DATA SCARCITY IN ROBOTICS: LEVERAGING STRUCTURAL PRIORS AND REPRESENTATION LEARNING

by

Artem Molchanov

A Dissertation Presented to the FACULTY OF THE USC GRADUATE SCHOOL UNIVERSITY OF SOUTHERN CALIFORNIA In Partial Fulfillment of the Requirements for the Degree DOCTOR OF PHILOSOPHY (Computer Science)

January 2021

Copyright 2021

Artem Molchanov

Acknowledgments

There are many people I would like to thank for their constant support and numerous contributions on my path towards my Ph.D.

First and foremost, I would like to thank my advisor Professor Gaurav S. Sukhatme, who showed me the path toward developing myself as a researcher. His patience and insightful advice has always helped me make the right decisions in important moments. Gaurav's approach to mentoring has always provided me with plenty of freedom to find and explore problems I wanted to follow. That is what made me an independent thinker and a researcher and that is what makes the time spent at RESL such a valuable experience.

I would like to express special gratitude toward my committee members: Professor Heather Culbertson, Professor Satyandra K. Gupta, Professor Nora Ayanian, Professor Stefan Schaal and Professor Joseph Lim for their constructive feedback, encouragement and overall help they kindly provided in my qualifying exams and the defense.

I am extremely grateful to my colleague, collaborator and dear friend Wolfgang Hönig. It is very hard to over-appreciate his contributions to my life at the hardest moments during my Ph.D. On numerous occasions, Wolfgang helped me as an advisor, as a software and hardware contributor to my projects and simply as a friend when I felt completely lost in the problem at hand. I'd also like to give special thanks to Karol Hausman and Yevgen Chebotar, who came to support me in the most critical stages of our projects under extremely stressful conditions. They presented me with an excellent example of a path that I have always aimed to replicate. I am very grateful to my first collaborator and mentor Andreas Breitenmoser, whose patience and great insights helped me to solve my first hard problem, that I am very proud of. I will never forget my numerous colleagues from RESL, ACT, CLVR, and CLMC laboratories. It would be impossible to mention everyone, but each of you made an exceptional social circle that I will greatly miss in the future steps of my life. I would also like to acknowledge my colleagues, friends and mentors I interacted with during my internships: Franziska Meier, Edward T. Grefenstette, Jan Kautz, Stan Birchfield, Jonathan Tremblay, Lutz Junge, Premkumar Natarajan, Marci Meingast, and many others. All of them helped me a lot during those short summer sprints that greatly enriched my life and professional experience.

And last, but not least, I would like to express deepest gratitude to my mother and my closest friend Philo Wells, who has become a family to me over these years. My family has always showed me infinite and unconditional support. It was a lot of ups and downs and without all of these amazing people I would not have been able to succeed in my Ph.D. Thank you all for your major contributions on this challenging path.

Table of Contents

| Acknowl | edgments | ii |
|--|---|--|
| List of Ta | ables | vii |
| List of Fi | igures | iii |
| Abstract | x | iii |
| 1 Introd 1.1 1.2 | duction | 1 3 6 |
| I Perce | ptual Data Scarcity | 10 |
| 2 Overa 2.1 2.2 2.3 2.4 2.5 | coming Scarcity of Point-of-Contact Estimation with Tactile SensorsRelated WorkBiomimetic Tactile SensorPoint-of-Contact Estimation2.3.1Contact Detection2.3.2Feature Extraction2.3.3Regression2.3.4ClassificationEvaluation2.4.1Data Collection Setup2.4.2Results and DiscussionConclusion | 11 14 17 18 19 19 20 22 23 26 30 |
| 3 Aggr 3.1 3.2 | egating Multiple Active Drifters Under Scarcity Of Current MeasurementsRelated WorkSingle Drifter Control3.2.1Regional Ocean Modeling System for Active Drifter Control3.2.2Predictions and Measurements of Currents | 32 35 36 37 38 |

| | | 3.2.3 | Controlling a Single Active Drifter 38 |
|----|-------|----------------------|---|
| | | 3.2.4 | Conclusion |
| | 3.3 | Multi-I | Orifter Control |
| | | 3.3.1 | Problem Formulation |
| | | 3.3.2 | Drifter Control Methods 4' |
| | | 3.3.3 | Performance Evaluation |
| | 3.4 | Conclu | sion |
| | | | |
| II | Task | Data So | carcity |
| 4 | Sim- | to-(Mult | i)-Real: Transfer of Low-Level Robust Control Policies to Multiple Quadro- |
| | tors | Deleter | |
| | 4.1 | Drobler | 1 WOIK |
| | 4.2 | Proble | in Statement |
| | 4.3 | | $\begin{array}{c} \text{Incs Simulation} \\ \text{Divid Deduce for Overlapter} \\ \end{array}$ |
| | | 4.3.1 | Negretalized Mater Threat Lengt |
| | | 4.3.2 | Simulation of Non-Ideal Motors |
| | | 4.3.3 | Simulation of Non-Ideal Motors |
| | | 4.5.4 | Numerical Matheda 77 |
| | 11 | 4.3.3 Loornii | Numerical Methods |
| | 4.4 | | Pandomization 7 |
| | | 4.4.7 | Policy Representation 70 |
| | | ч.ч.2 ДДЗ | Policy Learning 70 |
| | | т.т. <i>3</i> ДДД | Sim-to-Sim Verification |
| | | 445 | Sim-to-Real Verification 8 |
| | 45 | Experi | ments 8 |
| | 1.5 | 4 5 1 | Ablation Analysis on Cost Components |
| | | 452 | Sim-to-Real: Learning with Estimated Model |
| | | 4.5.2 | Sim-to-Multi-Real: Learning without Model |
| | | 454 | Control Policy Robustness and Recovery |
| | 4.6 | Conclu | sions |
| 5 | Tack | Spacifi | Learning with Searce Data via Mate learned Lesses |
| 5 | 1 ask | Palata | Work |
| | 5.1 | Moto I | arming via Learned Loss |
| | 5.2 | 5 2 1 | $MI^{3} \text{ for Supervised Learning} \qquad \qquad$ |
| | | 5.2.1 | ML ³ Poinforcement Learning |
| | | 5.2.2 | Shaping MI ³ loss by adding extra loss information during mata train 10^{4} |
| | 52 | J.2.J Export | Shaping with loss by adding extra loss information during meta-if $unit = 10^{-1}$ |
| | 5.5 | 5 3 1 | Learning to mimic and improve over known task losses |
| | | 537 | Shaping loss landscapes by adding extra information at mate train time 11 |
| | | 5.5.2 | Shaping loss landscapes by adding extra information at meta-train time . 11 |

| | 5.4 | Conclusions | 17 |
|-----|--------|----------------------|----|
| 6 | Conc | sion and Future Work | 9 |
| | 6.1 | Conclusion | 20 |
| | 6.2 | uture Work | 21 |
| Bil | bliogr | hy | 23 |

List of Tables

| 3.1 | Simulation parameters | 57 |
|-----|--|----|
| 3.2 | System performance over 100 simulations | 58 |
| 4.1 | Randomization variables and their distributions. | 78 |
| 4.2 | Robot Properties. | 82 |
| 4.3 | Ablation Analysis on Cost Components. α values not listed are 0 | 85 |
| 4.4 | Sim-to-Multi-Real Results. | 86 |

List of Figures

| 2.1 | Cross-sectional schematic of the BioTac sensor (adapted from Su et al. [2012]). | 18 |
|------|--|----|
| 2.2 | The 18 objects used for data collection | 22 |
| 2.3 | An example of a grasp used and contact samples collected during the experiment. | |
| | Left: an example grasp of an object. Right: an example of the detected contact | |
| | points and directions projected onto the <i>zx</i> plane | 23 |
| 2.4 | The normalized histogram of the peak tapping forces applied to the object during | |
| | the experiment. | 24 |
| 2.5 | The tapping rod. Left: an image of the rod. Right: rod coordinate frame with the | |
| | locations of the markers. | 25 |
| 2.6 | The wrist band with Vicon markers. Left: CAD model of the Barrett hand with | |
| | the band. Right: wrist coordinate frame with locations of the markers | 26 |
| 2.7 | Results of regression using different sensor modalities for NN and GP regressors. | |
| | <i>NN:Electr</i> - NN with electrodes only; <i>NN:Electr+PAC+PDC</i> - NN with the | |
| | full set of sensor modalities (electrodes, PAC, PDC); GP:Electr+PAC+PDC | |
| | - GP with the full set of sensor modalities. x, y, z, yaw, pitch are results for | |
| | independent dimensions. xyz represents results for combined error, i.e. the norms | |
| | of the error vectors in Cartesian space. Errors in Cartesian space are reported in | |
| | cm; angular errors are reported in degrees | 27 |
| 2.8 | Results of classification using different sets of sensor modalities: NN:Electr - | |
| | NN classifier with electrodes only; NN:Electr+PAC+PDC - NN classifier with | |
| | electrodes, PAC and PDC features; ST-HMP:Electr - ST-HMP feature learning | |
| | algorithm with SVM classifier applied to electrode features | 28 |
| 2.9 | Results of NN classification using different time lengths of features, starting | |
| | from the moment contact event is detected. | 29 |
| 2.10 | Results of NN classification using different grid steps | 30 |
| 2.11 | An example confusion matrix of the x dimension classification for one of the | |
| | objects | 31 |
| 3.1 | Lagrangian drifters. Left: A prototype of our passive drifter system, which was | |
| | deployed in the Southern California Bight to measure ocean currents. Right: The | |
| | schematic shows the main components and the mode of operation of an active | |
| | drifter | 33 |
| | | |

| 3.2 | Comparison of ocean current measurements and ROMS predictions: the cur- rent forecasts (green vectors) and nowcasts (red vectors) are plotted along the measured trajectory and estimated currents of the deployed passive drifter (blue dotted line and blue vectors). The black circle in the middle represents the center | |
|------|---|----|
| | of a cell of the grid that underlies ROMS. | 39 |
| 3.3 | Comparison of ocean current measurements and ROMS predictions: trajectories generated by the "drop a drifter" web page, and the real trajectory of the drifter | 40 |
| 3.4 | Comparison of ocean current measurements and ROMS predictions: absolute | 40 |
| 3.5 | values of the ocean current velocity vectors | 41 |
| 3.6 | opposite direction) | 42 |
| | pictures depict the evolution of trajectories over time (2 days, 30 days corre- spondingly). The colored crosses represent goal targets. The passive drifters managed to hit only 1 target out of 8. The simulation is run in a static ocean current field generated from ROMS data. | 43 |
| 3.7 | Simulated trajectories (colored curves) of active drifters starting from the same location (green circle). From left to right, the pictures depict the evolution of trajectories over time (2 days, 30 days correspondingly). The colored crosses represent goal targets. The color of the target corresponds to the color of the drifter it was assigned to. The active drifters managed to hit 5 targets out of 8. | |
| 3.8 | The simulation is run in a static ocean current field generated from ROMS data. The mean of the minimum distance to a destination vs. initial distance to a destination. Each graph is made of 6 discrete data points with different initial distances. Each point is the mean over 1000 simulated trials with the same initial distance. In each trial, a drifter was given the task to go from a random initial point to a destination picked randomly, but with the predefined distance between them. The minimum distance achieved by the drifter in each trial was recorded. | 43 |
| | The duration of one trial is 180 days. | 44 |
| 3.9 | The artificial potential energy function $U(d)$ for tuple a_u . The function defines | |
| | the quality of aggregation of a pair of drifters. | 52 |
| 3.10 | ROMS point-wise fitness function $f_a(\theta_{max1}, \theta_{max2})$. The function defines how well the plane is spanned at a point by a triple of vectors, based on the angles | |
| | between each pair of vectors. | 54 |
| 3.11 | Plots of ROMS fitness maps. Left: A map created from the full 12 layer set. Center: A map created from the best 3 layer set. Right: A map created from the | |
| | 6 layer set used in our simulations. | 55 |

| 3.12 | Evolution of the average aggregation metric over time based on 100 simulation runs. The green curve represents the performance of ideal drifters. The black curve represents the performance of realistic drifters with a low level controller | 50 |
|------|---|----|
| 3.13 | Performance of the realistic drifter system with the basic low-level controller under different tolerance angles. Every point is based on 100 simulation runs. The black curve represents the mean value and the green curves represent the | 59 |
| 3.14 | standard deviation of the aggregation metric at every point Performance of the realistic drifter system with the extended low-level controller under different T_{trust} . Every point is based on 100 simulation runs. The black curve represents the mean value and the green curves represent the standard | 60 |
| 3.15 | deviation of the aggregation metric at every point | 61 |
| 3.16 | deviation of the aggregation metric at every point | 62 |
| 4.1 | Three quadrotors of different sizes controlled by the same policy trained entirely in simulation. | 67 |
| 4.2 | Top-down view of our generalized quadrotor model. There are 5 components: baselink, payload, 4 arms, 4 motors, and 4 rotors. The model always assumes the \times configuration, with the front and left pointing at the positive X and Y directions respectively. Motors are indexed in the clockwise direction starting from the front-right motor. The front-right motor rotates counterclockwise and generates a thrust force f_1 | 73 |
| 4.3 | Trajectory tracking performance of a Crazyflie 2.0 using different controllers. The target trajectory is a figure-eight (to be executed in 5.5 s). The Mellinger controller has the lowest tracking error (0.11 m). Our baseline NN controller has lower tracking error (0.19 m) than the Mellinger with integral gains disabled (Mellinger no memory; 0.2 m). | 88 |
| | | |

| 4.4 | An example of a recovery trajectory from a random throw with an initial linear velocity of approximately 4 m/s. Trajectory colors correspond to quadrotor's |
|-----|--|
| | velocities. Arrow colors correspond to the names of the body-frame axes 91 |
| 5.1 | Framework overview: The learned meta-loss is used as a learning signal to optimize the optimize f_a which can be a regressor a classifier or a control policy 94 |
| 5.2 | Meta-learning for regression (top) and binary classification (bottom) tasks. (a) |
| 5.2 | meta-train task, (b) meta-test tasks $\dots \dots \dots$ |
| 5.3 | performance of the meta-network on the meta-train task as a function of (outer) |
| | meta-train iterations in blue, as compared to SGD using the task-loss directly in |
| | orange, (b) average performance of meta-loss on meta-test tasks as a function of |
| 54 | In number of gradient update steps $\dots \dots \dots$ |
| 5.7 | the ML ^{3} loss generalizes well the loss was trained on the blue trajectories and |
| | tested on the orange ones for the PointmassGoal task. ML^3 loss also significantly |
| | speeds up learning when compared to the task loss at meta-test time on the |
| | PointmassGoal (b) and the ReacherGoal (c) environments |
| 5.5 | ML^3 for model-free RL: results are averaged across 10 tasks. (a+b) Policy |
| | learning on new task with ML ³ loss compared to PPO objective performance |
| | during <i>meta-test</i> time. The learned loss leads to faster learning at meta-test |
| | time. (c+d) Using the same ML° loss, we can optimize policies of different |
| 56 | architectures, showing that our learned loss manifulns generality |
| 5.0 | ground-truth (extra) information at meta-train time: a) Comparison of learned |
| | ML^3 loss (top) and MSE loss (bottom) landscapes for fitting the frequency of |
| | a sine function. The red lines indicate the ground-truth values of the frequency. |
| | b) Comparing optimization performance of: ML ³ loss trained with (green), and |
| | without (blue) ground-truth frequency values; MSE loss (orange). The ML^3 loss |
| | learned with the ground-truth values outperforms both the non-shaped ML^3 loss |
| | and the MSE loss. c-d) Comparing performance of inverse dynamics model |
| | learning for ReacherGoal (c) and Sawyer arm (d). ML ³ loss trained with (green) |
| | and without (blue) ground-truth inertia matrix is compared to MSE loss (orange). The changed ML 3 loss outperforms the MSE loss in all asses |
| 57 | (a) MountainCar trajectory for policy optimized with iL OR compared to ML ³ |
| 5.7 | loss with extra information. (b) optimization performance during meta-test time |
| | for policies optimized with iLQR compared to ML^3 with and without extra |
| | information |

5.8 ReacherGoal with expert demonstrations available during meta-train time. (a) shows the targets in end-effector space. The four blue dots show the training targets for which expert demonstrations are available, the orange dots show the meta-test targets. In (b) we show the reaching performance of a policy trained with the shaped ML³ loss at meta-test time, compared to the performance of training simply on the behavioral cloning objective and testing on test targets. 116

Abstract

Recent advances in Artificial Intelligence have benefited significantly from access to large pools of data accompanied in many cases by labels, ground truth values, or perfect demonstrations. In robotics, however, such data are scarce or absent completely. Overcoming this issue is a major barrier to move robots from structured laboratory settings to the unstructured real world. In this dissertation, by leveraging structural priors and representation learning, we provide several solutions when data required to operate robotics systems is scarce or absent.

In the first part of this dissertation we study sensory feedback scarcity. We show how to use high-dimensional alternative sensory modalities to extract data when primary sensory sources are absent. In a robot grasping setting, we address the problem of contact localization and solve it using multi-modal tactile feedback as the alternative source of information. We leverage multiple tactile modalities provided by electrodes and hydro-acoustic sensors to structure the problem as spatio-temporal inference. We employ the representational power of neural networks to acquire the complex mapping between tactile sensors and the contact locations. We also investigate scarce feedback due to the high cost of measurements. We study this problem in a challenging field robotics setting where multiple severely underactuated aquatic vehicles need to be coordinated. We show how to leverage collaboration among the vehicles and the spatio-temporal smoothness of the ocean currents as a prior to densify feedback about ocean currents in order to acquire better controllability.

In the second part of this dissertation, we investigate scarcity of the data related to the desired task. We develop a method to efficiently leverage simulated dynamics priors to perform sim-to-real transfer of a control policy when no data about the target system is available. We investigate this problem in the scenario of sim-to-real transfer of low-level stabilizing quadrotor control policies.

We demonstrate that we can learn robust policies in simulation and transfer them to the real system while acquiring no samples from the real quadrotor. Finally, we consider the general problem of learning a model with a very limited number of samples using meta-learned losses. We show how such losses can encode a prior structure about families of tasks to create well-behaved loss landscapes for efficient model optimization. We demonstrate the efficiency of our approach for learning policies and dynamics models in multiple robotics settings.

Intelligence is not the ability to store information, but to know where to find it.

- Albert Einstein

Chapter 1

Introduction

In the past few years there has been significant progress in the field of artificial intelligence (AI): today AI programs are capable of recognizing images with accuracy exceeding human level [He et al., 2015], answer questions from visual content [Chen et al., 2015], provide high-quality language translations [Vaswani et al., 2017], generate synthetic voice closely resembling humans [Oord et al., 2016], play Atari and Go games at the level exceeding the best champions in the world [Mnih et al., 2015, Silver et al., 2016], and much more.

Such rapid advancement creates hope that robotics is about to take the next step and robots will soon acquire capabilities and skills resembling humans. Nonetheless, despite recent progress in robotics, we still have not managed to go much further than the level of very specialized systems that can only work in structured environments, such as research laboratories, industrial assembly lines, or paved roads. It is especially interesting to notice that we are still struggling with simple pick and place tasks that roboticists were working on 50 years ago [Feldman et al., 1971].

This discrepancy largely stems from the fact that the recent impressive results in AI are driven mainly by the availability of vast amounts of data [Krizhevsky et al., 2012, Abu-El-Haija et al., 2016] and computational power capable of utilizing large pools of information. These data are often accompanied by low-noise descriptive labels, annotations, measurements, and/or demonstrations

that directly (and in some cases, trivially) relate to the objective. In this dissertation, we aggregate these types of data under the term *high quality data*. In robotics, however, availability of such data is traditionally very limited. This is due to several factors: First, robots are diverse and expensive to run and operate both in terms of time and cost. Second, they often use very specific information, for example, motor commands or special properties of objects, such as orientation, friction coefficients of surfaces and materials, or grasping points and grasp configurations. As a result, robots often lack sensors to directly acquire task-relevant information in a straightforward way. Additionally, these factors are highly dependent on specific robot hardware and the experimental conditions making standardization of datasets problematic. Third, in order to operate in the real world, robots have to deal with a wide range of diverse conditions and adapt to new situations under noisy and incomplete information that can non-trivially relate to the task. Fourth, collecting data to allow adaptation to novel conditions on robots is often an unsafe process, because controllers can lose stability or behave unpredictably in novel situations – resulting in damaging hardware or surrounding structures, and even leading to injuries to humans involved in the process.

Realistically therefore, we can only assume *scarce* access to high quality data available to the robot during its operation. To overcome this challenge, we have to shift our focus to developing methods that can i) find and leverage alternative sources of information; ii) utilize sources of data, where the relationship between the data and useful quantities of interest to the control system designer is not straightforward and; iii) introduce useful, yet not too restrictive, forms of priors that drastically improve the utilization of the few samples of data available. Such a paradigm shift is the central theme of this dissertation. In this dissertation, we classify data scarcity problems into two groups:

Perceptual/Feedback data scarcity. This type of scarcity relates to the lack of sensors that directly measure the needed quantity or when the required quantity cannot be directly measured with an acceptable spatio-temporal density. An example is rare access to the GPS signal, when a robot has to operate inside of a building or underwater.

Task data scarcity. This type of scarcity refers to the case when the information about the task itself or the system executing the task is absent or not abundantly available. An example of such a scenario is a quadrotor that requires initial controller tuning. Availability of data, in this case, is very limited due to high probability of a crash and the resulting hardware damage.

1.1 Perceptual Data Scarcity

The main goal of a perceptual sub-system is estimation or inference of a state. The state represents a compressed (i.e. low-dimensional) world representation. Such representation is either interpretable or abstract. The latter is primarily used within the end-to-end learning pipelines that are growing in popularity. But despite the ongoing push toward using abstract representations [Bengio et al., 2013], there is still a lot of value in developing and augmenting robotic systems with interpretable low-dimensional feedback. Such feedback (and state) representation allows us to employ traditional components of the robotics pipeline, such as planning. For the less traditional (end-to-end) pipelines, interpretability is one of the key factors for better understanding processes in the system. Using compressed rather than high-dimensional representations enables a significant reduction in complexity of algorithms, both in terms of computations and collected experience, which is a key factor for compute-constrained robotic systems. Hence, in this portion of the dissertation, we focus on problems related to estimating interpretable low-dimensional states.

The most frequent problems of estimating states in robotics are i) high noise in measurements; ii) scarce availability of the direct measurement of the state; or iii) complete absence of a sensor measuring the desired quantity. The first two problems have been partially addressed in the areas of *data fusion* and *state estimation*. A majority of approaches are based on different forms of the Bayesian filter [Thrun et al., 2005] that aims to estimate a belief about the state. These approaches incorporate measurements, also referred to as observations, from different sensory modalities with a known probabilistic model relating those measurements to the state. While there is a large variety of different filters, they all assume that such an observation model is available. Here, we focus on problems ii) and iii) instead, which pose additional major challenges:

- **Complex observation-state mapping.** In many cases the engineering of the transformation between the observation and the state spaces is not feasible. This challenge is especially relevant for the scenario of the high-dimensional and possibly multi-modal sensors.
- **Time-energy cost.** This challenge arises when the act of measurement is expensive in terms of time and/or energy. That leads to the need for a careful trade-off between acquiring feedback and performing the desired task.

In the first part of this dissertation we address these two challenges considering two pertinent robotics applications that suffer from perceptual data scarcity.

The first application [Molchanov et al., 2016b,a] is concerned with the absence of a wellcharacterized mapping from sensory inputs to the estimated state. Specifically, we consider the difficult task of estimating the point of contact between an object handled by the robot and the environment (later called *indirect contacts* for simplicity) using high-dimensional and multimodal tactile sensors. This problem is particularly important for contact verification and failure monitoring during tool handling. Force-torque sensors are the traditional sensors for this purpose, but their high cost and relatively narrow application drives exploration of more general sensor, such as cameras. Unfortunately, vision-based perceptual systems have numerous fundamental issues inhibiting their utility for dexterous manipulation and contact handling in particularly. Among these issues are occlusions, changing lighting conditions, and fundamental ambiguity of depth perception [Manhardt et al., 2019, Makihara et al., 2003, Ishikawa et al., 1996]. These shortcoming make access to manipulation-relevant features very scarce. On the other hand there is a strong evidence in neuroscience that humans heavily rely on tactile perception [Johansson and Flanagan, 2009b, Goodwin and Wheat, 2004] for manipulation skills. Thus, tactile perception should provide us with an alternative source of information for contact handling, possibly even in the complete absence of information from other sources. State-of-the-art tactile sensors provide high-dimensional and multi-modal observations [Wettels et al., 2008]. The mapping between such observations and the required manipulation states can be arbitrary complex. To solve this problems we employ a two-fold approach. First, we leverage state-of-the-art multi-modal tactile sensors that allow structuring the problem as spatio-temporal inference. Second, we apply state-of-the-art learning approaches, including spatio-temporal hierarchical matching pursuit (ST-HMP) and deep neural networks to extract the unknown mapping between high-dimensional tactile sensor modalities directly to the target states (see Chapter 2).

The second application [Molchanov et al., 2015, 2014] demonstrates the data scarcity connected to time and energy cost of the measurements. In the more common version of this scenario one has to consider time or energy consumption of the system only, in other words, the robot can reduce uncertainty of the state without significant sacrifice to the performance on the task. An example could be a task where a robot has to recognize an object using multiple viewpoints for the purpose of adding the object to a semantic map. The viewpoint selection problem is often treated in isolation since it does not cause serious problems to collect more samples to reduce uncertainty up to a desired level. This type of problems has traditionally been tackled in the areas of active and interactive perception [Bohg et al., 2017]. To minimize the amount of interactions, many

approaches consider maximization of some information-theoretical criteria, such as information gain, to select the next available action. A more complicated version of this scenario arises when the act of measurement itself, besides taking effort, also directly and negatively impacts the goal behavior of the system, which creates a very hard exploration-exploitation dilemma. On one hand, the agent must perform measurements to reduce uncertainty of the state to select actions reasonably well. On the other hand, while reducing uncertainty it is forced to interrupt the ongoing task or redirect resources that significantly degrades the progress. In this dissertation, we consider such a scenario with an example of a swarm of highly underactuated vehicles, called drifters that are performing aggregation and coverage tasks in the open ocean. On top of the aforementioned challenges, our system is complicated by highly unpredictable dynamics and severely restricted actions available to the vehicles. We show how one can leverage collaboration among drifters and the spatio-temporal smoothness of the ocean currents as a prior to densify feedback about ocean currents to acquire controllability in this poorly controllable and unpredictable environment (see Chapter 3).

1.2 Task Data Scarcity

Another key difference between robotics and many traditional fields in AI, such as computer vision or natural language processing, comes from the fact that the purpose of robots is not only to perceive the world, but to change it in some meaningful way. Thus, in robotics, we need to find the appropriate action that the robot should perform at every state and/or moment of time. In practice, this requires us to communicate the task to the robot. Usually to solve a task we need two kinds of information i) information about dynamics model of the robot and the environment; and ii) the specification of the goal that the robot should achieve in the environment. The former is usually solved by one of the following approaches:

- Hand-engineering. This is the traditional method based on our understanding of physical systems. It provides meaningful and interpretable priors for the design of the controllers based on them. It is often restrictive to very specific systems and limited in how accurately one can express the dynamics.
- Learning. Learning-based approaches often operate on abstract models. Parameters of these models are learned for each system from the data collected on the system. Such models are very flexible and can potentially fit very complex dynamics. Yet, these methods are often data-hungry and require substantial coverage of the state space to avoid the problem of covariate shift, where the robot may be in a region of the state space it has not seen before and the model would fail to generalize.

As for the goal specification, despite a long history of intrinsically motivated systems [Oudeyer and Kaplan, 2009, Schmidhuber, 2010, Houthooft et al., 2016, Jaderberg et al., 2017, Forestier et al., 2017], the main source of the goal knowledge comes from human expertise that can be separated into two major sources:

- **Hand-engineering.** For example, it could be fully designed deterministic controllers, such as traditional linear control systems [Dorf and Bishop, 2000] or numerous trajectory planning algorithms [Gasparetto et al., 2012]. Often in these cases the system becomes very narrow-purpose built, i.e. a system that can not generalize and adapt under changing conditions.
- **Demonstrations and annotations.** This type of human expertise can also come in a variety of forms, such as labels, observations of skills, or the actions corresponding to certain skills.

Practically all of the aforementioned scenarios require scarcely available data. For example, hand engineering requires data collection for system identification that involves sophisticated setups and restricts applicability. Learning involves unsafe data collection, e.g. following a few trajectories before the expensive robotic system crashes and breaks. Demonstrations are also hard or (in some cases) impossible to acquire, for example, when a robot has to take care of elderly, when humans are simply not present, or, as in many cases, when it is impossible to provide demonstrations directly on the robot. Hence, in this part of the dissertation we address these shortcomings in the context of two different scenarios.

In the first scenario, we argue that simulators can help with i) (virtually) unlimited data for learning; ii) appropriate priors and help injecting already available knowledge into the learning system; and iii) problems with complexity, safety, and the cost instrumentation of training setups. We demonstrate [Molchanov et al., 2019a, Dawson, 2020, Molchanov et al., 2019b] that it is possible to solve complex robotics problems using a combination of reinforcement learning (RL) and simulators to match and even exceed performance on robotics problems with well-engineered controllers, while minimizing the engineering effort and requiring no samples from the target robotic system. Reinforcement learning is, arguably, one of the most contradictory tools in robotics research. On one hand, it has shown a lot of promise in finding complex control policies, while requiring minimal engineering for specific problems. On the other hand, due to extreme data inefficiency, it remains unclear if it can be successfully applied to robotic systems to improve upon existing (and more traditional) approaches. So far, the majority of demonstrations have shown results on either toy examples or have solved problems that have been solved half a century ago. One idea of overcoming the problem of data requirements for RL is to perform training in simulation and transfer the resulting policy to real robotic systems. But, so far, it has been challenging to show that sim-to-real transfer could work in practice for any interesting applications

since, in most cases, RL trained policies overfit to the specifics of the simulators. Hence, in our work, we would like to demonstrate the utility of combining RL and simulators. For this purpose, we provide a framework that learns robust stabilizing quadrotor policies transferable to multiple real quadrotors with varying dynamics properties while requiring no data collection from the physical systems (see Chapter 4).

Despite the compelling generalization capability that we showed in the previous scenario, it would be unrealistic to assume that our results will work for all robotic systems. In practice, we need to, and often can, collect some limited (scarce) amount of data from the target robotic system. Hence, there is a lot of value in developing general methods capable of learning tasks specifics of the target system under scarce data availability. The development of such a method is the topic of our second scenario. We tackle this problem by following the route of *meta-learning* [Bechtle et al., 2021] (or learning to learn). The meta-learning family of methods allows for the leveraging abundant meta-training data (presumably from a simulator). These methods learn important optimization priors using samples from a distribution of possible tasks provided at *meta-train* time. Such learned priors enable sample-efficient task identification from scarce data available at *meta-test* time, i.e. the time of adaptation to a specific system. Specifically, we propose a meta-learning algorithm that captures task priors in the form of a parametric loss functions. Such losses perform shaping of the optimization objective landscapes to improve optimization robustness and sample efficiency. Compared to other meta-learning approaches, our learned loss functions are policy architecture agnostic. Furthermore, our framework represents an orthogonal direction to the existing meta-learning body of work. We apply our approach to a diverse set of problems, thus, demonstrating its flexibility and generality (see Chapter 5).

Part I

Perceptual Data Scarcity

Chapter 2

Overcoming Scarcity of Point-of-Contact Estimation with Tactile Sensors

Many manipulations in unstructured environments require a robot to use a grasped object, i.e., a tool, to interact with other objects. Often, a specific part of the tool, such as the tool's tip, needs to make contact with another object to perform the manipulation. For example, a hammer should make contact with a nail on the flat surface of its head when performing a hammering task. By sensing if and where the tool has made contact, the robot can verify that it is performing the skill correctly and otherwise adapt the skill accordingly.

Detecting transient contacts between a handheld object and other objects is particularly important when considering obstacles in the environment. For example, when placing a box on a cluttered table top, the robot may detect an unexpected contact on the side of the box rather than on

Parts of this chapter appeared in:

[•] Artem Molchanov, Oliver Kroemer, Zhe Su, and Gaurav S. Sukhatme. Contact localization on grasped objects using tactile sensing. In 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS, 2016, pages 216–222. URL: https://doi.org/10.1109/IROS.2016.7759058.

Artem Molchanov, Oliver Kroemer, Zhe Su, and Gaurav S. Sukhatme. Model-free contact localization for manipulated objects using biomimetic tactile sensors. In Humanoids 2016 Workshop on Tactile Sensing for Manipulation, 2016.

the bottom. In this case, the robot should immediately stop executing the skill to prevent damaging the object and either place it at a different location or move the obstacles.

The task of localizing contacts has usually been approached by using either vision or wristmounted force-torque sensors or joint torque sensors. However, both of these modalities have numerous problems making contact feedback scarce: visual information relevant to the contact estimation is inherently hard to acquire due to occlusions, changing lighting conditions, and fundamental ambiguity of depth perception. Approaches based on force-torque sensors face other challenges such as bias drift and, in some cases, they require an accurate model of the robot and the grasped object. Besides that, high cost and the narrow application of the force-torque sensors make them scarcely available in robotics. Thus, contact localization is a prime example for perceptual data scarcity because there is no sensor that can directly measure the desired quantity and the observation-state mapping is too complex to be modeled by hand-engineering accurately and effectively.

On the other hand, there is substantial evidence in neuroscience that tactile sensors are one of the main sources of information for various manipulation tasks [Johansson and Flanagan, 2009b, Goodwin and Wheat, 2004]. Thus, it is natural to attempt applying tactile sensors for the purpose of external contact estimation thus using them as a source of data in absence or scarcity of feedback from other sensors.

In practice, application of tactile sensors face a few challenges:

Lack of understanding. We only have very rough understanding of the way human skin generates and represents tactile signals, making it hard to replicate similar measurements in robotics. On top of that it is often unclear if tactile modalities generate enough information about relevant states which is particularly evident for our problem.

- **Poor standardization and lack of datasets.** Partially due to the lack of understanding, there exist numerous approaches and modalities for tactile perception, making it hard to reuse and standardize data. Due to the infancy of the field the sensors are also rarely available on existing robotic arms, which further complicates the problem.
- **Fragility of sensors.** Tactile sensors have to withstand constant physical contact with different objects. Unfortunately, due to early development stage of the latest tactile technologies their rigidity requires considerable improvement. It further limits amount of collectable information, types of experiments and approaches that one can explore.

Despite aforementioned problems recent developments in the field have made significant step toward human-like tactile perception [Wettels et al., 2008], by combining information about pressure, vibrations, and temperature. The large amount of data provided by these sensors could potentially result in significantly more robust manipulation skills for robots. However, in order to fulfill this potential, the robots will also need suitable estimation methods to process the tactile data.

In this chapter, we address the problem of contact localization and solve it using multi-modal tactile feedback. We leverage multiple tactile modalities provided by electrodes and hydro-acoustic sensors to structure the problem as spatio-temporal inference. We employ machine learning techniques to acquire the complex mapping between tactile sensors and the contact locations. Specifically, we investigate estimating the positions and directions of contact points using neural network (NN) classification and regression, Gaussian process (GP) regression, and support vector machine (SVM) classification with features learned using spatio-temporal hierarchical matching pursuit (ST-HMP) [Madry et al., 2014]. We evaluate the methods using data collected from 18 objects with different shapes, sizes and materials. The experiments were designed to provide

accurate ground truth information of the contact events and the contact parameters. In our work, we rely on a few assumptions. First, we restrict our contacts to the transient type, i.e., short taps, which correspond to tactile events such as bumping into objects or making contact. We also restrict our investigation to a single contact point and perform object dependent learning.

The key contributions of this chapter are: a) a model-free approach for estimating contact positions and directions between the environment and a grasped object based on tactile signals, b) an accurate labelled dataset* for evaluating and benchmarking contact localization methods, and c) an evaluation of the presented approaches using real robot experiments.

2.1 Related Work

The problem of contact detection and localization has received significant attention in the literature. The majority of the related approaches use force-torque sensors and analytical models to estimate the contact location. These methods often rely on accurate and precise calibrations of the force-torque sensors. For example, Karayiannidis et al. [2014] estimate a point of contact using first order differential kinematics in combination with force-torque measurements. This approach is only applicable to rigid grasps where the object cannot slip in the hand, as it assumes knowledge of the static center of mass of the object. Likar and Zlajpah [2014] present an approach that estimates the force and the point of contact from joint torque information. Their method requires exact knowledge of the robot's dynamics, which are often difficult to obtain for real robots. Some works do not estimate contact locations explicitly but rather use force-torque measurements to estimate other task-relevant states, such as the alignment errors in the peg-in-hole problem [Bruyninckx]

^{*}The dataset is available at http://bicl.robotics.usc.edu

et al., 1995] or the transitions between contact states for assembly tasks [Hovland and McCarragher, 1997].

Other methods exploit geometric models of the robot to estimate the points of contact of the manipulator with the environment. For example, Petrovskaya et al. [2007] use compliant motions in order to simultaneously estimate geometric parameters of the robot's links and the points of contact. The link parameters are estimated using a least squares approach while the contact points are inferred using a Bayesian approach. Koonjul et al. [2011] propose two geometry-based approaches. The first approach extends the self posture changeability method [Kaneko and Tanie, 1994] to use multiple compliant joints, while the second approach is a model-free method that maps joint displacements directly to the point of contact.

Tactile sensors have also been used for contact estimation, but most of the work has focused on estimating contact locations and interaction forces on robotic digits [Ciobanu et al., 2014, Su et al., 2015]. Tactile sensors have also been used to estimate other object properties, e.g., the object's pose or material properties. For example, Corcoran and Platt [2010] use particle filters to estimate the pose of an object based on its contacts with a robot hand. The pose of the object during manipulation is estimated by a measurement model which integrates the likelihood of contact measurements over the space of all possible contact positions on the surface of the object. Li et al. [2014] use vision-based tactile sensor, GelSight, to localize objects in a robot hand by matching key points between object height maps with RANSAC. Su et al. [2016b] use the bio-inspired structure of BioTac sensors to achieve high-sensitivity in estimating the orientation of the contacted object. In contrast to previous work on estimating the pose of grasped objects or localizing contacts on the surface of robotic digits, our work focuses on estimating the locations of contacts between grasped objects and the environment. Contact points can also be estimated using vision-based methods. Bernabé et al. [2013] propose a method that combines robot motion with point-cloud-based object tracking to estimate the location of contacts on an occupancy grid map. Hu et al. [2000] use vision features, such as binocular disparity, shadows, and inter-reflections, to detect imminent contact for manipulation tasks. Due to poor accuracy and ambiguity in contact localization caused by visual occlusions, vision-based methods are generally not meant to be used alone, but rather in combination with other sensors.

Different sensor modalities for estimating contact parameters provide complimentary strengths and limitations [Kroemer et al., 2011]. A robot can therefore often obtain a more accurate estimate by fusing the data from multiple sources. Felip et al. [2014] present a method that fuses multiple hypotheses from different modalities, including force-torque, tactile, and range sensors, to compute the likelihood of contact points. In their framework, pressure-sensitive tactile sensors on the hand are used to generate hypothesis for contacts between the object and the hand. Ishikawa et al. [1996] propose estimating a contact point by intersecting a force line acquired using force-torque sensor measurements and the plane containing the contact point extracted using a task specific vision system.

In this chapter, we investigate the problem of estimating the contact point between a grasped object and the environment using tactile sensors. This problem is challenging as the contact point is not made directly with the tactile sensor. The robot must therefore estimate the contact point based on the forces and vibrations transferred through the grasped object. Rather than relying on analytical models, we propose a model-free data driven approach to the problem.

Object-environment contacts have usually been estimated using force-torque data from wristmounted sensors. However, tactile sensing also provides important information for detecting and estimating these contacts [Johansson and Flanagan, 2009a]. Our proposed approach thus provides an alternative sensor modality for estimating contacts, which could be combined with other modalities in a sensor fusion framework. Similar to other approaches for estimating objectenvironment contacts, we assume a single point of contact. We plan on extending the presented approach to multiple contacts in the future.

2.2 Biomimetic Tactile Sensor

In our experiments, we use a haptically-enabled Barrett robot arm with a three-fingered Barrett hand. Each finger is equipped with a biomimetic tactile sensor (BioTac) [Wettels et al., 2008]. Each BioTac consists of a rigid core containing an array of 19 electrodes surrounded by an elastic skin, as illustrated in Fig. 2.1. The skin is inflated with an incompressible and conductive liquid.

The BioTac provides three complementary sensory modalities: force, pressure, and temperature. When the skin is in contact with an object, the liquid is displaced, resulting in distributed impedance changes in the electrode array. The impedance of each electrode depends on the local thickness of the liquid between the electrode and the skin. Micro-vibrations in the skin propagate through the fluid and are detected by the hydro-acoustic pressure sensor. The high- and low-frequency pressure vibration signals are referred to as PAC and PDC respectively. Temperature and heat flow are transduced by a thermistor near the surface of the rigid core. Since temperature conditions do not change in our experiments we do not consider this modality for our analysis.

2.3 Point-of-Contact Estimation

In this section, we describe our pipeline for estimating the point-of-contact from tactile data. Our contact learning pipeline consists of two main parts: contact detection and contact localization. In Section 2.3.1, we describe how the robot estimates the time of contact. The contact localization



Figure 2.1: Cross-sectional schematic of the BioTac sensor (adapted from Su et al. [2012]).

is subsequently performed using the tactile data around this time point, as described in Section 2.3.2. Contact point localization can be framed as either a regression or a classification problem. We show how the contacts can be localized using neural network or Gaussian process regression in Section 2.3.3. In Section 2.3.4, we explain how the contact parameters can be estimated using neural network or support vector machine classifiers.

2.3.1 Contact Detection

Similar to fast afferents in human skin [Johansson and Flanagan, 2009a], we detect contact events using the high frequency vibration signals extracted from the BioTac sensor. The robot uses a 5-th order high-pass Butterworth filter with a cut-off frequency of 20 Hz to remove biases. Contact event candidates are then extracted by detecting when the filtered pressure signal passes a threshold value. In our experiments, we used a threshold value of 500. In order to remove jitter, closely located events are reduced into a single event using the density based spatial clustering of applications with noise (DBSCAN) [Ester et al., 1996] algorithm. The resulting clusters in time dimension define the beginning and the end of every contact event.

2.3.2 Feature Extraction

The BioTacs' electrodes and hydro-acoustic pressure sensors provide multimodal tactile data for localizing contacts. To form a feature vector for localizing the contact point, we extract the values of the tactile signals after a contact event is detected. We use a window of t = 25 consecutive time steps at a sampling rate of 100Hz from the beginning of the event. The size of the window corresponds to the mean duration of the detected contact events plus two standard deviations. The resulting feature vector consists of the PAC, PDC, and 19 electrode signals collected from the three fingers and concatenated over the time window resulting in $s \times t = (2 + 19) \times 3 \times 25 = 1575$ features. In order to reduce the influence of gravity on the tactile readings, we subtract the average signal values of the three time steps immediately before the contact event.

2.3.3 Regression

We parametrize the location of the contact point using its Cartesian position (x, y, z) relative to the robot's palm. The direction of the contact's direction is parametrized using the yaw (ψ - rotation around z axis) and the *pitch* (ν - rotation around y axis). The wrist coordinate frame is shown in Fig. 2.6. We exclude the rotation around the contact's direction, i.e., roll, from the parametrization, as it is ambiguous for a single point-of-contact.

Our regression approach learns a separate mapping from the tactile features $\mathbf{d} \in \mathbb{R}^{st}$, described in Section 2.3.2, to each of the five continuous contact point parameters $\{x, y, z, \psi, \nu\}$, which we denote as $\hat{c} \in \mathbb{R}$. We define the function as

$$\hat{c} = f(\mathbf{d}, \theta), \tag{2.1}$$

where θ is a vector of function parameters that define the mapping from features to contact parameters. We compare two different machine learning techniques for learning the contact estimation function: neural networks and Gaussian processes.

The NN architecture consists of two fully connected hidden layers with 900 neurons each. The training of the NN is performed in a supervised manner using stochastic gradient descend with the Euclidean quadratic loss function

$$L_e = \frac{1}{2N} \sum_{i=1}^{N} (\hat{c}_n - c_n)^2, \qquad (2.2)$$

where $\hat{c}_n \in \mathbb{R}$ is the predicted value of a contact parameter, $c_n \in \mathbb{R}$ is the ground truth value of the parameter, and N is the number of samples in the training set. To avoid over-fitting, we use dropout for the hidden layers with a dropout ratio of 0.5, and we pick the best snapshot on the validation set observed during the training.

We compare NN regression with GP regression, a state-of-the-art non-parametric Bayesian supervised learning approach with automatic relevance determination (ARD). For the Gaussian process model we use a zero mean function, a squared exponential ARD kernel, and Gaussian noise. We also initialize the length scales, signal variance, and likelihood hyperparameters to 0, 0, and -2.3, respectively. Due to the high computational cost of GPs, we compute the average signal values during contact event window and use them as features, such that $\mathbf{d} \in \mathbb{R}^{s}$.

2.3.4 Classification

Regression is the standard methodology for learning mappings from features to continuous variables. However, the prediction of the contact parameters from the tactile data may be ambiguous, with a multi-modal distribution over the contact parameters. The regression approach would result in the robot averaging over the multiple possible contact parameters. Instead, we want the robot to select the most likely set of contact parameters. This can be achieved by framing the contact localization as a classification problem. In the following, we introduce a classification approach for the contact parameter estimation, which is able to represent ambiguities by producing a distribution over the possible contact parameter values.

In the classification approach, we represent the point-of-contact in the form of a distribution over the discretized contact pose parameters:

$$h(c) = p(c|\mathbf{d},\theta), \tag{2.3}$$

where $c \in \mathbb{Z}$ is a one dimensional random variable corresponding to discretized contact parameter values, $\mathbf{d} \in \mathbb{R}^{s\dot{t}}$ and θ is the feature vector and classifier parameters as described in 2.3.3.

For NN we use the same architecture for the estimator, however for classification the neural net output is converted to the distribution over labels using soft-max. Furthermore, we use RMSprop adaptive learning rate with cross-entropy classification loss for training:

$$L = -\frac{1}{N} \sum_{n=1}^{N} \ln\left(\hat{h}_{n}^{l}\right),$$
(2.4)

where l is the true (target) label, N is the total number of training samples, and \hat{h}_n^l is the predicted probability for the target label l of the *n*-th sample. The target labels are created by discretizing the continuous ground-truth data and assigning all values to the corresponding discrete bins. If a value is inside a particular bin, the corresponding target label receives the probability value 1 and all other are assigned probability 0.


Figure 2.2: The 18 objects used for data collection

Since the combined discretized parameter space produces a significant number of classes, every class would have only a few training samples given a limited amount of training data. This creates a significant obstacle for learning, since neural networks are typically prone to over-fitting. In order to mitigate this problem, we use separate classifiers for each of the five contact pose parameters.

We also apply linear support vector machines (SVMs) to perform the classification. The robot uses spatio-temporal hierarchical matching pursuit (ST-HMP) [Madry et al., 2014] to compute suitable features for the linear SVM. This feature learning framework has been successfully used for other tactile sensing applications, such as accurately predicting grasp stability [Chebotar et al., 2016] and detecting sensory goals using both static and dynamic tactile signals [Su et al., 2016a]. Details on how to apply ST-HMP to BioTacs can be found in the paper of Su et al. [2016a].

2.4 Evaluation

In this section we present the evaluation of the contact localization using the proposed methods. We describe our experimental robotics setup and perform empirical comparative analysis.



Figure 2.3: An example of a grasp used and contact samples collected during the experiment. Left: an example grasp of an object. Right: an example of the detected contact points and directions projected onto the zx plane.

2.4.1 Data Collection Setup

The experiments were performed using a three-fingered Barrett hand. Each finger tip is equipped with a BioTac tactile sensor.

Eighteen objects with a variety of sizes and materials were chosen for this experiment. All of the objects are shown in Fig. 2.2. The grasped objects' textures were modified to prevent slippages during tapping. During the data collection, the robot held one of the objects (see Fig. 2.3 for an example grasp) while a person tapped the object with a plastic rod shown in Fig. 2.5. Fig. 2.4 shows a normalized histogram of the peak tapping forces applied to the objects, which were measured using an ATI force-torque plate attached to the tapping rod. The steady state grip forces applied



Figure 2.4: The normalized histogram of the peak tapping forces applied to the object during the experiment.

by the robot and measured by the normal force on the BioTacs [Su et al., 2015], were 11.55N ($\pm 5.84N$) among 18 objects.

In order to determine the location and direction of the contact event, the rod was tracked using a marker-based Vicon tracking system, as shown in Fig. 2.5. The wrist of the robot was also tracked using Vicon markers, as shown in Fig. 2.6 with the corresponding coordinate frame. Using the coordinate frames of the wrist and the rod, the robot computed the location and direction parameters of the contact points relative to the robot's wrist.

The data from the contact events was recorded as a continuous time series. The BioTac sensor readings were then extracted using the contact detection method described in Section 2.3.1. Using this approach, the robot detected ≈ 15100 samples for all 18 objects. The number of taps applied to a particular object varied depending on the object's size.

In order to evaluate the contact localization 80% of the collected data was randomly selected to train the estimators using the methods described in Section 2.3. The test set consisted of 12% of the collected data. Since the neural nets are prone to overfitting to the data if trained for too long, the remaining 8% of the data was used as a validation set to continuously evaluate the learners'



Figure 2.5: The tapping rod. Left: an image of the rod. Right: rod coordinate frame with the locations of the markers.

generalization performance. Training continued until no improvement in validation error was observed during 350 consecutive iterations. The network with the smallest validation error was selected and evaluated on the test data. This evaluation process was repeated for all 18 objects.

Since classification requires a finite number of classes, we define an estimation area around the wrist coordinate frame as a rectangular box with dimensions: x = -20..20 cm; y = -10..10 cm; z = 0..30 cm; and discretize these dimensions using different grid sizes: 3 cm, 2.5 cm, 2 cm, 1.5 cm, 1 cm, 0.5 cm and 15°, 12.5°, 10°, 7.5°, 5°, 2.5° for Cartesian and angular coordinates respectively, resulting in 13..80 classes for the *x* dimension; 6..40 classes for the *y* dimension; 10..60 classes for the *z* dimension; 24..144 classes for the *yaw* and 12..72 classes for the *pitch* dimension for different grid sizes. The dimensions of the estimation area were picked to accommodate sizes of all objects. The minimal step size for the grid was limited by the Vicon system's tracking error.



Figure 2.6: The wrist band with Vicon markers. Left: CAD model of the Barrett hand with the band. Right: wrist coordinate frame with locations of the markers.

2.4.2 Results and Discussion

Fig. 2.7 shows the mean absolute error (MAE) calculated from errors of all 18 objects. We evaluated NN regression using electrode features (*NN:Electr* in the figure) and using the full feature set (*NN:Electr+PAC+PDC* in the figure). We also evaluated using Gaussian Processes regression using the full feature set (Fig. 2.7 *GP:Electr+PAC+PDC*). The three regression approaches resulted in considerable errors that, in some cases, exceed 50% of the object's size, although we did not observe overfitting in either of cases. Such significant errors are probably caused by ambiguities in the mapping between features and the estimated contact parameters, which can not be represented properly by the regression. These results motivated us to approach our problem from the point of classification, which is also more suitable for application of neural networks, where they traditionally show superior results over other machine learning techniques [Krizhevsky et al., 2012, Simonyan and Zisserman, 2014].

Similar to the regression approach, we also evaluated NN classifiers using only the electrode features and using the full set of features. For this experiment, we pick a 1cm/5° grid with 25 time steps as a baseline parameter set. Fig. 2.8 presents the MAE across the test sample sets of



Figure 2.7: Results of regression using different sensor modalities for NN and GP regressors. *NN:Electr* - NN with electrodes only; *NN:Electr*+*PAC*+*PDC* - NN with the full set of sensor modalities (electrodes, PAC, PDC); GP:Electr+*PAC*+*PDC* - GP with the full set of sensor modalities. x, y, z, yaw, pitch are results for independent dimensions. xyz represents results for combined error, i.e. the norms of the error vectors in Cartesian space. Errors in Cartesian space are reported in cm; angular errors are reported in degrees.

all 18 objects using our classification approach described in Section 2.3.4. In this case, MAE is calculated as the average of the absolute value of the difference between the real contact point parameters and the middle of the bin predicted by the classifier. We also combine predictions of individual dimensions in Cartesian coordinates for every sample in order to calculate the average Euclidean norm of the error vector for all location predictions. This prediction is denoted as *xyz*. The results indicate that electrodes are the most relevant features for contact localization. Incorporation of PAC and PDC injects additional noise and leads to overfitting, which is a typical problem in learning [Domingos, 2012]. Thus, for further investigation we restrict our experiments to electrodes only and apply an SVM classifier with ST-HMP features extracted from the array of electrodes (see Fig. 2.8, *ST-HMP:Electr*). The results show that NN classifier outperforms ST-HMP by ≈ 0.9 cm for the overall Euclidean error and up to 3° for angular coordinates. Given these results, we decided to use the NN classifiers for the following investigation.



Figure 2.8: Results of classification using different sets of sensor modalities: *NN:Electr* - NN classifier with electrodes only; *NN:Electr+PAC+PDC* - NN classifier with electrodes, PAC and PDC features; *ST-HMP:Electr* - ST-HMP feature learning algorithm with SVM classifier applied to electrode features.

The robot's reaction time to contacts can play a crucial role for some applications of the contact localization algorithms, e.g., if they are used in a control loop. Thus, it is important to understand what contributes the most to the estimation delays. For our algorithms, the most significant source of delays is the data accumulation. For example, 25 time points with 100 Hz frequency causes a delay of 250 ms, whereas running forward pass of the NN requires less then 5 ms on a Core i5 CPU. It is therefore interesting to see how the prediction error varies with respect to the number of time steps of the features collected. Fig. 2.9 shows the classification MAE for the NN given different numbers of time steps retained from the moment the contact event is detected. As one can see from the figure, the algorithm is quite robust to the time window size, and its performance degrades gracefully until five time steps. Below five time steps, the performance starts dropping. The drop in performance can partially be attributed to imperfections in the detection pipeline, which introduce variable time shifting of the signals relative to the position inside the time window.



Figure 2.9: Results of NN classification using different time lengths of features, starting from the moment contact event is detected.

In all previous experiments we used a 1 cm/5° grid as a baseline parameter. Thus, it would be interesting to see how sensitive our classifiers are to the grid resolution used to discretize the contact parameter estimations. To investigate that effect, we vary our grid sizes from $0.5 \text{ cm}/2.5^{\circ}$ to 3 cm/15° with steps of $0.5 \text{ cm}/2.5^{\circ}$ and train NN classifiers. Our results, as shown in Fig. 2.10, indicate that the MAE does not change significantly and it does not exceed 3 cm/11°. This result means that even if the classifiers cannot guess the exact label they usually predict one of the adjacent bins, which indicates that they learn the underlying relations between different classes. This can also be seen from an example confusion matrix for the *x* dimension presented in Fig. 2.11. In the confusion matrix, the rows represent the ground truth (target) classes and the columns represent the predicted classes. In most cases the recall is quite high and misclassified labels cluster around the target class. The results also clearly indicate that classifiers produce almost no predictions for the classes that were not presented for the training, thus, making it safe to preallocate larger prediction areas if needed and keep the architecture of the classifiers the same while incorporating new data samples into learning.



Figure 2.10: Results of NN classification using different grid steps.

2.5 Conclusion

In this chapter, we demonstrate how we can overcome scarcity of point-of-contact estimation with external objects using BioTac tactile sensors. We leverage spatio-temporal structure of the problem and employ modern machine learning techniques to learn complex mapping from different tactile sensory modalities to the contact locations and orientations. In particular, our pipeline consists of contact detection and contact localization. The detection is performed by filtering and thresholding the BioTacs' pressure signals. Contact localization is performed by applying different machine learning methods, including neural networks, Gaussian processes, and support vector machines with ST-HMP feature learning, to tactile data. We frame the contact localization as both a regression and a classification problem and investigate sensitivity of our algorithms to time and component-wise changes in the input features, as well as various discretizations of the parameter space for our classification approach.

We evaluate our methods using hundreds of contact events from eighteen objects with different shapes and material properties. Our classification approach results in the best performance, with



Figure 2.11: An example confusion matrix of the *x* dimension classification for one of the objects.

expected localization errors less than 2.5 cm/10° for individual objects and poses. Our results clearly show that BioTac sensor is a rich source of indirect information for contact localization. It can provide valuable sensory modalities to overcome contact information scarcity even in the complete absence of alternative ways of contact estimation.

Chapter 3

Aggregating Multiple Active Drifters Under Scarcity Of Current Measurements

In many robotics scenarios, measurements are typically considered cost-free. However, in some cases acquiring measurements might induce an additional cost in terms of time or energy. For example, robots often have to apply actions, such as to change their configuration or even the scene in order to receive desired observations. Thus, this kind of perceptual data scarcity leads to a careful trade-off between acquiring measurements/observations and performing the desired task.

We look at this problem in the scope of an extremely challenging scenario of controlling Lagrangian drifters. Lagrangian drifters are monitoring devices that are used by oceanographers and biologists to track ocean currents and measure water characteristics. In this chapter, we are

Parts of this chapter appeared in:

Artem Molchanov, Andreas Breitenmoser, and Gaurav S. Sukhatme. Active drifters: Towards a practical multi-robot system for ocean monitoring. In IEEE International Conference on Robotics and Automation, ICRA, 2015, pages 545–552. URL: https://doi.org/10.1109/ICRA.2015.7139232

Artem Molchanov, Andreas Breitenmoser, and Gaurav S. Sukhatme. Active drifters: Sailing with the ocean currents. In RSS Workshop on Autonomous Control, Adaptation, and Learning for Underwater Vehicles, 2014.

concerned with drifters that are composed of a surface float and a tethered drogue that acts as an "underwater sail" * (see Fig. 3.1).



Figure 3.1: Lagrangian drifters. Left: A prototype of our passive drifter system, which was deployed in the Southern California Bight to measure ocean currents. Right: The schematic shows the main components and the mode of operation of an active drifter.

Drifters are underactuated, and usually passive systems. Positioning its drogue at a fixed depth causes a drifter to travel passively with the current at that depth. This is a common technique in oceanography to tag and track currents (and all that they carry with them).

To achieve better spatial resolution, multiple drifters are commonly deployed. There are two main challenges when it comes to deploying multiple drifters in coastal regions. The first is that they do not tend to provide uniform sampling resolution close to the shore. The second is that they are expensive (in terms of ship time) to retrieve at the end of the mission (being widely dispersed).

Here we study *active drifters* (see Sec. 3.2.1 for more details), specifically those with a single actuator that adjusts the drogue depth. Changing the drogue depth permits the *in situ* measurement

^{*}See http://www.pacificgyre.com for an example of commercially available drifters today.

and estimation of ocean current velocity at varying depths. This opens up the possibility of gaining (limited) control of drifter motion since a drifter with a depth-adjustable drogue can actively select the best ocean current for propulsion that achieves some high level mission goal (e.g., aggregation).

Controlling such vehicles is a very challenging problem for two reasons: First, changing the depth of the drogue is a slow and time consuming process. During this process the drifter is constantly involved in a random motion caused by the currents at intermediate depths. Thus, every decision that the drifter has to make comes at considerable price which results in scarcity of current measurements. Second, ocean currents are the only drifters' propulsive force. They are dynamic and extremely difficult to predict. Thus, active drifters are an example of a highly underactuated robotic system with external forcing. Other single and multi-robot systems, including underwater vehicles [Pereira et al., 2013] and aerial robots [Desaraju and Michael, 2014, Wolf et al., 2010], are exposed to external forces in real-world applications. In such settings, tasks such as navigation, station keeping, or formation maintenance are very challenging.

In this chapter, we consider a multi-drifter system, with two mission objectives. The first objective is *coverage* (drifters are required to spread out) and the second is *aggregation* (drifters are required to cluster near each other). We present a coupled control and current sampling method for active drifters that offers a solution to these two classical multi-robot problems under scarcity of current estimation and severe underactuation from external forcing induced by the ocean currents. Spreading out enhances the deployment process that drifters naturally experience in the ocean. Aggregation offers the practical benefit that a recovery vessel does not need to search for and pick up each individual drifter at disparate locations at the end of a monitoring mission. It can easily collect the aggregated drifters by visiting few clusters of drifters, which reduces ship operation cost. We report here on the control design and a simulation-based feasibility study to inform the design of a practical active multi-drifter system.

Our results suggest that a practical implementation of drifters with our method could be made to aggregate and disperse in the coastal ocean with relatively inexpensive components. We are able to show that after deployment a significant fraction of drifters can be aggregated in few clusters over a 90 day period (see 3.3.3), which greatly facilitates the recovery of the deployed drifters.

We organize this chapter as follows. First, in Section 3.2 we introduce Regional Ocean Modeling System (ROMS) and perform preliminary investigation of different control strategies for a single drifter in an idealistic scenario of no-cost current measurements. Then, in Section 3.3, we investigate our target problem of controlling multiple drifters under expensive ocean current measurements.

3.1 Related Work

Active drifter systems have been the subject of recent study. Regarding the mode of operation, the system we study here is closest to the system described in [Dunbabin, 2012a], which can raise and lower a drogue via a winch. An alternative approach uses a free-floating submerged drifter [Han et al., 2010] wherein the entire drifter body dives to a certain depth by changing its buoyancy instead of sitting at the surface and controlling a drogue on a tether. Argo floats [Gould et al., 2004] are larger Lagrangian profilers, which are in wide use in the ocean today. They also adjust their depth by buoyancy control and could theoretically be operated as active drifter systems [Smith and Huynh, 2013], though in practice they are not operated as such today.

In terms of active drifter control, there are two main approaches. The first is a predictive control approach which explicitly relies on predictions of the ocean currents based on an ocean model [Dunbabin, 2012a, Smith and Huynh, 2013]. The underlying assumption is that the predictions by the ocean models are reliable, which is not always the case 3.2.2. We follow an

alternative approach where each drifter uses in situ measurements of ocean currents to make control decisions.

A control strategy for rendezvous with multiple drifters is presented in [Ouimet and Cortés, 2014]. Although their targeted application of aggregating multiple drifters is similar, they use a different approach, where the ocean dynamics is represented by internal wave models with known parameters. In contrast, we assume no explicit knowledge about the dynamics of the ocean, except the vertical component of the flow, which we assume to be zero. In Jouffroy et al. [2013b] the target application is to control the absolute position of a drifter in a coastal scenario, where the drifter can anchor itself at the sea bottom if necessary. As an extension, the deep ocean scenario for a single drifter is examined under idealized conditions (the drifter is assumed to have instant estimates of currents at the present location, the ocean currents are assumed to be stationary and to span the plane positively). The external forcing by current flow fields and its impact on the control of underactuated robotic systems is furthermore addressed by Kwok and Martínez [2010] for a coverage task with self-propelled vehicles of bounded velocity in a river environment and by Michini et al. [2014] for tasks of tracking coherent structures on flows with autonomous underwater vehicles in the ocean.

3.2 Single Drifter Control

In this section we aim to provide preliminary investigation of the quality of ROMS model for long horizon predictions and possible control strategies of a single drifter in idealistic scenarios.

3.2.1 Regional Ocean Modeling System for Active Drifter Control

Any underwater or surface vehicle operating in the ocean is exposed to ocean currents. The currents are generally treated as noise which perturbs the vehicles' trajectories. In the case of underactuated and rather passive systems like passive or active drifters, the currents however act as the main driving force and can be seen as the primary component of a controller. In order to anticipate the way the vehicles are affected by the currents and design a control policy that actuates the drogue of an active drifter to reach a desired destination, it is essential to obtain good estimates of the ocean currents at all times.

Due to the large spatial scale of the ocean, it is hard to acquire such estimations at decent resolutions from standard ocean measurement tools, such as moorings, HF radars and satellite data, solely. Another route that has been explored in recent works (e.g., Smith and Huynh [2013]) is to utilize predictions based on ocean models, for example, using the Regional Ocean Modeling System (ROMS), possibly enhanced by the additional assimilation of real data [†]. However, many ocean phenomena are not yet completely understood, and as we show in this section, these models are often not very accurate and have a rather low resolution. In order to estimate the current vectors directly on site, our approach tries sampling the vectors locally at different drogue depths. For example, in our present drifter prototype, the vectors can be measured via the drifter's motion from two successive GPS locations. The approach does not require any prior information about the ocean, which may simplify the deployment of the system and potentially enable its use in arbitrary areas of the ocean.

[†]For details, see https://science.jpl.nasa.gov/projects/OurOcean/

3.2.2 Predictions and Measurements of Currents

In order to evaluate how well real and predicted ocean currents match, we deployed a passive drifter with drogue fixed at 3 m depth in the Southern California Bight near the coast of Los Angeles over 2 days (see Fig. 3.1 on the left). We operate at a local scale within kilometer range, which is below the minimum resolution of ROMS of 3 km × 3 km; hence a few relevant ROMS data points are available only. Although ROMS is a valuable tool at larger scales, the recorded data of Fig. 3.2 indicate that locally ROMS predictions often deviate significantly from the measured currents [‡]. This is especially well demonstrated by Fig. 3.3, which depicts the trajectory followed by the deployed drifter and the trajectories predicted by ROMS [§]. More detailed analysis of the data from the deployed drifter (see Fig. 3.4 and Fig. 3.5) shows that ROMS predictions and the real currents have weak positive correlation in the direction (the correlation coefficient is equal to -0.35). We see that as an indication of relatively poor consistency between ROMS predictions and the real in situ measurements. At the same time, nowcast and forecast predictions seem to be consistent with one another, i.e., they are highly correlated.

3.2.3 Controlling a Single Active Drifter

Given the above results of our field experiments, we aim at an approach that does not rely on ocean current predictions primarily, which is different from most former related works (Dunbabin [2012b], Smith and Huynh [2013] among others). The idea is to estimate currents in situ and use these estimates to design a reactive control policy that actuates the drogue to select a favorable

^{*}This presents an interesting direction for future research on how to combine on-line data from active drifters with ocean models like ROMS to further improve drifter navigation, as well as the ocean models themselves.

[§]The trajectories are generated using the "drop a drifter" web page. See http://www.cencoos.org/sections/models/ roms/ca/drifter/.



Figure 3.2: Comparison of ocean current measurements and ROMS predictions: the current forecasts (green vectors) and nowcasts (red vectors) are plotted along the measured trajectory and estimated currents of the deployed passive drifter (blue dotted line and blue vectors). The black circle in the middle represents the center of a cell of the grid that underlies ROMS.

current to drive the active drifter toward the destination. One can apply the following control

policies, which all use the measurements of the currents at the present drifter position \P :

Maximum projection (PRJ) . Select the current whose vector produces the largest projection

onto the axis of sight, which is the axis originating at the present position of the drifter and

passing through the destination.

[¶]Although we do not use ROMS predictions in the control policies, we rely on ROMS in our simulations as a generator of realistic ocean currents.



Figure 3.3: Comparison of ocean current measurements and ROMS predictions: trajectories generated by the "drop a drifter" web page, and the real trajectory of the drifter deployed in the Southern California Bight.

- Minimum distance (DIST) . Select the current for which the prediction of one step (or multiple steps) ahead results in a new drifter position with the smallest predicted distance to the destination.
- **Minimum angle (JF)** . Select the current whose direction is "closest" to the axis of sight, i.e., with the smallest angle between the current vector and the axis of sight. A similar idea was proposed by Jouffroy et al. [2013a].



Figure 3.4: Comparison of ocean current measurements and ROMS predictions: absolute values of the ocean current velocity vectors.

Since drifters do not have their own propulsion, the controllability of active drifters depends heavily on the set of currents present at a location. As shown in Jouffroy et al. [2013a], current vectors of different depth must span the plane positively everywhere to guarantee controllability of the drifter. Although an active drifter is only partially controllable and cannot reach every destination, it usually still performs better than a drifter with no control. For instance, Fig. 3.6 and Fig. 3.7 show the situation where 7 out of 8 destinations are unreachable for passive drifters (independent of their drogue depth), but the active drifter (here with minimum angle policy) manages to hit at least 5 of them.

Fig. 3.8 presents a comparison of the performance of the three different control policies using the minimum distance to a destination as a metric. The simulations were carried out in a dynamic ocean current field generated from ROMS data. The simulations show that the DIST and the PRJ controllers have very similar performance and statistically perform just slightly better than the JF



Figure 3.5: Comparison of ocean current measurements and ROMS predictions: direction angles of ocean current velocity vectors. The bottom part displays the cosine of the angle between the vectors of the ROMS forecast and the measurement, which visualizes their alignment (1: the same direction, 0: perpendicular, -1: opposite direction).

control policy. At the same time, all three controllers result in significant increase in performance compared to a passive drifter with a random choice of depth of the drogue. In this particular simulation, we observed 2 to 4 times better performance compared to a drifter with no control (a passive drifter) with the given metric.

3.2.4 Conclusion

In this section we provide preliminary investigation on the quality of ROMS model for long horizon predictions and possible control strategies of a single drifter in idealistic scenarios. Due to underactuation of the drifters and the ocean's chaotic and highly unpredictable dynamics, the problem of predicting and controlling the trajectories of drifters in the ocean is extremely challenging. Our field experiments reveal that ROMS is not sufficiently reliable for exact prediction



Figure 3.6: Simulated trajectories (colored curves) of passive drifters with varying drogue depths starting from the same location (green circle). From left to right, the pictures depict the evolution of trajectories over time (2 days, 30 days correspondingly). The colored crosses represent goal targets. The passive drifters managed to hit only 1 target out of 8. The simulation is run in a static ocean current field generated from ROMS data.



Figure 3.7: Simulated trajectories (colored curves) of active drifters starting from the same location (green circle). From left to right, the pictures depict the evolution of trajectories over time (2 days, 30 days correspondingly). The colored crosses represent goal targets. The color of the target corresponds to the color of the drifter it was assigned to. The active drifters managed to hit 5 targets out of 8. The simulation is run in a static ocean current field generated from ROMS data.



Figure 3.8: The mean of the minimum distance to a destination vs. initial distance to a destination. Each graph is made of 6 discrete data points with different initial distances. Each point is the mean over 1000 simulated trials with the same initial distance. In each trial, a drifter was given the task to go from a random initial point to a destination picked randomly, but with the predefined distance between them. The minimum distance achieved by the drifter in each trial was recorded. The duration of one trial is 180 days.

of currents thus making it hard to use ROMS for long horizon path planning. Given that we choose to utilize reactive control strategies based on in situ measurements and propose several approaches. In the next section, we will investigate the target scenario of a multi-drifter system in the presence of expensive current measurements.

3.3 Multi-Drifter Control

3.3.1 Problem Formulation

Our representation of the ocean relies on ocean current velocity vectors which are changing over time. This defines the time-varying flow field $\mathbf{f}: \mathbb{R}^3 \times \mathbb{R}_{\geq 0} \to \mathbb{R}^3$, with vectors

$$\mathbf{f}(\mathbf{x},t) = (f_{\mathbf{x}}(\mathbf{x},t), f_{\mathbf{y}}(\mathbf{x},t), f_{\mathbf{z}}(\mathbf{x},t))^{\mathrm{T}}.$$
(3.1)

 $f_{\mathbf{x}}(\mathbf{x}, t)$ and $f_{\mathbf{y}}(\mathbf{x}, t)$ are the horizontal components of the velocity vectors in east and north direction. *t* denotes the time in the model and $\mathbf{x} = (x, y, z)^{\mathrm{T}}$ is the 3D position. In what follows, we will neglect the vertical flow component and assume $f_{\mathbf{z}}(\mathbf{x}, t) = 0$ everywhere.

In order to simulate realistic dynamics of the ocean, we obtain the flow field from the Regional Ocean Modeling System (ROMS) [Shchepetkin and McWilliams, 2005]^{\parallel}. ROMS serves as our "ocean simulator": it generates velocity vectors for discrete time, with temporal resolution of 1 hour over a grid with spatial resolution of 3 km × 3 km. We interpolate the velocity vectors over the ROMS grid to retrieve continuity.

Given a group of N active drifters, each drifter's state is described by the 2D position of its surface float, $\mathbf{p}_i = (x_i, y_i)^T \in \mathbb{R}^2$, and the position of its tethered drogue at depth z_i , $i \in \{1, ..., N\}$. In this treatment we assume that the drogue has the same horizontal position as the surface float. In practice this will not be the case, but the horizontal offset between the two will not be significant relative to the size of the coverage area. We further assume absolute localization and global communication capabilities of the drifters within a centralized network (e.g., using GPS and satellite communication via a receiver on the surface float), such that the drifters can measure their

[&]quot;See also http://ourocean.jpl.nasa.gov/.

positions with sufficient accuracy and are able to exchange data among each other via a central base station if needed.

We model the drifters as first-order systems with Lagrangian dynamics,

$$\dot{\mathbf{p}}_i = \left(f_{\mathbf{x}}(\mathbf{x}_i(t), t), f_{\mathbf{y}}(\mathbf{x}_i(t), t)\right)^{\mathrm{T}} \text{ and } \dot{z}_i = u_i(t), \qquad (3.2)$$

with $\mathbf{x}_i(t) = (\mathbf{p}_i(t)^{\mathrm{T}}, z_i(t))^{\mathrm{T}}$ and control input

$$u_{i}(t) = \begin{cases} v & , z_{i}^{*} > z_{i}(t) \\ 0 & , z_{i}^{*} = z_{i}(t) \\ -v & , z_{i}^{*} < z_{i}(t) , \end{cases}$$
(3.3)

where z_i^* denotes the desired depth, and v represents the constant vertical speed for actuating the drogue.

Given the above ocean and drifter models, our goal is to control the drogues of N drifters in such a way that the drifters manage to 1) spread out inside a preference area A, and 2) aggregate within the overall mission area Ω , where $A \subset \Omega \subset \mathbb{R}^2$, despite the external forcing of the ocean currents. In other words, with respect to spreading out, the drifters should finally cover A in the limit by incrementally maximizing the minimum distance between drifters contained in A,

$$\max_{\mathbf{u}(t)} \left(\min_{\mathbf{p}_i, \mathbf{p}_j \in A, \ i, j \in \{1, \dots, N\}, \ i \neq j} \| \mathbf{d}_{i, j} \|_2 \right),$$
(3.4)

where $\mathbf{u}(t) = (u_1(t), \dots, u_i(t), \dots, u_N(t))$ and $\mathbf{d}_{i,j} = \mathbf{p}_j - \mathbf{p}_i$. Similarly, with respect to aggregation, the drifters are required to converge to clusters within Ω , i.e., the distances between the drifters $\|\mathbf{d}_{i,j}\|_2$ are minimized over time. This (reverse) process results in the maximization

of an *aggregation metric*. We delay the introduction of this metric till Section 3.3.3 where the quantitative evaluation of the active drifter system is discussed.

3.3.2 Drifter Control Methods

We design a control law that utilizes in situ measurements of ocean currents collected by the drifter. Our control method consists of two layers: a *high-level controller* that generates the desired instantaneous motion direction for each drifter and a *low-level controller* (a tracking controller) that selects a depth at which the ocean current best causes the drifter to move along the desired direction.

High-Level Control Law

The high-level controller is a potential field controller, which generates a unit direction vector \mathbf{F}_i for a drifter *i*,

$$\tilde{\mathbf{F}}_{i} = \frac{\mathbf{F}_{i}}{\|\mathbf{F}_{i}\|_{2}} , \text{ with } \mathbf{F}_{i} = \sum_{j=1, j \neq i}^{n} \mathbf{F}_{d_{i,j}} + \mathbf{C}_{i} , \qquad (3.5)$$

and

$$\mathbf{F}_{d_{i,j}} = \begin{cases} \xi \left(\frac{r_d}{\|\mathbf{d}_{i,j}\|_2}\right)^2 \tilde{\mathbf{d}}_{i,j} &, \|\mathbf{d}_{i,j}\|_2 > d_{\min} \\ 0 &, \|\mathbf{d}_{i,j}\|_2 \le d_{\min} \end{cases}$$
(3.6)

$$\mathbf{C}_{i} = \left(\frac{\|\mathbf{d}_{i,c}\|_{2}}{r_{c}}\right)^{w_{c}} \tilde{\mathbf{d}}_{i,c}, \qquad (3.7)$$

where $\mathbf{F}_{d_{i,j}}$ is the interaction force which attracts (or repels) drifter *i* toward (or from) drifter *j*, $\tilde{\mathbf{d}}_{i,j}$ and $\tilde{\mathbf{d}}_{i,c}$ are unit direction vectors, d_{\min} is the user-defined minimum distance of interaction between two drifters, and r_d is the radius at which the 2-norm of the interaction force of the drifters equals 1. $\xi = 1$ holds when a drifter is in the "aggregation" mode (attraction), and $\xi = -1$ when a drifter is in the "spreading" mode (repulsion). C_i is the force of attraction toward the center of the preference area A defined by the central point \mathbf{p}_c and the radius r_c . C_i should not disturb the interaction of drifters inside the preference area, but, at the same time, it should dominate outside the area. We achieve this by setting w_c to a high value. As the high-level controller generates only the direction for the motion, we are not interested in the magnitude of the force and normalize the vector \mathbf{F}_i to obtain the unit direction vector \mathbf{F}_i .

Low-Level Control Law

The purpose of the low-level controller is to select the best depth at which the propulsion due to the ocean currents will cause the drifter to move in the direction designated by the high-level control law. In our work, we consider a discrete set of depths. The choice of these depths is non-trivial and is explained in Section 3.3.3; it involves *sampling*, i.e., making measurements of currents at various depths, and *decision making*, i.e., choosing a particular depth.

With the drogue positioned at a particular depth for a small duration, the GPS on the surface is able to measure a change in drifter position and hence estimate motion locally. We ascribe this motion entirely due to the propulsive force at that depth, thereby estimating the ocean current at the present drogue depth. The low-level controller projects the ocean current vector onto the desired direction $\tilde{\mathbf{F}}_i$ (generated by the high-level controller) to obtain the *sample quality Q* of the ocean current that has effectively been sampled at the present depth. The sample quality thus measures the fitness of the ocean current at a particular depth.

The condition that triggers a control decision event is defined as

$$|\alpha - \alpha_{\text{prev}}| > \alpha_{\text{tol}}, \qquad (3.8)$$

where α is the angle between the latest estimate of the current the drifter is traveling with and the direction vector $\tilde{\mathbf{F}}_i$, α_{prev} corresponds to α for the drogue depth at which the previous control decision was made, and α_{tol} is the *tolerance angle* (a threshold value). Thus, the low-level controller makes a new decision every time the situation changes "significantly enough", which is defined by the tolerance angle.

The basic low-level controller employs a simple strategy. It periodically estimates the direction of the current it is traveling with. If a control decision event is triggered, the drogue starts sampling by cycling through a discrete range of depths. At each depth, the ocean current is estimated. Once all estimates are made, the controller picks the depth with the highest value of Q, and saves the present value of α (at the depth chosen by the controller) as α_{prev} .

This basic controller has two limitations. First, it does not reuse ocean current estimates. All the estimates from the previous control event are treated as outdated every time a new decision is required, even though, according to (3.8), in some cases a control decision event may be triggered due to frequent changes in the desired direction. For example, such behavior is typically observed when the drifters are aggregated in a cluster. Second, it does not leverage the aggregation of drifters. Accordingly, we extend the basic controller as follows to improve on performance.

In order to reuse previous ocean current estimates, we introduce the *trust time* T_{trust} and associate time stamps with ocean current estimates. T_{trust} defines how long an estimate can be used before it is outdated. The sampling strategy is modified as follows. Every time a drifter needs to make a new control decision, it re-samples only at the depths where its ocean current estimates are outdated. A decision is made once the estimates at all depths in the set are up to date.

The second extension uses samples acquired by other drifters in the cluster and induces *collaboration* among closely located drifters. Every time a new control decision is required, a drifter creates a pool of ocean current estimates, which consists of estimates acquired by itself and

those acquired from neighboring drifters that are within the data exchange distance. We choose the data exchange distance for our system to be the same as the minimum distance of interaction d_{\min} . For every depth value, the estimates are sorted according to their time stamp (the most recently acquired estimates are put first). The drifter puts the latest estimate for every depth value into the *best set*. Next, the drifter moves the drogue to only the depths for which the estimates are outdated in the best set. The best set is updated every time a new estimate is made. Once all estimates in the best set are up to date, the new control decision is made. This process is asynchronous and stochastic. It is asynchronous in the sense that current estimates can be shared as available—there are no preset communication slots. The clocks on the drifters are GPS equipped. The process is stochastic in the sense that there is no centralized current estimation task assignment, and, since all drifters act asynchronously, ocean current estimation happens randomly.

For the practical application, the high-level controller can be implemented as a centralized controller at the base station and a copy of the low-level controller resides on each drifter. Each drifter communicates its coordinates to the base station, and receives a desired direction vector back from the base station. Additionally, since drifters collaborate only when they are in close proximity to each other, we believe that the exchange of ocean current estimates can be done via RF modems in a practical system.

Both of these extensions we see as the way to exploit indirect information about ocean current measurements. It is indirect in the sense that in both cases we exploit measurements that are spatially and temporally separated from the true current values in the area.

3.3.3 Performance Evaluation

Aggregation Metric

In order to numerically evaluate the performance of the control algorithms, we propose a metric similar to a normalized pairwise potential energy of points,

$$M_a(D, n, N) = \frac{\left(\sum_{i=1}^n \sum_{j=1}^n U(\|\mathbf{d}_{i,j}\|_2)\right) - n}{N^2 - N},$$
(3.9)

where N is the initial number of drifters and n denotes the number of remaining drifters still left in the mission area Ω at a given time. $D \in \mathbb{R}^{n^2}$ is a square matrix of horizontal distances between the n remaining drifters, and $\|\mathbf{d}_{i,j}\|_2$ are elements of this matrix. Thus, since the normalization factor depends on N but the sum of the potential energies of individual pairs depends on n, the aggregation metric penalizes the fact that drifters leave the mission area (such drifters are considered to be lost). The function U is given as

$$U(d) = f_s(d, a_u), (3.10)$$

where the tuple a_u is a special instance of the parameter tuple $a = (a_{00}, a_{11}, a_{12}, a_{13}, a_{21}, a_{22}, a_{23})$ of the logistic function f_s for the potential field. The logistic function f_s is defined as

$$f_s(x,a) = \frac{1}{a_{00}} \left(\frac{a_{11}}{1 + e^{a_{12}(x - a_{13})}} + \frac{a_{21}}{1 + e^{a_{22}(x - a_{23})}} \right).$$
(3.11)

The parameters a_u are chosen such that the function has the highest drop approximately between 20 000 and 50 000, which corresponds to distance in meters. This zone is a transition area, once drifters get into this zone the metric starts growing faster, reflecting the fact that the drifters are



Figure 3.9: The artificial potential energy function U(d) for tuple a_u . The function defines the quality of aggregation of a pair of drifters.

approaching aggregation. Fig. 3.9 shows the graph of U(d) for our choice of a_u with parameter values $a_{00} = 1.77203$, $a_{11} = 1$, $a_{12} = 0.0002$, $a_{13} = 35\,000$, $a_{21} = 1$, $a_{22} = 0.000035$, $a_{23} = 35\,000$.

 M_a is normalized, thus it assumes 1 when all drifters are aggregated at one point (independent of the number of drifters). In order to determine the approximate number of clusters of drifters, n_{clust} , we use the heuristic rule

$$n_{\text{clust}} \approx 1/M_a \,, \tag{3.12}$$

which agrees with our observations in simulations.

ROMS Parameters

In our simulations, we use a discrete ROMS model with interpolation to the nearest grid point in order to generate a realistic ocean flow field. Since there is a cost associated with estimating currents at each depth, in practice we want the controller to work with as few depths as possible. An active drifter is controllable if the currents in the area span the plane positively Jouffroy et al. [2013b]. This criterion can be reformulated as follows

$$\max_{i,j}(\theta_{i,j}(\mathbf{p}_l)) < \pi, \qquad (3.13)$$

where $\theta_{i,j}(\mathbf{p}_l)$ is the angle (in the horizontal plane) between adjacent current vectors at depths i, j at the point \mathbf{p}_l . This means that if we project current vectors with the same horizontal initial point \mathbf{p}_l onto the horizontal plane, and the largest angle between these vectors is less than π , the system is controllable. This rule also implies that the minimum number of depths necessary is 3 (i.e., two vectors cannot span the plane positively). Thus, in order to find the appropriate depths we introduce the integral fitness function

$$R = \sum_{k=1}^{T} \sum_{l=1}^{P} f_a(\theta_{\max 1}(f'_{k,l}), \theta_{\max 2}(f'_{k,l})), \qquad (3.14)$$

where *R* depends on the ROMS dataset represented by *P* spatial points (in our case it is a mission area) and *T* time points (in our case it is a number of hours) in a given time period of the dataset. $f'_{k,l} = (\mathbf{f}((\mathbf{p}_l^T, z_1)^T, t_k), \mathbf{f}((\mathbf{p}_l^T, z_2)^T, t_k), \mathbf{f}((\mathbf{p}_l^T, z_3)^T, t_k)))$ is a set of three ocean current vectors at a given time t_k and given horizontal point \mathbf{p}_l , but at different depths z_1, z_2, z_3 . θ_{max1} is the largest angle of a triple of angles formed by the current vectors at different depths at the same $\mathbf{p}_l, \theta_{\text{max2}}$ is the second largest angle. The point-wise fitness $f_a(\theta_{\text{max1}}, \theta_{\text{max2}})$ is now given by

$$f_a(\theta_{\max 1}, \theta_{\max 2}) = f_s(\frac{\theta_{\max 1}}{\pi}, a_a) + (\frac{\pi - \theta_{\max 2}}{\pi}) f_s(\frac{\theta_{\max 2}}{\pi}, a_a) - \Delta f_s, \qquad (3.15)$$

where Δf_s is bias compensation, such that the point-wise fitness is 0 when all three vectors are collinear, a_a is a special instance of the parameter tuple a of the logistic function f_s presented



Figure 3.10: ROMS point-wise fitness function $f_a(\theta_{max1}, \theta_{max2})$. The function defines how well the plane is spanned at a point by a triple of vectors, based on the angles between each pair of vectors.

earlier in (3.11). Fig. 3.10 depicts the point-wise fitness function $f_a(\theta_{\max 1}, \theta_{\max 2})$ in the space of its arguments. As one can see, the function has its maximum value when $\theta_{\max 1} = \theta_{\max 2} = 120^\circ$, and has the highest derivative in the area where $\theta_{\max 1} \approx 180^\circ$ with minor influence of $\theta_{\max 2}$. f_a is also used to plot fitness maps, which are color maps (with a color corresponding to the value of the function at the point) acquired by summation of point-wise fitness values over time at each point. Such fitness maps allow visualization of the points where currents span the plane positively. These fitness maps are similar to the controllability heat maps introduced in Smith and Huynh [2013]. The fitness function presented above allows us to compare and pick the most appropriate combination of three depths. For this, we calculate the fitness value over the whole map. As an exemplar we calculate the best combination of depths for January 2014—the first month of data that is used in the simulations reported here. We also calculated the best combination for February



Figure 3.11: Plots of ROMS fitness maps. Left: A map created from the full 12 layer set. Center: A map created from the best 3 layer set. Right: A map created from the 6 layer set used in our simulations.

and March of 2014. They yielded the same combination of depths. Based on this, for purposes of the present study, all further investigation was conducted for the January 2014 dataset. The best depths are $\{0, 75, 400\}$ m. To evaluate the results we build the fitness map for these three depths and compare it to the best map H_{12} . We build H_{12} as follows: for a given hour k, we take the time slice of the full (12 depth) ROMS flow field (i.e., a multidimensional array representing this flow field at the given hour) and for each surface position \mathbf{p}_l of this slice, we compare all combinations of depth triples by their point-wise fitness value. For each position, we select the best depth triple and save its scalar fitness value in the matrix H_{12}^k . The target best map is obtained as the matrix

$$H_{12} = \sum_{k=1}^{T} H_{12}^{k}, \qquad (3.16)$$

where T depends on the time span used for calculations (as mentioned, we used one month of data). The results are shown in Fig. 3.11 on the left and in the center.

As one can see, three depths are not enough to approximate H_{12} . We add additional depths greedily, by adding the depth that gives the highest gain of fitness at each step, and terminate the

procedure once the map is represented "reasonably well". The resulting map built on a six depth set $\{0, 10, 30, 75, 200, 400\}$ m is shown in Fig. 3.11 on the right. The overall fitness of this map is $\approx 85\%$ of the fitness of the 12 depth map. We treat this as a reasonable approximation; these are the six discrete depths used in the simulations reported here ** .

Upper Bound of Performance

We find the upper bound by evaluating the performance of the high-level algorithm for an *ideal drifter*. An ideal drifter is capable of

- instantaneously changing the drogue depth;
- instantaneously estimating ocean currents, i.e., the drifter spends no time with the drogue at a particular depth to estimate currents at that depth;
- · making ideal measurements of ocean currents with no estimation noise; and
- perfect communication.

Thus, at the moment of taking a control decision, the ideal drifter knows all currents at its present location instantaneously. This scenario gives us an upper bound of what can be achieved by our algorithm in the best case. All further evaluations are performed with the simulation setup parameters given in Tab. 3.1, for 100 simulation runs.

As mentioned earlier, the aggregation area in the mission space Ω is defined by the central attraction point \mathbf{p}_c and the radius r_c of the aggregation high-level controller introduced in (3.5). Each of the 100 simulation runs performs a full aggregation scenario with random initial positions of drifters inside the preference area A. The results of the simulation for 90 and 180 days

^{**}Note that adjusting the set of depths online during a mission presents an interesting extension of the controller and is subject to future work.

| Table 3.1: | Simulation | parameters |
|------------|------------|------------|
|------------|------------|------------|

| ROMS dates | starting with January, 01, 2014 | |
|----------------------|---------------------------------------|--|
| Simulation time span | 90 / 180 days | |
| Depths set | { 0, 10, 30, 75, 200, 400 } m | |
| Number of drifters | 30 | |
| Number of runs | 100 | |
| r_d | 30 000 m | |
| d_{\min} | 1 000 m | |
| p_c | longitude = -124.75, latitude = 35.75 | |
| r _c | 250 000 m | |
| w _c | 10 | |
| Simulation time step | 300 s | |

are presented in Tab. 3.2 for the ideal drifters. Fig. 3.12 presents the average evolution of the aggregation metric (3.9) over time. As one can see, the ideal drifter scenario is almost saturated after 90 days of simulation and it has $\approx 90\%$ of its final value at this point. Thus, we decide to limit simulation time to 90 days for all further evaluations. This time cutoff has some advantages. First, late aggregation is penalized, which is important since in practice drifters may have limited life time. Second, we partially account for cases where saturation may happen earlier and the system starts deteriorating by losing drifters (i.e., they leave Ω), which happens inevitably when the time span is increased.

From the simulations and applying heuristic (3.12), one can see that for the ideal drifter case we obtain ≈ 2 clusters after the process is finished, which is a very good result, considering the
| | Averag | e metric | Average number of drifters lost | | | |
|-------------------|---------|----------|------------------------------------|----------|--|--|
| | 90 days | 180 days | 90 days | 180 days | | |
| Ideal drifter | 0.49 | 0.55 | 7 | 8.5 | | |
| Realistic drifter | 0.21 | 0.29 | 9 | 12.8 | | |

Table 3.2: System performance over 100 simulations

nature of external forcing and the fact that the algorithm utilizes no predictions but only in situ measurements. The choice of parameters for the *realistic drifter* scenario is explained next.

Performance under Estimation Limitations

In the previous section, we presented simulations that were performed under the assumption that the drifter can estimate currents instantaneously. This is never the case in reality. A drifter must sample currents at various depths before it makes an appropriate decision according to the low-level controller (see 3.3.2).

For our system, we assume that, in order to acquire a current estimate, the drogue has to submerge to a corresponding depth and drift for a certain time T_{est} with the ocean current. As an initial guess, we assume the estimation time to be $T_{est} = 10$ minutes per depth.

The vertical velocity of a realistic drogue is bounded. Our model uses the constant speed v for climbing and diving. The choice is based on parameters of existing buoyancy-driven vehicles [Reed et al., 2011, Han et al., 2010, Schwithal and Roman, 2009, Eriksen et al., 2001] that report achievable maximum vertical velocities in the range of 0.1 to 0.5 m/s. For our system, we chose a reasonable value in the middle of this range, and set it to 0.3 m/s.



Figure 3.12: Evolution of the average aggregation metric over time based on 100 simulation runs. The green curve represents the performance of ideal drifters. The black curve represents the performance of realistic drifters with a low-level controller.

As can be inferred from [Serrano et al., 2004], given a slowly moving object, such as the drifter, with estimation time on the order of minutes and almost constant speed, modern estimation approaches can reduce the velocity measurement error to the order of millimeters per second. Hence, given that the average velocity of the drifter is around 0.1 to 0.2 m/s, such estimation errors due to noise can be neglected even for realistic drifters at first.

Another parameter of the basic low-level control algorithm is the tolerance angle α_{tol} . Fig. 3.13 presents the graph of performance of the system under different tolerance angles. From this graph we can infer that our system has the best average performance at $\alpha_{tol} = 20^{\circ} \pm 5^{\circ}$, thus, for the selection of the remaining parameters and further simulations we fix the value to 20° .

The next set of parameters belongs to the extended version of the low-level control algorithm (see 3.3.2). The first is the trust time T_{trust} . Fig. 3.14 depicts the graph of the performance of the



Figure 3.13: Performance of the realistic drifter system with the basic low-level controller under different tolerance angles. Every point is based on 100 simulation runs. The black curve represents the mean value and the green curves represent the standard deviation of the aggregation metric at every point.

system for varying T_{trust} . As one can see, the performance deterioration starts approximately after the 3.5 hour mark. Thus, for further experiments we pick $T_{\text{trust}} = 3.5$ hours.

The final addition to the low-level control algorithm is the ability to share information about ocean current estimates. We implemented it as a sample pool—a set of current estimates that every drifter collects from drifters that are within a distance d_{\min} when the control decision is taken. Such collaboration helps drifters to save time by splitting the task of current estimation between nearby drifters stochastically, and, hence, taking advantage of clustering (see 3.3.2). The performance of the system in this setting is shown in Tab. 3.2 for the realistic drifter scenario. The average performance of 0.21 and 0.29 (for 90 and 180 days of simulation) corresponds to $\approx 3-5$ clusters (with around 6–10 drifters per cluster). Although it is lower than the average score for the ideal drifter scenario, we are encouraged by the performance and the implication for the practical use.



Figure 3.14: Performance of the realistic drifter system with the extended low-level controller under different T_{trust} . Every point is based on 100 simulation runs. The black curve represents the mean value and the green curves represent the standard deviation of the aggregation metric at every point.

Considering high ship operation costs, such clustering can significantly facilitate the process of collection of drifters.

Finally, we demonstrate an example of a complete deployment scenario, where drifters are dropped at one position (marked with a green cross), then spread all over the area during 90 days (Fig. 3.16 left) and finally aggregate over the last 90 days (Fig. 3.16 right). Although we have quite a good aggregation (more than a half of drifters assembled), the seemingly low score in this scenario (≈ 0.31) reflects the fact that the rest of the drifters are either lost or completely spread through the area due to external forcing (in the particular case of Fig. 3.16: 10 lost, 3 outliers).

Performance under Noisy Estimation

In this section, we finally evaluate the system's performance under estimation noise. We assume estimations of velocities are acquired from GPS data solely. As mentioned above in 3.3.3, modern



Figure 3.15: Performance of the realistic drifter system with the extended controller under different estimation noise levels. The level of the noise is defined through the upper limit of a triangular distribution. The system does not exhibit significant drop in performance up to the noise level of 5 cm/s. After this mark the performance degrades gracefully. Every point is based on 100 simulation runs. The black curve represents the mean value and the green curves represent the standard deviation of the aggregation metric at every point.

estimation approaches allow to obtain estimation errors on the order of millimeters per second. We now calculate the errors for a worst case scenario. For that, let us consider the simple velocity estimator

$$\hat{\mathbf{f}} = \frac{\mathbf{p}_k - \mathbf{p}_{k-1}}{T_{est}} \,. \tag{3.17}$$

For the GPS error, we assume simple additive noise with a uniform distribution symmetric around 0 with a lower limit $-e_{\text{max}}$ and upper limit e_{max} , and independence of noise components. For (3.17), we get a triangular distribution of the noise (as the distribution of the sum of uniformly distributed random variables) with half the distribution width $e_{\text{tr}} = \frac{2e_{\text{max}}}{T_{\text{est}}}$. For the value of e_{max} , we take the



Figure 3.16: Left: An example of positions of 30 drifters after 90 days of spreading. Right: An example of positions of 30 drifters after 90 days of aggregation following the spreading phase (180 days of simulation total). Blue triangles denote drifters and drifter clusters. Numbers near triangles denote cluster sizes. Triangles without numbers denote single drifters. Crossed out triangles on the borders (i.e., borders of Ω) denote lost drifters. The green cross marks the initial position of drifters and the green circle denotes the preference area A. The black thin curve denotes the coastal line.

error of a popular inexpensive GPS satellite messenger Spot Tracker^{††} $e_{\text{max}} = 6.4 \text{ m}$. Thus, for the triangular distribution we have $e_{\text{tr}} \approx 0.021 \text{ m/s}$.

Fig. 3.15 depicts the curve of the average aggregation metric for different e_{tr} . As one can see, the system is robust to the estimation noise. In the worst case scenario with e_{tr} and even beyond (with noise up to 5 cm/s) there is no significant drop in performance. More than that, the system exhibits critical problems (i.e., performance lower than 0.1) only at the upper limit of the noise distribution of 20 cm/s, which corresponds to the average velocity of the ocean currents.

^{††}For more details, see https://www.findmespot.com/downloads/SPOT2-SellSheet.pdf

3.4 Conclusion

We investigate a very challenging scenario of perceptual scarcity resulted from expensive measurements in a problem of aggregation and coverage of a multi-drifter system. First, we investigate prediction accuracy of the ROMS model and show that it can not be used reliably for long horizon predictions. Based on these results we propose using a simple controller utilizing in situ measurements of ocean currents. Second, we address the target problem of aggregation and coverage for a multi-drifter system. Since this problem is inherently scarce in terms of ocean current measurements, we propose a solution utilizing a combination of a hierarchical control, delayed decisions, and indirect information about ocean currents in the form of spatially and temporally separated current measurements. We show that our approach can significantly reduce the number of clusters of drifters during the aggregation stage. The reduction of clusters drastically reduces logistics complexity after completion of a mission.

Part II

Task Data Scarcity

Chapter 4

Sim-to-(Multi)-Real: Transfer of Low-Level Robust Control Policies to Multiple Quadrotors

Reinforcement learning (RL) has demonstrated capability of finding impressively complex policies in simulated environments, significantly outperforming more traditional control approaches. Unfortunately, when it comes to applying RL to real robotic systems the results have remained modest. The primary reason is that RL requires tremendous amounts of data, which is exceptionally hard to acquire on the real robots due to complexity of operating these systems and high chance of damaging hardware when executing an untrained policy. In other words, the key challenge in applying RL to robotic systems is data scarcity. One really challenging example is quadrotor control where the failure often results in a high speed impact with the ground, hence, at best, only a few trajectories can be acquired. Nonetheless, RL can provide significant reduction in

Parts of this chapter appeared in:

Artem Molchanov, Tao Chen, Wolfgang Hönig, James A. Preiss, Nora Ayanian, and Gaurav S. Sukhatme. Sim-to-(multi)-real: Transfer of low-level robust control policies to multi-ple quadrotors. In 2019 IEEE/RSJ International Conference on Intelligent Robots and Sys-tems, IROS 2019, Macau, China, 2019, pages 59–66. URL: https://doi.org/10.1109/IROS40897.2019.8967695.

[•] Artem Molchanov, Tao Chen, Wolfgang Hönig, James A. Preiss, Nora Ayanian, and Gaurav S. Sukhatme. Sim-to-(multi)-real: Transfer of low-level robust control policies to multiple quadrotors. In Southern California Robotics Symposium, SCR 2019, Pasadena, CA, US, 2019.



Figure 4.1: Three quadrotors of different sizes controlled by the same policy trained entirely in simulation.

engineering effort, since traditional control-theoretic approaches to stabilizing a quadrotor often require careful, model-specific system identification and parameter tuning to succeed. Particularly, we are interested in finding a single control policy that stabilizes any quadrotor and moves it to a goal position safely, without manual parameter tuning. Such a control policy can be very useful for testing of new custom-built quadrotors, and as a backup safety controller. Our primary objective for the controller is robustness to external disturbances and recovery from collisions. Our secondary objective is position control and trajectory tracking. To achieve our goal and overcome data scarcity in this problem, we combine RL with simulated quadrotor models as a task structural prior for learning a control policy. To close the gap between simulation and reality, we analyze the impact of various key parameters, such as modeling of sensor noise and using different cost functions. We also investigate how we can improve sim-to-real transfer (S2R) by applying domain randomization, a technique that trains over a distribution of the system dynamics to help trained policies be more resistant to a simulator's discrepancies from reality [Tobin et al., 2017, Sadeghi

and Levine, 2017, James et al., 2017, Tan et al., 2018]. We transfer our policy to three different quadrotor platforms and provide data on hover quality, trajectory tracking, and robustness to strong disturbances and harsh initial conditions.

To the best of our knowledge, ours is the first neural network (NN) based low-level quadrotor attitude-stabilizing (and trajectory tracking) policy trained completely in simulation that is shown to generalize to multiple quadrotors. Our contributions can be summarized as follows:

- A system for training model-free low-level quadrotor stabilization policies without any auxiliary pre-tuned controllers.
- Successful training and transfer of a single control policy from simulation to multiple real quadrotors.
- An investigation of important model parameters and the role of domain randomization for transferability.
- A software framework for flying Crazyflie 2.0 (CF) based platforms using NN controllers and a Python-based simulated environment compatible with OpenAI Gym [Brockman et al., 2016] for training transferable simulated policies*.

4.1 Related Work

Transferring from simulation to reality (S2R) is a very attractive approach to overcome the issues of safety and complexity of data collection for reinforcement learning on robotic systems. In the following we group the related work into different categories.

^{*}Available at https://sites.google.com/view/sim-to-multi-quad.

S2R with model parameter estimation. A substantial body of work considers closing the S2R gap by carefully estimating parameters of the real system to achieve a more realistic simulation. For example, Lowrey et al. [2018] transfer a non-prehensile manipulation policy for a system of three one-finger robot arms pushing a single cylinder. Tan et al. [2018] show transferability of agile locomotion gaits for quadruped robots. Antonova et al. [2017] learn a robust policy for rotating an object to a desired angle. While careful parameter estimation can provide a good estimate of the model, it often requires sophisticated setups [Förster, 2015]. In contrast, we transfer a learned policy to novel quadrotors for which we do not perform accurate model parameter estimation.

S2R with iterative data collection. An alternative way to overcome the problem of S2R gap is learning distributions of dynamics parameters in an iterative manner. For example, Christiano et al. [2016] learn inverse dynamics models from data gradually collected from a real robotic system, while transferring trajectory planning policy from a simulator. Chebotar et al. [2019], Zhu et al. [2018] transfer manipulation policies by iteratively collecting data on the real system and updating a distribution of dynamics parameters for the simulator physics engine. Similar principles work for the problem of humanoid balancing [Tan et al., 2016].

The common problem of these approaches is the necessity to execute untrained policies directly on the robot, which may raise safety concerns. In contrast to all of the works presented above, we are interested in a method that can i) transfer very low-level policies that directly control actuator forces, ii) transfer to multiple real robots with different dynamics, iii) control inherently unstable systems with dangerous consequences of failures, iv) avoid the need for data collection on the real system with a possibly unstable policy, and v) avoid complex setups for system parameter estimation.

Domain randomization. Domain randomization (DR) [Tobin et al., 2017] is a simple albeit promising domain adaptation technique that is well suited for S2R. It compensates for the

discrepancy between different domains by extensively randomizing parameters of the training (source) domain in simulation. In some cases, it can eliminate the need for data collection on the real robot completely. DR has been successfully applied for transferring visual features and high-level policies. Tobin et al. [2017] employ DR for training a visual object position predictor for the task of object grasping. The policy is trained in a simulation with random textures and lightning. A similar direction explores intermediate-level representations, such as object corners identification, and trained an interpretable high-level planning policy to stack objects using the Baxter robot [Tremblay et al., 2018].

S2R transfer of dynamics and low-level control policies is considered a more challenging task due to the complexity of realistic physics modeling. Nonetheless, there have been some promising works. Peng et al. [2018] apply DR to the task of pushing an object to a target location using a Fetch arm. The policy operates on joint angle positions instead of directly on torques. Mordatch et al. [2015] transfer a walking policy by optimizing a trajectory on a small ensemble of dynamics models. The trajectory optimization is done offline for a single Darwin robot. In contrast, our work does not require careful selection of model perturbations and is evaluated on multiple robots.

S2R for quadrotor control. S2R has also been applied to quadrotor control for transferring high-level visual navigation policies. Most of these works assume the presence of a low-level controller capable of executing high-level commands. Sadeghi and Levine [2017] apply CNN trained in simulation to generate a high-level controller selecting a direction in the image space that is later executed by a hand-tuned controller. Kang et al. [2019] look at the problem of visual navigation using a Crazyflie quadrotor and learn a high-level yaw control policy by combining simulated and real data in various ways.

The approaches most related to ours are the works of Koch et al. [2019] and Hwangbo et al. [2017]. In the former work, a low-level attitude controller is replaced by a neural network and

transferred to a real quadrotor [Koch et al., 2019]. In the latter work, a low-level stabilizing policy for the Hummingbird quadrotor that is trained in simulation is transferred to the Hummingbird. In contrast to their work i) we assume minimal prior knowledge about quadrotor's dynamics parameters; ii) we transfer a single policy to multiple quadrotor platforms; iii) we simplify the cost function used for training the policy; iv) we investigate the importance of different model parameters and the role of domain randomization for S2R transfer of quadrotor's low-level policies; and v) unlike Hwangbo et al. [2017] we do not use auxiliary pre-tuned PD controller in the learned policy.

4.2 **Problem Statement**

We aim to find a policy that directly maps the current quadrotor state to rotor thrusts. The quadrotor state is described by the tuple (\mathbf{e}_p , \mathbf{e}_v , R, \mathbf{e}_ω), where $\mathbf{e}_p \in \mathbb{R}^3$ is the position error, $\mathbf{e}_v \in \mathbb{R}^3$ is the linear velocity error of the quadrotor in the world frame, \mathbf{e}_ω is the angular velocity error in the body frame, and $R \in SO(3)$ is the rotation matrix from the quadrotor's body coordinate frame to the world frame. The objective is to minimize the norms of \mathbf{e}_p , \mathbf{e}_v , \mathbf{e}_ω and drive the last column of R to $[0,0,1]^T$ in the shortest time. The policy should be robust, i.e., it should be capable of recovering from different initial conditions, as well as transferable to other quadrotor platforms while retaining high performance. We make the following assumptions:

• We only consider quadrotors in the × configuration (Fig. 4.2). This configuration is the most widely used, since it allows convenient camera placement. It would not be possible for a single policy to control both + and × configurations without some input specifying the configuration of the current system, since the fundamental geometric relationship between the motors and the quadrotor's axes of rotation is altered.

- We assume access to reasonably accurate estimates of the quadrotor's position, velocity, orientation, and angular velocity. Similar assumptions are typical in the robotics literature. They are commonly satisfied in practice by fusion of inertial measurements with localization from one or more of the following: vision, LIDAR, GPS, or an external motion capture system.
- We consider quadrotors within a wide but bounded range of physical parameters [†]. The ranges we experiment with are presented in Table 4.1.

4.3 Dynamics Simulation

In this section we describe our dynamics model for simulation in detail.

4.3.1 Rigid Body Dynamics for Quadrotors

We treat the quadrotor as a rigid body with four rotors mounted at the corners of a rectangle. This rectangle lies parallel to the x - y plane of the body coordinate frame, as do the rotors. Each rotor is powered by a motor that only spins in one direction, hence the rotors can only produce positive thrusts in the z-direction of the body frame. The origin of the body frame is placed in the center of mass of the quadrotor. The dynamics are modeled using Newton-Euler equations:

$$m \cdot \ddot{\mathbf{x}} = m \cdot \mathbf{g} + R \cdot \mathbf{F} \tag{4.1}$$

$$\dot{\boldsymbol{\omega}} = I^{-1}(\boldsymbol{\tau} - \boldsymbol{\omega} \times (I \cdot \boldsymbol{\omega})) \tag{4.2}$$

$$\dot{R} = \boldsymbol{\omega}_{\times} R, \tag{4.3}$$

[†]This assumption is introduced to restrict our system to typical quadrotor shapes and properties. We are not considering edge cases, such as very high thrust-to-weight ratios, or unusual shapes with a large offset of the center of mass with respect to quadrotor's geometric center.



Figure 4.2: Top-down view of our generalized quadrotor model. There are 5 components: baselink, payload, 4 arms, 4 motors, and 4 rotors. The model always assumes the \times configuration, with the front and left pointing at the positive X and Y directions respectively. Motors are indexed in the clockwise direction starting from the front-right motor. The front-right motor rotates counterclockwise and generates a thrust force f_1 .

where $m \in \mathbb{R}_{>0}$ is the mass, $\ddot{\mathbf{x}} \in \mathbb{R}^3$ is the acceleration in the world frame, $\mathbf{g} = [0, 0, -9.81]^T$ is the gravity vector, $R \in SO(3)$ is the rotation matrix from the body frame to the world frame, $\mathbf{F} \in \mathbb{R}^3$ is the total thrust force in the body frame, $\boldsymbol{\omega} \in \mathbb{R}^3$ is the angular velocity in the body frame, $I \in \mathbb{R}^{3\times3}$ is the inertia tensor, $\boldsymbol{\tau} \in \mathbb{R}^3$ is the total torque in the body frame, and $\boldsymbol{\omega}_{\times} \in \mathbb{R}^{3\times3}$ is a skew-symmetric matrix associated with $\boldsymbol{\omega}$ rotated to the world frame.

The total torque au is calculated as:

$$\boldsymbol{\tau} = \boldsymbol{\tau}_p + \boldsymbol{\tau}_{th},\tag{4.4}$$

where τ_{th} is a thruster torque, produced by thrust forces [Martin and Salaün, 2010], τ_p is a torque along the quadrotor's z-axis, produced by difference in speed rotation of the propellers:

$$\boldsymbol{\tau}_{p} = r_{t2t} \cdot [+1, -1, +1, -1]^{T} \odot \mathbf{f}, \tag{4.5}$$

where r_{t2t} is a torque to thrust coefficient, -1 indicates the rotor turns clockwise, +1 indicates the rotor turns counterclockwise, and $\mathbf{f} = [f_1, f_2, f_3, f_4]^T$ is a vector representing the force generated by each rotor.

4.3.2 Normalized Motor Thrust Input

It is common in the literature of quadrotor control to assume that the motor speed dynamics are nearly instantaneous. This assumption allows us to treat each motor's thrust as a directly controlled quantity. In this paper, to facilitate transfer of the same policy to multiple quadrotors, we instead define a normalized control input $\hat{\mathbf{f}} \in [0,1]^4$ such that $\hat{\mathbf{f}} = \mathbf{0}$ corresponds to no motor power and $\hat{\mathbf{f}} = \mathbf{1}$ corresponds to full power. Note that the nominal value of $\hat{\mathbf{f}}$ for a hover state depends on the thrust-to-weight ratio of the quadrotor. By choosing this input, we expect the policy to learn a behavior that is valid on quadrotors with different thrust-to-weight ratios without any system identification needed. The input $\hat{\mathbf{f}}$ is derived from the policy action $\mathbf{a} \in \mathbb{R}^4$ by the affine transformation

$$\hat{\mathbf{f}} = \frac{1}{2}(\mathbf{a}+\mathbf{1}) \tag{4.6}$$

to keep the policy's action distribution roughly zero-mean and unit-variance. Since thrust is proportional to the square of rotor angular velocity, we also define $\hat{\mathbf{u}} = \sqrt{\hat{\mathbf{f}}}$ as the normalized angular velocity command associated with a given normalized thrust command.

4.3.3 Simulation of Non-Ideal Motors

The assumption of instantaneous motor dynamics is reasonable for the slowly-varying control inputs of a human pilot, but it is unrealistic for the noisy and incoherent control inputs from an untrained stochastic neural network policy. To avoid training a policy that exploits a physically implausible phenomenon of the simulator, we introduce two elements to increase realism: motor lag simulation and a noise process.

Motor lag We simulate motor lag with a discrete-time first-order low-pass filter:

$$\hat{\mathbf{u}}_{t}' = \frac{4dt}{T} (\hat{\mathbf{u}}_{t} - \hat{\mathbf{u}}_{t-1}') + \hat{\mathbf{u}}_{t-1}', \qquad (4.7)$$

where $\hat{\mathbf{u}}'_t \in \mathbb{R}^4$ is the vector of the filtered normalized rotor angular velocities, dt is the time between the inputs and $T \ge 4dt$ is the 2% *settling time*, defined for a step input as

$$T = dt \cdot \min\{t \in \mathcal{N} : \|\hat{\mathbf{u}}_{t'}' - \hat{\mathbf{u}}_t'\|_{\infty} < 0.02 \text{ for all } t' \ge t\}.$$

$$(4.8)$$

Motor noise We add motor noise $\epsilon_t^{\mathbf{u}}$ following a discretized Ornstein-Uhlenbeck process:

$$\epsilon_t^{\mathbf{u}} = \epsilon_{t-1}^{\mathbf{u}} + \theta(\mu - \epsilon_{t-1}^{\mathbf{u}}) + \sigma \mathcal{N}(0, 1), \tag{4.9}$$

where $\mu = 0$ is the process mean, θ is the decay rate, σ is the scale factor, and $\mathcal{N}(0, 1)$ is a random variable with four-dimensional spherical Gaussian distribution.

The final motor forces are computed as follows:

$$\mathbf{f} = f_{\max} \cdot \left(\hat{\mathbf{u}}_t' + \epsilon_t^{\mathbf{u}}\right)^2. \tag{4.10}$$

Here, f_{max} is found using the *thrust-to-weight ratio* r_{t2w} :

$$f_{max} = 0.25 \cdot g \cdot m \cdot r_{t2w},\tag{4.11}$$

where g is the gravity constant.

4.3.4 Observation Model

In addition to the motor noise, we also model sensor and state estimation noise. Noise in estimation of position and orientation as well as in linear velocity is modeled as zero-mean Gaussian. Noise in angular velocity measured by the gyroscope follows the methods presented by Furrer et al. [2016]. Sensor noise parameters were roughly estimated from data recorded while quadrotors were resting on the ground.

4.3.5 Numerical Methods

We use first-order Euler's method to integrate the differential equations of the quadrotor's dynamics. Due to numerical errors in integration, the rotational matrix R loses orthogonality over time. We re-orthogonolize R using the Singular Value Decomposition (SVD) method. We perform reorthogonolization every 0.5s of simulation time or when the orthogonality criteria fails:

$$\|RR^T - I_{3\times3}\|_1 \ge 0.01,\tag{4.12}$$

where $\|\cdot\|_1$ denotes the elementwise L_1 norm. If $U\Sigma V^T = R$ is the singular value decomposition of R, then $R_{\perp} = UV^T$ is the solution of the optimization problem

$$\underset{A \in \mathbb{R}^{3 \times 3}}{\text{minimize}} \|A - R\|_F \text{ subject to } A^T A = I,$$
(4.13)

where $\|\cdot\|_F$ denotes the Frobenius norm [Higham, 1989], making R_{\perp} a superior orthogonalization of R than, e.g., that produced by the Gram-Schmidt process.

4.4 Learning & Verification

In this section, we discuss the methodology we use to train a policy, including the domain randomization of dynamics parameters, the reinforcement learning algorithm, the policy class, and the basic structure of our experimental validations. More details of the experiments are given in Section 4.5.

| Variable | Unit | Nominal Randomization | Total Randomization |
|--------------|-----------------|-----------------------|---------------------------------|
| m | kg | 0.028 | ≤ 5 |
| $l_{body,w}$ | m | 0.065 | $\sim \mathscr{U}(0.05, 0.2)$ |
| T | S | 0.15 | $\sim \mathscr{U}(0.1, 0.2)$ |
| r_{t2w} | kg/N | 1.9 | $\sim \mathscr{U}(1.8, 2.5)$ |
| r_{t2t} | s ⁻² | 0.006 | $\sim \mathscr{U}(0.005, 0.02)$ |

Table 4.1: Randomization variables and their distributions.

4.4.1 Randomization

We investigate the role of domain randomization for generalization toward models with unknown dynamics parameters. The geometry of the generalized quadrotor is defined by variables shown in Fig. 4.2. For brevity, we omit the height variables $(l_{*,h})$, where * is the component name) in the figure. Table 4.1 shows the list of variables we randomize.

During training, we sample dynamics parameters for each individual trajectory. We experiment with two approaches for dynamics sampling:

- 1. Randomization of parameters around a set of nominal values assuming that approximate estimates of the parameters are available. We use existing Crazyflie 2.0 parameter estimates [Förster, 2015].
- 2. Randomization of parameters within a set of limits. The method assumes that the values of parameters are unknown but bound by the limits.

In the first scenario, we randomize all parameters describing our quadrotor around a set of nominal values, and in case a Gaussian distribution is used, we check the validity of the randomized values (mostly to prevent negative values of inherently positive parameters). In the second scenario, we start by sampling the overall width of the quadrotor $(l_{body,w})$ and the rest of the geometric parameters are sampled with respect to it. The total mass m of the quadrotor is computed by sampling densities of individual components. The inertia tensors of individual components with respect to the body frame are found using the parallel axis theorem.

4.4.2 Policy Representation

We use a fully-connected neural network to represent a policy. The neural network has two hidden layers with 64 neurons each and the tanh activation function, except for the output layer that has the linear activation. The network input is an 18-dimensional vector representing the quadrotor state presented in Section 4.2. Rather than inputting the current state and a goal state, we input only the error between the current and goal state, except for the rotation matrix which represents the current orientation. This reduces the input dimensionality, and trajectory tracking is still possible by shifting the goal state. In our policy we do not use any additional stabilizing PID controllers and directly control the motor thrusts, in contrast to existing approaches [Hwangbo et al., 2017]. Hence, our neural network policy directly outputs the normalized thrust commands a that are later converted to the normalized force commands $\hat{\mathbf{f}}$ (see Eq. (4.6)).

4.4.3 Policy Learning

The policy is trained using the Proximal Policy Optimization (PPO) algorithm [Schulman et al., 2017]. PPO has recently gained popularity for its robustness and simplicity. PPO is well-suited for RL problems with continuous state and action spaces where interacting with the environment is not considered expensive. We use the implementation available in $Garage^{\ddagger}$, a TensorFlow-based

[‡]https://github.com/rlworkgroup/garage

open-source reinforcement learning framework. This framework is an actively supported and growing reincarnation of the currently unsupported rllab framework [Duan et al., 2016a].

During training, we sample initial states uniformly from the following sets: orientation is sampled from the full SO(3) group, position within a 2 m box around the goal location, velocity with a maximum magnitude of 1 m/s, and angular velocity with a maximum magnitude of 2π rad/s. The goal state is always selected to hover at $[0, 0, 2]^T$ in the world coordinates. At execution time, we can translate the coordinate system to use the policy as a trajectory tracking controller. We parameterize the quadrotor's attitude with a rotation matrix instead of a quaternion because the unit quaternions double-cover the rotation group SO(3), meaning that a policy with a quaternion input must learn that the quaternions q and -q represent the same rotation.

The reinforcement learning cost function is defined as

$$c_{t} = \left(\|\mathbf{e}_{p}\|_{2} + \alpha_{v} \|\mathbf{e}_{v}\|_{2} + \alpha_{\omega} \|\mathbf{e}_{\omega}\|_{2} + \alpha_{a} \|\mathbf{a}\|_{2} + \alpha_{R} \cos^{-1}\left((\operatorname{tr}(R) - 1)/2\right)\right) dt,$$
(4.14)

where R is the rotation matrix and $\alpha_{\omega}, \alpha_{a}, \alpha_{R}, \alpha_{v}$ are non-negative scalar weights. The term $\cos^{-1}((\operatorname{tr}(R) - 1)/2)$ represents the angle of rotation between the current orientation and the identity rotation matrix. We investigate influence of different components in the experimental section.

4.4.4 Sim-to-Sim Verification

Before running the policy on real quadrotors, the policy is tested in a different simulator. This sim-to-sim transfer helps us verify the physics of our own simulator and the performance of policies in a more realistic environment. In particular, we transfer to the Gazebo simulator with the RotorS

package [Furrer et al., 2016] that has a higher-fidelity simulation compared to the one we use for training.

Gazebo by default uses the ODE physics engine, rather than our implementation of the Newton-Euler equations. RotorS models rotor dynamics with more details, e.g. it models drag forces which we neglect during learning. It also comes with various pre-defined quadrotor models, which we can use to test the performance of trained policies for quadrotors where no physical counterpart is available.

We found that using our own dynamics simulation for learning is faster and more flexible compared to using Gazebo with RotorS directly.

4.4.5 Sim-to-Real Verification

We verify our approach on various physical quadrotors that are based on the Crazyflie 2.0 platform. The Crazyflie 2.0 is a small quadrotor that can be safely operated near humans. Its light weight (27 g) makes it relatively crash-tolerant. The platform is available commercially off-the-shelf with an open-source firmware. We build heavier quadrotors by buying standard parts (e.g., frames, motors) and using the Crazyflie's main board as a flight computer.

We test policies by sequentially increasing quadrotor size (starting with the Crazyflie 2.0) for safety reasons. We quantify the performance of our policies using three different experiments. First, we evaluate the hover quality by tasking the quadrotor to hover at a fixed position and record its pose at a fixed frequency. For each sample, we compute the Euclidean position error $\|\mathbf{e}_p\|$ and the angular error ignoring yaw:

$$e_{\theta} = \arccos(R(:,3) \cdot [0,0,1]^T) = \arccos(R(3,3)),$$
 (4.15)

81

| Robot | CF | Small | Medium |
|-------------------------|-----|-------|--------|
| Weight [g] | 33 | 73 | 124 |
| $l_{body,w}$ [mm] | 65 | 85 | 90 |
| $l_{rotor,r}$ [mm] | 22 | 33 | 35 |
| r_{t2w} (approximate) | 1.9 | 2.0 | 2.7 |

Table 4.2: Robot Properties.

where R(:,3) is the last column of the rotation matrix R, and R(3,3) is its bottom-right element. We denote the mean of the position and angular errors over all collected hover samples as e_h (in m) and \bar{e}_{θ} (in deg), respectively. We characterize oscillations by executing a fast Fourier transform (FFT) on the roll and pitch angles, and report f_o (in Hz) – the highest frequency with a significant spike. Second, we evaluate the trajectory tracking capabilities by tasking the quadrotor to track a pre-defined figure-eight trajectory and record the position errors $\|\mathbf{e}_p\|$ (in m). We denote the mean of the errors during the flight as e_t . Finally, we disturb the quadrotors and check if they recover using our policies (an experiment that is difficult to quantify on a physical platform).

4.5 Experiments

We validate our control policies on three different quadrotors with varying physical properties: Crazyflie 2.0, small, and medium size as described in Table 4.2. All quadrotors use a similar control board with a STM32F405 microcontroller clocked at 168 MHz, executing the same firmware. We use the Crazyswarm testbed [Preiss et al., 2017] for our experiments. In particular, the state estimate is computed by an extended Kalman filter (EKF) that fuses on-board IMU data and motion capture information. For the experiments with trajectories, we upload them at the beginning of the flight and compute the moving goal states on-board. We make three major changes to the firmware: First, we add a control policy, which is an auto-generated C-function from the trained NN model in TensorFlow. Second, we remove the software low-pass filter of the gyroscope, and increase the bandwidth of its hardware low-pass filter. We found that the reduction in the gyroscope delay significantly reduces the quadrotor's physical oscillations when controlled by our policy. Third, we only use the motion capture system to estimate linear velocities using finite differences and ignore accelerometer readings. We found that the velocity estimates were noisier otherwise, which caused a large position offset when using our policy. Whenever we compare to a non-learned controller, we use the default Crazyswarm firmware without our modifications. Our motion capture system captures pose information at 100 Hz; all on-board computation (trajectory evaluation, EKF, control) is done at 500 Hz. Evaluating the neural network takes about 0.8 ms.

To train the policy, we collect 40 simulated trajectories with a duration of 7 s (i.e. 4.7 min of simulated flight) at each iteration of PPO. In simulation the policy runs at 100 Hz and the dynamics integration is executed at 200 Hz. Samples for training are collected with the policy rate. We train the majority of our policies for 3000 iterations, which we found sufficient for convergence. The exception is the scenarios with randomization, for which we run 6000 iterations due to slower convergence. In each scenario, we train five policies by varying the seed of the pseudorandom number generator used to generate the policy's stochastic actions and inject noise into the environment. For the test on the real system, we select the two best seeds according to the average (among trajectories) sum-over-trajectory Euclidean distance cost (i.e. $\|\mathbf{e}_p\|$) computed during policy training. After that, we visually inspect the performance of the two seeds in simulation and select the one that generates smoother trajectories and exhibits smaller attitude oscillations (a subjective measure).

4.5.1 Ablation Analysis on Cost Components

We analyze the necessity of different terms in the RL training cost function (4.14) on the flight performance in simulation and on a real quadrotor, because we are interested in a simpler cost function with fewer hyper-parameters. During training, we use approximate parameters of the Crazyflie 2.0 quadrotor model Förster [2015]. Here, we do not apply parameter randomization, but we incorporate sensor and thrust noise. We let the quadrotor hover at a fixed point and record its pose at 100 Hz for 10 s and report the mean position error e_h , mean angular error \bar{e}_{θ} , and oscillation frequency f_o , as defined in Section 4.4.5. Our results are shown in Table 4.3.

We notice that we can train a successful policy with a cost that only penalizes position, angular velocity, and actions, as long as α_{ω} is larger than 0.05 but smaller than 1 (see rows 1 – 6 in Table 4.3). The optimal value of α_{ω} differs slightly: in simulation $\alpha_{\omega} = 0.25$ achieves the lowest position and angular errors. On the real quadrotor, we notice that higher α_{ω} can result in significantly higher errors in position. Thus, we chose $\alpha_{\omega} = 0.1$, $\alpha_a = 0.05$ (and $\alpha_R = \alpha_v = 0$) as a baseline for our further experiments.

We can add a cost for rotational errors by setting $\alpha_R > 0$, which improves position and angular errors in simulation, but results in slightly larger position errors on the physical quadrotor (see rows 7 and 8 in Table 4.3) compared to the baseline. The major advantage of this added cost is that it also stabilizes yaw, which might be desired for takeoff or if the quadrotor is carrying a camera.

Finally, we compared our cost function with the cost function that is similar to the one previously introduced by Hwangbo et al. [2017]. It additionally includes cost on linear velocity (i.e. $\alpha_v > 0$; see row 9 in Table 4.3). This cost function is harder to tune because of the larger number of hyper-parameters. The learned policy showed slightly worse position and angular errors

| # | Cost Parameters | RotorS | | | CF | | | | | |
|---|---|-----------------|--------------------|-------|-------|-------------------|-------|--|--|--|
| | | e_h | \bar{e}_{θ} | f_o | e_h | $\bar{e}_{	heta}$ | f_o | | | |
| 1 | $\alpha_{\omega}: 0.00, \alpha_a: 0.05$ | Training failed | | | | | | | | |
| 2 | $\alpha_{\omega}: 0.05, \alpha_a: 0.05$ | No Takeoff | | | 0.14 | 3.52 | 1.0 | | | |
| 3 | $\alpha_{\omega}: 0.10, \alpha_a: 0.05$ | 0.05 0.84 | | 1.0 | 0.09 | 2.07 | 0.9 | | | |
| 4 | $\alpha_{\omega}: 0.25, \alpha_a: 0.05$ | 0.05 | 0.02 | 0.7 | 0.21 | 2.59 | 0.6 | | | |
| 5 | $\alpha_{\omega}: 0.50, \alpha_a: 0.05$ | 0.08 | 0.07 | 0.5 | 0.30 | 2.34 | 0.5 | | | |
| 6 | $\alpha_{\omega}: 1.00, \alpha_a: 0.05$ | Training failed | | | d | | | | | |
| 7 | $\alpha_{\omega}: 0.10, \alpha_a: 0.05 \ \alpha_R: 0.25$ | 0.06 | 0.02 | 1.1 | 0.14 | 1.67 | 1.0 | | | |
| 8 | $\alpha_{\omega}: 0.10, \alpha_a: 0.05 \ \alpha_R: 0.50$ | 0.04 | 0.01 | 0.8 | 0.14 | 1.51 | 0.8 | | | |
| 9 | $\alpha_{\omega}: 0.075, \alpha_a: 0.050, \alpha_R: 0.000, \alpha_v: 0.125$ | 0.04 | 4 31 | 12 | 0.14 | 3 73 | 10 | | | |
| , | (cmp. Hwangbo et al. [2017]) | | 1.51 | 1.2 | 0.11 | 5.75 | 1.0 | | | |

Table 4.3: Ablation Analysis on Cost Components. α values not listed are 0.

in simulation and on the physical robot. All policies did not show any significant oscillations $(f_o \leq 1.2 \text{ Hz}).$

4.5.2 Sim-to-Real: Learning with Estimated Model

Based on our findings in Section 4.5.1, we use the cost function with parameters $\alpha_{\omega} = 0.1$, $\alpha_a = 0.05$ and test the influence of noise and motor delays (settling time) in a trajectory tracking task on the Crazyflie 2.0. The task includes taking off, flying a figure-eight at moderate speeds (up to 1.6 m/s, 5.4 m/s², 24 deg roll angle; 5.5 s long), and landing. In all cases, we do not perform model randomization.

| # | Policy CF | | | Small | | | Medium | | | |
|----|------------------------------------|--------------------|-------|-------|--------------------|-------|--------|--------------------|-------|-------|
| | | \bar{e}_{θ} | f_o | e_t | \bar{e}_{θ} | f_o | e_t | \bar{e}_{θ} | f_o | e_t |
| 1 | Mellinger | 1.93 | 1.4 | 0.11 | 1.11 | 0.5 | 0.14 | 0.78 | 1.1 | 0.04 |
| 2 | Mellinger (no memory) | 3.68 | 6.2 | 0.20 | 1.23 | 0.8 | 0.16 | 5.53 | 5.5 | 0.07 |
| 3 | Mellinger (uniform) | 3.17 | 6.3 | 0.19 | 1.12 | 0.7 | 0.19 | 5.29 | 5.4 | 0.32 |
| 4 | NN CF | 1.53 | 0.9 | 0.19 | 1.13 | 1.0 | 0.21 | 2.99 | 1.0 | 0.47 |
| 5 | NN CF (w/o delay) | 1.90 | 1.0 | 0.21 | 0.93 | 0.8 | 0.21 | 2.25 | 1.0 | 0.42 |
| 6 | NN CF 10 % random | 1.61 | 1.1 | 0.30 | 1.34 | 1.0 | 0.22 | 3.51 | 1.0 | 0.47 |
| 7 | NN CF 20 % random | 1.53 | 1.1 | 0.20 | 1.12 | 1.0 | 0.21 | 1.67 | 0.9 | 0.33 |
| 8 | NN CF 30 % random | 2.65 | 1.4 | 0.23 | 4.33 | 1.0 | 0.24 | 1.96 | 1.0 | 0.33 |
| 9 | NN CF random t2w (1.5 – 2.5) | 1.68 | 1.1 | 0.23 | 1.44 | 1.0 | 0.33 | 1.22 | 0.9 | 0.39 |
| 10 | NN CF random t2w (1.8 – 2.5) | 2.32 | 1.9 | 0.21 | 1.91 | 1.0 | 0.26 | 1.70 | 1.8 | 0.49 |
| 11 | NN Fully random (t2w $1.5 - 2.5$) | 1.70 | 1.0 | 0.25 | 1.63 | 0.9 | 0.24 | 1.61 | 0.9 | 0.35 |

Table 4.4: Sim-to-Multi-Real Results.

As a baseline, we use the non-linear controller that is part of the Crazyswarm using default gains ("Mellinger"), where we observe an average Euclidean position error of 0.11 m. A second baseline is the same controller where we remove any computation that requires an additional state ("Mellinger (no memory)"), i.e., we set the gains for the integral terms to zero. As expected, the average position error increases – in this case to 0.2 m.

Our neural network with the motor settling time T = 0.15 has a mean position error of 0.19 m, which is similar to the hand-tuned baseline controller without memory. A network trained without motor delays (T nearly zero) overshoots frequently and has a larger mean position error of 0.21 m. If the network is trained without sensor and motor noise, we measure a mean position error of 0.24 m. The standard deviation of the norm of the position error for the neural networks is nearly twice as high as for the non-linear feedback controller (0.06 and 0.11, respectively). Plots for some controllers are shown in Fig. 4.3.

Note that none of our policies are explicitly trained for trajectory tracking. Nonetheless, they still show competitive tracking performance compared to the baseline trajectory-tracking controller specifically tuned for the Crazyflie.

4.5.3 Sim-to-Multi-Real: Learning without Model

We now investigate how well a single policy works across different quadrotor platforms. In all cases, we quantify the hover quality as well as trajectory tracking using the metrics defined in Section 4.4.5. For the medium quadrotor, we artificially limit the output RPM to 60% of its maximum, to compensate for its higher thrust-to-weight ratio[§]. The results are shown in Table 4.4.

We use three different baselines to provide an estimate on achievable performance. The first two baselines are identical to the Mellinger baselines in Section 4.5.2. A third baseline is used to

[§]Among all parameters the rough value of the thrust-to-weight ratio is relatively easy to estimate.



Figure 4.3: Trajectory tracking performance of a Crazyflie 2.0 using different controllers. The target trajectory is a figure-eight (to be executed in 5.5 s). The Mellinger controller has the lowest tracking error (0.11 m). Our baseline NN controller has lower tracking error (0.19 m) than the Mellinger with integral gains disabled (Mellinger no memory; 0.2 m).

test transferability, in which we find a uniform set of attitude and position gains for the Mellinger controller without any memory, but keep the manually tuned values for gravity compensation in place ("Mellinger (uniform)"). This baseline provides an estimate on how well a single policy might work across different quadrotors.

We compare these "Mellinger" baselines to different networks, including our baseline network (BN), BN when motor delays ignored, and various policies that use randomized models during training. We make the following observations:

- 1. All our policies show somewhat comparable performance to the Mellinger controllers on all platforms. There are no significant oscillations for the learned policies, whereas there are significant oscillations for some of the Mellinger baselines (see rows 2 and 3).
- Unsurprisingly, the network specifically trained for the Crazyflie works best on this platform. It also performs very well on the small quadrotor, but shows large position errors on the medium quadrotor (row 4). Surprisingly, modeling the motor delay during training has a very small impact on tracking performance (row 5).
- 3. Randomization around a set of nominal values can improve the performance, but it works best if the randomization is fairly small (20 % in our case), see rows 6 8. This improvement is not caused just by randomizing different thrust-to-weight ratios (rows 9 and 10).
- 4. Full randomization shows more consistent results over all platforms, but performs not as well as other policies (row 11).

4.5.4 Control Policy Robustness and Recovery

We perform recovery robustness tests by making repetitive throws of the quadrotors in the air. We use the baseline neural network that was trained for the Crazyflie specifically on all platforms. In

these tests we do not perform comparison with the "Mellinger" controller since it could not recover properly from deviations from its goal position larger than 0.5 m. This controller is mainly tuned for trajectory tracking with closely located points provided as full state vectors. Our policy, on the other hand, shows a substantial level of robustness on all three platforms. It performs especially well on the Crazyflie platform recovering from 80 % of all throws and up to 100 % throws with moderate attitude changes ($\leq 35^{\circ}$)[¶]. Interestingly, it can even recover in more than half of scenarios after hitting the ground.

The policy also shows substantial level of robustness on other quadrotors. Similar to the Crazyflie platform, the throws with moderate attitude change do not cause serious problems to any of the platforms and they recover in $\ge 90\%$ of trials. Stronger attitude disturbances are significantly harder and we observe roughly 50% recovery rate on average.

More mild tests, like light pushes and pulling from the hover state, do not cause failures. One observation that we make is that all policies learned some preference on yaw orientation that it tries to stabilize although we did not provide any yaw-related costs apart from the cost on angular velocities. We hypothesize that the policy seeks a "home" yaw angle because it becomes unnecessary to reason about symmetry of rotation in the xy plane if the quadrotor is always oriented in the same direction.

Another surprising observation is that the control policies can deal with much higher initial velocities than those encountered in training (≤ 1 m/s). In practice, initial velocities in our tests often exceed 3 m/s (see Fig. 4.4 for an example of a recovery trajectory). The policies can take-off from the ground, thus overcoming the near-ground airflow effects. They can also fly from distances far exceeding the boundaries of the position initialization box observed in the training. All these

[¶]Computed as a 95-th percentile in our experiments.



Figure 4.4: An example of a recovery trajectory from a random throw with an initial linear velocity of approximately 4 m/s. Trajectory colors correspond to quadrotor's velocities. Arrow colors correspond to the names of the body-frame axes.

factors demonstrate strong generalization of the learned policy to the out-of-distribution states. Our supplemental video shows examples on the achieved robustness with all three platforms.

4.6 Conclusions

In this chapter, we address the problem of data scarcity in learning complex quadrotor policies by combining reinforcement learning and simulated quadrotor models as task priors. We demonstrate how a single neural network policy trained completely in simulation for a task of recovery from harsh initial conditions can generalize to multiple quadrotor platforms with unknown dynamics parameters and no data collection from the target platform. We present a thorough study on the importance of many modeled quadrotor dynamics phenomena for the task of sim-to-real transfer. We investigate a popular domain adaptation technique, called domain randomization, for the purpose of reducing the simulation to reality gap.

Our experiments show the following interesting results. First, it is possible to transfer a single policy trained on a specific quadrotor to multiple real platforms, which significantly vary in sizes and inertial parameters. Second, the transferred policy is capable of generalizing to many out-of-distribution states, including much higher initial velocities and much more distant initial positions. Third, even policies that are trained when ignoring real physical effects (such as motor delays or sensor noise) work robustly on real systems. Modeling such effects explicitly during training improves flight performance slightly. Fourth, the transferred policies show high robustness to harsh initial conditions better than the hand-tuned nonlinear controller we used as a baseline. Fifth, domain randomization is capable of improving results, but the extent of the improvement is moderate in comparison to the baseline performance trained without parameter perturbations.

Our findings open exciting directions for future work. For example, one can incorporate a scarce (limited) number of samples collected from real systems with our policies to improve the trajectory tracking performance of our policy without manual tuning. Alternative path may include investigation of more complex network architectures aimed for explicit or implicit online systems identification. Another interesting avenue of work is exploring learning policies that are robust to failure by training with different failure cases such as broken motors. Finally, our work opens the possibility of investigating end-to-end policy learning and transfer for multi-drone systems.

Chapter 5

Task Specific Learning with Scarce Data via Meta-learned Losses

In Chapter 4, we demonstrate how we can achieve remarkably robust policies by transferring them from the simulation while collecting no samples from the target robotic system. But it would be unreasonable to assume that similar approach can work in many scenarios. In fact, we often can collect a limited (or scarce) amount of data samples from the target system, hence there is a lot of value in developing methods allowing sample efficient adaptation of a learner.

Inspired by the remarkable capability of humans to quickly learn and adapt to new tasks, the concept of learning to learn, or *meta-learning*, recently became popular within the machine learning community [Andrychowicz et al., 2016, Duan et al., 2016b, Finn et al., 2017]. This family of algorithms leverages an available prior on the task distribution at *meta-train* time to learn various components of the learning system allowing sample-efficient adaptation to a specific task at *meta-test* time. We can classify *learning to learn* methods into roughly two categories: approaches that learn representations that can generalize and are easily adaptable to new tasks [Finn et al.,

Parts of this chapter appeared in Sarah Bechtle*, **Artem Molchanov***, Yevgen Chebotar*, Edward Grefenstette, Ludovic Righetti, Gaurav Sukhatme, Franziska Meier. Meta-learning via learned loss. IEEE International Conference on Pattern Recognition, ICPR, 2021. URL: https://arxiv.org/pdf/1906.05374.pdf.


Figure 5.1: Framework overview: The learned meta-loss is used as a learning signal to optimize the optimize f_{θ} , which can be a regressor, a classifier or a control policy.

2017], and approaches that learn how to optimize models [Andrychowicz et al., 2016, Duan et al., 2016b].

In this chapter we investigate the second type of approach. We propose a learning framework that is able to learn any parametric loss function—as long as its output is differentiable with respect to its parameters. Such learned functions can be used to efficiently optimize models for new tasks.

Specifically, the purpose of this work is to encode learning strategies into a parametric loss function, or a *meta-loss*, which generalizes across multiple training contexts or tasks. Inspired by *inverse reinforcement learning* [Ng and Russell, 2000], our work combines the *learning to learn* paradigm of meta-learning with the generality of learning loss landscapes. We construct a unified, fully differentiable framework that can learn optimizee-independent loss functions to provide a strong learning signal for a variety of learning problems, such as classification, regression or reinforcement learning. Our framework involves an inner and an outer optimization loops. In the inner loop, a model or an *optimizee* is trained with gradient descent using the loss coming from

our learned meta-loss function. Fig. 5.1 shows the pipeline for updating the optimizee with the meta-loss. The outer loop optimizes the meta-loss function by minimizing a *task-loss*, such as a standard regression or reinforcement-learning loss, that is induced by the updated optimizee.

The contributions of this work are as follows: i) we present a framework for learning adaptive, high-dimensional loss functions through back-propagation that create the loss landscapes for efficient optimization with gradient descent. We show that our learned meta-loss functions improves over directly learning via the task-loss itself while maintaining the generality of the task-loss. ii) We present several ways our framework can incorporate extra information that helps shape the loss landscapes at *meta-train* time. This extra information can take on various forms, such as exploratory signals or expert demonstrations for RL tasks. After training the meta-loss function, the task-specific losses are no longer required since the training of optimizees can be performed entirely by using the meta-loss function alone, without requiring the extra information given at *meta-train* time. In this way, our meta-loss can find more efficient ways to optimize the original task loss.

We apply our meta-learning approach to a diverse set of problems demonstrating our framework's flexibility and generality. The problems include regression problems, image classification, behavior cloning, model-based and model-free reinforcement learning. Our experiments include empirical evaluation for each of the aforementioned problems.

5.1 Related Work

Meta-learning originates from the concept of learning to learn [Schmidhuber, 1987, Bengio and Bengio, 1990, Thrun and Pratt, 2012]. Recently, there has a been a wide interest in finding ways to improve learning speeds and generalization to new tasks through meta-learning. Let us consider

gradient based learning approaches, that update the parameters of an *optimizee* $f_{\theta}(x)$, with model parameters θ and inputs x as follows:

$$\theta_{\text{new}} = h_{\psi}(\theta, \nabla_{\theta} \mathscr{L}_{\phi}(y, f_{\theta}(x));$$
(5.1)

where we take the gradient of a loss function \mathscr{L} , parametrized by ϕ , with respect to the optimizee's parameters θ and use a gradient transform h, parametrized by ψ , to compute new model parameters θ_{new}^* . In this context, we can divide related work on meta-learning into learning model parameters θ that can be easily adapted to new tasks [Finn et al., 2017, Mendonca et al., 2019, Gupta et al., 2018, Yu et al., 2018], learning optimizer policies h that transform parameters updates with respect to known loss or reward functions [Maclaurin et al., 2015, Andrychowicz et al., 2016, Li and Malik, 2017, Franceschi et al., 2017, Meier et al., 2018, Duan et al., 2016b], or learning loss/reward function representations ϕ [Sung et al., 2017, Houthooft et al., 2018, Zou et al., 2019]. Alternatively, in unsupervised learning settings, meta-learning has been used to learn unsupervised rules that can be transferred between tasks [Metz et al., 2019, Hsu et al., 2019].

Our framework falls into the category of learning loss landscapes. Similar to works by Sung et al. [2017] and Houthooft et al. [2018], we aim at learning loss function parameters ϕ that can be applied to various optimizee models, e.g. regressors, classifiers or agent policies. Our learned loss functions are independent of the model parameters θ that are to be optimized, thus they can be easily transferred to other optimizee models. This is in contrast to methods that meta-learn model-parameters θ directly [e.g. Finn et al., 2017, Mendonca et al., 2019], which are are orthogonal and complementary to ours, where the learned representation θ cannot be separated from the original model of the optimizee. The idea of learning loss landscapes or reward functions

^{*}For simple gradient descent: $h(\theta, \nabla_{\theta} \mathscr{L}(y, f_{\theta}(x)) = \theta - \psi \nabla_{\theta} \mathscr{L}(y, f_{\theta}(x))$

in the reinforcement learning (RL) setting can be traced back to the field of inverse reinforcement learning [Ng and Russell, 2000, Abbeel and Ng, 2004, IRL]. However, in contrast to IRL we do not require expert demonstrations (however we can incorporate them). Instead we use task losses as a measure of the effectiveness of our loss function when using it to update an optimizee.

Closest to our method are the works on *evolved policy gradients* [Houthooft et al., 2018], *teacher networks* [Wu et al., 2018], *meta-critics* [Sung et al., 2017] and *meta-gradient RL* [Xu et al., 2018]. In contrast to using an evolutionary approach [e.g. Houthooft et al., 2018], we design a differentiable framework and describe a way to optimize the loss function with gradient descent in both supervised and reinforcement learning settings. Wu et al. [2018] propose that instead of learning a differentiable loss function directly, a teacher network is trained to predict parameters of a manually designed loss function, whereas each new loss function class requires a new teacher network design and training. In Xu et al. [2018], discount and bootstrapping parameters are learned online to optimize a task-specific meta-objective. Our method does not require manual design of the loss functions are learned entirely from data. Finally, in work by Sung et al. [2017] a *meta-critic* is learned to provide a task-conditional value function, used to train an actor policy. Although training a meta-critic in the supervised setting reduces to learning a loss function as in our work, in the reinforcement learning setting we show that it is possible to use learned loss functions to optimize directly with gradient descent.

5.2 Meta-Learning via Learned Loss

In this chapter, we aim to learn a loss function, which we call *meta-loss*, that is subsequently used to train an *optimizee*, e.g. a classifier, a regressor or a control policy. More concretely, we aim to

learn a meta-loss function \mathcal{M}_{ϕ} with parameters ϕ , that outputs the loss value $\mathcal{L}_{\text{learned}}$ which is used to train an optimize f_{θ} with parameters θ via gradient descent:

$$\theta_{\text{new}} = \theta - \alpha \nabla_{\theta} \mathscr{L}_{\text{learned}}, \tag{5.2}$$

where
$$\mathscr{L}_{\text{learned}} = \mathscr{M}_{\phi}(y, f_{\theta}(x))$$
 (5.3)

where y can be ground truth target information in supervised learning settings or goal and state information for reinforcement learning settings. In short, we aim to learn a loss function that can be used as depicted in Algorithm 2. Towards this goal, we propose an algorithm to learn meta-loss function parameters ϕ via gradient descent.

The key challenge is to derive a training signal for learning the loss parameters ϕ . In the following, we describe our approach to addressing this challenge, which we call Meta-Learning via Learned Loss (ML³).

5.2.1 ML³ for Supervised Learning

We start with supervised learning settings, in which our framework aims at learning a metaloss function $\mathcal{M}_{\phi}(y, f_{\theta}(x))$ that produces the loss value given the ground truth target y and the predicted target $f_{\theta}(x)$. For clarity purposes we constrain the following presentation to learning a meta-loss network that produces the loss value for training a regressor f_{θ} via gradient descent, however the methodology trivially generalizes to classification tasks.

Our meta-learning framework starts with randomly initialized model parameters θ and loss parameters ϕ . The current loss parameters are then used to produce loss value $\mathscr{L}_{\text{learned}} = \mathscr{M}_{\phi}(y, f_{\theta}(x))$. To optimize model parameters θ we need to compute the gradient of the loss value with respect to θ , $\nabla_{\theta}\mathscr{L} = \nabla_{\theta}\mathscr{M}_{\phi}(y, f_{\theta}(x))$. Using the chain rule, we can decompose the gradient computation into the gradient of the loss network with respect to predictions of model $f_{\theta}(x)$ times the gradient of model f with respect to model parameters[†],

$$\nabla_{\theta}\mathscr{M}_{\phi}(y, f_{\theta}(x)) = \nabla_{f}\mathscr{M}_{\phi}(y, f_{\theta}(x)) \nabla_{\theta} f_{\theta}(x).$$
(5.4)

Once we have updated the model parameters $\theta_{new} = \theta - \alpha \nabla_{\theta} \mathscr{L}_{learned}$ using the current meta-loss network parameters ϕ , we want to measure how much learning progress has been made with loss-parameters ϕ and optimize ϕ via gradient descent. Note, that the new model parameters θ_{new} are implicitly a function of loss-parameters ϕ , because changing ϕ would lead to different θ_{new} . In order to evaluate θ_{new} , and through that loss-parameters ϕ , we introduce the notion of a *task-loss* during *meta-train* time. For instance, we use the mean-squared-error (MSE) loss, which is typically used for regression tasks, as a task-loss $\mathscr{L}_{\mathscr{T}} = (y - f_{\theta_{new}}(x))^2$. We now optimize loss parameters ϕ by taking the gradient of $\mathscr{L}_{\mathscr{T}}$ with respect to ϕ as follows[†]:

$$\nabla_{\phi}\mathscr{L}_{\mathscr{T}} = \nabla_{f}\mathscr{L}_{\mathscr{T}} \nabla_{\theta_{\text{new}}} f_{\theta_{\text{new}}} \nabla_{\phi} \theta_{\text{new}} = \nabla_{f}\mathscr{L}_{\mathscr{T}} \nabla_{\theta_{\text{new}}} f_{\theta_{\text{new}}} \nabla_{\phi} \left[\theta - \alpha \nabla_{\theta} \mathbb{E}\left[\mathscr{M}_{\phi}(y, f_{\theta}(x))\right] \right]$$
(5.5)

where we first apply the chain rule and show that the gradient with respect to the meta-loss parameters ϕ requires the new model parameters θ_{new} . We expand θ_{new} as one gradient step on θ based on meta-loss \mathcal{M}_{ϕ} , making the dependence on ϕ explicit.

Optimization of the loss-parameters can either happen after *each inner gradient* step (where inner refers to using the current loss parameters to update θ), or after *M* inner gradient steps with the current meta-loss network \mathcal{M}_{ϕ} .

The latter option requires back-propagation through a chain of all optimizee update steps. In practice we notice that updating the meta-parameters ϕ after each inner gradient update step

[†]Alternatively this gradient computation can be performed using automatic differentiation

| Algorithm 1: ML | ³ at (<i>meta-train</i>) |
|-----------------|---------------------------------------|
|-----------------|---------------------------------------|

| 1: $\phi \leftarrow$ randomly initialize |
|--|
| 2: while not done do |
| 3: $\theta \leftarrow$ randomly initialize |
| 4: $x, y \leftarrow \text{Sample task samples from } \mathscr{T}$ |
| 5: $\mathscr{L}_{\text{learned}} = \mathscr{M}(y, f_{\theta}(x))$ |
| 6: $\theta_{\text{new}} \leftarrow \theta - \alpha \nabla_{\theta} \mathbb{E}_x \left[\mathscr{L}_{\text{learned}} \right]$ |
| 7: $\phi \leftarrow \phi - \eta \nabla_{\phi} \mathscr{L}_{\mathscr{T}}(y, f_{\theta_{\text{new}}})$ |
| 8: end while |
| |
| Algorithm 2: ML ³ at (<i>meta-test</i>) |
| 1: $M \leftarrow \#$ of optimizee updates |
| 2: $\theta \leftarrow$ randomly initialize |
| 3: for $j \in \{0,, M\}$ do |
| 4: $x, y \leftarrow \text{Sample task samples from } \mathscr{T}$ |
| 5: $\mathscr{L}_{\text{learned}} = \mathscr{M}(y, f_{\theta}(x))$ |
| 6: $\theta \leftarrow \theta - \alpha \nabla_{\theta} \mathbb{E}_x [\mathscr{L}_{\text{learned}}]$ |
| 7: end for |

works better. We reset θ after *M* inner gradient steps. We summarize the meta-train phase in Algorithm 1, with one inner gradient step.

5.2.2 ML³ Reinforcement Learning

In this section, we introduce several modifications that allow us to apply the ML³ framework to reinforcement learning problems. Let $\mathscr{M} = (S, A, P, R, p_0, \gamma, T)$ be a finite-horizon Markov Decision Process (MDP), where S and A are state and action spaces, $P : S \times A \times S \rightarrow \mathbb{R}_+$ is a state-transition probability function or system dynamics, $R : S \times A \rightarrow \mathbb{R}$ a reward function, $p_0 : S \rightarrow \mathbb{R}_+$ an initial state distribution, γ a reward discount factor, and T a horizon. Let $\tau = (s_0, a_0, \dots, s_T, a_T)$ be a trajectory of states and actions and $R(\tau) = \sum_{t=0}^{T-1} \gamma^t R(s_t, a_t)$ the trajectory return. The goal of reinforcement learning is to find parameters θ of a policy $\pi_{\theta}(a|s)$ that maximizes the expected discounted reward over trajectories induced by the policy: $\mathbb{E}_{\pi_{\theta}}[R(\tau)]$ where $s_0 \sim p_0, s_{t+1} \sim P(s_{t+1}|s_t, a_t)$ and $a_t \sim \pi_{\theta}(a_t|s_t)$. In what follows, we show how to train a meta-loss network to perform effective policy updates in a reinforcement learning scenario. To apply our ML³ framework, we replace the optimizee f_{θ} from the previous section with a stochastic policy $\pi_{\theta}(a|s)$. We present two applications of ML³ to RL.

ML³ for Model-Based Reinforcement Learning

Model-based RL (MBRL) attempts to learn a policy π_{θ} by first learning a dynamic model P. Intuitively, if the model P is accurate, we can use it to optimize the policy parameters θ . As we typically do not know the dynamics model a-priori, MBRL algorithms iterate between using the current approximate dynamics model P, to optimize the policy π_{θ} such that it maximizes the reward R under P, then use the optimized policy π_{θ} to collect more data which is used to update the model P. In this context, we aim to learn a loss function that is used to optimize policy parameters through our meta-network \mathcal{M} .

Similar to the supervised learning setting we use current meta-parameters ϕ to optimize policy parameters θ under the current dynamics model $P: \theta_{new} = \theta - \alpha \nabla_{\theta} \left[\mathscr{M}_{\phi}(\tau, g) \right]$, where $\tau = (s_0, a_0, \dots, s_T, a_T)$ is the sampled trajectory and the variable g captures some task-specific information, such as the goal state of the agent. To optimize ϕ we again need to define a task loss, which in the MBRL setting can be defined as $\mathscr{L}_{\mathscr{T}}(g, \pi_{\theta_{new}}) = -\mathbb{E}_{\pi_{\theta_{new}}, P}[R_g(\tau_{new})]$, denoting the reward that is achieved under the current dynamics model P. To update ϕ , we compute the gradient of the task loss $\mathscr{L}_{\mathscr{T}}$ wrt. ϕ , which involves differentiating all the way through the reward function, dynamics model and the policy that was updated using the meta-loss \mathscr{M}_{ϕ} . The pseudo-code in Algorithm 3 illustrates the MBRL learning loop. In Algorithm 5, we show the policy optimization procedure during meta-test time. Notably, we have found that in practice, the model of the dynamics P is not needed anymore for policy optimization at meta-test time. The meta-network learns to implicitly represent the gradients of the dynamics model and can produce a loss to optimize the policy directly.

ML³ for Model-Free Reinforcement Learning

Finally, we consider the model-free reinforcement learning (MFRL) case, where we learn a policy without learning a dynamics model. In this case, we can define a surrogate objective, which is independent of the dynamics model, as our task-specific loss [Williams, 1992, Sutton et al., 1999, Schulman et al., 2015]:

$$\mathscr{L}_{\mathscr{T}}(g, \pi_{\theta_{\text{new}}}) = -\mathbb{E}_{\pi_{\theta_{\text{new}}}}[R_g(\tau_{\text{new}})\log\pi_{\theta_{\text{new}}}(\tau_{\text{new}})]$$
(5.6)

— 4

$$= -\mathbb{E}_{\pi_{\theta_{\text{new}}}}\left[R_g(\tau_{\text{new}})\sum_{t=0}^{T-1}\log\pi_{\theta_{\text{new}}}(a_t|s_t)\right].$$
(5.7)

Similar to the MBRL case, the task loss is indirectly a function of the meta-parameters ϕ that are used to update the policy parameters. Although we are evaluating the task loss on full trajectory rewards, we perform policy updates from Eq. (5.2) using stochastic gradient descent (SGD) on the meta-loss with mini-batches of experience (s_i, a_i, r_i) for $i \in \{0, \dots, B-1\}$ with batch size B, similar to Houthooft et al. [2018]. The inputs of the meta-loss network are the sampled states, sampled actions, task information g and policy probabilities of the sampled actions: $\mathcal{M}_{\phi}(s, a, \pi_{\theta}(a|s), g)^{\ddagger}$. In this way, we enable efficient optimization of very high-dimensional policies with SGD provided only with trajectory-based rewards. In contrast to the above MBRL setting, the rollouts used for task-loss evaluation are real system rollouts, instead of simulated

[‡]We notice that in practice, including the policy's distribution parameters directly in the meta-loss inputs, e.g. mean μ and standard deviation σ of a Gaussian policy, works better than including the probability estimate $\pi_{\theta}(a|s)$, as it provides a direct way to update the distribution parameters using back-propagation through the meta-loss.

rollouts. At test time, we use the same policy update procedure as in the MBRL setting, see Algorithm 5.

| Algorithm 3: | ML^3 for N | MBRL (meta | a-train) |
|---------------------|--------------|------------|----------|
|---------------------|--------------|------------|----------|

- φ, ← randomly initialize parameters
 Randomly initialize dynamics model P
- 3: while not done do
- 4: $\theta \leftarrow$ randomly initialize parameters
- 5: $\tau \leftarrow \text{forward unroll } \pi_{\theta} \text{ using } P$
- 6: $\pi_{\theta_{\text{new}}} \leftarrow \text{optimize}(\tau, \mathscr{M}_{\phi}, g, R)$
- 7: $\tau_{\text{new}} \leftarrow \text{forward unroll } \pi_{\theta_{\text{new}}} \text{ using } P$
- 8: Update ϕ to maximize reward under *P*:
- 9: $\phi \leftarrow \phi \eta \nabla_{\phi} \mathscr{L}_{\mathscr{T}}(\tau_{\text{new}})$
- 10: $\tau_{real} \leftarrow \text{roll out } \pi_{\theta_{new}} \text{ on real system}$
- 11: $P \leftarrow$ update dynamics model with τ_{real}
- 12: end while

Algorithm 4: ML³ for MFRL (*meta-train*)

1: $I \leftarrow \#$ of inner steps 2: $\phi \leftarrow$ randomly initialize parameters 3: while not done do $\theta_0 \leftarrow$ randomly initialize policy 4: $\mathscr{T} \leftarrow$ sample training tasks 5: $\tau_0, R_0 \leftarrow \text{roll out policy } \pi_{\theta_0}$ 6: for $i \in \{0, ..., I\}$ do 7: $\pi_{\theta_{i+1}} \leftarrow \text{optimize}(\pi_{\theta_i}, \mathscr{M}_{\phi}, \tau_i, R_i)$ 8: $\tau_{i+1}, R_{i+1} \leftarrow \text{roll out policy } \pi_{\theta_{i+1}}$ 9: $\mathscr{L}^{i}_{\mathscr{T}} \leftarrow \text{compute task-loss } \mathscr{L}^{i}_{\mathscr{T}}(\tau_{i+1}, R_{i+1})$ 10: end for 11: $\mathscr{L}_{\mathscr{T}} \leftarrow \mathbb{E}\left[\mathscr{L}_{\mathscr{T}}^{i}\right]$ 12: $\phi \leftarrow \phi - \eta \nabla_{\phi} \hat{\mathcal{L}}_{\mathcal{T}}$ 13: 14: end while

Algorithm 5: ML^3 for RL (meta-test)

1: $\theta \leftarrow$ randomly initialize policy 2: **for** $j \in \{0, ..., M\}$ **do** 3: $\tau, R \leftarrow$ roll out π_{θ} 4: $\pi_{\theta} \leftarrow$ optimize $(\pi_{\theta}, \mathcal{M}_{\phi}, \tau, R)$ 5: **end for**

5.2.3 Shaping ML³ loss by adding extra loss information during *meta-train*

So far, we have discussed using standard task losses, such as MSE-loss for regression or reward functions for RL settings. However, it is possible to provide more information about the task at *meta-train* time, which can influence the learning of the loss-landscape. We can design our task-losses to incorporate extra penalties; for instance we can extend the MSE-loss with \mathcal{L}_{extra} and weight the terms with β and γ :

$$\mathscr{L}_{\mathscr{T}} = \beta (y - f_{\theta}(x))^2 + \gamma \mathscr{L}_{\text{extra}}$$
(5.8)

In our work, we experiment with 4 different types of extra loss information at *meta-train* time: for supervised learning we show that adding extra information through $\mathcal{L}_{extra} = (\theta - \theta^*)^2$, where θ^* are the optimal regression parameters, can help shape a convex loss-landscape for otherwise non-convex optimization problems; we also show how we can use \mathcal{L}_{extra} to induce a physics prior in robot model learning. For reinforcement learning tasks we demonstrate that by providing additional rewards in the task loss during meta-train time, we can encourage the trained meta-loss to learn exploratory behaviors; and finally also for reinforcement learning tasks, we show how expert demonstrations can be incorporated to learn loss functions which can generalize to new tasks. In all settings, the additional information shapes the learned loss function such that the environment does not need to provide this information during meta-test time.

5.3 Experiments

In this section we evaluate the applicability and the benefits of the learned meta-loss from two different view points. First, we study the benefits of using standard task losses, such as the



(a) Meta-Train Tasks

(b) Meta-Test Tasks

Figure 5.2: Meta-learning for regression (top) and binary classification (bottom) tasks. (a) meta-train task, (b) meta-test tasks

mean-squared error loss for regression, to train the meta-loss in Section 5.3.1. We analyze how a learned meta-loss compares to using a standard task-loss in terms of generalization properties and convergence speed. Second, we study the benefit of adding extra information at *meta-train* time to shape the loss landscape in Section 5.3.2.

5.3.1 Learning to mimic and improve over known task losses

First, we analyze how well our meta-learning framework can learn to mimic and improve over standard task losses for both supervised and reinforcement learning settings. For these experiments, the meta-network is parameterized by a neural network with two hidden layers of 40 neurons each.



Figure 5.3: Meta-learning for regression (top) and binary classification (bottom) tasks. (a) performance of the meta-network on the meta-train task as a function of (outer) meta-train iterations in blue, as compared to SGD using the task-loss directly in orange, (b) average performance of meta-loss on meta-test tasks as a function of the number of gradient update steps

Meta-Loss for Supervised Learning

In this set of experiments, we evaluate how well our meta-learning framework can learn loss functions \mathcal{M}_{ϕ} for regression and classification tasks. In particular, we perform experiments on sine function regression and binary classification of digits. At meta-train time, we randomly draw one task for meta-training (see Fig. 5.2 (a)), and at meta-test time we randomly draw 10 test tasks for regression, and 4 test tasks for classification (Fig. 5.2(b)). For the sine task at meta-train time, we draw 100 data points from function $y = \sin(x - \pi)$, with $x \in [-2.0, 2.0]$. For meta-test time we draw 100 data points from function $y = A \sin(x - \omega)$, with $A \sim [0.2, 5.0]$, $\omega \sim [-\pi, pi]$ and $x \in [-2.0, 2.0]$. We initialize our model f_{θ} to a simple feedforward NN with 2 hidden



Figure 5.4: ML^3 for MBRL: results are averaged across 10 runs. We can see in (a) that the ML^3 loss generalizes well, the loss was trained on the blue trajectories and tested on the orange ones for the PointmassGoal task. ML^3 loss also significantly speeds up learning when compared to the task loss at meta-test time on the PointmassGoal (b) and the ReacherGoal (c) environments.

layers and 40 hidden units each. For the binary classification task f_{θ} is initialized via the *LeNet* architecture [LeCun et al., 1998]. For both experiments we use a fixed learning rate $\alpha = \eta = 0.001$ for both inner (α) and outer (η) gradient update steps. We average results across 5 random seeds, where each seed controls the initialization of both initial model and meta-network parameters, as well as the the random choice of meta-train/test task(s), and visualize them in Fig. 5.2.

We compare the performance of using SGD with the task-loss \mathscr{L} directly (in orange) to SGD using the learned meta-network \mathscr{M}_{ϕ} (in blue), both using a learning rate $\alpha = 0.001$. In Fig. 5.3 (a) we show the average performance of the meta-network \mathscr{M}_{ϕ} as it is being learned, as a function of (outer) meta-train iterations in blue. In both regression and classification tasks, the meta-loss eventually leads to a better performance on the meta-train task as compared to the task loss. In Fig. 5.3 (b) we evaluate SGD using \mathscr{M}_{ϕ} vs SGD using \mathscr{L} on previously unseen (and out-ofdistribution) meta-test tasks as a function of the number of gradient steps. Even on these novel test tasks, our learned \mathscr{M}_{ϕ} leads to improved performance as compared to the task-loss.

Learning Reward functions for Model-based Reinforcement Learning

In the MBRL example, the forward model of the dynamics is represented in both cases by a neural network, the input to the network is the current state and action, the output is the next state of the environment. The tasks consist of a free movement task of a point mass in a 2D space, we call this environment PointmassGoal, and a reaching task with a 2-link 2D manipulator, which we call the ReacherGoal environment. The Pointmass state space is four-dimensional. For PointmassGoal (x, y, \dot{x}, \dot{y}) are the 2D positions and velocities, and the actions are accelerations (\ddot{x}, \ddot{y}) . The ReacherGoal environment for the MBRL experiments is a lower-dimensional variant of the MFRL environment. It has a four dimensional state, consisting of position and angular velocity of the joints $[\theta_1, \theta_2, \dot{\theta_1}, \dot{\theta_2}]$ the torque is two dimensional $[\tau_1, \tau_2]$ The dynamics model P is updated once every 100 outer iterations with the samples collected by the policy from the last inner optimization step of that outer optimization step, i.e. the latest policy. The task distribution $p(\mathcal{T})$ consists of different target positions that either the point mass or the arm should reach. During meta-train time, a model of the system dynamics, represented by a neural network, is learned from samples of the currently optimal policy. The task loss during meta-train time is $\mathscr{L}_{\mathscr{T}}(\theta) = \mathbb{E}_{\pi_{\theta}, P}[R(\tau)]$, where $R(\tau)$ is the final distance from the goal g, when rolling out $\pi_{\theta_{\text{new}}}$ in the dynamics model P. Taking the gradient $\nabla_{\phi} \mathbb{E}_{\pi_{\theta_{\text{new},P}}}[R(\tau)]$ requires the differentiation through the learned model P (see Algorithm 3). The input to the meta-network is the state-action trajectory of the current roll-out and the desired target position. The meta-network outputs a loss signal together with the learning rate to optimize the policy. Fig. 5.4a shows the qualitative reaching performance of a policy optimized with the meta loss during meta-test on PointmassGoal. The meta-loss network was trained only on tasks in the right quadrant (blue trajectories) and tested on the tasks in the left quadrant (orange trajectories) of the x, y plane, showing the generalization capability of the meta loss. Fig. 5.4b and Fig. 5.4c show a comparison in terms of final distance to the target position at test time. The performance of policies trained with the meta-loss is compared to policies trained with the task loss, in this case final distance to the target. The curves show results for 10 different goal positions (including goal positions where the meta-loss needs to generalize). When optimizing with the task loss, we use the dynamics model learned during the meta-train time, as in this case the differentiation through the model is required during test time. As mentioned in Section 5.2.2, this is not needed when using the meta-loss.

Learning Reward functions for Model-free Reinforcement Learning

In the following, we move to evaluating on model-free RL tasks. In our experiments we use two continuous control tasks based on OpenAI Gym MuJoCo environments [Todorov et al., 2012, Brockman et al., 2016]: ReacherGoal and AntGoal. The ReacherGoal environment is a 2-link 2D manipulator that has to reach a specified goal location with its end-effector. The task distribution (at meta-train and meta-test time) consists of an initial link configuration and random goal locations within the reach of the manipulator. The performance metric for this environment is the mean trajectory sum of negative distances to the goal, averaged over 10 tasks. As a trajectory reward $R_g(\tau)$ for the task-loss (see Eq. 5.6) we use $R_g(\tau) = -d + 1/(d + 0.001) - |a_t|$, where d is the distance of the end-effector to the goal g specified as a 2-d Cartesian position. The environment has eleven dimensions specifying angles of each link, direction from the end-effector to the goal, Cartesian coordinates of the target and Cartesian velocities of the end-effector.

The AntGoal environment requires a four-legged agent to run to a goal location. The task distribution consists of random goals initialized on a circle around the initial position. The performance metric for this environment is the mean trajectory sum of differences between the initial and the current distances to the goal, averaged over 10 tasks. Similar to the previous

environment we use $R_g(\tau) = -d + 5/(d + 0.25) - |a_t|$, where *d* is the distance from the center of the creature's torso to the goal *g* specified as a 2D Cartesian position. In contrast to the ReacherGoal this environment has 33 [§] dimensional state space that describes Cartesian position, velocity and orientation of the torso as well as angles and angular velocities of all eight joints. Note that in both environments, the meta-network receives the goal information *g* as part of the state *s* in the corresponding environments. Also, in practice, including the policy's distribution parameters directly in the meta-loss inputs, e.g. mean μ and standard deviation σ of a Gaussian policy, works better than including the probability estimate $\pi_{\theta}(a|s)$, as it provides a more direct way to update θ using back-propagation through the meta-loss.

Fig. 5.5 shows results for our tasks [¶] Fig. 5.5a and Fig. 5.5b show the results of the meta-test time performance for the ReacherGoal and the AntGoal environments respectively. We can see that ML^3 loss significantly improves optimization speed in both scenarios compared to PPO. In our experiments, we observed that on average ML^3 requires 5 times fewer samples to reach 80% of task performance in terms of our metrics for the model-free tasks.

To test the capability of the meta-loss to generalize across different architectures, we first meta-train \mathcal{M}_{ϕ} on an architecture with two layers and meta-test the same meta-loss on architectures with varied number of layers. Fig. 5.5 (c+d) show meta-test time comparison for the ReacherGoal and the AntGoal environments in a model-free setting for four different model architectures. Each curve shows the average and the standard deviation over ten different tasks in each environment. Our comparison clearly indicates that the meta-loss can be effectively re-used across multiple architectures with a mild variation in performance compare to the overall variance of the corresponding task optimization.

[§]In contrast to the original Ant environment we remove external forces from the state.

[¶]Our framework is implemented using open-source libraries *Higher* [Grefenstette et al., 2019] for convenient second-order derivative computations and *Hydra* [Yadan, 2019] for simplified handling of experiment configurations.



Figure 5.5: ML^3 for model-free RL: results are averaged across 10 tasks. (a+b) Policy learning on new task with ML^3 loss compared to PPO objective performance during *meta-test* time. The learned loss leads to faster learning at meta-test time. (c+d) Using the same ML^3 loss, we can optimize policies of different architectures, showing that our learned loss maintains generality.

5.3.2 Shaping loss landscapes by adding extra information at meta-train time

This set of experiments shows that our meta-learner is able to learn loss functions that incorporate extra information available only during meta-train time. The learned loss will be shaped such that optimization is faster when using the meta-loss compared to using a standard loss.

Illustration: Shaping loss

We start by illustrating the loss shaping on an example of sine frequency regression where we fit a single parameter for the purpose of visualization simplicity.

For this illustration we generate training data $\mathscr{D} = \{x_n, y_n\}^N$, N = 1000, by drawing data samples from the ground truth function $y = sin(\nu x)$, for x = [-1, 1]. We create a model $f_{\omega}(x) = sin(\omega x)$, and aim to optimize parameter ω on \mathscr{D} , with the goal of recovering value ν . Fig. 5.6a (bottom) shows the loss landscape for optimizing ω , when using the MSE loss. The target frequency ν is indicated by a vertical red line. As noted by Parascandolo et al. [2017], the landscape of this loss is highly non-convex and difficult to optimize with conventional gradient descent.

Here, we show that by utilizing additional information about the ground truth value of the frequency at meta-train time, we can learn a better shaped loss. Specifically, during meta-train time, our task-specific loss is the squared distance to the ground truth frequency: $(\omega - \nu)^2$ that we later call the *shaping loss*. The inputs of the meta-network $\mathcal{M}_{\phi}(y, \hat{y})$ are the training targets y and predicted function values $\hat{y} = f_{\omega}(x)$, similar to the inputs to the mean-squared loss. After meta-train time commences our learned loss function \mathcal{M}_{ϕ} produces a convex loss landscapes as depicted in Fig. 5.6a(top).

To analyze how the shaping loss impacts model optimization at meta-test time, we compare 3 loss functions: 1) directly using standard MSE loss (orange), 2) ML^3 loss that was trained via the MSE loss as task loss (blue), and 3) ML^3 loss trained via the shaping loss, Fig. 5.6b. When comparing the performance of these 3 losses, it becomes evident that without shaping the loss landscape, the optimization is prone to getting stuck in a local optimum.

112



Figure 5.6: Meta-test time evaluation of the shaped meta-loss (ML^3), i.e. trained with shaping ground-truth (extra) information at meta-train time: a) Comparison of learned ML^3 loss (top) and MSE loss (bottom) landscapes for fitting the frequency of a sine function. The red lines indicate the ground-truth values of the frequency. b) Comparing optimization performance of: ML^3 loss trained with (green), and without (blue) ground-truth frequency values; MSE loss (orange). The ML^3 loss learned with the ground-truth values outperforms both the non-shaped ML^3 loss and the MSE loss. c-d) Comparing performance of inverse dynamics model learning for ReacherGoal (c) and Sawyer arm (d). ML^3 loss trained with (green) and without (blue) ground-truth inertia matrix is compared to MSE loss (orange). The shaped ML^3 loss outperforms the MSE loss in all cases.

Shaping loss via physics prior for inverse dynamics learning

Next, we show the benefits of shaping our ML³ loss via ground truth parameter information for a robotics application. Specifically, we aim to learn and shape a meta-loss that improves sample efficiency for learning (inverse) dynamics models, i.e. a mapping $u = f(q, \dot{q}, \ddot{q}_{des})$, where: $q, \dot{q},$ \ddot{q}_{des} are vectors of joint angular positions, velocities and desired accelerations; u is a vector of joint torques.

Rigid body dynamics (RBD) provides an analytical solution to computing the (inverse) dynamics and can generally be written as:

$$M(q)\ddot{q} + F(q,\dot{q}) = u, \tag{5.9}$$

where the inertia matrix M(q), and $F(q, \dot{q})$ are computed analytically [Featherstone, 2014]. Learning an inverse dynamics model using neural networks can increase the expressiveness compared to RBD but requires many data samples that are expensive to collect. Here we follow the approach in [Lutter et al., 2019], and attempt to learn the inverse dynamics via a neural network that predicts the inertia matrix $M_{\theta}(q)$. To improve upon sample efficiency we apply our method by shaping the loss landscape during meta-train time using the ground truth inertia matrix M(q)provided by a simulator. Specifically, we use the task loss $\mathscr{L}_{\mathcal{T}} = (M_{\theta}(q) - M(q))^2$ to optimize our meta-loss network. During meta-test time we use our trained meta-loss shaped with the physics prior (the inertia matrix exposed by the simulator) to optimize the inverse dynamics neural network. In Fig. 5.6-c we show the prediction performance of the inverse dynamics model during meta-test time on new trajectories of the ReacherGoal environment. We compare the optimization performance during meta-test time when using the meta-loss trained with physics prior, the meta loss trained without physics prior (i.e via MSE loss) to the optimization with MSE loss. Fig. 5.6-d shows a similar comparison for the Sawyer environment - a simulator of the 7 degrees-of-freedom Sawyer anthropomorphic robot arm [Sawyer, 2012]. Inverse dynamics learning using the meta loss with physics prior achieves the best prediction performance on both robots. ML³ without physics prior performs worst on the ReacherGoal environment, in this case the task loss formulated only in the action space did not provide enough information to learn a $\mathcal{L}_{learned}$ useful for optimization. For the Sawyer training with MSE loss leads to a slower optimization, however the asymptotic performance of MSE and ML³ is the same. Only ML³ with shaped loss outperforms both.



Figure 5.7: (a) MountainCar trajectory for policy optimized with iLQR compared to ML^3 loss with extra information. (b) optimization performance during meta-test time for policies optimized with iLQR compared to ML^3 with and without extra information.

Shaping Loss via intermediate goal states for RL

We analyze loss landscape shaping on the MountainCar environment [Moore, 1990], a classical control problem where an under-actuated car has to drive up a steep hill. The propulsion force generated by the car does not allow steady climbing of the hill, thus greedy minimization of the distance to the goal often results in a failure to solve the task. The state space is two-dimensional

consisting of the position and velocity of the car, the action space consists of a one-dimensional torque. In our experiments, we provide intermediate goal positions during meta-train time, which are not available during the meta-test time. The meta-network incorporates this behavior into its loss leading to an improved exploration during the meta-test time as can be seen in Fig. 5.7-a, when compared to a classical iLQR-based trajectory optimization [Tassa et al., 2014]. Fig. 5.7-b shows the average distance between the car and the goal at last rollout time step over several iterations of policy updates with ML³ with and without extra information and iLQR. As we observe, ML³ with extra information can successfully bring the car to the goal in a small amount of updates, whereas iLQR and ML³ without extra information is not able to solve this task.



Figure 5.8: ReacherGoal with expert demonstrations available during meta-train time. (a) shows the targets in end-effector space. The four blue dots show the training targets for which expert demonstrations are available, the orange dots show the meta-test targets. In (b) we show the reaching performance of a policy trained with the shaped ML^3 loss at meta-test time, compared to the performance of training simply on the behavioral cloning objective and testing on test targets.

Shaping loss via expert information during meta-train time

Expert information, like demonstrations for a task, is another way of adding relevant information during meta-train time, and thus shaping the loss landscape. In *learning from demonstration* (*LfD*) [Pomerleau, 1991, Ng and Russell, 2000, Billard et al., 2008], expert demonstrations are used for initializing robotic policies. In our experiments, we aim to mimic the availability of an expert at meta-test time by training our meta-network to optimize a behavioral cloning objective at meta-train time. We provide the meta-network with expert state-action trajectories during train time, which could be human demonstrations or, as in our experiments, trajectories optimized using iLQR. During meta-train time, the task loss is the behavioral cloning objective $\mathscr{L}_{\mathscr{T}}(\theta) = \mathbb{E}\left[\sum_{t=0}^{T-1} [\pi_{\theta_{new}}(a_t|s_t) - \pi_{expert}(a_t|s_t)]^2\right]$. Fig. 5.8b shows the results of our experiments in the ReacherGoal environment.

5.4 Conclusions

In this chapter we present a framework that allows leveraging scarce data available on a target task at meta-test time by extracting relevant task priors from a distribution of tasks at meta-train time and encoding them into meta-learned losses. We showed how the meta-learned loss can become well-conditioned and suitable for an efficient optimization with gradient descent. When using the learned meta-loss we observe significant speed improvements in regression, classification and benchmark reinforcement learning tasks. Furthermore, we show that by introducing additional guiding information during training time we can train our meta-loss to develop exploratory strategies that can significantly improve performance during the meta-test time.

We believe that the ML³ framework is a powerful tool to incorporate prior experience and transfer learning strategies to new tasks. Future directions of our work include combining multiple

learned meta-loss functions in order to generalize over different families of tasks. Another interesting direction is further introduction of additional curiosity rewards during training time to improve the exploration strategies learned by the meta-loss.

Chapter 6

Conclusion and Future Work

In this dissertation, we provide multiple examples of how we can operate robotic systems under scarcity of data by leveraging learning and/or priors on the structure of the target task.

In the first part, we consider the scenarios of scarcity of perceptual (or feedback) data needed to close control loops. First, we show how one can utilize high-dimensional and multi-modal tactile feedback in combination with machine learning techniques to infer locations of external contacts when more traditional sources of measurements are completely absent or scarce. Second, we demonstrate a challenging scenario of data scarcity resulting from expensive measurements with an example of a swarm of highly underactuated vehicles called drifters. We show how one can leverage temporal and spatial smoothness of ocean currents and other drifters to densify measurements by adaptively delaying and stochastically distributing the work.

In the second part, we study the problems of scarcity of data required for task identification. First, we consider the problems of acquiring very robust low-level stabilizing controllers for multiple drones (with differing physical properties) without collecting any data from the target drone. We show how a combination of reinforcement learning and a structured prior in the form of a simulator can produce exceptionally robust and transferable control policies. Second, we address the problem of leveraging a few available samples from the target task, by encoding appropriate priors into the meta-learned losses. We demonstrate wide applicability of this approach on a family of control and classical machine learning problems.

6.1 Conclusion

There are a few key insights from our work. First, there is a lot of value in leveraging alternative sensory modalities and robot collaborations when it comes to feedback scarcity. If a single modality can not provide the full information about the state, other modalities, given a set of simple spatio-temporal assumptions about the problem's structure, can effectively complement the others.

Even for uniquely hard problems with a lot of uncertainty, effective solutions can be relatively simple and involve a set of heuristics. Although heuristics often do not produce optimal solutions, they retain better generalization capability since they do not rely on implicitly build models of environments that can overfit to the training state distribution and can be fragile in the presence of uncertainty. This is also seen from examples of human behavior. We tend to extract a set of simple heuristics [Gigerenzer and Brighton, 2009] from our experiences instead of building precise internal models of processes.

Simulators can provide an important prior for learning transferable policies using reinforcement learning. We demonstrate in our work that even in well-studied problems RL can find uniquely good solutions. We believe this comes from i) relaxation of assumptions in the hierarchical nature of controllers that are usually required for provably stable solutions; and ii) RL does not always find fragile controllers that rely on precise modelling, but can learn robust heuristics that generalize beyond the assumptions of the simulator.

6.2 Future Work

There is still a long path to traverse until robots will be able to operate in unstructured real world settings. One of the key challenges is incremental adaptation using scarce data available to the robot. Robots should be able to disentangle samples of useful data from the entire information flow, efficiently learn complex mappings from new sources, switch in a timely way between sources of information, and quickly adapt based on a few samples available. Our work advances the state-of-the-art a few steps toward this goal and leads to several important future work directions.

- **Perceptual Novelty Detection and Data Attention Models.** In the constantly changing conditions of the real world the robot must be able to effectively distinguish i) relevant data for state estimation, such as, low-noise data allowing high observability of states; and ii) novel and task-relevant observations to incorporate only important experiences. This mitigates the problem of catastrophic forgetting and allows better fitting of important mappings.
- **Effective Selection and Expansions of Priors.** As the robot discovers novel observations, it should not learn how to use them from scratch. Humans effectively re-use already learned representations to bootstrap learning. There is a great value in finding algorithms that together with the novelty detection models would enable to constantly grow and re-use learned perceptual mappings and motor skills.
- **Better Models of World Representation.** Besides the re-use of the proper priors for bootstrapping, we also need to incorporate structures into model architectures that allow better sample efficiency. Since humans tend to think in terms of objects and their relations, we believe that such inductive biases should be reflected in the structures of perceptual systems.

- Sim-to-Real and Real-to-Sim. As our computational resources and understanding of the problems develop, we will be able to run progressively more complex systems in realtime [Angles et al., 2019, Holden et al., 2019, Gao et al., 2018, Fei et al., 2018, Lee et al., 2019] and super-realtime. Nonetheless, initially provided simulators will always have some discrepancy with the real world and may not account for the constantly changing distribution of tasks that the robot will experience during its lifespan. Consequently, we need to create a feedback loop that enables adaptation of the simulator. This adaptation should reflect not only the changes in the goal distribution, but also the unmodeled dynamics and novel observations.
- **Learning of Safe Exploration Policies.** Robots should be able to explore in a safe manner to humans and themselves. Outside of ethical and financial concerns, this allows efficient acquisition of new task-relevant data and creates a good curriculum for learning of novel skills and representations.
- Meta Learning and Automatic Discovery of Inductive Biases. Automating machine learning engineering is the next step in modern artificial intelligence. It includes self discovering of more appropriate architectures for the learners and the structures of the learning algorithms themselves. Combined with aforementioned mechanisms, it creates unlimited potential for life-long adaptation of robotic systems. For robotic systems, this opens the possibility of working in unstructured and ever-changing conditions of the real world.

Bibliography

Pieter Abbeel and Andrew Y. Ng. Apprenticeship learning via inverse reinforcement learning. In Carla E. Brodley, editor, *Machine Learning, Proceedings of the Twenty-first International Conference (ICML 2004), Banff, Alberta, Canada, July 4-8, 2004*, volume 69 of *ACM International Conference Proceeding Series*. ACM, 2004. doi: 10.1145/1015330.1015430. URL https://doi.org/10.1145/1015330.1015430.

Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. *CoRR*, abs/1609.08675, 2016. URL http://arxiv.org/abs/1609.08675.

Marcin Andrychowicz, Misha Denil, Sergio Gomez Colmenarejo, Matthew W. Hoffman, David Pfau, Tom Schaul, and Nando de Freitas. Learning to learn by gradient descent by gradient descent. In *NeurIPS*, pages 3981–3989, 2016. URL http://dblp.uni-trier.de/db/conf/nips/nips2016. html#AndrychowiczDCH16.

Baptiste Angles, Daniel Rebain, Miles Macklin, Brian Wyvill, Loic Barthe, John P. Lewis, Javier von der Pahlen, Shahram Izadi, Julien P. C. Valentin, Sofien Bouaziz, and Andrea Tagliasacchi. VIPER: volume invariant position-based elastic rods. *Proc. ACM Comput. Graph. Interact. Tech.*, 2(2):19:1–19:26, 2019. doi: 10.1145/3340260. URL https://doi.org/10.1145/3340260.

Rika Antonova, Silvia Cruciani, Christian Smith, and Danica Kragic. Reinforcement learning for pivoting task. *CoRR*, abs/1703.00472, 2017. URL http://arxiv.org/abs/1703.00472.

Sarah Bechtle, Artem Molchanov, Yevgen Chebotar, Edward Grefenstette, Ludovic Righetti, Gaurav S. Sukhatme, and Franziska Meier. Meta-learning via learned loss. In *IEEE International Conference on Pattern Recognition, ICPR, Milan, Italy, January 10-15, 2021*. IEEE, 2021. URL http://arxiv.org/abs/1906.05374.

Yoshua Bengio and Samy Bengio. Learning a synaptic learning rule. Technical Report 751, Département d'Informatique et de Recherche Opérationelle, Université de Montréal, Montreal, Canada, 1990.

Yoshua Bengio, Aaron C. Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828, 2013. doi: 10.1109/TPAMI.2013.50. URL https://doi.org/10.1109/TPAMI.2013.50.

José A Bernabé, Javier Felip, Angel P Del Pobil, and Antonio Morales. Contact localization through robot and object motion from point clouds. In 2013 13th IEEE-RAS International Conference on Humanoid Robots (Humanoids), pages 268–273. IEEE, 2013.

Aude Billard, Sylvain Calinon, Rüdiger Dillmann, and Stefan Schaal. Robot programming by demonstration. In *Springer Handbook of Robotics*, pages 1371–1394. Springer, 2008.

Jeannette Bohg, Karol Hausman, Bharath Sankaran, Oliver Brock, Danica Kragic, Stefan Schaal, and Gaurav S. Sukhatme. Interactive perception: Leveraging action in perception and perception in action. *IEEE Trans. Robotics*, 33(6):1273–1291, 2017. URL https://doi.org/10.1109/TRO. 2017.2721939.

Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *CoRR*, abs/1606.01540, 2016. URL http://arxiv.org/abs/1606.01540.

Herman Bruyninckx, Stefan Dutre, and Joris De Schutter. Peg-on-hole: a model based solution to peg and hole alignment. In *Robotics and Automation, 1995. Proceedings., 1995 IEEE International Conference on*, volume 2, pages 1919–1924. IEEE, 1995.

Yevgen Chebotar, Karol Hausman, Zhe Su, Gaurav S. Sukhatme, and Stefan Schaal. Selfsupervised regrasping using spatio-temporal tactile features and reinforcement learning. In *IEEE Int. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, 2016.

Yevgen Chebotar, Ankur Handa, Viktor Makoviychuk, Miles Macklin, Jan Issac, Nathan D. Ratliff, and Dieter Fox. Closing the sim-to-real loop: Adapting simulation randomization with real world experience. In *IEEE Intl Conf. on Robotics and Automation (ICRA)*, 2019. URL http://arxiv.org/abs/1810.05687.

Kan Chen, Jiang Wang, Liang-Chieh Chen, Haoyuan Gao, Wei Xu, and Ram Nevatia. ABC-CNN: an attention based convolutional neural network for visual question answering. *CoRR*, abs/1511.05960, 2015. URL http://arxiv.org/abs/1511.05960.

Paul F. Christiano, Zain Shah, Igor Mordatch, Jonas Schneider, Trevor Blackwell, Joshua Tobin, Pieter Abbeel, and Wojciech Zaremba. Transfer from simulation to real world through learning deep inverse dynamics model. *CoRR*, abs/1610.03518, 2016. URL http://arxiv.org/abs/1610.03518.

Vlad Ciobanu, Decebal Popescu, and Adrian Petrescu. Point of contact location and normal force estimation using biomimetical tactile sensors. In *Complex, Intelligent and Software Intensive Systems (CISIS), 2014 Eighth International Conference on*, pages 373–378. IEEE, 2014.

Craig Corcoran and Robert Platt. A measurement model for tracking hand-object state during dexterous manipulation. In *Robotics and Automation (ICRA), 2010 IEEE International Conference* on, pages 4302–4308. IEEE, 2010.

Caitlin Dawson. Could this nearly invincible drone be the future of disaster relief? *News*, Feb 2020. URL https://viterbischool.usc.edu/news/2020/02/ could-this-nearly-invincible-drone-be-the-future-of-disaster-relief/.

Vishnu R. Desaraju and Nathan Michael. Hierarchical adaptive planning in environments with uncertain, spatially-varying disturbance forces. In 2014 IEEE International Conference on Robotics and Automation, ICRA 2014, Hong Kong, China, May 31 - June 7, 2014, pages 5171–5176. IEEE, 2014. doi: 10.1109/ICRA.2014.6907618. URL https://doi.org/10.1109/ICRA.2014.6907618.

Pedro M. Domingos. A few useful things to know about machine learning. *Commun. ACM*, 55(10): 78–87, 2012. doi: 10.1145/2347736.2347755. URL https://doi.org/10.1145/2347736.2347755.

Richard C. Dorf and Robert H. Bishop. *Modern Control Systems*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 9th edition, 2000. ISBN 0130306606.

Yan Duan, Xi Chen, Rein Houthooft, John Schulman, and Pieter Abbeel. Benchmarking deep reinforcement learning for continuous control. In *Intl Conf. on Machine Learning (ICML)*, pages 1329–1338, 2016a. URL http://jmlr.org/proceedings/papers/v48/duan16.html.

Yan Duan, John Schulman, Xi Chen, Peter L. Bartlett, Ilya Sutskever, and Pieter Abbeel. Rl²: Fast reinforcement learning via slow reinforcement learning. *CoRR*, abs/1611.02779, 2016b. URL http://dblp.uni-trier.de/db/journals/corr/corr1611.html#DuanSCBSA16.

M. Dunbabin. Optimal 4D Path-Planning in Strongly Tidal Coastal Environments: Application to AUVs and Profiling Drifters. In *RSS 2012 Workshop on Robotics for Environmental Monitoring, RSS 2012*, 2012a.

M. Dunbabin. Optimal 4D Path-Planning in Strongly Tidal Coastal Environments: Application to AUVs and Profiling Drifters. In *Proc. of the RSS 2012 Workshop on Robotics for Environmental Monitoring*, 2012b.

C. C. Eriksen, T. J. Osse, R. D. Light, T. Wen, T. W. Lehman, P. L. Sabin, J. W. Ballard, and A. M. Chiodi. Seaglider: a long-range autonomous underwater vehicle for oceanographic research. *IEEE Journal of Oceanic Engineering*, 26(4):424–436, Oct 2001. ISSN 0364-9059. doi: 10.1109/48.972073.

Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In Evangelos Simoudis, Jiawei Han, and Usama M. Fayyad, editors, *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96), Portland, Oregon, USA*, pages 226–231. AAAI Press, 1996.

Roy Featherstone. Rigid body dynamics algorithms. Springer, 2014.

Yun (Raymond) Fei, Christopher Batty, Eitan Grinspun, and Changxi Zheng. A multi-scale model for simulating liquid-fabric interactions. *ACM Trans. Graph.*, 37(4):51:1–51:16, 2018. doi: 10.1145/3197517.3201392. URL https://doi.org/10.1145/3197517.3201392.

Jerome A. Feldman, Karl K. Pingle, Thomas O. Binford, Gilbert Falk, Alan C. Kay, R. Paul, Robert F. Sproull, and Jay M. Tenenbaum. The use of vision and manipulation to solve the "instant insanity" puzzle. In *Proceedings of the 2Nd International Joint Conference on Artificial Intelligence*, IJCAI'71, pages 359–364, San Francisco, CA, USA, 1971. Morgan Kaufmann Publishers Inc. URL http://jcai.org/Proceedings/71/Papers/031.pdf.

Javier Felip, Antonio Morales, and Tamim Asfour. Multi-sensor and prediction fusion for contact detection and localization. In *2014 IEEE-RAS Int. Conf. on Humanoid Robots*, pages 601–607. IEEE, 2014.

Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017. URL http://dblp.uni-trier.de/db/conf/icml/icml2017.html# FinnAL17.

Sébastien Forestier, Yoan Mollard, and Pierre-Yves Oudeyer. Intrinsically motivated goal exploration processes with automatic curriculum learning. *CoRR*, abs/1708.02190, 2017. URL http://arxiv.org/abs/1708.02190.

Julian Förster. System identification of the crazyflie 2.0 nano quadrocopter. BA Thesis, ETH Zurich, 2015. URL https://doi.org/10.3929/ethz-b-000214143.

Luca Franceschi, Michele Donini, Paolo Frasconi, and Massimiliano Pontil. Forward and reverse gradient-based hyperparameter optimization. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1165–1173. JMLR. org, 2017.

Fadri Furrer, Michael Burri, Markus Achtelik, and Roland Siegwart. *RotorS—A Modular Gazebo MAV Simulator Framework*, pages 595–625. Springer Intl Publishing, 2016. doi: 10.1007/978-3-319-26054-9_23. URL https://doi.org/10.1007/978-3-319-26054-9_23.

Ming Gao, Andre Pradhana, Xuchen Han, Qi Guo, Grant Kot, Eftychios Sifakis, and Chenfanfu Jiang. Animating fluid sediment mixture in particle-laden flows. *ACM Trans. Graph.*, 37(4):149:1–149:11, 2018. doi: 10.1145/3197517.3201309. URL https://doi.org/10.1145/3197517.3201309.

Alessandro Gasparetto, Paolo Boscariol, Albano Lanzutti, and Renato Vidoni. Trajectory planning in robotics. *Mathematics in Computer Science*, 6(3):269–279, 2012. URL http://dblp.uni-trier.de/ db/journals/mics/.html#GasparettoBLV12.

Gerd Gigerenzer and Henry Brighton. Homo heuristicus: Why biased minds make better inferences. *Topics in Cognitive Science*, 1(1):107–143, 2009. doi: 10.1111/j.1756-8765.2008. 01006.x. URL https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1756-8765.2008.01006.x.

Antony W. Goodwin and Heather E. Wheat. Sensory signals in neural populations underlying tactile perception and manipulation. *Annual Review of Neuroscience*, 27:53–77, 2004.

John Gould, Dean Roemmich, Susan Wijffels, Howard Freeland, Mark Ignaszewsky, Xu Jianping, Sylvie Pouliquen, Yves Desaubies, Uwe Send, Kopillil Radhakrishnan, Kensuke Takeuchi, Kuh Kim, Mikhail Danchenkov, Phil Sutton, Brian King, Breck Owens, and Steve Riser. Argo profiling floats bring new era of in situ ocean observations. *Eos, Transactions American Geophysical Union*, 85(19):185–191, 2004. ISSN 2324-9250. doi: 10.1029/2004EO190002. URL http://dx.doi.org/10.1029/2004EO190002.

Edward Grefenstette, Brandon Amos, Denis Yarats, Phu Mon Htut, Artem Molchanov, Franziska Meier, Douwe Kiela, Kyunghyun Cho, and Soumith Chintala. Generalized inner loop metalearning. *arXiv preprint arXiv:1910.01727*, 2019. URL https://arxiv.org/abs/1910.01727.

Abhishek Gupta, Russell Mendonca, YuXuan Liu, Pieter Abbeel, and Sergey Levine. Metareinforcement learning of structured exploration strategies. In *Advances in Neural Information Processing Systems*, pages 5302–5311, 2018.

Younghee Han, Raymond A. de Callafon, Jorge Cortés, and J. Jaffe. Dynamic modeling and pneumatic switching control of a submersible drogue. In Joaquim Filipe, Juan Andrade-Cetto, and Jean-Louis Ferrier, editors, *ICINCO 2010, Proceedings of the 7th International Conference on Informatics in Control, Automation and Robotics, Volume 2, Funchal, Madeira, Portugal, June 15-18, 2010*, pages 89–97. INSTICC Press, 2010.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In 2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015, pages 1026–1034, 2015. URL https://doi.org/10.1109/ICCV.2015.123.

Nicholas J. Higham. Matrix nearness problems and applications. In M. J. C. Gover and S. Barnett, editors, *Applications of Matrix Theory*, pages 1–27. Oxford University Press, 1989.

Daniel Holden, Bang Chi Duong, Sayantan Datta, and Derek Nowrouzezahrai. Subspace neural physics: fast data-driven interactive simulation. In Sung-Hee Lee, Craig Schroeder, Stephen N. Spencer, Christopher Batty, and Jin Huang, editors, *Proceedings of the 18th annual ACM SIGGRAPH/Eurographics Symposium on Computer Animation, SCA 2019, Los Angeles, CA, USA, July 26-28, 2019*, pages 6:1–6:12. ACM, 2019. doi: 10.1145/3309486.3340245. URL https://doi.org/10.1145/3309486.3340245.

Rein Houthooft, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. VIME: variational information maximizing exploration. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, Barcelona, Spain*, pages 1109–1117, 2016. URL http://papers.nips.cc/paper/6591-vime-variational-information-maximizing-exploration.

Rein Houthooft, Yuhua Chen, Phillip Isola, Bradly C. Stadie, Filip Wolski, Jonathan Ho, and Pieter Abbeel. Evolved policy gradients. In *NeurIPS*, pages 5405–5414, 2018. URL http://dblp.uni-trier.de/db/conf/nips/nips2018.html#HouthooftCISWHA18.

Geir E Hovland and Brenan J McCarragher. Combining force and position measurements for the monitoring of robotic assembly. In *Intelligent Robots and Systems*, 1997. IROS'97., Proceedings of the 1997 IEEE/RSJ Int. Conf. on, volume 2, pages 654–660. IEEE, 1997.

Kyle Hsu, Sergey Levine, and Chelsea Finn. Unsupervised learning via meta-learning. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019. URL https://openreview.net/forum?id=r1My6sR9tX.

Helen H Hu, Amy A Gooch, William B Thompson, Brian E Smits, John J Rieser, and Peter Shirley. Visual cues for imminent object contact in realistic virtual environment. In *Proceedings of the conference on Visualization'00*, pages 179–185. IEEE Computer Society Press, 2000.

Jemin Hwangbo, Inkyu Sa, Roland Siegwart, and Marco Hutter. Control of a quadrotor with reinforcement learning. *IEEE Robotics and Automation Letters (RA-L)*, 2(4):2096–2103, 2017. doi: 10.1109/LRA.2017.2720851. URL https://doi.org/10.1109/LRA.2017.2720851.

Tomohiko Ishikawa, Shigeyuki Sakane, Tomomasa Sato, and Hideo Tsukune. Estimation of contact position between a grasped object and the environment based on sensor fusion of vision and force. *1996 IEEE/SICE/RSJ International Conference on Multisensor Fusion and Integration for Intelligent Systems (Cat. No.96TH8242)*, pages 116–123, 1996. URL https://ieeexplore.ieee.org/document/572167/authors#authors.

Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z. Leibo, David Silver, and Koray Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net, 2017. URL https://openreview.net/forum?id=SJ6yPD5xg.

Stephen James, Andrew J. Davison, and Edward Johns. Transferring end-to-end visuomotor control from simulation to real world for a multi-stage task. In *Conference on Robot Learning* (*CoRL*), pages 334–343, 2017.

Roland S Johansson and J Randall Flanagan. Coding and use of tactile signals from the fingertips in object manipulation tasks. *Nature Reviews Neuroscience*, 10(5):345–359, 2009a.

Roland S. Johansson and Randall J. Flanagan. Coding and use of tactile signals from the fingertips in object manipulation tasks. *Nature Reviews Neuroscience*, 10(5):345–359, April 2009b. ISSN 1471-003X. doi: 10.1038/nrn2621. URL http://dx.doi.org/10.1038/nrn2621.

J. Jouffroy, Q. Zhou, and O. Zielinski. On Active Current Selection for Lagrangian Profilers. *Modeling, Identification and Control*, 34(1):1–10, 2013a.

Jerome Jouffroy, Qiuyang Zhou, and Oliver Zielinski. On Active Current Selection for Lagrangian Profilers. *Modeling, Identification and Control*, 34(1):1–10, 2013b. doi: 10.4173/mic.2013.1.1.

Makoto Kaneko and Kazuo Tanie. Contact point detection for grasping an unknown object using self-posture changeability. *IEEE Transactions on Robotics and Automation*, 10(3):355–367, 1994.

Katie Kang, Suneel Belkhale, Gregory Kahn, Pieter Abbeel, and Sergey Levine. Generalization through simulation: Integrating simulated and real data into deep reinforcement learning for vision-based autonomous flight. In *IEEE Intl Conf. on Robotics and Automation (ICRA)*, 2019.

Yiannis Karayiannidis, Christian Smith, Francisco E Vina, and Danica Kragic. Online contact point estimation for uncalibrated tool use. In 2014 IEEE Int. Conf. on Robotics and Automation (ICRA), pages 2488–2494. IEEE, 2014.

William Koch, Renato Mancuso, and Azer Bestavros. Neuroflight: Next generation flight control firmware. *CoRR*, abs/1901.06553, 2019. URL http://arxiv.org/abs/1901.06553.

Gurdayal S Koonjul, Garth J Zeglin, and Nancy S Pollard. Measuring contact points from displacements with a compliant, articulated robot hand. In *Robotics and Automation (ICRA)*, 2011 IEEE International Conference on, pages 489–495. IEEE, 2011.

Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS 2012*, pages 1097–1105, 2012.

O. Kroemer, C.H. Lampert, and J. Peters. Learning dynamic tactile sensing with robust visionbased training. *IEEE Transactions on Robotics (T-Ro)*, 27(3):545–557, 2011.

A. Kwok and S. Martínez. A coverage algorithm for drifters in a river environment. In *Proceedings* of the 2010 American Control Conference, pages 6436–6441, June 2010. doi: 10.1109/ACC. 2010.5531467.

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, Nov 1998.

Seunghwan Lee, Moon Seok Park, Kyoung-Min Lee, and Jehee Lee. Scalable muscle-actuated human simulation and control. *ACM Trans. Graph.*, 38(4):73:1–73:13, 2019. doi: 10.1145/3306346.3322972. URL https://doi.org/10.1145/3306346.3322972.

Ke Li and Jitendra Malik. Learning to optimize. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net, 2017. URL https://openreview.net/forum?id=ry4Vrt5gl.

Rui Li, Robert Platt, Wenzhen Yuan, Andreas ten Pas, Nathan Roscup, Mandayam A Srinivasan, and Edward Adelson. Localization and manipulation of small parts using gelsight tactile sensing. In 2014 IEEE/RSJ International Conference on Intelligent Robots and Systems, pages 3988–3993. IEEE, 2014.

Nejc Likar and Leon Zlajpah. External joint torque-based estimation of contact information. *International Journal of Advanced Robotic Systems*, 11, 2014.
Kendall Lowrey, Svetoslav Kolev, Jeremy Dao, Aravind Rajeswaran, and Emanuel Todorov. Reinforcement learning for non-prehensile manipulation: Transfer from simulation to physical system. In *IEEE Intl Conf. on Simulation, Modeling, and Programming for Autonomous Robots* (*SIMPAR*), pages 35–42, 2018. doi: 10.1109/SIMPAR.2018.8376268. URL https://doi.org/10.1109/SIMPAR.2018.8376268.

Michael Lutter, Christian Ritter, and Jan Peters. Deep lagrangian networks: Using physics as model prior for deep learning. *CoRR*, abs/1907.04490, 2019. URL http://arxiv.org/abs/1907.04490.

Dougal Maclaurin, David Duvenaud, and Ryan Adams. Gradient-based hyperparameter optimization through reversible learning. In *International Conference on Machine Learning*, pages 2113–2122, 2015.

Marianna Madry, Liefeng Bo, Danica Kragic, and Dieter Fox. ST-HMP: Unsupervised Spatio-Temporal Feature Learning for Tactile Data. In *IEEE International Conference on Robotics and Automation (ICRA)*, may 2014.

Yasushi Makihara, Masao Takizawa, Yoshiaki Shirai, and Nobutaka Shimada. Object recognition under various lighting conditions. In Josef Bigün and Tomas Gustavsson, editors, *Image Analysis, 13th Scandinavian Conference, SCIA 2003, Halmstad, Sweden, June 29 - July 2, 2003, Proceedings*, volume 2749 of *Lecture Notes in Computer Science*, pages 899–906. Springer, 2003. doi: 10.1007/3-540-45103-X_119. URL https://doi.org/10.1007/3-540-45103-X_119.

Fabian Manhardt, Diego Martín Arroyo, Christian Rupprecht, Benjamin Busam, Tolga Birdal, Nassir Navab, and Federico Tombari. Explaining the ambiguity of object detection and 6d pose from visual data. In 2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019, pages 6840–6849. IEEE, 2019. doi: 10.1109/ICCV.2019.00694. URL https://doi.org/10.1109/ICCV.2019.00694.

Philippe Martin and Erwan Salaün. The true role of accelerometer feedback in quadrotor control. In *IEEE Intl Conf. on Robotics and Automation (ICRA)*, pages 1623–1629, 2010. doi: 10.1109/ROBOT.2010.5509980. URL https://doi.org/10.1109/ROBOT.2010.5509980.

Franziska Meier, Daniel Kappler, and Stefan Schaal. Online learning of a memory for learning rates. In 2018 IEEE International Conference on Robotics and Automation (ICRA), pages 2425–2432. IEEE, 2018.

Russell Mendonca, Abhishek Gupta, Rosen Kralev, Pieter Abbeel, Sergey Levine, and Chelsea Finn. Guided meta-policy search. *arXiv preprint arXiv:1904.00956*, 2019.

Luke Metz, Niru Maheswaranathan, Brian Cheung, and Jascha Sohl-Dickstein. Learning unsupervised learning rules. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=HkNDsiC9KQ.

Matthew Michini, M. Ani Hsieh, Eric Forgoston, and Ira B. Schwartz. Robotic tracking of coherent structures in flows. *IEEE Trans. Robotics*, 30(3):593–603, 2014. doi: 10.1109/TRO. 2013.2295655. URL https://doi.org/10.1109/TRO.2013.2295655.

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin A. Riedmiller, Andreas Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015. doi: 10.1038/nature14236. URL https://doi.org/10.1038/ nature14236.

Artem Molchanov, Andreas Breitenmoser, and Gaurav S. Sukhatme. Active drifters: Sailing with the ocean currents. In *RSS Workshop on Autonomous Control, Adaptation, and Learning for Underwater Vehicles, 2014*, 2014.

Artem Molchanov, Andreas Breitenmoser, and Gaurav S. Sukhatme. Active drifters: Towards a practical multi-robot system for ocean monitoring. In *IEEE International Conference on Robotics and Automation, ICRA 2015, Seattle, WA, USA, 26-30 May, 2015*, pages 545–552. IEEE, 2015. doi: 10.1109/ICRA.2015.7139232. URL https://doi.org/10.1109/ICRA.2015.7139232.

Artem Molchanov, Oliver Kroemer, Zhe Su, and Gaurav S. Sukhatme. Model-free contact localization for manipulated objects using biomimetic tactile sensors. In *Humanoids 2016 Workshop on Tactile Sensing for Manipulation, 2016*, 2016a.

Artem Molchanov, Oliver Kroemer, Zhe Su, and Gaurav S. Sukhatme. Contact localization on grasped objects using tactile sensing. In 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2016, Daejeon, South Korea, October 9-14, 2016, pages 216–222. IEEE, 2016b. doi: 10.1109/IROS.2016.7759058. URL https://doi.org/10.1109/IROS.2016.7759058.

Artem Molchanov, Tao Chen, Wolfgang Hönig, James A. Preiss, Nora Ayanian, and Gaurav S. Sukhatme. Sim-to-(multi)-real: Transfer of low-level robust control policies to multiple quadrotors. In 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2019, Macau, SAR, China, November 3-8, 2019, pages 59–66. IEEE, 2019a. doi: 10.1109/IROS40897.2019.8967695. URL https://doi.org/10.1109/IROS40897.2019.8967695.

Artem Molchanov, Tao Chen, Wolfgang Hönig, James A. Preiss, Nora Ayanian, and Gaurav S. Sukhatme. Sim-to-(multi)-real: Transfer of low-level robust control policies to multiple quadrotors. In *The Southern California Robotics Symposium, SCR 2019, Pasadena, CA, US, April, 2019*, 2019b.

Andrew Moore. Efficient memory-based learning for robot control. *PhD thesis, University of Cambridge*, 1990.

Igor Mordatch, Kendall Lowrey, and Emanuel Todorov. Ensemble-CIO: Full-body dynamic motion planning that transfers to physical humanoids. In *IEEE/RSJ Intl Conf. on Intelligent*

Robots and Systems (IROS), pages 5307–5314, 2015. doi: 10.1109/IROS.2015.7354126. URL https://doi.org/10.1109/IROS.2015.7354126.

Andrew Y. Ng and Stuart J. Russell. Algorithms for inverse reinforcement learning. In Pat Langley, editor, *Proceedings of the Seventeenth International Conference on Machine Learning (ICML 2000), Stanford University, Stanford, CA, USA, June 29 - July 2, 2000*, pages 663–670. Morgan Kaufmann, 2000.

Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.

Pierre-Yves Oudeyer and Frédéric Kaplan. What is intrinsic motivation? A typology of computational approaches. *Front. Neurorobot.*, 2009, 2009. URL https://doi.org/10.3389/neuro.12.006. 2007.

Michael Ouimet and Jorge Cortés. Coordinated rendezvous of underwater drifters in ocean internal waves. In *53rd IEEE Conference on Decision and Control, CDC 2014, Los Angeles, CA, USA, December 15-17, 2014*, pages 6099–6104. IEEE, 2014. doi: 10.1109/CDC.2014.7040344. URL https://doi.org/10.1109/CDC.2014.7040344.

Giambattista Parascandolo, Heikki Huttunen, and Tuomas Virtanen. Taming the waves: sine as activation function in deep neural networks. 2017.

Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *IEEE Intl Conf. on Robotics and Automation (ICRA)*, pages 1–8, 2018. doi: 10.1109/ICRA.2018.8460528. URL https://doi.org/10.1109/ICRA.2018.8460528.

Arvind Pereira, Jonathan Binney, Geoffrey A. Hollinger, and Gaurav S. Sukhatme. Risk-aware path planning for autonomous underwater vehicles using predictive ocean models. *J. Field Robotics*, 30(5):741–762, 2013. doi: 10.1002/rob.21472. URL https://doi.org/10.1002/rob.21472.

Anna Petrovskaya, Jaeheung Park, and Oussama Khatib. Probabilistic estimation of whole body contacts for multi-contact robot control. In *Proceedings 2007 IEEE International Conference on Robotics and Automation*, pages 568–573. IEEE, 2007.

Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neural Computation*, 3(1):88–97, 1991.

James A. Preiss, Wolfgang Hönig, Gaurav S. Sukhatme, and Nora Ayanian. Crazyswarm: A large nano-quadcopter swarm. In *IEEE Intl Conf. on Robotics and Automation (ICRA)*, pages 3299–3304. IEEE, 2017. doi: 10.1109/ICRA.2017.7989376. URL https://doi.org/10.1109/ICRA.2017.7989376.

Brooks L. Reed, Charles Ambler, Julio Guerrero, and Franz Hover. Vertical glider robots for subsea equipment delivery. In *IEEE International Conference on Robotics and Automation, ICRA 2011, Shanghai, China, 9-13 May 2011*, pages 2356–2361. IEEE, 2011. doi: 10.1109/ICRA. 2011.5980486. URL https://doi.org/10.1109/ICRA.2011.5980486.

Fereshteh Sadeghi and Sergey Levine. CAD2RL: real single-image flight without a single real image. In *Robotics: Science and Systems (RSS)*, 2017. doi: 10.15607/RSS.2017.XIII.034. URL http://www.roboticsproceedings.org/rss13/p34.html.

Sawyer. Rethink robotics. https://www.rethinkrobotics.com/sawyer/, 2012.

Juergen Schmidhuber. Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook. Institut für Informatik, Technische Universität München, 1987.

Jürgen Schmidhuber. Formal theory of creativity, fun, and intrinsic motivation (1990-2010). *Autonomous Mental Development*, 2(3):230–247, 2010. doi: 10.1109/TAMD.2010.2056368.

John Schulman, Nicolas Heess, Theophane Weber, and Pieter Abbeel. Gradient estimation using stochastic computation graphs. In *NeurIPS*, pages 3528–3536, 2015. URL http://dblp.uni-trier. de/db/conf/nips/nips2015.html#SchulmanHWA15.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL http://arxiv.org/abs/1707.06347.

A. Schwithal and C. Roman. Development of a new lagrangian float for studying coastal marine ecosystems. In *OCEANS 2009-EUROPE*, pages 1–6, May 2009. doi: 10.1109/OCEANSE.2009. 5278296.

L. Serrano, D. Kim, and R. B. Langley. A GPS Velocity Sensor: How Accurate Can It Be? – A First Look. In *Proc. of The Institute of Navigation. National Technical Meeting*, 2004, pages 875–885, 2004.

A. F. Shchepetkin and J. C. McWilliams. The regional oceanic modeling system (roms): a splitexplicit, free-surface, topography-following-coordinate oceanic model. *OCEAN MODELLING*, 9 (4):347–404, 2005. URL http://archipelago.uma.pt/pdf_library/Shchepetkin&McWilliams_2005_ OM.pdf.

David Silver, Aja Huang, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016. doi: 10.1038/nature16961. URL https://doi.org/10.1038/nature16961.

Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

Ryan N. Smith and Van T. Huynh. Controlling buoyancy-driven profiling floats for applications in ocean observation. *IEEE Journal of Oceanic Engineering*, PP(99):1 – 16, 2013. doi: http://dx.doi.org/10.1109/JOE.2013.2261895.

Z. Su, J.A. Fishel, T. Yamamoto, and G.E. Loeb. Use of tactile feedback to control exploratory movements to characterize object compliance. *Frontiers in neurorobotics*, 6, 2012.

Zhe Su, Karol Hausman, Yevgen Chebotar, Artem Molchanov, Gerald E Loeb, Gaurav S Sukhatme, and Stefan Schaal. Force estimation and slip detection/classification for grip control using a biomimetic tactile sensor. In *Humanoid Robots (Humanoids), 2015 IEEE-RAS 15th International Conference on*, pages 297–303. IEEE, 2015.

Zhe Su, Oliver Kroemer, Gerald E Loeb, Gaurav S. Sukhatme, and Stefan Schaal. Learning to switch between sensorimotorprimitives using multimodal haptic signals. In *International Conference on Simulation of Adaptive Behavior (SAB)*. Springer, 2016a.

Zhe Su, Stefan Schaal, and Gerald E Loeb. Surface tilt perception with a biomimetic tactile sensor. In 2016 6th IEEE RAS & EMBS international conference on biomedical robotics and biomechatronics (BioRob). IEEE, 2016b.

Flood Sung, Li Zhang, Tao Xiang, Timothy Hospedales, and Yongxin Yang. Learning to learn: Meta-critic networks for sample efficient learning. *arXiv preprint arXiv:1706.09529*, 2017.

Richard S. Sutton, David A. McAllester, Satinder P. Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In Sara A. Solla, Todd K. Leen, and Klaus-Robert Müller, editors, *Advances in Neural Information Processing Systems 12, [NIPS Conference, Denver, Colorado, USA, November 29 - December 4, 1999]*, pages 1057–1063. The MIT Press, 1999. URL http://papers.nips.cc/paper/ 1713-policy-gradient-methods-for-reinforcement-learning-with-function-approximation.

Jie Tan, Zhaoming Xie, Byron Boots, and C. Karen Liu. Simulation-based design of dynamic controllers for humanoid balancing. In *IEEE/RSJ Intl Conf. on Intelligent Robots and Systems (IROS)*, pages 2729–2736, 2016. doi: 10.1109/IROS.2016.7759424. URL https://doi.org/10. 1109/IROS.2016.7759424.

Jie Tan, Tingnan Zhang, Erwin Coumans, Atil Iscen, Yunfei Bai, Danijar Hafner, Steven Bohez, and Vincent Vanhoucke. Sim-to-real: Learning agile locomotion for quadruped robots. In *Robotics: Science and Systems (RSS)*, 2018. doi: 10.15607/RSS.2018.XIV.010. URL http://www.roboticsproceedings.org/rss14/p10.html.

Yuval Tassa, Nicolas Mansard, and Emo Todorov. Control-limited differential dynamic programming. In 2014 IEEE International Conference on Robotics and Automation, ICRA 2014, Hong Kong, China, May 31 - June 7, 2014, pages 1168–1175. IEEE, 2014. doi: 10.1109/ICRA.2014.6907001. URL https://doi.org/10.1109/ICRA.2014.6907001.

Sebastian Thrun and Lorien Pratt. Learning to learn. Springer Science & Business Media, 2012.

Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic Robotics (Intelligent Robotics and Autonomous Agents)*. The MIT Press, 2005. ISBN 0262201623.

Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *IEEE/RSJ Intl Conf. on Intelligent Robots and Systems (IROS)*, pages 23–30, 2017. doi: 10.1109/IROS.2017.8202133. URL https://doi.org/10.1109/IROS.2017.8202133.

Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2012, Vilamoura, Algarve, Portugal, October 7-12, 2012, pages 5026–5033. IEEE, 2012. doi: 10.1109/IROS.2012.6386109. URL https://doi.org/10.1109/IROS.2012.6386109.

Jonathan Tremblay, Thang To, Artem Molchanov, Stephen Tyree, Jan Kautz, and Stan Birchfield. Synthetically trained neural networks for learning human-readable plans from real-world demonstrations. In 2018 IEEE International Conference on Robotics and Automation, ICRA 2018, Brisbane, Australia, May 21-25, 2018. IEEE, 2018.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 6000–6010, 2017.

Nicholas Wettels, Veronica J. Santos, Roland S. Johansson, and Gerald E. Loeb. Biomimetic tactile sensor array. *Advanced Robotics*, 22(8):829–849, 2008. doi: 10.1163/156855308X314533. URL https://doi.org/10.1163/156855308X314533.

Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256, 1992.

Michael T. Wolf, Lars Blackmore, Yoshiaki Kuwata, Nanaz Fathpour, Alberto Elfes, and Claire Newman. Probabilistic motion planning of balloons in strong, uncertain wind fields. In *IEEE International Conference on Robotics and Automation, ICRA 2010, Anchorage, Alaska, USA, 3-7 May 2010*, pages 1123–1129. IEEE, 2010. doi: 10.1109/ROBOT.2010.5509135. URL https://doi.org/10.1109/ROBOT.2010.5509135.

Lijun Wu, Fei Tian, Yingce Xia, Yang Fan, Tao Qin, Jian-Huang Lai, and Tie-Yan Liu. Learning to teach with dynamic loss functions. In *NeurIPS*, pages 6467–6478, 2018. URL http://dblp.uni-trier.de/db/conf/nips/nips2018.html#WuTXFQLL18.

Zhongwen Xu, Hado van Hasselt, and David Silver. Meta-gradient reinforcement learning. In *NeurIPS*, pages 2402–2413, 2018.

Omry Yadan. A framework for elegantly configuring complex applications. Github, 2019. URL https://github.com/facebookresearch/hydra.

Tianhe Yu, Chelsea Finn, Annie Xie, Sudeep Dasari, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot imitation from observing humans via domain-adaptive meta-learning. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada,

April 30 - May 3, 2018, Workshop Track Proceedings. OpenReview.net, 2018. URL https://openreview.net/forum?id=Bk6p7YCLf.

Shaojun Zhu, Andrew Kimmel, Kostas E. Bekris, and Abdeslam Boularias. Fast model identification via physics engines for data-efficient policy search. In *Intl Joint Conf. on Ar-tificial Intelligence (IJCAI)*, pages 3249–3256, 2018. doi: 10.24963/ijcai.2018/451. URL https://doi.org/10.24963/ijcai.2018/451.

Haosheng Zou, Tongzheng Ren, Dong Yan, Hang Su, and Jun Zhu. Reward shaping via metalearning. *arXiv preprint arXiv:1901.09330*, 2019.