Deep learning in the visual cortex

Thomas Serre Brown University

- I. Fundamentals of primate vision
- II. <u>Computational mechanisms of rapid</u> recognition and feedforward processing
- III. <u>Beyond feedforward processing:</u> <u>Attentional mechanisms and cortical</u> <u>feedback</u>







Feedforward processing



- Coarse initial base representation
- Enables rapid object detection/recognition ('what is there?')
- Insufficient for object localization
- Sensitive to presence of clutter

### Feedforward processing



#### isolated objects



#### isolated objects



isolated objects







**Prediction**: recognition in clutter requires attentional mechanisms via cortical feedback

clutter



spatial attention

**Prediction**: recognition in clutter requires attentional mechanisms via cortical feedback

clutter



spatial attention

### 1. Monkey electrophysiology

Read-out of inferior temporal cortex population activity: Spatial attention eliminates clutter

with Zhang, Meyers, Bichot, Poggio & Desimone

# 2. Computational model of integrated attention and recognition

What and where: a Bayesian attention theory of attention

with Chikkerur, Tan & Poggio

### 1. Monkey electrophysiology

Read-out of inferior temporal cortex population activity: Spatial attention eliminates clutter

with Zhang, Meyers, Bichot, Poggio & Desimone

2. Computational model of integrated attention and recognition

What and where: a Bayesian attention theory of attention

with Chikkerur, Tan & Poggio

### The 'readout' approach



Zhang Meyers Bichot Serre Poggio Desimone PNAS'11



Zhang Meyers Bichot Serre Poggio Desimone unpublished data

























## Changes in the salience of distractor stimuli dominate over attention related enhancements



Aligned to the time when one of the

• Consistent with feedforward hierarchical models: In the absence of attention information about the identity of individual objects (and position) in clutter is greatly reduced relative to when objects are shown in isolation

- Consistent with feedforward hierarchical models: In the absence of attention information about the identity of individual objects (and position) in clutter is greatly reduced relative to when objects are shown in isolation
- Attention seems to restore the pattern of neural activity toward the vector representing the isolated object

- Consistent with feedforward hierarchical models: In the absence of attention information about the identity of individual objects (and position) in clutter is greatly reduced relative to when objects are shown in isolation
- Attention seems to restore the pattern of neural activity toward the vector representing the isolated object
- In spite of this nearly exclusive representation of the attended object, an increase in the salience of non-attended objects overrode these attentional enhancements

- Consistent with feedforward hierarchical models: In the absence of attention information about the identity of individual objects (and position) in clutter is greatly reduced relative to when objects are shown in isolation
- Attention seems to restore the pattern of neural activity toward the vector representing the isolated object
- In spite of this nearly exclusive representation of the attended object, an increase in the salience of non-attended objects overrode these attentional enhancements
- Results provide computational level explanation for how attention operates on neural representations to solve the problem of invariant recognition in clutter

 Monkey electrophysiology
Read-out of inferior temporal cortex population activity: Spatial attention eliminates clutter

with Zhang, Meyers, Bichot, Poggio & Desimone

# 2. Computational model of integrated attention and recognition

What and where: a Bayesian attention theory of attention

with Chikkerur, Tan & Poggio





### Perception as Bayesian inference





### Perception as Bayesian inference





### Perception as Bayesian inference

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





### Perception as Bayesian inference

### Hypothesis #1

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




Two independent streams of processing for object and location

#### Hypothesis #2



Objects encoded by collections of generic features (cond. ind. given an object and its location)

### Hypothesis #3



(position implicit)

### Perception as Bayesian inference





### Perception as Bayesian inference

#### p(L|I) saliency map





### Perception as Bayesian inference

- Goal of visual perception is 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 objects or locations (as well as on sensory input), and doing inference over the other latent variables
- Here we used belief propagation to solve the inference problem



Biologically-plausible implementations of belief propagation: Zemel et al. '98; Beck & Pouget '07; Deneve '08; George '08; Litvak & Ullman '09; Rao '04; Steimer et al. '09

Special case of the normalization model of attention by Reynolds & Heeger '09

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

0

V1







ventral / what

feedforward input

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

Bayesian inference and attention

Special case of the normalization model of attention by Reynolds & Heeger '09

sweep



Special case of the normalization model of attention by Reynolds & Heeger '09





### (parallel) feature-based attention

feedforward input

=

 $P(X^i|I)$ 



suppressive drive

Bayesian inference and attention



Consistent with data from V4 by Bichot et al '05



Consistent with data from V4 by Bichot et al '05



Consistent with data from V4 by Bichot et al '05



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







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

Predicted by Reynolds & Heeger '09

### Learning to localize cars and pedestrians in street scenes





## Learning to localize cars and pedestrians in street scenes









### The experiment

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

![](_page_57_Picture_9.jpeg)

![](_page_58_Figure_1.jpeg)

![](_page_59_Figure_1.jpeg)

Car Search

Car Search

![](_page_60_Picture_2.jpeg)

![](_page_60_Picture_3.jpeg)

Pedestrian Search

![](_page_60_Picture_5.jpeg)

![](_page_60_Picture_6.jpeg)

![](_page_60_Figure_7.jpeg)

![](_page_60_Figure_8.jpeg)

L

![](_page_61_Figure_1.jpeg)

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

1st three fixations

![](_page_62_Figure_2.jpeg)

![](_page_62_Picture_3.jpeg)

![](_page_63_Figure_1.jpeg)

![](_page_63_Figure_2.jpeg)

![](_page_64_Figure_1.jpeg)

![](_page_64_Figure_2.jpeg)

![](_page_65_Figure_1.jpeg)

![](_page_65_Figure_2.jpeg)

![](_page_66_Figure_1.jpeg)

Feature + contextual (spatial) priors Humans

![](_page_66_Figure_3.jpeg)

![](_page_67_Figure_1.jpeg)

Feature + contextual (spatial) priors Humans

![](_page_67_Figure_3.jpeg)

![](_page_68_Figure_1.jpeg)

Feature + contextual (spatial) priors Humans

![](_page_68_Figure_3.jpeg)

![](_page_69_Figure_1.jpeg)

Feature + contextual (spatial) priors Humans

![](_page_69_Figure_3.jpeg)

![](_page_70_Figure_1.jpeg)

Feature + contextual (spatial) priors Humans

![](_page_70_Figure_3.jpeg)

![](_page_71_Figure_1.jpeg)

Feature + contextual (spatial) priors Humans

![](_page_71_Figure_3.jpeg)


Feature + contextual (spatial) priors

Feature priors Humans



Predicting eye movements during searches for cars and pedestrians

Overall model accounts for 92% of inter-subject agreement!

\*similar (independent) results by Ehinger Hidalgo Torralba & Oliva (in press)

# Predicting eye movements during free viewing

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

human eye data from Bruce & Tsotsos

• Attention as part of the inference process that solves the visual recognition problem of 'what is where'

- Attention as part of the inference process that solves the visual recognition problem of 'what is where'
- Main goal of the visual system is to infer the identity and the position of objects in visual scenes:
  - Spatial attention emerges as a strategy to reduce the uncertainty in shape information while feature-based attention reduces the uncertainty in spatial information
  - Featural and spatial attention represent two distinct modes of a computational process solving the problem of recognizing and localizing objects, especially in difficult recognition tasks such as in cluttered natural scenes

- Attention as part of the inference process that solves the visual recognition problem of 'what is where'
- Main goal of the visual system is to infer the identity and the position of objects in visual scenes:
  - Spatial attention emerges as a strategy to reduce the uncertainty in shape information while feature-based attention reduces the uncertainty in spatial information
  - Featural and spatial attention represent two distinct modes of a computational process solving the problem of recognizing and localizing objects, especially in difficult recognition tasks such as in cluttered natural scenes
- Model agnostic about the specific algorithm for the inference process (i.e., no claim made about the brain computing probabilities explicitly)

#### Two modes of vision

## Two modes of vision

- Rapid bottom-up / feedforward processing during first 100-150ms of visual processing:
  - Coarse/initial base representation
  - Enables rapid object detection/recognition ('what is there?')

## Two modes of vision

- Rapid bottom-up / feedforward processing during first 100-150ms of visual processing:
  - Coarse/initial base representation
  - Enables rapid object detection/recognition ('what is there?')
- Top-down / re-entrant attentional processing
  - Enables recognition in clutter
  - Enables object localization

#### **Collaborators (MIT):**

- Narcisse Bichot
- Sharat Chikkerur
- Bob Desimone
- Ethan Meyers
- Tomaso Poggio
- Cheston Tan
- Ying Zhang



#### Robert J. and Nancy D. Carney Fund for Scientific Innovation

# Acknowledgments