

Introduction

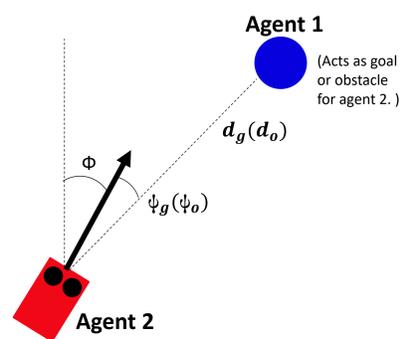
- Humans reliably attribute social interpretations to highly impoverished stimuli, such as interacting geometrical shapes (Heider and Simmel, 1944).
- Perception of animacy from such simple figures is dependent on a number of critical stimulus parameters (Tremoulet, Feldman 2000, 2006; Henrik et al., 2014).
- The perception of basic interactive actions, such as 'chasing' or 'fighting' has been addressed in several studies (Gao and Scholl 2013; Scholl and Tremoulet 2000; McAleer and Pollick 2000, Blythe et al., 1999); a set of six types of interactive movements has been repeatedly used in these studies.
- This perception of interaction has been explained by high-level cognitive processes, such as probabilistic reasoning and inference. (Baker et al., 2009)
- Building on classical biologically-inspired models for object and action perception (Riesenhuber and Poggio, 1999; Giese and Poggio, 2003), and a deep learning architecture (Simonyan, and Zisserman 2014) we propose a learning-based hierarchical neural network model that analyses such stimuli based shape and motion features directly from video sequences.
- The model has a simple feed-forward architecture and comprises two processing streams for form and object motion in the retinal frame of reference.
- The model contains only simple physiologically plausible operations.

Goal of our work

- Investigation if and how basic aspects of social and animacy perception can be accomplished by simple and physiologically plausible neural mechanisms, exploiting a hierarchical (deep) model of the visual pathway.

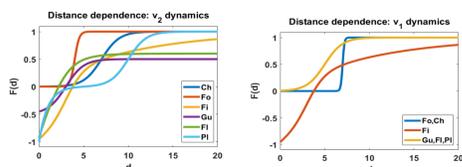
Generation of Stimuli

Modelling social interaction by modified human navigation model



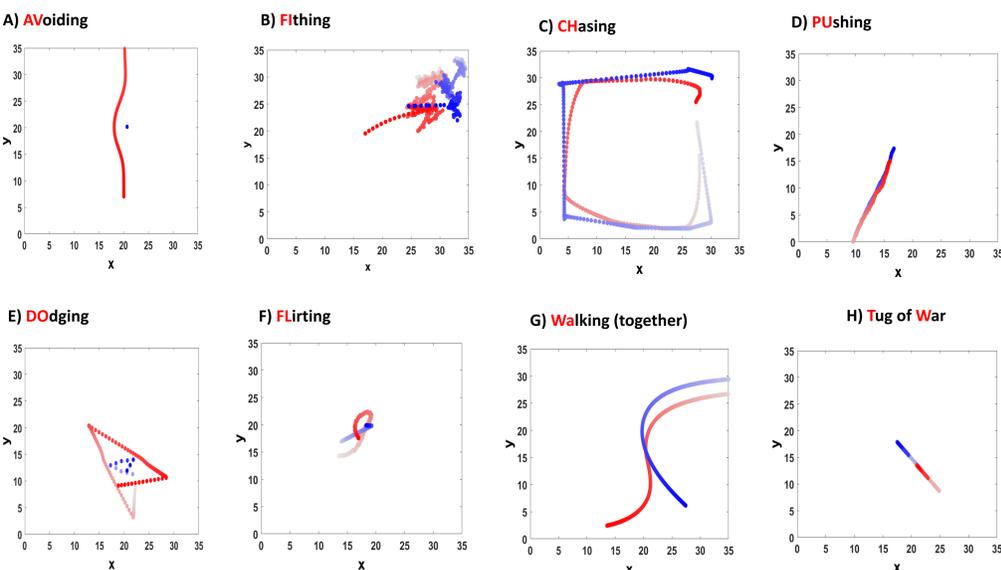
- Dynamics of heading direction (Fajen and Warren 2003):

$$\dot{\phi}_i = b\phi_i - k_g(\phi_i - \psi_{g,i})(e^{-c_1 d_{g,i}} + c_2) + k_o \sum_{n=1}^{N_{obst}} (\phi_i - \psi_{o,ni})(e^{-c_3 d_{o,ni}}) (e^{-c_4 d_{o,ni}})$$
- Dynamics of forward speed: $\tau \dot{v}_i = -v_i + F_i(d) + c_i \varepsilon_i(t)$
- Parameters fitted to movies by McAleer & Pollick (2008).



Sample trajectories from different intention categories (8 best recognized classes)

(Agent 1: blue, agent 2: red. Color saturation indicates time, the color fading with time.)



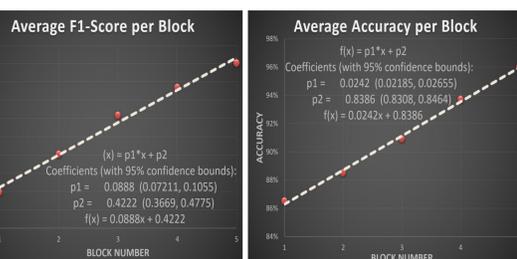
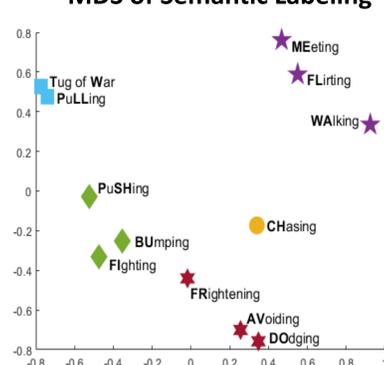
Psychophysical Experiment

- In the free labelling task, participants gave their own labels and interpreted each video freely.
- In the classification task, new subjects assigned labels to each video from the table of the most frequently proposed labels previously.
- In the last step, yet new subjects without watching videos, just gave pairwise rating to the semantic similarity of each label

Confusion Matrix

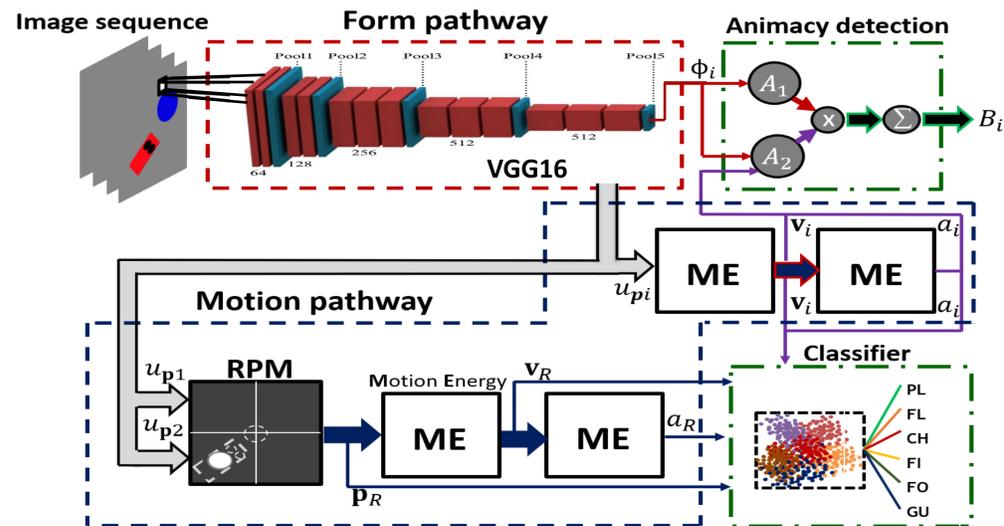
True Class	AV	BU	CH	DO	FI	FL	FR	ME	PLL	PSH	TW	WA
AV	43	2	8	1	2							1
BU	41			1	1	3	19					
CH	14	1	45	2	3	7	2					
DO	5	6	6	41	1	5						3
FI	6	4	17	47	1	7	4					
FL	1	5	7	1	44	4	4	1	2			1
FR	3			5	4	1	39	5	4	2	2	2
ME	1						9	2	45			6
PLL		8			2	7		39	12	2		
PSH		11			1	1	2	2	48			
TW		2		1	6				5	7	43	
WA		5		14	1		2		10			39
PPV	55.1%	51.2%	50.6%	71.9%	74.6%	72.1%	53.4%	65.2%	72.2%	47.5%	91.5%	79.6%
FDR	44.9%	48.7%	49.4%	28.1%	25.4%	27.9%	46.6%	34.8%	27.8%	52.5%	8.5%	20.4%

MDS of Semantic Labeling



- MDS result shows that misclassified labels are even semantically similar.
- Semantically similar actions can be distinguished from videos.
- F1-score and accuracy increase linearly with block numbers. Learning of categorization without explicit feedback!

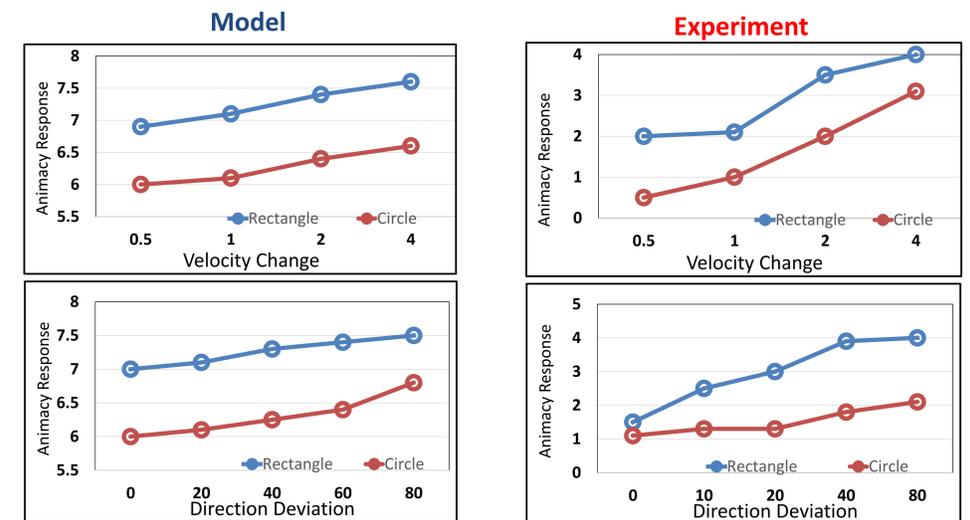
Model Architecture



- Hierarchical neural network with two pathways that analyze form and motion features.
- Mid-level features recognized by first layers of VGG16
- Two top levels that compute perceived animacy and classify perceived interaction.
- The choice of features for the computation of agency judgements was driven by results in the psychophysical literature.
- Critical features: absolute velocity and acceleration of agents, relative distance, velocity, and acceleration (McAleer and Pollick 2008).
- Testing multiple types of classifiers at the top level.

Results

Perception of animacy from the motion of a single object (Tremoulet, Feldman 2000)



- Consistent with the psychophysical results, the activity of the output 'agency neuron' increases with size of velocity and direction changes (testing trajectories where the agent followed a line and then suddenly changed direction or speed).
- Reproduction of increased animacy perception, compared to a moving circle (that does not have a body axis), if object has a body axis that is aligned with its velocity vector.

Social interaction classification

- Highest confusion rates between 'flirting' and 'chasing'; sometimes also 'playing' and 'guarding'.
- For all classifier types accuracy is at least 94 %, best classification result is obtained with linear support vector machine, reaching an accuracy of 99 %.
- All original videos from McAleer and Pollick (2008) were classified correctly, though they were not part of the training set.

Accuracy of different classifiers

Classifier	Accuracy
Linear SVM	99.0%
Gaussian kernel SVM	96.3%
LDA	94.7%
KNN	94.7%
Nonlinear LDA	94.3%
Neural Network	94.0%

Conclusions

- While our model is a quite simple but physiologically plausible it was able to reproduce several important characteristics of human perception of agency and of social interactions from strongly impoverished displays.
- Since the model includes a deep network for form with enough training data it can likely be extended for real video stimuli.

References

- Heider, F. and Simmel, M.: An Experimental Study of Apparent Behavior. The American Journal of Psychology (1944)
- Tremoulet, P.D., Feldman, J.: Perception of animacy from the motion of a single object. Perception 29, 943–951 (2000)
- Tremoulet, P. D. and Feldman, J.: The influence of spatial context and the role of intentionality in the interpretation of animacy from motion. Perception and psychophysics (2006)
- Hernik, M., Fearon, P., and Csibra, G.: Action anticipation in human infants reveals assumptions about anteroposterior body structure and action. Proceedings. Biological sciences (2014)
- Gao, T. and Scholl, B. J.: Perceiving animacy and intentionality. In Rutherford, M. D. and Kuhlmeier, V. A., editors, Social Perception. The MIT Press (2013)
- Scholl, B. J. and Tremoulet, P. D.: Perceptual causality and animacy. Trends in Cognitive Sciences, 4(8):299–309 (2000)
- McAleer, P., Pollick, F.E.: Understanding intention from minimal displays of human activity. Behavior Research Methods 40, 830–839 (2008)
- Fajen, B.R., Warren, W.H.: Behavioral dynamics of steering, obstacle avoidance, and route selection. Journal of Experimental Psychology: HPP (2003)
- Giese, M.A., Poggio, T.: Neural mechanisms for the recognition of biological movements. Nat Rev Neurosci 4, 179–192 (2003)
- Riesenhuber, M., Poggio, T.: Hierarchical models of object recognition in cortex. Nat. Neurosci. 2, 1019–1025 (1999)
- Blythe, P. W., Todd, P. M., & Miller, G. F.: How motion reveals intention: Categorizing social interactions. Simple heuristics that make us smart, pp. 257-285(1999)
- Baker, C.L., Saxe, R., Tenenbaum, J.B.: Action understanding as inverse planning. Cognition, Reinforcement learning and higher cognition 113, 329–349 (2009)
- Hovaidi-Ardestani M., Saini, N., Martinez, A. & Giese, M. A. (2018). Neural model for the visual recognition of animacy and social interaction. ICANN, Greece
- Simonyan, A. Zisserman., Very Deep CNNs for Large-Scale Visual Recognition.

Acknowledgements

This work was supported by: HFSP RGP0036/2016; the European Commission HBP FP7-ICT2013-FET-F/ 604102 and COGIMON H2020-644727, and the DFG GZ: GI 305/4-1 and KA 1258/15-1 as well as FP7-611909, DFG GZ: KA 1258/15-1, BMBF FKZ 01GQ1704, KONSENS BW Stiftung NEU007/1