

Modelli computazionali di attenzione visiva



Corso di Interazione uomo-macchina II

Prof. Giuseppe Boccignone

Dipartimento di Informatica
Università di Milano

boccignone@di.unimi.it
http://boccignone.di.unimi.it/IUM2_2014.html

Interazione fra organismi //tecnologie

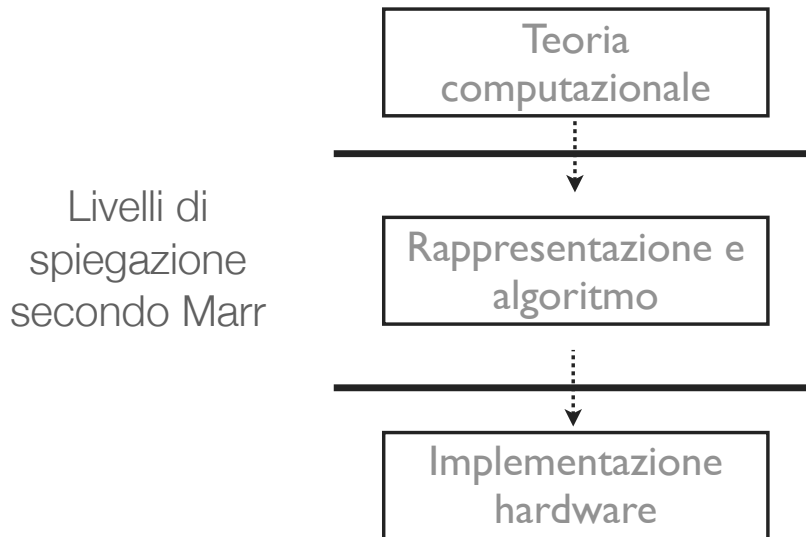
	Example Social Behaviours							Tech.		
	emotion	personality	status	dominance	persuasion	regulation	rapport	speech analysis	computer vision	biometry
Social Cues										
Physical appearance										
height			✓	✓					✓	✓
attractiveness		✓	✓	✓	✓		✓		✓	✓
body shape		✓		✓					✓	✓
Gesture and posture										
hand gestures	✓	✓			✓	✓	✓		✓	✓
posture	✓	✓	✓	✓	✓	✓	✓		✓	✓
walking		✓	✓	✓					✓	✓
Face and eyes behaviour										
facial expressions	✓	✓	✓	✓	✓	✓	✓		✓	✓
gaze behaviour	✓	✓	✓	✓	✓	✓	✓		✓	
focus of attention	✓	✓	✓	✓	✓	✓	✓		✓	
Vocal behaviour										
prosody	✓	✓		✓	✓		✓	✓		
turn taking	✓	✓	✓	✓		✓	✓	✓		
vocal outbursts	✓	✓		✓	✓	✓	✓	✓		
silence	✓		✓				✓	✓		
Space and Environment										
distance	✓	✓	✓		✓		✓		✓	
seating arrangement				✓	✓		✓		✓	



A. Vinciarelli, M. Pantic, H. Bourlard,
*Social Signal Processing: Survey of an
Emerging Domain,*
Image and Vision Computing (2008)

Modelli nelle scienze cognitive e nella percezione

//livelli di spiegazione



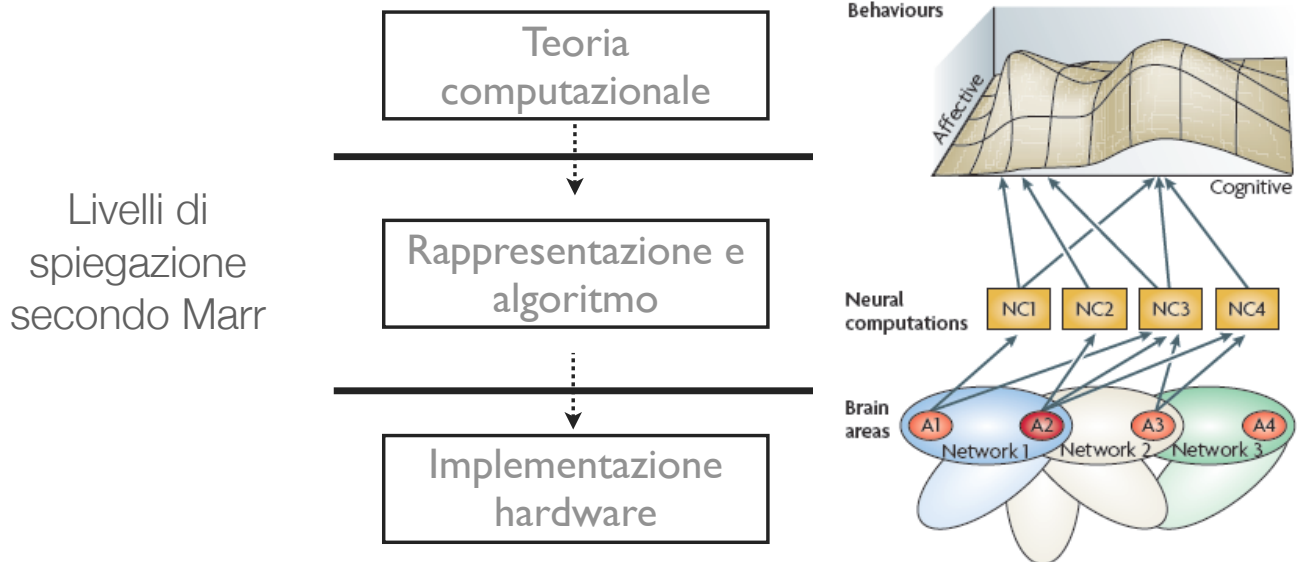
Modelli nelle scienze cognitive e nella percezione

//livelli di spiegazione



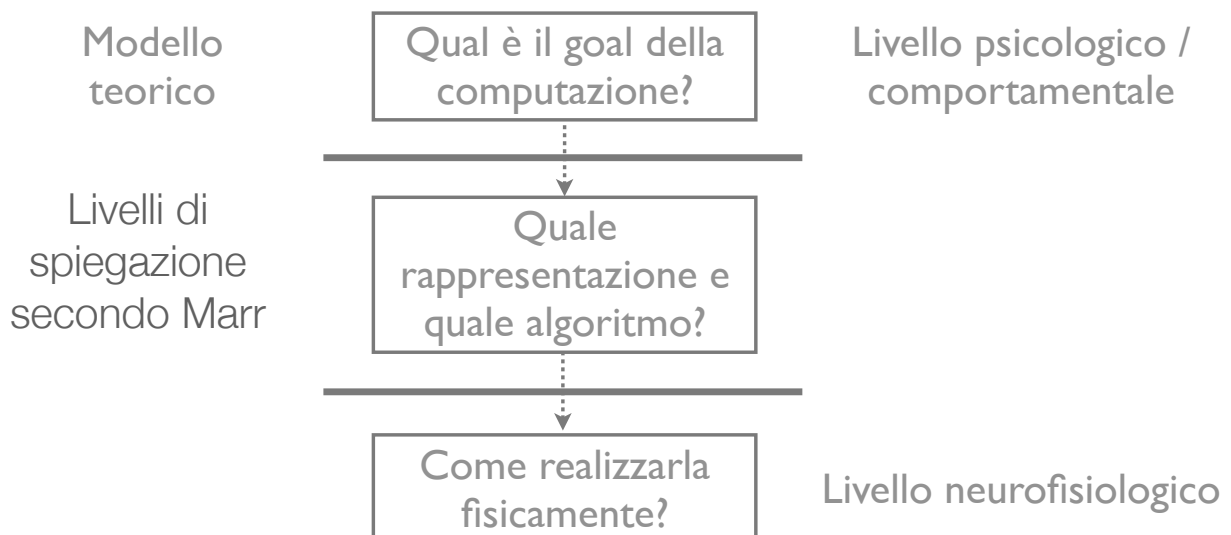
Modelli nelle scienze cognitive e nella percezione

//livelli di spiegazione



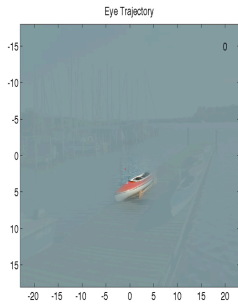
Modelli nelle scienze cognitive e nella percezione

//livelli di spiegazione



Modelli di attenzione visiva

//livelli di spiegazione



Qual è il goal della computazione?

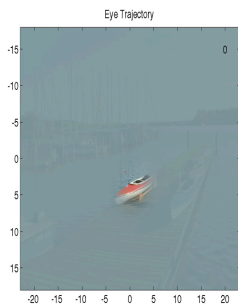
$$\mathbf{I} \mapsto \{\mathbf{r}_F(1), \mathbf{r}_F(2), \dots\}$$

Quale rappresentazione e quale algoritmo?

Come realizzarla fisicamente?

Modelli di attenzione visiva

//livelli di spiegazione



Qual è il goal della computazione?

che cosa guardo

$$\mathbf{I} \mapsto \mathcal{R}$$

$$\mathcal{R} \mapsto \{\mathbf{r}_F(1), \mathbf{r}_F(2), \dots\}$$

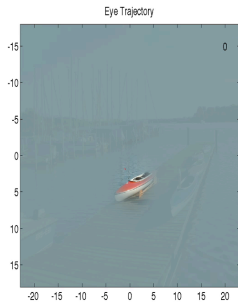
come guardo

Quale rappresentazione e quale algoritmo?

Come realizzarla fisicamente?

Modelli di attenzione visiva

//livelli di spiegazione



che cosa guardo
 $I \mapsto \mathcal{R}$

Qual è il goal della computazione?

$\mathcal{R} \mapsto \{r_F(1), r_F(2), \dots\}$
 come guardo

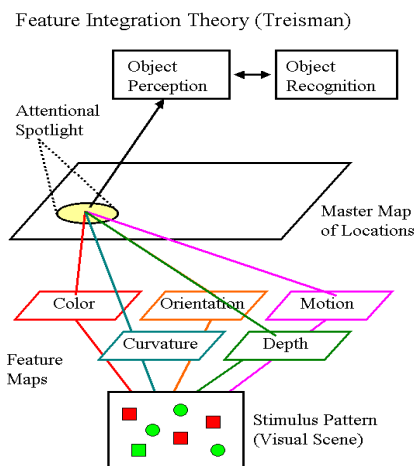
Quale rappresentazione e quale algoritmo?

Come realizzarla fisicamente?

Un semplice modello computazionale

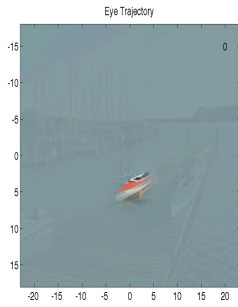
//Itti e Koch

- Nella sua formulazione originale è un modello bottom-up:
 - ha alla base il concetto di salienza degli stimoli fisici
- Basato sul modello psicologico della Treisman (FIT)



Modelli di attenzione visiva

//livelli di spiegazione



Qual è il goal della
computazione?

Quale
rappresentazione e
quale algoritmo?

Come realizzarla
fisicamente?

guardo i punti salienti

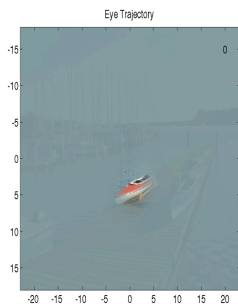
$$\mathbf{I} \mapsto \mathcal{R}$$

$$\arg \max \mathcal{R}$$

scelgo il più saliente

Modelli di attenzione visiva

//livelli di spiegazione



Qual è il goal della
computazione?

Quale
rappresentazione e
quale algoritmo?

Come realizzarla
fisicamente?

guardo i punti salienti

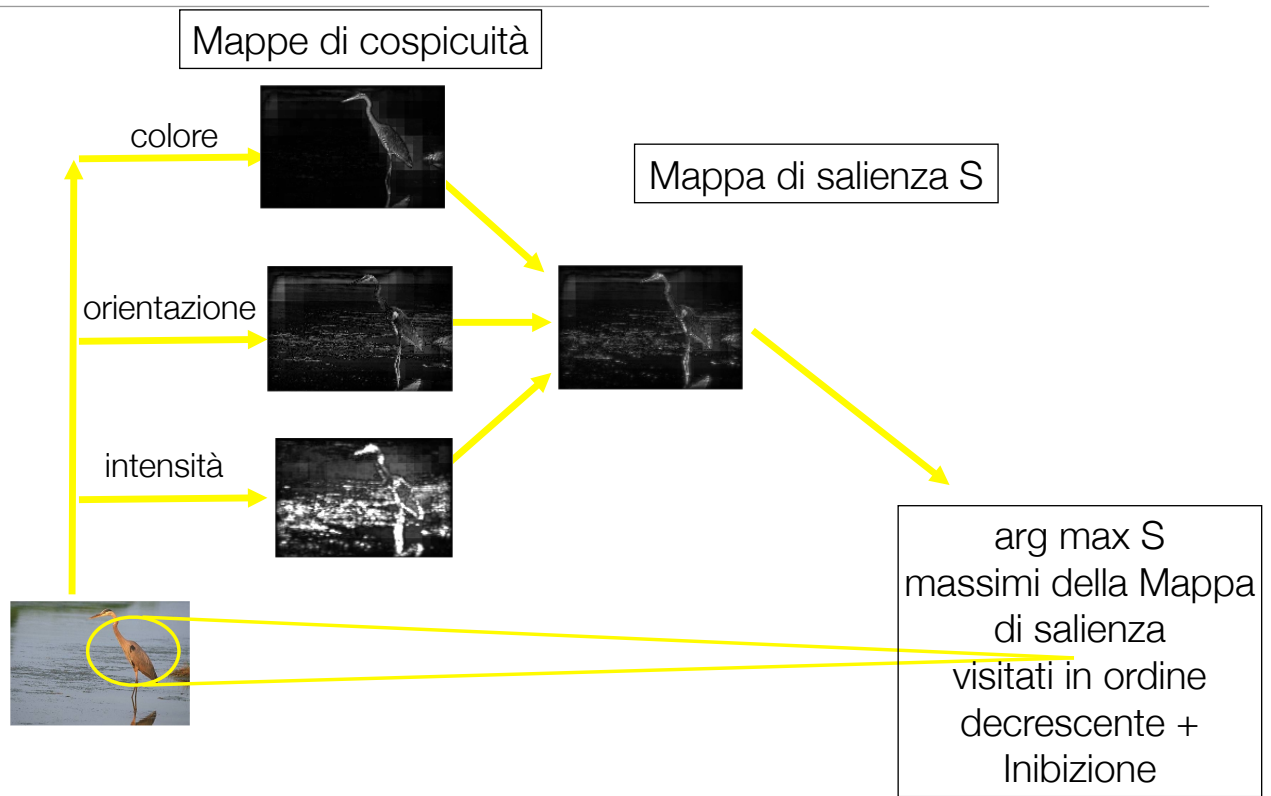
$$\mathbf{I} \mapsto \mathcal{R}$$

$$\arg \max \mathcal{R}$$

scelgo il più saliente

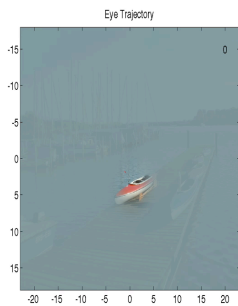
Un semplice modello computazionale

//Itti e Koch: algoritmo



Modelli di attenzione visiva

//livelli di spiegazione



Qual è il goal della
computazione?

guardo i punti salienti

$$\mathbf{I} \mapsto \mathcal{R}$$

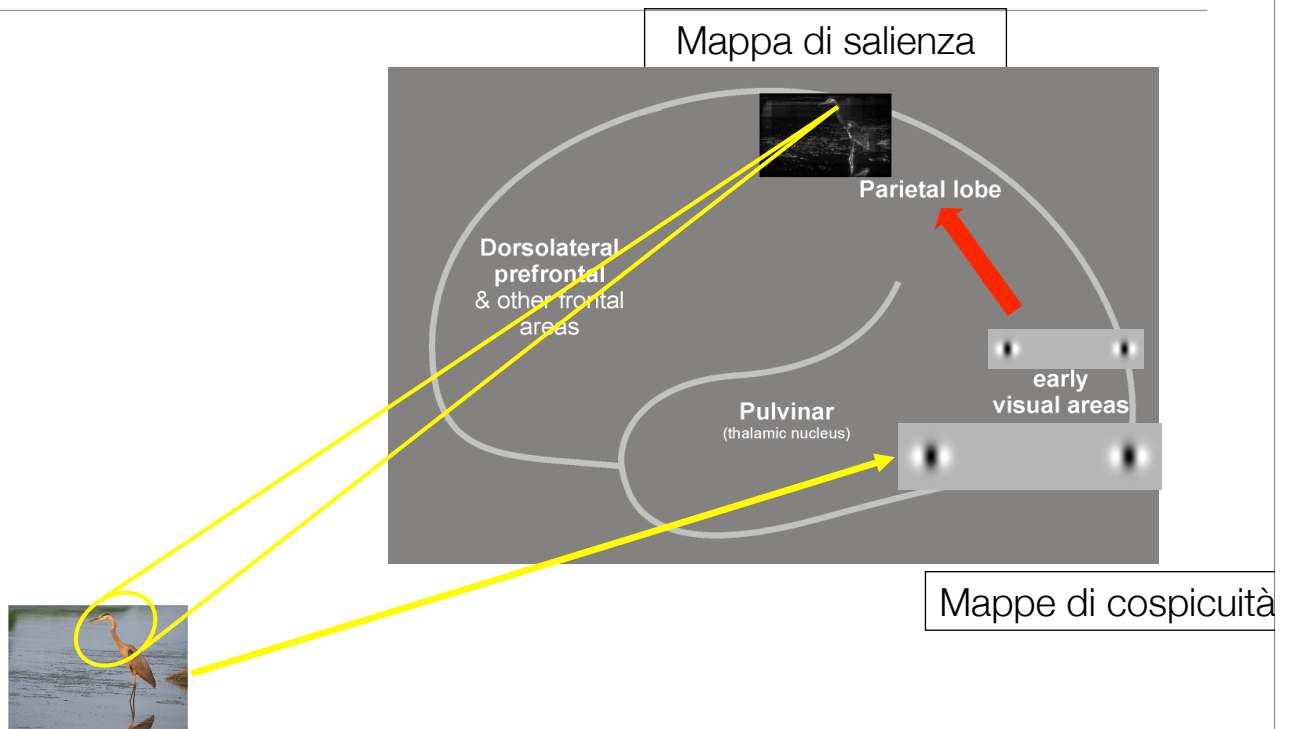
$$\arg \max \mathcal{R}$$

scelgo il più saliente

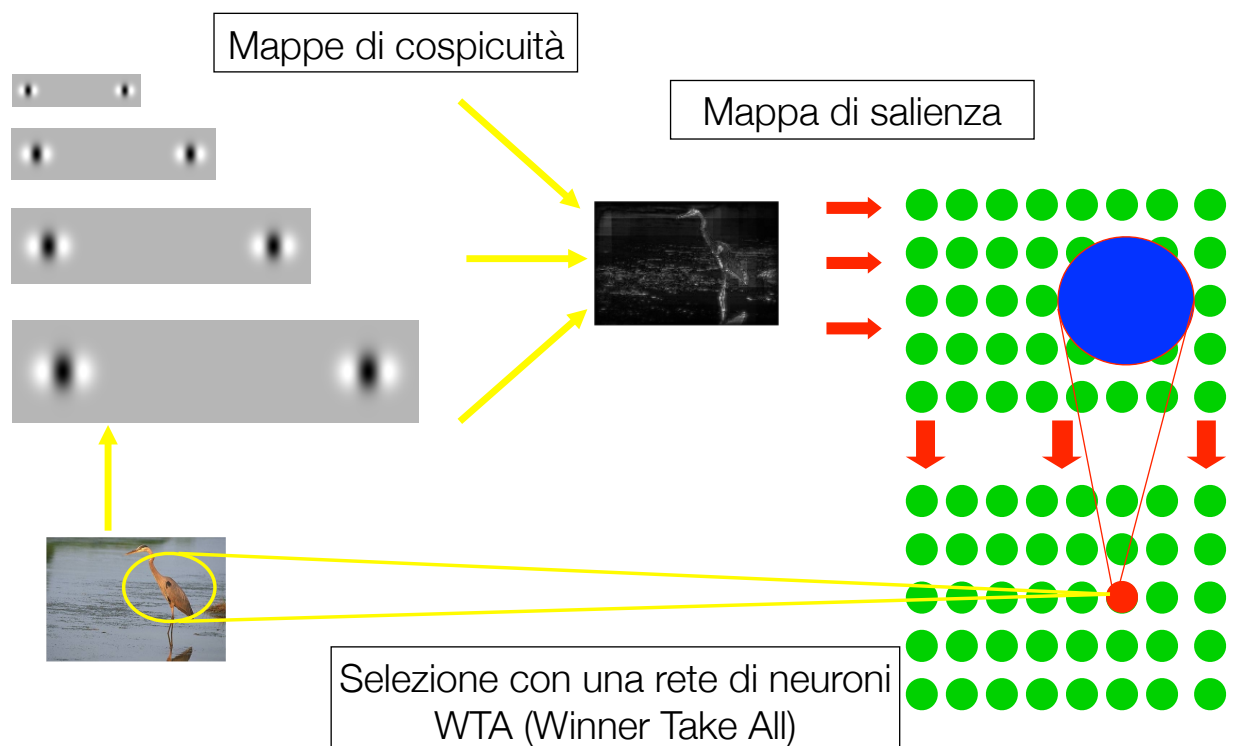
Quale
rappresentazione e
quale algoritmo?

**Come realizzarla
fisicamente?**

Un semplice modello computazionale //Itti e Koch: implementazione neurale

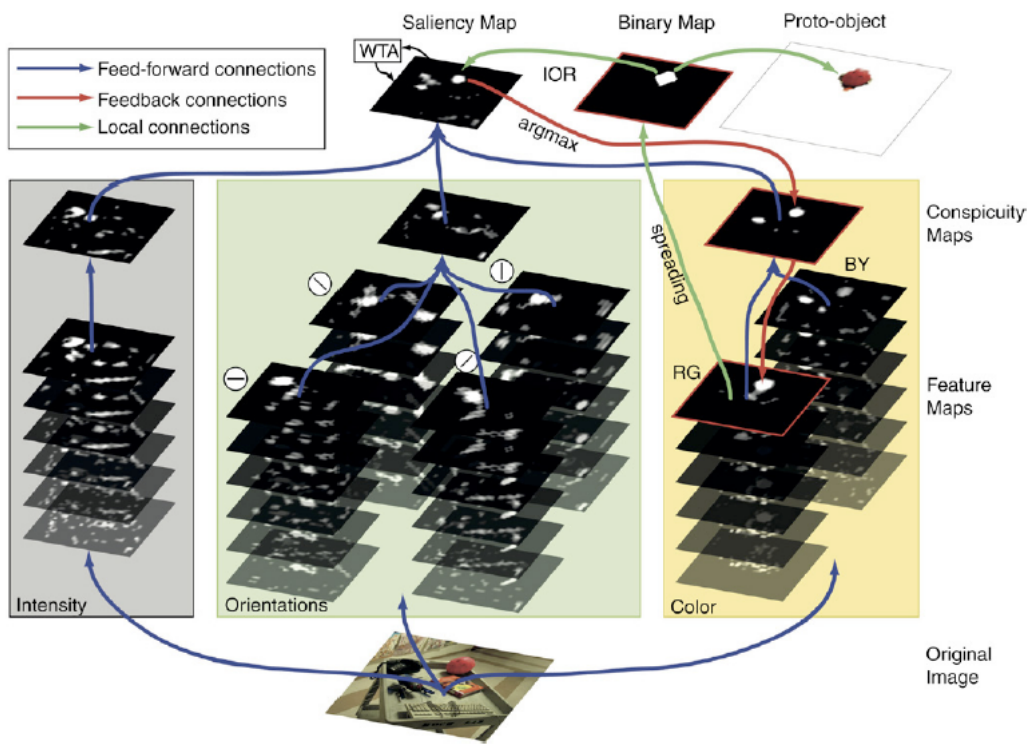


Un semplice modello computazionale //Itti e Koch: implementazione neurale

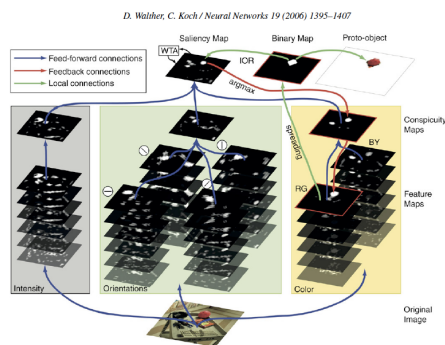


Un semplice modello computazionale //Itti e Koch

D. Walther, C. Koch / Neural Networks 19 (2006) 1395–1407



Un semplice modello computazionale //Itti e Koch



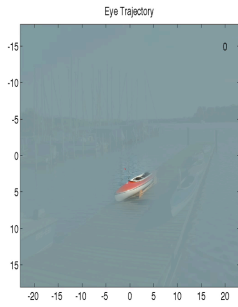
Free Matlab Code:

 toolbox ufficiale: <http://www.saliencytoolbox.net>

 diverse versioni: <http://www.klab.caltech.edu/~harel/share/gbvs.php>

Modelli di attenzione visiva

//livelli di spiegazione



Qual è il goal della
computazione?

che cosa guardo

$$\mathbf{I} \mapsto \mathcal{R}$$

$$\mathcal{R} \mapsto \{r_F(1), r_F(2), \dots\}$$

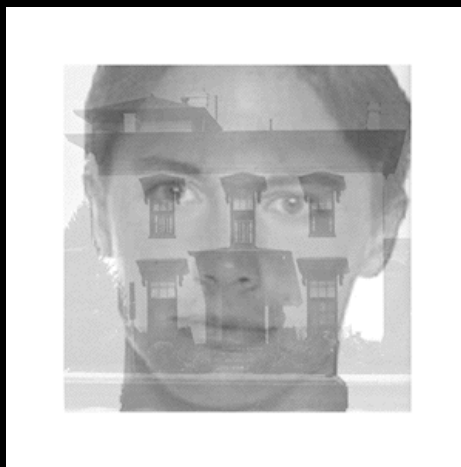
come guardo

Quale
rappresentazione e
quale algoritmo?

Come realizzarla
fisicamente?

What gets selected?

Object Based Attention



- O'Craven et al. (1999)

What gets selected?

Object Based Attention

- Summing up:
 - Importance of 'object-based' processes in dynamic tasks
 - Center-of-mass capture of attention
 - Physical constraints: objects vs. substances
 - Dynamic tuning of the Focus of Attention (FOA)

Towards event-based attention



People walking together



t_i

t_j

People fighting



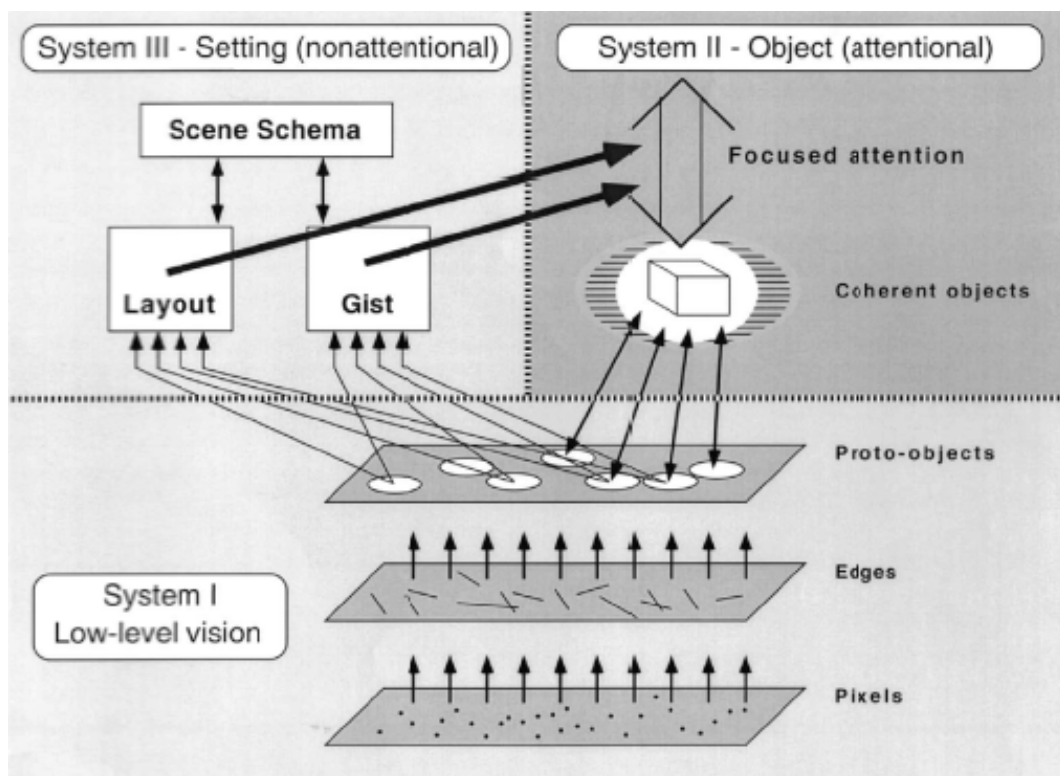
t_k

t_m

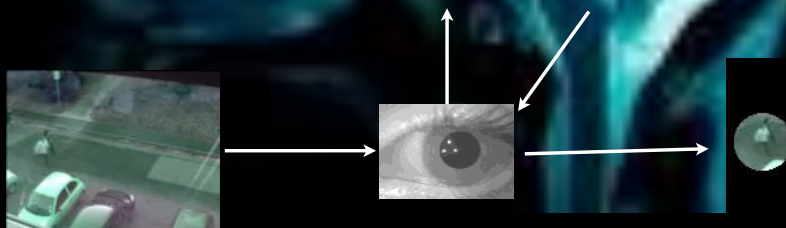
Towards event-based attention

- Nearly all of the work on attention concerns static objects, or moving objects
- Attention may also interact directly with information which is inherently dynamic, e.g. stereotypical motions

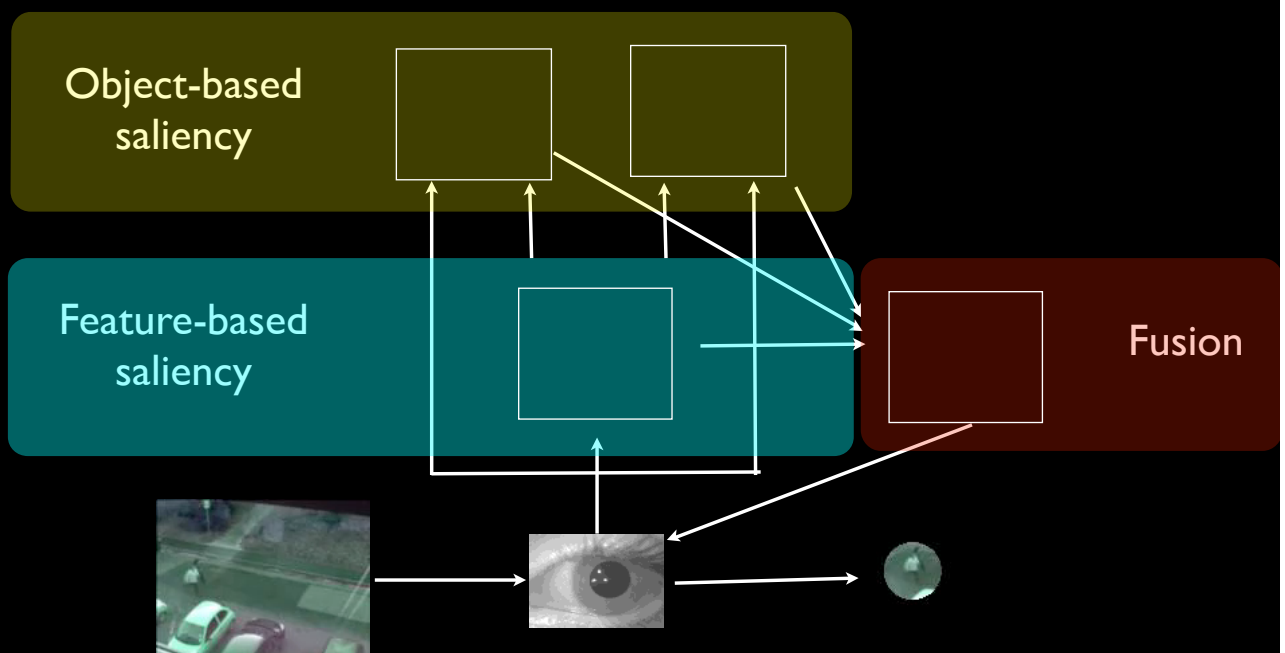
Livello di spiegazione psicologico
//rappresentazione dinamica di scene (Rensink)



Modelling visual attention

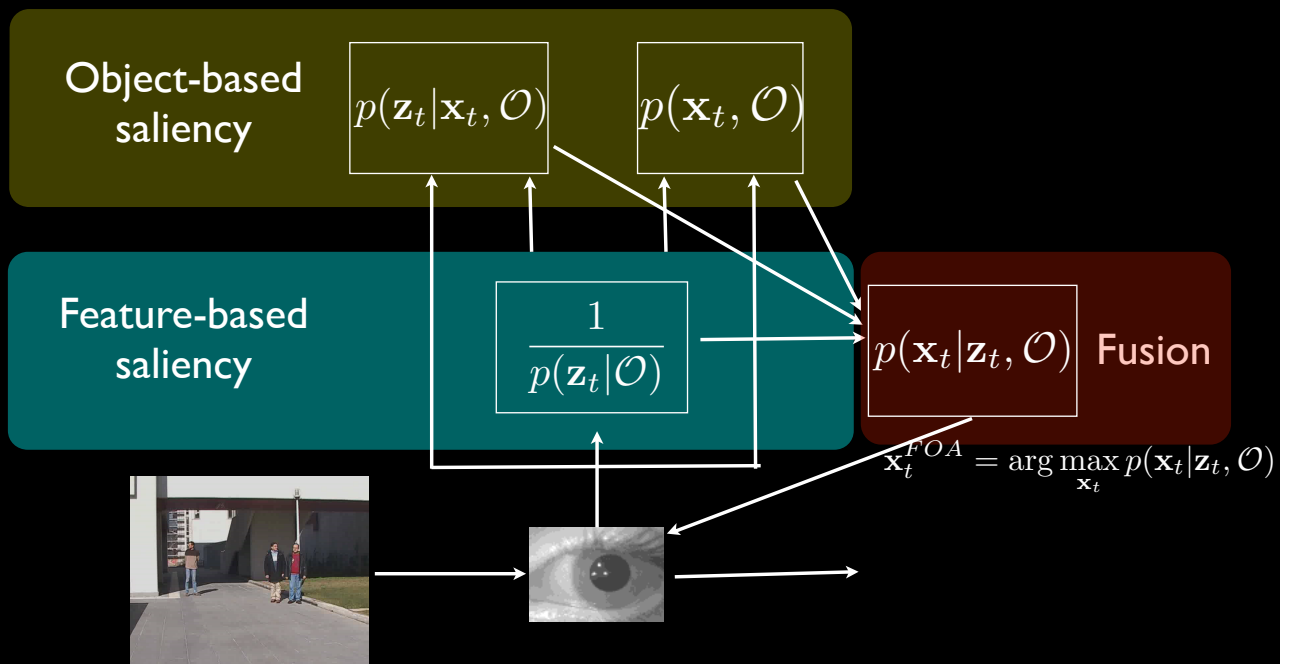


Modelling visual attention

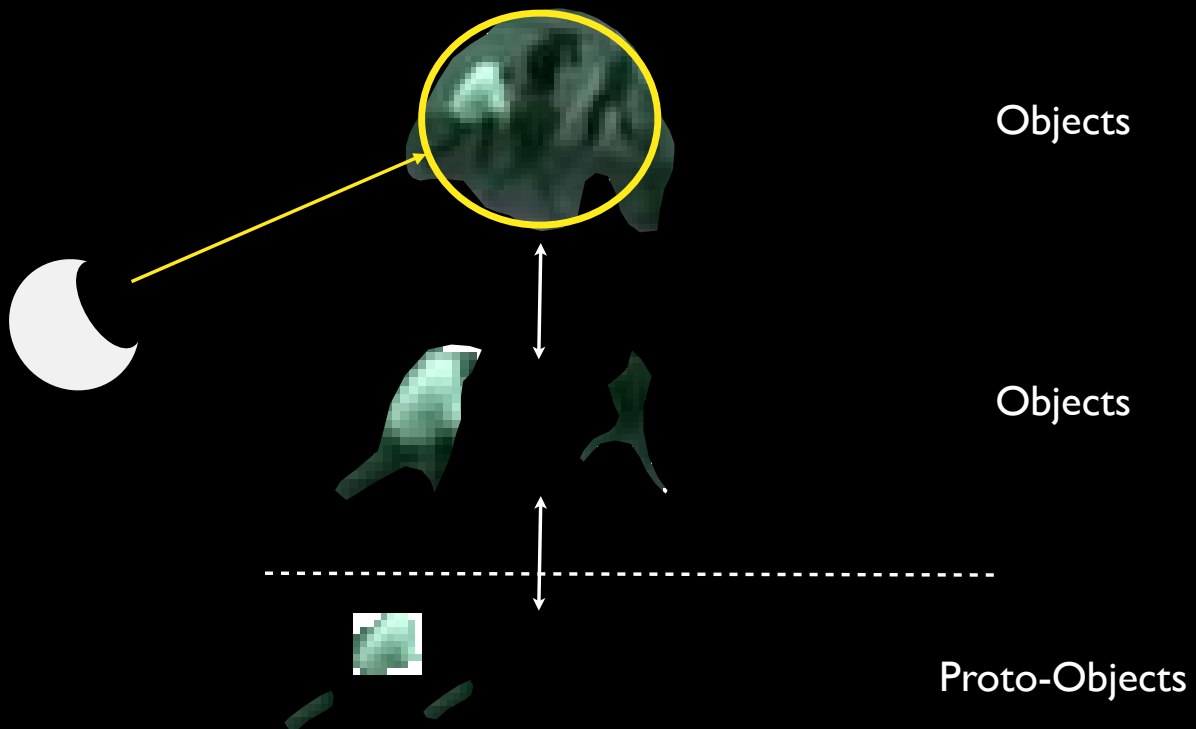


Modelling visual attention:

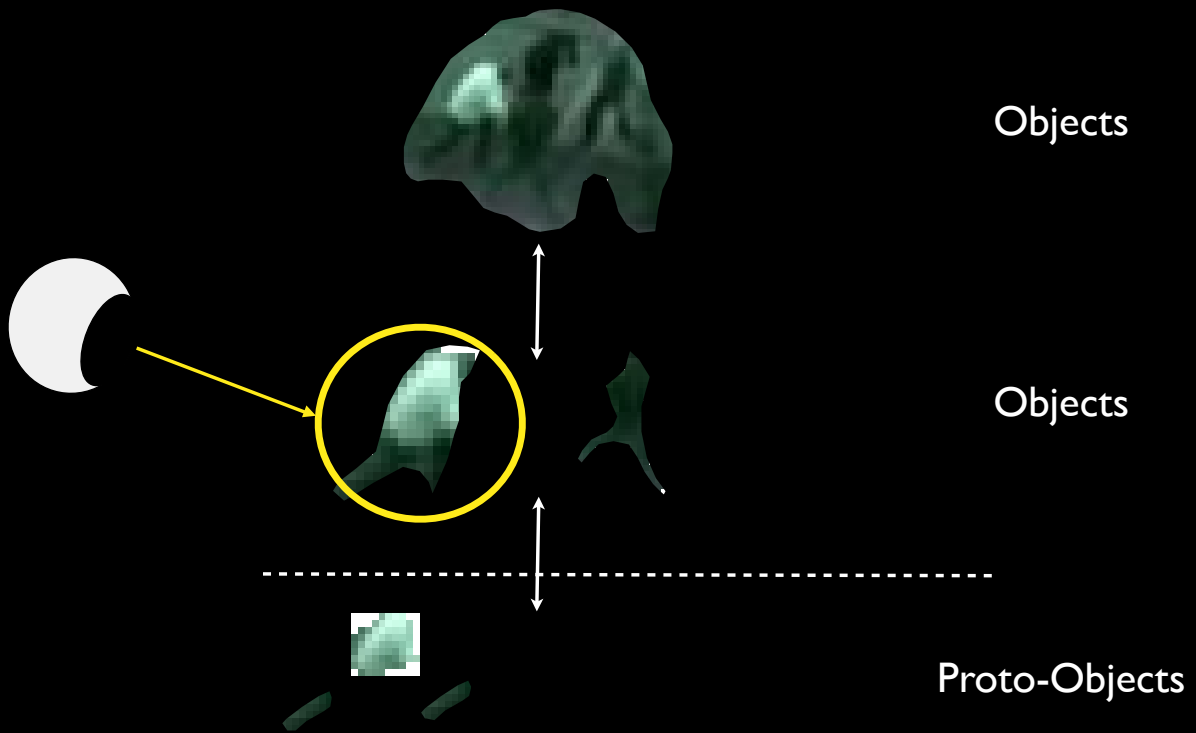
Vision as an inference process



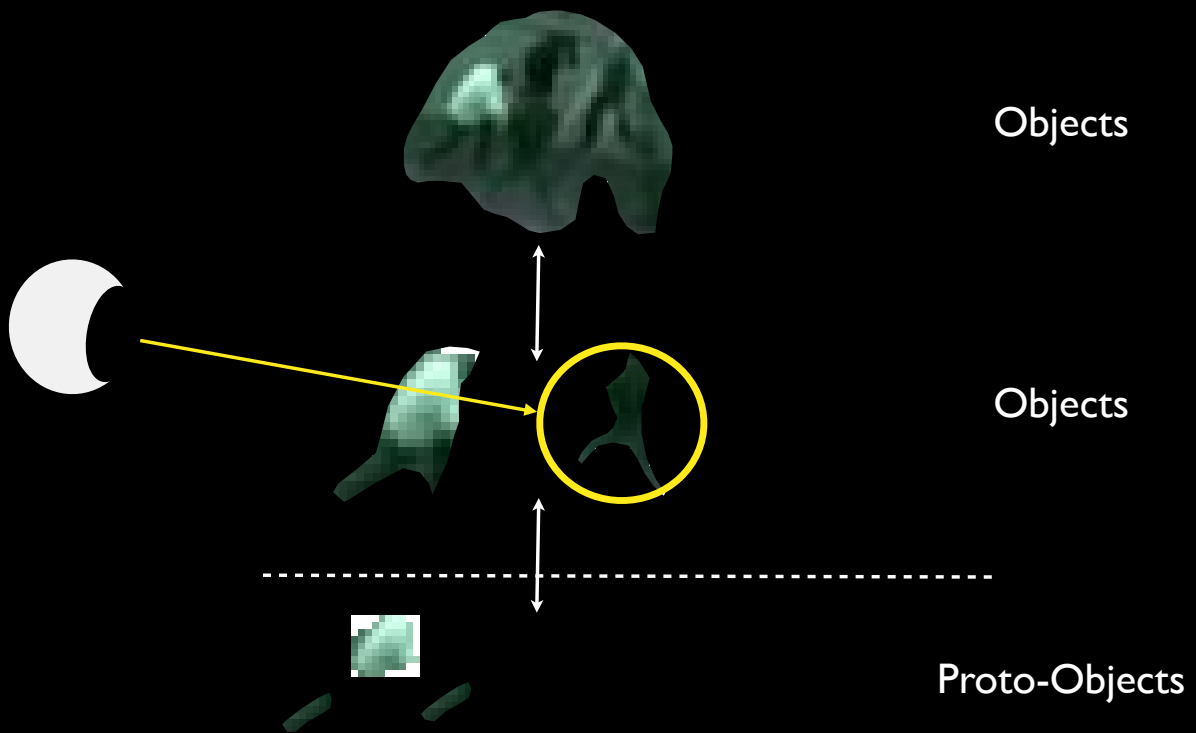
Towards event-based attention



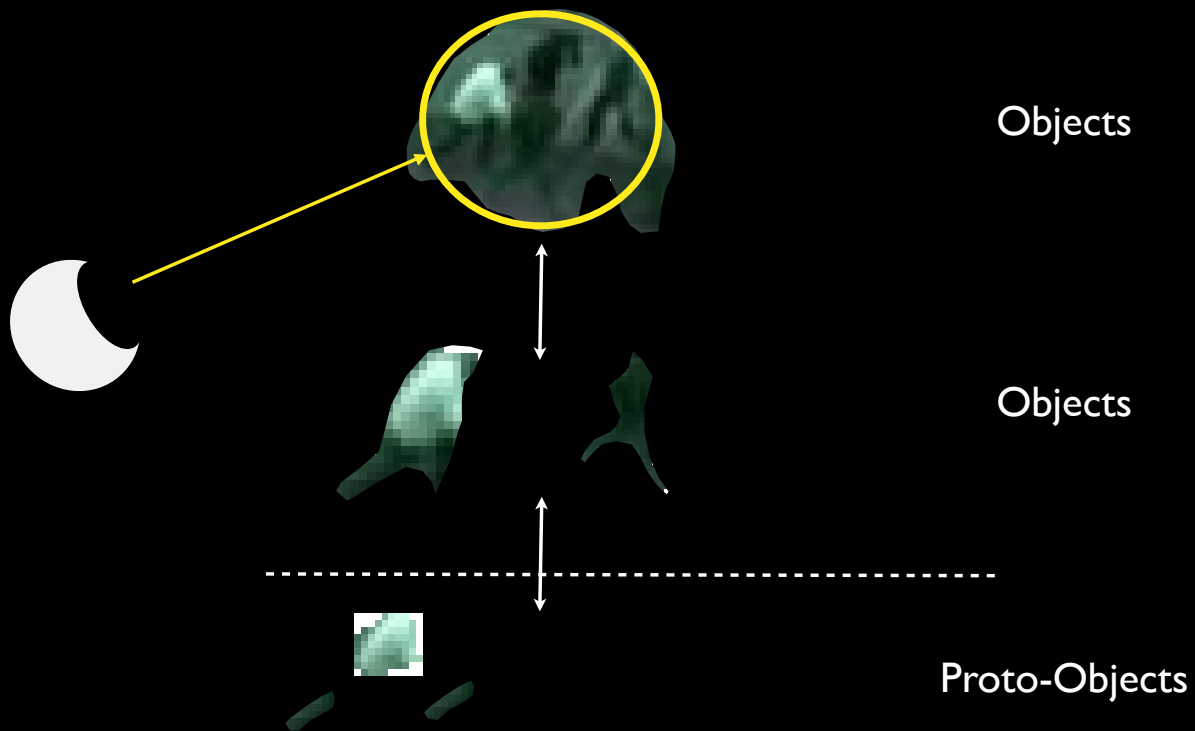
Towards event-based attention



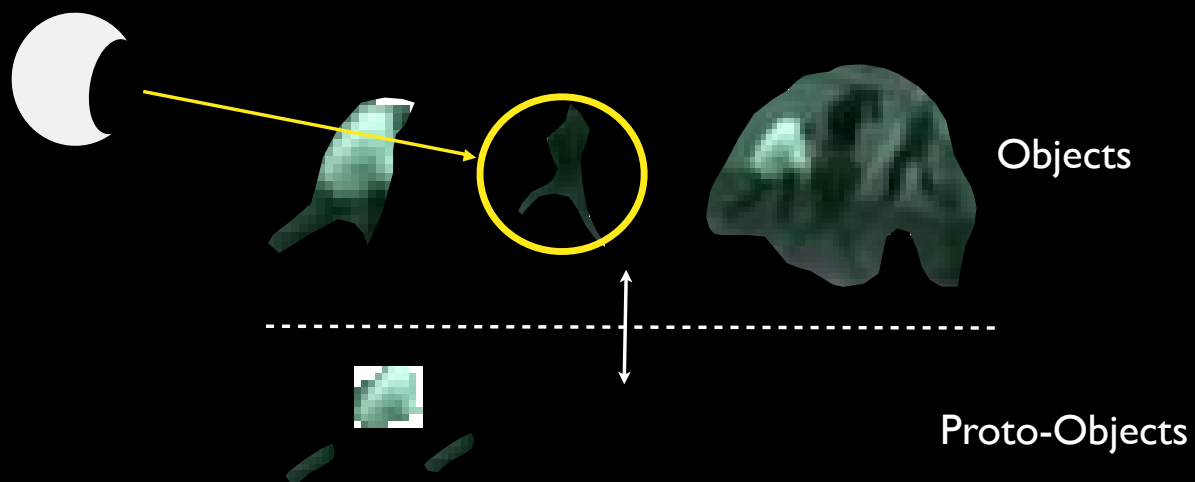
Towards event-based attention



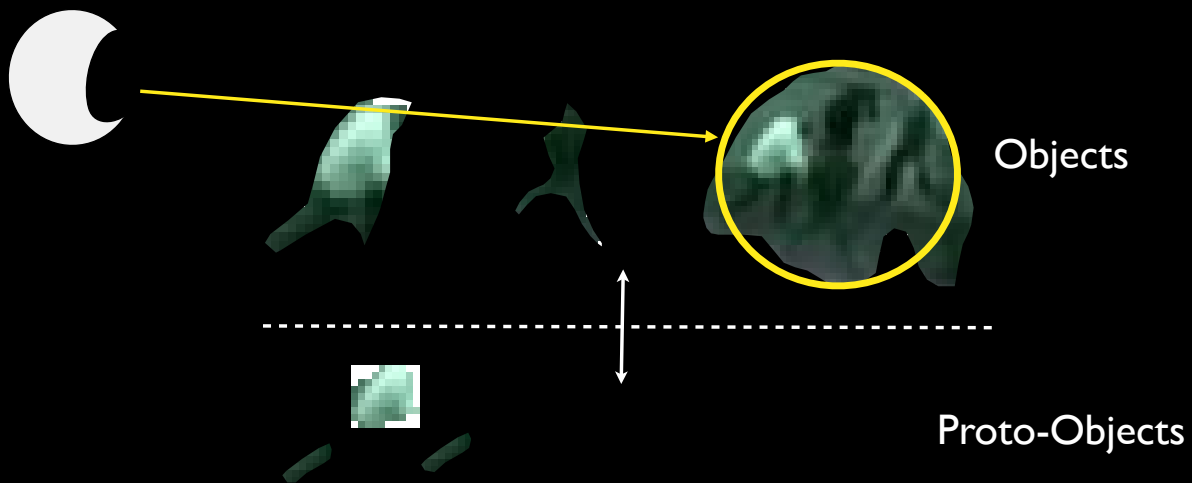
Towards event-based attention



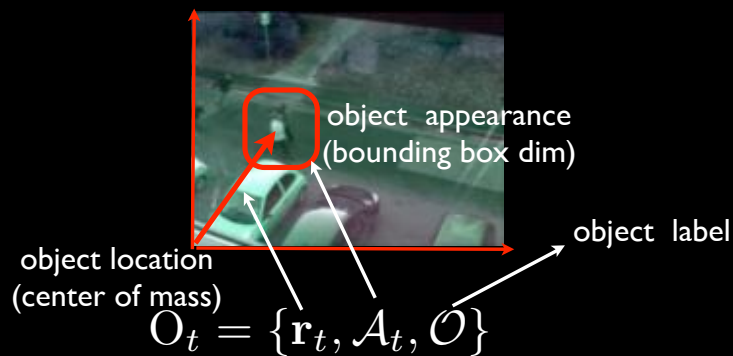
Towards event-based attention: may not be hierarchical



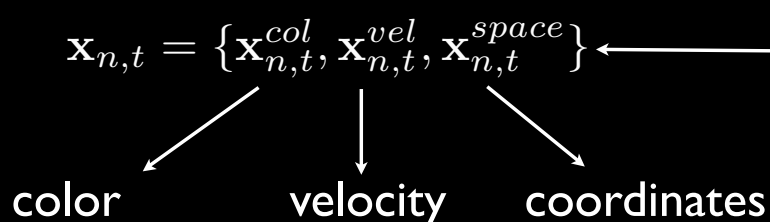
Towards event-based attention: may not be hierarchical



Towards event-based attention: a Bayesian model



Observed (spatio-temporal) features

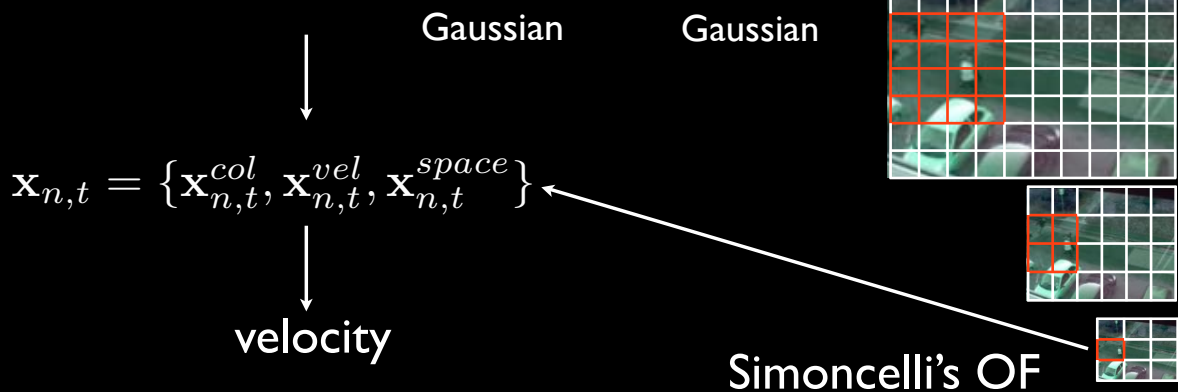


Towards event-based attention: a Bayesian model

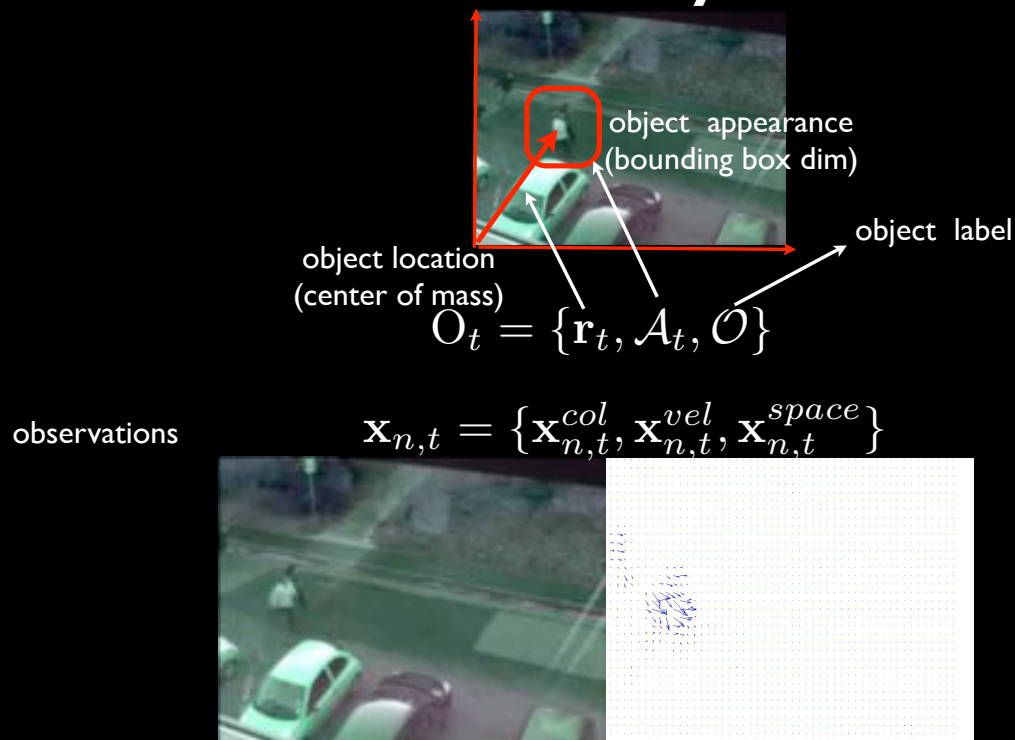
Observed (coarse) velocity features

$$\arg \max p(\mathbf{x}_t^{vel} | \frac{\partial x_t^L}{\partial \mathbf{r}}, \frac{\partial x_t^L}{\partial t}) \propto$$

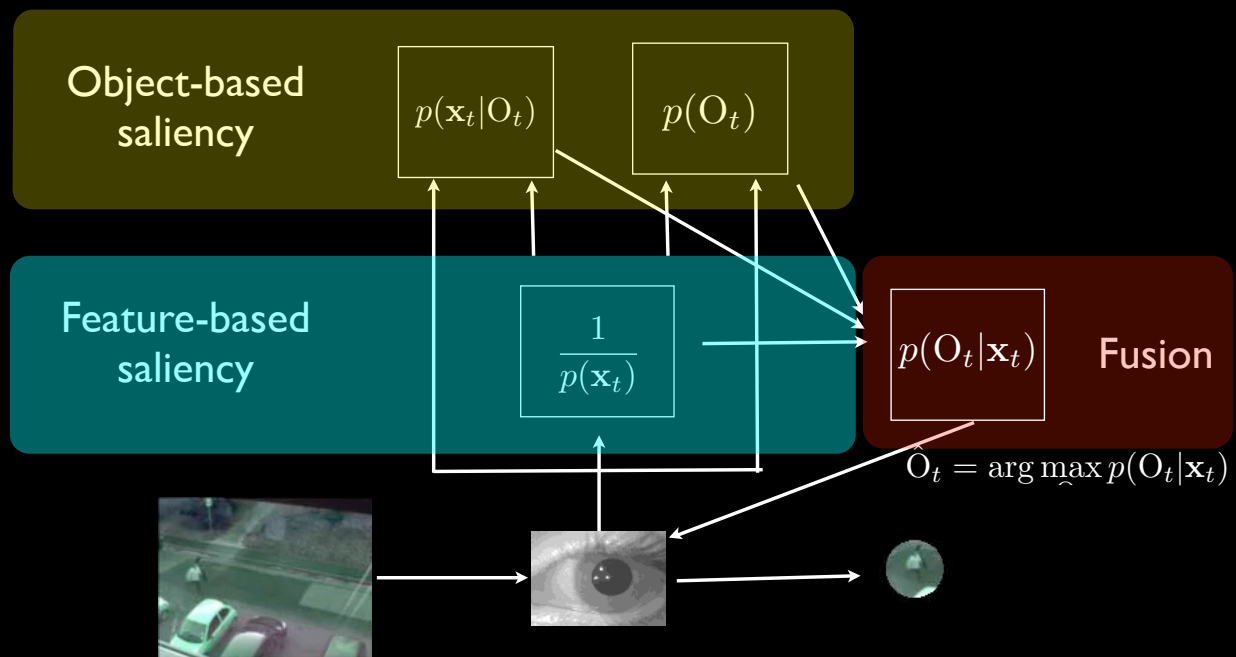
$$\arg \max p(\frac{\partial x_t^L}{\partial t} | \mathbf{x}_t^{vel}, \frac{\partial x_t^L}{\partial \mathbf{r}}) p(\mathbf{x}_t^{vel})$$



Towards event-based attention: a Bayesian model



Eye-movements and visual attention



Towards event-based attention: a Bayesian model

Objects/Groups of Objects features (e.g., velocity) Focus based on appearance & objects

$$p(O_t | \mathbf{x}_t) \simeq \frac{p(\mathbf{x}_t | \mathbf{r}_t, \mathcal{A}_t, \mathcal{O}) p(\mathbf{r}_t, | \mathcal{A}_t, \mathcal{O})}{p(\mathbf{x}_t)}$$

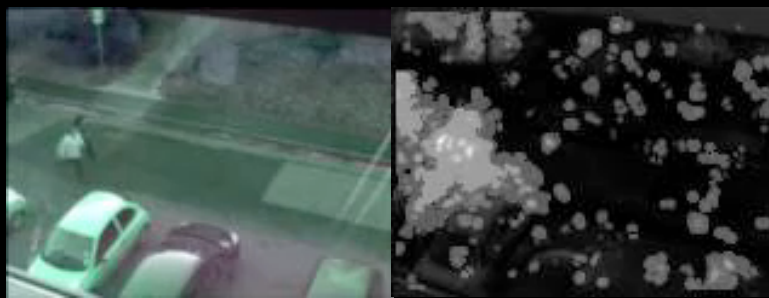
Low-level saliency & proto-object segmentation

A Bayesian model: proto-objects & saliency

$$p(O_t | \mathbf{x}_t) \simeq \frac{p(\mathbf{x}_t | \mathbf{r}_t, \mathcal{A}_t, \mathcal{O}) p(\mathbf{r}_t | \mathcal{A}_t, \mathcal{O})}{p(\mathbf{x}_t)}$$

Expectation-Maximization \longleftrightarrow Low-level saliency & proto-object segmentation

A Bayesian model: saliency



Feature-based
saliency

$$\frac{1}{p(\mathbf{x}_t)}$$

A Bayesian model: feature likelihood given proto-objects

Objects/Groups of Objects
features (e.g., velocity)

$$p(O_t | \mathbf{x}_t) \simeq \frac{p(\mathbf{x}_t | \mathbf{r}_t, \mathcal{A}_t, \mathcal{O}) p(\mathbf{r}_t | \mathcal{A}_t, \mathcal{O})}{p(\mathbf{x}_t)}$$

$$p(\mathbf{x}_t | \mathbf{r}_t, \mathcal{A}_t, \mathcal{O}) = \frac{\exp\{(\mathbf{x}_{n,t}^{vel} - \tilde{\mu}_k^{vel})^2\}}{\sum_n \exp\{(\mathbf{x}_{n,t}^{vel} - \tilde{\mu}_k^{vel})^2\}}$$

1) significant mean velocity

2) large spatial support, with respect to
spatiotemporal blobs that have non null velocity but
tiny space/time support

A Bayesian model: feature likelihood given proto-objects



$$p(\mathbf{x}_t | \mathbf{r}_t, \mathcal{A}_t, \mathcal{O}) = \frac{\exp\{(\mathbf{x}_{n,t}^{vel} - \tilde{\mu}_k^{vel})^2\}}{\sum_n \exp\{(\mathbf{x}_{n,t}^{vel} - \tilde{\mu}_k^{vel})^2\}}$$

A Bayesian model: fixation probability

Focus based on
appearance & objects

$$p(\mathcal{O}_t | \mathbf{x}_t) \simeq \frac{p(\mathbf{x}_t | \mathbf{r}_t, \mathcal{A}_t, \mathcal{O}) p(\mathbf{r}_t | \mathcal{A}_t, \mathcal{O})}{p(\mathbf{x}_t)}$$

$$p(\mathbf{r}_t | \mathcal{A}_t, \mathcal{O}) = w_i \mathcal{N}(\mathbf{r}_t; \mu_{R_i}, \Sigma_{R_i}), \mathbf{r}_t \in R_i$$

(width, height) of region R_i

importance weight of
region R_i

$$w_i = |R_i| / M_R$$

A Bayesian model: fixation probability



$$p(\mathbf{r}_t | \mathcal{A}_t, \mathcal{O}) = w_i \mathcal{N}(\mathbf{r}_t; \mu_{R_i}, \Sigma_{R_i}), \mathbf{r}_t \in R_i$$

A Bayesian model: inference



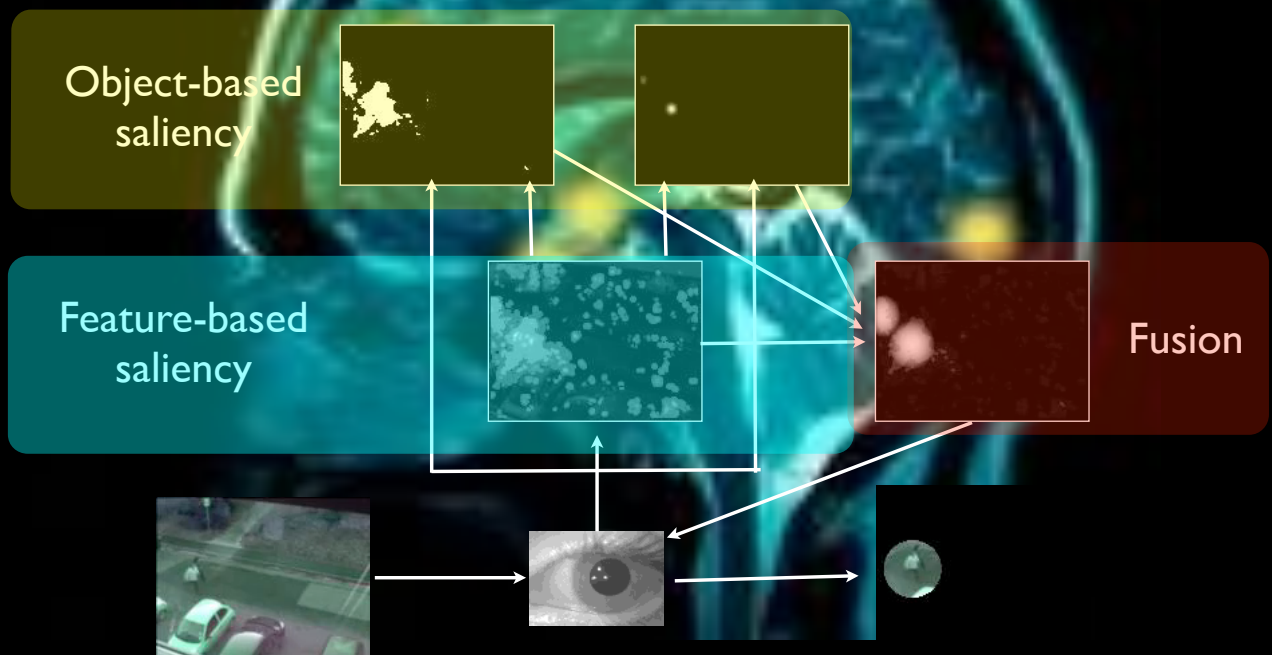
$$p(\mathcal{O}_t | \mathbf{x}_t) \simeq \frac{p(\mathbf{x}_t | \mathbf{r}_t, \mathcal{A}_t, \mathcal{O}) p(\mathbf{r}_t, | \mathcal{A}_t, \mathcal{O})}{p(\mathbf{x}_t)}$$

A Bayesian model: gaze selection



$$\hat{\mathcal{O}}_t = \arg \max_{\mathcal{O}_t} p(\mathcal{O}_t | \mathbf{x}_t)$$

The final picture



More examples



from: BEHAVE Interactions Test Case Scenarios
<http://groups.inf.ed.ac.uk/vision/BEHAVEDATA/>

More examples



from: BEHAVE Interactions Test Case Scenarios
<http://groups.inf.ed.ac.uk/vision/BEHAVEDATA/>

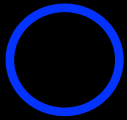
More examples



from: BEHAVE Interactions Test Case Scenarios
<http://groups.inf.ed.ac.uk/vision/BEHAVEDATA/>

Example 1 (unconstrained)

Average
eye-tracked
(36 observers)



Bayesian
Model

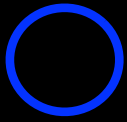


Example 2



Example 2 (unconstrained)

Average
eye-tracked
(36 observers)

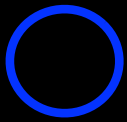


Bayesian
Model



Example 2 (constrained)

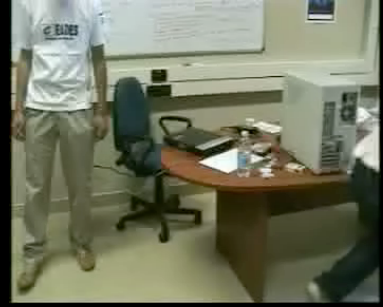
Average
eye-tracked
(36 observers)



Bayesian
Model



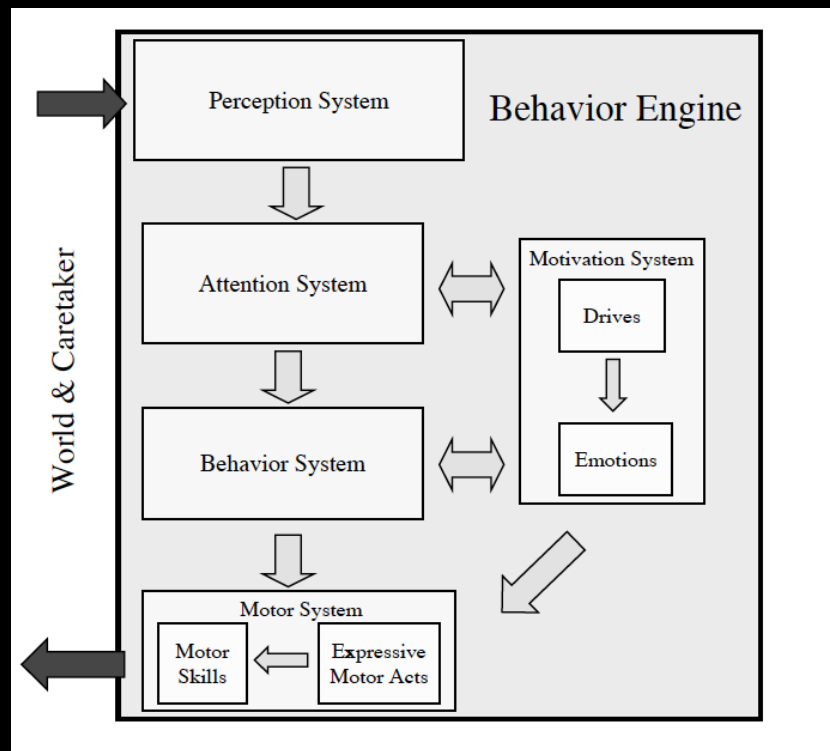
Situating vision in the world



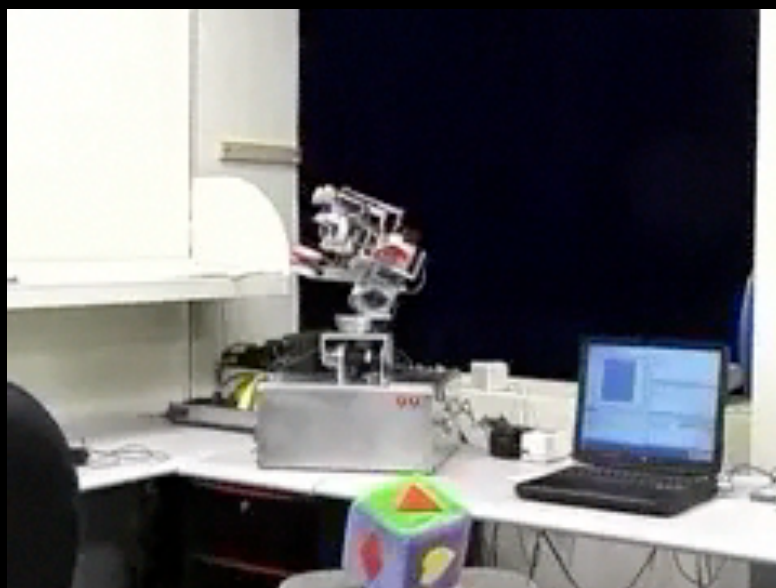
Situating vision in the world



Situating vision in the world

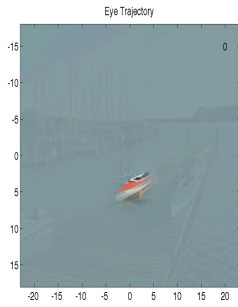


Situating vision in the world



Modelli di attenzione visiva

//livelli di spiegazione



Qual è il goal della computazione?

Quale rappresentazione e quale algoritmo?

Come realizzarla fisicamente?

che cosa guardo

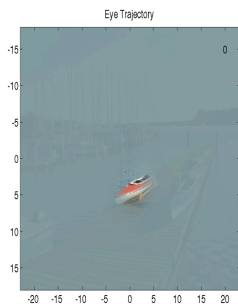
$$\mathbf{I} \mapsto \mathcal{R}$$

$$\mathcal{R} \mapsto \{r_F(1), r_F(2), \dots\}$$

come guardo

Modelli di attenzione visiva

//livelli di spiegazione: problema della variabilità



Qual è il goal della computazione?

Quale rappresentazione e quale algoritmo?

Come realizzarla fisicamente?

che cosa guardo

$$\mathbf{I} \mapsto \mathcal{R}$$

$$\mathcal{R} \mapsto \{r_F(1), r_F(2), \dots\}$$

come guardo

Modelli di attenzione visiva

//livelli di spiegazione: problema della variabilità



Modelli di attenzione visiva

//livelli di spiegazione: problema della variabilità



Modelli di attenzione visiva

//livelli di spiegazione: problema della variabilità



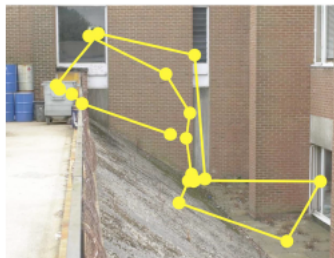
Modelli di attenzione visiva

//livelli di spiegazione: problema della variabilità

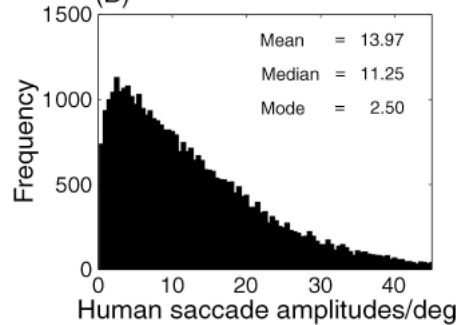
Journal of Vision (2011) 11(5):5, 1–23

Tatler, Hayhoe, Land, & Ballard

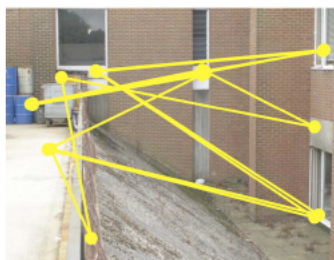
(A)



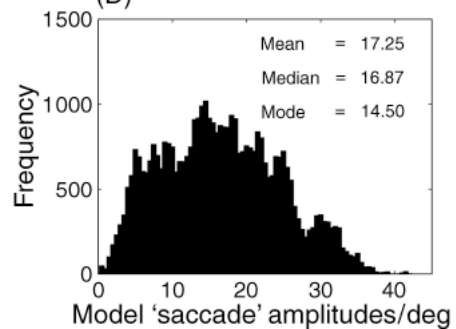
(B)



(C)



(D)



Modelli di attenzione visiva

//livelli di spiegazione: problema della variabilità

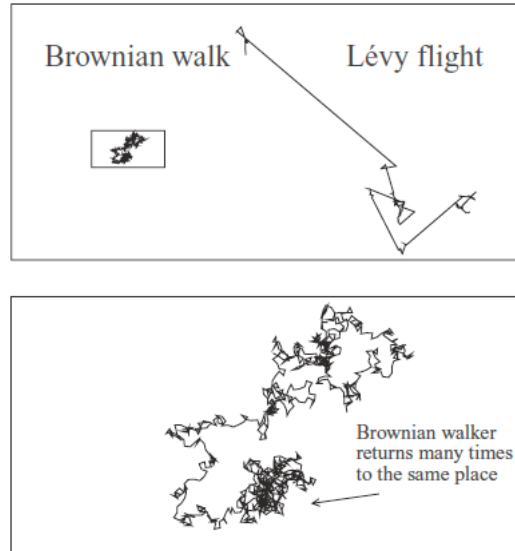
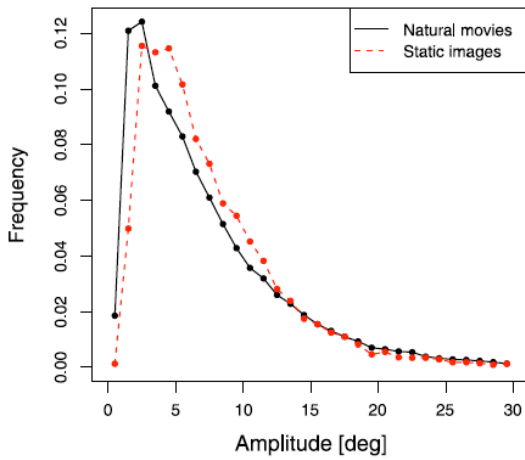
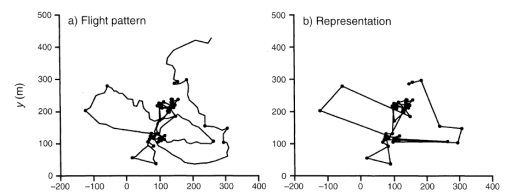
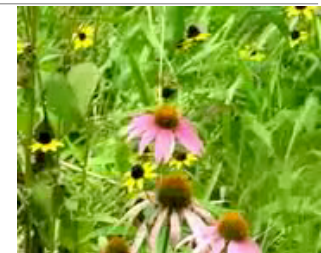
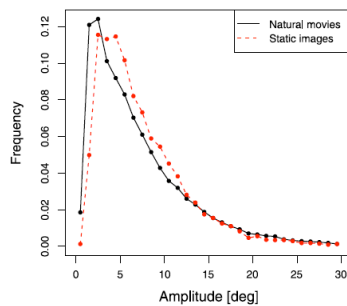
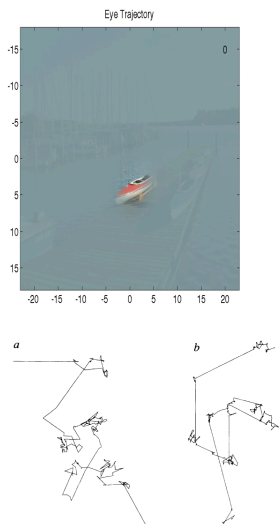


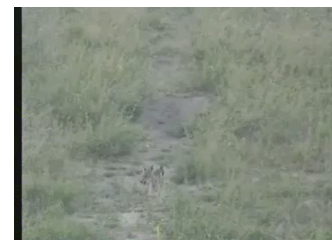
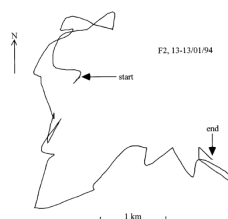
Figure 5.2 Two-dimensional Brownian random walk and Lévy flight of identical total length of 1000 units, shown to scale. (Bottom) Zoom of the Brownian walk. Note how the Brownian walker returns many times to previously visited locations, a phenomenon known as *oversampling*. In contrast, the Lévy flyer frequently takes ultralong jumps to virgin territory. This reduction in oversampling is part of the fundamental theoretical basis for interest in the Lévy flight foraging hypothesis.

Modelli di attenzione visiva

//livelli di spiegazione: problema della variabilità

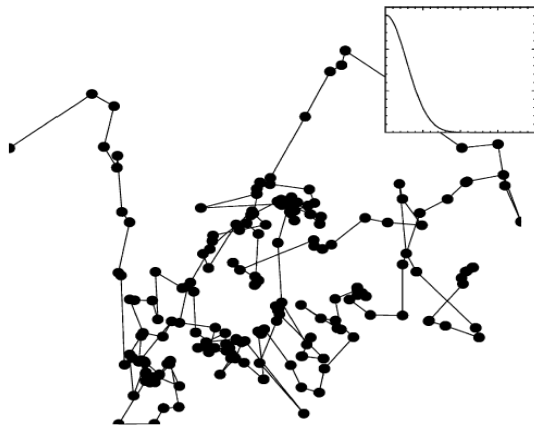


The Wandering Albatross....

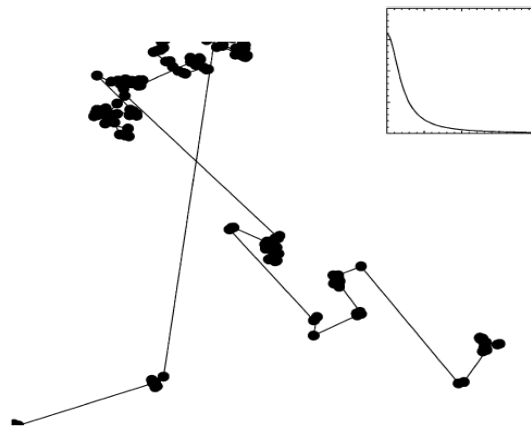


Modelli di attenzione visiva

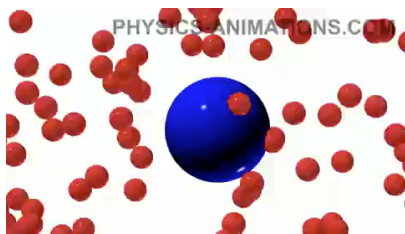
//livelli di spiegazione: problema della variabilità



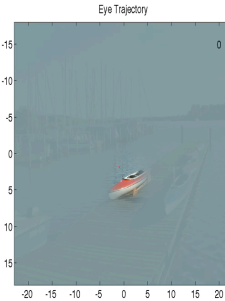
random walk Browniano



volo di Levy

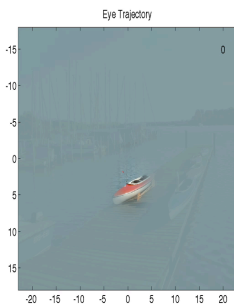


$$\mathbf{r}_{new}(t) = \mathbf{r}(t) - \nabla V + \eta$$



Modelli di attenzione visiva

//livelli di spiegazione: problema della variabilità



Qual è il goal della computazione?

Quale rappresentazione e quale algoritmo?

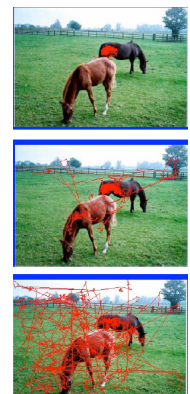
Come realizzarla fisicamente?

che cosa guardo

$$\mathbf{I} \mapsto \mathcal{R}$$

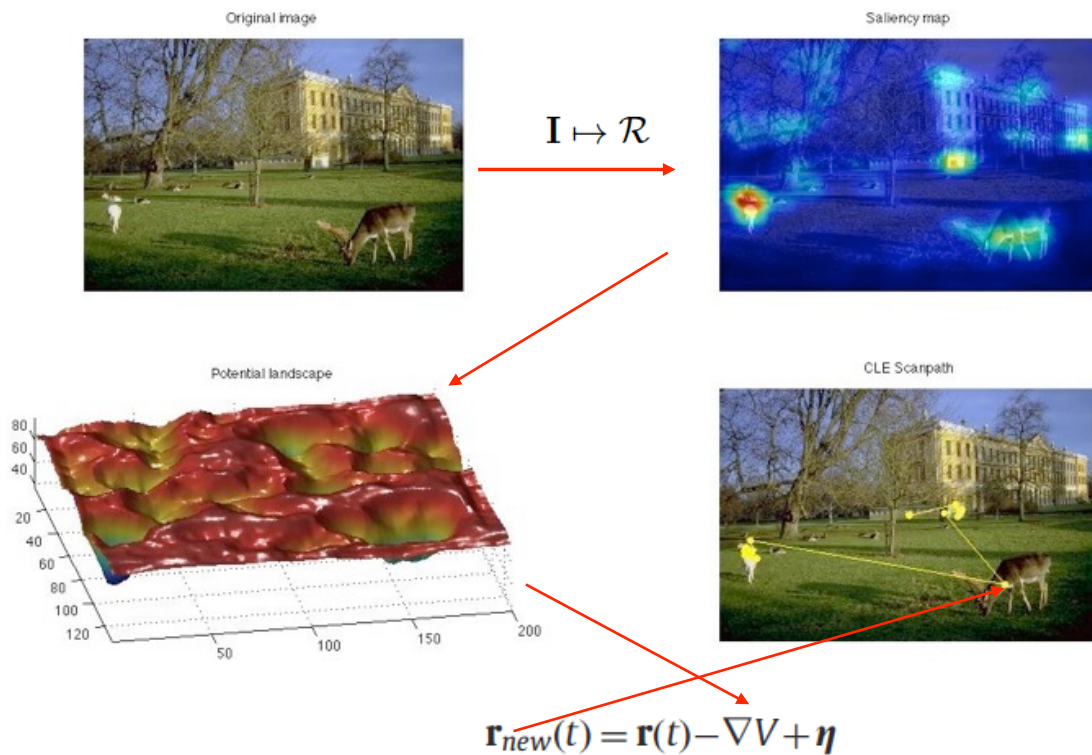
$\mathbf{r}_{new}(t) = \mathbf{r}(t) - \nabla V + \eta$
come guardo

simulazione Monte Carlo con accettazione Metropolis



Modelli di attenzione visiva

//livelli di spiegazione: problema della variabilità



Modelli di attenzione visiva

//livelli di spiegazione: problema della variabilità

Free Matlab Code:



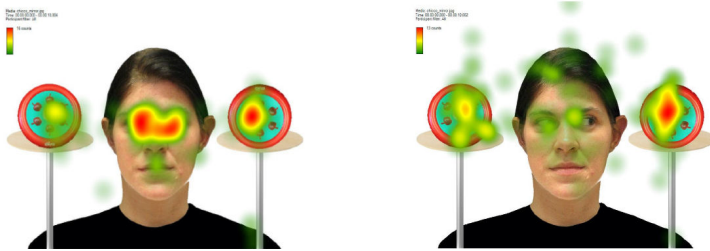
toolbox ufficiale su Matlab Central:

<http://www.mathworks.com/matlabcentral/fileexchange/38512-visual-scanpaths-via-constrained-levy-exploration-of-a-saliency-landscape>

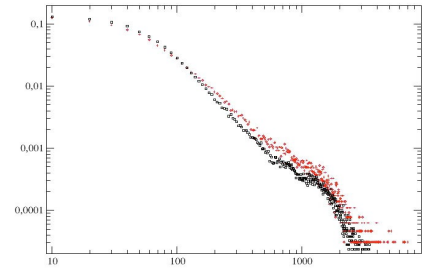
Modelli di attenzione visiva

//livelli di spiegazione: problema della variabilità

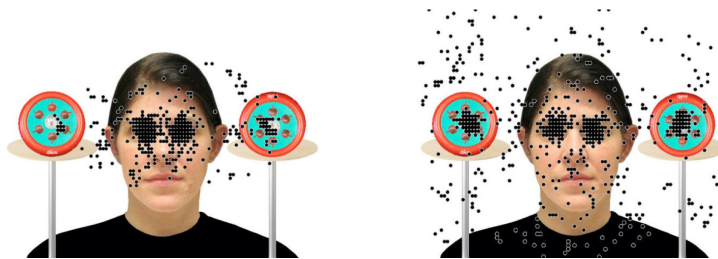
- Caratterizzazione di pazienti autistici



Experimental heat maps: TDs (left) and ASDs (right)



T is a sort of "attentional temperature", which seems to be able to describe the transition from TDs to ASDs visual behavior



1000 points generated by the model for $T=10$ (left) and $T=100$ (right) placed over a stimulus.

//the iCub

Modelli di attenzione visiva

//livelli di spiegazione: problema della variabilità

- Apprendimento stocastico nei robot

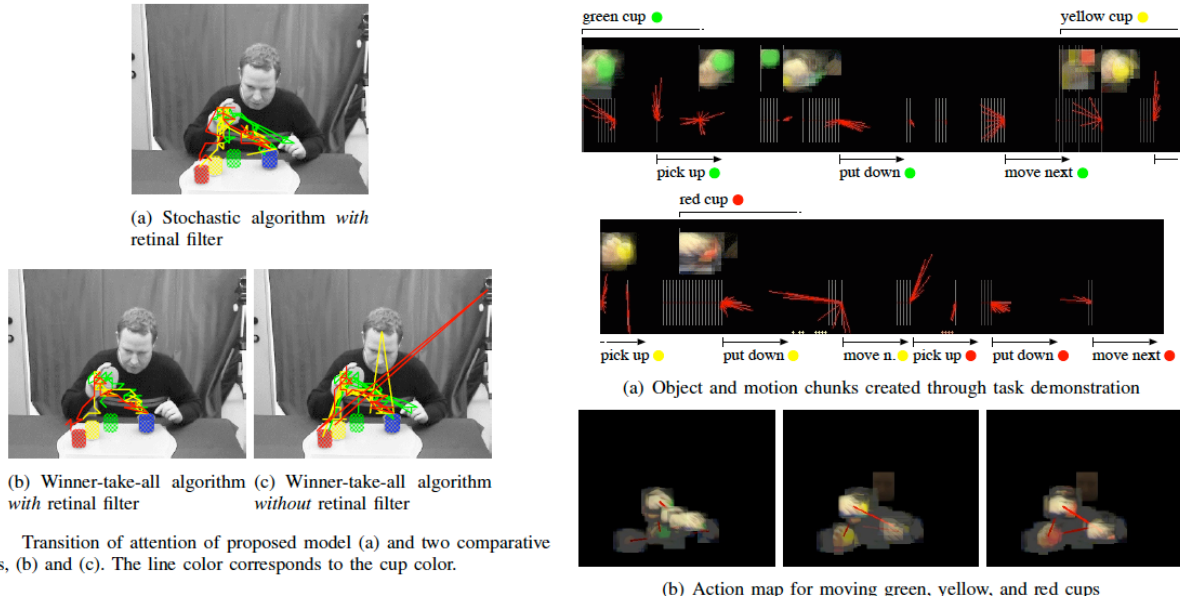
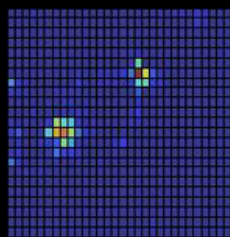
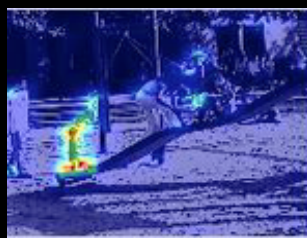


Fig. 7. Transition of attention of proposed model (a) and two comparative models, (b) and (c). The line color corresponds to the cup color.

IEEE TRANS. ON SMC-B

Ecological Sampling of Gaze Shifts

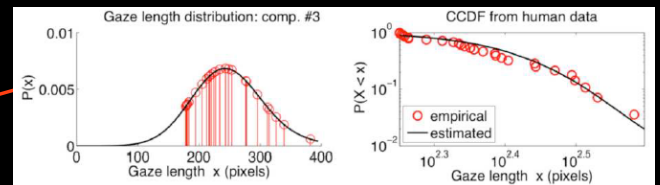
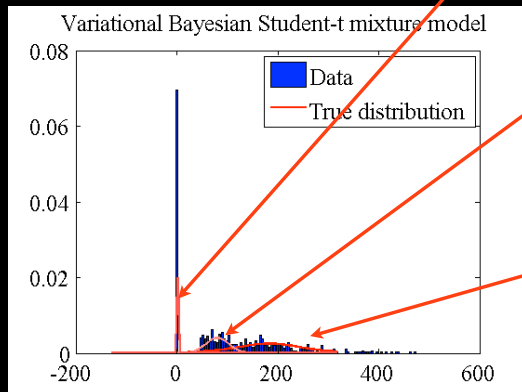
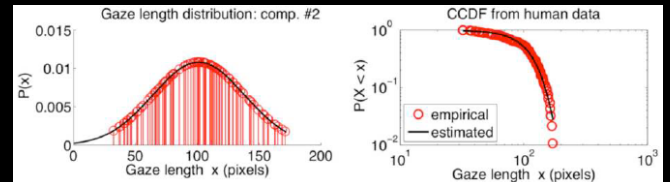
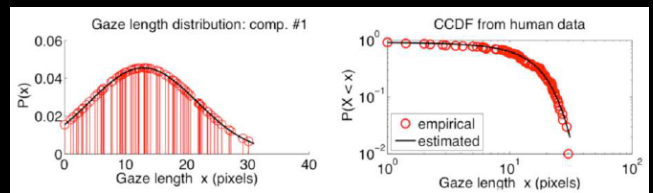
Giuseppe Boccignone and Mario Ferraro



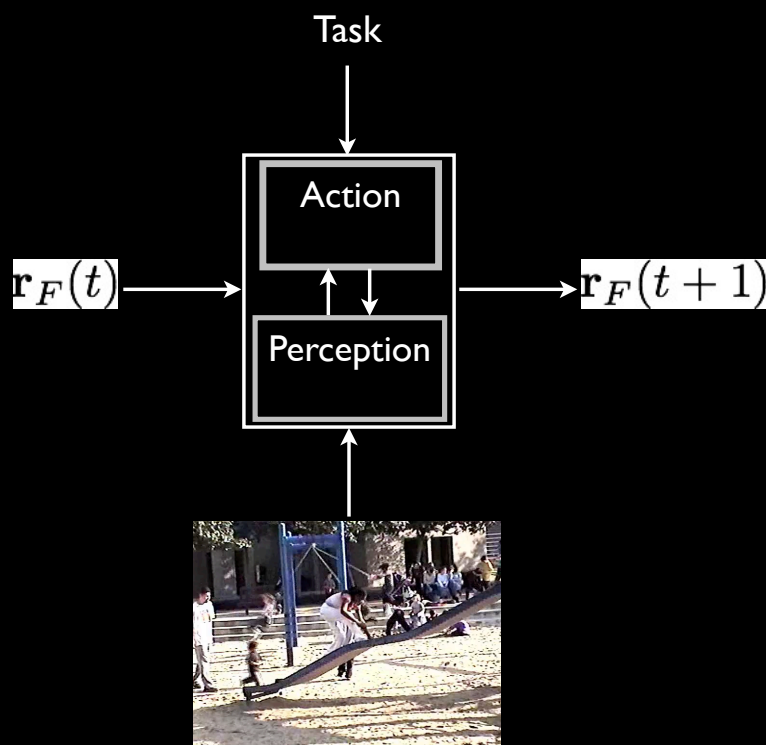
Matlab simulation code:

http://boccignone.di.unimi.it/Ecological_Sampling.html

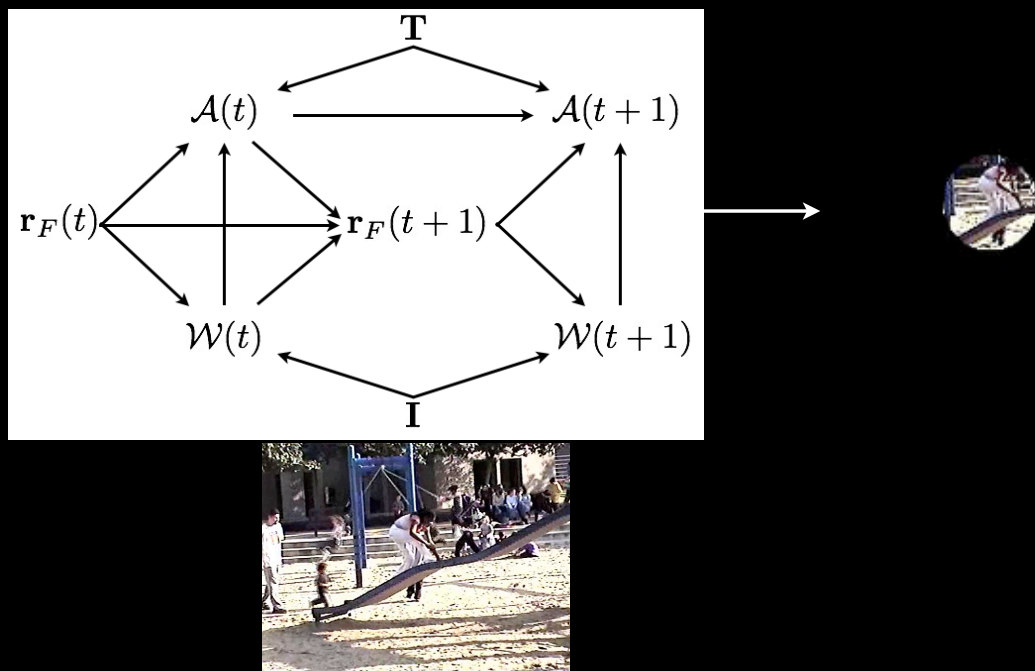
Ecological sampling: //data analysis



Ecological sampling of gaze shifts



Ecological sampling of gaze shifts //unfolding the action-perception loop



Ecological sampling of gaze shifts //sampling the action-perception loop

- Sampling the natural habitat

$$\mathcal{W}^*(t) \sim P(\mathcal{W}(t) | \mathbf{r}_F(t), \mathbf{F}(t), \mathbf{I}(t))$$

- Sampling the appropriate motor behavior

$$\mathcal{A}(t)^* \sim P(\mathcal{A}(t) | \mathcal{A}(t-1), \mathcal{W}^*(t))$$

- Sampling where to look next

$$\mathbf{r}_F(t+1) \sim P(\mathbf{r}_F(t+1) | \mathcal{A}(t)^*, \mathcal{W}^*(t), \mathbf{r}_F(t))$$

Ecological sampling of gaze shifts //sampling the natural habitat

- Sampling the saliency

$$\mathcal{S}^*(t) \sim P(\mathcal{S}(t) | \mathbf{F}(\widehat{\mathbf{I}}(t)))$$

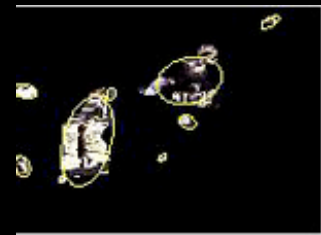
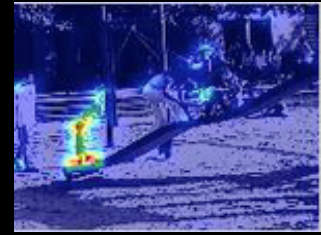
- Sampling patches...

$$\mathcal{M}^*(t) \sim P(\mathcal{M}(t) | \mathcal{S}^*(t))$$

- ...and their parametric descriptions

for $p = 1, \dots, N_P$

$$\theta_p^*(t) \sim P(\theta_p(t) | \mathcal{M}_p^*(t) = 1, \mathcal{S}^*(t))$$



Ecological sampling of gaze shifts //sampling the natural habitat

- Sampling “preys” from patches (food items / Interest Points)

$$O_p^*(t) \sim P(O_p(t) | \theta_p^*(t), \mathcal{M}_p^*(t) = 1, \mathcal{S}^*(t))$$

$$O(t) = \bigcup_{p=1}^{N_P} \{\mathbf{r}_{i,p}(t)\}_{i=1}^{N_{i,p}}$$

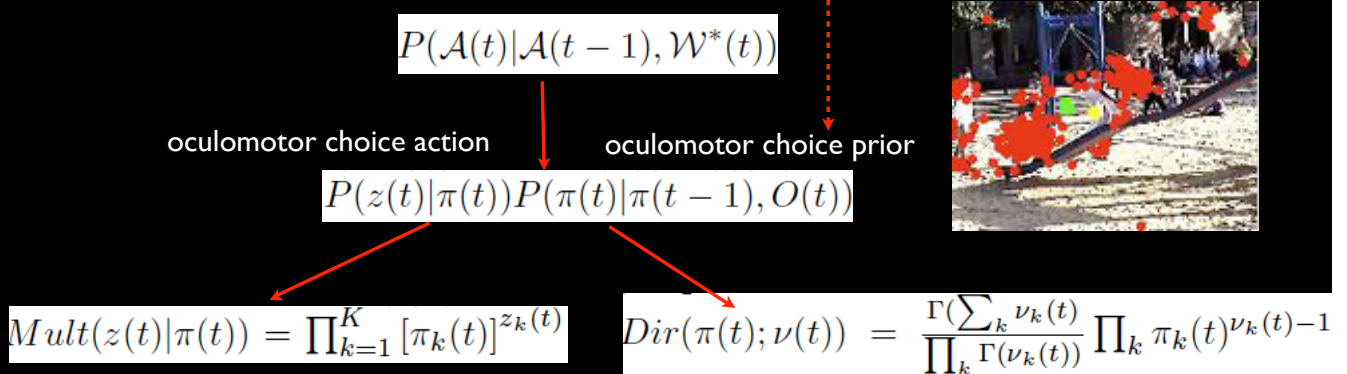
$$\mathbf{r}_{i,p} \sim \mathcal{N}(\mathbf{r}_p; \mu_p(t), \Sigma_p(t)), i = 1, \dots, N_{i,p}$$



Ecological sampling of gaze shifts

//sampling the motor behavior

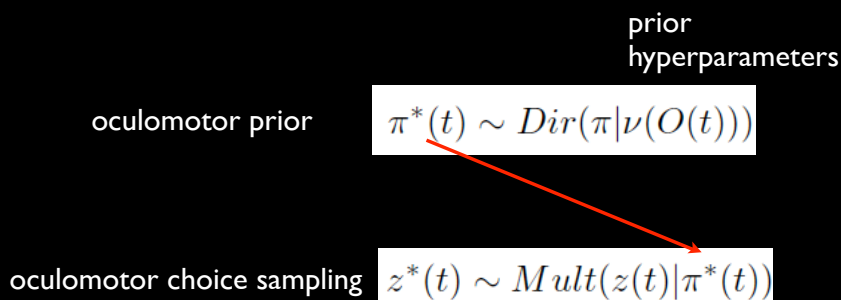
- Motor action $\mathcal{A}(t) = (z(t), \pi_t)$ determines the choice of oculomotor behavior
- Markov assumption $\mathcal{A}(t-1) \rightarrow \mathcal{A}(t)$



Ecological sampling of gaze shifts

//sampling the motor behavior

- Motor action $\mathcal{A}(t) = (z(t), \pi_t)$ determines the choice of oculomotor behavior
- Markov assumption $\mathcal{A}(t-1) \rightarrow \mathcal{A}(t)$



Ecological sampling of gaze shifts

//sampling the motor behavior

- Motor action $\mathcal{A}(t) = (z(t), \pi_t)$ determines the choice of oculomotor behavior
- Markov assumption $\mathcal{A}(t-1) \rightarrow \mathcal{A}(t)$

$$\pi^*(t) \sim \text{Dir}(\pi | \nu(O(t)))$$

$$\nu_k(t) = \nu_k(t-1) + [E_{\mathcal{C}(t)} = k], k = 1, \dots, K$$

prior
hyperparameters

$$\mathcal{C}(t) = \Delta(t) \cdot \Omega(t)$$



Ecological sampling of gaze shifts

//sampling the motor behavior

- Taking into account the complexity of the landscape in natural habitats

no cues



some cues



too many
cues



$$\mathcal{C}(t) = \Delta(t) \cdot \Omega(t)$$

Ecological sampling of gaze shifts

//sampling the motor behavior

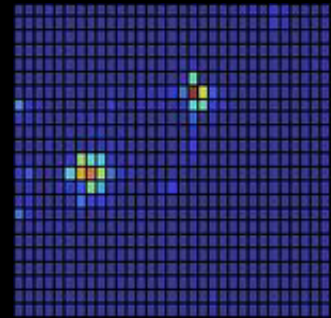
- Motor action $\mathcal{A}(t) = (z(t), \pi_t)$ determines the choice of oculomotor behavior
- Markov assumption $\mathcal{A}(t-1) \rightarrow \mathcal{A}(t)$

$$H(t) = -k_B \sum_{c=1}^{N_w} P(c, t) \log P(c, t)$$

$$\Delta \equiv H/H_{sup}$$

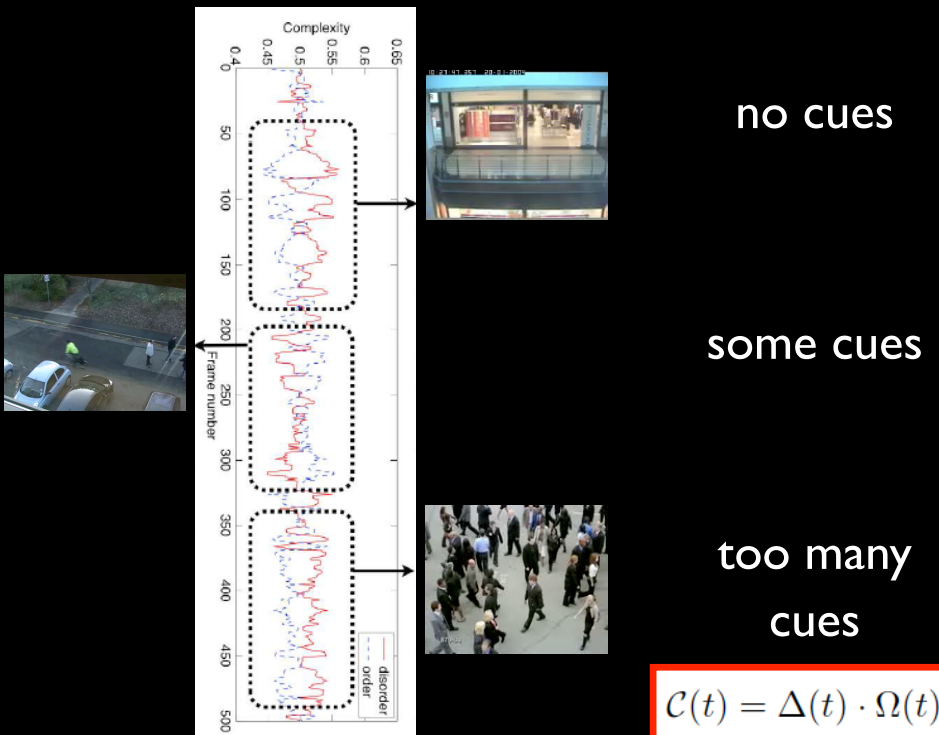
$$\Omega = 1 - \Delta$$

$$\mathcal{C}(t) = \Delta(t) \cdot \Omega(t)$$



Ecological sampling of gaze shifts

//sampling the motor behavior



Ecological sampling of gaze shifts

//sampling where to look next

- Motor action $\mathcal{A}(t) = (z(t), \pi_t)$ determines the choice of oculomotor behavior
- Sampling the next FOA

$$\mathbf{r}_F(t+1) \sim P(\mathbf{r}_F(t+1) | z^*(t) = k, \theta^*(t), \eta, \mathbf{r}_F(t))$$

- Langevin - Smoluchowski

$$d\mathbf{r}_F(t) = -\nabla V(\mathbf{r}_F, t)dt + \mathbf{D}(\mathbf{r}_F, t)\boldsymbol{\xi}_k(t)dt$$

Ecological sampling of gaze shifts

//sampling where to look next

- Motor action $\mathcal{A}(t) = (z(t), \pi_t)$ determines the choice of alpha stable parameters $f(\boldsymbol{\xi}; \alpha, \beta, \gamma, \delta)$
- Sampling the next FOA $\boldsymbol{\xi} \sim$

$$d\mathbf{r}_F(t) = -\nabla V(\mathbf{r}_F, t)dt + \mathbf{D}(\mathbf{r}_F, t)\boldsymbol{\xi}_k(t)dt$$

$$-\nabla V(\mathbf{r}_F, t) = -\sum_{p=1}^{N_V} (\mathbf{r}_F(t) - \mathbf{r}_p(t))$$

$$\mathbf{r}_F(t_{n+1}) \approx \mathbf{r}_F(t_n) - \sum_{p=1}^{N_V} (\mathbf{r}_F(t_n) - \mathbf{r}_p(t_n))\tau + \gamma_k \mathbb{I}\tau^{1/\alpha_k} \boldsymbol{\xi}_k.$$



Ecological sampling of gaze shifts

//sampling where to look next

no cues



#Levy flights >
#Gaussian flights

some cues



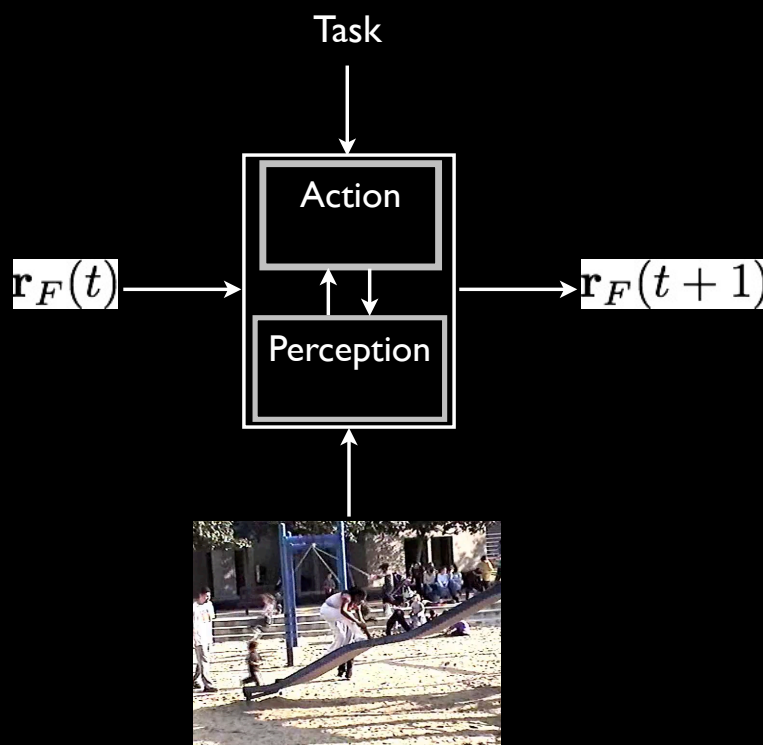
#Gaussian flights >
#Levy flights

too many
cues

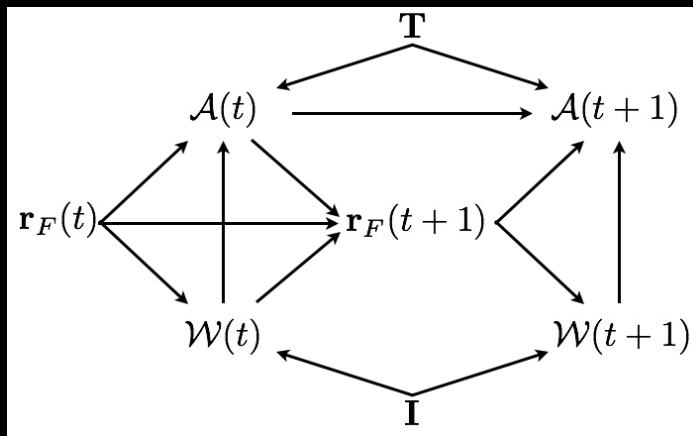


#Levy flights >
#Gaussian flights

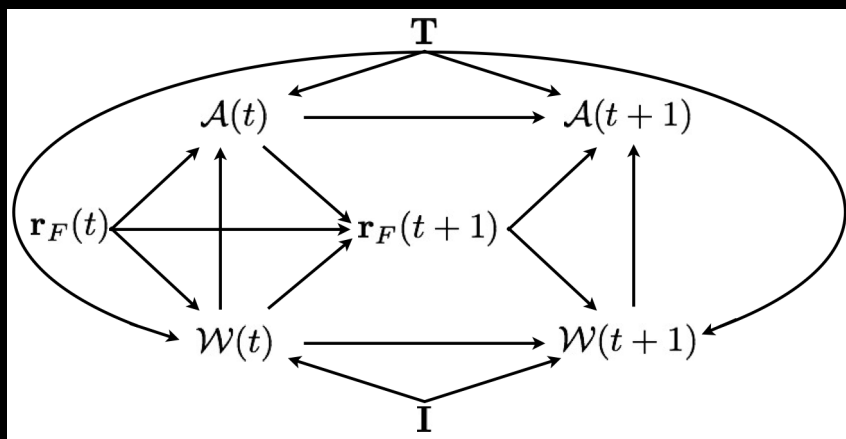
Ecological sampling of gaze shifts



Ecological sampling of gaze shifts //accounting for task

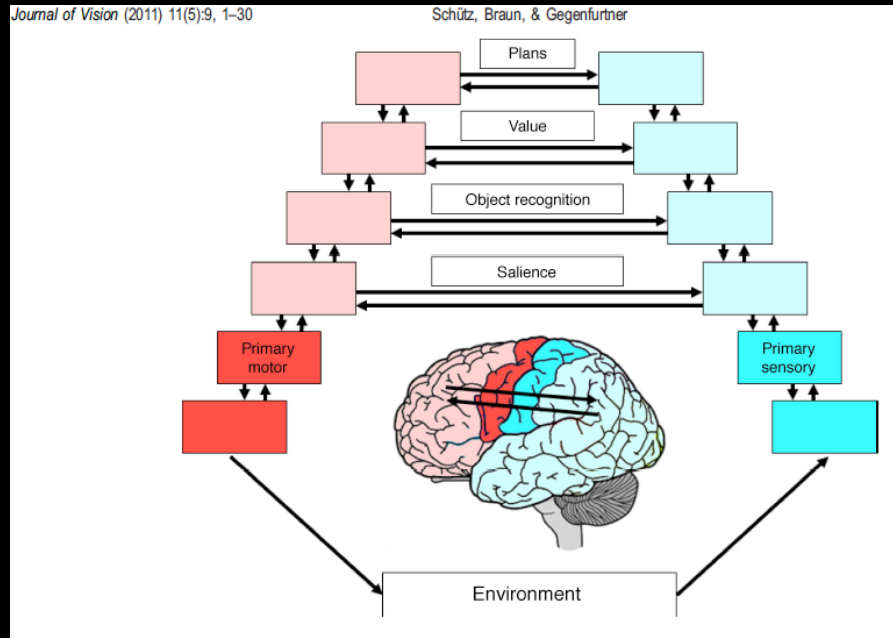


Ecological sampling of gaze shifts //accounting for task

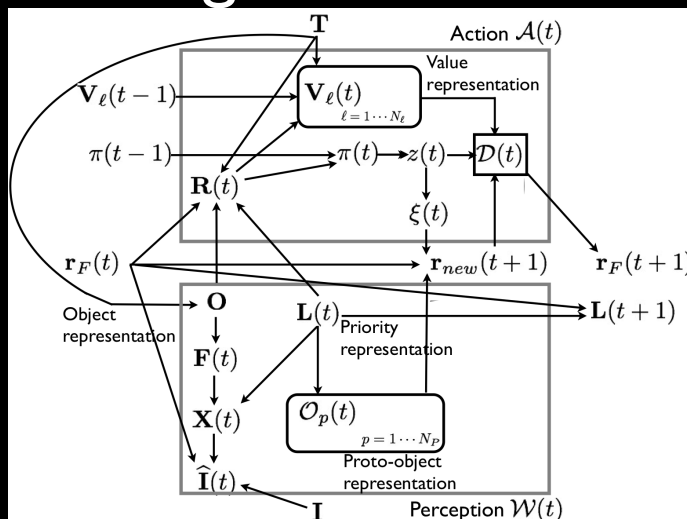


Ecological sampling of gaze shifts //accounting for task

- There are multiple levels of representation



Ecological sampling of gaze shifts //accounting for task



every gaze shift is a decision

