

# Attention and Memory in Deep Learning

Alex Graves

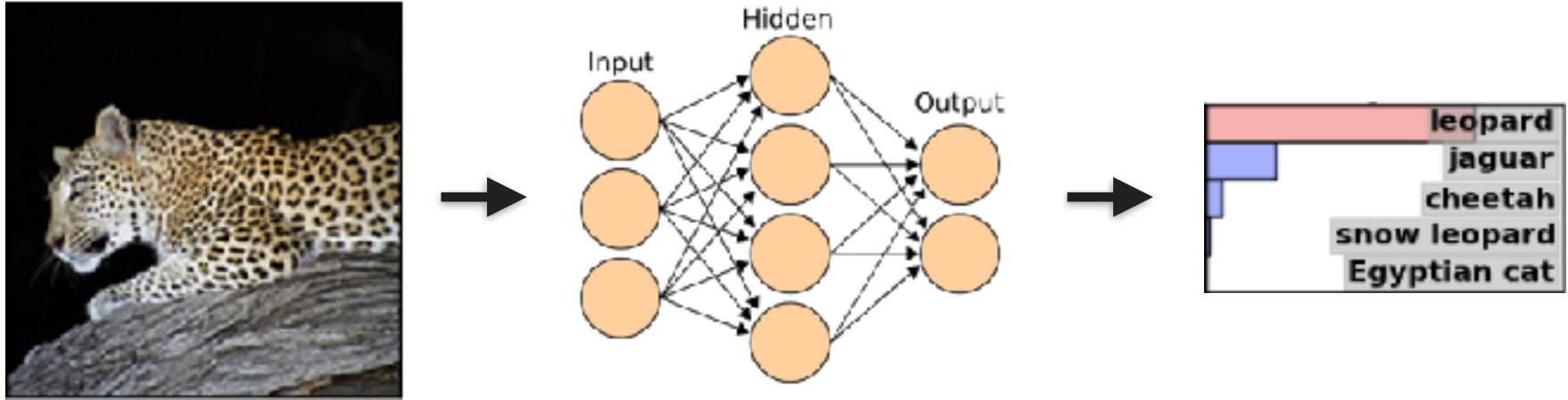
# Attention, memory and Cognition

The ability to focus on one thing and ignore others has a vital role in guiding cognition.

Not only does this allow us to pick out salient information from noisy data (the cocktail party problem) it also allows us to pursue one thought at a time, remember one event rather than all events...



# Neural Networks



Neural nets are parametric, nonlinear function approximations that can be fit to data to learn functions from input vectors (e.g. photographs) to output vectors (e.g. distributions over class labels)

What does that have to do with **attention**?

# Implicit Attention in Neural Networks

Deep nets naturally learn a form of **implicit attention** where they respond more strongly to some parts of the data than others

To a first approximation, we can visualise this by looking at the network **Jacobian** — sensitivity of the network outputs with respect to the inputs

# Neural Network Jacobian

$x$  = size  $k$  input vector

$y$  = size  $m$  output vector

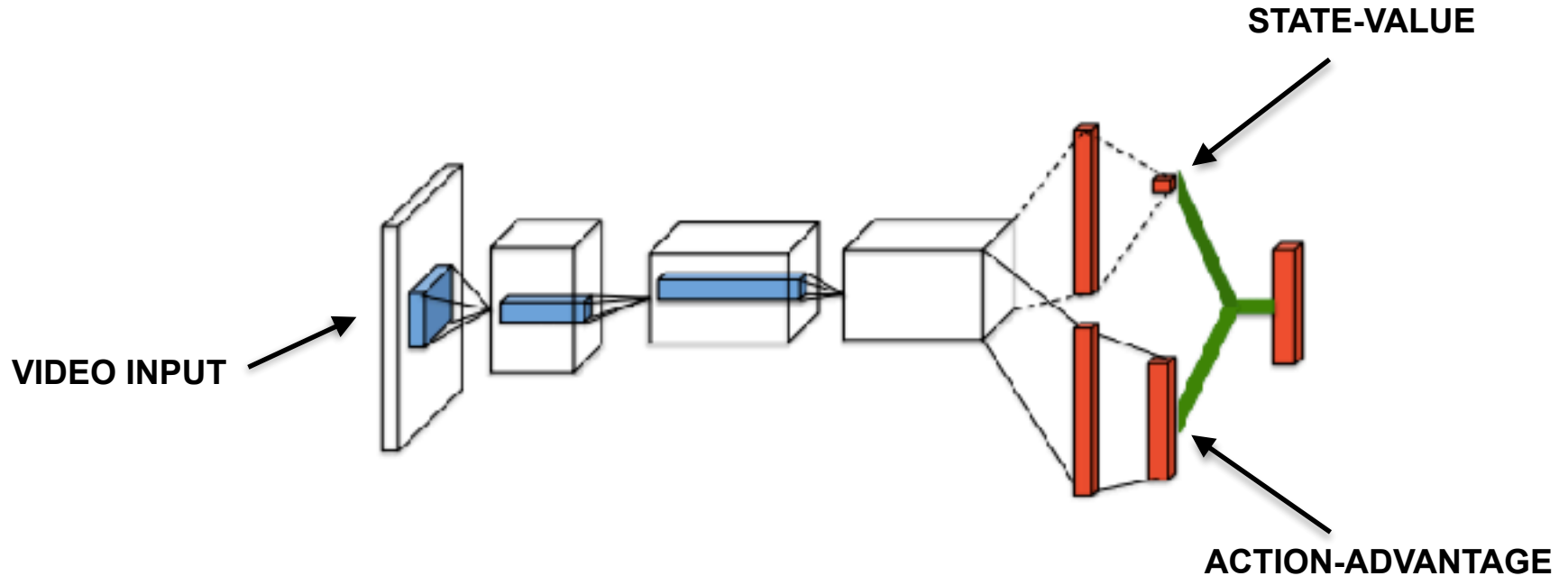
Jacobian  $J$  =  $m \times k$  matrix

$$J_{ij} = \frac{\partial y_i}{\partial x_j}$$

$$J = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_k} \end{bmatrix}$$

Can compute with ordinary backprop  
(just set output 'errors' = output activations)

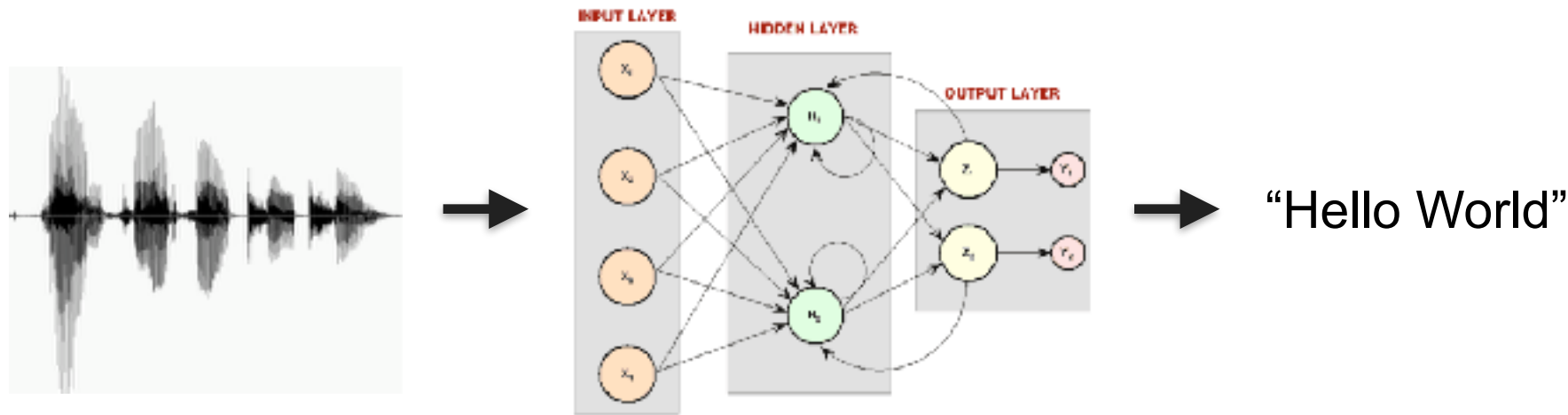
# Jacobian in Action: Duelling Network



*Dueling Network Architectures for Deep Reinforcement Learning, Wang et. al. (2015)*



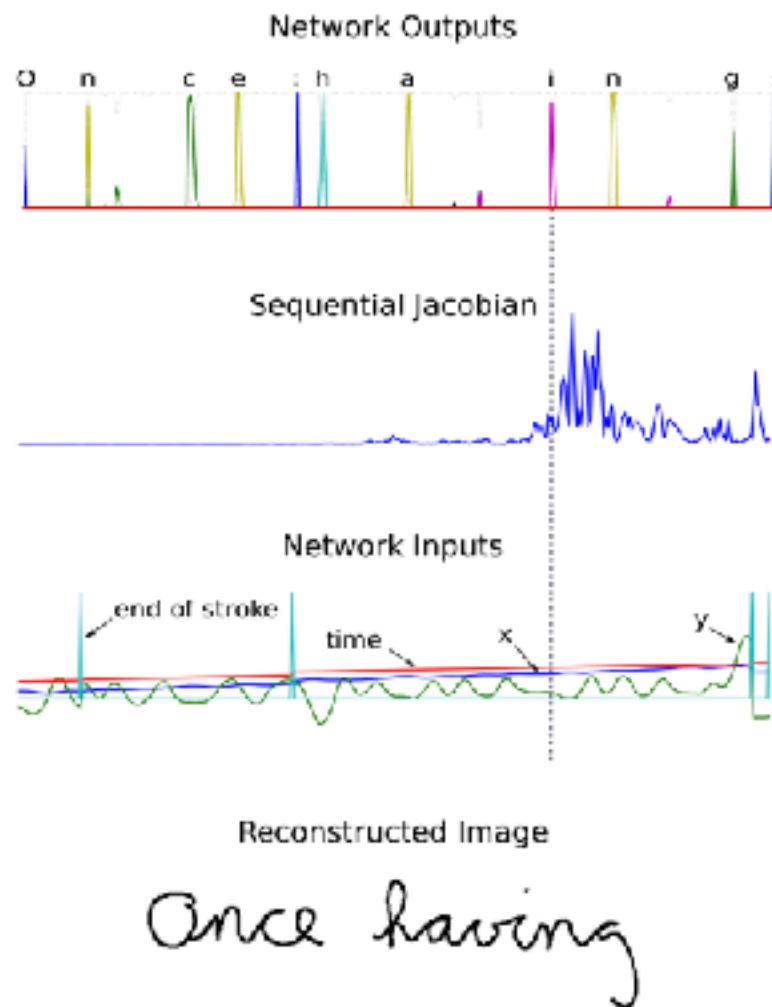
# Attention and memory in Recurrent Networks (RNNs)



**RNNs** contain a recursive hidden state and learn functions from sequences of inputs (e.g. a speech signal) to sequences of outputs (e.g. words)

The **sequential Jacobian** shows which past inputs they **remember** when predicting current outputs.

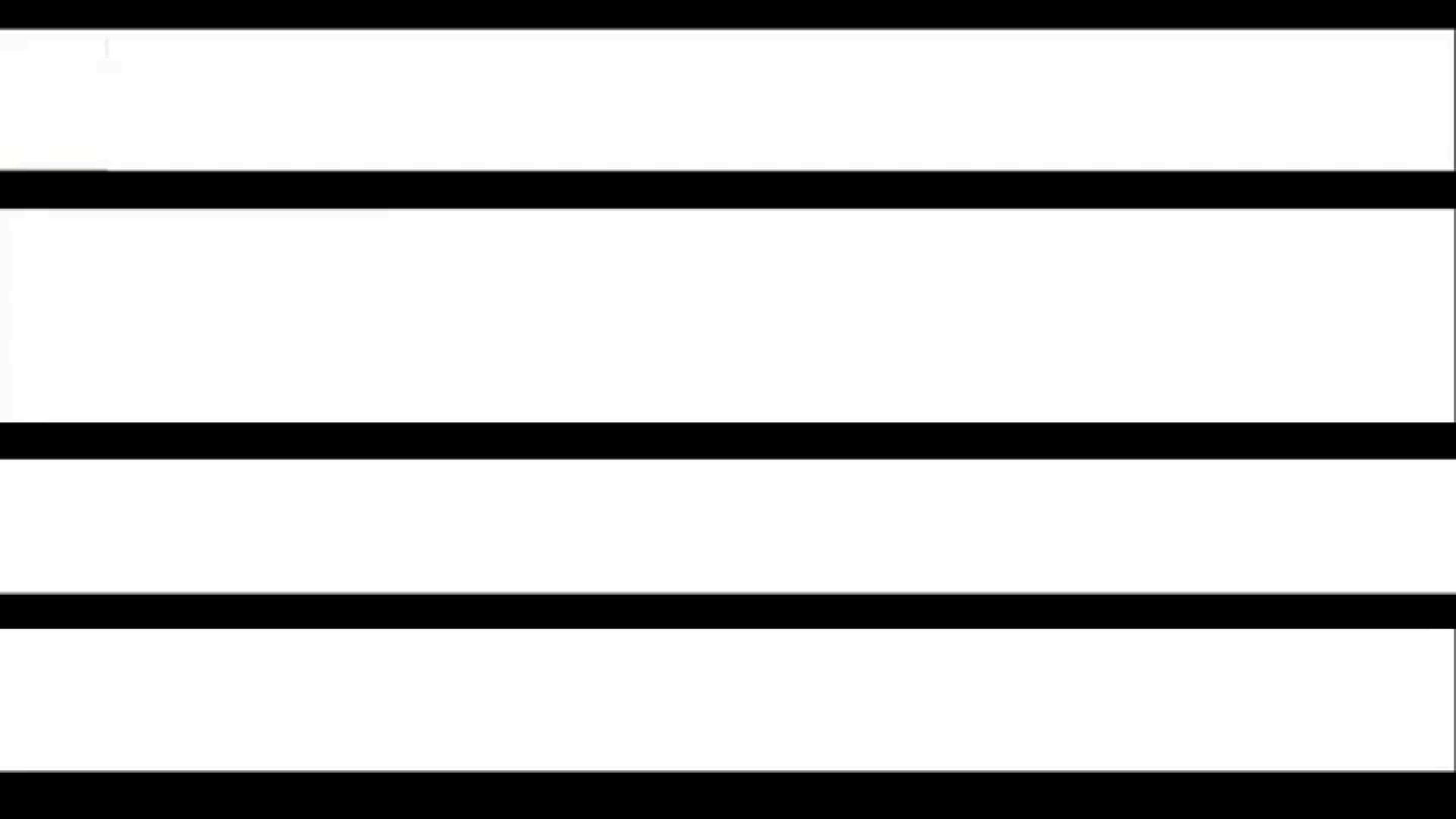




- ▶ The **Sequential Jacobian** is the set of derivatives of one network output with respect to all the inputs

$$J_k^t = \left( \frac{\partial y_k^t}{\partial \mathbf{x}^1}, \frac{\partial y_k^t}{\partial \mathbf{x}^2} \dots \right)$$

- ▶ It shows how the network responds to widely separated, but related, inputs, such as the **delayed dot** of the 'i' in 'having'



*“to reach” -> “zu erreichen”*

# Explicit Attention

Implicit attention is great, but there are still advantages to an **explicit attention** mechanism that limits the data presented to the network in some way:

- Computational **efficiency**
- **Scalability** (e.g. fixed sized glimpse for any size image)
- **Sequential processing** of static data (e.g. moving gaze)
- **Easier to interpret**

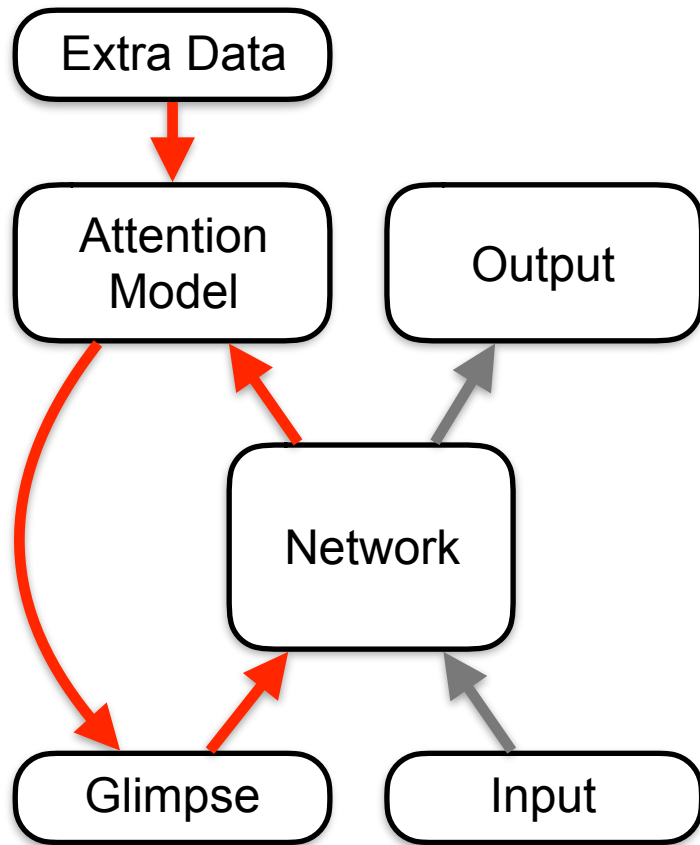
# Neural Attention Models

The **network** receives **input** and produces **output** as usual.

It also produces an extra set of outputs used to parameterise an **attention model**

The attention model then operates on some **extra data** (image, audio sample, text to be translated...) to create a fixed-size “**glimpse**” that is passed to the network as an extra input at the next time step

The complete system is **recurrent**, even if the network isn't



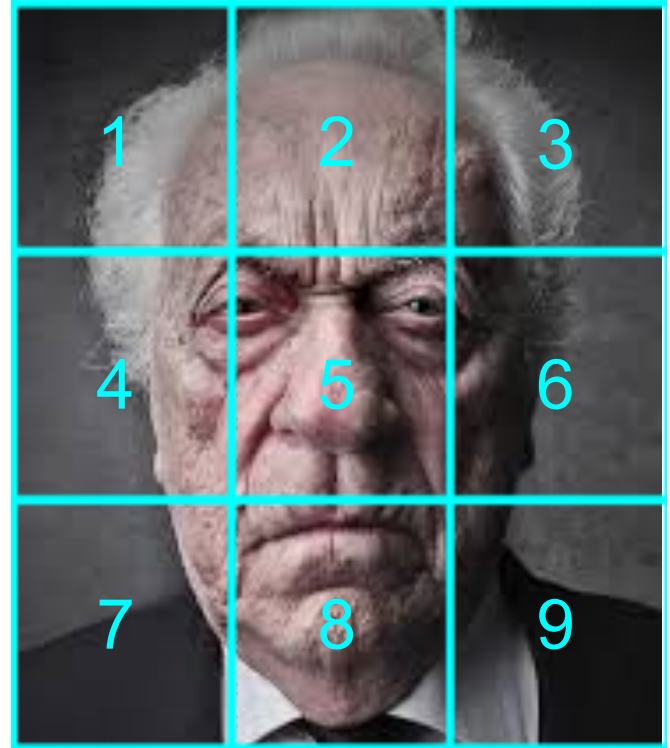
# Glimpse Distribution

Attention models generally work by defining a probability distribution over glimpses  $\mathbf{g}$  of the data  $\mathbf{x}$  and some set of attention outputs  $\mathbf{a}$  from the network:

$$\Pr(\mathbf{g}|\mathbf{a})$$

simplest case:  $\mathbf{a}$  just assigns probabilities to a set of discrete glimpses:

$$\Pr(\mathbf{g}_k|\mathbf{a}) = \frac{\exp(a_k)}{\sum_{k'} \exp(a_{k'})}$$



# Attention with RL

We can treat the distribution over glimpses  $\mathbf{g}$  as a **stochastic policy**  $\pi_{\mathbf{a}}$ , sample from it, and use **RL** techniques (with reward  $R = \text{task loss } L$  induced by the glimpse) to train the attention model

$$\pi_{\mathbf{a}} = \text{Pr}(\mathbf{g}_k | \mathbf{a})$$

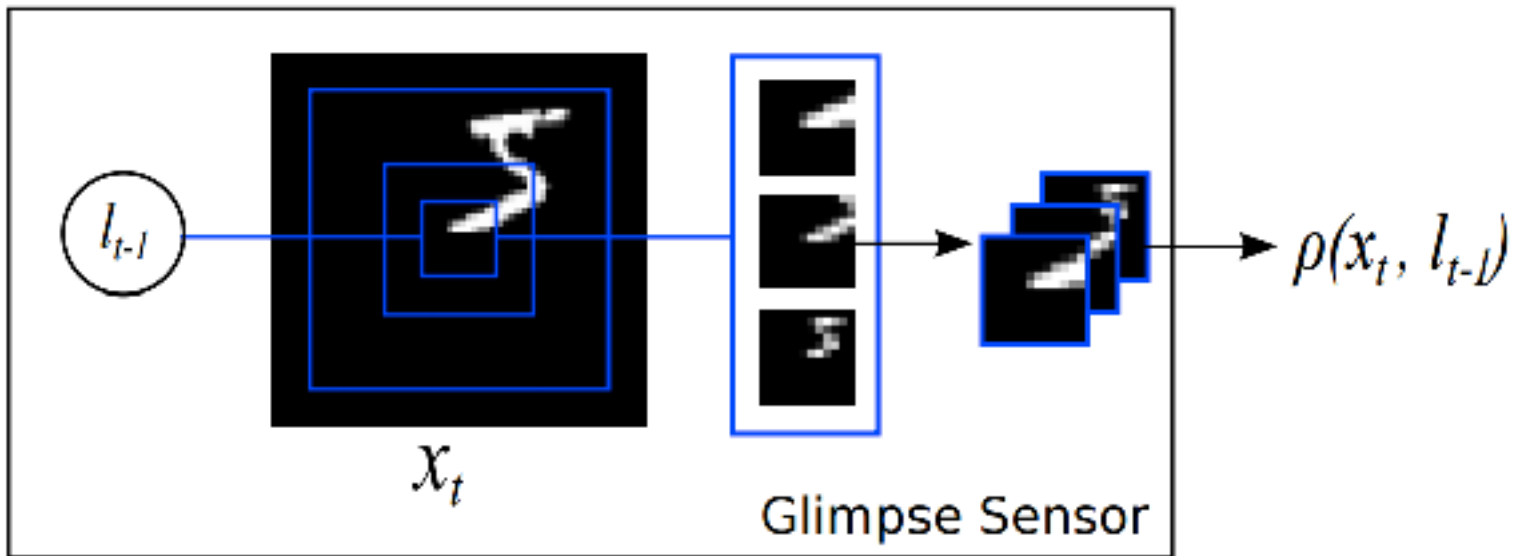
$$R = \mathbb{E}_{\mathbf{g} \sim \pi_{\mathbf{a}}} [\log \pi_{\mathbf{a}} L(\mathbf{g})]$$

$$\nabla_{\mathbf{a}} R = \mathbb{E}_{\mathbf{g} \sim \pi_{\mathbf{a}}} [\nabla_{\mathbf{a}} \log \pi_{\mathbf{a}} L(\mathbf{g})]$$

In general we can use RL methods for supervised tasks any time some module in the network is **non-differentiable**

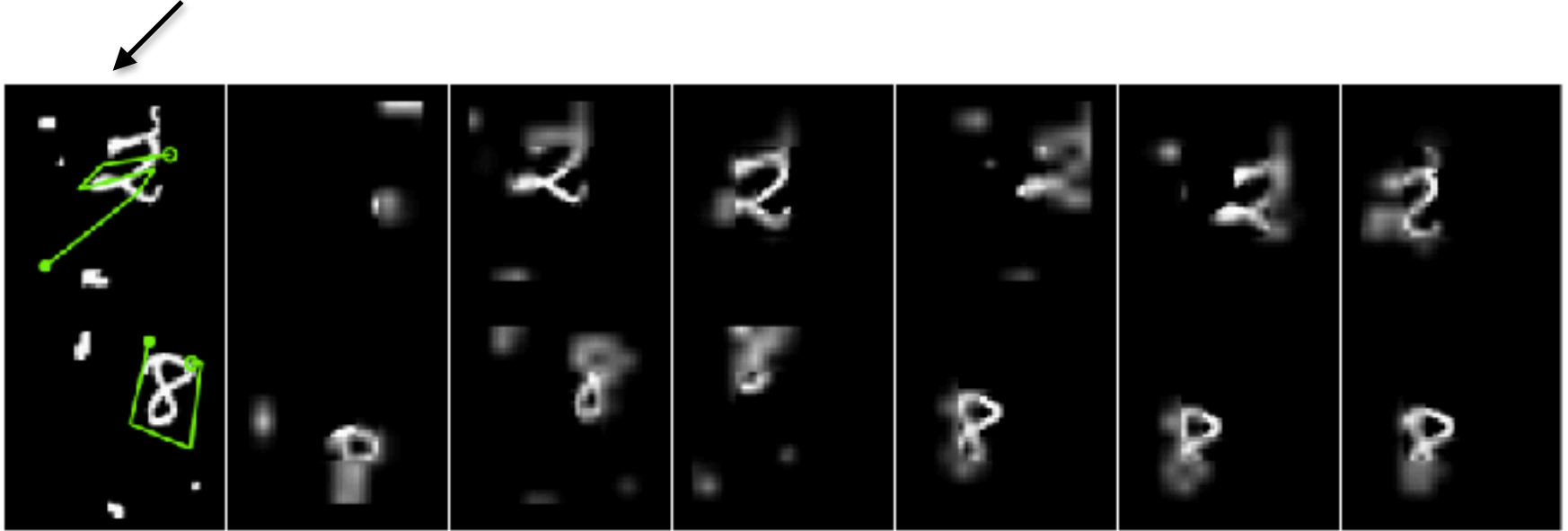
# Complex Glimpses

Generally the glimpse distribution is more complex than just a softmax (e.g. Gaussian over co-ordinates, width, height...) and the glimpses are more complex than image tiles (e.g. foveal models)



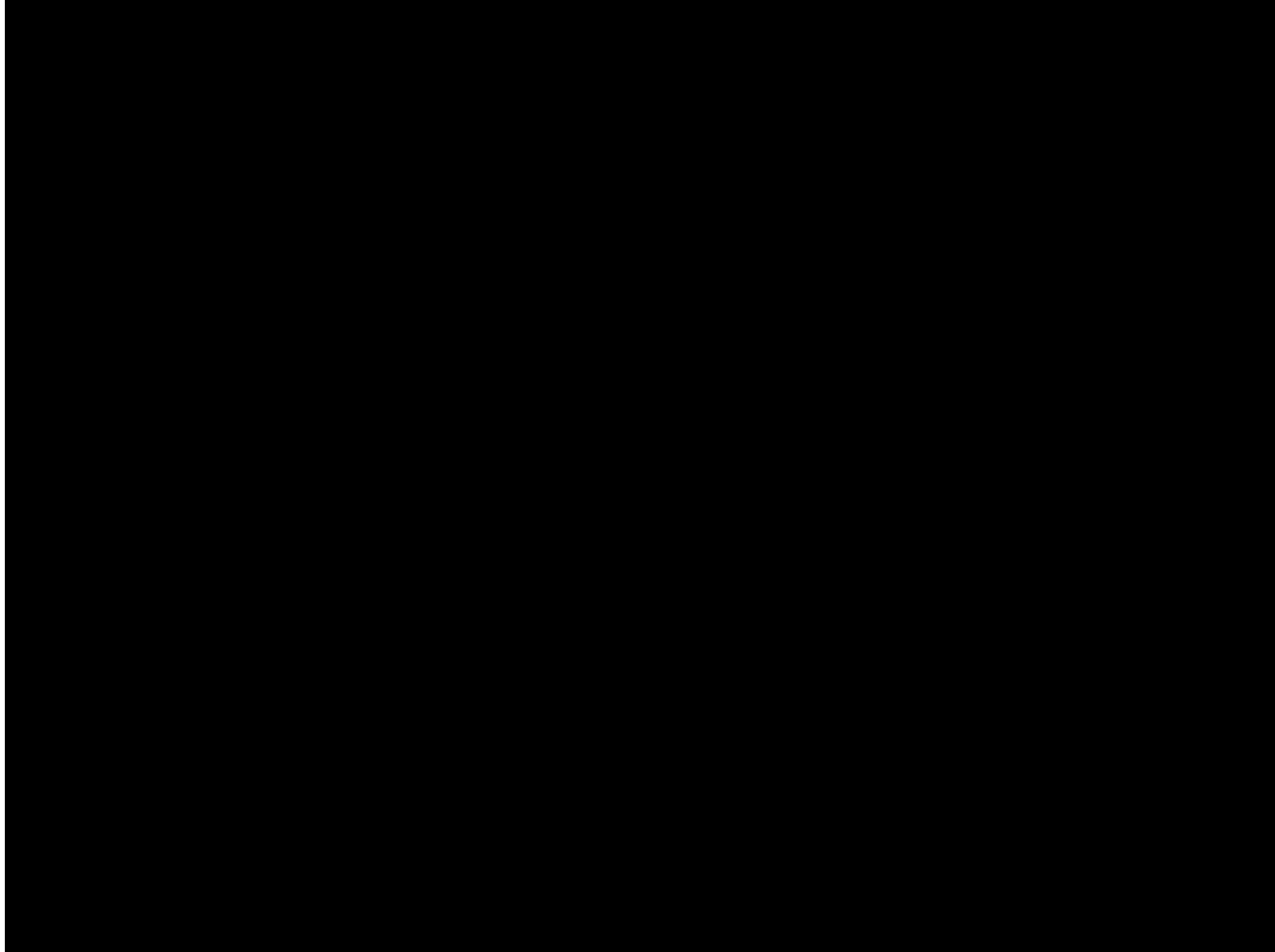


6 point “Glimpse path” (in green) while trying to classify image



6 Foveal Glimpses seen by the network

*Recurrent Models of Visual Attention*, Mnih et. al. (2014)



*Multiple Object Recognition with Visual Attention*, Ba et. al. (2014)

# Soft Attention

The last examples used **hard attention**: fixed size attention windows moved around the image, trained with RL techniques.

Robots have to look left or right, but in many cases attention doesn't need to be hard: we just want to focus more on certain regions and less on others.

If we do this in a differentiable way, we get **soft attention** which we can train **end-to-end** with **backprop**

Generally easier than using RL, but more expensive to compute

# Soft Attention

**Basic template:** we use the attention parameters  $\mathbf{a}$  to determine a distribution  $\Pr(\mathbf{g}|\mathbf{a})$  as before, only now we take an **expectation** over all possible glimpses instead of a **sample**

$$\mathbf{g} = \sum_{\mathbf{g}' \in \mathbf{x}} \mathbf{g}' \Pr(\mathbf{g}'|\mathbf{a})$$

This is differentiable w.r.t.  $\mathbf{a}$  as long as  $\Pr(\mathbf{g}|\mathbf{a})$  is:

$$\nabla_{\mathbf{a}} \mathbf{g} = \sum_{\mathbf{g}' \in \mathbf{x}} \mathbf{g}' \nabla_{\mathbf{a}} \Pr(\mathbf{g}'|\mathbf{a})$$

# Location Attention

Window vector (input to net)

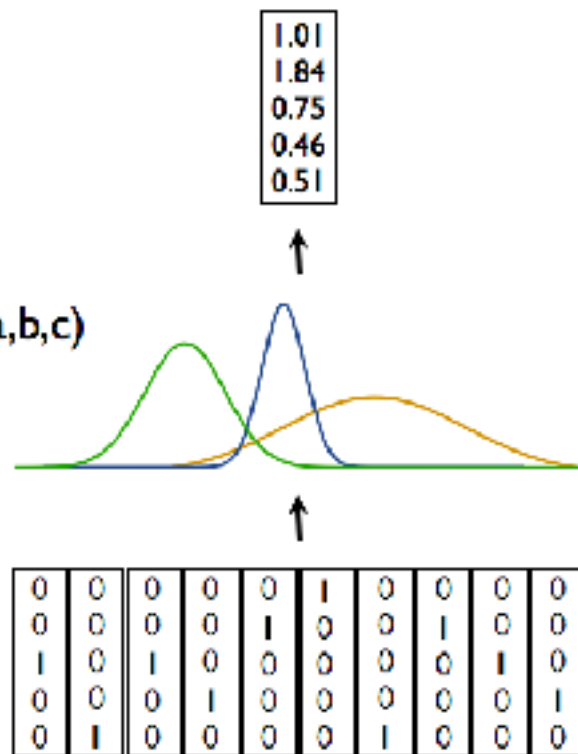
$$v^{t+1} = \sum_{i=1}^S w_i^t s_i$$

Window weights (net outputs for a,b,c)

$$w_i^t = \sum_{k=1}^K a_k^t \exp(-b_k^t [c_k^t - i]^2)$$

Input vectors (one-hot)

$(s_1, \dots, s_S)$



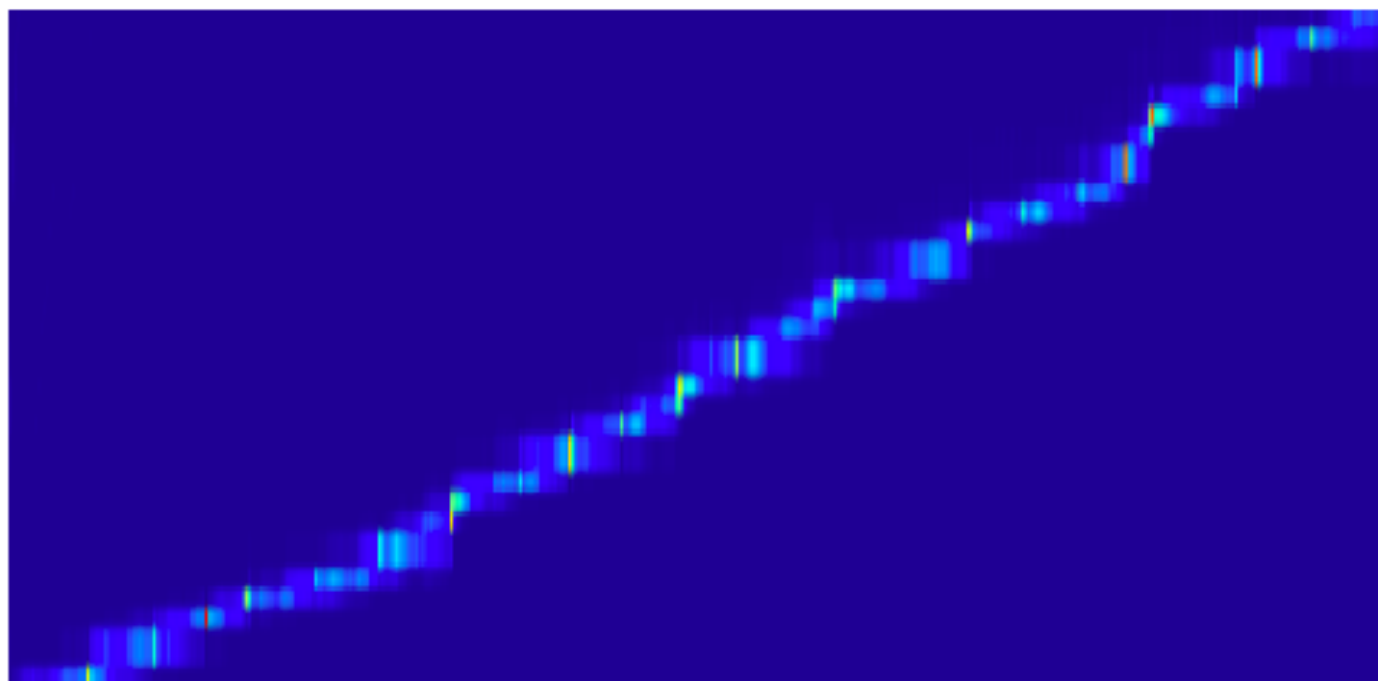
*Generating Sequences with Recurrent Neural Networks, Graves (2013)*

# Handwriting Synthesis With Soft Reading

these sequences were generated by  
picking samples at every star  
every line is a different style  
yes, real people write this badly

# Alignment

Thought that the muster from



Thought that the muster from

# Associative Attention

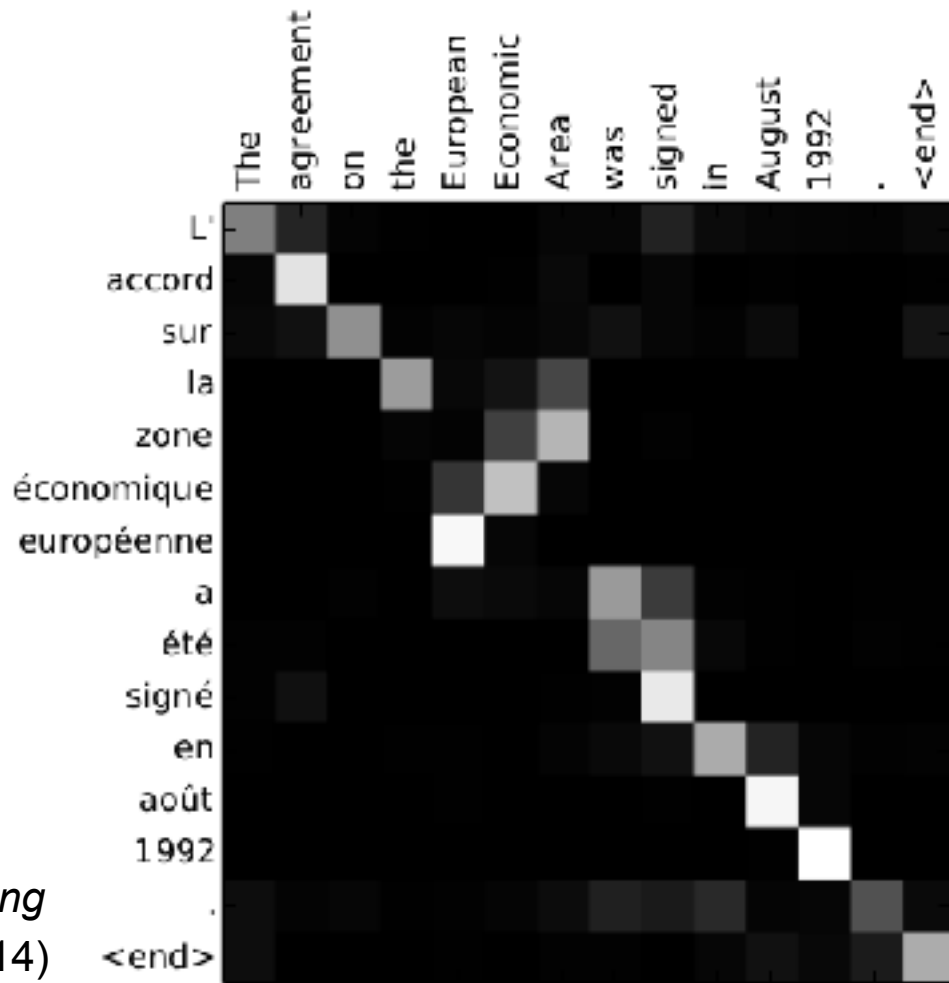
Instead of attending by position, we can attend by **content**. In this setting **a** = a **key vector** that is compared to *all* glimpses **g** using some **similarity function** *S*. The scores are normalised and used to define  $\text{Pr}(\mathbf{g}|\mathbf{a})$

$$\text{Pr}(\mathbf{g}|\mathbf{a}) = \frac{\exp(S(\mathbf{g}, \mathbf{a}))}{\sum_{\mathbf{g}' \in \mathbf{x}} \exp(S(\mathbf{g}', \mathbf{a}))}$$

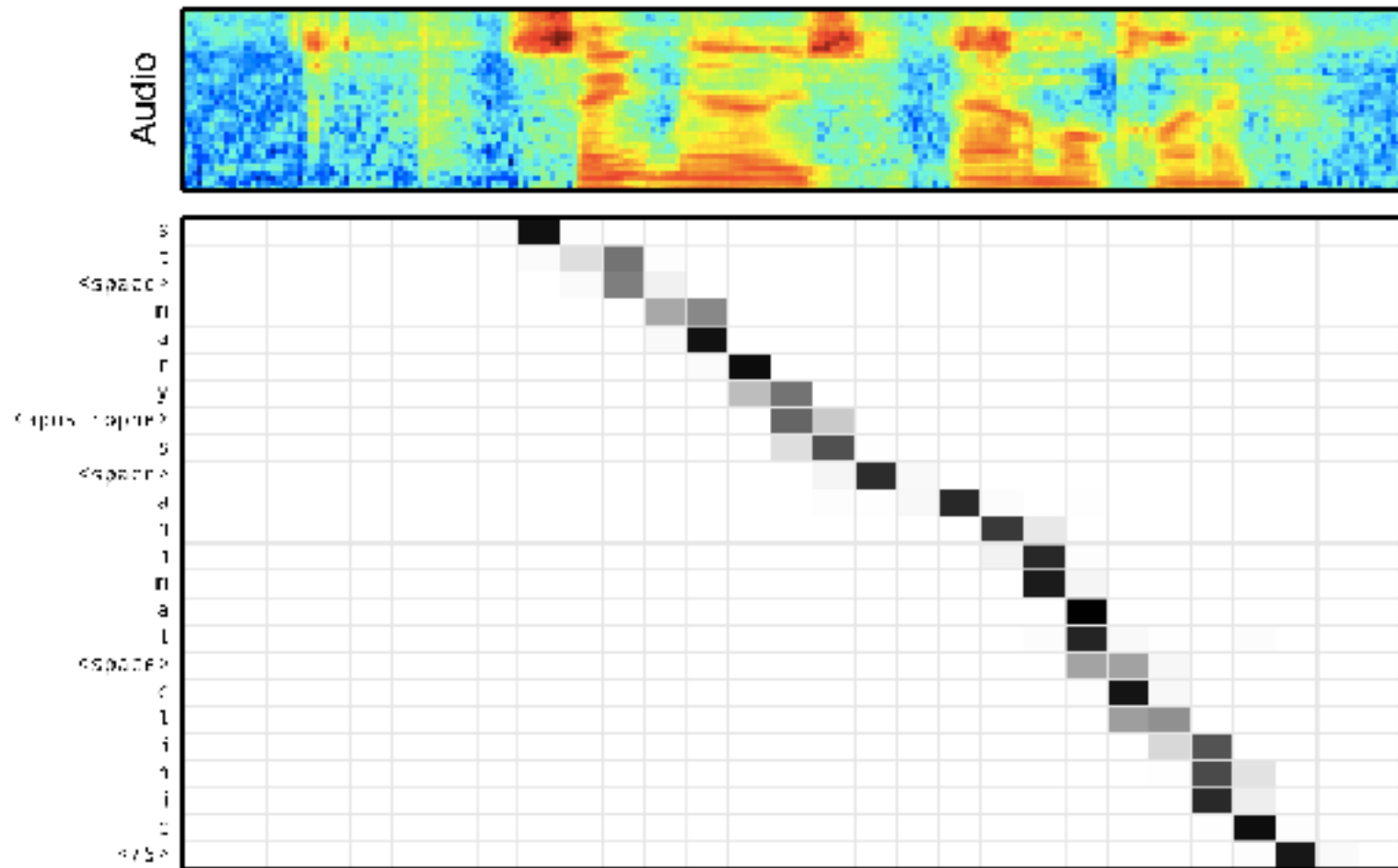
*S* can be learned (MLP, linear operator...) or fixed (dot product, cosine similarity...). Yields a **Multidimensional, feature-based** lookup: natural way to search data



# Reordering in machine translation using associative attention



*Neural Machine Translation by Jointly Learning to Align and Translate, Bahdanau et. al. (2014)*



*Listen, Attend and Spell*, Chan et. al. (2015)

by ent423 ,ent261 correspondent updated 9:49 pm et ,thu  
march 19 ,2015 ( ent261 ) a ent114 was killed in a parachute  
accident in ent45 ,ent85 ,near ent312 ,a ent119 official told  
ent261 on wednesday .he was identified thursday as  
special warfare operator 3rd class ent23 ,29 ,of ent187 ,  
ent265 .` ent23 distinguished himself consistently  
throughout his career .he was the epitome of the quiet  
professional in all facets of his life ,and he leaves an  
inspiring legacy of natural tenacity and focused

...

ent119 identifies deceased sailor as X ,who leaves behind  
a wife

by ent270 ,ent223 updated 9:35 am et ,mon march 2 ,2015  
( ent223 ) ent63 went famillal for fall at its fashion show in  
ent231 on sunday ,dedicating its collection to ``mamma"  
with nary a pair of ``mom jeans " in sight .ent164 and ent21 ,  
who are behind the ent196 brand ,sent models down the  
runway in decidedly feminine dresses and skirts adorned  
with roses ,lace and even embroidered doodles by the  
designers ' own nieces and nephews .many of the looks  
featured saccharine needlework phrases like ``I love you ,

...

X dedicated their fall fashion show to moms

# Introspective Attention

So far we have looked at attention to **external data**

Also useful to selectively attend to the network's internal state or memory: **introspective attention**

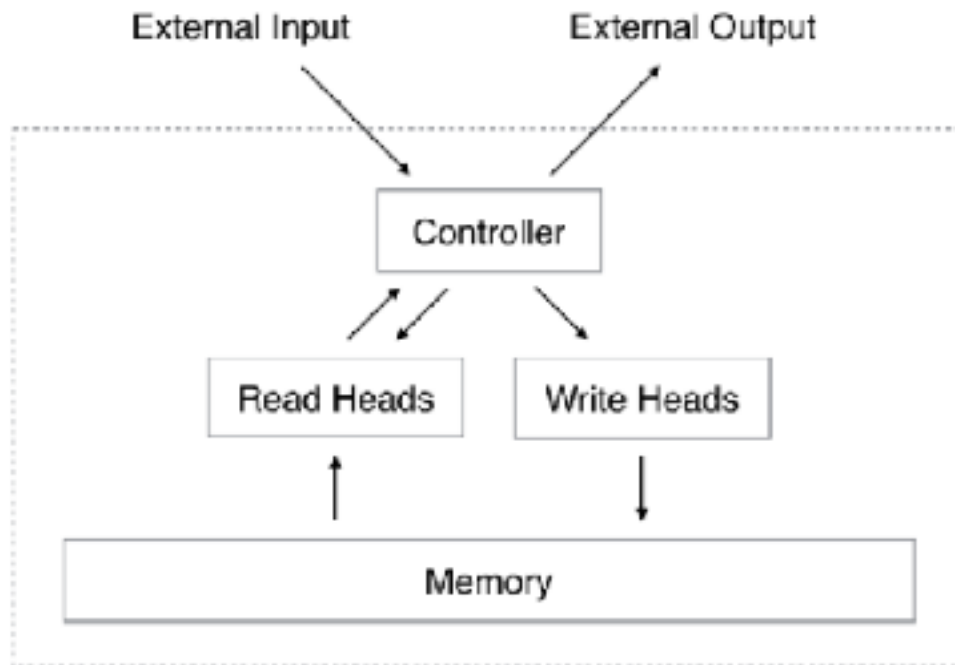
With internal information we can do selective **writing** as well as **reading**, allowing the network to **iteratively modify** its state

# Neural Turing Machines

The **Controller** is a **neural network** (recurrent or feedforward)

The **Heads** **select** portions of the memory and **read** or **write** to them

The **Memory** is a real-valued **matrix**



*Neural Turing Machines, Graves et. al. (2014)*

# Addressing by Content

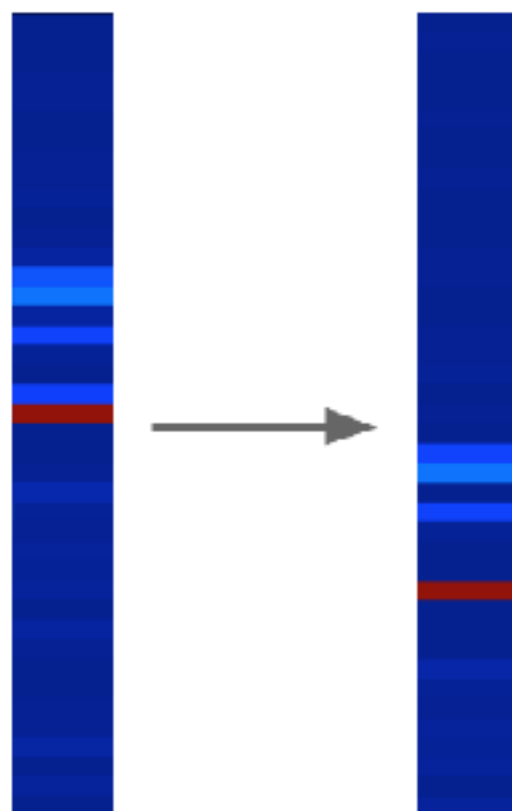
A **key vector**  $\mathbf{k}$  is emitted by the controller and compared to the content of each memory location  $\mathbf{M}[i]$  using a similarity measure  $S(\cdot, \cdot)$  (e.g. **cosine distance**) then normalised with a **softmax**. A '**sharpness**'  $\beta$  is used to narrow the focus. Finds the memories '**closest**' to the key

$$\mathbf{w}[i] = \frac{\exp(\beta S(\mathbf{k}, \mathbf{M}[i]))}{\sum_j \exp(\beta S(\mathbf{k}, \mathbf{M}[j]))}$$

# Addressing by Location

The controller outputs a **shift kernel**  $\mathbf{s}$  (e.g. a softmax on  $[-n,n]$ ) which is convolved with a weighting  $\mathbf{w}$  to produce a shifted weighting  $\hat{\mathbf{w}}$ .

$$\hat{\mathbf{w}}[i] = \sum_j \mathbf{w}[j] \mathbf{s}(i - j)$$



# Data Structure and Accessors

The combination of addressing mechanisms allows the controller to interact with the memory in several distinct modes, corresponding to different **data structures** and **accessors**.

**Content key only** — memory is accessed like an **associative map**

**Content and location** — key finds an **array**, shift **indexes** into it

**Location only** — shift **iterates** from the last focus



# Reading and Writing

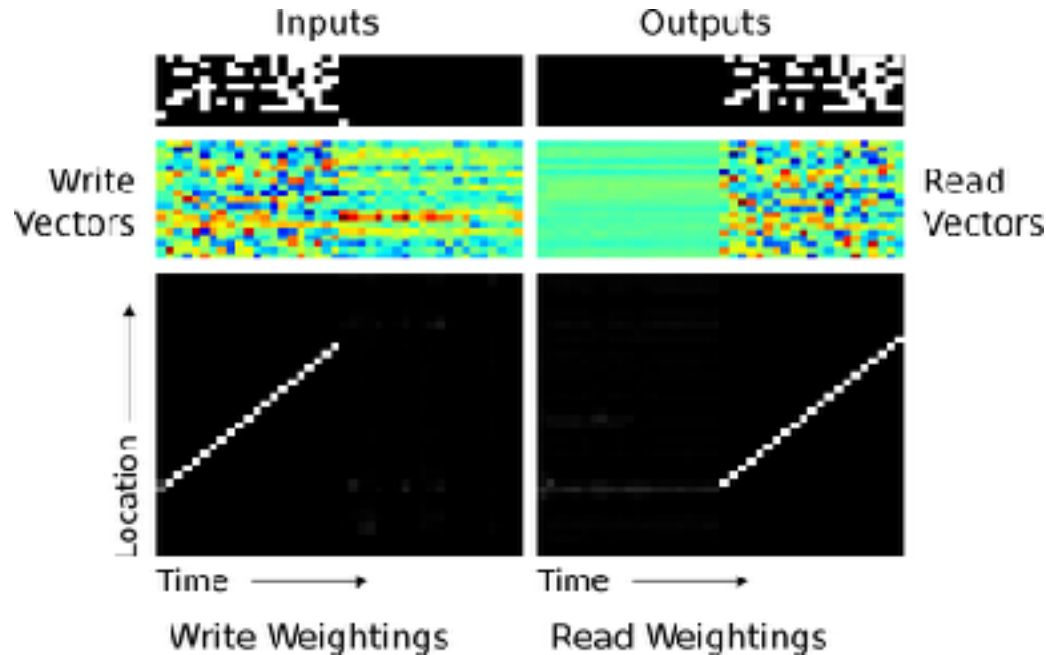
Once the weightings are defined, each **read head** returns a **read vector**  $\mathbf{r}$  as input to the controller at the next timestep

$$\mathbf{r} = \sum_i \mathbf{w}[i] \mathbf{M}[i]$$

Each **write head** receives an **erase vector**  $\mathbf{e}$  and an **add vector**  $\mathbf{a}$  from the controller and resets then writes to modify the memory (like **LSTM**)

$$\mathbf{M}[i] \leftarrow \mathbf{M}[i](\mathbf{1} - \mathbf{w}[i]\mathbf{e}) + \mathbf{w}[i]\mathbf{a}$$

# The NTM Copy Algorithm

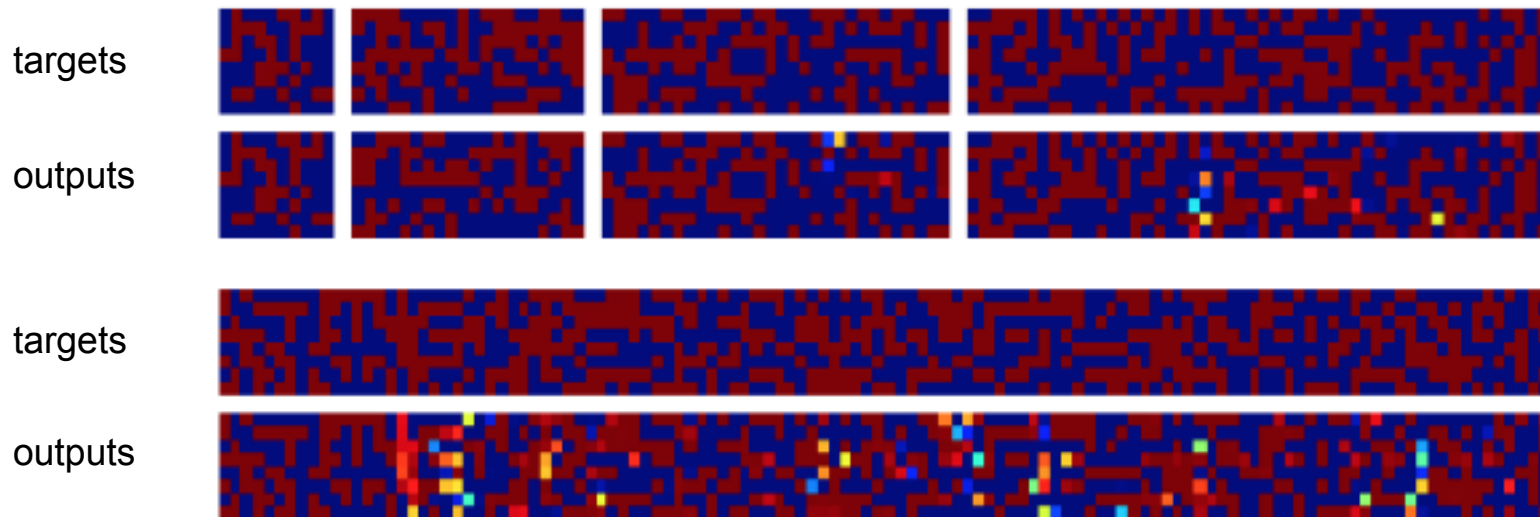


NTM++

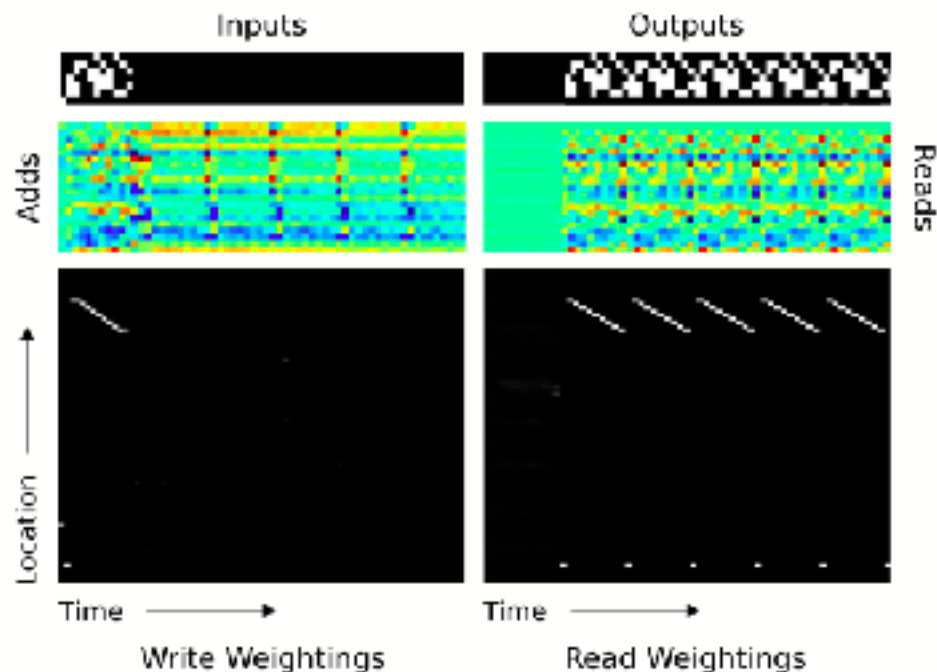
```
initialize: move head to start location
while input delimiter not seen do
    receive input vector
    write input to head location
    increment head location by 1
end while
return head to start location
while true do
    read output vector from head location
    emit output
    increment head location by 1
end while
```

pseudocode

# Copy Generalisation: length 10 to 120

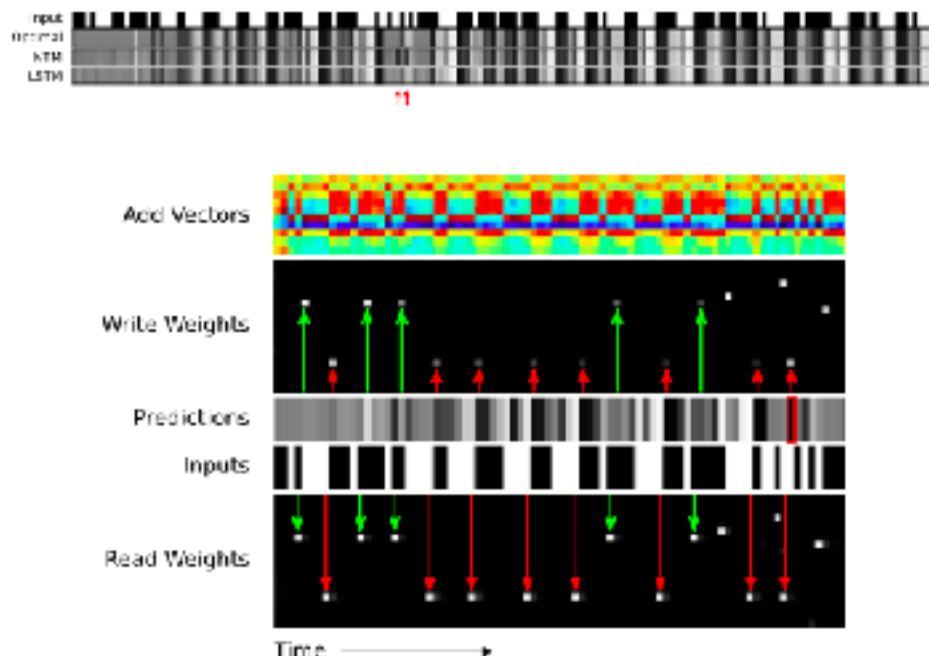


# Copy $N$ Times



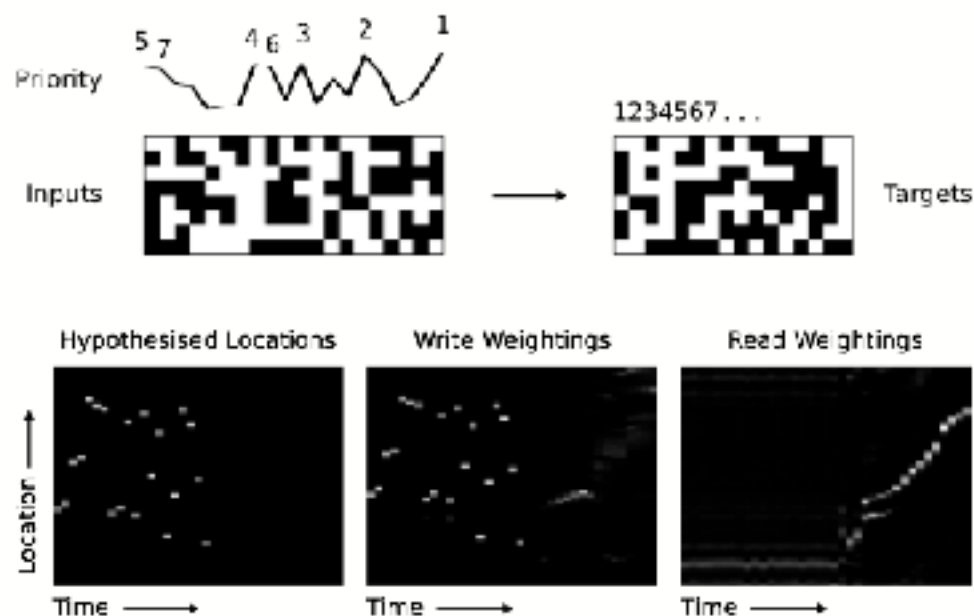
NTM learns its first **for-loop**, using **content** to jump, **iteration** to step, and a **variable** to **count** to  $N$ .

# N-Gram Inference



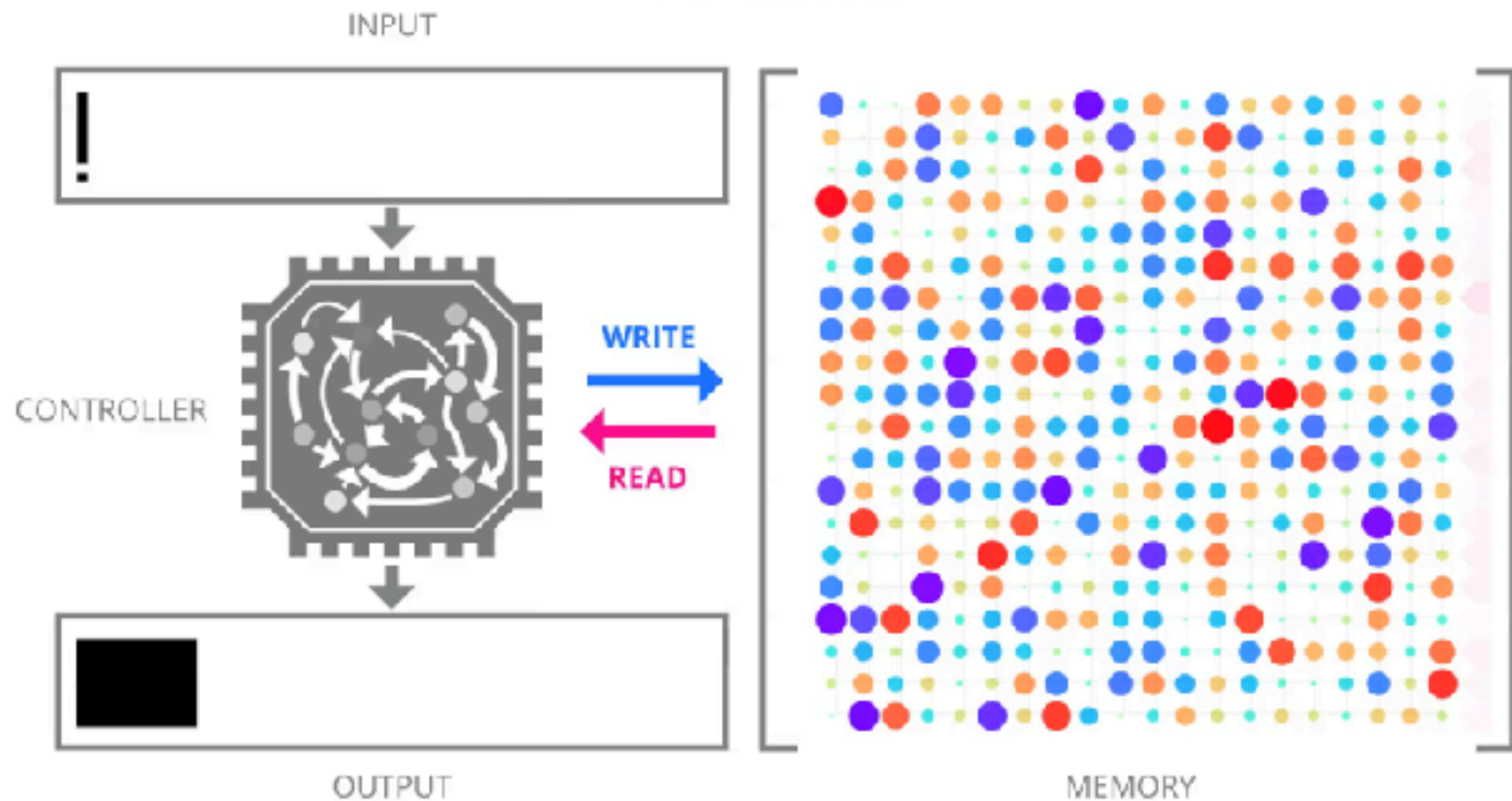
Specific memory locations store **variables** that **count** the **occurrences** of particular N-Grams

# Priority Sort



The network maps from **priorities** to **write locations**, then **iterates** through the memory to return the sorted list

# TRAINING





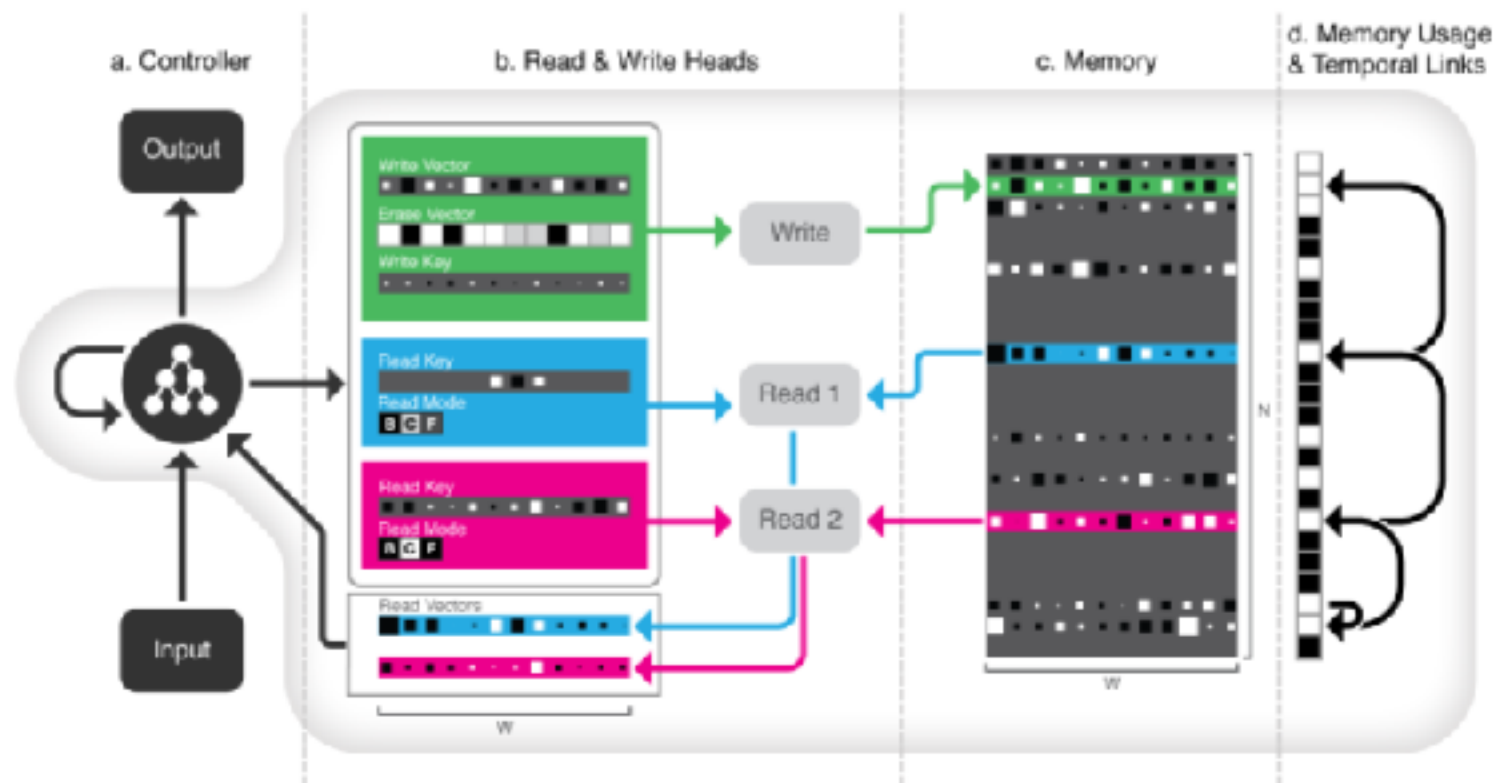
# Differentiable Neural Computers

**DNC** is a successor architecture to Neural Turing Machines with new attention mechanisms for memory access

*hybrid computing using a neural network with dynamic external memory, Graves et. al. (2016)*



# Overall Architecture



# Allocating Memory

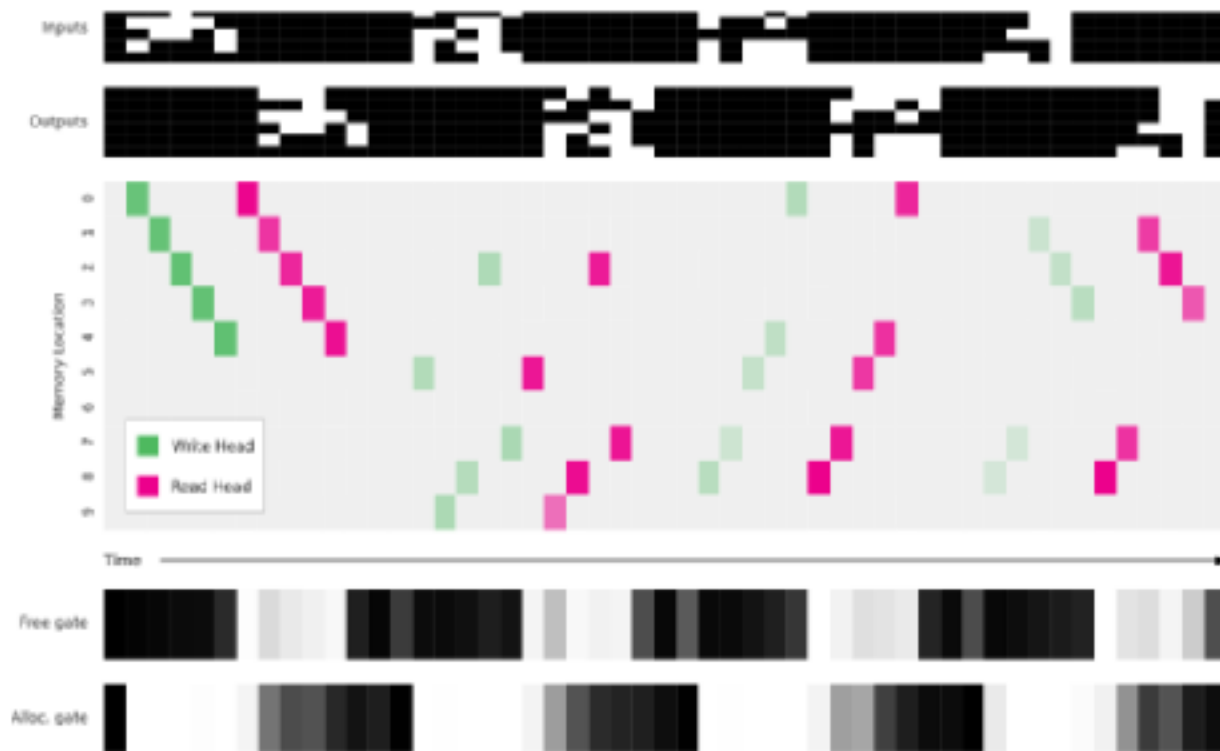
- NTM could only 'allocate' memory in contiguous blocks, leading to memory management problems
- DNC defines a differentiable *free list* tracking the *usage* ( $u_t$ ) of each memory location
- Usage is automatically *increased* after each write ( $\mathbf{w}_t^w$ ) and optionally *decreased* after each read ( $\mathbf{w}_t^r$ ) by *free gates* ( $f_t^i$ )

$$\mathbf{u}_t = (\mathbf{u}_{t-1} + \mathbf{w}_{t-1}^w - \mathbf{u}_{t-1} \circ \mathbf{w}_{t-1}^w) \circ \prod_{i=1}^R (1 - f_t^i \mathbf{w}_{t-1}^{r,i})$$

- The controller then uses an *allocation gate* ( $g_t^a$ ) to interpolate between writing to a newly allocated ( $\mathbf{a}_t$ ) location, or an existing one found by content ( $\mathbf{c}_t^w$ )

$$\mathbf{w}_t^w = g_t^w (g_t^a \mathbf{a}_t + (1 - g_t^a) \mathbf{c}_t^w)$$

# Memory Allocation Test



# Searching By Time

- NTM was able to retrieve memories in order of their *index* but not in order in which they were written
- Preserving temporal order is necessary for many tasks (e.g. sequence of instructions) and appears to play an important role in human cognition
- We wanted DNC to be able to iterate through memories in the order they were written (or rewritten)
- A *precedence weighting* ( $p_t$ ) keeps track of which locations were most recently written to:

$$p_t = \left(1 - \sum_i w_t^w[i]\right) p_{t-1} + w_t^w$$

# Searching By Time

$\mathbf{p}_t$  is then used to update a *temporal link matrix* ( $L_t$ ), defining the order locations were written in:

$$L_t[i,j] = (1 - \mathbf{w}_t^w[i] - \mathbf{w}_t^w[j])L_{t-1}[i,j] + \mathbf{w}_t^w[i]\mathbf{p}_{t-1}[j]$$

The controller can use  $L_t$  to retrieve the write before ( $\mathbf{b}_t^i$ ) or after ( $\mathbf{f}_t^i$ ) the last read location ( $\mathbf{w}_{t-1}^r$ ), allowing it to iterate forwards or backwards in time

$$\mathbf{b}_t^i = \hat{L}_t^\top \mathbf{w}_{t-1}^{r,i} \quad \mathbf{f}_t^i = \hat{L}_t \mathbf{w}_{t-1}^{r,i}$$

Finally three-way gates ( $\pi_t^i$ ) are used to interpolate among iterating forwards, backwards, or reading by content:

$$\mathbf{w}_t^{r,i} = \pi_t^i[1]\mathbf{b}_t^i + \pi_t^i[2]\mathbf{c}_t^{r,i} + \pi_t^i[3]\mathbf{f}_t^i$$

# Graph Experiments

## Training Data

a. Random Graph

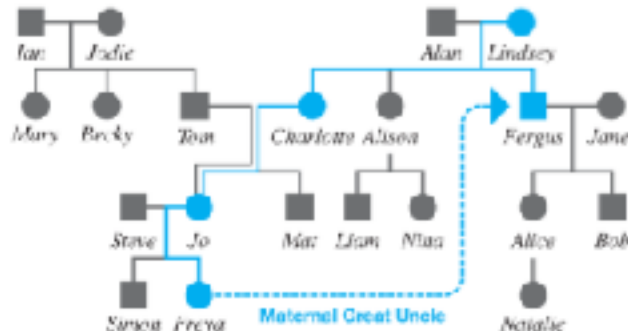


## Test Examples

b. London Underground



c. Family Tree



### Underground Input:

(Oxford Circus, Tottenham Court Road, Central)  
(Tottenham Court Road, Oxford Circus, Central)  
(Baker Street, Marylebone, Circle)  
(Baker Street, Marylebone, Bakerloo)  
(Baker Street, Oxford Circus, Bakerloo)

...

(Leicester Square, Charing Cross, Northern)  
(Tottenham Court Road, Leicester Square, Northern)  
(Oxford Circus, Piccadilly Circus, Bakerloo)  
(Oxford Circus, Tottenham Court Road, Central)  
(Oxford Circus, Euston, Victoria)

- 84 edges in total

### Traversal Question:

(Oxford Circus, Tottenham Court Road, Central)  
(Tottenham Court Road, Oxford Circus, Central)  
(Baker Street, Marylebone, Circle)  
(Baker Street, Marylebone, Bakerloo)  
(Baker Street, Oxford Circus, Bakerloo)

### Answer:

(Oxford Circus, Tottenham Court Road, Central)  
(Tottenham Court Road, Oxford Circus, Central)  
(Baker Street, Marylebone, Circle)  
(Baker Street, Marylebone, Bakerloo)  
(Baker Street, Oxford Circus, Bakerloo)

### Shortest Path Question:

(Moorgate, Piccadilly Circus, ...)

### Answer:

(Moorgate, Bank, Northern)  
(Bank, Tottenham Court Road, Central)  
(Tottenham Court Road, Leicester Square, Piccadilly)  
(Leicester Square, Piccadilly Circus, Piccadilly)

### Family Tree Input:

(Charlotte, Alan, Father)  
(Simon, Alison, Father)  
(Steve, Simon, Son)  
(Melanie, Alison, Mother)  
(Lindsey, Fergus, Son)

...

(Bob, Jane, Mother)  
(Natalie, Alice, Mother)  
(Mary, Ian, Father)  
(Jane, Alice, Daughter)  
(Mick, Charlotte, Mother)

- 54 edges in total

### Inference Question:

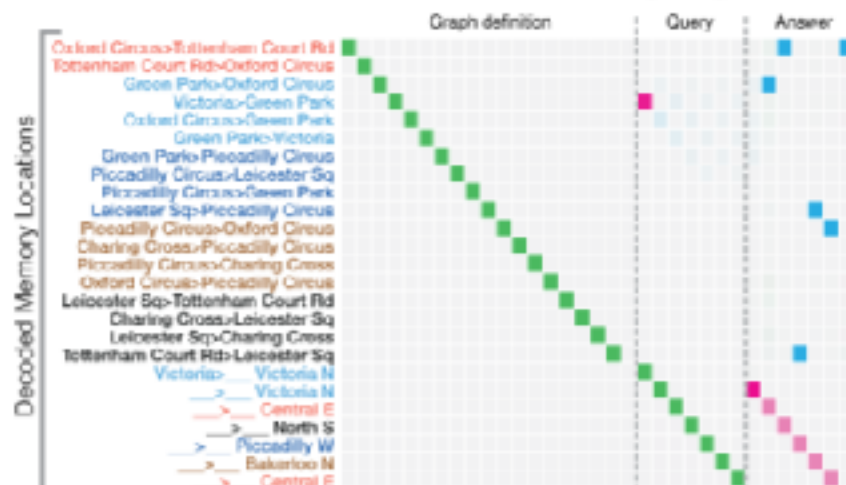
(Freya, Maternal Great Uncle)

### Answer:

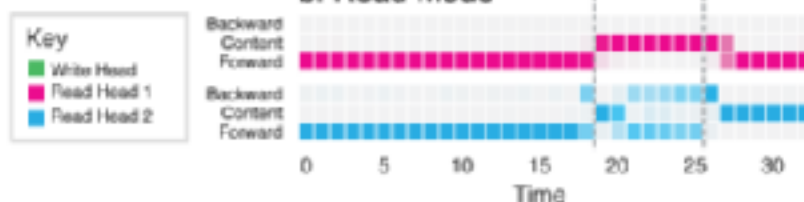
(Freya, Fergus, Maternal Great Uncle)

# How DNC Reads Graphs

a. Read and Write Weightings



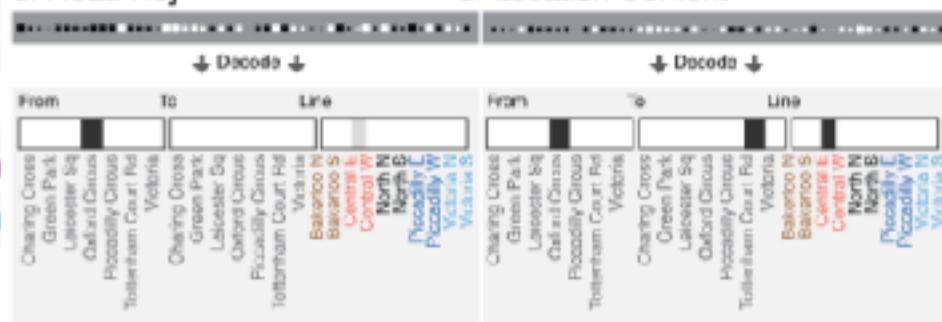
b. Read Mode



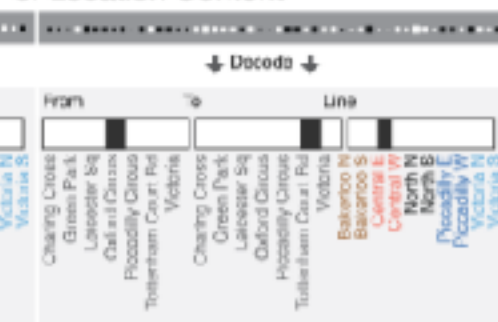
c. London Underground Map



d. Read Key



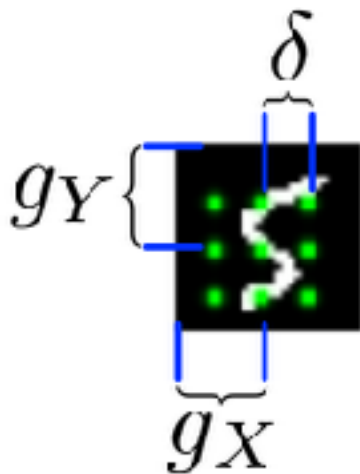
e. Location Content



# Differentiable Visual Attention

DRAW (Gregor et. al. 2015) uses a grid of Gaussian filters to **read** from input images and **draw** to a **canvas** image:

- $\sigma^2$  — filter variance
- $g_X, g_Y$  — grid centre
- $\delta$  — grid stride
- $\gamma$  — intensity



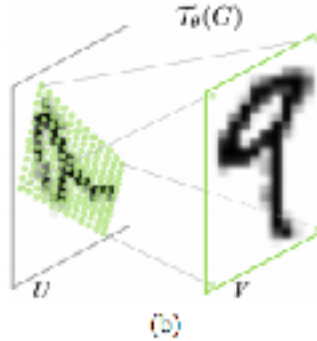


Reading MNIST

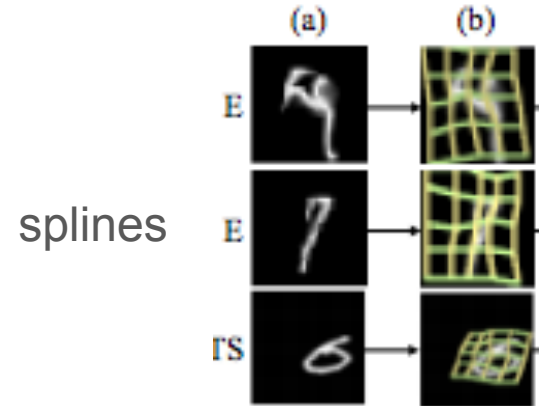
# Spatial Transformers

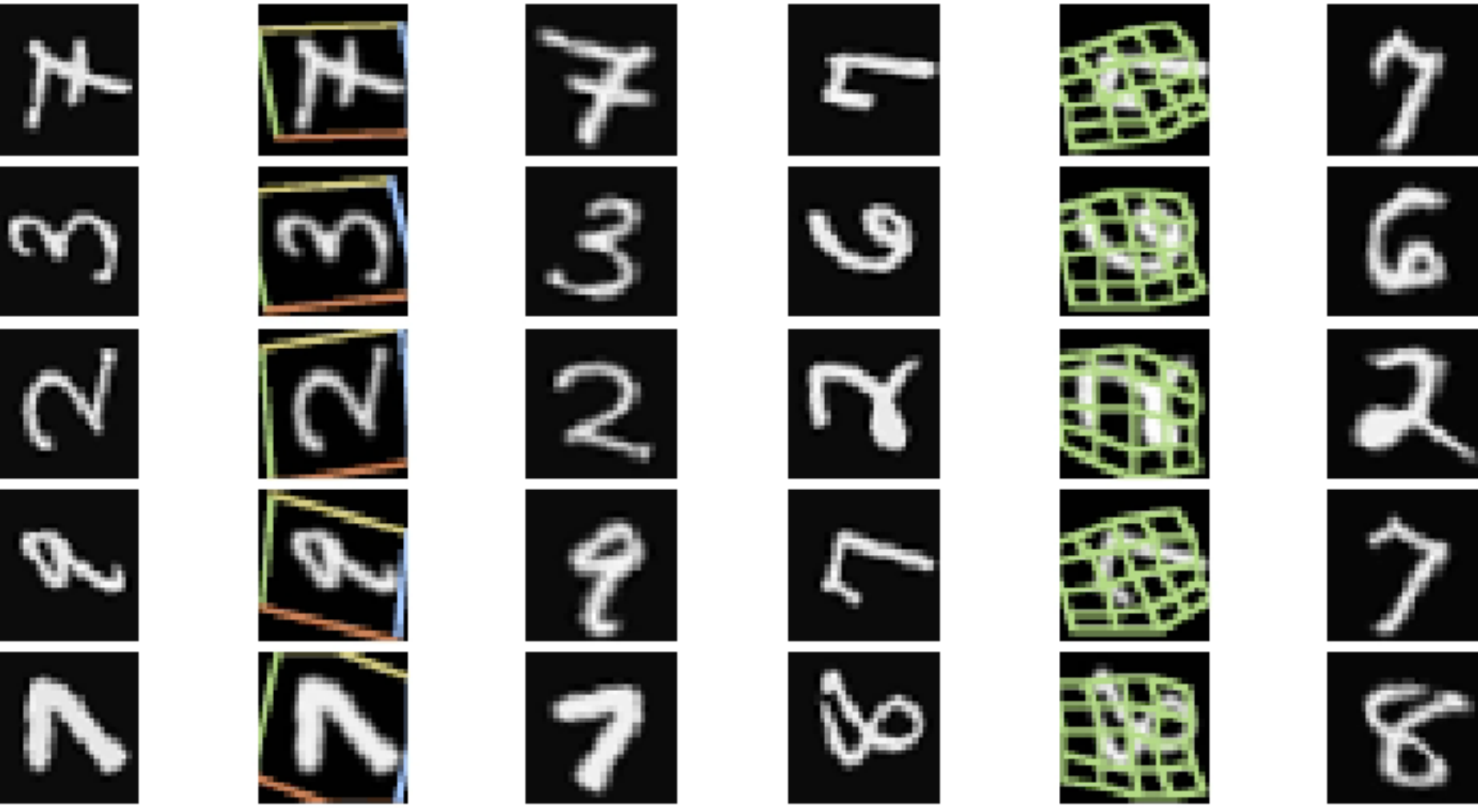
The attention mechanism in **Spatial Transformer Networks** (Jaderberg et. al. 2015) extends DRAW by allowing more complex transformations of the **sampling grid**

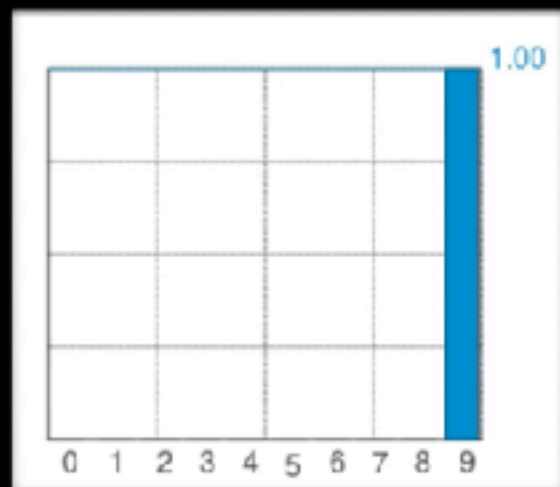
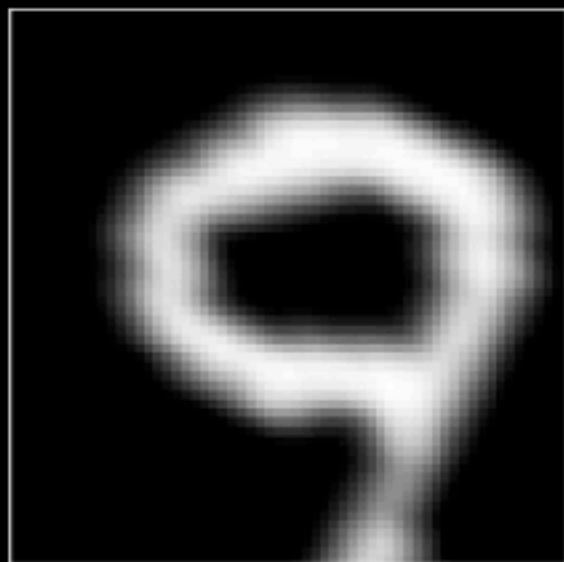
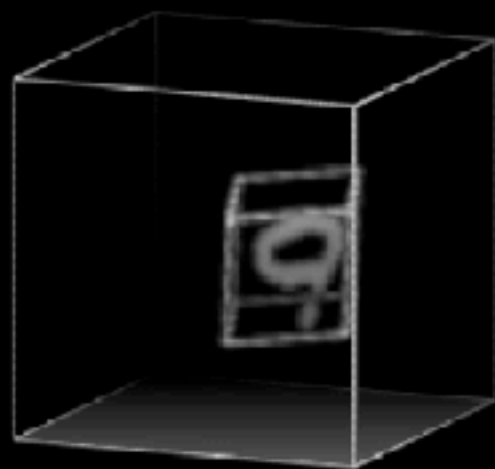
E.g. affine transformations



Can also extend to 3d sampling grids







# Summary

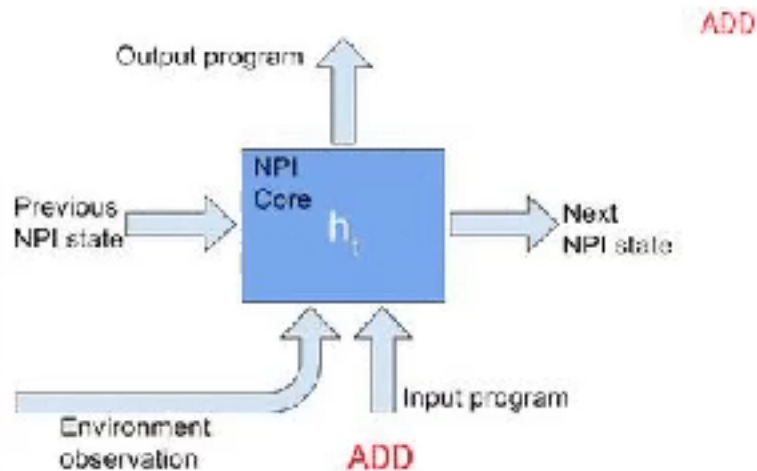
- Selective attention appears to be as useful for deep learning as it is for people
- Neural nets always have **implicit attention**, but we can also add **explicit attention** mechanisms
- These can be **stochastic** and trained with **reinforcement learning**
- Or **differentiable** and trained with ordinary **backprop**
- We can use attention to decide where to **write**, as well as where to **read**
- Many types of neural attention model (**content**, **location**, **visual**, **temporal**...) have been defined, and more are appearing all the time

### Addition scratch pad

	4	8	0	2	8	3	8	4	8	*
+	3	9	2	8	4	9	0	5	2	*
										*

### NPI inference

### Generated commands



*Neural Programmer-Interpreters*, Reed and de Freitas (2015)