

Encoding Longer-term Contextual Sensorimotor Information in a Predictive Coding Model

Junpei Zhong* and Tetsuya Ogata*[†] and Angelo Cangelosi[‡]

*National Advanced Institute of Industrial Science and Technology (AIST), Tokyo, Japan

[†] Lab for Intelligent Dynamics and Representation, Waseda University, Tokyo, Japan

[‡]School of Computer Science, University of Manchester, Manchester, UK

Abstract—Studies suggest that the difference of the sensorimotor events can be recorded with the fast- and slower-changing neural activities in the hierarchical brain areas, in which they have bi-directional connections. The slow-changing representations attempt to predict the activities on the faster level by relaying categorized sensorimotor events. On the other hand, the incoming sensory information corrects such event-based prediction on the higher level by the novel or surprising signal. From this motivation, we propose a predictive hierarchical artificial neural network model which is implemented the differentiated temporal parameters for neural updates. Also, both the fast- and slow-changing neural activities are modulated by the active motor activities. This model is examined in the driving dataset, recorded in various events, which incorporates the image sequences and the ego-motion of the vehicle. Experiments show that the model encodes the driving scenarios on the higher-level where the neuron recorded the long-term context.

I. INTRODUCTION

The predictive coding (PC) theory [1, 2, 3, 4] asserts a two-stream integrated theory for the sensorimotor loop: the top-down stream works as a predictive machine which utilizes both perception and action to minimize the prediction error, i.e. the difference between bottom-up/externally sensory stimuli. When encountering the sensorimotor events, the whole loop attempts to minimize the difference between the posterior estimation and the truth from its perception, by changing its internal learning model (“perceptual inference” (see also [5] and [6]) or by the action execution (“active inference”, see also [7] and [8]). Additionally, because of the integrative property of both perception and action, perceiving the world (perceptual inference) and acting on it (active inference) can be regarded as two aspects with the same aim: to minimize the prediction error.

As such, the integrative process of predictive model follows a bi-directional learning mechanism on each level of our hierarchical brain. It is suggested that within the hierarchical architecture, the topological higher level in the brain areas infers the prediction on the lower areas with a slower changing activities [9, 10]. This is done by its subsets of such prediction representations that are transmitted to the lower levels to predict the upcoming faster neural activities on the lower level. For instance, areas on the higher-level of our brain learn multiple world models and act as prior to explain the best descriptions of the upcoming percept. This continual top-down process acts as an “explain away” function (e.g. [11, 12]): the explanation on the higher-level offers the best parameters

based on the previous information to predict the most likely events of the sensory data on the lower levels, and explains away the other models. Such hierarchical function can be realised by the interaction of neural oscillations in different time-scales, which encode different temporal parameters of the world models.

Therefore, the higher level representation in a hierarchical model may implicitly represent the contextual information based on the understanding of the previous sensorimotor information. As such, the formed internal world model on the higher level has to be shaped by the statistical structure of the error. Based on such a hypothesis, we suggest that the importance of different time-scales should also an essential role in forming the internal models of the PC framework.

II. RELATED WORKS

The difference of the temporal scales of prediction results in different cognitive functions in embodied internal models. Some of the previous research focused on the short-term predictive function of the internal model. In most cases, such short-term prediction can act a compensation function of the sensorimotor integration (e.g. [13, 14, 15, 16]). Based on the PC framework, the PredNet model [17] is one of the hierarchical models that preserve both short- and mid-term memories for images. A recent work [18] proposed a model called AFA-PredNet which integrates both motor action and perception in the PC framework. In this network, the motor action is used as an attention model for the prediction from a couple of recurrent networks. However, the long-term prediction based on the understanding of the world model is still missing in both the PredNet and the AFA-PredNet models.

Indeed, when we think about the predictive functions in biological brains, there are no explicit boundaries between the short- and mid-term prediction and the long-term predictive: the short-term prediction is based on a long-term understanding and prediction of the world. For instance, [19] and [20] studied how to apply internal model to control the actual motor actions, mostly focusing on the predictive control of a motor action. [21] extended these models to imitation learning of the sensorimotor behaviours. The long-term planning behaviours can also emerge from internal simulation where the prediction occurs constantly (e.g. [22, 23]).

Specifically, we can consider the pre-symbolic representation as a understanding of the long-term context, which

is learnt in a unsupervised way. From this perspective, [24] reported an embodied experiment in which an association between the semantic meaning and the sensorimotor behaviours emerges by a recurrent architecture called Recurrent Neural Network with Parametric Bias Units (RNNPB). Based on the extension of this network, [25] discovered that the semantic representation about the object movements and object features also emerge in a recurrent neural network. Specifically, the network is able to predict the next probable position of the object movement, while to pre-symbolic representation is given.

When we regard the unification of different time-scales in a single predictive model with artificial recurrent connections, experiments based on the Multiple Timescale Neural Network (MTRNN) [26] offers an explanation from the view of the non-linear dynamical system for such phenomena. It can be regarded as another extended version of the RNNPB. The neurons on the higher-level of the MTRNN are with slower-changing neural activities, which modulates the neural activities on the lower-levels by the similar roles of the bias inputs. Thus, the whole network is able to work as a number of non-linear dynamic functions as a similar role of bifurcation. While the model is used to learn the temporal sequences such as the sensorimotor information of the robots, the model is able to represent different spatio-temporal scales of sensorimotor information, such as the language learning [27, 25] and object features/movements [28]. Similar concept of multiple time-scales has also been applied in Gated Recurrent Units for automatically context extraction [29, 30].

The multiple time-scales concept can also be extended in different modalities. For instance, the multiple spatio-temporal scales RNN (MSTRNN) [31] integrates the MTRNN and convolutional neural networks [32, 33], where both the spatial and temporal information are connected and associated on the higher level, where slower changing neurons represent the sensorimotor behaviours. The slower changing units on the higher level also makes the dynamics of the model easier to be interpreted, examined and changed. On the other hand, compared with MSTRNN, the PredNet [17] follows the definition of PC while using the difference as inputs on each layer. And it also uses the convolutional network to capture the local features of the visual streams. But the PredNet builds the temporal prediction in the top-down perception part, which

makes the model more biological plausible.

Building the PC embodied model with the concept of multiple time-scales would be beneficial for both engineering and cognitive studies. Firstly, it follows the results from the brain and cognitive studies that different response times while the neurons react to conscious/unconscious prediction. Second, the slower changing neurons in such a model would be easier for us to control and examine the dynamical behaviours of the model or the embodied systems. These are the main motivations why we are proposing the embodied model to examine the temporal scales in the predictive coding framework.

III. THE MODEL

The proposed MT-AFA-PredNet (Multiple Time-scale Action Formulated Predictive Network) is shown in Fig. 1. In general, the MT-AFA-PredNet is functionally organized as an integration with two networks: the left part is equivalent to a generative recurrent network, while the right part is a standard convolutional network.

In terms of architecture, it is similar as AFA-PredNet [18]. The important features of this architecture are:

- 1) There are a number of recurrent neural networks. (e.g. Convolutional LSTM) on each level of the model, which learn different possibilities of the prediction (a generative unit, *GU*, green)
- 2) The input of the motor action is used as an additional signal for the modulation of the prediction (the motor modulation unit, *MM*, grey). Specifically, it acts as an attention mechanism for the prediction from the upper level (top-down prediction);
- 3) The convolutional network in the bottom-up part capture the feature of the error on each level, (the discriminative unit, *DU*, blue);
- 4) The difference of the updating rate on different levels of the architecture determine different representation of the spatio-temporal properties of the sensorimotor behaviours (the error unit, *ER*, red).

The neural functions on each neural unit can be found in Eq. III. Although the main architecture of the MT-AFA-PredNet is the same as AFA-PredNet, the most important feature is that in the neural function of the generative unit (Eq. 4), the generated output is determined not only by the current neural status, but also its previous status. The fraction of the output is determined by the temporal parameter τ .

$$X_l(t) = \begin{cases} i(t), & \text{if } l = 0, \\ \text{MAXPOOL}(f(\text{Conv}(E_{l-1}(t))))), & \text{if } l > 0 \end{cases} \quad (1)$$

$$\hat{X}_l(t) = f(\text{Conv}(R_l(t))) \quad (2)$$

$$E_l(t) = [f(X_l(t) - \hat{X}_l(t)); f(\hat{X}_l(t) - X_l(t))] \quad (3)$$

$$R_l^d(t) = (1 - \frac{1}{\tau})R_l^d(t) + \frac{1}{\tau}\text{ConvLSTM}(E_l(t-1), R_l(t-1), \text{DevConv}(R_{l+1}(t))) \quad (4)$$

$$R_l(t) = \text{attention}(a(t)) \times R_l^d(t) \quad (5)$$

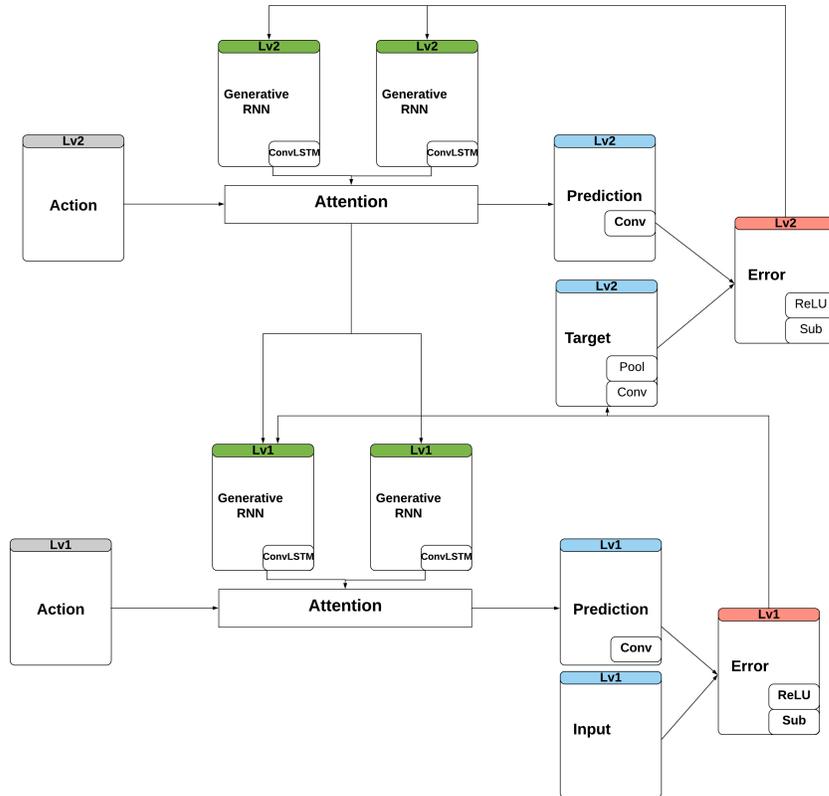


Fig. 1: A 2-layer AFA-PredNet

where $f(\cdot)$ is an activation function of the neurons, which we apply the ReLU function to ensure a faster learning in back-propagation, $X(\cdot)_i^t$ is the neural representation of the level l at time t . The representation on the EL layer l is $E(\cdot)_l$. The *MAXPOOL*, *Conv*, *ConvLSTM* and *MLP* are the corresponding neural algorithms. Specifically, to realize the time scale concept, Eq. 4 indicates that the predicted information in the *GU* unit should consider the previous state of the *ConvLSTM* outputs as well as the current output. This is determined by the time parameter τ .

The overall algorithm for learning a whole sequence is showed in Algorithm 1.

IV. EXPERIMENT

In this section, the performance of the network as well as the analysis of the neural activities will be conducted in driving dataset. We conducted the experiments to examine the performances of the multiple time-scale properties of the proposed AFA-PredNet network. The targets of the experiments are twofold:

- 1) The prediction of the incoming images by synthesising image sequences;
- 2) The representation of scenarios by the multi-time scales properties.

With consideration of this target, we chose the driving data-set¹. This data-set, provided by Daimler AG, contains five labelled driving scenarios, each of which contains 250 or 300 images. Additionally, the driving information is also included corresponding to the very time-stamp of the each image taken. The 5 units long of vector indicates the ego-motion information:

- 1) $\theta \in [-\pi, \pi]$: the angle of the steering wheel.
- 2) $\{v1, v2, v3, v4\} \in [0, 300]$: the velocities of each wheel.

A 3-layer MT-AFA-PredNet was used for training the sequence of both motor action vectors (i.e. the velocities of the wheels) and images, with the Adam optimizer [34]. Three different values of τ were applied in three different layers. With a larger *tau* on the upper levels, it indicates slower neural activities would be expected. Compared with the τ values selected in MTRNN works (e.g. [26, 28]), a much smaller τ values are chosen, because the LSTM networks performs longer term memories by themselves. The parameters are shown in the table:

A. Synthesis of Image

The epoch is set to be 300, each of which includes 500 iterations for each sequence. After the training, the RMS of

¹<https://ccv.wordpress.fos.auckland.ac.nz/eisats/set-1/>

```

Data:  $i(t) \& a(t) \in data$ 
while  $error > threshold$  or
 $iteration > maximum\_iteration$  do
  for  $t \leftarrow 0$  to  $T$  do
    for  $l \leftarrow 0$  to  $L$  do
      if  $l == L$  then
         $R_l^d(t) = (1 - 1/\tau)R_l^d(t-1) + 1/\tau \cdot$ 
         $ConvLSTM(E_l(t-1), R_l(t-1));$ 
      else
         $R_l^d(t) = (1 - 1/\tau)R_l^d(t-1) + 1/\tau \cdot$ 
         $ConvLSTM(E_l(t-1), R_l(t-1), DevConv(R_{l+1}(t)));$ 
      end
       $R_l(t) = MLP(a(t)) \times R_l^d(t);$ 
    end
    /* Generative (top-down) Process */
    for  $l \leftarrow L$  to  $0$  do
       $\hat{X}_l(t) = f(Conv(R_l(t))); E_l(t) =$ 
       $[f(X_l(t) - \hat{X}_l(t)); f(\hat{X}_l(t) - X_l(t));$ 
      /* Discriminative (bottom-up) Process */
    end
  end
end

```

Algorithm 1: MT-AFA-PredNet Computation

Parameters	Value
τ_0	1.0
τ_1	1.3
τ_2	2.0
Kernel	3×3
Padding	1
Pooling	2×2

TABLE I: parameters

each image are calculated as:

$$RMS = \frac{1}{T} \sqrt{\sum_{t \in (0, T]} \sum_{i \in pixels} (i(t) - \hat{i}(t))^2} \quad (6)$$

We compare the RMS error between the MT-AFA-PredNet and the LSTM in Tab. II. As a baseline, a single layer LSTM is used to predict the image sequences. As we can see, the MT-AFA-PredNet, which uses a multiple layer convolutional LSTM, performs better in prediction than the single layer LSTM, probably because the convolution calculation on each layer is beneficial to detect the image features.

The quantitative comparison between the ground-true and the synthesised images are shown below.

Fig. 2 and Fig. 4 show the comparison between some samples of the original and the predicted images in the scenario "crazy turn", and Fig. 3 and Fig. 5 show the scenario of "construction site". Since we did the normalisation after the prediction, the synthesized images show inverted colour.

	LSTM	MT-AFA-PredNet
Construction site	8.332	7.219
Crazy turn	10.315	8.287
Dancing light	9.834	7.314
Intern on bike	5.411	5.131
Safe turn	9.314	8.107
Squirrel	7.908	7.781

TABLE II: RMS between LSTM and MT-AFA-PredNet

B. Scenario Classification

We further visualise the neural activities on different layers to examine how time parameters τ affects the representation. Due to the page limit, in this subsection, only the quantitative results are shown: We first visualise the representation on the layers 1 and 2 in the first two scenarios ("crazy turn" and "construction site"). The quantitative comparison will be conducted to see whether the update on each layer has been differentiated. Then we will observe it has been categorized based on the representation on layer 2.

Corresponding to the prediction samples, the internal representations of 1st and 2nd layers of the GU units are shown respectively in Figs. 6, 7, 8 and 9. We can observe that the higher-level representation (Layer 2) remains more steady than the lower levels. And the representation seems is encoded in a sparse way. From this result, we can basically categorise different driving scenarios as shown in Fig. 10, where we can see there are different representation with different training sequences.

V. CONCLUSION

In this paper, we propose that the top-down prediction in the PC framework occurs based on the longer term understanding representing the contextual multi-modal information. A few neuroscience studies have suggested the temporal difference in neural activities can be found to be differentiated in the hierarchical brain areas. Therefore, the multiple time-scales concepts have been applied in an embodied PC model, which is called the Multiple Time-scale Action Formulated Predictive model (MT-AFA-PredNet). Specifically, in this model, the higher-level encodes slowly changing information of both perception and action, indicating the understanding of the sensorimotor event, resulting in the categorizing function of the model. The experiments has been conducted in the driving dataset with the ego-motion of the vehicle, which show that the update rates of different levels of neurons are differentiated, so that the information of the scenario (i.e. the longer-term context) is encoded on the higher level. Interesting results were observed, which leads us some questions for the next experiment:

- 1) The representation on the higher level updates slower than the lower level, but we have not figured out in which situations it updates mostly (e.g. changing of movement).
- 2) How is the representation of the GU units be different if different possibilities of motor actions are learnt with the same contextual information of perception?

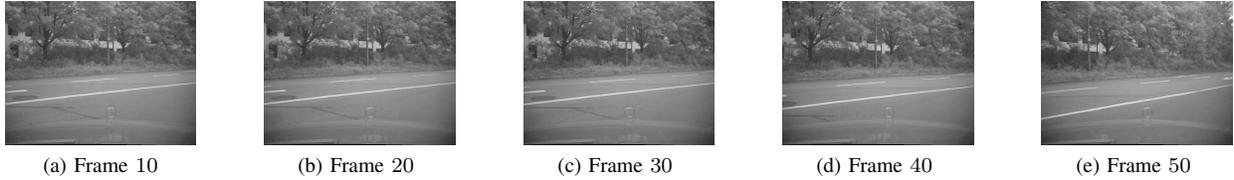


Fig. 2: Image Samples from the Left Camera (Crazy Turn)



Fig. 3: Image Samples from the Left Camera (Construction site)

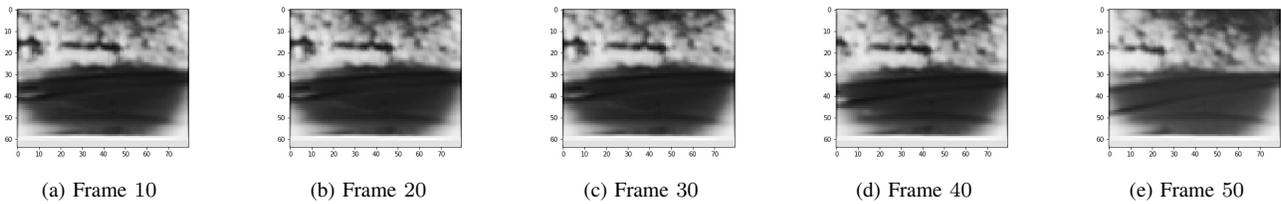


Fig. 4: Predicted Images after Training (Crazy turn)

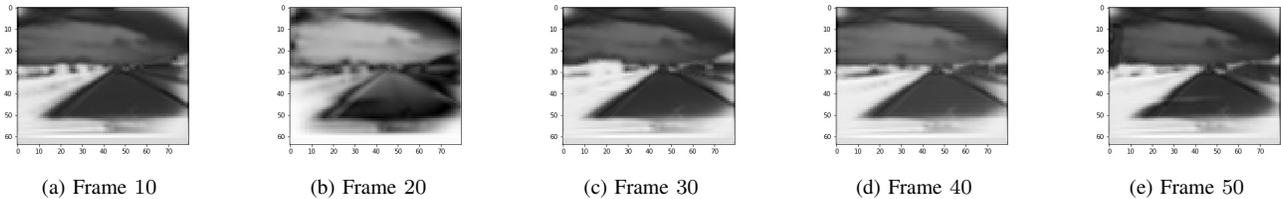


Fig. 5: Predicted Images after Training (Construction site)

3) How to fine-tune the sparseness of the representation?

Besides, at the next stage, we will examine the network performance in details with different hyper-parameters. Also, it would be interesting to explore the process of interaction between the short- and long-term prediction emerge in the neural representation.

ACKNOWLEDGEMENT

The research was partially supported by New Energy and Industrial Technology Development Organization (NEDO). A Pytorch implementation of MT-AFA-PredNet can be found on Github²

²https://github.com/jonizhong/mta_prednet.git

REFERENCES

- [1] A. Clark. “Whatever next? Predictive brains, situated agents, and the future of cognitive science”. In: *Behavioral Brain Sciences* (2012), pp. 1–86.
- [2] R. P. Rao and D. H. Ballard. “Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects”. In: *Nature neuroscience* 2.1 (1999), pp. 79–87.
- [3] K. Friston. “Learning and inference in the brain”. In: *Neural Networks* 16.9 (2003), pp. 1325–1352.
- [4] K. Friston. “A theory of cortical responses”. In: *Philosophical Transactions of the Royal Society B: Biological Sciences* 360.1456 (2005), pp. 815–836.
- [5] E. M. Segal and T. G. Halwes. “The influence of frequency of exposure on the learning of a phrase struc-

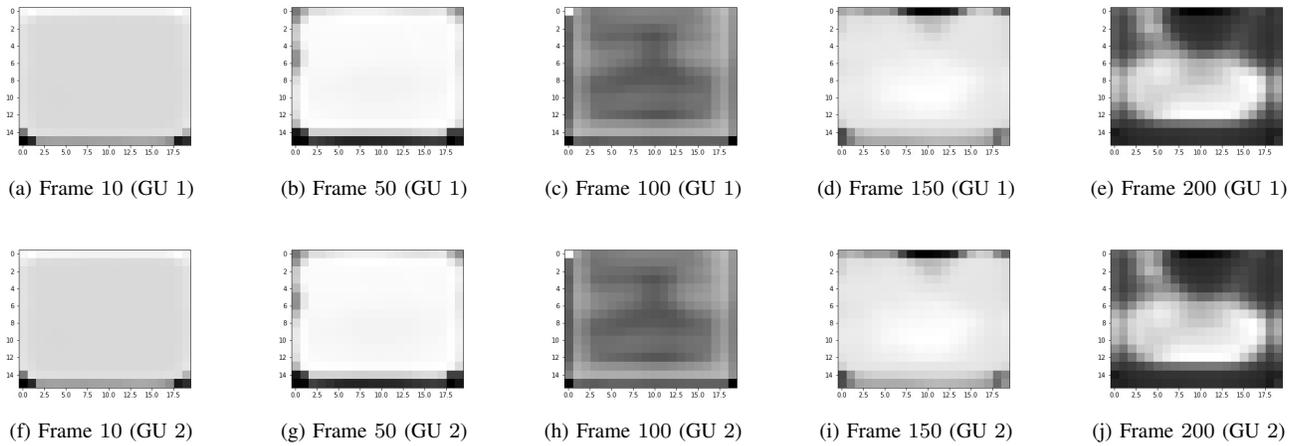


Fig. 6: Representation in GU Units (Layer 2) (Crazy turn)

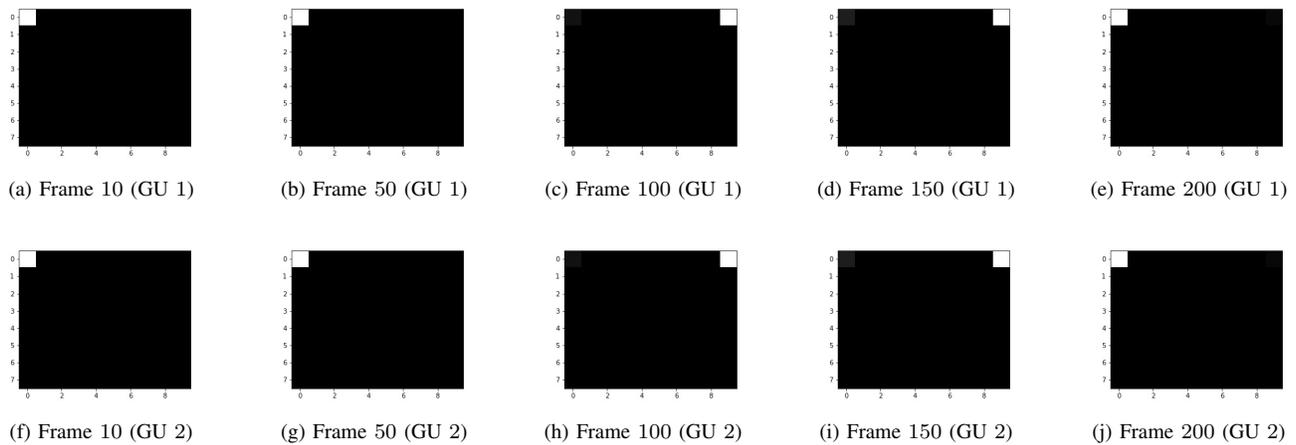


Fig. 7: Representation in GU Units (Layer 3) (Crazy turn)

- tural grammar”. In: *Psychonomic Science* 4.1 (1966), pp. 157–158.
- [6] K. Friston and S. Kiebel. “Cortical circuits for perceptual inference”. In: *Neural Networks* 22.8 (2009), pp. 1093–1104.
- [7] K. Friston, J. Mattout, and J. Kilner. “Action understanding and active inference”. In: *Biological cybernetics* 104.1 (2011), pp. 137–160.
- [8] G. Pezzulo, F. Rigoli, and K. Friston. “Active Inference, homeostatic regulation and adaptive behavioural control”. In: *Progress in Neurobiology* 134 (2015), pp. 17–35.
- [9] B. Han and R. VanRullen. “The rhythms of predictive coding? Pre-stimulus phase modulates the influence of shape perception on luminance judgments”. In: *Scientific reports* 7 (2017), p. 43573.
- [10] R. VanRullen. “Perceptual cycles”. In: *Trends in Cognitive Sciences* 20.10 (2016), pp. 723–735.
- [11] D. Kersten, P. Mamassian, and A. Yuille. “Object perception as Bayesian inference”. In: *Annual review of psychology* 55 (2004).
- [12] J. Hohwy, A. Roepstorff, and K. Friston. “Predictive coding explains binocular rivalry: An epistemological review”. In: *Cognition* 108.3 (2008), pp. 687–701.
- [13] E. von Holst and H. Mittelstaedt. “The reafference principle: Interaction between the central nervous system and the peripheral organs. Selected Papers of Erich von Holst: The Behavioural Physiology of Animals and Man”. In: (1950).
- [14] R. C. Miall and D. M. Wolpert. “Forward models for physiological motor control”. In: *Neural networks* 9.8 (1996), pp. 1265–1279.
- [15] N. L. Cerminara, R. Apps, and D. E. Marple-Horvat. “An internal model of a moving visual target in the lateral cerebellum”. In: *The Journal of physiology* 587.2 (2009), pp. 429–442.

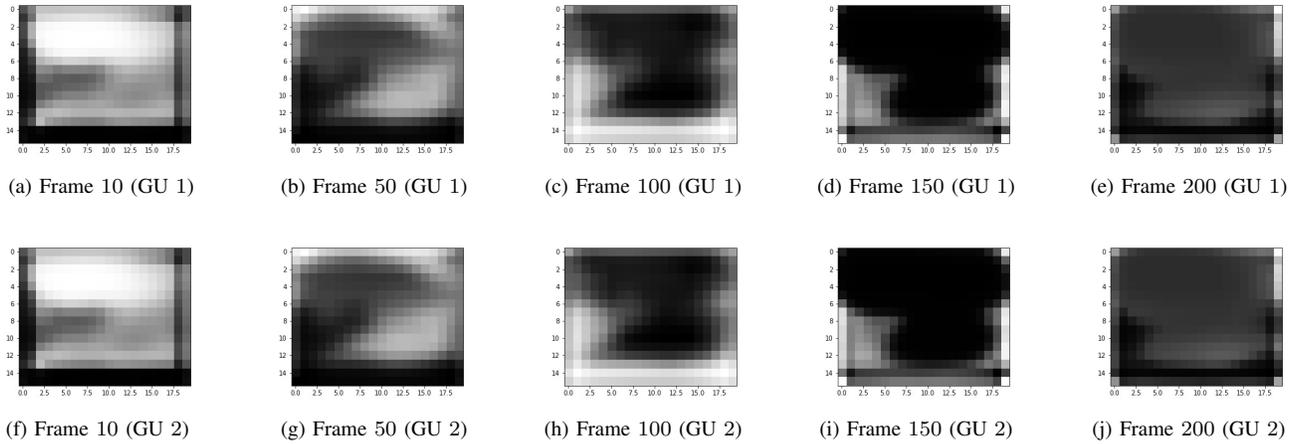


Fig. 8: Representation in GU Units (Layer 2) (Construction site)

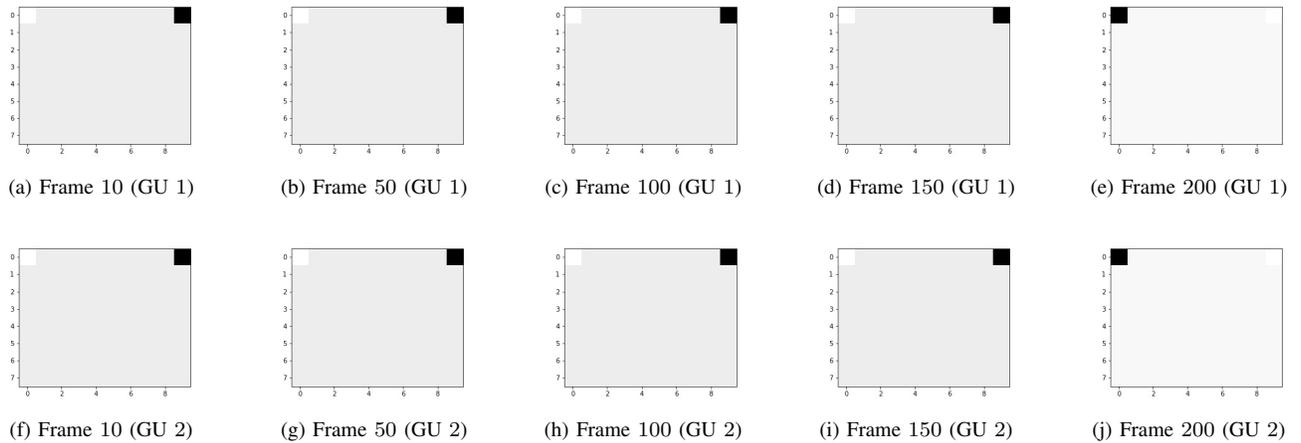


Fig. 9: Representation in GU Units (Layer 3) (Crazy turn)

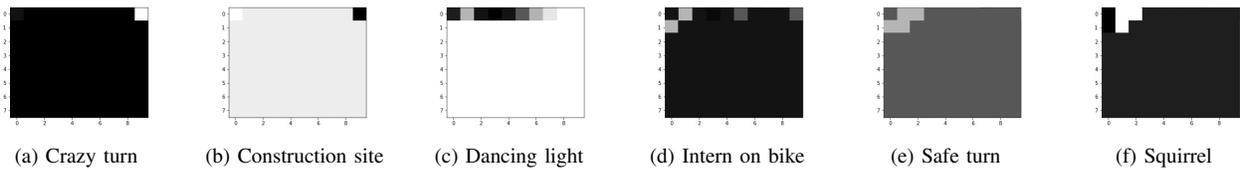


Fig. 10: Representation in GU Units (Layer 3, Frame 10) (6 Scenarios)

- [16] J. Zhong, C. Weber, and S. Wermter. “A Predictive Network Architecture for a Robust and Smooth Robot Docking Behavior”. In: *Paladyn. Journal of Behavioral Robotics* 3.4 (2012), pp. 172–180.
- [17] W. Lotter, G. Kreiman, and D. Cox. “Deep predictive coding networks for video prediction and unsupervised learning”. In: *arXiv preprint arXiv:1605.08104* (2016).
- [18] J. Zhong et al. “AFA-PredNet: The action modulation within predictive coding”. In: *International Joint Conference on Neural Networks (IJCNN)* (2018).
- [19] D. M. Wolpert, Z. Ghahramani, and M. I. Jordan. “An internal model for sensorimotor integration”. In: *Science* (1995), pp. 1880–1880.
- [20] D. M. Wolpert and M. Kawato. “Multiple paired forward and inverse models for motor control”. In: *Neural Networks* 11.7-8 (1998), pp. 1317–1329.
- [21] Y. Demiris and B. Khadhour. “Hierarchical attentive multiple models for execution and recognition of actions”. In: *Robotics and autonomous systems* 54.5 (2006), pp. 361–369.

- [22] H. Hoffmann. “Perception through visuomotor anticipation in a mobile robot”. In: *Neural Networks* 20.1 (2007), pp. 22–33.
- [23] R. Möller and W. Schenck. “Bootstrapping cognition from behavior: a computerized thought experiment”. In: *Cognitive Science* 32.3 (2008), pp. 504–542.
- [24] Y. Sugita and J. Tani. “Learning semantic combinatoriality from the interaction between linguistic and behavioral processes”. In: *Adaptive Behavior* 13.1 (2005), p. 33. ISSN: 1059-7123.
- [25] J. Zhong, A. Cangelosi, and S. Wermter. “Towards a self-organizing pre-symbolic neural model representing sensorimotor primitives”. In: *Frontiers in Behavioral Neuroscience* 8 (2014), p. 22.
- [26] Y. Yamashita and J. Tani. “Emergence of functional hierarchy in a multiple timescale neural network model: a humanoid robot experiment”. In: *PLoS Computational Biology* 4.11 (2008), e1000220.
- [27] T. Ogata and H. G. Okuno. “Integration of behaviors and languages with a hierarchical structure self-organized in a neuro-dynamical model”. In: *Robotic Intelligence In Informationally Structured Space (RiSS), 2013 IEEE Workshop on*. IEEE. 2013, pp. 89–95.
- [28] J. Zhong et al. “Sensorimotor Input as a Language Generalisation Tool: A Neurorobotics Model for Generation and Generalisation of Noun-Verb Combinations with Sensorimotor Inputs”. In: *arXiv preprint arXiv:1605.03261* (2016).
- [29] M. Kim, M. D. Singh, and M. Lee. “Towards Abstraction from Extraction: Multiple Timescale Gated Recurrent Unit for Summarization”. In: *arXiv preprint arXiv:1607.00718* (2016).
- [30] J. Zhong, A. Cangelosi, and T. Ogata. “Toward Abstraction from Multi-modal Data: Empirical Studies on Multiple Time-scale Recurrent Models”. In: *2017 International Joint Conference on Neural Networks (IJCNN)* (2017).
- [31] H. Lee, M. Jung, and J. Tani. “Recognition of visually perceived compositional human actions by multiple spatio-temporal scales recurrent neural networks”. In: *arXiv preprint arXiv:1602.01921* (2016).
- [32] Y. LeCun et al. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.
- [33] J. Donahue et al. “Long-term recurrent convolutional networks for visual recognition and description”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 2625–2634.
- [34] D. P. Kingma and J. Ba. “Adam: A method for stochastic optimization”. In: *arXiv preprint arXiv:1412.6980* (2014).