
This is an electronic reprint of the original article.
This reprint may differ from the original in pagination and typographic detail.

Hämäläinen, Perttu; Toikka, Juuso; Babadi, Amin; Liu, C Karen

Visualizing Movement Control Optimization Landscapes

Published in:
IEEE Transactions on Visualization and Computer Graphics

DOI:
[10.1109/TVCG.2020.3018187](https://doi.org/10.1109/TVCG.2020.3018187)

Published: 01/01/2022

Document Version
Publisher's PDF, also known as Version of record

Published under the following license:
CC BY

Please cite the original version:
Hämäläinen, P., Toikka, J., Babadi, A., & Liu, C. K. (2022). Visualizing Movement Control Optimization Landscapes. *IEEE Transactions on Visualization and Computer Graphics*, 28(3), 1648-1660.
<https://doi.org/10.1109/TVCG.2020.3018187>

Visualizing Movement Control Optimization Landscapes

Perttu Hämäläinen, Juuso Toikka, Amin Babadi, and C. Karen Liu

Abstract—

A large body of animation research focuses on optimization of movement control, either as action sequences or policy parameters. However, as closed-form expressions of the objective functions are often not available, our understanding of the optimization problems is limited. Building on recent work on analyzing neural network training, we contribute novel visualizations of high-dimensional control optimization landscapes; this yields insights into why control optimization is hard and why common practices like early termination and spline-based action parameterizations make optimization easier. For example, our experiments show how trajectory optimization can become increasingly ill-conditioned with longer trajectories, but parameterizing control as partial target states—e.g., target angles converted to torques using a PD-controller—can act as an efficient preconditioner. Both our visualizations and quantitative empirical data also indicate that neural network policy optimization scales better than trajectory optimization for long planning horizons. Our work advances the understanding of movement optimization and our visualizations should also provide value in educational use.

Index Terms—Visualization, Animation, Movement synthesis, Trajectory optimization, Policy optimization, Control optimization

1 Introduction

MUCH of computer animation research formulates animation as an optimization problem. It has been shown that complex movements can emerge from minimizing a cost function that measures the divergence from movement goals, such as moving to a specific pose or location while minimizing effort. In principle, this holds the promise of elevating an animator to the role of a choreographer, directing virtual actors and stuntmen through the definition of movement goals. However, solving the optimization problems can be hard in practice. It can require hours or even days of computing time, which is highly undesirable for interactive applications and the rapid iteration of movement goals; defining the goals can be a non-trivial design problem in itself, requiring multiple attempts to produce a desired aesthetic result.

Optimization problems can be divided into the four classes of increasing difficulty illustrated in Fig. 1:

- *Convex and well-conditioned* Convexity refers to the shape of the isocontours of the objective function, or level sets in a general d -dimensional case. In the ideal well-conditioned case, the isocontours are spherical, and simple gradient descent recovers a direct path to the optimum.
- *Convex and ill-conditioned* In ill-conditioned optimization, the isocontours are elongated instead of spherical. The gradient—coinciding with isocontour normals—no longer points towards the optimum, and numerical optimization may require more iterations.
- *Non-convex and unimodal* Non-convexity tends to make optimization even harder, but numerical opti-

mization usually still converges if there are no local optima to distract it.

- *Non-convex and multimodal* In this problem class, the landscape has local optima which can attract optimization. Gradient-free, sampling-based approaches like CMA-ES—common in animation research—may still find the global optimum [1], but this can be computationally expensive. Unfortunately, movement optimization can easily fall into this class, e.g. due to the discontinuities caused by colliding objects, and multiple options for going around obstacles [2], [3].

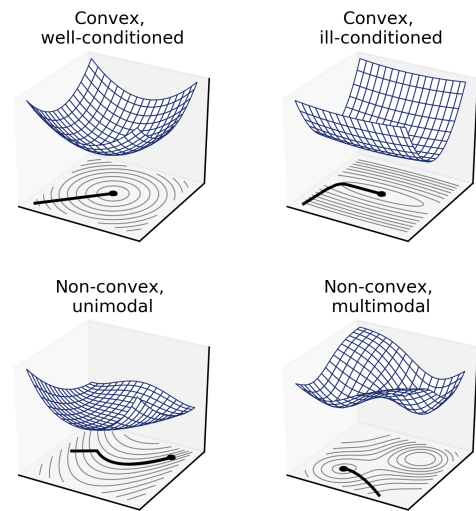


Fig. 1: Common types of optimization landscapes. The surfaces denote the values of 2-dimensional (bivariate) objective functions, with the isocontours displayed below the surface. The black curves show the progress of gradient descent optimization from an initial point.

- Hämäläinen, Toikka, and Babadi are with the Department of Computer Science at Aalto University.
E-mail: perttu.hamalainen@aalto.fi, juuso.toikka@aalto.fi, amin.babadi@aalto.fi
- C. Karen Liu is with the Department of Computer Science at Stanford University.
E-mail: karenliu@cs.stanford.edu

Landscape visualizations like those in Fig. 1 provide useful intuitions of optimization problems, and can help in reformulating a problem into a more tractable form. In general, one would like the objective function to be *more convex, well-conditioned, and unimodal*. We know that problem modifications such as the choice of action space can greatly affect movement optimization efficiency [4], but visualizing the effects on the optimization landscape is challenging because of high problem dimensionality. High dimensionality and/or non-differentiable physics simulators can also prevent analyzing problem conditioning through the eigenvalues of the Hessian.

The primary inspiration of this paper comes from recent work on visualizing neural network loss function landscapes by Li et al. [5]. Strikingly, the paper shows that *visualization of random 2D slices of a high-dimensional objective function can convey useful intuitions and predict the difficulty of optimization*, even with highly complex networks with millions of parameters. More specifically, the approach generates 3D landscape plots of the objective function $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$ by evaluating it along a plane (a 2D subspace in \mathbb{R}^d) defined by two random orthogonal basis vectors intersecting the optimum. Li et al. [5] use the approach to illustrate how deeper networks have more local optima, but adding skip-connections greatly helps in making the landscape more convex and unimodal.

Contribution: We contribute by showing that the random slice visualization approach of Li et al. [5] can be applied in the domain of movement optimization. Fig. 2 shows examples of this, visualizing the same objective function with different random basis vectors. Furthermore, we use the visualizations to investigate the following research questions:

- What is the effect of the number of timesteps – i.e., the planning horizon – on the optimization landscape? (Section 4)
- What is the effect of the choice of action space on the optimization landscape? (Section 5)
- What is the effect of converting instantaneous costs into rewards through exponentiation? (Section 6)
- What is the effect of early termination of movement trajectories or episodes, e.g., when deviating from a target state? (Section 7)
- Do the visualizations predict actual optimization performance, and generalize from simple to complex problems? (Sections 8 and 9).

Additionally, Section 10 provides a more theoretical investigation of the reliability and limitations of random 2D slice visualizations. We conclude that such visualization is a useful tool for diagnosing problems; the somewhat low sensitivity is compensated by high specificity. Our visualizations also explain why movement optimization best practices such as early termination and parameterizing actions as target angles work so well, which should make our work useful in teaching computer animation and movement optimization.

2 Related work

Spacetime optimization Much of the earlier work on animation as optimization focused on extensions of the seminal work on spacetime optimization by Witkin and Kass [6], [7],

[8], [9], [10], where the optimized variables included the root position and rotation as well as joint rotations for each animation frame. However, the synthesized motions were limited by the need for prior knowledge of contact information, such as when and which body parts should make contact with the ground. This limitation was overcome by [11], who introduced auxiliary optimized variables that specify the contact information. However, the number of colliding body parts was still limited.

Animation as simulation control In recent years, the focus of research has shifted towards animation as a simulation control problem. Typically, one optimizes simulation control parameters such as time-varying actuation torques of character joints, and an off-the-shelf physics simulator is used to realize the movement. While spacetime optimization can be performed with gradient-based optimization methods like Sequential Quadratic Programming [6] or L-BFGS [11], simulation control is commonly approached with sampling-based, gradient-free optimization due to non-differentiable dynamics and/or multimodality [2], [3], [12]. This is also what our work focuses on; it remains as future work to extend our visualizations to analyzing spacetime optimization.

Trajectory and policy optimization Two main classes of approaches include trajectory and policy optimization. In trajectory optimization, one optimizes the time-varying control parameters directly, which can be done both offline [13], [14], [15] or online, while the character moves and acts [2], [3], [16]. In policy optimization, one optimizes the parameters of a policy function such as a neural network that maps character state to (approximately) optimal control, typically independent of the current simulation time. This can be done both using neuroevolution [17], [18] or Reinforcement Learning (RL), which has recently proven powerful even with complex humanoid movements [19], [20], [21], [22], [23]. Unfortunately policy optimization/learning can be computationally expensive with large neural networks, and may require careful curriculum design [24]. On the other hand, it can produce controllers that require orders of magnitude less computing resources after training, compared to using trajectory optimization to solve each required movement in an interactive application such as a video game. Trajectory and policy optimization approaches can also be combined [25], [26], [27], which allows one to adjust the trade-off between training time and runtime expenses. In this paper, we provide analyses of both trajectory and policy optimization landscapes.

Visualizing optimization Many optimization visualizations are problem-specific, utilizing the semantics of optimized parameters [28]. Visualization is also used for letting a user interact and inform optimization [29]; this, however, falls outside the scope of this paper. Considering non-interactive, generic methods applicable to continuous-valued optimization, landscape visualizations like the ones in Fig. 1 are a standard textbook method. Although it is technically trivial to extend this to higher-dimensional problems by visualizing the objective function on a random plane (a 2D subspace), Li et al. [5] only recently demonstrated that such random slices can provide meaningful insights and, likewise, have some predictive power on the difficulty of very high-dimensional optimization. Inspired by [5], we test the random slice visualization approach in a new domain, and also provide additional

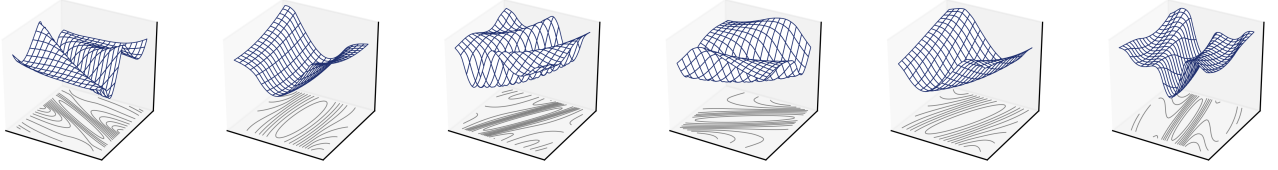


Fig. 2: Visualizing the inverted pendulum trajectory optimization objective of Equation 3 with different random basis vectors and $T = 100$. Although the plots are not exactly similar, they all exhibit the same overall structure, i.e., multimodality and an elongated, ill-conditioned optimum.

analyses of the method’s reliability and limitations. Other common methods for visualizing high-dimensional optimization include graphing the objective function along a straight line from the initial point to the found optimum [30], [31], [32], [33], or visualizing in a plane determined from the path taken during optimization [30]. There are also examples of visualizing movement optimization through a conversion to an interactive game or puzzle; in this case, players perform the optimization aided by predictive visualizations of how different actions affect the simulation state [34].

In computer animation and movement control research, objective functions are visualized occasionally, using various approaches to reduce the objectives to 2D. For example, Hämäläinen et al. [2] visualize contact discontinuities and multimodality in a 2D toy problem, and Sok et al. [35] visualize a high-dimensional multimodal objective with respect to two manually selected parameters. However, we know of no previous paper that focuses on visualizing movement optimization, or applies the 2D random slice approach of Li et al. [5] to movement optimization.

3 Test Problem: Inverted Pendulum Balancing

This section describes the inverted pendulum balancing problem that is used throughout Sections 4-8, before testing the visualization approach on the more complex simulated humanoid of Section 9. The pendulum is depicted in Fig. 3. Although it is simple, it offers the following benefits for analyzing movement optimization:

- In the trajectory optimization case, we know the true optimum. As the pendulum dynamics are differentiable, we can also compute the Hessian and its eigenvalues for further analysis. We implement this using Autograd [36].
- In the policy optimization case, we can use a simple P-controller as the parametric policy, which admits visualizing the full optimization landscape instead of only low-dimensional slices.

The dynamics governing the angle α_t , angular velocity ω_t , and control torque τ_t of the pendulum at timestep t are implemented as:

$$\omega_t = \omega_{t-1} + \delta(\tau_t + 0.5 l g \sin(\alpha_{t-1})), \quad (1)$$

$$\alpha_t = \alpha_{t-1} + \delta\omega_t, \quad (2)$$

where δ , l , and g are the simulation timestep, pendulum length, and force induced by gravity, respectively. We use $\delta = 0.1$, $l = 0.2$, and $g = 0.981$.

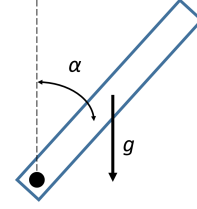


Fig. 3: The inverted pendulum model. The force exerted by gravity is denoted by g , and α denotes the angular deviation from an upright position.

3.1 Trajectory Optimization

Most of our trajectory optimization visualizations are generated from the problem of balancing a simple inverted pendulum, starting from an upward position such that the optimal torque sequence τ_1, \dots, τ_T is all zeros. The subscripts denote timestep indices and T is the planning horizon, i.e., the length of the simulated trajectory. The optimization objective is to minimize the trajectory cost \mathcal{C} computed as the sum of instantaneous costs:

$$\mathcal{C} = \sum_{t=1}^T (\alpha_t^2 + w\tau_t^2). \quad (3)$$

The cost is minimized when the pendulum stays upright at $\alpha = 0$ with zero torques. The relative importance of state cost α_t^2 and action cost τ_t^2 is adjusted by the multiplier w . We use $w = 1$ unless specified otherwise. The cost landscape is visualized in case of $T = 100$ in Fig. 2.

Some of our experiments convert the cost minimization problem into a reward maximization problem, computing trajectory reward \mathcal{R} using exponentiated costs as

$$\mathcal{R} = \sum_{t=1}^T (e^{-\alpha_t^2} + we^{-\tau_t^2}). \quad (4)$$

The reward formulation has been recently used with stellar results in the policy optimization of complex humanoid movements [20]. Exponentiation is also used in framing optimal control as estimation [37], [38].

3.2 Policy Optimization

In the case of policy optimization, we use the same pendulum simulation, with a minor adjustment. Instead of directly optimizing control torques, we use a policy $\tau_t = \pi_\theta(\mathbf{s}_t)$, parametrized by θ , where \mathbf{s} denotes pendulum state. The

optimization objective is to either minimize the expected trajectory cost, or maximize the expected reward, assuming that each trajectory or “episode” is started from a random initial state. As closed-form expressions for the expectations are not available, we replace them by averages. These are computed from 10 episodes, each started from a different initial pendulum angle.

In all the pendulum policy optimization visualizations, we use a simple P-controller as the policy:

$$\pi_{\theta}(\mathbf{s}_t) = \theta \alpha_t. \quad (5)$$

The benefit of this formulation is that we only have a single policy parameter to optimize; thus, we can visualize the full objective function, shown in Fig. 4. Despite the simplicity, this provides a multimodal optimization problem with properties similar to more complex problems. As shown in Fig. 4, there is global optimum at approximately $\theta = -0.1$, and also a false optimum near $\theta = 0$. At the false optimum, the action cost is minimized simply by letting the pendulum hang downwards. In a real-world case like controlling a simulated humanoid, the false optimum corresponds to resting on the ground with zero effort; readers familiar with humanoid control probably know that such a behavior is easy to elicit by having too large an effort cost.

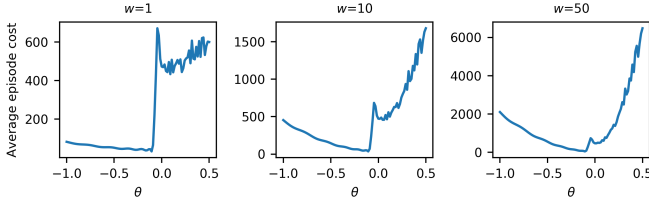


Fig. 4: Average episode cost (Equation 3) with $T = 200$, as a function of the P-controller policy parameter θ , and action cost weight w . A large w makes the local and true optima more equally good compared to the surrounding regions. This makes it more likely that a global Monte Carlo optimization method like CMA-ES will get attracted to the false optimum. The success of local gradient-based optimization depends on which cost basin the optimization is initialized in.

4 Effect of trajectory length

Figures 5 and 6 visualize the inverted pendulum trajectory and policy optimization landscapes, with different trajectory and episode lengths T . The figures yield two main insights:

- Trajectory optimization can become increasingly non-separable and ill-conditioned with large T .
- In both trajectory and policy optimization, the landscapes become more multimodal with large T .

4.1 Trajectory Optimization

The trajectory optimization landscapes of Fig. 5 are augmented with visualizations of the Hessian matrices of the cost function at the optimum. This allows further analysis of some important properties. A diagonal Hessian means that the optimization problem is separable, and the variables can

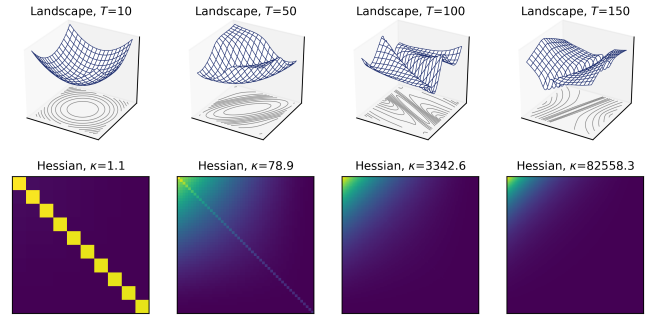


Fig. 5: Effect of trajectory length T on inverted pendulum trajectory optimization. The optimization problem becomes increasingly ill-conditioned and non-separable with longer action sequences. The bottom row shows the Hessian matrices at the optimal points with increasing T . κ denotes the condition number of the Hessian.

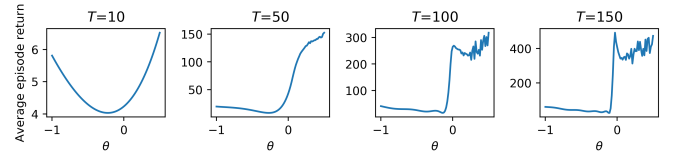


Fig. 6: Effect of trajectory length T on inverted pendulum trajectory policy optimization. Local optima become pronounced with large T , as the agent has more time to diverge and accumulate cost far from the desired states.

be optimized independent of each others. Strong off-diagonal elements imply that if one changes a variable, then one must also change some other variable to remain at the bottom of the valley in the landscape.

On the other hand, the eigenvalues of the Hessian measure curvature along the eigenvectors; the condition number κ , which denotes the ratio of the largest to smallest eigenvalues of the Hessian, is generally considered as a predictor of optimization difficulty. In the ideal case, the Hessian is a (scaled) identity matrix, i.e., a diagonal matrix where all the eigenvalues are the same; this indicates both separability and no ill-conditioning with $\kappa = 1$.

Intuitively, the ill-conditioning with a large T can be explained by the fact that perturbing an action of the optimal trajectory will lead to state divergence that accumulates over time. Thus, the total state cost is more sensitive to earlier actions, which leads to large differences in the Hessian eigenvalues. Non-separability stems from a need to adjust later actions to correct the state divergence. On the other hand, with a small T , the state has less time to diverge, and the τ^2 action cost dominates; this gives rise to the constant diagonal structure of the Hessian, as actions of all timesteps contribute equally and independently to the cost.

More formally, we may rewrite the cost function (Equation 3) of the trajectory optimization as:

$$\mathcal{C} = \sum_{t=1}^T \alpha_t(\tau_1, \dots, \tau_t)^2 + \sum_{t=1}^T w\tau_t^2, \quad (6)$$

emphasizing that the state α_t depends on the past actions τ_1

to τ_t . The Hessian of the cost function can then be expressed as:

$$\frac{\partial^2 \mathcal{C}}{\partial \tau^2} = 2 \sum_{t=1}^T (\nabla_{\tau} \alpha_t (\nabla_{\tau} \alpha_t)^T + \alpha_t \frac{\partial^2 \alpha_t}{\partial \tau^2}) + 2w\mathbf{I}. \quad (7)$$

The last term of Equation 7 is a diagonal matrix that does not depend on T , while the first term results in off-diagonal elements accumulated over time. As T grows larger, the first term begins dominating the second one, hence the decrease in separability (i.e., the fading diagonal shown in Fig. 5).

The increase of the condition number follows from the causality of the dynamics—future actions do not affect past states. This means that when $t = 1$, the first term of Equation 7 is a $T \times T$ matrix that has only one nonzero element at the upper-left corner, while, when $t = T$, the nonzero elements occupy the entire matrix. Summing up all the matrices from $t = 1$ to $t = T$, the first term of the Hessian will have larger values toward the upper-left corner, as evidenced by the bright areas in Fig. 5. The longer T is, the more disparate the upper-left and bottom-right corners become, leading to larger condition numbers.

Although it is well known that the difficulty of solving a trajectory optimization increases proportionally with the length of the planning horizon, our visualization shows that the difficulty is not solely caused by the increase in the size of the problem, but also by the increasingly ill-conditioned and non-separable objective function. Section 5 explains how the choice of action parameters can act as an efficient preconditioner.

4.2 Policy Optimization

In policy optimization, the optimized parameters have no similar dependency to the horizon T , as the effect of each policy parameter on the objective is averaged over all states. The separability of policy parameters depends on the function representation of the policy. For example, if the policy is represented as a multi-layer nonlinear neural network, the Hessian of the objective function is bound to be inseparable, as the neuron weights across layers are multiplied with each other when passing data through the network.

However, a longer horizon in policy optimization can induce a multimodal landscape in the objective function (Fig. 6). An explanation is that a longer time budget enables strategies that are not possible when the horizon is shorter. The “mollification” of the landscape at shorter horizons explains why policy optimization may benefit from a curriculum in which the planning horizon is gradually increased. The OpenAI Five Dota 2 bots provide a recent impressive example of this, gradually increasing the γ parameter of Proximal Policy Optimization [19] during training¹.

5 Effect of the Choice of Action Space

In the previous section, we saw that minimizing effort as the sum of squared actions gives rise to a strong constant diagonal in the Hessian, which in principle leads to better conditioned optimization. Parameterizing actions as torques, and having a squared torque cost term, is also common in continuous control benchmark problems such as OpenAI Gym MuJoCo

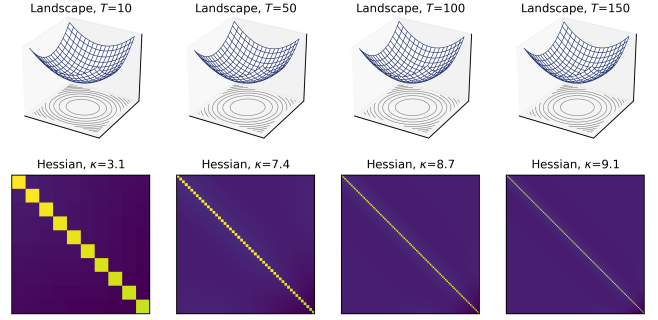


Fig. 7: Effect of episode length T on inverted pendulum trajectory optimization, when parameterizing actions as target angles. Remarkably, as opposed to the torque parameterization of Fig. 5, the landscape becomes much less ill-conditioned with large T .

[39]. However, [4] showed that parameterizing actions as target joint angles can make policy optimization more effective. The target poses that the policy outputs are converted to joint torques, typically using a P- or PD-controller. Such pose-based control is also common in earlier work on trajectory optimization [2], [3], [40], [41], often claimed to give better results than optimizing raw torques; despite this, comprehensive comparisons of control parameterizations are rare.

Fig. 7 shows that using actions as target angles $\bar{\alpha}$ does indeed scale better to long horizons. The strong constant diagonal of the Hessian does not fade out as much, and the condition number κ grows much slower with larger T . We compute the torques using a PD-controller: $\tau = k_p(\bar{\alpha} - \alpha) + k_d\omega$, using $k_p = 1$ and $k_d = -1$. The optimal $\bar{\alpha}$ sequence is all zeros, as the pendulum starts from an upright position with zero velocity.

An explanation for the effectiveness of the target angle parameterization is that *actions represent (partial) target states*. This makes the optimization of state cost or reward terms more separable and well-conditioned, and the Hessian closer to a scaled identity matrix. The optimal action at each timestep is more independent of preceding actions; a reasonable strategy is to always drive the pendulum towards the desired $\alpha = 0$. This explains why the Hessian is closer to diagonal. Furthermore, with the state-based control, perturbing early actions leads to less cumulative state divergence; hence, all actions contribute more equally to the cost. This explains why the spread of the diagonal values is low.

It should be noted that when using target joint angles with a character with an unactuated root, the parameterization cannot fully remove the dependencies between the actions of different timesteps. Deviations in initial actions can still lead to divergence, e.g. falling out of balance, which needs to be corrected by later actions.

We can also see the effect of action choice from the dynamics of the pendulum. Using torques as actions directly, the dynamic equation for the angle can be expressed as:

$$\alpha_t = \alpha_{t-1} + \delta\omega_{t-1} + \delta^2(\tau_t + 0.5 l g \sin(\alpha_{t-1})), \quad (8)$$

where the dependency to past states and actions is primarily due to the first two terms. When using target angles as

1. <https://blog.openai.com/openai-five/>

actions, we replace τ_t with $\bar{\alpha}_{t-1} - \alpha_{t-1} - \omega_{t-1}$ and arrive at a slightly different dynamic equation:

$$\alpha_t = (1 - \delta^2)\alpha_{t-1} + (\delta - \delta^2)\omega_{t-1} + \delta^2(\bar{\alpha}_{t-1} + 0.5 l g \sin(\alpha_{t-1})). \quad (9)$$

Now, the previous angle α_{t-1} and velocity ω_{t-1} have smaller multipliers $(1 - \delta^2)$ and $(\delta - \delta^2)$. This discount is applied recursively over time, such that the angle and velocity at n time steps ago are discounted by $(1 - \delta^2)^n$ and $(\delta - \delta^2)^n$. Using the same reasoning of causality as described in Section 4.1, the exponential reduction in dependency on previous states makes the Hessian of the cost function better conditioned (Equation 7).

5.1 Pose splines

In many papers, long action sequences are parameterized as "pose splines", i.e. parametric curves that define the target pose of the controlled character over time, which are then implemented using P- or PD-controllers [2], [15], [40]. The discussion above provides a strong motivation for this, as such parameterization achieves two things at once:

- Better separability and conditioning of the problem due to the choice of action space.
- Improved avoidance of ill-conditioning with long action sequences; in effect, the control points of a spline can be thought of as a shorter sequence of higher-level actions, each of which defines the instantaneous actions for multiple timesteps.

Naturally, using splines has aesthetic motivations as well, as they result in smoother motion with less frame-by-frame noise.

Experimental results comparing spline-based optimization with direct optimization of action sequences are provided in Section 9.

6 Effect of using rewards instead of costs

Fig. 8 shows the inverted pendulum trajectory optimization landscape using the reward function of Equation 4. Comparing this to Fig. 5, one notices the following:

- With large T , the landscape structures are essentially the same for both the cost function (Equation 3) and the reward function (Equation 4), with an elongated optimum and some local optima. The Hessian and $\kappa(T)$ at the optimal point are the same in both cases. Thus, in principle, the reward function should be as easy or hard to optimize as the cost function.
- On the other hand, at $T = 10$, the reward landscape shows some additional non-convexity. This suggests that in practice, the reward function may cause a performance hit in optimization. Section 8 provides evidence of this.

The main difference of the functions is that the sum of exponentiated costs is more tolerant to temporary deviations. A sum-of-squares cost function heavily penalizes the cost being large in even a single timestep, whereas the exponentiation clamps the rewards in the range $[0, 1]$. In principle, an agent could exploit the reward formulation by only focusing on some reward terms. We interpret the ridges in Fig. 8 ($T = 10$) as manifestations of this.

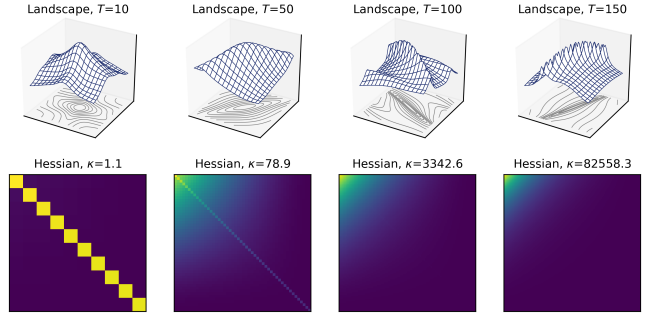


Fig. 8: Effect of trajectory length T on inverted pendulum trajectory optimization using the reward formulation of Equation 4. The landscape behaves similarly to the cost minimization in Fig. 5, except for the additional non-convexity at $T = 10$, and overall sharper ridges.

In Section 8, we shall see that the reward function can indeed result in worse optimization performance. However, the next section discusses how using rewards instead of costs can be more desirable when combined with the technique of early termination.

7 Effect of early termination

Standard continuous control policy optimization benchmark tasks [39], [42] utilize early termination of simulated episodes; this means that if the agent deviates from some desired region of the state space, e.g. a bipedal agent loses balance, the state is considered a terminal one; the agent stops receiving any rewards, and is reset to an initial state. Typically, termination greatly speeds up policy optimization, e.g. [20]. As far as we know, early termination has not been utilized in trajectory optimization, but there are no obstacles for this.

Figures 9 and 10 show the effect of terminating trajectories and episodes if $|\alpha| > 2.0$, i.e. if the pendulum deviates significantly from the desired upright pose. From the figures, one can gain two key insights:

- *Termination can greatly improve landscape convexity by removing false optima.* The pendulum cannot receive rewards from the local optimum of hanging downwards if termination prevents it from experiencing the corresponding states.
- *If the rewards are not strictly non-negative, new false optima are introduced to the landscape.* Naturally, if the agent is experiencing costs or negative rewards, a good strategy may be to terminate episodes as early as possible, which is what happens at the $\theta = 0.5$ local optimum in Fig. 10 (second subfigure from left, when termination is used with cost minimization). The problem can be mitigated by adding a termination penalty to the cost of the terminal simulation step, or a so-called alive bonus for all non-terminal states. This yields a clearly more convex landscape.

Although an alive bonus is used in many papers and benchmark tasks [24], [39], [43], we feel the danger of combining termination and negative rewards is not emphasized enough in previous literature. Our visualization highlights that the early termination can be a double-edged sword—it can be

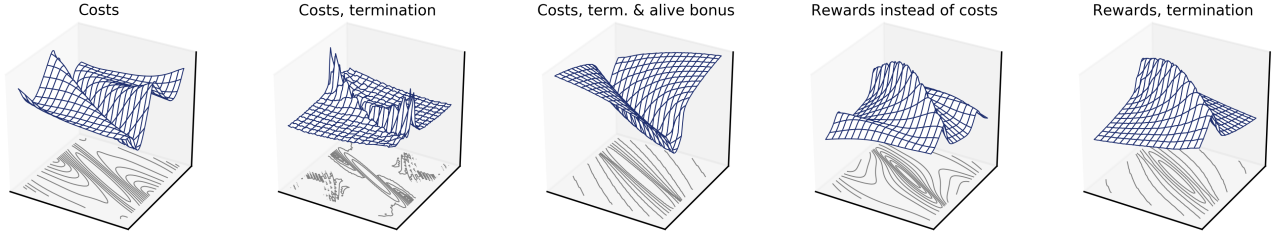


Fig. 9: Effect of early termination on inverted pendulum trajectory optimization landscape. Termination without an alive bonus increases multimodality of cost minimization, but makes reward maximization more convex.

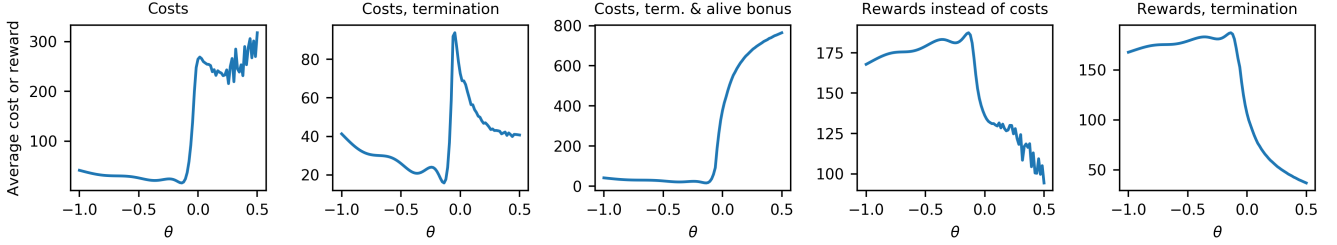


Fig. 10: Effect of termination on policy optimization. In cost minimization, termination creates a local minimum at $\theta = 0.5$, which drives the pendulum to termination to avoid accumulating more costs. Termination removes local optima when combined with an alive bonus or using rewards instead of costs.

harmful when the reward function consists of a mixture of positive and negative terms that are not well-balanced. Our visualization technique suggests that a good default strategy may be to use a combination of termination and rewards instead of costs. This way, one does not need to fine-tune the termination penalty or alive bonus. Furthermore, converting costs to rewards through exponentiation limits the results to the range $[0, 1]$, which is beneficial for most RL algorithms that feature some form of value function learning with neural networks. On the other hand, as shown in the next section, it can result in slower convergence.

8 Do visualizations predict optimization performance?

Fig. 11 tests how well the visualizations of the previous sections predict actual optimization performance in the inverted pendulum trajectory optimization. The figure compares quadratic costs to exponentiated rewards, and actions parameterized as torques to target angles implemented using a P-controller as in Section 5. To keep the figure readable, we did not include curves with and without termination; this does not make a big difference in the simple pendulum problem, and termination is further investigated in the next section. All the optimizations were performed using CMA-ES [44], [45], which is common in animation research, and is known to perform well with multimodal optimization tasks. A population size of 100 was used. To allow comparing both costs and rewards, progress is graphed as the Euclidian distance from the true optimum.

The results indicate that parameterizing actions as target angles is considerably more efficient, as predicted by the visualizations of Section 5. The exponential transform of costs

to rewards degrades performance, in line with the observations of Section 6.

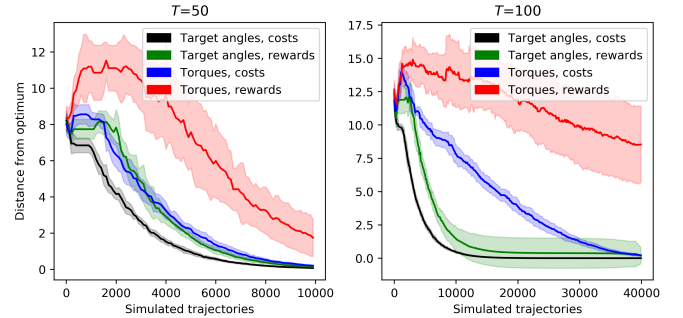


Fig. 11: Inverted pendulum trajectory optimization results, plotted as the mean of 10 runs of CMA-ES. We compare both cost minimization and reward maximization with actions parameterized both as torques and target angles. As predicted by our visualizations, the angle parameterization scales better for large T , and the reward maximization is less efficient.

9 Generalizability to more complex agents

This section tests the generalizability and usefulness of the visualization approach with a more complex agent, using both policy optimization and trajectory optimization. We use the 3D humanoid locomotion test from the Unity Machine Learning Agents framework [46] shown in Fig. 12. The test is designed for policy optimization; we modify it to also support trajectory optimization by starting each episode/trajectory from a fixed initial state, without randomization. We have also modified the code and environment to be fully determin-

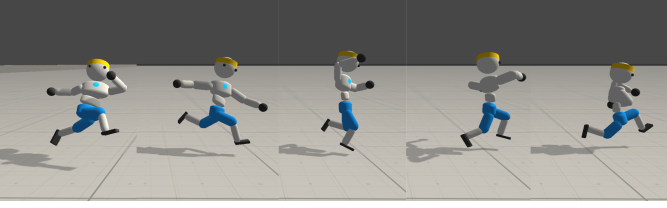


Fig. 12: 3D humanoid locomotion test from Unity Machine Learning Agents framework.

istic, i.e., non-reproducible simulation is not corrupting the landscape visualizations.

The action space is 39-dimensional, with actions defined as both joint target angles, as well as maximum torques that the joint motors are allowed to use for reaching their targets. We optimize with planning horizons of 1, 3, 5, and 10 seconds, resulting in 585, 1755, 2925, and 5850 optimized variables. The actions define target angles and maximum torques for the 16 joints of the humanoid. As the global optima are not known, we visualize the landscapes around the optima found in each test, following Li et al. [5]. We use a simulation timestep of 1/75 seconds, with control actions repeated for 5 timesteps, i.e., taking 15 actions per second.

As detailed below in Sections 9.1-9.4, the results support our earlier findings:

- 2D visualization slices reveal that trajectory optimization is highly multimodal, with optima that become narrower as the length of the planning horizon increases.
- As hypothesized in Section 5.1, utilizing a spline parameterization results in a more well-behaving landscape.
- Termination based on agent state reduces local optima in both trajectory and policy optimization.
- Policy optimization with neural network policies scales better for long planning horizons, with the landscape remaining essentially unchanged as the planning horizon grows. Notably, *policy optimization was more efficient than optimizing a single long trajectory*, even though our policy network has over 2M parameters to optimize, i.e., orders of magnitude more.

9.1 Trajectory Optimization

For trajectory optimization we use the recent highly scalable CMA-ES variant called LM-MA-ES [47]—as the physics simulator is not differentiable, gradient-based methods are not applicable. We visualize the landscapes around the optimum found by LM-MA-ES. Similar to CMA-ES, LM-MA-ES is a quasi-parameter-free method, typically only requiring adjustments to the iteration sampling budget (population size): this is increased from the recommended value for more difficult problems. The recommended budget—a logarithmic function of the number of variables—did not produce robust results. Instead, we used a 10 times larger budget in all the optimization runs.

We first tested trajectory optimization on action sequences of up to 5 seconds, i.e., up to 2925 optimized variables. This turned out to be a difficult task, resulting in somewhat unstable gaits even after hours of CPU time.

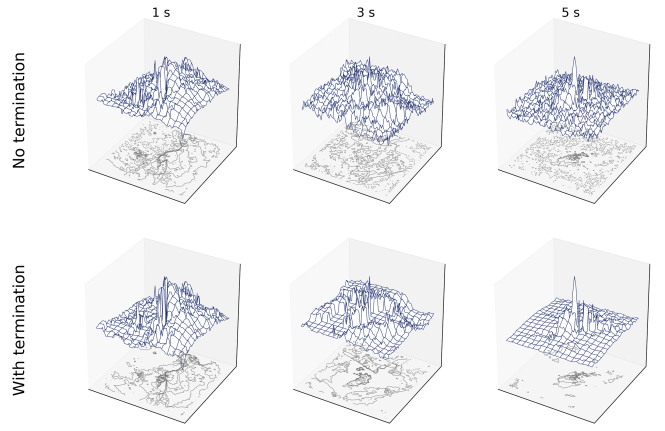


Fig. 13: Trajectory optimization landscapes of humanoid locomotion with different planning horizons. Termination removes local optima, producing a smoother landscape.

Here, 2D slice visualizations provided useful diagnostic information, revealing the difficult multimodality of the optimization problem. Fig. 13 shows the landscapes around the found local optimum for each tested planning horizon. Although there is a clear optimum in the center, the rest of the landscape is noisy and ill-conditioned. This is exacerbated with the longer planning horizons.

Fig. 13 also shows that termination helps remove local optima, and produces a smoother landscape. In the humanoid locomotion test, termination takes place when a body part other than the feet touches the ground.

9.2 Trajectory Optimization with Splines

In an effort to better handle long planning horizons, we tested trajectory optimization with spline-based parameterization. Instead of setting target angles and maximum torques for joints once every 5 timesteps, we interpolated the actions for each timestep using Catmull-Rom splines. The control points of the splines were optimized using LM-MA-ES, akin to the method presented by Hämäläinen et al. [2] for online optimization. Fig. 14 shows how this slightly improves the landscapes, with reduced noise and gentler slopes towards the optimum.

To provide a more direct comparison to optimizing action sequences without splines, Fig. 15 shows the non-spline trajectory optimization landscape around the action sequences resulting from the optimized splines of Fig. 14. Without the spline parametrization, the 3s and 5s landscapes degenerate, becoming noisier and more ill-conditioned.

Although trajectory optimization with splines results in cleaner landscapes and qualitatively better and smoother movement than non-spline trajectory optimization, finding good long trajectories required hours of CPU time. This motivated us to also test policy optimization, as explained below.

9.3 Policy Optimization

We trained a neural network policy for solving the humanoid locomotion task using Proximal Policy Optimization (PPO) [19] and different planning horizons. PPO utilizes episodic

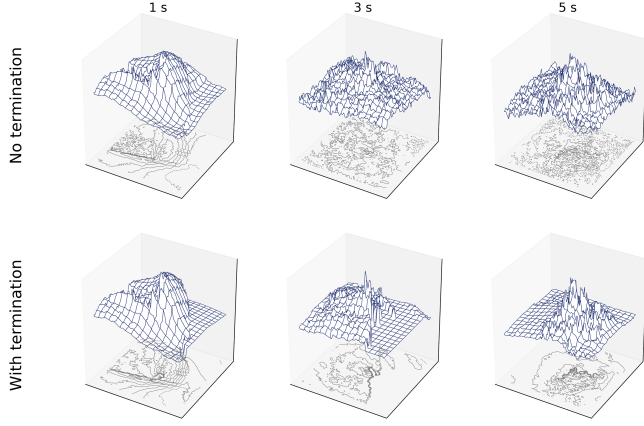


Fig. 14: Spline-parameterized trajectory optimization landscapes of humanoid locomotion. The landscapes exhibit less noise than the trajectory optimization landscapes of Fig. 13

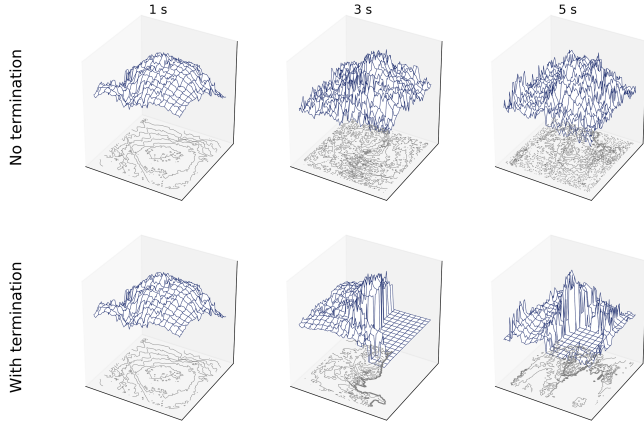


Fig. 15: Non-spline trajectory optimization landscapes around action sequences resulting from evaluating the optimal splines of Fig. 14. The landscapes become noisier and introduce more elongated ridges and valleys.

experience collection, i.e., the agent is started from some initial state, and explores states and actions until a terminal state, or maximum episode length, is reached. We used the planning horizon as the episode length.

The resulting landscapes, shown in Fig. 16, show a significant improvement in convexity, sphericity, and unimodality. The landscape also remains essentially unchanged with a longer planning horizon.

9.4 Scalability of Trajectory and Policy Optimization

The policy optimization visualizations suggest that as trajectory lengths increase, policy optimization should be more efficient than trajectory optimization. Fig. 17 provides evidence supporting this hypothesis. It compares optimization progress with different planning horizons and the three optimization strategies: LM-MA-ES, LM-MA-ES with splines, and PPO. With a trajectory length of 5 seconds, PPO already shows improved performance over the other methods, and the advantage is dramatically larger with 10 second trajectories. However, this comes with a caveat: with a neural network

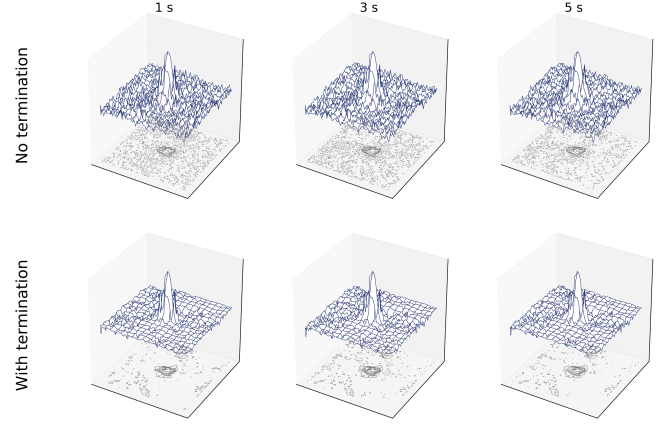


Fig. 16: Policy optimization landscapes of humanoid locomotion. As opposed to trajectory optimization, the task scales well with increasing planning horizon. Termination removes local optima.

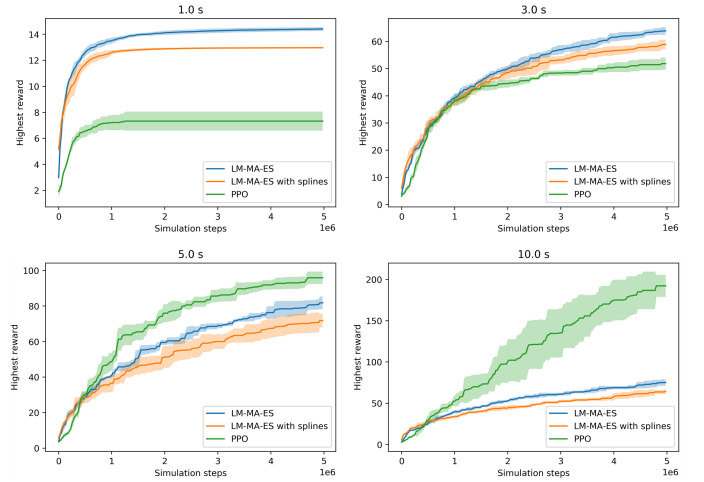


Fig. 17: Best trajectory or episode reward as a function of timesteps simulated during optimization, with different planning horizons and optimization approaches. The graphs show the mean and standard deviation of 5 independent optimization runs. Policy optimization with PPO scales significantly better for long planning horizons.

policy, there is more overhead per simulation step, as each optimization iteration requires training the policy and value networks using multiple minibatch gradient updates. Optimizing for 5 million timesteps with PPO was approximately 4 times slower than with LM-MA-ES in our tests, when measured in wall-clock time, using a single 4-core computer.

Interestingly, although the spline landscapes look slightly better in Fig. 14, spline trajectory optimization results in slightly lower rewards with a given simulation budget. A plausible explanation for this is that the Unity locomotion test has a very simple reward function that trajectory optimization can exploit with unnatural and jerky movements, as shown on the supplemental video² at 03:01. From a pure reward maximization perspective, a spline parameterization

2. https://youtu.be/5v_lsGCahSI

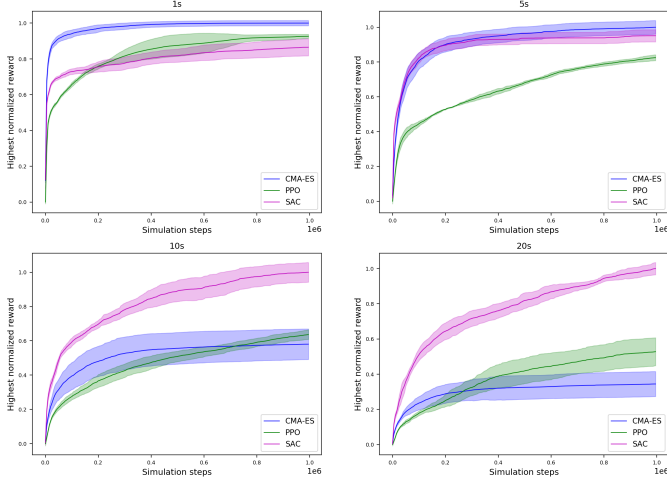


Fig. 18: Replicating the result of Fig. 17 in 4 OpenAI Gym MuJoCo tasks (Walker2d-v2, Hopper-v2, HalfCheetah-v2, Humanoid-v2), using 3 optimization methods (CMA-ES, PPO, and SAC), planning horizons 1s, 5s, 10s, and 20s, and 10 independent optimization runs per method and task. Each learning curve aggregates the results from all 4 tasks, using normalized episode/trajectory rewards.

that enforces a degree of smoothness may not be ideal and one may expect more gains if the reward function favors smooth movement. In general, best results are achieved when the action parameterization induces a useful prior for the optimization task.

Fig. 18 replicates the result of Fig. 17 with other agents and optimization methods, providing additional evidence that policy optimization scales better to long planning horizons. To generate the figure, we conducted trajectory and policy optimizations using 4 common OpenAI Gym [39] MuJoCo agents and locomotion tasks (or “environments” in the Gym lingo): 2D monopedal hopper (Hopper-v2), 2D bipedal walker (Walker2d-v2), 2D half quadruped (HalfCheetah-v2), and 3D humanoid (Humanoid-v2). In policy optimization, we tested both PPO and Soft Actor-Critic (SAC) [48], a more recent method that is growing in popularity. We used the Stable Baselines [49] PPO and SAC implementations with their default settings. In trajectory optimization, we used CMA-ES instead of LM-MA-ES, as it was easily available for the Python-based Gym framework. For each task and optimization method, we performed 10 independent training runs with different random seeds. To allow aggregating the convergence curves of all tasks in a single plot, we normalized the episode/trajectory rewards of each task over all optimization runs and methods to the range [0,1]. We used the default MuJoCo reward functions and episode termination, but removed the initial state randomization to allow direct comparison of trajectory and policy optimization, similar to the Unity humanoid tests above. We used a control frequency of 20Hz for all the tasks.

10 Properties and limitations of random 2D slice visualizations

So far, we have provided many examples of random 2D slice visualizations of high-dimensional objective functions. We have

also demonstrated that such visualizations have predictive power regarding the difficulty of optimization. However, as information is obviously lost in only evaluating the objective along a 2D slice, this section provides further analysis of the limitations of the approach. To allow investigating how the mathematical properties of objective functions manifest in the visualized slices, this section focuses on simple test functions with closed-form expressions. Such simple expressions are not available for real-life movement optimization objectives that depend on complex simulated dynamics, typically implemented using a black-box physics simulator.

10.1 2D Visualizations are Optimistic About Ill-conditioning

Consider the following cost function:

$$f(\mathbf{x}) = \sum_{i=1}^k x_i^2 + \epsilon \sum_{i=k+1}^d x_i^2 \quad (10)$$

$$= \|\mathbf{x}_{:k}\|^2 + \epsilon \|\mathbf{x}_{k:}\|^2, \quad (11)$$

where ϵ is a small constant, i.e., $f(\mathbf{x})$ mostly depends only on the first k optimized variables. $\mathbf{x}_{:k}$ denotes the projection of \mathbf{x} into the subspace of the first k dimensions. $\mathbf{x}_{k:}$ denotes the projection into the remaining dimensions. The Hessian of $f(\mathbf{x})$ is diagonal, containing the curvatures along the unit vectors as the diagonal elements. Curvature along the first k unit vectors equals 2 and curvature along the rest of the dimensions is 2ϵ . Thus, if $k \neq 0$ and $k \neq d$, the condition number $\kappa = 1/\epsilon$.

Geometrically, the isosurfaces of $f(\mathbf{x})$ are n -spheres elongated by a factor of $1/\sqrt{\epsilon}$ along the last $d - k$ dimensions. The visualized 2D isocontours correspond to the intersections of the isosurfaces with the visualization plane. This is illustrated in Fig. 19 and on the supplemental video for $d = 3$.

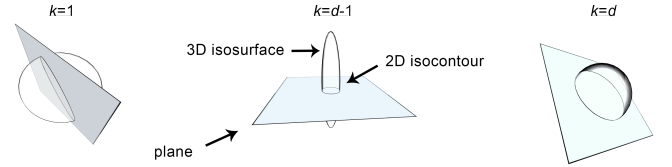


Fig. 19: The isosurfaces, random visualization planes, and 2D visualization isocontours for Equation 11 in the case of $d = 3$.

Investigating Fig. 19 reveals a basic property of the 2D visualizations: with the convex quadratic objective of Equation 11, the elongation of the 2D isocontours is less than or equal to the true elongation of the isosurfaces. *In other words*, $\kappa_{2D} \leq \kappa$.

This property follows from the isocontours corresponding to planar intersections of the isosurfaces. First, as illustrated in the middle of Fig. 19, it is possible to rotate the visualization plane such that the isocontours display less elongation. Second, the isocontours cannot display more than the real elongation; the visualized elongation is at maximum when one plane basis vector aligns with a direction of high elongation (vertical axis in the middle of Fig. 19), and the other basis vector aligns with a direction of low elongation, in which case the isocontours display the correct $\kappa_{2D} = \kappa = 1/\epsilon$.

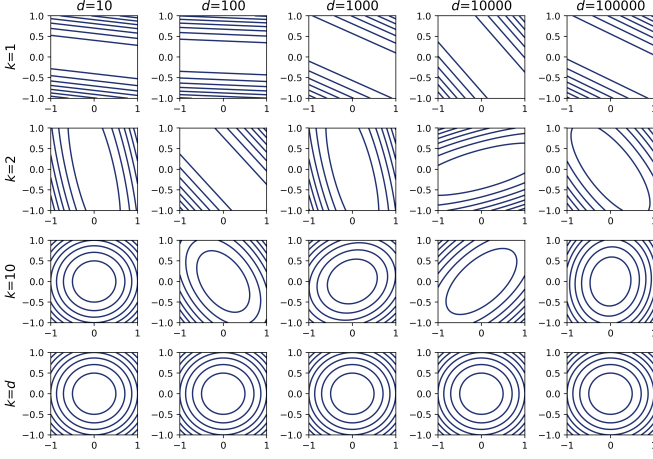


Fig. 20: Contour plots of random slices of $f(\mathbf{x})$ in Equation 12 with different k and d . Visualized ill-conditioning is more accurate with small k .

10.2 2D Visualizations Show Ill-conditioning More Accurately With Low Intrinsic Dimensionality

It turns out that *visualization accuracy depends on the intrinsic dimensionality k* . To analyze this, let us consider the extreme case of $\epsilon = 0$, i.e.,

$$f(\mathbf{x}) = f(\mathbf{x}_{:k}) = \|\mathbf{x}_{:k}\|^2. \quad (12)$$

Fig. 20 shows the contour plots of random 2D slices with different k and d . With low k , the visualized ill-conditioning is more accurate, independent of full problem dimensionality d . Deriving a closed-form expression of κ_{2D} as a function of d and k is beyond the scope of this paper. However, as shown below, $k = 1$ and $k = d$ result in the correct $\kappa_{2D} = \infty$ and $\kappa_{2D} = 1$, respectively.

Let $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ denote the slice basis vectors, with orthogonality $\mathbf{u}^T \mathbf{v} = 0$. On the visualization plane, $\mathbf{x} = p_1 \mathbf{u} + p_2 \mathbf{v} = [\mathbf{u} \ \mathbf{v}] \mathbf{p}$, where \mathbf{p} denotes the 2D position on the plane. Similarly, $\mathbf{x}_{:k} = [\mathbf{u}_{:k} \ \mathbf{v}_{:k}] \mathbf{p}$, and the objective can be expressed as:

$$f(\mathbf{x}_{:k}) = \|[\mathbf{u}_{:k} \ \mathbf{v}_{:k}] \mathbf{p}\|^2 \quad (13)$$

$$= \mathbf{p}^T \begin{bmatrix} \mathbf{u}_{:k}^T \mathbf{u}_{:k} & \mathbf{u}_{:k}^T \mathbf{v}_{:k} \\ \mathbf{v}_{:k}^T \mathbf{u}_{:k} & \mathbf{v}_{:k}^T \mathbf{v}_{:k} \end{bmatrix} \mathbf{p} \quad (14)$$

$$= \mathbf{p}^T \mathbf{A} \mathbf{p}. \quad (15)$$

Because \mathbf{A} is symmetric, the Hessian of the quadratic form w.r.t. \mathbf{p} is:

$$H(f(\mathbf{x}_{:k})) = \mathbf{A} + \mathbf{A}^T = 2\mathbf{A}. \quad (16)$$

The condition number $\kappa_{2D} = \kappa(H(f(\mathbf{x}_{:k}))) = \kappa(\mathbf{A})$, as the condition number is invariant to scaling the Hessian by a constant.

With $k=1$, the vectors $\mathbf{u}_{:k} = [u_1]$, $\mathbf{v}_{:k} = [v_1]$, and the determinant becomes zero:

$$\det(\mathbf{A}) = \mathbf{u}_{:k}^T \mathbf{u}_{:k} \mathbf{v}_{:k}^T \mathbf{v}_{:k} - \mathbf{u}_{:k}^T \mathbf{v}_{:k} \mathbf{v}_{:k}^T \mathbf{u}_{:k} \quad (17)$$

$$= u_1 u_1 v_1 v_1 - u_1 v_1 v_1 u_1 = 0. \quad (18)$$

This indicates at least one zero eigenvalue and, since \mathbf{A} is not a null matrix, some eigenvalue must also be nonzero, i.e., $\kappa(\mathbf{A}) = \max(\text{eig}(\mathbf{A})) / \min(\text{eig}(\mathbf{A})) = \infty$.

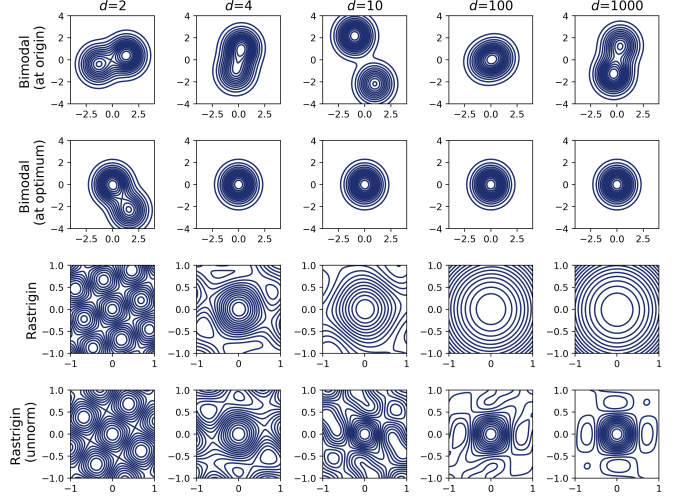


Fig. 21: Contour plots of random slices of multimodal test functions with different d . Large d can make multimodality less apparent, although it can be compensated by using unnormalized visualization basis vectors (bottom row). On the first row, the visualization plane is centered at the origin, between the optima. On the second row, it is centered at the optimum.

When k grows from 1 to d , the vectors $\mathbf{u}_{:k}, \mathbf{v}_{:k}$ gradually become closer to \mathbf{u}, \mathbf{v} , i.e., unit-length and orthogonal. Thus, the off-diagonal elements become zero and \mathbf{A} becomes the identity matrix, with $\kappa(\mathbf{A}) = 1$.

Although the quadratic $f(\mathbf{x})$ was chosen to be separable for easier mathematical analysis, the result generalizes to the arbitrarily rotated case.

10.3 2D Visualizations Can Be Optimistic About Multimodality

Fig. 21 shows contour plots of random visualization slices with different dimensionality d , using two multimodal test functions. Rastrigin's function is a standard multimodal optimization test function with the global minimum at the origin and infinitely many local minima:

$$f_{\text{Rastrigin}}(\mathbf{x}) = 10d + \sum_{i=1}^d [x_i^2 - 10 \cos(2\pi x_i)] \quad (19)$$

Additionally, we use the following bimodal function:

$$f_{\text{Bimodal}}(\mathbf{x}) = e^{-\frac{1}{2}\|\mathbf{x}-\mathbf{1}\|^2} + 0.8e^{-\frac{1}{2}\|\mathbf{x}+\mathbf{1}\|^2}, \quad (20)$$

where $\mathbf{1}$ denotes a vector of ones. We visualize this function both around the origin and around the dominant mode at $\mathbf{1}$.

Fig. 21 reveals two key insights:

- *Multimodality becomes less apparent with increasing dimensionality d .* The exception is the first row, where visualized multimodality only depends on plane rotation independent of d . This is because the visualization plane intersects the origin; in this case, the optima are always equally far from the plane and thus have similar influence on the visualized $f_{\text{Bimodal}}(\mathbf{x})$. However, when the plane intersects the optimum at $\mathbf{1}$, the other optimum tends to lie increasingly far from the

plane with increasing d , having a negligible effect on the visualization.

- The visualization of Rastrigin’s function illustrates how *landscape features may scale differently with dimensionality*. Rastrigin’s central mode becomes more dominant with large d . At the bottom of Fig. 21, we demonstrate how this can be compensated by omitting the unit-length normalization of the slice basis vectors, and instead normalizing them to the mean of their sampled lengths. Each basis vector element is sampled uniformly in the range $[-1, 1]$.

10.4 If 2D Visualizations Show Problems, There Really Are Problems

Visualized ill-conditioning is real The results above indicate that 2D visualizations have limited sensitivity as a diagnostic tool for detecting ill-conditioning and multimodality. Fortunately, $\kappa_{2D} \leq \kappa$ also means that the visualizations have high specificity, i.e. if the 2D isocontours are elongated, the problem is indeed ill-conditioned.

Visualized non-convexity indicates real non-convexity For a convex unimodal objective, the 2D visualization is likewise convex and unimodal. This follows from the intersection of two convex sets being convex. Each 2D isocontour encloses a set that is the intersection of two convex subsets of \mathbb{R}^d , i.e., the visualization plane and the volume enclosed by the corresponding isosurface. However, other non-convexity can be confused with multimodality. Consider a curved, banana-shaped 3D isosurface. It is possible to intersect this with a plane such that the resulting isocontours comprise two ellipses.

10.5 Summary of Limitations

In summary, random 2D visualization slices of high-dimensional objectives tend to be optimistic about both ill-conditioning and multimodality. However, this limitation is mitigated by the visualizations not showing illusory non-convexity or ill-conditioning. In other words, *as a diagnostic tool for detecting problems, random 2D visualizations have low sensitivity compensated by high specificity*.

Mitigating the low sensitivity is a potential topic for future work. For instance, if computing eigenvectors and eigenvalues of the Hessian or its low-rank approximation (e.g., [50]) is not too expensive, one could visualize using a 2D basis formed by the eigenvectors with lowest and highest eigenvalues. This way, the visualization would be in line with the condition number and show ill-conditioning with higher sensitivity. In this paper, however, we have focused on random visualization slices due to their simplicity and prior success in visualizing neural network loss landscapes [5]. Furthermore, at least for large policy networks with millions of parameters, computing eigenvector approximations is quite expensive.

11 Conclusion

We have presented several novel visualizations of continuous control trajectory and policy optimization landscapes, demonstrating the usefulness of the random 2D slice visualization approach of Li et al. [5] in this domain. We have also presented a mathematical analysis of the limitations of the random 2D slice visualizations.

Our visualizations provide new intuitions of movement optimization problems and explain why common best practices are powerful. The visualization approach can be used as a diagnostic tool for understanding why optimization does not converge, or progresses slowly. Even when a global optimum is not known, as in Section 9, it can be useful to plot the landscape around a found optimum: If the optimized movement is not satisfactory or one optimization approach performs worse than another, visualization can provide insights on why this is the case, e.g., due to ill-conditioning or multimodality.

We acknowledge that some of our results—e.g. the efficiency of episode termination—are already known to experienced readers. We do, however, provide novel visual evidence of the underlying reasons; for example, we show how termination based on agent state removes local optima in the space of optimized actions or policy parameters. This contributes to the understanding of movement optimization, and, as representative images are known to increase understanding and recall [51], it should also have pedagogical value in educating new researchers and practitioners.

To conclude, the key insights from our work can be summarized as:

- Random 2D slice visualizations are useful in analyzing high-dimensional movement optimization landscapes, and can predict movement optimization efficiency.
- The curse of dimensionality hits trajectory optimization hard, as it can become increasingly ill-conditioned with longer planning horizons. Policy optimization scales better in this regard. Perhaps counterintuitively, optimizing a neural network policy can be more efficient than optimizing a single action trajectory with orders of magnitude less parameters.
- Parameterizing actions as (partial) target states—e.g. target angles that the character’s joints are driven towards—is strongly motivated, as opposed to optimizing raw control torques. It can make trajectory optimization more well-conditioned and separable.
- Combining the two points above, one can explain the power of the common practice of optimizing splines that define time-varying target poses; pose parameterization leads to more well-conditioned optimization, and the spline control points define a shorter sequence of macro actions, which further counteracts ill-conditioning caused by sequence length. However, the smoothness constraints that splines impose on movements may not be ideal for all reward functions.
- Using early termination appears strongly motivated, as it typically results in a more convex landscape in both trajectory and policy optimization. However, combining termination with costs or negative rewards is dangerous, which our visualizations clearly illustrate.

In our future work, we aim to investigate the parameterization of actions as target states for complex controlled agents with unactuated roots. We hypothesize that this can be implemented for both trajectory and policy optimization, using a general-purpose neural network controller trained for reaching the target states.

Acknowledgements

This research has been supported by Academy of Finland grant 299358.

References

- [1] N. Hansen and S. Kern, "Evaluating the cma evolution strategy on multimodal test functions," in *International Conference on Parallel Problem Solving from Nature*. Springer, 2004, pp. 282–291.
- [2] P. Hämäläinen, S. Eriksson, E. Tanskanen, V. Kyrki, and J. Lehtinen, "Online Motion Synthesis Using Sequential Monte Carlo," *ACM Transactions on Graphics*, vol. 33, no. 4, p. 51, 2014.
- [3] P. Hämäläinen, J. Rajamäki, and C. K. Liu, "Online Control of Simulated Humanoids Using Particle Belief Propagation," *ACM Transactions on Graphics*, vol. 34, no. 4, p. 81, 2015.
- [4] X. B. Peng and M. van de Panne, "Learning locomotion skills using deeprl: does the choice of action space matter?" in *Proceedings of the ACM SIGGRAPH/Eurographics Symposium on Computer Animation*. ACM, 2017, p. 12.
- [5] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein, "Visualizing the loss landscape of neural nets," in *Advances in Neural Information Processing Systems*, 2018, pp. 6389–6399.
- [6] A. Witkin and M. Kass, "Spacetime constraints," in *Proc. SIGGRAPH '88*. New York, NY, USA: ACM, 1988, pp. 159–168. [Online]. Available: <http://doi.acm.org/10.1145/54852.378507>
- [7] M. F. Cohen, "Interactive spacetime control for animation," in *Proc. SIGGRAPH '92*. New York, NY, USA: ACM, 1992, p. 293–302. [Online]. Available: <http://doi.acm.org/10.1145/133994.134083>
- [8] A. C. Fang and N. S. Pollard, "Efficient synthesis of physically valid human motion," *ACM Trans. Graph.*, vol. 22, no. 3, pp. 417–426, 2003. [Online]. Available: <http://doi.acm.org/10.1145/1201775.882286>
- [9] A. Safonova, J. K. Hodgins, and N. S. Pollard, "Synthesizing physically realistic human motion in low-dimensional, behavior-specific spaces," *ACM Trans. Graph.*, vol. 23, no. 3, pp. 514–521, 2004. [Online]. Available: <http://doi.acm.org/10.1145/1186562.1015754>
- [10] K. Wampler and Z. Popović, "Optimal gait and form for animal locomotion," in *ACM Transactions on Graphics (TOG)*, vol. 28, no. 3. ACM, 2009, p. 60.
- [11] I. Mordatch, E. Todorov, and Z. Popović, "Discovery of complex behaviors through contact-invariant optimization," *ACM Trans. Graph.*, vol. 31, no. 4, p. 43:1–43:8, Jul. 2012. [Online]. Available: <http://doi.acm.org/10.1145/2185520.2185539>
- [12] L. Liu, K. Yin, M. van de Panne, T. Shao, and W. Xu, "Sampling-Based Contact-Rich Motion Control," *ACM Transactions on Graphics*, vol. 29, no. 4, pp. 128:1–128:10, Jul. 2010. [Online]. Available: <http://doi.acm.org/10.1145/1778765.1778865>
- [13] J. T. Ngo and J. Marks, "Spacetime constraints revisited," in *Proc. SIGGRAPH '93*. New York, NY, USA: ACM, 1993, p. 343–350. [Online]. Available: <http://doi.acm.org/10.1145/166117.166160>
- [14] M. Al Borno, M. de Lasa, and A. Hertzmann, "Trajectory optimization for full-body movements with complex contacts," *IEEE Transactions on Visualization and Computer Graphics*, vol. 19, no. 8, pp. 1405–1414, 2013.
- [15] K. Naderi, J. Rajamäki, and P. Hämäläinen, "Discovering and synthesizing humanoid climbing movements," *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, p. 43, 2017.
- [16] Y. Tassa, T. Erez, and E. Todorov, "Synthesis and stabilization of complex behaviors through online trajectory optimization," in *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*. IEEE, 2012, pp. 4906–4913.
- [17] T. Geijtenbeek, M. van de Panne, and A. F. van der Stappen, "Flexible Muscle-Based Locomotion for Bipedal Creatures," *ACM Transactions on Graphics*, vol. 32, no. 6, 2013.
- [18] F. P. Such, V. Madhavan, E. Conti, J. Lehman, K. O. Stanley, and J. Clune, "Deep neuroevolution: genetic algorithms are a competitive alternative for training deep neural networks for reinforcement learning," *arXiv preprint arXiv:1712.06567*, 2017.
- [19] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [20] X. B. Peng, P. Abbeel, S. Levine, and M. van de Panne, "Deepmimic: Example-guided deep reinforcement learning of physics-based character skills," *arXiv preprint arXiv:1804.02717*, 2018.
- [21] S. Lee, M. Park, K. Lee, and J. Lee, "Scalable muscle-actuated human simulation and control," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 4, pp. 1–13, 2019.
- [22] K. Bergamin, S. Clavet, D. Holden, and J. R. Forbes, "Drecon: data-driven responsive control of physics-based characters," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 6, pp. 1–11, 2019.
- [23] S. Park, H. Ryu, S. Lee, S. Lee, and J. Lee, "Learning predict-and-simulate policies from unorganized human motion data," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 6, pp. 1–11, 2019.
- [24] W. Yu, G. Turk, and C. K. Liu, "Learning symmetric and low-energy locomotion," *ACM Transactions on Graphics (TOG)*, vol. 37, no. 4, p. 144, 2018.
- [25] S. Levine and V. Koltun, "Guided policy search," in *International Conference on Machine Learning*, 2013, pp. 1–9.
- [26] I. Mordatch and E. Todorov, "Combining the benefits of function approximation and trajectory optimization," in *Robotics: Science and Systems*, 2014, pp. 5–32.
- [27] J. J. Rajamäki and P. Hämäläinen, "Continuous control monte carlo tree search informed by multiple experts," *IEEE transactions on visualization and computer graphics*, 2018.
- [28] C. V. Jones, "Visualization and optimization," *ORSA Journal on Computing*, vol. 6, no. 3, pp. 221–257, 1994.
- [29] D. Meignan, S. Knust, J.-M. Frayret, G. Pesant, and N. Gaud, "A review and taxonomy of interactive optimization methods in operations research," *ACM Transactions on Interactive Intelligent Systems (TiiS)*, vol. 5, no. 3, p. 17, 2015.
- [30] I. J. Goodfellow, O. Vinyals, and A. M. Saxe, "Qualitatively characterizing neural network optimization problems," *arXiv preprint arXiv:1412.6544*, 2014.
- [31] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang, "On large-batch training for deep learning: Generalization gap and sharp minima," *arXiv preprint arXiv:1609.04836*, 2016.
- [32] L. Dinh, R. Pascanu, S. Bengio, and Y. Bengio, "Sharp minima can generalize for deep nets," *arXiv preprint arXiv:1703.04933*, 2017.
- [33] L. N. Smith and N. Topin, "Exploring loss function topology with cyclical learning rates," *arXiv preprint arXiv:1702.04283*, 2017.
- [34] P. Hämäläinen, X. Ma, J. Takatalo, and J. Togelius, "Predictive physics simulation in game mechanics," in *Proceedings of the Annual Symposium on Computer-Human Interaction in Play*. ACM, 2017, pp. 497–505.
- [35] K. W. Sok, M. Kim, and J. Lee, "Simulating biped behaviors from human motion data," in *ACM SIGGRAPH 2007 Papers*, ser. SIGGRAPH '07. New York, NY, USA: Association for Computing Machinery, 2007, p. 107-es. [Online]. Available: <https://doi.org/10.1145/1275808.1276511>
- [36] D. Maclaurin, D. Duvenaud, and R. P. Adams, "Autograd: Effortless gradients in numpy," in *ICML 2015 AutoML Workshop*, 2015.
- [37] E. Todorov, "General duality between optimal control and estimation," in *Decision and Control, 2008. CDC 2008. 47th IEEE Conference on*. IEEE, 2008, pp. 4286–4292.
- [38] E. Todorov, "Efficient computation of optimal actions," *Proceedings of the national academy of sciences*, vol. 106, no. 28, 2009.
- [39] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," *arXiv preprint arXiv:1606.01540*, 2016.
- [40] M. Al Borno, M. De Lasa, and A. Hertzmann, "Trajectory optimization for full-body movements with complex contacts," *IEEE transactions on visualization and computer graphics*, vol. 19, no. 8, pp. 1405–1414, 2013.
- [41] J. Rajamäki and P. Hämäläinen, "Augmenting sampling based controllers with machine learning," in *Proceedings of the ACM SIGGRAPH / Eurographics Symposium on Computer Animation*, ser. SCA '17. New York, NY, USA: ACM, 2017, pp. 11:1–11:9. [Online]. Available: <http://doi.acm.org/10.1145/3099564.3099579>

- [42] Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. d. L. Casas, D. Budden, A. Abdolmaleki, J. Merel, A. Lefrancq *et al.*, “Deepmind control suite,” *arXiv preprint arXiv:1801.00690*, 2018.
- [43] A. Rajeswaran, S. Ghotra, B. Ravindran, and S. Levine, “Epopt: Learning robust neural network policies using model ensembles,” *arXiv preprint arXiv:1610.01283*, 2016.
- [44] N. Hansen and A. Ostermeier, “Completely derandomized self-adaptation in evolution strategies,” *Evolutionary computation*, vol. 9, no. 2, pp. 159–195, 2001.
- [45] N. Hansen, “The cma evolution strategy: A tutorial,” *arXiv preprint arXiv:1604.00772*, 2016.
- [46] A. Juliani, V.-P. Berges, E. Vckay, Y. Gao, H. Henry, M. Mattar, and D. Lange, “Unity: A general platform for intelligent agents,” *arXiv preprint arXiv:1809.02627*, 2018.
- [47] I. Loshchilov, T. Glasmachers, and H.-G. Beyer, “Large scale black-box optimization by limited-memory matrix adaptation,” *IEEE Transactions on Evolutionary Computation*, 2018.
- [48] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” in *International Conference on Machine Learning*, 2018, pp. 1861–1870.
- [49] A. Hill, A. Raffin, M. Ernestus, A. Gleave, A. Kanervisto, R. Traore, P. Dhariwal, C. Hesse, O. Klimov, A. Nichol, M. Plappert, A. Radford, J. Schulman, S. Sidor, and Y. Wu, “Stable baselines,” <https://github.com/hill-a/stable-baselines>, 2018.
- [50] K.-C. Li, “On principal hessian directions for data visualization and dimension reduction: Another application of stein’s lemma,” *Journal of the American Statistical Association*, vol. 87, no. 420, pp. 1025–1039, 1992.
- [51] R. N. Carney and J. R. Levin, “Pictorial illustrations still improve students’ learning from text,” *Educational psychology review*, vol. 14, no. 1, pp. 5–26, 2002.



Perttu Hämäläinen received an M.Sc.(Tech) degree from Helsinki University of Technology in 2001, an M.A. degree from the University of Art and Design Helsinki in 2002, and a doctoral degree in computer science from Helsinki University of Technology in 2007. Presently, Hämäläinen is an associate professor at Aalto University, publishing on human-computer interaction, computer animation, machine learning, and game research. Hämäläinen is passionate about human movement in its many forms, ranging from analysis and simulation to first-hand practice of movement arts such as parkour or contemporary dance.



Juuso Toikka received his M.Sc.(Tech) degree in Computer Science from Aalto University in 2019. While this manuscript was in preparation, Toikka worked as a research assistant at the Department of Computer Science at Aalto University, Finland, but he has since moved on to Ubisoft RedLynx, pursuing a game industry career. His professional interests include animation tools, procedural animation, movement control optimization, reinforcement learning, and emergence in games.



Amin Babadi is a doctoral candidate at the Department of Computer Science, Aalto University, Finland. His research focuses on developing efficient, creative movement artificial intelligence for physically simulated characters in multi-agent settings. Babadi has previously worked on three commercial games, developing AI, animation, gameplay, and physics simulation systems.



C. Karen Liu is an associate professor in the Department of Computer Science at Stanford University. She received her Ph.D. degree in Computer Science from the University of Washington. Liu’s research interests are in computer graphics and robotics, including physics-based animation, character animation, optimal control, reinforcement learning, and computational biomechanics. She developed computational approaches to modeling realistic and natural human movements, learning complex control policies for humanoids and assistive robots, and advancing fundamental numerical simulation and optimal control algorithms. The algorithms and software developed in her lab have fostered interdisciplinary collaboration with researchers in robotics, computer graphics, mechanical engineering, biomechanics, neuroscience, and biology. Liu received a National Science Foundation CAREER Award, an Alfred P. Sloan Fellowship, and was named Young Innovators Under 35 by Technology Review. In 2012, Liu received the ACM SIGGRAPH Significant New Researcher Award for her contribution in the field of computer graphics.

Paper Supplement: Visualizing Movement Control Optimization Landscapes

This supplementary document presents additional results to augment the paper’s Section 9.4. Recall that the central result of Section 9.4 is that policy optimization scales better to long trajectories/episodes, although a policy neural network typically has orders of magnitude more parameters to optimize than a single trajectory (even a long one). This was tested with multiple locomotion tasks and optimizers: A Unity Machine Learning Agents 3D humanoid (optimized using LM-MA-ES and PPO) and four different MuJoCo agents: A 2D monopedal hopper (Hopper-v2), 2D bipedal walker (Walker2d-v2), 2D half quadruped (HalfCheetah-v2), and a 3D humanoid (Humanoid-v2), optimized using CMA-ES, PPO, and SAC. The optimization landscape visualizations of section 9.4—from the Unity humanoid locomotion case—agree on the result, displaying much less multimodality and ill-conditioning in the policy optimization case.

Similar landscape plots of the MuJoCo agents are included below. All landscape visualizations are centered around the found optima (a vector of control torques for each time step in trajectory optimization, or a vector of neural network parameters in policy optimization). The visualizations were computed using grids of 100×100 points, computing the mean return of 10 trajectories/episodes for each grid point. To improve visual clarity, all landscapes were also filtered using Gaussian blur with $\sigma = 1.0$.

These MuJoCo landscapes support the results of Section 9.4. This is clearest in the hopper landscapes shown in Fig. 22.

CMA-ES trajectory optimization landscapes become increasingly ill-conditioned with long trajectories, with narrow ridges where an optimizer typically zigzags back and forth over the ridge, making very slow progress along the ridge. In contrast, the policy optimization landscapes show almost spherical optima. Note that the policy networks for PPO and SAC have different parameter counts (as per the default parameters of the Stable Baselines implementations that we used). Thus, the landscapes cannot display exactly same optima.

The plots for the other MuJoCo environments (Fig. 23-25) exhibit similar qualities, although less clearly. The trajectory optimization landscapes also become increasingly multimodal and/or noisy with longer trajectories. It should be noted that each landscape’s vertical axis is normalized to show maximal detail, i.e., the heights of the optima in different landscapes cannot be directly compared.

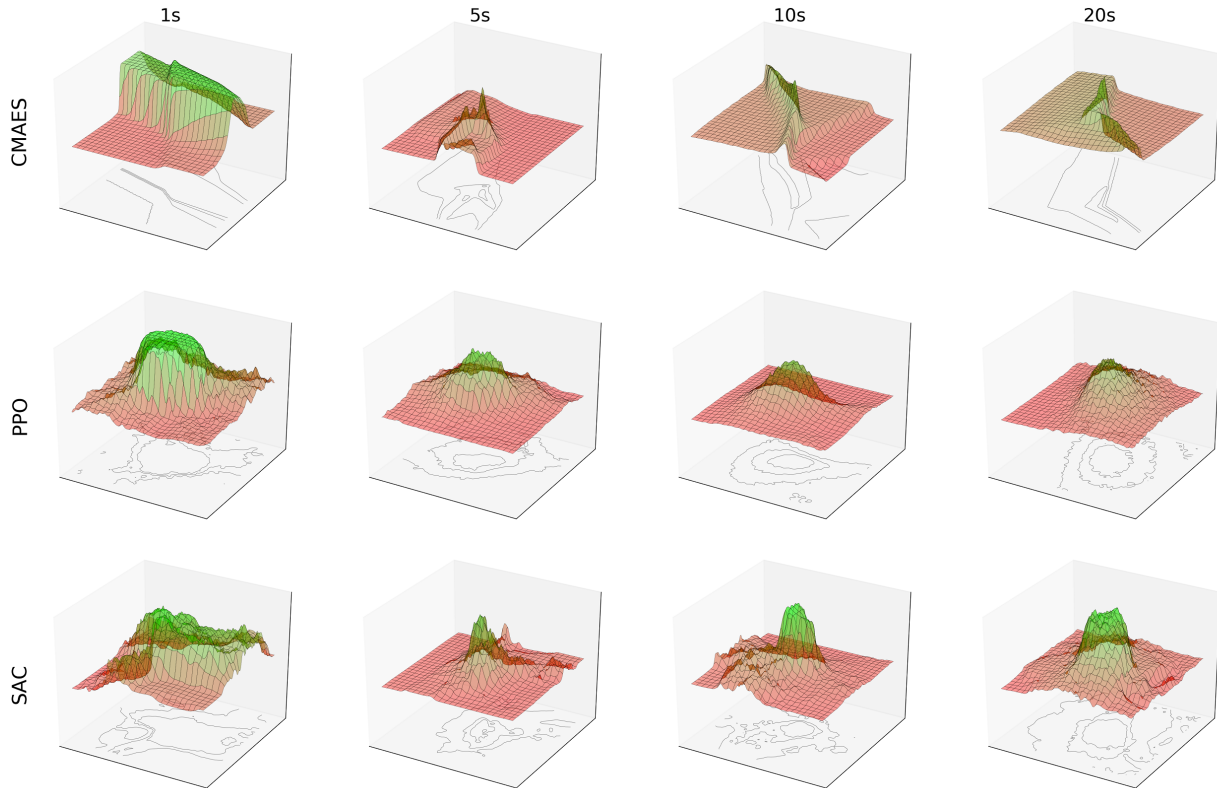


Fig. 22: Trajectory optimization (CMA-ES) and policy optimization (PPO, SAC) landscapes for the Hopper-v2 MuJoCo environment, with trajectory/episode lengths ranging from 1 to 20 seconds. Trajectory optimization becomes highly ill-conditioned for long trajectories.

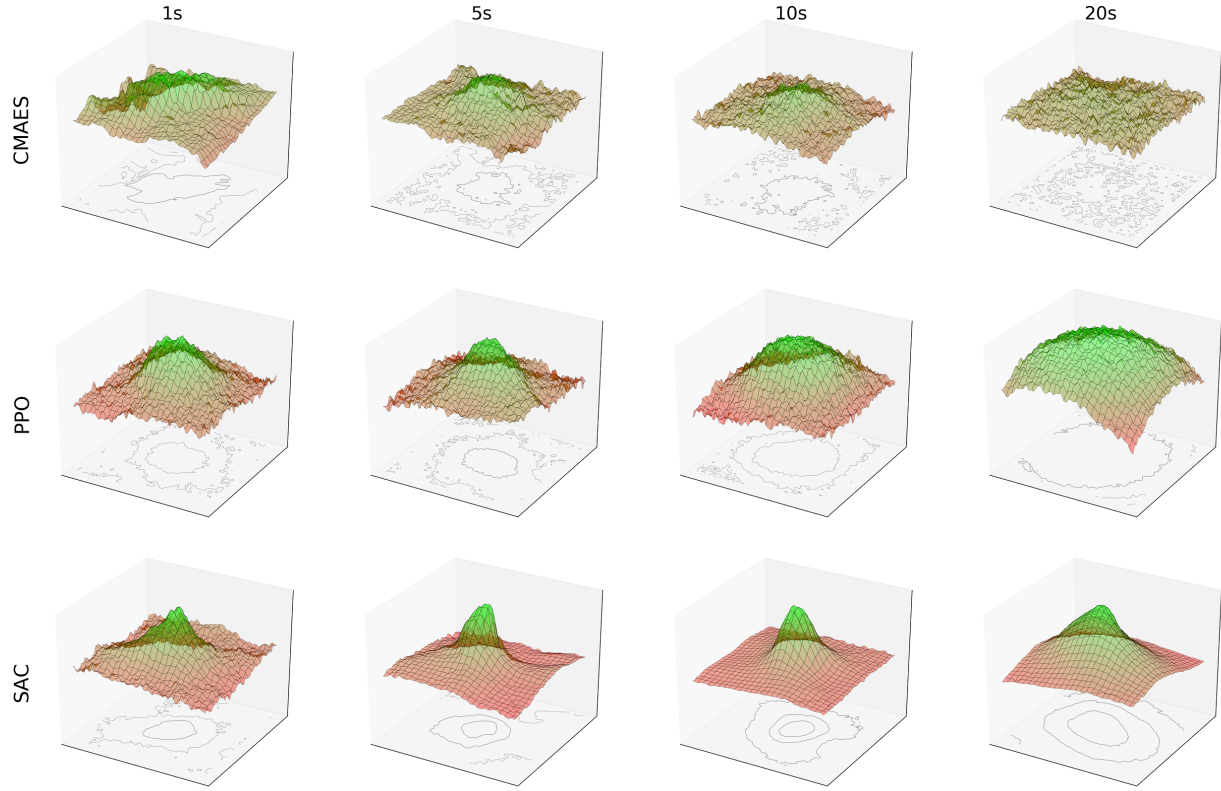


Fig. 23: Trajectory optimization (CMA-ES) and policy optimization (PPO, SAC) landscapes for the HalfCheetah-v2 MuJoCo environment, with trajectory/episode lengths ranging from 1 to 20 seconds.

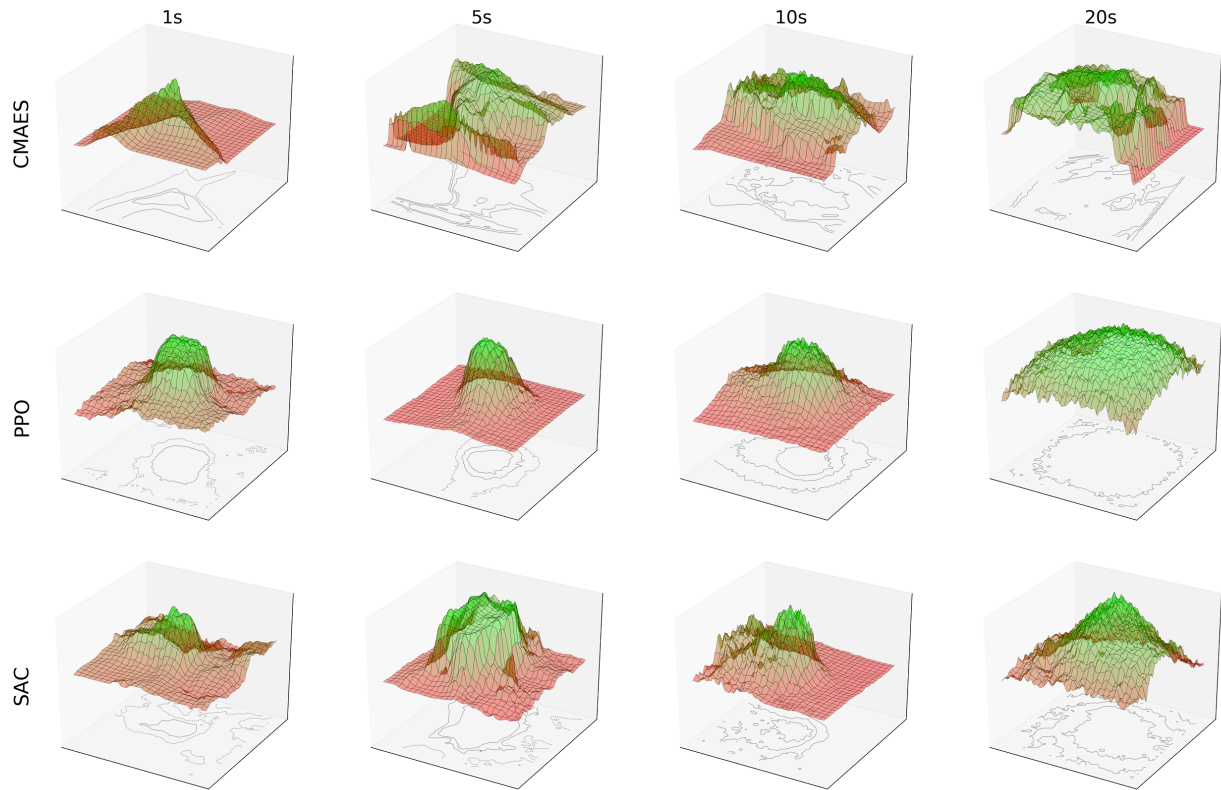


Fig. 24: Trajectory optimization (CMA-ES) and policy optimization (PPO, SAC) landscapes for the Walker2d-v2 MuJoCo environment, with trajectory/episode lengths ranging from 1 to 20 seconds.

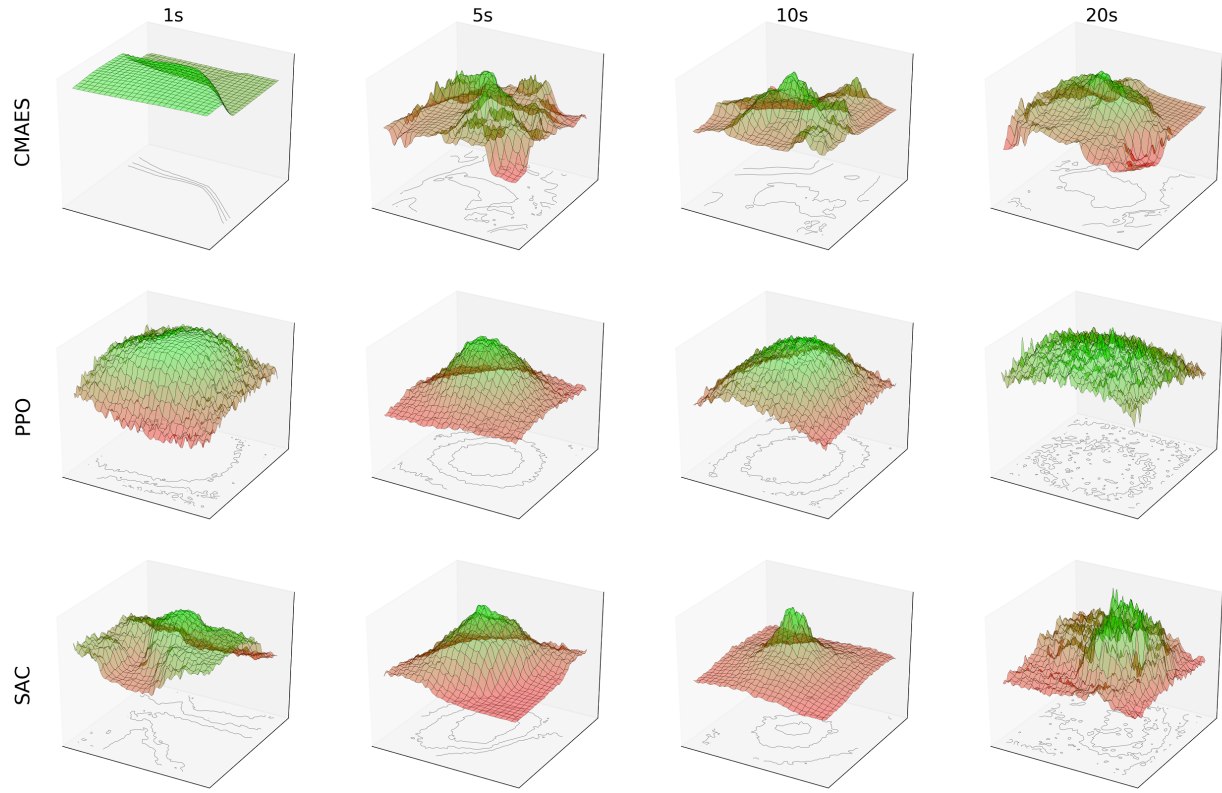


Fig. 25: Trajectory optimization (CMA-ES) and policy optimization (PPO, SAC) landscapes for the Humanoid-v2 MuJoCo environment, with trajectory/episode lengths ranging from 1 to 20 seconds.