

---

# Neural Networks with Memory

Memory Networks, End-to-End Memory Networks,  
Condensed Memory Networks, Neural Turing  
Machine & Differential Neural Computer

MemNN, E2E, C-MemNN, NTM, DNC, OMG !!

---

Aaditya Prakash  
March 27, 2017

---

# Neural Networks with Memory

Memory Networks, End-to-End Memory Networks,  
Condensed Memory Networks, Neural Turing  
Machine & Differential Neural Computer

MemNN, E2E, C-MemNN, NTM, DNC, OMG !!

Aaditya Prakash  
March 27, 2017



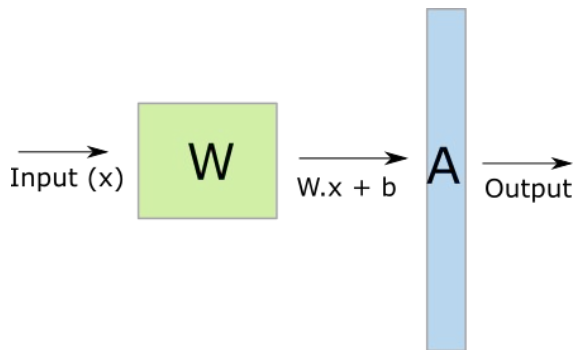
# Why?

Siri, play me a good song  
Sorry, I couldn't find a 'good song' in your music.



Because AI is good only with context !

# Neural Network - limitation



1. Not all problems can be mapped to  $y = f(x)$
2. Need to remember external context
3. Need to know where to look for in the context
4. Need to know what to look for in the context
5. Need to know how to reason using external context
6. Need to handle potentially changing external context

Text source: Sumit Chopra FAIR, DLSS 2016

# Why not RNN ?

1. Inherent temporal structure (not all knowledge is temporal)
2. Very slow (because of sequential learning, not embarrassingly parallel)
3. Cannot scale to large network (overhead in terms of learning)
4. Prone to “vanishing/exploding gradient” (multiply gradient many times)
5. Back-propagation in time is hard (always needs truncation and clipping)

# Reasoning is hard

Bilbo travelled to the cave. Gollum dropped the ring there. Bilbo took the ring.  
Bilbo went back to the Shire. Bilbo left the ring there. Frodo got the ring.  
Frodo journeyed to Mount-Doom. Frodo dropped the ring there. Sauron died.  
Frodo went back to the Shire. Bilbo travelled to the Grey-havens. The End.

Where is the ring?

Where is Bilbo now?

Where is Frodo now?

# Reasoning is hard

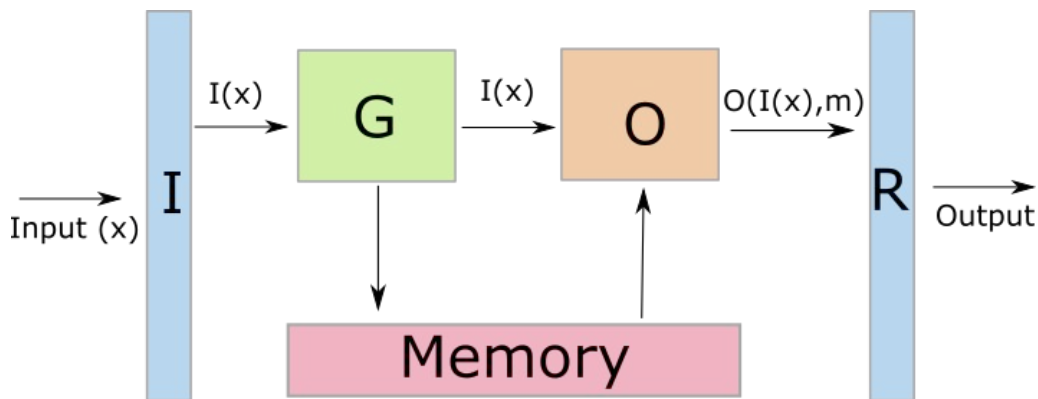
Bilbo travelled to the cave. Gollum dropped the ring there. Bilbo took the ring.  
Bilbo went back to the Shire. Bilbo left the ring there. Frodo got the ring.  
Frodo journeyed to Mount-Doom. Frodo dropped the ring there. Sauron died.  
Frodo went back to the Shire. Bilbo travelled to the Grey-havens. The End.

Where is the ring?      A: Mount-Doom

Where is Bilbo now?    A: Grey-havens

Where is Frodo now?    A: Shire

# Memory Networks



Convert the input (x) to internal representation (aka feature space)  
- learn embeddings (bow/gru/lstm), use pre-trained vectors (glove/w2v)

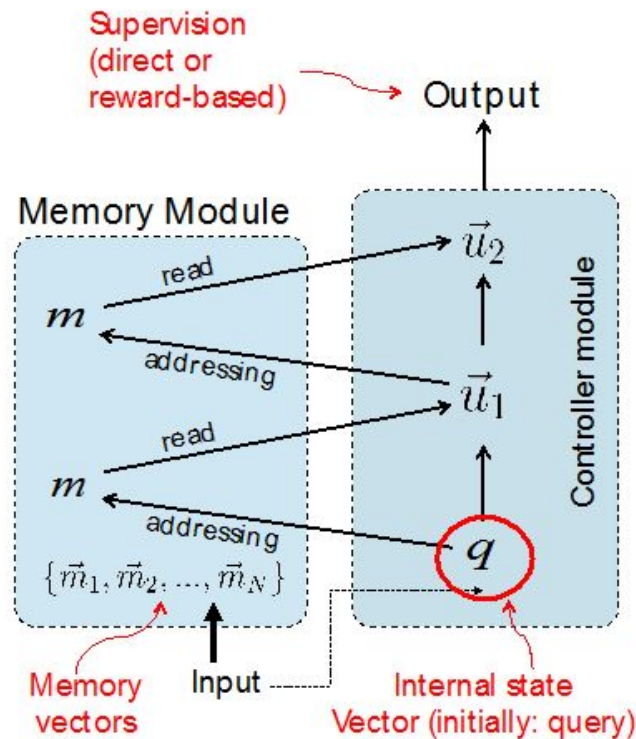
Make a memory state (u), write initial data to the memory, update when necessary

Generate the output state, which is combination of memory state and input

Convert the output (O) into response as desired by the model



# Memory Networks

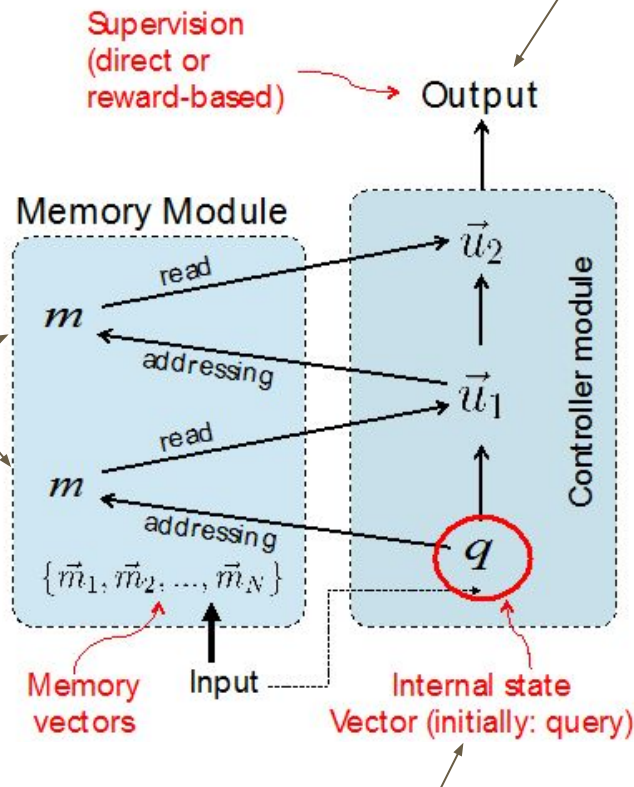


# Memory Networks

bABI tasks ---

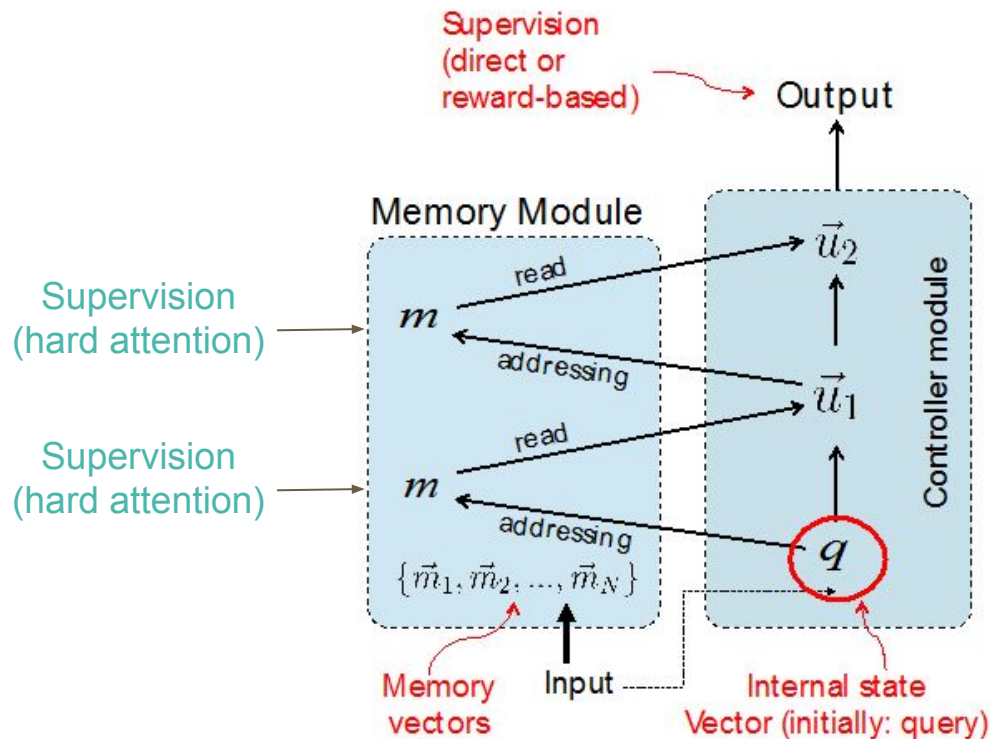
John is in the playground.  
Bob is in the office.  
John picked up the football.  
Bob went to the kitchen.

Where is the football? A:playground  
Where was Bob before the kitchen? A:office

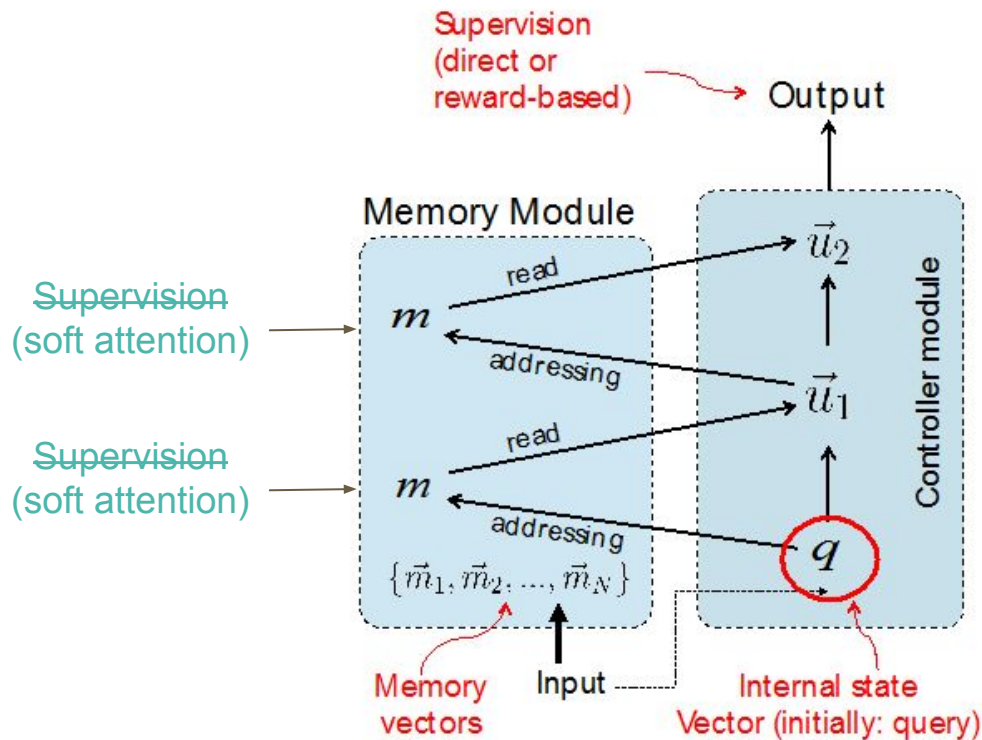


Where is the football?  
Where was Bob before the kitchen?

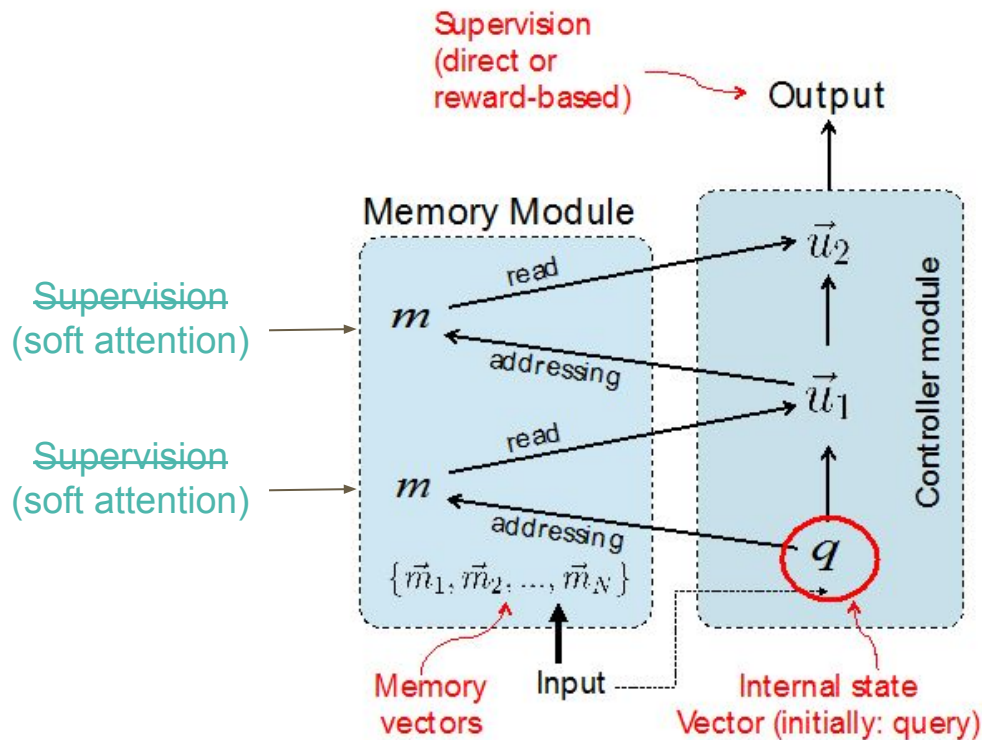
# Hard Attention (Strong supervision)



# Soft Attention (Weak supervision)

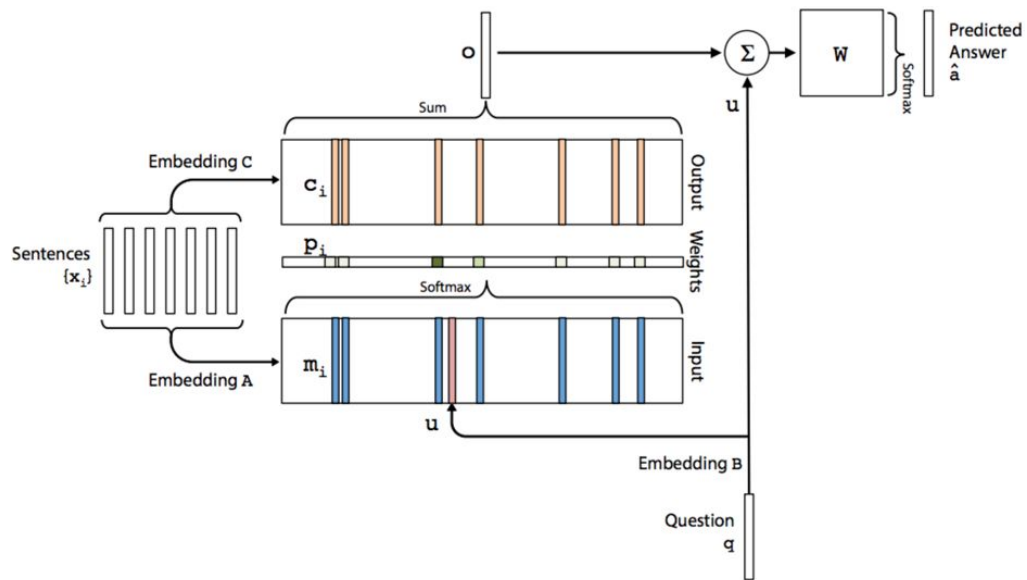


# Soft Attention (Weak supervision)

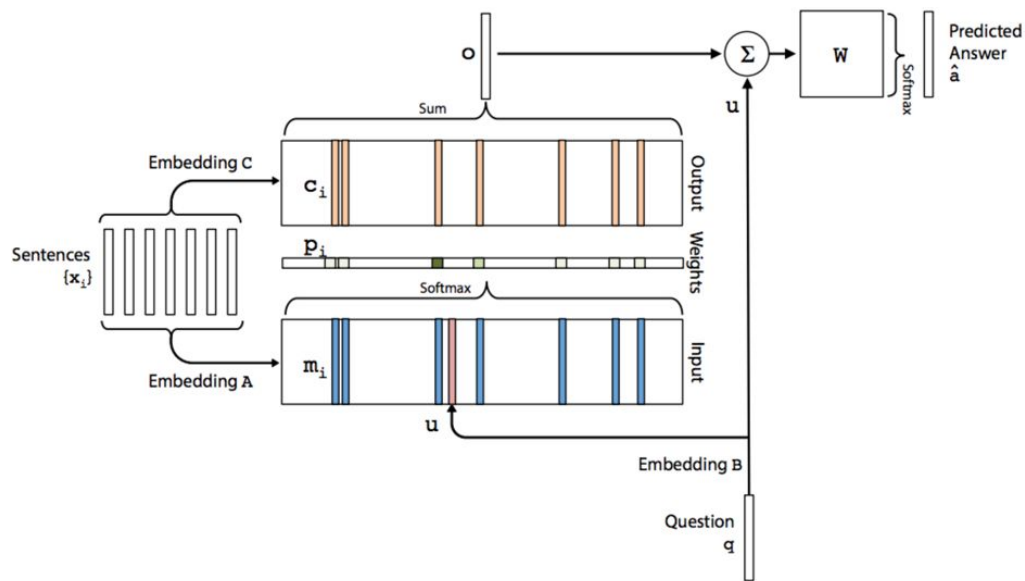


Now applicable to more problems where it is hard to get strong supervision

# End to End Memory Networks



# End to End Memory Networks



$$x_i = \{x_{i1}, x_{i2}, \dots, x_{in}\}$$

$$m_i = \sum_j A x_{ij}$$

$$p_i = \text{softmax}(u^T m_i)$$

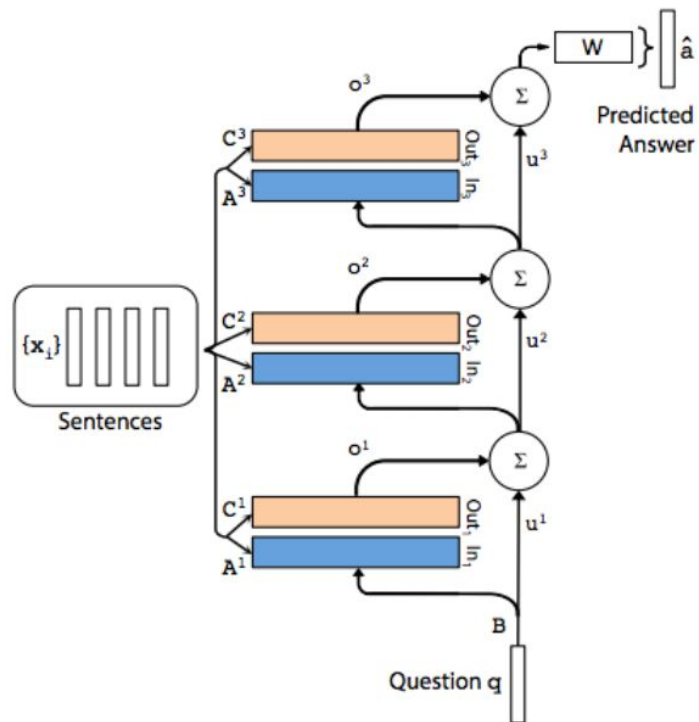
$$= \text{softmax}(q^T B^T \sum_j A x_{ij})$$

$$c_i = \sum_j C x_{ij}$$

$$o = \sum_i p_i c_i$$

$$a = \text{softmax}(W(o + u))$$

# One hop is just not sufficient



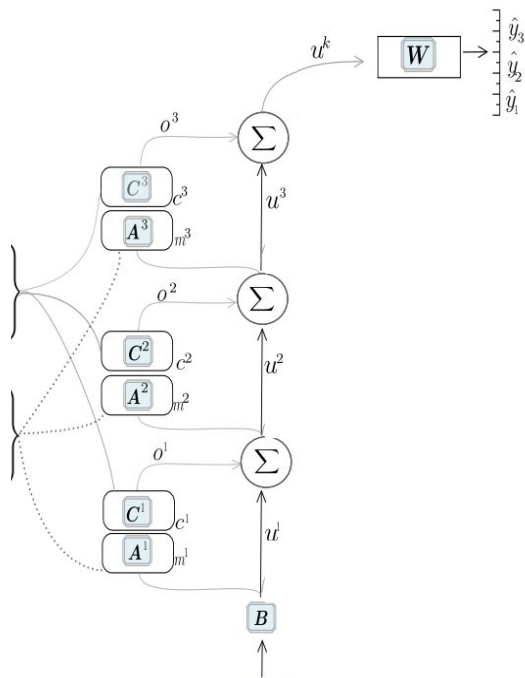
$$u^{k+1} = u^k + o^k$$

Each layer has its own embedding matrices  $A^k, C^k$

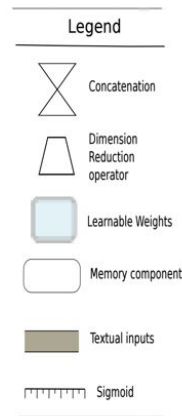
$$\hat{a} = \text{Softmax}(W u^{K+1}) = \text{Softmax}(W(o^K + u^K))$$



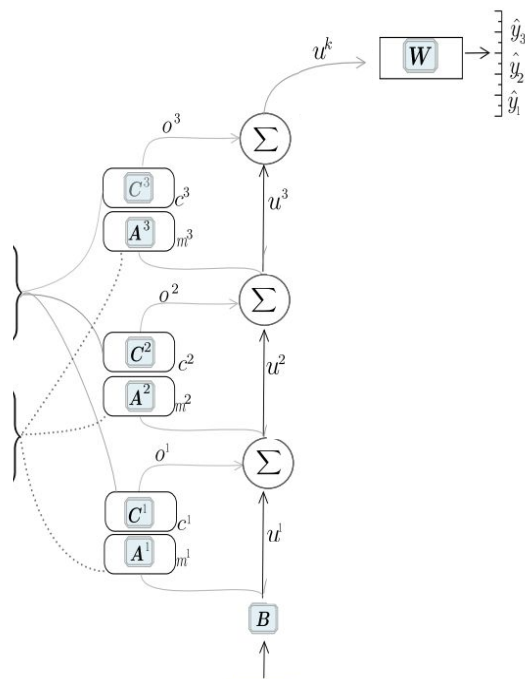
# Condensed Memory Network [ C-MemNN ]



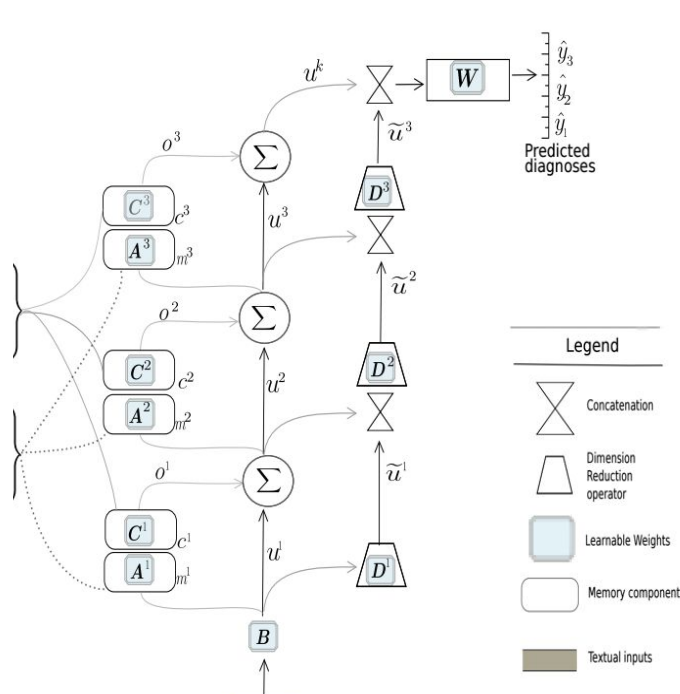
MemNN



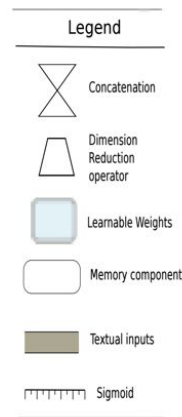
# Condensed Memory Network [ C-MemNN ]



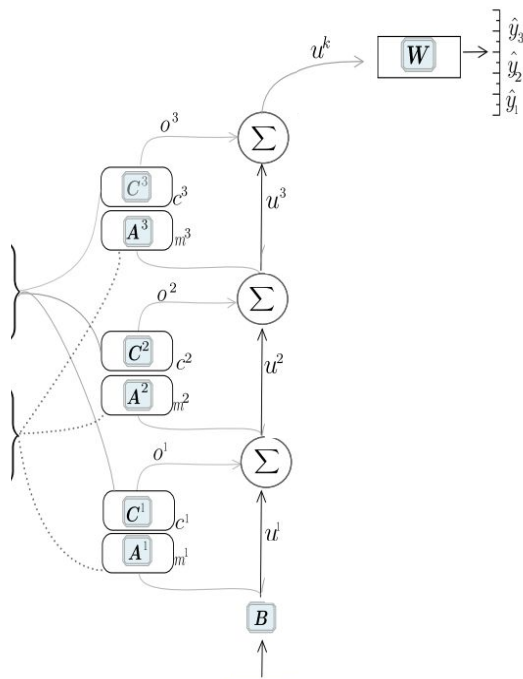
MemNN



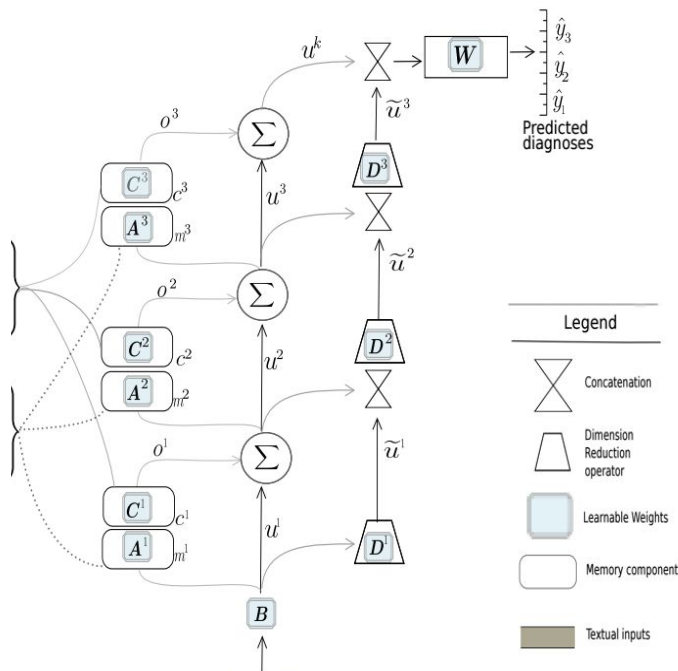
C-MemNN



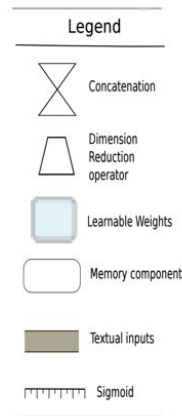
# Condensed Memory Network [ C-MemNN ]



MemNN



C-MemNN



## C-MemNN

$$o^k = \sum_i \text{Addressing}(u^k, m_i^k) \cdot c_i^k$$

$$u^{k+1} = u^k + o^k$$

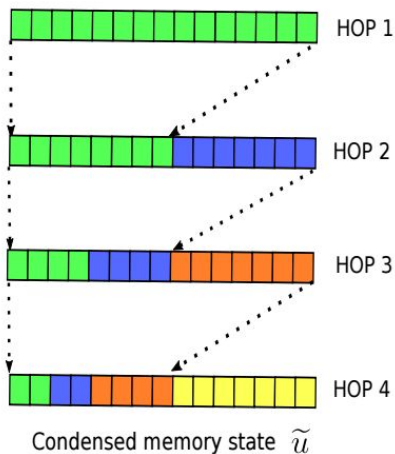
$$\tilde{u}^{k+1} = u^{k+1} \oplus D_1 \cdot \tilde{u}^k$$

$$\hat{y}_r = \arg \max_{r \in R} \frac{1}{1 + e^{-1 * (\tilde{u}^{k+1} \cdot W)}}$$

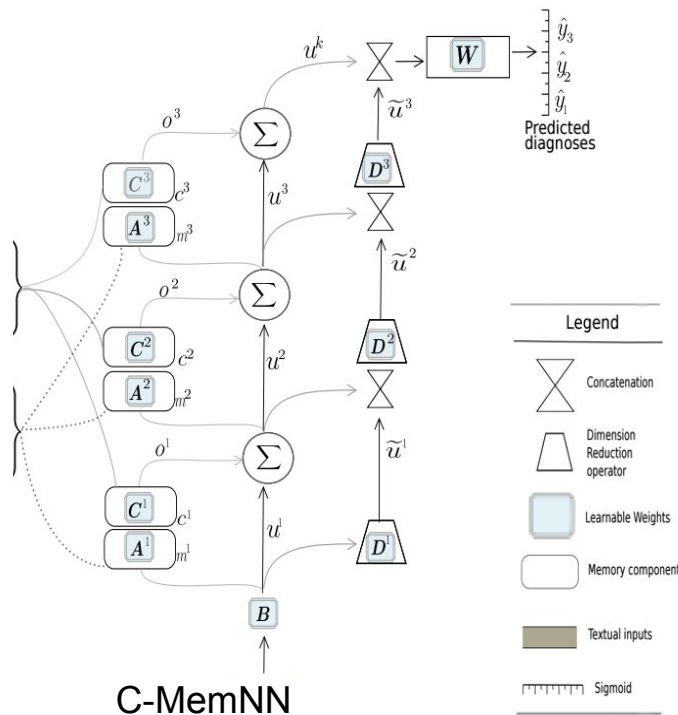
## A-MemNN

$$\tilde{u}^{k+1} = \tilde{u}^k + \frac{\tilde{u}^{k-1}}{2} + \frac{\tilde{u}^{k-2}}{4} + \dots$$

# Condensed Memory Network [ C-MemNN ]



Condensation of memory state - continuous hierarchy.



## C-MemNN

$$o^k = \sum_i \text{Addressing}(u^k, m_i^k) \cdot c_i^k$$

$$u^{k+1} = u^k + o^k$$

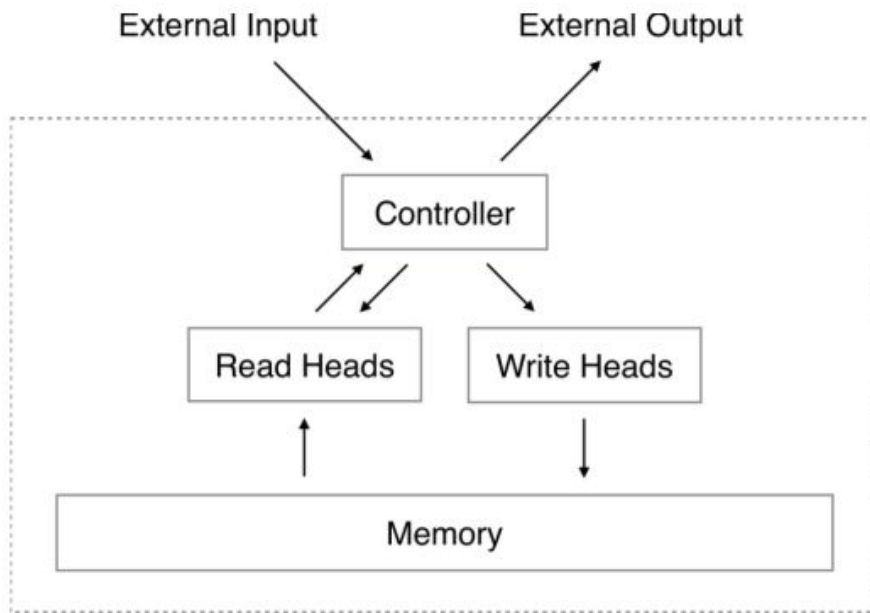
$$\tilde{u}^{k+1} = u^{k+1} \oplus D_1 \cdot \tilde{u}^k$$

$$\hat{y}_r = \arg \max_{r \in R} \frac{1}{1 + e^{-1 * (\tilde{u}^{k+1} \cdot W)}}$$

## A-MemNN

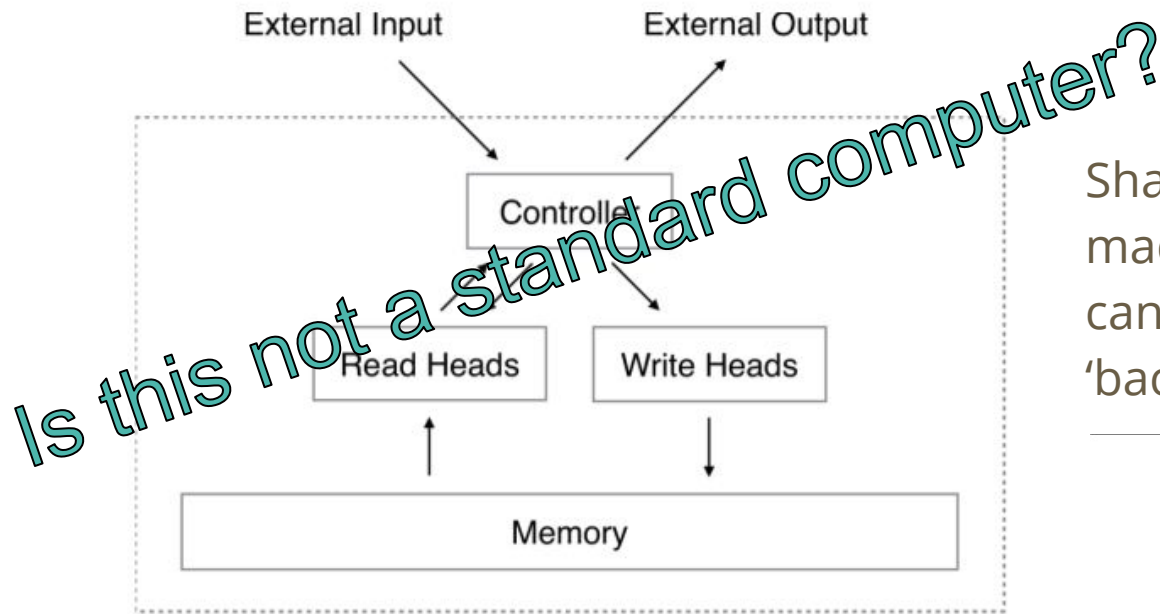
$$\tilde{u}^{k+1} = \tilde{u}^k + \frac{\tilde{u}^{k-1}}{2} + \frac{\tilde{u}^{k-2}}{4} + \dots$$

# Neural Turing Machine

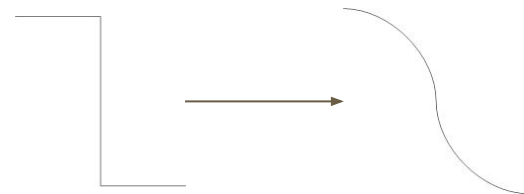


NTM can learn basic algorithms like sorting.

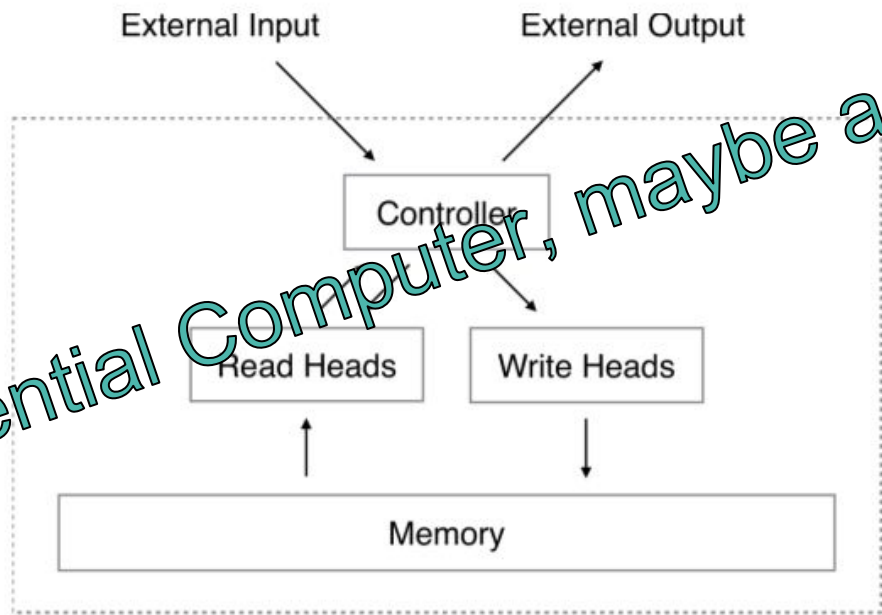
# Neural Turing Machine



Sharp functions  
made smooth. Now  
can be trained with  
'backpropagation'

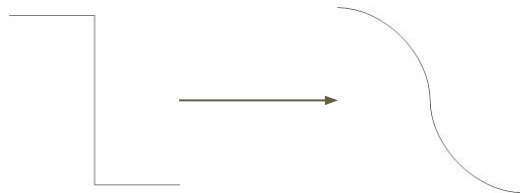


# Neural Turing Machine



So Differential Computer, maybe a better name?

Sharp functions  
made smooth. Now  
can be trained with  
'backpropagation'

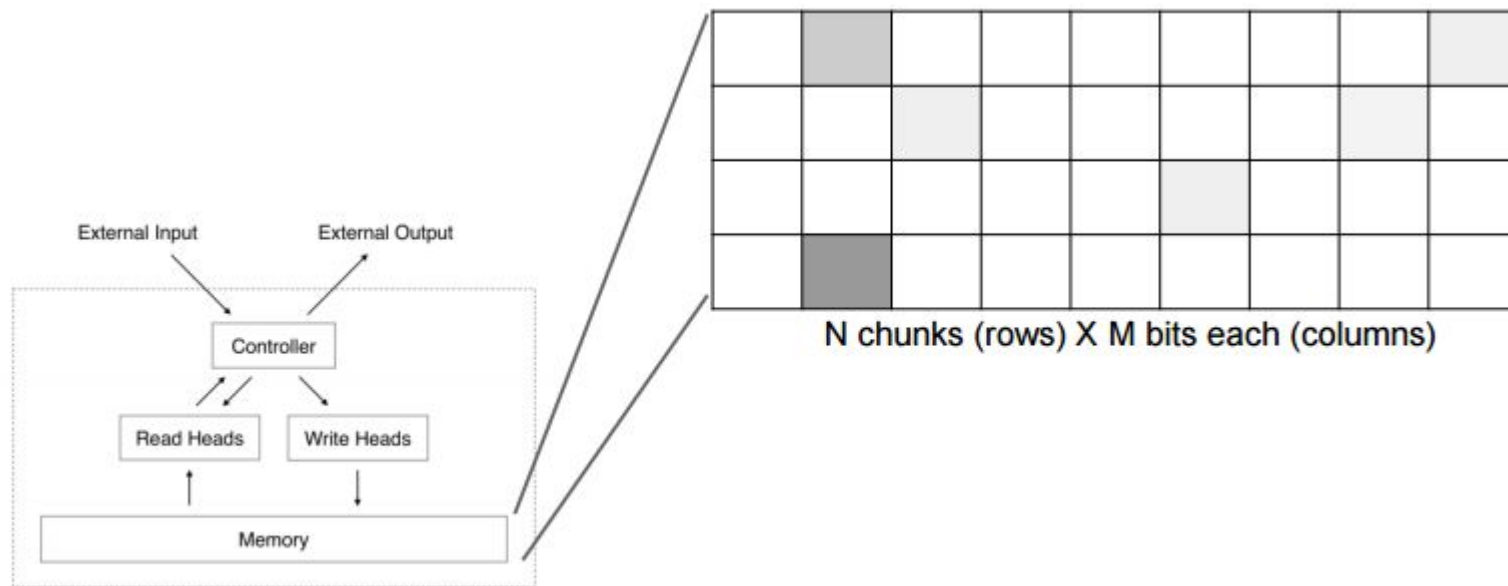


# Neural Turing Machine - Selective Attention

- Focus on the parts of the memory the network will read and write to:  
‘An **introspective attention** model’
- Use the controller outputs to **parameterise a distribution** (weights) over the rows (slots) in the memory matrix
- Weights are defined by two main attention mechanism:
  - One based on **content**
  - One based on **location**



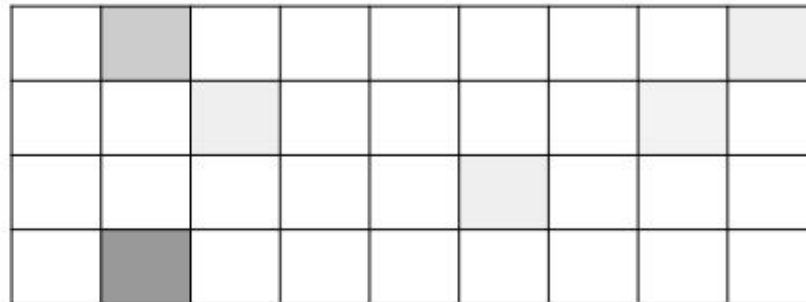
# Neural Turing Machines - Memory



# Neural Turing Machines - Memory

Read from memory (“blurry”)

$$\mathbf{r}_t \longleftarrow \sum_i w_t(i) \mathbf{M}_t(i),$$



N chunks (rows) X M bits each (columns)

Write to memory (“blurry”)

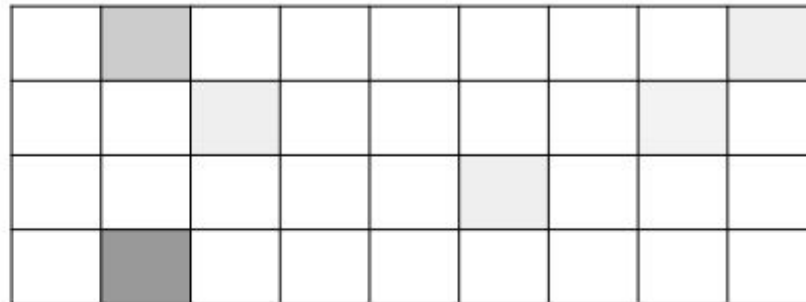
$$\tilde{\mathbf{M}}_t(i) \longleftarrow \mathbf{M}_{t-1}(i) [\mathbf{1} - w_t(i) \mathbf{e}_t].$$

$$\mathbf{M}_t(i) \longleftarrow \tilde{\mathbf{M}}_t(i) + w_t(i) \mathbf{a}_t.$$

# Neural Turing Machines - Memory

Read from memory ("blurry")

$$\mathbf{r}_t \leftarrow \sum_i w_t(i) \mathbf{M}_t(i),$$



N chunks (rows) X M bits each (columns)

Write to memory ("blurry")

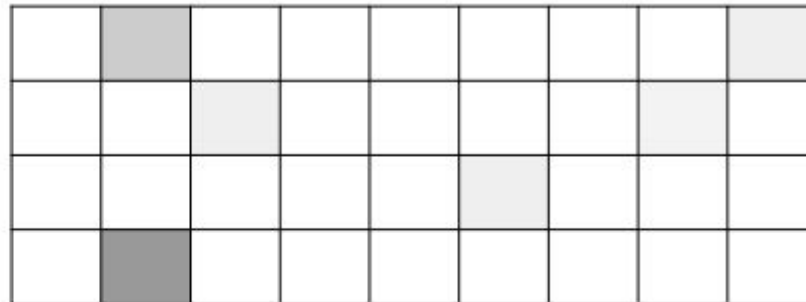
$$\begin{aligned} \tilde{\mathbf{M}}_t(i) &\leftarrow \mathbf{M}_{t-1}(i) [\mathbf{1} - w_t(i) \mathbf{e}_t] \\ \mathbf{M}_t(i) &\leftarrow \tilde{\mathbf{M}}_t(i) + w_t(i) \mathbf{a}_t. \end{aligned}$$

Forget gate !!  
Déjà vu

# Neural Turing Machines - Memory

Read from memory ("blurry")

$$\mathbf{r}_t \leftarrow \sum_i w_t(i) \mathbf{M}_t(i),$$



N chunks (rows) X M bits each (columns)

Write to memory ("blurry")

$$\tilde{\mathbf{M}}_t(i) \leftarrow \mathbf{M}_{t-1}(i) [\mathbf{1} - w_t(i) \mathbf{e}_t]$$

Forget gate !!

Déjà vu

$$\mathbf{M}_t(i) \leftarrow \tilde{\mathbf{M}}_t(i) + w_t(i) \mathbf{a}_t.$$

# LSTM

# Neural Turing Machines - Memory

Addressing by content (similarity)

$$w_t^c(i) \leftarrow \frac{\exp\left(\beta_t K[\mathbf{k}_t, \mathbf{M}_t(i)]\right)}{\sum_j \exp\left(\beta_t K[\mathbf{k}_t, \mathbf{M}_t(j)]\right)}$$

$$K[\mathbf{u}, \mathbf{v}] = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \cdot \|\mathbf{v}\|}$$

Addressing by location (shift)

$$\tilde{w}_t(i) \leftarrow \sum_{j=0}^{N-1} w_t^g(j) s_t(i-j) \quad w_t(i) \leftarrow \frac{\tilde{w}_t(i)^{\gamma_t}}{\sum_j \tilde{w}_t(j)^{\gamma_t}}$$


N chunks (rows) X M bits each (columns)

# Neural Turing Machines - Memory

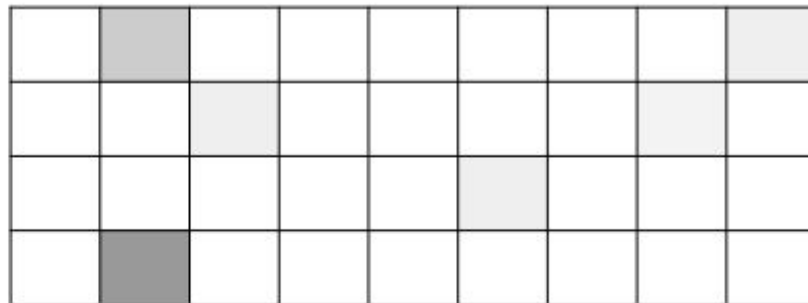
Addressing by content (similarity)

$$w_t^c(i) \leftarrow \frac{\exp\left(\beta_t K[\mathbf{k}_t, \mathbf{M}_t(i)]\right)}{\sum_j \exp\left(\beta_t K[\mathbf{k}_t, \mathbf{M}_t(j)]\right)}$$

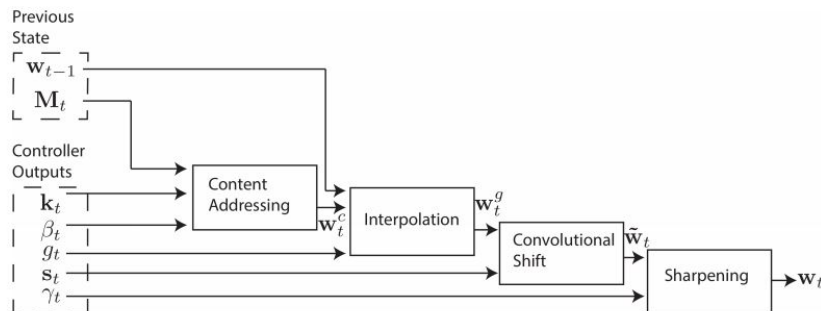
$$K[\mathbf{u}, \mathbf{v}] = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \cdot \|\mathbf{v}\|}$$

Addressing by location (shift)

$$\tilde{w}_t(i) \leftarrow \sum_{j=0}^{N-1} w_t^g(j) s_t(i-j) \quad w_t(i) \leftarrow \frac{\tilde{w}_t(i)^{\gamma_t}}{\sum_j \tilde{w}_t(j)^{\gamma_t}}$$



N chunks (rows) X M bits each (columns)

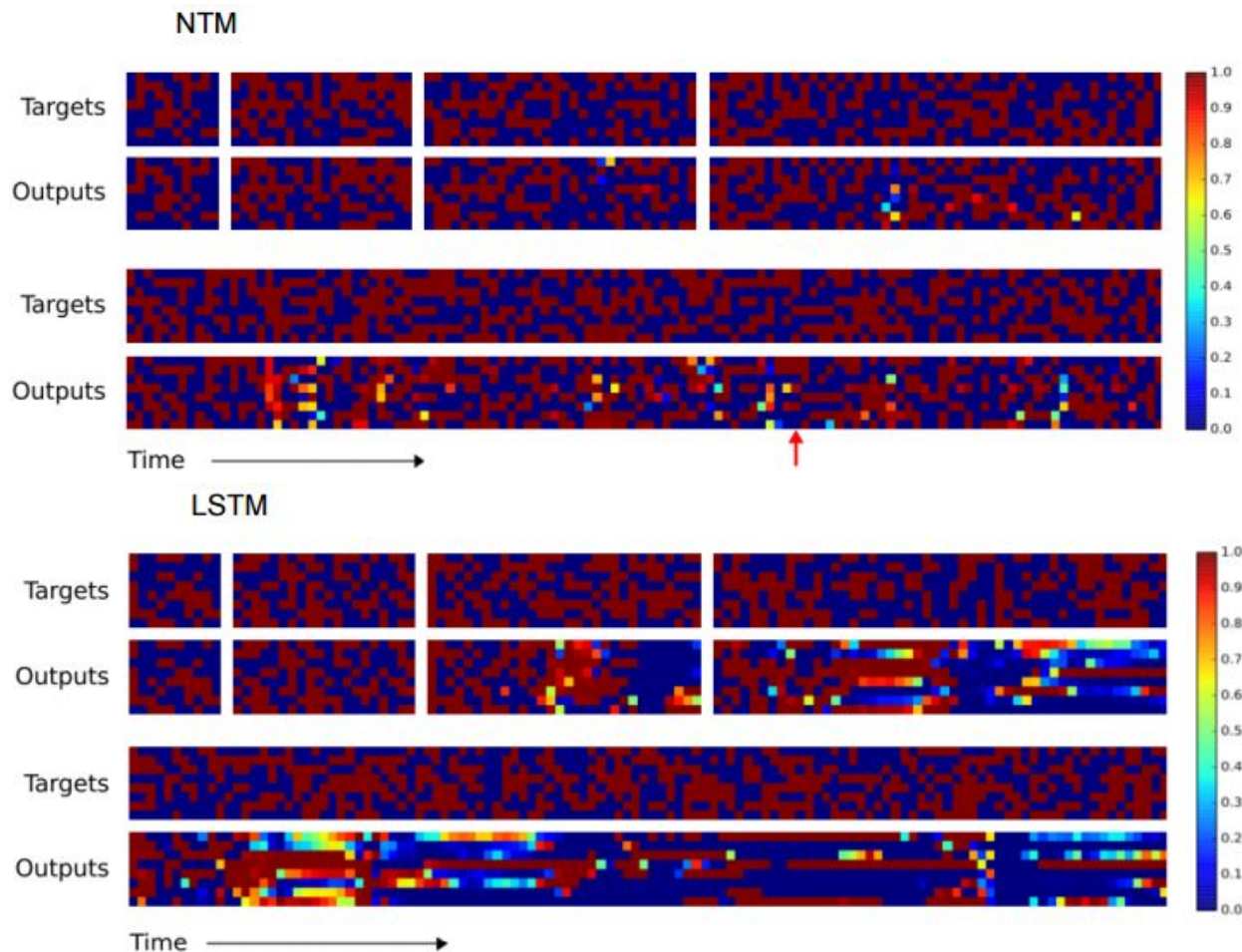


The key vector  $\mathbf{k}_t$  & key strength  $\beta_t$  are used to perform content-based addressing of the memory matrix  $\mathbf{M}_t$ . The resulting content-based weighting is interpolated with the weighting from the previous time step based on the value of the interpolation gate  $g_t$ . The shift weighting  $s_t$  determines whether & by how much the weighting is rotated. Depending on  $\gamma_t$ , the weighting is sharpened and used for memory access.

# NTM - Copy

- NTM can learn basic algorithms like copy, loop, sort, associative recall, N-gram inference.

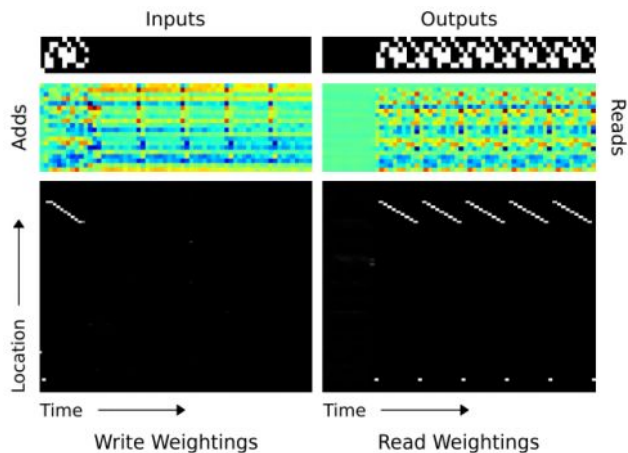
- A simple task of 'copying' the input back can be hard, when the time period becomes very long.





# NTM - Mult Copy

Copying the same sequence many times.



## NTM

Length 10, Repeat 20

Targets



Outputs



Length 20, Repeat 10

Targets



Outputs



## LSTM

Length 10, Repeat 20

Targets



Outputs

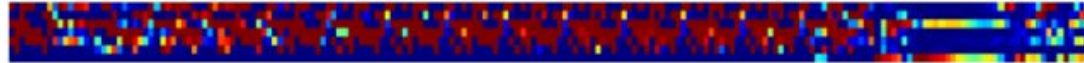


Length 20, Repeat 10

Targets



Outputs

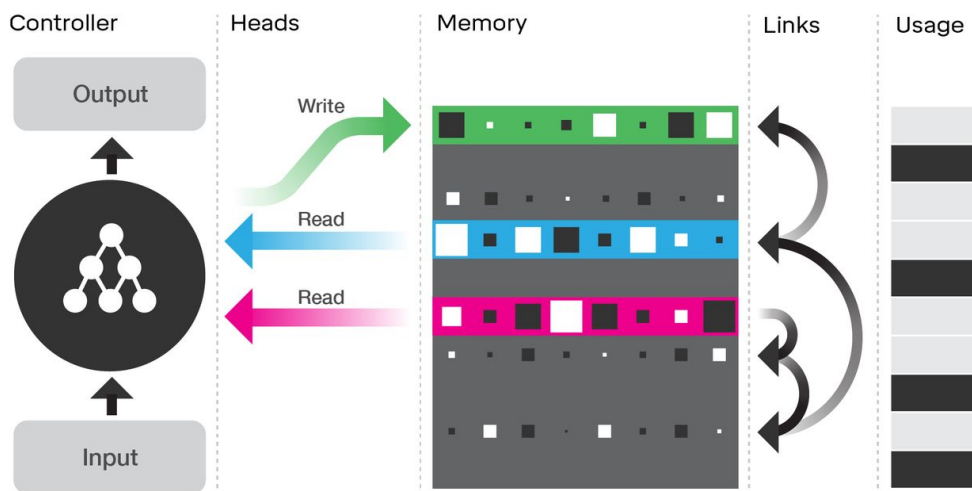


Time →

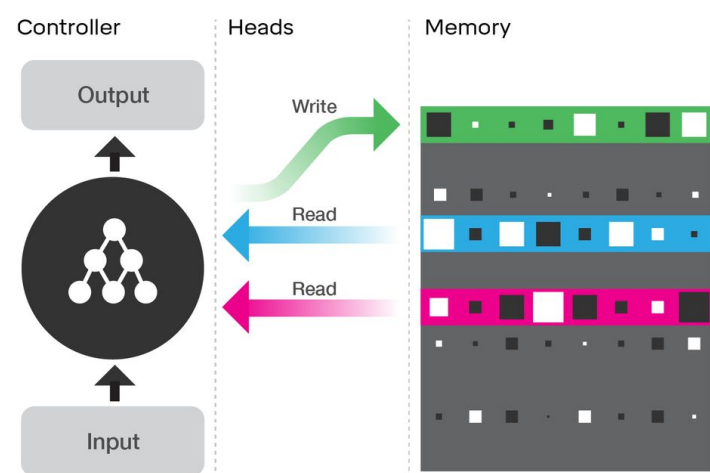


# Differential Neural Computer

Architecture of DNC

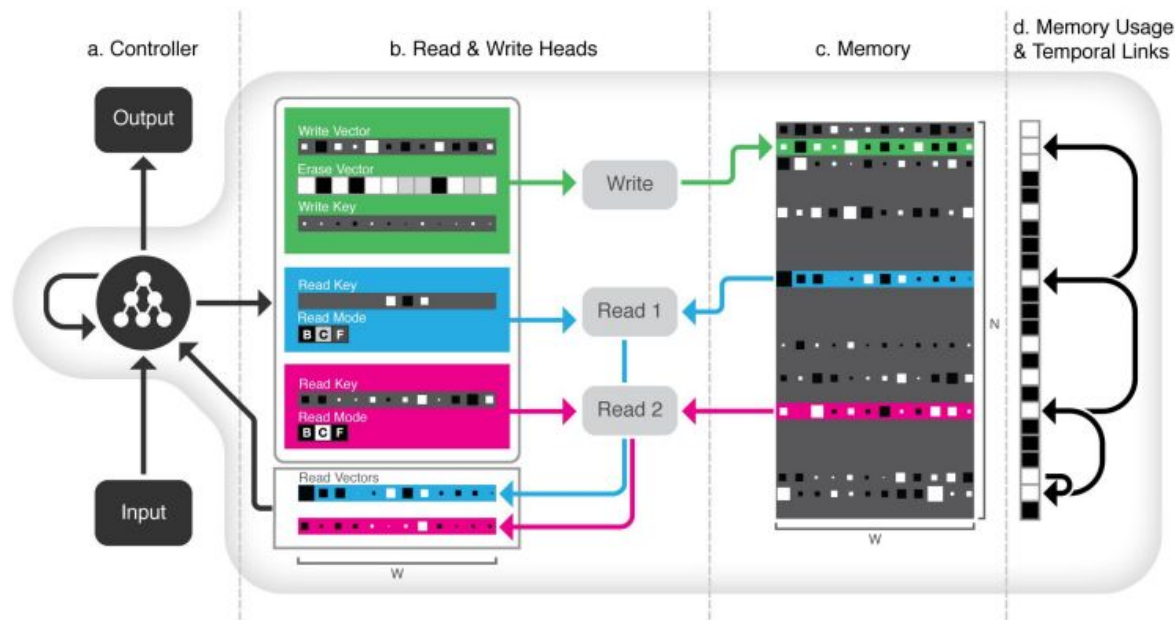


Architecture of NTM



# Differential Neural Computer

Architecture of DNC



- 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
- DNC tries 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:

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

# Differential Neural Computer

$\mathbf{p}_t$  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 ( $\mathbf{w}_{t-1}^{r,i}$ ), allowing it to iterate forwards or backwards in time

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

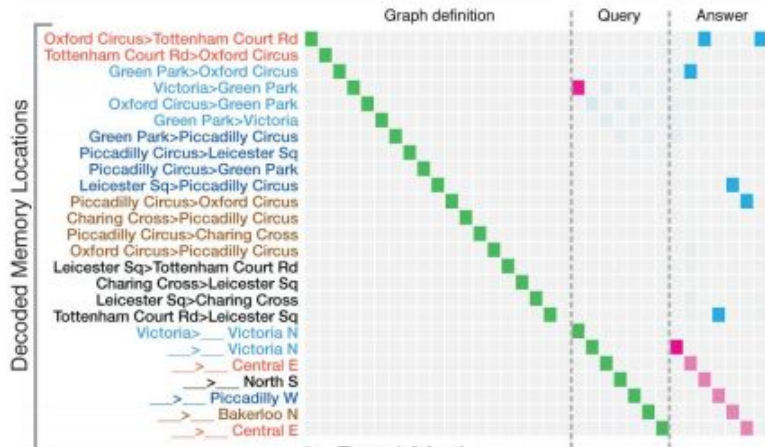
Finally three-way gates ( $\pi_t^i$ ) are used to interpolate among iterating forwards, backwards, or reading by content:

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

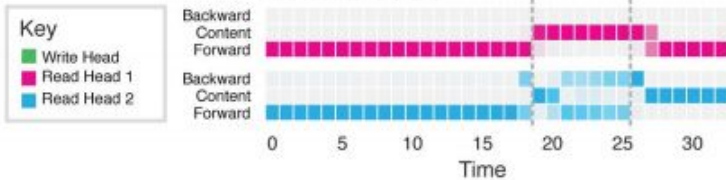
- RNNs are great at processing sequences: text, audio, time-series...
- But graph-structured data is even more general: maps, molecules, parse-trees, knowledge graphs, social networks...
- DNC is able to interpret and answer questions about graphs, even if presented in sequential form
- Started with explicit graphs, but ultimately interested in implicit graphs: relations in natural language, scene parsing, agent's surroundings...

# How DNC Reads Graphs

### a. Read and Write Weightings



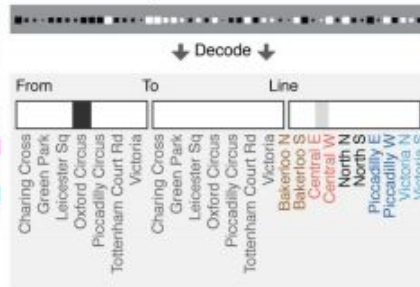
### b. Read Mode



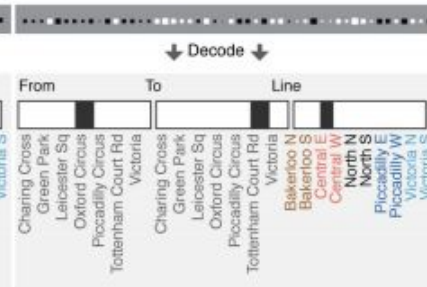
### c. London Underground Map



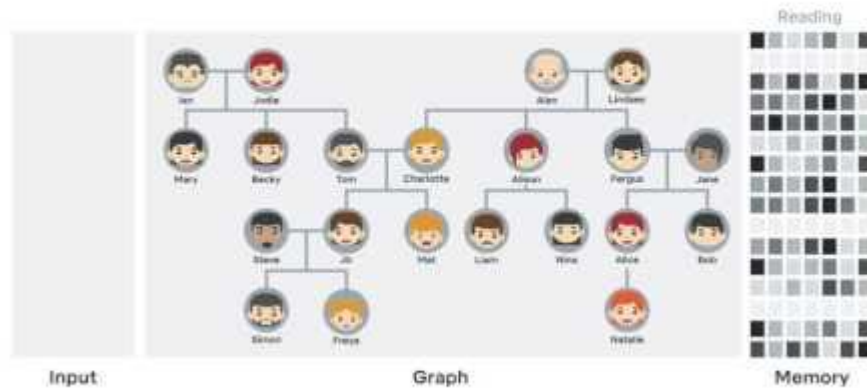
d. Read Key



### e. Location Content



# Differential Neural Computer



# Relevant papers

1. J. Weston, S. Chopra, A. Bordes. [Memory Networks](#). ICLR 2015 (and arXiv:1410.3916).
2. S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus. [End-To-End Memory Networks](#). NIPS 2015
3. J. Weston, et al. [Towards AI-Complete Question Answering: A Set of Prerequisite Toy Tasks](#)
4. A. Bordes, N. Usunier, S. Chopra, J. Weston. [Large-scale Simple Question Answering with Memory Networks](#)
5. Alex Graves, et al [Neural Turing Machine](#). Nature 2015
6. Prakash, et al. [Condensed Memory Networks](#). AAAI 2017
7. Alex Graves, et al [Differential Neural Computers](#). Nature 2017

# Source materials (also good for tutorials)

1. [Jason Weston, Memory Networks](#)
2. [Daniel Shank, NTM](#)
3. [DeepMind, DNC](#)
4. [J. Schmidhuber, How to learn an algorithm](#)
5. [Alex Graves, RNN Symposium on DNC](#)