

$$\Omega_t^{(m)} = [t_{m-1}, t_m] \times \Omega, \quad \Omega_t = \bigcup_{m=1}^M \Omega_t^{(m)}, \quad \Omega_t^{(m)} \cap \Omega_t^{(m+1)} = \Gamma_t^{(m)}$$

where $\Gamma_t^{(m)}$ denotes the shared interface between adjacent regions. Within each subdomain, the governing physical law is expressed as

$$\mathcal{F}^{(m)}(\mathbf{x}, t, u, \nabla u, \nabla^2 u, a(\mathbf{x}, t), p^{(m)}) = 0, \quad (\mathbf{x}, t) \in \Omega_t^{(m)}.$$

Here, $u(\mathbf{x}, t)$ is the state field (e.g., temperature, concentration, or pressure), $a(\mathbf{x}, t)$ represents the input or forcing field, and $p^{(m)}$ denotes physical parameters associated with the m^{th} subdomain. The DD-PINN jointly optimizes all subnetworks through a composite physics-constrained loss enforcing local equation satisfaction, initial or boundary compliance, and interface continuity:

$$\begin{aligned} \min_{\{\theta^{(m)}\}_{m=1}^M} \sum_{m=1}^M \mathbb{E}_{(a,b,p) \sim \mathcal{D}^{(m)}} & \left[w_{\text{eq}} \left\| \mathcal{F}^{(m)}(\mathcal{N}_{\theta^{(m)}}(a, b, p)) \right\|_2^2 \right. \\ & + w_{\text{ic}} \left\| \mathcal{I}(\mathcal{N}_{\theta^{(1)}}(a, b, p)) \right\|_2^2 \\ & \left. + w_{\text{int}} \left\| \mathcal{C}^{(m)}(\mathcal{N}_{\theta^{(m)}}, \mathcal{N}_{\theta^{(m+1)}}) \right\|_2^2 \right]. \end{aligned}$$

Here, (a, b, p) collectively denote the *input field*, *boundary specification*, and *physical parameters* that characterize the governing operator within each subdomain, drawn from a distribution $\mathcal{D}^{(m)}$ of admissible physical configurations. Each local neural operator $\mathcal{N}_{\theta^{(m)}}$ learns a mapping

$$(a, b, p) \mapsto u^{(m)}(\mathbf{x}, t),$$

enabling generalization across varying input conditions, boundary profiles, and parameter regimes. The interface condition

$$\mathcal{C}^{(m)} = \alpha_m \hat{u}_{\theta^{(m)}}(\Gamma_t^{(m)}) - \beta_m \hat{u}_{\theta^{(m+1)}}(\Gamma_t^{(m)}) = 0$$

ensures smooth transitions between adjacent subdomains, where coefficients α_m and β_m specify field or flux continuity. Each local operator minimizes a residual-based objective:

$$\mathcal{L}_{\text{res}}^{(m)} = w_{\text{eq}} \mathcal{L}_{\text{eq}}^{(m)} + w_{\text{ic}} \mathcal{L}_{\text{ic}}^{(m)} + w_{\text{int}} \mathcal{L}_{\text{int}}^{(m)},$$

with

$$\mathcal{L}_{\text{eq}}^{(m)} = \frac{1}{N^{(m)}} \sum_{i=1}^{N^{(m)}} \left| \mathcal{R}^{(m)}(\hat{u}_{\theta^{(m)}}(\mathbf{x}_i, t_i), a(\mathbf{x}_i, t_i), p^{(m)}) \right|^2,$$

$$\mathcal{R}^{(m)}(\mathbf{x}, t) = \frac{\partial u}{\partial t} + \mathcal{L}^{(m)}(u, a, p^{(m)}),$$

where $\mathcal{L}^{(m)}(\cdot)$ represents the local differential operator acting on u . Interface and initial-condition losses are defined as

$$\mathcal{L}_{\text{int}}^{(m)} = \frac{1}{N_{\Gamma}^{(m)}} \sum_{(\mathbf{x}, t) \in \Gamma_t^{(m)}} \left| \hat{u}_{\theta^{(m)}}(\mathbf{x}, t) - \hat{u}_{\theta^{(m+1)}}(\mathbf{x}, t) \right|^2,$$

$$\mathcal{L}_{\text{ic}}^{(m)} = \frac{1}{N_0} \sum_{\mathbf{x}_k \in \Omega} \left| \hat{u}_{\theta^{(1)}}(\mathbf{x}_k, 0) - u_0(\mathbf{x}_k) \right|^2.$$

To capture multi-frequency behavior, each subnetwork may employ Fourier feature mappings of the form

$$\phi_F(\mathbf{x}, t) = [\sin(B[\mathbf{x}, t]), \cos(B[\mathbf{x}, t])], \quad B \in \mathbb{R}^{d_t \times r},$$

where r is the number of Fourier modes and entries of B are sampled from a Gaussian distribution scaled by σ_B , enabling the encoding of fine-scale spatial and temporal variations. The DD-PINN naturally extends to multi-physics systems by allowing each subdomain to follow a distinct operator $\mathcal{F}^{(m)}$ and parameters $p^{(m)}$, with interface constraints $\mathcal{C}^{(m)}$ ensuring compatibility across heterogeneous regions. This unified residual-based formulation enforces local physical laws, continuity across subdomains, and global boundary consistency within a single optimization framework. The Diffusion Operator (DO) extends neural operator learning to a probabilistic and physics-constrained setting capable of modeling deterministic and stochastic systems governed by ordinary or partial differential equations across spatial dimensions $d \in \{0, 1, 2, 3\}$ and time. Unlike deterministic mappings, the DO learns a stochastic generative operator defined through a forward–reverse diffusion process that evolves along an auxiliary diffusion time variable $\tau \in [0, 1]$. Let $u_0 = u(\mathbf{x}, t)$ denote the normalized physical state field over the spatio-temporal domain $\Omega_t = \Omega \times [0, T]$, where $\Omega \subseteq \mathbb{R}^d$ is the spatial region and t represents physical time. The forward diffusion process $q(u_\tau | u_0)$ introduces a stochastic perturbation to u_0 :

$$q(u_\tau | u_0) = \mathcal{N}(\sqrt{\bar{\alpha}_\tau} u_0, (1 - \bar{\alpha}_\tau) \mathbf{I}), \quad \bar{\alpha}_\tau = \prod_{s=1}^{\tau} (1 - \beta_s),$$

where $\beta_s \in (0, 1)$ specifies the diffusion variance schedule and \mathbf{I} is the identity covariance. This yields the diffused field

$$u_\tau = \sqrt{\bar{\alpha}_\tau} u_0 + \sqrt{1 - \bar{\alpha}_\tau} \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}),$$

representing a continuous interpolation between the structured state u_0 and isotropic Gaussian noise ϵ . The family $\{u_\tau\}_{\tau \in [0, 1]}$ thus forms a continuum of realizations whose entropy increases with τ . The reverse process reconstructs u_0 from u_τ by learning a parameterized inverse transition $p_\theta(u_{\tau-\Delta\tau} | u_\tau, c)$, where θ are trainable parameters and $c = (a, b, p)$ represents conditioning variables: $a(\mathbf{x}, t)$ denotes external forcing, $b(\mathbf{x}, t)$ defines boundary or initial conditions, and $p \in \mathbb{R}^k$ captures physical parameters such as diffusivities or reaction coefficients. The neural operator \mathcal{G}_θ approximates the true reverse kernel,

$$p_\theta(u_{\tau-\Delta\tau} | u_\tau, c) \approx q(u_{\tau-\Delta\tau} | u_\tau, c),$$

through score matching, which links the known forward diffusion q with the learnable operator. The score function satisfies

$$\nabla_{u_\tau} \log q(u_\tau | c) = -\frac{u_\tau - \sqrt{\bar{\alpha}_\tau} u_0}{1 - \bar{\alpha}_\tau} \simeq -\frac{1}{\sqrt{1 - \bar{\alpha}_\tau}} \mathcal{G}_\theta(u_\tau, \tau, \phi(c)),$$

where $\phi(c)$ is an embedding of the conditioning context. The operator $\mathcal{G}_\theta(u_\tau, \tau, \phi(c))$ thus approximates the noise component in the forward process. Model parameters are optimized by minimizing

$$\mathcal{L}_{\text{diff}}(\theta) = \mathbb{E}_{\tau, u_0, \epsilon} [\|\epsilon - \mathcal{G}_\theta(u_\tau, \tau, \phi(c))\|_2^2],$$

which enforces accurate denoising across all diffusion times. To evaluate \mathcal{G}_θ , the field u_τ is encoded into a latent feature representation $h^{(0)}$ that captures spatial, temporal, and physical context:

$$h^{(0)}(\mathbf{x}, t, \tau) = \text{Linear}([u_\tau(\mathbf{x}, t), \mathbf{x}, t, \phi(c)]) + \psi_\tau(\tau),$$

where $h^{(0)} \in \mathbb{R}^{D_h}$ is the latent embedding, $\psi_\tau(\tau)$ encodes diffusion time, and D_h is the latent dimension. The features evolve hierarchically as

$$h^{(l)} = \text{Conv}(\text{SiLU}(\text{GN}(h^{(l-1)}))) + \psi_c(\phi(c)) + \text{Attn}(h^{(l-1)}),$$

where $\text{Conv}(\cdot)$ denotes spatial convolution (1D/2D/3D depending on dimension d), $\text{GN}(\cdot)$ is group normalization, $\text{SiLU}(\cdot)$ the smooth activation, and $\text{Attn}(\cdot)$ a multi-head attention mechanism capturing nonlocal dependencies. The final feature state $h^{(L)}$ is projected to the diffusion score as

$$\mathcal{G}_\theta(u_\tau, \tau, \phi(c)) = \text{Linear}(h^{(L)}),$$

yielding the predicted perturbation field. Physical consistency is enforced using a residual operator derived from the governing equation $\frac{\partial u}{\partial t} = \mathcal{F}(u, a, p)$:

$$\mathcal{R}(\mathbf{x}, t) = \frac{\partial u}{\partial t} + \mathcal{L}(u, a, p),$$

leading to the combined optimization objective

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{diff}} + w_{\text{phys}} \|\mathcal{R}(\hat{u}_\theta)\|_2^2,$$

where w_{phys} controls the relative weight of the physics constraint. During inference, the reverse diffusion reconstructs the clean field iteratively:

$$u_{\tau-\Delta\tau} = \frac{1}{\sqrt{\alpha_\tau}} \left(u_\tau - \frac{1-\alpha_\tau}{\sqrt{1-\bar{\alpha}_\tau}} \mathcal{G}_\theta(u_\tau, \tau, \phi(c)) \right) + \sigma_\tau \epsilon,$$

with $\sigma_\tau = \sqrt{\tilde{\beta}_\tau}$ as the stochastic step variance. After convergence ($\tau \rightarrow 0$), the reconstructed field is obtained by projecting onto the physical coordinates,

$$\hat{u}_\theta(\mathbf{x}, t) = \mathcal{P}(\mathcal{G}_\theta(u_\tau, \tau, \phi(c))),$$

where $\mathcal{P}(\cdot)$ denotes the projection operator. The reconstruction accuracy is measured by the relative L_2 error,

$$\epsilon_{\text{rel}} = \frac{\|\hat{u}_\theta - u^*\|_2}{\|u^*\|_2},$$

where $u^*(\mathbf{x}, t)$ is the reference physical solution.

In this formulation, the latent representations $h^{(l)}$ serve as hierarchical encodings of the diffused state u_τ , embedding spatial-temporal correlations and physical dependencies before projection through \mathcal{G}_θ . The DO thus establishes a coherent mapping between stochastic diffusion dynamics and neural operator learning, yielding a dimension-agnostic, physics-aware, and time-conditioned framework that unifies probabilistic generative modeling and physics-informed reasoning within a single formulation.

Deep Operator Network (DeepONet) The *Deep Operator Network (DeepONet)* is a neural framework that learns mappings between input functions, boundary-initial specifications, and physical parameters to predict the corresponding solution fields. It achieves this by decomposing the operator into a branch network, which encodes input functions, and a trunk network, which represents the solution basis over the spatio-temporal domain. It represents the solution operator

$$\mathcal{G} : (a, b, p) \mapsto u,$$

that maps an input or forcing function $a(\mathbf{x}, t)$, boundary/initial condition specification b , and physical parameters p to the corresponding state field $u(\mathbf{x}, t)$ defined over a domain $\Omega_t = \Omega \times [0, T]$, where $\Omega \subseteq \mathbb{R}^d$ with spatial dimension $d \in \{0, 1, 2, 3\}$. DeepONet decomposes the operator into two neural sub-networks: (i) the *branch network* \mathcal{B}_θ encodes the input functional $a(\mathbf{x}, t)$, boundary/initial conditions b , and parameters p into a latent representation in \mathbb{R}^r , and (ii) the *trunk network* \mathcal{T}_θ represents the coordinate-dependent basis functions over (\mathbf{x}, t) in \mathbb{R}^r . The predicted solution field is expressed as

$$\hat{u}_\theta(\mathbf{x}, t; a, b, p) = \mathcal{B}_\theta(a, b, p) \cdot \mathcal{T}_\theta(\mathbf{x}, t) + \mathbf{b}_\theta,$$

where the inner product couples the input-dependent latent space with the coordinate-dependent basis, and \mathbf{b}_θ denotes learnable offsets for each physical state variable.

The governing physical process is represented by a generalized differential operator

$$\mathcal{F}(\mathbf{x}, t, u, \nabla u, \nabla^2 u, \frac{\partial u}{\partial t}, a(\mathbf{x}, t), p) = 0, \quad (\mathbf{x}, t) \in \Omega_t,$$

where $\mathcal{F}(\cdot)$ may describe transport, diffusion, reaction, or other physical mechanisms. For steady-state PDEs, the temporal term $\partial u / \partial t$ vanishes, while for purely temporal ODEs, the spatial derivatives ∇u and $\nabla^2 u$ are absent. The residual operator that enforces the physical constraint is given by

$$\mathcal{R}(\mathbf{x}, t) = \frac{\partial \hat{u}_\theta}{\partial t} + \Phi(\hat{u}_\theta, \nabla \hat{u}_\theta, \nabla^2 \hat{u}_\theta, a(\mathbf{x}, t), p),$$

where $\Phi(\cdot)$ encapsulates the differential form of the governing equations. All necessary spatial and temporal derivatives of \hat{u}_θ are computed analytically via automatic differentiation through the trunk network. The learned solution must satisfy auxiliary conditions defined as: (i) *boundary conditions* $\mathcal{B}(u) = 0$ on $\partial\Omega$, enforcing prescribed Dirichlet, Neumann, or Robin constraints on the spatial boundary; and (ii) the *initial condition* $\mathcal{I}(u(\mathbf{x}, 0)) = u_0(\mathbf{x})$ in Ω , specifying the state of the system at the initial time $t = 0$. These conditions are incorporated into the training objective through additional loss terms that penalize deviation from the prescribed boundary and initial values. The overall training objective minimizes the weighted sum of the residual, boundary, and initial condition losses:

$$\min_{\theta} \mathbb{E}_{(a, b, p) \sim \mathcal{D}} \left[w_{\text{eq}} \|\mathcal{R}(\mathbf{x}, t)\|_2^2 + w_{\text{bc}} \|\mathcal{B}(\hat{u}_\theta)\|_2^2 + w_{\text{ic}} \|\mathcal{I}(\hat{u}_\theta)\|_2^2 \right],$$

where $w_{\text{eq}}, w_{\text{bc}}, w_{\text{ic}}$ control the relative importance of enforcing the physical equation, boundary conditions, and initial conditions. The optimal model parameters are obtained by

$$\theta^* = \arg \min_{\theta} \mathcal{L}_{\text{res}}(\theta), \quad \mathcal{L}_{\text{res}} = w_{\text{eq}} \mathcal{L}_{\text{eq}} + w_{\text{bc}} \mathcal{L}_{\text{bc}} + w_{\text{ic}} \mathcal{L}_{\text{ic}}.$$

DeepONet generalizes seamlessly from purely temporal systems to spatio-temporal domains and extends to multi-physics interactions by assigning separate trunk networks to different fields while sharing a common branch encoder. Thus, DeepONet provides a dimension-agnostic, physics-constrained, and time-aware operator-learning framework capable of learning continuous mappings between functional inputs, boundary–initial specifications, and solution fields across diverse dynamical systems.

Bayesian Deep Operator Network (B-DeepONet) The *Bayesian Deep Operator Network (B-DeepONet)* extends the deterministic DeepONet framework by introducing probabilistic inference into the operator-learning paradigm. It aims to learn mappings between input functions, boundary–initial specifications, and physical parameters while quantifying epistemic uncertainty in the learned operator. The model decomposes the stochastic operator into a probabilistic branch network, which encodes input functions, and a probabilistic trunk network, which represents the solution basis over the spatio-temporal domain. Each neural layer is modeled with Bayesian weights and biases, enabling the propagation of uncertainty through the operator space.

Formally, the Bayesian operator is defined as

$$\mathcal{G}_\theta : (a, b, p) \mapsto u,$$

where $a(\mathbf{x}, t)$ denotes the input or forcing function, b represents the boundary and initial specifications, and p contains the physical parameters. The corresponding stochastic prediction $\hat{u}_\theta(\mathbf{x}, t; a, b, p)$ represents a distribution over possible solution fields. The model approximates the true operator posterior $p(\theta|\mathcal{D})$ using a variational distribution $q_\phi(\theta)$, leading to a tractable formulation via variational inference. The network architecture consists of two Bayesian sub-networks: (i) a *Bayesian branch network* \mathcal{B}_θ , which encodes the functional inputs, boundary–initial conditions, and parameters into a latent random representation in \mathbb{R}^r , and (ii) a *Bayesian trunk network* \mathcal{T}_θ , which constructs coordinate-dependent stochastic basis functions over $(\mathbf{x}, t) \in \Omega_t$. The output is expressed as

$$\hat{u}_\theta(\mathbf{x}, t; a, b, p) = \mathcal{B}_\theta(a, b, p) \cdot \mathcal{T}_\theta(\mathbf{x}, t) + \mathbf{b}_\theta,$$

where each weight in \mathcal{B}_θ and \mathcal{T}_θ is sampled from a Gaussian posterior with trainable mean and variance, i.e.,

$$w = \mu + \sigma \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I),$$

allowing Monte Carlo sampling for uncertainty propagation during inference. The governing physics is described by a generalized differential operator

$$\mathcal{F}(\mathbf{x}, t, u, \nabla u, \nabla^2 u, \frac{\partial u}{\partial t}, a(\mathbf{x}, t), p) = 0, \quad (\mathbf{x}, t) \in \Omega_t,$$

subject to boundary and initial conditions

$$\mathcal{B}(u) = 0 \quad \text{on } \partial\Omega, \quad \mathcal{I}(u(\mathbf{x}, 0)) = u_0(\mathbf{x}) \quad \text{in } \Omega.$$

The residual operator enforcing the physical constraint is defined as

$$\mathcal{R}(\mathbf{x}, t) = \frac{\partial \hat{u}_\theta}{\partial t} + \Phi(\hat{u}_\theta, \nabla \hat{u}_\theta, \nabla^2 \hat{u}_\theta, a(\mathbf{x}, t), p),$$

where $\Phi(\cdot)$ captures the underlying physical relationships, and all spatial–temporal derivatives are obtained via automatic differentiation. The Bayesian training objective minimizes the variational free energy (negative evidence lower bound), which combines the physics-informed residual loss with a Kullback–Leibler (KL) divergence term regularizing the variational posterior:

$$\mathcal{L}_{\text{Bayes}}(\theta) = \mathbb{E}_{q_\phi(\theta)} \left[w_{\text{eq}} \|\mathcal{R}(\mathbf{x}, t)\|_2^2 + w_{\text{bc}} \|\mathcal{B}(\hat{u}_\theta)\|_2^2 + w_{\text{ic}} \|\mathcal{I}(\hat{u}_\theta)\|_2^2 \right] + \beta_{\text{KL}}$$

where $p(\theta)$ denotes the prior distribution over network weights and β_{KL} controls the trade-off between data fitting and regularization. The optimization objective is

$$\theta^* = \arg \min_{\theta} \mathcal{L}_{\text{Bayes}}(\theta).$$

At inference time, the Bayesian DeepONet generates a predictive distribution by sampling multiple realizations of the operator:

$$p(\hat{u}|a, b, p) = \int p(\hat{u}|a, b, p, \theta) q_\phi(\theta) d\theta,$$

from which the predictive mean $\mathbb{E}[\hat{u}]$ and uncertainty (variance or credible intervals) can be estimated via Monte Carlo sampling. This allows explicit quantification of both epistemic and aleatoric uncertainty in the learned operator space. The Bayesian formulation thus extends the deterministic DeepONet into a probabilistic operator-learning paradigm that captures model uncertainty, improves robustness under limited or noisy data, and generalizes to unseen parameter regimes.

Galerkin Transformer (GT) The *Galerkin Transformer (GT)* establishes a unified, dimension-agnostic, and physics-constrained operator-learning paradigm capable of solving ordinary and partial differential equations, as well as coupled multi-physics systems that evolve over time. The central objective is to learn a mapping between functional spaces that satisfies the governing physical laws globally, rather than pointwise, by embedding the Galerkin weak formulation into the neural operator training process. This ensures global consistency, stability, and strict adherence to boundary and initial constraints. Let the physical process be described by a generalized differential operator

$$\mathcal{F}(\mathbf{x}, t, u, \nabla u, \nabla^2 u, \frac{\partial u}{\partial t}, a, p) = 0,$$

where $u(\mathbf{x}, t)$ denotes the state field defined on a spatial domain $\Omega \subseteq \mathbb{R}^d$ and temporal domain $t \in [0, T]$; $a(\mathbf{x}, t)$ represents external inputs or forcing fields; and $p \in \mathbb{R}^k$ denotes a vector of physical parameters such as diffusivities, reaction constants, or conductivities. The system satisfies the boundary and initial constraints

$$\mathcal{B}(u) = 0 \quad \text{on } \partial\Omega, \quad \mathcal{I}(u(\mathbf{x}, 0)) = u_0(\mathbf{x}) \quad \text{in } \Omega,$$

where $\mathcal{B}(\cdot)$ and $\mathcal{I}(\cdot)$ denote the boundary and initial operators, respectively. Instead of enforcing the governing equations at discrete collocation points, the Galerkin formulation requires the residual of the governing operator to be orthogonal to a finite set of basis functions in the chosen function

space. Using orthogonal Legendre polynomials $\{P_k(\xi)\}_{k=1}^{N_w}$ as the basis, the weak form is expressed as

$$\int_{\Omega_t} P_k(\xi(t)) \mathcal{R}(\hat{u}_\theta(\mathbf{x}, t), a(\mathbf{x}, t), p) d\mathbf{x} dt = 0, \quad k = 1, \dots, N_w,$$

where $\Omega_t = \Omega \times [0, T]$ denotes the spatio-temporal domain and $\xi(t) = 2(t - t_0)/(T - t_0) - 1$ maps the temporal interval to $[-1, 1]$. The residual operator

$$\mathcal{R}(\mathbf{x}, t) = \frac{\partial \hat{u}_\theta}{\partial t} + \mathcal{L}(\hat{u}_\theta, a(\mathbf{x}, t), p),$$

quantifies the local deviation of the predicted solution $\hat{u}_\theta(\mathbf{x}, t)$ from the governing physical law, where $\mathcal{L}(\cdot)$ encapsulates the spatial or parametric components of the physics such as diffusion, advection, or reaction. The corresponding weak-form loss penalizes the L_2 norm of these projected residuals:

$$\mathcal{L}_{\text{wf}} = \frac{1}{N_w} \sum_{k=1}^{N_w} \left\| \int_{\Omega_t} P_k(\xi(t)) \mathcal{R}(\hat{u}_\theta(\mathbf{x}, t), a(\mathbf{x}, t), p) d\mathbf{x} dt \right\|_2^2,$$

where $P_k(\xi(t)) \mathcal{R}(\cdot) d\mathbf{x} dt$ represents the inner product between the k^{th} Legendre polynomial basis and the physics residual integrated over space and time. The Galerkin Transformer parameterizes the trial solution $\hat{u}_\theta(\mathbf{x}, t)$ as a neural operator that maps the forcing fields and parameters to the solution field:

$$\hat{u}_\theta(\mathbf{x}, t) = \mathcal{G}_\theta(a(\mathbf{x}, t), p) = \text{Linear}(h^{(L)}(\mathbf{x}, t)),$$

where the latent representation $h^{(L)}(\mathbf{x}, t)$ is obtained through a hierarchy of attention-based layers initialized by

$$h^{(0)}(\mathbf{x}, t) = \text{Linear}(\phi(\mathbf{x}, t, a(\mathbf{x}, t), p)) + \text{PE}(\mathbf{x}, t).$$

Here, $\phi(\mathbf{x}, t, a, p)$ denotes the concatenated feature vector containing spatial coordinates \mathbf{x} , temporal coordinate t , input or forcing values $a(\mathbf{x}, t)$, and physical parameters p , while $\text{PE}(\mathbf{x}, t)$ encodes positional information. The hidden states are recursively updated through

$$h^{(l)} = \text{MHA}(h^{(l-1)}) + \text{FFN}(h^{(l-1)}), \quad l = 1, \dots, L.$$

The Transformer thus parameterizes the function space within which Galerkin orthogonality is enforced: $h^{(0)}(\mathbf{x}, t)$ provides the neural embedding from which the trial solution \hat{u}_θ is constructed, and $\mathcal{R}(\mathbf{x}, t)$ quantifies its deviation from the governing dynamics. Boundary and initial constraints are softly enforced via auxiliary losses that complement the weak-form residual:

$$\mathcal{L}_{\text{bc}} = \frac{1}{N_b} \sum_{\mathbf{x}_j \in \partial\Omega} |\mathcal{B}(\hat{u}_\theta(\mathbf{x}_j, t_j))|^2, \quad \mathcal{L}_{\text{ic}} = \frac{1}{N_0} \sum_{\mathbf{x}_k \in \Omega} |\hat{u}_\theta(\mathbf{x}_k, 0) - u_0(\mathbf{x}_k)|^2.$$

The total optimization objective becomes

$$\mathcal{L}_{\text{total}} = w_{\text{wf}} \mathcal{L}_{\text{wf}} + w_{\text{bc}} \mathcal{L}_{\text{bc}} + w_{\text{ic}} \mathcal{L}_{\text{ic}},$$

where $w_{\text{wf}}, w_{\text{bc}}, w_{\text{ic}}$ control the contribution of weak, boundary, and initial constraints, respectively.

The projection term $P_k(\xi(t)) \mathcal{R}(\hat{u}_\theta(\mathbf{x}, t), a(\mathbf{x}, t), p) d\mathbf{x} dt$ directly depends on the trial solution predicted by the Transformer, establishing a continuous feedback loop between

the neural representation and the governing equations. The Transformer generates \hat{u}_θ through its encoded features $h^{(0)}(\mathbf{x}, t)$, while the Galerkin projection evaluates and minimizes the global physics residual across the domain. Through this coupling, the network learns to represent operators that satisfy the governing system in a weak, energy-minimizing sense. By enforcing orthogonality of the residual to the Legendre polynomial basis and incorporating boundary and initial conditions within the training objective, the Galerkin Transformer achieves globally consistent, physically faithful, and numerically stable solutions. This framework unifies time-dependent ODEs, steady and transient PDEs, and coupled multi-physics problems under a single physics-constrained operator-learning formulation.

Fourier Neural Operator (FNO) The Fourier Neural Operator (FNO) is a fully spectral operator-learning framework designed to learn mappings between continuous function spaces by representing input and solution fields in the Fourier domain, capturing spatial-temporal dependencies through a compact set of spectral modes. The underlying Fourier Neural Operator (FNO) defines a mapping between functional spaces,

$$\mathcal{G}_\theta : a(\mathbf{x}, t) \mapsto u(\mathbf{x}, t),$$

which is learned through iterative updates in Fourier space as

$$v_{t+1}(\mathbf{x}) = \sigma(W v_t(\mathbf{x}) + \mathcal{F}^{-1}(R_\theta(k) \mathcal{F}[v_t](k))),$$

where v_t denotes the hidden representation at iteration t , \mathcal{F} and \mathcal{F}^{-1} denote the Fourier and inverse Fourier transforms, $R_\theta(k)$ represents learnable complex-valued weights acting on a truncated set of spectral modes, W is a local linear transformation, and $\sigma(\cdot)$ is a nonlinear activation function. This formulation enables global communication through spectral convolution while retaining local adaptability via residual connections. The learnable operator $\mathcal{G}_\theta : (a, b, p) \mapsto u$ is parameterized using sequential spectral convolutional layers and residual transformations. Each layer transforms the field $u(\mathbf{x}, t)$ into its frequency representation $\tilde{u}_k = \mathcal{F}[u](k)$, applies learnable weights W_k to selected Fourier coefficients, and reconstructs the signal via $u' = \mathcal{F}^{-1}[W_k \tilde{u}_k] + W u$. This formulation preserves global smoothness while enabling local corrections through residual blending. Training is performed by minimizing a physics-constrained loss function that enforces the governing equations, along with initial and boundary conditions. The total loss is defined as

$$\mathcal{L}_{\text{res}}(\theta) = w_{\text{eq}} \sum_{m=1}^M \frac{1}{N_t} \sum_{i=1}^{N_t} |\mathcal{R}_m(\hat{u}_\theta^m(\mathbf{x}_i, t_i), a(\mathbf{x}_i, t_i), p)|^2 + w_{\text{ic}} \sum_{m=1}^M |\hat{u}_\theta^m$$

where m indexes each coupled state variable. The residual operator for each field is expressed as

$$\mathcal{R}_m(\mathbf{x}, t) = \frac{\partial u_m}{\partial t} + \Phi_m(u_1, \dots, u_M, \nabla u_m, \nabla^2 u_m, a(\mathbf{x}, t), p),$$

with temporal derivatives computed directly in the spectral domain, ensuring continuous and differentiable gradi-

ent evaluation without discretization errors. To ensure stable and efficient convergence, the framework employs adaptive learning strategies. The initial-condition weight w_{ic} follows a cosine annealing schedule, allowing early alignment with the initial state before focusing on physics residuals. A curriculum-learning mechanism progressively increases the temporal range of residual evaluation, improving stability for long-horizon training. The Fourier architecture generalizes naturally to higher-dimensional problems, where spectral convolutions operate as $\text{SpectralConv}_d : \mathbb{R}^{B \times N_1 \times \dots \times N_d \times C_{in}} \rightarrow \mathbb{R}^{B \times N_1 \times \dots \times N_d \times C_{out}}$, enabling a unified treatment of coupled physical systems within a single formulation. By embedding physical laws directly into the optimization process and performing continuous operator approximation through spectral differentiation, the PINO establishes a dimension-agnostic, physics-constrained, and time-aware framework for solving ODEs, PDEs, and multi-physics systems under one unified formulation.

Factorized Fourier Neural Operator (F-FNO) The Factorized Fourier Neural Operator (F-FNO) generalizes spectral operator learning by introducing a separable decomposition of Fourier transforms along each coordinate axis, enabling efficient representation of high-dimensional and anisotropic physical processes. Let the governing system be defined by a generalized differential operator

$$\mathcal{F}(\mathbf{x}, t, u, \nabla u, \nabla^2 u, \frac{\partial u}{\partial t}, a(\mathbf{x}, t), p) = 0, \quad (\mathbf{x}, t) \in \Omega_t,$$

where $u(\mathbf{x}, t)$ is the state field, $a(\mathbf{x}, t)$ is an input or forcing field, and p denotes physical parameters. The objective is to approximate the solution operator $\mathcal{G} : (a, b, p) \mapsto u$ with a learnable neural operator \mathcal{G}_θ parameterized by θ . Each F-FNO layer performs independent spectral transformations along spatial and temporal coordinates. For an input field $x(\mathbf{x}, t)$, the factorized spectral mapping is

$$\mathcal{S}_\theta(x) = \sum_{d=1}^D \mathcal{F}_d^{-1}(R_d(k_d) \mathcal{F}_d[x]),$$

where \mathcal{F}_d and \mathcal{F}_d^{-1} are the Fourier and inverse Fourier transforms along dimension d , and $R_d(k_d)$ are learnable complex-valued filters truncated to low-frequency modes. This separable construction learns independent correlations along each coordinate while combining them additively to capture cross-dimensional interactions. Each spectral block is coupled with a residual connection and nonlinear activation,

$$x' = \sigma(\text{Norm}(\mathcal{S}_\theta(x) + Wx)) + x,$$

where W is a local linear transformation, $\text{Norm}(\cdot)$ ensures stability, and $\sigma(\cdot)$ is a smooth activation. Training minimizes a physics-constrained objective that enforces the governing equations and auxiliary conditions:

$$\mathcal{L}_{\text{res}}(\theta) = w_{\text{eq}} \frac{1}{N} \sum_{i=1}^N |\mathcal{R}(\hat{u}_\theta(\mathbf{x}_i, t_i), a(\mathbf{x}_i, t_i), p)|^2 + w_{ic} \|\hat{u}_\theta(\mathbf{x}, 0) - u_0(\mathbf{x})\|_2^2 + w_{bc} \|\hat{u}_\theta(\partial\Omega, t) - b(\partial\Omega, t)\|_2^2,$$

where $\mathcal{R}(\cdot)$ is the residual of the physical operator, and $w_{\text{eq}}, w_{ic}, w_{bc}$ balance the respective terms. Spatial and temporal derivatives are obtained through automatic differentiation, ensuring continuous and differentiable constraint enforcement. This formulation extends naturally to any spatial

dimension $D \in \{0, 1, 2, 3\}$, with spectral convolutions acting independently as

$$\text{SpectralConv}_d : \mathbb{R}^{B \times N_1 \times \dots \times N_d \times C_{in}} \rightarrow \mathbb{R}^{B \times N_1 \times \dots \times N_d \times C_{out}},$$

enabling scalable modeling of complex dynamical systems. By embedding physical consistency directly within the optimization and factorizing spectral learning across all axes, the F-FNO achieves efficient, stable, and generalizable operator learning under a single unified framework.

References