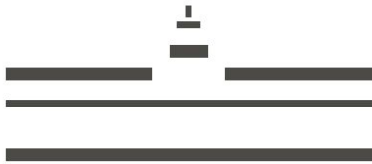


EYESHOTS
Kick-off Meeting

Attention and 3D-Object recognition

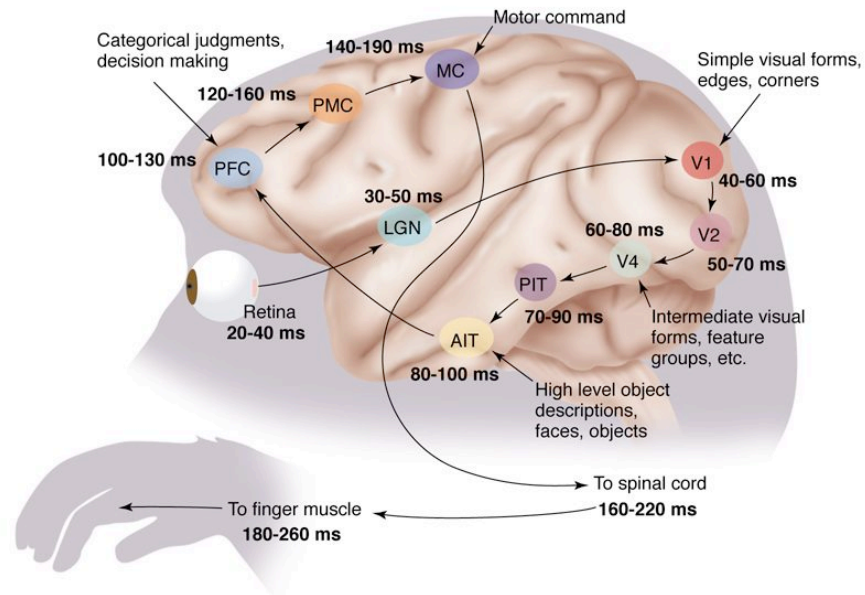
PD Dr. Fred Hamker



Westf. Wilhelms-University Münster, Germany
Department of Psychology
Otto-Creutzfeldt Center for Cognitive and Behavioral
Neuroscience

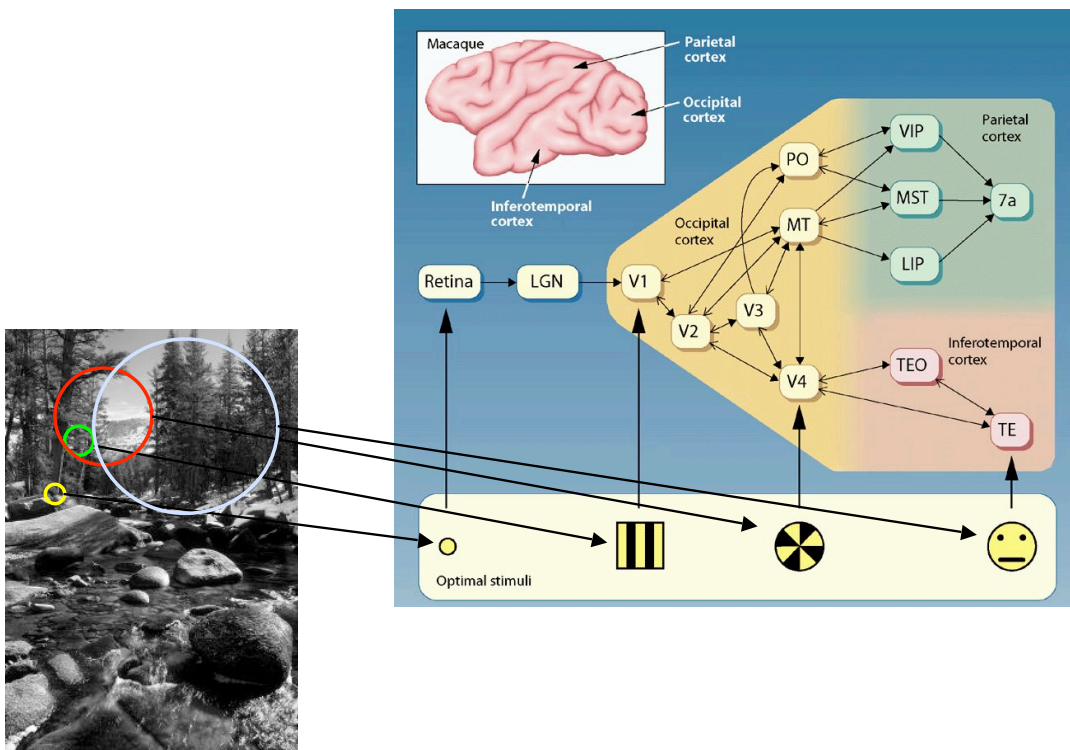


Object recognition in primates



Thorpe, S.J. & Fabre-Thorpe, M. (2001) Seeking categories in the brain. *Science*, 291:260-263.

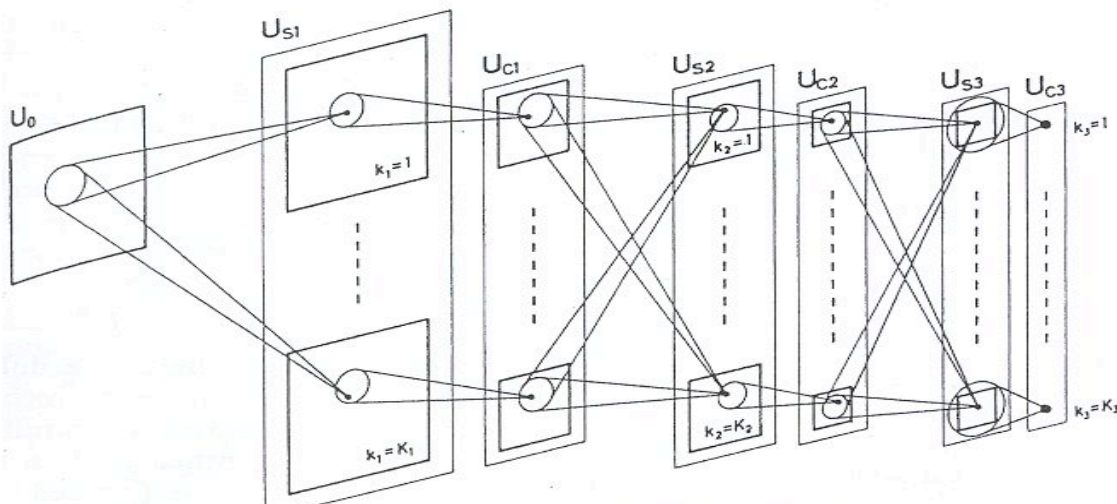
The ventral and dorsal pathway



Models of object recognition

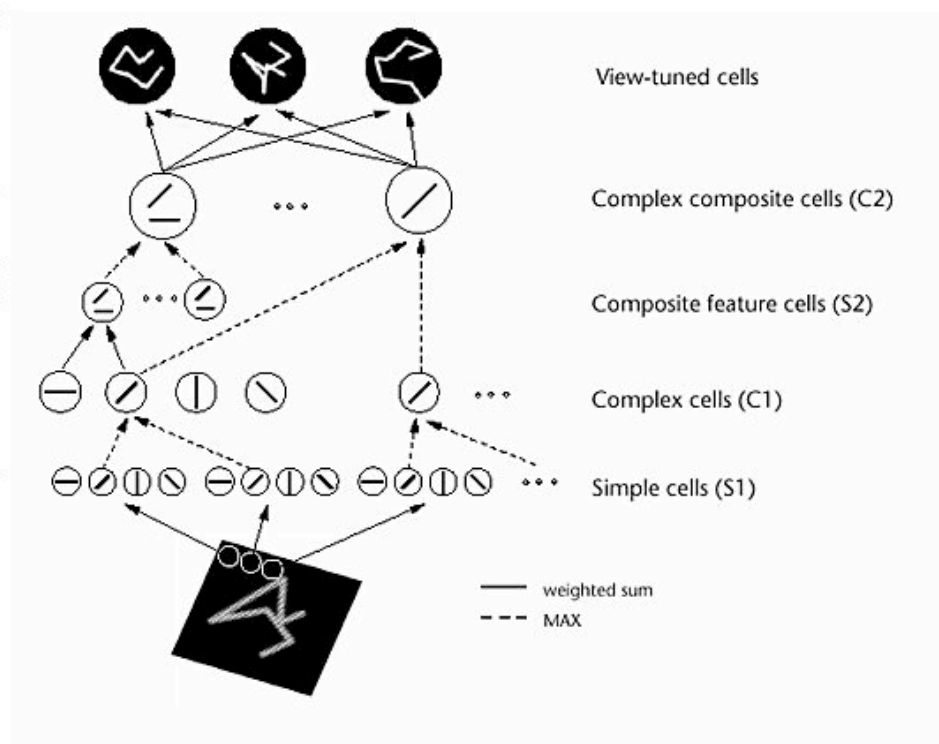
- Fukushima, K (1980) Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position, Biol. Cybern. 36:193-202.
- Riesenhuber, M & Poggio, T (1999) Hierarchical models of object recognition in cortex, Nat. Neurosci. 2:1019-1025.

Position invariant recognition in the Neocognitron (Fukushima 1980)

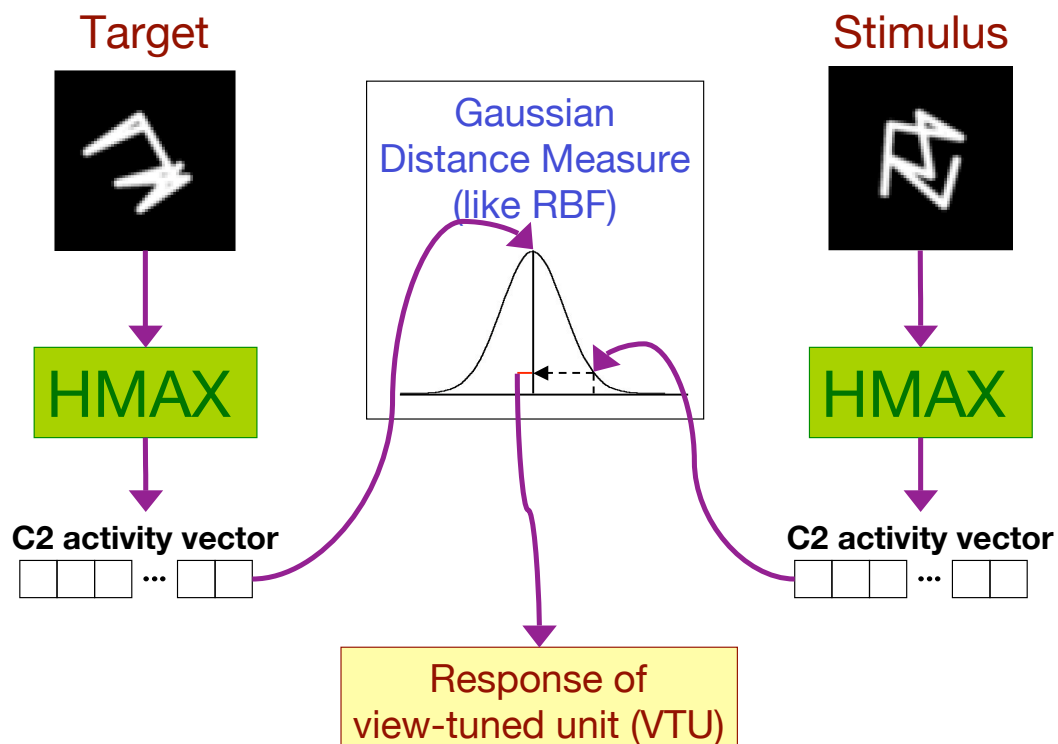


Several processing layers, comprising simple (S) and complex (C) cells.
S-cells in one layer respond to conjunctions of C-cells in previous layer.
C-cells in one layer are excited by small neighborhoods of S-cells.

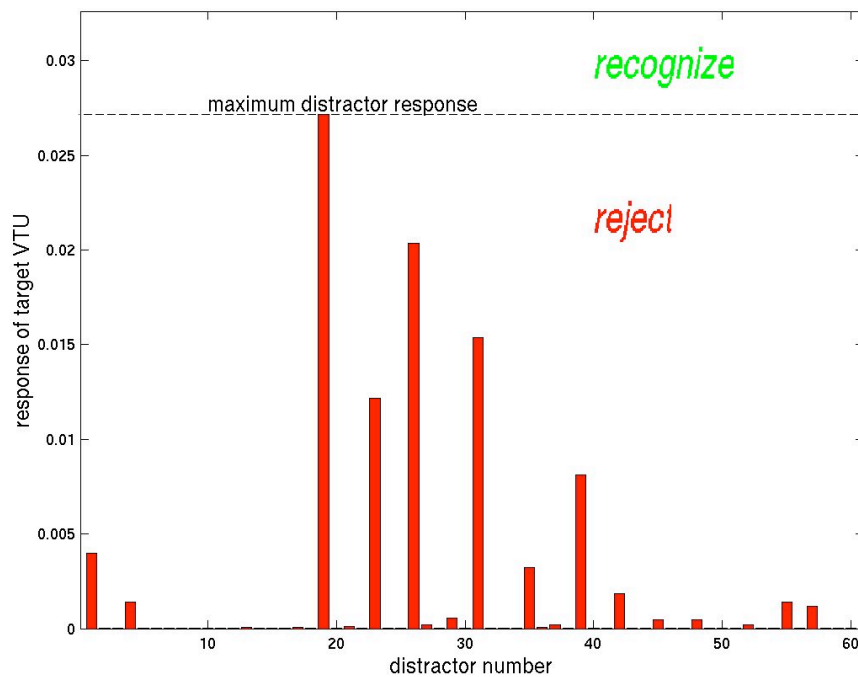
HMAX (Riesenhuber and Poggio, 1999)



Details of recognition in HMAX

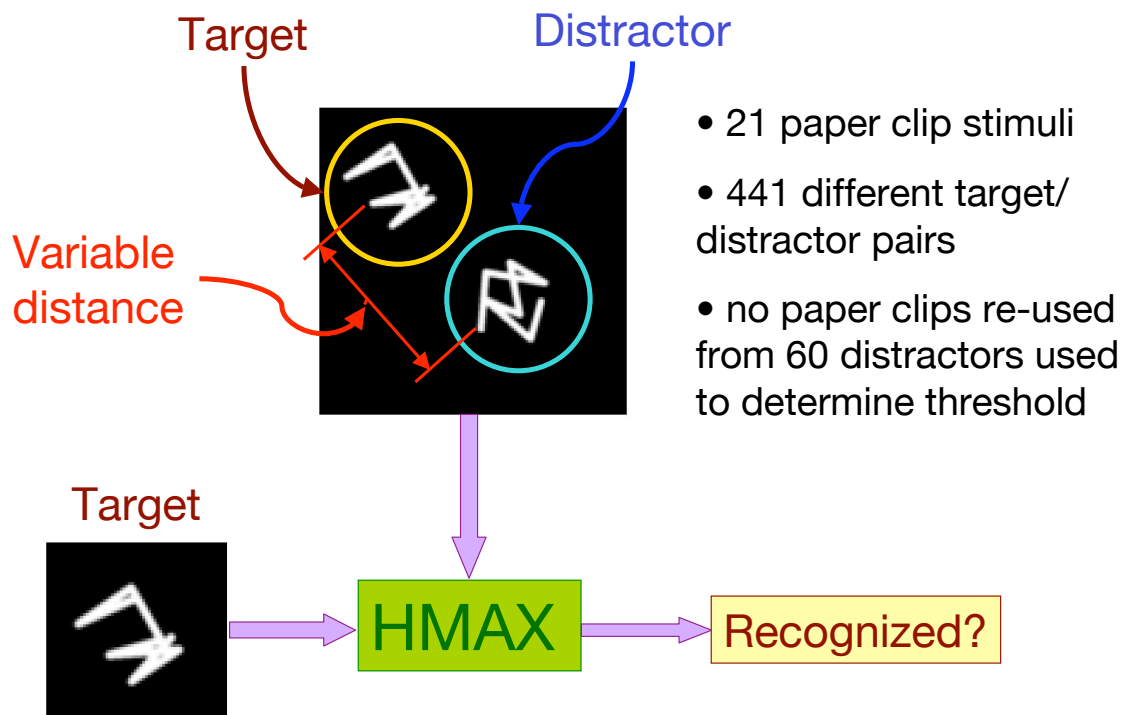


When is a stimulus recognized ?

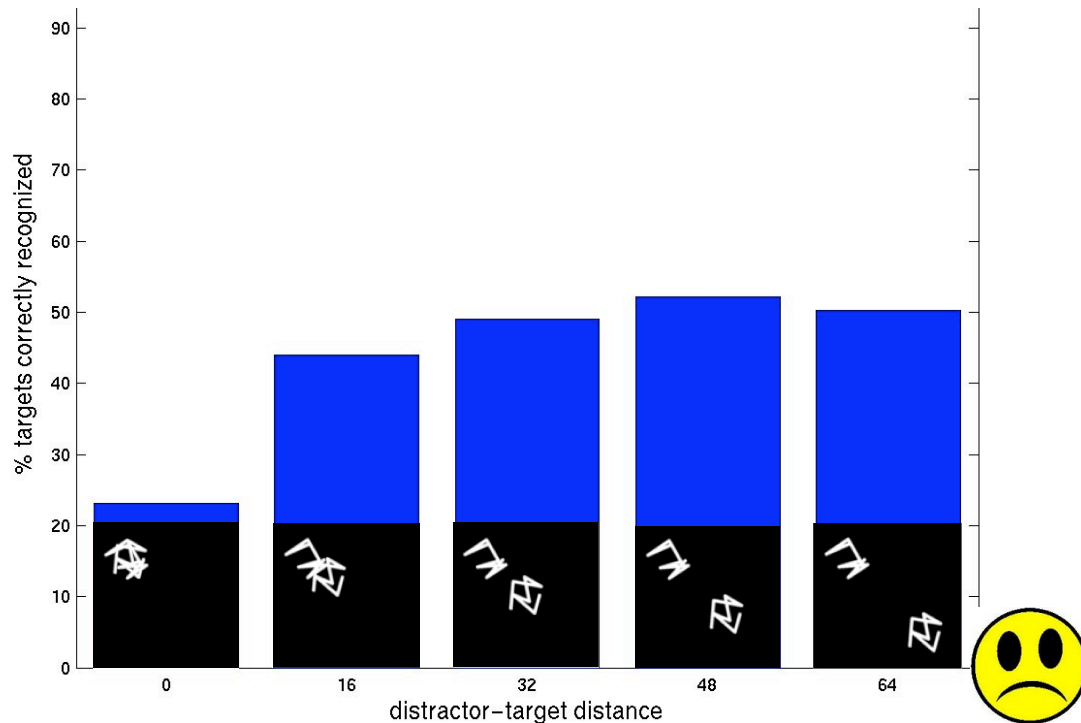


60 randomly chosen distracter paper clips

Two paper clips



Two paper clips - Results

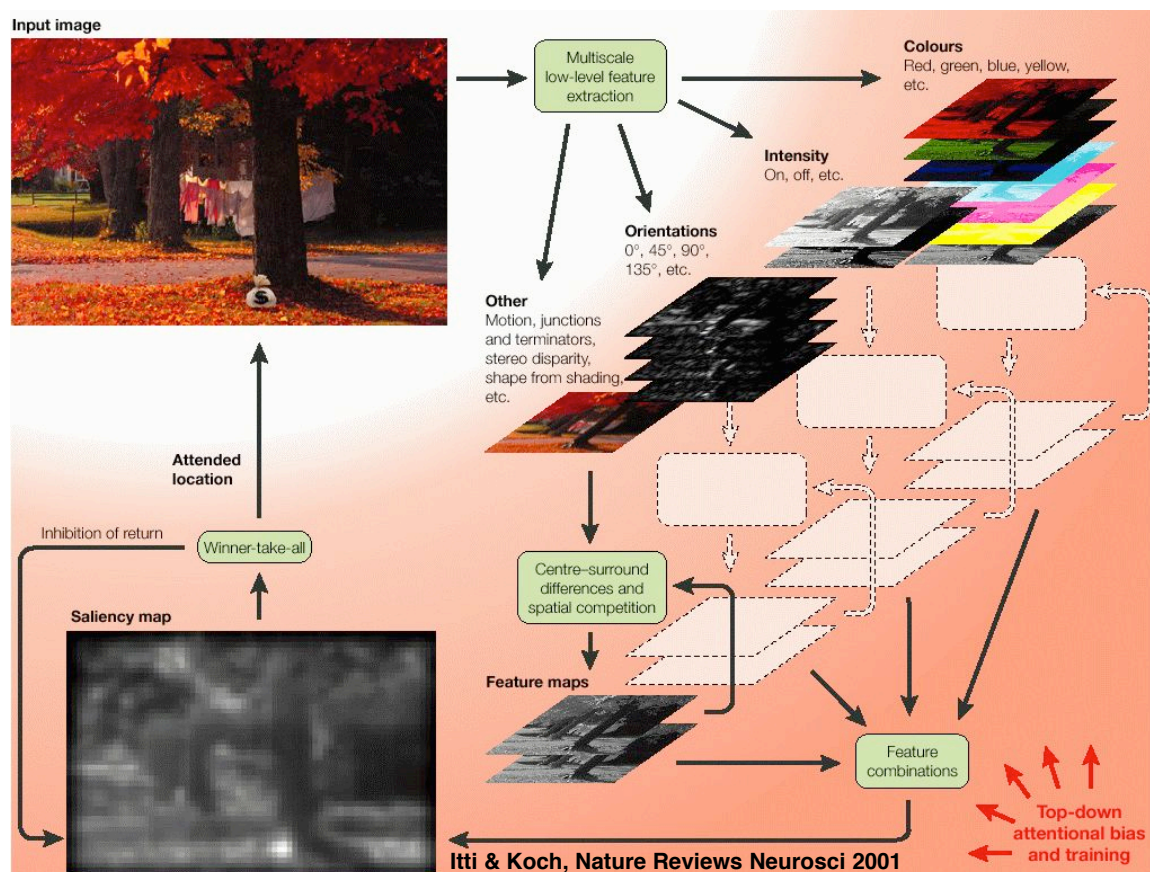


Hierarchical Template Matching for Object Recognition

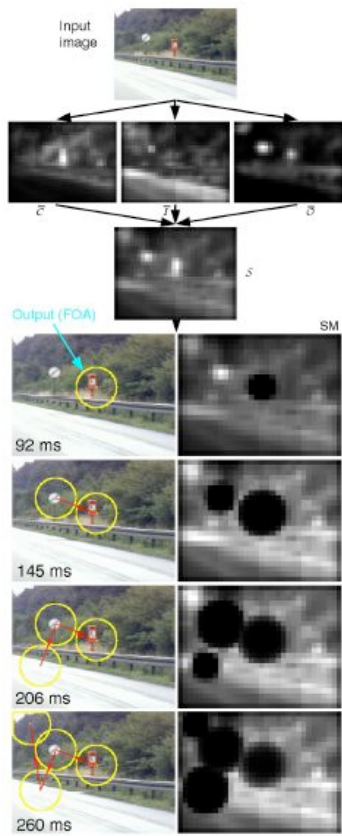
- Image passed through layers of units with progressively more complex features at progressively less specific locations.
- Hierarchical in that features at one stage are built from features at earlier stages.
- Processing hierarchy yields activation of view-tuned units.
- A collection of view-tuned units is associated with one object.
- Object recognition is severely impaired in the presence of clutter
- At present, no learning algorithm for tuning the weights has been developed (but see Wersing and Körner 2003 and LeChun, 1998).

The saliency map model of attention

- Itti, L., Koch, C., Niebur, E. (1998) A model of saliency-based visual attention for rapid scene analysis. IEEE Transactions on Pattern Analysis and Machine Intelligence, 20:1254-1259.
- Itti, L., Koch, C. (2000) A saliency-based search mechanism for overt and covert shifts of visual attention. Vision Res., 40:1489-1506.



Example

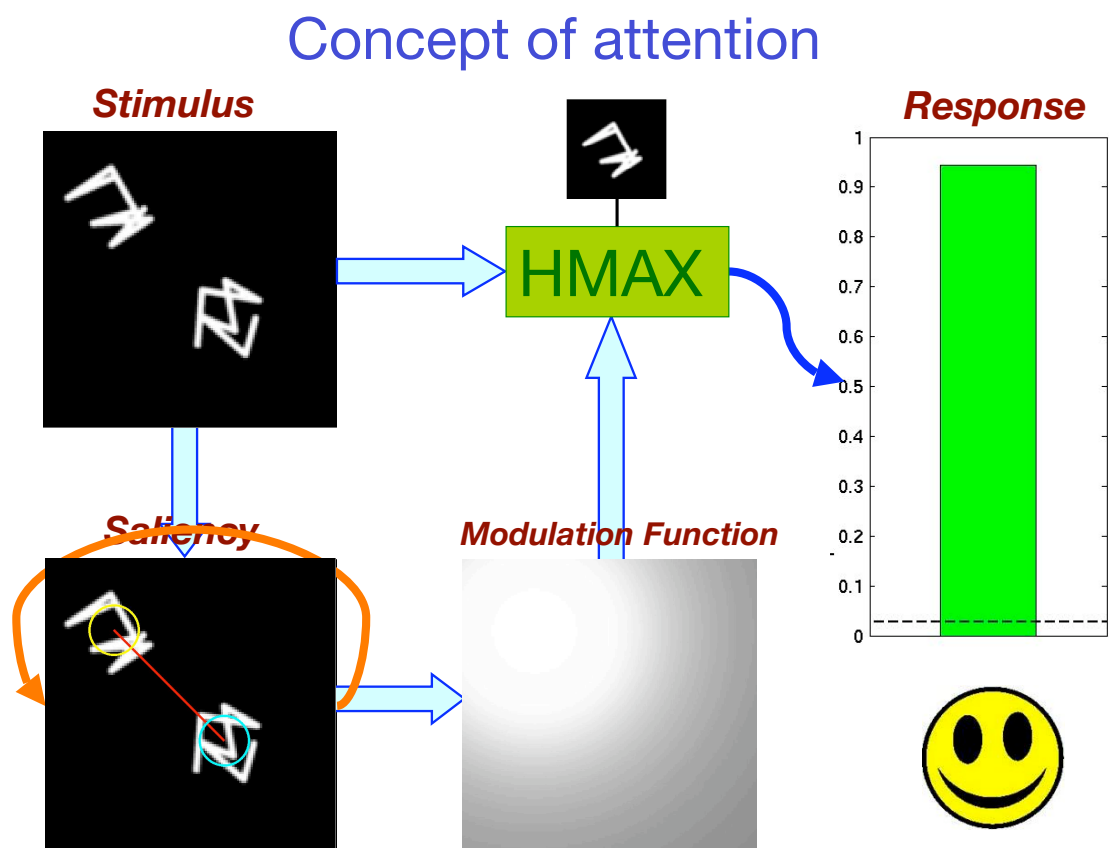
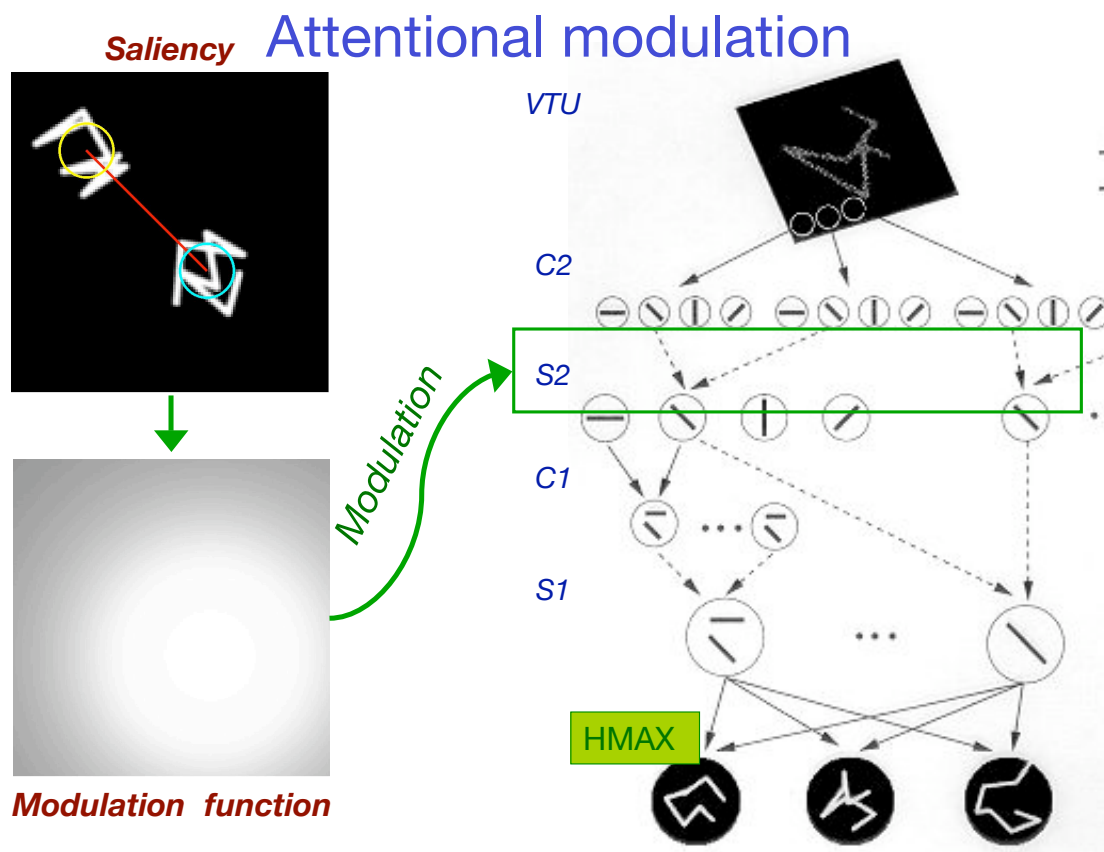


Discussion

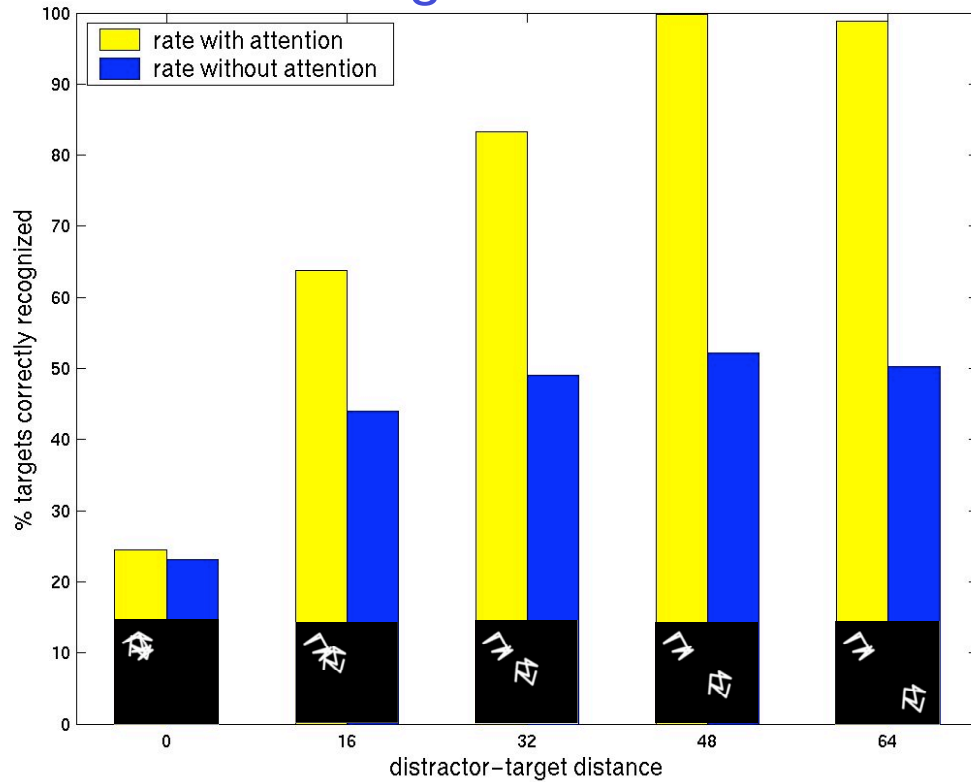
- The saliency model offers a fast algorithm for guiding vision to potentially meaningful parts of a scene.
- It selects only a point in space, as compared to an object or region. Region selection has to be added by a separate mechanism.
- Saliency is restricted to simple features.
- Attention is defined solely as the selection in space (no, or only indirect feature-based selection).
- The advantage of this mechanism for object recognition is limited, since a selection in space does not necessarily promote object-recognition.

Combining the saliency map with hierarchical models of object recognition

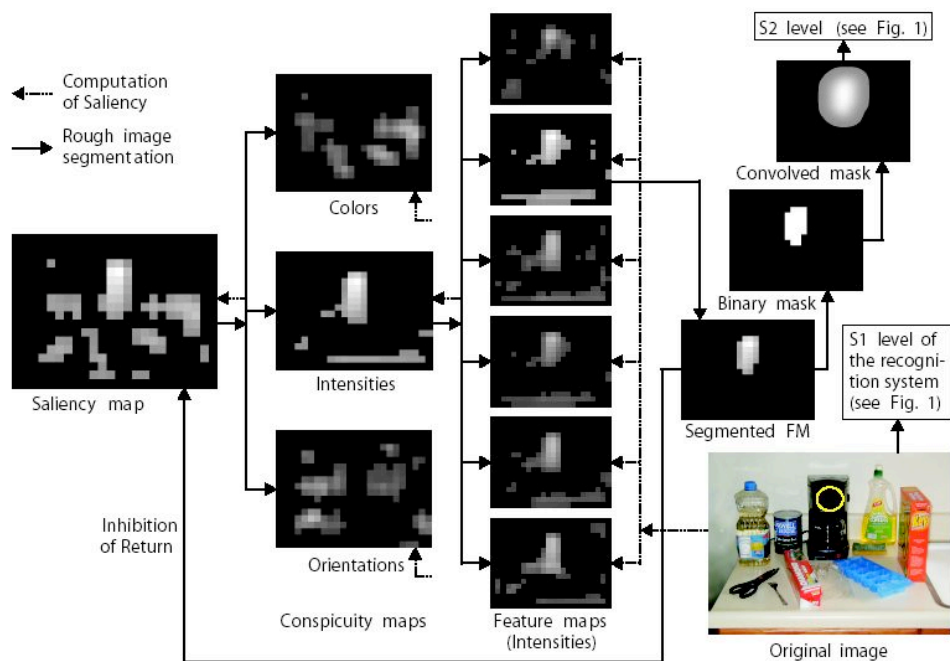
- Walther, Itti, Riesenhuber, Poggio, Koch (2002) Attentional Selection for Object Recognition - a Gentle Way. In: Biologically Motivated Computer Vision. Lecture Notes in Computer Science. Berlin, Heidelberg, New York: Springer Verlag, 472-479.



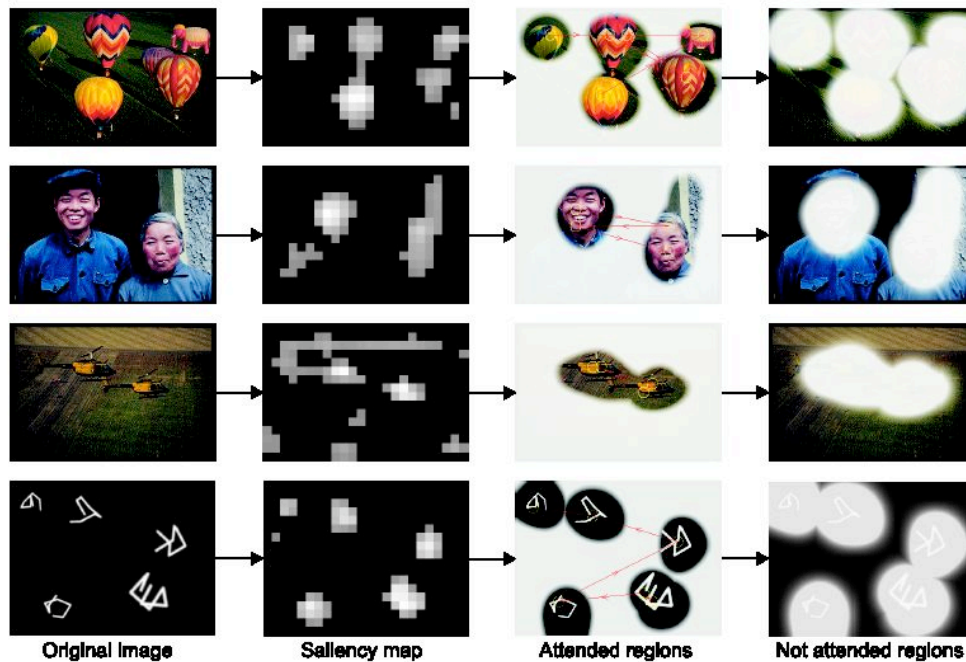
Recognition Rate



Towards attention and object recognition in natural scenes



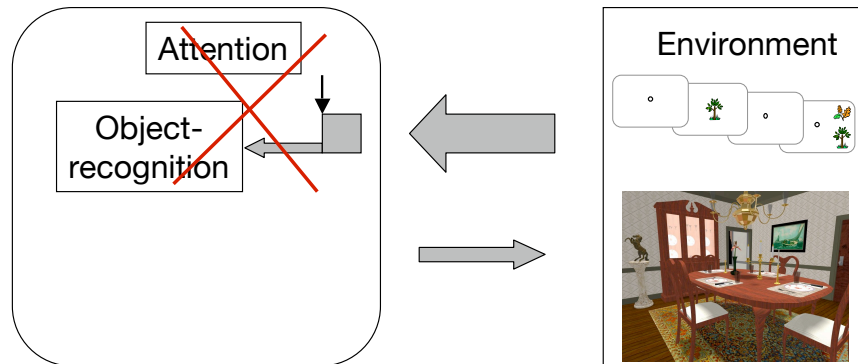
Towards attention and object recognition in natural scenes



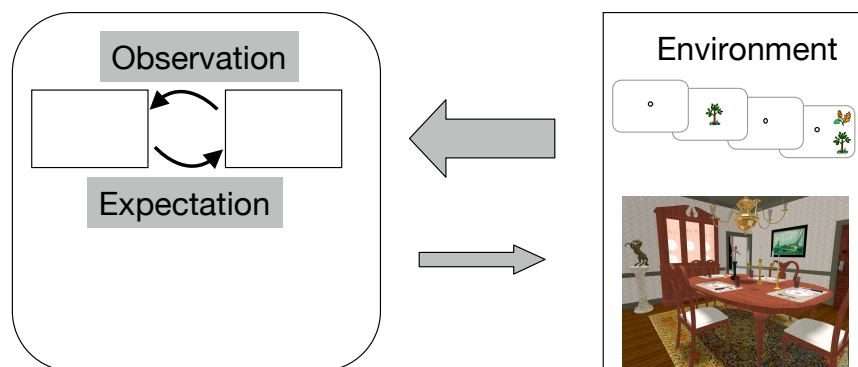
Discussion

- The combination of the Saliency-Map model with a spatially selective focus and a hierarchical model for object recognition appears to be a straightforward way to go.
- Recognition depends on the quality of the focus.
- The focus is not determined by the recognition task.
- The model predicts that prior selection is necessary for object recognition, which appears to be a contradiction to the ability of category detection in dual-task situations.

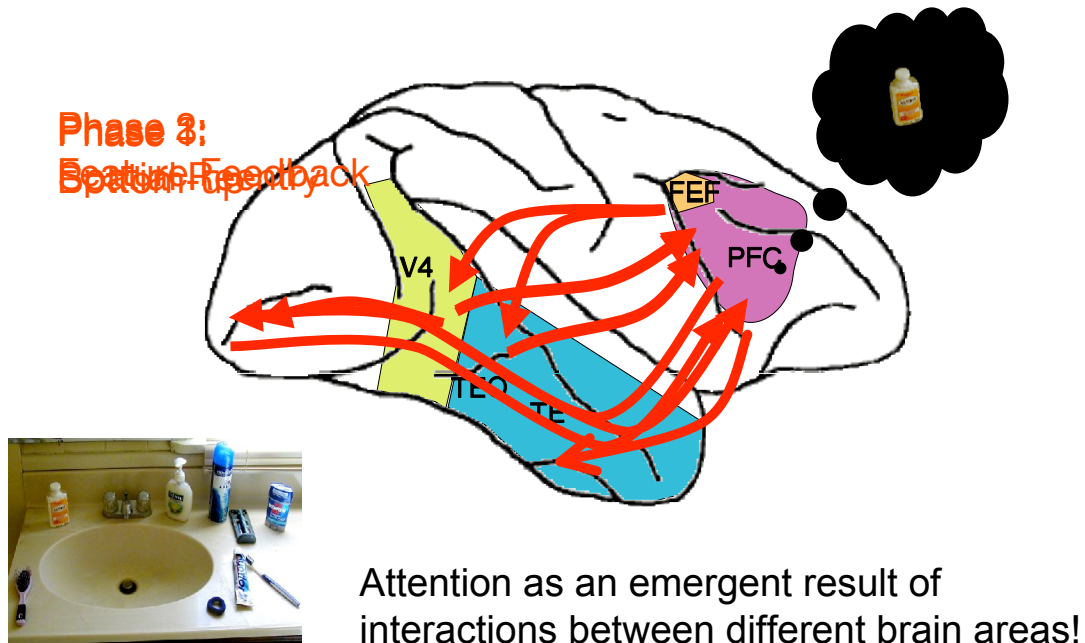
Classical approach of visual attention and object recognition



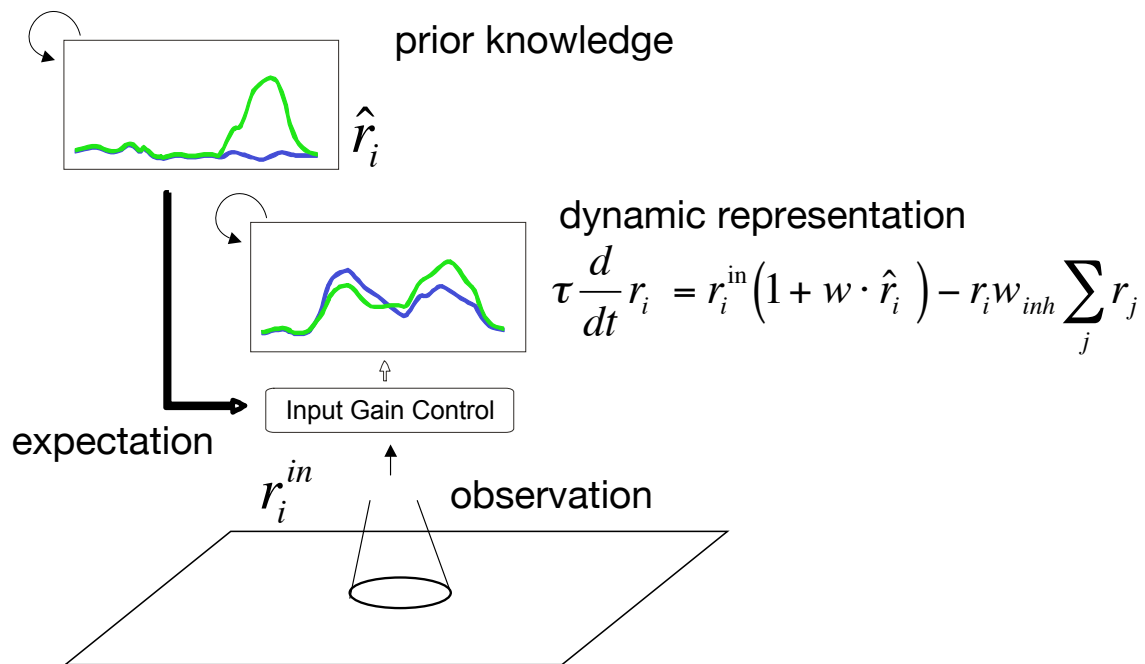
Alternative approach of visual attention and object recognition



The three phases of visual perception



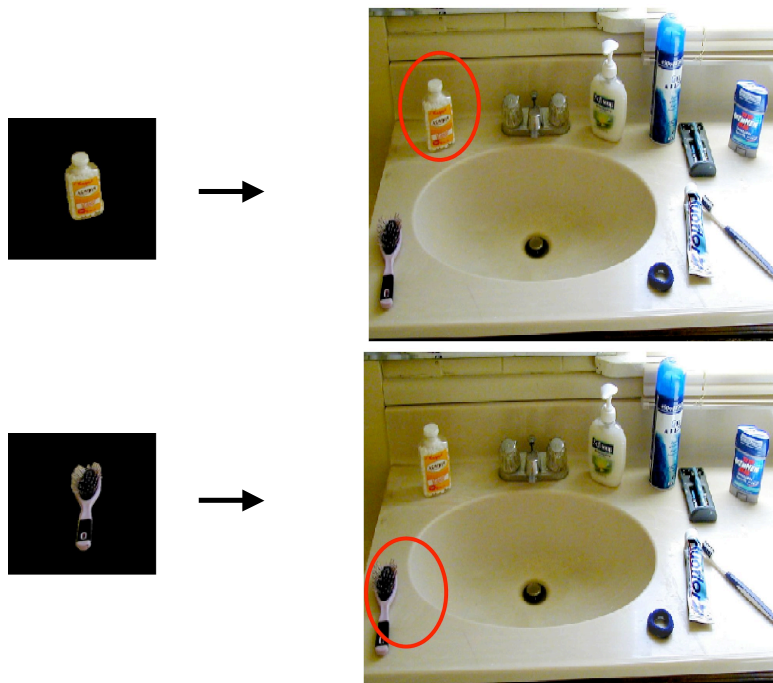
The concept of population-based inference

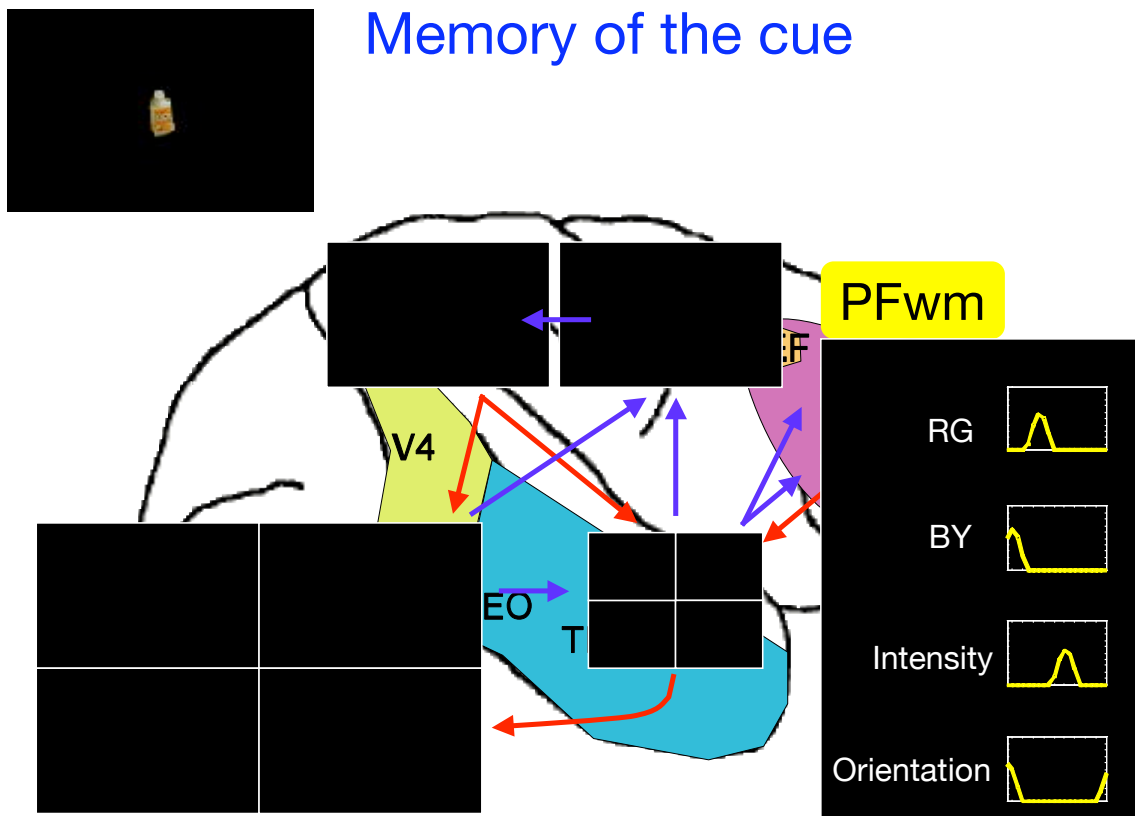
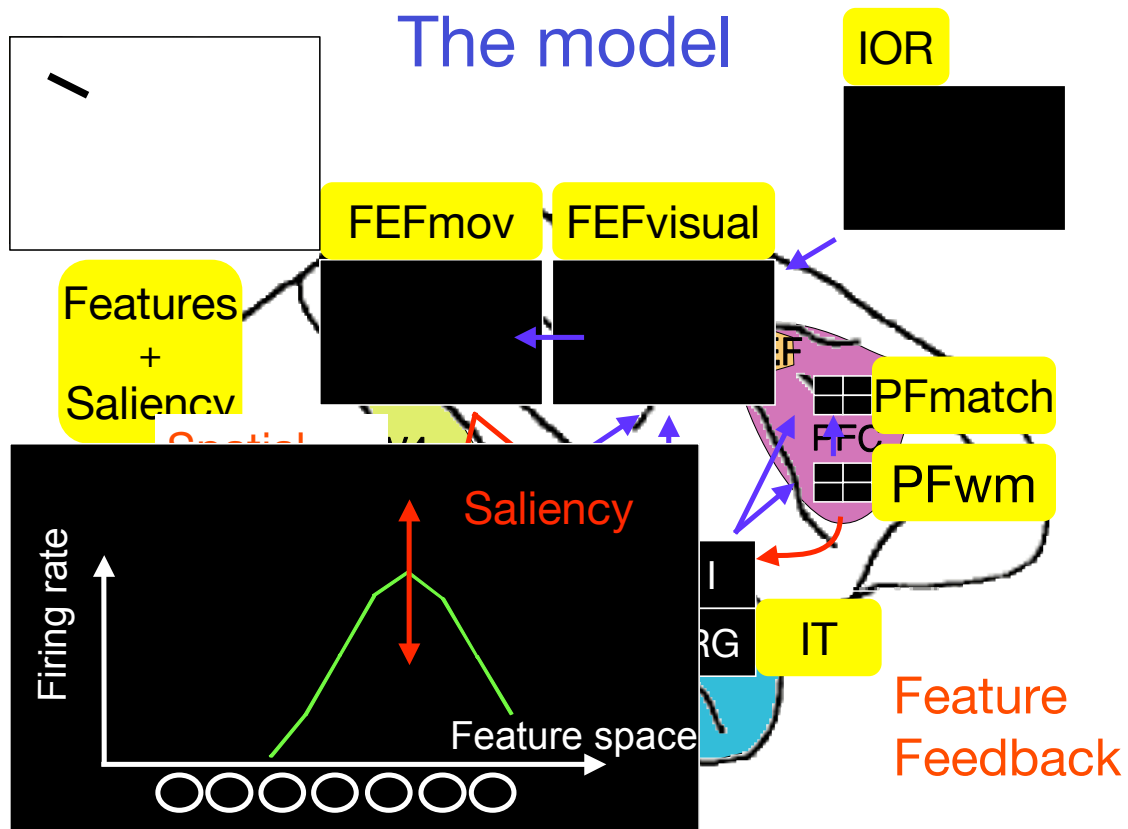


Search “without Target”

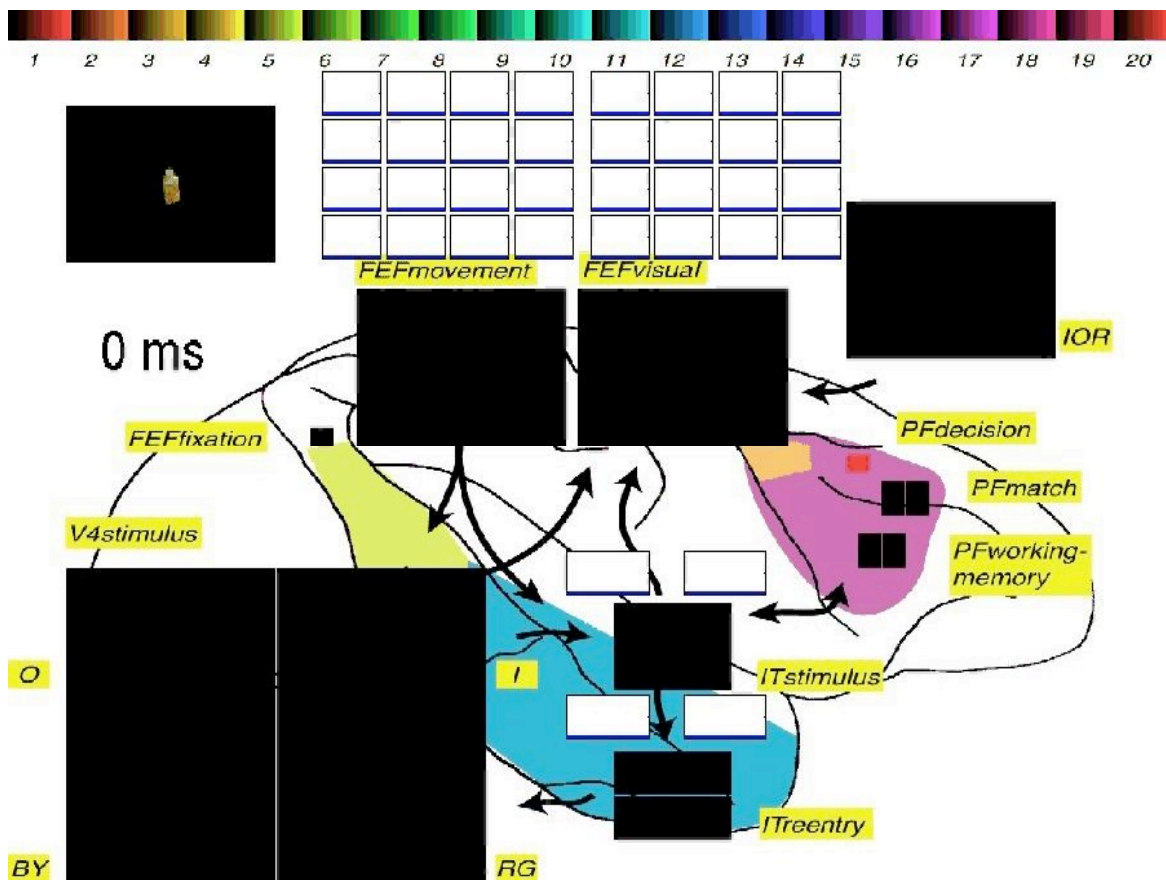
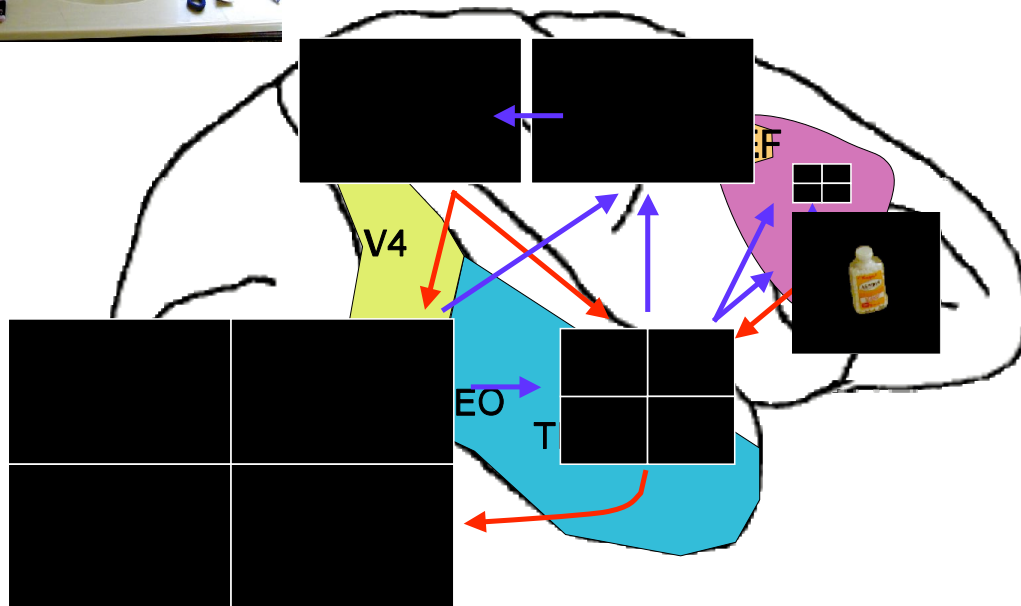


Search with Target

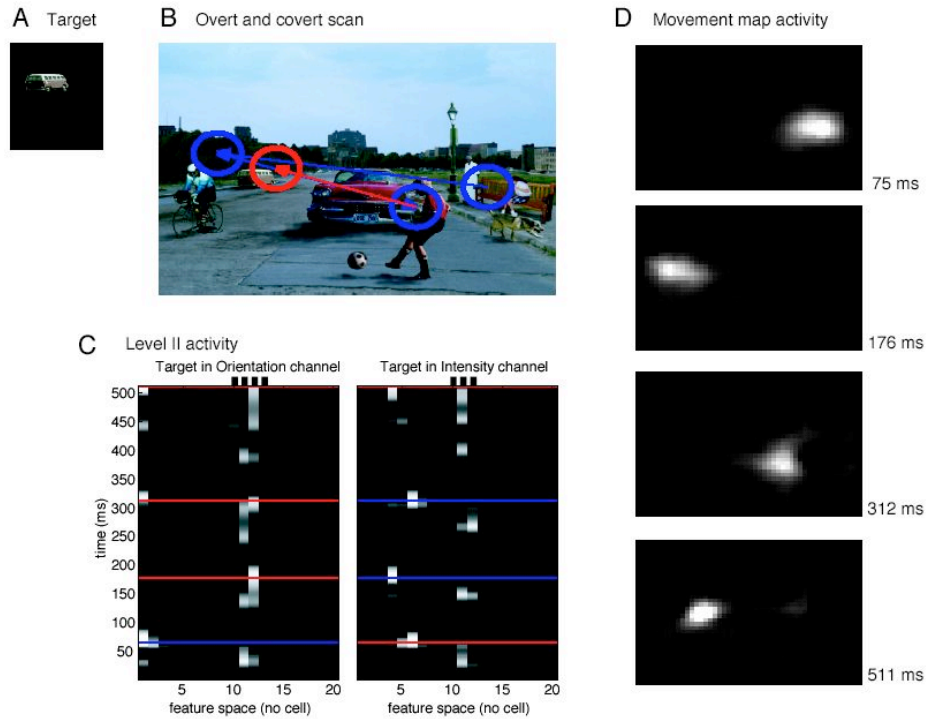




Goal-directed perception

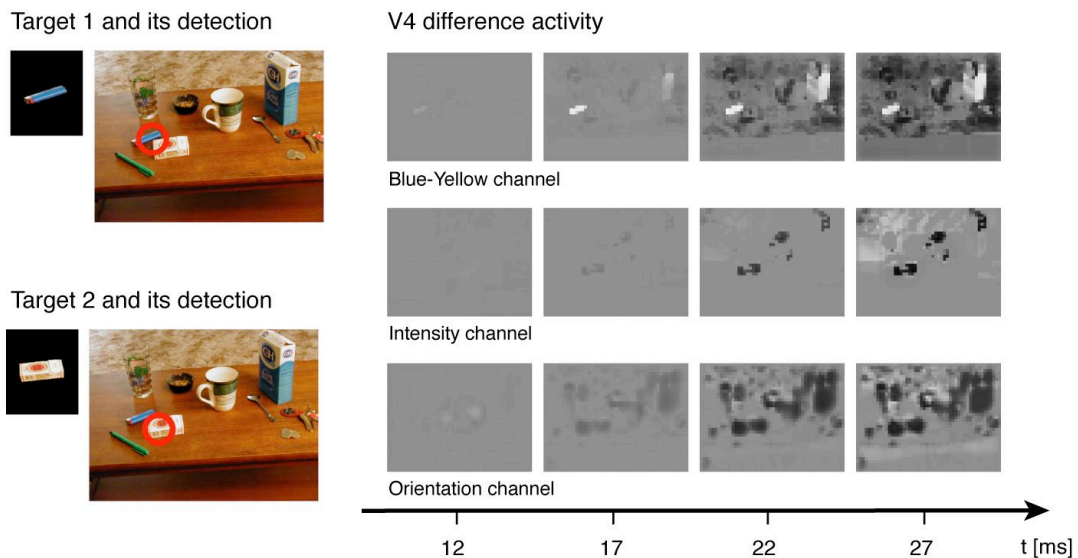


Overt and covert attention



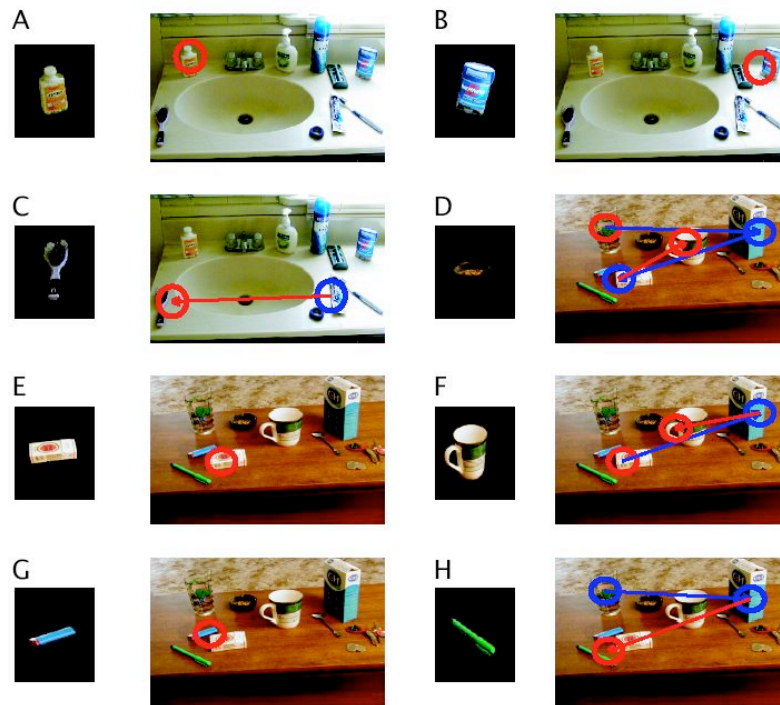
Hamker, FH (2005) Journal for Computer Vision and Image Understanding.

Feature-based attention in natural scenes



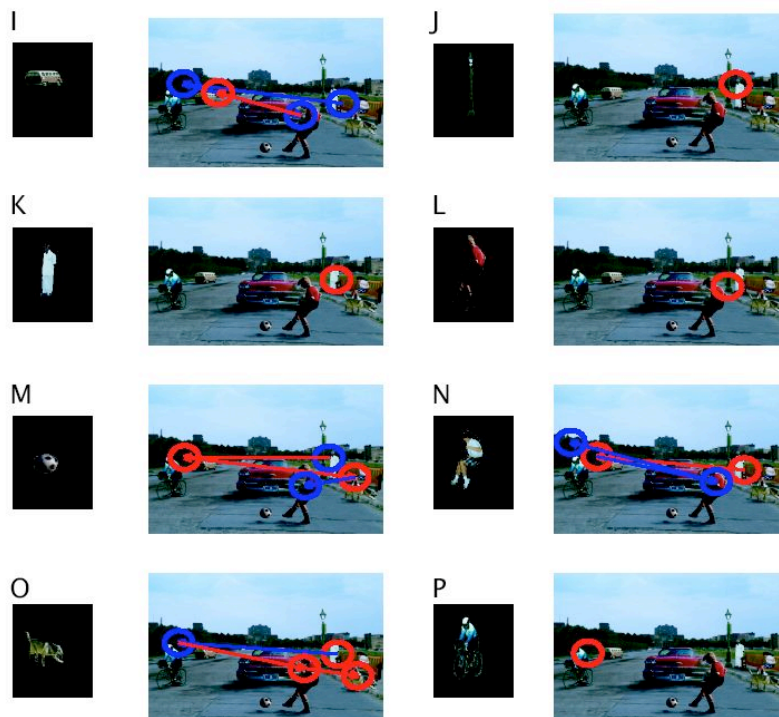
The model predicts that prior to any spatial selection, V4 contains information about potential target objects - feature-based attention.

Visual search examples



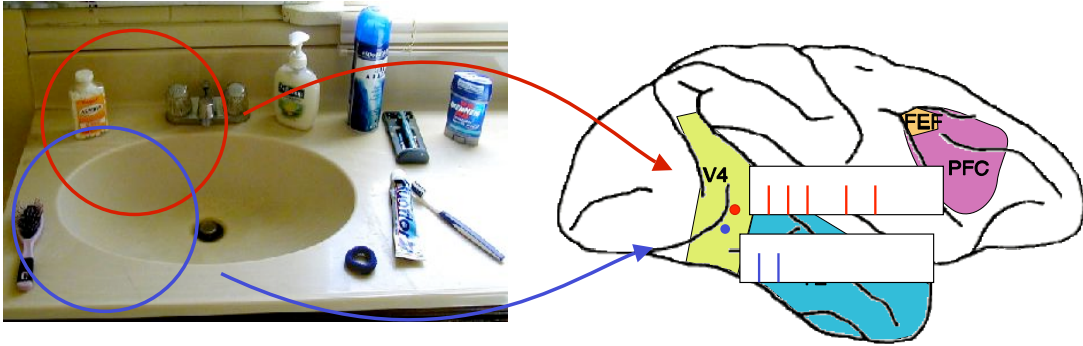
Hamker, FH (2005) Journal for Computer Vision and Image Understanding.

Visual search examples

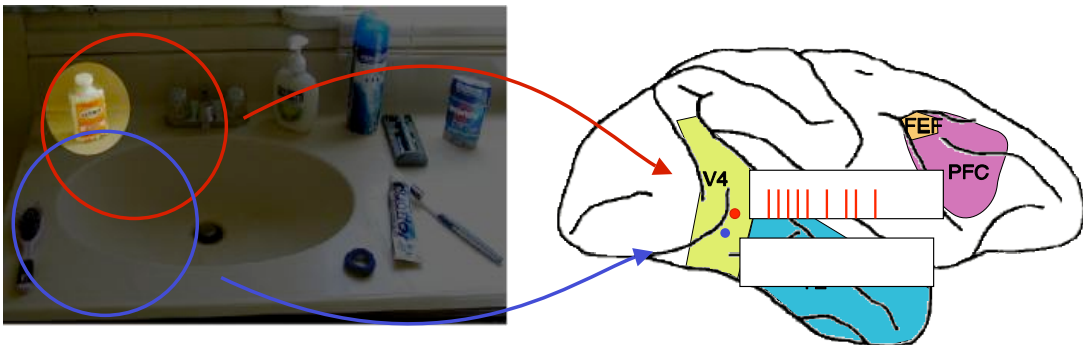


Hamker, FH (2005) Journal for Computer Vision and Image Understanding.

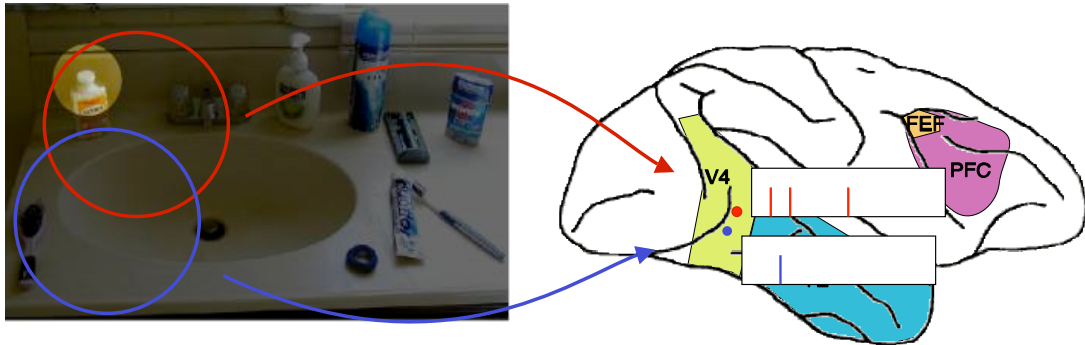
How does attention facilitate
object recognition ?



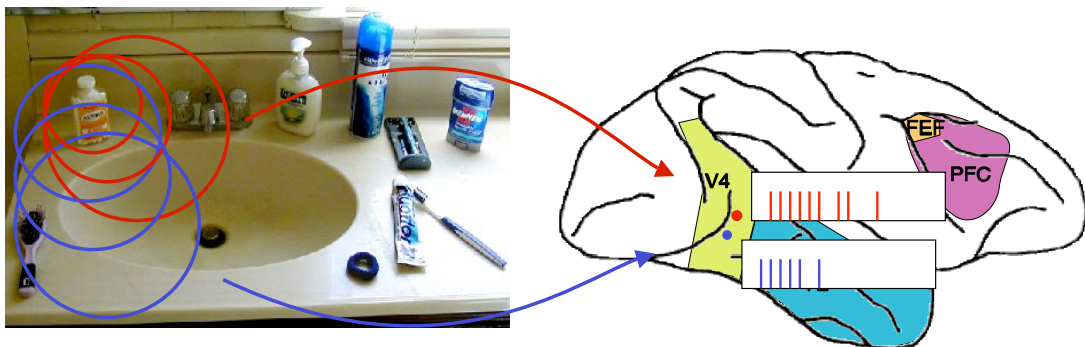
Spotlight Metaphor

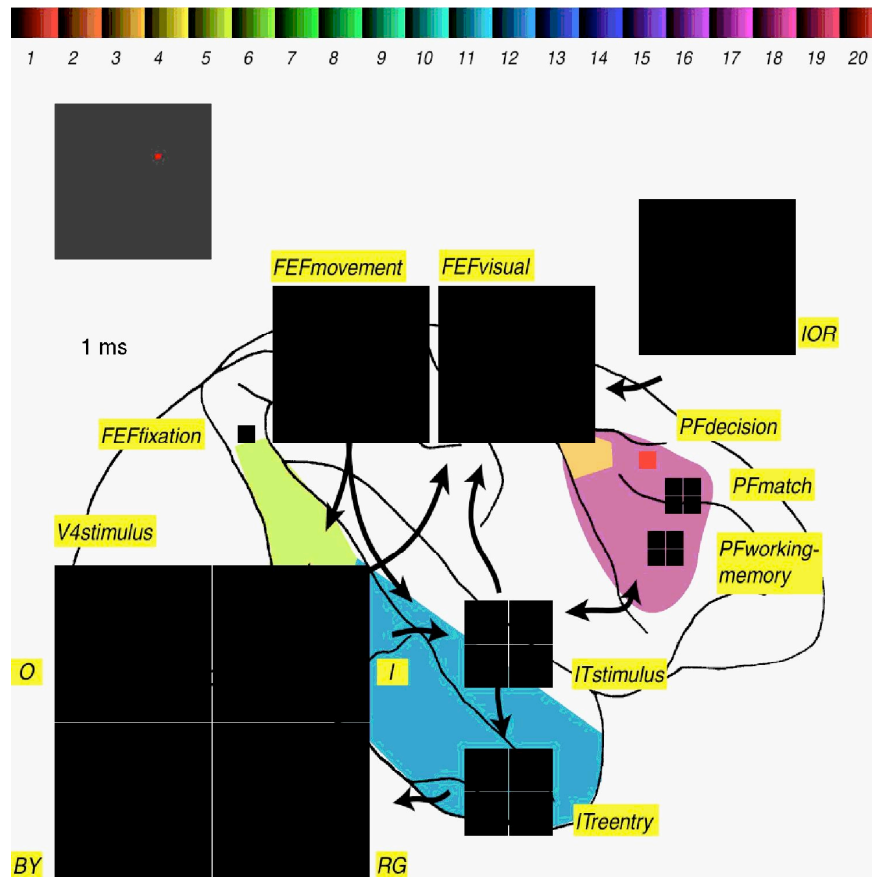


Problem of the Spotlight Metaphor

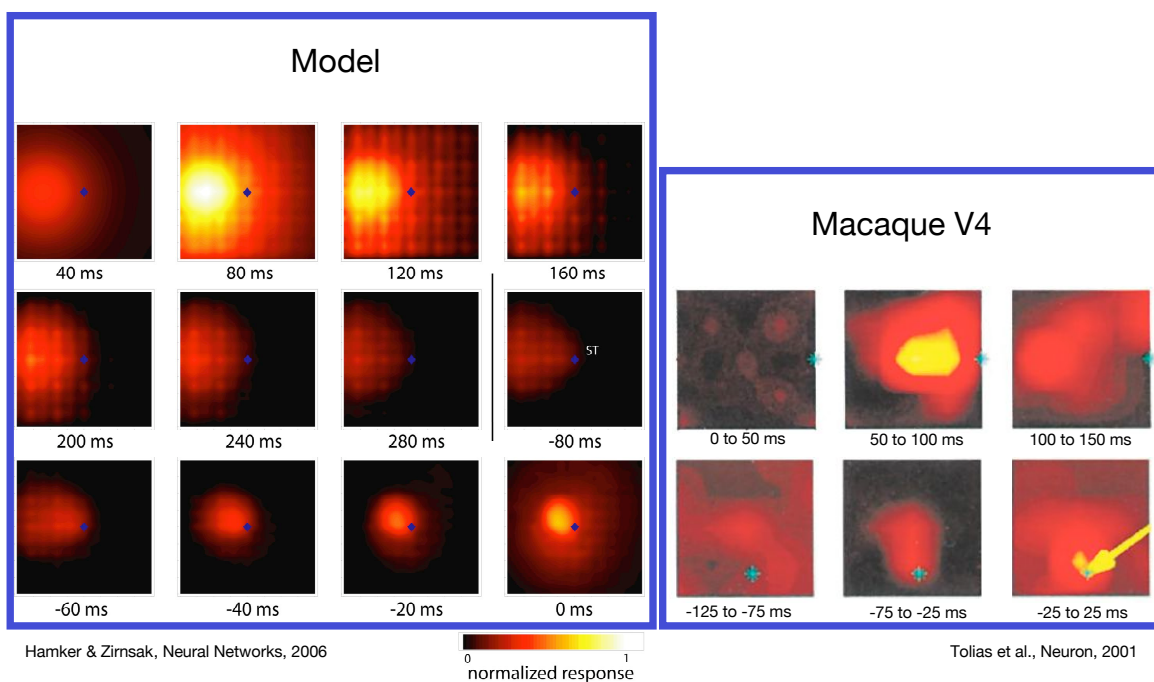


Attention tunes the RF properties

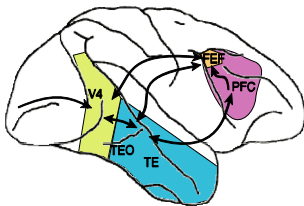
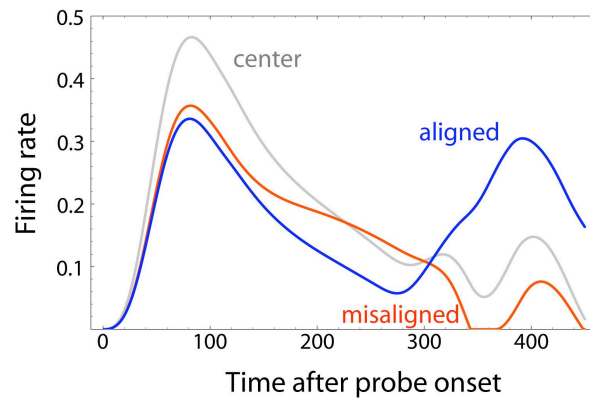
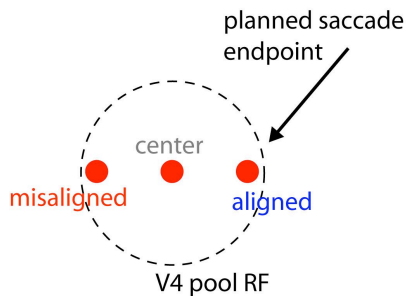




Model RF dynamics compared to V4



Effects of a saccade on the neural firing rate of model V4 cells



Hypotheses

Attention is a network property

It emerges since high level task descriptions have to be connected to low level scene descriptions

The planning of an eye movement provides a reentry signal which influences perception

Feature-specific feedback within the object recognition pathway, gain control and competitive interactions directly enhance the features of interest and guide spatial attention to the object of interest.

I propose that the direction of attention and recognition must be an iterative process to be effective.

Limitations of the present approach

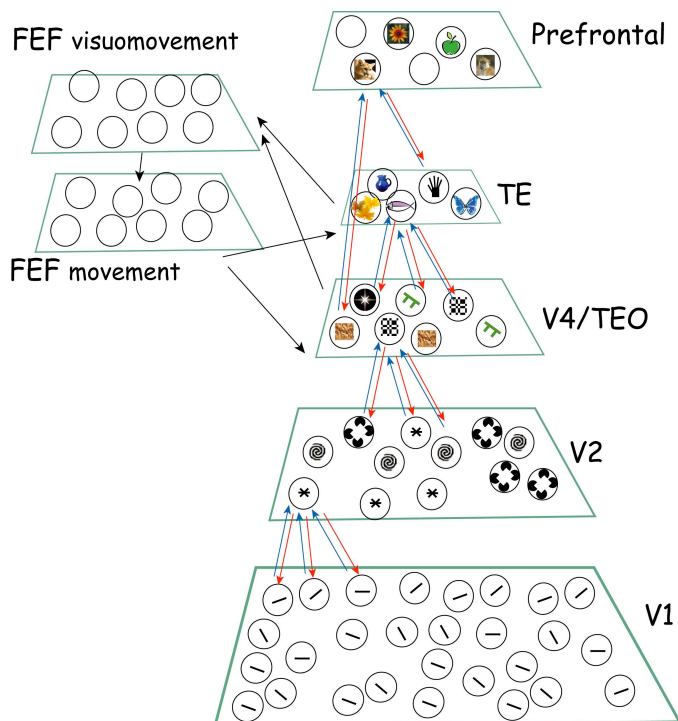
The representation on which object detection is made does hardly allow for real object recognition tasks and the guidance of vision can only be based on simple color and orientation cues.

Extend the present approach by learning feedback and feedforward transformations within the feature spaces of different complexity considering image statistics.

The model does not know what too look for.

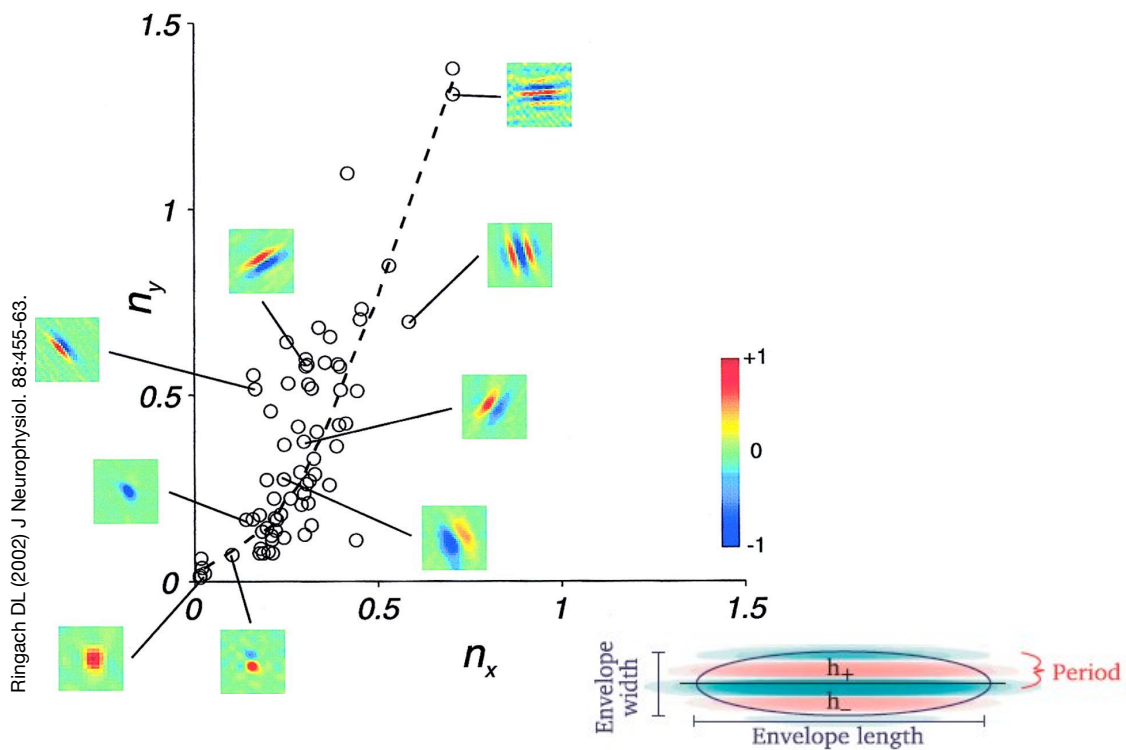
Develop a computational approach to learn the recall of target features on the task at hand in a reward based scenario for guiding visual perception.

Vision as an active, constructive process

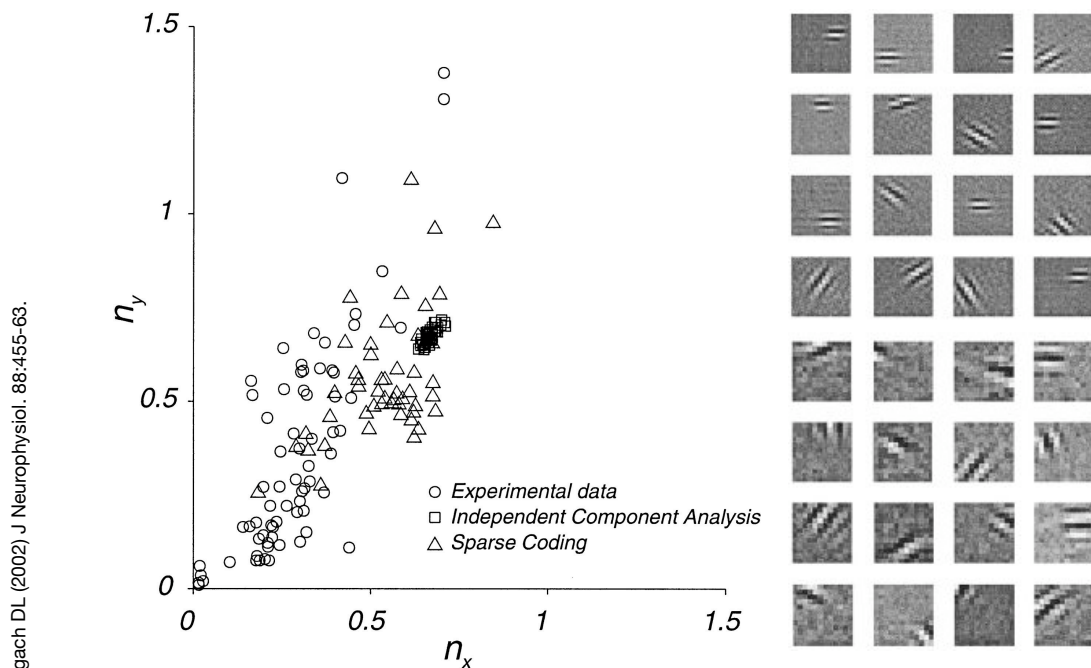


1. **Top-down guidance** – vision is guided by templates
2. **Enhancement of features of interest** – rather than only selecting just the location
3. **Parallel pattern matching** – switch into a serial search if parallel search does not discriminate the target
4. **Fast bottom-up recognition** – despite recurrent interactions in the system
5. **Flexible pattern matching** – matching process is flexible and indicates the similarity of target and object

Data of V1 receptive fields



Models of V1 receptive fields



The diagram illustrates a hierarchical visual processing model. At the bottom right, an input image of a rocky shore with a log is shown. A red square on the log indicates the receptive field of an ON/OFF cell. This cell's response, $r_i^{On/Off}$, is shown as two grayscale images: one for the ON response (positive) and one for the OFF response (negative). These responses are fed into the LGN (Lateral Geniculate Nucleus), represented by a horizontal row of cells. The LGN output, r_i^I , is then processed by a V1 (Primary Visual Cortex) cell, represented by a vertical column of cells. The V1 cell's response, r_j^{II} , is shown as a grayscale image of the receptive field. The diagram also includes labels for weights w_{ij} and a_{ji} connecting the LGN and V1 stages.

The diagram illustrates a V1 model. At the bottom, the LGN (Lateral Geniculate Nucleus) is shown with a horizontal row of neurons. Some are labeled 'ON' and others 'OFF'. Two specific receptive fields are highlighted: an 'ON' field and an 'OFF' field. The 'ON' field is represented by a dark, noisy image, and the 'OFF' field is represented by a light, noisy image. Arrows point from these fields to the corresponding neurons in the LGN. The output of the LGN neurons is represented by a vertical column of neurons in the V1 layer. The activity of these V1 neurons is labeled r_i^I and r_j^I . The weight connecting the LGN neuron i to the V1 neuron j is labeled w_{ij} . The activity of the V1 neuron j is also labeled r_j^I . The diagram also shows the 'Anti-Hebb' learning rule: $\Delta c_{ij} = -\eta \cdot r_i^I \cdot r_j^I$. The receptive fields (RF 1 and RF 2) are shown as two sets of images. The 'Diff' (difference) between the two RFs is shown as a dark, noisy image. The 'Anti-Hebb' result is shown as a light, noisy image.

Results of learning in model V1

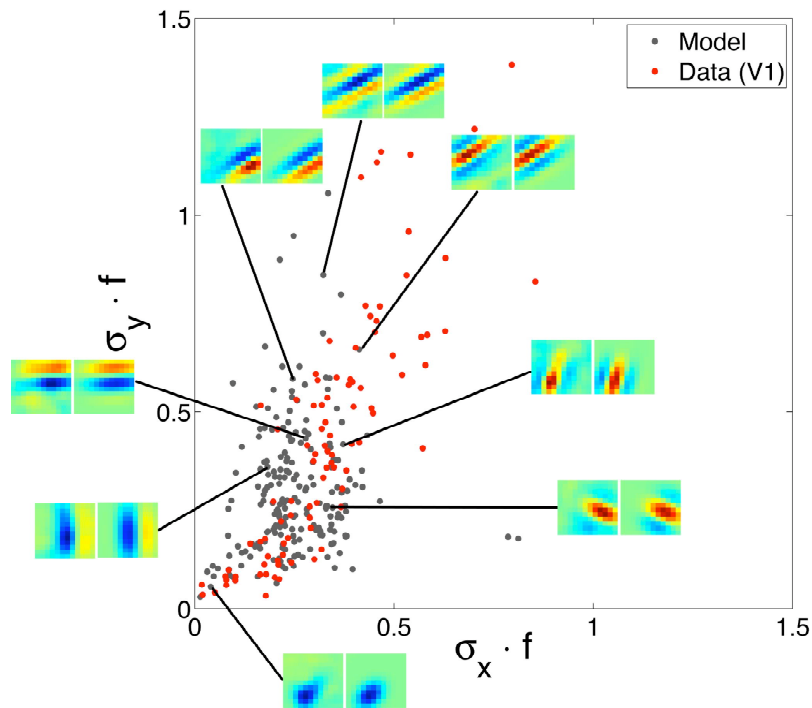


Image reconstruction

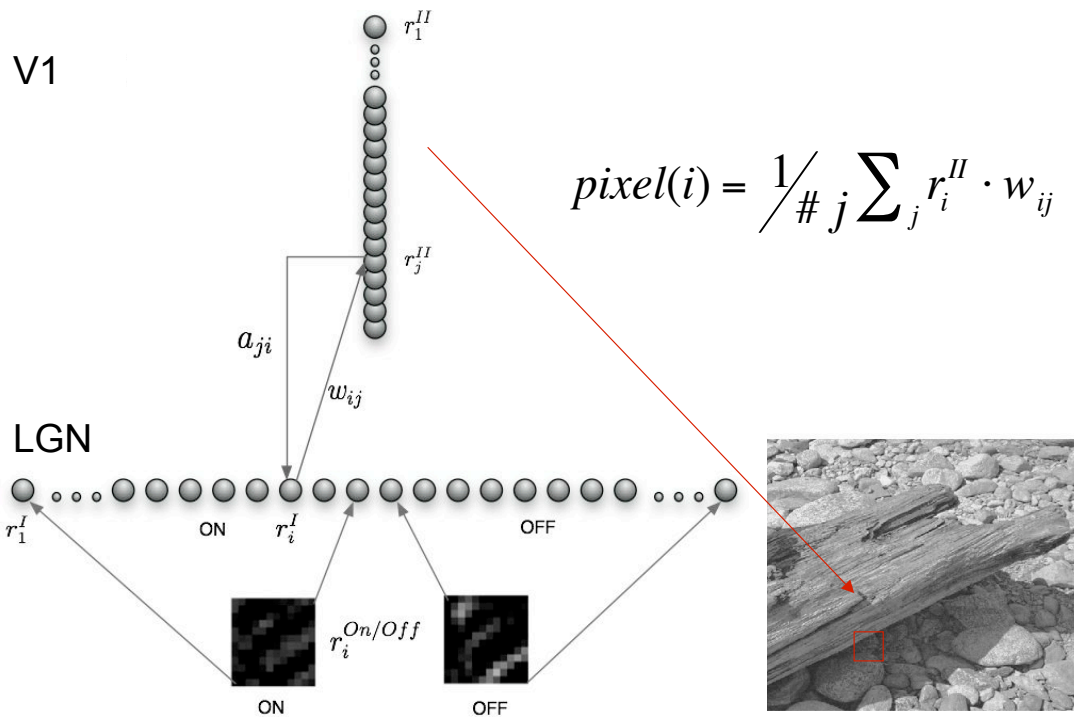
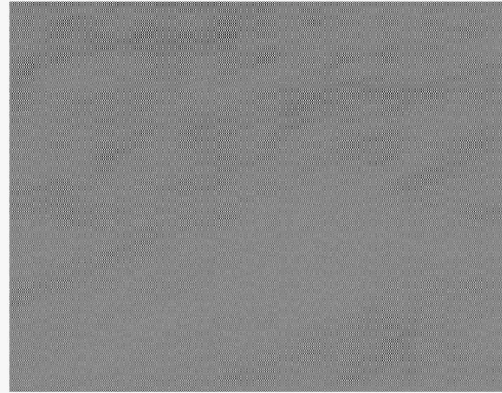


Image reconstruction



Original



Reconstruction

Challenges and planned work in Eyeshots

- Learning of joint feature and disparity information in model V1 receptive fields, by training the model with stereo images (requires images sequences).
- Expand approach to the next higher level to learn more complex features (including disparity).
- Implement attentional dynamics within this network.
- Learn when to look for a particular object through reward.

Deliverables in Eyeshots

Deliverables

D3.1a: Demonstration of learning disparity-tuned feature selective cells. Software module. (Month 12).

D3.1b: Demonstration of object selective cells at intermediate complexity showing properties of disparity. Technical report. (Month 24)

D3.2: Object-based top-down selection using learned bi-directional connections between feature detectors to localize the object of interest in a cluttered 3D scene. Software module. (Month 36)

D3.3a: A model of working memory that allows to activate context information for the task at hand based on the association of previous events leading to reward. We will use the learned feature responses (WP2) on real world scenes if the feature-detectors are available, otherwise artificial representations will be used. Software module. (Month 24).

D3.3b: Final, fully tested version of the Working Memory Model. Technical report. (Month 36).