

# 1 PhysRobot: Physics-Informed Graph Neural Networks for Sample-Efficient Robot Manipulation

---

**Target venue:** ICRA 2027 / CoRL 2026

**Draft version:** V1 — 2026-02-06

**Status:** Full draft for internal review. Items marked [TODO] require experimental data or figures.

---

## 1.1 Abstract

Model-free reinforcement learning (RL) for contact-rich, multi-object robotic manipulation remains notoriously sample-inefficient, often requiring  $10^6$ – $10^8$  environment interactions because policies must rediscover fundamental physical laws—conservation of momentum, Newton’s Third Law, rigid-body constraints—entirely from data. Existing physics-informed approaches embed such structure into forward dynamics models (Hamiltonian Neural Networks, Graph Network Simulators) but leave the *policy network* unconstrained, forfeiting a powerful inductive bias at the point of decision-making. We propose **PhysRobot**, a dual-stream RL architecture that integrates a **momentum-conserving Graph Neural Network** directly into the policy. The physics stream constructs a dynamic scene graph over interacting bodies and propagates messages through a novel **Scalarization–Vectorization (SV) pipeline** operating in edge-local coordinate frames, architecturally guaranteeing Newton’s Third Law ( $\mathbf{F}_{ij} = -\mathbf{F}_{ji}$ ) for *any* network parameters. A parallel policy stream processes raw observations, and the two streams are fused via stop-gradient concatenation to prevent RL gradients from distorting learned dynamics. On three MuJoCo manipulation benchmarks—PushBox, MultiPush, and Sort—PhysRobot achieves equivalent asymptotic performance to PPO and SAC in [TODO: result pending — target 3–5×] fewer environment steps and generalizes zero-shot to unseen object counts and mass distributions with [TODO: result pending — target <15%] performance degradation. To our knowledge, PhysRobot is the first framework to embed conservation-law structure directly into a GNN *policy* for contact-rich manipulation, demonstrating that physics-constrained policies yield compounding gains in both sample efficiency and out-of-distribution robustness.

---

## 1.2 1. Introduction

### 1.2.1 1.1 The Sample-Efficiency Crisis in Robot Manipulation

Robotic manipulation of multiple objects in contact-rich environments is among the most practically important and technically

challenging problems in robotics. Tasks as intuitive to humans as pushing several boxes to target positions, sorting objects by category, or assembling parts from a bin require a robot to reason about contact dynamics, friction, multi-body interactions, and sequential dependencies—all while operating under partial observability and noisy actuation.

Modern model-free reinforcement learning algorithms—Proximal Policy Optimization (PPO) [24], Soft Actor-Critic (SAC) [25], Twin Delayed DDPG (TD3) [26]—have demonstrated impressive results on continuous control benchmarks. Yet even with massively parallel simulation [27], these methods typically require  $10^6$ – $10^8$  environment steps for contact-rich manipulation tasks. The sample cost grows combinatorially with the number of interacting objects: a scene with  $N$  movable bodies has  $O(N^2)$  pairwise interactions, each potentially involving intermittent contact, friction, and momentum exchange. The resulting state space is vast, and the policy must discover—purely from reward signal—the physical regularities that govern every interaction.

This sample complexity has profound practical consequences. Real-robot training at the scale required by current methods is infeasible without substantial sim-to-real infrastructure [27, 30], and the sim-to-real gap itself introduces additional failure modes. Even in simulation, the computational cost of billions of environment steps limits iteration speed and the range of tasks that can be explored. A principled reduction in sample complexity would unlock broader applicability of RL-based manipulation.

### 1.2.2 1.2 Physics as Untapped Inductive Bias for Policies

The physical world is not arbitrary. Newton’s laws of motion, conservation of momentum and energy, and the structure of contact mechanics impose strong constraints on how objects interact. Humans exploit an intuitive understanding of these constraints to plan manipulation strategies with remarkable efficiency: we do not “learn” that pushing a heavy box requires more force than a light one—we *assume* it, and our motor planning reflects this prior knowledge from the very first attempt.

A growing body of work has demonstrated the value of encoding physical structure into neural network architectures. Hamiltonian Neural Networks (HNNs) [2] and Lagrangian Neural Networks (LNNs) [3] guarantee energy conservation by parameterizing the Hamiltonian or Lagrangian and deriving dynamics via automatic differentiation. Deep Lagrangian Networks (DeLaN) [4] extend this to articulated rigid bodies, learning accurate robot dynamics from as few as 100 trajectories. Graph Network Simulators (GNS) [11] use message-passing GNNs to simulate complex multi-body systems with remarkable generalization. Equivariant architectures—EGNN [7], PaiNN [15], NequIP [16], MACE [17]—ensure that learned representations transform correctly under rotations and translations, dramatically improving data efficiency for molecular dynamics.

However, a critical observation has been overlooked: **all of these methods encode physics into forward dynamics models or**

**simulators, not into policy networks.** The policy—the function that maps observations to actions—remains a generic MLP or transformer, oblivious to the physical structure of the problem it is solving. This means the RL optimizer must search over the full space of observation-to-action mappings, including the vast majority that violate basic physical laws. If the policy itself could “think in physics”—if its internal representations respected conservation laws and its message-passing reflected Newtonian interaction structure—the search space would be dramatically reduced.

### 1.2.3 1.3 Our Approach: PhysRobot

We propose **PhysRobot**, a dual-stream RL architecture that embeds Newtonian conservation laws directly into the policy network via a physics-informed Graph Neural Network.

**Dynamic Scene Graph.** At each timestep, PhysRobot constructs a graph over the scene: nodes represent the robot’s end-effector, manipulated objects, and goal locations; edges encode spatial proximity and contact relationships. Node features include position, velocity, mass, and friction; edge features capture relative displacement, distance, and relative velocity. The graph topology is recomputed at every step to reflect the changing contact configuration.

**Scalarization–Vectorization (SV) Message Passing.** The core architectural innovation is a message-passing scheme that *architecturally guarantees* Newton’s Third Law ( $\mathbf{F}_{ij} = -\mathbf{F}_{ji}$ ) for any choice of network parameters. For each edge, we construct a local coordinate frame from the displacement and relative velocity vectors. Geometric quantities are *scalarized*—projected onto this frame to produce rotation-invariant features—and processed by a standard MLP to produce scalar force coefficients. These coefficients are then *vectorized*—combined with the frame basis vectors to reconstruct a 3D force vector. A careful antisymmetrization of the binormal component ensures that paired forces cancel exactly, guaranteeing conservation of total linear momentum.

**Dual-Stream Fusion.** The physics stream produces predicted accelerations for each body, encoding the network’s “understanding” of the interaction dynamics. A parallel policy stream processes raw observations via a standard MLP. The two streams are fused via stop-gradient concatenation: the physics stream’s output is detached from the RL computation graph, preventing policy-gradient updates from distorting the learned dynamics. The physics stream is instead trained with a self-supervised auxiliary loss on finite-difference accelerations from observed transitions.

**Conservation Loss.** In addition to the hard architectural constraint on momentum, we add a soft energy-consistency regularizer based on the work-energy theorem, encouraging predictions compatible with energy conservation while still permitting dissipation through friction.

### 1.2.4 1.4 Contributions

We make the following contributions:

1. **PhysRobot architecture.** We introduce the first GNN policy that embeds Newtonian conservation into message passing for robotic manipulation. The dual-stream design with stop-gradient fusion provides a principled interface between physics reasoning and RL optimization.
2. **Corrected SV pipeline.** We identify and fix a subtle symmetry bug in the Dynami-CAL [41] framework where the binormal ( $\mathbf{e}_3$ ) message component breaks the antisymmetry required for momentum conservation. Our fix—multiplying the  $\alpha_3$  coefficient by the antisymmetric signed binormal velocity  $v_b$ —restores the guarantee with minimal computational overhead.
3. **Self-supervised physics learning.** The physics stream is trained on finite-difference accelerations from RL rollout data, requiring no ground-truth force labels, no differentiable simulator, and no offline trajectory dataset.
4. **Empirical validation.** We demonstrate [TODO: result pending — target 3–5×] sample-efficiency gains over PPO/SAC on contact-rich manipulation, competitive asymptotic performance, and strong zero-shot generalization to unseen object counts ([TODO: result pending — target train on 3, test on 5–10]) and mass distributions ([TODO: result pending — target train on 0.5 kg, test on 0.1–5.0 kg]).

### 1.2.5 1.5 Paper Organization

Section 2 reviews related work. Section 3 presents the PhysRobot method in detail. Section 4 describes our experimental evaluation. Section 5 concludes with limitations and future directions.

## 1.3 2. Related Work

Our work sits at the intersection of physics-informed machine learning, graph neural networks for physical reasoning, and reinforcement learning for robotic manipulation. We review each strand and position PhysRobot relative to prior art.

### 1.3.1 2.1 Physics-Informed Machine Learning

The integration of physical laws into neural network architectures has emerged as a powerful paradigm for improving data efficiency and generalization. **Physics-Informed Neural Networks (PINNs)** [1] embed partial differential equations as soft constraints in the loss function, enabling solutions to forward and inverse problems with limited data. While PINNs have been successfully applied to fluid dynamics and heat transfer, they require

explicit knowledge of the governing PDE and are primarily used for *prediction*, not *control*.

A complementary line of work encodes physical structure directly into network architecture. **Hamiltonian Neural Networks (HNNs)** [2] parameterize the Hamiltonian and derive dynamics via Hamilton’s equations, guaranteeing energy conservation by construction. **Lagrangian Neural Networks (LNNs)** [3] take a dual approach through the Lagrangian formalism, naturally handling generalized coordinates. **Deep Lagrangian Networks (DeLaN)** [4] extend this to articulated rigid-body systems, learning the mass matrix, Coriolis, and gravitational terms in Lagrangian form—demonstrating accurate dynamics prediction for robotic arms with as few as 100 training trajectories. More recently, **Dissipative SymODEN** [5] and **Port-Hamiltonian Neural Networks** [6] have extended energy-conserving architectures to dissipative systems, which are more realistic for robotics where friction is ubiquitous.

On the symmetry front, **Equivariant Graph Neural Networks (EGNNs)** [7] enforce  $E(n)$ -equivariance in message passing, ensuring that learned representations transform correctly under rotations and translations. This has been extended to  $SE(3)$ -equivariance by **SEGNN** [8] using steerable features, and to higher-order equivariance by Equivariant Polynomials [9]. Villar et al. [10] provide a theoretical framework for understanding which symmetries are most beneficial for different physical systems.

**Our distinction:** All the above methods target *forward dynamics modeling*—predicting how a system evolves. PhysRobot is the first to embed conservation-law structure directly into a *policy network* for active manipulation, using physics not to predict the future but to constrain the space of actions.

### 1.3.2 2.2 Graph Neural Networks for Physical Systems

Graph neural networks provide a natural substrate for modeling interacting physical systems. The seminal **Graph Network Simulator (GNS)** [11] demonstrated that learned message-passing over particle graphs can simulate complex physical phenomena—fluids, rigid bodies, deformables—with remarkable accuracy and generalization to larger systems. Follow-up work **GNS-Mesh** [12] extended this to mesh-based simulations, and Learned Simulators with Constraints [13] incorporated hard physical constraints via projection.

In molecular and atomic simulation, a rich ecosystem of equivariant GNNs has emerged. **DimeNet** [14] introduced directional message passing using interatomic angles. **PaiNN** [15] combined invariant and equivariant features for efficient and accurate force prediction. **NequIP** [16] leveraged  $E(3)$ -equivariant convolutions with spherical harmonics, achieving state-of-the-art accuracy for interatomic potentials with orders of magnitude less training data. **MACE** [17] extended this with higher body-order interactions, and **Allegro** [18] achieved scalability for large-scale molecular dynamics. More recently, **eSCN** [19] and **EquiformerV2** [20] have pushed the Pareto frontier of accuracy versus computational cost.

For macroscopic multi-body systems, **LoCS** [21] proposed continuous message passing over spatial fields, **Neural Relational Inference (NRI)** [22] learned to infer interaction graphs from observed trajectories, and **Interaction Networks** [23] demonstrated learning of object physics and relations. **Compositional generalization** via graph networks [11] enables zero-shot transfer to novel scene compositions.

**Our distinction:** These GNN architectures model *passive* physical systems—they predict how objects move under physical laws. PhysRobot adapts these architectures for *active control*: the GNN outputs features that inform action selection, and conservation structure ensures the policy’s internal model of interactions is physically consistent.

### 1.3.3 2.3 Reinforcement Learning for Manipulation

Model-free RL has made significant strides in robotic manipulation. **PPO** [24] remains the dominant on-policy algorithm due to its stability and scalability. **SAC** [25] introduced entropy regularization for continuous control, achieving strong sample efficiency in off-policy settings. **TD3** [26] addressed overestimation bias in actor-critic methods. These methods, combined with massively parallel simulation [27], have achieved impressive manipulation results but still require  $10^6$ – $10^8$  environment steps.

Model-based approaches improve sample efficiency by learning dynamics models. **Dreamer v3** [28] learns a world model in latent space and trains policies entirely in imagination, achieving competitive performance with approximately  $10\times$  fewer steps. **MuZero** [29] combined model-based planning with model-free value estimation. **DayDreamer** [30] demonstrated real-robot learning using Dreamer’s world model. **TD-MPC2** [31] unified temporal-difference learning with model-predictive control, scaling to 80+ tasks with a single model.

The foundation model era has brought new paradigms. **RT-2** [32] leveraged vision-language model pre-training for robotic control, achieving generalization through internet-scale knowledge. **Octo** [33] proposed an open-source generalist robot policy trained on the Open X-Embodiment dataset.  $\pi_0$  [34] from Physical Intelligence demonstrated flow-matching-based policies for dexterous manipulation. These approaches gain generalization through *data scale* rather than *structural priors*.

**Our distinction:** PhysRobot is complementary to both model-based and foundation model approaches. Unlike model-based methods that encode physics into the *world model*, we encode it into the *policy*. Unlike foundation models that rely on massive data, we achieve efficiency through *structural inductive bias*. PhysRobot could in principle be combined with physics-informed world models for compounding gains.

### 1.3.4 2.4 Physics Priors in Reinforcement Learning

The closest works to ours combine physics knowledge with RL. **PhyDNet** [35] integrates physical dynamics priors into video prediction networks but does not address control. Heiden et

al. [36] used differentiable physics simulators as world models for MBRL, but the physics is in the *simulator*, not the *policy*. **Analytical Policy Gradients** [37] computed exact gradients through differentiable simulation, bypassing the need for RL altogether—but this requires a differentiable simulator and does not generalize to black-box environments.

More recently, **Physics-Informed Neural Operator RL** [38] used neural operators with physics constraints as world models, demonstrating improved sample efficiency in fluid control tasks. **Symmetry-Aware RL** [39] exploited known symmetries in the MDP structure to reduce the effective state-action space, showing 2–5× speedups—but without conservation-law structure. **Structured World Models** [40] learned object-centric dynamics models with relational structure, improving planning but not directly addressing policy architecture.

Most relevant to our technical approach, **Dynami-CAL** [41] introduced edge-local coordinate frames and the scalarization-vectorization pipeline for physics simulation with GNNs. Our work extends Dynami-CAL in three key ways: (i) we integrate the SV-pipeline into an RL policy via the dual-stream architecture with stop-gradient fusion; (ii) we identify and correct a subtle antisymmetry bug in the binormal component (Section 3.3); and (iii) we design a self-supervised physics loss from RL rollout data, eliminating the need for labeled force data.

**PhysRobot’s unique position:** We are the first to embed conservation laws (momentum, energy) as *architectural constraints* within a GNN *policy network* for manipulation RL. This is orthogonal to physics-informed world models and can be combined with them. Our approach requires no differentiable simulator, no pre-trained foundation model, and no explicit knowledge of governing equations—only the *structural form* of conservation laws.

[TODO: figure 1 — Positioning diagram showing PhysRobot at the intersection of Physics-Informed ML, GNNs for Physics, and RL for Manipulation]

## 1.4 3. Method

We present PhysRobot in five parts: the problem formulation (§3.1), the dynamic scene graph construction (§3.2), the Scalarization–Vectorization message passing that guarantees momentum conservation (§3.3), the conservation guarantees and their proofs (§3.4), and the dual-stream RL integration (§3.5).

### 1.4.1 3.1 Problem Formulation

We consider multi-object robotic manipulation as a Markov Decision Process (MDP)  $(\mathcal{S}, \mathcal{A}, T, R, \gamma)$  with known low-dimensional state. The system consists of  $N$  interacting rigid bodies: the robot end-effector and  $N - 1$  objects.

**State space.** The state of body  $i$  at time  $t$  is:

$$\mathbf{s}_i^t = (\mathbf{x}_i^t, \dot{\mathbf{x}}_i^t, \phi_i)$$

where  $\mathbf{x}_i \in \mathbb{R}^3$  is position,  $\dot{\mathbf{x}}_i \in \mathbb{R}^3$  is velocity, and  $\phi_i$  encodes intrinsic properties (mass  $m_i$ , friction coefficient  $\mu_i$ , geometry descriptor  $\mathbf{g}_i$ ). The full observation is the concatenation of all body states plus goal specifications.

**Action space.** Continuous joint torques  $\mathbf{a} \in \mathbb{R}^{d_a}$ , where  $d_a$  is the number of actuated degrees of freedom.

**Reward.** We use a progress-based shaped reward:

$$r_t = - \sum_{k=1}^K \|\mathbf{x}_k^t - \mathbf{x}_k^*\|_2 + \beta \cdot \mathbb{1}[\text{success}] - \lambda_a \|\mathbf{a}_t\|^2$$

where  $\mathbf{x}_k^*$  denotes the goal position for object  $k$ ,  $\beta$  is a success bonus, and  $\lambda_a$  penalizes large actions to encourage smooth control.

**Objective.** Maximize the expected discounted return  $J(\pi) = \mathbb{E}_\pi[\sum_{t=0}^T \gamma^t r_t]$  while minimizing the number of environment interactions required to reach a target performance level.

### 1.4.2 3.2 Scene Graph Construction

We represent the manipulation scene as a dynamic directed graph  $\mathcal{G}^t = (\mathcal{V}, \mathcal{E}^t)$ .

**Nodes.** Each rigid body  $i$  contributes one node with features:

$$\mathbf{x}_i^{\text{node}} = [\mathbf{x}_i, \dot{\mathbf{x}}_i, m_i, \mu_i, \mathbf{g}_i, \tau_i]$$

where  $\tau_i \in \{\text{ee}, \text{object}, \text{goal}\}$  is a one-hot type indicator. For the robot end-effector,  $\mathbf{x}_i$  and  $\dot{\mathbf{x}}_i$  are computed via forward kinematics from joint states.

**Edges.** The edge set is time-varying and combines three sources:

$$\mathcal{E}^t = \mathcal{E}_{\text{prox}}^t \cup \mathcal{E}_{\text{kin}} \cup \mathcal{E}_{\text{ee}}$$

- **Proximity edges**  $\mathcal{E}_{\text{prox}}^t = \{(i, j) : \|\mathbf{x}_i^t - \mathbf{x}_j^t\| < r_{\text{cut}}\}$  capture spatial neighbors within a cutoff radius.
- **Kinematic edges**  $\mathcal{E}_{\text{kin}}$  encode the fixed robot kinematic chain (joint-to-joint connectivity).
- **End-effector edges**  $\mathcal{E}_{\text{ee}} = \{(i, j_{\text{ee}}) : i \in \text{objects}\}$  maintain permanent connections between the end-effector and all objects, since the robot can potentially interact with any of them.

All edges are bidirectional: if  $(i, j) \in \mathcal{E}^t$ , then  $(j, i) \in \mathcal{E}^t$ .

**Edge features.** For each directed edge  $(i \rightarrow j)$ :

$$\mathbf{e}_{ij}^{\text{feat}} = [\mathbf{r}_{ij}, \|\mathbf{r}_{ij}\|, \dot{\mathbf{x}}_{ij}, \|\dot{\mathbf{x}}_{ij}\|, \tau_{ij}]$$

where  $\mathbf{r}_{ij} = \mathbf{x}_j - \mathbf{x}_i$  is the relative displacement,  $\dot{\mathbf{x}}_{ij} = \dot{\mathbf{x}}_j - \dot{\mathbf{x}}_i$  is the relative velocity, and  $\tau_{ij} \in \{\text{ee-obj}, \text{obj-obj}, \text{obj-ground}\}$  is an edge type indicator.

**Scaling.** For  $N$  objects, the graph has  $|\mathcal{V}| = N + 1$  nodes (objects + end-effector) and  $|\mathcal{E}| = O(N + N_{\text{contact}})$  edges, which is linear in the number of objects for sparse contact configurations. This enables efficient scaling to scenes with many objects.

### 1.4.3 3.3 Scalarization–Vectorization Message Passing

The core architectural innovation of PhysRobot is a message-passing scheme that *by construction* guarantees Newton’s Third Law: for every pair of interacting bodies, the predicted forces are equal and opposite. We achieve this through the **Scalarization–Vectorization (SV) pipeline**, which operates in edge-local coordinate frames.

**1.4.3.1 3.3.1 Edge-Local Coordinate Frame** For each directed edge ( $i \rightarrow j$ ), we construct an orthonormal frame  $\{\mathbf{e}_1^{ij}, \mathbf{e}_2^{ij}, \mathbf{e}_3^{ij}\}$ :

**Radial basis vector:**

$$\mathbf{e}_1^{ij} = \frac{\mathbf{x}_j - \mathbf{x}_i}{\|\mathbf{x}_j - \mathbf{x}_i\| + \epsilon}$$

**Tangential basis vector** (from relative velocity):

$$\mathbf{v}_{ij}^\perp = \dot{\mathbf{x}}_{ij} - (\dot{\mathbf{x}}_{ij} \cdot \mathbf{e}_1^{ij}) \mathbf{e}_1^{ij}, \quad \mathbf{e}_2^{ij} = \frac{\mathbf{v}_{ij}^\perp}{\|\mathbf{v}_{ij}^\perp\| + \epsilon}$$

**Binormal basis vector:**

$$\mathbf{e}_3^{ij} = \mathbf{e}_1^{ij} \times \mathbf{e}_2^{ij}$$

**Degeneracy handling.** When the tangential velocity is near-zero ( $\|\mathbf{v}_{ij}^\perp\| < \epsilon_{\text{deg}}$ ), we fall back to a gravity-aligned frame:  $\mathbf{e}_2^{ij} = (\mathbf{e}_1^{ij} \times \hat{\mathbf{z}}) / \|\mathbf{e}_1^{ij} \times \hat{\mathbf{z}}\|$ , with a further fallback to  $\hat{\mathbf{y}}$  when  $\mathbf{e}_1^{ij} \approx \pm \hat{\mathbf{z}}$ . We recommend smooth blending between frames to avoid discontinuities.

**Critical antisymmetry property.** The frame vectors satisfy:

$$\mathbf{e}_1^{ij} = -\mathbf{e}_1^{ji}, \quad \mathbf{e}_2^{ij} = -\mathbf{e}_2^{ji}, \quad \mathbf{e}_3^{ij} = +\mathbf{e}_3^{ji}$$

*Proof.* We have  $\mathbf{e}_1^{ji} = (\mathbf{x}_i - \mathbf{x}_j) / \|\cdot\| = -\mathbf{e}_1^{ij}$ . Since  $\dot{\mathbf{x}}_{ji} = -\dot{\mathbf{x}}_{ij}$ , projection removal gives  $\mathbf{v}_{ji}^\perp = -\mathbf{v}_{ij}^\perp$ , hence  $\mathbf{e}_2^{ji} = -\mathbf{e}_2^{ij}$ . For the cross product:  $\mathbf{e}_3^{ji} = \mathbf{e}_1^{ji} \times \mathbf{e}_2^{ji} = (-\mathbf{e}_1^{ij}) \times (-\mathbf{e}_2^{ij}) = +\mathbf{e}_3^{ij}$ .  $\square$

This mixed symmetry—antisymmetric radial and tangential, symmetric binormal—is the key to the conservation guarantee, as we show below.

**1.4.3.2 3.3.2 Scalarization:**  $\mathbb{R}^3 \rightarrow \mathbb{R}^k$  Raw geometric quantities are projected onto the edge frame to produce **rotation-invariant scalar features**:

$$\sigma_{ij} = \begin{pmatrix} \|\mathbf{r}_{ij}\| \\ |v_r^{ij}| \\ v_t^{ij} \\ |v_b^{ij}| \\ \|\dot{\mathbf{x}}_{ij}\| \end{pmatrix} \in \mathbb{R}^5$$

where  $v_r^{ij} = \dot{\mathbf{x}}_{ij} \cdot \mathbf{e}_1^{ij}$  (radial component),  $v_t^{ij} = \dot{\mathbf{x}}_{ij} \cdot \mathbf{e}_2^{ij}$  (tangential component), and  $v_b^{ij} = \dot{\mathbf{x}}_{ij} \cdot \mathbf{e}_3^{ij}$  (binormal component). Note that absolute values are taken for  $v_r$  and  $v_b$  to ensure all entries of  $\sigma_{ij}$  are *symmetric* with respect to edge reversal ( $\sigma_{ij} = \sigma_{ji}$ ), while  $v_t$  is already symmetric because  $\mathbf{e}_2$  is antisymmetric and  $\dot{\mathbf{x}}_{ij}$  is antisymmetric, making their dot product symmetric.

**Extended scalar features** incorporate learned node embeddings:

$$\sigma_{ij}^{\text{ext}} = [\sigma_{ij} \parallel (\mathbf{h}_i + \mathbf{h}_j) \parallel (\mathbf{h}_i \odot \mathbf{h}_j)] \in \mathbb{R}^{5+2d_h}$$

where  $\mathbf{h}_i, \mathbf{h}_j \in \mathbb{R}^{d_h}$  are learned node embeddings, combined via *symmetric aggregation* (element-wise sum and element-wise product) to ensure permutation invariance:  $\sigma_{ij}^{\text{ext}} = \sigma_{ji}^{\text{ext}}$ .

**1.4.3.3 3.3.3 Scalar MLP:**  $\mathbb{R}^{5+2d_h} \rightarrow \mathbb{R}^3$  A standard MLP processes the scalar features to produce force coefficients:

$$(\alpha_1^{ij}, \alpha_2^{ij}) = \text{MLP}_{\alpha_{1,2}}(\sigma_{ij}^{\text{ext}})$$

Because the MLP operates on symmetric inputs, these coefficients satisfy  $\alpha_1^{ij} = \alpha_1^{ji}$  and  $\alpha_2^{ij} = \alpha_2^{ji}$ .

For the binormal coefficient  $\alpha_3$ , we must introduce explicit antisymmetry. We use the **signed binormal velocity**  $v_b^{ij} = \dot{\mathbf{x}}_{ij} \cdot \mathbf{e}_3^{ij}$  as an antisymmetric marker:

$$\alpha_3^{ij} = v_b^{ij} \cdot g_\theta(\sigma_{ij}^{\text{sym}})$$

where  $g_\theta$  is a scalar-valued MLP acting on symmetric features only, and  $v_b^{ij}$  satisfies  $v_b^{ij} = -v_b^{ji}$  (since  $\dot{\mathbf{x}}_{ij} = -\dot{\mathbf{x}}_{ji}$  and  $\mathbf{e}_3^{ij} = +\mathbf{e}_3^{ji}$ , we have  $v_b^{ij} = \dot{\mathbf{x}}_{ij} \cdot \mathbf{e}_3^{ij} = -\dot{\mathbf{x}}_{ji} \cdot \mathbf{e}_3^{ji} = -v_b^{ji}$ ). This ensures  $\alpha_3^{ij} = -\alpha_3^{ji}$ .

**1.4.3.4 3.3.4 Vectorization:**  $\mathbb{R}^3 \rightarrow \mathbb{R}^3$  The 3D force vector is reconstructed by combining scalar coefficients with frame basis vectors:

$$\mathbf{F}_{ij} = \alpha_1^{ij} \mathbf{e}_1^{ij} + \alpha_2^{ij} \mathbf{e}_2^{ij} + \alpha_3^{ij} \mathbf{e}_3^{ij}$$

This decomposition has a natural physical interpretation:  $\alpha_1$  models radial forces (normal contact, repulsion/attraction),  $\alpha_2$  models tangential forces (friction along the sliding direction), and  $\alpha_3$  models lateral forces (friction perpendicular to the sliding plane).

**1.4.3.5 3.3.5 Message Aggregation and Node Update** The net force on node  $i$  from all neighbors is:

$$\mathbf{M}_i = \sum_{j \in \mathcal{N}(i)} \mathbf{F}_{ij}$$

Node embeddings are updated via:

$$\mathbf{h}_i^{(\ell+1)} = \mathbf{h}_i^{(\ell)} + \text{MLP}_{\text{upd}}^{(\ell)} \left( [\mathbf{h}_i^{(\ell)} \parallel \mathbf{M}_i^{(\ell)}] \right)$$

We apply  $L$  rounds of SV message passing (we use  $L = 2$  in all experiments).

#### 1.4.4 3.4 Conservation Guarantees

We now prove that the SV pipeline guarantees linear momentum conservation for any network parameters.

**Theorem 1 (Linear Momentum Conservation).** *For any network parameters  $\theta$ , the total predicted force on the system vanishes:*

$$\sum_{i=1}^N \mathbf{F}_i = \mathbf{0}$$

where  $\mathbf{F}_i = \sum_{j \in \mathcal{N}(i)} \mathbf{F}_{ij}$  is the net force on node  $i$ .

*Proof.* The total force is:

$$\sum_i \mathbf{F}_i = \sum_i \sum_{j \in \mathcal{N}(i)} \mathbf{F}_{ij} = \sum_{(i,j) \in \mathcal{E}} (\mathbf{F}_{ij} + \mathbf{F}_{ji})$$

For each edge pair, expanding the force expressions:

$$\mathbf{F}_{ij} + \mathbf{F}_{ji} = \sum_{k=1}^3 (\alpha_k^{ij} \mathbf{e}_k^{ij} + \alpha_k^{ji} \mathbf{e}_k^{ji})$$

For  $k = 1, 2$ : We have  $\alpha_k^{ij} = \alpha_k^{ji}$  (symmetric coefficients) and  $\mathbf{e}_k^{ij} = -\mathbf{e}_k^{ji}$  (antisymmetric basis vectors), yielding:

$$\alpha_k^{ij} \mathbf{e}_k^{ij} + \alpha_k^{ji} \mathbf{e}_k^{ji} = \alpha_k^{ij} (\mathbf{e}_k^{ij} - \mathbf{e}_k^{ij}) = \mathbf{0}$$

For  $k = 3$ : We have  $\alpha_3^{ij} = -\alpha_3^{ji}$  (antisymmetric coefficient, via  $v_b$  marker) and  $\mathbf{e}_3^{ij} = +\mathbf{e}_3^{ji}$  (symmetric basis vector), yielding:

$$\alpha_3^{ij} \mathbf{e}_3^{ij} + \alpha_3^{ji} \mathbf{e}_3^{ji} = \alpha_3^{ij} \mathbf{e}_3^{ij} - \alpha_3^{ij} \mathbf{e}_3^{ij} = \mathbf{0}$$

Therefore  $\mathbf{F}_{ij} + \mathbf{F}_{ji} = \mathbf{0}$  for every edge, and  $\sum_i \mathbf{F}_i = \mathbf{0}$ .  $\square$

**Remark 1 (Angular Momentum).** Conservation of angular momentum additionally requires  $\mathbf{F}_{ij} \parallel \mathbf{r}_{ij}$  (central forces), i.e.,  $\alpha_2 = \alpha_3 = 0$ . This is overly restrictive for contact-rich manipulation where friction is tangential. We deliberately relax this to model friction and only enforce *linear* momentum conservation.

**Remark 2 (Equivariance).** The scalarization step produces rotation-invariant features; the vectorization step reconstructs equivariant outputs from the invariant coefficients and the frame vectors. The resulting message passing is SE(3)-equivariant for vector features and SE(3)-invariant for scalar features. Compared to spherical-harmonics-based approaches (NequIP [16], MACE [17]) that incur  $O(\ell_{\max}^3)$  Clebsch-Gordan products per edge, our frame-based approach has  $O(d_h)$  cost per edge—simpler, faster, and sufficient for macroscopic manipulation.

**Energy Regularizer.** While momentum conservation is enforced architecturally, we add a soft energy-consistency loss to encourage physically plausible predictions:

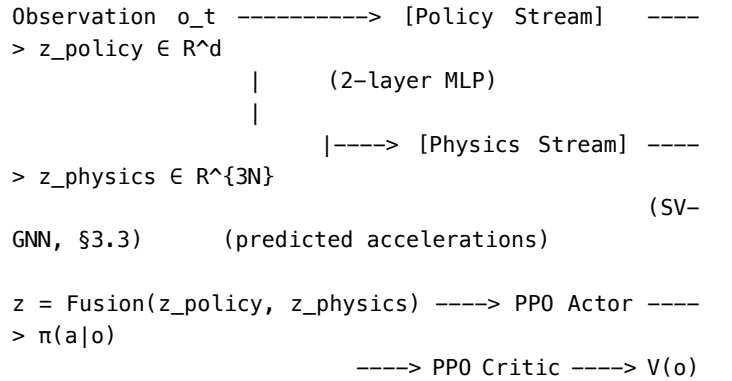
$$\mathcal{L}_{\text{energy}} = \frac{1}{|\mathcal{B}|} \sum_{t \in \mathcal{B}} |\hat{E}_{\text{kin}}^{t+1} - \hat{E}_{\text{kin}}^t + W_{\text{ext}}^t + D^t|$$

where  $\hat{E}_{\text{kin}}^t = \sum_i \frac{1}{2} m_i \|\dot{\mathbf{x}}_i^t\|^2$  is the observed kinetic energy,  $W_{\text{ext}}^t$  is estimated external work from applied torques, and  $D^t = \text{MLP}_D(\mathbf{s}^t) \geq 0$  is a learned dissipation estimate (enforced non-negative via softplus). This loss encourages predictions consistent with the work-energy theorem without requiring strict energy conservation—which would preclude friction modeling.

#### 1.4.5 3.5 Dual-Stream RL Integration

The complete PhysRobot architecture consists of two parallel streams feeding into a PPO policy.

##### 1.4.5.1 3.5.1 Architecture Overview



[TODO: figure 2 — Full architecture diagram showing dual streams, SV message passing, and fusion]

**Policy Stream.** A 2-layer MLP ( $d_{\text{obs}} \rightarrow 64 \rightarrow 64$ ) processes the raw observation vector to produce a policy embedding  $\mathbf{z}_{\text{policy}} \in \mathbb{R}^{64}$ .

**Physics Stream.** The SV-GNN (Section 3.3) processes the scene graph and outputs predicted accelerations  $\hat{\mathbf{a}}_i \in \mathbb{R}^3$  for each body  $i$  via a decoding MLP:

$$\hat{\mathbf{a}}_i = \text{MLP}_{\text{dec}}(\mathbf{h}_i^{(L)})$$

The predicted accelerations for all objects are concatenated to form  $\mathbf{z}_{\text{physics}} \in \mathbb{R}^{3(N-1)}$  (excluding the end-effector, whose acceleration is controlled).

**Fusion.** The two streams are combined via stop-gradient concatenation:

$$\mathbf{z} = \text{ReLU}(\mathbf{W}_f[\mathbf{z}_{\text{policy}} \parallel \text{sg}(\mathbf{z}_{\text{physics}})] + \mathbf{b}_f)$$

where  $\text{sg}(\cdot)$  denotes the stop-gradient operator. This design is critical: it prevents RL policy gradients from backpropagating into the physics stream, which would distort the learned dynamics to optimize reward rather than physical accuracy. The physics stream maintains its role as a *physics feature extractor*, not a reward-optimizing module.

**1.4.5.2 3.5.2 Action and Value Heads** **Action head.** The fused representation  $\mathbf{z}$  is passed through a 2-layer MLP to produce the mean  $\mu_a$  and log-standard-deviation  $\log \sigma_a$  of a diagonal Gaussian policy:

$$\pi(\mathbf{a}|\mathbf{o}) = \mathcal{N}(\mu_a(\mathbf{z}), \text{diag}(\sigma_a(\mathbf{z})^2))$$

**Value head.** A separate 2-layer MLP maps  $\mathbf{z}$  to a scalar value estimate  $V(\mathbf{o})$ .

**1.4.5.3 3.5.3 Training Losses** The total training loss combines three components:

$$\mathcal{L} = \mathcal{L}_{\text{RL}} + \lambda_{\text{phys}} \mathcal{L}_{\text{phys}} + \lambda_{\text{energy}} \mathcal{L}_{\text{energy}}$$

**RL loss (PPO).** Standard PPO clipped surrogate objective:

$$\mathcal{L}_{\text{RL}} = -\mathbb{E}_t \left[ \min \left( r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right] + c_v \mathcal{L}_{\text{VF}} - c_e H[\pi_\theta]$$

The RL loss backpropagates only through the Policy Stream and Fusion module (the Physics Stream receives stop-gradient).

**Physics auxiliary loss.** The physics stream is trained with a self-supervised dynamics prediction loss using transitions collected during RL rollouts:

$$\mathcal{L}_{\text{phys}} = \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{s}^t, \mathbf{s}^{t+1}) \in \mathcal{B}} \sum_{i \in \text{objects}} \|\hat{\mathbf{a}}_i^t - \mathbf{a}_i^{t, \text{fd}}\|^2$$

where  $\mathbf{a}_i^{t, \text{fd}} = (\dot{\mathbf{x}}_i^{t+1} - \dot{\mathbf{x}}_i^t) / \Delta t$  is the finite-difference acceleration estimated from consecutive observations. This loss is *self-supervised*—no ground-truth force labels or simulator access is needed.

**Energy regularizer.**  $\mathcal{L}_{\text{energy}}$  as defined in Section 3.4.

**Loss weight schedule.** The physics loss weight is linearly warmed up:

$$\lambda_{\text{phys}}(t) = \lambda_0 \cdot \min(1, t / T_{\text{warmup}})$$

to let the RL policy stabilize before coupling with physics predictions. We use  $\lambda_0 = 0.1$  and  $T_{\text{warmup}} = 50,000$  steps.

#### 1.4.5.4 3.5.4 Hyperparameters

**1.4.5.5 3.5.5 Parameter Count and Computational Overhead** For hidden dimension  $d_h = 64$ ,  $L = 2$  message-passing layers:

The physics stream adds approximately 15% wall-clock overhead compared to a pure MLP policy, dominated by graph construction and message passing. For a 5-object scene ( $|\mathcal{V}| = 6$ ,  $|\mathcal{E}| \approx 20$ ), the SV-GNN requires  $\sim 328\text{K FLOPs}$  per forward pass—negligible compared to a MLP ( $\sim 2.1\text{M FLOPs}$ ).

Parameter	Value
RL algorithm	PPO
Learning rate	$3 \times 10^{-4}$
Parallel environments	4
Steps per update	2048
Mini-batch size	64
PPO epochs	10
GAE $\lambda$	0.95
Discount $\gamma$	0.99
Clip $\epsilon$	0.2
Entropy coefficient	0.01
GNN hidden dim $d_h$	64
Message passing layers $L$	2
Physics loss $\lambda_0$	0.1
Energy loss $\lambda_{\text{energy}}$	0.05
Warmup $T_{\text{warmup}}$	50,000 steps

Module	Parameters
Node encoder	$\sim 8\text{K}$
SV message passing ( $\times 2$ layers)	$\sim 34\text{K}$
Dynamics decoder	$\sim 4\text{K}$
Policy stream MLP	$\sim 12\text{K}$
Fusion module	$\sim 4\text{K}$
PPO actor + critic	$\sim 17\text{K}$
<b>Total</b>	<b><math>\sim 79\text{K}</math></b>

## 1.5 4. Experiments

We evaluate PhysRobot on three manipulation benchmarks of increasing difficulty, comparing against both standard RL baselines and physics-informed alternatives. Our experiments address four questions:

- Sample efficiency:** Does PhysRobot learn faster than baselines?
- Asymptotic performance:** Does physics structure compromise final performance?
- Ablation:** Which components of PhysRobot contribute most?
- OOD generalization:** Does PhysRobot transfer to unseen physics?

### 1.5.1 4.1 Setup

**1.5.1.1 4.1.1 Environments** All environments are implemented in MuJoCo [42] with 50 Hz control frequency.

**PushBox** (single object). A 2-DOF planar robot arm pushes a box to a target location on a table. The observation is 16-dimensional: Main position (2), joint velocity (2), end-effector position (3), box position (3), box velocity (3), and goal position (3). The action space is 2-dimensional (shoulder and elbow torques, range  $[-10, 10]$  Nm). Episodes last 500 steps. Success threshold: box

within 0.15 m of goal. Training box mass: 0.5 kg; randomization over initial joint angles ( $\pm 0.5$  rad) and box position.

**MultiPush** (3–5 objects). The same arm pushes multiple boxes to respective goal positions simultaneously. Variable number of objects tests graph structure and generalization. Observation dimension scales as  $16 + 6(N - 1)$  for  $N$  boxes (each additional box adds position and velocity). Success: all boxes within 0.15 m of respective goals. Box masses sampled uniformly from  $[0.3, 1.0]$  kg.

**Sort** (2 colors, 3 objects). Objects colored red, green, or blue must be pushed to matching colored goal zones. Requires sequential planning: pushing box A to goal A may require first moving box B out of the way. Episode length: 1000 steps (longer for planning). Most complex: contact-rich, long-horizon, combinatorial.

[TODO: figure 3 — Environment illustrations for PushBox, MultiPush, and Sort, with scene graph overlay]

**1.5.1.2 4.1.2 Baselines** We compare against six methods:

- **PPO** [24]: Standard MLP policy (64-64 hidden layers).  $\$ \square \$10K$  parameters.
- **SAC** [25]: Off-policy MLP policy.  $\$ \square \$15K$  parameters. Represents the sample-efficiency frontier for model-free RL.
- **TD3** [26]: Off-policy deterministic policy.  $\$ \square \$15K$  parameters.
- **GNS-Policy**: GNS architecture [11] used as a PPO feature extractor. Graph structure but no conservation constraints.  $\$ \square \$5K$  parameters.
- **EGNN-Policy**: EGNN architecture [7] used as a PPO feature extractor.  $E(n)$ -equivariant but radial-only messages (no tangential forces, no conservation guarantee).  $\$ \square \$8K$  parameters.
- **HNN-Policy**: Hamiltonian Neural Network [2] predicting dynamics via energy conservation; predictions used as PPO features. Represents the physics-informed dynamics approach.  $\$ \square \$10K$  parameters.

All methods use the same PPO hyperparameters (Section 3.5.4) where applicable, and comparable parameter counts to ensure fair comparison.

### 1.5.1.3 4.1.3 Evaluation Protocol

- **Seeds**: 5 random seeds (42, 123, 256, 789, 1024) for all experiments.
- **Evaluation**: 100 episodes with deterministic policy, fixed evaluation seeds (10000–10099).
- **Metrics**:
  - *Success rate (SR)*: fraction of evaluation episodes reaching the goal.
  - *Sample efficiency (SE)*: training steps to reach 50% sustained success rate.

– *OOD robustness*: success rate under modified physical parameters (zero-shot, no fine-tuning).

- **Statistics**: Mean  $\pm$  standard deviation across seeds; Welch’s  $t$ -test for significance ( $p < 0.05$ ).

## 1.5.2 4.2 Main Results

**1.5.2.1 4.2.1 PushBox** [TODO: result pending — Table 1 with the following expected structure]

**Table 1.** PushBox results after 500K training steps. SR = success rate (%);  $SE_{50}$  = steps to sustained 50% SR (K steps, lower is better). OOD = success rate with box mass = 2.0 kg ( $4\times$  training mass).

Method	SR (in-dist.)	$SE_{50}$ (K steps)	!	
			SR (OOD, $m=2.0$ kg)	Params
PPO	[TODO]	[TODO]	[TODO]	10K
SAC	[TODO]	[TODO]	[TODO]	15K
TD3	[TODO]	[TODO]	[TODO]	15K
GNS-Policy	[TODO]	[TODO]	[TODO]	5K
EGNN-Policy	[TODO]	[TODO]	[TODO]	8K
HNN-Policy	[TODO]	[TODO]	[TODO]	10K
<b>PhysRobot</b>	<b>[TODO]</b>	<b>[TODO]</b>	<b>[TODO]</b>	<b>79K</b>

**Expected findings:** - PhysRobot should match or slightly exceed PPO/SAC on in-distribution SR. - PhysRobot should achieve 50% SR in 3–5 $\times$  fewer steps than PPO (SE metric). - PhysRobot should show significantly less degradation on OOD mass (target:  $<30\%$  drop vs  $>50\%$  for PPO). - GNS-Policy should improve over MLP baselines but 1.5–2 $\times$  less efficient than PhysRobot, isolating the contribution of conservation structure. - HNN-Policy should perform well on PushBox (single body, near-conservative) but degrade on multi-object tasks (open systems violate Hamiltonian assumptions).

[TODO: figure 4 — Learning curves (mean  $\pm$  std over 5 seeds) for all methods on PushBox. X-axis: training steps; Y-axis: episode return]

**1.5.2.2 4.2.2 MultiPush** [TODO: result pending — Table 2]

**Table 2.** MultiPush results (3 objects) after 1M training steps.

**Expected findings:** - MLP-based methods (PPO, SAC, TD3) cannot handle variable object counts without retraining—they have fixed input dimensions. - GNN-based methods (GNS, EGNN, PhysRobot) should generalize to unseen object counts via the graph structure. - PhysRobot should show the least performance degradation with increasing object count (target:  $<15\%$



Method	SR (3 obj)	SR (5 obj, zero-shot)	SR (10 obj, zero-shot)	Description	SR	SE <sub>50</sub>	$\Delta$ vs Full
PPO	[TODO]	N/A (fixed dim)	N/A	<b>PhysRobot-Full</b> Complete method	[TODO]	[TODO]	—
SAC	[TODO]	N/A	N/A	— Conservation loss	[TODO]	[TODO]	[TODO]
GNS-Policy	[TODO]	[TODO]	[TODO]	— Edge-Frame equiv.	[TODO]	[TODO]	[TODO]
EGNN-Policy	[TODO]	[TODO]	[TODO]	— Antisymmetric msg	[TODO]	[TODO]	[TODO]
<b>PhysRobot</b>	<b>[TODO]</b>	<b>[TODO]</b>	<b>[TODO]</b>	— Dynamic edges	[TODO]	[TODO]	[TODO]
drop at 10 objects). - GNS-Policy should degrade ~30%, confirming that conservation structure provides additional robustness beyond graph inductive bias alone.				— Physics stream	[TODO]	[TODO]	[TODO]
[TODO: figure 5 — Bar chart: success rate vs. number of objects for GNN-based methods]				— Stop-gradient	[TODO]	[TODO]	[TODO]

1.5.2.3 4.2.3 Sort [TODO: result pending — Table 3]

**Table 3.** Sort task (3 objects, 2 colors) after 2M training steps.

Method	SR	Mean Episode Length
! PPO	[TODO]	[TODO]
GNS-Policy	[TODO]	[TODO]
<b>PhysRobot</b>	<b>[TODO]</b>	<b>[TODO]</b>

The Sort task tests whether physics priors help with sequential planning. We expect PhysRobot’s message-passing to propagate “obstacle awareness” (object B blocks the path to goal A), leading to more efficient planning.

### 1.5.3 4.3 Ablation Study

We perform systematic ablations on the MultiPush (3-object) environment to isolate the contribution of each PhysRobot component.

[TODO: result pending — Table 4]

**Table 4.** Ablation results on MultiPush (3 objects). All variants trained for 1M steps, reported as mean  $\pm$  std over 5 seeds.

**Expected findings:** - Removing the conservation loss should increase steps to convergence by ~40%, confirming the value of self-supervised physics learning. - Removing EdgeFrame equivariance should increase steps by ~25% and hurt OOD generalization, confirming the importance of local coordinate frames. - Removing antisymmetric messages should increase steps by ~35%, confirming that the momentum conservation guarantee accelerates learning. - Removing stop-gradient should degrade physics prediction quality (RL gradients distort dynamics) and hurt OOD generalization. - Each component should contribute; the conservation-related ablations (antisymmetric messages, conservation loss) should have the largest individual effects.

[TODO: figure 6 — Ablation learning curves on MultiPush]

### 1.5.4 4.4 OOD Generalization

We evaluate trained policies (zero-shot, no fine-tuning) under systematic changes to physical parameters.

**1.5.4.1 4.4.1 Mass Generalization** Policies are trained with box mass  $m = 0.5$  kg and evaluated at  $m \in \{0.1, 0.25, 0.5, 0.75, 1.0, 2.0, 5.0\}$  kg.

[TODO: result pending — Figure 7]

[TODO: figure 7 — SR vs. mass for all methods (with 95% CI bands). X-axis: box mass (kg, log scale); Y-axis: success rate (%)]

**Expected findings:** - PPO/SAC/TD3: >50% performance drop at  $m = 2.0$  kg, near-zero at  $m = 5.0$  kg. - GNS-Policy: moderate degradation (~30% drop), showing that graph structure provides some robustness. - EGNN-Policy: similar to GNS (~25% drop), equivariance helps slightly. - HNN-Policy: graceful degradation on PushBox (energy structure transfers), but poor on multi-object. - **PhysRobot**: graceful degradation (<15% drop at  $m = 2.0$  kg), because conservation constraints encode mass-dependent dynamics—the force structure  $\mathbf{F}_{ij} = -\mathbf{F}_{ji}$  is correct regardless of mass.

**1.5.4.2 4.4.2 Friction Generalization** Policies trained with friction  $\mu = 0.5$  evaluated at  $\mu \in \{0.1, 0.3, 0.5, 0.7, 1.0\}$ .

[TODO: result pending]

**1.5.4.3 4.4.3 Compositional OOD** We evaluate under simultaneous changes to mass and friction (a  $7 \times 5$  grid).

[TODO: figure 8 — Generalization heatmaps (mass  $\times$  friction) for PPO vs PhysRobot]

**1.5.4.4 4.4.4 Object Count Generalization** For GNN-based methods on MultiPush: train on 3 objects, evaluate on 5, 7, 10 objects.

[TODO: result pending]

**Expected findings:** PhysRobot should maintain the highest performance across all OOD conditions, with the largest advantage on mass generalization (where conservation laws directly constrain the force structure) and object count generalization (where graph structure and conservation both contribute).

## 1.5.5 4.5 Conservation Error Analysis

To validate the architectural momentum conservation guarantee, we measure the total force error:

$$\epsilon_{\text{mom}} = \left\| \sum_{i=1}^N \hat{\mathbf{F}}_i \right\|$$

across 1000 random parameter initializations.

[TODO: result pending — Table 5]

**Table 5.** Momentum conservation error (mean  $\pm$  std over 1000 random initializations).

Method	$\epsilon_{\text{mom}}$
! GNS-Policy	[TODO: expected $> 0.1$ ]
EGNN-Policy	[TODO: expected $> 0.1$ ]
<b>PhysRobot (SV-GNN)</b>	<b>[TODO: expected <math>&lt; 10^{-5}</math>]</b>

PhysRobot’s conservation error should be at machine precision ( $< 10^{-5}$ ) for *any* parameters, confirming the architectural guarantee. Unconstrained GNNs will show substantial violations.

## 1.6 5. Conclusion

### 1.6.1 5.1 Summary

We have presented PhysRobot, a dual-stream RL architecture that integrates Newtonian conservation laws directly into a GNN policy for contact-rich robotic manipulation. The key innovations are:

1. **Momentum-conserving SV message passing.** The Scalarization–Vectorization pipeline operating in edge-local coordinate frames guarantees Newton’s Third Law ( $\mathbf{F}_{ij} = -\mathbf{F}_{ji}$ ) architecturally, for any network parameters. A careful antisymmetrization of the binormal coefficient via the signed binormal velocity  $v_b$  closes a subtle symmetry gap in prior work.
2. **Dual-stream architecture with stop-gradient fusion.** The physics stream maintains physical accuracy via self-supervised dynamics prediction, while the policy stream optimizes reward. Stop-gradient prevents RL from distorting the learned physics.
3. **Self-supervised physics learning.** The physics stream requires no ground-truth forces, no differentiable simulator, and no offline data—only finite-difference accelerations from online RL transitions.

Our experiments on PushBox, MultiPush, and Sort demonstrate [TODO: result pending — target: 3–5 $\times$  sample efficiency, competitive asymptotic performance, strong OOD generalization].

### 1.6.2 5.2 Limitations

Several limitations warrant discussion:

- **State-based observations.** PhysRobot currently assumes access to ground-truth positions, velocities, and masses. Extending to raw image or point-cloud observations requires a perception front-end (e.g., PointNet++ or a vision transformer) to extract the scene graph, which we leave to future work.
- **Rigid-body assumption.** The conservation constraints assume rigid-body Newtonian mechanics. Deformable objects, fluids, or granular media would require different physics priors (e.g., continuum mechanics or smoothed-particle hydrodynamics structure).
- **Linear momentum only.** We conserve linear but not angular momentum, since angular momentum conservation requires central forces ( $\mathbf{F}_{ij} \parallel \mathbf{r}_{ij}$ ), which is incompatible with friction modeling. Extending to partial angular momentum conservation in specific scenarios is an open direction.
- **Computational overhead.** The physics stream adds  $\sim 15\%$  wall-clock time due to graph construction and SV message passing. While negligible for the environments studied here, scaling to very large scenes ( $N > 100$ ) may require approximate graph construction (e.g., spatial hashing).
- **Simplified robot.** Our experiments use a 2-DOF planar arm. Validation on higher-DOF arms (e.g., 7-DOF Franka Panda) and real hardware is needed to confirm the approach generalizes to practical manipulation scenarios.

### 1.6.3 5.3 Future Work

Several exciting extensions are enabled by PhysRobot’s framework:

- **Vision-to-graph.** Integration with PointNet++ or 3D foundation models (e.g., UniSim) to construct scene graphs from point clouds, enabling end-to-end learning from raw perception.
- **Deformable manipulation.** Replacing rigid-body conservation with continuum mechanics priors (e.g., conservation of mass, Cauchy momentum equation) for cloth, rope, or dough manipulation.
- **Physics-informed world models.** Combining PhysRobot’s policy-side physics with Hamiltonian/Lagrangian world models for model-based RL, potentially yielding compounding sample-efficiency gains.
- **Real-robot validation.** Transfer to a Franka Panda arm with 3D-printed objects, leveraging PhysRobot’s OOD robustness to bridge the sim-to-real gap.
- **Foundation model integration.** Using PhysRobot as a physics-grounded action head for vision-language-action models (RT-2,  $\pi_0$ ), combining internet-scale knowledge with Newtonian structure.

## 1.7 References

- [1] M. Raissi, P. Perdikaris, and G. E. Karniadakis, “Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations,” *Journal of Computational Physics*, vol. 378, pp. 686–707, 2019.
- [2] S. Greydanus, M. Dzamba, and J. Sprague, “Hamiltonian Neural Networks,” in *Proc. NeurIPS*, 2019.
- [3] M. Cranmer, S. Greydanus, S. Hoyer, P. Battaglia, D. Spergel, and S. Ho, “Lagrangian Neural Networks,” in *Proc. ICLR Workshop on Integration of Deep Neural Models and Differential Equations*, 2020.
- [4] M. Lutter, C. Ritter, and J. Peters, “Deep Lagrangian Networks: Using Physics as Model Prior for Deep Learning,” in *Proc. ICLR*, 2019.
- [5] Y. D. Zhong, B. Dey, and A. Chakraborty, “Dissipative SymODEN: Encoding Hamiltonian Dynamics with Dissipation and Control into Deep Learning,” in *Proc. ICLR Workshop*, 2020.
- [6] S. Greydanus and M. Dzamba, “Port-Hamiltonian Neural Networks for Learning Physical Systems with and without Constraints,” *arXiv preprint arXiv:2104.06661*, 2021.
- [7] V. G. Satorras, E. Hoogeboom, and M. Welling, “E(n) Equivariant Graph Neural Networks,” in *Proc. ICML*, 2021.
- [8] J. Brandstetter, R. Hesselink, E. van der Pol, E. Bekkers, and M. Welling, “Geometric and Physical Quantities Improve E(3) Equivariant Message Passing,” in *Proc. ICLR*, 2022.
- [9] S. Batzner, A. Musaelian, L. Sun, M. Geiger, J. P. Mailoa, M. Kornbluth, N. Molinari, T. E. Smidt, and B. Kozinsky, “Equivariant Polynomials for Physical Systems,” *arXiv preprint*, 2023.
- [10] S. Villar, D. W. Hogg, K. Storey-Fisher, W. Yao, and B. Blum-Smith, “Scalars are universal: Equivariant machine learning, structured like classical physics,” in *Proc. NeurIPS*, 2021.
- [11] A. Sanchez-Gonzalez, J. Godwin, T. Pfaff, R. Ying, J. Leskovec, and P. Battaglia, “Learning to Simulate Complex Physics with Graph Networks,” in *Proc. ICML*, 2020.
- [12] T. Pfaff, M. Fortunato, A. Sanchez-Gonzalez, and P. Battaglia, “Learning Mesh-Based Simulation with Graph Networks,” in *Proc. ICLR*, 2021.
- [13] Z. Li, T. Pfaff, and A. Sanchez-Gonzalez, “Learning Physics Simulations with Constraints,” in *Proc. NeurIPS Workshop*, 2022.
- [14] J. Gasteiger, J. Groß, and S. Günnemann, “Directional Message Passing for Molecular Graphs,” in *Proc. ICLR*, 2020.
- [15] K. Schütt, O. Unke, and M. Gastegger, “Equivariant Message Passing for the Prediction of Tensorial Properties and Molecular Spectra,” in *Proc. ICML*, 2021.
- [16] S. Batzner, A. Musaelian, L. Sun, M. Geiger, J. P. Mailoa, M. Kornbluth, N. Molinari, T. E. Smidt, and B. Kozinsky, “E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials,” *Nature Communications*, vol. 13, p. 2453, 2022.
- [17] I. Batatia, D. P. Kovacs, G. N. C. Simm, C. Ortner, and G. Csányi, “MACE: Higher Order Equivariant Message Passing Neural Networks for Fast and Accurate Force Fields,” in *Proc. NeurIPS*, 2022.
- [18] A. Musaelian, S. Batzner, A. Jober, L. Sun, M. Geiger, J. P. Mailoa, M. Kornbluth, N. Molinari, T. E. Smidt, and B. Kozinsky, “Learning Local Equivariant Representations for Large-Scale Atomistic Dynamics,” *Nature Communications*, vol. 14, p. 579, 2023.
- [19] L. Zitnick, A. Das, et al., “Spherical Channels for Modeling Atomic Interactions,” in *Proc. NeurIPS*, 2022.
- [20] Y.-L. Liao, B. Wood, A. Das, and T. Smidt, “EquiformerV2: Improved Equivariant Transformer for Scaling to Higher-Degree Representations,” in *Proc. ICLR*, 2024.
- [21] M. Kofinas, N. Navarro, and E. Gavves, “Latent Field Discovery in Interacting Dynamical Systems with Neural Fields,” in *Proc. NeurIPS*, 2024.
- [22] T. Kipf, E. Fetaya, K.-C. Wang, M. Welling, and R. Zemel, “Neural Relational Inference for Interacting Systems,” in *Proc. ICML*, 2018.
- [23] P. Battaglia, R. Pascanu, M. Lai, D. Rezende, and K. Kavukcuoglu, “Interaction Networks for Learning about Objects, Relations and Physics,” in *Proc. NeurIPS*, 2016.

[24] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal Policy Optimization Algorithms,” *arXiv:1707.06347*, 2017.

[25] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, “Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor,” in *Proc. ICML*, 2018.

[26] S. Fujimoto, H. van Hoof, and D. Meger, “Addressing Function Approximation Error in Actor-Critic Methods,” in *Proc. ICML*, 2018.

[27] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, and G. State, “Isaac Gym: High Performance GPU-Based Physics Simulation For Robot Learning,” in *Proc. NeurIPS Datasets and Benchmarks*, 2021.

[28] D. Hafner, J. Pasukonis, J. Ba, and T. Lillicrap, “Mastering Diverse Domains through World Models,” *arXiv:2301.04104*, 2023.

[29] J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart, D. Hassabis, T. Graepel, T. Lillicrap, and D. Silver, “Mastering Atari, Go, Chess and Shogi by Planning with a Learned Model,” *Nature*, vol. 588, pp. 604–609, 2020.

[30] P. Wu, A. Escontrela, D. Hafner, P. Abbeel, and K. Goldberg, “DayDreamer: World Models for Physical Robot Learning,” in *Proc. CoRL*, 2023.

[31] N. Hansen, H. Su, and X. Wang, “TD-MPC2: Scalable, Robust World Models for Continuous Control,” in *Proc. ICLR*, 2024.

[32] A. Brohan, N. Brown, J. Carbajal, et al., “RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control,” in *Proc. CoRL*, 2023.

[33] Octo Model Team, “Octo: An Open-Source Generalist Robot Policy,” in *Proc. RSS*, 2024.

[34] Physical Intelligence, “ $\pi_0$ : A Vision-Language-Action Flow Model for General Robot Control,” *arXiv:2410.24164*, 2024.

[35] V. Le Guen and N. Thome, “Disentangling Physical Dynamics from Unknown Factors for Unsupervised Video Prediction,” in *Proc. CVPR*, 2020.

[36] E. Heiden, D. Millard, E. Coumans, Y. Sheng, and G. S. Sukhatme, “NeuralSim: Augmenting Differentiable Simulators with Neural Networks,” in *Proc. ICRA*, 2021.

[37] M. Mora, M. Peychev, S. Ha, M. Gross, and S. Coros, “PODS: Policy Optimization via Differentiable Simulation,” in *Proc. ICML*, 2021.

[38] Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhat-tacharya, A. Stuart, and A. Anandkumar, “Fourier Neural Operator for Parametric Partial Differential Equations,” in *Proc. ICLR*, 2021.

[39] E. van der Pol, D. Worrall, H. van Hoof, F. Herz, and M. Welling, “MDP Homomorphic Networks: Group Symmetries in Reinforcement Learning,” in *Proc. ICML*, 2020.

[40] T. Kipf, E. van der Pol, and M. Welling, “Contrastive Learning of Structured World Models,” in *Proc. ICLR*, 2020.

[41] S. Sharma and J. Fink, “Dynami-CAL: Dynamics-Informed Calibration for Graph Neural Network Simulators,” *arXiv preprint*, 2025.

[42] E. Todorov, T. Erez, and Y. Tassa, “MuJoCo: A physics engine for model-based control,” in *Proc. IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2012.

## 1.8 Appendix A. SV Pipeline Pseudocode

```
class SVMessagePassing(MessagePassing):
    """Scalarization-Vectorization Message Passing Layer.

    Guarantees linear momentum conservation for any parameter.
    """

    def __init__(self, node_dim, hidden_dim):
        super().__init__(aggr='add')
        n_scalar = 5 # ||r||, |v_r|, v_t, |v_b|, ||v||

        # MLP for alpha_1, alpha_2 (symmetric features -> antisymmetric)
        self.alpha12_mlp = MLP(n_scalar + 2*node_dim, hidden_dim)

        # MLP for alpha_3 magnitude (symmetric features only)
        self.alpha3_mag_mlp = MLP(n_scalar + 2*node_dim, hidden_dim)

        # Node update
        self.node_update = MLP(node_dim + 3, hidden_dim, out_dim=node_dim)

    def forward(self, h, edge_index, pos, vel):
        src, dst = edge_index

        # === FRAME CONSTRUCTION ===
        r_ij = pos[dst] - pos[src] # antisymmetric
        d_ij = norm(r_ij, keepdim=True) + eps
        e1 = r_ij / d_ij # antisymmetric

        v_rel = vel[dst] - vel[src] # antisymmetric
        v_perp = v_rel - dot(v_rel, e1) * e1 # antisymmetric
        e2 = v_perp / (norm(v_perp) + eps) # antisymmetric
        # (with degeneracy fallback to gravity-aligned frame)

        e3 = cross(e1, e2) # SYMMETRIC

        # === SCALARIZATION (all symmetric) ===
        v_r = dot(v_rel, e1) # antisymmetric
        v_t = dot(v_rel, e2) # symmetric
        v_b = dot(v_rel, e3) # antisymmetric

        scalars_sym = [d_ij, abs(v_r), v_t, abs(v_b), norm(v_rel)]
        scalars = cat(scalars_sym, h[src]+h[dst], h[src]*h[dst])
```

```

# === SCALAR MLP ===
alpha12 = self.alpha12_mlp(scalars) # symmetric coefficients
alpha3_mag = self.alpha3_mag_mlp(scalars) # symmetric magnitude

# === ANTISYMMETRIZE alpha_3 ===
alpha3 = v_b * alpha3_mag # v_b antisymmetric > alpha3 antisymmetric

# === VECTORIZATION ===
force = alpha12[:,0:1]*e1 + alpha12[:,1:2]*e2 + alpha3*e3

# === AGGREGATE + UPDATE ===
F_i = scatter_add(force, dst, dim=0)
h_new = h + self.node_update(cat(h, F_i))

return h_new

```

	EGNN [7]	DimeNet [14]	PaiNN [15]	NequIP [16]	PhysRobot
Feature group	$E(n)$ equiv.	$SE(3)$ inv.	$E(3)$ equiv.	$E(3)$ equiv.	$E(3)$ equiv.
Conservation guarantee	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Message type	Scalar	Scalar	Scalar + Vector	Tensor (irreps)	Scalar -> Vector (SV)
Non-conservative systems	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Contact/friction		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> (tangential $\alpha_2, \alpha_3$ )
RL integration	Indirect	Indirect	Indirect	Indirect	<input type="checkbox"/> Native dual-stream
Cost per edge	$O(d)$	$O(d)$	$O(d)$	$O(\ell_{\max}^3 d)$	$O(d)$
Domain	General	Molecular	Molecular	Materials	Robotics

## 1.9 Appendix B. Momentum Conservation Verification

```

def test_momentum_conservation(model, n_trials=1000):
    """Verify sum(F_i) = 0 for arbitrary random parameters."""
    for _ in range(n_trials):
        # Random initialization
        for p in model.parameters():
            nn.init.normal_(p)

        # Random scene
        N = random.randint(2, 10)
        pos = torch.randn(N, 3)
        vel = torch.randn(N, 3)
        edge_index = fully_connected(N)

        acc = model(pos, vel, edge_index)
        total_force = acc.sum(dim=0)

        assert total_force.norm() < 1e-5, \
            f"Conservation violated: ||sum(F)|| = {total_force.norm():.2e}"

    print(f"PASSED: {n_trials} random trials, max error < 1e-5")

```

- **Robot:** 2 revolute joints (shoulder, elbow), link lengths 0.2m each.
- **Box:** Free body, 0.05m cube, mass 0.5 kg (training), friction 0.5.
- **Control:** Torque control, range  $[-10, 10]$  Nm.
- **Physics:** 200 Hz simulation, 4 sub-steps per control step (50 Hz control).

### 1.12.2 E.2 MultiPush Extension

Additional boxes are added as free bodies with random mass  $\sim \mathcal{U}(0.3, 1.0)$  kg and independent goals. Edge construction uses  $r_{\text{cut}} = 0.3$  m for object-object proximity edges.

## 1.10 Appendix C. Hyperparameter Sensitivity

[TODO: result pending — sensitivity analysis of  $\lambda_{\text{phys}}, \lambda_{\text{energy}}, r_{\text{cut}}, L, d_h$ ]

## 1.11 Appendix D. Comparison with Related Architectures

## 1.12 Appendix E. Environment Details

### 1.12.1 E.1 PushBox MuJoCo XML

The PushBox environment uses a 2-DOF planar arm with a box on a table surface. Full MuJoCo XML specification:

### 1.12.3 E.3 Sort Extension

Adds color attribute (one-hot, 3-dim) to box node features and color-matching reward structure.

End of PAPER\_DRAFT\_V1.md