

Dynamics-Adaptive Continual Reinforcement Learning via Progressive Contextualization

Tiantian Zhang, Zichuan Lin, Yuxing Wang, Deheng Ye, Qiang Fu,
Wei Yang, Xueqian Wang, Bin Liang, Bo Yuan, and Xiu Li

Abstract—A key challenge of continual reinforcement learning (CRL) in dynamic environments is to promptly adapt the RL agent’s behavior as the environment changes over its lifetime, while minimizing the catastrophic forgetting of the learned information. To address this challenge, in this article, we propose DaCoRL, i.e., dynamics-adaptive continual RL. DaCoRL learns a context-conditioned policy using progressive contextualization, which incrementally clusters a stream of stationary tasks in the dynamic environment into a series of contexts and opts for an expandable multihead neural network to approximate the policy. Specifically, we define a set of tasks with similar dynamics as an environmental context and formalize context inference as a procedure of online Bayesian infinite Gaussian mixture clustering on environment features, resorting to online Bayesian inference to infer the posterior distribution over contexts. Under the assumption of a Chinese restaurant process prior, this technique can accurately classify the current task as a previously seen context or instantiate a new context as needed without relying on any external indicator to signal environmental changes in advance. Furthermore, we employ an expandable multihead neural network whose output layer is synchronously expanded with the newly instantiated context, and a knowledge distillation regularization term for retaining the performance on learned tasks. As a general framework that can be coupled with various deep RL algorithms, DaCoRL features consistent superiority over existing methods in terms of the stability, overall performance and generalization ability, as verified by extensive experiments on several robot navigation and MuJoCo locomotion tasks.

Index Terms—Dynamic environment, continual reinforcement learning (CRL), incremental context detection, adaptive network expansion.

I. INTRODUCTION

REINFORCEMENT learning (RL) [1] is a major learning paradigm in machine learning for sequential decision making tasks. It aims to train a competent policy for an agent that properly maps states to actions to maximize the

cumulative reward by interacting with an environment in a trial-and-error manner. Traditional RL algorithms, such as Q-learning [2] and SARSA [3], have been widely studied as tabular methods and successfully applied to Markov decision processes (MDPs) with finite discrete state-action spaces. The emergence of advanced function approximation techniques based on deep neural networks (DNNs) enables RL to have a higher level of understanding of the physical world [4], and solve high-dimensional tasks ranging from playing video games directly from pixels [5], [6] to making real-time decisions on continuous robot control tasks [7]–[11].

The progresses of RL have been predominantly focused on learning a single task with the assumption of a stationary¹ and fully-explorable environment for sampling observations. Nevertheless, in the real-world, environments are often non-stationary and characterized by ever-changing dynamics such as shifts in the terrain or weather conditions, changes of the target position in robot navigation [12], different traffic inflow rates and demand patterns at different times of a day (e.g., peak and off-peak hours) in vehicular traffic signal control [13], and the variation in coexisting agents in multiagent systems [14]. These scenarios demand competent RL agents that can continually adapt to new environmental dynamics while retaining performance on all previously encountered environmental conditions.

Unfortunately, the above requirements are difficult for existing RL methods to fulfill. On the one hand, storing all past experiences may result in a constant growth in memory consumption and computational power. On the other hand, if the size of the replay buffer is limited, the agent may inevitably suffer from the phenomenon known as catastrophic forgetting [15]–[17], incapable of retaining the knowledge and skills learned in previously encountered situations.

Recently, Continual RL (CRL) [13], [18]–[22] has been investigated as an effective solution for adaptation to dynamic environments and mitigation of catastrophic forgetting. In this setting, the dynamic environment can be considered as a stream of stationary tasks on a certain timescale where each task corresponds to the specific environmental dynamics during the associated time period. As shown in Fig. 1, the previously learned policy (e.g., $\pi_{\theta_{t-1}^*}$ over $M_1 \sim M_{t-1}$) is used for the initialization of the new policy (e.g., π_{θ_t} on M_t), and it is subsequently updated to fit in the current task during the learning period in a continual fashion, retaining

¹A stationary environment is an environment whose dynamics normally represented by the reward and state transition functions of the MDP do not change over time.

This work was partly supported by the Science and Technology Innovation 2030-Key Project under Grant 2021ZD0201404 and Tencent Rhino-Bird Research Elite Program. (Corresponding author: Bo Yuan and Xiu Li.)

Tiantian Zhang is with the Department of Automation, Tsinghua University, Beijing 100084, China, and also with the Tencent AI Lab, Shenzhen 518000, China (e-mail: ztt19@mails.tsinghua.edu.cn).

Zichuan Lin, Deheng Ye, Qiang Fu, and Wei Yang are with the Tencent AI Lab, Shenzhen 518000, China (e-mail: zichuanlin@tencent.com; deri-ye@tencent.com; leonfu@tencent.com; willyang@tencent.com).

Yuxing Wang, Xueqian Wang and Xiu Li are with the Shenzhen International Graduate School, Tsinghua University, Shenzhen 518055, China (e-mail: wyx20@mails.tsinghua.edu.cn; wang.xq@sz.tsinghua.edu.cn; li.xiu@sz.tsinghua.edu.cn).

Bin Liang is with the Department of Automation, Tsinghua University, Beijing 100084, China (e-mail: liangbin@mail.tsinghua.edu.cn).

Bo Yuan is with Qianyuan Institute of Sciences, Hangzhou 310000, China (e-mail: boyuan@ieee.org).

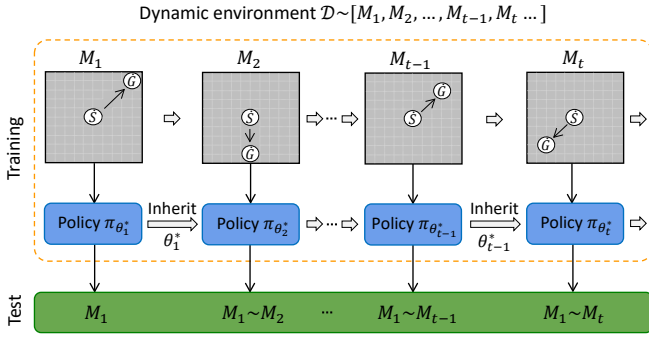


Fig. 1. Continual reinforcement learning (CRL) in dynamic environments. $M_t \in \mathcal{M}, t = 1, 2, \dots$ denotes the specific MDP/task in time period t ; \mathcal{D} denotes the dynamic environment over the MDPs space \mathcal{M} ; θ are the learning parameters of policy; $\pi_{\theta_t^*}$ represents the approximate optimal policy over all learned tasks $M_1 \sim M_t$.

previously learned abilities. In other words, CRL is capable of developing proper behaviors for new tasks while keeping the overall performance across all learned tasks. Such features of continual learning are highly desirable for intelligent systems in real-world applications where the environments are subject to consistent changes.

Existing CRL methods [18]–[21] assume a decomposition of the original problem into disjoint sub-domains of similar dynamics (also called “*context*”) and their boundaries are known in advance. Consequently, previous studies mainly focus on how to construct effective mechanisms to mitigate the catastrophic forgetting among contexts, largely ignoring the challenge of automatic context inference during the learning process. For the alleviation of catastrophic forgetting/interference caused by data distribution drift in the single-task RL, Zhang *et al.* [22] employ sequential K-means clustering to achieve automatic context inference, given the number of contexts (k) in advance. This method works well because it is feasible to acquire an approximate estimate of the state distribution with sufficient exploration to determine the value of k in a single task. However, significant challenges are expected in dynamic environments, where it is impractical to accurately determine the number of environmental contexts in advance as the changes of environmental dynamics are usually infinite and highly uncertain. As a result, it is more rational for the agent to infer and instantiate environmental contexts in a fully online and incremental manner during the CRL process.

In this article, we investigate CRL in dynamic environments to achieve continual inference of environmental contexts and necessary adaptation. The ultimate objective is to ensure the overall performance of the agent in the whole environment. This work is a significant and essential extension to our previous framework for single-task RL in [22], which relies on the prior knowledge about the number of contexts.

To this end, we propose a novel dynamics-adaptive continual reinforcement learning scheme (DaCoRL) with progressive contextualization, which incrementally clusters a stream of stationary tasks in dynamic environment into a series of contexts and opts for an expandable multihead neural network to learn a context-conditioned policy. The progressive contextualization

contains two core modules: The first one is the incremental context detection procedure for automatically detecting the changes in environmental dynamics and clustering a set of tasks with similar dynamics into an environmental context. The second one is the joint optimization procedure to train the policy online for each unique context using an expandable multihead neural network and a knowledge distillation regularization term.

To detect the changes of environmental dynamics over time, we introduce the online Bayesian infinite Gaussian mixture model to cluster environment features in a latent context space, where each cluster corresponds to a separate context. We employ the online Bayesian inference to update the model of contexts in a fully incremental manner, assuming that the prior distribution over the contexts is a Chinese restaurant process (CRP) [23]. With this online incremental clustering technique, DaCoRL can incrementally instantiate new contexts according to the concentration parameter² of CRP as needed, without requiring any external information of environmental changes such as the predetermined number of contexts in [22]. During the joint optimization procedure, we introduce an expandable multihead neural network whose output heads can be adaptively expanded according to the number of instantiated contexts. Compared with the fixed structure neural networks, this approach can eliminate unnecessary redundancy of network structure and improve learning efficiency.

The contributions of this article are summarized as follows.

- 1) An incremental context detection strategy is introduced for CRL in dynamic environments. It formalizes context inference in dynamic environments as an online Bayesian infinite Gaussian mixture clustering procedure on environmental features, enabling the agent to properly identify changes in environmental dynamics in an online manner without any prior knowledge of contexts.
- 2) A novel dynamics-adaptive CRL training scheme called DaCoRL is proposed for dynamic environments in continuous spaces. By employing an expandable multihead neural network in which an output head is added synchronously with the newly instantiated context, and a knowledge distillation regularization term, DaCoRL can effectively alleviate the catastrophic forgetting and ensure competitive capacity of RL agents for continual learning in dynamic environments.
- 3) Extensive experiments on a suite of continuous control tasks ranging from robot navigation to MuJoCo locomotion are conducted to validate the overall superiority of our method over baselines in terms of the stability, overall performance and generalization ability.

In the rest part of this article, Section II reviews the related work and Section III introduces the problem statement of CRL and relevant concepts and notations. In Section IV, the framework of DaCoRL is presented, with details on its mechanism and implementation. Experimental results and analyses on several robot navigation and MuJoCo locomotion tasks are presented in Section V to provide comprehensive evidences

²It is a hyperparameter in CRP denoted by α in this article to control the likelihood of new contexts.

on the superiority of DaCoRL over existing techniques. This article is concluded in Section VI with some discussions and directions for future work.

II. RELATED WORK

Continual learning is conceptually related to incremental learning [24]–[29] and online learning [30], [31] as they all assume that tasks or samples are presented in a sequential manner. On the one side, although incremental learning and continual learning are frequently used interchangeably in the literature, they are not always the same. In some studies [24]–[26], incremental learning is used to describe a learning process where a sequence of incremental tasks are learned in a continual manner — in this case, continual learning can be referred to as incremental learning. Nevertheless, several other studies on incremental learning [27]–[29] concentrate on how to incrementally adjust the previously learned policy to facilitate the fast adaptation to a new task, while ignoring how the agent performs on old tasks. In this sense, continual learning is different from incremental learning.

On the other side, online learning aims to fit a single model to *a single task* over a sequence of data instances without any adaptation to new tasks or concerns for the mitigation of catastrophic forgetting [30], [31]. By contrast, continual learning considers how to learn *a sequence of tasks* while maintaining the performance on previously learned tasks. Another field related to continual learning is multitask learning [32]–[34], which tries to train a single model on multiple tasks simultaneously. The key difference is that, in multitask learning, the training data for all tasks are simultaneously available while continual learning is usually based on the assumption that the training data of a previous task cannot be readily used for training on the current task.

For CRL in dynamic environments, the greatest challenge comes from detecting the changes of environment autonomously. RL-CD [35] is a model-based approach that estimates a set of partial models (containing the transition probability and reward function for each underlying MDP) for predicting environmental dynamics. It detects context changes by continuously evaluating the prediction quality of each partial model on the given experience transitions in the current environment. The model with the highest prediction quality is designated as the current active model, and when there is no model with quality higher than a predefined threshold, a new one is created. Although no prior knowledge is required, RL-CD is computationally and memory intensive since it needs to build an MDP for each context.

Recently, some model-free methods [21], [22], [36] perform context division directly based on experienced transitions to achieve continual reinforcement learning. Context QL [21], a tabular CRL method, applies an online parametric Dirichlet change point (ODCP) algorithm to state-reward sequences collected during learning to detect changes in the dynamic environments. According to the results of detection, Context QL may learn a new Q table for the newly detected context, or improve the policy learned if the current environmental context has been previously experienced. However, there is

a key assumption that the pattern of context changes is known and the number of such changes is finite, since ODCP can only determine whether the environmental context has changed, instead of the specific context. Meanwhile, it is only applicable to RL problems with a small and discrete state-action space due to the Q-tables in use.

For dynamic environments in continuous state spaces, CRL-Unsup [36] is an end to end model-free strategy that detects distributional shifts by tracking the ability of the agent to perform the task and then consolidates the memory when the change is detected. In practice, when the difference between the short-term and long-term moving averages of rewards goes below a certain threshold, the memory elastic weight consolidation (EWC [18]) procedure is triggered to prevent catastrophic forgetting. This method requires storing all policy weights learned before each EWC process is triggered and can be problematic in the case of a positive forward transfer followed by a possible negative backward transfer, as the non-decreasing training cumulative reward curve is used to detect distributional changes. By contrast, our method only needs to store an extra policy model learned up to the beginning of the current task as the teacher network for knowledge distillation, and the incremental context detection module is more effective and timely in the detection of environmental changes.

Recently, IQ [22] shows that performing context division by online clustering and training a multihead neural network with a knowledge distillation regularization term can alleviate the catastrophic interference caused by data distribution drift in the single-task RL. However, it requires the prior knowledge of the number of contexts, which limits its applicability in dynamic environments. In this article, we lift this restriction by introducing an incremental context detection module, which can instantiate contexts incrementally as needed without requiring any prior knowledge of contexts.

In addition to the aforementioned CRL efforts, LLIRL [29] is a recently proposed lifelong adaptation approach for dynamic environments, which employs the EM algorithm, together with a CRP prior on the context distribution, to learn an infinite mixture model to cluster the tasks in dynamic environments incrementally over time. In this way, LLIRL can build upon previous experiences and selectively retrieve the necessary experience to facilitate the adaptation to the current task. It should be noted that LLIRL needs to optimize two separate sets of network parameters to train the behavior policy and parameterize the environment for each instantiated context, respectively, which greatly increases the training and storing cost. Furthermore, it does not consider the issue of catastrophic interference among in-context tasks. By contrast, our proposed method avoids the burden of multiple neural network training by performing policy learning using an expandable multihead neural network. Furthermore, the knowledge distillation technique can effectively reduce the interference among both between-context and in-context tasks, making it possible to conduct effective learning in dynamic environments with a single policy network.

In this article, we aim to design an efficient online CRL method that can automatically detect and identify environmental changes without any prior knowledge of environmental

contexts and achieve continual and stable learning in dynamic environments with continuous state-action spaces, avoiding the adverse effect of catastrophic forgetting.

III. PRELIMINARIES AND NOTATIONS

The formulation of the CRL problem in the domain of dynamic environments and the related key concepts are introduced in this section. The notations used in this article are summarized in Table I.

A. Problem Formulation

1) *Reinforcement Learning in Continuous Spaces*: RL is commonly studied following the MDP framework [1], which is defined as a tuple $M = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$, where \mathcal{S} is the set of states; \mathcal{A} is the set of actions; $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the environment transition probability function; $\mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the reward function, and $\gamma \in [0, 1]$ is the discount factor. At each time step $t \in \mathbb{N}$, the agent moves from s_t to s_{t+1} with probability $p(s_{t+1}|s_t, a_t)$ after it takes action a_t , and receives instant reward r_t .

Most RL algorithms rely on the mechanism of policy gradient to handle tasks with continuous state and action spaces [37]. In such cases, the policy is defined as a function $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, mapping each state to a probability distribution of actions, with $\sum_{a \in \mathcal{A}} \pi(a|s) = 1, \forall s \in \mathcal{S}$. If only the state space is continuous, the policy can use the Boltzmann distribution to select discrete actions, while when the action space is also continuous, the Gaussian distribution is commonly used. The goal of policy-based RL is to find an optimal policy π^* with internal parameter $\theta \in \Theta$ that maximizes the expected long-term discount return

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [R(\tau)] = \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] \quad (1)$$

where, the expectation is over the complete trajectory $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots)$ generated following π_{θ} until the end of the agent's lifetime.

In basic policy gradient methods (e.g., REINFORCE [38]), the action-selection distribution is usually parameterized by a deep neural network trained by taking the gradient ascent with the partial derivative of the objective (i.e., maximize the expected return) with respect to the policy parameters. The policy gradient can be approximated expressed as

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [R(\tau)] \\ &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [R(\tau) \nabla \log \pi_{\theta}(\tau)] \\ &\approx \frac{1}{N} \sum_{n=1}^N R(\tau^n) \nabla_{\theta} \log \pi_{\theta}(\tau^n) \end{aligned} \quad (2)$$

where $(\tau^1, \tau^2, \dots, \tau^N)$ is a batch of learning trajectories sampled from policy π_{θ} to estimate the return expectation.

2) *CRL in Dynamic Environments*: Following the convention in [29], in this article, we consider the dynamic environment as an infinite sequence of stationary tasks where each task corresponds to the specific environmental dynamics within its

TABLE I
NOTATIONS AND THEIR DESCRIPTIONS

Notation	Description
\mathcal{M}	space of MDPs
\mathcal{D}	dynamic environment over time in \mathcal{M}
M_t	stationary task in the t^{th} time period
x_t	feature vector of M_t
θ_t	weights of policy network in time period t
π_{θ^*}	approximate optimal policy over $M_1 \sim M_t$
π_r	uniform random policy
z_t, z_t^*	latent and assigned context label of task M_t
$z_{1:t}^*$	assigned context labels (z_1, z_2, \dots, z_t) for $[M_1, M_2, \dots, M_t]$
α	concentration parameter in CRP
$m_k^{(t)}$	number of assignments to context k up to the t^{th} time period
K_t	number of instantiated contexts up to the t^{th} time period
$\varphi_k^{(t)}$	estimation of parameters of context k in time period t
$\mu_k^{(t)}$	estimation of centroid vector of context k in time period t
θ_S	weights of the shared representation layers in policy network
$\theta_{\mathcal{H},k}$	weights of the k^{th} output head in policy network
β	learning rate for policy update
λ	coefficient of distillation regularization

time period and the same dynamics may recur more than once across different time periods. The time period is assumed to be long enough for the agent to get sufficient experience samples to finish policy learning for the associated task. Suppose that there is a space of MDPs denoted as \mathcal{M} , and a dynamic environment \mathcal{D} changing over time in \mathcal{M} . The CRL agent interacts with $\mathcal{D} = [M_1, M_2, \dots, M_{t-1}, M_t, \dots]$, where each $M_t \in \mathcal{M}$ is a specific MDP/task that is stationary in the t^{th} time period, and the identify of each task $M_t, t \in [1, 2, \dots]$ is *unknown* to the agent.

To represent the policies for multiple tasks with a single model, we assume that the observations of the dynamic environment contain the discriminative information of the tasks. Under this assumption, the goal of CRL in dynamic environments in time period t is to extend the acquired knowledge, accumulated from previously learned tasks $M_1 \sim M_{t-1}$, to the current task M_t , to learn a single policy to achieve the maximum return on all learned tasks $M_1 \sim M_t$

$$\theta_t^* = \arg \max_{\theta} \sum_{i=1}^t J_{M_i}(\theta) \quad (3)$$

where J_{M_i} is the expected return on task M_i . Note that the agent is expected to learn a sequence of tasks one by one strategically so that it can retain the previously acquired knowledge when learning new tasks. In other words, the policy π with parameter θ_t^* in (3) is an approximate optimal policy over all learned tasks $M_1 \sim M_t$. In this article, the overall performance of the learned policy across all tasks is the primary metric for judging the performance of CRL agents.

B. Chinese Restaurant Process

The CRP [23] is a discrete-time stochastic process that defines a prior distribution over the cluster structures. The term CRP arises from an analogy of seating a sequence of customers (equivalent to a stream of observations) in a Chinese restaurant with an infinite number of tables (equivalent to clusters) and

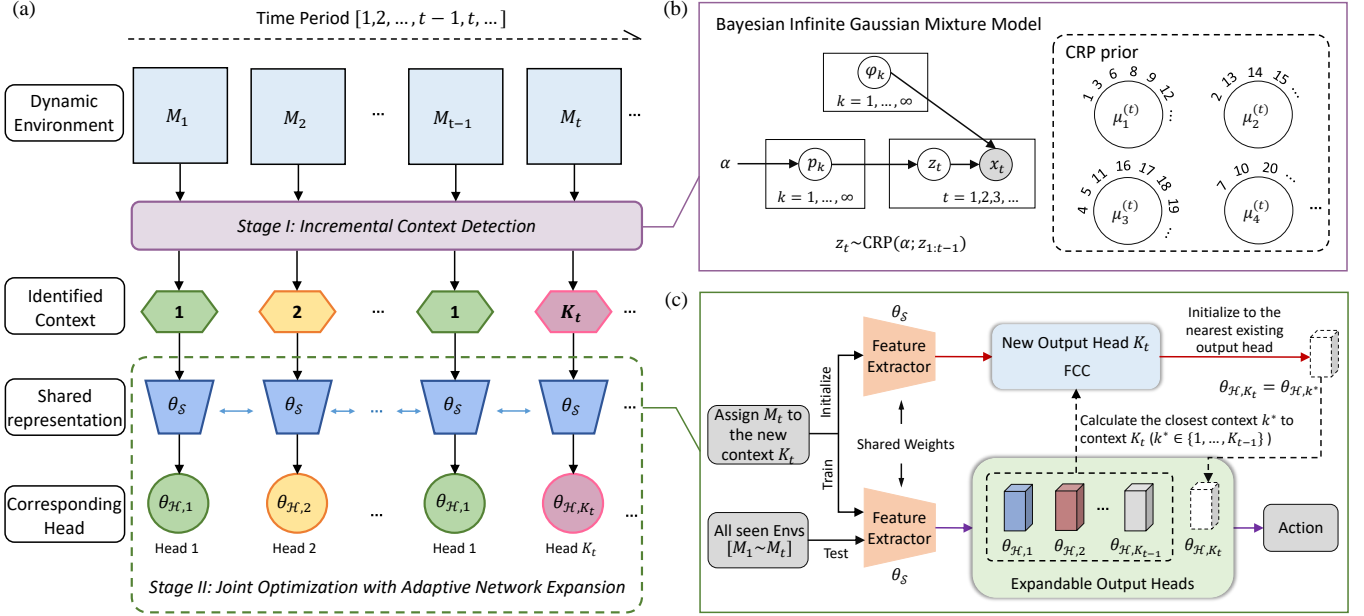


Fig. 2. An overview of DaCoRL in dynamic environments. (a) The general framework of DaCoRL. The dynamic environment is represented by a sequence of stationary tasks $[M_1, M_2, \dots, M_{t-1}, M_t, \dots]$ with which the CRL agent interacts sequentially. The incremental context detection module either associates an existing context with the current task (e.g., M_{t-1} belongs to context 1) or instantiates a new context as needed (e.g., M_t belongs to a new context K_t) online. When training sequentially on different tasks, a policy network with shared feature extractor (blue) and a set of expandable output heads corresponding to different contexts is maintained. (b) The Bayesian infinite Gaussian mixture model with the CRP prior for context inference. (c) Multihead neural network expansion. For task M_t assigned to the new context K_t , we add a new output head and initialize it to the nearest existing output head, and then train the whole network on M_t while keeping the performance on learned tasks unchanged.

each table has infinity capacity. Each customer sits randomly at an occupied table with probability proportional to the number of current customers at that table, or at an unoccupied table with probability proportional to a hyperparameter α ($\alpha > 0$). The conditional probability for the t^{th} customer sitting at the k^{th} table is

$$p(z_t = k | z_{1:t-1}^*, \alpha) = \begin{cases} \frac{m_k^{(t-1)}}{t-1+\alpha}, & k \leq K_{t-1} \\ \frac{\alpha}{t-1+\alpha}, & k = K_{t-1} + 1 \end{cases} \quad (4)$$

where z_t is the latent cluster label of the t^{th} customer; k is the cluster label and $z_{1:t-1}^* = \{z_1^*, z_2^*, \dots, z_{t-1}^*\}$ represents the cluster labels assigned to the $t-1$ customers, respectively; $m_k^{(t-1)}$ counts the number of assignments to cluster k up to time $t-1$ and K_{t-1} denotes the number of occupied clusters after the t^{th} customer is seated; α serves as the concentration parameter that controls the likelihood of new clusters.

IV. PROPOSED METHOD

In this section, we give a detailed description of DaCoRL whose overview is shown in Fig. 2. We first introduce the key components including incremental context detection and adaptive network expansion. Then, we present a joint optimization scheme combining the proposed techniques with the policy-based RL to achieve competent CRL in dynamic environments. As DaCoRL can be incorporated into any canonical policy-based RL methods, for the sake of clarity, we present an instantiation of DaCoRL using the vanilla policy gradient RL algorithm REINFORCE [38].

A. Incremental Context Detection

In dynamic environments, it is important to automatically detect and identify the changes of the environment, which can enable specialized policy learning for different tasks, alleviating the catastrophic forgetting. In this article, to characterize the environmental dynamics, we inherit and extend the concept of “context” in [22] where a set of tasks with similar dynamics is regarded as the same context. Under this setting, the detection of environmental changes can be transformed into a context division procedure in the latent space.

However, compared with the single stationary environment investigated in [22], context division in dynamic environments brings significantly more challenges. In particular, the changes are usually infinite and highly uncertain and, without prior knowledge, it is hard to determine in advance the number of contexts that need to be instantiated.

To address this issue, we propose an incremental context detection module in DaCoRL, which can instantiate contexts when necessary in an incremental manner without any prior knowledge. Specifically, we formalize the task of context detection as an online Bayesian infinite Gaussian mixture clustering procedure. To do so, we define a set of environment features to represent the environmental dynamics, and then perform clustering on the features by assuming a CRP prior distribution over the clustering structure. In CRP, the instantiation of new contexts is controlled by the concentration parameter α : the larger the value of α , the more contexts instantiated in general. On the one hand, when $\alpha = 0$, only a single context is instantiated during the learning process. On the other hand, when $\alpha \rightarrow \infty$, a new context is instantiated for

each task encountered. The complete procedure of incremental context detection with CRP is described in Algorithm 1, which mostly entails the next two steps.

1) *Environment Feature Construction (lines 1-2)*: With the assumption that the task information is contained in the observation space, we can implicitly detect the changes of environmental dynamics by tracking the variation of the observation distribution. In practical implementations, we construct the features for a specific task from the collected observations. Specifically, before the start of RL training in each time period, the agent is required to explore the current environment using a random policy

$$\pi_r(a|s) = \text{Uniform}(A(s)) \quad (5)$$

where $A(s)$ is the set of available actions in the state s and $\text{Uniform}(\cdot)$ is the uniform distribution function. Subsequently, we construct the feature vector (x_t) for M_t so that x_t can represent the dynamics of the current environment from the observations $\mathcal{T}_\varepsilon = \{s_0^i, s_1^i, s_2^i, \dots\}_{i=1}^m$, where m is the number of trajectories and i denotes the trajectory index.

Here, we approximate the set of collected observations as a Gaussian distribution and use its mean vector as the feature of the current environment. Naturally, we can obtain x_t directly from the original observation space

$$x_t = \frac{1}{N} \sum_{i=1}^N s_i \quad (6)$$

for the vector inputs, or from the embedding space

$$x_t = \frac{1}{N} \sum_{i=1}^N \phi(s_i) \quad (7)$$

for the visual inputs, where N is the number of states contained in \mathcal{T}_ε and $\phi(s)$ can be represented by a random encoder [39] or a standard GAN model [40].

2) *Context Detection via Online Incremental Clustering (lines 3-12)*: To enable incremental context detection in dynamic environments, based on the constructed environment features, we employ online Bayesian inference to update the context models in a fully online fashion, avoiding the necessity for storing previously seen samples. The key step is to estimate the context parameters $\{\varphi_k\}_{k=1}^{K_t}$, which can build up a mapping between specific tasks and context variables in the latent space. In our implementations, the context model (predictive likelihood function) represents environment features using diagonal Gaussian distributions

$$p(x_t|\varphi_k) = \mathcal{N}(x_t; \mu_k, \sigma^2) \quad (8)$$

where μ_k is the mean of the Gaussian used to denote context k and σ^2 is a constant indicating the variance. Under this setting, the context centroids are just the mean vectors of the Gaussian distributions to be estimated: $\varphi_k = \{\mu_k\}$.

For a sequence of tasks $[M_1, M_2, \dots, M_{t-1}, M_t]$, the first task is assigned to the first context by default. For M_t in the t^{th} time period, we instantiate a new potential context $K_{t-1} + 1$, and initialize the new context model parameter $\varphi_{K_{t-1}+1}^{(t-1)}$ as the feature vector x_t , regardless of whether M_t has been encountered before (line 3). After that, the posterior probabilities

Algorithm 1 Incremental Context Detection with CRP

Input: The learning task M_t in the t^{th} time period;

The parameter set $\{\varphi_k^{(t-1)}\}_{k=1}^{K_{t-1}}$ of contexts already instantiated before the t^{th} time period.

Parameter: Concentration parameter α .

Output: Task-to-context assignment z_t^* for M_t .

The parameters $\{\varphi_k^{(t)}\}_{k=1}^{K_t}$ of instantiated contexts.

- 1: Sample m trajectories from M_t using a uniform policy π_r : $\mathcal{T}_\varepsilon = \{\tau_i\}_{i=1}^m, \tau_i \sim \pi_r$.
 - 2: Construct the feature vector x_t of M_t from \mathcal{T}_ε .
 - 3: Initialize $\varphi_{K_{t-1}+1}^{(t-1)}$ of a new potential context as x_t .
 - 4: Compute CRP prior $p(z_t = k|z_{1:t-1}^*, \alpha)$ for x_t .
 - 5: Compute posterior probabilities of task-to-context assignment $p(z_t = k|z_{1:t-1}^*, x_t)$.
 - 6: **if** $p(z_t = K_{t-1} + 1|z_{1:t-1}^*, x_t) > p(z_t = k|z_{1:t-1}^*, x_t), \forall k \leq K_{t-1}$ **then**
 - 7: $K_t = K_{t-1} + 1$, add $\varphi_{K_{t-1}+1}^{(t-1)}$ to $\{\varphi_k^{(t-1)}\}_{k=1}^{K_{t-1}}$.
 - 8: **else**
 - 9: $K_t = K_{t-1}$.
 - 10: **end if**
 - 11: Update parameters $\{\varphi_k^{(t-1)}\}_{k=1}^{K_t}$.
 - 12: Calculate $z_t^* = \arg \max_k p(x_t|z_t = k, \varphi_k^{(t)})$, $\forall k \leq K_t$ to obtain the final assignment of M_t .
 - 13: **return** $z_t^*, \{\varphi_k^{(t)}\}_{k=1}^{K_t}$.
-

of the task-to-context assignment over all $K_{t-1} + 1$ contexts are estimated (lines 4-5) by

$$p(z_t = k|z_{1:t-1}^*, x_t, \varphi_k^{(t-1)}, \alpha) \propto p(x_t|\varphi_k^{(t-1)})p(z_t = k|z_{1:t-1}^*, \alpha). \quad (9)$$

By combining the definition of the predictive likelihood in (8) and the CRP prior distribution in (4), the posterior distribution can be rewritten as

$$p(z_t = k|z_{1:t-1}^*, x_t, \varphi_k^{(t-1)}, \alpha) \propto \begin{cases} m_{k,t-1} \mathcal{N}(x_t; \mu_k, \sigma^2), & k \leq K_{t-1} \\ \alpha \mathcal{N}(x_t; \mu_k, \sigma^2), & k = K_{t-1} + 1 \end{cases} \quad (10)$$

which is abbreviated to $p(z_t = k|z_{1:t-1}^*, x_t)$ in subsequent derivations for simplicity.

With the estimated $p(z_t = k|z_{1:t-1}^*, x_t)$ for each context, we can determine whether to keep the new context for M_t (lines 6-10). If the posterior probability of the potential context is greater than those of the K_{t-1} existing contexts, this new context is instantiated for M_t .

Next, based on the inferred posterior probabilities, we update context parameters by optimizing the expected log-likelihood

$$\mathcal{L}(x_t|\varphi_k^{(t-1)}) = \mathbb{E}_{M_t} \log p(x_t|z_t = k, \varphi_k^{(t-1)}) \quad (11)$$

where $M_t \sim p(z_t = k|z_{1:t-1}^*, x_t)$. Suppose that all context models with the prior parameters $\varphi_k^{(t-1)}$ have been optimized up to the $(t-1)^{\text{th}}$ time period. The estimation of $\{\varphi_k^{(t)}\}_{k=1}^{K_t}$ can be updated based on the gradient (line 11)

$$\varphi_k^{(t)} \leftarrow \varphi_k^{(t-1)} + \eta_{k,t} \nabla_{\varphi_k^{(t-1)}} \mathcal{L}(x_t|\varphi_k^{(t-1)}), \forall k \leq K_t \quad (12)$$

where $\eta_{k,t} = \frac{1}{m_{k,t-1} + p(z_t=k|z_{1:t-1}^*, x_t)}$ is the learning rate for the k^{th} context in time period t . The gradient term can be derived as $p(z_t=k|z_{1:t-1}^*, x_t) \nabla_{\varphi_k} \log p(x_t|z_t=k, \varphi_k^{(t-1)})$.

Based on the updated context parameters, the identity z_t^* of M_t can be finally obtained by computing an MAP estimate on the predictive likelihood (line 12), selecting the context model that best fits the current environment.

B. Adaptive Network Expansion

To minimize the interference among contexts in CRL, we opt for an expandable neural network as the policy network in DaCoRL, in which an output head is added synchronously with the newly instantiated context. Each output head specializes on a specific context, and the representation layers are shared among different contexts. This adaptive and expandable network structure can parameterize a dedicated policy for each context without any unnecessary redundancy. In Fig. 2(a), the set of weights of the expandable multihead neural network is denoted by $\theta = \{\theta_S, \theta_{\mathcal{H},1}, \theta_{\mathcal{H},2}, \dots, \theta_{\mathcal{H},K_t}, \dots\}$, where θ_S is a set of weights for shared representation layers, and $\theta_{\mathcal{H},k}$, $k \in \{1, 2, \dots\}$ are the parameters of the k^{th} output head.

The policy network is initialized as a canonical single-head neural network for the forthcoming learning in the first context. When a new context is instantiated, the neural network is adaptively expanded by adding an output head whose structure is consistent with that of the existing output heads. For the newly added output head K_t , we propose the following three practical implementations for parameter initialization.

- 1) *Random Initialization*: The weights of the newly added output head are initialized to random values. In this case, the policy of the newly instantiated context inherits the representation module of the learned policies, while its output layer is trained from scratch.
- 2) *Random Trained Head Initialization*: It randomly selects one of the trained output heads and then initializes the weights of the newly added output head to those of the selected one

$$\theta_{\mathcal{H},K_t} = \theta_{\mathcal{H},k}, \quad k = \text{Uniform}(\{1, 2, \dots, K_{t-1}\}). \quad (13)$$

This method enables the new policy to inherit knowledge from the learned context. However, it is likely that the cloned policy may hinder the CRL agent's ability to properly explore the task space corresponding to the current context, especially when there are significant differences between the two contexts.

- 3) *Nearest Trained Head Initialization*: It initializes the newly added output head to a specific trained one whose associated context is nearest to the current instantiated context in the latent context space.

$$\theta_{\mathcal{H},K_t} = \theta_{\mathcal{H},k^*}, \quad k^* = \arg \min_k \text{dist}(\varphi_k, \varphi_{K_t}) \quad (14)$$

where $k \in \{1, 2, \dots, K_{t-1}\}$ and $\text{dist}(\varphi_k, \varphi_{K_t})$ denotes the distance between contexts k and K_t .

Intuitively, the third implementation is most likely to encourage forward transfer. The experimental results and analysis in Section V-F provide further elaboration on this point.

C. Joint Optimization Scheme

In the policy optimization stage, we integrate the above components with policy-based RL algorithms to achieve efficient CRL in dynamic environments. Furthermore, we employ the knowledge distillation technique [22] to mitigate the interference in the learning in different contexts (corresponding to different output heads in the policy network) caused by the shared low-level representation. This technique can effectively lessen the catastrophic interference caused by distribution drifts among contexts as well as the minor interference due to the differences among tasks within the same context. Taking REINFORCE as the underlying policy-based RL algorithm as an example, the optimization objective of our proposed DaCoRL is derived as follows.

First of all, we rewrite the original loss function of REINFORCE in (1) with the context label variable z_t^* as

$$\mathcal{L}_{\text{ori}}(\theta_{z_t^*}) = \mathbb{E}_{\tau \sim \pi_{\theta_{z_t^*}}(\tau)} [R(\tau)] \quad (15)$$

where $\theta_{z_t^*} = \{\theta_S, \theta_{\mathcal{H},z_t^*}\}$ and $\pi_{\theta_{z_t^*}}$ is the policy corresponding to the context associated with the current task M_t .

As a representative of model compression, knowledge distillation can work well for encouraging the outputs of one network to approximate the outputs of another [41], [42]. In DaCoRL, we use it as a regularization term in the probability distribution estimation of actions to preserve the previously learned policies. To construct the distillation regularization term, we regard the policy network from in the last time period as the teacher network, expressed as π_{θ^-} , and the current policy network to be trained as the student network, expressed as π_{θ} . We use the Kullback-Leibler (KL) divergence to constrain the difference between the policies corresponding to the two networks. Thus, the distillation loss of the output head for context k is defined as

$$\mathcal{L}_{\mathcal{D}_k}(\theta_k) = \mathbb{E}_{\tau \sim \pi_{\theta_{z_t^*}}(\tau)} \left[\text{KL}[\pi_{\theta_k}(\cdot|s), \pi_{\theta_k^-}(\cdot|s)] \right] \quad (16)$$

where $\theta_k = \{\theta_S, \theta_{\mathcal{H},k}\}$. In the t^{th} time period, considering all K_t output heads in the policy network, the distillation loss term sums up as

$$\mathcal{L}_{\mathcal{D}}(\theta) = \sum_{k=1}^{K_t} \mathcal{L}_{\mathcal{D}_k}(\theta_k) \quad (17)$$

where $\theta = \{\theta_S, \{\theta_{\mathcal{H},k}\}_{k=1}^{K_t}\}$ denotes the set of all weights of the policy network.

Finally, to optimize a policy network that can guide the agent to make proper decisions in dynamic environments without being adversely affected by catastrophic forgetting, we combine (15) and (17) to form a joint optimization scheme. Namely, we solve the CRL problem in dynamic environments by the following optimization objective

$$\begin{aligned} & \max_{\theta_S, \theta_{\mathcal{H}}} \mathcal{L}_{\text{ori}}(\theta_{z_t^*}) - \lambda \mathcal{L}_{\mathcal{D}}(\theta) \\ & \theta_{z_t^*} = \{\theta_S, \theta_{\mathcal{H},z_t^*}\} \\ & \theta = \{\theta_S, \{\theta_{\mathcal{H},k}\}_{k=1}^{K_t}\} \end{aligned} \quad (18)$$

where $\lambda \in [0, 1]$ is a coefficient to control the tradeoff between learning the new policy and preserving the learned policies.

Algorithm 2 DaCoRL

Input: Dynamic environment $\mathcal{D} = [M_1, \dots, M_t, \dots, M_T]$.
 Single-head policy network π_θ with random weights
 $\theta = \{\theta_S^{(0)}, \theta_{\mathcal{H},1}^{(0)}\}$.

Parameter: Concentration parameter α ; Learning rate β ;
 Distillation regularization coefficient λ .

Output: Parameter set $\{\varphi_k\}_{k=1}^{K_T}$ of instantiated contexts;
 Approximate optimal policy parameters θ^* for \mathcal{D} .

- 1: **for** each time period $t \in \{1, 2, \dots, T\}$ **do**
- 2: **if** $t = 1$ **then**
- 3: Sample randomly and construct the feature vector x_1
 for the task M_1 .
- 4: Instantiate M_1 as the first context: $\varphi_1^{(1)} = x_1, z_1^* = 1$.
- 5: Set $K_1 = 1, m_1^{(1)} = 1$ in the CRP model.
- 6: Update the policy network from scratch using the
 canonical policy gradient method to obtain θ^* :

$$\theta \leftarrow \theta + \beta \nabla_\theta \mathcal{L}_{ori}.$$
- 7: **else**
- 8: Infer the task-to-context assignment z_t^* for M_t using
 incremental context detection with CRP:

$$z_t^*, \{\varphi_k^{(t)}\}_{k=1}^{K_t} \leftarrow \text{Algorithm 1}(M_t, \{\varphi_k^{(t-1)}\}_{k=1}^{K_{t-1}}, \alpha).$$
- 9: Update the CRP model according to z_t^* :

$$m_k^{(t)} = m_k^{(t-1)}, \forall k \leq K_t; m_{z_t^*}^{(t)} = m_{z_t^*}^{(t-1)} + 1.$$
- 10: **if** $z_t^* = K_{t-1} + 1$ **then**
- 11: Expand the policy network π_θ by adding an output
 head with the nearest trained head initialization,
 and add $\theta_{\mathcal{H},K_t}$ to θ .
- 12: **end if**
- 13: Set $\pi_{\theta^-} = \pi_\theta$ for distillation.
- 14: Update the policy network using the policy gradient
 method in (18) to obtain θ^* :

$$\theta \leftarrow \theta + \beta \nabla_\theta (\mathcal{L}_{ori} - \lambda \mathcal{L}_{\mathcal{D}}).$$
- 15: **end if**
- 16: **end for**
- 17: **return** $\{\varphi_k\}_{k=1}^{K_T}, \theta^*$.

The complete procedure of DaCoRL is summarized in Algorithm 2 where the agent interacts with a dynamic environment $\mathcal{D} = [M_1, \dots, M_t, \dots, M_T]$. In the first time period $t = 1$, the task M_1 is instantiated as the first context with parameter $\varphi_1^{(1)}$ (line 3). We initialize the CRP model with $K_1 = 1, m_1^{(1)} = 1$ (line 4) and train the single-head policy network from scratch using the canonical policy gradient method (line 5). In the t^{th} ($t \geq 2$) time period, we first apply the incremental context detection module to identify the context to which the current task belongs (line 7), and update the CRP model based on the identity of M_t for future context detection (line 8). Then, we employ an expandable multihead neural network for policy optimization. Specifically, when a new context is instantiated, the policy network is synchronously expanded by adding an output head initialized with the nearest trained head (lines 9-12). Next, the knowledge distillation regularization term is integrated into the loss of the original RL algorithm to reduce the catastrophic forgetting of learned tasks, and the parameters

of the policy network are updated till convergence (lines 13-14). Finally, K_T contexts with parameters $\{\varphi_k\}_{k=1}^{K_T}$ and the optimal policy π_{θ^*} are obtained for all learned tasks.

V. EXPERIMENTS AND EVALUATIONS

In this section, we conduct comprehensive experiments on several continuous control tasks from robot navigation to MuJoCo locomotion to demonstrate the effectiveness of our method. We design a variety of sequential learning tasks with diverse changes in the underlying dynamics. These problem settings are expected to be representative of the dynamic environments that RL agents may encounter in real-world scenarios. The following are the overarching questions that we aim to answer from our experiments and analysis.

- Q1 Does DaCoRL successfully achieve better continual reinforcement learning in various dynamic environments compared with existing methods?
- Q2 How does the initialization strategy of the newly added output head affect the performance of DaCoRL?
- Q3 How does the number of instantiated contexts in the latent space affect the performance of DaCoRL?
- Q4 Can DaCoRL achieve a positive forward transfer during the learning process?
- Q5 How is DaCoRL’s generalization ability to previously unseen tasks?

A. Datasets

1) *Robot Navigation* [28], [29]: It contains three types of dynamic environments with parametric variation across tasks. In Fig. 3, Types I, II, and III indicate that the dynamic environments are created in terms of parametric variation in the goal position (changes in the reward function), the puddles positions (changes in the state transition function), and both the goal and puddles positions (changes in both the reward function and state transition functions), respectively. In each navigation task, a robot agent needs to move to a goal position within a unit square. The state consists of the agent’s current 2-D position, and the action corresponds to the 2-D velocity commands in the range of $[-0.1, 0.1]$. The reward is equal to the negative squared distance to the goal position minus a small control cost that is proportional to the action’s scale. Each learning episode always starts from a given position and terminates when the agent is within 0.01 from the goal or when the episode length is greater than 100. We choose these commonly used domains as they are well-understood, suitable for highlighting the mechanism and verifying the effectiveness of our proposed method in a straightforward manner.

2) *MuJoCo Locomotion* [29]: It contains three locomotion tasks with parametric variation and growing dimensions of state-action spaces, as shown in Fig. 4. These continuous control tasks require a one-legged hopper, a planar cheetah or a 3-D quadruped ant robot to run at a particular velocity along the positive x -direction. The reward is an alive bonus plus a regular part that is negatively correlated with the absolute value between the current velocity of the agent and a preset target velocity. The dynamic environment is designed to apply parametric variations in the target velocity within a range:

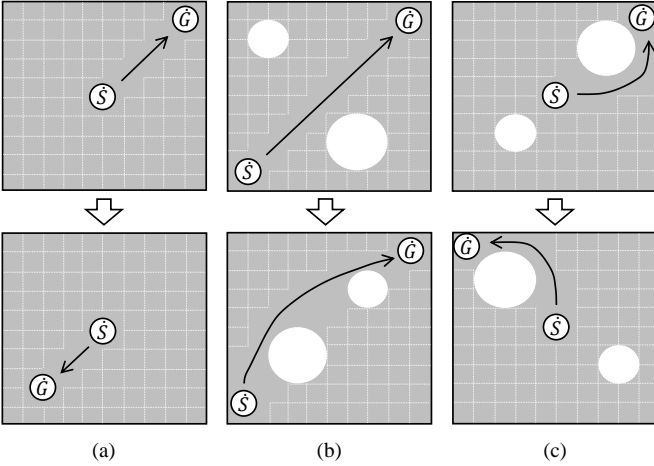


Fig. 3. Examples of three types of dynamic environments in the robot navigation tasks [28], [29]. S is the start point and G is the goal point. Puddles are shown in white. (a) Type I: the goal changes. (b) Type II: the puddles change. (c) Type III: both the goal and the puddles change.

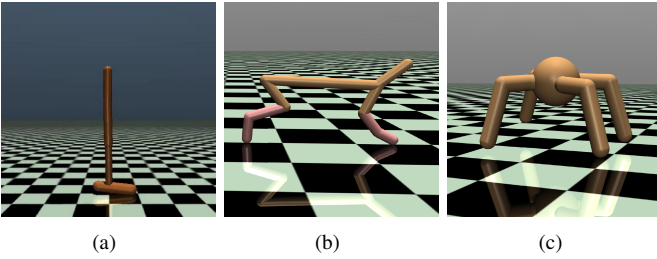


Fig. 4. Representative MuJoCo locomotion tasks [29] with growing dimensions of state-action spaces including (a) Hopper, $|\mathcal{S}| = 11$, $|\mathcal{A}| = 3$, and $r = 1 - 4 \cdot |v_x - x_g|$. (b) HalfCheetah, $|\mathcal{S}| = 20$, $|\mathcal{A}| = 6$, and $r = -|v_x - x_g|$. (c) Ant, $|\mathcal{S}| = 111$, $|\mathcal{A}| = 8$, and $r = 1 - 3 \cdot |v_x - x_g|$. v_x is the agent’s velocity in the positive x-direction and v_g is the target velocity.

$[0.0, 1.0]$ for Hopper, $[0.0, 2.0]$ for HalfCheetah, and $[0.0, 0.5]$ for Ant. Each learning episode always starts from a given physical status of the agent and terminates when the agent falls down or when the episode length is greater than 100. We choose these domains to further evaluate the efficiency of our method on more sophisticated domains.

In our experiments, to constitute the continual learning process, we uniformly generate four Gaussian clusters (representing the expected contexts in this article) for each type of the dynamic environment in its parametric variation space. Each of the first three clusters consists of 12 samples and the fourth cluster contains 14 samples, with each sample containing the values of the variable parameters of a specific task. We sequentially arrange these $T = 50$ tasks in a random order, resulting in a dynamic environment $\mathcal{D} = [M_1, M_2, \dots, M_T]$. For DaCoRL (Oracle), the supervised version of DaCoRL mentioned in the next section, we make available the task-to-context assignments based on the generation process of dynamic environments. By contrast, the correspondence between tasks and contexts is unknown and needs to be identified by DaCoRL itself. In addition, by reference to the benchmarks in the scenario of learning multiple tasks with a single model in environments with parametric variation [43], we construct ob-

servations spaces with extra dimensions of variable parameters for the above dynamic environments to meet the assumption in this article that the observations contain the discriminative information of the tasks.

B. Baselines

We evaluate our method in comparison to the following four state-of-the-art baseline methods and one supervised version of our proposed DaCoRL in dynamic environments.

- 1) *Naive*: It refers to canonical RL methods (e.g., REINFORCE [38], Proximal Policy Optimization (PPO) [44]) that simply train a policy model during the learning process, without paying attention to any possible environmental changes or forgetting.
- 2) *CRLUnsup* [36]: It detects environmental changes by observing the ability of the agent to perform the task (i.e., the difference between the actual reward and the expected one). When this difference goes below a certain threshold, it triggers the memory consolidation procedure employing Elastic Weight Consolidation (EWC) [18].
- 3) *CDKD*: It is adapted from the IQ [22] framework that can alleviate catastrophic interference for value-based RL in stationary environments. In this article, we extend IQ to policy-based RL and perform context division and knowledge distillation (renamed as CDKD), and set the true number of contexts in advance so that it can conduct continual learning in dynamic environments.
- 4) *LLIRL* [29]: It develops and maintains a library that contains an infinite mixture of parameterized environment models for lifelong adaptation in dynamic environments, which focuses on building upon the prior knowledge accumulated during previous learning to optimize the learning parameters to achieve the maximum return in the current environment, *without* considering the forgetting of previously learned policies.
- 5) *DaCoRL (Oracle)*: This approach can be regarded as a supervised version of our proposed DaCoRL, and the key difference is that the agent is informed in each time period of the specific task-to-context mapping (i.e., all context identifications are made available).

In the experiments, we use policy search with nonlinear function approximation to handle continuous control tasks in dynamic environments. For the robot navigation tasks, we perform gradient updates using the vanilla policy gradient RL algorithm REINFORCE. In addition, we employ PPO as the base RL algorithm to learn in the more challenging MuJoCo locomotion tasks.

C. Implementation

We adopt a similar network architecture for all tasks of robot navigation and MuJoCo locomotion. The policy of DaCoRL is approximated by a feed-forward neural network that contains a fully connected hidden layer (with 200 units) used as the feature extractor and an expandable output head module used as the action distribution predictor, where each output head consists of a 200-unit fully connected hidden layer and a

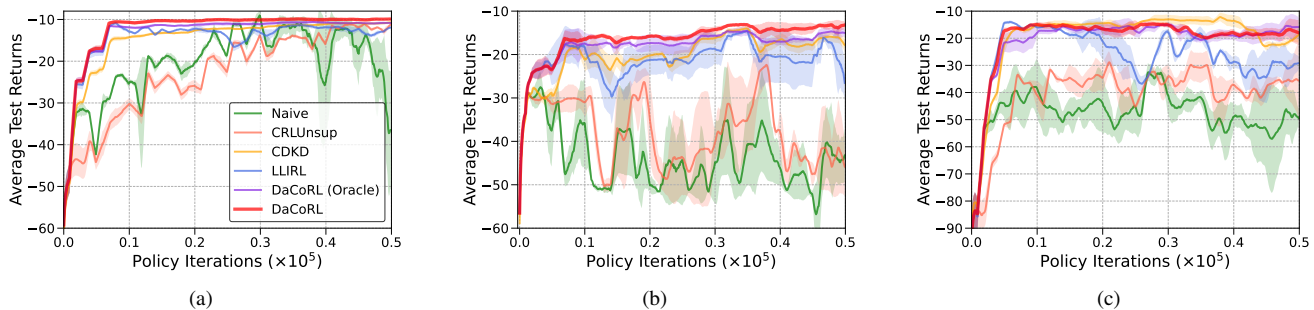


Fig. 5. Testing curves of the average returns over all tasks in the dynamic environments (\mathcal{R}_{ave}) of the robot navigation tasks. Here and in related figures in the following, K_T is the number of instantiated contexts by DaCoRL and the solid lines and shaded regions denote the means and standard deviations of average test returns, respectively, across five runs. (a) Type I, $K_T = 5$. (b) Type II, $K_T = 6$. (c) Type III, $K_T = 4$.

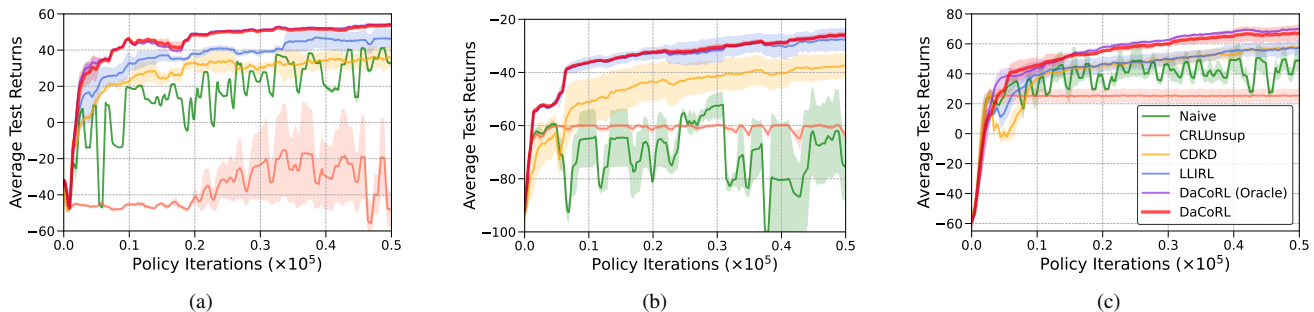


Fig. 6. Testing curves of the average returns over all tasks in the dynamic environments (\mathcal{R}_{ave}) of the MuJoCo locomotion tasks. $K_T = 4$ contexts are instantiated for DaCoRL in all tasks. (a) Hopper. (b) HalfCheetah. (c) Ant.

fully connected output layer. The hidden layers are connected by ReLU nonlinearity, following the network configuration for these tasks in [29]. For a fair comparison, the network architecture of CDKD is set to the same as that of DaCoRL and the number of contexts predetermined for CDKD is also set to be consistent with that automatically detected by DaCoRL. For Naive, CRLUnsup, and LLIRL, each policy network consists of two 200-unit hidden layers connected by ReLU nonlinearity and a fully connected output layer.

During the learning process, we train the model for 1k policy iterations on each task, and evaluate the policy performance by testing the current policy on all tasks in the dynamic environment every 100 iterations. All results reported are the average performance over five independent runs.

D. Evaluation Metrics

Following the convention in previous studies [28], [29], we define two performance metrics to systematically evaluate our proposed method. The first one is the average return over a batch of test episodes on all tasks in the dynamic environment, which is used to evaluate the overall performance of the model on all tasks after a fixed number of policy iterations during training in real time. It is defined as

$$\mathcal{R}_{ave} = \frac{1}{Tm} \sum_{i=1}^T \sum_{j=1}^m R(\tau_{ij}) \quad (19)$$

where T is the number of tasks in the dynamic environment; m is the number of episodes tested for each task; $R(\tau_{ij})$ is the cumulative reward obtained in the j^{th} test episode of task

i . The other is the average return over all test episodes, which evaluates the average performance of the model over the entire training process. It is defined as

$$\bar{\mathcal{R}}_{ave} = \frac{1}{J} \sum_{i=1}^J \mathcal{R}_{ave}^{(i)} \quad (20)$$

where J is the number of \mathcal{R}_{ave} evaluations in the learning process, and $\mathcal{R}_{ave}^{(i)}$ is \mathcal{R}_{ave} in the i^{th} evaluation.

Remark 1: For task M , the test procedure of DaCoRL is conducted in two steps: 1) Policy selection. It firstly constructs the feature vector x of M following the procedure in lines 1-2 of Algorithm 1, and the identity z^* of M can be obtained by the MAP estimate of the predictive likelihood of the K_T contexts with parameters $\{\varphi_k\}_{k=1}^{K_T}$ obtained by training. The policy corresponding to the output head z^* of the trained neural network is the final policy selected for M . 2) Policy execution. The agent applies the selected policy on M to evaluate its performance in terms of the cumulative reward in each episode.

E. Results

To investigate the effectiveness of our proposed method (Q1), we present the results of DaCoRL and all baselines in three types of dynamic environments of the robot navigation tasks. Fig. 5 shows the average episodic returns over all tasks during training according to (19), and Table II reports the numerical results in terms of the average returns over $1000/100 * 50$ tests throughout the whole training process according to (20). For DaCoRL, the numbers of instantiated

TABLE II
NUMERICAL RESULTS OF $\bar{\mathcal{R}}_{ave}$ OF ALL METHODS IN THE ROBOT NAVIGATION TASKS (BASED ON THE RESULTS IN FIG. 5. HERE AND IN RELATED TABLES, THE CONFIDENCE INTERVALS ARE STANDARD DEVIATIONS. THE BEST PERFORMANCE IS MARKED IN BOLDFACE.)

Task	Type I	Type II	Type III
Naive	-20.44 ± 1.07	-43.55 ± 2.25	-47.28 ± 3.98
CRLUnsup	-22.43 ± 0.97	-37.73 ± 1.82	-40.32 ± 1.02
CDKD	-14.34 ± 0.30	-19.86 ± 0.69	-19.13 ± 0.95
LLIRL	-14.24 ± 0.24	-21.08 ± 1.92	-25.92 ± 0.41
DaCoRL	-11.93 ± 0.43	-16.15 ± 0.30	-19.58 ± 0.08
DaCoRL (Oracle)	-12.79 ± 0.26	-17.61 ± 0.03	-20.46 ± 1.19

TABLE III
NUMERICAL RESULTS OF $\bar{\mathcal{R}}_{ave}$ OF ALL METHODS IN THE MUJoCo LOCOMOTION TASKS (BASED ON THE RESULTS IN FIG. 6.)

Task	Hopper	HalfCheetah	Ant
Naive	16.18 ± 0.32	-69.98 ± 6.44	36.39 ± 5.87
CRLUnsup	-35.87 ± 9.38	-61.08 ± 0.55	23.37 ± 4.40
CDKD	26.47 ± 3.73	-46.36 ± 6.63	41.24 ± 1.55
LLIRL	35.36 ± 0.75	-34.84 ± 3.10	42.76 ± 4.25
DaCoRL	44.07 ± 0.46	-34.11 ± 0.62	52.18 ± 1.34
DaCoRL (Oracle)	44.57 ± 0.45	-34.31 ± 0.58	54.09 ± 1.87

contexts are $K_T = 5$, $K_T = 6$, and $K_T = 4$ for the three types of navigation tasks, respectively. In Fig. 5, it is clear that DaCoRL is significantly superior to all baselines in terms of the average test return and the stability. Naive shows the worst forgetting and performance fluctuations since it does not employ any mechanism to overcome context drifts in dynamic environments. CRLUnsup is slightly better than Naive, possibly due to the EWC technique. However, in Fig. 5(b), it still suffers from wild fluctuations in performance, as the changes of environment cannot be detected just by looking at the cumulative reward curve, and learning the new task can overwrite the learned policy without triggering the memory consolidation procedure. CDKD and LLIRL perform clearly better than the above two baselines. Nevertheless, the fixed policy network structure may lead to an over-constrained optimization objective at the early stage of CDKD training, and although LLIRL employs separate neural networks to learn tasks from different context clusters, it does not consider the interference among tasks within the same cluster.

By contrast, DaCoRL achieves superior performance compared with all baselines. Due to the effective context division of environments in the latent space, along with an expandable multihead policy network and the knowledge distillation technique, DaCoRL can achieve competent continual learning of sequential tasks in dynamic environments. In particular, DaCoRL is on a par with [see Fig. 5(c)] or even outperforms [see Figs. 5(a) and 5(b)] DaCoRL (Oracle) where the task-to-context assignments are provided in advance. This results reveal that the incremental context detection module may produce more rational and valuable context division for continual learning than the supervised version with known task-to-context assignments. Furthermore, the average results in

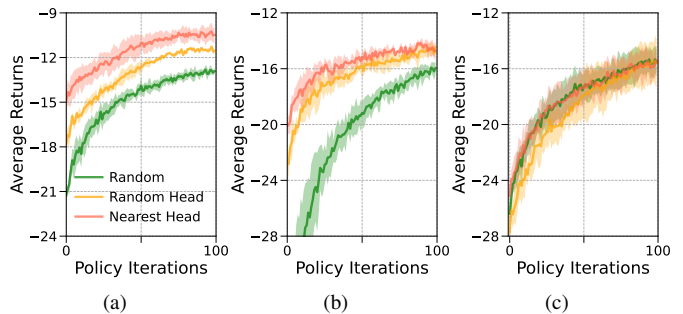


Fig. 7. Average training episodic return over all tasks per iteration of DaCoRL with different output head initialization implementations in the robot navigation tasks. (a) Type I. (b) Type II. (c) Type III.

Table II indicate that DaCoRL receives significantly larger or comparable average returns over all test episodes during the entire learning process than all baselines. Meanwhile, the statistical results also indicate that DaCoRL generally features smaller standard deviations in performance than the baselines.

To further demonstrate the scalability and flexibility of our method, we evaluate DaCoRL and all baselines on MuJoCo locomotion tasks. The average episodic returns over all tasks during training according to (19) are shown in Fig. 6 and the corresponding average returns over all tests throughout the whole learning process according to (20) are summarized in Table III. It shows that DaCoRL consistently exhibits better and more stable performance than all baseline methods, even in these complex dynamic environments. More specifically, DaCoRL instantiates totally four contexts ($K_T = 4$) for each dynamic environment, which is consistent with the given contexts in DaCoRL (Oracle). Since the learning curves of these two methods are largely coincident in all tasks in Fig. 6, it confirms that DaCoRL can accurately identify environmental context changes in a fully self-adaptive manner and can achieve comparable performance with the case where the task-to-context assignments are known in advance.

F. Analysis

1) *Influence of the Initialization of Output Heads*: To address Q2, we evaluate DaCoRL in the robot navigation tasks with different initialization strategies proposed in Section IV-B. The average returns over all tasks during the first 100 policy iterations are shown in Fig. 7. Here, we shorten the three initialization implementations to “Random”, “Random Head”, “Nearest Head”, respectively, to ensure the readability.

From Fig. 7, it is clear that the initialization by Nearest Head can enable DaCoRL to attain a better initial policy and more positive forward transfer to the new task, which is consistent with our analysis in Section IV-B. Meanwhile, Random Head exhibits obvious positive forward transfer on Type I and Type II environments and negative forward transfer on the Type III environment, compared with Random initialization. Referring to [27], [28], a possible explanation is that, in the Type III environment, Random Head often chooses an initial policy that hinders the exploration of the new task such that the agent needs to spend more time in the early training stage to counter the old policy, resulting in slow performance improvement.

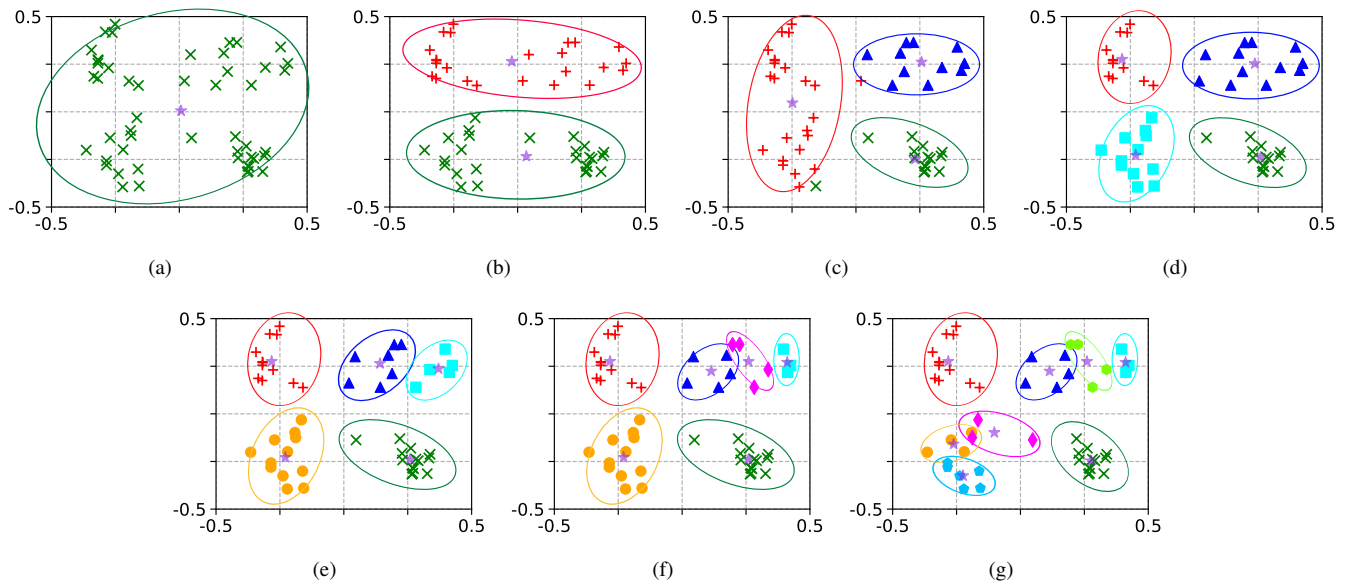


Fig. 8. The clustering patterns of contexts in DaCoRL with different numbers of instantiated contexts in the Type I navigation environment. (a) One context. (b) Two contexts. (c) Three contexts. (d) Four contexts. (e) Five contexts. (f) Six contexts. (g) Eight contexts.

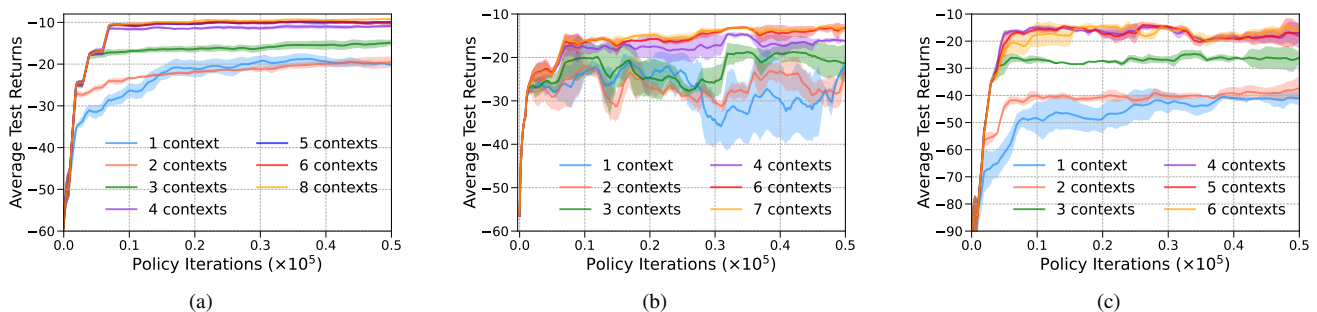


Fig. 9. The average returns over all tasks in the dynamic environments (\mathcal{R}_{ave}) of DaCoRL with different numbers of instantiated contexts in the robot navigation tasks. (a) Type I. (b) Type II. (c) Type III.

2) Influence of the Number of Instantiated Contexts (K_T):

Intuitively, instantiating a separate context for each task contained in the dynamic environment (i.e., $K_T = T$, each output head of the policy network corresponds to the policy for an individual task) is likely to result in an optimal policy for the agent. However, it is hard to operate in practice, especially when the dynamic environment comprises a large number of tasks due to the complicated network structure and challenging model training process. Thus, it is necessary to investigate the influence of the number of instantiated contexts on the performance of DaCoRL (Q3).

We vary the concentration parameter α in the CRP prior distribution to change the numbers of instantiated contexts in DaCoRL in the robot navigation tasks. Since the environmental features are highly correlated with their variation parameters, that is, intuitively, taking the Type I environment for example, the tasks with adjacent goal positions are more similar to each other and tend to belong to the same context. Hence, we use the goal position in the 2-D coordinate as a visualization to reveal the clustering patterns of contexts for the Type I navigation environment with different numbers of instantiated contexts. The results are shown in Fig. 8, where the star

TABLE IV
NUMERICAL RESULTS OF \mathcal{R}_{ave} OF DaCoRL WITH DIFFERENT NUMBERS OF INSTANTIATED CONTEXTS IN THE ROBOT NAVIGATION TASKS (BASED ON THE RESULTS IN FIG. 9.)

K_T	Type I	Type II	Type III
1	-22.86 ± 1.10	-27.34 ± 3.42	-47.16 ± 3.69
2	-22.32 ± 0.34	-26.80 ± 0.43	42.08 ± 1.35
3	-17.09 ± 0.88	-22.59 ± 2.10	-28.74 ± 0.89
4	-12.89 ± 0.30	-18.00 ± 0.96	-19.58 ± 0.08
5	-11.93 ± 0.43	—	-19.80 ± 0.42
6	-11.89 ± 0.43	-16.15 ± 0.30	-19.86 ± 0.56
7	—	-15.79 ± 0.21	—
8	-11.48 ± 0.19	—	—

shapes represent the centroids of instantiated contexts obtained from the incremental context detection procedure, and other markers represent tasks ($M_i, i \in [1, 2, \dots, T]$) in the dynamic environment. It is clear that our incremental context detection module can capture the appropriate context patterns under different parameter settings in a fully online manner, regardless of how many contexts are ultimately instantiated.

TABLE V
NUMERICAL RESULTS IN TERMS OF THE AVERAGE INITIAL PERFORMANCE OVER ALL SEQUENTIAL TASKS DURING TRAINING IN THE ROBOT NAVIGATION AND MUJoCo LOCOMOTION TASKS

Task	Type I	Type II	Type III	Hopper	HalfCheetah	Ant
Naive	-22.17 ± 1.83	-43.63 ± 2.42	-50.86 ± 5.18	8.67 ± 0.47	-75.11 ± 4.57	30.64 ± 4.39
CRLUnsup	-23.07 ± 1.23	-39.09 ± 0.58	-44.78 ± 0.90	-39.40 ± 9.42	-63.55 ± 0.47	20.91 ± 4.56
CDKD	-20.11 ± 0.60	-31.34 ± 1.77	-25.96 ± 0.86	24.42 ± 3.93	-49.66 ± 7.64	34.57 ± 2.43
LLIRL	-19.08 ± 0.51	-23.80 ± 1.71	-33.47 ± 1.36	28.16 ± 0.77	-37.63 ± 3.55	40.53 ± 3.70
DaCoRL	-14.92 ± 0.74	-20.11 ± 0.78	-25.14 ± 0.67	37.28 ± 0.46	-36.02 ± 0.13	49.07 ± 1.38

TABLE VI
NUMERICAL RESULTS IN TERMS OF THE AVERAGE TEST RETURN OVER 50 DIFFERENT TASKS THAT HAVE NOT BEEN SEEN DURING TRAINING IN ROBOT NAVIGATION AND MUJoCo LOCOMOTION TASKS

Task	Type I	Type II	Type III	Hopper	HalfCheetah	Ant
Naive	-33.78 ± 19.51	-39.82 ± 9.78	-50.19 ± 8.02	39.36 ± 0.90	-69.84 ± 12.17	55.26 ± 6.01
CRLUnsup	-11.29 ± 0.59	-42.52 ± 9.09	-38.24 ± 16.41	-43.41 ± 12.73	-62.92 ± 1.57	34.13 ± 4.34
CDKD	-11.86 ± 0.24	-22.46 ± 2.06	-18.77 ± 3.93	34.25 ± 9.38	-36.14 ± 4.05	63.26 ± 1.86
LLIRL	-14.52 ± 1.28	-39.84 ± 1.26	-35.94 ± 0.15	48.76 ± 3.39	-27.67 ± 3.85	59.39 ± 2.25
DaCoRL	-11.28 ± 0.22	-17.58 ± 0.89	-23.38 ± 4.46	53.45 ± 0.90	-26.36 ± 1.20	67.65 ± 4.24

Additionally, we show the performance of DaCoRL with different numbers of contexts in the three types of navigation environments in Fig. 9 and Table IV according to (19) and (20), respectively. Overall, the results are consistent with the intuition that DaCoRL tends to achieve better performance with more contexts instantiated. Meanwhile, maintaining fewer contexts means that tasks with large differences may be clustered into the same context, leading to serious interference during policy network training. Nevertheless, Fig. 9 and Table IV indicate that the performance gains become relatively minor when the number of contexts exceeds a certain value. Consequently, for the sake of model complexity and training cost, we recommend choosing a moderate number of instantiated contexts by adjusting the value of α . For instance, in our experiments, $K_T = 5$ for Type I, $K_T = 6$ for Type II and $K_T = 4$ for Type III are sufficient to obtain reasonable performance in these robot navigation tasks. Simultaneously, we empirically found that the concentration parameter α is insensitive for the number of instantiated contexts, and the same context detection results can be obtained by all α values over a continuous interval. For instance, in our experiments, all values of α in $[0.64, 0.93]$ for Type I, in $[0.85, 0.89]$ for Type II and in $[0.19, 0.81]$ for Type III can obtain the same number of contexts as expected above.

3) *Forward Transfer*: In DaCoRL, tasks with similar dynamics are grouped into the same context and share the same policy throughout the CRL process. In each time period, DaCoRL retrieves the policy corresponding to the most similar context as the initial policy for the current learning process, which can enable the agent to get better startup performance on each new task. To validate this feature (Q4), we calculate the average training episodic return for each algorithm at the first policy iteration on each task, and the results are shown in Table V. In comparison to all baselines, DaCoRL always achieves better initial startup performance in all experimental dynamic environments, demonstrating significant positive forward transfer.

4) *Generalization Results*: To investigate the performance of DaCoRL in unseen tasks (Q5), we randomly generate 50 different tasks that have not been seen during training in each dynamic environment, and test the policies learned by DaCoRL and all baselines in these tasks, where each policy is tested on a specific task for 100 times. The average test returns over all tasks are shown in Table VI. In general, DaCoRL features superior generalization ability in most unseen tasks compared with all baselines, except on Type III robot navigation tasks, where its generalization performance is slightly inferior to CDKD but still significantly better than other baselines. This phenomenon is explainable since CDKD is the most similar algorithm to DaCoRL and is likely to be comparable to DaCoRL especially when the number of contexts is known in advance.

VI. CONCLUSION AND FUTURE WORK

In this article, we present a continual reinforcement learning framework named DaCoRL as a viable solution to the challenge of dynamics adaptability in continuous spaces. The goal of DaCoRL is to continuously adapt the RL agent’s behavior towards the changing environment, which can be regarded as a sequence of tasks, and to minimize the catastrophic forgetting of previously learned tasks. To this end, DaCoRL employs an incremental context detection module to categorize the stream of tasks into a set of distinct contexts using online Bayesian infinite Gaussian mixture clustering. It also simultaneously optimizes a context-conditioned policy with an expandable multihead neural network, in conjunction with a knowledge distillation regularization term, to avoid the interference both between and within contexts. A key advantage of our method is that it can achieve incremental context instantiation in a fully online manner without requiring any external information to explicitly signal environmental changes. Meanwhile, it only relies on a single policy network to accomplish effective continual learning in dynamic environments. Experiments on

several continuous control tasks confirm that DaCoRL can significantly outperform state-of-the-art algorithms in terms of the stability, overall performance and generalization ability.

While REINFORCE and PPO are used as the underlying RL algorithms in the current experimental studies, our proposed framework is versatile and can be easily coupled with any other policy-based RL algorithms. A promising direction for future work is to develop an efficient strategy that can facilitate the positive backward transfer during the CRL process to make the current learning in turn further improve the performance on the learned tasks. Another direction is to address more challenging and realistic cases, such as learning in intensively changing environments where environment changes may happen between consecutive episodes [28].

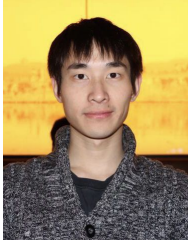
REFERENCES

- [1] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. MIT press, 2018.
- [2] C. J. C. H. Watkins and P. Dayan, “Q-learning,” *Machine Learning*, vol. 8, no. 3-4, pp. 279–292, 1992.
- [3] G. A. Rummery and M. Niranjan, “On-line q-learning using connectionist systems,” *Technical Report*, 1994.
- [4] H. Li, Q. Zhang, and D. Zhao, “Deep reinforcement learning-based automatic exploration for navigation in unknown environment,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 6, pp. 2064–2076, 2019.
- [5] V. Mnih *et al.*, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [6] O. Vinyals *et al.*, “Grandmaster level in starcraft II using multi-agent reinforcement learning,” *Nature*, vol. 575, no. 7782, pp. 350–354, 2019.
- [7] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, “Trust region policy optimization,” in *Proceedings of the International Conference on Machine Learning*. PMLR, 2015, pp. 1889–1897.
- [8] A. Faust *et al.*, “PRM-RL: Long-range robotic navigation tasks by combining reinforcement learning and sampling-based planning,” in *Proceedings of the International Conference on Robotics and Automation*, 2018, pp. 5113–5120.
- [9] H.-T. L. Chiang, A. Faust, M. Fiser, and A. Francis, “Learning navigation behaviors end-to-end with autolr,” *IEEE Robotics and Automation Letters*, vol. 4, no. 2, pp. 2007–2014, 2019.
- [10] A. Francis *et al.*, “Long-range indoor navigation with PRM-RL,” *IEEE Transactions on Robotics*, vol. 36, no. 4, pp. 1115–1134, 2020.
- [11] M. G. Bellemare *et al.*, “Autonomous navigation of stratospheric balloons using reinforcement learning,” *Nature*, vol. 588, no. 7836, pp. 77–82, 2020.
- [12] M. A. K. Jaradat, M. Al-Rousan, and L. Quadan, “Reinforcement based mobile robot navigation in dynamic environment,” *Robotics and Computer-Integrated Manufacturing*, vol. 27, no. 1, pp. 135–149, 2011.
- [13] K. Kheterpal, M. Riemer, I. Rish, and D. Precup, “Towards continual reinforcement learning: A review and perspectives,” *arXiv preprint arXiv:2012.13490*, 2020.
- [14] Y. Zheng, Z. Meng, J. Hao, Z. Zhang, T. Yang, and C. Fan, “A deep bayesian policy reuse approach against non-stationary agents,” in *Proceedings of the Conference on Neural Information Processing Systems*, vol. 31, 2018.
- [15] M. McCloskey and N. J. Cohen, “Catastrophic interference in connectionist networks: The sequential learning problem,” in *Psychology of Learning and Motivation*, 1989, vol. 24, pp. 109–165.
- [16] R. M. French, “Catastrophic forgetting in connectionist networks,” *Trends in Cognitive Sciences*, vol. 3, no. 4, pp. 128–135, 1999.
- [17] R. Hadsell, D. Rao, A. A. Rusu, and R. Pascanu, “Embracing change: Continual learning in deep neural networks,” *Trends in Cognitive Sciences*, 2020.
- [18] J. Kirkpatrick *et al.*, “Overcoming catastrophic forgetting in neural networks,” in *Proceedings of the National Academy of Sciences*, vol. 114, no. 13, 2017, pp. 3521–3526.
- [19] C. Fernando *et al.*, “Pathnet: Evolution channels gradient descent in super neural networks,” *arXiv preprint arXiv:1701.08734*, 2017.
- [20] S. Kessler, J. Parker-Holder, P. Ball, S. Zohren, and S. J. Roberts, “UNCLEAR: A straightforward method for continual reinforcement learning,” in *Proceedings of the International Conference on Machine Learning*, 2020.
- [21] S. Padakandla, K. Prabuchandran, and S. Bhatnagar, “Reinforcement learning algorithm for non-stationary environments,” *Applied Intelligence*, vol. 50, no. 11, pp. 3590–3606, 2020.
- [22] T. Zhang, X. Wang, B. Liang, and B. Yuan, “Catastrophic interference in reinforcement learning: A solution based on context division and knowledge distillation,” *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [23] D. J. Aldous, “Exchangeability and related topics,” in *École d’Été de Probabilités de Saint-Flour XIII—1983*. Springer, 1985, pp. 1–198.
- [24] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, “iCaRL: Incremental classifier and representation learning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2001–2010.
- [25] A. Rosenfeld and J. K. Tsotsos, “Incremental learning through deep adaptation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 3, pp. 651–663, 2018.
- [26] M. Hersche, G. Karunaratne, G. Cherubini, L. Benini, A. Sebastian, and A. Rahimi, “Constrained few-shot class-incremental learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 9057–9067.
- [27] Z. Wang, C. Chen, H.-X. Li, D. Dong, and T.-J. Tarn, “Incremental reinforcement learning with prioritized sweeping for dynamic environments,” *IEEE/ASME Transactions on Mechatronics*, vol. 24, no. 2, pp. 621–632, 2019.
- [28] Z. Wang, H.-X. Li, and C. Chen, “Incremental reinforcement learning in continuous spaces via policy relaxation and importance weighting,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 6, pp. 1870–1883, 2019.
- [29] Z. Wang, C. Chen, and D. Dong, “Lifelong incremental reinforcement learning with online bayesian inference,” *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- [30] L. Bottou *et al.*, “Online learning and stochastic approximations,” *Online Learning in Neural Networks*, vol. 17, no. 9, p. 142, 1998.
- [31] S. Shalev-Shwartz *et al.*, “Online learning and online convex optimization,” *Foundations and Trends® in Machine Learning*, vol. 4, no. 2, pp. 107–194, 2012.
- [32] R. Caruana, “Multitask learning,” *Machine learning*, vol. 28, no. 1, pp. 41–75, 1997.
- [33] M. Crawshaw, “Multi-task learning with deep neural networks: A survey,” *arXiv preprint arXiv:2009.09796*, 2020.
- [34] Y. Zhang and Q. Yang, “A survey on multi-task learning,” *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [35] B. C. Da Silva, E. W. Basso, A. L. Bazzan, and P. M. Engel, “Dealing with non-stationary environments using context detection,” in *Proceedings of the International Conference on Machine Learning*, 2006, pp. 217–224.
- [36] V. Lomonaco, K. Desai, E. Culurciello, and D. Maltoni, “Continual reinforcement learning in 3d non-stationary environments,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 248–249.
- [37] Y. Duan, X. Chen, R. Houthoofd, J. Schulman, and P. Abbeel, “Benchmarking deep reinforcement learning for continuous control,” in *Proceedings of the International Conference on Machine Learning*. PMLR, 2016, pp. 1329–1338.
- [38] R. J. Williams, “Simple statistical gradient-following algorithms for connectionist reinforcement learning,” *Machine learning*, vol. 8, no. 3, pp. 229–256, 1992.
- [39] Y. Seo, L. Chen, J. Shin, H. Lee, P. Abbeel, and K. Lee, “State entropy maximization with random encoders for efficient exploration,” in *Proceedings of the International Conference on Machine Learning*. PMLR, 2021, pp. 9443–9454.
- [40] I. Goodfellow *et al.*, “Generative adversarial nets,” *Proceedings of the Conference on Neural Information Processing Systems*, vol. 27, 2014.
- [41] J. Gou, B. Yu, S. J. Maybank, and D. Tao, “Knowledge distillation: A survey,” *International Journal of Computer Vision*, vol. 129, no. 6, pp. 1789–1819, 2021.
- [42] L. Wang and K.-J. Yoon, “Knowledge distillation and student-teacher learning for visual intelligence: A review and new outlooks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [43] T. Yu *et al.*, “Meta-World: A benchmark and evaluation for multi-task and meta reinforcement learning,” in *Proceedings of the Conference on Robot Learning*. PMLR, 2020, pp. 1094–1100.
- [44] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 2017.



Tiantian Zhang received the B.Sc. degree in automation from the Department of Information Science and Technology, Central South University, Changsha, China, in 2015, and the M.Sc. degree in control engineering from the Department of Automation, Tsinghua University, Beijing, China, in 2018, where she is currently pursuing the Ph.D. degree in control science and engineering.

Her research interests include data science, decision-making, and reinforcement learning.



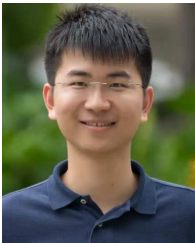
Zichuan Lin finished his Ph.D. degree from the department of Computer Science and Technology, Tsinghua University, Beijing, China, in 2021.

He is now a Researcher at Tencent, Shenzhen, China, working on sample-efficient deep reinforcement learning algorithms. He is interested in machine learning, reinforcement learning as well as their applications on game AI and natural language processing. He has been serving as a PC for NeurIPS, ICML, ICLR and AAAI.



Yuxing Wang received the B.Sc. degree in communication engineering from the Department of Electronic Information, Southwest Minzu University, Chengdu, China, in 2020. He is currently pursuing the M.Sc. degree in electronic information with the Shenzhen International Graduate School, Tsinghua University, Shenzhen, China.

His research interests include evolutionary computation, reinforcement learning and embodied AI.



Deheng Ye finished his Ph.D. from the School of Computer Science and Engineering, Nanyang Technological University, Singapore, in 2016.

He is now a Principal Researcher and Team Manager with Tencent, Shenzhen, China, where he leads a group of engineers and researchers developing large-scale learning platforms and intelligent AI agents. He is interested in applied machine learning, reinforcement learning, and software engineering. He has been serving as a PC/SPC for NeurIPS, ICML, ICLR, AAAI, and IJCAI.



Qiang Fu received the B.S. and M.S. degrees from the University of Science and Technology of China, Hefei, China, in 2006 and 2009, respectively.

He is the Director of the Game AI Center, Tencent AI Lab, Shenzhen, China. He has been dedicated to machine learning, data mining, and information retrieval for over a decade. His current research focus is game intelligence and its applications, leveraging deep learning, domain data analysis, reinforcement learning, and game theory.



Wei Yang received the M.S. degree from the Huazhong University of Science and Technology, Wuhan, China, in 2007.

He is currently the General Manager of the Tencent AI Lab, Shenzhen, China. He has pioneered many influential projects in Tencent in a wide range of domains, covering Game AI, Medical AI, search, data mining, large-scale learning systems, and so on.



Xueqian Wang (Member, IEEE) received the M.Sc. and Ph.D. degrees in control science and engineering from the Harbin Institute of Technology (HIT), Harbin, China, in 2005 and 2010, respectively.

From June 2010 to February 2014, he was a Post-Doctoral Researcher with HIT. From March 2014 to November 2019, he was an Associate Professor with the Division of Informatics, Shenzhen International Graduate School, Tsinghua University, Shenzhen, China. He is currently a Professor and the Leader of the Center for Artificial Intelligence and Robotics, Shenzhen International Graduate School, Tsinghua University. His research interests include robot dynamics and control, teleoperation, intelligent decision-making and game playing, and fault diagnosis.



Bin Liang (Senior Member, IEEE) received the B.Sc. and M.Sc. degrees in control engineering from the Honors College, Northwestern Polytechnical University, Xi'an, China, in 1989 and 1991, respectively, and the Ph.D. degree in control engineering from the Department of Precision Instrument, Tsinghua University, Beijing, China, in 1994.

From 1994 to 2003, he held his positions as a Post-Doctoral Researcher, an Associate Researcher, and a Researcher with the China Academy of Space Technology (CAST), Beijing. From 2003 to 2007, he held his positions as a Researcher and an Assistant Chief Engineer with the China Aerospace Science and Technology Corporation, Beijing. He is currently a Professor with the Research Center for Navigation and Control, Department of Automation, Tsinghua University. His research interests include modeling and control of intelligent robotic systems, teleoperation, and intelligent sensing technology.



Bo Yuan (Senior Member, IEEE) received the B.E. degree in computer science from the Nanjing University of Science and Technology, Nanjing, China, in 1998, and the M.Sc. and Ph.D. degrees in computer science from The University of Queensland (UQ), St Lucia, QLD, Australia, in 2002 and 2006, respectively.

From 2006 to 2007, he was a Research Officer on a project funded by the Australian Research Council, UQ. From 2007 to 2021, he was a faculty member (Lecturer: 2007-2009 and Associate Professor: 2009-2021) in the Division of Informatics, Tsinghua Shenzhen International Graduate School, P.R. China, and served as the Deputy Director of Office of Academic Affairs (2013-2020). In July 2021, he co-founded Shenzhen Wisdom & Strategy Technology Co., Ltd., an innovative K-12 AI education solution provider. He has authored or coauthored more than 110 papers in refereed international conferences and journals. His research interests include data science, evolutionary computation, and reinforcement learning.



Xiu Li (Member, IEEE) received her Ph.D. degree in computer integrated manufacturing from the Nanjing University of Aeronautics and Astronautics, Nanjing, China, in 2000.

From 2003 to 2010, she was an Associate Professor with the Department of Automation, Tsinghua University, Beijing, China. From 2010 to 2016, she was an Associate Professor with the Division of Informatics, Shenzhen International Graduate School, Tsinghua University, Shenzhen, China. She is currently a Professor and the director of the Division of Informatics, Shenzhen International Graduate School, Tsinghua University. Her research interests include intelligent system, pattern recognition, and data mining.