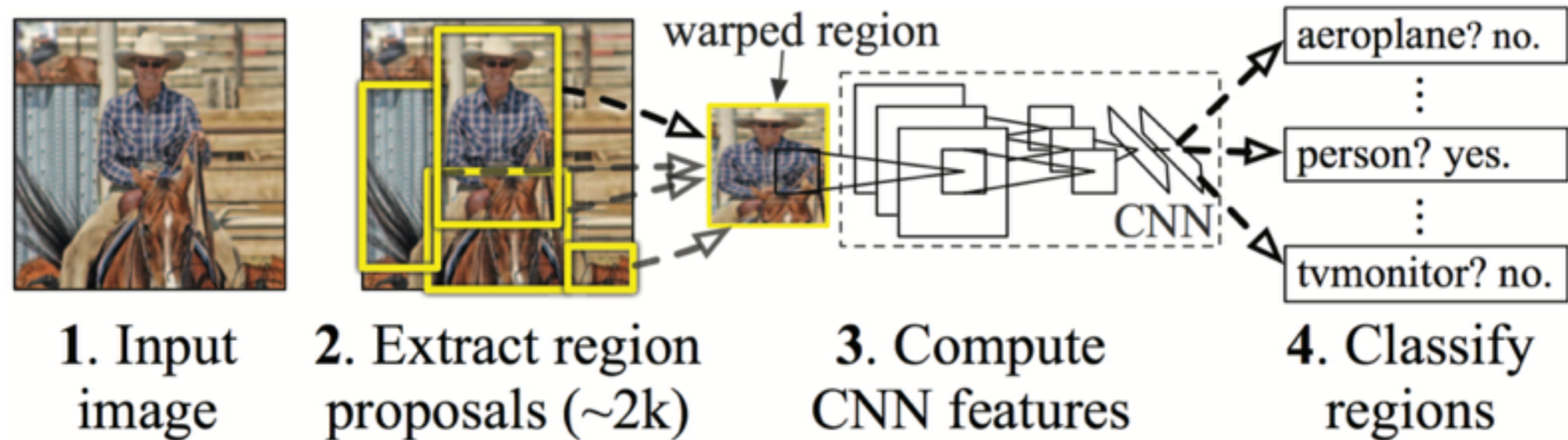


# Visual Attention With Neural Networks

Main Paper: Recurrent Models of Visual Attention

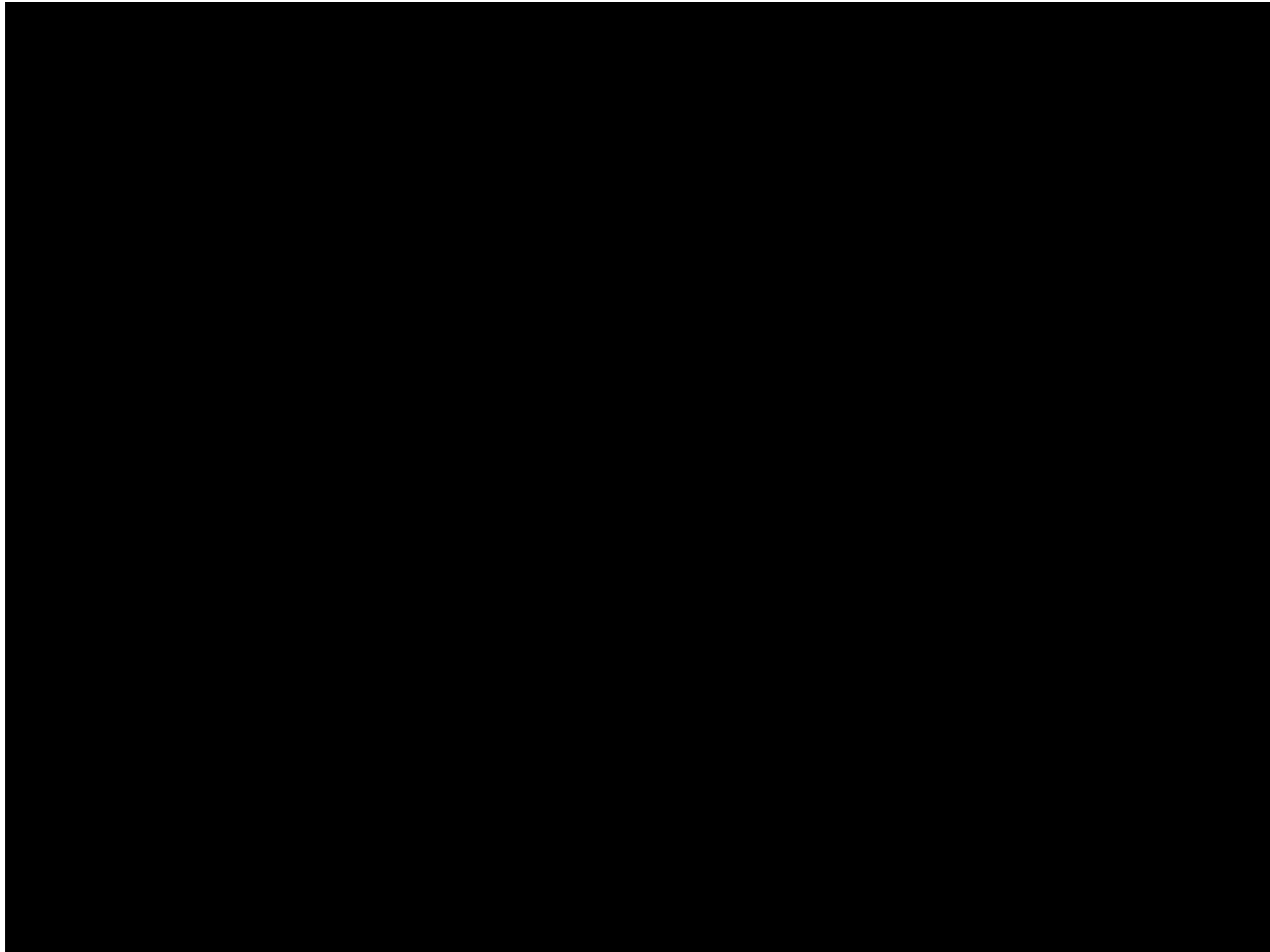
Presentation by Matthew Shepherd

# Full image processing is computationally expensive



Regions are can be selected intelligently but time still scales with the size of the image

Humans focus on specific  
regions in their FOV

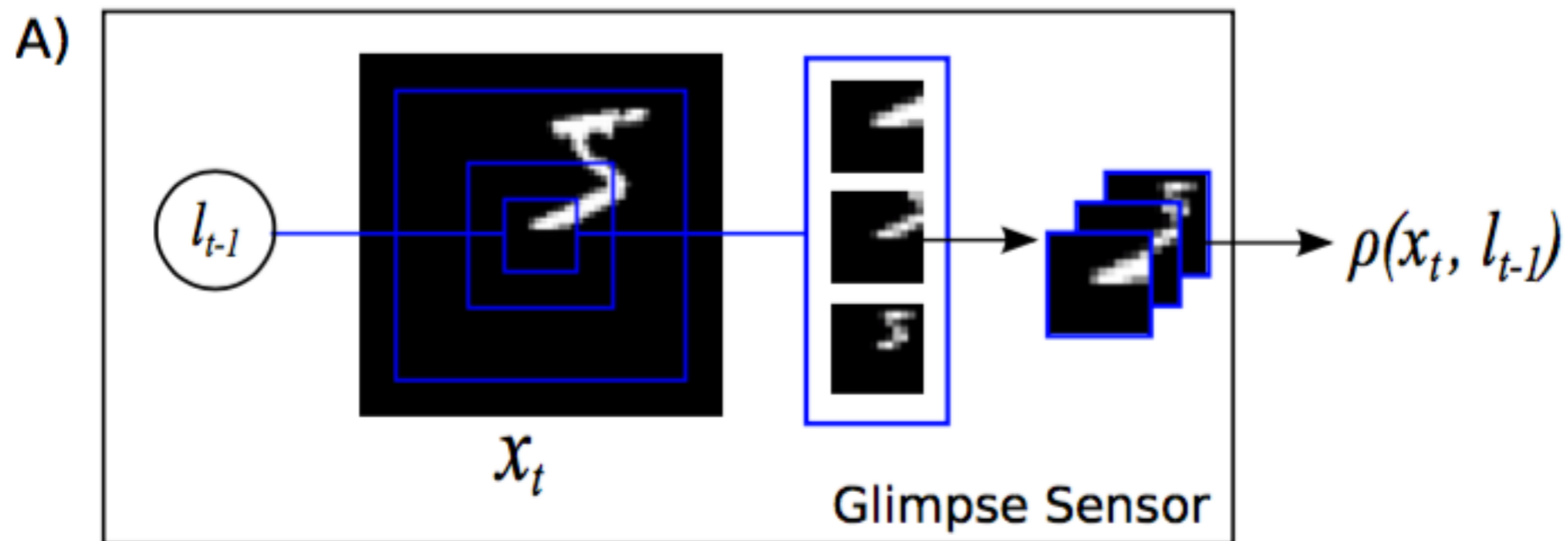


# “Vision as a sequential decision task”

The model sequentially chooses small windows of information on the data

Integrates information from all past windows to make its next decision.

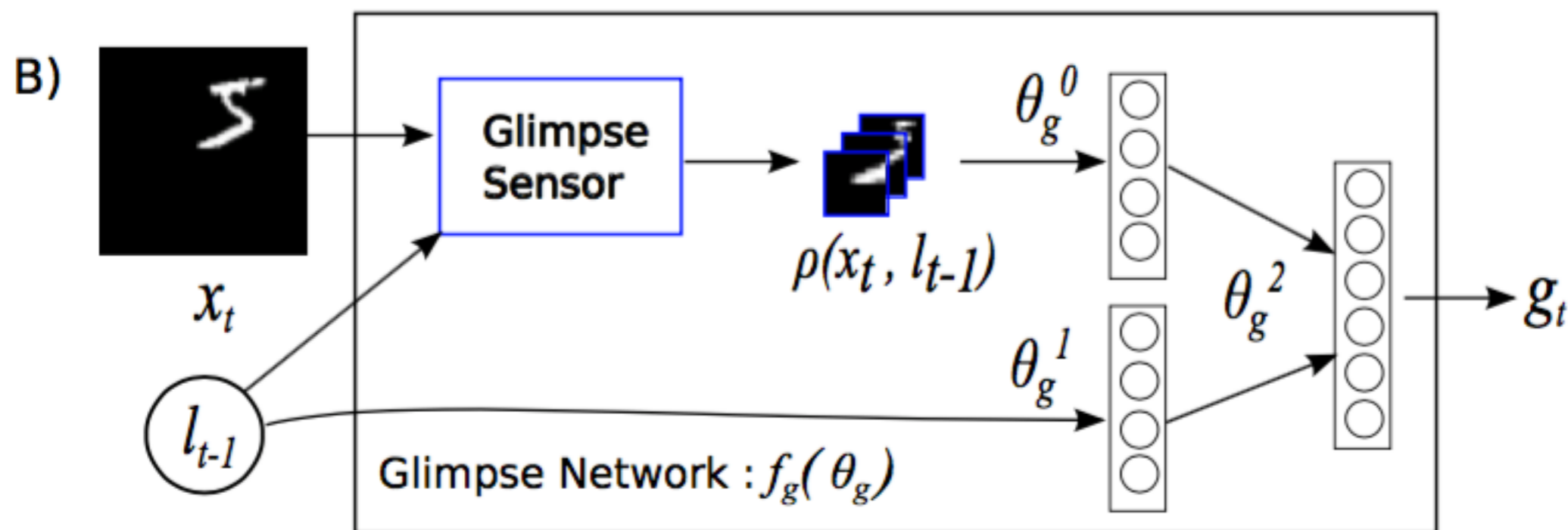
The sensor only provides limited information about the scene,  $X_t$ , focused at a location,  $l_{t-1}$



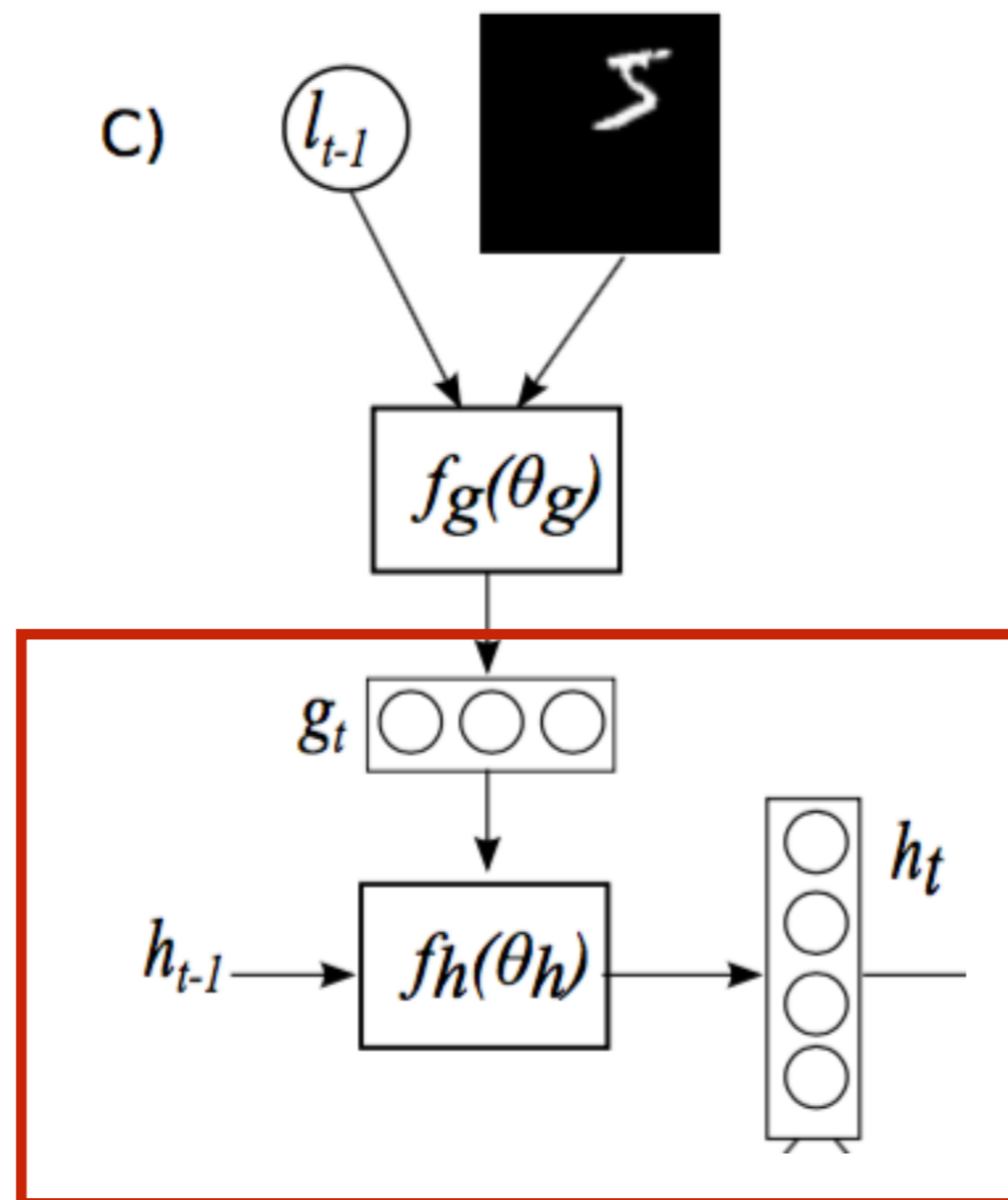
Referred to as a  
“Glimpse”

retina-like  
representation

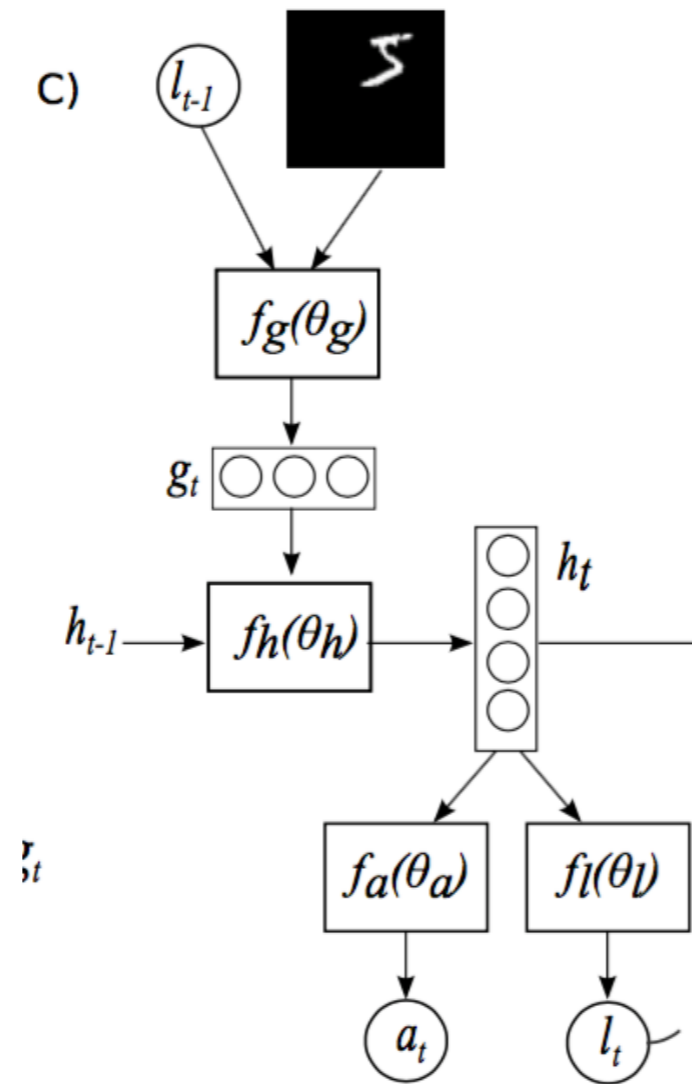
The “glimpse network” incorporates sensor and location information into a single vector



The core network  $f_h(\theta_h)$  incorporates the sensor information as well as past information

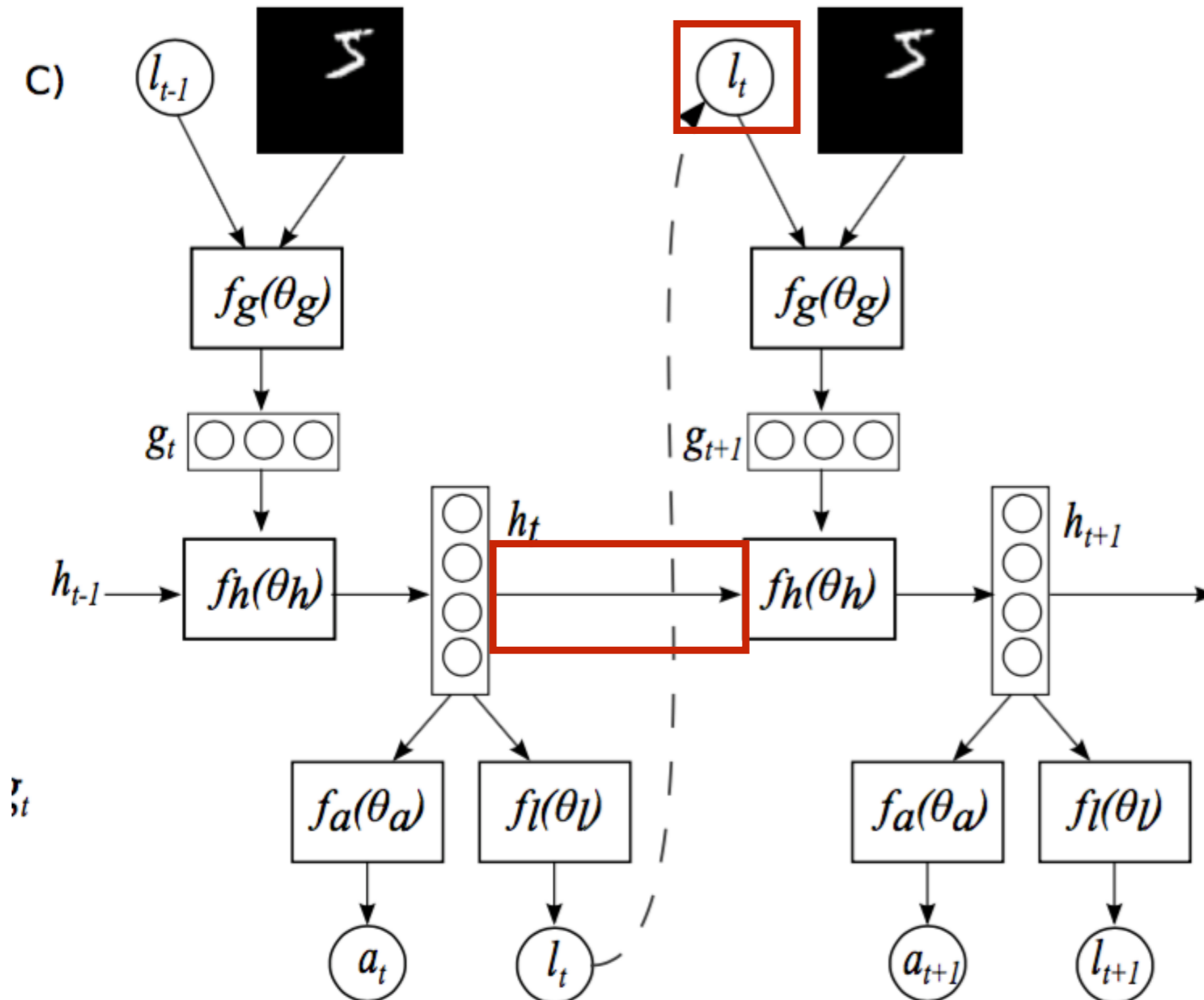


The output of the core network is then used to deploy the sensor and make a classification





# The network is recurrent



# Partially Observed Markov Decision Problem (POMDP)

- Glimpse can be seen as a partial view of the state
- Network must learn a policy

$$\pi(\mathbf{l}_t, \mathbf{a}_t | \mathbf{s}_{1:t}; \theta)$$

- The policy is determined by the NN
- State history is encapsulated by the hidden state of the network

# So how do we train it?

$$J(\theta) = \mathbb{E}_{p(s_{1:T};\theta)} \left[ \sum r_t \right]$$

$$\nabla_{\theta} J = \sum_{t=1}^T \mathbb{E}_{p(s_{1:T};\theta)} [\nabla_{\theta} \log \pi(u_t | s_{1:t}; \theta) R] \approx \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \nabla_{\theta} \log \pi(u_t^i | s_{1:t}^i; \theta) R^i$$

The REINFORCE rule

# Additional training details

Useful for determining  $f_l$  but  $f_a$  can be determined more directly by minimizing cross-entropy loss.

A baseline value,  $b_t$ , is added to the gradient approximation to reduce variance

$$\frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \nabla_{\theta} \log \pi(u_t^i | s_{1:t}^i; \theta) (R_t^i - b_t)$$

# Experiments

# RAM performs well on translated MNIST



Model	Error
FC, 2 layers (64 hidden each)	6.42%
FC, 2 layers (256 hidden each)	2.63%
Convolutional, 2 layers	1.62%
RAM, 4 glimpses, $12 \times 12$ , 3 scales	1.54%
RAM, 6 glimpses, $12 \times 12$ , 3 scales	<b>1.22%</b>
RAM, 8 glimpses, $12 \times 12$ , 3 scales	<b>1.2%</b>

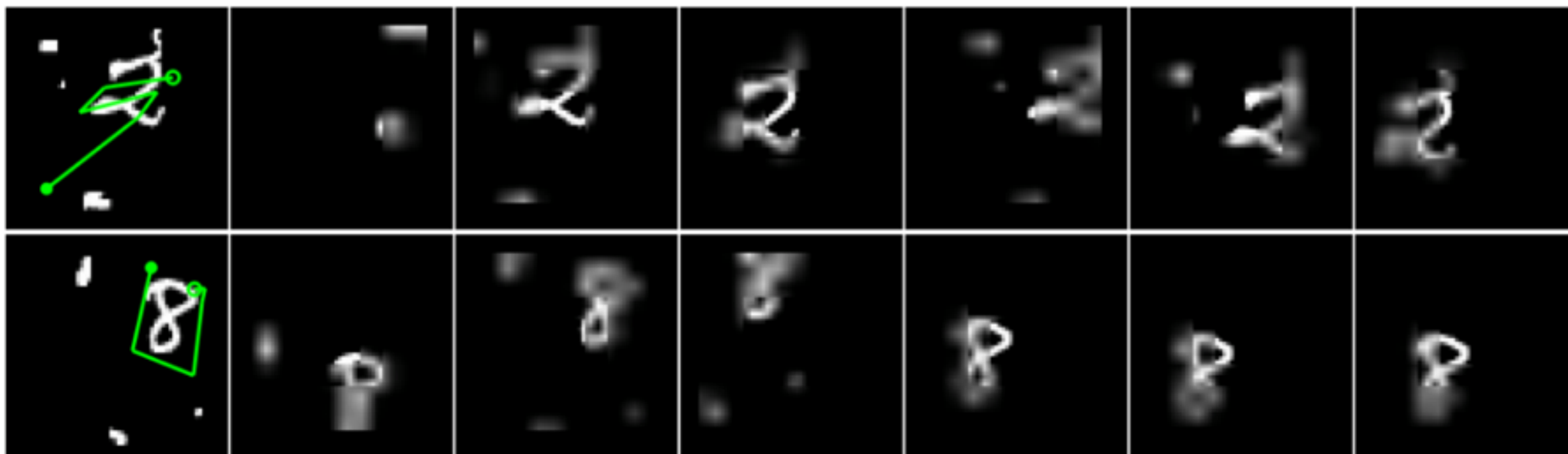
# RAM performs very well on cluttered translated images



**(b) 100x100 Cluttered Translated MNIST**

Model	Error
Convolutional, 2 layers	14.35%
RAM, 4 glimpses, $12 \times 12$ , 4 scales	9.41%
RAM, 6 glimpses, $12 \times 12$ , 4 scales	8.31%
RAM, 8 glimpses, $12 \times 12$ , 4 scales	8.11%
RAM, 8 random glimpses	28.4%

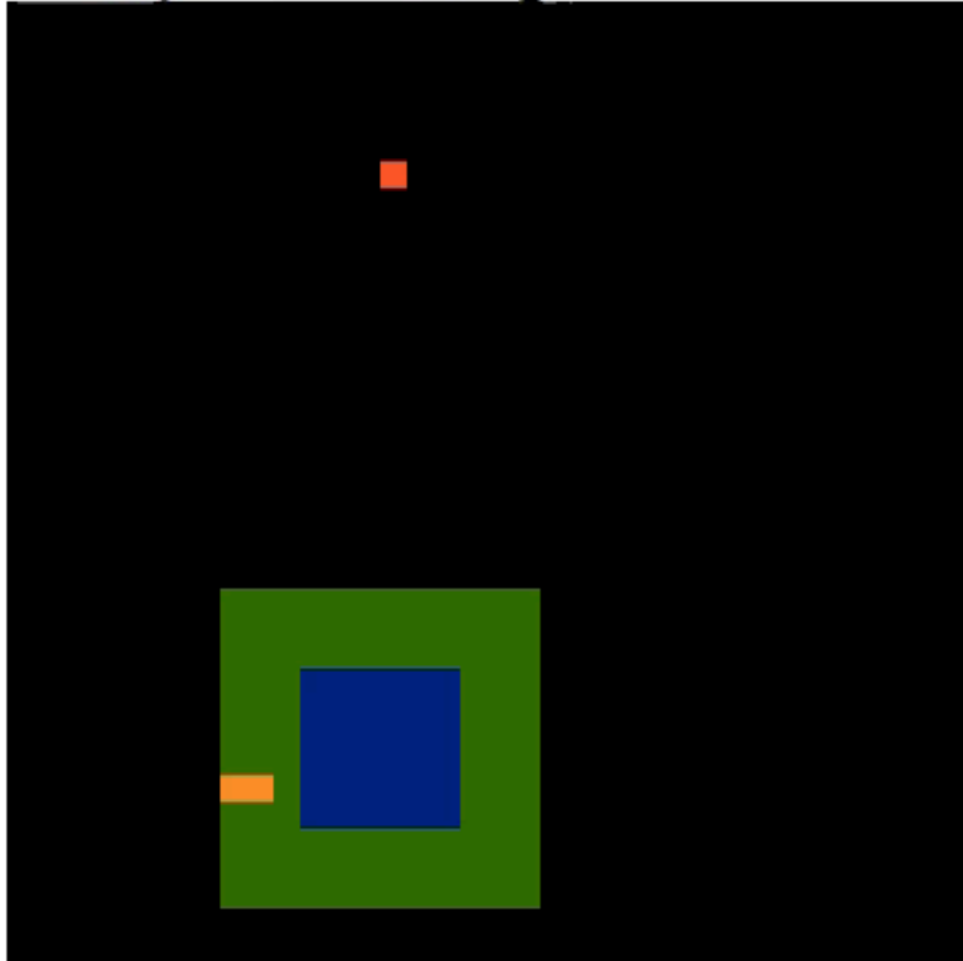
# Meaningful policies are learned



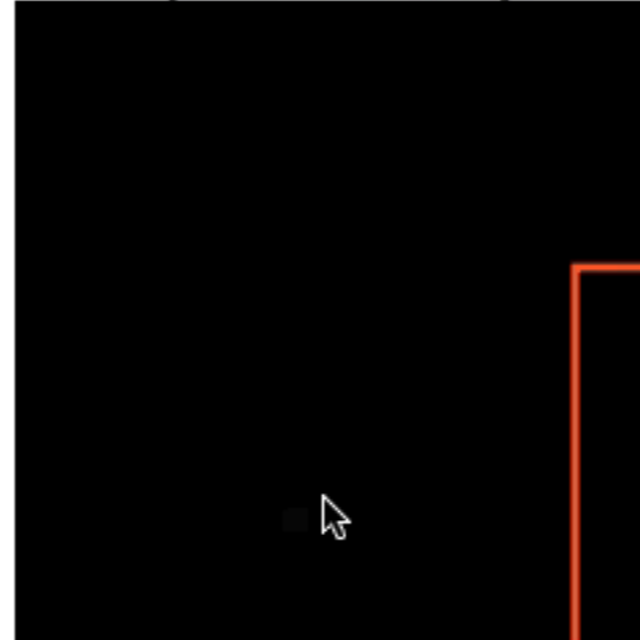


# RAM performs in a dynamic environment

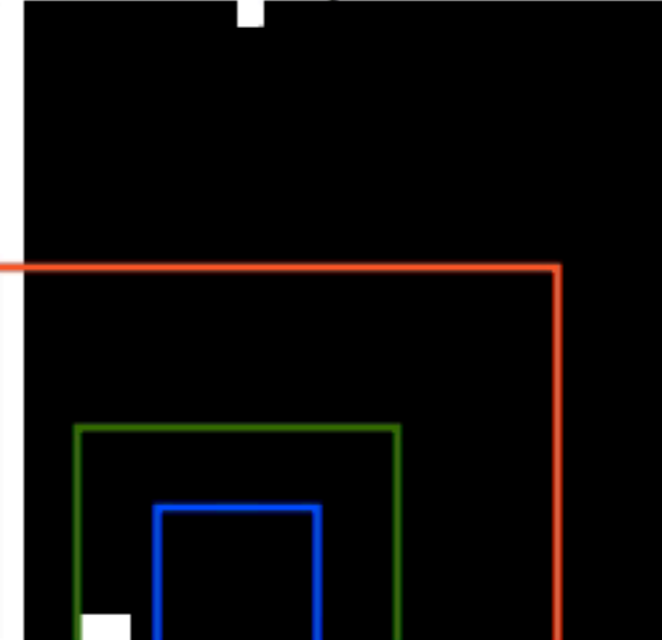
Glimpse Coverage



Glimpse History



Game Play

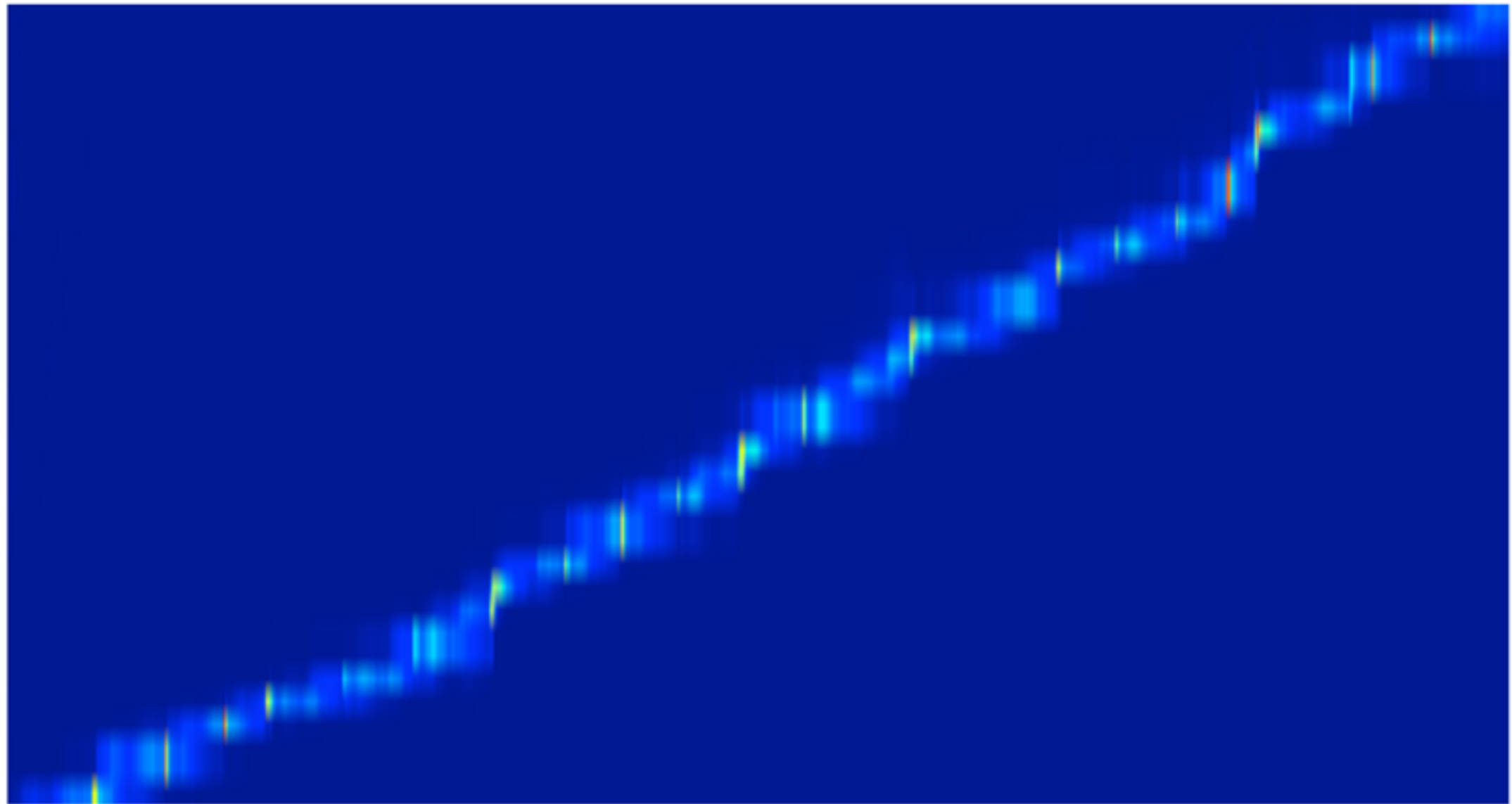


Individual Glimpses



# Other models of attention

Thought that the muster from

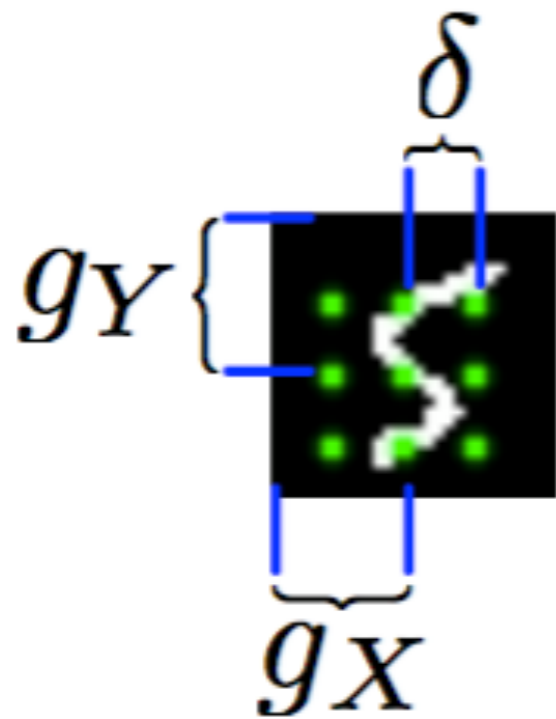


Thought that the muster from

# DRAW: A Recurrent Neural Network For Image Generation

- Combines an attention mechanism with a sequential variational auto-encoder
- Reading *and* writing are now both sequential tasks

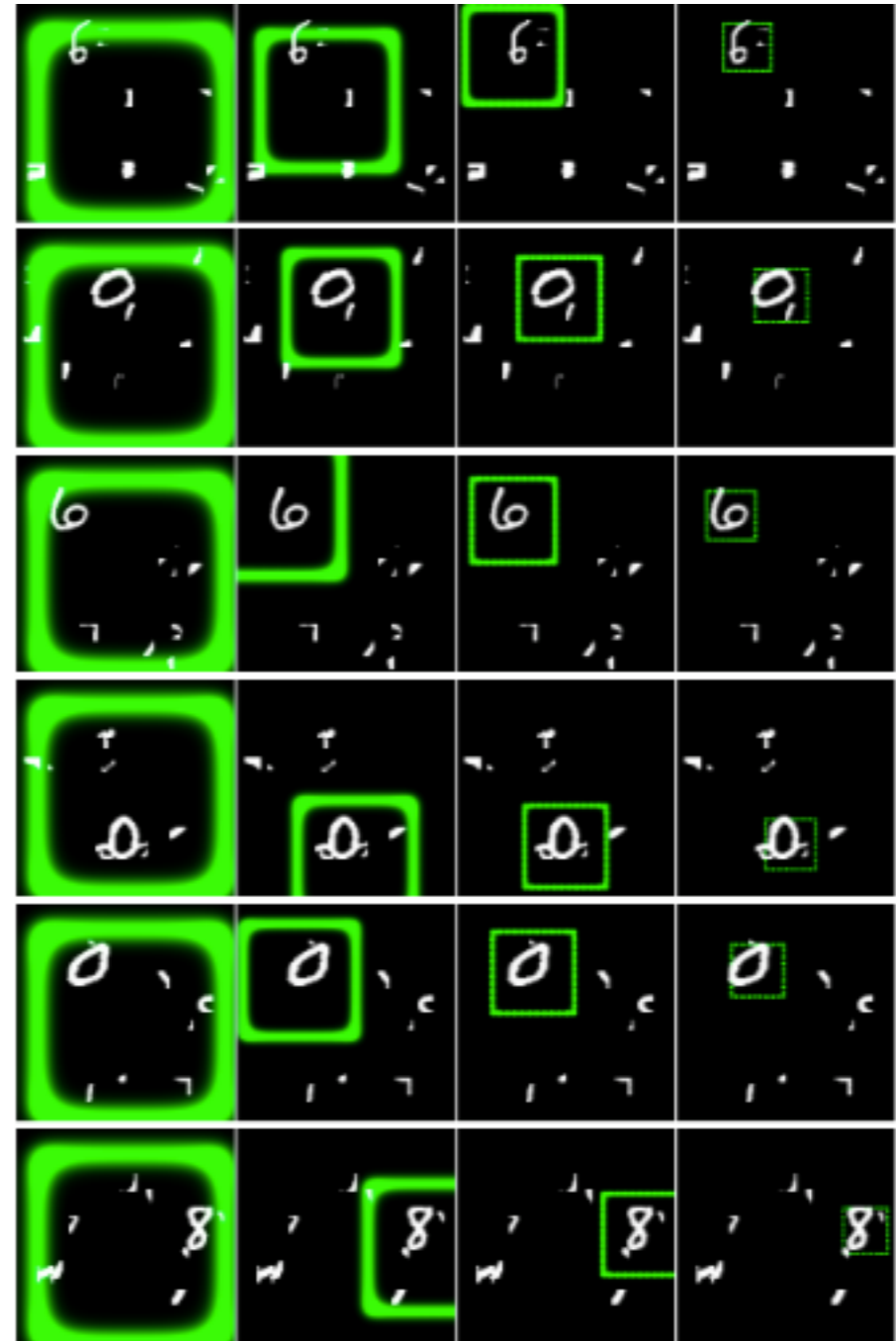
# Differentiable RAM



$$\mu_X^i = g_X + (i - N/2 - 0.5) \delta$$

$$\mu_Y^j = g_Y + (j - N/2 - 0.5) \delta$$

$$(\tilde{g}_X, \tilde{g}_Y, \log \sigma^2, \log \tilde{\delta}, \log \gamma) = W(h^{dec})$$



# Differentiable RAM performance

*Table 1.* Classification test error on  $100 \times 100$  Cluttered Translated MNIST.

Model	Error
Convolutional, 2 layers	14.35%
RAM, 4 glimpses, $12 \times 12$ , 4 scales	9.41%
RAM, 8 glimpses, $12 \times 12$ , 4 scales	8.11%
Differentiable RAM, 4 glimpses, $12 \times 12$	4.18%
Differentiable RAM, 8 glimpses, $12 \times 12$	<b>3.36%</b>

# DRAW-ing with attention



Reading MNIST