

Composable Energy Policies for Reactive Motion Generation and Reinforcement Learning

Julen Urain*, Anqi Li†, Puze Liu*, Carlo D’Eramo*, Jan Peters*

*IAS, TU Darmstadt, †University of Washington

Email: {julen, puze, carlo}@robot-learning.de, anqi14@cs.washington.edu, jan.peters@tu-darmstadt.de

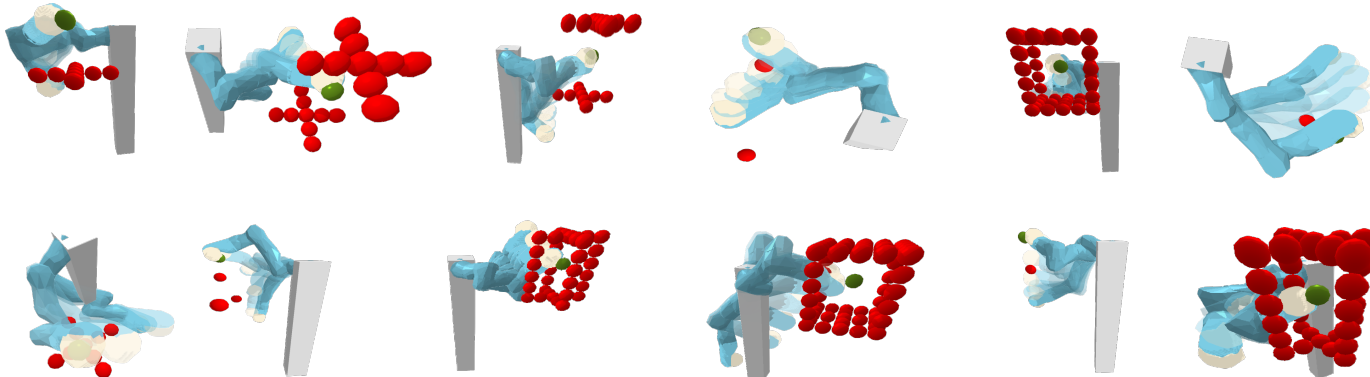


Fig. 1. Composable Energy Policies generate trajectories reactively in a wide set of environments, as the ones in the figure. The generated trajectories are attracted to a target pose, visualized by the green sphere, and avoid a set of obstacles, shown as red spheres.

Abstract—Reactive motion generation problems are usually solved by computing actions as a sum of policies. However, these policies are independent of each other and thus, they can have conflicting behaviors when summing their contributions together. We introduce Composable Energy Policies (CEP), a novel framework for modular reactive motion generation. CEP computes the control action by optimization over the product of a set of stochastic policies. This product of policies will provide a high probability to those actions that satisfy all the components and low probability to the others. Optimizing over the product of the policies avoids the detrimental effect of conflicting behaviors between policies choosing an action that satisfies all the objectives. Besides, we show that CEP naturally adapts to the Reinforcement Learning problem allowing us to integrate, in a hierarchical fashion, any distribution as prior, from multimodal distributions to non-smooth distributions and learn a new policy given them. Video in <https://sites.google.com/view/composable-energy-policies/home>

I. INTRODUCTION

Many robotic tasks deal with finding a control action satisfying multiple objectives. An apparently simple task such as watering some plants requires satisfying multiple objectives to perform it properly. The robot should reach the targets (the plants) with the watering can, avoid pouring water on the floor while approaching, and avoid colliding and break plant’s branches with its arms. In contrast with more sequential tasks [42, 17, 18, 39], in which the objectives are satisfied concatenating them in time, in the presented work, we consider tasks in which multiple objectives must be satisfied in parallel.

The problem has been faced with a spectrum of solutions that balance between global optimality and computational complexity. Path planning methods [26, 25, 20] find a global

trajectory from start to goal by a computationally intense Monte-Carlo sampling process. Trajectory optimization methods [45, 32, 19, 37] reduce the computational burden of planning methods by learning the global trajectory given an initial trajectory candidate. These methods reshape the global trajectory to satisfy the objectives. However, they still require solving an optimization problem over long temporal horizon trajectories. Reactive Motion Generators, such as Artificial Potential Fields (APF) methods [21, 23, 11, 22, 33, 4] have a very low computational cost, but lack any guarantees of finding a global trajectory satisfying the objectives. These methods propose a modular approach. Each component proposes a deterministic policy to satisfy one of the objectives, and the control action is obtained by the sum of the policy contributions.

There is a clear gap in methodology. The first (path planning and trajectory optimization), approach the problem as an optimization or inference problem over the combined objectives, while the second (reactive motion generation) assumes a complete independence between objectives, lacking any optimality guarantees. We frame Reactive Motion Generation as an optimization problem over a product of expert policies [14, 31, 30, 12, 43]

$$\mathbf{a}^* = \arg \max_{\mathbf{a}} \prod_k \pi_k(\mathbf{a} | \mathbf{s}). \quad (1)$$

Most of the previous approaches for reactive motion generation, compute each action maximizing a certain objective independently and then, sum the actions together [33, 21, 11, 15]. This independence between policies can end up in conflicting behaviors leading to oscillations or local minima. In our

approach, we first compute the product of a set of stochastic policies, and then, we compute the action maximizing the product of them. This approach can be understood as a probabilistic logical conjunction (AND operator) [8, 43] between the stochastic policies (see Fig. 2 for visual representation).

Nevertheless, there are exceptions. Dynamic Window Approach (DWA) [10, 47] solves the reactive motion generation for a 2D planar robot in a two steps optimization algorithm. First, the search space of possible actions is reduced given a set of constraints. Then, given an objective function, the optimal action is computed solving an optimization problem. In contrast with DWA, in our work, we compute the optimal action only through an optimization phase. Additionally, we consider 7-dof robots instead of planar robots and allow to integrate objectives defined in arbitrarily different task spaces. Finally, given the flexibility of our method, we can integrate policies from multiple sources in a single model. We could integrate obstacle avoidance policies with data-driven learned policies, compute a solution by trajectory optimization and add reaction to unexpected situations or integrate priors in a Reinforcement Learning (RL) policy to accelerate learning and ensure safety.

Contribution: We present Composable Energy Policies (CEP), a new framework for modular robot control. In contrast with previous methods [21, 33, 35] that solves reactive motion generation with a sum of deterministic policies, our method solves reactive motion generation by maximum likelihood estimation over a product of expert policies. We claim that optimizing over the composed policy distributions will increase the guarantees of satisfying all the objectives jointly w.r.t. adding actions coming from different policies. To validate our claims, we derive the reactive motion generation problem from control as inference view and evaluate the optimality guarantees of applying a product of expert policies. In (Appendix A), we show the deviation between the optimal policy and CEP w.r.t. the planning horizon.

In our work, we aim for flexible composition of modular policy distributions. We propose a framework to integrate stochastic policies represented in arbitrary state-action spaces. Additionally, these policies might have multiple sources from data-driven learned multimodal policies to computationally costly control as inference solutions. In section V, we show how to integrate a RL policy with a set of prior policies to improve exploration and guarantee safety exploration.

In the following, we focus on presenting the method’s main blocks while in the Appendix, we focus on the theoretical interpretation of the method. In section VI, we validate our algorithm in both reactive motion generation and prior based RL problems.

Notation: As our discussion will involve a set of policies and a set of spaces in which these policies are represented, we will use superscript (π^x) to represent the space in which the policy is and subscript (π_x) to represent the policy index. $\mathbf{f}_x^z(\cdot)$ represent a transformation map from space \mathcal{X} to space \mathcal{Z} , \mathbf{s}^x is the state in space \mathcal{X} and \mathbf{s}^z , the state in space \mathcal{Z} .

II. PRELIMINARIES

A. Artificial Potential Fields

Reactive motion generation deals with the problem of generating local and quick motion. The developed methods require to have low computational burden to give a fast response reacting to unexpected situations. APF [21] propose to solve the problem by a weighted sum of a set of dynamic systems represented in a set of task spaces. Let’s assume a set of transformations maps $\mathbf{f}_x^{z_0}, \dots, \mathbf{f}_x^{z_K}$, that transform a point in the configuration space, \mathcal{X} , to a set of task spaces, \mathcal{Z}_k , $(\mathbf{z}_k, \dot{\mathbf{z}}_k) = (\mathbf{f}_x^{z_k}(\mathbf{x}), \mathbf{J}^{z_k}(\mathbf{x})\dot{\mathbf{x}})$, with $\mathbf{J}^{z_k} = \frac{\partial \mathbf{f}_x^{z_k}(\mathbf{x})}{\partial \mathbf{x}}$; and a set of deterministic second-order dynamic systems in the task space, $\mathbf{g}^{z_0}, \dots, \mathbf{g}^{z_k}$. The APF represents the dynamic system in the configuration space \mathcal{X}

$$\mathbf{g}^x = \sum_{k=0}^K \mathbf{J}^{z_k^\dagger} \mathbf{\Lambda}_k(\mathbf{z}_k) \mathbf{g}^{z_k} \quad (2)$$

with a weighted ($\mathbf{\Lambda}_k(\mathbf{z}_k)$) sum of the of the projected dynamics and $\mathbf{J}^{z_k^\dagger}$ the pseudoinverse Jacobian. Each dynamic component represent the policy to satisfy a particular objective. Given that we are looking for the motion that satisfy multiple objectives, the sum of the dynamics might not be enough as the different components might have conflicting behaviours and thus lead to local minima or oscillations. In Appendix C, we show that APF can be rewritten as a particular case of product of experts and provide an interpretation of their failures.

We aim to solve the conflicts between dynamics by the maximum likelihood estimation over the product of a set of dynamic distributions. In our work, we model each dynamic component by a distribution. These distributions can be understood as a weighted set of possible dynamics for each component. In the combination of these distributions, our method will search for the dynamics that better satisfy the combined set of objectives and thus, avoid conflicting behaviors, as shown in Figure 2.

III. COMPOSABLE ENERGY POLICIES

Let us assume a set of independent stochastic policies $\pi_1(\mathbf{a}|\mathbf{s}), \dots, \pi_K(\mathbf{a}|\mathbf{s})$ modeled by a Boltzmann distribution

$$\pi_i(\mathbf{a}|\mathbf{s}) = \exp(E_i(\mathbf{a}, \mathbf{s})) \frac{1}{Z_i(\mathbf{s})} \quad (3)$$

with $E : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is an arbitrarily represented energy function and $Z(\mathbf{s}) = \int_{\mathbf{a}} \exp(E(\mathbf{a}; \mathbf{s})) d\mathbf{a}$ is the normalization factor. Choosing a Boltzmann distribution is not an arbitrary choice. Boltzmann distribution allows representing an arbitrary distribution by a suitable definition of the energy function (E) [27]. Additionally, computing the product of experts

$$\pi(\mathbf{a}|\mathbf{s}) = \prod_{k=0}^K \pi_k(\mathbf{a}|\mathbf{s})^{\beta_k} \propto \exp\left(\sum_k \beta_k E_k(\mathbf{a}, \mathbf{s})\right). \quad (4)$$

will end up in a weighted sum over the individual energy components in the log-space. Linearly combined energies are beneficial for computing the energies contribution in parallel

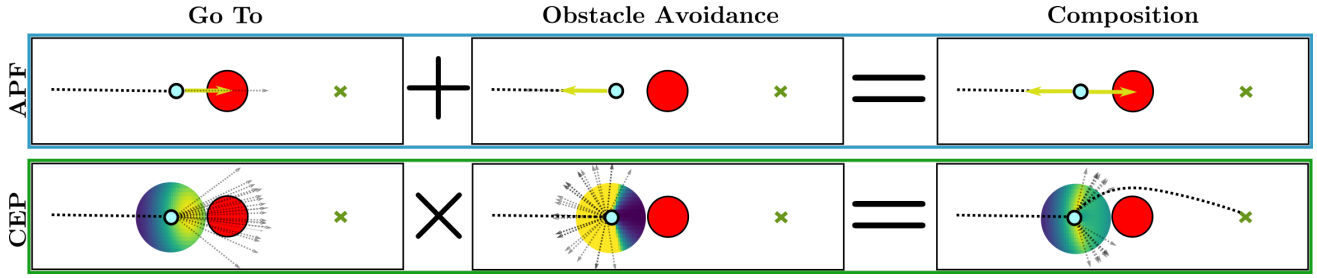


Fig. 2. Visual Representation of modular control for Goto + Obstacle Avoidance. In the top box, we show Artificial Potential Field (APF) [21] policy. In the bottom box, we show Composable Energy Policies (CEP). In contrast, with artificial potential field method, that sums deterministic actions (goto, avoid obstacle), CEP computes the product of the policy distributions. The composition will provide high probability to those actions that satisfy both components and low to the rest. This approach helps avoiding conflicts between different components. Robot: blue circle, obstacle: red circle and target: green cross. Thick dotted line: performed trajectory, light dotted line: possible future trajectories.

and thus, in practise, parallelize the computation by multi-processing, increasing the control frequency. Product of experts is a natural choice for our policy; given a set of energies E_1, \dots, E_K , product of experts is the posterior of an inference problem maximizing the energies (See Appendix B).

A. Energy Tree

In the composition proposed in (4), each energy function is considered to be in the same state-action space. However, in most of the robotics scenarios, we might be interested in composing together energies defined in different task spaces. Navigation of the robot's end-effector towards a target while avoiding the obstacles, composes skills defined in different task spaces. Each robot link should avoid obstacles and the end effector should reach a certain target position. Inspired by APF [21] and Riemannian Motion Policies (RMP) [33], we propose to model the composition of energies in different task spaces. We introduce the probabilistic graphical model for our policy in Fig.3.

Our architecture is composed of two main components. First, we have a set of policies $\pi^z(\mathbf{a}^z | \mathbf{s}^z)$, defined in different state-action latent spaces (task space) $(\mathbf{s}^z, \mathbf{a}^z) \in \mathcal{Z}$. Second, we consider a set of deterministic mappings that transform the state-actions in the common space (configuration space) $(\mathbf{s}^x, \mathbf{a}^x) \in \mathcal{X}$ to the latent state-action spaces \mathcal{Z} , $\mathbf{f}_x^z: \mathcal{X} \rightarrow \mathcal{Z}$.

Applying the change of variable rule in probabilistic density functions [2, 3], given the mapping \mathbf{f}_x^z is bijective in the actions ($\mathbf{a}^z \leftrightarrow \mathbf{a}^x$), we can write the policy in the common space, π^x , in terms of the latent space policy π^z [46]

$$\begin{aligned} \pi^x(\mathbf{a}^x | \mathbf{s}^x) &= \pi^z(\mathbf{a}^z | \mathbf{s}^z) \sqrt{\det(\mathbf{J}^\top \mathbf{J})} \\ &= \pi^z(\mathbf{f}_x^z(\mathbf{s}^x, \mathbf{a}^x)) \sqrt{\det(\mathbf{J}^\top \mathbf{J})}. \end{aligned} \quad (5)$$

with $\mathbf{J} = \partial \mathbf{f}_x^z(\mathbf{a}^x, \mathbf{s}^x) / \partial \mathbf{a}^x$ is the Jacobian of \mathbf{f}_x^z in the action \mathbf{a}^x . The equality in (5) does not hold for injective or surjective transformations. In practice, we apply the equality for any possible map \mathbf{f}_x^z (bijective, surjective, and injective), and we leave for future work the study of the policy distribution under surjective or injective mappings.

From (5), we can write $\pi^x(\mathbf{a}^x | \mathbf{s}^x)$ w.r.t. the energy function defined in the latent space $E^z(\mathbf{a}^z; \mathbf{s}^z)$

$$\pi^x(\mathbf{a}^x | \mathbf{s}^x) = \exp(E^z(\mathbf{f}_x^z(\mathbf{a}^x, \mathbf{s}^x))) \frac{1}{Z(\mathbf{s}^x)} \quad (6)$$

with

$$Z = \frac{\int_{\mathbf{a}^z} \exp(E^z(\mathbf{a}^z, \mathbf{s}^z)) d\mathbf{a}^z}{\sqrt{\det(\mathbf{J}^\top \mathbf{J})}}. \quad (7)$$

This policy representation allows to define the energy function in an easy to represent task space. Then, for any $\mathbf{s}^x, \mathbf{a}^x \in \mathcal{X}$, we can map them to the latent space \mathcal{Z} , $(\mathbf{s}^z, \mathbf{a}^z) = \mathbf{f}_x^z(\mathbf{s}^x, \mathbf{a}^x)$ and compute the un-normalized log-probability of the action in a given latent policy $E^z(\mathbf{a}^z, \mathbf{s}^z)$.

Our energy tree is inspired by the APF formulation, but it has important differences. In APF formulation, there is an explicit function from state \mathbf{s}^x to action \mathbf{a}^x

$$\mathbf{a}^x = \mathbf{f}_x^{z-1}(\mathbf{a}^z) = \mathbf{f}_x^{z-1}(\mathbf{g}(\mathbf{s}^z)) = \mathbf{f}_x^{z-1}(\mathbf{g}(\mathbf{f}_x^z(\mathbf{s}^x))) \quad (8)$$

where \mathbf{f}_x^z is the function mapping from common space to latent space and \mathbf{g} is a deterministic policy in the latent space. In this approach, \mathbf{f}_x^z is required to be invertible. In contrast, CEP assumes an implicit function from \mathbf{s}^x to \mathbf{a}^x

$$\mathbf{a}^x = \arg \max_{\mathbf{a}^x} E^z(\mathbf{f}_x^z(\mathbf{a}^x, \mathbf{s}^x)). \quad (9)$$

In this approach, we do not need to compute the inverse of \mathbf{f}_x^z . Additionally, E^z represents the energy in the latent space, so we can combine distributions instead of combining deterministic dynamics. In Sec. IV we show that having an implicit function might be beneficial for reactive motion generation.

In case a set of policies exist in different task spaces $\mathcal{Z}_1, \dots, \mathcal{Z}_N$, we can compose the policies together

$$\pi^x(\mathbf{a}^x | \mathbf{s}^x) = \exp\left(\sum_{k=1}^K \beta_k E^{z_k}(\mathbf{f}_x^{z_k}(\mathbf{a}^x, \mathbf{s}^x))\right) \frac{1}{Z(\mathbf{s}^x)} \quad (10)$$

Shown in Fig. 3, CEP has a recursive architecture. Our model allows us to represent the latent policies $\pi^z(\mathbf{a}^z | \mathbf{s}^z)$ with a set of policies in an even deeper latent space \mathcal{W} , and still compute the optimization in the same way.

B. Maximum Likelihood Estimation on CEP

In CEP, the optimal action is obtained by Maximum Likelihood Estimation over the product of expert policies. Given the current state \mathbf{s}^x

$$\mathbf{a}_{\text{ML}}^x = \arg \max_{\mathbf{a}^x} \prod_{k=0}^K \pi_k(\mathbf{a}^x | \mathbf{s}^x)^{\beta_k}, \beta_k > 0. \quad (11)$$

the weighted sum of deterministic policies (Artificial Potential Fields methods), can be framed as a particular case of CEP.

Go To a target: The simplest distribution model for going to a target position is a normal distribution

$$\pi(\ddot{\mathbf{x}}|\mathbf{x}, \dot{\mathbf{x}}) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}, \dot{\mathbf{x}}), \boldsymbol{\Sigma}(\mathbf{x})) \quad (16)$$

with

$$\boldsymbol{\mu}(\mathbf{x}, \dot{\mathbf{x}}) = -\mathbf{K}_p(\mathbf{x} - \mathbf{x}^*) - \mathbf{K}_v\dot{\mathbf{x}} \quad (17)$$

$$\boldsymbol{\Sigma}(\mathbf{x}) = \alpha(\mathbf{x} - \mathbf{x}^*)^\top(\mathbf{x} - \mathbf{x}^*). \quad (18)$$

The distribution proposes as mean $\boldsymbol{\mu}$, a PD-controller with \mathbf{x}^* the target state and $\mathbf{K}_p > 0$ and $\mathbf{K}_v > 0$ gains and a covariance matrix $\boldsymbol{\Sigma}$ represented as the quadratic distance to the target, scaled by $\alpha > 0$. The variance will shrink closer to the target and, thus, the relevance of this energy component will increase the closer we are to the desired state.

Obstacle Avoidance: We represent obstacle avoidance energy in the unidimensional space represented by the vector between the cartesian robot position \mathbf{x}_r in task space and the cartesian obstacle position \mathbf{x}_o .

We first compute the vector pointing to the obstacle

$$\mathbf{v}_{r,o} = \mathbf{x}_r - \mathbf{x}_o \quad (19)$$

Then, we project the velocity $\dot{\mathbf{x}}$ and the acceleration $\ddot{\mathbf{x}}$ vectors in the task space

$$\begin{aligned} \dot{\mathbf{x}}_p &= \dot{\mathbf{x}} \frac{\dot{\mathbf{x}} \cdot \mathbf{v}_{r,o}}{\|\dot{\mathbf{x}}\| \|\mathbf{v}_{r,o}\|} \\ \ddot{\mathbf{x}}_p &= \ddot{\mathbf{x}} \frac{\ddot{\mathbf{x}} \cdot \mathbf{v}_{r,o}}{\|\ddot{\mathbf{x}}\| \|\mathbf{v}_{r,o}\|} \end{aligned} \quad (20)$$

Finally, we compute the energy for the acceleration $\ddot{\mathbf{x}}_p$, $\dot{\mathbf{x}}_p$ and $\mathbf{v}_{r,o}$ as

$$E(\ddot{\mathbf{x}}_p, \dot{\mathbf{x}}_p, \mathbf{v}_{r,o}) = \begin{cases} 0 & \text{if } \dot{\mathbf{x}}_p < 0 \text{ or } \|\mathbf{v}_{r,o}\| > \gamma \\ -\infty & \text{if } \ddot{\mathbf{x}}_p > -\alpha\dot{\mathbf{x}}_p - \beta \\ 0 & \text{otherwise} \end{cases} \quad (21)$$

with a distance threshold parameter $\gamma > 0$ to activate the energy component. If the projected robot velocity is not pointing to the obstacle or if the robot's task space position is too far from the robot, the distribution is a uniform distribution over the whole acceleration space. In case the robot is close to the obstacle and the velocity vector points to the obstacle, the distribution then becomes a uniform distribution between $-\infty$ and $-\alpha\dot{\mathbf{x}}_p - \beta$, $\mathcal{U}_p(-\infty, -\alpha\dot{\mathbf{x}}_p - \beta)$ with $\alpha > 0$ and $\beta > 0$.

Avoid Joint Limits: Similarly to the collision avoidance energy, we model the joint limits avoidance with a uniform distribution that provides high probability to those accelerations pointing in the opposite direction to the joint limit as long as the joint is close to the limits. Given the distance to the joint limit

$$d_l = q - q_l \quad (22)$$

with q_l the joint limit, the energy can be represented as

$$E(\ddot{q}; \dot{q}, q) = \begin{cases} 0 & \text{if } \frac{d_l}{|d_l|}\dot{q} < 0 \text{ or } |d_l| > \gamma \\ -\infty & \text{if } \frac{d_l}{|d_l|}\dot{q} > -\alpha\frac{d_l}{|d_l|}\dot{q} - \beta \\ 0 & \text{otherwise} \end{cases} \quad (23)$$

with $\alpha > 0$ and $\beta > 0$.

Imitation Learning, Control-as-Inference and RL: CEP architecture is not only limited to heuristic energies. The architecture allows to integrate any form distributions. CEP architecture allows to integrate multimodal policies learned by behavioural cloning [46, 30], posterior distributions computed by control as inference [24, 49, 50] or a RL policy.

In the following, we propose a novel hierarchical reinforcement learning policy to integrate an RL agent with a set of prior distributions.

V. CEP FOR HIERARCHICAL REINFORCEMENT LEARNING

Priors have been widely used in RL, accelerating the learning process and therefore improving the sample efficiency. Most popular approaches can be categorized in two groups. On the one hand, methods that first approximate a parameterized policy to mimic the priors

$$\min_{\theta} D_{\text{KL}}(\pi_{\theta}(\mathbf{a}|\mathbf{s}) || \pi_q(\mathbf{a}|\mathbf{s}))$$

and then, update this policy to maximize for a new task [29, 46].

$$\max_{\theta} \mathbb{E}_{p_{\pi_{\theta}}(\mathbf{s}, \mathbf{a})} [\mathcal{R}(\mathbf{s}, \mathbf{a})].$$

Even if this approach is the most popular one, if the parameterized policy is not sufficiently expressive is not going to cover all the prior information. Moreover, the parameterized model might forget important information coming from the prior after some RL updates. On the other hand, Hierarchical Reinforcement Learning (HRL) methods [16, 40, 41, 6, 31] consider a two-layered policy.

$$\pi^L(\mathbf{a}|\mathbf{s}) = \int_{\mathbf{a}_H} \pi_q(\mathbf{a}|\mathbf{s}, \mathbf{a}_H) \pi_{\theta}^H(\mathbf{a}_H|\mathbf{s}) d\mathbf{a}_H \quad (24)$$

Given some priors, a policy π_q is modeled conditioned on a higher-order action \mathbf{a}_H . For sampling an action, we first sample an action from the high-level policy, $\mathbf{a}_H \sim \pi_{\theta}^H(\mathbf{a}_H|\mathbf{s})$ and then, we sample from the conditioned low-level policy, $\mathbf{a} \sim \pi_q(\mathbf{a}|\mathbf{s}, \mathbf{a}_H)$. Most of the previous HRL policies assume rather a mixture of skills [31, 13, 29, 9] or a simple sum of the prior actions [16, 40] as HRL policies. These policy architectures might not satisfy the desired properties in a real robot RL problem, rather because they are not safe enough or because they are so limited in their exploration space.

In a robot learning scenario, a desired HRL policy should consider (1) priors that guarantee hard constraints to avoid collisions, (2) multimodal/non-normal priors to provide a rich set of possible exploring regions, (3) allow to integrate multiple prior sources in the same policy and (4) do not limit the exploration region to a few set of given skills. In the

following, we show that CEP allows integrating all the desired properties in the HRL policy.

In our work, we propose to model the low level policy with a CEP

$$\begin{aligned} \mu_H, \Sigma_H &\sim \pi_\theta^H(\mu_H, \Sigma_H | s) \\ \pi^L(\mathbf{a} | s, \mu_H, \Sigma_H) &= \pi_q(\mathbf{a} | s) \mathcal{N}(\mathbf{a} | \mu_H, \Sigma_H) \end{aligned} \quad (25)$$

In our model, the action is computed by the maximization over the low-level policy

$$\mathbf{a} = \arg \max_{\mathbf{a}} \pi^L(\mathbf{a} | s, \mu_H, \Sigma_H). \quad (26)$$

Our proposed policy combines the given prior policy distribution with a normal distribution parameterized with the high-level policy. This structure allows integrating any type of prior distribution with a learnable policy. The prior could impose hard constraints to avoid actions that lead to collisions by setting $\pi_q(\mathbf{a} | s) = 0$ or encourage the robot to explore in a mixture of regions by setting a mixture of distributions as prior.

The variance Σ_H control controls the importance of the learning policy w.r.t. the priors. The bigger Σ_H is the less relevant the learning policy is and thus, the higher the influence of the priors in the action; in the limit $\lim_{\Sigma_H \rightarrow 0} \pi^L(\mathbf{a} | s, \mu_H, \Sigma_H) = \delta(\mu_H - \mathbf{a})$, the prior effect disappear and we only trust in the learning policy.

It is interesting to remark that our model applies RL in the high-level policy and solves a stochastic search optimization problem in the low-level policy. With the low-level stochastic search optimization, we can guarantee that some conditions (obstacle avoidance) will be satisfied in the robot action, while, a simple residual control could not.

VI. EXPERIMENTAL EVALUATION

The model evaluation is split into two parts. First, we study CEP performance for reactive motion generation. We evaluate how the different energy components perform by observing the success rate and the collisions in a set of obstacle avoidance environments. We compare CEP w.r.t. previous modular reactive controllers. We exclude trajectory optimization and path planning algorithms as our evaluation is interested in reactive controllers that can provide a control response in high control frequencies ($\geq 500\text{Hz}$).

Second, we evaluate CEP in an RL problem. We want to observe if using CEP as prior boosts the learning performance of an RL agent. Additionally, we want to evaluate if the proposed priors are sufficiently strong to impose safety priors in the policy and avoid collisions. We compare our approach to previous methods applying both *Behavioural Cloning + RL* and *Hierarchical RL*.

Reactive Motion Generation

We consider the problem of controlling a robot manipulator in a set of obstacle avoidance environments. We use a 7 dof KUKA-LWR robot manipulator and the objective is to reach a target pose (3D and 6D) while avoiding colliding

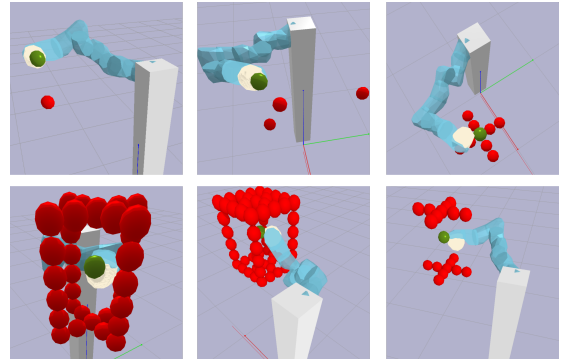


Fig. 4. Reactive Motion Generation Environments. In the top, from left to right: 1 Obstacle, 3 Obstacles, Cross. In the bottom, from left to right: Cage I, Cage II, Double Cross. Green Sphere: target pose, red spheres: obstacles.

with the obstacles. We consider increasingly difficult environments presented in Fig. 4. The robot is initialized in a random joint configuration and its motion is generated by a set of modular reactive motion generators, without additional global path planning algorithms. We evaluate the success rate and the collision rate of our method and compare it w.r.t. deterministic modular reactive motion controllers (Artificial Potential Fields [21] and Riemannian Motion Policies [33]) as baselines.

We build the CEP model with

- Go-To energy in the end effector’s task space.
- Obstacle Avoidance energy in each robot link’s task space.
- Joint limits avoidance energy in the configuration space.

To guarantee high frequency control actions ($\geq 500\text{Hz}$), we solve the optimization problem (11) in a Nvidia GPU RTX-2800. Each energy component is parallelized and the action samples are evaluated in batch. With a sufficient amount of samples per optimization step (10.000 samples), the optimization is solved in two steps.

Results and analysis: We summarize the obtained results for the 3D Goto problem in table I. Few obstacles environments are easily solved by all the methods. We observe a success rate of almost 100% for the first three environments in the three cases. This result shows that simple scenarios can be easily solved with local reactive controllers and it is not required to solve a global trajectory planning problem. In complex scenarios, CEP performs better than the baselines. We can obtain an 89% success rate in the first cage environment and 46% in the second cage. We think that these results might be related to the way we compute the robot’s acceleration. While the baselines methods compute the accelerations provided by the different components (Go-To, Avoid Joint Limits, Obstacle Avoidance), and do not take into account any possible conflict of interests between each action, our method first computes the distribution obtained by the product of all the components and then, computes the action that maximizes this composed distribution. The built composition will provide high probabilities to those actions that have high probability in each component and low to the rest. Thus, the optimal action from

| Methods | 1 Obstacle | | 3 Obstacles | | Cross | | Double Cross | | Cage I | | Cage II | |
|---|----------------|--------------|----------------|--------------|---------------|--------------|---------------|--------------|---------------|--------------|---------------|--------------|
| | Success | Collide | Success | Collide | Success | Collide | Success | Collide | Success | Collide | Success | Collide |
| Riemannian Motion Policies [33] | 100/100 | 0/100 | 98/100 | 0/100 | 94/100 | 0/100 | 84/100 | 0/100 | 55/100 | 0/100 | 5/100 | 0/100 |
| Artificial Potential Fields [21] | 100/100 | 0/100 | 100/100 | 0/100 | 90/100 | 0/100 | 76/100 | 0/100 | 43/100 | 0/100 | 4/100 | 2/100 |
| Composable Energy Policies(Ours) | 100/100 | 0/100 | 100/100 | 0/100 | 98/100 | 0/100 | 91/100 | 0/100 | 89/100 | 0/100 | 46/100 | 0/100 |

TABLE I
RESULTS FOR 3D GoTo + OBSTACLE AVOIDANCE TASK

| Methods | 1 Obstacle | | 3 Obstacles | | Cross | | Double Cross | | Cage I | | Cage II | |
|---|----------------|--------------|----------------|--------------|---------------|--------------|---------------|--------------|---------------|--------------|---------------|--------------|
| | Success | Collide | Success | Collide | Success | Collide | Success | Collide | Success | Collide | Success | Collide |
| Riemannian Motion Policies [33] | 100/100 | 0/100 | 91/100 | 0/100 | 75/100 | 0/100 | 63/100 | 0/100 | 18/100 | 0/100 | 1/100 | 0/100 |
| Artificial Potential Fields [21] | 100/100 | 0/100 | 90/100 | 0/100 | 78/100 | 0/100 | 60/100 | 0/100 | 21/100 | 0/100 | 1/100 | 0/100 |
| Composable Energy Policies(Ours) | 100/100 | 0/100 | 100/100 | 0/100 | 91/100 | 0/100 | 87/100 | 0/100 | 43/100 | 0/100 | 14/100 | 0/100 |

TABLE II
RESULTS FOR 6D GoTo + OBSTACLE AVOIDANCE TASK

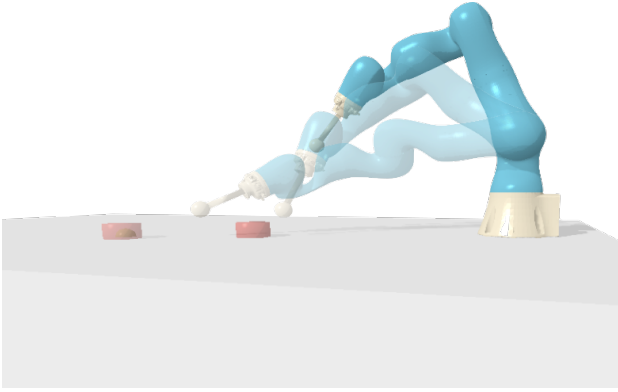


Fig. 5. Puck hitting environment. CEP is tested in a RL environment where the robot needs to learn how to hit the puck. The reward is defined by the euclidean distance between the puck and the green target. An additional reward penalizes if the robot hits the table.

the composed distribution is expected to be the one that better satisfies all the individual components.

The performance worsens for the 6D GoTo problem (Table II). The orientation sets an important constraint in the possible final configurations and reduces the set of trajectories that solves the problem. CEP was able to perform relatively better than the baselines but it got less than 50% success rate in both cage environments, suggesting that in complex scenarios an additional global path planner should be integrated with CEP.

Prior Based Policies for RL

In the following experiment, we evaluate the performance of CEP as a policy that combines prior policies with a new learnable policy to solve a new task. We consider the problem of hitting a puck and place it in a target position Fig. 5. We use a 7 dof LBR-IIWA robot, the action space is defined by the end-effector’s task space cartesian velocity $\mathbf{a} \in \mathbb{R}^3$ and the state $\mathbf{s} \in \mathbb{R}^{18}$ is represented by the end-effector’s cartesian position \mathbf{x}_{ee} , puck’s position \mathbf{x}_{puck} , their relative position $r_{p-ee} = \mathbf{x}_{ee} - \mathbf{x}_{puck}$, the puck’s velocity \mathbf{v}_{puck} , the end effector’s velocity \mathbf{v}_{ee} and the target position \mathbf{x}_{target} . One episode has 1400 steps. We define the reward as

$$r = -\|\mathbf{x}_{puck} - \mathbf{x}_{target}\|^2 + r_T(\mathbf{x}_{ee}) \quad (27)$$

with $r_T(\mathbf{x}_{ee}) = -1000$ if the robot’s end effector position \mathbf{x}_{ee} collides the table. The first term $-\|\mathbf{x}_{puck} - \mathbf{x}_{target}\|^2$ will encourage the robot to put the puck close to the target, minimizing the euclidean distance. The second term will encourage the robot to avoid to collide against the table. We compare CEP approach presented in (25) with other approaches that combines prior information with RL. We consider Behavioural Cloning + RL [29] and Residual Policy Learning [16, 40] as baselines.

To guarantee a fair comparison, we apply the same prior policy in the three scenarios. We build a prior policy with two components in a CEP. We consider

- Go-To energy in the robot’s end effector’s task space.
- Table avoidance in the robot’s end effector’s task space.

The prior is combined in different forms depending on the applied algorithm:

Behavioural Cloning + RL. We first fit a policy to the prior policy

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{p(\mathbf{s})} [D_{KL}(\pi_q(\mathbf{a}|\mathbf{s})||\pi_{\theta}(\mathbf{a}|\mathbf{s}))] \quad (28)$$

and then directly sample from the fit policy

$$\mathbf{a} \sim \pi_{\theta}(\mathbf{a}|\mathbf{s}). \quad (29)$$

Residual Policy Learning. We consider a hierarchical structure. We first sample, from the learnable policy and add the optimal solution from the prior

$$\mathbf{a} = \mathbf{a}^* + \mathbf{a}_H, \mathbf{a}_H \sim \pi_{\theta}(\mathbf{a}_H|\mathbf{s}) \quad (30)$$

with

$$\mathbf{a}^* = \arg \max_{\mathbf{a}} \pi_q(\mathbf{a}|\mathbf{s}). \quad (31)$$

Composable Energy Policies. We also consider a hierarchical structure. In this case, we first sample from the learnable policy

$$\mathbf{a}_H \sim \pi_{\theta}(\mathbf{a}_H|\mathbf{s}) \quad (32)$$

and then, optimize over the composition

$$\mathbf{a} = \arg \max_{\mathbf{a}} \pi_q(\mathbf{a}|\mathbf{s}, \mathbf{a}_H). \quad (33)$$

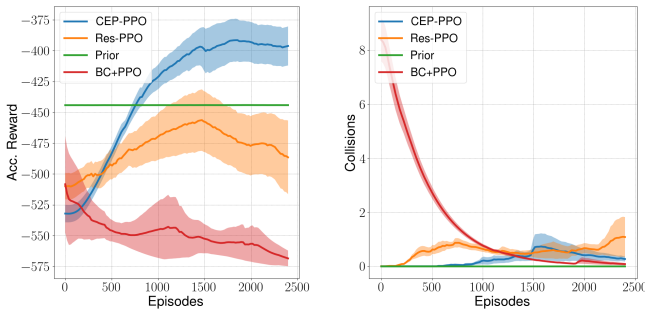


Fig. 6. Training curves for Hitting a puck environment. CEP-PPO performs consistently better than other prior + RL methods. Collision avoidance prior is stronger in CEP-PPO and provides a safer training.

We apply the hierarchical policy introduced in (25). For all the experiments, we used the Mushroom-RL [7] PPO [38] implementation. We include additional details in Appendix D.

a) Results and analysis: We show the obtained results in Fig. 6. Our method, CEP-PPO performs on average better than Residual Learning and Behavioural Cloning + RL methods. The obtained results show that while CEP can consistently improve its performance, residual policy learning has a slower learning curve and behavioral cloning + RL decays in performance. We remark that the three methods start performing worse than the prior due to the stochasticity of the policies in contrast with the deterministic prior. The obtained results might be related to the collisions per episode. CEP based policy allows imposing hard constraints to avoid collisions. Both, behavioral cloning + RL and residual learning adds white noise around the optimal prior action and then, increase the chances of table collision. Fewer collisions mean that $r_T(\mathbf{x}_{ee})$ is less activated during the training and then, the learning can focus on moving the puck close to the target. In the Behavioural Cloning case, we see that the table collisions push the robot to avoid the table, and then, it forgets how to reach the target. The Behavioural Cloning + RL policy ends with a policy that makes the robot neither collide with the table and neither to hit the puck. In the residual learning case, the policy is updated by taking both hitting the puck and avoid the table into consideration, while in CEP, the collisions against the table are less, and thus, the training is more focused on improving the hitting tactic.

CEP ablation study: When applying CEP-PPO, the Σ_H parameter from (25) controls how big is the influence of the learnable policy in the overall policy. The smaller Σ_H is; the smaller the entropy for the learnable part and the higher the influence in the overall policy. In order to study the influence of the Σ_H parameter, we compare CEP-PPO with fixed Σ_H and CEP-PPO with learnable Σ_H . We present the obtained results in Fig. 7 As expected having the freedom to set the Σ_H parameter allows the robot to learn a better behavior to hit the puck. If the Σ_H parameter is fix and too big, then the learned policy has a low impact on the overall policy and it can not improve a lot. With a very small Σ_H , the learned policy will have a very high impact in the overall policy and then, it might not get benefited by the priors. When learning Σ_H

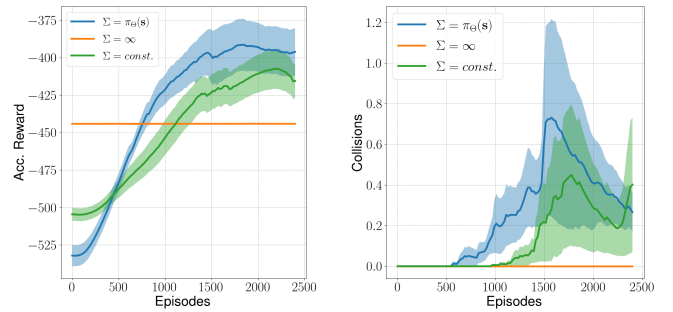


Fig. 7. Training curves for Hitting a puck environment. Learning the variance of the policy improves the learning curve, but also increases the probability of collisions.

parameter in the high-level policy, the robot decides when to trust more in the priors and when less. This allows the robot to get help from priors only if required, but also to decouple from them to learn a better policy. However, shown in the collisions plot from Fig. 7, reducing the effect of the priors might increase the probability to go into dangerous regions and collide more often against the table. Nevertheless, learning the Σ_H parameter seems to provide better results than keeping it fixed.

VII. CONCLUSIONS

We have presented Composable Energy Policies, a novel approach for modular robot control. In contrast with the previous method for reactive motion generation, which computes the optimal action for each component independently; our method solves an optimization problem over the sum of components. We have shown that APF policies can be rewritten as a CEP and provide reasoning of why their method might fail. Through an evaluation phase, we have also shown that our method can solve a set of reactive motion generation problems more efficiently than artificial potential field methods.

Our method is flexible and allows us to represent each component in an arbitrary set of state-action spaces. Additionally, CEP allows integrating policies from multiple sources and samples from the product of them. In our work, we show how to integrate a RL agent with a set of priors. The resulting algorithm combines a high-level RL policy with a low-level optimization problem. This low-level optimization allows to guarantee that the policy is safe (if an obstacle avoidance prior is added) or to explore in informative regions (if guiding priors are added). We have compared our method with two methods that also combine prior information with an RL agent. Given the low-level optimization of our method, we can guarantee that some priors such as obstacle avoidance were satisfied and the learning is more stable.

Additionally, we have introduced in the Appendix, the relations between optimal control and reactive motion generation. Setting their connections allows us to compute how the optimal policy for a multi-objective optimal control problem diverges from the product of the optimal policies for each objective. These insights allow us to better understand the policy composition and improve the design of our components.

ACKNOWLEDGMENTS

This project has received funding from the European Union’s Horizon 2020 research and innovation programmes under grant agreement No. #820807 (SHAREWORK).

REFERENCES

- [1] Abbas Abdolmaleki, Sandy Huang, Leonard Hasenclever, Michael Neunert, Francis Song, Martina Zambelli, Murilo Martins, Nicolas Heess, Raia Hadsell, and Martin Riedmiller. A distributional view on multi-objective policy optimization. In *International Conference on Machine Learning*, pages 11–22. PMLR, 2020.
- [2] Adi Ben-Israel. The change-of-variables formula using matrix volume. *SIAM Journal on Matrix Analysis and Applications*, 21(1):300–312, 1999.
- [3] Adi Ben-Israel. An application of the matrix volume in probability. *Linear Algebra and its Applications*, 321(1-3):9–25, 2000.
- [4] Ching-An Cheng, Mustafa Mukadam, Jan Issac, Stan Birchfield, Dieter Fox, Byron Boots, and Nathan Ratliff. Rmpflow: A computational graph for automatic motion policy generation. In *International Workshop on the Algorithmic Foundations of Robotics*, 2018.
- [5] Marco Da Silva, Frédo Durand, and Jovan Popović. Linear bellman combination for control of character animation. *ACM SIGGRAPH*, 2009.
- [6] Christian Daniel, Gerhard Neumann, and Jan Peters. Hierarchical relative entropy policy search. In *Artificial Intelligence and Statistics*, pages 273–281. PMLR, 2012.
- [7] Carlo D’Eramo, Davide Tateo, Andrea Bonarini, Marcello Restelli, and Jan Peters. Mushroomrl: Simplifying reinforcement learning research. <https://github.com/MushroomRL/mushroom-rl>, 2020.
- [8] Yilun Du, Shuang Li, and Igor Mordatch. Compositional visual generation and inference with energy based models. In *Conference on Neural Information Processing Systems*, 2020.
- [9] Carlos Florensa, Yan Duan, and Pieter Abbeel. Stochastic neural networks for hierarchical reinforcement learning. In *International Conference on Learning Representations*, 2017.
- [10] Dieter Fox, Wolfram Burgard, and Sebastian Thrun. The dynamic window approach to collision avoidance. *IEEE Robotics & Automation Magazine*, 4(1):23–33, 1997.
- [11] Shuzhi Sam Ge and Yun J Cui. Dynamic motion planning for mobile robots using potential field method. *Autonomous robots*, 13(3):207–222, 2002.
- [12] Tuomas Haarnoja, Vitchyr Pong, Aurick Zhou, Murtaza Dalal, Pieter Abbeel, and Sergey Levine. Composable deep reinforcement learning for robotic manipulation. In *IEEE International Conference on Robotics and Automation*, pages 6244–6251. IEEE, 2018.
- [13] Nicolas Heess, Greg Wayne, Yuval Tassa, Timothy Lillicrap, Martin Riedmiller, and David Silver. Learning and transfer of modulated locomotor controllers. *arXiv preprint arXiv:1610.05182*, 2016.
- [14] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002.
- [15] Andrew Howard, Maja J Matarić, and Gaurav S Sukhatme. Mobile sensor network deployment using potential fields: A distributed, scalable solution to the area coverage problem. In *Distributed Autonomous Robotic Systems 5*, pages 299–308. Springer, 2002.
- [16] Tobias Johannink, Shikhar Bahl, Ashvin Nair, Jianlan Luo, Avinash Kumar, Matthias Loskyll, Juan Aparicio Ojea, Eugen Solowjow, and Sergey Levine. Residual reinforcement learning for robot control. In *International Conference on Robotics and Automation*, pages 6023–6029. IEEE, 2019.
- [17] Leslie Pack Kaelbling and Tomás Lozano-Pérez. Hierarchical task and motion planning in the now. In *IEEE International Conference on Robotics and Automation*, pages 1470–1477. IEEE, 2011.
- [18] Leslie Pack Kaelbling and Tomás Lozano-Pérez. Integrated task and motion planning in belief space. *The International Journal of Robotics Research*, 32(9-10):1194–1227, 2013.
- [19] Mrinal Kalakrishnan, Sachin Chitta, Evangelos Theodorou, Peter Pastor, and Stefan Schaal. Stomp: Stochastic trajectory optimization for motion planning. In *IEEE international conference on robotics and automation*, pages 4569–4574. IEEE, 2011.
- [20] Lydia E Kavraki, Petr Svestka, J-C Latombe, and Mark H Overmars. Probabilistic roadmaps for path planning in high-dimensional configuration spaces. *IEEE transactions on Robotics and Automation*, 12(4):566–580, 1996.
- [21] O. Khatib. Real-time obstacle avoidance for manipulators and mobile robots. In *IEEE International Conference on Robotics and Automation*, volume 2, pages 500–505, 1985. doi: 10.1109/ROBOT.1985.1087247.
- [22] Oussama Khatib. A unified approach for motion and force control of robot manipulators: The operational space formulation. *IEEE Journal on Robotics and Automation*, 1987.
- [23] Yoram Koren and Johann Borenstein. Potential field methods and their inherent limitations for mobile robot navigation. In *IEEE International Conference on Robotics and Automation*, 1991.
- [24] Alexander Lambert, Adam Fishman, Dieter Fox, Byron Boots, and Fabio Ramos. Stein variational model predictive control. *Conference on Robot Learning*, 2020.
- [25] Steven M LaValle. *Planning Algorithms*. Cambridge University Press, Cambridge, U.K., 2006.
- [26] Steven M LaValle, James J Kuffner, BR Donald, et al. Rapidly-exploring random trees: Progress and prospects. *Algorithmic and computational robotics: new directions*, 5:293–308, 2001.
- [27] Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, and F Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.
- [28] Sergey Levine. Reinforcement learning and control

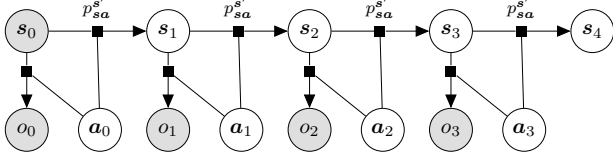
- as probabilistic inference: Tutorial and review. *CoRR*, abs/1805.00909, 2018.
- [29] Katharina Mülling, Jens Kober, Oliver Kroemer, and Jan Peters. Learning to select and generalize striking movements in robot table tennis. *The International Journal of Robotics Research*, 32(3):263–279, 2013.
- [30] Alexandros Paraschos, Christian Daniel, Jan R Peters, and Gerhard Neumann. Probabilistic movement primitives. In *Advances in Neural Information Processing Systems*, pages 2616–2624, 2013.
- [31] Xue Bin Peng, Michael Chang, Grace Zhang, Pieter Abbeel, and Sergey Levine. Mcp: Learning composable hierarchical control with multiplicative compositional policies. *Advances in Neural Information Processing Systems* 32, 2019.
- [32] Nathan Ratliff, Matt Zucker, J Andrew Bagnell, and Siddhartha Srinivasa. Chomp: Gradient optimization techniques for efficient motion planning. In *IEEE International Conference on Robotics and Automation*, pages 489–494, 2009.
- [33] Nathan D Ratliff, Jan Issac, Daniel Kappler, Stan Birchfield, and Dieter Fox. Riemannian motion policies. *arXiv preprint arXiv:1801.02854*, 2018.
- [34] Konrad Rawlik, Marc Toussaint, and Sethu Vijayakumar. On stochastic optimal control and reinforcement learning by approximate inference. *Proceedings of Robotics: Science and Systems VIII*, 2012.
- [35] Elon Rimon and Daniel E Koditschek. Exact robot navigation using artificial potential functions. *Departmental Papers (ESE)*, page 323, 1992.
- [36] Reuven Rubinfeld. The cross-entropy method for combinatorial and continuous optimization. *Methodology And Computing In Applied Probability*, 1999.
- [37] John Schulman, Yan Duan, Jonathan Ho, Alex Lee, Ibrahim Awwal, Henry Bradlow, Jia Pan, Sachin Patil, Ken Goldberg, and Pieter Abbeel. Motion planning with sequential convex optimization and convex collision checking. *The International Journal of Robotics Research*, 33(9):1251–1270, 2014.
- [38] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [39] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- [40] Tom Silver, Kelsey Allen, Josh Tenenbaum, and Leslie Kaelbling. Residual policy learning. *arXiv preprint arXiv:1812.06298*, 2018.
- [41] Avi Singh, Huihan Liu, Gaoyue Zhou, Albert Yu, Nicholas Rhinehart, and Sergey Levine. Parrot: Data-driven behavioral priors for reinforcement learning. *International Conference on Learning Representations*, 2021.
- [42] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- [43] Geraud Nangue Tasse, Steven James, and Benjamin Rosman. A boolean task algebra for reinforcement learning. *Conference on Neural Information Processing Systems*, 2020.
- [44] Emanuel Todorov. Compositionality of optimal control laws. *Advances in neural information processing systems*, 22:1856–1864, 2009.
- [45] Marc Toussaint. Robot trajectory optimization using approximate inference. In *international conference on machine learning*, pages 1049–1056, 2009.
- [46] Julen Urain, Michele Ginesi, Davide Tateo, and Jan Peters. Imitationflows: Learning deep stable stochastic dynamic systems by normalizing flows. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2020.
- [47] Jur Van Den Berg, Stephen J Guy, Ming Lin, and Dinesh Manocha. Reciprocal n-body collision avoidance. In *Robotics research*, pages 3–19. Springer, 2011.
- [48] Benjamin Van Niekerk, Steven James, Adam Earle, and Benjamin Rosman. Composing value functions in reinforcement learning. In *International Conference on Machine Learning*, pages 6401–6409, 2019.
- [49] Joe Watson, Hany Abdulsamad, and Jan Peters. Stochastic optimal control as approximate input inference. In *Conference on Robot Learning*, pages 697–716. PMLR, 2020.
- [50] Grady Williams, Andrew Aldrich, and Evangelos A Theodorou. Model predictive path integral control: From theory to parallel computation. *Journal of Guidance, Control, and Dynamics*, 40(2):344–357, 2017.
- [51] Brian D Ziebart, J Andrew Bagnell, and Anind K Dey. Modeling interaction via the principle of maximum causal entropy. In *International Conference on Machine Learning*, 2010.

APPENDIX A

A CONTROL AS INFERENCE VIEW FOR REACTIVE MOTION GENERATION

In the following section, we want to highlight the connections between reactive motion generation and control as inference to evaluate the optimality guarantees of composing energies in a multi-objective optimal control problem.

Fig. 8. Graphical model for the Optimal Control problem. s_t denoted the state, \mathbf{a}_t denotes the action and o_t is an additional variable representing the optimality of the state and action for a given reward.



We frame optimal control as Bayesian inference problem [34, 28] over the sequence of actions $\mathbf{a}_{0:T}$. The optimal control problem is visualized as a graphical model in Fig. 8. The problem is formulated introducing an auxiliary variable $o_{0:T}$ that represents the optimality of s_t and \mathbf{a}_t under a certain reward function, $p(o_t | s_t, \mathbf{a}_t) \propto \exp(r(s_t, \mathbf{a}_t))$. Given a certain prior distribution $q(\mathbf{a} | s_0)$ and given s_0 is known, the inference problem is

$$p(A_0^T | s_0, O_0^T) = \frac{p(O_0^T | A_0^T, s_0) q(A_0^T | s_0)}{p(O_0^T | s_0)} \quad (34)$$

with

$$p(O_0^T | A_0^T, s_0) = \int_{\mathbf{s}_{1:T}} p(O_0^T | S_0^T, A_0^T) p(S_1^T | A_0^T, s_0) dS_1^T \quad (35)$$

where $A_0^T : \{\mathbf{a}_0, \dots, \mathbf{a}_T\}$, $O_0^T : \{o_0, \dots, o_T\}$ and $S_0^T : \{s_0, \dots, s_T\}$. There are two main directions to solve the posterior in (34). First, methods that frame the problem as an Hidden Markov Model (HMM) and solve it in an Expectation-Maximization approach. Second, methods that compute the posterior in the trajectory level A_0^T . The first, are computationally demanding as they require several forward and backward message passing to compute the posterior. The second, needs to solve the problem in the trajectory level and thus the dimension of the variables grows linearly with the trajectory length, T .

Reactive motion generation instead, solves a one-step-ahead problem. The solutions in reactive motion generation are local rather than global and thus, they are computationally more efficient.

We frame reactive motion generation as a one-step control as inference problem

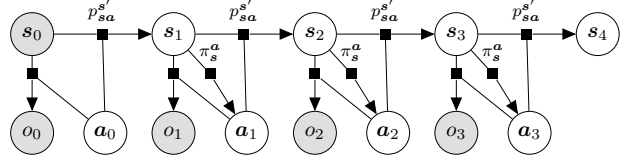
$$p(\mathbf{a}_0 | s_0, O_0^T) = \frac{p(O_0^T | \mathbf{a}_0, s_0) q(\mathbf{a}_0 | s_0)}{p(O_0^T | s_0)} \quad (36)$$

with

$$p(O_0^T | \mathbf{a}_0, s_0) = \int_{S_{1:T}} \int_{A_{1:T}} p(O_0^T | S_0^T, A_0^T) p(S_1^T | \mathbf{a}_0, s_0) \pi(A_1^T | S_1^T) dS_1^T dA_1^T. \quad (37)$$

In contrast with the classical stochastic optimal control problem that finds the posterior for the whole trajectory A_0^T , in one-step-ahead control as inference problem, we aim to find the posterior for only the instant next control action, \mathbf{a}_0 . Computing the posterior only for \mathbf{a}_0 requires the likelihood to be defined as the marginal of not only the state trajectory S_1^T , but also the action trajectory A_1^T . The graphical model for one-step-ahead control as inference is presented in Fig. 9. In

Fig. 9. Graphical model for one-step ahead optimal control problem. In this approach, the actions A_1^T are dependant on S_1^T given a policy π .



one-step ahead control as inference, we introduce an additional distribution π that provides us the probability of A_1^T given S_1^T . This additional policy π is interpreted as the policy the agent will apply in the future. In this context, given we know the policy the agent will apply in the future, we are looking for the instant action that maximized the cumulative reward in the long horizon trajectory. (37) can be rewritten as the expectation over $p(O_0^T | S_0^T, A_0^T)$

$$\begin{aligned} \mathbb{E}_{S_0^T, A_0^T \sim p^\pi(S_0^T, A_0^T | s_0, \mathbf{a}_0)} [p(O_0^T | S_0^T, A_0^T)] &\propto \\ \mathbb{E}_{S_0^T, A_0^T \sim p^\pi(S_0^T, A_0^T | s_0, \mathbf{a}_0)} \left[\exp\left(\sum_{t=0}^T r(s_t, \mathbf{a}_t)\right) \right] &= \\ \exp(Q_r^\pi(s_0, \mathbf{a}_0)) &\quad (38) \end{aligned}$$

From what follows, the likelihood for our inference problem is proportional to the $\exp(Q)$ defined over a certain policy π . Given a π , the posterior for our inference problem is given by

$$p(\mathbf{a}_0 | s_0, O_0^T) \propto \exp(Q_r^\pi(s_0, \mathbf{a}_0)) q(\mathbf{a}_0 | s_0). \quad (39)$$

The provided Q function depends on π and then, the quality of our reactive motion generator to solve a long horizon optimal control problem directly depends on the quality of π . In the optimal case, for Q^* , the reactive motion generator can find the optimal trajectory, even if the optimization is done locally.

A. Optimality Guarantees

In CEP, we propose to model the Q function as the sum of a set of optimal Q_k^* functions. Instead, we are aware that $Q_\Sigma = \frac{1}{K} \sum_{k=0}^K Q_k^*$ is not the optimal function Q^* for the sum of the rewards $r = \frac{1}{K} \sum_{k=1}^K r_k$. The closer Q_Σ is from the

optimal Q^* , the closer the product of experts policy would be from the optimal policy. In this section, we study, given a certain reward $r = \frac{1}{2}(r_1 + r_2)$, how much the sum of the individual components Q_Σ diverge from the optimal Q^* . In [28] is shown, that the optimal Q function for the control as inference problem can be computed by recursively solving the soft-value iteration [51]

$$\begin{aligned} Q(\mathbf{s}, \mathbf{a}) &= r(\mathbf{s}, \mathbf{a}) + \mathbb{E}_{p(\mathbf{s}'|\mathbf{s}, \mathbf{s})} [V(\mathbf{s}')] \\ V(\mathbf{s}) &= \log \int_{\mathcal{A}} \exp(Q(\mathbf{a}, \mathbf{s})). \end{aligned} \quad (40)$$

We assume a finite horizon control as inference problem and evaluate how much Q_Σ diverges from Q^* when increasing the control horizon

For $t = T$,

$$Q^{*T} = \frac{1}{2}(r_1 + r_2) \quad (41)$$

and

$$Q_1^{*T} = r_1, Q_2^{*T} = r_2. \quad (42)$$

Then,

$$Q^{*T} = \frac{1}{2}(Q_1^{*T} + Q_2^{*T}) \quad (43)$$

and the divergence

$$\Delta Q^T = Q^{*T} - Q_\Sigma^T = 0. \quad (44)$$

From (43), for $T = 1$ horizon optimal control problems, the sum of the optimal components Q_Σ is equal to the optimal Q . For longer horizon control as inference problems, the optimal Q is computed by recursively solving the soft-bellman update backward in time. For computing $t = T - 1$,

$$\begin{aligned} Q^{*T-1}(\mathbf{s}, \mathbf{a}) &= r(\mathbf{s}, \mathbf{a}) + \mathbb{E}_{p(\mathbf{s}'|\mathbf{s}, \mathbf{s})} [V^{*T}(\mathbf{s}')] \\ V^{*T}(\mathbf{s}) &= \log \int_{\mathcal{A}} \exp(Q^{*T}(\mathbf{a}, \mathbf{s})). \end{aligned} \quad (45)$$

we do a soft bellman update. The difference between Q^{*T-1} and Q_Σ^{T-1}

$$\Delta Q^{T-1} = Q^{*T-1} - Q_\Sigma^{T-1} = \mathbb{E}_{p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [V^{*T} - V_\Sigma^T] \quad (46)$$

with

$$\begin{aligned} V^{*T} - V_\Sigma^T &= \log \frac{\int_{\mathcal{A}} \exp(Q^{*T})}{\int_{\mathcal{A}} (\exp(Q_1^{*T}) \int_{\mathcal{A}} \exp(Q_2^{*T}))^{\frac{1}{2}}} \\ &= \log \frac{\int_{\mathcal{A}} \exp(\frac{1}{2}Q_1^{*T}) \exp(\frac{1}{2}Q_2^{*T})}{(\int_{\mathcal{A}} \exp(Q_1^{*T}) \int_{\mathcal{A}} \exp(Q_2^{*T}))^{\frac{1}{2}}}. \end{aligned} \quad (47)$$

Then,

$$\begin{aligned} Q^{*T-1} &= Q_\Sigma^{T-1} + \Delta Q^{T-1} \\ &= Q_\Sigma^{T-1} + \mathbb{E}_{p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[\log \frac{\int_{\mathcal{A}} \exp(\frac{1}{2}Q_1^{*T}) \exp(\frac{1}{2}Q_2^{*T})}{(\int_{\mathcal{A}} \exp(Q_1^{*T}) \int_{\mathcal{A}} \exp(Q_2^{*T}))^{\frac{1}{2}}} \right] \end{aligned} \quad (48)$$

From (46) and (47), we can obtain the recurrence relation for the divergence error

$$\Delta Q^{t-1} = \mathbb{E}_{p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[\log \frac{\int_{\mathcal{A}} \exp(\frac{1}{2}Q_1^{*t}) \exp(\frac{1}{2}Q_2^{*t}) \exp(\Delta Q^t)}{(\int_{\mathcal{A}} \exp(Q_1^{*t}) \int_{\mathcal{A}} \exp(Q_2^{*t}))^{\frac{1}{2}}} \right]. \quad (49)$$

From (46) and (47), we can see that if $Q_1^* = Q_2^*$, then $\Delta Q = 0$ and the more they differ, the bigger the divergence error to the optima Q^* .

We can also extrapolate interesting results from (49). For any Q_1^* and Q_2^* ,

$$\left(\int_{\mathcal{A}} \exp(Q_1^{*t}) \int_{\mathcal{A}} \exp(Q_2^{*t}) \right)^{\frac{1}{2}} > \int_{\mathcal{A}} \exp(\frac{1}{2}Q_1^{*t}) \exp(\frac{1}{2}Q_2^{*t}) \quad (50)$$

and then,

$$\Delta Q < 0, \forall Q_1^*, Q_2^* \quad (51)$$

And then, from (49), ΔQ monotonically decreases. The longer the horizon T, the bigger the divergence error ΔQ .

In practice, these insights are beneficial for modeling our Q functions. We can expect the performance of our reactive motion generation to decay in those environments in which long-horizon planning is required as the ΔQ will increase over long horizons. On the other hand, the more overlapping regions the Q components have, the smaller the divergence error would be and the closer we would be to the optima.

Similar theoretical studies have been already developed [12, 48, 44, 1]. In Optimal Control, composable optimality guarantees were proven for linear dynamics, by the linearly-solvable Markov Decision Processes (LMDP) approach. In [44, 5], a weighted policy sum was proven to be optimal for the sum of the rewards, as long as the rewards differ only in the terminal reward.

In [12, 43], the optimality of the composition is studied in a maximum entropy reinforcement learning problem

$$J(\pi) = \mathbb{E}_{\pi, p(\mathbf{s}_0)} \left[\sum_{t=0}^T r(\mathbf{s}_t, \mathbf{a}_t) + \alpha \mathcal{H}(\pi) | \mathbf{s}_{t+1} \sim p(\cdot | \mathbf{s}_t, \mathbf{a}_t) \right] \quad (52)$$

with $\mathcal{H}(\pi)$ the entropy of the policy. In [12] is proven, that optimal soft-Q function for $r = \frac{1}{2}r_1 + r_2$ is bounded

$$Q_\Sigma(\mathbf{s}, \mathbf{a}) \geq Q^*(\mathbf{s}, \mathbf{a}) \geq Q_\Sigma(\mathbf{s}, \mathbf{a}) - \mathcal{C}^*(\mathbf{s}, \mathbf{a}) \quad (53)$$

with $Q_\Sigma = \frac{1}{2}Q_1^* + Q_2^*$ and \mathcal{C}^* is the fixed point of

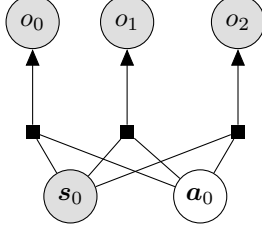
$$\mathcal{C}(\mathbf{s}, \mathbf{a}) \leftarrow \gamma \mathbb{E}_{p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[\mathcal{D}_{\frac{1}{2}}(\pi_1^*(\mathbf{a}|\mathbf{s}) || \pi_2^*(\mathbf{a}|\mathbf{s})) + \max_{\mathbf{a}'} \mathcal{C}(\mathbf{s}', \mathbf{a}') \right] \quad (54)$$

and $\mathcal{D}_{\frac{1}{2}}$ is the Renyi divergence of order $\frac{1}{2}$.

APPENDIX B
COMPOSABLE ENERGY POLICIES AS INFERENCE

The product of experts is the natural expression for the multi-objective inference problem presented in Fig. 10. We introduce a set of auxiliary variables o_0, \dots, o_K that represents the optimality of \mathbf{s}_0 and \mathbf{a}_0 given a certain energy function $p(o_k = 1 | \mathbf{s}_0, \mathbf{a}_0) \propto \exp(E_k(\mathbf{s}, \mathbf{a}))$. Then, the likelihood for

Fig. 10. Graphical model for Composable Energy Policies. o_k is an auxiliary variable that represents the optimality of \mathbf{s}_0 and \mathbf{a}_0 for a particular energy.



the graphical model in Fig. 10 can be computed as the product of the terms

$$p(\mathbf{s}_0, \mathbf{a}_0, o_{0:2}) = q(\mathbf{a}_0)p(\mathbf{s}_0) \prod_{k=0}^2 p(o_k | \mathbf{s}_0, \mathbf{a}_0). \quad (55)$$

We are interested on computing the maximum a posteriori for \mathbf{a}_0 given \mathbf{s}_0 and $o_k = 1, \forall k = (0, 1, 2)$

$$\begin{aligned} \mathbf{a}_0^* &= \arg \max_{\mathbf{a}_0} p(\mathbf{a}_0 | \mathbf{s}_0, o_{0:2} = 1) \\ &\propto \exp \left(\sum_{k=0}^2 E_k(\mathbf{a}_0, \mathbf{s}_0) \right) q(\mathbf{a}_0). \end{aligned} \quad (56)$$

From (56), we can see that the maximum a posteriori will search for the optimal \mathbf{a}_0 given an exponentiated sum of energies and the prior distribution $q(\mathbf{a}_0)$. If we assume a uniform distribution for the prior, we recover (4). Thus, the product of experts represents the distribution that satisfies best a set of components E_0, \dots, E_k . We highlight that the CEP graphical model in Fig. 10 and in Fig. 3 are equivalent. The model in Fig. 10, represent the CEP as the distribution maximizing a set of rewards E_0, \dots, E_k , while the model in Fig. 3 represent the CEP as the distribution maximizing the combination of a set of policies. Depending on the problem, we might be interested to frame the problem with one or the other approach.

APPENDIX C
COMPOSABLE ENERGY POLICIES AND ARTIFICIAL
POTENTIAL FIELDS METHOD

In the following section, we will show that the Artificial Potential Fields Method can be written as a particular case of Composable Energy Policies.

Artificial Potential Field Method assumes a weighted sum of the accelerations generated by a set of dynamic components

$$\ddot{\mathbf{q}} = \sum_{k=0}^K \Lambda_k \mathbf{g}_k(\mathbf{q}, \dot{\mathbf{q}}). \quad (57)$$

Let's consider a particular case of CEP were all the components are represented with a normal distribution centered in $\mathbf{g}_k(\mathbf{q}, \dot{\mathbf{q}})$,

$$\pi_k(\ddot{\mathbf{q}} | \mathbf{q}, \dot{\mathbf{q}}) = \mathcal{N}(\mathbf{g}_k(\mathbf{q}, \dot{\mathbf{q}}), \Sigma_k). \quad (58)$$

Then, given K components

$$\pi(\ddot{\mathbf{q}} | \mathbf{q}, \dot{\mathbf{q}}) = \prod_{k=0}^K \mathcal{N}(\mathbf{g}_k(\mathbf{q}, \dot{\mathbf{q}}), \Sigma_k) = \mathcal{N}(\mathbf{g}_\Sigma(\mathbf{q}, \dot{\mathbf{q}}), \Sigma_\Sigma) \quad (59)$$

with

$$\mathbf{g}_\Sigma = \left(\sum_{k=0}^K \Sigma_k^{-1} \right)^{-1} \left(\sum_{k=0}^K \Sigma_k^{-1} \mathbf{g}_k \right) \quad (60)$$

and

$$\Sigma_\Sigma = \left(\sum_{k=0}^K \Sigma_k^{-1} \right)^{-1}. \quad (61)$$

In CEP, the acceleration is obtained by an optimization of

$$\ddot{\mathbf{q}} = \arg \max_{\ddot{\mathbf{q}}} \pi(\ddot{\mathbf{q}} | \mathbf{q}, \dot{\mathbf{q}}). \quad (62)$$

The maximum of 59 is known; given the policy is a gaussian, the maximum is the mean in 60

$$\ddot{\mathbf{q}} = \left(\sum_{k=0}^K \Sigma_k^{-1} \right)^{-1} \left(\sum_{k=0}^K \Sigma_k^{-1} \mathbf{g}_k(\mathbf{q}, \dot{\mathbf{q}}) \right). \quad (63)$$

We can rewrite (63) to (57)

$$\Lambda_k = \left(\sum_{k=0}^K \Sigma_k^{-1} \right)^{-1} \cdot \Sigma_k^{-1} \quad (64)$$

We show that APF policy can be rewritten as a CEP, modelling the dynamics by a product of normal distributions. Modelling all the components by normal distributions might be a bad choice for a proper modules integration. normal distributions assume that (1) there is a unique optimal action (the mean) to solve the task and (2) the actions quality is related to the mahalanobis distance to the optima. While tasks like reaching a target might satisfy (1) and (2); tasks like obstacle avoidance might require richer representations to properly solve the task.

APPENDIX D
EXPERIMENTS

A. *Reactive Motion Generation*

The modular components (Go-to, Obstacle Avoidance and joint limits avoidance) for both baselines APF and RMP were modeled based on [21] and [4] respectively. CEP were modeled following the energies presented in section IV. We introduce a table with the hyperparameters for all CEP energies

B. *Prior based policies for RL*

We introduce a Table with the Hyperparameters for the PPO in Mushroom-RL [7]. We used the same hyperparameters for the three methods (Behavioural Cloning + RL, Residual RL and CEP for RL).

| Energy Modules | Parameters | | |
|------------------------|----------------|---------------|----------------|
| Reach Target | $K_p = 20.$ | $K_v = 30.$ | $\alpha = 10.$ |
| Obstacle Avoidance | $\gamma = 0.2$ | $\alpha = 4.$ | $\beta = 0.1$ |
| Joint Limits avoidance | $\gamma = 0.3$ | $\alpha = 4.$ | $\beta = 0.1$ |

TABLE III

COMPONENT PARAMETERS FOR COMPOSABLE ENERGY POLICIES IN REACTIVE MOTION GENERATION

| Hyperparameters | |
|-------------------|--------------|
| horizon | 1400 |
| Policy Net | 18-128-128-3 |
| policy std0 | 1. |
| policy batch | 256 |
| n_epochs_policy | 2 |
| policy learn rate | 1e-4 |
| Critic Net | 18-256-256-1 |
| critic batch | 256 |
| critic learn rate | 3e-4 |
| n_steps_per_fit | horizon x 30 |
| Discount factor | .997 |

TABLE IV

COMPONENT PARAMETERS FOR COMPOSABLE ENERGY POLICIES IN REACTIVE MOTION GENERATION