Exploring the transition from novice to expert in RL policies for motor skill

Oussama Gabouj oussama.gabouj@epfl.ch Ahmed Aziz Ben Haj Hmida ahmed.benhajhmida@epfl.ch Salim Boussofara salim.boussofara@epfl.ch

Abstract

1 This study investigates the use of reinforcement learning for motor control, building on the curriculum-based RL approach introduced by Chiappa et al [1]. It explores 2 the evolution from novice to expert reinforcement learning policies in motor skill 3 acquisition, focusing on the manipulation of Baoding balls, which is a task requiring 4 complex motor skills. We employ recurrent PPO [2] with LSTM layers to handle 5 partial observability. Our study, also, simplifies the learning process by distilling 6 complex expert strategies into novice agents using various dimensionality reduction 7 techniques involving both static and dynamic reductions of observation and action 8 spaces with the aim to find a balance that maintains essential information while 9 enhancing learning efficiency. This report provides an overview of the methods 10 and insights gained into the effectiveness of dimensionality reduction in training 11 RL agents for intricate motor control tasks. 12

13 1 Introduction

14 **1.1 Context and motivation**

Understanding biological motor control is a major problem facing neuroscience today. The complex coordination of muscles required for various tasks ranging from daily activities to athletic achievements demonstrates the remarkable capabilities of biological systems. Using computer-based tools like musculoskeletal simulators and RL algorithms can help us learn about these processes and create better artificial motor control systems.

20 1.2 Background

Our project build on the foundational work achieved by Chiappa et al. [1]. They presented a new way 21 of using curriculum-based RL for motor control. Their method, the Static to Dynamic Stabilization 22 (SDS) curriculum, won the NeurIPS MyoChallenge for its effectiveness in training a model to 23 manipulate Baoding balls, a task requiring complex motor skills. This innovative approach mirrors 24 human learning by progressively teaching an RL agent to stabilize static configurations before moving 25 on to dynamic transitions, thereby enhancing learning efficiency and performance. The study by 26 Chiappa et al. [1] highlights the potential of combining physiologically-detailed simulators with 27 28 powerful RL algorithms to tackle complex motor control challenges.

29 An important aspect of their work is to verify the hypothesis that it is possible to extract synergies

from artificial agents as in biological muscles via dimensionality reduction in motor control. Indeed, The human motor system has numerous degrees of freedom, making the control problem highly

Submitted to 38th Conference on Neural Information Processing Systems (NeurIPS 2024). Do not distribute.

complex. The SDS curriculum implicitly reduces this complexity by breaking down the learning
 process into manageable stages. However, understanding the intrinsic dimensionality of motor control
 tasks can further enhance learning efficiency. Identifying the key components or synergies that govern
 effective motor control could help to develop more efficient RL algorithms that operate in a reduced

³⁶ dimensional space, leading to faster learning and better generalization.

37 2 Experimental environment

The environment used for the experiments is the *Myosuite*, specifically the baoding balls task which a challenging motor control problem, divided into two distinct phases:

• **Phase I:** It focuses on counter-clockwise rotations with fixed task parameters.

• **Phase II:** It introduces additional complexities such as clockwise rotations, hold conditions, and random variations in task parameters like rotation period, ball size, and friction.

The agent receives a comprehensive set of observations, described by an 86-dimensional observation 43 vector, which provides a detailed representation of the system's state. This vector includes 23 joint 44 angles, positions and velocities of the balls, positions and distances of the targets and previous 45 activations of the 39 muscles. These observations enable the agent to understand the current state 46 and make informed decisions about subsequent actions. Indeed in response to the observations, the 47 agent generates actions within a 39-dimensional action space, corresponding to the activations of 48 the muscle-tendon units controlling the human forearm model. These actions are influenced by the 49 internal state representations formed by the LSTM layers, which accumulate information over time. 50

⁵¹ Therefore, the interaction between the control policy and the MuJoCo physics simulator is formulated

as a Partially-Observable Markov Decision Process, represented as $M = \langle S, A, O, T, R, \gamma \rangle$. This process comprises the following components:

- **State Space** (*S*): The complete set of possible states of the system.
- Observation Function (*O*): Maps the state to an 86-dimensional observation vector (*O* : $S \to \mathbb{R}^{86}$).
- Action Space (A): Represents the possible muscle activations ($A \subset \mathbb{R}^{39}$).
- **Transition Function** (*T*): Defines how the environment evolves ($T : S \times A \rightarrow S$).
- **Reward Function** (*R*): Associates rewards with state transitions $(R : S \times A \times S \rightarrow \mathbb{R})$.
- **Discount Factor** (γ): Balances immediate and future rewards.

61 **3** Methodology

The primary task of our project is to analyze the components of motor synergies and understand how variance behaves with the aim to find ways to distill expert motor strategies and effectively transfer them to novice agents to achieve comparable performance in the second phase. The reinforcement algorithm used to train the novice agent is the recurrent proximal policy optimization which balances exploration and exploitation effectively. The PPO configuration includes a recurrent neural network architecture with LSTM layer to handle the partial observability of the environment.

To analyze the motor synergies of the novice and the expert agents, we conducted a series of structured
 experiments.

70 3.1 Train an agent on phase I using the architecture used for phase II

To ensure consistency, we began by training an agent on Phase I of the Baoding balls task using the architecture designed for Phase II. The training process starts with the agent having no prior knowledge, learning solely from rewards received during interactions. The agent trains in the Phase I environment, focusing on counter-clockwise rotations of the Baoding balls over multiple episodes, each lasting 200 time steps (5 seconds). Early termination is used if the balls fall below the palm to focus learning on productive interactions. The agent's performance is based on rewards for maintaining the balls' proximity to the target trajectory, with cumulative rewards updating the policy through PPO. Using the advanced Phase II architecture in Phase I training is expected to improve learning efficiency and performance, providing a baseline for further analysis.

80 3.2 Extract the principal components of the expert agent (Phase II)

Next, we extracted the PCs from the expert agent trained in Phase II to identify key components
driving effective performance and understand underlying motor synergies. This step is foundational
for the subsequent distillation processes, as it highlights the most significant features and actions that
contribute to successful task execution. PCA is performed on the expert agent's features and actions
from Phase II of the Baoding balls task.

- Feature Space: Extract PCs from the last hidden layer and map into a smaller feature space.
- Action Space: Extract PCs from the action space (i.e muscles activation's space) and map into a smaller action space.

Using an analysis of the cumulative explained variance, which we will discuss later in the report, we decomposed the high-dimensional feature and action spaces of the expert agent into a small set of orthogonal components capturing the most of the variance in the agent's behavior. This allows us to identify the key factors of the expert's performance, representing motor synergies [3] and efficient strategies that can be transferred to novice agents. This understanding helps in reducing the dimensionality of observation and action spaces, combining hand movements and simplifying the learning process for novice agents.

96 **3.3** Policy distillation strategies

To explore the transition from novice to expert in RL policies for motor skill acquisition, we employ two primary strategies for policy distillation: reducing the observation mapping space and reducing the action space dimensionality.

100 3.3.1 Feature space dimensionality reduction

In the first part, we focus on reducing the observation mapping space to guide the agent in taking the
 best possible actions. By reducing the dimensionality of the observation space, we aim to constrain
 the agent's exploration, helping it to learn more efficiently. This approach addresses the curse of
 dimensionality by simplifying the complex observation mapping space.

Static PCA: We project the features of the novice agent into a lower-dimensional space using the principal components (PCs) of the observation mapping space derived from the expert agent as illustrated in the Figure 1. The PCA layer as indicated by red box is a frozen layer. By reducing the feature space to these key components, we create a simplified representation that the novice agent can use to learn the task. The agent is, then, trained to take actions given the reduced observation mapping space, effectively focusing its learning on the most relevant aspects of the task as determined by the expert agent's experience. The reduced feature space helps limit the exploration of the novice

agent to the most promising regions.



Figure 1: Static PCA architecture

Bottleneck: Instead of using a static PCA, we employ an additional fully connected layer to the policy 113 neural network to dynamically learn the best feature mapping before action selection. The Figure 2 114 illustrates the architecture of the modified network. It involves adding a bottleneck layer, just before 115 the output layer, which forces the network to compress the information into a lower-dimensional 116 representation. It should allow the network to adaptively determine the most important features, 117 potentially capturing more nuanced relationships than static PCA. This bottleneck architecture serves 118 as a baseline for comparison. It allows us to evaluate the effectiveness of PCA-reduced observation 119 space from the expert against a dynamically learned reduced observation mapping space. 120



Figure 2: Bottleneck architecture

121 **3.3.2** Action space dimensionality reduction

The second approach focuses on directly reducing the action space dimensionality. We aim to constrain the agent's exploration so that it only explores the most probable and relevant action space derived from the expert agent's experience. By reducing the dimensionality of the action space, we add constraints to guide the agent's exploration towards the best possible actions, captured by the PC of the expert.

Projection and Back-Projection: In this strategy, actions are projected into a reduced space formed by the PCs that were extracted from the expert. The obtained vectors are, then, projected back to the original high-dimensional action space. This process, illustrated in Figure 3, involves splitting the projection and back-projection phases to reduce small variations and noise. The idea is that by reconstructing from the projected space, the basic actions remain consistent, but the intensity of the activated muscles is smoothed, removing small variations. This smoothing effect helps the agent to capture the overall behavior more effectively.



Figure 3: Projection and Back-Projection architecture

Reduced Latent Space Exploration: This approach enforces exploration constraints by projecting
 actions into a reduced latent space, conducting exploration within this space, and then projecting back
 to the original action space. The key difference from the previous method is the focus on latent space,
 which represents a more abstract and potentially more informative reduced space for exploration.
 By operating within this latent space, the agent can explore variations around the high-variance

139 components identified from the expert's policy. This method ensures that the agent's exploration

is primarily within the most critical regions of the action space, promoting efficient learning and potentially faster convergence to an optimal policy. This approach is illustrated in Figure 4, where

- potentially faster convergence to an optimal policy. This approach is illustrated in Figure 4, where the exploration is conducted after PCA and then projected back to ensure the agent focuses on the
- most informative components.



Figure 4: Latent Space Exploration architecture

143

144 **4 Results Analysis**

145 4.1 PCA Analysis

After analyzing the figures showcasing the cumulative explained variance by principal components 146 in both the action space and feature space, we can conclude that a suitable number of principal 147 components (PCs) for each domain can be determined by considering their cumulative contribution 148 to the variance. In the action space, a suitable number of PCs to be extracted is 16, as their cumulated 149 contribution to the variance is greater than 90%. For the feature space, a suitable number of PCs to 150 be extracted is 40, accounting for a substantial 88% of the variance. We can note that achieving a 151 cumulative variance of over 0.9 requires adding at least 10 components, covering only 0.02 of the 152 variance. Keeping only 40 PCs seems, therefore, to be a good tradeoff. 153

By selecting these 16 PCs for the action space and 40 PCs for the feature space, we effectively capture the most significant features, enabling proper dimensionality reduction while preserving the essential characteristics of both spaces for our analysis.



Figure 5: Cumulative explained variance of the PCs

157 4.2 Experiments

All four models were evaluated under stochastic conditions. The mean reward indicates the average reward that a model achieved throughout training, and it varies within the range specified by the

standard deviation STD Reward. Similarly, the mean episode length, which measures the duration 160

before failure (e.g., balls dropping) also fluctuates within the range defined by its standard deviation 161

STD Episode Length. 162

Model	Exp 1	Exp 2	Exp 3	Exp4	Expert 2 Baseline				
Mean Reward	485	321	41	39	990.4				
STD Reward	7.3	7.8	2	2.1	-				
Mean Episode Length	121.4	111.3	22.5	21.5	200				
STD Episode Length	3.5	4.2	1.5	0.9	-				
Table 1: Summary of Neural Network Experiment Results									

able	1:	Summary	of Neura	al Netwo	ork Exp	erimer	t Results
aore	. .	ounnur,	01 1 10 011	AI I 100011	orn Dap	CI IIIICI.	it recourto

Feature space dimensionality reduction 4.2.1 163

Experiment 1: It involved integrating a PCA layer into the original neural network architecture 164 to reduce the feature space from 256 dimensions 40 which is the number of principal components 165 (PCs), and then mapping these PCs to the final 39 features. This model achieved a mean reward 166 of 485 with a standard deviation of 7.3 indicates that this model performed robustly under varied 167 conditions, consistently earning high rewards. The relatively stable mean episode length of 121.4 168 with a standard deviation of 3.5 suggests that the model was effective in sustaining performance over 169 time before failure occurred. The integration of a PCA layer demonstrated a substantial improvement 170 in performance compared to the more complex models in subsequent experiments. This suggests that 171 feature reduction at this level of dimensionality effectively balances complexity and performance, 172 offering a sweet spot for this specific task. 173

Experiment 2: Similar to the first experiment, but replacing the static PCA layer with a fully 174 connected layer that dynamically learns the best feature mapping for action selection. The mean 175 reward decreased to 321 with a standard deviation of 7.8, and the mean episode length also reduced to 176 111.3 with a standard deviation of 4.2. These changes suggest that, while the network was adapting, it 177 might have faced challenges in stabilizing its performance, possibly due to the increased complexity 178 or overfitting to the training, as the network adjusts the feature representation based on feedback from 179 180 the environment, potentially capturing more complex patterns beneficial for the task.

4.2.2 Action space dimensionality reduction 181

Experiment 3: This retained the original neural network structure but added a PCA layer that 182 projected the 39 output features into a smaller feature space of only 16 dimensions, previously defined 183 by the expert, before projecting them back to the 39 output features. The results showed a significant 184 drop in performance, with a mean reward of 41 and a standard deviation of 2. The episode length was 185 also much shorter at 22.5 with a standard deviation of 1.5, indicating a quicker failure rate and less 186 robustness in maintaining performance. The drastic reduction in feature space seems to have been 187 too severe, potentially omitting necessary information for making effective decisions. This suggests 188 that there is a critical threshold below which further reduction in dimensionality adversely affects the 189 model's capability to function effectively. 190

Experiment 4: This setup is built on the third experiment by adding exploration of the reduced 191 feature space before mapping back to the original action space. The slight decrease in mean reward 192 to 39 and an episode length of 21.5, both with small standard deviations of 2.1 and 0.9, indicates a 193 continued struggle in leveraging the reduced feature space effectively, even with the added exploration 194 phase. The results suggest that simply spending more time within this space is not sufficient to 195 compensate for the loss of critical information due to excessive compression, although the exploration 196 197 of the reduced feature space was a logical step.

4.2.3 Evaluation: 198

The first two experiments demonstrated significantly higher rewards and episode lengths compared 199 to the last two. The higher performance suggests that, in these initial experiments, the models 200 were at least able to maintain control over the task (e.g., holding the balls) for longer periods, 201 which is the initial step of the curriculum learning SDS according to the project paper, although 202 they struggled with more complex manipulations like correctly rotating the balls, as indicated in 203

the skeletal simulations. The significantly lower rewards and episode durations in the latter two experiments suggest difficulties in basic task retention, such as holding the balls, which was visually confirmed through the simulations.

207 5 Discussion

208 5.1 Evaluating Task-Specific Performance Through PCA

The series of experiments underscores the critical role of selecting an appropriate number of principal components (PCs) in feature reduction strategies. This choice must align with the complexity of the tasks to ensure effective performance. In these experiments, the reconstruction from the reduced feature space maintains consistent basic actions but smooths out the intensity of muscle activation, eliminating minor variations. This smoothing occurs because the agent struggles to fully execute all its tasks within the constricted feature space

For instance, while the agent learns to hold the balls, a task that does not require fine variations, it fails to effectively rotate the balls, which demands maximal muscle extension. This limitation in the model's learning process induced by the smoothing prevents it from extending the muscles sufficiently to perform complex tasks like ball rotation. Consequently, while this effect allows the agent to capture overall behavior more efficiently, it overlooks crucial, nuanced movements essential for complete task execution.

Given that Experiment 1 demonstrated a capability to achieve a reward of 400 just by holding the balls, we suggest a potential to sequentially train the model on more complex tasks such as ball rotation and then boarding. Since training to hold the balls was accomplished in just 4 hours, it is projected that training the model to perform all three tasks, holding, rotating, and boarding, could be completed in less than 24 hours. This is a substantial reduction from the 400 hours initially cited in the literature, indicating a significant efficiency improvement in the training process.

227 5.2 Optimizing Dimensionality Reduction Strategies for Enhanced Model Training

The strategic selection of 16 principal components (PCs) for the action space and 40 PCs for the 228 feature space of the last hidden layer was intended to maintain explained variances above 0.9 and 0.88 229 respectively. However, the model's persistent struggles, particularly with tasks involving ball rotation, 230 suggest that the lower variations excluded from the model (those accounting for the remaining 0.1 231 232 and 0.12 of variance) are indeed critical for complete training success. Since the explained variance is highly task-specific and our model manages intricate details like training 39 distinct muscles, it 233 appears necessary to include more variance by increasing the explained variance percentage for both 234 action and feature spaces. 235

In light of this, discussions with the supervisor have showed that it was recommended to adjust the number of components for the action space to 24 instead of 16. This adjustment aims to preserve more of the essential variations needed for comprehensive model training.

Moreover, the use of Principal Component Analysis (PCA) itself may need reconsideration. While 239 PCA is effective for reducing dimensionality by capturing maximum variance within fewer dimen-240 sions, it might not be the most suitable method for tasks that involve complex, dependent actions like 241 the ones our model is trained on. Alternative dimensionality reduction techniques might be more 242 appropriate like Independent Component Analysis (ICA) which, unlike PCA, that prioritizes direc-243 244 tions that maximize variance, focuses on finding components that are statistically independent from each other. This could be particularly advantageous in tasks where independent features contribute 245 uniquely to performance outcomes. 246

Additionally, experimenting with the positioning of the PCA layer within the neural network architecture might yield improvements. Shifting the PCA layer to different points between other layers could potentially optimize the variance explained, ensuring that essential information is not lost during dimension reduction. This approach may validate whether an explained variance of 0.9, for example, is sufficient to maintain the necessary components for the model to learn effectively without significant loss of crucial details.

253 6 Conclusion

The experiments conducted in this study reveal significant insights into the efficiency of dimensionality 254 reduction and policy distillation strategies for training RL agents in motor control tasks. While 255 reducing the observation and action spaces aids in simplifying the learning environment, it is crucial 256 to balance the extent of reduction to avoid losing critical information necessary for task execution. 257 The results indicate that a careful selection of the number of principal components and strategically 258 placing the projection layer within the neural network architecture are keys to maintaining an effective 259 learning process. Moreover, transitioning from novice to expert performance can be optimized by 260 adjusting the dimensional reduction techniques to enhance learning outcomes. Future work should 261 explore alternative dimensionality reduction techniques and different configurations of the learning 262 architecture to further improve the efficiency and effectiveness of training RL agents in complex 263 motor tasks. 264

265 **References**

- [1] Nisheet Patel Abigaïl Ingster Alexandre Pouget Alexander Mathis Alberto Silvio Chiappa,
 Pablo Tano. Acquiring musculoskeletal skills with curriculum-based reinforcement learning.
 2024.
- [2] Frank Zimmer Mike Preuss Marco Pleines, Matthias Pallasch. Generalization, mayhems and
 limits in recurrent proximal policy optimization. 2022.
- [3] Zhang Z Atzori M Müller H Xie Z Scano A Zhao K, Wen H. Evaluation of methods for the extraction of spatial muscle synergies reinforcement learning. 2022.