

---

# SAMBA: Safe Model-Based & Active Reinforcement Learning

---

**Alexander I. Cowen-Rivers\***

Huawei R&D UK  
alexander.cowen.rivers@huawei.com

**Daniel Palenicek\***

Huawei R&D UK  
Technical University Darmstadt  
daniel.palenicek@stud.tu-darmstadt.de

**Vincent Moens\***

Huawei R&D UK  
vincent.moens@huawei.com

**Mohammed Amin ABDULLAH**

Huawei R&D UK  
mohammed.abdullah@huawei.com

**Aivar Sootla**

Huawei R&D UK  
aivar.sootla@huawei.com

**Jun Wang**

Huawei R&D UK  
University College London  
w.j@huawei.com

**Haitham Bou-Ammar<sup>†</sup>**

Huawei R&D UK  
University College London  
haitham.ammar@huawei.com

## Abstract

In this paper, we propose SAMBA, a novel framework for safe reinforcement learning that combines aspects from probabilistic modelling, information theory, and statistics. Our method builds upon PILCO to enable active exploration using novel (semi-)metrics for out-of-sample Gaussian process evaluation optimised through a multi-objective problem that supports conditional-value-at-risk constraints. We evaluate our algorithm on a variety of safe dynamical system benchmarks involving both low and high-dimensional state representations. Our results show orders of magnitude reductions in samples and violations compared to state-of-the-art methods. Lastly, we provide intuition as to the effectiveness of the framework by a detailed analysis of our active metrics and safety constraints.

## 1 Introduction

Reinforcement learning (RL) has seen successes in domains such as video and board games [35, 48, 49], and control of simulated robots [4, 42, 43]. Though successful, these applications assume idealised simulators and require tens of millions of agent-environment interactions typically performed by randomly exploring policies. In real-world safety-critical applications, however, such an idealised framework of random exploration with the ability to gather samples at ease falls short, partly due to the catastrophic costs of failure and the high operating costs. Hence, if RL algorithms are to be applied in the real world, safe agents that are sample-efficient and capable of

---

\*The first three authors are joint first authors.

<sup>†</sup>Honorary position at UCL.

mitigating risk need to be developed. To this end, different works adopt varying safety definitions, where some are interested in safe learning, i.e., safety *during* the learning process, while others focus on acquiring safe policies *eventually*. Safe learning generally requires some form of safe initial policy as well as regularity assumptions on dynamics, all of which depend on which notion of safety is considered. For instance, safety defined by constraining trajectories to safe regions of the state-action space is studied in, e.g., [2], which assumes partially-known control-affine dynamics with Lipschitz regularity conditions, as well as in [5, 32], both of which require strong assumptions on dynamics and initial control policies to give theoretical guarantees of safety. Safety in terms of Lyapunov stability [31] is studied in, e.g., [13, 14] (model-free), which require a safe initial policy, and [7] (model-based), which requires Lipschitz dynamics and a safe initial policy. The work in [54] (which builds on [52]) considers deterministic dynamics and attempts to optimise expected return while not violating a pre-specified safety threshold. Several papers attempt to keep expectation constraints satisfied during learning, e.g., [1] extends ideas of [29], [16] adds a safety layer which makes action corrections to maintain safety. When it comes to safe final policies (defined in terms of corresponding constraints), on the other hand, some works [11, 37] considered risk-based constraints and developed model-free solvers.

Unfortunately, most of these methods are sample-inefficient and make a large number of visits to unsafe regions. Given our interest in algorithms achieving safe final policies while reducing the number of visits to unsafe regions, we pursue a safe model-based framework that assumes no safe initial policies nor *a priori* knowledge of the transition model. As we believe that sample efficiency is key for effective safe solvers, we choose Gaussian processes [39] for our model. Of course, any such framework is hampered by the quality and assumptions of the model’s hypothesis space. Aiming at alleviating these restrictions, we go further and integrate *active* learning (discussed in Section 2.2), wherein the agent influences where to query/sample from in order to generate new data so as to reduce model uncertainty, and so ultimately to learn more efficiently. Successful application of this approach is very much predicated upon the chosen method of quantifying potential uncertainty reduction, i.e., what we refer to as the *(semi-)active metric*. Common (semi-)active metrics are those which identify points in the model with large entropy or large variance, or where those quantities would be most reduced in a posterior model if the given point were added to the data set [20, 33, 34, 46]. However, our desire for safety adds a complication, since a safe learning algorithm will likely have greater uncertainty in regions where it is unsafe by virtue of not exploring those regions [19, 30]. Indeed our experiments (Figure 1) support this claim.

Attacking the above challenges, we propose two novel out-of-sample (semi-)metrics for Gaussian processes that allow for exploration in novel areas while remaining close to training data, thus avoiding unsafe regions. To enable effective and grounded introduction of active exploration and safety constraints, we define a novel constrained bi-objective formulation of RL and provide a policy multi-gradient solver that is proven effective on a variety of safety benchmarks. In short, our contributions can be stated as follows: 1) novel constrained bi-objective formulation enabling exploration and safety consideration, 2) safety-aware active (semi-)metrics for exploration, and 3) policy multi-gradient solver trading off cost minimisation, exploration maximisation, and constraint feasibility. We test our algorithm on three stochastic dynamical systems after augmenting these with safety regions and demonstrate a significant reduction in sample and cost complexities compared to the state-of-the-art.

## 2 Background and notation

### 2.1 Reinforcement learning

We consider Markov decision processes (MDPs) with continuous states and action spaces;  $\mathcal{M} = \langle \mathcal{X}, \mathcal{U}, \mathcal{P}, c, \gamma \rangle$ , where  $\mathcal{X} \subseteq \mathbb{R}^{d_{\text{state}}}$  denotes the state space,  $\mathcal{U} \subseteq \mathbb{R}^{d_{\text{act}}}$  the action space,  $\mathcal{P} : \mathcal{X} \times \mathcal{U} \times \mathcal{X} \rightarrow [0, 1]$  is a transition density function,  $c : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$  is the cost function and  $\gamma \in [0, 1]$  is a discount factor. At each time step  $t = 0, \dots, T$ , the agent is in state  $\mathbf{x}_t \in \mathcal{X}$  and chooses an action  $\mathbf{u}_t \in \mathcal{U}$  transitioning it to a successor state  $\mathbf{x}_{t+1} \sim \mathcal{P}(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t)$ , and yielding a cost  $c(\mathbf{x}_t, \mathbf{u}_t)$ . Given a state  $\mathbf{x}_t$ , an action  $\mathbf{u}_t$  is sampled from a policy  $\pi : \mathcal{X} \times \mathcal{U} \rightarrow [0, 1]$ , where we write  $\pi(\mathbf{u}_t | \mathbf{x}_t)$  to represent the conditional density of an action. Upon subsequent interactions, the agent collects a trajectory  $\tau = [\mathbf{x}_{0:T}, \mathbf{u}_{0:T}]$ , and aims to determine an optimal policy  $\pi^*$  by minimising total expected cost:  $\pi^* \in \arg \min_{\pi} \mathbb{E}_{\tau \sim p_{\pi}(\tau)} [\mathcal{C}(\tau) := \sum_{t=0}^{T-1} \gamma^t c(\mathbf{x}_t, \mathbf{u}_t)]$ , where  $p_{\pi}(\tau)$  denotes the trajectory density defined as:  $p_{\pi}(\tau) = \mu_0(\mathbf{x}_0) \prod_{t=0}^{T-1} \mathcal{P}(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) \pi(\mathbf{u}_t | \mathbf{x}_t)$ , with

$\mu_0(\cdot)$  being an initial state distribution.

**Constrained MDPs:** The above can be generalised to include various forms of constraints, often motivated by the desire to impose some form of safety measures. Examples are expectation constraints [1, 3] (which have the same form as the objective, i.e., expected discounted sum of costs), constraints on the variance of the return [38], chance constraints (a.k.a. Value-at-Risk (VaR)) [12], and Conditional Value-at-Risk (CVaR) [11, 12, 37]. The latter is the constraint we adopt in this paper for reasons that will be elucidated upon below. Adding constraints means we can't directly apply standard algorithms like policy gradient [51], and different techniques are required, e.g., via Lagrange multipliers [8], as was done in [11, 12, 37] besides many others. Further, current methods only consider cost minimisation with no regard to exploration as we do in this paper.

**Model-Based Reinforcement Learning:** Current solutions to the problem described above (constrained or unconstrained) can be split into model-free and model-based methods. Though effective, model-free algorithms are highly sample inefficient [26]. For sample-efficient solvers, we follow model-based strategies that we now detail. To reduce the number of interactions with the real environments, model-based solvers build surrogate models,  $\mathcal{P}_{\text{surr}}$ , to determine optimal policies. These methods, typically, run two main loops. The first gathers traces from the real environment to update  $\mathcal{P}_{\text{surr}}$ , while the second improves the policy using  $\mathcal{P}_{\text{surr}}$  [19, 24]. Among various candidate models, e.g., world models [23], in this paper, we follow PILCO [19] and adopt Gaussian processes (GPs) as we believe that uncertainty quantification and sample efficiency are key for real-world considerations of safety. In this construction, one places a Gaussian process prior on a latent function  $f$  to map between input-output pairs. Such a prior is fully specified by a mean,  $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$ , and a covariance function  $k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$  [39]. We write  $f \sim \mathcal{GP}(m(\cdot), k(\cdot, \cdot))$  to emphasize that  $f$  is sampled from a GP [39]. Given a data-set of input-output pairs  $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{n_1}$ , corresponding, respectively, to state-action and successor state tuples, one can perform predictions on a query set of  $n_2$  test data points  $\{\mathbf{x}_*^{(j)}\}_{j=1}^{n_2}$ . Such a distribution is Gaussian with predictive mean-vectors and covariance matrices given by:  $\boldsymbol{\mu}_* = \mathbf{K}_{n_2, n_1} \mathbf{A}_{n_1, n_1} \mathbf{y}_{n_1}$  and  $\boldsymbol{\Sigma}_{n_2, n_2} = \mathbf{K}_{n_2, n_2} - \mathbf{K}_{n_2, n_1} \mathbf{A}_{n_1, n_1} \mathbf{K}_{n_1, n_2}$ , where  $\mathbf{A}_{n_1, n_1} = [\mathbf{K}_{n_1, n_1} + \sigma_\omega^2 \mathbf{I}]^{-1}$  with  $\sigma_\omega$  being the noise covariance that is assumed to be Gaussian. In the above, we also defined  $\mathbf{y}_{n_1}$  as a vector concatenating all training labels,  $\mathbf{K}_{n_1, n_1} = K(\mathbf{X}, \mathbf{X}) = [k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})]_{ij}$ ,  $\mathbf{K}_{n_1, n_2} = \mathbf{K}_{n_2, n_1}^\top = K(\mathbf{X}, \mathbf{X}_*) = [k(\mathbf{x}^{(i)}, \mathbf{x}_*^{(j)})]_{ij}$ , and  $\mathbf{K}_{n_2, n_2}(\mathbf{X}_*, \mathbf{X}_*) = [k(\mathbf{x}_*^{(i)}, \mathbf{x}_*^{(j)})]_{ij}$ , where  $\mathbf{X}$  and  $\mathbf{X}_*$  are feature matrices with  $\#\text{input-dim} \times n_1$  and  $\#\text{input-dim} \times n_2$  sizes respectively. We executed training in GPyTorch [21], and used multi-output-GPs as defined in [18].

## 2.2 Active learning in dynamical systems

In *active learning* [20, 46], an agent chooses points to sample/query that best improve learning or model updates. This is often performed by optimising an *acquisition function*, which gives some quantification of how much a model would improve if a given data point were queried, e.g., points where the model has high entropy or where variance can be most reduced. Active learning with GPs has been studied in the static case, where points can be selected at will (see, e.g., [33, 34]). In the context of dynamical systems, however, added complications arise as one is not always able to directly drive the system into a desired state. Recent work has attempted to resolve this problem, e.g., in [10] and [44], receding horizon optimisation is used to iteratively update a model, and in [10], actions are favoured that maximise the sum of differential entropy terms at each point in the mean trajectory predicted to occur by those actions. Moreover, in [44], a sum of variance terms is optimised to improve Bayesian linear regression. Again, for computational tractability, the predicted mean of states is used as propagating state distributions in the model is difficult. Different to our paper, neither of these works deal with safety, nor do they have additional objectives to maximise/minimise avoiding a bi-objective formulation. In [28] a GP model that is used for MPC is updated by greedily selecting points which maximise *information gain*, i.e., reduction in entropy, as is done in [34]. Only very recently, the authors in [6] proposed an active learning approach coupled with MBRL. Similar to SAMBA, they use an adaptive convex combination of objectives, however their exploration metric is based on reward variance computed from a (finite) collection of models increasing the burden on practitioners who now need to predefine the collection of dynamics. They do not use GPs as we do, and do not consider safety. Compared to [6], we believe SAMBA is more flexible supporting model-learning from scratch and enabling principled exploration coupled with safety consideration. Further afield from our work, active learning has been recently studied in the context of GP time-series in [55], and for pure exploration in [47], which uses a finite collection of models.

Our (semi-)metrics generalise the above to consider safe-regions and future information trade-off as we detail in Section 3.2.

### 3 SAMBA: Framework & solution

In designing SAMBA, we take PILCO [19] as a template and introduce two novel ingredients allowing for active exploration and safety. Following PILCO, SAMBA runs a main loop that gathers traces from the real environment and updates a surrogate model,  $\mathcal{P}_{\text{GP}}(\cdot) : \mathcal{X} \times \mathcal{U} \times \mathcal{X} \rightarrow [0, 1]$ , encoded by a Gaussian process. Given  $\mathcal{P}_{\text{GP}}(\cdot)$ , PILCO and other model-based methods [50] attempt to obtain a policy that minimises total-expected cost with respect to traces,  $\tau$ , sampled from the learnt model by solving  $\min_{\pi} \mathbb{E}_{\tau \sim p_{\text{sur}}(\tau)} [\mathcal{C}(\tau)]$  with  $p_{\text{sur}}(\tau) = \mu_0(\mathbf{x}_0) \prod_{t=0}^{T-1} \mathcal{P}_{\text{GP}}(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) \pi(\mathbf{u}_t | \mathbf{x}_t)$ . The updated policy is then used to sample new traces from the real system where the above process repeats. During this sampling process, model-based algorithms consider various metrics in acquiring transitions that reveal novel information, which can be used to improve the surrogate model’s performance. PILCO, for instance, makes use of the GP uncertainty, while ensemble models [41, 53] explore by their aggregated uncertainties. With sufficient exploration, this allows policies obtained from surrogate-models to control real-systems. Our safety considerations mean we would prefer agents that learn well-behaving policies with minimal sampling from unsafe regions of state-action spaces; a property we achieve later by incorporating CVaR constraints as we detail in Section 3.1. Requiring a reduced number of visits to unsafe regions, hence, lessens the amount of “unsafe” data gathered in such areas by definition. Therefore, model entropy is naturally increased in these territories and algorithms following such exploration strategies are, as a result, encouraged to sample from hazardous states. As such, a naive adaptation of entropy-based exploration can quickly become problematic by contradicting safety requirements. To circumvent these problems, we introduce two new active exploration (semi-)metrics in Section 3.2, that assess information beyond training-data availability and consider input-output data distributions. Our (semi-)metrics operate under the assumption that during any model update step, “safe” transition data (i.e., a set of state-action-successor states sampled from safe regions) is more abundant in number when compared to “unsafe” triplets. Considering such a skew between distributions, our (semi-)metrics yield increased values on test queries close to high-density training data. Given such (semi-)metrics, we enable novel model-based algorithms that solve a bi-objective optimisation problem that attempts to minimise cost, while maximising active values. In other words, during this step, we not only update policies to be well-behaving in terms of the total cost but also to actively explore safe transitions that allow for improved models in successor iterations.

Of course, the assumption of having skew towards safe regions in training data distribution is generally not true since solving the above only ensures good expected returns. To frequently sample safe regions, we augment cost minimisation with a safety constraint that is encoded through the CVaR of a user-defined safety cost function with respect to *model* traces. Hence, SAMBA solves a bi-objective constrained optimisation problem (§3.1) aimed at minimising cost, maximising active exploration, and meeting safety constraints.

#### 3.1 Bi-objective constrained MDPs

Given a (GP) model of an environment, we formalise our problem as a generalisation of constrained MDPs to support bi-objective losses. We define a bi-objective MDP by a tuple  $\mathcal{M}_{\text{BiO}} = \langle \mathcal{X}, \mathcal{U}, \mathcal{P}_{\text{GP}}, c, l, \zeta, \gamma \rangle$  consisting of state space  $\mathcal{X} \subseteq \mathbb{R}^{d_{\text{state}}}$ , action space  $\mathcal{U} \subseteq \mathbb{R}^{d_{\text{act}}}$ , Gaussian process transition model  $\mathcal{P}_{\text{GP}}$ , cost function  $c : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ , constraint cost function (used to encode safety)  $l : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ , additional objective function  $\zeta : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ , and discount factor  $\gamma \in [0, 1]$ .<sup>3</sup> In our setting, for instance,  $l$  encodes the state-action’s risk by measuring the distance to an unsafe region, while  $\zeta$  denotes an active (semi-)metric from these described in Section 3.2. To finalise the problem definition of bi-objective MDPs, we need to consider an approachable constraint to describe safety considerations. In incorporating such constraints, we are chiefly interested in those that are flexible (i.e., can support different user-designed safety criteria) and allow us to quantify events occurring in tails of cost distributions. When surveying single-objective constrained MDPs, we realise that the literature predominately focuses on expectation-type constraints [1, 40] – a not so flexible approach restricted to being safe on average. Others, however, make use of conditional-value-at-risk (CVaR); a coherent risk measure [12] that provides powerful and flexible notions of safety

<sup>3</sup>It is worth noting that the (semi-)metrics we develop in Section 3.2 map to  $\mathbb{R}_+$  instead of  $\mathbb{R}$ .

(i.e., can support expectation, tail-distribution, or hard – unsafe state visitation – constraints) and quantifies tail risk in the worst  $(1 - \alpha)$ -quantile. Formally, given a random variable  $Z$ ,  $\text{CVaR}_\alpha(Z)$  is defined as:  $\text{CVaR}_\alpha(Z) = \min_\nu [\nu + \frac{1}{1-\alpha} \mathbb{E}[(Z - \nu)^+]]$ , where  $(Z - \nu)^+ = \max(Z - \nu, 0)$ . With such a constraint, we can write the optimisation problem of our bi-objective MDP as:

$$\min_{\pi} \mathbb{E}_{\tau \sim p_{\text{surf}}(\tau)} [\mathcal{L}(\tau), \zeta(\tau)]^T \quad \text{s.t.} \quad \text{CVaR}_\alpha(\mathcal{L}(\tau)) \leq \xi, \quad (1)$$

with  $\mathcal{L}(\tau) = \sum_{t=0}^T \gamma^t l(\mathbf{x}_t, \mathbf{u}_t)$  being total accumulated safety cost along  $\tau$ , and  $\xi \in \mathbb{R}_+$  a safety threshold. Of course, Equation 1 describes a problem not standard to reinforcement learning. In Section 3.3, we devise *policy multi-gradient* updates to determine  $\pi$ .

### 3.2 $\zeta$ -functions for safe active exploration

In general,  $\zeta$  can be any bounded objective that needs to be maximised/minimised in addition to standard cost. Here, we choose one that enables active exploration in safe state-action regions. To construct  $\zeta$ , we note that a feasible policy – i.e., one abiding by CVaR constraints – of the problem in Equation 1 samples tuples that mostly reside in safe regions. As such, the training data distribution is skewed in the sense that safe state-action pairs are more abundant than unsafe ones. Exploiting such skewness, we can indirectly encourage agents to sample safe transitions by maximising information gain (semi-)metrics that only grow in areas close-enough to training data.

**$\zeta_{\text{LOO}}$ : Leave-One-Out semi-metric** Consider a GP dynamics model,  $\mathcal{P}_{\text{GP}}$ , that is trained on a state-action-successor-state data set  $\mathcal{D} = \{(\tilde{\mathbf{x}}^{(i)}, y^{(i)})\}_{i=1}^{n_1}$  with  $\tilde{\mathbf{x}}^{(i)} = (\mathbf{x}^{(i)}, \mathbf{u}^{(i)})$ .<sup>4</sup> Such a GP induces a posterior allowing us to query predictions on  $n_2$  test points  $\tilde{\mathbf{x}}_* = \{(\mathbf{x}_*^{(j)}, \mathbf{u}_*^{(j)})\}_{j=1}^{n_2}$ . As noted in Section 2.1, the posterior is also Gaussian with the following mean vector and covariance matrix:<sup>5</sup>  $\boldsymbol{\mu}_* = \mathbf{K}_{n_2, n_1} \mathbf{A}_{n_1, n_1} \mathbf{y}_{n_1}$  and  $\boldsymbol{\Sigma}_{n_2, n_2} = \mathbf{K}_{n_2, n_2} - \mathbf{K}_{n_2, n_1} \mathbf{A}_{n_1, n_1} \mathbf{K}_{n_1, n_2}$ . Our goal is to design a measure that increases in regions with dense training-data (due to the usage of CVaR constraint) to aid agents in exploring novel yet safe tuples. To that end, we propose using an expected leave-one-out semi-metric between two Gaussian processes defined, for a one query data point  $\tilde{\mathbf{x}}_*$ , as:  $\zeta_{\text{LOO}}(\tilde{\mathbf{x}}_*) = \mathbb{E}_{i \sim \text{Uniform}[1, n_1]} [\text{KL}(p(\mathbf{f}_* | \mathcal{D}_{-i}) || p(\mathbf{f}_* | \mathcal{D}))]$  with  $\mathcal{D}_{-i}$  being  $\mathcal{D}$  with point  $i$  left-out. Importantly, such a measure will only grow in regions which are close-enough to sampled training data, as posterior mean and covariance of  $p(\mathbf{f}_* | \mathcal{D}_{-i})$  shift by a factor that scales linearly and quadratically, respectively, with the total covariance between  $\tilde{\mathbf{x}}_*$  and  $\mathbf{X}_{-i}$  where  $\mathbf{X}_{-i}$  denotes a feature matrix with the  $i^{\text{th}}$  row removed.<sup>6</sup> In other words, such a semi-metric fulfils our requirement in the sense that if a test query is distant (in distribution) from all training input data, it will achieve low  $\zeta_{\text{LOO}}$  score. Though appealing, computing a full-set of  $\zeta_{\text{LOO}}$  can be highly computationally expensive, of the order of  $\mathcal{O}(n_1^4)$  – computing  $\mathbf{A}_{n_1-i, n_1-i}$  requires  $\mathcal{O}(n_1^3)$  and this has to be repeated  $n_1$  times. A major source contributing to this expense, well-known in GP literature, is related to the need to invert covariance matrices. Rather than following variational approximations (which constitute an interesting direction), we prioritise sample-efficiency and focus on exact GPs. To this end, we exploit the already computed  $\mathbf{A}_{n_1, n_1}$  during the model-learning step and make-use of the matrix inversion lemma [36] to recursively update the mean and covariances of  $p(\mathbf{f}_* | \mathcal{D}_{-i})$  for all  $i$  (see appendix):  $\boldsymbol{\mu}_*^{(i)} = \boldsymbol{\mu}_* - \mathbf{K}_{n_2, n_1} \frac{\mathbf{a}_i^\top \mathbf{a}_i}{a_{i,i}} \mathbf{y}_{n_1}$  and  $\boldsymbol{\Sigma}_{n_2, n_2}^{(i)} = \boldsymbol{\Sigma}_{n_2, n_2} + \mathbf{K}_{n_2, n_1} \frac{\mathbf{a}_i^\top \mathbf{a}_i}{a_{i,i}} \mathbf{K}_{n_1, n_2}$ , with  $\mathbf{a}_i$  being the  $i^{\text{th}}$  row of  $\mathbf{A}_{n_1, n_1}$ . Hence, updating the inverse covariance matrix only requires computing and adding the outer product of the  $i^{\text{th}}$  row  $\mathbf{A}_{n_1, n_1}$ , divided by the  $i^{\text{th}}$  diagonal element. This, in turn, reduces complexity from  $\mathcal{O}(n_1^4)$  to  $\mathcal{O}(n_1^3)$ .

**$\zeta_{\text{Bootstrap}}$ : Bootstrapped symmetric metric** We also experimented with another metric that quantifies posterior sensitivity to bi-partitions of the data as measured by symmetric KL-divergence,<sup>7</sup> averaged over possible bi-partitions:  $\zeta_{\text{Bootstrap}}(\tilde{\mathbf{x}}_*) = \mathbb{E}_{\langle \mathcal{D}_1, \mathcal{D}_2 \rangle} [\text{KL}_{\text{sym}}(p(\mathbf{f}_* | \mathcal{D}_1) || p(\mathbf{f}_* | \mathcal{D}_2))]$ , where  $\langle \mathcal{D}_1, \mathcal{D}_2 \rangle$  is a random bi-partition of the data  $\mathcal{D}$ . In practice, we randomly split the data in half, and do this  $K$  times (where  $K$  is a tuneable hyper-parameter) to get a collection of  $K$  bi-partitions. We then average over that collection. Similar to  $\zeta_{\text{LOO}}$ ,  $\zeta_{\text{Bootstrap}}$  also assigns low importance to query points far

<sup>4</sup>As  $\mathcal{D}$  changes at every outer iteration, we simply concatenate all data in one larger data set; see Algorithm 1.

<sup>5</sup>For clarity, we describe a one-dimensional scenario. We extend to multi-output GPs in our experiments.

<sup>6</sup>Though intuitive, we provide a formal treatment of the reason behind such growth properties in the appendix.

<sup>7</sup>A symmetric KL-divergence between two distributions  $p$ , and  $q$  is defined as:  $(\text{KL}(p||q) + \text{KL}(q||p))/2$ .

from the training inputs, and hence, can be useful for safe-decision making. In our experiments,  $\zeta_{\text{LOO}}$  provided better-behaving exploration strategy, see Section 4.

**Transforming  $\zeta(\tilde{\mathbf{x}}_*)$  to  $\zeta(\boldsymbol{\tau})$**  Both introduced functions are defined in terms of query test points  $\tilde{\mathbf{x}}_*$ . To incorporate in Equation 1, we define trajectory-based expected total information gain as  $\zeta_{\text{LOO}}(\boldsymbol{\tau}) = \sum_{t=0}^T \gamma^t \zeta_{\text{LOO}}(\langle \mathbf{x}_t, \mathbf{u}_t \rangle)$  and  $\zeta_{\text{Bootstrap}}(\boldsymbol{\tau}) = \sum_{t=0}^T \gamma^t \zeta_{\text{Bootstrap}}(\langle \mathbf{x}_t, \mathbf{u}_t \rangle)$ . Interestingly, this characterisation trades off long-term versus short-term information gain similar to how cost trades-off optimal greedy actions versus long-term decisions. In other words, it is not necessarily optimal to seek an action that maximises immediate information gain since such a transition can ultimately drive the agent to unsafe states (i.e., ones that exhibit low  $\zeta(\tilde{\mathbf{x}})$  values). In fact, such horizon-based definitions have also recently been shown to improve modelling of dynamical systems [10, 47]. Of course, our problem is different in the sense that we seek safe policies in a safe decision-making framework, and thus require safely exploring (semi-)metrics.

### 3.3 Solution method

We now provide a solver to the problem in Equation 1. We operate using  $\zeta_{\text{LOO}}(\boldsymbol{\tau})$  and note that our derivations can exactly be repeated for  $\zeta_{\text{Bootstrap}}(\boldsymbol{\tau})$ . Since we maximise exploration, we use  $-\zeta_{\text{LOO}}(\boldsymbol{\tau})$  in the minimisation problem in Equation 1. Effectively, we need to overcome two hurdles for an implementable algorithm. The first, relates to the bi-objective nature of our problem, while the second is concerned with the CVaR constraint that requires a Lagrangian-type solution.

**From Bi- to Single objectives:** We transform the bi-objective problem into a single objective one through a linear combination of  $\mathcal{C}(\boldsymbol{\tau})$  and  $\zeta_{\text{LOO}}(\boldsymbol{\tau})$ . This relaxed yet constrained version is given by  $\min_{\pi} \lambda_{\pi} \mathbb{E}_{\boldsymbol{\tau} \sim p_{\text{surr}}(\boldsymbol{\tau})}[\mathcal{C}(\boldsymbol{\tau})] - (1 - \lambda_{\pi}) \mathbb{E}_{\boldsymbol{\tau} \sim p_{\text{surr}}(\boldsymbol{\tau})}[\zeta_{\text{LOO}}(\boldsymbol{\tau})]$  s.t.  $\text{CVaR}_{\alpha}(\mathcal{L}(\boldsymbol{\tau})) \leq \xi$ ,<sup>8</sup> where  $\lambda_{\pi}$  is a policy dependent weighting. The choice of  $\lambda_{\pi}$ , however, can be difficult as not any arbitrary combination is acceptable as it has to ultimately yield a Pareto-efficient solution. Fortunately, the authors in [45] have demonstrated that a theoretically-grounded choice for such a weighting in a stochastic multi-objective problem is one that produces a common descent direction that points opposite to the minimum-norm vector in the convex hull of the gradients, i.e., one solving:  $\lambda_{\pi}^* = \arg \min_{\lambda_{\pi}} \|\lambda_{\pi} \nabla_{\pi} \mathbb{E}_{\boldsymbol{\tau}}[\mathcal{C}(\boldsymbol{\tau})] - (1 - \lambda_{\pi}) \nabla_{\pi} \mathbb{E}_{\boldsymbol{\tau}}[\zeta_{\text{LOO}}(\boldsymbol{\tau})]\|_2^2$ . Luckily, solving for  $\lambda_{\pi}^*$  is straight-forward and can be encoded using a rule-based strategy that compares gradient norms, see appendix for further details.

**From Constrained to Unconstrained Objectives:** We write an unconstrained problem using a Lagrange multiplier  $\lambda_{\text{CVaR}}$ :  $\min_{\pi} \lambda_{\pi}^* \mathbb{E}_{\boldsymbol{\tau} \sim p_{\text{surr}}(\boldsymbol{\tau})}[\mathcal{C}(\boldsymbol{\tau})] - (1 - \lambda_{\pi}^*) \mathbb{E}_{\boldsymbol{\tau} \sim p_{\text{surr}}(\boldsymbol{\tau})}[\zeta_{\text{LOO}}(\boldsymbol{\tau})] + \lambda_{\text{CVaR}}[\text{CVaR}_{\alpha}(\mathcal{L}(\boldsymbol{\tau})) - \xi]$ . Due to non-convexity of the problem, we cannot assume strong duality holds, so in our experiments, we schedule  $\lambda_{\text{CVaR}}$  proportional to gradients using a technique similar to that in [43] that has proven effective.<sup>9</sup> To solve the above optimisation problem, we first fix  $\nu$  and perform a policy gradient step in  $\pi$ .<sup>10</sup> To minimise the variance in the gradient estimator of  $\mathcal{C}(\boldsymbol{\tau})$  and  $\zeta_{\text{LOO}}(\boldsymbol{\tau})$ , we build two neural network critics that we use as baselines. The first attempts to model the value of the standard cost, while the second learns information gain values. For the CVaR’s gradient, we simply apply policy gradients. As CVaR is non-Markovian, it is difficult to estimate its separate critic. In our experiments, a heuristic where discounted safety losses as unbiased baselines was used and proved effective. In short, our main update equations when using a policy parameterised by a neural network with parameters  $\boldsymbol{\theta}$  can be written as:

$$\boldsymbol{\theta}^{[j][k+1]} = \boldsymbol{\theta}^{[j][k]} - \eta_k \left( \mathbb{E}_{\boldsymbol{\tau}} \left[ \sum_{t \geq 1} \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{u}_t | \mathbf{x}_t) \left( \lambda_{\pi}^* (Q_{\mathcal{C}}(\mathbf{x}_t, \mathbf{u}_t) - V_{\zeta}^{\phi_1}(\mathbf{x}_t)) - (1 - \lambda_{\pi}^*) (Q_{\zeta}(\mathbf{x}_t, \mathbf{u}_t) - V_{\zeta}^{\phi_2}(\mathbf{x}_t)) \right) \right] + \lambda_{\text{CVaR}} \nabla_{\boldsymbol{\theta}} \text{CVaR}_{\alpha}(\mathcal{L}(\boldsymbol{\tau})) \right) \quad (2)$$

where  $\eta_k$  is a learning rate, and  $V_{\zeta}^{\phi_1}(\mathbf{x}_t)$ ,  $V_{\zeta}^{\phi_2}(\mathbf{x}_t)$  are neural network critics with parameters  $\phi_1$  and  $\phi_2$ . We present the main steps in Algorithm 1 and more details in the appendix.

<sup>8</sup>Please note that the negative sign in the linear combination is due to the fact that we used  $-\zeta_{\text{LOO}}(\boldsymbol{\tau})$ .

<sup>9</sup>Note that a primal dual-method as in [15] is not applicable due to non-convexity. In the future, we plan to study approaches from [22] to ease determining  $\lambda_{\text{CVaR}}$ .

<sup>10</sup>We resorted to policy gradients for two reason: 1) cost functions are not necessarily differentiable, and 2) better experimental behaviour when compared to model back-prop especially on OpenAI’s safety gym tasks.

---

**Algorithm 1** SAMBA: Safe Model-Based & Active Reinforcement Learning

---

- 1: **Inputs:**  $\lambda_{\text{CVaR}}$  initialisation, initial random policy  $\pi_1$ ,  $\mathcal{D}_0 = \{\}$ ,  $\alpha$ -quantile, safety threshold  $\xi$
  - 2: **for**  $j = 1 : \#\text{env-iterations}$  **do:**
  - 3:     Sample traces from the real environment using  $\pi_j$ , concatenate all data and update  $\mathcal{P}_{\text{GP}}^{[j]}(\cdot)$
  - 4:     **for**  $k = 1 : \#\text{control-iterations}$  **do:**
  - 5:         Sample traces from  $\mathcal{P}_{\text{GP}}^{[j]}(\cdot)$  and compute  $\zeta_{\text{LOO}}^{[j][k]}(\tau)$  (§ 3.2)
  - 6:         Solve for  $\lambda_{\pi}^{[j][k],*} = \arg \min_{\lambda_{\pi}} \|\lambda_{\pi} \nabla_{\theta} \mathbb{E}_{\tau}[\mathcal{C}(\tau)] - (1 - \lambda_{\pi}) \nabla_{\theta} \mathbb{E}_{\tau}[\zeta_{\text{LOO}}(\tau)]\|_2^2$
  - 7:         Update  $\phi_1^{[j][k]}$  and  $\phi_2^{[j][k]}$  with  $\mathbb{E}_{\tau}[\mathcal{C}(\tau)]$  and  $\mathbb{E}_{\tau}[\zeta_{\text{LOO}}(\tau)]$  as targets, and  $\theta^{[j][k]}$  (§ 3.3)
  - 8:         Set  $\pi_{j+1} = \pi_{\#\text{control-iterations}}$
  - 9: **Output:** Policy  $\pi_{\#\text{env-iterations}}$
- 

## 4 Experiments

We assess SAMBA in terms of both *safe learning* (train) and *safe final policies* (test) on three dynamical systems, two of which are adaptations of standard dynamical systems for MBRL (Safe Pendulum and Safe Cart-Pole Double Pendulum), while the third (Fetch Robot – optimally control end-effector to reach a 3D goal) we adapt from OpenAI’s robotics environments [9]. In each of these tasks, we define unsafe regions as areas in state spaces and design the safety loss (i.e.,  $\mathcal{L}(\tau)$ ) to correspond to the (linearly proportional) distance between the end-effector’s position (when in the hazard region) to the centre of the unsafe region. SAMBA implemented a more stable proximal update of Equation 2 following a similar method to [43]. We compare against algorithms from both model-free and model-based literature. Model-free comparisons against TRPO [42], PPO [43], CPO [1], STRPO (safety-constrained TRPO) [40] and SPPO [40] (safety-constrained PPO) enable us to determine if SAMBA improves upon the following: sample complexities during training; total violations (TV), that is, the total number of timesteps spent inside the unsafe region; total accumulated safety cost (TC).<sup>11</sup> Comparison with model-based solvers (e.g., PlaNet [25], (P)PILCO [19]) sheds light on the importance of our active exploration metrics. It is important to note that when implementing PILCO, we preferred a flexible solution that does not assume moment-matching and specific radial-basis function controllers. Hence, we adapted PILCO to support proximal policy updates, referred to as PPILCO in our experiments, and similarly, SPPIILCO (safety-constrained PPILCO), which also proved effective. As SAMBA introduces exploration components to standard model-based learners, we analysed these independently before combining them and reporting TV and TC <sup>11</sup> (see Table 1). All policies are represented by two-hidden-layer (32 units each) neural networks with *tanh* non-linearities. Each requires under 12 hours of training on an NVIDIA GeForce RTX 2080 Ti which has a power consumption of 225W, yielding an estimated training cost of £ 0.1 per model, per seed. Due to space constraints, all hyper-parameters to reproduce our results can be found in the appendix.

**(Semi-)metrics Component:** Evaluating our (semi-)metrics, we conducted an analysis that reports a 2D projection view of the state space at various intervals in the data-collection process. We compare  $\zeta_{\text{LOO}}(\tau)$  and  $\zeta_{\text{Bootstrap}}(\tau)$  against an entropy-based exploration metric and report the results on two systems in Figure 1. It is clear that both our (semi-)metrics encourage agents to explore safe regions in the state space as opposed to entropy that mostly grows in unsafe regions. Similar results are demonstrated with the Fetch Reach robot task (in the appendix). It is also worth noting that due to the high-dimensional nature of the tasks, visual analysis can only give indications. Still, empirical analysis supports our claim and performance improvements are clear; see Table 1.

**Learning and Evaluation:** Having concluded that our (semi-)metrics indeed provide informative signals to the agent for safe exploration, we then conducted learning and evaluation experiments comparing SAMBA against state-of-the-art methods. Results reported in Table 1 demonstrate that SAMBA reduces the amount of training TC, <sup>11</sup> and samples by orders of magnitude compared to others. Interestingly, during safe evaluation (deploying learnt policy and evaluating performance), we see SAMBA’s safety performance competitive with (if not significantly better than) policies trained for safety in terms of TC and TV. <sup>11</sup> Such results are interesting as SAMBA was never explicitly designed to specifically minimise test TV, <sup>11</sup> but it was still able to acquire significant reductions.

---

<sup>11</sup>Note, we report safe learning process TC, which is the total incurred safety cost throughout all training environment interactions, and safe evaluation TC and TV, which similarly is the total incurred safety cost during evaluation, and the total violations from the timesteps spent inside the unsafe region during evaluation.

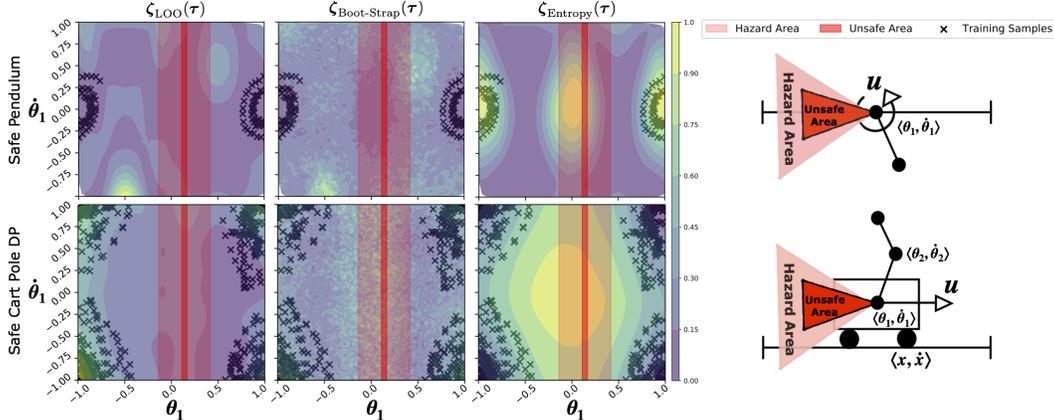


Figure 1: Comparing active learning metric values across the entire state space. Note, policies used within this framework are encouraged to visit regions with higher (yellow) values. The dynamics model for Safe Pendulum was trained on 100 samples and for Safe Cart Pole Double Pendulum on 500 samples respectively. More plots for different numbers of samples can be found in the appendix.

Of course, one might argue that these results do not convey a safe final policy, as violations are still non-zero. We remind, however, that we define safety in terms of CVaR constraints, which do not totally prohibit violations, rather, limit the average cost of excess violations beyond a (user-defined) risk level  $\alpha$ . Indeed, as mentioned above, it is not possible to guarantee total safety without strong assumptions on dynamics and/or an initial safe policy (and of course, none of the algorithms in Table 1 have zero violations).

Table 1: Safe learning and safe evaluation for Safe Pendulum, Safe Cart Pole Double Pendulum, and Safe Fetch Reach. We averaged each policy’s safe evaluation over three random seeds, and collected 10000 evaluation samples per seed. <sup>11</sup>

Learner	Safe Pendulum				Safe Cart Pole Double Pendulum				Safe Fetch Reach			
	Learning		Evaluation		Learning		Evaluation		Learning		Evaluation	
	Samples	TC	TV	TC	Samples	TC	TV	TC	Samples	TC	TV	TC
TRPO	1000	110	4.3	4.9	1000	310	10	23	1000	790	1.8	2.7
PPO	1000	81	3.6	5	1000	140	5.5	35	1000	670	1.6	3.1
PlanNet	100	56	4.5	5.2	40	2.6	7.7	29	100	110	2.3	2.9
PPILCO	2	0.04	2.8	4.5	6	0.097	7.1	33	5	0.49	1.9	3
PlaNet w RS	100	51	3.5	4.2	40	1.3	1.6	2.5	-	-	-	-
CPO	1000	58	4.7	4.4	1000	95	1.1	5	1000	160	1.9	2.7
STRPO	1000	38	2.0	2.1	1000	83	2.1	2.7	1000	170	1.7	2.1
SPPO	1000	65	1.7	2.4	1000	68	1	1.8	1000	290	1.5	2
SPPILCO	1.8	0.02	1.7	2.1	6	0.062	1.2	1.6	5	0.37	1.4	2.2
<b>SAMBA</b>	<b>1.6</b>	<b>0.01</b>	<b>1.5</b>	<b>2.0</b>	<b>5</b>	<b>0.054</b>	<b>0.85</b>	<b>1.4</b>	<b>5</b>	<b>0.23</b>	<b>1.2</b>	<b>1.9</b>

The data is scaled by  $10^3$ . Note: PlaNet w RS on Safe Fetch Reach diverged during training.

To evaluate risk performance, we conducted an in-depth study evaluating all methods on the two most challenging environments, Safe Cart Pole Double Pendulum and Safe Fetch Reach, using cost limits  $\xi = 0.25$  and  $\xi = 0.75$ . Table 2 demonstrates that SAMBA achieves lower safety cost quartiles and lower expected safety cost. Therefore, we conclude that, indeed, SAMBA produces safe-final policies in terms of its objective in Equation 1.

## 5 Conclusion and future work

We proposed SAMBA, a safe and active model-based learner that makes use of GP models and solves a bi-objective constraint problem to trade-off cost minimisation, exploration, and safety. We evaluated our method on three benchmarks, including ones from Open AI’s safety gym and demonstrated significant reduction in training cost and sample complexity, as well as safe-final policies in terms of CVaR constraints compared to the state-of-the-art.

Table 2: Constraint satisfaction (✓) or constraint violation (✗) on Safe Cart Pole Double Pendulum and Safe Fetch of model-based (orange) and model-free (blue) solvers. We evaluated expectation constraints (Exp.) and CVaR. Results show SAMBA satisfied safety constraints, outperforming others in different quartiles. There is also a clear trade-off between returns and safety for all algorithms. Interestingly, SAMBA is safer yet acquiring acceptable expected returns – close to SPPO for instance.

Learner	Safe Cart Pole Double Pendulum						Safe Fetch Reach					
	Quartile			Constraint			Quartile			Constraint		
	0.25	0.5	0.75	Exp.	CVaR $_{\alpha}$	$\mathbb{E}_{\tau}[\mathcal{C}(\tau)]$	0.25	0.5	0.75	Exp.	CVaR $_{\alpha}$	$\mathbb{E}_{\tau}[\mathcal{C}(\tau)]$
PPO	5.0	5.7	9.4	✗	✗	-19	0.95	1.25	1.88	✗	✗	-0.26
PPILCO	5.0	5.8	8.6	✗	✗	-21	0.91	1.22	1.45	✗	✗	-0.36
PlaNet	5.2	6.0	8.8	✗	✗	-21	0.64	1.09	1.88	✗	✗	-0.39
TRPO	5.1	5.9	7.0	✗	✗	-22	0.69	0.79	0.92	✗	✗	-0.58
CPO	0.41	0.47	0.96	✗	✗	-23	0.81	0.92	1.02	✗	✗	-0.43
PlaNet w RS	0.83	0.94	1.11	✗	✗	-24	-	-	-	-	-	-
SPILCO	0.21	0.28	0.33	✗	✗	-25	0.45	0.77	1.13	✓	✗	-0.57
STRPO	0.22	0.27	0.32	✓	✗	-28	0.14	0.46	0.86	✓	✗	-0.98
SPPO	0.19	0.24	0.29	✓	✗	-27	0.27	0.44	0.71	✓	✓	-1.51
<b>SAMBA</b>	<b>0.15</b>	<b>0.21</b>	<b>0.24</b>	✓	✓	-27	<b>0.00</b>	<b>0.05</b>	<b>0.19</b>	✓	✓	-2.27

Note: PlaNet w RS on Safe Fetch Reach diverged during training.

In future, we plan to generalise to variational GPs [17], and to apply our method in real-world robotics. Additionally, we want to study the theoretical guarantees of SAMBA to demonstrate convergence to well-behaving, exploratory, and safe policies.

## Broader Impact

Not applicable.

## References

- [1] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In *International Conference on Machine Learning*, pages 22–31, 2017.
- [2] Anayo K Akametalu, Jaime F Fisac, Jeremy H Gillula, Shahab Kaynama, Melanie N Zeilinger, and Claire J Tomlin. Reachability-based safe learning with gaussian processes. In *IEEE Conference on Decision and Control*, pages 1424–1431, 2014.
- [3] Eitan Altman. *Constrained markov decision processes*, volume 7. CRC Press, 1999.
- [4] Haitham Bou Ammar, Eric Eaton, Paul Ruvolo, and Matthew Taylor. Online multi-task learning for policy gradient methods. In *International Conference on Machine Learning*, pages 1206–1214, 2014.
- [5] Anil Aswani, Humberto Gonzalez, S Shankar Sastry, and Claire Tomlin. Provably safe and robust learning-based model predictive control. *Automatica*, 49(5):1216–1226, 2013.
- [6] Philip Ball, Jack Parker-Holder, Aldo Pacchiano, Krzysztof Choromanski, and Stephen Roberts. Ready policy one: World building through active learning. *arXiv preprint arXiv:2002.02693*, 2020.
- [7] Felix Berkenkamp, Matteo Turchetta, Angela Schoellig, and Andreas Krause. Safe model-based reinforcement learning with stability guarantees. In *Advances in Neural Information Processing Systems*, pages 908–918, 2017.
- [8] Dimitri P Bertsekas. Nonlinear programming. *Journal of the Operational Research Society*, 48(3):334–334, 1997.
- [9] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.

- [10] Mona Buisson-Fenet, Friedrich Solowjow, and Sebastian Trimpe. Actively learning gaussian process dynamics. *arXiv preprint arXiv:1911.09946*, 2019.
- [11] Yinlam Chow and Mohammad Ghavamzadeh. Algorithms for cvar optimization in mdps. In *Advances in Neural Information Processing Systems*, pages 3509–3517, 2014.
- [12] Yinlam Chow, Mohammad Ghavamzadeh, Lucas Janson, and Marco Pavone. Risk-constrained reinforcement learning with percentile risk criteria. *The Journal of Machine Learning Research*, 18(1):6070–6120, 2017.
- [13] Yinlam Chow, Ofir Nachum, Edgar Duenez-Guzman, and Mohammad Ghavamzadeh. A lyapunov-based approach to safe reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 8092–8101, 2018.
- [14] Yinlam Chow, Ofir Nachum, Aleksandra Faust, Mohammad Ghavamzadeh, and Edgar Duenez-Guzman. Lyapunov-based safe policy optimization for continuous control. *arXiv preprint arXiv:1901.10031*, 2019.
- [15] Yinlam Chow, Aviv Tamar, Shie Mannor, and Marco Pavone. Risk-sensitive and robust decision-making: a cvar optimization approach. In *Advances in Neural Information Processing Systems*, pages 1522–1530, 2015.
- [16] Gal Dalal, Krishnamurthy Dvijotham, Matej Vecerik, Todd Hester, Cosmin Paduraru, and Yuval Tassa. Safe exploration in continuous action spaces. *arXiv preprint arXiv:1801.08757*, 2018.
- [17] Andreas Damianou, Michalis K Titsias, and Neil D Lawrence. Variational gaussian process dynamical systems. In *Advances in Neural Information Processing Systems*, pages 2510–2518, 2011.
- [18] Taco de Wolff, Alejandro Cuevas, and Felipe Tobar. Mogptk: The multi-output gaussian process toolkit. *arXiv preprint arXiv:2002.03471*, 2020.
- [19] Marc Deisenroth and Carl E Rasmussen. Pilco: A model-based and data-efficient approach to policy search. In *International Conference on Machine Learning*, pages 465–472, 2011.
- [20] Valerii Vadimovich Fedorov. *Theory of optimal experiments*. Elsevier, 2013.
- [21] Jacob Gardner, Geoff Pleiss, Kilian Q Weinberger, David Bindel, and Andrew G Wilson. Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration. In *Advances in Neural Information Processing Systems*, pages 7576–7586, 2018.
- [22] CJ Goh and XQ Yang. Nonlinear lagrangian theory for nonconvex optimization. *Journal of Optimization Theory and Applications*, 109(1):99–121, 2001.
- [23] David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- [24] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. In *International Conference on Learning Representations*, 2020.
- [25] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In *International Conference on Machine Learning*, 2019.
- [26] Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. In *AAAI Conference on Artificial Intelligence*, 2018.
- [27] Hjmshi. *hjmshi/pytorch-lbfgs*. <https://github.com/hjmshi/PyTorch-LBFGS>, 2018.
- [28] Achin Jain, Truong Nghiem, Manfred Morari, and Rahul Mangharam. Learning and control using gaussian processes. In *ACM/IEEE International Conference on Cyber-Physical Systems*, pages 140–149, 2018.

- [29] Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *International Conference on Machine Learning*, pages 267–274, 2002.
- [30] Sanket Kamthe and Marc Peter Deisenroth. Data-efficient reinforcement learning with probabilistic model predictive control. In *International Conference on Artificial Intelligence and Statistics*, 2018.
- [31] Hassan K Khalil and Jessy W Grizzle. *Nonlinear systems*, volume 3. Prentice hall Upper Saddle River, NJ, 2002.
- [32] Torsten Koller, Felix Berkenkamp, Matteo Turchetta, and Andreas Krause. Learning-based model predictive control for safe exploration. In *IEEE Conference on Decision and Control*, pages 6059–6066, 2018.
- [33] Andreas Krause and Carlos Guestrin. Nonmyopic active learning of gaussian processes: an exploration-exploitation approach. In *International Conference on Machine Learning*, pages 449–456, 2007.
- [34] Andreas Krause, Ajit Singh, and Carlos Guestrin. Near-optimal sensor placements in gaussian processes: Theory, efficient algorithms and empirical studies. *Journal of Machine Learning Research*, 9(Feb):235–284, 2008.
- [35] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [36] KB Petersen, MS Pedersen, et al. The matrix cookbook, vol. 7. *Technical University of Denmark*, 15, 2008.
- [37] LA Prashanth. Policy gradients for cvar-constrained mdps. In *International Conference on Algorithmic Learning Theory*, pages 155–169. Springer, 2014.
- [38] LA Prashanth and Mohammad Ghavamzadeh. Actor-critic algorithms for risk-sensitive mdps. In *Advances in Neural Information Processing Systems*, pages 252–260, 2013.
- [39] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press, 2005.
- [40] Alex Ray, Joshua Achiam, and Dario Amodei. Benchmarking safe exploration in deep reinforcement learning. 2019.
- [41] Rohan Saphal, Balaraman Ravindran, Dheevatsa Mudigere, Sasikanth Avancha, and Bharat Kaul. Seerl: Sample efficient ensemble reinforcement learning. *arXiv preprint arXiv:2001.05209*, 2020.
- [42] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International Conference on Machine Learning*, pages 1889–1897, 2015.
- [43] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [44] Matthias Schultheis, Boris Belousov, Hany Abdulsamad, and Jan Peters. Receding horizon curiosity. In *Conference on Robot Learning*, 2019.
- [45] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. In *Advances in Neural Information Processing Systems*, pages 527–538, 2018.
- [46] Burr Settles. Active learning literature survey. Technical report, University of Wisconsin-Madison Department of Computer Sciences, 2009.
- [47] Pranav Shyam, Wojciech Jaśkowski, and Faustino Gomez. Model-based active exploration. In *International Conference on Machine Learning*, 2019.

- [48] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- [49] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354–359, 2017.
- [50] Aravind Srinivas, Michael Laskin, and Pieter Abbeel. Curl: Contrastive unsupervised representations for reinforcement learning. *arXiv preprint arXiv:2004.04136*, 2020.
- [51] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [52] Matteo Turchetta, Felix Berkenkamp, and Andreas Krause. Safe exploration in finite markov decision processes with gaussian processes. In *Advances in Neural Information Processing Systems*, pages 4312–4320, 2016.
- [53] Joost van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarín Gal. Simple and scalable epistemic uncertainty estimation using a single deep deterministic neural network. *arXiv preprint arXiv:2003.02037*, 2020.
- [54] Akifumi Wachi, Yanan Sui, Yisong Yue, and Masahiro Ono. Safe exploration and optimization of constrained mdps using gaussian processes. In *AAAI Conference on Artificial Intelligence*, 2018.
- [55] Christoph Zimmer, Mona Meister, and Duy Nguyen-Tuong. Safe active learning for time-series modeling with gaussian processes. In *Advances in Neural Information Processing Systems*, pages 2730–2739, 2018.

## A Algorithm

In this section we provide a more detailed description of Algorithm 1. We shall use the following short-hand for the advantage function:  $A_C^{\phi_1}(\mathbf{x}, \mathbf{u}) := Q_C(\mathbf{x}, \mathbf{u}) - V_C^{\phi_1}(\mathbf{x})$ ,  $A_\zeta^{\phi_2}(\mathbf{x}, \mathbf{u}) := Q_\zeta(\mathbf{x}, \mathbf{u}) - V_\zeta^{\phi_2}(\mathbf{x})$ .

---

### Algorithm 2 SAMBA: Safe Model-Based & Active Reinforcement Learning

---

- 1: **Inputs:** Risk-level  $\alpha$ , safety threshold  $\xi$ ,
  - 2: **Hyper-parameters:** Data batch size  $B$ , model traces batch size  $N$ , #env-iterations  $J$ , #control-iterations  $K$
  - 3: **Initialisation:** GP hyper-parameters (§ 2.1),  $\mathcal{D} = \emptyset$ ,  $\lambda_{\text{CVaR}} \geq 0$  arbitrarily, policy parameter  $\theta^{[0][0]}$  and critic parameters  $\phi_1^{[0][0]}$ ,  $\phi_2^{[0][0]}$  randomly
  - 4: **for**  $j = 0 : J - 1$  **do:**
  - 5:   Sample batch of  $B$  new traces  $\mathcal{N}_j = \{\tau_\ell\}_{\ell=1}^B$  from the real environment using  $\theta^{[j][0]}$ ,
  - 6:   Transform the traces  $\mathcal{N}_j$  and append to data set  $\mathcal{D}$ :
  - 7:     **for each**  $\tau_\ell \in \mathcal{N}_j$  **do**
  - 8:       **for each**  $(\mathbf{x}_t, \mathbf{u}_t, \mathbf{x}_{t+1}) \in \tau_\ell$  **do**
  - 9:          $\mathcal{D} = \mathcal{D} \cup \{\langle \tilde{\mathbf{x}}^{(\ell,t)} = (\mathbf{x}_t, \mathbf{u}_t), \mathbf{y}^{(\ell,t)} = \mathbf{x}_{t+1} \rangle\}$  (§ 3.2)
  - 10:   Update  $\mathcal{P}_{\text{GP}}^{[j]}(\cdot)$  using  $\mathcal{D}$  and GP equations (§ 2.1)
  - 11:   **for**  $k = 0 : K - 1$  **do:**
  - 12:     Sample  $N$  traces  $\mathcal{S} = \{\tau_\kappa\}_{\kappa=1}^N$  from  $\mathcal{P}_{\text{GP}}^{[j]}(\cdot)$  using  $\theta^{[j][k]}$
  - 13:     Compute estimates of  $\nabla_{\theta} \mathbb{E}_{\tau}[\mathcal{C}(\tau)]$  and  $\nabla_{\theta} \mathbb{E}_{\tau}[\zeta_{\text{LoO}}(\tau)]$  using  $\mathcal{S}$ :
  - 14:       
$$\nabla_{\theta} \mathbb{E}_{\tau}[\mathcal{C}(\tau)] \approx \frac{1}{N} \sum_{\tau_\kappa \in \mathcal{S}} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_t^{(\kappa)} | \mathbf{x}_t^{(\kappa)}) A_C^{\phi_1^{[j][k]}}(\mathbf{x}_t^{(\kappa)}, \mathbf{u}_t^{(\kappa)})$$
  - 15:       
$$\nabla_{\theta} \mathbb{E}_{\tau}[\zeta_{\text{LoO}}(\tau)] \approx \frac{1}{N} \sum_{\tau_\kappa \in \mathcal{S}} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_t^{(\kappa)} | \mathbf{x}_t^{(\kappa)}) A_{\zeta}^{\phi_2^{[j][k]}}(\mathbf{x}_t^{(\kappa)}, \mathbf{u}_t^{(\kappa)})$$
  - 16:     Solve for  $\lambda_{\pi}^{[j][k],*} = \arg \min_{\lambda_{\pi}} \|\lambda_{\pi} \nabla_{\theta} \mathbb{E}_{\tau}[\mathcal{C}(\tau)] - (1 - \lambda_{\pi}) \nabla_{\theta} \mathbb{E}_{\tau}[\zeta_{\text{LoO}}(\tau)]\|_2^2$ :
  - 17:     **if**  $\nabla_{\theta} \mathbb{E}_{\tau}[\mathcal{C}(\tau)]^{\top} \nabla_{\theta} \mathbb{E}_{\tau}[\zeta_{\text{LoO}}(\tau)] \geq \|\nabla_{\theta} \mathbb{E}_{\tau}[\mathcal{C}(\tau)]\|_2^2$  **then**  $\lambda_{\pi}^{[j][k],*} = 1$
  - 18:     **else if**  $\nabla_{\theta} \mathbb{E}_{\tau}[\mathcal{C}(\tau)]^{\top} \nabla_{\theta} \mathbb{E}_{\tau}[\zeta_{\text{LoO}}(\tau)] \geq \|\nabla_{\theta} \mathbb{E}_{\tau}[\zeta_{\text{LoO}}(\tau)]\|_2^2$  **then**  $\lambda_{\pi}^{[j][k],*} = 0$
  - 19:     **else**
  - 20:       
$$\lambda_{\pi}^{[j][k],*} = \frac{(\nabla_{\theta} \mathbb{E}_{\tau}[\zeta_{\text{LoO}}(\tau)] - \nabla_{\theta} \mathbb{E}_{\tau}[\mathcal{C}(\tau)])^{\top} \nabla_{\theta} \mathbb{E}_{\tau}[\zeta_{\text{LoO}}(\tau)]}{\|\nabla_{\theta} \mathbb{E}_{\tau}[\zeta_{\text{LoO}}(\tau)] - \nabla_{\theta} \mathbb{E}_{\tau}[\mathcal{C}(\tau)]\|_2^2}$$
  - 21:     Update  $\phi_1^{[j][k]}$
  - 22:       
$$\phi_1^{[j][k+1]} = \phi_1^{[j][k]} - \frac{\eta_k^{\phi_1}}{NT} \sum_{\tau_\kappa \in \mathcal{S}} \sum_{t=0}^T \left( V_C^{\phi_1^{[j][k]}}(\mathbf{x}_t^{(\kappa)}) - \sum_{t'=t}^T \gamma^{t'-t} c(\mathbf{x}_{t'}^{(\kappa)}, \mathbf{u}_{t'}^{(\kappa)}) \right) \nabla_{\phi_1} V_C^{\phi_1^{[j][k]}}(\mathbf{x}_t^{(\kappa)})$$
  - 23:     Update  $\phi_2^{[j][k]}$
  - 24:       
$$\phi_2^{[j][k+1]} = \phi_2^{[j][k]} - \frac{\eta_k^{\phi_2}}{NT} \sum_{\tau_\kappa \in \mathcal{S}} \sum_{t=0}^T \left( V_{\zeta}^{\phi_2^{[j][k]}}(\mathbf{x}_t^{(\kappa)}) - \sum_{t'=t}^T \gamma^{t'-t} \zeta_{\text{LoO}}(\mathbf{x}_{t'}^{(\kappa)}, \mathbf{u}_{t'}^{(\kappa)}) \right) \nabla_{\phi_1} V_{\zeta}^{\phi_2^{[j][k]}}(\mathbf{x}_t^{(\kappa)})$$
  - 25:     Update  $\theta^{[j][k]}$  (§ 3.3)
  - 26:       
$$\theta^{[j][k+1]} = \theta^{[j][k]} - \frac{\eta_k^{\theta}}{N} \sum_{\tau_\kappa \in \mathcal{S}} \left[ \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_t^{(\kappa)} | \mathbf{x}_t^{(\kappa)}) \left( \lambda_{\pi}^* A_C^{\phi_1^{[j][k]}}(\mathbf{x}_t^{(\kappa)}, \mathbf{u}_t^{(\kappa)}) \right. \right. \\ \left. \left. - (1 - \lambda_{\pi}^*) A_{\zeta}^{\phi_2^{[j][k]}}(\mathbf{x}_t^{(\kappa)}, \mathbf{u}_t^{(\kappa)}) \right) + \lambda_{\text{CVaR}} \nabla_{\theta} \text{CVaR}_{\alpha}(\mathcal{L}(\tau_{\kappa})) \right]$$
  - 27:   Set  $\theta^{[j+1][0]} = \theta^{[j][K]}$ ,  $\phi_1^{[j+1][0]} = \phi_1^{[j][K]}$ ,  $\phi_2^{[j+1][0]} = \phi_2^{[j][K]}$
  - 28: **Return**  $\theta^{[J][0]}$
-

## B Active metrics

### B.1 Active (semi-)metrics intuition

We implement our code in GPyTorch [21], which computes a Cholesky decomposition of  $\mathbf{A} = \mathbf{L}\mathbf{L}^\top$  during training. The use of the triangular matrix  $\mathbf{L}$  allows the problem to be divided in such a way that memory requirements are reduced, as most matrix multiplications can be pre-computed and re-used when needed. Note also that the diagonal elements of  $\mathbf{A}$  can be readily computed as  $a_{i,i} = \sum_{j=1}^i L_{i,j}^2$ . All these operations considered, the total computational burden of a full-set LOO distribution is reduced from  $\mathcal{O}(N^4)$  to  $\mathcal{O}(N^3)$ . The space complexity of our batch operations, however, can increase due to the need of storing a large tensor. Rather than doing so, we implement a batched version that eliminates such a need and achieves desirable results.

To illustrate what information the LOO semi-metric brings, we analyse under which conditions its value is null. Coming back to our definition, we can write explicitly the KL-divergence from  $p(\mathbf{f}_\star|\mathcal{D}_{-i})$  to  $p(\mathbf{f}_\star|\mathcal{D})$ , for a single test data-point (i.e.,  $n_2 = 1$ ):

$$\text{KL}(p(\mathbf{f}_\star|\mathcal{D}_{-i})||p(\mathbf{f}_\star|\mathcal{D})) = \frac{1}{2} \left( \underbrace{\left( \frac{\mu_\star - \mu_\star^{(i)}}{\sigma_{n_2, n_2}} \right)^2}_{\geq 0} + \underbrace{\left( \frac{\sigma_{n_2, n_2}^{(i)}}{\sigma_{n_2, n_2}} \right)^2 - \log \left( \frac{\sigma_{n_2, n_2}^{(i)}}{\sigma_{n_2, n_2}} \right)^2 - 1}_{\geq 0} \right), \quad (3)$$

where we use non-bold and non-capitalised fonts in order to signify that  $\sigma_{n_2, n_2}^{(i)}$ ,  $\sigma_{n_2, n_2}$ ,  $\mu_\star^{(i)}$  and  $\mu_\star$  are scalars. Note that the second expression is non-negative since  $\log(x) \leq x - 1$  for all  $x$  and the equality is achieved if and only if  $x = 1$ . Therefore for the KL divergence in Equation 3 to be zero, two main conditions need to be met:  $\mu_\star^{(i)} = \mu_\star$  and  $\sigma_{n_2, n_2}^{(i)} = \sigma_{n_2, n_2}$ . Since in general  $\sigma_{n_2, n_2}^{(i)} = \sigma_{n_2, n_2} + \mathbf{K}_{n_2, n_1} \mathbf{a}_i^\top \mathbf{a}_i \mathbf{K}_{n_2, n_1} / a_{i,i}$  we obtain another necessary condition for the KL divergence being equal to zero:  $\mathbf{K}_{n_2, n_1} \mathbf{a}_i^\top = 0$ . Now we will show that  $\mathbf{K}_{n_2, n_1} \mathbf{a}_i^\top = 0$  (and hence  $\sigma_{n_2, n_2}^{(i)} = \sigma_{n_2, n_2}$ ) is also sufficient. Developing the condition  $\mu_\star^{(i)} = \mu_\star$  we have  $\mathbf{K}_{n_2, n_1} \mathbf{a}_i^\top \mathbf{a}_i \mathbf{y}_{n_1} = 0$  since  $\mu_\star^{(i)} = \mu_\star - \mathbf{K}_{n_2, n_1} \mathbf{a}_i^\top \mathbf{a}_i \mathbf{y}_{n_1}$ . As both  $\mathbf{K}_{n_2, n_1} \mathbf{a}_i^\top$  and  $\mathbf{a}_i \mathbf{y}_{n_1}$  are scalars this equivalently amounts to  $\mathbf{K}_{n_2, n_1} \mathbf{a}_i^\top = 0$  or  $\mathbf{a}_i \mathbf{y}_{n_1} = 0$ . Therefore if  $\mathbf{K}_{n_2, n_1} \mathbf{a}_i^\top = 0$  (or equivalently  $\sigma_{n_2, n_2}^{(i)} = \sigma_{n_2, n_2}$ ) then  $\mu_\star = \mu_\star^{(i)}$  and the KL divergence in Equation 3 is equal to zero. In other words, the KL divergence between the two distribution will be zero if and only if the covariance of the test data and the target training data is zero when all other points are observed. As a consequence, the LOO semi-metric can be seen as a proxy to the expected reduction of uncertainty about training output values conditioned on the observation of the queried datapoint.

The intuition behind these conditions is that the expected KL divergence between LOO and the full distributions answers a completely different question when compared to entropy bonus: whereas entropy bonus indicates how much ‘‘absolute’’ uncertainty there is about a *test* output – which we always expect to be high for points situated at a distance of the training inputs – the LOO semi-metric indicates how much, on average, we can learn (i.e. how big would be the reduction in uncertainty would be) about the *training* outputs if we were to observe the queried test data. Therefore, the LOO semi-metric truly answers the question of how much we can learn about the current model when observing the queried point.

### B.2 LOO computation

Without loss of generality, assume that we leave the last sample out, i.e., we will set  $i = n_1$ . We have the following expressions for the means:

$$\begin{aligned} \mu_\star &= \mathbf{K}_{n_2, n_1} \mathbf{A}_{n_1, n_1} \mathbf{y}_{n_1}, \\ \mu_\star^{(n_1)} &= \mathbf{K}_{n_2, n_1-1} \mathbf{A}_{n_1-1, n_1-1} \mathbf{y}_{n_1-1}, \end{aligned} \quad (4)$$

and the covariance matrices:

$$\begin{aligned} \Sigma_{n_2, n_2} &= \mathbf{K}_{n_2, n_2} - \mathbf{K}_{n_2, n_1} \mathbf{A}_{n_1, n_1} \mathbf{K}_{n_1, n_2}, \\ \Sigma_{n_2, n_2}^{(n_1)} &= \mathbf{K}_{n_2, n_2} - \mathbf{K}_{n_2, n_1-1} \mathbf{A}_{n_1-1, n_1-1} \mathbf{K}_{n_1-1, n_2}, \end{aligned} \quad (5)$$

where

$$\begin{aligned}\mathbf{A}_{n_1, n_1} &= [\mathbf{K}_{n_1, n_1} + \sigma_\omega^2 \mathbf{I}]^{-1}, \\ \mathbf{A}_{n_1-1, n_1-1} &= [\mathbf{K}_{n_1-1, n_1-1} + \sigma_\omega^2 \mathbf{I}]^{-1}.\end{aligned}$$

The main difficulty in computing the updates  $\boldsymbol{\mu}_*^{(n_1)} - \boldsymbol{\mu}_*$ ,  $\boldsymbol{\Sigma}_{n_2, n_2}^{(n_1)} - \boldsymbol{\Sigma}_{n_2, n_2}$  is actually dealing with the update involving the matrices  $\mathbf{A}_{n_1-1, n_1-1}$  and  $\mathbf{A}_{n_1, n_1}$ .

Note that there exist  $\mathbf{b}_0, b_1, c_0$  such that:

$$\begin{aligned}\mathbf{K}_{n_2, n_1} &= [\mathbf{K}_{n_2, n_1-1} \quad c_0] \\ \mathbf{K}_{n_1, n_1} &= \begin{bmatrix} \mathbf{K}_{n_1-1, n_1-1} & \mathbf{b}_0^\top \\ \mathbf{b}_0 & b_1 \end{bmatrix}\end{aligned}\quad (6)$$

We have

$$\mathbf{A}_{n_1, n_1} = [\mathbf{K}_{n_1, n_1} + \sigma_\omega^2 \mathbf{I}]^{-1} = \begin{bmatrix} \mathbf{A}_{\neg n_1, \neg n_1} & \mathbf{a}_{n_1, \neg n_1}^\top \\ \mathbf{a}_{n_1, \neg n_1} & a_{n_1, n_1} \end{bmatrix}$$

As  $\mathbf{A}_{n_1, n_1}^{-1} = \mathbf{K}_{n_1, n_1} + \sigma_\omega^2 \mathbf{I}$ , using the block inversion lemma it is straightforward to show that:

$$\mathbf{K}_{n_1-1, n_1-1} + \sigma_\omega^2 \mathbf{I} = [\mathbf{A}_{\neg n_1, \neg n_1} - \mathbf{a}_{n_1, \neg n_1}^\top a_{n_1, n_1}^{-1} \mathbf{a}_{n_1, \neg n_1}]^{-1},$$

and consequently:

$$\mathbf{A}_{n_1-1, n_1-1} = [\mathbf{K}_{n_1-1, n_1-1} + \sigma_\omega^2 \mathbf{I}]^{-1} = \mathbf{A}_{\neg n_1, \neg n_1} - \mathbf{a}_{n_1, \neg n_1}^\top a_{n_1, n_1}^{-1} \mathbf{a}_{n_1, \neg n_1},$$

For simplicity, let us introduce the following matrix:

$$\begin{aligned}\boldsymbol{\Delta} &= \mathbf{A}_{n_1, n_1} - \begin{bmatrix} \mathbf{A}_{n_1-1, n_1-1} & 0 \\ 0 & 0 \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{A}_{\neg n_1, \neg n_1} & \mathbf{a}_{n_1, \neg n_1}^\top \\ \mathbf{a}_{n_1, \neg n_1} & a_{n_1, n_1} \end{bmatrix} - \begin{bmatrix} \mathbf{A}_{\neg n_1, \neg n_1} - \mathbf{a}_{n_1, \neg n_1}^\top a_{n_1, n_1}^{-1} \mathbf{a}_{n_1, \neg n_1} & 0 \\ 0 & 0 \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{a}_{n_1, \neg n_1}^\top a_{n_1, n_1}^{-1} \mathbf{a}_{n_1, \neg n_1} & \mathbf{a}_{n_1, \neg n_1}^\top \\ \mathbf{a}_{n_1, \neg n_1} & a_{n_1, n_1} \end{bmatrix} = \begin{bmatrix} \mathbf{a}_{n_1, \neg n_1}^\top \\ a_{n_1, n_1} \end{bmatrix} [\mathbf{a}_{n_1, \neg n_1} \quad a_{n_1, n_1}] / a_{n_1, n_1} = \frac{\mathbf{a}_{n_1, \neg n_1}^\top \mathbf{a}_{n_1, \neg n_1}}{a_{n_1, n_1}},\end{aligned}$$

where  $\mathbf{a}_{n_1} = [\mathbf{a}_{n_1, \neg n_1} \quad a_{n_1, n_1}]$  is the  $n_1$ -th row of the matrix  $\mathbf{A}_{n_1, n_1}$ .

Now:

$$\begin{aligned}\boldsymbol{\Sigma}_{n_2, n_2}^{(n_1)} - \boldsymbol{\Sigma}_{n_2, n_2} &= \mathbf{K}_{n_2, n_1} \mathbf{A}_{n_1, n_1} \mathbf{K}_{n_1, n_2} - \mathbf{K}_{n_2, n_1-1} \mathbf{A}_{n_1-1, n_1-1} \mathbf{K}_{n_1-1, n_2} = \\ &= [\mathbf{K}_{n_2, n_1-1} \quad c_0] \boldsymbol{\Delta} \begin{bmatrix} \mathbf{K}_{n_1-1, n_2} \\ c_0 \end{bmatrix} = \mathbf{K}_{n_2, n_1} \frac{\mathbf{a}_{n_1, \neg n_1}^\top \mathbf{a}_{n_1, \neg n_1}}{a_{n_1, n_1}} \mathbf{K}_{n_1, n_2}.\end{aligned}$$

Similarly,

$$\begin{aligned}\boldsymbol{\mu}_*^{(n_1)} - \boldsymbol{\mu}_* &= \mathbf{K}_{n_2, n_1-1} \mathbf{A}_{n_1-1, n_1-1} \mathbf{y}_{n_1-1} - \mathbf{K}_{n_2, n_1} \mathbf{A}_{n_1, n_1} \mathbf{y}_{n_1} = \\ &= -\mathbf{K}_{n_2, n_1} \boldsymbol{\Delta} \mathbf{y}_{n_1} = -\mathbf{K}_{n_2, n_1} \frac{\mathbf{a}_{n_1, \neg n_1}^\top \mathbf{a}_{n_1, \neg n_1}}{a_{n_1, n_1}} \mathbf{y}_{n_1}\end{aligned}$$

Hence, we have that

$$\begin{aligned}\boldsymbol{\mu}_*^{(i)} &= \boldsymbol{\mu}_* - \mathbf{K}_{n_2, n_1} \frac{\mathbf{a}_i^\top \mathbf{a}_i}{a_{i, i}} \mathbf{y}_{n_1}, \\ \boldsymbol{\Sigma}_{n_2, n_2}^{(i)} &= \boldsymbol{\Sigma}_{n_2, n_2} + \mathbf{K}_{n_2, n_1} \frac{\mathbf{a}_i^\top \mathbf{a}_i}{a_{i, i}} \mathbf{K}_{n_1, n_2}\end{aligned}$$

for any  $i$ .

Table 3: Experiment hyper-parameters

	Safe Pendulum	Safe Cartpole Double Pendulum	Safe Fetch Reacher	
Dynamics Model	N samples	30	30	300
	N init samples	30	30	1000
	Optimizer	FullBatchLBFGS [27]	FullBatchLBFGS [27]	FullBatchLBFGS [27]
	Predict delta states	True	True	True
	Normalize states	True	True	True
	GP model	Exact	Exact	Exact
	Objective function	Exact Marginal Log Likelihood	Exact Marginal Log Likelihood	Exact Marginal Log Likelihood
	kernel	RBF	RBF	RBF
	CG tolerance	0.00001	0.00001	0.00001
	Max preconditioner size	200	30	150
	Max cg iterations	15000	15000	15000
	Optimization iterations	300	300	100
	Learning rate	0.1	0.1	0.0001
	Agent	N update epochs	80	80
Mini batch size		30000	30000	90000
N neurons		32	32	32
Policy learning rate		3e-4	3e-4	3e-4
Value function learning rate		1e-3	1e-3	1e-3
Penalty learning rate		5e-2	5e-2	5e-2
Clipping		0.2	0.2	0.2
Max gradient norm		0.5	0.5	1.0
Gamma		0.99	0.95	0.99
Lambda		0.97	0.97	0.95
Value function targets		Monte Carlo	Monte Carlo	Monte Carlo
Value function loss fn		Squared loss	Squared loss	Clipped loss
CVaR cost limit		0.025	0.025	0.0166666667
CVaR risk		$\alpha \in (0.5, 1.0)$	$\alpha \in (0.5, 1.0)$	$\alpha \in (0.5, 1.0)$
Runner	Max trajectory length	30	30	30
	N training epochs	100	100	100
	N epoch samples	30000	30000	90000

## C Environmental settings

We implemented a safety cost function with the following settings; The unsafe region started at  $USR_{min} = 20\pi/180$  and ended at  $USR_{max} = 30\pi/180$ , therefore if  $\theta_1$  was in  $USR_{min} \leq \theta_1 \leq USR_{max}$  we would refer to this as a violation and safety cost would be incurred. The Hazard Region contained within  $HZ_{min} = USR_{min} - \pi/4 \leq HZ \leq USR_{max} + \pi/4 = HZ_{max}$ . Where  $\theta_1$  is transformed to always remain in the region  $\theta_1 \in [-\pi, \pi]$ . Therefore, safety cost is linearly proportional to the distance from the edge of the hazard region ( $HZ_{min}$  or  $HZ_{max}$ ) to the centre of the hazard region  $(HZ_{max} + HZ_{min})/2$ . We also implemented batched versions of the pendulum reward function as well as batched versions of the pendulum safety function using torch. We also implement the unsafe region centre as the middle between the fetchers position and its goal (with respect to the  $x$  position only), then the hazard region expands around this point with respect to  $1/4$  of the total distance between the fetchers start state and goal state. Note, we terminated Safe Pendulum once the last 5 rewards in a trajectory were all  $\leq -0.01$ , as this provides an adequately stabilised pendulum.

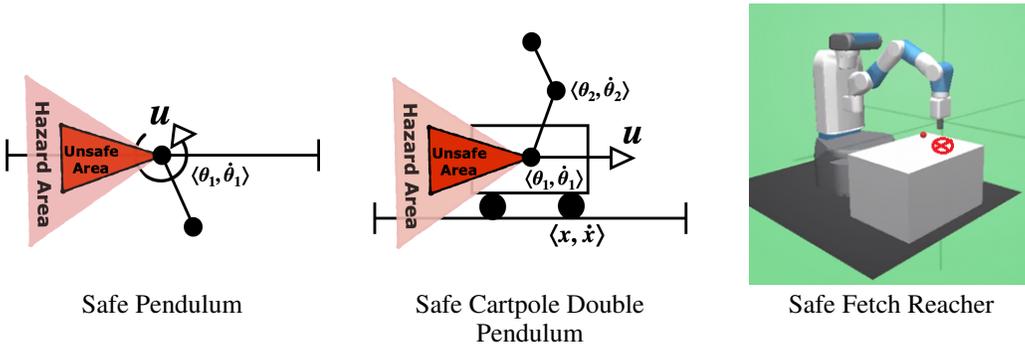


Figure 2: Environments with their respective unsafe regions visualized.

## D Experimental Evaluation

### D.1 LOO performance analysis

This section empirically compares the typical entropy-based method with LOO. For comparison, we collected  $N$  random rollouts from Safe Pendulum to train our GP. Next, we calculated the entropy/LOO metric for unseen random rollouts totaling 20,000 test samples. Figure 3 shows the result of this evaluation, where we compared Entropy vs. LOO for differing sizes of training data  $N$  for the GP ( $N = 20$ ,  $N = 40$  and  $N = 100$ ). The red box contains the USR (Unsafe Region). The experiment suggests that LOO is much more conservative and will aim to guide the policy away from the dangerous area. At the same time, entropy encourages exploration into the unsafe area, additionally encouraged with the increase of  $N$  - a highly undesirable property when safe active exploration is required.

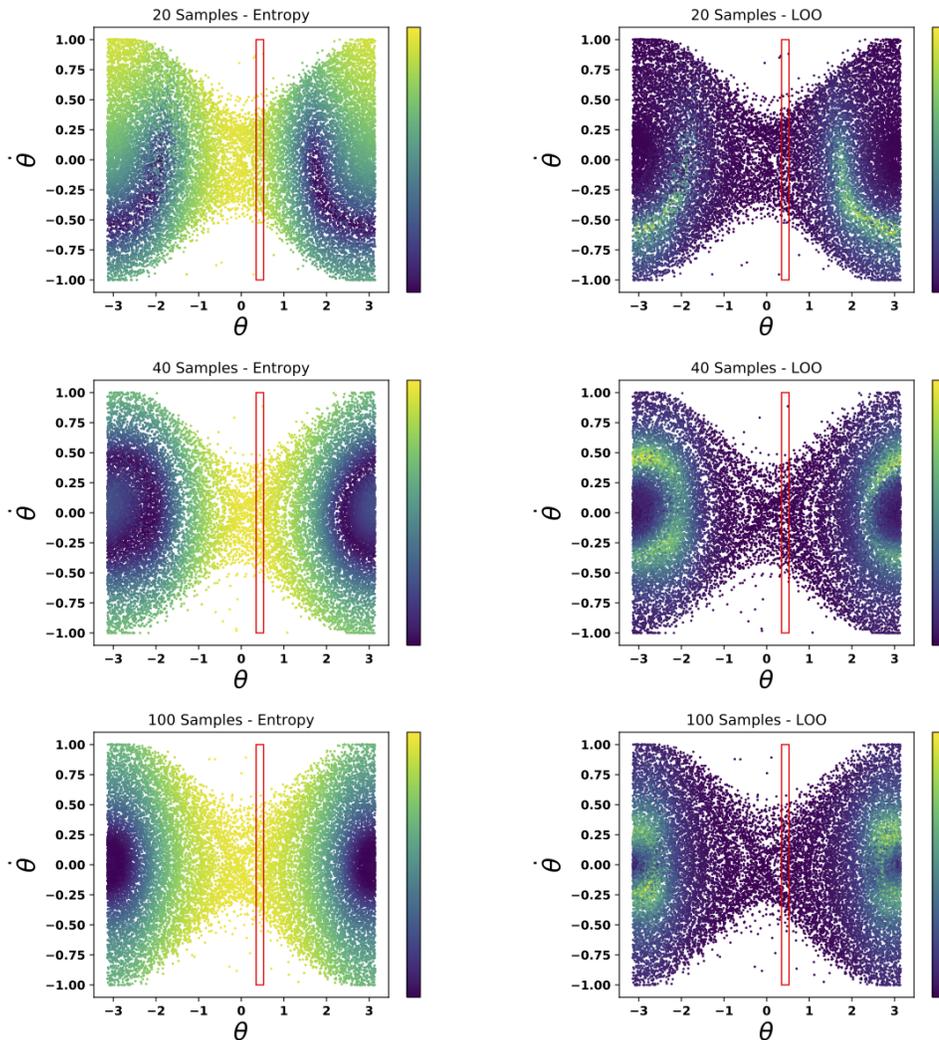


Figure 3: Comparison between Entropy vs. LOO values across the entire state space of Safe Pendulum across true sampled data points. Yellow corresponds to a higher metric values, while blue to a lower. Note that PPO (or any policy used within this framework) is encouraged to visit regions with higher (yellow) weight. The red box indicates the unsafe region of the state space.

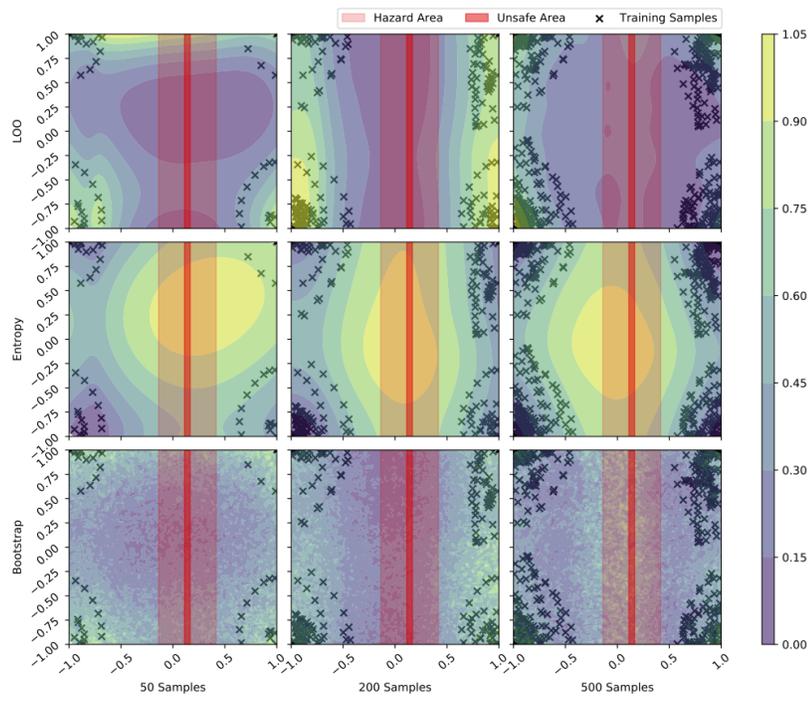
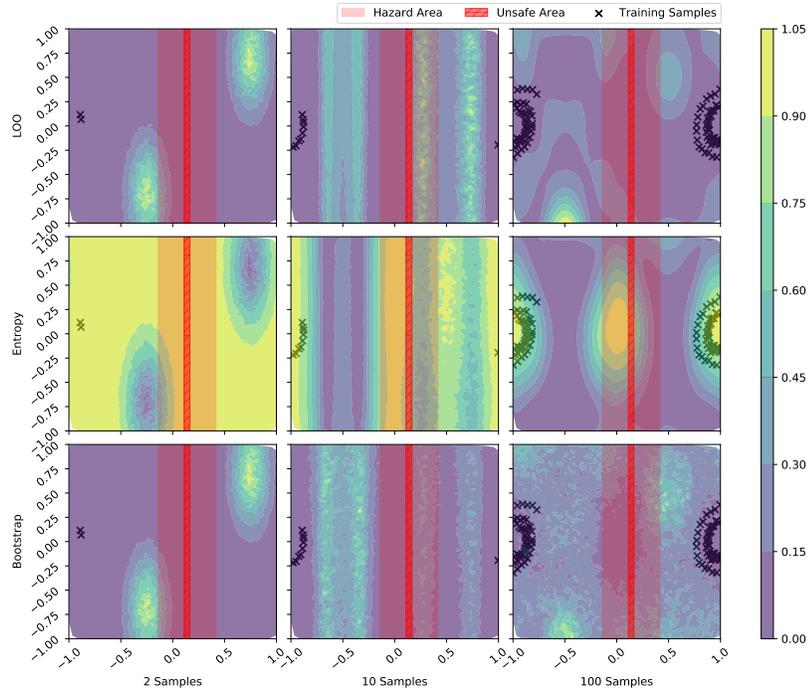


Figure 4: Comparison Entropy vs. LOO values across the entire state space of Safe Pendulum and Safe Cartpole Double Pendulum. Yellow corresponds to a higher metric value, while blue to a lower. Note that PPO (or any policy used within this framework) is encouraged to visit regions with higher (yellow) weight. The red box indicates the unsafe region of the state space.

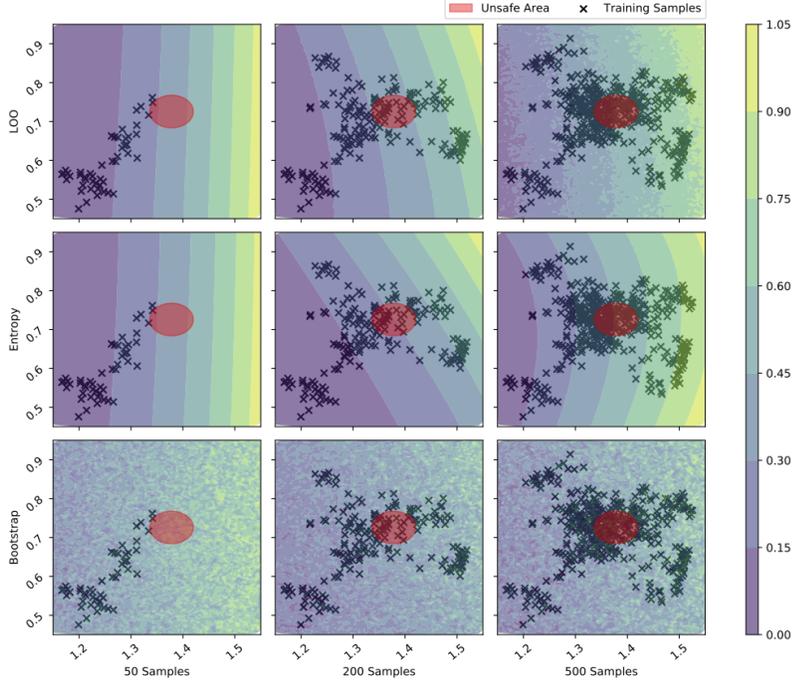


Figure 5: Comparison Entropy vs. LOO values across the entire state space of Safety Fetcher. Yellow corresponds to a higher metric value, while blue to a lower. Note that PPO (or any policy used within this framework) is encouraged to visit regions with higher (yellow) weight. The red box indicates the unsafe region of the state space.

## D.2 Dynamics model analysis

This section presents and compares sample traces between different GP dynamics models trained on a different number of data points (20, 40, and 100 samples). The procedure for the comparison is the following: We rollout the real environment for one full trajectory with random actions and record states and actions. For each GP, we show three sample traces of length 100 time steps (dashed blue line). The sampled traces are compared to the true trajectories by concatenating the true action taken (from the true trajectory) to the open-loop predicted states. We can see in Fig 6 how even when the dynamical system differences to the model, there is still many similarities to the transition dynamics.

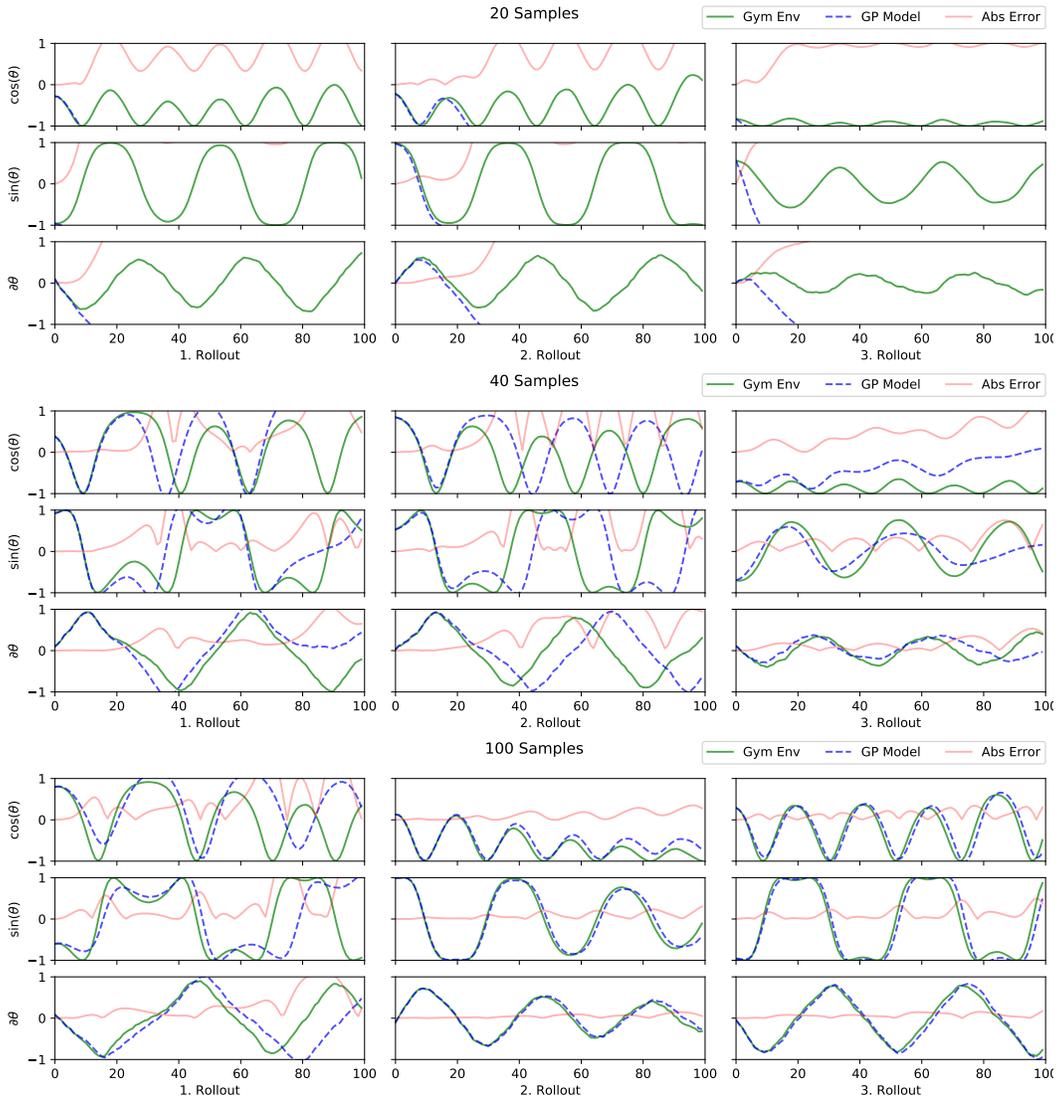


Figure 6: Samples trace from Pendulum.

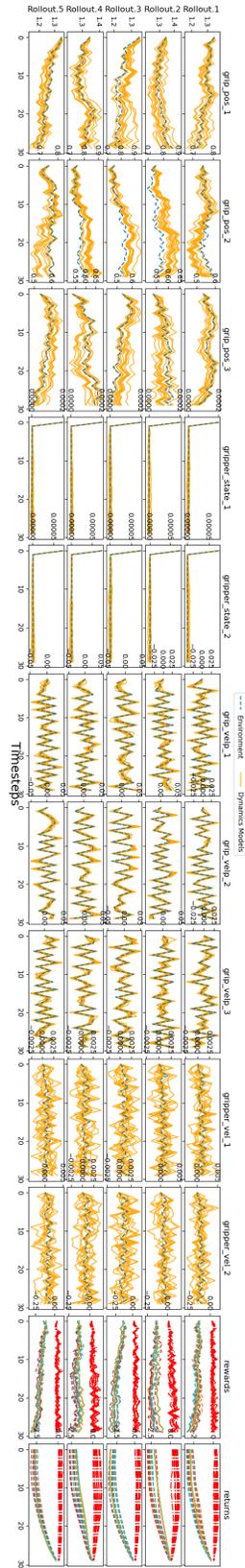


Figure 7: Sampled traces from Safe Fetch Reacher.