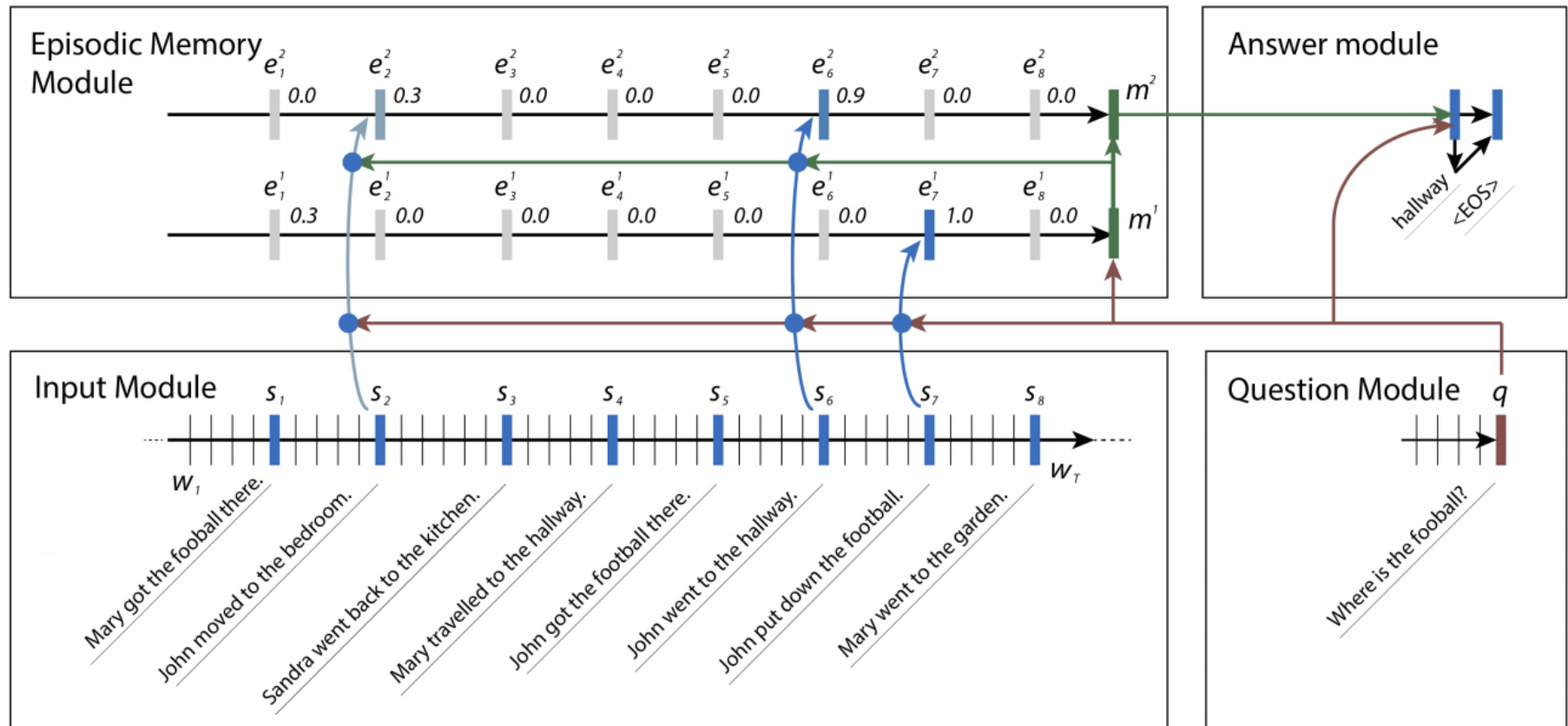


# Neural Architectures with Memory

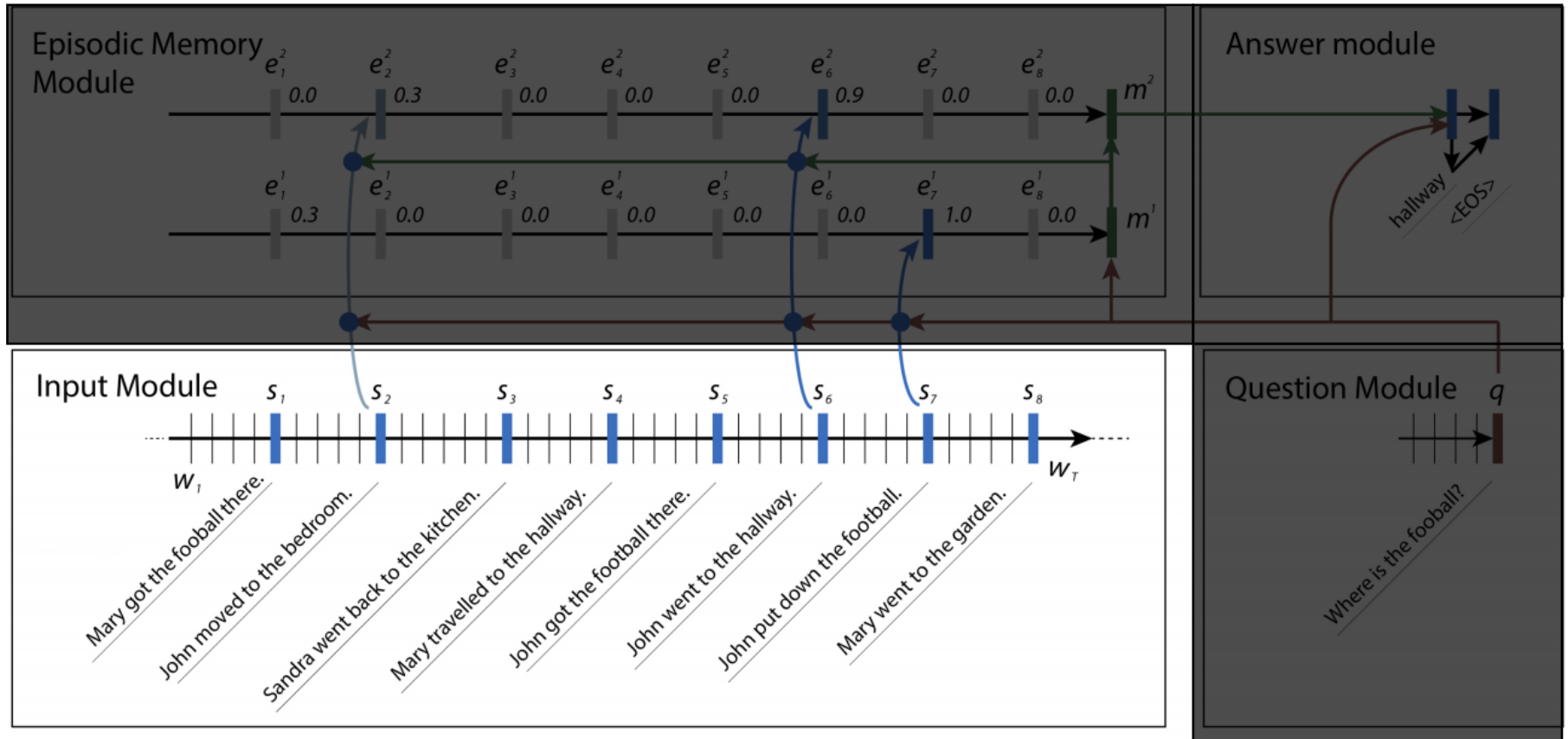
By Mohit Kumar(DeepMind)

# Dynamic Memory Networks - The Beast



Use RNNs, specifically GRUs for every module

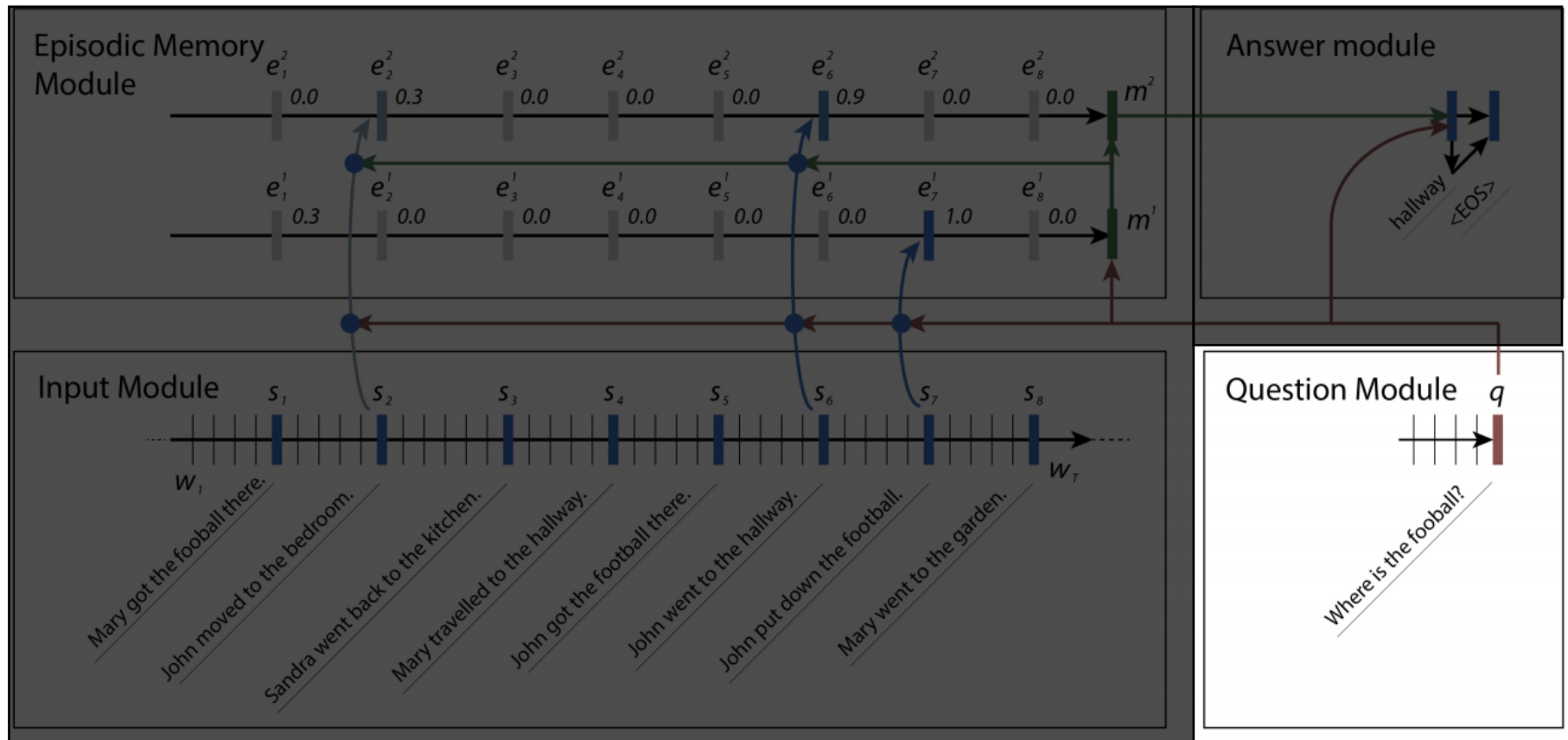
# DMN



Final GRU  
Output for  
sentence

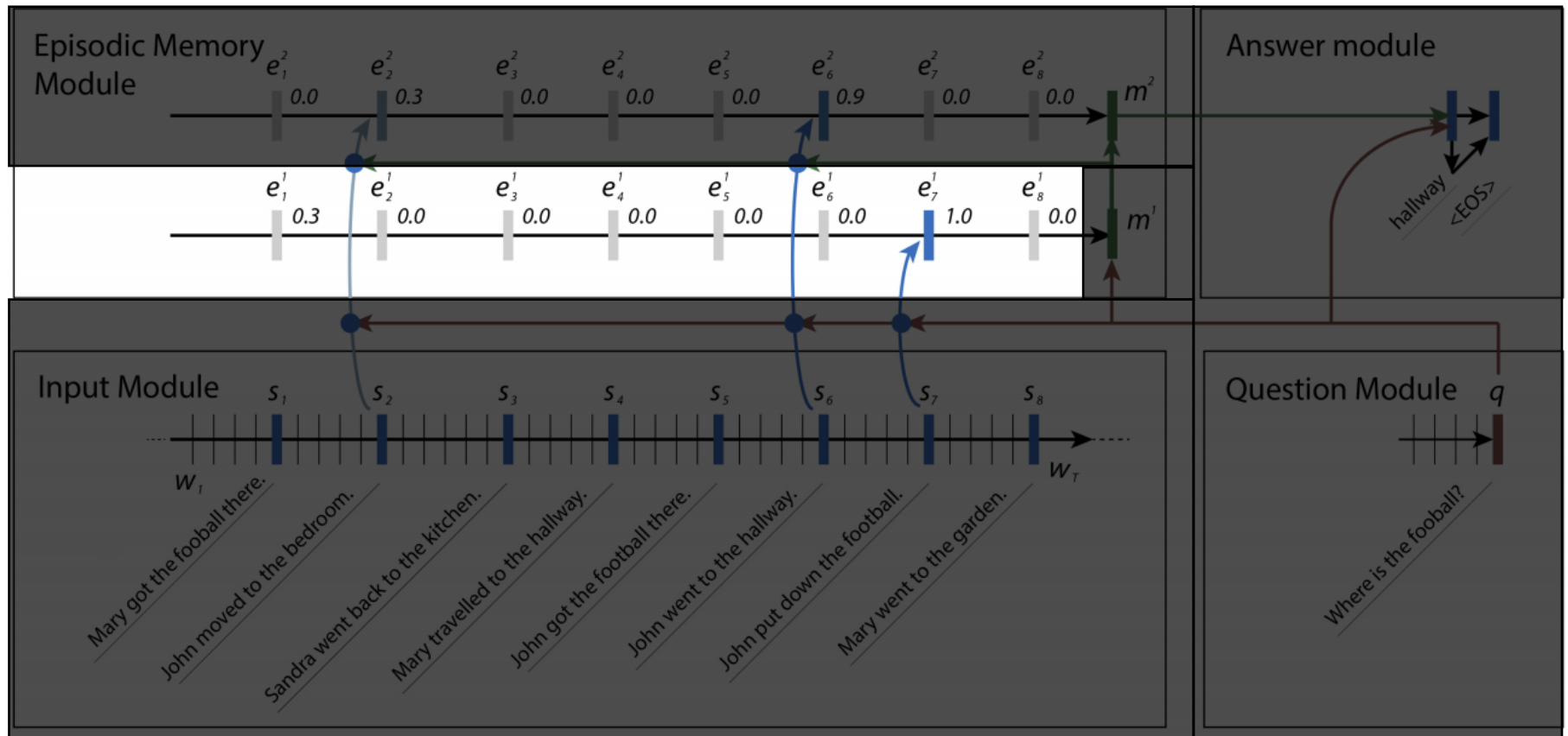
$$c_t = \text{GRU}(w_t^i, c_t^{i-1})$$

# DMN



$$q = \text{GRU}(q_w^i, q^{i-1})$$

# DMN



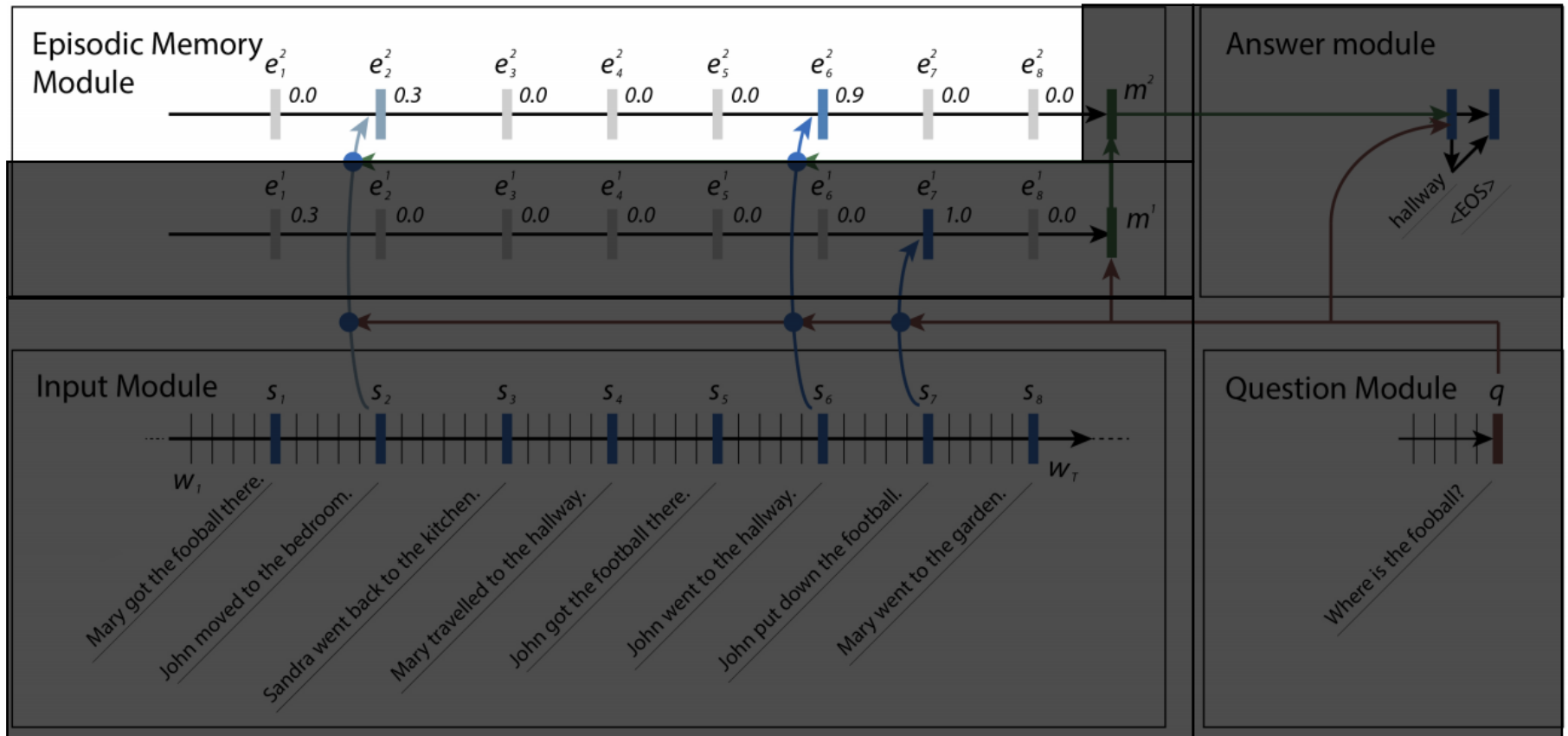
$$\text{Hop} = i$$

$$i = 1$$

$$h_t^i = g_t^i \text{GRU}(c_t, h_{t-1}^i) + (1 - g_t^i) h_{t-1}^i$$

$$e^i = h_{T_C}^i$$

# DMN



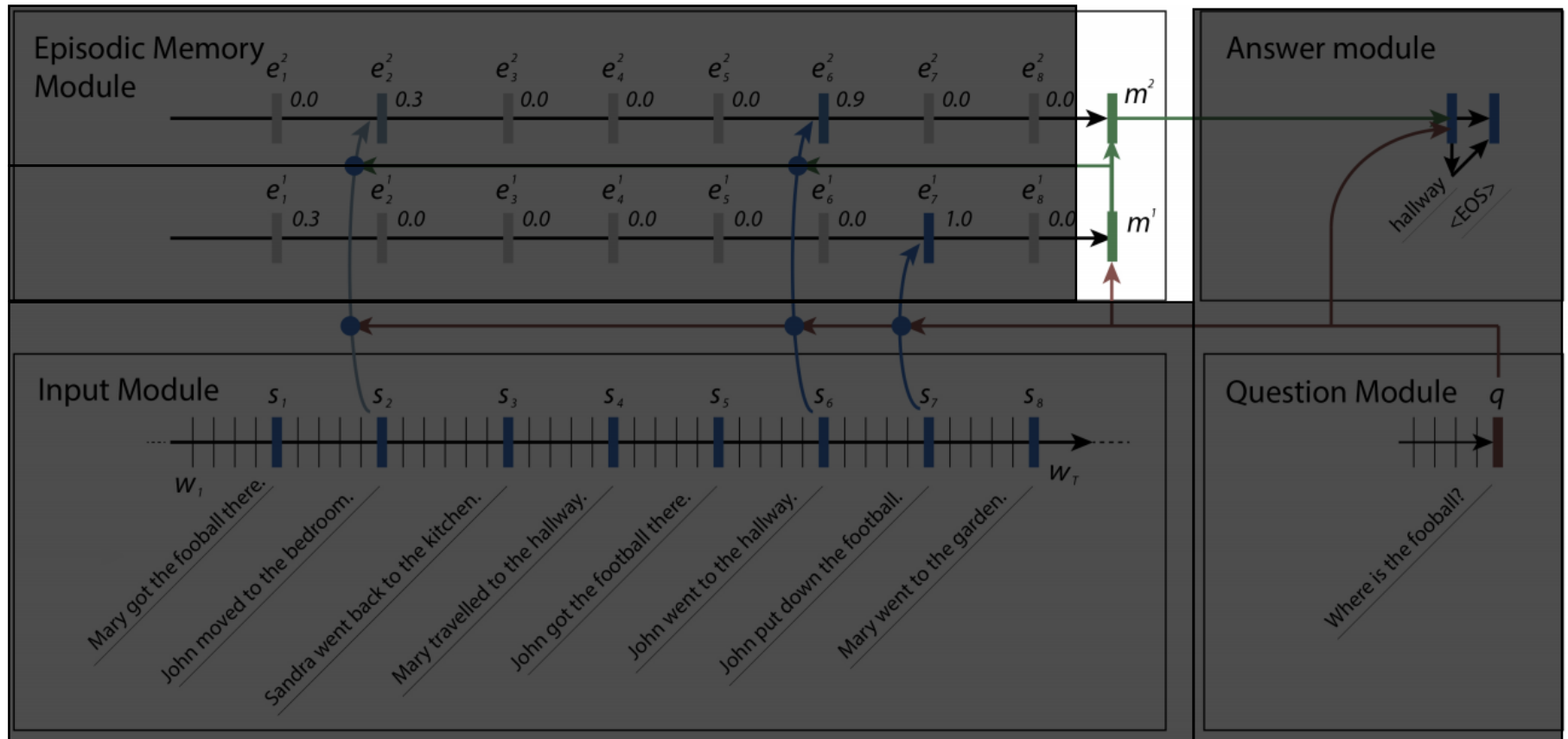
$$\text{Hop} = i$$

$$i = 2$$

$$h_t^i = g_t^i \text{GRU}(c_t, h_{t-1}^i) + (1 - g_t^i) h_{t-1}^i$$

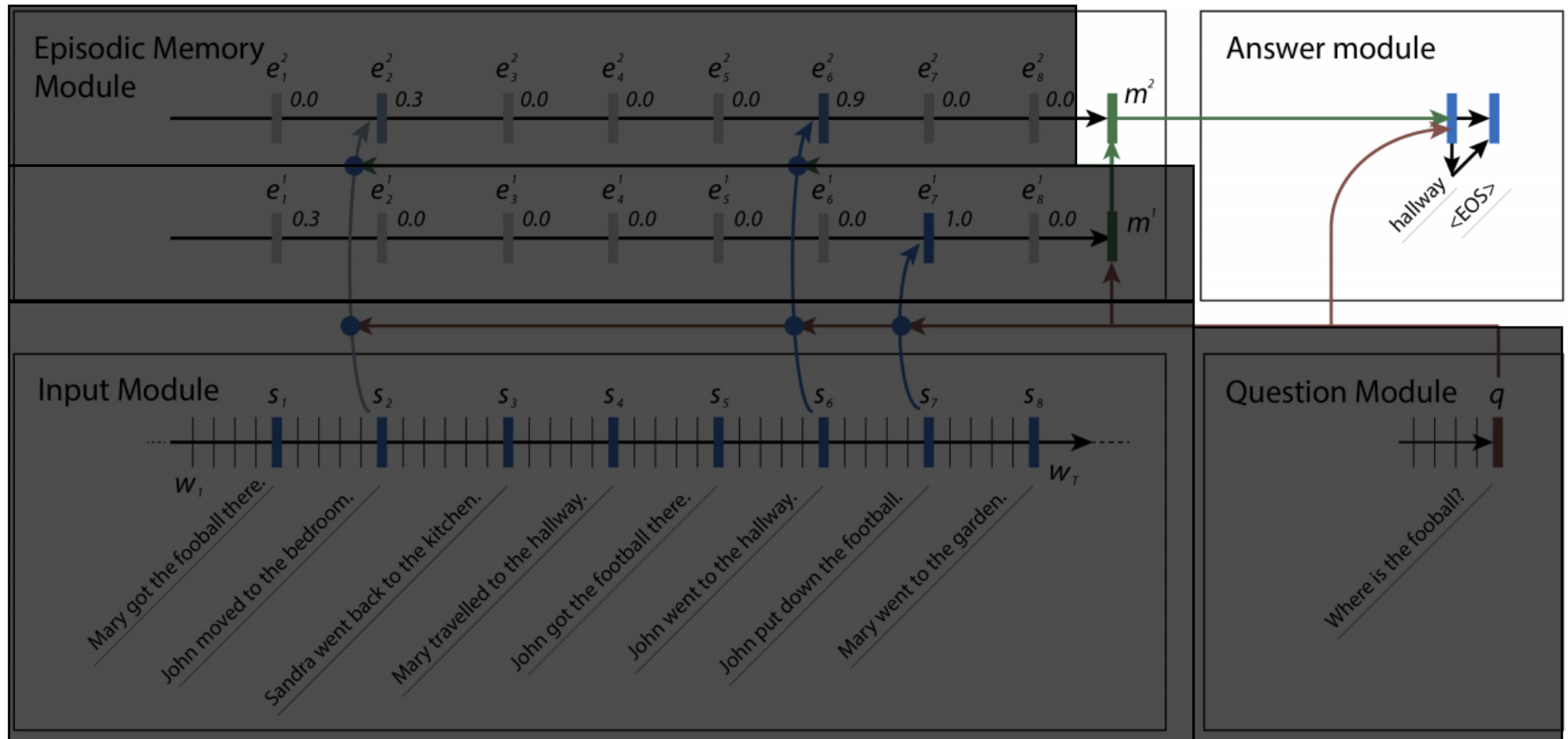
$$e^i = h_{T_C}^i$$

# DMN



$$m^i = \text{GRU}(e^i, m^{i-1})$$

# DMN



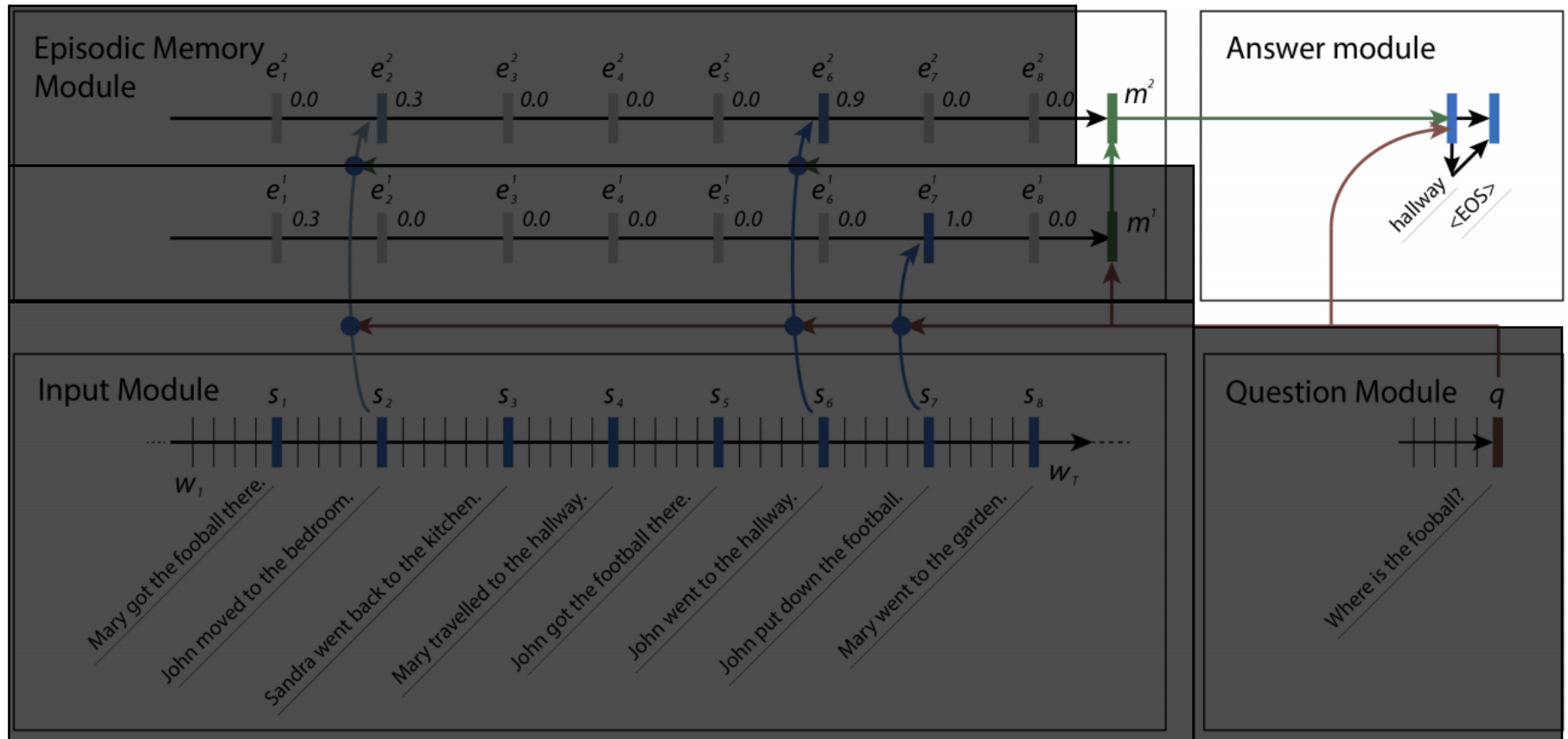
$$y_t = \text{Softmax}(W^{(a)} \alpha_t)$$

$$\alpha_0 = m^{T_m}$$

$$\alpha_t = \text{GRU}([y_{t-1}, q], \alpha_{t-1})$$



# DMN

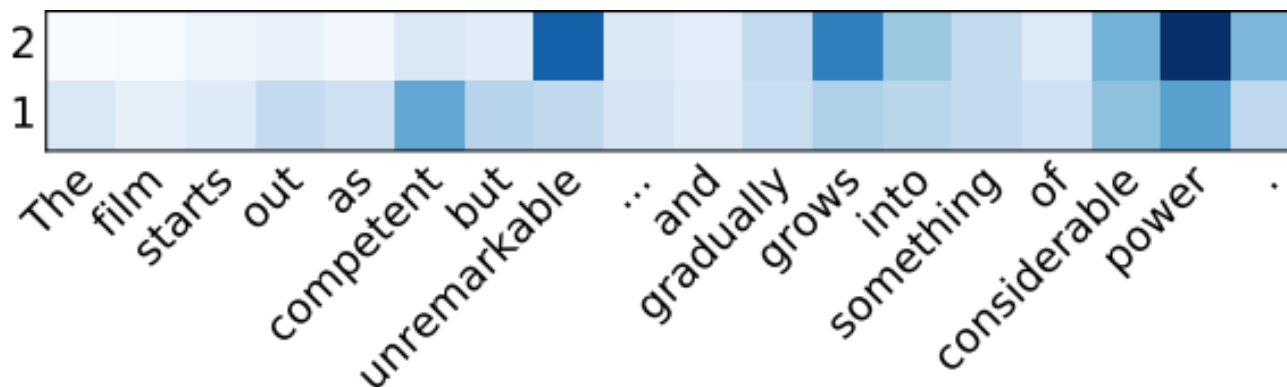


How many GRUs were used with 2 hops?

# DMN – Qualitative Results

**Question:** Where was Mary before the Bedroom?  
**Answer:** Cinema.

Facts	Episode 1	Episode 2	Episode 3
Yesterday Julie traveled to the school.			
Yesterday Marie went to the cinema.			
This morning Julie traveled to the kitchen.			
Bill went back to the cinema yesterday.			
Mary went to the bedroom this morning.			
Julie went back to the bedroom this afternoon.			
[done reading]			



# Algorithm Learning

# Neural Turing Machine

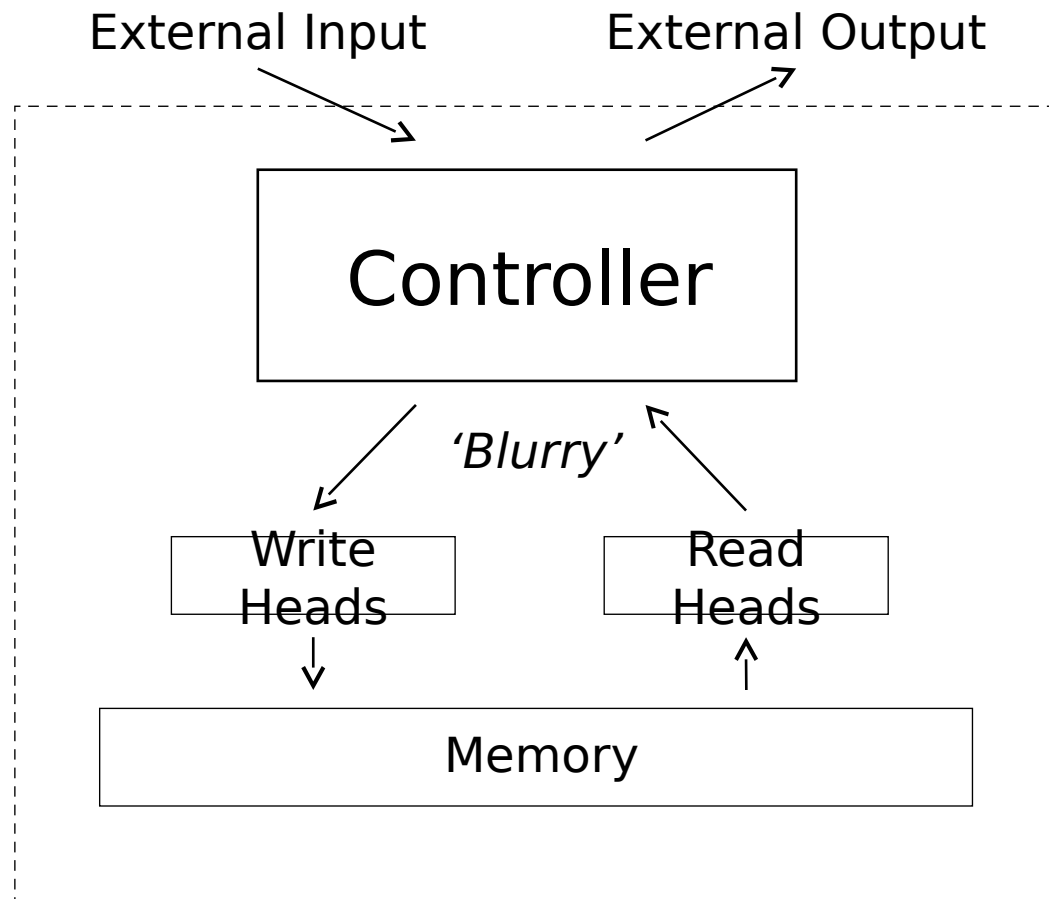
## Copy Task: Implement the Algorithm

Given a list of numbers at input, reproduce the list at output

### Neural Turing Machine Learns:

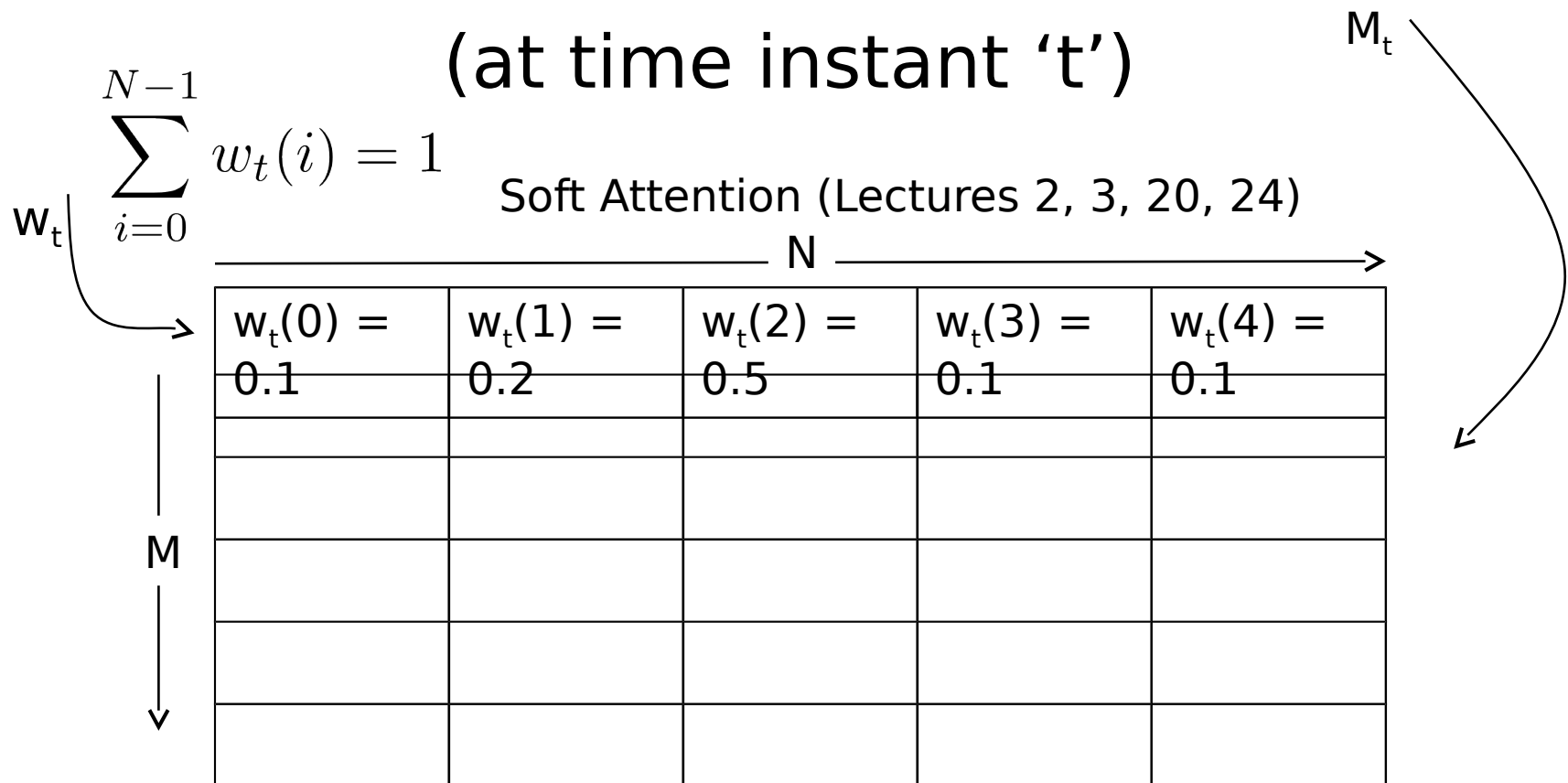
1. What to write to memory
2. When to write to memory
3. When to stop writing
4. Which memory cell to read from
5. How to convert result of read into final output

# Neural Turing Machines



# Neural Turing Machines

## ‘Blurry’ Memory Addressing (at time instant ‘t’)



# Neural Turing Machines

More formally,

## **Blurry Read Operation**

Given:  $M_t$  (memory matrix) of size  $N \times M$

$w_t$  (weight vector) of length  $N$

$t$  (time index)

$$r_t = \sum_{i=0}^{N-1} w_t(i) \mathbf{M}_t(i)$$

# Neural Turing Machines: Blurry Writes

## Blurry Write Operation

Decomposed into blurry erase + blurry add

Given:  $M_t$  (memory matrix) of size  $N \times M$

$w_t$  (weight vector) of length  $N$

$t$  (time index)

$e_t$  (erase vector) of length  $M$

$$M_t(i) \mathbf{a}_t = \underbrace{M_{t-1}(i) (1 - w_t(i) e_t)}_{\text{Erase Component}} + \underbrace{w_t(i) \mathbf{a}_t}_{\text{Add Component}}$$



# Neural Turing Machines: Erase

$$\mathbf{M}_t(i) = \mathbf{M}_{t-1}(i)(1 - w_t(i)\mathbf{e}_t)$$

1xN

$\mathbf{M}_0 \Rightarrow$

$w_1(0) =$	$w_1(1) =$	$w_1(2) =$	$w_1(3) =$	$w_1(4) =$
0.1	0.2	0.5	0.1	0.1
5	7	9	2	12
11	6	3	1	2
3	7	3	10	6
4	2	5	9	9
3	5	12	8	4

$\mathbf{e}_1$

1.0
0.7
0.2
0.5
0.0

Mx1

# Neural Turing Machines: Erase

$$\mathbf{M}_t(i) = \mathbf{M}_{t-1}(i)(1 - w_t(i)\mathbf{e}_t)$$

$w_1(0) =$	$w_1(1) =$	$w_1(2) =$	$w_1(3) =$	$w_1(4) =$
0.1	0.2	0.5	0.1	0.1
4.5	5.6	4.5	1.8	10.8
10.23	5.16	1.95	0.93	1.86
2.94	6.72	2.7	9.8	5.88
3.8	1.8	3.75	8.55	8.55
3	5	12	8	4

# Neural Turing Machines: Addition

$$\mathbf{M}_t(i) = \mathbf{M}_{t-1}(i)(1 - w_t(i)\mathbf{e}_t) + w_t(i)\mathbf{a}_t$$

$w_1(0) =$	$w_1(1) =$	$w_1(2) =$	$w_1(3) =$	$w_1(4) =$	$\mathbf{a}_1$
0.1	0.2	0.5	0.1	0.1	
4.5	5.6	4.5	1.8	10.8	3
10.23	5.16	1.95	0.93	1.86	4
2.94	6.72	2.7	9.8	5.88	-2
3.8	1.8	3.75	8.55	8.55	0
3	5	12	8	4	2

# Neural Turing Machines: Blurry Writes

$$\mathbf{M}_t(i) = \mathbf{M}_{t-1}(i)(1 - w_t(i)\mathbf{e}_t) + w_t(i)\mathbf{a}_t$$

$\mathbf{M}_1 \Rightarrow$

4.8	6.2	6	2.1	11.1
10.63	5.96	3.95	1.33	2.26
2.74	6.32	1.7	9.6	5.68
3.8	1.8	3.75	8.55	8.55
3.2	5.4	13	8.2	4.2

# Neural Turing Machines: Attention Model

## Generating $w_t$

### Content Based

Example: QA Task

- Score sentences by similarity with Question
- Weights as softmax of similarity scores

### Location Based

Example: Copy Task

- Move to address  $(i+1)$  after writing to index  $(i)$
- Weights  $\approx$  Transition probabilities

# Neural Turing Machine: Attention Model

Prev. State

$\mathbf{M}_t$

$w_{t-1}$

Controller

Outputs

$\mathbf{k}_t$

$\beta_t$

$g_t$

$s_t$

$\gamma_t$

$\mathbf{C}$   
 $\mathbf{A}$

$w_t^c$

$\mathbf{I}$

$w_t^g$

$\mathbf{CS}$

$\hat{w}_t$

$\mathbf{S}$

$w_t$

## Steps for generating

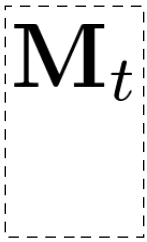
$\mathbf{W}_t$

1. Content Addressing
2. Peaking
3. Interpolation
4. Convolutional Shift (Location Addressing)
5. Sharpening

# Neural Turing Machine: Attention Model

Prev. State

$\mathbf{M}_t$

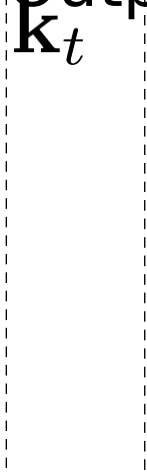


$\mathbf{k}_t$  : Vector (length M) produced by Controller

Controller

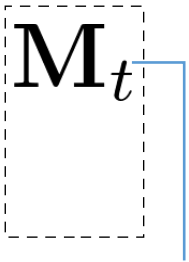
Outputs

$\mathbf{k}_t$



# Neural Turing Machine: Attention Model

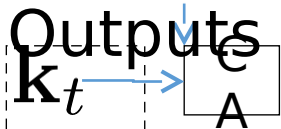
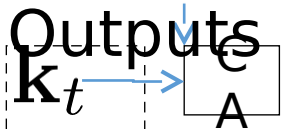
Prev. State



Step 1: Content Addressing

$$\frac{(CA)}{w_t^c(i)} = \frac{\exp \langle \mathbf{M}_t(i), \mathbf{k}_t \rangle}{\sum_i \exp \langle \mathbf{M}_t(i), \mathbf{k}_t \rangle}$$

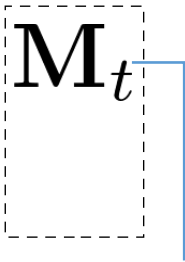
Controller





# Neural Turing Machine: Attention Model

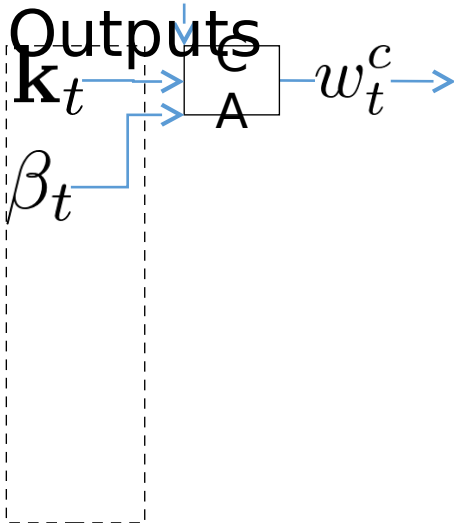
Prev. State



Step 2: Peaking

$$w_t^c(i) = \frac{\exp(\beta_t(\langle \mathbf{M}_t(i), \mathbf{k}_t \rangle))}{\sum_i \exp(\beta_t(\langle \mathbf{M}_t(i), \mathbf{k}_t \rangle))}$$

Controller



# Neural Turing Machine: Attention Model

Prev. State

$\mathbf{M}_t$

$w_{t-1}$

Controller

Outputs

$\mathbf{k}_t$

$\beta_t$

$g_t$

$\mathbf{C}$   
 $\mathbf{A}$

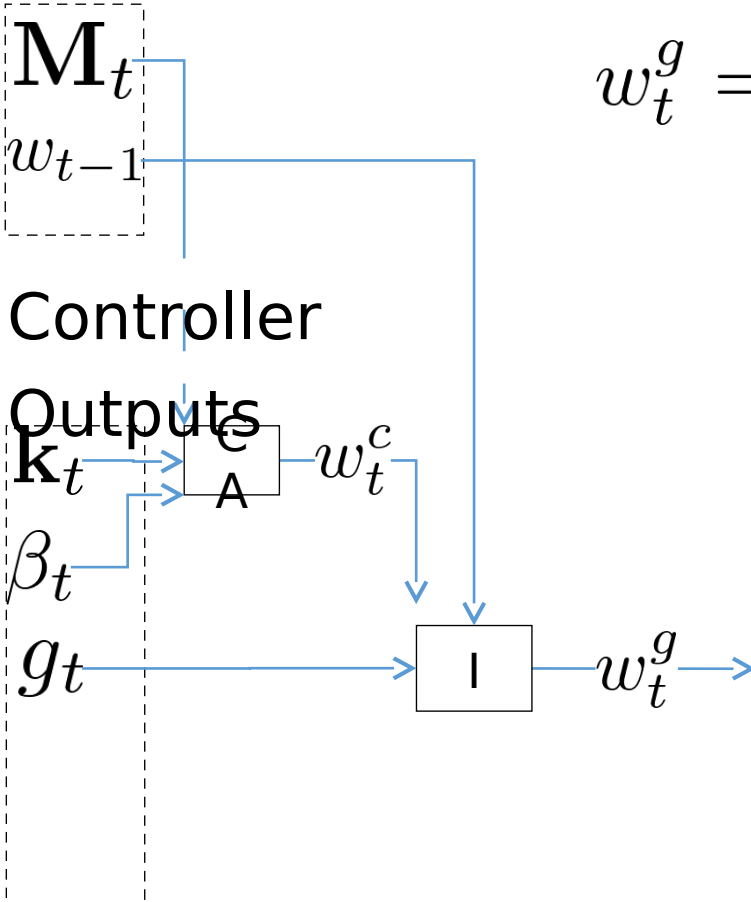
$w_t^c$

$\mathbf{I}$

Step 3: Interpolation (I)

$$w_t^g = g_t w_t^c + (1 - g_t) w_{t-1}$$

$w_t^g$



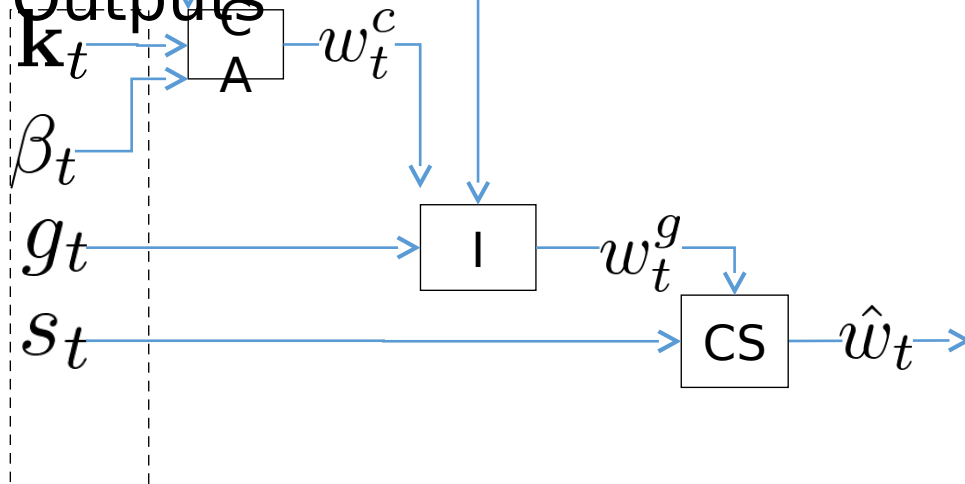
# Neural Turing Machine: Attention Model

Prev. State

$\mathbf{M}_t$   
 $w_{t-1}$

Controller

Outputs



Step 4: Convolutional Shift  
(CS)  $s_t$

- Controller outputs  $w_t^c$ , a normalized distribution over all N possible shifts
- Rotation-shifted weights computed as  $\sum_{j=0}^{N-1} w_t^c(j) \beta_t(j - i)$

# Neural Turing Machine: Attention Model

Prev. State

$\mathbf{M}_t$   
 $w_{t-1}$

Controller

Outputs

$\mathbf{k}_t$   
 $\mathbf{A}$

$\beta_t$

$g_t$

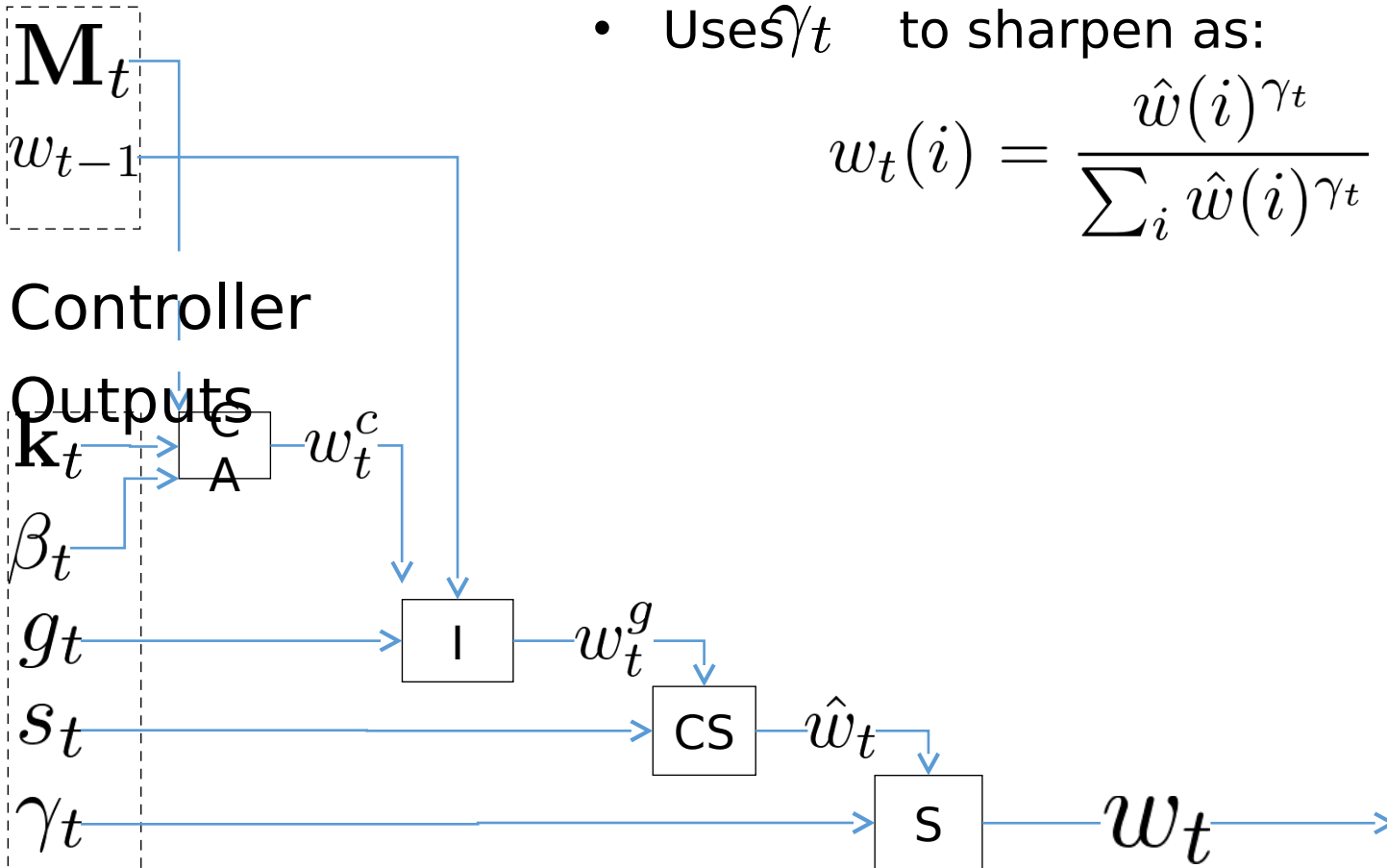
$s_t$

$\gamma_t$

Step 5: Sharpening (S)

- Uses  $\hat{w}_t$  to sharpen as:

$$w_t(i) = \frac{\hat{w}(i)^{\gamma_t}}{\sum_i \hat{w}(i)^{\gamma_t}}$$



# Neural Turing Machine: Controller Design

- Feed-forward: faster, more transparency & interpretability about function learnt
- LSTM: more expressive power, doesn't limit the number of computations per time step

Both are end-to-end differentiable!

1. Reading/Writing -> Convex Sums
2.  $w_t$  generation -> Smooth
3. Controller Networks

# Neural Turing Machine: Network Overview

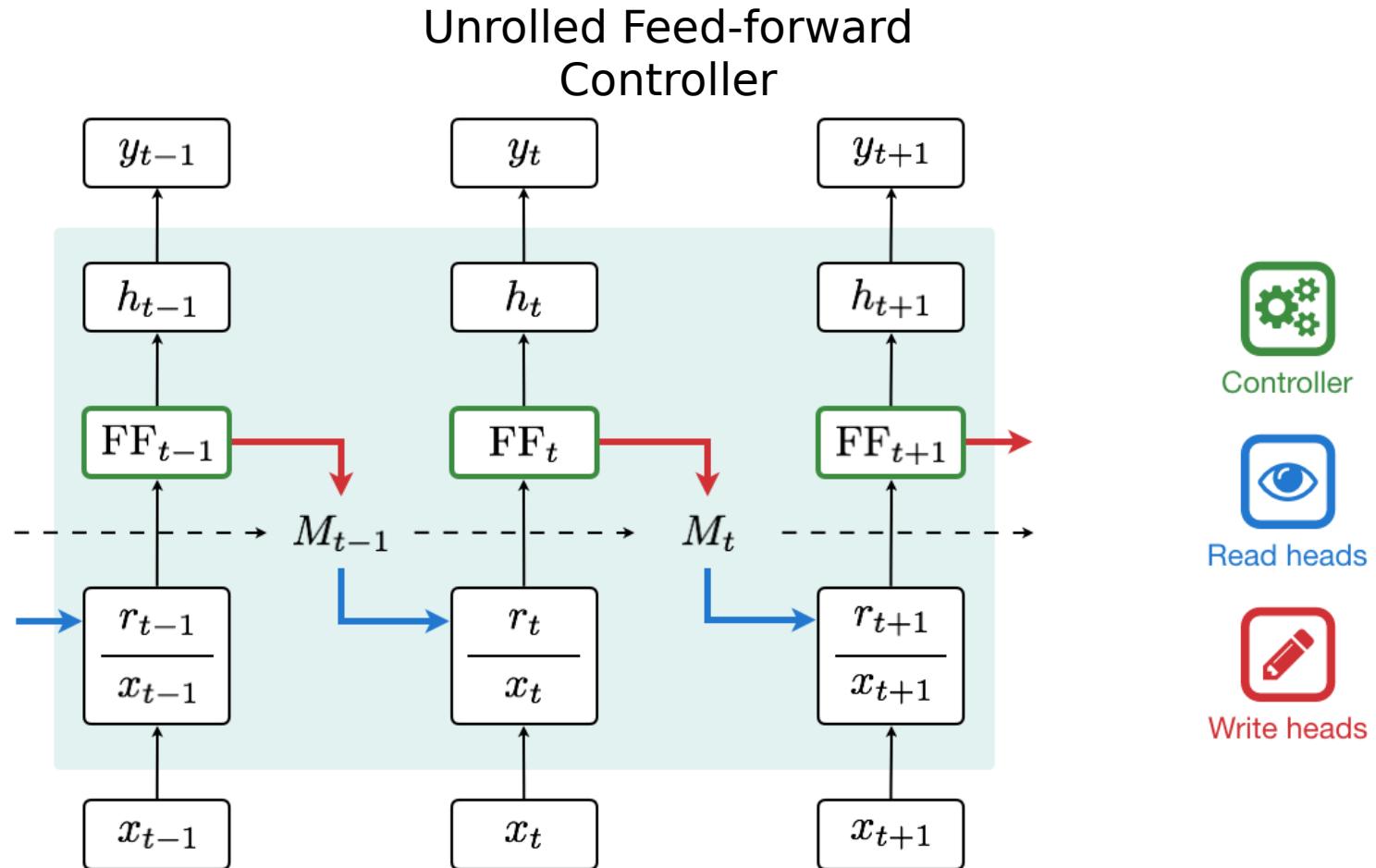


Figure from  
[Snips AI's Medium Post](#)

# Neural Turing Machines vs. MemNNs

## MemNNs

- Memory is static, with focus on retrieving (reading) information from memory

## NTMs

- Memory is continuously written to and read from, with network learning when to perform memory read and write

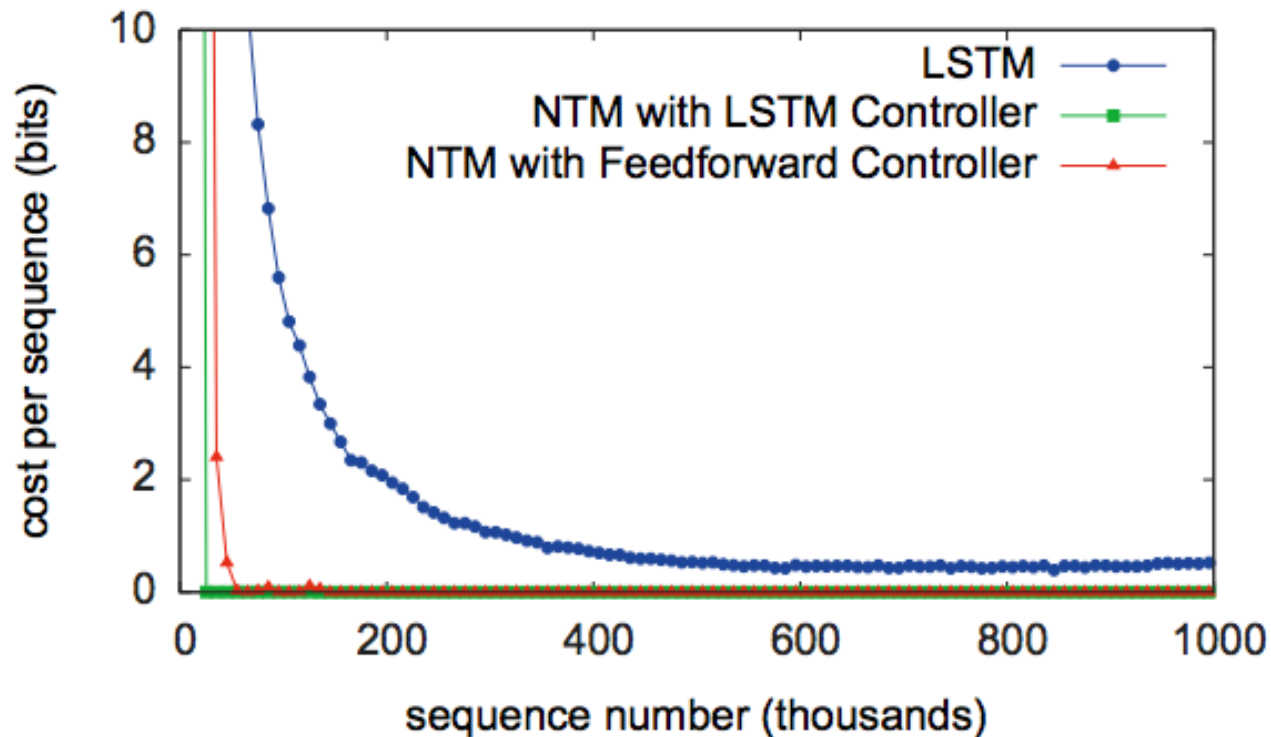
# Neural Turing Machines: Experiments

<b>Task</b>	<b>Network Size</b>		<b>Number of Parameters</b>	
	<b>NTM w/ LSTM*</b>	<b>LSTM</b>	<b>NTM w/ LSTM</b>	<b>LSTM</b>
Copy	3 x 100	3 x 256	67K	1.3M
Repeat Copy	3 x 100	3 x 512	66K	5.3M
Associative	3 x 100	3 x 256	70K	1.3M
N-grams	3 x 100	3 x 128	61K	330K
Priority Sort	2 x 100	3 x 128	269K	385K

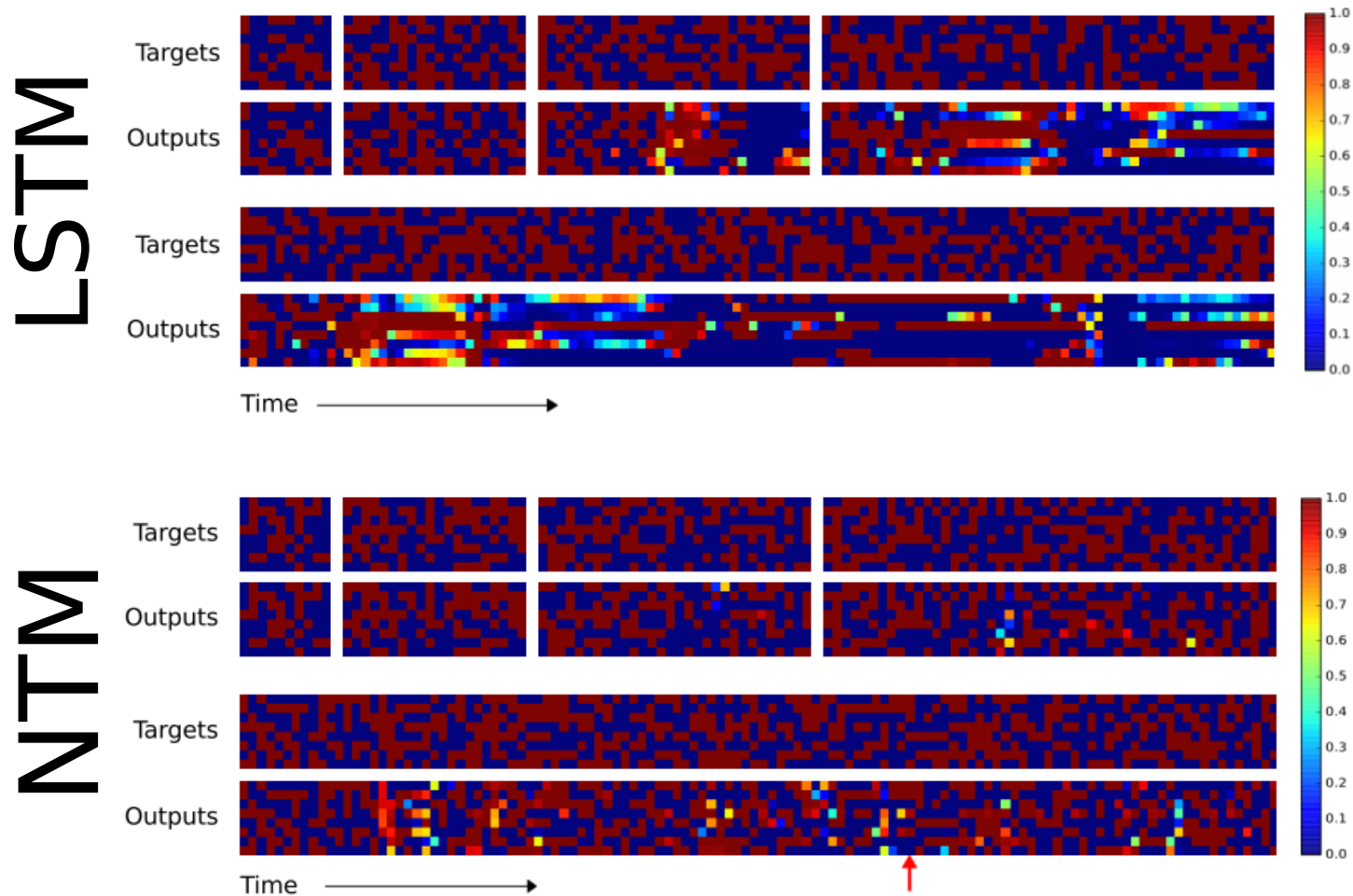


# Neural Turing Machines: 'Copy' Learning Curve

Trained on 8-bit sequences,  $1 \leq \text{sequence length} \leq 20$



# Neural Turing Machines: 'Copy' Performance



Neural Turing Machines  
triggered an outbreak of  
Memory Architectures!

# Dynamic Neural Turing Machines

## Experimented with addressing schemes

- Dynamic Addresses: Addresses of memory locations learnt in training – allows non-linear location-based addressing
- Least recently used weighting: Prefer least recently used memory locations + interpolate with content-based addressing
- Discrete Addressing: Sample the memory location from the content-based distribution to obtain a one-hot address
- Multi-step Addressing: Allows multiple hops over memory

	<b>Location NTM</b>	<b>Content NTM</b>	<b>Soft DNTM</b>	<b>Discrete DNTM</b>
1- step	31.4%	33.6%	29.5%	27.9%
3- step	32.8%	32.7%	24.2%	21.7%

# Stack Augmented Recurrent Networks

Learn algorithms based on stack implementations (e.g. learning fixed sequence

⌋

Sequence generator	Example
$\{a^n b^n \mid n > 0\}$	aab <b>b</b> aaab <b>bb</b> ab <b>aaaa</b> ab <b>bbbb</b>
$\{a^n b^n c^n \mid n > 0\}$	aaab <b>bbccc</b> ab <b>caaaa</b> ab <b>bbbbcccc</b>
$\{a^n b^n c^n d^n \mid n > 0\}$	aab <b>bccdd</b> aaab <b>bbcccd</b> ddab <b>cd</b>
$\{a^n b^{2n} \mid n > 0\}$	aab <b>bbba</b> aaab <b>bbbbbb</b> abb
$\{a^n b^m c^{n+m} \mid n, m > 0\}$	aab <b>ccca</b> aaab <b>cccc</b> ab <b>cc</b>
$n \in [1, k], X \rightarrow nXn, X \rightarrow =$	$(k = 2) 12=$ <b>21</b> $2122=$ <b>2212</b> $11121=$ <b>12111</b>

Uses a stack data structure to store memory  
(as opposed to a memory matrix)

# Stack Augmented Recurrent Networks

- Blurry ‘push’ and ‘pop’ on stack. E.g.:

$$s_t[0] = a[\text{Push}](h_t) + a[\text{Pop}]s_{t-1}[1]$$

- Some results:

method	$a^n b^n$	$a^n b^n c^n$	$a^n b^n c^n d^n$	$a^n b^{2n}$	$a^n b^m c^{n+m}$
RNN	25%	23.3%	13.3%	23.3%	33.3%
LSTM	100%	100%	68.3%	75%	100%
List RNN 40+5	100%	33.3%	100%	100%	100%
Stack RNN 40+10	100%	100%	100%	100%	43.3%
Stack RNN 40+10 + rounding	100%	100%	100%	100%	100%

# Differentiable Neural Computers

## Advanced addressing mechanisms:

- Content Based Addressing
- Temporal Addressing
  - Maintains notion of sequence in addressing
  - Temporal Link Matrix  $L$  (size  $N \times N$ ),  $L[i, j]$  = degree to which location  $i$  was written to after location  $j$ .
- Usage Based Addressing

# DNC: Usage Based Addressing

- Writing increases usage of cell, reading decreases usage of cell
- Least used location has highest usage-based weighting
- Interpolate b/w usage & content based weights for final write weights



# DNC: Example



Hybrid computing using a neural network with dynamic external memory, Graves et. al., Nature vol. 538

# DNC: Improvements over NTMs

## NTM

- Large contiguous blocks of memory needed
- No way to free up memory cells after writing

## DNC

- Memory locations non-contiguous, usage-based
- Regular de-allotment based on usage-tracking

# DNC: Experiments

## Graph Tasks

Graph Representation: (source, edge, destination) tuples

Types of tasks:

- Traversal: Perform walk on graph given source, list of edges
- Shortest Path: Given source, destination
- Inference: Given source, relation over edges; find destination

# DNC: Experiments

## Graph Tasks

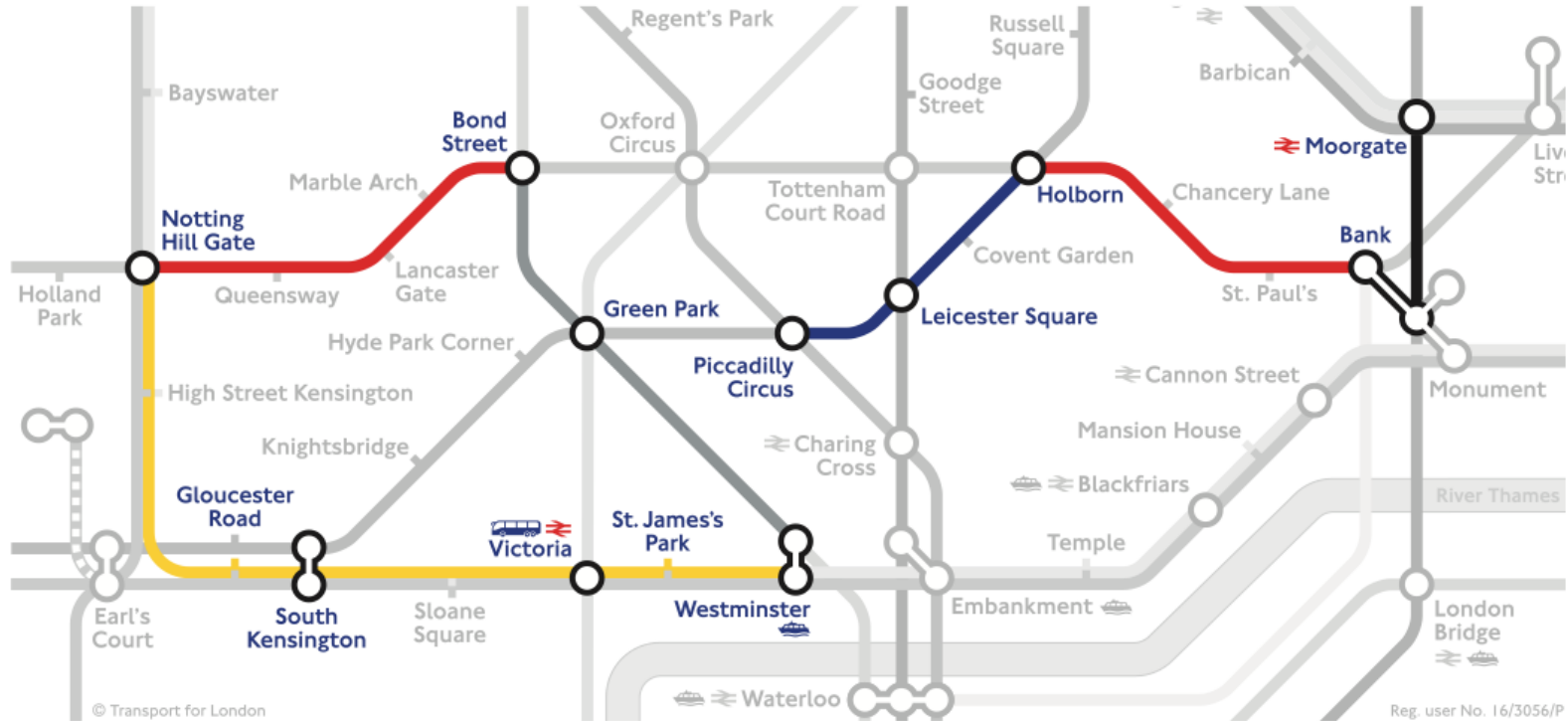
Training over 3 phases:

- Graph description phase: (source, edge, destination) tuples fed into the graph
- Query phase: Shortest path (source, \_\_\_\_, destination), Inference (source, hybrid relation, \_\_\_\_), Traversal (source, relation, relation ..., \_\_\_\_)
- Answer phase: Target responses provided at output

Trained on random graphs of maximum size 1000

# DNC: Experiments

## Graph Tasks: London Underground



# DNC: Experiments

## Graph Tasks: London Underground

Input Phase

(OxfordCircus, TottenhamCtRd, Central)  
(TottenhamCtRd, OxfordCircus, Central)  
(BakerSt, Marylebone, Circle)  
(BakerSt, Marylebone, Bakerloo)  
(BakerSt, OxfordCircus, Bakerloo)  
⋮  
(LeicesterSq, CharingCross, Northern)  
(TottenhamCtRd, LeicesterSq, Northern)  
(OxfordCircus, PiccadillyCircus, Bakerloo)  
(OxfordCircus, NottingHillGate, Central)  
(OxfordCircus, Euston, Victoria)

# DNC: Experiments

## Graph Tasks: London Underground **Traversal Task**

(BondSt, \_, Central),  
(\_, \_, Circle), (\_, \_, Circle),  
(\_, \_, Circle), (\_, \_, Circle),  
(\_, \_, Jubilee), (\_, \_, Jubilee),

Query Phase

(BondSt, NottingHillGate, Central)  
(NottingHillGate, GloucesterRd, Circle)  
⋮  
(Westminster, GreenPark, Jubilee)  
(GreenPark, BondSt, Jubilee)

Answer Phase

# DNC: Experiments

Graph Tasks: London Underground

## Shortest Path Task

(Moorgate, PiccadillyCircus, \_)

(Moorgate, Bank, Northern)

(Bank, Holborn, Central)

(Holborn, LeicesterSq, Piccadilly)

(LeicesterSq, PiccadillyCircus, Piccadilly)

Query Phase

Answer Phase



# DNC: Experiments

Graph Tasks: Freya's Family Tree



# Conclusion

- Machine Learning models require memory and multi-hop reasoning to perform AI tasks better
- Memory Networks for Text are an interesting direction but very simple
- Generic architectures with memory, such as Neural Turing Machine, limited applications shown
- Future directions should be focusing on applying generic neural models with memory to more AI Tasks.

# Reading List

- Karol Kurach, Marcin Andrychowicz & Ilya Sutskever **Neural Random-Access Machines**, ICLR, 2016
- Emilio Parisotto & Ruslan Salakhutdinov **Neural Map: Structured Memory for Deep Reinforcement Learning**, ArXiv, 2017
- Pritzel et. al. **Neural Episodic Control**, ArXiv, 2017
- Oriol Vinyals, Meire Fortunato, Navdeep Jaitly **Pointer Networks**, ArXiv, 2017
- Jack W Rae et al., **Scaling Memory-Augmented Neural Networks with Sparse Reads and Writes**, ArXiv 2016
- Antoine Bordes, Y-Lan Boureau, Jason Weston, **Learning End-to-End Goal-Oriented Dialog**, ICLR 2017
- Junhyuk Oh, Valliappa Chockalingam, Satinder Singh, Honglak Lee, **Control of Memory, Active Perception, and Action in Minecraft**, ICML 2016
- Wojciech Zaremba, Ilya Sutskever, **Reinforcement Learning Neural Turing Machines**, ArXiv 2016