

Robots that Anticipate Pain: Anticipating Physical Perturbations from Visual Cues
through Deep Predictive Models

by

Indranil Sur

A Thesis Presented in Partial Fulfillment
of the Requirements for the Degree
Master of Science

Approved April 2017 by the
Graduate Supervisory Committee:

Heni Ben Amor, Chair
Georgios Fainekos
Yezhou Yang

ARIZONA STATE UNIVERSITY

May 2017

ABSTRACT

To ensure system integrity, robots need to proactively avoid any unwanted physical perturbation that may cause damage to the underlying hardware. In this thesis work, we investigate a machine learning approach that allows robots to anticipate impending physical perturbations from perceptual cues. In contrast to other approaches that require knowledge about sources of perturbation to be encoded before deployment, our method is based on experiential learning. Robots learn to associate visual cues with subsequent physical perturbations and contacts. In turn, these extracted visual cues are then used to predict potential future perturbations acting on the robot. To this end, we introduce a novel deep network architecture which combines multiple sub-networks for dealing with robot dynamics and perceptual input from the environment. We present a self-supervised approach for training the system that does not require any labeling of training data. Extensive experiments in a human-robot interaction task show that a robot can learn to predict physical contact by a human interaction partner without any prior information or labeling. Furthermore, the network is able to successfully predict physical contact from either depth stream input or traditional video input or using both modalities as input.

ACKNOWLEDGEMENTS

I would like to express my gratitude to my thesis advisor, Dr. Heni Ben Amor for his wisdom and friendship, for providing me this incredible research opportunity, for the countless conversations and whiteboard discussions, and for financial support through Teaching Assistanship. I would also like to thank Dr. Georgios Fainekos for the opportunity to work in his ‘Cyber Physical Systems Lab’ and showing me the ways of conducting proper research.

I would like to thank David Vogt at TU Bergakademie Freiberg University for exciting and thought provoking early discussions on my research and other topics. I would also like to thank my labmates Joseph Campbell, Nikhil Kalige, Kevin Sebastian Luck, Trevor Barron at ‘Interactive Robotics Lab’ and Erkan Tuncali, Kangjin Kim, Bardh Hoxha at ‘Cyber Physical Systems Lab’ for wonderful discussions and timely help. I would also like to thank Mark Strickland, Simon Stepputtis, Ashish Kumar, Trevor Richardson, Geoffrey Clark, Ramsundar Kalpagam Ganesan, Dax McDonald, David Kish for taking time out and taking part in the experiments.

Finally, I owe gratitude to my family and close friends for their continued love, support and patience.

TABLE OF CONTENTS

| | Page |
|---|------|
| LIST OF TABLES | iv |
| LIST OF FIGURES | v |
| CHAPTER | |
| 1 INTRODUCTION | 1 |
| 2 RELATED WORK | 5 |
| 3 DEEP PREDICTIVE MODELS OF PHYSICAL PERTURBATIONS ... | 7 |
| 3.1 Deep Dynamics Model | 8 |
| 3.2 Perceptual Anticipation Model | 13 |
| 4 EXPERIMENTAL RESULTS | 18 |
| 4.1 Simulation Experiment | 19 |
| 4.1.1 Simulation Collision Example | 20 |
| 4.1.2 Prediction Error on Simulation Training and Test Sets | 20 |
| 4.1.3 Test Error for each predictor | 21 |
| 4.1.4 Generalization Test | 22 |
| 4.2 Human-Robot Interaction | 24 |
| 4.2.1 Example Interaction | 28 |
| 4.2.2 Visualizing Network Saliency | 29 |
| 4.2.3 Prediction Error on Training and Test Sets | 30 |
| 4.2.4 Test Error after Model Selection | 33 |
| 5 FUTURE WORK | 36 |
| 6 CONCLUSIONS | 37 |
| REFERENCES | 38 |
| APPENDIX A | 41 |

LIST OF TABLES

| Table | Page |
|--|------|
| 4.1 Comparison of Different Configurations of DPM on Simulation Train Data | 20 |
| 4.2 Comparison of Different Configurations of DPM on Simulation Test Data | 21 |
| 4.3 Test Error of Different Predictors on Depth & Frame Channel | 22 |
| 4.4 Generalization Test Error for Model Trained with Depth Data | 25 |
| 4.5 Generalization Test Error for Model Trained with Frame Data | 26 |
| 4.6 Generalization Test Error for Model Trained with Depth & Frame Data | 26 |
| 4.7 Comparison of Different Configurations of DPM on Physical Train Data | 30 |
| 4.8 Comparison of Different Configurations of DPM on Physical Test Data | 31 |
| 4.9 Test Errors after Model Selection with Depth Data | 32 |
| 4.10 Test Errors after Model Selection with Frame Data | 33 |

LIST OF FIGURES

| Figure | Page |
|--|------|
| 1.1 Baxter Robot Anticipates Physical Contact by Human..... | 2 |
| 1.2 Overview Architecture of the Deep Predictive Model | 3 |
| 3.1 Visualization of the Nature of Different Motor Babbling Policies..... | 10 |
| 3.2 End-Effector Coordinates for Different Motor Babbling Policies..... | 11 |
| 3.3 Deep Dynamics Neural Network | 12 |
| 3.4 Analysis of Predictive Error Functions..... | 13 |
| 3.5 Perceptual Anticipation Model..... | 15 |
| 4.1 Simulation Scenario to Test Perceptual Anticipation Model | 18 |
| 4.2 Example Collision Visualization..... | 19 |
| 4.3 ROC Plot for Test Data on Depth & Frame Channel | 21 |
| 4.4 Simulation Scenario to Test Generalization capability..... | 23 |
| 4.5 Spike Visualization for Different Test Scenarios..... | 24 |
| 4.6 Example Collision Visualization for Generalization Test..... | 25 |
| 4.7 Example Interaction Visualization..... | 27 |
| 4.8 Saliency Map Visualization | 29 |
| 4.9 Spike Visualization of Ground Truth and Prediction of Perturbation ... | 31 |
| 4.10 ROC Plots for Test Data after Model Selection..... | 32 |
| 4.11 Quantitative Analysis on Different Test Participants..... | 34 |

Chapter 1

INTRODUCTION

According to Isaac Asimov’s third law of robotics [2], “a robot must protect its own existence as long as such protection does not conflict with the First or Second Law”, i.e., as long as it does not harm a human. Aspects of safety and self-preservation are tightly coupled to autonomy and longevity of robotic systems. For robots to explore their environment and engage in physical contact with objects and humans, they need to ensure that any such interaction may not lead to tear, damage, or irreparable harm to the underlying hardware. Situations that jeopardize the integrity of the system need to be detected and actively avoided. This determination can be performed in either a reactive way, e.g., by using sensors to identify forces acting on the robot, or in a pro-active way, e.g., by detecting impending collisions. In recent years, a plethora of safety methods have been proposed that are based on reactive strategy. Approaches for compliant control and, in particular, impedance control techniques have been shown to enable safe human-robot interaction [13] in a variety of close-contact tasks. Such methods are typically referred to as *post-contact* approaches, since they react to forces exchanged between the system and its environment after they first occur.

In many application domains, however, robots need to pro-actively reason about impending damage before it occurs. Methods that tackle these scenarios, have mostly focused on proximity detection and collision avoidance. Motion tracking or other sensing devices are used to identify nearby humans and obstacles in order determine whether they intersect the robot’s path. To this end, a human expert has to reason about the expected obstacles before deployment and, in turn, hand-code methods for object recognition, human tracking, or collision detection. Such methods are

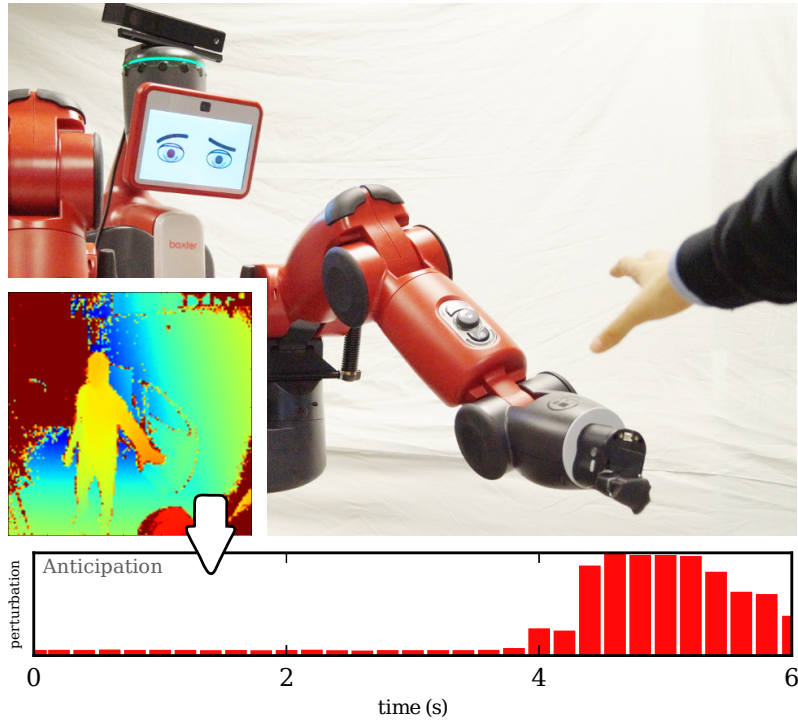


Figure 1.1: Baxter robot anticipates physical contact by human.

largely based on the status quo of the environment and do not take in account the expected future states. For example, the behavior of a human interaction partner may already provide cues whether or not physical contact is to be expected. In addition, such methods suffer from limited adaptivity – the robotic system is not able to incorporate other or new sources of physical perturbations that have not been considered by the human expert. Changes in the application domain typically require the expert to specify the set of obstacles/perturbants and how they can be identified from sensor data. However, for robots to autonomously explore new goals, tasks, and environments they cannot be constantly relying on human intervention and involvement. In order to increase resilience of robotic systems, the processes responsible for ensuring safety should inherently be (a) adaptive to changes that occur during the cycle of operation and (b) anticipative in nature, so as to proactively avoid

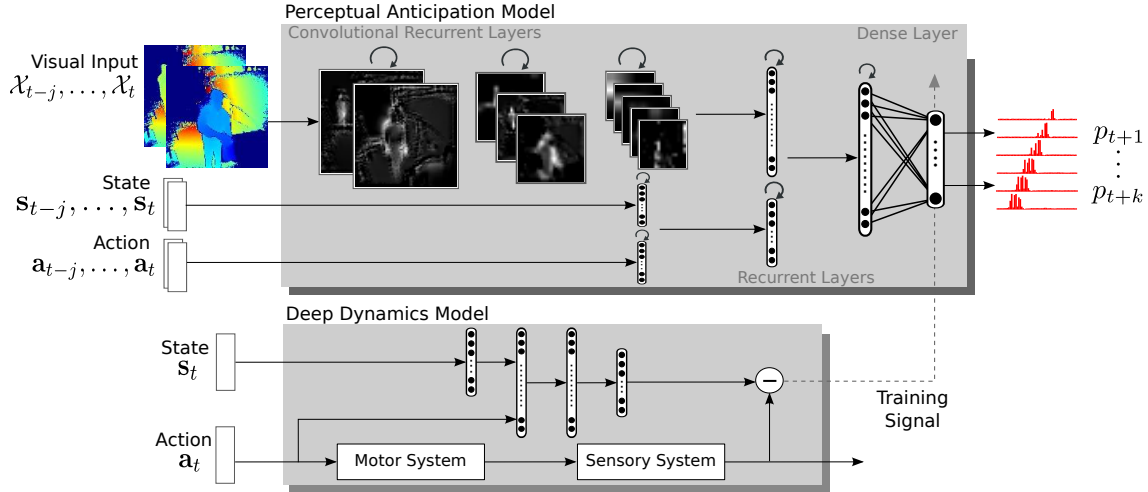


Figure 1.2: Overview of architecture of the Deep Predictive Model for perturbation prediction from visual cues.

damage.

In biology such processes are common place: humans and animals experience *pain* as the guiding signal which ensures self-preservation and safe, “open ended” exploration of the world. Over time, biological creatures *learn to anticipate* impending pain sensations from environmental and proprioceptive cues, e.g., from the velocity and direction of an obstacle, or an imbalance in posture. The relationship between pain and learning is bi-directional in humans and vertebrate animals. Our perception of pain shapes the way we learn, but the perception itself is also shaped by the experiences we undergo during learning. In particular, repeated exposure to a sensory cue (e.g. a specific image or object) preceding a negative outcome (e.g. electric shock) will turn the cue into a conditioned stimulus that helps anticipate the shock in future trials. This ability to associate a certain stimuli with negative future sensations is considered to be essential to survival.

Inspired by the relationship between pain and learning in biological creatures, we propose in this paper a new method for learning to anticipate physical perturbations

in robotic systems. Robot safety is realized through an adaptive process in which prior experiences are used to predict impending harm to the system. These predictions are based on (1) environmental cues, e.g., camera input, (2) proprioceptive information, and (3) the intended actions of the robot. We introduce a *deep predictive model* (DPM) that leverages all three sources of information in order to anticipate contact and physical perturbations in a window of future time steps. An important characteristic of this deep predictive model, is that it differentiates between forces and perturbations that are the result of an executed robot behavior, e.g., fast arm movements, and external perturbations that are the result of outside influence. After learning, the DPM can be used to anticipate future states that result in forces being applied to the robot. In turn, this information can be used to provide visual feedback to the user, or actively avoid the predicted states, or change the impedance of the system.

We will evaluate the introduced system in a set of experiments in which a robot has to anticipate physical perturbations caused by a human interaction partner. We will show that the introduced model effectively anticipates contact using either RGB or RGB-D camera sensors. While the introduced method can be used to actively avoid noxious states, we will focus our analysis in this paper to the detection of perturbations only.

Chapter 2

RELATED WORK

Safety plays a critical role in the field of robotics. Recently, various methods have been put forward in order to protect a human interaction partner from harm. The work in [8], for example, uses proprioceptive sensing to identify collisions and, in turn, execute a reactive control strategy that enables the robot to retract from the location of impact. In a similar vein, the work in [14] describes methods for rapid collision detection and response through trajectory scaling. Many approaches to safe human-robot interaction are based on rapid sensing through force-torque sensors or tactile sensors [17]. Another approach is to generate estimates of external forces acting on a robot using disturbance observers [11]. These approaches, however, require a model of the underlying plant, which in the case of complex humanoid robots can become challenging to derive. In addition, nonlinearities underlying the current robot or task can often lead to instabilities in the system [21]. To circumvent such challenges, several approaches have been proposed for learning perturbation filters using a data-driven machine learning method [3, 4, 5]. An alternative, bio-inspired approach was proposed in [19]. In particular, an “Artificial Robot Nervous System” was introduced which emulates the human skin structure, as well as the spikes generated whenever an impact occurs.

All of the above techniques are reactive in nature. Only after contact with its environment, can a robot detect a physical impact or perturbation and react to it. However, many critical situations require a proactive avoidance strategy. To this end, various methods for human motion anticipation have been proposed [20, 23, 1]. Given a partially observed human motion, a robot can anticipate the goal location and inter-

mediate path and accordingly generate a collision free navigation strategy. However, such approaches are brittle in that they require some form of human tracking. If the source of the perturbation is non-human, e.g., a moving object then no anticipation can occur. In this paper, we are interested in dynamic approaches to anticipation of perturbations. Robots learn to predict contacts or impact by associating them with visual cues. This is similar in spirit to [6] where dashcam videos were used to predict car collisions. However, an important difference is that our predictions are based on both visual cues, proprioceptive sensors, as well as next robot actions.

DEEP PREDICTIVE MODELS OF PHYSICAL PERTURBATIONS

Our goal is to enable an intelligent system to anticipate undesired external perturbations before they occur. Following the biological inspiration, repeated exposure to a sensory cue preceding a physical perturbation will turn the cue into a predictive variable that helps anticipate the perturbation. To this end, we propose the deep predictive model seen in Fig. 1.2. A key insight of our approach, is that a robot can learn to correlate perceptual data, e.g., camera image, to its future internal states of the robot. By establishing a mapping between perceptual information and internal states, the it can anticipate harm before it occurs.

The DPM is based on a modular architecture, in which two modules generate the necessary relationship in order to map perceptual data to future anticipated noxious signals. The first module of the DPM is a *deep dynamics module* that learns to discriminate between external perturbations caused by outside events and natural variations of sensor readings due to the currently executed behavior and sensor noise. The deep dynamics module is trained to generate a signal, whenever the recorded sensor values cannot be explained by the actions of the robot. This produces a dense training signal for self-supervised learning of external perturbations. Prediction in the deep dynamics model is performed within a probabilistic framework in order to estimate the model uncertainties.

The second module of the DPM is the *perceptual anticipation module*– a deep network that takes visual percepts, proprioceptive states, and actions taken as inputs and generates predictions over future noxious signals. It learns to associate specific visual and proprioceptive cues to harmful states. The perceptual anticipation module

contains convolutional recurrent layers [24], [10] that process the visual input in both space and time. In addition, it features recurrent and dense layers, that combine the processed visual information with information about the robot state and actions to produce multiple predictions over expected perturbations. Learning is performed using the paradigm of self-supervised learning. More specifically, the training signal produced by the deep dynamics module is used as a target signal. No human intervention is needed in order to either provide new training data or label existing data.

Subsequently, we will explain the elements of the introduced network architecture and describe the underlying training and inference process in more detail.

3.1 Deep Dynamics Model

The task of the deep dynamics model is to identify exogenous perturbations affecting the robot. Detecting such perturbations in the sensor stream can be challenging when the robot is performing dynamic movements that by themselves cause fluctuation in the sensor readings. To discriminate between exogenous and endogenous perturbations, we will use a strategy inspired by the human motor system.

Before sending an action $\mathbf{a}_t \in \mathbb{R}^Q$ to the actuators, the robot creates a copy of \mathbf{a}_t , the so-called efference copy, and predicts the expected sensory stimuli after execution. This prediction is performed by the deep dynamics model, a neural network that maps the current state \mathbf{s}_t and intended action \mathbf{a}_t onto the expected next state \mathbf{s}_{t+1}^* .

After executing action \mathbf{a}_t , we can measure the discrepancy between the expected sensations and the measured sensor values, also called the refference. The degree of discrepancy is an indicator for external physical perturbations.

This methodology is related to the concept of disturbance observers [11]. However, in contrast to disturbance observers, no explicit model of the system needs to be

provided. Instead, the deep dynamics model is entirely learned, which is particularly important for compliant and cable-driven robots, for which analytical models can often be hard to devise and difficult to calibrate.

Next, we will describe how to collect data and perform learning in a deep dynamics models. Then, we will show how the learned model can be used for probabilistic inference and perturbation estimation.

Data Collection: Without loss of generality, we define for the remainder of the paper the system state to be $\mathbf{s}_t = [\boldsymbol{\theta}_t, \dot{\boldsymbol{\theta}}_t, \ddot{\boldsymbol{\theta}}_t, \mathbf{p}_t]^T \in \mathbb{R}^P$ which includes joint angles $\boldsymbol{\theta}_t$, joint velocities $\dot{\boldsymbol{\theta}}_t$, accelerations $\ddot{\boldsymbol{\theta}}_t$ and end-effector pose \mathbf{p}_t .

To collect training data for training the deep dynamics model, we use motor babbling [25, 9]. To this end, small changes are applied as the control action $\mathbf{a}_t \in \mathbb{R}^Q$, where Q is the number of degrees of freedom for the robot, i.e., the number of joints. Considering a naive implementation of motor babbling, the action can be sampled from an isotropic Gaussian distribution $\mathbf{a}_t \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. However, empirically we have observed that such a naive sampling approach does not effectively cover the state-action space by unnecessarily focusing on samples in the center of the distribution.

To alleviate this problem, we use a bi-modal distribution and incorporate a simple momentum term

$$\begin{aligned} \pi &\sim \mathcal{B}(0.5) \\ \mathbf{u} &\sim \pi \mathcal{N}(\mu \mathbf{1}, \sigma^2 \mathbf{I}) + (1 - \pi) \mathcal{N}(-\mu \mathbf{1}, \sigma^2 \mathbf{I}) \\ \mathbf{a}_t &= (\mathbf{u} + \mathbf{a}_{t-1})/2 \end{aligned} \tag{3.1}$$

Actions sampled using the above strategy effectively cover the state-action space and generate trajectories without causing wear and tear. The result of the motor babbling phase, is a dataset for training which consists of N triplets $(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1})$ consisting of the current state, current action, and the next state. The individual matrices storing all states and actions are denoted by $\mathbf{S}_t, \mathbf{S}_{t+1} \in \mathbb{R}^{N \times P}$ and $\mathbf{A}_t \in$

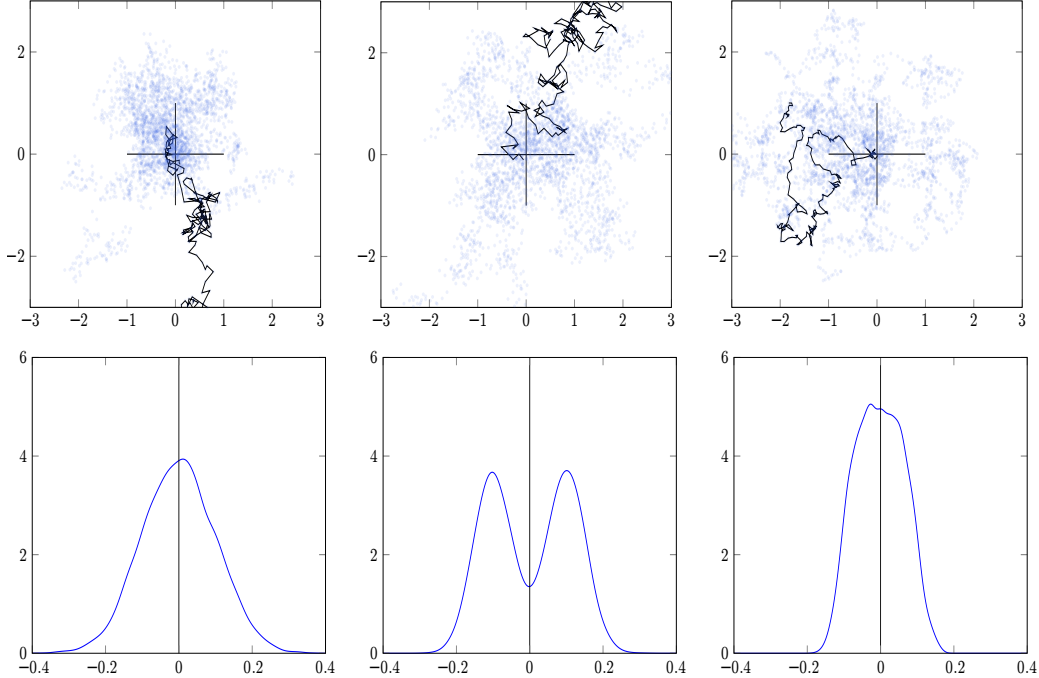


Figure 3.1: Visualization of 2-D state-space coverage and the action distributions underlying the motor babbling process. Left: Gaussian distribution resulting in Brownian motion in action space. Middle: Bi-modal distribution. Right: Incorporating momentum yields smoother motions and better coverage of the state-action space.

$\mathbb{R}^{N \times Q}$.

The effects of using momentum are demonstrated in Fig. 3.1 (right) and Fig. 3.2 (right).

Model Learning: The deep dynamics model is an artificial neural network that maps a state \mathbf{s}_t and action \mathbf{a}_t on to the expected next state \mathbf{s}_{t+1}^* . Training is performed using Dropout [15] and data collected in the motor babbling process. Euclidean loss function, $\mathcal{L} = \frac{1}{2M} \sum_{i=1}^M \|\mathbf{s}_i - \hat{\mathbf{s}}_i\|$, where M is the mini-batch size, is used to identify the error. To identify the optimal network parameters, we used a grid search over the network architectures and hyper-parameters yielding the network shown in Fig. 3.3.

The neural network generates a point estimate for any set of inputs. However,

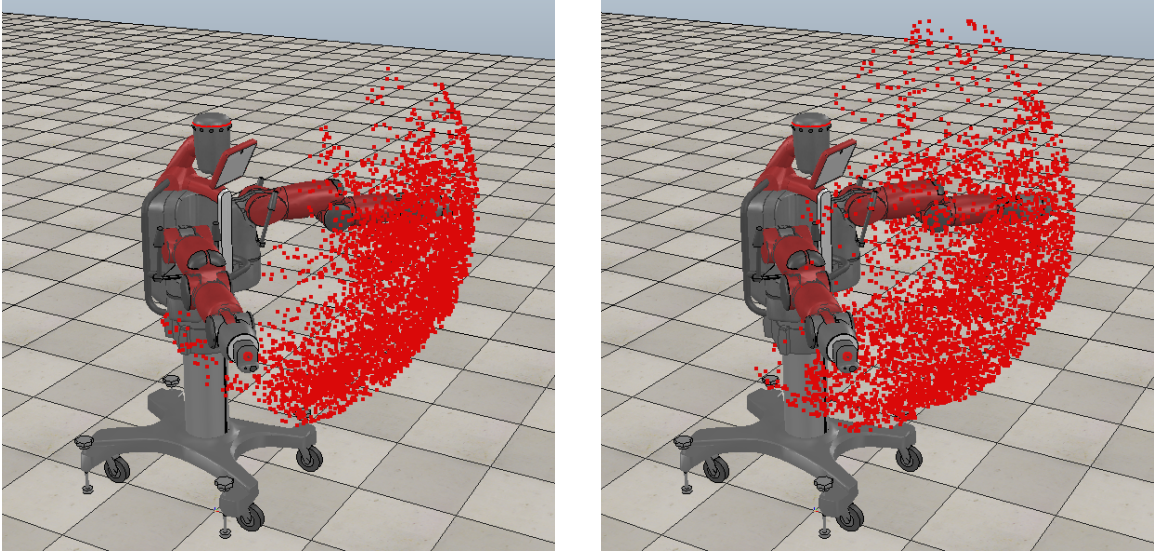


Figure 3.2: By visualizing the end-effector coordinates we can also get a perception of state-action space coverage. On left with gaussian sampled action we can see how clustered the end-effector coordinates are while using bimodal distribution and using momentum we can see the state-action space is fairly sampled.

due to the non-determinism and noise underlying such tasks, it is important to reason about uncertainties when making predictions. To this end, we leverage recent theoretical insights in order to generate probabilistic predictions from a trained neural network. In particular, it was shown in [12] that neural network learning using the Dropout method [15] is equivalent to a Bayesian approximation of a Gaussian Process modeling the training data. Following this insight, we generate a set of predictions $\{\hat{\mathbf{s}}_1, \dots, \hat{\mathbf{s}}_T\}$ from a trained network through T stochastic forward passes. The generated predictions form a possibly complex distribution represented as a population of solutions. We then extract an approximate parametric form of the underlying

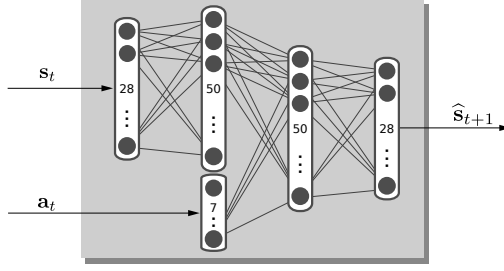


Figure 3.3: Deep Dynamics Neural Network

distribution through moment matching

$$\begin{aligned} \mathbb{E}(\mathbf{s}_{t+1}^*) &\approx \frac{1}{T} \sum_{i=1}^T \hat{\mathbf{s}}_i \\ \text{Var}(\mathbf{s}_{t+1}^*) &\approx \frac{1}{T} \sum_{i=1}^T \hat{\mathbf{s}}_i^T \hat{\mathbf{s}}_i - \mathbb{E}(\mathbf{s}_{t+1}^*)^T \mathbb{E}(\mathbf{s}_{t+1}^*) \end{aligned} \quad (3.2)$$

Given the above distribution moments we can reason about the uncertainty underlying the prediction process. We will see in the next section, that this information can be incorporated into any distance measure used for identifying incongruent sensations.

Perturbation as Predictive Error: As described above and depicted in Fig. 1.2, we generate in every time step a prediction of the next sensory state \mathbf{s}_{t+1}^* based on the current state and action. Consequently, after executing the action \mathbf{a}_t , we measure the *true* sensor values of the robot and compare them to the prediction, i.e.,

$$\boldsymbol{\delta} = \mathbb{E}(\mathbf{s}_{t+1}^*) - \mathbf{s}_{t+1}^* \quad (3.3)$$

Taking the norm of vector $\boldsymbol{\delta}$ we get an estimate of the discrepancy between robot expectation and reality. Assuming a reasonably accurate model, this discrepancy is an estimate of exogenous perturbations that affected the system dynamics. To correct for the inherent model uncertainty and the probabilistic nature of the predictions, we

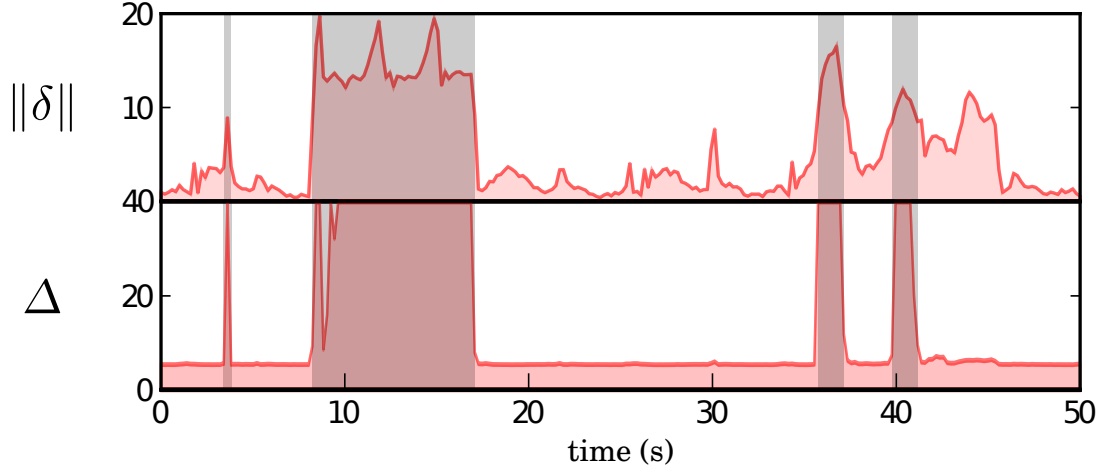


Figure 3.4: Analysis of predictive error functions

can take an exponential at the confidence bound according to Equ. ??.

$$\begin{aligned} \tau_i &= \exp(\delta_i - 2 \text{Var}(\mathbf{s}_{t+1}^*)_i), \\ \Delta &= \|\boldsymbol{\tau}\| \end{aligned} \tag{3.4}$$

This approach effectively incorporates the confidence bounds of the prediction into perturbation estimation and the exponential scaling simplifies the process of discerning nominal behavior from abnormal behavior due to external physical influence.

Fig. 3.4 depicts the effect the exponential scaling at confidence bound. The ground-truth highlighted in gray was hand-labeled using a video stream as reference. We can see both predictive error functions produce elevated responses during moments of contact. Yet, by incorporating an exponential scaling as performed in Δ we can reduce spurious activations and false-positives.

3.2 Perceptual Anticipation Model

In this section, we will describe how experienced perturbations can be correlated to predictive visual cues. In turn, these visual cues can later be used to anticipate the occurrence of physical contact. For example, perceiving an approaching wall may

indicate an impending collision. Still, whether or not an external perturbation will occur is also dependent on the next actions of the robot, e.g., whether or not the robot will stop its course. Hence, states and actions need to be included in the prediction process.

In our approach, a dedicated model – the Perceptual Anticipation Model – learns to predict impending perturbations from a sequence of visual input data, robot actions, and sensory states. Visual input representing the environment is given by an $R \times C$ grid. Each cell of the grid stores F measurements, i.e., depth or color channels that may vary over time. Thus, the visual input corresponds to a tensor $\mathbf{X} \in \mathbb{R}^{R \times C \times F}$.

Training is based on repeated physical interaction with the environment, e.g., a human or static objects. Throughout this process, a stream of visuals is recorded using either a traditional RGB camera or an RGB-D depth camera. The result is a sequence of observations \mathbf{X}_t . As the same time, states \mathbf{s}_t and actions \mathbf{a}_t of the robot are recorded. In addition, at every time step, the previously introduced deep dynamics model is run, in order to generate estimates p_t of perturbations currently acting on the robot.

Given the above data sets, the goal of training the Perceptual Anticipation Model is to approximate the distribution

$$\begin{aligned}
 P(p_{t+1}, \dots, p_{t+k} \mid & \mathbf{X}_{t-j}, \dots, \mathbf{X}_t, \\
 & \mathbf{s}_{t-j}, \dots, \mathbf{s}_t, \\
 & \mathbf{a}_{t-j}, \dots, \mathbf{a}_t)
 \end{aligned}
 \tag{3.5}$$

where j defines a window of past time steps. More specifically, we are interested in a predictive model that *generates expected future perturbations* p_{t+1}, \dots, p_{t+k} conditioned on *past inputs* \mathbf{X}_{t-j} , as well as current and previous robot states and actions.

The presented task requires both attention to spatial patterns within the visual

input, as well as attention to temporal patterns and behavioral dynamics. Hence, special care has to be taken to ensure that the model architecture used for learning can effectively identify and incorporate both sources of information when making a prediction.

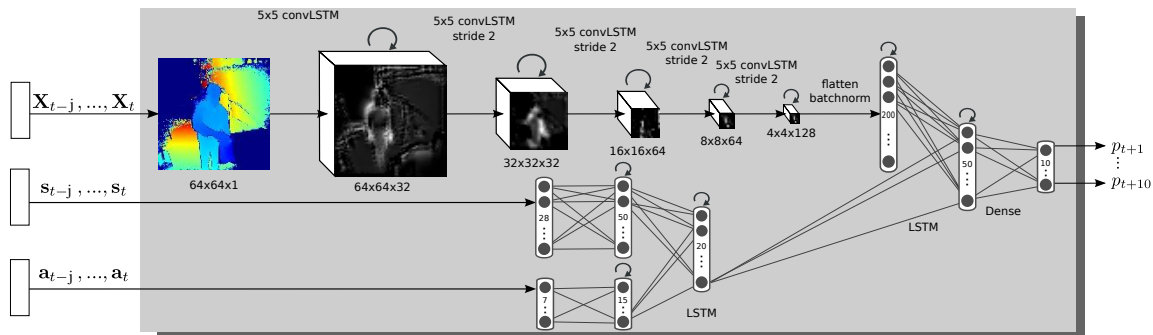


Figure 3.5: Perceptual Anticipation Model

Network Architecture: In order to base all predictions on both spatial and temporal information, we employ a *deep convolutional recurrent network* to learn the anticipation model. The input of the network is a sequence of visual data, the robot states, and actions. The output of the network is a vector $\mathbf{p}_t = [p_{t+1}, \dots, p_{t+k}]^T \in \mathbb{R}^k$ that describes the likelihood of a perturbation at each of the time steps $\{t+1, \dots, t+k\}$. In order to incorporate temporal dynamics, recurrent units are used according to either the Long Short Term Memory (LSTM) [16] model or the Gated Recurrent Units (GRU) [7] model. These recurrent units keep track of a hidden state. In turn, the next state of the network is calculated as a function of the hidden state and the new inputs. In the case of using convLSTM units [27], the network output is governed by

the following set of equations

$$\begin{aligned}
\mathbf{I}_t &= \sigma(\mathbf{W}_{xi} * \mathbf{X}_t + \mathbf{W}_{hi} * \mathbf{H}_{t-1} + \mathbf{b}_i) \\
\mathbf{J}_t &= \tanh(\mathbf{W}_{xc} * \mathbf{X}_t + \mathbf{W}_{hc} * \mathbf{H}_{t-1} + \mathbf{b}_c) \\
\mathbf{F}_t &= \sigma(\mathbf{W}_{xf} * \mathbf{X}_t + \mathbf{W}_{hf} * \mathbf{H}_{t-1} + \mathbf{b}_f) \\
\mathbf{O}_t &= \sigma(\mathbf{W}_{xo} * \mathbf{X}_t + \mathbf{W}_{ho} * \mathbf{H}_{t-1} + \mathbf{b}_o) \\
\mathbf{C}_t &= \mathbf{F}_t \odot \mathbf{C}_{t-1} + \mathbf{I}_t \odot \mathbf{J}_t \\
\mathbf{H}_t &= \mathbf{O}_t \odot \tanh(\mathbf{C}_t)
\end{aligned} \tag{3.6}$$

where \mathbf{X}_t is the input at time t , \mathbf{C}_t is the memory cell output, and \mathbf{H}_t is the hidden state. Using memory cells is a critical element of LSTM and ensures that the network output is conditioned on previous activations of the network. The gate variables \mathbf{I}_t , \mathbf{F}_t , \mathbf{O}_t of the convLSTM model denote 3D tensors with 2 dimensions as spatial dimensions and one dimension for the convolution filter dimension. σ denotes the sigmoid function, $*$ is convolution operation and \odot is Hadamard multiplication. In the case of using convGRU units [18], the neural network outputs are defined by equations

$$\begin{aligned}
\mathbf{Z}_t &= \sigma(\mathbf{W}_{xz} * \mathbf{X}_t + \mathbf{W}_{hz} * \mathbf{H}_{t-1} + \mathbf{b}_z) \\
\mathbf{R}_t &= \sigma(\mathbf{W}_{xr} * \mathbf{X}_t + \mathbf{W}_{hr} * \mathbf{H}_{t-1} + \mathbf{b}_r) \\
\widehat{\mathbf{H}}_t &= \Phi(\mathbf{W}_x * \mathbf{X}_t + \mathbf{W}_h * (\mathbf{R}_t \odot \mathbf{H}_{t-1}) + \mathbf{b}) \\
\mathbf{H}_t &= (1 - \mathbf{Z}_t) \odot \mathbf{H}_{t-1} + \mathbf{Z}_t \odot \widehat{\mathbf{H}}_t
\end{aligned} \tag{3.7}$$

with inputs $\mathbf{X}_1, \dots, \mathbf{X}_t$, cell outputs/hidden states $\mathbf{H}_1, \dots, \mathbf{H}_t$ and gates \mathbf{Z}_t , \mathbf{R}_t of convGRU are 3D tensors with 2 dimensions as the spatial dimensions and one dimension for the convolution filter dimension. The symbol Φ describes the activation function used, e.g., tanh or Rectified Linear Unit (ReLU) [22]. Note that the equations governing the convGRU have fewer parameters, which reduced both training and execution time. Faster forward passes can be crucial for real-time robotics application

as the one described here.

Loss Function: In order to train all network parameters, we use a weighted binary cross-entropy as a loss function

$$\mathcal{L} = -\frac{1}{T} \sum_{t=1}^T \left(w \mathbf{p}_t \log(\hat{\mathbf{p}}_t) + (1 - \mathbf{p}_t) (1 - \log(\hat{\mathbf{p}}_t)) \right) \quad (3.8)$$

where w is weight penalty, \mathbf{p}_t is the ground truth, and $\hat{\mathbf{p}}_t$ is the likelihood of perturbation generated by the network.

EXPERIMENTAL RESULTS

As simulation result we have investigated the Perceptual Anticipation Model in a dynamic environment interacting with the robot. As real world experiment we conducted a human-robot interaction study to investigate the validity of the introduced approach. In all of the following experiments the prediction rate was 5Hz. System state dimension P is 28 and degrees of freedom Q is 7. Mean and standard deviation for action sampling is $\mu = 0.1, \sigma = 0.05$. Dimension of environment data: $R, C = 64$ and $F = 1$ (gray-scale frame channel or depth channel). The parameters of DPM are set to: $k = 10, w = 3, j = 9$.

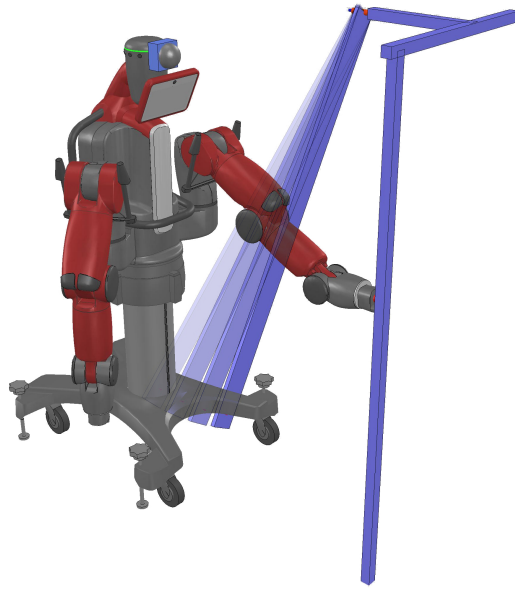
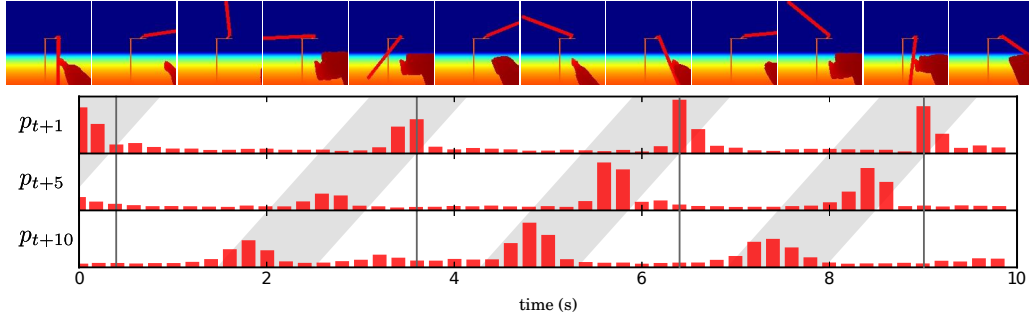
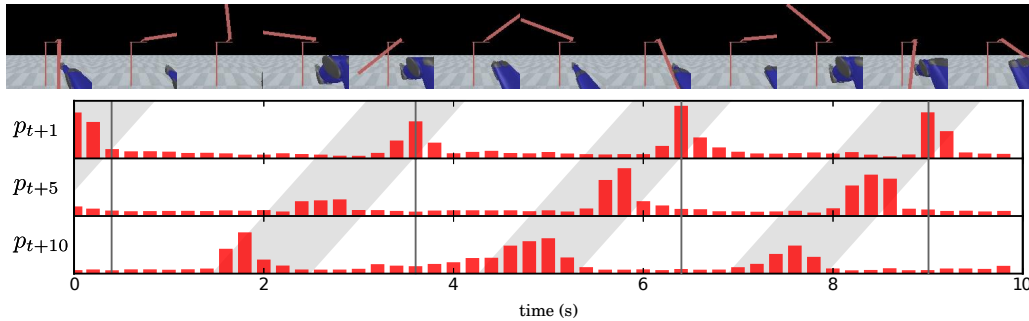


Figure 4.1: Dynamic scenario in V-REP simulation to test the anticipation capability of the Perceptual Anticipation Model



(a) Predictions of the likelihood of physical collision generated from a perceptual anticipation model which was trained on depth input. The predictors fire at different moments ahead of the actual moment of contact.



(b) Predictions of the likelihood of physical collision generated from a perceptual anticipation model which was trained on grayscale video input. The predictors fire at different moments ahead of the actual moment of contact.

Figure 4.2: Example Collision Visualization

4.1 Simulation Experiment

To demonstrate the anticipation capability of the Perceptual Anticipation Model, we have created a dynamic scenario in V-REP simulation as shown in Fig. 4.1. Here a rod is rotating and also the robot is performing the constant policy. The Perceptual Anticipation Model has to anticipate the collision between the rotating rod and the robot arm performing its own actions. The scenario is simple but its non-linear in nature as the rod is rotating and is non-deterministic in nature as the rod is not

moving in constant angular velocity.

4.1.1 Simulation Collision Example

A collision scenario and the anticipation capability of the Perceptual Anticipation Model is demonstrated in Fig. 4.10. Here we can see that the predictions for p_{t+1} , p_{t+5} and p_{t+10} are predicting the likelihood of collision 1,5,10 timesteps ahead correspondingly. We also can see that the models are equally adept in extracting visual cues from both depth and frame data.

4.1.2 Prediction Error on Simulation Training and Test Sets

The prediction error on the training set can be seen in Tab. 4.1 and test set can be seen in Tab. 4.2. For training different data configurations were used and tested, namely depth data vs. grayscale frame data and both depth and frame data. The predictions are thresholded to get optimal MCC (Matthews correlation coefficient) in each of the configuration.

| Config | true+ | pred+ | precision | recall | MCC |
|----------------|-------|-------|-----------|--------|--------|
| ConvLSTM-SA-D | 14693 | 30887 | 0.4757 | 0.7689 | 0.5675 |
| ConvLSTM-SA-F | 14093 | 28554 | 0.4936 | 0.7375 | 0.5669 |
| ConvLSTM-SA-DF | 13346 | 26081 | 0.5117 | 0.6984 | 0.5622 |

Table 4.1: Comparison between different configurations of the DPM on Train dataset (19110 perturbation points in 27000)

This shows when using both modalities of depth and frame, the model has better capability to generalize and hence has better MMC on test dataset.

| Config | true+ | pred+ | precision | recall | MCC |
|----------------|-------|-------|-----------|--------|--------|
| ConvLSTM-SA-DF | 3114 | 7146 | 0.4358 | 0.7670 | 0.5389 |
| ConvLSTM-SA-D | 3003 | 6703 | 0.4480 | 0.7397 | 0.5370 |
| ConvLSTM-SA-F | 3127 | 7530 | 0.4153 | 0.7702 | 0.5243 |

Table 4.2: Comparison between different configurations of the DPM on Test dataset (4060 perturbation points in 60000)

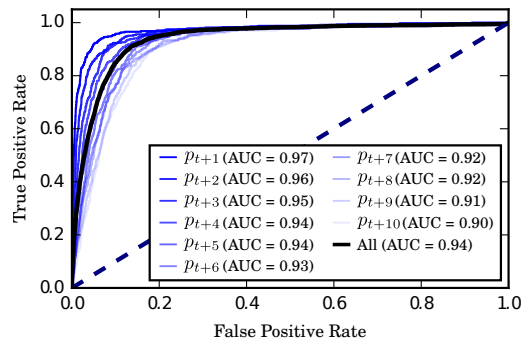


Figure 4.3: ROC(Receiver Operating Characteristic) plot for test data with convLSTM-SA-DF

4.1.3 Test Error for each predictor

Given the above results, we selected a state-action conditioned ConvLSTM with both depth and frame data(ConvLSTM-SA-DF) as the underlying network model for the prediction. Consequently, we analyzed the prediction errors for the generated predictions p_{t+1}, \dots, p_{t+k} separately. Tab. 4.3 show the results of this analysis. We can see precision and recall deteriorate, as the network produces predictions for longer horizons.

Similarly in Fig. 4.3 we can see the ROC(Receiver Operating Characteristic) plot for test data with our selected configuration. Decreasing AUC(area under curve) indicating how well the different predictors are performing. We can see that as we go

| pred no | true+ | pred+ | precision | recall | MCC |
|----------------|-------|-------|-----------|--------|--------|
| p_{t+1} | 335 | 453 | 0.7395 | 0.8251 | 0.7644 |
| p_{t+2} | 293 | 419 | 0.6993 | 0.7217 | 0.6890 |
| p_{t+3} | 287 | 483 | 0.5942 | 0.7069 | 0.6203 |
| p_{t+4} | 316 | 637 | 0.4961 | 0.7783 | 0.5878 |
| p_{t+5} | 302 | 677 | 0.4461 | 0.7438 | 0.5373 |
| p_{t+6} | 318 | 817 | 0.3892 | 0.7833 | 0.5083 |
| p_{t+7} | 297 | 769 | 0.3862 | 0.7315 | 0.4863 |
| p_{t+8} | 297 | 783 | 0.3793 | 0.7315 | 0.4807 |
| p_{t+9} | 335 | 1075 | 0.3116 | 0.8251 | 0.4538 |
| p_{t+10} | 332 | 1129 | 0.2941 | 0.8177 | 0.4339 |
| Overall | 311 | 715 | 0.4358 | 0.7670 | 0.5389 |

Table 4.3: ConvLSTM-SA Depth, Frame Test dataset (406 perturbation points in 6000)

towards higher horizon predictors the AUC decreases, indicating lower predictors are performing better.

4.1.4 Generalization Test

Here we are checking the generalization capabilities of Perceptual Anticipation Model trained with data from the stimulation scenario where the bar is just rotating in the clockwise direction. We test the models learnt on data from scenarios like the bar rotating in the anticlockwise direction, two bars rotating in different combinations of rotations as shown in Fig. 4.4

In Fig. 4.5 we can see the spike visualization for different test scenario. The

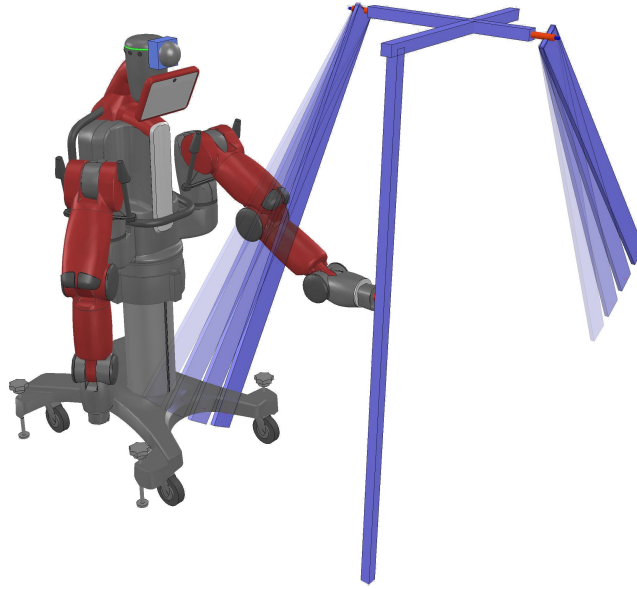


Figure 4.4: Dynamic scenario in V-REP simulation to test the generalization capability of the Perceptual Anticipation Model

configuration being tested is convLSTM-SA-DF. We can clearly see that the model is performing very well when its tested on data its trained on. Not only that, when there is a different dynamic object in the form of another rotating, the model is anticipating very well if the closer bar is rotating in the same direction it was trained on, irrespective of the rotation of the other bar. But when the closer bar is rotating towards the other direction the model can not anticipate the collision. Also, Fig. 4.6 show an example excerpt showing the anticipation of collision with the new scenario.

In Tab. 4.4, 4.5, 4.6 we can see the quantified results of the test errors for different models for different scenarios. We can see for all the models whenever the first bar is rotating clockwise, the model predictions are very good, while for whenever the first bar is rotating anti-clockwise the the models are not able to pick up any cues and hence its not able to anticipate the collision properly.

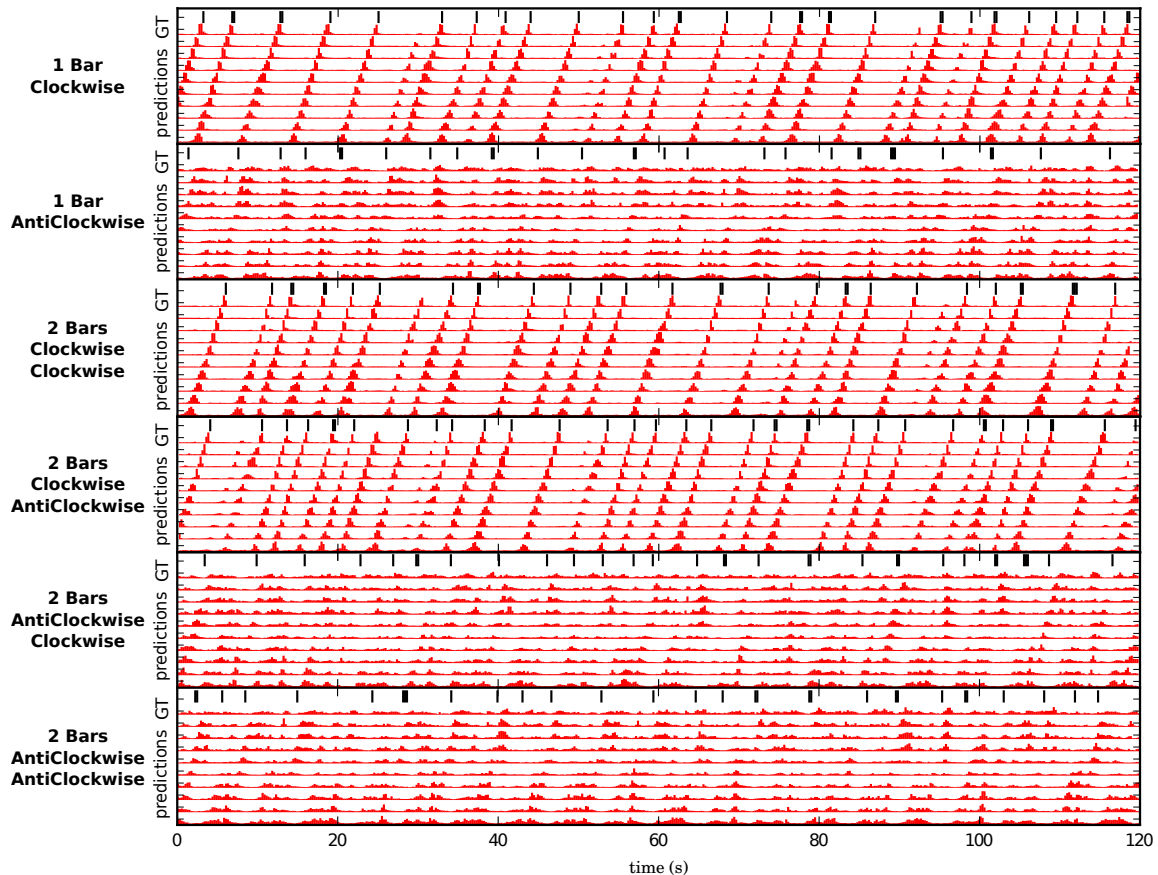


Figure 4.5: Spike visualization for different test scenarios of convLSTM-SA-DF configuration

4.2 Human-Robot Interaction

In the real-world scenario we have trained and tested the Perceptual Anticipation Model on anticipation of humans interacting with the robot while its performing a specific task. For training the anticipation model, we have created an interaction dataset involving 10 participants. Each participant interacted with the robot for about ten minutes. This resulted in a data set of 3000 data points per person. The participants were instructed to move towards the robot and touch its arm at any location. Throughout this process, the arm of the robot was continuously performing

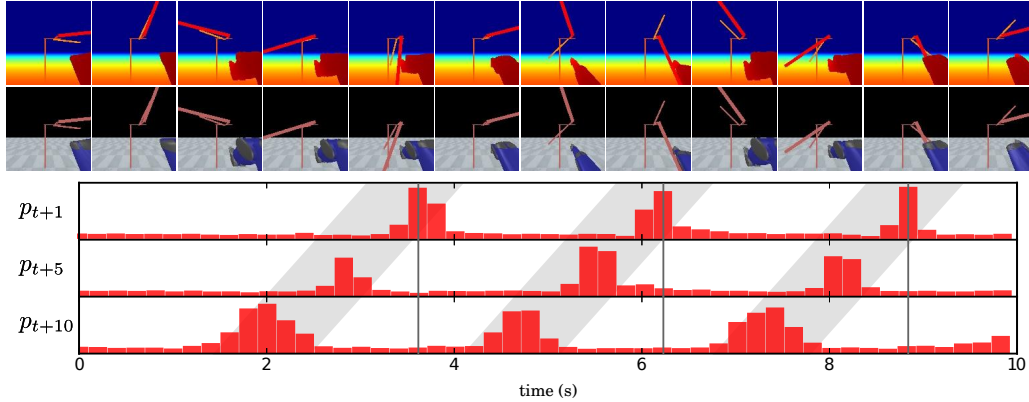


Figure 4.6: Predictions of the likelihood of physical perturbation generated from a perceptual anticipation model on a new scenario where another dynamic object present.

| Scenario | precision | recall | MCC |
|---------------------------------------|-----------|--------|--------|
| 1-Bar Clockwise | 0.4628 | 0.7649 | 0.5578 |
| 1-Bar Anti-Clockwise | 0.0707 | 0.8552 | 0.0466 |
| 2-Bars Clockwise, Clockwise | 0.4764 | 0.7662 | 0.5667 |
| 2-Bars Clockwise, Anti-Clockwise | 0.4268 | 0.7166 | 0.5101 |
| 2-Bars Anti-Clockwise, Clockwise | 0.0735 | 0.8894 | 0.0402 |
| 2-Bars Anti-Clockwise, Anti-Clockwise | 0.0741 | 0.8816 | 0.0431 |

Table 4.4: Test Error for different scenarios for model ConvLSTM-SA-D

a forward-backward movement.

In addition, in 40-50% of the interactions, the participants were instructed to pretend approaching the robot and then turning back without any contact. Incorporating such *feigning moves* on the part of the human interaction partner, as well as a constant movement of the arm, ensures that the robot has to constantly monitor its state and the environment in order to appropriately update its belief about impend-

| Scenario | precision | recall | MCC |
|---------------------------------------|-----------|--------|--------|
| 1-Bar Clockwise | 0.5122 | 0.6356 | 0.5355 |
| 1-Bar Anti-Clockwise | 0.0757 | 0.6247 | 0.0480 |
| 2-Bars Clockwise, Clockwise | 0.4636 | 0.7390 | 0.5461 |
| 2-Bars Clockwise, Anti-Clockwise | 0.3747 | 0.7730 | 0.4904 |
| 2-Bars Anti-Clockwise, Clockwise | 0.0774 | 0.8329 | 0.0561 |
| 2-Bars Anti-Clockwise, Anti-Clockwise | 0.0803 | 0.6087 | 0.0467 |

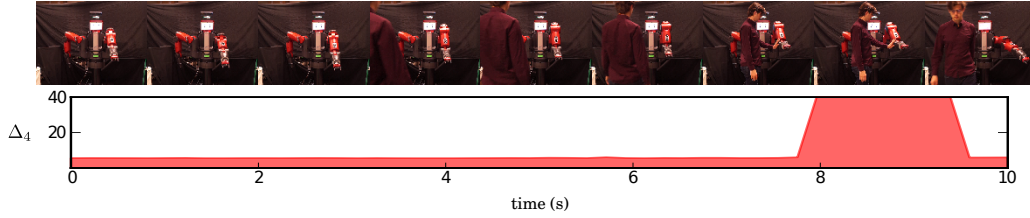
Table 4.5: Test Error for different scenarios for model ConvLSTM-SA-F

| Scenario | precision | recall | MCC |
|---------------------------------------|-----------|--------|--------|
| 1-Bar Clockwise | 0.4633 | 0.7634 | 0.5576 |
| 1-Bar Anti-Clockwise | 0.0661 | 0.9031 | 0.0158 |
| 2-Bars Clockwise, Clockwise | 0.4830 | 0.7662 | 0.5715 |
| 2-Bars Clockwise, Anti-Clockwise | 0.4011 | 0.7664 | 0.5097 |
| 2-Bars Anti-Clockwise, Clockwise | 0.0694 | 0.9923 | 0.0140 |
| 2-Bars Anti-Clockwise, Anti-Clockwise | 0.0691 | 0.9942 | 0.0066 |

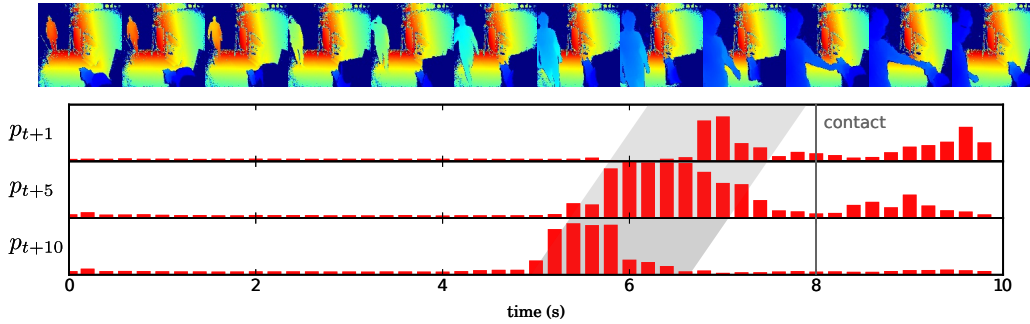
Table 4.6: Test Error for different scenarios for model ConvLSTM-SA-DF

ing physical contact. Data recorded from 9 of the participants was used for training, while 1 participant was used for validation. Data from a separate set of 3 participants was used as test data to evaluate the generalization capabilities of the model.

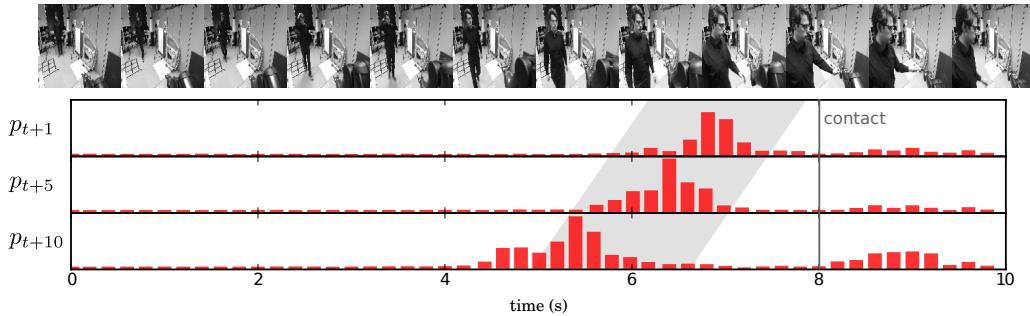
Two different versions of the training set were generated. First, we created a setup in which depth images were used as input. In the second setup, we used grayscale video images as input.



(a) Detection of physical perturbations using the Deep Dynamics Model. A human subject approaches the robot, pushes the arm, and moves back. Detected perturbations correspond to the moment of contact and release.



(b) Predictions of the likelihood of physical perturbation generated from a perceptual anticipation model which was trained on depth input. The predictors fire at different moments ahead of the actual moment of contact.



(c) Predictions of the likelihood of physical perturbation generated from a perceptual anticipation model which was trained on grayscale video input. The predictors fire at different moments ahead of the actual moment of contact.

Figure 4.7: Example Interaction Visualization

4.2.1 Example Interaction

Fig. 4.7a shows an example interaction between the robot and a human subject applying forces on its arm. The perturbation identified by the Deep Dynamics Model accurately reflects the moment of contact between the human and the robot. These perturbations are used to train the anticipation model, so that in subsequent interactions the robot can predict the occurrence of a perturbation by only observing the human approaching and lifting his hand.

This process can be seen in Fig. 4.7b. After training the perceptual anticipation model, we provide current observations in form of depth images as input. In turn, the network generates estimates for the next k time steps which reflect the likelihood of a future perturbation at $t + k$. The moment of contact is depicted in Fig. 4.7b by a vertical line. We can see that the predictions for p_{t+1}, p_{t+5} and p_{t+10} are aligned along a diagonal. The predictor with a larger horizon, i.e., looking ten time steps into the future, activates early to indicate a high likelihood of a contact. The predictors with a shorter horizon activate at later time steps, based on the horizon they have been trained on. Note the attenuation of the activations, once the robot is outside the envelope for which the predictors have been trained to fire.

Fig. 4.7c shows the same process with a perceptual anticipation model trained on typical grayscale video images. Again, the predictors fire mostly within the envelope (gray) imposed by the respective horizons, with p_{t+10} firing first. However, the responses are less accentuated when compared with the activations generated from depth images.



Figure 4.8: Saliency map visualization of p_{t+5} predictor. It anticipates any interactions into the future. The images show, highlighted in red, the pixels the network is paying attention when generating an output. In the beginning, almost no region is active, but as the human starts approaching the robot, the network focuses on the human and also on the robot arm. Once the subject is in contact with the robot, the activations attenuate.

4.2.2 Visualizing Network Saliency

To better understand the decision process by which the perceptual anticipation model generates predictions, we visualized the underlying saliency using the method introduced in [26]. Fig. 4.8 shows an example of saliency maps at different moments in time during a human-robot interaction. Regions highlighted in red or yellow correspond to pixels with a strong influence on the output of the network. Originally, the network is not focusing on any particular region within the image. However, as the human subject starts walking in direction of the robot, the network approximately starts focusing on pixels around the body of both the human subject, as well as the right arm of the robot. As the human comes closer, the saliency scores for each pixel become larger, as indicated by the yellow coloring. At the onset of a contact, the saliency scores attenuate.

The latter analysis of saliency, is in line with our expectations of features relevant to the task, i.e., the shape and location of the human and robot arm. It is important to note that the activations for this particular predictor (p_{t+5}) quickly attenuate as the human comes close to the robot. We noticed that this pattern is dependent on the specific predictor. Consequently, the saliency maps for the p_{t+1} attenuate at a later moment, while the saliency maps for p_{t+10} attenuate at an earlier moment.

4.2.3 Prediction Error on Training and Test Sets

The prediction error on the training set can be seen in Tab. 4.7 and test set can be seen in Tab. 4.8. Different models configurations were tested, namely combinations of (a) ConvLSTM vs. ConvGRU, (b) state-action conditioned networks vs only visual data, and (c) depth data vs. grayscale frame data. The predictions are thresholded to get optimal MCC (Matthews correlation coefficient) in each of the configuration.

| Config | true+ | pred+ | precision | recall | MCC |
|---------------|-------|-------|-----------|--------|--------|
| ConvLSTM-SA-F | 7907 | 12509 | 0.6321 | 0.8511 | 0.7228 |
| ConvLSTM-F | 7803 | 13534 | 0.5765 | 0.8399 | 0.6833 |
| ConvLSTM-D | 7624 | 13027 | 0.5852 | 0.8207 | 0.6804 |
| ConvGRU-D | 7711 | 13451 | 0.5733 | 0.8300 | 0.6769 |
| ConvGRU-SA-D | 7724 | 13680 | 0.5646 | 0.8314 | 0.6720 |
| ConvLSTM-SA-D | 7513 | 13114 | 0.5729 | 0.8087 | 0.6675 |
| ConvGRU-SA-F | 7782 | 14307 | 0.5439 | 0.8377 | 0.6612 |
| ConvGRU-F | 7134 | 13262 | 0.5379 | 0.7679 | 0.6278 |

Table 4.7: Comparison between different configurations of the DPM on Train dataset (9290 perturbation points in 270000)

| Config | true+ | pred+ | precision | recall | MCC |
|---------------|-------|-------|-----------|--------|--------|
| ConvLSTM-SA-D | 1907 | 4607 | 0.4139 | 0.6421 | 0.4953 |
| ConvLSTM-D | 2087 | 5999 | 0.3479 | 0.7027 | 0.4711 |
| ConvGRU-D | 1686 | 4060 | 0.4153 | 0.5677 | 0.4651 |
| ConvGRU-SA-D | 1876 | 5206 | 0.3604 | 0.6316 | 0.4541 |
| ConvGRU-F | 2129 | 6696 | 0.3180 | 0.7168 | 0.4522 |
| ConvGRU-SA-F | 2040 | 6360 | 0.3208 | 0.6869 | 0.4442 |
| ConvLSTM-SA-F | 1993 | 6190 | 0.3220 | 0.6710 | 0.4396 |
| ConvLSTM-F | 2127 | 8039 | 0.2646 | 0.7162 | 0.4060 |

Table 4.8: Comparison between different configurations of the DPM on Test dataset (2970 perturbation points in 90000)

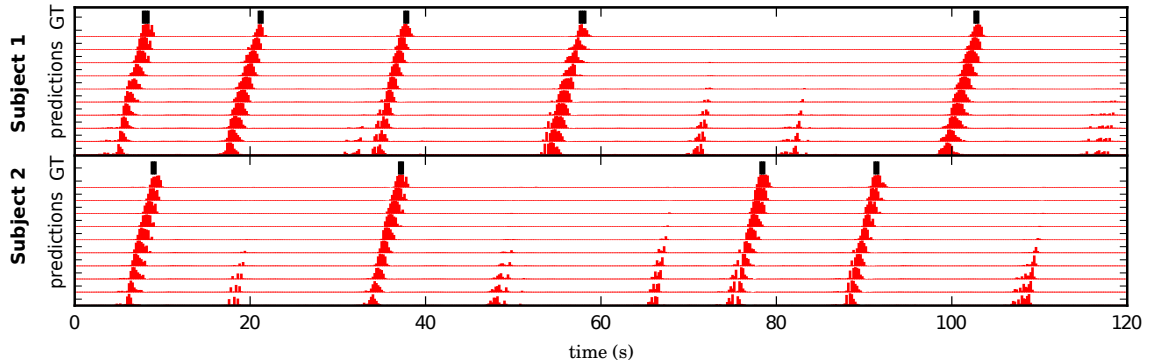
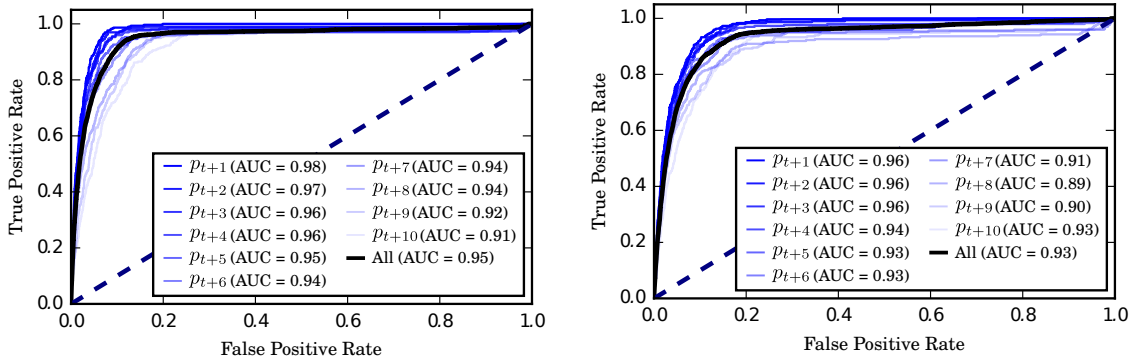


Figure 4.9: Visualization of the ground truth timing (black) of a physical perturbation and the activations (red) of the anticipation model for two test subjects.

Comparing the performance of different architectures we can say convLSTM type architectures are performing slightly better than convGRU type architectures. Also, models trained on depth data are generalizing better than those trained on frame data.

| pred no | true+ | pred+ | precision | recall | MCC |
|----------------|-------|-------|-----------|--------|--------|
| p_{t+1} | 239 | 526 | 0.4544 | 0.8047 | 0.5877 |
| p_{t+2} | 227 | 487 | 0.4661 | 0.7643 | 0.5799 |
| p_{t+3} | 233 | 522 | 0.4464 | 0.7845 | 0.5742 |
| p_{t+4} | 224 | 531 | 0.4218 | 0.7542 | 0.5451 |
| p_{t+5} | 215 | 481 | 0.4470 | 0.7239 | 0.5507 |
| p_{t+6} | 192 | 424 | 0.4528 | 0.6465 | 0.5226 |
| p_{t+7} | 212 | 576 | 0.3681 | 0.7138 | 0.4905 |
| p_{t+8} | 166 | 443 | 0.3747 | 0.5589 | 0.4352 |
| p_{t+9} | 187 | 647 | 0.2890 | 0.6296 | 0.3989 |
| p_{t+10} | 194 | 790 | 0.2456 | 0.6532 | 0.3691 |
| Overall | 191 | 461 | 0.4139 | 0.6421 | 0.4953 |

Table 4.9: ConvLSTM-SA Depth Test dataset (297 perturbation points in 9000)



(a) ROC(Receiver Operating Characteristic) plot with convLSTM-SA-D (b) ROC(Receiver Operating Characteristic) plot with convLSTM-SA-F

Figure 4.10: ROC Plots for Test Data after Model Selection

| pred no | true+ | pred+ | precision | recall | MCC |
|----------------|-------|-------|-----------|--------|--------|
| p_{t+1} | 202 | 492 | 0.4106 | 0.6801 | 0.5083 |
| p_{t+2} | 208 | 539 | 0.3859 | 0.7003 | 0.4986 |
| p_{t+3} | 203 | 530 | 0.3830 | 0.6835 | 0.4901 |
| p_{t+4} | 196 | 516 | 0.3798 | 0.6599 | 0.4788 |
| p_{t+5} | 182 | 486 | 0.3745 | 0.6128 | 0.4567 |
| p_{t+6} | 193 | 545 | 0.3541 | 0.6498 | 0.4564 |
| p_{t+7} | 209 | 641 | 0.3261 | 0.7037 | 0.4543 |
| p_{t+8} | 171 | 493 | 0.3469 | 0.5758 | 0.4230 |
| p_{t+9} | 174 | 535 | 0.3252 | 0.5859 | 0.4113 |
| p_{t+10} | 185 | 755 | 0.2450 | 0.6229 | 0.3592 |
| Overall | 199 | 619 | 0.3220 | 0.6710 | 0.4396 |

Table 4.10: ConvLSTM-SA Frame Test dataset (297 perturbation points in 9000)

4.2.4 Test Error after Model Selection

Given the above results, we selected a state-action conditioned ConvLSTM (ConvLSTM-SA) as the underlying network model for the prediction. Consequently, we analyzed the prediction errors for the generated predictions p_{t+1}, \dots, p_{t+k} separately. Tab. 4.9 and Tab. 4.10 show the results of this analysis for depth data and frame data respectively. In both cases precision and recall deteriorate, as the network produces predictions for longer horizons. However, using depth input generates significantly better results for short horizon predictors.

Similarly in Fig. 4.10a and Fig. 4.10b we can see the ROC(Receiver Operating Characteristic) plot for test data with different configurations. Decreasing AUC(area under curve) indicating how well the different predictors are performing. We can see

that as we go towards higher horizon predictors the AUC decreases, indicating lower predictors are performing better.

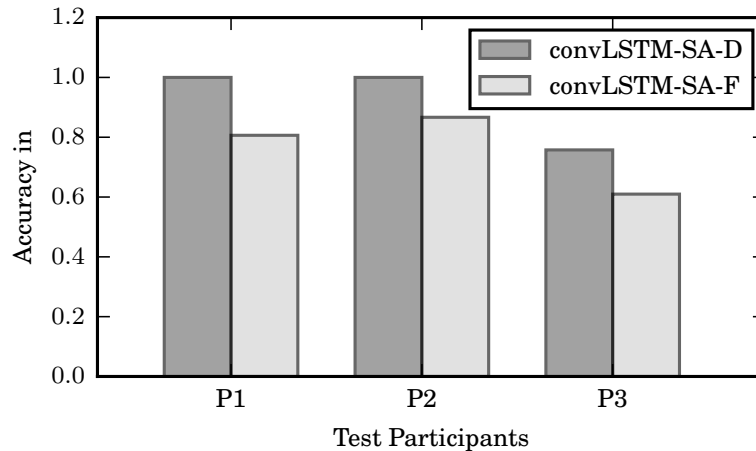


Figure 4.11: Percentage of times all ten predictors activate ahead of ground truth activation.

The above tables show a relatively low precision on the test set. To investigate this phenomenon, we visualized the ground truth versus the generated predictions by the network. An excerpt of this can be seen in Fig. 4.9. The figure shows the network activations for two different subjects over time. The black marks indicate the moments of contact, i.e., the ground truth (GT). The ten plots in red underneath the ground truth indicate the output predictions of the network. We can see that in the majority of cases the network performs the right classification. We can also see that the activations are aligned along a diagonal, since they outputs fire for different horizons. In some cases (Subject 2, time step 18) the earlier predictors start to fire but immediately cease thereafter. An analysis of the video showed that these cases correspond to the feigning moves the subjects were asked to perform from time to time. This shows that the network performs according to our expectations: the early predictors fire, while the late predictors are awaiting more evidence. Hence, the low

accuracy is mostly caused by the network sometimes firing slightly ahead of time.

To get a better estimate of the overall prediction accuracy we performed another evaluation in which we counted the number of times all ten predictors activate ahead of a ground truth activation for different participants, see Fig. 4.11. We can see that the networks performed well for a two out of three subjects, i.e., accuracy of 90% using depth images. For participant P3, accuracies dropped to about 75% (depth) and 60% (grayscale frame data).

Chapter 5

FUTURE WORK

For future work, we aim at extending the framework, such that the predictions can be conditioned on entire control policies of the robot and not only a single next action. Furthermore, we will incorporate prediction framework into reinforcement learning, in order to also autonomously learn preventive motions for avoiding perturbations.

Chapter 6

CONCLUSIONS

In this thesis work, we introduced a novel methodology that allows robots to associate perceptual cues with impending physical impact to the body. By learning a mapping between observed visual features and future perturbations, robots can anticipate upcoming hazardous or unwanted states. To this end, we introduced a complex neural network model that combines spatial, temporal, and probabilistic reasoning in order to generate predictions over a horizon of next time steps. The network learns to focus on visual features that are most indicative of future exogenous perturbations.

Our experiments show that a humanoid robot can effectively learn when physical contact with a human interaction partner will occur. The method autonomously learns to anticipate physical perturbations from any source and is not restricted to human-robot interactions. One interesting insight is that the network learned to roughly identify the human shape, as well as the shape of the robot without any supervised training data or labelling.

REFERENCES

- [1] Amor, H. B., G. Neumann, S. Kamthe, O. Kroemer and J. Peters, “Interaction primitives for human-robot cooperation tasks”, in “2014 IEEE International Conference on Robotics and Automation (ICRA)”, pp. 2831–2837 (2014).
- [2] Asimov, I., *I, Robot* (Fawcett Publications, 1950).
- [3] Berger, E., D. Müller, D. Vogt, B. Jung and H. B. Amor, “Transfer entropy for feature extraction in physical human-robot interaction: Detecting perturbations from low-cost sensors”, in “2014 IEEE-RAS International Conference on Humanoid Robots”, pp. 829–834 (IEEE, 2014).
- [4] Berger, E., M. Sastuba, D. Vogt, B. Jung and H. B. Amor, “Dynamic mode decomposition for perturbation estimation in human robot interaction”, in “The 23rd IEEE International Symposium on Robot and Human Interactive Communication”, pp. 593–600 (IEEE, 2014).
- [5] Berger, E., M. Sastuba, D. Vogt, B. Jung and H. Ben Amor, “Estimation of perturbations in robotic behavior using dynamic mode decomposition”, *Advanced Robotics* **29**, 5, 331–343 (2015).
- [6] chan, f.-h., y.-t. chen, y. xiang and m. sun, “Anticipating accidents in dashcam videos”, in “asian conference computer vision (accv)”, (????).
- [7] Cho, K., B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk and Y. Bengio, “Learning phrase representations using rnn encoder-decoder for statistical machine translation”, arXiv preprint arXiv:1406.1078 (2014).
- [8] De Luca, A., A. Albu-Schaffer, S. Haddadin and G. Hirzinger, “Collision detection and safe reaction with the dlr-iii lightweight manipulator arm”, in “International Conference on Intelligent Robots and Systems (IROS), 2006 IEEE/RSJ”, pp. 1623–1630 (IEEE, 2006).
- [9] Demiris, Y. and A. Dearden, “From motor babbling to hierarchical learning by imitation: a robot developmental pathway”, (2005).
- [10] Donahue, J., L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko and T. Darrell, “Long-term recurrent convolutional networks for visual recognition and description”, in “The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)”, (2015).
- [11] Eom, K. S., I. H. Suh, W. K. Chung and S. R. Oh, “Disturbance observer based force control of robot manipulator without force sensor”, in “Proceedings. 1998 IEEE International Conference on Robotics and Automation (Cat. No.98CH36146)”, vol. 4, pp. 3012–3017 vol.4 (1998).

- [12] Gal, Y. and Z. Ghahramani, “Dropout as a bayesian approximation: Representing model uncertainty in deep learning”, arXiv preprint arXiv:1506.02142 **2** (2015).
- [13] Haddadin, S., *Towards Safe Robots: Approaching Asimov’s 1st Law* (Springer Publishing Company, Incorporated, 2013).
- [14] Haddadin, S., A. Albu-Schaffer, A. De Luca and G. Hirzinger, “Collision detection and reaction: A contribution to safe physical human-robot interaction”, in “2008 IEEE/RSJ International Conference on Intelligent Robots and Systems”, pp. 3356–3363 (IEEE, 2008).
- [15] Hinton, G. E., N. Srivastava, A. Krizhevsky, I. Sutskever and R. R. Salakhutdinov, “Improving neural networks by preventing co-adaptation of feature detectors”, arXiv preprint arXiv:1207.0580 (2012).
- [16] Hochreiter, S. and J. Schmidhuber, “Long short-term memory”, *Neural computation* **9**, 8, 1735–1780 (1997).
- [17] Iwata, H., H. Hoshino, T. Morita and S. Sugano, “Force detectable surface covers for humanoid robots”, in “2001 IEEE/ASME International Conference on Advanced Intelligent Mechatronics. Proceedings (Cat. No.01TH8556)”, vol. 2, pp. 1205–1210 vol.2 (2001).
- [18] Kaiser, L. and I. Sutskever, “Neural gpu learn algorithms”, arXiv preprint arXiv:1511.08228 (2015).
- [19] Kuehn, J. and S. Haddadin, “An artificial robot nervous system to teach robots how to feel pain and reflexively react to potentially damaging contacts”, *IEEE Robotics and Automation Letters* **2**, 1, 72–79 (2017).
- [20] Mainprice, J., R. Hayne and D. Berenson, “Goal set inverse optimal control and iterative replanning for predicting human reaching motions in shared workspaces”, *IEEE Transactions on Robotics* **32**, 4, 897–908 (2016).
- [21] Murakami, T., F. Yu and K. Ohnishi, “Torque sensorless control in multidegree-of-freedom manipulator”, *IEEE Transactions on Industrial Electronics* **40**, 2, 259–265 (1993).
- [22] Nair, V. and G. E. Hinton, “Rectified linear units improve restricted boltzmann machines”, in “Proceedings of the 27th international conference on machine learning (ICML-10)”, pp. 807–814 (2010).
- [23] Pérez-D’Arpino, C. and J. A. Shah, “Fast target prediction of human reaching motion for cooperative human-robot manipulation tasks using time series classification”, in “2015 IEEE International Conference on Robotics and Automation (ICRA)”, (2015).
- [24] Pinheiro, P. H. and R. Collobert, “Recurrent convolutional neural networks for scene labeling.”, in “ICML”, pp. 82–90 (2014).

- [25] Saegusa, R., G. Metta, G. Sandini and S. Sakka, “Active motor babbling for sensorimotor learning”, in “Robotics and Biomimetics, 2008. ROBIO 2008. IEEE International Conference on”, pp. 794–799 (IEEE, 2009).
- [26] Simonyan, K., A. Vedaldi and A. Zisserman, “Deep inside convolutional networks: Visualising image classification models and saliency maps”, arXiv preprint arXiv:1312.6034 (2013).
- [27] Xingjian, S., Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong and W.-c. Woo, “Convolutional lstm network: A machine learning approach for precipitation nowcasting”, in “Advances in Neural Information Processing Systems”, pp. 802–810 (2015).

APPENDIX A

IRB APPROVAL FOR HUMAN-ROBOT INTERACTION DATA COLLECTION

APPROVAL: EXPEDITED REVIEW

[Hani Ben Amor](#)
[Computing, Informatics and Decision Systems Engineering, School of \(CIDSE\)](#)

-

Hani.Benamor@asu.edu

Dear [Hani Ben Amor](#):

On 3/21/2017 the [ASU IRB](#) reviewed the following protocol:

| | |
|---------------------|--|
| Type of Review: | Initial Study |
| Title: | Anticipating Touch from Visual Cues |
| Investigator: | Hani Ben Amor |
| IRB ID: | STUDY00005986 |
| Category of review: | (7)(b) Social science methods, (7)(a) Behavioral research |
| Funding: | None |
| Grant Title: | None |
| Grant ID: | None |
| Documents Reviewed: | <ul style="list-style-type: none"> • Informed_Consent-2-2.pdf, Category: Consent Form; • Recruitment Flyer.pdf, Category: Recruitment Materials; • anticipation-benamor-2.docx, Category: IRB Protocol; |

The IRB approved the protocol from 3/21/2017 to 3/20/2018 inclusive. Three weeks before 3/20/2018 you are to submit a completed Continuing Review application and required attachments to request continuing approval or closure.

If continuing review approval is not granted before the expiration date of 3/20/2018 approval of this protocol expires on that date. When consent is appropriate, you must use final, watermarked versions available under the “Documents” tab in ERA-IRB.

In conducting this protocol you are required to follow the requirements listed in the INVESTIGATOR MANUAL (HRP-103).

Sincerely,

IRB Administrator

cc:

Indranil Sur