

Tutorial at KDD, August 14<sup>th</sup> 2021

# From Deep Learning to Deep Reasoning

## Part A: Learning to reason

Truyen Tran, Vuong Le, Hung Le and Thao Le

{truyen.tran,vuong.le,hung.le,thao.le}@deakin.edu.au

<https://bit.ly/37DYQn7>

# Logistics



Truyen Tran



Vuong Le



Hung Le



Thao Le

<https://bit.ly/37DYQn7>

# Agenda

- **Introduction**
- Part A: Learning-to-reason framework
- Part B: Reasoning over unstructured and structured data
- Part C: Memory | Data efficiency | Recursive reasoning

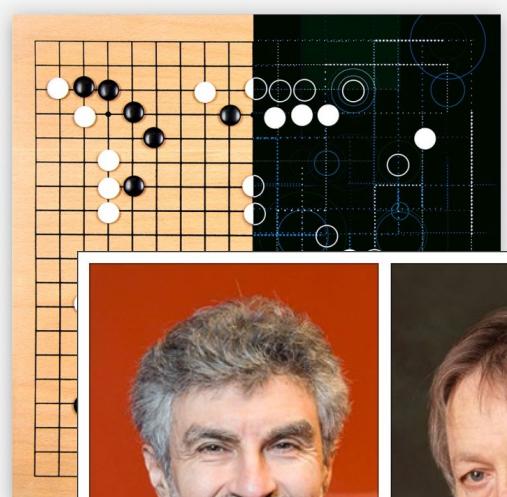


2012

# DL: 8 years snapshot



# AlphaGo



Yoshua Bengio



Ge



Geoffrey Hinton  
@geoffreyhinton

## Turing A

Extrapolating the spectacular performance of GPT3 into the future suggests that the answer to life, the universe and everything is just 4.398 trillion parameters.

6:26 AM · Jun 11, 2020 · Twitter Web App

1 Retweets

75 Quote Tweets

3.8K Likes

# GPT-3 2020

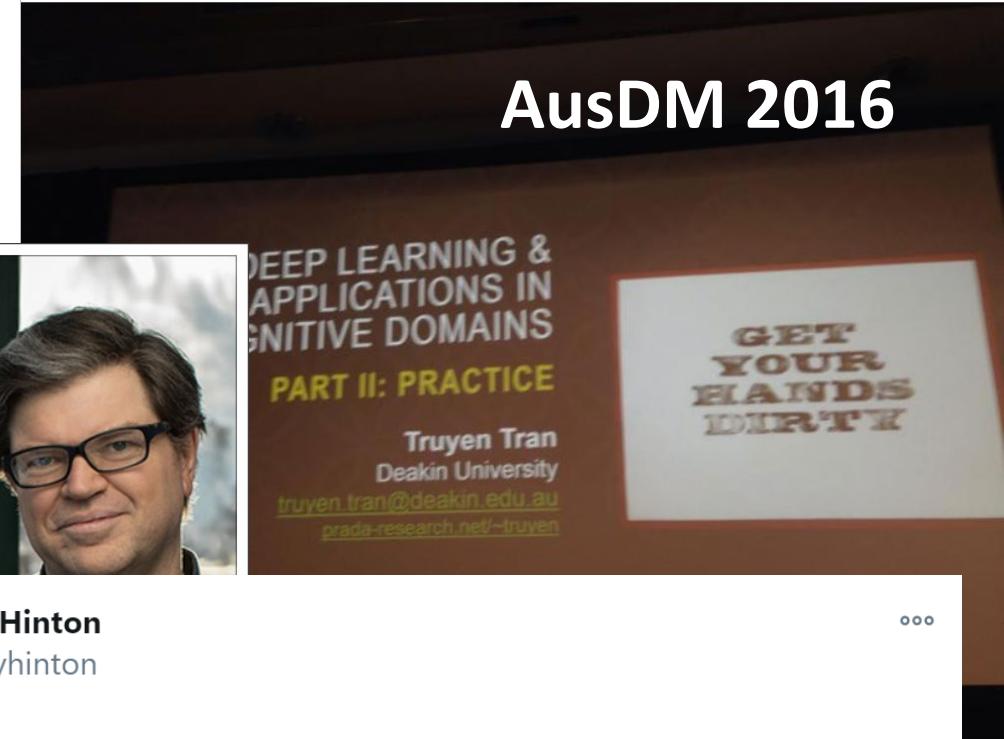


Xingjun (Daniel) Ma • 1st

Lecturer@Deakin University, Machine Learning and AI Scientist  
4yr

...

Excellent deep learning tutorial in AusDM2016 by Truyen Tran!



# AusDM 2016

# DL has been fantastic, but ...

- It is great at interpolating
  - → data hungry to cover all variations and smooth local manifolds
  - → little systematic generalization (novel combinations)
- Lack of human-perceived reasoning capability
  - Lack natural mechanism to incorporate prior knowledge, e.g., common sense
  - No built-in causal mechanisms
  - → Have trust issues!
- To be fair, many of these problems are common in statistical learning!

# Why still DL in 2021?

## Theoretical

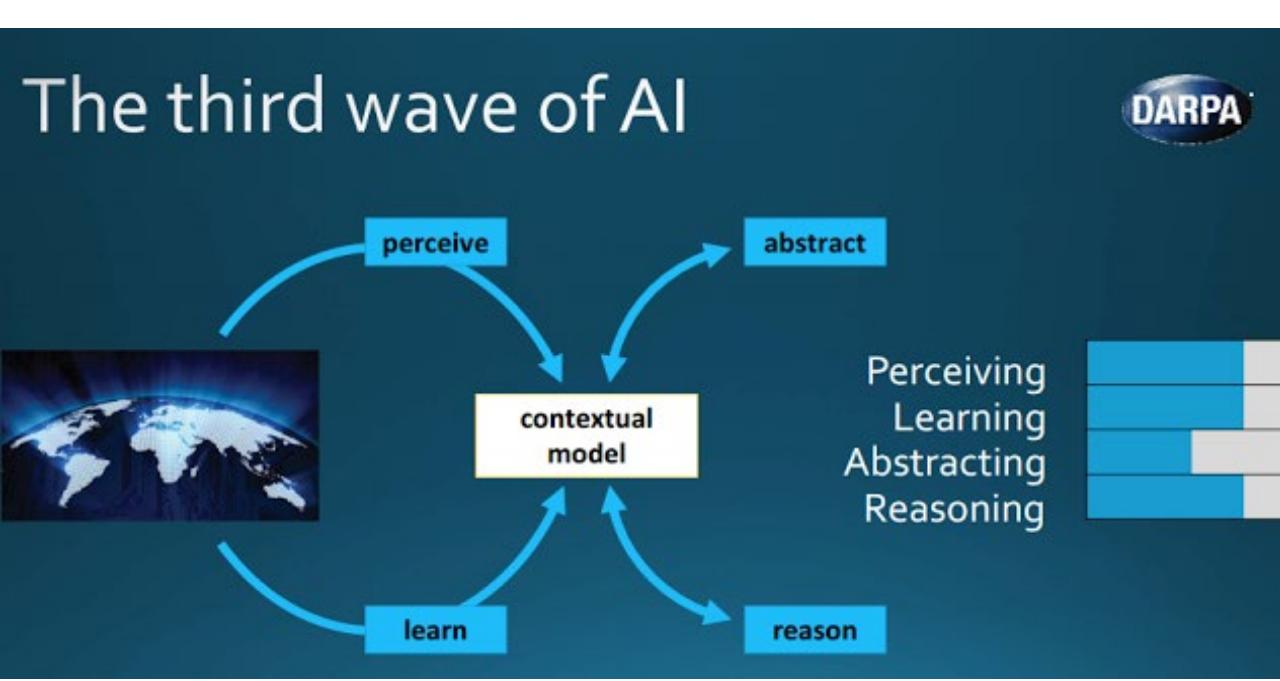
- **Expressiveness:** Neural nets can approximate any function.
- **Learnability:** Neural nets are trained easily.
- **Generalisability:** Neural nets generalize surprisingly well to unseen data.

## Practical

- **Generality:** Applicable to many domains.
- **Competitive:** DL is hard to beat as long as there are data to train.
- **Scalability:** DL is better with more data, and it is very scalable.

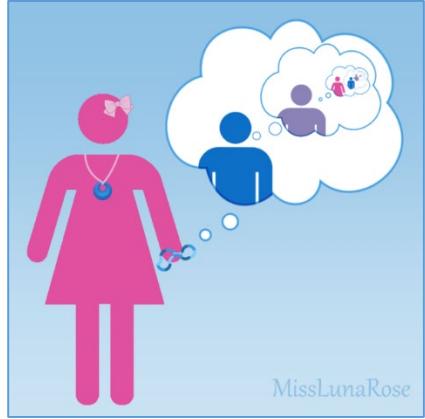


# The next AI/ML challenge



2020s-2030s

- Learning + reasoning, general purpose, human-like
- Has contextual and common-sense reasoning
- Requires less data
- Adapt to change
- Explainable



Theory of mind  
Recursive reasoning

Memory

Facts  
Semantics  
Events and relations  
Working space



Perception

Multiple

System 1:  
Intuitive

- Fast
- Implicit/automatic
- Pattern recognition
- Multiple

Single

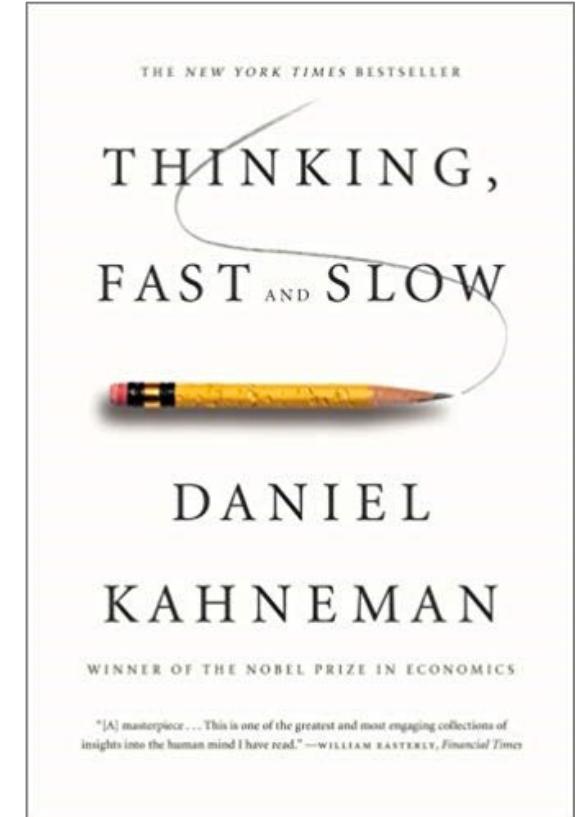
System 2:  
Analytical

- Slow
- Deliberate/rational
- Careful analysis
- Single, sequential

# Toward deeper reasoning

# System 2

- Holds hypothetical thought
- Decoupling from representation
- Working memory size is not essential.  
Its attentional control is.



# Reasoning in Probabilistic Graphical Models (PGM)

- Assuming models are fully specified (e.g., by hand or learnt)

- Estimate MAP as energy minimization
- Compute marginal probability
- Compute expectation & normalisation constant

- Key algorithm: **Pearl's Belief Propagation**, a.k.a Sum-Product algorithm in factor graphs.

- Known result in 2001-2003: BP minimises **Bethe free-energy minimization**.

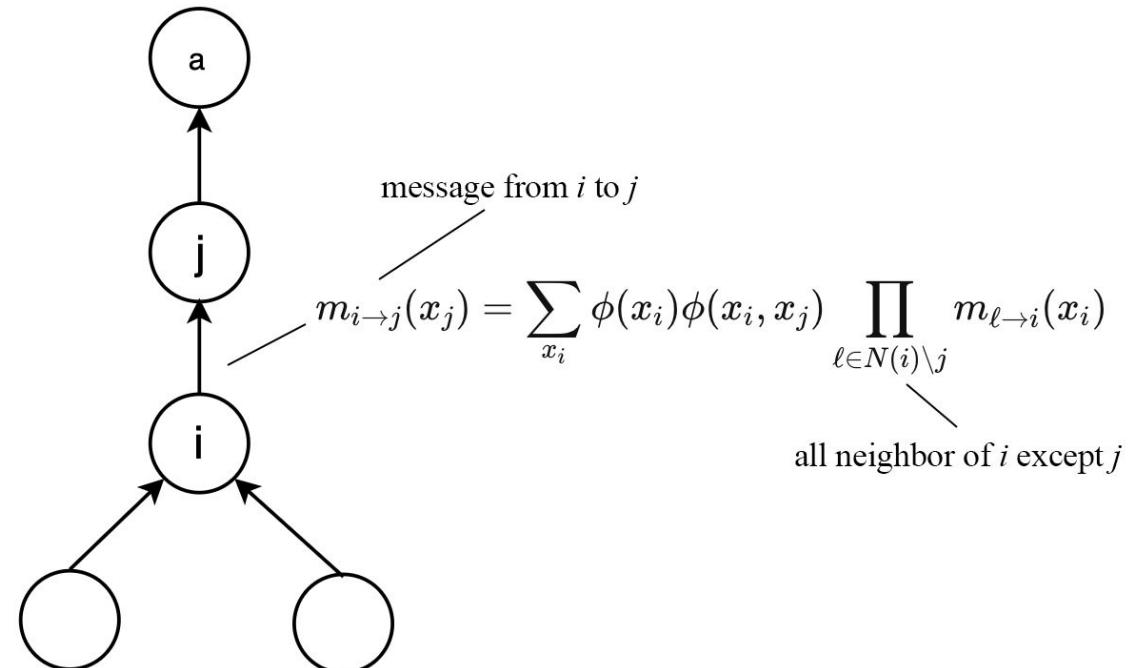


Figure credit: Jonathan Hui

**Can we learn to infer directly from data  
without full specification of models?**

# Agenda

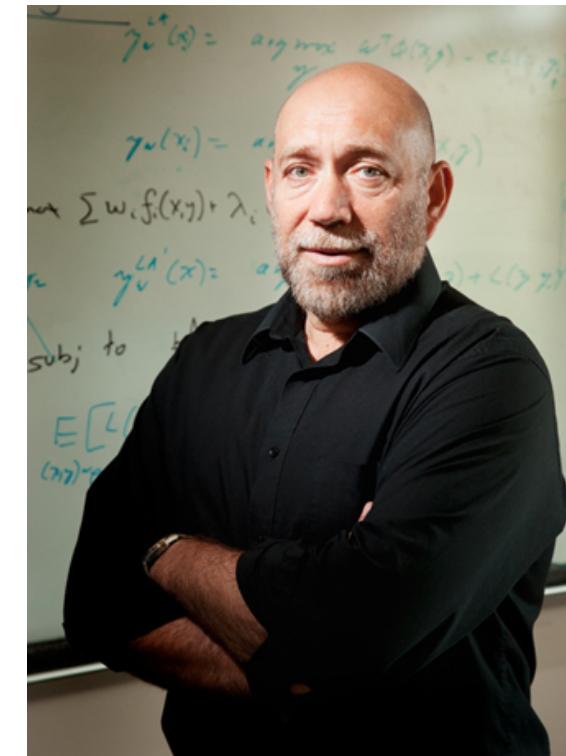
- Introduction
- **Part A: Learning-to-reason framework**
- Part B: Reasoning over unstructured and structured data
- Part C: Memory | Data efficiency | Recursive reasoning

# Part A: Sub-topics

- Reasoning as a prediction skill that can be learnt from data.
  - Question answering as zero-shot learning.
- Neural network operations for learning to reason:
  - Concept-object binding.
  - Attention & transformers.
  - Dynamic neural networks, conditional computation & differentiable programming.
- Reasoning as iterative representation refinement & query-driven program synthesis and execution
  - Compositional attention networks.
  - Neural module networks.
- Combinatorics reasoning

# Learning to reason

- Learning is to self-improve by experiencing ~ acquiring knowledge & skills
- Reasoning is to deduce knowledge from previously acquired knowledge in response to a query (or a cues)
- Learning to reason is to improve the ability to decide if a knowledge base entails a predicate.
  - E.g., given a video  $f$ , determines if the person with the hat turns before singing.
- Hypotheses:
  - Reasoning as just-in-time **program synthesis**.
  - It employs **conditional computation**.



(Dan Roth; ACM Fellow; IJCAI  
John McCarthy Award)

Khardon, Roni, and Dan Roth. "Learning to reason." *Journal of the ACM* (JACM) 44.5 (1997): 697-725.

# Learning to reason, a definition

*Definition 2.1.1.* An algorithm  $A$  is an *exact reasoning* algorithm for the reasoning problem  $(\mathcal{F}, \mathcal{Q})$ , if for all  $f \in \mathcal{F}$  and for all  $\alpha \in \mathcal{Q}$ , when  $A$  is presented with input  $(f, \alpha)$ ,  $A$  runs in time polynomial in  $n$  and the size of  $f$  and  $\alpha$ , and answers “yes” if and only if  $f \vDash \alpha$ .

E.g., given a video  $f$ , determines if the person with the hat turns before singing.

# Practical setting: (query, database, answer) **triplets**

- This is very general:
  - Classification: Query = *what is this?* Database = *data*.
  - Regression: Query = *how much?* Database = *data*.
  - QA: Query = *NLP question*. Database = *context/image/text*.
  - Multi-task learning: Query = *task ID*. Database = *data*.
  - Zero-shot learning: Query = *task description*. Database = *data*.
  - Drug-protein binding: Query = *drug*. Database = *protein*.
  - Recommender system: Query = *User (or item)*. Database = *inventories (or user base)*;

# Can neural networks reason?

Reasoning is not necessarily achieved by making logical inferences

There is a continuity between  
[algebraically rich inference] and  
[connecting together trainable learning systems]

Central to reasoning is composition rules to guide the combinations of modules to address new tasks



"When we observe a visual scene, when we hear a complex sentence, we are able to explain in formal terms the relation of the objects in the scene, or the precise meaning of the sentence components. However, there is no evidence that such a formal analysis necessarily takes place: we see a scene, we hear a sentence, and we just know what they mean. **This suggests the existence of a middle layer, already a form of reasoning, but not yet formal or logical.**"

# Hypotheses

- Reasoning as just-in-time program synthesis.
- It employs conditional computation.
- Reasoning is recursive, e.g., mental travel.

# Two approaches to neural reasoning

- **Implicit chaining of predicates through recurrence:**
  - Step-wise query-specific attention to relevant concepts & relations.
  - Iterative concept refinement & combination, e.g., through a working memory.
  - Answer is computed from the last memory state & question embedding.
- **Explicit program synthesis:**
  - There is a set of modules, each performs an pre-defined operation.
  - Question is parse into a symbolic program.
  - The program is implemented as a computational graph constructed by chaining separate modules.
  - The program is executed to compute an answer.

# In search for basic neural operators for reasoning

- Basics:
  - Neuron as feature detector → Sensor, filter
  - Computational graph → Circuit
  - Skip-connection → Short circuit
- Essentials
  - Multiplicative gates → AND gate, Transistor, Resistor
  - Attention mechanism → SWITCH gate
  - Memory + forgetting → Capacitor + leakage
  - Compositionality → Modular design
- ..

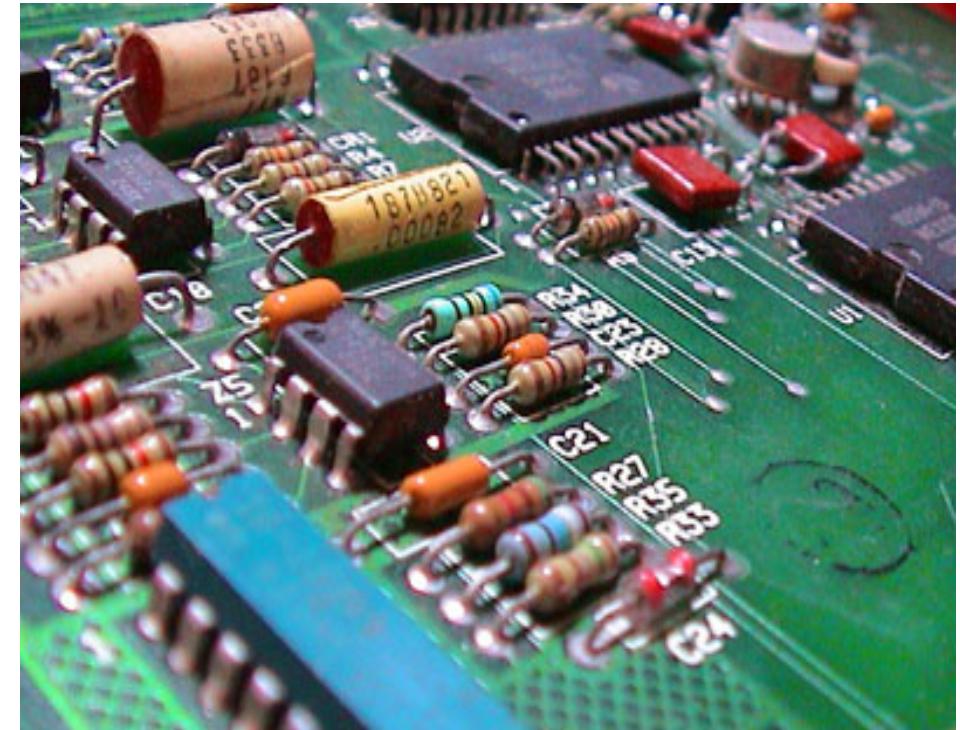


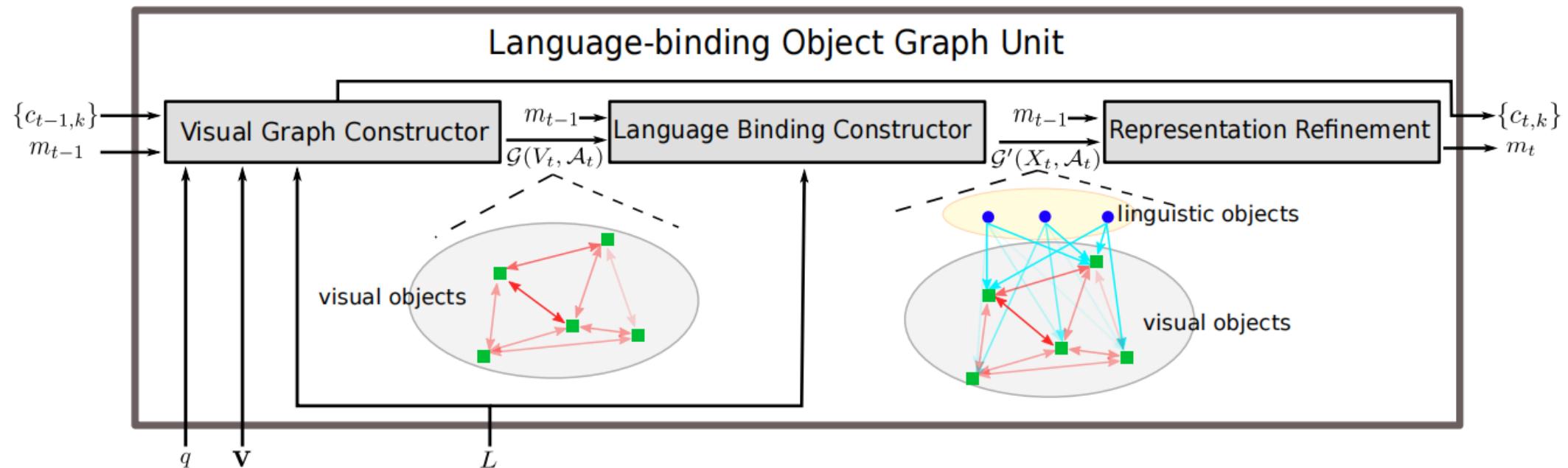
Photo credit: [Nicola Asuni](#)

# Part A: Sub-topics

- Reasoning as a prediction skill that can be learnt from data.
- Question answering as zero-shot learning.
- **Neural network operations for learning to reason:**
  - Concept-object binding.
  - Attention & transformers.
  - Dynamic neural networks, conditional computation & differentiable programming.
- Reasoning as iterative representation refinement & query-driven program synthesis and execution.
  - Compositional attention networks.
  - Reasoning as Neural module networks.
- Combinatorics reasoning

# Concept-object binding

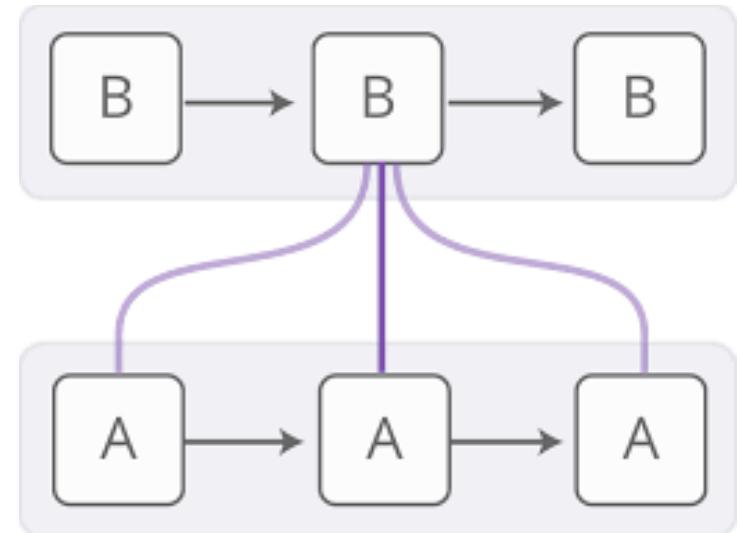
- Perceived data (e.g., visual objects) may not share the same semantic space with high-level concepts.
- Binding between concept-object enables reasoning at the concept level



Example of concept-object binding in LOGNet (Le et al, IJCAI'2020)

# Attentions: Picking up only what is needed at a step

- Need attention model to select or ignore certain **computations** or **inputs**
  - Can be “soft” (differentiable) or “hard” (requires RL)
  - Needed for selecting predicates in reasoning.
- Attention provides a short-cut → long-term dependencies
  - Needed for long chain of reasoning.
  - Also encourages sparsity if done right!



<http://distill.pub/2016/augmented-rnns/>

# Fast weights | HyperNet – the multiplicative interaction

- Early ideas in early 1990s by Juergen Schmidhuber and collaborators.
- **Data-dependent weights** | Using a controller to generate weights of the main net.

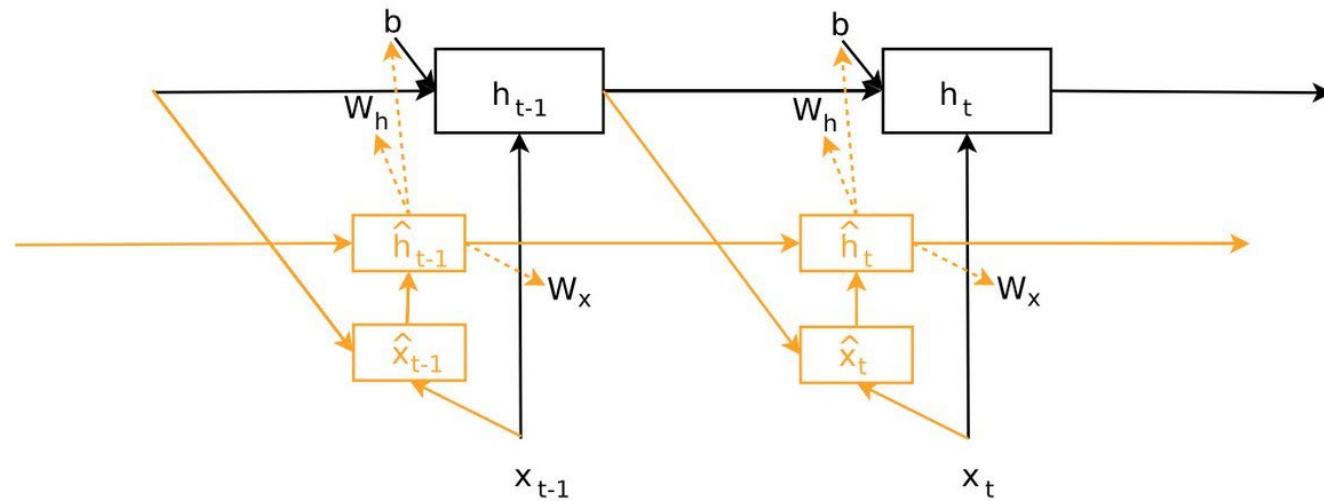
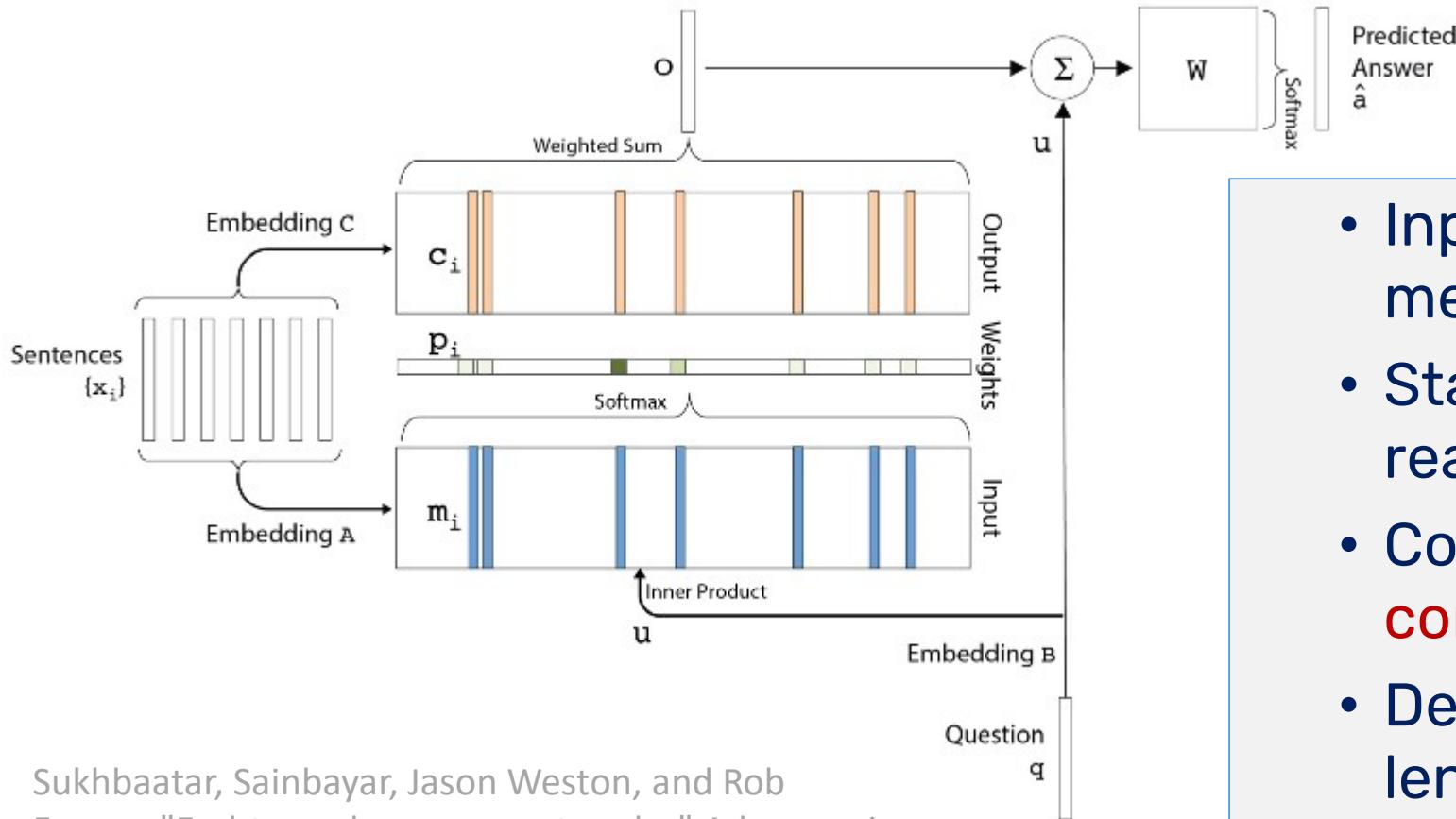


Figure: The HyperRNN system. The black system represents the main RNN while the orange system represents the weight-generating HyperRNN cell.

Ha, David, Andrew Dai, and Quoc V. Le. "Hypernetworks." *arXiv preprint arXiv:1609.09106* (2016).

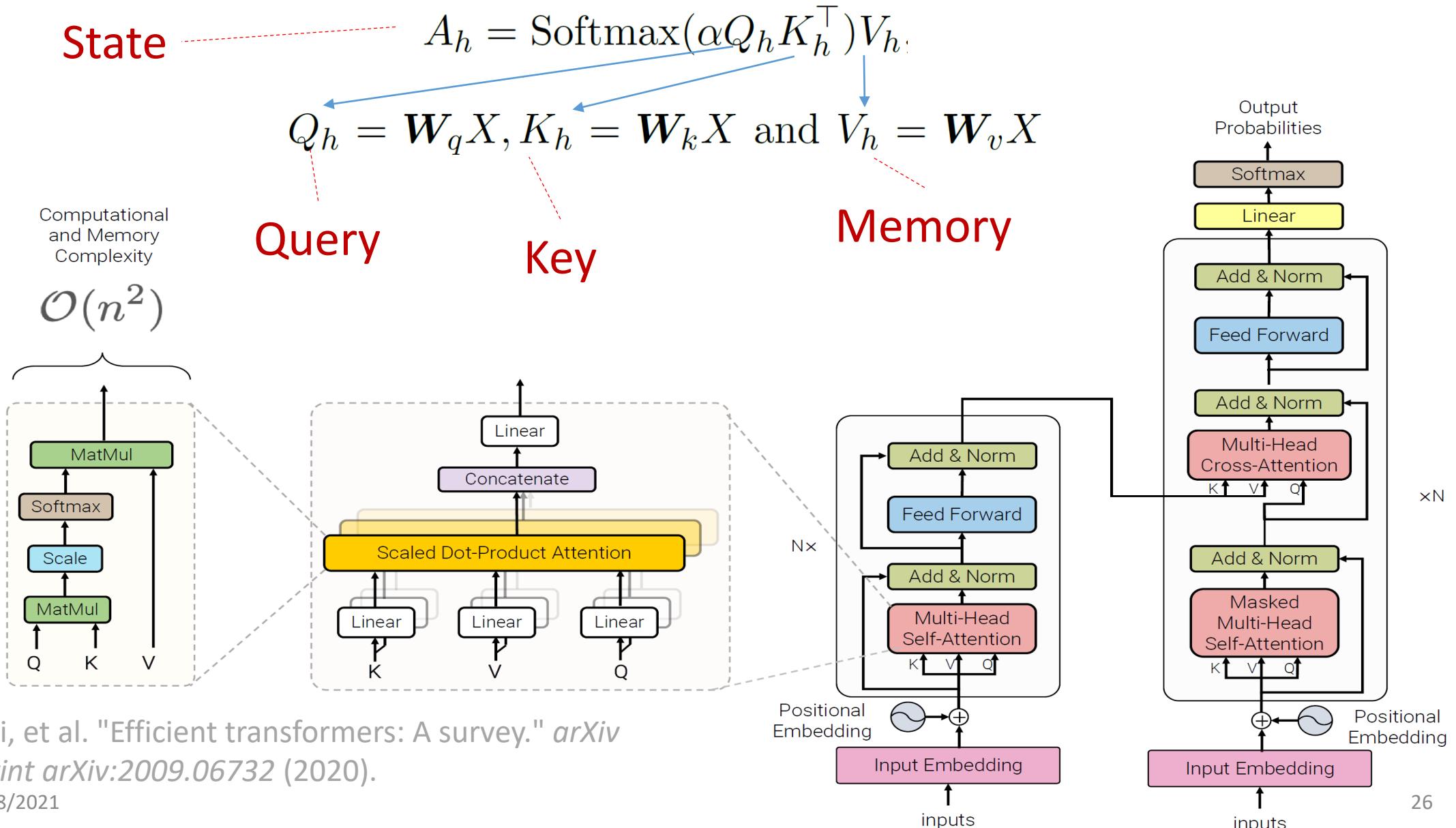
# Memory networks: Holding the data ready for inference



Sukhbaatar, Sainbayar, Jason Weston, and Rob Fergus. "End-to-end memory networks." Advances in neural information processing systems. 2015.

- Input is a set → Load into memory, which is NOT updated.
- State is a RNN with attention reading from inputs
- Concepts: **Query**, **key** and **content** + **Content addressing**.
- Deep models, but constant path length from input to output.
- Equivalent to a RNN with shared input set.

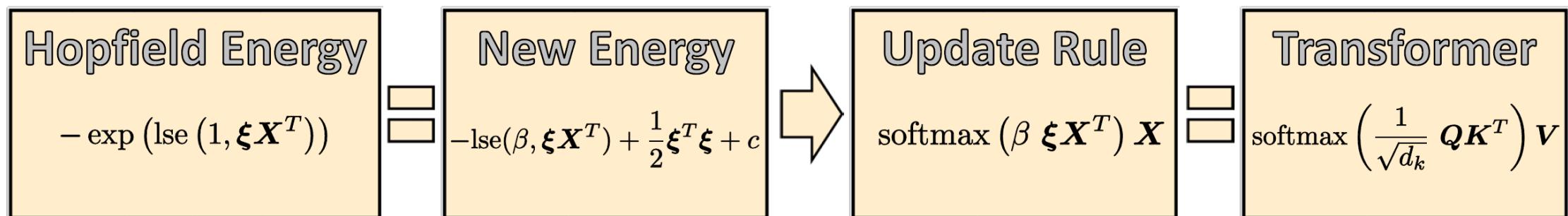
# Transformers: Analogical reasoning through self-attention



Tay, Yi, et al. "Efficient transformers: A survey." *arXiv preprint arXiv:2009.06732* (2020).

# Transformer as implicit reasoning

- Recall: Reasoning as (free-) energy minimisation
  - The classic Belief Propagation algorithm is minimization algorithm of the Bethe free-energy!
- Transformer has relational, iterative state refinement makes it a great candidate for implicit relational reasoning.



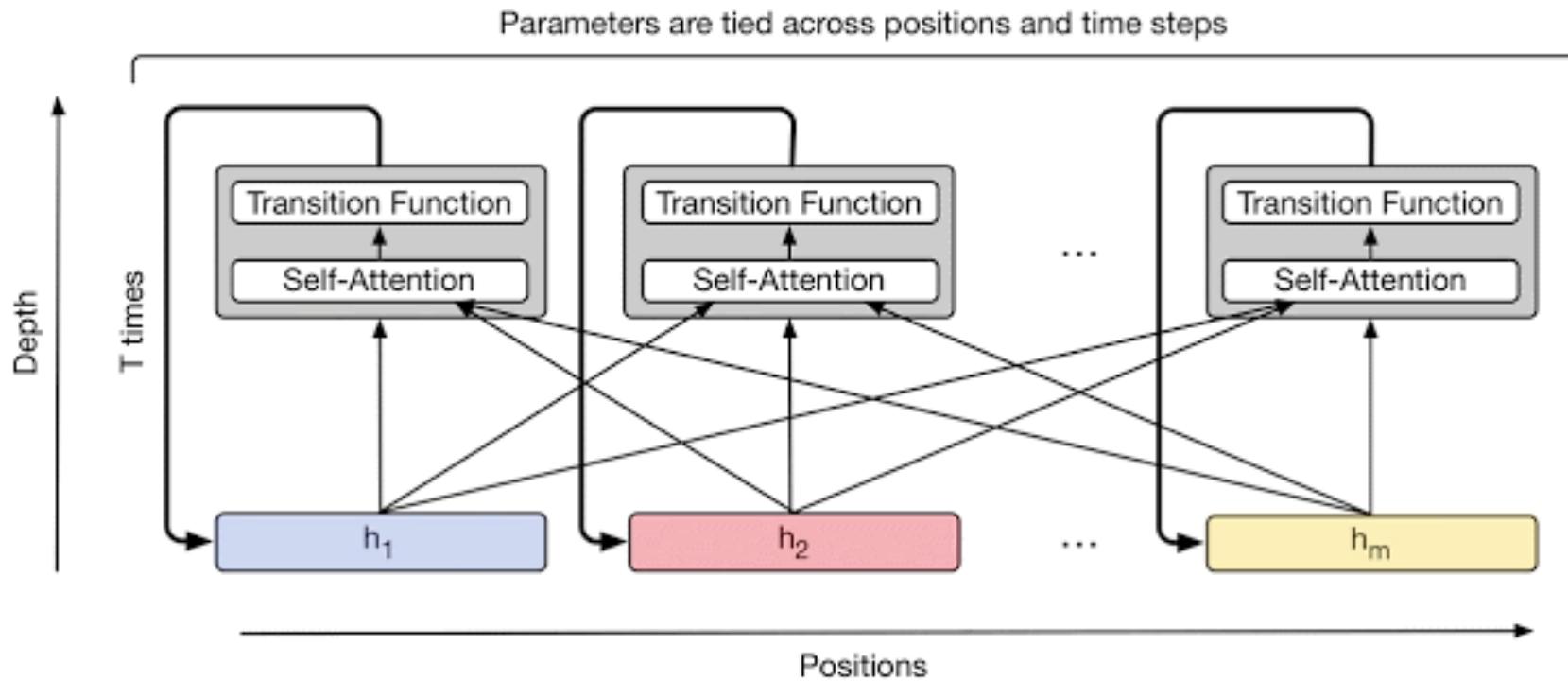
Ramsauer, Hubert, et al. "Hopfield networks is all you need." *arXiv preprint arXiv:2008.02217* (2020).

# Transformer v.s. memory networks

- Memory network:
  - Attention to input set
  - One hidden state update at a time.
  - Final state integrate information of the set, conditioned on the query.
- Transformer:
  - Loading all inputs into working memory
  - Assigns one hidden state per input element.
  - All hidden states (including those from the query) to compute the answer.

# Universal transformers

Dehghani, Mostafa, et al. "Universal Transformers." *International Conference on Learning Representations*. 2018.



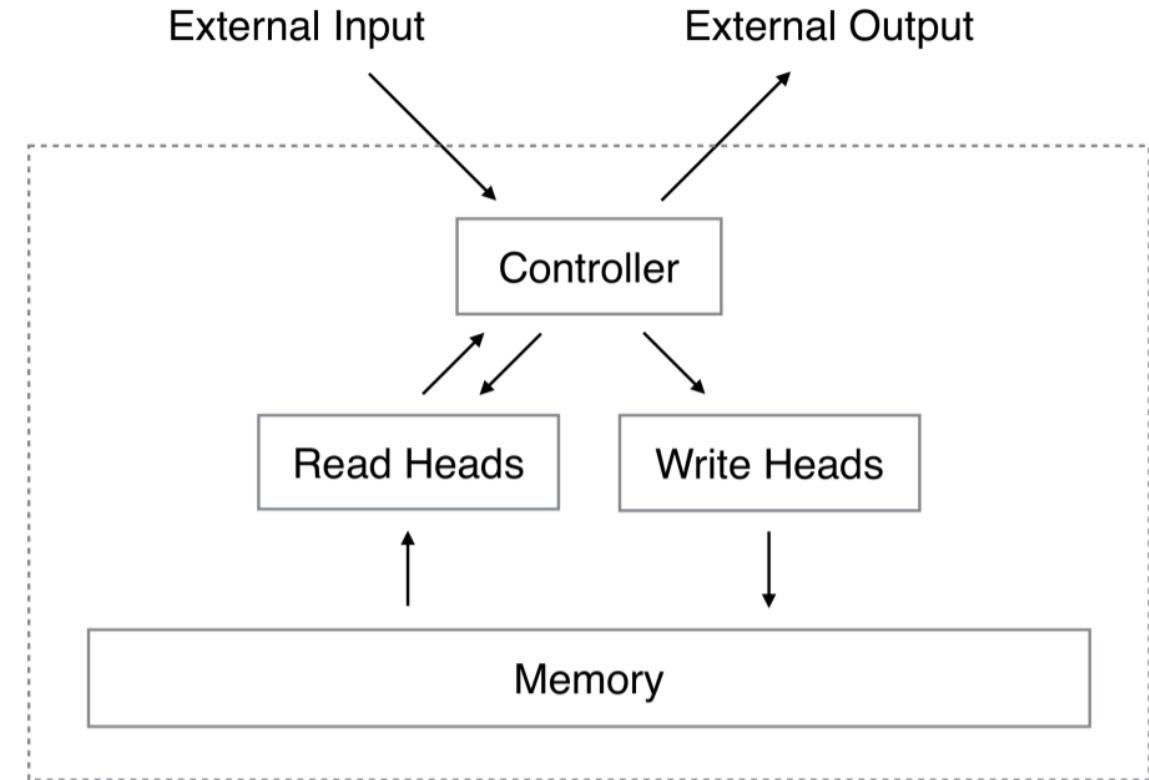
# Dynamic neural networks

- Memory-Augmented Neural Networks
- Modular program layout
- Program synthesis

# Neural Turing machine (NTM)

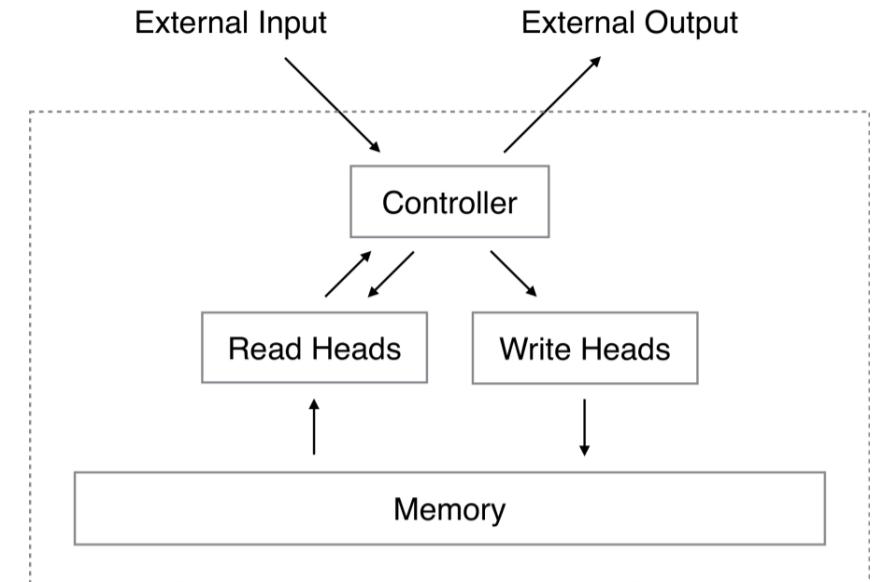
## A memory-augmented neural network (MANN)

- A controller that takes input/output and talks to an external memory module.
- Memory has read/write operations.
- The main issue is where to write, and how to update the memory state.
- All operations are differentiable.



# MANN for reasoning

- Three steps:
  - Store data into memory
  - Read query, process sequentially, consult memory
  - Output answer
- Behind the scene:
  - Memory contains data & results of intermediate steps
  - LOGNet does the same, memory consists of object representations
- Drawbacks of current MANNs:
  - No memory of controllers → Less modularity and compositionality when query is complex
  - No memory of relations → Much harder to chain predicates.

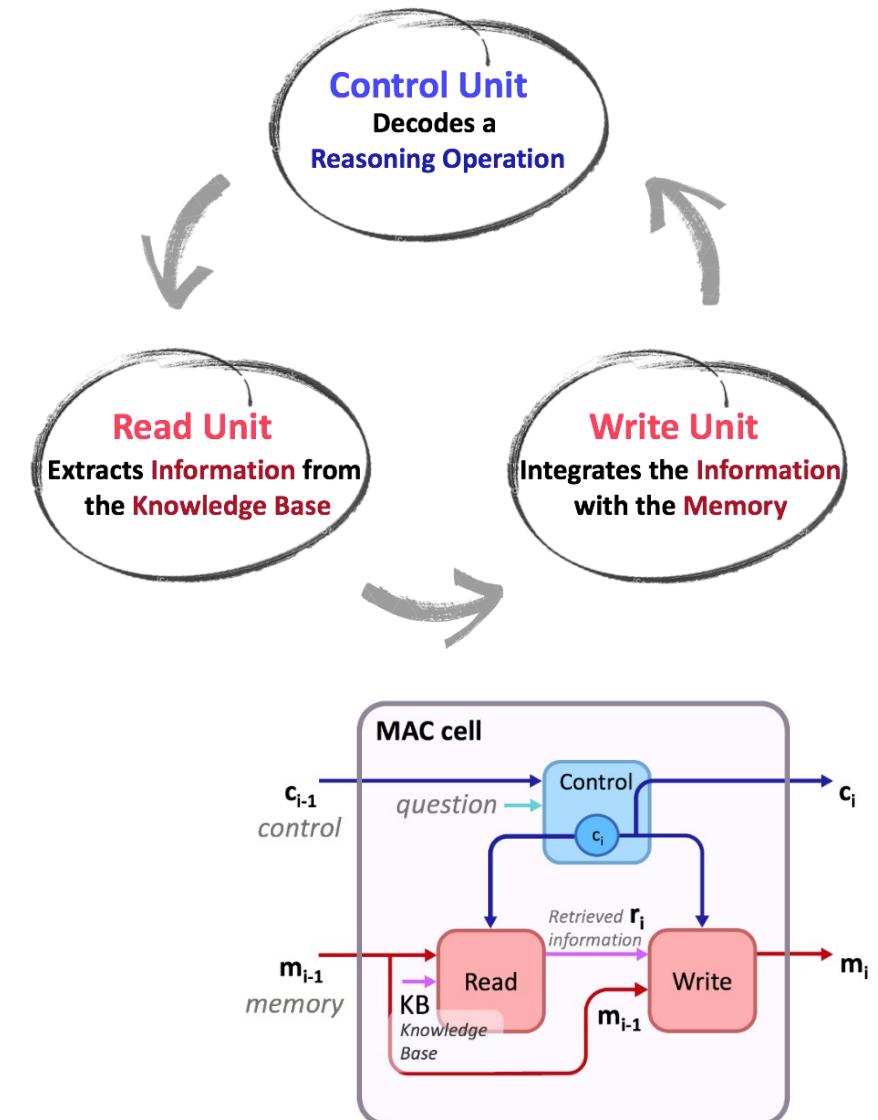
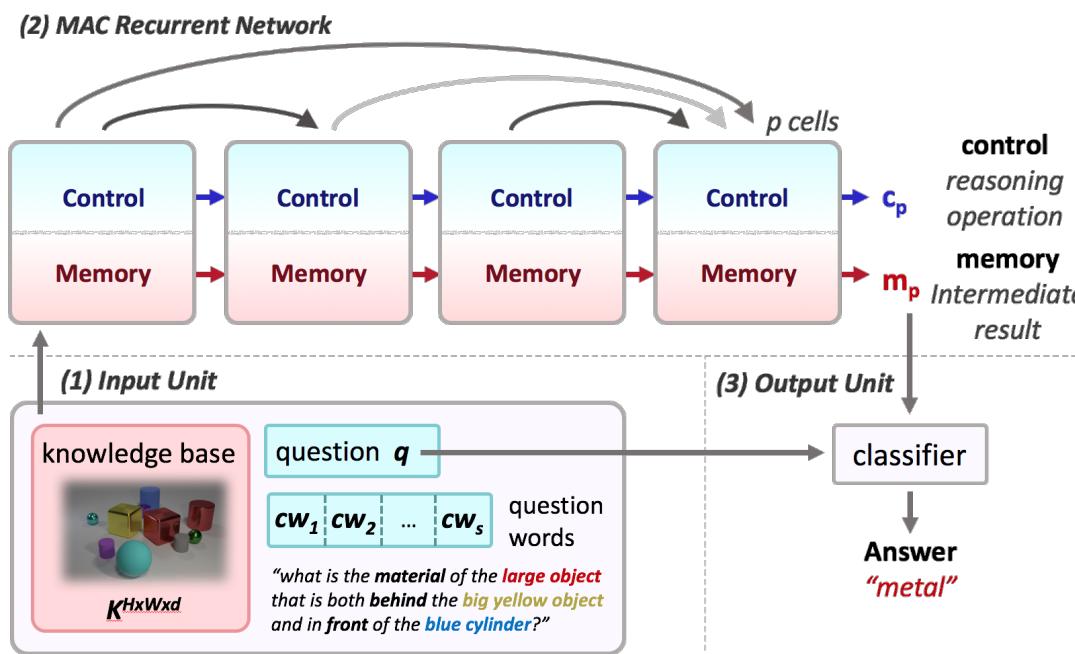


Source: [rylanschaeffer.github.io](https://rylanschaeffer.github.io)

# Part A: Sub-topics

- Reasoning as a prediction skill that can be learnt from data.
  - Question answering as zero-shot learning.
- Neural network operations for learning to reason:
  - Concept-object binding.
  - Attention & transformers.
  - Dynamic neural networks, conditional computation & differentiable programming.
- **Reasoning as iterative representation refinement & query-driven program synthesis and execution.**
  - Compositional attention networks.
  - Reasoning as Neural module networks.
- Combinatorics reasoning

# MAC Net: Recurrent, iterative representation refinement



Hudson, Drew A., and Christopher D. Manning. "Compositional attention networks for machine reasoning." *ICLR* 2018.

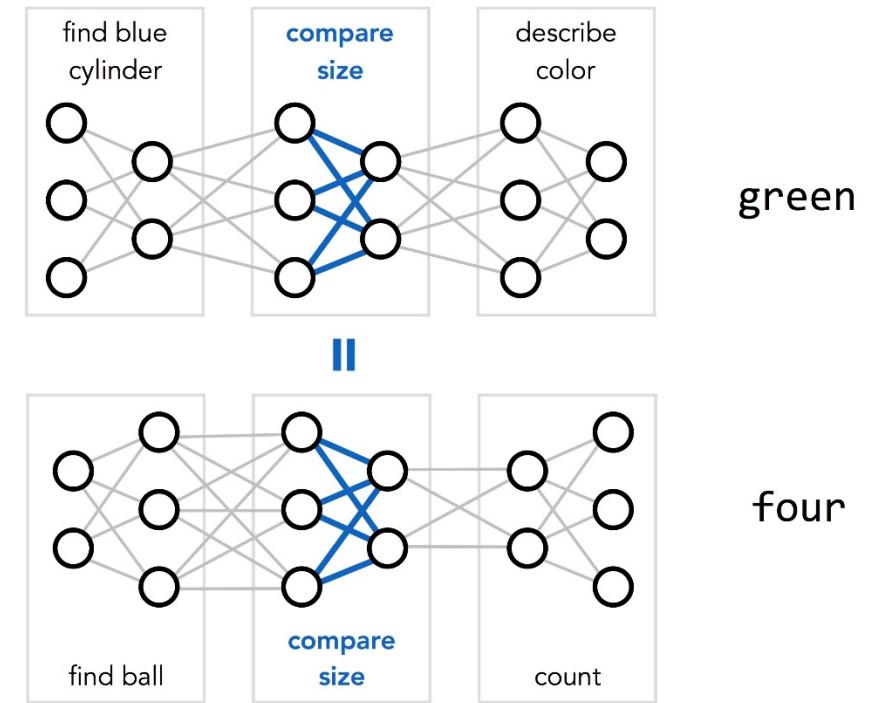
# Module networks

## (reasoning by constructing and executing neural programs)

- Reasoning as laying out modules to reach an answer
- Composable neural architecture → question parsed as program (layout of modules)
- A module is a function ( $x \rightarrow y$ ), could be a sub-reasoning process  $((x, q) \rightarrow y)$ .

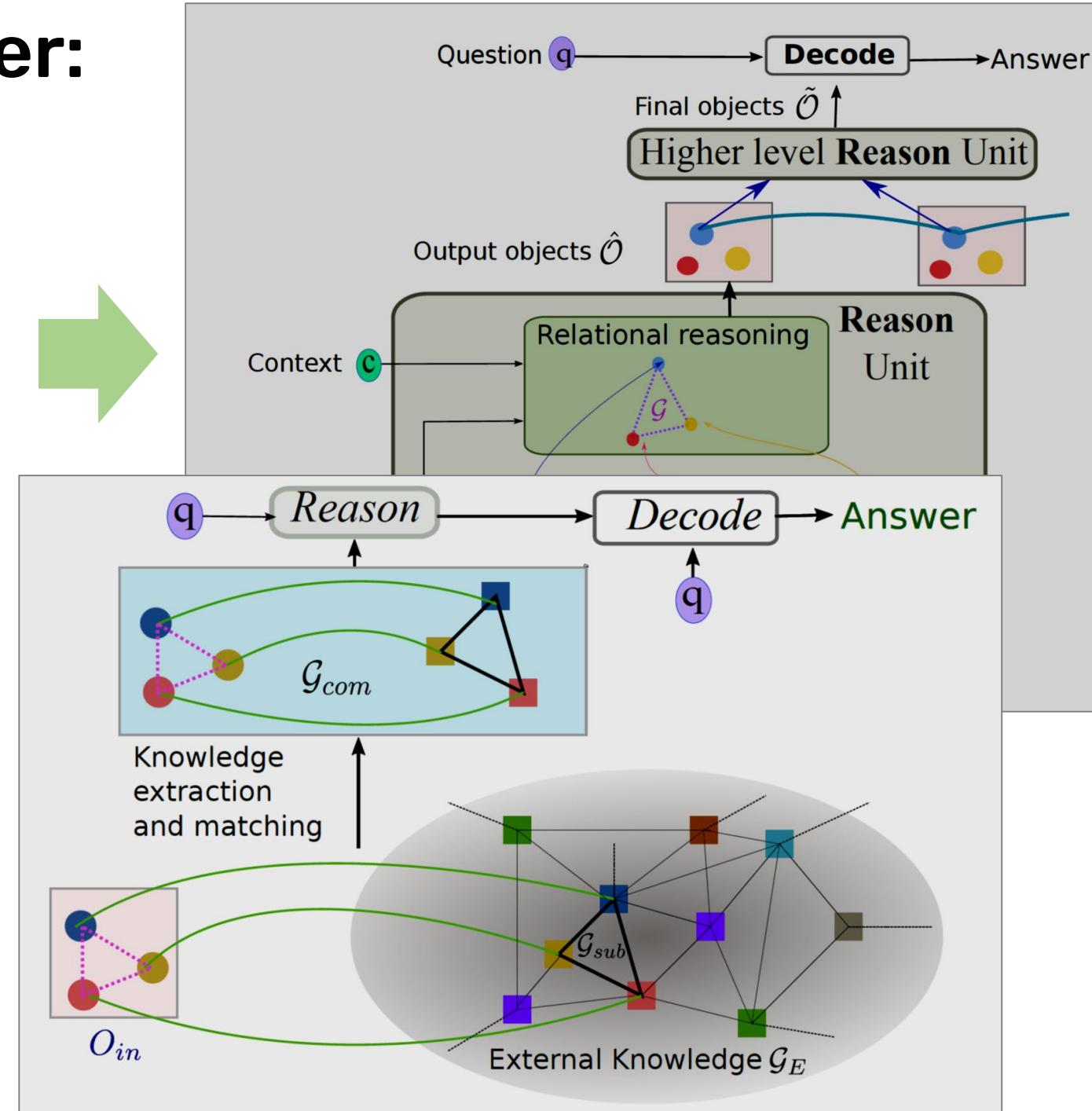
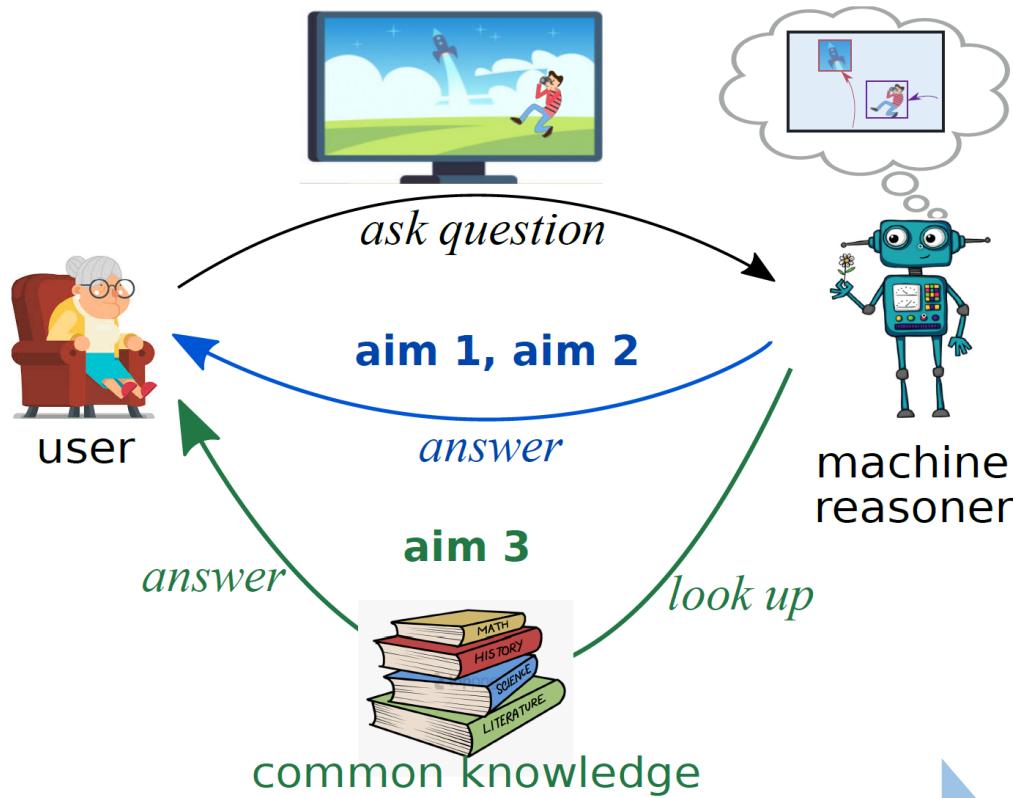


What color is the thing with the same size as the blue cylinder?



How many things are the same size as the ball?

# Putting things together: A framework for visual reasoning



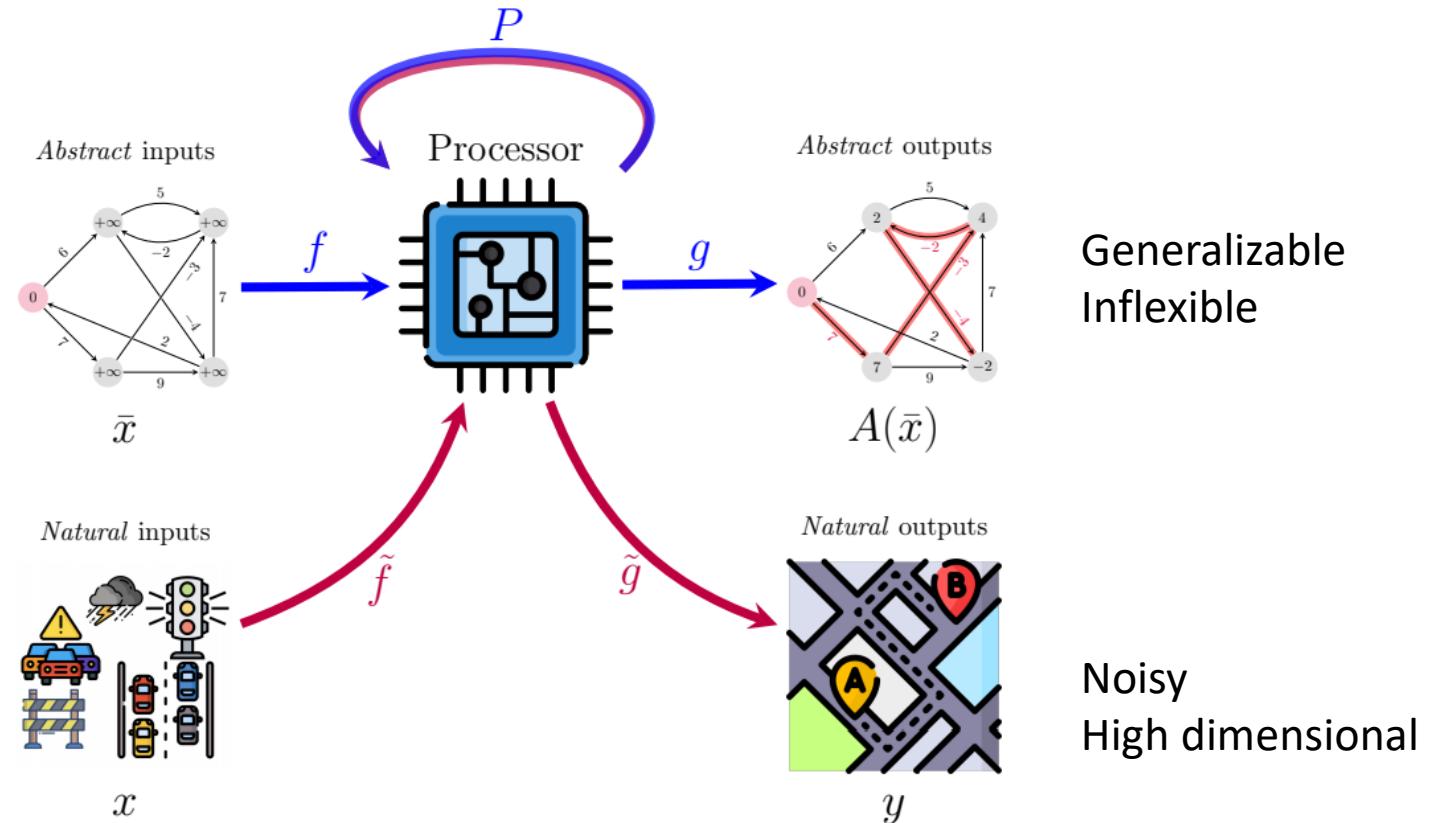
# Part A: Sub-topics

- Reasoning as a prediction skill that can be learnt from data.
  - Question answering as zero-shot learning.
- Neural network operations for learning to reason:
  - Concept-object binding.
  - Attention & transformers.
  - Dynamic neural networks, conditional computation & differentiable programming.
- Reasoning as iterative representation refinement & query-driven program synthesis and execution.
  - Compositional attention networks.
  - Reasoning as Neural module networks.
- **Combinatorics reasoning**

# Implement combinatorial algorithms with neural networks

- 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

Train neural processor P to imitate algorithm A



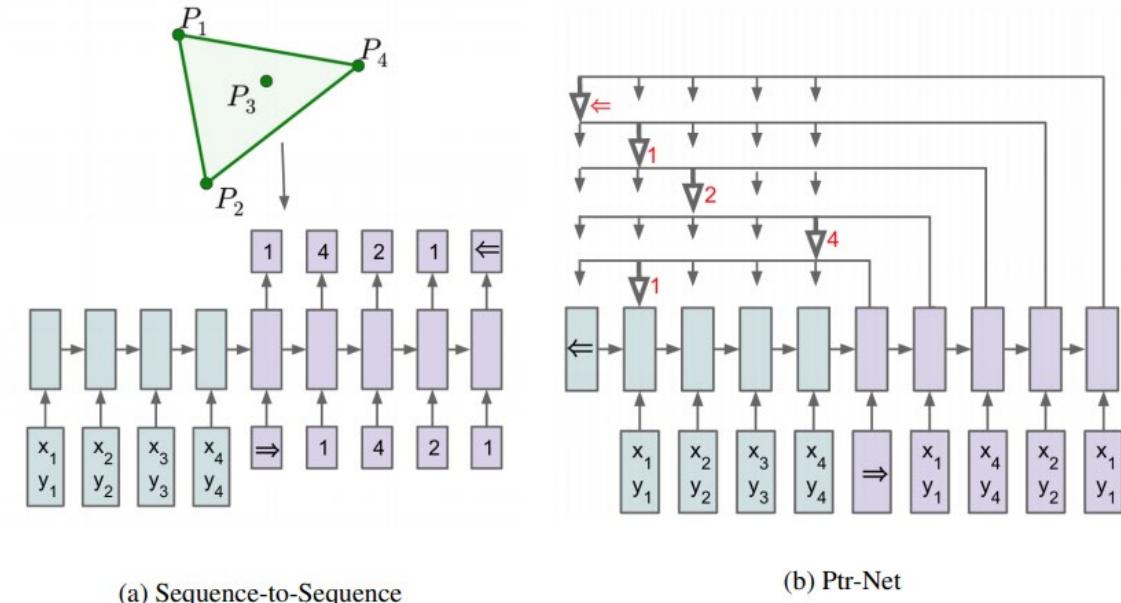
# Processor as RNN

- 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

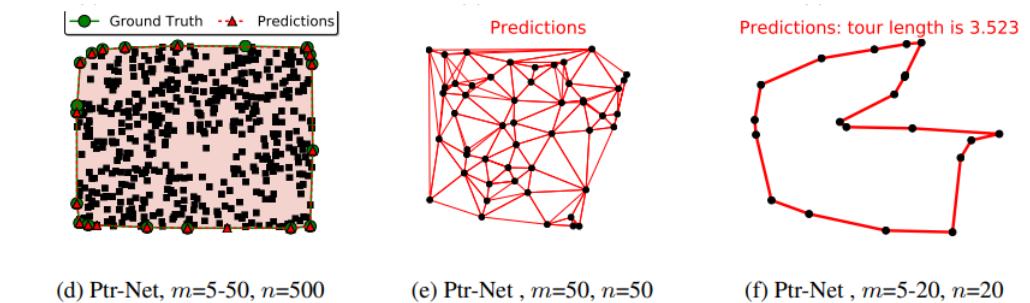
$$u_j^i = v^T \tanh(W_1 e_j + W_2 d_i) \quad j \in (1, \dots, n)$$

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



(a) Sequence-to-Sequence

(b) Ptr-Net

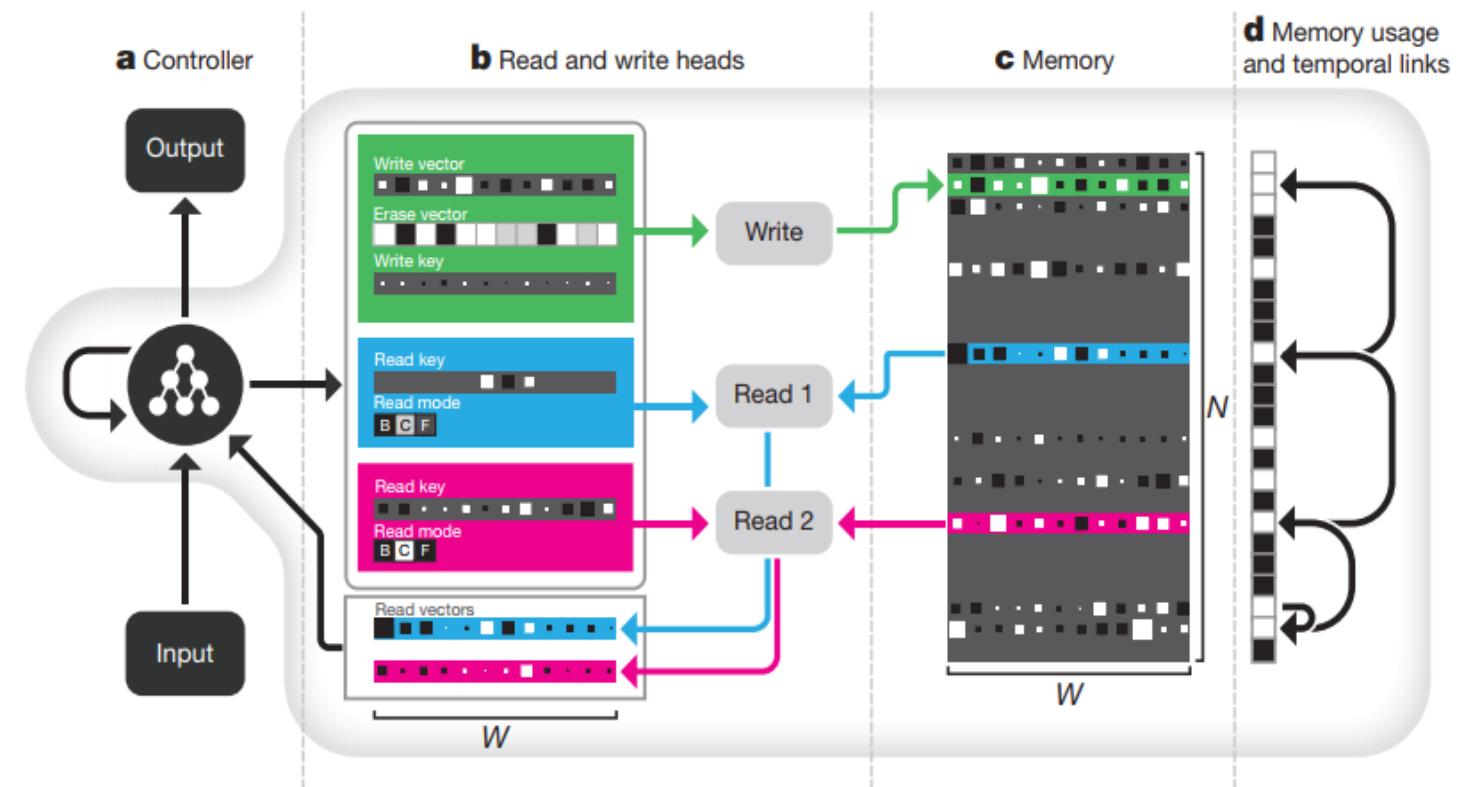


Vinyals, Oriol, Meire Fortunato, and Navdeep Jaitly. "Pointer networks."

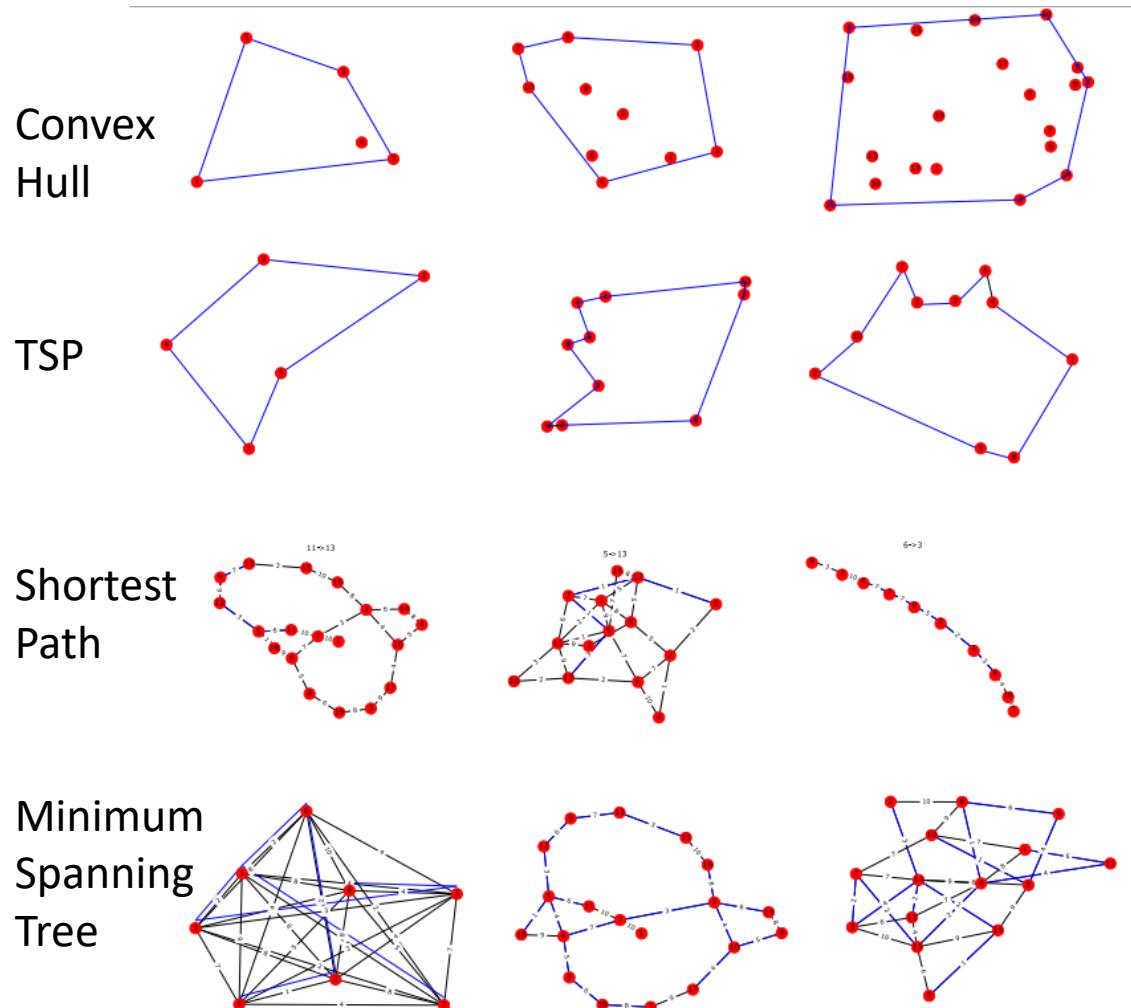
In *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 2*, pp. 2692-2700. 2015.

# 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



# Sequential encoding of graphs



- Each node is associated with random one-hot or binary features
- Output is the features of the solution

Geometry

[x<sub>1</sub>,y<sub>1</sub>, feature1],  
[x<sub>2</sub>,y<sub>2</sub>, feature2],  
...



[feature4],  
[feature2],  
...

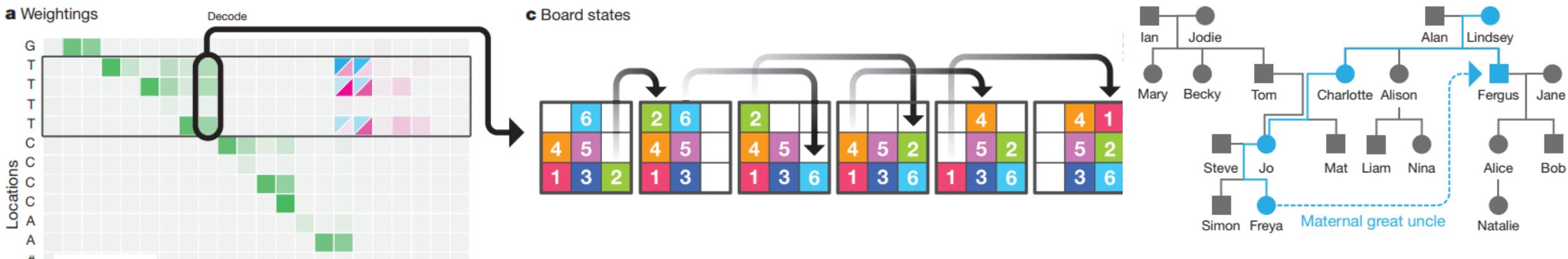
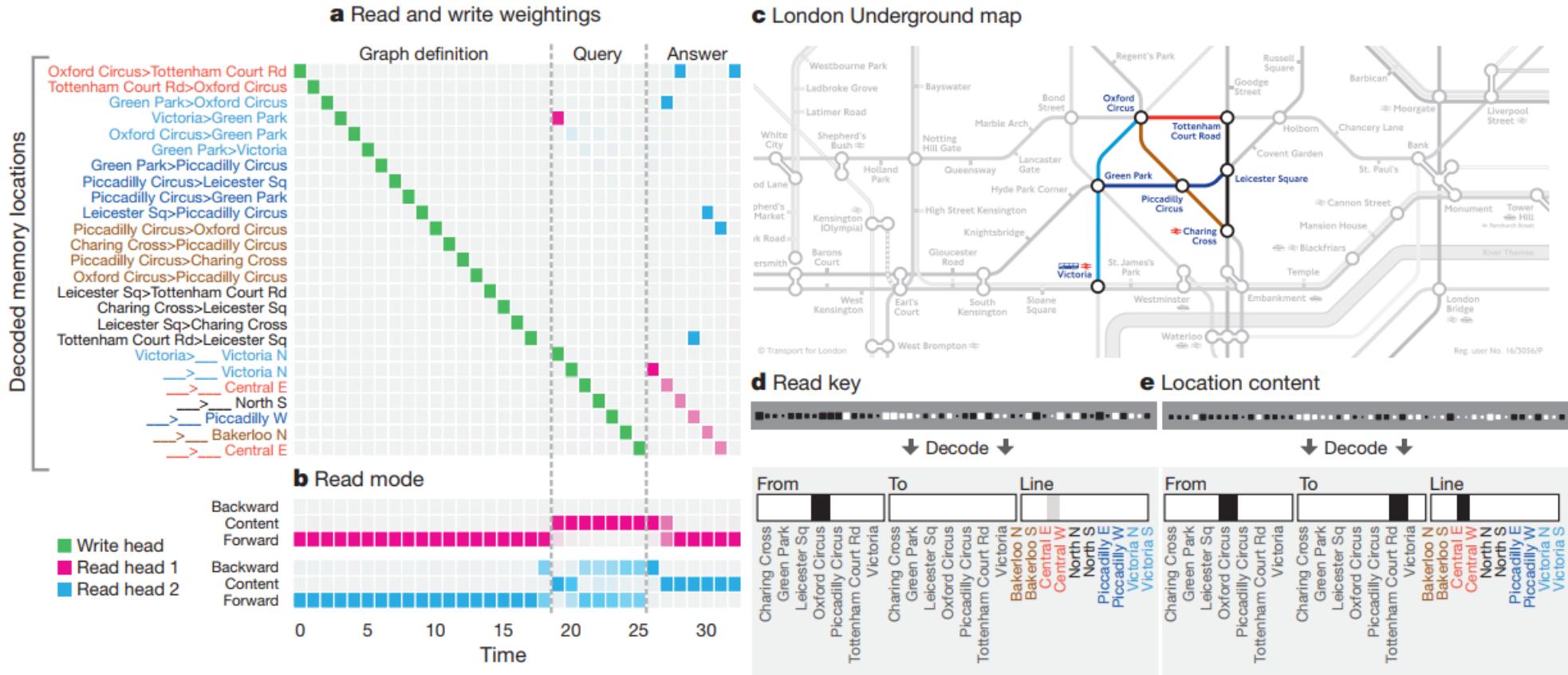
Graph

[node\_feature1, node\_feature2, edge12],  
[node\_feature1, node\_feature2, edge13],  
...

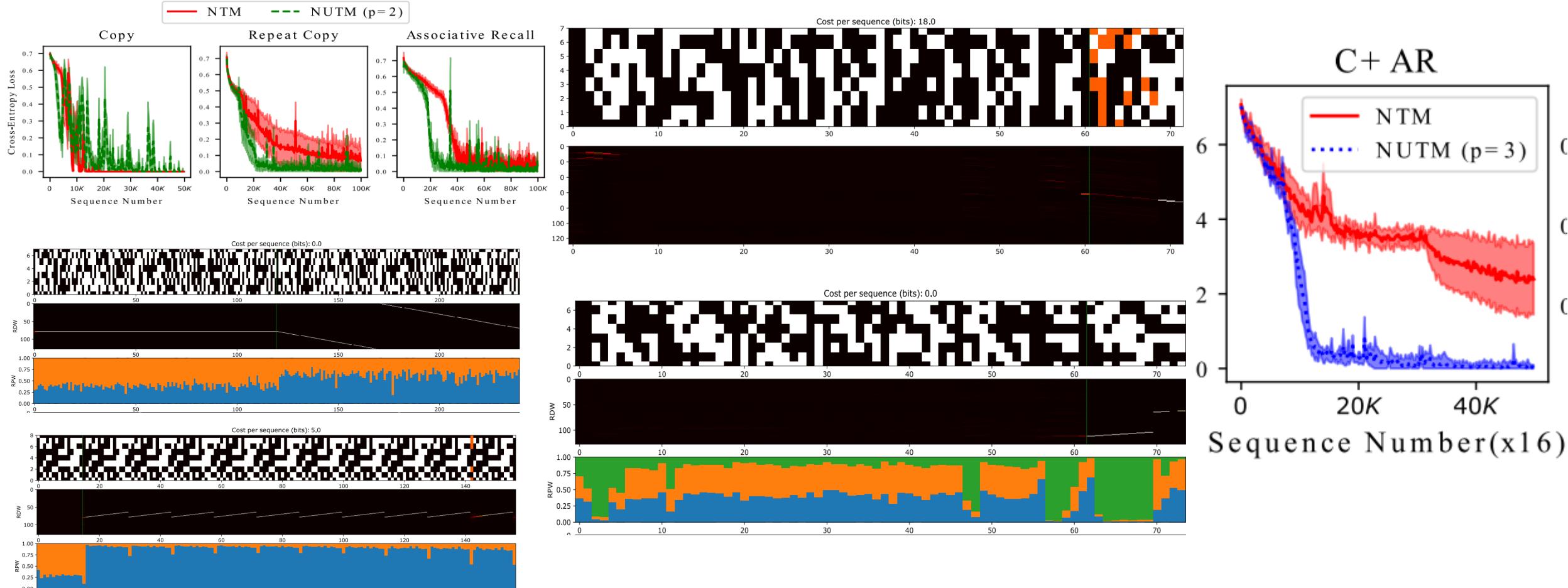


[node\_feature4],  
[node\_feature2],  
...

# DNC: graph reasoning

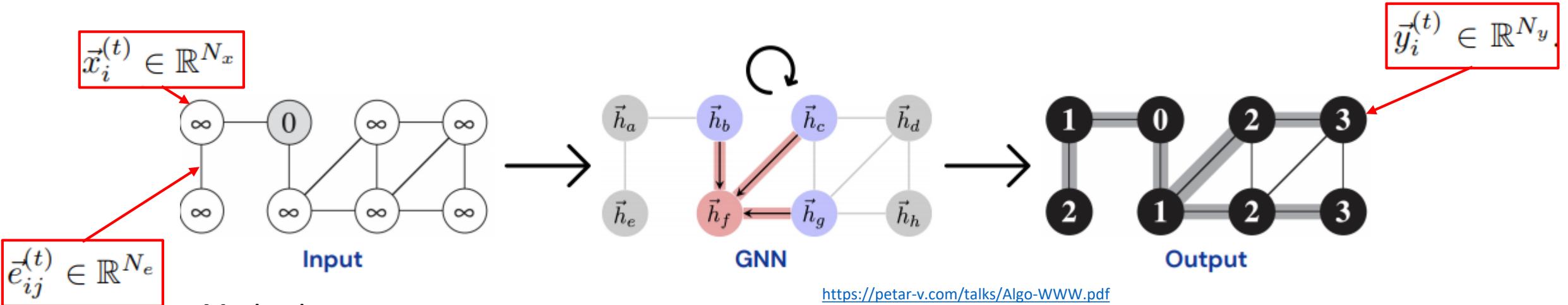


# NUTM: learning multiple algorithms at once



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

# Processor as graph neural network (GNN)



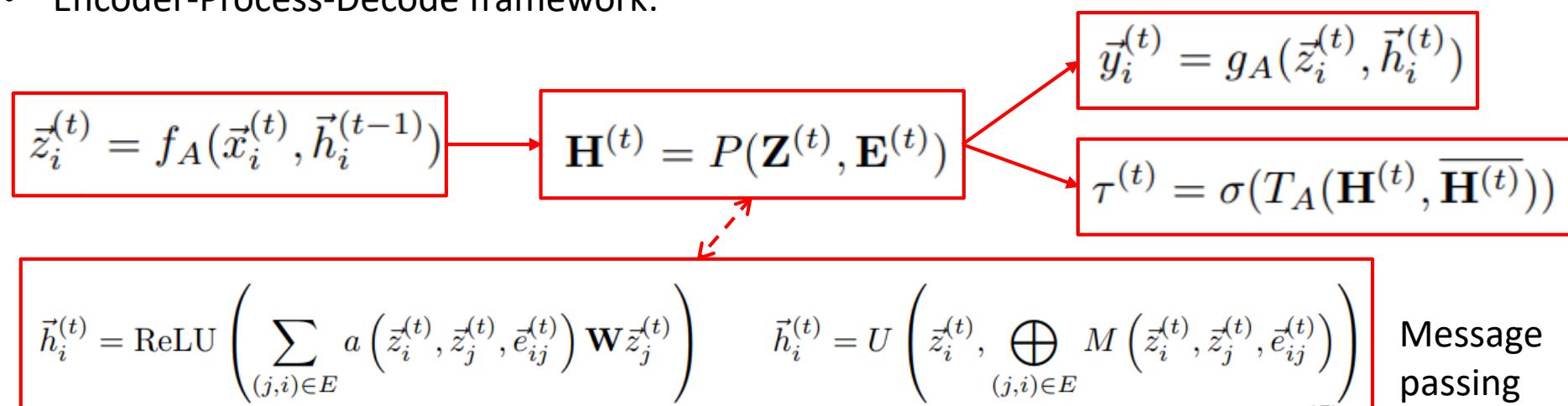
Motivation:

- 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 Padovano, 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 on output-level

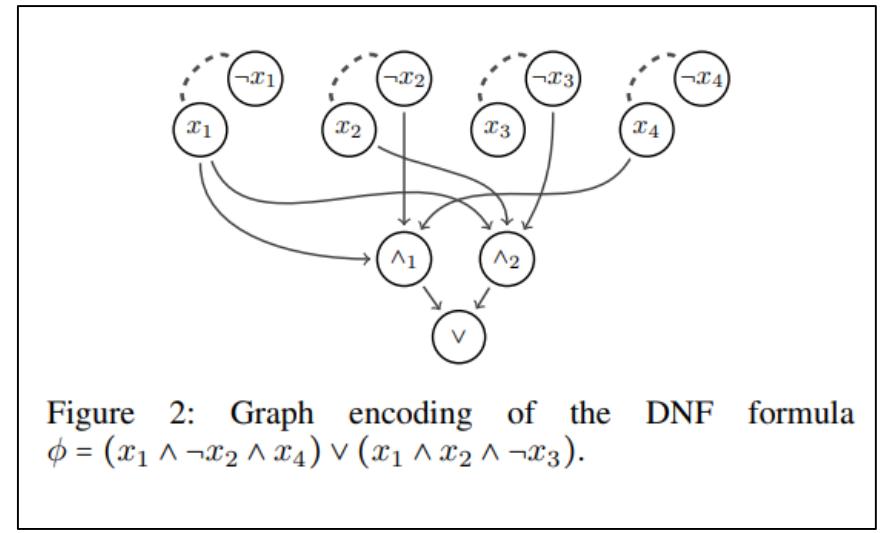
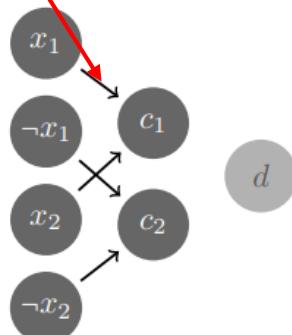


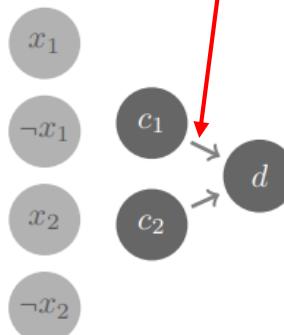
Figure 2: Graph encoding of the DNF formula  $\phi = (x_1 \wedge \neg x_2 \wedge x_4) \vee (x_1 \wedge x_2 \wedge \neg x_3)$ .

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



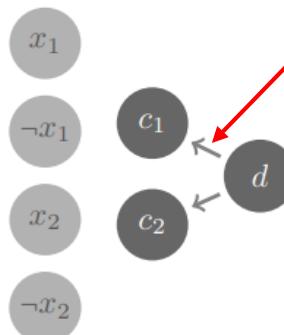
(a)

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



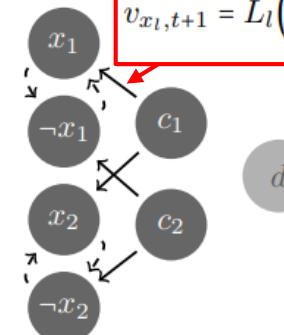
(b)

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



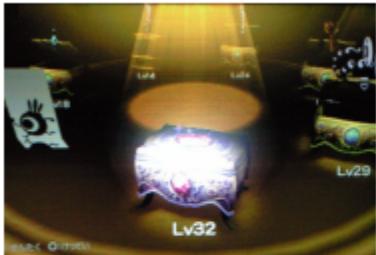
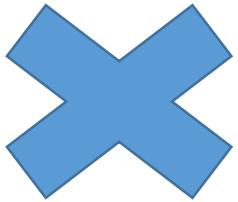
(c)

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



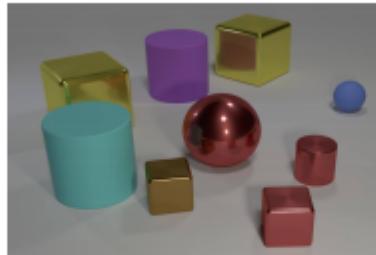
Best:  $O(m+n)$

# Neural networks and algorithms alignment



*Summary statistics*

What is the maximum value difference among treasures?



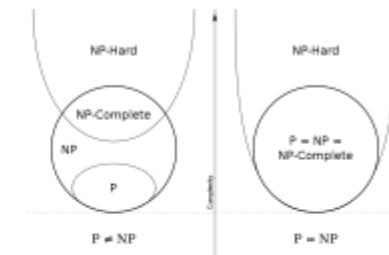
*Relational argmax*

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?

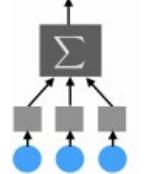


*NP-hard problem*

Subset sum: Is there a subset that sums to 0?



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

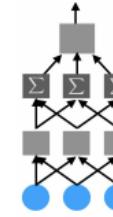


**Deep Sets** (Zaheer et al.)  
~ summary statistics

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



**GNNs**  
~ (pairwise) relations



$$\begin{aligned} h_s^{(k)} &= \sum_{t \in S} \text{MLP}_1^{(k)} \left( h_s^{(k-1)}, h_t^{(k-1)} \right) \\ y &= \text{MLP}_2 \left( \sum_{s \in S} h_s^{(K)} \right) \end{aligned}$$

Neural exhaustive search

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

# GNN is aligned with Dynamic Programming (DP)

## Graph Neural Network

```
for k = 1 ... 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

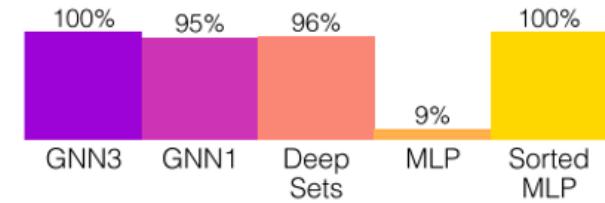
```
for k = 1 ... ISI - 1:  
  for u in S:  
     $d[k][u] = \min_v d[k-1][v] + \text{cost}(v, u)$ 
```

*Learns a simple reasoning step*

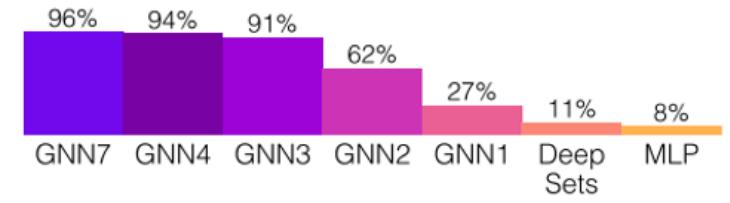
 $h_u^{(k)}$ 
 $d[k][u]$ 
 $\sum_v$ 
 $\min_v$ 
 $\text{MLP}(h_v^{(k-1)}, h_u^{(k-1)})$ 
 $d[k-1][v] + \text{cost}(v, u)$ 

Neural exhaustive search

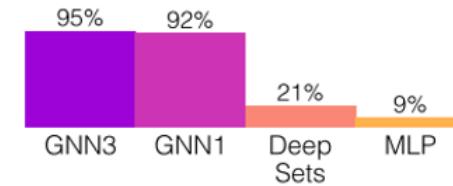
$\text{MLP}_2(\max_{\tau \subseteq S} \text{MLP}_1 \circ \text{LSTM}(X_1, \dots, X_{|\tau|} : X_1, \dots, X_{|\tau|} \in \tau)).$



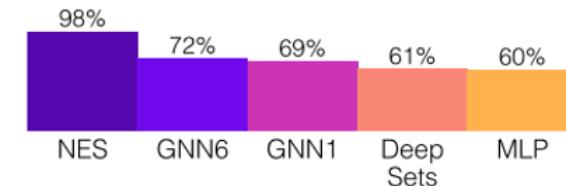
(a) Maximum value difference.



(c) Monster trainer.



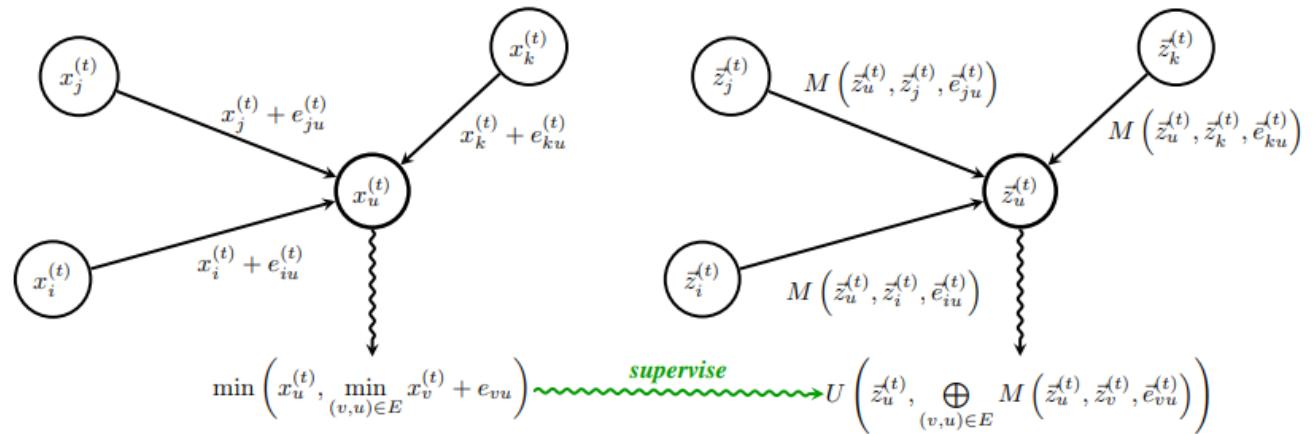
(b) Furthest pair.



(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

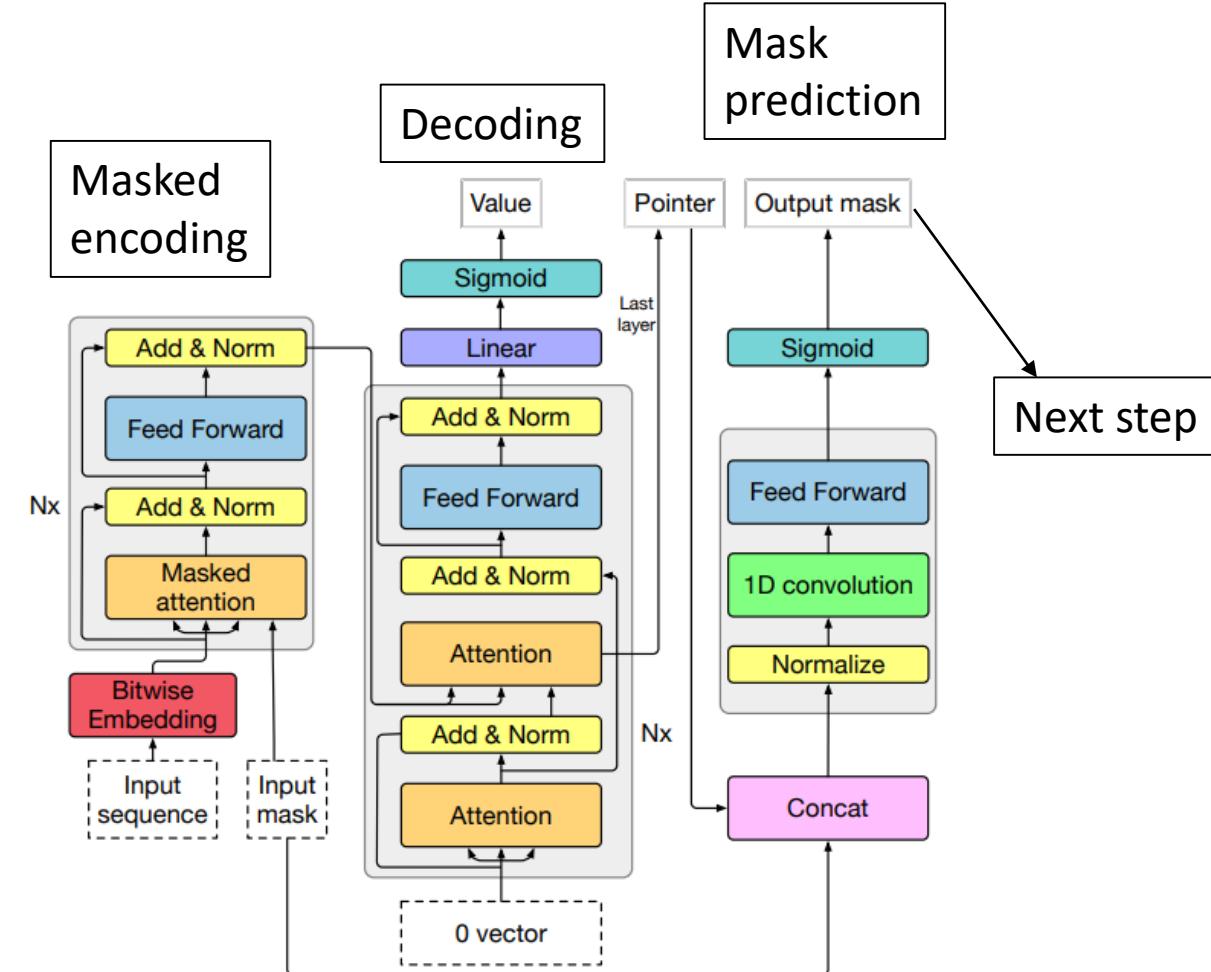


Algorithm	Inputs	Supervision signals
Breadth-first search	$x_i^{(t)}$ : is $i$ reachable from $s$ in $\leq t$ hops?	$x_i^{(t+1)}$ , $\tau^{(t)}$ : has the algorithm terminated?
Bellman-Ford	$x_i^{(t)}$ : shortest distance from $s$ to $i$ (using $\leq t$ hops)	$x_i^{(t+1)}$ , $\tau^{(t)}$ , $p_i^{(t)}$ : 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

# 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)

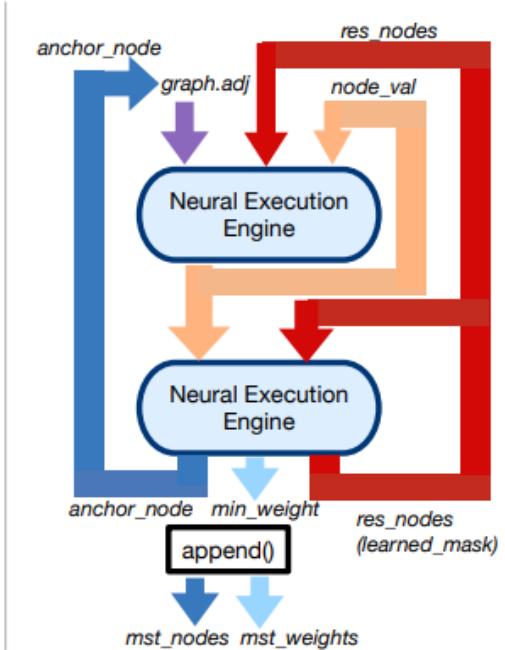
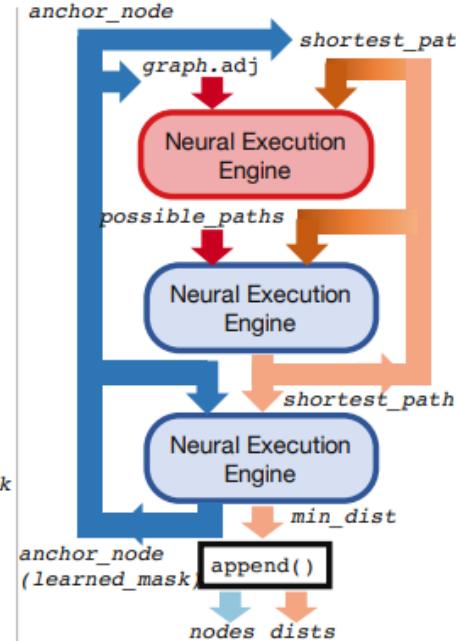
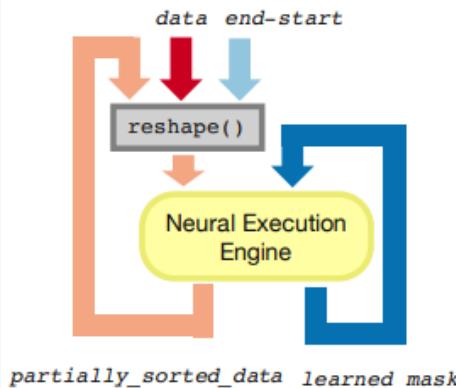
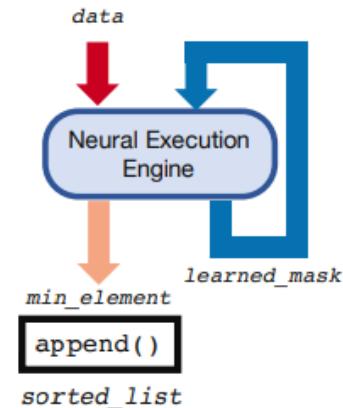
        node_val(res_nodes) = min(node_val(res_nodes),
                                 adj_list(res_nodes))

        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

```



End of part A

<https://bit.ly/37DYQn7>