

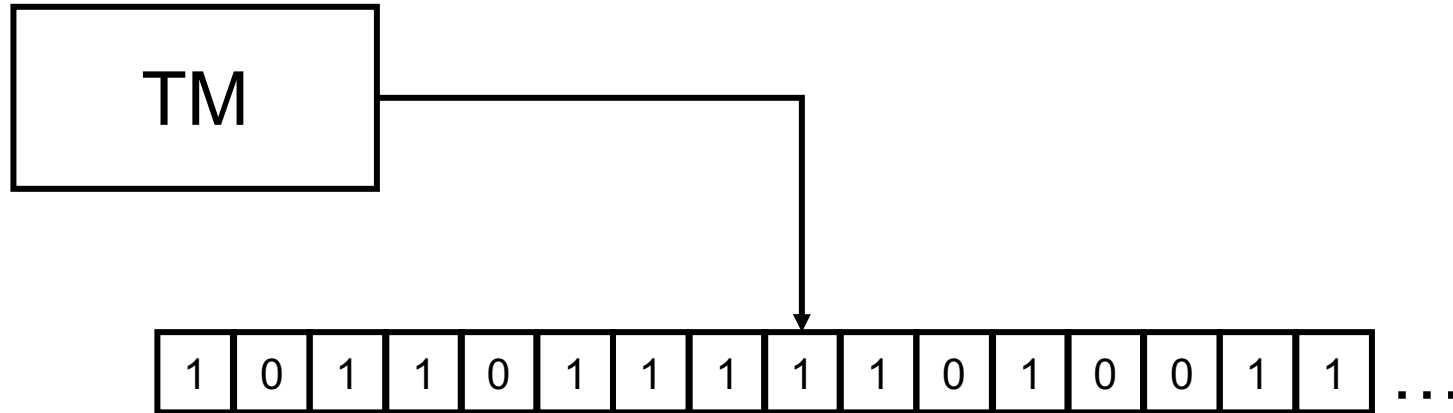
# Computational Power and Computational Limits of Neural Networks

Seminar in Deep Neural Networks FS2022

17 May 2022

Martynas Noruišis  
Mentor: Peter Belcák

# Turing Machines



# Simulating Turing Machines with Neural Networks

# Construction

- A Recurrent Neural Network of 886 processors (nodes)
- Rational weights
- *No slow down (real time simulation)*

# Sigmoid

*Saturated linear function*

$$\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } 0 \leq x \leq 1 \\ 1 & \text{if } x > 1 \end{cases}$$

1-tape TM = 2-stack Push-Down  
Automaton (PDA) &  
1-tape TM = 2-stack PDA  $\Rightarrow$   
 $p$ -tape TM =  $2p$ -stack PDA

# Encoding values

$$\omega = \omega_1 \omega_2 \dots \omega_n$$

$$\delta[\omega] = \delta[\omega_1 \omega_2 \dots \omega_n] = \sum_{i=1}^n \frac{\omega_i}{2^i}$$

1 0 0 0 0 1 0 1

1 0 0 0 0 1 0 1 0



1 0 0 0 0 0 0 0 0 0 0 0  
0 1 1 1 1 1 1 1 1 1 1

# Encoding binary strings

$$\omega = \omega_1 \omega_2 \dots \omega_n$$

$$\delta[\omega] = \delta[\omega_1 \omega_2 \dots \omega_n] = \sum_{i=1}^n \frac{2\omega_i + 1}{4^i}$$

# Stack operations

$$\text{stack } q = \sum_{i=1}^n \frac{2\omega_i + 1}{4^i}$$

peek

$$\text{top}(q) = \sigma(4q - 2)$$

push 0

$$\frac{q}{4} + \frac{1}{4}$$

push 1

$$\frac{q}{4} + \frac{3}{4}$$

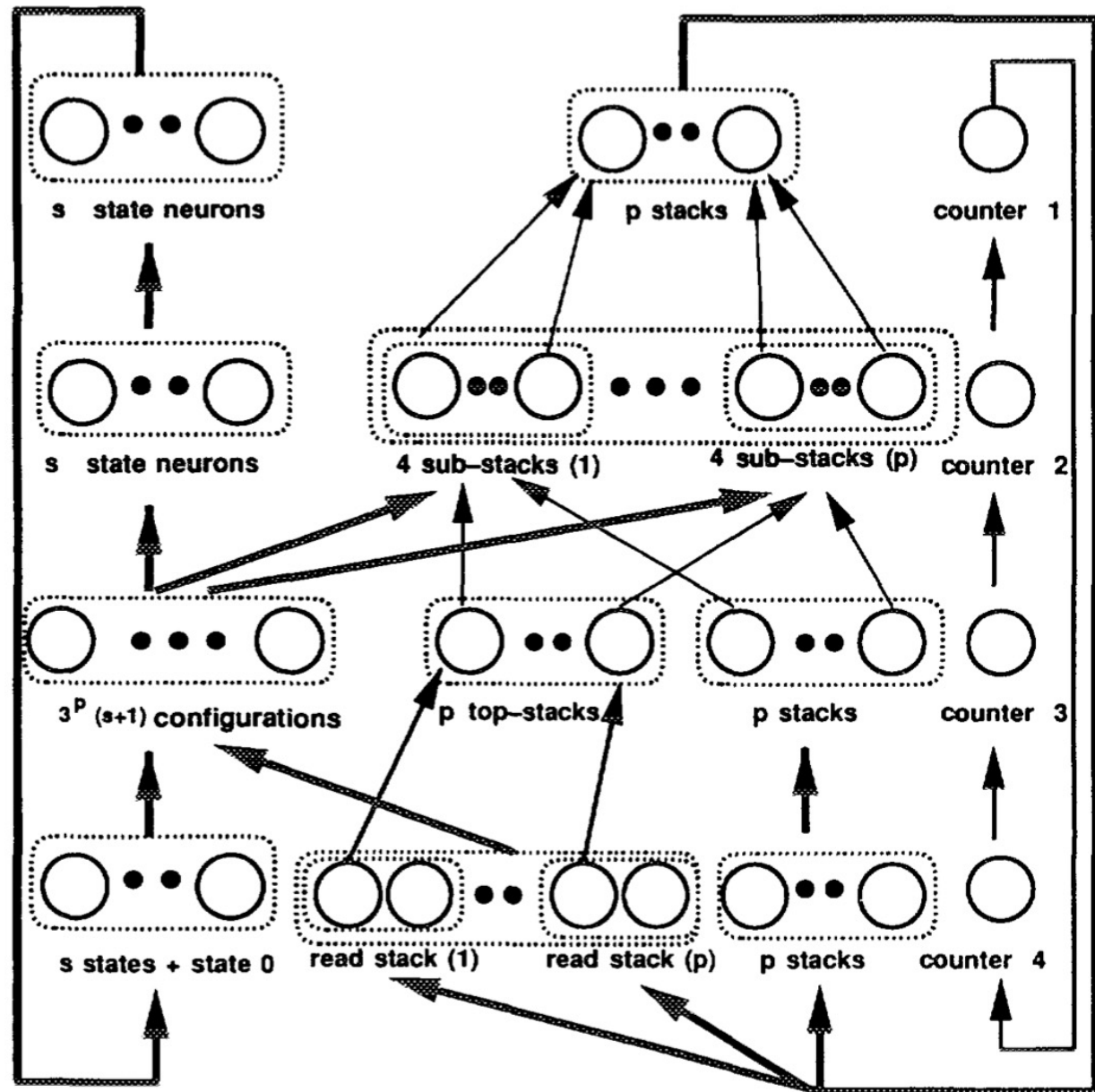
pop

$$4q - (2 \text{top}(q) + 1)$$

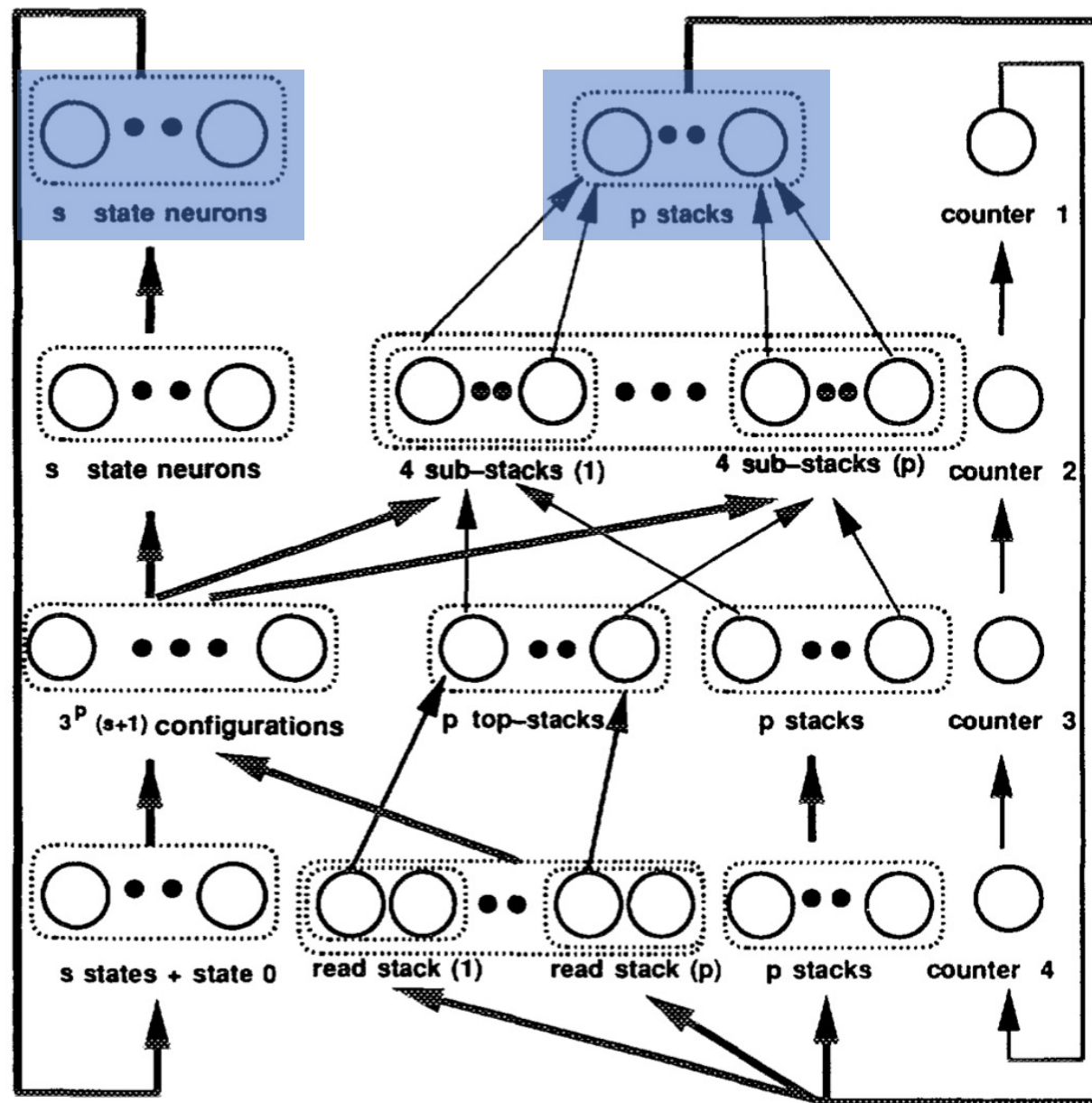
non-empty stack

$$\sigma(4q)$$

# Layout of the construction



$$(x_1, \dots, x_p, q_1, \dots, q_p) \in \mathbb{Q}^{s+p}$$



$F_1$

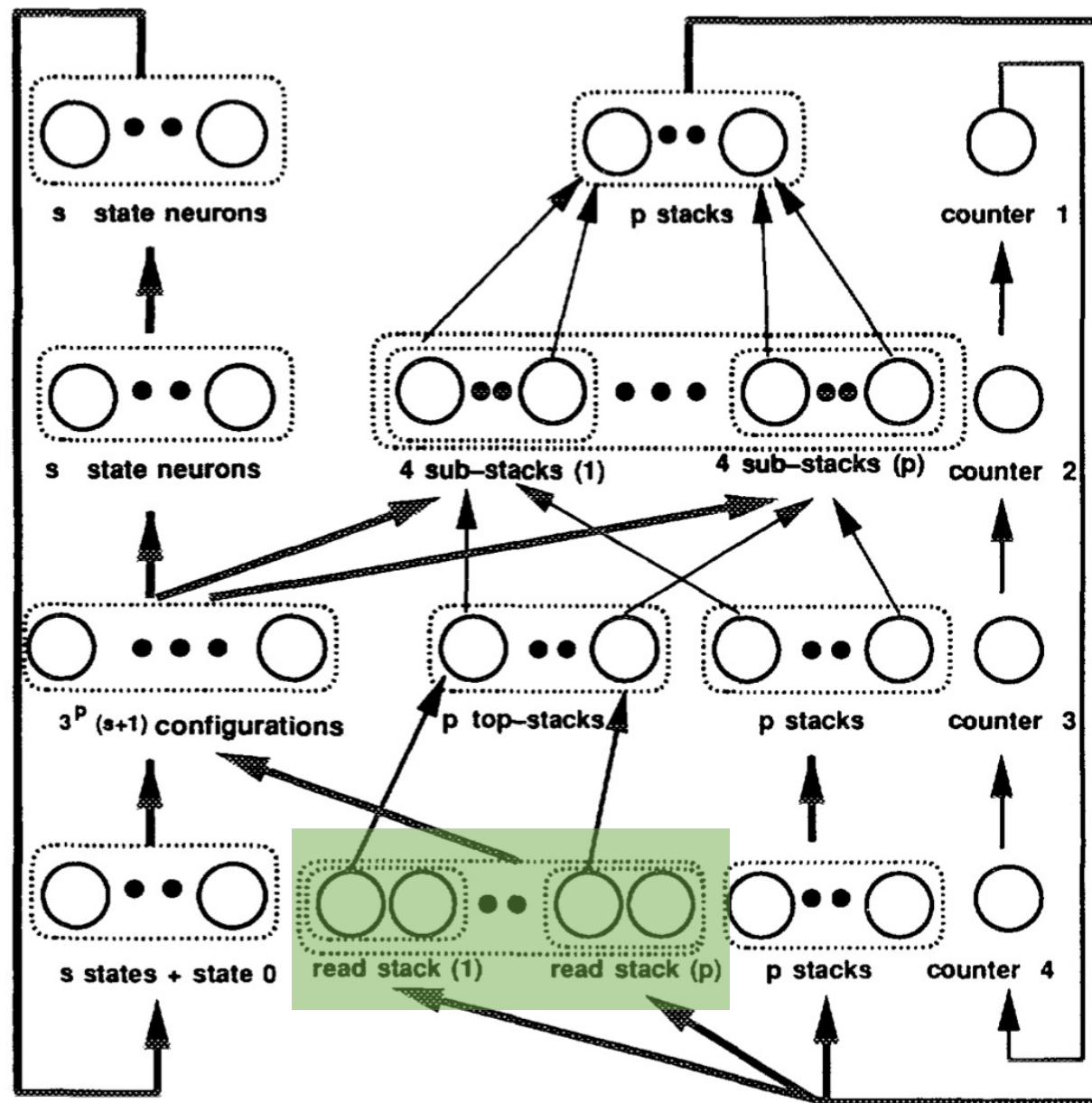
$F_2$

$F_3$

$F_4$

$$(\zeta[q_i], \tau[q_i]) \in \{0, 1\} \times \{0, 1\}$$

(top( $q_i$ ), non\_empty( $q_i$ ))



$F_1$

$F_2$

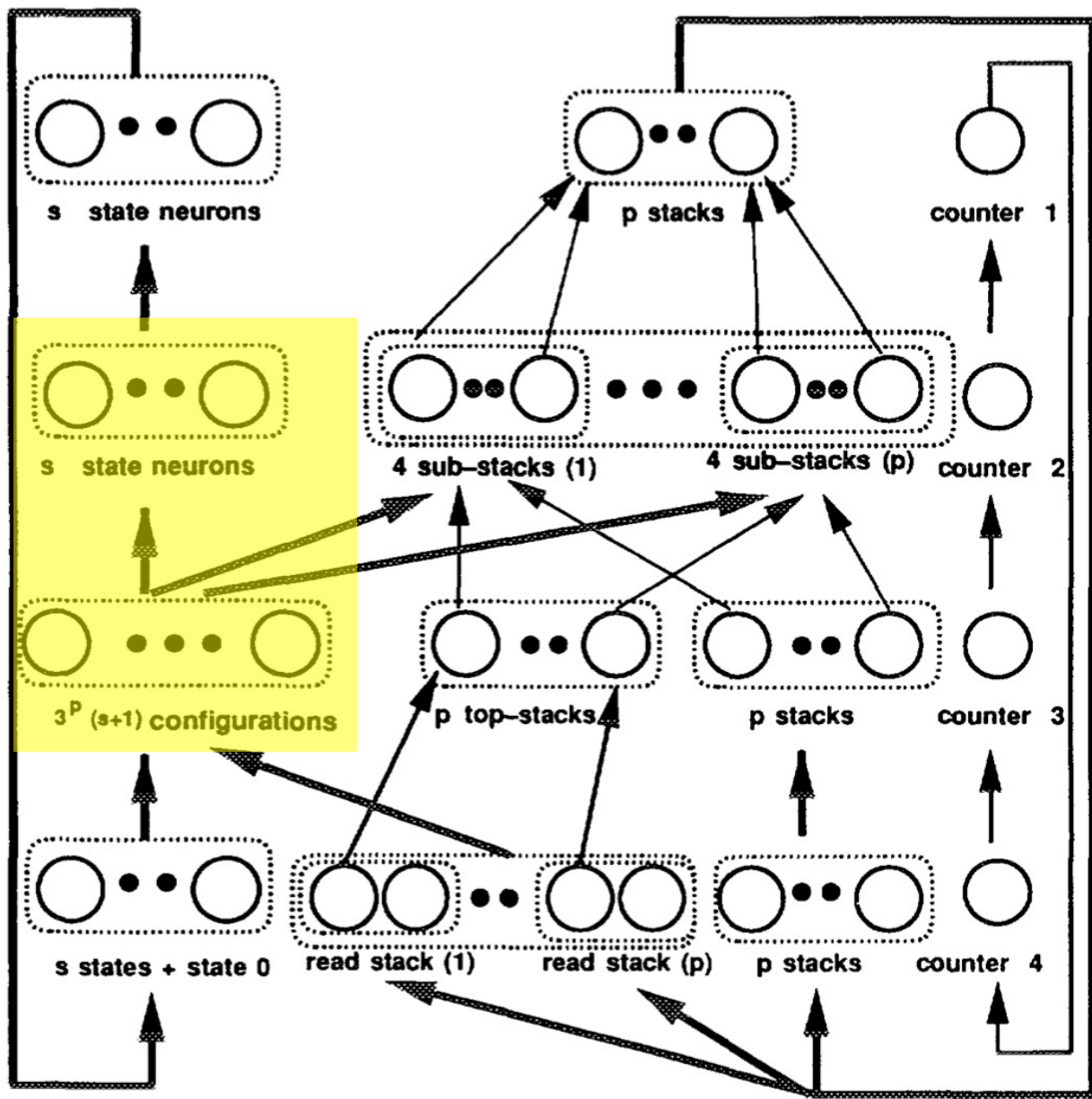
$F_3$

$F_4$

$$x_i^+ := \sum_{j=0}^s \beta_{ij}(\zeta[q_1], \dots, \zeta[q_1], \tau[q_i], \dots, \tau[q_p])x_j$$

$\beta_{ij}(\zeta[q_1], \dots, \zeta[q_1], \tau[q_i], \dots, \tau[q_p]) = 1$   
 $\Leftrightarrow$  there is a transition from  $j$  to  $i$  with  
the configuration

$$\beta(d_1, d_2, \dots, d_t)x = \sum_{r=1}^{2^t} c_r \sigma(v_r \cdot \mu)$$





$$q_i^+ := \left( \sum_{j=0}^s y_{ij}^1(a_1, \dots, a_p, b_1, \dots, b_p) x_j \right) \times q_i$$

$$+ \left( \sum_{j=0}^s y_{ij}^2(a_1, \dots, a_p, b_1, \dots, b_p) x_j \right)$$

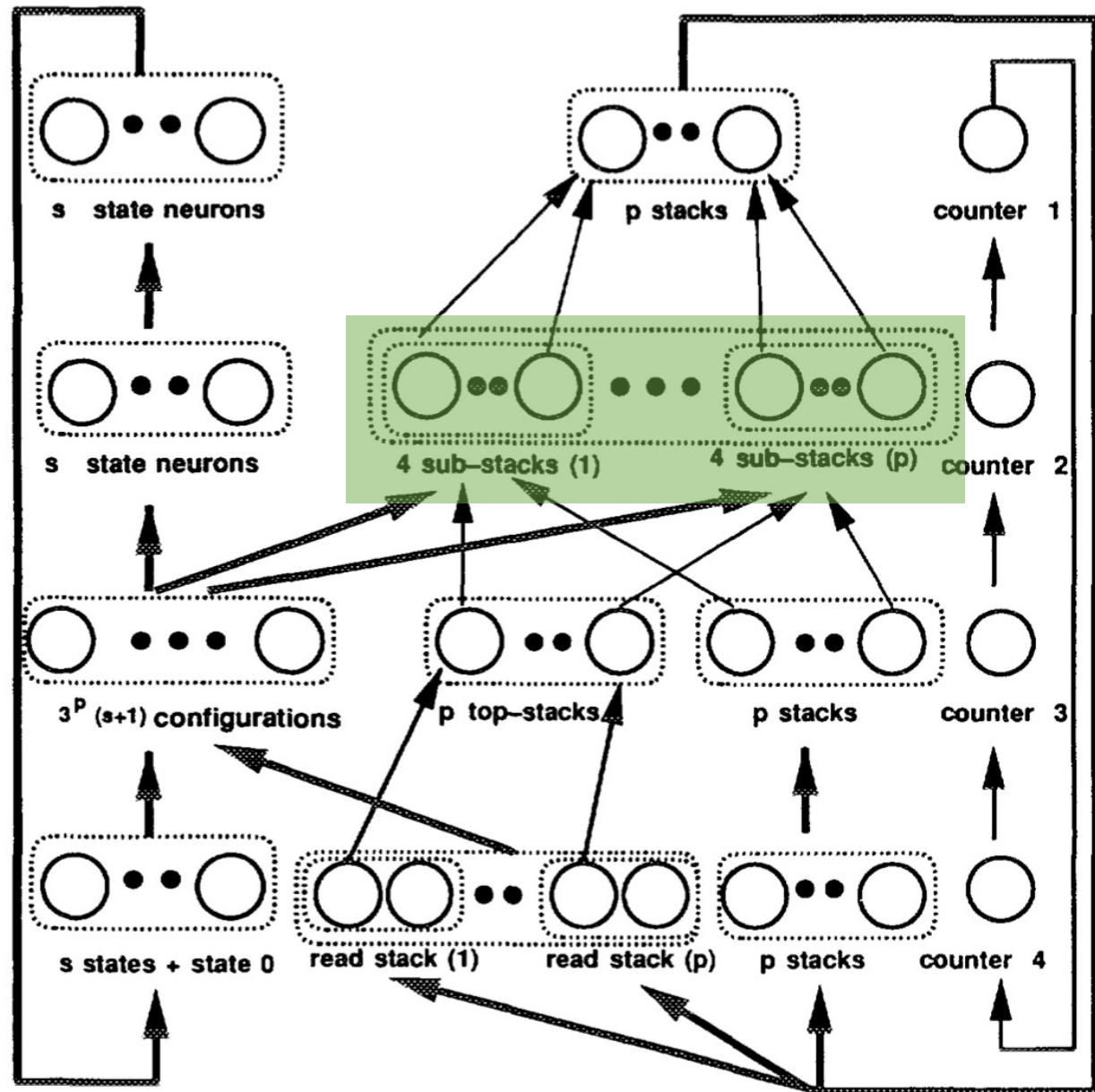
$$\times \left( \frac{1}{4} q_i + \frac{1}{4} \right)$$

$$+ \left( \sum_{j=0}^s y_{ij}^3(a_1, \dots, a_p, b_1, \dots, b_p) x_j \right)$$

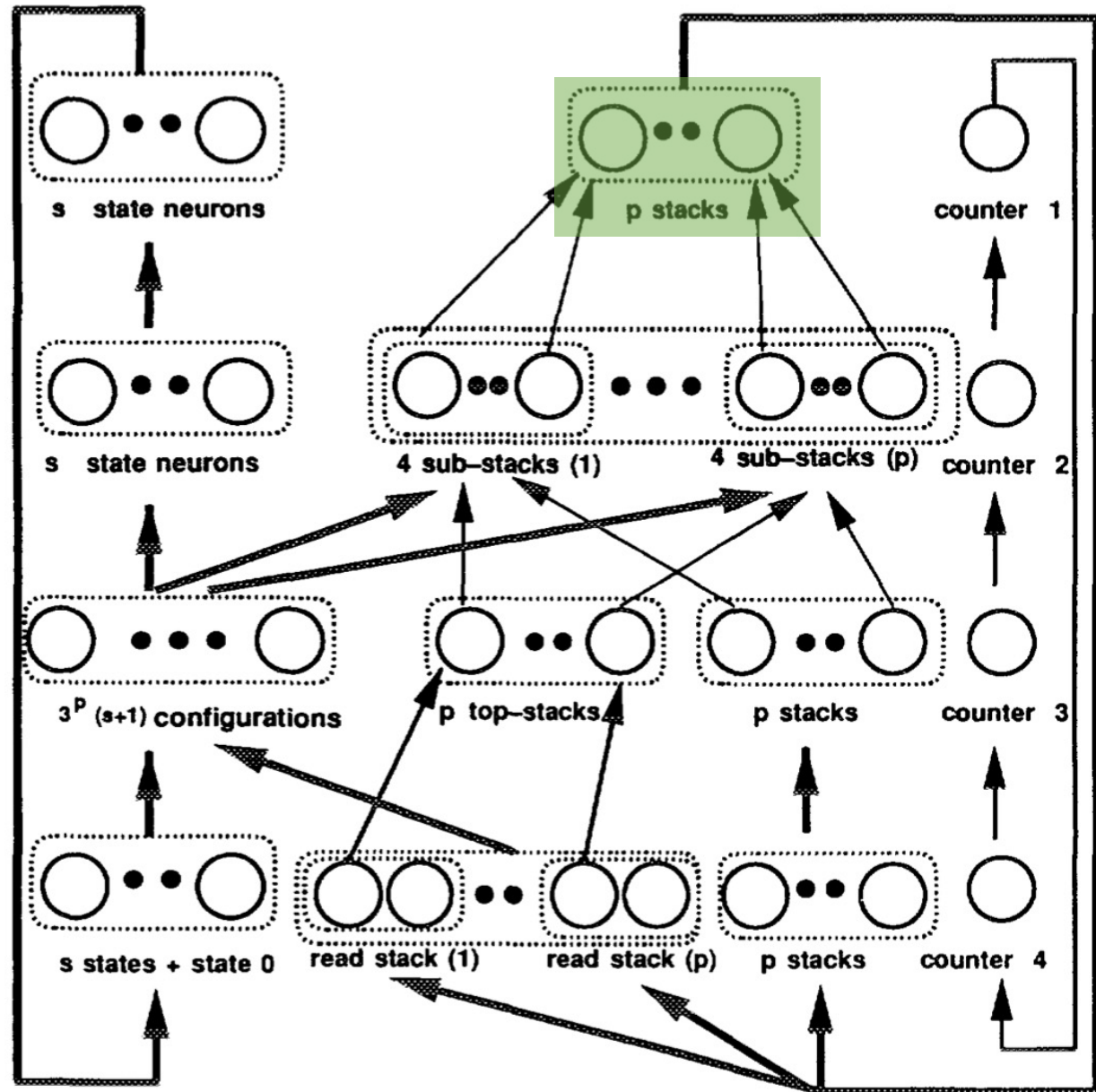
$$\times \left( \frac{1}{4} q_i + \frac{1}{4} \right)$$

$$+ \left( \sum_{j=0}^s y_{ij}^4(a_1, \dots, a_p, b_1, \dots, b_p) x_j \right)$$

$$\times (4q_i - 2\zeta[q_i] - 1)$$



$$\begin{aligned}
q_i^+ &:= \left( \sum_{j=0}^s y_{ij}^1(a_1, \dots, a_p, b_1, \dots, b_p) x_j \right) \times q_i \\
&+ \left( \sum_{j=0}^s y_{ij}^2(a_1, \dots, a_p, b_1, \dots, b_p) x_j \right) \\
&\times \left( \frac{1}{4} q_i + \frac{1}{4} \right) \\
&+ \left( \sum_{j=0}^s y_{ij}^3(a_1, \dots, a_p, b_1, \dots, b_p) x_j \right) \\
&\times \left( \frac{1}{4} q_i + \frac{1}{4} \right) \\
&+ \left( \sum_{j=0}^s y_{ij}^4(a_1, \dots, a_p, b_1, \dots, b_p) x_j \right) \\
&\times (4q_i - 2\zeta[q_i] - 1)
\end{aligned}$$



# Construction

- An RNN of 886 processors (nodes)
- Rational weights
- No slow down (real time simulation)

# Are Neural Networks more powerful than TMs?

H. T. Siegelmann and E. D. Sontag, 'On the computational power of neural nets', *JCSS*, 1995

H. T. Siegelmann and E. D. Sontag, 'Analog computation via neural networks', *Theoretical Computer Science*, vol. 131, 1994

J. L. Balcazar, R. Gavalda, and H. T. Siegelmann, 'Computational power of neural networks: a characterization in terms of Kolmogorov complexity', *IEEE Trans. Inform. Theory*, vol. 43, Jul. 1997

# Super-Turing computation

- NNs with rational weights = as powerful as Turing Machines
- NNs with real weights = more powerful than Turing Machines (Super-Turing!)
  - Can decide any language in exponential time
  - In polynomial time accept languages in P/poly

H. T. Siegelmann and E. D. Sontag, 'On the computational power of neural nets', *JCSS*, 1995

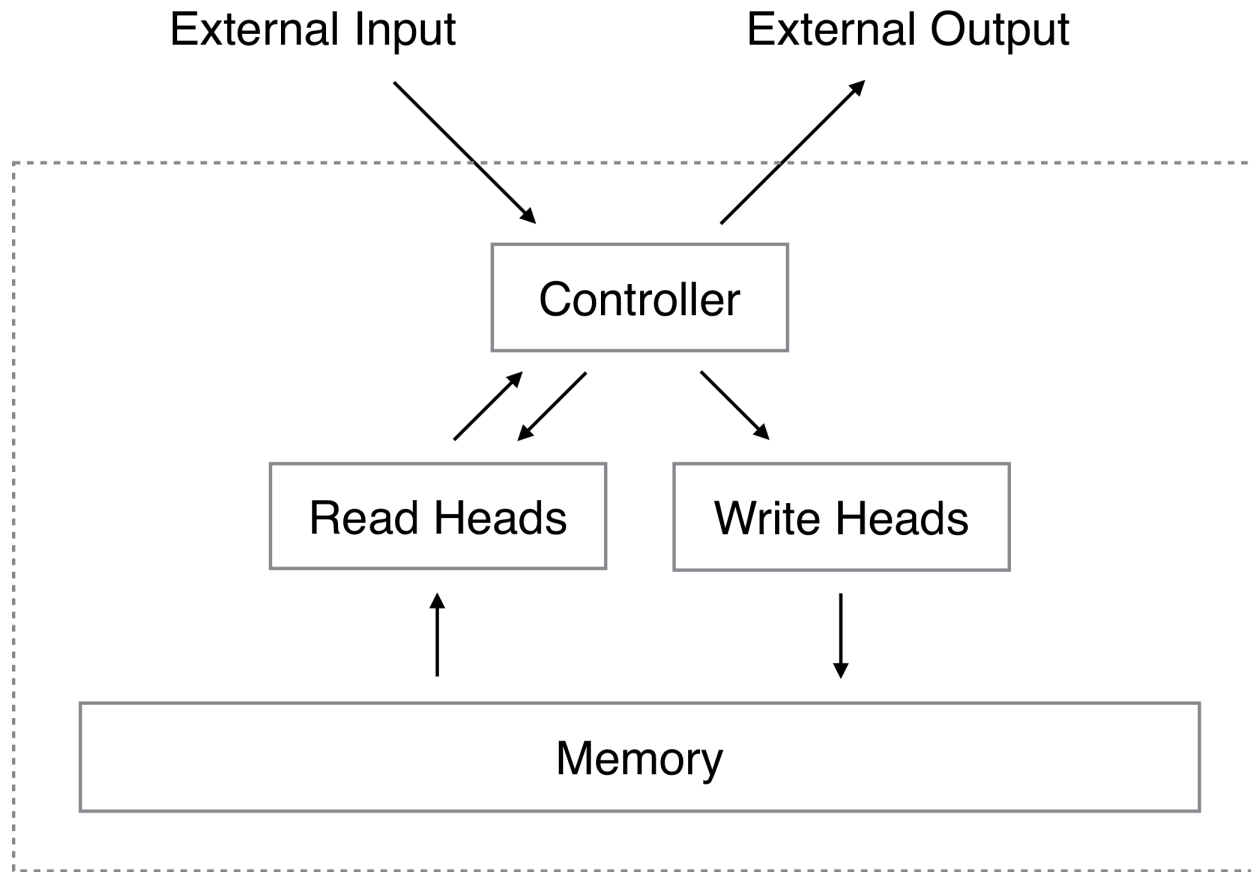
H. T. Siegelmann and E. D. Sontag, 'Analog computation via neural networks', *Theoretical Computer Science*, vol. 131, 1994

J. L. Balcazar, R. Gavalda, and H. T. Siegelmann, 'Computational power of neural networks: a characterization in terms of Kolmogorov complexity', *IEEE Trans. Inform. Theory*, vol. 43, Jul. 1997

# Mimicking Turing Machines

# Neural Turing Machines

# Neural Turing Machines





# Reading

read vector

$$\mathbf{r}_t \leftarrow \sum_i w_t(i) \mathbf{M}_t(i)$$

# Writing

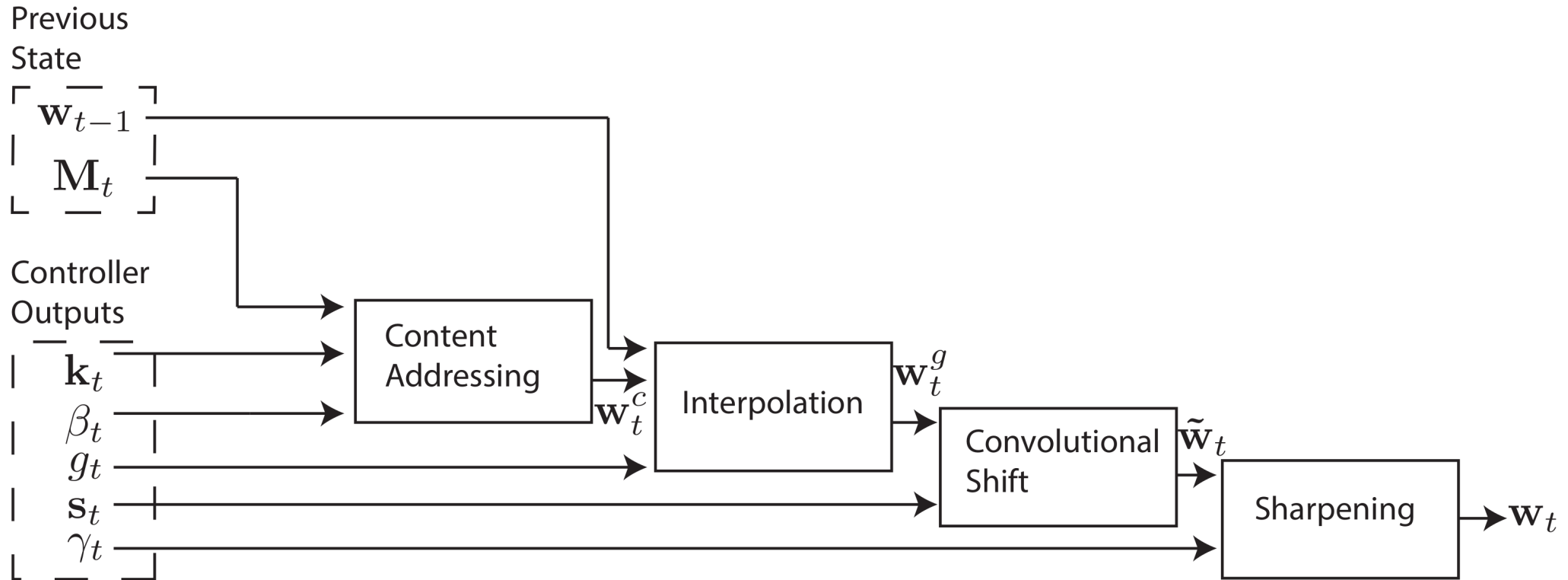
erase

$$\tilde{\mathbf{M}}_t(i) \leftarrow \mathbf{M}_{t-1}(i) [\mathbf{1} - w_t(i) \mathbf{e}_t]$$

add

$$\mathbf{M}_t(i) \leftarrow \tilde{\mathbf{M}}_t(i) + w_t(i) \mathbf{a}_t$$

# Addressing



# Focusing by Content

$$w_t^c(i) \leftarrow \frac{\exp(\beta_t K[\mathbf{k}_t, \mathbf{M}_t(i)])}{\sum_j \exp(\beta_t K[\mathbf{k}_t, \mathbf{M}_t(j)])}$$

# Focusing by Location

interpolation

$$\mathbf{w}_t^g \leftarrow g_t \mathbf{w}_t^c + (1 - g_t) \mathbf{w}_{t-1}$$

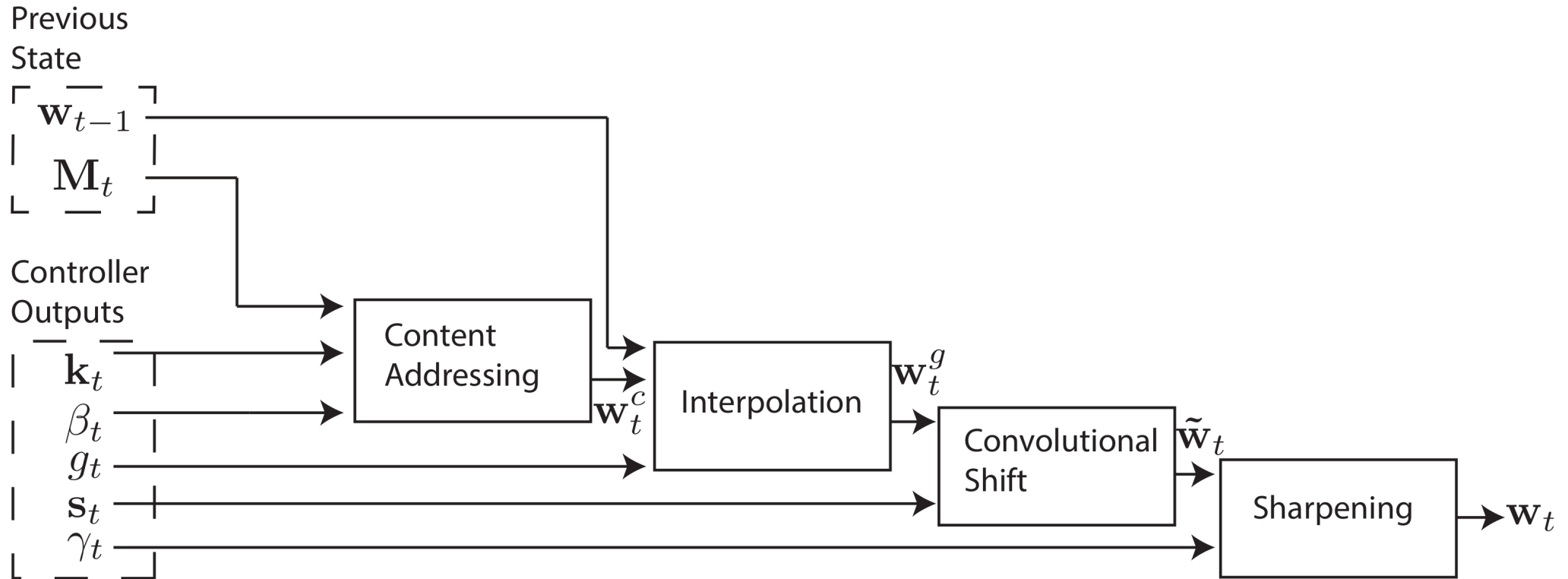
convolutional shift

$$\tilde{w}_t(i) \leftarrow \sum_{j=0}^{N-1} w_t^g(j) s_t(i - j)$$

sharpening

$$w_t(i) \leftarrow \frac{\tilde{w}_t(i)^{\gamma t}}{\sum_j \tilde{w}_t(j)^{\gamma t}}$$

# Addressing



# LSTM

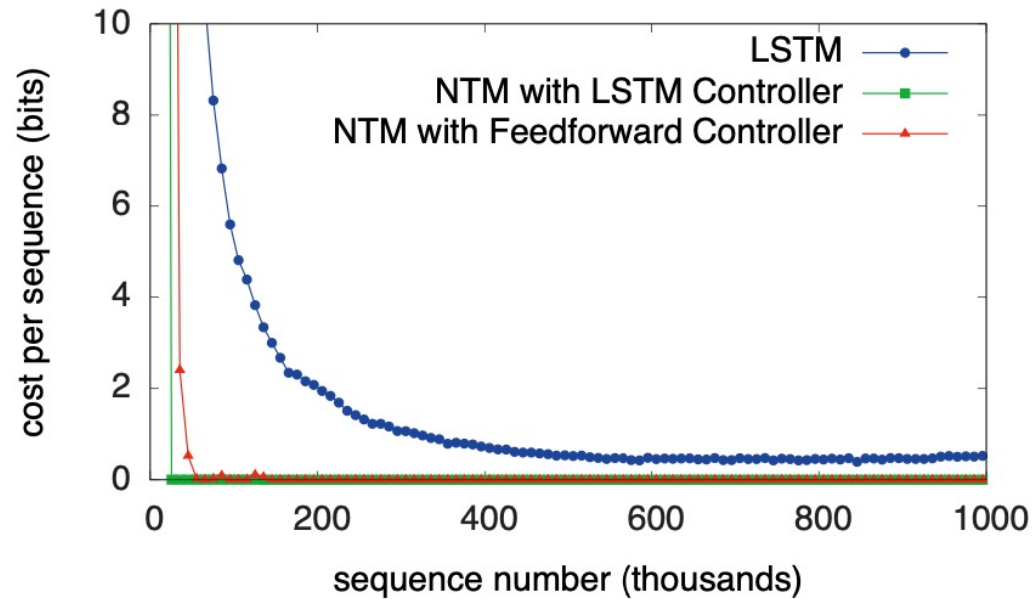
- Small internal memory
- Less interpretable

# Feed-forward

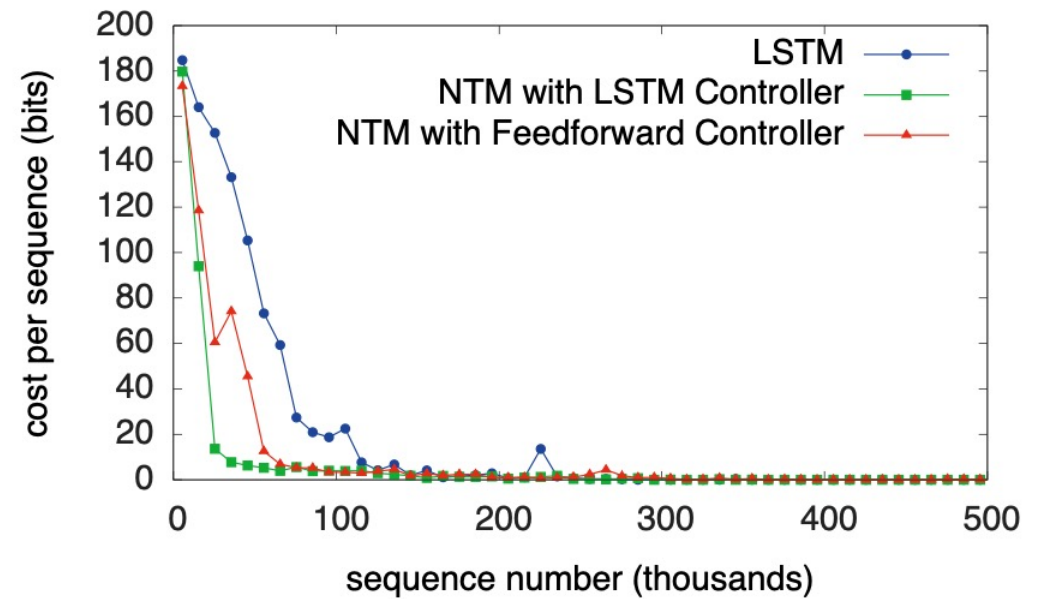
- No internal memory
- More interpretable

# Experiments

# Copy and Repeat Copy



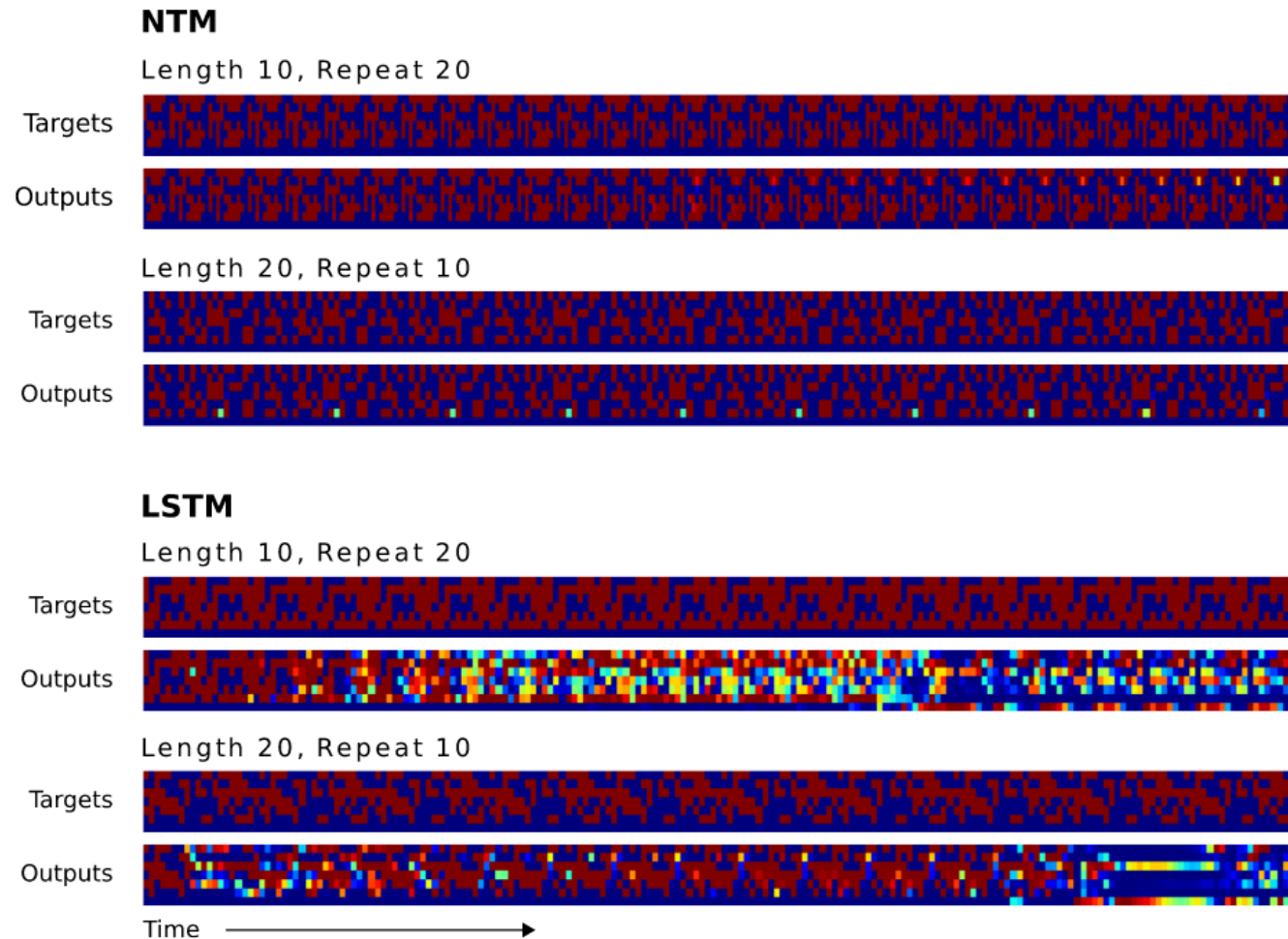
**copy**

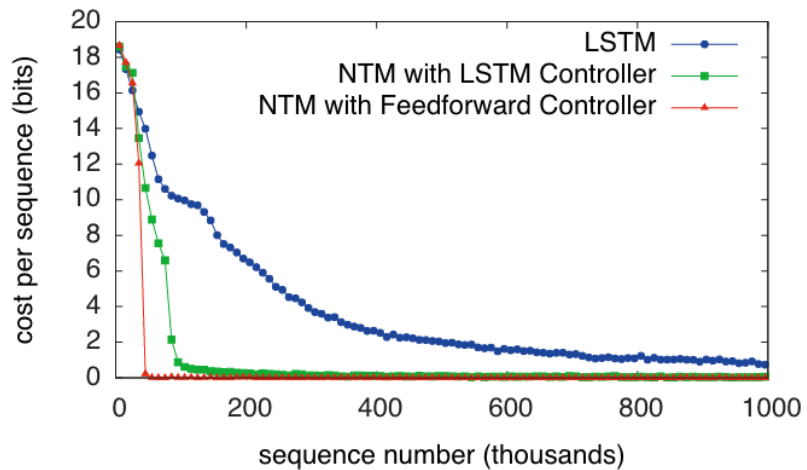


**repeat copy**

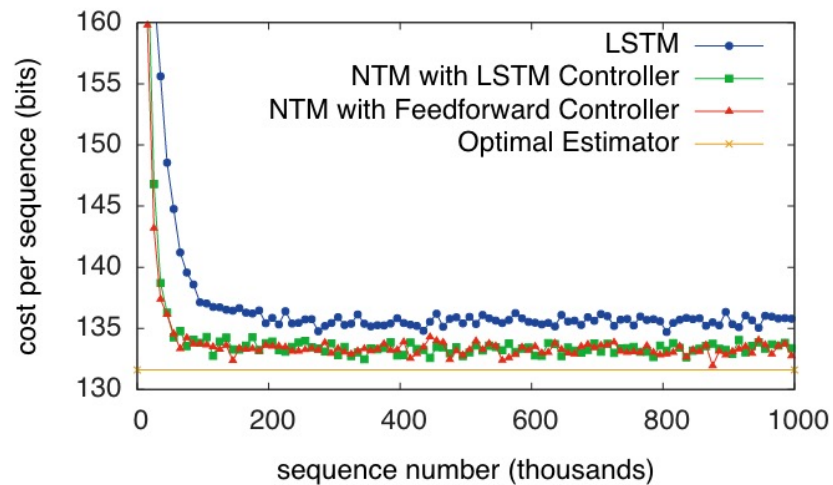


# Repeat Copy Generalisation

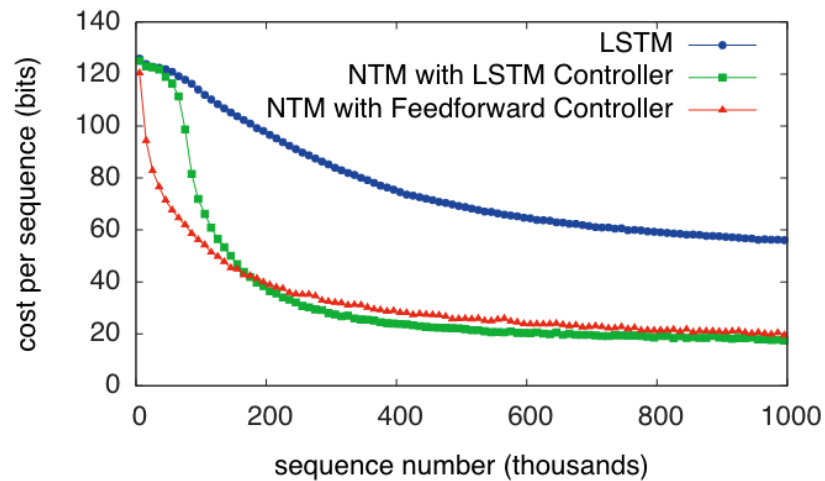




**associative recall**



**dynamic n-gram**



**priority sort**

# Memory Networks

# Components

- **$I$**  (input feature map)
  - convert  $x$  to an internal representation  $I(x)$
- **$G$**  (generalisation)
  - $m_i = G(m_i, I(x), m) \forall i$
- **$O$**  (output feature map)
  - $o = O(I(x), m)$
- **$R$**  (response)
  - $r = R(o)$

# Basic model

memories stored at the next available slot

$$\mathbf{m}_N = x, N = N + 1$$

$k$  relevant memories chosen by  $O$

$$o_1 = O_1(x, \mathbf{m}) = \operatorname{argmax}_{i=1, \dots, N} s_O(x, \mathbf{m}_i)$$

$$o_2 = O_2(x, \mathbf{m}) = \operatorname{argmax}_{i=1, \dots, N} s_O([x, \mathbf{m}_{o_1}], \mathbf{m}_i)$$

...

$$o_k = O_k(x, \mathbf{m}) = \operatorname{argmax}_{i=1, \dots, N} s_O([x, \mathbf{m}_{o_1}, \dots, \mathbf{m}_{o_k}], \mathbf{m}_i)$$

# Basic model

choose best word to output

$$r = \operatorname{argmax}_{w \in W} s_R([x, \mathbf{m}_{O_1}, \dots, \mathbf{m}_{O_k}], w)$$

# Training

minimise (for  $k = 2$ )

$$\begin{aligned} & \sum_{\bar{f} \neq \mathbf{m}_{o_1}} \max(0, \gamma - s_O(x, \mathbf{m}_{o_1}) + s_O(x, \bar{f})) \\ & + \sum_{\bar{f}' \neq \mathbf{m}_{o_2}} \max(0, \gamma - s_O([x, \mathbf{m}_{o_1}], \mathbf{m}_{o_2}) + s_O([x, \mathbf{m}_{o_1}], \bar{f}')) \\ & + \sum_{\bar{r} \neq r} \max(0, \gamma - s_R([x, \mathbf{m}_{o_1}, \mathbf{m}_{o_2}], r) + s_R([x, \mathbf{m}_{o_1}, \mathbf{m}_{o_2}], \bar{r})) \end{aligned}$$

# Extensions



# Efficient memory via hashing

# Efficient memory via hashing

The **capital** of **France** is **Paris**

# Efficient memory via hashing

1. Train embedding matrix  $U_O$
2. Use K-Means to cluster word vectors  $(U_O)_i$

# Modelling write time

$$s_{O_t}(x, y, y') = \Phi_x(x)^T U_{O_t}(\Phi_y(y) - \Phi_y(y') + \Phi_t(x, y, y'))$$

# Experiments

# Large-scale QA

- Open question answering (QA) on a 14M statement dataset (Fader et al. (2013))

Method	F1
(Fader et al., 2013)	0.54
(Bordes et al., 2014b)	0.73
MemNN (embedding only)	0.72
MemNN (with BoW features)	0.82

Method	Embedding F1	Embedding + BoW F1	Candidates (speedup)
MemNN (no hashing)	0.72	0.82	14M (0x)
MemNN (word hash)	0.63	0.68	13k (1000x)
MemNN (cluster hash)	0.71	0.80	177k (80x)

# Simulated world QA

Joe went to the kitchen. Fred went to the kitchen. Joe picked up the milk.

Joe travelled to the office. Joe left the milk. Joe went to the bathroom.

Where is the milk now? **A: office**

Where is Joe? **A: bathroom**

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

# Simulated world QA

Method	Difficulty 1			Difficulty 5	
	actor w/o before	actor	actor+object	actor	actor+object
RNN	100%	60.9%	27.9%	23.8%	17.8%
LSTM	100%	64.8%	49.1%	35.2%	29.0%
MemNN $k = 1$	97.8%	31.0%	24.0%	21.9%	18.5%
MemNN $k = 1$ (+time)	99.9%	60.2%	42.5%	60.8%	44.4%
MemNN $k = 2$ (+time)	100%	100%	100%	100%	99.9%





# References

1. A. M. Turing, 'Systems of logic based on ordinals', *Proceedings of the London Mathematical Society*, Series 2, 45, pp. 161-228, 1939
2. H. T. Siegelmann and E. D. Sontag, 'On the computational power of neural nets', *Journal of Computer and System Sciences*, 50(1), pp. 132-150, 1995, doi: [10.1006/jcss.1995.1013](https://doi.org/10.1006/jcss.1995.1013)
3. H. T. Siegelmann and E. D. Sontag, 'Analog computation via neural networks', *Theoretical Computer Science*, vol. 131, no. 2, pp. 331–360, Sep. 1994, doi: [10.1016/0304-3975\(94\)90178-3](https://doi.org/10.1016/0304-3975(94)90178-3)
4. J. L. Balcazar, R. Gavalda, and H. T. Siegelmann, 'Computational power of neural networks: a characterization in terms of Kolmogorov complexity', *IEEE Trans. Inform. Theory*, vol. 43, no. 4, pp. 1175–1183, Jul. 1997, doi: [10.1109/18.605580](https://doi.org/10.1109/18.605580)
5. J. Kilian and H. T. Siegelmann, 'The Dynamic Universality of Sigmoidal Neural Networks', *Information and Computation*, vol. 128, no. 1, pp. 48–56, Jul. 1996, doi: [10.1006/inco.1996.0062](https://doi.org/10.1006/inco.1996.0062)
6. S. Chung and H. T. Siegelmann, 'Turing Completeness of Bounded-Precision Recurrent Neural Networks', in *Advances in Neural Information Processing Systems*, 2021, vol. 34, pp. 28431–28441. Accessed: May 08, 2022. [Online]. Available: <https://proceedings.neurips.cc/paper/2021/hash/ef452c63f81d0105dd4486f775adec81-Abstract.html>
7. A. Graves, G. Wayne, and I. Danihelka, 'Neural Turing Machines', *arXiv:1410.5401 [cs]*, Dec. 2014, Accessed: Mar. 05, 2022. [Online]. Available: <http://arxiv.org/abs/1410.5401>
8. J. Weston, S. Chopra, and A. Bordes, 'Memory Networks', *arXiv:1410.3916 [cs, stat]*, Nov. 2015, Accessed: Mar. 05, 2022. [Online]. Available: <http://arxiv.org/abs/1410.3916>
9. A. Fader, L. Zettlemoyer, and O. Etzioni, 'Paraphrase-driven learning for open question answering', in *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1608-1618, Aug. 2013.