# Learning Fast and Precise Pixel-to-Torque Control

Steffen Bleher[1], Steve Heim[1,2], Sebastian Trimpe[1,3]

[1] Intelligent Control Systems Group, Max Planck Institute for Intelligent Systems, Stuttgart, Germany;
[2] Biomimetic Robotics Lab, Massachusetts Institute of Technology, Cambridge, USA;
[3] Institute for Data Science in Mechanical Engineering, RWTH Aachen University, Aachen, Germany

## I. INTRODUCTION

In the field, robots often need to operate in unknown and unstructured environments, where accurate sensing and state estimation (SE) becomes a major challenge. Cameras have been used to great success in mapping and planning in such environments [1], as well as complex but quasi-static tasks such as grasping [2], but are rarely integrated into the control loop for unstable systems. Learning pixel-to-torque control promises to allow robots to flexibly handle a wider variety of tasks. Although they do not present additional theoretical obstacles, learning pixel-to-torque control for unstable systems that that require precise and high bandwidth control still poses a significant practical challenge, and best practices have not yet been established. Part of the reason is that many of the most auspicious tools, such as deep neural networks (DNN), are opaque: the cause for success on one system is difficult to interpret and generalize.

The machine learning community has alleviated this problem by establishing standard data sets and standardized simulation environments that allow different approaches to be easily benchmarked against each other. This trend is not well established in the robotics community, as there are many more hurdles to reproduce a system in hardware than purely in simulation. To help drive reproducible research on the practical aspects of learning pixel-to-torque control, we propose a platform that can flexibly represent the entire process, from lab to deployment, for learning pixel-to-torque control on a robot with fast, unstable dynamics: the vision-based Furuta pendulum. The platform, shown in Figure 1 and detailed in "Reproducible Platform", can be reproduced with either off-the-shelf or custom-built hardware. We expect that this platform will allow researchers to quickly and systematically test different approaches, as well as reproduce and benchmark case studies from other labs.

We also present a first case study on this system using DNNs which, to the best of our knowledge, is the first demonstration of learning pixel-to-torque control on an unstable system with update rates faster than $100\,\mathrm{Hz}$. A video synopsis can be found online at https://youtu.be/S2llScfG-8E, and in the supplementary material.

## II. RELATED WORK

DNNs combined with RL have had tremendous success recently in a variety of robotics applications [2]–[5], though many open challenges remain [6]. One of the most critical challenges in all these approaches is sample efficiency—or



Fig. 1: The vision-based Furuta pendulum: a platform for reproducible research on learning fast and precise pixel-to-torque control.

rather, sample inefficiency. In most cases, learning is done in simulation only, which adds the challenge of transferring the learned policy from simulation to the actual hardware. To make this transfer successfully, a lot of effort is typically put into modeling and system identification [4], [7], such that the gap between simulation and reality is 'small' in some sense. For certain dynamics, such as turbulent flows or soft matter [8], [9], accurate models are unavailable or prohibitive to simulate. Overcoming the sample-efficiency challenge would not only allow learning to be leveraged on these systems, but also alleviate the reality gap in general: policies can be refined on the real hardware after initial training in a low-fidelity simulation. Indeed, Lee *et al.* [3] point out that, while they needed a high-fidelity model of the robot dynamics to train a DNN policy, they could successfully make the transfer from simulation to reality using only low-fidelity terrain models.

One of the key concepts they leverage is *privileged information*: ground-truth data is often available during training, and can be leveraged to substantially reduce training time. Chen *et al.* [10] coined this term, and use it to improve imitation learning by first training an autonomous driving policy with a

birds-eye view of the environment, then using its evaluations as training examples for the final policy, which only has access to a regular car-mounted camera as input. Lee *et al.* [3] use the same concept to infer terrain properties from a history of proprioceptive data and thus avoid the need for external sensing entirely. In both these studies, learning is done in simulation, and privileged information can be directly accessed from the simulator. Privileged information is also available when learning directly in hardware, especially if it takes place in a controlled lab setting. For example, Srinivasan *et al.* [11] learn accurate SE for a racing car from only IMU and wheel encoder readings. Training targets are generated with a mixed Kalman filter that has access to two additional velocity sensors, which are very accurate but also expensive. Previously, Levine *et al.* [12] used this concept to learn to estimate the position of a target object from images, using supervised learning. During this phase, the object is placed in the robot's gripper, so the robot can directly estimate its position relative to the camera through joint-position measurements and forward kinematics.

Levine *et al.* [12] also report significantly better performance with full end-to-end learning, that is, training a single DNN for both the state estimator and controller. However, separating SE and control also has benefits, such as improved sample efficiency and more targeted development. For example, Srinivasan *et al.* [11] rely on existing methods for perception, mapping, and control, and focus on learning a convolutional DNN specifically to estimate velocities, which can be challenging for a Kalman filter during aggressive maneuvers. Hoeller *et al.* [13] implement a full, highly modular learning pipeline, which separately tackles state-representation[1] and motion planning. This pipeline is trained to high performance with remarkable sample efficiency, requiring on the order of 70'000 depth-images and 17 hours worth of trajectories, using a mixture of simulated and real-world data.

Despite recent successes in learning vision-based controllers for grasping [2], [14], [15], learning pixel-to-torque control remains elusive, especially for fast, unstable systems. An important difference is that for tasks such as grasping, a conventional low-level controller can be relied on to stabilize the dynamics. Instead, the challenge is to generate appropriate desired kinematics such as grasping positions [2], [15], or kinematic trajectories, often called primitives [14], [15]. In other words, learning is used for *planning*, rather than for *control*.

Since torque control is usually required when systems have fast and unstable dynamics, high control bandwidth is often a concern when learning pixel-to-torque control. This need for fast and precise feedback is a key characteristic of the proposed platform, which distinguishes it from more common platforms for research on vision-based learning.

Lambert *et al.* [16] use DNNs to learn a predictive low-level controller of a hovering quadcopter using onboard sensors as input, and manage to obtain stable hovering for several seconds at a time after training on only a few minutes of data. Despite running a relatively simple DNN architecture on a powerful,

---

[1] State representation only differs from state estimation in that it includes relevant state of the environment, such as moving obstacles.

offboard GPU, evaluation time is the bottleneck: to obtain sufficiently long prediction horizons requires multiple evaluations of the DNN, which limits their control bandwidth to 50 Hz. This bottleneck is greatly exasperated when vision is used for feedback since the high dimensions and complexity of vision typically require larger and more sophisticated DNNs, which are consequently slower to evaluate. For example, Mattner *et al.* [17] use an auto-encoder architecture to learn pixel-to-torque control for balancing an inverted pendulum with minimal domain knowledge. The entire DNN size is kept manageable in a number of ways, including down-sampling the input image to 40×30 pixels and limiting the output torque to only three values. Nonetheless, the control bandwidth is limited to roughly 10 Hz. Fast evaluation becomes even more challenging for autonomous robots that rely on onboard computation, since more powerful computers add strain to a limited payload and battery supply. To run vision-based control fully onboard a small quadcopter, Kaufmann *et al.* [18] only learn a DNN for part of the control pipeline, which runs at a lower frequency of 10 Hz. Similar to the grasping examples above, stability is maintained by a conventional controller running at a higher frequency, in this case a minimum-jerk planner. To learn the full pixel-to-torque control pipeline of fast systems, the trade-off between the learned SE's precision and its evaluation speed takes a central role in designing the learning pipeline, as we will explore in our case study.

## III. LEARNING PIPELINE

We found the two central criteria for designing the learning pipeline to be sample efficiency, and simultaneously fulfilling the precision and control bandwidth requirement. For the Furuta pendulum used in this study, the minimum control frequency translates to a time budget for the entire vision-based control of roughly 8 ms (see "Reproducible Platform"). We found it essential to split up the learning pipeline into four steps (see Figure 2): in step *A*, online RL of a control policy using privileged information, in step *B*, policy analysis and sample-collection using privileged information, in step *C*, offline supervised learning of the SE, and finally in step *D*, online adaptation learning of the control policy to the SE.

We rely on privileged information in the form of rich and accurate state measurements, which are often available in the lab setting via external sensing such as motion capture. We also rely on a means of automating sample collection with a specified distribution. In the case of the Furuta pendulum, state measurements are readily available from the joint encoders, and samples can be easily gathered using a standard combination of energy-pumping and LQR controllers. An important benefit of automating sample collection is that it makes it possible to quickly and easily collect new data sets. As is always the case, development is an iterative process, and speeding up this process is critical yet seldomly discussed in literature [7].

### A. Learning the Control Policy

To focus on a reliable and sample efficient training process, we train a Proximal Policy Optimization (PPO) [24] RL agent

## Reproducible Platform

To foster research on learning pixel-to-torque control, an ideal platform should represent the entire learning pipeline while also allowing different challenges to be addressed independently. It should be simple to reproduce, allow fast and easy iterations, and clear benchmarking. In order to push the boundary of learning and vision-based control in highly dynamic settings, we propose to combine the above requirements in a representative setup with high demands for precision and control bandwidth of larger than 100 Hz.

### The Vision-based Furuta pendulum

The Furuta pendulum is a well-studied system that combines important challenges, being nonlinear, unstable, and underactuated [19]. A minimal state space can be described as $\boldsymbol{x} = [\theta, \alpha, \dot{\theta}, \dot{\alpha}]^T$, where $\theta$ refers to the arm angle and $\alpha$ to the pendulum angle. The pendulum has a single control input $u$: the voltage applied to a DC-motor actuating the arm joint $\theta$. Access to its built-in encoders serves as a simple proxy for a controlled lab environment with external sensing, without the need for additional expensive equipment such as a motion capture setup.

Well-established controllers, such as energy-based swing-up [20] and LQR controllers [21], provide a reliable and rigorous baseline to compare against. Because the required control frequency directly depends on the pendulum's natural frequency, it can be easily increased or decreased by simple modifications to the pendulum length. Any desktop-sized Furuta pendulum can thus be easily modified to have comparable dynamics to the one used in this study, including platforms commonly used in classrooms [22].

### Hardware Details

We use the off-the-shelf Quanser Qube Servo 2 Furuta pendulum; as mentioned above, other Furuta pendulums with similar dynamics can be used instead. For our setup, a minimum sample frequency of 80 Hz is required to stabilize the pole with an LQR controller. From experiments, we found a sampling frequency of 120 Hz to be a good compromise and use this frequency for the case study presented.

In our experience, the learning pipeline is not very sensitive to changes in system dynamics, but more sensitive to changes in the camera and lighting setup. We use a FLIR BlackFly 3 high-speed camera with a resolution of 0.3 MP and a sample rate of 522 Hz. We found a high frame rate to be important to avoid additional latency in the control cycle. To reduce the effect of ambient light, we mounted the camera and pendulum in a white box. We also use a dimmable light source with a maximal illuminance of 4600 lx/0.5 m, which can both provide a steady illumination, and also simulate transfer to other lighting conditions in a controlled manner.

We ran experiments on a standard desktop computer with an Intel Xeon W-2123 CPU and an Nvidia Quadro P620 GPU, but also used a more powerful desktop to accelerate offline training (reported training time corresponds to the described computer). The hardware interface is based on code from [23] and contains a flexible setup based on OpenAI Gym and PyTorch, facilitating the reproducibility and reusability of all components in the proposed learning pipeline.

Our entire setup, including off-the-shelf pendulum, camera, and desktop PC, was purchased for a total cost of less than 10'000 €.

### Baseline Controllers & Automated Data Collection

For the swing-up baseline, we use the energy-pumping controller described in [20]. For balancing, we use a standard LQR controller.

To automate data generation, these baseline controllers are modified with reference trajectories for the pendulum arm and a sinusoidal input signal on the input to ensure sufficient coverage of the state space. This is critical because, although the dynamics are invariant to the arm angle $\theta$, the vision-based SE is not. For details on the baseline controllers, choice of gains, and modifications, see the Appendix B.

### Instructions and Code Repositories

Detailed instructions on how to set up the hardware platform are provided at https://git.rwth-aachen.de/quanser-vision/vision-based-furuta-pendulum.

All code needed to reproduce the pixel-to-torque case study presented here is provided at https://git.rwth-aachen.de/quanser-vision/furuta-pixel-to-torque-control.
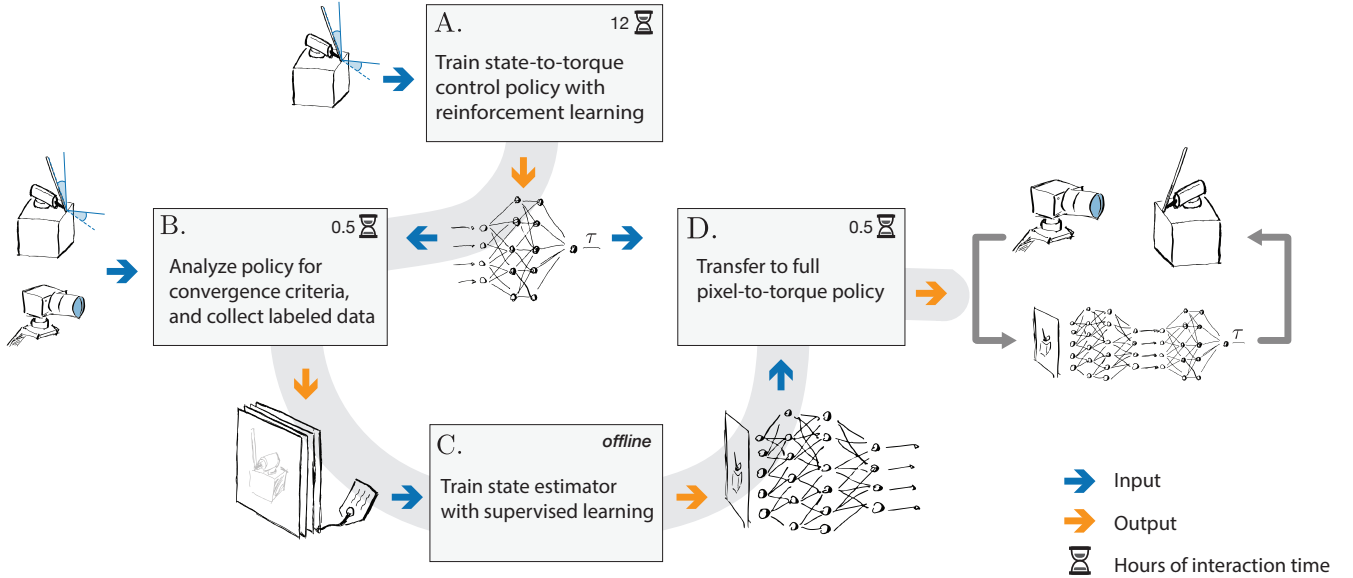
Fig. 2: The four steps of the pipeline can be completed in less than 18 hours, with only 13 hours of hardware interaction time, which is automated largely automated. Note that step *C* is done completely offline and can be accelerated by running the training on a dedicated cluster. Privileged information is used in steps *A* and *B*.

using privileged information as input. In the case of the Furuta platform, the agent learns to swing up and balance the pole in approximately $12\,\mathrm{h}$ of interaction time, which is equivalent to $8\,\mathrm{h}$ worth of samples gathered for learning and $4\,\mathrm{h}$ for resetting. The entire process is automated and could be run in a single session without any intervention.

To enable the agent to learn this task reliably, it was important to tune the reward function and to adjust hyperparameters based on knowledge about the system. We use a continuous reward function, which accelerates training by providing a reliable, steady increase in the accumulated reward. For the Furuta pendulum, we use a quadratic reward penalizing the angle positions of the pendulum with

$$r_t = \left(1 - \frac{4}{5}\frac{|\alpha_t|}{180°} - \frac{1}{5}\frac{|\theta_t|}{180°}\right)^2. \tag{1}$$

We train the agent with a small learning rate and a small clipping factor (compare Table I), which also helps to reliably increase the reward over training episodes. Agents with a large learning rate learned the swing-up task more quickly, but were not able to learn to balance the pendulum reliably: they were susceptible to 'fatal forgetting', or sudden large drops in reward. We surmise this is because balancing requires very precise control inputs, and therefore also a smaller learning rate.

### B. Policy Analysis and Data Collection

Based on the control policy trained on privileged information, we empirically identify minimum precision requirements by injecting noise on the state until the task can no longer be fulfilled. This threshold is then used as the convergence criteria for training the SE in step III-C. For the Furuta

pendulum, we add zero-mean Gaussian noise on the angles, and propagate it via finite differences to the angular velocities. At a sampling frequency of $120\,\mathrm{Hz}$, the agent can tolerate noise with a standard deviation $< 1°$. We also noticed that this level of precision is only necessary to balance the pendulum near the equilibrium point; the policy is able to swing up the pendulum even with higher noise. Based on this observation, we separately collect data for the swing-up portion of the task, and the balancing portion (see "Reproducible Platform"). The convergence criteria is then only tested on images relevant for balancing, which we heuristically determined as $|\alpha| < 10°$. As we will see in Section III-C, converging to high precision over the entire state space not only requires more training time, but a larger DNN.

### C. Learning precise State Estimation

Precise predictions require an unknown minimum network capacity, which makes it difficult to reduce execution time on limited computational resources. We balance this trade-off with a deliberate choice of the DNN architecture, a biased data set, and data augmentation methods.

To increase precision, we simplify the learning task and train a DNN using standard convolutional layers to estimate only the pose from a single image. Velocities are then computed from a buffer of previously estimated positions and velocities via finite differences and a first-order low pass filter (see Figure 3). This structure reduces the SE's prediction error to roughly a fifth compared to a recurrent neural network architecture similar in size, which we speculate is due to the freed capacity being available for higher accuracy on a simpler task. Alternatively, velocities could be estimated by using a history of images as input, but again this would significantly

increase the network size, which we need to reduce as much as possible.

We also downsample the input image from 540×720 to 220×220 pixels, which allows the DNN depth to be increased; we found this was more important for precision than a higher image resolution. To compensate for the downsampling, we add a very small stride of 1-pixel per step. With a depth of 12 layers, the SE reaches a precision that that is able to distinguish individual pixels.

Despite these measures, the limited network size makes it difficult for the DNN to converge to a low error everywhere. Precise state estimates are often not needed throughout the entire state space, and we can evaluate where the SE should be more precise based on the policy analysis conducted in step III-B. For the Furuta pendulum, we bias the training data set to be more densely sampled around the upper equilibrium point. An SE trained on a very biased data set can meet our convergence criteria after just four episodes of training. Due to its reliably low prediction error for small angles (compare Figure 4), the RL agent could also adapt much faster.

To avoid overfitting to the training data set, and to increase the SE's robustness, we also apply data augmentation methods [25] during training. The input images are randomly zoomed, rotated, shifted, and modulated in brightness (compare Table II). While augmenting training data did not increase the accuracy on the validation data set, it substantially increased the prediction performance while testing on the hardware setup.

### D. Policy Transfer

Although the state-to-torque policy does not perform well 'out of the box' with the DNN-based SE, we found that it can be quickly and easily transferred with additional training, without adjusting any learning hyperparameters. With an additional training time of approximately $30\,\mathrm{min}$, the policy adapts to the new input and achieves performance comparable to the policy relying on privileged information. A typical run is shown in Figure 5, where the pendulum reaches the upright position in a single swing.

## IV. Discussion

Our case study demonstrates learning fast and precise pixel-to-torque control of an unstable system, with deep neural networks trained exclusively on real-world data. Although this task does not pose any theoretical obstacles, the practical challenges are substantial, and to the best of our knowledge, this is the first demonstration on a system requiring a control bandwidth of $100\,\mathrm{Hz}$ or higher. We hope this case study will serve as a starting point for further reproducible and comparable studies.

In addition to the usual challenges of using DNNs in control, such as sample efficiency and robustness, we found that achieving both high bandwidth *and* precision becomes a critical challenge for fast and unstable systems. This is especially the case when using function approximators such as neural networks. While large, sophisticated DNN architectures have tremendous representation power, these networks are not only slower to train but also slower to execute at runtime. The high bandwidth requirements severely limit the network size, putting speed and precision of the DNN in direct conflict. In fact, we found that making this trade-off was the dominating factor in designing the learning pipeline.

In order to achieve this while also keeping sample requirements reasonable, we resorted to separating the problems of control and state estimation, and enforcing a state representation based on first principles. Deep learning literature often advocates against a clear separation [5], [12] in favor of allowing the learning process to converge to a latent space representation, which may be more parsimonious and task-invariant. However, in our initial attempts using spatial autoencoders [5], we found that a DNN small enough to meet our execution time requirements was by far not expressive enough to learn a meaningful representation. On a data set of more than 400'000 samples, the autoencoder could not detect features precise enough for the RL agent to noticeably increase the reward after an interaction time of over 24 hours. Learning control and SE separately allowed us to use sample-efficient algorithms, and more easily leverage domain knowledge and privileged information to systematically reduce the DNN size. Furthermore, for the Furuta system, we can directly compare our learning algorithms against conventional controllers as a clear baseline, which is particularly helpful when debugging unexpected learning outcomes.

One of the advantages we did not initially anticipate is that, by learning the control policy first, we could quantify the robustness of this policy to SE noise. This provides clear convergence criteria for supervised learning of the SE and is particularly helpful as a quantitative measure for balancing the trade-off in precision and bandwidth. To cope with limited representation power, we bias the SE training data to regions that require high precision, and compromise in the regions that do not. With the Furuta pendulum, these regions were simply chosen based on system knowledge, but this is not always straightforward to do for more complex systems. An alternative we find promising is to not only test the robustness of a standard control policy, but deliberately train robust control policies with a curriculum [3], for example with progressively noisier environments. Such a policy would further relax the requirements on the SE, which will be critical for more complex tasks or if even higher bandwidth is required. To this end, an important direction of research is to quantify the robustness of a control policy [26].

In the presented case study, we put minimal effort into making the SE robust to changes in the environment [27], and this is certainly an important topic for further studies. Nonetheless, we found the applied data augmentation methods were crucial to increase the SE's performance on the hardware system. To our surprise, the final policy was quite robust to lighting changes: the LED lamp could be dimmed by $30\,\%$ before the policy failed. Instead, we believe that the robustness of the learning pipeline itself is an important aspect that is rarely discussed in the literature. Indeed, we have so far reproduced these results multiple times ourselves, with mixed results. The entire pipeline was first developed with TensorFlow, then re-implemented in PyTorch essentially un-
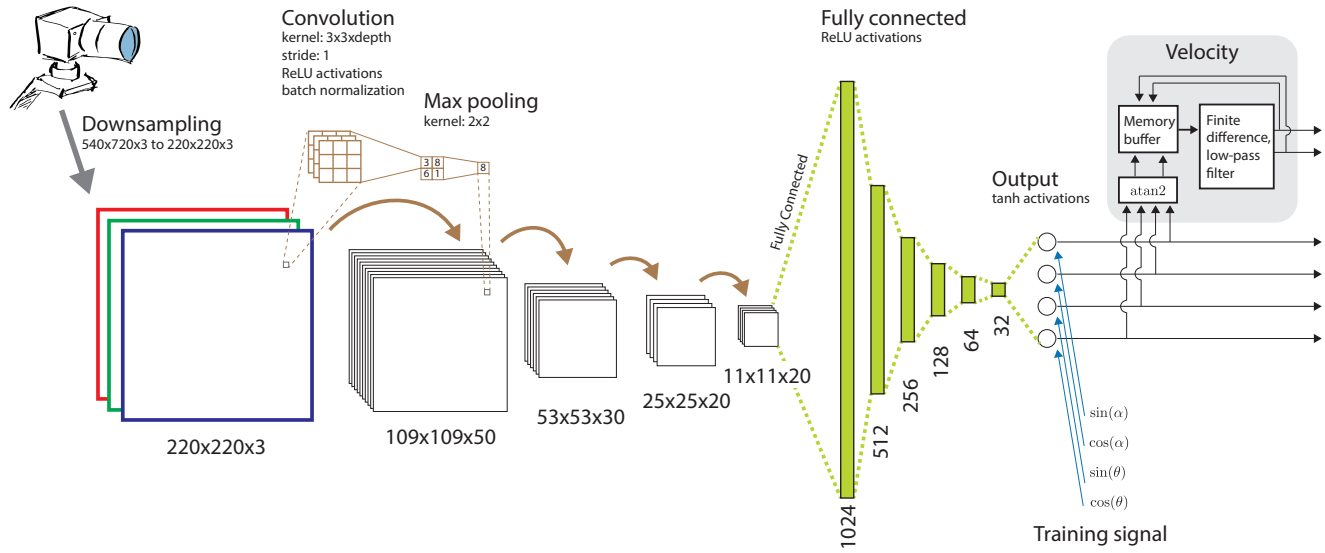
Fig. 3: The architecture of the SE. Instead of predicting the angles directly, we train the DNN to predict $[\cos\theta, \sin\theta, \cos\alpha, \sin\alpha]^T$. We thereby ensure that samples that are very close to each other in the input space are also close in the output space by avoiding jumps for the angles around $\pm 180°$. During training, we apply neuron dropouts with a probability of 0.1 after every max-pooling layer and every fully connected layer to prevent overfitting to the training data set.
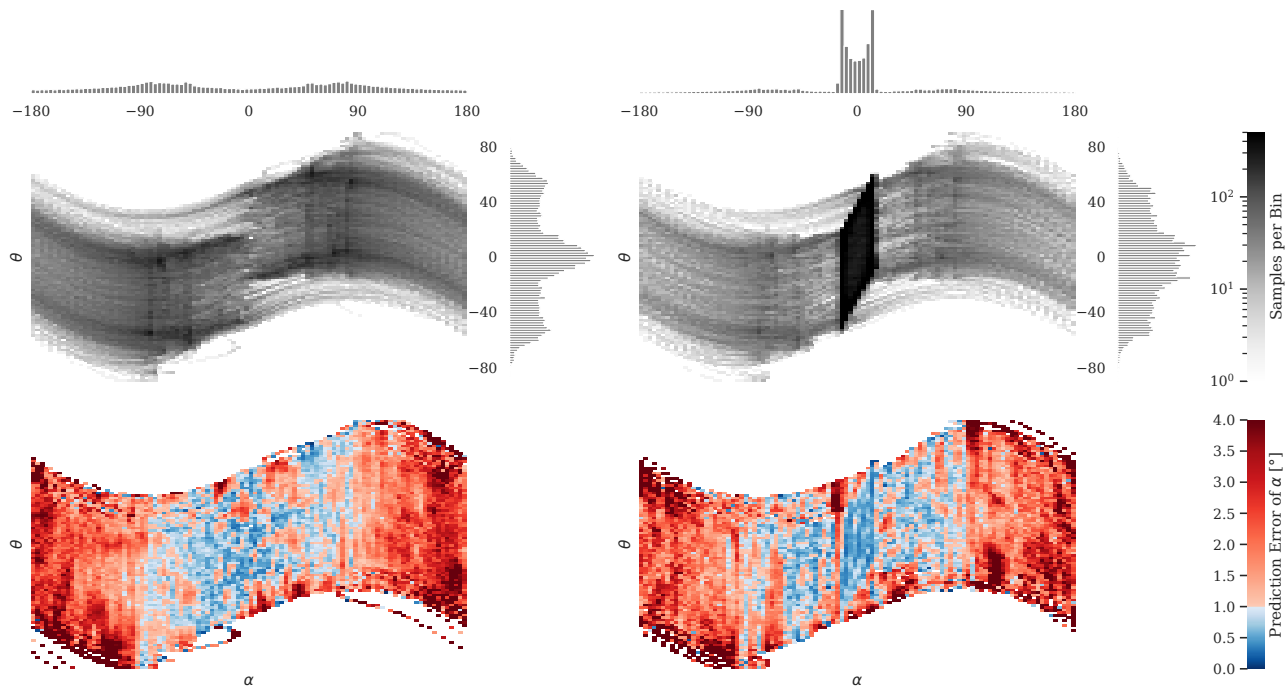


Fig. 4: Comparison of an unbiased (left) and a data set biased around the upright pendulum position (right), with the absolute data distribution over $\theta$ and $\alpha$ (top), and the resulting prediction error distribution after training of the SE (bottom). Both data sets contain 336'000 images and were trained on the same DNN architecture. While the unbiased data set met the convergence criteria after 48 episodes, the unbiased converged after just four episodes. We also found that the RL agent could adapt much more quickly and more reliably to the SE trained on the biased data set, further reducing the interaction time on the hardware system.
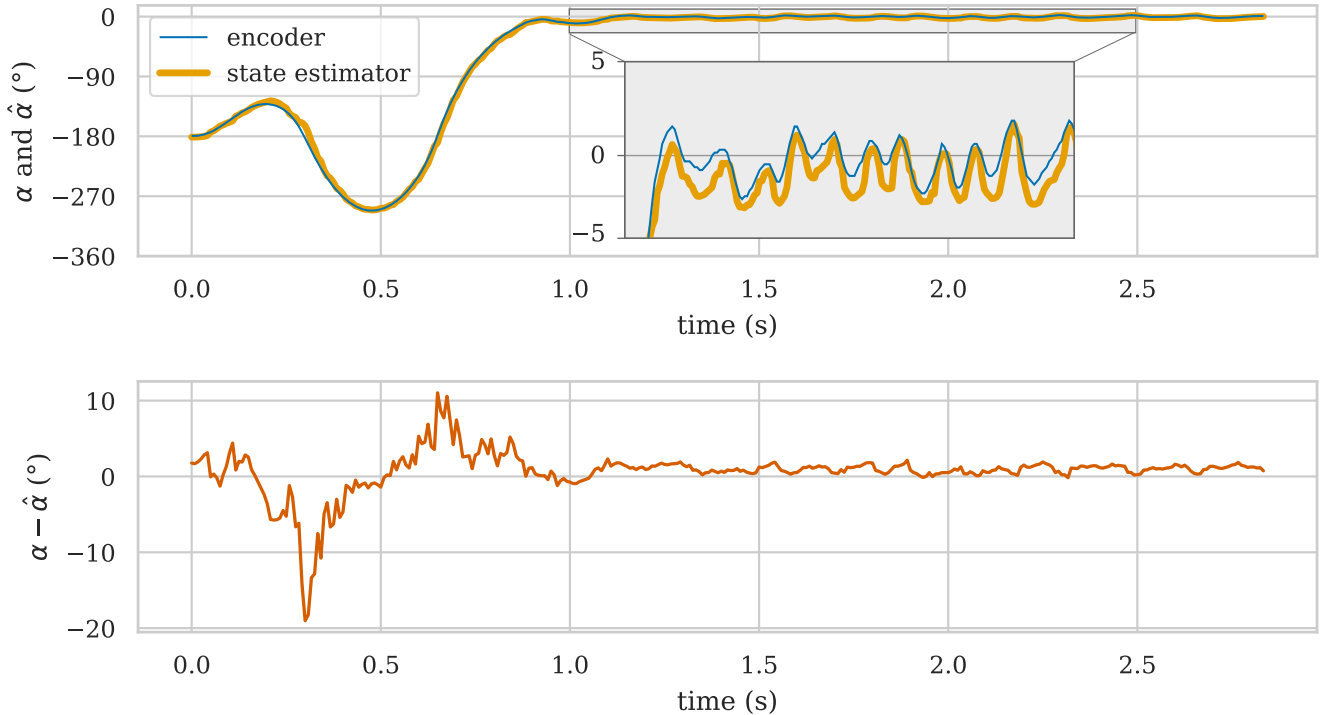
Fig. 5: Swing-up and balance trajectory of the RL agent with the state estimator as input. Despite much larger errors during swing-up, the policy is able to reliably swing up with only one or two swings, which is faster than a typical run using a standard energy-pumping controller.

changed; while learning the SE worked out of the box, the RL agent typically converged to much more aggressive policies, and required additional parameter tuning and testing. After the pandemic started, the entire hardware setup was moved to Steffen's apartment; here, we were pleasantly surprised that the same pipeline, without any algorithmic changes or hyperparameter tuning, reproduced our results. However, once the setup was moved back to the lab, it took an unexpected three days of debugging, recollecting training samples, and training from scratch in order to once again reproduce these results. To better understand how to create reliable learning pipelines, reproducible studies—and reproducing them—are sorely needed. We believe the vision-based Furuta platform we have presented is ideal for such studies. Not only does it capture important challenges for fast and unstable systems, it is simple enough that development iterations can be made quickly, and conventional controllers provide not only a clear baseline to compare against, but also a tool to debug and validate different parts of the learning pipeline. For example, we often determined whether to put more effort into step *A* or step *C* by comparing the DNN-based SE coupled with an LQR controller against the DNN-based policy with encoder readings. We believe the required effort to recreate the vision-based Furuta platform is reasonable, and we look forward to studies that reproduce, and improve on, the results we have presented.

### REFERENCES

[1] C. J. Ostafew, A. P. Schoellig, T. D. Barfoot, and J. Collier, "Learning-based nonlinear model predictive control to improve vision-based mobile robot path tracking," *Journal of Field Robotics*, vol. 33, no. 1, 2016.

[2] S. Levine, P. Pastor, A. Krizhevsky, J. Ibarz, and D. Quillen, "Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection," *The International Journal of Robotics Research*, vol. 37, no. 4-5, 2018.

[3] J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning quadrupedal locomotion over challenging terrain," *Science robotics*, vol. 5, no. 47, 2020.

[4] OpenAI, I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino, M. Plappert, G. Powell, R. Ribas, J. Schneider, N. Tezak, J. Tworek, P. Welinder, L. Weng, Q. Yuan, W. Zaremba, and L. Zhang, *Solving rubik's cube with a robot hand*, 2019. arXiv: 1910.07113.

[5] C. Finn, Xin Yu Tan, Yan Duan, T. Darrell, S. Levine, and P. Abbeel, "Deep spatial autoencoders for visuomotor learning," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2016.

[6] J. Ibarz, J. Tan, C. Finn, M. Kalakrishnan, P. Pastor, and S. Levine, "How to train your robot with deep reinforcement learning: Lessons we have learned," *The International Journal of Robotics Research*, vol. 40, no. 4-5, 2021.

[7] Z. Xie, P. Clary, J. Dao, P. Morais, J. Hurst, and M. Panne, "Learning locomotion skills for cassie: Iterative design and sim-to-real," in *Conference on Robot Learning (CoRL)*, vol. 100, Proceedings of Machine Learning Research, 2020.

[8] J. W. Roberts, L. Moret, J. Zhang, and R. Tedrake, "Motor learning at intermediate reynolds number: Experiments with policy gradient on the flapping flight of a rigid wing," in *From Motor Learning to Interaction Learning in Robots*. 2010.

[9] A. von Rohr, S. Trimpe, A. Marco, P. Fischer, and S. Palagi, "Gait learning for soft microrobots controlled by light fields," in *IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2018.

[10] D. Chen, B. Zhou, V. Koltun, and P. Krähenbühl, "Learning by cheating," in *Conference on Robot Learning (CoRL)*, vol. 100, Proceedings of Machine Learning Research, 2019.

[11] S. Srinivasan, I. Sa, A. Zyner, V. Reijgwart, M. I. Valls, and R. Siegwart, "End-to-end velocity estimation for autonomous racing," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, 2020.

[12] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *The Journal of Machine Learning Research*, vol. 17, no. 1, 2016.

[13] D. Hoeller, L. Wellhausen, F. Farshidian, and M. Hutter, "Learning a state representation and navigation in cluttered and dynamic environments," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, 2021.

[14] A. Zeng, S. Song, J. Lee, A. Rodriguez, and T. Funkhouser, "Tossingbot: Learning to throw arbitrary objects with residual physics," in *Proceedings of Robotics: Science and Systems (RSS)*, 2019.

[15] H. Ha and S. Song, "Flingbot: The unreasonable effectiveness of dynamic manipulation for cloth unfolding," in *Conference on Robotic Learning (CoRL)*, vol. 155, Proceedings of Machine Learning Research, 2021.

[16] N. O. Lambert, D. S. Drew, J. Yaconelli, R. Calandra, S. Levine, and K. S. J. Pister, "Low-level control of a quadrotor with deep model-based reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, 2019.

[17] J. Mattner, S. Lange, and M. Riedmiller, "Learn to swing up and balance a real pole based on raw visual input data," in *International Conference on Neural Information Processing (NIPS)*, Springer, 2012.

[18] E. Kaufmann, A. Loquercio, R. Ranftl, A. Dosovitskiy, V. Koltun, and D. Scaramuzza, "Deep drone racing: Learning agile flight in dynamic environments," in *Conference on Robot Learning (CoRL)*, vol. 87, Proceedings of Machine Learning Research, 2018.

[19] B. S. Cazzolato and Z. Prime, "On the dynamics of the furuta pendulum," *Journal of Control Science and Engineering*, 2011.

[20] K. J. Åström and K. Furuta, "Swinging up a pendulum by energy control," *Automatica*, vol. 36, no. 2, 2000.

[21] I. Paredes, M. Sarzosa, M. Herrera, P. Leica, and O. Camacho, "Optimal-robust controller for furuta pendulum based on linear model," in *2017 IEEE Second Ecuador Technical Chapters Meeting (ETCM)*, 2017.

[22] *Build-its blog*, Accessed: 2020-06-23. [Online]. Available: https://build-its-inprogress.blogspot.com/search/label/Pendulum.

[23] K. Polzounov, R. Sundar, and L. Redden, *Blue River Controls: A toolkit for Reinforcement Learning Control Systems on Hardware*, Accepted at the Workshop on Deep Reinforcement Learning at the 33rd Conference on Neural Information Processing Systems (NeurIPS 2019), Vancouver, Canada. 2019. eprint: arXiv:2001.02254.

[24] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," 2017. [Online]. Available: http://arxiv.org/abs/1707.06347.

[25] J. Wang, L. Perez, *et al.*, "The effectiveness of data augmentation in image classification using deep learning," *Convolutional Neural Networks Visual Recognition*, vol. 11, 2017.

[26] S. Heim, A. Rohr, S. Trimpe, and A. Badri-Spröwitz, "A learnable safety measure," in *Conference on Robot Learning (CoRL)*, vol. 100, Proceedings of Machine Learning Research, 2019.

[27] M. Laskin, K. Lee, A. Stooke, L. Pinto, P. Abbeel, and A. Srinivas, "Reinforcement learning with augmented data," in *Advances in Neural Information Processing Systems (NIPS)*, vol. 33, 2020.

[28] *Quanser qube servo 2 manual*, Accessed: 2020-06-23. [Online]. Available: https://www.quanser.com/courseware-resources/?fwp_resource_related_products=1472.

## APPENDIX

### A. Training Parameters

### B. Baseline controllers and data generation

The parameter values are specific to the Quanser Qube Servo 2 pendulum and may need to be adjusted for other Furuta pendulum systems. The control signal is saturated to stay within motor constraints by simple thresholding.

*1) Swing-Up Control:* As a baseline, we use an energy-based control law from [20].

$$u_{\text{swing-up}} = \mu(E_0 - E)\,\text{sign}(\dot{\alpha}\cos\alpha)$$

where $\mu$ is a tunable gain, $E_0 = mgl$ is the potential energy of the pendulum at the upright equilibrium, with mass $m$ and length $l$, and $E$ is the current total energy. The sign operator applies a bang-bang control signal, which results in better performance.

TABLE I: Training Parameters of the PPO Reinforcement Learning Agent

| General | |
| --- | --- |
| Interaction time, Training on States | 11 h 50 min |
| (corresponding sample time) | 8 h 28 min |
| Interaction time, Adaption to SE | 0 h 28 min |
| (corresponding sample time) | 0 h 16 min |
| Sample Frequency | 120 Hz |
| Horizon | 2048 |
| Minibatch Size | 32 |
| Epochs | 10 |
| Learning Rate | $2 \times 10^{-4}$ |
| Generalized Advantage Estimation | 0.98 |
| Discount Factor | 0.995 |
| Value Function Coefficient | 0.5 |
| Entropy Coefficient | 0.0 |
| Clipping | 0.1 |
| **Policy Network** | |
| Type | Multi-Layer-Perceptron |
| Neurons per Layer | [64, 64, 12] |
| Activation Function | tanh |

TABLE II: Training Parameters of the State Estimator

| Data Set | |
| --- | --- |
| Interaction Time | 28 min |
| Sample Frequency | 200 Hz |
| **General** | |
| Episodes | 4 (until convergence criteria is met) |
| Batch Size | 16 |
| Loss | Mean Squared Error |
| Optimizer | ADAM |
| **Data Augmentation** | |
| Zoom | (1.0, 1.02) |
| x-y-Translation | (-0.01, 0.01) |
| Brightness | (0.9, 1.1) |

For data collection, the baseline controller is modified to sweep large areas of the state space. We add a PID controller $k_{\text{PID},\theta}$ on the arm angle $\theta$ to follow a reference trajectory $\theta_{\text{ref}}(t) = A_{\text{data},1} \sin(f_{\text{data},1}t)$:

$$u_{\text{data},1} = u_{\text{swing-up}} + k_{\text{PID},\theta} \left(\theta - \theta_{\text{ref}}(t)\right)$$

The PID parameters $k_p$, $k_i$ and $k_d$ and trajectory parameters $A_{data,1}$ and $f_{data,1}$ are listed in Table III.

*2) Balancing Control:* As a baseline, we use a Linear Quadratic Regulator (LQR) controller

$$u_{\text{balance}} = -\boldsymbol{K}\boldsymbol{x}.$$

To design the feedback gain $\boldsymbol{K}$ we solve the Ricatti equation of the linear dynamic model of the pendulum at the upright equilibrium $\boldsymbol{x}_0 = [0,0,0,0]^T$ and $\boldsymbol{u}_0 = 0$. The system matrices of the linear model are provided by Quanser [28] as

$$\boldsymbol{A} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 149.3 & -0.01004 & 0 \\ 0 & 261.6 & 1 & -0.0103 \end{bmatrix}$$

and $\boldsymbol{B} = [0, 0, 49.73, 49.15]^T$.

We found the LQR weight matrices

$$\boldsymbol{Q} = \begin{bmatrix} 12 & 0 & 0 & 0 \\ 0 & 5 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

and $R = 1$ to be suitable for reliable and stable data generation trajectories.

To generate data around the upper equilibrium point we perturb the controller with two input signals:

$$u_{\text{data},2} = -\boldsymbol{K}\left(\boldsymbol{x} - \boldsymbol{x}_{\text{ref}}(t)\right) + u_{\text{oscillation}}(t)$$

where $u_{\text{oscillation}}(t) = A_{\text{data},1} \sin(f_{\text{data},1}t)$ is a fast oscillation and serves to gather more samples in the region $|\alpha| < 15°$, and $\boldsymbol{x}_{\text{ref}}(t) = [\theta_{\text{ref}}(t), 0, 0, 0]^T$ is a slow oscillation with $\theta_{\text{ref}}(t) = A_{\text{data},2} \sin(f_{\text{data},2}t)$, and serves to cover large areas of $\theta$. All control and signal parameters are listed in Table III.

TABLE III: Parameters for Data Generation

| **Energy-based swing-up control** | |
| --- | --- |
| $k_p$ | 0.5 |
| $k_i$ | 0.5 |
| $k_d$ | 0.05 |
| $A_{\text{data},0}$ | 60° |
| $f_{\text{data},0}$ | 0.05 Hz |
| **LQR control for balancing** | |
| $A_{\text{data},1}$ | 28 V |
| $f_{\text{data},1}$ | 2.4 Hz |
| $A_{\text{data},2}$ | 30° |
| $f_{\text{data},2}$ | 0.03 Hz |