

GETTING TUNED TO OTHERS' INTENTIONS

A model for hierarchical spatio-temporal pattern recognition

Masters Thesis

of

Floran Stuijt

June 20, 2010

Supervisors:

Dr. W.F.G. Haselager

Dr. I. van Rooij

ing. M. Vlietstra

Radboud University Nijmegen



Faculty of Social Sciences
Institute of Artificial Intelligence

In collaboration with



Contents

Summary	vi
Acknowledgements	vii
1 Introduction	1
1.1 The Social Skill of Mindreading	1
1.2 An Argument for Hierarchical Organization	3
2 Towards an Improved Model	7
2.1 Foundations	7
2.2 Formal Description	9
2.3 Proposal for Learning Method	10
2.4 Selection of a Prediction Mechanism	11
3 Simulation Experiments	13
3.1 Implementation Specific Details	13
3.2 Experiment 1: Construction of a Pattern Detection Filter	14
3.3 Experiment 2: Detecting Hierarchical Organized Sequences	15
3.4 Experiment 3: Goal recognition in a joint-dial experiment	17
4 Discussion & Conclusion	21
4.1 Discussion	21
4.1.1 Implementation Specific Aspects	21
4.1.2 Fundamental Aspects	22
4.1.3 Relation to Mirror Neurons	24
4.1.4 Biological Evidence for Existence of Prediction Errors	25
4.2 Conclusion	25
A MATLAB Code Listing	31

Summary

The human ability to understand other's actions and attribute mental states to those actions plays a prominent role in social interaction between humans. An understanding of this mindreading ability may result in new techniques which for example can be applied to the detection of threatening or abusive behavior in surveillance applications, or can be used to take human-robot interaction to a higher level.

It is believed that visual perception in the human makes use of hierarchical decomposition, an idea that has gained significant popularity since the findings of simple and complex cells in the visual cortex of cats (Hubel and Wiesel, 1959) and macaque monkeys (Hubel and Wiesel, 1968). In this thesis it is investigated whether a previous approach to mental state inference, the Mental State Inference (MSI) model by Oztop et al. (2005), can be improved by adding hierarchical structure to the model.

A novel model is presented which is based on the detection of biological motion in order to infer the mental states of others. This is achieved by a series of so called complex 'tuning forks' each of is tuned resonate with a specific input signal. The behavior of these tuning forks is defined so that these forks 'resonate' during the time the pattern, to which the forks have been 'tuned' to, is present in the incoming signal. I.e., tuning forks behave as pattern recognizers, and their output is defined as the degree to which the input corresponds to or resonates with their intrinsic pattern. It is these pattern recognizers that form the basic building block of the presented model. The model is based on hierarchical decomposition since the output of the pattern recognizers within the model, referred to as responsibility signals, is further analyzed by another series of pattern recognizers. This additional level of analysis allows us to detect primitive biological motions and detect the order in which these occur, which in turn allows us to combine primitive behaviors to create a representation of more complex behavior, and more specifically allows us to detect hierarchical sequential ordered behavior.

In this study it is hypothesized that the introduction of the above-mentioned hierarchical decomposition principle can be used to infer the intentions of others, given a sufficient amount of pattern recognizer layers as described above. In order to verify whether such an approach is fruitful three experiments have been set up. In the first experiment it is shown that the model is able to detect patterns embedded in a random time series. The second experiment shows that the responsibility values computed in the first experiment can be used to detect a specific sequence of patterns, and at the same time shows that the model is not limited to pattern recognition in one dimensional time series but can also be used to detect patterns in multidimensional time series. The third experiment shows that the model is not only able to detect patterns in synthetically generated signals but can also be used to detect patterns in empirical data.

Accumulating evidence shows that 'action understanding may precede, rather than fol-

low from, action mirroring' (Csibra, 2007). In line with this argument, the model presented in this thesis demonstrates that mental state inference is possible without the need of mental simulation, and is therefore compatible with recent findings as part of the functional role of the mirror neuron system and action understanding in general.

Acknowledgements

After 9 months of keeping my nose to the grindstone it is finally time to thank all the people that were involved in my attempt to reach my academic objectives. Spending almost a year on writing a thesis in combination with a 16 hour job is difficult but thanks to the people who helped me I was able to keep on going.

As this thesis was written during an internship at Logica, I'm happy to thank Logica for providing me the resources necessary to perform the study. The Working Tomorrow program in which I was allowed to participate was both inspiring and enjoyable at the same time. In particular, I would like to thank project leader Martijn Vlietstra for encouraging me and convincing me that I was doing my job properly. I also would like to thank all students who were graduating in the same period as I did, for taking my mind off things when I could use some distraction. I hereby wish them good luck with their future career.

The experimental part of my research could have taken a much longer time if it was not Ruud Meulenbroek and Oliver Herbort who were so kind to provide me with the data required to perform one of the simulation experiments, and significantly reduced the time required to get obtain the results presented in this thesis.

Furthermore, thanks go out to Iris van Rooij and Pim Haselager for their excellent supervision. Every progress meeting gave me a huge boost of energy to continue the work on my research.

And finally, I would like to thank Boukje for her patience and for taken the time listening to my moaning on the difficulties I was experiencing during my research.

Introduction

1.1. The Social Skill of Mindreading

As technology in robotics advances, the field of human-robot interaction is gaining an increasing research interest. Researchers in this field are concerned with the development of methods which improve interaction between humans and robots, and one way to improve this interaction is by making it more similar to human-human interaction. Human-human interaction is characterized by its extensive use of social skills, including for instance emotion detection, action understanding, and learning by demonstration, and the development of methods allowing one to equip a robot with these skills is considered as one of the biggest challenges in this field.

Understanding other people's minds is a social skill we humans are remarkably good at. To illustrate this, consider for example a person sitting in front of a glass of water on a table. Imagine this person grasping the glass with the palm of his hand faced against the side of the glass. We immediately infer that this person is probably thirsty and wants to drink. On the other hand, if the person grasps the glass with the palm of his hand on top of the glass, we expect the person to move the glass to another location on the table. Understanding these differences is considered part of a 'theory of mind', and providing a robot the capability to read other's minds can pave the way for taking human-robot interaction to a higher level. Other fields can benefit from this ability as well, consider for example an application which automatically detects abusive or threatening behavior in surveillance applications. Alternatively, this ability may also be exploited by entertainment industries, for example, by adjusting a computer game to the current intentional state of the game player in order to improve the gaming experience.

Another application for reading other's minds can be found in imitation learning by robots. Significant progress has been made in this field (for an overview, see Breazeal and Scassellati, 2002), however, it is still an open question how we can give a robot the ability to determine whether an observed action is relevant to a specific task. Clearly, imitating everything an observer observes is not useful. Obviously, imitation is most valuable when the observed movements coincide with the goals of the observer. In that case, the observed actions are presumably relevant to the achievement of these goals, and a straightforward solution would be to associate these movements with the current intentional state of the robot. For example, if a robot 'intends' that a glass should be emptied, and a demonstrator performs the actions required to empty that same glass, the robot may learn how to achieve his intentions by mimicking the observed movements later on.

Earlier work in the field of artificial intelligence mainly focused on providing a solution

to the action understanding problem, also known as plan and intent recognition, by taking a symbolic approach. This work was centered around the principle of rationality, the assumption that rational agents tend to achieve their desires as optimally as possible, given their beliefs (Baker et al., 2007). Consider for example the work by Appelt and Pollack (1992), in which a logical abduction method for plan recognition is described. Abduction can be seen as the logical equivalent of inference to the best explanation, or reasoning from effect to cause. However, the method described in this work is built on the assumption that information is either true or false, a notion which does not fit well in our inherently uncertain world. Moreover, all rules have to be defined by hand, for the reason that no learning mechanism is provided. Later work changed focus towards probabilistic models which are able to handle uncertainty in information. An example of such a model is the Bayesian model of plan recognition by Charniak and Goldman (1993).

A more recent and particular interesting approach to intention understanding is put forward by the mental state inference (MSI) model (Oztop et al., 2005). This biologically inspired model aims to provide a solution to this problem by selecting the mental state that is associated with the motor plan which best predicts the observed actions. This is accomplished by extracting a so called control variable from the sensory input, which is used for comparison with the outcome of the mental movement simulation which in turn results in a control variable prediction. In order to analyze this model, several simulation experiments have been conducted. In one of these experiments, the task consisted of inferring the goal of a demonstrated tool-use. For this simulation, three types of stimuli sequences were recorded, namely holding a hammer, driving a nail, and prying a nail with the hammer. Results show that the model is able to successfully infer (i.e. correctly classify) the underlying goals of these sequences even after limited exposure to each of these sequences.

Although simulation results are promising, the MSI model lacks at least three properties which I consider essential in practical situations. First, no description for a learning mechanism is provided, as motor plans are hardwired in the model. Second, for imitation learning the MSI model is impractical since it assumes that the observer has a series of planning modules in place before observed actions can be recognized, while obviously the purpose of imitation *is* to learn the motor plans associated with these same planning modules. Third, one motor controller is often insufficient to encode complex movements such as non-monotonic grasp movements (e.g. moving away from an object first before moving closer to the object). A similar argument holds for control variables; it is unlikely that one control variable is sufficient to encode the details necessary to infer the mental state of another person. To illustrate this, consider again the glass grasping example. In addition to the type of grasp, another important source of information is the amount of water in the glass. If the glass is empty and a precision grip is performed, the demonstrator does not want to drink. However, the amount of water in the glass may not always be visible. Information regarding the amount of water in the glass can increase the likelihood of correct intention inference, but it is not strictly required to do intention inference in the first place. This example illustrates that the information required to infer a mental state can be divided over multiple control variables, but the MSI model does not offer an explanation on how input from multiple control variables should be integrated.

In this thesis a model for intention inference based on the MSI model is presented (Section 2) and investigated. This model is built upon the hypothesis that introducing a hierarchical structure offers a natural way to overcome the above-mentioned issues with the MSI model, a notion which I will elaborate in more detail in Section 1.2. In order to validate

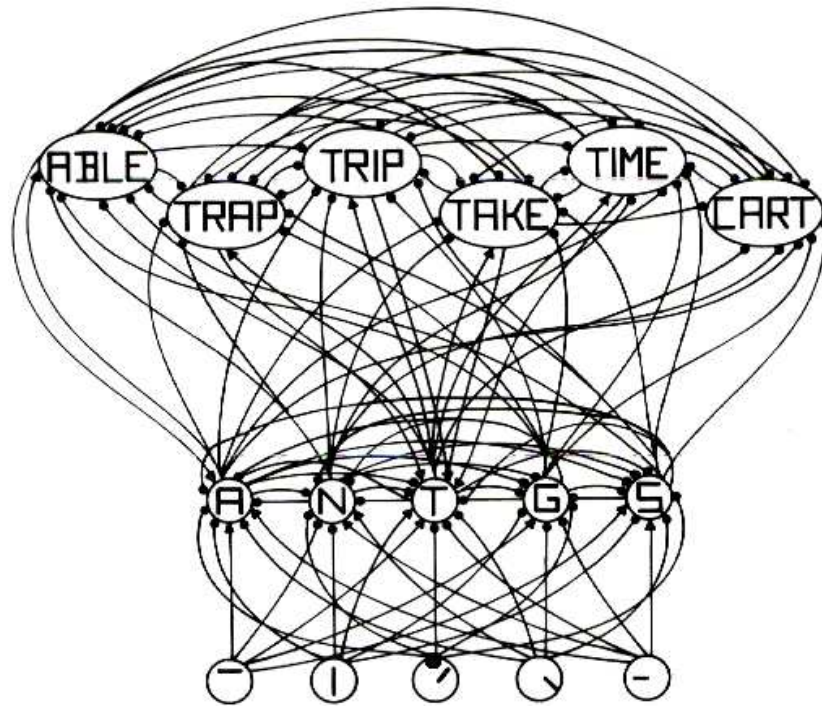


Figure 1.1: The interactive activation (IA) model for word recognition. This model is able to explain various phenomena in human reading, including the ability of recognizing words from degraded stimuli, and the word superiority effect (McClelland and Rumelhart, 2002).

this hypothesis three experiments have been set up, in which the goal of the first two experiments is to show that the model is capable of detecting patterns in synthetically generated multidimensional time series, and thus has the potential capacity of mental state inference by detecting patterns associated with a mental state and labeling them as such. The third experiment is conducted to show that usage of the model is not limited to synthetically generated data but can be applied to empirical data as well.

1.2. An Argument for Hierarchical Organization

For humans it is natural to structure their environment following a hierarchical organization. Our environment abounds in hierarchical structure, consider for example a house, consisting of multiple rooms, each of which might contain multiple pieces of furniture. Neurophysiological studies show that our brain might apply hierarchical decomposition in perception, a notion which is supported by the discovery of ‘simple cells’ and ‘complex cells’ in the visual cortex of cats (Hubel and Wiesel, 1959) and macaque monkeys (Hubel and Wiesel, 1968). In these studies it was demonstrated that some of these simple cells were behaving as edge detectors in visual stimuli, since these cells responded selectively to the presentation of edges with a specific orientation. Similar cells were found, responding to increasingly complex features, such as edges moving in a particular direction. It is believed that the behavior of these cells, termed complex cells, is built up by a proper combination of simple cells.

Hierarchical organization is not limited to our environment, but can also be found in spoken language. Sentences are constructed by combining words or morphemes, each of

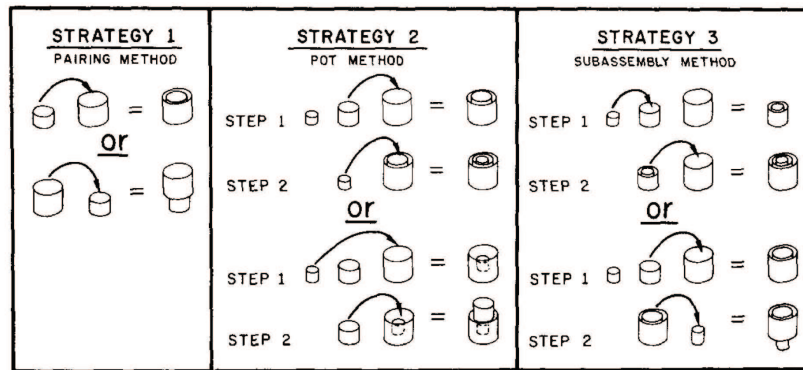


Figure 1.2: During their development infants follow an increasing progress in hierarchical complexity of their manipulative capabilities. In the first two strategies, chain-like sequences are employed to manipulate the cups. The third strategy which is later in development, employs the utilization of a subassembly to achieve the final nesting of the cups (Greenfield et al., 1972).

which in turn can be segmented into a series of phonemes. A similar organization can also be found in written language, and the interactive activation (IA) model (McClelland and Rumelhart, 2002), a computational model for word recognition, demonstrates this by having a structure similar to the one discovered in animal brains by Hubel and Wiesel (See Figure 1.1). The IA model is able to explain various robust phenomena in human reading, including the human ability of recognizing words from degraded stimuli, and the word superiority effect.

The findings by Hubel and Wiesel, and McClelland, both support the notion that the brain represents our inherently hierarchical environment in a hierarchical way. Hierarchical structure can also be found in behavior of human beings. It has, for example, been shown that during their development infants follow an increasing progress in hierarchical complexity of their manipulative capabilities. To illustrate this, consider the cup nesting strategies shown in Figure 1.2. This figure presents the order by which cup-nesting strategies develop in infants from the age of 11 months (Greenfield et al., 1972). In the first two strategies, chain-like sequences are employed to manipulate the cups. The third strategy, which develops later, makes use of a subassembly to achieve the final nesting of the cups. Contrary to the chain-like sequences, the third strategy shows hierarchically organized sequential behavior, as one cup is placed inside the other to form a subassembly, which is followed by combining the resulting subassembly with the third cup. It is believed that the development of such a strategy forms the basis for symbol combination and tool use, and that speech and the capacity of manual object combination both depend on a common neural substrate whose function is that of a supramodal hierarchical processor (Greenfield et al., 1991). The results of these studies exemplify the inherent hierarchical organization in human behavior, and it is likely that the brain uses hierarchical decomposition here as well in order to interpret and represent the behavior of human beings.

The IA model is built on the assumption that stimulus information is static, i.e. it is assumed that stimuli do not change over time. However, temporal information can be important as well, as it may convey a substantial amount of information. Johansson (1973) demonstrated this by showing that people can easily infer biological motion types from a

series of moving bright spots. Light emitting spots were placed on the major joints of a demonstrator's body, and the motion patterns of these spots were recorded during several motion types, such as walking and running. However, by randomly rearranging these bright spots relative to each other after recording, subjects become unable to attribute any biological motion to the moving spots. This demonstrates that the combination of both spatial and temporal information is essential in order to perform a meaningful analysis of biological motion.

The MSI model discussed earlier exploits temporal information by employing a predictive component. However, the model does not adopt a hierarchical decomposition strategy, which is, as emphasized above, believed to play a significant role in order to be able to recognize hierarchically organized sequential behavior. Furthermore, the model can be improved by enabling it to integrate information coming from multiple sensory inputs. In Section 2, a conceptual model with these improvements in mind is presented.

Towards an Improved Model

In this section the foundations of a novel conceptual model for mental state inference are presented. These foundations are based on the principle that, in addition to static features, hierarchical decomposition can also be applied to features changing over time, also referred to as dynamical features (Jaeger, 2007). This approach aims to provide an explanation on how information coming from multiple sensory inputs can be integrated such that the mental state which underlies this sensory information can be inferred from observation. To be more precise, it aims to provide an explanation on how we can *learn* to infer the mental state from the observation of a hierarchically organized behavior. In Section 2.1 a functional and informal description of the model is given, which is followed by a formal description of the model presented in Section 2.2.

2.1. Foundations

A schematic overview of the presented conceptual model is given in Figure 2.1. The model illustrated in this figure consists of multiple sets of ‘complex tuning forks’ or pattern recognizers. The behavior of these tuning forks is defined so that these forks ‘resonate’ during the time the pattern, to which the forks have been ‘tuned’ to, is present in the incoming signal. I.e., tuning forks behave as pattern recognizers, and their output is defined as the degree to which the input corresponds to or resonates with their intrinsic pattern. It is these pattern recognizers that form the basic building block of the presented model.

At the bottom of the model, each set of pattern recognizers receives input in the form of a feature, such as the distance from an end-effector to a target object or the orientation of an object. The output of a pattern recognizer will represent the conformity of the input to its intrinsic pattern. This output, referred to as a pattern recognizer’s *responsibility* or *vote*, indicates the intervals on which specific patterns have been present in the input signal, and does not contain any information about the characteristics of the patterns themselves when considered in isolation.

Ideas similar to those presented so far have also been proposed in the context of motor control, e.g. in a model by Wolpert and Kawato (1998), and in the MOSAIC model by Haruno et al. (2001). The models presented in these studies are primarily based on a number of ‘neural experts’, trained in order to become an expert at predicting input signals. The prediction error of each neural expert indicates if the pattern the expert has been trained on is currently present as input. However, these models do not apply hierarchical decomposition to the incoming signals. In fact, Haruno et al. (2001) address this issue in their study and propose to further analyze the MOSAIC model’s analogues of the responsibility signals

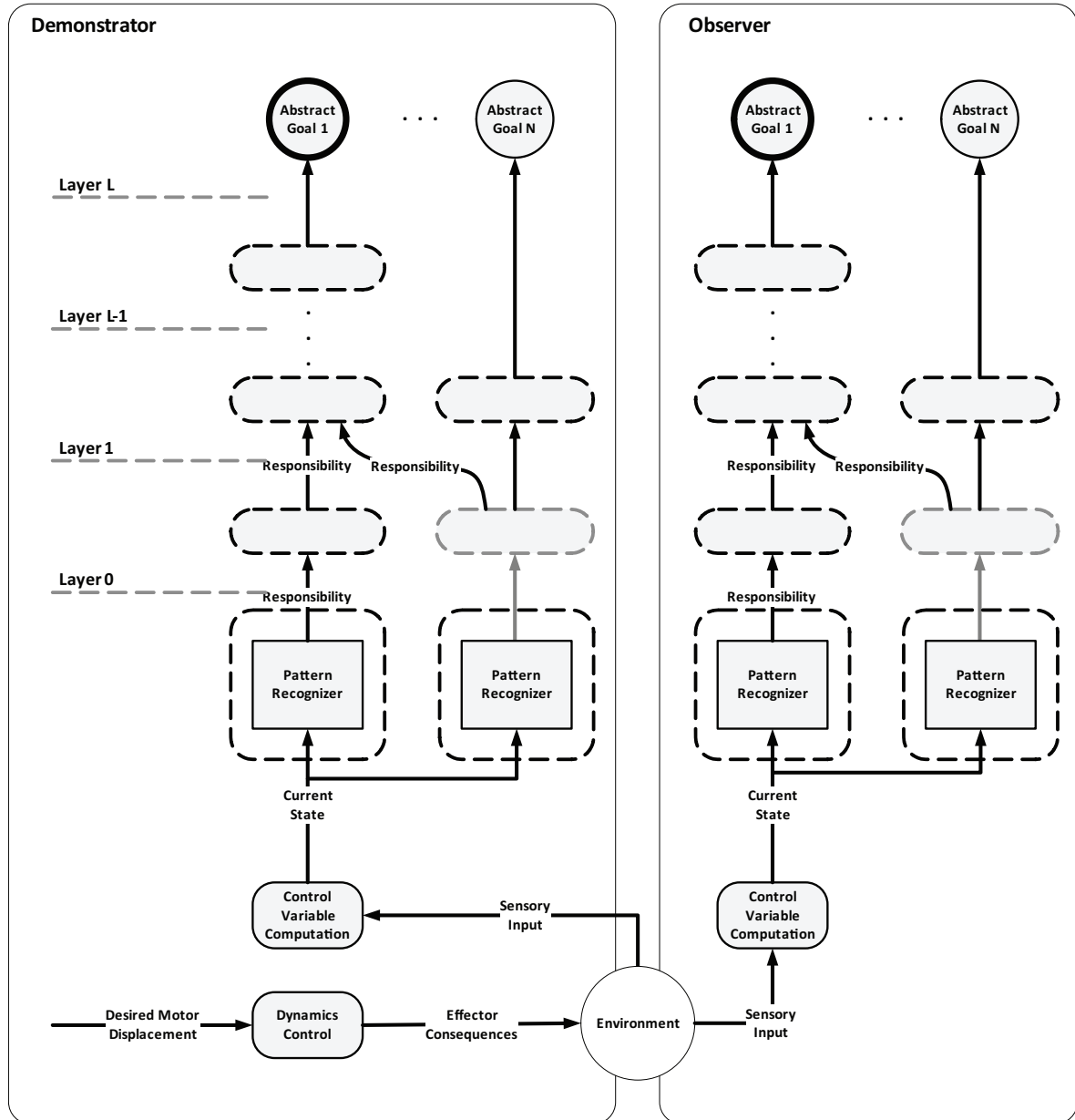


Figure 2.1: A conceptual model for goal recognition. See text for details.

by another series of experts.

At this point it has to be emphasized that features can change over time, and as a consequence computed responsibilities can also be regarded as time varying quantities. As such, these responsibilities can also be interpreted as dynamical features, and can, as proposed by Haruno et al. (2001), be put through another level of analysis. Therefore, in the presented model responsibility signals are sent to a higher layer, again consisting of multiple pattern recognizers, similar to the ones in the subordinate layer. This higher layer is responsible for detecting patterns in responsibility signals, and the output in turn consists of a more abstract description of its input. A similar strategy can be repeated for several times, and it is assumed that responsibility signals computed at a specific level can be regarded as abstract representations of the input signal. This brings us to the essence of the ideas outlined above: can we relate sufficiently abstract representations of the input signal to mental states, and can we use this relation to infer the mental state of others?

The approach presented above is based on the assumption that mental states are predefined in the model, and that their activity is representative for the mental state of the actor. This allows the association strength between mental states and the output of a specific level in the model to be modulated by coactivation, i.e. the model is capable of learning by modulating the association strengths between one or more pattern recognizers in a specific layer and a mental state, as a result of simultaneous activity. This method does not only allow learning by self observation, but also allows learning by imitation since observations do not have to be caused by self generated actions in order to be associated with a mental state of the actor. Moreover, as soon as the relation between responsibility values of a specific layer and corresponding mental states has been learnt, similar observations which are not the result of self generated actions might yield similar responsibility values as self generated actions and thus activate the corresponding mental state of the actor.

In order to complete the description of this conceptual model, a description is needed on how to determine which patterns the pattern recognizers are tuned to. As mentioned above, pattern recognizers can be regarded as experts, and these experts represent a model of a part of the model's potential input. Following this line of thought, one can intuitively conclude that experts which seldom correctly predict the next input are of no avail. In fact, a similar approach is proposed in this study, by retraining pattern recognizers that rarely or never correctly predict their input signals.

The ideas introduced above define the basis of the model presented in this research, i.e. can one infer a mental state by relating an abstract responsibility quantity to a mental state, and does the addition of abstraction layers provide a means for inferring the mental state associated with actions which are hierarchical by nature?

2.2. Formal Description

In this Section, I will elaborate on how the model in Section 2.1 can be implemented, using a mathematical formalization supported by a block scheme of the model which is shown in Figure 2.2.

The model shown consists of L layers, and each layer l consists of K_l pattern recognizers. The set of pattern recognizers at layer l is denoted by $\mathbf{PR}^{(l)}$, where

$$\mathbf{PR}^{(l)} = \{\mathbf{PR}_1^{(l)}, \mathbf{PR}_2^{(l)}, \dots, \mathbf{PR}_{K_l}^{(l)}\} \quad (2.1)$$

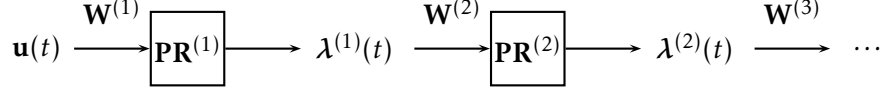


Figure 2.2: A block scheme notation of the model as shown in Figure 2.1, see text for details.

Each pattern recognizer $\text{PR}_i^{(l)}$ produces a responsibility signal $\lambda_i^{(l)}(t)$ as a function of its input at the current time step t and the complete input history, thus

$$\lambda_i^{(l)}(t) = \text{PR}_i^{(l)}(\mathbf{W}_i^{(l)} \lambda^{(l-1)}(t), \mathbf{W}_i^{(l)} \lambda^{(l-1)}(t-1), \dots, \mathbf{W}_i^{(l)} \lambda^{(l-1)}(1)) \quad (2.2)$$

with

$$\lambda_i^{(0)}(t) = u_i(t) \quad (2.3)$$

where $u_i(t)$ denoted the model's i th input at time step t . Furthermore, $\mathbf{W}_i^{(l)}, 1 \leq l \leq L$, denotes a $K_l \times K_{l-1}$ transformation matrix which is used to specify the spatial relationship between the input signal u or λ and the input of a pattern recognizer $\text{PR}_i^{(l)}$.

The output of $\text{PR}^{(l)}$ can be written as a vector $\lambda^{(l)}(t)$, where

$$\lambda^{(l)}(t) = \{\lambda_1^{(l)}(t), \lambda_2^{(l)}(t), \dots, \lambda_{K_l}^{(l)}(t)\} \quad (2.4)$$

For convenience, I will adopt the following notation as a short hand for (2.2) and (2.4):

$$\lambda^{(l)}(t) = \tilde{\mathbf{P}}\mathbf{R}^{(l)}(\mathbf{W}^{(l)} \lambda^{(l-1)}(t)) \quad (2.5)$$

It has to be emphasized that a PR model a dynamical system of arbitrary order, and as a consequence repetitive computation of $\tilde{\mathbf{P}}\mathbf{R}(\cdot)$ does not necessarily result in the same outcome.

The modules denoted by PR form the basic building blocks of the model. The inner working of these blocks is illustrated in Figure 2.3. The weight multiplied input is fed both into a so called ESN module, which I will discuss in detail shortly. For now it suffices to say that an ESN module is responsible for the prediction of a future input signal $\hat{\mathbf{x}}(t+k)$, given the current input signal $\mathbf{x}(t)$. The output of an ESN module is turn fed into a k -unit time delay module denoted by Z^{-k} , from which the output is compared to the current input of the ESN module, which results in a prediction or error signal $\varepsilon(t)$. Thus,

$$\varepsilon(t) = Z^{-k}[\text{ESN}(x(t))] - x(t) = \hat{x}(t) - x(t) \quad (2.6)$$

Next, $\varepsilon(t)$ is fed into a scoring function $\sigma(\cdot)$ which transforms the error signal to a responsibility signal.

2.3. Proposal for Learning Method

In addition to the computation of a responsibility signal, the error signal $\varepsilon(t)$ can also utilized to assess the predictive quality of a pattern recognizer module. The rationale behind

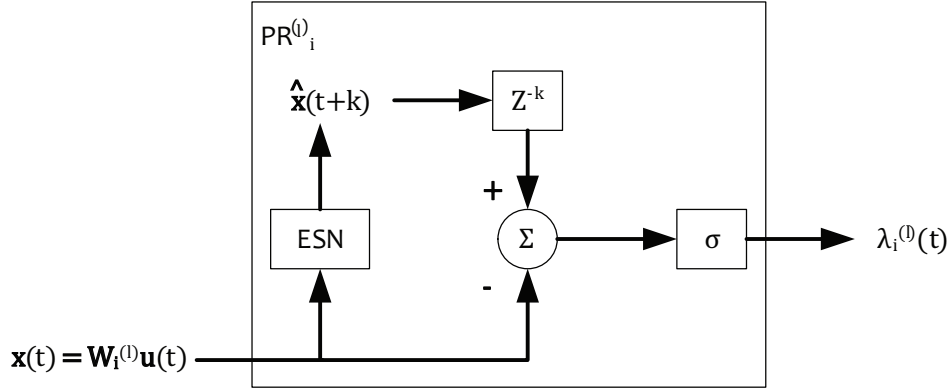


Figure 2.3: Pattern recognizers form the basic building block of the model. The ESN component is responsible for predicting future input signals, denoted by $\hat{x}(t+k)$. These predictions are compared using the Z^{-k} delay operator with the actual input signal $x(t)$, resulting in a prediction error ε which can be transformed to a proper responsibility value using the appropriate transformation function.

this approach is that PR's which seldom correctly predict the next input are of no avail, and should be retrained so that these modules can be reused as predictors of more frequently recurring patterns. This idea can be implemented by a simple bookkeeping mechanism that keeps track of the ratio R between erroneous predictions and correct predictions, where a prediction is correct if $\varepsilon(t) < T_{false}$. As soon as $R > T_{retrain}$, the module can be flagged for retraining and the bookkeeping mechanism can restart counting.

PR modules which are flagged for retraining should be trained to predict a signal representative for a part of its incoming signals. An example implementation for this mechanism is that the module, as soon as it has been flagged for retraining, uses the signals from the next N time steps as a training set for learning. Simulations have to prove whether such mechanism is useful, but are left out the scope of this study due to time constraints.

2.4. Selection of a Prediction Mechanism

As discussed above, the ESN module is responsible for predicting a future input signal $\hat{x}(t+k)$ given the current input signal $x(t)$. Such modules are typically implemented by drawing samples from a generative model, like for example a Hidden Markov Model (HMM). Another option is to use recurrent neural networks (RNNs), but one of the major disadvantages of RNNs is that these generally suffer from slow convergence during training, and convergence to a global optimum in the error surface is not guaranteed.

Recently, so called Echo State Networks (ESN) were introduced by Jaeger (2001), which do not suffer from the drawbacks mentioned in the context of recurrent neural networks. A detailed description of ESNs is outside the scope of this thesis (cf. Jaeger, 2002, for an in depth tutorial), but an important benefit in favor of ESNs is that these offer a biologically plausible method for implementing a prediction mechanism, and are fairly easy to implement in a numerical computing environment such as MATLAB. However, it has to be emphasized that the choice for the selection of such a predictive module is not of major importance within the scope of this study. Indeed, the implementation of a pattern recognizer

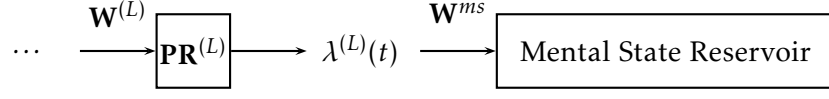


Figure 2.4: Continuation of the block scheme as shown in Figure 2.2. This figure shows the relation between responsibility signals computed at layer L and mental states. During learning the model has to determine the values of connection matrix \mathbf{W}^{ms} which specifies which responsibilities represent a given mental state.

module might also be realized by using other techniques suitable to model their behavior as a (non-linear) adaptive filter.

Finally, responsibilities at level L have to be associated with mental states, which is illustrated in Figure 2.4. The set of possible mental states is referred to as the mental state reservoir, and during learning the model has to determine the values of connection matrix \mathbf{W}^{ms} which specifies which responsibilities represent a given mental state. A solution to this problem would be to increase the weight values in \mathbf{W}^{ms} as a result of coactivation. It is straightforward to implement such a Hebbian like learning mechanism, but again simulations have to point out whether such a mechanism is useful in practice.

Simulation Experiments

In this chapter we will demonstrate how the conceptual model proposed in Section 2 can be implemented and applied to an intention inference task. This demonstration consists of a number of experiments of which the results can provide insight on the capabilities and characteristics of the proposed model. Before we will continue with the details of the experiments, we will first provide the reader with information regarding implementation of the model.

3.1. Implementation Specific Details

In order to perform the above-mentioned experiments the model as described in Section 2 has been implemented using MATLAB¹. For implementation of the echo state network (ESN) component the ESN Tools library² has been used.

Although a formal description of the model is given in Section 2, it is still left unspecified how to implement the scoring function responsible for transforming a prediction error to a quantity which behaves as a responsibility value. Intuitively, an ideal implementation of such a filter would not assign high responsibility values to accidentally correct predictions, but only assign high responsibility values to prediction errors that remain close to zero for a longer time.

To achieve the afore-mentioned behavior, in our experiments a bounded leaky integrator filter is used which is defined in terms of the differential equation

$$\frac{dy}{dt} = \min(-\beta y + \|\mathbf{c}\|_1, \gamma - \|\mathbf{c}\|_1), \quad (3.1)$$

where \mathbf{c} is a vector representing the prediction error, γ is the integrator's upper bound, and β is the leakage rate. Furthermore, the sum of absolute errors $\|\mathbf{c}\|_1$, also known as Manhattan or L_1 distance, is defined as

$$\|\mathbf{c}\|_1 = \sum_{i=1}^n |c_i|. \quad (3.2)$$

A plot of (3.1) for $\beta = 0.1$ and $\gamma = 1$ is given in Figure 3.1. As can be seen in this figure the mechanism described above prohibits the model from assigning low prediction errors to accidentally correct predictions since by definition the leaky integrator filter approaches zero output only after a sufficiently long series of correct predictions. On the other hand,

¹<http://www.mathworks.com/products/matlab/>

²<http://reservoir-computing.org/node/129>

due to its slow dynamics small prediction errors occurring after a series of correct predictions do not immediately result in a rapid decrease of the filter's output.

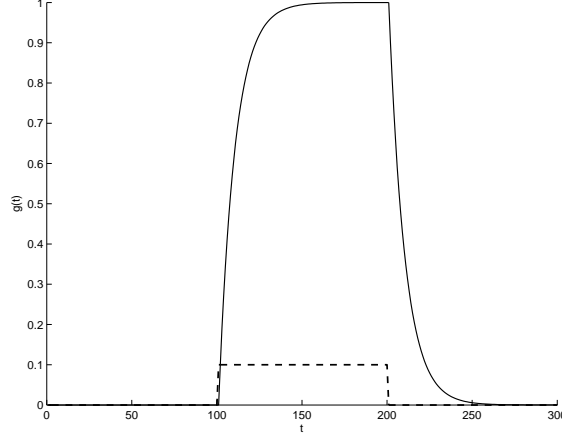


Figure 3.1: This figure shows the output of the bounded leaky integrator defined in (3.1). The dashed line represents the input x at time step t , where the solid line denotes the filter's output for the input signal with $\beta = 0.1$ and $\gamma = 1$. As can be seen in this figure the mechanism described above prohibits the model from assigning low prediction errors to accidentally correct predictions since the leaky integrator filter by definition only approaches zero output after a long enough series of correct predictions. On the other hand, due to its slow dynamics small prediction errors occurring after a series of correct predictions do not immediately result in a sudden increase of its output.

In order to transform the outcome of the leaky integrator filter such that its output lies on the interval $[0, 1]$ the resulting quantity is transformed using the error function $\psi(x)$ which is defined as

$$\psi(x; \sigma) = e^{-x^2/2\sigma}, \quad 0 \leq \psi(x; \sigma) \leq 1. \quad (3.3)$$

A plot of (3.3) is shown in Figure 3.2. As can be seen in this figure, this function only attributes high responsibility values to sufficiently small prediction errors. The spread parameter σ defines the tolerance of the error function, which for this plot is set to $\sigma = 0.005$.

3.2. Experiment 1: Construction of a Pattern Detection Filter

For the first experiment the model is manually trained to detect the occurrence of different patterns embedded in a distractor time series. The distractor time series is constructed by using a Gaussian random walk model to generate a time series S , in which each value at a given time step t is dependent on the value at time step $t - 1$ according to the equation

$$S(t + 1) = S(t) + \varphi^{-1}(z, \mu, \sigma), \quad (3.4)$$

where $0 \leq z \leq 1$ is a uniform distributed random number, and φ is a Gaussian distribution with mean $\mu = 0$ and standard deviation $\sigma = 2$.

In this experiment a single layered model is used, in which the first layer consists of three pattern recognition modules PR_1, PR_2 , and PR_3 . These modules have been trained with time shift parameter $k = 20$ on the following three signals respectively:

$$p_1(t) = \frac{1}{2}(\sin(2\pi \frac{1}{3}t) + \frac{1}{2}), \quad (3.5)$$

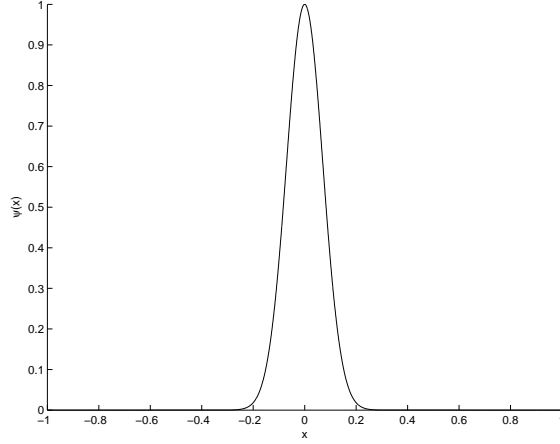


Figure 3.2: A plot of the error function $\psi(x; \sigma)$, see (3.3) for definition. As can be seen in this figure, this function only attributes near one responsibility values to a small prediction errors. The spread parameter σ defines the tolerance of the error function, which for this plot is set to $\sigma = 0.005$.

$$p_2(t) = P(3t - 6), \text{ and} \quad (3.6)$$

$$p_3(t) = 1 - P(3t - 6), \quad (3.7)$$

where t is time in seconds and P is the logistic function defined by

$$P(t) = \frac{1}{1 + e^{-t}}. \quad (3.8)$$

Next, signals p_1, p_2 , and p_3 are embedded at predefined time offsets in the distractor time series S , resulting in a time series u which is used to assess the model's performance on detection of p_1, p_2 , and p_3 . The results of this experiment are shown in Figure 3.3. As can be seen in this figure the model is able to accurately classify and detect the presence of (3.5), (3.6), and (3.7) even in the presence of distractor data.

As can be observed by the responsibility signals the model requires some time in order to recognize the patterns correctly. This is due to the choice for predictive components, as the used echo state network components need some time to stabilize before they are able to produce the correct output. Furthermore, during some timesteps patterns detection errors are small due to the 'accidental' similarity of the output of the modules to the time shifted input signal. However, since a bounded leaky integrator filter has been applied (see Section 3.1), these short drops in prediction error do not lead to noticeably peaks in the responsibility signal.

3.3. Experiment 2: Detecting Hierarchical Organized Sequences

In this second experiment our goal is to demonstrate whether the model is able to detect patterns in responsibility signals as well. Here we use the responsibility signals from our first experiment which, as discussed in Section 3.2, represent the presence of pattern P_1 followed by P_3 , which will be denoted by $\{P_1, P_3\}$. Since the responsibility signals computed in our first experiment are generated by a series of three pattern recognizer modules, the

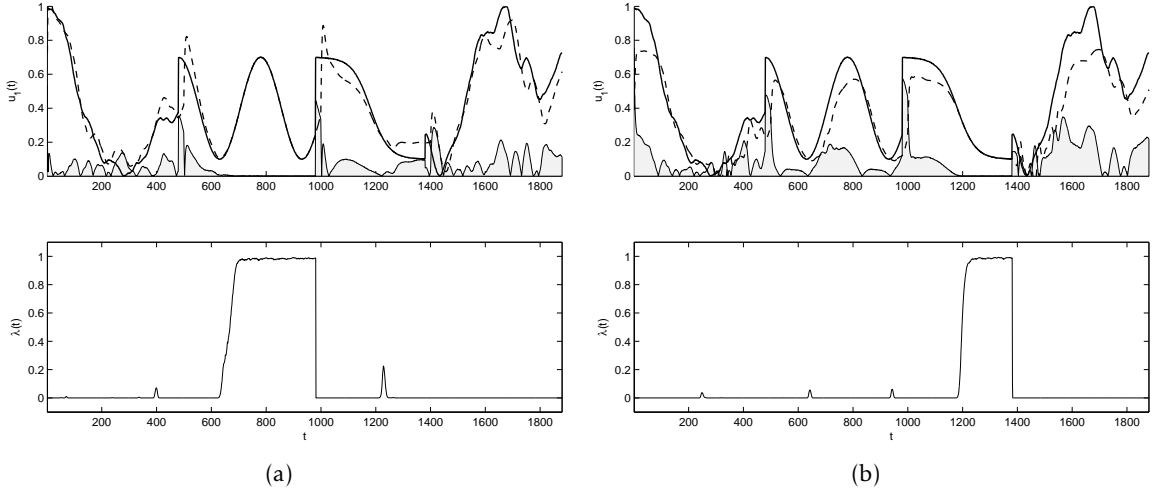


Figure 3.3: (a) provides an overview of the predictive process taking place in a pattern recognizer module trained on the signal defined by (3.5). The solid line in the top row shows the input signal of the model, while the dashed line shows the output of the predictive component. The gray area denotes the absolute prediction error. The bottom row displays the responsibility values associated with the prediction errors in the top row. (b) is similar to (a), but now for another pattern recognizer module trained on the signal defined by (3.7). As can be seen in both figures, both pattern recognizer respond properly since activity of responsibility signals is limited to the time steps at which the corresponding patterns are present in time.

output of the model is a three dimensional time series $\lambda^{(1)}$. In order to recognize $\{P_1, P_3\}$, a layer consisting of two pattern recognizers is added to the model, where the first pattern recognizer is trained for recognizing the three dimensional time series computed in our first experiment, i.e. responsibility signals belonging to detection of $\{P_1, P_3\}$ in the first layer, and the second pattern recognizer is trained on the responsibility signals belonging to the detection of $\{P_2, P_3\}$.

As can be seen in Figure 3.3, the above-mentioned responsibility signals show block like characteristics, and experience shows that the echo state networks used for prediction in the pattern recognizers are difficult to train on prediction of signals with such characteristics. Because of this, these responsibility signals $\lambda(t)$ are transformed using a sinc convolution filter acting as a low pass filter resulting in a more smooth time series $\lambda'(t)$:

$$\lambda'(t) = \lambda(t) * \text{sinc}(x), -2 \leq x \leq 2. \quad (3.9)$$

where $*$ denotes the mathematical convolution operation and sinc is defined as follows:

$$\text{sinc}(x) = \frac{\sin(\pi x)}{\pi x}. \quad (3.10)$$

The results of this experiment are shown in Figure 3.4. As can be seen in Figure 3.4(a) the model generates acceptable predictions of the input signal in the interval $500 < t < 1500$, and as a result positive responsibility values are computed for this interval. To illustrate the behavior of a pattern recognizer during presentation of a pattern that is unknown to a specific pattern recognizer, the input signal also includes the responsibility values associated

with detection of $\{P_2, P_3\}$ as computed by the subordinate layer. What can be seen is that the predictive component is highly sensitive to the presentation of signals it cannot predict, and as a result the responsibility values during the corresponding time intervals remain zero. However, as soon as the second part of the pattern (i.e. P_3) is presented, the model quickly recognizes the second half of the pattern since it is part of $\{P_1, P_3\}$. This phenomenon can also be observed in 3.4(b), which shows the computation process within the second pattern recognizer trained on $\{P_2, P_3\}$. Similar to the first pattern recognizer, P_3 is part of the pattern the second pattern recognizer has been trained on, resulting in correct predictions of the input signal during the intervals on which P_3 is present, thus yielding positive responsibility values.

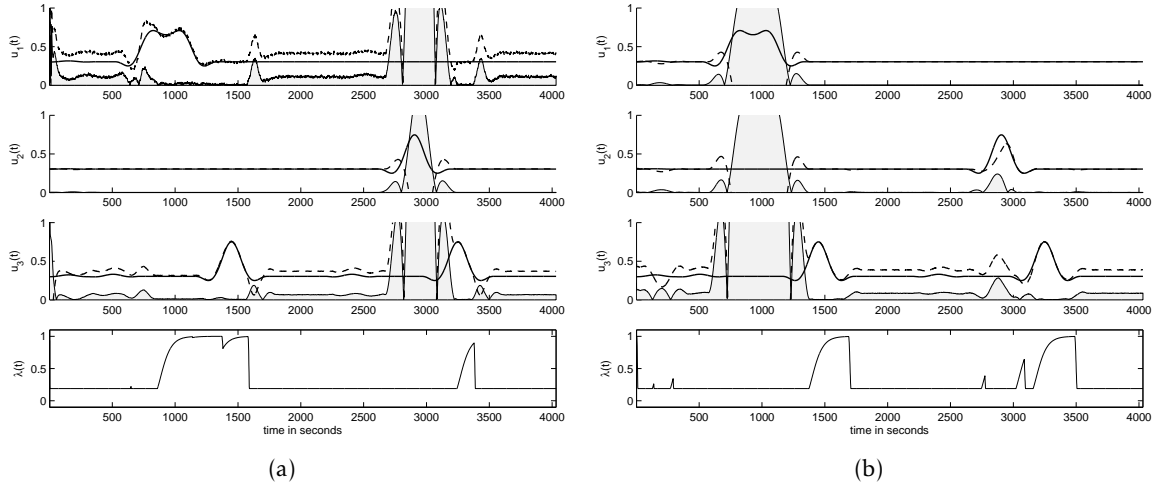


Figure 3.4: As can be seen in Figure 3.4(a) the model generates acceptable predictions of the input signal in the interval $500 < t < 1500$, and as a result positive responsibility values are computed for this interval. To illustrate the behavior of a pattern recognizer during presentation of a pattern that is unknown to a specific pattern recognizer, the input signal also includes the responsibility values associated with detection of $\{P_2, P_3\}$ as computed by the subordinate layer. What can be seen is that the predictive component is highly sensitive to the presentation of signals it cannot predict, and as a result the responsibility values during the corresponding time intervals remain zero. However, as soon as the second part of the pattern (i.e. P_3) is presented, the model quickly recognizes the second half of the pattern since it is part of $\{P_1, P_3\}$. This phenomenon can also be observed in 3.4(b), which shows the computation process within the second pattern recognizer trained on $\{P_2, P_3\}$. Similar to the first pattern recognizer, P_3 is part of the pattern the second pattern recognizer has been trained on, resulting in correct predictions of the input signal during the intervals on which P_3 is present, thus yielding positive responsibility values.

3.4. Experiment 3: Goal recognition in a joint-dial experiment

The experiments discussed above all were conducted using synthetically generated data from mathematical models. To address the question whether the model is able to infer intentions from behavior it is useful to apply the model on real world data as well.

In a recent study Herbert and Butz (2009) demonstrated that human participants adjust the grasp of a dial to the direction and extent of an upcoming rotation of the dial. This an-

ticipatory movement executed before touching the dial might be used by observers to infer the direction and extent of the upcoming dial rotation. Differences in direction and extent of the upcoming rotation yield different trajectories and can therefore be used to infer the goal of the actor. To validate this latter claim, a joint-dial turn experiment has been set up by Herbort (2010) in which one participant, referred to as the actor, was instructed to rotate a dial from its neutral position to either 45° , 90° , -45° or 90° , marked by a LED exclusively visible to the actor. Another participant, referred to as the observer, was instructed to rotate the dial to its neutral position afterwards. During rotation of the dial pronation and supination (p/s) angles of the actor and observer's forearm have been recorded, as well as the distance of the actor's hands to the dial and the rotation angle of the dial.

Details on which variables were manipulated during the course of the experiment are outside the scope of this thesis. However, the p/s angles measured during execution of the rotation movement by the actor are of great interest within the scope of the present study since this information might offer sufficient cues to infer the instructions provided to the actor by the experimenter, even before the actor actually rotates the dial. This information thus can be applied to the model in order to infer the goal of the actor given the anticipatory p/s angles before the actor actually reaches the dial.

In this experiment data from the joint-dial experiment is used for simulation of the model, such that we can assess the model's ability to infer the upcoming rotation of the dial given the p/s angles of an actor, even before the actor touches the dial.

For this simulation several prototypical trials of one subject were selected for each dial target type. For each of these trials all 50 samples before the actor turns the dial were extracted, averaged, and normalized resulting in a prototypical p/s time series for each type of dial target. Next, a model was constructed having one layer, consisting of a pattern recognizer for each of the four target types. These pattern recognizers were in turn trained on the prototypical time series obtained from the selection procedure above.

Visual inspection of the data already revealed that the p/s angle of the underarm is likely to be insufficient in order to differentiate between different extents of rotation, i.e. differentiating between 45° and 90° and between their counter rotations. In order to provide an additional cue for differentiating the extent to which the actor turns the dial, context information was added to the input signal by adding an inverted and exact copy of the underarm's angle respectively in the 45° and 90° signal, and their counter rotations. An appurtenant benefit of augmenting the source signal is that this results in a multi dimensional input signal, which is of great benefit for illustrative reasons.

An example of a training signal is shown in Figure 3.5. This figure shows the training signal corresponding to the pronation/supination (p/s) angles of the forearm executed prior to the 45° turn of a dial in the joint-dial turn experiment. The first dimension, denoted by $u_1(t)$, represents a normalized average of the p/s angle as recorded during the experiment. The other dimension $u_2(t)$ consists of an inverted copy of the first dimension which can be used as a cue for differentiating the time series with the time series corresponding to a 90° turn of a dial.

Results of a simulation run using the joint dial experiment signals are shown in Figure 3.6. In this a detailed overview is given of the recognition process within a single pattern recognizer module executed in the context of the joint dial turning experiment. The two top rows in Figure 3.6(a) represent respectively the normalized pro/supination angle $u_1(t)$ of the actor's underarm during the experiment and contextual information $u_2(t)$ which was added to make it possible to discern between 45° and 90° , and their counter rotations. In

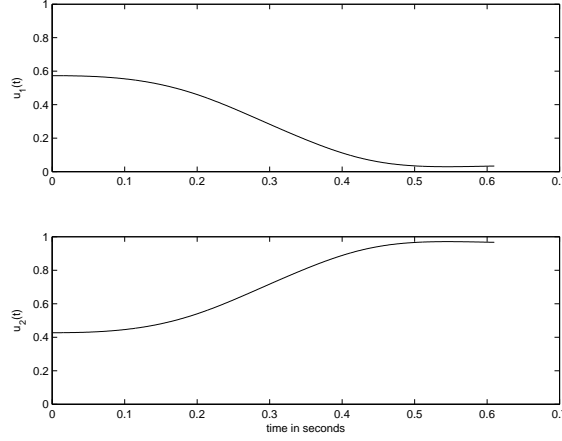


Figure 3.5: Example of a training signal corresponding to the pronation/supination (p/s) angles of the forearm executed prior to the 45° turn of a dial in the joint-dial turn experiment. The first dimension, denoted by $u_1(t)$, represents a normalized average of the p/s angle as recorded during the experiment. The other dimension $u_2(t)$ consists of an inverted copy of the first dimension which can be used as a cue for differentiating the time series with the time series corresponding to a 90° turn of a dial.

this case, since in this trial the target was 45° , $u_2(t)$ consists of an inverted copy of $u_1(t)$. The pattern recognizer shown in this figure has been trained on a prototypical underarm movement corresponding to the anticipatory movement executed before a 45° turn of the dial in the experiment, augmented with an inverted copy of this signal in the form of contextual information (see Figure 3.5). The actual moment at which the dial is turned during the trial is indicated by the vertical line marked with TK. The dashed lines represent the output $\hat{u}_1(t)$ and $\hat{u}_2(t)$ of the pattern recognizer's predictive component. The grey areas correspond to the absolute prediction errors which in turn are used to produce the responsibility signal $\lambda(t)$ shown in the bottom row.

A similar run of the model is shown in Figure 3.6(b). However, this pattern recognizer module is trained on recognizing a prototypical turn of -45° . Here, prediction errors indicate that the module is unable to correctly predict its subsequent input and therefore the corresponding responsibility signal $\lambda(t)$ remains close to zero as expected.

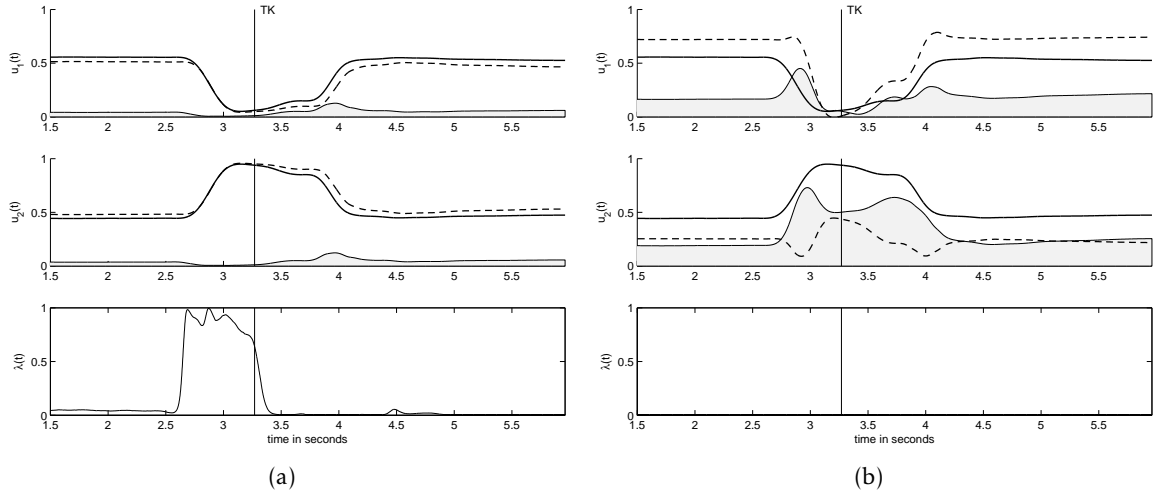


Figure 3.6: (a) This figure provides a detailed overview of the recognition process within a single pattern recognizer module executed in the context of the joint dial turning experiment. The black lines in the two top rows represent respectively the normalized p/s angle $u_1(t)$ of the actor's underarm during the experiment and contextual information $u_2(t)$ (*see text for details*). The pattern recognizer shown in this figure has been trained on a prototypical underarm movement corresponding to the anticipatory movement executed before a 45° turn of the dial. The actual moment at which the dial is turned is indicated by the vertical line marked as TK. Dashed lines represent the output $\hat{u}_1(t)$ and $\hat{u}_2(t)$ of the pattern recognizer's predictive component. Grey areas denote absolute prediction errors which are transformed to the responsibility signal $\lambda(t)$ shown in the bottom row. (b) A similar run but now with a pattern recognizer module trained on recognizing a prototypical turn of -45° . Here, prediction errors indicate that the module is unable to correctly predict its subsequent input and therefore the corresponding responsibility signal $\lambda(t)$ remains approximately zero, as expected.

Discussion & Conclusion

4.1. Discussion

In Section 3 the results of three experiments are described. The results of the first experiment show that the presented model is, like the MSI model, capable of detecting patterns in time series. The second experiment, however, shows that the outcome of this pattern recognition process can in turn be used to detect higher level or more abstract patterns in these same time series. Contrary to the MSI model the presented model analyzes the outcome of this pattern recognition process allowing the model to recognize mental states associated with hierarchical organized behavior. The third experiment shows that, in addition to the first two experiments, the model is not limited to detecting patterns in synthetic data but can also be applied to empirical data.

The simulations conducted in this study reveal some interesting aspects of the presented model, and also offer a number of insights into the use of predictive coding models in general. Some of these aspects are problematic from an engineering perspective since here we are interested in solving the intention understanding problem as best as possible. However, from a scientific perspective these findings might prove to be interesting since these could correspond to similar findings in e.g. humans, increasing the model’s biological plausibility and thus can provide additional insights in the workings of human mindreading capabilities. In the next sections both of these aspects will be discussed in more detail.

4.1.1. Implementation Specific Aspects

In its rudimentary form the presented model does not impose any restrictions on the implementation of the pattern recognition modules it consists of. In order to perform simulations with the model some choices have to be made regarding the implementation of these modules. Some of these choices have a significant effect on the findings that are presented in this thesis, and the effects of these choices will be discussed in this section.

As mentioned before the basic building block for the presented model are pattern recognizer modules. These modules have been implemented using echo state networks, each of which fulfills the role as a predictor of a part of its input space. Echo state networks are a class of dynamical systems, and one property of a dynamical system is that these systems generally have so called attractor dynamics which describes the typical behavior of a dynamical system given a specific driving force and dissipation of the system. The behavior which is a result of this attractor dynamics can be perturbed by changing the driving force, and this might result in chaotic behavior.

In this study the ESN component is tuned or trained in order to output a time shifted copy of the learnt pattern given the pattern as a driving force or input of the ESN component. If the input of the ESN component does not correspond to the pattern it has been trained to, the system quickly departs from its typical behavior and will produce chaotic behavior. Obviously, during this chaotic behavior the predictive quality of the ESN is poor and the corresponding responsibility values will remain close to zero. However, if the pattern which these pattern recognizers have been trained to is present in the incoming signal the ESN will quickly revert to its typical behavior and as a result responsibility values will increase as prediction errors will diminish since the attractor dynamics of the system coincide with the time shifted variant of the input signal.

Utilizing the attractor dynamics of a dynamical system for pattern recognition allows us to quickly respond to the presence of its signal of interest, and thus allows fast recognition of time series in its input space. However, due to the selection of this particular predictive component the model is very sensitive to noise in its incoming signals. Furthermore, experiments show that ESNs do not perform well on the prediction of block-like signals. A possible reason is that echo state networks assume a reservoir of rich dynamics, and these dynamics are formed by the dynamics of the input signal. If the input signal does not contain these dynamics, the ESN is unable to generate a proper output signal.

During learning of predictors we must take some precautions into account. First of all, we must be certain that the predictors do predict the signal and do not just output the current input signal, which for a signal S with slow dynamics results in reasonable predictions (i.e. if $|S(t) - S(t+k)| < \epsilon$). Furthermore, if we train the predictive component on predicting a periodic input for timestep $t + \varphi$ where φ is equal to the input signal's period, this might also result in reasonable predictions if the predictor simply outputs the input signals.

Some approaches to solve the problems mentioned above have been undertaken in this study. In one of these approaches frequency modulation was applied to the responsibility signals in order to increase the dynamics of the signal without loss of information. However, this approach was not useful since small differences in the onset of subparts of a pattern result in phase differences between responsibility signals, which in turn makes it harder for superordinate layers to view these signals as being similar. Another approach which shows significant improvement in recognition is the transformation of responsibility signals using a leaky integrator function. This transformation results in increased dynamics in the resulting signal, and also has a positive effect on the model's ability in distinguishing signals only differing in the order of their subparts. This latter finding can be explained by the fact that the effects of earlier patterns are 'smeared out', i.e. slowly diminish over time due to the leaky characteristic of the transformation function. However, the method which works best for the experiments described in this study is the use of a *sinc* transformation filter as described in Section 3.3. The choice for a low pass filter affects the general performance of the model and will therefore be discussed in the next section.

4.1.2. Fundamental Aspects

Responsibility signals resulting from prediction errors do not contain any information on the incoming signal itself when considered in isolation. These signals by nature contain less information than the input signals, and therefore have different characteristics than the original time series. Most important, these signals tend to contain only low frequency components, and thus less information is lost when these signals are considered in a more

coarse time scale (i.e. lower sampling frequency). Time series prediction requires a certain amount of short term memory (STM), since its output is not solely dependent on its current input but also on its preceding input. STM demands can be reduced by compressing these signals over time, reducing their lengths and therefore also the demand for required STM capacity. In addition to signal compression the model also reduces its STM demands by employing the hierarchical decomposition principle, since pattern recognizers are trained to recognize small parts of the input space and therefore do not need high STM capacities themselves.

The presented model can be extended with the ability to automatically form abstract representations by implementing the suggested learning method. Although this method has not been investigated in this thesis, it is possible to provide some predictions. By implementing the suggested learning method abstract representations evolve within each pattern recognition layer since only representations which are part of the input space will remain in the long run. However, since the proposed model does not include any feedback connections from the mental state reservoir, it is not guided in developing useful representations, and this might result in the formation of meaningless abstractions being of no avail for mental state inference. Furthermore, pattern recognizers might develop overlapping representations, i.e. in theory each pattern recognizer can become responsible for recognizing the same part of the input space. One way to solve this problem is to extend the model by adding lateral competition between the modules within each layer, offering the opportunity to disable training for modules during the presence of a time series which is already recognized by one of the other modules within a particular layer. Furthermore, the addition of feedback signals from the mental state reservoir to subordinate layers can be used to prevent the model from forming meaningless abstractions by, analogous to the previous solution, disabling training during the presence of time series if no mental state is active during observation of that time series.

In Section 4.1.1 it was already mentioned that the choice for a low pass filter of the responsibility signals affects the performance of the model. To be more specific, it affects the degree to which the model is able to distinguish patterns which are composed of the same subparts, but are ordered in a different way. To illustrate this, let us define two patterns, firstly P_1 which consists of the subpatterns A , B , and C , for which we will use the notation $P_1 = \{A, B, C\}$, and similarly $P_2 = \{C, B, A\}$. Suppose we have two pattern recognizers, R_1 and R_2 , which are trained to recognize P_1 and P_2 respectively. Presuming that both R_1 and R_2 are capable of recognizing these patterns after exposure to just one sample, these recognizers will by definition not be able to distinguish P_1 and P_2 , as both patterns are composed of the same parts, and each recognizer recognizes its intrinsic pattern in both time series. This is caused by the fact that these time series are by definition part of their intrinsic pattern, and as a consequence equal responsibility signals are computed. Introducing reverberations by low pass filtering the responsibility signals has the positive effect of a decrease in confusion between time series consisting of the same but differently ordered subparts.

Another important aspect of the pattern recognizers is their capacity to generalize over different time series, i.e. the amount to which these recognizers are invariant over different properties of the time series these are able to recognize, such as noise sensitivity, and differences in speed. These properties may influence the model's discernment negatively, since overgeneralization can result in the model not being able to distinguish between time series related to different mental states only differing in the properties mentioned above.

4.1.3. Relation to Mirror Neurons

Although the model proposed in this thesis is definitely designed from an engineering perspective, it is certainly inspired by biological models for intention understanding. It is therefore worthwhile to discuss its relevance to other fields, in particular from the perspective of cognitive science.

The discovery of ‘mirror neurons’ led to a significant increase of interest in action understanding. These neurons do not only become active during performance of an action, but are also active during observation of that same movement. Cells exhibiting these properties were first discovered in macaque monkeys, and regions possessing similar properties were discovered later in humans. Although the function of the mirror neuron system (MNS) is subject of fierce debate, one popular explanation for the function of the MNS is that it is responsible for mentalizing, or mindreading. This idea is rooted in simulation theory, which states that mindreading is achieved by replicating others mental states, like putting one in other’s shoes. Because of its mirror properties, the MNS is regarded as evidence in favor of simulation theory, and can as such be considered as the neural correlate of these mentalizing capabilities.

Fogassi et al. (2005) reported that mirror neurons have been located which respond selectively to motor acts as belonging to an action sequence, and thus predicting the intended goal of a complex action. In their study, monkeys were tested in two conditions, namely food placement in mouth and placement in container near the mouth. Neurons were found that responded selectively to the execution as well as the observation of movements followed by bringing food to the mouth. Activity in the same cells was absent when the executed movement was the same but was followed by putting food in the container next to the mouth. About two third of neurons discharged preferentially when grasping was embedded in a specific motor action, and since it is assumed that the kinematics of the movement is the same in both conditions it is concluded that the activity of these neurons allows the prediction of intentions.

A common conception is that the MNS is associated with both action and perception. This is reflected in the abundance of computational models for the MNS system and their relation to motor control. This relation is evident when considering the RNNPB model (Tani et al., 2004), which was initially devised as a model for imitation learning but exhibits mirror properties as well. However, imitation has not been proven to exist in primates, and it is therefore debatable whether the MNS’s primary goal is to serve this purpose. In this thesis we would like to point out that the role for the MNS as a mentalizing module does not imply the existence of a motor system. The model presented in this thesis shows that action understanding is possible even in the absence of motor control.

Later evidence partly contradicts with some MNS ideas, namely that the MNS might be involved in action understanding but not inferring in higher level intentions. Another explanation is that it is only involved in familiar scenarios where it acts as a matching system, but does not differentiate actions belonging to higher level intentions unfamiliar to the observant (Brass et al., 2007). In Kilner and Frith (2008) it is argued that this contradiction is in fact an instance of hierarchical organization and that it is best considered within a predictive coding framework, like the model that is presented in this thesis.

4.1.4. Biological Evidence for Existence of Prediction Errors

In the presented model information is transmitted by signals derived from prediction errors. Neurophysiological evidence shows the existence of analogues of these signals in the animal brain, in particular in monkey dopamine cells (Schultz and Dickinson, 2000). Several functions for these neuronal prediction error signals have been suggested. In one of these suggestions it is stated that these prediction errors might play a role in guidance of learning, which is related to the Rescorla-Wagner model (Rescorla and Wagner, 1972). This model is able to predict robust phenomena in associative learning, in particular the blocking effect. Blocking occurs when an conditioned stimulus (CS1) has been paired with an unconditioned stimulus (US). As soon as another conditioned stimulus CS2 in compound with CS1 is presented before the US, in theory CS2 could be used to predict the US. However, in practise this association is not learnt, a phenomenon which is called the blocking effect. The Rescorla-Wagner model predicts this effect since it incorporates a prediction error for learning, i.e. as soon as the CS predicts the US, no learning takes place. If prediction performance is bad, the error is high and learning will take place.

4.2. Conclusion

In this thesis it has been investigated whether the Mental State Inference (MSI) model by Oztot et al. (2005) could be improved by adding hierarchical structure to the model. The MSI model is built upon multiple visuomanual feedback controllers, of which forward models form the key ingredients. The role of these forward models is twofold, namely i) to predict sensory consequences of motor commands eliminating sensory processing delays and ii) to be used in mental simulation, allowing an observer to infer the mental state associated with the forward model which best predicts other's actions.

Oztot and colleagues justify the use of mental simulation in order to infer others' mental state as follows: 'Mental state inference using mental simulation allows interpretation without prior experience, since as long as a movement or behavior is in the repertoire of an organism, it will be interpretable without any training'. However, this approach does not allow interpretation of mental states corresponding to behavior which is not producible by the observing organism itself (e.g. behavior produced by organisms having a different kinematic structure). Their claims are substantiated by experimental evidence in favor of the involvement of the motor system in action observation/mental state inference, in particular the involvement of mirror neurons, since experiments show that these neurons discharge during action observation as well as during execution of the same action. These findings are compatible with an action understanding mechanism based on a mental simulation loop involving mirror neurons, and are compatible with simulation theories and in particular with the direct matching hypothesis as proposed by Rizzolatti et al. (2001). However, recent studies dispute this direct matching hypothesis, and suggest a more restrained role of a mirror neuron system in action understanding (see Figure 4.1). In fact, accumulating evidence shows that 'action understanding may precede, rather than follow from, action mirroring' (Csibra, 2007). In line with Csibra and colleagues' argument, the model presented in this thesis demonstrates that mental state inference is possible without the need of mental simulation, and is therefore compatible with recent findings as part of the functional role of the mirror neuron system and action understanding in general. Furthermore, by eliminating the involvement of the motor system in the mental state inference process, the presented

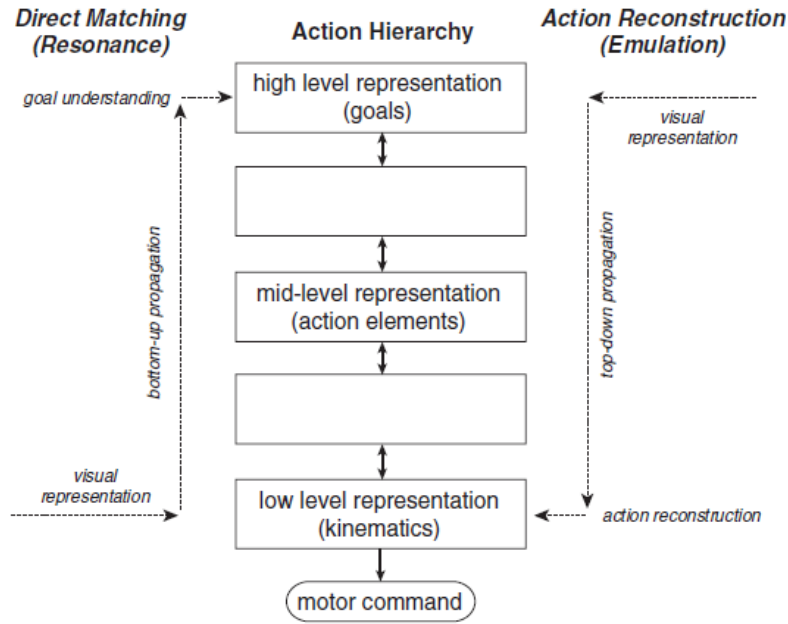


Figure 4.1: Direct matching vs. action reconstruction (Csibra, 2007)

model does not suffer from the self-sustaining issues associated with the idea of using the MSI model for imitation learning as discussed in the introduction.

Oztop and colleagues argue that ‘a fixed amount of resources is required to implement the mental simulation circuit as opposed to a dedicated neural circuit that requires an ever-growing storage requirement with the increasing number of behaviors to interpret’. However, their model requires controllers to be as complex as the behavior associated with the mental state, and therefore the amount controllers required to interpret complex behavior in terms of mental states results in a combinatorial explosion. For example, in the MSI model two different controllers are required for recognizing the difference in e.g. moving towards a cup and performing a precision grip, and moving towards the same cup and performing a power grip, which both might be associated with a different mental state or higher level goal (resp. drinking and moving the cup).

The model presented in this thesis overcomes this combinatorial explosion by using hierarchical decomposition, which is based on the idea that in order to recognize complex behavior we must account for different levels of action recognition. It is assumed that the recognition of more abstract actions depends on the recognition of hierarchical ordered sequential behavior by detecting the occurrence of physical movements. This information can be combined such that higher order descriptions (e.g. behavioral primitives) are formed, which again can be combined in order to form even more abstract descriptions (e.g. goals). Results from simulations show that the model is indeed capable of recognizing complex hierarchical organized patterns in multidimensional time series. Assuming that these time series are the result of a specific mental state either originating in the observer or the demonstrator, the presence of these time series can be associated with a particular mental

state allowing one to infer mental states from observations later on.

Finally, in addition to the MSI model the presented model includes a description of a learning method which can be used to automatically develop abstract representations of patterns in the input space of the model. Although model learning has not been investigated in this thesis, the ideas presented here show that our ability to ‘mindread’ might very well be rooted in a primitive ability of humans to become ‘tuned’ to the world surrounding us, resulting in the development of a hierarchical information processing system allowing us not only to relate observed behavior to its underlying mental states but also to predict and anticipate to this same behavior.

Bibliography

- Appelt, D. and Pollack, M. (1992). Weighted abduction for plan ascription. *User Modeling and User-Adapted Interaction*, 2(1):1–25.
- Baker, C., Tenenbaum, J., and Saxe, R. (2007). Goal inference as inverse planning. In *Proceedings of the 29th annual meeting of the cognitive science society*.
- Brass, M., Schmitt, R., Spengler, S., and Gergely, G. (2007). Investigating action understanding: inferential processes versus action simulation. *Current Biology*.
- Breazeal, C. and Scassellati, B. (2002). Robots that imitate humans. *Trends in Cognitive Sciences*, 6(11):481–487.
- Charniak, E. and Goldman, R. (1993). A Bayesian model of plan recognition. *Artificial Intelligence*, 64(1):53–79.
- Csibra, G. (2007). Action mirroring and action interpretation: An alternative account. *Sensorimotor foundations of higher cognition: Attention and performance XXII*, ed. P. Haggard, Y. Rosetti & M. Kawato. Oxford University Press.[aPC].
- Fogassi, L., Ferrari, P., Gesierich, B., Rozzi, S., Chersi, F., and Rizzolatti, G. (2005). Parietal lobe: from action organization to intention understanding. *Science*, 308(5722):662–667.
- Greenfield, P. et al. (1991). Language, tools and brain: The ontogeny and phylogeny of hierarchically organized sequential behavior. *Behavioral and brain sciences*, 14(4):531–551.
- Greenfield, P., Nelson, K., and Saltzman, E. (1972). The development of rulebound strategies for manipulating seriated cups: A parallel between action and grammar* 1. *Cognitive psychology*, 3(2):291–310.
- Haruno, M., Wolpert, D., and Kawato, M. (2001). Mosaic model for sensorimotor learning and control. *Neural Computation*, 13(10):2201–2220.
- Herbort, O. (2010). Can an Observer Predict the Goal of a Dial-turn by Observing the Grasp? *Unpublished manuscript*.
- Herbort, O. and Butz, M. (2009). Anticipatory planning and control of hand orientation in grasping movements. *Manuscript submitted for publication*.
- Hubel, D. and Wiesel, T. (1959). Receptive fields of single neurones in the cat’s striate cortex. *The Journal of Physiology*, 148(3):574.

- Hubel, D. and Wiesel, T. (1968). Receptive fields and functional architecture of monkey striate cortex. *The Journal of Physiology*, 195(1):215.
- Jaeger, H. (2001). Short term memory in echo state networks. *German National Institute for Computer Science*.
- Jaeger, H. (2002). Tutorial on training recurrent neural networks, covering BPPT, RTRL, EKF and the "echo state network" approach. *Fraunhofer Institute for Autonomous Intelligent Systems (AIS), International University Bremen*.
- Jaeger, H. (2007). Discovering multiscale dynamical features with hierarchical Echo State Networks. *Jacobs University Bremen, Tech. Rep.*
- Johansson, G. (1973). Visual perception of biological motion and a model for its analysis. *Perceiving events and objects*.
- Kilner, J. and Frith, C. (2008). Action observation: inferring intentions without mirror neurons. *Current Biology*, 18(1):32–33.
- McClelland, J. and Rumelhart, D. (2002). An interactive activation model of context effects in letter perception. *Psycholinguistics: critical concepts in psychology*, 88(5):422.
- Oztop, E., Wolpert, D., and Kawato, M. (2005). Mental state inference using visual control parameters. *Cognitive Brain Research*, 22(2):129–151.
- Rescorla, R. and Wagner, A. (1972). A theory of Pavlovian conditioning: Variations in the effectiveness of reinforcement and nonreinforcement. *Classical conditioning II: Current research and theory*, pages 64–99.
- Rizzolatti, G., Fogassi, L., and Gallese, V. (2001). Neurophysiological mechanisms underlying the understanding and imitation of action. *Nature Reviews Neuroscience*, 2(9):661–670.
- Schultz, W. and Dickinson, A. (2000). Neuronal coding of prediction errors. *Annual Review of Neuroscience*, 23(1):473–500.
- Tani, J., Ito, M., and Sugita, Y. (2004). Self-organization of distributedly represented multiple behavior schemata in a mirror system: reviews of robot experiments using RNNPB. *Neural Networks*, 17(8-9):1273–1289.
- Wolpert, D. and Kawato, M. (1998). Multiple paired forward and inverse models for motor control. *Neural Networks*, 11(7-8):1317–1329.

MATLAB Code Listing
