

GeoDT: Geometry-Aware Decision Transformer for Robust Safe Multi-Task Offline Reinforcement Learning

Anonymous authors
Paper under double-blind review

Abstract

Scaling offline reinforcement learning across heterogeneous tasks remains challenging, especially under safety constraints. In multi-task settings, features processed by a shared model may play different semantic roles across tasks, leading to semantic inconsistency, conflicting optimization signals, and performance degradation as task diversity increases. While prior multi-task and safe offline RL methods address parts of this challenge, few provide a unified framework that is both effective and safety-aware. We propose GeoDT (Geometry-Aware Decision Transformer), a framework for safe multi-task offline RL that biases cross-task sharing toward geometry-related trajectory structure to mitigate semantic inconsistency. GeoDT learns to separate geometry-related structure from task-specific semantics, constructs geometry-aware context from prompt trajectories through relational structure induction and prototype memory, and incorporates safety by using cost signals to shape the feasible region of geometric reuse. The resulting context is fused with task-specific semantic features to condition a cost-aware Decision Transformer. To better assess behavior as task diversity increases, we further introduce the Task Scaling Robustness Score (TSRS) and Inter-Task Balance Score (ITBS), which measure performance retention and cross-task balance as the number of tasks increases. Experiments on multi-task safe offline RL benchmarks show that GeoDT achieves strong reward-cost trade-offs, improved robustness under increasing task diversity, and zero-shot adaptation to unseen safety budgets compared with competitive baselines. These results suggest that geometry-related trajectory structure can provide an effective basis for safe multi-task offline reinforcement learning.

1 Introduction

Real-world decision-making systems, from robotic control to safety-critical automation, increasingly require agents that can operate across multiple tasks while respecting explicit safety constraints (Wu et al., 2024). Offline reinforcement learning (RL) is a natural framework for such settings, since it learns from fixed datasets without unsafe online exploration. However, learning a single offline policy across heterogeneous tasks remains challenging, especially when a single policy must jointly handle varying objectives, contextual observations, and safety requirements.

A key obstacle in heterogeneous multi-task offline RL is semantic inconsistency under shared representation learning, where a common encoder must process features whose semantic roles may differ across tasks, often leading to task interference or negative transfer (Yu et al., 2021; Ishfaq et al., 2024). In heterogeneous multi-task settings, observations from different tasks are often aligned into a shared input format and processed by a common encoder, but aligned input dimensionality does not imply aligned semantics. Features handled by the same model can still correspond to different physical quantities, contextual cues, or safety-relevant signals across tasks. As a result, naive parameter sharing can induce representation conflict, conflicting optimization signals, and negative transfer as task diversity increases.

This challenge is particularly pronounced in the offline safe RL setting. Because training is restricted to static datasets, the agent cannot rely on additional exploration to disambiguate misleading shared representations or recover from harmful cross-task interference, while distributional shift further amplifies errors under joint

sharing (Fujimoto et al., 2019). At the same time, safety introduces an additional layer of difficulty: beyond learning task-relevant behavior, the agent must also determine which trajectory patterns remain feasible under task-dependent cost budgets (Lin et al., 2023). In other words, multi-task offline safe RL requires not only effective cross-task sharing, but also selective sharing that respects safety.

Existing approaches address only parts of this problem. Multi-task offline RL methods often reduce interference through task identifiers, task embeddings, or modular conditioning (Yang et al., 2020; Hu et al., 2024; Yu et al., 2021; Ishfaq et al., 2024), but such approaches either rely on explicit task identity or leave under-specified what structure should be shared across heterogeneous tasks. Prompt- and context-based policy learning methods improve adaptation and transfer (Laskin et al., 2022; Huang et al., 2024; Yoo et al., 2025), yet they primarily focus on performance generalization rather than semantic misalignment under joint training. Meanwhile, safe offline RL methods have shown strong results in single-task settings (Zheng et al., 2024; Koirala et al., 2024b; Yao et al., 2024; Gong et al., 2025; Chemingui et al., 2025b; Liu et al., 2023b; Zhang et al., 2023), but they do not explicitly address the representation conflict introduced by heterogeneous multi-task learning. As a result, few existing methods jointly tackle semantic inconsistency, robustness under task diversity, and safety in a unified offline multi-task framework.

We argue that a more stable basis for cross-task sharing lies in the *geometry- and dynamics-related structure* of trajectories. Although task-specific semantics, reward functions, and contextual cues may vary substantially, agent motion and environment interaction often still exhibit reusable geometric regularities. These regularities offer a more reliable anchor for cross-task sharing than raw semantics alone. Moreover, safety should not be treated merely as a decision-time constraint appended after representation learning; it should also shape which trajectory relations and structural patterns are feasible to reuse across tasks.

Motivated by this view, we propose **GeoDT** (Geometry-Aware Decision Transformer), a framework for safe multi-task offline RL. GeoDT learns to separate geometry-related structure from task-specific semantics, regularizes the geometric branch with a self-supervised trajectory consistency objective, and constructs geometry-aware context through relational structure induction and prototype memory. To incorporate safety without collapsing it into the representation itself, GeoDT uses cost information to modulate the feasible region of geometric reuse, so that safety influences which relational patterns and structural memories are emphasized. The resulting context is fused with task-specific semantic features to condition a cost-aware Decision Transformer (Chen et al., 2021), enabling robust and safety-aware policy learning without explicit task identifiers. To further support this motivation, we provide a simplified theoretical analysis and a targeted empirical gradient analysis in Appendix B, illustrating how semantic inconsistency under shared training can induce cross-task optimization difficulties and how geometry–semantics decomposition helps alleviate them.

Meanwhile, standard average reward and cost can hide two failure modes that are central to multi-task scaling: substantial degradation relative to single-task reference performance, and uneven performance concentrated on only a subset of tasks. To better assess robustness as task diversity increases, we further introduce two diagnostic metrics: the Task Scaling Robustness Score (TSRS) and the Inter-Task Balance Score (ITBS). TSRS and ITBS are designed to capture these two aspects, namely performance retention and cross-task balance, as the number of jointly trained tasks increases.

Our contributions are summarized as follows:

- We identify semantic inconsistency under shared encoding as a key obstacle in safe multi-task offline RL with heterogeneous tasks, and formalize a practical setting where tasks differ in objectives, contextual semantics, and safety constraints, with additional zero-shot evaluation under varying target safety budgets.
- We propose GeoDT, a geometry-aware decision transformer that biases cross-task sharing toward geometry- and dynamics-related trajectory structure through learnable state decomposition, self-supervised consistency learning, relational structure induction, prototype memory, and safety-aware feasible reuse.
- We introduce TSRS and ITBS as diagnostic measures of robustness under increasing task diversity, and conduct extensive experiments demonstrating strong reward–cost trade-offs, robustness under increasing task diversity, and zero-shot generalization to unseen safety budgets.

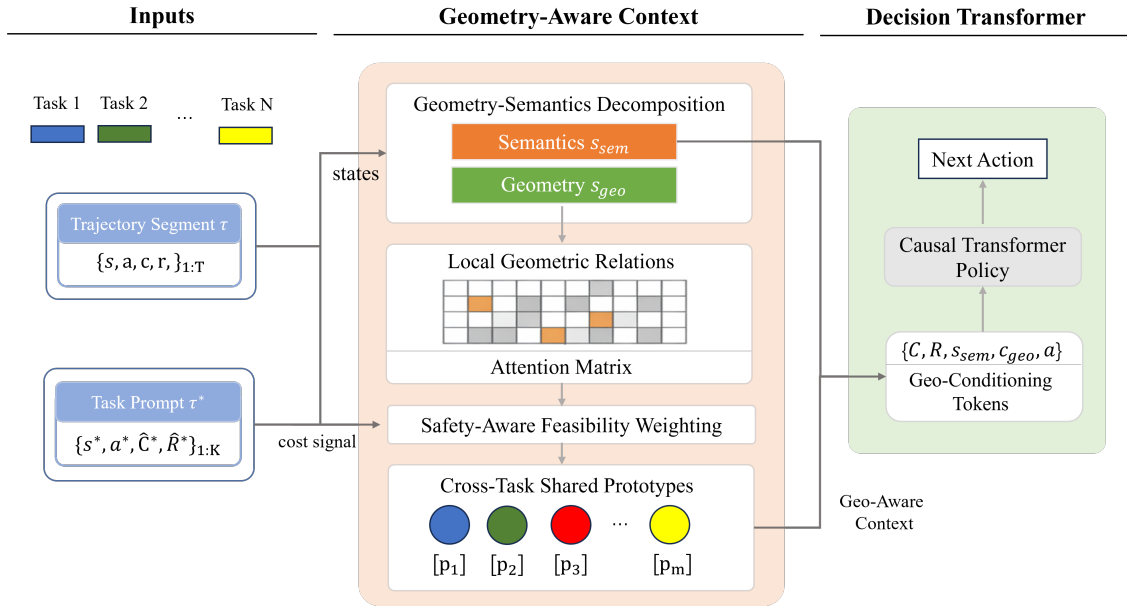


Figure 1: Overview of GeoDT. To mitigate semantic inconsistency across heterogeneous tasks, GeoDT separates geometry-related structure from task-specific semantics. Unlike prior methods that directly share task semantics or rely on task identifiers, GeoDT performs cross-task sharing through geometry-aware trajectory structure. A prompt trajectory τ^* is used to build the geometry-aware context c_{geo} through relational structure induction, safety-aware feasibility reweighting, and shared prototype retrieval. The resulting context is fused with semantic state features and fed into a cost-aware Decision Transformer for reward- and cost-conditioned action prediction.

The remainder of this paper is organized as follows. Section 2 reviews related work on safe offline RL, multi-task offline RL, and structure-aware sequence modeling. Section 3 introduces the problem formulation for multi-task safe offline RL. Section 4 presents GeoDT, including learnable state decomposition, geometry-aware structure induction, prototype memory, and safety-aware decision conditioning. Section 5 describes the experimental setup and empirical results, including robustness evaluation under increasing task diversity and ablation studies. Finally, Section 6 concludes with limitations and future directions.

2 Related Work

Safe offline reinforcement learning. Safe offline reinforcement learning studies how to learn policies from fixed datasets while satisfying safety constraints. Early work introduced constrained batch policy optimization and conservative policy learning under offline settings (Le et al., 2019; Fujimoto et al., 2019). Subsequent approaches improved safety through policy projection (Polosky et al., 2022), value regularization against unsafe out-of-distribution actions (Xu et al., 2022a), and stationary-distribution correction under cost constraints (Lee et al., 2022). More recent methods have adopted generative and sequence-modeling paradigms. CDT (Liu et al., 2023b) and SaFormer (Zhang et al., 2023) extend Transformer-based decision modeling to constrained settings, while FISOR (Zheng et al., 2024), OASIS (Yao et al., 2024), LSPC (Koirala et al., 2024b), FAWAC (Koirala et al., 2024a), TraC (Gong et al., 2025), and CAPS (Chemingui et al., 2025b) improve reward–cost trade-offs through diffusion modeling, latent-space optimization, feasibility-aware regression, unsafe-trajectory filtering, or adaptive constraint handling. These methods have demonstrated strong performance in single-task safe offline RL, but most assume fixed task semantics and do not address the additional representation conflict that arises when heterogeneous tasks are trained jointly.

Multi-task and context-based offline reinforcement learning. Offline multi-task RL aims to learn a shared policy or sequence model across multiple tasks from static data. Existing methods typically

mitigate task interference through sample routing, factorized representations, or task-conditioned modulation. Conservative data sharing (Yu et al., 2021) selectively routes transitions to reduce negative transfer, while MORL (Ishfaq et al., 2024) learns low-rank shared structure across tasks. Other approaches improve cross-task generalization through task-dependent architectures, latent skill abstractions, and prompt- or context-conditioned sequence modeling (Yang et al., 2020; Yoo et al., 2022; Laskin et al., 2022; Lee et al., 2023; Wang et al., 2025; Xu et al., 2022b; Hu et al., 2024; Yoo et al., 2025). However, most of these methods primarily target performance transfer and do not explicitly model the semantic inconsistency that emerges when heterogeneous observations are aligned into a shared representation but retain different semantic roles across tasks. Moreover, safety constraints are usually absent or handled only implicitly.

Structure-aware and geometry-aware trajectory modeling. A growing line of work suggests that effective generalization in sequential decision making depends not only on task identity, but also on how structural information is extracted from trajectories and interactions. Context-based sequence models leverage prompts or retrieved histories to infer task-relevant behavior patterns (Laskin et al., 2022; Xu et al., 2022b; Huang et al., 2024; Yoo et al., 2025). More broadly, representation learning methods for control often seek to separate shared dynamical structure from task-specific variation, for example through latent abstractions, temporal consistency, or structured memory (Gelada et al., 2019; Schwarzer et al., 2020; Khetarpal et al., 2024). Geometry-aware reinforcement learning has also been explored in other settings, such as robotic manipulation, where object geometry, interaction structure, or symmetry is used to improve policy generalization (Bellemare et al., 2019; Hoang et al., 2025).

Our work is most closely related to this structure-aware perspective, but differs from existing geometry-aware and representation-learning approaches in three important ways. First, we study *heterogeneous multi-task offline RL*, where the central challenge is semantic inconsistency under shared encoding rather than online exploration or embodiment-specific generalization. Second, we use *trajectory geometry* as the anchor for cross-task sharing, rather than explicit task identifiers, purely semantic context, or object-specific geometric priors. Third, we incorporate safety not merely as an auxiliary decision-time condition, but as a signal that shapes the feasible region of geometric reuse across tasks. This allows GeoDT to jointly address semantic inconsistency and safety while maintaining robust performance across diverse tasks in a unified multi-task offline RL framework.

3 Preliminaries

3.1 Constrained Markov Decision Processes and Safe Multi-Task Offline RL

Safe reinforcement learning is commonly formulated under the constrained Markov decision process (CMDP) framework (Altman, 2021). A CMDP is defined as

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

where \mathcal{S} is the state space, \mathcal{A} is the action space, $P : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$ is the transition distribution, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, $c : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}_{\geq 0}$ is the safety cost, and $\gamma \in [0, 1)$ is the discount factor. The goal is to learn a policy π that maximizes the expected discounted return while keeping the expected cumulative cost below a budget κ :

$$\max_{\pi} \mathbb{E}_{\tau \sim \pi}[R(\tau)] \quad \text{s.t.} \quad \mathbb{E}_{\tau \sim \pi}[C(\tau)] \leq \kappa, \quad (1)$$

where $R(\tau) = \sum_{t=1}^T \gamma^{t-1} r_t$ and $C(\tau) = \sum_{t=1}^T \gamma^{t-1} c_t$.

In this work, we consider a *multi-task safe offline RL* setting with a collection of N tasks, indexed by $i \in \{1, \dots, N\}$:

$$\{\mathcal{M}_i\}_{i=1}^N, \quad \mathcal{M}_i = (\mathcal{S}_i, \mathcal{A}, P_i, r_i, c_i, \gamma).$$

The tasks share the same action space \mathcal{A} but may differ in state semantics, transition dynamics, reward functions, and cost functions. In particular, even when state vectors are represented in a shared input format or have the same dimensionality across tasks, corresponding features may encode different physical meanings, giving rise to semantic inconsistency and cross-task interference during joint training.

For each task, we assume access only to an offline dataset collected by unknown behavior policies, with no additional environment interaction during training. Let $\mathcal{D} = \bigcup_{i=1}^N \mathcal{D}_i$ denote the union of task datasets, where each transition sequence in \mathcal{D}_i consists of tuples $(s_t, a_t, r_t, c_t, s_{t+1})$. At test time, the agent is required to produce actions that achieve favorable reward–cost trade-offs under a specified target safety budget, while generalizing across heterogeneous tasks without explicit task identifiers.

Our goal is to learn a single offline policy for heterogeneous safe RL tasks that remains robust as task diversity increases and is adaptive to different safety budgets. The main challenge is that naive parameter sharing across tasks can be ineffective when state dimensions are semantically misaligned. GeoDT addresses this challenge by identifying geometry-related structure that is more stable across tasks, using it to construct structured context for action prediction, and incorporating safety as a signal that influences which shared geometric patterns are feasible to reuse under a given budget.

3.2 Decision Transformer

Decision Transformer (DT) (Chen et al., 2021) formulates offline RL as a sequence modeling problem. Instead of learning value functions explicitly, DT models actions autoregressively conditioned on trajectory context and desired outcomes using a causal Transformer. A standard DT takes as input an interleaved sequence

$$(\hat{R}_1, s_1, a_1, \hat{R}_2, s_2, a_2, \dots),$$

where $\hat{R}_t = \sum_{i=t}^T r_i$ denotes the return-to-go at step t . Given a finite context window, the model predicts the next action by

$$a_t = f_\theta(\hat{R}_{t-L:t}, s_{t-L:t}, a_{t-L:t-1}),$$

where f_θ is implemented as a causal Transformer decoder.

Several subsequent works have extended DT to richer conditioning signals, such as prompts, context trajectories, and safety-related objectives (Lee et al., 2023; Xu et al., 2022b; Huang et al., 2024; Liu et al., 2023b; Zhang et al., 2023). Our method follows this sequence-modeling paradigm, but augments DT with geometry-aware structured context to better handle semantic inconsistency across heterogeneous tasks under safety constraints.

4 Method

Figure 1 illustrates the overall architecture of **GeoDT**. Instead of relying on explicit task identifiers or textual task descriptions, GeoDT uses short trajectory prompts sampled from the offline dataset to provide compact contextual information about task behavior and safety preference. For task \mathcal{T}_i , a prompt is defined as

$$\tau_i^* = (\hat{R}_{t:t+K}^*, \hat{C}_{t:t+K}^*, s_{t:t+K}^*, a_{t:t+K}^*), \quad (2)$$

where $\hat{R}_t^* = \sum_{j=t}^T r_j$, $\hat{C}_t^* = \sum_{j=t}^T c_j$, and K is the prompt length. Since K is much shorter than the full horizon, the prompt provides contextual information about task behavior and safety preferences without degenerating into direct imitation.

GeoDT is built on a simple principle: trajectory prompts contain both task-specific semantics and cross-task geometric regularities, and these two sources of information should be modeled differently. To this end, GeoDT first decomposes each state into a geometric component and a semantic component, and regularizes the geometric branch using a self-supervised consistency objective. It then constructs a geometry-aware context by combining prompt-level relational structure with a small set of learned geometric prototypes. Finally, safety is incorporated by using cost signals to reweight geometric relations, so that the resulting context emphasizes feasible trajectory patterns before being passed to a cost-aware Decision Transformer.

4.1 Geometry–Semantics Decomposition and Consistency Learning

A central challenge in multi-task offline RL is *semantic inconsistency*: even when observations are aligned into a shared input format, corresponding features may encode different meanings across tasks, leading

to conflicting gradients under shared training. To mitigate this issue, we decompose each state into two complementary components: a geometric part that captures more stable cross-task structure, and a semantic part that preserves task-specific information.

Given a state $s_t \in \mathbb{R}^d$, we compute a learnable gating vector

$$m_t = \sigma(\text{MLP}(s_t)) \in (0, 1)^d, \quad (3)$$

and decompose the state as

$$s_{geo}^t = m_t \odot s_t, \quad s_{sem}^t = (1 - m_t) \odot s_t, \quad (4)$$

where \odot denotes element-wise multiplication. This decomposition is fully data-driven and does not rely on task identifiers.

To bias the geometric branch toward stable dynamical structure, we introduce a one-step self-supervised consistency objective. Specifically, we learn a predictor f_ψ such that

$$\hat{s}_{geo}^{t+1} = f_\psi(s_{geo}^t, a_t), \quad (5)$$

and minimize

$$\mathcal{L}_{geo} = \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \mathcal{D}} \left[\|f_\psi(s_{geo}^t, a_t) - s_{geo}^{t+1}\|_2^2 \right]. \quad (6)$$

This objective does not enforce exact invariance, but provides a self-supervised bias toward representations that respect shared transition structure. We further provide a simplified theoretical analysis and an empirical gradient analysis in Appendix B to illustrate how semantic inconsistency under shared training can induce cross-task optimization difficulties and why geometry–semantics decomposition helps alleviate them. In the following, s_{geo} is used to construct cross-task geometric context, while s_{sem} is reserved for task-specific decision making.

4.2 Geometry-Aware Safe Context Construction

Given a prompt τ_i^* , we first map each geometric component into a latent embedding space,

$$z_t^{geo} = \phi(s_{geo}^t), \quad (7)$$

which serves as the basis for structured context construction.

Relational structure induction. To capture dependencies among prompt elements, we learn a data-dependent relational structure end-to-end. Specifically, we compute pairwise affinities and normalize them as

$$\tilde{A}_{ij} = \frac{(W_q z_i^{geo})^\top (W_k z_j^{geo})}{\sqrt{d}}, \quad A_{ij} = \text{softmax}_j(\tilde{A}_{ij}). \quad (8)$$

This yields a fully connected but input-adaptive relational structure. Using this structure, we aggregate local information and summarize the prompt as

$$h_i = \sum_j A_{ij} (W_v z_j^{geo}), \quad h_{\text{prompt}} = \text{Pool}(\{h_i\}), \quad (9)$$

where $\text{Pool}(\cdot)$ denotes a permutation-invariant readout.

Safety-aware structure reweighting. Not all geometric relations are equally desirable under safety constraints. Rather than directly entangling cost with the representation, we treat cost-to-go as a feasibility signal that reweights relational structure:

$$w_t = \sigma(-\alpha \hat{C}_t), \quad A_{ij}^{\text{safe}} = \text{softmax}_j(w_i w_j \cdot \tilde{A}_{ij}). \quad (10)$$

We then recompute the prompt representation using the safety-aware structure,

$$h_i^{\text{safe}} = \sum_j A_{ij}^{\text{safe}} (W_v z_j^{geo}), \quad h_{\text{prompt}}^{\text{safe}} = \text{Pool}(\{h_i^{\text{safe}}\}). \quad (11)$$

In this way, safety affects *which* geometric relations are emphasized, rather than *how* geometry itself is represented.

Prototype-based global context. To complement prompt-level structure with global regularities shared across tasks, we introduce M learnable geometric prototypes $P = \{p_1, \dots, p_M\}$, where $p_m \in \mathbb{R}^d$. Given the safety-aware prompt summary, we retrieve prototype memory via

$$q = \phi_q(h_{\text{prompt}}^{\text{safe}}), \quad \beta_m = \text{softmax}(q^\top p_m), \quad h_{\text{proto}} = \sum_{m=1}^M \beta_m p_m. \quad (12)$$

The prototypes act as global anchors for recurring geometric patterns across tasks.

Final context construction. Finally, we combine the safety-aware prompt representation with the retrieved prototype memory to obtain the geometry-aware context

$$c_{\text{geo}} = \text{MLP}([h_{\text{prompt}}^{\text{safe}}, h_{\text{proto}}]). \quad (13)$$

4.3 Cost-Aware Decision Transformer

To generate actions, we combine the geometry-aware context with the semantic branch. We first encode the semantic component by

$$z_t^{\text{sem}} = \phi_{\text{sem}}(s_{\text{sem}}^t), \quad (14)$$

and fuse it with the geometry-aware context:

$$z_t = \text{MLP}([z_t^{\text{sem}}, c_{\text{geo}}]). \quad (15)$$

Following prior work on safe sequence modeling (Liu et al., 2023b; Zhang et al., 2023), we further condition action generation on both return-to-go and cost-to-go. Given a context window of length H , the transformer input at timestep t is

$$\mathbf{o}_t = \left\{ \hat{R}_{t-H:t}, \hat{C}_{t-H:t}, z_{t-H:t}, a_{t-H:t-1} \right\}, \quad (16)$$

where $\hat{R}_t = \sum_{j=t}^T r_j$ and $\hat{C}_t = \sum_{j=t}^T c_j$.

We parameterize the policy as a stochastic Gaussian distribution:

$$\pi_\theta(\cdot | \mathbf{o}_t) = \mathcal{N}(\mu_\theta(\mathbf{o}_t), \Sigma_\theta(\mathbf{o}_t)), \quad (17)$$

where μ_θ and Σ_θ are predicted by a causal transformer decoder. This design allows GeoDT to leverage shared geometry for cross-task generalization while adapting behavior to different reward and safety requirements.

4.4 Training Objective

We train GeoDT end-to-end on offline trajectories using maximum likelihood estimation with entropy regularization (Haarnoja et al., 2018). The policy loss is

$$\mathcal{L}_{\text{pol}} = \mathbb{E}_{\mathbf{o} \sim \mathcal{D}} \left[-\log \pi_\theta(\mathbf{a} | \mathbf{o}) - \lambda H(\pi_\theta(\cdot | \mathbf{o})) \right], \quad (18)$$

where $\lambda \geq 0$ controls the trade-off between conservativeness and stochasticity.

The final training objective combines the sequence modeling loss with the geometry consistency regularizer:

$$\mathcal{L} = \mathcal{L}_{\text{pol}} + \lambda_{\text{geo}} \mathcal{L}_{\text{geo}}. \quad (19)$$

During training, trajectory prompts are sampled from the offline dataset to construct context, and no additional environment interaction is required.

5 Experiment

We evaluate GeoDT on a suite of multi-task safe offline RL benchmarks to assess its ability to (i) remain robust across diverse tasks with heterogeneous semantics, (ii) maintain favorable reward–cost trade-offs under varying safety constraints, and (iii) generalize via trajectory-based context without explicit task identifiers. Implementation details are provided in Appendix C.3.

Experiment setup. Following prior work (Liu et al., 2023b; Achiam et al., 2017; Zhang et al., 2020), we consider four continuous-control safe locomotion tasks: **Run**, **Circle**, **Reach**, and **Gather**. Each task is instantiated with four robot morphologies: **Car**, **Ball**, **Drone**, and **Ant**, resulting in a diverse multi-task setting with varying dynamics, objectives, and safety constraints. All experiments are conducted on the BulletSafetyGym benchmark (Gronauer, 2022). Additional environment details and hyperparameters are provided in Appendix D.

Offline datasets. For the **Run** and **Circle** environments, we use the publicly available DSRL dataset (Liu et al., 2023a), which is specifically designed for offline safe RL. For **Reach** and **Gather**, we collect trajectories using a PPO-Lag agent (Ray et al., 2019) under progressively varying cost thresholds to ensure diverse reward–cost trade-offs. All datasets follow the DSRL format, including state, action, reward, cost, and terminal signals, ensuring consistency across tasks.

Baselines. We compare GeoDT against a set of strong baselines covering both safe offline RL and multi-task learning. All baselines are adapted to the same multi-task training protocol for fair comparison:

- *CDT* (Liu et al., 2023b): a Decision Transformer variant with cost conditioning.
- *Prompt-DT* (Xu et al., 2022b): a prompt-based Decision Transformer that uses task prompts to support multi-task generalization.
- *LSPC* (Koirala et al., 2024b): a latent-space method balancing reward and cost.
- *CPQ* (Xu et al., 2022a): a Q-learning approach with explicit constraint penalties.
- *BCQ-Lag* (Fujimoto et al., 2019): a Lagrangian extension of BCQ for safe offline RL.
- *SoftMod* (Yang et al., 2020): a modular multi-task architecture, augmented with cost-to-go inputs.
- *FISOR* (Zheng et al., 2024): a diffusion-based method with feasibility guidance.
- *O3SRL* (Chemingui et al., 2025a): an offline safe RL method with online optimization for constraint satisfaction.

Evaluation metrics. We evaluate all methods using normalized reward return and normalized cost return, following standard practice in offline RL (Fu et al., 2020; Liu et al., 2023b). The normalized reward is defined as

$$R_{\text{norm}} = \frac{R_{\pi} - r_{\min}(\mathcal{T})}{r_{\max}(\mathcal{T}) - r_{\min}(\mathcal{T})}, \quad (20)$$

and the normalized cost as

$$C_{\text{norm}} = \frac{C_{\pi}}{\kappa + \epsilon}, \quad (21)$$

where R_{π} and C_{π} denote the return and cost of policy π , and κ is the task-specific cost threshold.

To further evaluate robustness under increasing task diversity, we additionally report the *Task Scaling Robustness Score (TSRS)* and *Inter-Task Balance Score (ITBS)*, which measure performance degradation and reward–cost balance as the number of tasks increases (see Section 5.2).

5.1 Main Multi-Task Performance

Table 1 reports the task-wise performance of a *single jointly trained multi-task policy* for each method. For each agent family, all methods are trained once on the full set of tasks simultaneously, and the resulting policy is then evaluated separately on each task. This setting is substantially more challenging than training independent single-task policies, since the model must resolve cross-task interference while maintaining favorable reward–cost trade-offs under heterogeneous dynamics, objectives, and safety constraints.

Table 1: Task-wise performance of *jointly trained multi-task policies*. For each agent family, every method is trained once on all tasks simultaneously, and the resulting single policy is evaluated on each task separately. Higher reward and lower cost are better.

Method	Car multi-task setting								Ball multi-task setting								Average	
	Run		Circle		Reach		Gather		Run		Circle		Reach		Gather		Reward \uparrow	Cost \downarrow
	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost		
GeoDT	0.95	0.77	0.73	0.61	0.70	0.67	0.41	0.43	0.74	0.74	0.75	0.63	0.87	0.79	0.48	0.63	0.704	0.659
CDT	0.94	0.79	0.68	0.56	0.69	1.27	0.43	0.86	0.73	0.80	0.76	1.21	0.83	0.73	0.45	0.55	0.686	0.836
PDT	0.93	4.42	0.68	1.96	0.63	0.65	0.40	0.63	0.88	0.97	0.69	1.25	0.84	0.94	0.45	0.50	0.685	1.414
CPQ	0.90	1.18	0.59	0.84	0.49	7.10	0.14	0.29	0.37	1.71	0.58	1.05	0.62	0.87	0.24	0.68	0.491	1.715
BCQL	0.89	1.32	0.61	3.57	0.53	0.58	0.31	0.54	0.29	0.28	0.79	1.77	0.85	0.79	0.46	0.56	0.593	1.177
LSPC	0.80	0.64	0.42	1.87	0.28	0.29	0.22	0.55	0.91	1.10	0.62	2.38	0.44	1.29	0.25	0.45	0.492	1.069
SoftMod	0.95	2.68	0.70	4.43	0.60	0.62	0.37	0.55	1.14	1.17	0.82	2.52	0.82	0.94	0.43	0.50	0.729	1.675
FISOR	0.22	0.33	0.28	0.26	0.51	0.69	0.40	0.15	0.02	0.03	0.37	0.01	0.74	0.97	0.51	0.45	0.379	0.360
O3SRL	0.93	0.83	0.81	2.95	0.45	0.89	0.40	0.90	1.04	1.14	0.86	2.55	0.85	0.64	0.43	0.50	0.722	1.299

Note: Results are reported for a *single jointly trained multi-task policy* evaluated separately on each task. **Bold** denotes safe policies, gray denotes unsafe policies, and **blue bold** marks the best safe result.

Overall, GeoDT achieves the strongest average reward–cost balance among the safe methods considered. In particular, it attains the best average normalized reward (0.704) while also yielding the lowest average normalized cost (0.659) among high-performing methods. This indicates that GeoDT is able to improve task performance without sacrificing safety, rather than simply trading lower cost for substantially reduced reward. Compared with CDT and Prompt-DT, which are the strongest sequence-modeling baselines in terms of reward, GeoDT produces more consistent safety behavior and avoids the large cost violations observed in several tasks. Compared with conservative safe RL baselines such as FISOR, GeoDT achieves substantially higher reward while remaining within a competitive safety range.

A more fine-grained view reveals that GeoDT is consistently strong across individual tasks rather than relying on a few outlier wins. It achieves the best safe result on multiple tasks, including Car Run, Car Circle, Car Reach, Ball Circle, and Ball Reach, and remains competitive on the remaining tasks. In contrast, several baselines exhibit a sharper reward–cost imbalance: some methods achieve high reward at the expense of severe safety violations (e.g., O3SRL and SoftMod on several tasks), while others remain safe but become overly conservative and underperform in reward (e.g., FISOR). This pattern suggests that the main difficulty in multi-task safe offline RL is not only maximizing return or minimizing cost in isolation, but balancing the two under cross-task semantic mismatch.

These results are consistent with the central motivation of GeoDT. By separating geometry-related structure from task-specific semantics and constructing geometry-aware context before decision making, GeoDT is better able to share useful information across tasks without amplifying semantic inconsistency. The resulting benefit is not merely higher average reward, but a more favorable and stable reward–cost trade-off across heterogeneous multi-task settings.

While Table 1 demonstrates the overall effectiveness of GeoDT, average task performance alone does not fully characterize how well a method maintains performance as task diversity increases. We therefore next analyze cross-task robustness more directly using TSRS and ITBS.

5.2 Scaling Robustness under Increasing Task Diversity

While Table 1 shows that GeoDT achieves strong overall reward–cost trade-offs, average task performance alone does not fully characterize how well a method maintains performance. In multi-task safe offline RL, a method may achieve strong average performance while still degrading sharply relative to its single-task counterpart, or exhibit large variation across tasks as the number of jointly trained tasks grows. To better characterize robustness under increasing task diversity, we introduce two diagnostic metrics: the **Task Scaling Robustness Score** (TSRS) and the **Inter-Task Balance Score** (ITBS). Our design is motivated by robustness-oriented statistical measures used to summarize performance degradation and variability in other domains, such as ecological sensing (Xu, 2006) and image quality assessment (Wang et al., 2004).

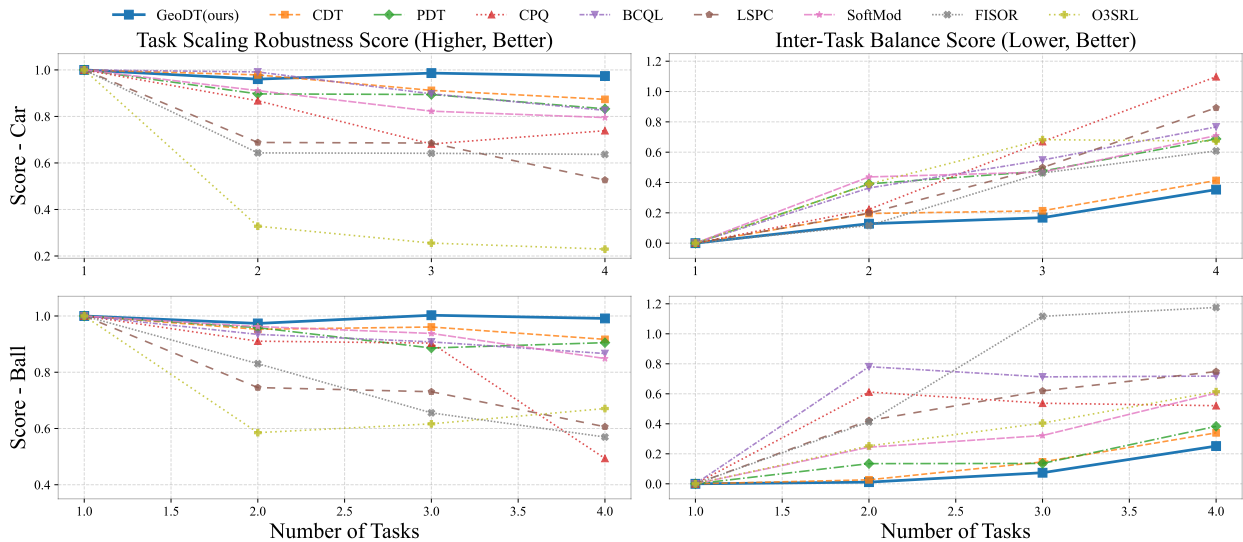


Figure 2: Scaling robustness as the number of jointly trained tasks increases. We report the Task Scaling Robustness Score (TSRS, higher is better) and the Inter-Task Balance Score (ITBS, lower is better). GeoDT consistently maintains stronger robustness and more balanced task-wise behavior than competing methods as task diversity increases.

Task Scaling Robustness Score (TSRS). TSRS measures how well a method preserves its overall reward–cost performance relative to the corresponding single-task setting as the number of jointly learned tasks increases. For an n -task setting, we define

$$\text{TSRS}(n) = 1 + \frac{1}{\sqrt{n}} \sum_{i=1}^n [(1 - \lambda)\delta R_i - \lambda\delta C_i], \quad (22)$$

where

$$\delta R_i = \frac{R_i - R_{\text{ref}}}{R_i + R_{\text{ref}} + \varepsilon}, \quad \delta C_i = \frac{C_i - C_{\text{ref}}}{C_i + C_{\text{ref}} + \varepsilon}. \quad (23)$$

Here, R_i and C_i denote the reward and cost on task i under multi-task training, while R_{ref} and C_{ref} are the corresponding single-task reference values. ε is a small constant for numerical stability, and λ controls the reward–cost trade-off; we set $\lambda = 0.5$ in all experiments. A higher TSRS indicates that the method better preserves its single-task performance as task diversity increases.

Inter-Task Balance Score (ITBS). TSRS evaluates overall performance retention, but does not indicate whether performance is evenly distributed across tasks. To quantify balance, we define

$$\text{ITBS}(n) = \text{CV}(R_1, \dots, R_n) + \lambda \cdot \text{CV}(C_1, \dots, C_n), \quad (24)$$

where CV denotes the coefficient of variation. Lower ITBS indicates more balanced behavior across tasks, with less task-to-task fluctuation in both reward and cost.

Figure 2 reports TSRS and ITBS as the number of jointly trained tasks increases. GeoDT consistently maintains a higher TSRS and a lower ITBS than competing methods, especially in the most challenging four-task setting. Although several baselines remain competitive when only two or three tasks are trained jointly, their performance degrades more sharply as task diversity increases. In contrast, GeoDT maintains performance more effectively, preserving a larger fraction of its single-task performance while maintaining more balanced reward–cost behavior across tasks.

These results are consistent with the design of GeoDT. By separating geometry-related structure from task-specific semantics and constructing safety-aware geometric context before decision making, GeoDT

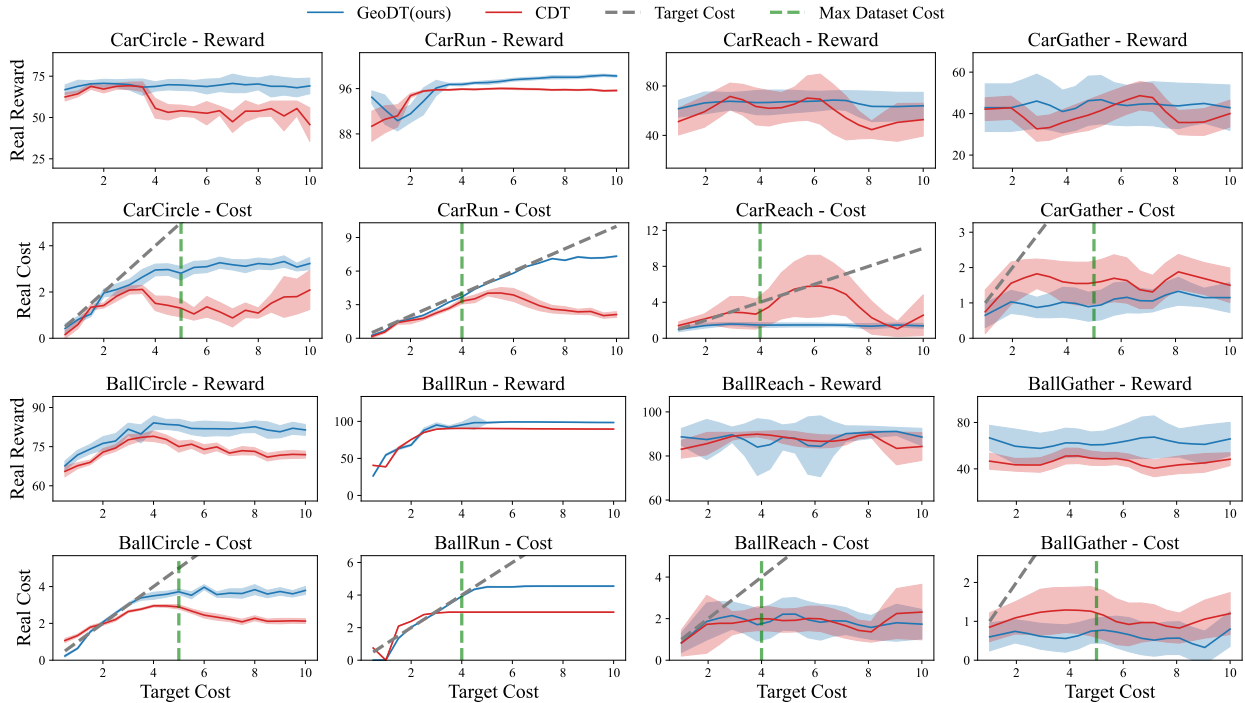


Figure 3: Zero-shot adaptation under varying target cost budgets. Each column corresponds to one task, and the x-axis denotes the target cost. Results are averaged over 20 evaluation episodes and three random seeds. Solid lines show the mean and shaded regions indicate ± 1 standard deviation. The dashed vertical line marks the maximum offline cost observed in the dataset, and the gray diagonal indicates the ideal relation $y = \text{target cost}$.

reduces cross-task interference and avoids over-specializing to a subset of tasks. As task diversity grows, this yields a more robust and balanced multi-task policy, rather than one that achieves strong performance on only a few tasks at the expense of others.

5.3 Zero-Shot Adaptation to Unseen Safety Budgets

An important advantage of return- and cost-conditioned sequence models is their ability to adapt behavior at inference time without retraining. In safe multi-task offline RL, this capability is particularly valuable: a single policy may need to operate under different safety budgets depending on deployment requirements, even when no additional interaction or task-specific fine-tuning is allowed. Following prior work on return- and cost-conditioned decision transformers (Xu et al., 2022b; Liu et al., 2023b), we therefore evaluate whether GeoDT can generalize zero-shot to *unseen target cost budgets* in the multi-task setting.

We compare GeoDT against a multi-task-trained CDT on four tasks for both the **Car** and **Ball** agents. For each task, we sweep 20 target cost budgets at evaluation time. In addition, obstacle layouts and goal positions are randomly regenerated to produce unseen environment configurations, making the evaluation more challenging than simply replaying the offline training distribution. The results are shown in Figure 3.

Overall, GeoDT achieves a more reliable reward–cost trade-off across a wide range of target budgets. As the target cost increases, GeoDT typically adjusts its behavior smoothly, improving reward when additional cost budget is available while remaining substantially better aligned with the desired cost level. This trend is especially clear in tasks such as **Run** and **Circle**, where reward is strongly coupled with exploratory or aggressive behavior. In these settings, GeoDT adjusts performance more smoothly with the budget, whereas CDT often becomes unstable or violates the intended safety level under tighter constraints.

Table 2: Ablation of key components in GeoDT. We report average normalized reward and cost over all Car tasks and all Ball tasks, respectively, as well as the overall average. Higher reward and lower cost are better.

Method	Car Avg		Ball Avg		Overall Avg	
	Reward \uparrow	Cost \downarrow	Reward \uparrow	Cost \downarrow	Reward \uparrow	Cost \downarrow
GeoDT	0.697	0.620	0.710	0.698	0.704	0.659
GeoDT (w/o feasibility weighting)	0.613	0.922	0.665	0.814	0.639	0.868
GeoDT (w/o prototype)	0.610	0.783	0.660	0.691	0.635	0.737
GeoDT (w/o decomposition)	0.678	0.710	0.675	0.712	0.676	0.711
GeoDT (w/o \mathcal{L}_{geo})	0.665	0.773	0.635	0.763	0.659	0.768

Note: All variants are trained under the same multi-task protocol. Car Avg and Ball Avg denote averages over the four corresponding tasks for each agent family.

GeoDT is also more robust when extrapolating beyond the effective training range. Even when target costs approach or exceed the maximum offline cost observed in the dataset, GeoDT maintains comparatively stable behavior and avoids the sharp degradation seen in CDT on several tasks. This suggests that the geometry-aware context learned by GeoDT provides a more transferable representation of safe behavior patterns, allowing the model to adjust to new safety requirements without retraining.

These results complement the multi-task performance and scaling analyses above. Beyond achieving strong average reward–cost trade-offs and robust performance under increasing task diversity, GeoDT also supports zero-shot adaptation to unseen safety budgets, which is essential for practical deployment in safety-critical multi-task settings.

5.4 Ablation Study

We perform ablations to evaluate the contribution of the main components in GeoDT, including geometry–semantics decomposition, prototype memory, and feasibility-aware structure reweighting. Table 2 reports the performance of a single jointly trained multi-task policy for each ablated variant.

Overall, Table 2 shows that all major components of GeoDT contribute to the final performance. Removing feasibility-aware reweighting causes the largest increase in average cost, especially on Car tasks, indicating that safety-aware structure selection is essential for maintaining cost compliance. Removing the prototype memory also degrades both reward and cost, showing that global geometric memory provides useful structural information beyond prompt-level relations. Removing the geometry–semantics decomposition leads to a smaller but still consistent drop in the overall reward–cost trade-off, supporting our claim that disentangling shared geometry from task-specific semantics helps reduce cross-task interference. Notably, removing the geometry consistency loss \mathcal{L}_{geo} leads to an even larger degradation than removing decomposition itself. This suggests that decomposition alone is insufficient: without an explicit geometric consistency bias, the learned partition between geometric and semantic features becomes weakly constrained and may even be less effective than not decomposing the state at all. Together, these results indicate that the benefit of GeoDT comes from the interaction between decomposition and geometry consistency, rather than from either design choice in isolation.

We further study sensitivity to the number of geometric prototypes M on two representative agents, **Car** and **Ant**. As shown in Figure 4, GeoDT remains relatively stable across $M \in \{4, 8, 16\}$, indicating that the prototype memory is not overly sensitive to the exact choice of M . For **Car**, the average reward varies only slightly across different prototype counts, while the cost remains within a relatively narrow range, suggesting that the learned global geometric memory is robust once a moderate capacity is provided. A similar trend is observed for **Ant**, where different values of M lead to comparable reward–cost trade-offs without dramatic degradation. Overall, these results suggest that the gains of GeoDT do not rely on careful tuning of the prototype count; instead, a moderate number of prototypes is sufficient to capture useful cross-task geometric regularities. We therefore use $M = 8$ in the main experiments as a balanced default choice.

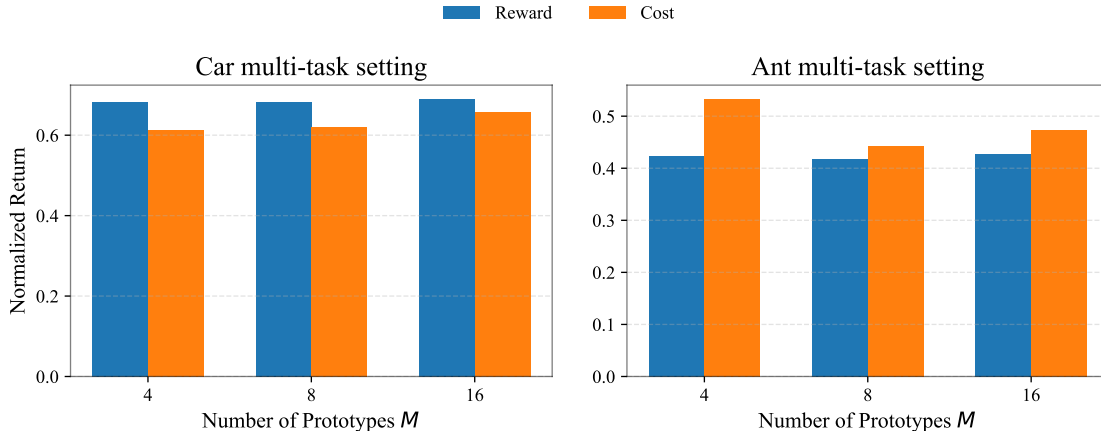


Figure 4: Sensitivity to the number of geometric prototypes M on representative agents. We report the average normalized reward and cost over all tasks for the Car and Ant multi-task settings. GeoDT remains stable across a moderate range of M , indicating that the prototype memory is not overly sensitive to this hyperparameter.

Additional experiments on two more agents, as well as more ablation studies and key hyperparameters, are presented in the Appendix E.

6 Conclusion

In this work, we presented **GeoDT**, a geometry-aware framework for robust, safe multi-task offline reinforcement learning. GeoDT addresses a central challenge in multi-task offline RL—*semantic inconsistency* across heterogeneous tasks—by separating geometry-related trajectory structure from task-specific semantics and using the resulting geometry-aware context to guide decision making. In particular, GeoDT combines learnable geometry–semantics decomposition, self-supervised geometry consistency learning, prototype-augmented context construction, and safety-aware structure reweighting within a cost-conditioned Decision Transformer.

Empirically, GeoDT achieves strong reward–cost trade-offs across diverse safe multi-task benchmarks and remains more robust as task diversity increases, and generalizes zero-shot to unseen safety budgets without requiring explicit task identifiers. These results suggest that exploiting shared trajectory geometry is a promising way to improve robustness and safety in offline multi-task decision making.

Limitations and future work. Despite these promising results, several limitations remain. First, GeoDT is still evaluated within task families seen during training; although it generalizes across heterogeneous tasks and unseen safety budgets, transferring to entirely novel task families may still require additional adaptation. Second, our method improves safety through cost-aware structure learning and cost-conditioned decision making, but does not provide hard safety guarantees during deployment. Finally, while the prototype memory is empirically stable, its interpretation remains implicit and could be further studied.

Future work will explore stronger forms of safety assurance, including integrating offline safe RL with verifiable control tools such as reachability-based methods (Zheng et al., 2024). Another promising direction is to extend geometry-aware representation learning to broader cross-domain and cross-embodiment settings, where task semantics vary even more substantially across environments.

References

Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In *International Conference on Machine Learning*, pp. 22–31. PMLR, 2017.

- Eitan Altman. *Constrained Markov decision processes*. Routledge, 2021.
- Marc Bellemare, Will Dabney, Robert Dadashi, Adrien Ali Taiga, Pablo Samuel Castro, Nicolas Le Roux, Dale Schuurmans, Tor Lattimore, and Clare Lyle. A geometric perspective on optimal representations for reinforcement learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yassine Chemingui, Aryan Deshwal, Alan Fern, Thanh Nguyen-Tang, and Janardhan Rao Doppa. Online optimization for offline safe reinforcement learning. In *Advances in Neural Information Processing Systems*, 2025a.
- Yassine Chemingui, Aryan Deshwal, Honghao Wei, Alan Fern, and Jana Doppa. Constraint-adaptive policy switching for offline safe reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 15722–15730, 2025b.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in Neural Information Processing Systems*, 34:15084–15097, 2021.
- Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In *International Conference on Machine Learning*, pp. 2052–2062. PMLR, 2019.
- Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G Bellemare. Deepmdp: Learning continuous latent space models for representation learning. In *International Conference on Machine Learning*, pp. 2170–2179. PMLR, 2019.
- Ze Gong, Akshat Kumar, and Pradeep Varakantham. Offline safe reinforcement learning using trajectory classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 16880–16887, 2025.
- Sven Gronauer. Bullet-safety-gym: A framework for constrained reinforcement learning. Technical report, mediaTUM, 2022.
- 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. PMLR, 2018.
- Tai Hoang, Huy Le, Philipp Becker, Vien Anh Ngo, and Gerhard Neumann. Geometry-aware RL for manipulation of varying shapes and deformable objects. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=7BLXhmWvwF>.
- Shengchao Hu, Ziqing Fan, Li Shen, Ya Zhang, Yanfeng Wang, and Dacheng Tao. Harmodt: Harmony multi-task decision transformer for offline reinforcement learning. *arXiv preprint arXiv:2405.18080*, 2024.
- Sili Huang, Jifeng Hu, Hechang Chen, Lichao Sun, and Bo Yang. In-context decision transformer: Reinforcement learning via hierarchical chain-of-thought. *arXiv preprint arXiv:2405.20692*, 2024.
- Haque Ishfaq, Thanh Nguyen-Tang, Songtao Feng, Raman Arora, Mengdi Wang, Ming Yin, and Doina Precup. Offline multitask representation learning for reinforcement learning. *Advances in Neural Information Processing Systems*, 37:70557–70616, 2024.
- Khimya Khetarpal, Zhaohan Daniel Guo, Bernardo Avila Pires, Yunhao Tang, Clare Lyle, Mark Rowland, Nicolas Heess, Diana Borsa, Arthur Guez, and Will Dabney. A unifying framework for action-conditional self-predictive reinforcement learning. *arXiv preprint arXiv:2406.02035*, 2024.
- Prajwal Koirala, Zhanhong Jiang, Soumik Sarkar, and Cody Fleming. Fawac: Feasibility informed advantage weighted regression for persistent safety in offline reinforcement learning. *arXiv preprint arXiv:2412.08880*, 2024a.

- Prajwal Koirala, Zhanhong Jiang, Soumik Sarkar, and Cody Fleming. Latent safety-constrained policy approach for safe offline reinforcement learning. *arXiv preprint arXiv:2412.08794*, 2024b.
- Michael Laskin, Luyu Wang, Junhyuk Oh, Emilio Parisotto, Stephen Spencer, Richie Steigerwald, DJ Strouse, Steven Hansen, Angelos Filos, Ethan Brooks, et al. In-context reinforcement learning with algorithm distillation. *arXiv preprint arXiv:2210.14215*, 2022.
- Hoang Le, Cameron Voloshin, and Yisong Yue. Batch policy learning under constraints. In *International Conference on Machine Learning*, pp. 3703–3712. PMLR, 2019.
- Jonathan Lee, Annie Xie, Aldo Pacchiano, Yash Chandak, Chelsea Finn, Ofir Nachum, and Emma Brunskill. Supervised pretraining can learn in-context reinforcement learning. *Advances in Neural Information Processing Systems*, 36:43057–43083, 2023.
- Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. Coptidice: Offline constrained reinforcement learning via stationary distribution correction estimation. *arXiv preprint arXiv:2204.08957*, 2022.
- Qian Lin, Bo Tang, Zifan Wu, Chao Yu, Shangqin Mao, Qianlong Xie, Xingxing Wang, and Dong Wang. Safe offline reinforcement learning with real-time budget constraints. In *International Conference on Machine Learning*, pp. 21127–21152. PMLR, 2023.
- Zuxin Liu, Zijian Guo, Haohong Lin, Yihang Yao, Jiacheng Zhu, Zhepeng Cen, Hanjiang Hu, Wenhao Yu, Tingnan Zhang, Jie Tan, et al. Datasets and benchmarks for offline safe reinforcement learning. *arXiv preprint arXiv:2306.09303*, 2023a.
- Zuxin Liu, Zijian Guo, Yihang Yao, Zhepeng Cen, Wenhao Yu, Tingnan Zhang, and Ding Zhao. Constrained decision transformer for offline safe reinforcement learning. In *International Conference on Machine Learning*, pp. 21611–21630. PMLR, 2023b.
- Zuxin Liu, Zijian Guo, Haohong Lin, Yihang Yao, Jiacheng Zhu, Zhepeng Cen, Hanjiang Hu, Wenhao Yu, Tingnan Zhang, Jie Tan, and Ding Zhao. Datasets and benchmarks for offline safe reinforcement learning. *Journal of Data-centric Machine Learning Research*, 2024.
- Nicholas Polosky, Bruno C Da Silva, Madalina Fiterau, and Jithin Jagannath. Constrained offline policy optimization. In *International Conference on Machine Learning*, pp. 17801–17810. PMLR, 2022.
- Alex Ray, Joshua Achiam, and Dario Amodei. Benchmarking safe exploration in deep reinforcement learning. *arXiv preprint arXiv:1910.01708*, 7(1):2, 2019.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- Max Schwarzer, Ankesh Anand, Rishab Goel, R Devon Hjelm, Aaron Courville, and Philip Bachman. Data-efficient reinforcement learning with self-predictive representations. *arXiv preprint arXiv:2007.05929*, 2020.
- Zhe Wang, Haozhu Wang, and Yanjun Qi. Hierarchical prompt decision transformer: Improving few-shot policy generalization with global and adaptive guidance. In *Companion Proceedings of the ACM on Web Conference 2025*, pp. 520–529, 2025.
- Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.
- Jingda Wu, Chao Huang, Hailong Huang, Chen Lv, Yuntong Wang, and Fei-Yue Wang. Recent advances in reinforcement learning-based autonomous driving behavior planning: A survey. *Transportation Research Part C: Emerging Technologies*, 164:104654, 2024.
- Hanqiu Xu. Modification of normalised difference water index (ndwi) to enhance open water features in remotely sensed imagery. *International Journal of Remote Sensing*, 27(14):3025–3033, 2006.

- Haoran Xu, Xianyuan Zhan, and Xiangyu Zhu. Constraints penalized q-learning for safe offline reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 8753–8760, 2022a.
- Mengdi Xu, Yikang Shen, Shun Zhang, Yuchen Lu, Ding Zhao, Joshua Tenenbaum, and Chuang Gan. Prompting decision transformer for few-shot policy generalization. In *International Conference on Machine Learning*, pp. 24631–24645. PMLR, 2022b.
- Ruihan Yang, Huazhe Xu, Yi Wu, and Xiaolong Wang. Multi-task reinforcement learning with soft modularization. *Advances in Neural Information Processing Systems*, 33:4767–4777, 2020.
- Yihang Yao, Zhepeng Cen, Wenhao Ding, Haohong Lin, Shiqi Liu, Tingnan Zhang, Wenhao Yu, and Ding Zhao. Oasis: Conditional distribution shaping for offline safe reinforcement learning. *Advances in Neural Information Processing Systems*, 37:78451–78478, 2024.
- Minjong Yoo, Sangwoo Cho, and Honguk Woo. Skills regularized task decomposition for multi-task offline reinforcement learning. *Advances in Neural Information Processing Systems*, 35:37432–37444, 2022.
- Minjong Yoo, Woo Kyung Kim, and Honguk Woo. In-context policy adaptation via cross-domain skill diffusion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 22191–22199, 2025.
- Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Sergey Levine, and Chelsea Finn. Conservative data sharing for multi-task offline reinforcement learning. *Advances in Neural Information Processing Systems*, 34:11501–11516, 2021.
- Qin Zhang, Linrui Zhang, Haoran Xu, Li Shen, Bowen Wang, Yongzhe Chang, Xueqian Wang, Bo Yuan, and Dacheng Tao. Saformer: A conditional sequence modeling approach to offline safe reinforcement learning. *arXiv preprint arXiv:2301.12203*, 2023.
- Yiming Zhang, Quan Vuong, and Keith Ross. First order constrained optimization in policy space. *Advances in Neural Information Processing Systems*, 33:15338–15349, 2020.
- Yinan Zheng, Jianxiong Li, Dongjie Yu, Yujie Yang, Shengbo Eben Li, Xianyuan Zhan, and Jingjing Liu. Safe offline reinforcement learning with feasibility-guided diffusion model. *arXiv preprint arXiv:2401.10700*, 2024.

A Discussion on State Structure and Semantic Inconsistency

In the safe locomotion benchmarks considered in this work, the observation vector typically contains two qualitatively different types of information: (1) *agent-centric dynamical variables*, which describe the agent’s internal physical state such as position, velocity, orientation, and actuator status, and (2) *task-dependent environmental signals*, which encode goals, obstacles, boundaries, or other task-specific cues perceived through external sensors.

The first type is generally more stable across tasks for a fixed embodiment. For example, in BulletSafetyGym, the agent-centric part includes variables such as position, velocity, orientation, angular velocity, and motor-related states. These quantities preserve consistent physical semantics across tasks. In contrast, the second type depends strongly on the task definition and observation design, and may vary substantially in both dimensionality and meaning.

For instance, the agent-centric observations are structured as follows:

- **Car:** a 7-dimensional vector including position (x, y) , velocity (v_x, v_y) , yaw (expressed as sine and cosine), and yaw rate.
- **Ball:** a 7-dimensional vector including position (x, y) , velocity (v_x, v_y) , and 3-dimensional angular velocity.

- **Drone**: a 17-dimensional vector including position (x, y, z) , velocity (v_x, v_y, v_z) , a 4-dimensional orientation quaternion, 3-dimensional angular velocity, and motor-related measurements.
- **Ant**: a 33-dimensional vector including position (x, y, z) , velocity (v_x, v_y, v_z) , a 4-dimensional orientation quaternion, 3-dimensional angular velocity, 16-dimensional motor relative positions, and 4-dimensional leg contact states.

Task-dependent observations are defined as:

- **Run**: no explicit task-specific observations, since the task only requires forward motion under a speed constraint.
- **Circle**: a 1-dimensional value representing the distance to the circle boundary.
- **Reach**: for **Car**, a 26-dimensional vector combining a 2-dimensional goal distance and a 24-dimensional LIDAR scan; for **Ball**, **Drone**, and **Ant**, a 50-dimensional vector combining goal distance with two 24-dimensional LIDAR scans.
- **Gather**: a 32-dimensional vector from two 16-ray LIDAR scans, one detecting goals and the other detecting obstacles.

These examples illustrate why semantic inconsistency naturally arises in multi-task learning. Even when observation vectors have similar dimensionality across tasks, features handled by a shared representation may correspond to different physical quantities or serve different functional roles. For example, one dimension may reflect proximity to a goal in one task but proximity to an obstacle or boundary in another. As a result, naively sharing representations across tasks can introduce ambiguity and conflicting optimization signals.

This observation motivates the design of GeoDT. Rather than manually fixing which dimensions should be shared across tasks, GeoDT learns to separate more stable geometry-related structure from task-dependent semantics in a data-driven way. The state structure described above therefore serves as intuition for why such a decomposition is beneficial, rather than as a hand-crafted rule used directly by the method.

B A Simplified Analysis of Cross-Task Gradient Conflict

B.1 Simplified Theoretical Analysis

This appendix provides a simplified theoretical analysis to illustrate why semantic inconsistency under shared training can induce cross-task gradient conflict, and why separating geometry-related structure from task-specific semantics can help alleviate this effect. The goal is not to fully characterize the optimization dynamics of GeoDT, but to formalize the intuition behind the decomposition used in Section 4.1.

Setup. Consider two tasks i and j trained with a shared representation parameterized by θ . We assume that the input to task i can be decomposed as

$$x_i = z + u_i, \quad x_j = z + u_j,$$

where z denotes a geometry- or dynamics-related component that is shared across tasks, while u_i and u_j denote task-specific semantic components. Intuitively, z captures the part of the input that is useful for cross-task sharing, whereas u_i and u_j represent semantic variation that may differ across tasks even when the inputs are processed by a common encoder.

Let the task losses be $\mathcal{L}_i(\theta; x_i)$ and $\mathcal{L}_j(\theta; x_j)$, and let the corresponding gradients on the shared parameters be

$$g_i = \nabla_{\theta} \mathcal{L}_i(\theta; x_i), \quad g_j = \nabla_{\theta} \mathcal{L}_j(\theta; x_j).$$

To quantify cross-task gradient conflict, we consider the expected gradient alignment

$$A_{ij} := \mathbb{E}[\langle g_i, g_j \rangle].$$

Larger values of A_{ij} indicate better gradient agreement, whereas smaller or negative values indicate stronger conflict.

Linearized shared training model. Suppose the shared encoder directly processes the full input $x_i = z + u_i$. Linearizing the task gradients around the shared component z , we write

$$g_i \approx \bar{g}_i + J_i u_i, \quad g_j \approx \bar{g}_j + J_j u_j,$$

where

$$\bar{g}_i := \nabla_{\theta} \mathcal{L}_i(\theta; z), \quad \bar{g}_j := \nabla_{\theta} \mathcal{L}_j(\theta; z),$$

and J_i, J_j denote the sensitivity of the gradients to task-specific semantic perturbations.

Expanding the inner product yields

$$\langle g_i, g_j \rangle \approx \langle \bar{g}_i, \bar{g}_j \rangle + \langle \bar{g}_i, J_j u_j \rangle + \langle J_i u_i, \bar{g}_j \rangle + \langle J_i u_i, J_j u_j \rangle. \quad (25)$$

The first term reflects alignment induced by the shared geometry-related structure z , while the remaining terms arise from task-specific semantic variation.

Assumptions. For simplicity, assume that:

1. $\mathbb{E}[u_i] = \mathbb{E}[u_j] = 0$;
2. the perturbations are independent of the shared component in the linearized regime;
3. $\mathbb{E}[u_i u_j^{\top}]$ is small when the task-specific semantics differ substantially across tasks.

Under these assumptions, the cross terms involving one perturbation and one shared term vanish in expectation or remain small, so the dominant perturbation effect comes from the last term in Eq. equation 25.

Proposition 1. Under the linearized model above, the expected gradient alignment can be decomposed as

$$A_{ij} \approx \mathbb{E}[\langle \bar{g}_i, \bar{g}_j \rangle] + \mathbb{E}[\langle J_i u_i, J_j u_j \rangle].$$

Moreover,

$$|\mathbb{E}[\langle J_i u_i, J_j u_j \rangle]| \leq \|J_i\| \|J_j\| \mathbb{E}[\|u_i\| \|u_j\|].$$

Interpretation. Proposition 1 shows that naive shared encoding mixes two effects: a shared-structure term, which promotes gradient agreement, and a task-specific perturbation term, which scales with both the magnitude of semantic mismatch and the sensitivity of the shared encoder to that mismatch. As task-specific semantics become more heterogeneous, the perturbation term can dominate and reduce gradient alignment, thereby inducing stronger optimization conflict across tasks.

Proof sketch. The decomposition follows by taking expectation in Eq. equation 25 and using the assumptions above to suppress the mixed terms. The bound is obtained by Cauchy–Schwarz:

$$|\langle J_i u_i, J_j u_j \rangle| \leq \|J_i u_i\| \|J_j u_j\| \leq \|J_i\| \|J_j\| \|u_i\| \|u_j\|,$$

followed by expectation on both sides.

Effect of geometry–semantics decomposition. Now consider a decomposition

$$x_i \mapsto (x_i^{\text{geo}}, x_i^{\text{sem}}),$$

where the shared branch is encouraged to depend primarily on the geometry-related component z , while task-specific semantic variation is routed to the semantic branch. In the decomposed model, let the shared gradients be

$$\tilde{g}_i \approx \bar{g}_i + \tilde{J}_i u_i, \quad \tilde{g}_j \approx \bar{g}_j + \tilde{J}_j u_j,$$

where \tilde{J}_i and \tilde{J}_j measure the residual sensitivity of the shared branch to task-specific semantic perturbations after decomposition.

Table 3: Empirical gradient analysis on the shared geometry-aware branch during joint training. We report the average pairwise gradient cosine similarity, the fraction of task pairs with negative cosine similarity, and the average gradient norm. Higher cosine similarity and gradient norm are better, while lower negative ratio is better.

Method	Car			Ball		
	Cosine \uparrow	Neg. Ratio \downarrow	Grad Norm \uparrow	Cosine \uparrow	Neg. Ratio \downarrow	Grad Norm \uparrow
GeoDT	0.0420	0.4113	5.386	0.0308	0.3300	2.526
w/o decomposition	-0.0300	0.7500	0.079	0.0013	0.5000	0.145

Proposition 2. If the decomposition reduces the sensitivity of the shared branch to task-specific semantic variation, i.e.,

$$\|\tilde{J}_i\| \leq \|J_i\|, \quad \|\tilde{J}_j\| \leq \|J_j\|,$$

then the perturbation component of the expected gradient alignment satisfies

$$|\mathbb{E}[\langle \tilde{J}_i u_i, \tilde{J}_j u_j \rangle]| \leq \|\tilde{J}_i\| \|\tilde{J}_j\| \mathbb{E}[\|u_i\| \|u_j\|] \leq \|J_i\| \|J_j\| \mathbb{E}[\|u_i\| \|u_j\|].$$

Therefore, geometry–semantics decomposition reduces the magnitude of the task-specific perturbation term in shared gradients and improves expected gradient alignment relative to naive shared encoding.

Interpretation. Proposition 2 formalizes the role of decomposition: it does not need to eliminate task-specific semantic influence entirely. It is enough to reduce the sensitivity of the shared branch to semantic mismatch, thereby shrinking the perturbation term that degrades cross-task gradient agreement.

Role of the geometry consistency objective. In GeoDT, the geometry consistency loss further biases the shared branch toward transition-stable structure. In the simplified model above, this can be interpreted as further reducing the sensitivity of the shared branch to task-specific semantic perturbations, i.e., decreasing $\|\tilde{J}_i\|$ and $\|\tilde{J}_j\|$.

Corollary 1. If the geometry consistency objective further reduces the sensitivity of the shared branch to task-specific semantic perturbations, then it tightens the perturbation bound in Proposition 2 and further improves expected gradient alignment across tasks.

Discussion. This analysis suggests that semantic inconsistency hurts multi-task training by injecting task-specific perturbations into gradients on shared parameters. Geometry–semantics decomposition helps by preserving a cleaner shared optimization signal, while the geometry consistency objective further suppresses task-specific sensitivity in the shared branch. Although simplified, this view provides a simplified optimization-based explanation for why geometry–semantics decomposition can reduce cross-task interference under heterogeneous shared training.

B.2 Empirical Gradient Analysis

To complement the simplified analysis above, we further examine optimization behavior on the shared geometry-aware branch during joint training. Specifically, we compare GeoDT with the variant without decomposition on the 4-task Car and Ball settings, and compute task-wise gradients on the shared geometry-aware parameters. This analysis is intended to test whether geometry–semantics decomposition alleviates the optimization difficulties caused by semantic inconsistency under naive shared training.

We report three quantities. First, we measure the average pairwise gradient cosine similarity across tasks, where higher values indicate better gradient alignment on the shared branch. Second, we report the fraction of task pairs with negative cosine similarity, where lower values indicate fewer conflicting task updates. Third, we report the average gradient norm on the shared geometry-aware branch, which reflects the strength of the effective optimization signal received by the shared parameters. Specifically, gradients are measured on the

decomposition, geometry encoder, relational structure, and prototype-context parameters, and the reported values are averaged over the last 10 logging points and 3 seeds.

As shown in Table 3, GeoDT achieves higher gradient cosine similarity, a lower fraction of conflicting task pairs, and substantially larger shared-branch gradient norms than the variant without decomposition on both Car and Ball. In particular, the negative-pair ratio is reduced by 44.9% on Car and 34.0% on Ball. These empirical observations are consistent with the simplified analysis above: semantic inconsistency under naive shared training can induce more conflicting task updates and weaken the effective optimization signal on the shared branch, whereas geometry–semantics decomposition helps preserve more stable and informative shared updates.

C Implementation Details

C.1 Data Collection Algorithm

We construct our dataset using the PPO-Lag algorithm (Ray et al., 2019), an extension of the PPO method to constrained reinforcement learning.

The learning problem is formulated as a Constrained Markov Decision Process (CMDP), where the objective is to maximize expected returns while satisfying predefined safety constraints. To incorporate the constraint, a non-negative Lagrange multiplier λ is introduced, yielding the following unconstrained Lagrangian objective:

$$\mathcal{L}(\theta, \lambda) = J_R(\theta) - \lambda (J_C(\theta) - d). \quad (26)$$

where d is the cost budget. PPO-Lag performs on-policy updates by optimizing a clipped surrogate objective with a Lagrangian-modified advantage function:

$$L^{\text{PPO-Lag}}(\theta) = \mathbb{E}_t \left[\min(r_t(\theta) \hat{A}_t^{\text{Lag}}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^{\text{Lag}}) \right], \quad (27)$$

$$\hat{A}_t^{\text{Lag}} = \hat{A}_{R,t} - \lambda \hat{A}_{C,t}.$$

Here, $\hat{A}_{R,t}$ and $\hat{A}_{C,t}$ denote empirical estimates of the advantage functions with respect to the reward and cost signals, respectively. These quantify the relative benefit of an action compared to the policy’s average behavior under each signal, and are typically estimated using a value baseline or generalized advantage estimation (GAE) (Schulman et al., 2015).

The Lagrange multiplier is updated via projected gradient ascent:

$$\lambda \leftarrow [\lambda + \alpha_\lambda (J_C(\theta) - d)]_+ \quad (28)$$

where $[\cdot]_+$ denotes projection onto the non-negative orthant. The code implementation is adapted from the FSRL (Liu et al., 2024) framework.

C.2 Dataset Specifications

The statistics of the dataset for each task are summarized in Table 4.

Table 4: Dataset statistics across domains and tasks.

Domain	Task	Trajectories	Average Reward	Average Cost
Car	Circle	1450	369.9 ± 141.9	49.8 ± 30.5
	Run	651	508.3 ± 66.6	14.2 ± 11.6
	Reach	1404	57.1 ± 17.5	9.0 ± 8.7
	Gather	1260	30.3 ± 9.4	1.6 ± 1.2
Ball	Circle	886	591.6 ± 220.9	43.1 ± 22.2
	Run	940	618.6 ± 270.1	48.7 ± 23.1
	Reach	1456	381.1 ± 48.8	16.2 ± 11.6
	Gather	1253	29.6 ± 8.4	1.1 ± 1.0
Drone	Circle	1923	768.8 ± 149.0	53.8 ± 32.9
	Run	1990	396.4 ± 101.3	50.5 ± 43.2
	Reach	936	119.3 ± 93.9	15.6 ± 11.7
	Gather	751	12.6 ± 5.0	0.5 ± 0.7
Ant	Circle	5728	239.0 ± 108.3	85.3 ± 54.9
	Run	1816	601.6 ± 170.6	68.1 ± 47.1
	Reach	847	10.4 ± 12.8	9.3 ± 12.1
	Gather	656	12.7 ± 4.4	0.4 ± 0.6

Note: Values are reported as mean \pm standard deviation.

The offline datasets exhibit substantial return variability under the imposed collection cost limits. For instance, the standard deviation of episodic returns reaches ± 141.9 for Car–Circle, ± 270.1 for Ball–Run, and ± 149.0 for Drone–Circle, indicating that policy performance fluctuates widely as it strives to satisfy the cost constraint. By contrast, the simpler Reach and Gather tasks yield markedly lower cost dispersion—approximately ± 8 – 12 for Reach and only ± 0.6 – 1.2 for Gather, reflecting their reduced complexity and hence diminished sensitivity to the cost limit, which results in more stable average performance.

C.3 Training Details

We provide additional implementation details for GeoDT in this appendix. The implementation and the associated datasets are available at <https://anonymous.4open.science/r/hdt-D512>. For each dataset, 10% of the trajectories are randomly reserved for evaluation, while the remaining 90% are used for training. During both training and evaluation, prompts are sampled from the corresponding split. All experiments are conducted on eight NVIDIA RTX 6000 Ada Generation GPUs, each with 48 GB of memory. Unless otherwise stated, all reported results are averaged over 3 random seeds, and each evaluation is performed over 20 episodes per seed.

GeoDT is trained end-to-end without a separate pretraining stage. In each iteration, a trajectory segment is sampled for policy learning together with a short prompt for geometry-aware context construction. Algorithm 1 summarizes how the prompt is converted into geometry-aware context, and Algorithm 2 describes the overall training procedure.

Algorithm 1 Geometry-Aware Context Construction

Require: Prompt trajectory $\tau^* = (\hat{R}^*, \hat{C}^*, s^*, a^*)$, prototype memory $P = \{p_1, \dots, p_M\}$ **Ensure:** Geometry-aware context c_{geo}

- 1: Decompose prompt states into geometric and semantic parts:

$$(s_{geo}, s_{sem}) \leftarrow \text{DECOMPOSE}(s^*)$$

- 2: Encode geometric prompt tokens:

$$z^{geo} \leftarrow \phi(s_{geo})$$

- 3: Induce prompt-level relational structure and summarize local geometry:

$$h_{\text{prompt}} \leftarrow \text{STRUCTUREINDUCE}(z^{geo})$$

- 4: Reweight relations using cost-to-go feasibility signals:

$$h_{\text{prompt}}^{\text{safe}} \leftarrow \text{FEASIBLEREWIGHT}(h_{\text{prompt}}, \hat{C}^*)$$

- 5: Retrieve global prototype memory:

$$h_{\text{proto}} \leftarrow \text{PROTOTYPERETRIEVE}(h_{\text{prompt}}^{\text{safe}}, P)$$

- 6: Fuse local and global geometry:

$$c_{geo} \leftarrow \text{MLP}([h_{\text{prompt}}^{\text{safe}}, h_{\text{proto}}])$$

- 7:
- return**
- c_{geo}
-

Algorithm 2 GeoDT Training Procedure

Require: Offline multi-task dataset $\mathcal{D} = \{\mathcal{T}_1, \dots, \mathcal{T}_N\}$, training iterations L , prompt length K , trajectory length H

- 1:
- for**
- $\ell = 1$
- to
- L
- do**

- 2: Sample trajectory segments
- τ
- and prompts
- τ^*
- from each task

- 3:
- for all**
- (τ, τ^*)
- in the minibatch
- do**

- 4: Construct geometry-aware context:

$$c_{geo} \leftarrow \text{GEOCONTEXT}(\tau^*)$$

- 5: Fuse
- c_{geo}
- with semantic state features in
- τ

- 6: Predict action distribution with the cost-aware Decision Transformer

- 7: Compute geometry consistency loss
- \mathcal{L}_{geo}

- 8:
- end for**

- 9: Compute sequence-modeling policy loss
- \mathcal{L}_{DT}

- 10: Form the total objective:

$$\mathcal{L} = \mathcal{L}_{\text{DT}} + \lambda_{geo} \mathcal{L}_{geo}$$

- 11: Update model parameters by gradient descent

- 12:
- end for**
-

When the optional prototype balancing regularizer is enabled in the appendix study, we further add $\lambda_{\text{bal}} \mathcal{L}_{\text{bal}}$ to the total objective.

D Experiment Setting and Hyperparameters

D.1 Environment Task Description

We use BulletSafetyGym (Gronauer, 2022) environments for training and evaluation. In the **Circle** tasks, the agent is expected to move on a circle in a clockwise direction and is rewarded based on its velocity and proximity to the boundary of the circle. The reward and cost are defined as:

$$\begin{aligned} r(s_t) &= 0.1 \times \frac{-y_t v_x + x_t v_y}{1 + \|x_t\|_2 - r} + 0.01 \times r_{\text{agent}}(s_t) \\ c(s_t) &= \mathbf{1}(\|x_t\|_2 > x_{\text{lim}}) \end{aligned} \quad (29)$$

where r is the radius of the circle, x_{lim} is the safety region, and $r_{\text{agent}}(\cdot)$ represents additional rewards based on the agent state, such as electricity costs, speed costs, joint cost, etc.

In the **Run** tasks, the agent is expected to run through an avenue between two safety boundaries. The reward and cost are defined as:

$$\begin{aligned} r(s_t) &= (\phi_t - \phi_{t-1}) + r_{\text{agent}}(s_t) \\ c(s_t) &= \min(1, \mathbf{1}(\|y_t\|_2 > y_{\text{lim}}) + \mathbf{1}(\|v_t\|_2 > v_{\text{lim}})) \end{aligned} \quad (30)$$

where $\phi_t = -60 \times \|(x_t, y_t) - (g_x, 0)\|_2$ is the task potential to a fictitious target position $(g_x, 0)$, y_{lim} is the safety region, and v_{lim} is the speed limit.

In the **Reach** tasks, the agent is supposed to move towards a series of goals. Obstacles are placed to hinder the agent from trivial solutions. The agent is rewarded for moving closer to the goal and entering the goal zone. The reward and cost are defined as:

$$\begin{aligned} r(s_t) &= (\phi_t - \phi_{t-1}) + r_{\text{agent}}(s_t) \\ c(s_t) &= \min(1, N_t + \mathbf{1}(z_t > z_{\text{lim}})) \end{aligned} \quad (31)$$

where $\phi_t = -\|(x_t, y_t) - (g_x, g_y)\|_2$ is the task potential to the current goal (g_x, g_y) , N_t is the number of collisions and z_{lim} is the height limit.

In the **Gather** tasks, the agent is expected to navigate and collect as many green apples as possible while avoiding red bombs. In contrast to the other tasks, agents in the gather tasks receive only sparse rewards when reaching apples. Costs are also sparse and are received when touching bombs. Let the agent’s 2D position be $\mathbf{p}_t = (x_t, y_t)$ and let $\mathcal{O}_{\text{apple}}$ denote the set of all apple-type obstacles. Let $\delta_t(o) = \mathbb{I}[\text{visible}_t(o) \wedge d_{o,t} < d_{\text{det}}]$. The reward is defined as:

$$r(s_t) = \begin{cases} R_{\text{dead}}, & \text{if the agent is dead at } t, \\ R_{\text{apple}} \sum_{o \in \mathcal{O}_{\text{apple}}} \delta_t(o), & \text{otherwise} \end{cases} \quad (32)$$

where $d_{o,t} = \|\mathbf{p}_{o,t} - \mathbf{p}_t\|$ is the agent distance to apple. The cost is defined as:

$$c(s_t) = \min(1, N_t \cdot B + \mathbf{1}(z_t > z_{\text{lim}})) \quad (33)$$

where N_t is the number of bombs touched at time t and B the per-bomb penalty.

D.2 Hyperparameters

For all baseline methods, we adopt their default hyperparameter configurations. To ensure a fair comparison across all methods, we set the rollout length for each task to match the maximum number of interaction steps used by GeoDT and the baselines. The cost threshold for the baselines is set to 10 for all tasks, except for the **Gather** task, where it is set to 1. The key hyperparameters used for the Causal Transformer-based models (GeoDT, CDT and PDT) are provided in Table 5. The key hyperparameters used for other baselines are shown in Table 6. We apply the same setting for multi-task and single-task cases.

Table 5: Causal Transformer-based model configuration parameters

Parameter	GeoDT	CDT	PDT
<i>Common Settings:</i>			
Embedding dimension		128	
Batch size		512	
Dropout		0.1	
<i>Architecture:</i>			
Number of layers	3	3	3
Number of attention heads	8	8	1
Context length	20	10	20
Prompt length	20	–	20
<i>Optimization:</i>			
Learning rate ($\times 10^{-4}$)	1.0	1.0	1.0
Grad norm clip	0.25	0.25	0.50

Table 6: Other baseline model configuration parameters

Parameter	CPQ	BCQ-L	LSPC	SoftMod	FISOR	O3SRL
<i>Common Settings:</i>						
Training steps			1×10^6			
Batch size			512			
Dropout			0.1			
<i>Algorithm-Specific Settings:</i>						
Hidden layer size	256	256	256	64	256	256
Soft update rate (τ)	0.005	0.005	0.005	0.005	0.001	0.005
<i>Learning Rates ($\times 10^{-3}$):</i>						
Actor learning rate	1.0	1.0	0.3	0.1	0.3	0.3
Critic learning rate	1.0	1.0	0.3	0.1	0.3	0.3
Activation function	ReLU	ReLU	ReLU	ReLU	<i>mish</i>	ReLU

E Additional Experiments Results

E.1 Experiment Results on Agent Drone and Ant

We further evaluate the effectiveness of GeoDT on the more complex **Drone** agent, which features higher-dimensional observations and actions across tasks. We conduct multi-task experiments with the number of tasks to be learned ranging from 2 to 4. Table 7 presents the performance of various methods on the full set of four tasks: **Run**, **Circle**, **Reach**, and **Gather**.

Compared with agent **Drone**, **Ant** is even more complex, which has the highest state dimension over all the **SafetyBulletGym** environments. Following the same evaluation protocol of previous experiments, we compare the performance of GeoDT and other baselines on the full set of four tasks: **Run**, **Circle**, **Reach**, and **Gather**, which results are shown in Table 8.

For both agents, despite the high-dimensional state space, GeoDT demonstrates exceptional multi-task performance across all baselines, maintaining safety satisfaction in all tasks, highlighting the effectiveness of the GeoDT framework.

Table 7: Performance comparison with the Drone agent across tasks.

Methods	Drone-Run		Drone-Circle		Drone-Reach		Drone-Gather		Average	
	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost
GeoDT	0.67	0.89	0.75	0.86	0.33	0.92	0.16	0.15	0.477	0.704
CDT	0.65	0.38	0.72	0.79	0.40	1.29	0.13	0.30	0.474	0.690
PDT	0.43	0.45	0.38	0.72	0.13	0.79	0.12	0.76	0.262	0.678
CPQ	0.85	1.42	0.51	1.78	0.21	1.32	0.12	0.23	0.419	1.188
BCQL	0.79	1.20	0.83	3.41	0.35	0.85	0.09	0.33	0.517	1.449
LSPC	0.61	0.75	0.85	1.834	0.30	1.43	0.21	0.15	0.493	1.041
SoftMod	0.91	1.51	0.69	4.93	0.28	0.90	0.13	0.15	0.502	1.872
FISOR	0.02	0.91	0.37	0.07	0.32	0.78	0.14	0.08	0.211	0.459
O3SRL	0.47	0.74	0.82	1.83	0.43	1.15	0.15	0.50	0.469	1.054

Table 8: Performance comparison with the Ant agent across tasks.

Methods	Ant-Run		Ant-Circle		Ant-Reach		Ant-Gather		Average	
	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost
GeoDT	0.86	0.25	0.54	0.87	0.27	0.31	0.15	0.05	0.455	0.370
CDT	0.91	0.61	0.81	1.63	0.07	1.79	0.10	0.05	0.472	1.020
PDT	0.23	0.03	0.10	0.31	0.05	1.13	0.02	0.03	0.100	0.375
CPQ	0.10	0.00	0.02	0.00	-0.10	0.95	-0.28	0.21	-0.066	0.245
BCQL	1.26	7.74	0.68	3.12	0.07	0.82	-0.10	0.03	0.479	2.924
LSPC	0.86	2.19	1.02	3.11	0.22	0.26	0.20	0.00	0.573	1.388
SoftMod	0.71	4.18	0.87	3.61	0.18	0.46	0.10	0.05	0.463	2.073
FISOR	0.17	0.01	0.22	0.00	0.17	0.58	0.15	0.03	0.178	0.153
O3SRL	0.50	0.19	0.90	4.32	0.36	0.23	0.10	0.00	0.464	1.183

E.2 How Does the Prompt Length Affect the Performance?

Our method uses trajectory prompts both to identify the current task and to guide the agent’s behaviour. In this section, we evaluate how prompt length influences performance by comparing multi-task results for the Car and Ball agents with prompt lengths of 10, 20, and 40. The results are summarized in Table 9 and Table 10.

Table 9: Prompt length effect on the Car agent across tasks.

Prompt Len	Car-Run		Car-Circle		Car-Reach		Car-Gather		Average	
	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost
10	0.951	1.345	0.717	0.813	0.666	0.638	0.358	0.600	0.673	0.849
20	0.947	0.770	0.726	0.614	0.701	0.673	0.413	0.420	0.697	0.620
40	0.943	0.645	0.731	0.608	0.650	0.630	0.367	0.625	0.673	0.627

At the beginning of each episode, the trajectory prompt is randomly selected from the offline dataset, without bias toward high-reward or low-cost segments. Although one might expect that selecting such prompts could enhance performance, we did not observe clear improvements in preliminary experiments. However, it is noticed that the prompt length does make a difference, and its effectiveness depends critically on the quality and diversity of the offline datasets and must be interpreted on a per-task basis—it is by no means “the longer, the better.” For dynamically simple tasks (e.g., Run and Circle), the dataset typically contains high-quality, highly repetitive trajectories; longer prompts can thus leverage more safe-and-efficient behavior examples, enabling the Transformer to better capture safety patterns and achieve lower costs or higher rewards. Indeed, as prompt length increases, the safety of Ball-Run, Car-Run, and Car-Circle improves, while the reward for Ball-Circle and Car-Circle rises markedly. By contrast, for dynamically complex tasks (e.g., Reach

Table 10: Prompt length effect on the Ball agent across tasks.

Prompt Len	Ball-Run		Ball-Circle		Ball-Reach		Ball-Gather		Average	
	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost
10	0.804	1.128	0.737	0.650	0.895	0.890	0.525	0.400	0.740	0.767
20	0.740	0.743	0.745	0.627	0.874	0.785	0.480	0.630	0.710	0.696
40	0.753	0.903	0.762	0.635	0.869	0.658	0.492	0.575	0.720	0.693

and **Gather**), where obstacle layouts and goal positions vary randomly and data noise is significant, overly long contexts fail to mitigate distributional shift and instead introduce irrelevant trajectory fragments. The resulting noise in task identification or policy generation confuses the model and degrades performance. Consequently, shorter or medium-length prompts (10 or 20) prove more robust for these tasks.

In summary, prompt length should be tailored to task complexity to balance information richness against noise interference. Simple tasks benefit from longer prompts that provide more positive examples, whereas complex tasks require restraint to avoid diluting critical context. In our experiments, a prompt length of 20 consistently achieved the lowest cost across most tasks, suggesting it strikes an optimal trade-off between capturing safety patterns and minimizing redundant interference.

E.3 Optional Prototype Balancing Regularization

In GeoDT, the prototype memory is retrieved through a soft assignment mechanism, which may in principle become overly concentrated on a small subset of prototypes. To examine whether a lightweight balancing regularizer can improve memory utilization, we additionally consider an optional prototype balancing loss.

Let $\beta \in \mathbb{R}^{B \times K}$ denote the prototype retrieval weights for a batch, where K is the number of prototypes. We compute the average prototype usage

$$\bar{\beta}_k = \frac{1}{B} \sum_{b=1}^B \beta_k^{(b)}, \quad (34)$$

and define a balancing regularizer as

$$\mathcal{L}_{\text{bal}} = \sum_{k=1}^K \bar{\beta}_k \log \bar{\beta}_k, \quad (35)$$

which encourages a less degenerate usage distribution over the prototype memory. The overall objective becomes

$$\mathcal{L} = \mathcal{L}_{\text{DT}} + \lambda_{\text{geo}} \mathcal{L}_{\text{geo}} + \lambda_{\text{bal}} \mathcal{L}_{\text{bal}}. \quad (36)$$

We evaluate this variant with $K = 8$, which is the default setting used in the main experiments. As shown in Table 11, adding prototype balancing leads to only minor changes in reward, but tends to increase the realized cost. In other words, encouraging a more even prototype usage distribution does not improve the overall reward–cost trade-off in our setting, even though it can slightly alter the retrieval behavior of the memory module.

This observation is consistent with the design of GeoDT. The prototype memory is intended to provide a compact set of reusable geometric anchors, but we do not require these prototypes to be maximally distinct or uniformly utilized. Some overlap in prototype usage is natural, since different tasks may still share partially similar geometric structure. We therefore treat prototype balancing as an optional regularizer rather than a core component of the method, and keep the default GeoDT formulation as simple as possible in the main paper. Indeed, GeoDT does not rely on enforcing prototypes to be maximally dissimilar; rather, it only requires them to provide a useful and reusable set of geometric reference patterns.

Table 11: Effect of optional prototype balancing regularization with $K = 8$. While balancing slightly changes average reward, it tends to increase realized cost, leading to no clear improvement in the overall reward–cost trade-off.

Method	Car Avg		Ant Avg		Overall Avg	
	Reward	Cost	Reward	Cost	Reward	Cost
GeoDT ($K = 8$)	0.697	0.620	0.418	0.443	0.558	0.532
GeoDT + balance ($K = 8$)	0.695	0.718	0.428	0.52	0.562	0.619

E.4 How Does the Safe Coefficient Affect TSRS and ITBS?

In section 5.2, we initially evaluated all methods using a fixed safe coefficient $\lambda = 0.5$ to balance reward and cost in a principled way. To better understand how this balance influences comparative performance, we conducted a sensitivity analysis by sweeping $\lambda \in \{0.0, 0.25, 0.5, 0.75, 1.0\}$ and computing TSRS and ITBS for each method at $n = 4$ in both the Car and Ball environments.

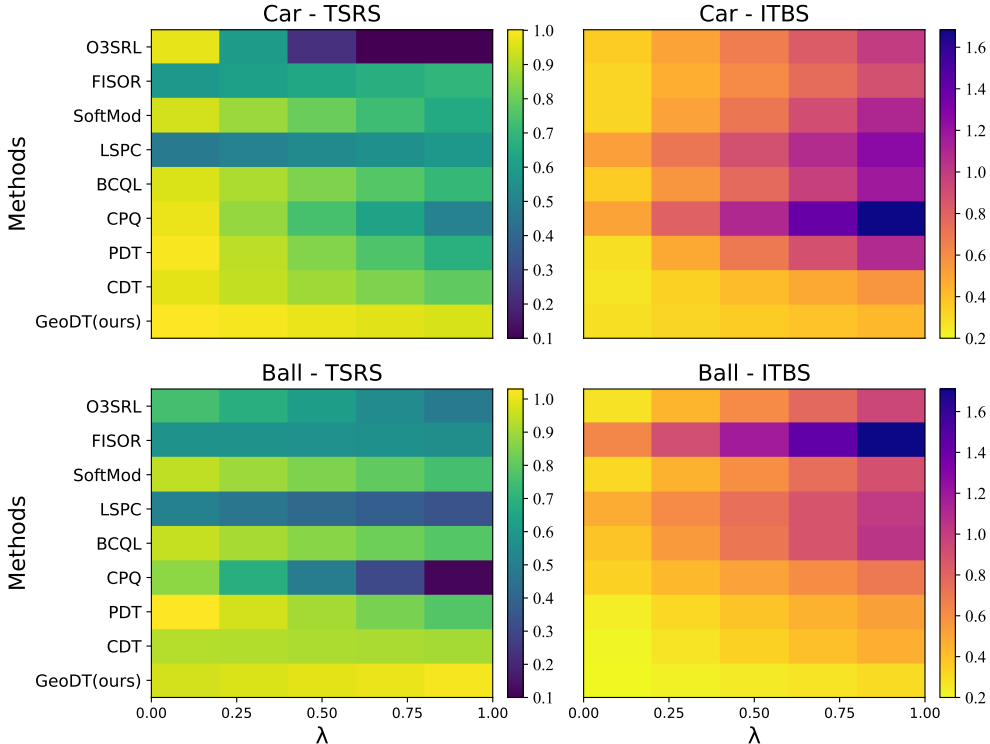


Figure 5: Effect of the safety coefficient sweep on TSRS and ITBS. As the safety coefficient λ increases, TSRS tends to decrease while ITBS increases. A moderate value, such as $\lambda = 0.5$, offers a reasonable trade-off between model separation and ranking stability. Notably, GeoDT consistently achieves the highest TSRS and lowest ITBS across all settings, demonstrating its robustness to varying safety weights and, more importantly, its balanced response to reward and cost trade-offs.

As shown in Figure 5, the results demonstrate that our method yields the smoothest TSRS curves and the smallest ITBS fluctuations across all λ values and in both environments, consistently tracking close to the ideal reference. This remarkable insensitivity to λ confirms that GeoDT’s internal reward–cost balancing mechanism and safety control remain robust even under extreme weightings. Moreover, we find that $\lambda \approx 0.5$ offers an optimal compromise between model separation and ranking stability: at lower λ , TSRS gradually

decreases (improving reward alignment) while ITBS increases (focusing greater on cost variability), yet GeoDT maintains its characteristic flat profile around this point.

For applications prioritizing maximal differentiation in the `Car` environment, raising λ toward the 0.75–1.0 range is effective; conversely, in the `Ball` environment, selecting $\lambda \in [0.25, 0.5]$ better suppresses inter-task volatility. These findings not only justify GeoDT’s default choice of $\lambda = 0.5$ but also provide practical guidance for dynamically tuning λ in other methods to balance discriminative power against performance consistency.

F Single-task performance

We report the single-task performance of all methods on four agents for reference in Table 12, 13, 14 and 15. Although GeoDT is not the best-performing model in every individual task—several baselines achieve top scores on specific tasks—it consistently delivers competitive results across all settings. As previously discussed, GeoDT exhibits strong robustness under increasing task diversity: as the number of tasks increases, it maintains performance close to the single-task upper bound, demonstrating its robustness and adaptability in multi-task environments.

Table 12: Single-Task Performance on `Car` agent.

Methods	Car-Run		Car-Circle		Car-Reach		Car-Gather		Average	
	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost
GeoDT	0.95	0.71	0.74	0.77	0.70	0.57	0.37	0.58	0.690	0.658
CDT	0.95	0.74	0.70	0.60	0.67	0.88	0.48	0.51	0.700	0.683
PDT	0.93	1.53	0.68	1.26	0.67	0.90	0.38	0.50	0.665	1.048
CPQ	0.75	0.52	0.65	0.53	0.55	1.33	0.15	0.55	0.525	0.733
BCQL	0.93	0.70	0.64	1.80	0.61	0.64	0.30	0.52	0.620	0.915
LSPC	0.71	0.22	0.78	0.23	0.71	0.73	0.48	0.55	0.671	0.430
SoftMod	0.93	0.64	0.75	2.56	0.70	0.70	0.38	0.70	0.690	1.150
FISOR	0.76	0.05	0.40	0.19	0.54	0.48	0.47	0.38	0.543	0.275
O3SRL	0.94	0.00	0.76	0.00	0.57	0.14	0.38	0.50	0.603	0.160

Table 13: Single-Task Performance on `Ball` agent.

Methods	Ball-Run		Ball-Circle		Ball-Reach		Ball-Gather		Average	
	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost
GeoDT	0.79	0.84	0.77	0.63	0.89	0.84	0.48	0.55	0.733	0.715
CDT	0.88	0.86	0.79	0.82	0.88	0.64	0.45	0.55	0.750	0.718
PDT	0.79	0.81	0.73	0.90	0.72	0.75	0.49	0.43	0.683	0.723
CPQ	0.56	0.68	0.68	0.23	0.59	0.51	0.24	0.26	0.518	0.420
BCQL	0.35	0.20	0.79	1.20	0.88	0.75	0.46	0.52	0.620	0.668
LSPC	0.75	0.84	0.48	0.24	0.48	0.62	0.38	0.40	0.522	0.526
SoftMod	1.00	1.04	0.99	1.99	0.85	0.72	0.50	0.35	0.835	1.025
FISOR	0.31	0.00	0.44	0.01	0.65	0.78	0.44	0.73	0.460	0.380
O3SRL	0.66	0.02	0.63	0.02	0.88	0.94	0.37	0.50	0.633	0.369

Table 14: Single-Task Performance on **Drone** agent.

Methods	Drone-Run		Drone-Circle		Drone-Reach		Drone-Gather		Average	
	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost
GeoDT	0.72	0.79	0.77	0.87	0.32	0.78	0.11	0.30	0.480	0.685
CDT	0.71	0.60	1.14	0.78	0.46	1.27	0.15	0.30	0.615	0.738
PDT	0.54	0.36	0.42	0.83	0.20	1.16	0.12	0.44	0.320	0.698
CPQ	0.26	0.44	0.56	0.56	0.15	0.53	0.06	0.08	0.258	0.403
BCQL	0.65	0.71	0.68	1.19	0.36	1.07	0.13	0.05	0.455	0.755
LSPC	0.72	0.69	0.66	0.67	0.45	0.90	0.25	0.25	0.518	0.627
SoftMod	0.56	0.23	0.96	1.95	0.30	0.77	0.13	0.25	0.488	0.800
FISOR	0.60	0.00	0.41	0.35	0.36	0.66	0.11	0.28	0.370	0.323
O3SRL	0.38	0.00	0.56	0.74	0.44	0.66	0.60	0.10	0.495	0.375

Table 15: Single-Task Performance on **Ant** agent.

Methods	Ant-Run		Ant-Circle		Ant-Reach		Ant-Gather		Average	
	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost	Reward	Cost
GeoDT	0.90	0.43	0.44	0.74	0.24	0.68	0.20	0.05	0.445	0.475
CDT	0.92	0.62	0.85	1.34	0.14	1.03	-0.02	0.06	0.472	0.763
PDT	0.38	0.21	0.10	0.18	0.07	0.85	-0.08	0.02	0.118	0.315
CPQ	0.13	0.01	0.00	0.00	0.04	0.83	-0.09	0.06	0.020	0.225
BCQL	1.02	4.63	0.61	1.42	0.17	0.72	-0.13	0.02	0.418	1.698
LSPC	0.94	1.46	0.40	0.78	0.23	0.18	0.20	0.10	0.440	0.630
SoftMod	1.21	8.34	0.74	2.48	0.15	0.41	0.10	0.10	0.550	2.833
FISOR	0.60	0.38	0.41	0.00	0.18	0.29	0.15	0.05	0.335	0.180
O3SRL	0.37	0.54	0.50	0.00	0.20	0.81	0.30	0.00	0.343	0.338