

UNIVERSITÄT HAMBURG

**Artificial Neural Models for Feedback  
Pathways for Sensorimotor  
Integration**

by

Junpei Zhong

Dissertation with the aim of achieving a doctoral degree  
at the Faculty of Mathematics, Informatics and Natural Sciences  
Department of Computer Science  
of Universität Hamburg

Submitted by Junpei Zhong  
2015 in Hamburg

**Evaluators:**

Prof. Dr. Angelo Cangelosi

Prof. Dr. Stefan Wermter

Prof. Dr. Jianwei Zhang

Date of Disputation: 17th April, 2015

*“Biology gives you a brain. Life turns it into a mind.”*

Jeffrey Eugenides

# *Abstract*

The brain comprises hierarchical modules on various physiological levels. Neural feedback signals (including lateral and top-down connections) modulate the neural activities via inhibitory or excitatory connections within/between these levels. They have predictive and filtering functions influencing the neuronal population coding of the bottom-up sensory-driven signals in the perception-action system.

In this thesis, we propose that the predictive role of the feedback pathways at most levels of action and perception can be modelled by the recurrent connections in different artificial cognitive platforms (simulation and humanoid robots). This will be examined by three recurrent neural network models. Furthermore, the three models and experiments with them show that the recurrent neural networks are able to model feedback pathways and to exhibit the feedback-related sensorimotor predictive functions.

In the first model, inspired by the study of neurobiology, we emphasize that the feedback connections facilitate a predictive mechanism to compensate for the neural delay in the two streams (ventral and dorsal) of the visual system. We model this with a novel recurrent network with a horizontal product. In the simulation, the recurrent connections give rise to the fast- and slow-changing neural activations in the dorsal- and ventral-like hidden layer. Particularly the recurrent connections build a feedback channel to predict the upcoming neural activity in the dorsal-like hidden layer, while another feedback channel maintains stable neural encoding in the ventral-like hidden layer.

In the second part of the thesis, a sensorimotor integration model with visual prediction is implemented, whose visual perception part is considered to be the dorsal stream representation of the first model. This further augments the visual prediction with its role of guiding motor action. Together with the action module which adopts a continuous reinforcement learning algorithm, this model allows a smooth and faster docking behaviour for a humanoid robot.

In the third experiment, we propose that the source of the feedback pathway could be the high-level cognitive processes, such as pre-symbolic representations. Furthermore, the emergence of these cognitive processes and feedback-related sensorimotor functions are not independent processes but they integrate and assist each other in a hierarchical way. Therefore, we augment the first horizontal product

model with additional units, called parametric bias (PB) units, as a pre-symbolic representation. In the robot experiments, we show that during the learning process of observing sensorimotor primitives, the pre-symbolic representation is self-organized in the parametric units; during prediction, these representational units act as a prior expectation which guides the robot to recognize and to expect various pre-learned sensorimotor primitives.

These three experiments demonstrate that implementation of the feedback pathways with recurrent connections can realize predictive sensorimotor functions. The emergence of these feedback pathways also accounts for the pre-symbolic representation in cognitive systems. Furthermore, we claim that the recurrent connections can be one of possible neural structures to build up the feedback pathways on the sensorimotor integration in artificial cognitive systems.

# *Zusammenfassung*

Das Gehirn besteht aus hierarchisch angeordneten Modulen auf verschiedenen physiologischen Ebenen. Neuronale Rückkopplungen (einschließlich lateraler und hierarchischer Verbindungen) modulieren die neuronalen Aktivitäten über hemmende oder anregende Verbindungen innerhalb sowie zwischen diesen Ebenen. Die Rückkopplungen haben Prädiktions- und Filterfunktionen bezogen auf die neuronale Codierung der sensorisch getriebenen Signale im Wahrnehmungs-Aktionssystem.

In dieser Arbeit stellen wir die Hypothese auf, dass die prädiktive Rolle der Rückkopplungen auf den meisten Ebenen von Wahrnehmung und Handlung durch die rekurrente Verbindungen in verschiedenen künstlichen kognitiven Plattformen (Simulation und humanoide Roboter) modelliert werden kann. Dies wird anhand von drei rekurrenten neuronalen Netzwerkmodellen untersucht. Darüber hinaus zeigen unsere Experimente mit den drei Modellen, dass die rekurrenten neuronalen Netze Rückkopplungen modellieren und rückkopplungsbezogene sensomotorisch prädiktive Funktionen aufweisen.

Im ersten Modell, inspiriert durch das Studium der Neurobiologie, betonen wir, dass die rekurrenten Verbindungen einen prädiktiven Mechanismus ermöglichen, der neuronale Verzögerung in den beiden Strömen (ventral und dorsal) des visuellen Systems kompensiert. Wir modellieren dies mit einem neuartigen rekurrenten neuronalen Netzwerk als horizontales Produkt. In der Simulation führen die rekurrenten Verbindungen zu sich schnell und langsam ändernden neuronalen Aktivierungen in der verborgenen ventralen und dorsalen Schicht. Dabei bilden die rekurrenten Verbindungen einen Rückkopplungskanal, um die kommende neuronale Aktivität in der dorsalen Schicht vorherzusagen, während ein anderer Rückkopplungskanal eine stabile neuronale Codierung in der ventralen Schicht aufrechterhält.

Im zweiten Teil der Arbeit wird ein sensomotorisches Integrationsmodell mit visueller Vorhersage implementiert, dessen visueller Teil als weitergehende Implementierung des dorsalen Pfads des ersten Modells verstanden werden kann. Diese erweitert die visuelle Vorhersage durch die Funktion, motorische Aktionen auszuführen. Zusammen mit dem Aktionsmodul, das einen Algorithmus zum kontinuierlichen Verstärkungslernen einsetzt, erlaubt dieses Modell ein reibungsloses und schnelleres Dockingverhalten für einen humanoiden Roboter.

Im dritten Versuch stellen wir die Hypothese auf, dass die Quelle der Rückkopplungen höhere kognitive Prozesse, wie zum Beispiel präsymbolische Repräsentationen sein können. Darüber hinaus sind die Entstehung dieser kognitiven Prozesse und rückkopplungs-bezogenen sensomotorische Funktionen nicht unabhängige Prozesse, sondern sie integrieren und unterstützen sich gegenseitig auf eine hierarchische Weise. Deshalb erweitern wir das erste horizontale Produktmodell um zusätzliche Einheiten, die so sogenannten parametrischen Bias (PB) Einheiten, als präsymbolische Repräsentation. In Roboterexperimenten zeigen wir, dass während des Lernprozesses, bei dem sensomotorische Primitive beobachtet werden, sich die präsymbolische Repräsentation in den PB Einheiten selbstorganisiert. Während der Vorhersage wirken diese Darstellungseinheiten als Vorerwartung, die den Roboter dazu führt, verschiedene vorher gelernte sensomotorische Primitive zu erwarten und zu erkennen.

Diese drei Versuche zeigen, dass die Implementierung der Rückkopplungen mit rekurrenten Verbindungen eine prädiktive Sensomotorik verwirklichen kann. Die Entstehung dieser Rückkopplungen ist auch für die präsymbolische Repräsentationen in kognitiven Systemen verantwortlich. Außerdem behaupten wir, dass rekurrente Verbindungen mögliche neuronale künstliche Netzwerkstrukturen zum Aufbau der Rückkopplungen für die sensomotorische Integration in künstlichen kognitiven Systeme sein können.

## *Acknowledgements*

The past four years have been quite a marathon both academically and personally. First and foremost I would like to thank my supervisor, Prof. Stefan Wermter. His serious and diligent attitude toward research will be a good example for my future career. Besides, I have also benefited from his plenty support to my daily life in Hamburg, especially at the initial stage of my doctoral years.

Furthermore, I am extremely grateful to Dr. Cornelius Weber, for his intensive and valuable discussions about my research work and project ideas. He also showed me in detail how to conduct research and design an engineering product preciously as a (German) scientist and a (German) engineer.

I also have to thank Prof. Angelo Cangelosi of Plymouth University, who gave continuous support to me and to other RobotDoC fellows. He is a good example about how to enjoy research.

Thanks also go to my colleagues of the knowledge technology (WTM) group who always gave me discussions and suggestions through my doctoral career. Special thanks to Katja Kösters and Erik Strahl who often helped me to solve a lot of problems and to Nicolás Navarro-Guerrero for his fruitful discussion about work and other stuff.

I will be remembering the days with other RobotDoC fellows across Europe. We always had a great moment when we gathered (and travelled) every half year.

I would also like to express my gratitude to Dr. Lola Cañamero and other colleagues in her group for their care and support when I was finishing this thesis at University of Hertfordshire.

Applause should also be sent to Kira Chow, who did two of the lovely drawings (Figs. 2.6 and 8.2) for this thesis. I also acknowledge the language improvement by Stefan Wermter, Katja Kösters, Cornelius Weber, Alex Smith, Junlei Yu and Alexandra Lindqvist.

Finally, I would like to thank the EU project RobotDoC under 235065 from the FP7, Marie Curie Action ITN, for its financial support.

I dedicate this work to all of my friends in China, Hong Kong and Europe, for their continuous supports and endless asking when I am going to finish.

The end of a marathon is not an end of training. I guess an ultra-marathon will be the next target.





# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Cognitive Robotics . . . . .	1
1.2	Modelling Feedback Pathways on Cognitive Robotic Systems . . .	4
1.3	Research Questions and Contributions . . . . .	5
1.4	Structure of the Thesis . . . . .	6
<b>2</b>	<b>Feedback Pathways in a Hierarchical Modularity Brain</b>	<b>8</b>
2.1	Sensorimotor Integration . . . . .	8
2.2	Sensorimotor System . . . . .	9
2.2.1	Sensory and Motor Cortices . . . . .	9
	Sensory Cortices . . . . .	9
	Motor System . . . . .	12
2.3	Functional Modularity of Sensorimotor System . . . . .	14
2.4	Multi-dimensional Hierarchical Modularity . . . . .	19
2.5	Neural Feedback Connections . . . . .	21
2.5.1	Embodied Feedback Pathways . . . . .	23
2.6	Phenomena from Feedback Pathways . . . . .	24
2.6.1	Binocular Rivalry . . . . .	24
2.6.2	Retina Prediction . . . . .	26
2.6.3	Feedback Pathways on Predictive Action . . . . .	27
2.7	Representation as Bayesian Inference . . . . .	27
2.8	Hypotheses of Feedback Integration . . . . .	30
2.8.1	Predictive Coding . . . . .	30
2.8.2	Biased Competition . . . . .	32
2.9	Discussion . . . . .	33
2.10	Summary . . . . .	35
<b>3</b>	<b>Artificial Recurrent Neural Network Models of Feedback Pathways</b>	<b>36</b>
3.1	Introduction . . . . .	36
3.2	Multi-Layer Perceptron . . . . .	37
3.2.1	A Single Perceptron . . . . .	37
3.2.2	Multi-Layer Perceptron Network . . . . .	38

3.2.3	Back-propagation . . . . .	40
3.3	Artificial Recurrent Neural Networks (ARNNs) . . . . .	42
3.3.1	Recurrent Connections . . . . .	42
3.3.2	Back-propagation through Time . . . . .	43
3.3.3	Variants of ARNN . . . . .	47
	RNN with Parametric Biases . . . . .	47
	Echo State Network . . . . .	50
	Long Short Term Memory . . . . .	51
3.4	Discussion . . . . .	52
3.5	Summary . . . . .	53
<b>4</b>	<b>Perception-Action Model with Hierarchical Feedback Pathways</b>	<b>55</b>
4.1	Perception-Action Model . . . . .	55
4.2	Neural Basis of Architecture . . . . .	56
4.2.1	Somatotopic Arrangement of Motor and Visual Cortices . . . . .	56
4.2.2	Two Visual Streams for Action . . . . .	58
4.2.3	Mirror Neurons and Ideomotor Principle . . . . .	59
4.3	Sensorimotor Integration Architecture . . . . .	59
4.4	Summary . . . . .	62
<b>5</b>	<b>Feedback-influenced Motion-coding in Visual Cortex</b>	<b>63</b>
5.1	The Visual System . . . . .	63
5.1.1	Two-stream Theory . . . . .	63
5.1.2	Identification and Tracking . . . . .	66
5.1.3	Motivation . . . . .	68
5.2	Horizontal Recurrent Model . . . . .	69
5.2.1	Horizontal Product . . . . .	69
5.2.2	Algorithm . . . . .	71
	Training . . . . .	71
	Intrinsic Plasticity . . . . .	74
5.3	Experiments . . . . .	75
5.4	Summary . . . . .	78
<b>6</b>	<b>Feedback Signals on a Predictive Sensorimotor System</b>	<b>79</b>
6.1	The Sensorimotor System . . . . .	79
6.1.1	Sensory Prediction . . . . .	79
6.1.2	Tracking and Prediction . . . . .	80
6.1.3	Motivation . . . . .	81
6.1.4	Experiments Setting . . . . .	82
6.2	Recurrent Prediction Sensorimotor Model . . . . .	84
6.2.1	Landmark-based Detection . . . . .	84
6.2.2	Algorithm . . . . .	85
	Visual Prediction via Recurrent Connections . . . . .	85
	Smooth Action Generation . . . . .	86
6.3	Experiments . . . . .	88
6.3.1	Training Scheme . . . . .	88

---

6.3.2	Experimental Results . . . . .	90
	Approaching based on Reinforcement Learning with- out Prediction . . . . .	90
	Approaching based on Reinforcement Learning and Predictive Sensory System . . . . .	93
	Docking Trials with Different Connections . . . . .	94
6.4	Summary . . . . .	94
<b>7</b>	<b>Pre-symbolic Representation Emerged from Sensorimotor Feed- back</b>	<b>95</b>
7.1	Language Acquisition . . . . .	95
	7.1.1 Pre-symbolic Communication . . . . .	95
	7.1.2 Motivation . . . . .	98
7.2	Horizontal Recurrent Network with Parametric Bias . . . . .	101
	7.2.1 Algorithm . . . . .	102
	Learning Mode . . . . .	103
	Recognition Mode . . . . .	105
	Generation Mode . . . . .	105
7.3	Experiments . . . . .	105
	7.3.1 Preliminary Experiments . . . . .	105
	7.3.2 Pre-symbolic Learning via Interaction . . . . .	111
	Learning . . . . .	111
	Recognition . . . . .	114
	Generation . . . . .	118
	Generalization in Recognition . . . . .	118
	PB Representation with Different Speeds . . . . .	120
7.4	Summary . . . . .	120
<b>8</b>	<b>Discussion and Conclusion</b>	<b>122</b>
8.1	Discussion . . . . .	122
	8.1.1 Hierarchical Action Control . . . . .	122
	8.1.2 Language Acquisition from Sensorimotor Integration . . . . .	124
	8.1.3 Predictive Perception . . . . .	125
	8.1.4 Robotics as Synthetic Methodology and Neuro-robotics . . . . .	127
8.2	Future Work . . . . .	129
	8.2.1 Conceptor Representation and Mirror Neuron System . . . . .	129
	8.2.2 Deep Learning and Predictive Coding . . . . .	132
8.3	Conclusion . . . . .	133
	<b>Bibliography</b>	<b>137</b>
	<b>Declaration of Authorship</b>	<b>161</b>

# List of Figures

2.1	Ventral and Dorsal Streams in Visual System . . . . .	11
2.2	Motor Control Flowchart . . . . .	13
2.3	Layers of Cortical Columns . . . . .	16
2.4	Cortices in Sensorimotor System . . . . .	17
2.5	Multi-dimensional Feedback Pathways . . . . .	18
2.6	Example of Neuronal Feedback from Cognitive Processes: Goalkeeper . . . . .	22
2.7	Binocular Rivalry . . . . .	25
2.8	Schematic of Predictive Coding . . . . .	31
2.9	Schematic of Biased Competition . . . . .	33
3.1	A Perceptron Unit . . . . .	38
3.2	A Multi-layer Perceptron (MLP) . . . . .	39
3.3	Elman RNN Network . . . . .	44
3.4	Jordan RNN Network . . . . .	45
3.5	Unfolding an RNN for BPTT ( $\tau = 3$ ) . . . . .	46
3.6	RNNPB with Elman-like Connections . . . . .	47
3.7	Three Modes of RNNPB . . . . .	49
3.8	Echo State Network (ESN) . . . . .	50
3.9	A LSTM Block . . . . .	52
4.1	Homunculus Organization in Motor and Somatosensory Cortices . . . . .	57
4.2	A Hierarchical Perception-Action Model with Action and Visual Dorsal Stream . . . . .	61
5.1	Anatomy of Two-stream Theory . . . . .	64
5.2	Example of Pooling of Positional Variation . . . . .	67
5.3	Horizontal Product Recurrent Network architecture . . . . .	70
5.4	Partial Sample of Input Data . . . . .	75
5.5	Input of Activation Movement . . . . .	76
5.6	Corresponding Output from Fig. 5.5 . . . . .	76
5.7	Network Activations in Hidden Layers . . . . .	77
5.8	Output Error through Iterations . . . . .	77
6.1	Overall Architecture Combining Sensory Prediction and Sensorimotor Action . . . . .	83
6.2	Sample of the Landmark . . . . .	84
6.3	Shelf Installation / Docking Station . . . . .	89

---

6.4	Training Curves of Two Modules . . . . .	91
6.5	Comparison of Trajectories . . . . .	92
7.1	Diagram of Sensorimotor Integration with the Object Interaction. .	99
7.2	The HoRNNPB Network Architecture . . . . .	100
7.3	Prediction of Three Curves . . . . .	107
7.3	Prediction of Three Curves (cont.) . . . . .	108
7.4	PB Values in Recognition . . . . .	109
7.5	Prediction of Three Untrained Curves . . . . .	110
7.5	Prediction of Three Untrained Curves (cont.) . . . . .	111
7.6	Experimental Scenario . . . . .	112
7.7	Values of PB Units in Two Streams . . . . .	114
7.8	Update of the PB Values in Recognition Mode . . . . .	115
7.8	Generated Values from HoRNNPB . . . . .	117
7.9	Update of the PB Values in Recognition Mode with an Untrained Feature (Circle) . . . . .	119
7.10	PB Values with Faster Speed . . . . .	120
8.1	Hierarchical Perception-Action Model . . . . .	123
8.2	Cognitive Development . . . . .	128
8.3	Neural and Cognitive Model by RNNPB with Deep Structure . . .	133

# List of Tables

2.1	Examples of Functional Modularity in Different Physiology Scales .	19
5.1	Network Parameters (Horizontal Product RNN) . . . . .	76
6.1	Time delay in the Camera-arm Cycle of NAO (in milliseconds) . .	81
6.2	Network Parameters (Predictive Sensorimotor Integration) . . . . .	87
6.3	Docking Trials by Reinforcement Learning with and without Prediction . . . . .	90
6.4	Docking Trials in Different Connections . . . . .	94
7.1	Root Mean Square Error of Two Curve Set Predictions . . . . .	108
7.2	Network Parameters (HoRNNPB) . . . . .	113
7.3	Prediction Error . . . . .	118

# Abbreviations

<b>ARNN</b>	Artificial Recurrent Neural Network
<b>AMG</b>	Amygdala
<b>BPTT</b>	Back-propagation through Time
<b>CACLA</b>	Continuous Actor-Critic Learning Automaton
<b>CNS</b>	Central Nervous System
<b>ESN</b>	Echo State Network
<b>FFA</b>	Fusiform Face Area
<b>fMRI</b>	Functional Magnetic Resonance Imagine
<b>HMM</b>	Hidden Markov Model
<b>HoRNNPB</b>	Horizontal Recurrent Network with Parametric Bias
<b>ICA</b>	Independent Component Analysis
<b>IFG</b>	Inferior Frontal Gyrus
<b>IT</b>	Inferotemporal Area
<b>LGN</b>	Lateral Geniculate Nucleus
<b>LSTM</b>	Long Short Term Memory
<b>M Cell</b>	Magnocellular cell
<b>MLP</b>	Multi-Layer Perceptron
<b>MNS</b>	Mirror Neuron System
<b>MT</b>	Medial Temporal Lobe
<b>NMF</b>	Non-negative Matrix Factorization
<b>OFC</b>	Orbitofrontal Cortex
<b>P Cell</b>	Parvocellular Cell
<b>PAM</b>	Perception-Action Model
<b>PFC</b>	Pre-frontal Cortex



---

<b>POMDP</b>	Partially Observable Markov Decision Process
<b>PPC</b>	Posterior Parietal Cortex
<b>PMA</b>	Premotor Cortex
<b>RBC</b>	Recognition-by-components
<b>RBF</b>	Radial Basis Function
<b>RBM</b>	Restricted Boltzmann Machine
<b>RNNPB</b>	Recurrent Neural Network with Parametric Bias
<b>SC</b>	Superior Colliculus
<b>SIFT</b>	Scale-invariant Feature Transformation
<b>SMA</b>	Supplementary Motor Area
<b>SMC</b>	Sensorimotor Contingency
<b>SRN</b>	Simple Recurrent Network
<b>SURF</b>	Speeded Up Robust Features
<b>SSE</b>	Summed Squared Error
<b>STS</b>	Superior Temporal Sulcus
<b>SVM</b>	Support Vector Machine
<b>V1</b>	Primary Visual Area

*To my parents and my grandmother: I will be eternally thankful for your support and encouragement. Without you, this work would not have been possible.*

# Chapter 1

## Introduction

The main aim of this thesis is to implement feedback neural connections using artificial recurrent neural networks (ARNN) on the sensorimotor integration of cognitive robotics systems. This is based on the cognitive finding that neural feedback transmits signals from high-level cognitive functions to the lower-level neuron activities in the sensorimotor systems. These processes also account for several cognitive functions and phenomena. This thesis proposes that the feedback information is partially originated from cognitive processes, such as (pre-)symbolic representations. These feedback mechanisms in sensorimotor processes are implemented with recurrent neural models in this thesis. Furthermore, such recurrent connections-based neural models will be examined in cognitive robotic systems.

### 1.1 Cognitive Robotics

Since as early as 1920 when the Czech writer Čapek invented the word ‘robot’, different kinds of robotic systems have been designed and deployed in various fields. As a branch of technology, a robotic system is usually built to provide solutions for a single or a set of task(s) or problem(s) with configurations of mechanical systems, electrical systems and control systems. Due to the fact that robots are built with tireless and (mostly) faultless manipulators, they are suitable to aid us with repetitive, dangerous and demanding situations. Therefore, these systems have been extensively used in the fields of industrial manipulation, space exploration, etc. For instance, a number of industrial robots (mostly articulated robots) have been deployed in the auto-mobile industry. The ratio of robots to the employed human workers has increased to ten to one [Gates, 2007]. The semi-autonomous

robots ‘Robonaut’ have been sent to space stations to handle different types of tasks in space with their dexterous manipulation skills [Lovchik & Diftler, 1999, Ambrose et al., 2000, Bluethmann et al., 2003].

Nevertheless, the industry robots and the ‘Robonaut’ robots are not fully autonomous robotic systems. Instead, there are only a small number of fully autonomous service robotic systems which are used in the real world and share the same working space with humans. The reason why only such a small number of robotic systems is employed to accompany a human being in daily life is that such a robotic system needs to be accurate enough to understand the situation by perception, to be adaptive enough to deal with the changing environment with noisy sensors and to be knowledgeable enough to communicate with humans. These problems (among others) are still yet to be fully solved by engineers and researchers.

One solution to solve the problem regarding full autonomy is to develop a robot that possesses its own cognitive capabilities. Thus, the topics of cognitive systems and robotic systems are correlated and overlapping with the subject ‘cognitive robotics’. Inspired by multiple disciplines such as cognitive science, neuroscience and psychology, research in cognitive robotics mainly concerns how to develop cognitive ability by designing architectures and algorithms in hardware and software for robot systems so that they can execute intelligent behaviours in terms of human-like perception and motor action as well as high-level cognition. Therefore, these systems are able to be operated in a dynamic environment driven to accomplish one or several complex goal(s).

Also, designing a cognitive robotic system is different from designing architectures and algorithms to merely provide an ability for a machine to plan, to reason and to deliberate a solution according to symbolic rules, although this artificially intelligent method has been proven to be successful in applications that can be reduced as a formulation of symbol manipulation, such as playing chess; for instance, the renowned supercomputer ‘Deep Blue’ won a chess match against the human world champion [Schaeffer & Plaat, 1997]. Such machines still belong to disembodied devices, which can only passively receive and process symbolic data as a calculator does. Despite it has intelligent abilities based on learning, it cannot carry out a large diversity of tasks because it does not really learn from information perceived from sensors, and does not own a body for interacting with the environment. Intelligence and understanding are not only symbol manipulation. Therefore, cognitive robotics is much related to provide an ability of ‘intelligent’

thinking to a robotic system. The study of cognitive robotics is not only about learning to ‘think’. Derived but different from the conventional ‘artificial intelligence’, a ‘cognitive robot’ that possesses intelligence should also own ‘a group of operations of the mind by which reasoning is performed, to give (an embodied) expression to them in the symbolical language of calculus’ [Boole, 1854]. So, the embodied knowledge (the ‘mind’ proposed by Boole) coming from different cognitive systems differs when they have various configurations and constraints: birds (or flying robots) that float in a six-dimensional (i.e. three-dimensional position and three-dimensional rotation) free space (i.e. sky) have a different perceptual world than ants (or vacuum-cleaning robots) which are restricted to planar surfaces; human beings (or humanoid robots that have legs and arms) need a more complex control scheme for learning dexterous hand and finger movements.

Thus, the behaviours of the embodied ‘acting’, i.e. taking into account the physical body for learning, are crucial for cognitive systems. In line with the learning of ‘acting’ in human and other biological cognitive systems, it actually takes up the majority of the cerebral cortex in the brain to accomplish such related sensorimotor tasks<sup>1</sup>. This can also answer Moravec’s paradox<sup>2</sup>: the more abstract knowledge obtained from sensorimotor integration can consequently be utilized as a description of the machine’s reasoning. To sum up, to realise a full ‘artificial cognitive system’, the software and hardware design of a cognitive system should first take into account ‘the embodiment of intelligence’, which advocates that there is no clear distinction between the representation of thinking and the way of perceiving and acting. Perceiving and acting are major ways to acquire knowledge of thinking; the cognition ability is acquired during interaction with the environment and presented with the body. Besides, there are no differences in terms of basic cognitive mechanisms between biological agents and artificial agents, although they have different manifestations.

In terms of artificial systems, some of the researchers focused on building artificial machines to realize embodied behaviour-based intelligence by demonstrating low level sensorimotor behaviours (e.g. navigation [Cordeschi, 2002, Holland, 2003, Webb, 2002]). But most of them have merely done simple formulation of embodiment, as they ignored either higher-level cognition or its strong link grounded in bodily activity and experience.

---

<sup>1</sup>The human temporal lobe, occipital lobe and parietal lobe, where various sensory cortices and motor cortices locate, occupy approximately 60-65 % of the cerebral hemispheres.

<sup>2</sup>According to Moravec’s paradox, it is contrary to traditional assumptions that the high-level reasoning requires very little computation, but low-level sensorimotor skills require much more computational resources.

Alternatively, some researchers started to identify possible cognitive mechanisms in biological systems that endow intelligence to emerge from dynamic interaction with the environment through sensorimotor experience. These findings are then transferred by building artificial cognitive systems. For instance, researchers investigate the developmental mechanisms of infants and program them in robotic systems (see also [Asada et al., 2001, Weng et al., 2001]). These systems deliberate about the surrounding by understanding and interaction; the acquisition of ‘cognition’, especially the way to act, is also done with the accumulating process of sensorimotor skills instead of being programmed by the human designer. Furthermore, high-level cognition, such as symbolic representation for language acquisition and reasoning, is also obtained by the interaction process which involves the physical body, its sensorimotor process and the environment [R. A. Wilson & Foglia, 2011]. So, cognitive learning is not only a knowledge acquisition process mediated by the physical body, but also an abstraction process grounding intelligence and skill development at the same time. In terms of artificial systems, this is achieved by constituting a structural coupling between mind, neural structures and the physical body, by a small set of pre-programmed learning rules. By this mean, the system acquires a representation from sensorimotor knowledge: it becomes capable of interaction using more complex verbal and non-verbal expressions.

## 1.2 Modelling Feedback Pathways on Cognitive Robotic Systems

In the context of sensorimotor integration, it is straight-forward to regard that the perceptual world is directly obtained by a series of sensory-driven information flows that come from various kinds of perceptual receptors (a.k.a. the bottom-up influences); they directly represent the physical characteristics of the stimulus-driven perception, which indicates that the percept is a *true* representation of reality. For instance, in the visual system, the electrical signals from the retina carrying visual information are sent across a series of cortical areas until they reach a high-level neural representation in which the visual scene is understood. Therefore, the higher level represents more abstract information; this continues until the highest cognitive representations, such as language acquisition, decision making and reasoning, are formed on the highest levels.

However, our percept is also affected by the information about previous experience, concerning how the percept/action might appear [Gilbert & Li, 2013], which is transmitted in the feedback neural pathways. These schemata are affected by the existing attention, expectation, perceptual tasks, working memories and proposed motor commands. Existing in the feedback neural pathways, the complete schema is constituted by the information provided from the top-down and lateral connections. Physiologically, these influences can be represented as a type of neural transmission that originates from the high-level cortical areas and exerts neuronal influences to the low-level neuronal activities, such as chemical and electrical transmission. It consequently affects the existing response patterns of neuronal population in order to become better suited to the (anticipated) environment. This is especially useful, as its existence increases processing speed and accuracy, and reduces the bandwidth of sensorimotor processing by interpreting a percept or motor action, filtering sensory information or predicting the upcoming sensory information.

In this thesis, we specifically focus on designing embodied neural networks to model the feedback pathways on the *sensorimotor* integration of cognitive robotic systems. This is inspired by biological cognitive systems, whose sensorimotor control is ubiquitously affected or modulated by the feedback pathways. Generally speaking, the feedback pathways link the cortical areas, which constitute various levels of representation in a hierarchical sensorimotor system. Generally, the feedback pathways include:

- the top-down projections from high-level to low-level representation;
- the lateral connections within the same level.

In agreement with our embodied cognition theory, the emergence of such feedback information is also accomplished by learning knowledge in an embodied way through sensorimotor integration, which refers to a dynamic process involving action, perception and interaction.

### 1.3 Research Questions and Contributions

In short, we will investigate the feedback pathways on sensorimotor functions of cognitive robotic systems, modelled by different kinds of recurrent neural networks in this thesis. We will answer the following questions:

- What kind of information do the feedback pathways deliver?
- What are the mathematical descriptions of these feedback signals? What is the model to describe them under these descriptions?
- From an engineering perspective, is it possible that the feedback pathways also give rise to some of the sensorimotor functions in artificial cognitive systems as they do in biological systems? Do the artificial feedback pathways also facilitate the sensorimotor learning and adaptation in a dynamic environment for a cognitive robotic system?

## 1.4 Structure of the Thesis

The remainder of this thesis is organised as follows:

Chap. 2 sees an overview of the feedback pathways in the context of the neuroscience and cognitive science. The concepts of functional modularity and hierarchical organization in the sensorimotor areas also provide the theoretical foundation of the existence of the feedback pathways. Additionally, cognitive evidence is introduced to support the existence of the feedback pathways. Two theories (predictive coding and biased competition) are introduced which attempt to explain how the information from feedback pathways affect the sensory-driven bottom-up neural stimuli on the cortical functions.

In Chap. 3, we focus on the implementation of feedback pathways with an emphasis on their feedback signals. This can be implemented by recurrent artificial neural networks on situated agents. In this chapter, the fundamental structures and algorithms of recurrent neural networks are introduced in the context of computer science. We also present a few variants of recurrent network models.

In Chap. 4, an overall cognitive architecture of this thesis is proposed, which also gives a general framework of this thesis. This framework is based on the perception-action model, where the hierarchical perception and action share the common coding representation. This representation also plays a role in exerting feedback pathways in both perception and action. The following chapters describe different parts of this architecture.

From Chap. 5 to Chap. 7, we begin to propose the neural network models on embodied systems. Firstly, in Chap. 5, we concentrate on the modelling of different temporal-encoding requirements in the two-stream theory within the primary visual cortex. These requirements, specifically, mean that the fast-changing units



are needed for encoding dorsal information, while the slow-changing units are needed for encoding ventral information. We claim the maintenance of such units is contributed by the feedback pathways on the neural structure level in the visual system. To model this, two homogeneous recurrent connections are used to maintain the encoding of object features and movements. These two recurrent networks are linked by a horizontal product where the information about the features and movements of an object are able to be separated successfully in the two hidden layers.

Since the dorsal pathway also exerts the motor-relevant information which allows smooth sensorimotor integration, the filtering and predictive functions of the feedback information ensure a stable and rapid motor action. Therefore, we examine the predictive function on the sensorimotor system based on feedback pathway models in a robotic system in Chap. 6.

In Chap. 7, a hypothesis that the pre-symbolic representation emerges from sensorimotor integration is modelled and examined by robot-object-interaction. This representation relates to a high-level language acquisition and conversely it operates as one of the sources of the feedback influences. Here, the feedback signal is across neuronal and cognitive levels as well as motor and sensory cortices. The model is realised by a novel horizontal product recurrent neural network model featuring a recurrent network with parametric biases.

Finally, a discussion and conclusions are given in Chap. 8.

## Chapter 2

# Feedback Pathways in a Hierarchical Modularity Brain

In this chapter, the feedback pathways of the sensorimotor system will be reviewed. The sections introduce the sensorimotor system from a lower (neurobiology) to a higher (cognition) level: we first introduce the anatomical organization of the sensorimotor system of biological cognitive systems. The hierarchically organized cortices are extensively linked by reciprocal connections of the sensorimotor system, which is the prerequisite of the existence of the feedback pathways. These feedback pathways convey prior knowledge guided perception/action into neuronal and even lower levels. Together with the sensory-driven bottom-up influences, the information from the feedback pathways maintain asymmetric information flows on the artificial/biological sensorimotor system. In addition, two hypotheses about how the feedback pathways influence low-level neural activities are introduced.

### 2.1 Sensorimotor Integration

Perception and action hold a synergistic relationship, in which these two parts communicate and coordinate with each other. In particular, the sensory representation as well as sensory awareness emerges from the changes of the perceptual world in the sensory input, which results from the active execution of certain sensorimotor skills through cognitive processes, rather than an internal representation merely from sensory signals. Therefore, perception can simply be conceived as a process of probing the external world by action (e.g. moving the arm and touching, and perceive what has been changed). This relationship can be depicted as

sensorimotor regularities, according to the framework of sensorimotor contingency (SMC) [O'Regan & Noë, 2001]. This law-like relationship between perception and action rejects the traditional theory that perception is fully composed by cognitive processes in the brain. Also, this law allows perception to be acquired by engaging a full set of skillful action and establishes cognitive processes from sensorimotor interaction.

Within the SMC framework, Prinz [1984, 1992, 2003] proposed that perception and action systems are represented as a *common-coding* mechanism, stating that perception and action share the same representation which reflects the perceptual events that actions produce, rather than that there are cognitive processes in-between. This means that perceiving an action triggers the same representations as the perception system does when it receives sensory input. These representations are called 'common-coding'. These shared representations do not encode explicit actions, but rather do they encode the information which is perceived as a consequence of the corresponding actions.

Although there is still not a convincing conclusion in how perception and motor actions are related, the common-coding theory may be one answer to the key part in most, if not all aspects of cognition. As the common representation may be able to be scaled up to represent a perceptual symbol, this theory may be able to conclude a unified perception representation that comes from the exploration of the environment and from the immediate effect of the sensorimotor contingencies. This representation can be further exploited for one's planning, reasoning, and speech behaviours. From this theory, our prior experience and knowledge come from both the sensory input channels (sensory cortices) and the motor action channels (motor cortices). In the next section, we will review the anatomy of the sensorimotor system.

## 2.2 Sensorimotor System

### 2.2.1 Sensory and Motor Cortices

**Sensory Cortices** The sensory cortices mainly include the visual cortex, the auditory cortex and the somatosensory cortex. These brain areas, and the receptors constitute a hierarchical system for sensory systems. We will briefly introduce them from low to high levels.

The lateral geniculate nucleus (LGN) is the first processor of visual information as it receives its information from the retina. Receiving the majority of input from LGN, the primary visual area (V1) is the first stage of the cortex that processes visual information in all visual areas. It forms a mapping of the whole visual field in a topographical way: anatomical locations in V1 maintain the similar spatial arrangement of the visual field. Inside V1, neurons with different receptive fields are functionally considered to be local feature detectors for preferred orientation in the visual stimuli. The neural encoding of orientation perception, subject movement as well as feature recognition usually begins on this level of processing. So, the V1 area is the first brain area that receives the electrical signals from the eyes and transmits it for further perception and motor action. After V1, visual information is transmitted to the surrounding visual areas such as V2 and the associative visual areas: V3, V4, V5 (or MT), etc. Although these areas are not as comprehensively studied as V1, the main functions of these areas have been identified:

- **Visual area V2**, also called prestriate cortex, is the second major area in the visual cortex. It receives strong connections from V1, and is mostly considered as a mapping of V1, but its neurons possess more complex receptive fields. Also strong connections are sent from V2 to V3, V4, and V5.
- **Visual area V3** receives the majority of inputs from V2, and projects to the area MT and V4. Part of the V3 normally contains a representation of the dynamic shape of visual stimuli.
- **Visual area V4** receives feed-forward connections which are from V1 via V2. Different from the complete mapping of V1 and V2 in the visual field, the receptive field of V4 is only sensitive to both colour and orientation.
- **The inferotemporal area (IT)** is located anterior to V4. Its neurons also activate to a wide range of colours and an object feature of intermediate complexity.
- **Visual area V5/The medial temporal lobe (MT)** is part of the extrastriate visual cortex. With connections to area V3, it is considered to play an essential role in the perception of motion.

From V1 via V2 to other associative cortices, the receptive fields have a relatively increasing complexity. Also, these areas appear to form two major cortical systems for processing visual information: a ventral visual stream begins with V1, into

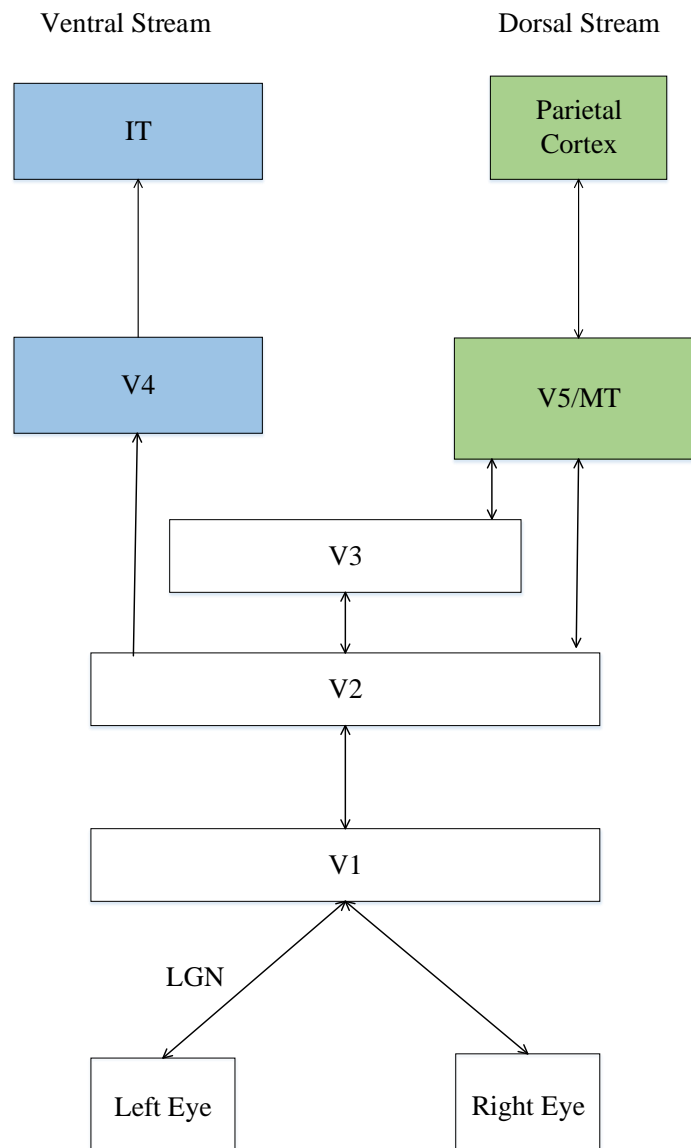


FIGURE 2.1: Ventral and Dorsal Streams in Visual System  
Areas in ventral streams are denoted in blue, while areas in dorsal streams are denoted in green. Common areas are in white. We can see a hierarchical organization in both layers from the eyes to higher cortices for visual information understanding.

V2 and V4, and terminates at the IT cortex. This stream processes the local features of a visual stimuli, which is utilized to process form recognition as well as object representation and to keep a long-term memory. All this information is used to identify, recognize, and remember objects. The other pathway is the dorsal stream, which begins with V1, goes through V2 and V5 until the posterior parietal cortex. The dorsal stream is essential to provide spatial information of stimuli. This information further relates to motor action, location of the object, and control of the saccades.

From the above-mentioned connections in the visual system (Fig. 2.1), we can conclude that various cortical areas of the visual system are connecting and functioning in a hierarchical manner related to their receptive fields' complexities and abstractness: from low- to high-level, the sensory representation in the neural activities are becoming more and more abstract.

**Motor System** A complete motor system includes the motor cortex, the central nervous system and the manipulator.

Particularly, the motor cortex can be divided into several parts:

- **The primary motor cortex (M1)** controls the execution of movement by neural impulses through the central nervous system. It is located at the frontal lobe, forming a somatotopic representation of different parts of the body.
- **The premotor cortex (PMA)** controls some aspects of motor action, such as torso muscles of the body. Also, it is involved in the preparation for movement and the planning of a movement. The so-called mirror neurons are also located in area F5 of the premotor cortex.
- **The supplementary motor area (SMA)** controls body movement, which has a direct connection to the spinal cord.
- **The posterior parietal cortex (PPC)** is responsible for multi-sensory information to motor commands.

The control process of motor action is simply following a top-down control scheme from the cognition level to the execution of a voluntary action. This control scheme is enacted hierarchically (Fig. 2.2).

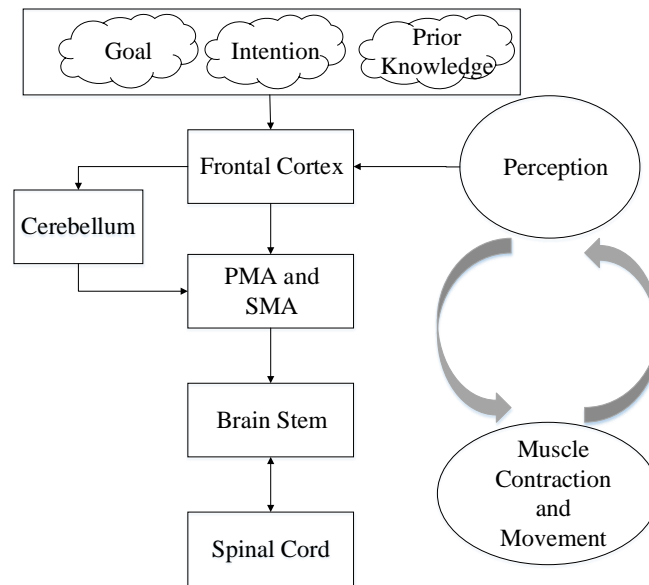


FIGURE 2.2: Motor Control Flowchart

- Intention Description** First, as any voluntary movement is jointly determined, mediated or affected by the motor cortex and numerous other neural systems, cognitive processes determine motor strategies according to goals, objectives or intentions of the intended movement at the topmost level. This is anatomically involved in the pre-frontal cortex (PFC).
- Vision-for-action** The parietal cortex projects spatial perception to the frontal cortex. The function of the frontal cortex involves the analysis of the position of body, from which further motor actions can be accordingly determined. The basal ganglia are also involved in this process.
- Pre-structured Actions** The secondary motor areas (PMA and SMA), together with the cerebellum, further augment the goal-directed pre-structured motor programs into motor synergies, which calculate the precise force produced by different muscles in agreement with the principle of redundancy.
- Muscles Contractions** According to the outcome of secondary motor areas, the primary motor cortex, the brain stem and the spinal cord further specify and generate the contractions of all the muscles needed for the motor actions.

The voluntary motor action is determined by the cognitive processes with consideration of the sensory information obtained from the perceptual world, the

current state of the body and the goal (objective or intention). The whole process is executed in a top-down way: the cognitive process augments the goal into a set of motor primitives, which are physically executed by muscles. More specifically, this execution process requires a top-down control scheme from the motor cortices and the central nervous system to the musculoskeletal system.

This section briefly introduced the anatomy of the visual and motor cortices. To summarize, the cortices are organized in a hierarchical way. Among these cortices, besides of the main stream signals (i.e. sensory-driven signal in visual cortex, top-down signal in motor cortex), there are also feedback signals (c.f. Fig. 2.1 and Fig. 2.2). We will address this problem in the following section.

### 2.3 Functional Modularity of Sensorimotor System

The human body is physiologically constituted by clustering on various levels: the neuronal cells are formed by clusters of molecules, while the neurons constitute various structures of the brain. Also, accomplishment of a cognitive function involves a series of structures. To sum up, these physiological parts which own similar or related functions usually interact, cluster and function as a whole and form another physiological network on a higher level. In this way these parts become another unique level in the physiological organization. That is the way neurons, neural structure and cognitive functions emerge. Literally, we call the clustering phenomenon to form cells, neurons, brain issues and networks with similar/related functions *functional modularity*.

Although this is a common phenomenon to form various parts of the physiological body, we only depict a few examples in the human brain in this section. On the neuronal level, for instance, neurons are segregated into different layers within an individual cortical column (also known as the hyper-column or cortical module) with different kinds of connectivities in the cerebral cortex. These modules carry out specific cellular functions, such as signal transmission, by the interaction of cells within the same modularity. Those cells own various physiological characteristics, all of which are linked by the connections. These connections endow cells with different characteristics to interact efficiently and to perform cellular functions jointly [Mountcastle, 1957].

Within each of these cortical columns, the neurons have almost identical receptive fields, which implies that they have similar firing activities while a similar stimuli is presented. In the V1 area, specifically, cells with the same eye preference as well



as the same orientation of line stimuli are grouped into the same cortical columns. In other words, each column typically responds to a sensory stimulus representing a certain feature of sound, vision or other sensory modalities.

Such columnar organisation can be mostly found in the sensory and motor areas of the neocortex, in which the cells in neighbouring cortical columns have similar functions (e.g. orientation selection, eye preference), so the cortical columns are considered to be one of the most basic repeating functional units on the neuronal level [Ng, 2009]. This has been recorded from micro-electrode mapping experiments, metabolic studies and nerve regeneration experiments (e.g. [Kaas, 1987, Mountcastle, 1957, 1997, O. Favorov & Whitsel, 1988, O. V. Favorov & Diamond, 1990, Tommerdahl et al., 1993]).

As a basic functional unit, each cellular module contains afferent excitatory and inhibitory connections from and to other modules as well as intra-cortical connections. This is realised by cells with different functions within the same module. For instance, there are six different layers of different neurons within the cortical column. Each of them has distinct functions (Fig. 2.3): the neocortical neurons and pyramidal cells retain excitatory connections which are grouped into separate bundles with dendritic cells at their centres, while in another layer the basket cells and stellate cells form local inhibitory connections to exert strong intra-modular lateral connections to other modules [M. E. Newman, 2004, 2006]. The six layers of cells constitute a module as an informational encapsulation and limited central accessibility, except with some kind of input and output channels. On the other hand, due to the encapsulation, these modules require information flows to communicate out of the module and form a larger system via inter-modular connections such as axons and dendrites.

Similar to the formulation of the cellular level, the quasi-independent neural modules with similar functions are integrated within themselves. They exhibit some degrees of interdependency among other modules and form another level of modules, from which they assemble the basic networks on a higher-level: structural level.

A structure module in the human brain is one tissue to accomplish several cognitive functions in the human body, which usually is anatomically referred to as the Brodmann areas (such as sensory areas of visual cortex, somatosensory cortex, motor cortex, premotor cortex, etc.) in the cerebral cortex. A bit different from the modularity on the cellular level, the organization of structural areas suggests that both functional segregation and functional integration are happening

simultaneously in the brain to some extent (for a detailed review see [Bullmore & Sporns, 2009]).

For instance, the functional modularity in brain structures can be indicated during detection of neural activation in the fusiform face area (FFA) and other structural areas, in response to different facial recognition tasks [Sergent et al., 1992, Kanwisher et al., 1997]. These activations of neural coupling are dynamically changing according to various types of stimuli and context [Ishai, 2008] depending on stimulus and tasks. This suggests that the mapping between brain structures and functions are not static but a dynamic changing process according to the conscious/subconscious cognitive processes. Furthermore, the modularity results in a prompt interaction between neurons within the same structure if they are involved in a similar task. Fig. 2.4 shows the main brain structures which are involved in visual and auditory perception. We can see the visual and auditory cortical areas are located closer to their linking areas in the occipital lobe and the temporal

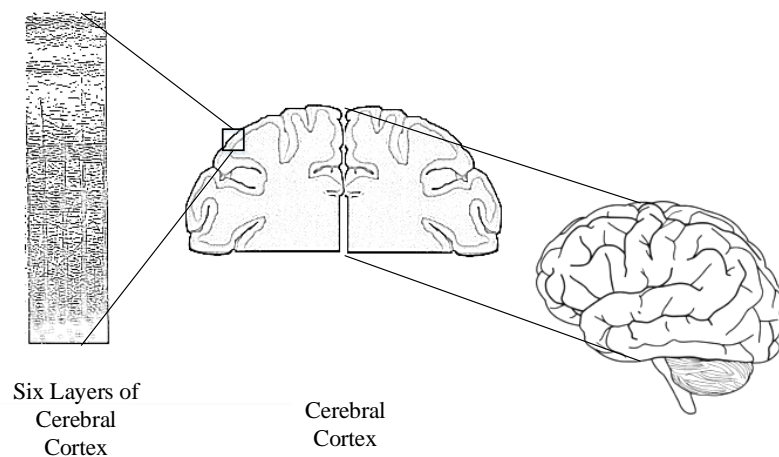


FIGURE 2.3: Layers of Cortical Columns

The cerebral cortex is the outermost of the mammalian brain. It constitutes of up to six horizontal layers. These layers have a different composition in terms of neurons and connectivity<sup>1</sup>.

<sup>1</sup>This image is a derivative work based on images from [http://commons.wikimedia.org/wiki/File:Human\\_cerebral\\_cortex.png](http://commons.wikimedia.org/wiki/File:Human_cerebral_cortex.png), [http://commons.wikimedia.org/wiki/File:Cajal\\_cortex\\_drawings.png](http://commons.wikimedia.org/wiki/File:Cajal_cortex_drawings.png) and <http://pixabay.com/en/brain-human-anatomy-body-155655/>. All of the images are licensed under the Public Domain license.

Systems	Neuron	Structure	Cognition
Visual System	[Tootell et al., 1998], [Fujita et al., 1992]	[Felleman & Van Essen, 1991], [Ungerleider & Pessoa, 2008]	[Ward, 2008], [Zeki & Bartels, 1998]
Motor System	[Eisenberg et al., 2010], [Takei et al., 2001]	[Rizzolatti et al., 1988], [Alexander et al., 1986]	[Kuppuswamy & Harris, 2013], [Carruthers, 2002]

TABLE 2.1: Examples of Functional Modularity in Different Physiology Scales

lobe, respectively.

This could be explained by the fact that the organization of modularity saves a lot of time for neural information transmission, although the speed of neural impulse can reach  $100 \text{ m/s}^2$ . Also, it is an advantage in evolution to reduce energy required for information transmission. Research has been conducted to investigate how neuronal modularities emerge in complex networks (e.g. [Kashtan & Alon, 2005, Kashtan et al., 2007, Chen et al., 2006, Clune et al., 2013]). It is widely believed that the modular structure of complex brain networks plays a critical role in their functionality to make the transmission of information within cognitive functions more efficient by shortening the nerve length for a fast information transmission, as the brain structures have various information routes in different contexts.

Besides the functional modularity examples above, there are more examples of multiple levels of functional modularity, which are shown in Tab. 2.1. We can see that modularity is a pervasive phenomenon in the whole physiological body on various levels. From these examples of the brain, we can conclude that a cluster on one physiological level allows them to process one function at a time without changing too much information within one module. The whole brain system is composed of specialised-function cognitive modules, which are localised from the low-level peripheral neurons to high-level brain structures. Because of the dynamic connectivity between the structures in the facial recognition tasks, modularity also requires a series of dynamic inter-modular interactions. Some of them are transmitted via feedback pathways.

## 2.4 Multi-dimensional Hierarchical Modularity

Since we have introduced the functional modularity of the brain (especially in the cellular and structural levels), the question might arise how do these functionally

<sup>2</sup>For instance, some myelinated neurons conducting at speeds up to  $120 \text{ m/s}$  ( $432 \text{ km/h}$ ).

coherent modules combine into larger, less cohesive subsystems as well as the complete cognitive functional sensorimotor system?

Here, we conclude that both within and across various physiological scales, there is a hierarchical organisation to control the body, to carry out metabolism and to restructure its own physiological function in each dimension. This conclusion is based on the anatomy and physiology studies on the human body and brain. As shown in Fig. 2.5, the feedback pathways have two directions along two physiological dimensions.

We call the first dimension hierarchical organization within the same level (scale) of physiology. Along this dimension, modules are with the similar physiological composition (i.e. the same physiological scale). The studies of this dimension can be extensively conducted on neural structural level, such as in the visual systems (e.g. [Gilbert & Li, 2013]) and the auditory system (e.g. [Polley et al., 2006]). For instance, the feedback pathways on the neural structure level can be found in electrophysiological studies of the visual system, especially the visual information extraction through the ‘what’ and ‘where’ pathways. In general, the visual

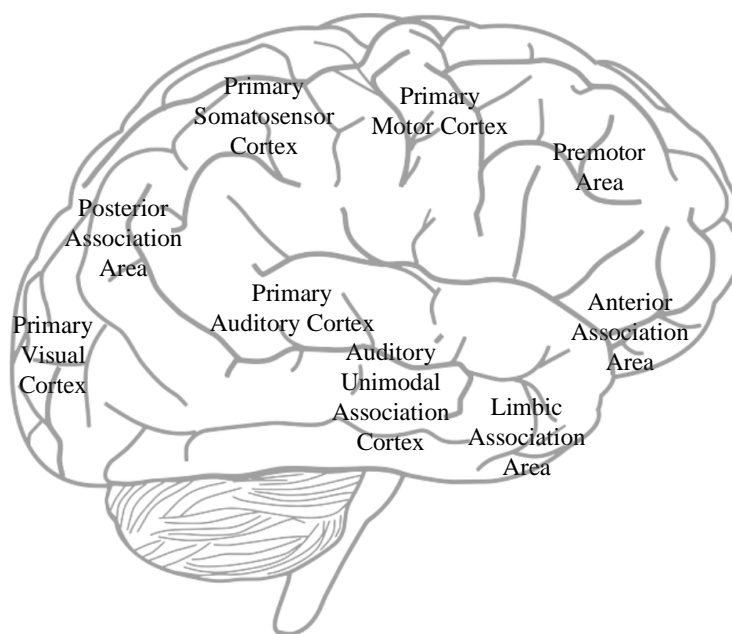


FIGURE 2.4: Cortices in Sensorimotor System

Nearby cortices are grouped and linked in a hierarchical way: from primary sensory areas, association sensory areas to higher-order areas.

(Redrawn from [Saper et al., 2000]<sup>3</sup>)

<sup>3</sup>This image is a derivative work based on image <http://pixabay.com/en/brain-human-anatomy-body-155655/>, which is licensed under the Public Domain license.

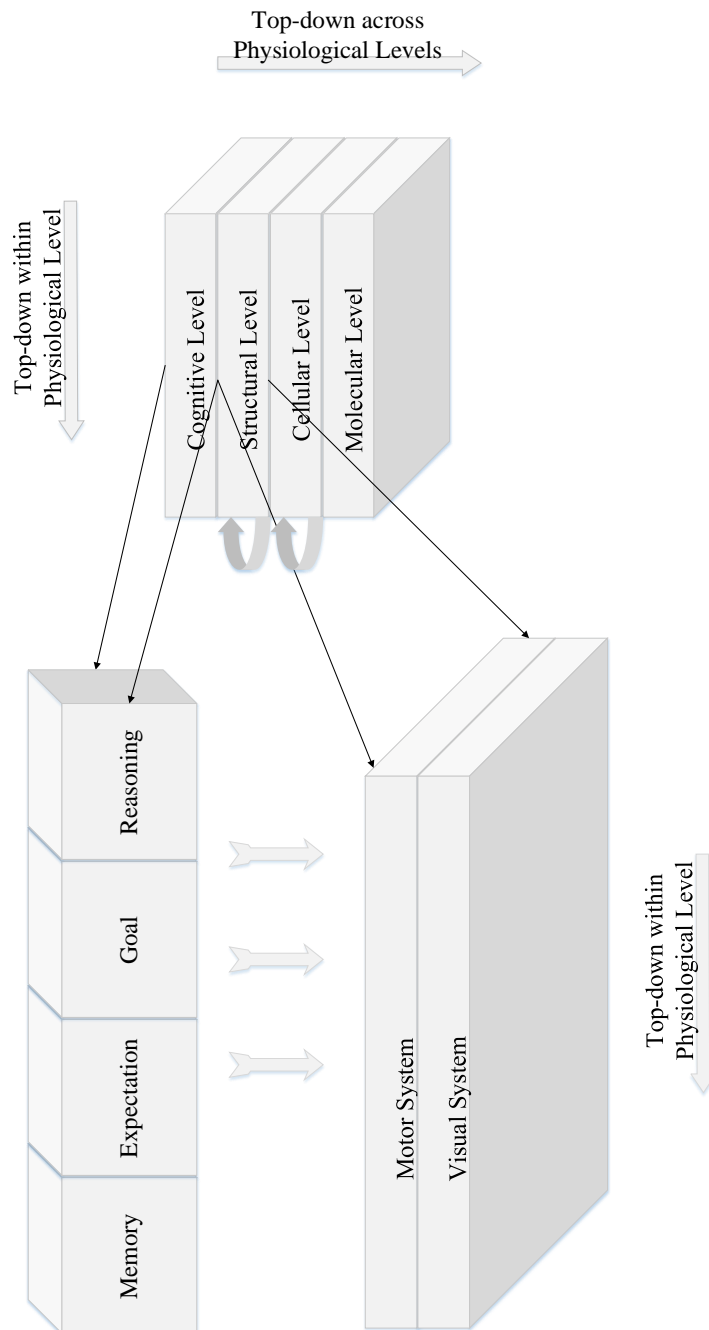


FIGURE 2.5: Multi-dimensional Feedback Pathways

The feedback pathways (including top-down and lateral connections) spread across various physiological levels and within the same physiological level, while we can also find numerous lateral connections within neuronal structural and cellular levels.

processing starts from early to late vision system, which is physiologically corresponding roughly to the pathway from LGN to the visual cortex hierarchically. Neurons on each level of brain area in this hierarchy extract a more abstract information, forming non-random structural features in the ventral and dorsal streams [Van Essen et al., 1992], which correspond approximately neural activities in the spatially posterior-anterior and to the dorsal-ventral brain areas. In addition to the sensory-driven neural signals which are transmitted successively from magnocellular retinal cell (M-cells) to V1, V2 and MT areas, there also exist feedback signals from MT to V3 and V2 (Fig. 2.1). Neuron population in the MT is sensitive to a moving stimulus (i.e. a spatial information), which considered to be accounted for motor action such as saccades for fixation, rapid eye movement (see also [Milner et al., 2006]). The bottom-up influences on dorsal stream reaches part of the frontal cortex, which is associated with high-level control of cognition. On the other hand, the position change information of the visual stimuli in the MT area also mediates low-level perception, resulting in visual illusions, such as ‘flash-lag effect’, which is a visual illusion wherein a flash or a moving object that perceived is displaced from the actual position (usually perceived as motion extrapolation [Nijhawan, 1994]).

At the other dimension, a new *decomposable* module can be formed from lower physiological level modules in a self-organizing way by multiple, sparsely interconnections. In other words, modules on one level form a specialised-function system (i.e. network) in some extents, which becomes another module on a higher physiological level. In this way, modules on various levels are segregated and nested in a hierarchical way: the macro-molecular networks are composed from molecules, the macro-molecular networks constitute neuron circuits, and neuron circuits constitute neuronal networks, which further build the whole brain system of neuronal networks. The whole system is assembled from various neuronal structures. They operate as an organized computational system in which a kind of high-level processing with multiple ‘syndromic response’ functions<sup>4</sup> is formed. That is, it becomes a physiological fact that various cells combine to form tissues which then organize into larger units called organs. On the highest level of the hierarchical system, there are other cognitive processes controlling the neural processes via attention and consciousness. Among those levels, the information transmission from high-level to low-level becomes one of the top-down influences, a kind of feedback signals. Examples can also be found in the visual system: there

---

<sup>4</sup>Syndromic response means that the response of the neuronal system is tuned in a form which should be described as a multi-facet function among various neurons.

are top-down influences from frontal areas which modulate the receptive fields of V1 by predicting the possible position of the flanking stimulus to the relative centre. It implies that neural representations of attention in frontal regions are at the top of the hierarchy, voluntarily assisting a prompt spatial processing in area V1 [Ozaki, 2011]. There are still quite a few feedback pathways controlling the physiological function from the brain, such as cellular functions controlled by the autonomic nervous system and motor movement controlled by the central nervous system. These belong to the cross-level dimension feedback pathways in Fig. 2.5. Although there are various kinds of feedback pathways in physiology, in the following text, we focus only on the feedback pathways in the brain which mainly happen within the neural structure level and cognitive level. Such feedback pathways are considered to represent a kind of subjective experience depending not only on sensory information from the environment but also on parts of cognitive processes, such as our prior knowledge or expectations (Fig. 2.6).

Thus, this suggests that the feedback pathways can spread among multiple physiology scales. The high-level module of the hierarchy generalises the low-level statistics, which includes structural functions, and in turn it controls some parts of the neural activities, cell activities or molecular activity. For instance, the neurotransmitter dopamine modulation (such as arousal) are affected by emotion, which is why body states can express emotion to some extents, e.g. the fear emotion causes trembling hands quickly. On the neural structure level of perception, the encoding of high-level neuron populations generalize the low-level stimuli, so the maintenance of such neural activity predicts part of the low-level sensory driven inputs, according to the previous neural activity from motor action and contextual stimuli as a prior. On the neural structure level of the motor action, the higher-level control includes, but is not limited to, the upcoming motor action command according to the perception and action information as well as the goal reinforcement, etc. (Fig. 2.6). Generally speaking, we can further represent these feedback pathways on the neural and cognitive level in a Bayesian framework, which will be stated in next section. This Bayesian inference on each neural structure level should take into account the stimuli from the low-level sensory input which reflects the structure of the perceptual world around us, with consideration of the other kinds of prior.

In this section, we have concluded that the physical parts constitute the complete cognitive system which can be considered as function modules. Particularly, modules are systematically organized as layers in a *hierarchical* manner in the brain,

as well as in part of the biological body (cognitive, neural structure and neuronal, etc.).

## 2.5 Neural Feedback Connections

In this section, we focus our discussion on how the feedback connections affect the neural structure level. Despite of the fact that the reciprocal connections dominate the neural transmission in the brain (e.g. [Felleman & Van Essen, 1991, Coogan & Burkhalter, 1993]) as it is more efficient and robust for control and information transmission, among these levels and modules there are also complex asymmetric interactions formed by the feedback connections. Specifically, the top-down influences form a series of feed-back projections from a high-level neuron population to a low-level one, while there are also numerous lateral connections

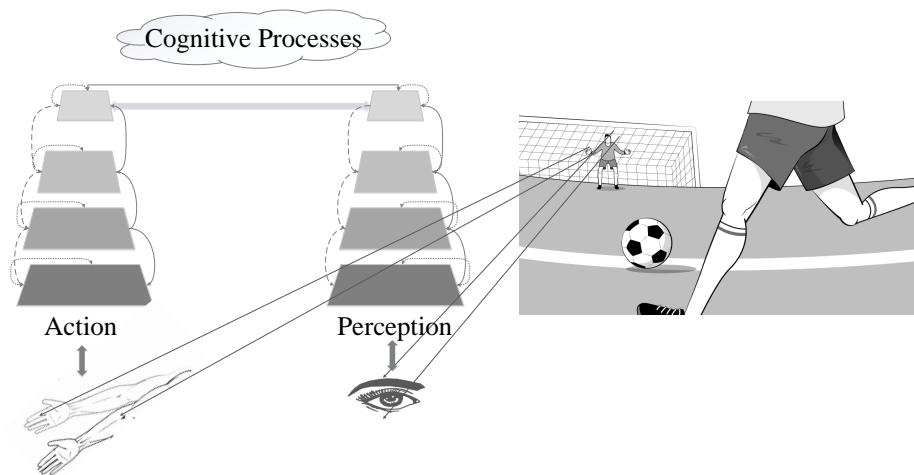


FIGURE 2.6: Example of Neuronal Feedback from Cognitive Processes: Goal-keeper

The perception predicts the object movement, which may account for the ‘flash-lag effect’ (we will introduce later). The motor action part is affected by the predictive perception too. It predicts the object movement from sensory information based on prior knowledge (e.g. experience) in order to accomplish a certain task (e.g. to save the goal)<sup>5</sup>. After that, if a high-level cognitive process, such as reasoning, makes a decision and concludes that there is no need to be fearful, such motor actions will slowly dissipate. In this example, emotion as a process on the cognitive level can affect the dopamine neurons; it also controls the motor cortex, the central nervous system (CNS) as well as muscles. Part of these control flows (from emotion process to neural and dopamine modulation and from reasoning process to motor cortex) can be included in two feedback pathways in Fig. 2.5.

<sup>5</sup>Copyright by Kira Chow.



within one cortical area. As the V1 cortex has been well studied, an example is made in this particular area.

Generally speaking, the visual information at the primary visual system is specialized to process static and moving objects and their patterns, while the higher-level areas receive the basic visual information from the primary visual system and generalize more abstract information, such as visual object identity and movement. Although there exist varieties between the exact connectivities of different species of animals in their sensory or motion cortices, the bottom-up influences in different biological cognitive systems generally play a similar role: they receive information from the raw sensory input and extract information in a hierarchical way. Each level in this hierarchical organization processes one specific feature within the bottom-up influences; this processed feature proceeds to the next higher level [Hubel & Wiesel, 1963].

For instance, the feedback connections ubiquitously link the cortical areas, modulating the neural activities by the top-down and recurrent influences.

- There are direct feedback projections from higher-level (V2, V3, V4, V5) to V1 (e.g. [Ungerleider & Desimone, 1986a,b, Shipp & Zeki, 1989, Rockland & Van Hoesen, 1994]).
- It has also be found that direct feedback projecting signals from V1 to Superior Colliculus (SC) and Lateral Geniculate Nucleus (LGN) and pulvinar (e.g. [Lund et al., 1975, Graham, 1982, Fries, 1990]).
- Lateral connections also ubiquitously exist within the same neural structure. (e.g. [Raiguel et al., 1989, Kisvarday et al., 1997, Angelucci et al., 2002]).

A similar case can also be found in the motor system, where the motor primitives ascend through various areas [Calais-Germain & Lamotte, 1996]. To sum up, the bottom-up influences on the sensorimotor process provide a raw sensory input to construct a representation to perception and motor actions by analysing its raw inputs and building up an impartial source of the perceptual world or action commands.

If we trace down the roots of the feedback pathways, except the lateral inhibition connections within V1, the feedback signals are partially derived from the cognitive processes such as memory, experience and expectations. Experiments done by Bar et al. [2001] suggested that the predictive perception and action should come from memory. Therefore, our perceived world is constituted from both a proactive

prediction based on the prior experience of sensory input as well as the actual incoming information [Herwig & Schneider, 2014]. Furthermore, predictions are suggested to account for the construction of a static visual world despite of the fact that the saccade is happening continuously (e.g., [Colby et al., 1992, McConkie & Currie, 1996, Rolfs et al., 2011, Sommer & Wurtz, 2006]).

### 2.5.1 Embodied Feedback Pathways

Since environment-agent interaction is the fundamental part to develop a complete autonomous (artificial or biological) agent [Brooks, 1991, Beer, 1995], the same rule applies to the development of the feedback pathways on the sensorimotor integration of robotic systems too. However, a few differences between them should be addressed:

- First, the configurations of the agents are different from each other. As we are already aware, the feedback pathways are mostly rooted in the past knowledge which is obtained from environment-agent interaction, encoded in genes or computer codes and stored in the brain neurons or storage media. For instance, what an artificial agent perceives from its camera is RGB pixels, which is fundamentally different from what a biological agent perceives with its eyes.
- This difference further relates to the cues and coding in the knowledge source of the feedback pathways, which is stored by the construction of new memory proteins in a biological system, or change of patterns of magnetization in an artificial system. Thus, the differences in anatomy and physiology of processing units (e.g. the brain) and memory also result in the diversity of knowledge representation and learning schemes on feedback pathways.

Nevertheless, the basic form of the feedback pathways (i.e. it is being transmitted from a high level to a lower one or within the same level) should be consistent. Although quite a few previous works have been focused on models of top-down influences (e.g. [Li et al., 2004, Gilbert & Sigman, 2007, Itti & Koch, 2001]) and lateral connections (e.g. [Amari, 1977, Sirosh & Miikkulainen, 1997]) of biological systems, few studies have been done on an artificial sensorimotor systems.

## 2.6 Phenomena from Feedback Pathways

Feedback pathways influence human sensorimotor integration. This can be found in various sensorimotor phenomena. Additionally, they are also able to indirectly demonstrate that the existence of the feedback pathways is beneficial to the brain cognitive functions.

### 2.6.1 Binocular Rivalry

Binocular rivalry depicts a visual phenomenon that when two distinct images are shown to each eye simultaneously: instead of the two images being overlapped or merged, these images are perceived for a few moments one by one in the perception, as they are competing in a consecutive and ‘bi-stable’ manner (Fig. 2.7). This usually happens when sufficiently dissimilar stimuli are presented to the two eyes; the stimuli can be as simple as different gratings in orientation or as complex as pictures of a human face or a horse.

This can be explained by the feedback pathways on the neural structure level, which correlates for the conscious visual experience. It changes stable visual stimuli into fluctuations in perception [Alais & Blake, 2005]. According to S.-H. Lee et al. [2004], such top-down influences may start hierarchically from high-level ‘expectation’ and descend to V1. When they are integrated with the sensory information perceived from two eyes, the percept is represented as the images competing with each other.

Furthermore, the functional mechanism of the feedback pathways is not only distributed in a hierarchical manner over visual cortices, but also forms part of the cognitive processes. For example, experiments done by Blake [2001] found that Jewish and Catholic believers judged the relative predominance of symbols mostly according to their two religions: Jewish believers are likely to see Judaist symbols during binocular rivalry in a longer visual dominance duration, while the Catholic believers are likely to see Catholic religious symbols in a longer visual dominance duration. Similarly, an upright human face tends to predominate over an inverted face, which could be accounted for by the feedback pathways of facial recognition area.

To summarize, binocular rivalry is a result from an imposition or stamping of the high-level temporal cortical areas back onto the V1 cortex. Thus, the feedback pathways may determine dominance appearances of the binocular rivalry. The

binocular rivalry dynamics could be influenced by multiple cognitive processes, such as sensory, cognitive, working-memory and affective factors [Tong et al., 2006].

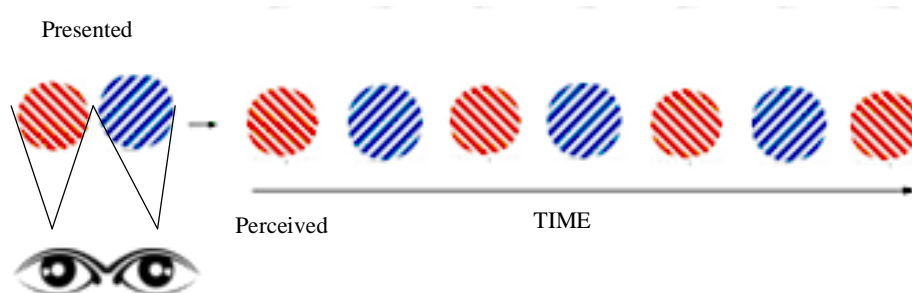


FIGURE 2.7: Binocular Rivalry

When two distinct stimuli are presented, in perceptual awareness over time two different stimuli compete for perceptual dominance<sup>6</sup>.

## 2.6.2 Retina Prediction

The vertebrate retina is a tissue in the visual system that converts light energy into electrical signals in a form of nerve impulses, by its ganglion cells. As part of the visual system, the ganglion cells inside the retina have the off-centre and on-centre properties: the off-centre cells with excitatory synapses are hyper-polarized by light, while the on-centre cells having inhibitory relationships with synapses are suppressed without the light.

However, this property of the on-centre-off-surround receptive field (or reverse) is not always held with constant lighting. They also have an ability to adapt according to background illumination, intensity and duration of stimulation and other factors of the stimuli within the receptive fields [Kuffler et al., 1953]. The firing rate of both types of ganglion cells may be altered if there is a prolonged presence of contrast or luminance stimuli presented. This could not only account for natural selection from development, but also the relatively short adaptation time (only a few seconds) indicates that this is also the result from the adaptation

<sup>6</sup>This image is a derivative work based on <https://openclipart.org/detail/856/eyes-by-molumen>, which is licensed under the Public Domain license.

of spatio-temporal organization of the receptive fields of the ganglion cells, with the predictive vision learnt from prior knowledge.

Besides, the majority of ganglion cells are also sensitive to temporal movement patterns, so that they are able to follow the detection of temporal patterns in the environment in a predictive manner [Barlow, 1953]; for example, measurement of different periodic waveforms has shown that the operation of the ganglion cells is modulated by both luminance and chromate. Hosoya et al. [2005] proposed that the local encoding in the ganglion cells also have differences in raw image intensity in a biphasic temporal sequence.

Therefore, the adaptation of ganglion cells in the form of spatial antagonism and temporal antagonism should be a result of the feedback pathways. To sum up, the information provided from feedback pathways to the retina encode the possible image intensities in the perception, assuming that nearby spots of the receptive field will display similar image intensities. This adaptation can be achieved by ‘anti-Hebbian’ learning<sup>7</sup>, which serves as a ‘novelty filter’ that learns to suppress the sensitivity of ganglion cells which correspond to the predictable elements. At the same time, it increases the correlation between the ‘unexpected forth-coming events’ and the visual stimuli.

### 2.6.3 Feedback Pathways on Predictive Action

We mainly introduced the feedback pathways on perception. However, due to the homogeneous hierarchical organization of perception and action, the feedback pathways should exist in both parts (perception and action). Also, those influences are not independent but may influence each other, which results in the feedback pathways integrating the most updated sensory stimuli as well as the current motor action on the whole sensorimotor integration.

Brown et al. [2011] asserted that the peripheral motor action is a kind of active inference in estimating the attention execution in order to compensate predicted sensory signals. This assumption accounts for the fact that the spatial attention is mediated by feedback pathways, which usually predict the possible visual stimuli with consideration of the action inference. In the context of embodied intelligence, if an agent actuates an action, it may also tend to react with the environment changes according to this very action, so that the sensation matches the mental

---

<sup>7</sup>Anti-Hebbian learning is a learning process in contrast to ‘Hebbian learning’ which proposes an algorithm that the corresponding synaptic weight increases if a repeat firing of one cell contributes to the firing of another cell connected.

prediction. Such action mediated feedback pathways can account for the visual stability [Wurtz, 2008], which means that our visual perception remains stable even though the movements of our eyes, head and body create a stochastic yet predictable motion pattern. Here this action inference may include a scene or goal movement understanding, long-term memory and expectation. This facilitates quicker processing of the incoming visual information when an agent has prior knowledge of this kind of advanced information, which biases the processing of incoming visual information. This mediated attention also makes our perceptual world stable by utilising an unfolded sensory prediction to cause the subsequent motor action [Hawkins, 2004]. The representation of this mechanism may be similar to the predictive sensorimotor integration framework, such as Wolpert & Kawato [1998] and Kawato et al. [2003]. Also, this can be modelled as a partially observable Markov decision process (POMDP) model in the perspective of optimal feedback control theory [Todorov & Jordan, 2002].

## 2.7 Representation as Bayesian Inference

In terms of the mathematical representation of the feedback pathways, the Bayesian Inference may be one common link between the biological and the artificial systems. Von Helmholtz, as a pioneer to interpret perception within a Bayesian framework, proposed a general rule: that perception is actually composed of visual statistical representations, which are determined by the previous perceptions themselves [von Helmholtz et al., 1909]. He stated that previous perception must be inferred in order to fully understand the pattern perception appearance, so that a single percept can be regarded as a result of a full description of the prior as well as the raw sensory inputs.

This hypothesis inspired a group of theorists proposing inverse inference and ‘analysis-by-synthesis’ (e.g. [Neisser, 1967, D. MacKay, 1956, Gregory, 1980]), which infers an up-coming prediction from a pre-learnt knowledge from low-level sensory information (short-term prior), and a learnt internal model (long-term prior). Eqs. 2.1 and 2.2 show how a visual system infers the most probable representation according to the Bayesian perspective if we already know the prior probability of visual knowledge and action [T. Lee & Mumford, 2003]. It is achieved by the posterior  $S_i$  given a particular sensory evidence ( $E$ ), motor action ( $A$ ) and other prior information we have already known ( $I$ ).

$$P(S_i|E, I) \propto P(E|S_i)P(S_i) \quad (2.1)$$

$$P(S_i|E, A, I) \propto P(E|S_i)P(A|S_i)P(S_i) \quad (2.2)$$

where  $P(S_i|E, I)$  and  $P(S_i|E, A, I)$  are the conditional probabilities given the scene  $E$ , motor action  $A$  and the prior information  $I$ .  $P(S_i|I)$  and  $P(S_i|A, I)$  are the prior probability.

Generally, this inference process includes two steps: at the first step, the perception is obtained from inverse Bayesian Inference, which presupposes an integration of the past input, as well as a model used to determine how the input should be estimated. At the second step, the final inference of the perception is obtained from an integration from a top-down generative neuronal representation and the bottom-up influence. This can be formulated mathematically using a sequential hypothesis to compute the posterior by mixed probability distributions, each of which represents a single prior, depending on a causal relation inferred from sensory information.

Apart from the earlier works proposing to construct non-probabilistic generative models (e.g. [D. MacKay, 1956, Pece, 1992]), the theory from von Helmholtz inspired the development of a family of probabilistic models called Helmholtz machines [Dayan et al., 1995, Dayan & Hinton, 1996]. These models attempt to learn a new uniform representation of deep regularities through temporal cycling of perceptual sensing, thus creating a succinct internal model without any prior knowledge of pre-classified samples. This kind of generalisation methods performs stochastic recognition and reconstruction through the interaction of bottom-up sensory data and top-down expectation. Moreover, this kind of generative model is also akin to the information flow in the cortex. This explains the asymmetric information flows between the top-down and bottom-up influences, which encode with a more ‘experienced’ perceptual and an up-to-date sensory correction on their activity, respectively [Hohwy, 2007]. This procedure combines ‘top-down’ and ‘bottom-up’ influences in a delicate and potent fashion, and explains how the non-intrinsic activity in the perception area developed with interactions across levels [Friston, 2003].

Hawkins [2004] also proposed that the two basic principles of sensorimotor integration are ‘hierarchy’ and ‘integration’: when one’s own behaviour is involved,

the associated motor actions not only ‘precede sensation’ but they can also ‘determine sensation’ by changing the experience of sensation. This is partially realised by simulating the changes in the perceptual world prior to any actual movements. During this process, the nervous system transforms the difference between the current and the desired sensory coordinates into the motor system’s coordinates, so that the motor system is able to generate the necessary motor commands to move the actuator and thus reach the target state. At the same time, the efferent copy of the motor command forms a feedback signal that is able to improve the sensory perception by predicting the next time-step. A minimised representation of perception requires at least one single feedback (or bias) to deliver the error between the expectation and the sensory information hierarchically [Fletcher & Frith, 2008, McMains & Kastner, 2011], by learning from reciprocal interactions on various hierarchical levels. Although the Helmholtz machine is identical from observation in a neuroimaging study that the top-down influences actually activate the neural activities in the perception by providing ‘surprise’ signals [Egner et al., 2010], a single-layer Helmholtz machine cannot be used as a universal function approximator. Instead, a more simple but still Bayesian compatible method is learning by back-propagation [D. J. MacKay, 1996]. This may be practical to be implemented and adopted [McClelland & Rumelhart, 1981, Rumelhart & McClelland, 1986]. Learning mechanisms in the field of artificial neural networks will be discussed in detail later in Chap. 3.

To summarise, one role of feedback pathways (especially in the perception) can be explained by a Bayesian model which is able to estimate the statistical dependencies on a high-level cognition from the temporal perception inputs by an inverted generative process (i.e., the top-down influences deliver a prior knowledge). This model justifies the intertwining relation between neuronal activities on the high-level of perception and motor action and on the low-level sensory motor action cortices.

## 2.8 Hypotheses of Feedback Integration

We have been aware that part of the feedback pathways can be formulated as a Bayesian inference. However, investigations about how the brain integrate these two sources of information seamlessly in a dynamic and rapidly changing environment are still in progress. A few theories have been proposed concerning how the



integration of top-down and bottom-up influences happen in the perception and action cortices.

### 2.8.1 Predictive Coding

As argued by von Helmholtz et al. [1909], what is perceived by the brain is not exactly what is sensed, as the brain itself continuously predicts the upcoming percepts and corrects a certain kind of error in various hierarchies of the cortices. As shown in Fig. 2.8, predictive coding theory proposes that these cortices work in a cascade so that the high-level system attempts to predict the up-stream from the low-level statistics according to the innate or learnt models (e.g. Bayesian model) [Friston, 2005, Rao & Ballard, 1999, Clark, 2012]. In other words, *only* the errors are accumulated and transmitted from the lower level of sensory input to the higher-level cognition. This is how the bottom-up perception information adapts the environment changes and reduces the error between top-down expectations and bottom-up raw sensory inputs. Thus, the predictive coding theory asserts that the brain is always working in ‘error-correction’ mode.

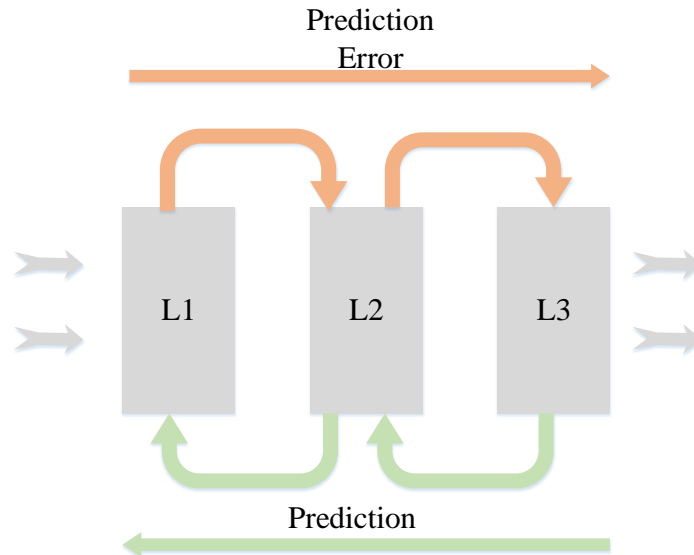


FIGURE 2.8: Schematic of Predictive Coding

According to predictive coding, prediction about the incoming sensory information comes from each level of this hierarchical architecture. If the expected information is different to the prediction, error signals or surprising information is incorporated in the feedback. Also, the predictive model is updated.

The hierarchical predictive coding theory is suggested by a number of experimental studies. For example, functional magnetic resonance imaging (fMRI) data has shown that illusory contours remain in the early visual cortex even if the contours from sensory input disappear [Muckli et al., 2005], which can be explained by the encoding at the local lateral interactions within V1 as well as the top-down predictive influences mediate from higher-level motion-sensitive areas such as MT/V5. The predictive coding theory can also explain other low-level adaptation phenomena, such as relative attenuation of neural signals (a.k.a. repetition suppression) [Summerfield et al., 2008], which means a reduction of neural response when stimuli are presented repeatedly. To further explain the predictive coding theory on the visual motion processing in a neuroanatomical context, Bar [2004] proposed that the medial frontal regions should encode predictive templates (constructed by individual objects) that are learnt associatively in contextual scenes at a higher visual cortex, which acts as a low-frequency adaptive filter. Similarly, this has also been observed in electrophysiological recording of auditory cortices; while the participants listened to auditory stimuli with varying pitch strength, the neural activities between the adjacent and primary auditory cortical areas can be explained with the principle of predictive coding [Kumar et al., 2011].

To summarize, all of these phenomena can be interpreted by the theory that there is an innate mechanism that compares expected and actual inputs (previous percepts in perception, or both percept and actuated actions in sensorimotor processing) in information processing, during which the expected information acts with a predictive coding-like mechanism. If the prediction of perception is perfectly identical to the raw sensory input, successful perception, cognition and action transmit an identical suppression which ‘explains away’ prediction error so that no neural suppressions (responses) happen. If the predictive coding is related to action, it is also considered to be a type of suppression that always attempts to minimise the expected percept, similar to the ‘error correction’ in predictive coding theory.

The coding mechanism is similar to the data compression technique used in audio signal processing and speech processing that encode the signal at the current time-step by using a weighted representation from the previous time-steps. That is where the term ‘predictive coding’ was taken from.

## 2.8.2 Biased Competition

An alternative theory, biased competition, can also explain the role of feedback pathways on the sensorimotor coordination. This theory claims that the sensorimotor process is the result of competitive interactions among large assemblies of neuronal processes, including the top-down and bottom-up influences as well as local lateral neuronal dynamics. This biased role from top-down and lateral influences on the bottom-up sensory input may be similar to negative feedback (Fig. 2.9); but the feedback signal comes from more than one process on a single level. Among those processes, there is no explicit selection for specific processes which results to behaviourally relevant stimuli. Instead, the final spatio-temporal sequences in perception are biased by all of those stimuli [Desimone, 1998].

In the single-cell activity study conducted by Kastner & Ungerleider [2001], when multiple but simultaneous stimuli are shown within the same receptive field, the neural response to the paired stimuli was reduced compared to a single stimulus. It suggests that these stimuli are not processed independently in the visual process but that they are mutually suppressed. Therefore, it also implies that the larger the number of stimuli, the smaller chance that attention is routed to a specific object due to the increase of neural suppression. Other neurological studies, such as neuronal spiking recording by Yilmaz [2012], have also suggested that attention may result from a biasing routing by the feedback pathways.

This theory can explain the formulation of the attention by regarding that it is an obligatory competing process for tracking multiple objects. For an attention process, one reason of employing such a biased mechanism is that a visual system (as well as other perception systems) has only a limited information capacity to focus on the object (stimuli) of interest. Therefore, when multiple objects are presented simultaneously in the visual field, the stimuli will compete in the neural representation due to the limited routing and processing resources of the visual system. According to the biased competition theory, a final attention might be caused by a biased effect from some other mental processes, which results in the tracking of the features which are previously attended in the visual field and prioritise the feature-driven attention. In addition, this biases the attention to move toward the object which is the most relevant to the agent's behaviour, or is the most interesting in conscious/unconscious processes.

Although it seems that the theories of the biased competition and the predictive coding are incompatible, in principle they both depict that the feedback

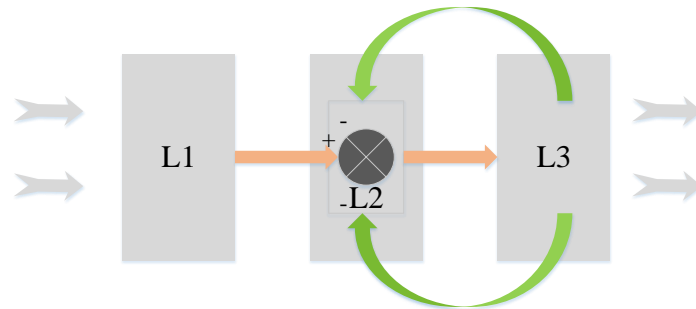


FIGURE 2.9: Schematic of Biased Competition

On each level, competition among multiple stimuli is biased by mutual influences.

pathways are integrated with bottom-up influences in various manners. Physiologically speaking, this is realised by the competition of neuron populations via lateral inhibition in a cortical region. Globally it becomes neural suppression. Spratling [2008] unified these two theories mathematically as well, with which he successfully explained the single-cell electrophysiological mechanism in visual attention.

## 2.9 Discussion

How the feedback signals and the sensory-driven influences integrate may be explained by the free-energy principle. This principle explains that a (biological or other dynamical) system tends to minimise the free-energy function of their finite degrees of internal states by maximizing the similarity of their internal mental states to environment orders. In other words, the free-energy can in-turn be measured by a *surprisal* (which represents the quantity of ‘surprise’ of seeing an outcome). This *surprisal*, in the context of perception, is delivered by the feedback pathways. Therefore, reducing the information-theoretic free energy inside a system’s world model means to reduce prediction errors, i.e. reduces *surprisal*. This is how the statistical regularities are employed to infer sensory perception by

eliminating the redundant energy in the uniformity of the spatial and temporal domains of the sensory signals.

Furthermore, a similar principle also applies to any biological system (from single-cell organisms to social networks) that resists a tendency to become disorderly, which is probably a result of evolution to efficiently save energy to perceive sensory inputs, as well as survive in a dynamically changing environment. This could also account for the topological characterisation of cortical areas, where a complex network is formed based on the best optimisation result of the structural and functional organisation in order to ensure the shortest routes for frequent information transmission. Therefore, neurons with similar functions are clustered, structures that communicate often are close. Connections are mainly reciprocal. This process of modularity is naturally formed with the most efficient information presentation [Clune et al., 2013]. The minimisation of energy is also identical to the basic training rule of machine learning techniques, such as the Boltzmann machine [Hinton & Sejnowski, 1983] or back-propagation learning [Rumelhart et al., 1986].

The hierarchical control of physiology partially contributes to the formation of the feedback pathways in the cognitive processes, in which smaller numbers of parameters (e.g. a symbolic command) determines the high degree of freedom (DOF) dynamics in the lower layer. This has also been discovered by motor control in biological systems. As Bernstein [1967] proposed, the generation of certain movements is not so trivial, as the body in most biological systems is a highly-redundant manipulator, which is comparable to solving the control in extreme abundance of DOF systems in a stable manner. Kuppuswamy & Harris [2013] recently also claimed that only smaller number of variables are needed for the acquisition of motor control. Such an idea can be extended from the central nervous system (CNS) to other hierarchical control of cognitive processes (e.g. the grandmother cell theory [Gross, 2002]). Conversely, the sensory-driven bottom-up processes relieve the high-dimension of neuro-mechanical redundancy in the body of organisms. This is realized by extracting relatively persistent ‘profiles’ of sensory data (e.g. identity of a moving visual object, an emotion-driven or goal-directed behaviour) to construct high-level cognitive processes. Such ‘profiles’ may result in the emergence of language acquisition too (Chap. 6).

These feedback and sensory-driven influences are highly integrated, which results in an agile movement and a stable perceptual world of the cognitive agents. Therefore, in this thesis, we also advocate this idea and build up a model to interpret

how the ‘familiarity’ is derived from the previously learnt knowledge on the high-level mental state and affects the low-level sensorimotor processes by top-down influences.

## 2.10 Summary

This chapter begins with a review about the hierarchical modularity of physiology, from which we conclude that there are multi-dimensional feedback pathways within the same level of physiology, and across various levels. Due to this fact, the feedback pathways affect neuron activities, control motor actions, carry out metabolism and restructure physiological functions. Specifically, in this chapter, the feedback pathways on the sensorimotor integration of the brain at the neuronal and cognitive levels are investigated. The roots of the feedback pathways here are cognitive processes like goal movement understanding, long-term memory and expectation, or neural processes like first-order movement encoding.

The hierarchical modularity organization of sensory and motor cortices are connected with long-distance inter-cortical connections, which are physiological paths of feedback pathways. On each level, these feedback pathways are distributed and integrated with the sensory-driven stimuli. In this way, they affect the perception and action in various ways and can be observed as a few phenomena and illusions, such as binocular rivalry, prediction in vertebrate retina and feedback pathways on motor actions.

As the feedback pathways can be formulated as a Bayesian inference, which is identical to the hypothesis by von Helmholtz, the integration of the bottom-up and feedback influences can be explained by two main theories: namely predictive coding and biased competition. Both of them emphasize the fact that these two influences should integrate on each level and affect each other.

It inspires a few models in cognitive modelling, computational neuroscience and machine learning (e.g. back-propagation in multi-layer perceptron (MLP), Helmholtz machines and other variant models). Since the feedback pathways existing on different levels of the hierarchical cortical areas (especially visual and motor cortices) endow predictive sensorimotor functions for a cognitive system, it also encourages us to design an architecture to realize such mechanisms in artificial cognitive systems.

## Chapter 3

# Artificial Recurrent Neural Network Models of Feedback Pathways

### 3.1 Introduction

The mechanisms of the neural signal feedback transmission motivates us to implement it in an artificial cognitive agent. In this chapter, we investigate the possibility of using an artificial recurrent neural network (ARNN) to be one of the possible techniques to achieve this target by reviewing its background.

From the perspective of machine learning, ARNN are a class of artificial neural networks which usually include directed connections between units and which contain cycles in the graphical model. These connections establish feedback connections and maintain a network activation in a temporal loop. This chapter firstly introduces the basic component of an artificial neural network (ANN): the perceptron model. After the introduction of the multi-layer perceptron (MLP) networks, we describe how the recurrent connections are constituted and trained, and investigate how they contribute to the network dynamics in a simple ARNN and other ARNN variants. A comparison between an ARNN and the neural feedback mechanism will be also given.

## 3.2 Multi-Layer Perceptron

An MLP network is an artificial neural network that consists of multiple layers of simple units called perceptrons. These units are connected by directed weighted connections without any cycles or loops in the network.

### 3.2.1 A Single Perceptron

The perceptron model designed by Rosenblatt [1958] computes an output from a non-linear transfer function with a weighted summation of all inputs and a bias value. Mathematically, the output of the perceptron can be written as:

$$z = \varphi\left(\sum_{i=1}^n w_i x_i + b\right) = \varphi(\mathbf{w}^T \mathbf{x} + b) \quad (3.1)$$

where  $\mathbf{w}$  denotes the vector of weights,  $\mathbf{x}$  is the vector of inputs,  $b$  is the bias.  $\varphi$  represents a binary classifier:

$$\varphi = \begin{cases} 0 & \text{if } \mathbf{w}^T \mathbf{x} + b > 0 \\ 1 & \text{otherwise} \end{cases} \quad (3.2)$$

As shown in Fig. 3.1, the actual output of the perceptron can be written as

$$\varphi(\mathbf{w}^T \mathbf{x} + b) = z \quad (3.3)$$

Eq. 3.3 indicates that the output of a single perceptron is a linear function of all the inputs. In the simplest case of two-dimensional problems, it means that the two inputs  $x_1$  and  $x_2$  are separated by a straight line which is determined by two weights and the bias of the perceptron:

$$w_1 x_1 + w_2 x_2 + b = 0 \quad (3.4)$$

Similarly, in a higher-dimensional space it means that the data points can be classified by a hyperplane.

Training of a perceptron is done by adjusting the weighting matrix connecting to the inputs with feature vectors correlated to the error between desired and actual outputs at one iteration (Eq. 3.5). This correction step is executed iteratively



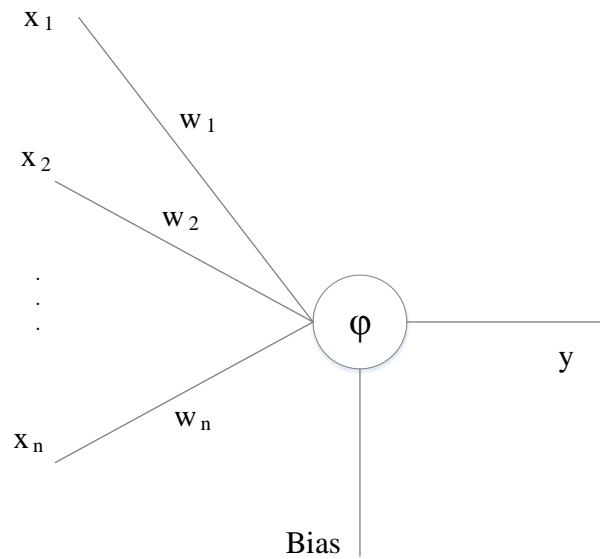


FIGURE 3.1: A Perceptron Unit

until the network learns to reproduce the desired response or the error is within a certain threshold.

$$w_i(t + 1) = w_i(t) + \eta x_i \delta \quad (3.5)$$

where  $x_i$  is the  $i$ -th input, and  $w_i$  is the corresponding weight,  $\eta$  is the learning rate and  $\delta$  denotes the output error:

$$\delta = d - z \quad (3.6)$$

which means that the update of the weights depends on the difference between expected (target) output  $d$  and the actual output  $z$ .

### 3.2.2 Multi-Layer Perceptron Network

A single perceptron is a linear classifier; it cannot solve a simple XOR problem<sup>1</sup>. Since the nature of the decision boundaries varies with the network topology, it has

<sup>1</sup>XOR is true whenever an odd number of inputs is true.

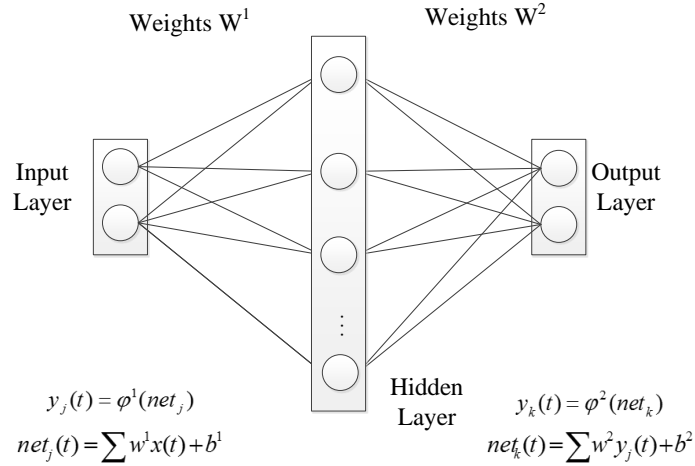


FIGURE 3.2: A Multi-layer Perceptron (MLP)

The transfer functions  $\varphi^1(\cdot)$  and  $\varphi^2(\cdot)$  in an MLP should be differentiable non-linear functions.

been proved that a neural network with three layers, i.e. an MLP network, is able to generate arbitrary decision boundaries. This network is constituted by a group of perceptron units that are being connected by a set of weighted connections. In this way they constitute an input layer, a hidden layer and an output layer. As all the connections are directed, the input signal actually propagates through the network layer by layer (Fig. 3.2).

Fig. 3.2 presents how an output can be calculated from forward-propagation layer-by-layer in an MLP; the function of an MLP can be interpreted as a non-linear mapping  $f : R^D \rightarrow R^P$ , where  $D$  is the size of input vector  $x$  and  $P$  is the size of the output vector  $f(x)$ . Generally, a full mapping function in matrix notation is given as:

$$z = f(x) = \varphi^2(b^{(2)} + \mathbf{W}^{(2)}(\varphi^1(b^{(1)} + \mathbf{W}^{(1)}x))) \quad (3.7)$$

where  $b^{(1)}$  and  $b^{(2)}$  are bias vectors,  $\mathbf{W}^{(1)}$  and  $\mathbf{W}^{(2)}$  are weighting matrices, and  $\varphi^1(\cdot)$  and  $\varphi^2(\cdot)$  are transfer functions in hidden and output layers, respectively.  $\mathbf{W}^{(1)} \in \mathbb{R}^{D \times H}$  and  $\mathbf{W}^{(2)} \in \mathbb{R}^{H \times P}$  are the two weighting matrices between three layers.

Particularly, the transfer function in a perceptron unit can be realised with any differentiable non-linear function, such as the logistic sigmoid function:

$$z = \varphi(y) = \frac{1}{1 + e^{-y}} \quad (3.8)$$

or the hyperbolic tangent function:

$$z = \varphi(y) = \frac{1 - e^{-2y}}{1 + e^{-2y}} \quad (3.9)$$

Such a non-linear function is used as a transfer function in the perceptron model, because it prevents rapid saturations and the learning makes the output to regress into a certain range until the cost function of the network reaches the local optima (usually the cost function reaches the local minimum value) after training. This allows the output of the perceptron unit to be constrained within a certain pre-defined interval, thus enabling it to contribute to a stable mapping to the desired output values. In practice, such selection of transfer functions often depends on the experience of the designer based on what kind of data (inputs and outputs) the model has to cope with.

We use vector  $\theta$  to define all the parameters of an MLP that are to be learnt,

$$\theta = \{W^{(2)}, b^{(2)}, W^{(1)}, b^{(1)}\} \quad (3.10)$$

Similar to what we introduced in the single perceptron section, a logistic sigmoid, the hyperbolic tangent function or other non-linear functions can be utilised as the transfer function of a layer of an MLP. The update of this vector is basically based on the error back-propagated from the outputs.

### 3.2.3 Back-propagation

The training of an MLP is more complicated than training of a single perceptron, because the error in neuronal units of the hidden layer is difficult to define. Usually it is derived from a cost function, which is a direct mapping to the difference between the actual and the desired output patterns; using this difference, the gradient of a cost (or error) is used to modify the weighting matrices according to the conventional gradient decent method. The most frequently used cost function is the summed squared error (SSE), which is defined as the summation of

the squared differences between each actual output and its corresponding desired output.

$$C = \frac{1}{2} \sum_t^T \sum_k^P (d_k(t) - z_k(t))^2 \quad (3.11)$$

where  $d$  is the desired output vector,  $T$  is the total number of training samples and  $P$  is the size of the output vector. According to the gradient descent, each weight change in the network should be proportional to the negative gradient of this cost with respect to the specific weight that is going to be modified. Intuitively, the larger the neural activation is, the bigger error it contributes, so it should be corrected more by the training process. Mathematically, it can be written as

$$\Delta w = -\eta \frac{\partial C}{\partial w} \quad (3.12)$$

where  $\eta$  is a learning rate.

In order to obtain the partial derivative, the exact weight change with respect to the cost function can be rewritten as the product of the internal error of each neuron  $\delta = -\partial C / \partial net$  and the network output with respect to the specific weight  $\partial net / \partial w$ .

$$-\eta \frac{\partial C}{\partial w} = -\eta \frac{\partial C}{\partial net} \frac{\partial net}{\partial w} \quad (3.13)$$

in which  $\partial C / \partial net$  is the internal error of the output layer. It can be derived as,

$$\delta_k = \frac{\partial C}{\partial net_k} = -\frac{\partial C}{\partial \varphi^2} \frac{\partial \varphi^2}{\partial net_k} = (d_k - y_k) \varphi^2(y_k)' \quad (3.14)$$

where  $\varphi^2(\cdot)$  is the transfer function of the output layer.  $\varphi^2(\cdot)'$  represents the first derivative of the output layer's transfer function.

Similarly, the error of the hidden layer is given by,

$$\delta_j = -\sum_k^P \frac{\partial C}{\partial \varphi^2} \frac{\partial \varphi^2}{\partial net_k} \frac{\partial net_k}{\partial \varphi^1} \frac{\partial \varphi^1}{\partial net_j} = \sum_k^P \delta_k w_{kj} \varphi^1(y_j)' \quad (3.15)$$

where  $w_{kj}$  is the element of  $k$ -th row and  $j$ -th column in the weighting matrix  $\mathbf{W}^{(2)}$ ,  $\varphi^1(\cdot)$  is the transfer function of the hidden layer.  $\varphi^1(\cdot)'$  is the first derivative of the hidden layer.

For a first-order polynomial,  $\partial net / \partial W^{(2)}$  equals the outputs of the hidden layer. Denoting the  $k$ -th row and  $j$ -th column of element in matrix  $\mathbf{W}^{(2)}$  as  $w_{kj}$ , the weighting matrix  $\mathbf{W}^{(2)}$  between the hidden and output layers is updated according to

$$\Delta w_{kj} = \eta \delta_k y_j \quad (3.16)$$

Also  $\partial net / \partial W^{(1)}$  equals the network inputs. Denoting the  $j$ -th row and  $i$ -th column of element in matrix  $\mathbf{W}^{(1)}$  as  $w_{ji}$ , the weight change  $\mathbf{W}^{(1)}$  between the input and hidden layers is given by:

$$\Delta w_{ji} = \eta \delta_j x_i \quad (3.17)$$

where  $x_i$  is the activation of the  $i$ -th element of the input layer.

Thus, combining Eqs. 3.14 and 3.16, we can derive the update of the weights between hidden layer and output layer  $\mathbf{W}^{(2)}$ , and combining Eqs. 3.14, 3.15 and 3.17, we can derive the update of the weights between input layer and hidden layer  $\mathbf{W}^{(1)}$ .

### 3.3 Artificial Recurrent Neural Networks (ARNNs)

When additional weighting matrices are connecting either two layers in one network with a directed cycle, a recurrent neural network is established. With such a cycle, the outputs of an RNN become input functions of next states, thus the internal states of an RNN can theoretically affect the network dynamics with an infinite time-length.

#### 3.3.1 Recurrent Connections

Recurrent connections offer a feedback loop in the whole network, which enables signals from one layer to be fed back to a previous layer. With this directed flow, neural signals can affect *future* network dynamics, so that arbitrary temporal sequences can be represented as a function of the previous internal network states. The ARNN is able to learn any arbitrary dynamical system with arbitrary precision, whereas an MLP can merely learn static non-linear models.

Typically, a feedback connection within one layer or between two layers forms a simple recurrent connection. Two simple recurrent networks (SRNs) are illustrated in Fig. 3.3 (Elman Network [Elman, 1990]) and Fig. 3.4 (Jordan Network [Jordan, 1997]). In both of them, the hidden layer is updated not only from the external input of the network, but also with activation of certain layers from the previous states. Specifically, since the connections are recurrently connected within the network (i.e. they feed-back onto the network itself), the local recurrence results in a compensation of the decrease of the internal neural activities by a constant  $\tau$  (a time constant that defines how many time-steps the network is unfolded during training) as they are fed back by building a short-term memory. This constant also determines the depth of the short-term memory (i.e. how long a given value fed to the context unit will be stored).

Like the feed-forward connections, this feedback is also learnt by modifying the weights which enable an adaptation of the temporal effects of the internal states. However, since the weighting matrix connection is in the temporal domain (i.e. the past activity of one layer is the input of the current activity of another/the same layer), it is necessary to modify the back-propagation algorithm.

### 3.3.2 Back-propagation through Time

Elman [1990] proposed an approximation learning rule based on truncated back-propagation, in which the time-delayed input  $x_i(t-1)$  is regarded as an additional input (i.e. time constant  $\tau = 1$ ), so that error from the output patterns was also back-propagated to the weighting matrix between the hidden layer and the additional input layer. However, it was found that this approach is not sufficiently accurate to find the optimal weight change according to the gradient descent, because the effect of error should further propagate even further in the temporal domain. Therefore, back-propagation should also be applied in the temporal domain, which leads to the so-called ‘back-propagation through time’ (BPTT) approach [Rumelhart et al., 1988].

The time-constant  $\tau$  defines the number of time-steps within which the error is back-propagated. Then all the recurrent connections are duplicated spatially within these time-steps, thus it is a mapping from the temporal dynamics to the spatial dynamics. We take an Elman network as an example, as shown in Fig. 3.3; each hidden layer sends its activation (either directly or indirectly) to the current output by recurrent connections, with  $\tau$  numbers of copies of the neuron activations. Therefore, all of the internal states within  $\tau$  time-steps contribute to

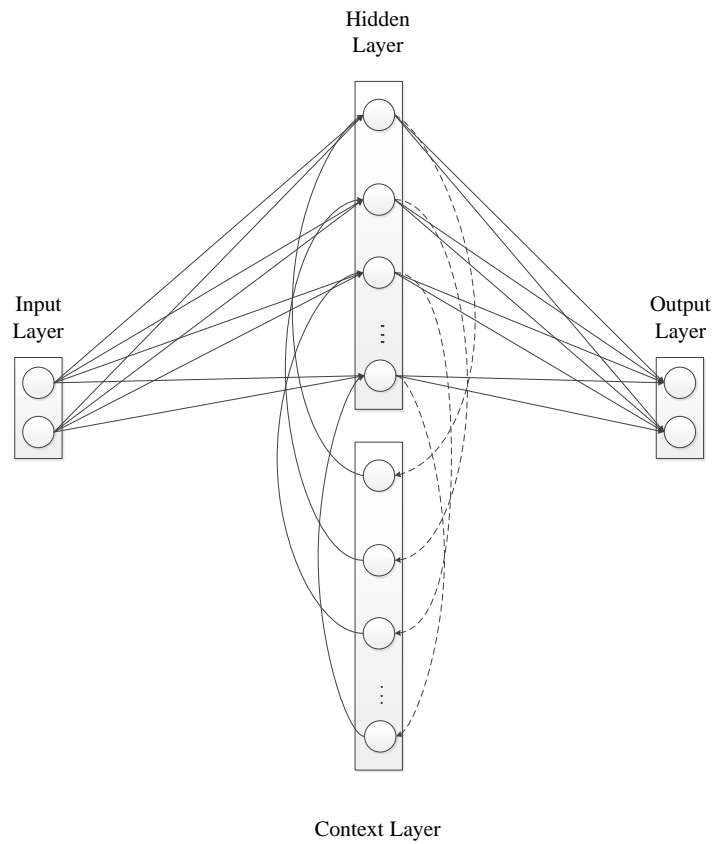


FIGURE 3.3: Elman RNN Network

Weights with dash arrows are fixed, i.e. the context layer is the copy of the hidden layer at the previous time-step.

the output. Conversely the output error can also be back-propagated along these unfolded connections.

From Eq. 3.15, the error on the hidden layer at the previous one time-step, which can be obtained by back-propagation from the current state of the hidden layer by the recurrent weights can be written as

$$\delta_j(t-1) = \sum_j^H \delta_j(t) W^3 \varphi^1(y_j(t-1))' \quad (3.18)$$

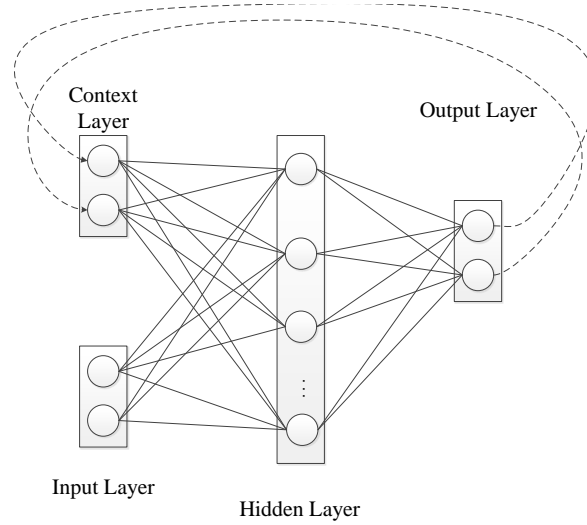


FIGURE 3.4: Jordan RNN Network

Weights with dash arrows are fixed, i.e. the context layer is the copy of the output layer at the previous time-step.

where  $h$  is the index of the vector that is from time-step  $t$ , and  $j$  is the index of the vector that is from the previous time-step  $t - 1$ .

If the error on the hidden layer is calculated back to the  $T$ th step, it can be written as a recursive form:

$$\delta_j(t - T) = \sum_j^H \delta_j(t - T + 1) W^3 \varphi^1(y_j(t - T))' \quad (3.19)$$

where  $T \in \{1, 2, \dots, \tau\}$ . Hence, the recurrent connection is updated (for totally  $\tau$  times for each sample) by

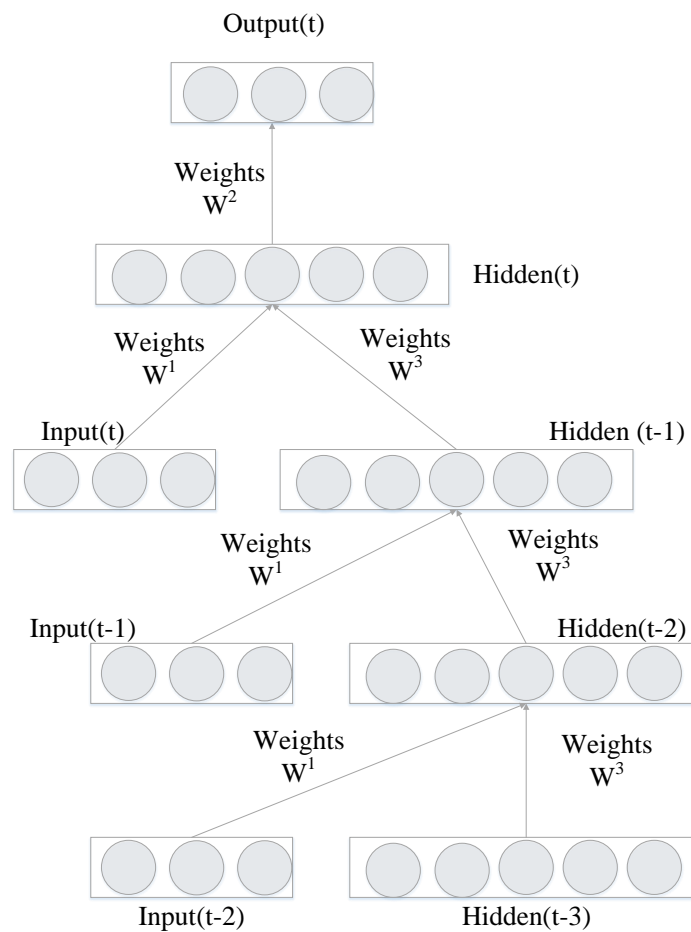
$$\Delta w_{jj'} = \eta \delta_j(t - T) y_{j'}(t - T + 1) \quad (3.20)$$

Note that the update of the weights  $\Delta w_{jj'}$  can be done within several time-steps, depending on the time-constant  $\tau$ . Therefore, Eq. 3.20 is applied  $\tau$  times for each update.

Generally, this equation allows a calculation of the weight update according to the error at previous time-steps from  $t - \tau$  to  $t$  with arbitrary length. Further deduction can be made according to Eqs. 3.18 and 3.20.

However, it is important to note that error  $\delta$  contributes to each weight in Fig. 3.5



FIGURE 3.5: Unfolding an RNN for BPTT ( $\tau = 3$ )

which is unfolded in a spatial domain. Clearly, this requires a lot of memory to store both previous errors and activations if we set  $\tau$  to be too large. Additionally, in practice, the error from previous time-steps with too large  $\tau$  makes too small contribution to the weight update due to a ‘vanishing gradient effect’ [Bengio et al., 1994], which means that at each time-step (each layer in Fig. 3.5), when the error is back-propagated through one time-step, it gets smaller and smaller until it quantitatively vanishes. On the other hand, the error from previous time-steps with too small  $\tau$  will result in another truncated back-propagation.

### 3.3.3 Variants of ARNN

As we mentioned, the generic ARNNs only own a short-term memory which results from the simple recurrent connections. However, in some circumstances, we need a memory that can sustain longer. Besides, multiple spatiotemporal training sequences usually result in training divergence when they attempt to implement multiple attractor dynamics if the network memory is too small. Therefore, more state-of-the-art recurrent networks have been developed to address these problems.

**RNN with Parametric Biases** A recurrent neural network with parametric bias (RNNPB) [Tani & Ito, 2003] is capable of learning different sequences with different parametric biases. These sequences are learnt as non-linear dynamic attractors while the parametric biases are represented as bifurcation parameters. Interestingly, the parametric bias learnt by BPTT can also be regarded as a high-level representation in the neural architecture.

Comparing with the generic RNN networks, with which it is difficult to implement multiple attractor dynamics, the RNNPB is able to generate and recognize multiple temporal sequence patterns by its self-organizing property within an additional layer, called parametric bias units (PB Units). As shown in Fig. 3.6, an RNNPB is essentially a recurrent network (Elman [1990] or Jordan [1997] types) with a set of bias units with adjustable values.

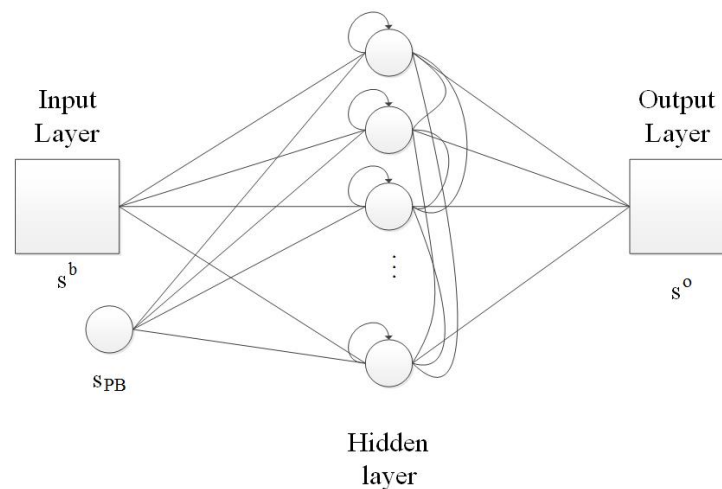


FIGURE 3.6: RNNPB with Elman-like Connections

The parametric bias (PB) units in this recurrent network are connected to the hidden layer as ordinary bias units, but the internal values of them are also updated through back-propagation. Comparing with the generic RNN, the RNNPB owns

the additional PB variables which act as bifurcation parameters for the non-linear dynamics. There exist three running modes in RNNPB. In the learning mode, all connection weights and PB values are updated by BPTT. In the recognition mode, only the PB units are updated. The PB value is manually set in the generation mode.

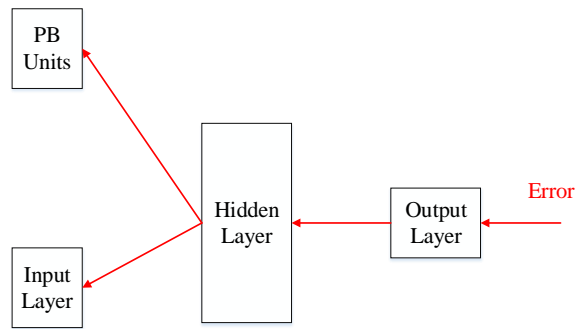
Furthermore, according to Cuijpers et al. [2009], a trained RNNPB can successfully retrieve and recognize different types of pre-learnt, non-linear oscillation dynamics. Thus, this bifurcation function can be regarded as an expansion of the storage capability of working-memory within the sensory system. Furthermore, it adds the generalization ability of the PB units in terms of recognizing and generating non-linear dynamics.

Three running modes (learning, recognition and generation) can functionally simulate different stages between sensorimotor sequences and high-level representation of these sequences. The illustrative demonstration of the model with three modes is shown in Fig. 3.7, where parameters that are modified are denoted in red and the constant weights are denoted in blue.

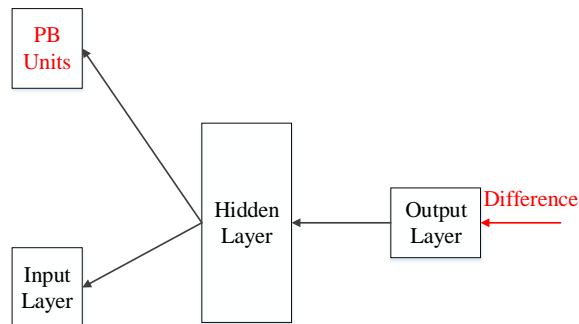
*Action learning mode* (Fig. 3.7(a)): The learning is performed off-line, but the internal values in PB units are learnt *unsupervised*. When providing the training stimulus for each movement pattern, the weights are updated with BPTT (back-propagation through time). Similarly, the internal values of the PB units are also updated in a self-organizing way from back-propagation. If we refer to one entire learning cycle (all sequences) as an epoch  $e$ , in each epoch, the  $k$ th PB unit  $u$  updates its internal value based on the summation of the back-propagated error from one complete sequence.

*Action recognition mode* (Fig. 3.7(b)): This mode recognizes the types of behaviour sequences by updating the PB units according to the past observation. The information flow in this running mode is mostly the same as in the learning mode, i.e. back-propagation, except that the synaptic weights are not updated; rather, the error between target and prediction is only back-propagated and updated into the PB units. If a trained sequence is presented to the network, the activation of the PB units will converge to the values that were previously shown in the learning mode in order to recover the PB values trained before.

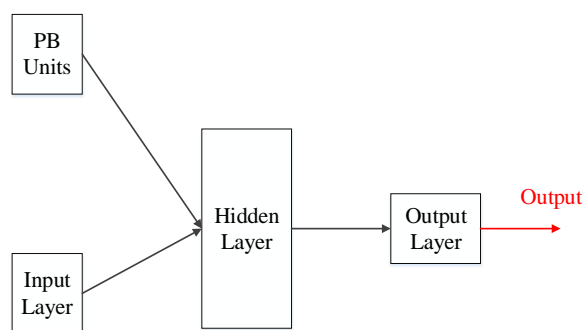
*Action generation mode* (Fig. 3.7(c)): After learning and after the synaptic weights are determined, the RNNPB can act in a closed-loop way, in which the output prediction can be applied as an input for the next time step. In principle, the



(a) Learning Mode



(b) Recognition Mode



(c) Generation Mode

FIGURE 3.7: Three Modes of RNNPB

The figures show different information flow in three modes of RNNPB. Internal values/weights in red will be updated in each mode.

network can generate a trained sequence by providing initial value of the input and externally setting the PB values.

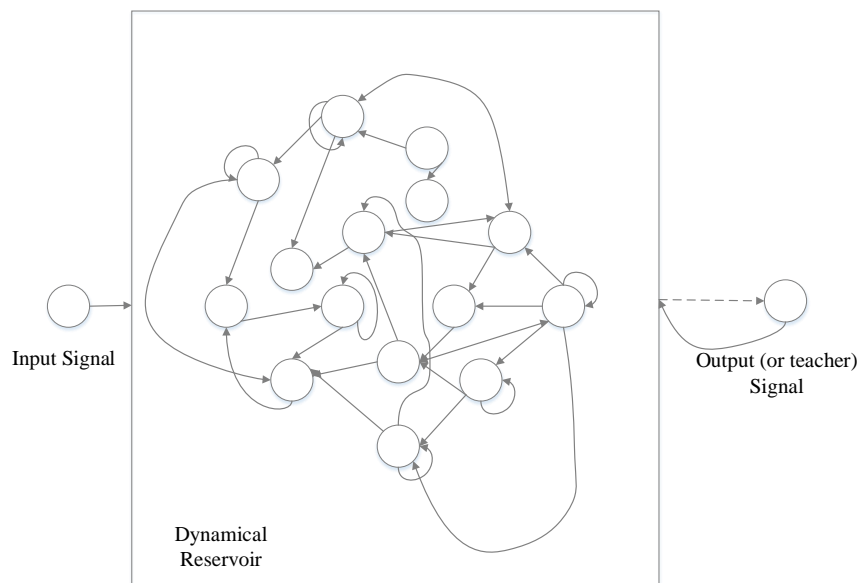


FIGURE 3.8: Echo State Network (ESN)

Hidden neurons inside dynamical reservoir are usually sparsely connected (with typically 1% connectivity). The weights of hidden neurons are fixed. Output weights (dashed) can be trained. Modified from [Jaeger, 2001]

**Echo State Network** Echo state networks (ESN) [Jaeger, 2002] is a kind of recurrent neural networks that have a sparsely and randomly connected hidden layer. Among all the connections, only the weights of output neurons can be changed and trained. Generally, an ESN is created using the following steps:

- A set of neurons with random connections is created to constitute a dynamic reservoir. These artificial neurons can employ any neuron model (e.g., non-spiking leaky integrator neurons are used in the frequency generator example [Jaeger, 2001]). The input neurons are added to the reservoir with randomly assigned all-to-all connections. If the feedback from outputs is required, another set of output neurons must be connected to the reservoir to provide a teacher signal, also with randomly assigned all-to-all connections.
- After that we need to train the reservoir states. If the output signals are presented, a ‘teacher forcing’ is applied to do the training: the error between

the output and the teacher signals are used to modify the output weights between the reservoir and the output. If there are no teacher signals, the neural activity in the reservoir is only driven by the input, which becomes the reservoir states.

- Finally, the output weights are updated by the linear regression (dotted arrows in Fig. 3.8).

Due to the randomly-assigned connecting feedback (weights) within the neurons, the reservoir states require certain asymptotic properties in the weighting matrix. These properties lead to the echo state property which builds up the short-term memory to the internal states. This echo state property is beneficial to some applications, for instance, it has been shown that the ESNs are able to reproduce certain time series as a frequency generator [Jaeger, 2001] or to predict chaotic dynamics in wireless communication [Jaeger & Haas, 2004]. There are numerous studies in the literature concerning the necessary conditions to achieve the appropriate echo state property and maintain the short-term memory (e.g., [Jaeger, 2001, Buehner & Young, 2006]).

**Long Short Term Memory** A long short term memory (LSTM) model [Hochreiter & Schmidhuber, 1997] uses memory blocks to replace the hidden units in a conventional RNN. These memory blocks are called LSTM block (Fig. 3.9). Each of these blocks constitutes an LSTM layer in a recurrent network. This block has the property that can store an internal value for an arbitrary length of time with control by the gated signals. As shown in Fig. 3.9, the recurrent weight of the linear unit is set to be 1.0. If there are no other inputs, this connection serves to preserve the block's current state to the next time-step. Additionally, there are several gates to realise 'memory' function of the block as follows:

- The 'input gate', which is directly connected to the input unit with a product operation, i.e. when the input gate is set to be zero, it wipes out the value from the input unit.
- The 'forget gate', which if it is set to zero, stops the internal memory function. Therefore, it will forget whatever value it was remembering.
- The 'output gate' which determines when the unit should output the value in its memory.

As a single LSTM block can keep an arbitrary long memory, it is suitable to be applied in classification and prediction time series with unknown length where important events are hidden. For instance, an LSTM performs well in handwriting recognition when it requires to scan meaningful characters in the whole writing space with context information [Graves et al., 2009].

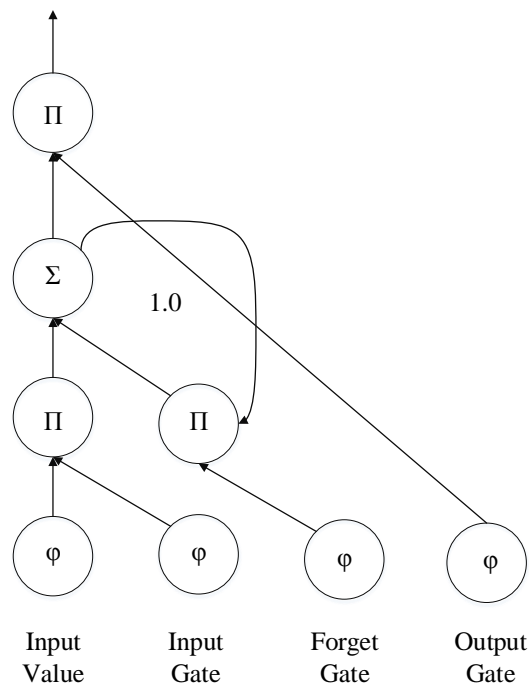


FIGURE 3.9: A LSTM Block

### 3.4 Discussion

Feedback is ubiquitous in the brain. Simple recursive loops and circuit elements in neurobiology also provide the ability to generate patterns of inhibition or excitation.

Basically, a feedback circuit is formed when some neurons' outputs are fed back to the input. There are two types of feedback in neuronal circuits: positive and negative ones. In positive feedback circuits, the effect of the output is to sustain or to increase the activity of the initial input firing. A few neural mechanisms are accounted for by the positive feedback. For instance, a cell firing is caused by the opening of the ion channels that allows the positive feedback and therefore it

causes an increase of the membrane potential. In turn, the increase of membrane potential also causes more current to activate the channels. Such mechanism was first characterised by the well-known Hodgkin-Huxley model [Hodgkin & Huxley, 1952].

The negative feedback circuits aim to inhibit the activity of the initial input realized by a loop of axon, through which the neuron sends suppressing signals to itself. For example, dopamine neurons label sensory stimuli with appropriate values, which are predicted and detected as their own reward signals [Schultz, 1998]. Furthermore, the feedback loop can also be formed between different neurons (interneurons). For instance, the Renshaw cells send inhibitory axon to synapse of the alpha motor neurons.

Comparing with the biological recurrent neural networks (e.g., [Grossberg, 1982, H. R. Wilson & Cowan, 1972, Pilly & Grossberg, 2012]) which attempt to quantitatively interpret the exact electro-chemical activities in neurons and synapses, the ARNNs mainly focus on building a simple model of recurrent connections that can easily encode information of sensory percepts and motor outputs, but also keep the basic ‘neural activity maintenance’ role of recurrent connections. Moreover, although there are several ARNN models as we introduced, in this thesis, we mainly concentrate on the Elman-like models and its variant (RNNPB), because these models are easier to implement in artificial systems and are to simulate sensorimotor functions related to the feedback signals. Also, the short-term memory inside the Elman-like ARNNs is comparable to the working memory in the sensorimotor integration.

### 3.5 Summary

At the neuron and cognitive levels of the brain, feedback signals are utilised to constitute a form of recurrent pathways that interact with and adjust the adaptive filtering function of the bottom-up influence. Followed by the Bayesian inference description of the feedback pathways, in this chapter we apply recurrent connections to model such feedback signals. Recurrent connections in artificial neural networks constitute a directed cycle in the temporal domain. They have been demonstrated to provide a successful short-term memory function in spatio-temporal learnt sequences. Usually, training is performed using back-propagation through time by minimising the energy function in order to reduce the training error. Compared to biologically realistic recurrent network models that focus on



biological mechanisms, the implementation and training of artificial recurrent neural networks are simple and robust. Therefore, in the following chapters of this thesis we will describe the implementation of the feedback pathways in artificial cognitive agents with artificial recurrent connections.

## Chapter 4

# Perception-Action Model with Hierarchical Feedback Pathways

In this chapter, we present a cognitive system framework based on the principle of perception-action model. This framework determines the whole organization of the thesis; the models from Chap. 5 to Chap. 7 describe different parts of this architecture with various kinds of implementations of the feedback information.

### 4.1 Perception-Action Model

The framework of Perception-Action Model (PAM) is based on the common coding theory which advocates that action and perception are intertwined by sharing the same representational basis [Prinz, 1997]. This model asserted that this common representation is simply formed by either the mapping from perception or the perceptual events that actions produce. Note that the representation does not explicitly represent actions; instead, there is an encoding of the possible *future* percept which is learnt from prior sensorimotor knowledge. This perception-action framework is derived from the ideomotor principle [James, 1890], which advocates that actions are represented with prediction of their perceptual consequences, i.e. it encodes the forthcoming perception that is going to happen when an action is executed (i.e. motor imagery) [Greenwald, 1970].

Different from the conventional view that a prior mapping rule should be acquired before the linkage between action and perception (e.g. sensorimotor contingency)

is created, the common coding theory asserts that perception and action may modulate each other directly via the shared coding by a similarity-based matching of common codes. Such a matching does not require a prior knowledge of preceding rules, but it needs only the primitives of sensorimotor knowledge. Therefore, the pairing of perception and action, i.e. the acquisition of ‘common coding’, emerges from prior sensorimotor knowledge. For instance, assuming that one person (called ‘presenter’) is facing another person (called ‘observer’) doing a certain kind of hand movement, according to the PAM model, the corresponding representational domain in the observer about the hand movement should activate, either when the hand movement is observed or the action is executed by the observer itself. Here, both of the current afferent information (referring to the perceived event) and predictive efferent information (referring to intended events from actions) have the same format and structure of a perceptual representation. Specifically, the action being executed is determined by the predictive effects in perception which is caused by the intended action. Thus, in a long term, the acquisition of the common coding from sensorimotor processes is also a learning process for action planning.

## 4.2 Neural Basis of Architecture

### 4.2.1 Somatotopic Arrangement of Motor and Visual Cortices

Based on the evidence of common coding representation [Buccino et al., 2001], our proposed architecture account for the somatotopic arrangement between primary motor and somatosensory cortices. The primary motor cortex (M1) and the somatosensory cortex (S1) are located next to each other in the frontal lobe. These two areas have different functions: the M1 is mainly involved in the execution of voluntary movements by generating neural impulses to activate skeletal muscles, while the S1 is part of the sensory system and is mainly involved with the conscious perception of various sensory modalities, such as touch, pressure and pain. Also, the primary motor cortex has plenty of afferent and efferent connections with the somatosensory cortex, which indicates that the motor commands are partially integrated with the ongoing somatic sensory state of the body from the somatosensory cortex.

Furthermore, the motor control representation in the brain is similar to the perception representation. The motor representation from toe to mouth in brain area M1 is arranged from the top of the cerebral hemisphere to the bottom as

an inverted person. This representation in the cortex is also called ‘humunculus’ (which means ‘little person’ in Latin). From Fig. 4.1, we can see the motor cortex has a large motor representation in dealing with speech and manipulation of objects by the hands, so the humunculus has a large mouth and large hands. This somatosensory cortex also has a similar homunculus representation. Although the exact proportions of the homunculus organization may vary (Fig. 4.1), the general sequences of the stimulation associated with motor actions and sensory perception are similar. The shortest distance between somatosensory cortex and motor cortex forms a convenient link between certain modalities of perception and motor action.

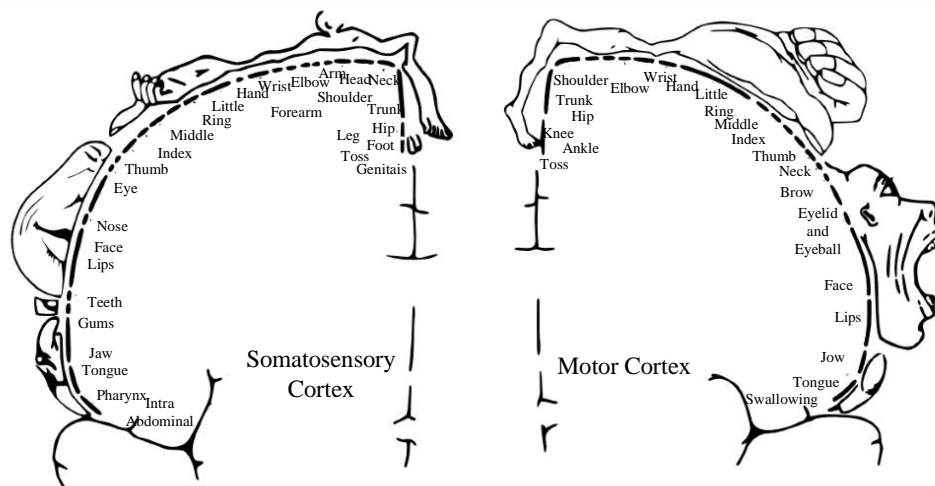


FIGURE 4.1: Homunculus Organization in Motor and Somatosensory Cortices  
Both the motor and somatosensory cortices corresponds point-to-point to one area in the body, in which the fingertip representations in both areas are relatively large, indicating that they are one of the most sensitive/dexterous parts of the body<sup>1</sup>.

Therefore, we propose that parts of the neurons in perception (not first order sensory perception neurons) and part of the motor action should be closely connected, so that the somatosensory cortex maps a continuous representation in some modalities of the primary motor cortex by the neural projection. This is proven by the neurological finding that some neurons in the same brain region

<sup>1</sup>This image is a derivative work based on image <http://commons.wikimedia.org/wiki/File:Homunculus-ja.png>, which is licensed under the Public Domain license.

fire for both execution and observation [Keysers et al., 2010]. This somatotopic arrangement also accounts for the phonetic learning which may involve motor cortex and somatosensory cortex. Some of the motor neurons project to the neurons of the somatosensory cortex with sound receptors, indicating that the executions of articulatory movements of the lips and tongue are closely related to the reproduction of particular lip- or tongue-related phonemes [Pulvermüller et al., 2006].

### 4.2.2 Two Visual Streams for Action

As for the visual system, the representation of the common coding of the dorsal stream emerges during the process of the real-time guidance of action. On the other hand, the low-level visual process is affected by the asymmetric feedback pathways, which guides the predictive dorsal processes. For instance, the neural activities in the visual cortex are predictively modulated by the attention as well as the intended action [Reynolds & Chelazzi, 2004]. Meanwhile, the less direct route is through the ventral stream which maintains visual patterns, allowing the development of visual memories by exploring the novel environment.

Thus, the development of system functions of the visual streams involves the execution of action, which emerges a prior perceptual knowledge about certain actions from the end-manipulator; conversely, this common coding improves both perception and actions by feedback pathways. On the other hand, vision-for-perception and vision-for-action must interact on some levels to accomplish a certain action. For instance, grasping an object with suitable muscle force and movements needs an estimation of the object weights and property, possibly in a semantic form, which are from the ventral stream and the dorsal stream, respectively. Another distinguishing feature provided by the ventral and dorsal streams is that they use a different frames of reference: the ventral stream uses an object-centred frame, while the dorsal stream uses various forms of egocentric frames [Committeri et al., 2004]. Nevertheless, the ventral and dorsal streams cooperate on some levels to voluntarily accomplish a certain action: vision-for-perception (ventral stream) forms an object awareness, mostly according to their features. The actual action planning mainly involves information from the dorsal stream (vision-for-action), but it is also modulated by the representation of object awareness [Schenk & McIntosh, 2010]. Therefore, in our proposed architecture, we also assert the fact that two streams in the visual system are integrated on different levels of the motor hierarchical representation.

### 4.2.3 Mirror Neurons and Ideomotor Principle

As we introduced, the ideomotor principle assumes that actions should be represented as the upcoming perceptual when the actions are taken. The neuroscience finding of mirror neurons in the monkey's premotor cortex (F5) [Gallese et al., 1996, Rizzolatti et al., 1996] and inferior parietal lobule (IPL) [Rizzolatti et al., 2001, Fogassi et al., 2005] also implicitly endorsed this principle by verifying the statement from James [1890]:

*'Every mental representation of a movement awakens to some degree the actual movement which is its object.'*

Indeed, the mirror neuron theory provides the neuroscience evidence of the ideomotor principle in terms of its functional logic in the brain. Moreover, the importance of discovery of mirror neurons indicates that there is a mapping from one's action into cognitive knowledge in an automatic way [Rizzolatti & Craighero, 2004]. The mirror neuron can be considered as a root of language development when it encodes the meaning of action-related words and controls the execution of those actions (e.g. [Hauk et al., 2004, Liberman & Mattingly, 1985]). Also, the discovery of the mirror neurons in premotor and somatosensory cortices supports the ideomotor principle. Based on the PAM model, our architecture also defends the mirror neuron theory; when the cognitive process is consciously intending to execute a motor action, it forms a loop involving perceptual knowledge to drive the muscle movement. Later, the (predictive) perceptual world is involved to maintain the action; together with the bottom-up sensory-driven perception, it determines and updates the perceptual knowledge and the next motor action. Thereby, the mirror neurons fill the gap in a sensorimotor loop, which establishes a dynamical equilibrium between various entities: the mind, the body and the environment [Case et al., 2013].

## 4.3 Sensorimotor Integration Architecture

From the above-mentioned sensorimotor functions, our proposed model is shown in Fig. 4.2. It is mainly based on the common representation framework of perception and action, where the information of the feedback pathways are formed through various levels in the hierarchical organization of perception and action. As a source of the feedback pathways, the common representation domain coding also represents a perception-action linkage between perception, motor imagery, and

action planning. To establish the links between movements ( $a$ ) and their sensory effects ( $e$ ), one needs continuous learning throughout childhood. For instance, the object-directed reaching [Woodward, 1998] and grasping [Rochat, 1987] during the early stages of infant development are considered to be the learning of movements and sensory effects, with consideration of object affordances. Once this link has been established, these perception-action associations in this architecture allow the following operations:

- First, these associations allow to predict the perceptual outcome of given actions by means of the forward models (e.g. Bayesian Model) ( $a \rightarrow e$ ). In the formulation of Bayesian inference which we introduced in Chap. 2, it can be written as

$$P(E|A, I) \propto P(A|E)P(E|I) \quad (4.1)$$

where  $E$  estimates the upcoming perception evidence given an executed action  $A$  and other prior information you have already known ( $I$ ). The term  $P(A|E)$  suggests a pre-learnt model representing the possibility of a motor action  $A$  will be executed given a (possible) resulting sensory evidence ( $E$ ) is perceived (backward computation).

This perceptual prediction also affects low-level activities such as neural activities, which account for the phenomena in perception we mentioned before. The kind of sensorimotor integration proposed in Chap. 5 shows the integration of perception and action in one modality. They share the common predictive representation in form of the recurrent weights which are used to explicitly represent the upcoming visual percept or the percept that is caused from actions.

- Second, these associations allow to select an appropriate movement given an intended perceptual representation. From the backward computations introduced in Eq. 4.2 ( $e \rightarrow a$ ), a predictive sensorimotor integration occurs (Chap. 6).

$$P(A|E, G) \propto P(E|A)P(A|G) \quad (4.2)$$

where  $A$  indicates a particular action selected given the (intended) sensory information  $E$  and a goal  $G$ . Here we assume that one's action is only determined by the current sensory input and the goal.

Note that the above Bayesian inferences are not independent, but they incrementally calculate (deduce) the forthcoming motor action and perception,

and correct its internal model simultaneously, by the linkage of common coding. The complete sensorimotor integration is hereby learnt.

- In terms of its hierarchical organization, it also allows this operation: with bidirectional information pathways, a low level perception representation can be expressed on a higher level, with a more complex receptive field, and vice versa ( $e_{low} \leftrightarrow e_{high}$ ). This can be realised by deep architectures or by our proof-concept model shown in Chap. 7. These operations can be achieved by extracting statistical regularity.

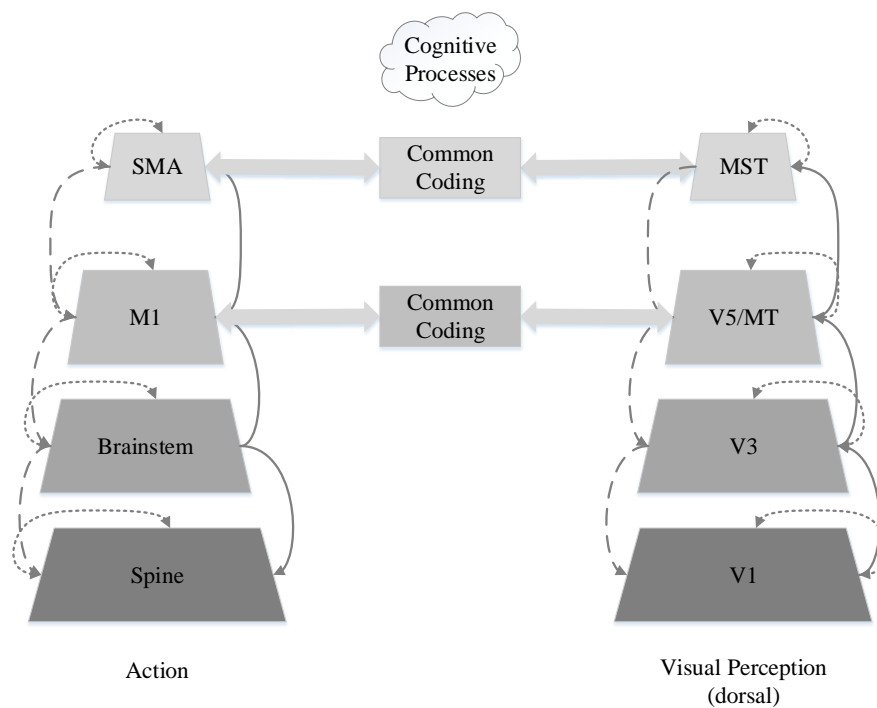


FIGURE 4.2: A Hierarchical Perception-Action Model with Action and Visual Dorsal Stream

To conclude, this framework proposes that as a source of feedback pathways, the common coding domain represents a linkage (the same representation) between perception, motor, and action planning on a higher cognitive level, while there are also feedback pathways that maintain the representation on various levels in both perception and action. We will discuss the details in the following chapters.



## 4.4 Summary

In this chapter, we proposed a cognitive framework with which the different models are going to be discussed in various perspectives in the next chapters. This hierarchical architecture is based on the perception-action model, in which the feedback pathways transmit information in both perception and actions parts. The execution of an action is sharing the same representation with perception, while this common coding is also modulated by the hierarchical feedback signals. Therefore, the representation of common coding is integrated from both sensory-driven and feedback signals. Particularly, the highest-level of perception and action encodes prior knowledge and offers a source of feedback information, which is acquired by the segregated emergence of perception and action.

## Chapter 5

# Feedback-influenced Motion-coding in Visual Cortex

In this chapter, we mainly focus on the modelling of the feedback pathways on two streams in the visual system, with emphasis on its role on the dorsal stream which deals with motion perception. We firstly review the neurobiological background of this theory. Then a brief introduction of the relevant techniques for implementing the two streams for computer vision using neural learning is given. Different from those models, we propose a horizontal product recurrent network model to encode an object's identity and its movement. For the recurrent connections, we claim that the neural activity in the hidden layer is comparable to the observed activities in neurobiological studies. Since there exists a significant neural delay which is caused by the transmission of electrical and chemical signals, our hypothesis is that recurrent connections compensate such neural delay, e.g. by predicting neural activities for motion perception. This chapter is based on our published paper [Zhong et al., 2012b,a].

### 5.1 The Visual System

#### 5.1.1 Two-stream Theory

Mishkin et al. [1983] established the hypothesis that there are two parallel and independent streams in the visual system of humans, in which the 'dorsal pathway' encodes spatial information, invariant of stimulus-specific properties, while the 'ventral pathway' encodes object feature identity, invariant of positions and sizes. These ventral and dorsal streams can also be called the 'what' and 'where' streams

or ‘perception’ and ‘action’ streams [Goodale & Milner, 1992] to some extent. Generally, the visual areas along both the dorsal and ventral streams are organized hierarchically [Livingstone & Hubel, 1988], where the abstractness and complexity increase from lower to higher stages.

As shown in Fig. 5.1, these two streams convey the ‘what’ and ‘where’ information from the visual stimuli in the following specific neuroanatomical areas:

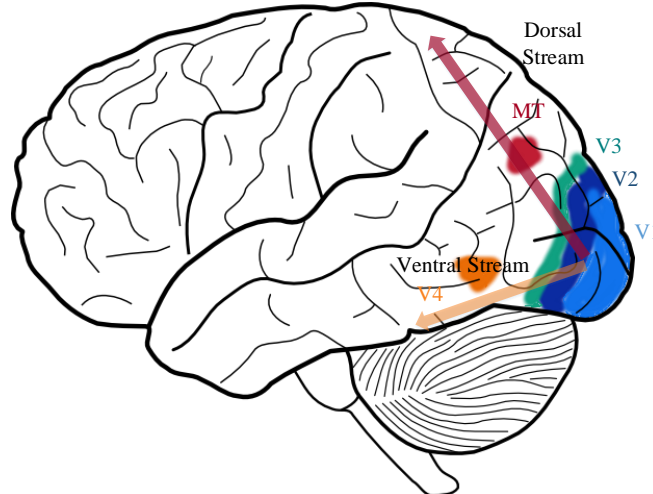


FIGURE 5.1: Anatomy of Two-stream Theory

From V1 in the occipital lobe, visual processing continues in two streams: one into the temporal lobe (ventral stream), one into the parietal lobe (dorsal stream)<sup>1</sup>.

- The two streams first originate from the retina, which turns the lighting signals into electric activity. Specifically, two kinds of retinal ganglion cells exhibit different responses to different properties of visual information: parvocellular cells (P cells) sustain colour-opponent responses while magnocellular cells (M cells) exhibit transient responses to a stimulus. These two types of cells, together with the corresponding parvocellular pathway and magnocellular pathway in LGN, are considered to be the beginning of ventral and dorsal streams.
- The V1 locating in the occipital lobe has a topological arrangement of the visual image. The ventral and dorsal streams start from here and receive information directly from the LGN. There are simple and complex cells in the V1 area. The classification of these cells is based on their responses to

<sup>1</sup>This image is a derivative work based on image <http://pixabay.com/en/brain-human-anatomy-body-155655/>, which is licensed under the Public Domain license.

drifting visual stimuli: it has been widely accepted that both simple and complex cells are orientation selective. Particularly, recent research [Priebe et al., 2006] discovered that the complex cells in V1 are tuned by not only *orientation* but also *speed* of movement. This implies that the dorsal visual stream starts to segregate from V1, encoding both orientation feature and movement information.

- Both streams go through the V2 area, and they further link into higher order visual cortices.
- The ventral stream goes through the higher order visual cortices (V4, IT) in the temporal area, and further conveys information to the inferior temporal cortex. In this cortex, neurons are mostly tuned to be responsive for object shape of intermediate complexity. These areas are considered to have a strong modulation in attention control by visual memory of prior visual salience (e.g. [Desimone, 1996, Bussey & Saksida, 2007]).
- The dorsal stream is transmitted to the dorsomedial area (V3) and visual area MT (also known as V5 in humans) and to the posterior parietal cortex. Specially, MT in the middle temporal lobe is involved in visual motion processing. It is also related to the functions of relaying local motion signals and controlling eye movements (e.g. [Luna et al., 1998, M. Corbetta, 1998]).

Generally, the average receptive field size increases from a lower cortex to a higher one in both streams. Besides, neural response latencies vary in the two streams. In Schmolesky et al. [1998], it was found that most of the cortical visual areas in the dorsal stream show a nearly simultaneous onset of activity for flashed stimuli while the visual response latencies in occipito-temporal areas V2 and V4 are significantly higher. This property of neurons suggests that the variance of visual response latencies in these two streams may facilitate the difference of encoding in two streams. In other words, the shorter latency neural response encodes motion information from the stimuli [Priebe et al., 2006], and the longer latency response encodes visual features from the stimuli. The latency can also be regarded to be significantly involved in the formulation of the recurrent influences by assembling the prior knowledge of the speed and direction (dorsal stream) or features (ventral stream) of the visual stimuli.

Furthermore, the feedback pathways that convey the top-level encoding arouse or sustain the neural activities on the lower neural levels. This could account for various visual functions, such as the compensation of the neural delay in the

visual dorsal stream [Nijhawan, 1994]. Since the neural delay exists in the visual cortex, it could be crucial to eliminate it in certain circumstances. For example, neural delays of one-tenth of a second will cause a large bias in visual perception if a high-speed object is in the visual field. Therefore, the feedback pathways on the dorsal stream play the role of a prediction mechanism, which can account for some visual illusions such as the ‘flash-lag effect’. Also the feedback in the ventral stream transmits the encoding of the features of the stimuli, facilitating the formulation of visual memory. Therefore, the neural responses on higher levels of the two streams can be regarded as the source of such feedback signals, which affects the neural activities on lower levels by means of sustainment or prediction.

### 5.1.2 Identification and Tracking

In the field of computer vision, object identification and pattern recognition have been active topics for decades. However, little attention has been given to encode the object location (as well as movement) as it is straightforward to do the tracking by sliding window. In this section, we will motivate our proposed horizontal recurrent network model by reviewing the related techniques and models of learning ‘what’ and ‘where’ in both computer vision and computational neuroscience.

In the computer vision community, invariant object recognition has been an active topic, which is fundamental in scene recognition, autonomous driving, etc. Existing techniques for invariant object recognition are mainly based on selecting features from visual input (called feature-based techniques) or matching a template (appearance-based techniques) to identify specific objects in the pictures or video clips. For instance, scale-invariant feature transformation (SIFT) [Lowe, 1999] and Speeded Up Robust Features (SURF) [Bay et al., 2006], such techniques search and identify features in the target frame and compare to the pre-learned (pre-defined) ones. From a psychologist’s point of view, these feature-based techniques are compatible to the theory of RBC (Recognition-by-components) [Biederman, 1987], which asserts that the representation of a set of combinations includes the basic elements in the visual fields (called geons, such as cubes, cylinders, wedges, etc). The geons, together with their interrelations (e.g. the relative positions, size of the geons), compose the concept of an object in the brain. The combinations of these basic representations may yield millions of components to represent a real visual object, which are stored in the brain by using structural descriptions, which is a kind of semantic representation. These descriptions are further used for matching when the brain needs to identify an object.

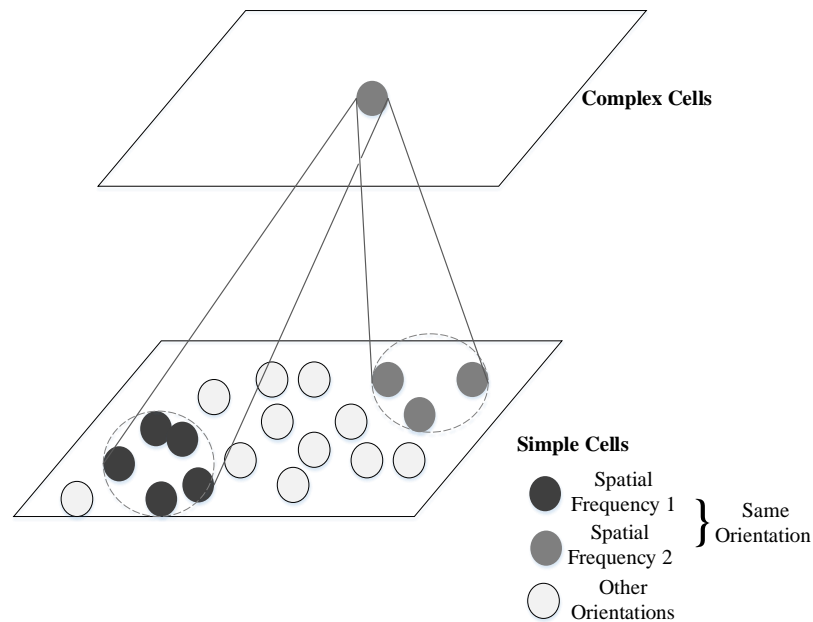


FIGURE 5.2: Example of Pooling of Positional Variation

The pooling method over a group of simple cells allows a small variation from one dimension of the visual stimuli, which may include different positions, orientations, spatial frequencies, etc. Hence it results in the tolerance of such dimension of the visual information. For instance, if neurons with similar orientation and spatial frequency in space, but with different position preferences are separated in the simple cell map, for a complex cell to achieve a position invariant response, it pools the neural responses of simple cells of corresponding positional variation located in the simple cell map.

In the field of computational neuroscience, this problem is formulated as transformation invariance. In terms of transformation invariance modelling, layered networks with simple and complex cell representations [Fukushima, 1980, Földiák, 1991, Hyvärinen & Hoyer, 2000] attempt to learn the transformation invariant perception of objects through self-organizing or by constraints of maximizing sparseness, in which simple and complex cells represent local features and transformation invariant features. These simple and complex cells are arranged in layers one after another so that the degree of transformation invariance gets higher and higher. During learning, a set of images with small transformations are presented, so that similar and localised features in the visual stimulus can be pooled in the complex-cell-like layer(s) (Fig. 5.2). A Hebbian-learning-like method can be adopted to learn transformation invariance by encouraging the neurons to fire invariantly while transformations are performed in their input stimuli [Földiák, 1991, Wiskott & Sejnowski, 2002]. Objects without positional transformation can also be learnt

by a statistical representation of features, such as Restricted Boltzmann Machine (RBM) [Hinton, 2002] and its deep hierarchical architectures [Norouzi et al., 2009, Salakhutdinov & Hinton, 2009].

To conclude, the common feature of these models/techniques for object identification is that most of them utilise localized, oriented or edge filters. Furthermore, some of them employ the idea of hierarchical structures to extract features in order to understand certain visual scenarios, which are considered to be similar to the hierarchical feature extraction in the visual cortex [Livingstone & Hubel, 1988]. Nevertheless, most of these object identification techniques/models ignore the information about the object location.

However, in terms of building an artificial cognitive system based on the principles of biological systems, it is better to preserve all kinds of information gained on various levels of visual perception, especially those involved in action and further cognitive functions. Therefore, it is essential to keep the concurrence of ventral and dorsal streams as these two sorts of information may intertwine again in higher cognitive brain parts, such as the hippocampus in the medial temporal lobe (ML). Thus the techniques which simply disregard object locations are deficient. Therefore, it becomes more attractive for researchers in both neuroscience and computer science to encode object identity and transformation simultaneously. For instance, using bilinear multiplication [Freeman & Tenenbaum, 1997, C. Anderson et al., 2005] it is possible to separate the invariant features and their transformations separately by sparse coding. This is based on an assumption that a transformed image can be represented as a bilinear model of a standard transformation invariant representation of features and the control units representing transformation parameters. In bilinear models of visual routing, a set of control neurons dynamically modifies the weights of the ‘what’ pathway on a short time scale. The control units, encoding the object’s position, thereby route the visual information from any retinal position to an object-centred reference frame on the top-most level of the ‘what’ pathway [Olshausen et al., 1993, Bergmann & von der Malsburg, 2011, Memisevic & Hinton, 2007].

### 5.1.3 Motivation

Considering the above network models, as well as the finding in V1 area of macaque monkey [Priebe et al., 2006], which revealed that some of the complex cells are ‘speed tuned’ in V1, we advocate that the neural modelling of the V1 area should not only constitute the information of both ‘what’ and ‘where’, but also encode

the first-order movement percept. Furthermore, the difference of visual response latencies may be accounted for by the horizontal connections within areas or feedback signals from higher cortical areas [Lamme & Roelfsema, 2000], which is identical to our proposed recurrent connection based on implementation of feedback pathways.

With the idea of solving the ‘what’ and ‘where’ problem jointly, in this chapter we propose a recurrent neural model that can extract two or more components of information into separate pathways from visual stimuli. Unlike previous approaches, our model uses the horizontal product model, which efficiently reduces the computational complexity. This model addresses the following problems:

- to encode *motion direction* and to predict the next visual stimuli;
- to mimic the feedback pathways on the ventral and dorsal streams. Both pathways incorporate recurrent connections to capture different latencies of neural responses of these streams.

## 5.2 Horizontal Recurrent Model

### 5.2.1 Horizontal Product

One way of integrating both of the two separate pathways, as well as motion perception and prediction, is using a horizontal product model. For instance, Köster et al. [2009] applied the horizontal product model together with Independent Component Analysis (ICA). The model was successful to learn to separate the localized image feature and the transformation. One feature of using ICA model to solve bilinear formulation is to reduce the computational effort by decreasing connections because it isolates two pathways. For instance, if there are  $N$  numbers of locations, considering  $M$  possibilities of transformations and  $F$  features,  $F^2 \cdot M \cdot N^2$  connections are needed based on bilinear multiplication. We only need  $F^2N + FMN^2$  connections if we employ a horizontal product (Fig. 5.3). Therefore, inspired by the functional properties of dorsal and ventral pathway neurons, our model segregates the hidden layers into two groups:

- dorsal-like units that encode the (fast changing) current and future object position, and more specifically, the object movement;
- ventral-like units that encode the (slow changing) object feature(s).



We specify a three-layer network with recurrent connections and a horizontal product (Fig. 5.3). The input layer corresponds to the simple cells in V1 (cortical layer IV), while the hidden layer corresponds to complex cells in V1 (cortical layer II and III). The output layer holds a similar object representation of the input, but with a one-step prediction. The hidden layer contains two independent sets of units representing dorsal-like ‘ $d$ ’ and ventral-like ‘ $v$ ’ neurons respectively, inspired by the two visual streams and their functional properties we introduced in the last section. The recurrent connection in the hidden layers helps to predict movement in layer  $d$  and maintain a persistent representation of an object in layer  $v$ . The horizontal product brings both pathways together again in the output layer with one-step ahead predictions. Let us denote the output layer’s input from layer  $d$  and layer  $v$  as  $x^d$  and  $x^v$ , respectively. The network output  $s^o$  is obtained via the horizontal product as

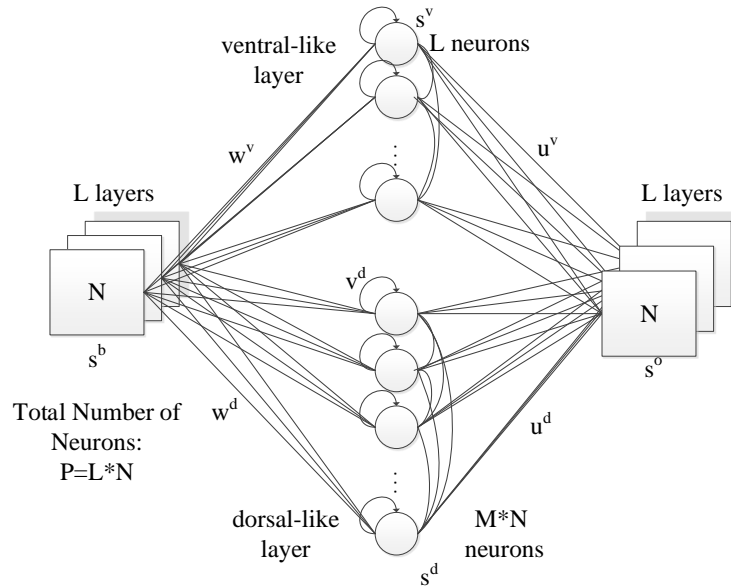


FIGURE 5.3: Horizontal Product Recurrent Network architecture

Due to the complexity, the one-step delayed input is not included in this figure, but it is fully connected to the two hidden layers as the original input  $s^b$ .

$$s^o = x^d \odot x^v \quad (5.1)$$

where  $\odot$  indicates element-wise multiplication, so each pixel is defined by the product of two independent parts, i.e. for unit  $i$  it is given that  $s_i^o = x_i^d \cdot x_i^v$ .

### 5.2.2 Algorithm

**Training** We use  $s_i^b(t)$  to represent the activation of the input unit  $i$  at the  $t$ -th time-step. In some of the following equations, we will omit the time-index  $t$  if all activations are in the same time-step. The hidden units' inputs  $y_j^v$  in the ventral pathway and  $y_j^d$  in the dorsal pathway are defined as

$$y_j^v(t) = \sum_i s_i^b(t)w_{ji}^v + \sum_i s_i^b(t-1)\bar{w}_{ji}^v + \sum_{j'} s_{j'}^v(t-1)v_{jj'}^v \quad (5.2)$$

$$y_l^d(t) = \sum_i s_i^b(t)w_{li}^d + \sum_i s_i^b(t-1)\bar{w}_{li}^d + \sum_{l'} s_{l'}^d(t-1)v_{ll'}^d \quad (5.3)$$

where  $w_{li}^d/w_{ji}^v$  represent the weighting matrices between dorsal/ventral layers and the input layer,  $\bar{w}_{li}^d/\bar{w}_{ji}^v$  represent the weighting matrices between a one-step delayed input and the two hidden layers and  $v_{ll'}^d/v_{jj'}^v$ , indicate the recurrent weighting matrices within the hidden layers (see Fig. 5.3). The incorporation of the time-delayed inputs directly from  $s_i^b$  can introduce more stable input signals in both layers regardless of the short-time changes of object features.

The transfer functions in both hidden layers employ a logistic function and a soft-max function:

$$z_j^v = \frac{1}{1 + \exp(-a_j y_j^v + b_j)} \quad ; \quad z_l^d = \frac{1}{1 + \exp(-a_l y_l^d + b_l)} \quad (5.4)$$

$$s_j^v = \frac{\exp(z_j^v)}{\sum_{j'} \exp(z_{j'}^v)} \quad ; \quad s_l^d = \frac{\exp(z_l^d)}{\sum_{l'} \exp(z_{l'}^d)} \quad (5.5)$$

The logistic function has two local modifiable parameters  $a$  and  $b$ , leading to the intrinsic plasticity of neurons, which we will discuss in the following paragraphs. Together with the soft-max function, the logistic functions ensure the regularity of firing rate on the hidden layer.

The terms of the horizontal products of both pathways can be presented as follows:

$$x_k^v = \sum_j s_j^v u_{kj}^v; \quad x_k^d = \sum_l s_l^d u_{kl}^d \quad (5.6)$$

Again, the output of the network composes a horizontal product from two hidden layers:

$$s^o = x^d \odot x^v \quad (5.7)$$

The training progress is determined by a cost function:

$$C = \frac{1}{2} \sum_t^T \sum_k^P (s_k^b(t+1) - s_k^o(t))^2 \quad (5.8)$$

where  $s_i^b(t+1)$  is the one-step ahead input, as well as the desired output,  $s_k^o(t)$  is the current output,  $T$  is the total number of available time-step samples and  $P$  is the number of output nodes (i.e.  $P = L \cdot N$ ), which equals to the number of input nodes. Following gradient descent, each weight update in the network is proportional to the negative gradient of the cost with respect to the specific weight  $w$  that will be modified:

$$\Delta w = -\eta \frac{\partial C}{\partial w} \quad (5.9)$$

We use the lower indices  $i \in \{1, 2, \dots, P\}$  indicating neurons in the input layer,  $j \in \{1, 2, \dots, F\}$  indicates neurons in the ventral-like hidden layer,  $l \in \{1, 2, \dots, MN\}$  indicates neurons in the dorsal-like hidden layer and  $k \in \{1, 2, \dots, P\}$  indicates neurons in the output layer. The back-propagation training progress is then modified according to the horizontal product as follows.

- Weights from ventral-like hidden layer to output layer:

$$\Delta u_{kj}^v(t) = \eta \delta_k \cdot s_j^v(t) \{(1 - s_k^o(t)) s_k^o(t)\} * x_k^d \cdot a_j^v(t) \quad (5.10)$$

- Weights from dorsal-like hidden layer to output layer:

$$\Delta u_{kl}^d(t) = \eta \delta_k \cdot s_l^d(t) \{(1 - s_k^o(t)) s_k^o(t)\} * x_k^v \cdot a_l^d(t) \quad (5.11)$$

- Weights from current input to ventral-like hidden layer:

$$\Delta w_{ji}^v(t) = \eta s_i^b(t) \cdot s_j^v(t) (1 - s_j^v(t)) * \sum_{k=1}^P x_k^d \left\{ \delta_k (1 - s_k^o(t)) s_k^o(t) \right\} u_{kj}^v \cdot a_j^v(t) \quad (5.12)$$

- Weights from current input to dorsal-like hidden layer:

$$\Delta w_{li}^d(t) = \eta s_i^b(t) \cdot s_l^d(t) (1 - s_l^d(t)) * \sum_{k=1}^P x_k^v \left\{ \delta_k (1 - s_k^o(t)) s_k^o(t) \right\} u_{kl}^d \cdot a_l^d(t) \quad (5.13)$$

- Weights from delayed input to ventral-like hidden layer:

$$\Delta \bar{w}_{ji}^v(t) = \eta s_i^b(t-1) \cdot s_j^v(t)(1 - s_j^v(t)) * \sum_{k=1}^P x_k^d \left\{ \delta_k (1 - s_k^o(t)) s_k^o(t) \right\} u_{kj}^v \cdot a_j^v(t) \quad (5.14)$$

- Weights from delayed input to dorsal-like hidden layer:

$$\Delta \bar{w}_{li}^d(t) = \eta s_i^b(t-1) \cdot s_l^d(t)(1 - s_l^d(t)) * \sum_{k=1}^P x_k^v \left\{ \delta_k (1 - s_k^o(t)) s_k^o(t) \right\} u_{kl}^d \cdot a_l^d(t) \quad (5.15)$$

- Recurrent weights of ventral-like hidden layer:

$$\begin{aligned} \Delta v_{jj'}^v(t) = & \eta s_j^v(t-1) \cdot s_j^v(t)(1 - s_j^v(t)) * \sum_{k=1}^P \delta_k x_k^d u_{kj}^v \cdot a_j^v(t) \\ & + \eta s_j^v(t-2) \cdot s_j^v(t-1)(1 - s_j^v(t-1)) * v_{jj'}^v \left\{ s_j^v(t-1) \cdot \right. \\ & \left. s_j^v(t)(1 - s_j^v(t)) * \sum_{k=1}^P \sum_{l=1}^{MN} s_l^d(t) u_{kj}^d(t) \left\{ \delta_k (1 - s_k^o(t)) s_k^o(t) \right\} u_{kj}^v \cdot a_j^v(t) \right\} \end{aligned} \quad (5.16)$$

- Recurrent weights of dorsal-like hidden layer:

$$\begin{aligned} \Delta v_{ll'}^d(t) = & \eta s_l^d(t-1) \cdot s_l^d(t)(1 - s_l^d(t)) * \sum_{k=1}^P \delta_k x_k^v u_{kl}^d \cdot a_l^d(t) \\ & + \eta s_l^d(t-2) \cdot s_l^d(t-1)(1 - s_l^d(t-1)) * v_{ll'}^d \left\{ s_l^d(t-1) \cdot \right. \\ & \left. s_l^d(t)(1 - s_l^d(t)) * \sum_{k=1}^P \sum_{j=1}^F s_j^v(t) u_{kl}^v(t) \left\{ \delta_k (1 - s_k^o(t)) s_k^o(t) \right\} u_{kl}^d \cdot a_l^d(t) \right\} \end{aligned} \quad (5.17)$$

where  $\delta$  denotes the output error, i.e. the difference between expected (target) output  $d$  and the actual output  $s^o$ :

$$\delta_k = d_k - s_k^o \quad (5.18)$$

The object identity and position information from the input data is distinguished and extracted by the two pathways during training. This distinctive information

coding is induced by the following arrangements: The activations in layer  $v$  are first determined by Eq. 5.4 and 5.5; after that, a constraint is set to the ventral-like units  $v$  so that the states in the following time-steps are forced to be equal to the first time-step as long as the identity of the object remains unchanged. That is, the ventral-like units' activations remain the same until the appearance of a new object. The dorsal-like units, which do not have such a constraint, can update quickly according to the current position of the object.

To ensure a connection between two nodes could be excitatory, the weight matrices between input and hidden layers, between delayed input and hidden layers, and between hidden layers and output, are set to contain only non-negative elements. In a mathematical perspective, it equals to non-negative matrix factorization (NMF), which only allows additive operation.

**Intrinsic Plasticity** When a neuron in the hidden layers employs non-linear transformation function, it is possible for it to maintain equilibrium of an exponentially distributed firing rate resulting in a regular firing in the hidden layer. That is the reason we adopt the intrinsic plasticity in the neurons of the hidden layers based on Weber & Triesch [2008]. With this model, the transfer function of a neuron can basically adapt to fit a sparse exponential regime by adjusting its parameters, slope and threshold. The logistic transfer function here is adjusted with respect to parameters  $a$  and  $b$ . They are updated in order to minimize  $d(f_z || f_{exp})$  which represents the Kullback Leibler divergence between the hidden layer neuron's firing rate distribution  $f_z(z_i)$  and a sparse target distribution  $f_{exp}(z_i) \approx e^{-z_i}$

$$\Delta a_i = \eta_a \left( \frac{1}{a_i} + y_i - 2y_i z_i - \frac{1}{\mu} y_i z_i + \frac{1}{\mu} y_i z_i^2 \right) \quad (5.19)$$

$$\Delta b_i = \eta_b \left( 1 - 2z_i - \frac{1}{\mu} z_i + \frac{1}{\mu} z_i^2 \right) \quad (5.20)$$

where  $\mu$  is the mean for the exponential defined over the positive half-axis. The learning of parameters  $a$  and  $b$  leads to different shapes of the transfer function. The parameter  $a$  controls the gain of the input, changing the slope of sigmoid function, while the parameter  $b$  shifts the non-zero point threshold of the function.

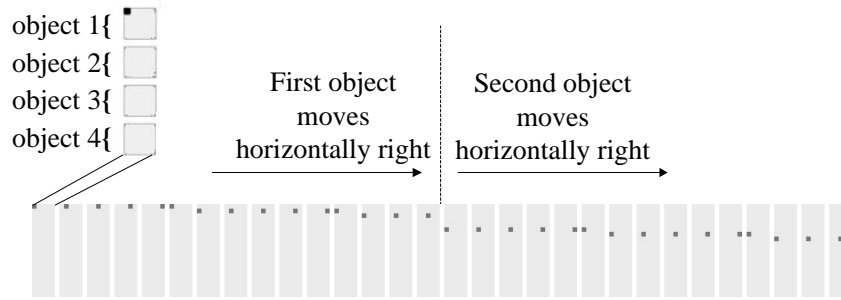


FIGURE 5.4: Partial Sample of Input Data

Objects 1 and 2 (in 1st and 2nd squares) are moving horizontally rightwards.

### 5.3 Experiments

In the experiment here, we present artificially generated input data to the network to simulate the neural coding in the cortical column; here the data set is generated in a formulation of various two-dimensional matrices. One and only one element in only one matrix is activated, which represents one object feature. The moving of the activation mimics the moving objects, i.e. their positions change quickly, but their features change rarely. In this way, the learning of this network is demonstrated through showing different objects in various layers of inputs. Generally, assume that the input images are composed of  $k$  layers which represent  $k$  various kinds of objects. In each layer, there are  $m \times n$  positions, in which the object moves horizontally or vertically. The training data set comprises a complete horizontal moving activation for both directions. It covers all of the possible horizontal movements from all of the positions including all objects. For instance, the first data set contains an activation in the first layer moving from coordinates  $(1, 1)$  to  $(1, 2), (1, 3), \dots$ . These movements vary in different starting points and different layers. Fig. 5.4 demonstrates an example of one part of the input sequence, in which there are two kinds of objects, represented in the first and second squares of the input block. Both objects are moving horizontally right. During the training process, the target data is the one-step ahead instance from the training data.

The network parameters are shown in Tab. 5.1, the maximum iteration is set to 100,000. In order to learn movement appropriately, activation in both of the hidden layers is reset to zero after the activation of changes between objects.

With the input in Fig. 5.5, Fig. 5.6 shows the corresponding one-step ahead prediction. We can generally observe that the output layer predicts the one-step ahead

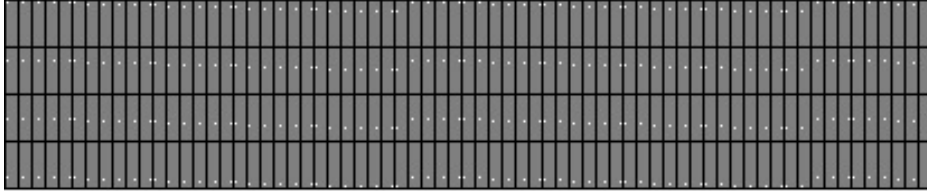


FIGURE 5.5: Input of Activation Movement

Each row represents a complete movement of one object in one direction. In this partial view of input, we only show horizontal movement from left to right.

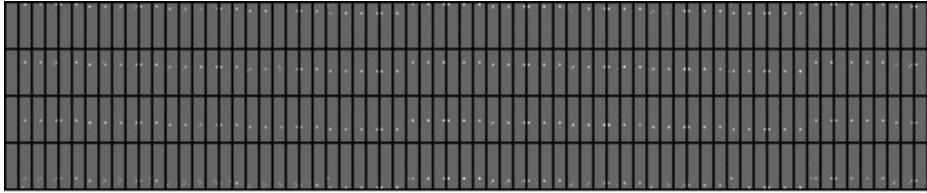


FIGURE 5.6: Corresponding Output from Fig. 5.5

Parameters	Description	Value
$L \times N_1 \times N_2$	Size of Input Layer	$4 \times 5 \times 5$
$M$	Number of Movement Directions	4
$\mu_v$	Parameter of Intrinsic Plasticity in Ventral Stream	0.1
$\mu_d$	Parameter of Intrinsic Plasticity in Ventral Stream	0.01
$\eta_a$	Learning Rate of $a$ in Intrinsic Plasticity	0.0001
$\eta_b$	Learning Rate of $b$ in Intrinsic Plasticity	0.0001
$\eta$	Learning Rate of Weights	0.01
$\epsilon$	Minimum Error Decreasing as a Stopping Criteria	$10^{-8}$

TABLE 5.1: Network Parameters (Horizontal Product RNN)

movement. Note that the output is inactive in every first time-step since the recurrent and time-delayed connections require the previous inputs which are not available in the first time-step. As depicted in Fig. 5.7, activations in the corresponding activations of hidden layer  $v$  stay stable when one object appears, while we can distinguish various patterns in the dorsal-like layer  $d$  representing perceptions of different positions. The training error over the course of learning is depicted in Fig. 5.8. The stopping criterion ( $\text{output error}(t-1) - \text{output error}(t) < \epsilon$  or output error is increasing) was achieved by around iteration 3100.

On the other hand, the activations in the dorsal pathway, including the hidden layer  $d$ , horizontal product with weighting matrix, fluctuate with the changing of the object position. In particular, we can tell there are different patterns in the hidden layer  $d$  while the object is moving horizontally right or left.

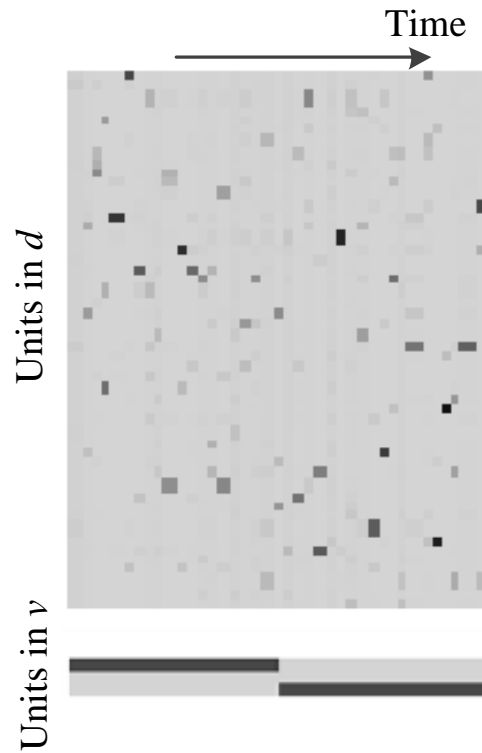


FIGURE 5.7: Network Activations in Hidden Layers

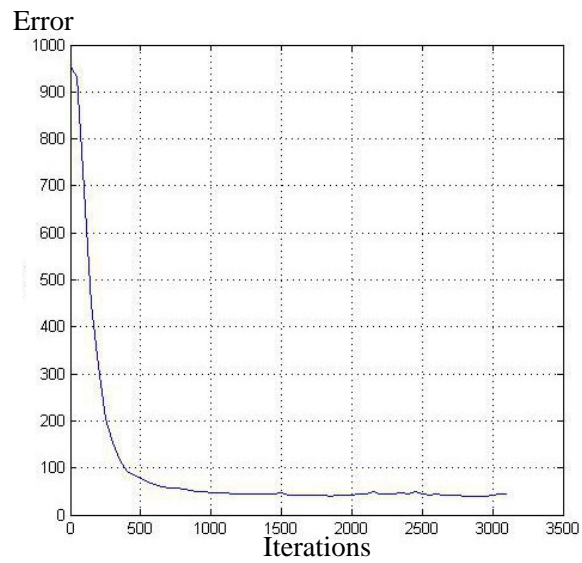


FIGURE 5.8: Output Error through Iterations



## 5.4 Summary

In this chapter we introduce a recurrent network architecture of modelling feedback pathways on two streams in the visual system. Advocating the concept of isolating ‘what’ and ‘where’ in ventral and dorsal pathways, a prediction in the dorsal pathway is also encoded due to the recurrent connection.

The experimental results show that information of object identity and position have successfully been separated: the activation of the ventral pathway remains stable when presenting the same object, while in the activation of the dorsal pathway, especially, there are different patterns appearing in the hidden layer  $d$  indicating different kinds of movement directions. This result can be comparable to the finding of the difference of latency of neural responses in two streams.

This model demonstrates an example of how the feedback signals (e.g. lateral connections) affect neural activities in different parts of the visual system. The role of recurrent connections in the visual system has been stressed in terms of their transmission as a kind of feedback signals. In the modelling perspective, it is inevitable to apply recurrent connections within the dorsal pathway because its short-term memory stores the past movement information. Furthermore, conventional back-propagation training through the horizontal product model is able to isolate the ‘what’ and ‘where’ information into two pathways, to represent the object identity and movement direction respectively, and then to couple them again together in the output layer by the horizontal product, being able to predict movement in the next time-step.

## Chapter 6

# Feedback Signals on a Predictive Sensorimotor System

In the last chapter, we introduced the prediction in the perception of motion. In this chapter, we augment this idea with a motor action module and propose that this kind of predictive mechanism is also beneficial in motion, particularly in artificial cognitive systems.

The prediction is achieved by tracking a moving object by observing saliency and predictively coding the evidence of preferred visual evidence, so that the upcoming sensory data is predicted by the feedback pathways based on the prior sensory information. We claim that with the predicted sensory information, the sensorimotor integration reacts smoother and faster.

This chapter is mostly modified from our published paper [Zhong et al., 2012c].

### 6.1 The Sensorimotor System

#### 6.1.1 Sensory Prediction

In the sensorimotor cycle of an artificial cognitive system, especially when the artificial system is a robot which physically interacts with the environment, there usually exists a temporal delay mainly contributed by the processing time of sensors, transmission time of signals and mechanical latency. For example, because few object recognition programs can recognize the identity of human faces from

visual inputs in real time, the running speed of human-following behaviour based on object recognition should be slower than a normal walking speed of human beings. It is difficult for a robot to keep searching a face in consecutive images within a short time scale. A simple predictive mechanism, such as Kalman filters, can solve this problem by predicting the movement of a person as soon as he/she has been identified (e.g. [Foresti, 1999]). However, since a Kalman Filter is based on a linearity assumption, it does not consider very complicated movements with e.g. non-linear influence from context. Such problems may be tackled by neural networks which can learn to predict percepts in a general dynamic environment. Sensory prediction is of great benefit to dynamic robot behaviours such as obstacle avoidance, visually guided navigation, reaching, visual search, and rapid decision-making under uncertainty, since these kinds of behaviour highly rely on current sensory information. In these scenarios with non-linear dynamics, a developmental sensory prediction is needed to learn to compensate for the latency of the sensorimotor cycle.

A second reason for employing the predictive mechanism is that the sensory inputs of artificial cognitive systems are often noisy and inaccurate in a real environment, which may lead to failure of robot behaviours. In that case, a predictive sensory module can compare its prediction based on previous short-term sensory percepts to the current sensory value. A noisy sensor value can thereby be identified and an action adjustment executed. A severe case may result from sensor failure caused by hardware problems or a change of the environment (e.g. lighting conditions). In such cases, an embodied predictive sensory module can act as a filter to recursively estimate the incoming percepts.

### **6.1.2 Tracking and Prediction**

A well-known method for tracking and prediction used for artificial systems are the Kalman filters [Kalman, 1960]. They are a set of equations built on linear operators and model a POMDP to estimate the state of a process with Gaussian noise. Despite that the conventional Kalman filters are based on the assumption of a linear dynamic system, other improved/adaptive Kalman filters have been proposed [Bonato et al., 2009, Klein et al., 2012] to avoid such a limitation.

Particle filters are also one set of the techniques to deal with non-linear/non-Gaussian tracking problems. As their name implies, they apply a set of particles to represent the posterior density in the state space by recursively calculating the Bayesian inference. They have been extensively used in object tracking (e.g.

Connection	Mean	Variance
LAN	366.0	84.5
WIFI	797.7	187.1

TABLE 6.1: Time delay in the Camera-arm Cycle of NAO (in milliseconds)

[Schulz et al., 2001, Yan et al., 2011]) as well as other robot perception problems (e.g. [Thrun, 2002, Grisetti et al., 2007, Zhong & Fung, 2012, Zhong et al., 2010]).

Neural networks can learn universal function approximation and thereby optimally predict non-linear data [Schaefer et al., 2008, Möller, 2012, Hirel et al., 2011, Saegusa et al., 2007]. For example, a simple feed-forward network approximates a time step occurring in a Hidden Markov model (HMM) to estimate an agent’s position along a motion sequence [Thrun, 1998]. Furthermore, a neural network with recurrent connections is able to predict the movement trajectory recursively, since it also represents the recent data in its short-term memory, it exhibits smooth and stable neural dynamics.

### 6.1.3 Motivation

We have argued in the last chapter that any cortical area should compensate for its own processing delay via the feedback pathways. In this chapter, we will implement this idea on a humanoid robot: the NAO robot<sup>1</sup>. It is designed as an autonomous, programmable humanoid robot. However, as a disadvantage also found in other robots, a sensorimotor cycle latency also exists in NAOs. To quantify this latency, an experiment was conducted by capturing two time-stamps between issuing a movement command and perceiving this movement from the robot camera. Based on ten recordings in each case, the delay of NAO’s sensorimotor cycle is around 0.8s with a wireless connection and 0.36s with a LAN cable between the control PC and the robot (Tab. 6.1). This long delay was observed despite the fact that the used QVGA resolution allowed up to 30 frames per second in the LAN modus. This further motivates to realize sensory prediction to compensate for this delay, which is realized by a predictive sensorimotor model:

- Prediction within an autonomous cognitive robot can happen in both perception and action parts, but in this chapter we only consider the prediction in the perception, i.e. a system predicting sensory signals given the current and previous sensory states. Nevertheless, such prediction is also valuable

<sup>1</sup><http://www.aldebaran.com/en/humanoid-robot/nao-robot>

for the motor action part since it results in a smoother and faster motor action.

- Such local sensory prediction may be easier to implement than prediction of the entire system response, and enable the mixed hierarchical and parallel processing in the (visual) cortex, including short-cut connections, while keeping the representations in all areas temporally aligned.
- Prediction of restricted sensory percepts (instead of e.g. the next action) may be generalizable to many contexts and actions.

#### 6.1.4 Experiments Setting

Since prediction of the complete raw sensory percept (e.g. all pixels of the camera image) is not desirable and would be very difficult for an autonomous robot, it is advisable to predict only few features extracted from the sensory percepts as human perception does [Darrin et al., 2004, Natale et al., 2007]. This can be implemented as a learnt non-linear mapping of sensory representations to predict the forthcoming sensory flow.

Again, this work shows that motion perception in early visual processing is affected by the feedback pathways by modelling it with recurrent neural connections in a robotic system. Furthermore, such implementation is also practical to be applied as a predictive sensorimotor model as part of a developmental robotic system. In our proposed design, we also emphasize that the motor action of the robot should be executed based on the predicted sensory percepts, reacting to the forthcoming sensory data. In this way, it enables faster and smoother robot behaviours. Prediction of motion generally includes both active and passive movement, i.e. movement of sensory events generated by action of the agent itself and those generated by the movement of the observed objects [Cullen, 2004]. In this chapter, we concentrate on the prediction of active movement, inferred from the perceived visual motion of a fixed landmark. However, since the prediction does not use additional input from any behaviours or motor commands, this architecture can predict sensory percepts, no matter whether they are caused by active or passive movement. In statistical notation, we would regard the prediction process as an HMM rather than a partially observable Markov decision process (POMDP), since the action is unknown.

In the work by Navarro-Guerrero et al. [2012], the authors have realized a goal-directed behaviour based on reinforcement learning by a NAO humanoid robot.

The resulting behaviour, however, did not look natural compared to the walking patterns of humans. In this chapter, therefore, we make several new contributions. First, instead of a discrete representation of state- and action spaces, we use continuous representations in both, which results in more sophisticated walking behaviours. This is made possible by the Continuous Actor-Critic Learning Automaton (CACL) [van Hasselt & Wiering, 2007]. The continuous action space facilitates generalization of learnt actions to unlearnt regions in the state space, which could help speed up learning and optimize the action selection. This reinforcement learning algorithm is similar to motor adaptation according to error-prediction in biology [Shadmehr et al., 2010, Izawa & Shadmehr, 2011].

Another property of our previous work was the robot's slowness. Instead of using prediction, the visual percept was retrieved after a short waiting time in which the robot stood still to obtain a clear camera image. A behavior that now takes little more than a second with our proposed model (see results below) took roughly half a minute due to the still standing periods [Navarro-Guerrero et al., 2012].

Corresponding to the integration of predictive visual information and smooth motor action, two modules, the sensory and motor modules, are incorporated in this architecture (Fig. 6.1). Within the sensory module, an Elman network is applied to predict the upcoming sensory information, while the motor-action module uses a network trained by CACL [van Hasselt & Wiering, 2007]. Integration of these two modules drives a NAO robot to approach a target position with smooth and continuous behaviour in an autonomous manner. Also, the NAO robot hardware and control commands allow continuous walking without stopping, so that the robot can change the parameters of walking, i.e. speed, walking direction and torso orientation, at any given time while walking. Besides, the light weight of the robot guarantees that it reacts fairly fast.

## 6.2 Recurrent Prediction Sensorimotor Model

### 6.2.1 Landmark-based Detection

The NAO robot is endowed with circular NAOMarks landmarks and a built-in detection routine. However, this closed-source recognition program costs quite a lot of computation power and causes significant delay. Besides it leads to over-estimation of the landmark size if the images are blurry, which leads to wrongly estimated poses. Therefore, we designed an own landmark to identify the position

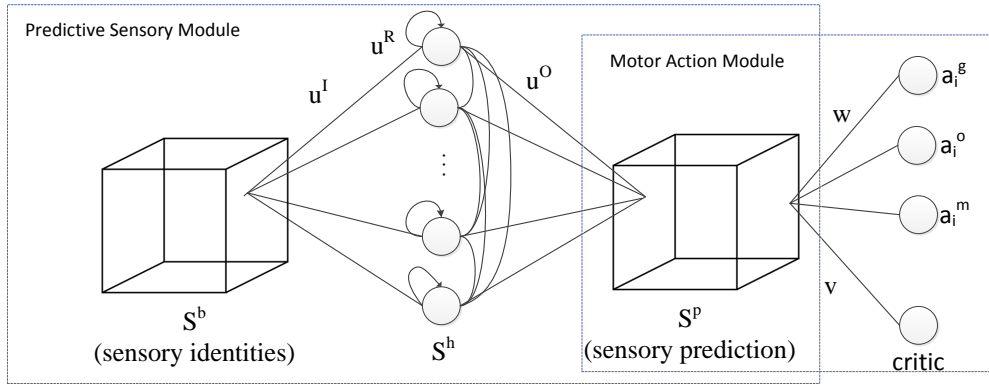


FIGURE 6.1: Overall Architecture Combining Sensory Prediction and Sensorimotor Action

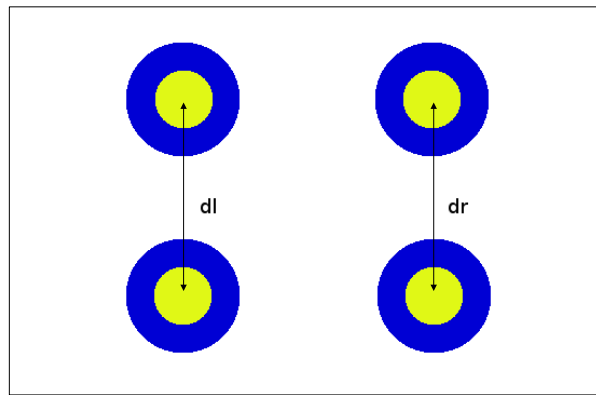


FIGURE 6.2: Sample of the Landmark

The perceived vertical distances between two circle centres from the NAO camera are denoted as  $d_l$  and  $d_r$ , from which we calculate three identities based on Eq. 6.1.

of the approach target. It consists of four circles, each including a large blue ring with a smaller yellow circle inside (Fig. 6.2). We detect the positions of the circle centres in the robot's visual field using 2D Gaussian Fourier filters in RGB channels. The colour combination of the landmark is different from the wooden docking station to be easily distinguishable. Then, Hough circle detection finds circles within a certain radius range [Ballard, 1981]. Our routine is faster than the previous NAOMarks, the total pre-processing and searching time of the landmark data is less than 0.01s.

Then the position and orientation of the robot with respect to the landmark can be defined by the following three values:

$$I_1 = \frac{d_l + d_r}{2}, \quad I_2 = \frac{d_l - d_r}{I_1}, \quad I_3 = \sum_x pix_x \quad (6.1)$$

where the measurements  $d_l$  and  $d_r$  are defined according to Fig. 6.2.  $pix_x$  is the summation of the  $x$  coordinate of all the four circle centres within the robot's visual field. Referring to the overall installation of the shelf in Fig. 6.3, the first identity  $I_1$  correlates with the proximity between the robot and the landmark. The second identity  $I_2$  correlates with the angle of the robot's position w.r.t. the landmark, while the third identity  $I_3$  informs about the robot's orientation w.r.t. the horizontal direction to the landmark.

The three values of Eq. 6.1 contain all position/orientation information relevant for the robot's approach behaviour. Neurons in the input layer are arranged as a three-dimensional cube, in which they are activated with a Gaussian activation blob that is centred around  $(s_1^c, s_2^c, s_3^c)$  defined as:

$$s_n^c = \frac{I_n - I_n^{min}}{I_n^{max} - I_n^{min}} \times N_n \quad (6.2)$$

where  $n$  ( $n = 1, 2, 3$ ) is the  $n$ -th dimension of the cube, corresponding to the  $n$ -th identity of Eq. 6.1.  $I_n^{min}$  and  $I_n^{max}$  are the minimum and maximum value of the identities data, respectively.  $N_n$  is the number of neurons in the  $n$ -th dimension. The activation of neighbouring neurons  $s^b(x_1, x_2, x_3)$  is distributed according to a Gaussian:

$$s_i^b(x_1, x_2, x_3) \sim \mathcal{N}(s_n^c, \delta_m) \quad (6.3)$$

These values define the activation on the sensory input layer.

### 6.2.2 Algorithm

**Visual Prediction via Recurrent Connections** Similar to Chap. 5, the predictive module consists of a three-layer Elman network. Inputs to this network are the three observed identities of the landmark from the perceived images  $\{s^b\}$ . Outputs are the one-step ahead predictions  $\{s^p\}$  (c.f. Fig. 6.1).



$s_i^b(t)$  denotes the state of input neuron  $i$  at the  $t$ -th time-step. The activations of the hidden units  $y_j$  at time  $t$  are defined as

$$y_j(t) = \sum_i^{N_1 \times N_2 \times N_3} s_i^b(t) u_{ji}^I + \sum_i^{N_1 \times N_2 \times N_3} s_i^b(t-1) \bar{u}_{ji}^I + \sum_{j'}^{N_h} s_{j'}^h(t-1) u_{jj'}^R, \quad (6.4)$$

where  $u_{ji}^I$  represents the weight matrix between input layer and hidden layer,  $\bar{u}_{ji}^I$  represents the weight matrix between the time-delayed input and the hidden layer and  $u_{jj'}^R$  indicates the recurrent weight matrix within the hidden layer. The above equation shows that the hidden layer is connected to the weighted stimuli of the current input and the delayed input, while there are additional lateral connections.

The transfer function of the hidden layer is the logistic function,

$$s_j^h(t) = \frac{1}{1 + \exp(-\beta y_j(t))} \quad (6.5)$$

The  $k$ -th output for sensory prediction  $s_k^p(t)$  at time  $t$  is

$$s_k^p(t) = \sum_j^{N_h} s_j^h(t) u_{kj}^O \quad (6.6)$$

where  $u_{kj}^O$  are the weights to the output layer (cf. Fig. 6.1).

**Smooth Action Generation** In the CACLA-trained reinforcement learning network, the input layer encodes the predicted sensory states  $s^p$  and the output layer encodes a critic value and three robot action commands. Two of these represent the moving direction and torso orientation change, and are activated linearly as

$$a^{m/o/g}(t) = \sum_k s_k^p(t) w_k^a \quad (6.7)$$

where  $w_k^a$  is the weight matrix between the sensory state units  $s_k^p(t)$  and the action units  $a^{m/o/g}(t)$ . The third action unit signals the robot to stop walking (to initiate a possible grasping action) and is activated by a sigmoid function:

$$grasp(t) = \frac{1}{1 + \exp(-a^g)} \quad (6.8)$$

A value of  $grasp(t) > 0.5$  in this unit causes the robot to stop, while a value  $< 0.5$  lets it continue the docking behaviour.

Parameters	Description	Value
$\delta_s$	Variance of Gaussian Distribution in Sensory Module	1.0
$N_1 \times N_2 \times N_3$	Size of Input Layer in Sensory Module	$10 \times 10 \times 10$
$N_h$	Size of Hidden Layer in Sensory Module	1,500
$\eta$	Learning Rate in Sensory Module	0.1
$\beta$	Slope in Logistic Function	1.0
$\delta_m$	Variance of Gaussian Distribution in Motor Module	1.0
$\eta_w$	Learning Rate in Motor Action Module	0.4
$\gamma$	Discount Factor of Reinforcement Learning	0.8
$\epsilon$	Decay Rate of Reinforcement Learning	0.5

TABLE 6.2: Network Parameters (Predictive Sensorimotor Integration)

A critic unit guides reinforcement learning. It is activated as

$$c(t) = \sum_k s_k^p(t) v_k \quad (6.9)$$

where  $v_k$  are the weights between the sensory state units  $s_k^p(t)$  and the critic.

The action weights  $w^a$  are updated by the following rule:

$$w_j^a(t+1) = w_j^a(t) + \eta_w \delta_a a^{m/o/g} s_j \quad (6.10)$$

where  $\delta_a$  is the action output error. According to CACLA by van Hasselt & Wiering [2007], this update is only performed while the critic value (Eq. 6.9) is increasing. The updating of the critic weights is defined by:

$$v_{ij}(t+1) = v_{ij}(t) + \epsilon \delta_p s_j \quad (6.11)$$

where  $\delta_p$  is the prediction error. When the reward is achieved, it is defined as

$$\delta_p = r - c(t) \quad (6.12)$$

while the reward is not yet achieved,

$$\delta_p = r + \gamma c(t+1) - c(t) \quad (6.13)$$

where  $\gamma$  is the temporal discount factor. An important difference of CACLA and conventional actor-critic learning is that the action weights  $w^a$  are only updated if the state value is increased, but not if it is decreased, since the optimum may exist between the current selection and the executed one due to the continuous encoding. Tab. 6.2 shows the training parameters used in the two network modules.

## 6.3 Experiments

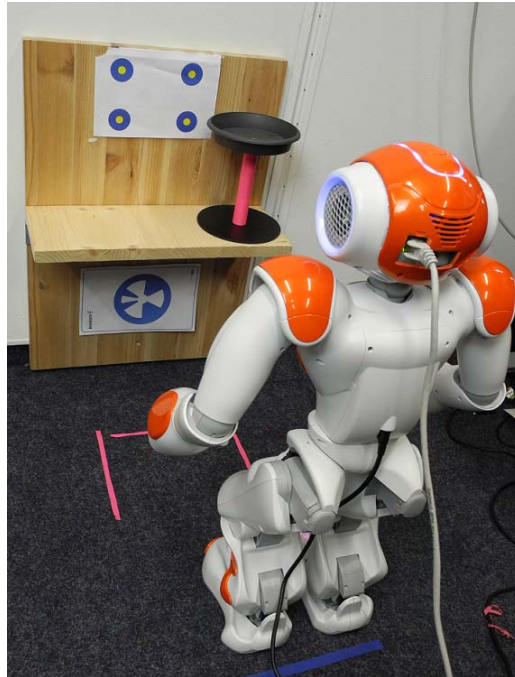
Previously Navarro-Guerrero et al. [2012] have studied that, autonomous robot approaching based on reinforcement learning can serve as a basis for different kinds of robot behaviours, such as grasping, human-robot-interaction, re-charging, etc. We test our predictive sensory model in our home lab, where a shelf with a landmark is installed (Fig. 6.3). The goal of autonomous docking is to approach a narrow area which allows the robot to execute the grasping behaviour afterwards. Grasping will be controlled by a self-organizing map with supervised control output, which ensures robustness and tolerance towards the position and pose reached by the docking: as long as the object (a plastic goblet) is visible in a certain area of the visual field, the robot is able to grasp it.

In our scenario, we define an area of  $22\text{cm}$  by  $13\text{cm}$  in front of the shelf as the optimal stopping area within which both feet of the robot must halt after approaching (red square in Fig. 6.3(b)) based on the kinematics and dimensions of the NAO robot for the grasping behaviour. The largest distance is limited by the robot's arm length, while a too short distance to the shelf leads to the robot's arm being blocked by the shelf when raised. A larger area for the start of the approach is constrained by the requirement that the landmark must be visible within the robot's visual field.

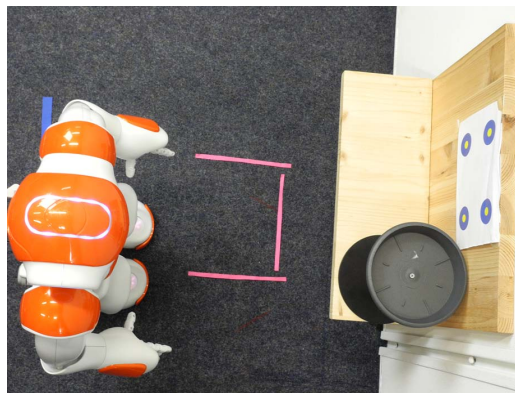
### 6.3.1 Training Scheme

First, training sequences of the robot were collected through manual control in real world experiments, which leads to supervised learning of the forward model and a form of supervised reinforcement learning for the action model [Navarro-Guerrero et al., 2012]. We avoided the use of a robot simulator since it is very difficult to configure it to reflect the exact physical parameters of the real world, such as the camera optics or friction on various carpets. With real world training, these factors are learnt without explicit model. The Gaussian activations (Eq. 6.3) speed up training and a large training set can be avoided.

Since the NAO robot has many degrees of freedom, for simplicity we keep the robot pose constant except for leg movements to keep the number of action units small. The training sequences are not recorded under continuous walking, but step-wise; in every step, the following data is recorded: the sensory identities from the landmark, the robot movement direction, a change in torso orientation and the stop action for grasping. Based on the hardware constraints of the NAO robot,



(a) Back view of the docking station



(b) Top view of the docking station

FIGURE 6.3: Shelf Installation / Docking Station

The robot tries to approach the square mark. After docking, the robot could grasp some object from the shelf.

the walking direction is between  $-\pi$  and  $+\pi$  and the torso orientation change is from  $-0.1$  to  $+0.1$  (both in radians). Furthermore, to keep the dimensionality of the continuous action space limited, we try to keep the consistency of the step-walk-distance and the continuous-walk-distance covered by the sampling time, i.e.  $\Delta D = 3cm$ .

Methods	# of Trials	# of Success	Success Rate	Avg. Time	Variance
Without	25	16	0.640	1.57	0.75
With	25	20	0.800	1.22	0.51

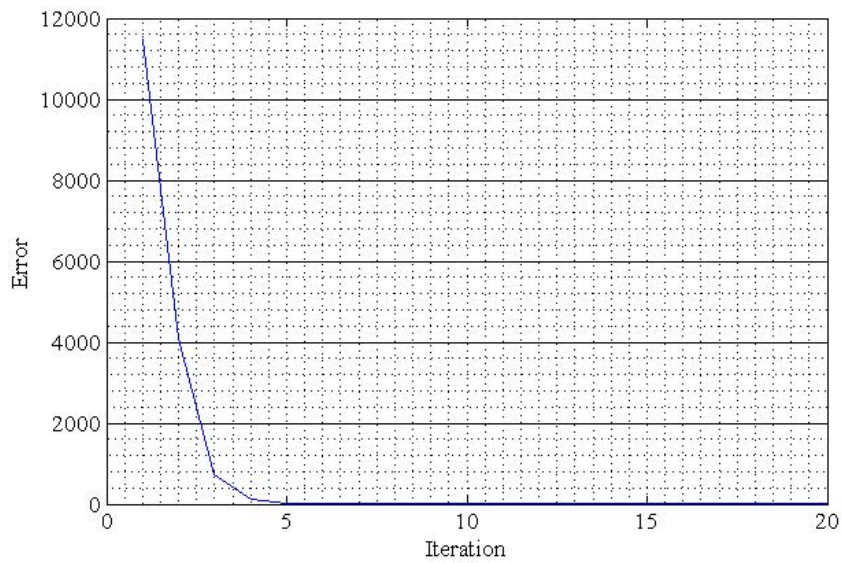
TABLE 6.3: Docking Trials by Reinforcement Learning with and without Prediction

The approaching can start at any point within the approach area. For the reason of spatially balancing weights representing the walking data in the training sequence, we carefully selected the starting points as four points in the middle, six points in the left half and six points in the right half of the approaching area. We chose more training sequences of walking sideways because they are more challenging.

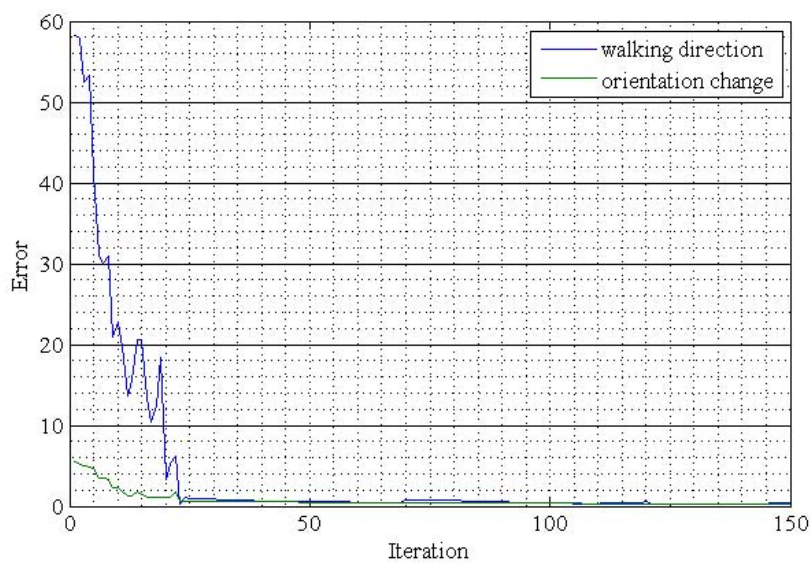
The recorded training sequences were used to train both modules off-line using the rules described above. The training in the sensory predictive module is identical to the conventional back-propagation through time (BPTT) algorithm. Fig. 6.4 shows the learning curves representing the output error in both two modules. The figures show that both modules quickly converge before the 50th iteration.

### 6.3.2 Experimental Results

**Approaching based on Reinforcement Learning without Prediction** As soon as the first training procedure of CACLA was done, the NAO robot was able to approach the shelf in a continuous way. Connecting the sensory input directly to the motor action module, we conducted twenty-five approach trials without a predictive model with LAN connection. Results are shown in Tab. 6.3. We count a trial successful if the subsequent grasping behaviour leads to a successful grasp of the goblet.



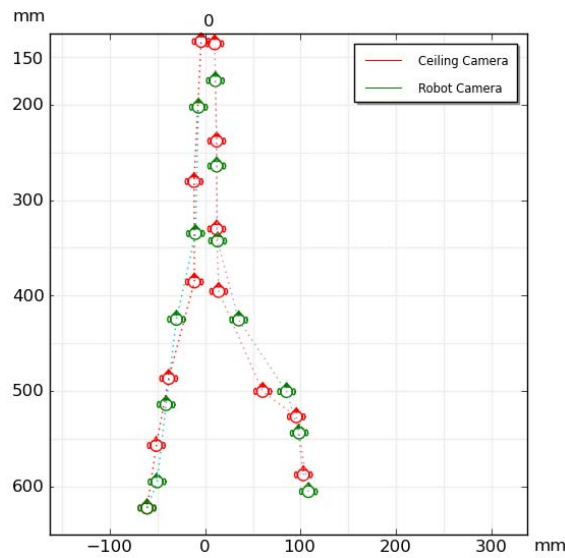
(a) Training Curve of the Predictive Sensory Module. The error is calculated by the mean squared difference between the target and the output of each unit in the output cube.



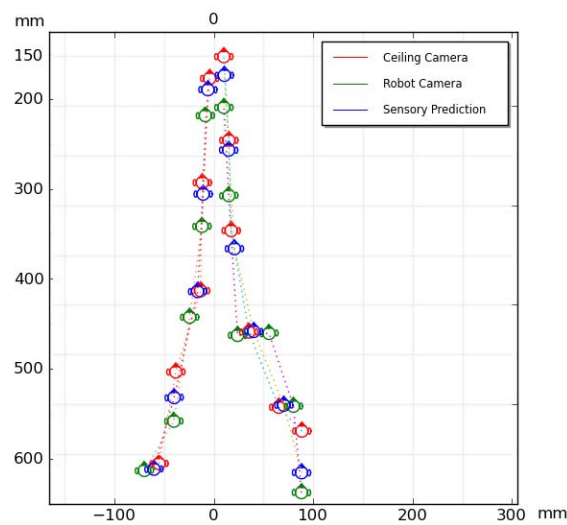
(b) Training curve of the motor action module.

FIGURE 6.4: Training Curves of Two Modules

From the on-site observations of docking trials, the robot sometimes reacted late when observing the landmark. For instance, it usually stopped too close to the docking station, which can lead to failure of grasping. In some cases the robot



(a) Robot Trajectories without Sensory Prediction



(b) Robot Trajectories with Sensory Prediction

FIGURE 6.5: Comparison of Trajectories

Trajectories of docking trials observed by a ceiling camera (true positions, red) and inferred from the robot's current sensory perception (green). (a) Without sensory prediction, sensorimotor delay causes the grasp signal to be produced at positions beyond the optimal position at  $(0\text{cm}, 15\text{cm})$ , which often results in failure of grasping. (b) The sensory prediction (blue) matches the true position better than the current estimate. With prediction, the robot correctly stops and gives the grasping command before the optimal position is inferred from the camera.

performed corrective actions following such responses, which led to longer average docking durations, as evident in Tab. 6.3.

Two trials from both sides are shown in Fig. 6.5. To objectively compare the offsets of the robot camera, the trajectories were also tracked from a ceiling camera. With 0.2s sampling rate, we measured the x-y position of the robot. After synchronizing the ceiling camera and the robot observation data, we observe that a delay happens during the whole docking process, manifested by the observed inferred position being offset from the true position. This latency is likely to be the cause of the observed late NAO reaction, which led to a longer approaching time, when it attempted to produce additional back and forth movements, and even failure of the docking and the following grasping, specifically when NAO elicited the grasping signal when it was already too close to the station (see the final stopping/grasping points were not at (0cm, 150cm) but beyond it in Fig. 6.5).

### **Approaching based on Reinforcement Learning and Predictive Sensory System**

In the following experiment of the integration of both modules, the predictive sensory module should build up an internal mental model to predict the upcoming sensory signal sequence based on the previous sensory experience. We used both the predictive sensory value and the real sensory value by averaging them and applying the averaged value as the input of the units  $a^m$  and  $a^o$ . This method can compensate the sensorimotor cycle latency, filter noisy sensory information and therefore enhance the success-rate and speed while the robot is walking. But the stopping/grasping signal should arrive even earlier due to command delay in mechanics, so only the predictive sensory value is fed as input of the grasping command unit  $a^g$  to further compensate this delay. These trial results are compared to the former results as shown in Tab. 6.3.

To visualize the effect of the predictive sensory model, we also synchronized the predictive sensory percepts of x-y distances and the observed ones in Fig. 6.5(b) with the ceiling camera data. Since the grasping signal only depends on the predictive sensory information, it solves the problem of the grasping signal coming too late and the robot stopping too close to the docking station, as it occurred in the previous experiment. Besides, in Fig. 6.5 we can see that the predicted trajectory is smoother than the one obtained from current robot vision, which hints to a denoising function of the predictive sensory module.



Methods	# of Approaches	# of Success	Success Rate	Avg. Time	Variance
LAN	20	16	0.800	1.25	0.59
WIFI	20	15	0.750	1.35	0.58

TABLE 6.4: Docking Trials in Different Connections

**Docking Trials with Different Connections** To test the model predicting sensory percepts in a general context, an experiment with a different network connection was conducted. Due to different network delays (c.f. Tab. 6.1), the robot received different kinds of sensory percept sequences when we used the slower wireless (WIFI) connection. As mentioned before, the predictive sensory module estimates an HMM process. To train the predictive sensory module for the different connection delays, we adjusted training samples to match the delay time of different connections. Tab. 6.4 shows that the results of docking trials under different connections are similar, implying that the predictive sensory architecture adapts its prediction. Hence motor responses will be different in a different sensory percept context.

## 6.4 Summary

In this chapter, we presented a predictive sensory architecture that predicts the visually retrieved coordinates. It specifically extends the dorsal-like stream in the visual cortex based on recurrent connections, which was shown in the last chapter and models its integration with a motor action module. Together with a continuous reinforcement-learned action strategy based on these predicted sensory values, the predictive architecture resulted in a smoother and faster approaching behaviour in the case study of robot approaching. The filtering function of the predictive sensory module provided a smooth sensory signal leading to a smooth and robust behaviour. We also showed that the predictive sensory model effectively compensated the latency of the sensorimotor cycle of the robot, which led to less errors being made and to faster executed behaviour.

## Chapter 7

# Pre-symbolic Representation Emerged from Sensorimotor Feedback

In last two chapters, we investigated the prediction function of feedback signals on perception, and how it affects the sensorimotor integration. In this chapter, we continue to examine such influences coming from various cognitive processes.

Cognitive processes are learnt by sensory-driven signals. For instance, the acquisition of the symbolic and linguistic representations of sensorimotor behaviour is done by an agent when it is executing and/or observing own and others' actions. Conversely, this cognitive process accomplishes the sensory prediction function via the feedback pathways. This chapter is mostly based on our published work [Zhong et al., 2011, 2014].

### 7.1 Language Acquisition

#### 7.1.1 Pre-symbolic Communication

Although infants are not supposed to acquire a symbolic representational system at the sensorimotor stage, based on Piaget's definition of infant development, the preparation of language development, such as a pre-symbolic representation for conceptualization, has been set at the time when the infant starts babbling [Mandler, 1999]. Experiments have shown that infants have established the concept of animate and inanimate objects, even if they have not yet seen the objects before

[Gelman & Spelke, 1981]. Similar phenomena also include the conceptualization of object affordances such as the conceptualization of containment [Bonniec, 1985]. This conceptualization mechanism is developed at the sensorimotor stage to represent sensorimotor primitives and other object-affordance related properties.

During an infant's development at the sensorimotor stage, one way to learn affordances is to interact with objects using tactile perception, to observe the object from visual perception and thus learn the causal relation between the visual features, affordances and movements as well as to conceptualize them. This learning starts with the basic ability to move an arm towards the visual-fixated objects in new-born infants [von Hofsten, 1982], continues through object-directed reaching at the age of 4 months [Streri et al., 1993, D. Corbetta & Snapp-Childs, 2009], and can also be found during the object exploration of older infants (c.f. [Mandler, 1992, Ruff, 1984]). From these interactions leading to visual and tactile percepts, infants gain experience through the instantiated 'bottom-up' knowledge about object affordances and sensorimotor primitives. Building on this, infants at the age of around 8-12 months gradually expand the concept of object features, affordances and the possible causal movements in the sensorimotor context [Gibson, 1988, C. Newman et al., 2001, Rocha et al., 2006]. For instance, they realize that it is possible to pull a string that is tied to a toy car to fetch it instead of crawling towards it. An associative rule has also been built that connects conceptualized visual feature inputs, object affordance and the corresponding frequent auditory inputs of words, across various contexts [Romberg & Saffran, 2010]. At this stage, categories of object features are particularly learnt in different contexts due to their affordance-invariance [Bloom et al., 1993].

Therefore, the integrated learning process of the object's features, movements according to the affordances, and other knowledge is a globally conceptualized process through visual and tactile perception. This conceptualized learning is a precursor of a pre-symbolic representation of language development. This learning is the process to form an abstract and simplified representation for information exchange and sharing<sup>1</sup>. To conceptualize from visual perception, it usually includes a planning process: first the speaker receives and segments visual knowledge in the perceptual flow into a number of states on the basis of different criteria, then the speaker selects essential elements, such as the units to be verbalized, and last the speaker constructs certain temporal perspectives when the events have to be anchored and linked (c.f. [Habel & Tappe, 1999, von Stutterheim & Nuse,

---

<sup>1</sup>For comparison of conceptualization between engineering and language perspectives, see [Gruber & Olsen, 1994, Bowerman & Levinson, 2001].

2003]). Assuming this planning process is distributed between ventral and dorsal streams, the conceptualization process should also emerge from the visual information that is perceived in each stream, associating the distributed information in both streams. As a result, the candidate concepts of visual information are statistically associated with the input stimuli. For instance, they may represent a particular visual feature with a particular class of label (e.g. a particular visual stimuli with an auditory wording ‘circle’) [Chemla et al., 2009]. Furthermore, the establishment of such links also strengthens the high-order associations that generate sensory predictions and generalize to novel visual stimuli [Yu, 2008]. Once the infants have learnt a sufficient number of words, they begin to detect a particular conceptualized cue with a specific kind of wording. At this stage, infants begin to use their own conceptualized visual ‘database’ of known words to identify a novel meaning class and possibly to extend their wording vocabulary [Smith et al., 2002]. Thus, this associative learning process enables the acquisition and the extension of the concepts of domain-specific information (e.g. features and movements in our experiments) with the visual stimuli.

This conceptualization will further result in a pre-symbolic way for infants to communicate when they encounter a conceptualized object and intend to execute a correspondingly conceptualized well-practised sensorimotor action towards that object. For example, behavioural studies showed that when 8-to-11-month-old infants are unable to reach and pick up an empty cup, they may point it out to the parents and execute an arm movement intending to bring it to their lips. The conceptualized shape of a cup reminds infants of its affordance and thus they can communicate in a pre-symbolic way. Thus, the emergence from the conceptualized visual stimuli to the pre-symbolic communication also gives further rise to the different periods of learning nouns and verbs in infancy development (c.f. [Gentner, 1982, Tardif, 1996, Bassano, 2000]). This evidence supports that the production of verbs and nouns are not correlated to the same modality in sensory perception: experiments performed by Kersten [1998] suggest that nouns are more related to the movement orientation caused by the intrinsic properties of an object, while verbs are more related to the trajectories of an object. Thus, we argue that such differences of acquisitions in lexical classes also relate to the conceptualized visual ventral and dorsal streams. The finding is consistent with Damasio & Tranel [1993]’s hypothesis that verb generation is modulated by the perception of conceptualization of movement and its spatio-temporal relationship.

For this reason, we propose that the conceptualized visual information, which is

a prerequisite for the pre-symbolic communication, is also modulated by perception in two visual streams. As we introduced in the last chapter, there have been studies of modelling the functional modularity in the development of ventral and dorsal streams (e.g. [Jacobs et al., 1991, Mareschal et al., 1999]) and the bilinear models of visual routing (e.g. [Olshausen et al., 1993, Bergmann & von der Malsburg, 2011, Memisevic & Hinton, 2007]). However, a model which explains the development of conceptualization from both streams and results in an explicit representation of conceptualization of both streams, while the visual stimuli are presented, is still missing in the literature. This conceptualization should be able to encode the same category for information flows in both ventral and dorsal streams like ‘object files’ in the visual understanding [Fields, 2011] so that they could be discriminated in different contexts during language development.

On the other hand, this conceptualized representation that is distributed in two visual streams is also able to predict the tendency of appearance of an action-oriented object in the visual field via feedback pathways, which causes some sensorimotor phenomena such as object permanence [Tomasello & Farrar, 1986] showing that the infants’ attention is usually driven by the object’s features and movements. For instance, when infants are observing the movement of the object, recording showed an increase of the looking times when the visual information after occlusion is violated in either surface features or location [Mareschal & Johnson, 2003]. Also, the words and sounds play a top-down role in the early infants’ visual attention [Sloutsky & Robinson, 2008]. This could hint at the different development stages of the ventral and dorsal streams and their effect on the conceptualized prediction mechanism in the infant’s consciousness.

Accordingly, the model we propose about the conceptualized visual information should also be able to explain the emergence of a predictive function in the sensorimotor system, e.g. the ventral stream attempts to track the object and the dorsal stream processes and predicts the object’s spatial location, when the sensorimotor system is involved in an object interaction. We have been aware of that this build-in predictive function in a forward sensorimotor system is essential: neuroimaging research has revealed the existence of internal forward models in the parietal lobe and the cerebellum that predict sensory perception from efference copies of motor commands [Kawato et al., 2003] and supports fast motor reactions (e.g. [Hollerbach, 1982]). Since the probable position and the movement pattern of the action should be predicted on a short time scale, sensory feedback produced by a forward model with negligible delay is necessary in this sensorimotor loop.

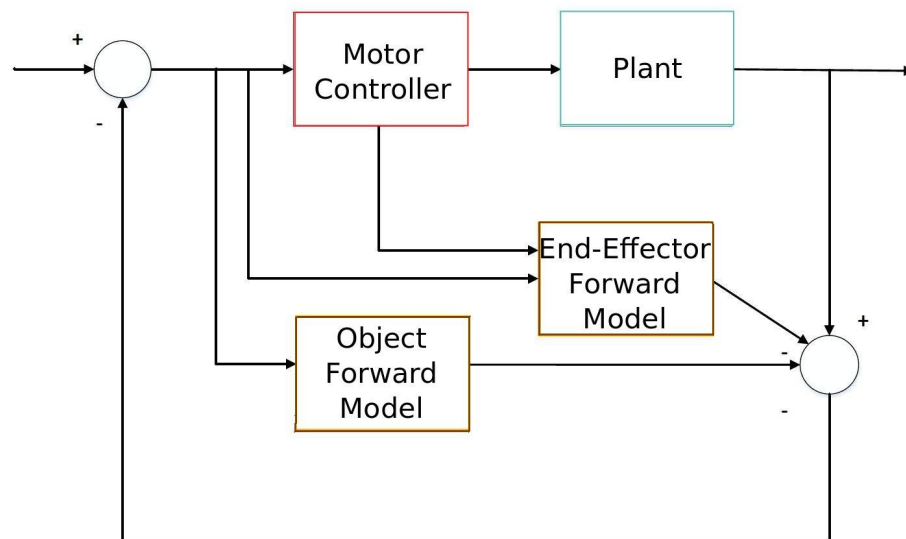


FIGURE 7.1: Diagram of Sensorimotor Integration with the Object Interaction. The lower forward model predicts the object movement, while the upper forward model extracts the end-effector movement from sensory information in order to accomplish a certain task (e.g. object interaction).

### 7.1.2 Motivation

The sensorimotor model to explain the predictive role of the feedback pathways is suitable to work as one of the building modules that takes into account the predictive object movement in a forward sensorimotor system to deal with object interaction from visual stimuli input as Fig. 7.1 shows. This system is similar to Wolpert et al. [1995]’s sensorimotor integration, but it includes an additional sensory estimator (the lower brown block) which takes into account the visual stimuli from the object so that it is able to predict the dynamics of both the end-effector (which is accomplished by the upper brown block) and the sensory input of the object. This object-predictive module is essential in a sensorimotor system to generate sensorimotor actions like tracking and avoiding when dealing with fast-moving objects, e.g. in ball sports. We also assert that the additional inclusion of forward models in the visual perception of the objects can explain some predictive developmental sensorimotor phenomena, such as object permanence.

In summary, we propose a model which should accomplish the following tasks:

- Links should be established between the development of ventral/dorsal visual streams and the emergence of the conceptualization in visual streams, which further leads to the feedback predictive function of a sensorimotor system;

- As such conceptualization may become a source of the feedback signals in a forward (predictive) sensorimotor model, we verify the hypothesis that a higher-level cognition process (e.g. symbolic representation) and a predictive sensorimotor process are integrated.

To validate this proof-of-concept model, we also conducted experiments in a simplified robotics scenario. Two NAO robots were employed in the experiments: one of them was used as a ‘presenter’ and moved its arm along pre-programmed trajectories as motion primitives. A ball was attached at the end of the arm so that another robot could obtain the movement by tracking the ball. Our neural network was trained and run on the other NAO, which was called the ‘observer’. In this way, the observer robot perceived the object movement from its vision passively, so that its network took the object’s visual features and the movements into account. Though we could also use one robot and a human presenter to run the same tasks, we used two identical robots, due to the following reasons: 1. the object movement trajectories can be done by a pre-programmed machinery so that the types and parameters of it can be adjusted; 2. the use of two identical robots allows to interchange the roles of the presenter and observer in an easier manner.

## 7.2 Horizontal Recurrent Network with Parametric Bias

A similar forward model exhibiting sensory prediction for visual object perception has been proposed in the previous chapter (Chap. 5) where we suggested an RNN implementation of the sensory forward model. Together with a CACLA trained multi-layer network as a controller model, the forward model embodied in a robot receiving visual landmark percepts enabled a smooth and robust robot behaviour. However, one drawback of this work was its inability to store multiple sets of spatial-temporal input-output mappings, i.e. the learning did not converge if there appeared several spatial-temporal mapping sequences in the training. Consequently, a simple RNN network was not able to predict different sensory percepts for different reward-driven tasks. Another problem was that it assumed that only one visual feature would appear in the robot’s visual field, and this one was the only visual cue it could learn during development.

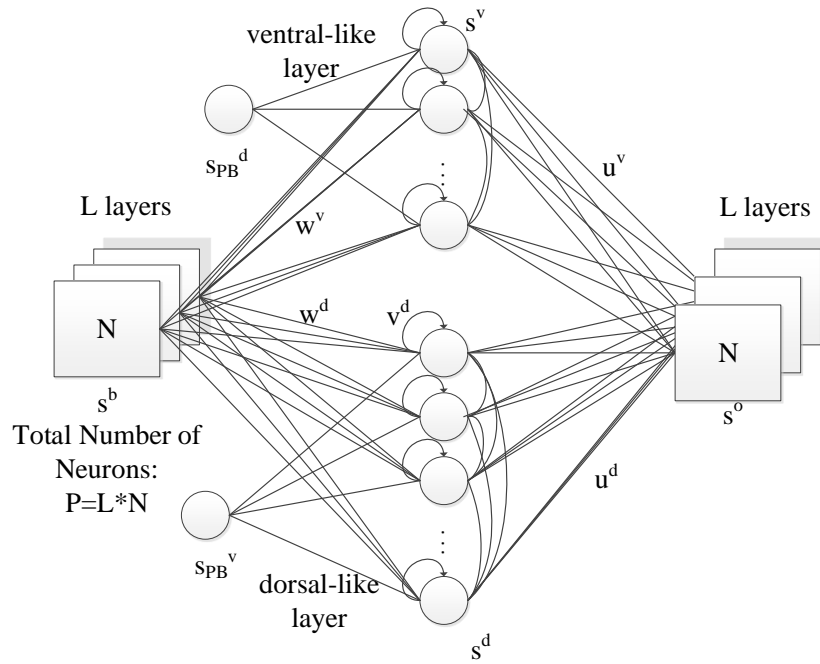


FIGURE 7.2: The HoRNNPB Network Architecture

In this network,  $L$  layers represent  $L$  different types of features. Size of  $N$  indicates the positional information of the object.

To solve the first problem, we will consider to use the RNNPB model we introduced in Chap. 3. Merging the ideas of RNNPB and the forward model, in the context of sensorimotor integration of object interaction, the PB units can be considered as a small set of high-level conceptualized units that describe various types of non-linear dynamics of visual percepts, such as features and movements. This representation is more related to the ‘natural prototypes’ from visual perception, for instance, than a specific language representation [Rosch, 1973]. The development of PB units can also be seen as the pre-symbolic communication that emerges during sensorimotor learning. The conceptualization, on the other hand, could also result in the prediction of future visual percepts of moving objects in sensorimotor integration via feedback pathways.

In this novel horizontal recurrent model with parametric bias (HoRNNPB) (Fig. 7.2), we propose a three-layer, horizontal-product Elman network with PB units. Similar to the original RNNPB model, the network is capable of being executed under three running modes, according to the pre-known conditions of inputs and outputs: learning, recognition and prediction. In learning mode, the representation of object features and movements are first encoded in the weights of both streams, while the bifurcation parameters with a smaller number of dimensions are encoded



in the PB units. This is consistent with the emergence of the conceptualization at the sensorimotor stage of infant development.

Apart from the PB units, another novelty in the network is that the visual object information is encoded in two neural streams and is further conceptualized in PB units. Two streams share the same set of input neurons, where the coordinates of the object in the visual field are used as identities of the perceived images. The appearance of values in different layers represents different visual features: in our experiment, the colour of the object detected by the yellow filter appears in the first layer whereas the colour detected by the green filter appears in the second layer; the other layer remains zero. For instance, the input  $((0, 0), (x, y))$  represents a green object at  $(x, y)$  coordinates in the visual field. The hidden layer contains two independent sets of units representing dorsal-like ‘ $d$ ’ and ventral-like ‘ $v$ ’ neurons respectively. Similar to the model we showed in Chap. 4, these two sets of neurons are inspired by the functional properties of dorsal and ventral streams:

- fast responding dorsal-like units predict object position and hence encode movements;
- slow responding ventral-like units represent object features.

The recurrent connection in the hidden layers also helps to predict movements in layer  $d$  and to maintain a persistent representation of an object’s feature in layer  $v$ . The horizontal product brings both pathways together again in the output layer with one-step ahead predictions. Let us denote the output layer’s input from layer  $d$  and layer  $v$  as  $x^d$  and  $x^v$ , respectively. The network output  $s^o$  is obtained via the horizontal product as

$$s^o = x^d \odot x^v \tag{7.1}$$

where  $\odot$  indicates element-wise multiplication, so each pixel is defined by the product of two independent parts, i.e. for output unit  $k$  it is  $s_k^o = x_k^d \cdot x_k^v$ .

### 7.2.1 Algorithm

Similar as in the previous two chapters, we use  $s^b(t)$  to represent the activation and  $PB^{d/v}(t)$  to represent the activation of the dorsal/ventral PB units at time-step  $t$ . In some of the following equations, the time-index  $t$  is omitted if all activations

are from the same time-step. The inputs to the hidden units  $y_j^v$  in the ventral stream and  $y_j^d$  in the dorsal stream are defined as

$$y_l^d(t) = \sum_i s_i^b(t)w_{li}^d + \sum_{l'} s_l^d(t-1)v_{ll'}^d + \sum_{n_2} PB_{n_2}^v(t)\bar{w}_{ln_2}^d \quad (7.2)$$

$$y_j^v(t) = \sum_i s_i^b(t)w_{ji}^v + \sum_{j'} s_j^v(t-1)v_{jj'}^v + \sum_{n_1} PB_{n_1}^d(t)\bar{w}_{jn_1}^v \quad (7.3)$$

where  $w_{li}^d, w_{ji}^v$  represent the weighting matrices between dorsal/ventral layers and the input layer,  $\bar{w}_{li}^d, \bar{w}_{ji}^v$  represent the weighting matrices between PB units and the two hidden layers, and  $v_{ll'}^d$  and  $v_{jj'}^v$  indicate the recurrent weighting matrices within the hidden layers.

The transfer functions in both hidden layers and the PB units all employ the sigmoid function recommended by LeCun et al. [1998],

$$s_{l/j}^{d/v} = 1.7159 \cdot \tanh\left(\frac{2}{3}y_{l/j}^{d/v}\right) \quad (7.4)$$

$$PB_{n_1/n_2}^{d/v} = 1.7159 \cdot \tanh\left(\frac{2}{3}\rho_{n_1/n_2}^{d/v}\right) \quad (7.5)$$

where  $\rho^{d/v}$  represent the internal values of the PB units.

The terms of the horizontal products of both pathways can be presented as follows:

$$x_k^v = \sum_j s_j^v u_{kj}^v; \quad x_k^d = \sum_l s_l^d u_{kl}^d \quad (7.6)$$

The output of the two streams composes a horizontal product for the network output as we defined in Eq. 7.1.

**Learning Mode** The training progress is basically determined by the cost function:

$$C = \frac{1}{2} \sum_t^T \sum_k^N (s_k^b(t+1) - s_k^o(t))^2 \quad (7.7)$$

where  $s_i^b(t+1)$  is the one-step ahead input (as well as the desired output),  $s_k^o(t)$  is the current output,  $T$  is the total number of available time-step samples in a complete sensorimotor sequence and  $N$  is the number of output nodes which is equal to the number of input nodes. Following gradient descent, each weight update in the network is proportional to the negative gradient of the cost with

respect to the specific weight  $w$  that will be updated:

$$\Delta w_{ij} = -\eta_{ij} \frac{\partial C}{\partial w_{ij}} \quad (7.8)$$

where  $\eta_{ij}$  is the adaptive learning rate of the weights between neuron  $i$  and  $j$ , which is adjusted in every epoch [Kleesiek et al., 2013]. To determine whether the learning rate has to be increased or decreased, we compute the changes of the weight  $w_{i,j}$  in consecutive epochs:

$$\sigma_{i,j} = \frac{\partial C}{\partial w_{i,j}}(e-1) - \frac{\partial C}{\partial w_{i,j}}(e) \quad (7.9)$$

The update of the learning rate is

$$\eta_{i,j}(e) = \begin{cases} \min(\eta_{i,j}(e-1) \cdot \xi^+, \eta_{max}) & \text{if } \sigma_{i,j} > 0, \\ \max(\eta_{i,j}(e-1) \cdot \xi^-, \eta_{min}) & \text{if } \sigma_{i,j} < 0, \\ \eta_{i,j}(e-1) & \text{else.} \end{cases} \quad (7.10)$$

where  $\xi^+ > 1$  and  $\xi^- < 1$  represent the increasing/decreasing rate of the adaptive learning rates, with  $\eta_{min}$  and  $\eta_{max}$  as lower and upper bounds, respectively. Thus, the learning rate of a particular weight increases by  $\xi^+$  to speed up the learning when the changes of that weight from two consecutive epochs have the same sign, and vice versa.

Besides the usual weight update according to back-propagation through time, the accumulated error over the whole time-series also contributes to the update of the PB units. The update for the  $i$ -th unit in the PB vector for a time-series of length  $T$  is defined as:

$$\rho_i(e+1) = \rho_i(e) + \gamma_i \sum_{t=1}^T \delta_{i,j}^{PB} \quad (7.11)$$

where  $\delta^{PB}$  is the error back-propagated to the PB units,  $e$  is  $e$ -th time-step in the whole time-series (e.g. epoch),  $\gamma_i$  is PB units' adaptive updating rate which is proportional to the absolute mean value of the back-propagation error at the  $i$ -th PB node over the complete time-series of length  $T$ :

$$\gamma_i \propto \frac{1}{T} \sum_{t=1}^T \delta_{i,j}^{PB} \quad (7.12)$$

The reason for applying the adaptive technique is that it was realized that the PB units converge with difficulty. Usually a smaller learning rate is used in the generic version of RNNPB to ensure the convergence of the network. However, this results in a trade-off in convergence speed. The adaptive learning rate is an efficient technique to overcome this trade-off [Kleesiek et al., 2013].

**Recognition Mode** The recognition mode is executed with a similar information flow as the learning mode: given a set of spatio-temporal sequences, the error between the target and the real output is back-propagated through the network to the PB units. However, the synaptic weights remain constant and only the PB units will be updated, so that the PB units are self-organized as the pre-trained values after certain epochs. Assuming the length of the observed sequence is  $a$ , the update rule is defined as:

$$\rho_i(e+1) = \rho_i(e) + \gamma \sum_{t=T-a}^T \delta_{i,j}^{PB} \quad (7.13)$$

where  $\delta^{PB}$  is the error back-propagated from a certain sensory information sequence to the PB units and  $\gamma$  is the updating rate of PB units in recognition mode, which should be larger than the adaptive rate  $\gamma_i$  at the learning mode.

**Generation Mode** The values of the PB units can also be manually set or obtained from recognition, so that the network can generate the upcoming sequence with one-step prediction.

## 7.3 Experiments

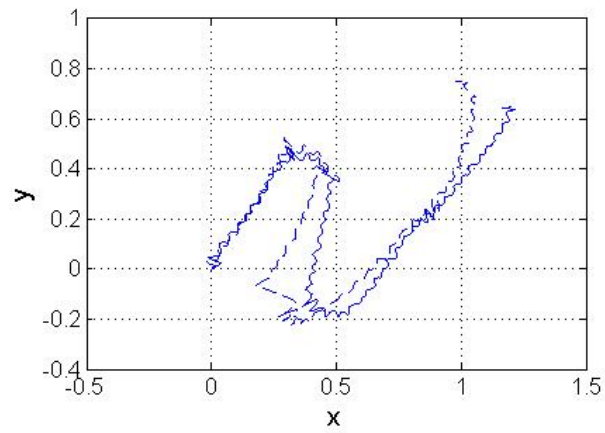
### 7.3.1 Preliminary Experiments

To introduce the usage of RNNPB and test if it is useful to recognize and to predict multiple perception sequences, in this part of the experiment, we will first conduct preliminary experiments on action recognition and prediction using a generic RNNPB model. Due to PB units' property of recognizing temporal input sequences, PB units should be beneficial to robot action understanding. With a prototype architecture of PB units, we will focus on robot trajectory prediction and recognition in the following sections, and test their ability to understand *what* the robot is doing.

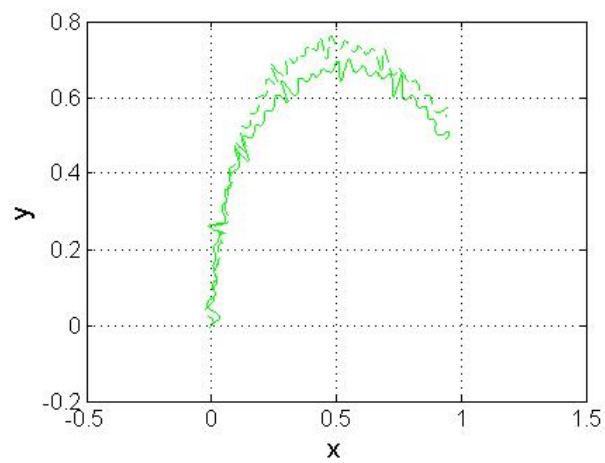
As a foundation of robot action understanding, the recognition and prediction of robot walking trajectories are the objectives of the following experiments. For effectiveness, we use the Webots simulator [Michel, 2004] to collect the trajectory data. Our NAO robot is controlled in the Webots simulator to walk along pre-defined trajectories. From the supervisor function within Webots, three kinds of trajectories, that is a straight line, a sine curve and a half circle were recorded using x and y coordinates. Different combinations of these curves with different parameters make the robot walk in various trajectories. As an initial controlled experiment with known trajectories, we select three trajectories: 1. a sine curve:  $y = 50\sin(\frac{2\pi}{3} * x)cm$ ; 2. a half circle curve with 50cm radius; 3. a straight line.

The reason why we use these kinds of trajectories is that they can constitute different kinds of trajectories, e.g. a walking trajectory when doing obstacle avoidance, by changing their parameters, i.e. frequency and amplitude in sine curve and radius in half circle. We train the network with three types of input sequences. The expectation is that the generalization ability of PB units can recognize similar trajectories with different parameters. For all simulations we use the same network: 2 input nodes, 10 hidden nodes, 10 context nodes, 2 output nodes. Additionally, we use 3 PB nodes in the experiment. The empirically determined network parameters are:  $l = 30$ ,  $\eta_l = 0.01$ ,  $\eta_r = 0.5$ , and the learning rate of connection weights in back-propagation is defined by  $\eta_{BP} = 0.01$ .

After training, we input the walking records from the above three trajectories respectively and attempt to predict the position one step ahead given the previous inputs. As shown in Figure 7.3, after several time-steps, the network can basically predict the learnt trajectories. Furthermore, we inspect the values in the PB units while we continuously feed three types of input sequences into it. As shown in Figure 7.4 internal values in the PB units can reflect different input patterns of the whole learning sequence.

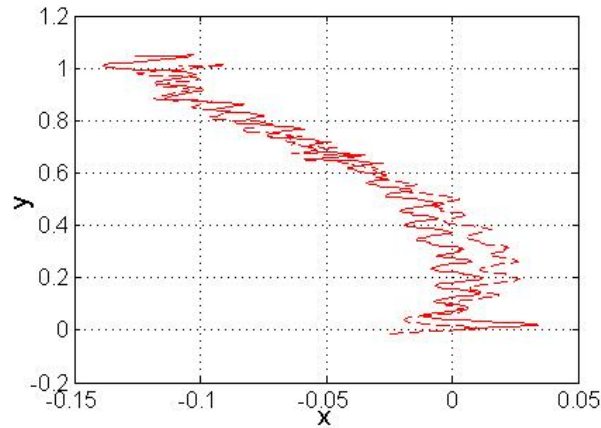


(a) Prediction of sin curve



(b) Prediction of circle curve

FIGURE 7.3: Prediction of Three Curves



(a) Prediction of straight line

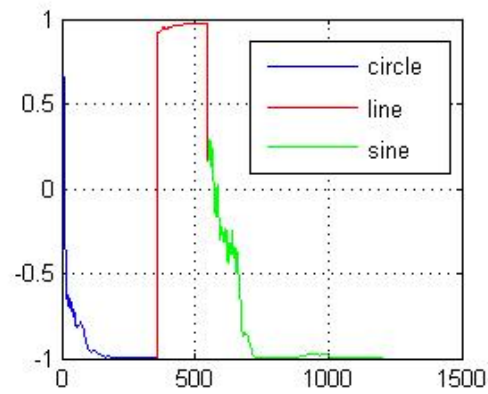
FIGURE 7.3: Prediction of Three Curves (cont.)

Prediction experiments were done for three types of trained sequences. Solid lines represent the true positions and the dashed line represent the predictions. It can be observed that the predicted sequence and the target sequence were quite close in the above figures.

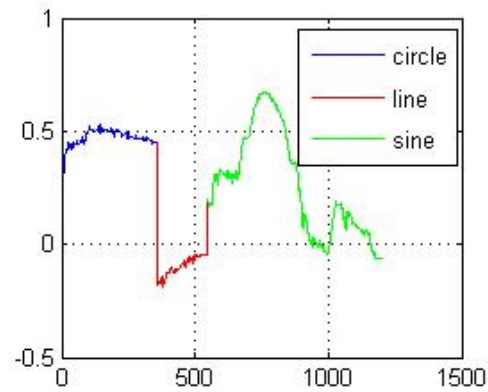
RMSE	sine	line	circle	sine2	line2	circle2
x coordinate	0.0714	0.0052	0.0077	0.2655	0.0187	0.0427
y coordinate	0.0829	0.0066	0.0108	0.1884	0.0094	0.0744

TABLE 7.1: Root Mean Square Error of Two Curve Set Predictions

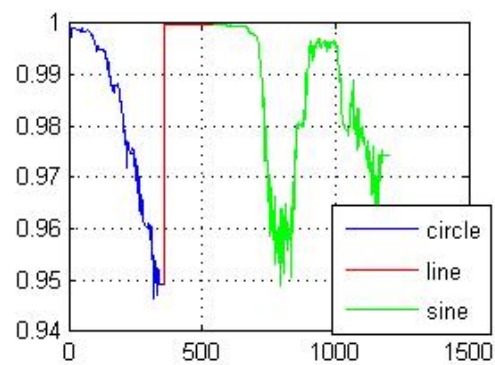
Secondly, we attempt to input another three different types of similar patterns, but with different parameters in order to test the generalization ability for other untrained trajectories. Fig. 7.5 shows the prediction results of the network. Although some errors occur, the generalization of the network still successfully predicts the trend of the curves: 1. a sine curve:  $y = 100\sin(\frac{\pi}{2} * x)cm$ ; 2. a half circle curve with  $30cm$  radius; 3. a straight line. Tab. 7.1 shows the RMS error between the true value and prediction.



(b) Output in PB Unit 1



(c) Output of PB Unit 2

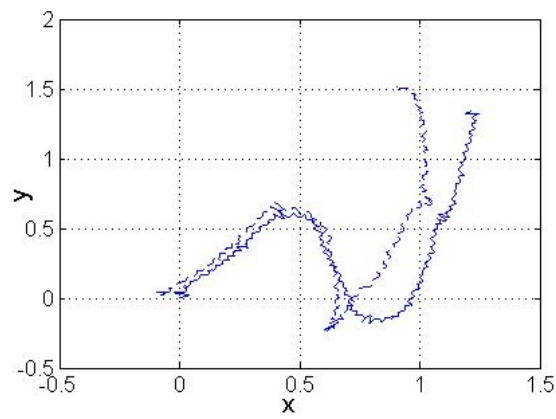


(d) Output of PB Unit 3

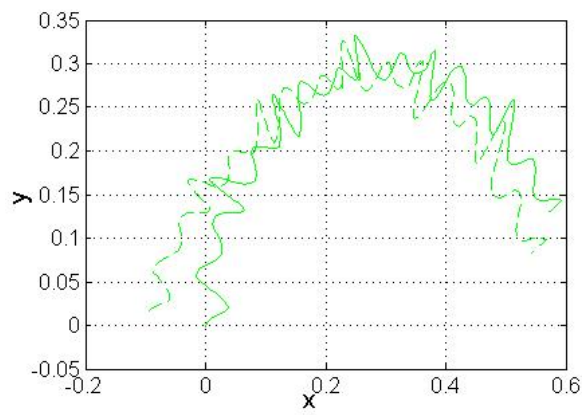
FIGURE 7.4: PB Values in Recognition

Three sequences were fed into the network to demonstrate the recognition in PB units.



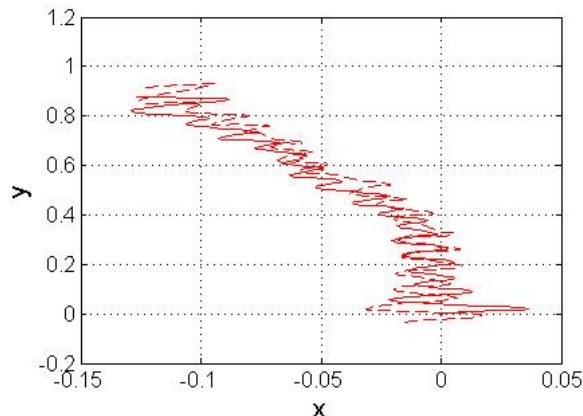


(a) Prediction of sin curve 2



(b) Prediction of circle curve 2

FIGURE 7.5: Prediction of Three Untrained Curves



(a) Prediction of straight line 2

FIGURE 7.5: Prediction of Three Untrained Curves (cont.)

The errors between the predicted curves and untrained curves were larger than those in Fig. 7.3, but the trend of the similar curves can also be predicted.

### 7.3.2 Pre-symbolic Learning via Interaction

In this experiment, we examined our HoRNNPB network with robotic experiments. Two NAO robots were placed face-to-face in a rectangle box of  $61.5\text{cm} \times 19.2\text{cm}$  as shown in Fig. 7.6. These distances were carefully adjusted so that the observer was able to keep track of movement trajectories in its visual field during all experiments using the images from the lower camera. The NAO robot has two cameras. We use the lower one to capture the images because its installation angle is more suitable to track the balls when they are held in the other NAO's hand.

Two  $3.8\text{cm}$ -diameter balls with yellow/green colour were used for the following experiments. The presenter consecutively held each of the balls to present the object interaction. The original image, received from the lower camera of the observer, was pre-processed with threshold in HSV colour-space and the coordinates of its centroid in the image moment were calculated. Here we only considered two different colours as the only feature to be encoded in the ventral stream, as well as two sets of movement trajectories encoded in the dorsal stream. Although we have only tested a few categories of trajectories and features, we believe the results can be predicted to multiple categories in future applications.

**Learning** The two different trajectories (in  $\text{cm}$ ) are defined as below,

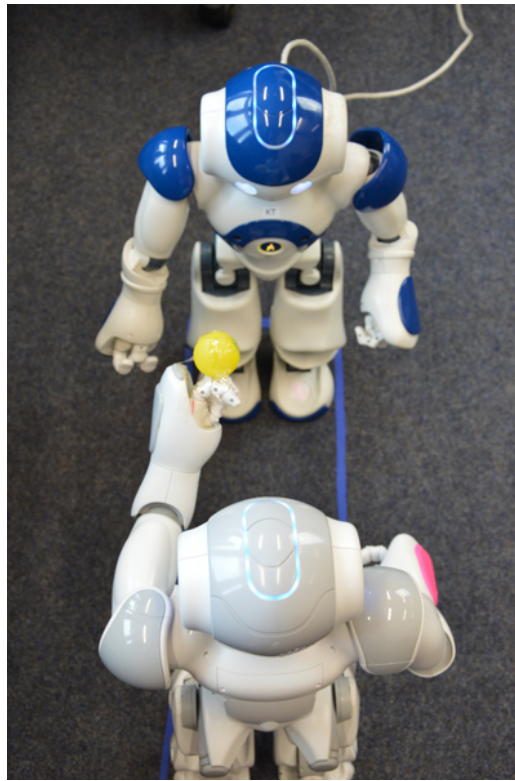


FIGURE 7.6: Experimental Scenario  
Two NAOs are standing face-to-face with in a rectangle box.

the *cosine* curve,

$$x = 12 \tag{7.14}$$

$$y = 8 \cdot \left(-\frac{t}{2}\right) + 0.04 \tag{7.15}$$

$$z = 4 \cdot \cos(2t) + 0.10 \tag{7.16}$$

and the *square* curve,

$$x = 12 \quad (7.17)$$

$$y = \begin{cases} 0 & t \leq -\frac{3\pi}{4} \\ \frac{16}{\pi}t + 12 & -\frac{3\pi}{4} < t \leq -\frac{\pi}{4} \\ 8 & -\frac{\pi}{4} < t \leq \frac{\pi}{4} \\ -\frac{16}{\pi}t + 12 & \frac{\pi}{4} < t \leq \frac{3\pi}{4} \\ 0 & t > \frac{3\pi}{4} \end{cases} \quad (7.18)$$

$$z = \begin{cases} \frac{16}{\pi}t + 20 & t \leq -\frac{3\pi}{4} \\ 14 & -\frac{3\pi}{4} < t \leq -\frac{\pi}{4} \\ -\frac{16}{\pi}t + 10 & -\frac{\pi}{4} < t \leq \frac{\pi}{4} \\ 6 & \frac{\pi}{4} < t \leq \frac{3\pi}{4} \\ \frac{16}{\pi}t - 6 & t > \frac{3\pi}{4} \end{cases} \quad (7.19)$$

where the 3-dimension tuple  $(x, y, z)$  are the coordinates (centimetres) of the ball w.r.t the torso frame of the NAO presenter.  $t$  loops between  $(-\pi, \pi]$ . In each loop, we calculated 20 data points to construct trajectories with 4s sleeping time between every two data points. Note that although we have defined the optimal desired trajectories, the arm movement was not ideally identical to the optimal trajectories due to the noisy position control of the end-effector of the robot. On the observer side, the  $(x, y)$  coordinates of the colour-filtered moment of the ball in the visual field were recorded to form a trajectory with sampling time of 0.2s. Five trajectories, in the form of tuple  $(x, y, z)$  w.r.t the torso frame of the NAO observer were recorded with each colour and each curve, so in total 20 trajectories were available for training.

Parameters	Parameter's Descriptions	Value
$\eta_{ventral}$	Learning Rate in Ventral Stream	$1.0 \times 10^{-5}$
$\eta_{dorsal}$	Learning Rate in Dorsal Stream	$1.0 \times 10^{-3}$
$\eta_{max}$	Maximum Value of Learning Rate	$1.0 \times 10^{-1}$
$\eta_{min}$	Minimum Value of Learning Rate	$1.0 \times 10^{-7}$
$M_\gamma$	Proportionality Constant of PB Units Updating Rate	$1.0 \times 10^{-2}$
$n_1$	Size of PB Unit 1	1
$n_2$	Size of PB Unit 2	1
$n_v$	Size of Ventral-like Layer	50
$n_d$	Size of Dorsal-like Layer	50
$\xi^-$	Decreasing Rate of Learning Rate	0.999999
$\xi^+$	Increasing Rate of Learning Rate	1.000001

TABLE 7.2: Network Parameters (HoRNNPB)

In each training epoch, these trajectories, in the form of tuples, were fed into the input layer one after another for training, with the tuples of the next time-step serving as a training target. The empirically-obtained parameters are listed in Tab. 7.2. The final PB values were examined after the training was done, and the values were shown in Fig. 7.7. It can be seen that the first PB unit, along with the dorsal stream, was approximately self-organized with the colour information, while the second PB unit, along with the ventral stream, was self-organized with the movement information.

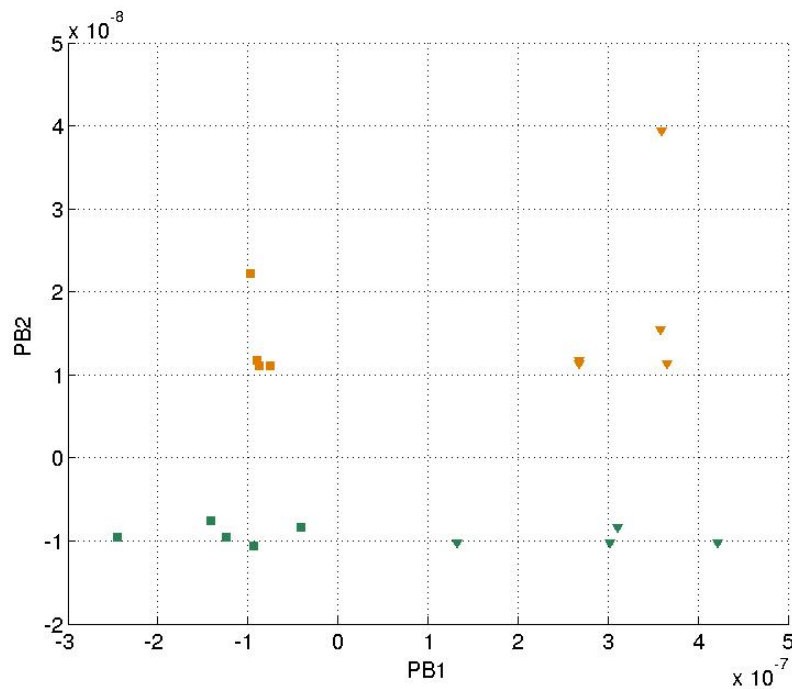
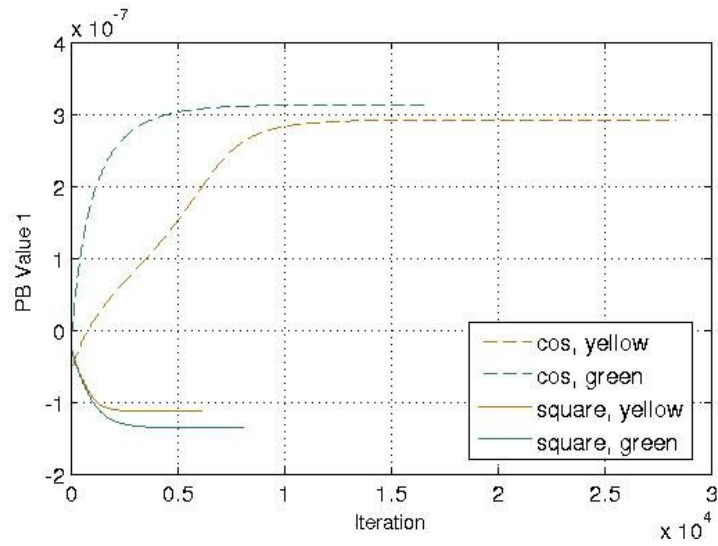


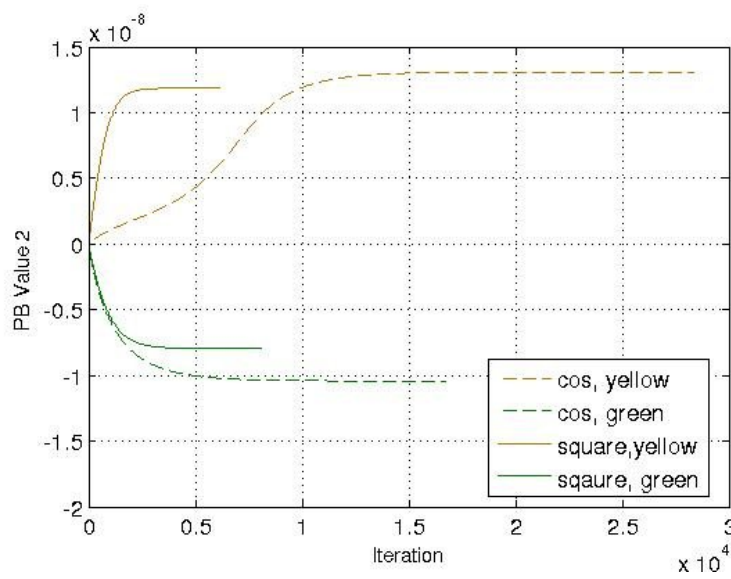
FIGURE 7.7: Values of PB Units in Two Streams

The square markers represent those PB units after the *square* curves training and the triangle markers represent those of the *cosine* curves training. The colours of the markers, yellow and green, represent the colours of the balls used for training.

**Recognition** Another four trajectories were presented in the recognition experiment, in which the length of the sliding-window is equal to the length of the whole time-series, i.e.  $T = a$  in Eq. 7.13. The update of the PB units are shown in Fig. 7.8. Although we used the complete time-series sequence for the recognition, it should also be possible to use only part of the sequence, e.g. through the sliding-window approach with a smaller number of  $a$  to fulfil the real-time requirement in the future.

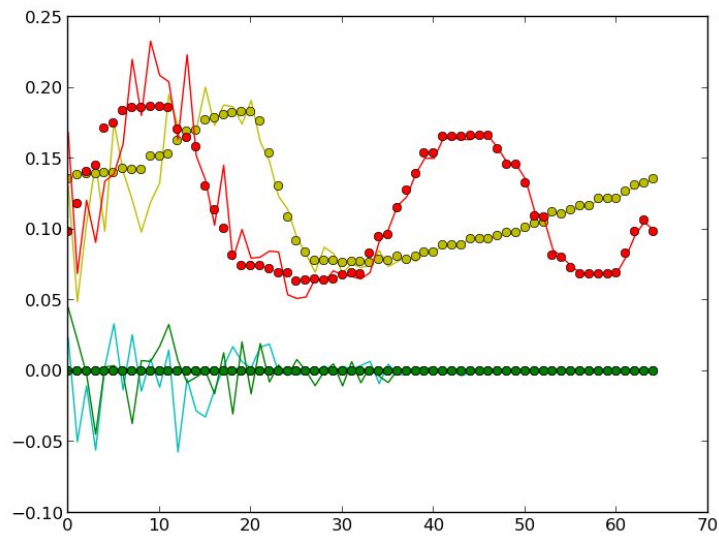


(a) PB value 1

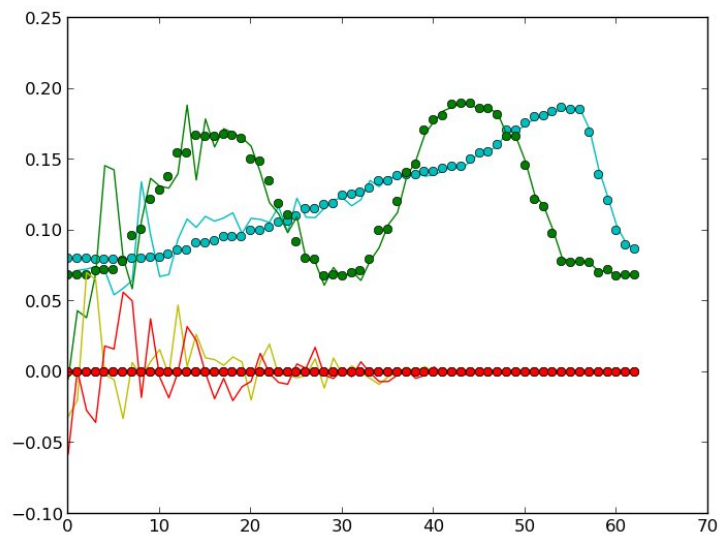


(b) PB value 2

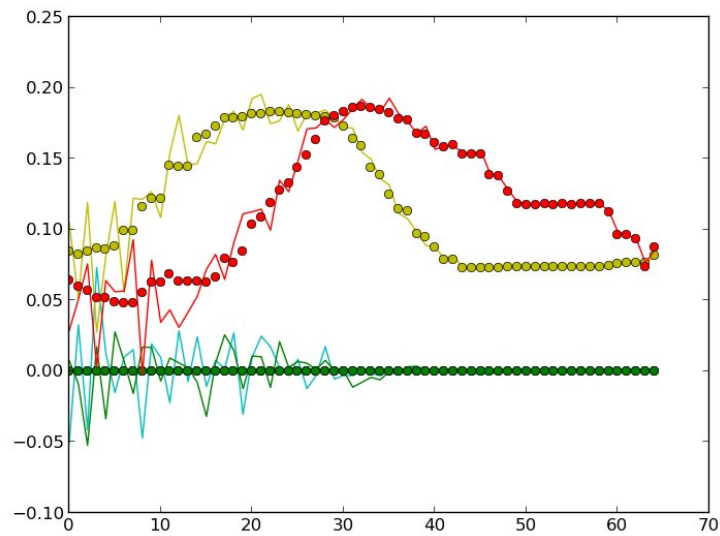
FIGURE 7.8: Update of the PB Values in Recognition Mode



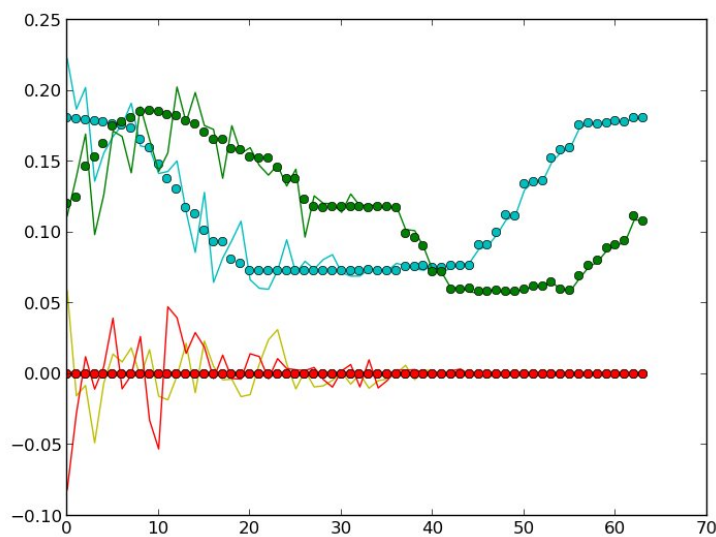
(a) Cosine curve, yellow ball



(b) Cosine curve, green ball



(c) Square curve, yellow ball



(d) Square curve, green ball

FIGURE 7.8: Generated Values from HoRNNPB

The dots denote the true values for comparison, curves show the estimated ones. Yellow and red colours represent the values of the two neurons in the first layer (yellow), the colours green and cyan represent those in the second layer (green).



Error of Outputs	Unit 1	Unit 2	Unit 3	Unit 4
cosine, yellow	$2.28 \times 10^{-4}$	$8.09 \times 10^{-5}$	$7.29 \times 10^{-4}$	$8.63 \times 10^{-4}$
cosine, green	$8.34 \times 10^{-4}$	$7.04 \times 10^{-4}$	$1.50 \times 10^{-4}$	$2.01 \times 10^{-4}$
square, yellow	$3.91 \times 10^{-4}$	$9.64 \times 10^{-5}$	$1.74 \times 10^{-3}$	$3.23 \times 10^{-4}$
square, green	$1.40 \times 10^{-3}$	$3.27 \times 10^{-4}$	$3.54 \times 10^{-4}$	$2.60 \times 10^{-4}$

TABLE 7.3: Prediction Error

**Generation** In this simulation, the obtained PB units from the previous recognition experiment were used to generate the predicted movements using the prior knowledge of a specific object. Then, the one-step prediction from the output units were again applied to the input at the next time-step, so that the whole time-series corresponding to the object’s movements and features were obtained. Fig. 7.8 presents the comparisons between the true values (the same as used in recognition) and the predicted ones.

From Fig. 7.8, it can be observed that the estimation was biased quite largely to the true value within the first few time-steps, as the RNN needs to accumulate enough input values to access its short-term memory. However, the error became smaller and it kept track of the true value in the following time-steps. Considering that the curves are automatically generated given the PB units and the values at the first time-step, the error between the true values and the estimated ones are acceptable. Moreover, this result shows clearly that the conceptualization affects the (predictive) visual perception.

**Generalization in Recognition** To testify whether our new computational model has the generalization ability as Cuijpers et al. [2009] proposed, we recorded another set of sequences of a circle trajectory. The trajectory (in *cm*) is defined as:

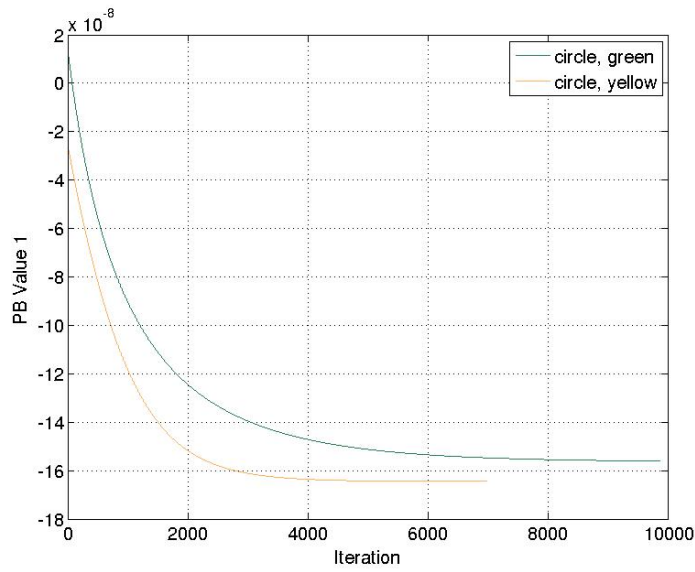
$$x = 12 \tag{7.20}$$

$$y = 4 \cdot \sin(2t) + 0.04 \tag{7.21}$$

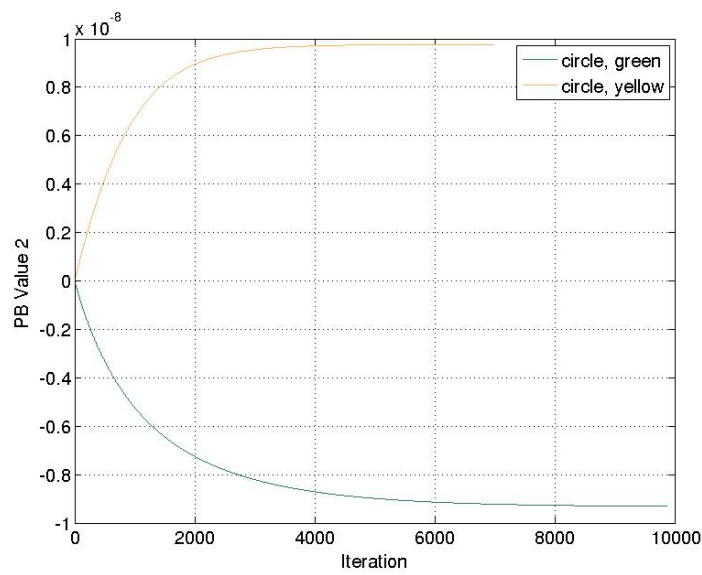
$$z = 4 \cdot \cos(2t) + 0.10 \tag{7.22}$$

The yellow and green balls were still used. We ran the recognition experiment again with the weight previously trained. The update of the PB units are shown in Fig.7.9. Comparing Fig. 7.7 and Fig.7.9, we can observe that the positive and negative signs of PB values are similar to the square trajectory. This is

probably because the visual perception of circle and square movements have more similarities than those between circle and cosine movements.



(e) PB value 1



(f) PB value 2

FIGURE 7.9: Update of the PB Values in Recognition Mode with an Untrained Feature (Circle)

**PB Representation with Different Speeds** We further generated 20 trajectories with the same data functions (Eqs. 7.14 - 7.19) but with a slower sampling time. In other words, the movement of the balls seemed to be faster with the robot’s observation. The final PB values after training were shown in Fig. 7.10.

Comparing with Fig. 7.7, It can be seen that generally the PB values were smaller in Fig.7.9, which was probably because there was less error being propagated during training. Moreover, the corresponding PB values corresponding to colours (green and yellow) and movements (cosine and square) were interchanged within the same PB unit (i.e. along the same axis) due to the difference of random initial parameters of the network. But the PB unit along with the dorsal stream still encoded colour information, while the PB unit along with the ventral stream encoded movement information. The network was still able to show properties of spatio-temporal sequences data in the PB units’ representation.

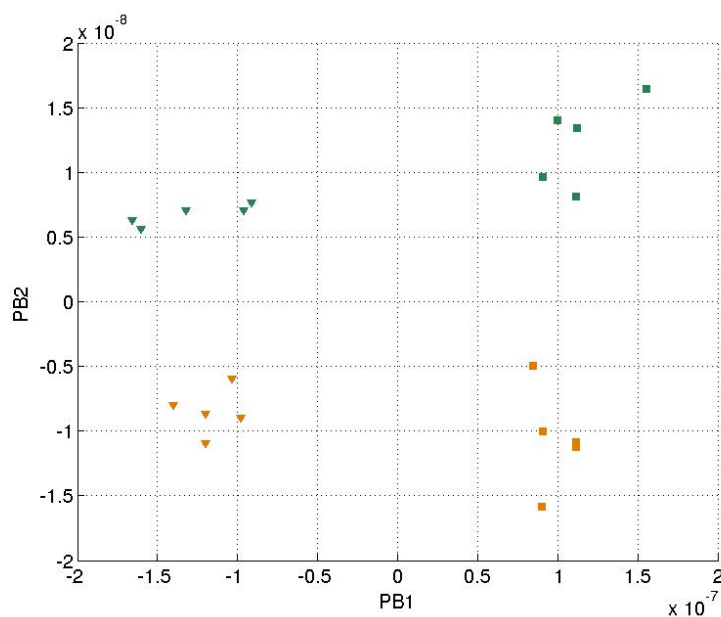


FIGURE 7.10: PB Values with Faster Speed

Values of two sets of PB units in the two streams after training with faster speed. The representation of the markers is the same as in Fig. 7.7.

## 7.4 Summary

In this chapter a recurrent network architecture integrating the RNNPB model and the horizontal product model has been presented, which shows that it is possible

to link the conceptualization of ventral/dorsal visual streams, the emergence of pre-symbolic communication, and the predictive sensorimotor system.

Based on the horizontal product model, the information in the dorsal and ventral streams is separately encoded in two network streams and the predictions of both streams are integrated via the horizontal product while the PB units act as a conceptualization of both streams. These PB units allow for storing multiple sensory sequences. After training, the network is able to recognize the pre-learned conceptualized information and to predict the up-coming visual perception. The network also shows generalization abilities in both ventral and dorsal streams. Therefore, our approach offers preliminary concepts for a similar development of conceptualized language in pre-symbolic communication and further in infants' sensorimotor-stage learning.

## Chapter 8

# Discussion and Conclusion

In this chapter, we will discuss the related issues of this thesis in the context of computer science, cognitive science and neuroscience. Some potential research related to the topic of feedback pathway modelling based on our proposed models will also be discussed. Finally, a summary of this thesis will be given.

### 8.1 Discussion

#### 8.1.1 Hierarchical Action Control

As it has been widely accepted that perception is constituted in a hierarchical way (e.g. [Van Essen & Maunsell, 1983]), we argue that action is executed in a similar way (as e.g. [Rosenbaum, 2009] proposed). From Fig. 8.1 which depicts the general framework of the sensorimotor integration in this thesis, we claim that:

- Execution of an action activates the same representation as if the action is perceived.
- Each upper level has top-down influences to the lower one.
- Lateral connections exist within each level.

As we can see in Fig. 4.2, according to the common coding theory, the action is represented as a homogeneous predictive percept. As we discussed in Chap. 2, Latash et al. [1996] claimed that the control of a specific action leads not merely to a simple movement but to a series of movements. Execution of a single motor primitive spans multiple levels in the neural dimension, with the increasing complexity of the receptive fields from muscle over spine to brain.

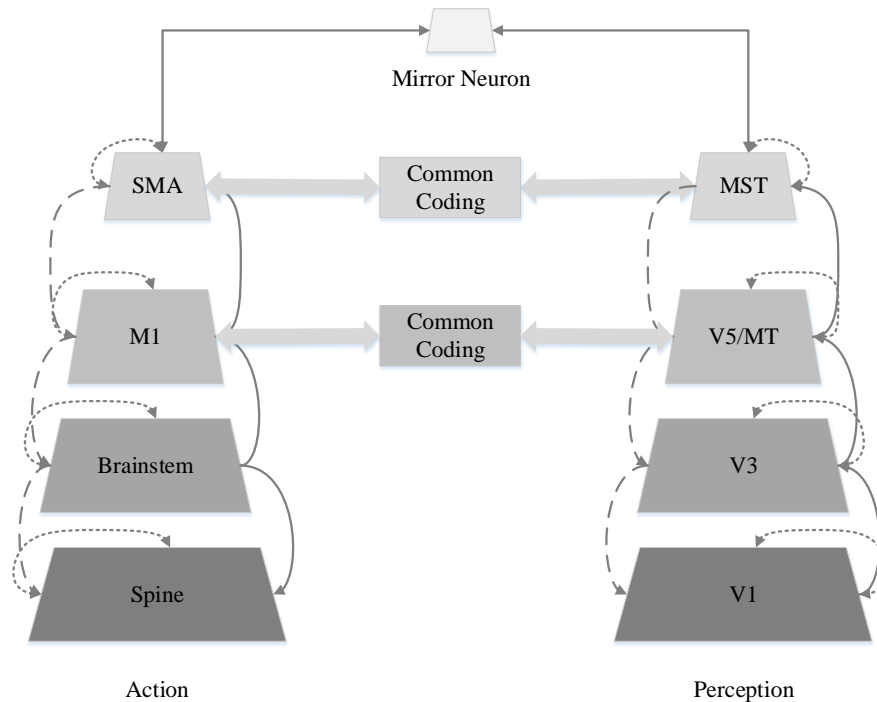


FIGURE 8.1: Hierarchical Perception-Action Model

Generally, there are five essential rules to accomplish a hierarchical action representation [Grafton & de C. Hamilton, 2007]:

- An action includes a series of movements to accomplish a final, temporally distal goal;
- Action (i.e. a series of movements) should be adaptive to the environment;
- The basic movement (motor primitives) can be learnt and retrieved;
- The highest representation is to achieve a desired goal and to solve a problem;
- It is possible to integrate the basic motor elements into a single module.

Despite of the fact that an actual action execution relates to multiple factors, such as musculoskeletal form and function, biomechanics, observations of goal-driven behaviour, a single execution can generally be grouped into musculoskeletal and neural parts.

### 8.1.2 Language Acquisition from Sensorimotor Integration

The visual conceptualization and perception are intertwined processes. Besides the infant development we introduced in Chap. 7, these intertwined processes can also happen during adult learning. As experiments in Schyns & Oliva [1999] showed, when the visual observation is not clear, the brain automatically predicts the visual percept and updates the categorization labels on various levels according to what has been gained from the visual field. On the other hand, this conceptualization also affects the immediate visual perception in a predictive manner. For instance, the conceptualization of a human face predictively spreads conceptualizations on other levels (e.g. face emotion). These feedback pathways propagate from object identity to other local conceptualizations, such as object affordance, motion, edge detection and other processes at the early stages of visual processing. This can be tested by classic illusions, such as ‘figure–ground vase illusion’, where perception depends largely on prior knowledge derived from experiences rather than direct observation.

Therefore, our pre-symbolic representation model in Chap. 6 to some extent also demonstrates the integrated process between conceptualization and spatio-temporal visual perception. This predictive perception may also arouse other visual-based predictive behaviours such as those arising from object permanence. The model explains that, in a hierarchical structure of a sensorimotor system, the high-level conceptualization representation is continuously updated with the partial sensory information perceived in a short-time scale from sensory-driven inputs. Conversely, the feedback pathways accomplish some of the sensorimotor functions by the conceptualized high-level representation of visual perception which is identical to the integration conceptualization and (predictive) visual perception.

The models that we proposed in this thesis are proof-of-concept models focusing on the feasibility of using recurrent connections to model feedback pathways. Therefore, the networks we used can be augmented in two different ways:

- the memory can be enlarged by increasing the number of hidden unit(s);
- more layers can be added to extract more abstract features of the spatio-temporal sequences to form deep-learning architectures.

The language acquisition and mirror neuron theory are also closely related, in the sense that language evolves while the mirror neurons are learning to understand

the intention/meaning of certain motor actions [Hurford, 2004]. It is based on the following two mechanisms.

- Mirror neurons may explain speech imitation during the infant development stages; they imitate the speech from the words they hear [Studdert-Kennedy, 2002]. Also, the firing of mirror neurons trigger the corresponding parts of the somatosensory cortex to configure the vocal tract by executing the muscles' contraction [Goldstein et al., 2006].
- Mirror neurons may support concept representation, especially the mental representation (or simulation) of certain actions. This coding may be represented as the perception anticipation according to the outcome of these actions, based on the common coding theory [Prinz, 1984].

In terms of the model we proposed in Chap. 6, we claim that it was supported by the neuroscience study about the role of the mirror neuron system (MNS) in executing object-oriented-actions (e.g. grasping). This property is similar to the 'data-driven' models such as MNS [Oztop & Arbib, 2002] and MNS2 [Bonaiuto et al., 2007, Bonaiuto & Arbib, 2010], although the main hypothesis in our model is not based on the MNS theory. In the MNS review paper by Oztop et al. [2006], the action generation mode of the RNNPB model was considered to be excessive as there has not been found evidence yet to show that the MNS participates in action generation. However, in our model (as well as the overall architecture) the generation mode has a key role of conceptualized PB units in the sensorimotor integration of object interaction. Nevertheless, the similar network architecture (RNNPB) used in modelling mirror neurons [Tani et al., 2004] and our pre-symbolic sensorimotor integration models may imply a close relationship between language (pre-symbolic) development, object-oriented actions, and the mirror neuron theory.

### 8.1.3 Predictive Perception

The feedback affecting sensory input can be regarded as a kind of predictive information retrieved from the internal memory [J. Anderson & Schooler, 2000]. Based on predictive coding theory, in the hierarchical architecture, the feedback signals (especially the top-down signals) predict the forthcoming sensory input, while the sensory-driven bottom-up signals only deliver the error of the estimation.



The predictive function of feedback pathways is essential as they have the following benefit on the lower-level peripheral perception functions:

- The target of the feedback pathways in perception is applied for sensory prediction. It is realized by extracting cues from multimodal or amodal perception via feature extraction (e.g. by the early visual system) which becomes a prior. Then, the posterior estimation is applied to the next predictive perception.
- If there is a difference between the posterior estimation and the current receptor signals, the percept may be derived as a combination of the two to avoid the fluctuation caused by neuronal or receptor noise. On the other hand, the error signals are also transmitted from bottom-up signals to further act as a prior in the perception cues.

These functions are not independent; instead they are processes that happen at the same time and integrate with each other. They are performed with the similar Bayesian inference and are always interchanging prior knowledge on the cognitive processes level. On various levels of perception, the feedback signals play a role of significant modulation to the lower-levels, indicating that the top-down processes can affect early perceptual cortex and play a role of sensory prediction. Shulman et al. [1997] showed that the top-down feedback selectively modulates attention based on prior knowledge about the feature or feature-related-analysis of the object. Therefore, perception is constructed not only by ‘bottom-up’ external stimuli, but also formed by the internal priori constraints such as expectation, memory or the current goal, although these constraints are also initially learnt from ‘bottom-up’ sensory signals. That is how different kinds of a prior work as a forward model to form the perception prediction via feedback pathways.

Although our models (Chaps. 5 and 6) with lateral connections based on the Elman networks only dealt with one-dimensional motion prediction, it showed that the recurrent connections should be a necessary part to build the predictive role of the feedback connections. Specifically, the recurrent connections had the following roles in these models:

- In both models, recurrent connections record the possible movements by exciting the corresponding units and inhibiting the others. Therefore, the weighting matrix of the recurrent connections deliver a predictive perception in the feedback signals.

- In the model shown in Chap 6, the inhibition of the recurrent connections also plays a role of filtering perception noise. The recurrent weights are learnt from the perception caused from action execution, and are also used for predictive action. Thereby, this model builds a common-coding model with filtering functions.

#### 8.1.4 Robotics as Synthetic Methodology and Neuro-robotics

In the previous chapters, we implemented three models based on recurrent connections. As we mentioned, the main target of these chapters is to prove that they can be used to model feedback pathways in the sensorimotor cortices of artificial cognitive systems. Training of these neural feedback is done by interacting with environment, called embodiment (Fig. 8.2). Based on the embodied cognition, any embodied cognitive system should regard the sensorimotor learning as an interaction process that is shaped by the environment, while the motor action itself also changes the environment in various ways.

Also, this is how we can regard the construction of artificial cognitive systems as a way to examine different hypotheses in neuroscience and cognitive science. This is usually realized by building up models based on such hypotheses and implementing these models in a robotic system. Therefore, using robot as a synthetic methodology is an important way to understand biological principles, although testing such computational models is often constrained by different configurations in executors, sensors and environments.

Furthermore, the embodiment theory also indicates the following procedures of using robots as synthetic methodology:

- Hypotheses are formulated based on findings of biological studies.
- Basic principles should be extracted based on these hypotheses.
- Biologically-inspired but simple enough models can be built to realize these principles.
- In terms of the sensorimotor interaction, it would be more convenient to extract only the relevant information for the model processing, rather than using the raw inputs.

Among all modelling methods for feedback pathways we have shown that the artificial neural network (ANN) is one of the structures to be adopted as it is

comparable to the feedback structures. Also, its information and coding are also similar to the existing feedback pathways in the neural structures, according to cognitive and neuroscience experiments. As a result, the robots which are controlled by ANNs are called ‘neuro-robots’. Nowadays, neuro-robotics research has covered modelling of visual cortex (e.g. [Orabona et al., 2005, Vijayakumar et al., 2001]), auditory localisation (e.g. [Webb, 1995, Liu et al., 2009]) to higher cognition modelling (e.g. [Yamashita & Tani, 2008]) and other sensorimotor behaviours (see e.g. [Wermter et al., 2005] and [Wermter et al., 2014] for edited collections presenting broad samples of models in neuro-robotics). These models offered a



FIGURE 8.2: Cognitive Development  
Cognition is achieved through sensorimotor interaction with the environment<sup>1</sup>.

<sup>1</sup>Copyright by Kira Chow.

complementary insight into understanding various levels of the cognitive and neural systems, by providing a grounded demonstration of neural activities, language acquisition and motor behaviours which have emerged from real-world interaction. Furthermore, from the robotics engineering point of view, an ideally learnt sensorimotor model on a robot can assist to robustly provide flexible interaction between human and environment, to estimate the incoming sensory data in a stable way for a changing sensory environment and even to have its own thoughts. Although some engineering methods could also mimic the process of the sensorimotor interaction, e.g. the potential field method that allows the robot to have obstacle avoidance and path planning [Khatib, 1985, Huang, 2009, Pradhan et al., 2011], their limitations include local minima, linearity assumptions or non-adaptability to environmental changes.

Due to the complexity of building a sensorimotor system that can interact with a fully dynamic environment, we cannot build a complete sensorimotor model with all of the sensorimotor functions. However, from Chap. 5 to Chap. 7, we gradually built three models with different constraints in perception and action. With the framework introduced in Chap. 4, these models are able to interact together with bottom-up sensory-driven signals and feedback signals, by which a sensorimotor integration system is built.

## 8.2 Future Work

### 8.2.1 Conceptor Representation and Mirror Neuron System

As introduced before, with the increasing complexity of the receptive fields from low-level to high-level in the sensorimotor hierarchy, it is possible that only a small number of neurons are activated by a specific percept or a particular motor action on the highest level. This also leads to the hypothesis that the brain may contain ‘grandmother cells’ which is a set of neurons that have a specific receptive field that responds to only one precise stimulus: e.g. the face of the person’s grandmother.

In terms of modelling, interestingly, the PB units in the RNNPB model, which we applied, exhibit a similar property of grandmother cell-like encoding in the sensorimotor system. Likewise, the hypothesis model of ‘conceptor’ is also proposed by Jaeger [2014]. A small number of ‘conceptors’ may control complex human sensorimotor behaviours.

Besides the theory of ‘conceptors’, this high-level representation of action is also related to firing of mirror neurons in the mirror neuron theory. Specifically, the mirror neurons in this theory indicate some neurons in the pre-motor cortex not only responding to the execution of an action, but also firing during the observation of that same action that is executed by others [Rizzolatti et al., 1996]. This function of action recognition as well as action mirroring is embedded in perception, incorporating the action context and the interacting objects and mapping into various parts of the mirror neurons system [Oztop & Arbib, 2002]. From the studies on macaque monkeys, there are three cortical areas forming a hierarchical architecture of the MNS with the reciprocal connections: area F5 and area PF form a premotor-parietal MNS system [Luppino et al., 1999], with the bottom-up inputs about the object features from the STS area as well as a feedback influence [Harries & Perrett, 1991, Seltzer & Pandya, 1994]. In this way, it is possible that the proposed MNS performs action understanding on multiple levels via reciprocal connections [Hamilton et al., 2007, Kilner et al., 2007]:

- The intention level that determines the long-term goals from the agent itself or anticipates them from observing others;
- The goal level that describes the short-term goals related to achievement of the long-term intention;
- The kinematic level that controls the spatio-temporal inverse kinematics of the actuators;
- The muscle level that governs the meta-patterns of the cells of the muscle tissues required to execute the actions.

The whole system is thus always attempting to understand one’s short-term goals as well as the long-term intention by observing its movements, or it is attempting to control and to adjust its own movement by changing short-term goals contextually according to the long-term intention by the integration of the feedback and bottom-up influences. As Iacoboni et al. [2005] suggested: the mirror neurons encode goal-directed actions; when an action is observed by another individual, the feedback information comes from the inferior parietal lobule (IPL), encoding the final goals of the actions, i.e. ‘why’ the actor is doing it. For example, when an agent is grasping an apple with a certain series of motor actions, it can gradually be understood that it is grasping-for-eating but not grasping-for-placing, i.e. its intention of actions, according to the trivial difference actions and objects.

Specifically, this action understanding role is accomplished by both bottom-up mechanisms arising in early visual areas and feedback mechanisms arising in pre-frontal cortex. In this way, the feedback pathways encoding certain intention modulates the forthcoming action understanding and recognition. Moreover, the functions of the MNS may also include intention understanding [Iacoboni et al., 2005] as well as other action-related functions such as action imitation [Rizzolatti et al., 2001] and even language acquisition [Rizzolatti & Arbib, 1998].

The findings and hypotheses introduced above about the mirror neuron theory conclude the properties of the related theories about ‘conceptor’ neurons:

- These neurons representing some kinds of ‘profile’ of sensorimotor primitives work in a small-dimensional space to encode the sensorimotor sequences. The units can be regarded to code high-level representations such as action goals. The dimension of this space should not be as large as that in musculoskeletal movement, since it is more convenient to execute the reasoning to select an action given the a specific goal, when it is a low dimensional space to encode goals, motor commands, etc.
- On the other hand, the representation on the lower level becomes more and more simple. On the lowest level, the sensorimotor sequences being encoded are composed from musculoskeletal movements, which are explicitly encoded in their corresponding motor cortices with a high number of dimensions;
- Corresponding to the continuous space of sensorimotor primitives, such an expression space of ‘profile’ should also be continuous.
- Importantly, the ideas of ‘conceptor’, mirror neurons, etc. are attempting to connect the missing link between the symbolic (conceptualized) representation and the non-linear sensorimotor sequences, because the highest level of abstract representation of sensorimotor sequences should be essential for reasoning, memory and other cognitive processes.

As a proof-of-concept model, the RNNPB can fulfil the above requirements given that

- There is a limited number of sensorimotor sequences.
- Only two levels of representations (i.e. conceptual and sensorimotor representations) are needed in the whole architecture.

Although it has constraints, the RNNPB can realize several hierarchical sensorimotor functions due to its flexibility in neural dynamics:

- It can generate multiple sequences with the implication of high-level representation: since a recurrent network has the ability to approximate the exact Bayesian processes [Freaun et al., 2006], it is theoretically feasible for a recurrent network with a fixed number of hidden layers to train with a finite set of spatio-temporal sequences of movement data.
- Likewise, the low-dimensional PB space can be trained as a continuous expression space with behaviour data.

Admittedly, other models can also be employed, if there are different constraints in the process of building the hierarchical sensorimotor architecture or in the artificial system itself. Modelling the brain of a fly only needs a smaller capacity of memory (e.g. [Greenspan & Van Swinderen, 2004]). On the other hand, a long-term memory is also needed for modelling a human brain if memory is involved in the feedback signal. Thus, it implies that models that are similar to the RNNPB, HoRNNPB and the conceptor models can be further developed as a hierarchical architecture, where a higher-level representation (e.g. PB units) extracts information from lower levels.

## 8.2.2 Deep Learning and Predictive Coding

In this thesis, we mainly developed three-layer recurrent networks to examine the functions of recurrent connections. Particularly, in terms of the sensorimotor levels they represent, the models in Chaps. 5 and 6 constitute feedback (lateral) connections within the same level in the sensorimotor hierarchy, while the HoRNNPB model in Chap 7 constructs a two-level architecture in the sensorimotor hierarchy.

Admittedly, three-layer networks could not extract enough sensorimotor representations. They cannot cover all the variability in realization of predictive perception either. However, the three models we proposed suggest that feedback connections with a more sophisticated deep architecture can be beneficial in terms of extracting sensorimotor features from complex structures [Bengio, 2009]. Similar to the hierarchical organization of cortices we mentioned in Chap. 2, the deep structure cannot only extract complex structure and build internal representation from the regularities of the sensory inputs, but also plays a role in transforming the neuronal spike waveforms to a very abstract level of representation. Nevertheless, we

assert that the PB units in the RNNPB, or models with similar neural dynamics, can also realize the function of extracting a statistical regularity on a higher level.

Future experiments can be pursued by stacking several similar statistical models into a deep architecture which is learnt in a self-organized way by the feedback error (for detailed reviews see also [Bengio, 2009, Bengio et al., 2013]). The self-organization can be inspired from the cognitive science theory that the perceptual world is actually constructed by the error signals plus our top-down prediction; the error signals are caused by the mismatch and are transmitted and distributed on various levels from the ‘predictive coding’ theory.

Also, it should be promising to build a hybrid architecture by stacking ARNN models (especially the PB models) with a deep learning structure, where the abstract representation can form another symbolic representation on a cognitive level in a hierarchical way (Fig. 8.3).

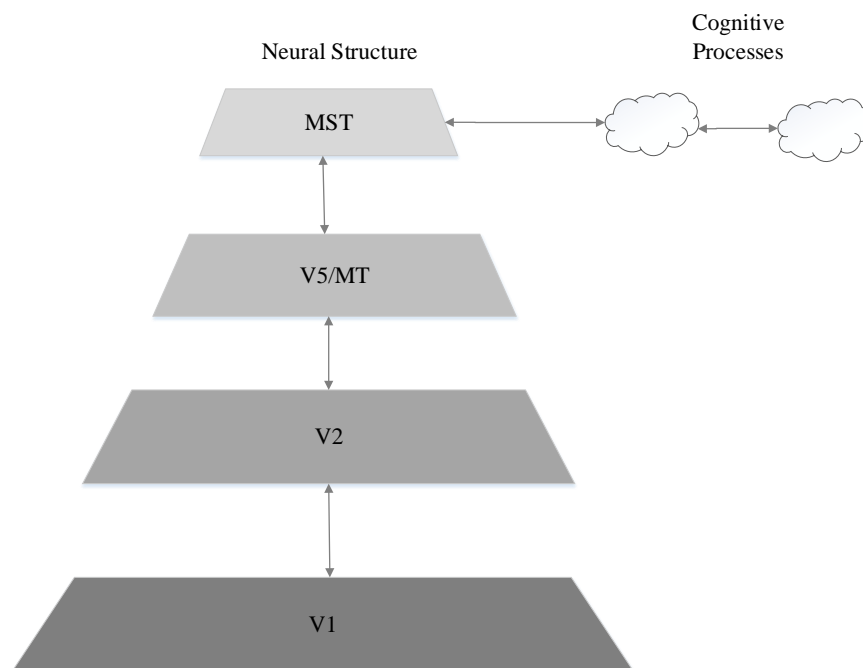


FIGURE 8.3: Neural and Cognitive Model by RNNPB with Deep Structure

Each layer can be regarded as an abstract representation of the lower one. At the uppermost layer there are different kinds of abstract concept representation, which also relate to various cognitive processes.



### 8.3 Conclusion

Feedback pathways exist on different levels of the hierarchical cortical areas (especially visual and motor cortices). Concerning their functions in the sensorimotor system, we raised four main questions at the beginning of this thesis (Chap. 1).

To answer these questions, we investigated the feedback signals in biological systems. Having it formulated as Bayesian inference, we model those signals in artificial recurrent connections. Consequently, we conducted three experiments using various types of recurrent neural networks to examine their feasibility and performance in artificial systems.

In the first experiment, we focused on the predictive encoding of information in the primary visual cortex. To model this, we designed a recurrent predictive network with a horizontal product where the information of object feature and object movement becomes successfully separated in its two hidden layers. This experiment demonstrates how the recurrent network constructs lateral connections which allow to accomplish predictive function in both visual pathways.

Since the dorsal pathway also allows for motor-relevant representations, in the second experiment, we also developed such recurrent connections for sensory latency compensation which also supports smoother and faster behaviours in sensorimotor integration tasks. For example, the latency of the sensorimotor cycle of a robot may affect the response time for the motor action. We expanded the use of recurrent connections in the sensorimotor system, particularly in the sensory prediction part, so that the recurrent connections can compensate the delay in the sensory percepts. A continuous actor-critic automaton (CACLA) was developed for the generation of smooth behaviours corresponding to the predictive sensory percepts. Experiments showed that the predictive sensorimotor architecture successfully increases the speed and robustness of the robot docking experiment.

The recognition and prediction functions are not independent processes, but they also integrate and assist each other in a hierarchical way. Furthermore, we propose that they result in the development of pre-symbolic communication. Therefore, in the last experiment, we proposed that the learning of the visual pathways also leads to the conceptualization of visual information. This was realized by a horizontal recurrent network with parametric bias (HoRNNPB). We examined this model through a robot passively observing an object to learn its features and movements. During the learning process of observing sensorimotor primitives,

i.e. observing a set of trajectories of arm movements and its oriented object features, the pre-symbolic representation was self-organized in the parametric units. These representational units acted as bifurcation parameters guiding the robot to recognize and predict various learnt sensorimotor primitives.

The above three experiments also examine our proposed feedback common-coding sensorimotor framework (Chap. 4) in different perspectives:

- The bidirectional pathways maintain the perception representation in the neural dynamics, within a single level (Chap. 5) and across two levels (Chap. 7). In Chap. 5, we focus on the feedback pathways modelling in the dorsal and ventral streams on a single level. We only model the lateral connections with the Elman networks, but these feedback pathways can be also modelled by other recurrent networks, in which the higher level represents memory, expectation and other top-down prior experience too. The model in Chap. 7 further elaborates the model in Chap. 5 in a hierarchical way.
- As we mentioned in Chap. 4, the posteriors of action motion and perception are inferred incrementally. At the same time, in the loop of sensorimotor inferences, the linkage of common coding is learnt. Chap.6 showed the basic principle of this learning loop ( $e \rightarrow a \rightarrow e$ ) by a proof-of-concept recurrent model.
- The common coding is learnt by either perception or action or both in an unsupervised way; in this thesis, it is represented as shared weighted connections, which are also modulated by the feedback information. In Chap. 6, we build a linkage between perception and action with the PAM model. Specifically, the learning of action-oriented visual perception is driven by the motor actions.

To sum up, all of the three recurrent connections based models we used in these experiments are built to demonstrate the role of feedback information in the context of predictive sensorimotor integration. This information transmitted in the feedback pathways comes from higher-levels of neural activities as well as different kinds of cognitive processes.

The main scientific contribution is that in order to answer the questions about feedback pathways in biological and artificial cognitive systems we raised in Chap. 1, we designed three novel artificial neural models for artificial cognitive systems based on the investigations about feedback signals on biological systems. After

having conducting experiments based on these models, we answer the questions as follows,

- First, the feedback pathways in a hierarchical sensorimotor structure account for the predictive functions in perception and action on each layer. Since such prediction cannot only account for some phenomena such as flash-lag effect and retina prediction, but also result in a smoother and robust sensorimotor integration, we propose that the feedback signals deliver a prior knowledge to mediate the neural activities. Some of the feedback pathways may come from representation of cognitive processes, such as memory and expectation. These processes can be regarded as conceptualised representation, or a symbolic representation, which is learnt during hierarchical sensorimotor processes as prior knowledge. For instance, such a higher-level symbolic representation can also be considered as an immature stage for language acquisition.
- Second, the feedback signal can be modelled as a Bayesian inference within the common coding framework. By implementing Elman recurrent connections to model feedback pathways, these models exhibit some of complementary cognitive functions seen in biological systems. Although the properties of feedback pathways vary in the brain, various types of recurrent connections that are similar to Elman networks may also be able to form a Bayesian inference from prior knowledge.
- Third, several sensorimotor functions which can be considered related to the feedback pathways can be realised on cognitive robots by Elman connections, such as sensory prediction, its corresponding action prediction and pre-symbolic emergence. Experiments on real robots or in simulation demonstrate that recurrent connections either improve the performance in adaptivity or provide cognitive capabilities such as pre-symbolic representation. Feedback pathways are beneficial to artificial cognitive systems, also because they allow a flexible and fast sensorimotor integration to the cognitive systems.

# Bibliography

- Alais, D., & Blake, R. (2005). *Binocular rivalry*. MIT press.
- Alexander, G. E., DeLong, M. R., & Strick, P. L. (1986). Parallel organization of functionally segregated circuits linking basal ganglia and cortex. *Annual Review of Neuroscience*, 9(1), 357–381.
- Amari, S.-i. (1977). Dynamics of pattern formation in lateral-inhibition type neural fields. *Biological Cybernetics*, 27(2), 77–87.
- Ambrose, R. O., Aldridge, H., Askew, R. S., Burridge, R. R., Bluethmann, W., Diftler, M., ... Rehnmark, F. (2000). Robonaut: NASA's space humanoid. *IEEE Intelligent Systems*, 15(4), 57–63.
- Anderson, C., Van Essen, D., & Olshausen, B. (2005). Directed visual attention and the dynamic control of information flow. *Neurobiology of Attention.*, 11–17.
- Anderson, J., & Schooler, L. (2000). The adaptive nature of memory. *The Oxford Handbook of Memory*.
- Angelucci, A., Levitt, J. B., Walton, E. J., Hupe, J.-M., Bullier, J., & Lund, J. S. (2002). Circuits for local and global signal integration in primary visual cortex. *The Journal of Neuroscience*, 22(19), 8633–8646.
- Asada, M., MacDorman, K. F., Ishiguro, H., & Kuniyoshi, Y. (2001). Cognitive developmental robotics as a new paradigm for the design of humanoid robots. *Robotics and Autonomous Systems*, 37(2), 185–193.
- Ballard, D. (1981). Generalizing the hough transform to detect arbitrary shapes. *Pattern Recognition*, 13(2), 111–122.
- Bar, M. (2004). Visual objects in context. *Nature Reviews Neuroscience*, 5(8), 617–629.

- Bar, M., Tootell, R. B., Schacter, D. L., Greve, D. N., Fischl, B., Mendola, J. D., ... Dale, A. M. (2001). Cortical mechanisms specific to explicit visual object recognition. *Neuron*, *29*(2), 529–535.
- Barlow, H. B. (1953). Summation and inhibition in the frog's retina. *The Journal of Physiology*, *119*(1), 69–88.
- Bassano, D. (2000). Early development of nouns and verbs in French: Exploring the interface between lexicon and grammar. *Journal of Child Language*, *27*(3), 521–559.
- Bay, H., Tuytelaars, T., & Van Gool, L. (2006). SURF: Speeded up robust features. In *Computer Vision—ECCV 2006* (pp. 404–417). Springer.
- Beer, R. D. (1995). A dynamical systems perspective on agent-environment interaction. *Artificial Intelligence*, *72*(1), 173–215.
- Bengio, Y. (2009). Learning deep architectures for AI. *Foundations and Trends in Machine Learning*, *2*(1), 1–127.
- Bengio, Y., Courville, A., & Vincent, P. (2013). Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *35*(8), 1798–1828.
- Bengio, Y., Simard, P., & Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, *5*(2), 157–166.
- Bergmann, U., & von der Malsburg, C. (2011). Self-organization of topographic bilinear networks for invariant recognition. *Neural Computation*, 1–28.
- Bernstein, N. (1967). Co-ordination and regulation of movements.
- Biederman, I. (1987). Recognition-by-components: a theory of human image understanding. *Psychological Review*, *94*(2), 115.
- Blake, R. (2001). A primer on binocular rivalry, including current controversies. *Brain and Mind*, *2*(1), 5–38.
- Bloom, L., Tinker, E., & Margulis, C. (1993). The words children learn: Evidence against a noun bias in early vocabularies. *Cognitive Development*, *8*(4), 431–450.

- Bluethmann, W., Ambrose, R., Diftler, M., Askew, S., Huber, E., Goza, M., ... Magruder, D. (2003). Robonaut: A robot designed to work with humans in space. *Autonomous Robots*, *14*(2-3), 179–197.
- Bonaiuto, J., & Arbib, M. (2010). Extending the mirror neuron system model, II: what did I just do? a new role for mirror neurons. *Biological Cybernetics*, *102*(4), 341–359.
- Bonaiuto, J., Rosta, E., & Arbib, M. (2007). Extending the mirror neuron system model, I. *Biological Cybernetics*, *96*(1), 9–38.
- Bonato, V., Marques, E., & Constantinides, G. (2009). A floating-point extended Kalman filter implementation for autonomous mobile robots. *Journal of Signal Processing Systems*, *56*, 41-50.
- Bonniec, G. P.-L. (1985). From visual-motor anticipation to conceptualization: Reaction to solid and hollow objects and knowledge of the function of containment. *Infant Behavior and Development*, *8*(4), 413–424.
- Boole, G. (1854). *An investigation of the laws of thought: on which are founded the mathematical theories of logic and probabilities* (Vol. 2). Walton and Maberly.
- Bowerman, M., & Levinson, S. (2001). *Language acquisition and conceptual development* (Vol. 3). Cambridge University Press.
- Brooks, R. A. (1991). Intelligence without representation. *Artificial Intelligence*, *47*(1), 139–159.
- Brown, H., Friston, K., & Bestmann, S. (2011). Active inference, attention, and motor preparation. *Frontiers in Psychology*, *2*, 1-10.
- Buccino, G., Binkofski, F., Fink, G. R., Fadiga, L., Fogassi, L., Gallese, V., ... Freund, H.-J. (2001). Action observation activates premotor and parietal areas in a somatotopic manner: an fMRI study. *European Journal of Neuroscience*, *13*(2), 400–404.
- Buehner, M., & Young, P. (2006). A tighter bound for the echo state property. *IEEE Transactions on Neural Networks*, *17*(3), 820–824.
- Bullmore, E., & Sporns, O. (2009). Complex brain networks: graph theoretical analysis of structural and functional systems. *Nature Reviews Neuroscience*, *10*(3), 186–198.

- Bussey, T., & Saksida, L. (2007). Memory, perception, and the ventral visual-perirhinal-hippocampal stream: Thinking outside of the boxes. *Hippocampus*, *17*(9), 898–908.
- Calais-Germain, B., & Lamotte, A. (1996). *Anatomy of movement exercises*. Eastland Press.
- Carruthers, P. (2002). The cognitive functions of language. *Behavioral and Brain Sciences*, *25*(06), 657–674.
- Case, L. K., Gosavi, R., & Ramachandran, V. S. (2013). Heightened motor and sensory (mirror-touch) referral induced by nerve block or topical anesthetic. *Neuropsychologia*, *51*(10), 1823–1828.
- Chemla, E., Mintz, T., Bernal, S., & Christophe, A. (2009). Categorizing words using ‘frequent frames’: what cross-linguistic analyses reveal about distributional acquisition strategies. *Developmental Science*, *12*(3), 396–406.
- Chen, B. L., Hall, D. H., & Chklovskii, D. B. (2006). Wiring optimization can relate neuronal structure and function. *Proceedings of the National Academy of Sciences of the United States of America*, *103*(12), 4723–4728.
- Clark, A. (2012). Whatever next? Predictive brains, situated agents, and the future of cognitive science. *Behavioral Brain Sciences*, 1–86.
- Clune, J., Mouret, J.-B., & Lipson, H. (2013). The evolutionary origins of modularity. *Proceedings of the Royal Society B: Biological sciences*, *280*(1755), 20122863.
- Colby, C. L., Duhamel, J.-R., & Goldberg, M. E. (1992). Posterior parietal cortex and retinocentric space. *Behavioral and Brain Sciences*, *15*(04), 727–728.
- Committeri, G., Galati, G., Paradis, A., Pizzamiglio, L., Berthoz, A., & LeBihan, D. (2004). Reference frames for spatial cognition: different brain areas are involved in viewer-, object-, and landmark-centered judgments about object location. *Journal of Cognitive Neuroscience*, *16*(9), 1517–1535.
- Coogan, T. A., & Burkhalter, A. (1993). Hierarchical organization of areas in rat visual cortex. *The Journal of Neuroscience*, *13*, 3749–3749.
- Corbetta, D., & Snapp-Childs, W. (2009). Seeing and touching: the role of sensory-motor experience on the development of infant reaching. *Infant Behavior and Development*, *32*(1), 44–58.

- Corbetta, M. (1998). Frontoparietal cortical networks for directing attention and the eye to visual locations: identical, independent, or overlapping neural systems? *Proceedings of the National Academy of Sciences*, *95*(3), 831–838.
- Cordeschi, R. (2002). *The discovery of the artificial: Behavior, mind and machines before and beyond cybernetics* (Vol. 28). Springer.
- Cuijpers, R., Stuijt, F., & Sprinkhuizen-Kuyper, I. (2009). Generalisation of action sequences in RNNPB networks with mirror properties. In *Proceedings of the European Symposium on Neural Networks (ESANN)* (p. 251-256).
- Cullen, K. (2004). Sensory signals during active versus passive movement. *Current Opinion in Neurobiology*, *14*(6), 698.
- Damasio, A. R., & Tranel, D. (1993). Nouns and verbs are retrieved with differently distributed neural systems. *Proceedings of the National Academy of Sciences*, *90*(11), 4957–4960.
- Darrin, C., Christopher, G., Aleš, U., & Cheng, G. (2004). Learning to act from observation and practice. *International Journal of Humanoid Robotics*, *1*(04), 585–611.
- Dayan, P., & Hinton, G. (1996). Varieties of Helmholtz machine. *Neural Networks*, *9*(8), 1385–1403.
- Dayan, P., Hinton, G., Neal, R., & Zemel, R. (1995). The Helmholtz machine. *Neural Computation*, *7*(5), 889–904.
- Desimone, R. (1996). Neural mechanisms for visual memory and their role in attention. *Proceedings of the National Academy of Sciences*, *93*(24), 13494–13499.
- Desimone, R. (1998). Visual attention mediated by biased competition in extrastriate visual cortex. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, *353*(1373), 1245–1255.
- Egner, T., Monti, J. M., & Summerfield, C. (2010). Expectation and surprise determine neural population responses in the ventral visual stream. *The Journal of Neuroscience*, *30*(49), 16601–16608.
- Eisenberg, M., Shmuelof, L., Vaadia, E., & Zohary, E. (2010). Functional organization of human motor cortex: directional selectivity for movement. *The Journal of Neuroscience*, *30*(26), 8897–8905.



- Elman, J. (1990). Finding structure in time. *Cognitive Science*, *14*(2), 179–211.
- Favorov, O., & Whitsel, B. L. (1988). Spatial organization of the peripheral input to area 1 cell columns. I. The detection of ‘segregates’. *Brain Research Reviews*, *13*(1), 25–42.
- Favorov, O. V., & Diamond, M. E. (1990). Demonstration of discrete place-defined columns ‘segregates’ in the cat SI. *Journal of Comparative Neurology*, *298*(1), 97–112.
- Felleman, D. J., & Van Essen, D. C. (1991). Distributed hierarchical processing in the primate cerebral cortex. *Cerebral Cortex*, *1*(1), 1–47.
- Fields, C. (2011). Trajectory recognition as the basis for object individuation: a functional model of object file instantiation and object-token encoding. *Frontiers in Psychology*, *2*.
- Fletcher, P. C., & Frith, C. D. (2008). Perceiving is believing: A Bayesian approach to explaining the positive symptoms of schizophrenia. *Nature Reviews Neuroscience*, *10*(1), 48–58.
- Fogassi, L., Ferrari, P. F., Gesierich, B., Rozzi, S., Chersi, F., & Rizzolatti, G. (2005). Parietal lobe: from action organization to intention understanding. *Science*, *308*(5722), 662–667.
- Földiák, P. (1991). Learning invariance from transformation sequences. *Neural Computation*, *3*, 194–200.
- Foresti, G. (1999). Object recognition and tracking for remote video surveillance. *IEEE Transactions on Circuits and Systems for Video Technology*, *9*(7), 1045–1062.
- Frean, M., Lilley, M., & Boyle, P. (2006). Implementing Gaussian process inference with neural networks. *International Journal of Neural Systems*, *16*(05), 321–327.
- Freeman, W. T., & Tenenbaum, J. B. (1997). Learning bilinear models for two-factor problems in vision. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 1997* (pp. 554–560).
- Fries, W. (1990). Pontine projection from striate and prestriate visual cortex in the macaque monkey: An anterograde study. *Visual Neuroscience*, *4*(03), 205–216.

- Friston, K. (2003). Learning and inference in the brain. *Neural Networks*, *16*(9), 1325–1352.
- Friston, K. (2005). A theory of cortical responses. *Philosophical Transactions of the Royal Society B: Biological Sciences*, *360*(1456), 815–836.
- Fujita, I., Tanaka, K., Ito, M., & Cheng, K. (1992). Columns for visual features of objects in monkey inferotemporal cortex. *Nature*, *360*(6402), 343–346.
- Fukushima, K. (1980). Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological Cybernetics*, *36*(4), 193–202.
- Gallese, V., Fadiga, L., Fogassi, L., & Rizzolatti, G. (1996). Action recognition in the premotor cortex. *Brain*, *119*(2), 593–609.
- Gates, B. (2007). A robot in every home. *Scientific American*, *296*(1), 58–65.
- Gelman, R., & Spelke, E. (1981). The development of thoughts about animate and inanimate objects: Implications for research on social cognition. *Social Cognitive Development: Frontiers and Possible Futures*, 43–66.
- Gentner, D. (1982). *Why nouns are learned before verbs: Linguistic relativity versus natural partitioning* (Tech. Rep.).
- Gibson, E. (1988). Exploratory behavior in the development of perceiving, acting, and the acquiring of knowledge. *Annual Review of Psychology*.
- Gilbert, C. D., & Li, W. (2013). Top-down influences on visual processing. *Nature Reviews Neuroscience*, *14*(5), 350–363.
- Gilbert, C. D., & Sigman, M. (2007). Brain states: top-down influences in sensory processing. *Neuron*, *54*(5), 677–696.
- Goldstein, L., Byrd, D., & Saltzman, E. (2006). The role of vocal tract gestural action units in understanding the evolution of phonology. *Action to Language via the Mirror Neuron System*, 215–249.
- Goodale, M. A., & Milner, A. D. (1992). Separate visual pathways for perception and action. *Trends in Neurosciences*, *15*(1), 20–25.
- Grafton, S. T., & de C. Hamilton, A. F. (2007). Evidence for a distributed hierarchy of action representation in the brain. *Human Movement Science*, *26*(4), 590–616.

- Graham, J. (1982). Some topographical connections of the striate cortex with subcortical structures in macaca fascicularis. *Experimental Brain Research*, 47(1), 1–14.
- Graves, A., Liwicki, M., Fernández, S., Bertolami, R., Bunke, H., & Schmidhuber, J. (2009). A novel connectionist system for unconstrained handwriting recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(5), 855–868.
- Greenspan, R. J., & Van Swinderen, B. (2004). Cognitive consonance: complex brain functions in the fruit fly and its relatives. *Trends in Neurosciences*, 27(12), 707–711.
- Greenwald, A. G. (1970). Sensory feedback mechanisms in performance control: with special reference to the ideo-motor mechanism. *Psychological Review*, 77(2), 73.
- Gregory, R. L. (1980). Perceptions as hypotheses. *Philosophical Transactions of the Royal Society of London. B, Biological Sciences*, 290(1038), 181–197.
- Grisetti, G., Tipaldi, G. D., Stachniss, C., Burgard, W., & Nardi, D. (2007). Fast and accurate slam with rao-blackwellized particle filters. *Robotics and Autonomous Systems*, 55(1), 30–38.
- Gross, C. G. (2002). Genealogy of the grandmother cell. *The Neuroscientist*, 8(5), 512–518.
- Grossberg, S. (1982). Contour enhancement, short term memory, and constancies in reverberating neural networks. In *Studies of Mind and Brain* (pp. 332–378). Springer.
- Gruber, T., & Olsen, G. (1994). An ontology for engineering mathematics. In *Fourth International Conference on Principles of Knowledge Representation and Reasoning* (Vol. 94, pp. 258–269).
- Habel, C., & Tappe, H. (1999). *Processes of segmentation and linearization in describing events*.
- Hamilton, A. F. d. C., Grafton, S. T., & Hamilton, A. (2007). The motor hierarchy: from kinematics to goals and intentions. *Sensorimotor Foundations of Higher Cognition*, 381–408.

- Harries, M., & Perrett, D. (1991). Visual processing of faces in temporal cortex: Physiological evidence for a modular organization and possible anatomical correlates. *Journal of Cognitive Neuroscience*, *3*(1), 9–24.
- Hauk, O., Johnsrude, I., & Pulvermüller, F. (2004). Somatotopic representation of action words in human motor and premotor cortex. *Neuron*, *41*(2), 301–307.
- Hawkins, J. (2004). *On intelligence*. Macmillan.
- Herwig, A., & Schneider, W. X. (2014). Predicting object features across saccades: Evidence from object recognition and visual search. *Journal of Experimental Psychology: General*, *143*(5), 1903–1922.
- Hinton, G. E. (2002). Training products of experts by minimizing contrastive divergence. *Neural Computation*, *14*(8), 1771–1800.
- Hinton, G. E., & Sejnowski, T. J. (1983). Optimal perceptual inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 448–453).
- Hirel, J., Gaussier, P., & Quoy, M. (2011). Biologically inspired neural networks for spatio-temporal planning in robotic navigation tasks. In *IEEE International Conference on Robotics and Biomimetics, ROBIO* (pp. 1627–1632).
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, *9*(8), 1735–1780.
- Hodgkin, A. L., & Huxley, A. F. (1952). A quantitative description of membrane current and its application to conduction and excitation in nerve. *The Journal of Physiology*, *117*(4), 500.
- Hohwy, J. (2007). Functional integration and the mind. *Synthese*, *159*(3), 315–328.
- Holland, O. (2003). The first biologically inspired robots. *Robotica*, *21*(04), 351–363.
- Hollerbach, J. M. (1982). Computers, brains and the control of movement. *Trends in Neurosciences*, *5*, 189–192.
- Hosoya, T., Baccus, S. A., & Meister, M. (2005). Dynamic predictive coding by the retina. *Nature*, *436*(7047), 71–77.

- Huang, L. (2009). Velocity planning for a mobile robot to track a moving target: a potential field approach. *Robotics and Autonomous Systems*, 57(1), 55–63.
- Hubel, D., & Wiesel, T. (1963). Shape and arrangement of columns in cat's striate cortex. *The Journal of Physiology*, 165(3), 559–568.
- Hurford, J. R. (2004). Language beyond our grasp: what mirror neurons can, and cannot, do for the evolution of language. *Evolution of Communication Systems*, 297–314.
- Hyvärinen, A., & Hoyer, P. (2000). Emergence of phase- and shift-invariant features by decomposition of natural images into independent feature subspaces. *Neural Computation*, 12(7), 1705–1720.
- Iacoboni, M., Molnar-Szakacs, I., Gallese, V., Buccino, G., Mazziotta, J. C., & Rizzolatti, G. (2005). Grasping the intentions of others with one's own mirror neuron system. *PLoS Biology*, 3(3), e79.
- Ishai, A. (2008). Let's face it: It's a cortical network. *NeuroImage*, 40(2), 415–419.
- Itti, L., & Koch, C. (2001). Computational modelling of visual attention. *Nature Reviews Neuroscience*, 2(3), 194–203.
- Izawa, J., & Shadmehr, R. (2011). Learning from sensory and reward prediction errors during motor adaptation. *PLoS Computational Biology*, 7(3), e1002012.
- Jacobs, R., Jordan, M., & Barto, A. (1991). Task decomposition through competition in a modular connectionist architecture: The what and where vision tasks. *Cognitive Science*, 15(2), 219–250.
- Jaeger, H. (2001). The 'echo state' approach to analysing and training recurrent neural networks—with an erratum note. *Bonn, Germany: German National Research Center for Information Technology GMD Technical Report*, 148.
- Jaeger, H. (2002). *Tutorial on training recurrent neural networks, covering BPPT, RTRL, EKF and the 'echo state network' approach*. GMD-Forschungszentrum Informationstechnik.
- Jaeger, H. (2014). Controlling recurrent neural networks by conceptors. *CoRR*, abs/1403.3369.
- Jaeger, H., & Haas, H. (2004). Harnessing nonlinearity: Predicting chaotic systems and saving energy in wireless communication. *Science*, 304(5667), 78–80.

- James, W. (1890). The consciousness of self. *The Principles of Psychology*, 8.
- Jordan, M. I. (1997). Serial order: A parallel distributed processing approach. *Advances in Psychology*, 121, 471–495.
- Kaas, J. H. (1987). The organization of neocortex in mammals: Implications for theories of brain function. *Annual Review of Psychology*, 38(1), 129–151.
- Takei, S., Hoffman, D. S., & Strick, P. L. (2001). Direction of action is represented in the ventral premotor cortex. *Nature Neuroscience*, 4(10), 1020–1025.
- Kalman, R. E. (1960). A new approach to linear filtering and prediction problems. *Journal of basic Engineering*, 82(1), 35–45.
- Kanwisher, N., McDermott, J., & Chun, M. M. (1997). The fusiform face area: A module in human extrastriate cortex specialized for face perception. *The Journal of Neuroscience*, 17(11), 4302–4311.
- Kashtan, N., & Alon, U. (2005). Spontaneous evolution of modularity and network motifs. *Proceedings of the National Academy of Sciences of the United States of America*, 102(39), 13773–13778.
- Kashtan, N., Noor, E., & Alon, U. (2007). Varying environments can speed up evolution. *Proceedings of the National Academy of Sciences*, 104(34), 13711–13716.
- Kastner, S., & Ungerleider, L. G. (2001). The neural basis of biased competition in human visual cortex. *Neuropsychologia*, 39(12), 1263–1276.
- Kawato, M., Kuroda, T., Imamizu, H., Nakano, E., Miyauchi, S., & Yoshioka, T. (2003). Internal forward models in the cerebellum: fMRI study on grip force and load force coupling. *Progress in Brain Research*, 142, 171–188.
- Kersten, A. W. (1998). An examination of the distinction between nouns and verbs: Associations with two different kinds of motion. *Memory & cognition*, 26(6), 1214–1232.
- Keysers, C., Kaas, J. H., & Gazzola, V. (2010). Somatosensation in social perception. *Nature Reviews Neuroscience*, 11(6), 417–428.
- Khatib, O. (1985). Real-time obstacle avoidance for manipulators and mobile robots. In *Proc. IEEE International Conference on Robotics and Automation* (Vol. 2, pp. 500–505).

- Kilner, J. M., Friston, K. J., & Frith, C. D. (2007). Predictive coding: an account of the mirror neuron system. *Cognitive Processing*, 8(3), 159–166.
- Kisvarday, Z., Toth, E., Rausch, M., & Eysel, U. (1997). Orientation-specific relationship between populations of excitatory and inhibitory lateral connections in the visual cortex of the cat. *Cerebral Cortex*, 7(7), 605–618.
- Kleesiek, J., Badde, S., Wermter, S., & Engel, A. K. (2013). Action-driven perception for a humanoid. In *Agents and Artificial Intelligence* (pp. 83–99). Springer.
- Klein, T., Jeka, J., Kiemel, T., & Lewis, M. (2012). Navigating sensory conflict in dynamic environments using adaptive state estimation. *Biological Cybernetics*, 1–14.
- Köster, U., Lindgren, J., Gutmann, M., & Hyvärinen, A. (2009). Learning natural image structure with a horizontal product model. *Independent Component Analysis and Signal Separation*, 507–514.
- Kuffler, S. W., et al. (1953). Discharge patterns and functional organization of mammalian retina. *Journal of Neurophysiology*, 16(1), 37–68.
- Kumar, S., Sedley, W., Nourski, K. V., Kawasaki, H., Oya, H., Patterson, R. D., ... Griffiths, T. D. (2011). Predictive coding and pitch processing in the auditory cortex. *Journal of Cognitive Neuroscience*, 23(10), 3084–3094.
- Kuppuswamy, N., & Harris, C. M. (2013). Developing learnability—the case for reduced dimensionality. In *2013 IEEE Third Joint International Conference on Development and Learning and Epigenetic Robotics (ICDL)* (pp. 1–7).
- Lamme, V., & Roelfsema, P. (2000). The distinct modes of vision offered by feedforward and recurrent processing. *Trends in Neurosciences*, 23(11), 571–579.
- Latash, M. L., Turvey, M. T., & Bernshtein, N. A. (1996). *Dexterity and its development*. Lawrence Erlbaum.
- LeCun, Y., Bottou, L., Orr, G. B., & Müller, K. (1998). Efficient backprop. In *Neural Networks: Tricks of the Trade* (pp. 9–50). Springer.
- Lee, S.-H., Blake, R., & Heeger, D. J. (2004). Traveling waves of activity in primary visual cortex during binocular rivalry. *Nature Neuroscience*, 8(1), 22–23.

- Lee, T., & Mumford, D. (2003). Hierarchical bayesian inference in the visual cortex. *Journal of Optical Society of America A*, *20*(7), 1434–1448.
- Li, W., Piëch, V., & Gilbert, C. (2004). Perceptual learning and top-down influences in primary visual cortex. *Nature Neuroscience*, *7*(6), 651–657.
- Liberman, A. M., & Mattingly, I. G. (1985). The motor theory of speech perception revised. *Cognition*, *21*(1), 1–36.
- Liu, J., Perez-Gonzalez, D., Rees, A., Erwin, H., & Wermter, S. (2009). Multiple sound source localisation in reverberant environments inspired by the auditory midbrain. In *Artificial Neural Networks–ICANN 2009* (pp. 208–217). Springer.
- Livingstone, M., & Hubel, D. (1988). Segregation of form, color, movement, and depth: Anatomy, physiology, and perception. *Science*, *240*(4853), 740.
- Lovchik, C., & Diftler, M. A. (1999). The robonaut hand: A dexterous robot hand for space. In *Proceedings. 1999 IEEE International Conference on Robotics and Automation, 1999.* (Vol. 2, pp. 907–912).
- Lowe, D. G. (1999). Object recognition from local scale-invariant features. In *The Proceedings of the Seventh IEEE International Conference on Computer Vision, 1999.* (Vol. 2, pp. 1150–1157).
- Luna, B., Thulborn, K., Strojwas, M., McCurtain, B., Berman, R., Genovese, C., & Sweeney, J. (1998). Dorsal cortical regions subserving visually guided saccades in humans: an fmri study. *Cerebral Cortex*, *8*(1), 40–47.
- Lund, J. S., Lund, R. D., Hendrickson, A. E., Bunt, A. H., & Fuchs, A. F. (1975). The origin of efferent pathways from the primary visual cortex, area 17, of the macaque monkey as shown by retrograde transport of horseradish peroxidase. *Journal of Comparative Neurology*, *164*(3), 287–303.
- Luppino, G., Murata, A., Govoni, P., & Matelli, M. (1999). Largely segregated parietofrontal connections linking rostral intraparietal cortex (areas AIP and VIP) and the ventral premotor cortex (areas F5 and F4). *Experimental Brain Research*, *128*(1-2), 181–187.
- MacKay, D. (1956). *The epistemological problem for automata.* Automata Studies, Princeton, NJ: Princeton University Press.
- MacKay, D. J. (1996). Bayesian methods for backpropagation networks. In *Models of Neural Networks III* (pp. 211–254). Springer.



- Mandler, J. M. (1992). The foundations of conceptual thought in infancy. *Cognitive Development*, 7(3), 273–285.
- Mandler, J. M. (1999). Preverbal representation and language. *Language and Space*, 365.
- Mareschal, D., & Johnson, M. H. (2003). The what and where of object representations in infancy. *Cognition*, 88(3), 259–276.
- Mareschal, D., Plunkett, K., & Harris, P. (1999). A computational and neuropsychological account of object-oriented behaviours in infancy. *Developmental Science*, 2(3), 306–317.
- McClelland, J. L., & Rumelhart, D. E. (1981). An interactive activation model of context effects in letter perception: I. An account of basic findings. *Psychological Review*, 88(5), 375.
- McConkie, G. W., & Currie, C. B. (1996). Visual stability across saccades while viewing complex pictures. *Journal of Experimental Psychology: Human Perception and Performance*, 22(3), 563.
- McMains, S., & Kastner, S. (2011). Interactions of top-down and bottom-up mechanisms in human visual cortex. *The Journal of Neuroscience*, 31(2), 587–597.
- Memisevic, R., & Hinton, G. (2007). Unsupervised learning of image transformations. In *2007 IEEE Conference on Computer Vision and Pattern Recognition* (pp. 1–8).
- Michel, O. (2004). Webots: Professional mobile robot simulation. *Journal of Advanced Robotics Systems*, 1(1), 39–42.
- Milner, A. D., Goodale, M. A., & Vingrys, A. J. (2006). *The visual brain in action* (Vol. 2). Oxford University Press Oxford.
- Mishkin, M., Ungerleider, L., & Macko, K. (1983). Object vision and spatial vision: Two cortical pathways. *Trends in Neurosciences*, 6, 414–417.
- Möller, R. (2012). A model of ant navigation based on visual prediction. *Journal of Theoretical Biology*, 305, 118–130.
- Mountcastle, V. B. (1957). Modality and topographic properties of single neurons of cat's somatic sensory cortex. *Journal of Neurophysiology*, 20(4), 408–434.

- Mountcastle, V. B. (1997). The columnar organization of the neocortex. *Brain*, *120*(4), 701–722.
- Muckli, L., Kohler, A., Kriegeskorte, N., & Singer, W. (2005). Primary visual cortex activity along the apparent-motion trace reflects illusory perception. *PLoS Biology*, *3*(8), e265.
- Natale, L., Nori, F., Sandini, G., & Metta, G. (2007). Learning precise 3D reaching in a humanoid robot. In *IEEE 6th International Conference on Development and Learning, ICDL* (pp. 324–329).
- Navarro-Guerrero, N., Weber, C., Schroeter, P., & Wermter, S. (2012). Real-world reinforcement learning for autonomous humanoid robot docking. *Robotics and Autonomous Systems*, *60*(11), 1400–1407.
- Neisser, U. (1967). *Cognitive psychology*. Appleton-Century-Crofts.
- Newman, C., Atkinson, J., & Braddick, O. (2001). The development of reaching and looking preferences in infants to objects of different sizes. *Developmental Psychology*, *37*(4), 561.
- Newman, M. E. (2004). Detecting community structure in networks. *The European Physical Journal B-Condensed Matter and Complex Systems*, *38*(2), 321–330.
- Newman, M. E. (2006). Modularity and community structure in networks. *Proceedings of the National Academy of Sciences*, *103*(23), 8577–8582.
- Ng, G. W. (2009). *Brain-mind machinery: brain-inspired computing and mind opening*. World Scientific.
- Nijhawan, R. (1994). Motion extrapolation in catching. *Nature*.
- Norouzi, M., Ranjbar, M., & Mori, G. (2009). Stacks of convolutional restricted Boltzmann machines for shift-invariant feature learning. In *IEEE Conference on Computer Vision and Pattern Recognition, 2009*. (pp. 2735–2742).
- Olshausen, B., Anderson, C., & Van Essen, D. (1993). A neurobiological model of visual attention and invariant pattern recognition based on dynamic routing of information. *The Journal of Neuroscience*, *13*(11), 4700–4719.
- Orabona, F., Metta, G., & Sandini, G. (2005). Object-based visual attention: a model for a behaving robot. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition-Workshops, 2005. CVPR Workshops*. (pp. 89–89).

- O'Regan, J. K., & Noë, A. (2001). A sensorimotor account of vision and visual consciousness. *Behavioral and Brain Sciences*, *24*(5), 939–972.
- Ozaki, T. J. (2011). Frontal-to-parietal top-down causal streams along the dorsal attention network exclusively mediate voluntary orienting of attention. *PLoS One*, *6*(5), e20079.
- Oztop, E., & Arbib, M. (2002). Schema design and implementation of the grasp-related mirror neuron system. *Biological Cybernetics*, *87*(2), 116–140.
- Oztop, E., Kawato, M., & Arbib, M. (2006). Mirror neurons and imitation: A computationally guided review. *Neural Networks*, *19*(3), 254–271.
- Pece, A. (1992). Redundancy reduction of a Gabor representation: A possible computational role for feedback from primary visual cortex to lateral geniculate nucleus. *Artificial Neural Networks*, *2*, 865–868.
- Pilly, P. K., & Grossberg, S. (2012). How do spatial learning and memory occur in the brain? Coordinated learning of entorhinal grid cells and hippocampal place cells. *Journal of Cognitive Neuroscience*, *24*(5), 1031–1054.
- Polley, D. B., Steinberg, E. E., & Merzenich, M. M. (2006). Perceptual learning directs auditory cortical map reorganization through top-down influences. *The Journal of Neuroscience*, *26*(18), 4970–4982.
- Pradhan, N., Burg, T., & Birchfield, S. (2011). Robot crowd navigation using predictive position fields in the potential function framework. In *IEEE American Control Conference, ACC* (pp. 4628–4633).
- Priebe, N., Lisberger, S., & Movshon, J. (2006). Tuning for spatiotemporal frequency and speed in directionally selective neurons of macaque striate cortex. *The Journal of Neuroscience*, *26*(11), 2941–2950.
- Prinz, W. (1984). Modes of linkage between perception and action. In *Cognition and Motor Processes* (pp. 185–193). Springer.
- Prinz, W. (1992). Why don't we perceive our brain states? *European Journal of Cognitive Psychology*, *4*(1), 1–20.
- Prinz, W. (1997). Perception and action planning. *European Journal of Cognitive Psychology*, *9*(2), 129–154.
- Prinz, W. (2003). Experimental approaches to action. *Agency and Self-awareness*, 165–187.

- Pulvermüller, F., Huss, M., Kherif, F., del Prado Martin, F. M., Hauk, O., & Shtyrov, Y. (2006). Motor cortex maps articulatory features of speech sounds. *Proceedings of the National Academy of Sciences*, *103*(20), 7865–7870.
- Raiguel, S. E., Lagae, L., Gulyàs, B., & Orban, G. A. (1989). Response latencies of visual cells in macaque areas v1, v2 and v5. *Brain Research*, *493*(1), 155–159.
- Rao, R. P., & Ballard, D. H. (1999). Predictive coding in the visual cortex: A functional interpretation of some extra-classical receptive-field effects. *Nature Neuroscience*, *2*(1), 79–87.
- Reynolds, J. H., & Chelazzi, L. (2004). Attentional modulation of visual processing. *Annual Reviews of Neuroscience*, *27*, 611–647.
- Rizzolatti, G., & Arbib, M. A. (1998). Language within our grasp. *Trends in Neurosciences*, *21*(5), 188–194.
- Rizzolatti, G., Camarda, R., Fogassi, L., Gentilucci, M., Luppino, G., & Matelli, M. (1988). Functional organization of inferior area 6 in the macaque monkey. *Experimental Brain Research*, *71*(3), 491–507.
- Rizzolatti, G., & Craighero, L. (2004). The mirror-neuron system. *Annual Review of Neuroscience*, *27*(1), 169–192.
- Rizzolatti, G., Fadiga, L., Gallese, V., & Fogassi, L. (1996). Premotor cortex and the recognition of motor actions. *Cognitive Brain Research*, *3*(2), 131–141.
- Rizzolatti, G., Fogassi, L., & Gallese, V. (2001). Neurophysiological mechanisms underlying the understanding and imitation of action. *Nature Reviews Neuroscience*, *2*(9), 661–670.
- Rocha, N., Silva, F., & Tudella, E. (2006). The impact of object size and rigidity on infant reaching. *Infant Behavior and Development*, *29*(2), 251–261.
- Rochat, P. (1987). Mouthing and grasping in neonates: Evidence for the early detection of what hard or soft substances afford for action. *Infant Behavior and Development*, *10*(4), 435–449.
- Rockland, K. S., & Van Hoesen, G. W. (1994). Direct temporal-occipital feedback connections to striate cortex (V1) in the macaque monkey. *Cerebral Cortex*, *4*(3), 300–313.
- Rolfs, M., Jonikaitis, D., Deubel, H., & Cavanagh, P. (2011). Predictive remapping of attention across eye movements. *Nature Neuroscience*, *14*(2), 252–256.

- Romberg, A., & Saffran, J. (2010). Statistical learning and language acquisition. *Wiley Interdisciplinary Reviews: Cognitive Science*, 1(6), 906–914.
- Rosch, E. (1973). Natural categories. *Cognitive Psychology*, 4(3), 328–350.
- Rosenbaum, D. A. (2009). *Human motor control*. Academic Press.
- Rosenblatt, F. (1958). The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6), 386.
- Ruff, H. A. (1984). Infants' manipulative exploration of objects: Effects of age and object characteristics. *Developmental Psychology*, 20(1), 9.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088), 533–536.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1988). Learning internal representations by error propagation. In *Neurocomputing: Foundations of research* (pp. 673–695). Cambridge, MA, USA: MIT Press.
- Rumelhart, D. E., & McClelland, J. L. (1986). *Parallel distributed processing: Explorations in the microstructure of cognition. volume 1. foundations*. MIT Press, Cambridge, MA, USA.
- Saegusa, R., Nori, F., Sandini, G., Metta, G., & Sakka, S. (2007). Sensory prediction for autonomous robots. In *7th IEEE-RAS International Conference on Humanoid Robots* (pp. 102–108).
- Salakhutdinov, R., & Hinton, G. (2009). Deep Boltzmann machines. *Artificial Intelligence*, 5(2).
- Saper, C. B., Iversen, S., & Frackowiak, R. (2000). Integration of sensory and motor function: The association areas of the cerebral cortex and the cognitive capabilities of the brain. *Principles of Neural Science*, 4, 349–380.
- Schaefer, A., Udluft, S., & Zimmermann, H. (2008). Learning long-term dependencies with recurrent neural networks. *Neurocomputing*, 71, 2481–2488.
- Schaeffer, J., & Plaat, A. (1997). Kasparov versus deep blue: The rematch. *ICCA Journal*, 20(2), 95–101.
- Schenk, T., & McIntosh, R. D. (2010). Do we have independent visual streams for perception and action? *Cognitive Neuroscience*, 1(1), 52–62.

- Schmoleky, M. T., Wang, Y., Hanes, D. P., Thompson, K. G., Leutgeb, S., Schall, J. D., & Leventhal, A. G. (1998). Signal timing across the macaque visual system. *Journal of Neurophysiology*, *79*(6), 3272–3278.
- Schultz, W. (1998). Predictive reward signal of dopamine neurons. *Journal of Neurophysiology*, *80*(1), 1–27.
- Schulz, D., Burgard, W., Fox, D., & Cremers, A. B. (2001). Tracking multiple moving targets with a mobile robot using particle filters and statistical data association. In *Proceedings of IEEE International Conference on Robotics and Automation, 2001*. (Vol. 2, pp. 1665–1670).
- Schyns, P. G., & Oliva, A. (1999). Dr. Angry and Mr. Smile: When categorization flexibly modifies the perception of faces in rapid visual presentations. *Cognition*, *69*(3), 243–265.
- Seltzer, B., & Pandya, D. N. (1994). Parietal, temporal, and occipital projections to cortex of the superior temporal sulcus in the rhesus monkey: A retrograde tracer study. *Journal of Comparative Neurology*, *343*(3), 445–463.
- Sergent, J., Signoret, J.-L., Bruce, V., & Rolls, E. (1992). Functional and anatomical decomposition of face processing: Evidence from prosopagnosia and PET study of normal subjects [and discussion]. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, *335*(1273), 55–62.
- Shadmehr, R., Smith, M., & Krakauer, J. (2010). Error correction, sensory prediction, and adaptation in motor control. *Annual Review of Neuroscience*, *33*, 89–108.
- Shipp, S., & Zeki, S. (1989). The organization of connections between areas v5 and v2 in macaque monkey visual cortex. *European Journal of Neuroscience*, *1*(4), 333–354.
- Shulman, G. L., Corbetta, M., Buckner, R. L., Raichle, M. E., Fiez, J. A., Miezin, F. M., & Petersen, S. E. (1997). Top-down modulation of early sensory cortex. *Cerebral Cortex*, *7*(3), 193–206.
- Sirosh, J., & Miikkulainen, R. (1997). Topographic receptive fields and patterned lateral interaction in a self-organizing model of the primary visual cortex. *Neural Computation*, *9*(3), 577–594.

- Sloutsky, V. M., & Robinson, C. W. (2008). The role of words and sounds in infants' visual processing: From overshadowing to attentional tuning. *Cognitive Science*, *32*(2), 342–365.
- Smith, L., Jones, S., Landau, B., Gershkoff-Stowe, L., & Samuelson, L. (2002). Object name learning provides on-the-job training for attention. *Psychological Science*, *13*(1), 13–19.
- Sommer, M. A., & Wurtz, R. H. (2006). Influence of the thalamus on spatial visual processing in frontal cortex. *Nature*, *444*(7117), 374–377.
- Spratling, M. W. (2008). Reconciling predictive coding and biased competition models of cortical function. *Frontiers in Computational Neuroscience*, *2*, 1-8.
- Streri, A., Pownall, T., & Kinglerlee, S. (1993). *Seeing, reaching, touching: The relations between vision and touch in infancy*. Harvester Wheatsheaf Oxford.
- Studdert-Kennedy, M. (2002). Mirror neurons, vocal imitation, and the evolution of particulate speech. *Mirror Neurons and the Evolution of Brain and Language*, 207–227.
- Summerfield, C., Trittschuh, E. H., Monti, J. M., Mesulam, M.-M., & Egner, T. (2008). Neural repetition suppression reflects fulfilled perceptual expectations. *Nature Neuroscience*, *11*(9), 1004–1006.
- Tani, J., & Ito, M. (2003). Self-organization of behavioral primitives as multiple attractor dynamics: A robot experiment. *IEEE Transactions on Systems, Man and Cybernetics, Part A: Systems and Humans*, *33*(4), 481–488.
- Tani, J., Ito, M., & Sugita, Y. (2004). Self-organization of distributedly represented multiple behavior schemata in a mirror system: reviews of robot experiments using RNNPB. *Neural Networks*, *17*(8-9), 1273–1289.
- Tardif, T. (1996). Nouns are not always learned before verbs: Evidence from mandarin speakers' early vocabularies. *Developmental Psychology*, *32*(3), 492.
- Thrun, S. (1998). Bayesian landmark learning for mobile robot localization. *Machine Learning*, *33*(1), 41–76.
- Thrun, S. (2002). Particle filters in robotics. In *Proceedings of the Eighteenth Conference on Uncertainty in Artificial Intelligence* (pp. 511–518).
- Todorov, E., & Jordan, M. I. (2002). Optimal feedback control as a theory of motor coordination. *Nature Neuroscience*, *5*(11), 1226–1235.

- Tomasello, M., & Farrar, M. J. (1986). Object permanence and relational words: A lexical training study. *Journal of Child Language*, *13*(03), 495–505.
- Tommerdahl, M., Favorov, O., Whitsel, B., Nakhle, B., & Gonchar, Y. (1993). Minicolumnar activation patterns in cat and monkey SI cortex. *Cerebral Cortex*, *3*(5), 399–411.
- Tong, F., Meng, M., & Blake, R. (2006). Neural bases of binocular rivalry. *Trends in Cognitive Sciences*, *10*(11), 502–511.
- Tootell, R. B., Hadjikhani, N. K., Vanduffel, W., Liu, A. K., Mendola, J. D., Sereno, M. I., & Dale, A. M. (1998). Functional analysis of primary visual cortex (V1) in humans. *Proceedings of the National Academy of Sciences*, *95*(3), 811–817.
- Ungerleider, L. G., & Desimone, R. (1986a). Cortical connections of visual area MT in the macaque. *Journal of Comparative Neurology*, *248*(2), 190–222.
- Ungerleider, L. G., & Desimone, R. (1986b). Projections to the superior temporal sulcus from the central and peripheral field representations of V1 and V2. *Journal of Comparative Neurology*, *248*(2), 147–163.
- Ungerleider, L. G., & Pessoa, L. (2008). What and where pathways. *Scholarpedia*, *3*(11), 5342.
- Van Essen, D. C., Anderson, C. H., Felleman, D. J., et al. (1992). Information processing in the primate visual system: an integrated systems perspective. *Science*, *255*(5043), 419–423.
- Van Essen, D. C., & Maunsell, J. H. (1983). Hierarchical organization and functional streams in the visual cortex. *Trends in Neurosciences*, *6*, 370–375.
- van Hasselt, H., & Wiering, M. (2007). Reinforcement learning in continuous action spaces. In *IEEE International Symposium on Approximate Dynamic Programming and Reinforcement Learning, 2007*. (pp. 272–279).
- Vijayakumar, S., Conradt, J., Shibata, T., & Schaal, S. (2001). Overt visual attention for a humanoid robot. In *Proceedings. 2001 IEEE/RSJ International Conference on Intelligent Robots and Systems, 2001*. (Vol. 4, pp. 2332–2337).
- von Helmholtz, H., et al. (1909). *Handbuch der Physiologischen Optik*. Hamburg: Voss.



- von Hofsten, C. (1982). Eye–hand coordination in the newborn. *Developmental Psychology*, *18*(3), 450.
- von Stutterheim, C., & Nuse, R. (2003). Processes of conceptualization in language production: language-specific perspectives and event construal. *Linguistics*, *41*(5; ISSU 387), 851–882.
- Ward, L. M. (2008). Attention. *Scholarpedia*, *3*(10), 1538.
- Webb, B. (1995). Using robots to model animals: a cricket test. *Robotics and Autonomous Systems*, *16*(2), 117–134.
- Webb, B. (2002). Robots in invertebrate neuroscience. *Nature*, *417*(6886), 359–363.
- Weber, C., & Triesch, J. (2008). A sparse generative model of V1 simple cells with intrinsic plasticity. *Neural Computation*, *20*(5), 1261–1284.
- Weng, J., McClelland, J., Pentland, A., Sporns, O., Stockman, I., Sur, M., & Thelen, E. (2001). Autonomous mental development by robots and animals. *Science*, *291*(5504), 599–600.
- Wermter, S., Palm, G., & Elshaw, M. (Eds.). (2005). *Biomimetic neural learning for intelligent robots*. Springer.
- Wermter, S., et al. (Eds.). (2014). *Artificial neural networks and machine learning–ICANN 2014*. Springer Heidelberg.
- Wilson, H. R., & Cowan, J. D. (1972). Excitatory and inhibitory interactions in localized populations of model neurons. *Biophysical Journal*, *12*(1), 1–24.
- Wilson, R. A., & Foglia, L. (2011). Embodied cognition. In E. N. Zalta (Ed.), *The Stanford Encyclopedia of Philosophy* (Fall 2011 ed.).
- Wiskott, L., & Sejnowski, T. (2002). Slow feature analysis: Unsupervised learning of invariances. *Neural Computation*, *14*(4), 715–770.
- Wolpert, D., Ghahramani, Z., & Jordan, M. I. (1995). An internal model for sensorimotor integration. *Science*, 1880–1880.
- Wolpert, D., & Kawato, M. (1998). Multiple paired forward and inverse models for motor control. *Neural Networks*, *11*(7-8), 1317–1329.
- Woodward, A. L. (1998). Infants selectively encode the goal object of an actor’s reach. *Cognition*, *69*(1), 1–34.

- Wurtz, R. H. (2008). Neuronal mechanisms of visual stability. *Vision Research*, 48(20), 2070–2089.
- Yamashita, Y., & Tani, J. (2008). Emergence of functional hierarchy in a multiple timescale neural network model: a humanoid robot experiment. *PLoS Computational Biology*, 4(11), e1000220.
- Yan, W., Weber, C., & Wermter, S. (2011). A hybrid probabilistic neural model for person tracking based on a ceiling-mounted camera. *Journal of Ambient Intelligence and Smart Environments*, 3(3), 237–252.
- Yilmaz, O. (2012). Oscillatory synchronization model of attention to moving objects. *Neural Networks*, 29, 20–36.
- Yu, C. (2008). A statistical associative account of vocabulary growth in early word learning. *Language Learning and Development*, 4(1), 32–62.
- Zeki, S., & Bartels, A. (1998). The autonomy of the visual systems and the modularity of conscious vision. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 353(1377), 1911–1914.
- Zhong, J., & Fung, Y.-F. (2012). Case study and proofs of ant colony optimisation improved particle filter algorithm. *Control Theory & Applications, IET*, 6(5), 689–697.
- Zhong, J., Fung, Y.-F., & Dai, M. (2010). A biologically inspired improvement strategy for particle filter: Ant colony optimization assisted particle filter. *International Journal of Control, Automation and Systems*, 8(3), 519–526.

# Publications Arising from this Thesis

- Zhong, J., Weber, C., & Wermter, S. (2011). Robot trajectory prediction and recognition based on a computational mirror neurons model. In *Artificial Neural Networks and Machine Learning–ICANN 2011* (Vol. 2, pp. 333–340). Espoo, Finland: Springer.
- Zhong, J., Weber, C., & Wermter, S. (2012a). Learning features and predictive transformation encoding based on a horizontal product model. In *Artificial Neural Networks and Machine Learning–ICANN 2012* (pp. 539–546). Lausanne, Switzerland: Springer.
- Zhong, J., Weber, C., & Wermter, S. (2012b). Learning features and transformation encoding based on a generative horizontal product model. In *Proceedings of the Sixteenth International Conference on Cognitive and Neural Systems (IC-CNS 2012)*. Boston, MA, USA.
- Zhong, J., Weber, C., & Wermter, S. (2012c). A predictive network architecture for a robust and smooth robot docking behavior. *Paladyn. Journal of Behavioral Robotics*, 3(4), 172–180.
- Zhong, J., Cangelosi, A., & Wermter, S. (2014). Towards a self-organizing pre-symbolic neural model representing sensorimotor primitives. *Frontiers in Behavioral Neuroscience*, 8, 22.

# Erklärung über die Eigenständigkeit der Dissertation

Hiermit erkläre ich an Eides statt, dass ich die vorliegende Dissertationsschrift selbst verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe.

## Declaration of Authorship

I hereby declare, on oath, that I have written the present dissertation by my own and have not used other than the acknowledged resources and aids.

Hamburg, den

Datum/Date:

---

Unterschrift/Signed:

---

Name/Name:

---