

# MODELS OF SALIENCE AND EYE MOVEMENTS

NICHOLAS BUTKO,  
JAVIER MOVELLAN, GARRISON COTTRELL,  
TERRY SEJNOWSKI, MATT TONG,  
CHRIS KANAN, LINGYUN ZHANG,  
LEANNE CHUKOSKIE, MIKE MOZER,  
MIKE ARNOLD

# MATHEMATICAL FORMULATION OF INFORMATION

☼ “Information” has two (related) mathematical meanings:

☼ 1) How much did you expect something you experience (“it is going to rain”)?

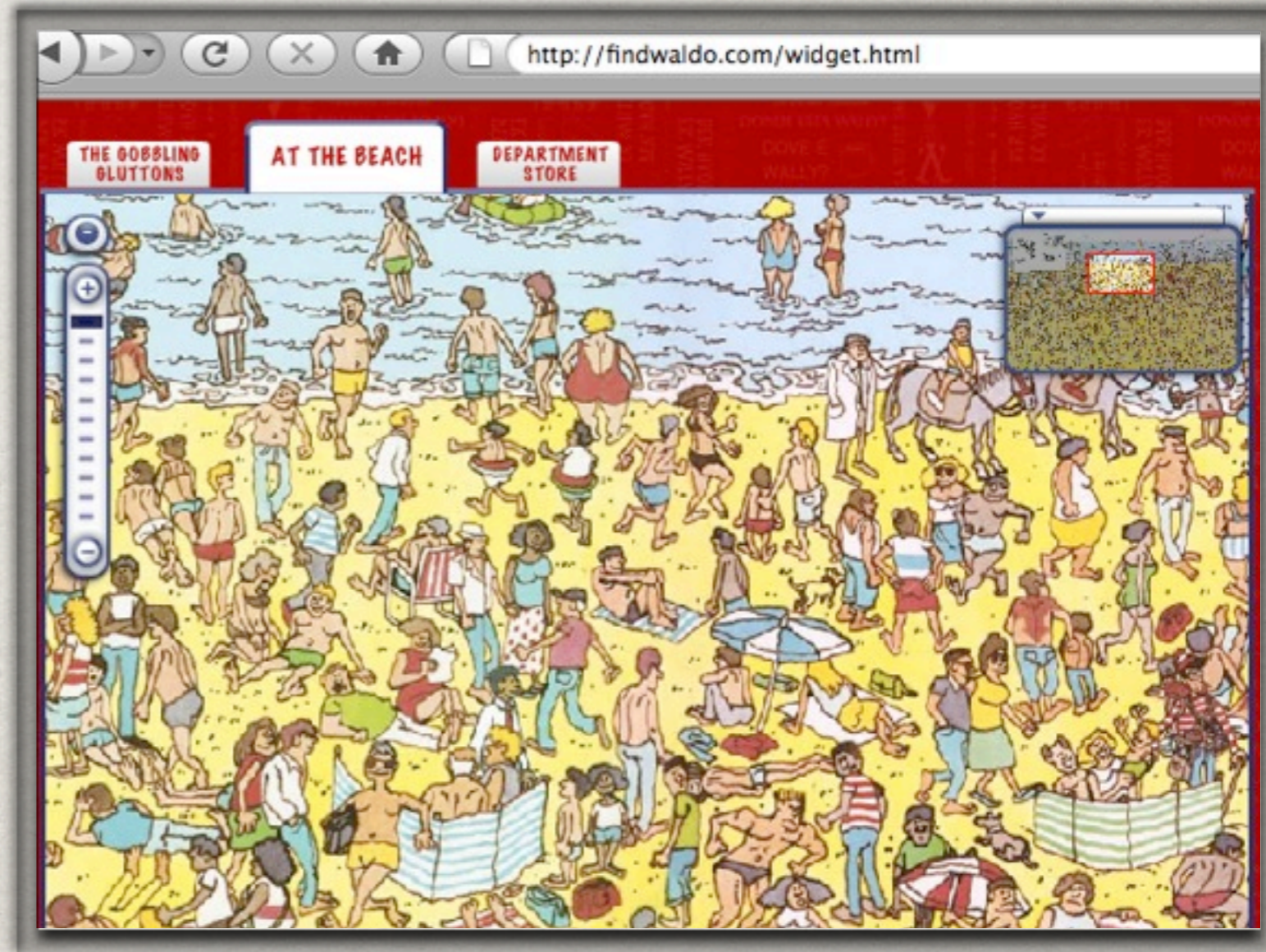
$$-\log p(x)$$

☼ 2) How unsure are you about some aspect of the world (e.g. “is it going to rain?”)?

$$-\sum_{\text{Possibilities}} p(x) \log p(x)$$

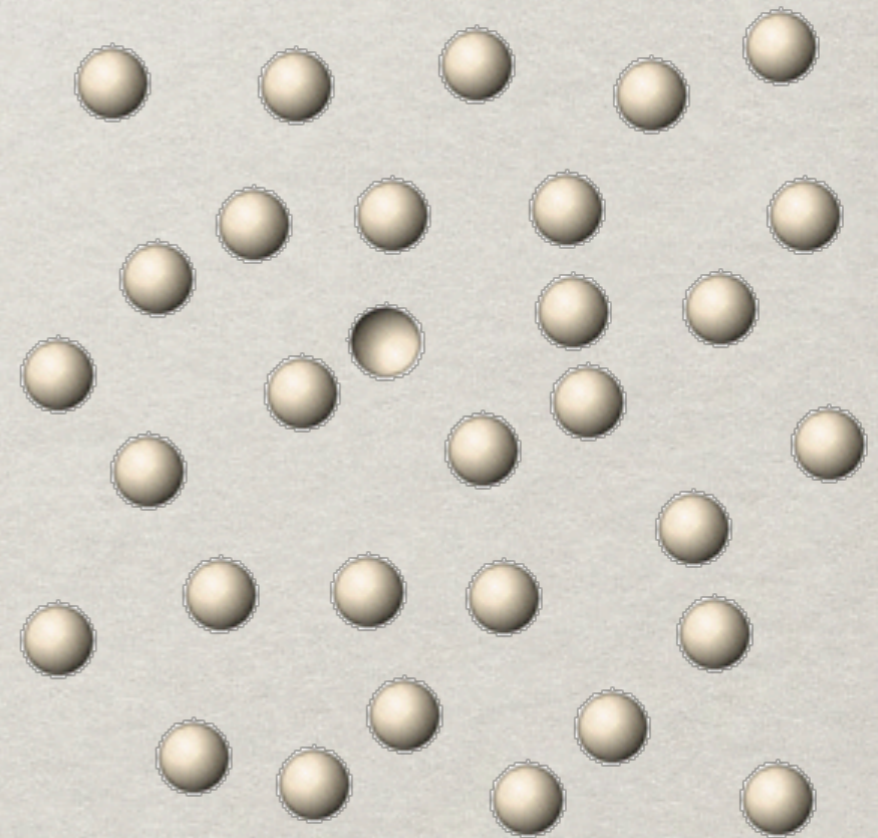
# HOW HUMANS SEARCH SCENES

- ☼ People don't closely examine every inch of the world.
- ☼ Eye-movements are tuned to optimally gather information (Sense 1 and Sense 2).
- ☼ These different notions of information led to very different models:
  - ☼ 1) Visual Saliency
  - ☼ 2) Digital Retina



# VISUAL SALIENCY

- ✻ Salient objects “pop out” of visual scenes.
  - ✻ Simple preprocessing step directs computational resources.
  - ✻ Rare (improbable) image features are more salient than common (probable ones)
  - ✻ Improbable events carry more *information* (Sense 1).
- ✻ We developed an efficient way to model the statistics of a video stream, and analyze it for salient “pop out”.



# A PROMISING FRAMEWORK

- ✱ A common framework is shared by several authors.
- ✱ Claim: The **goal** of eye-movements is to **find visual targets**.
- ✱ Approach: Attend to regions  $x$  of the visual plane which **contain visual targets with high probability**.
  - ✱ In open-ended tasks, drop class-specific terms.

$$\begin{aligned} \text{Saliency}(x) &= \log[p(C_x = 1 | \text{ImgFeats}_x)] \\ &= \log[p(\text{ImgFeats}_x | C_x = 1)] + \log[p(C_x = 1)] - \log[p(\text{ImgFeats}_x)] \end{aligned}$$

← Estimate using current  
image histogram  
(Torralba et al.)

← Estimate using local  
region histogram  
(Bruce & Tsotsos)

↓ Estimate using natural  
image histograms  
(Zhang et al.)\*

\*Best suited to real-time implementation

# ZHANG'S SUN MODEL

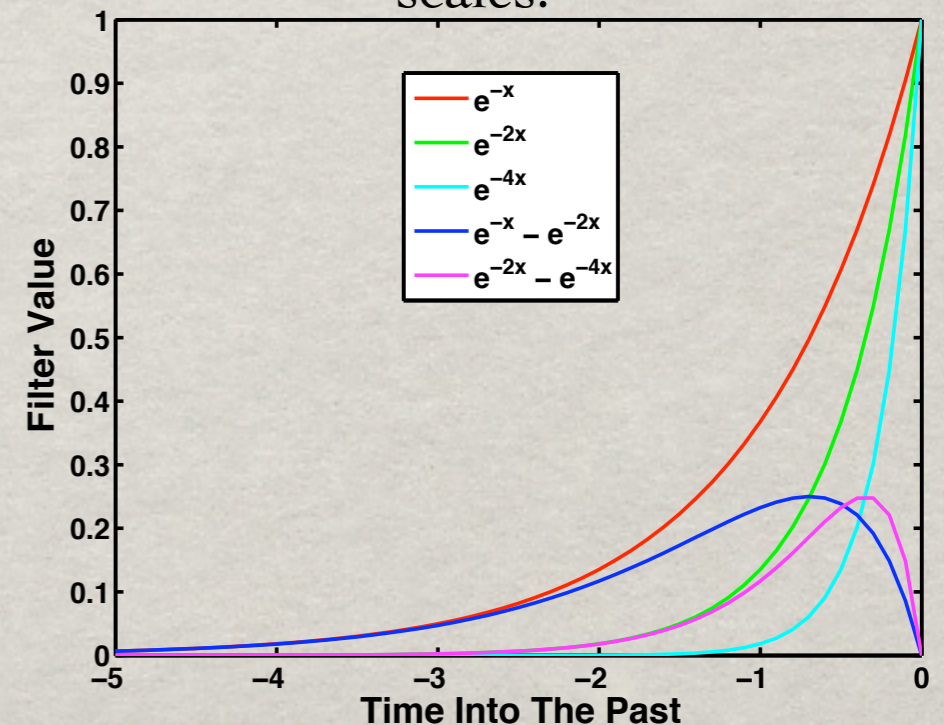
- ☼ Zhang et al. created the “Saliency Using Natural-statistics” model of visual saliency.

$$\text{Saliency}(x) = -\log[p(\text{ImgFeats}_x)]$$

**Space:** Difference Of Gaussians filters at increasing spatial scales.



**Time:** Difference Of Exponentials filters at increasing temporal scales.



**Probability:** Generalized Gaussian

$$p(\text{ImgFeats}_x) = \prod_i C_i \exp(-|\text{ImgFeats}_x^i / \sigma_i|^{\theta_i})$$

\*Parameters  $\sigma$  and  $\theta$  are estimated from Image Features in natural images.

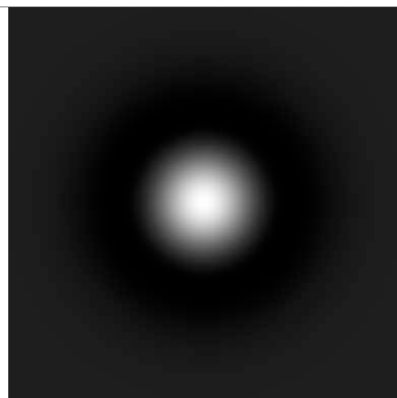
# SUN ALGORITHM SKETCH

1. Grab a new video frame
2. Filter frame with  $N$  Difference-of-Gaussian filters at increasing spatial scales.
3. Integrate each  $DoG$  filter with  $M+1$  previous Exponential filters at increasing time-scales.  
 $[\tau_j/(1+\tau_j) DoG_k + 1/(1+\tau_j) OldExponential_{kj}]$
4. Compute  $NM$  Difference-of-Exponential temporal filters.
5. Compute  $-\log p(DoE)$  for each pixel  $x$  of each  $DoE$  filter  $i$ :  
 $-\log p(DoE_i) = |DoE_i / \sigma_i|^{\theta_i}$ ; for  $\theta_i$  and  $\sigma_i$  fit to spatiotemporal scale in natural images.
6. Sum all  $NM$   $-\log p(DoE)$  to get salience for each pixel  $x$ .



# EFFICIENT APPROXIMATION

- ✱ Our goal is a much-faster-than-real-time algorithm.
- ✱ We achieve this with two approximations.
  1. Instead of Difference-of-Gaussian spatial filters, we use Difference-of-Box Haar-style filters. [Efficient convolution]
  2. Instead of Generalized Gaussian probability model, we use Laplacian probability model. [Efficient inference]

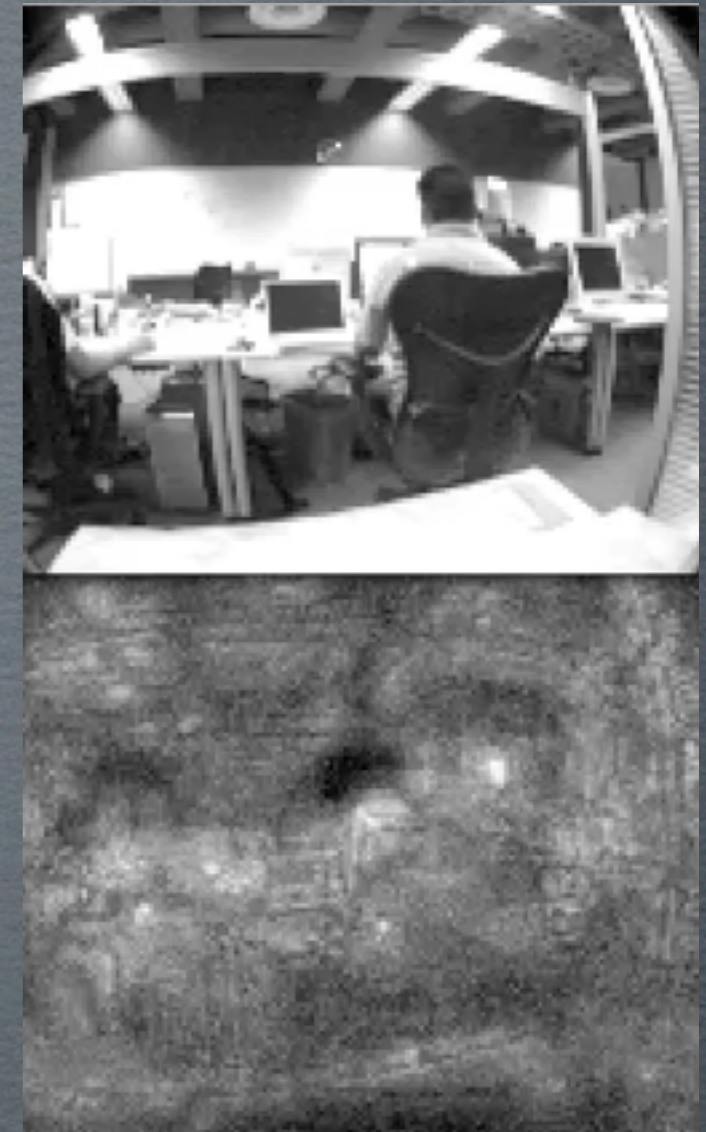


\*Butko, Zhang, Cottrell, Movellan



Online: Camera Control

Offline: Video Analysis



TWO EXAMPLES

# “POP OUT” HELPS TRACK PEOPLE



**Salience Tracking Condition**

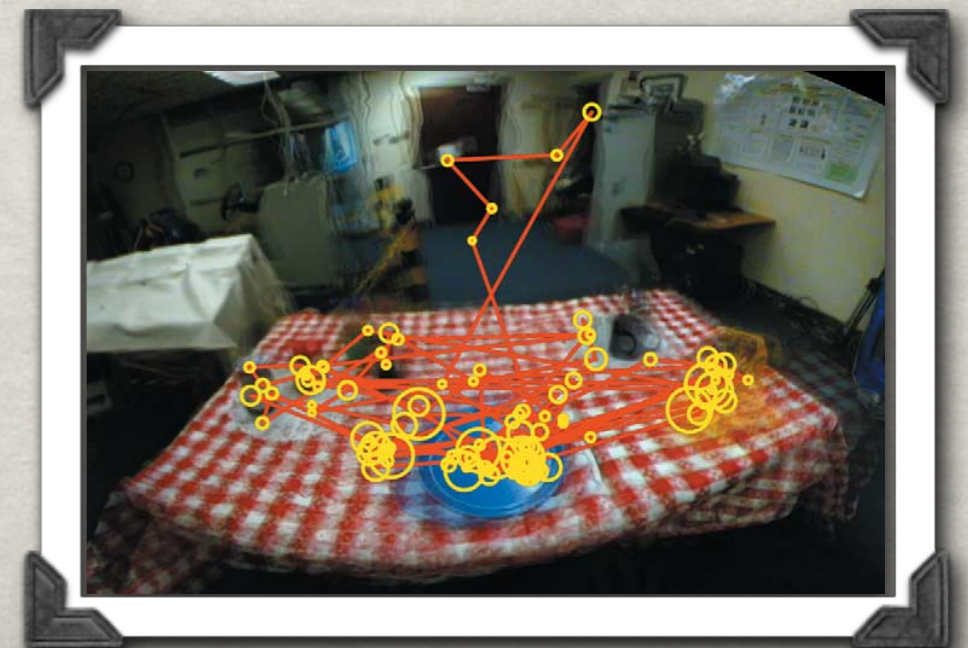
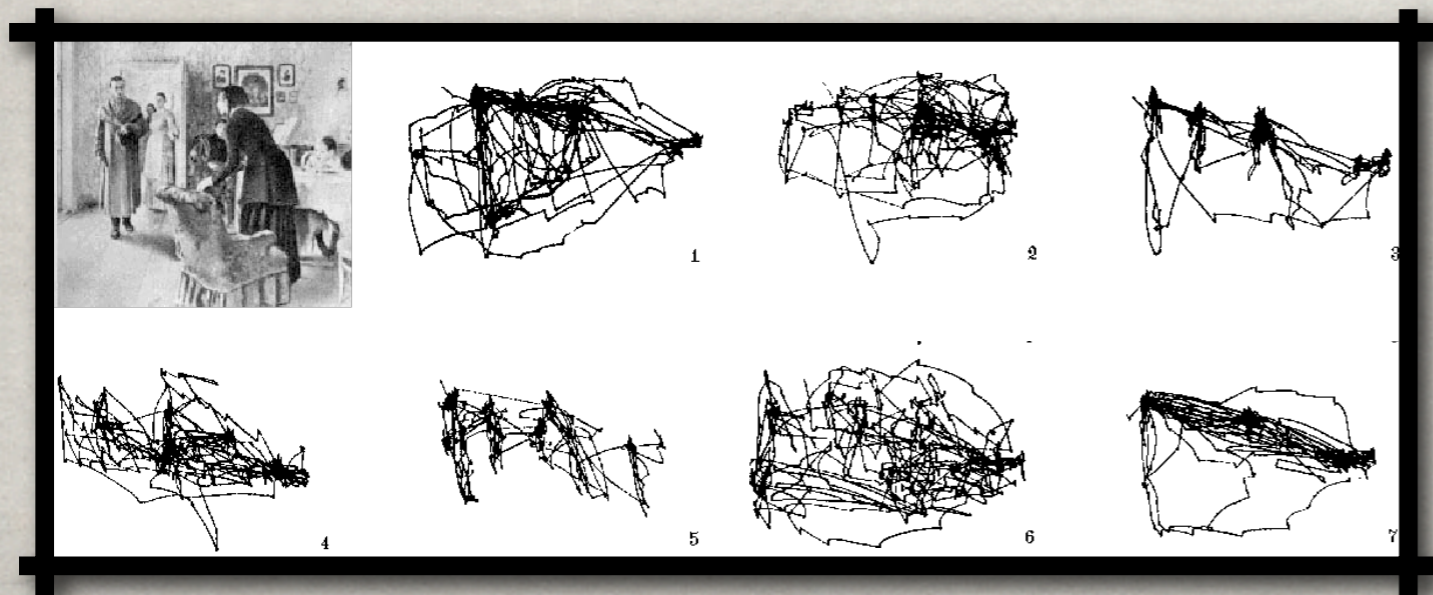


**Playback Condition**



# TASK DIRECTED LOOKING BEHAVIOR

- ✿ Visual Popout can be useful for robots, and it seems to be important in people, but it can't account for task-specific looking behavior.
- ✿ It has long been known that where people look depends on what information they are trying to gather [Yarbus 1967]
- ✿ Current studies have difficulty making quantitative claims: "Fixations are tightly linked in time to the evolution task. Very few irrelevant regions are fixated." [Hayhoe & Ballard 2005]



# A PROMISING FRAMEWORK

- ✱ A common framework is shared by several authors.
- ✱ Claim: The **goal** of eye-movements is to **find visual targets**.
- ✱ Approach: Attend to regions  $x$  of the visual plane which **contain visual targets with high probability**.
  - ✱ In open-ended tasks, drop class-specific terms.

$$\begin{aligned} \text{Saliency}(x) &= \log[p(C_x = 1 | \text{ImgFeats}_x)] \\ &= \log[p(\text{ImgFeats}_x | C_x = 1)] + \log[p(C_x = 1)] - \log[p(\text{ImgFeats}_x)] \end{aligned}$$

Object Appearance  
Information

Location Prior

Image Channel  
Information

“Mutual Information” between object presence and image features.

\*Tong, Kanan, Cottrell

# QUALITATIVE RESULTS (MUG SEARCH)

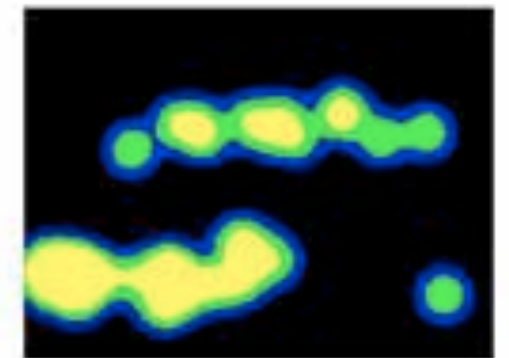
- Where we disagree the most with Torralba et al. (2006)

Gist

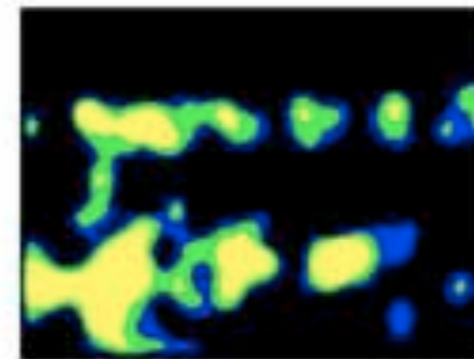
SUN



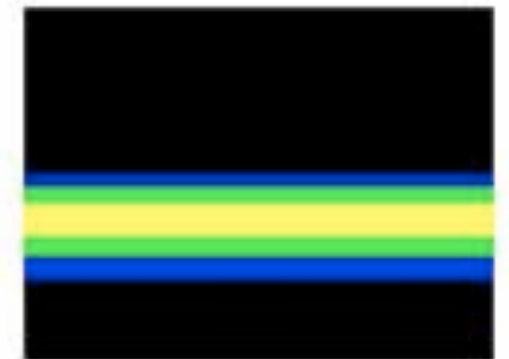
Targets: mugs



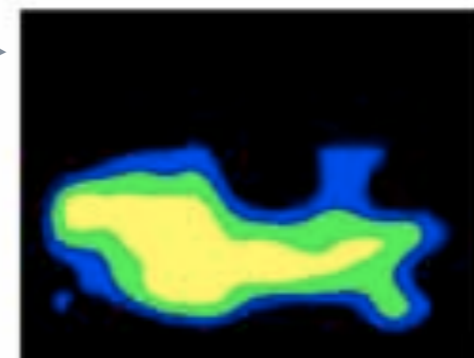
Subject Consistency



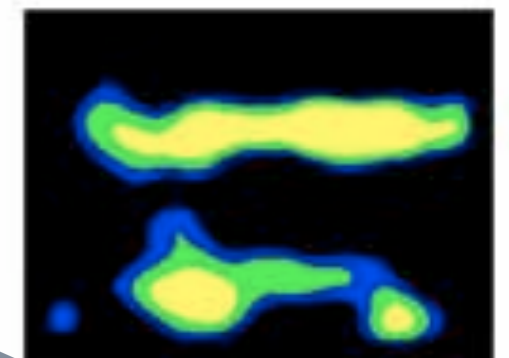
Bottom-Up ( $1/p(F|G)$ )



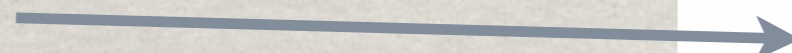
Contextual Modulation ( $p(L|C,G)$ )



Contextual Guidance ( $p(L|C,G)/p(F|G)$ )



Appearance ( $p(C|F)$ )



# QUALITATIVE RESULTS (PICTURE SEARCH)

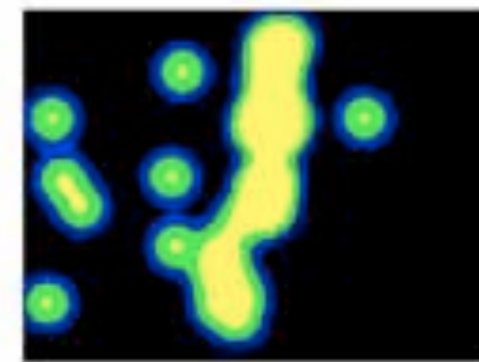
- Where we disagree the most with Torralba et al. (2006)

Gist

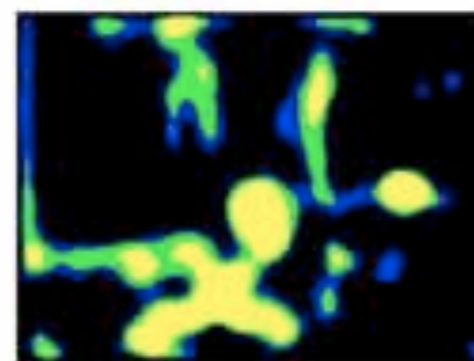
SUN



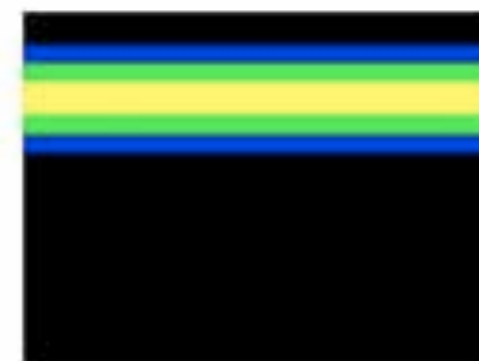
Targets: paintings



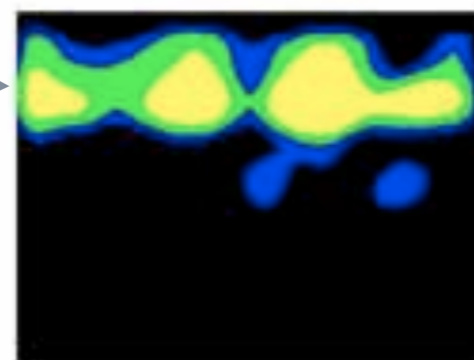
Subject Consistency



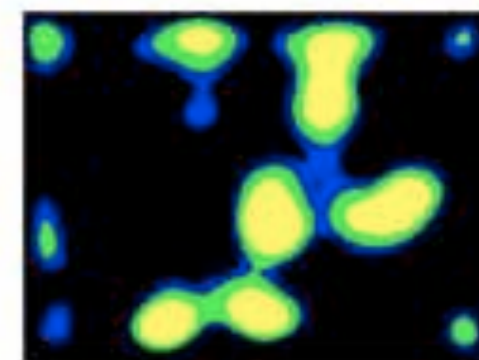
Bottom-Up ( $1/p(F|G)$ )



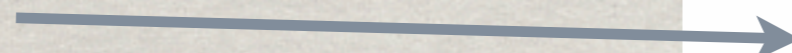
Contextual Modulation ( $p(L|C,G)$ )



Contextual Guidance ( $p(L|C,G)/p(F|G)$ )



Appearance ( $p(C|F)$ )

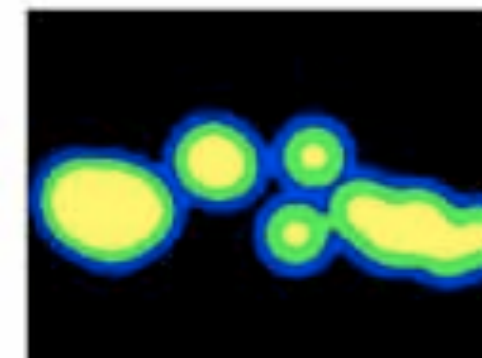


# QUALITATIVE RESULTS (PEOPLE SEARCH)

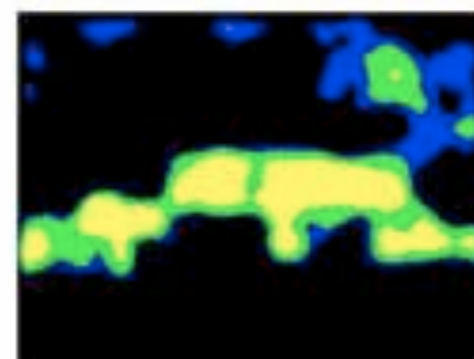
- Where we agree the most with Torralba et al. (2006)



Targets: people



Subject Consistency

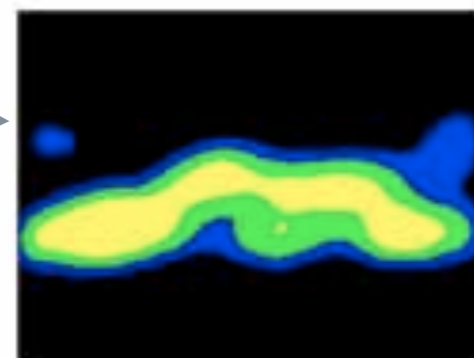
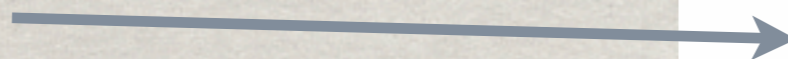


Bottom-Up ( $1/p(F|G)$ )

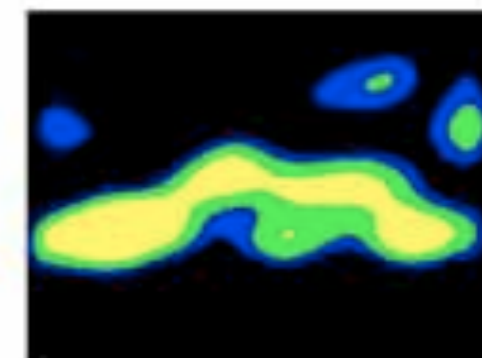


Contextual Modulation ( $p(L|C,G)$ )

Gist



Contextual Guidance ( $p(L|C,G)/p(F|G)$ )



Appearance ( $p(C|F)$ )

SUN



# QUALITATIVE RESULTS (PAINTING SEARCH)

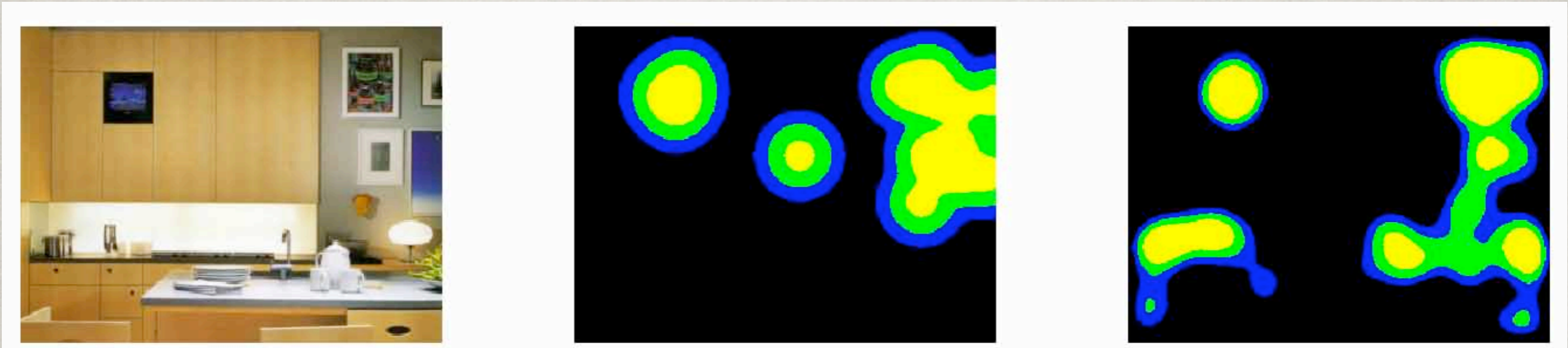


Image Humans SUN

- ☼ This is an example where SUN and humans make the same mistake due to the similar appearance of TV's and pictures (the black square in the upper left is a TV!).



# MATHEMATICAL FORMULATION OF INFORMATION

☼ “Information” has two (related) mathematical meanings:

☼ 1) How much did you expect something you experience (“it is going to rain”)?

$$-\log p(x)$$

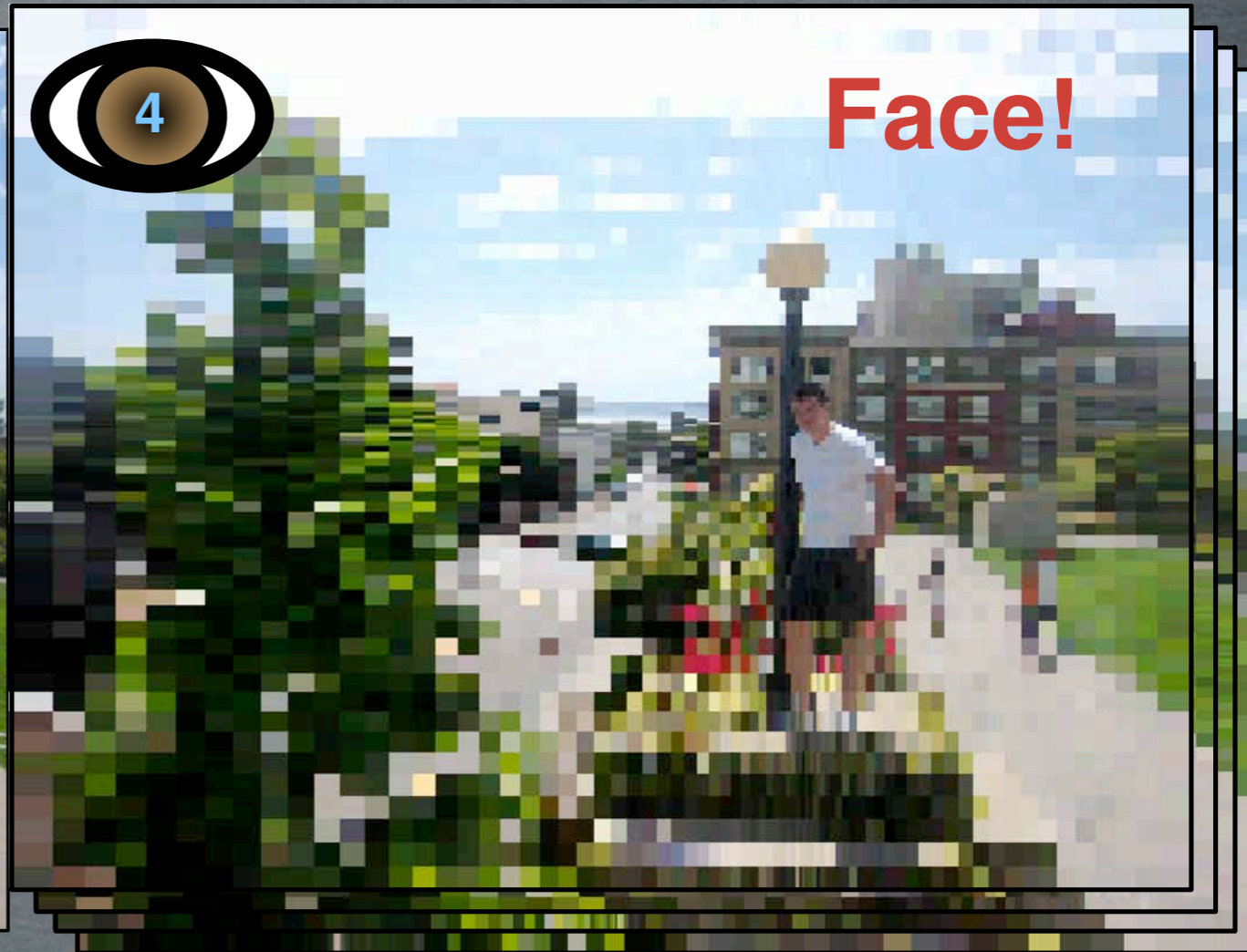
☼ 2) How unsure are you about some aspect of the world (e.g. “is it going to rain?”)?

$$-\sum_{\text{Possibilities}} p(x) \log p(x)$$

# GATHERING INFORMATION

$$- \sum_{\text{Possibilities}} p(x) \log p(x)$$

- ✱ Question:
  - ✱ “Where is a face?”
- ✱ Possibilities:
  - ✱ Top-left, Middle, Bottom-right, etc...
  - ✱ Or, nowhere.
- ✱ Information (above) says how much information we have left to gather about the face location.
  - ✱ Once we have gathered the maximum amount of information, we will know where the face is.



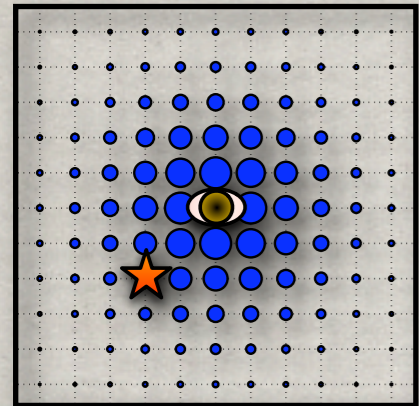
# SEARCHING FOR FACES

# MODELING THE RETINA

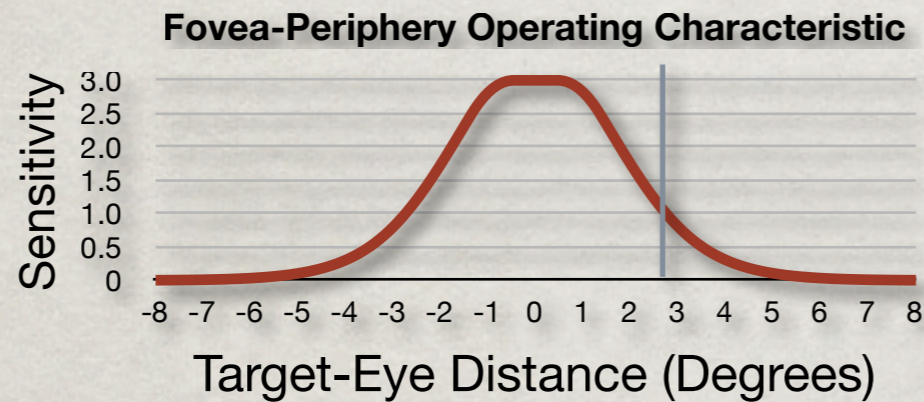
[ADAPTED FROM NAJEMNIK & GEISLER 2005]

Signal+Noise [N(0,1)]  
(Observation)

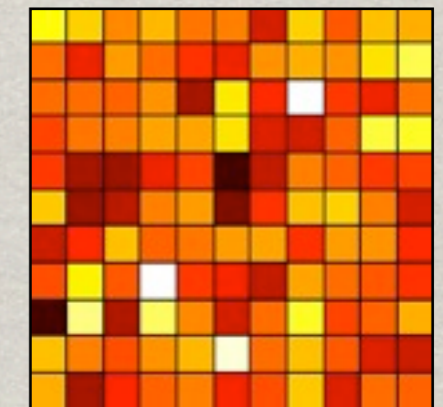
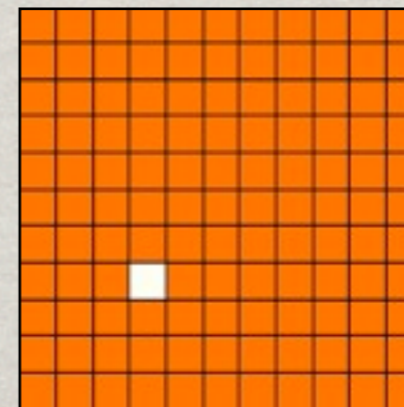
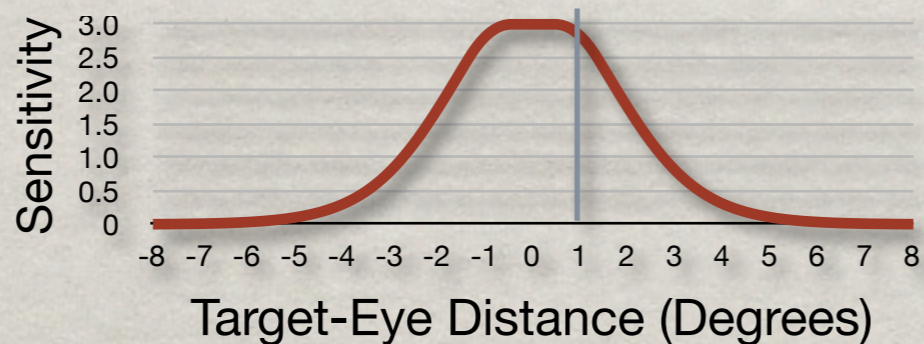
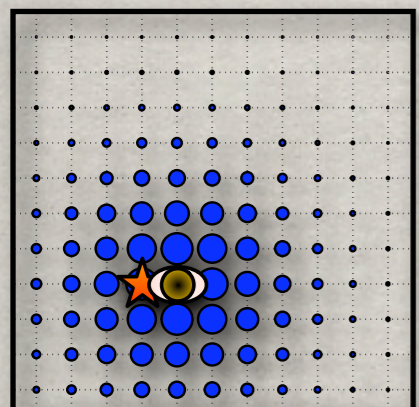
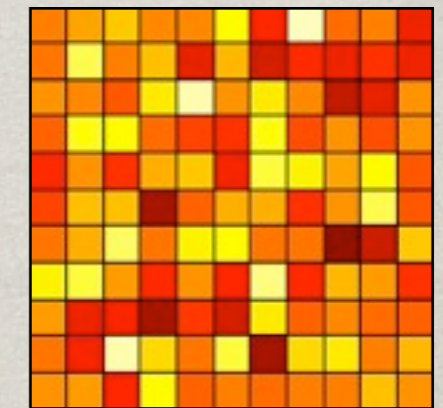
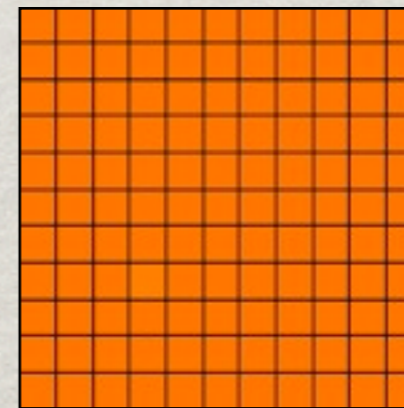
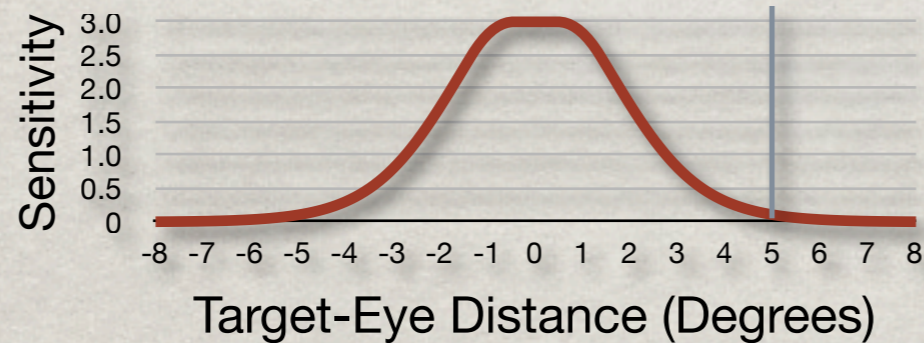
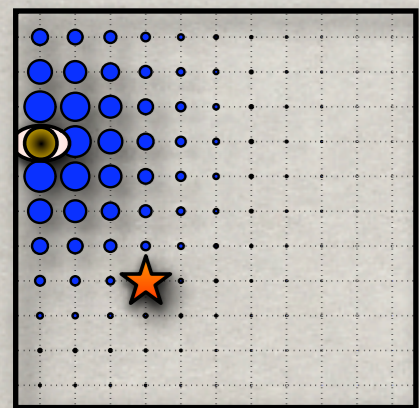
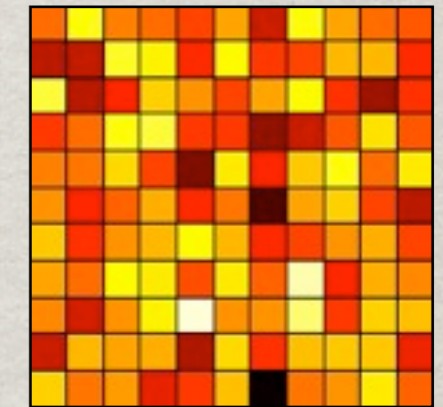
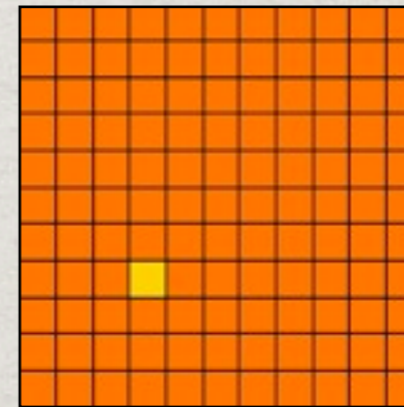
State / Action



Target Signal Strength



Signal

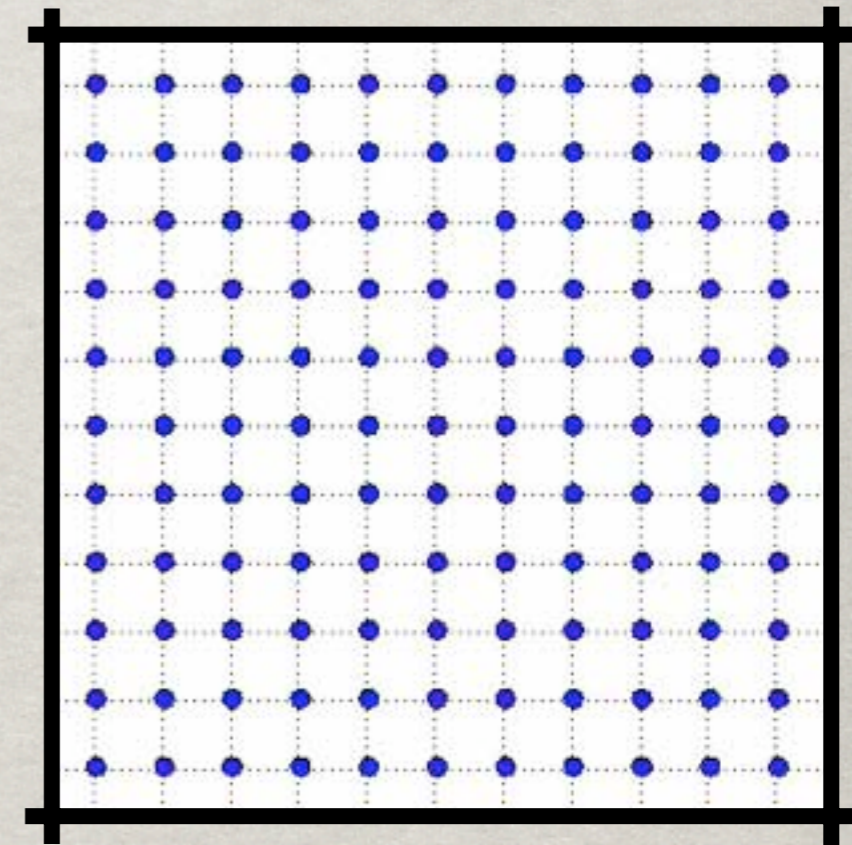


\*Apply Infomax principle to learn optimal eye-movement behavior!

# INTEGRATING OBSERVATIONS

- ✱ Characterized the noise properties of the sensory system.
- ✱ The POMDP framework specifies how to *infer* the likelihood that the target is at each location:

$$\begin{aligned}
 p(S = i | A_{1:t}, \vec{O}_{1:t}) &\propto p(\vec{O}_t = \vec{o} | S = i, A_t) p(S = i | \vec{O}_{1:t-1}, A_{1:t-1}) \\
 p(\vec{O}_t = \vec{o} | S = i, A_t = k) &= \prod_{j=1}^N p(o_j | S = i, A_t = k) \\
 &= 1/\sqrt{2\pi} \exp(-(o_i - d_{i,k})^2/2) \prod_{j \neq i} 1/\sqrt{2\pi} \exp(-(o_j)^2/2) \\
 &= \frac{\exp(-(o_i - d_{i,k})^2/2)}{\exp(-(o_i)^2/2)} Z \\
 &= \exp(\alpha_{i,k} d_{i,k}) Z; \quad \alpha \equiv (o_i - d_{i,k})/2 \\
 B_t^i &\propto \exp(\alpha_{i,k} d_{i,k}) B_{t-1}
 \end{aligned}$$



- ✱ I-POMDP Bayesian analysis:
  - ✱ Online Learning
  - ✱ Local update rule



# DIGITAL RETINA IN ACTION

# INFOMAX APPROACH IMPROVES STATE OF THE ART AI

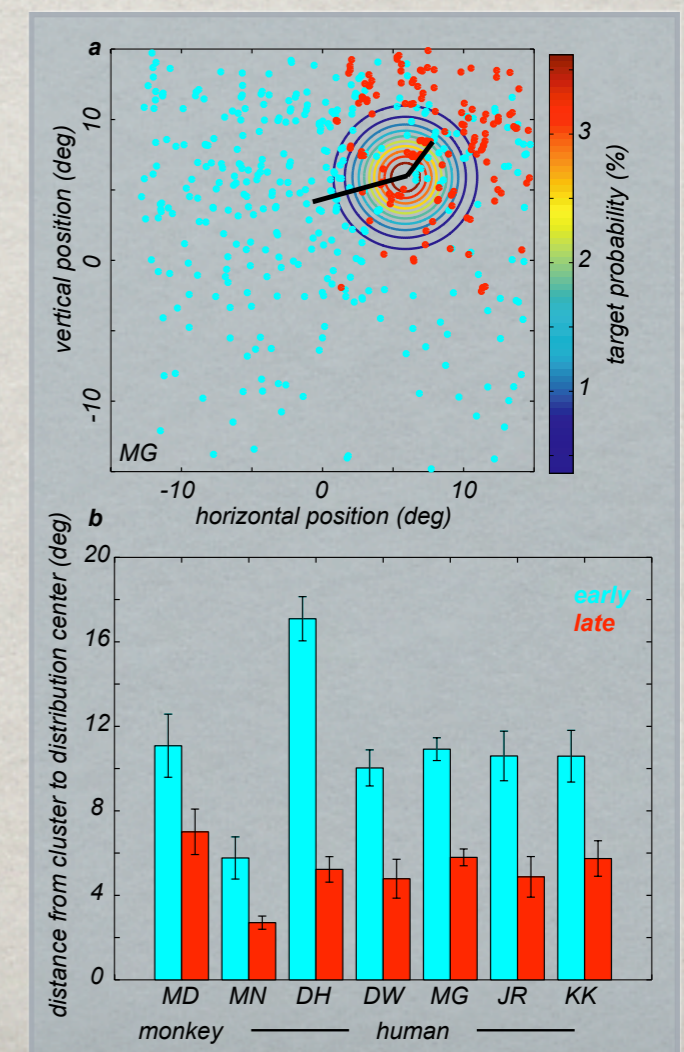
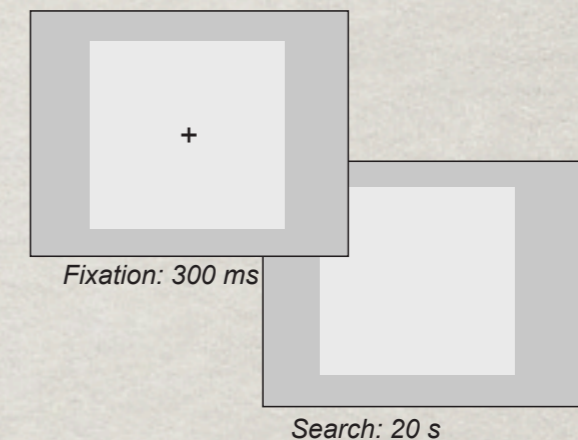
- ✱ Apply digital retina sequentially to a static image Vs. search for faces using a standard face detector.
  - ✱ Achieve two-fold speed increase with minimal loss in accuracy.
- ✱ Optimal Information (Sense 2) gathering.

	Digital Retina	Full Image
Runtime (ms/1000px)	<b>0.57</b>	1.25
Displacement (% Width)	7.6%	<b>6%</b>

\*Butko, Movellan

# IS INFORMATION THE ONLY MODEL FOR EYE-MOVEMENTS?

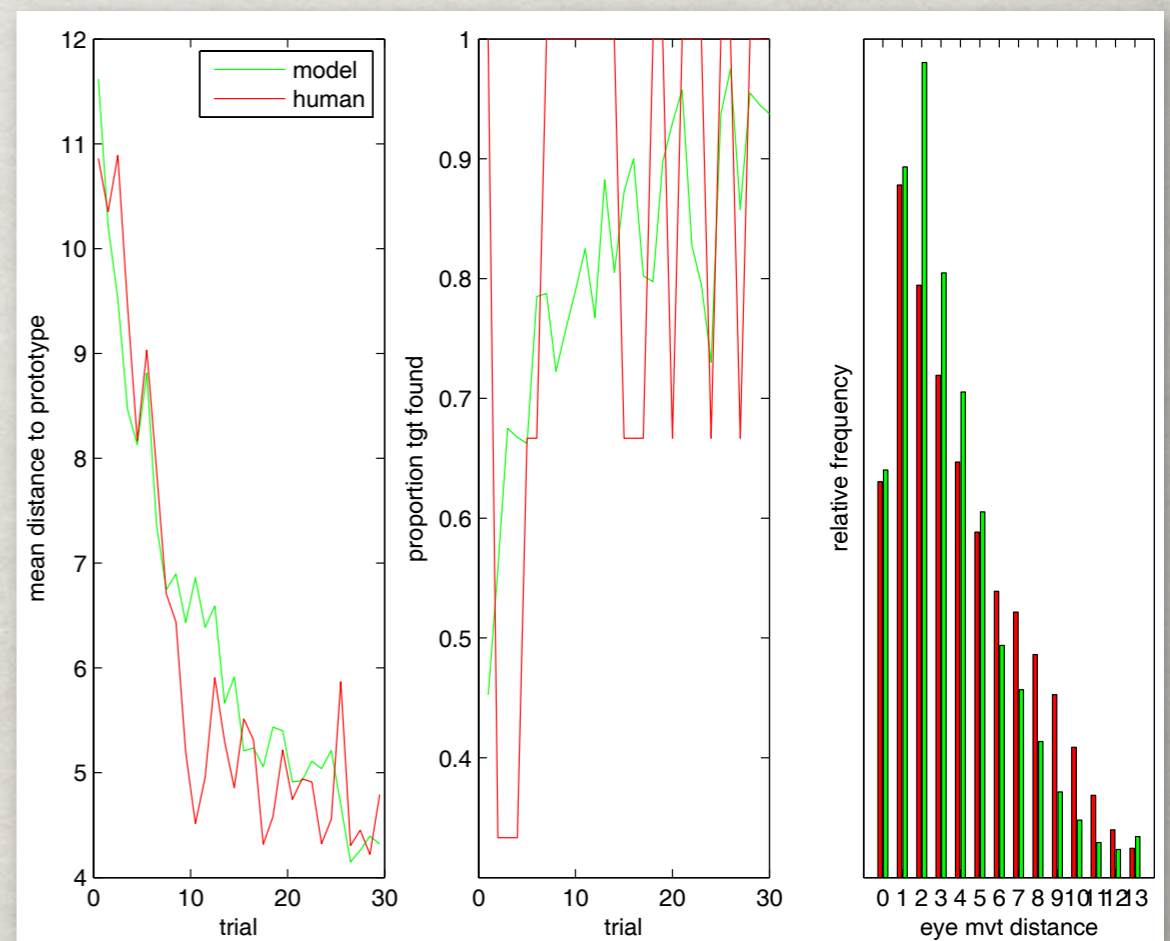
- Given a set of eye-movement data, how should we model it?
- Experiments in “top-down” effects of eye-movement:
  - “Hidden” target, have to look at something invisible to end trial.
  - No “bottom up” visual information to aid eye-movement.
  - After many trials, learn where target is likely to be, move eyes in absence of visual cues.
- How can we model this study?





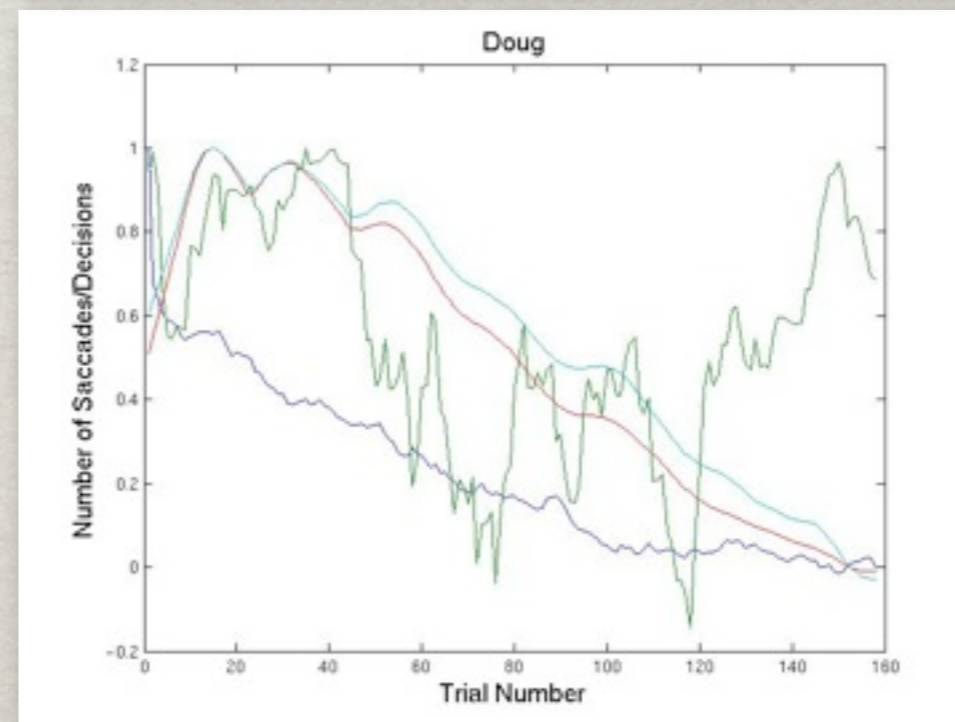
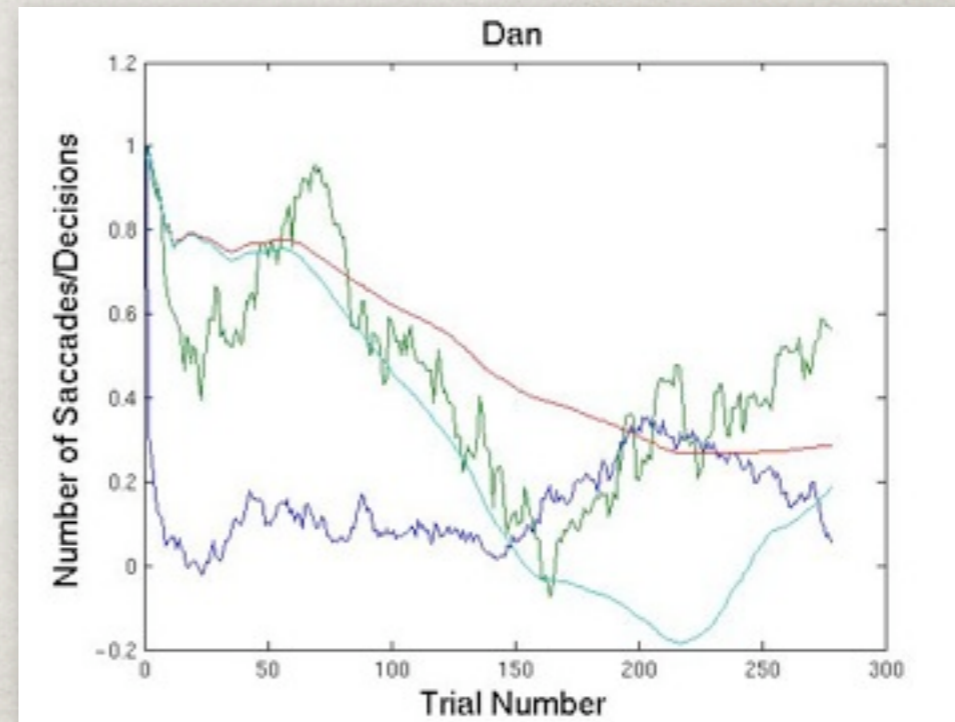
# MODEL 1: SPATIAL Q-LEARNING (RL)

- ✱ A subject doing this task is “rewarded” for finding the search target, “penalized” for moving their eyes too much (wasting energy).
  - ✱ Can leverage reinforcement learning.
    - ✱ Actions that are close (in retrospect, L1 distance) to the target are rewarded more than far away ones.
    - ✱ Movements that are far (L1 distance) from the last fixation are penalized.
  - ✱ Learn reward structure.
- ✱ Choose an eye-movement as a soft-max over the reward of each state.



# MODEL 2: BAYESIAN OPTIMAL OBSERVER

- ✿ Use Bayesian parameter estimation to estimate target location distribution.
- ✿ After each found-target, update estimates of [x-mean, x-variance, y-mean, y-variance].
- ✿ For each trial, sample eye-movements from location distribution.
  - ✿ If not-found, set probability for that location to zero.
  - ✿ Renormalize and resample to generate next fixation.
- ✿ Very similar to Infomax Approach (earlier)



Subject Performance

Model Performance

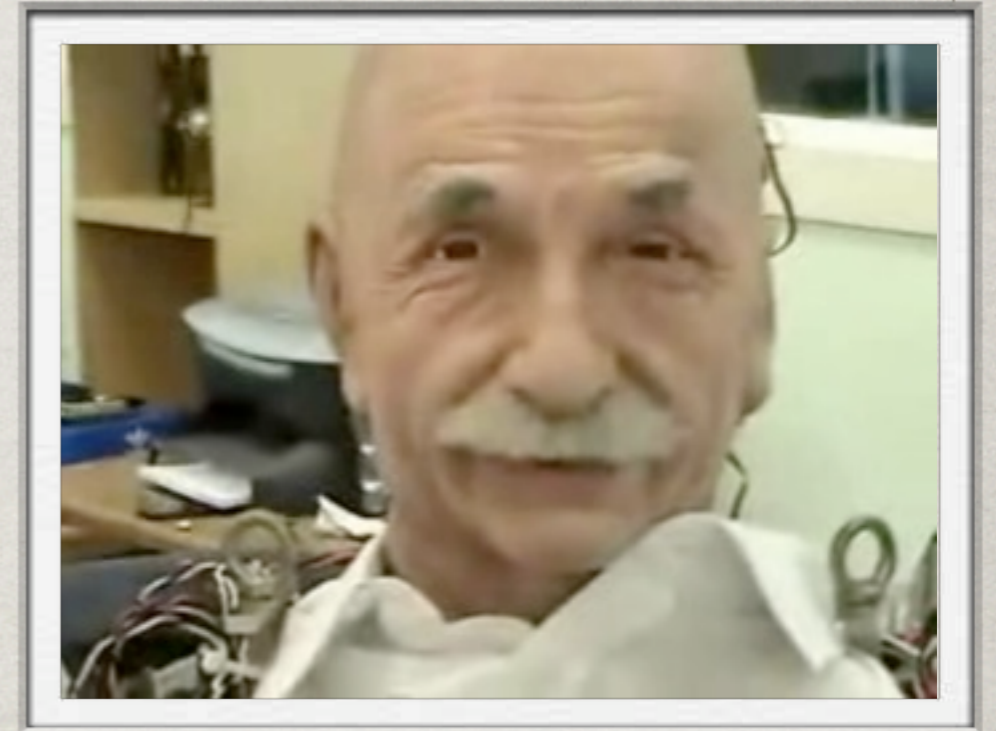
KL Distance Data Distr. to Target Distr.

KL Distance Recent Data to Target Distr.

# NEXT STEPS

- ✿ Einstein Tutor

- ✿ In order to teach, you need to effectively gather information about the mental state of your pupil.
- ✿ Attentive? Confused? Bored?



- ✿ Project One

- ✿ Robotic Platform to simulate developmental processes and learning during the first year of life.

