Bottom-up and top-down processing in visual perception

Thomas Serre

Brown University

Department of Cognitive & Linguistic Sciences Brown Institute for Brain Science Center for Vision Research

The 'what' problem monkey electrophysiology



Willmore Prenger & Gallant '10

The 'what' problem

monkey electrophysiology



Yamane Carlson Bowman Wang & Connor '08







The 'what' problem

human psychophysics

Thorpe et al '96; VanRullen & Thorpe '01; Fei-Fei et al '02 '05; Evans & Treisman '05; Serre Oliva & Poggio '07



Vision as 'knowing what is where'



'What' and 'where' cortical pathways

Ungerleider & Mishkin '84



'What' and 'where' cortical pathways



cortical pathways

How does the visual system combine information about the identity and location of objects?



cortical pathways

How does the visual system combine information about the identity and location of objects?

Central thesis: visual attention

(see also Van Der Velde and De Kamps '01; Deco and Rolls '04)



cortical pathways

Part I: A 'Bayesian' theory of attention that solves the 'what is where' problem

• Predictions for neurophysiology and human eye movement data

- Sharat Chikkerur
- Cheston Tan
- Tomaso Poggio

Part II: Readout of IT population activity

Attention eliminates clutter

- Ying Zhang
- Ethan Meyers
- Narcisse Bichot
- Tomaso Poggio
- Bob Desimone

Part I: A 'Bayesian' theory of attention that solves the 'what is where' problem

• Predictions for neurophysiology and human eye movement data

- Sharat Chikkerur
- Cheston Tan
- Tomaso Poggio



- Part II: Readout of IT population activity
 - Attention eliminates clutter

- Ying Zhang
- Ethan Meyers
- Narcisse Bichot
- Tomaso Poggio
- Bob Desimone



Perception as Bayesian inference

$P(S|I) \propto P(I|S)P(S)$



Perception as Bayesian inference

$P(S|I) \propto P(I|S)P(S)$





Perception as Bayesian inference



$P(S|I) \propto P(I|S)P(S)$





Perception as Bayesian inference



Perception as Bayesian inference

 To recognize and localize objects in the scene, the visual system selects objects, one object at a time



P(O, L, I)



Assumption #1: Attentional spotlight

Broadbent '52 '54; Treisman '60; Treisman & Gelade '80; Duncan & Desimone '95; Wolfe '97; and many others • Object location *L* and identity *O* are independent







Assumption #2: 'what' and 'where' pathways • Object location *L* and identity *O* are independent

P(O, L, I) = P(O)P(L)P(I|L, O)



Assumption #2: 'what' and 'where' pathways

- Objects encoded by collections of 'universal' features (of intermediate complexity)
 - Either present or absent
 - Conditionally independent given the object and its location



Assumption #3: universal features

see Riesenhuber & Poggio '99; Serre et al '05 '07

- Objects encoded by collections of 'universal' features (of intermediate complexity)
 - Either present or absent
 - Conditionally independent given the object and its location



Assumption #3: universal features

see Riesenhuber & Poggio '99; Serre et al '05 '07













Model: Serre et al '05 Experimental data: Hung* Kreiman*et al '05

Serre Oliva & Poggio '07



Serre Oliva & Poggio '07





Serre Oliva & Poggio '07





Serre Oliva & Poggio '07





- Goal of visual perception: to estimate posterior probabilities of visual features, objects and their locations in an image
- Attention corresponds to conditioning on high-level latent variables representing particular features or locations (as well as on sensory input), and doing inference over the other latent variables





$$P(X^{i}|I) = \frac{P(I|X^{i}) \sum_{F^{i},L} P(X^{i}|F^{i},L)P(L)P(F^{i})}{\sum_{X^{i}} \left\{ P(I|X^{i}) \sum_{F^{i},L} P(X^{i}|F^{i},L)P(L)P(F^{i}) \right\}}$$










John H. Reynolds^{1,*} and David J. Heeger²

¹Salk Institute for Biological Studies, La Jolla, CA 92037-1099, USA

²Department of Psychology and Center for Neural Science, New York University, New York, NY 10003, USA *Correspondence: reynolds@salk.edu

DOI 10.1016/j.neuron.2009.01.002

$$R(x,\theta) = \frac{A(x,\theta)E(x,\theta)}{S(x,\theta) + \sigma}$$



Bayesian inference and attention



$$P(X^{i} = x|I) \propto \sum_{F^{i},L} P(X^{i} = x|F^{i},L)P(I|X^{i})P(F^{i})P(L)$$



Multiplicative scaling of tuning curves by spatial attention

$$P(X^{i}|I) = \frac{P(I|X^{i}) \sum_{F^{i},L} P(X^{i}|F^{i},L)P(L)P(F^{i})}{\sum_{X^{i}} \left\{ P(I|X^{i}) \sum_{F^{i},L} P(X^{i}|F^{i},L)P(L)P(F^{i}) \right\}}$$

Trujillo and Treue '02

Mc Adams and Maunsell '99





Contrast vs. response gain

$$P(X^{i} = x|I) \propto \sum_{F^{i},L} P(X^{i} = x|F^{i},L)P(I|X^{i})P(F^{i})P(L)$$





Feature-based attention

neural data from Bichot et al '05











- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker

- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker

Uniform priors (bottom-up) Feature priors Feature + contextual (spatial) priors Humans

- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker

Uniform priors (bottom-up) Feature priors Feature + contextual (spatial) priors Humans

Uniform priors (bottom-up)
Feature priors
Feature + contextual (spatial) priors
Humans

- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker





- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker



Uniform priors (bottom-up)
Feature priors
Feature + contextual (spatial) priors
Humans

The experiment

• Eye movements as proxy for attention

• Dataset:

- 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker

Uniform priors (bottom-up)
Feature priors
Feature + contextual (spatial) priors
Humans

- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker



Uniform priors (bottom-up)
Feature priors
Feature + contextual (spatial) priors
Humans

- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker



Uniform priors (bottom-up) Feature priors Feature + contextual (spatial) priors Humans

- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker



Uniform priors (bottom-up) Feature priors Feature + contextual (spatial) priors Humans

- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker



Uniform priors (bottom-up)
Feature priors
Feature + contextual (spatial) priors
Humans

- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker



Uniform priors (bottom-up)
Feature priors
Feature + contextual (spatial) priors
Humans

- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker



Uniform priors (bottom-up)
Feature priors
Feature + contextual (spatial) priors
Humans

- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker

Explains 92% of the inter-subject agreement!



1st three fixations

Uniform priors (bottom-up)
Feature priors
Feature + contextual (spatial) priors
Humans

The experiment

- Eye movements as proxy for attention
- Dataset:
 - 100 street-scenes images with cars & pedestrians and 20 without
- Experiment
 - 8 participants asked to count the number of cars/pedestrians
 - block design for cars and pedestrians
 - eye movements recorded using an infra-red eye tracker

*similar (independent) results by Ehinger Hidalgo Torralba & Oliva '10

Method	ROC area
Bruce & Tsotos '06	72.8%
Itti et al '01	72.7%
Proposed	77.9%





Bottom-up saliency and free-viewing

human eye data from Bruce & Tsotsos

- Goal of vision:
 - To solve the problem of 'what is where' problem

- Goal of vision:
 - To solve the problem of 'what is where' problem
- Key assumptions:
 - 'Attentional spotlight' \rightarrow recognition done sequentially, one object at a time
 - 'What' and 'where' independent pathways

- Goal of vision:
 - To solve the problem of 'what is where' problem
- Key assumptions:
 - 'Attentional spotlight' \rightarrow recognition done sequentially, one object at a time
 - 'What' and 'where' independent pathways
- Attention as the inference process implemented by the interaction between ventral and dorsal areas
 - Integrates bottom-up and top-down (feature-based and context-based) attentional mechanisms
 - Seems consistent with known physiology

- Goal of vision:
 - To solve the problem of 'what is where' problem
- Key assumptions:
 - 'Attentional spotlight' \rightarrow recognition done sequentially, one object at a time
 - 'What' and 'where' independent pathways
- Attention as the inference process implemented by the interaction between ventral and dorsal areas
 - Integrates bottom-up and top-down (feature-based and context-based) attentional mechanisms
 - Seems consistent with known physiology
- Main attentional effects in the presence of clutters
 - Spatial attention reduces uncertainty over location and improves object recognition performance over first bottom-up (feedforward) pass

Part I: A 'Bayesian' theory of attention that solves the 'what is where' problem

• Predictions for neurophysiology and human eye movement data

- Sharat Chikkerur
- Cheston Tan
- Tomaso Poggio

Part II: Readout of IT population activity

Attention eliminates clutter

- Ying Zhang
- Ethan Meyers
- Narcisse Bichot
- Tomaso Poggio
- Bob Desimone




















The 'readout' approach



The 'readout' approach





train readout classifier on isolated object



test generalization in clutter









































- Robert Desimone (MIT)
- Christof Koch (CalTech)
- Laurent Itti (USC)
- Tomaso Poggio (MIT)
- David Sheinberg (Brown)





Thank you!