# Article: **Hybrid Computing Using a Neural Network with Dynamic External Memory**

Speaker: Valeria Bubnova

# Prepare to me amazed…

1) What's the idea?

2) Does it even work?
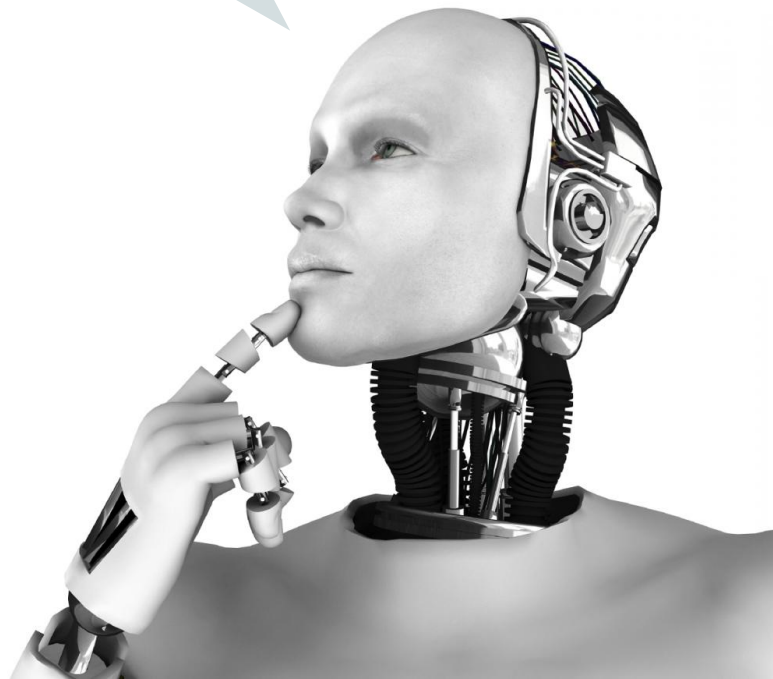
3) Details (you'll need them!)

# Neural Networks

What about… me?

Store weights

Learn from data

Learn some distributions

Weak with algorithms & data  structures

# Neural Networks

# Computers

| Neural Networks | Computers |
|---|---|
| Store weights | Processor and RAM |
| Learn from data | Passively store data |
| Learn some distributions | All data stored equally |
| Weak with algorithms & data structures | Can process data structures |

# Differentiable Neural Computers

Store weights

Learn from data

Learn some distributions

~~Weak with algorithms & data structures~~

~~Processor and RAM~~  Controller and Memory matrix

~~Passively~~ store data

~~All data stored equally~~  Data is stored with some distribution

Can process data structures

#NeuralTuringMachine

# How?
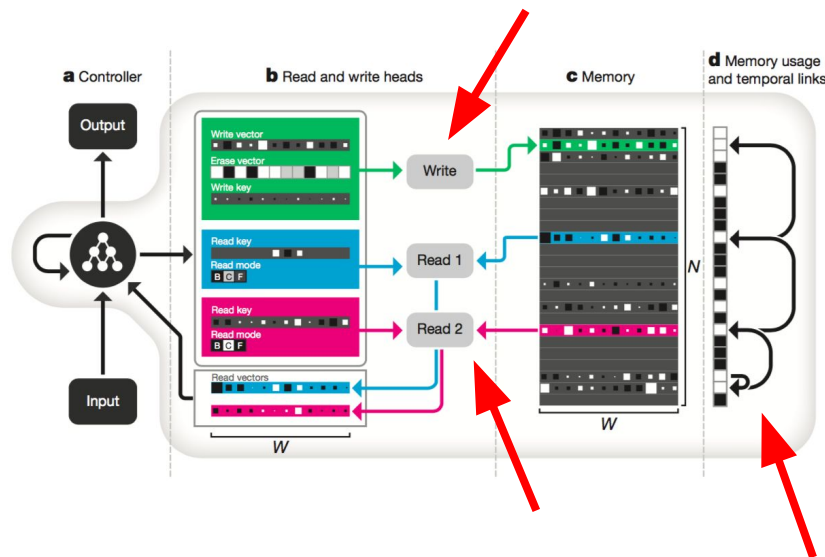
Differentiable everything!

Yes, and memory as well....

# 3 forms of <differentiable> attention

- Content Lookup

- Temporal Linking

- Memory Allocation

# 3 forms of \<differentiable\> attention

#Hippocampus
#CA3_synapses
#CA1_synapses

| Attention mechanism | Computational considerations |
|---|---|
| Content Lookup | Formation of associative data structures |
| Temporal Linking | Sequential retrieval of input sequences |
| Memory Allocation | Provides the write head with unused locations |

**a** Random graph

**b** London Underground

Traversal

Shortest-path

Underground input:
(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)

84 edges in total

Traversal question:
(BondSt, _, Central),
(_, _, Circle), (_, _, Circle),
(_, _, Circle), (_, _, Circle),
(_, _, Jubilee), (_, _, Jubilee),

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

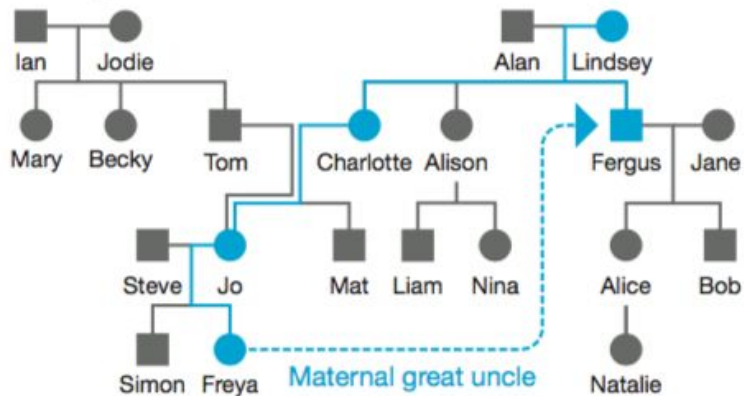Shortest-path question:
(Moorgate, PiccadillyCircus, _)

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

LSTM: 37% accuracy
    after 2 million examples

DNC: 98.8% accuracy
    after 1 million examples

**c** Family tree

Maternal great uncle

Family tree input:
(Charlotte, Alan, Father)
(Simon, Steve, Father)
(Steve , Simon, Son1)
(Nina, Alison, Mother)
(Lindsey, Fergus, Son1)
⋮
(Bob, Jane, Mother)
(Natalie, Alice, Mother)
(Mary, Ian, Father)
(Jane, Alice, Daughter1)
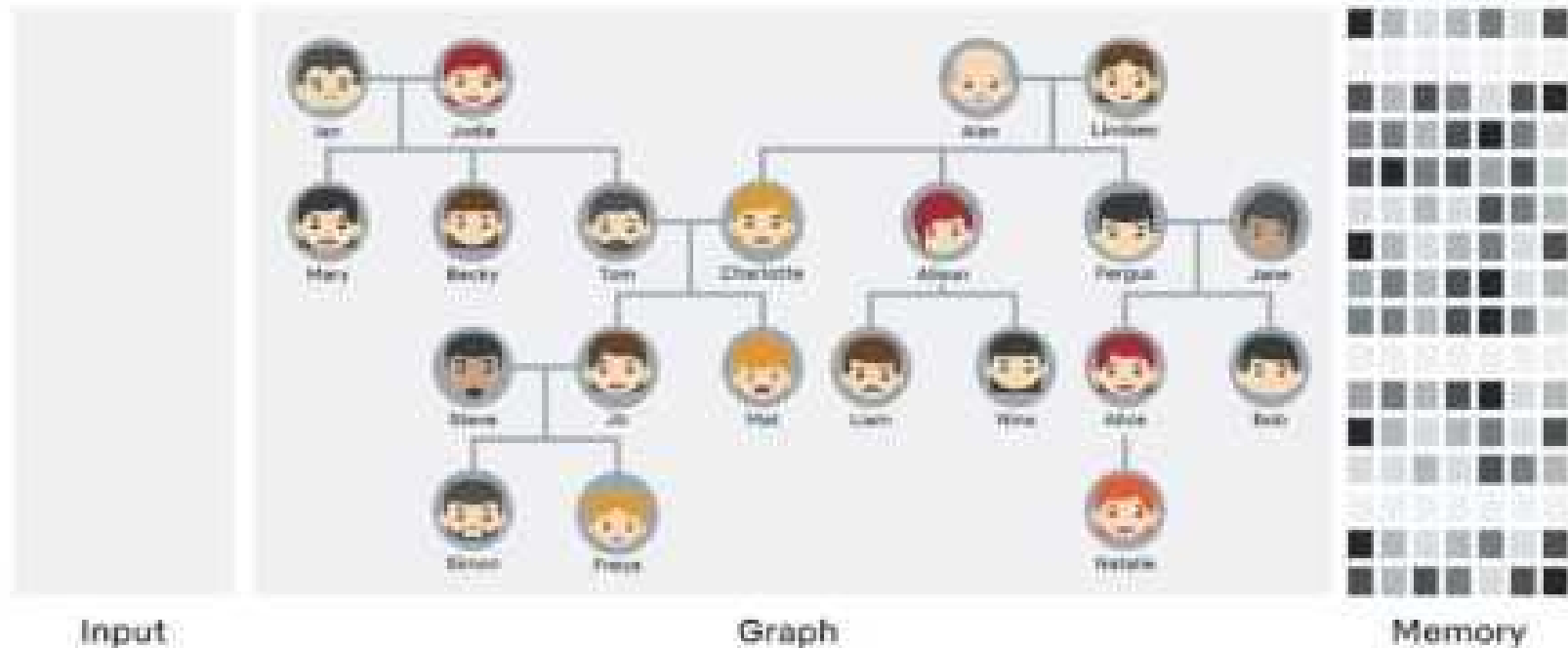(Mat, Charlotte, Mother)

54 edges in total

Inference question:
(Freya, _, MaternalGreatUncle)

Answer:
(Freya, Fergus, MaternalGreatUncle)

DNC: 81.8% accuracy
on four-step relations

Input                                   Graph                                   Memory
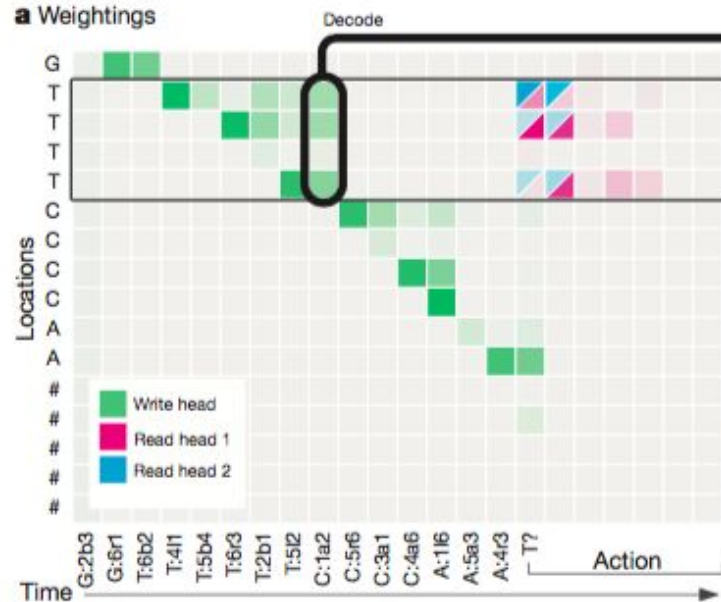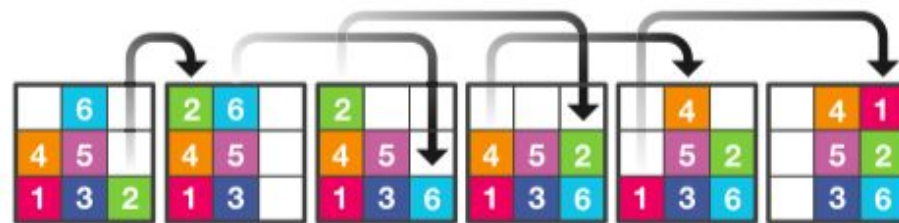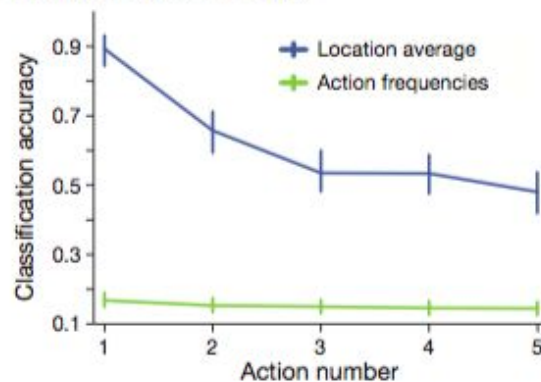
Sheep are afraid of wolves.

Gertrude is a sheep.

Mice are afraid of cats.

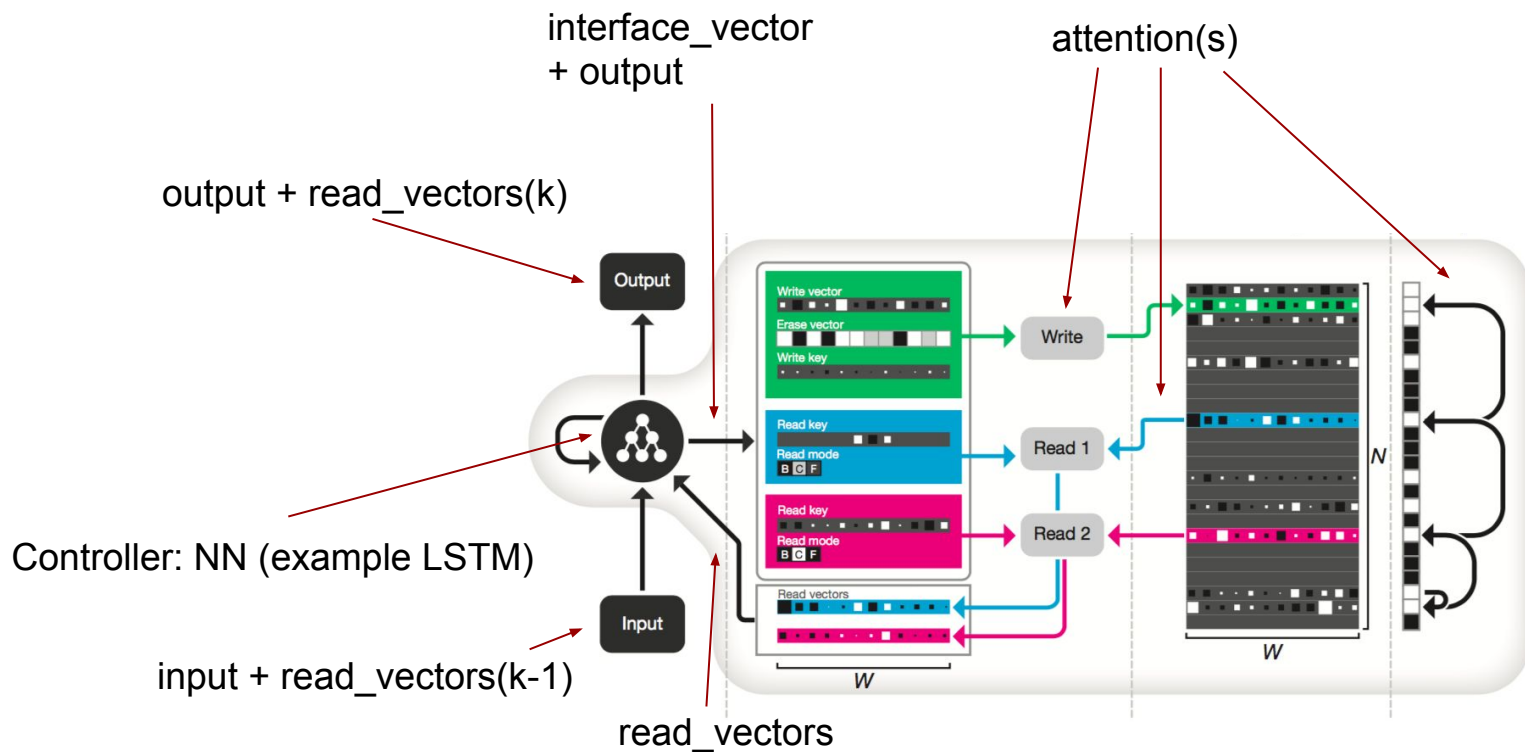What is Gertrude afraid of?

LSTM: 7.5% mean test error

DNC: 3.8% mean test error

**a** Weightings

Decode

Locations

- Write head
- Read head 1
- Read head 2

Time

G:2b3  G:6r1  T:6b2  T:4l1  T:5b4  T:6r3  T:2b1  T:5l2  C:1a2  C:5r6  C:3a1  C:4a6  A:1l6  A:5a3  A:4r3  T?

Action

**b** Goal T constraints

| 2 | | 4 | | 4 | | 3 | | 1 | | 5 |
| 6 | | 4 1 | | 5 | | 3 6 | | 2 | | 5 2 |
| 6b2 | | 4l1 | | 5b4 | | 6r3 | | 2b1 | | 5l2 |

**c** Board states

**d** Planned action decodings

Classification accuracy

- Location average
- Action frequencies

0.9
0.7
0.5
0.3
0.1

1    2    3    4    5

Action number

**e** t-SNE location goal labels

# DETAILS

# DETAILS



Controller: NN (example LSTM)

input
$$i_t^l = \sigma(W_i^l[\chi_t; h_{t-1}^l; h_t^{l-1}] + b_i^l)$$

forget
$$f_t^l = \sigma(W_f^l[\chi_t; h_{t-1}^l; h_t^{l-1}] + b_f^l)$$

state
$$s_t^l = f_t^l s_{t-1}^l + i_t^l \tanh(W_s^l[\chi_t; h_{t-1}^l; h_t^{l-1}] + b_s^l)$$

output gate activation
$$o_t^l = \sigma(W_o^l[\chi_t; h_{t-1}^l; h_t^{l-1}] + b_o^l)$$

hidden
$$h_t^l = o_t^l \tanh(s_t^l)$$

# DETAILS



interface_vector + output

attention(s)

output + read_vectors(k)

Controller: NN (example LSTM)

input + read_vectors(k-1)

read_vectors

interface_vector
+ output

# interface_vector

**Interface parameters.** Before being used to parameterize the memory interactions, the interface vector $\boldsymbol{\xi}_t$ is subdivided as follows:

$$\boldsymbol{\xi}_t = \left[ \boldsymbol{k}_t^{\mathrm{r},1}; \ldots; \boldsymbol{k}_t^{\mathrm{r},R}; \hat{\beta}_t^{\mathrm{r},1}; \ldots; \hat{\beta}_t^{\mathrm{r},R}; \boldsymbol{k}_t^{\mathrm{w}}; \hat{\beta}_t^{\mathrm{w}}; \hat{\boldsymbol{e}}_t; \boldsymbol{v}_t; \hat{f}_t^1; \ldots; \hat{f}_t^R; \hat{g}_t^{\mathrm{a}}; \hat{g}_t^{\mathrm{w}}; \hat{\boldsymbol{\pi}}_t^1; \ldots; \hat{\boldsymbol{\pi}}_t^R \right]$$

- the write key $\boldsymbol{k}_t^{\mathrm{w}} \in \mathbb{R}^W$;

- the write strength $\beta_t^{\mathrm{w}} = \mathrm{oneplus}(\hat{\beta}_t^{\mathrm{w}}) \in [1, \infty)$;

- the erase vector $\boldsymbol{e}_t = \sigma(\hat{\boldsymbol{e}}_t) \in [0,1]^W$;

- the write vector $\boldsymbol{v}_t \in \mathbb{R}^W$;

- $R$ free gates $\{f_t^i = \sigma(\hat{f}_t^i) \in [0,1]; 1 \le i \le R\}$;

- the allocation gate $g_t^{\mathrm{a}} = \sigma(\hat{g}_t^{\mathrm{a}}) \in [0,1]$;

- the write gate $g_t^{\mathrm{w}} = \sigma(\hat{g}_t^{\mathrm{w}}) \in [0,1]$; and

- $R$ read modes $\{\boldsymbol{\pi}_t^i = \mathrm{softmax}(\hat{\boldsymbol{\pi}}_t^i) \in \mathcal{S}_3; 1 \le i \le R\}$.
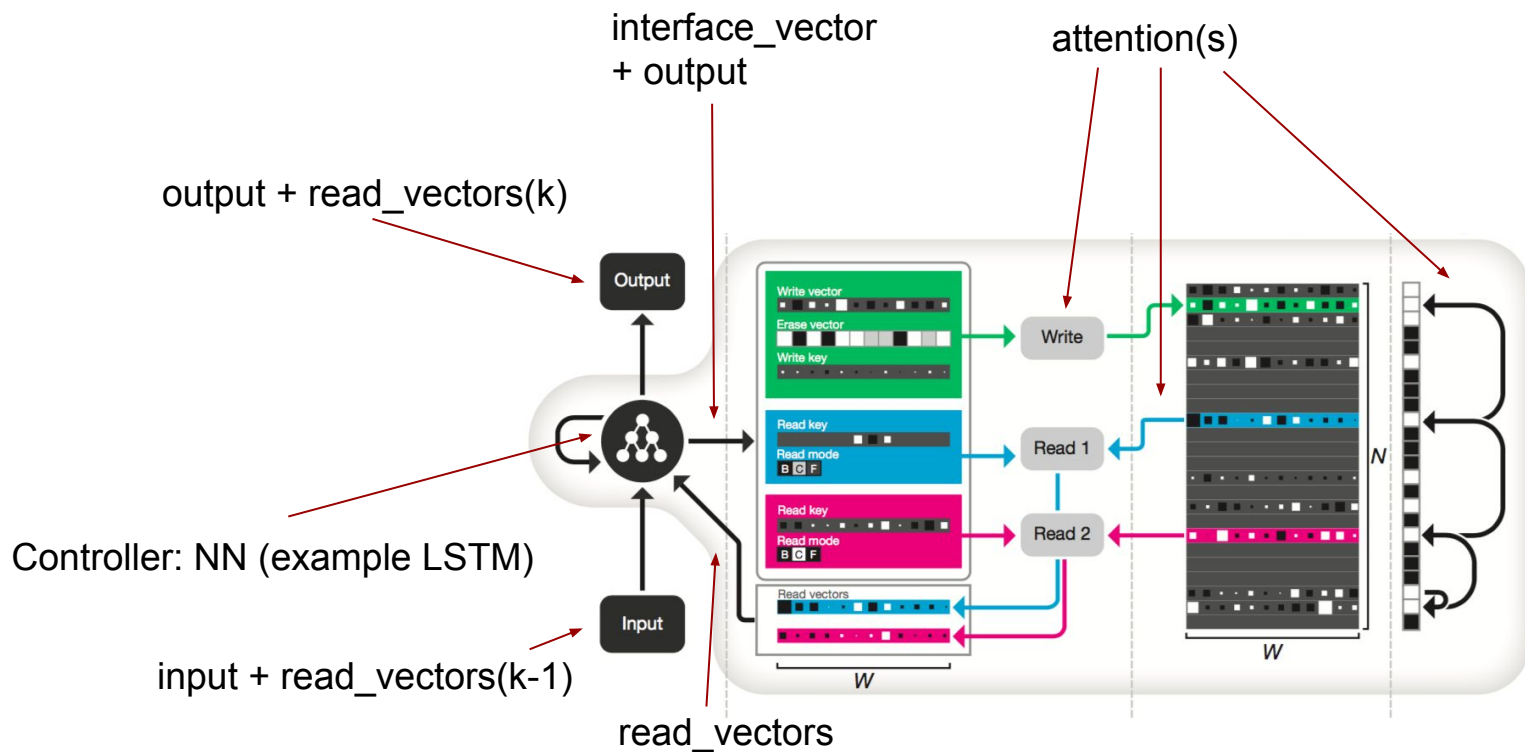
# interface_vector

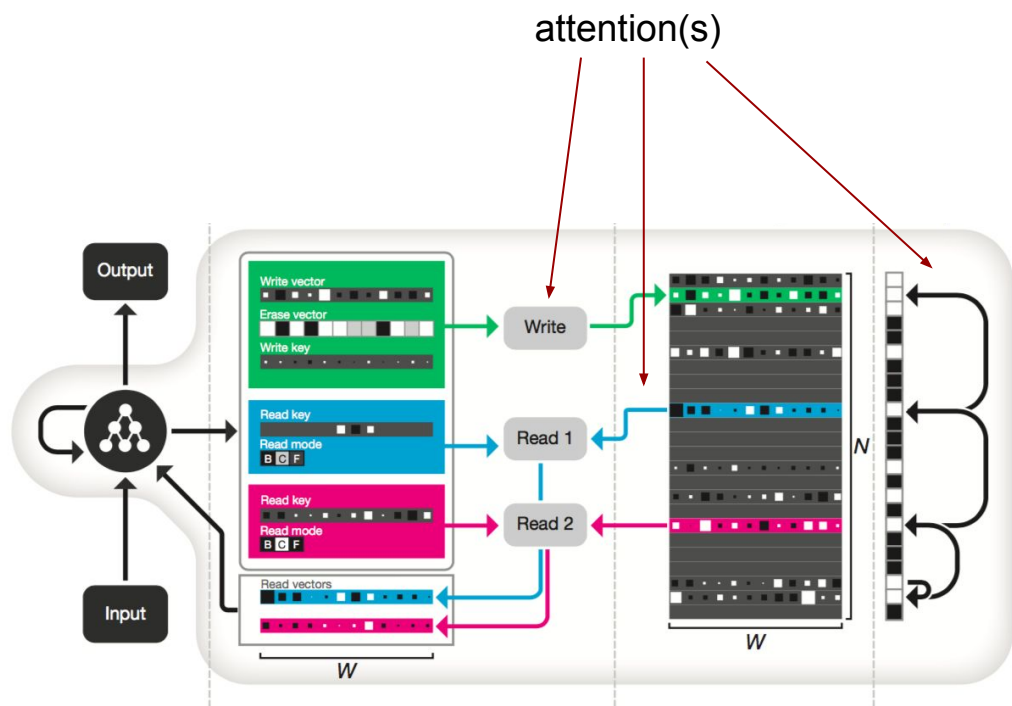**Interface parameters.** Before being used to parameterize the memory interactions, the interface vector $\boldsymbol{\xi}_t$ is subdivided as follows:

$$\boldsymbol{\xi}_t = \left[ \boldsymbol{k}_t^{\text{r},1}; \ldots; \boldsymbol{k}_t^{\text{r},R}; \hat{\beta}_t^{\text{r},1}; \ldots; \hat{\beta}_t^{\text{r},R}; \boldsymbol{k}_t^{\text{w}}; \hat{\beta}_t^{\text{w}}; \hat{\boldsymbol{e}}_t; \boldsymbol{v}_t; \hat{f}_t^1; \ldots; \hat{f}_t^R; \hat{g}_t^{\text{a}}; \hat{g}_t^{\text{w}}; \hat{\boldsymbol{\pi}}_t^1; \ldots; \hat{\boldsymbol{\pi}}_t^R \right]$$

$$\boldsymbol{v}_t = W_y [\boldsymbol{h}_t^1; \ldots; \boldsymbol{h}_t^L]$$

$$\boldsymbol{\xi}_t = W_\xi [\boldsymbol{h}_t^1; \ldots; \boldsymbol{h}_t^L]$$

# DETAILS



interface_vector + output

attention(s)

output + read_vectors(k)

Controller: NN (example LSTM)

input + read_vectors(k-1)

read_vectors

attention(s)

Content Lookup

Memory Allocation

Temporal Linking

# Content lookup (attention 1)

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

where $\boldsymbol{k} \in \mathbb{R}^W$ is a lookup key, $\beta \in [1, \infty)$ is a scalar representing key strength and $\mathcal{D}$ is the cosine similarity:

$$\mathcal{D}(\boldsymbol{u}, \boldsymbol{v}) = \frac{\boldsymbol{u} \cdot \boldsymbol{v}}{|\boldsymbol{u}||\boldsymbol{v}|}$$

# Memory allocation (attention 2)

$$\boldsymbol{\psi}_t = \prod_{i=1}^{R} \left( \mathbf{1} - f_t^i \mathbf{w}_{t-1}^{\mathrm{r},i} \right)$$

$$\boldsymbol{u}_t = (\boldsymbol{u}_{t-1} + \boldsymbol{w}_{t-1}^{\mathrm{w}} - \boldsymbol{u}_{t-1} \circ \boldsymbol{w}_{t-1}^{\mathrm{w}}) \circ \boldsymbol{\psi}_t$$

$$\boldsymbol{a}_t[\boldsymbol{\phi}_t[j]] = (1 - \boldsymbol{u}_t[\boldsymbol{\phi}_t[j]]) \prod_{i=1}^{j-1} \boldsymbol{u}_t[\boldsymbol{\phi}_t[i]]$$

# Memory allocation (attention 2)

Free gates
(from interface vector)

Retention vector

$$\boldsymbol{\psi}_t = \prod_{i=1}^{R} \left( \mathbf{1} - f_t^i \mathbf{w}_{t-1}^{\mathbf{r},i} \right)$$

Weights
(learnable)

Element-wise product

Memory usage vector

$$\boldsymbol{u}_t = (\boldsymbol{u}_{t-1} + \boldsymbol{w}_{t-1}^{\mathbf{w}} - \boldsymbol{u}_{t-1} \circ \boldsymbol{w}_{t-1}^{\mathbf{w}}) \circ \boldsymbol{\psi}_t$$

Free list

$\phi_t[1]$ is the index
of the least used location

$$\boldsymbol{a}_t[\phi_t[j]] = (1 - \boldsymbol{u}_t[\phi_t[j]]) \prod_{i=1}^{j-1} \boldsymbol{u}_t[\phi_t[i]]$$

Allocation weighting

*Write weighting.* The controller can write to newly allocated locations, or to locations addressed by content, or it can choose not to write at all. First, a write content weighting $c_t^w \in \mathcal{S}_N$ is constructed using the write key $k_t^w$ and write strength $\beta_t^w$:

$$c_t^w = \mathcal{C}(M_{t-1}, k_t^w, \beta_t^w)$$

$c_t^w$ is interpolated with the allocation weighting $a_t$ defined in equation (1) to determine a write weighting $w_t^w \in \Delta_N$:

$$w_t^w = g_t^w \left[ g_t^a a_t + (1 - g_t^a) c_t^w \right] \qquad (2)$$

where $g_t^a \in [0,1]$ is the allocation gate governing the interpolation and $g_t^w \in [0,1]$ is the write gate. If the write gate is 0, then nothing is written, regardless of the other write parameters; it can therefore be used to protect the memory from unnecessary modifications.

# Temporal linking (attention 3)

$$p_0 = 0$$

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

Precedence weights (writing order)

$$L_0[i,j] = 0 \quad \forall i, j$$

$$L_t[i,i] = 0 \quad \forall i$$

$$L_t[i,j] = (1 - w_t^{\mathrm{w}}[i] - w_t^{\mathrm{w}}[j])L_{t-1}[i,j] + w_t^{\mathrm{w}}[i]p_{t-1}[j]$$

Temporal linking matrix

*Read weighting.* Each read head $i$ computes a content weighting $\mathbf{c}_t^{r,i} \in \Delta_N$ using a read key $\boldsymbol{k}_t^{r,i} \in \mathbb{R}^W$:
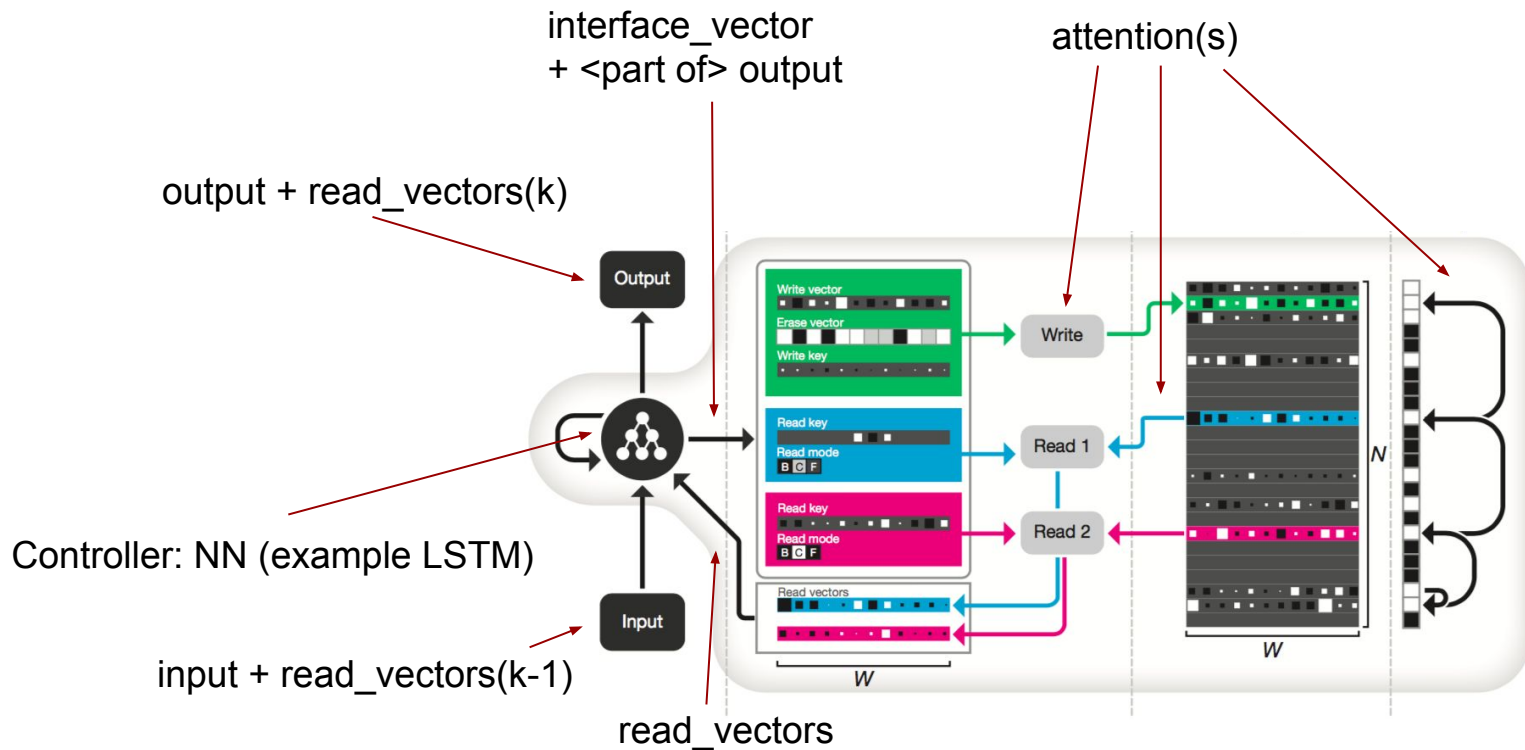
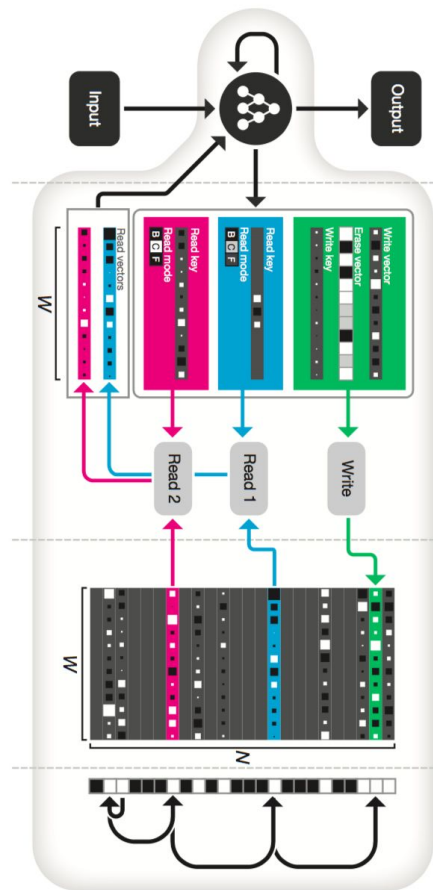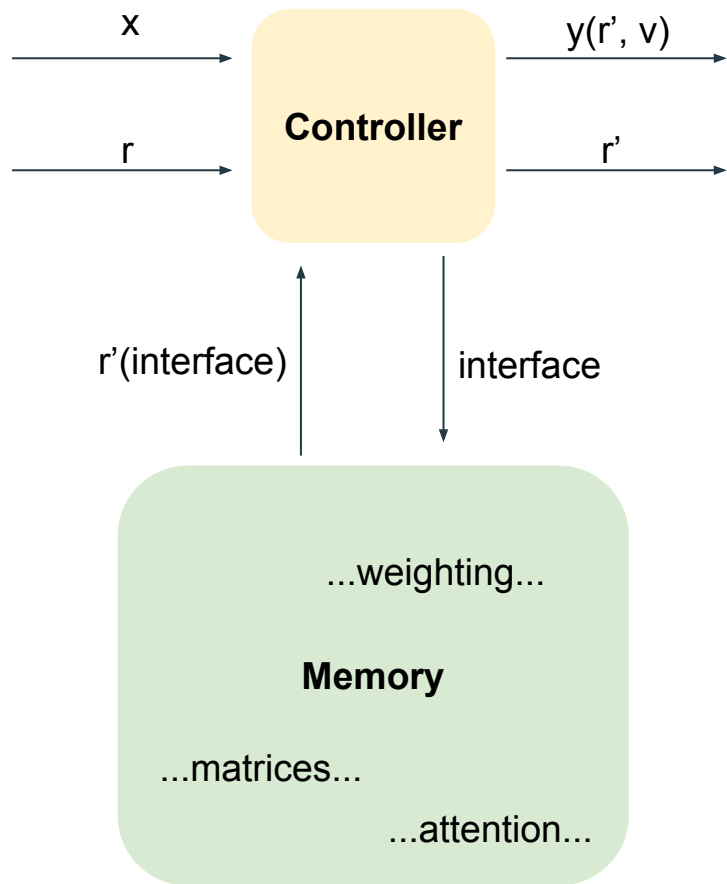$$\mathbf{c}_t^{r,i} = \mathcal{C}(M_t, \boldsymbol{k}_t^{r,i}, \beta_t^{r,i})$$

Each read head also receives a read mode vector $\boldsymbol{\pi}_t^i \in \mathcal{S}_3$, which interpolates among the backward weighting $\boldsymbol{b}_t^i$, the forward weighting $\boldsymbol{f}_t^i$ and the content read weighting $\boldsymbol{c}_t^{r,i}$, thereby determining the read weighting $\boldsymbol{w}_t^{r,i} \in \mathcal{S}_3$:

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

If $\boldsymbol{\pi}_t^i[2]$ dominates the read mode, then the weighting reverts to content lookup using $\boldsymbol{k}_t^{r,i}$. If $\boldsymbol{\pi}_t^i[3]$ dominates, then the read head iterates through memory locations in the order they were written, ignoring the read key. If $\boldsymbol{\pi}_t^i[1]$ dominates, then the read head iterates in the reverse order.
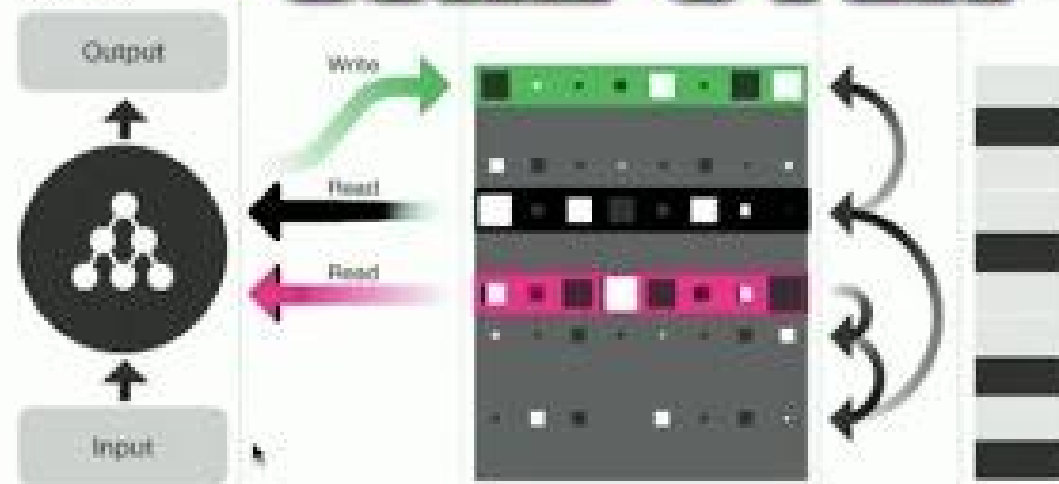
# DETAILS



interface_vector
+ <part of> output

attention(s)

output + read_vectors(k)

Controller: NN (example LSTM)

input + read_vectors(k-1)

read_vectors

# DIFFERENTIABLE NEURAL COMPUTER

Controller

Output

Write

Read

Read

Input

this consists of a neural network that can read from and write to an external memory analogous to the random-access memory in a conventional computer.