

3. Reasoning as Memory

Introduction

Item memory

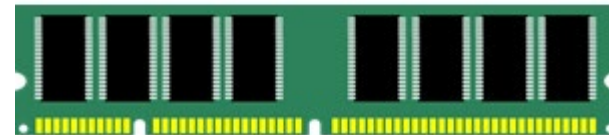
Relational memory

Program memory

Introduction

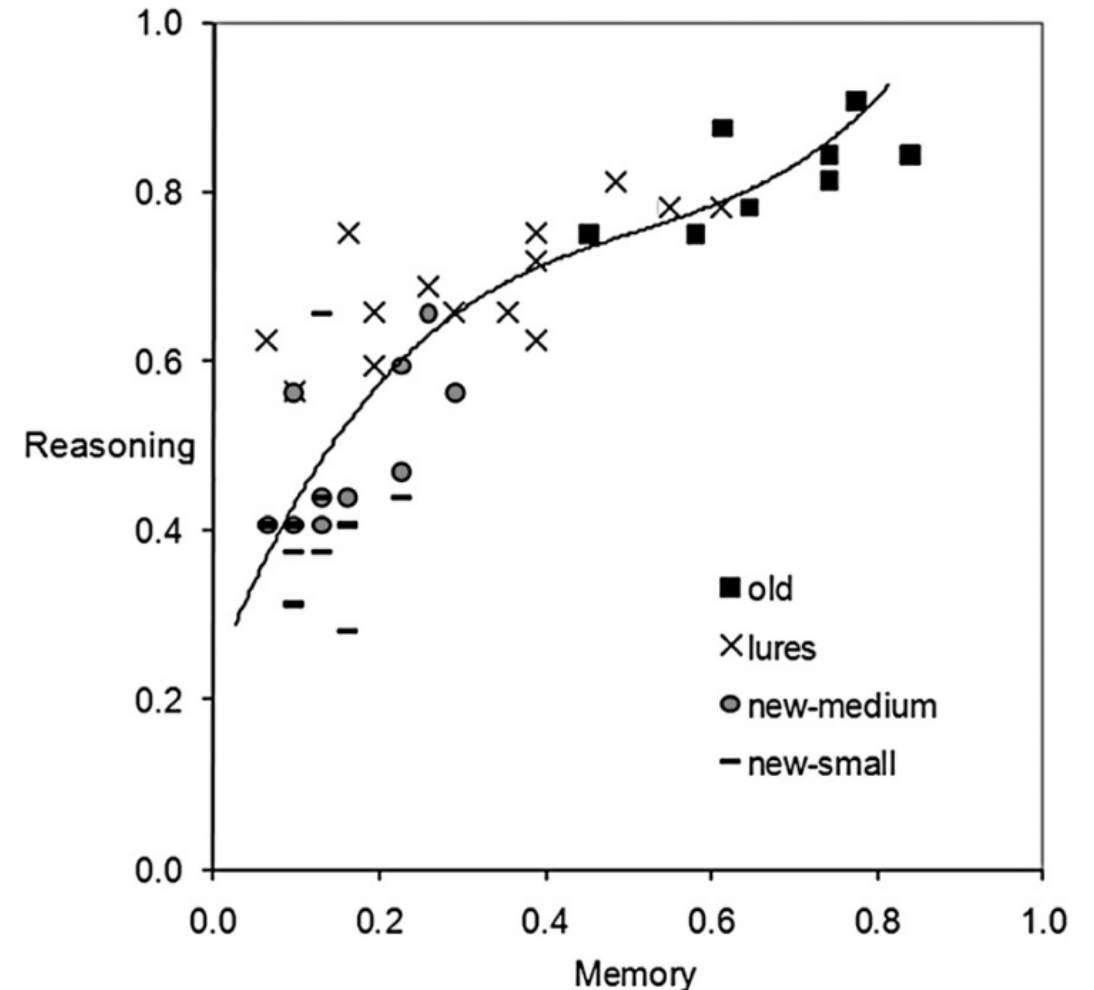
Memory is part of intelligence

- Memory is the ability to **store**, **retain** and **recall** information
- Brain memory stores items, events and high-level structures
- Computer memory stores data and temporary variables



Memory-reasoning analogy

- 2 processes: fast-slow
 - Memory: familiarity-recollection
- Cognitive test:
 - Corresponding reasoning and memorization performance
 - Increasing # premises, inductive/deductive reasoning is affected



Common memory activities

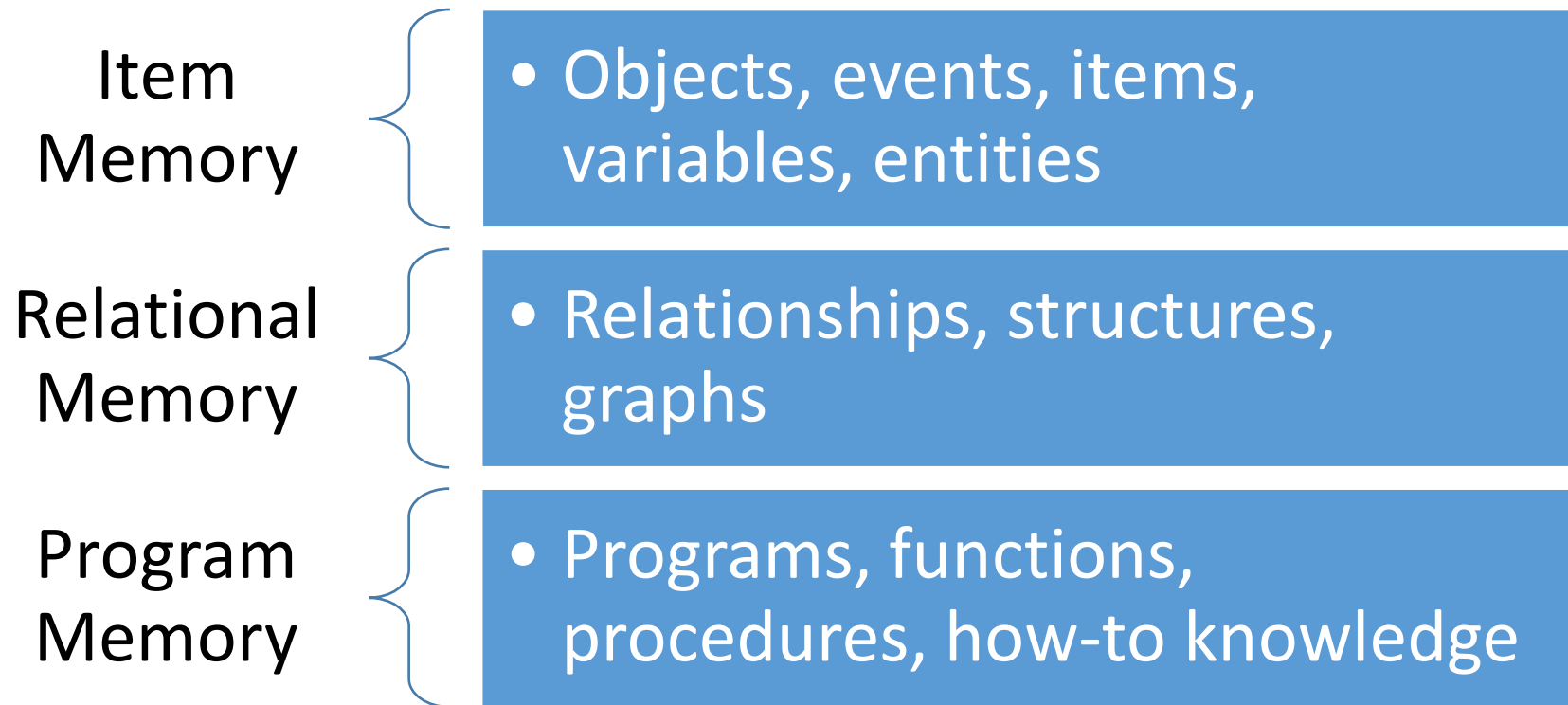
- **Encode**: write information to the memory, often requiring compression capability
- **Retain**: keep the information overtime. This is often assumed in machinery memory
- **Retrieve**: read information from the memory to solve the task at hand

Encode

Retain

Retrieve

Memory taxonomy based on memory content



Item memory

Associative memory

RAM-like memory

Independent memory

Distributed item memory as associative memory

Language

"Green" means
"go," but what
does "red" mean?

Semantic
memory

Time

birthday party on
30th Jan

Episodic
memory

Object

Where is my pen?
What is the
password?

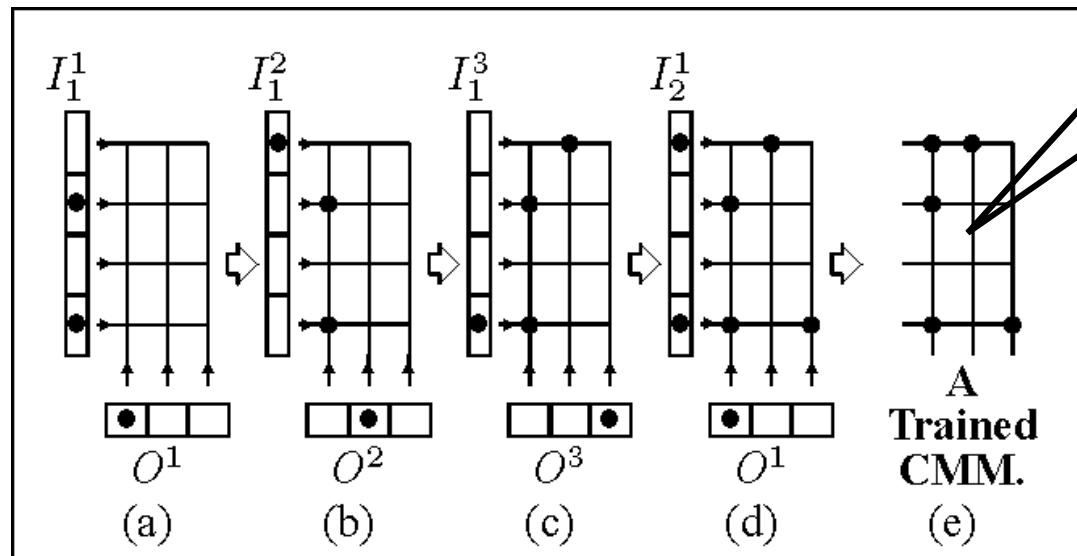
Working
memory

Behaviour



Motor
memory

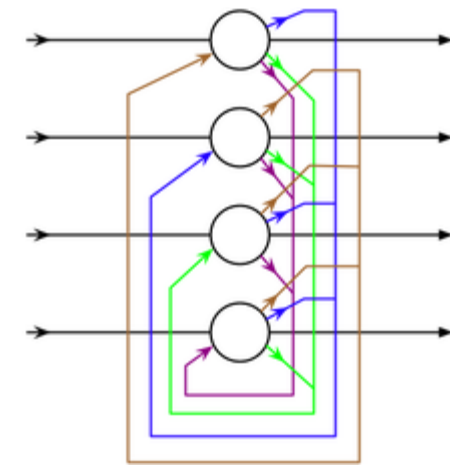
Associate memory can be implemented as Hopfield network



Correlation matrix memory

Encode $\hat{\mathbf{M}} = \sum_{k=1}^q \mathbf{b}_k \mathbf{a}_k^T$ Retrieve $\mathbf{b} = \hat{\mathbf{M}} \mathbf{a}_j$

Feed-forward retrieval



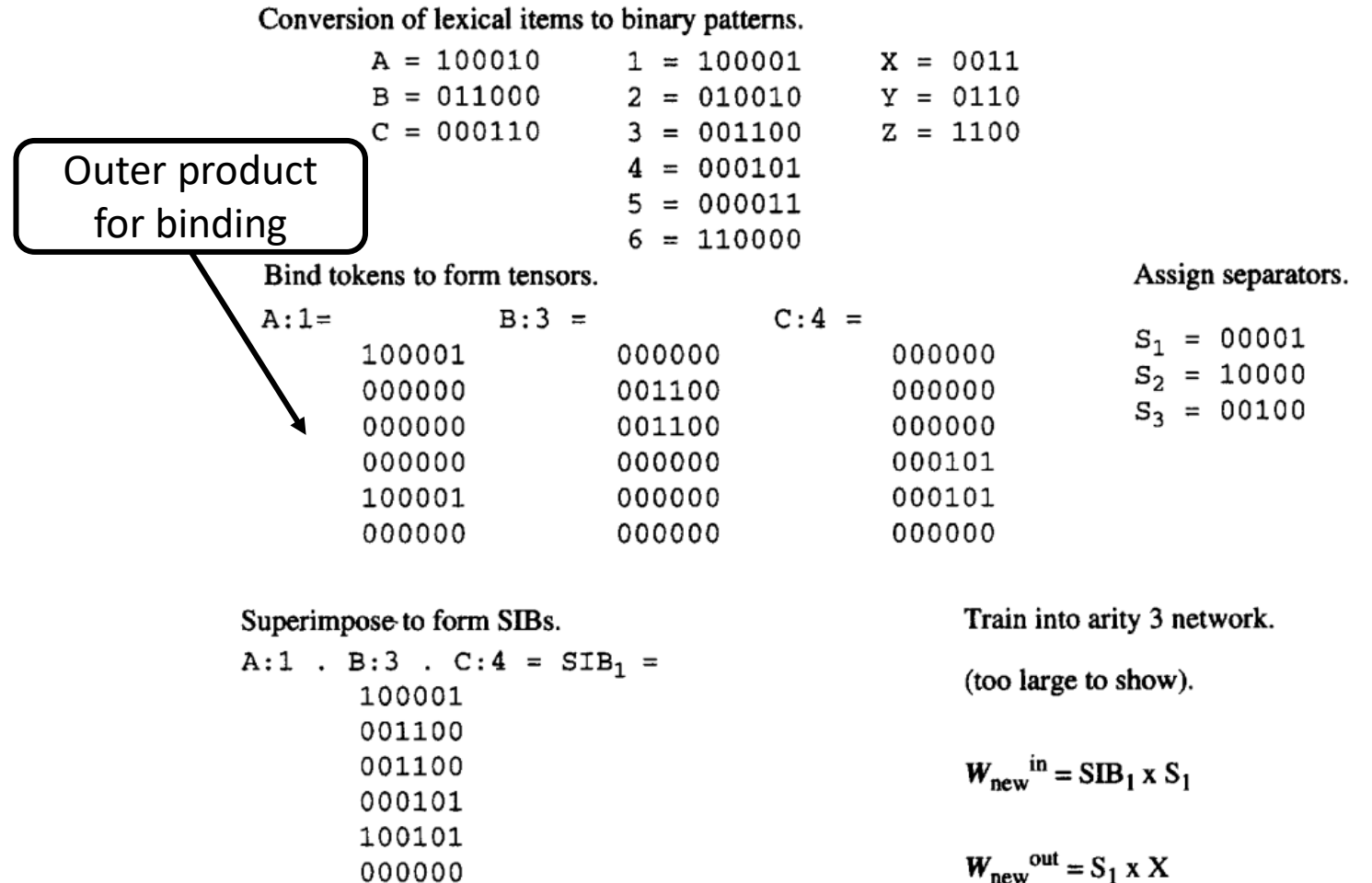
Hopfield network

Retrieve $x_i(t+1) = \text{sign} \left(\sum_{j=1}^N w_{ij} x_j(t) \right)$

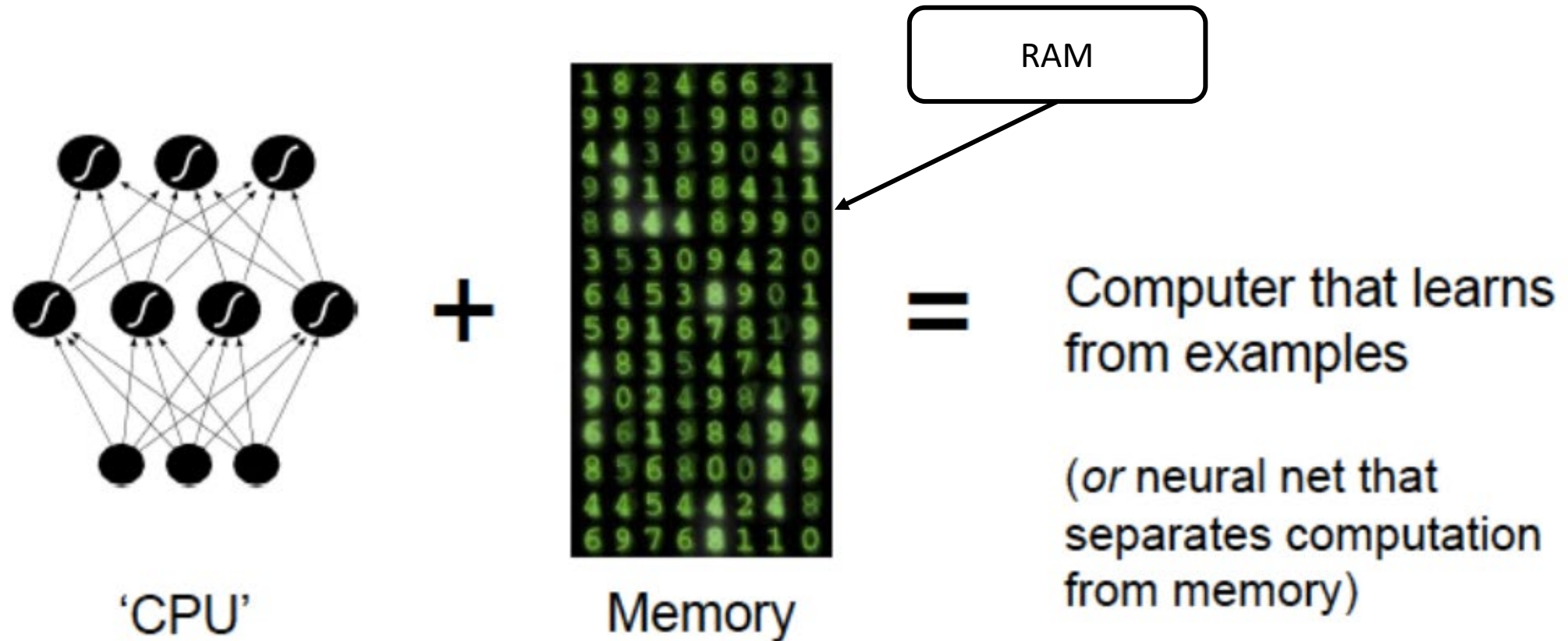
Recurrent retrieval

Rule-based reasoning with associative memory

- Encode a set of rules:
“pre-conditions
→ post-conditions”
- Support variable
binding, rule-conflict
handling and partial
rule input
- Example of encoding
rule “A:1,B:3,C:4→X”

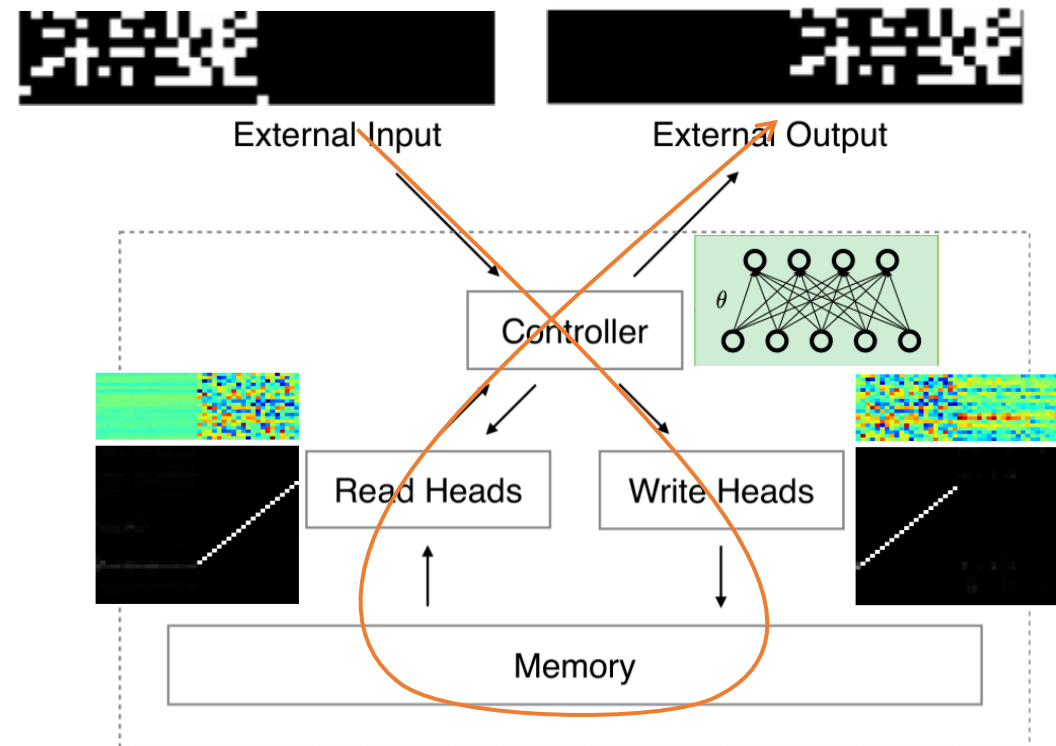


Memory-augmented neural networks: computation-storage separation

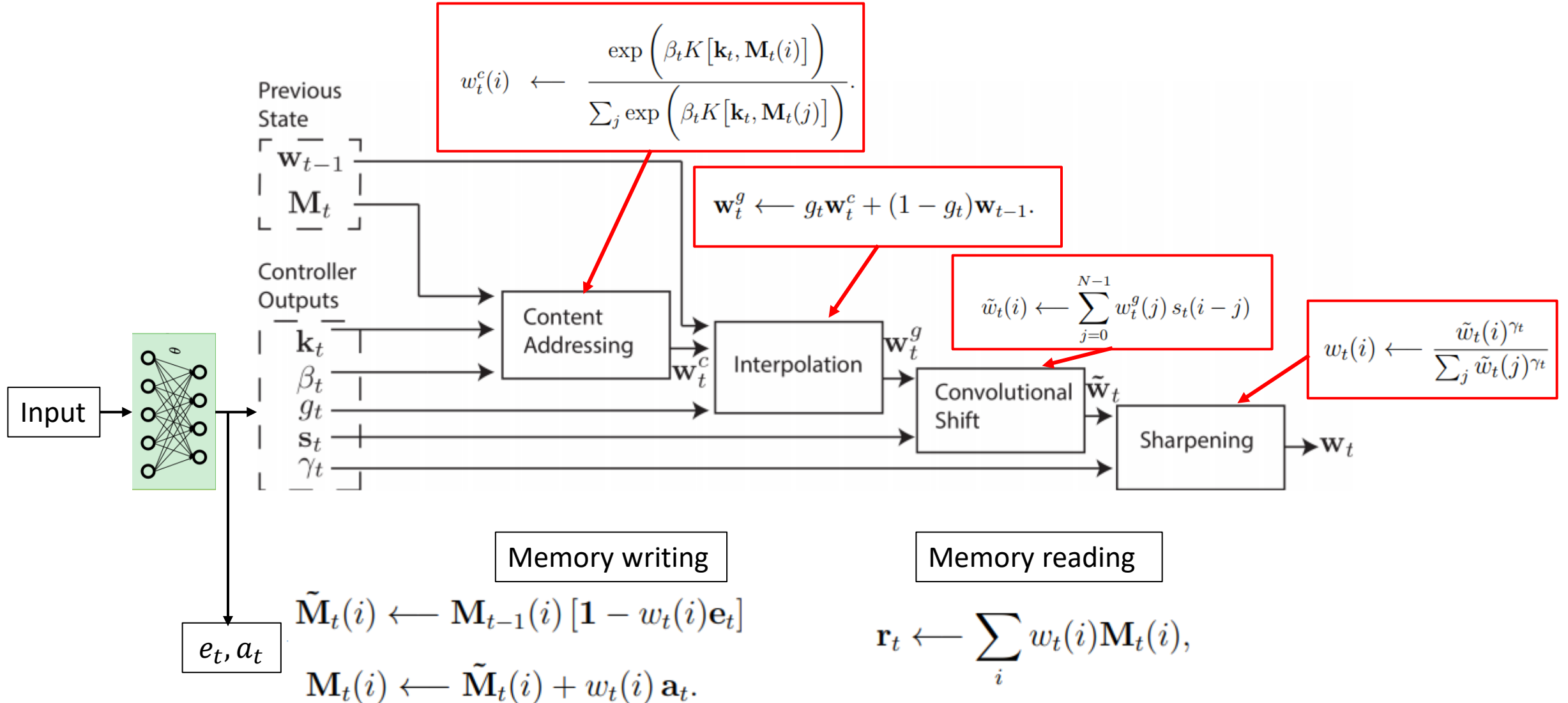


Neural Turing Machine (NTM)

- Memory is a **2d matrix**
- Controller is a **neural network**
- The controller read/writes to memory at certain addresses.
- Trained **end-to-end**, differentiable
- Simulate Turing Machine
→ support symbolic reasoning, algorithm solving

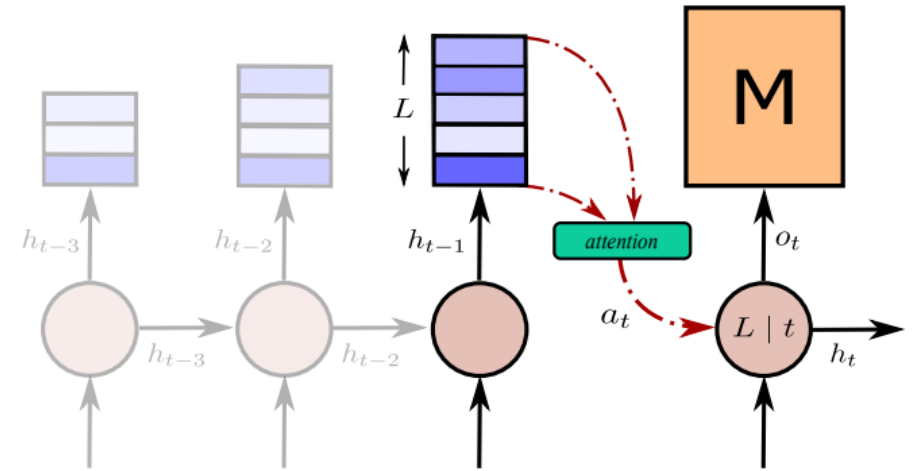


Addressing mechanism in NTM



Optimal memory writing for memorization

- Simple finding: writing too often deteriorates memory content (not retainable)
- Given input sequence of length T and only D writes, *when should we write to the memory?*



Theorem 3. Given D memory slots, a sequence with length T , a decay rate $0 < \lambda \leq 1$, then the optimal intervals $\{l_i \in \mathbb{R}^+\}_{i=1}^{D+1}$ satisfying $T = \sum_{i=1}^{D+1} l_i$ such that the lower bound on the average contribution $I_\lambda = \frac{C}{T} \sum_{i=1}^{D+1} f_\lambda(l_i)$ is maximized are the following:

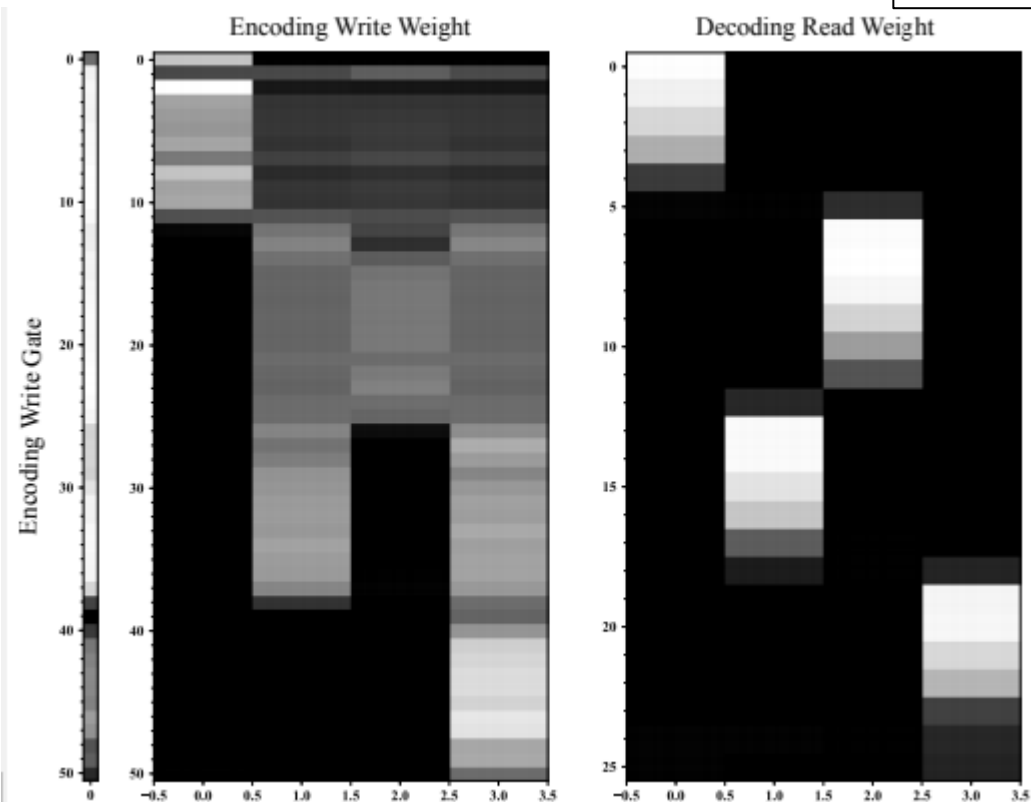
$$l_1 = l_2 = \dots = l_{D+1} = \frac{T}{D+1} \quad (7)$$

Uniform writing is optimal for memorization

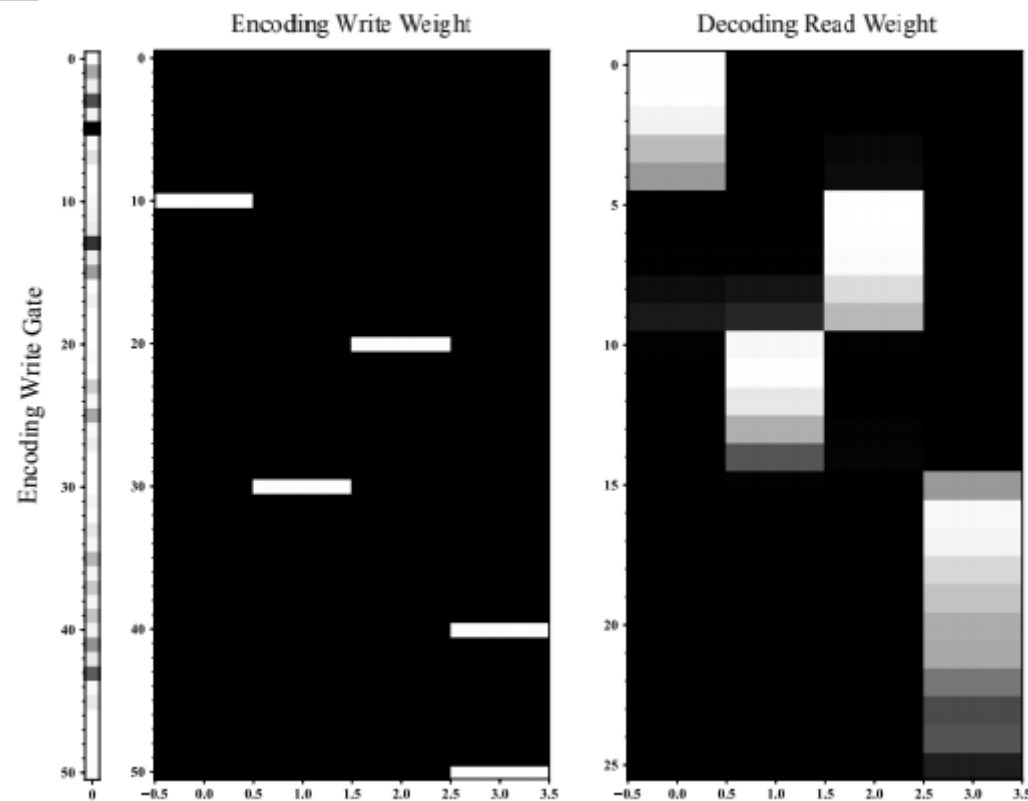
Better memorization means better algorithmic reasoning

Max	$x_1 x_2 \dots x_T$	$\max(x_1, x_2) \max(x_3, x_4) \dots \max(x_{T-1}, x_T)$
-----	---------------------	--

T=50, D=5



Regular



Uniform (cached)

Memory of independent entities

Weston, Jason, Bordes, Antoine, Chopra, Sumit, and Mikolov, Tomas.
Towards ai-complete question answering: A set of prerequisite toy tasks. CoRR, abs/1502.05698, 2015.

- Each slot store one or some entities
 - Memory writing is done separately for each memory slot
- each slot maintains the life of one or more entities
- The memory is a set of N parallel RNNs

Task 3: Three Supporting Facts

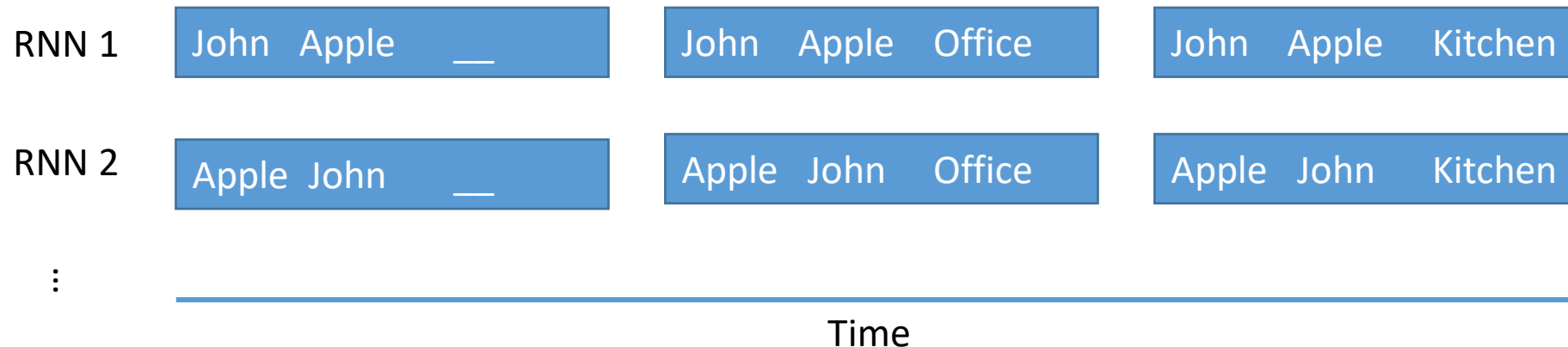
John picked up the apple.

John went to the office.

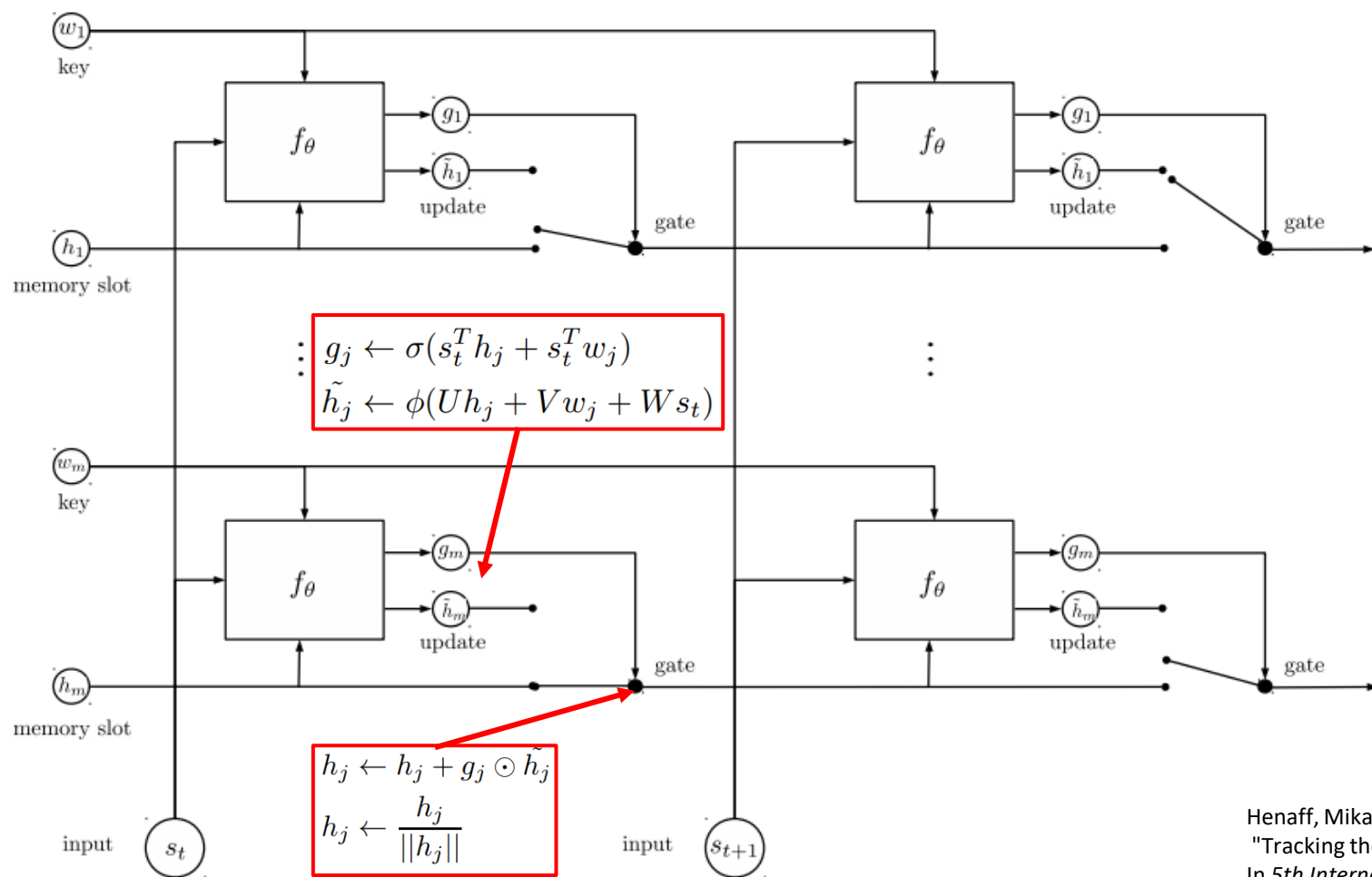
John went to the kitchen.

John dropped the apple.

Where was the apple before the kitchen? **A:office**



Recurrent entity network



“Where is the ball?”

$$p_j = \text{Softmax}(q^T h_j)$$

$$u = \sum_j p_j h_j$$

$$y = R\phi(q + Hu)$$

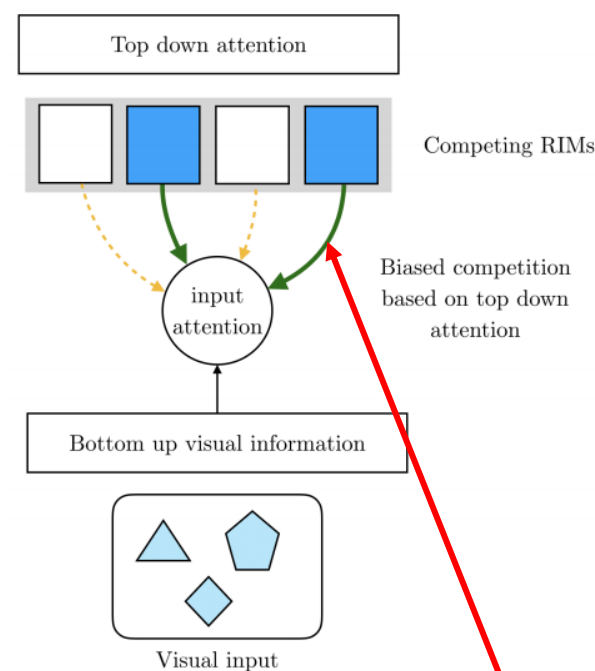
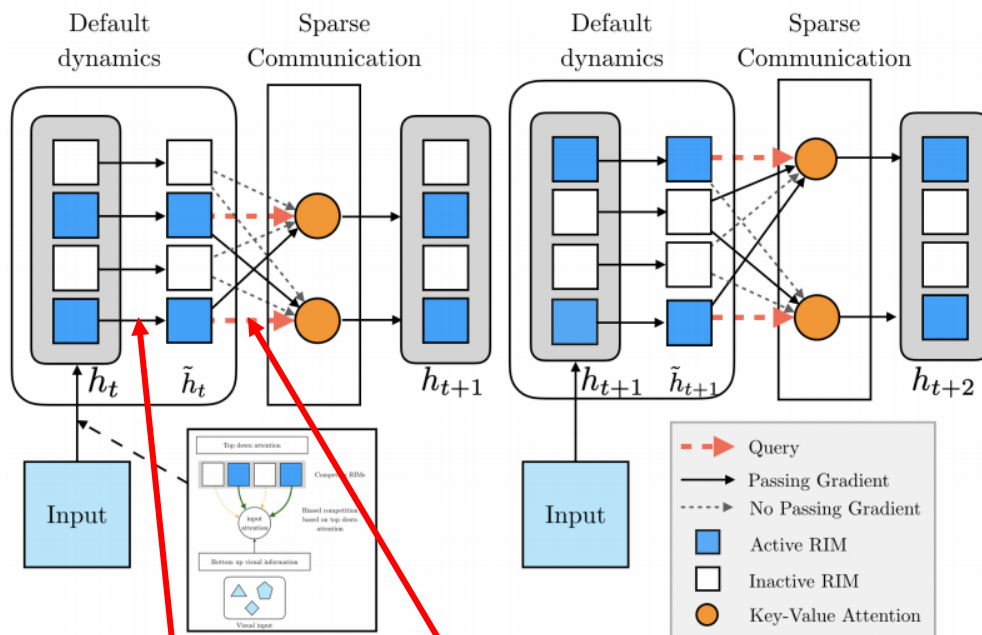
Garden

Mary picked up the ball.

Mary went to the garden.

Henaff, Mikael, Jason Weston, Arthur Szlam, Antoine Bordes, and Yann LeCun. "Tracking the world state with recurrent entity networks." In *5th International Conference on Learning Representations, ICLR 2017*. 2017.

Recurrent Independent Mechanisms



$$\tilde{h}_{t,k} = D_k(h_{t,k}) = LSTM(h_{t,k}, A_k^{(in)}; \theta_k^{(D)}) \quad \forall k \in \mathcal{S}_t$$

$$Q_{t,k} = \tilde{W}_k^q \tilde{h}_{t,k}, \forall k \in \mathcal{S}_t \quad K_{t,k} = \tilde{W}_k^e \tilde{h}_{t,k}, \forall k \quad V_{t,k} = \tilde{W}_k^v \tilde{h}_{t,k}, \forall k$$

$$h_{t+1,k} = \text{softmax} \left(\frac{Q_{t,k} (K_{t,:})^T}{\sqrt{d_e}} \right) V_{t,:} + \tilde{h}_{t,k} \quad \forall k \in \mathcal{S}_t.$$

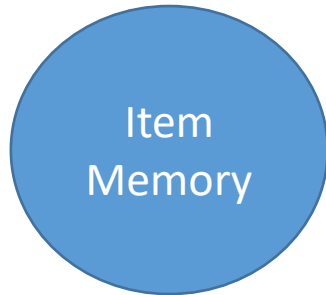
$$A_k^{(in)} = \text{softmax} \left(\frac{h_t W_k^q (X W^e)^T}{\sqrt{d_e}} \right) X W^v,$$

Relational memory

Graph memory

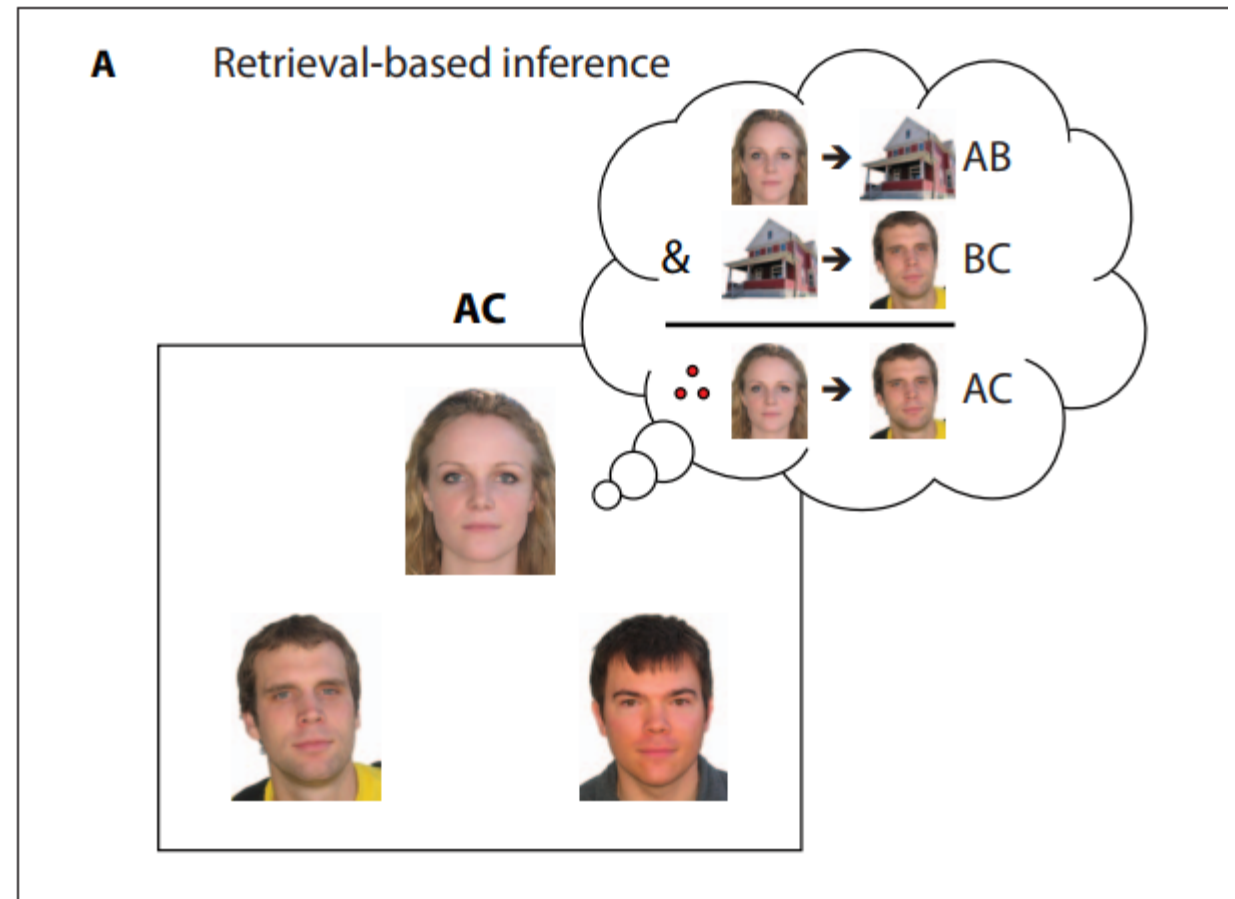
Tensor memory

Motivation for relational memory: item memory is weak at recognizing relationships

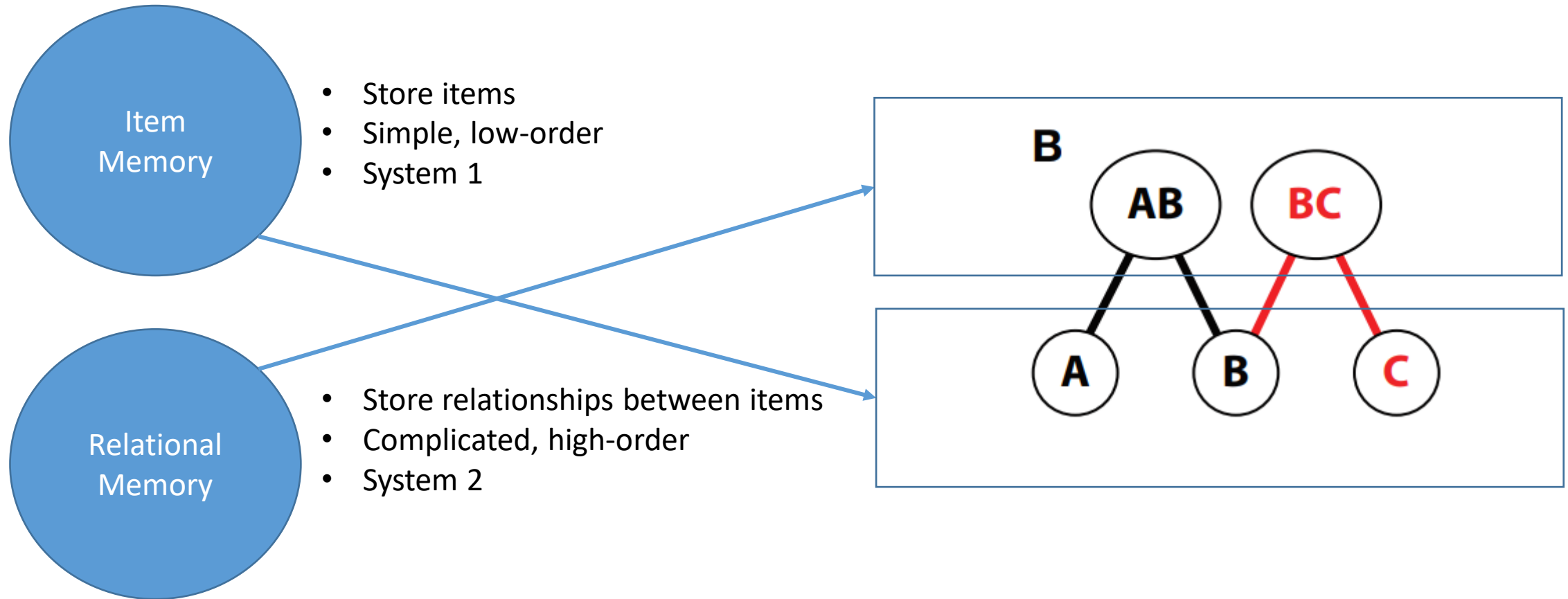


- Store and retrieve individual items
- Relate pair of items of the same time step
- Fail to **relate temporally distant items**

$$\hat{\mathbf{M}} = \sum_{k=1}^q \mathbf{b}_k \mathbf{a}_k^T$$



Dual process in memory

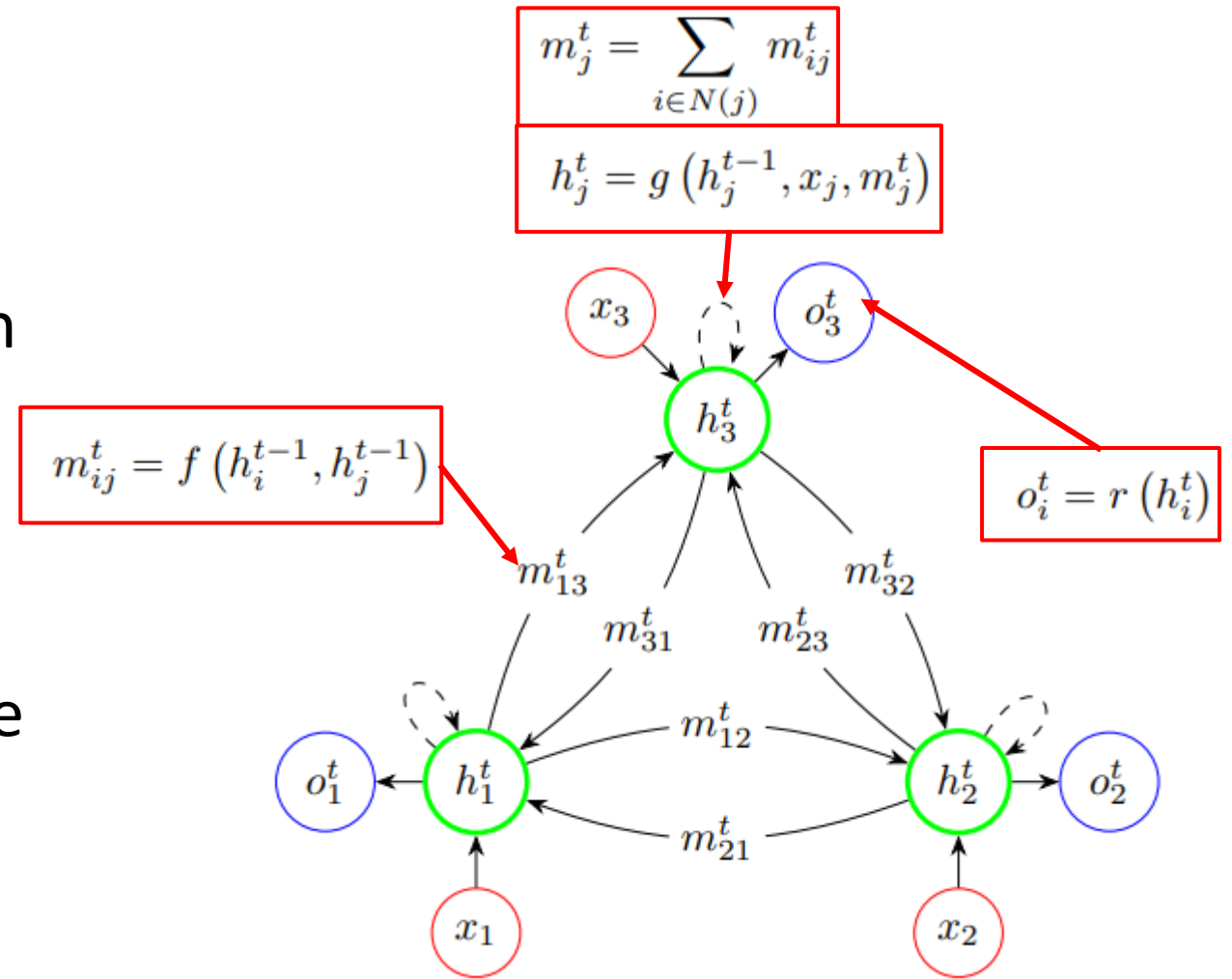


Howard Eichenbaum, *Memory, amnesia, and the hippocampal system* (MIT press, 1993).

Alex Konkelt and Neal J Cohen, "Relational memory and the hippocampus: representations and methods", *Frontiers in neuroscience* 3 (2009).

Memory as graph

- Memory is a **static graph** with fixed nodes and edges
- Relationship is somehow known
- Each memory node stores the state of the graph's node
- Write to node via message passing
- Read from node via MLP



bAbI

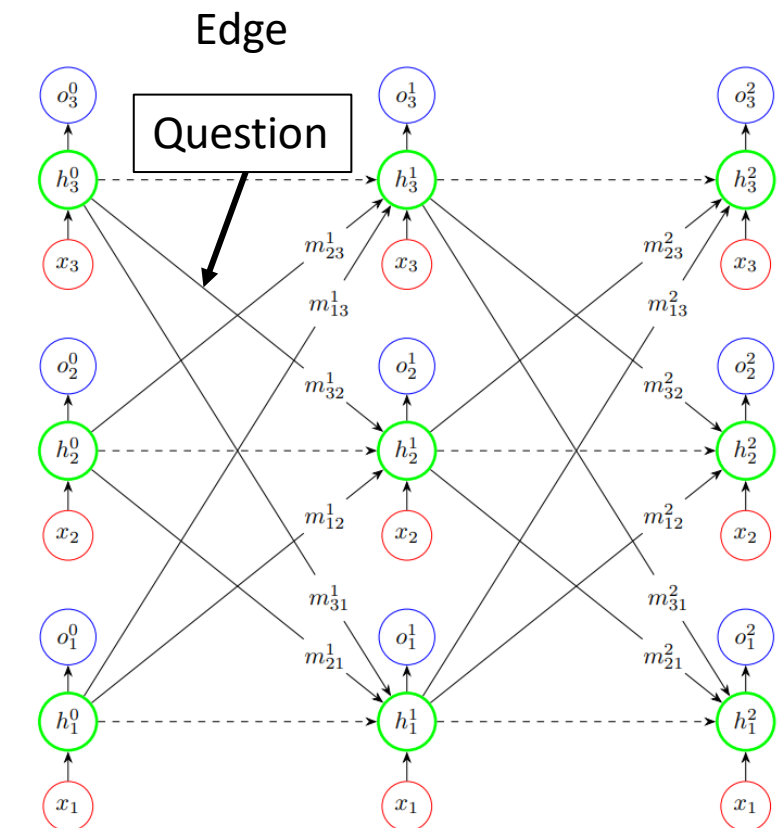
$$x_i = \text{MLP}(\text{concat}(\text{last}(\text{LSTM}_S(s_i)), \text{last}(\text{LSTM}_Q(q)), \text{onehot}(p_i + o)))$$

Node

Fact 1

Fact 2

Fact 3



$$o^t = r(\sum_i h_i^t)$$

Answer

$$m_{ij}^t = f(h_i^{t-1}, h_j^{t-1}, e_{ij})$$

Method	N	Mean Error (%)	Failed tasks (err. >5%)
RRN* (this work)	15	0.46 ± 0.77	0.13 ± 0.35

CLEVER

$$o_i = \text{concat}(p_i, \text{onehot}(c_i), \text{onehot}(m_i))$$

$$q = \text{concat}(\text{onehot}(s), \text{onehot}(n))$$

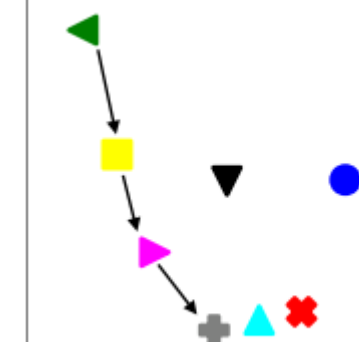
$$x_i = \text{MLP}(\text{concat}(o_i, q))$$

Node

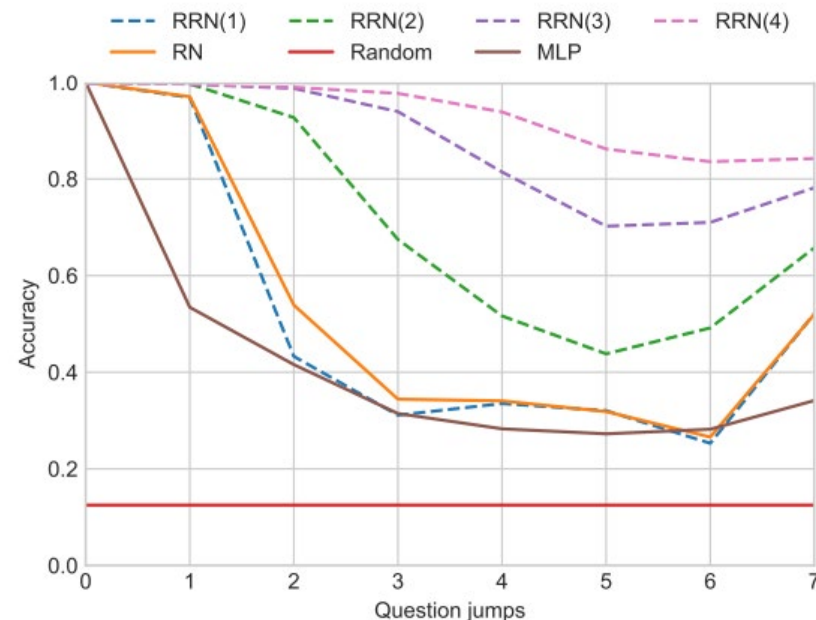
(colour, shape, position)

Edge

(distance)

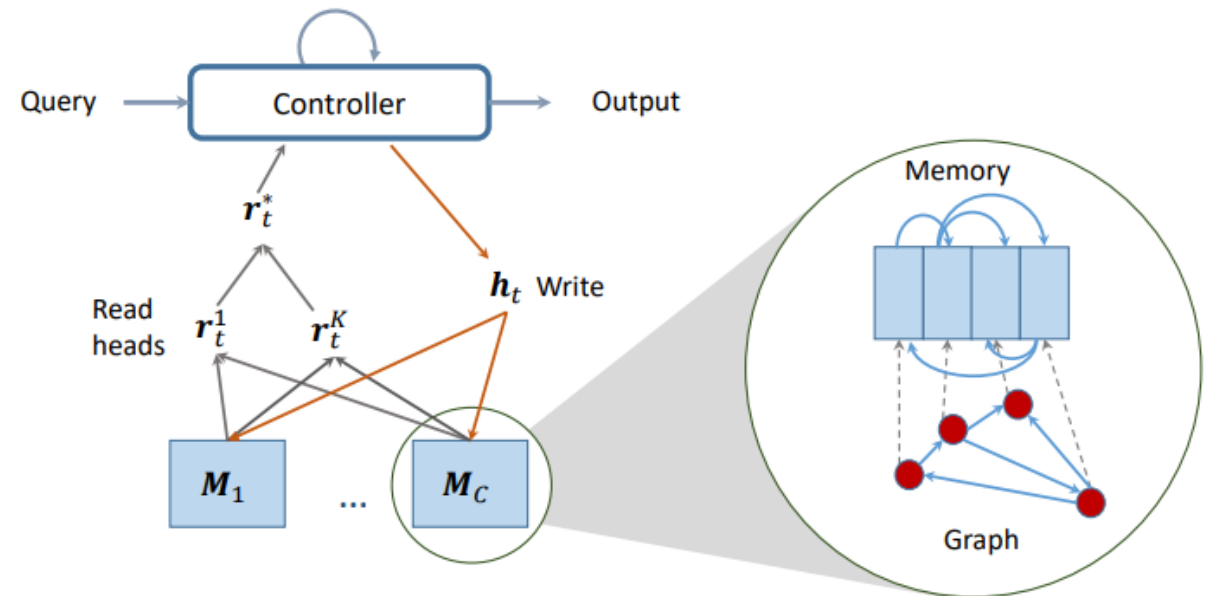


solution to the question: “green, 3 jumps”, which is “plus”,



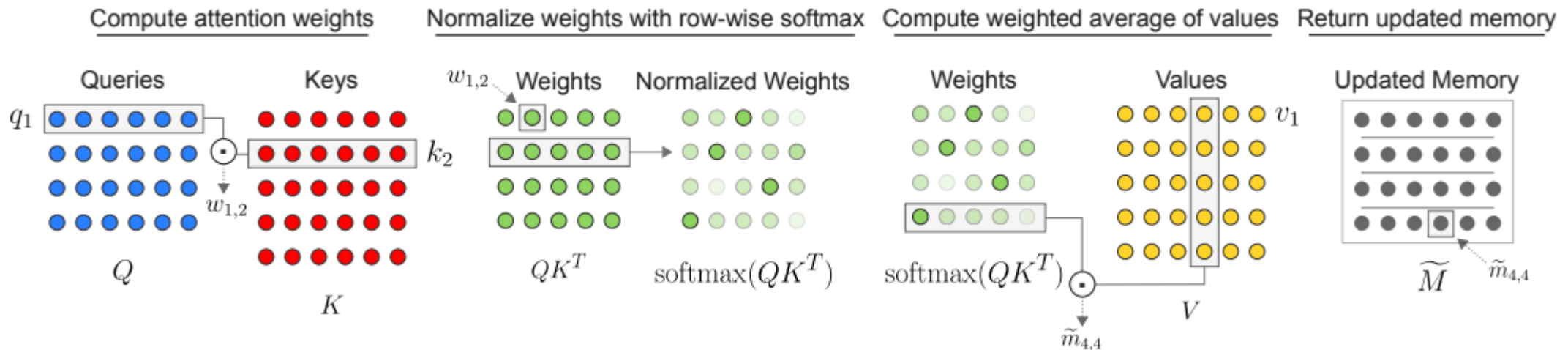
Memory of graphs access conditioned on query

- Encode multiple graphs, **each graph is stored in a set of memory row**
- For each graph, the controller read/write to the memory:
 - Read uses content-based attention
 - Write use message passing
- Aggregate read vectors from all graphs to create output

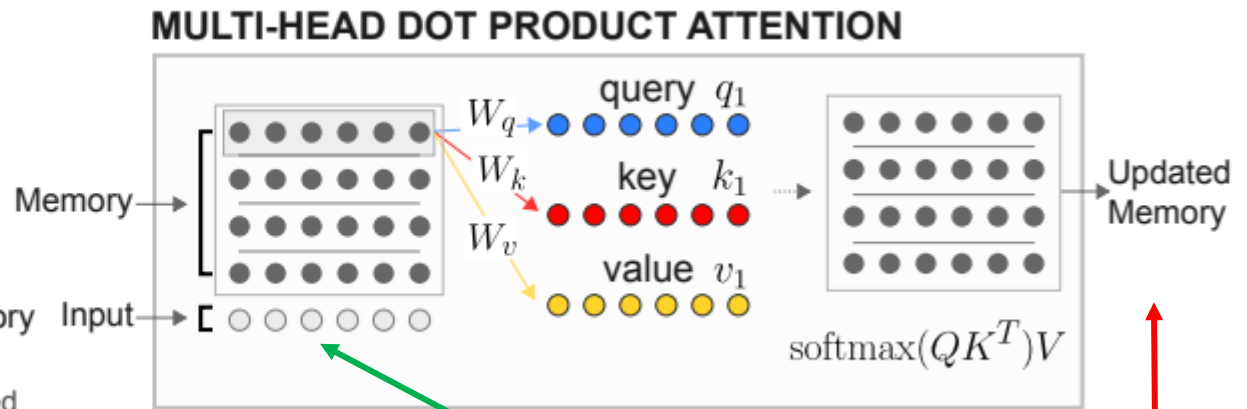
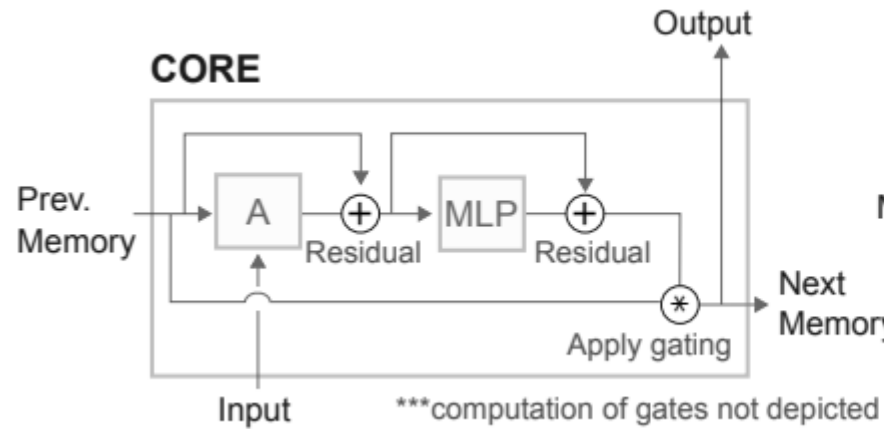


Capturing relationship can be done via memory slot interactions using attention

- Graph memory needs customization to an explicit design of nodes and edges
- Can we automatically learn structure with a 2d tensor memory?
- Capture relationship: each slot interacts with all other slots (self-attention)



Relational Memory Core (RMC) operation



$$s_{i,t} = (h_{i,t-1}, m_{i,t-1})$$

$$f_{i,t} = W^f x_t + U^f h_{i,t-1} + b^f$$

$$i_{i,t} = W^i x_t + U^i h_{i,t-1} + b^i$$

$$o_{i,t} = W^o x_t + U^o h_{i,t-1} + b^o$$

$$m_{i,t} = \sigma(f_{i,t} + \tilde{b}^f) \circ m_{i,t-1} + \sigma(i_{i,t}) \circ \underbrace{g_\psi(\tilde{m}_{i,t})}_{\text{row/memory-wise MLP with layer normalisation}}$$

$$h_{i,t} = \sigma(o_{i,t}) \circ \tanh(m_{i,t})$$

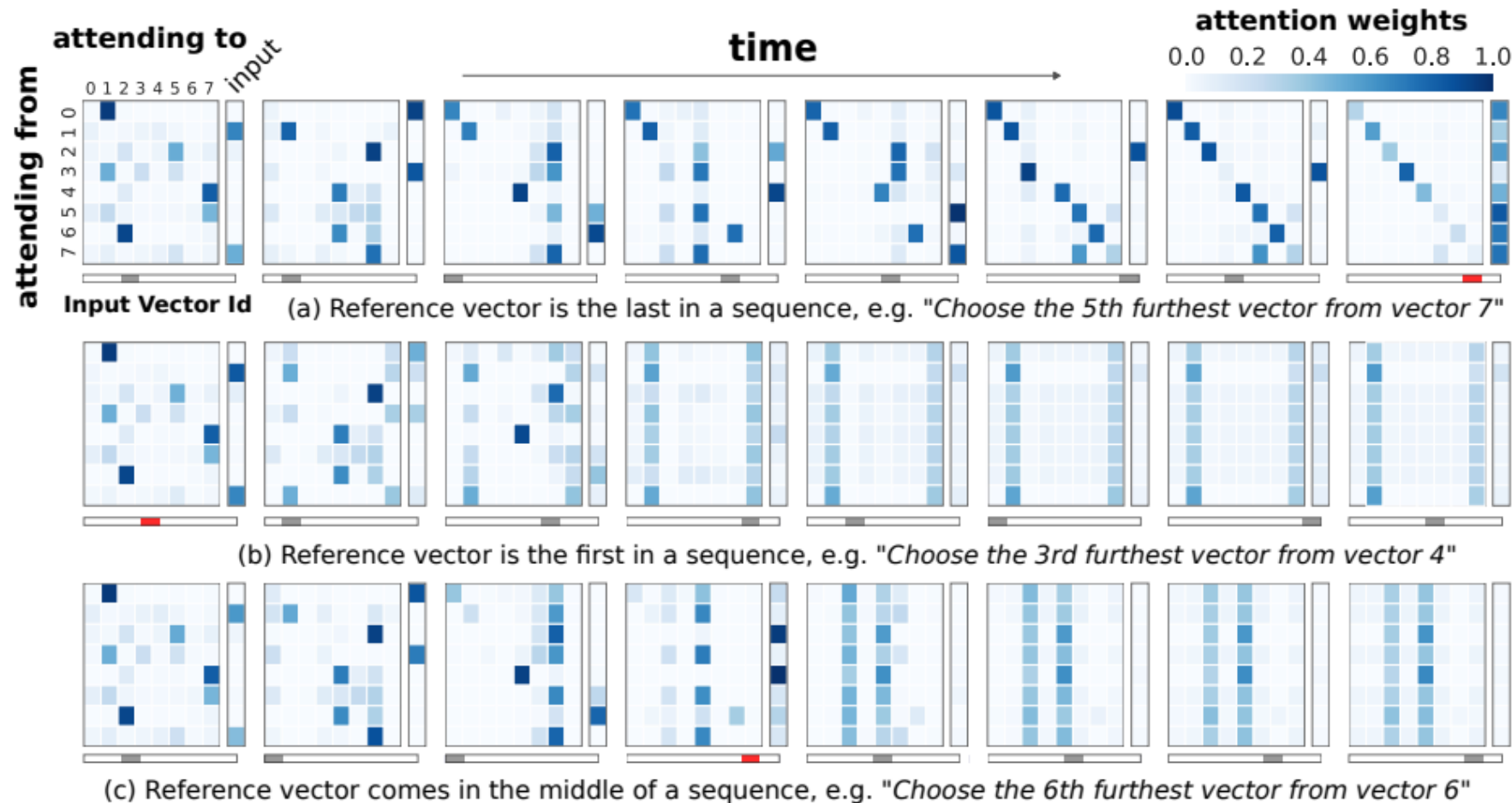
row/memory-wise MLP with layer normalisation

$$s_{i,t+1} = (m_{i,t}, h_{i,t})$$

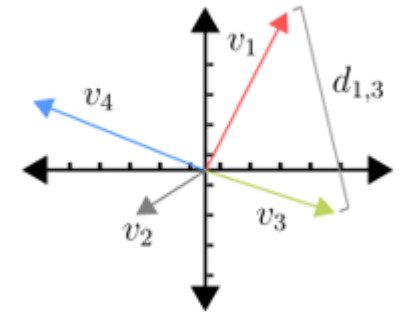
$$\tilde{M} = \text{softmax} \left(\frac{MW^q ([M; x] W^k)^T}{\sqrt{d^k}} \right) [M; x] W^v,$$

RNN-like
Interface

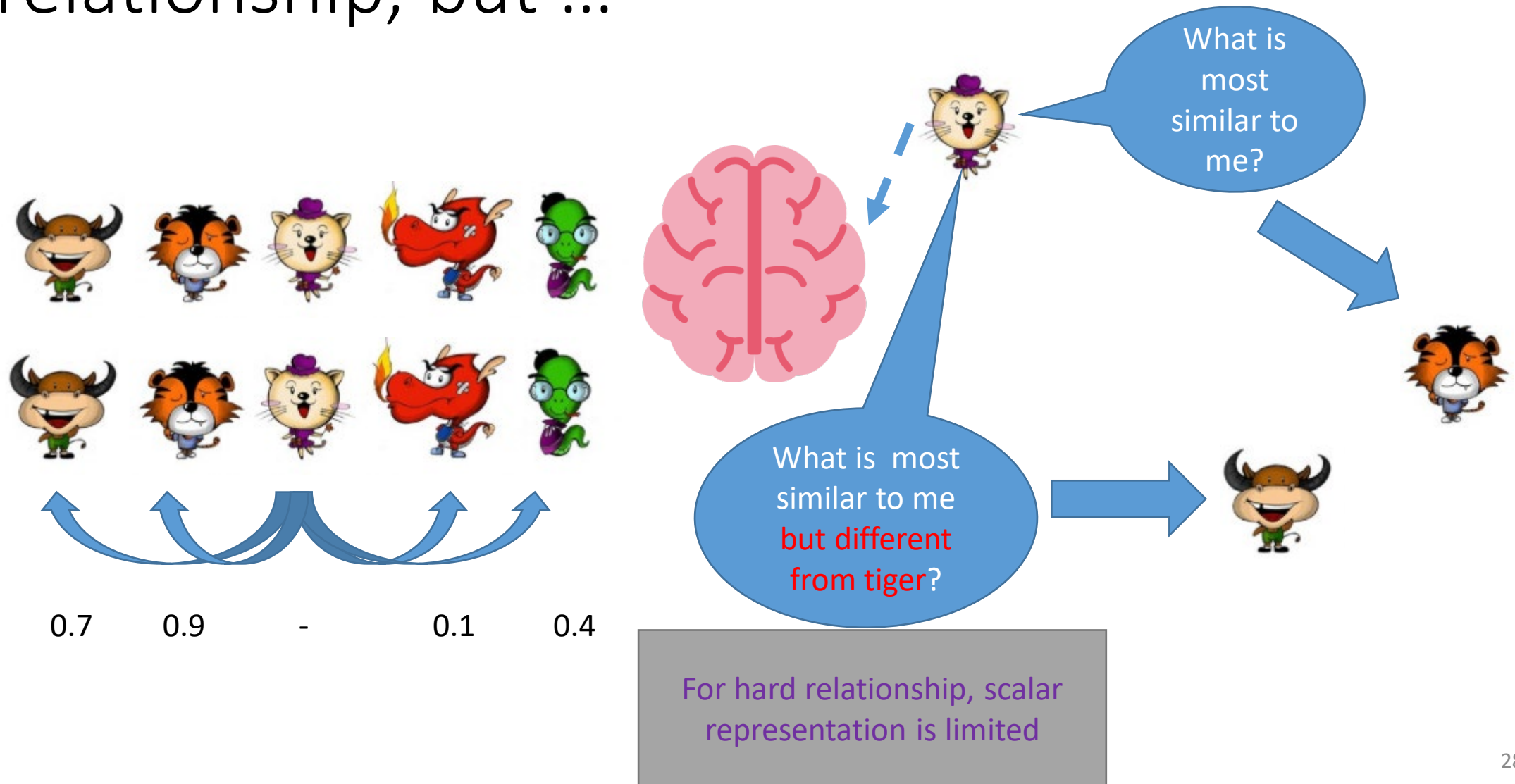
Allowing pair-wise interactions can answer questions on temporal relationship



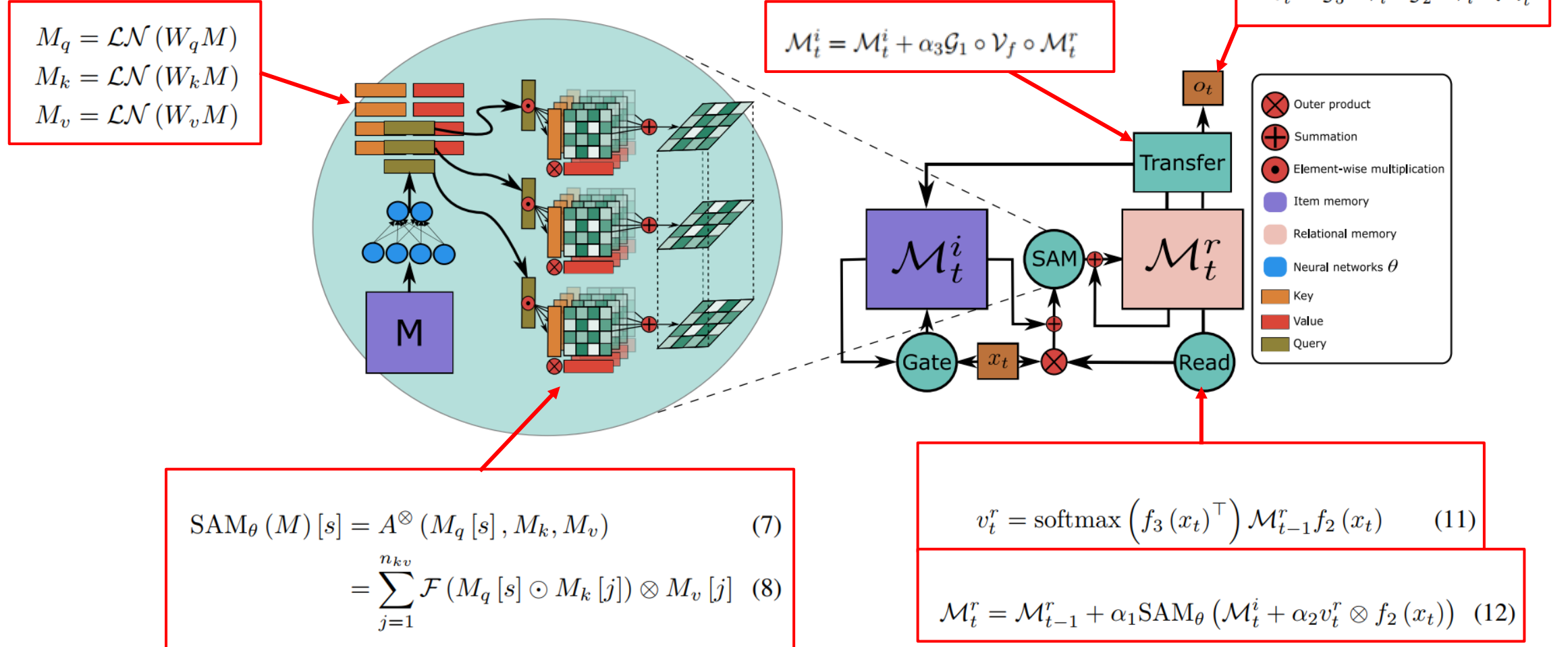
What is the N^{th} furthest from vector m ?



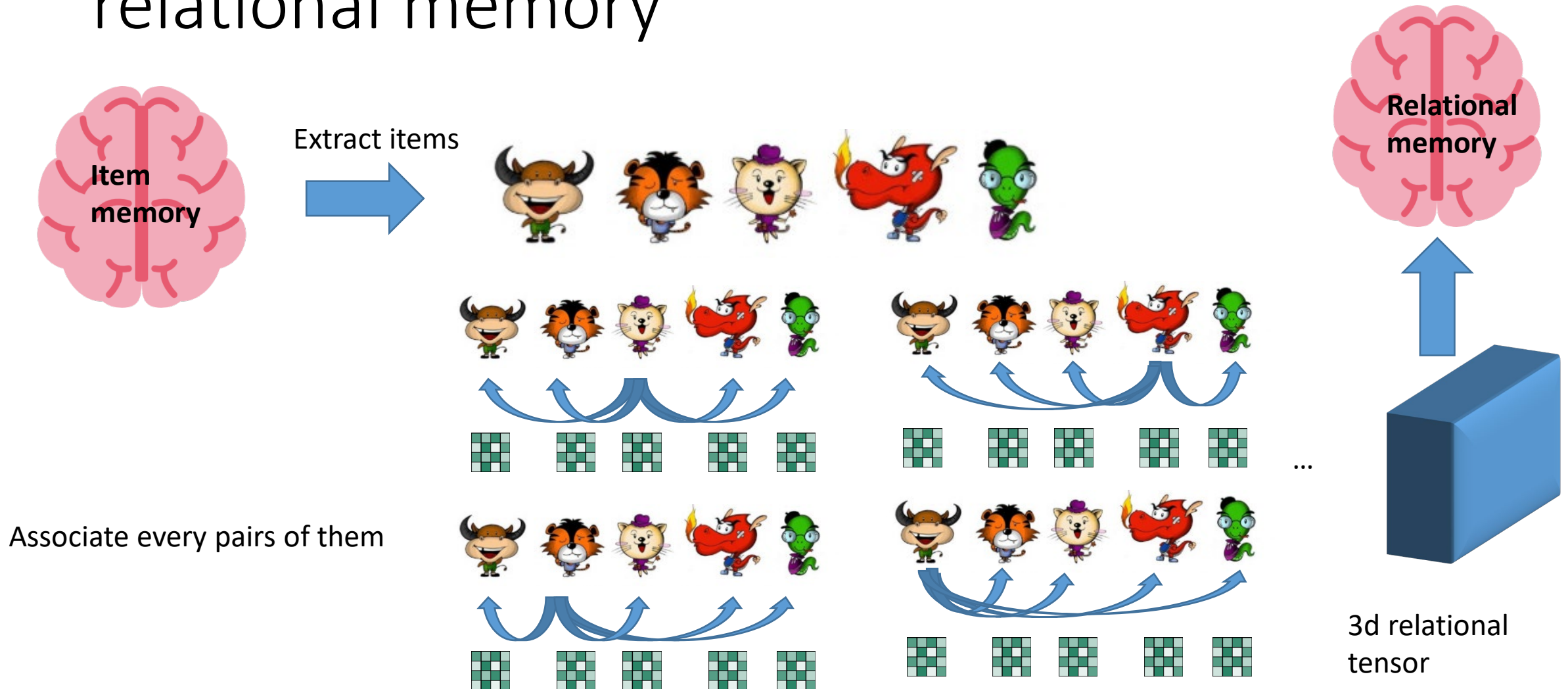
Dot product attention works for simple relationship, but ...



Self-attentive associative memory



Complicated relationship needs high-order relational memory



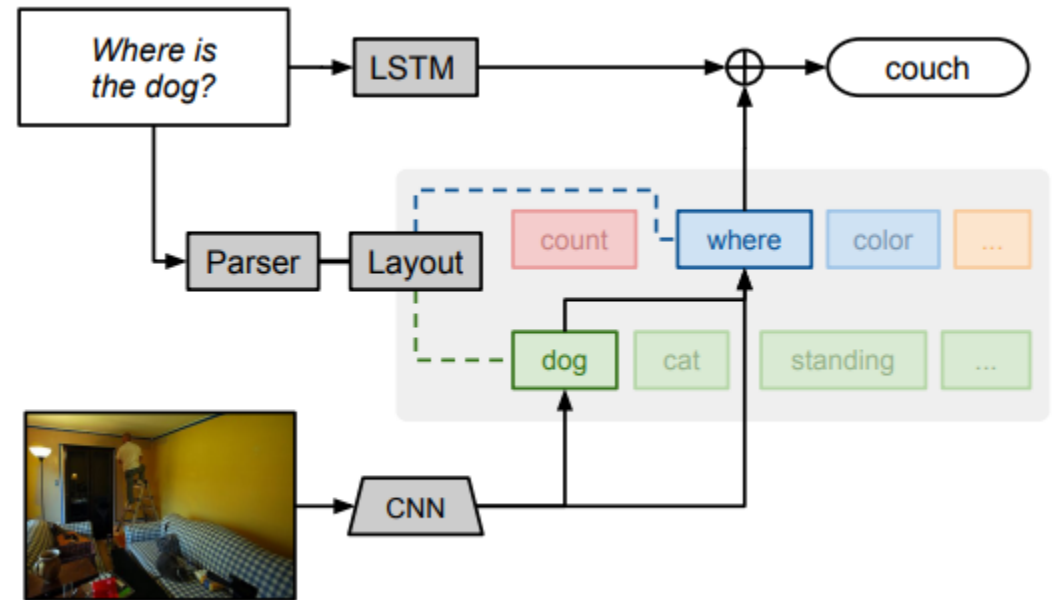
Program memory

Module memory

Stored-program memory

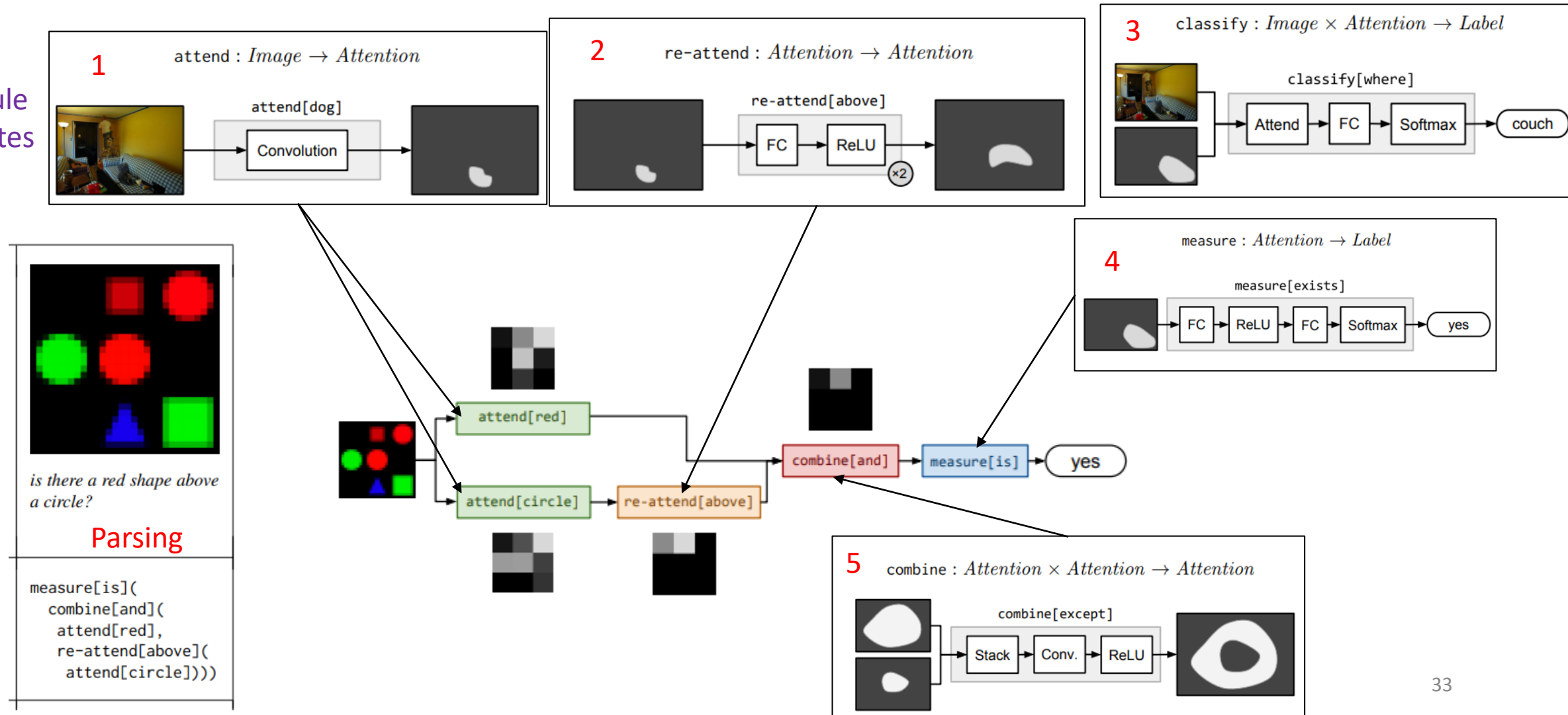
Predefining program for subtask

- A program designed for a task becomes a **module**
- Parse a question to module layout (order of program execution)
- Learn the weight of each module to master the task



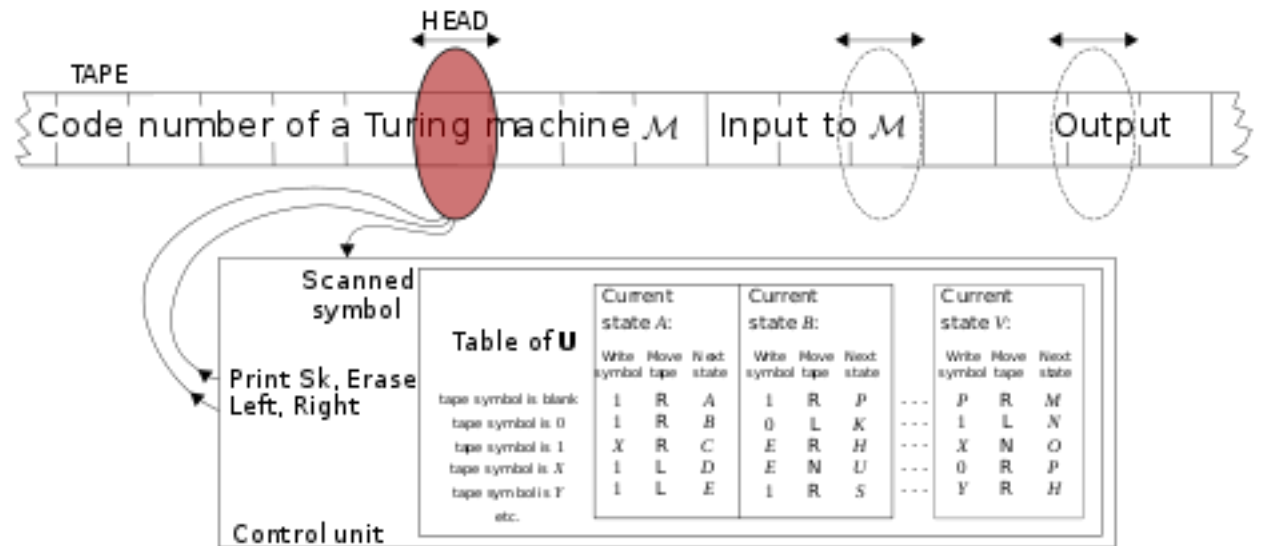
Program selection is based on parser, others are end2end trained

5 module templates



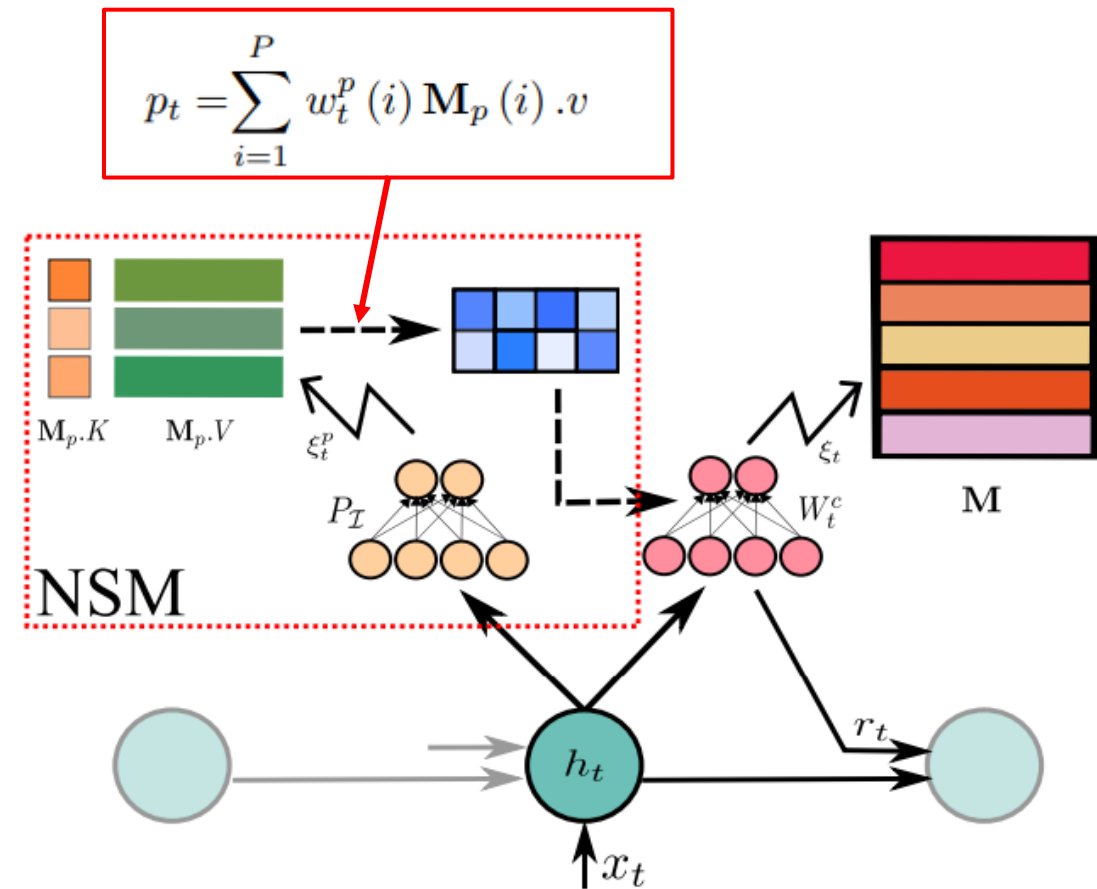
The most powerful memory is one that stores both program and data

- Computer architecture: Universal Turing Machines/Harvard/VNM
- **Stored-program principle**
- Break a big task into subtasks, each can be handled by a TM/single purposed program stored in a program memory



NUTM: Learn to select program (neural weight) via program attention

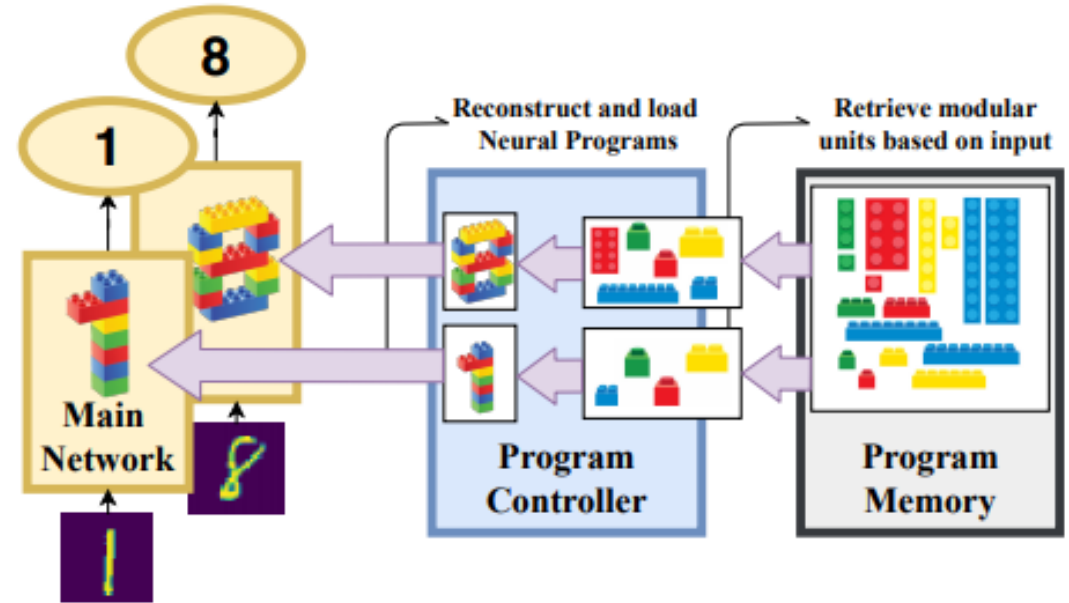
- Neural stored-program memory (NSM) stores **key (the address)** and **values (the weight)**
- The weight is selected and loaded to the controller of NTM
- The stored NTM weights and the weight of the NUTM is learnt end-to-end by backpropagation



Le, Hung, Truyen Tran, and Svetha Venkatesh. "Neural Stored-program Memory." In *International Conference on Learning Representations*. 2019.

Scaling with memory of mini-programs

- Prior, 1 program = 1 neural network (millions of parameters)
- Parameter inefficiency since the programs do not share common parameters
- Solution: store sharable mini-programs to compose infinite number of programs

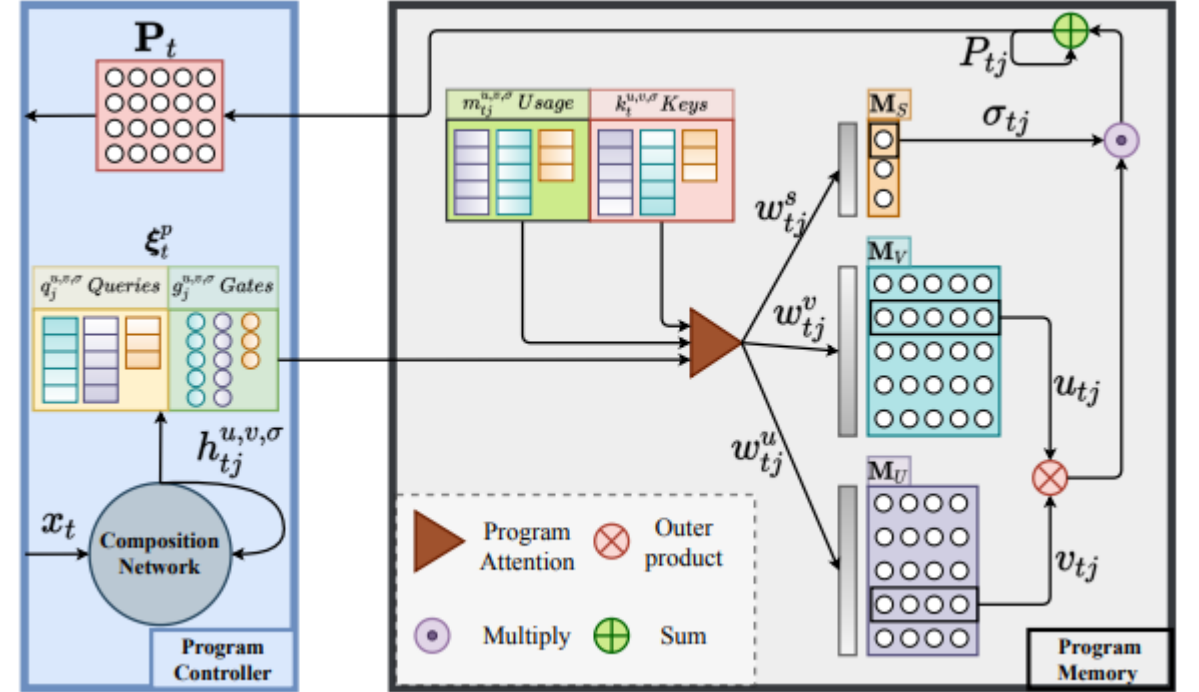


it is analogous to building Lego structures corresponding to inputs from basic Lego bricks.

Recurrent program attention to retrieve singular components of a program

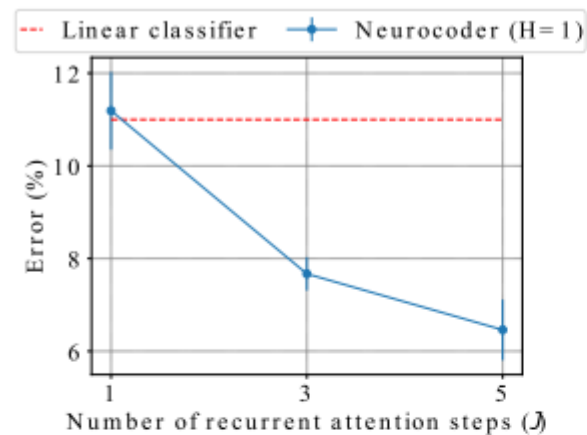
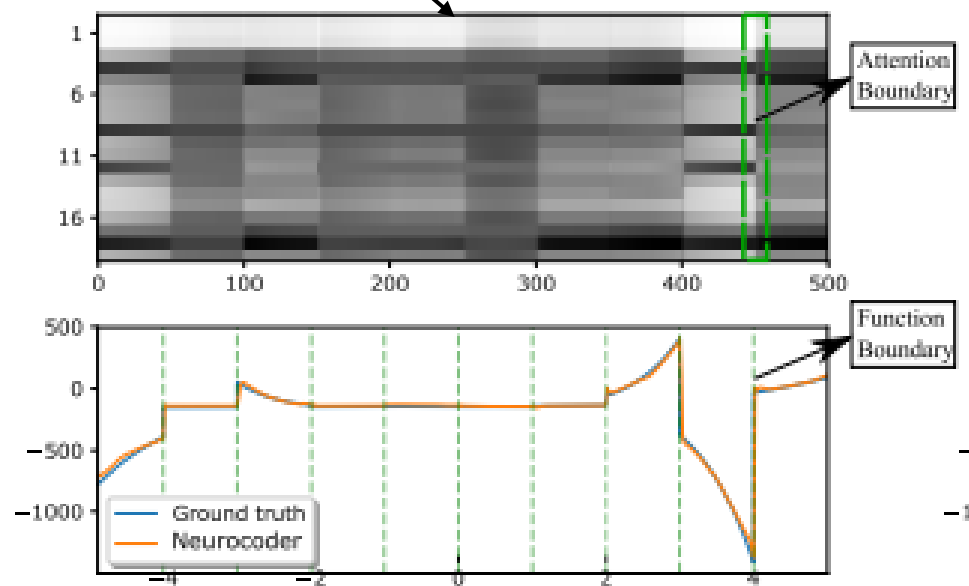
$$\begin{aligned} \mathbf{P}_t &= \mathbf{U}\mathbf{S}\mathbf{V}^T \\ &= \sum_n \sigma_{tn} u_{tn} v_{tn}^T \\ u_{tn} &= \sum_{i=1}^{P_u} w_{tin}^u \mathbf{M}_U(i) \\ v_{tn} &= \sum_{i=1}^{P_v} w_{tin}^v \mathbf{M}_V(i) \end{aligned}$$

$$\sigma_{tn} = \begin{cases} \text{softplus} \left(\sum_{i=1}^{P_s} w_{tin}^\sigma \mathbf{M}_S(i) \right) & n = r_m \\ \sigma_{tn+1} + \text{softplus} \left(\sum_{i=1}^{P_s} w_{tin}^\sigma \mathbf{M}_S(i) \right) & n < r_m \end{cases}$$

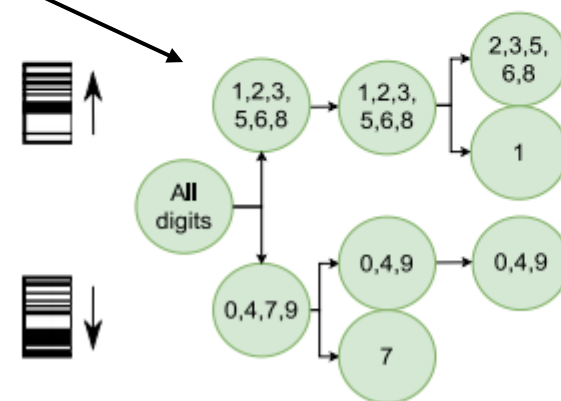


Program attention is equivalent to binary decision tree reasoning

Recurrent program attention auto detects task boundary

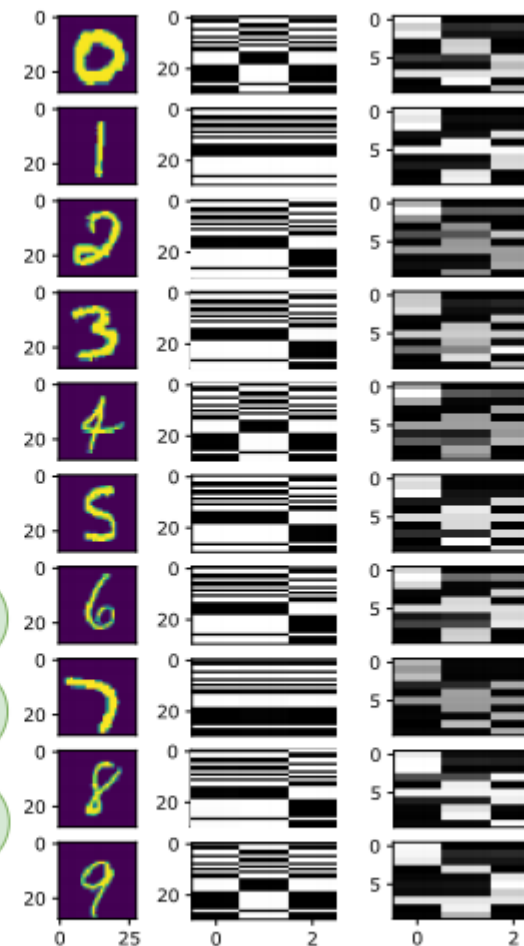


(a)



(c)

(d)



(b)

QA

10. Combinatorics reasoning

RNN

MANN

GNN

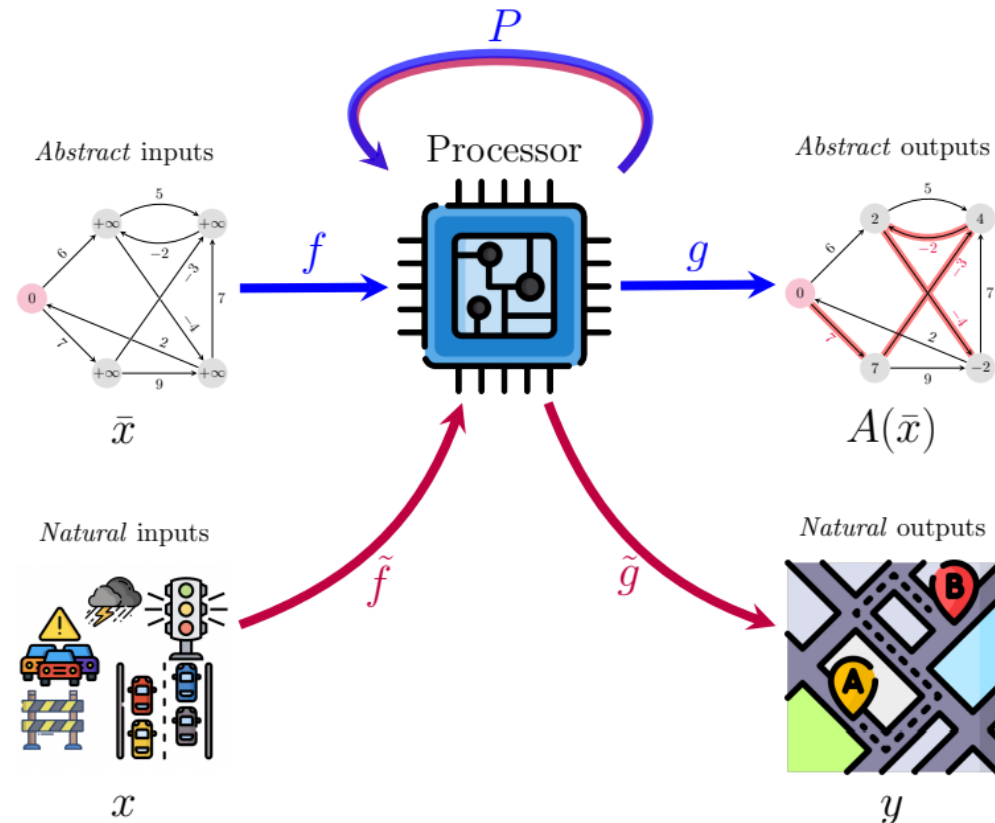
Transformer

Implement combinatorial algorithms with neural networks

Train neural processor P to imitate algorithm A

Processor P:

- (a) aligned with the computations of the target algorithm;
- (b) operates by matrix multiplications, hence natively admits useful gradients;
- (c) operates over high-dimensional latent spaces



Generalizable
Inflexible

Noisy
High dimensional

Processor as RNN

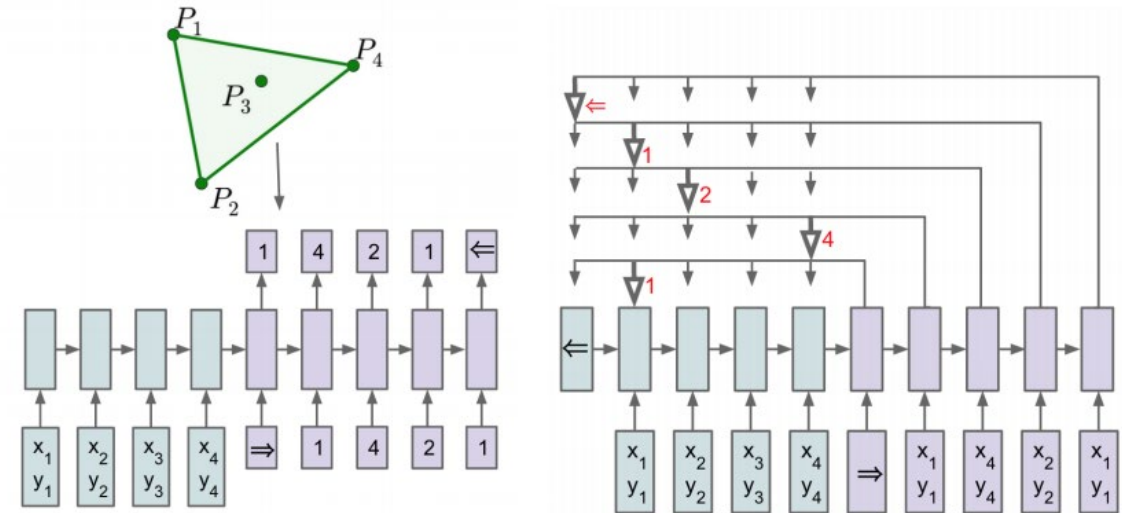
- Do not assume knowing the structure of the input, input as a sequence
→ not really reasonable, harder to generalize
- RNN is Turing-complete → can simulate any algorithm
- But, it is not easy to learn the simulation from data (input-output)

→ Pointer network

Assume $O(N)$ memory
And $O(N^2)$ computation
 N is the size of input

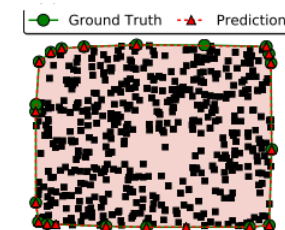
$$u_j^i = v^T \tanh(W_1 e_j + W_2 d_i) \quad j \in (1, \dots, n)$$

$$p(C_i | C_1, \dots, C_{i-1}, \mathcal{P}) = \text{softmax}(u^i)$$

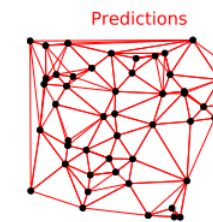


(a) Sequence-to-Sequence

(b) Ptr-Net



(d) Ptr-Net, $m=5-50, n=500$



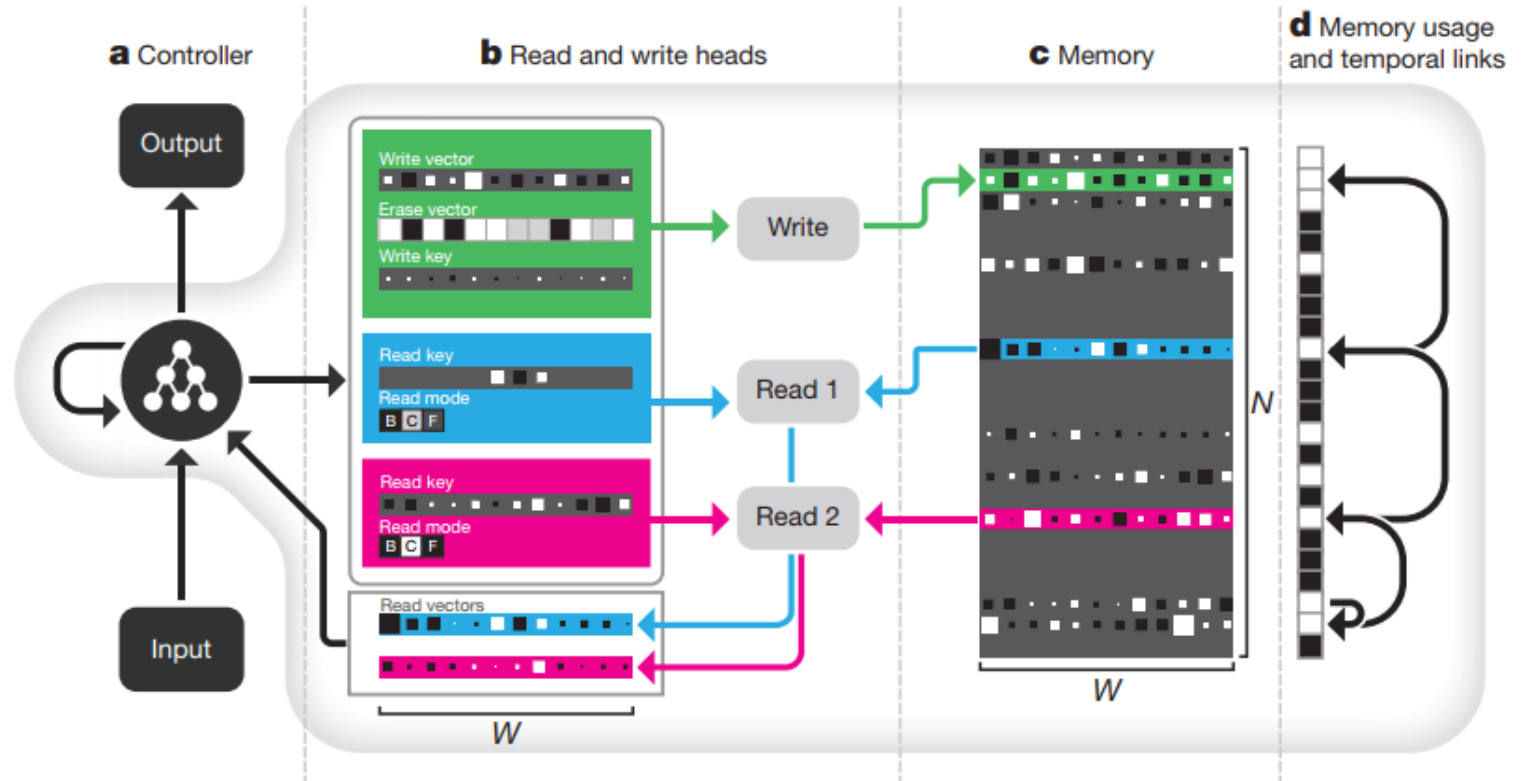
(e) Ptr-Net, $m=50, n=50$



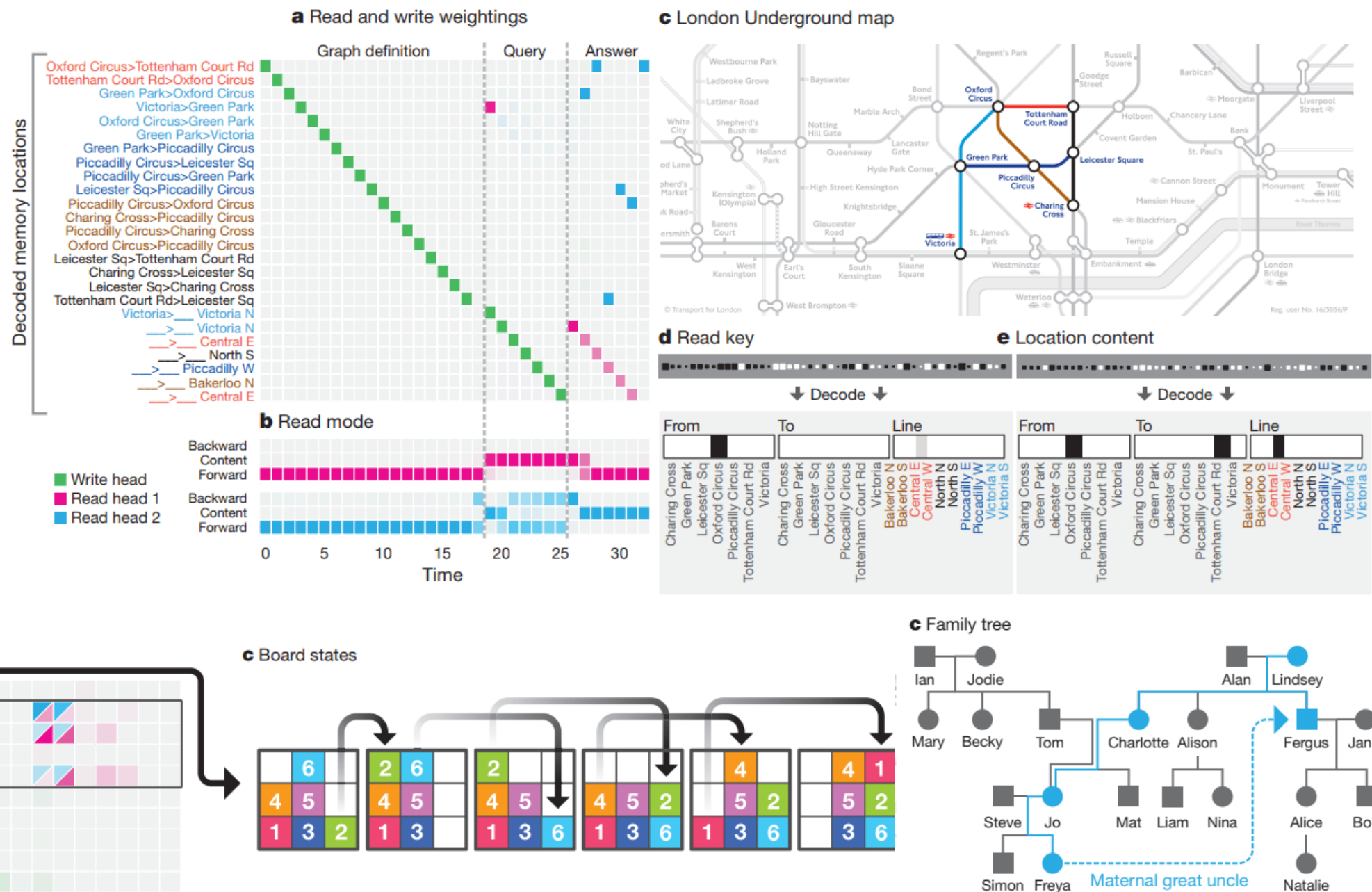
(f) Ptr-Net, $m=5-20, n=20$

Processor as MANN

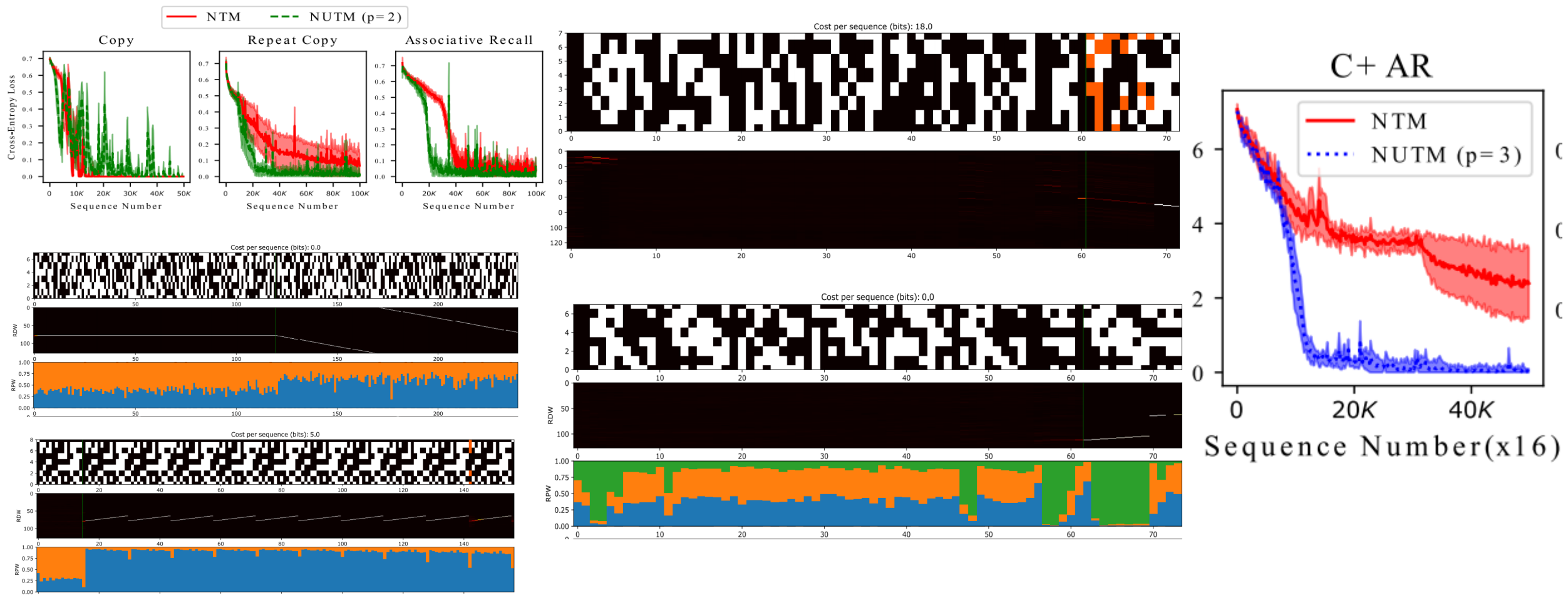
- MANN simulates neural computers or Turing machine → ideal for implement algorithms
- Sequential input, no assumption on input structure
- Assume $O(1)$ memory and $O(N)$ computation



DNC: item memory for graph reasoning



NUTM: implementing multiple algorithms at once



Le, Hung, Truyen Tran, and Svetha Venkatesh. "Neural Stored-program Memory."
In *International Conference on Learning Representations*. 2019.

STM: relational memory for graph reasoning

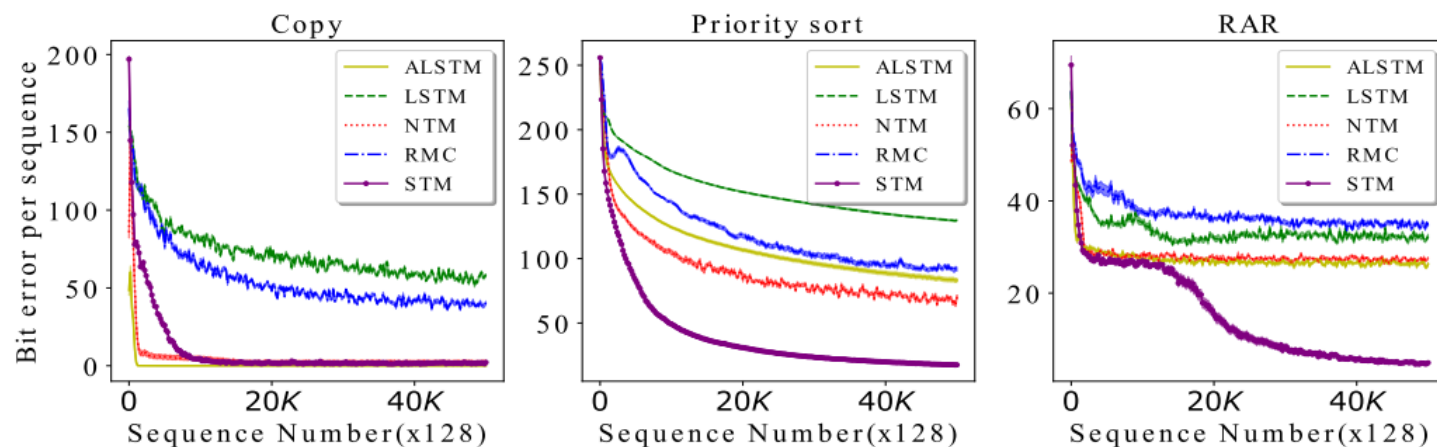
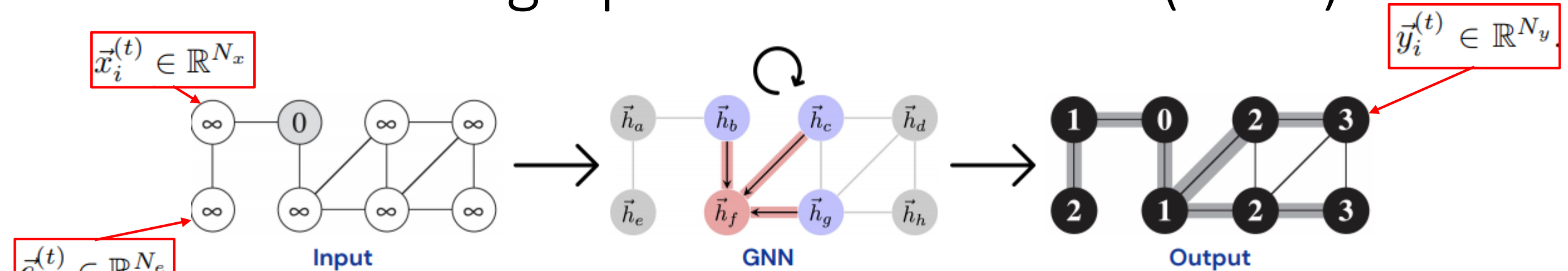


Figure 2. Learning curves on algorithmic synthetic tasks.

Model	#Parameters	Convex hull		TSP		Shortest path	Minimum spanning tree
		$N = 5$	$N = 10$	$N = 5$	$N = 10$		
LSTM	4.5 M	89.15	82.24	73.15 (2.06)	62.13 (3.19)	72.38	80.11
ALSTM	3.7 M	89.92	85.22	71.79 (2.05)	55.51 (3.21)	76.70	73.40
DNC	1.9 M	89.42	79.47	73.24 (2.05)	61.53 (3.17)	83.59	82.24
RMC	2.8 M	93.72	81.23	72.83 (2.05)	37.93 (3.79)	66.71	74.98
STM	1.9 M	96.85	91.88	73.96 (2.05)	69.43 (3.03)	93.43	94.77

Table 3. Prediction accuracy (%) for geometry and graph reasoning tasks with random *one-hot* associated features. Italic numbers are tour length–additional metric for TSP. Average optimal tour lengths found by brute-force search for $N = 5$ and 10 are 2.05 and 2.88, respectively.

Processor as graph neural network (GNN)



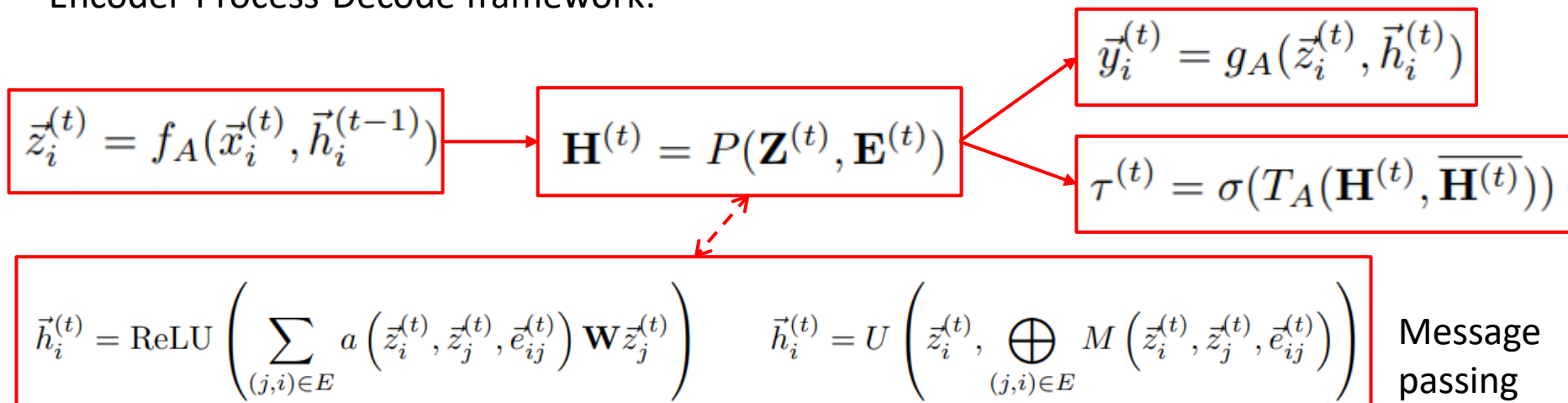
Motivation:

- Many algorithm operates on graphs
- Supervise graph neural networks with algorithm operation/step/final output
- Encoder-Process-Decode framework:

<https://petar-v.com/talks/Algo-WWW.pdf>

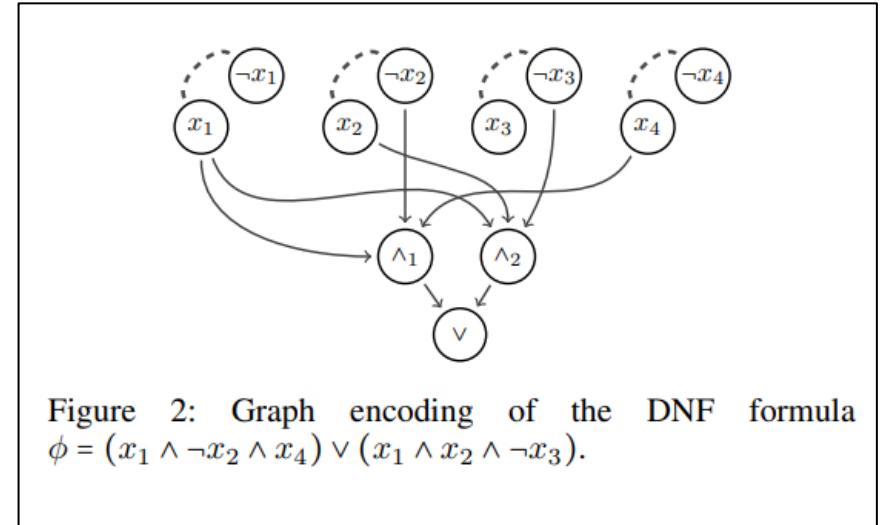
Veličković, Petar, Rex Ying, Matilde Padova, Raia Hadsell, and Charles Blundell.

"Neural Execution of Graph Algorithms." In *International Conference on Learning Representations*. 2019.

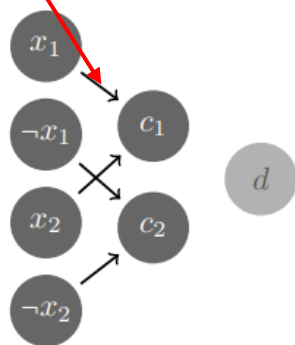


Example: GNN for a specific problem (DNF counting)

- Count #assignments that satisfy disjunctive normal form (DNF) formula
- Classical algorithm is P-hard $O(mn)$
- m : #clauses, n : #variables
- Supervised training

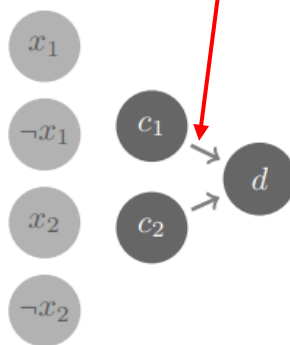


$$\hat{v}_{x_c, t+1} = L_{c_1} \left(v_{x_c, t}, \sum_{x_l \in N(x_l)} M_l(v_{x_l, t}) \right)$$



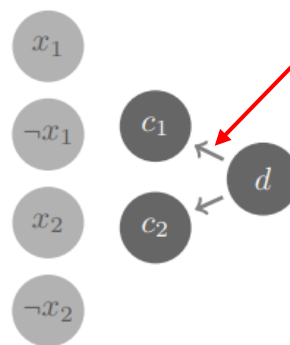
(a)

$$v_{x_d, t+1} = L_d \left(v_{x_d, t}, \sum_{x_c \in N(x_d)} M_c(\hat{v}_{x_c, t+1}) \right)$$



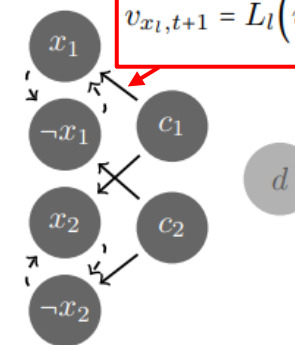
(b)

$$v_{x_c, t+1} = L_{c_2} \left(\hat{v}_{x_c, t+1}, M_d(v_{x_d, t+1}) \right)$$



(c)

$$v_{x_l, t+1} = L_l \left(v_{x_l, t}, \left(\sum_{x_c \in N(x_l)} M_c(v_{x_c, t+1}) \parallel M_l(v_{\neg x_l, t}) \right) \right)$$

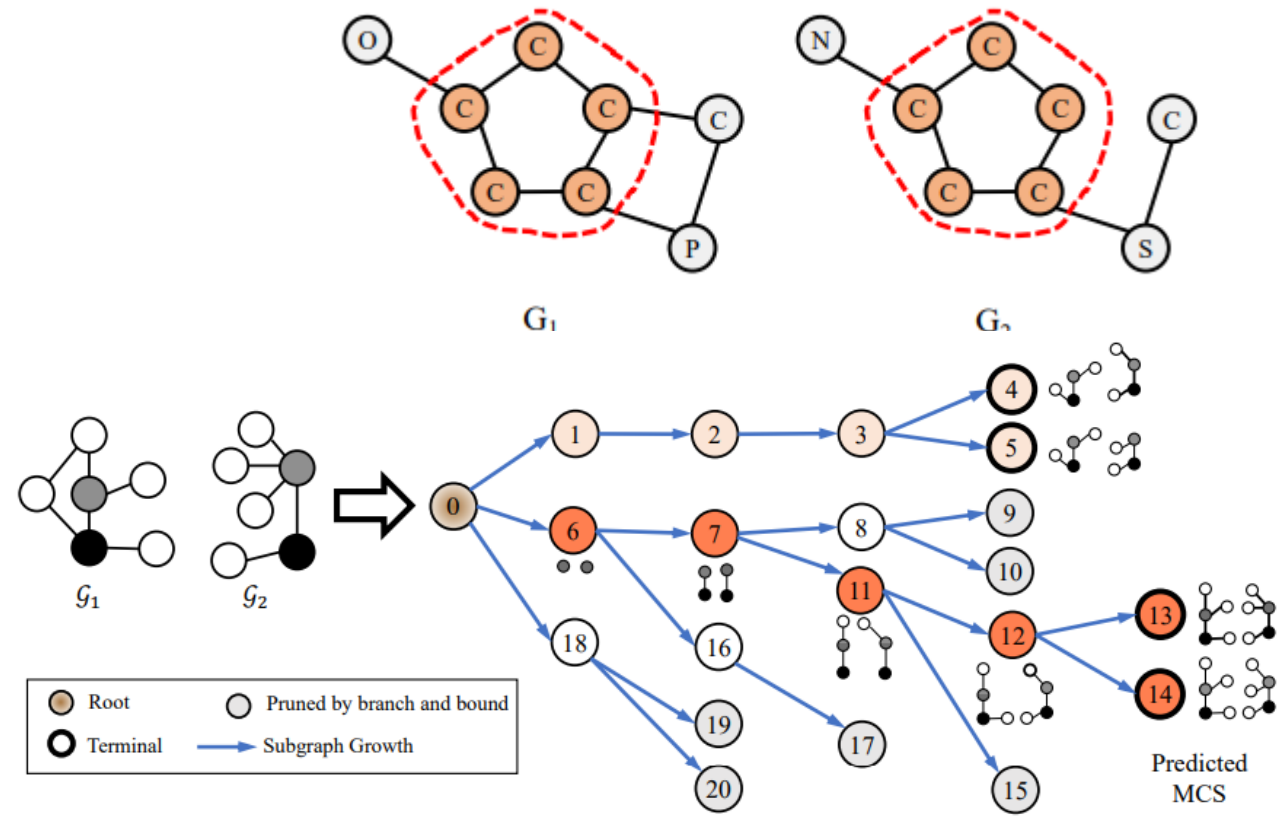


(d)

Best: $O(m+n)$

Example: GNN trained with reinforcement learning (maximum common subgraph)

- Maximum common subgraph (MCS) is NP-hard
- Search for MCS:
 - BFS then pruning
 - Which node to visit first?
- Cast to RL:
 - State:
 - Current subgraph
 - Node-node mapping
 - Input graph
 - Action: Node pair or edge will be visited
 - Reward: +1 if a node pair is selected
 - $Q(s,a)=\text{largest common subgraph size}$



Learning state representation with GNN

Bidomain representation

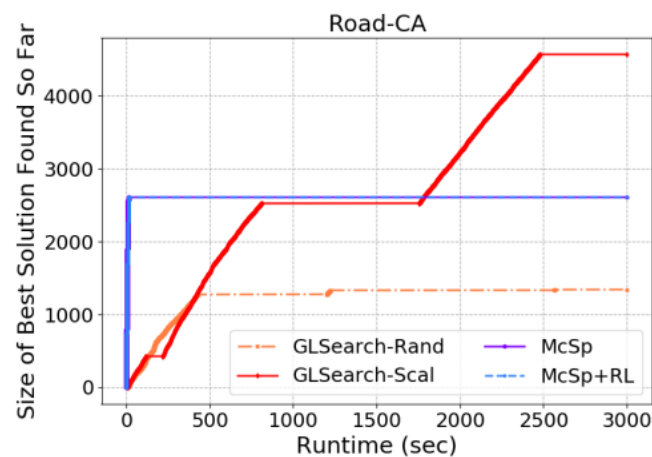
$$Q^*(s_t, a_t), \text{ as } r_t + \gamma V^*(s_{t+1})$$

$$Q(s_t, a_t) = 1 + \gamma \text{MLP} \left(\text{CONCAT} \left(\text{INTERACT}(\mathbf{h}_{\mathcal{G}_1}, \mathbf{h}_{\mathcal{G}_2}), \right. \right. \\ \left. \left. \text{INTERACT}(\mathbf{h}_{s1}, \mathbf{h}_{s2}), \mathbf{h}_{\mathcal{D}_c}, \mathbf{h}_{\mathcal{D}_0} \right) \right).$$

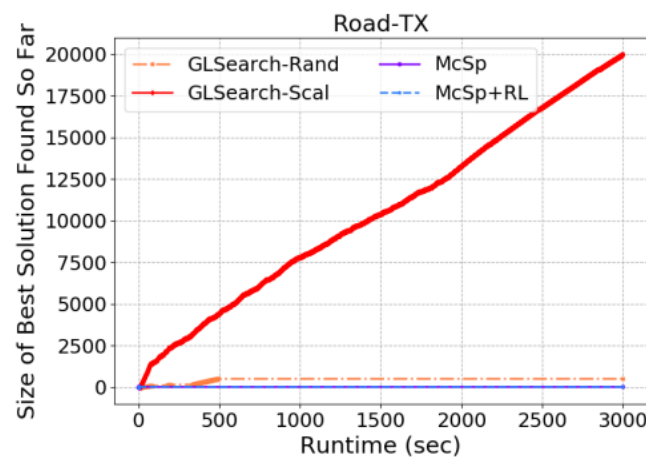
Pretrain with ground-truth Q or expert estimation
Then train as DQN

$$\mathbf{h}_{D_k} = \text{INTERACT} \left(\text{READOUT}(\{\mathbf{h}_i | i \in \mathcal{V}'_{k1}\}), \right. \\ \left. \text{READOUT}(\{\mathbf{h}_j | j \in \mathcal{V}'_{k2}\}) \right).$$

$$\text{READOUT}(\{\mathbf{h}_{D_k} | k \in \mathcal{D}^{(c)}\})$$



(a) Result on ROAD-CA with 978513 nodes.



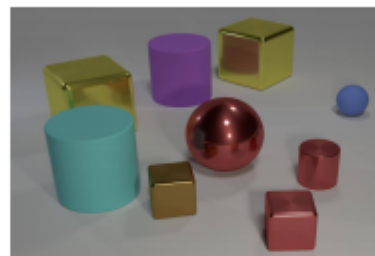
(b) Result on ROAD-TX with 1080909 nodes.

Method	ROAD	DBEN
GLSEARCH (no $\mathbf{h}_{\mathcal{G}}$)	0.977	0.878
GLSEARCH (no \mathbf{h}_s)	1.000	0.874
GLSEARCH (no $\mathbf{h}_{\mathcal{D}_c}$)	0.803	0.780
GLSEARCH (no $\mathbf{h}_{\mathcal{D}_0}$)	0.576	0.856
GLSEARCH (SUM interact)	0.902	0.913
GLSEARCH (unfactored)	0.447	0.807
GLSEARCH (unfactored-i)	0.500	0.789
GLSEARCH	0.992	1.000
BEST SOLUTION SIZE	132	508

Neural networks and algorithms alignment



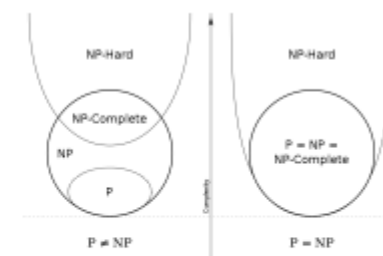
Summary statistics
What is the maximum value difference among treasures?



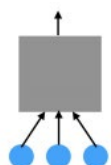
Relational argmax
What are the colors of the furthest pair of objects?



Dynamic programming
What is the cost to defeat monster X by following the optimal path?

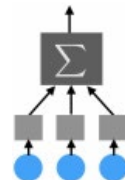


NP-hard problem
Subset sum: Is there a subset that sums to 0?



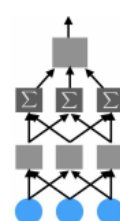
MLPs
~ feature extraction

$$y = \text{MLP}(\|_{s \in S} X_s)$$



Deep Sets (Zaheer et al.,
~ summary statistics

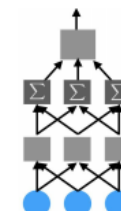
$$y = \text{MLP}_2 \left(\sum_{s \in S} \text{MLP}_1(X_s) \right)$$



GNNs
~ (pairwise) relations

$$h_s^{(k)} = \sum_{t \in S} \text{MLP}_1^{(k)} \left(h_s^{(k-1)}, h_t^{(k-1)} \right)$$

$$y = \text{MLP}_2 \left(\sum_{s \in S} h_s^{(K)} \right)$$



GNNs
~ (pairwise) relations

Neural exhaustive search

GNN is aligned with Dynamic Programming (DP)

Graph Neural Network

for $k = 1 \dots$ GNN iter:

for u in S :

No need to learn for-loops

$$h_u^{(k)} = \sum_v \text{MLP}(h_v^{(k-1)}, h_u^{(k-1)})$$

Bellman-Ford algorithm

for $k = 1 \dots |S| - 1$:

for u in S :

$$d[k][u] = \min_v d[k-1][v] + \text{cost}(v, u)$$

Learns a simple reasoning step

$h_u^{(k)}$

$d[k][u]$

\sum_v

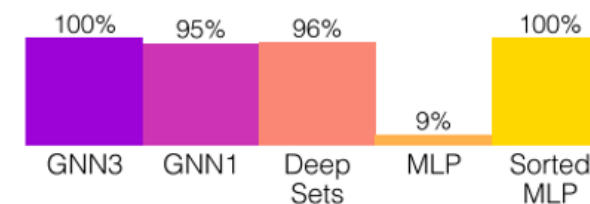
\min_v

$\text{MLP}(h_v^{(k-1)}, h_u^{(k-1)})$

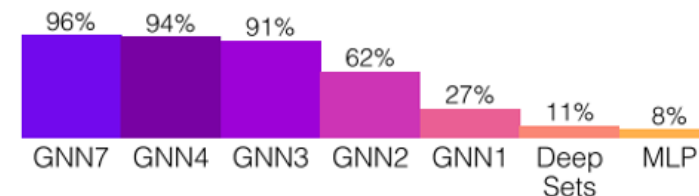
$d[k-1][v] + \text{cost}(v, u)$

Neural exhaustive
search

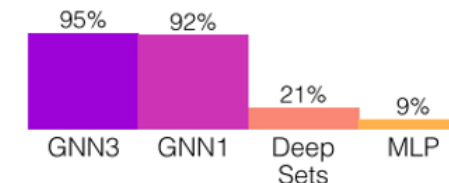
$$\text{MLP}_2(\max_{\tau \subseteq S} \text{MLP}_1 \circ \text{LSTM}(X_1, \dots, X_{|\tau|} : X_1, \dots, X_{|\tau|} \in \tau)).$$



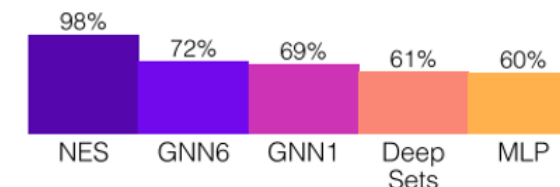
(a) Maximum value difference.



(c) Monster trainer.



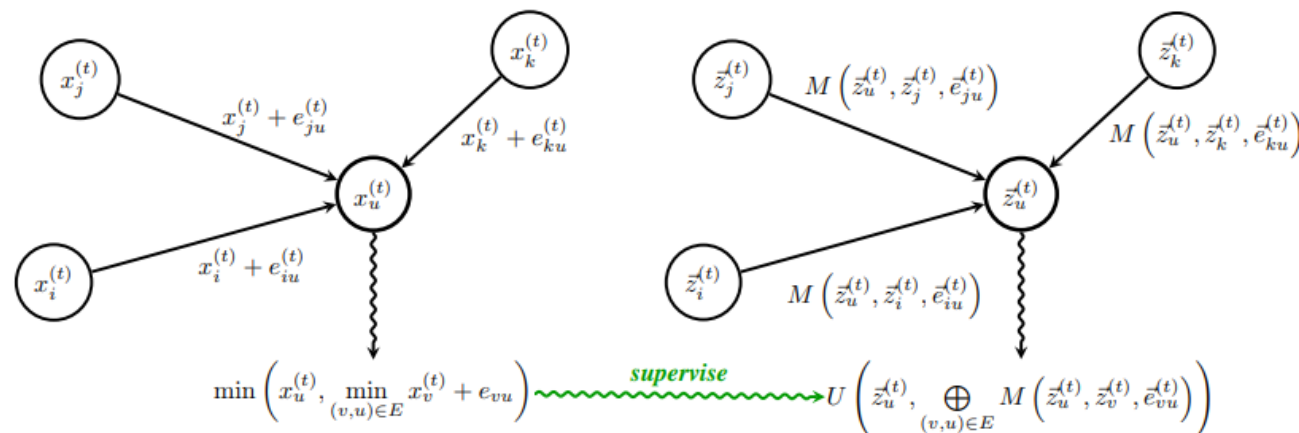
(b) Furthest pair.



(d) Subset sum. Random guessing yields 50%.

If alignment exists \rightarrow step-by-step supervision

- Merely simulate the classical graph algorithm, generalizable
- No algorithm discovery



Algorithm	Inputs	Supervision signals
Breadth-first search	$x_i^{(t)}$: is i reachable from s in $\leq t$ hops?	$x_i^{(t+1)}$, $\tau^{(t)}$: has the algorithm terminated?
Bellman-Ford	$x_i^{(t)}$: shortest distance from s to i (using $\leq t$ hops)	$x_i^{(t+1)}$, $\tau^{(t)}$, $p_i^{(t)}$: predecessor of i in the shortest path tree (in $\leq t$ hops)
Prim's algorithm	$x_i^{(t)}$: is node i in the (partial) MST (built from s after t steps)?	$x_i^{(t+1)}$, $\tau^{(t)}$, $p_i^{(t)}$: predecessor of i in the partial MST

Joint training is encouraged

Table 1: Accuracy of predicting reachability at different test-set sizes, trained on graphs of 20 nodes. GAT* correspond to the best GAT setup as per Section 3 (GAT-full using the full graph).

Model	Reachability (mean step accuracy / last-step accuracy)		
	20 nodes	50 nodes	100 nodes
LSTM (Hochreiter & Schmidhuber, 1997)	81.97% / 82.29%	88.35% / 91.49%	68.19% / 63.37%
GAT* (Veličković et al., 2018)	93.28% / 99.86%	93.97% / 100.0%	92.34% / 99.97%
GAT-full* (Vaswani et al., 2017)	78.40% / 77.86%	85.76% / 91.83%	88.98% / 91.51%
MPNN-mean (Gilmer et al., 2017)	100.0% / 100.0%	61.05% / 57.89%	27.17% / 21.40%
MPNN-sum (Gilmer et al., 2017)	99.66% / 100.0%	94.25% / 100.0%	94.72% / 98.63%
MPNN-max (Gilmer et al., 2017)	100.0% / 100.0%	100.0% / 100.0%	99.92% / 99.80%

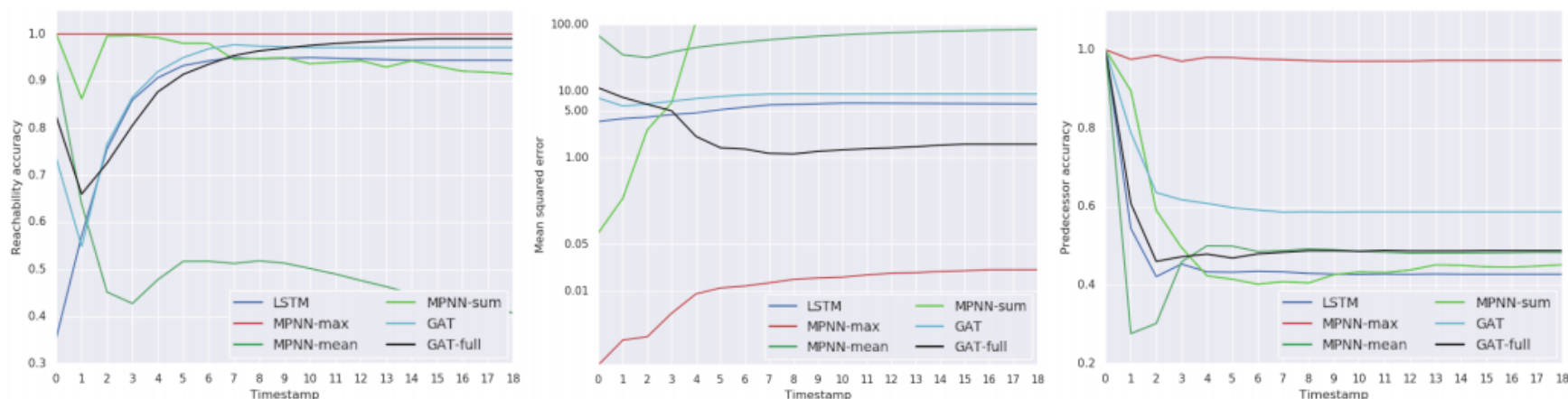
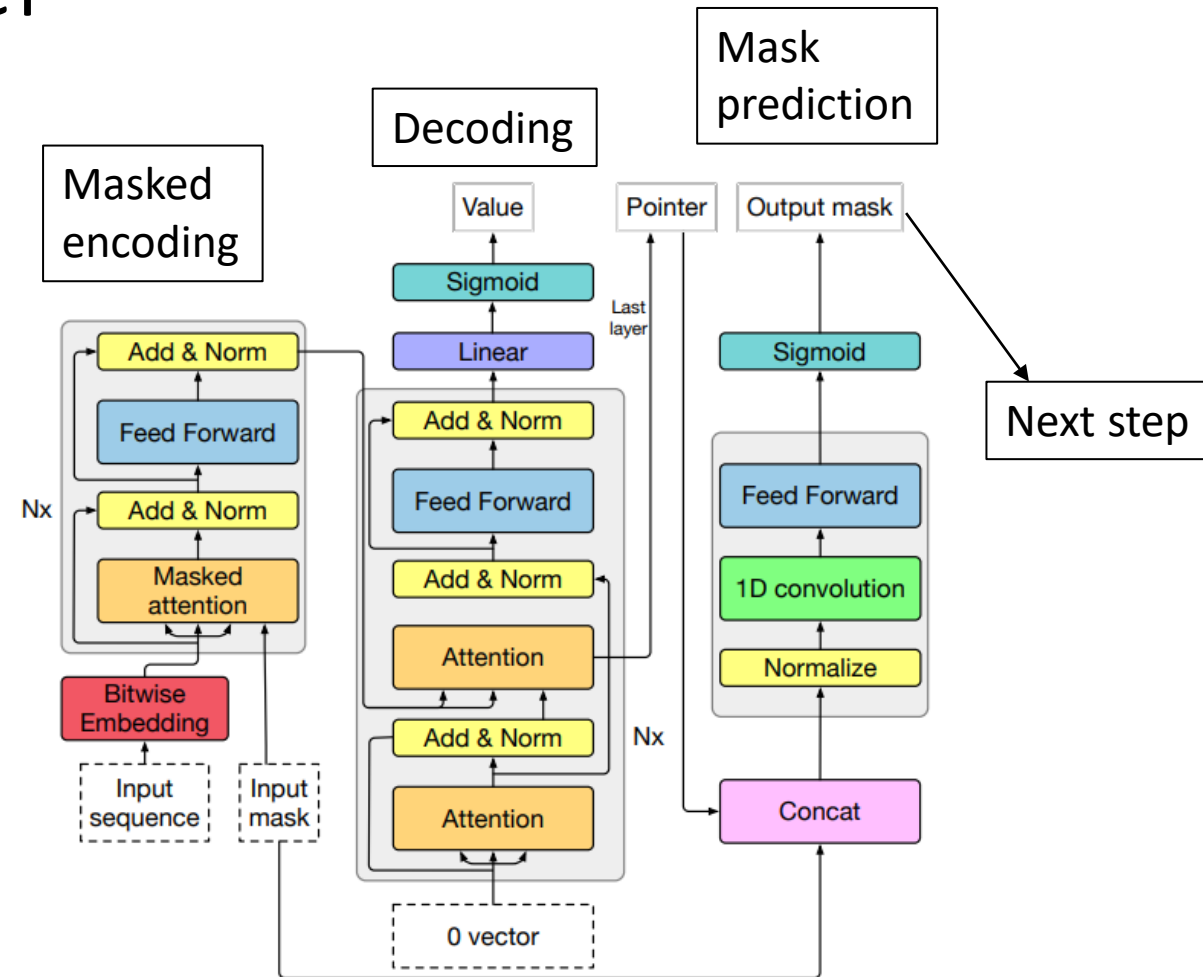


Figure 3: The per-step algorithm execution performances in terms of reachability accuracy (**left**), distance mean-squared error (**middle**) and predecessor accuracy (**right**), tested on 100-node graphs after training on 20-node graphs. Please mind the scale of the MSE plot.

Processor as Transformer

- Back to input sequence (set), but stronger generalization
- Transformer with encoder mask ~ graph attention
- Use Transformer with:
 - Binary representation of numbers
 - Dynamic conditional masking

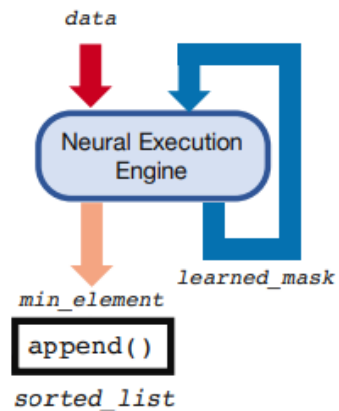


Training with execution trace

```

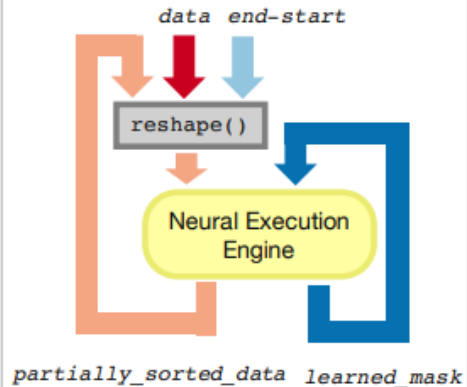
selection_sort(data):
    sorted_list = []
    while (len(data) > 0):
        min_index, min_element = find_min(data)
        data.delete(min_index)
        sorted_list.append(min_element)
    return sorted_list

find_min(data):
    min_element = -1
    min_index = -1
    for index, element in enumerate(data):
        if (element < min_element):
            min_element = element
            min_index = index
    return [min_index, min_element]
    
```



```

merge_sort(data, start, end):
    if (start < end):
        mid = (start + end) / 2
        merge_sort(data, start, mid)
        merge_sort(data, mid+1, end)
    return merge(data, start, mid, end)
    
```



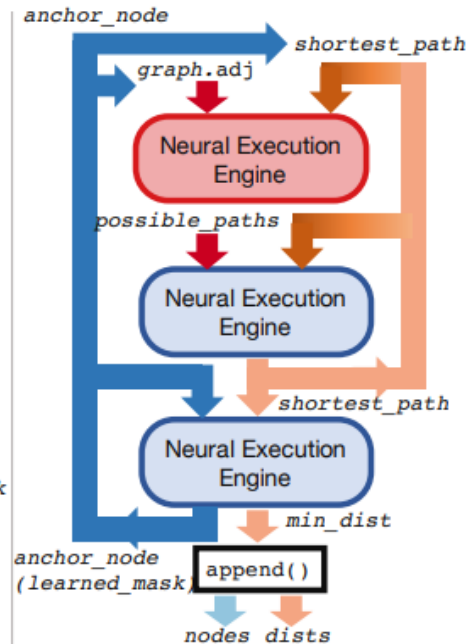
```

shortest_path(graph, source_node, shortest_path):
    dists = []
    nodes = []
    anchor_node = source_node
    node_list = graph.get_nodes()

    while node_list:
        possible_paths = sum(graph.adj(anchor_node),
                               shortest_path(anchor_node))
        shortest_path = min(possible_paths, shortest_path)
        anchor_node, min_dist = min(shortest_path(node_list))

        node_list.delete(anchor_node)
        nodes.append(anchor_node)
        dists.append(min_dist)

    return dists, nodes
    
```



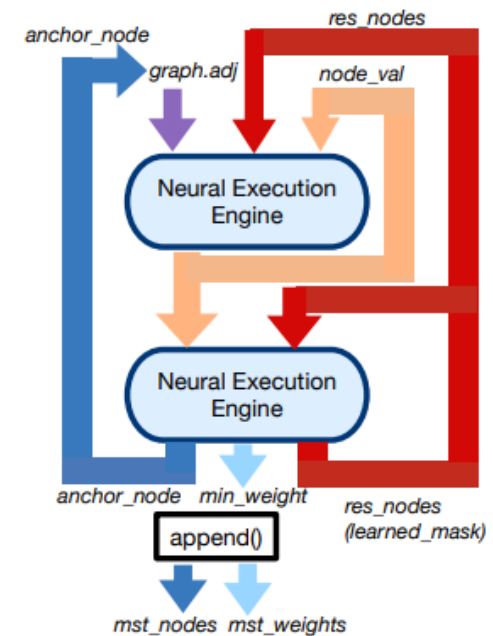
```

minimum_spanning_tree(graph, source_node, node_val):
    mst_nodes = []
    mst_weights = []
    anchor_node = source_node
    res_nodes = graph.get_nodes()

    while node_list:
        adj_list = graph.adj(anchor_node)
        node_val(res_nodes) = min(node_val(res_nodes),
                                   adj_list(res_nodes))
        anchor_node, min_weight = min(node_val(res_nodes))

        mst_nodes.append(anchor_node)
        mst_weights.append(min_weight)
        res_nodes.delete(anchor_node)

    return mst_nodes, mst_weights
    
```



The results show strong generalization

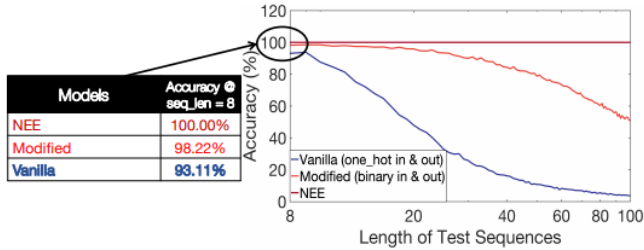


Figure 3: Sorting performance of transformers trained on sequences of up to length 8.

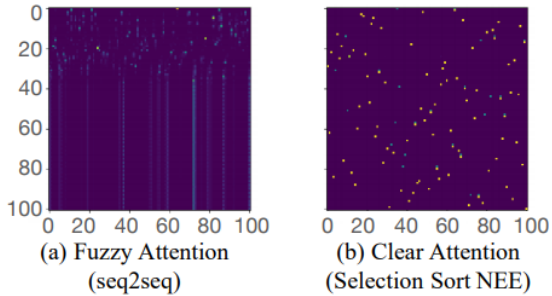


Figure 4: Visualizing decoder attention weights. Attention is over each row. Transformer attention saturates as the output sequence length increases, while NEE maintains sharp attention.

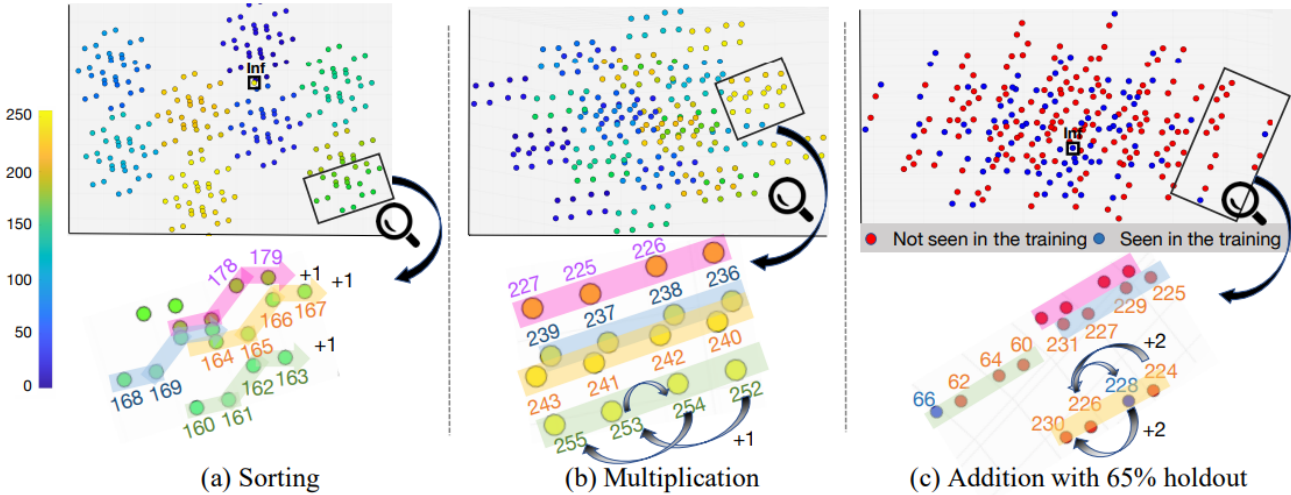


Figure 5: 3D PCA visualization of learned bitwise embeddings for different numeric tasks. The embeddings exhibit regular, task-dependent structure, even when most numbers have not been seen in training (c).

	Sizes	25	50	75	100
Accuracy					
Selection sort		100.00	100.00	100.00	100.00
Merge sort		100.00	100.00	100.00	100.00
Shortest path		100.00	100.00	100.00	100.00*
Minimum spanning tree		100.00	100.00	100.00	100.00

QA