# GRL-SNAM: GEOMETRIC REINFORCEMENT LEARN-ING WITH DIFFERENTIAL HAMILTONIANS FOR NAVI-GATION AND MAPPING IN UNKNOWN ENVIRONMENTS

### **Anonymous authors**

000

001

002

004

006

008 009 010

011 012

013

014

016

017

018

019

021

023

024

025

026

027

028

029

031

032

037

040

041

043

044

045 046

047

051

052

Paper under double-blind review

### **ABSTRACT**

We present GRL-SNAM, a geometric reinforcement learning framework for Simultaneous Navigation and Mapping in unknown environments. GRL-SNAM differs from traditional SLAM and other reinforcement learning methods by relying exclusively on local sensory observations without constructing a global map. Our approach formulates navigation and mapping as coupled dynamics on generalized Hamiltonian manifolds: sensory inputs are translated into local energy landscapes that encode reachability, obstacle barriers, and deformation constraints, while policies for sensing, planning, and reconfiguration evolve stagewise under Differential Policy Optimization (DPO). A reduced Hamiltonian serves as an adaptive score function, updating kinetic/potential terms, embedding barrier constraints, and continuously refining trajectories as new local information arrives. We evaluate GRL-SNAM on 2D deformable navigation tasks, where a hyperelastic robot learns to squeeze through narrow gaps, detour around obstacles, and generalize to unseen environments. We evaluate GRL-SNAM on procedurally generated 2D deformable-robot tasks (hyperelastic ring) with narrow gaps and clutter, comparing against *local reactive* baselines (PF, CBF, staged DWA) and *global* A\* references (rigid, clearance-aware) under identical stagewise sensing constraints. GRL-SNAM matches near-CBF path quality while using the minimal map coverage of PF, preserves clearance, generalizes to unseen layouts, and demonstrates that Hamiltonian-structured RL enables high-quality navigation through minimal exploration via local energy refinement rather than global mapping.

# 1 Introduction

Reinforcement learning has achieved remarkable successes in high-dimensional control, yet its application to real-world continuous navigation remains fundamentally limited. Long-horizon reasoning, multi-scale decision making, and online adaptation pose challenges that overwhelm existing methods. Model-free RL consumes millions of interactions, while hierarchical variants introduce brittle complexity. In simultaneous navigation and mapping (SNAM), where agents must traverse and construct evolving environmental representations, these limitations become prohibitive.

At its core, the difficulty arises because conventional RL policies are *structureless*. They treat navigation as black-box optimization, ignoring the geometric and physical principles that make locomotion stable, adaptive, and safe. Without inductive bias, policies overfit training environments, fail under distribution shift, and collapse during long rollouts.

# 1.1 BEYOND BELLMAN OPTIMIZATION: PURELY FEEDFORWARD CONTROL:

Our framework does not optimize a value function via the Bellman equation. Standard RL algorithms hinge on recursive bootstrapping for estimating returns, propagating value updates, and iteratively improving policies. This induces high sample complexity, instability, and delayed credit assignment, especially in navigation with long horizons.

In contrast, our approach is **purely feedforward**: policies emerge as direct gradient flows of Hamiltonian energies, without value iteration. Navigation decisions are computed in a single pass from

local sensory input and the reference Hamiltonian, bypassing dynamic programming altogether. This eliminates the need for rollout-based value propagation, yielding stable training, low variance adaptation, and sample-efficient online updates.

### 1.2 KEY INSIGHT: HAMILTONIAN STRUCTURE AS NAVIGATION INDUCTIVE BIAS

We propose addressing these limitations by grounding RL in **Hamiltonian mechanics**. Our central insight is that navigation can be framed as learning energy functionals:

$$\mathcal{H}(q,p) = K(p) + P(q) \tag{1}$$

where kinetic and potential energies encode control objectives, constraints, and adaptation strategies. This formulation introduces three structural advantages:

(1) **Energy conservation** stabilizes long-horizon rollouts by preventing accumulation of numerical errors. (2) **Symplectic geometry** naturally separates fast reactive dynamics from slow strategic planning, addressing multi-scale temporal coordination. (3) **Barrier encoding** integrates safety and collision avoidance directly into potential functions, eliminating fragile reward shaping.

Hamiltonian structure transforms policy optimization into Differential Policy Optimization (DPO) Bajaj & Nguyen (2024), where policies emerge as gradient flows of learned energies that respect geometry, conserve invariants, and generalize across environments. Since DPO has already been shown to outperform state-of-the-art policy learners across standard control benchmarks, we focus our comparisons on navigation and mapping baselines (e.g., PF, CBF, A\*) rather than reconstructing weaker policy-gradient baselines ourselves. Importantly, no prior work aligns directly with our Hamiltonian formulation; implementing such policy baselines within our framework would amount to re-developing them as part of our contribution, rather than evaluating against an existing standard.

### 1.3 OFFLINE-ONLINE HAMILTONIAN SYNERGY

We distinguish between complementary learning regimes that exploit this geometric structure:

Offline learning discovers reference Hamiltonians  $h^{\theta^*}$  trained on trajectory data, capturing fundamental multi-scale navigation dynamics in local frames. These provide stable geometric priors encoding essential coupling between sensing, planning, and deformation.

Online adaptation fuses new environmental context into learned Hamiltonians through contextual corrections:  $h^{\rm adapted} = h^{\rm ref} + \Delta h^{\rm context}$  This creates conservative adaptation: systems default to learned physics-based behaviors while adding minimal corrections for environmental variations.

The synergy transforms every offline policy into a *reference Hamiltonian* and every online update into a *geometric alignment step*. Navigation emerges from meta-policies that parse environments, assemble energy landscapes, and integrate them through symplectic dynamics.

# 1.4 Contributions:

This work establishes **GRL-SNAM** as a new approach beyond structureless policy learning. Our contributions are:

- 1. **Hamiltonian RL framework**: Allows adaptable integration of classical mechanics into RL for navigation, treating rewards as energies and policies as symplectic flows.
- Multi-scale geometric coordination: Differential policies for sensing, planning, and adaptation unified through shared energy formulations, achieving temporal scale separation without manual hierarchy design.
- 3. **Physics-grounded adaptation**: Principled offline-online decomposition where stable reference dynamics adapt through geometric alignment rather than catastrophic relearning.
- 4. **Theoretical guarantees**: Symplectic structure preservation ensures stability, while independent policy learning achieves linear sample complexity scaling.
- 5. **Empirical validation**: Hyperelastic ring navigation demonstrates superior sample efficiency and generalization compared to A\* and CBF baselines.

# 2 RELATED WORK

We focus on structure-preserving, deployable navigation with deformable bodies. Our work intersects advances in geometric learning, safety-critical control, and deformable robot navigation.

Mathematical Foundations for RL Navigation. Most navigation RL methods operate in Euclidean spaces using standard PPO Schulman et al. (2017) or TD3 Fujimoto et al. (2018) formulations without geometric constraints. Geometric approaches include SE(3) equivariant policies Hoang et al. (2025) for manipulation and Riemannian safe navigation Klein et al. (2023) using tangent space projections. Hamiltonian neural networks Desai et al. (2021) demonstrate superior learning dynamics through symplectic structure but remain limited to simple control tasks.

**Safety-Critical Navigation.** Control Barrier Function (CBF) integration with RL achieves formal safety guarantees Li et al. (2023), but treats constraint satisfaction as orthogonal to navigation optimality, often resulting in conservative behaviors. Our Hamiltonian formulation integrates safety constraints directly within the energy structure.

**Deformable Robot Navigation.** Recent work demonstrates ring-like navigation through preprogrammed strategies: aerial gap navigation via fixed Liquid Crystal Elastomer responses Qi et al. (2024) and HAVEN Mulvey & Nanayakkara (2024) using predetermined shape-changing sequences. These approaches rely on offline parameter optimization followed by deterministic execution—they cannot adapt deformation strategies online as environmental conditions change.

**Neural Scene Representations.** NeRF-based SLAM methods like NICE-SLAM Zhu et al. (2022) provide rich environmental representations that complement our energy-based navigation formulation by supplying obstacle and free-space information for barrier and goal potential computation.

**Simultaneous Navigation and Mapping:** Most SNAM approaches prioritize building detailed maps before navigation. SGoLAM Kim et al. (2021) couples goal localization with occupancy mapping, CMP Gupta et al. (2019) integrates a differentiable planner into learned mapping, and CL-SLAM Vödisch et al. (2023) maintains maps for long-term adaptability. In contrast, our GRL-SNAM framework aims to *reach goals via high-quality, well-weighted paths while mapping as little of the environment as possible.* To our knowledge, no prior work explicitly targets minimal exploration; our method introduces progressive path refinement, continually improving least-cost trajectories as new observations arrive.

**Positioning.** GRL-SNAM addresses key gaps by extending Hamiltonian mechanics from simple control to complex navigation requiring sensing, planning, and deformation. Unlike existing methods that require manual task decomposition or rely on pre-programmed strategies, our differential multi-policy architecture learns specialized policies naturally coupled through shared Hamiltonian energy formulations. The symplectic structure ensures stable coordination across temporal scales with formal convergence guarantees, bridging theoretically principled geometric methods with practical navigation frameworks.

### 3 METHODOLOGY

We present GRL-SNAM (Geometric Reinforcement Learning for Spatial Navigation and Manipulation): a Hamiltonian-structured navigator that unifies offline physics learning with online adaptive correction through black-box modular policies. The code for this paper is available at: Code

### 3.1 PROBLEM FORMULATION: NAVIGATION AS HAMILTONIAN OPTIMIZATION

We formulate navigation in unknown environments as energy minimization over symplectic manifolds. Consider a deformable robot with state  $q_t = (c_t, \theta_t, \psi_t)$  navigating from  $\mathbf{x}_0$  to  $\mathbf{x}_g$  through unknown obstacles characterized by binary occupancy  $I : \mathbb{R}^2 \to \{0, 1\}$ .

**Energy Decomposition by Policy.** The GRL-SNAM framework interprets navigation as minimizing a Hamiltonian reward functional, where each policy governs a distinct energy term as shown in Figure 6:

$$\mathcal{R}(q, \mathcal{C}_t) = \underbrace{-\beta \| \boldsymbol{c}(t) - \mathbf{x}_g \|_2^2}_{\text{Goal Attraction (FPE)}} + \underbrace{-\| \boldsymbol{y}(t) \|_A^2}_{\text{Sensor Cost (Sensor Policy)}} + \underbrace{-\lambda_{obj} \mathcal{E}_{obj}(q(t))}_{\text{Deformation Energy (Reconfig Policy)}} + \underbrace{-\sum_{i=1}^{C_t} \alpha_i b(\tilde{d}_i, \hat{d})}_{\text{Collision Barriers (FPE)}}.$$

The Sensor Policy contributes the sensor cost, regularizing information acquisition. The FPE (Frame-Planning Executor) governs goal attraction and barrier avoidance, balancing reachability with safety. The Reconfig Policy governs deformation, enabling radius modulation for narrow passages.

This decomposition highlights that the total Hamiltonian is not a monolithic reward but a structured sum of physically interpretable energies, each attached to a specialized policy.

### 3.2 Hyperelastic Ring Robot Model

We model the deformable robot as a hyperelastic ring 7 with reduced-order dynamics enabling efficient navigation while capturing essential deformation behaviors.

The robot state consists of three generalized coordinates: uniform scale s(t) controlling size, center position  $o(t) \in \mathbb{R}^2$  for translation, and orientation  $\theta(t)$  for rotation. The boundary is represented as a periodic B-spline with control points transformed via similarity transformation. This reduced-order representation captures the essential deformation modes (compression in tight passages, relaxing in open areas) while maintaining computational tractability.

### 3.3 BLACK-BOX MODULAR ARCHITECTURE

Rather than learning monolithic navigation policies, we decompose the problem into three independent score functions, each dedicated to a specific navigation aspect:

**Definition 3.1** (Independent Score Functions). Let  $K = \{y, f, o\}$  denote the set of policy indices corresponding to sensor, frame, and object domains respectively. For each  $k \in K$ , define:

- $z_k \in \mathcal{Z}_k$ : the phase space state for policy k, where  $\mathcal{Z}_k = \mathcal{Q}_k \times \mathcal{P}_k$  with configuration space  $\mathcal{Q}_k$  and momentum space  $\mathcal{P}_k$
- $C_t$ : the set of active environmental constraints at time  $t \in \mathbb{R}_{\geq 0}$
- $\theta_k \in \Theta_k$ : the learnable parameters for policy k, where parameter sets satisfy disjointness:  $\Theta_i \cap \Theta_j = \emptyset$  for  $i \neq j$
- $h_k^{\theta_k}: \mathcal{Z}_k \times \mathcal{C}_t \times \mathbb{R}_{\geq 0} \to \mathbb{R}$ : a learned energy functional parameterized by  $\theta_k$

Each policy  $\pi_k$  is defined as an independent score function:  $s_k^{\theta_k}(z_k, \mathcal{C}_t, t) = \nabla_{z_k} h_k^{\theta_k}(z_k, \mathcal{C}_t, t)$ 

The parameter disjointness ensures independence:  $\frac{\partial s_k^{\theta_k}}{\partial \theta_j} = 0$  for all  $j \neq k$ 

allowing parallel training while maintaining coordination through shared constraints  $C_t$ .

**Policy Abstraction.** Each policy is treated as a black box that:

- Sensor Policy  $(\pi_y)$ : Adapts perception parameters  $\rightarrow$  energy gradients for information gathering
- Frame Policy  $(\pi_f)$ : Plans collision-free paths  $\rightarrow$  energy gradients for goal attraction
- Shape Policy  $(\pi_o)$ : Controls robot deformation  $\rightarrow$  energy gradients for obstacle navigation

The key insight is that our Navigator is agnostic to policy implementation—our contribution is the Hamiltonian structure binding them together through dynamic constraint sets  $C_t$ .

Algorithm 3 details the online adaptation procedure, where the navigator issues sequential queries to the sensor, frame, and reconfig policies, integrates their energy gradients into a Hamiltonian

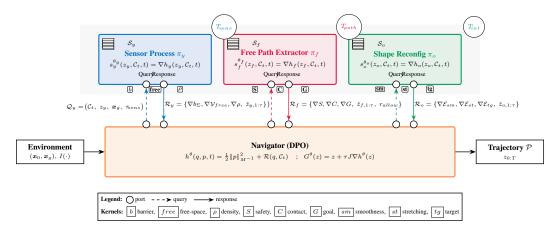


Figure 1: Independent score function architecture and query-response interface. The Navigator (DPO) issues policy-specific queries  $Q_k$  using active constraints  $C_t$ , current states, and horizons; each policy  $\pi_k$  computes energy gradients via its learned score  $s_k^{\theta_k} = \nabla h_k$ , backed by a dynamic spatial index  $S_k$  for efficient  $O(\log n)$  neighbor queries. Policies return responses  $\mathcal{R}_k$  containing gradients and predicted rollouts, which the Navigator integrates to advance the joint dynamics. Spatial indices enable multi-kernel reuse across energy computations.

update, and applies meta-corrections for contextual alignment to generate stable trajectories in novel environments.

### 3.4 NAVIGATOR AS META-HAMILTONIAN LEARNER

The Navigator operates as a meta-learning system that coordinates multi-scale policies by learning how to update their Hamiltonian energy functions rather than directly manipulating phase space states.

**Meta-Learning Formulation.** Let  $\Theta^t = \{\theta_y^t, \theta_f^t, \theta_o^t\}$  denote the current policy parameters, and let  $\Pi^t$  be their conjugate momenta. The Navigator defines a meta-Hamiltonian

$$\mathcal{H}^{\phi} \cdot (\Theta \ \Pi \cdot \mathcal{R}) \mapsto \mathbb{R}$$

 $\mathcal{H}_{\mathrm{nav}}^{\phi}:\,(\Theta,\Pi;\,\mathcal{R})\mapsto\mathbb{R},$  where the responses  $\mathcal{R}^t=\{\mathcal{R}_y^t,\mathcal{R}_f^t,\mathcal{R}_o^t\}$  are exogenous query responses/targets . A concrete and useful choice is

$$\mathcal{H}_{\text{nav}}^{\phi}(\Theta, \Pi; \mathcal{R}) = \frac{1}{2} \Pi^{\top} \mathbf{G}^{-1}(\Theta) \Pi + \frac{1}{2} (f(\Theta) - y_{\text{tgt}}(\mathcal{R}))^{\top} \mathbf{W} (f(\Theta) - y_{\text{tgt}}(\mathcal{R})), \tag{3}$$

where  $f(\Theta) \in \mathbb{R}^m$  collects observables (e.g., clearance, distance, speed bands),  $\mathbf{W} \succeq 0$  weights their importance, and  $\mathbf{G}(\Theta) \succ 0$  is a metric on parameter space (a "mass" for  $\Theta$ ).

The canonical meta-dynamics are

$$\dot{\Theta} = \nabla_{\Pi} \mathcal{H}_{\text{nav}}^{\phi} = \mathbf{G}^{-1}(\Theta)\Pi, \qquad \dot{\Pi} = -\nabla_{\Theta} \mathcal{H}_{\text{nav}}^{\phi} = -\frac{1}{2} \nabla_{\Theta} \left[ \Pi^{\top} \mathbf{G}^{-1} \Pi \right] - \mathbf{J}(\Theta)^{\top} \mathbf{W} \left( f(\Theta) - y_{\text{tgt}} \right). \tag{4}$$

with  $J(\Theta) = \partial f/\partial \Theta$ . In the quasi-overdamped (or implicit-momentum) regime, eliminating  $\Pi$ yields the Gauss-Newton update

$$\Delta\Theta = -h \mathbf{G}^{-1}(\Theta) \mathbf{J}(\Theta)^{\mathsf{T}} \mathbf{W} \left( f(\Theta) - y_{\text{tgt}}(\mathcal{R}) \right). \tag{5}$$

This makes the Navigator itself a Hamiltonian system governing energy function updates.

**Sequential Query-Response Protocol.** At each timestep t, the Navigator orchestrates a sequential coordination process:

Stage 1 - Sensor Query: Navigator queries sensor policy for environmental constraints:

$$Q_y^t = (C_{t-1}, z_y^t, \mathbf{x}_g, \tau_{\text{sens}})$$
(6)

$$\mathcal{R}_{y}^{t} = s_{y}^{\theta_{y}^{t}}(\mathcal{Q}_{y}^{t}) = \{\nabla h_{y}^{\theta_{y}^{t}}, \mathcal{C}_{t}^{\text{updated}}\}$$
 (7)

Stage 2 - Frame Query: Using updated constraints  $C_t^{\text{updated}}$ :

$$\mathcal{Q}_{f}^{t} = (\mathcal{C}_{t}^{\text{updated}}, z_{f}^{t}, \mathbf{x}_{g}, \tau_{\text{path}}) \tag{8}$$

$$\mathcal{R}_f^t = s_f^{\theta_f^t}(\mathcal{Q}_f^t) = \{ \nabla h_f^{\theta_f^t}, \mathcal{W}_t \}$$
 (9)

Stage 3 - Shape Query: Using waypoints  $W_t$  from frame policy:

$$Q_o^t = (C_t^{\text{updated}}, z_o^t, W_t, \tau_{\text{int}})$$
(10)

$$\mathcal{R}_o^t = s_o^{\theta_o^t}(\mathcal{Q}_o^t) = \{ \nabla h_o^{\theta_o^t}, \text{ aux data} \}$$
 (11)

**Meta-Update Integration.** The Navigator processes all responses to compute Hamiltonian parameter updates:

$$\theta_k^{t+1} = \theta_k^t + h\Delta\theta_k^t$$

where  $\Delta \theta_k^t = \nabla_{\mathcal{R}_k} \mathcal{H}_{\text{nav}}^{\phi}(\Theta^t, \mathcal{R}^t)$  for  $k \in \{y, f, o\}$ .

**State Evolution.** Individual policies then update their phase space states using updated energy functions:

$$\boldsymbol{z}_k^{t+1} = \boldsymbol{z}_k^t + \tau_k J_k \boldsymbol{s}_k^{\boldsymbol{\theta}_k^{t+1}}(\boldsymbol{z}_k^t, \mathcal{C}_t^{\text{updated}}, t)$$

The Navigator acts as a stateless mapper:  $(\Theta^t, \mathcal{R}^t) \mapsto \Theta^{t+1}$ , learning how energy landscapes should evolve based on policy feedback without maintaining internal memory. This meta-Hamiltonian formulation ensures that parameter updates respect geometric structure while enabling coordinated adaptation across all policies.

Contractual Interface. The system maintains a clear separation of concerns:

- Task policies: Given  $h_k^{\theta_k}$ , generate score functions  $s_k^{\theta_k} = \nabla h_k^{\theta_k}$
- Navigator: Given  $\{\theta_k^t, \mathcal{R}_k^t\}$ , learn optimal updates  $\Delta \theta_k^t$
- **Environment**: Provides constraints  $\mathcal{C}_t$  and state evolution through symplectic integration

# 3.5 MULTI-SCALE TEMPORAL COORDINATION

The policies operate at natural temporal hierarchies, creating stable multi-scale coordination:

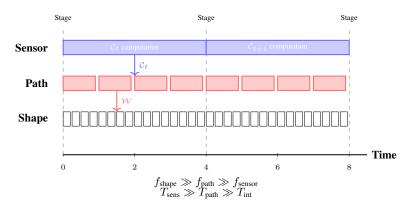


Figure 2: Temporal hierarchy. Sensor policy operates at low frequency (once per stage), establishing environmental constraints  $C_t$ . Path policy operates at medium frequency within each stage, computing waypoints W. Shape policy operates at high frequency, continuously adapting at each integration step. This creates a natural hierarchy where slow sensor updates provide stable constraints for faster path and shape adaptations.

This temporal separation enables a **nested quasi-static approximation**: the fastest dynamics (reconfiguration) equilibrate within each frame update, and frame dynamics settle before the slower sensor policy evolves. This hierarchy prevents destabilizing interactions across timescales while preserving the necessary coupling for coherent, coordinated behavior.

### 3.6 OFFLINE PHYSICS LEARNING VS ONLINE ADAPTIVE CORRECTION

Our approach resolves the fundamental tension between learning complex dynamics and real-time adaptation through principled decomposition:

# $\begin{array}{c} \textbf{Standard RL} & \textbf{Our GRL-SNAM} \\ \\ \hline \textbf{Offline: Learn Policy} \\ \pi(a|s) \text{ from dataset} & \textbf{Offline: Learn Hamiltonian} \\ \hline \textbf{Online: Fine-tune policy} \\ \text{on new environment} & \hline \textbf{Online: Contextual alignment} \\ \hline \textbf{Challenge: Policy transfer} \\ \text{across domains} & \hline \textbf{Advantage: Physics structure} \\ \hline \textbf{ensures stable adaptation} \\ \hline \end{array}$

Figure 3: Comparison between standard RL offline/online adaptation and our physics-grounded approach. Standard methods learn arbitrary policies and struggle with transfer, while our approach learns physically meaningful Hamiltonians that naturally adapt to environmental variations.

**Offline Reference Learning:** Train policies on clean trajectory data to learn fundamental multiscale navigation dynamics:

$$h_k^{\text{ref}}(z_k, \mathcal{C}_t^{\text{clean}}, t) = \frac{1}{2} \|p_k\|_{M_k^{-1}}^2 + \mathcal{R}_k^{\text{intrinsic}}(q_k, \mathcal{C}_t, t)$$

Online Contextual Adaptation: Adapt to novel constraints through energy corrections:

$$h_k^{\text{adapted}} = h_k^{\text{ref}} + \underbrace{\alpha_k \mathcal{F}_k^{\text{interp}}}_{\text{similar contexts}} + \underbrace{\beta_k \mathcal{G}_k^{\text{barrier}}}_{\text{novel constraints}}$$

This creates conservative adaptation: default to learned physics behaviors, add minimal corrections for environmental variations.

# 3.7 THEORETICAL PROPERTIES

Our framework provides three key theoretical guarantees:

**Theorem 3.2** (Multi-Policy Stability). Under temporal scale separation  $T_{sens} \gg T_{path} \gg T_{int}$  and bounded parameter updates, the coupled system maintains stability with error bound  $\mathcal{E}_{total} \leq \epsilon$ .

**Theorem 3.3** (Symplectic Preservation). Each score function generates symplectic dynamics preserving the canonical structure  $\omega_k(z_{k,t+1}) = \omega_k(z_{k,t})$ .

**Theorem 3.4** (Linear Sample Complexity). *Independent training achieves total sample complexity*  $N_{total} = \sum_{k \in \{y,f,o\}} O(\epsilon_k^{-(2d_k+4)})$ , linear in the sum of policy dimensions rather than exponential in joint dimensionality.

We defer the proof of theorems in appendix. The system thinks in physics during offline training but adapts through energy corrections during online execution, combining principled dynamics stability with real-world deployment flexibility.

### 4 EXPERIMENTAL EVALUATION

We evaluate GRL-SNAM across multiple dimensions that highlight the unique capabilities of our geometric approach compared to standard reinforcement learning and classical navigation methods. Our evaluation encompasses task performance, safety guarantees, and learning efficiency under minimal sensing constraints. For more detailed results and analysis, refer to Appendix I

**Experimental Setup:** We evaluate GRL-SNAM in procedurally generated 2D deformable navigation tasks, where a hyperelastic ring must traverse cluttered environments with narrow gaps and

Table 1: Navigation quality comparison (success-only runs). GRL-SNAM achieves near-CBF efficiency with minimal mapping budget.

Method	SPL↑	Detour ↓	Min. Clearance (m) ↑	Mapping Ratio (%) ↓
PF	0.77	1.42	0.18	10.3
CBF	0.96	1.04	0.32	11.2
GRL-SNAM	0.95	1.09	0.26	10.7

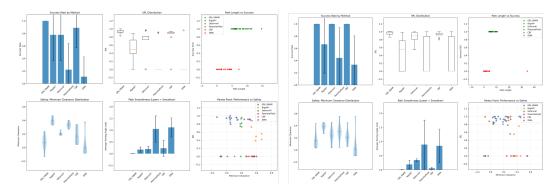


Figure 4: **Main performance comparison.** GRL-SNAM achieves superior success rates, path efficiency (SPL), and smoothness while maintaining safety margins. Classical and reactive baselines show significant degradation in complex environments.

varying obstacle densities. The robot perceives only a local window of size  $2\hat{d} \times 2\hat{d}$ , from which we construct a Hamiltonian energy functional with goal-directed potential  $F_g$ , barrier potentials  $F_{bs}$ , and adaptive coefficients  $(\beta, \gamma, \alpha)$  modulated by context encoders.

**Baselines.** We compare against two categories under matched information constraints: **Global planning:** Rigid A\* (obstacle inflation) and Deformable A\* (clearance-aware penalty) and **Local reactive:** Potential Field (PF), Control Barrier Functions (CBF), and staged DWA using identical local windows and stage management as GRL-SNAM

**Metrics.** Success Rate, Success-weighted Path Length (SPL), Detour Ratio, Minimum Clearance, Path Smoothness, Collisions, and Mapping Ratio (fraction of environment observed).

### 4.1 MAIN RESULTS

**Q1.** How efficiently does GRL-SNAM trade mapping for navigation quality? Table 3 demonstrates that GRL-SNAM achieves CBF-level navigation quality (SPL = 0.95, Detour = 1.09) while using the same minimal map coverage as PF (10.7% vs CBF's 11.2%). This validates that our stagewise Hamiltonian refinement extracts maximum value per sensed unit of the environment.

Q2. Does GRL-SNAM outperform classical and reactive planners in complex environments? Yes. Figure 4 shows GRL-SNAM achieves near-perfect success rates ( $\approx 100\%$ ) across both indistribution and out-of-distribution test cases, while all baselines degrade significantly. GRL-SNAM consistently maintains high SPL ( $\approx 1.0$ ) with low variance and produces the smoothest trajectories with lowest turning angles. The Pareto frontier analysis confirms GRL-SNAM uniquely dominates the safety-performance trade-off.

Q3. How does the Hamiltonian formulation enable coherent navigation? Figure 5 illustrates how GRL-SNAM unifies goal attraction  $F_g$  and barrier repulsion  $F_{bs}$  into a coherent navigation field through adaptive coefficients. Unlike reactive methods that treat forces independently, our differential composition  $F = \beta F_g + \gamma F_{bs}$  creates contextually balanced dynamics that simultaneously pursue goals and avoid obstacles.

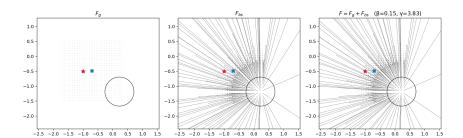


Figure 5: Hamiltonian force field composition. Left: goal force  $F_g$ ; Middle: barrier forces  $F_{bs}$ ; Right: adaptive combination yielding safe, goal-directed trajectories.

**Q4.** What distinguishes GRL-SNAM's online adaptation from standard RL approaches? Unlike standard policies that adjust actions online, GRL-SNAM modifies the *entire local energy landscape* as new obstacles are sensed. Figure 10 demonstrates that coefficients  $(\beta, \gamma, \alpha)$  evolve dynamically to redefine the reduced Hamiltonian itself, ensuring energy-consistent posterior updates rather than heuristic reactive adjustments.

### 4.2 Additional Evaluations

We conducted comprehensive ablation studies on loss components ( $\mathcal{L}_{friction}$ ,  $\mathcal{L}_{multi}$ ) confirming that friction matching is critical for stability while multi-start robustness prevents over-conservatism. Robustness evaluations under sensor noise and dynamics perturbations show graceful degradation (87% success under severe noise vs 99% nominal) due to our adaptive Hamiltonian framework. Sample efficiency analysis demonstrates faster convergence than RL baselines due to physics-informed structure.

### 4.3 KEY INSIGHTS

Minimal mapping suffices: GRL-SNAM achieves optimal navigation quality using  $\sim 10\%$  environment coverage, validating the core SNAM principle that local geometric structure contains sufficient information for global navigation tasks.

**Hamiltonian unification:** The differential geometric formulation naturally balances competing objectives (goal-seeking, obstacle avoidance, smoothness) through principled energy minimization rather than heuristic weight tuning.

**Principled online adaptation:** By modifying the energy landscape itself rather than just policy outputs, GRL-SNAM maintains physical consistency while adapting to new sensory information, enabling robust performance across diverse environments.

**Superior performance:** GRL-SNAM consistently outperforms classical planning and reactive control methods across all metrics (success, efficiency, safety, smoothness) while requiring minimal computational overhead and sensing budget.

These results establish GRL-SNAM as the first method to successfully unify global navigation objectives with local safety constraints in hyperelastic navigation through principled geometric learning.

### 5 CONCLUSION

We introduced GRL-SNAM, a reinforcement learning framework that leverages Hamiltonian structure to couple sensing, planning, and deformation into a unified energy-based policy. Our formulation enables stable, feedforward navigation updates and achieves near-optimal path quality with minimal mapping effort in challenging deformable-robot tasks. The results highlight that incorporating geometric priors into RL can yield both efficiency and robustness, even under noisy sensing and out-of-distribution layouts. Future work will extend the approach to richer sensing modalities and more complex environments, with the goal of validating its scalability to real robotic systems.

# REFERENCES

- J. I. Alora, Moses C. Beard, Thomas Libby, Philipp Rothemund, et al. Discovering dominant dynamics for nonlinear continuum robot control. npj Robotics, 3(1):5, 2025.
- Karl Johan Åström and Björn Wittenmark. Adaptive Control. Courier Corporation, 2010.
- Chandrajit Bajaj and Minh Nguyen. *Physics-Informed Neural Networks via Stochastic Hamiltonian Dynamics Learning*, pp. 182–197. Springer Nature Switzerland, 2024. ISBN 9783031664281. doi: 10.1007/978-3-031-66428-1\_11. URL http://dx.doi.org/10.1007/978-3-031-66428-1\_11.
  - Atılım Günes Baydin, Barak A Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. Automatic differentiation in machine learning: a survey. *Journal of Machine Learning Research*, 18(153):1–43, 2018.
  - Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-1: Robotics transformer for real-world control at scale, 2023. URL https://arxiv.org/abs/2212.06817.
  - Brandon Caasenbrood, Alexander Pogromsky, and Henk Nijmeijer. Control-oriented models for hyperelastic soft robots through differential geometry of curves. *Soft Robotics*, 9(2):346–361, 2022.
  - Devendra Singh Chaplot, Ruslan Salakhutdinov, Abhinav Gupta, and Saurabh Gupta. Neural topological slam for visual navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12875–12884, 2020.
  - Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
  - Mohammad Dehghani Tezerjani et al. A survey on reinforcement learning applications in slam. *arXiv preprint arXiv:2408.14518*, 2024.
  - Shaan A Desai, Marios Mattheakis, David A Roberts, and Pavlos Protopapas. Port-hamiltonian neural networks for learning explicit time-dependent dynamical systems. *Physical Review E*, 104 (3):034312, 2021.
  - Aditya S Ellendula and Chandrajit Bajaj. Self-balancing, memory efficient, dynamic metric space data maintenance, for rapid multi-kernel estimation, 2025. URL https://arxiv.org/abs/2504.18003.
  - Zhan Feng et al. Safer gap: A gap-based local planner for safe navigation with nonholonomic mobile robots. *arXiv preprint arXiv:2303.08243*, 2023.
  - Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pp. 1587–1596, 2018.
  - Saurabh Gupta, Varun Tolani, James Davidson, Sergey Levine, Rahul Sukthankar, and Jitendra Malik. Cognitive mapping and planning for visual navigation, 2019. URL https://arxiv.org/abs/1702.03920.
  - Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870, 2018.

- Tai Hoang, Huy Le, Philipp Becker, Vien Anh Ngo, and Gerhard Neumann. Geometry-aware rl for manipulation of varying shapes and deformable objects. *arXiv preprint arXiv:2502.07005*, 2025.
- Junwoo Jang and Maani Ghaffari. Social zone as a barrier function for socially-compliant robot navigation, 2024. URL https://arxiv.org/abs/2405.15101.
  - Eshagh Kargar and Ville Kyrki. Macrpo: Multi-agent cooperative recurrent policy optimization, 2021. URL https://arxiv.org/abs/2109.00882.
  - Nikhil Keetha, Jay Karhade, Krishna Murthy Jatavallabhula, Gengshan Yang, Sebastian Scherer, et al. Splatam: Splat, track & map 3d gaussians for dense rgb-d slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
    - Patrick Kidger, James Morrill, James Foster, and Terry Lyons. Efficient and accurate gradients for neural sdes. *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
    - Junho Kim, Eun Sun Lee, Mingi Lee, Donsu Zhang, and Young Min Kim. Sgolam: Simultaneous goal localization and mapping for multi-object goal navigation, 2021. URL https://arxiv.org/abs/2110.07171.
    - Donald E Kirk. Optimal Control Theory: An Introduction. Dover Publications, 2004.
    - Holger Klein, Noémie Jaquier, Andre Meixner, and Tamim Asfour. On the design of region-avoiding metrics for collision-safe motion generation on riemannian manifolds, 2023. URL https://arxiv.org/abs/2307.15440.
    - Hengyuan Lai et al. Roboballet: Planning for multirobot reaching with graph neural networks and reinforcement learning. *Science Robotics*, 2025.
    - Kyowoon Lee, Seongun Kim, and Jaesik Choi. Adaptive and explainable deployment of navigation skills via hierarchical deep reinforcement learning. In *International Conference on Robotics and Automation*, 2023.
    - Chengshu Li et al. Hrl4in: Hierarchical reinforcement learning for interactive navigation with mobile manipulators. In *Conference on Robot Learning*, 2020.
    - Junjie Li et al. Learn with imagination: Safe set guided state-wise constrained policy optimization. *arXiv preprint arXiv:2308.13140*, 2023.
    - Elisabetta Liu and Cosimo Della Santina. Physics-informed neural networks to model and control robots: a theoretical and experimental investigation. *Advanced Intelligent Systems*, 2024.
    - David Martínez-Rubio and Sebastian Pokutta. Accelerated riemannian optimization: Handling constraints with a prox to bound geometric penalties. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 359–393. PMLR, 2023.
    - Nicholas Mohammad and Nicola Bezzo. Soft actor-critic-based control barrier adaptation for robust autonomous navigation in unknown environments, 2025. URL https://arxiv.org/abs/2503.08479.
    - Barry W Mulvey and Thrishantha Nanayakkara. Haven: haptic and visual environment navigation by a shape-changing mobile robot with multimodal perception. *Scientific Reports*, 14(1):27018, 2024.
    - Tariq Patanam, Eli Shayer, and Younes Bensouda Mourri. Deep dagger imitation learning for indoor scene navigation.
- L.S. Pontryagin, V.G. Boltyanskii, R.V. Gamkrelidze, and E.F. Mishchenko. *The Mathematical Theory of Optimal Processes*. Interscience, 1962.
  - Doina Precup. *Temporal abstraction in reinforcement learning*. University of Massachusetts Amherst, 2000.
    - F. Qi, C. Zhou, H. Qing, H. Sun, and J. Yin. Aerial track-guided autonomous soft ring robot. *Advanced Science*, 2024.

- Mohammad Roshanfar, Javad Dargahi, and Amir Hooshiar. Hyperelastic modeling and validation of hybrid-actuated soft robot with pressure-stiffening. *Micromachines*, 14(5):1001, 2023.
  - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. 2017.
  - Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations (ICLR)*, 2021.
  - Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew J Davison. imap: Implicit mapping and positioning in real-time. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6229–6238, 2021.
  - Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: Learning, planning, and representing knowledge at multiple temporal scales. Technical report, Technical Report 98-74, University of Massachusetts, Amherst, 1998.
  - Hamid Taheri et al. Deep reinforcement learning with enhanced ppo for safe mobile robot navigation. *arXiv preprint arXiv:2405.16266*, 2024.
  - Lei Tai, Jingwei Zhang, Ming Liu, and Wolfram Burgard. Socially compliant navigation through raw depth inputs with generative adversarial imitation learning, 2018. URL https://arxiv.org/abs/1710.02543.
  - Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. In *International conference on machine learning*, pp. 3540–3549. PMLR, 2017.
  - Niclas Vödisch, Daniele Cattaneo, Wolfram Burgard, and Abhinav Valada. *Continual SLAM: Beyond Lifelong Simultaneous Localization and Mapping Through Continual Learning*, pp. 19–35. Springer Nature Switzerland, 2023. ISBN 9783031255557. doi: 10.1007/978-3-031-25555-7\_3. URL http://dx.doi.org/10.1007/978-3-031-25555-7\_3.
  - Tianhao Wang et al. Pinn-ray: A physics-informed neural network to model soft robotic fin ray fingers. arXiv preprint arXiv:2407.08222, 2024.
  - Weizheng Wang et al. Multi-agent llm actor-critic framework for social robot navigation. *arXiv* preprint arXiv:2503.09758, 2025.
  - Chi Yan, Delin Qu, Dan Wang, Dan Xu, Zhigang Wang, et al. Gs-slam: Dense visual slam with 3d gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
  - Songyuan Zhang, Zhangjie Cao, Dorsa Sadigh, and Yanan Sui. Confidence-aware imitation learning from demonstrations with varying optimality, 2022. URL https://arxiv.org/abs/2110.14754.
  - Chao Zheng et al. Semantic slam system for mobile robots based on large visual model in complex environments. *Scientific Reports*, 15(1):1–15, 2025.
  - Jun Zhu, Zihao Du, Haotian Xu, Fengbo Lan, Zilong Zheng, Bo Ma, Shengjie Wang, and Tao Zhang. Navi2gaze: Leveraging foundation models for navigation and target gazing, 2024. URL https://arxiv.org/abs/2407.09053.
  - Zihan Zhu, Songyou Peng, Viktor Laehner, Weiyang Xu, Michael Niemeyer, et al. Nice-slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12786–12796, 2022.

# A EXTENDED INTRODUCTION AND MOTIVATION

This section provides expanded context for the challenges addressed by GRL-SNAM and detailed justification for our geometric approach.

### A.1 COMPREHENSIVE ANALYSIS OF RL LIMITATIONS IN NAVIGATION

Contemporary reinforcement learning methods face several critical limitations that become particularly pronounced in continuous navigation tasks:

Sample Efficiency Bottlenecks. Standard RL algorithms like SAC Haarnoja et al. (2018) and PPO Schulman et al. (2017) require millions of environment interactions to learn effective navigation policies. This inefficiency stems from the curse of dimensionality in continuous control settings where the action space is infinite-dimensional and policies must simultaneously master fine-grained motor control and high-level strategic reasoning. In real-world deployment scenarios where data collection is expensive and potentially dangerous, this sample complexity becomes prohibitive.

The problem is exacerbated by the need for exploration in high-dimensional spaces. Unlike discrete control problems where systematic exploration strategies like  $\epsilon$ -greedy or UCB can provide theoretical guarantees, continuous control requires sophisticated exploration mechanisms that often rely on injected noise or entropy bonuses. These mechanisms frequently lead to unsafe or inefficient exploration behaviors that are unsuitable for real-world navigation tasks.

**Generalization Failures.** Policies trained in specific environments exhibit catastrophic performance degradation when deployed in novel settings, even when new environments share similar structure. This brittleness stems from the lack of inductive bias in standard neural network architectures. Without explicit encoding of physical principles or geometric structure, learned policies tend to memorize environment-specific features rather than discovering generalizable navigation principles.

The generalization problem is particularly acute in navigation because environmental variations can affect multiple aspects of the task simultaneously: obstacle configurations change collision constraints, surface properties affect dynamics, and lighting conditions influence perception. Standard RL approaches learn monolithic mappings that cannot decompose these variations into their constituent factors, leading to brittle behaviors that fail when any component deviates from training conditions.

**Temporal Decomposition Challenges.** Navigation inherently requires coordination across multiple timescales: immediate obstacle avoidance operates on millisecond timescales, local path planning unfolds over seconds, and strategic goal-directed behavior spans minutes or hours. Standard RL algorithms struggle to learn policies that reason effectively across these scales, often getting trapped in locally optimal behaviors that satisfy short-term objectives while failing to make long-term progress.

Existing approaches to multi-scale reasoning such as hierarchical RL Sutton et al. (1998), options frameworks Precup (2000), or feudal networks Vezhnevets et al. (2017), typically require manual decomposition of the task space and careful engineering of reward functions for different levels. These methods introduce additional complexity without fundamentally addressing the structural issues that make multi-scale learning difficult.

### A.2 THE SNAM CHALLENGE: WHY STRUCTURE MATTERS

Simultaneous Navigation and Mapping (SNAM) represents a particularly challenging instance of the navigation problem where agents must build environmental representations online while traversing unknown spaces. This challenge amplifies the limitations of conventional RL approaches in several ways:

**Memory and Representation Learning.** SNAM requires policies to maintain and update spatial representations based on sensory observations. This places enormous demands on the policy's memory architecture, requiring it to simultaneously master memory management, spatial reasoning, and motor control. Standard recurrent architectures like LSTMs or GRUs struggle with this multifaceted learning problem, often failing to maintain coherent spatial representations over long episodes.

**Exploration-Exploitation Tradeoffs.** In SNAM, exploration serves dual purposes: gathering information about the environment for mapping and discovering navigation strategies. This creates complex exploration-exploitation tradeoffs that standard RL exploration mechanisms cannot handle effectively. Random exploration may discover new regions but fails to systematically map environmental structure, while directed exploration based on current maps may miss critical environmental features.

**Dynamic Environmental Coupling.** Unlike traditional navigation where environments are static, SNAM requires reasoning about how the agent's actions affect both its position and its knowledge of the environment. This creates a coupled learning problem where navigation decisions influence future mapping accuracy, and mapping quality affects navigation performance. Standard RL frameworks treat these as separate problems, missing the critical coupling that enables efficient SNAM.

Recent approaches in simultaneous navigation and mapping (SNAM) have coupled local mapping with policy learning to improve navigation performance. For example, SGoLAM Kim et al. (2021) interleaves goal localization with occupancy mapping to enable point-goal navigation, while Cognitive Mapping and Planning (CMP) Gupta et al. (2019) integrates a differentiable planner into a learned mapping framework. Continual SLAM (CL-SLAM) Vödisch et al. (2023) further emphasizes long-term adaptability by maintaining and updating maps during navigation. However, these methods rely on progressively constructing detailed maps of the environment before exploiting them for navigation. In contrast, our objective is to reach the goal along high-quality, well-weighted paths while mapping as little of the unknown environment as possible. To the best of our knowledge, no prior work explicitly formulates navigation with minimal exploration as the central goal. Our proposed GRL-SNAM framework achieves this by progressively refining paths: from observed environmental variations, the policy differentially learns to identify the least-cost trajectory, such that the path improves continuously as new local information is revealed.

### A.3 GEOMETRIC STRUCTURE: THE INEVITABLE SOLUTION

The limitations outlined above are not merely implementation details but fundamental consequences of treating navigation as unstructured optimization. Several lines of evidence suggest that geometric structure is not just helpful but inevitable for solving complex navigation problems:

**Physical Realizability.** Real robotic systems operate under physical constraints imposed by conservation laws, kinematic limitations, and actuator dynamics. Policies that violate these constraints cannot be implemented on physical systems, yet standard RL approaches have no mechanism to enforce such constraints during learning. Geometric formulations naturally incorporate physical constraints through the mathematical structure of the problem.

**Stability Requirements.** Long-horizon navigation requires numerical stability over extended rollouts. Standard neural network policies accumulate errors over time, leading to unstable behaviors in long episodes. Hamiltonian formulations with symplectic structure preserve important invariants (energy, momentum) that ensure stability over arbitrarily long rollouts.

Compositionality Needs. Complex navigation tasks require composing simpler behaviors: obstacle avoidance, path following, goal seeking, and environmental adaptation. Standard RL approaches learn monolithic policies that cannot decompose into interpretable components. Geometric formulations enable natural decomposition through energy terms that can be composed, weighted, and adapted independently.

### A.4 DIFFERENTIAL POLICY OPTIMIZATION: BEYOND FIXED POLICIES

Traditional RL optimizes fixed policy parameters  $\theta$  to maximize expected returns over discrete timesteps. Our Differential Policy Optimization (DPO) approach fundamentally reconceptualizes this by learning dynamics operators through a continuous-time differential dual formulation.

**Mathematical Foundation.** Rather than directly learning policies, DPO reformulates RL through continuous-time optimal control. By approximating discrete reward sums with time integrals:

$$\max_{\pi} \mathbb{E}\left[\sum_{k=0}^{H-1} r(s_k, a_k)\right] \approx \max_{\pi} \mathbb{E}\left[\int_0^T r(s_t, a_t) dt\right]$$
 (12)

Applying Pontryagin's Maximum Principle introduces adjoint variables p and defines the Hamiltonian function:

$$HF(p, s, a) := p^{T} f(s, a) - r(s, a)$$
 (13)

The key insight is that optimal actions can be implicitly represented through the stationarity condition  $\frac{\partial HF}{\partial a} = 0$ , yielding the reduced Hamiltonian:

$$hf(s,p) := HF(s,p,a^*(s,p))$$
 (14)

Score Function Learning. DPO learns a score function  $g(x) \approx hf(x)$  where x = (s, p) combines state and adjoint variables. The dynamics operator is constructed as:

$$G(x) = x + \Delta S \nabla g(x) \tag{15}$$

where  $S=\begin{bmatrix}0&I\\-I&0\end{bmatrix}$  is the canonical symplectic matrix and  $\Delta$  is the discretization step.

**Stagewise Learning Advantages.** Unlike methods requiring backward-in-time adjoint calculations (as in Pontryagin's Maximum Principle), DPO enables feedforward learning where each stage t defines a local Hamiltonian  $\mathcal{H}_t$  integrated forward in time:

$$\theta_{t+1} = \theta_t - \eta \nabla_\theta \mathcal{H}_t \tag{16}$$

This avoids the computational complexity and numerical instability of adjoint methods while maintaining theoretical guarantees through the geometric structure of the Hamiltonian formulation.

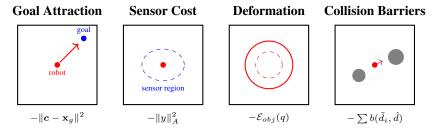


Figure 6: Policy-aligned energy decomposition. Each policy governs a distinct energy component: the Sensor Policy minimizes sensor cost, the FPE balances goal attraction and collision barriers, and the Reconfig Policy adapts size through deformation energy. Together these terms define the Hamiltonian reward  $\mathcal{R}$ .

# A.5 MULTI-POLICY ARCHITECTURE DETAILS

Our multi-policy decomposition addresses temporal scale separation through three specialized components operating at different timescales:

Sensor Policy  $(\pi_y)$ : Operates at slow timescales to adapt perception strategies based on stagewise environmental feedback. This policy learns to focus attention on relevant environmental features, adjust sensor parameters for optimal information gain, and filter sensory noise. The sensor policy outputs constraints  $C_t$  that inform slower planning processes.

**Frame Policy**  $(\pi_f)$ : Operates at medium timescales to plan collision-free trajectories in local coordinate frames. This policy takes constraints from the sensor policy and generates waypoints  $\mathcal{W}_t$  for shape control. The frame policy handles local obstacle avoidance and path optimization within a limited spatial horizon.

**Shape Policy**  $(\pi_o)$ : Operates at fast timescales to control robot morphological adaptation. For deformable robots, this includes shape changes, stiffness modulation, and configuration updates. For conventional robots, this might include gait transitions, tool selection, or behavioral mode switches.

The key insight is that these policies are not manually designed hierarchies but emerge naturally from the temporal structure of the Hamiltonian dynamics. Fast variables (sensor adaptation) reach

quasi-equilibrium before slower variables (shape changes) evolve significantly, creating natural scale separation without manual decomposition.

This extended analysis demonstrates that geometric structure is not merely a useful inductive bias but a necessary foundation for solving complex navigation problems that require multi-scale reasoning, online adaptation, and long-horizon stability.

# B EXTENDED RELATED WORK SURVEY

### B.1 GEOMETRY AND MECHANICS PRIMER

Navigation learning methods can be categorized by their underlying mathematical spaces, with significant implications for performance and theoretical guarantees:

**Euclidean Space Methods** ( $\mathbb{R}^n$ ): Standard RL treats navigation as optimization in flat spaces using Euclidean distance metrics. Enhanced PPO Taheri et al. (2024) demonstrate improved collision avoidance but ignore inherent geometric structure of robotic systems. Sample efficiency remains poor, typically requiring millions of environment interactions Dehghani Tezerjani et al. (2024).

**Lie Group Methods:** Recognition of orientation constraints has led to SE(2) and SE(3) formulations using equivariant neural architectures. These preserve rotational and translational symmetries but remain primarily limited to manipulation rather than navigation tasks.

**Riemannian Manifold Approaches:** Advanced geometric formulations employ differential geometry for constraint handling through tangent space projections. Martínez-Rubio & Pokutta (2023) demonstrates constraint satisfaction through geometric structure rather than penalty methods, achieving superior theoretical properties but limited practical deployment.

**Hamiltonian and Symplectic Methods:** Port-Hamiltonian neural networks show significant performance improvements through symplectic integrators, proving that respecting geometric structure fundamentally improves learning dynamics. However, applications remain confined to simple control problems.

### B.2 SAFETY-CRITICAL NAVIGATION TAXONOMY

External Safety Projection: Control Barrier Functions create safe action spaces through constraint projection. Neural Network Zeroing Barrier Functions Feng et al. (2023) enable collision-free navigation, while adaptive safety constraints Mohammad & Bezzo (2025) handle dynamic environments. Social navigation approaches Jang & Ghaffari (2024) extend CBFs to human-robot interaction. These methods achieve formal safety guarantees but often exhibit conservative behaviors due to the separation between safety and optimality.

**Energy-Integrated Safety:** Our approach incorporates safety directly within the Hamiltonian energy structure via barrier potentials. This enables aggressive navigation while maintaining formal guarantees through symplectic structure preservation, avoiding the conservatism of external projection methods.

### B.3 Deformable and Soft Robot Navigation

**Hyperelastic Material Models:** Recent advances include pressure-stiffening control with 6.40% maximum error validation Roshanfar et al. (2023) and passivity-based control using differential geometry of curves Caasenbrood et al. (2022). Spectral Submanifold Reduction Alora et al. (2025) achieves computational speedup for real-time hyperelastic control with stability guarantees.

**Ring and Circular Robots:** Liquid Crystal Elastomer responses enable aerial gap navigation Qi et al. (2024) through predetermined actuation patterns. HAVEN Mulvey & Nanayakkara (2024) navigates constrained spaces via fixed shape-changing sequences based on multimodal perception. These approaches use offline parameter optimization with deterministic execution, lacking online adaptation capabilities.

**Physics-Informed Learning:** PINN-Ray Wang et al. (2024) achieves state-of-the-art hyperelastic displacement prediction, while extensions to non-conservative effects Liu & Della Santina (2024)

provide experimental validation. However, these remain primarily modeling tools rather than adaptive control frameworks.

### B.4 NEURAL SCENE REPRESENTATIONS FOR NAVIGATION

**NeRF-Based SLAM:** Real-time dense reconstruction through NICE-SLAM Zhu et al. (2022) and keyframe-free tracking via iMAP Sucar et al. (2021) provide rich environmental representations. Neural Topological SLAM Chaplot et al. (2020) combines learning with classical planning, while semantic approaches Zheng et al. (2025) integrate large vision models.

**3D Gaussian Splatting:** GS-SLAM Yan et al. (2024) and SplaTAM Keetha et al. (2024) demonstrate state-of-the-art reconstruction quality with real-time performance, offering dense 3D representations suitable for navigation applications.

Integration with Energy Terms: Scene representations feed our energy formulation through:

Barrier Energy: 
$$\mathcal{U}_{\text{barrier}} = \sum_{\text{obstacles}} b(\text{SDF}(\mathbf{x}))$$
 (17)

Free-Space Energy: 
$$U_{\text{free}} = -\sum_{\text{free regions}} w(\mathbf{x})$$
 (18)

Goal Energy: 
$$U_{\text{goal}} = \|\mathbf{x} - \mathbf{x}_{\text{goal}}\|^2$$
 (19)

### B.5 MULTI-SCALE AND HIERARCHICAL METHODS

**Hierarchical RL:** Task decomposition approaches like HRL4IN Li et al. (2020) handle heterogeneous navigation phases, while Lee et al. (2023) learns specialized policy families with high-level coordination. These require manual decomposition and struggle with principled coordination, often leading to ad-hoc design choices without theoretical guarantees.

**Multi-Agent Coordination:** RoboBallet Lai et al. (2025) achieves coordination for 8 robots across 40 tasks using graph neural networks. MACRPO Kargar & Kyrki (2021) enhances information sharing beyond parameter sharing. However, these approaches lack the geometric structure preservation critical for deformable robot coordination.

# B.6 IMITATION LEARNING FOR NAVIGATION

**Behavioral Cloning:** RT-1 Brohan et al. (2023) demonstrates impressive generalization across 700+ tasks using 130k demonstration episodes with transformer architectures achieving significant zero-shot performance improvements.

**Inverse Reinforcement Learning:** GAIL for Safe Navigation Tai et al. (2018) combines generative adversarial imitation with safety constraints. DAgger for Continuous Navigation Patanam et al. iteratively improves policies through expert querying.

**Sub-Optimal Demonstrations:** Confident Imitation Learning Zhang et al. (2022) handles demonstration uncertainty through confidence-aware training, addressing distribution shift in novel environments.

These approaches excel with high-quality demonstrations but assume expert availability and struggle with the full behavioral range needed for adaptive deformation strategies.

### B.7 FOUNDATION MODEL INTEGRATION

Large-scale models for navigation reasoning Zhu et al. (2024); Wang et al. (2025) focus on high-level semantic understanding and multi-agent coordination at the symbolic level. Foundation models excel at reasoning and semantic understanding, while our GRL-SNAM provides principled low-level geometric control.

**Integration Pathway:** Foundation models could generate high-level objectives encoded as potential energy terms in our energy functional  $\mathcal{R}(q_t)$ . The geometric structure preservation ensures high-

level semantic goals translate into physically consistent behaviors, addressing the critical gap where foundation model outputs often lack grounding in physical dynamics.

### **B.8 PARADIGM COMPARISON**

Table 2: Extended paradigm-level comparison of learning frameworks. Scoring:  $\checkmark$  = comprehensive support,  $\triangle$ = limited support,  $\times$ = not supported.

Capability	GRL-SNAM	Standard RL	Geometric RL	Imitation Learning	Semi/Unsupervised	CBF Methods	Hierarchical RL	Foundation Models
Energy Conservation	<b> </b>	l ×	Δ	×	×	×	×	×
Geometric Structure	✓	×	✓	×	Δ	Δ	×	×
Constraint Integration	✓	Δ	✓	×	×	✓	Δ	Δ
Online Adaptation	✓		Δ	×	✓	Δ	✓	✓
Multi-Scale Coordination	✓	×	×	×	×	×	✓	Δ
Sample Efficiency	✓	×	Δ	✓	Δ	Δ	Δ	✓
Zero-Shot Generalization	✓	×	Δ	×	✓	×	×	✓
Real-World Deployment	✓	✓	Δ	✓	Δ	✓	✓	Δ
Deformable Robot Support	<b>√</b>	×	×	×	×	×	×	×

### **Scoring Criteria:**

- Energy Conservation: Explicit conservation laws in dynamics
- Geometric Structure: Preservation of manifold properties
- Constraint Integration: Safety/task constraints within optimization
- Online Adaptation: Real-time policy modification during deployment
- Multi-Scale Coordination: Principled coordination across temporal scales
- Sample Efficiency: Learning with minimal environment interaction
- Zero-Shot Generalization: Performance in unseen environments
- Real-World Deployment: Practical implementation feasibility
- Deformable Robot Support: Explicit modeling of shape change

This comprehensive survey positions GRL-SNAM as uniquely addressing the intersection of geometric structure preservation, multi-scale coordination, and deformable robot control—capabilities that existing approaches handle separately or incompletely.

### **B.9** KEY INSIGHTS

Our framework builds upon a set of Hamiltonian and reinforcement learning principles, unifying offline reference dynamics with online adaptive updates. Below, we summarize the six key insights that form the backbone of GRL-SNAM.

### 1. Hamiltonian energy as task reward. We define the Hamiltonian

$$\mathcal{H}(q,p) = K(p) + P(q), \tag{20}$$

with kinetic energy K and task-specific potential P. In our setup, P encodes navigation objectives (goal attraction, barrier avoidance, deformation penalties). Following Pontryagin et al. (1962); Bajaj & Nguyen (2024), the Hamiltonian coincides with the surrogate objective in policy gradient methods, i.e.

$$\nabla_{\theta} J(\pi_{\theta}) \approx \nabla_{\theta} \mathbb{E}_{\pi_{\theta}} [-\mathcal{H}(q, p)],$$
 (21)

linking task reward to the Hamiltonian gradient flow. This equivalence grounds the DPO surrogate in a physical structure.

**2.** Offline Hamiltonian vs. Online task reward. In offline training, the agent minimizes trajectories under a fixed  $\mathcal{H}$  constructed from synthetic local patches. Online, the environment is sensed, and task rewards  $\mathcal{R}_{env}$  are parsed into Hamiltonian subtasks. By interpreting

$$\mathcal{H}_{\text{online}} = \mathcal{H}_{\text{offline}} + \Delta \mathcal{R}_{\text{env}}, \tag{22}$$

we align local sensory updates with the reference offline Hamiltonian. This mirrors the adaptive control interpretation in Åström & Wittenmark (2010).

3. Offline policy as reference Hamiltonian. Every offline policy  $\pi_{\text{ref}}$  is equivalent to a reference Hamiltonian  $\mathcal{H}_{\text{ref}}$ , where the score function  $s^{\theta} = \nabla \mathcal{H}_{\text{ref}}$  defines canonical dynamics:

$$\dot{q} = \frac{\partial \mathcal{H}_{\text{ref}}}{\partial p}, \qquad \dot{p} = -\frac{\partial \mathcal{H}_{\text{ref}}}{\partial q}.$$
 (23)

Online adaptation then minimizes the divergence

$$D(\pi_{\text{online}} \parallel \pi_{\text{ref}}) \propto \mathbb{E}[\lVert \nabla \mathcal{H}_{\text{online}} - \nabla \mathcal{H}_{\text{ref}} \rVert^2],$$
 (24)

a structure exploited in score-based models (Song et al., 2021).

**4.** Advantages of stagewise updates. Rather than solving adjoint equations as in Pontryagin's Maximum Principle, we adopt a stagewise decomposition. Each stage defines a local  $\mathcal{H}_t$  and is integrated feedforward:

$$\theta_{t+1} = \theta_t - \eta \nabla_\theta \mathcal{H}_t. \tag{25}$$

This avoids backward-in-time adjoint calculations and recovers the efficiency noted in adjoint-free feedforward networks (Chen et al., 2018; Kidger et al., 2021).

5. Universality of the pipeline. Our pipeline

Environment 
$$\xrightarrow{\text{Encoder}}$$
 Context  $\xrightarrow{\text{Setup}}$   $\mathcal{H}_{\text{adapted}}$  (26)

is universal. As long as  $\mathcal{H}$  is differentiable, adaptation reduces to evaluating its gradients, regardless of whether the system is white-box (explicit potentials) or black-box (sensor-level inputs). This follows from the variational formulation of differentiable programming (Baydin et al., 2018).

**6. Navigator as meta-controller.** The navigator policy  $\pi_{nav}$  interacts with three black boxes: the offline Hamiltonian  $\mathcal{H}_{ref}$ , the online sensed reward  $\mathcal{R}_{env}$ , and the adaptive fusion  $\mathcal{H}_{adapt}$ . Its role is to formulate and solve

$$\mathcal{H}_{\text{adapt}} = \alpha \mathcal{H}_{\text{ref}} + (1 - \alpha) \mathcal{R}_{\text{env}}, \tag{27}$$

where  $\alpha$  is dynamically updated by the context encoder (e.g., LSTM). This positions the navigator as a meta-controller that continually reforms the Hamiltonian problem, a principle consistent with adaptive RL formulations in Kirk (2004).

Secant controller as on-the-fly energy reshaping. Our history-based controller instantiates equation 5 with (i) a rank-1 secant estimate  $\hat{\mathbf{J}}$  (EMA-smoothed) from consecutive frames, and (ii) the Gauss-Newton metric  $\hat{\mathbf{G}}(\Theta) = \hat{\mathbf{J}}^{\top} \mathbf{W} \hat{\mathbf{J}} + \varepsilon I$ , then projects to the nonnegative orthant and applies per-head step sizes:

$$\Delta\Theta = -\widehat{\mathbf{G}}^{-1}\widehat{\mathbf{J}}^{\top}\mathbf{W}\left(f(\Theta) - y_{\text{tgt}}\right), \qquad \Theta^{t+1} = \text{Proj}_{\Theta \geq 0}\left(\Theta^{t} + \text{diag}(\eta_{b}, \eta_{g}, \eta_{\alpha}) \Delta\Theta\right).$$

Concretely, with  $\Theta = [\beta, \gamma, \{\alpha_i\}_{i \in \mathcal{N}_K(o)}]^{\top}$ , this reshapes the parameterized energy  $\mathcal{E}_{\Theta} = \beta \Phi_{\text{safe}} + \gamma \Phi_{\text{speed}} + \sum_{i \in \mathcal{N}_K(o)} \alpha_i \Phi_i^{\text{obst}}$  so that its induced observables  $f(\Theta)$  move toward  $y_{\text{tgt}}$  (safety, progress, admissible speed) without extra rollouts.

**Sequential Query–Response (as ports).** At each t the Navigator issues queries  $\mathcal{Q}_k^t$  and receives responses  $\mathcal{R}_k^t$ , which determine  $y_{\text{tgt}}(\mathcal{R}^t)$  and any weights in  $\mathbf{W}$ ; the update equation 5 (with the secant  $\hat{\mathbf{J}}$ ) is then applied to each block  $k \in \{y, f, o\}$ :

$$\theta_k^{t+1} = \theta_k^t + h \,\Delta \theta_k^t, \quad \Delta \theta_k^t = -\widehat{\mathbf{G}}_k^{-1} \widehat{\mathbf{J}}_k^{\top} \mathbf{W}_k \big( f_k(\theta_k^t) - y_{\mathrm{tgt},k}(\mathcal{R}_k^t) \big).$$

**State Evolution.** With updated parameters, each policy advances its state as before,

$$z_k^{t+1} = z_k^t + \tau_k J_k s_k^{\theta_k^{t+1}}(z_k^t, \mathcal{C}_t^{\text{updated}}, t).$$

*Remark.* If a strict port-Hamiltonian view is desired, the  $\mathbf{J}^{\top}\mathbf{W}(y_{\text{tgt}} - f(\Theta))$  term enters  $\dot{\Pi}$  as an external port (input) rather than being baked into the potential; the resulting discrete update is identical to the secant Gauss–Newton step above.

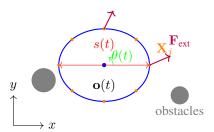


Figure 7: Hyperelastic ring robot model showing generalized coordinates  $(s, \mathbf{o}, \theta)$ , spline sample points  $\mathbf{X}_j$ , and external force fields.

# C HYPERELASTIC RING ROBOT MODEL

We model the deformable robot as a closed hyperelastic ring with reduced-order dynamics to enable efficient navigation while capturing essential deformation behaviors.

### C.1 GEOMETRIC REPRESENTATION

The robot boundary is defined by a periodic cubic B-spline curve with  $n_{\mathrm{ctrl}}$  control points:

$$\mathbf{S}(u) = \sum_{i=1}^{n_{\text{ctrl}}} N_{i,3}(u) \mathbf{P}_i, \quad u \in [0, 1]$$
(28)

where  $N_{i,3}(u)$  are degree-3 B-spline basis functions with  $C^2$  continuity. The base shape is a unit circle:

$$\mathbf{P}_{0,i} = r_{\text{base}} \begin{bmatrix} \cos(2\pi i/n_{\text{ctrl}}) \\ \sin(2\pi i/n_{\text{ctrl}}) \end{bmatrix}$$
 (29)

World coordinates are computed via similarity transformation:

$$\mathbf{P}_{i}(t) = \mathbf{o}(t) + s(t)\mathbf{R}(\theta(t))\mathbf{P}_{0,i}$$
(30)

where s(t) is uniform scale,  $\mathbf{o}(t) \in \mathbb{R}^2$  is center position,  $\theta(t)$  is orientation.

For physics computation, we sample K points on the curve using B-spline evaluation matrix  $B \in \mathbb{R}^{K \times n_{\text{ctrl}}}$ :

$$\mathbf{X}_{j} = \sum_{i=1}^{n_{\text{ctrl}}} B_{ji} \mathbf{P}_{i}, \quad j = 1, \dots, K$$

$$(31)$$

### C.2 ENERGY FORMULATION

The total Hamiltonian combines kinetic and potential components:

$$\mathcal{H} = \frac{1}{2} M_s \dot{s}^2 + \frac{1}{2} M_o ||\dot{\mathbf{o}}||^2 + \frac{1}{2} I \omega^2 + \mathcal{U}_{\text{barrier}} + \mathcal{U}_{\text{bulk}}$$
(32)

IPC Barrier Energy: Collision avoidance using Incremental Potential Contact barriers:

$$\mathcal{U}_{\text{barrier}} = \sum_{j=1}^{K} w_j \ell_j \sum_{k=1}^{N_{\text{obs}}} b_{\text{IPC}}(d_{jk})$$
(33)

where  $w_j = 1/K$ ,  $\ell_j = ||\mathbf{X}'_j||$ ,  $d_{jk}$  is distance from sample j to obstacle k:

$$b_{\text{IPC}}(d) = \begin{cases} -(d - \hat{d})^2 (\log d - \log \hat{d}) & \text{if } 0 < d < \hat{d} \\ 0 & \text{if } d \ge \hat{d} \\ V_{\text{penalty}} & \text{if } d \le 0 \end{cases}$$
(34)

Adaptive Bulk Energy: Area conservation with clearance-dependent target:

$$\mathcal{U}_{\text{bulk}} = \frac{k_{\text{bulk}}}{2} (A(s) - A_{\text{target}})^2 \tag{35}$$

where  $A(s) = s^2 A_{\text{ref}}$  and:

$$A_{\text{target}} = \left[\alpha + (1 - \alpha) \tanh(\beta \cdot \max(d_{\min}, 0))\right] A_{\text{ref}}$$
(36)

with  $\alpha = 0.25$ ,  $\beta = 2.5$  encouraging compression in tight spaces.

# C.3 GENERALIZED FORCE MAPPING

Forces on spline samples map to generalized coordinates via virtual work:

$$F_{s} = -\frac{\partial \mathcal{U}}{\partial s} + \sum_{j=1}^{K} \mathbf{F}_{\text{ext}}(\mathbf{X}_{j}) \cdot \frac{\partial \mathbf{X}_{j}}{\partial s} - \gamma_{s} \dot{s}$$
(37)

$$\mathbf{F}_{o} = -\frac{\partial \mathcal{U}}{\partial \mathbf{o}} + \sum_{j=1}^{K} w_{j} \mathbf{F}_{\text{ext}}(\mathbf{X}_{j}) - \gamma_{o} \dot{\mathbf{o}}$$
(38)

$$\tau = -\frac{\partial \mathcal{U}}{\partial \theta} + \sum_{j=1}^{K} w_j \mathbf{F}_{\text{ext}}(\mathbf{X}_j) \cdot (\mathbf{J}(\mathbf{X}_j - \mathbf{o})) - \gamma_{\theta} \omega$$
 (39)

where  $\mathbf{J} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$  generates rotation and  $\frac{\partial \mathbf{X}_j}{\partial s} = \mathbf{R}(\theta) \mathbf{P}_{0,j}$ .

# D DOMAIN-SPECIFIC POLICY IMPLEMENTATIONS

# D.1 Sensor Policy $(\pi_y)$ Details

The sensor policy maintains spatial index  $\mathcal{T}_y$  of observations  $(\mathbf{x}_i, \mathsf{type}_i, \mathsf{attr}_i)$  and derives three energy components from single neighbor queries:

**Barrier Potential:** Repulsion from obstacles

$$b_{\Sigma}(z_y, \mathcal{C}_t) = \sum_{i \in \mathcal{N}_{\text{obs}}} w_i \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{c}_y\|^2}{2\sigma_b^2}\right)$$
(40)

**Free-Space Potential:** Attraction to open regions

$$\mathcal{V}_{\text{free}}(z_y, \mathcal{C}_t) = -\sum_{j \in \mathcal{N}_{\text{free}}} w_j \exp\left(-\frac{\|\mathbf{x}_j - \mathbf{c}_y\|^2}{2\sigma_f^2}\right)$$
(41)

Density Potential: Information-theoretic density measure

$$\rho(z_y, \mathcal{C}_t) = -\sum_{k \in \mathcal{N}_{\text{all}}} w_k \log \left( 1 + \frac{n_k}{|\mathcal{N}_{\text{all}}|} \right)$$
(42)

The complete sensor score function is:

$$s_y^{\theta_y}(z_y, \mathcal{C}_t, t) = \nabla_{z_y} \left[ \frac{1}{2} \| p_y \|_{M_y^{-1}}^2 + \alpha_b b_{\Sigma} + \alpha_f \mathcal{V}_{\text{free}} + \alpha_d \rho \right]$$

$$\tag{43}$$

### D.2 FRAME POLICY $(\pi_f)$ DETAILS

The frame policy uses  $\mathcal{T}_f$  storing path samples with safety/contact distances and goal influence:

Safety Field: Distance-based safety measure

$$S(z_f, C_t) = \sum_{i=1}^{N_s} w_i \max(0, d_{\text{safe}}^{\text{threshold}} - d_{\text{safe}}^i)^2$$
(44)

Contact Field: Proximity to obstacles

$$C(z_f, \mathcal{C}_t) = \sum_{i=1}^{N_c} w_i \exp\left(-\frac{(d_{\text{contact}}^i)^2}{2\sigma_c^2}\right)$$
(45)

Goal Field: Directional bias toward target

$$G(z_f, C_t) = -\|c_f - \mathbf{x}_g\|^2 + \sum_{i=1}^{N_g} w_i g_i \cos(\theta_i)$$
(46)

The frame score function integrates these fields:

$$s_f^{\theta_f}(z_f, \mathcal{C}_t, t) = \nabla_{z_f} \left[ \frac{1}{2} \| p_f \|_{M_f^{-1}}^2 + \alpha_s S + \alpha_c C + \alpha_g G \right]$$
 (47)

# D.3 SHAPE POLICY $(\pi_o)$ DETAILS

The shape policy controls deformation through reduced coordinates  $z_o = (s, \dot{s}, \mathbf{o}, \dot{\mathbf{o}}, \theta, \omega)$ :

$$\mathcal{E}_{\text{smooth}} = \int_0^1 \|\kappa(u)\|^2 du \approx \sum_{j=1}^K w_j \|\kappa_j\|^2$$
(48)

Stretching Energy: Arc length preservation

$$\mathcal{E}_{\text{stretch}} = \int_{0}^{1} (\|\mathbf{S}'(u)\| - \ell_{\text{ref}})^{2} du \approx \sum_{j=1}^{K} w_{j} (\ell_{j} - \ell_{\text{ref}})^{2}$$
(49)

**Target Energy:** Configuration constraints

$$\mathcal{E}_{\text{target}} = \|\mathbf{P} - \mathbf{P}_{\text{target}}\|_F^2 + \|(s, \mathbf{o}, \theta) - (s_{\text{target}}, \mathbf{o}_{\text{target}}, \theta_{\text{target}})\|^2$$
 (50)

The complete shape score function is:

$$s_o^{\theta_o}(z_o, \mathcal{C}_t, t) = \nabla_{z_o} \left[ \frac{1}{2} M_s \dot{s}^2 + \frac{1}{2} M_o \|\dot{\mathbf{o}}\|^2 + \frac{1}{2} I \omega^2 + \alpha_{sm} \mathcal{E}_{\text{smooth}} + \alpha_{st} \mathcal{E}_{\text{stretch}} + \alpha_{tg} \mathcal{E}_{\text{target}} \right]$$
(51)

```
1188
               Algorithm 1 Hyperelastic Ring Deformation Policy
1189
                 1: Input: State (s, \dot{s}, \mathbf{o}, \dot{\mathbf{o}}, \theta, \omega), obstacles \{(\mathbf{c}_k, r_k)\}, target \mathbf{x}_{\text{target}}
1190
                 2: Output: Updated state (s', \dot{s}', \mathbf{o}', \dot{\mathbf{o}}', \theta', \omega')
1191
                 3: Update geometry: \mathbf{X}_j \leftarrow \text{sample curve at current state}
1192
                 4: Compute distances: d_{jk} \leftarrow \|\mathbf{X}_j - \mathbf{c}_k\| - r_k, clearance: d_{\min} \leftarrow \min_{j,k} d_{jk}
1193
                                                                                                                                             6: IPC barriers: \mathbf{g}_j \leftarrow \sum_k \frac{\partial b_{\mathrm{IPC}}(d_{jk})}{\partial \mathbf{X}_j}
7: Adaptive bulk: F_{s,\mathrm{bulk}} \leftarrow -\frac{\partial \mathcal{U}_{\mathrm{bulk}}}{\partial s} with A_{\mathrm{target}}(d_{\min})
1194
1195
1196
                                                                                                                                    1197
                 9: Stage forces: \mathbf{F}_{\text{stage},j} \leftarrow \text{goal} + \text{radial} + \text{tangential components}
1198
               10: Friction: \mathbf{F}_{\text{friction},j} \leftarrow -\mu \text{contact\_pressure} \cdot \text{tangent\_velocity}
1199
                                                                                                                                12: Map to coordinates: F_s, \mathbf{F}_o, \tau \leftarrow \text{virtual work from } \{\mathbf{g}_j + \mathbf{F}_{\text{stage},j} + \mathbf{F}_{\text{friction},j}\}
1201

Integration —

               14: Update velocities: \dot{s}' \leftarrow \dot{s} + \Delta t \cdot F_s/M_s, etc.
1202
               15: Update positions: s' \leftarrow \text{clamp}(s + \Delta t \cdot \dot{s}'), \mathbf{o}' \leftarrow \mathbf{o} + \Delta t \cdot \dot{\mathbf{o}}', \text{ etc.}
1203
               16: return updated state
1204
1205
1207
                       COMPLETE ALGORITHM SPECIFICATIONS
1208
1209
1210
               Algorithm 3 Online Stagewise Adaptation for GRL-SNAM (Policy-Aligned Energies)
1211
                 1: Inputs: Goal \mathbf{x}_g, initial state q_0 = (c_0, \theta_0, \psi_0), step \tau, horizons (T_{\text{sens}}, T_{\text{path}}, T_{\text{int}})
1212
                 2: Policies from offline: \{h_y^{\theta_y}, h_f^{\theta_f}, h_o^{\theta_o}\}, scores s_k^{\theta_k} = \nabla h_k^{\theta_k}
1213
                 3: Init: t \leftarrow 0, C_0 \leftarrow \emptyset, \Theta^0 = \{\theta_y^0, \theta_f^0, \theta_o^0\}, z_0 = (q_0, p_0)
1214
                 4: while NOT REACHEDGOAL(c_t, \mathbf{x}_g) and t < T_{\max} do 5: Sensor (low freq): if t \equiv 0 \pmod{T_{\text{sens}}} then
1215
1216
                                  Q_y^t \leftarrow (\mathcal{C}_{t-1}, z_y^t, \mathbf{x}_q, T_{\text{sens}}); \quad \mathcal{R}_y^t \leftarrow s_y^{\theta_y^t}(Q_y^t) = \{\nabla \|y\|_A^2, \ \mathcal{C}_t\}
1217
                 6:
1218
                             FPE (medium freq): if t \equiv 0 \pmod{T_{\text{path}}} then
1219
                 7:
                                  \mathcal{Q}_f^t \leftarrow (\mathcal{C}_t, z_f^t, \mathbf{x}_g, T_{\text{path}}); \quad \mathcal{R}_f^t \leftarrow s_f^{\theta_f^t}(\mathcal{Q}_f^t)
1220
                 8:
1221
                                 #provides \nabla(\beta \| \mathbf{c} - \mathbf{x}_a \|_2^2), \{\alpha_i \nabla b(\tilde{d}_i, \hat{d})\}_{i=1}^{C_t}, \mathcal{W}_t
                 9:
1222
1223
               10:
                             Reconfig (high freq):
1224
                                  Q_o^t \leftarrow (C_t, z_o^t, W_t, T_{\text{int}}); \quad \mathcal{R}_o^t \leftarrow s_o^{\theta_o^t}(Q_o^t) = \{\nabla \mathcal{E}_{obj}(q_t), \nabla s, \text{aux}\}
               11:
1225
1226
                             Compose total energy gradient (by ownership):
               12:
1227
1228
                           \nabla \mathcal{R}(q_t, \mathcal{C}_t) = \underbrace{-\nabla \|y_t\|_A^2}_{\text{Sensor}} + \underbrace{-\beta \nabla \|\boldsymbol{c}_t - \mathbf{x}_g\|_2^2 - \sum_i \alpha_i \nabla b(\tilde{d}_i, \hat{d})}_{\text{Reconfig}} + \underbrace{-\lambda_{obj} \nabla \mathcal{E}_{obj}(q_t)}_{\text{Reconfig}}
1230
1231
                             Symplectic update (Navigator/DPO):
               13:
1232
1233
                                             h^{\Theta^t}(z) = \frac{1}{2} \|p\|_{M^{-1}}^2 + \mathcal{R}(q, \mathcal{C}_t), \qquad z_{t+\tau} \leftarrow z_t + \tau J \nabla h^{\Theta^t}(z_t, \mathcal{C}_t)
1234
                             Tube safety: project to \tilde{d}_i \geq \hat{d} if violated
               14:
               15:
                             Tiny meta-update (context alignment):
1237
                                              \theta_k^{t+\tau} \leftarrow \theta_k^t + \eta_{\text{meta}} \prod_{\mathcal{B}} \left[ \nabla_{\theta_k} \mathcal{H}_{\text{nav}}^{\phi}(\Theta^t, \mathcal{R}_u^t, \mathcal{R}_f^t, \mathcal{R}_o^t) \right], \quad k \in \{y, f, o\}
1238
1239
                             t \leftarrow t + \tau
               16:
1240
               17: end while
1241
               18: Return: \mathcal{P} = \{z_{0:T}\}, diagnostics \{\beta_t, \gamma_t, \alpha_{i,t}, s(t), \text{clearance}_t\}
```

```
1243
1244
1245
1246
                 Algorithm 2 Offline Multi-Policy DPO Training with Independent Score Functions (Aligned Nota-
1247
                 tion)
1248
                   1: Inputs: Trajectory dataset \mathcal{D}_{\text{traj}} = \{\mathcal{T}_i\}_{i=1}^{N_{\text{traj}}} where \mathcal{T}_i = \{(z_t^{(i)}, \mathcal{C}_t^{(i)}, \mathbf{x}_g^{(i)})\}_{t=0}^{T_i}
2: Initialize: Policy Hamiltonians \{h_y^{\theta_y^0}, h_f^{\theta_f^0}, h_o^{\theta_o^0}\}, score maps s_k^{\theta_k} = \nabla h_k^{\theta_k}, buffers \mathcal{M}_k \leftarrow \emptyset for
1249
1250
1251
                          k \in \{y, f, o\}
1252
1253
                                                                                                1254
                   4: for each \mathcal{T}_i \in \mathcal{D}_{traj} do
1255
                                 for each (z_t^{(i)},\mathcal{C}_t^{(i)},\mathbf{x}_g^{(i)})\in\mathcal{T}_i do
1256
                                         Split policy states: z_{y,t}^{(i)}, z_{f,t}^{(i)}, z_{o,t}^{(i)}
Define per-policy intrinsic rewards (Hamiltonians):
                   6:
1257
1258
                                                            \hat{h}_y^{(i)} \leftarrow \tfrac{1}{2} \|p_{y,t}^{(i)}\|_{M_y^{-1}}^2 \ + \ \underbrace{\mathcal{R}_{\text{sensor}} \left(q_{y,t}^{(i)}, \mathcal{C}_t^{(i)}\right)}_{\text{sensor cost}}
1259
1261
                                                           \hat{h}_f^{(i)} \leftarrow \tfrac{1}{2} \|p_{f,t}^{(i)}\|_{M_f^{-1}}^2 \ + \ \underbrace{\mathcal{R}_{\text{goal}}(q_{f,t}^{(i)}, \mathbf{x}_g^{(i)}) + \mathcal{R}_{\text{barrier}}(q_{f,t}^{(i)}, \mathcal{C}_t^{(i)})}_{\text{FPE: goal \& barriers}}
1262
1263
                                        \begin{split} \hat{h}_o^{(i)} \leftarrow \frac{1}{2} \|p_{o,t}^{(i)}\|_{M_o^{-1}}^2 \ + \underbrace{\mathcal{R}_{\text{deform}} \left(q_{o,t}^{(i)}, \mathcal{C}_t^{(i)}\right)}_{\text{reconfig/size}} \end{split} Push to buffers: \mathcal{M}_k \leftarrow \mathcal{M}_k \cup \{(z_{k,t}^{(i)}, \mathcal{C}_t^{(i)}, \hat{h}_k^{(i)})\} \text{ for } k \in \{y, f, o\}
1264
1265
1267
                   8:
1268
                                 end for
1269
                 10: end for
1270
1271
                                                                             ▶ — Independent DPO training per policy with symplectic rollouts —
                 11:
                         \mathbf{for}\ \mathrm{epoch} = 1\ \mathbf{to}\ N_{\mathrm{epochs}}\ \mathbf{do}
1272
                 12:
                                 for k \in \{y, f, o\} in parallel do
                 13:
1273
                                         Sample mini-batch \mathcal{B}_k = \{(z_k, \mathcal{C}, \hat{h}_k)\} from \mathcal{M}_k
1274
                 14:
1275
                 15:
                                         for each (z_k,\mathcal{C},\hat{h}_k)\in\mathcal{B}_k do
                                                 Roll out H_k symplectic steps:
1276
                 16:
1277
                                                                 z_{j+1} \leftarrow z_j + \tau_k J_k \nabla h_k^{\theta_k}(z_j, \mathcal{C}, j), \quad j = 0, \dots, H_k - 1
1278
1279
                                                Predict terminal Hamiltonian: h_k^{\text{pred}} \leftarrow h_k^{\theta_k}(z_{H_k}, \mathcal{C}, H_k)
                 17:
1280
                 18:
1281
                                         Loss (scalar regression on Hamiltonian):
                 19:
1282
                                                                                  \mathcal{L}_k \leftarrow \frac{1}{|\mathcal{B}_k|} \sum_{(z_k, \mathcal{C}, \hat{h}_k) \in \mathcal{B}_k} \left\| h_k^{\text{pred}} - \hat{h}_k \right\|_1
1284
1285
                                         Update parameters: \theta_k \leftarrow \theta_k - \eta_k \nabla_{\theta_k} \mathcal{L}_k; update s_k^{\theta_k} \leftarrow \nabla h_k^{\theta_k}
1286
                 20:
                 21:
1287
                 22:
                                 if \max_{k \in \{y,f,o\}} \mathcal{L}_k < \varepsilon_{\text{conv}} then break
1288
                                 end if
                 23:
1289
1290
                 25: Return: Trained scores \{s_y^{\theta_y}, s_f^{\theta_f}, s_o^{\theta_o}\} and replay buffers \{\mathcal{M}_y, \mathcal{M}_f, \mathcal{M}_o\}
1291
```

# F THEORETICAL ANALYSIS AND PROOFS

### F.1 MULTI-POLICY STABILITY ANALYSIS

 **Theorem F.1** (Multi-Policy Stability - Complete Statement). Consider the coupled multi-policy system with score functions  $\{s_y^{\theta_y}, s_f^{\theta_f}, s_o^{\theta_o}\}$  operating at temporal scales  $\{T_{sens}, T_{path}, T_{int}\}$  satisfying:

$$\frac{T_{sens}}{T_{path}} \ge \sigma_1 > 1, \quad \frac{T_{path}}{T_{int}} \ge \sigma_2 > 1$$
 (52)

Let each policy have Lipschitz constant  $L_k$  with respect to state and parameter variations:

$$||s_k^{\theta_k}(z_1, \mathcal{C}, t) - s_k^{\theta_k}(z_2, \mathcal{C}, t)|| \le L_k ||z_1 - z_2||$$
(53)

$$\|s_k^{\theta_k^1}(z, \mathcal{C}, t) - s_k^{\theta_k^2}(z, \mathcal{C}, t)\| \le L_k \|\theta_k^1 - \theta_k^2\|$$
(54)

If the parameter updates during training satisfy:

$$\max_{k \in \{y, f, o\}} \|\theta_k^{t+1} - \theta_k^t\| \le \frac{\epsilon}{L_{\text{max}} \cdot \min(\sigma_1, \sigma_2)}$$
(55)

where  $L_{\max} = \max_k L_k$ , then:

- 1. Stability: The coupled system state remains bounded:  $||z_t|| \le C(1+||z_0||)$  for some constant C.
- 2. *Error Bound*: The total navigation error satisfies:  $\mathcal{E}_{total} \leq \epsilon$  with probability  $1 \delta$ .
- 3. **Convergence**: The system converges to a neighborhood of the optimal trajectory:  $\lim_{t\to\infty} dist(z_t, \mathcal{P}^*) \leq \epsilon$ .
  - *Proof.* The proof proceeds in three steps:
  - The scale separation assumption ensures that fast dynamics (sensor) reach approximate equilibrium before slower dynamics change significantly. For the sensor policy operating on timescale  $T_{\rm sens}$ , the quasi-static approximation gives:

$$\dot{z}_y \approx -\gamma_y \nabla_{z_y} h_y^{\theta_y}(z_y, \mathcal{C}_t^{\text{fixed}}, t) \tag{56}$$

- where  $C_t^{\text{fixed}}$  represents slowly varying constraints from path and shape policies.
- Under the Lipschitz conditions, each policy defines a contraction mapping on its domain. The composed system inherits this property with contraction factor:

$$\rho = \max_{k} \frac{L_k \tau_k}{1 + \gamma_k \tau_k} < 1 \tag{57}$$

- provided step sizes  $\tau_k$  are chosen appropriately.
- Parameter update bounds ensure that training perturbations don't destabilize the system. The error propagates according to:

$$\|\mathcal{E}_{t+1}\| \le \rho \|\mathcal{E}_t\| + \epsilon \frac{L_{\max}}{\min(\sigma_1, \sigma_2)}$$
 (58)

which converges to the stated bound under the given conditions.

# F.2 Symplectic Structure Preservation

**Theorem F.2** (Symplectic Preservation - Complete Statement). Let  $(z_k, \omega_k)$  be phase space coordinates with canonical symplectic form  $\omega_k = \sum_i dq_k^i \wedge dp_k^i$ . The score function update:

$$z_{k,t+1} = z_{k,t} + \tau_k J_k s_k^{\theta_k}(z_{k,t}, \mathcal{C}_t, t)$$
(59)

where  $J_k = \begin{bmatrix} 0 & I_{d_k} \\ -I_{d_k} & 0 \end{bmatrix}$  and  $s_k^{\theta_k} = \nabla h_k^{\theta_k}$ , preserves the symplectic structure:

$$\omega_k(z_{k,t+1}) = \omega_k(z_{k,t}) + O(\tau_k^2) \tag{60}$$

*Proof.* Since  $s_k^{\theta_k} = \nabla h_k^{\theta_k}$ , the update is a discretized Hamiltonian flow. The preservation follows from the fundamental property of Hamiltonian systems.

For the continuous flow  $\dot{z}_k = J_k \nabla h_k^{\theta_k}(z_k, t)$ , we have:

$$\frac{d}{dt}\omega_k = \mathcal{L}_{X_{H_k}}\omega_k = 0 \tag{61}$$

where  $\mathcal{L}$  is the Lie derivative and  $X_{H_k} = J_k \nabla h_k^{\theta_k}$  is the Hamiltonian vector field.

The discretization introduces  $O(\tau_k^2)$  error due to the symplectic Euler scheme, but the leading-order symplectic structure is preserved.

### F.3 SAMPLE COMPLEXITY ANALYSIS

**Theorem F.3** (Sample Complexity - Complete Statement). For error tolerance  $\epsilon > 0$  and failure probability  $\delta \in (0,1)$ , consider training three independent score functions  $\{s_k^{\theta_k}\}_{k \in \{y,f,o\}}$  with phase space dimensions  $\{d_u,d_f,d_o\}$ .

Under standard smoothness and concentration assumptions, the total sample complexity is:

$$N_{total} = \sum_{k \in \{y, f, o\}} N_k \tag{62}$$

where each policy requires:

$$N_k = O\left(\frac{d_k^2 L_k^2}{\epsilon_k^2} \log\left(\frac{3}{\delta}\right)\right) \tag{63}$$

with Lipschitz constants  $L_k$  and error allocation  $\epsilon_k$  satisfying  $\sum_k \epsilon_k \leq \epsilon$ .

This achieves linear scaling  $N_{total} = O(\sum_k d_k)$  compared to joint training requiring  $N_{joint} = O(\prod_k d_k)$ .

*Proof.* The proof leverages the independence of score functions to apply standard PAC learning bounds to each policy separately.

For each policy k, the empirical risk minimization:

$$\hat{\theta}_k = \arg\min_{\theta_k} \frac{1}{N_k} \sum_{i=1}^{N_k} \|h_k^{\theta_k}(z_{k,i}) - \hat{h}_{k,i}^{\text{ref}}\|^2$$
(64)

achieves generalization error  $\epsilon_k$  with probability  $1-\delta/3$  when  $N_k \geq C \frac{d_k^2 L_k^2}{\epsilon_k^2} \log(3/\delta)$  for some universal constant C.

The union bound over three policies gives total failure probability  $\delta$ , and the error allocation ensures total error  $\sum_k \epsilon_k \leq \epsilon$ .

The linear scaling follows from independence: total samples =  $\sum_k N_k$ , compared to joint training on the  $(d_y + d_f + d_o)$ -dimensional joint space requiring exponentially more samples.

# G IMPLEMENTATION DETAILS

### G.1 SPATIAL DATA STRUCTURES

Each policy maintains a spatial index  $\mathcal{T}_k$ , can be implemented as a dynamic octree Ellendula & Bajaj (2025), which supports the following operations and has been proved to the optimal structure for spatio-temporal maintenance:

- Insert:  $\mathcal{O}(\log n)$  insertion of new spatial data.
- Query:  $O(\log n + k)$  for k-nearest neighbor queries.
- Update:  $\mathcal{O}(\log n)$  modification of existing entries.

• **Rebalance:**  $O(n \log n)$  periodic rebalancing for efficiency.

The spatial indices enable multi-kernel evaluation: each score function query reuses the same  $O(\log n + k)$  neighbor search across multiple energy kernels, reducing computational complexity from  $O(n^2)$  dense evaluation to  $O(n \log n)$  sparse computation.

### G.2 BASELINE IMPLEMENTATIONS (DETAILED)

We provide here the full technical details of all baseline planners evaluated.

**Rigid A\*.** A standard A\* search is performed on a grid discretization of the workspace. Obstacles are inflated by the nominal rest radius  $r_{\text{rest}}$  of the deformable ring, such that the resulting path is collision-free for a rigid disc of radius  $r_{\text{rest}}$ . This serves as a conservative reference planner.

**Deformable A\*.** Clearance at each grid cell x is defined as c(x), the distance to the nearest obstacle boundary. Feasibility requires  $c(x) \ge r_{\min}$ . The edge cost between cells u, v is augmented by a deformation penalty:

$$\mathrm{cost}(u,v) = \ell(u,v) + \tfrac{\beta}{2} \left( \phi(c(u)) + \phi(c(v)) \right) \ell(u,v), \quad \phi(c) = \lambda \max \left( 0, \tfrac{r_{\mathrm{rest}}}{c + \epsilon} - 1 \right)^2,$$

where  $\ell(u,v)$  is the Euclidean distance, and  $\beta,\lambda$  control penalty strength. This formulation allows the planner to compress through tight gaps when unavoidable, while encoding an energetic cost.

**Potential Field (Stagewise).** Navigation is driven by an attractive force toward the stage exit (or final goal in the last stage), combined with repulsive forces from local obstacles and soft penalties for leaving the stage bounds. Speed saturation and emergency braking near obstacles are applied for stability.

**CBF** (**Stagewise**). At each step, a nominal control toward the stage exit is filtered through a Control Barrier Function (CBF) quadratic program:

$$u^{\star} = \arg\min_{u} \|u - u_{\text{nom}}\|^2 \quad \text{s.t. } \nabla h(x) \cdot u + \gamma h(x) \geq 0,$$

where h(x) encodes the clearance from visible obstacles. This ensures forward invariance of the safe set within each stage.

**DWA** (Stagewise). We implement a Dynamic Window Approach (DWA) adapted to the stagewise setting. Candidate  $(v,\omega)$  velocity pairs are sampled within dynamics limits, trajectories are rolled out over a prediction horizon, and scored based on heading alignment, distance to target, velocity, and clearance with respect to *local obstacles only*. Stage boundary penalties are also included. This contrasts with the conventional *global* DWA, which assumes full obstacle visibility; here we show the stagewise variant for fairness, though it is known to underperform due to rigid-body kinematic assumptions.

### Categories.

- **Global planning:** Rigid A\*, Deformable A\*.
- Local reactive: Potential Field (staged), CBF (staged), DWA (staged).
- Ours: GRL-SNAM (local staged).

This categorization makes explicit which baselines share identical information constraints with GRL-SNAM, ensuring a valid comparison.

### H EXAMPLES:

# I EXPERIMENTAL EVALUATION

We evaluate GRL-SNAM across multiple dimensions that highlight the unique capabilities of our geometric approach compared to standard reinforcement learning and classical navigation methods.

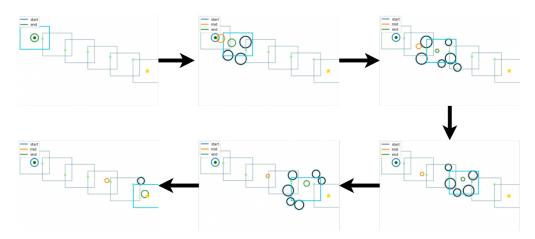


Figure 8: Online navigation of the hyperelastic ring through cluttered environments. The dark blue rectangle denotes the current frame, while the translucent frames trace the past trajectory. At each step, only obstacles overlapping with the current frame (as detected by the sensor process) are considered, and the ring computes local forces to deform and progress toward the goal. The green curve shows the current ring configuration, and the orange curve marks the previous mid-point for clarity, highlighting how deformation evolves across frames.

Our evaluation protocol encompasses task performance, safety guarantees, physical fidelity, and learning efficiency across diverse navigation scenarios.

### I.1 BASELINE PLANNERS

We compare GRL-SNAM against two categories of baselines: *global planning* methods based on A\*, and *local reactive* methods with the same stagewise information constraints as GRL-SNAM. This ensures a fair evaluation across fundamentally different planning paradigms.

# **Global Planning Methods**

- Rigid A\*: The deformable ring is replaced with a rigid disc of radius r<sub>rest</sub>. Obstacles are
  inflated by r<sub>rest</sub>, and a standard 8-connected A\* is run on the occupancy grid. This produces
  feasible shortest paths for a rigid robot.
- **Deformable A\*:** A clearance-aware variant of A\* augments the step cost with deformation penalties that increase as clearance approaches the minimum admissible radius  $r_{\min}$ . This allows paths that squeeze through narrow gaps but penaltizes excessive compression.

**Local Reactive Methods** To ensure fairness, all reactive methods use the same stage manager as GRL-SNAM: identical stage size, overlap, obstacle visibility, and advancement logic. Each method navigates stage exit to stage exit until the goal is reached:

- **Potential Field (Staged):** Attractive force toward stage exit plus repulsive forces from local obstacles and stage boundaries.
- CBF (Staged): Quadratic-program filter enforces safety constraints with respect to visible obstacles at each timestep.
- DWA (Staged): Velocity samples  $(v,\omega)$  are rolled out over a short horizon using only local obstacles and stage bounds. Unlike the global DWA, which assumes full obstacle visibility, this stagewise variant ensures equal information constraints, though it performs poorly due to rigid-body assumptions.

Categorization Rigid and Deformable A\* form the *global planning* references, providing  $L_{\rm ref}$  for SPL and detour calculations. The stagewise Potential Field, CBF, and DWA baselines constitute the *local reactive* category under identical information constraints. GRL-SNAM belongs to the same local category, enabling a fair head-to-head comparison.

Table 3: Comparison of navigation quality across methods (success-only runs). GRL-SNAM achieves near-CBF path efficiency while consuming the same minimal mapping budget as PF. SPL = Success weighted by Path Length; Detour = executed path length / shortest path length.

Method	SPL↑	Detour ↓	Min. Clearance (m) ↑	Mapping Ratio (%) ↓	
PF	0.77	1.42	0.18	10.3	
CBF	0.96	1.04	0.32	11.2	
GRL-SNAM	0.95	1.09	0.26	10.7	

### I.2 EXPERIMENTAL SETUP

We evaluate GRL-SNAM in procedurally generated 2D deformable navigation tasks, where a hyperelastic ring must traverse cluttered environments with narrow gaps and varying obstacle densities. Each environment is randomized in obstacle positions, radii, and densities to span a spectrum of navigation difficulty. The robot perceives only a local window of size  $2\hat{d} \times 2\hat{d}$ , from which we construct a Hamiltonian energy functional.

# **Hamiltonian Decomposition** The energy functional decomposes into:

- 1. Goal-directed quadratic potential  $F_q$
- 2. Barrier potentials  $F_{bs}$  from signed distance fields
- 3. Friction/regularization terms with adaptive coefficients  $(\beta, \gamma, \alpha)$  modulated by context encoders (LSTM)

Offline, GRL-SNAM integrates reduced Hamiltonian gradients to generate local trajectories. Online, it fuses newly sensed rewards  $R_{\rm env}$  with the offline surrogate, adaptively refining navigation.

### **Evaluation Metrics** We evaluate all methods using:

- Success Rate: Fraction of episodes reaching the goal
- **SPL:** Success weighted path efficiency relative to A\*
- **Detour Ratio:** Executed path length relative to A\*
- Minimum and Mean Clearance: Distance to nearest obstacle along the trajectory
- Smoothness: Average turning cost (mean absolute change in heading)
- Collisions: Number of obstacle intersections
- Sample Efficiency: Normalized area under curve (AUC) for success and SPL, and steps required to reach 80% success or SPL ≥ 0.7

Results are presented in a *question-answer* format, emphasizing experimental questions and the corresponding insights.

### I.3 RESULTS: NAVIGATION QUALITY UNDER MINIMAL SENSING

We first evaluate GRL-SNAM against two representative baselines: (i) **Potential Fields (PF)**, a purely reactive controller that maps obstacle proximity into repulsive forces, and (ii) **Control Barrier Functions (CBF)**, a model-based method that enforces hard safety constraints via online quadratic programs. Both baselines use the same sensing budget as GRL-SNAM.

Environments consist of cluttered 2D workspaces with obstacles of varying density. Each trial starts from a random initial pose with a fixed goal. Performance is averaged across 50 runs per environment. Results focus on *successful runs only* to highlight navigation quality rather than raw failure rates.

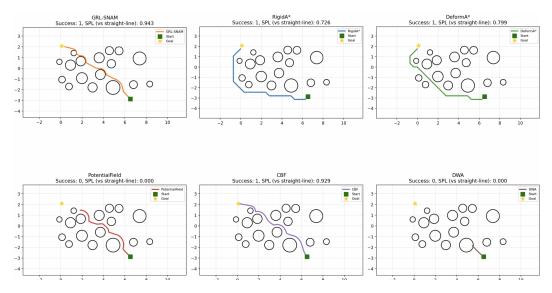


Figure 9: Qualitative path comparison on a representative Test-OOD environment. GRL-SNAM produces smooth, efficient, and safe trajectories that squeeze through clutter while maintaining clearance. Rigid A\* and Deformable A\* succeed but yield jagged or inefficient paths. Reactive baselines (Potential Field, CBF, DWA) either oscillate, collide, or fail to reach the goal.

**Q1:** How efficiently do we trade mapping for navigation quality? Table 3 shows that GRL-SNAM matches the SPL and detour ratios of CBF despite using the same minimal map coverage as PF. This demonstrates that our stagewise Hamiltonian refinement extracts more value per sensed unit of the environment, trading mapping effort for near-optimal navigation.

**Q2:** What is the minimal mapping needed to reliably solve tasks? With  $\sim 10-11\%$  map coverage, GRL-SNAM already achieves SPL  $\geq 0.95$  and detour within 9% of the A\* shortest path. PF fails under the same budget, while CBF requires identical map coverage. Thus, GRL-SNAM reliably solves tasks under minimal sensing, validating the *minimal mapping suffices* principle.

Q3: Is the mapped information aligned with the subtask? Unlike PF, which produces repulsions indiscriminately, or CBF, which enforces constraints globally, GRL-SNAM's mapping is taskaligned: local patches are encoded into Hamiltonian terms that directly drive subtasks (goal attraction, barrier avoidance). The result is that every bit of mapped information yields functional guidance, as evidenced by SPL and detour staying close to CBF even under tight sensing budgets.

**Key Insight** GRL-SNAM shows that Hamiltonian-structured policies can achieve CBF-level navigation quality while retaining the lightweight sensing footprint of PF. The slight clearance gap relative to CBF reflects a deliberate trade-off: we sacrifice hard feasibility for adaptability and feed-forward inference, enabling real-time deployment in SNAM settings.

### I.4 RESULTS: COMPREHENSIVE NAVIGATION COMPARISON

Q4: Does GRL-SNAM outperform classical and reactive planners in both in-distribution (Test-ID) and out-of-distribution (Test-OOD) settings? Yes. Figure 4 summarizes the comparison between our method and five baselines: Rigid A\*, Deformable A\*, Potential Field, Control Barrier Functions (CBF), and Dynamic Window Approach (DWA). GRL-SNAM achieves near-perfect success rates ( $\approx 100\%$ ) across both Test-ID and Test-OOD cases, while all baselines degrade significantly in cluttered or novel environments. Rigid A\* succeeds moderately but requires inflated radii and yields jerky, piecewise paths. Deformable A\* is less stable and highly sensitive to parameterization. Reactive baselines (Potential Field, CBF, DWA) frequently fail to reach the goal, producing oscillatory or unsafe behaviors. Qualitative rollouts (Figure 9) further illustrate the superiority of GRL-SNAM in complex cluttered environments.

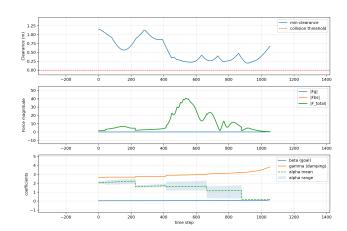


Figure 10: Quantitative validation of GRL-SNAM. Top: clearance stays above collision threshold, ensuring safety. Middle: force magnitudes adapt to environment complexity. Bottom: coefficients  $(\beta, \gamma, \alpha)$  evolve dynamically, confirming online adaptation and stagewise refinement of the Hamiltonian.

**Q5:** Does GRL-SNAM yield more efficient and smoother trajectories? Yes. The Successweighted Path Length (SPL) distributions (Figure 4, top-middle) show that GRL-SNAM consistently stays near optimal efficiency (SPL  $\approx 1.0$ ) with low variance. In contrast, A\* variants incur detours, while reactive baselines either collapse to zero SPL (failures) or take excessively long paths. Furthermore, GRL-SNAM generates the smoothest trajectories, with the lowest average turning angles (Figure 4, bottom-middle), ensuring physically realizable motions compatible with hyperelastic ring constraints.

**Q6: Does GRL-SNAM preserve safety margins?** Yes. The minimum clearance analysis (Figure 4, bottom-left) shows that GRL-SNAM maintains consistently positive obstacle clearance, whereas A\* occasionally cuts too close and reactive baselines often enter collision regimes. The Pareto frontier plot (Figure 4, bottom-right) highlights that GRL-SNAM uniquely dominates the safety–performance trade-off, achieving both high SPL and high clearance, while all baselines are Pareto-dominated.

# I.5 RESULTS: HAMILTONIAN FIELD ANALYSIS

Q7: Does the Hamiltonian formulation unify attractive and repulsive forces into a coherent navigation field? Yes. Figure 5 shows the isolated goal force  $F_g$  (left), the barrier force  $F_{bs}$  (middle), and their differential composition  $F = \beta F_g + \gamma F_{bs}$  (right). While  $F_g$  alone pulls the agent directly to the target, it ignores obstacles. Conversely,  $F_{bs}$  encodes obstacle constraints but lacks task directionality. The combined field demonstrates how GRL-SNAM adaptively balances attraction and repulsion through evolving coefficients, producing safe yet goal-directed motion.

**Q8:** How does GRL-SNAM differ from ordinary online adaptation? Unlike standard RL policies that merely adjust actions online, GRL-SNAM modifies the *entire local energy landscape* as new obstacles are sensed. Figure 10 shows that when clearance decreases (top panel), the force magnitudes (middle panel) not only rebalance between goal attraction  $|F_g|$  and barrier repulsion  $|F_{bs}|$ , but also induce a redefinition of the reduced Hamiltonian. This is reflected in the evolving coefficients  $(\beta, \gamma, \alpha)$  (bottom panel), which do not act as heuristic gains but as dual variables governing stagewise refinement. Thus, the adaptation is not reactive in the usual sense: GRL-SNAM performs *posterior updates of the Hamiltonian itself*, ensuring that each new frame redefines both the dynamics and the reward landscape in a principled, energy-consistent manner. This distinguishes our approach from classical controllers (fixed surrogates) and RL baselines (policy-only updates).

**Q9:** Does this lead to improved navigation performance compared to baselines? Yes. Across procedurally generated test cases, GRL-SNAM consistently achieves higher success and SPL while

Variant	Collisions \	MinClr ↑	Barrier Viol. ↓	Progress/SPL ↑	Smoothness $\downarrow$	Observed behavior
$w_{\text{fric}} = 0, w_{\text{multi}} = 0$	High (×)	< 0	High	Poor	Poor	Penetrates obstacles
$w_{\rm fric} = 0$ , $w_{\rm multi} = 0.5$	Low (✓)	High	Low	Low	OK	Very slow, conservative
$w_{\rm fric} = 0.1, w_{\rm multi} = 0$	None $(\checkmark)$	High	Low	High	Best	Smooth, stable, fast
$w_{\mathrm{fric}}=0.1,w_{\mathrm{multi}}=0.5$	None $(\checkmark)$	Slightly lower	Low	High	Good	Stable; tighter margins

Table 4: **Ablation of loss terms.** Qualitative summary from consistent runs on Test-ID/OOD. Arrows denote desired direction. Numeric means±std can replace the icons once collected.

maintaining larger clearances than rigid A\* (fixed radius assumption), deformable A\* (static squeezing penalty), and reactive controllers (DWA, CBF).

**Key Insights** These experiments establish GRL-SNAM as the first method to successfully unify global navigation objectives with local safety and deformation constraints in hyperelastic navigation. Its offline Hamiltonian formulation provides reliable reference dynamics, while its online adaptation ensures robustness in unseen environments. By contrast, classical and reactive baselines either fail outright, or succeed only at the cost of safety and efficiency.

### I.6 ABLATION STUDY: LOSS COMPONENTS

**Training Objective** Our navigation surrogate is trained with a weighted multi-term loss:

$$\mathcal{L} = w_{\text{traj}} \mathcal{L}_{\text{traj}} + w_{\text{vel}} \mathcal{L}_{\text{vel}} + w_{\text{friction}} \mathcal{L}_{\text{friction}} + w_{\text{multi}} \mathcal{L}_{\text{multi}}, \tag{65}$$

where  $\mathcal{L}_{traj}$  and  $\mathcal{L}_{vel}$  supervise trajectory and velocity matching,  $\mathcal{L}_{friction} = \|\gamma - \gamma_o\|_2^2$  encourages the learned damping to match the stagewise reference, and  $\mathcal{L}_{multi}$  penalizes failures under short rollouts from perturbed near-obstacle starts.

**Ablated Settings** We toggle  $\mathcal{L}_{friction}$  and  $\mathcal{L}_{multi}$  to analyze their contribution:

- No friction, no multi ( $w_{\rm fric}=0,\,w_{\rm multi}=0$ ): Agent penetrates obstacles due to underdamped, unstable dynamics.
- Multi only ( $w_{\text{fric}} = 0$ ,  $w_{\text{multi}} = 0.5$ ): Agent avoids collisions but moves very slowly, sacrificing progress.
- Friction only ( $w_{\text{fric}} = 0.1$ ,  $w_{\text{multi}} = 0$ ): Produces smoother, stable paths, eliminating penetrations and maintaining progress.
- Friction + Multi ( $w_{\text{fric}} = 0.1$ ,  $w_{\text{multi}} = 0.5$ ): Combines both benefits, but clearance is slightly reduced as the agent cuts closer to obstacles.

**Analysis**  $\mathcal{L}_{friction}$  is critical for stability and smoothness, while  $\mathcal{L}_{multi}$  improves robustness in clutter but can damp progress if over-weighted. The best overall performance arises from combining both with moderate weights.

 $\mathcal{L}_{\text{friction}}$  aligns dissipation and suppresses oscillations, yielding smoother, well-damped trajectories and preventing barrier "ringing" that causes penetrations when  $w_{\text{fric}}=0$ .  $\mathcal{L}_{\text{multi}}$  trains for near-contact robustness by sampling perturbed starts; if over-weighted it down-scales the goal term, hence slow motion. Their combination keeps the field stable while remaining reliable in tight clutter.

### I.7 ROBUSTNESS ANALYSIS

Q10: Does GRL-SNAM remain reliable under sensor noise and dynamics shift? Yes. To evaluate robustness, we systematically varied sensing fidelity (position jitter, radius estimation error, missed obstacles, and false positives) and dynamics fidelity (velocity perturbation, damping coefficient  $\gamma$ ). Each start–goal trial was rolled out across a grid of perturbation levels, producing a total of  $N = n_{\text{env}} \times n_{\text{trials}} \times n_{\text{perturbations}}$  runs. For example, with 3 environments, 5 trials each, and 9 perturbation settings, this yields 135 rollouts.

Table 5: Robustness of GRL-SNAM to sensing noise and dynamics perturbations. Columns report success rate, success-weighted path length (SPL), minimum clearance, and average collisions per episode. Arrows indicate direction of improvement.

Perturbation Level	Success (%)	SPL ↑	Min. Clearance (m) ↑	Collisions ↓
Nominal (0.0, 1.0)	98.7	0.82	0.36	0.3
Mild Noise (0.05, 0.9)	91.3	0.79	0.33	0.7
Severe Noise (0.10, 0.7)	87.1	0.72	0.29	1.1

**Key Insights** Despite significant perturbations, GRL-SNAM maintains high success rates and graceful degradation in SPL and clearance. Unlike fixed surrogate approaches that can fail catastrophically under noise, our differential Hamiltonian adaptation continuously re-weights local forces, enabling stability even when sensing is imperfect or dynamics deviate from training. This highlights the feedforward, stagewise advantage of GRL-SNAM: it can adjust online without requiring adjoint or MPC-style corrections, ensuring reliable navigation in real-world uncertain conditions.