# VendiRL: A Framework for Self-Supervised Reinforcement Learning of Diversely Diverse Skills

**Erik M. Lintunen**
Department of Computer Science
Aalto University, Finland
erik.lintunen@aalto.fi

## Abstract

In self-supervised reinforcement learning (RL), one of the key challenges is learning a diverse set of skills to prepare agents for unknown future tasks. Despite impressive advances, scalability and evaluation remain prevalent issues. Regarding scalability, the search for meaningful skills can be obscured by high-dimensional feature spaces, where relevant features may vary across downstream task domains. For evaluating skill diversity, defining what constitutes "diversity" typically requires a hard commitment to a specific notion of what it means for skills to be diverse, potentially leading to inconsistencies in how skill diversity is understood, making results across different approaches hard to compare, and leaving many forms of diversity unexplored. To address these issues, we adopt a measure of sample diversity that translates ideas from ecology to machine learning—the *Vendi Score*—allowing the user to specify and evaluate any desired form of diversity. We demonstrate how this metric facilitates skill evaluation and introduce *VendiRL*, a unified framework for learning diversely diverse sets of skills. Given distinct similarity functions, VendiRL motivates distinct forms of diversity, which could support skill-diversity pretraining in new and richly interactive environments where optimising for various forms of diversity may be desirable.
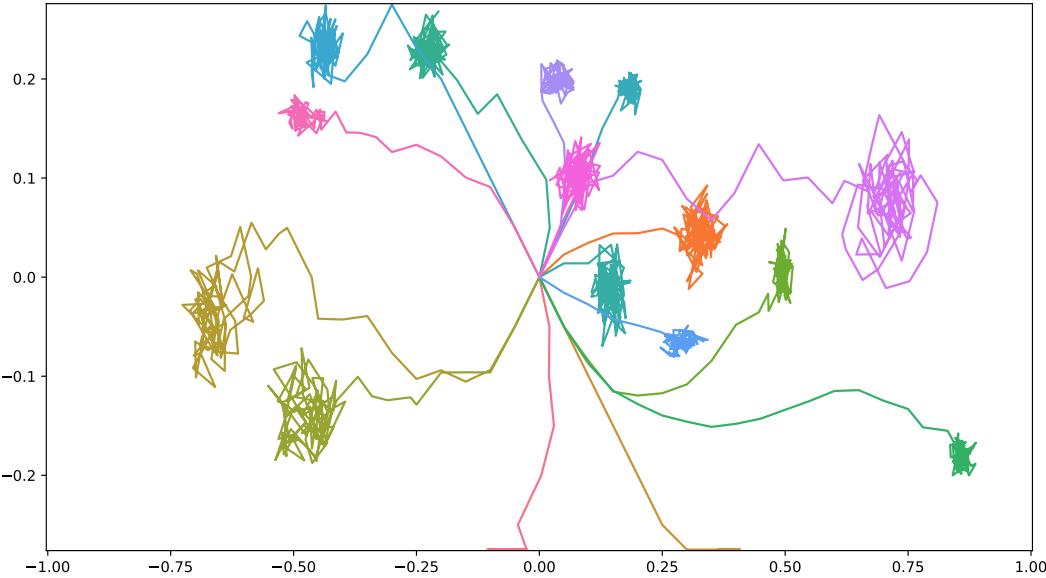
Figure 1: A playful example of *VendiRL*-motivated skills in a 2D environment. An agent optimised its skills for diversity using a linear combination of two distinct skill-similarity measures, illustrating our framework's flexibility in targeting diversity—not only individual notions, but also combinations.

Preprint.

# 1 Introduction

In deep reinforcement learning (RL), it is common to use various pretraining tasks to learn transferable knowledge and skills. Central to learning transferable skills is the challenge of acquiring a *diverse* set of skills. This area of research has been popularised by impressive results showing that diverse skills can empower agents in adapting to a wide range of future tasks (e.g., Gregor et al., 2016, p. 11; Eysenbach et al., 2019, pp. 6–7; Hansen et al., 2020, pp. 6–7; Sharma et al., 2020, pp. 8–10; Laskin et al., 2021, p. 4; Baumli et al., 2021, p. 6737; Shafiullah and Pinto, 2022, pp. 8–9; Yang et al., 2023, pp. 7–8; Park et al., 2024, pp. 8–10). Additionally, some of these successes have been examined theoretically, with recent work suggesting that such approaches can be used to learn an optimal initialisation for unknown reward functions (Eysenbach et al., 2022; cf., Yang et al., 2024) by facilitating the discovery of the ground-truth features of an environment (Reizinger et al., 2025).

Despite these advances, it remains unclear how well existing methods scale in complex and richly interactive environments—both virtual and physical—that are characteristic of many applications. Due to resource constraints, such as training time and compute, these methods often rely on various forms of user supervision. That is, users apply prior knowledge of downstream tasks to more effectively discover meaningfully diverse skills. Useful priors include determining the appropriate number of skills to learn, identifying subsets of the feature space known to facilitate learning a set of target behaviours, and choosing a skill-diversity objective that motivates diversity in a desirable way.

We focus on the latter: defining an objective for diversity often requires a hard commitment to a specific understanding of what it means for skills to be "diverse", such as those that are discriminable (Gregor et al., 2016), costly to transform into one another (He et al., 2022), or temporally distant (Park et al., 2024). Such variance can lead to inconsistencies in how skill diversity is conceptualised and measured—a construct lacking objective evaluation metrics.[1] Moreover, these definitions may fail to represent the wide spectrum of possible interpretations of diversity in the real world, many of which could turn out to be beneficial for modelling and evaluating aspects of open-ended learning.

To this end, we introduce *VendiRL*, a framework inspired by the *Vendi Score* (Friedman and Dieng, 2023)—a measure of sample diversity that connects ideas from ecology to machine learning. The Vendi Score enables the specification and evaluation of any desired form of diversity as defined by a similarity function. Our contributions are three-fold: (1) we adapt the Vendi Score to measuring skill diversity, resulting in an interpretable and intuitive evaluation metric for diversity in open-ended RL; (2) a novel framework for the acquisition of diverse skills through plug-in objectives, specifically via distinct similarity functions and their combinations (as exemplified in Figure 1), enabling the discovery of *diversely* diverse skills under a single reward formalism; and (3) this "pick-and-mix" approach to defining diversity could potentially support the scalability of skill-diversity pretraining to new and richly interactive environments by offering a variety of diversity rewards to choose from.

# 2 Preliminaries

To support the understanding of our work, we introduce the following concepts: reward functions (Section 2.1), skills (Section 2.2), intrinsic rewards (Section 2.3), and the Vendi Score (Section 2.4). Put succinctly, a task for a learning agent is defined by a reward function; skills are conceptualised as a goal-conditioned policy; an agent can self-supervise the process of learning diverse skills by generating its own intrinsic reward signal; and the Vendi Score offers machine learning researchers and practitioners an interpretable, ecologically inspired evaluation metric for any form of diversity.

## 2.1 Tasks are defined by reward functions

Conceptually, at the core of our work is the idea of tasking an agent to learn a diverse set of skills. To this end, we first elaborate on our understanding of a "task" in the context of RL. Central to RL is the *reward function*, mapping from interactions involving states, actions, and successor states to scalar rewards. This construct is anchored in the *reward hypothesis*—a belief strongly influencing the field—suggesting that "all of what we mean by goals and purposes can well be thought of as the maximization of the expected value of the cumulative sum of a received scalar signal (called reward)"

---

[1]The lack of benchmarks and objective evaluation metrics is not unique to skill diversity; it is a prevalent issue in the broader field of open-ended learning. For a discussion on this topic, see Colas et al. (2022, pp. 1167–1169).

(Sutton and Barto, 2018, p. 53; cf., Silver et al., 2021, p. 4). In this sense, a reward function can be understood as defining a *task* for a learning agent (Ng et al., 1999, p. 1).

## 2.2 Skills are policies conditioned on goals

In RL, "skills" can be conceptualised as policies trained to accomplish specific tasks. The challenge of learning a set of multiple skills can be formalised as a goal-augmented Markov decision process (GAMDP) within the framework of goal-conditioned reinforcement learning (GCRL; for a review see Liu et al., 2022). This extends the standard definition of a reward function to be conditioned on goals:

$$r : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \times \mathbb{G} \to \mathbb{R}.$$

Henceforth, $\mathbf{s}, \mathbf{s}' \in \mathbb{S}$ represent a state and successor state, respectively, from the set of possible states, $\mathbf{a} \in \mathbb{A}$ represents an action from the set of possible actions, and $\mathbf{g} \in \mathbb{G}$ represents a *goal* from the set of possible goals. This goal, or *goal-defining variable*, is a parameter to the reward function (cf., Colas et al., 2022, p. 1165; Aubret et al., 2023, p. 6); it indicates which reward function the agent is aiming to maximise. Then, a *skill*, $\pi(\mathbf{a} \mid \mathbf{s}, \mathbf{g})$, is a policy given a goal, optimising for some notion of cumulative reward according to the goal-conditioned reward. In this sense, goals can be viewed as "a set of *constraints* ... that the agent seeks to respect" (Colas et al., 2022, p. 1165, emphasis in original).

While the most immediate intuition of a goal is often as a desired state for the agent to reach (e.g., Kaelbling, 1993, p. 1), the formalism allows for a more general set of constraints on behaviour. In effect, any behaviour that can be defined by attempting to maximise some reward function on the environment can be formulated as a goal–skill pairing. For a concise typology of goal representations in the intrinsically motivated GCRL literature, see, for example, Colas et al. (2022, p. 1171–1174).

## 2.3 Intrinsic rewards enable self-supervised learning

Intrinsically motivated RL deviates from traditional RL by employing *intrinsic* reward functions to generate pseudo-rewards (Lidayan et al., 2025, p. 3). These intrinsic reward functions evaluate agent-internal variables (Oudeyer and Kaplan, 2008, p. 3; cf., Berlyne, 1965, p. 246 and Oudeyer and Kaplan, 2007), which enables their applicability across a diverse range of environments. Various intrinsic rewards have been explored, such as novelty, learning progress, and empowerment, primarily due to their effectiveness in supporting open-ended development, task-agnostic learning, and the ability to deal with sparse rewards (Colas et al., 2022, p. 1161). Models of intrinsic motivation have also been used to explain various aspects of "wet" (biological) RL (e.g., Brändle et al., 2023; Molinaro et al., 2024; Modirshanechi et al., 2025), results potentially applicable for the development of truly open-ended and human-like RL agents. For reviews of intrinsically motivated RL, see Oudeyer and Kaplan (2007); Linke et al. (2020); Colas et al. (2022); Aubret et al. (2023); Lidayan et al. (2025).

A typical feature of intrinsic rewards is their non-stationary nature, which induces a partially observable Markov decision process (POMDP) where the dynamics of the reward distribution are unobserved by the agent. For example, consider rewards computed as a function of the parameters, $\phi$, of a learned neural network; then, the reward $r(\mathbf{s}, \mathbf{a}, \mathbf{s}', \mathbf{g}, \phi_t)$ is likely to differ from $r(\mathbf{s}, \mathbf{a}, \mathbf{s}', \mathbf{g}, \phi_{t+1})$. This characteristic is also prevalent in approaches for learning diverse skills, exemplified in Appendix D.1.

## 2.4 A measure of sample diversity inspired by ecology: the Vendi Score

A long line of work in ecology concerns a fundamental conceptual problem: how can diversity be quantified in a meaningful way? Some interpretations emphasise the importance of *abundance* while others the importance of *balance* (Leinster, 2021, pp. 4–5). One choice is to quantify how *different* the species in a community are. The Vendi Score (Friedman and Dieng, 2023) bridges this understanding to machine learning, representing "the effective number of dissimilar elements in a sample" (p. 6).

The Vendi Score is defined as the exponential of the Shannon entropy of the eigenvalues of a kernel (sample-similarity) matrix. A useful feature of the quantity is that it is interpretable: zero entropy results in the effective number 1 ($\Leftrightarrow$ all elements are equal) and maximum entropy returns the number of elements $n$ ($\Leftrightarrow$ all elements are effectively unique). Formally (Friedman and Dieng, 2023, pp. 5–6):

**Definition 1 (Vendi Score).** Let $\mathbf{x}_1, \cdots, \mathbf{x}_n \in \mathbb{X}$ denote a collection of $n$ samples, let $k : \mathbb{X} \times \mathbb{X} \to \mathbb{R}$ be a positive semidefinite similarity function, with $k(\mathbf{x}, \mathbf{x}) = 1$ for all $\mathbf{x}$, and let $\mathbf{K} \in \mathbb{R}^{n \times n}$ denote a kernel matrix with entry $K_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$. Denote by $\lambda_1, \cdots, \lambda_n$ the eigenvalues of $\mathbf{K}/n$. The

Vendi Score ($VS_k$) is defined as the exponential of the Shannon entropy of the eigenvalues of $\mathbf{K}/n$:

$$VS_k(\mathbf{x}_1, \cdots, \mathbf{x}_n) := \exp\left(-\sum_{i=1}^{n} \lambda_i \log \lambda_i\right), \tag{1}$$

following convention $0 \log 0 = 0$. The Vendi Score is computed from the eigenvalues of $\mathbf{K}/n$, instead of $\mathbf{K}$, so that its eigenvalues sum to one, and therefore the Shannon entropy is well-defined.

## 3 VendiRL

First, we adapt the Vendi Score to measuring skill diversity (Section 3.1). Then, we introduce our framework for learning diversely diverse skills (Section 3.2). Lastly, we demonstrate that VendiRL can motivate different forms of diversity in skills according to a variable similarity function (Section 3.3).

### 3.1 A new way of thinking about skill diversity

**Definition 2 (Effective number of unique skills).** We denote skills as specific configurations of policy parameters: let $\boldsymbol{\theta}_i$ represent a particular skill corresponding to goal $i$, a specific configuration of the skills being learned, $\boldsymbol{\theta}$, from the set of all possible skills $\Theta$. A user-specified $k : \Theta \times \Theta \to \mathbb{R}$ defines a similarity function according to which two skills are compared, such that $k(\boldsymbol{\theta}_i, \boldsymbol{\theta}_i) = 1$ for all $i \in \{1, \cdots, n\}$, where $n$ represents the number of skills being learned, and $\mathbf{K} \in \mathbb{R}^{n \times n}$ denotes a kernel matrix with entry $K_{i,j} = k(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j)$. Then, $VS_{k^t}(\boldsymbol{\theta}_1, \cdots, \boldsymbol{\theta}_n)$ represents the diversity of an agent's skills at a particular point in time, $t$. We call this quantity *the effective number of unique skills*.

In our experiments, two skills are compared based on the trajectories of observations they induce. Depending on the choice of $k$, we may use either full trajectories or summary statistics. For instance, the mean of a trajectory over time captures information about the region of observation space the trajectory predominantly occupies, while the determinant of a trajectory's covariance matrix captures information about the trajectory's volume in the observation space. In Figure 2, we used full trajectories, using a function $k$ that estimates *overlap* of skills in the observation space. This function estimates the $F_1$ score using a method originally intended for the finite approximation of data manifolds in the generative modelling literature (Kynkäänniemi et al., 2019), which we adapted to RL for approximating "skill manifolds" (see Appendix A). Other choices of $k$ are explored in Section 3.3, where we provide examples of the sets of diverse skills they can motivate using VendiRL.



(a) Random skills (no training).  (b) Skills trained with MISL (see Section 4.1).
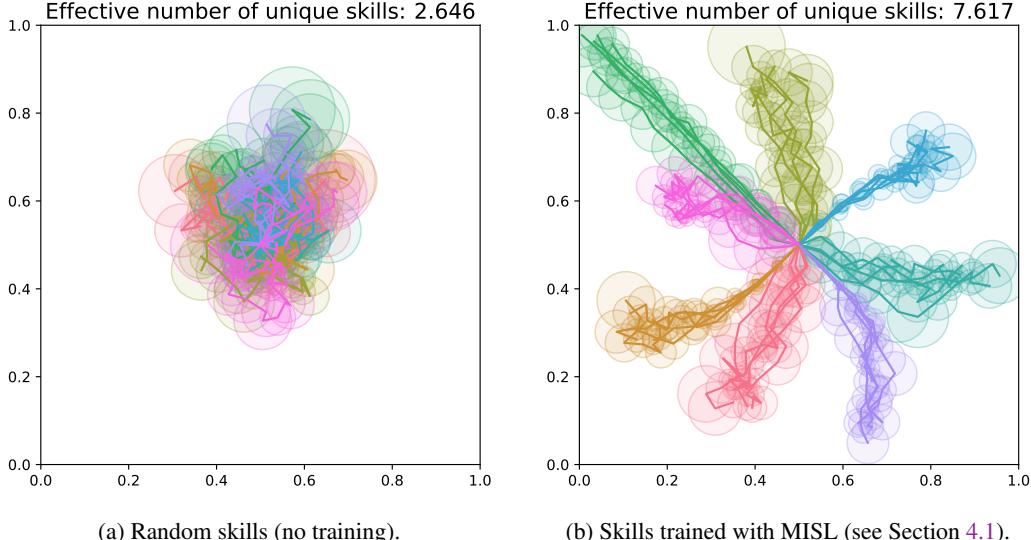
Figure 2: Quantifying skill diversity with the Vendi Score; $F_1$ score is used as the measure of similarity. The similarity function estimates overlap between two skills in the feature space (see Appendix A). Each of the eight skills is represented by five i.i.d. trajectories, and skill are differentiated by colour.

Figure 2 illustrates the use of the Vendi Score in measuring the diversity of two distinct sets of skills. Random skills (Figure 2a) result in extensive overlap when rolled out, so the effective number of unique skills is low (2.646). In contrast, training skills with mutual information skill learning (MISL; see Section 4.1) results in skills that are discriminable in the observation space—clearly reflected in the Vendi Score (7.617), close to the maximum of eight (the agent has learned a set of eight skills).

Because the user can specify the similarity function used to compare skills, the Vendi Score lends itself to measuring any form of diversity that can be described mathematically. This flexibility is also fundamental to our approach for motivating diversely diverse skills, which we elaborate on next.

## 3.2 How VendiRL works



**Algorithm 1:** VendiRL
Define: similarity function $k$, number of skills $n$.
Initialise: skills $\boldsymbol{\theta}$, skill memory $\mathbf{M}$, kernel matrix $\mathbf{K}$.
**while** *training* **do**
    Refill $\mathbf{M}$ by generating skill-trajectories from $\boldsymbol{\theta}$.
    Update $\mathbf{K}$ with $k$ (for all $n$ skills).
    Sample goal uniformly at random: $\mathbf{g} \sim \mathcal{U}\{0, n-1\}$.
    **for** *steps in epoch* **do**
        Sample action: $\mathbf{a}_t \sim \pi_{\boldsymbol{\theta}}(\mathbf{a}_t \mid \mathbf{s}_t, \mathbf{g})$.
        Step environment: $\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)$.
        Store observation in $\mathbf{M}$ (for current skill).
        Update $\mathbf{K}$ with $k$ (for current skill).
        Compute reward: $r_t := \exp\left(-\sum_{i=1}^n \lambda_i \log \lambda_i\right)$.
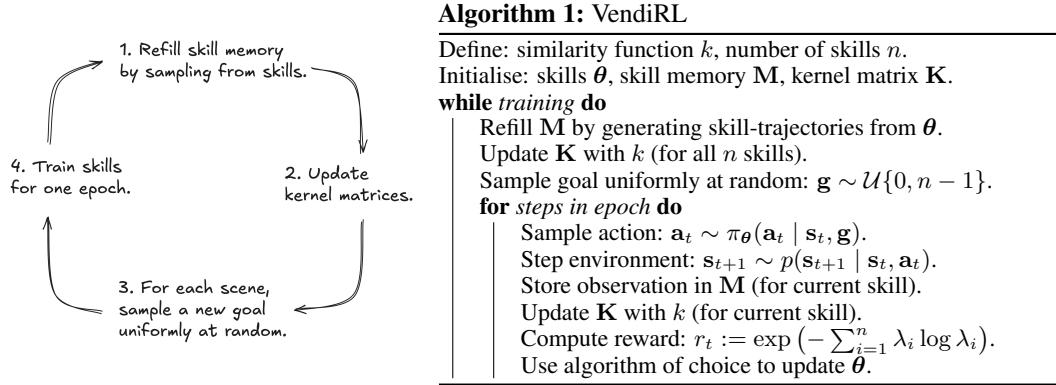        Use algorithm of choice to update $\boldsymbol{\theta}$.

Figure 3: A schematic illustrating the high-level training loop (left) and our algorithm (right).

The high-level training loop for an agent is as follows. Following Definition 2 (Section 3.1), we fix a similarity function $k$ and a number of skills to learn $n$. Before any learning takes place, the agent generates a trajectory from each of its skills with the randomly initialised parameters $\boldsymbol{\theta}$ and stores the observations in a skill memory $\mathbf{M}$. Following this, the agent computes pairwise similarities between its skills using $k$, which results in the kernel matrix $\mathbf{K}$. Then, the agent selects a goal, $\mathbf{g}$, uniformly at random, following and updating the corresponding skill, $\boldsymbol{\theta}_{\mathbf{g}}$, for a fixed number of steps.
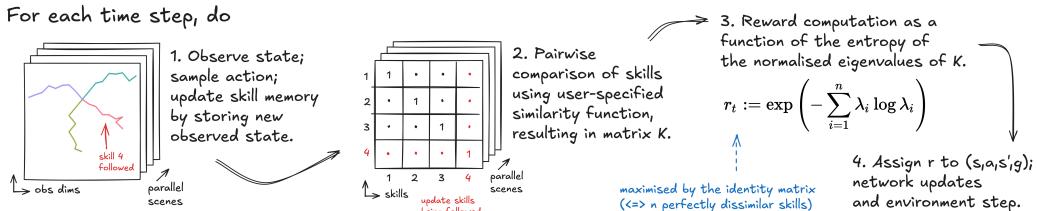


Figure 4: A schematic illustrating a single *VendiRL* training step within a parallelised training setup.

At each time step, the agent observes the current state, samples an action from its skill, and steps in the environment according to the environment's dynamics. Observations are stored in the skill memory $\mathbf{M}$—which, in our experiments, is implemented as a queue with a maximum length equal to one episode. This design allows the agent to track the most recent trajectory induced by each skill. Consequently, the agent always replaces the observation at time $t$ from the previous episode with the observation at time $t$ from the current episode. After storing the observation, the skill-similarity values in the kernel matrix $\mathbf{K}$ are updated for the elements corresponding to the skill currently being followed, $\boldsymbol{\theta}_{\mathbf{g}}$. Specifically, the skill is compared pairwise with all other skills, $\{\boldsymbol{\theta}_1, \cdots, \boldsymbol{\theta}_n\} \setminus \{\boldsymbol{\theta}_{\mathbf{g}}\}$. Following this update, the agent computes its goal-conditioned transition reward as defined below.

**Definition 3 (VendiRL reward).** Given the kernel matrix, $\mathbf{K} \in \mathbb{R}^{n \times n}$, consisting of pairwise similarity measurements of $n$ skills, $\boldsymbol{\theta}_1, \cdots, \boldsymbol{\theta}_n$, at time $t+1$, such that $K_{i,j} = k(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j)$, the goal-conditioned transition reward is defined as the exponential of the Shannon entropy of the eigenvalues

of $\mathbf{K}/n$ (following Equation 1, Section 2.4). That is, the Vendi Score, $VS_k$, at time $t+1$:

$$r_t(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}, \mathbf{g}_t) := VS_{k^{t+1}}(\boldsymbol{\theta}_1, \cdots, \boldsymbol{\theta}_n). \tag{2}$$

Empirically, we found that transforming the reward worked well with some RL algorithms. Alternative formulations of the reward function, such as a *time derivative* and *penalty* are defined in Appendix B.

After computing the transition reward, skills ($\boldsymbol{\theta}$) are updated according to an optimiser of choice. This motivates the agent to learn diverse skills according to $k$, thus driving a variable form of diversity. After the skill has been followed for a fixed number of steps, the skill memory is refilled with new trajectories generated from $\boldsymbol{\theta}$,[2] the kernel matrix updated, a new goal selected, and the process repeated. Our algorithm can be found in Figure 3, and we illustrate a single training step in Figure 4.
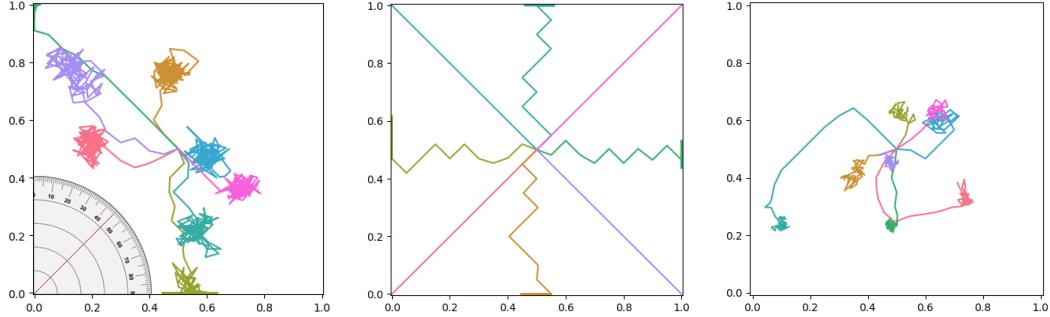
## 3.3 What skills are learned



Figure 5: Examples of *VendiRL*-motivated skills in a $[0, 1]^2$-bounded environment. In each case, an agent has learned eight skills that are differentiated by colour. Similarity between skills is measured as a function of (from left to right): cosine similarity, maximum mean discrepancy, covariance structure.

To showcase the capabilities of VendiRL, we tested the framework using several distinct similarity functions. As a proof of concept, we analysed the skills learned by agents in a simple 2D environment without extrinsic rewards, meaning that our agents were exclusively motivated by VendiRL rewards.

Figure 5 presents results from training three distinct sets of eight skills, each utilising a different similarity function. The first function (left) is defined as cosine similarity between trajectory means, which drives the skills to angular separation around the origin. The second function (middle) employs maximum mean discrepancy with a linear kernel, encouraging the skills to push away from one another.[3] The third function (right) measures similarity in covariance structure by computing absolute differences in the determinants of the skills' covariance matrices, thereby driving diversity in how widely the skills spread. These similarity functions are described mathematically in Appendix C.

Beyond inducing distinct patterns of skill variation, Figure 1 shows that VendiRL can target multiple notions of diversity simultaneously. The 16 skills shown were trained using a similarity function that linearly combines cosine similarity and similarity in covariance structure. This compositionality makes VendiRL a practical tool for designing composite diversity objectives and evaluation metrics.

## 3.4 Challenges in scaling up

We follow Freeman et al. (2021) by leveraging the auto-vectorisation capabilities of JAX (Bradbury et al., 2018) to parallelise data collection across a batch of independent environment scenes. This changes the learning dynamics in several ways. For instance, because goals are sampled independently per scene, an agent can pursue different skills in different scenes (and the same skill in multiple scenes). Consequently, a single training epoch can expose the agent to a broad set of transitions spanning many skills, possibly including multiples of each skill that vary in state-action coverage.

---

[2]Refilling the skill memory is a step taken to synchronise reward distributions across parallel scenes of an environment, a design choice discussed in Section 3.4.

[3]This form of diversity is similar to that targeted by many state-of-the-art skill-diversity methods, visualised for comparison in Appendix D.2, Figure 6.

Parallelisation massively accelerates the acquisition of diverse skills, but it introduces a challenge: different scenes induce different reward distributions. To understand why this is the case and how we address it, we turn to the mechanisms underlying reward computation.

Each scene is allocated its own skill memory $\mathbf{M}$ and kernel matrix $\mathbf{K}$ rather than pooling them across scenes. Pooling complicates credit assignment because: (1) the same skill can be active in multiple scenes simultaneously; (2) transition rewards depend on the current contents of the skill memory (i.e., the diversity among stored observations); and (3) skills are stochastic, so the same skill may visit different regions of the state space across scenes, especially early in training when the skills are near-random. With a pooled memory, multiple observations induced by the same skill could be inserted into the memory at the same step, and the resulting reward would be broadcast to all those transitions, regardless of whether the underlying $(\mathbf{s}, \mathbf{a}, \mathbf{s}', \cdot)$ are comparable.

Using independent per-scene memories resolves this credit-assignment issue but creates a new effect: a skill may be reinforced to increase diversity *locally* in one scene while inadvertently encouraging similarity in another. In other words, distinct memories induce distinct reward distributions. To mitigate drift between reward distributions, the agent periodically synchronises memories: between epochs, each scene's memory is refilled with new trajectories generated from the shared skills $\boldsymbol{\theta}$.

More broadly, we view the design of $\mathbf{M}$ as central to scaling up the task of learning diverse skills. Alternative parametrisations—for example, running averages of skills or learned skill representations against which new experience is compared—could amortise memory costs and thereby improve scaling. Likewise, varying how skill information is shared across scenes provides a natural axis for experimentation, allowing us to study the trade-offs between reward stability and training throughput.

### 3.5 Summary of contributions

Our work advances open-ended RL in three main ways:

- Unified, domain-agnostic evaluation: we adapt the Vendi Score to RL, enabling comparable measures of skill diversity across algorithms, environments, and representations;
- Plug-in diversity objectives: a similarity function directly defines the reward, thus pursuing a different form of diversity does not require changing the system's architecture; and
- Pick-and-mix diversity: because the framework supports composable forms of diversity, it is a versatile tool for training and evaluating many forms of skill diversity.

## 4 Related work

Recent progress in self-supervised learning of diverse skills is heavily skewed towards a single methodological family; accordingly, we centre our review on that literature, then survey approaches to evaluating skill diversity, and close with related work outside the dominant paradigm.

### 4.1 Mutual information skill learning

Mutual information skill learning (MISL) is a well-studied subset of those intrinsically motivated skill acquisition methods focused on learning diverse sets of skills. In MISL, an agent is tasked to maximise behavioural mutual information (BMI), that is, the mutual information between some representation of actions and some representation of states following those actions (Choi et al., 2021, p. 3). Most often, this is done between skills (macro-actions) and some function of the agent's (skill-induced) trajectory, for example, the latest observation of the state. To make the task tractable, most methods employ some form of *variational* MISL; in this context, BMI is approximated using a variational lower bound (e.g., Barber and Agakov, 2003, p. 2), and the problem of learning a diverse set of skills is formulated as a cooperative game between an actor and a learned discriminator model.[4]

Choi et al. (2021, p. 3) provides an overview of variational approaches to MISL, also known as variational empowerment maximisation (for empowerment see Klyubin et al., 2005). Many different objectives for variational MISL have been proposed; examples of relevant empirical work include

---

[4]The formal details underlying variational MISL are given in Appendix D.1. For the interested reader, Eysenbach (2025) provides an accessible tutorial on the topic.

Gregor et al. (2016), Warde-Farley et al. (2019), Eysenbach et al. (2019), Hansen et al. (2020), Sharma et al. (2020), Laskin et al. (2021), Baumli et al. (2021), Yang et al. (2023), and Zheng et al. (2025). The works of Eysenbach et al. (2022) and Reizinger et al. (2025) explain how MISL can be used to learn a "universal" representation for solving downstream tasks (Reizinger et al., 2025, p. 5).

One issue with scaling MISL is its reliance on a well-constructed feature space.[5] Without it, agents may end up learning trivially diverse skills that neither result in observably distinct behaviours in the environment nor transfer effectively to downstream tasks—rendering the pretraining less worthwhile. Although VendiRL does not fully resolve the issue, it alleviates it by enabling the integration of an additional layer of prior knowledge. Rather than partitioning the feature space into uniform subregions, VendiRL can target variable notions of diversity better suited to a downstream task domain.

## 4.2 Evaluating skill diversity

The dominant way to evaluate skill-diversity methods is via their utility: measure how well the skills learned during pretraining transfer to a selection of downstream tasks. For example, some methods are evaluated by fine-tuning a set of learned skills on a set of downstream tasks (e.g., Gregor et al., 2016, p. 11; Eysenbach et al., 2019, pp. 6–7; Hansen et al., 2020, pp. 6–7; Laskin et al., 2021, p. 4; Yang et al., 2023, pp. 7–8), and others using a hierarchical controller to control the repertoire of learned skills on a set of downstream tasks (e.g., Eysenbach et al., 2019, pp. 7–8, Sharma et al., 2020, pp. 9–10; Baumli et al., 2021, p. 6737; Shafiullah and Pinto, 2022, p. 8; Park et al., 2024, pp. 8–10).

Beyond utility, other common approaches for evaluation include qualitative analyses visualising skill trajectories (e.g., Eysenbach et al., 2019, p. 5; Sharma et al., 2020, pp. 6–8; Gu et al., 2021, Appendix A, pp. 19–23; Park et al., 2024, pp. 7–8), measuring state space coverage (e.g., Zheng et al., 2024, p. 8; Park et al., 2024, pp. 8–9), and scoring on an agent on the diversity reward used for training (e.g., Gregor et al., 2016, p. 8, Eysenbach et al., 2019, Appendix D.1, p. 15; Choi et al., 2021, pp. 6–8).

Other proposed proxies include: the effective number of skills (Eysenbach et al., 2019, p. 6; Appendix D, pp. 17–18), measured as the exponential of the Shannon entropy of *the goal-selection policy*, resulting in a quantity representing the effective number of skills *being considered for selection by an agent at a given time*; Latent Goal Reaching (e.g., Choi et al., 2021, pp. 6–8; Gu et al., 2021, Appendix B, p. 25); distributions over task reward (Eysenbach et al., 2019, Appendix D.3, p. 16); skill-separability metrics (e.g., Yang et al., 2024, pp. 5, 7); per-dimension goal achievement (Warde-Farley et al., 2019, pp. 7–10); Diffusion Time (e.g., Machado et al., 2017, Appendix B, p. 13); System Neural Diversity (Bettini et al., 2023); and Quality-Diversity (QD) scores (see Pugh et al., 2016, p. 9).

While each measure usefully captures some facet of skill diversity, most make a hard commitment to a single notion and thus do not support evaluating variable forms of diversity. With the exception of utility-based and QD-like measures, they rarely enable the user to specify what "diversity" should mean for a particular domain. Several evaluation metrics can, in principle, reflect different forms of diversity by changing the feature representation, but in practice the representation is constrained by the domain and determined by what works well for learning the task. In contrast, our use of the Vendi Score brings a unified, domain-agnostic notion of skill diversity to RL, enabling consistent evaluation of behavioural diversity under a variable similarity function specified by the user.

## 4.3 Other approaches for learning diverse skills

**Regularity-based rewards** Sancaktar et al. (2023) use regularity as an intrinsic reward. In their approach, the state space is factorised into objects, mapped to a multiset of symbols (e.g., discretised positions, colours). Varying this mapping induces distinct types of structured behaviour by changing what is considered "regular", though the aim is not explicitly to support variable notions of diversity.

**Automatic curricula** OpenAI et al. (2021) use asymmetric self-play in a goal-conditioned setting: one agent aims to propose challenging tasks, the other aims to solves them. This interaction results in a curriculum that uncovers complex tasks and the diverse skills to solve them. Colas et al. (2019) partition the sensory space into subspaces ("modules") and design a system for agents to self-organise

---

[5]An extended discussion of the challenges MISL approaches face in scaling can be found in Appendix D.2.

curricula over modules, prioritising subspaces of goals they are increasingly or decreasingly successful in reaching. Thereby, agents can learn the skills needed to tackle a diverse range of goals.

**Incremental expansion of capacities**   Shafiullah and Pinto (2022) grow a diverse skill repertoire sequentially (each its own policy), optimising new skills for high state entropy w.r.t. existing skills and low entropy within-skill, promoting diversity in skills while preserving controllability. Pong et al. (2020) learn a maximum-entropy goal distribution via importance-weighted training of a goal generator such that rarely visited states are given more weight. This can result in uniform coverage of states, inducing a diverse range of skills when reaching diverse goals requires diverse behaviours.

**Laplacian-based approaches**   Eigenoptions (Machado et al., 2017) are skills derived from temporal properties of the state space though the eigendecomposition of a matrix representation of its underlying graph. The decomposition yields intrinsic reward functions for learning skills that operate at different time scales. Chen et al. (2023) unify Laplacian-based methods with MISL, extracting the benefits of the two into a powerful framework for learning skills that are diverse and achieve high state coverage.

**Language-guided discovery**   Rho et al. (2025) use language models to guide skill discovery, using *language-distance* as a proxy for semantic distance to learn "semantically diverse" skills.

**Broader context**   The landscape of self-supervised RL for diverse skills is rapidly evolving; Park et al. (2024, Appendix A, pp. 17–18) provide an extensive list of recent approaches.

## 5   Conclusion and future work

We introduced VendiRL, a self-supervised RL framework that centres both skill learning and evaluation on a user-specified similarity function and the Vendi Score. Our proof-of-concept shows that different similarity functions induce distinct behaviours, enabling "pick-and-mix" diversity without having to redesign the mechanism for generating rewards. This decouples what counts as diverse from how rewards are computed, offering a path to scaling diversity-driven pretraining in richer, more open-ended settings where being able to motivate variable forms of diversity may be desirable.

Adapting the Vendi Score to RL further supports evaluation and benchmarking: the same metric can quantify outcomes across algorithms, representations, and domains, assuming a comparable number of skills. This opens up new uses, for example: direct comparisons of existing skill-learning methods under shared notions of diversity, use of diversity as a reward signal for unsupervised environment design (e.g., generating environment variants based on the skill diversity they afford), and construction of skill-diversity benchmarks with variable and composable diversity targets.

Looking ahead, several clear directions follow. First, in assessing skill transfer: how useful are VendiRL-motivated skills in solving tasks across downstream domains, and can properties of those domains be mapped to effective similarity functions or their mixtures? Second, can we find utility in agents that autonomously select and weight different forms of diversity via a meta-controller instead of relying on human priors? Third, studying the scalability and stability of our framework by experimenting with its key components: the choice of skill representation, the design of skill memory, mechanisms to mitigate drift between reward distributions under massive parallelism, and the underlying RL algorithm. Lastly, perhaps most importantly, we aim to deepen our understanding of how different similarity functions shape the behavioural diversity that emerges. This, in turn, can support a widely acknowledged need in RL: the evaluation and fair comparison of the capabilities of open-ended systems.

In short, VendiRL aims to turn what is considered diverse into a variable, optimisable design choice.

## Acknowledgements and disclosure of funding

## References

Arthur Aubret, Laetitia Matignon, and Salima Hassas. An information-theoretic perspective on intrinsic motivation in reinforcement learning: A survey. *Entropy*, 25(2), 2023. URL https://doi.org/10.3390/e25020327.

David Barber and Felix Agakov. The IM algorithm: a variational approach to information maximization. In *Advances in Neural Information Processing Systems*, 2003. URL https://dl.acm.org/doi/10.5555/2981345.2981371.

Kate Baumli, David Warde-Farley, Steven Hansen, and Volodymyr Mnih. Relative variational intrinsic control. In *AAAI Conference on Artificial Intelligence*, 2021. URL https://doi.org/10.1609/aaai.v35i8.16832.

Daniel Berlyne. *Structure and Direction in Thinking*. John Wiley & Sons, 1965.

Matteo Bettini, Ajay Shankar, and Amanda Prorok. System neural diversity: Measuring behavioral heterogeneity in multi-agent learning. *arXiv preprint arXiv:2305.02128*, 2023. URL https://doi.org/10.48550/arXiv.2305.02128.

James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL http://github.com/jax-ml/jax.

Franziska Brändle, Lena J. Stocks, Joshua B. Tenenbaum, Samuel J. Gershman, and Eric Schulz. Empowerment contributes to exploration behaviour in a creative video game. *Nature Human Behaviour*, 7(9):1481–1489, 2023. URL https://doi.org/10.1038/s41562-023-01661-2.

Jiayu Chen, Vaneet Aggarwal, and Tian Lan. A unified algorithm framework for unsupervised discovery of skills based on determinantal point process. In *Advances in Neural Information Processing Systems*, 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/d6938c8e88ef62394d2f4f3fd428e036-Paper-Conference.pdf.

Jongwook Choi, Archit Sharma, Honglak Lee, Sergey Levine, and Shixiang Shane Gu. Variational empowerment as representation learning for goal-conditioned reinforcement learning. In *International Conference on Machine Learning*, 2021. URL https://proceedings.mlr.press/v139/choi21b.html.

Cédric Colas, Pierre Fournier, Mohamed Chetouani, Olivier Sigaud, and Pierre-Yves Oudeyer. CURIOUS: Intrinsically motivated modular multi-goal reinforcement learning. In *International Conference on Machine Learning*, 2019. URL https://proceedings.mlr.press/v97/colas19a.html.

Cédric Colas, Tristan Karch, Olivier Sigaud, and Pierre-Yves Oudeyer. Autotelic agents with intrinsically motivated goal-conditioned reinforcement learning: a short survey. *Journal of Artificial Intelligence Research*, 74:1159–1199, 2022. URL https://doi.org/10.1613/jair.1.13554.

Danica. An overview of MMD, 2017. URL https://stats.stackexchange.com/a/276618.

Benjamin Eysenbach. Self-supervised representations and reinforcement. Tutorial given at the 2025 Multi-disciplinary Conference on Reinforcement Learning and Decision Making (RLDM), 2025. URL https://ben-eysenbach.github.io/self-supervised-rl/.

Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=SJx63jRqFm.

Benjamin Eysenbach, Ruslan Salakhutdinov, and Sergey Levine. The information geometry of unsupervised reinforcement learning. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=3wU2UX0voE.

Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier Bachem. Brax - a differentiable physics engine for large scale rigid body simulation. In *Neural Information Processing Systems*, 2021. URL https://datasets-benchmarks-proceedings.neurips.cc/paper_files/paper/2021/file/d1f491a404d6854880943e5c3cd9ca25-Paper-round1.pdf.

Dan Friedman and Adji Bousso Dieng. The vendi score: A diversity evaluation metric for machine learning. *Transactions on Machine Learning Research*, 2023. URL https://openreview.net/forum?id=g97OHbQyk1.

Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *arXiv preprint arXiv:1611.07507*, 2016. URL https://doi.org/10.48550/arXiv.1611.07507.

Shixiang Shane Gu, Manfred Diaz, Daniel C. Freeman, Hiroki Furuta, Seyed Kamyar Seyed Ghasemipour, Anton Raichuk, Byron David, Erik Frey, Erwin Coumans, and Olivier Bachem. Braxlines: Fast and interactive toolkit for rl-driven behavior engineering beyond reward maximization. *arXiv preprint arXiv:2110.04686*, 2021. URL https://doi.org/10.48550/arXiv.2110.04686.

Steven Hansen, Will Dabney, André Barreto, David Warde-Farley, Tom Van de Wiele, and Volodymyr Mnih. Fast task inference with variational intrinsic successor features. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=BJeAHkrYDS.

Shuncheng He, Yuhang Jiang, Hongchang Zhang, Jianzhun Shao, and Xiangyang Ji. Wasserstein unsupervised reinforcement learning. In *AAAI Conference on Artificial Intelligence*, 2022. URL https://ojs.aaai.org/index.php/AAAI/article/view/20645.

Antti Honkela. Chapter 3: Multivariate normal distributions and numerical linear algebra. In *Computational Statistics I*, 2020. URL https://www.cs.helsinki.fi/u/ahonkela/teaching/compstats1/book/.

Leslie Pack Kaelbling. Learning to achieve goals. In *International Joint Conference on Artificial Intelligence*, 1993. URL https://api.semanticscholar.org/CorpusID:5538688.

A.S. Klyubin, D. Polani, and C.L. Nehaniv. Empowerment: a universal agent-centric measure of control. *IEEE Congress on Evolutionary Computation*, 2005. URL https://doi.org/10.1109/CEC.2005.1554676.

Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. In *Advances in Neural Information Processing Systems*, 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/0234c510bc6d908b28c70ff313743079-Paper.pdf.

Michael Laskin, Denis Yarats, Hao Liu, Kimin Lee, Albert Zhan, Kevin Lu, Catherine Cang, Lerrel Pinto, and Pieter Abbeel. URLB: Unsupervised reinforcement learning benchmark. In *Advances in Neural Information Processing Systems*, 2021. URL https://openreview.net/forum?id=lwrPkQP_is.

Tom Leinster. *Entropy and Diversity: The Axiomatic Approach*. Cambridge University Press, 2021.

Aly Lidayan, Michael D Dennis, and Stuart Russell. BAMDP shaping: a unified theoretical framework for intrinsic motivation and reward shaping. In *International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=tijmpS9Vy2.

Cam Linke, Nadia M. Ady, Martha White, Thomas Degris, and Adam White. Adapting behavior via intrinsic reward: A survey and empirical study. *Journal of Artificial Intelligence Research*, 69: 1287–1332, 2020. URL https://doi.org/10.1613/jair.1.12087.

Erik M. Lintunen, Nadia M. Ady, and Christian Guckelsberger. Diversity progress for goal selection in discriminability-motivated RL. In *Intrinsically Motivated Open-ended Learning Workshop at NeurIPS 2024*, 2024. URL https://openreview.net/forum?id=nz9iquQEJF.

Minghuan Liu, Menghui Zhu, and Weinan Zhang. Goal-conditioned reinforcement learning: Problems and solutions. In *International Joint Conference on Artificial Intelligence*, 2022. URL https://doi.org/10.24963/ijcai.2022/770. Survey Track.

Marlos C. Machado, Marc G. Bellemare, and Michael Bowling. A Laplacian framework for option discovery in reinforcement learning. In *International Conference on Machine Learning*, 2017. URL https://proceedings.mlr.press/v70/machado17a.html.

Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=B1QRgziT-.

Alireza Modirshanechi, Peter Dayan, and Eric Schulz. An integrative framework for the human sense of control. *PsyArXiv preprint 10.31234/osf.io/cnkyz_v1*, 2025. URL https://doi.org/10.31234/osf.io/cnkyz_v1.

Gaia Molinaro, Cédric Colas, Pierre-Yves Oudeyer, and Anne Collins. Latent learning progress drives autonomous goal selection in human reinforcement learning. In *Advances in Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=GbqzN9HiUC.

Andrew Y. Ng, Daishi Harada, and Stuart J. Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *International Conference on Machine Learning*, 1999. URL https://dl.acm.org/doi/10.5555/645528.657613.

OpenAI OpenAI, Matthias Plappert, Raul Sampedro, Tao Xu, Ilge Akkaya, Vineet Kosaraju, Peter Welinder, Ruben D'Sa, Arthur Petron, Henrique P. d. O. Pinto, Alex Paino, Hyeonwoo Noh, Lilian Weng, Qiming Yuan, Casey Chu, and Wojciech Zaremba. Asymmetric self-play for automatic goal discovery in robotic manipulation. *arXiv preprint arXiv:2101.04882*, 2021. URL https://arxiv.org/abs/2101.04882.

Pierre-Yves Oudeyer and Frederic Kaplan. What is intrinsic motivation? A typology of computational approaches. *Frontiers in Neurorobotics*, 1:1–14, 2007. URL https://doi.org/10.3389/neuro.12.006.2007.

Pierre-Yves Oudeyer and Frederic Kaplan. How can we define intrinsic motivation? In *International Conference on Epigenetic Robotics*, 2008. URL https://inria.hal.science/inria-00420175/document.

Seohong Park, Oleh Rybkin, and Sergey Levine. Metra: Scalable unsupervised rl with metric-aware abstraction. In *International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=c5pwL0Soay.

Vitchyr H. Pong, Murtaza Dalal, Steven Lin, Ashvin Nair, Shikhar Bahl, and Sergey Levine. Skew-fit: State-covering self-supervised reinforcement learning. In *International Conference on Machine Learning*, 2020. URL https://proceedings.mlr.press/v119/pong20a.html.

Justin K. Pugh, Lisa B. Soros, and Kenneth O. Stanley. Quality diversity: A new frontier for evolutionary computation. *Frontiers in Robotics and AI*, 3, 2016. URL https://doi.org/10.3389/frobt.2016.00040.

Patrik Reizinger, Bálint Mucsányi, Siyuan Guo, Benjamin Eysenbach, Bernhard Schölkopf, and Wieland Brendel. Skill learning via policy diversity yields identifiable representations for reinforcement learning. *arXiv preprint arXiv:2507.14748*, 2025. URL https://doi.org/10.48550/arXiv.2507.14748.

Seungeun Rho, Laura Smith, Tianyu Li, Sergey Levine, Xue Bin Peng, and Sehoon Ha. Language guided skill discovery. In *International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=i3e92uSZCp.

Cansu Sancaktar, Justus Piater, and Georg Martius. Regularity as intrinsic reward for free play. In *Advances in Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=BHHrX3CRE1.

Nur Muhammad Shafiullah and Lerrel Pinto. One after another: Learning incremental skills for a changing world. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=dg79moSRqIo.

Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware unsupervised discovery of skills. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=HJgLZR4KvH.

David Silver, Satinder Singh, Doina Precup, and Richard S. Sutton. Reward is enough. *Artificial Intelligence*, 299, 2021. URL https://doi.org/10.1016/j.artint.2021.103535.

Matthew J. A. Smith, Jelena Luketina, Kristian Hartikainen, Maximilian Igl, and Shimon Whiteson. Learning skills diverse in value-relevant features. In Sarath Chandar, Razvan Pascanu, and Doina Precup, editors, *Conference on Lifelong Learning Agents*, pages 1174–1194. PMLR, 2022. URL https://proceedings.mlr.press/v199/smith22a.html.

Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018.

David Warde-Farley, Tom Van de Wiele, Tejas Kulkarni, Catalin Ionescu, Steven Hansen, and Volodymyr Mnih. Unsupervised control through non-parametric discriminative rewards. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=r1eVMnA9K7.

Rushuai Yang, Chenjia Bai, Hongyi Guo, Siyuan Li, Bin Zhao, Zhen Wang, Peng Liu, and Xuelong Li. Behavior contrastive learning for unsupervised skill discovery. In *International Conference on Machine Learning*, 2023. URL https://proceedings.mlr.press/v202/yang23a.html.

Yucheng Yang, Tianyi Zhou, Qiang He, Lei Han, Mykola Pechenizkiy, and Meng Fang. Task adaptation from skills: Information geometry, disentanglement, and new objectives for unsupervised reinforcement learning. In *International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=zSxpnKh1yS.

Chongyi Zheng, Jens Tuyls, Joanne Peng, and Benjamin Eysenbach. Can a misl fly? analysis and ingredients for mutual information skill learning. In *International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=xoIeVdFO7U.

Xiang Zheng, Xingjun Ma, Chao Shen, and Cong Wang. CIM: Constrained intrinsic motivation for reinforcement learning, 2024. URL https://openreview.net/forum?id=UnuSBQjgqK.

# A Skill diversity as minimal overlap in the feature space

As one example of measuring skill diversity with the Vendi Score, we employ a skill-similarity measure derived from a $k$-nearest neighbors formulation of *precision* and *recall* developed by Kynkäänniemi et al. (2019). Note that $k$ here denotes a scalar variable, deviating from the use of $k$ in the main body of the paper to represent a kernel function.

To compute the overlap between two stochastic skills, $\boldsymbol{\theta}_a$ and $\boldsymbol{\theta}_b$, we collect two samples of $N$ trajectories, $\mathbf{X}_a \sim \boldsymbol{\theta}_a$ and $\mathbf{X}_b \sim \boldsymbol{\theta}_b$, respectively. In our case, each trajectory is a $\mathbb{R}^{T \times D}$ matrix, where $T$ represents the length of the trajectory in time steps and $D$ the number of observation dimensions in the given environment. For each skill, the samples are concatenated, so we end up with an $\mathbb{R}^{NT \times D}$ matrix consisting of $NT$ observation vectors each denoted $\mathbf{x}$. Then, for each set of observation vectors, $\mathbf{X} \in \{\mathbf{X}_a, \mathbf{X}_b\}$, we compute pairwise Euclidean distances between all observations in the set and, for each observation vector, form a hypersphere with radius equal to the distance to its $k$th nearest neighbour. Together, these hyperspheres form an estimate of the true skill manifold in the observation space (illustrated in Figure 2). Following Kynkäänniemi et al. (2019, p. 3), we use a binary function to determine whether a given observation is located within this manifold:

$$f(\mathbf{x}, \mathbf{X}) := \begin{cases} 1, & \text{if } ||\mathbf{x} - \mathbf{x}'|| \leq ||\mathbf{x}' - \mathrm{NN}_k(\mathbf{x}', \mathbf{X})|| \text{ for at least one } \mathbf{x}' \in \mathbf{X} \\ 0, & \text{otherwise,} \end{cases} \tag{3}$$

where $\mathrm{NN}_k(\mathbf{x}', \mathbf{X})$ represents the $k$th nearest observation from $\mathbf{x}' \in \mathbf{X}$. Then, as in Kynkäänniemi et al. (2019, p. 3), precision and recall are defined respectively as:

$$\mathrm{pr}(\mathbf{X}_a, \mathbf{X}_b) := \frac{1}{|\mathbf{X}_b|} \sum_{\mathbf{x}_b \in \mathbf{X}_b} f(\mathbf{x}_b, \mathbf{X}_a), \quad \mathrm{re}(\mathbf{X}_a, \mathbf{X}_b) := \frac{1}{|\mathbf{X}_a|} \sum_{\mathbf{x}_a \in \mathbf{X}_a} f(\mathbf{x}_a, \mathbf{X}_b). \tag{4}$$

Then, the overlap between $\mathbf{X}_a$ and $\mathbf{X}_b$ is given by the harmonic mean of precision and recall:

$$F_1(\mathbf{X}_a, \mathbf{X}_b) := \frac{2 \times \mathrm{pr}(\mathbf{X}_a, \mathbf{X}_b) \times \mathrm{re}(\mathbf{X}_a, \mathbf{X}_b)}{\mathrm{pr}(\mathbf{X}_a, \mathbf{X}_b) + \mathrm{re}(\mathbf{X}_a, \mathbf{X}_b)}, \tag{5}$$

such that $F_1 = 1 \iff \mathbf{X}_a = \mathbf{X}_b$, and $F_1 \to 0$ indicates little to no overlap.

Precision measures the proportion of observations drawn from skill $\boldsymbol{\theta}_b$ that fall within the estimated support of $\boldsymbol{\theta}_a$, and recall measures the proportion of observations drawn from skill $\boldsymbol{\theta}_a$ that fall within the estimated support of $\boldsymbol{\theta}_b$. The $F_1$ score weights the two measures equally. After computing pairwise scores between skills, their overall diversity—or, the effective number of unique skills—is given by computing the Vendi Score from the kernel matrix (as detailed in Section 3.1).

# B Transforming VendiRL rewards

## B.1 Time derivative

To reward an agent for *increasing* the diversity of its skills between two consecutive time steps, the VendiRL rewards can be defined as $r_t(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}, \mathbf{g}_t) := VS_{k^{t+1}}(\boldsymbol{\theta}_1, \cdots, \boldsymbol{\theta}_n) - VS_{k^t}(\boldsymbol{\theta}_1, \cdots, \boldsymbol{\theta}_n)$.

## B.2 Penalty

The Vendi Score lies in $[1, n]$. Thus, a VendiRL reward can be easily transformed into a *penalty* with $r_t(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}, \mathbf{g}_t) := VS_{k^{t+1}}(\boldsymbol{\theta}_1, \cdots, \boldsymbol{\theta}_n) - n$. We also found $(f \circ g)(r_t)$, where $f(x) := \log(x)$ and $g(x) := x/n$, such that $g : [1, n] \to [1/n, 1]$ and $f : [1/n, 1] \to [\log(1/n), 0]$, to work well.

# C Similarity functions included in this paper

As before, we denote skills as specific configurations of policy parameters: let $\boldsymbol{\theta}_i$ represent a particular skill corresponding to goal $i$, a specific configuration of the skills being learned, $\boldsymbol{\theta}$, from the set of all possible skills $\Theta$. Then, $k(\boldsymbol{\theta}_a, \boldsymbol{\theta}_b)$ measures the similarity between $\boldsymbol{\theta}_a$ and $\boldsymbol{\theta}_b$ such that $k : \Theta \times \Theta \to \mathbb{R}$.

Notably, in case one wishes to enforce $k(\boldsymbol{\theta}_i, \boldsymbol{\theta}_i) = 1$ for all $i \in \{1, \cdots, n\}$, where $n$ denotes the number of skills being learned, outputs from some choices of $k$ will have to be scaled accordingly.

## C.1 Cosine similarity

In addition to the above, let $\boldsymbol{\mu}_i$ represent the mean of a trajectory, over time, induced by skill $\boldsymbol{\theta}_i$. Then, the cosine similarity between skills $\boldsymbol{\theta}_a$ and $\boldsymbol{\theta}_b$ is given by

$$k(\boldsymbol{\theta}_a, \boldsymbol{\theta}_b) := \frac{\boldsymbol{\mu}_a \cdot \boldsymbol{\mu}_b}{\|\boldsymbol{\mu}_a\| \, \|\boldsymbol{\mu}_b\|}. \tag{6}$$

In an unbounded feature space, the output of $k$ is in $[-1, 1]$. Left unscaled, $k(\boldsymbol{\theta}_a, \boldsymbol{\theta}_b) = 1$ when the two means have no angular separation and $k(\boldsymbol{\theta}_a, \boldsymbol{\theta}_b) = -1$ when opposite one another.

## C.2 Maximum mean discrepancy

Again, where $\boldsymbol{\mu}_i$ denotes the mean of a trajectory induced by $\boldsymbol{\theta}_i$, the maximum mean discrepancy (with linear kernel) between skills $\boldsymbol{\theta}_a$ and $\boldsymbol{\theta}_b$ is given by

$$k(\boldsymbol{\theta}_a, \boldsymbol{\theta}_b) := (f \circ g)(\boldsymbol{\mu}_a, \boldsymbol{\mu}_b), \quad \text{where}$$

$$f(x) := \frac{1}{\exp(x)} \quad \text{and} \quad g(\boldsymbol{\mu}_a, \boldsymbol{\mu}_b) := \|\boldsymbol{\mu}_a - \boldsymbol{\mu}_b\|. \tag{7}$$

The function $g$ is derived as follows (Danica, 2017). Maximum mean discrepancy (MMD) is defined based on a feature map $\varphi : \mathbb{X} \to \mathbb{H}$, where $\mathbb{H}$ is some Hilbert space. Then, when $\mathbb{X} = \mathbb{H} = \mathbb{R}^d$ and $\varphi(x) = x$, corresponding to a linear kernel, we have

$$\begin{aligned} \mathrm{MMD}(P, Q) &= \|\mathbb{E}_{X \sim P}\left[\varphi(X)\right] - \mathbb{E}_{Y \sim Q}\left[\varphi(Y)\right]\|_{\mathbb{H}} \\ &= \|\mathbb{E}_{X \sim P}\left[X\right] - \mathbb{E}_{Y \sim Q}\left[Y\right]\|_{\mathbb{R}^d} \\ &= \|\mu_P - \mu_Q\|_{\mathbb{R}^d}. \end{aligned}$$

Computing the distance ($g$) and converting the output into a similarity ($f$) results in $k(\boldsymbol{\theta}_a, \boldsymbol{\theta}_b) = 1$ when the two means are equal and a quantity in $(0, 1)$ otherwise (tending to zero as the distance grows in magnitude).

## C.3 Covariance structure

Given a trajectory of observations induced by skill $\boldsymbol{\theta}_i$, we denote by $\boldsymbol{\Sigma}_i$ the corresponding sample covariance matrix. Then, the similarity in covariance structure is defined as

$$k(\boldsymbol{\theta}_a, \boldsymbol{\theta}_b) := (f \circ g)(\boldsymbol{\Sigma}_a, \boldsymbol{\Sigma}_b), \quad \text{where}$$

$$f(x) := \frac{1}{\exp(x)} \quad \text{and} \quad g(\boldsymbol{\Sigma}_a, \boldsymbol{\Sigma}_b) := |\det(\boldsymbol{\Sigma}_a) - \det(\boldsymbol{\Sigma}_b)|. \tag{8}$$

The determinant $\det(\boldsymbol{\Sigma}_i)$ is evaluated using Cholesky decomposition as follows (Honkela, 2020). The symmetric positive definite matrix $\boldsymbol{\Sigma}_i$ can be represented as $\boldsymbol{\Sigma}_i = \mathbf{L}\mathbf{L}^T$, where $\mathbf{L}$ is a lower-triangular matrix. Then, using basic properties of the determinant and the logarithm, we obtain

$$\begin{aligned} \log \det \boldsymbol{\Sigma}_i &= \log\left(\det(\mathbf{L}\mathbf{L}^T)\right) = \log\left(\det \mathbf{L} \det \mathbf{L}^T\right) \\ &= \log\left((\det \mathbf{L})^2\right) = 2 \log(\det \mathbf{L}) \\ &= 2 \log\left(\prod_{i=1}^{d} l_{ii}\right) = 2 \sum_{i=1}^{d} \log(l_{ii}). \end{aligned}$$

The determinant of the skill's covariance matrix, as the product of its eigenvalues—each of which represents the magnitude of the skill's spread on a principal axis—captures information about the skill's *volume* in feature space. Taking the absolute difference ($g$) and converting the output into a similarity ($f$) results in $k(\boldsymbol{\theta}_a, \boldsymbol{\theta}_b) = 1$ when the two determinants are equal and a quantity in $(0, 1)$ otherwise (tending to zero as the absolute difference grows in magnitude).

# D  Extended discussion of MISL

## D.1  Variational MISL

Formally, variational MISL approximates the lower bound on mutual information using a discriminator $q$ with parameters $\phi$, between the goal-defining variable, $\mathbf{g}$, and some function of the trajectory induced by the corresponding skill. For concreteness, we represent the output of this function by a common choice, the observed successor state, $\mathbf{s}'$, determined by the state-transition distribution, $p_{\boldsymbol{\theta}}(\mathbf{s}' \mid \mathbf{s}, \mathbf{g})$, conditional on the skill-defining parameters $(\boldsymbol{\theta}, \mathbf{g})$. Then, the objective is to maximise

$$\mathcal{F}(\boldsymbol{\theta}, \phi) := \mathbb{E}_{\mathbf{g} \sim p(\mathbf{g}), \mathbf{s}' \sim p_{\boldsymbol{\theta}}(\mathbf{s}'|\mathbf{s},\mathbf{g})} [\underbrace{\log q_{\phi}(\mathbf{g} \mid \mathbf{s}')}_{(\alpha)} - \underbrace{\log p(\mathbf{g})}_{(\beta)}]. \tag{9}$$

($\alpha$) Given the agent's ($\boldsymbol{\theta}$) behaviour, the discriminator ($\phi$) tries to predict which skill the agent is following. The agent is rewarded for learning predictable and thus diverse skills: successfully discriminating skills requires the agent to observe distinct regions of the feature space.

($\beta$) This term is maximised in expectation when skills are selected uniformly at random. If the agent does not learn $p(\mathbf{g})$, it is a common choice to fix it to a uniform distribution, resulting in a constant term. For more detail on fixing versus learning $p(\mathbf{g})$, see Lintunen et al. (2024).

## D.2  Problems in scaling MISL



(a) *Diversity is All You Need* (Eysenbach et al., 2019), a prototypical MISL approach, as part of which an agent maximises the mutual information between skills and the states they induce.

(b) The corresponding decision boundary of the learned classifier used to predict skills from observations. To maximise rewards, the agent has has to be correct in and certain of its predictions.

(c) Spectral normalisation can be used to enforce smoothness in the decision landscape (Choi et al., 2021, pp. 5–6). This makes the task harder to solve, but its use helps avoid overfitting to noise.

(d) *Contrastive Successor Features* (Zheng et al., 2025, visualisation from cited paper with no changes) maximises the mutual information between transitions and skills, fixing dithering issues.
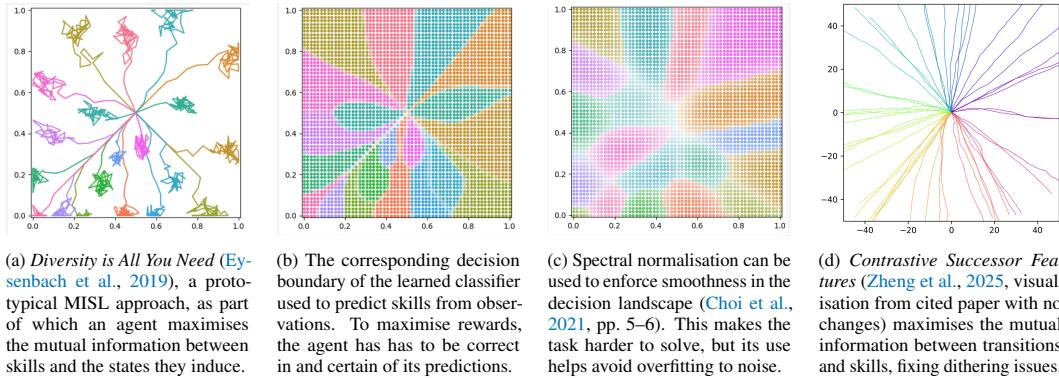
Figure 6: Evolution of mutual information skill learning. Some problems and subsequent solutions.

### D.2.1  Having to determine the number of skills

Learning *the* optimal initialisation for the set of all unknown reward functions in an environment would typically require an agent to learn a large number of skills, making the pretraining computationally demanding if not intractable (to do optimally).[6] Instead, it is common to try to approximate the optimal initialisation by fixing some small number of skills, $n$, and learning $n$ skills from the set of all possible skills. However, not all methods motivate agents to go beyond what is required to discriminate the skills. If $n$ is too small, the skills can start dithering in discriminable regions (e.g., Figure 6a). Further, if the dimensionality of the feature space is high, finding discriminable regions in that space can become trivial—without increasing the number of skills—due to its increase in volume. Discriminability-motivated learning can thus be challenging to implement effectively in complex environments. That said, some MISL methods have been developed recently to address this problem (e.g., Figure 6d), but they still rely on the user's intuition of how many skills an agent should learn.

### D.2.2  The discriminator can overfit to random skills

Since the skills are randomly initialised, and the discriminator is typically a neural network with high expressive power, the neural network can easily overfit to the random skills. This can lead to a

---

[6]Reizinger et al. (2025, p. 8) show that the canonical way of fixing a small number of skills and drawing skills uniformly at random during training is insufficient for learning the ground-truth features of an environment.

non-smooth decision landscape of the discriminator (e.g., Figure 6b), which can result in the skills converging to a suboptimal solution with respect to their diversity. One proposed solution is to use spectral normalisation for regularising the discriminator (Choi et al., 2021, pp. 5–6, referring to the work of Miyato et al., 2018). This enforces smoothness in the discriminator's decision landscape (e.g., Figure 6c); making the task harder to solve, but helping to prevent overfitting to near-random skills. While spectral normalisation clearly offers a solution to the problem of overfitting to random skills, we lack research on the trade-offs between different levels of discriminator expressivity (learning to discriminate skills effectively) and regularisation (preventing overfitting).

### D.2.3 Relying on the assumption of a well-structured feature space

The use of MISL often relies on the user's supervision for structuring the feature space. In practice, this means that when implementing a MISL agent, some dimensions of the feature space are ignored by design. With some specific subset of downstream reward functions in mind, the user determines the most appropriate subset using their prior knowledge of task-relevant features. In such cases, agents maximise diversity in some low-dimensional subset of the space of behaviours (e.g., Eysenbach et al., 2019, p. 7). Without such supervision, agents can end up learning trivially diverse skills that neither translate to observably diverse behaviours in the environment, nor transfer to downstream tasks effectively enough to make the pretraining worthwhile. While automated solutions to this problem have been proposed, such as learning a feature representation that supports transfer learning (e.g., Smith et al., 2022), they still rely on external supervision to the extent of specifying the set—or at the very least, the distribution—of future tasks.