

Grounded Action Transformation for Robot Learning in Simulation

Josiah P. Hanna, Peter Stone

Dept. of Computer Science
 The University of Texas at Austin
 Austin, TX 78712 USA
 {jphanna,pstone}@cs.utexas.edu

Abstract

Robot learning in simulation is a promising alternative to the prohibitive sample cost of learning in the physical world. Unfortunately, policies learned in simulation often perform worse than hand-coded policies when applied on the physical robot. Grounded simulation learning (GSL) promises to address this issue by altering the simulator to better match the real world. This paper proposes a new algorithm for GSL – Grounded Action Transformation – and applies it to learning of humanoid bipedal locomotion. Our approach results in a 43.27% improvement in forward walk velocity compared to a state-of-the-art hand-coded walk. We further evaluate our methodology in controlled experiments using a second, higher-fidelity simulator in place of the real world. Our results contribute to a deeper understanding of grounded simulation learning and demonstrate its effectiveness for learning robot control policies.

Introduction

Manually designing control policies for every possible situation a robot could encounter is impractical. Reinforcement learning (RL) provides a promising alternative to hand-coding skills. Recent applications of RL to high dimensional control tasks have seen impressive successes within simulation (Schulman et al. 2015; Lillicrap et al. 2015). Unfortunately, a large gap exists between what is possible in simulation and the reality of robot learning. State-of-the-art learning methods require thousands of episodes of experience which is impractical for a physical robot. Aside from the time it would take, collecting the required training data may lead to substantial wear on the robot. Furthermore, as the robot explores different policies it may execute unsafe actions which could damage the robot.

An alternative to learning directly on the robot is learning in simulation (Cutler and How 2015; Koos, Mouret, and Doncieux 2010). Simulation is a valuable tool for robotics research as execution of a robotic skill in simulation is comparatively easier than real world execution. However, the value of simulation learning is limited by the inherent inaccuracy of simulators in modeling the dynamics of the physical world (Kober, Bagnell, and Peters 2013). As a result, learning that takes place in a simulator is unlikely to improve real world performance.

Copyright © 2017, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

Grounded Simulation Learning (GSL) is a framework for learning with a simulator in which the simulator is modified with data from the physical robot, learning takes place in simulation, the new policy is evaluated on the robot, and data from the new policy is used to further modify the simulator (Farchy et al. 2013). The paper introducing GSL demonstrates the effectiveness of the method in a single, limited experiment, by increasing the forward walking velocity of a slow, stable bipedal walk by 26.7%.

This paper introduces a new algorithm – Grounded Action Transformation (GAT) – for simulator grounding within the GSL framework. The algorithm is fully implemented and evaluated using a high-fidelity simulator as a surrogate for the real world. The results of this study contribute to a deeper understanding of transfer from simulation methods and the effectiveness of GAT. In contrast to prior work, our real-world evaluation of GAT starts from a state-of-the-art walk engine as the base policy, and nonetheless is able to improve the walk velocity by over 43%, leading to what may be the fastest known walk on the SoftBank NAO robot.

Background

Preliminaries

We formalize robot skill learning as a *reinforcement learning* (RL) problem (Sutton and Barto 1998). At discrete time-step t the robot takes action $A_t \sim \pi(\cdot|S_t)$ according to a policy, π , which is a distribution over actions, $A_t \in \mathcal{A}$, conditioned on the current state, $S_t \in \mathcal{S}$. The environment, E , responds with $S_{t+1} \sim P(\cdot|S_t, A_t)$ according to the dynamics, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}_{\geq 0}$ which is a probability density function over next states conditioned on the current state and action. A trajectory of length L is a sequence of states and actions, $\tau := S_0, A_0, \dots, S_L, A_L$. We also define a cost function, c , which assigns a scalar cost to each (s, a) pair. We denote the value of policy, π , as $J(\pi) := \mathbb{E}_{\tau \sim Pr(\cdot|\pi)} \left[\sum_{t=0}^L c(S_t, A_t) \right]$ where $Pr(\tau|\pi)$ is the probability of τ when selecting actions according to π .

We assume π is parameterized by a vector θ and denote the parameterized policy as π_θ . Since θ determines the policy distribution we overload notation by referring to π_θ as θ . Given initial parameters θ_0 , the goal of policy improvement is to find θ' such that $J(\theta') < J(\theta_0)$. In this work, θ is a deterministic policy that maps state observations to an action

vector \mathbf{a}_t . Each component of \mathbf{a}_t , \mathbf{a}_t^i , is the desired joint angle for degree-of-freedom i . At each time-step, t , θ selects \mathbf{a}_t according to the robot's configuration in joint space, \mathbf{x}_t , high level intention commands (e.g., walk forward at 75% of maximum velocity), ω_t , and a sensor vector, ψ_t , of inertial measurement (IMU) and foot sensor readings. The state of the robot can be fully described as $\mathbf{s}_t = \langle \mathbf{x}_t, \dot{\mathbf{x}}_t, \ddot{\mathbf{x}}_t, \omega_t, \psi_t \rangle$ where $\dot{\mathbf{x}}_t$ and $\ddot{\mathbf{x}}_t$ are the time derivatives of \mathbf{x}_t . Since the robot only observes \mathbf{x}_t , ω_t , and ψ_t , E is partially observable. Note that while the high level intention commands are determined by the robot, from the point-of-view of θ , these directional commands are state features.

In this paper, learning takes place in a simulator which is an environment, E_{sim} , that models E . Specifically E_{sim} has the same state-action space as E but inevitably a different dynamics distribution, P_{sim} . In many robotics settings, c is user defined and thus is identical for E and E_{sim} . However, the different dynamics distribution mean that for any policy, θ , $J(\theta) \neq J_{\text{sim}}(\theta)$ since θ induces a different trajectory distribution in E than in E_{sim} . Thus θ' with $J_{\text{sim}}(\theta') < J_{\text{sim}}(\theta_0)$ does not imply $J(\theta') < J(\theta_0)$ – in fact $J(\theta')$ could even be worse than $J(\theta_0)$. In practice and in the literature, learning in a simulator frequently leads to catastrophic degradation of J . This paper explores methods for learning in E_{sim} that result in lower policy cost.

Related Work

This section surveys previous research in simulation-transfer methods. Our work also relates to model-based reinforcement learning, however, we restrict our attention here to methods directly concerned with learning in simulation.

One approach to simulation-transfer is using experience in simulation to reduce learning time on the physical robot. Cutler et al. (2014) use lower fidelity simulators to narrow the action search space for faster learning in higher fidelity simulators or the real world. Cully et al. (2015) use a simulator to estimate fitness values for low-dimensional robot behaviors which gives the robot prior knowledge of how to adapt its behavior if it becomes damaged during real world operation. Cutler et al. (2015) use experience in simulation to estimate a prior for a Gaussian process model to be used with the PILCO learning algorithm (Deisenroth and Rasmussen 2011). Rusu et al. (2016a; 2016b) introduce progressive neural network policies which are initially trained in simulation before a final period of learning in the true environment. In contrast to these methods, our method performs all learning in a grounded simulator and only uses the physical robot for policy evaluation.

Another class of simulation-transfer methods optimize an objective besides J_{sim} in simulation. Rajeswaran et al. (2016) use a set of different simulators to learn robust policies that can perform well in a variety of environments. Koos et al. (2010) use multi-objective optimization to find policies that trade off between optimizing $J_{\text{sim}}(\pi)$ and *transferability*. Iocchi et al. (2007) compute a correction function from $J(\pi)$ such that $J_{\text{sim}}(\pi) = J(\pi)$ and then optimize the corrected objective. In contrast, we directly optimize J_{sim} within a modified simulator.

Christiano et al. (2016) turn simulation policies into real

world policies by transforming policy actions so that they produce the same effect that they did in simulation. This approach is complementary to ours although it also requires a simulator which can be reset to an arbitrary state during policy execution.

Grounded Simulation Learning

In this section we introduce the Grounded Simulation Learning (GSL) framework as presented by Farchy et al. (2013). GSL improves robot learning in simulation by using state transition data from the physical system to modify E_{sim} such that the modified E_{sim} is a higher fidelity model of E . The process of making the simulator more like the real world is referred to as *grounding*.

The GSL framework assumes the following:

1. There is an imperfect simulator, $E_{\text{sim}} = \langle \mathcal{S}, \mathcal{A}, P_{\text{sim}}, c \rangle$, that models the environment of interest, E . Furthermore, E_{sim} must be *modifiable*. A modifiable simulator has parameterized transition probabilities $P_\phi(\cdot|s, a) := P_{\text{sim}}(\cdot|s, a; \phi)$ where the vector ϕ can be changed to produce, in effect, a different simulator.
2. Additionally, GSL assumes the physical robot can record a data set, \mathcal{D} of trajectories when executing any policy θ .
3. Finally, GSL requires a policy improvement routine, *optimize*, that searches for θ which decreases $J_{\text{sim}}(\theta)$. The *optimize* routine returns a set of candidate policies, Π_c to evaluate on the physical robot.

Let $d(p, q)$ be a measure of similarity between probabilities p and q . GSL grounds E_{sim} by finding ϕ^* such that:

$$\phi^* = \underset{\phi}{\operatorname{argmin}} \sum_{\tau \in \mathcal{D}} d(\operatorname{Pr}(\tau|\theta), \operatorname{Pr}_{\text{sim}}(\tau|\theta, \phi)) \quad (1)$$

where $\operatorname{Pr}(\tau|\theta)$ is the probability of observing trajectory τ on the physical robot and $\operatorname{Pr}_{\text{sim}}(\tau|\theta, \phi)$ is the probability of τ in the simulator with dynamics parameterized by ϕ . For instance, if $d(p, q) := \log p - \log q$ for two probabilities p and q then ϕ^* minimizes the Kullback-Leibler divergence between $\operatorname{Pr}(\cdot|\theta)$ and $\operatorname{Pr}_{\text{sim}}(\cdot|\theta, \phi^*)$.

Assuming GSL can solve (1), the framework is as follows:

1. Execute policy θ_0 on the physical robot to collect a data set of trajectories, \mathcal{D} .
2. Use \mathcal{D} to find ϕ^* that satisfies Equation 1.
3. Use *optimize* with J_{sim} and P_{ϕ^*} to learn a set of candidate policies Π_c in simulation which are expected to perform well on the physical robot.
4. Evaluate each proposed $\theta_c \in \Pi_c$ on the physical robot and return the policy, θ_1 , with minimal J .

GSL can be applied iteratively with θ_1 being used to collect more trajectories to ground the simulator again before learning θ_2 . The re-grounding step is necessary since changes to θ result in changes to the distribution of states that the robot observes. When this happens, a simulator that has been modified with data from the state distribution of θ_0 may be a poor model under the state distribution of θ_1 .

Farchy et al. proposed a GSL algorithm, which we refer to as GUIDED GSL. GUIDED GSL blends simulator modification with human guided policy optimization. This algorithm demonstrated the efficacy of GSL in optimizing forward walk velocity of a bipedal robot. The robot in that research began learning with a slow but stable walk policy. Four iterations of GUIDED GSL led to an improvement of over 26% on this baseline walk. Two limitations of GUIDED GSL are that:

1. It makes the assumption that actions in simulation achieve their desired effect instantaneously.
2. It required a human expert to manually select which components of θ were allowed to change between iterations of the optimize routine.

In this work, we introduce an enhanced GSL methodology that eliminates both assumptions and optimizes one of the fastest available NAO walk engines.

Robot Platform and Task

Before presenting GAT, our new GSL algorithm, we describe the robot platform, initial walk policy, and simulators which we use for evaluation. While our method is applicable to other robots, tasks, and simulators we describe these components first in order to ground the algorithm’s presentation.

Robot Platform: Our target task is bipedal walking using the SoftBank NAO.¹ The NAO is a humanoid robot with 25 degrees of freedom (See Figure 2a). For walking, our NAO uses an open source walk engine developed at the University of New South Wales (UNSW) (Ashar et al. 2015; Hall et al. 2016). This walk engine has been used by at least one team in each of the past three RoboCup Standard Platform League (SPL) championship games in which teams of five NAOs compete in soccer matches. To the best of our knowledge, it is one of the fastest open source walks available.² The walk is a zero moment point walk based on an inverted pendulum model. The walk is closed loop, using the inertial measurement (IMU) sensors and a learned ankle control policy to improve stability. The UNSW walk engine has 15 parameters that determine features of the walk (e.g., step height, pendulum model height). The values of the parameters from the open source release constitute θ_0 . For full information on the UNSW walk see (Hengst 2014).

Simulators: We use two different simulators in this work. The first is the open source SimSpark simulator.³ The simulator simulates realistic physics with the Open Dynamics Engine.⁴ This simulator was also used in the work introducing GSL (Farchy et al. 2013). We also use the Gazebo simulator⁵ provided by the Open Source Robotics Foundation.⁶

¹<https://www.ald.softbankrobotics.com/en>

²While other regular RoboCup participants also have competitive walks, we are not aware of any published result that confirms a speed as much as 42% faster than UNSW’s, as we achieve in this paper.

³<http://simspark.sourceforge.net>

⁴<http://www.ode.org/>

⁵<http://gazebo.org/>

⁶<http://www.osrfoundation.org/>

Gazebo is an open source simulator that is commonly used in robotics research. SimSpark enables fast simulation but is a lower fidelity model of the real world. Gazebo enables high fidelity simulation with an additional computational cost. In one of our experiments we consider the more realistic Gazebo environment as E and SimSpark as E_{sim} . The main difference between these simulators and the physical robot is how actions change the robot’s configuration. In SimSpark, actions achieve the desired command angle almost instantaneously. On the physical robot there is a delay. Gazebo also models joints as more reactive than the real world although less reactive than in SimSpark.

Grounded Action Transformation

We now introduce our principle algorithmic contribution, GSL with *Grounded Action Transformation* (GAT) which improves the grounding step (step 2) of the GSL framework. GAT improves grounding by correcting differences in the effects of actions between E and E_{sim} . GSL with GAT is presented in Algorithm 1.

The GSL framework grounds the simulator by finding ϕ^* that satisfies (1). Since it is often easier to minimize error in the one step dynamics distribution than error in the trajectory distributions, GAT uses:

$$\phi^* = \underset{\phi}{\operatorname{argmin}} \sum_{\tau_i \in \mathcal{D}} \sum_{t=0}^L d(P(\mathbf{s}_{t+1}^i | \mathbf{s}_t^i, \mathbf{a}_t^i), P_{\phi}(\mathbf{s}_{t+1}^i | \mathbf{s}_t^i, \mathbf{a}_t^i)) \quad (2)$$

as a surrogate loss function which can be minimized with transitions observed in \mathcal{D} . GSL with GAT begins by collecting the dataset \mathcal{D} (Algorithm 1, line 4).

Physics-based simulators (such as SimSpark and Gazebo) have a large number of parameters determining the physics of the simulated environment (e.g., friction coefficients). However, using these parameters as ϕ is not amenable to numerical optimization of (2). To find ϕ^* efficiently, GAT uses a parameterized *action transformation* function which takes the agent’s state and action as input and outputs a new action which – when taken in simulation – will result in the robot transitioning to the same next state it would have in E . We denote this function, $g_{\phi} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{A}$; the parameters of g serve as ϕ under the GSL framework. GAT constructs g with parameterized models of the robot’s dynamics and the simulator’s *inverse* dynamics. Assuming the simulated robot can record a dataset \mathcal{D}_{sim} of experience like the physical robot, GAT reduces (2) to a supervised learning problem.

GAT defines $g := g_{\phi}$ by a deterministic forward model of the robot’s dynamics, f and a deterministic model of the simulator’s inverse dynamics, f_{sim}^{-1} . The function, f maps $(\mathbf{s}_t, \mathbf{a}_t)$ to the maximum likelihood estimate of \mathbf{x}_{t+1} under P . The function f_{sim}^{-1} maps $(\mathbf{s}_t, \mathbf{s}_{t+1})$ to the action that is most likely to produce this transition in simulation. When executing θ in simulation, the robot selects $\mathbf{a}_t \sim \pi_{\theta}(\cdot | \mathbf{s}_t)$ and then uses f to predict what the resulting configuration, \mathbf{x}_{t+1} , would be in E . Then \mathbf{a}_t is replaced with $\hat{\mathbf{a}}_t := f_{\text{sim}}^{-1}(\mathbf{s}_t, f(\mathbf{s}_t, \mathbf{a}_t))$. The result is that the robot achieves the exact \mathbf{x}_{t+1} it would have on the physical robot.⁷

⁷GAT subsumes GUIDED GSL which makes the additional as-

Algorithm 1 Grounded Action Transformation (GAT) Pseudo code. Input: An initial policy, θ , the environment, E , a simulator, E_{sim} , smoothing parameter α , and a policy improvement method, `optimize`. The function `rolloutN(θ , N)` executes N trajectories with θ and returns the observed state transition data. The functions `trainForwardModel` and `trainInverseModel` estimate models of the forward and inverse dynamics respectively.

```

1: function GAT
2:    $\theta_0 \leftarrow \theta$ 
3:    $\mathcal{D}_{\text{robot}} \leftarrow \text{RolloutN}(E, \theta_0, N)$ 
4:    $\mathcal{D}_{\text{sim}} \leftarrow \text{RolloutN}(E_{\text{sim}}, \theta_0, N)$ 
5:    $f \leftarrow \text{trainForwardModel}(\mathcal{D}_{\text{robot}})$ 
6:    $f_{\text{sim}}^{-1} \leftarrow \text{trainInverseModel}(\mathcal{D}_{\text{sim}})$ 
7:    $g(s, a) \leftarrow \alpha f_{\text{sim}}^{-1}(s, f(s, a)) + (1 - \alpha) \cdot a_t$ 
8:    $\Pi \leftarrow \text{optimize}(E_{\text{sim}}, \theta, g)$ 
9:   return  $\text{argmin}_{\theta \in \Pi} J(\theta)$ 
10: end function

```

In practice f and f_{sim}^{-1} are represented with supervised regression models and learned from \mathcal{D} and \mathcal{D}_{sim} respectively (Algorithm 1 lines 5-6). The approximation of g introduces noise into the robot’s motion. To ensure stable motions, GAT uses a smoothing parameter α . The action transformation function (Algorithm 1 line 8) is then defined as:

$$g(s_t, \mathbf{a}_t) := \alpha f_{\text{sim}}^{-1}(s_t, f(s_t, \mathbf{a}_t)) + (1 - \alpha) \mathbf{a}_t$$

In our experiments, we set α as high as possible subject to the walk remaining stable. Figure 1 illustrates the GAT-modified E_{sim} . GAT then proceeds to improve θ with `optimize` within the grounded simulator (lines 8-9).

GAT Implementation

In this work, GAT uses two neural networks to approximate f and f_{sim}^{-1} . Each function is a three layer network with 200 hidden units in the first layer and 180 hidden units in the second. During simulator modification, the f network receives s_t and \mathbf{a}_t as input and the f_{sim}^{-1} network receives s_t and the output of f as input. The final output from f_{sim}^{-1} is the replacement action $\hat{\mathbf{a}}_t$. While \mathbf{a}_t is a vector of desired joint angles, the action input to f and action output of f_{sim}^{-1} is encoded as a desired change in \mathbf{x}_t which was found to improve prediction. We follow Fu et al. (2015) by having f predict the joint acceleration caused by executing \mathbf{a}_t in s_t instead of directly predicting s_{t+1} . The accelerations can then be integrated and added to s_{t+1} to produce the resulting next state. Additionally, f and f_{sim}^{-1} regress to the sine and cosine of the target angular accelerations and then pass the

assumption $f_{\text{sim}}^{-1}(s_t, \mathbf{x}_{t+1}) = \mathbf{x}_{t+1}$, in other words that requesting joint angles means the robot achieves those joint angles at the next time-step. GAT removes all human guidance from policy optimization. GAT also removes the need for human expertise in selecting the components of θ which are allowed to change during policy optimization.

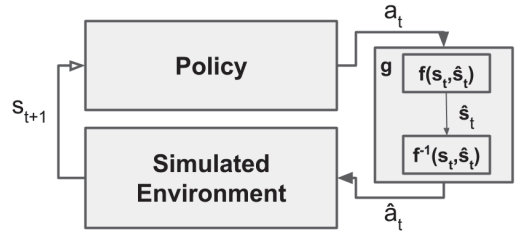


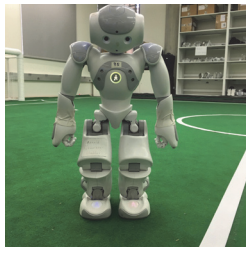
Figure 1: An illustration of the modifiable simulator GAT induces. At each time-step the robot takes an action, \mathbf{a}_t , and passes \mathbf{a}_t to a modification function, g . The modification function uses a deterministic model of the real robot’s dynamics, f , to predict the effect of executing \mathbf{a}_t on the physical robot. Then, a model of the simulated robot’s inverse dynamics uses the prediction, $\hat{\mathbf{x}}_t$, to predict the action $\hat{\mathbf{a}}_t$ which will achieve $\hat{\mathbf{x}}_t$ in simulation. Finally, $\hat{\mathbf{a}}_t$ is passed to the environment, E_{sim} and the resulting state transition will be similar to the transition that would have occurred in E .

outputs through the arctan function to produce the final angular acceleration. Passing the network outputs through the arctan function ensures f and f_{sim}^{-1} produce valid joint angles. The true state, s_t , is estimated by concatenating joint configurations $\mathbf{x}_t, \mathbf{x}_{t-1}, \dots, \mathbf{x}_{t-4}$ and past actions $\mathbf{a}_{t-1}, \dots, \mathbf{a}_{t-4}$. The state estimate $\langle \mathbf{x}_t, \dots, \mathbf{x}_{t-4}, \mathbf{a}_{t-1}, \dots, \mathbf{a}_{t-4} \rangle$ improves the predictions of f and f_{sim}^{-1} because it implicitly captures the unobserved $\dot{\mathbf{x}}_t$ and $\ddot{\mathbf{x}}_t$ state variables. The length of the configuration history was chosen to balance predictive accuracy with keeping the number of inputs to the networks small. Both networks are trained with backpropagation.

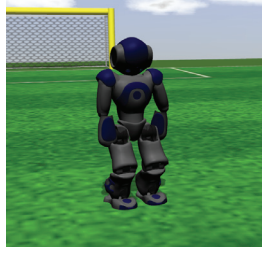
Empirical Results

Experimental Set-up

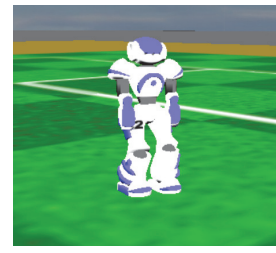
We evaluate GAT on the task of bipedal robot walking. The walk takes a target forward velocity parameter in the range $[0, 1]$. We set this parameter to 0.75 which we found led to the fastest walk that was reliably stable. The robot walks forward towards a target at this velocity. If the robot’s angle to the target becomes greater than five degrees it turns back towards the target while continuing to walk forward. In all environments, $J(\theta)$ is the negative of the average forward walk velocity while executing θ . On the physical robot a trajectory terminates once the robot has walked four meters or falls. A trajectory generated with θ_0 lasts ≈ 20.5 seconds on the robot. In simulation a trajectory terminates after a fixed time interval (7.5 seconds in SimSpark and 10 seconds in Gazebo) or when the robot falls. Previous work has shown that when using a model estimated from data it is better to use shorter action trajectories to avoid overfitting to an inaccurate model (Jiang et al. 2015). Even for identical s_t and \mathbf{a}_t , s_{t+1} in E_{sim} will most likely be different than s_{t+1} in E . This error compounds over the course of a trajectory since state error at s_t is propagated forward into s_{t+1} . This observation motivates using shorter trajectory lengths as simulator fidelity decreases.



(a) A Softbank NAO Robot



(b) A simulated NAO in Gazebo



(c) A simulated NAO in SimSpark

Figure 2: The three robotic environments used in this work. The Softbank NAO is our target physical robot. The NAO is simulated in the Gazebo and SimSpark simulators.

Optimizing θ We use the Covariance Matrix Adaption-Evolutionary Strategy (CMA-ES) algorithm (Hansen 2006) as the `optimize` routine. CMA-ES is a policy search method which samples a population of candidate θ from a Gaussian distribution over parameters. The top k parameter vectors are used to update the sampling distribution so that the mean is closer to an optimal policy. We modify $J_{\text{sim}}(\theta)$ for the optimization by adding a cost of 15 for any trajectory in which the robot falls. The added penalty encourages CMA-ES to strongly favor stable policies over faster, less stable ones. To clarify terminology, a generation refers to a single update of CMA-ES; an iteration refers to a complete cycle of GAT.

SimSpark to Gazebo: Since a large number of trials are difficult to obtain on a physical robot, we present a study of GAT using Gazebo as a surrogate for the real world. In this setting we evaluate the effectiveness of GAT compared to learning with no grounding and grounding E_{sim} by injecting noise into the robot’s actions. Adding an “envelope” of noise has been used before to minimize simulation bias by preventing the policy improvement algorithm from overfitting to the simulator’s dynamics (Jakobi, Husbands, and Harvey 1995). We refer to this baseline as Noise-Envelope. Since GAT with function approximation introduces noise into the robot’s actions we wish to verify that GAT offers benefits over such methods. Noise-Envelope adds standard Gaussian noise to the robot’s actions within the ungrounded simulator. We also attempted to evaluate GUIDED GSL but preliminary experiments showed that the assumption that actions achieve their desired effect instantaneously did not hold for this setting.

We run 10 trials of each method. For GAT we collect 50 trajectories of robot experience to train f and 50 trajectories of simulated experience to train f_{sim}^{-1} . For each method, we optimize θ for 10 generations of the CMA-ES algorithm. In each generation, 150 policies are sampled and evaluated in simulation. CMA-ES estimates J_{sim} with 20 trajectories from each policy. Overall, the CMA-ES optimization requires 30,000 simulated trajectories for each trial.

Table 1 gives the average improvement in stable walk policies for each method and the number of trials in which a method failed to produce a stable improvement. This table only considers policies found after the first generation of CMA-ES. The reason for this is that the policies in the first generation were randomly sampled from an initial search

distribution. We consider learning to begin once CMA-ES has used the J_{sim} values of the first generation to update the search distribution. Using J_{sim} to identify θ which improve J in the first generation is a *policy evaluation* problem while improvement afterwards is a *policy improvement* problem.

Simulation to NAO Set-up: We evaluate GAT for transferring policies learned in Simspark or Gazebo to the physical NAO. The data set \mathcal{D} consists of 15 trajectories collected with θ_0 on the physical NAO. For each iteration, we optimize θ for 10 generations of the CMA-ES algorithm. We evaluate the best policy from each generation with 5 trajectories on the physical robot. If the robot falls in any of the 5 trajectories the policy is considered unstable. While the number of evaluations is small, in practice stable policies had small variance in walk velocity and unstable policies fell in almost all trajectories.

One challenge in this setting is the simulators receive robot actions at 50 Hz while the physical NAO receives robot actions at 100 Hz. The discrepancy in action frequency poses a problem for using real world data to modify how joints move in simulation. By skipping every other measurement to get an effectively 50Hz data trace, we are able to model how the physical robot’s joints move at the simulator’s frequency.

Experimental Results

SimSpark to Gazebo Results: Table 1 shows that GAT maximizes improvement in J while minimizing iterations without improvement. NOISE-ENVELOPE improves upon no grounding in both improvement and number of iterations without improvement. Adding noise to the simulator encourages CMA-ES to propose robust policies which are more likely to be stable. However, GAT further improves over NOISE-ENVELOPE. This result demonstrates that action transformations are grounding the simulator in a more effective way than simply injecting noise.

Table 1 also shows that on average GAT finds an improved policy within the first few policy updates after grounding. When learning with no grounding finds an improvement it is also usually in an early generation of CMA-ES. As policy improvement progresses, the best policies in each generation begin to overfit to the dynamics of E_{sim} . Without grounding overfitting happens almost immediately. NOISE-ENVELOPE is shown to be more robust to overfitting since any policy it

Method	% Improve	Failures	Best Gen.
No Ground	11.094	7	1.33
Noise-Envelope	18.93	5	6.6
GAT	22.48	1	2.67

Table 1: This table compares Grounded Action Transformation (GAT) with baseline approaches for transferring learning between SimSpark and Gazebo. The first column displays the average maximum improvement found by each method after the first policy update made by CMA-ES. The second column is the number of times a method failed to find a stable walk. The third column gives the average generation of CMA-ES when the best policy was found. No Ground refers to learning done in the unmodified simulator.

proposes achieved good cost in a noisy E_{sim} . The grounding done by GAT is inherently local to the trajectory distribution of θ_0 . Thus as θ changes in policy improvement, the action transformation function fails to produce a more realistic simulator. Noise modification methods can mitigate overfitting by emphasizing robust policies although it is also limited in finding as strong of an improvement as GAT.

Simulator to Physical NAO Results: Our final empirical evaluation applies GAT to learning bipedal walks on a physical NAO. Table 2 gives the walk velocity of policies learned in simulation with GAT. The physical robot walks at a velocity of 19.52 cm/s with θ_0 . Two iterations of GAT with SimSpark increased the walk velocity of the NAO to 27.97 cm/s — an improvement of 43.27% compared to θ_0 .⁸ GAT with SimSpark and GAT with Gazebo both improved walk velocity by over 30%. This result demonstrates generality of our approach across different simulators.

As in the previous experiment, policy improvement with CMA-ES required 30,000 trajectories per iteration to find the 10 policies that were evaluated on the robot. In contrast the total number of trajectories executed on the physical robot is 65 (15 trajectories in \mathcal{D} and 5 evaluations per $\theta_c \in \Pi_c$). This result demonstrates GAT can use sample-intensive simulation learning to optimize real world skills with a low number of trajectories on the physical robot.

Farchy et al. demonstrated the benefits of re-grounding and further optimizing θ . We reground using the 15 trajectories collected with the best policy found by GAT with SimSpark and optimize for a further 10 generations in simulation. The second iteration results in a walk, θ_2 , which averages 27.97 cm/s for a total improvement of 43.27% over θ_0 .

Finally, we evaluate θ_0 with 100% of maximum velocity requested (i.e., the forward velocity request parameter set to 1.0). The average velocity of this walk across five stable trajectories is 24.3 cm / s. This result shows that GAT can learn walk policies that outperform one of the best hand coded walks available.

⁸A video of the learned walk policies is available at https://www.cs.utexas.edu/users/AustinVilla/?p=research/real_and_sim_walk_learning.

Method	Velocity (cm/s)	% Improve
θ_0	19.52	0.0
GAT SimSpark θ_1	26.27	34.58
GAT SimSpark θ_2	27.97	43.27
GAT Gazebo θ_1	26.89	37.76

Table 2: This table gives the maximum learned velocity and percent improvement for each method starting from θ_0 (top row).

Discussion and Future Work

Our proposed algorithm, GAT, has some limitations that we discuss here. The decision to learn an action modification function, g , makes the assumption that $\forall(s_t, a_t, s_{t+1})$ that could be observed on the physical robot there exists an action \hat{a}_t that will produce the same transition when used in place of a_t . We posit that this assumption is often true since \hat{a}_t can usually be executed with more or less force in simulation to achieve the desired response. However, this assumption is likely to break down under contact dynamics where external forces resist the robot’s actions. Other tasks may introduce other forms of simulator bias that GAT is currently limited in handling. An important direction for future work is to characterize the settings where this approach is limited and to identify alternatives.

We specify simulator grounding as a supervised learning problem however the distribution of inputs to the learned models, f and f_{sim}^{-1} , during policy execution differ from the training distributions. The distribution change is likely to weaken the quality of action modification under g . To see why the change happens consider that f is trained under (s_t, a_t, s_{t+1}) sampled from executing θ on the real robot while f_{sim}^{-1} is trained on (s_t, a_t, s_{t+1}) sampled in simulation. During simulator modification, both models are evaluated on data sampled from the distribution of θ in the grounded simulator. We have started to explore methods similar to the Dataset Aggregation (DAgger) algorithm (Ross and Bagnell 2012) which collects new data to retrain f_{sim}^{-1} as the state distribution of θ in E_{sim} changes with grounding.

Finally, this paper considers action modification but a complementary approach could consider sensor modification. Sensor modification would add another layer to the modified simulator such that the environment returns a state and the sensor modification function predicts what the robot would observe in that state. The two methods have an interesting interaction since sensor modification changes the action θ selects while action modification changes the next state which changes the observed sensor values as well.

Conclusion

This paper proposed and evaluated the Grounded Action Transformation (GAT) algorithm for grounded simulation learning. Our method led to a 43.27 % improvement in the walk velocity of a state-of-the-art bipedal robot. We conducted an empirical study that analyzed the benefits of GAT compared to a pair of baseline simulation-transfer methods. This experiment demonstrates GAT enhances learning in simulation in comparison to other methods.

Acknowledgments

We would like to thank Matthew Hausknecht, Patrick MacAlpine, and Garrett Warnell for insightful discussions and the anonymous reviewers for helpful comments. This work has taken place in the Learning Agents Research Group (LARG) at UT Austin. LARG research is supported in part by NSF (CNS-1330072, CNS-1305287, IIS-1637736, IIS-1651089), ONR (21C184-01), and AFOSR (FA9550-14-1-0087). Josiah Hanna is supported by an NSF Graduate Research Fellowship. Peter Stone serves on the Board of Directors of, Cogitai, Inc. The terms of this arrangement have been reviewed and approved by the University of Texas at Austin in accordance with its policy on objectivity in research.

References

- Ashar, J.; Ashmore, J.; Hall, B.; and Harris, S. e. a. 2015. Robocup spl 2014 champion team paper. In *RoboCup 2014: Robot World Cup XVIII*, volume 8992 of *Lecture Notes in Computer Science*. Springer International Publishing, 70–81.
- Christiano, P.; Shah, Z.; Mordatch, I.; Schneider, J.; Blackwell, T.; Tobin, J.; Abbeel, P.; and Zaremba, W. 2016. Transfer from simulation to real world through learning deep inverse dynamics model. *arXiv preprint arXiv:1610.03518*.
- Cully, A.; Clune, J.; Tarapore, D.; and Mouret, J.-B. 2015. Robots that can adapt like animals. *Nature*.
- Cutler, M., and How, J. P. 2015. Efficient reinforcement learning for robots using informative simulated priors. In *IEEE International Conference on Robotics and Automation, ICRA*.
- Cutler, M.; Walsh, T. J.; and How, J. P. 2014. Reinforcement learning with multi-fidelity simulators. In *IEEE Conference on Robotics and Automation, ICRA*.
- Deisenroth, M. P., and Rasmussen, C. E. 2011. Pilco: A model-based and data-efficient approach to policy search. In *International Conference on Machine Learning, ICML*.
- Farchy, A.; Barrett, S.; MacAlpine, P.; and Stone, P. 2013. Humanoid robots learning to walk faster: From the real world to simulation and back. In *Twelfth International Conference on Autonomous Agents and Multiagent Systems, AAMAS*.
- Fu, J.; Levine, S.; and Abbeel, P. 2015. One-shot learning of manipulation skills with online dynamics adaptation and neural network priors. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*.
- Hall, B.; Harris, S.; Hengst, B.; Liu, R.; Ng, K.; Pagnucco, M.; Pearson, L.; Sammut, C.; and Schmidt, P. 2016. Robocup spl 2015 champion team paper.
- Hansen, N. 2006. The CMA evolution strategy: a comparing review. In Lozano, J.; Larranaga, P.; Inza, I.; and Bengoetxea, E., eds., *Towards a new evolutionary computation. Advances on estimation of distribution algorithms*. Springer. 75–102.
- Hengst, B. 2014. runswift walk2014 report robocup standard platform league. Technical report, The University of New South Wales.
- Iocchi, L.; Libera, F. D.; and Menegatti, E. 2007. Learning humanoid soccer actions interleaving simulated and real data. In *Second Workshop on Humanoid Soccer Robots*.
- Jakobi, N.; Husbands, P.; and Harvey, I. 1995. Noise and the reality gap: The use of simulation in evolutionary robotics. In *European Conference on Artificial Life*, 704–720. Springer.
- Jiang, N.; Kulesza, A.; Singh, S.; and Lewis, R. 2015. The dependence of effective planning horizon on model accuracy. In *International Conference on Autonomous Agents and Multiagent Systems, AAMAS*.
- Kober, J.; Bagnell, J. A.; and Peters, J. 2013. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*.
- Koos, S.; Mouret, J.-B.; and Doncieux, S. 2010. Crossing the reality gap in evolutionary robotics by promoting transferable controllers. In *Proceedings of the 12th annual conference on Genetic and evolutionary computation*, 119–126. ACM.
- Lillicrap, T. P.; Hunt, J. J.; Pritzel, A.; Heess, N.; Erez, T.; Tassa, Y.; Silver, D.; and Wierstra, D. 2015. Continuous control with deep reinforcement learning. *CoRR abs/1509.02971*.
- Rajeswaran, A.; Ghotra, S.; Levine, S.; and Ravindran, B. 2016. Epopt: Learning robust neural network policies using model ensembles. *arXiv preprint arXiv:1610.01283*.
- Ross, S., and Bagnell, J. A. 2012. Agnostic system identification for model-based reinforcement learning. In *29th International Conference on Machine Learning, ICML*.
- Rusu, A. A.; Rabinowitz, N. C.; Desjardins, G.; Soyer, H.; Kirkpatrick, J.; Kavukcuoglu, K.; Pascanu, R.; and Hadsell, R. 2016a. Progressive neural networks. *arXiv preprint arXiv:1606.04671*.
- Rusu, A. A.; Vecerik, M.; Rothörl, T.; Heess, N.; Pascanu, R.; and Hadsell, R. 2016b. Sim-to-real robot learning from pixels with progressive nets. *arXiv preprint arXiv:1610.04286*.
- Schulman, J.; Moritz, P.; Levine, S.; Jordan, M. I.; and Abbeel, P. 2015. High-dimensional continuous control using generalized advantage estimation. In *International Conference on Learning Representations*.
- Sutton, R. S., and Barto, A. G. 1998. *Reinforcement Learning: An Introduction*. MIT Press.