



COMENIUS UNIVERSITY IN BRATISLAVA
FACULTY OF MATHEMATICS, PHYSICS AND INFORMATICS
DEPARTMENT OF APPLIED INFORMATICS

GROUNDING THE MEANING
IN SENSORIMOTOR COGNITION:
A CONNECTIONIST APPROACH

Dissertation thesis

RNDr. Kristína Rebrová

Study programme: Informatics
Field of study: 9.2.1 Informatics
Supervisor: doc. Ing. Igor Farkaš, PhD.

Bratislava, 2013



THESIS ASSIGNMENT

Name and Surname: RNDr. Kristína Rebrová
Study programme: Computer Science (Single degree study, Ph.D. III. deg., full time form)
Field of Study: 9.2.1. Computer Science, Informatics
Type of Thesis: Dissertation thesis
Language of Thesis: English

Title: Grounding the meaning in sensorimotor cognition: a connectionist approach
Aim: Design and evaluate a connectionist model for bidirectional mapping between sensory and motor representations, that could serve as a basis for the underlying mirror neuron system based action understanding.
Annotation: Computational models of motor resonance and mirror neuron system are a plausible approach towards learning bidirectional sensorimotor mappings that can help explain action understanding from the grounded cognition perspective. Especially appealing is the direct matching hypothesis that emphasizes the role of the close link between sensory and motor representations mediating the understanding of the observed behavior.
Keywords: sensory-motor representations, mirror neuron system, neural network, grounding, response (in)variance

Tutor: doc. Ing. Igor Farkaš, PhD.
Department: FMFI.KAI - Department of Applied Informatics
Head of department: doc. PhDr. Ján Rybár, PhD.
Assigned: 20.10.2010
Approved: 20.10.2010
prof. RNDr. Branislav Rován, PhD.
Guarantor of Study Programme

Student

Tutor



Univerzita Komenského v Bratislave
Fakulta matematiky, fyziky a informatiky

ZADANIE ZÁVEREČNEJ PRÁCE

- Meno a priezvisko študenta:** RNDr. Kristína Rebrová
Študijný program: informatika (Jednoodborové štúdium, doktorandské III. st., denná forma)
Študijný odbor: 9.2.1. informatika
Typ záverečnej práce: dizertačná
Jazyk záverečnej práce: anglický
- Názov:** Ukotvenie významov v sensorimotorickej kognícii: konekcionistický prístup
Cieľ: Navrhnete a vyhodnotíte konekcionistický model na obojsmerné prepojenie medzi senzoricými a motorickými reprezentáciami, ktorý bude slúžiť ako základ pre porozumenie akciám prostredníctvom systému zrkadliacich neurónov.
Anotácia: Výpočtové modely motorickej rezonancie a systému zrkadliacich neurónov sú prijateľným prístupom k učeniu obojsmerných senzomotorických prepojení, ktoré môžu vysvetliť porozumenie akciám z perspektívy ukotvenej kognície. Obzvlášť atraktívnou je hypotéza priameho mapovania, ktorá vyzdvihuje úlohu úzkeho prepojenia medzi senzoricými a motorickými reprezentáciami, sprostredkúceho porozumenie pozorovanému správaniu.
Kľúčové slová: senzo-motorické reprezentácie, systém zrkadliacich neurónov, neurónová sieť, ukotvenie, (in)variantnosť odozvy
- Školiteľ:** doc. Ing. Igor Farkaš, PhD.
Katedra: FMFI.KAI - Katedra aplikovanej informatiky
Vedúci katedry: doc. PhDr. Ján Rybár, PhD.
- Spôsob prístupnosti elektronickej verzie práce:** bez obmedzenia
- Dátum zadania:** 20.10.2010
- Dátum schválenia:** 20.10.2010
- prof. RNDr. Branislav Rován, PhD.
garant študijného programu

.....
študent

.....
školiteľ

Acknowledgment

I would like to thank my supervisor Igor Farkaš for the guidance and help he has provided me. I would also like to thank my undergraduate supervisor and colleague Martin Takáč for his guidance, knowledge and enthusiasm which inspired me to follow a scientific career. Regarding my research, I would like to thank to my younger colleague Matej Pecháč and also to my other students of cognitive science who inspired me during the whole course of my studies. Great many thanks belong to my parents and my brother for they boundless support and love. Last, but not least, I would like to thank my partner Ľudovít Malinovský for all his support and advice in both personal and professional domain.

Abstrakt

Teória ukotvenej kognície predpokladá, že koncepty sú v mozgu/mysli zakódované na základe senzorických a motorických reprezentácií, ktoré s nimi súvisia. V doméne porozumenia akciám bolo mnohokrát empiricky overené, že pozorovanie pohybu vyvoláva v mozgu aktivitu nielen vo vizuálnych oblastiach, ale aj v motorickej kôre. Tieto poznatky možno aplikovať v doméne kognitívnej robotiky pri tvorbe riadiacich architektúr pre humanoidné roboty. Prepojenie percepcie a akcie na vyššej úrovni umožní agentovi oddeliť vykonávanie akcie od jej pomenovania, čo je problém, ktorým trpí mnoho modelov v tejto problematike. Na úvod opisujem konceptuálnu a empirickú teóriu k problematike prepojenia percepcie a akcie, významné modely a trérovacie algoritmy umelých neurónových sietí a súčasný stav problematiky v oblasti výpočtového modelovania systému zrkadliacich neurónov a ukotvenia významu v senzomotorickej interakcii.

Jadro tejto dizertačnej práce tvorí konekcionistický model systému zrkadliacich neurónov zostavený z rôznych, vzájomne prepojených, umelých neurónových sietí navrhnutý pre kognitívneho robota iCub. Aktiváciu zrkadliacich neurónov je modelované pomocou obojsmernej asociácie vysokoúrovňových vizuálnych a motorických reprezentácií. Pre účely tohto modelu sme spolu s mojimi vedúcimi vytvorili nový obojsmerný učiaci algoritmus pre viacvrstvové neurónové siete ako alternatívu k biologicky neplauzibilnému spätnému šíreniu chyby.

Výsledky z experimentov s neurónovými sieťami tvoriacimi modulárnu architektúru navrhovaného modelu ukázali, že v modeli úspešne vznikajú koherentné vysokoúrovňové reprezentácie a tiež, že v ňom vzniká želaná zrkadliaca aktivita. Dôležitým prínosom tejto práce sú aj experimenty s novým učiacim algoritmom BAL, ktoré naznačujú, že tento algoritmus je schopný skonvergovať k dobrým riešeniam a úspešne plní funkciu obojsmerného mapovania arbitrárnych ako aj reálnych (napr. robotických) vstupných vzorov.

Kľúčové slová: výpočtové modelovanie, umelé neurónové siete, senzomotorická kognícia, ukotvenie významu, porozumenie akciám, zrkadliace neuróny, kognitívna robotika

Abstract

According to grounded theories of cognition, concepts in the mind/brain are stored in terms of their sensory and motor consequences. In the domain of action understanding, recent empirical evidence suggests, that to understand the observed action the brain exploits areas that are involved in the production of the observed movement patterns. Such principle might be very useful also for a cognitive robot, allowing it to link the observed action with its own motor repertoire in order to understand the observed scene. Using this mechanism, the robot should be able to name actions without a necessity to produce the action while naming it (as many other models have to).

This thesis provides a conclusive overview of empirical and conceptual evidence on a common framework for action and perception in the brain, it describes related ANN architectures and learning algorithms, and summarizes a state-of-the-art in computational modeling of mirror neuron circuitry and grounded meaning acquisition.

In this thesis I present a connectionist model of the mirror neuron system built from various interconnected artificial neural networks for a cognitive robot iCub. In this model, the firing of mirror neurons is modeled through bidirectional association of high-level sensory and motor representations. To achieve this mapping I propose a new supervised learning algorithm for multi-layer bidirectional networks designed to overcome the biological implausibility of standard error back-propagation.

Results from experiments with neural networks comprising this model showed that the model is successfully able to form coherent high-level representations and form associations between them, based on which the model displays the desired mirroring property. Another outcome of this thesis are the results on the novel BAL algorithm, which suggest that it can converge to good solutions., i.e. learn fully bidirectional mappings between arbitrary and real-world input patterns.

Keywords: computational modeling, artificial neural network, sensorimotor cognition, meaning grounding, action understanding, mirror neurons, cognitive robotics

Contents

1	Introduction	1
2	Theoretical and empirical background	3
2.1	Action understanding: broad perspective	4
2.1.1	Ideomotor theory of action	4
2.1.2	Common coding theory	6
2.1.3	Neural correlates of action understanding	10
2.2	Mirror neurons	10
2.2.1	Mirror neurons in monkeys	10
2.2.2	Mirror neurons in humans	13
2.2.3	Mirror neurons and action understanding	15
2.2.4	Action understanding as a continuum	17
2.2.5	Mirror neurons and imitation	18
2.2.6	Mirror neurons and goal understanding	20
2.2.7	Other possible roles of mirror neurons in cognition	21
2.2.8	The origin of mirror neurons	22
2.2.9	Mirror neurons and high-level representations	26
2.3	Motor resonance in language comprehension	28
2.3.1	Neural evidence	28
2.3.2	Behavioral evidence	29
2.3.3	Language and mirror neurons	30
2.4	Summary	31
3	ANN architectures and learning	33
3.1	Error-driven learning	34
3.1.1	Multi-layer perceptron	34

3.1.2	Reinforcement learning and CACLA	36
3.1.3	Recurrent neural network	39
3.1.4	Recurrent neural network with parametric biases	40
3.1.5	Generalized Recirculation	44
3.1.6	Bidirectional Activation-based Learning	47
3.2	Unsupervised learning	49
3.2.1	Kohonen’s self-organising map	49
3.2.2	Merge self-organising map	51
4	Computational modeling	53
4.1	Cognitive robotics	54
4.1.1	The iCub robot	55
4.2	Mirror neuron system models	56
4.2.1	FARS	57
4.2.2	Mirror neuron system 1	59
4.2.3	Extending Mirror Neuron System	61
4.2.4	Mental State Inference	65
4.2.5	Hebbian Account	67
4.2.6	Higher-order Hopfield network	68
4.2.7	Tessitore’s model	69
4.2.8	Knott’s model	71
4.2.9	RNNPB model	73
4.2.10	SRNPB model	75
4.2.11	MOSAIC	76
4.2.12	Wiedermann’s finite cognitive agents	78
4.2.13	Bayesian approach	81
4.3	Grounding meaning in action	82
4.3.1	Direct grounding of language in action	82
4.3.2	RNNPB and language	85
4.3.3	Grounding in sensorimotor behavior using RL	89
4.3.4	The MirrorBot project	91
4.3.5	Linking language to motor chains	94
4.3.6	Connectionist model of symbol grounding transfer	96

4.3.7	Neural theory of language	99
4.3.8	TWIG	103
5	Towards robotic MNS	107
5.1	Architecture of the model	109
5.1.1	Executive and perceptual modules	110
5.1.2	Higher associative areas	111
5.1.3	PF pathway	112
5.1.4	AIP pathway	113
5.1.5	Model function and learning	115
5.2	Experiments with BAL	115
5.2.1	4-2-4 encoder	116
5.2.2	Simple binary vector association	118
5.2.3	Complex binary vector association	119
5.2.4	Conclusion	121
5.3	Experiments with robotic MNS	122
5.3.1	Level 2: self-organization of sensory and motor inputs .	122
5.3.2	Level 3: bidirectional sensorimotor association	128
5.3.3	Level 4: towards invariance	135
5.4	Discussion and future work	136
6	Conclusion	139
	Bibliography	140

List of Shortcuts

MNS	Mirror neurons
MNS	Mirror neuron system
EEG	Electroencepalography
fMRI	Functional magnetic resonance
STS	Superior temporal sulcus
ACE	Action–sentence compatibility effect
ANN	Artificial neural network
MLP	Multi-layer percetron
BP	Error back-propagation
CACLA	Continuous actor-critic learning-automaton
TD	Temporal difference
RNN	Recurrent neural network
SRN	Simple recurrent network
BPTT	Back propagation through time
RTRL	Real-time recurrent learning
GeneRec	Generalized recirculation
BAL	Bidirectional activation-based learning
RNNPB	Recurrent neural network with parametric biases
SOM	Self-organizing map
MSOM	Merge self-organizing map
CR	Cognitive robotics
DoF	Degree of freedom
ODE	Open dynamics engine

List of Figures

2.1	Macaque monkey's mirror neuron responses	11
2.2	MNS in the macaque brain	13
2.3	Experiment on mirror neurons in monkeys with two types of plies	21
2.4	Perspective variant neurons in STS of a macaque monkey	27
2.5	Illustration of three perspectives eliciting different mirror neuron responses	28
2.6	Brain activation resulting from listening to various action verbs	29
3.1	Schematic depiction of a two-layer perceptron.	34
3.2	Scheme of actor-critic learning paradigm.	38
3.3	Schematic depiction unfolding a generic RNN in time.	41
3.4	Schematic depiction of two generic RNNPB nets	44
3.5	Schematic depiction of two activation phases in GeneRec	45
3.6	Schematic depiction of the BAL model	48
3.7	Schematic depiction of a self-organizing map.	49
3.8	Scheme of merge self-organizing map.	51
4.1	ICub robot and its simulator	56
4.2	Schematic depiction of the most commonly modeled brain areas related to MNS.	58
4.3	Schematic depiction of FARS model	59
4.4	Schematic depiction of MNS1 model	60
4.5	Schematic depiction of MNS2 model	62
4.6	Schematic depiction of MSI model	66
4.7	Schematic depiction of Knott's model	72
4.8	Scheme of two interconnected RNNPB	74
4.9	Schematic depiction of the SRNPB model	76
4.10	MOSAIC model in action control and action observation mode	77

4.11	Schematic depiction of a finite cognitive agent	79
4.12	Schematic depiction of an extended finite cognitive agent	80
4.13	Illustration of two interconnected RNNPB	86
4.14	RNNPB architecture with a meta-level network	89
4.15	The modular architecture of grounded RL model	90
4.16	Schematic depiction of the modular hierarchical self-organizing memory .	93
4.17	Schematic depiction neuronal pools and connections between them . . .	95
5.1	The sketch of our robotic MSN model.	109
5.2	Examples of three grasp types from the observer’s perspective	111
5.3	4-2-4 encoder: performance as function of λ	117
5.4	Encoder 4-2-4: development of network convergence	118
5.5	Bidirectional associator: network performance as a function of λ and n_H	119
5.6	Bidirectional associator: development of network performance over time .	120
5.7	Random complex data for BAL experiments.	120
5.8	Bidirectional associator with complex data: network performance as a function of λ and n_H	121
5.9	Bidirectional associator with complex data: development of network per- formance	122
5.10	Bidirectional associator with complex data: pattern match visualization .	123
5.11	Contour plots of MSOM quantitative measures as a function of different α and β	124
5.12	Examples of the trained motor and visual maps.	125
5.13	Number of unique patterns after k -WTA binarization as a function of k . .	127
5.14	Disambiguated dataset from visual MSOM 16×16 and motor MSOM 12×12	128
5.15	Perfomance of BAL as a function of α	130
5.16	Perfomance of BAL as a function of n_H	131
5.17	BAL with first-perspective robotic data: network performance in time . .	131
5.18	BAL with first-perspective robotic data: pattern match.	132
5.19	BAL performance in two learning phases	133
5.20	BAL after two learning phases: pattern match.	134
5.21	Results from STSa module.	135

List of Tables

3.1	Equilibrium network variables in GeneRec model.	46
3.2	Activation phases and states in BAL model.	47
5.1	Optimal parameters for MSOM-based modules.	124
5.2	Optimal parameters for MSOM response binarization.	127
5.3	BAL performance in two learning phases (50 nets).	133

Chapter 1

Introduction

The main topic of this dissertation thesis is grounding of the meaning in sensorimotor interaction. In short, it has been theorized and empirically proven that when we categorize, reason or talk about actions – voluntary goal-driven sequences of movements – the brain exploits the areas that are responsible for the execution of these movements. A partial activation of the motor circuitry based on observation of the movement of another without producing any movement is called the motor resonance. Embodied and grounded theories of meaning, like perceptual symbol system (PSS) hypothesis (Barsalou, 1999), assume that concepts in the mind/brain are stored in terms of their sensory consequences and representations from other modalities as a sort of “hubs“ of multimodal information. Using such a mechanism, when we predict outcomes of an action or reason about our future plans, we exploit these concepts to “simulate” the result. Motor resonance, which also appears when we only imagine a movement, is assumed to support the existence of this simulation. Regarding understanding in the movement domain, a prominent topic are the mirror neurons. These are special cells discovered in monkeys that are active not only when fulfilling their regular function, i.e. execution of grasping movement, but also when the monkey observes such movements.

The primary goal of my dissertation is to propose, implement, and evaluate a mirror neuron system (MNS) model in the framework of cognitive robotics. The MNS model forms a control architecture for a simulated iCub

robot which learns to grasp objects, to remember and represent these grasping actions, and tries to understand (recognize) the observed grasping behavior. A robot endowed with a mirror neuron circuitry will, in line with the empirical evidence in this thesis, be able to separately produce or observe and name actions, without necessity to produce the action when naming it (as many models have to, more in Sec. 4).

In Chapter 2, I present an overview of conceptual and empirical basis on the interconnection of action and perception in human (and animal) cognition. I present the concepts and principles of ideomotor theory, common coding theory and motor resonance. Subsequently, I discuss a broad spectrum of topics regarding the mirror neurons, from their discovery, their emergence, up to their assumed role in high cognition. Finally, I discuss the evidence on involvement of the motor modality in understanding of the linguistic meaning.

Chapter 3 provides the background on selected artificial neural network (ANN) architectures and learning algorithms frequently mentioned in this thesis. Within this chapter I also introduce a new learning algorithm BAL (Farkaš and Rebrová, 2013).

Chapter 4 is devoted to computational models mainly based on ANN. After a brief introduction including a definition of the cognitive robotics framework central to this thesis, I provide overview and critical evaluation of selected models of MNS. Later on, I discuss also models of grounding of the (linguistic) meaning.

Finally, in Chapter 5, I describe the proposed robotic MNS model. The model consists of four logical parts (levels) which I describe and evaluate separately. This chapter also discusses experiments and results of various experiments with the model.

The closing chapter of this thesis is the conclusion with a overview on the whole thesis and its outcomes.

Chapter 2

Theoretical and empirical background

In this chapter I provide a broad empirical overview of the concept of connecting perception and action in the common framework providing a basis for grounding of the (linguistic) meaning. First, I introduce historical and conceptual overview of how action and perception can be linked to provide multi-modal grounded representations. I describe the concept of motor resonance – the partial activation of neural circuitry triggered by observation of a movement without the production, and its possible role in action understanding. I discuss this particular topic in more detail in (Rebrova, 2012).

Subsequently I describe the discovery, anatomy, functionality, and origin of mirror neurons that were theorized to serve as the neural basis for action understanding. I explain differences between visual and motor theories of understanding, emphasizing that the motor modality involvement might be very useful in assessing complex visual scenes, predicting the outcomes of an action, and understanding the goals of the observed agent. Finally, I summarize the evidence on motor-based understanding of linguistic meaning. This evidence suggest that the linguistic modality is most likely linked to perception and action.

2.1 Action understanding: broad perspective

2.1.1 Ideomotor theory of action

In cognitive science the term action control refers to a mechanism that causes people to express motor behavior, the action. There are two main theories of action control (Iacoboni, 2009). First, the so-called sensorimotor framework¹ states that actions emerge only as behavioral responses to external stimuli. Perception and action (motor control) are separate from each other. Stimuli are translated into motor commands by a special stimulus–response mapping. Ideas (if existent) do not cause actions. Regarding imitation and understanding of actions this approach suffers from the *correspondence problem*: a transformation of perception of other’s movement to the observer’s (imitator’s) movement is needed. This kind of transformation is from the computational perspective very demanding, and it is still not well known, whether brain employs detailed coordinate translations, or there is a different, more simple mechanism.

Ideomotor framework of action, on the other hand, avoids the correspondence problem easily, hence it primarily assumes that perception and actions (motor domain) share common neural codes. The ideomotor theory of action has its roots in the work of several German and British scholars (Stock and Stock, 2004). However, it became well known more than 50 years later, from James’ *Principles of Psychology* (1890). According to the ideomotor principle, human actions are initiated by sensory consequences that typically result from them, in other words, by the anticipation of their effects. Actions are represented on the basis of perceptual aspects that are usually present during their execution. The association between actions and perception indeed has to be learned through experience.

¹In the case of history of representing and explaining motor behavior the term sensorimotor framework refers to opposite meaning to its meaning in the usual context in this thesis. I encourage the reader to note this distinction carefully.

History of ideomotor theory

The ideomotor principle was formed by Herbart (1816, 1825) as a simple solution to the mind–body problem. Herbart considered the process of action initiation through ideas as a basic principle for intentional behavior. Followed by his successors and colleagues, his account gained a form very similar to the contemporary ideomotor theory. Lotze (1852) added more emphasis to the learning process in which *ideagenic* and *kinetic* substrates representing actions are formed. He also added that this process is inevitable. Harless (1861) formulated the theory of the *apparatus of will*, in which he assumed that “the organic bases of sensation have functional coherence with organic basis of movements”. His theoretical model can be considered the first model for common coding of perception and action. The second, independent root of the ideomotor theory, were English scholars. Namely it was Laycock (1845) who adopted medical-physiological perspective and Carpenter (1852) who coined the term *ideomotor*. However, his account on occult phenomena did not enter the reformulation of theory by James (1890).

In his *Principles* James integrated the two roots. He took the most suitable term from Carpenter (but only the term) and adopted the full theory from Lotze and Harless. He added that the ideomotor reaction happens “unhesitatingly and immediately”. He pointed out that everyday actions have reflex-like nature. There are many behavioral routines we execute subconsciously, for instance when we eat raisins out of a cake, fully engaged in a conversation, without noticing the final consequences of exposing a bad habit we have. On the behalf of the automaticity of an ideomotor reaction, James explained that the effect induced during sole observation of the movement might not necessarily lead to action. It can be inhibited by a “competing idea” that blocks the execution. It is important to note the conceptual contribution of this theory, formed nearly hundred years before the actual empirical evidence from neuroscience and neuropsychology emerged. James’ account was heavily criticized by Thorndike (1913), who was at that time the president of American Psychological Association (APA). Thorndike’s main argument against the ideomotor theory, from the position of the sensorimotor

framework, was that an idea can never initiate action, only sensory stimuli are responsible for any behavior. He considered ideomotor theory a kind of mystification and rejected it completely. It is possible that this act caused a diminished interest in this theory for another six decades.

Contemporary ideomotor theory

The ideomotor theory was revisited later by Greenwald (1970), who reformulated it in more empirically verifiable terms. He characterized three basic elements of ideomotor action: stimulus (S), response (R), and (sensory) effect (E), and proposed an experimental methodology to study the relationship between them. Such experiments form a core of contemporary ideomotor theory (Shin et al., 2010) as well as experiments on the *theory of event coding* (TEC) (Hommel et al., 2001). TEC assumes that perceived and to-be-produced events are represented in a common domain, in a distributed fashion, and in a hierarchical structure. There were various ideomotor phenomena observed, including stimulus–response compatibility (facilitation of reaction on the basis of congruence with the stimulus) and ideomotor action, an involuntary movement that tends to arise when observing another’s performance.

As suggested already by James (1890), ideomotor principle evokes an action automatically, but the movement might not be executed. From modern simulationist view (Barsalou, 1999; Jeannerod, 2001; Borghi et al., 2010) unexecuted ideomotor action represented in a common coding framework might as well serve as a mental simulation of the observed action resulting in the understanding of the observed behavior. In the next section I review various sources of empirical evidence on common coding of perception and action. Then I introduce mirror neurons as a possible substrate of common coding in the brain.

2.1.2 Common coding theory

The *common coding theory* (Prinz, 1997; Hommel et al., 2001) suggests that there is a common representational base for perception and action (motor performance). The perception of action automatically activates its motor

component and vice versa. The common coding framework might also be considered the means for sensorimotor simulation (Barsalou, 1999; Jeannerod, 2001; Wolpert et al., 2003). As shown by Ehrsson et al. (2003), the same neural mechanisms are involved in mental imagery of a motor act as in its execution. It is likely that the understanding of the observed action works on the basis of motor resonance as well. Van der Wel et al. (2013) summarize various sources of empirical evidence and several problems for the common coding theory. I will briefly review some of them.

Behavioral evidence

A powerful source of evidence is the motor resonance based understanding. The main idea is that *the more closely the observed action maps onto the observer's own motor repertoire, the more accurate will be the observer's prediction of the course and the result of the action.* One source of such experiments is the comparison between the observation of self and others. Knoblich and Flach (2001) first filmed the participants throwing darts. A week later they observed their and others' dart throws and predicted to which part of the dart board will the dart land. After a short phase for adaptation to an unusual situation of observing oneself, participants displayed significantly better results for predicting results of their actions compared to the actions of others.

Another class of experiments on motor resonance can be named "professionals versus novices". For instance, Repp and Knoblich (2007) employed a special perceptual illusion, the Triton paradox. According to it, there are perceptually bistable pairs of tones, the perception of the pitch going up or down is equally likely. Repp and Knoblich tested two groups of participants: professional pianists and a group of unskilled controls. The task of the participants was to listen to (bistable) tones and to press keys on keyboard simultaneously. Unlike the controls, the pianists always perceived the change in the tone accordingly to the movement on the keyboard in a way that the piano works. When they pressed a series of keys from left to right, the tone seemed to go up and vice versa.

The effect of anticipatory performance enhancement was well demonstrated on experiments with professional sportsmen. For instance, Aglioti et al. (2008) compared how professional basketball players, expert watchers (coaches, sport journalists), and novices predicted the success of throw shots presented in video clips. Indeed, expert watchers outperformed novices, but the performance of expert players was even better. On the other hand, both target group displayed increased motor-evoked potentials. As suggested by James, the motor resonance is an automatic and inevitable phenomenon, although the motor response might not match well the observed behavior. Interestingly, in their experiments with judging the faking behavior in basketball, Sebanz and Shiffrar (2009) discovered that experts outperformed novices only when they observed full videos, but not when they viewed static images.

Neural evidence

Another source of empirical evidence for the common coding and its possible implementation in motor resonance are results from neuropsychology, specifically from studies on the suppression of EEG mu-rhythm (approximately in the same range as alpha-frequency band). The mu-rhythm is characteristic of the motor rest and vanishes when the subject begins to move. Already in 1950s, Cohen-Seat et al. (1954) and Gastaut and Bert (1954) observed the desynchronization of the EEG mu-rhythm in subjects that did no move, but observed actions performed by others. This phenomenon was later revisited by Oberman and Ramachandran (2007), who extended the effect of observation of human demonstrators to similar inanimate effectors like a robotic arm. In addition, motor resonance is not limited to adult observers. Van Elk et al. (2008) discovered the desynchronization of mu- and beta-frequency bands in 14–16 month old infants observing other infants crawl and walk. Importantly, the desynchronization was stronger for crawling than for walking. These findings suggest that some resonance appears even for movements out of the observer’s motor repertoire, but since it is “mapped” only loosely, the resonance is weaker.

Motor resonance can be modulated by the effect of learning. Cross et al. (2006) reported changes in brain activation during action observation while experts in modern dance learned new dance patterns. The comparison of results of scanning sessions at the beginning and in the middle of 4 weeks practice showed an increased activity in inferior parietal lobule and ventral premotor areas during the observation of the newly acquired movement in comparison with observation of the same movement before training. Moreover, Cross et al. (2009) extended these findings to pure perceptual expertise. They concluded that training on the physical and on the perceptual basis produced similar changes in the brain activation. It is then possible, that motor resonance can to some extent be enhanced by visual stimuli.

Problem of agency

The crucial problem for the common coding theory is the sense of agency (Sato and Yasuda, 2005). The question is: If perceived and executed actions activate the same action representations, how do we differentiate who actually executes the action? Van der Wel et al. (2013) distinguish between two basic accounts on the agency in common coding theory. First there are sensorimotor accounts that propose existence of a forward model (Wolpert and Kawato, 1998) that generates predictions, and an inverse model that reversely activates actions that could possibly lead to the observed situation. The brain always generates multiple forward models. If there is forward model that matches the observed sensory consequences, the agent experiences himself as the source of the action. The sensorimotor account in principle evaluates, whether the motor command was executed or whether only motor resonance took place.

The second group are perceptual accounts based solely on the evaluation of perceptual information. The empirical evidence for perceptual involvement in motor control comes from experiments of Fournieret and Jeannerod (1998), who observed that people adjust actions to perceived visual feedback without explicit awareness and ignoring sensorimotor cues arising from movements. Additionally to perceptual cues, Flach et al. (2003) emphasize the timing

of the events. Georgieff and Jeannerod (1998) proposed that additionally to *where* and *what* systems a third, *who* system exists in posterior parietal cortex. It generates time-tags linked to common codes for perception and action, hence providing the agent with the clue whether they are the source of the perceived action or not. Although these two views can be considered separate, there is a good possibility that the brain employs mechanisms from both of them. As in various other cognitive capacities, attention may play a key role in determining agency as well (Knott, 2012).

2.1.3 Neural correlates of action understanding

Among the theories aiming to explain the nature of action understanding together with its neural correlates, two rival theories can be distinguished (Rizzolatti et al., 2001). According to *the visual hypothesis*, the observed action is assessed on the basis of sole visual processing and is mediated by visual areas of the brain such as the Superior Temporal Sulcus (STS), which is sensitive to a large class of biological movements. On the contrary, *the direct matching hypothesis* emphasizes involvement of motor modality, concretely the mapping of the observed action onto an action in one's own motor repertoire. This second theory shares the common ground with the evidence provided in this section. The property of matching the observation with the execution has been found in the so-called mirror neurons. In the following section I describe the discovery of the mirror neuron system, its anatomy and properties in monkeys and humans and its possible roles in action understanding, and also in other cognitive capacities.

2.2 Mirror neurons

2.2.1 Mirror neurons in monkeys

Mirror neurons (MN) were originally discovered in area F5 in the ventral premotor cortex of a macaque monkey by Pellegrino et al. (1992). This area of the monkey's brain is characteristic with neurons that become active dur-

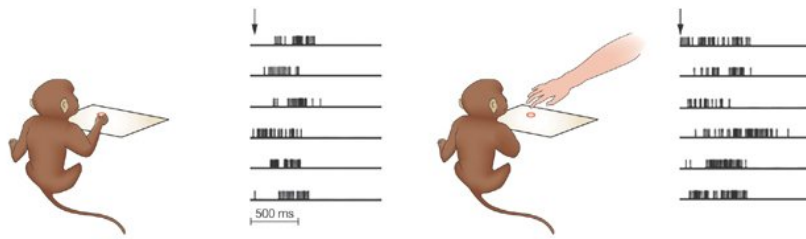


Figure 2.1: Illustration of macaque monkey’s mirror neuron responses (Rizzolatti et al., 2009).

ing particular hand and mouth movements, such as grasping, holding and tearing. Many of F5 neurons react only to very specific types of actions, for instance only to a precision grip, and some neurons in this area are activated also by visual stimuli like graspable food (Rizzolatti et al., 1988). Accidentally, during the preparation of an experiment, Pellegrino et al. (1992) discovered that some of the F5 neurons discharged not only during the execution of a certain motor act, but also when the monkey observed the particular motor act performed by the experimenter. This mirroring activity occurred only when the target of the motor act was present on the scene and when the motor act finished completely, for instance when the experimenter reached for a piece of food and grasped it successfully. Mirror neurons did not respond to meaningless actions or to the presentation of an object alone, even if it was food, which the monkey regards as the most interesting.

Subsequently, Gallese et al. (1996) and Rizzolatti et al. (1996) provided detailed reports on mirror neurons and their properties. Mirror neurons were also found to respond to auditory stimuli typical for the particular action, like nut-cracking (Köhler et al., 2002). Interestingly, Umiltà et al. (2001) realized that mirror neurons fire also when the object acted upon is not visible, but only if the monkey has sufficient clues to “figure out” that the object is hidden (e.g. behind a fence). Therefore mirror neuron activity can be considered strongly oriented towards the object and towards the goal of the action (see Sec. 2.2.6).

Mirror neuron system

On the basis of these findings, mirror neurons were theorized to be involved in action understanding, since they connect perception of an action with its representation in the observer's motor repertoire. Rizzolatti and Sinigaglia (2010) summarize that this function was attributed to the so-called *parieto-frontal action-observation action-execution brain circuit* also called the mirror neuron system (MNS). This circuit and other relevant areas are schematically depicted in Fig. 2.3. It consists of the above mentioned area F5, area PFG in rostral part of inferior parietal lobule (IPL) between areas PF and PG, and the anterior intraparietal area (AIP). The two parietal areas are both connected with F5 and both receive high-order visual information from areas located inside the superior temporal sulcus (STS) and the inferior temporal lobe (IT) providing input to frontal motor-control area F5. STS, similarly to F5, encodes biological motion, but it lacks motor properties and therefore cannot be considered a true part of the MNS, which applies for IT as well. The parieto-frontal circuit is also connected with the area F6 (pre-supplementary motor area) and the ventral prefrontal cortex, which are the higher-order areas that control it.

In addition to the parieto-frontal circuit, neurons with mirror properties were found in other areas of the parietal lobe, the lateral intraparietal area (LIP), which contributes to joint attention (Shepherd et al., 2009), and ventral intraparietal area (VIP). Neurons in VIP encode tactile and visual stimuli occurring in peripersonal space and might be responsible for encoding body-directed motor acts rather than object-directed motor acts represented by mirror neurons in F5 (Ishida et al., 2010). Mirror neurons were discovered also in frontal areas of the monkey brain, in primary motor (M1) and dorsal premotor (PMd) cortices (Tkach et al., 2007).

Canonical neurons

Together with mirror neurons, another interesting multi-modal type of neurons has been discovered. Unlike the MN, the canonical neurons respond only to presentation of a sole graspable object. The activity of the canon-

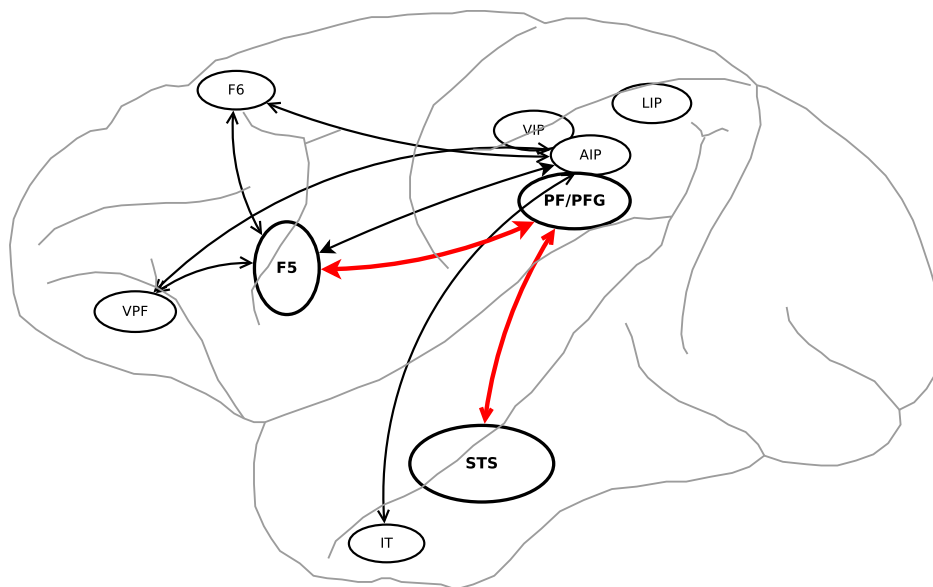


Figure 2.2: Schematic depiction of MNS in the macaque brain. The red arrows indicate brain connections of the highest interest.

ical neurons is bound to the agent's ability to apply the mirrored action to the object (e.g. when a neuron for power grip is firing in presence of an apple). These neurons were theorized to represent the object's affordances (Gibson, 1977), i.e. the set of actions that the agent can afford to apply to the object. According to Jeannerod et al. (1995), canonical neurons in AIP–F5 circuit mediate the transformation of object properties into appropriate hand postures. Canonical neurons (as well as the MNs) were also discovered in human brain (Grezes et al., 2003) on the basis of functional magnetic resonance (fMRI). Interestingly, a cooperation of mirror and canonical neurons has been found and the areas identified (area 44) roughly reflect the position of F5 in the macaque brain.

2.2.2 Mirror neurons in humans

Although the direct evidence for mirror neurons in humans emerged only recently, the background for assuming their existence is more than 50 years old. As summarized in Sec. 2.1.2, first empirical evidence for motor resonance comes from studies on the desynchronization of the EEG mu rhythm (Cohen-

Seat et al., 1954; Gastaut and Bert, 1954) and was recently confirmed and extended by Oberman and Ramachandran (2007). Rizzolatti and Craighero (2004) summarize a variety of EEG, MEG, and TMS studies that (indirectly) confirm the existence of mirror neurons in humans. According to Rizzolatti and Craighero, the core of the human mirror neuron system consists of the rostral part of the IPL, the lower part of the precentral gyrus, and the posterior part of the inferior frontal gyrus (IFG).

The first evidence of mirror neurons in humans on a single-cell level was provided by Mukamel et al. (2010). They recorded data from 21 patients with pharmacologically intractable epilepsy, who executed or observed hand grasping actions and facial emotional expressions. The observed regions of the brain were chosen according to clinical criteria only, so the main areas of interest, namely the Broca's area, which is a possible homologue of F5 (Rizzolatti and Arbib, 1998), were not examined. Significant proportions of cells responding to both action observation and action execution were found in the medial frontal lobe (SMA) and, interestingly, in the medial temporal lobe (namely hippocampus, the parahippocampal gyrus and the entorhinal cortex). According to Mukamel and colleagues, the mirroring activity, recorded during observation of an action in the medial temporal lobe, might correspond to the reactivation of the memory traces formed previously during its execution.

Unlike the mirror neurons in monkeys, which are triggered only by meaningful actions such as reaching for and grasping food, motor resonance in humans appears also based on intransitive meaningless arm movements (Fadiga et al., 1995). Another important difference between humans and monkeys is, that the postulated human mirror system, unlike that of the monkey, responds to actions regardless of the effector used to perform them, be it an animal, a human, or a robotic arm, or even a tool, with or without the presence of the target object (Peeters et al., 2009). The mirror system in monkey, on the other hand, requires an interaction between a biological effector (the hand or the mouth) and an object, and will not respond to an agent mimicking an action (Rizzolatti and Craighero, 2004). This suggests that the human

mirror system is more general and more abstract, and possibly involved in a larger variety of cognitive functions.

2.2.3 Mirror neurons and action understanding

Based on the execution-observation matching property, mirror neurons were theorized to be involved in understanding the actions of others. In line with the ideomotor theory of action and common coding theory, mirror neurons could be identified as a part of the multimodal representations of actions, since their activity can be triggered both by perceptual stimuli and motor activity. The role of mirror neurons in action understanding is still a subject of a vivid debate between neuroscientists and psychologists. Regarding motor involvement in action understanding two camps (similarly to two basis hypotheses) can be identified. Firstly, there are proponents of motor theories of understanding (direct-matching hypothesis). Secondly, there are those defending sole visual assessment (visual hypothesis) to whom we (Farkaš et al., 2011a) refer as opponents of motor-based (or mirror neuron based) understanding.

Regarding the opposition against motor understanding, a very careful assessment is in place. Unfortunately, opponents often restate the motor hypothesis as avoiding the visual component of understanding or declining a possibility of sole visual assessment. On the other hand, proponents of mirror neuron based theories do not state visual understanding impossible, only emphasize that the motor component might be useful in assessing the visual input. For instance, Tessitore et al. (2010) emphasize the bidirectional flow between motor and visual areas that might be very helpful in processing complex and broad visual information. Since this study is using the computational modeling methodology, I will describe it in more detail in the next section dedicated specially to computational models. According to wide-range analysis of Molenberghs et al. (2012), who applied clustering methods on 125 different fMRI studies on human mirror systems, the core structures activated during observation and comprehension of movement include not only mirror neuron areas, but also visual cortices, cerebellum and

parts of the limbic system. In line with the provided evidence, I believe that the interaction between all of these areas is necessary to describe the whole process of action understanding and the gradedness of this capacity (discussed in the next section).

One of the strongest opponents of motor-based understanding are Hickok (2008) and Hickok and Hauser (2010), who based their opposition on the claim that if mirror neurons mediate understanding, then the recognition of the observed actions and the motor ability to produce it should not dissociate. To falsify this claim, they use results by Mahon and Caramazza (2005), who reported various cases of patients with different motor impairments able to recognize and to name actions normally. Taking into account such findings, one can expect a sort of dichotomy between motor and non-motor understanding. On the other hand, in this study there is no indication whether the subjects were deprived of their motor abilities during their life or from birth. This might make a crucial distinction, since the motor capabilities lost during life might still leave the conceptual system of representation based on prior motor activity.

Another famous opponent of the involvement of MNs in understanding is Heyes (2010). Similarly to Hickok and Hauser, she states that the activation of mirror neurons is only an epiphenomenon emerging on the basis of associative learning. Interestingly, Heyes and her team (Catmur et al., 2007) showed, that mirror neurons in humans can be temporarily re-trained to respond to different action (i.e. one finger movement associated with completely different finger movement). However, unless we are able to do single-cell recording, we can only speak in terms of motor resonance, which is more general and less specific, since it is bound to non-invasive measuring techniques. It is well known that EEG has very good time resolution, but very poor spacial resolution. Thus, to claim that mirror neurons can be retrained only on the basis of such weak evidence might be considered an overstatement.

2.2.4 Action understanding as a continuum

An important question, which is often avoided is the exact definition of action understanding. Farkas et al. (2011a,b) adopt the definition of Gallese et al. (1996) who see action understanding as “the capacity to recognize that an individual is performing an action, to differentiate this action from other analogous to it, and to use this information in order to act appropriately”. This action has to be produced on the basis of some reasoning processes on the internal states of the agent and the assessment of this state through a behavioral response. Since an appropriate response should be beneficial to the agent, we can expect and evaluate it in advance.

Opponents of MN based action understanding often freely interchange the term “understanding” and “recognition”. We (Farkas et al., 2011a,a) do not regard them the same concepts, but as two different degrees of the understanding process. For us understanding is a continuous graded phenomenon, ranging from mere recognition – a simple categorical judgment, i.e. “this is grasping” – to the capacity to anticipate the internal state of the observed agent and to make predictions about the outcome of the observed action. This gradedness of understanding has been studied in context of the strength of motor resonance. As described in Sec. 2.1.2, neural and behavioral evidence shows that stronger activation in motor cortices strongly correlates with better assessment of the observed action. Interestingly, the difference between the performance of professional players and professional observers only appears when the stimulus material is a video, not when assessing static images (which do not trigger motor resonance).

In line with the action–continuum hypothesis, recognition might be mediated solely by categorical visual assessment, but motor resonance and action–perception matching still remains a crucial component of full (deep) understanding of the observed action and its consequences or goals (for goal-encoding see Sec. 2.2.6). One of Hickock’s arguments in his eight problems for mirror neurons (Hickok, 2008) is that “musically untrained people can recognize, say, saxophone playing even if they’ve never touched the instrument”. However, one can recognize that someone is playing the saxophone only be-

cause of recognizing the instrument itself. Regarding cognition and brain processes, everything runs in parallel. Therefore in this example we cannot clearly say that the observer recognized the saxophone play, because object recognition might have preceded action recognition. An interesting challenge for Hickock's theory would be to assess whether people can still recognize that someone is playing the saxophone in case of pantomime, so the brain could not rely on assessing the object itself, only the movement. In line with our theory, people who have no experiences with musical instruments whatsoever might display difficulties in this task.

This prediction should be testable using a behavioral experiment. At first, one group of participants will learn how to execute a simple, but novel movement, which is not part of their motor repertoire. During this time, the second group will observe videos of the same movement and receive a linguistic description of it. At the end, both groups will have to answer various questions regarding this movement, for example, whether it is possible to grasp an apple with it. Our prediction is that the second group will perform worse than the first group because mere observation of an action prevents the subjects from reaching the same level of motor proficiency, as in the case of subjects who could actively learn the task.

2.2.5 Mirror neurons and imitation

Rizzolatti and Arbib (1998) postulated that mirror neurons are involved in two crucial cognitive abilities – understanding of actions and their imitation. However, only few animal species are able to imitate (Iacoboni, 2009), so this claim is mostly restricted to humans and some higher animals. Heyes (2001) refers to imitation as to “copying by an observer of a feature of the body movement of a model” implying a causal link between observation of the movement and its execution. The imitation of the particular movement feature must happen after observation of this (and not other) feature, and it must not occur by chance. A similar position is taken by Rizzolatti et al. (2001) who emphasize the “response facilitation” – an automatic tendency to reproduce the observed movement, which may occur with or without under-

standing. They postulate that imitation with understanding is present only in adult humans and that it is possibly mediated by mirror neurons.

A more detailed description and, more importantly, different distinctions were introduced by Hurley (2008). She distinguishes four types of social learning: stimulus enhancement, goal emulation², movement priming and true imitation, while the latter is only present in humans. True imitation in this sense does not only require a proper copying of movements (means), but also a successful copying of goals (ends). Hurley explains how the ability to imitate, but also deliberation and mind reading, can be enabled by subpersonal (functional, not conscious, nor neural) mechanisms of control, mirroring, and simulation in her multi-layer Shared Circuits Model (SCM). In this model, consistently with the common coding theory and the direct-matching hypothesis, perception and action share the same neural resources. Hurley's aim is to reconcile ideomotor and associative theories of action understanding that are often posed against each other, emphasizing both the sharing of resources and associative learning.

An interesting question arises considering the development of imitative abilities in humans. Both infants and their parents are known to imitate facial gestures due to overt (contagious) imitation processes. As proposed by Heyes (2010), the imitative behavior of parents might as well serve as a natural mirror for infants to associate their gestures with their visual representations, providing then the motor resonance underlying action understanding. This can be considered a "imitation-first" view. On the other hand, "understanding-first" view emphasizes that understanding is needed in order to imitate. In Hurley's view, which I find quite plausible, these two skills develop hand in hand. At first, even if the imitation is not successful, for instance when the movement is somehow copied but its goal is not achieved, neural representation of the new movement starts to form. Following mul-

²The term "goal emulation" refers to a situation when an animal (or a human) sees an action and then produces a slightly different action to produce the same goal. It is known that children imitate both means and goals "blindly", even when the means are ineffective. On the other hand, it was found that children, given a context, will copy the goal, but using a simpler movement to achieve it, than the one executed by the experimenter (Bekkering et al., 2000).

multiple attempts, this representation gains its strength as well as the motor resonance evoked during the re-observation needed for further attempts to imitate, since a similar motor plan is already in the repertoire. In this sense, imitation and action understanding (mediated by mirroring processes) might develop together in a mutually beneficial way.

2.2.6 Mirror neurons and goal understanding

Another important property of the MNS is that “mirror neurons may encode the goal of the motor acts of another individual in an observer-centered spatial framework, thus providing the observer with crucial information for organizing their own future behavior in cooperation or competition with the observed individuals” (Rizzolatti and Sinigaglia, 2010). The intuition behind this theory is the existence of two types of mirror neurons according to their congruence with the observed action. The *strictly congruent* mirror neurons react only to certain type of motor act, for instance only to a precision grip. The *broadly congruent* mirror neurons, on the other hand, may react to a whole category of motor acts sharing the same goal (Gallese et al., 1996).

A compelling evidence for the goal-encoding role in monkeys was provided by Umiltà et al. (2008) who recorded single-cell activity in monkeys using and observing the usage of two types of pliers (normal and reverse pliers with opposite mechanisms for opening and closing). Results of this study showed the same pattern of some mirror neurons’ activity during the observation of both types of grasp. An example from human subjects supporting the goal-encoding hypothesis was provided in an interesting fMRI study. Gazzola et al. (2007) presented aplasic individuals (born without arms) with actions performed by hands, feet, and mouth and measured the motor resonance. The results of this experiment showed that the aplasic individuals’ motor cortices responded also to hand actions, which were obviously impossible for them to execute. Interestingly, these were only those actions, which the subjects were able to accomplish by mouth or feet. These findings suggest that there exist broadly congruent goal-encoding neurons, which connect not

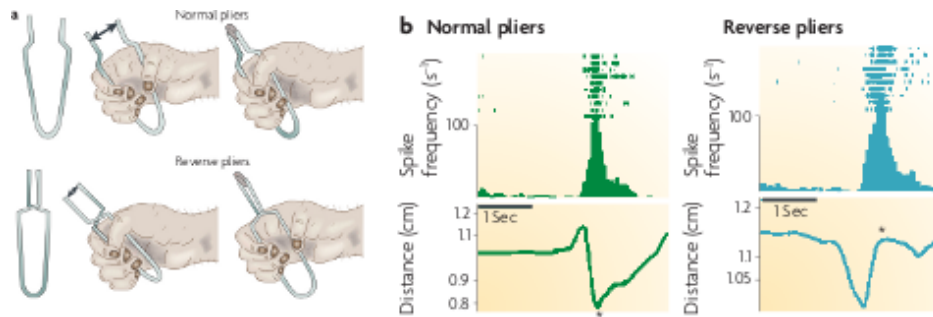


Figure 2.3: Schematic depiction of an experiment by Umiltà et al. (2008). Figures on the left illustrate two types of pliers and on the right the firing of one broadly congruent mirror neuron.

only various types of an action (various grasp types), but may also bind certain actions according to their goal.

2.2.7 Other possible roles of mirror neurons in cognition

Similarly to motor mirror neurons possibly involved in understanding of actions, mirroring mechanisms have also been considered to be involved in understanding of emotions. Gallese et al. (2004) propose these as a basic functional mechanism that provides an insight into other minds. As recently discussed by Rizzolatti and Sinigaglia (2010), a possible dichotomy between resonance-based and cognitive interpretation of the visual stimuli may exist. It is also quite likely that the latter type of processing is of a quick and superficial nature. On the other hand, “deep” understanding and empathy might require a mirroring mechanism to endow the individual with the ability to “relive” the present state of the observed agent.

These and various other findings about mirror neurons, or more precisely about the mirroring mechanism, suggest that its function is not limited to action understanding. Most important is the matching property, which provides the mapping from perception to a brain area that mediates the concrete function. The primary function of mirroring mechanism and perceptually triggered resonance in the particular effector brain area can be hypothesized

to mediate the simulation-based understanding and reasoning in line with simulationist theories of cognition like Barsalou’s PPS (Barsalou, 1999).

Lastly, one of the most intriguing conclusions drawn from the existence of mirror neurons is their possible role as a “missing link” between animal communication and human language (Arbib, 2005). The theory and the evidence in favor of it was provided already by Rizzolatti et al. (1996), who state that area F5 and Broca’s area might not only be anatomical homologues,³ but could also share functional properties crucial for development, production and understanding of communication gestures, which gave rise to the evolution of language. Similarly, Rizzolatti and Arbib (1998) state that the way to the evolution of the open vocalization system present in humans (speech) was paved by the evolution of the manual gestural system, facilitated by the action-execution action-observation matching property of neurons in Broca’s area. Interestingly, there has been no such evolution in monkeys. A great challenge in this field remains to explain the evolutionary changes of mirror neuron system and the related circuitry in humans leading to emergence of language.

2.2.8 The origin of mirror neurons

Regarding the origin of mirror neurons, there are several accounts, which tend to oppose one another. Firstly, there is another criticism of the approach of Rizzolatti’s and Arbib’s groups and of the evolutionary favoring of the MNS assumed by Rizzolatti and Arbib (1998). Heyes (2010) points out that according to the *adaptation hypothesis* (i.e. the evolutionary account of Rizzolatti and Arbib, 1998), experience plays a relatively minor role in the development of mirror neurons (it triggers or facilitates it), and that the capacity of mirror neurons to match observed with executed actions is genetically inherited. She also proposes an alternative, opposing account – the *association hypothesis*, which states that mirror neurons are merely a

³However, the relationship may be more complicated, as suggested by Grezes et al. (2003).

byproduct of associative learning and that their existence is not caused by any evolutionary mechanisms.

Association versus evolutionary account

According to the association hypothesis, the motor resonance during action observation occurs due to memory retrieval of the execution of the observed action. The memory, triggered by a visual stimulus, was formed in the past, when the observer executed the particular action with visual guidance. These memory-triggered mirror neurons are a product of associative learning in the sense of Pavlovian conditioning and are extensively trained by correlated experience, even if the executed action is different from a simultaneously observed action. Using this correlation account, Heyes explains differences between humans and monkeys. She claims that “humans receive a great deal of more correlated experience of observing and executing similar actions” and so the human mirror system can react to a greater variety of stimuli. We (Farkaš et al., 2011b) agree that correlated experience plays a crucial role in the development of mirror neurons. However, we believe they are not an insignificant byproduct of some other processes or a random phenomenon, but a functional piece of a larger mechanism underlying action understanding and (dependent on the location in the brain) understanding in general.

In (Farkaš et al., 2011b) we claim that both genetic factors and sensorimotor experience are *crucial* for emergence of mirror neurons. We argue that genetics expresses itself largely in terms of cortical wiring at various levels of granularity, while experience manifests itself in tuning synaptic efficiencies and potentially also in some degree of synaptic rewiring. Architectural constraints are much more credible as supported by numerous neuroscience evidence. These constraints relate to various levels of granularity ranging from neuron-level constraints, such as specification of neuron types and their associated characteristics, via local neural circuits (e.g. layered organization of the cortex at various parts of the brain, degree of interconnectivity in terms of “fan-in” and “fan-out”) to global architectural constraints in the brain, like

the thalamo-cortical or cortical-cortical pathways.⁴ These constraints, however, do not exclude mechanisms for possible rewiring, evoked by changing experience (brain plasticity). Lastly, chronotopic constraints are reflected in the timing of events in the developmental process.

Various stages of cognitive development, for instance in language acquisition, or the “theory of mind”, are known to take place at a typical age of a maturing child. This may have some genetic basis which can be slightly modulated by individual characteristics and the environment. Therefore, a plausible account for the formation of the mirror neuron system will necessarily depend on obtaining a right balance between nature and nurture factors so that they interact correctly. Oberman and Ramachandran (2008) present a similar view (in an open peer commentary in Hurley, 2008) and in addition, they propose an experiment on newborn monkeys that could help to disentangle the inborn capabilities and learning. They also propose an alternative experiment with adult monkeys that might shed light on whether mirroring ability could be achieved by Hebbian associations.

Hebbian account

A similar account to the one of Heyes (2010) is the Hebbian account of Keyesers and Perrett (2004), who claim that the existence of mirror neurons can be explained on the basis of anatomical connections between core circuits, i.e. STS, PF and F5, and the Hebbian learning rule. In her debate, Heyes (2010) criticizes also the Hebbian account and emphasizes the distinction between it and her association hypothesis. She states that the Hebbian learning only implies contiguity, whereas the associative account requires both contiguity – the closer the two events occur in time, the stronger the association, and contingency – required correlation or predictive relationship between the

⁴We argue that it is possible (and necessary) to distinguish between representational constraints and architectural constraints even if they have similar consequences. An instance of a representational constraint can be that a concrete neuron in F5 is connected to the neuron in F1 which is connected to the, say, thumb on the left foot. An instance of an architectural constraint can be the forced growth of neural synapses between F5 neurons and F1 neurons to follow a specific distribution; what in turn causes that some neurons will really (statistically) be connected to the left-foot-thumb neuron in F1.

events. Farkaš et al. (2011a) believe there is no significant difference between the two accounts and that also the Hebbian learning requires both contiguity and contingency, although the latter was not explicitly mentioned in the original Hebb's postulate.

Although the computational mechanisms operate at the level of neurons, it should be kept in mind that the both the associative and the Hebbian accounts of mirror neurons are actually a high-level psychological explanation that links “spatial” (non-sequential) sensory and motor patterns. As argued by Knott (2012), in the context of associations between STS and F5, mediated by PF, the Hebbian account assumes that sensory and motor representations, that have inherently natural sequential structure, are first independently integrated in time to yield static representations, which can then be associated in one-to-one fashion. For instance, a particular sensory representation corresponding to a concrete grasp type is associated with a corresponding underlying motor representation (Keysers and Perrett, 2004). It is an open question whether we can rely on such a type of association, or whether we need to consider the sequential nature of these learned associations. I will address the sequential nature of actions and the way of representing them in higher-order level in Chap. 5.

Counter-mirror neurons

An interesting question emerging from recent single-cell evidence on both monkeys (Kraskov et al., 2009) and humans (Mukamel et al., 2010) is that mirror neurons as such might not gain their firing properties due to their own setup, but due to the structure of brain areas. Both above mentioned studies report neurons which have mirroring properties, but with opposite pattern of firing. These neurons are hypothesized to inhibit self-movement during action observation. This is consistent with the ideomotor and common coding theories from the previous chapter. Basically, if we observe, we always have a tendency to mirror – repeat the observed action. Counter-mirroring mechanism might help us to distinguish between the situation when we produce movement, and the situation when we observe it.

2.2.9 Mirror neurons and high-level representations

As described in Sec. 2.2.6, different types of mirror neurons react to different categories of stimuli. There are strictly congruent mirror neurons that react to particular type of grasp, but there are also broadly congruent mirror neurons that tend to respond more to the goal of an action rather than to specific movement which has to be done to fulfill it. Therefore a hierarchy inside the mirror neuron circuitry can be assumed, which might contribute to both specific, but also broad high-level representations.

An important role in action understanding and MNS functionality is played by the superior temporal sulcus (STS) mentioned in Sec. 2.2.1. This visual area is sensitive to a large variety of biological movements, but it lacks multi-modal properties displayed by mirror neurons so it is not considered a true part of MNS. On the other hand, STS is one of the primary sources of visual information for the mirror neuron areas. STS is connected with F5 through two distinct pathways (Nelissen et al., 2011). Firstly, the posterior part of STS (STSp) is connected with F5c through PF (PFG). Secondly, the anterior part of STS (STSa) is connected with F5a through AIP. Borra and Rockland (2011) suggest that there are even more sources of visual information for area F5, including prefrontal area BA12.

STS is a very diverse and interesting brain area with topologically separated classes of neurons. For instance the lower banks of STS (TEa) encode the biological movement, but neurons in upper banks (TPO a PGa) encode rather the identity of the observed individual and classes of movements observed. Based on the location in STS, there are neurons that are sensitive to viewpoint from which the object is observed (e.g. front view, side view, etc.), but also neurons that are invariant to it (object-centered), what applies to perception of faces (Perrett et al., 1989, 1991), but also movements (Jellema and Perrett, 2006). This phenomenon is illustrated in Fig. 2.4.

Perrett et al. (1991) assume, that this phenomenon might be an outcome of hierarchical organization. All variant STS neurons representing different views of an object/action feed information to invariant neurons. This allows the invariant neurons to react to the object/action regardless of the viewers

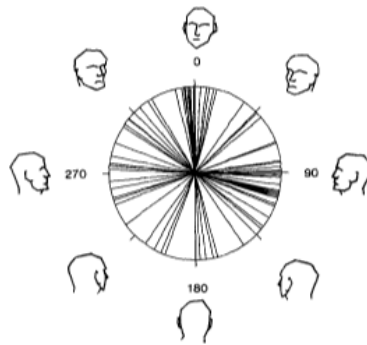


Figure 2.4: Schematic depiction of perspective variant neurons in STS of a macaque monkey (Perrett et al., 1989).

perspective and provide a high-level (categorical) representation. This phenomenon is referred to as the response pooling of lower-level units towards higher-level units. Jellema and Perrett (2006) also show, that these invariant or more precisely object-centered neurons have different anatomical location inside the STS. Neurons that are selective to viewpoint are located in posterior areas (STSp) and viewpoint-invariant in anterior banks (STSa), which are close to frontal cortices. It is well known that representation frontal cortices are of a very high-level nature. This anatomical tendency of having more general representation in locations closer to frontal cortex might not apply only in STS, but also in other parts of the parietal and prefrontal lobes, such as motor (mirror neuron) areas.

In addition, variant and invariant properties have been recently discovered by Caggiano et al. (2009) also in responses of the mirror neurons in monkeys' area F5. In their experiments, monkeys observed grasping actions filmed from three different perspectives, namely the self-observing view (0°), the side view (90°) and the opposite view (180°), which is illustrated in Fig. 2.5. Caggiano et al. found both variant and invariant mirror neurons (roughly in the 3:1 ratio).

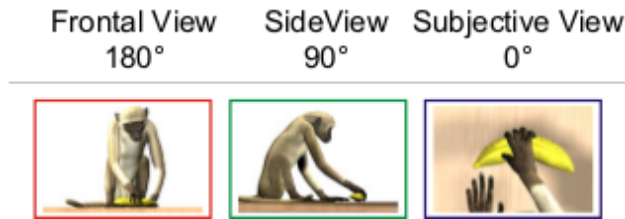


Figure 2.5: Illustration of three perspectives eliciting different mirror neuron responses (Caggiano et al., 2009).

2.3 Motor resonance in language comprehension

We introduced the ideomotor theory of action initiation and the motor-based theory of action understanding. The main assumption in both of them is the common coding for perception and action. In the following section I will describe empirical evidence on the involvement of motor resonance in language comprehension, suggesting a possible amendment to the multimodal sensorimotor codes.

2.3.1 Neural evidence

One of the richest sources of evidence for motor involvement in language comprehension are the results of the contemporary neuropsychology. For instance, Pulvermüller et al. (2001; 2005) or Hauk et al. (2004) measured activity in motor areas of the brain during comprehension of simple action verbs connected to different effectors, like “kick” executed with leg, “pick” with hand, and “lick” with mouth. Results from various experiments showed somatotopic activation in motor cortex⁵ only 250 ms after the stimulus onset. This means, that language evokes motor resonance in an automatic, involun-

⁵Somatotopic organization in the motor cortex was described already by Penfield and Rasmussen (1950). It means that different parts of the body are represented in different locations of the cortex, similarly to a map. The mouth and articulators are represented close to the sylvian fissure, the arms and hand at dorsolateral sites and the foot and leg in the vertex and interhemispheric sulcus.

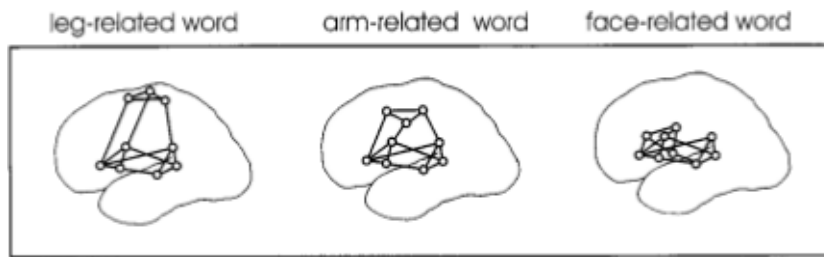


Figure 2.6: Schematic depiction of brain activation resulting from listening to various action verbs from Pulvermüller et al. (2001).

tary fashion (before the word reaches consciousness), similarly to what James (1890) expected about ideas of actions evoked by their anticipation.

2.3.2 Behavioral evidence

Another important source of empirical evidence comes from psycholinguistic experiments. One of the most influential findings was the *action–sentence compatibility effect* (ACE) (Glenberg and Kaschak, 2002). It demonstrates the presence of motor resonance in language comprehension on the basis of interference. Glenberg and Kaschak measured reaction times of participants judging the sensibility of a special type of transitive sentences called transfer sentences. Such sentences always include two agents (agent and patient) and an object that was transferred from one to another, while the participant is always one of them. For example a sentence “Jane handed you a book” includes Jane transferring the book to the addressee of the sentence. In the experiment, participants had to provide their answers using two buttons, one near and one far away from their body. Results showed that the reaction time was significantly shorter when participants had to make a move congruent with the direction implied by the sentence, in comparison with the incongruent direction. For instance, when the sentence was “Jane handed you the book”, implying the direction from “Jane” to “you”, the reaction was quicker when the “yes” button was near the body, implying movement from Jane to the recipient. Subsequently, Glenberg et al. (2008) showed that ACE applies not only to concrete, but also to abstract transfer sentences, for instance “Jane told you a story”.

In a similar experiment, Zwaan and Taylor (2006) tested the influence of both motor response interference and additional visual stimulus, moving congruently or incongruently with the movement implied by the sentence. They focused on a different movement, concretely the rotation. An example sentence with rotation would be “Jane turned down the volume”. Participants were instructed to make sensibility judgments rotating a knob either clockwise or counterclockwise. The results of various different experiments showed the difference in reaction time for congruent and incongruent movements. The comprehension of sentences was also influenced by the sight of a rotating object (on a computer screen). This visual input influenced not only processing of the linguistic input, but also the hand movement itself.

2.3.3 Language and mirror neurons

Gallese and Lakoff (2005) suggested a direct relationship between language comprehension and mirror neurons. They proposed the sensorimotor system as the most likely neural substrate for representation of concepts, including abstract concepts. In the center of their theory is the so-called *neural exploitation*, the adaptation of mechanisms for perception and action (like motor resonance) to mediate also “higher” cognitive functions as language use or reasoning. Gallese and Lakoff claim that there is nothing like a specialized language center in the (human) brain. This theory suggests something similar to my note in the beginning of this section: language might be integrated into common sensorimotor action codes and to sensorimotor simulation. A weak version of this claim would be that language is not integrated, but can trigger these common codes.

Recently, Knott (2012) proposed a similar, but less radical account, proposing that not necessarily brain circuits, but the general mechanisms for action execution and language about action might be quite similar. He points out, that sentences share the deep syntactic structure across languages. This *logical form* of sentences might be viewed as descriptions of sequences of “attentional operations” of moving one’s attention from one part of the observed scene (actor) to another (to-be-manipulated object).

2.4 Summary

Summarizing the given evidence on ideomotor action and common coding theory it seems that motor resonance plays an important role in understanding the observed actions and in evaluating their outcomes. Motor resonance can be viewed as a sensorimotor simulation of the observed action. The greater the activation of the motor circuitry, the better prediction about the observed action can be produced. The common coding theory was formulated as a theoretical framework without specific predictions on its neural implementation. However, the discovery of *mirror neurons* in early 1990's provided an interesting view on a possible neural mechanism behind the motor resonance (see Pineda, 2005) and formed a likely candidate for (at least a partial) neural implementation of common coding for action and perception.

The execution–observation matching property of mirror neurons gave rise to a revival of motor-based theories of action understanding. The existence of mirroring activity during action observation cannot be doubted. However, the exact meaning and full functional value of this phenomenon is yet to be assessed. Action understanding indeed involves visual cortices of the brain such as STS. Nevertheless, the motor component might indeed be needed, especially in cases when the observer has to judge the possible outcomes of the observed actions and understand the goal followed by the observed agent. Mirroring mechanisms might also quite naturally explain social phenomena such as empathy.

From the growing amount of empirical evidence on the involvement of motor resonance in language comprehension it seems that language capacity is rather distributed across the brain than limited to a single specialized area. It seems that grounding of concepts and language is accomplished using sensorimotor representations and sensorimotor simulation. This direction of research on motor resonance and common coding suggests a possible amendment to the multimodal sensorimotor codes in form of their linguistic labels.

In the next chapter, I will introduce a thorough overview of computational models that implement the mirror neuron system, sensorimotor interaction

and prediction, as well as models that ground meaning in sensorimotor loops. Subsequently I present my account on this topic.

Chapter 3

Artificial neural network architectures and learning

In this chapter, I describe architectures and learning algorithms for most of the artificial neural networks (ANN) mentioned in the following chapter this thesis, which summarizes various computational models of mirror neurons and the grounding of meaning grounding. The chapter is divided into two sections based on the learning principle: error driven and unsupervised learning techniques.

Error-driven or supervised learning in artificial neural network uses a teacher signal from the environment and adapts the network weights based on difference between the desired outcome and the network's actual outcome. Supervised techniques are most suitable for learning various input-output mappings. An exception in this section is the reinforcement learning and CACLA subsection (3.1.2), which basically does not belong to supervised learning techniques, but forms a new, third paradigm of learning. However, since I focus on describing artificial neural networks, which are the building block of related CACLA-based models, I introduce this approach along with them, rather than dividing the chapter into three parts.

Unsupervised techniques, on the other hand, work without any external signals (except the input data) and aim at finding statistical regularities in the processed data. Thus unsupervised learning and architectures are

most suitable for data clustering, categorization, and pattern encoding and recovery. The most prominent unsupervised learning principle is the Hebbian learning based on the Hebb's postulate about learning in the brain, i.e. that “neurons that fire together, wire together” meaning that neurons that are concurrently activated, gain stronger synaptic weights than other neurons.

Regarding the biological plausibility, unsupervised learning is evidently prominent in the brain, but also supervised and reinforcement learning can be found (O'Reilly and Munakata, 2000; Doya, 2000).

3.1 Error-driven learning

3.1.1 Multi-layer perceptron

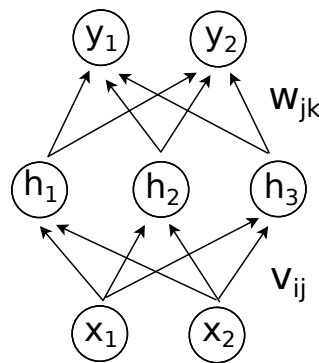


Figure 3.1: Schematic depiction of a two-layer perceptron.

Since their emergence in 1980's, the multi-layer perceptrons (MLP) have gained a firm position among various applications of computational modeling. They are quite simple and computationally powerful, able to solve various tasks such as function approximation, pattern classification or (in some cases) prediction. A standard MLP consists of an input layer \mathbf{x} , one or more hidden layers \mathbf{h} , and an output layer \mathbf{y} with weight matrices $\mathbf{v}(\mathbf{x} \rightarrow \mathbf{h})$ and $\mathbf{w}(\mathbf{h} \rightarrow \mathbf{y})$. Each projecting layer contains a trainable bias input fed with constant input -1, so when computing the layer activation the input vector has $k + 1$ nodes (where k is the actual size of the projecting layer). A generic MLP architecture is displayed in Fig. 3.1. Units in the network compute a weighted

sum of an incoming activation, which is then subjected to the activation function f (usually sigmoid, e.g. logistic) according to:

$$h_j = f\left(\sum_{i=1}^{n+1} v_{ij}x_i\right), \quad (3.1)$$

and

$$y_k = f\left(\sum_{j=1}^{q+1} w_{jk}h_j\right). \quad (3.2)$$

The MLPs are usually trained in a supervised manner using error back-propagation (BP) (Rumelhart et al., 1986) which functions in two phases. First, the network produces an estimate at the output layer (forward pass). Subsequently, the error computed as the difference between desired and estimated values on the output layer is propagated through the network in the backward direction (backward pass) and weights are updated according to:

$$\Delta w_{jk} = \alpha \delta_k h_j, \text{ where } \delta_k = (d_k - y_k) f'_k, \quad (3.3)$$

and

$$\Delta v_{ij} = \alpha \delta_j x_i, \text{ where } \delta_j = \left(\sum_k w_{jk} \delta_k\right) f'_j, \quad (3.4)$$

where $\alpha > 0$ is the learning rate.

As summarized by Haykin (2007), the computing power of BP lies in its two main attributes: it is (1) a local method for updating weights and biases, and (2) an efficient method for computing all partial derivatives of the cost function with respect to free parameters. Interestingly, despite its simplicity, a simple MLP can serve also as a forward model (Wolpert et al., 2003) (see definition in Sec. 4.2.4), which is usually modeled using more suitable recurrent neural networks. In this setup, the input and output layers of the model share their dimensionality and interpretation, since the output layer represents the estimate of the network about the next state of a (sensorimotor) sequence. The network can be taught the target behavior simply using the BP. However, BP method is considered biologically implausible, since it imposes the target difference on the whole network. Metaphorically, using

the BP for learning a robot to reach and grasp an object is like pulling its arm to let it remember the correct movement. However, when a child learns a movement it exploits the so-called motor baling, meaning that it makes a lot of different movement and afterwards evaluates their beneficiality. This process is captured in a more ecologically valid method called the *reinforcement learning*.

3.1.2 Reinforcement learning and CACLA

Reinforcement learning (RL) (Sutton and Barto, 1998) is an ecologically valid method of training artificial systems based on notions from (psychological) behaviorism. Like in supervised learning techniques, the agent (or a network) learns on the basis of numerical signal. However, in RL the agent is evaluated on arbitrary basis, i.e. the reward or punishment signal is given to the agent with an arbitrary time delay Maly (2013). In this way, the agent has no immediate feedback about appropriateness of the action it just performed. Therefore, the evaluation of the current performance has to be estimated by the agent on the basis of previous experiences. The overall goal of an agent is to maximize the long-term reward. Basically, RL is a form of trial and error learning.

The RL task can be viewed as Markov Decision Process (MDP), which operates on a discrete set of states. A MDP has the Markov property, i.e. the probability of the next state depends only upon the present state. The agent can visit states and gain reward according to the *reward function*, which stands for the feedback of the environment for agent's actions. The agent's actions are chosen according to an adaptable *policy*. The goal of an agent is to find the optimal policy to gain maximal reward from its actions. During its course of action an agent holds estimates about the states by adapting the *value function*. The value function represents what is good for the agent "in the long run" (Sutton and Barto, 1998). Although most of the traditional reinforcement learning algorithms were designed for small finite sets of states and actions, algorithms for finding good policies in continuous domains are studied as well.

Van Hasselt and Wiering (2007) introduced a new class of learning algorithms for continuous spaces and actions, named the *continuous actor-critic learning automaton* (CACL). It is based on the *actor-critic architecture*, which consists of two entities, the actor and the critic (Fig. 3.2). The executive part is the actor, which predicts the next best action of an agent. The actor behavior is modulated by the critic which evaluates the performance of the actor and provides it with internal rewards and punishments for its immediate actions. In continuous spaces, instead of a standard discrete transition tables, ANNs can be used as function approximators representing the actor and the critic (typically, standard two-layer perceptrons are used). The actor maps the agent's state onto the agent's action and the critic maps the same input onto evaluation of the current state (one value). Learning in CACL is summarized in Algorithm 1. The value function $V(s)$, computed by the critic, is updated based on temporal differences between agent's subsequent states, using

$$V_{t+1}(s_t) = V_t(s_t) + \alpha_t \delta_t, \quad (3.5)$$

where $\delta_t = r_{t+1} + \gamma V_t(s_{t+1}) - V_t(s_t)$ is the temporal-difference (TD) error, $0 \leq \alpha_t \leq 1$ denotes the learning rate and $0 \leq \gamma \leq 1$ is the discount factor. The reward r_{t+1} is received by the agent immediately after executing the action, which results in a state transition from s_t to s_{t+1} . It is known that using the update given by Eq. 3.5 for the discrete RL will result in convergence of the values to the actual expected rewards for a fixed policy (Sutton and Barto, 1998). CACL extends the usability of this update in continuous RL by yielding accurate function approximators.

Unlike other actor-critic methods, in CACL it is not desirable to update the policy in the opposite direction when the sign of the TD error is negative. An extreme case would be considering an actor that already outputs the optimal action in each state for some deterministic Markov decision processes, so for most exploring actions, the TD error would be negative. If the actor was updated away from such an optimal action, its output would no longer be optimal (van Hasselt, 2012, p.31). Using only positive values of δ_t can hence be seen as a drive for improvement, and the improper behavior is

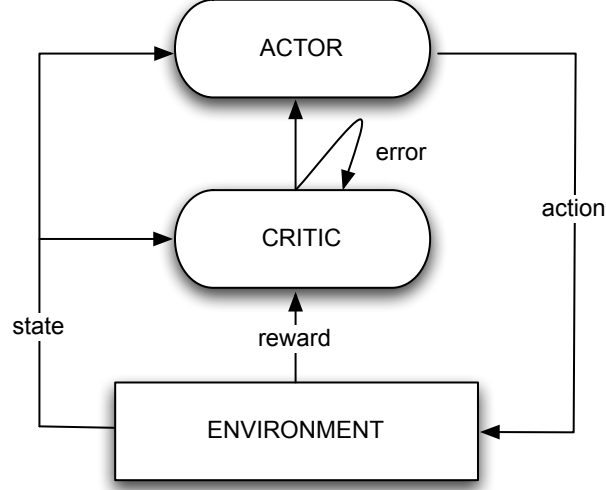


Figure 3.2: Scheme of actor-critic learning paradigm.

Algorithm 1 CACLA learning algorithm

$s_0 \leftarrow$ initial state

Initialise parameters $\theta_{i,0}^A$ and $\theta_{i,0}^C$ randomly (to small values)

for $t = 0, 1, 2 \dots$ **do**

$\tilde{a}_t \leftarrow a_t(s_t)$ using exploration

 perform action \tilde{a}_t and move to s_{t+1} , get new r_{t+1} and $V_t(s_{t+1})$

 update critic's parameters:

$$\theta_{i,t+1}^C = \theta_{i,t}^C + \alpha \delta_t \frac{\partial V_t(s_t)}{\partial \theta_i^C(t)} \quad (3.6)$$

if $V_{t+1}(s_t) > V_t(s_t)$ **then**

 update actor's parameters:

$$\theta_{i,t+1}^A = \theta_{i,t}^A + \alpha (\tilde{a}_t - a_t) \frac{\partial a_t(s_t)}{\partial \theta_{i,t}^A} \quad (3.7)$$

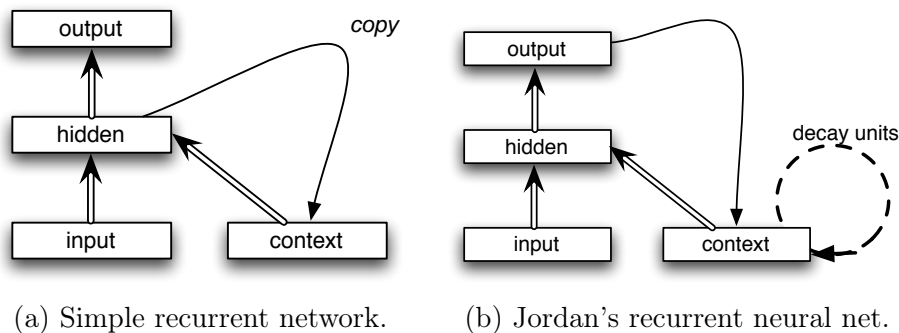
end if

end for

implicitly unlearned by learning a new behavior. In other words, CACLA only updates the actor when actual improvements have been observed. This principle helps to avoid slow learning when there are plateaus in the value space and the TD errors are small. Additionally, van Hasselt and Wiering (2009) have shown, that CACLA is comparable to commonly used discrete TD methods such as SARSA or Q-learning.

3.1.3 Recurrent neural network

Recurrent neural networks (RNN) are characteristic with their connectivity enhancing the network's weights in time. The most common are the partially recurrent ANN which are, basically, multi-layer perceptrons enhanced with context layer (or layers) which contain the information from previous time step. The temporal information from the context layer allows recurrent networks to manage tasks that MLPs cannot, like prediction and generation of sequences or modeling finite state automata. A prominent RNN model is the simple recurrent network (SRN) developed by Elman (1990). In SRN the context layer represents a copy of the network hidden layer activation from the previous step. The context layer has weighted connections to the hidden layer as displayed in Fig. 3.3a. Elman (1990) originally trained the SRN to learn a simplified grammar based on its ability to predict the next character or the next word in a sentence.



Another similar and widely used architecture is the Jordan's net (Jordan, 1986) (Fig. 3.3b), in which the context layer represents the previous activation on the output layer. Additionally, the decay units allow the con-

text to contain information from desired number of steps, i.e. $c_i(t+1) = y_i(t) + \alpha c_i(t-1)$. Unlike SRN, the Jordan's net is not only able to recognize input sequences, but also to generate valid sequences on the output.

Although RNN can be trained using standard BP algorithm. Although, better results are obtained using algorithms that account also for the temporal factor. For this purpose, the back-propagation through time (BPTT) algorithm was created (Rumelhart et al., 1986). BPTT works on the same basis as standard BP, but it takes into account also the activation from the previous time steps. The algorithm literally copies the network in time $1 \leq t \leq T$, where T indicates the length of the input sequence resulting in a network with $2(T+1)$ neurons. This “unfolding” process is illustrated in Fig. 3.3. Note that no weights are adapted during the unfolding, i.e. they are time-independent. To adapt weights the algorithm first computes partial derivatives $\partial E / \partial w_{ij}^{(t)}$ for each t . Subsequently, weights in all weight matrices are updated according to:

$$\Delta w_{ij} = -\alpha \sum_{t=1}^T \partial E / \partial w_{ij}^{(t)}, \quad (3.8)$$

where $\alpha > 0$ is the learning rate. This method is quite space-demanding (complexity $\theta(m+n)h$, where n is the number of neurons, m is the dimensionality of the input and h is the length of an epoch) and less computationally demanding ($\theta(n^2)$) (Williams and Zipser, 1995).

An alternative to BPTT is the real-time recurrent learning (RTRL) (Williams and Zipser, 1995). In short, RTRL tracks the influence of each weight on the activity of each output neuron and uses this information when adapting weights. Although this algorithm does not require the unfolding of network activations from the past (space complexity $\mathcal{O}(n^2)$, where n is the number of neurons), it is still very computationally demanding (complexity $\mathcal{O}(n^4)$).

3.1.4 Recurrent neural network with parametric biases

Recurrent neural network with parametric biases (RNNPB) (Tani, 2003; Tani and Ito, 2003) is a modified Jordan RNN with one hidden layer. An addition

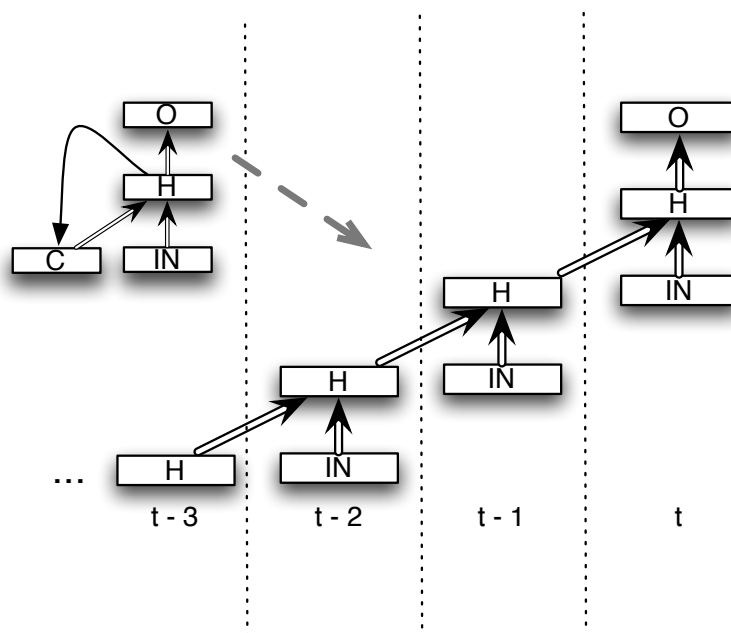


Figure 3.3: Schematic depiction unfolding a generic RNN in time.

to Jordan's architecture are the parametric bias nodes (PB), which are a part of the input layer and serve as bifurcation parameters. A scheme of two generic RNNPB is depicted in Fig. 3.4. The basic function of the model is to predict the sensorimotor activation vector for the next time step and thus control the movement of an agent (a robot). All layers except the input layer have sigmoid activation functions. The values of PB nodes during concrete behaviors can be stored and subsequently (manually) fed as an input to the network, which will then produce the learned behavior. Interestingly, the PB vectors are self-determined by the network in an unsupervised manner.

In a standard setting, the network works as a forward model. The input layer consists of sensory input neurons, the PB neurons, and context neurons that receive input from the context output neurons from the previous time step. Next, there is a hidden layer whose size depends on the input size, and an output layer consisting of prediction neurons (the same number as input nodes) and context neurons. All layers except the input layer have sigmoid activation functions.

During the learning of various sequences of inputs, two different mechanisms apply, one for weight modification, the other one for the PB vector modification. While the weight matrix is the same for all sequences, the PBs differ in each case. Both are computed simultaneously using BPTT to minimize the value of the learning error function E over all training sequences q_k using:

$$E(W, p_0, \dots, p_{s-1}) = \sum_{k=0}^{s-1} E_k(W, p_k), \quad (3.9)$$

and

$$E_k(W, p_k) = \sum_{t=0}^{l_k-1} \sum_{n \in ONodes} (r_{kn}(t) - o_{kn}(W, p_k, t))^2, \quad (3.10)$$

where E_k is the learning error of the training sequence q_k , l_k is its length, p_k is its corresponding PB vector, s is the number of training sequences, W is the whole network's weight matrix, N is the number of output nodes, $o_{kn}(W, p_k, t)$ is the output of the output node n (from the set of all output nodes $ONodes$) at the time step t , and $r_{kn}(t)$ is the desired output of this unit.

Connection weights are initialized randomly and then iteratively updated according to:

$$\delta^2 w_{nm}^{(T)} = - \frac{\partial E(W^{(T)}, p_0^{(T)}, \dots, p_{s-1}^{(T)})}{\partial w_{nm}}, \quad (3.11)$$

$$\delta w_{nm}^{(T)} = \eta_w \cdot \delta w_{nm}^{(T-1)} + \epsilon_w (1 - \eta_w) \cdot \delta^2 w_{nm}^{(T)}, \quad (3.12)$$

and

$$w_{nm}^{(T)} = w_{nm}^{(T-1)} + \delta w_{nm}^{(T)}, \quad (3.13)$$

where the delta error $\delta^2 w_{nm}^{(T)}$ at a training iteration T is computed and back-propagated to the connection weight from node n to node m , η_w stands for the learning rate and ϵ_w represents the time constant of the update modification.

Unlike connection weights, the PB vectors encompass the whole sequence. Therefore in each PB vector p_k each j -th element is initially set to 0.5 and

then updated every training iteration T according to:

$$\delta^2 p_{kj}^{(T)} = -\frac{\partial E(W^{(T)}, p_0^{(T)}, \dots, p_{s-1}^{(T)})}{\partial p_{kj}}, \quad (3.14)$$

$$\delta p_{kj}^{(T)} = \eta_p \cdot \delta p_{kj}^{(T-1)} + \epsilon_p (1 - \eta_p) \cdot \delta^2 p_{kj}^{(T)}, \quad (3.15)$$

and

$$p_{kj}^{(T)} = p_{kj}^{(T-1)} + \delta p_{kj}^{(T)}, \quad (3.16)$$

where $\delta^2 p_{kj}^{(T)}$ is the delta error, η_w and ϵ_w are learning parameters with the same function as in the case of weights.

After the learning phase, the network can generate learned sequences according to the PB vector set as an input. The PB vector can be taken from another RNNPB network or from the set of all self-determined PB vectors from the network. To achieve this, PB vectors can be computed through regression of the past sequence pattern (using the regression window), for each sequence stored, as well as for newly encountered sensory data. RNNPB can also operate in the so called *closed-loop mode*, in which the input of the network consists only of predictive output. In this way it can generate imaginary sequences without receiving the actual input from the environment.

To employ the task of learning in two separate domains (like behavior and language) described in 4.3.2 two RNNPB can be trained simultaneously with their PB vectors “bound together”. Subsequently, the PB vectors from each net can be fed on the other’s input and elicit activity without stimulus input. To accomplish this, Sugita and Tani (2005) designed the method of *PB binding*. PB binding introduces the so-called *interaction term* to the learning equations for the PB vectors in both networks:

$$p_{s_k}^{(T)} = p_{s_k}^{(T-1)} + \delta p_{s_k}^{(T)} + \gamma_L \cdot (p_{b_k}^{(T-1)} - p_{s_k}^{(T-1)}) \quad (3.17)$$

and

$$p_{b_k}^{(T)} = p_{b_k}^{(T-1)} + \delta p_{b_k}^{(T)} + \gamma_B \cdot (p_{s_k}^{(T-1)} - p_{b_k}^{(T-1)}), \quad (3.18)$$

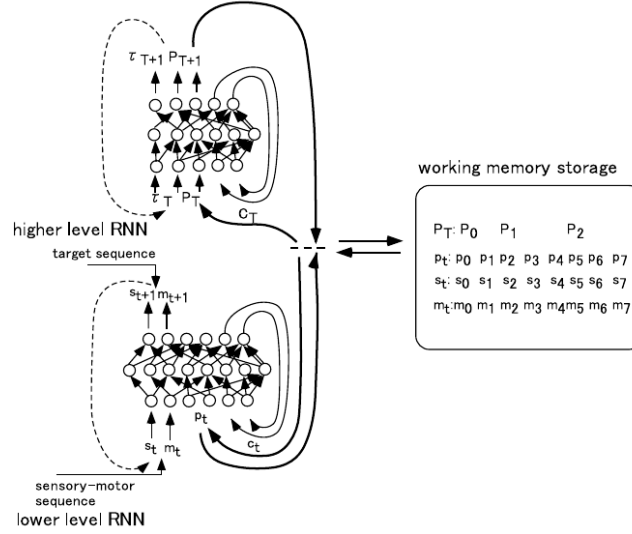


Figure 3.4: Schematic depiction of two generic RNNPB nets connected through their PB vectors using PB-binding (Tani, 2003).

where $p_{s_k}^{(T)}$ denotes the value of linguistic network PB vector for sentence s_k at time T , similarly $p_{b_k}^{(T)}$ stands for behavioral net's PB vector for behavioral sequence b_k , and γ_L and γ_B are positive coefficients of strength of the binding.

In sum, the PB allow the RNNPB architecture not only to react to particular classes of inputs, but also trigger the behavior without any input, forming a sort of procedural memory. Using PB binding, two separate RNNPB can be trained simultaneously, so the resulting PB encode the activity patterns of both networks. In this way, one network can trigger the activity in the other and vice versa. This property makes RNNPB a great candidate for modeling of both mirror neurons and language grounding, as described in Sec. 4.2.9 and Sec. 4.3.2.

3.1.5 Generalized Recirculation

The standard BP learning algorithm is known to be biologically implausible because it requires the mechanism of error propagation and it does not use locally available, activation-based variables. With this in mind, O'Reilly (1996a) designed Generalized Recirculation (GeneRec) algorithm that avoids

the computation of error derivatives, yet can lead to error minimization. GeneRec is an extension the recirculation model (Hinton and McClelland, 1988), which was restricted to autoassociation.

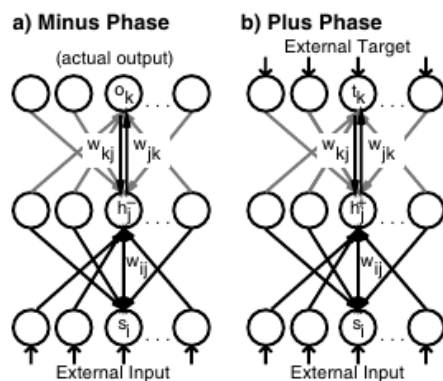


Figure 3.5: Schematic depiction of two activation phases in GeneRec model (O'Reilly and Munakata, 2000).

The GeneRec model is a three-layer network with full connectivity between layers whose activation rules are described in Table 3.1 reproduced from O'Reilly (1996a). Unlike original recirculation algorithm, which uses a four-stage activation update process, GeneRec runs in just two phases as displayed in Fig. 3.5. The model has reciprocal connectivity between hidden and output layer with symmetric weights. The activation flow starts in minus phase, when the stimulus s_i is presented. Note that the net input term at the hidden layer includes the input from both visible layers before applying the sigmoid activation function $\sigma(\eta) = 1/(1 + \exp(-\eta))$. Output units produce activations o_k^- in minus phase but can also be clamped to target activations o_k^+ at the onset of plus phase. Input units can only deliver stimuli s_i at the onset of minus phase. This model was developed in the Leabra framework (O'Reilly, 1996b), which uses dynamic units approximating the behavior of biological neurons. O'Reilly (1996a) has shown that, under certain conditions, GeneRec computes the same error derivatives as Almeida-Pineda recurrent back-propagation (Almeida, 1987; Pineda, 1987).

Table 3.1: Equilibrium network variables in GeneRec model.

Layer	Phase	Net Input	Activation
Input (s)	–	–	$s_i = \text{stimulus input}$
Hidden (h)	–	$\eta_j^- = \sum_i w_{ij} s_i + \sum_k w_{kj} o_k^-$	$h_j^- = \sigma(\eta_j^-)$
	+	$\eta_j^+ = \sum_i w_{ij} s_i + \sum_k w_{kj} o_k^+$	$h_j^+ = \sigma(\eta_j^+)$
Output (o)	–	$\eta_k^- = \sum_j w_{jk} h_j$	$o_k^- = \sigma(\eta_k^-)$
	+	–	$o_k^+ = \text{target output}$

The basic weight update rule in GeneRec is:

$$\Delta w_{pq} = \lambda a_p^- (a_q^+ - a_q^-), \quad (3.19)$$

where a_p^- denotes the presynaptic and a_q^- denotes the postsynaptic unit activation in minus phase, a_p^+ is the presynaptic activation from plus phase (in output-to-hidden direction) and λ denotes the learning rate. The learning rule given in Eq. 3.19 is applied to both input-hidden and hidden-output weights as well as to bias weights.

In his work, O'Reilly experimented with several modifications of GeneRec, determined by the weight update rules. For instance, he showed that the symmetry-preserving version of GeneRec, i.e. with symmetric hidden-to-output weights and symmetric weight update according to:

$$\frac{1}{\epsilon} \Delta w_{ij} = a_i^- (a_j^+ - a_j^-) + a_j^- (a_i^+ - a_i^-) \quad (3.20)$$

The symmetric learning rule combined with the so-called midpoint method according to:

$$\frac{1}{\alpha} \Delta w_{ij} = (a_i^+ a_j^+) - (a_i^- a_j^-) \quad (3.21)$$

results in weight update equivalent to Contrastive Hebbian learning (CHL) for training Boltzmann machines (both in stochastic and deterministic versions). GeneRec as well as CHL are based on differences between two activation phases. Forward (minus) phase involves activation propagation from inputs toward outputs producing the network estimate of the output values. Subsequent backward (plus) phase flows in the opposite direction propa-

Table 3.2: Activation phases and states in BAL model.

Layer	Phase	Net Input	Activation
\mathbf{x}	F	-	$x_i^F = \text{stimulus}$
\mathbf{h}	F	$\eta_j^F = \sum_i w_{ij}^{IH} x_i^F$	$h_j^F = \sigma(\eta_j^F)$
\mathbf{y}	F	$\eta_k^F = \sum_j w_{jk}^{HO} h_j^F$	$y_k^F = \sigma(\eta_k^F)$
\mathbf{y}	B	-	$y_k^B = \text{stimulus}$
\mathbf{h}	B	$\eta_j^B = \sum_k w_{kj}^{OH} y_k^B$	$h_j^B = \sigma(\eta_j^B)$
\mathbf{x}	B	$\eta_i^B = \sum_j w_{ji}^{HI} h_j^B$	$x_i^B = \sigma(\eta_i^B)$

gating the desired output throughout the network. In the next section I introduce a bidirectional activation-based learning (BAL), which is based on the GeneRec model, but is completely symmetrical regarding the activation propagation and the weight update rules.

3.1.6 Bidirectional Activation-based Learning

With a motivation to create biologically plausible and fully bidirectional algorithm, (Farkaš and Rebrová, 2013) derived the Bidirectional Activation-based Learning algorithm (BAL) from GeneRec. The aim for creating BAL was to mediate the high-level associations between sensory and motor representations in the proposed MNS model (Chap. 5) in more biologically plausible way than using standard BP. I present results from preliminary experiments with BAL in Sec. 5.2.

BAL shares with GeneRec the two activation phases, but differs from it by the connectivity that allows completely bidirectional associations to be established (GeneRec focuses on input-to-output mapping). Unlike GeneRec, BAL uses two pairs of weight matrices for each activation phase and computes and uses activation values for all layers in both phases. In addition, BAL does not use dynamical settling process, but computes the activations in one step as described in Table 3.2.

We avoid input-output notation of layers as used in GeneRec, because in our case not only output can be evoked by input presentation, but also vice versa. Hence, we label the two outer (visible) layers \mathbf{x} and \mathbf{y} and the

hidden layer \mathbf{h} . Let forward activation be denoted by superscript F and backward activation by superscript B. Then during the forward pass, the \mathbf{x} units are clamped to \mathbf{x}^F and we get the activations $\mathbf{x}^F \rightarrow \mathbf{h}^F \rightarrow \mathbf{y}^F$. During the backward pass, the \mathbf{y} units are clamped to \mathbf{y}^B and we get the activations $\mathbf{y}^B \rightarrow \mathbf{h}^B \rightarrow \mathbf{x}^B$ (Fig. 3.6).

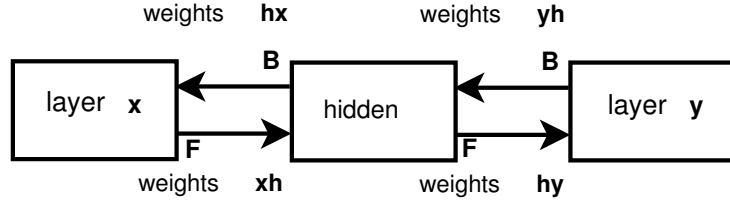


Figure 3.6: Schematic depiction of the BAL model

The mechanism of weights update partially matches that of GeneRec. Each weight in BAL network (i.e. belonging to one of the four weight matrices) is updated using the same learning mechanism, according to which the weight difference is proportional to the product of the presynaptic (sending) unit activation a_p and the difference of postsynaptic (receiving) unit activations a_q , corresponding to two activation phases (F and B, in particular order). Weights in \mathbf{x} -to- \mathbf{y} direction (belonging to \mathbf{h} and \mathbf{y} units) are updated as

$$\Delta w_{pq}^F = \lambda a_p^F (a_q^B - a_q^F), \quad (3.22)$$

where, as in the GeneRec algorithm, a_p^F denotes the presynaptic activity, a_q^F is the postsynaptic activity, and a_q^B denotes the postsynaptic activity from the opposite phase (\mathbf{y} -to- \mathbf{h}). Analogically, the weights in \mathbf{y} -to- \mathbf{x} direction (belonging to \mathbf{h} and \mathbf{x} units) are updated as

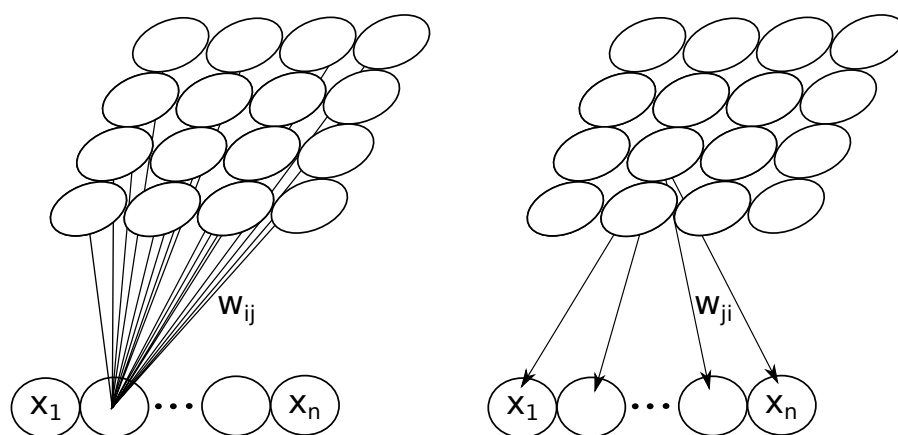
$$\Delta w_{pq}^B = \lambda a_p^B (a_q^F - a_q^B) \quad (3.23)$$

All units have trainable thresholds (biases) that are updated in the same way as regular weights and are fed a constant input 1.

3.2 Unsupervised learning

3.2.1 Kohonen's self-organising map

Self-organising maps (SOM) were introduced by Kohonen (1997) (hence they are called Kohonen networks) and became well known and popular. Extensions of SOM are still a state-of-the-art topic in artificial neural network modeling. Unlike perceptrons, SOMs acquire knowledge in unsupervised manner. A SOM consists of input nodes, one multi-dimensional layer of neurons (usually 2D) and a matrix of weight connections from all input nodes to all map nodes (Fig. 3.7a).



(a) SOM: capturing the pattern (connection weights from the input to the map).

(b) SOM: retrieving and comparing the pattern (connection weights from the map unit to input nodes).

Figure 3.7: Schematic depiction of a self-organizing map.

Each neuron can be seen as a data model – a representation of a data category, usually computed as a local average of the data that it encompasses. From the networks point of view, each neuron reacts to a given input with a continuous value that represents the distance between the input and the data model, which, basically, stands for the prototype of the data class. Unlike similar methods, SOMs have topographical organization, so the neurons that are close together represent models that resemble each other. Thanks to this

property, the SOMs are vastly used for discovering clustering structures in high-dimensional data.

Learning in SOMs resides in two basic principles: (1) competition among all nodes of the network, and (2) cooperation among neighboring nodes. When an input is presented to the network, the winner i^* is found whose distance from the input is minimum according to:

$$i^* = \arg \min_i \|\mathbf{x}(t) - \mathbf{w}_i\|, \quad (3.24)$$

where $\mathbf{x}(t)$ is the current input to the network, \mathbf{w}_i the weight vector of the i -th neuron, and $\|\cdot\|$ denotes the Euclidean norm. Subsequently, the weights of the winner and of the surrounding neurons are adapted according to:

$$\Delta \mathbf{w}_i = \alpha(t) h(i, i^*) (\mathbf{x}(t) - \mathbf{w}_i), \quad (3.25)$$

where $\alpha(t)$ is a decreasing learning rate at time t . To encompass the limitation of adapting only the weights of the neighbors to a certain degree of proximity to the winner, the neighborhood function $h(x, y)$ is introduced to the learning equation. Gaussian (Eq. 3.26) and Manhattan (Eq. 3.27) distance function are the most prominent examples.

$$h(i, i^*) = e^{-d(i, i^*)^2 / \sigma^2(t)} \quad (3.26)$$

$$h(i, i^*) = \begin{cases} 1 & \text{if } d_M(i, i^*) \leq \lambda(t) \\ 0 & \text{if otherwise} \end{cases} \quad (3.27)$$

After successful training, SOM forms a map with clusters of similar data. The vector of synaptic weights connecting any neuron j from the map with all input nodes represent a prototype of a certain class of the input data (Fig. 3.7b). SOM has been shown not only to successfully find data clusters among high-dimensional sequences, but also to generalize over new data. Additionally, prototypes of the classes of the data can be retrieved in form of the connection weights. Since SOM maintains topographic organization,

central neurons of larger clusters can be easily identified producing more general category prototypes.

3.2.2 Merge self-organising map

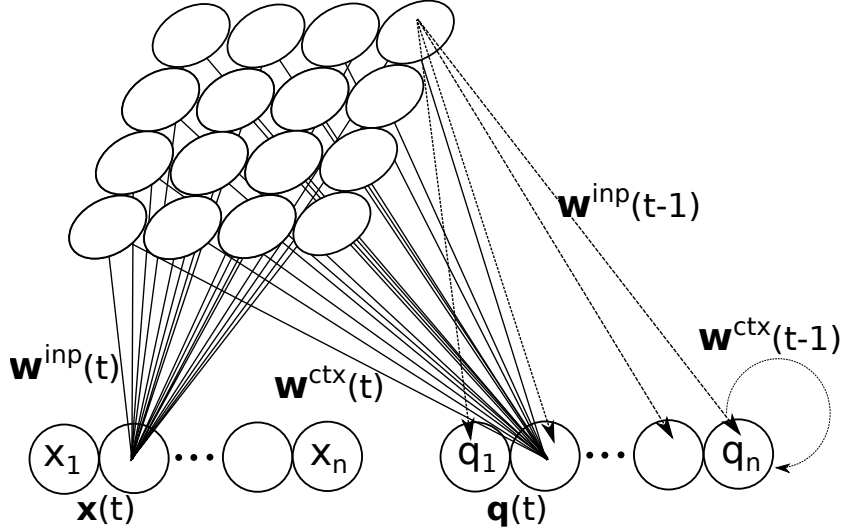


Figure 3.8: Scheme of merge self-organizing map.

The merge self-organising map (MSOM) model is based on the well-known SOM, but it has recurrent architecture (see Fig. 3.8), so it can be used for processing sequential data. In MSOM a standard SOM is enhanced with the context layer (descriptor) in a similar way as a recurrent neural network. Therefore each neuron $i \in \{1, 2, \dots, N\}$ in the map has not one, but two weight vectors associated with it:

1. $\mathbf{w}_i^{\text{inp}} \in \mathcal{R}^n$ – linked with an n -dimensional input vector $\mathbf{s}(t)$ feeding the network at time t , and
2. $\mathbf{w}_i^{\text{ctx}} \in \mathcal{R}^n$ – linked with the so-called context descriptor $\mathbf{q}(t)$ specified below.

The output of unit i at time t is computed as:

$$y_i(t) = \exp(-d_i(t)), \quad (3.28)$$

where

$$d_i(t) = (1 - \alpha) \cdot \|\mathbf{s}(t) - \mathbf{w}_i^{\text{inp}}(t)\|^2 + \alpha \cdot \|\mathbf{q}(t) - \mathbf{w}_i^{\text{ctx}}(t)\|^2 \quad (3.29)$$

Parameter $0 < \alpha < 1$ trades off the effect of the context and the current input on the neuron's profile and $\|\cdot\|$ denotes the Euclidean norm. The context descriptor is calculated based on the affine combination of both weight vectors of the previous winner according to:

$$\mathbf{q}(t) = (1 - \beta) \cdot \mathbf{w}_{i^*(t-1)}^{\text{inp}}(t) + \beta \cdot \mathbf{w}_{i^*(t-1)}^{\text{ctx}}(t), \quad (3.30)$$

where $i^*(t-1) = \arg \min_i \{d_i(t-1)\}$ is the previous winner and parameter $0 < \beta < 1$ trades off the impact of the context and the current input on the context descriptor.

The training sequences are presented in natural order, one input vector a time, and in each step both weight vectors are updated using the same form of Hebbian rule:

$$\Delta \mathbf{w}_i^{\text{inp}}(t) = \gamma \cdot h_{ii^*} \cdot (\mathbf{s}(t) - \mathbf{w}_i^{\text{inp}}(t)), \quad (3.31)$$

$$\Delta \mathbf{w}_i^{\text{ctx}}(t) = \gamma \cdot h_{ii^*} \cdot (\mathbf{q}(t) - \mathbf{w}_i^{\text{ctx}}(t)), \quad (3.32)$$

where i^* is the winner index at time t and $0 < \gamma < 1$ is the learning rate. Neighborhood function h_{ii^*} is a Gaussian (of width σ) on the distance $d(i, k)$ of units i and i^* in the map: $h_{ii^*} = \exp(-d(i, i^*)^2/\sigma^2)$. The 'neighborhood width', σ , linearly decreases in time to allow for forming topographic representation of input sequences.

As a result, the units (i.e. their responsiveness) get organized according to sequence characteristics, biased towards their suffixes (most recent inputs). Since the MSOM model shares properties with basic SOM we can assume it to be able to generalize new inputs and retrieve the data prototypes. However, to account for the sequential nature of the input data a specific method has to be developed.

Chapter 4

Computational modeling

This chapter is dedicated to computational modeling of selected cognitive phenomena. Under the term computational modeling I generally assume a methodology in which a hypothesis is evaluated using abstraction and computer simulation. In general, a computational model of this sort processes some kind of input data, for instance, stimuli received by participants of a psychological experiment, and learns to return some desired values. The aim is, especially in cognitive science, that the model behaves like expected. For instance, in case of modeling a cognitive phenomenon based on some empirical data, the model is not only desired to produce the correct answers, but also to make mistakes as the original subjects of an experiment did. Computational modeling in context of neuroscience requires the model to resemble brain areas and the connections between them.

On the other hand, in the context of robotics, a computational model plays a role of a control architecture for the robot. Unlike standard control architectures, nature-inspired models allow a robot to learn autonomously, adapt to changes in the environment, and to become intelligent in a sense. Although the “thinking machines” are yet unreachable, cognitive robotics and embodied approaches put forward very promising means to understand natural intelligence and mimic it in machines.

In this chapter, I first briefly introduce the framework of cognitive robotics. Subsequently, I describe some of the most prominent models of the mirror

neurons system. Finally, I also enter the topic of grounding the (linguistic) meaning in sensorimotor interaction within the field of cognitive robotics.

4.1 Cognitive robotics

Cognitive robotics (CR) is currently a very dynamic and one of the most attractive fields in cognitive sciences. It is characteristic with use of physical and simulated robots, usually operating in a simplified environment. CR models and experiments usually focus on a concrete problem, like acquisition of simple motor skills, sensorimotor interaction in particular tasks, acquisition of grounded lexicon, or learning a simplified language and other. Simple and small skills are built and generally assumed to be scalable to other skills and domains. Advances in human-robot interaction are also prominent aims of CR.

One of its core concepts of CR is the *synthetic methodology*, that can be characterized as “understanding by building” (Pfeifer and Scheier, 1999). Its aim is to study the cognitive development and various cognitive phenomena through their realization and emergence in artificial embodied agents, physical or simulated robots. Another core concept of CR is the embodiment and groundedness of the agent (Pfeifer and Scheier, 1999). In contrast with classic designer’s approaches (such as GOFAI), which create intelligent behavioral primarily in a top-down manner, cognitive robotics focuses on the agent itself. The cognitive robots are required to fulfill these two requirements:

1. *Embodiment*: the agent is required to have a physical body (which also applies to simulated agents), which is subject to physical powers and provides the robot with sensory information and the means to interact with the environment;
2. *Situatedness*: the agent has to be situated in the environment and learn on the basis of its interaction with the environment and other agents.

CR emphasizes that cognitive processes modeled in the agent should account for the whole agent with its body and environments. From this point of view,

a control architecture for one robot type (species) should not be transferable to another robot. Additionally, the robot’s knowledge base and skills should be acquired during its “life”, not be given to the agent previously to its interaction with environment. The development of a cognitive robot should resemble the development of human and animal infants.

Cognitive (developmental) robotics offers a stable platform to study grounding in separate cognitive abilities (Asada et al., 2009). A typical cognitive robotics methodology in the meaning grounding domain comprises a simple perceptual categorization followed by a verbal command leading to the execution of a simple action. In a different scenario the robot first produces an action and simultaneously perceives its consequences. Control architectures for these robots are usually based on artificial neural networks (see Chap. 3). Brain-inspired control architectures for cognitive robots process the sensory inputs and drive the agent-environment interaction.

Robots in the CR framework are built specifically for academic purposes. Cognitive robots strongly resemble human infants of various age or have at least humanoid shape and effectors. One of the most prominent example of a cognitive robot is the iCub.

4.1.1 The iCub robot

The iCub is a small-size humanoid robot (figure 4.1) created within the European project “Robotcub” (robotcub.org; Metta et al., 2008), specifically for cognitive robotics purposes. It is designed to resemble a 2.5-year-old child, it is 90 cm tall and has a weight of 23 kg. The iCub is endowed with 53 degrees of freedom (DoF) in joints distributed all over its body in proportion similar to human effectors (e.g. 9 DoFs for hands). In comparison with other humanoid robots of its size and type, the iCub has the highest kinematic complexity providing a very accurate model of an actual child’s body and effectors. It even has movable eyes each with a color camera.

iCub’s perception and motor control are centralized in one control system, which serves as an interface between the robot internal state and the external world. The robot communicates through an Ethernet network protocol and

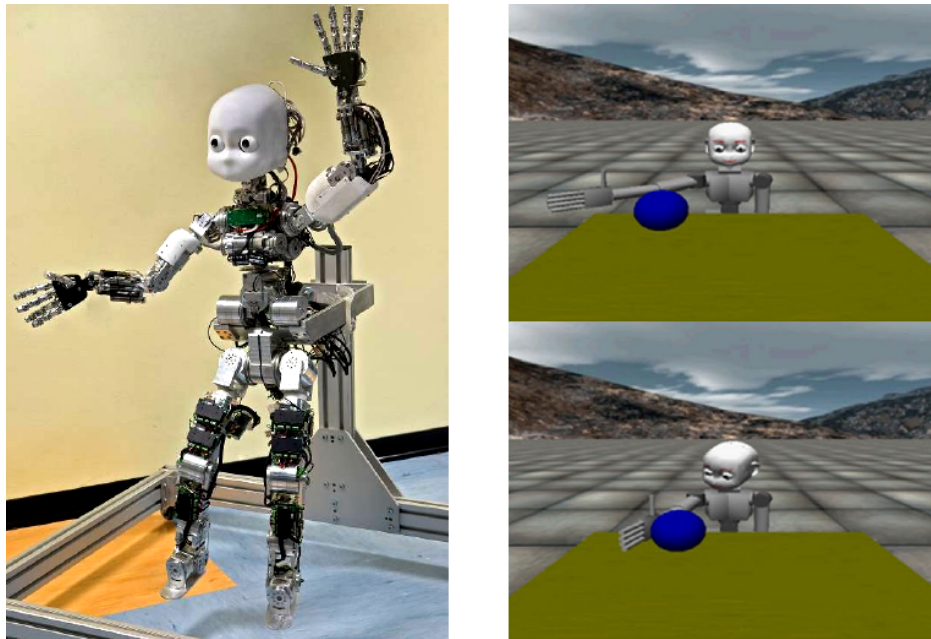


Figure 4.1: ICub robot and its simulator, figures from Metta et al. (2008), and from Marocco et al. (2010) respectively.

can be controlled (provide sensory readings and execute actions) using the integrated software platform YARP (Metta et al., 2006). This platform is also used for control of the iCub simulator (Tikhanoff et al., 2008), which is a carefully designed virtual form of iCub robot in the open dynamics engine (ODE, www.ode.org), providing a safe, yet realistic, environment for testing the control architectures before implementing them to the real robot.

4.2 Computational and conceptual models of the Mirror Neuron System

Since the discovery of mirror neurons in 1990's, a large variety of computational and conceptual models of the *mirror neuron system* (MNS) has been proposed and implemented using mainly artificial neural networks. Oztop et al. (2013) claim that computational models of MNS are the most powerful tool for explaining various aspects of mirror neuron function and emergence. Since the actual empirical research is truly very time-demanding and in most

cases in humans nearly impossible to do (single-cell recording and opening the skull in general is executed only if absolutely necessary), artificial neural networks might provide us with answers to various questions. Computational models are capable of virtually unlimited number of experiments in which various hypothesis on the function, emergence, and neural wiring of mirror neurons might be tested.

A majority of MNS models are aimed at modeling the actual neural circuitry. These models are built of modules or components that directly represent particular parts of the monkey's brain (Fig. 4.2), for instance the FARS model (Fagg and Arbib, 1998), MNS1 (Oztop and Arbib, 2002), MSI (Oztop et al., 2005), or MNS2 (Bonaiuto et al., 2007). I will briefly describe and review these models (for a very thorough review see Oztop et al., 2006). I will also present some alternative computational models aimed at explaining the emergence and function of mirror neurons.

On the other hand, there are models that are closer to the paradigm of cognitive robotics, which use properties of the mirror neurons in a specific architecture. Such models do not aim to encompass the neural circuitry, but rather endow the agent with some special capabilities. A good example is the RNNPB model (Tani et al., 2004), which combines temporal BP learning and unsupervised formation of special codes (parametric biases) that are able to trigger specific behavior of the network without an actual sensory input. The RNNPB model will be discussed in this section as well.

4.2.1 FARS

The very first modeling effort that is related to MNS is the FARS model (Fagg and Arbib, 1998). Rather than directly approaching the mirroring properties, the model focuses on the task of visually guided grasping. Mainly it analyzes how the canonical part of the MNS circuitry, centered on the AIP \rightarrow F5_{can} pathway in the macaque brain may account for this ability (see Fig. 4.2 also regarding other brain areas involved in MNS models). In the computational model, AIP represents the grasps afforded by the object while F5_{can} selects and drives the execution of an appropriate schemas. More details

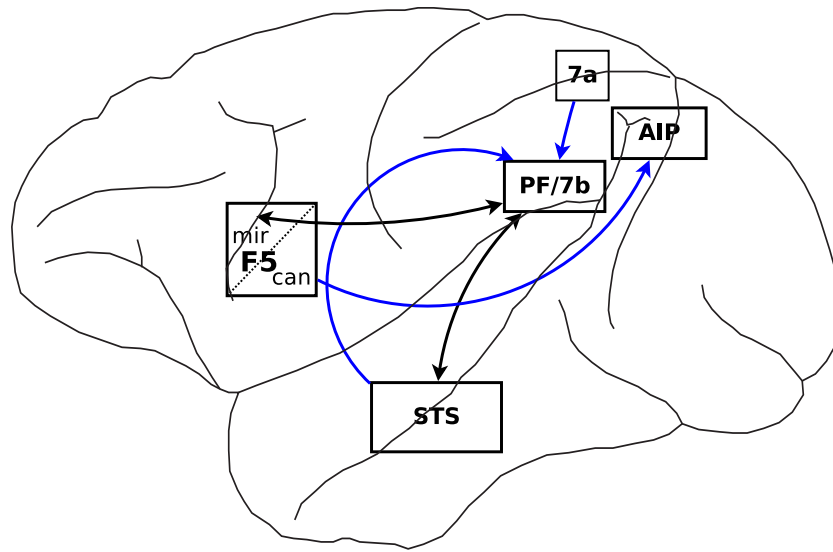


Figure 4.2: Schematic depiction of the most commonly modeled brain areas related to MNS.

on the model are contained in the dissertation thesis of Fagg (1996). Interestingly, the model also shows how task information and other constraints may resolve the problem of selecting and action from various opportunities provided by multiple affordances. The authors emphasize the contribution of various other brain areas. This complex system is depicted in Fig. 4.3. As a modeling approach, FARS uses the so-called μ -schemas, which consist of leaky-integrator neurons, each of which has a specialized function. These μ -schemas encompass movement programs as “discrete and encapsulated entities” (Fagg and Arbib, 1998) in order to separate conceptual and implementation level.

Evaluation

This approach is particularly influential for computational neuroscience. Fagg (1996) evaluated his model also on the basis of comparison with human-subject data and succeeded. Unfortunate for this model is the heuristically prewiring of the visuomotor object-to-grasp transformation. However, as an interesting and important feature of the model, this transformation can be learned using RL. The reward signal provides feedback about the success and

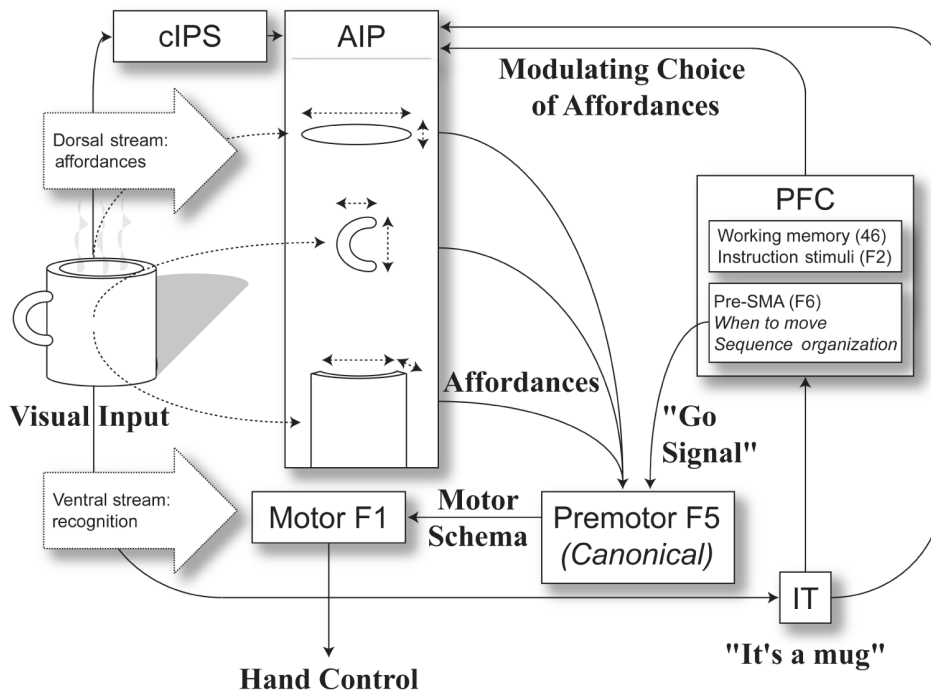


Figure 4.3: Schematic depiction of FARS model (Arbib, 2005).

efficiency of the chosen grasping action and guides the parameter tuning for better grasp configuration (Fagg, 1996).

4.2.2 Mirror neuron system 1

The MNS1 model (Oztop and Arbib, 2002) extends the modeling perspective of FARS by assuming two agents, the actor and the observer. Unlike FARS, from which it developed, the MNS1 model might be considered a true model of mirror neuron circuitry. In MNS1, the capacity to select a proper grasping movement on the basis of object affordances is used to form a capacity to recognize the observed action. The basis for action recognition are the so-called *hand states*, which encode the relation between the unfolding shape and trajectory of a hand and the affordances of an object, as well as the relative position of the hand with respect to the target object. Oztop and Arbib (2002) show that a feedforward two-layer neural network with perceptron units, representing a $7b \rightarrow F5_{\text{mir}}$ mapping (a core mirror circuit), could be

trained to recognize the grasp type from the hand-state trajectory, often achieving correct classification at a time before the hand reaches the object.

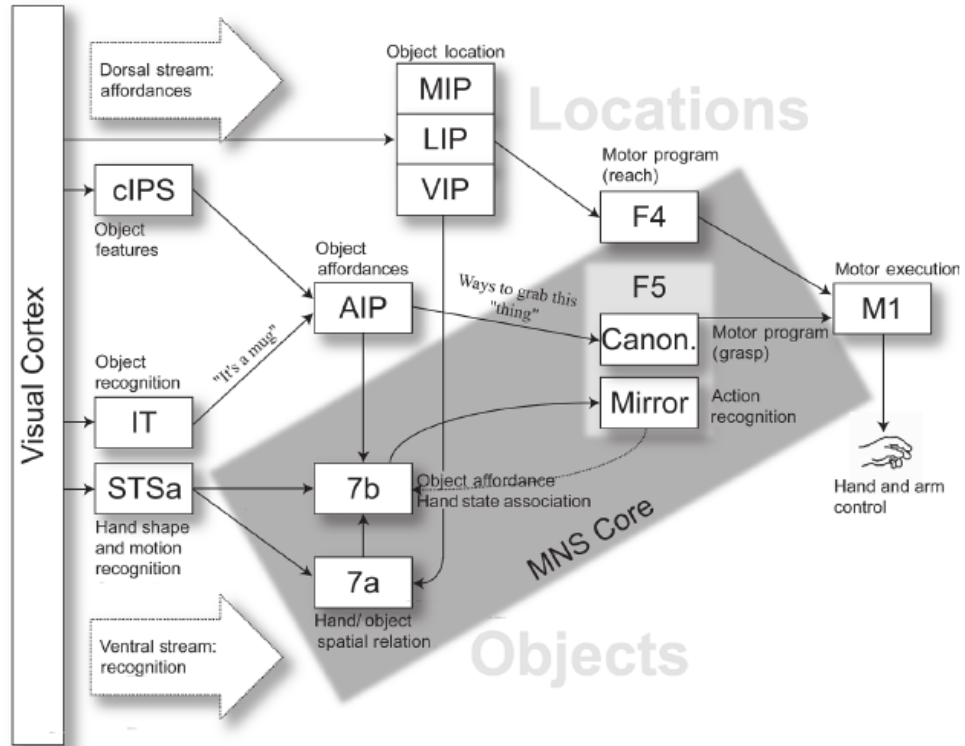


Figure 4.4: Schematic depiction of MNS1 model (Arbib, 2005).

The network takes as inputs pre-processed hand state representations assumed to be formed in STS and 7a (see Fig. 4.4), and as an output it generates action recognition signals in the $F5_{\text{mir}}$ layer (i.e. the action is recognized by the mirror neurons). The hidden layer (7b) mediates the sensorimotor mapping through the object affordance–hand state association. The model is situated in a virtual scenario that allows to generate the training sequences as well as the target responses. The network is trained using standard BP algorithm. The activity of the $F5_{\text{can}}$ neurons serves as a teaching signal for the $F5_{\text{mir}}$ neurons enabling the agent to learn which hand-object trajectories correspond to the canonically encoded grasps. As a result, the appropriate mirror neurons come to fire in response to viewing the appropriate trajectories even when the trajectory is not accompanied by $F5_{\text{can}}$ firing. The infor-

mation provided by the hand state is preprocessed, using the object-centered frame of reference, to yield an invariant representation (with respect to the agent of the action), allowing action recognition in $F5_{\text{mir}}$.

Evaluation

The MNS1 model surely serves its purpose of creating a mirroring mechanism, yet it is based on two questionable assumptions. Firstly, actions that are being learned by $F5_{\text{mir}}$ neurons are already in the monkey’s motor repertoire, so what the model learns is to associate the known motor representation with the known visual representation. However, in real life situation, actions are learned simultaneously with their perceptual qualities and affordances that drove their execution. Secondly, and more importantly, the hand-state trajectory conversion to the invariant representation is not yet firmly grounded in the empirical evidence (described in Sec. 2.2). The model does not even respect the $F5_{\text{mir}} \rightarrow \text{PF}(7b) \rightarrow \text{STS}$ pathway (discovered as first), which connects perspective-variant area STSp with area F5c. Although the model contains AIP area, it misses also the $F5_{\text{mir}} \rightarrow \text{PF}(7b) \rightarrow \text{STS}$ pathway, which was re-discovered with respect to mirror neurons only recently (Nelissen et al., 2011). This model is also unable to account for perspective variant mirror neurons in F5 (Caggiano et al., 2009).

4.2.3 Extending Mirror Neuron System

The Mirror neuron system 2 (MNS2) model of Bonaiuto et al. (2007) is an extension of MNS1 both in terms of architecture and representations. The central part of the model is a two-layer recurrent network (see Sec. 3.1.3) with modified Jordan architecture that slightly resembles the RNNPB model (see Sec. 4.2.9), illustrated in Fig. 4.5. The model learns to classify three types of hand trajectories that represent three different object grasps. Inputs to the model are the hand state representations from previous MNS1 model (see Sec. 4.2.2). The network is trained by BPTT (Rumelhart et al., 1986) using “self-observation” signals. In this case the categorical information about the type of the grasp is pre-given to the network as the object’s affordances. After training the model predicts the unfolding trajectory for a target object as well

as it activates the proper code in $F5_{\text{mir}}$ representing the action-recognition in mirror neurons. The output of the network is of two types (see Fig. 4.5). The *recurrent output layer* represents hand state information used during training and serves as a context for the network in the next step. The *external output layer* provides categorical information about the observed movement in one-hot encoding.

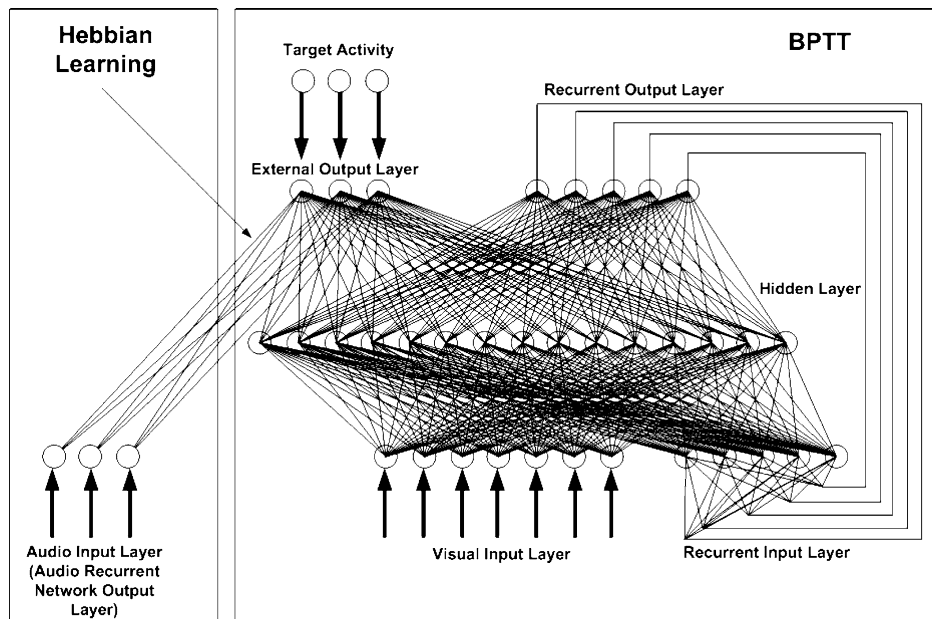


Figure 4.5: Schematic depiction of MNS2 model (Arbib, 2005).

Using the recurrent loop, MNS2 gained also an ability to predict correct hand trajectories even when the hand becomes temporarily invisible. Bonaiuto et al. (2007) refer to this as to the working memory. This capacity of the model reflects on the findings of Umiltà et al. (2001) (see Sec. 2.2.1). Another addition to the model is the auditory subsystem. It is the same Jordan-type network as the primary visuomotor system. The output neurons of the auditory subsystem project directly to the external output layer. These projections are trained using Hebbian learning and provide the association between actions and their characteristic noises (such as peanut breaking or paper ripping) on the basis of study by Köhler et al. (2002) (see Sec. 2.2.1). The result of this interconnection is that the mirror neurons respond to vi-

sual as well as to solely to specific auditory stimuli. Bonaiuto et al. (2007) consider the auditory-based representations inherently actor-invariant.

In the following work, Bonaiuto and Arbib (2010) came with an idea that the mirror neuron system might mediate another function, which they called "what did I just do". The motivation came from experiments with a cat re-learning to grasp food in a new way due to spinal lesions (Alstermark et al., 1981). To simulate this, Bonaiuto and Arbib (2010) designed an extended architecture using the MNS2 model. Interestingly, they switched to RL paradigm. Concretely, the model combines RL, action affordances, and competitive queuing in a mechanism called *augmented competitive queuing (ACQ)*. The model uses representations of the external state (distances between important objects) and the internal state (hunger) as inputs into two separate systems: the actor system and the mirror system. The actor system chooses an action based on its desirability, computed from the internal state, and its executability, computed from the external state. The model assumes a limited repertoire of meaningful actions (such as grasp-with-paw, grasp-with-jaws, rake, raise/lower neck etc.) and a variable repertoire of meaningless, "irrelevant" actions, to enlarge the search space for useful actions after the lesion. The final action is chosen according to its highest priority, computed as a product of its executability and desirability.

RL is applied after the recognition of action for updating both executability and desirability of the action. If the action was recognized by $F5_{\text{mir}}$, the reinforcement signal for the executability is positive (action is reinforced). If the activation in $F5_{\text{mir}}$ was below a minimal level, the reinforcement signal is negative. The executability is reinforced by the perceived ability to perform the recognized or the intended action. RL signal for desirability is formed by the presence of food in the animal's mouth and processed by the adaptive critic module that compares the predicted desirability of recognized action with the primary reinforcement. Effective RL signal is computed as a sum of primary reinforcement and the error in current prediction of desirability. The RL signal represents the difference between discounted predicted desirability of the current action and the predicted desirability of the previous action. Thus, if the next action is more desirable than the previous one, the RL sig-

nal for recognized action is slightly positive, as it brings the cat closer to the goal. If the food ends in the cat's mouth, the RL signal is maximum and the desirability of the performed action is strongly increased. The desirability of the successful end of an action will later serve to reinforce previous actions that lead to it.

Once the agent learns to perform all actions well, the lesion is induced in form of noise and inaccuracies at the input and the agent's motor schemas, and the desirability and the executability of actions have to reorganize. If the MNS is disabled, the reinforcement is markedly slower because the agent has no clue as to what action was actually performed. In this case only the primary reinforcement (food) is available. On the other hand, involving a functional MNS is shown to mediate rapid reorganization of successful behavior to compensate for the lesion.

Evaluation

MNS2 model shares with MNS1 two assumptions: (1) the use of an object-centered frame of reference, and (2) the setup in which $F5_{\text{can}}$ mirroring activity precedes activity in $F5_{\text{mir}}$ neurons, so it can serve as target activity for learning. From the empirical point of view, canonical neurons are not supposed to fire when the action of another agent is observed. On the other hand, they fire consistently when the graspable object is observed and when it is grasped by the monkey. Using this property to train the $F5_{\text{mir}}$ activity is a nice computational mechanism. On the other hand, it does not explain the emergence of the mirroring function as such, since it presupposes $F5_{\text{can}}$ activity. The extensions of MNS2 are quite fascinating and the model is very complex. In line with ecological validity requirements I appreciate experiments with ecologically valid RL. The capability of the mirror neurons in MNS2 extension to re-organize the animal's behavior fits the prediction. An interesting question would be the scalability of this model. The problem of hand-state preprocessing and absolute perspective invariance presupposed by the model remains to cloud the plausibility of the model. However, such phenomena might be accounted for using more artificial neural network modules, and special architecture design.

4.2.4 Mental State Inference

Mental State Inference (MSI) model developed Oztop et al. (2005) tries to encompass the capacity to “put oneself into the shoes of the observed agent”, which helps us understand his actions. Similarly to previous models of Oztop and Arbib, the MSI model is applied to the context of grasping actions. In the first step, the model learns to grasp an object using visual feedback, exploiting both inverse and forward models. As firstly mentioned in Sec. 2.1.2, a forward model (Wolpert and Kawato, 1998; Kawato, 1999) is a brain mechanism that generates predictions about the next possible action. It is coupled with an inverse model that searches for actions that could possibly lead to the observed situation. Forward and inverse models were proposed as one of the solutions for the problem of agency arising in the common coding theory (Sec. 2.1.2). Similarly, in the MSI model, the forward model is used for simulating the mental state of the observed actor and decoupling the mental movement from the actual movement.

The schema of the model with both actor and observer is displayed in Fig. 4.6. In case of voluntary action, the parietal cortex (possibly area 7b) extracts the control variables from visual information. These control variables refer to aspects relevant for the execution of a particular action. For instance, a control variable for reaching is the distance between the hand and the object. The premotor cortex receives this information and computes the motor signals (inverse model) to match the parietal cortex output to the desired neural code relayed by the prefrontal cortex. The premotor cortex also includes the forward model that mediates inverse model learning for establishing a feedforward control strategy. Once the agent’s own sensorimotor feedback loop is established, it is extended to be used for observing and understanding grasping behavior. Oztop et al. (2005) call this a “mental simulation loop”, built around a forward model, which in turn is used by a “mental state inference loop” to estimate the goal of the observed agent. During action observation process, the premotor cortex of the observer is inhibited to avoid any overt movement.

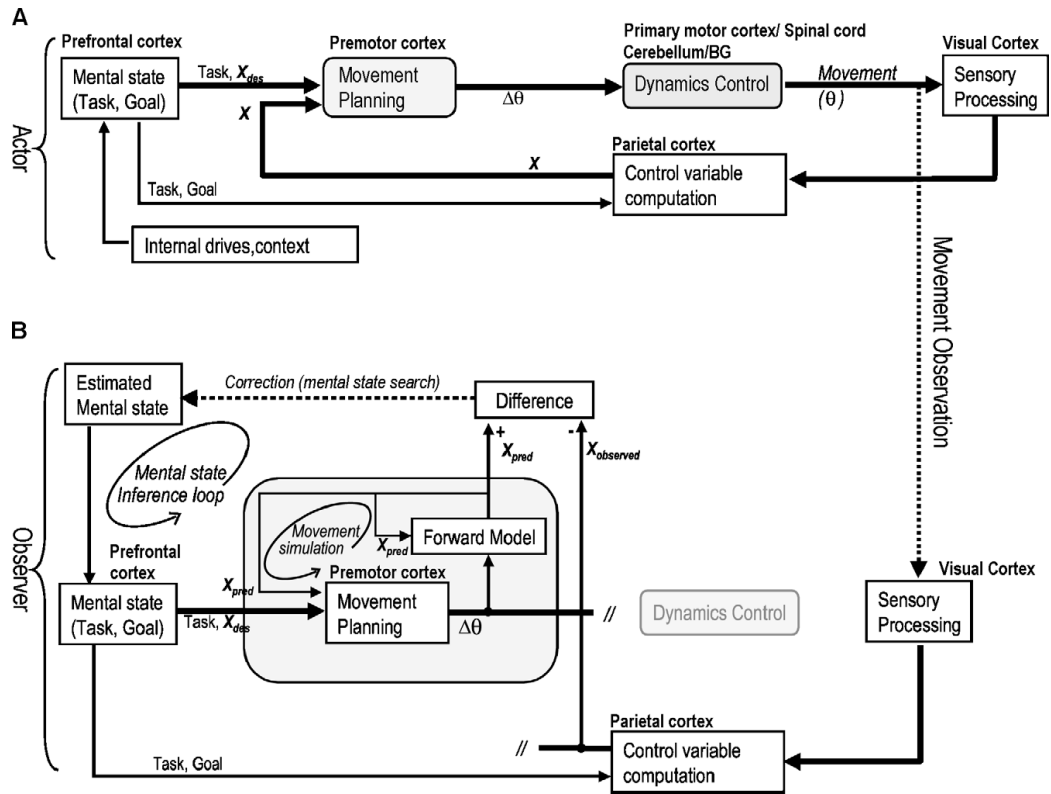


Figure 4.6: Schematic depiction of MSI model (Oztop et al., 2005).

The mental state inference loop is trained using a gradient-based method that minimizes the error between the predicted sensory outcome (generated by the forward model) and the observed sensory outcome. If the error is minimized, the observer can infer actor's mental state and hence accurately predict the trajectory of the observed action. The crucial assumption, necessary for functioning of the model, is the use of object-centered frame of reference for both executed and observed grasping movements (extracted control variables are invariant under translation and rotation).

Evaluation

The MSI model, in line with the direct matching hypothesis, shows that action understanding requires motor simulation of the observed action. This model, as well as its predecessors, does not reflect to the problem of perspective, but expects that mirror neurons receive perspective-invariant infor-

mation from the STS. However, recent evidence points out that perspective-variant neurons are found in F5, and that perspective-variant sites of STS (STSp) are connected with F5 (F5c) through PF, and that perspective-invariant sites of STS (STSa) are also connected with F5 (F5a), but through AIP area. Recently, Frischen et al. (2009) have shown that change in the current dominant perspective taken by an agent can occur subconsciously. They suggest that witnessing an action of a conspecific leads the observer to simulate the allocentric selective attention mechanisms such that they effectively perceive their surroundings from the other person's perspective. The question is whether MNS-based action understanding also requires a change of perspective or whether this ability can be achieved without it.

4.2.5 Hebbian Account

Keyesers and Perrett (2004) provide a conceptual model on Hebbian-learning-based emergence of mirror neurons. They focus on STS–PF–F5_{mir} network in which the PF and F5_{mir} are supposed to develop mirroring properties as a direct consequence of the anatomical connections between STS, PF and F5_{mir}. The connections between these areas become first associated during self-observation while executing an action (also theorized by Heyes, 2010), so the activations in STS neurons responding to the sight of this action (e.g. precision grasp) overlap in time with activity in the PF and F5 neurons that fire when the agent performs that action. This time relatedness poses a prerequisite for Hebbian learning. The same logic applies to learning to understand actions of others. Thanks to invariant properties of (object-centered) STS neurons, the observation of someone else performing a similar action will activate the same neurons in PF and F5. This way, mirror properties are supposed to emerge. The authors also explain how in a similar way mirroring of actions that we cannot see (e.g. one's own mouth movements) can emerge, and how many other forms of social learning can be conceived with Hebbian learning. In addition, this conceptual model includes the explanation for mechanisms that could inhibit STS activations during action execution.

Evaluation

The Hebbian view is a quite clear and relatively simple account of the emergence of mirror neurons. From the computational perspective it may be thought of as a higher-level account because it abstracts away from the sequential nature of sensory and motor signals. Hence, it associates activations at two different locations. Despite its psychological plausibility it is not clear whether this account encompasses the right level of association at the neural level. It is possible, that whole movement sequences have to be associated. Sequence association has been approached in the two models described here, specifically in the MSI (Sec. 4.2.4) and RNNPB (Sec. 4.2.9) models, using different learning paradigms. In addition, in the Hebbian view it is not clear, why only a subset of motor neurons gain mirror properties. Basically, neurons that respond also to action observation form about 20% of area F5. Probably, there are some (lateral) competitive mechanisms coupled with self-organization that need to be included to provide an explanation of this phenomenon. A theoretical option is that some neurons could be pre-determined to become mirror neurons. However, this rather strong nativist assumption has been challenged by many views (see Sec. 2.2.8).

4.2.6 Higher-order Hopfield network

Chaminade et al. (2008) designed and implemented an alternative recurrent neural network model (HHOC) that learns visuomotor associations using a robotic hand. The aim of this model is to account for the capacity to imitate originating from self-observation. The HHOC model is an extension of Hopfield's associative memory, a one-layer recurrent network with full connectivity and symmetric weights. The Hopfield network is known to be able to store patterns and, as a crucial property, also retrieve the pattern based on a partial representation or noisy data (Hopfield, 1982). In this case it was used to associate the visual signals presented on the retina and the motor signals that were generated by the agent during the "motor babbling" stage of its development. Compared to the Hebbian account (4.2.5), the pattern association of signals in the HHOC model is achieved in a more complex way. In

HHOC each visuomotor pattern is stored as an attractor with its surrounding basin of attraction that enables convergence to the attractor (representing a pattern-recall process). Specifically, the artificial input patterns consist of all possible hand postures with four fingers (all but thumb) up or down, and the expected retinal images for the posture. Interestingly, to compensate for the high overlap between (visuomotor) patterns to be stored, the HHOC uses the second-order units (rather than standard units) whose synaptic weight changes are proportional to a product of three unit activations (rather than two).

The HHOC model was shown to be quite robust to noise, to generalize across patterns by inducing a correct motor pattern when cued by a novel visual representation, and to generalize across agents by inducing an expected motor pattern when processing a visual input generated by a different artificial/human hand. The authors interpret the latter ability as imitation, because it leads to correctly evoked motor pattern. Since the visual and motor patterns are mediated by induced attractors, the availability of both pattern types can be safely assumed. Testing for generalization to other agents resides in using altered visual stimuli (hands) that evoke slightly different retinal images to be processed.

Evaluation

In other models, generalization is related to view-point relative processing in STS yielding object-centered representations. The use of the second-order units in the HHOC model, which leads to higher computational power, solves the modeling problem, and has also been argued to be biologically plausible (Mel and Koch, 1990). The HHOC model has not been linked to any anatomical areas but, in the context given above, one can assume that it is meant to link STS with $F5_{\text{mir}}$ in the form of long-range attractors, ignoring the mediating stages of processing performed by the parietal areas.

4.2.7 Tessitore's model

A slightly different theoretical approach to mirror neuron system function and modeling was adopted by Tessitore et al. (2010), who emphasize the

bidirectional nature of the flow of information between visual and motor areas. Note that in most of previous computational attempts the information flow was directed from visual areas to motor areas where triggering mirror response. The main assumption of Tessitore et al. (2010) is that mirror neurons facilitate action recognition and control processes, since they provide a simplified motor representation that narrows down a wide search space of the visual input. The model of Tessitore et al. represents a mapping function from visual representation (an image of the hand in grasping action) to motor representation (created with a special recording glove). The core of the model is preprocessing of the data using PCA and a Mixture Density Network (MDN) trained using standard BP algorithm.

The model is based on empirical findings which suggest that any grasping action can be expressed using a small set of hand-joint configuration parameters (Mason et al., 2001). Tessitore et al. (2010) encode these descriptions of hand postures using the so-called *action subspaces* which stand for principal components of hand postures, also called eigenpostures. In this way each hand configuration can be encoded as a weighted linear combination of eigenpostures. The goal of the model is to encompass the functional mapping between visual representation of a grasp and its motor representation using action subspaces. For this mapping, two MDNs were used, with standard BP algorithm. Results from experiments confirmed that motor information indeed simplifies the processing of visual inputs and improves visual recognition. In sum, Tessitore and colleagues model the firing of the mirror neurons during action execution and action observation as the action subspace selection process (i.e. applying the mapping trained function) stands here for

Another interesting property of this model is that it directly solves the problem of translation of the perspective, i.e. from the observer to oneself. It contrasts with most of the classic MNS models, such as those mentioned above, which do not address this property, but rather assume visual preprocessing in STS.

Evaluation

Tessitore et al. (2010) compare their model with the MOSAIC and RN-NPB models (described further in this text), that also allow for iterative interactions between sensory processing and mirror activity. However, in these models the iterative interaction occurs via coupled forward-inverse models and the nature of mirror neurons is somewhat different. Unlike the other two models, Tessitore’s model encodes the motor information as action classes, rather than action sequences, without requiring a precise reference to action kinematic parameters. Thus, in this model the motor information involved in action understanding operates on a higher level (the categorical representation and simplified features of the movement). Tessitore and colleagues support this aspect of their model using empirical evidence by Craighero et al. (2002), who have shown that mirror neurons respond to the kinematic characteristics of executed/observed actions in quite unselective manner, i.e. mirror neurons generally react to various movement setups. This level of modeling precision resembles the Hebbian model (Sec. 4.2.5).

4.2.8 Knott’s model

As a part of his theory about a tight link between sensorimotor cognition and natural language syntax, Knott (2012) proposes a quite elaborated (albeit not yet implemented) model of the mirror neuron circuit (section 2.7.5 in his book). This model (Fig. 4.7) combines the ideas of the Hebbian account (Iacoboni et al., 2001; Keysers and Perrett, 2004) and the forward modeling (Wolpert et al., 2003; Oztop and Arbib, 2002; Oztop et al., 2005). The model focuses on the STS–PF–F5_{mir} circuit whose parts are assumed to learn to bidirectionally associate visual representations in STS with motor representations in F5. As a core assumption, shared with earlier models, STS can form action representations, that are perspective invariant with respect to the agent. As an important part, the inclusion of the forward model aims at capturing the temporal nature of the signals to be associated as well as simplifying the matching to be learned between F5 and STS. The forward model $F5_{\text{mir}} \rightarrow PF \rightarrow STS$ converts motor signals into anticipated sensory

consequences (in line with Miall, 2003), which may be easier to match with STS signals than the motor representations from which they derive (because they are in the same domain).

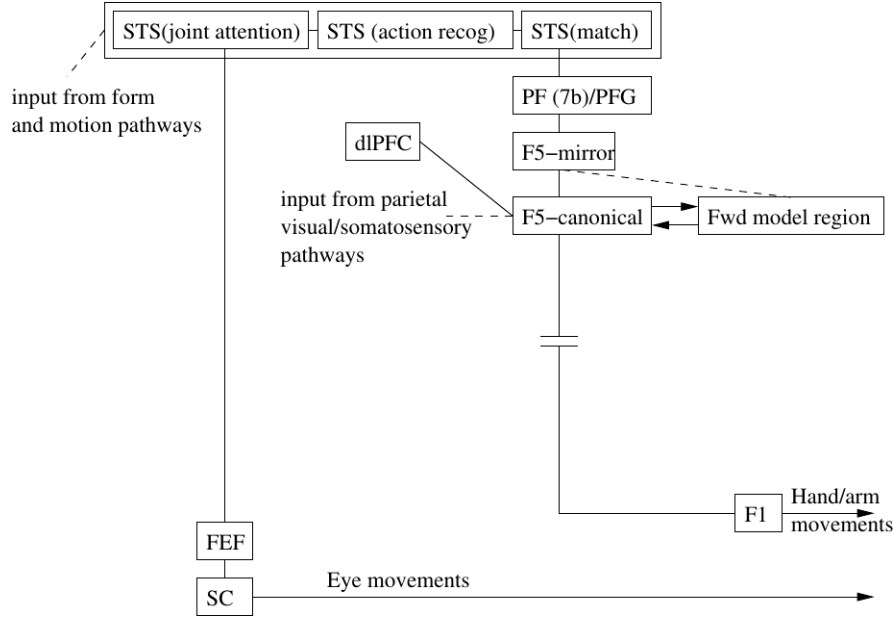


Figure 4.7: Schematic depiction of Knott's (2012) model of MNS.

Knott (2012) points out several significant differences between his model and several earlier models. For instance, the model of Oztop and Arbib (2002) assumes that STS is directly involved during the execution of reach/grasp actions, while in Knott's model it is not. However, the STS 'match' region, discovered by Iacoboni et al. (2001) to be involved in both execution and observation, is active in Knott's model during action recognition, as it receives signals from the motion recognition pathway, namely MT and MST.

Evaluation

Comparing the Knott's model to a typical MNS model, e.g. the MSI (Oztop et al., 2005) there is a difference in two main aspects. First, representations of the observed agent's hand are primarily computed in the parietal area, and then are sent to STS for matching with hypothesized representations, which reduces the role of STS. In Knott's model, STS receives inputs from the grasp/reach pathway (MT and MST) but also from the form

classification pathway, making STS more autonomous in generating visual representations of actions (analogical to object classification pathway). The second difference is that in the MSI model, the learned association between F5 and STS only runs in one direction, from F5 to STS, even during action recognition. This, in fact, reduces the role of STS in action recognition as well. Knott's model, on the other hand, assumes that STS generates its own invariant visual representations that can trigger activity in F5, hence making the links bidirectional. This is in line with the evidence on invariant representations in STS as well as with the model of Tessitore et al. (2010) discussed in Sec. 4.2.9. Knott's model is definitely worth implementing (using neural networks) such that its functionality and predictions could be better appreciated.

4.2.9 RNNPB model

Adopting the dynamic systems perspective, the model of Tani et al. (2004) represents a shift towards recurrent architectures, which allow learning of sequences. The recurrent neural network with parametric biases (RNNPB) was designed to allow learning, imitation and autonomous sequence generation. A generic architecture is depicted in Fig. 4.8, computational detail on architecture and training of RNNPBs are provided in Sec. 3.1.4. The most intriguing novel feature of RNNPB are the so-called parametric biases (PB), which enable the network to recognize and categorize various spatiotemporal patterns and thus modulate its own dynamic function. PB vectors emerge on the basis of self-organization, therefore they are considered self-determined by the network in an unsupervised manner. Apart from PB vectors, the network is trained using BPTT method.

In a standard setting, the network works as a forward model (Wolpert et al., 2003), generating predictions about the next state of the world from the current state when performing a certain action. The values of PB nodes during concrete behaviors can be stored and subsequently (manually) fed as an input to the network to produce the learned behavior. Thus, the network can initiate appropriate action without a need for sensory stimulation

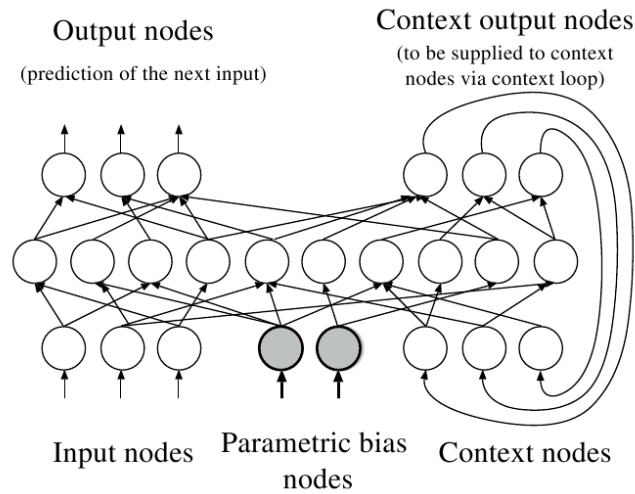


Figure 4.8: A scheme with two generic RNNPB nets connected through their PB vectors using PB-binding (Sugita and Tani, 2005).

(target). Regarding this capacity, Tani et al. (2004) propose that the parametric biases share properties and play similar role as the mirror neurons, because their activation both characterizes and triggers various sensorimotor sequences. Since the input layer can consist of sensory inputs of both visual and motor type, the PB vectors can bind information and produce estimates in both modalities simultaneously.

Tani and colleagues tested the RNNPB architecture in three robot experiments, each with physically embodied robots. In the first experiment, the Sony humanoid QRIO robot learns to observe and simultaneously imitate hand gestures of the human demonstrator. The state space analysis of the memory patterns in the RNNPB has shown that each pattern learned by the network forms a distinct cycling attractor. The authors concluded that the attractor dynamics implemented with the bifurcation of parameters makes the system manipulable by the users and robust against possible perturbations. In the second experiment, they used a different setup, a robotic arm, which was able not only to learn simple movement patterns, but also to differentiate classes of patterns. Concretely they were end-point and circular movements, which naturally have a quite different dynamics. Results of this experiment suggested that the RNNPB can produce quite nonlinear map-

ping between parametric bias and the generated behaviors and thus embed different types of attractor dynamics. In the third example, an experiment with linguistic behavior binding is introduced which is similar (of not the same) as the model and experiment described in Sec. 4.3.2.

Evaluation

An interesting aspect of the PB vectors is their proposed analogy to mirror neurons. However, the relationship between mirror neurons in F5 and PB vectors is more of a metaphorical nature. Although the PB vectors both characterize and trigger various sensorimotor sequences, like mirror neurons, their nature is rather multimodal than motor (mirror neurons are motor neurons with visual properties). Since the input layer can consist of sensory inputs of various types, here of visual and motor information, the PB vectors can bind information and produce estimates in any modalities depending on the interpretation of the input layer. From my point of view, PB vectors are more like pointers to the brain areas, which can trigger actions. Along with recent evidence on mirror neurons in humans (Mukamel et al., 2010) (see Sec. 2.2.2), PB vectors might be considered computational counterparts of mirror neurons in hippocampus or neighboring structures, which serve as pointers to the cortical brain memory areas (O'Reilly and Munakata, 2000).

4.2.10 SRNPB model

Following the work of Tani and Sugita, Zhong et al. (2011) proposed the SRNPB model, i.e. a SRN (Elman, 1990) with parametric biases (Fig. 4.9). The goal of their model was to gain an ability to predict the following trajectory of a walking humanoid robot based on its previous trajectory. Since the PB have basically the same properties as in RNPPB, they can be seen as categorical representations of trajectory types, but also interpreted as mirror neurons for walking actions. Results from experiments with SRNPB have shown that it is not only able to successfully predict and recognize sequences, but also generalize over new sequences. Since the model successfully implements a forward model, it can also be used for generating next step to be executed within a pre-learned action.

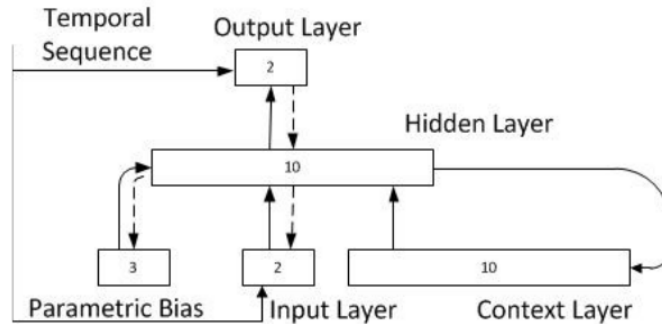


Figure 4.9: Schematic depiction of the SRNPB model (Zhong et al., 2011).

Evaluation

In this case the same questions can be raised as in the RNPPB model. The PB vectors are not modality specific, therefore they cannot be definitely identified with mirror neurons. On the other hand, with respect to cognitive robotics paradigm PB-based networks are a very powerful tool to build multimodal representations which can elicit explicit behavior as well as internal simulation (forward and inverse model). A crucial property is the ability to encompass sequences as single codes (PBs), because such mechanism is most probably mediated by mirror neurons. In this case, however, I would rather speak of common codes (Sec. 2.1.2) than mirror neurons. Interestingly, the prediction ability of the model can be put in parallel with motor-resonance based theories of understanding (Sec. 2.1.2). Basically, the more familiar the observed action and its actor, the more motor resonance will emerge, and (probably causally) the better the prediction about the observed action would be. It would be interesting to be able to measure motor resonance in SRNPB model, if possible, or to extend the model to account also for this phenomenon.

4.2.11 MOSAIC

The MOSAIC model is a rather sophisticated architecture based on Hidden-Markov models (HMM) that was originally designed for motor control (?). Wolpert et al. (2003) have shown that this model also encompasses mirror

neuron function in action recognition and imitation. MOSAIC is modular and allows a distributed cooperation and competition of the internal models (see the top of Fig. 4.10). The basic functional units of MOSAIC are multiple competing predictor–controller pairs (forward–inverse models). Competition drives the process of selecting controllers with better predicting forward models to be more influential to the overall control.

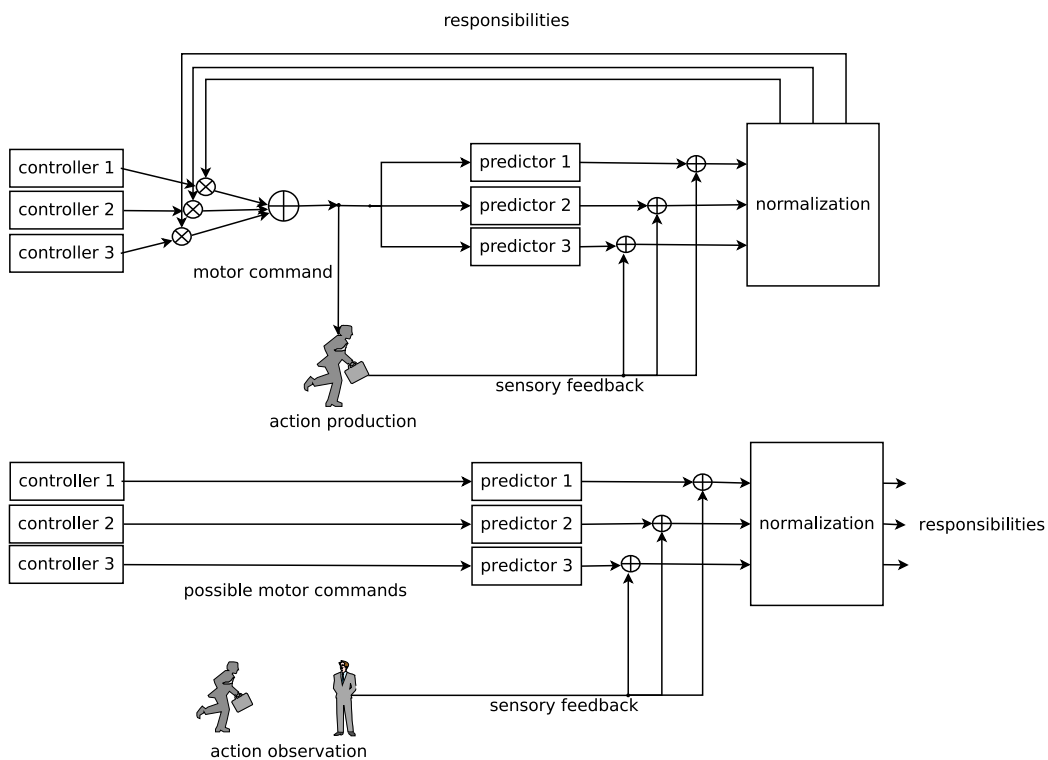


Figure 4.10: MOSAIC model in action control mode (upper part) and in action observation mode (lower part) (Farkaš et al., 2011a).

The use of MOSAIC for mirroring in case of action recognition or imitation requires three stages (Wolpert et al., 2003): First, the visual information of the actor’s movement must be converted into a format that can be used as inputs to the motor system of the observer. This conversion requires that the visual processing system extracts variables related to the agent’s state (e.g. joint angles). These are fed to the observer’s MOSAIC as the “desired state” of the actor. The second stage is that each controller generates the motor command, which represents the desired trajectory obtained from the

observation. In this “observation mode”, no movement generation occurs but the outputs of the controllers are used as inputs to the predictors paired with the controllers. Hence, the outputs of the forward predictions represent possible internal states of the observer. These predictions can then be compared with the actor’s actual next state. The difference – prediction error – indicates, which of the controller modules of the imitator must be active to generate the observed movement. The outcome of this stage is a symbolic sequence composed of indexes of the controllers selected (one at a time) during action observation. In the third stage, this sequence can either be compared in memory (action recognition) or used for repeating the action (imitation). The output of predictors might be considered analogous to the mirror neuron activity.

Evaluation

The distribution of control in the MOSAIC model is an interesting idea that allows efficient control. It is possible that such multiple controllers specialized for various behaviors (e.g. grasping) are also implemented in the brain. However, the need for explicit pre-specification of the number of controllers seems less plausible. The controllers and their division of labor might rather be formed during experience.

4.2.12 Wiedermann’s finite cognitive agents

Unlike the previously described models, consisting of artificial neural network, Wiedermann (2003) proposes a model based on finite state machines. He defines a finite cognitive agent (FCA), which consists of perceptual-motor units (PMU) representing the agent’s sensors and effectors and a finite-state transducer (FST), representing the mind/brain. Finite cognitive agents can interact with their environment through their PMU and select actions, reason, etc. using FST. Given an FCA A with k PMUs, a finite set Q of the agent’s states, a finite set M of agent’s actions, and set P of proprioceptive information, the activity of an FCA (Fig. 4.11) can be formalized using a

transition function δ :

$$\delta : Q \times (M \times S \times P)^k \rightarrow Q \times M^k. \quad (4.1)$$

The model works in the so-called *standard mode* as a control model for agent's behavior, as well as in an *observation mode* when the agent observes another agent's behavior and the appropriate motor plan is retrieved (mirror activity). The model shares this property with various other MNS models. It is very close to the common coding theory (see Sec. 2.1.2), since it can retrieve also information from sensory modalities typically accompanying the particular action. The FST of the FCA also works as a forward-model endowing the agent with motor imagery.

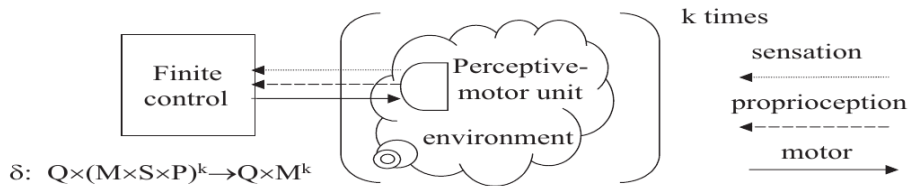


Figure 4.11: Schematic depiction of a finite cognitive agent (Wiedermann, 2003).

Later, Wiedermann (2009) refined his model in order to capture abstract semantics and high-level cognitive capacities as acquisition, processing and exploiting of semantic knowledge. In this model the agent's control unit consists of two internal world models, which complement each other in their functioning. The *mirror net* consists of multi-modal representations of actions (again resembling the common coding) that are connected to agent's sensors and effectors (S-M units). The *control unit* is divided into two sub-units representing embodied concepts and abstract concepts, derived from embodied ones (Fig. 4.12). Unlike the previous model, the control unit here is based on neural networks. The author emphasizes a possibility of the role of mirror neurons in higher cognition and in language evolution proposed mainly by Arbib (see Sec. 2.2.7).

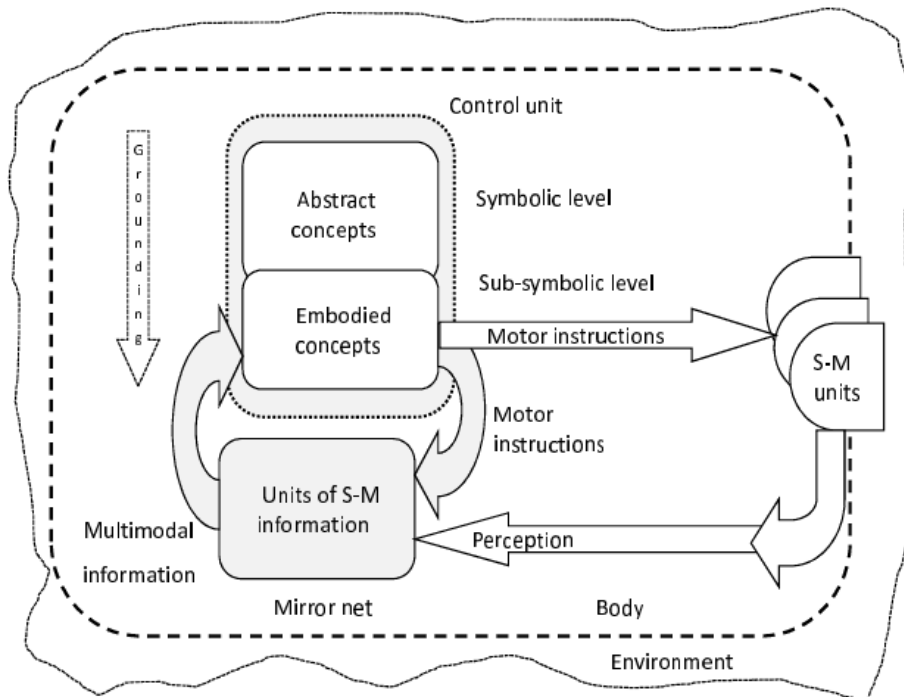


Figure 4.12: Schematic depiction of an extended finite cognitive agent (Wierdmann, 2009).

Evaluation

Although the model look quite interesting, its function has not been yet empirically tested, nor has the model been implemented. As a theoretical construct it is a fine illustration of how both symbolic and subsymbolic processes might work together in the human brain/mind. However, to implement such non-trivial model would require adopting a suitable framework, such as cognitive robotics. An architecture based on mirror neurons encompassing the common coding might be very beneficial in driving the robot's behavior. Such mechanism can provide the robot with the ability to categorize (or name) the observed behavior, and, based on the robot's experiment, to make estimates about the observed actions. Apart from grounding the model in a robotic platform there are many uncertain details, such as a description of actual neural network design and learning to encompass the control unit, or, in case of the FST the problem with the determinism of such machine (the

model assumes that there is only a finite number of actions and resulting perceptual states).

4.2.13 Bayesian approach

Kilner et al. (2007) have adopted a Bayesian perspective to modeling the MNS. They identify a the role of MNS with the agent's ability to infer intentions from the observed agent's movement. They focus on the problem of inferring the cause of an observed action and suggest that the problem could be solved by the MNS using predictive coding on the basis of a statistical approach known as empirical Bayesian inference. This means that the most likely cause of an observed movement can be inferred when minimizing the prediction error during action observation by comparing the predicted the kinematics on the basis of agent's own action system and the observed kinematics.

The computational mechanisms for this generative approach (as opposed to one-way, recognition approaches) are proposed to exist in two directions. The posterior-frontal pathway $STS \rightarrow PF \rightarrow F5_{\text{mir}}$ implements the prediction error of the motor action, and the frontal-posterior pathway $STS \leftarrow PF \leftarrow F5_{\text{mir}}$ implements the generative model that computes the sensory consequences of the performed action (forward model). The authors suggest that the forward model be an integral part of the motor function, since the same model becomes exploited both for action execution and action observation (similarly to the MSI model).

Evaluation

Bayesian framework is well principled and exploits bidirectional connectivity. It is consistent with the theory that the brain performs all kinds of predictions (expectations) at various levels of its organization (Friston, 2003). On the other hand, it could be argued that probabilistic modeling is a higher-level account that does not provide neurally-inspired learning mechanisms, characteristic of connectionist approaches.

4.3 Grounding meaning in action

In this section I will describe selected models of grounding meaning in action, most of them from the field of cognitive robotics (Sec. 4.1), implemented in physically embodied robots (or their simulated counterparts). For each model I will provide its description, properties, empirical motivation, and a critical evaluation. When selecting model for this chapter I focused on the process in which meanings are acquired and grounded in the agent's perception and action. Agents endowed with such models usually learn to recognize and name the actions they currently make or are able to initiate actions on verbal command. The models in this chapter are partially ordered based on the complexity of language they model, and also on the basis of the computational paradigm used (e.g. ANN versus formal logic).

Most of the models I describe are based on artificial neural networks, since they are the primary topic of this thesis. This chapter mentions various types of architectures and learning, most of which can be found in Chap. 3. The models I present implement the common coding theory, the involvement of the motor resonance in language learning and comprehension, and also functionally mimic mirror neurons. For comparison, I include models based on other mechanisms, like symbol systems, at the beginning and at the end of the chapter. After summarizing the models that influenced my work described in Chap. 5.

4.3.1 Direct grounding of language in action

Marocco et al. (2010) recently proposed a model of direct grounding of language in action (motor control). This model was designed and tested using the simulated iCub robot (see Sec. 4.1.1). The control architecture for iCub consists of a recurrent neural network.

In their current study, Marocco and colleagues used only a small subset of all iCub's DoFs, specifically one joint on the shoulder to push objects on a table and two joints in the neck, to express the position of the object in the visual field relative to the robot, providing the visual input. Next,

the encoder value of the shoulder joint is used as a proprioceptive feedback. The set of inputs for the model also contains a coarse information about the shape of the manipulated object called the *roundness* parameter and, of course, linguistic input. The sensorimotor state of the robot is updated every 500 ms. When the robot receives a target joint angle as input, it automatically generates a movement corresponding to the target angle using a preprogrammed proportional-integral-derivative (PID) controller. There are three different objects in the robot’s environment, which will always produce the same “response” to the pushing action. There is a sphere that rolls away, a box that slides, and a cylinder, which is fixed to the table, so it is unmovable.

The neural architecture is a fully connected recurrent neural network with 10 hidden, 8 input and 8 output units. It works as a standard forward model (Wolpert et al., 2003) predicting the next state of sensors, actuators and linguistic input. The input layer consists of three joint encoders, one neuron for tactile input (active when the hand touches the object), one neuron for roundness and three linguistic units locally encoding three words. Interestingly, the three words could be interpreted not only as names for actions, but rather as names of the properties of the objects the iCub encounters (e.g. rolling, sliding, unmovable). The activation value of linguistic units can vary between the training phase and the testing phase, respectively.

The network was trained using BPTT (Sec. 3.1.3). During the training phase the robot received sequences of 30 activation patterns each representing 15 seconds of the robot’s activity plus the linguistic input provided at the beginning, and computed by the anticipation of the network for the rest of the sequence. To facilitate the training process the iCub learns in a *closed-loop* mode. Subsequently, the model was tested in the same conditions as trained, but in an *open-loop* and in online manner. Results of various experiments showed that the robot was able to name its experiences with objects correctly. Interestingly, the roundness information, computed from the visual input allowed the agent to recognize and name the object (action) correctly, even before it came to contact with it.

Evaluation

An important aspect of this model, similarly to the works of Sugita and Tani (2005), is that language and action are fully interconnected. Execution of an action automatically activates its linguistic description and vice versa. Another feature which is very important for modeling of the embodied cognition is the interconnection of the robot’s cognition and its body. Meanings in this model are specific to the robot. They do not constitute any feature lists or other designer’s expectations about the learning and recognition process. Words are encoded in terms of the robot–object interactions. Regarding the roundness parameter, which partially suffices for activation of the correct word, this effect is a result of sensorimotor binding as well and could not emerge separately.

When evaluating the model’s biological plausibility, a small drawback can be found – the BPTT training method. It is well known that standard BP is not biologically plausible. It is not natural for the action to be literally forced to the agent, as if there was some hand dragging the child’s hand to learn to reach for an object. A more biologically plausible method is RL described in Sec. 3.1.2 with emphasis on continuous environments, which are typical for robots.

Using RL the network gets information about its performance in form of either reward or penalty (or both). An important difference between RL and BP is that the information about the correctness of the output is given to the network as a whole, not to separate units. This type of learning is from the “network’s view” quite difficult and might take a much longer time to converge to a good solution. Combined with the recurrent network architecture the task might seem impossible to accomplish. However, Tikhanoff et al. (2011) successfully implemented associative RL (Barto and Jordan, 1987) in an RNN.

Another small objection against this model, and also against few further mentioned models, could be made about the inseparability of naming and execution of actions. In order to name or recognize the action, the iCub has to do it as well.

4.3.2 RNNPB and language

One of the most intriguing of recent connectionist models of embodied language acquisition is the RNNPB based architecture by Sugita and Tani (2005), which was used in various setups where embodied robots learned repertoires of actions, and/or action names. Computational detail of this model can be found in Sec. 3.1.4. Using RNNPB architecture, Sugita and Tani provide a novel scheme for learning, in which no representations, symbolic nor structural, are implemented in the system. In comparison with various previous models like the Bailey's model described in Sec. 4.3.7, the language is directly grounded in the motor control component. There is no heavy preprogramming nor designing needed.

Sugita and Tani (2005) used a mobile robot, equipped with two wheels, one-joint arm and color vision, which operates in an environment (micro-world). In the robot's environment there are three colored objects placed at three different positions (on the left, center and on the right) in front of the robot. To simplify behavioral learning and account for compositionality, the position of the objects remains always (nearly) the same, for instance the red object is always on the right. The robot is capable of three actions: pointing, pushing and hitting, so its repertoire will consist of nine behavioral categories, each of which can be labeled by two different sentences, indicating either the location of the object to be manipulated or its color. The synonyms are introduced to observe the relationship between the behavioral similarity and the acquired linguistic semantic structure. If the behavioral structures correlate with the representation of syntactic structure, the system can be considered truly embodied.

To employ the task of simultaneous behavioral and language learning, the model consists of two interconnected parts, linguistic module and behavioral module, each consisting of one RNNPB (Fig. 4.13). These two networks are connected through their PB vectors, which can be fed to each other's input, or, more importantly, be trained to minimize the difference between them, virtually forming one PB vector that represents both the lexical label and the actual behavioral pattern of the concrete action. To bring these

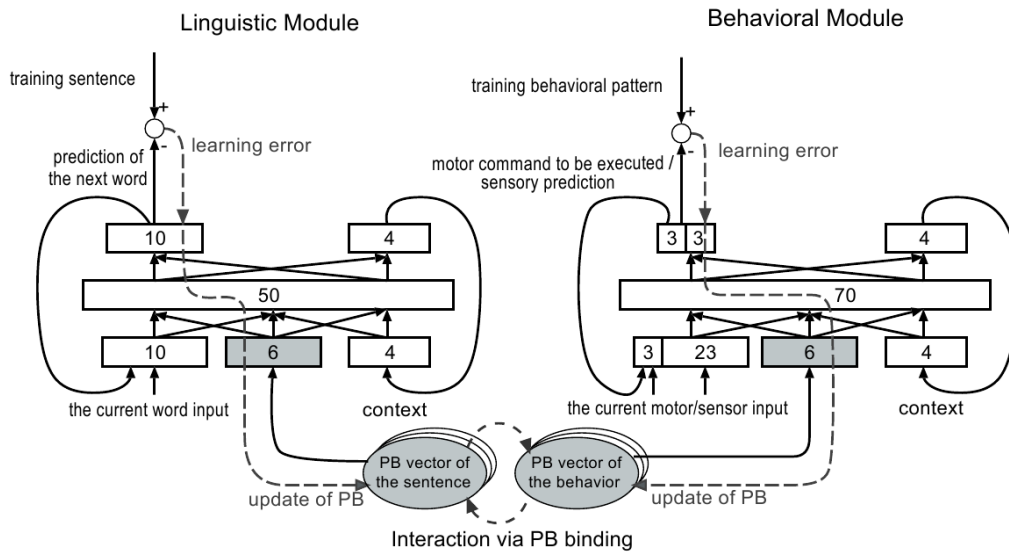


Figure 4.13: Illustration of two interconnected RNNPB (Sugita and Tani, 2005).

representation close to one another, Sugita and Tani (2005) used the so-called PB binding (see Sec. 3.1.4).

The linguistic module has 10 input nodes, 9 of which locally encode words (i.e. three action words, three locations, and three colors) and one stands for a fixed starting symbol denoting the start of the new sequence. The network task is to predict the next word in the sentence. The behavioral module is a bit more complex, with 26-dimensional sensorimotor vector on the input (as well as 6 PB nodes and 4 context nodes) and processes input sequences of various lengths, comprising 15 to 75 steps, unlike the fixed three steps in the linguistic module. Similarly to the other module, this network creates prediction of the next input, which is used as the actual motor command for the robot at the next time step. The two modules work simultaneously during the training phase, but the capabilities are tested separately. First the linguistic module is fed with a test sentence, the recognition takes place and a PB vector is computed. In the next step the obtained PB vector is fed to the behavioral module and the generated behavior is evaluated.

Sugita and Tani tested the model in three different experiments. First, only the linguistic module was trained and tested separately. In this case,

syntax was learned successfully by extracting combinatorial properties of the training set of sentences. An interesting phenomenon emerged during the training, that values propagated from the word input nodes to hidden nodes did not carry information about themselves, but only about their functional class (i.e. verb or noun), pointing to generalization in learning. The second experiment aimed to test the behavioral module and its robustness. The results indicated that the robot can perform the target behavior if the object remains in its sight even when it is slightly shifted from its position. Unlike the first experiment, in purely behavioral task no clear combinatorial structures emerged in the PB space.

Finally, the third experiment employed both modules. In the testing phase, as mentioned before, only the language-to-behavior route was examined. The robot was able to perform correctly for all linguistic inputs, including novel inputs not present in the training set. It was able to successfully generalize to learn the syntactic structure, its mapping to behavioral patterns and correctly react to verbal command. Regarding the PB vectors, it has been observed that even without being bound together during the learning, the PB vectors of novel sentences coincided with proper behavioral sentences.

Acquisition of a compositional representation

In a subsequent work Sugita and Tani (2008) focused on the compositional character of semantics¹. The model was created and embedded in a simulated agent to account for the emergence of reusable units (Skinner, 1957), resulting from generalization of multiple examples and for the usage-based account of language development (Tomasello, 2003). In general, the aim of the model was to explain how conventional symbolic usage-based models can be implemented on a subsymbolic level.

¹The compositionality of language in this sense refers to the human ability to understand sentences from the meanings of its constituents, but also from the way they are put together. This ability then mediates generalization of meaning through compositional semantics, for instance the generalization of roles of nouns and verbs in the sentences.

The experiment used a simulated agent based on a rather simple mobile robot, equipped with a color camera and two wheels, able to perceive objects in its close proximity, turn to them and reach them. The agent received commands in the form of concatenated labels for actions and their parameters. For example, *turnto-blue+18* means that the agent should turn 18 degrees to the blue object. The scene always contained the target object, occasionally accompanied by a dummy object (other than target) serving as a distractor.

An internal structure of the agent consists of a neural network with slightly different architecture than typical RNNPB mentioned above depicted in Fig. 4.14. The model has two parts. The first part, the *base-level network*, which transforms visual inputs into motor outputs (velocities of two wheels), is a conventional four-layered feed-forward network with one modification. Between the last two layers there are second-order connections to the second part of the network called *meta-level network*. This mechanism enables the base network to switch its function accordingly to meta-level network activation in the presence of the same visual input. The meta-level network has two layers, an input layer for PB vectors, a hidden layer with second-order connections to the base network, and a storage space for PB vectors (pseudo-linguistic input). During the course of each action, this network outputs the activation from the corresponding PB vector constantly, ensuring the binding between the action and its vector representation.

Before learning, a set of data of different sparseness was generated by a special teaching program. The network was then trained in a batch manner using the standard BP algorithm. Two capabilities of the agent contributing to its ability to generalize were studied, specifically the ability to transfer the skills into novel environments and the ability to combine learned actions into a novel action. Results of four different experiments showed that the agent successfully acquired both capabilities. Moreover, the principal component analysis (PCA), applied to the representative vectors of the targets (computed from PB vectors), showed that the 36 actions clustered together in the action space, both according to the color of the target and according to the operation applied. Since each of the targets was represented uniquely in the subspace regardless its surrounding context, the authors conclude that

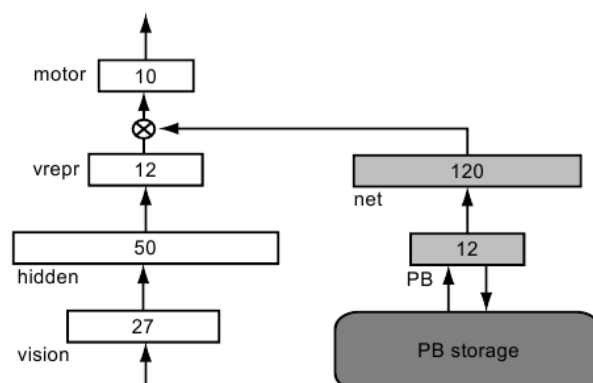


Figure 4.14: RNNPB architecture with a meta-level network (Sugita and Tani, 2005).

the agent was able to generalize from the given input sentences and form reusable concepts.

Evaluation

The power and capacity of the RNNPB model is quite astonishing. It is also in line with the ideomotor and common coding theories, as well as with recent findings on action and language understanding. The PB vectors could be interpreted directly as the common codes for perception and action. In this case, I can again mention the above mentioned objections against the BP-based learning (see Sec. 4.3.6). On the other hand, taking into account that training of RNNPB, especially of two synchronized networks must be quite difficult, the BP approach might be the best way. I find this model very inspiring. However I incline to a simpler mechanism for association of perception and action based on Hebbian learning and a slightly different architecture.

4.3.3 Grounding the meanings in sensorimotor behavior using reinforcement learning

Unlike many other cognitive robotic models, the work of Farkaš et al. (2012) aims at introducing an ecologically plausible mechanism of grounding con-

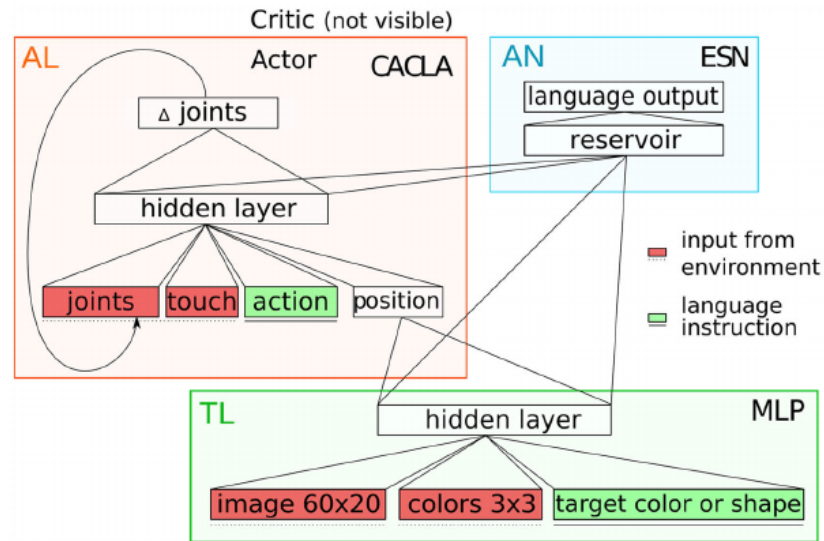


Figure 4.15: The modular architecture of grounded RL model (Farkaš et al., 2012).

cepts in sensorimotor behavior. The proposed control architecture for a simulated iCub robot allows it to learn the meanings of actions (point, touch and push) oriented towards objects in its peripersonal space. The model consists of three neural-network-based modules that are trained in different ways (Fig. 4.15). The first module (target localizer, TL) is a two-layer perceptron trained using standard BP algorithm. Its role is to attend to the target position in the visual scene, given the low-level visual information and the feature-based target information. The second module (action learning, AL) is the most distinguishing part of the whole model. It is an executive module trained using continuous actor-critic RL method (Sec. 3.1.2) to execute actions as sequences, based on a linguistic command. The third (action naming, AN) module consists of an echo-state network. The aim of this module is to provide the linguistic description of the executed actions.

From the experimental results Farkaš et al. (2012) conclude that the model successfully learned generalize in case of novel action-target combinations with randomized initial arm positions. It can also promptly adapt its behavior if the action/target suddenly changes during motor execution. Together with the other two modules, the robot was able to produce and

name all the action-target combinations (also the previously unknown). Interestingly, activations in the hidden layer of the AL module projected on a 2-dimensional SOM tended to form clusters characterizing various actions that the agent is able to produce. This information might be of use in further modeling, for instance when searching for a higher-level representation of the particular action.

A similar model was recently proposed by Tikhanoff et al. (2011), who created a neural network model for reaching and grasping together with the linguistic component. The architecture encompasses a feed-forward network for reaching and a recurrent neural network for grasping. The reaching task, trained by BP, is approached as one-step process, for which the training data is first acquired during the motor babbling stage when the required joint position are stored for various target positions in robot's peripersonal space. In contrast, the grasping behavior is viewed as a sequential task and employs an associative RL algorithm.

Evaluation

One drawback of this model of Farkaš et al. (2012), often present in similar works, is that in order to name the action, the agent has to execute it as well. On the other hand, in a real life one may need to name the action that is observed. To enhance the model with the capability to produce an action name when just observing it, a mirroring mechanism could be introduced. Although the model uses ecologically valid RL method for movement generation, the TL module is still based on standard MLP with BP. To enhance the model's plausibility this module might be exchanged for more suitable and possibly more sophisticated module for object recognition and target localization.

4.3.4 The MirrorBot project

The aim of the MirrorBot project (Wermter and Elshaw, 2003) was to produce a life-like model of perception system in which semantic representations of actions, percepts and concepts emerge in a robot endowed with cortical assemblies and mirror neurons. One step towards such a model was to pro-

vide the robot with self-organizing memory. In context of interconnection between learning of language and actions, Wermter and Elshaw (2003) introduce a model based on self-organizing maps (Sec. 3.2.1). In this model, verbally labeled actions are categorized – clustered according to the body parts they are associated with.

The neural architecture in MirrorBot was strongly influenced by recent neuropsychological findings, particularly by findings of Pulvermüller and colleagues (e.g. Pulvermüller et al., 2001, described in Sec. 2.3). In short, this evidence shows that various different brain regions are somatotopically activated in accordance with the modality which was the verb related to (e.g. kick - foot). The model consists of multiple interconnected SOMs, each representing either a body part or an association area. Its task is to extract and associate semantic features from sensory inputs that represent an action with a representation of a word.

The architecture is schematically depicted in Fig. 4.16. In the bottom left part, the primary sensory area (Body Part SOM) categorizes incoming sensory inputs and sends it to three separate SOMs for three different body parts – head, hand and leg. These maps, as well as the map responsible for processing the linguistic input (Word SOM), project to the association area, which finally generates the proper motor and linguistic behavior in the robot. The association area, inspired by such areas in the brain, is also a SOM. It projects to two separate modules for action and language in a way that was not described in the current paper, so in this case it might be only hypothetical. Most importantly, the system can receive inputs from sensors in two forms: an action representation that stands for visual perception of an action, and a word representation that stands for the linguistic label of an action. Likewise, the association map can produce both the motor output, a copy of the perceived action as well as the linguistic output, the name of the action.

The architecture was trained on sensor readings from Mirror-neuron Robot Agent (MIRA), which is equipped with a microphone, speakers, a pan-tilt camera, IR sensors, a 2-degree gripper with sensors that detects objects, and a wireless connection to a PC in which its control is implemented. MIRA

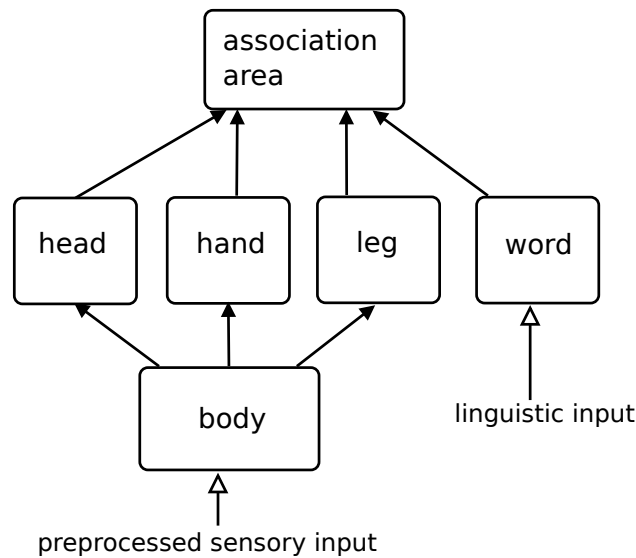


Figure 4.16: Schematic depiction of the modular hierarchical self-organizing memory (Wermter and Elshaw, 2003).

was set up to perform various human-like actions, like turning to the sides (leg actions), moving the head and manipulating the objects (hand actions). The actual input was preprocessed and normalized in order to reduce its size and make it suitable for the neural architecture. The preprocessed input fed to the primary body network comprised 120 units, each encoding one of the preprocessed sensor readings. This part of the network was experimentally set to various sizes and trained on the preprocessed data as well as the second-order body part maps, which received only inputs from the appropriate body parts.

To evaluate the network, Wermter and Elshaw tested the primary Body Part SOM and the three separate SOMs. The primary SOM was able to distinguish the body parts successfully when its size was 12×12 neurons. In case of separate maps, 8×8 neurons sufficed. These results show that the architecture is able to process and categorize the sensor reading, i.e. low semantic features, in a somatotopic manner. However, the linguistic module and the association areas, which are crucial for a real association between language and action remained only at a theoretical level of description.

Evaluation

Unfortunately, the authors only mention results from partial testing of the architecture. To fully evaluate the functionality and properties of the model, the whole architecture should be employed and, in the best scenario, also tested online in a physical robot. Since the association area can be activated in the same way when the agent perceives an action, executes an action and comprehends its linguistic label, the model can be considered an implementation of the common coding theory. This model is a good example of a biologically motivated architecture.

4.3.5 Linking language to motor chains

An interesting account on the neural substrate of common coding of perception, action and language is the *chain model* of Chersi et al. (2010). The empirical background of such a model draws on the evidence for goal-understanding function of mirror neurons (see Sec. 2.2.7). More precisely, of neuropsychological experiments (Fogassi et al., 2005; Bonini et al., 2010) that showed that the activation patterns of mirror neurons differ with the goal of the action. Chersi and colleagues propose that these differences account for a chain organization of motor and mirror neurons in parietal and premotor cortices. Such chains of neurons encode short habitual action sequences. Consequently, the execution and comprehension (mental simulation) of actions correspond to the propagation of activity within the specific chains of actions.

Similarly, perception of action-related language involves motor resonance in the form of activation of particular pools of mirror neurons. The aim of this model is to explain different experimental results on the interaction of motor resonance and language comprehension (interference versus facilitation). It also proposes a varying degree of crosstalk between neuronal populations. This co-activation of neurons in different pools depends on whether they encode the same motor act, the same effector or the same action goal. For instance neurons encoding reaching will be active for different types of grasps.

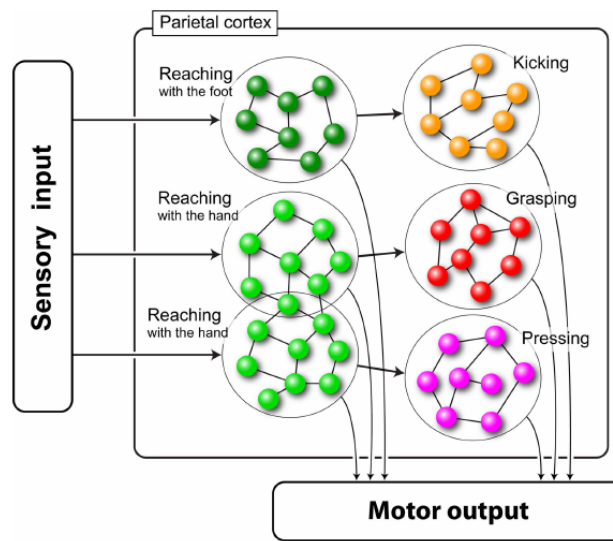


Figure 4.17: Schematic depiction neuronal pools and connections between them (Chersi et al., 2010).

Chersi and colleagues implemented their theory in a computational model based on spiking neurons. They modeled a hypothetical subject whose task was to read two sentences involving grasping and kicking actions, and subsequently reach and press a button after a “go” signal. The model consists of six neuron pools (one pool for one motor act, see Fig. 4.17), the behavior of which is described by a firing rate model with time-dependent synaptic currents (Dayan and Abbott, 2001). The model encompasses complex interactions between excitatory and inhibitory neurons within the pools, the dynamics of ionic currents and neurotransmitters, which are beyond the scope of this thesis.

The results of experiments show that even a simple model can reproduce various experimental results by exploiting only “low-level” properties of neurons. This confirms the idea that motor-resonance-based interaction effects might occur due to neurodynamical factors within the mirror neuron circuit rather than due to some high-level cognitive processes. Concluding the results of their experiment, Chersi and colleagues claim that “showing that interference and facilitation are actually two manifestations of the same process greatly strengthens the embodied view according to which the re-

cruitment of the motor system is fundamental for sentence comprehension”. Results of such experiments with biologically plausible computational representations of neurons confirm that there still exists a need for explorations in the field of actual motor control and its neural correlates. For comparison consider the “preprocessed” models with strong assumptions as the one described in Sec. 4.3.7. One of the aims of the cognitive robotics is to find a faithful compromise between the concrete and abstract levels of description.

4.3.6 Connectionist model of symbol grounding transfer

Cangelosi and Riga (2006) examined how grounded meanings of words can be combined to form the meaning of a new word. This process is called symbol grounding transfer, since the meaning is not grounded directly, but in its simpler components. For instance, words like “horse” and “horn” can be combined to a new concept “unicorn”. It is very likely that high-level abstract concepts, which cannot be directly associated with sensorimotor experience, are grounded using low-level concrete concepts (Barsalou and Wiemer-Hastings, 2005).

The model of Cangelosi and Riga is based on learning through imitation, which has been considered fundamental for acquisition of language (Tomasello, 2003). In their experiment they used two robots simulated in ODE.² These robots consist of two 3-segment arms and a torso with 4 wheels. One of the robots is manually preprogrammed and serves as demonstrator. The other one (imitator), endowed with neural architecture for motor control and language, learns from the demonstrator to perform and name simple actions (movements). To approximate the observed movement of the demonstrator the imitator uses a special “imitation algorithm”, which computes an estimation of the to-be-imitated motor output in the form of forces that applied to the joints of the robot’s body. The imitation algorithm is based on a hyperbolic tangent function and uses the states of joints of the demonstrator,

²ODE is an open source library for simulating rigid body dynamics with advanced joint types and integrated collision detection with friction.

as if the imitator had not only visual information, but something like a whole 3D model of the demonstrator in the current moment. This approximation is then used for training of the neural controller of the robot.

The neural controller of the robot is a standard multi-layer perceptron. It has 26 input units that encode the names of all possible actions in the environment. They are connected to 8 hidden units. The last layer consists of 8 output units, each representing one motorized joint. The network is trained online, using standard BP algorithm. During training the imitator watches and listens to the demonstrator and tries to repeat the action. The actual motor output of the imitator is compared to the ideal outcome computed using the imitation algorithm and the difference (error) is back-propagated through the network. What is actually learned, is the association between the body states and words on the input. The ideal outcome of the learning is the imitator's ability to produce actions correctly on the basis of verbal description of the input.

Cangelosi and Riga divided the training into three stages depending on the level of grounding and its transfer of the learned words. First the robot learns directly grounded actions. In the next two phases the robot learns to ground higher-level descriptions (words) based on verbal description from the user, so it only learns from the peer agent in the beginning. The higher-level description always has a form: new word (is a combination of) action1 (and) action2. An example of such a description could be "grab close-left-arm close-right-arm". In the last phase, the agent learns to combine basic actions with higher-level actions from the previous stage. The results showed that the robot was indeed able to learn composite higher-level words and execute the actions correctly (to a certain extent).

In a subsequent work, Cangelosi et al. (2007) extended this simple architecture to not only comprehend, but also to produce verbal descriptions of actions. Using the same robotic design and imitative learning paradigm, the novelty the authors introduced was a modification in the neural architecture. It remained a MLP, but visual (a simplified retina) and motor information was added to the input layer, the hidden layer was enlarged and neurons for linguistic descriptions were added to the output layer. The architecture

was trained in three phases in the same way as the original model. First the robot learned to execute different actions without language, but in association with their visual form (seen in the demonstrator). Then the robot learned to name the actions and produce them on verbal command, but without vision. Lastly, the prelearned components were integrated. In the final stage of learning, the robot learned also the higher-level words solely from linguistic inputs. Cangelosi and colleagues point out that the learning of such words without visual or desired motor input works on the basis of sensorimotor mental simulation (Barsalou, 1999), which was in the focus of the previous section of my thesis.

Evaluation

Although both the basic model of Cangelosi and colleagues and its extension are nice examples of grounded acquisition of words, and especially a demonstration of possible mechanisms for grounding of abstract words, they have certain drawbacks. First, it is quite impossible in real-life imitation for the imitator to know exactly all the positions and forces in all joints of the demonstrator. However the authors try to compensate for this shortcut introducing an algorithm to slightly modify these values for the imitator agent, so it seems that the training data might have been produced in observation. Such simplification might be explained on the basis of common coding theory as well. If the perception and action are coded in a common domain, one can find a motor plan for the observed action on the basis of matching the observed action with the perceptual aspects of action (or actions) stored in memory. Cangelosi et al. (2007) implemented something similar in the second version of the model, which uses both visual and motor information, apart from linguistic input.

In line with other models in which BP was used, this model might be as well considered a little implausible. Standard BP might be exchanged with continuous RL method such as CACLA, especially for a quite simple multi-layer perceptron.

One might make similar claims regarding the learning as in the previous models. Although BPTT is a very effective method, it is not biologically

plausible. A more plausible form of learning, especially in case of such simple model (regarding the amounts of layers and neurons in them), would be welcomed. As mentioned earlier in the text, RL is a most suitable candidate.

4.3.7 Neural theory of language

The neural theory of language (NTL) (Feldman and Narayanan, 2004; Feldman, 2006) is a very influential account on the neural implementation of language learning and understanding. This theory served as a background for the claim of Gallese and Lakoff (2005) on the neural exploitation mechanism for higher cognitive functions (described in Sec. 2.3 of this thesis). The NTL indeed accounts for the distributed nature of language comprehension and involvement of various modalities in it. The main concept here is the term *synergy* that stands for the coordinated movement of a range of muscles (for instance grasping), which is stereotyped, but also has to be parameterized (e.g. different hand postures according to different objects to be grasped). Feldman and Narayanan (2004) claim that the complex synergy is the core semantics of the word. They regard mirror neurons as a multimodal neural substrate for actions and action words. The matching properties of mirror neurons can account for variable meanings of action words, so word “grasping” can cover all situations: grasping, being grasped or observing grasping at the same time. This is in line with the three basic perspectives (observer, agent, and experiencer), which the child encounters during language learning (Bailey et al., 1998).

The understanding of language in NTL works on the basis of multimodal integrative representations, as suggested by the common coding theory. To contrast the cognitivist account, namely Chomsky’s *generative grammar* (Chomsky, 1966), Feldman and Narayan propose a modeling paradigm based on the *construction grammar* (CG). The aim of CG is to describe linguistic construction forms and link them to their embodied meaning. The CG approach generally assumes that there already exists a fully functional model for perception trained separately. This assumption seems to extend also to actual motor control. In CG, all linguistic elements, from simple words to

whole discourses are modeled in the form of $\langle form, meaning \rangle$ pairs. In the case of complex utterances, these meaning pairs are decomposed to constructional composition of meanings of their parts. As an example of NTL based model, Feldman and Narayan put forward a model of acquisition of simple action verbs (Bailey, 1997).

Model of acquisition of action verbs

In his dissertation thesis, Bailey (1997) created a model of early acquisition of action verbs, which operates on an abstract computational level. Bailey et al. (1997) differentiate between four levels of discourse. The highest is the cognitive level that comprises words and concepts. This level is implemented at the computational level, which is the focus of this study. The bottom two are the connectionist level (ANNs) and the neural level, the latter being still implicit (it requires further knowledge of the brain to implement). Such a hierarchy is quite interesting, however there is no direct evidence for this kind of division to be implemented in the brain. It thus remains in the top-down and high-level area of modeling.

Similarly to young children, Bailey's model has to deal with a *correlation problem* of how to realize what features and actions is the teacher (parent) talking about. To implement this, the model uses two types of representational entities: *action (x-) schemas* and *feature (f-) structures*. It is important to note, that the model takes a great advantage of some interesting properties of human cognition and children's language acquisition. First, it implements the gestalt perception, our ability to perceive separate features more as a whole than separately, using conjunctive representations of features (described in the following text). Then it takes for granted an assumption that children learn language without explicit negative feedback (only from positive examples). And finally, it draws on the fact that children are capable of the so called fast mapping (Carey, 1978). Sometimes a child learns the meaning of a word from just one example. I will discuss these properties and their plausibility at the end of this section.

Executing (x-) schemas, which describe actions (synergic activity) and their aspects, are implemented using Petri nets (Murata, 1989). Petri nets consist of places, transitions, and directed arcs between them, which always connect places to transitions (not to the same-type entities). Places are predicates that may contain a natural number of tokens. These tokens sort of “flow” in the network and cause the transitions to fire. When they fire, transitions consume tokens and place another tokens on their other ends. An important aspect of a Petri net is the nondeterminism. An enabled transition may fire, but it does not have to. An example of a net for a grasping action could comprise some information from the world state, like the size of the to-be-grasped object and transitions representing movements which form the action (reaching, various preshaping hand postures, final execution of the grasp, etc.)

The x-schemas are bidirectionally connected to and receive information from the f-structs. These can have probabilistic values that are consistent with prototype theory of categorization (Heider[Rosch], 1972). Apart from representing different sets of features describing actions of the agent and the current state of the world, f-structs also encode multiple senses of verbs. The meaning of a verb is then a set of f-structs that describe various categories of movements named by the word. The desired behavior of the agent endowed with this model is to react to verbal commands appropriately. The interpretation of a verbal command is a process of choosing a motor action and its parameters in accordance with the world state allowing context dependency. The best matching f-struct is selected and used to guide the execution of the movement (set parameters for the x-schema).

The agent learns action verbs as an additional input while it is executing an action. The ability to produce various movements is already hard-wired. To accomplish learning, the model uses Bayesian model-merging (Omohundro, 1993). At first, it assumes that each example is a separate word sense. Secondly, after a batch of presentations (or after one in online case) the model evaluates all word senses (f-structs) and potentially merges some of them to provide a better description of the training set. In practice learning means to apply Bayes’ law to yield $P(L)P(T|L)$ and maximizing this product, where

the higher probability $P(L)$ is assigned to more compact languages L , and contrasted with the likelihood probability $P(T|L)$ of L and the training set T . The results of various experiments with this models showed, that it is able to successfully encompass the training set of 18 different verbs (Feldman and Narayanan, 2004) in English and a similar amount of verbs in some other languages.

Evaluation

The NTL has two parallel implications. First, there is the theory on the involvement of various neural mechanisms in language comprehension and learning. Since I discussed this aspect and possible neural substrate for it in case of language about action (mirror neurons), I will not continue the discussion on this topic. On the other hand, there is the CG-based paradigm of top-down modeling based on the abstraction from movement synergies. The routine movements might be represented as parameterized schemes. However, such models take a vast amount of mechanisms for granted. The whole perceptual and motor component is omitted, leaving the representations to be specified by the designer of the model. Regardless of how extensive is the designer's effort to describe actions of either humans, robotic or artificial agents, this description might never be complete, and indeed will not be embodied in the strict sense. Similarly the agents knowledge of the world state is not derived from its perception but predefined. Although the motivation for the input data which does not contain bad examples is plausible, there is no noise introduced to the learning.

An important question is, whether such models are or are not grounded. From one point of view, we can create abstract agents for which the abstractions like x-schemas and f-structs are the elementary representations of actions. However, such models might only be able to learn and function properly in highly abstract simulated environments. In the case of embodied agents, a certain "conversion mechanism" to these higher-level representations is needed. The continuous nature of perception and action is partially represented by the probabilistic values. However, representations in Bailey's model still gives a strong impression of the ungroundedness. Since percep-

tion and action are assumed to be preprocessed by an unspecified mechanism, there emerges a problem similar to the unidentified nature of the *transduction* in Fodorian amodal symbols, criticized by Barsalou (1999).

The main advantage of this type of modeling lays in the abstraction. A CG-based model is probably much less computationally demanding, so it can provide a larger variety of functions. It also seems useful for testing of some aspects of language learning. However, it does not explain as much as some selected processes, leaving the gap between actual execution of actions and their naming. The connection between perception, actual execution of an action and language that describes it is crucial, and should be studied in connection, not separately. The direct grounding of language in motor control was studied by Marocco et al. (2010) (Sec. 4.3.1).

4.3.8 TWIG

Lastly, I briefly present an example of non-neural based embodied language learning. Transportable Word Intension Generator (TWIG) (Gold et al., 2009) is an interesting example of a formal logic based system that learns compositional meanings grounded in its sensory experience without supervision. Inspired by developmental psychology, TWIG models some of the basic strategies that children use to acquire new meanings of words in the absence of direct referents.³ First of these strategies is the usage of the grammatical context of the sentence. For instance, in English one can infer the meaning of an unknown word from its position in the sentence and relationship with other words (Brown, 1957). In sentence “A pig *foo* the ball” it is obvious, from the surrounding context that *foo* has to be a verb. The second heuristic is the principle of contrast (Clark, 1987), according to which children contrast the unknown new word with the known words resulting in acquisition of a new meaning and possibly in a change of the meanings of the already known words or concepts.

³Children often learn language from conversations which they do not take part in. Thus the referents of the actors and objects are not directly labeled, by pointing or gazing, and therefore must be inferred (O’Grady and O’Grady, 2005).

To successfully accomplish capabilities as the finding of the referent in the sentence, using language to learn language, understanding of deictic pronouns (e.g. I, you or this), and finally language production, TWIG uses a combination of two formal techniques. First, the *extension inference* is used to infer the real world referent of a new word in sentence. In the second step, definitions of words are created in the form of special types of decision trees, called *definition trees*, which describe the process in which the speaker chooses the word, and encompass the syntax as well as the meaning of words and sentences. These two mechanisms are used to acquire meanings of new words and produce grammatically and semantically correct sentences about the world the agents experiences.

In order to realize experiments, Gold and colleagues implemented TWIG in a humanoid non-mobile robot, equipped with video cameras, dual-channel microphones, and an ultra-sound communication system for localizing objects on the scene and their distances. Unfortunately, the estimation of distances requires all objects that the agent can evaluate to have a special device, which could cause significant difficulties when scaling the environment of the robot. However, this setting was not particularly specific to TWIG, which actually only requires to be connected with a robot which is able to generate predicates about the environment, such as locations of objects and people, and some basic relationships between them. Since the sensory information is inherently continuous, the system is designed to allow such values in the predicates, and subsequently generate thresholds for these values. Although, the paper does not clearly state, how the system generates these thresholds, or whether they must be initialized by the experimenter, what casts some shades on the unsupervised character of the learning in this system.

Evaluation

An interesting example of different approach to grounding. However, from the previous relationship of embodied cognitive science and formal approaches one has to maintain caution. Albeit plausible from our conscious experience, formal systems such as definition trees or inference mechanisms have yet not been empirically proven to be implemented in the brain. From

the cognitive robotics point of view, this model represents a good example of how the robot might exploit its environment. However, the specificity of the design might cause problem when scaling the model, or even more problems when embedding the agent into a real-world tasks and environments.

Chapter 5

Towards a robotic model of the mirror neuron system

This chapter is devoted to my thesis project. The main contribution of this thesis is the proposed robotic MNS model and the novel learning algorithm for ANN.

In the previous chapter, I described and evaluated various models of the mirror neuron system. As already mentioned, computational models of MNS are considered a quite prominent tool in mirror neuron research, complementary to various neurophysiological and neuropsychological studies. Oztop et al. (2013) emphasize that “computational models provide sufficient and causal explanations for observed phenomena involving mirror systems and the learning processes which form them”. Computational models also point out to the need of accounting for additional circuitry to “lift up” (escalate or evolve) the MNS of a monkey to maintain the quite diverse and complex nature of MNS in humans (see Sec. 2.2.2).

Most of the above mentioned MNS models aim at explaining high-level aspects of action understanding (such as understanding goals or emotions of the observed agent) and the role of mirror neurons in it. However, these models abstract from the problem of translating the perspective of the observed action. Typically such models take for granted that perspective-invariant coding emerges in the STS area feeding to area F5 exploiting the PF path-

way (Sec. 5.1.3). However, as revealed by recent evidence (e.g. Nelissen et al., 2011), invariant representations emerge in anterior bank of STS (STSa), which is indeed connected to F5 (F5a), but through a different pathway including the AIP area. Additionally, perspective variant representations have also been found in the core mirror neuron area F5 (Caggiano et al., 2011).

A new approach to exploring and modeling mirror neurons was adopted by Tessitore et al. (Sec. 4.2.7) who used the computational paradigm to account for results of psychological experiments. In short, the main assumption of Tessitore et al. (2010) is that mirror neurons facilitate action recognition and control processes, since they provide a simplified motor representation that narrows a wide search space of the visual input. Their model represents a mapping function from visual representation (an image of the hand during grasping action) to motor representation (created with a special recording glove). An interesting property of this model is that it directly solves the problem of translation of perspective, i.e. from the observer to oneself. As mentioned above, mirror neurons react to different perspectives similarly to neurons in STSp. Therefore, the perspective information should be encompassed by a computational model of mirror neurons as well. The aim of our model is to account for the existence of view-dependency of neurons in both STS and F5 as a possible outcome of their bidirectional connectivity.

In the beginning of Sec. 4.2, I divided MNS models to two classes: (1) models that aim at direct modeling of brain areas, and (2) models whose function was inspired by mirror neuron function and enable the cognitive robots to associate motor and perceptual stimuli in a way resembling the common coding theory (Sec. 2.1.2). A typical member of class 1 would be for instance the MSI model and of class 2 the RNNPB model. Our model is a sort of hybrid between these two classes; it roughly models brain areas at certain level of abstraction. On the other hand, the final aim of our model is to allow the robot (simulated iCub) to associate the observed action with actions in its own motor repertoire allowing it to make categorical judgments on the observed movement and to a certain degree to “understand” the observed movement. The modular architecture of our model is described in the next section.

5.1 Architecture of the model

Our robotic MNS model (Fig. 5.1) consists of several modules. Based on the network architecture and function it can be divided into four layers:

1. executive and perceptual modules at the bottom
2. high-level representations of motor and visual sequences (areas STSp and F5_{mir})
3. PF pathway, which connects STSp and F5_{mir}
4. AIP pathway, which connects STSa and F5_{mir}

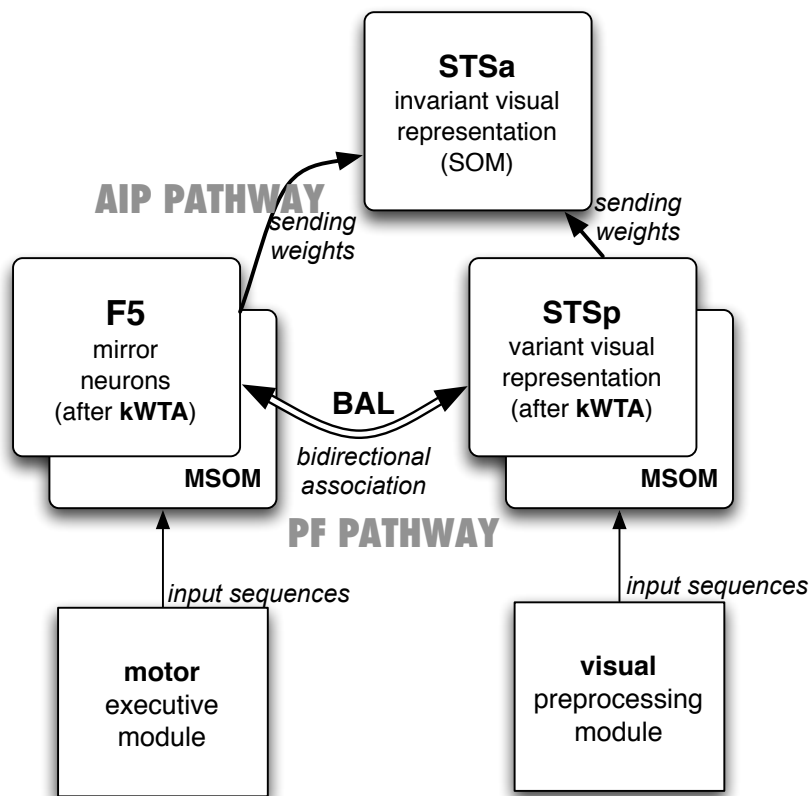


Figure 5.1: The sketch of our robotic MSN model.

The sketch of our model is displayed in Fig. 5.1 and was recently published in Rebrová et al. (2013). The model assumes that sensory-motor links are

established between higher level representations, rather than directly between low-level representations of the movement as a temporal sequence of the robot arm’s state. Note that we label the F5 area as $F5_{\text{mir}}$ because it represents only the mirror neurons. Since full representations in this area are linked with full representations in STSp, we can only call it a subset of $F5^1$. Also, in the present state, we do not divide the $F5_{\text{mir}}$ area into F5c and F5a according to neural evidence. Neither do we explicitly model the intermediate areas PF and F5. Area PF forms a hidden layer of neurons in BAL network, which we do not access directly, hence its activation is self-organized by the network behavior. AIP pathway is represented by the F5-to-STSa part of the model.

For encoding sequences in a high-level fashion we have chosen the MSOM model by Strickert and Hammer (2005) (see Sec. 3.2.2). In line with Thivierge and Marcus (2007) we consider topographic maps as a ubiquitous organizing principle in the brain. Although it is not known whether STS and F5 areas are organized this way, we have chosen the map organization because it provides compact distributed representations. In the following, text I define the function, architecture, and learning algorithms at all four levels of our models separately.

5.1.1 Executive and perceptual modules

On the lowest level of our robotic MNS model resides a control architecture for a simulated iCub robot (Zdechovan, 2012). More information on iCub robot and its simulator can be found in Sec. 4.1.1. The task of the robot is to learn three types of grasps, power, side and precision grasp, illustrated in Fig. 5.2. The executive module consisted of two neural networks, one for reaching and one for grasping, trained using RL. Since the robot’s state space and actions are continuous in nature, CACLA algorithm was used (see Sec. 3.1.2). Our previous successful implementation of a similar motor module can be found in the work of Farkaš et al. (2012) described in Sec. 4.3.3.

¹In biological systems, mirror neurons form only a subset (around 30%) of the macaques F5 (see Sec. 2.2).

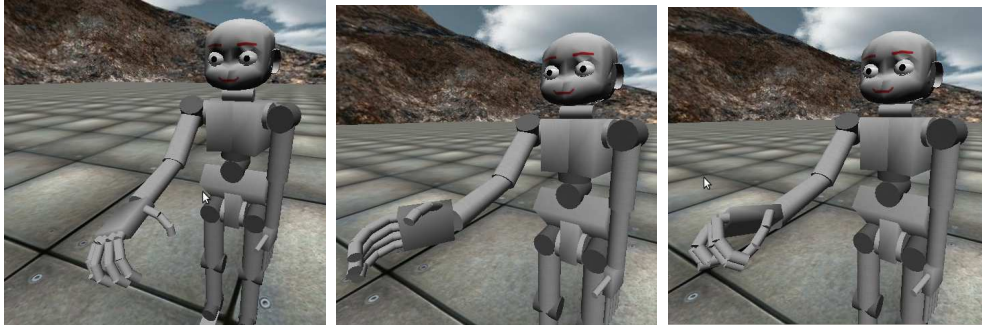


Figure 5.2: Examples of three grasp types from the observer’s perspective. Left to right: power grasp, side grasp and precision grasp.

After the training, motor and visual sequences were extracted from this module and stored to be fed to higher areas at the upper level. Motor sequences comprise the all joint angles from 16 degrees of freedom (DoF) in robot’s right arm stored during the execution of both reaching and grasping. Since these values are given in degrees, we rescaled them to interval $\langle -1, 1 \rangle$, independently for each DoF. The corresponding sensory representations are merely visual, provided by robot’s camera in its right eye (for simplicity, we used monocular information).

Visual information was taken in the form of 3D coordinates of all 16 arm joints (48 values), plus 3D coordinates of four finger tips (total 60) and projected onto the right camera yielding 2D coordinates. As well as the motor data, the values were rescaled to $\langle -1, 1 \rangle$, independently for each coordinate. The primary data sequences represent the self-observing view. To generate visual representations for other perspectives (90° , 180° , and 270°), not directly available from iCub simulator, we used self-observed trajectories (0°) and rotated them correspondingly using appropriate mathematical apparatus. Afterwards, all trajectories were projected onto 2D retina and rescaled.

5.1.2 Higher associative areas

At this level of our model there are two modules that process the low-level motor and visual information, and form high-level representation of movement in F5 and STSp, respectively. These modules are both implemented as

MSOM (Strickert and Hammer, 2005), the self-organizing maps with recurrent context, able to process sequential data. Details on MSOM architecture and learning algorithm are in Sec. 3.2.2. Importantly, after training, responsiveness of MSOM units get organized according to sequence characteristics, biased towards their suffixes. A bias toward sequence suffixes is most suitable, since the whole grasping movement sequences differ mostly at the end (reaching is very similar for all grasp types and the position at the end of the sequence is most characteristic of the grasp type).

In line with our modeling assumptions, we considered the responses of MSOMs in form of sparse distributed representations. The motivation for using sparse codes representations comes from considering the mapping between two domains that could lend itself to generalization and robustness. In biological networks, distributed representations are typically achieved by lateral inhibition. A computational shortcut to achieve such organization is the k -WTA (k winner-take-all) mechanism.

Similarly to standard WTA algorithm with one winner, this mechanism evaluates the response of the map units. In our case, based on the distance between the input and the weights of each neuron the mechanism selects k winning units with the closest distance to the sample. The winner units are then set to 1.0 and all the remaining units (bits in the output vector) are reset to zero. These units and their position on the map form a pattern which represents the concrete motor or visual sequence in F5_{mir} and STSp. Apart from biological plausibility, the advantage of binarization is in facilitation of the training and in simplification and clarity in assessing the network performance.

In conclusion, the activity on the MSOMs binarized with the k -WTA algorithm represents the activity in F5 and STSp areas.

5.1.3 PF pathway

The topmost part of the model includes a three-layer perceptron network as an abstraction of the F5c–PF–STSp circuit. This pathway forms a bidirectional link between invariant motor information from F5 module with

variant perceptual information in STSp via the parietal area (PF). For this pathway specifically, we (Farkaš and Rebrová, 2013) designed the bidirectional activation-based algorithm (BAL) introduced in Sec. 3.1.6. This algorithm was derived from biologically plausible GeneRec created by O’Reilly (1996a) (see Sec. 3.1.5). Like the GeneRec model, BAL is based on plausible two-phase activation spreading and activation-based weight adaptation. The main difference from GeneRec is, that BAL forms completely bidirectional associations. Given unproblematic data pairs, a trained BAL network can produce the desired activation pattern on both input-output layers given the activation of one of them.

In the learning process the network first forms mapping from each motor representation onto visual representation from the self-observation perspective. Subsequently the motor representation is trained to be associated with all possible visual representations, i.e. in 1:4 fashion. We do not expect the model to produce a good outcome on the ambiguous visual stimuli, meaning that when one motor pattern is associated with 4 visual patterns the network will always be confused (as any human would). However, our aim is to form this association in order to trigger the proper motor representations using all visual stimuli. We identify the activity in F5 module triggered through associated sensory information with the activity of mirror neurons.

5.1.4 AIP pathway

The newest part of our model is the F5a–AIP–STSa circuit, which links mirror neuron firing in $F5_{\text{mir}}$ (which is in this stage still invariant) with invariant representations in STSa. As suggested by Tessitore et al. (2010), simplified motor information originating in mirror neurons might facilitate the processing of complex visual inputs in STS. In line with this we assume that the motor information from $F5_{\text{mir}}$ might facilitate the process of forming invariant representations in STSa. One of the aims of our model is to account for this hypothesis.

Results from experiments with MSOMs and robotic data (Sec. 5.3.1) show that while motor MSOM ($F5_{\text{mir}}$) gets organized according to grasp types

(movement is generally perspective invariant), the visual MSOM (STSp) gets organized strongly according to perspective and subsequently according to grasp type. Therefore the organization in STSp is variant as we expected. However, to form invariant representations, variant maps do not suffice. As a first step in our modeling of invariant representations we wanted to connect the $F5_{\text{mir}}$ and STSp modules with the STSa area. For this we have chosen the SOM model (Kohonen, 1997) (see Sec. 3.2.1).

Preliminary experiments with the SOM model we have shown that not even the explicit motor information was able to drive the network to form absolutely invariant representation (i.e. the same network response to one class of movement observed from all perspectives). Therefore, an additional mechanism was necessary to drive the process of forming invariant representations. Competitive and cooperative lateral interactions are quite typical for the brain. Finally, the proposed STSa module consists of a SOM with a cooperative lateral mechanism. Specifically, it is a new set of weights \mathbf{w}^{lat} between all neurons of the map (except self-loops). Together with the lateral interaction the final output activation a of a particular neuron i is:

$$a_i = a_i^{\text{SOM}} + \sum_{j=0}^N \mathbf{w}_{ij}^{\text{lat}} a_j, \quad (5.1)$$

where the activation a^{SOM} of a neuron j from the SOM network is computed according to:

$$a_j^{\text{SOM}} = \exp(-d_j), \quad (5.2)$$

where d is the distance of the particular input from the prototype stored in neuron i computed using standard SOM equations (Sec. 3.2.1).

The strength of the synaptic connection between a neuron i and all other neurons in the map is adapted using:

$$\Delta w_{ij}^{\text{lat}} = \alpha(\lambda d(i, j) + (1 - \lambda)a_b)a_j, \quad (5.3)$$

where $d(i, j)$ is the Gaussian distance (Eq. 3.26) between neurons i and j , α is a non-zero learning rate, and λ is a trade-off factor indicating the influence

of the neurons distance over the influence of the activation of the receiving neuron. Combining this two factors should be beneficial to allow the active neuron to trigger their neighbors in overall more active areas of the map, rather than on the borders between clusters.

5.1.5 Model function and learning

In the process of acquiring the whole MNS functionality, the robot first learns to produce the three grasps. The information from the motor module is processed with the higher level F5c module (MSOM, Sec. 3.2.2) and gets organized on the resulting map as clusters of instances of the same movements. During the production of the movement, the motor information and the visual information from the self-observation perspective gets associated bidirectionally using the BAL algorithm (Sec. 3.1.6). At the same time, we assume that the robot observes another robot producing the same actions and creates visual representations of those actions from different perspectives (self, 90°, 180°, and 270°) in STSp and associates them with the motor representations as well (using BAL). Then, if the robot observes an action from various perspectives, the motor representation of the action is triggered as well. This motor representation, which is basically invariant, then projects to STSa module together with visual information from STSp. In line with Tessitore et al. (2010), motor information helps to form the view-independent representations in the visual areas, thus forming categorical representations potentially used for distinguishing and understanding of the movement as such.

5.2 Experiments with BAL algorithm

In this part of my thesis I summarize preliminary experiments with the BAL first to be presented in Farkaš and Rebrová (2013). As mentioned above, our main motivation for designing BAL was to implement it in our robotic MNS model, to mediate the bidirectional mapping between F5_{mir} and STSp modules (see Sec. 5.1.3).

BAL is based on biologically plausible algorithm GeneRec created by O'Reilly (1996a). Both algorithms are described in Chap. 3. In short, BAL provides a fully bidirectional association between two sets of (arbitrary) patterns in a supervised manner. In line with our modeling assumptions, we train BAL mostly on sparse binary patterns (see Sec. 5.1.2). Before implementing BAL in our MNS model, we evaluated its learning properties in three different experiments with artificial binary data sets. The network task was the same in all cases - to learn bidirectional associations between two datasets. Based on the experimental datasets, the three experiments can be characterized as follows:

1. the standard 4-2-4 encoder task
2. bidirectional association of high-dimensional sparse binary patterns
3. bidirectional association of complex binary patterns

For assessing the network performance, we used three quantitative measures (separately for F and B directions):

1. pattern success (patSucc), which indicates the proportion of output patterns that completely match targets,
2. bit success (bitSucc), the proportion of units matching their target, and
3. mean squared error (MSE) per neuron.

Based on preliminary experiments, we initialize the weights in all tests to small values from the normal distribution $\mathcal{N}(0; 1/\sqrt{n_I + 1})$, where n_I denotes the input data dimension.

5.2.1 4-2-4 encoder

To compare the performance of BAL with GeneRec, we ran tests using the well-known 4-2-4 encoder task. We investigated the convergence of BAL and the number of required training epochs as a function of the learning rate. Fig. 5.3 shows the convergence success for 100 networks and the average

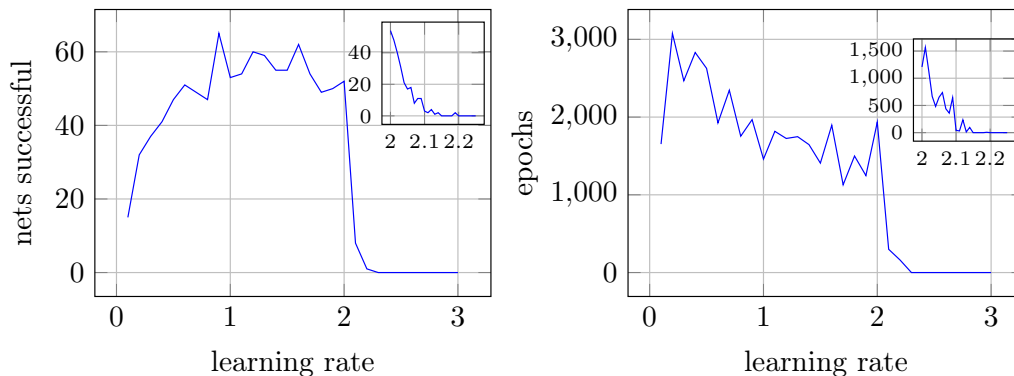


Figure 5.3: 4-2-4 encoder: results for 100 nets, number of successful runs (left), average number of training epochs needed for convergence (right), both as a function of λ . Details for critical values are shown in inset plots.

numbers of epochs needed. The simulations showed that convergence of BAL depends on the learning rate, with the highest number of 65% successful runs achieved for $\lambda = 0.9$. For comparison, O'Reilly (1996a) reports 90% success for basic GeneRec algorithm and 56% for a symmetric modification of GeneRec (Eq. 3.20) and its modification equivalent to CHL (Eq. 3.21).

In sum, probability of BAL convergence is lower than that of basic GeneRec rule, but comparable to its symmetric versions. We expect that the smaller number of successful runs is in both cases influenced by the bidirectional nature of the weight update. Another difference to be accounted is that unlike GeneRec, which provides input-output mapping, BAL is required to manage both input-output and output-input mapping. Even though the values are the same in this case (since it is an encoder), slight disturbances in one direction might cause the other direction to fail as well (remember, that weight update in BAL uses all network activation values in both phases).

BAL was observed to require a higher number of training epochs than GeneRec, with very high variability (and skewed distribution), ranging from 100 to thousands of epochs. On the contrary, O'Reilly reports only 418 epochs for GeneRec to converge, and less than 100 epochs for symmetric versions of GeneRec. An interesting property of BAL is that convergence probability sharply drops to zero beyond certain range of values of the learning rate, for 4-2-4 task at $\lambda = 2$. BAL convergence in 4-2-4 task and sensi-

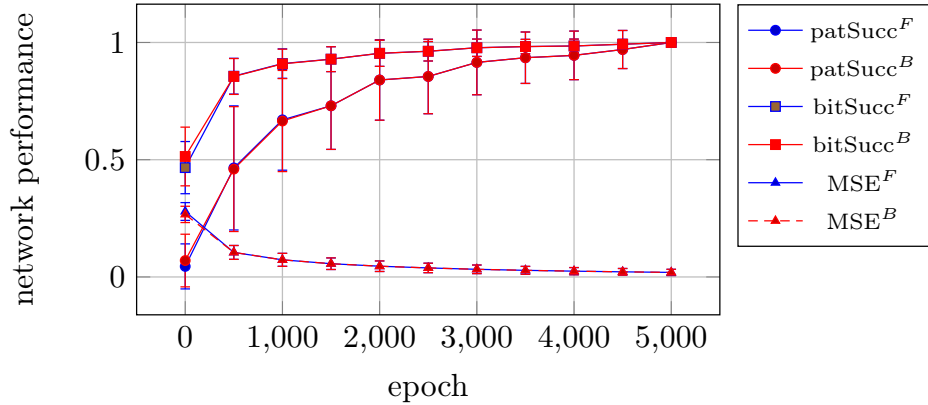


Figure 5.4: Encoder 4-2-4: development of network convergence (50 successful nets).

tivity to learning rate deserves further investigation. Fig. 5.4 illustrates the learning process of 50 successful networks during 5000 epochs using $\lambda = 0.9$. We conclude that MSE drops to minimum values satisfying error-free performance of the network as indicated by all success-based measures (converging to one) in both directions. If the network converges, it masters the encoder task perfectly.

5.2.2 Simple binary vector association

For this experiment we created a sparse binary data set with high dimensionality, which resembles sensory-motor patterns from the F5_{mir} and STSp modules (Sec. 5.1.2). Specifically, we used two sets of 100 binary vectors with 144 bits of which 12 were positive (as if they resulted from k -WTA algorithm with $k = 12$) arbitrarily associated to form one-to-one mapping. Unlike the patterns from the robotic MNS model, these random data do not form clusters of positive bits.

Similarly to the previous experiment, we tested the network performance with various values of learning rate using 144–120–144 architecture. Fig. 5.5 displays the results. The network learns the mapping well up to a certain value of the learning rate ($\lambda = 0.3$), beyond which it is again observed to quickly deteriorate (Fig. 5.5 left). Subsequently, using the estimated optimal

learning rate ($\lambda = 0.2$), we also tested selected sizes of the hidden layer n_H (Fig. 5.5 right). We can conclude that n_H has significant influence only on the amount of training epochs needed to reach 100% success (see the inset figure in Fig. 5.5 right).

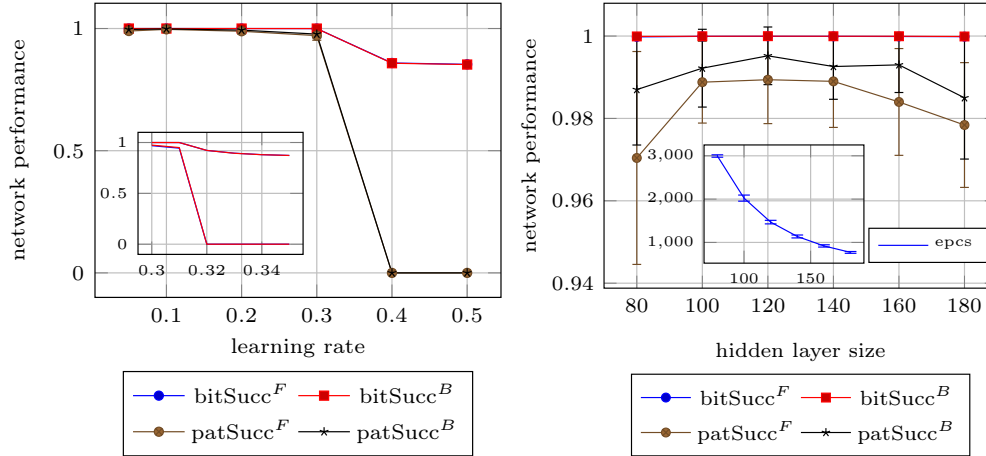


Figure 5.5: Bidirectional associator: network performance as a function of λ (on the left, detail for critical values in the inset plot) and n_H (on the right, with the number of epochs needed in the inset plot).

To demonstrate the network training process, we computed performance measures for 50 nets trained for 2500 epochs using optimized parameters $\lambda = 0.2$ and $n_H = 120$. Results in Fig. 5.6 show that the networks reliably converge to successful mappings between sparse patterns. To understand the network behavior we also examined the hidden layer. We observed that \mathbf{h}^F and \mathbf{h}^B activations have a tendency to move closer to each other, as could be expected from BAL (and also from GeneRec) learning rule. Interestingly, activations of \mathbf{h} units in both directions converged roughly to 0.5, so no tendency towards binary internal representations was observed. This property of internal coding is also worth further investigation.

5.2.3 Complex binary vector association

Motivated by the mappings between invariant motor and variant visual representations, we evaluated the network performance on 1:4 data associa-

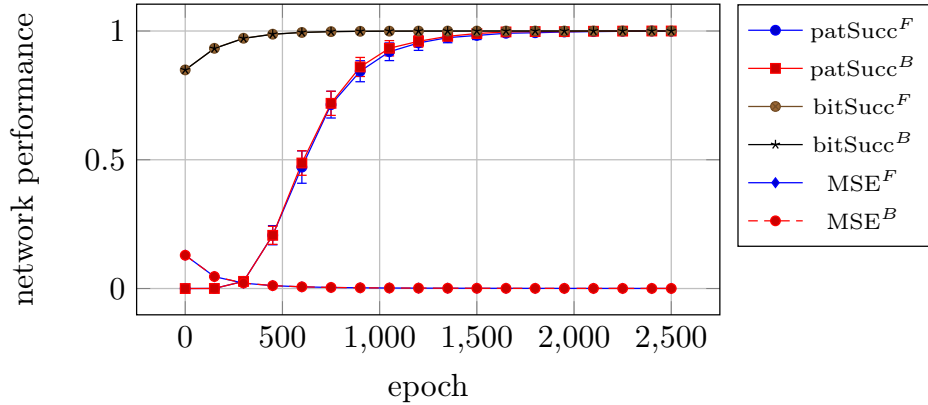


Figure 5.6: Bidirectional associator: development of network performance over time (50 nets).

tions. For this purpose we created low-dimensional sparse binary codes, 16-dimensional vectors (4×4 map) with $k = 3$ active units displayed in Fig. 5.7. For each target (\mathbf{y}), these four patterns (\mathbf{x}) were assumed to have nonzero overlap (1 pixel is shared by all patterns). In this experiment we, again, searched for optimal λ and n_H as displayed in Fig. 5.8. The best performance was achieved using $\lambda \approx 1$. The size of the hidden layer n_H does not seem to influence the network performance.

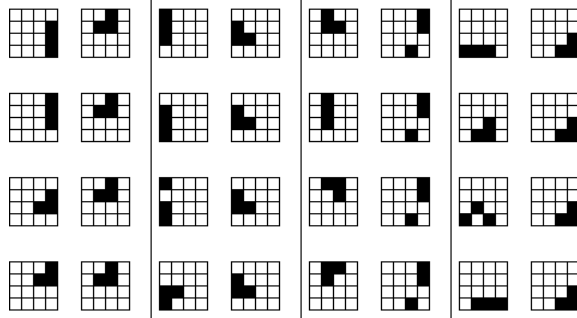


Figure 5.7: Random complex data for BAL experiments.

As in other two experiments, we computed performance measures for 50 nets trained for 2500 epochs using optimized parameters $\lambda = 1.0$ and a default value of $n_H = 14$. Results in Fig. 5.9 show that the network performance is optimal in the unambiguous F direction and worse in B direction. For the best λ the networks yielded $\text{patSucc}^B \approx 4\%$ and $\text{bitSucc}^B \approx 86\%$, which

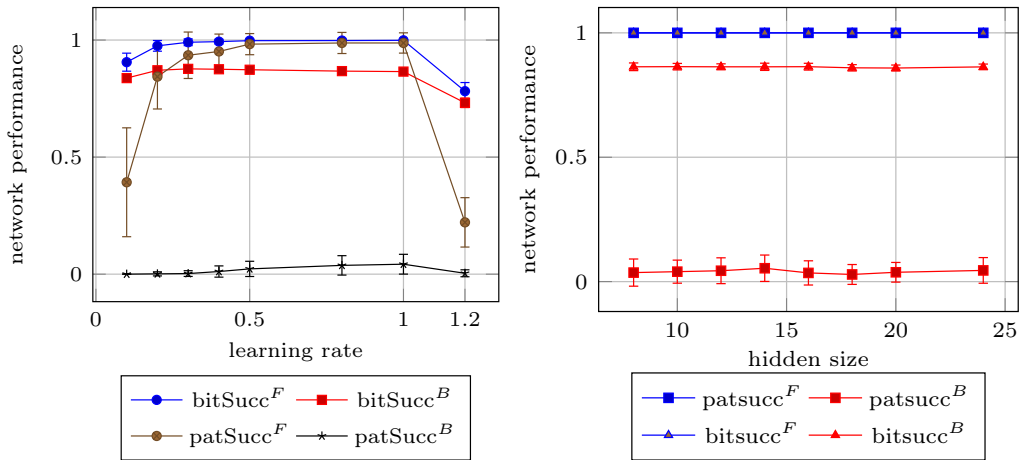


Figure 5.8: Bidirectional associator with complex data: network performance as a function of λ (left) and n_H (right).

means that the networks made some small errors in most of the patterns. This could be expected since the network cannot know which of the four associated patterns is to be reconstructed. It is known, that a network trained to associate more binary target patterns with one pattern tends to produce a mesh of outputs, weighed by their frequency of occurrence in the training set.

Examples of network outputs are illustrated in Fig. 5.10. Only positive-value pixels (i.e. neurons with activation from $\langle 0.5, 1 \rangle$) are filled with color. Green color indicates a match of the target and the estimated value, blue indicates the target activation that was not matched by the output, and red indicates false-positive activation on output. As expected, in the ambiguous B direction the network always tries to produce a mesh of all associated patterns.

5.2.4 Conclusion

In this section I presented results from a new training algorithm BAL for bidirectional mappings. The preliminary experiments have shown that using an appropriate learning rate, the BAL model can converge, albeit requiring more training epochs than GeneRec. In particular, for 4-2-4 encoder task

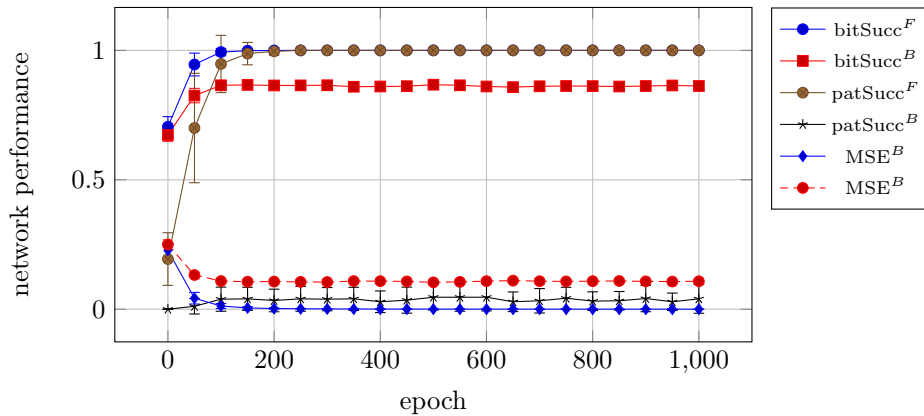


Figure 5.9: Bidirectional associator with complex data: development of network performance over time (50 nets).

the convergence is not guaranteed, which was observed also in the GeneRec model. The next step in this research should be to investigate the reasons for these performance discrepancies. Experiments with BAL also revealed that hidden unit activations tend to converge to similar values for F and B phases. They do not tend to binarize, which is probably not necessary for learning the task. Further experiments and a more detailed analysis of BAL are required to better understand this biologically motivated bidirectional learning algorithm.

5.3 Experiments with robotic MNS

5.3.1 Level 2: self-organization of sensory and motor inputs

These results will be (first) published in our recent work (Rebrová et al., 2013). In this experiment we searched for optimal MSOM parameters and evaluated its suitability for the task. To do this, we used data from the trained iCub as described in Sec. 5.1.1. More details on MSOM architecture and learning can be found in Sec. 3.2.2. Further information and detailed results can be found in master thesis of Pecháč (2013).

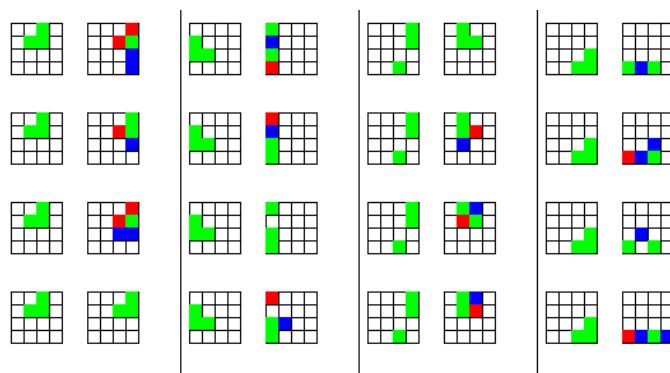


Figure 5.10: Bidirectional associator with complex data: pattern match visualization. Each pixel filled with color represents a neuron activated to 1.0 in the resulting and desired pattern projected through each other. Green color represents a match between target and output, blue pixels indicate positive bits unmatched by the network, and red color indicates the false positive outputs.

Finding the optimal maps

In our experiments we first searched for optimal MSOM parameters. We aimed at getting the map that would optimally distribute its resources (units) for best discrimination of input data. Following the methods from the recent work of Vančo and Farkaš (2010), we calculated three quantitative measures:

1. winner discrimination (WD), which stands for the proportion of different winners to all units at the end of training;
2. entropy (Ent), which evaluates how often various units become winners, so the highest entropy means most balanced unit participation in the competition process;
3. quantization error (QE), that calculates the average error at the unit as a result of quantization process.

To get the best MSOMs, we systematically varied parameters α (Eq. 3.29) and β (Eq. 3.30) in the interval (0,1) and selected the configuration with highest WD and Ent and possibly minimal QE. Optimal parameters α and β as well as other parameters used to create resulting maps (learning rate

γ and number of epochs ϵ) are displayed in Table 5.1. Fig. 5.11 displays contour plots of three tested quantitative measures as a function of α and β .

Table 5.1: Optimal parameters for MSOM-based modules.

module/parameter	α	β	γ	ϵ
STSp	0.3	0.7	0.1	300
F5 _{mir}	0.3	0.5	0.02	300

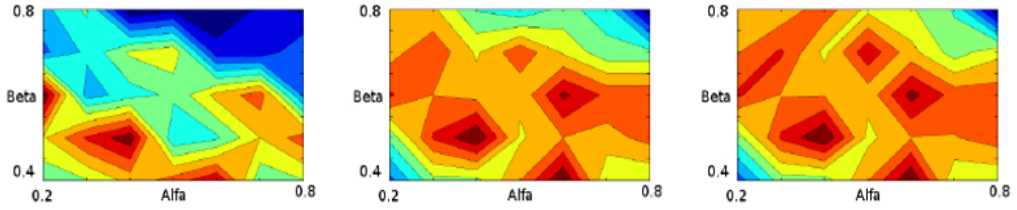


Figure 5.11: Contour plots of MSOM quantitative measures as a function of different α and β . From left to right: WD, Ent, and QE (from Pecháč, 2013).

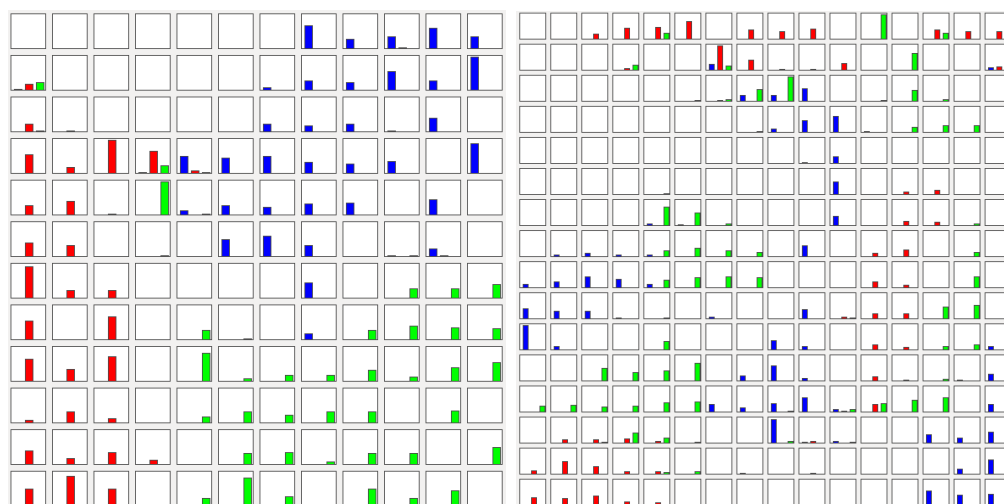
On the basis of educated guess and further experiments, we have chosen to implement maps with following dimensions:

- visual maps: 12×12 , 14×14 , 16×16 , 18×18 , and 20×20 ;
- motor maps: 8×8 , 10×10 , 12×12 , and 14×14 ;

Based on further experiments (described later in this text), we have chosen the map sizes of interest, specifically visual maps with 16×16 units and motor maps with 12×12 .

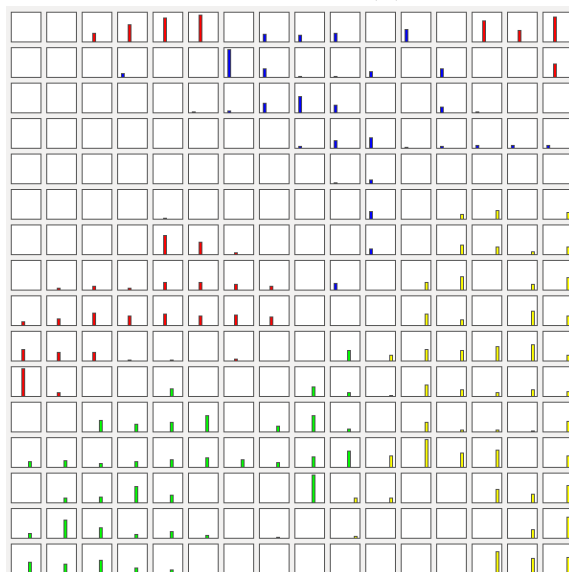
Using these parameters, we trained the MSOMs and evaluated winner hits (i.e. the number of times the particular unit became the winner), for three categories of grasp type for both motor (Fig. 5.12a) and visual dataset (Fig. 5.12b), and additionally for four perspectives for the visual dataset (Fig. 5.12c). Topographic organization of unit's sensitivity is evident in all cases. For visual map, the organization on a coarse level is arranged according to the perspectives, and on a more fine-grained level according to the grasp types. Although some units never became winners, they participated in distributed map representations as well. Topographic order reflects the

natural separability of classes (types of grasps) both in terms of their motor and visual features. Visual maps clearly reveal that perspective is a more strongly distinguishing parameter than the type of grasp.



(a) Motor map: grasp types

(b) Visual map: grasp types



(c) Visual map: perspectives

Figure 5.12: Examples of the trained motor and visual maps.

Generating the MSOM output responses

As mentioned in Sec. 5.1.2, the final output activation from $F5_{\text{mir}}$ and STSp modules is to be a sparse binary code representing a distributed brain representation of either movement or visual percept. We obtain these representations using the k -WTA algorithm. Up to this point we have not realized possible difficulties with the dataset. However, evaluating the binarized responses of MSOMs and subsequently exploring the original datasets from the executive module we realized that our input data are not consistent. Probably, it is caused by the method of obtaining the data. In fact, the training sequences were merely the same action with some additional random noise in the movement. Therefore it is questionable to call this a dataset of movement instances. Therefore, the responses of MSOMs were very similar in all movement classes and specially inside each class. Binarization of such responses caused many patterns to repeat unexpectedly. Unfortunately, when combined with the visual MSOM data to form pairs the associations between them gain m-to-n nature, which is rather ambiguous and impossible to learn for any known neural network.

To find an appropriate small number of k active neurons, I have decided to evaluate the following properties of the resulting data:

1. the amount of unique patterns in the dataset,
2. the average distance between centers of data patterns on the 2D map.

Numbers of unique patterns for visual and motor maps for $1 \leq k \leq 40$ are displayed in Fig. 5.13. Note that the motor dataset contains together 55 instances of three grasp types, and the visual dataset contains 4-times more, i.e. 256 different sequences. If the dataset was not ambiguous, we can expect up to 256 different patterns after binarization with $k = 1$. Unfortunately, the highest amounts of unique patterns with only one positive bit within visual MSOM responses was 117 for map size 20×20 and for motor MSOM responses 42 for map size 14×14 .

To overcome the complications with ambiguous dataset, I have decided to select unique binarized motor patterns and visual patterns bound with

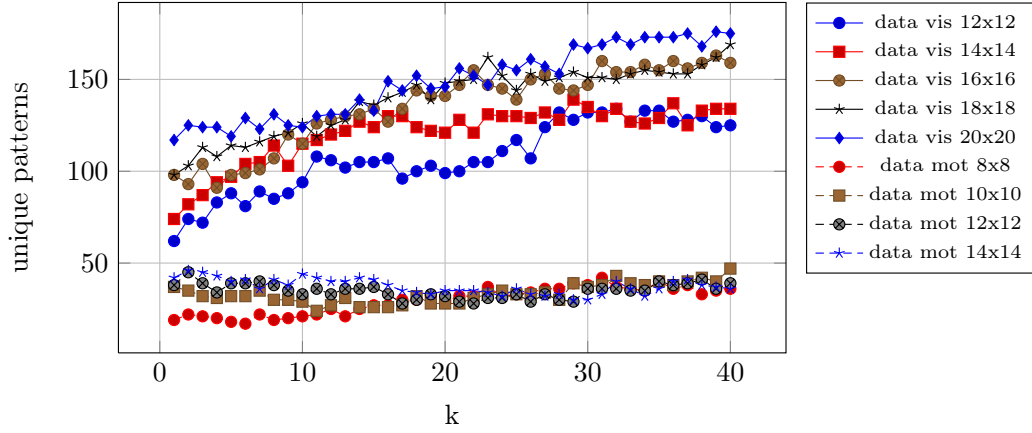


Figure 5.13: Number of unique patterns after k -WTA binarization as a function of k .

them. Subsequently we discarded the multiple occurrences of binary visual patterns forming a “disambiguated” dataset. On the basis of batch testing of above mentioned criteria and also experiments with BAL algorithm on different combinations of maps sizes and values of k , I have found the best dataset with final parameters listed in Table 5.2. From each grasp type I managed to gain 10 unique instances with exception of the second grasp type, for which I found only 9 well-formed instances. However, the original dataset for the second grasp type consisted of 17 instances unlike the first and the third, consisting of 19 instances. The dataset of disambiguated map responses binarized with k -WTA, displayed in Fig. 5.14, served for training the BAL model described in the next section.

Table 5.2: Optimal parameters for MSOM response binarization.

	map size	k
STSp	16×16	16
F5 _{mir}	12×12	16

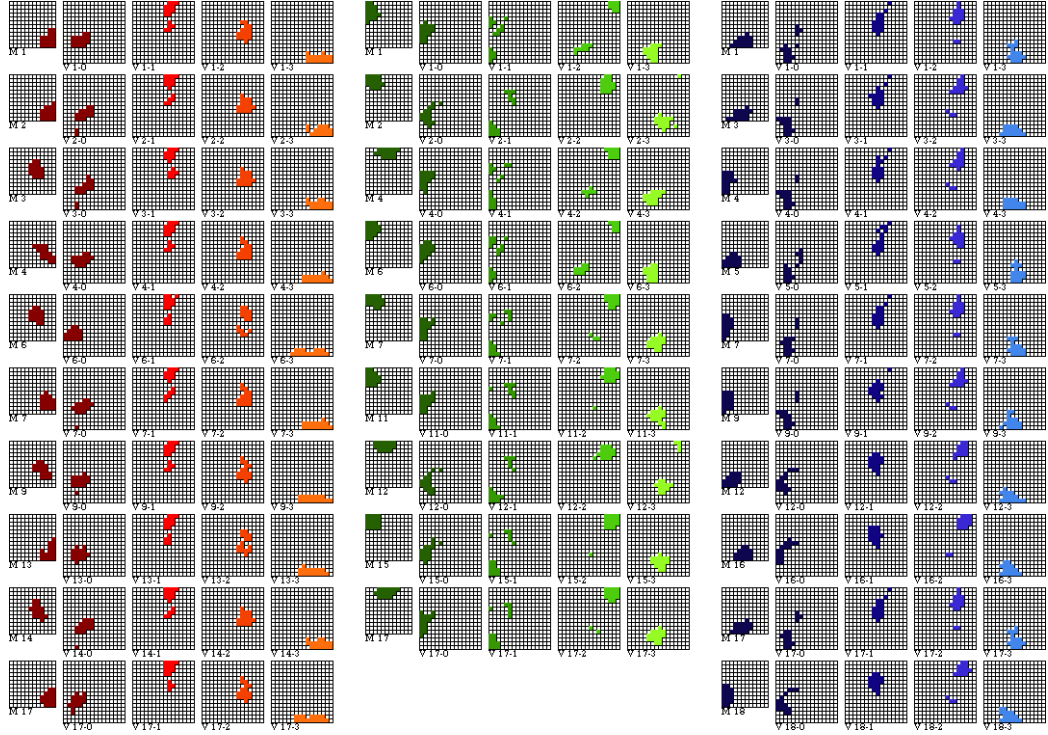


Figure 5.14: Disambiguated dataset from visual MSOM 16×16 and motor MSOM 12×12 , both binarized with $k = 16$. Leftmost side of each column depicts the resulting motor pattern. Four visual perspective-variant representations follow. Each binary motor pattern is associated with four visual patterns.

5.3.2 Level 3: bidirectional association of motor and visual representations

In this section I present results from the three-layer network trained using BAL (Sec. 3.1.6) to form bidirectional associations between sparse binary codes representing the activity in the $F5_{\text{mir}}$ and STSp layers. At the end of Sec. 5.3.1 I have shown that the data I intended to use were somewhat faulty. Therefore I pruned the original dataset of MSOM responses to 29 grasp instances (and 4-times more visual instances). Visualization of this dataset is in Fig. 5.14. As mentioned above, our robotic MNS model first forms a bidirectional association between the agent’s movements and their visual representation from the self-observed perspective. I present results from

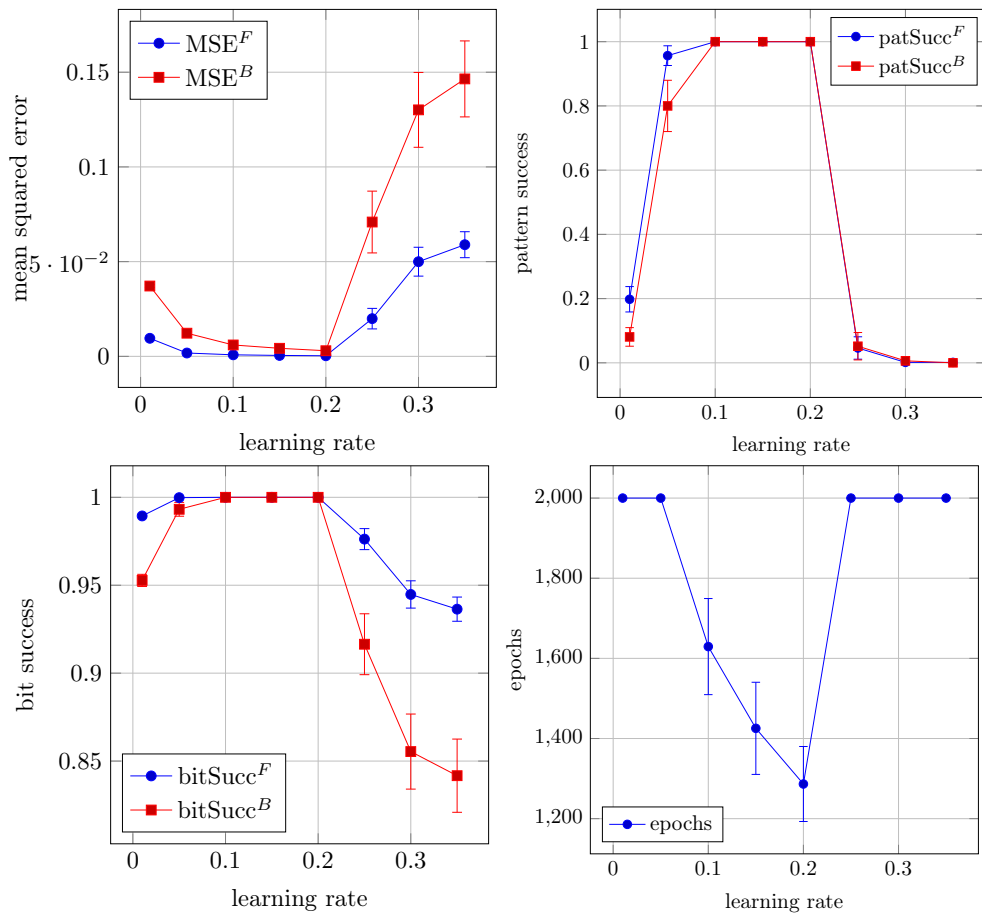
experiments with network parameters, an overview of network performance and a visualization of the network responses to acquired pattern associations. Subsequently, I present results from a *learning experiment* that ran in two phases. First, the self-observing perspective mapping is established. Next, the model is trained data from all perspectives.

The result presented in this section come from the same evaluation basis as preliminary experiments with BAL model (Sec. 5.2), concretely the mean squared error (MSE), the pattern success measure (patSucc), and the bit success measure (bitSucc) for both F and B directions. In testing the map responses, the output of each unit is considered correct if it lies in the correct half of the $\langle 0, 1 \rangle$ interval.

Establishing mapping from self-observing perspective

In this experiment I first evaluated various sizes of the hidden layer (n_H) and various values of the learning rate (α) in the same way as in previous BAL experiments. Additionally to above mentioned performance measures I evaluate also the number of epochs it took the network to reach 100% pattern success in both directions. The results are displayed in Fig. 5.15 and 5.16. From the results I conclude that α is the most influential parameter. Similarly to results in Sec. 3.1.6 from BAL trained on artificial binary data, the network converges to good solutions with $\alpha = 0.2$. Results from experiments with n_H were also similar to previous experiments, concretely I conclude that n_H does not influence the network convergence, only the speed (i.e. number of epochs) of converging to the solution.

To demonstrate the network training process, I computed performance measures for 50 nets trained for 1300 epochs using optimized parameters $\alpha = 0.2$ and $n_H = 240$. Results in Fig. 5.17 show that the networks reliably converge to successful mappings between sparse patterns. To visualize the final output of the network to the training data I used the same method as in Sec. 5.2. Each framed pixel on the map represents a unit in $F5_{\text{mir}}$ or $STSp$. The motor and visual maps are projected with only positive-value units filled with color. Green color indicates a match of the target and the estimated

Figure 5.15: Performance of BAL as a function of α

value (i.e. both the target input and the network output have a positive outcome from this unit), blue indicates the target activation that was not matched by the output, and red indicates false-positive activation on output. The resulting visualization is displayed in Fig. 5.18. Note, that since the network manages the mapping perfectly, all positive units are green. I use the same technique in the next section to evaluate the net responses to all data.

Learning the other perspectives

The learning experiment consists of two phases. In phase 1, the agent associates its own movements with their appearance on the basis of high-level

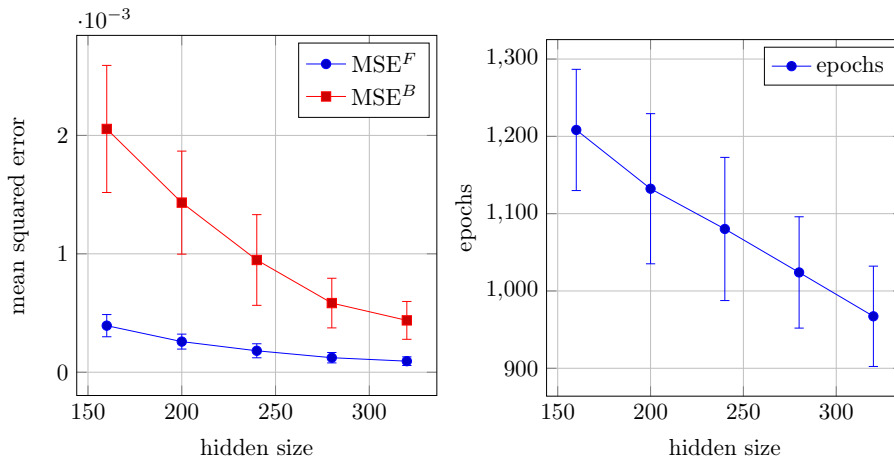
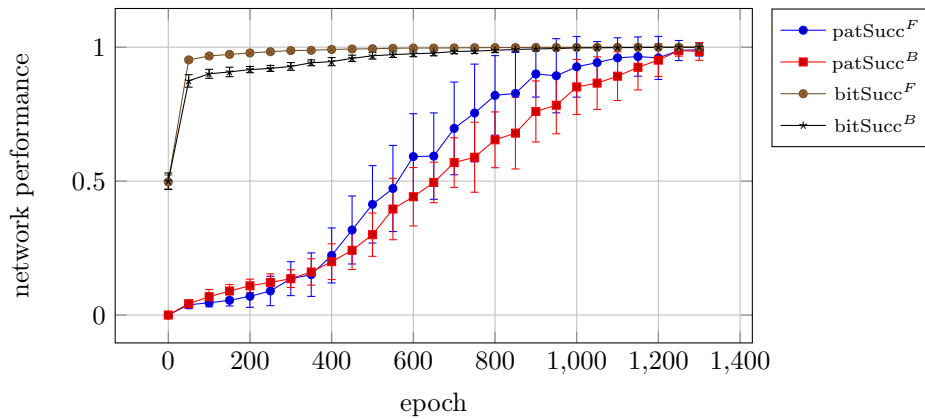
Figure 5.16: Performance of BAL as a function of n_H 

Figure 5.17: BAL with first-perspective robotic data: network performance in time (50 nets).

representations from MSOM binarized using the k -WTA mechanism. In phase 2, the agent forms bidirectional associations also with the movement seen from the other perspectives. Here we use a sort of “scenario-based shortcut”. The robot first produces the self movement, while observing its own arm. Right after it, while the generated motor pattern is assumed to be still residually active, the robot observes the same movement from another perspective (as if it was playing an educational game with its parent). It is known that parents often imitate children’s immediate behavior providing

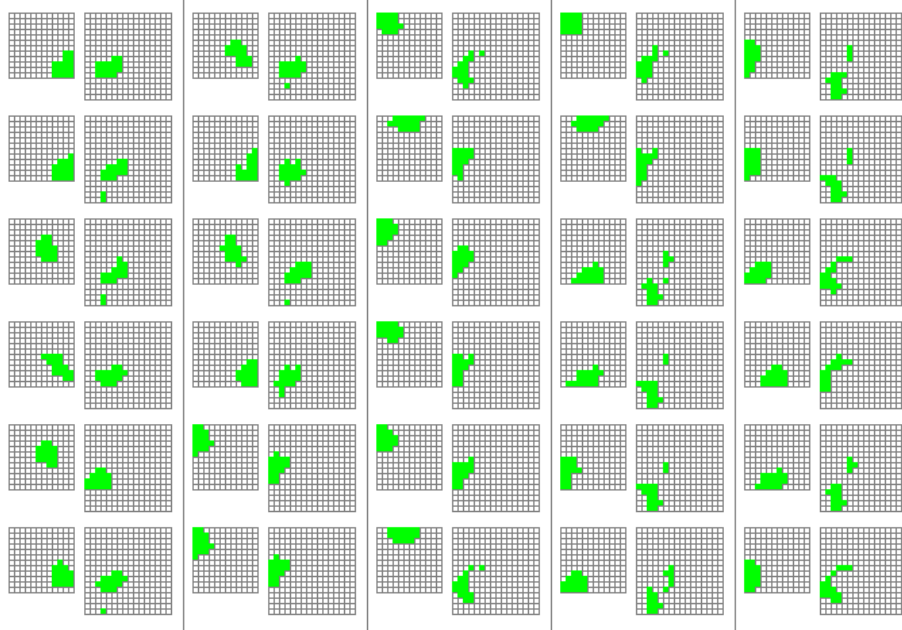


Figure 5.18: BAL with first-perspective robotic data: pattern match.

them with something like a mirror, which may explain how mirror neurons could emerge as a product of associative learning (Heyes, 2010).

The parameters in this experiment are the same as in the previous experiment ($\alpha = 0.2$ and $n_H = 240$) as well as the number of nets (50). Phase 1 ran in 1200 epochs and phase 2 in 3600 epochs. The decision to use a higher number of training epochs for the second phase comes from my observation, that allowing BAL to train for a high number of epochs helps to tweak its final performance (cf. the high amount of epochs in the 4-2-4 experiment Sec. 5.2). Therefore I allowed the second learning phase to take up to 3600 epochs. Note that since the network learns an impossible task from one direction, the algorithm will not converge and reach the stopping criterion. Results from this experiment are displayed in Tab.5.3 and in Fig. 5.19.

The results from phase 1 of learning indicate that BAL algorithm is able to form error-free associations between visual and motor representations. Regarding phase 2, it is clear that the task of bidirectional association of ambiguous data (1-to-4 associations) cannot be accomplished in one direction (compare patSucc^F and patSucc^B). The 50% of completely reconstructed

Table 5.3: BAL performance in two learning phases (50 nets).

measure	phase 1		phase 2	
	F	B	F	B
mse	0.0 ± 0.0	0.0 ± 0.0	0.091 ± 0.003	0.049 ± 0.003
pattern success	1.0 ± 0.0	1.0 ± 0.0	0.006 ± 0.005	0.502 ± 0.023
bit success	1.0 ± 0.0	1.0 ± 0.0	0.905 ± 0.003	0.947 ± 0.004

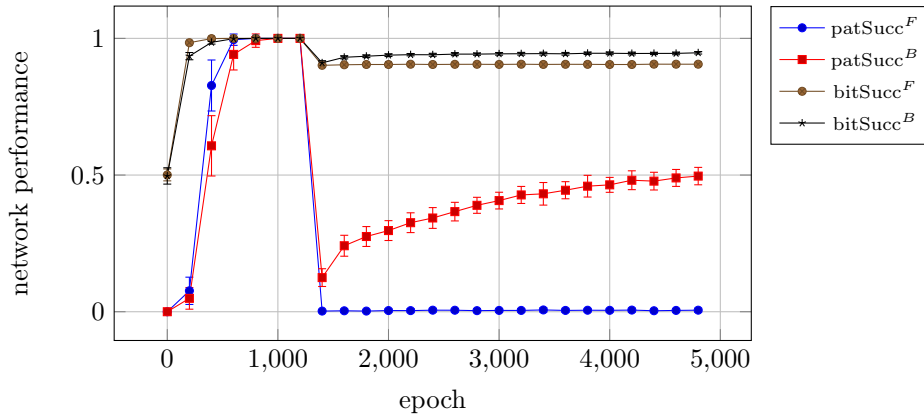


Figure 5.19: BAL performance in two learning phases (50 nets).

motor patterns (i.e. near 100% bitSucc) does not necessarily mean that the network does not manage the task. Although the network matches only 50% of patterns perfectly, the mistakes it make appear to be in close range of the desired pattern. This confirms also quite high bit success (up to 95%).

The visualization of network responses displayed in Fig. 5.20 also suggests (note every left column) that the errors the network makes are close to the target pattern. These errors might be a consequence of very similar, yet not the same, representations of grasp instances. Note that the patterns in one grasp category tend to overlap highly. It seems that the network tends to produce more “typical” representations of the grasp category. Therefore, I assume that the errors the network makes will always remain inside the category, i.e. that the network does not confuse instance of one category for an instance of other category. In conclusion, the representation of a correct movement would be triggered.

Hence, we can conclude that any visual representation of a particular movement will trigger a proper motor representation of this movement, representing the role of mirror neuron activity. The motor information activated on the basis of the visual input can further be used to facilitate the process of forming invariant representation of the movement in STSa. In addition, a mechanism similar to STSa lateral excitation could be implemented here, if more invariant representation of the movement were desirable (i.e. forming broadly congruent mirror neurons from Sec. 2.2.6).

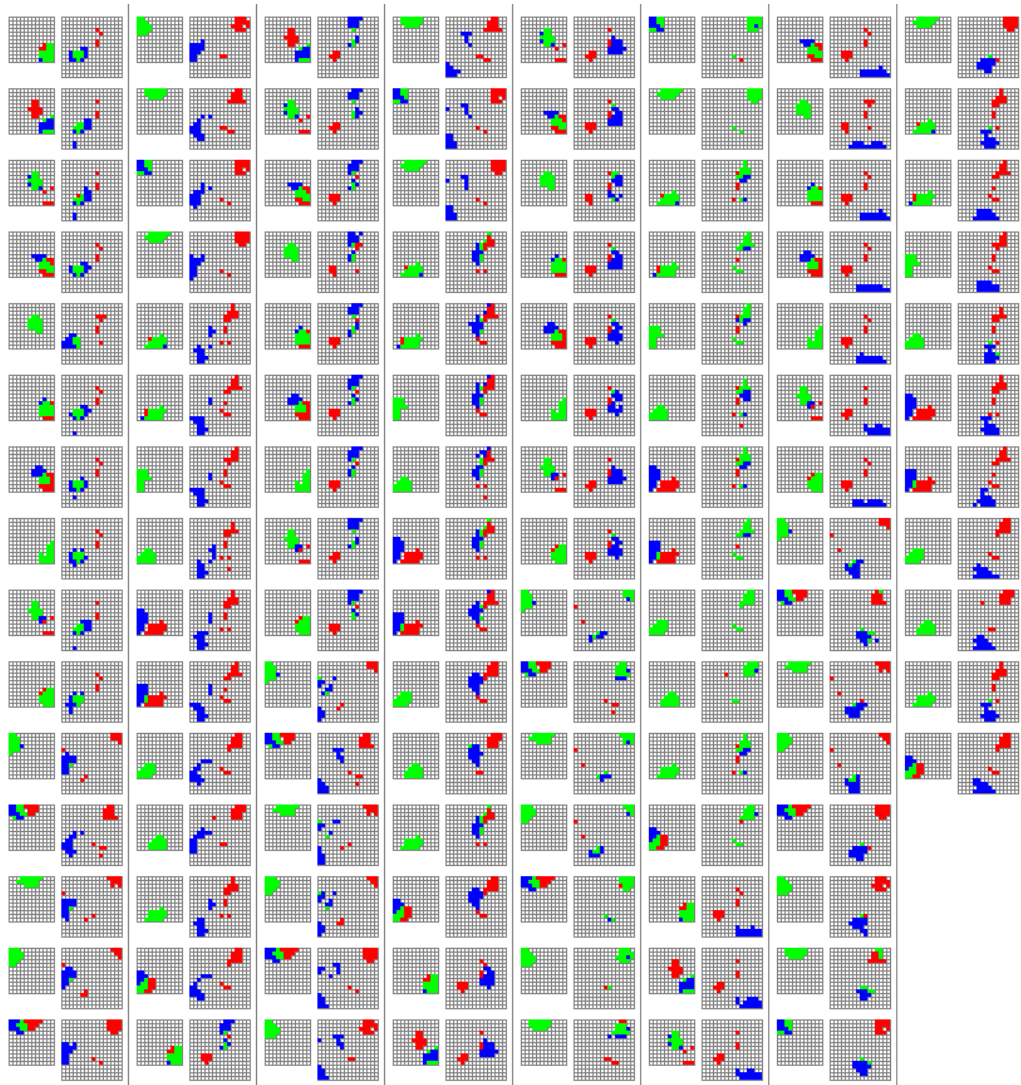


Figure 5.20: BAL after two learning phases: pattern match.

5.3.3 Level 4: self-organization and competition leading to invariance

In this part I briefly report on the results of preliminary experiments with STSa module. As described in Sec. 5.1.4, the module consists of a SOM fed with responses of both STSp and $F5_{\text{mir}}$. After the winner is computed in the map, activation of all units is adapted according to Eq. 5.3. The final organization of STSa area according to grasps and perspectives is displayed in Fig. 5.21b. The SOM part of this model was trained using learning rate $\alpha' = 0.1$ and the lateral interaction was computed using parameters $\alpha^{\text{lat}} = 0.5$, $\gamma = 0.01$, and variance σ for the Gaussian distance (Eq. 3.26).

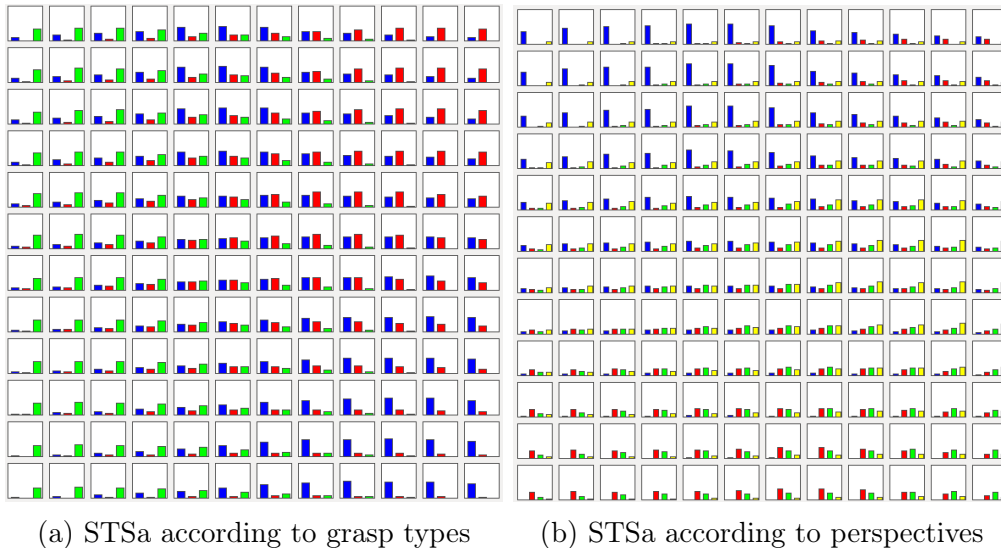


Figure 5.21: Results from STSa module.

Results from this preliminary experiment suggest that the lateral excitation mechanism (Eq. 5.3) helps to create invariant representations. Note that neurons in the central part of the map tend to react to all perspectives and in the corners there are neurons with prominent perspectives. There are also neurons that react to two or three perspectives only. In line with empirical evidence (Jellema and Perrett, 2006) the central representation tends to be more invariant and a topological organization of this process can be observed. On the other hand, we expected the organization to be clustered

around movements with invariant centers in the centuries of gravity of the particular grasp representation. However, neurons in the central part of the map tend to react to all stimuli. These omni-reactive neurons might be interpreted as the highest-level of representation of object grasping, regardless of the grasp type.

5.4 Discussion and future work

The robotic MNS model proposed in this chapter comprises various diverse modules. As depicted in Fig. 5.1, the model retrieves sequential motor and visual information from the lowest level modules. At the middle level, high-level representations of these sequences emerge based on self-organization in $F5_{\text{mir}}$ and STSp. These visual and motor representation are binarized and associated using a supervised BAL algorithm. Bidirectional association provides the model with the means to trigger motor representations with visual stimuli and thus to mimic the function of mirror neurons. Motor representation of the movement is projected back to visual areas through the second pathway. Together with variant visual representations, the motor representations are self-organized on the highest level. To help the emergence of invariant representations, an excitatory local mechanism is introduced to the STSa output map.

Regarding the results from the MSOM-based modules, I can conclude that, although the organization appeared according to expectations, further research is in place. Especially, new experiments should be made with the simulated iCub and the grasping module encompassing a more valid task and environment (unlike the current data which consisted of the same movement instances with noise). For instance, objects of different shapes and sizes, and in different locations might be used to create a diverse set of instances of the learned grasps. New data might shed light on the problem of m-to-n associations and might also cause a slight different organization on the maps. Last but not least, the task with various “stimuli” for the robot is definitely more ecologically valid than just adding noise to existing routines.

Results from the bidirectional activation-based algorithm imply that it indeed converges to right solutions. However, the training time required to converge is higher than that of the original GeneRec algorithm. Our experiments also revealed that hidden unit activations tend to converge to similar values for F and B phases. They do not tend to binarize, which is probably not necessary for learning the task. Further experiments and a more detailed analysis of BAL are required to account for its convergence properties and the hidden-layer representations, which are particularly interesting in the context of the MNS model.

When forming the association between self-observing perspective and motor representation, BAL has no difficulties to converge to 100% success. On the other hand, when we present it with the task, which is impossible to learn in one direction (1-to-4 mapping), the net does not converge even in the other direction, reaching maximum 50% of fully correct outcomes. However, this might be also a consequence of the training data problem mentioned above. Representations of particular grasp instances overlap to a varying degree, causing the net to erroneously produce the more prominent patterns (i.e. activate units that are taking part in most of the representations as false positives).

As to the original goal of this model, which was to account also for variant responses of mirror neurons in F5, further research is desirable. It is yet not known, why different mirror neurons fire in response to presentation of a movement from different perspectives. Neither is known the nature of the information transferred in the two pathways connecting STS and F5. What is certain is that mirror neurons in the first place represent the movement they encode as their primary function during its execution. Therefore, mirror neurons will always be invariant in their nature. These interesting firing properties might, on the other hand, have more complex meaning, for instance reflecting the attenuation to an actor or an object. Although I did not aim at modeling variant cells in $F5_{\text{mir}}$ in this stage of the model development, evaluating this model design and its function, I consider this task very difficult and high-level for the context of cognitive robotics. The main posi-

tive feature of this model is that it associates the information from different modalities, and also that it finally models the firing of mirror neurons.

Evidently, the mirror neuron circuitry and the mirror mechanism still gives birth to many new open questions. Computational models of mirror neurons might not only provide a powerful tool for explaining (Oztop et al., 2013) the neurophysiological and neuropsychological data on mirror neurons, but also challenge the empirical evidence and raise new questions to be experimentally evaluated.

In the future, the STS_a part of the model should be studied in more detail. In fact, the whole architecture should be proven to be scalable and possibly also transferable (the high levels) to a different humanoid robot. As mentioned above, a first improvement might be made collecting more diverse data from the iCub in a more plausible manner. A very interesting direction would be to explore the self-organizing maps with lateral excitatory and inhibitory connections such as proposed in Sec. 5.1.4. This principle can be applied also in forming high-level organization in motor areas, such as the broadly-congruent mirror neurons, which encode the goal of an action rather than the precise method to reach it.

Chapter 6

Conclusion

My thesis discusses how a control architecture (be it a brain or an ANN) might utilize its motor systems to mediate understanding of the actions of others. The claim that perception and action are on a high level represented in a common framework is well rooted in empirical evidence. Embodied cognitive science suggests that motor representations in the brain which can be activated without the actual movement production might serve as simulation mechanisms allowing us to “step into the shoes” of the observed agent. Such mechanism is assumed to be a primary role of mirror neurons in area F5 of the macaque brain.

The discovery of mirror neurons gave rise to various computational models, mainly based on artificial neural networks. Since the computational models of MNS are considered a prominent tool for explaining the mirror neuron function and emergence, most of the models aim at capturing the actual neural circuitry, by having components that directly represent particular parts of the monkey’s brain. On the other hand, there are models that are closer to the paradigm of cognitive robotics, which use properties of the mirror neurons in a specific architecture. Such models do not aim to encompass the neural circuitry, but rather endow the agent with some special capabilities. The robotic MNS model proposed in this thesis also belongs to this category. However, it also aims on faithfully encompassing the MNS and its core brain areas.

The proposed model consists of various mutually interconnected neural networks. To provide a way to form bidirectional mappings between sensory and motor representations in a more biologically plausible way than standard BP, a new learning algorithm BAL was designed. Results from experiments with neural networks comprising this model showed that the model is successfully able to form coherent high-level representations and associate them in 1-to-1 manner. Since bidirectional mapping of 1-to-4 patterns is an impossible task, adding more visual representations to the model reduces the success rate. However, in this way, various visual representations can trigger motor representations hence fulfilling the desired emergence of mirror neuron activity.

In the future, the model should be further evaluated and experimented with. The emergence of invariant representations on the basis of lateral interactions should be studied and possibly implemented in other parts of the model. Also the BAL algorithm should be further studied and possibly enhanced to drive the representations on the hidden layer to reorganize and display something like a mirroring activity, since mirror neurons were found also in the PF area. Last, but not least, a model providing a sensorimotor binding might be scaled to mediate simple language grounding. Concepts connected to action should be connected with motor areas as well. I believe that the common coding theory is most beneficial from the viewpoint of utilizing a theory in cognitive robotics. A common representational ground for action, perception, and meaning might extend to benefit the human-robot interaction.

Bibliography

- Aglioti, S., P. Cesari, M. Romani, and C. Urgesi (2008). Action anticipation and motor resonance in elite basketball players. *Nature Neuroscience* 11(9), pp. 1109–1116.
- Almeida, L. (1987). A learning rule for asynchronous perceptrons with feedback in a combinatorial environment. In M. Caudil and C. Butler (Eds.), *Proceedings of the IEEE First International Conference on Neural Networks, San Diego, CA*, pp. 609–618.
- Alstermark, B., A. Lundberg, U. Norrsell, and E. Sybirska (1981). Integration in descending motor pathways controlling the forelimb in the cat: 9 Differential behavioural defects after spinal cord lesions interrupting defined pathways from higher centres to motoneurons. *Experimental Brain Research* 42, pp. 299–318.
- Arbib, M. (2005). From monkey-like action recognition to human language: An evolutionary framework for neurolinguistics. *Behavioral and Brain Sciences* 28(02), pp. 105–124.
- Asada, M., K. Hosoda, Y. Kuniyoshi, H. Ishiguro, T. Inui, Y. Yoshikawa, M. Ogino, and C. Yoshida (2009). Cognitive developmental robotics: A survey. *IEEE Transactions on Autonomous Mental Development* 1(1), pp. 12–34.
- Bailey, D. (1997). *When push comes to shove: A computational model of the role of motor control in the acquisition of action verbs*. Ph. D. thesis, Computer Science Division, University of California, Berkeley.
- Bailey, D., N. Chang, J. Feldman, and S. Narayanan (1998). Extending embodied lexical development. In *Proceedings of the Twentieth Annual Conference of the Cognitive Science Society, 1998, University of Wisconsin-Madison*, pp. 84–89. Lawrence Erlbaum Associates.
- Bailey, D., J. Feldman, S. Narayanan, and G. Lakoff (1997). Modeling embodied lexical development. In *Proceedings of the nineteenth annual conference of the Cognitive Science Society, 1997, Stanford University, Stanford, CA*, pp. 19–24. Lawrence Erlbaum.
- Barsalou, L. (1999). Perceptual Symbol Systems. *Behavioral and Brain Sciences* 22(04), pp. 577–660.

- Barsalou, L. and K. Wiemer-Hastings (2005). Situating abstract concepts. In D. Pecher and R. Zwaan (Eds.), *Grounding Cognition: The Role of Perception and Action in Memory, Language, and Thinking*, pp. 129–163. Cambridge University Press.
- Barto, A. G. and M. I. Jordan (1987). Gradient following without back-propagation in layered networks. In *1st International Conference on Neural Networks, San Diego*, Volume 2.
- Bekkering, H., A. Wohlschläger, and M. Gattis (2000). Imitation of gestures in children is goal-directed. *The Quarterly Journal of Experimental Psychology Section A* 53(1), pp. 153–164.
- Bonaiuto, J. and M. Arbib (2010). Extending the mirror neuron system model, II: what did I just do? A new role for mirror neurons. *Biological Cybernetics* 102, pp. 341–359.
- Bonaiuto, J., E. Rosta, and M. Arbib (2007). Extending the mirror neuron system model, I: Audible actions and invisible grasps. *Biological Cybernetics* 96, pp. 9–38.
- Bonini, L., S. Rozzi, F. Serventi, L. Simone, P. Ferrari, and L. Fogassi (2010). Ventral premotor and inferior parietal cortices make distinct contribution to action organization and intention understanding. *Cerebral Cortex* 20(6), pp. 1372–1385.
- Borghi, A., C. Gianelli, and C. Scorolli (2010). Sentence Comprehension: Effectors and Goals, Self and Others An Overview of Experiments and Implications for Robotics. *Frontiers in Neurorobotics* 4. doi: 103389/fnbot201000003.
- Borra, E. and K. S. Rockland (2011). Projections to early visual areas v1 and v2 in the calcarine fissure from parietal association areas in the macaque. *Frontiers in neuroanatomy* 5. doi: 103389/fnana201100035.
- Brown, R. (1957). Linguistic determinism and the part of speech. *Journal of Abnormal and Social Psychology* 55(1), pp. 1–5.
- Caggiano, V., L. Fogassi, G. Rizzolatti, J. K. Pomper, P. Thier, M. A. Giese, and A. Casile (2011). View-based encoding of actions in mirror neurons of area f5 in macaque premotor cortex. *Current Biology* 21(2), pp. 144–148.
- Caggiano, V., L. Fogassi, G. Rizzolatti, P. Thier, and A. Casile (2009). Mirror neurons differentially encode the peripersonal and extrapersonal space of monkeys. *Science* 324(5925), pp. 403–406.
- Cangelosi, A. and T. Riga (2006). An embodied model for sensorimotor grounding and grounding transfer: Experiments with epigenetic robots. *Cognitive Science* 30(4), pp. 673–689.
- Cangelosi, A., V. Tikhanoff, J. Fontanari, and E. Hourdakis (2007). Integrating language and cognition: A cognitive robotics approach. *Computational Intelligence Magazine* 2(3), pp. 65–70.
- Carey, S. (1978). The child as word learner. *Linguistic Theory and Psychological Reality* .

- Carpenter, W. (1852). On the influence of suggestion in modifying and directing muscular movement, independently of volition. In *Proceedings of the Royal Institution of Great Britain*, Volume 1, pp. 147–153.
- Catmur, C., V. Walsh, and C. Heyes (2007). Sensorimotor learning configures the human mirror system. *Current Biology* 17(17), pp. 1527–31.
- Chaminade, T., E. Oztop, C. Gordon, and M. Kawato (2008). From self-observation to imitation: visuomotor association on a robotic hand. *Brain Research Bulletin* 75(6), pp. 775–784.
- Chersi, F., S. Thill, T. Ziemke, and A. Borghi (2010). Sentence processing: Linking language to motor chains. *Frontiers in Neurobotics* 4. 103389/fnbot201000004.
- Chomsky, N. (1966). *Topics in the theory of generative grammar*. Paris: Mouton De Gruyter.
- Clark, E. (1987). The principle of contrast: A constraint on language acquisition. *Mechanisms of language acquisition* , pp. 1–33.
- Cohen-Seat, G., H. Gastaut, J. Faure, and G. Heuyer (1954). Etudes expérimentales de l'activité nerveuse pendant la projection cinématographique. *Réflexions sur le Congrès International de Filmologie* 5, pp. 7–64.
- Craigheo, L., A. Bello, L. Fadiga, and G. Rizzolatti (2002). Hand action preparation influences the responses to hand pictures. *Neuropsychologia* 40(5), pp. 492–502.
- Cross, E., A. Hamilton, and S. Grafton (2006). Building a motor simulation de novo: observation of dance by dancers. *Neuroimage* 31(3), pp. 1257–1267.
- Cross, E., D. Kraemer, A. Hamilton, W. Kelley, and S. Grafton (2009). Sensitivity of the action observation network to physical and observational learning. *Cerebral Cortex* 19(2), pp. 315.
- Dayan, P. and L. Abbott (2001). *Theoretical neuroscience: Computational and mathematical modeling of neural systems*. Cambridge, MA: The MIT Press.
- Doya, K. (2000). Reinforcement learning in continuous time and space. *Neural Computation* 12(1), pp. 219–245.
- Ehrsson, H., S. Geyer, and E. Naito (2003). Imagery of voluntary movement of fingers, toes, and tongue activates corresponding body-part-specific motor representations. *Journal of Neurophysiology* 90(5), pp. 3304–3316.
- Elman, J. (1990). Finding structure in time. *Cognitive science* 14(2), pp. 179–211.
- Fadiga, L., L. Fogassi, G. Pavesi, and G. Rizzolatti (1995). Motor facilitation during action observation: a magnetic stimulation study. *Journal of Neurophysiology* 73(6), pp. 2608–2611.
- Fagg, A. (1996). *A Computational Model of The Cortical Mechanisms Involved in Primate Grasping*. Ph. D. thesis, University of Southern California, Computer Science

Department.

- Fagg, A. and M. Arbib (1998). Modeling parietal-premotor interactions in primate control of grasping. *Neural Networks* 11, pp. 1277–1303.
- Farkaš, I., T. Malík, and K. Rebrová (2012). Grounding the meanings in sensorimotor behavior using reinforcement learning. *Frontiers in Neurobotics* 6(1). doi: 103389/fnbot201200001.
- Farkaš, I., M. Malý, and K. Rebrová (2011a). Mirror neurons – theoretical and computational issues. Technical report, (TR-2011-28) Comenius University in Bratislava.
- Farkaš, I., M. Malý, and K. Rebrová (2011b). Porozumenie motorickým akciám – hypotéza kontinua. pp. 61–68. Opava, ČR: Slezská univerzita v Opavě.
- Farkaš, I. and K. Rebrová (2013). Bidirectional activation-based neural network learning algorithm. In *Proceedings of the International Conference on Artificial Neural Networks (ICANN), Sofia, Bulgaria*.
- Feldman, J. (2006). *From molecule to metaphor: A neural theory of language*. MIT Press.
- Feldman, J. and S. Narayanan (2004). Embodied meaning in a neural theory of language. *Brain and Language* 89(2), pp. 385–392.
- Flach, R., G. Knoblich, and W. Prinz (2003). Off-line authorship effects in action perception. *Brain and Cognition* 53(3), pp. 503–513.
- Fogassi, L., P. Ferrari, B. Gesierich, S. Rozzi, F. Chersi, and G. Rizzolatti (2005). Parietal lobe: from action organization to intention understanding. *Science* 308(5722), pp. 662–667.
- Fourneret, P. and M. Jeannerod (1998). Limited conscious monitoring of motor performance in normal subjects. *Neuropsychologia* 36(11), pp. 1133–1140.
- Frischen, A., D. Loach, and S. Tipper (2009). Seeing the world through another person’s eyes: Simulating selective attention via action observation. *Cognition* 111, pp. 212–218.
- Friston, K. (2003). Learning and inference in the brain. *Neural Networks* 16, pp. 1325–1352.
- Gallese, V., L. Fadiga, L. Fogassi, and G. Rizzolatti (1996). Action recognition in the premotor cortex. *Brain: A Journal of Neurology* 119, pp. 593–609.
- Gallese, V., C. Keysers, and G. Rizzolatti (2004). A unifying view of the basis of social cognition. *Trends in Cognitive Sciences* 8(9), pp. 396–403.
- Gallese, V. and G. Lakoff (2005). The brain’s concepts: The role of the sensory-motor system in conceptual knowledge. *Cognitive neuropsychology* 22(3-4), pp. 455–479.
- Gastaut, H. and J. Bert (1954). EEG changes during cinematographic presentation; moving picture activation of the EEG. *Electroencephalography and Clinical Neurophysiology* 6(3), pp. 433–444.

- Gazzola, V., H. van der Worp, T. Mulder, B. Wicker, G. Rizzolatti, and C. Keysers (2007). Aphasics born without hands mirror the goal of hand actions with their feet. *Current Biology* 17(14), pp. 1235–1240.
- Georgieff, N. and M. Jeannerod (1998). Beyond consciousness of external reality: A. *Consciousness and Cognition* 7(3), pp. 465–477.
- Gibson, J. (1977). The Theory of Affordances. *Perceiving, Acting, and Knowing: Toward an Ecological Psychology*, pp. 67–82.
- Glenberg, A. and M. Kaschak (2002). Grounding language in action. *Psychonomic Bulletin & Review* 9(3), pp. 558.
- Glenberg, A., M. Sato, L. Cattaneo, L. Riggio, D. Palumbo, and G. Buccino (2008). Processing abstract language modulates motor system activity. *The Quarterly Journal of Experimental Psychology* 61(6), pp. 905–919.
- Gold, K., M. Doniec, C. Crick, and B. Scassellati (2009). Robotic vocabulary building using extension inference and implicit contrast. *Artificial Intelligence* 173(1), pp. 145–166.
- Greenwald, A. (1970). Sensory feedback mechanisms in performance control: With special reference to the ideomotor mechanism. *Psychological Review* 77(2), pp. 73–99.
- Grezes, J., J. Armony, J. Rowe, and R. Passingham (2003). Activations related to “mirror” and “canonical” neurones in the human brain: an fMRI study. *Neuroimage* 18(4), pp. 928–937.
- Harless, E. (1861). Der Apparat des Willens. *Zeitschrift für Philosophie und philosophische Kritik* 38, pp. 50–73.
- Hauk, O., I. Johnsrude, and F. Pulvermüller (2004). Somatotopic representation of action words in human motor and premotor cortex. *Neuron* 41(2), pp. 301–307.
- Haykin, S. (2007). *Neural networks: A comprehensive foundation* (2nd ed).
- Heider[Rosch], E. (1972). Universals in color naming and memory. *Journal of Experimental Psychology* 93(1), pp. 10–20.
- Herbart, J. (1816). *Lehrbuch zur Psychologie*. Königsberg, Germany: Unzer.
- Herbart, J. (1825). *Psychologie als Wissenschaft neu gegründet auf Erfahrung, Metaphysik und Mathematik Zweiter, analytischer Teil*. Königsberg, Germany: Unzer.
- Heyes, C. (2001). Causes and consequences of imitation. *Trends in Cognitive Sciences* 5(6), pp. 253–261.
- Heyes, C. (2010). Where do mirror neurons come from? *Neuroscience and Biobehavioral Reviews* 34(4), pp. 575–83.
- Hickok, G. (2008, July). Eight problems for the mirror neuron theory of action understanding in monkeys and humans. *Journal of Cognitive Neuroscience* 21(7), pp. 1229–43.

- Hickok, G. and M. Hauser (2010). (Mis)understanding mirror neurons. *Current Biology* 20(14), pp. R593–4.
- Hinton, G. E. and J. L. McClelland (1988). Learning representations by recirculation. In *Neural Information Processing Systems*, pp. 358–366. New York: American Institute of Physics.
- Hommel, B., J. Müsseler, G. Aschersleben, and W. Prinz (2001). The theory of event coding (TEC): A framework for perception and action planning. *Behavioral and Brain Sciences* 24(05), pp. 849–878.
- Hopfield, J. (1982). Neural networks and physical systems with emergent collective computational properties. *Proceedings of the National Academy of Sciences* 79, pp. 2554–2558.
- Hurley, S. (2008). The shared circuits model (SCM): How control, mirroring, and simulation can enable imitation, deliberation, and mindreading. *Behavioral and Brain Sciences* 31(01), pp. 1–22.
- Iacoboni, M. (2009). Imitation, empathy and mirror neurons. *Annual Review of Psychology* 60, pp. 653–670.
- Iacoboni, M., L. Koski, M. Brass, H. Bekkering, R. Woods, M. Dubeau, J. Mazziotta, and G. Rizzolatti (2001). Re-afferent copies of imitated actions in the right superior temporal cortex. *Proceedings of the National Academy of Sciences* 98, pp. 13995–9.
- Ishida, H., K. Nakajima, M. Inase, and A. Murata (2010). Shared mapping of own and others' bodies in visuotactile bimodal area of monkey parietal cortex. *Journal of Cognitive Neuroscience* 22(1), pp. 83–96.
- James, W. (1890). *The Principles of Psychology, Vols I, II*. Cambridge, MA: Harvard University Press.
- Jeannerod, M. (2001). Neural simulation of action: a unifying mechanism for motor cognition. *Neuroimage* 14(1), pp. S103–S109.
- Jeannerod, M., M. A. Arbib, G. Rizzolatti, and H. Sakata (1995). Grasping objects: the cortical mechanisms of visuomotor transformation. *Trends in Neurosciences* 18(7), pp. 314–320.
- Jellema, T. and D. Perrett (2006). Neural representations of perceived bodily actions using a categorical frame of reference. *Neuropsychologia* 44, pp. 1535–1546.
- Jordan, M. (1986). Attractor dynamics and parallelism in a connectionist sequential network. In *Proceedings of the Eighth Annual Conference of the Cognitive Science Society*.
- Kawato, M. (1999). Internal models for motor control and trajectory planning. *Current Opinion in Neurobiology* 9(6), pp. 718–727.
- Keysers, C. and D. Perrett (2004). Demystifying social cognition: a Hebbian perspective. *Trends in Cognitive Sciences* 8(11), pp. 501–507.

- Kilner, J., K. Friston, and C. Frith (2007). The mirror-neuron system: a Bayesian perspective. *NeuroReport* 18(6), pp. 619–623.
- Knoblich, G. and R. Flach (2001). Predicting the effects of actions: Interactions of perception and action. *Psychological Science* 12(6), pp. 467–472.
- Knott, A. (2012). *Sensorimotor Cognition and Natural Language Syntax*. MIT Press Cambridge, MA.
- Köhler, E., C. Keysers, M. A. Umiltà, L. Fogassi, V. Gallese, and G. Rizzolatti (2002). Hearing sounds, understanding actions: action representation in mirror neurons. *Science* 297(5582), pp. 846–848.
- Kohonen, T. (1997). *Self-organizing Maps*. Springer.
- Kraskov, A., N. Dancause, M. Quallo, S. Shepherd, and R. Lemon (2009). Corticospinal neurons in macaque ventral premotor cortex with mirror properties: a potential mechanism for action suppression? *Neuron* 64(6), pp. 922–930.
- Laycock, T. (1845). *On the Reflex Function of the Brain*. London, England: Adlard.
- Lotze, R. (1852). *Medizinische Psychologie oder Physiologie der Seele*. Weidmann'sche Buchhandlung.
- Mahon, B. and A. Caramazza (2005). The orchestration of the sensory-motor systems: Clues from neuropsychology. *Cognitive Neuropsychology* 22(3/4), pp. 480–494.
- Maly, M. (2013). *Reinforcement Learning with Abstraction*. Ph. D. thesis, Comenius University in Bratislava.
- Marocco, D., A. Cangelosi, K. Fischer, and T. Belpaeme (2010). Grounding Action Words in the Sensorimotor Interaction with the World: Experiments with a Simulated iCub Humanoid Robot. *Frontiers in Neurorobotics* 4. doi: 103389/fnbot201000007.
- Mason, C., J. Gomez, and T. Ebner (2001). Hand synergies during reach-to-grasp. *Journal of Neurophysiology* 86(6), pp. 2896–2910.
- Mel, B. and C. Koch (1990). Sigma-pi learning: on radial basis functions and cortical associative learning. In D. Touretzky (Ed.), *Advances in NIPS-2*, pp. 474–481. The MIT Press.
- Metta, G., P. Fitzpatrick, and L. Natale (2006). YARP: yet another robot platform. *International Journal on Advanced Robotics Systems* 3(1), pp. 43–48.
- Metta, G., G. Sandini, D. Vernon, L. Natale, and F. Nori (2008). The iCub humanoid robot: an open platform for research in embodied cognition. In *Proceedings of the 8th Workshop on Performance Metrics for Intelligent Systems*, pp. 50–56. ACM.
- Miall, M. (2003). Connecting mirror neurons and forward models. *NeuroReport* 14(17), pp. 2135–2137.
- Molenberghs, P., R. Cunnington, and J. B. Mattingley (2012). Brain regions with mirror properties: a meta-analysis of 125 human fmri studies. *Neuroscience & Biobehavioral*

- Reviews* 36(1), pp. 341–349.
- Mukamel, R., A. Ekstrom, J. Kaplan, M. Iacoboni, and I. Fried (2010). Single-neuron responses in humans during execution and observation of actions. *Current Biology* 20(8), pp. 750–756.
- Murata, T. (1989). Petri nets: Properties, analysis and applications. *Proceedings of the IEEE* 77(4), pp. 541–580.
- Nelissen, K., E. Borra, M. Gerbella, S. Rozzi, G. Luppino, W. Vanduffel, G. Rizzolatti, and G. A. Orban (2011). Action observation circuits in the macaque monkey cortex. *Journal of Neuroscience* 31(10), pp. 3743–3756.
- Oberman, L. and V. Ramachandran (2008). How do shared circuits develop? *Behavioral and Brain Sciences: Open Peer Commentary* 31, pp. 34–35.
- Oberman, L. M. and V. S. Ramachandran (2007). The simulating social mind: the role of the mirror neuron system and simulation in the social and communicative deficits of autism spectrum disorders. *Psychological bulletin* 133(2), pp. 310–327.
- O’Grady, W. and W. O’Grady (2005). *How children learn language*. Cambridge University Press.
- Omohundro, S. (1993). Best-first model merging for dynamic learning and recognition. *Advances in Neural Information Processing Systems* , pp. 958–958.
- O’Reilly, R. (1996a). Biologically plausible error-driven learning using local activation differences: The generalized recirculation algorithm. *Neural Computation* 8(5), pp. 895–938.
- O’Reilly, R. (1996b). *The Leabra model of neural interactions and learning in the neocortex*. Ph. D. thesis, Carnegie Mellon University.
- O’Reilly, R. and Y. Munakata (2000). *Computational Explorations in Cognitive Neuroscience: Understanding the Mind by Simulating the Brain*. Cambridge, MA: The MIT Press.
- Oztop, E. and M. Arbib (2002). Schema design and implementation of the grasp-related mirror neuron system. *Biological Cybernetics* 87, pp. 116–140.
- Oztop, E., M. Kawato, and M. Arbib (2006). Mirror neurons and imitation: A computationally guided review. *Neural Networks* 19(3), pp. 254–271.
- Oztop, E., M. Kawato, and M. Arbib (2013). Mirror neurons: Functions, mechanisms and models. *Neuroscience Letters* .
- Oztop, E., D. Wolpert, and M. Kawato (2005). Mental state inference using visual control parameters. *Cognitive Brain Research* 22, pp. 129–151.
- Pecháč, M. (2013). Self-organization of sensorimotor representations and its use in grasping of objects (in Slovak). Master’s thesis, Comenius University in Bratislava.

- Peeters, R., L. Simone, K. Nelissen, M. Fabbri-Destro, W. Vanduffel, G. Rizzolatti, and G. Orban (2009). The representation of tool use in humans and monkeys: common and uniquely human features. *Journal of Neuroscience* 29(37), pp. 11523.
- Pellegrino, G., L. Fadiga, L. Fogassi, V. Gallese, and G. Rizzolatti (1992). Understanding motor events: a neurophysiological study. *Experimental Brain Research* 91(1), pp. 176–180.
- Penfield, W. and T. Rasmussen (1950). *The Cerebral Cortex of Man*. Macmillan.
- Perrett, D., M. Harries, R. Bevan, S. Thomas, P. Benson, A. Mistlin, A. Chitty, J. Hietanen, and J. Ortega (1989). Frameworks of analysis for the neural representation of animate objects and actions. *Journal of Experimental Biology* 146(1), pp. 87–113.
- Perrett, D., M. Oram, M. Harries, R. Bevan, J. Hietanen, P. Benson, and S. Thomas (1991). Viewer-centred and object-centred coding of heads in the macaque temporal cortex. *Experimental Brain Research* 86(1), pp. 159–73.
- Pfeifer, R. and C. Scheier (1999). *Understanding Intelligence*. Cambridge, MA: The MIT Press.
- Pineda, F. (1987). Generalization of back-propagation to recurrent neural networks. *Physics Review Letters* 59(19), pp. 2229–2232.
- Pineda, J. (2005). The functional significance of mu rhythms: translating “seeing” and “hearing” into “doing”. *Brain Research Reviews* 50(1), pp. 57–68.
- Prinz, W. (1997). Perception and action planning. *European Journal of Cognitive Psychology* 9(2), pp. 129–154.
- Pulvermüller, F. (2005). Brain mechanisms linking language and action. *Nature Reviews Neuroscience* 6(7), pp. 576–582.
- Pulvermüller, F., M. Härle, and F. Hummel (2001). Walking or Talking?: Behavioral and Neurophysiological Correlates of Action Verb Processing. *Brain and Language* 78(2), pp. 143–168.
- Rebrová, K. (2012). Stelesnené porozumenie a ideomotorická teória. In J. Rybár (Ed.), *Kognitívne paradigmy*, pp. 127–150. Vydavateľstvo Európa.
- Rebrová, K., M. Pecháč, and I. Farkaš (2013). Towards a robotic model of the mirror neuron system. In *Proceedings of the 3rd joint IEEE International Conference on Development and Learning and on Epigenetic Robotics*. Osaka, Japan.
- Repp, B. and G. Knoblich (2007). Action can affect auditory perception. *Psychological Science* 18(1), pp. 6–7.
- Rizzolatti, G. and M. Arbib (1998). Language within our grasp. *Trends in neurosciences* 21(5), pp. 188–194.
- Rizzolatti, G., R. Camarda, L. Fogassi, M. Gentilucci, G. Luppino, and M. Matelli (1988). Functional organization of inferior area 6 in the macaque monkey. *Experimental Brain*

- Research* 71(3), pp. 491–507.
- Rizzolatti, G. and L. Craighero (2004). The mirror-neuron system. *Annual Review of Neuroscience* 27, pp. 169–92.
- Rizzolatti, G., M. Fabbri-Destro, and L. Cattaneo (2009). Mirror neurons and their clinical relevance. *Nature Clinical Practice Neurology* 5(1), pp. 24–34.
- Rizzolatti, G., L. Fadiga, V. Gallese, and L. Fogassi (1996). Premotor cortex and the recognition of motor actions. *Cognitive Brain Research* 3(2), pp. 131–141.
- Rizzolatti, G., L. Fogassi, and V. Gallese (2001). Neurophysiological mechanisms underlying the understanding and imitation of action. *Nature Reviews Neuroscience* 2, pp. 661–670.
- Rizzolatti, G. and C. Sinigaglia (2010). The functional role of the parieto-frontal mirror circuit: interpretations and misinterpretations. *Nature Reviews Neuroscience* 11(4), pp. 264–74.
- Rumelhart, D., G. Hinton, and R. Williams (1986). *Learning internal representations by error propagation*, pp. 318–362. Number 1. Cambridge, MA: The MIT Press.
- Sato, A. and A. Yasuda (2005). Illusion of sense of self-agency: discrepancy between the predicted and actual sensory consequences of actions modulates the sense of self-agency, but not the sense of self-ownership. *Cognition* 94(3), pp. 241–255.
- Sebanz, N. and M. Shiffrar (2009). Detecting deception in a bluffing body: The role of expertise. *Psychonomic Bulletin & Review* 16(1), pp. 170–175.
- Shepherd, S., J. Klein, R. Deaner, and M. Platt (2009). Mirroring of attention by neurons in macaque parietal cortex. *Proceedings of the National Academy of Sciences* 106(23), pp. 9489–9494.
- Shin, Y., R. Proctor, and E. Capaldi (2010). A Review of Contemporary Ideomotor Theory. *Psychological Bulletin* 136(6), pp. 943–974.
- Skinner, B. (1957). *Verbal Behavior*. New York: Appleton-Century-Crofts.
- Stock, A. and C. Stock (2004). A short history of ideo-motor action. *Psychological Research* 68(2), pp. 176–188.
- Strickert, M. and B. Hammer (2005). Merge SOM for temporal data. *Neurocomputing* 64, pp. 39–71.
- Sugita, Y. and J. Tani (2005). Learning semantic combinatoriality from the interaction between linguistic and behavioral processes. *Adaptive Behavior* 13(1), pp. 33–52.
- Sugita, Y. and J. Tani (2008). A sub-symbolic process underlying the usage-based acquisition of a compositional representation: Results of robotic learning experiments of goal-directed actions. In *Proceeding of 7th IEEE International Conference on Development and Learning ICDL 2008*, pp. 127–132. IEEE.
- Sutton, R. and A. Barto (1998). *Reinforcement Learning: An Introduction*. MIT Press.

- Tani, J. (2003). Learning to generate articulated behavior through the bottom-up and the top-down interaction processes. *Neural Networks* 16(1), pp. 11–23.
- Tani, J. and M. Ito (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), pp. 481–488.
- Tani, J., M. Ito, and Y. Sugita (2004). Self-organization of distributedly represented multiple behavior schemata in a mirror system: reviews of robot experiments using RNNPB. *Neural Networks* 17(8-9), pp. 1273–1289.
- Tessitore, G., R. Prevede, E. Catanzariti, and G. Tamburrini (2010). From motor to sensory processing in mirror neuron computational modelling. *Biological Cybernetics* 103(6), pp. 471–485.
- Thivierge, J.-P. and G. Marcus (2007). The topographic brain: from neural connectivity to cognition. *Trends in Neurosciences* 30(6), pp. 251–259.
- Thorndike, E. (1913). Ideo-motor action. *Psychological Review* 20, pp. 91–106.
- Tikhanoff, V., A. Cangelosi, P. Fitzpatrick, G. Metta, L. Natale, and F. Nori (2008). An open-source simulator for cognitive robotics research: The prototype of the icub humanoid robot simulator. In *Proceedings of the 8th Workshop on Performance Metrics for Intelligent Systems*, pp. 57–61. ACM.
- Tikhanoff, V., A. Cangelosi, and G. Metta (2011). Integration of speech and action in humanoid robots: icub simulation experiments. *IEEE Transactions on Autonomous Mental Development* 3(1), pp. 17–29.
- Tkach, D., J. Reimer, and N. G. Hatsopoulos (2007). Congruent activity during action and action observation in motor cortex. *Journal of Neuroscience* 27(48), pp. 13241–50.
- Tomasello, M. (2003). *Constructing a language: A usage-based theory of language acquisition*. Harvard University Press.
- Umiltà, M., E. Kohler, V. Gallese, L. Fogassi, L. Fadiga, C. Keysers, G. Rizzolatti, et al. (2001). I know what you are doing: A neurophysiological study. *Neuron* 31(1), pp. 155–166.
- Umiltà, M. a., L. Escola, I. Intskirveli, F. Grammont, M. Rochat, F. Caruana, a. Jezzini, V. Gallese, and G. Rizzolatti (2008, February). When pliers become fingers in the monkey motor system. *Proceedings of the National Academy of Sciences of the United States of America* 105(6), pp. 2209–13.
- Van der Wel, R., N. Sebanz, and G. Knoblich (2013). Action Perception From a Common Coding Perspective. In K. Johnson and M. Shiffrar (Eds.), *People Watching: Social, Perceptual, and Neurophysiological Studies of Body Perception*, pp. 101. Oxford University Press.
- Van Elk, M., H. Van Schie, S. Hunnius, C. Vesper, and H. Bekkering (2008). You’ll never crawl alone: neurophysiological evidence for experience-dependent motor resonance in

- infancy. *Neuroimage* 43(4), pp. 808–814.
- van Hasselt, H. (2012). Reinforcement learning in continuous state and action spaces. *Reinforcement Learning*, pp. 207–251.
- Van Hasselt, H. and M. Wiering (2007). Reinforcement learning in continuous action spaces. In *IEEE International Symposium on Approximate Dynamic Programming and Reinforcement Learning (ADPRL)*, pp. 272–279.
- van Hasselt, H. and M. A. Wiering (2009). Using continuous action spaces to solve discrete problems. In *Neural Networks, 2009 IJCNN 2009 International Joint Conference on*, pp. 1149–1156. IEEE.
- Vančo, P. and I. Farkaš (2010). Experimental comparison of recursive self-organizing maps for processing tree-structured data. *Neurocomputing* 73(7-9), pp. 1362–1375.
- Wermter, S. and M. Elshaw (2003). Learning robot actions based on self-organising language memory. *Neural Networks* 16(5-6), pp. 691–699.
- Wiedermann, J. (2003). Mirror neurons, embodied cognitive agents and imitation learning. *Computing and Informatics* 22(6), pp. 545–559.
- Wiedermann, J. (2009). A high level model of a conscious embodied agent. In *Proceeding of 8th IEEE International Conference on Cognitive Informatics ICCI'09*, pp. 448–456. IEEE.
- Williams, R. J. and D. Zipser (1995). Gradient-based learning algorithms for recurrent networks and their computational complexity. *Back-propagation: Theory, architectures and applications*, pp. 433–486.
- Wolpert, D., K. Doya, and M. Kawato (2003). A unifying computational framework for motor control and social interaction. *Philosophical Transactions of the Royal Society of London B* 358, pp. 593–602.
- Wolpert, D. and M. Kawato (1998). Multiple paired forward and inverse models for motor control. *Neural Networks* 11(7-8), pp. 1317–1329.
- Zdechovan, L. (2012). Modeling the object grasping using the neural networks and icub robotic simulator (in Slovak). Master's thesis, Comenius University in Bratislava.
- Zhong, J., C. Weber, and S. Wermter (2011). Robot trajectory prediction and recognition based on a computational mirror neurons model. In T. Honkela, W. Duch, M. Girolami, and S. Kaski (Eds.), *Proceedings of the 21st International Conference on Artificial Neural Networks ICANN 2011*, Volume 2, pp. 333–340. Espoo, Finland: Springer.
- Zwaan, R. and L. Taylor (2006). Seeing, acting, understanding: Motor resonance in language comprehension. *Journal of Experimental Psychology-General* 135(1), pp. 1–11.