# Advances in
# Neural Turing Machines

**Truyen Tran**
Deakin University

**CafeDSL, Aug 2018**

✉ truyen.tran@deakin.edu.au

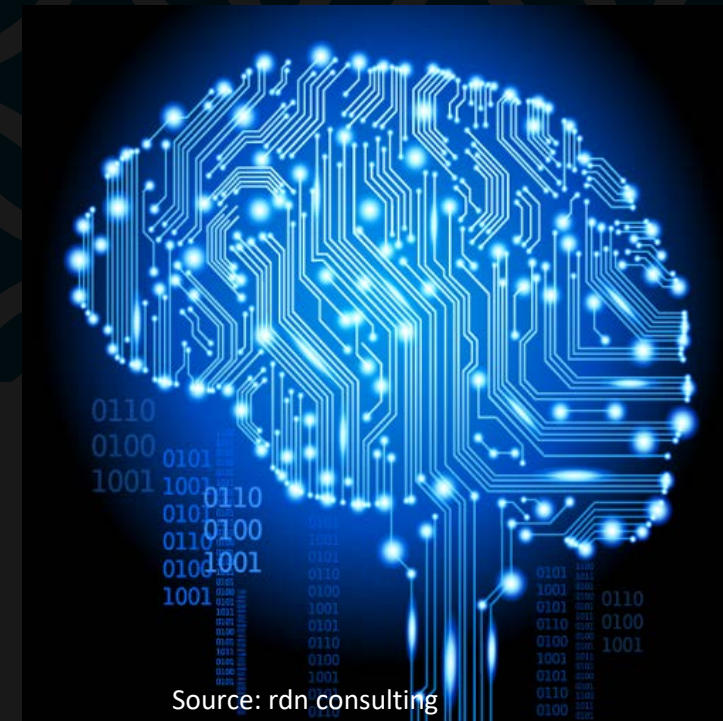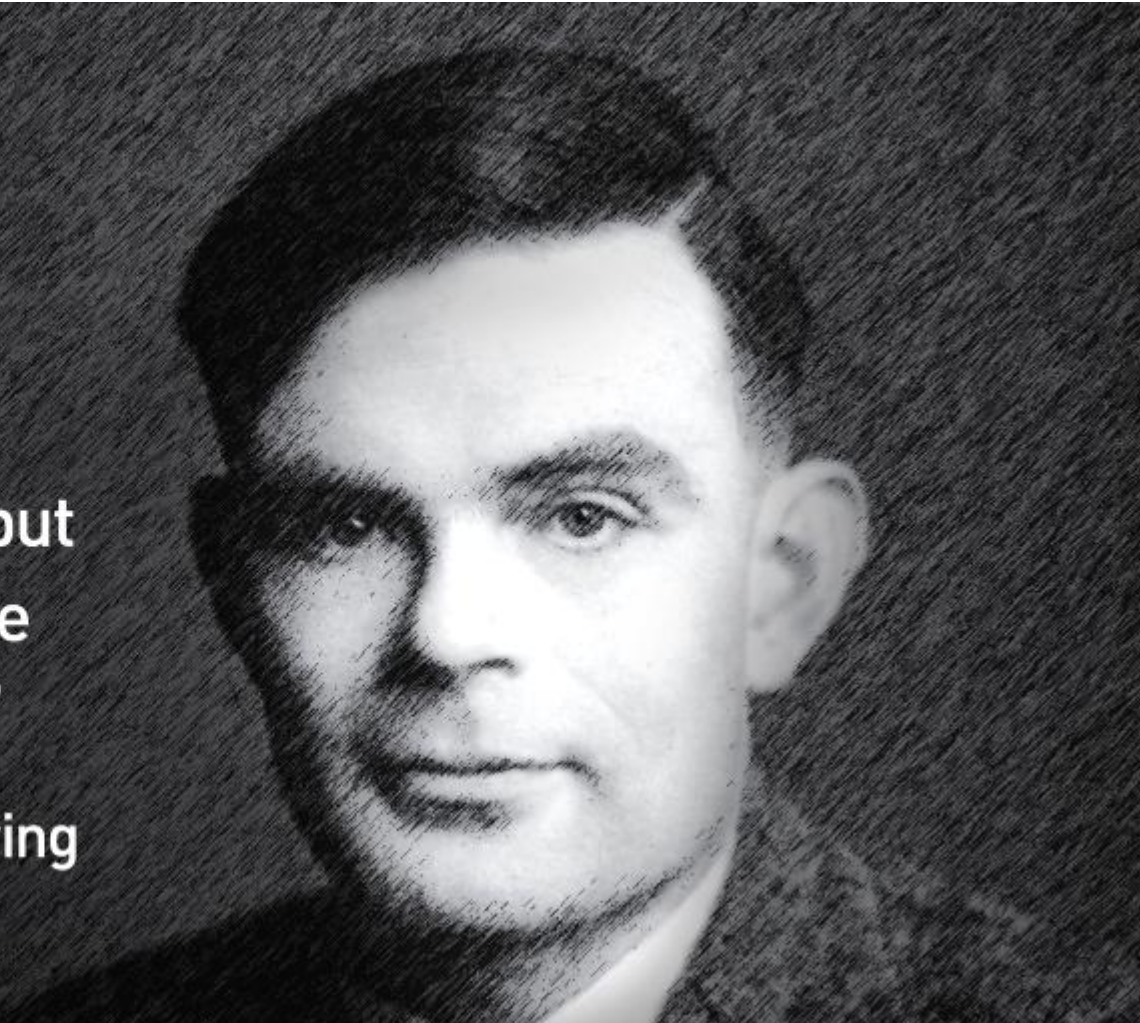🏠 truyentran.github.io

🐦 @truyenoz

📝 letdataspeak.blogspot.com

g⁺ goo.gl/3jJ1O0

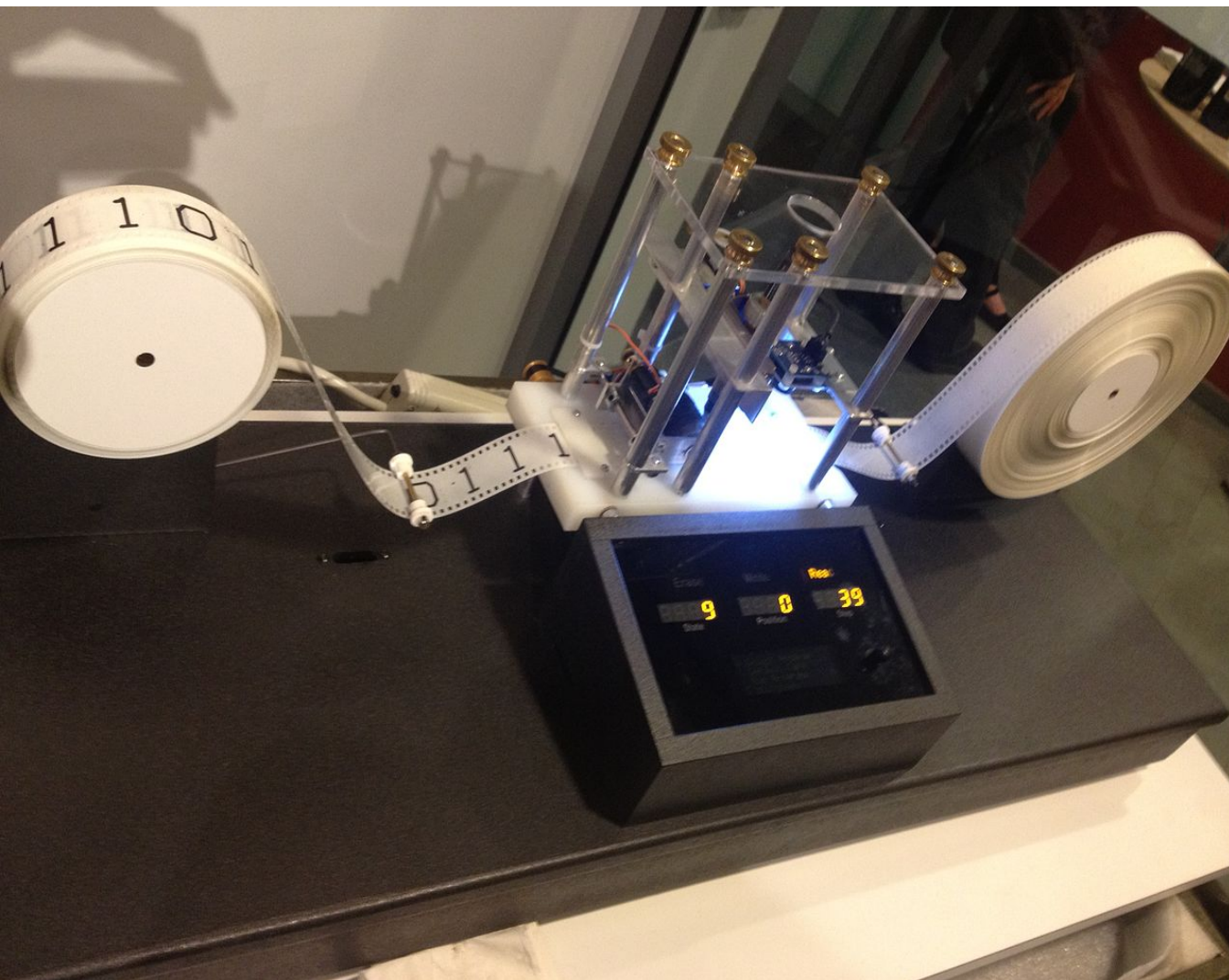Source: rdn consulting

"We can only see a short distance ahead, but we can see plenty there that needs to be done."

— Alan Turing

https://twitter.com/nvidia/status/1010545517405835264

# (Real) Turing machine

It is possible to invent a *single machine* which can be used to compute *any* computable sequence. If this machine **U** is supplied with the tape on the beginning of which is written the string of quintuples separated by semicolons of some computing machine **M**, then **U** will compute the same sequence as **M**.

Wikipedia

Can we learn from data a model that is as powerful as a Turing machine?

# Agenda

Brief review of deep learning

Neural Turing machine (NTM)

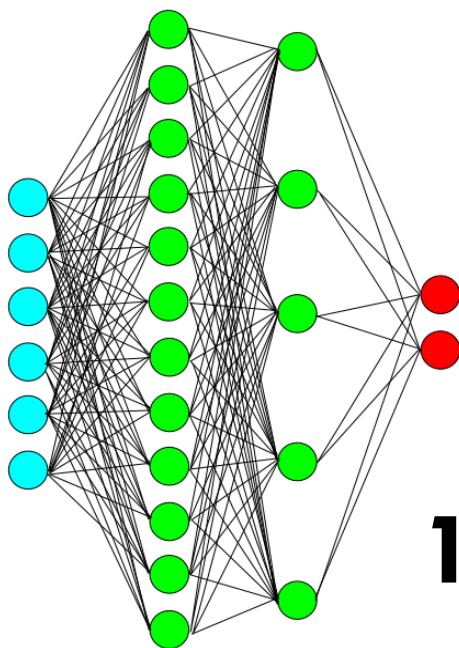Dual-controlling for read and write (PAKDD'18)

Dual-view in sequences (KDD'18)

Bringing variability in output sequences  (NIPS'18 ?)

Bringing relational structures into memory (IJCAI'17 WS+)
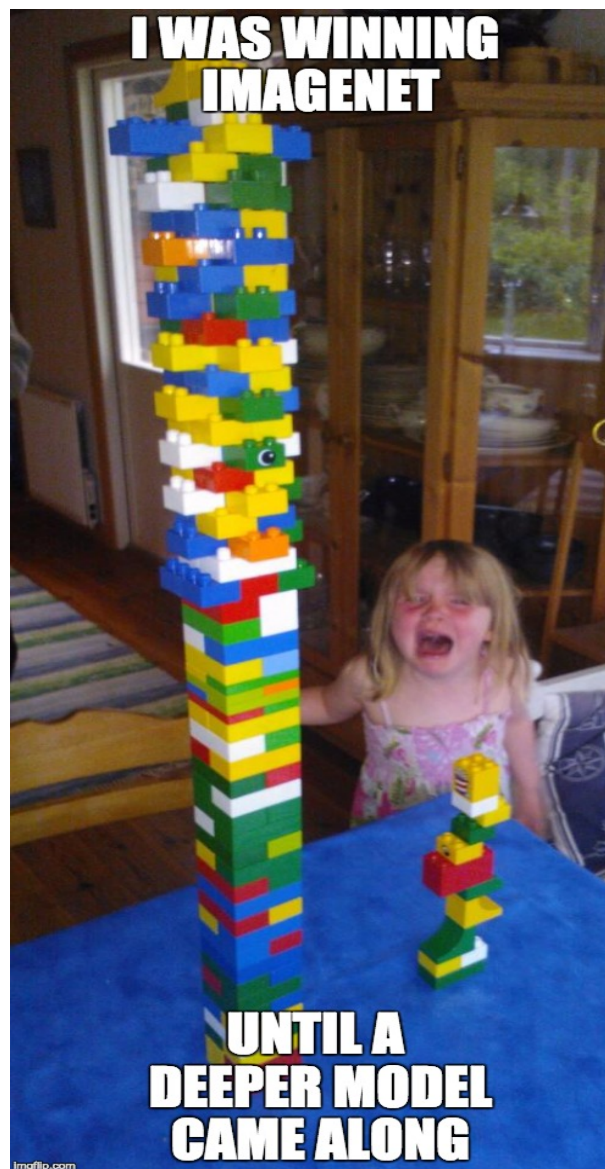
# Deep learning in a nutshell



Input layer   Hidden Layers   Output Layer

**1986**

http://blog.refu.co/wp-content/uploads/2009/05/mlp.png



I WAS WINNING IMAGENET

UNTIL A DEEPER MODEL CAME ALONG

imgflip.com

**2012**



Convolution
Pooling
Softmax
Other

**2016**

# Let's review current offerings

Feedforward nets (FFN)

Recurrent nets (RNN)

Convolutional nets (CNN)

Message-passing graph nets (MPGNN)

Universal transformer

…..

Work surprisingly well on LOTS of important problems

Enter the age of differentiable programming

**BUTS …**

No storage of intermediate results.

Little choices over what to compute and what to use

Little support for complex chained reasoning

Little support for rapid switching of tasks
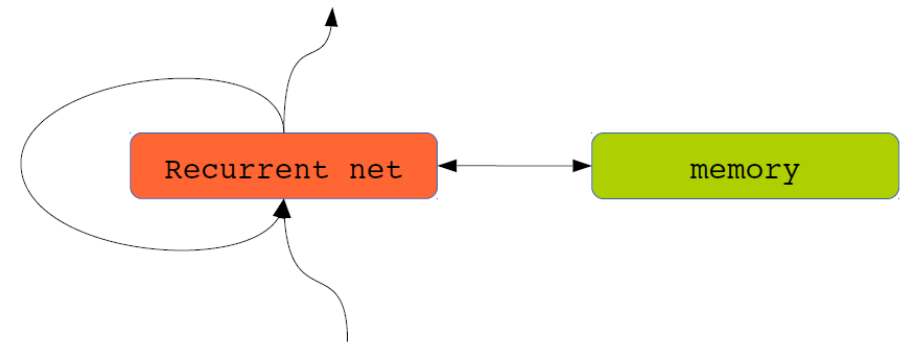
# Searching for better priors

Translation invariance in CNN

Recurrence in RNN

Permutation invariance in attentions and graph
neural networks

**Memory for complex computation**

   **→ Memory-augmented neural networks
(MANN)**



(LeCun, 2015)

# What is missing? A memory

Use multiple pieces of information

Store intermediate results (RAM like)

Episodic recall of previous tasks (Tape like)

Encode/compress & generate/decompress long sequences

Learn/store programs (e.g., fast weights)

Store and query external knowledge

Spatial memory for navigation

Rare but important events (e.g., snake bite)

Needed for complex control

Short-cuts for ease of gradient propagation = constant path length

Division of labour: program, execution and storage

Working-memory is an indicator of IQ in human

# Example: Code language model

```
FileWriter writer = new FileWriter(file);
writer.write(''This is an example'');
int count = 0;
System.out.prinltln(''Long gap'');
    . . . . . .

writer.flush();
writer.close();
```
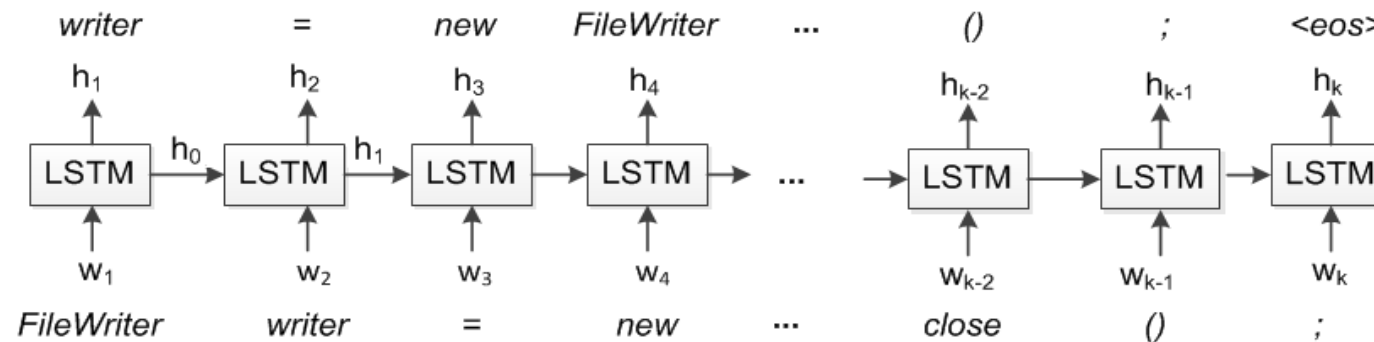
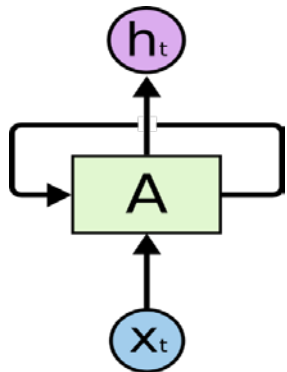**Still needs a better memory for:**

Repetitiveness
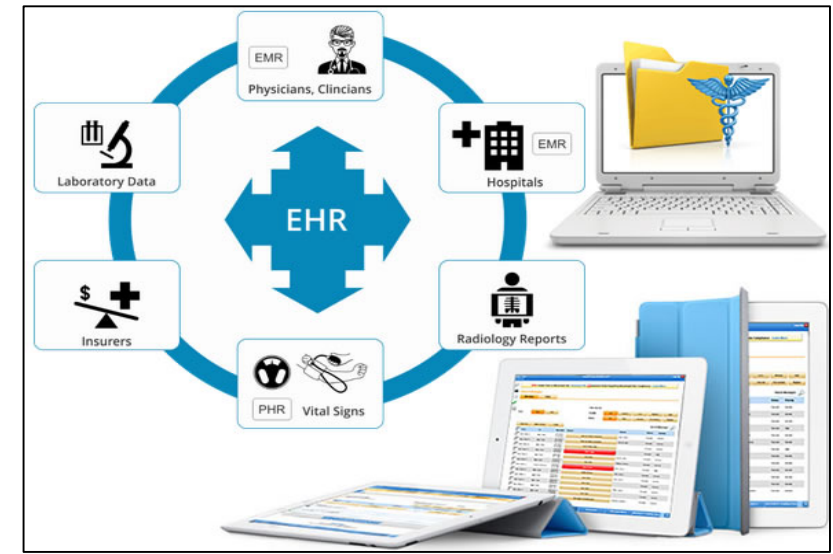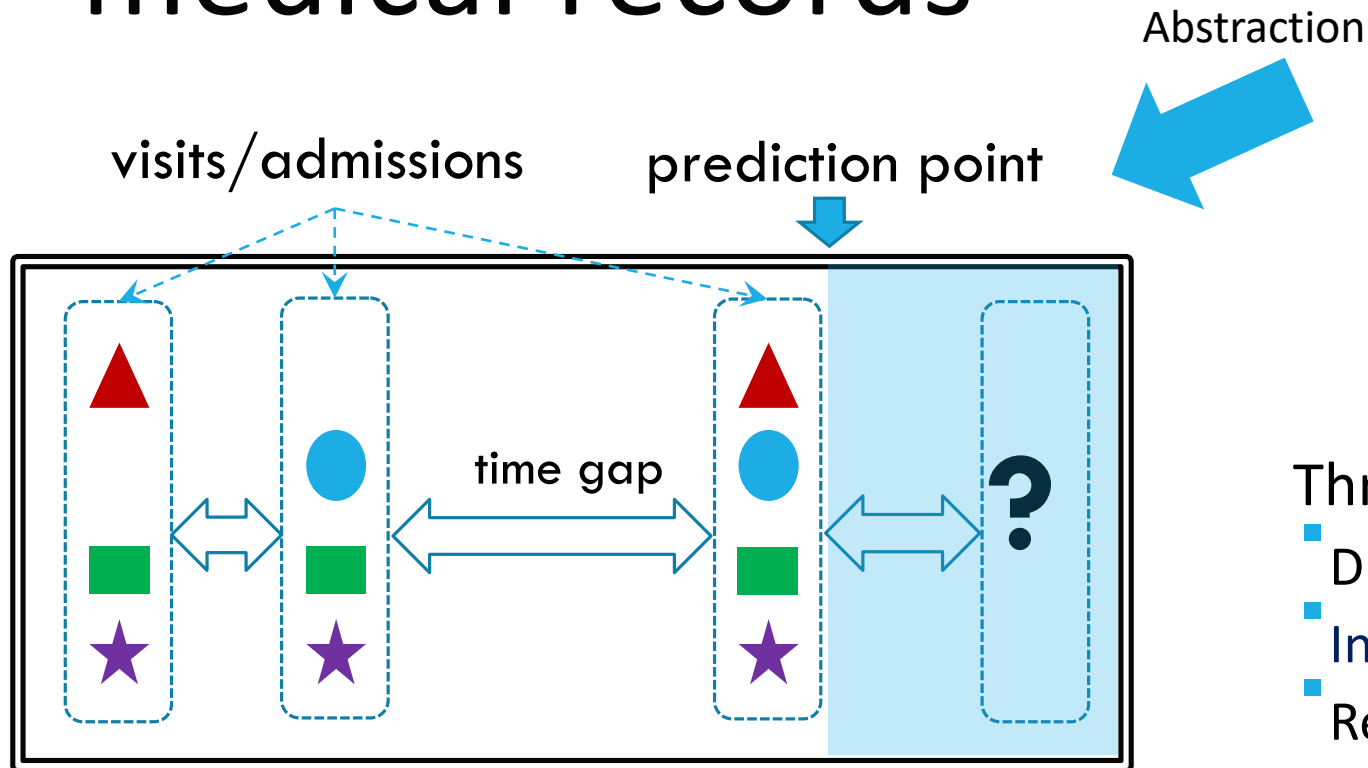    E.g. for (int i = 0; i < n; i++)
Localness
    E.g. *for (int size* may appear more often
    that *for (int i* in some source files.

Very long sequence (big file, or char level)

$$P(s) = P(w_1) \prod_{t=2}^{k} P\left(w_t \mid \boldsymbol{w}_{1:t-1}\right)$$

# Example: Electronic medical records



Source: medicalbillingcodings.org

Abstraction

Modelling

visits/admissions   prediction point

time gap

?

Three interwoven processes:
- Disease progression
- Interventions & care processes
- Recording rules

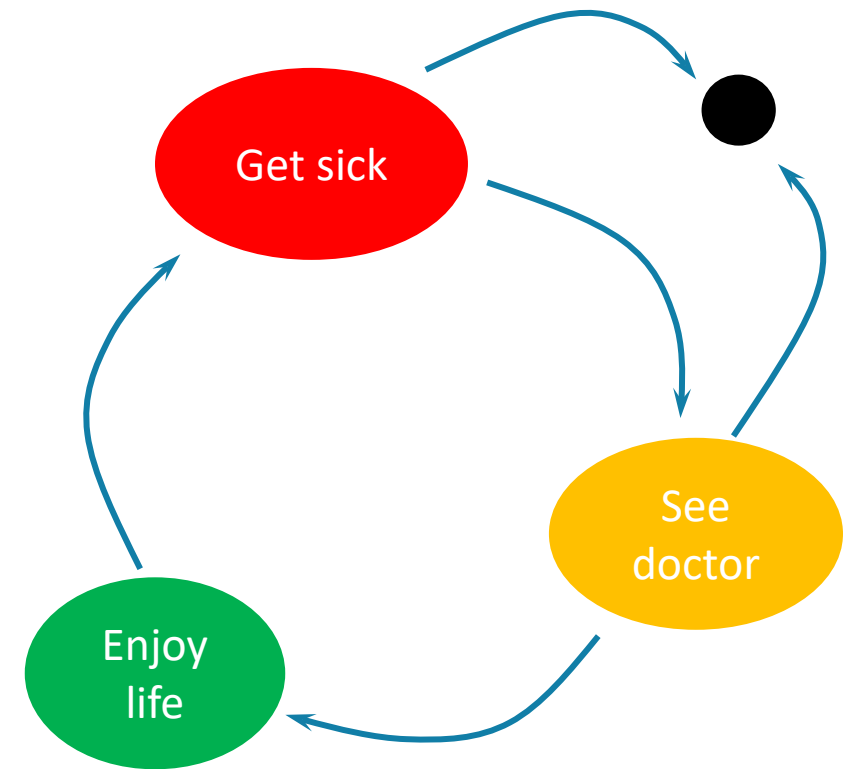Need memory to handle thousands of events

# EMR visualisation

A prototype system developed iHops (our spin-off)

# **Conjecture**: Healthcare is Turing computational

Healthcare processes as executable computer program obeying hidden "grammars"

The "grammars" are learnable through observational data

With "generative grammars", entire health trajectory can be simulated.

# Other possible applications of memory

Video captioning

QA, VQA

Machine translation

Machine reading (stories, books, DNA)

Business process continuation

Software execution

Code generation

Graph as sequence of edges
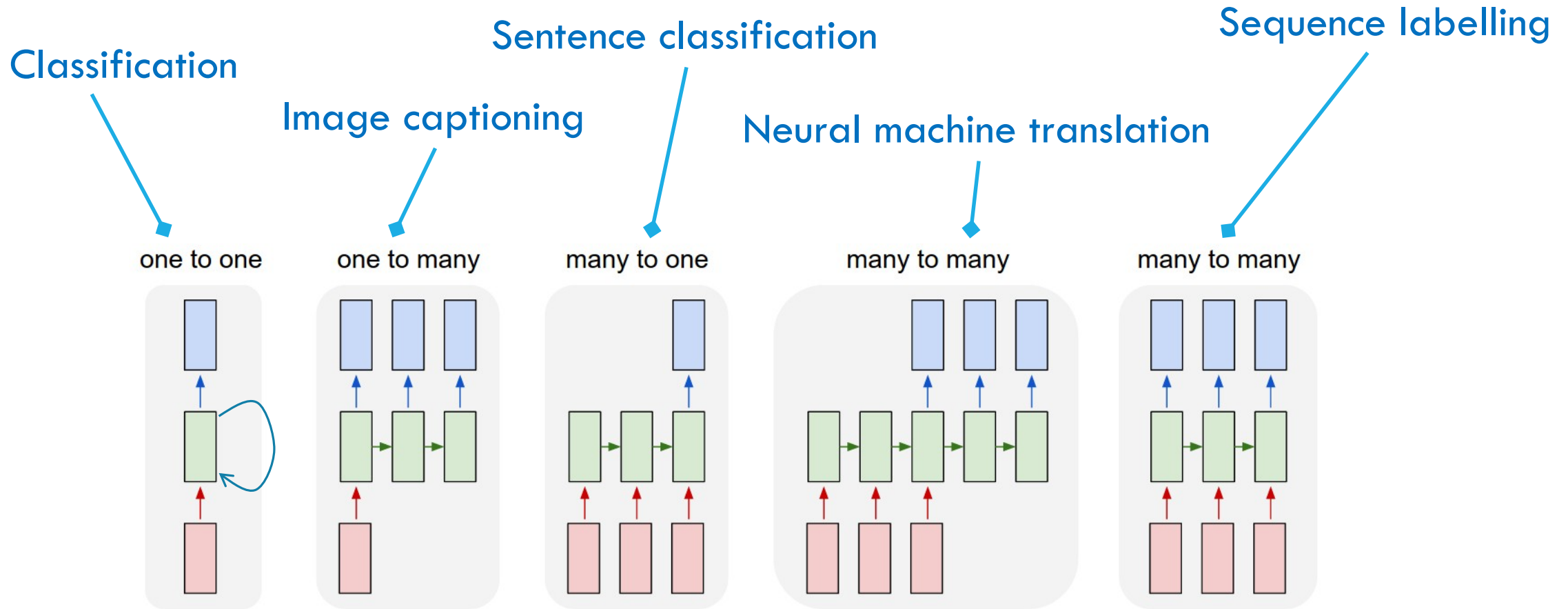
Event sequences

Graph traversal

Algorithm learning (e.g., sort)

Dialog systems (e.g., chat bots)

Reinforcement learning agents

# Neural Turing machine (NTM)

# RNN: theoretically powerful, practically limited

Classification

Sentence classification

Sequence labelling

Image captioning

Neural machine translation

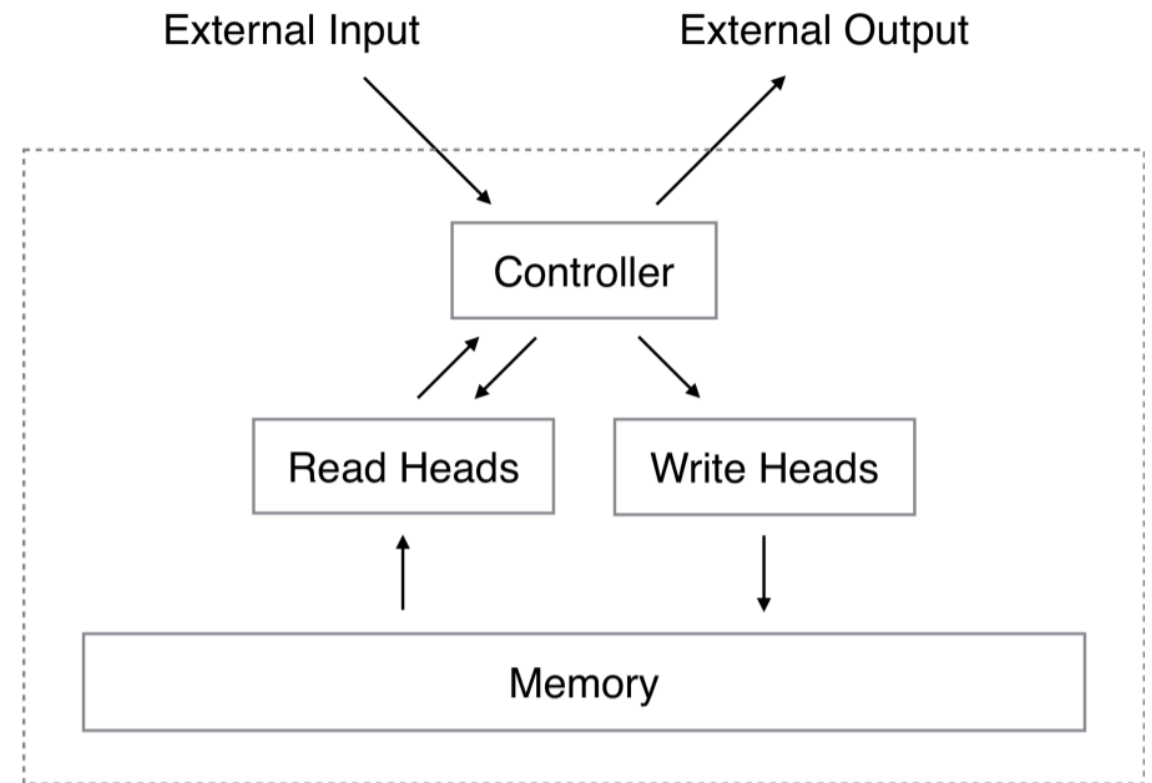| one to one | one to many | many to one | many to many | many to many |

# Neural Turing machine (NTM)

A controller that takes input/output and talks to an external memory module.

Memory has read/write operations.
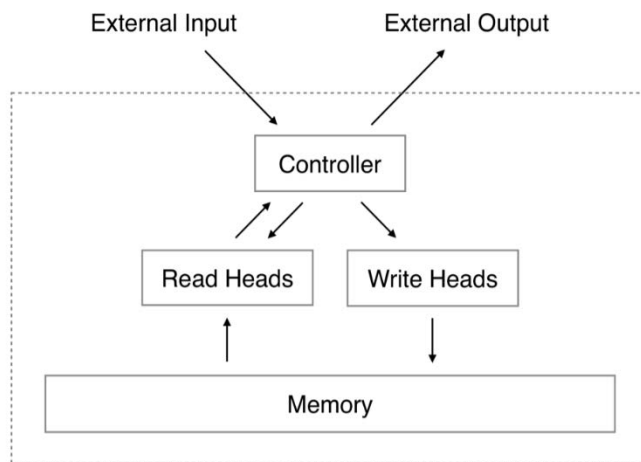
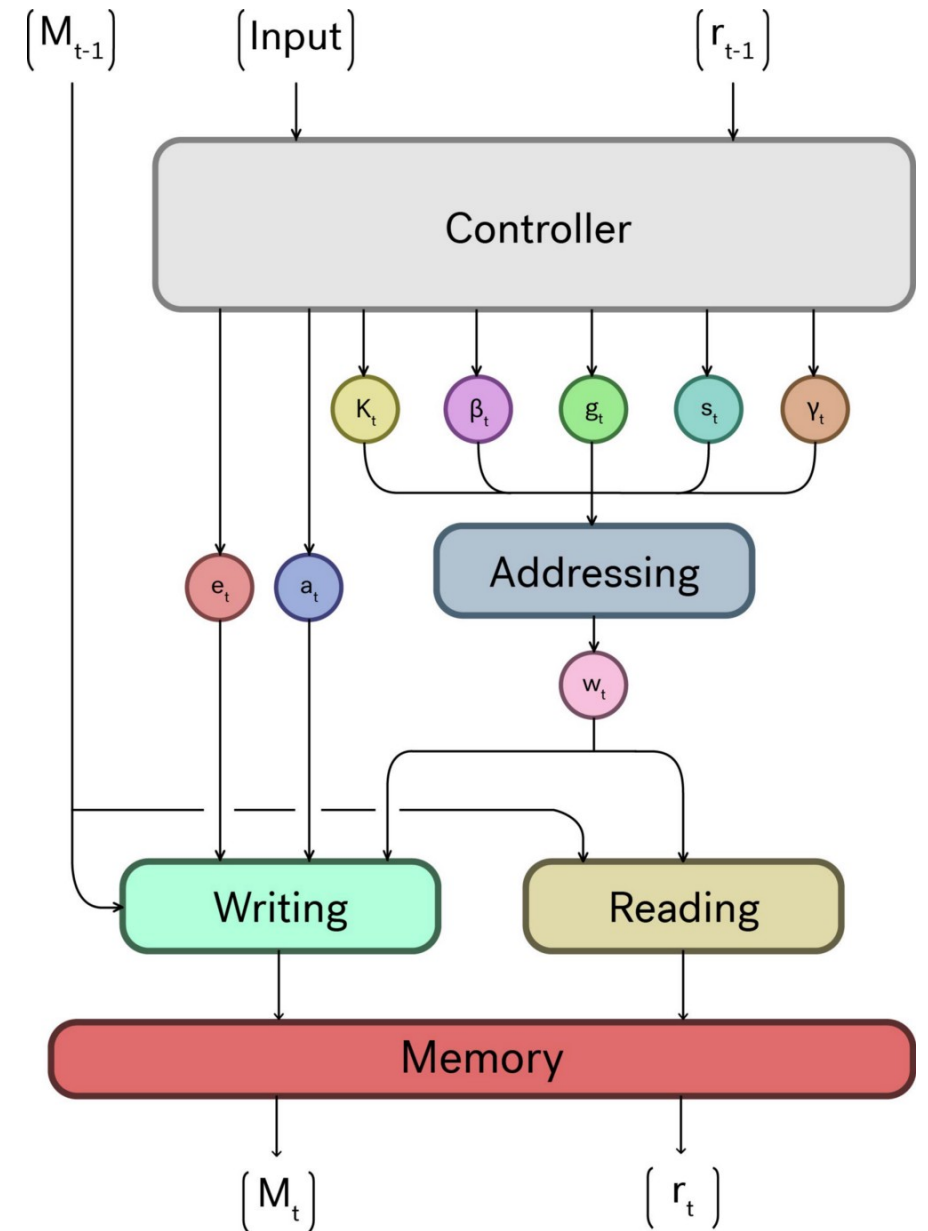The main issue is where to write, and how to update the memory state.

All operations are differentiable.

# NTM operations



https://rylanschaeffer.github.io



https://medium.com/@aidangomez/the-neural-turing-machine-79f6e806c0a1

NTM unrolled in time with LSTM as controller

#Ref: https://medium.com/snips-ai/ntm-lasagne-a-library-for-neural-turing-machines-in-lasagne-2cdce6837315

# Differentiable neural computer (DNC)

**2014**

**2016**

Illustration of the DNC architecture



https://rylanschaeffer.github.io

Source: deepmind.com

#REF: Graves, Alex, et al. "Hybrid computing using a neural network with dynamic external memory." *Nature* 538.7626 (2016): 471-476.

# Dual-controlling for read and write

Hung Le, Truyen Tran & Svetha Venkatesh

*PAKDD'18*

# MANN with dual control (DC-MANN)

Two controllers, for input & output

The encoder reads the input sequence is encoded into memory

The decoder reads the memory and produces a sequence of output symbols

**During decoding, the memory is write-protected (DCw-MANN)**



#REF: Hung Le, Truyen Tran, and Svetha Venkatesh. "Dual Control Memory Augmented Neural Networks for Treatment Recommendations", PAKDD18.

# DC-MANN

#Ref: https://medium.com/snips-ai/ntm-lasagne-a-library-for-neural-turing-machines-in-lasagne-2cdce6837315

# Result: Odd-Even Sequence Prediction

- Input: a sequence of random odd numbers → output: a sequence of even numbers

- Output:

$$y_n = \begin{cases} 2x_n & n \leq \lfloor \frac{L}{2} \rfloor \\ y_{n-1} + 2 & n > \lfloor \frac{L}{2} \rfloor \end{cases}$$
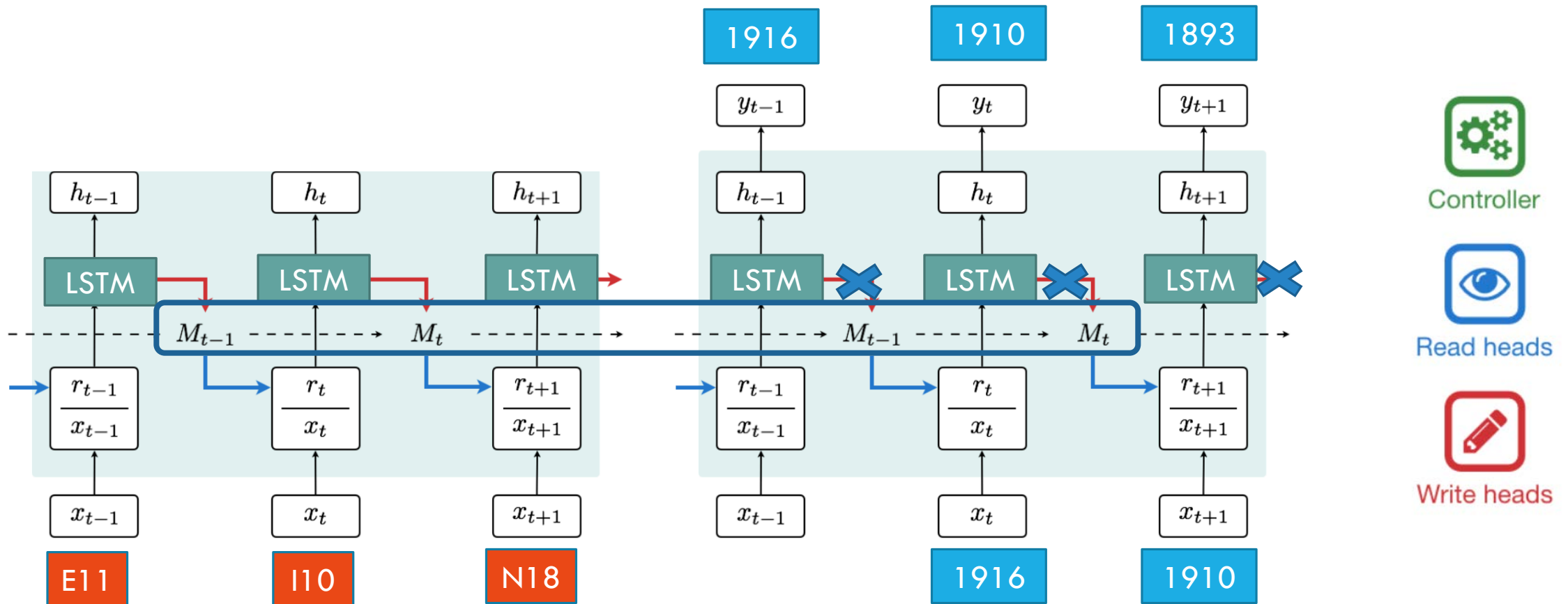
| Model | NLD |
|---|---|
| Seq2Seq | 0.679 |
| Seq2Seq with attention | 0.637 |
| DNC | 0.267 |
| DNC (write-protected) | 0.250 |
| DC-MANN | 0.161 |
| **DCw-MANN** | **0.082** |

Without memory, LSTMs fail the task

Write-protected helps

# Treatment recommendation



Admission 1 ··· Admission N-1 Admission N (current)

E11  I10  N18  1916  1910  ···  Z86  E11  A08  1952  1893  E11  T81  A08

Predict output sequence:
Treatments for current admission

?  ?  ?

# Result: Medicine prescription



MEDICINE PRESCRIPTION

| | Jaccard | Precision |
|---|---|---|
| DCw-MANN | 0.556 | 0.598 |
| DNC | 0.529 | 0.577 |
| Seq2Seq+attention | 0.142 | 0.224 |
| Seq2Seq | 0.138 | 0.22 |
| Random Forest | 0.405 | 0.491 |
| Logistic Regression | 0.311 | 0.412 |

# Compared to DNC



**Fig. 5.** Training Loss of Drug Prescription Task

**Fig. 6.** Testing Loss of Drug Prescription Task
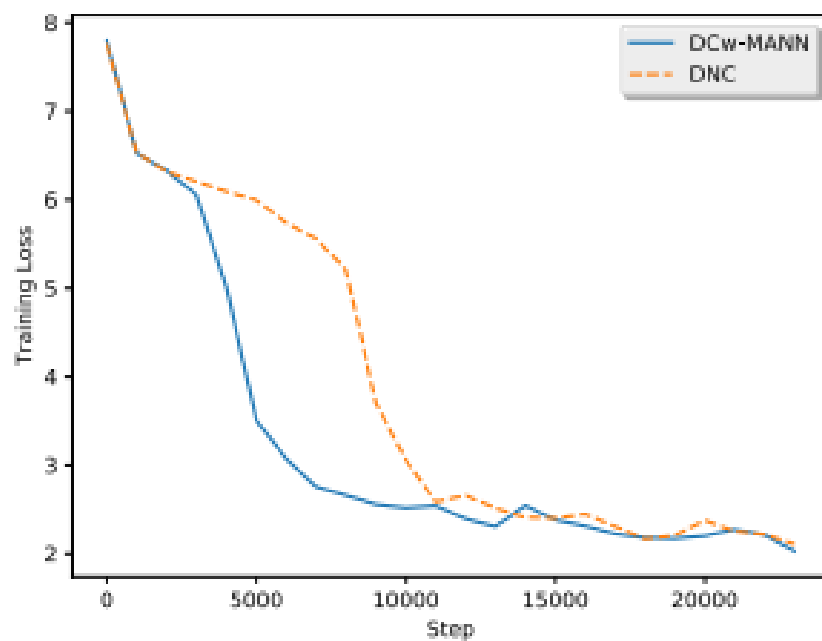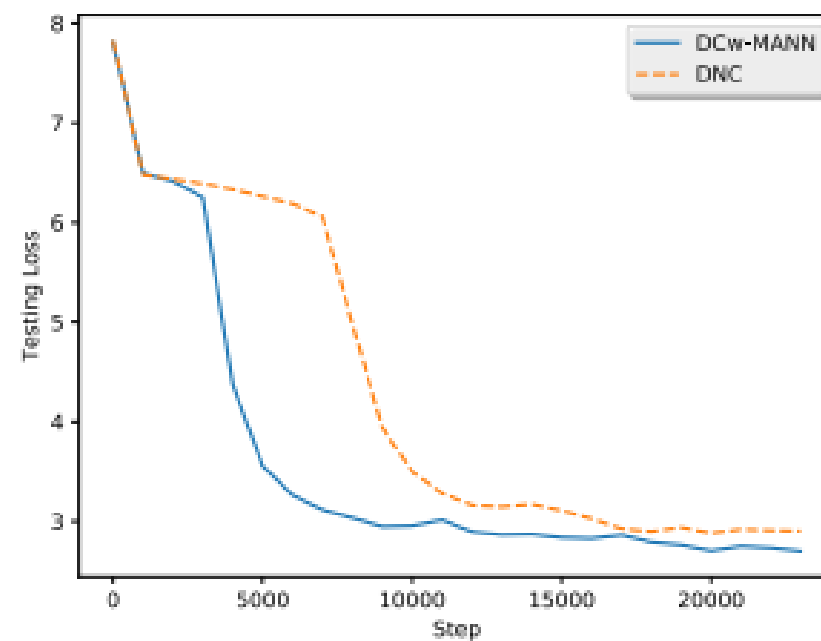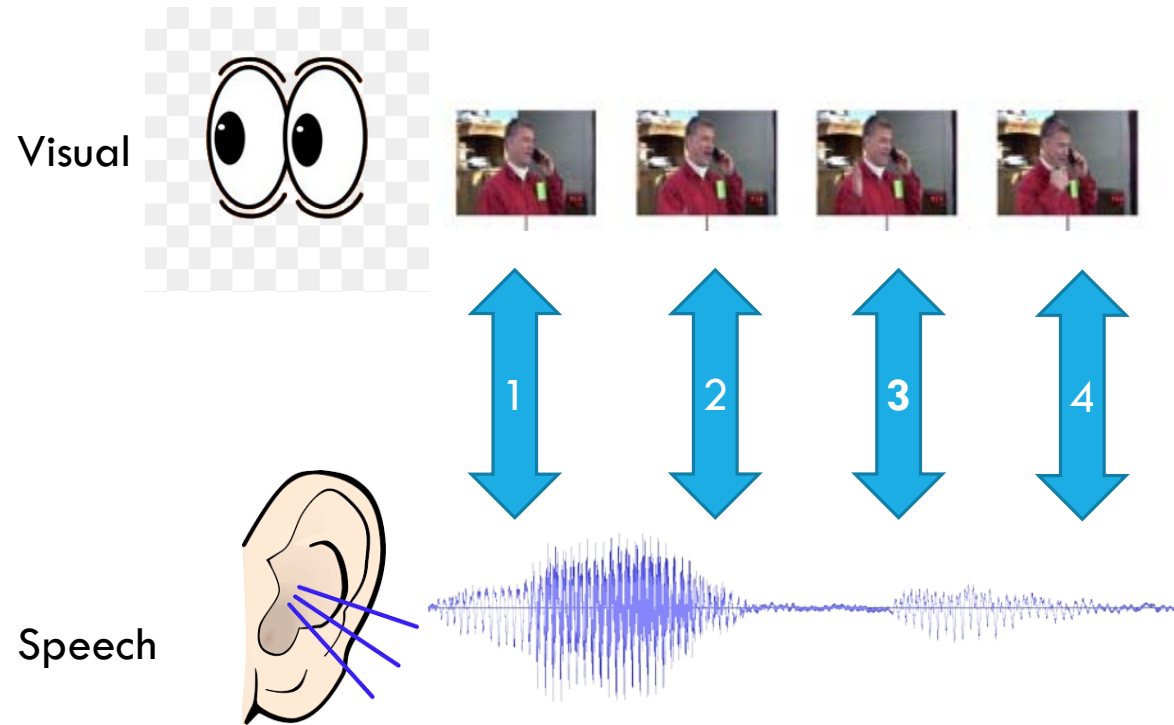
# Dual-view sequential problems

Hung Le, Truyen Tran & Svetha Venkatesh

*KDD'18*

# Synchronous two-view sequential learning

# Asynchronous two-view sequential learning Healthcare: medicine prescription



Diagnoses

E11 | I10 | N18 | Z86 | E11

Procedures

1916 | 1910 | 1952 | 1893

Medicines

DOCU100L | ACET325

# Asynchronous two-view sequential learning Healthcare: disease progression

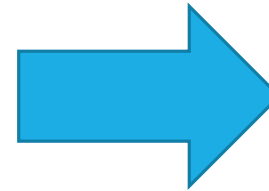Previous diagnoses

| E11 | I10 | N18 | Z86 | E11 |

Future diagnoses ???

| 1916 | 1910 | ACET325 | DOCU100L |

Previous interventions

# Intra-view & inter-view interactions

# Dual architecture



*Dual Memory Neural Computer (DMNC). There are two encoders and one decoder implemented as LSTMs. The dash arrows represent cross-memory accessing in early-fusion mode*

Simple sum, but distant, asynchronous

$$\{y_i = x_i^1 + x_{L+1-i}^2\}_{i=1}^{L}$$



Learning curve

Accuracy

DMNC   Others
≈ 99%   <55%

Medicine prescription performance
(data: MIMIC-III)

Disease progression performance
(data: MIMIC-III)

DMNC    WLAS    DeepCare

| | P@1 Dieabies | P@2 Dieabies | P@3 Dieabies | P@1 Mental | P@2 Mental | P@3 Mental |
|---|---|---|---|---|---|---|
| DMNC | 67.6 | 61.3 | 57 | 53.6 | 50 | 47.1 |
| WLAS | 65.9 | 60.8 | 56.5 | 51.8 | 48.9 | 45.7 |
| DeepCare | 66.2 | 59.6 | 53.7 | 52.7 | 49.4 | 46.2 |

# Bringing variability in output sequences

Hung Le, Truyen Tran & Svetha Venkatesh

*Submitted to NIPS'18*

# Motivation: Dialog system

A dialog system needs to maintain the history of chat (e.g., could be hours)

- → Memory is needed

The generation of response needs to be flexible, adapting to variation of moods, styles

- Current techniques are mostly based on LSTM, leading to "stiff" default responses (e.g., "I see").

There are many ways to express the same thought

- → Variational generative methods are needed.

# Variational Auto-Encoder (VAE)
## (Kingma & Welling, 2014)

Two separate processes: generative (hidden → visible) versus recognition (visible → hidden)



http://kvfrans.com/variational-autoencoders-explained/

# Variational memory encoder-decoder (VMED)

generated     context



$p(y|x,z)$

$p_\phi(z|x)$

$q_\theta(z|x,y)$

latent variables

**Conditional Variational Auto-Encoder**

generated     context

$q_\theta(z|x,y,r)$

reads

$p(y|x,z)$

$p_\phi(z|x,r)$

r

M

latent variables

memory

**VMED**

# Sample response

| Input context | Response |
|---|---|
| **Reddit comment:** What is your favorite scene in film history ? Mine is the restaurant scene in the Godfather. | **Seq2Seq:** The scene in <br> **Seq2Seq-att:** The final <br> **DNC:** The scene in <br> **CVAE:** Inception god! Not by a shark /*/ Amour great /*/ Pro thing you know 3 dead <br> **VMED (K=3):** The opening scene from history movie /*/ The scene in a shot nights! Robin movie /*/ The psycho scene in fight from |
| **Reddit comment:** What actor will win an Oscar in the next 10 years ? | **Seq2Seq:** Colin <br> **Seq2Seq-att:** Liam Neeson <br> **DNC:** Tom Gyllenhaal <br> **CVAE:** Daryl and Aaron /*/ Carefully count Alfred Deniro /*/ Ponyo Joker posible <br> **VMED (K=3):** Edward or Leo Dicaprio goes on /*/ Dicaprio will /*/ Dicaprio Tom has actually in jack on road |

# Sample response (2)

| | |
|---|---|
| **LJ user:** Your very pretty and your is very cute. May i add ya to my friends list ? | **Seq2Seq:** I! Add and too you back<br>**Seq2Seq-att:** I did may! Away back<br>**DNC:** Sure you added<br>**CVAE:** Sure ! I can add you /*/ Thank cool of surely /*/ Thank hi of sure!<br>**VMED (K=3):** Thank cool for sure /*/ Add definitely! Sure surely friends /*/ Thank girl go ahead |
| **Cornell movies:** Nobody can murder someone with their feelings . I mean that is crazy ? | **Seq2Seq:** Yes you are<br>**Seq2Seq-att:** Really not is it<br>**DNC:** Managing the truth<br>**CVAE:** Sure not to be in here /*/ Oh yes but i know! /*/ That to doing with here and<br>**VMED (K=4):** You are right! /*/ That is crazy /*/ You can't know Jimmy |
| **OpenSubtitle:** I'm obliged for your hospitality. I appreciate it and your husband too. | **Seq2Seq:** That is have got coming about these<br>**Seq2Seq-att:** May you not what nothing nobody<br>**DNC:** Yes i am taking this<br>**CVAE:** No . that for good! And okay /*/ All in the of two thing /*/ Sure. Is this!<br>**VMED (K=3):** I know. I can afford /*/ I know nothing to store for you pass /*/ I know. Doing anymore you father |

# Bringing relational structures into memory

Trang Pham, Truyen Tran & Svetha Venkatesh

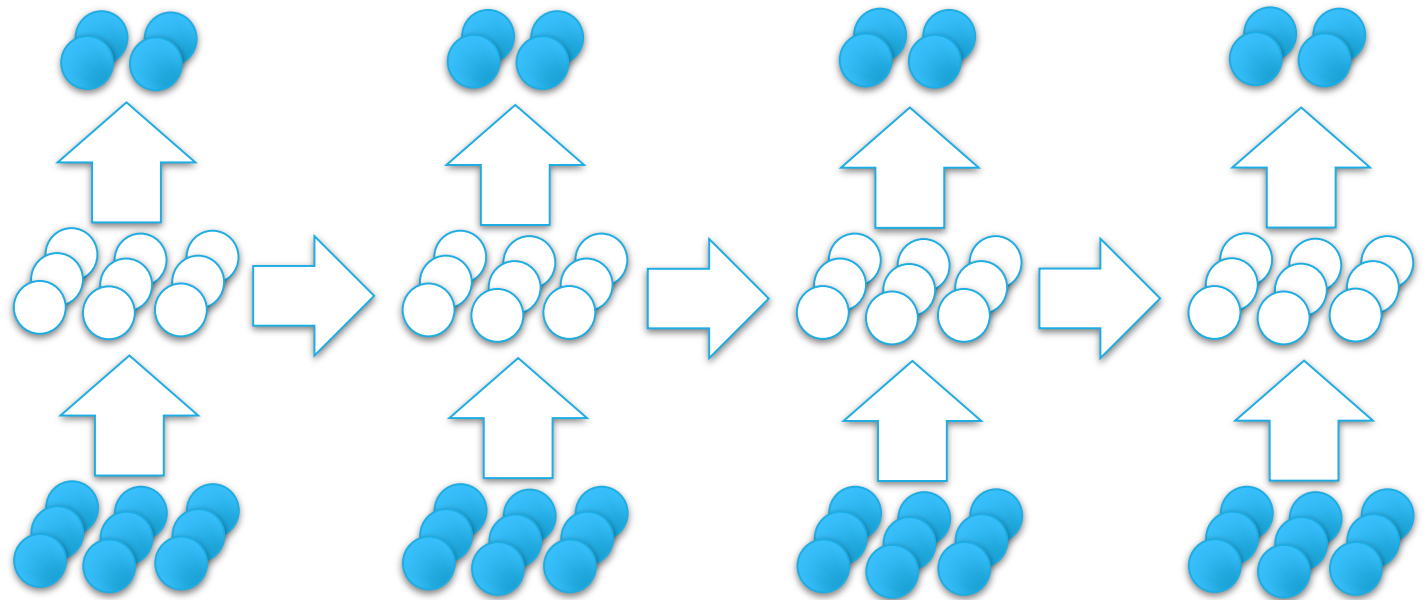*IJCAI'17 WS+*

# NTM as matrix machine

Controller and memory operations can be conceptualized as matrix operations

- **Controller is a vector changing over time**

- **Memory is a matrix changing over time**

#REF: Kien Do, Truyen Tran, Svetha Venkatesh, "Learning Deep Matrix Representations", *arXiv preprint arXiv:*1703.01454

$$H_t = \sigma(U_x^\mathsf{T} X_t V_x + U_h^\mathsf{T} H_{t-1} V_h + B)$$

# Idea: Relational memory

Independent memory slots not suitable for relational reasoning

Human working memory sub-processes seem inter-dependent

Transformation

**Relational structure**

$$H_t = \sigma(U_x^\intercal X_t V_x + U_h^\intercal H_{t-1} V_h + B)$$
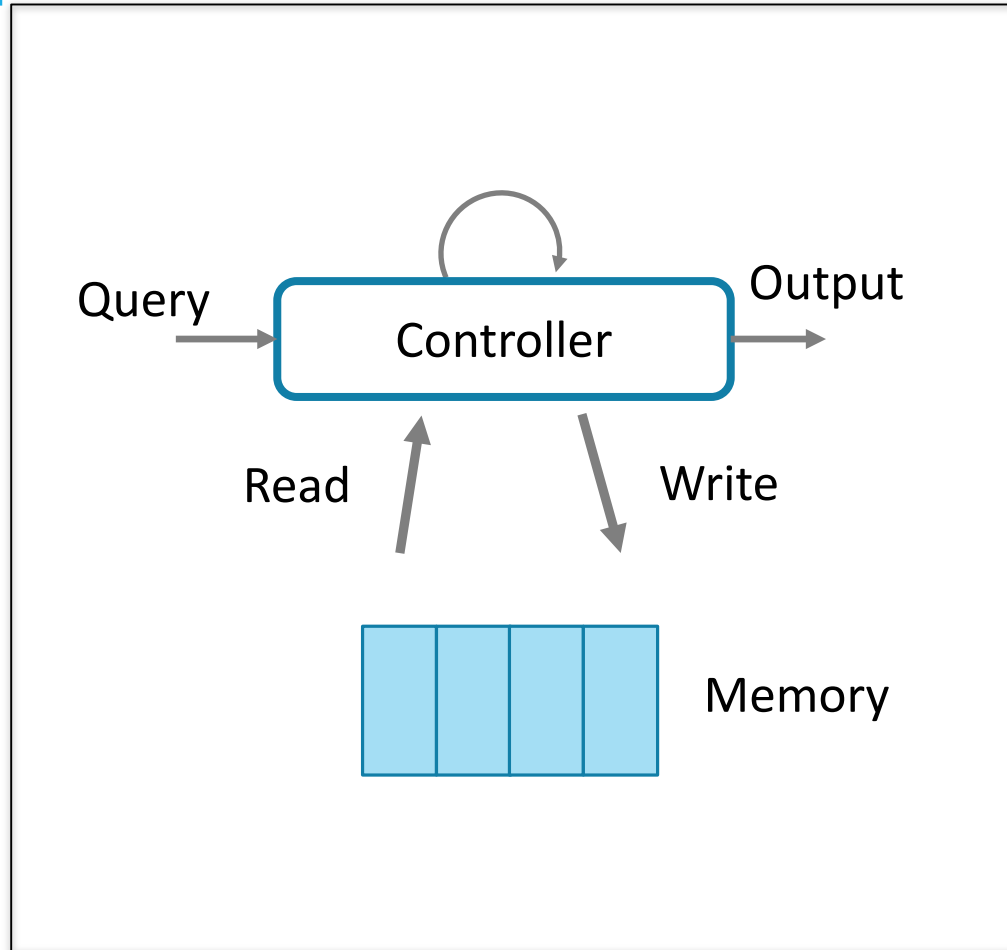
New memory proposal

New information
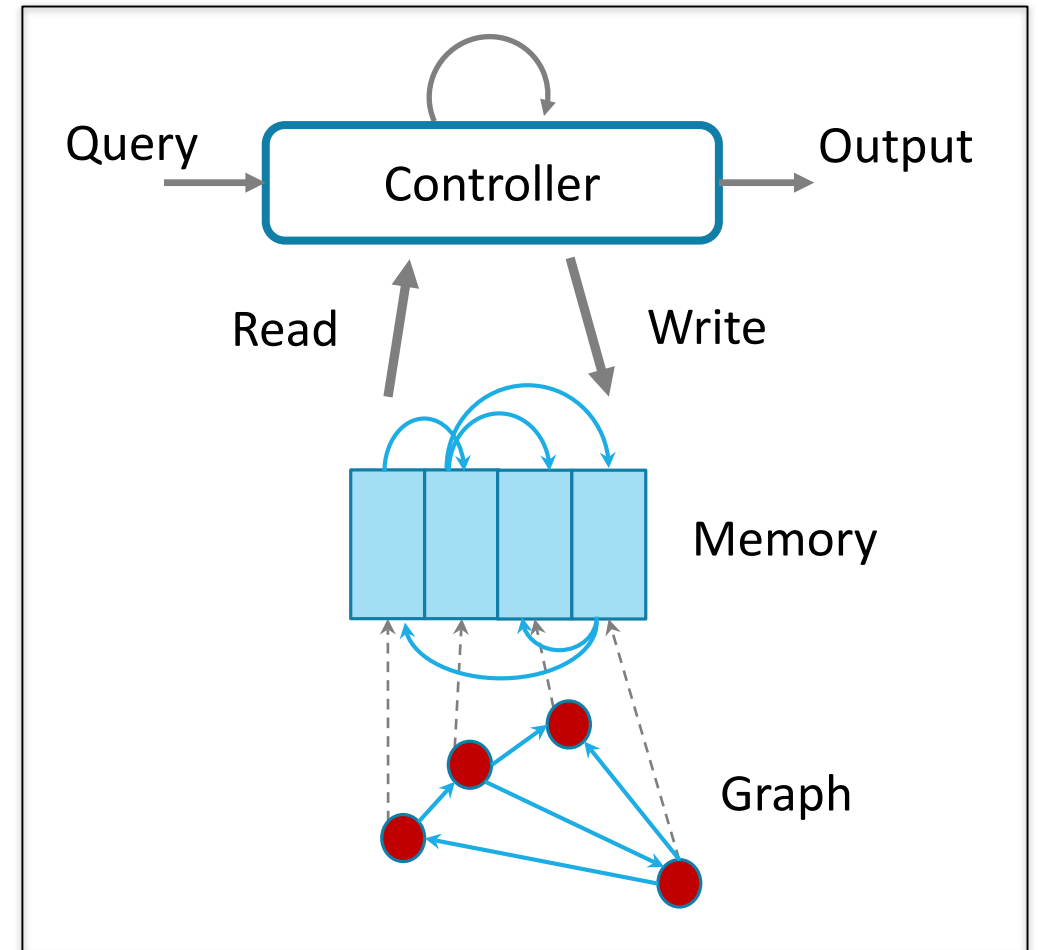
Old memory

Time-aware bias

# Relational Dynamic Memory Network (DMNN)



NTM

Relational Dynamic Memory Network

# RDMN unrolled



Output process

Memory process

Controller process

Message passing

Input process

# Drug-disease response

Molecule → Bioactivity

| Model | MicroF1 | MacroF1 | Average AUC |
|---|---|---|---|
| SVM | 66.4 | 67.9 | 85.1 |
| RF | 65.6 | 66.4 | 84.7 |
| GB | 65.8 | 66.9 | 83.7 |
| NeuralFP [19] | 68.2 | 67.6 | 85.9 |
| MT-NN [51] | 75.5 | 78.6 | 90.4 |
| RDMN | **77.8** | **80.3** | **92.1** |

# Chemical reaction

Molecules → Reaction

| | CCI900 | | CCI800 | |
|---|---|---|---|---|
| | AUC | F1-score | AUC | F1-score |
| Random Forests | 94.3 | 86.4 | 98.2 | 94.1 |
| Highway Networks | 94.7 | 88.4 | 98.5 | 94.7 |
| DeepCCI [38] | 96.5 | 92.2 | 99.1 | 97.3 |
| RDMN | 96.6 | 92.6 | 99.1 | 97.4 |
| RDMN+multiAtt | 97.3 | 93.4 | 99.1 | 97.8 |
| RDMN+FP | 97.8 | 93.3 | 99.4 | 98.0 |
| RDMN+multiAtt+FP | 98.0 | 94.1 | 99.5 | 98.1 |
| RDMN+SMILES | 98.1 | 94.3 | 99.7 | 97.8 |
| RDMN+multiAtt+SMILES | **98.1** | **94.6** | **99.8** | **98.3** |

# Team @ Deakin (A2I2)



**Thanks to many people who have created beautiful graphics & open-source programming frameworks.**

# References

Memory–Augmented Neural Networks for Predictive Process Analytics, A Khan, H Le, K Do, T Tran, A Ghose, H Dam, R Sindhgatta, *arXiv preprint* arXiv:1802.00938

Learning deep matrix representations, K Do, T Tran, S Venkatesh, *arXiv preprint* arXiv:1703.01454

Variational memory encoder-decoder, H Le, T Tran, T Nguyen, S Venkatesh, *arXiv preprint* arXiv:1807.09950

Relational dynamic memory networks, Trang Pham, Truyen Tran, Svetha Venkatesh, *arXiv preprint* arXiv:1808.04247

Dual Memory Neural Computer for Asynchronous Two-view Sequential Learning, H Le, T Tran, S Venkatesh, *KDD'18*

Dual control memory augmented neural networks for treatment recommendations, H Le, T Tran, S Venkatesh, *PAKDD'18*.