

# What if Neural Networks had explicit memory?

Batuhan Tömekçe

24.05.2021

# Current NNs fail at long sequences

- MLP can't generalize to arbitrary sequence length
- RNNs are conceptually powerful but have problems in training.
  - Gradient vanishing and explosion
- LSTMs solve gradient problems but hard to train and still don't scale well with long sequences

# What about learning algorithms?

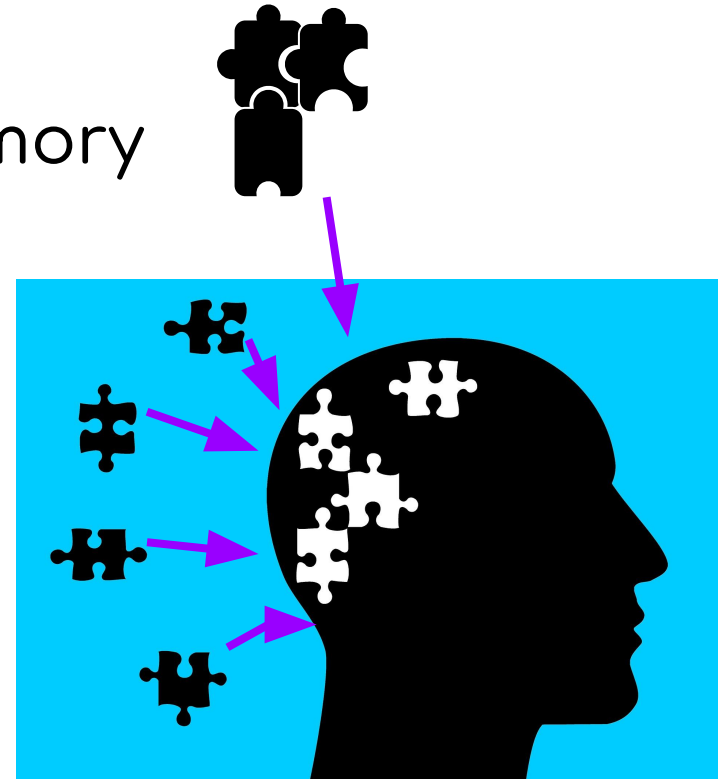
- Most of the use-cases of neural networks are regression or classification
- No functionality for rule based systems (branches, loops)
- The learnt functions aren't interpretable
- Memory component is only implicit in weights or a single state vector

# Outline

1. Inspirations from brains and computers
2. Where do neural networks come from?
3. Memory Networks - [extending neural networks with its own knowledge base](#)
4. Memory-Augmented Neural Network - [how to build one?](#)
5. Building a neural computer through MANNs
6. Discussion

# Brain uses both short-term and long-term memory

- Humans use their short term memory with chunks
- Global workspace theory - retrieve necessary knowledge per decision



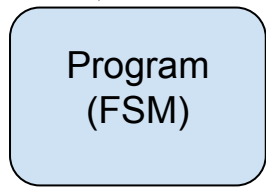
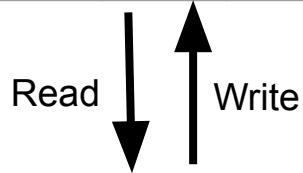
Source: [www.pixabay.com](http://www.pixabay.com)

# FSM with memory makes a Turing Machine

A Turing Machine can be represented as a 4-tuple of states, symbols, program and initial state



$$(Q, \Gamma, \delta, q_0)$$



$$\delta : Q \times \Gamma \rightarrow \Gamma \times \{-1, 1\} \times Q$$

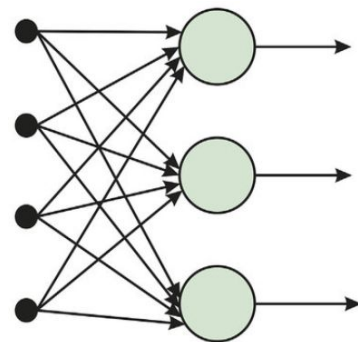
# Computers can implement algorithms

- Road to computers
  - Finite state machine to Turing Machine goes through Memory
- RNNs are turing complete
  - What does it mean?
- Universal Turing Machine
- Contrasting Turing Machines with statistical approach
  - MLPs are static function approximators

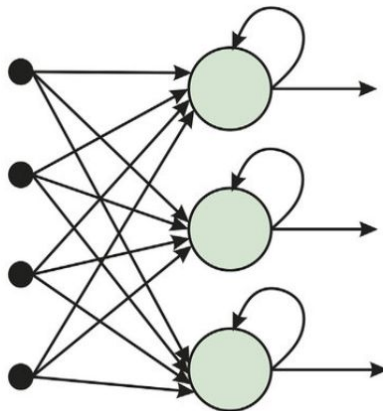
# Everything is connected and distributed

- Connectionism
  - a cognitive theory
- Parallel distributed processing
  - implementation of the theory

Feed-forward



Recurrent







Using Memory to store input knowledge

# Memory Networks equip neural networks with custom knowledge base

## 4 key ingredients

- **I** → input feature map
  - e.g. word to vector
- **G** → generalization
  - update memory
- **O** → output feature map
  - produce output
- **R** → response
  - e.g. output to word

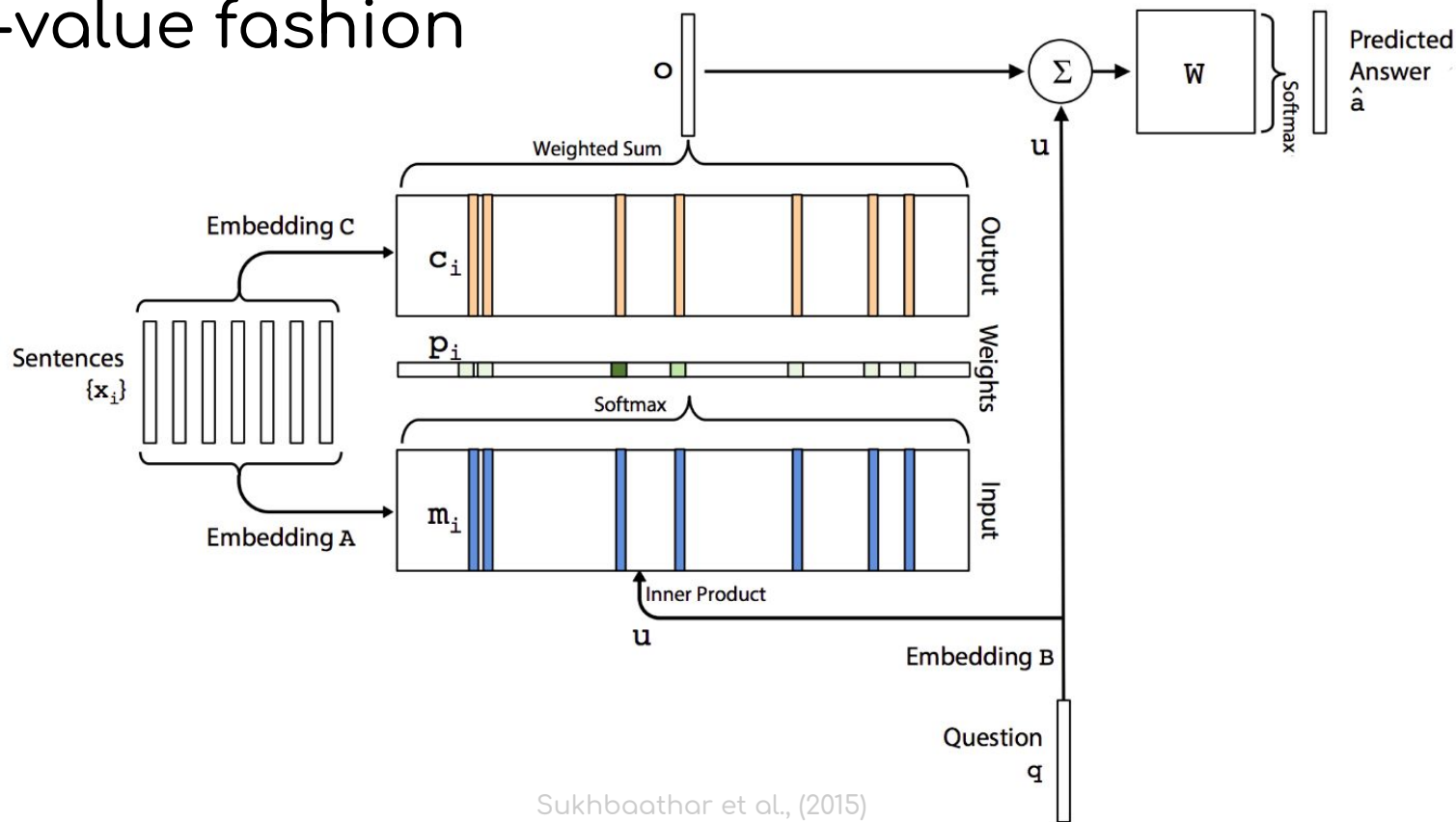
# Memory Networks equip neural networks with custom knowledge base

- Main goal is to reason about multiple dependent connections
  - e.g. question answering, reasoning about connections between sentences
- Stores incoming knowledge
- Retrieves sequentially

# Neurons fire together wire together

- End-to-end training improves memory networks
- Used in language modeling and questions answering
- Multiple hops is important for reasoning tasks

# The knowledge is encoded into memory in key-value fashion

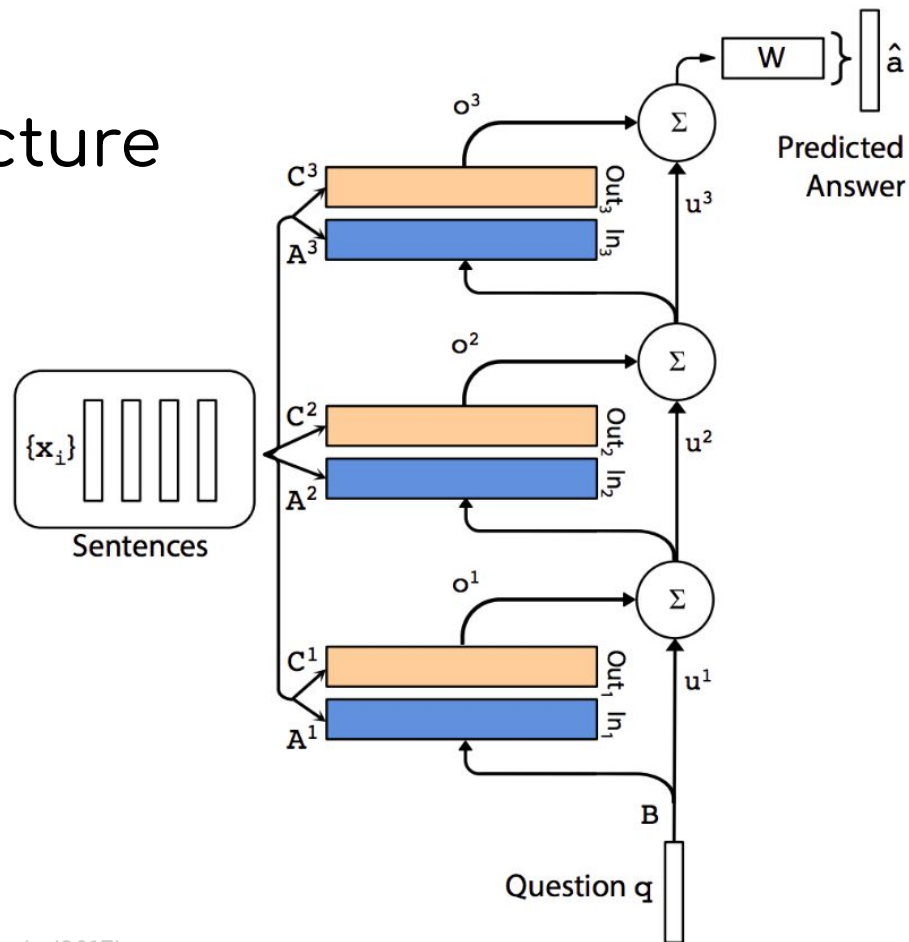


# Multiple hops enable reasoning in a chain structure

Sentences are e.g. Wikipedia

Who is presenting **neural network architectures for algorithms in Deep Neural Networks seminar** at **ETH**?

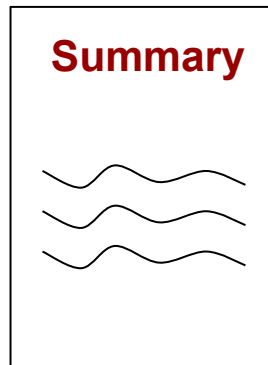
The question will be answered in hops



# Memory Networks have static memory

- Mostly used in real datasets (language, vision)
- Mostly fixes the memory at test time
- Mostly big memory size and works on discrete sets, multiple hops (multiple memory)

Memory  
Networks



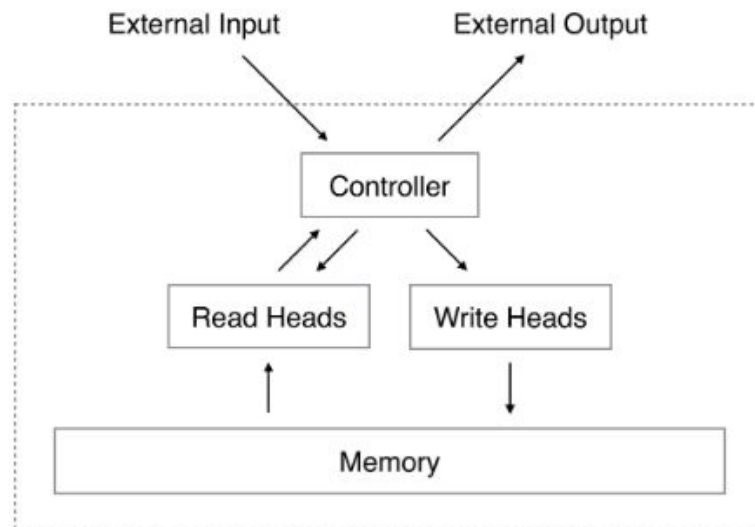


Enhancing a neural network with  
memory to solve algorithms

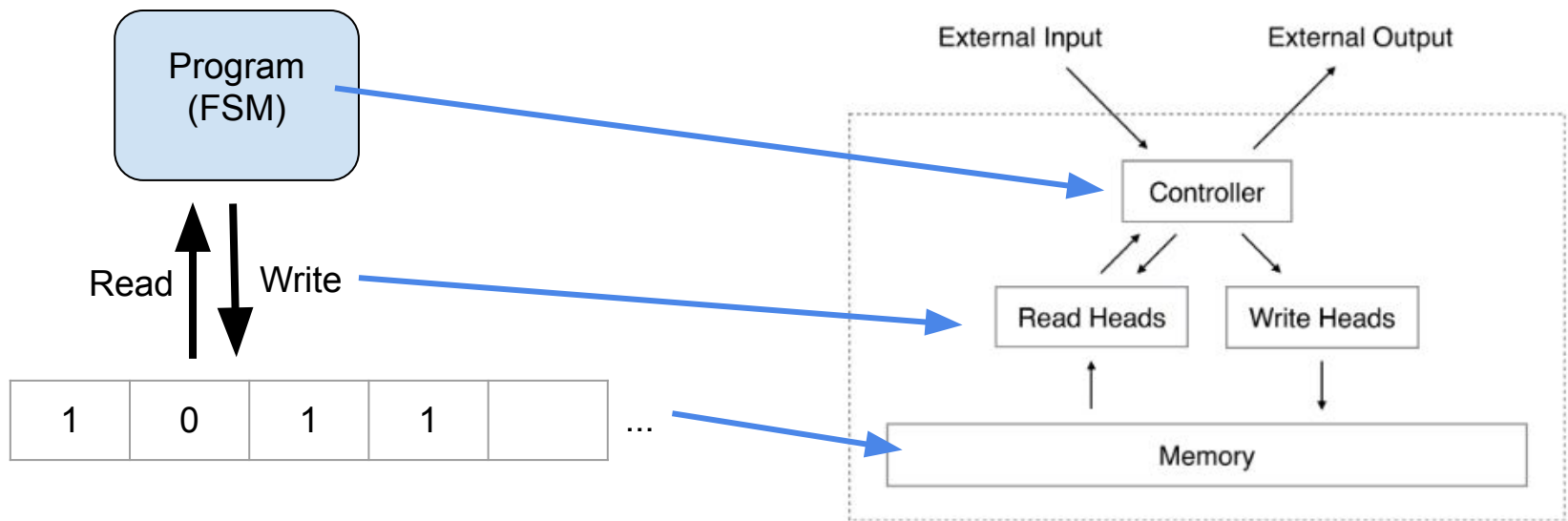


# Neural Turing Machine augments a neural network with memory

- Controller consists of neural network(s)
- First read then write
- Main method of communication is attention



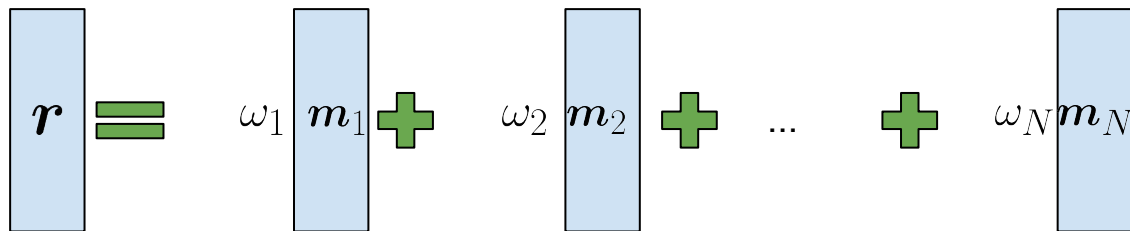
# Why it is called Neural Turing Machine?



# Read by convex combination of the memory cells

One read head at each time step computes

$$\mathbf{r}_t = \mathbf{M}_t^T \boldsymbol{\omega}_t$$



Read head generates the normalized weight vector

$$\sum \omega_i = 1$$

$$\mathbf{M} =$$

Memory bank

$m_1^T$
$m_2^T$
...
$m_N^T$

# Write to memory by erasing and adding

Write head generates erase, add and weight vectors

$$e_t \quad a_t \quad \omega_t$$

The  $i$ 'th memory cell is erased by

$$\tilde{m}_{t,i} = m_{t-1,i} - \omega_i m_{t-1,i} \odot e_t$$

The  $i$ 'th memory is added

$$m_{t,i} = \tilde{m}_{t,i} + \omega_i a_t \quad M =$$

Memory bank

$m_1^T$
$m_2^T$
...
$m_N^T$

# Content-based addressing is used to communicate with memory

A weighted softmax distribution is used for content-based addressing using a key vector

$$\omega_t^c = \frac{\exp \beta_t \mathcal{S}(\mathbf{k}_t, \mathbf{m}_{t,i})}{\sum_i^N \exp \beta_t \mathcal{S}(\mathbf{k}_t, \mathbf{m}_{t,i})}$$

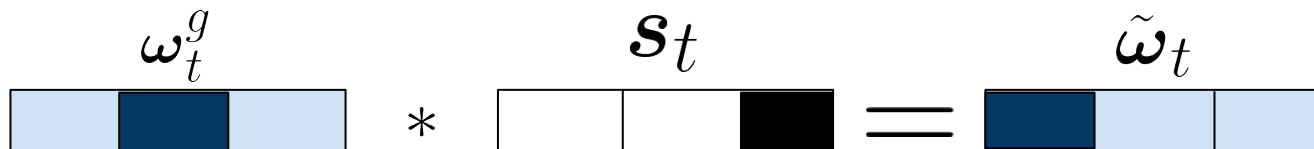
$$\omega_t^c = \frac{\exp \beta_t \mathcal{S}(\mathbf{k}_t, \mathbf{m}_{t,i})}{\sum_i^N \exp \beta_t \mathcal{S}(\mathbf{k}_t, \mathbf{m}_{t,i})}$$

# Location-based addressing is used for variable binding

First a convex combination of previous and content-based weight is taken

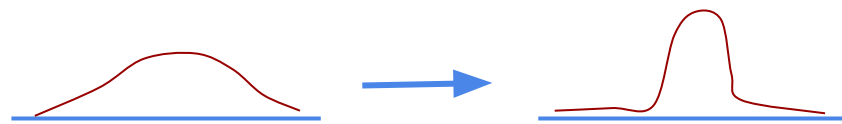
$$\omega_t^g \leftarrow (1 - g_t) \omega_t^c + g_t \omega_{t-1} \quad g_t \in (0, 1)$$

Then a circular convolution is taken

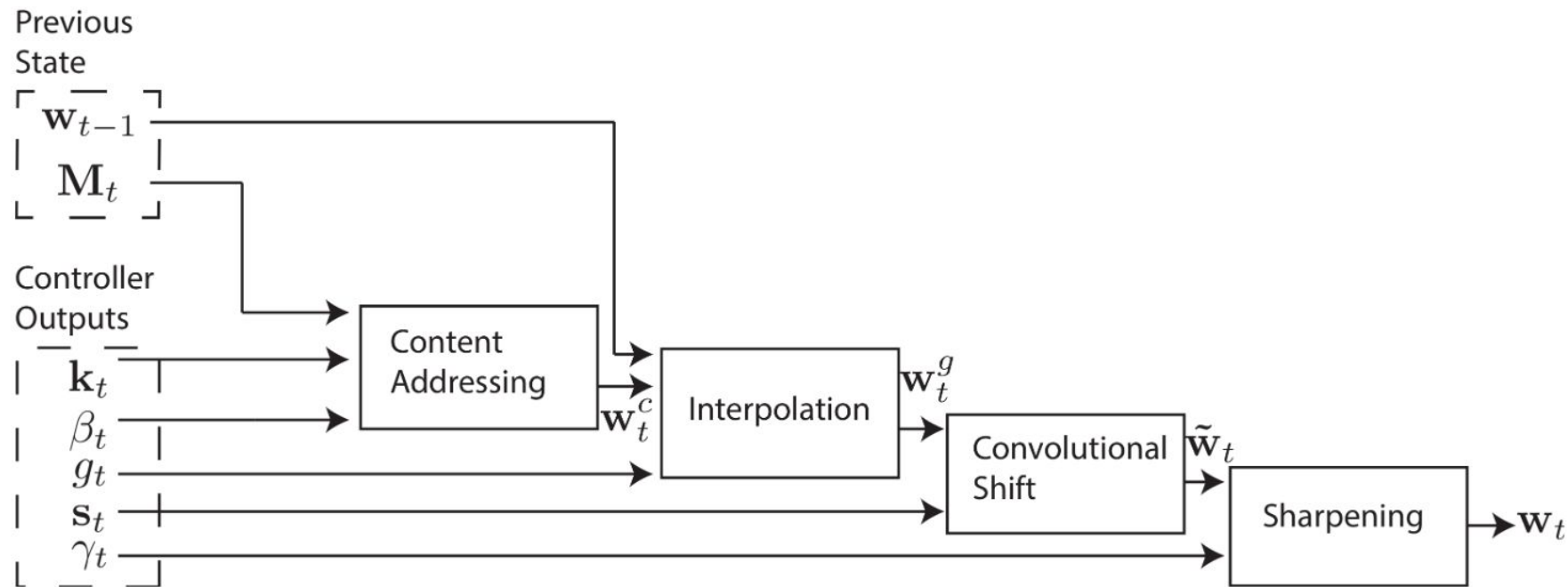


Then sharpen the distribution against accumulation errors

$$\omega_t = \text{softmax}(\gamma_t \tilde{\omega}_t)$$

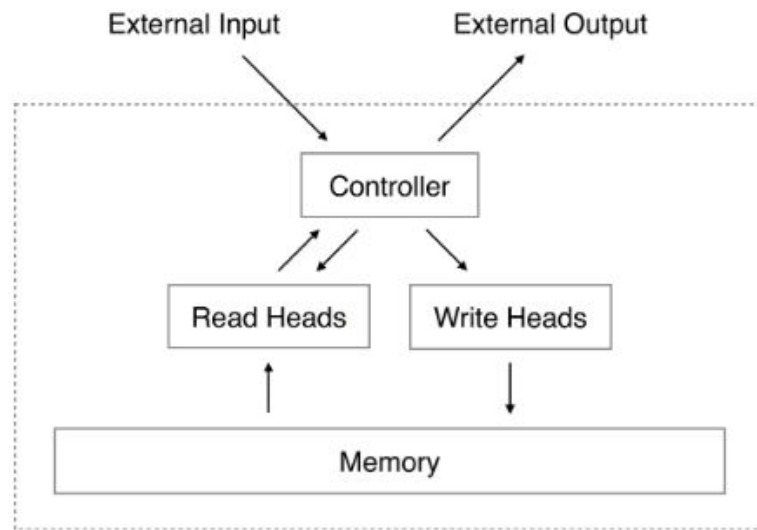


# Content based and location based addressing are employed together



# Controller consists of interface and state networks

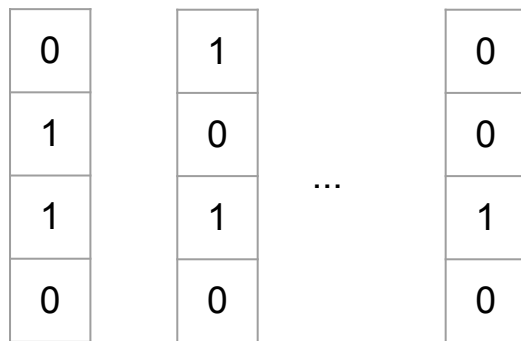
- State network is LSTM
- Interface network is MLP





# Evaluation is done through algorithmic rule based tasks that measure generalization

- Copy
- Repeat Copy

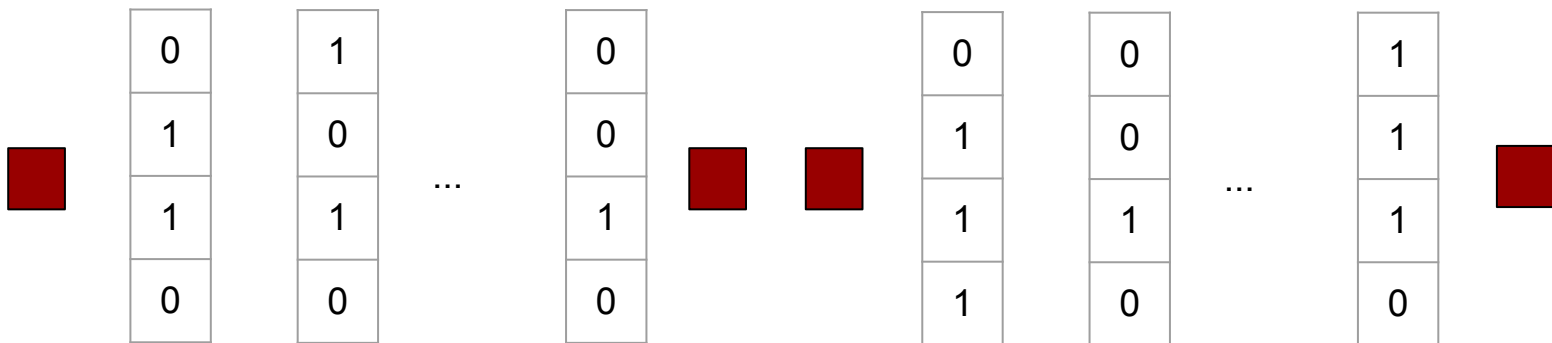


Sequence of binary vectors

# Evaluation is done through algorithmic rule based tasks that measure generalization

## Associative Recall

- Given list of items and a query item from the list, the model predicts the item next to the query



# Evaluation is done through algorithmic rule based tasks that measure generalization

## Dynamic N-gram

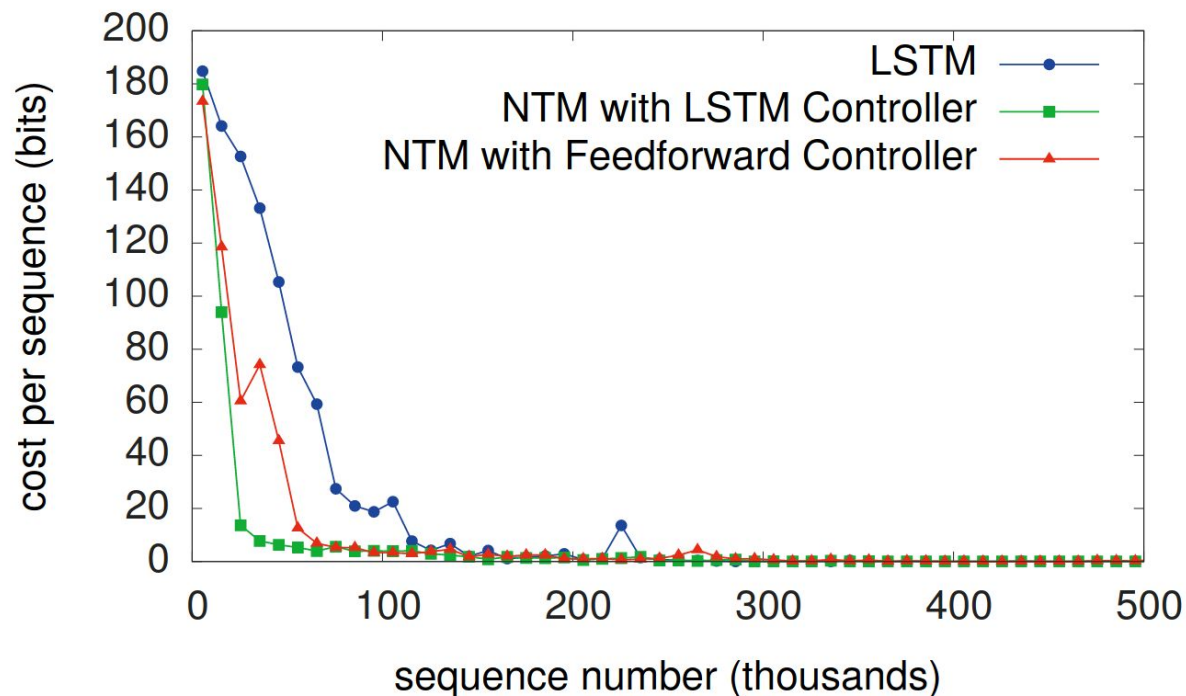
- Learn a distribution using the memory

## Priority Sorting

- Given vectors with preferences, sort them according to their preference

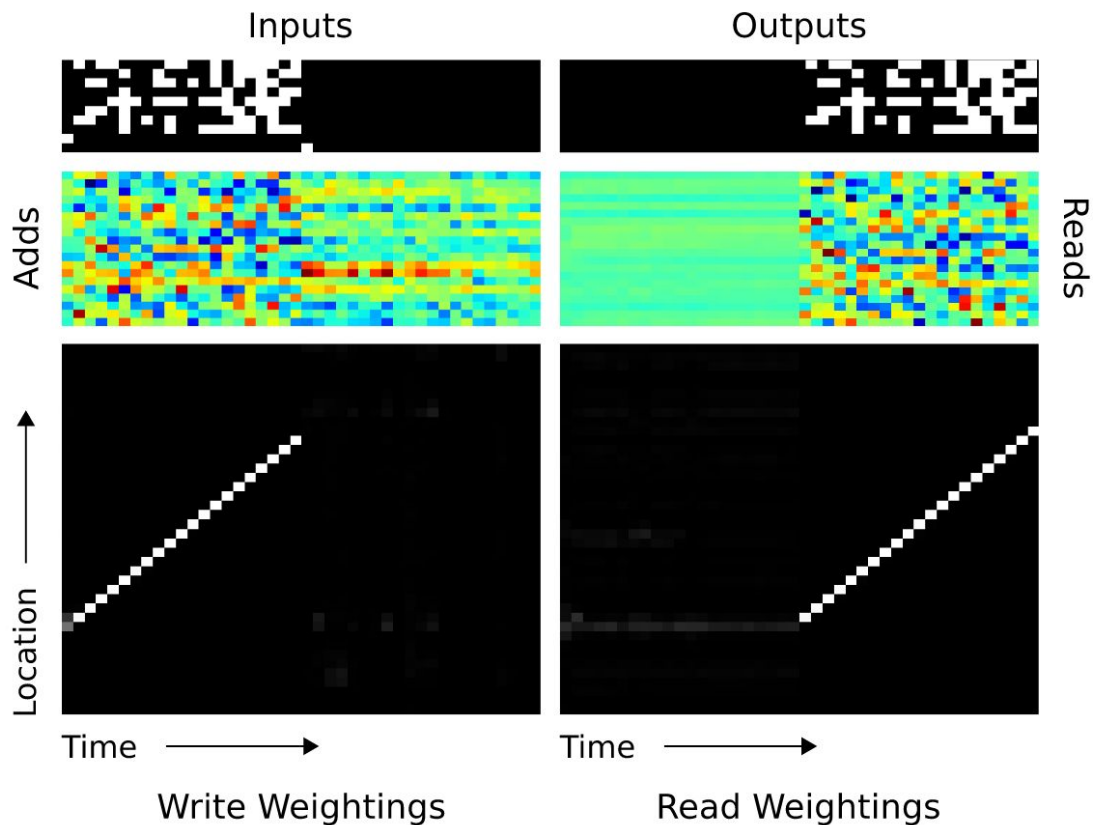
# NTM has more representational power and performs better

- First the input sequence is fed
- Then the the models produce output
- Multi-label binary classification
- The weights are resetted after each sequence



# NTM learns an algorithm

The sequence is stored in memory and then read from memory



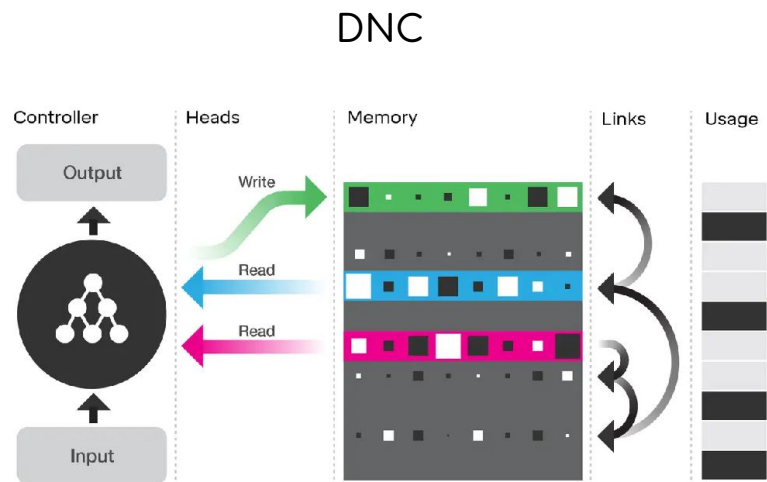
# Other examples of Memory-Augmented Neural Networks specialize for different task domains

[Differentiable Neural Computer](#) also stores the order of memory writes as a linked list

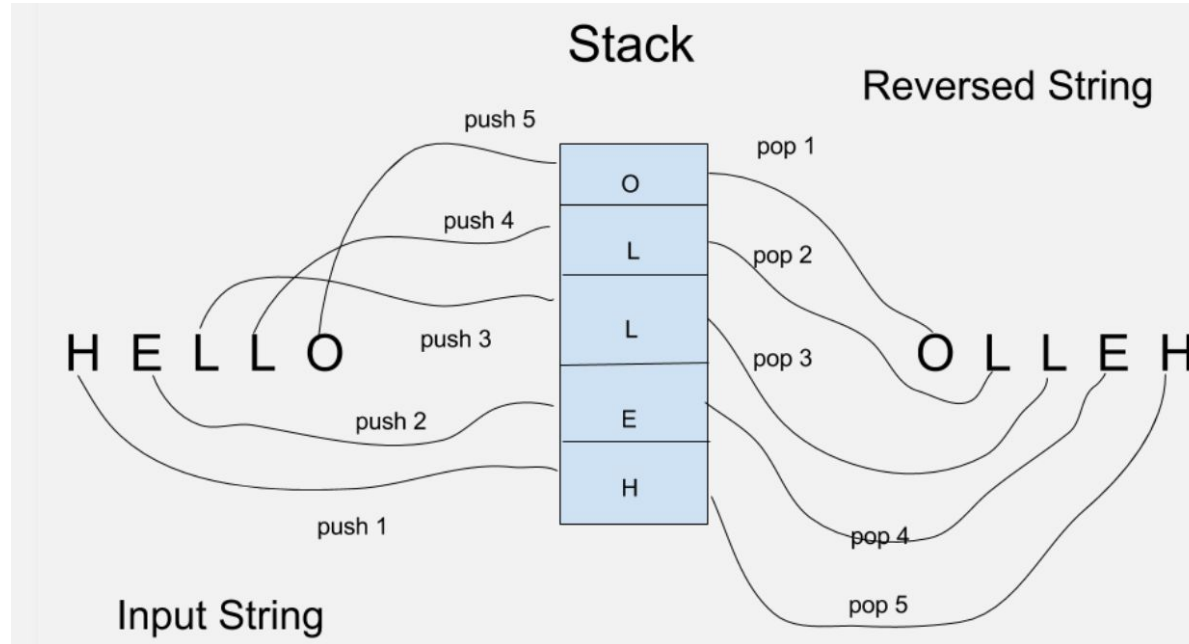
- used for complex data structures

[Least Recently Used Access](#) employs a content based addressing as the name suggests

- used in few-shot meta learning



# What about modifications to the memory bank?



Source: <https://jasdeep04.github.io/posts/Neural-Stacks/>

# Memory Networks have static memory whereas MANNs have dynamic memory

## Memory Networks

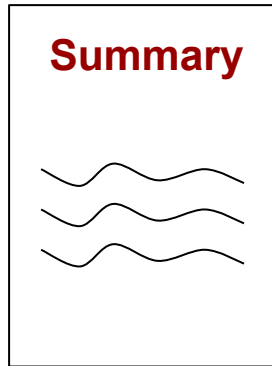
- Mostly used in real datasets (language, vision)
- Mostly fixes the memory at test time
- Mostly bigger memory size and works on discrete sets, multiple hops (multiple memory)

## Memory-Augmented Neural Networks

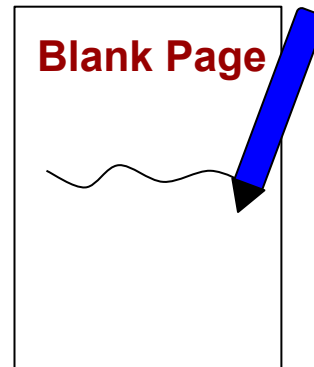
- Mostly used in simulated algorithmic tasks
- Use its memory to store objects at test time for algorithmic purposes
- Smaller external memory (different version are possible) provides variable binding



Memory  
Networks



NTM





Using memory to store neural networks

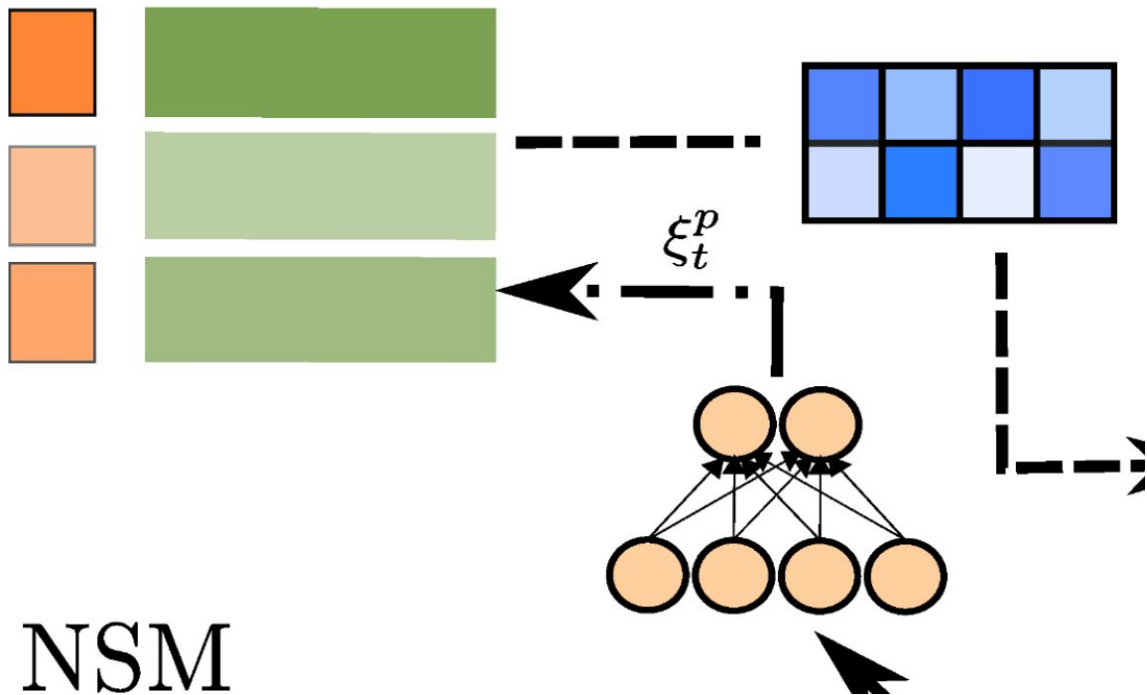
# MANNs can be improved with program memory

- Computers store different programs and also the data in RAM
- Selection of different program/models
- Different programs can be used for meta learning and multi-task learning
- Going in the direction of a neural computer

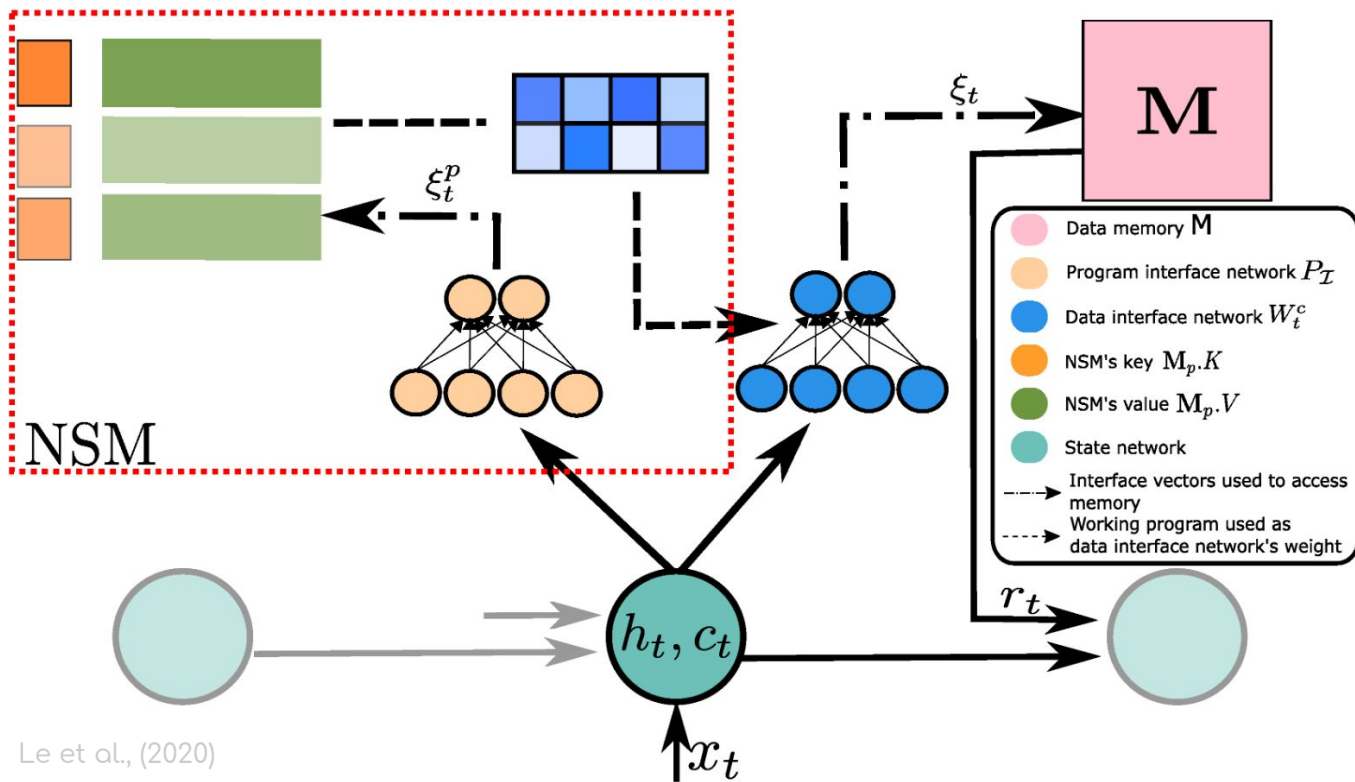
# Neural Stored Program Memory uses key-value attention to retrieve weights

A memory bank to store the weights of the controller

A meta network emits the keys for weight retrieval

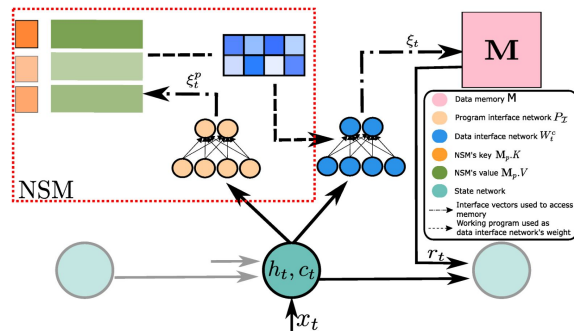


# Neural Universal Turing Machine is built by equipping a MANN with NSM

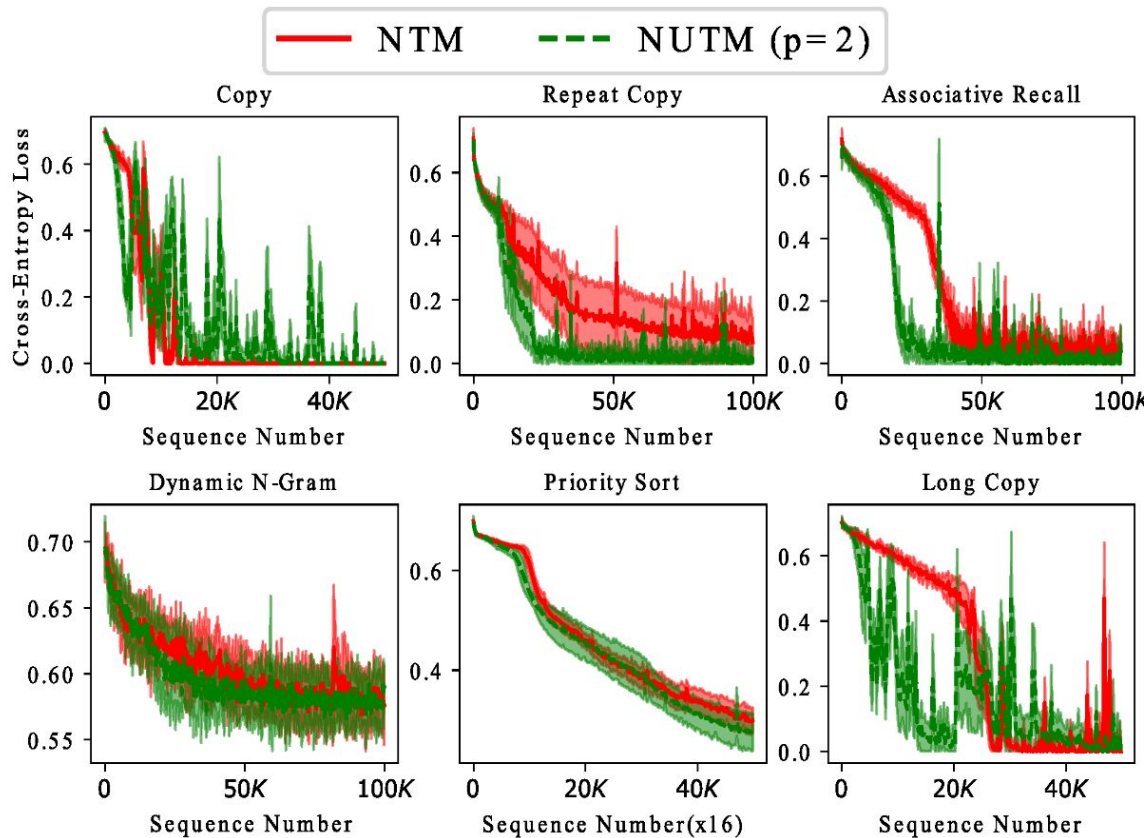


# How does NSM combined with NTM makes a Neural Universal Turing Machine?

- We can build Universal Turing Machines by putting the Turing Machines into the tape
- NSM only stores the weights of MLP interface
- one state Universal Turing Machine

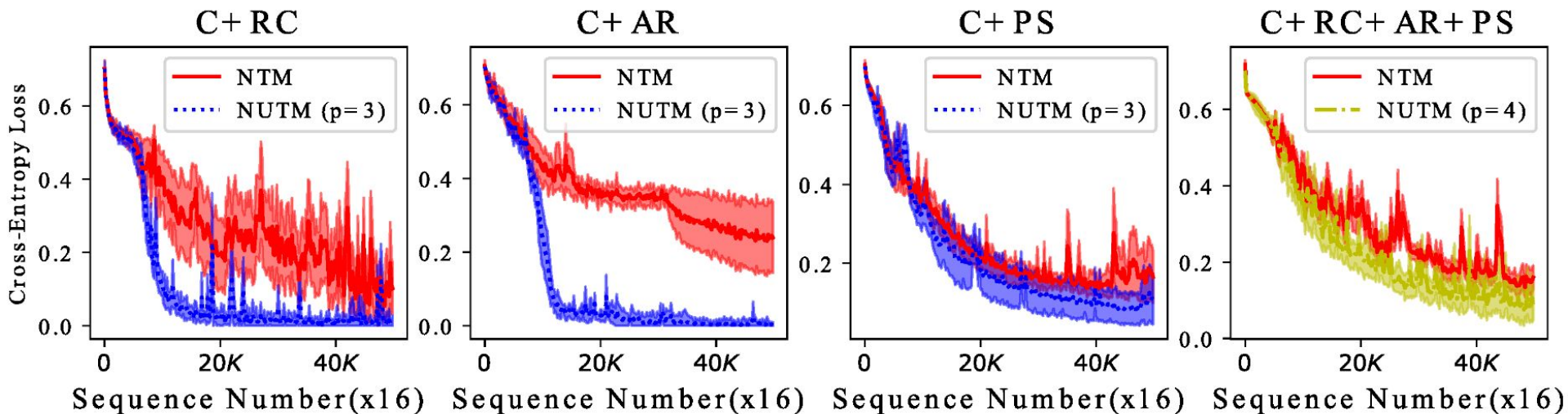


# NUTM converges faster than NTM?



# NUTM can learn different programs at once

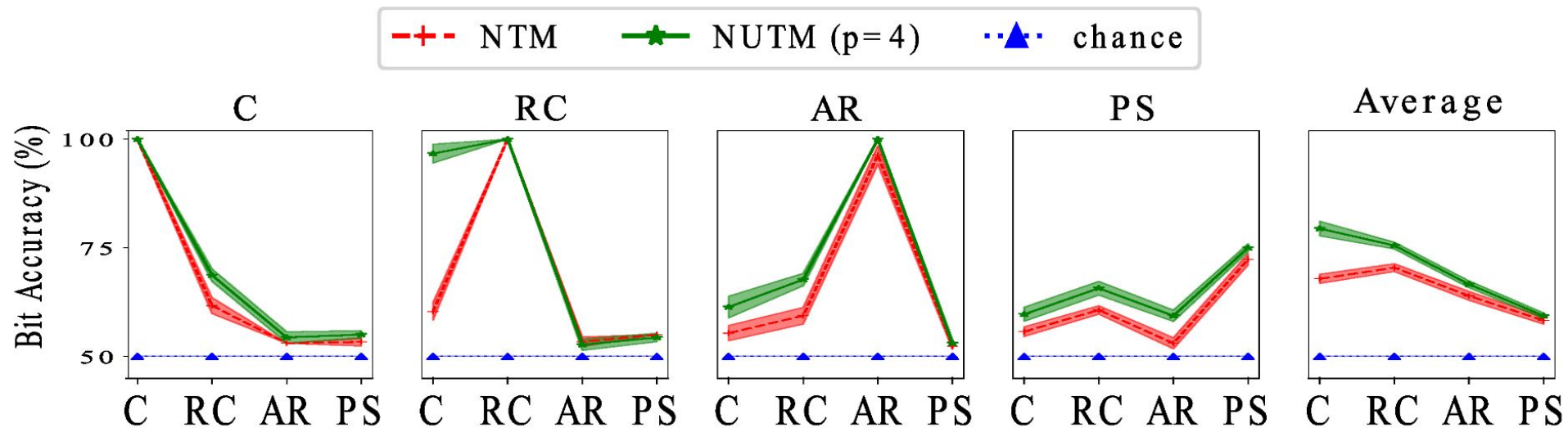
Combination of the atomic tasks after the other





# NUTM forgets less

How much of the task does the model remember?

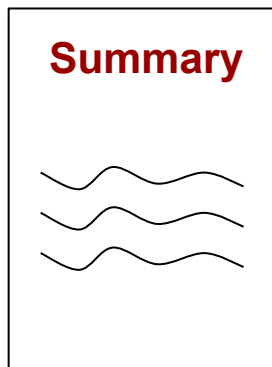


# NSM with MANN implements a form of fast and slow weights

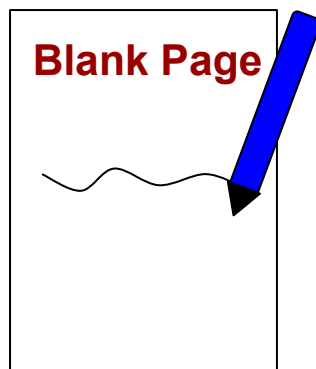
- A form of meta-learning
- slow weights are through backpropagation
- fast weights are through interpolation of programs

# Looking into the models as different exam types

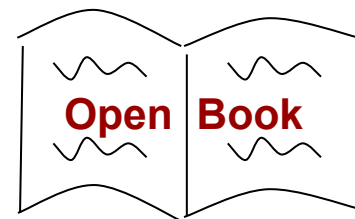
Memory  
Networks



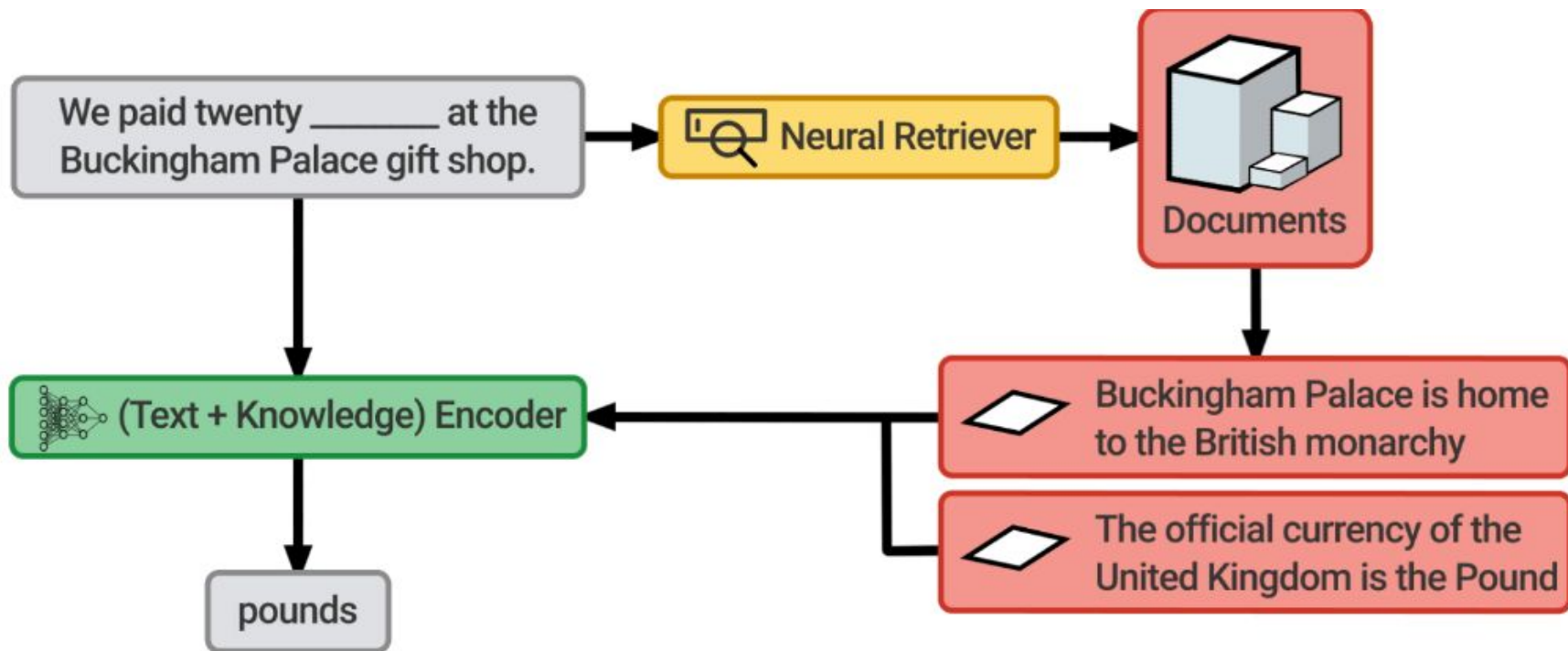
NTM



?



# REALM can cite while answering questions



Source: <https://ai.googleblog.com/2020/08/realms-integrating-retrieval-into.html>

# Recap

- Importance of memory and differentiability for current intelligence systems
- Memory networks
- Memory-augmented neural networks
- Meta learning perspective, storing neural networks in memory

# Discussion

- Memory in RL
- Neural stored-program memory for different application domains
- Learning algorithms with seq2seq + attention
- Trainability and reproducibility of MANNs
- What is the goal?



MLP is all you need

# References

- Bernhard J. Baars, "In the theatre of consciousness: Global workspace theory, a rigorous scientific theory of consciousness" (1997)
- Miller, G. A. "The magical number seven, plus or minus two: some limits on our capacity for processing information" (1956)
- Siegelman and Sontag, "On the computational power of neural nets" (1995)
- Weston et al., "Memory networks" (2014)
- Sukhbaatar et al., "End-to-end memory networks" (2015)
- Graves et al., "Neural Turing machine" (2014)
- Graves and Wayne et al., "Hybrid computing using a neural network with dynamic external memory" (2016)
- Santoro et al., "Meta-learning with memory-augmented neural networks" (2016)
- Grefenstette et al., "Learning to Transduce with Unbounded Memory" (2015)
- Le et al., "Neural stored-program memory" (2020)
- von der Magsburg, "The correlation theory of brain function" (1981)
- Guu et al., "REALM: Retrieval-Augmented Language Model Pre-Training" (2020)