

Neural networks with dynamic external memory

Differentiable neural computer

Maciej Żelazczyk

December 13, 2017

PhD Student in Computer Science

Division of Artificial Intelligence and Computational Methods

Faculty of Mathematics and Information Science

m.zelazczyk@mini.pw.edu.pl

**Warsaw University
of Technology**

Recurrent neural networks

- Feedforward nets process one input at a time.

Recurrent neural networks

- Feedforward nets process one input at a time.
- Order might be important (e.g. text, sound, video).

Recurrent neural networks

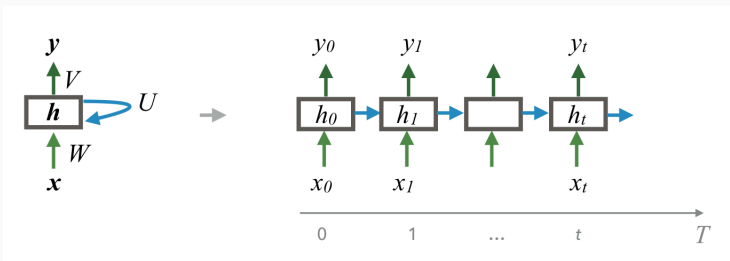
- Feedforward nets process one input at a time.
- Order might be important (e.g. text, sound, video).
- Need to divide data into chunks and process it in sequence.

Recurrent neural networks

- Feedforward nets process one input at a time.
- Order might be important (e.g. text, sound, video).
- Need to divide data into chunks and process it in sequence.
- Adapt feedforward architecture.

Recurrent neural networks

- Feedforward nets process one input at a time.
- Order might be important (e.g. text, sound, video).
- Need to divide data into chunks and process it in sequence.
- Adapt feedforward architecture.

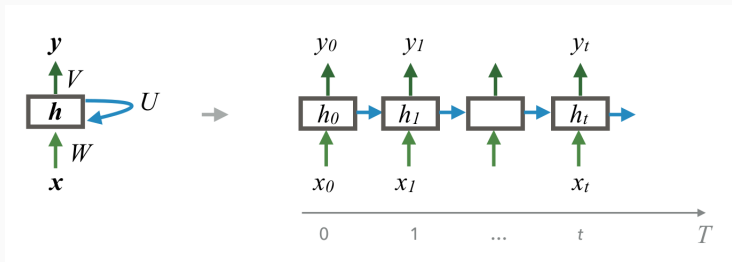


Source: Gakhov, A., *Recurrent Neural Networks. Part 1: Theory*

Vanilla RNNs

- Classic RNN architecture.

$$\mathbf{h}_t = \tanh \left(\begin{bmatrix} \mathbf{U} & \mathbf{W} \end{bmatrix} \begin{bmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{bmatrix} \right)$$

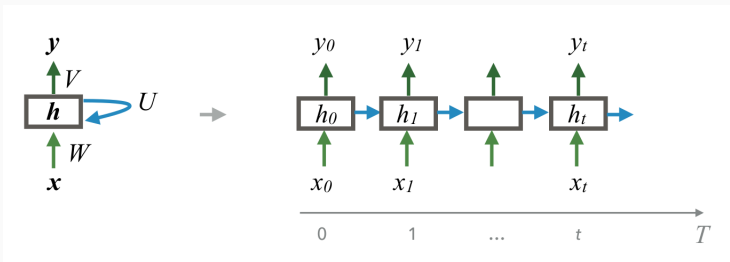


Vanilla RNNs

- Classic RNN architecture.

$$\mathbf{h}_t = \tanh \left(\begin{bmatrix} \mathbf{U} & \mathbf{W} \end{bmatrix} \begin{bmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{bmatrix} \right)$$

- Possible to think of \mathbf{h}_t as of internal memory.

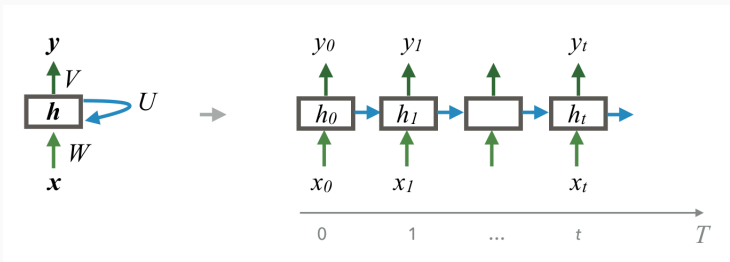


Vanilla RNNs

- Classic RNN architecture.

$$\mathbf{h}_t = \tanh \left(\begin{bmatrix} \mathbf{U} & \mathbf{W} \end{bmatrix} \begin{bmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{bmatrix} \right)$$

- Possible to think of \mathbf{h}_t as of internal memory.
- In practice, this only works for a couple of steps.

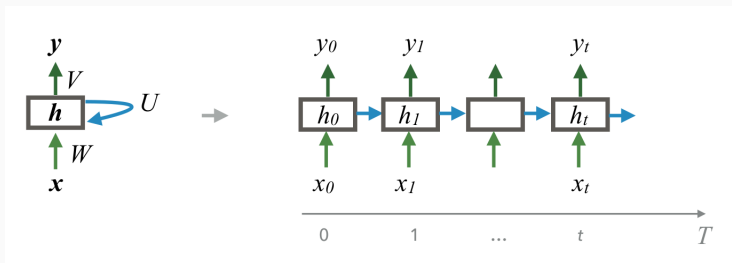


Vanilla RNNs

- Classic RNN architecture.

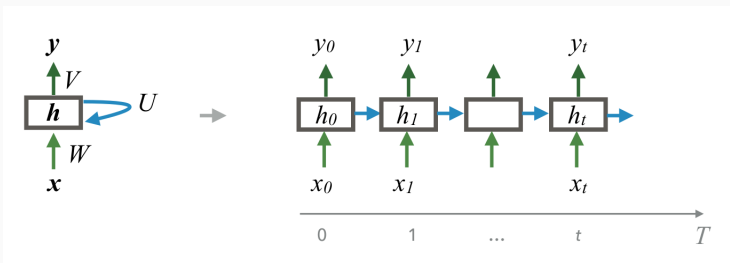
$$\mathbf{h}_t = \tanh \left(\begin{bmatrix} \mathbf{U} & \mathbf{W} \end{bmatrix} \begin{bmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{bmatrix} \right)$$

- Possible to think of \mathbf{h}_t as of internal memory.
- In practice, this only works for a couple of steps.
- Gradient either vanishes or explodes during training.



Vanishing gradient problem

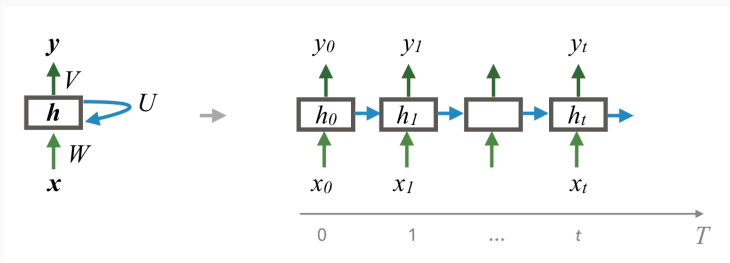
- First explanation of unstable gradients in [Hochreiter, 1991].



Source: Gakhov, A., *Recurrent Neural Networks. Part 1: Theory*

Vanishing gradient problem

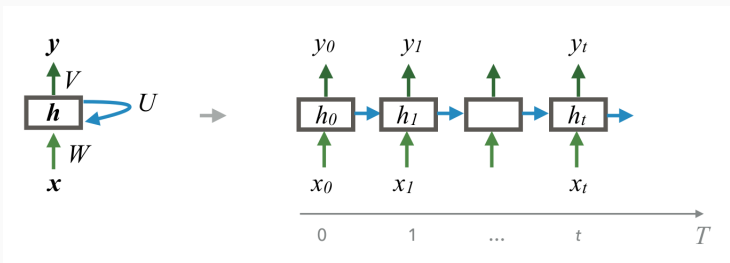
- First explanation of unstable gradients in [Hochreiter, 1991].
- General idea: multiplying by $\frac{d}{dx} \tanh x = 1 - \tanh^2 x \in (0, 1]$.



Source: Gakhov, A., *Recurrent Neural Networks. Part 1: Theory*

Vanishing gradient problem

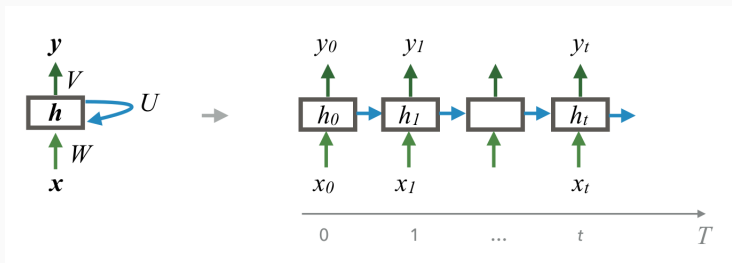
- First explanation of unstable gradients in [Hochreiter, 1991].
- General idea: multiplying by $\frac{d}{dx} \tanh x = 1 - \tanh^2 x \in (0, 1]$.
- Formal argument: based on eigenvalues.



Source: Gakhov, A., *Recurrent Neural Networks. Part 1: Theory*

Vanishing gradient problem

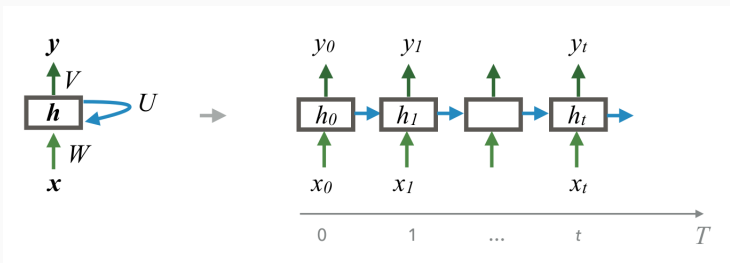
- First explanation of unstable gradients in [Hochreiter, 1991].
- General idea: multiplying by $\frac{d}{dx} \tanh x = 1 - \tanh^2 x \in (0, 1]$.
- Formal argument: based on eigenvalues.
- Vanilla RNNs are inherently unstable in training.



Source: Gakhov, A., *Recurrent Neural Networks. Part 1: Theory*

Vanishing gradient problem

- First explanation of unstable gradients in [Hochreiter, 1991].
- General idea: multiplying by $\frac{d}{dx} \tanh x = 1 - \tanh^2 x \in (0, 1]$.
- Formal argument: based on eigenvalues.
- Vanilla RNNs are inherently unstable in training.
- Memory is limited to < 10 steps.



Source: Gakhov, A., *Recurrent Neural Networks. Part 1: Theory*

Long-short term memory

- Specifically design an architecture to circumvent vanishing gradients [Hochreiter and Schmidhuber, 1997].

Long-short term memory

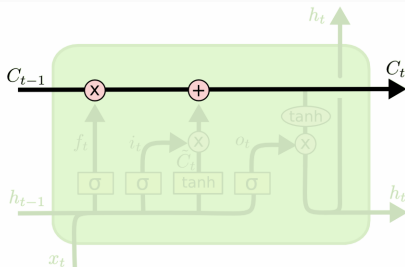
- Specifically design an architecture to circumvent vanishing gradients [Hochreiter and Schmidhuber, 1997].
- General idea: additive interactions transport gradients better.

Long-short term memory

- Specifically design an architecture to circumvent vanishing gradients [Hochreiter and Schmidhuber, 1997].
- General idea: additive interactions transport gradients better.
- Add a separate state cell \mathbf{c}_t .

Long-short term memory

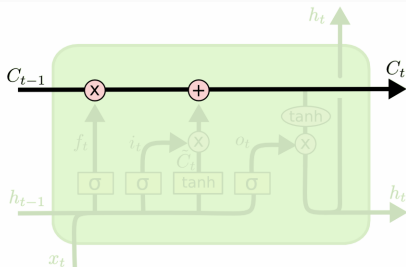
- Specifically design an architecture to circumvent vanishing gradients [Hochreiter and Schmidhuber, 1997].
- General idea: additive interactions transport gradients better.
- Add a separate state cell c_t .



Source: Olah, C., *Understanding LSTM Networks*

Long-short term memory

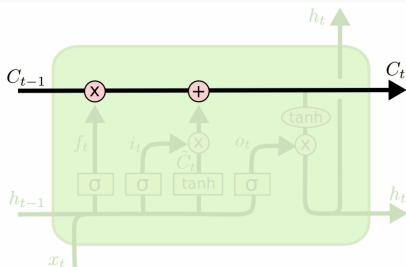
- In practice, this works relatively well (text classification, translation etc.).



Source: Olah, C., *Understanding LSTM Networks*

Long-short term memory

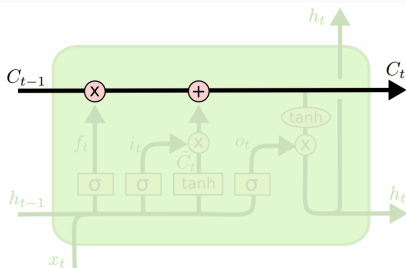
- In practice, this works relatively well (text classification, translation etc.).
- Memory persists for ≈ 100 steps.



Source: Olah, C., *Understanding LSTM Networks*

Long-short term memory

- In practice, this works relatively well (text classification, translation etc.).
- Memory persists for ≈ 100 steps.
- State cell was not designed as memory in traditional sense.



Source: Olah, C., *Understanding LSTM Networks*

Limits of LSTMs

Despite their undeniable success, LSTMs suffer from a number of limitations:

Source: Graves, A., *IJCNN 2017 Plenary Talk: Frontiers in Recurrent Neural Network Research*

Limits of LSTMs

Despite their undeniable success, LSTMs suffer from a number of limitations:

1. 100 steps is not how human memory works.

Source: Graves, A., *IJCNN 2017 Plenary Talk: Frontiers in Recurrent Neural Network Research*

Limits of LSTMs

Despite their undeniable success, LSTMs suffer from a number of limitations:

1. 100 steps is not how human memory works.
2. In practice, hidden state \mathbf{h}_t is modified at each time step.

Source: Graves, A., *IJCNN 2017 Plenary Talk: Frontiers in Recurrent Neural Network Research*

Limits of LSTMs

Despite their undeniable success, LSTMs suffer from a number of limitations:

1. 100 steps is not how human memory works.
2. In practice, hidden state \mathbf{h}_t is modified at each time step.
3. Increasing the size of memory is equivalent to expanding the vector \mathbf{h}_t and the whole network. No. of weights grows at least linearly with required memory.

Source: Graves, A., *IJCNN 2017 Plenary Talk: Frontiers in Recurrent Neural Network Research*

Limits of LSTMs

Despite their undeniable success, LSTMs suffer from a number of limitations:

1. 100 steps is not how human memory works.
2. In practice, hidden state \mathbf{h}_t is modified at each time step.
3. Increasing the size of memory is equivalent to expanding the vector \mathbf{h}_t and the whole network. No. of weights grows at least linearly with required memory.
4. Memory might become "hard-coded." Specific parts of the network might be used to detect given features. Location and content are intertwined.

Source: Graves, A., *IJCNN 2017 Plenary Talk: Frontiers in Recurrent Neural Network Research*

Memory/location entanglement

Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.

Cell that turns on inside quotes:

"You mean to imply that I have nothing to eat out of... on the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

Source: Karpathy, A., *The Unreasonable Effectiveness of Recurrent Neural Networks*

Networks with external memory

Adding an external memory source mitigates some of the mentioned problems:

Networks with external memory

Adding an external memory source mitigates some of the mentioned problems:

1. Not all of the memory is interacted with all the time. Specific parts of memory are accessed at each time step. Memory is "protected."

Networks with external memory

Adding an external memory source mitigates some of the mentioned problems:

1. Not all of the memory is interacted with all the time. Specific parts of memory are accessed at each time step. Memory is "protected."
2. Computational cost is not necessarily scaling up with the size of the memory. In theory, memory can be very large. Analogy: increase amount of RAM without changing the CPU.

Networks with external memory

Adding an external memory source mitigates some of the mentioned problems:

1. Not all of the memory is interacted with all the time. Specific parts of memory are accessed at each time step. Memory is "protected."
2. Computational cost is not necessarily scaling up with the size of the memory. In theory, memory can be very large. Analogy: increase amount of RAM without changing the CPU.
3. Content is separated out from location. Computation separated from memory.

Networks with external memory

Adding an external memory source mitigates some of the mentioned problems:

1. Not all of the memory is interacted with all the time. Specific parts of memory are accessed at each time step. Memory is "protected."
2. Computational cost is not necessarily scaling up with the size of the memory. In theory, memory can be very large. Analogy: increase amount of RAM without changing the CPU.
3. Content is separated out from location. Computation separated from memory.
4. Easier to deal with variables, linked lists, etc. Abstraction comes in handy.

Differentiable neural computer

Cast abstract idea in concrete architecture: Differentiable neural computer (DNC) [Graves et al., 2016]. Design principles:

Differentiable neural computer

Cast abstract idea in concrete architecture: Differentiable neural computer (DNC) [Graves et al., 2016]. Design principles:

- Differentiable end-to-end.

Differentiable neural computer

Cast abstract idea in concrete architecture: Differentiable neural computer (DNC) [Graves et al., 2016]. Design principles:

- Differentiable end-to-end.
- Read-write memory.

Differentiable neural computer

Cast abstract idea in concrete architecture: Differentiable neural computer (DNC) [Graves et al., 2016]. Design principles:

- Differentiable end-to-end.
- Read-write memory.

Overview of DNC:

Differentiable neural computer

Cast abstract idea in concrete architecture: Differentiable neural computer (DNC) [Graves et al., 2016]. Design principles:

- Differentiable end-to-end.
- Read-write memory.

Overview of DNC:

(a) Controller - neural network, e.g. deep LSTM.

Differentiable neural computer

Cast abstract idea in concrete architecture: Differentiable neural computer (DNC) [Graves et al., 2016]. Design principles:

- Differentiable end-to-end.
- Read-write memory.

Overview of DNC:

- (a) Controller - neural network, e.g. deep LSTM.
- (b) Read and write heads.

Differentiable neural computer

Cast abstract idea in concrete architecture: Differentiable neural computer (DNC) [Graves et al., 2016]. Design principles:

- Differentiable end-to-end.
- Read-write memory.

Overview of DNC:

- (a) Controller - neural network, e.g. deep LSTM.
- (b) Read and write heads.
- (c) Memory matrix.

Differentiable neural computer

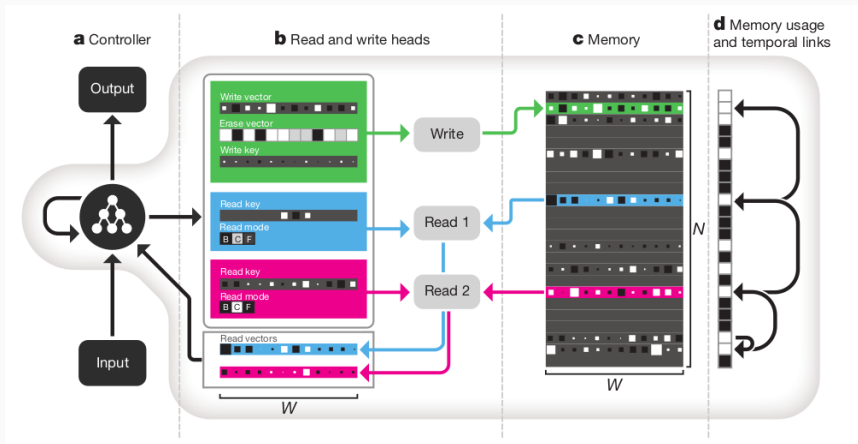
Cast abstract idea in concrete architecture: Differentiable neural computer (DNC) [Graves et al., 2016]. Design principles:

- Differentiable end-to-end.
- Read-write memory.

Overview of DNC:

- (a) Controller - neural network, e.g. deep LSTM.
- (b) Read and write heads.
- (c) Memory matrix.
- (d) Memory usage vector and temporal link matrix.

Differentiable neural computer



Source: [Graves et al., 2016]

Controller

Neural network \mathcal{N} . Let's use a deep LSTM architecture, which carries a hidden state vector $[\mathbf{h}_t^1 \ \dots \ \mathbf{h}_t^L]$.

Input:

Neural network \mathcal{N} . Let's use a deep LSTM architecture, which carries a hidden state vector $\left[\mathbf{h}_t^1 \quad \dots \quad \mathbf{h}_t^L \right]$.

Input:

- External input $\mathbf{x}_t \in \mathbb{R}^X$.

Controller

Neural network \mathcal{N} . Let's use a deep LSTM architecture, which carries a hidden state vector $[\mathbf{h}_t^1 \ \dots \ \mathbf{h}_t^L]$.

Input:

- External input $\mathbf{x}_t \in \mathbb{R}^X$.
- R read vectors $\mathbf{r}_{t-1}^1, \dots, \mathbf{r}_{t-1}^R \in M_{t-1} \in \mathbb{R}^{N \times W}$.

Controller

Neural network \mathcal{N} . Let's use a deep LSTM architecture, which carries a hidden state vector $[\mathbf{h}_t^1 \ \dots \ \mathbf{h}_t^L]$.

Input:

- External input $\mathbf{x}_t \in \mathbb{R}^X$.
- R read vectors $\mathbf{r}_{t-1}^1, \dots, \mathbf{r}_{t-1}^R \in M_{t-1} \in \mathbb{R}^{N \times W}$.

Output:

Controller

Neural network \mathcal{N} . Let's use a deep LSTM architecture, which carries a hidden state vector $\begin{bmatrix} \mathbf{h}_t^1 & \dots & \mathbf{h}_t^L \end{bmatrix}$.

Input:

- External input $\mathbf{x}_t \in \mathbb{R}^X$.
- R read vectors $\mathbf{r}_{t-1}^1, \dots, \mathbf{r}_{t-1}^R \in M_{t-1} \in \mathbb{R}^{N \times W}$.

Output:

- Controller output vector $\mathbf{v}_t = W_y \begin{bmatrix} \mathbf{h}_t^1 & \dots & \mathbf{h}_t^L \end{bmatrix} \in \mathbb{R}^Y$.

Controller

Neural network \mathcal{N} . Let's use a deep LSTM architecture, which carries a hidden state vector $\begin{bmatrix} \mathbf{h}_t^1 & \dots & \mathbf{h}_t^L \end{bmatrix}$.

Input:

- External input $\mathbf{x}_t \in \mathbb{R}^X$.
- R read vectors $\mathbf{r}_{t-1}^1, \dots, \mathbf{r}_{t-1}^R \in M_{t-1} \in \mathbb{R}^{N \times W}$.

Output:

- Controller output vector $\mathbf{v}_t = W_y \begin{bmatrix} \mathbf{h}_t^1 & \dots & \mathbf{h}_t^L \end{bmatrix} \in \mathbb{R}^Y$.
- Interface vector $\hat{\xi}_t = W_\xi \begin{bmatrix} \mathbf{h}_t^1 & \dots & \mathbf{h}_t^L \end{bmatrix} \in \mathbb{R}^{(W \times R) + 3W + 5R + 3}$.

Controller

Neural network \mathcal{N} . Let's use a deep LSTM architecture, which carries a hidden state vector $\begin{bmatrix} \mathbf{h}_t^1 & \dots & \mathbf{h}_t^L \end{bmatrix}$.

Input:

- External input $\mathbf{x}_t \in \mathbb{R}^X$.
- R read vectors $\mathbf{r}_{t-1}^1, \dots, \mathbf{r}_{t-1}^R \in M_{t-1} \in \mathbb{R}^{N \times W}$.

Output:

- Controller output vector $\mathbf{v}_t = W_y \begin{bmatrix} \mathbf{h}_t^1 & \dots & \mathbf{h}_t^L \end{bmatrix} \in \mathbb{R}^Y$.
- Interface vector $\hat{\xi}_t = W_\xi \begin{bmatrix} \mathbf{h}_t^1 & \dots & \mathbf{h}_t^L \end{bmatrix} \in \mathbb{R}^{(W \times R) + 3W + 5R + 3}$.
- Memory-augmented output vector $\mathbf{y}_t = \mathbf{v}_t + W_r \begin{bmatrix} \mathbf{r}_t^1 & \dots & \mathbf{r}_t^R \end{bmatrix} \in \mathbb{R}^Y$.

- Interface vector before processing $\hat{\xi}_t =$
 $\left[\mathbf{k}_t^{r,1} \dots \mathbf{k}_t^{r,R}; \hat{\beta}_t^{r,1} \dots \hat{\beta}_t^{r,R}; \mathbf{k}_t^w; \hat{\beta}_t^w; \hat{\mathbf{e}}_t; \mathbf{v}_t; \hat{f}_t^1 \dots \hat{f}_t^R; \hat{\mathbf{g}}_t^a; \hat{\mathbf{g}}_t^w; \hat{\pi}_t^1 \dots \hat{\pi}_t^R \right]$

Interface vector

- Interface vector before processing $\hat{\xi}_t =$
 $\left[\mathbf{k}_t^{r,1} \dots \mathbf{k}_t^{r,R}; \hat{\beta}_t^{r,1} \dots \hat{\beta}_t^{r,R}; \mathbf{k}_t^w; \hat{\beta}_t^w; \hat{\mathbf{e}}_t; \mathbf{v}_t; \hat{f}_t^1 \dots \hat{f}_t^R; \hat{\mathbf{g}}_t^a; \hat{\mathbf{g}}_t^w; \hat{\pi}_t^1 \dots \hat{\pi}_t^R \right]$
- Define: $\text{oneplus}(x) = 1 + \ln(1 + e^x) \in [1, \infty)$.

Interface vector

- Interface vector before processing $\hat{\xi}_t =$
 $\left[\mathbf{k}_t^{r,1} \dots \mathbf{k}_t^{r,R}; \hat{\beta}_t^{r,1} \dots \hat{\beta}_t^{r,R}; \mathbf{k}_t^w; \hat{\beta}_t^w; \hat{\mathbf{e}}_t; \mathbf{v}_t; \hat{f}_t^1 \dots \hat{f}_t^R; \hat{g}_t^a; \hat{g}_t^w; \hat{\pi}_t^1 \dots \hat{\pi}_t^R \right]$
- Define: $\text{oneplus}(x) = 1 + \ln(1 + e^x) \in [1, \infty)$.
- Define: $\text{softmax}(\mathbf{x})_j = \frac{e^{x_j}}{\sum_i e^{x_i}}$.

Interface vector

- Interface vector before processing $\hat{\xi}_t =$
$$\left[\mathbf{k}_t^{r,1} \dots \mathbf{k}_t^{r,R}; \hat{\beta}_t^{r,1} \dots \hat{\beta}_t^{r,R}; \mathbf{k}_t^w; \hat{\beta}_t^w; \hat{\mathbf{e}}_t; \mathbf{v}_t; \hat{f}_t^1 \dots \hat{f}_t^R; \hat{g}_t^a; \hat{g}_t^w; \hat{\pi}_t^1 \dots \hat{\pi}_t^R \right]$$
- Define: $\text{oneplus}(x) = 1 + \ln(1 + e^x) \in [1, \infty)$.
- Define: $\text{softmax}(\mathbf{x})_j = \frac{e^{x_j}}{\sum_i e^{x_i}}$.
- $\beta_t^{r,i} = \text{oneplus}(\hat{\beta}_t^{r,i}), \beta_t^w = \text{oneplus}(\hat{\beta}_t^w)$

Interface vector

- Interface vector before processing $\hat{\xi}_t =$
 $\left[\mathbf{k}_t^{r,1} \dots \mathbf{k}_t^{r,R}; \hat{\beta}_t^{r,1} \dots \hat{\beta}_t^{r,R}; \mathbf{k}_t^w; \hat{\beta}_t^w; \hat{\mathbf{e}}_t; \mathbf{v}_t; \hat{f}_t^1 \dots \hat{f}_t^R; \hat{g}_t^a; \hat{g}_t^w; \hat{\pi}_t^1 \dots \hat{\pi}_t^R \right]$
- Define: $\text{oneplus}(x) = 1 + \ln(1 + e^x) \in [1, \infty)$.
- Define: $\text{softmax}(\mathbf{x})_j = \frac{e^{x_j}}{\sum_i e^{x_i}}$.
- $\beta_t^{r,i} = \text{oneplus}(\hat{\beta}_t^{r,i}), \beta_t^w = \text{oneplus}(\hat{\beta}_t^w)$
- $\mathbf{e}_t = \sigma(\hat{\mathbf{e}}_t), f_t^i = \sigma(\hat{f}_t^i), g_t^a = \sigma(\hat{g}_t^a), g_t^w = \sigma(\hat{g}_t^w)$

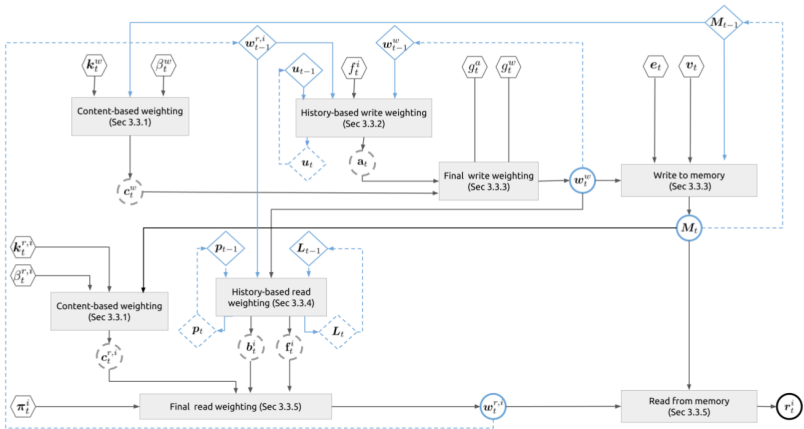
Interface vector

- Interface vector before processing $\hat{\xi}_t =$
 $\left[\mathbf{k}_t^{r,1} \dots \mathbf{k}_t^{r,R}; \hat{\beta}_t^{r,1} \dots \hat{\beta}_t^{r,R}; \mathbf{k}_t^w; \hat{\beta}_t^w; \hat{\mathbf{e}}_t; \mathbf{v}_t; \hat{f}_t^1 \dots \hat{f}_t^R; \hat{g}_t^a; \hat{g}_t^w; \hat{\pi}_t^1 \dots \hat{\pi}_t^R \right]$
- Define: $\text{oneplus}(x) = 1 + \ln(1 + e^x) \in [1, \infty)$.
- Define: $\text{softmax}(\mathbf{x})_j = \frac{e^{x_j}}{\sum_i e^{x_i}}$.
- $\beta_t^{r,i} = \text{oneplus}(\hat{\beta}_t^{r,i}), \beta_t^w = \text{oneplus}(\hat{\beta}_t^w)$
- $\mathbf{e}_t = \sigma(\hat{\mathbf{e}}_t), f_t^i = \sigma(\hat{f}_t^i), g_t^a = \sigma(\hat{g}_t^a), g_t^w = \sigma(\hat{g}_t^w)$
- $\pi_t^i = \text{softmax}(\hat{\pi}_t^i)$

Interface vector

- Interface vector before processing $\hat{\xi}_t =$
$$\left[\mathbf{k}_t^{r,1} \dots \mathbf{k}_t^{r,R}; \hat{\beta}_t^{r,1} \dots \hat{\beta}_t^{r,R}; \mathbf{k}_t^w; \hat{\beta}_t^w; \hat{\mathbf{e}}_t; \mathbf{v}_t; \hat{f}_t^1 \dots \hat{f}_t^R; \hat{g}_t^a; \hat{g}_t^w; \hat{\pi}_t^1 \dots \hat{\pi}_t^R \right]$$
- Define: $\text{oneplus}(x) = 1 + \ln(1 + e^x) \in [1, \infty)$.
- Define: $\text{softmax}(\mathbf{x})_j = \frac{e^{x_j}}{\sum_i e^{x_i}}$.
- $\beta_t^{r,i} = \text{oneplus}(\hat{\beta}_t^{r,i}), \beta_t^w = \text{oneplus}(\hat{\beta}_t^w)$
- $\mathbf{e}_t = \sigma(\hat{\mathbf{e}}_t), f_t^i = \sigma(\hat{f}_t^i), g_t^a = \sigma(\hat{g}_t^a), g_t^w = \sigma(\hat{g}_t^w)$
- $\pi_t^i = \text{softmax}(\hat{\pi}_t^i)$
- Interface vector after processing $\xi_t =$
$$\left[\mathbf{k}_t^{r,1} \dots \mathbf{k}_t^{r,R}; \beta_t^{r,1} \dots \beta_t^{r,R}; \mathbf{k}_t^w; \beta_t^w; \mathbf{e}_t; \mathbf{v}_t; f_t^1 \dots f_t^R; g_t^a; g_t^w; \pi_t^1 \dots \pi_t^R \right]$$

Interacting with memory



Source: Hsin, C., *Implementation and Optimization of Differentiable Neural Computers*

1. Content-based addressing:

1. Content-based addressing:

- $$C(M, \mathbf{k}, \beta)[i] = \frac{\exp\{\mathcal{D}(\mathbf{k}, M[i, \cdot])\beta\}}{\sum_j \exp\{\mathcal{D}(\mathbf{k}, M[j, \cdot])\beta\}}$$

1. Content-based addressing:

- $\mathcal{C}(M, \mathbf{k}, \beta)[i] = \frac{\exp\{\mathcal{D}(\mathbf{k}, M[i, \cdot])\beta\}}{\sum_j \exp\{\mathcal{D}(\mathbf{k}, M[j, \cdot])\beta\}}$
- cosine similarity $\mathcal{D}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \in [-1, 1]$

1. Content-based addressing:

- $\mathcal{C}(M, \mathbf{k}, \beta)[i] = \frac{\exp\{\mathcal{D}(\mathbf{k}, M[i, \cdot])\beta\}}{\sum_j \exp\{\mathcal{D}(\mathbf{k}, M[j, \cdot])\beta\}}$
- cosine similarity $\mathcal{D}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{|\mathbf{u}| |\mathbf{v}|} \in [-1, 1]$
- $\mathbf{c}_t^w = \mathcal{C}(M_{t-1}, \mathbf{k}_t^w, \beta_t^w) \in \mathcal{S}_N$

1. Content-based addressing:

- $\mathcal{C}(M, \mathbf{k}, \beta)[i] = \frac{\exp\{\mathcal{D}(\mathbf{k}, M[i, \cdot])\beta\}}{\sum_j \exp\{\mathcal{D}(\mathbf{k}, M[j, \cdot])\beta\}}$
- cosine similarity $\mathcal{D}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \in [-1, 1]$
- $\mathbf{c}_t^w = \mathcal{C}(M_{t-1}, \mathbf{k}_t^w, \beta_t^w) \in \mathcal{S}_N$

2. Dynamic memory allocation:

Writing to memory

1. Content-based addressing:

- $\mathcal{C}(M, \mathbf{k}, \beta)[i] = \frac{\exp\{\mathcal{D}(\mathbf{k}, M[i, \cdot])\beta\}}{\sum_j \exp\{\mathcal{D}(\mathbf{k}, M[j, \cdot])\beta\}}$
- cosine similarity $\mathcal{D}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \in [-1, 1]$
- $\mathbf{c}_t^w = \mathcal{C}(M_{t-1}, \mathbf{k}_t^w, \beta_t^w) \in \mathcal{S}_N$

2. Dynamic memory allocation:

- memory retention vector $\psi_t = \prod_{i=1}^R (\mathbf{1} - f_t^i \mathbf{w}_{t-1}^{r,i}) \in [0, 1]^N$

Writing to memory

1. Content-based addressing:

- $\mathcal{C}(M, \mathbf{k}, \beta)[i] = \frac{\exp\{\mathcal{D}(\mathbf{k}, M[i, \cdot])\beta\}}{\sum_j \exp\{\mathcal{D}(\mathbf{k}, M[j, \cdot])\beta\}}$
- cosine similarity $\mathcal{D}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \in [-1, 1]$
- $\mathbf{c}_t^w = \mathcal{C}(M_{t-1}, \mathbf{k}_t^w, \beta_t^w) \in \mathcal{S}_N$

2. Dynamic memory allocation:

- memory retention vector $\psi_t = \prod_{i=1}^R (\mathbf{1} - f_t^i \mathbf{w}_{t-1}^{r,i}) \in [0, 1]^N$
- memory usage vector
 $\mathbf{u}_t = (\mathbf{u}_{t-1} + \mathbf{w}_{t-1}^w - \mathbf{u}_{t-1} \circ \mathbf{w}_{t-1}^w) \circ \psi_t \in [0, 1]^N$

Writing to memory

1. Content-based addressing:

- $\mathcal{C}(M, \mathbf{k}, \beta)[i] = \frac{\exp\{\mathcal{D}(\mathbf{k}, M[i, \cdot])\beta\}}{\sum_j \exp\{\mathcal{D}(\mathbf{k}, M[j, \cdot])\beta\}}$
- cosine similarity $\mathcal{D}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \in [-1, 1]$
- $\mathbf{c}_t^w = \mathcal{C}(M_{t-1}, \mathbf{k}_t^w, \beta_t^w) \in \mathcal{S}_N$

2. Dynamic memory allocation:

- memory retention vector $\psi_t = \prod_{i=1}^R (\mathbf{1} - f_t^i \mathbf{w}_{t-1}^{r,i}) \in [0, 1]^N$
- memory usage vector
 $\mathbf{u}_t = (\mathbf{u}_{t-1} + \mathbf{w}_{t-1}^w - \mathbf{u}_{t-1} \circ \mathbf{w}_{t-1}^w) \circ \psi_t \in [0, 1]^N$
- sort indices of memory locations in ascending order of usage,
 $\phi_t \in \mathbb{N}^+, \phi_t[1]$ is the least used location

Writing to memory

1. Content-based addressing:

- $\mathcal{C}(M, \mathbf{k}, \beta)[i] = \frac{\exp\{\mathcal{D}(\mathbf{k}, M[i, \cdot])\beta\}}{\sum_j \exp\{\mathcal{D}(\mathbf{k}, M[j, \cdot])\beta\}}$
- cosine similarity $\mathcal{D}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \in [-1, 1]$
- $\mathbf{c}_t^w = \mathcal{C}(M_{t-1}, \mathbf{k}_t^w, \beta_t^w) \in \mathcal{S}_N$

2. Dynamic memory allocation:

- memory retention vector $\psi_t = \prod_{i=1}^R (\mathbf{1} - f_t^i \mathbf{w}_{t-1}^{r,i}) \in [0, 1]^N$
- memory usage vector
 $\mathbf{u}_t = (\mathbf{u}_{t-1} + \mathbf{w}_{t-1}^w - \mathbf{u}_{t-1} \circ \mathbf{w}_{t-1}^w) \circ \psi_t \in [0, 1]^N$
- sort indices of memory locations in ascending order of usage,
 $\phi_t \in \mathbb{N}^+, \phi_t[1]$ is the least used location
- allocation weighting
 $\mathbf{a}_t[\phi_t[j]] = (1 - \mathbf{u}_t[\phi_t[j]]) \prod_{i=1}^{j-1} \mathbf{u}_t[\phi_t[i]] \in \Delta_N$

Writing to memory

1. Content-based addressing:

- $\mathcal{C}(M, \mathbf{k}, \beta)[i] = \frac{\exp\{\mathcal{D}(\mathbf{k}, M[i, \cdot])\beta\}}{\sum_j \exp\{\mathcal{D}(\mathbf{k}, M[j, \cdot])\beta\}}$
- cosine similarity $\mathcal{D}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \in [-1, 1]$
- $\mathbf{c}_t^w = \mathcal{C}(M_{t-1}, \mathbf{k}_t^w, \beta_t^w) \in \mathcal{S}_N$

2. Dynamic memory allocation:

- memory retention vector $\psi_t = \prod_{i=1}^R (\mathbf{1} - f_t^i \mathbf{w}_{t-1}^{r,i}) \in [0, 1]^N$
- memory usage vector
 $\mathbf{u}_t = (\mathbf{u}_{t-1} + \mathbf{w}_{t-1}^w - \mathbf{u}_{t-1} \circ \mathbf{w}_{t-1}^w) \circ \psi_t \in [0, 1]^N$
- sort indices of memory locations in ascending order of usage,
 $\phi_t \in \mathbb{N}^+, \phi_t[1]$ is the least used location
- allocation weighting
 $\mathbf{a}_t[\phi_t[j]] = (1 - \mathbf{u}_t[\phi_t[j]]) \prod_{i=1}^{j-1} \mathbf{u}_t[\phi_t[i]] \in \Delta_N$

3. Write weighting:

Writing to memory

1. Content-based addressing:

- $\mathcal{C}(M, \mathbf{k}, \beta)[i] = \frac{\exp\{\mathcal{D}(\mathbf{k}, M[i, \cdot])\beta\}}{\sum_j \exp\{\mathcal{D}(\mathbf{k}, M[j, \cdot])\beta\}}$
- cosine similarity $\mathcal{D}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \in [-1, 1]$
- $\mathbf{c}_t^w = \mathcal{C}(M_{t-1}, \mathbf{k}_t^w, \beta_t^w) \in \mathcal{S}_N$

2. Dynamic memory allocation:

- memory retention vector $\psi_t = \prod_{i=1}^R (\mathbf{1} - f_t^i \mathbf{w}_{t-1}^{r,i}) \in [0, 1]^N$
- memory usage vector
 $\mathbf{u}_t = (\mathbf{u}_{t-1} + \mathbf{w}_{t-1}^w - \mathbf{u}_{t-1} \circ \mathbf{w}_{t-1}^w) \circ \psi_t \in [0, 1]^N$
- sort indices of memory locations in ascending order of usage,
 $\phi_t \in \mathbb{N}^+, \phi_t[1]$ is the least used location
- allocation weighting
 $\mathbf{a}_t[\phi_t[j]] = (1 - \mathbf{u}_t[\phi_t[j]]) \prod_{i=1}^{j-1} \mathbf{u}_t[\phi_t[i]] \in \Delta_N$

3. Write weighting:

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

Writing to memory

1. Content-based addressing:

- $\mathcal{C}(M, \mathbf{k}, \beta)[i] = \frac{\exp\{\mathcal{D}(\mathbf{k}, M[i, \cdot])\beta\}}{\sum_j \exp\{\mathcal{D}(\mathbf{k}, M[j, \cdot])\beta\}}$
- cosine similarity $\mathcal{D}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \in [-1, 1]$
- $\mathbf{c}_t^w = \mathcal{C}(M_{t-1}, \mathbf{k}_t^w, \beta_t^w) \in \mathcal{S}_N$

2. Dynamic memory allocation:

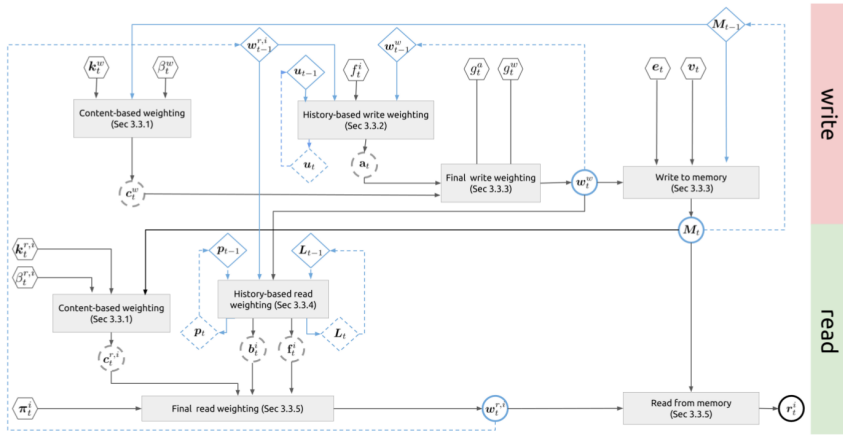
- memory retention vector $\psi_t = \prod_{i=1}^R (\mathbf{1} - f_t^i \mathbf{w}_{t-1}^{r,i}) \in [0, 1]^N$
- memory usage vector
 $\mathbf{u}_t = (\mathbf{u}_{t-1} + \mathbf{w}_{t-1}^w - \mathbf{u}_{t-1} \circ \mathbf{w}_{t-1}^w) \circ \psi_t \in [0, 1]^N$
- sort indices of memory locations in ascending order of usage,
 $\phi_t \in \mathbb{N}^+, \phi_t[1]$ is the least used location
- allocation weighting
 $\mathbf{a}_t[\phi_t[j]] = (1 - \mathbf{u}_t[\phi_t[j]]) \prod_{i=1}^{j-1} \mathbf{u}_t[\phi_t[i]] \in \Delta_N$

3. Write weighting:

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

4. Actual write operation: $M_t = M_{t-1} \circ (E - \mathbf{w}_t^w \mathbf{e}_t^T) + \mathbf{w}_t^w \mathbf{v}_t^T$

Interacting with memory



Source: Hsin, C., *Implementation and Optimization of Differentiable Neural Computers*

1. Content-based addressing:

1. Content-based addressing:

- $\mathbf{c}_t^{r,i} = \mathcal{C}(M_t, \mathbf{k}_t^{r,i}, \beta_t^{r,i}) \in \mathcal{S}_N$

1. Content-based addressing:

- $\mathbf{c}_t^{r,i} = \mathcal{C}(M_t, \mathbf{k}_t^{r,i}, \beta_t^{r,i}) \in \mathcal{S}_N$

2. Temporal memory linkage:

1. Content-based addressing:

- $\mathbf{c}_t^{r,i} = \mathcal{C}(M_t, \mathbf{k}_t^{r,i}, \beta_t^{r,i}) \in \mathcal{S}_N$

2. Temporal memory linkage:

- temporal link matrix $L_t \in [0, 1]^{N \times N}$, $L_t[i, \cdot] \in \Delta_N$, $L_t[\cdot, j] \in \Delta_N$

1. Content-based addressing:

- $\mathbf{c}_t^{r,i} = \mathcal{C}(M_t, \mathbf{k}_t^{r,i}, \beta_t^{r,i}) \in \mathcal{S}_N$

2. Temporal memory linkage:

- temporal link matrix $L_t \in [0, 1]^{N \times N}$, $L_t[i, \cdot] \in \Delta_N$, $L_t[\cdot, j] \in \Delta_N$
- precedence weighting $\mathbf{p}_t = (1 - \sum_i \mathbf{w}_t^w[i]) \mathbf{p}_{t-1} + \mathbf{w}_t^w \in \Delta_N$

Reading from memory

1. Content-based addressing:

- $\mathbf{c}_t^{r,i} = \mathcal{C}(M_t, \mathbf{k}_t^{r,i}, \beta_t^{r,i}) \in \mathcal{S}_N$

2. Temporal memory linkage:

- temporal link matrix $L_t \in [0, 1]^{N \times N}$, $L_t[i, \cdot] \in \Delta_N$, $L_t[\cdot, j] \in \Delta_N$
- precedence weighting $\mathbf{p}_t = (1 - \sum_i \mathbf{w}_t^w[i]) \mathbf{p}_{t-1} + \mathbf{w}_t^w \in \Delta_N$
- linkage logic: $\forall i : L_t[i, i] = 0$, $L_t[i, j] = (1 - \mathbf{w}_t^w[i] - \mathbf{w}_t^w[j]) L_{t-1}[i, j] + \mathbf{w}_t^w[i] \mathbf{p}_{t-1}[j]$

Reading from memory

1. Content-based addressing:

- $\mathbf{c}_t^{r,i} = \mathcal{C}(M_t, \mathbf{k}_t^{r,i}, \beta_t^{r,i}) \in \mathcal{S}_N$

2. Temporal memory linkage:

- temporal link matrix $L_t \in [0, 1]^{N \times N}$, $L_t[i, \cdot] \in \Delta_N$, $L_t[\cdot, j] \in \Delta_N$
- precedence weighting $\mathbf{p}_t = (1 - \sum_i \mathbf{w}_t^w[i]) \mathbf{p}_{t-1} + \mathbf{w}_t^w \in \Delta_N$
- linkage logic: $\forall i : L_t[i, i] = 0$, $L_t[i, j] = (1 - \mathbf{w}_t^w[i] - \mathbf{w}_t^w[j]) L_{t-1}[i, j] + \mathbf{w}_t^w[i] \mathbf{p}_{t-1}[j]$
- forward weighting: $\mathbf{f}_t^i = L_t \mathbf{w}_{t-1}^{r,i} \in \Delta_N$

Reading from memory

1. Content-based addressing:

- $\mathbf{c}_t^{r,i} = \mathcal{C}(M_t, \mathbf{k}_t^{r,i}, \beta_t^{r,i}) \in \mathcal{S}_N$

2. Temporal memory linkage:

- temporal link matrix $L_t \in [0, 1]^{N \times N}$, $L_t[i, \cdot] \in \Delta_N$, $L_t[\cdot, j] \in \Delta_N$
- precedence weighting $\mathbf{p}_t = (1 - \sum_i \mathbf{w}_t^w[i]) \mathbf{p}_{t-1} + \mathbf{w}_t^w \in \Delta_N$
- linkage logic: $\forall i : L_t[i, i] = 0$, $L_t[i, j] = (1 - \mathbf{w}_t^w[i] - \mathbf{w}_t^w[j]) L_{t-1}[i, j] + \mathbf{w}_t^w[i] \mathbf{p}_{t-1}[j]$
- forward weighting: $\mathbf{f}_t^i = L_t \mathbf{w}_{t-1}^{r,i} \in \Delta_N$
- backward weighting: $\mathbf{b}_t^i = L_t^T \mathbf{w}_{t-1}^{r,i} \in \Delta_N$

Reading from memory

1. Content-based addressing:

- $\mathbf{c}_t^{r,i} = \mathcal{C}(M_t, \mathbf{k}_t^{r,i}, \beta_t^{r,i}) \in \mathcal{S}_N$

2. Temporal memory linkage:

- temporal link matrix $L_t \in [0, 1]^{N \times N}$, $L_t[i, \cdot] \in \Delta_N$, $L_t[\cdot, j] \in \Delta_N$
- precedence weighting $\mathbf{p}_t = (1 - \sum_i \mathbf{w}_t^w[i]) \mathbf{p}_{t-1} + \mathbf{w}_t^w \in \Delta_N$
- linkage logic: $\forall i : L_t[i, i] = 0$, $L_t[i, j] = (1 - \mathbf{w}_t^w[i] - \mathbf{w}_t^w[j]) L_{t-1}[i, j] + \mathbf{w}_t^w[i] \mathbf{p}_{t-1}[j]$
- forward weighting: $\mathbf{f}_t^i = L_t \mathbf{w}_{t-1}^{r,i} \in \Delta_N$
- backward weighting: $\mathbf{b}_t^i = L_t^T \mathbf{w}_{t-1}^{r,i} \in \Delta_N$

3. Read weighting:

Reading from memory

1. Content-based addressing:

- $\mathbf{c}_t^{r,i} = \mathcal{C}(M_t, \mathbf{k}_t^{r,i}, \beta_t^{r,i}) \in \mathcal{S}_N$

2. Temporal memory linkage:

- temporal link matrix $L_t \in [0, 1]^{N \times N}$, $L_t[i, \cdot] \in \Delta_N$, $L_t[\cdot, j] \in \Delta_N$
- precedence weighting $\mathbf{p}_t = (1 - \sum_i \mathbf{w}_t^w[i]) \mathbf{p}_{t-1} + \mathbf{w}_t^w \in \Delta_N$
- linkage logic: $\forall i : L_t[i, i] = 0$, $L_t[i, j] = (1 - \mathbf{w}_t^w[i] - \mathbf{w}_t^w[j]) L_{t-1}[i, j] + \mathbf{w}_t^w[i] \mathbf{p}_{t-1}[j]$
- forward weighting: $\mathbf{f}_t^i = L_t \mathbf{w}_{t-1}^{r,i} \in \Delta_N$
- backward weighting: $\mathbf{b}_t^i = L_t^T \mathbf{w}_{t-1}^{r,i} \in \Delta_N$

3. Read weighting:

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

Reading from memory

1. Content-based addressing:

- $\mathbf{c}_t^{r,i} = \mathcal{C}(M_t, \mathbf{k}_t^{r,i}, \beta_t^{r,i}) \in \mathcal{S}_N$

2. Temporal memory linkage:

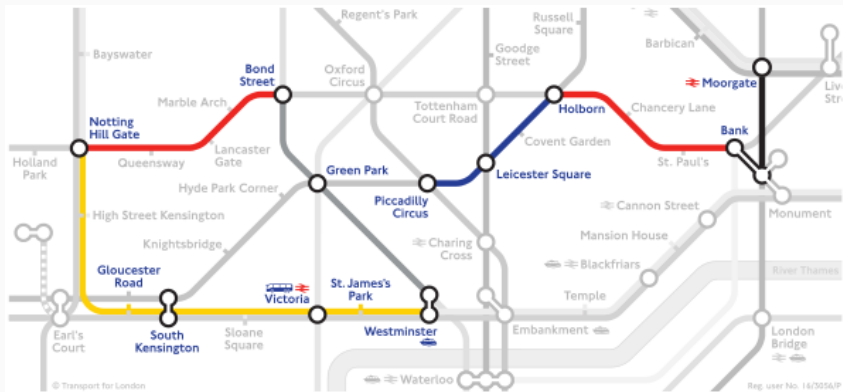
- temporal link matrix $L_t \in [0, 1]^{N \times N}$, $L_t[i, \cdot] \in \Delta_N$, $L_t[\cdot, j] \in \Delta_N$
- precedence weighting $\mathbf{p}_t = (1 - \sum_i \mathbf{w}_t^w[i]) \mathbf{p}_{t-1} + \mathbf{w}_t^w \in \Delta_N$
- linkage logic: $\forall i : L_t[i, i] = 0$, $L_t[i, j] = (1 - \mathbf{w}_t^w[i] - \mathbf{w}_t^w[j]) L_{t-1}[i, j] + \mathbf{w}_t^w[i] \mathbf{p}_{t-1}[j]$
- forward weighting: $\mathbf{f}_t^i = L_t \mathbf{w}_{t-1}^{r,i} \in \Delta_N$
- backward weighting: $\mathbf{b}_t^i = L_t^T \mathbf{w}_{t-1}^{r,i} \in \Delta_N$

3. Read weighting:

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

4. Actual read operation: $\mathbf{r}_t^i = M_t^T \mathbf{w}_t^{r,i}$.

Traversing London Underground



Source: [Graves et al., 2016]

Traversing London Underground

- London Underground as a graph.

Traversing London Underground

- London Underground as a graph.
- Explicit vector representation of an edge:
 $\left[\text{station}_1 \quad \text{station}_2 \quad \text{line} \right]$

Traversing London Underground

- London Underground as a graph.
- Explicit vector representation of an edge:
[station₁ station₂ line]
- Queries: traversal, shortest path.

Traversing London Underground

- London Underground as a graph.
- Explicit vector representation of an edge:
 $\left[\text{station}_1 \quad \text{station}_2 \quad \text{line} \right]$
- Queries: traversal, shortest path.
- Training: graphs with random nodes and connections.

Traversing London Underground

- London Underground as a graph.
- Explicit vector representation of an edge:
 $\left[\text{station}_1 \quad \text{station}_2 \quad \text{line} \right]$
- Queries: traversal, shortest path.
- Training: graphs with random nodes and connections.
- Curriculum learning with increasing complexity of graphs and queries.

Traversing London Underground

- London Underground as a graph.
- Explicit vector representation of an edge:
 $\left[\text{station}_1 \quad \text{station}_2 \quad \text{line} \right]$
- Queries: traversal, shortest path.
- Training: graphs with random nodes and connections.
- Curriculum learning with increasing complexity of graphs and queries.
- Tested without re-training on the London Underground graph.

Traversal

Traversal question:

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

Answer:

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

Shortest-path

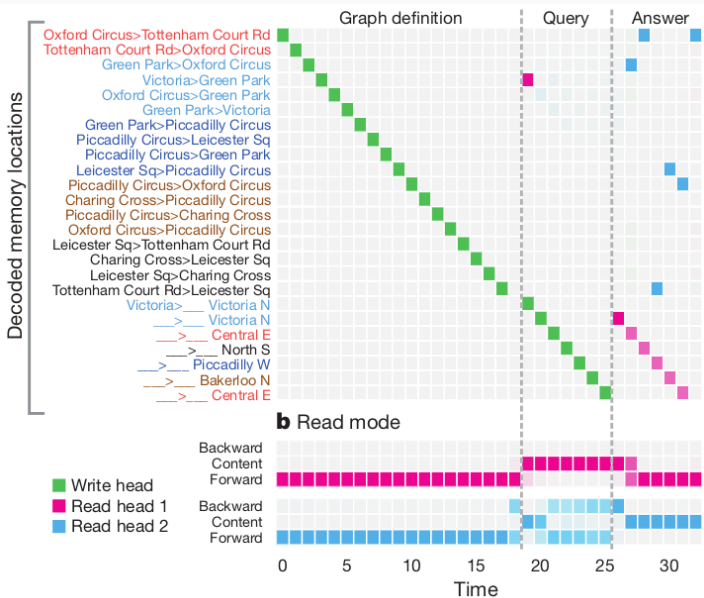
Shortest-path question:

(Moorgate, PiccadillyCircus, _)

Answer:

(Moorgate, Bank, Northern)
(Bank, Holborn, Central)
(Holborn, LeicesterSq, Piccadilly)
(LeicesterSq, PiccadillyCircus, Piccadilly)

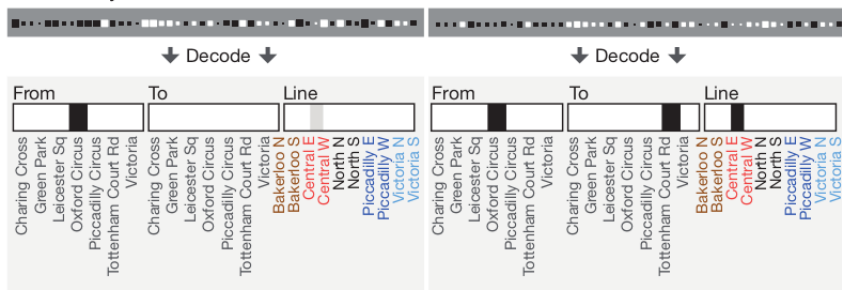
Traversal



Traversal

d Read key

e Location content



Source: [Graves et al., 2016]

Further research

- Synthetic gradients [Jaderberg et al., 2016].

Further research

- Synthetic gradients [Jaderberg et al., 2016].
- Speed up training.

Further research

- Synthetic gradients [Jaderberg et al., 2016].
- Speed up training.
- DNC with other types of neural networks.

Further research

- Synthetic gradients [Jaderberg et al., 2016].
- Speed up training.
- DNC with other types of neural networks.
- Scale up.

Further research

- Synthetic gradients [Jaderberg et al., 2016].
- Speed up training.
- DNC with other types of neural networks.
- Scale up.
- Tasks beyond graphs.



Graves, A., Wayne, G., et al. (2016).

Hybrid computing using a neural network with dynamic external memory.

Nature, 538:471–476.



Hochreiter, S. (1991).

Untersuchungen zu dynamischen neuronalen netzen.

Diploma thesis, Technical University Munich.



Hochreiter, S. and Schmidhuber, J. (1997).

Long short-term memory.

Neural Computation, 9(8):1735–1780.



Jaderberg, M., Czarnecki, W. M., et al. (2016).

Decoupled neural interfaces using synthetic gradients.

arXiv.