

# Differentiable Neural Computers

HYBRID COMPUTING USING A NEURAL NETWORK WITH  
DYNAMIC EXTERNAL MEMORY (GRAVES ET AL. 2016)

---

Konstantinos Kogkalidis

May 28, 2018

Logic and Computation

# Overview: Probabilistic Programming

## Cross-domain

- Data Flow Programming
- Bayesian Reasoning
- Machine Learning
- Functional Programming

# Overview: Probabilistic Programming

## Cross-domain

- Data Flow Programming
- Bayesian Reasoning
- Machine Learning
- Functional Programming

Intuition: *"Rather than explicitly write a program, write some **constraints** on the behavior of the desired program and use computational tools to search the program space for **models** satisfying these constraints."*

# Overview: Probabilistic Programming

## Cross-domain

- Data Flow Programming
- Bayesian Reasoning
- Machine Learning
- Functional Programming

Intuition: *"Rather than explicitly write a program, write some **constraints** on the behavior of the desired program and use computational tools to search the program space for **models** satisfying these constraints."*

PROGRAM	MODEL
Discrete	Continuous
Deterministic	Stochastic
Static	Adaptive

## Overview: DNC

### Differentiable Neural Computer

A recurrent neural network coupled with an external memory.

# Overview: DNC

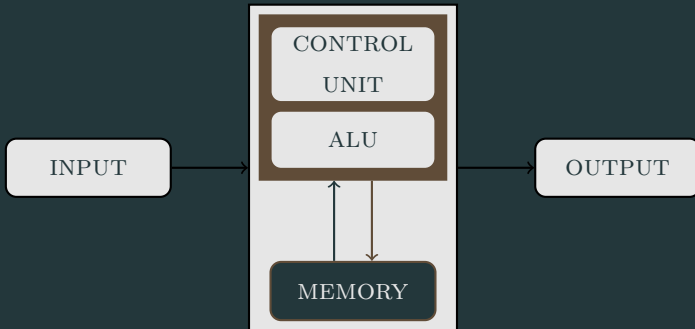
## Differentiable Neural Computer

A recurrent neural network coupled with an external memory.

- Inspired by biological memory
- Reimagining of classical concepts of computation
- Extension of NTMs
  - End-to-end differentiable
  - Auto-associative memory
  - Turing complete
- + Stronger memory management

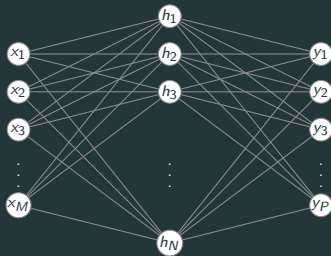
# Introduction: Classic Computation

## Von Neumann architecture



# Introduction: Neural Networks

## Simple Neural Net



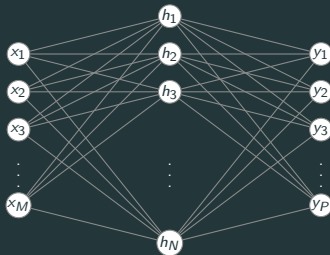
$$NN : \mathbf{x}_{t_i} \mapsto \mathbf{y}_{t_i}$$

No memory



# Introduction: Neural Networks

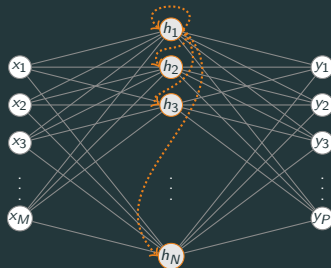
## Simple Neural Net



$$NN : \mathbf{x}_{t_i} \mapsto \mathbf{y}_{t_i}$$

No memory

## Simple Recurrent Net



$$RNN : \mathbf{x}_{t_0} \otimes \mathbf{x}_{t_1} \otimes \dots \otimes \mathbf{x}_{t_i} \mapsto \mathbf{y}_{t_i}$$

Finite memory

# Introduction: Neural Networks

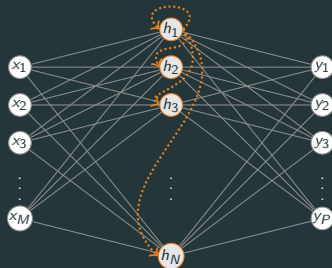
## Simple Neural Net



$$NN : \mathbf{x}_{t_i} \mapsto \mathbf{y}_{t_i}$$

No memory

## Simple Recurrent Net



$$RNN : \mathbf{x}_{t_0} \otimes \mathbf{x}_{t_1} \otimes \cdots \otimes \mathbf{x}_{t_i} \mapsto \mathbf{y}_{t_i}$$

Finite memory

*"If training vanilla neural nets is optimization over functions,  
training recurrent nets is **optimization over programs**."*

## Approach

Train a RNN to act as the **controller** of a memory matrix  $M$  of  $N$  addresses through  $R$  **read heads** and one **write head**.

# Approach

Train a RNN to act as the **controller** of a memory matrix  $M$  of  $N$  addresses through  $R$  **read heads** and one **write head**.

## 1. Content Lookup

- **Attention** over memory defined by weightings  $w \in \mathcal{S}^N$
- Compare controller output with memory objects (**auto-associative memory**)
- Allow partial matches (**pattern completion**)

# Approach

Train a RNN to act as the **controller** of a memory matrix  $M$  of  $N$  addresses through  $R$  **read heads** and one **write head**.

## 1. Content Lookup

- **Attention** over memory defined by weightings  $w \in \mathcal{S}^N$
- Compare controller output with memory objects (**auto-associative memory**)
- Allow partial matches (**pattern completion**)

## 2. Sequential Retrieval

- Fill  $L \in [0, 1]^{N \times N}$  indexing **temporal transitions**
- **Shift** operations defined by  $Lw$ ,  $L^\top w$

# Approach

Train a RNN to act as the **controller** of a memory matrix  $M$  of  $N$  addresses through  $R$  **read heads** and one **write head**.

## 1. Content Lookup

- **Attention** over memory defined by weightings  $w \in \mathcal{S}^N$
- Compare controller output with memory objects (**auto-associative memory**)
- Allow partial matches (**pattern completion**)

## 2. Sequential Retrieval

- Fill  $L \in [0, 1]^{N \times N}$  indexing **temporal transitions**
- **Shift** operations defined by  $Lw$ ,  $L^T w$

## 3. Dynamic Allocation

- Mark memory locations with  $\{0, 1\}$  to **signal usage**
- Manipulate signals during R/W operations to enable **reallocation**
- Generalization to **unbounded memory**

## Controller: Overview

A deep long short-term memory network receiving input:

$$\mathbf{x}_t = [\mathbf{x}_t; \mathbf{r}_{t-1}^1; \dots; \mathbf{r}_{t-1}^R] \quad (\text{timestep } t)$$

and producing output:

$$(\mathbf{v}_T, \boldsymbol{\xi}_T) = \mathcal{N}([\mathbf{x}_1; \dots; \mathbf{x}_T]; \vartheta) \quad (\text{entire sequence})$$

where  $\mathcal{N}$  a set of state equations and  $\vartheta$  their trainable parameters.

## Controller: State Equations

A more detailed look into  $\mathcal{N}$ .



## Controller: State Equations

A more detailed look into  $\mathcal{N}$ .

LSTM layer equations

Input:  $[\mathbf{x}_t; \mathbf{h}'_{t-1}; \mathbf{h}'_{t-1}]$

Output:  $\mathbf{h}'_t$

## Controller: State Equations

A more detailed look into  $\mathcal{N}$ .

### LSTM layer equations

Input:  $[\mathcal{X}_t; \mathbf{h}'_{t-1}; \mathbf{h}'_{t-1}]$

Output:  $\mathbf{h}'_t$

$$\mathbf{i}'_t = \zeta_i([\mathcal{X}_t; \mathbf{h}'_{t-1}; \mathbf{h}'_{t-1}]) \quad (\text{input gate})$$

$$\mathbf{f}'_t = \zeta_f([\mathcal{X}_t; \mathbf{h}'_{t-1}; \mathbf{h}'_{t-1}]) \quad (\text{forget gate})$$

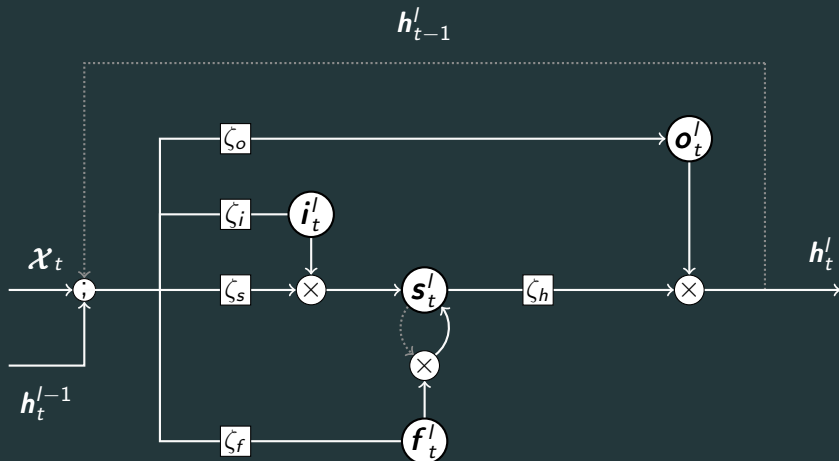
$$\mathbf{s}'_t = \mathbf{f}'_t \mathbf{s}'_{t-1} + \mathbf{i}'_t \zeta_s([\mathcal{X}_t; \mathbf{h}'_{t-1}; \mathbf{h}'_{t-1}]) \quad (\text{state update})$$

$$\mathbf{o}'_t = \zeta_o([\mathcal{X}_t; \mathbf{h}'_{t-1}; \mathbf{h}'_{t-1}]) \quad (\text{output gate})$$

$$\mathbf{h}'_t = \mathbf{o}'_t \zeta_h(\mathbf{s}'_t) \quad (\text{hidden layer})$$

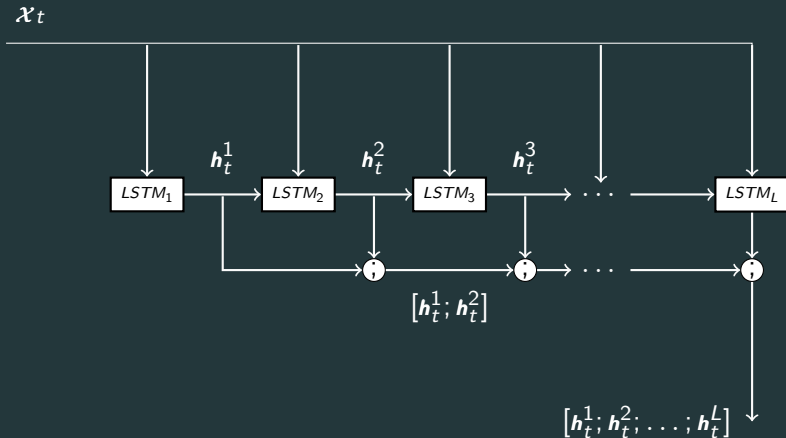
## Controller: Signal-Flow (1/2)

### Single LSTM layer



## Controller: Signal-Flow (2/2)

### LSTM Network (multiple layers)



## Controller: Outputs

$$(\boldsymbol{v}_t, \boldsymbol{\xi}_t) = \mathcal{N}([\boldsymbol{x}_1; \dots; \boldsymbol{x}_T]; \boldsymbol{\vartheta})$$

## Controller: Outputs

$$(\boldsymbol{v}_t, \boldsymbol{\xi}_t) = \mathcal{N}([\boldsymbol{x}_1; \dots; \boldsymbol{x}_T]; \boldsymbol{\vartheta})$$

Intermediate output  $\boldsymbol{v}_t = W_y[\boldsymbol{h}_t^1; \dots; \boldsymbol{h}_t^L]$

$$\boldsymbol{y}_t = \boldsymbol{v}_t + W_R[\boldsymbol{r}_t^1; \dots; \boldsymbol{r}_t^R] \quad (\text{Memory-conditioning})$$

## Controller: Outputs

$$(\mathbf{v}_t, \xi_t) = \mathcal{N}([\mathbf{x}_1; \dots; \mathbf{x}_T]; \vartheta)$$

Intermediate output  $\mathbf{v}_t = W_y[\mathbf{h}_t^1; \dots; \mathbf{h}_t^L]$

$$\mathbf{y}_t = \mathbf{v}_t + W_R[\mathbf{r}_t^1; \dots; \mathbf{r}_t^R] \quad (\text{Memory-conditioning})$$

Interface vector  $\xi_t = W_\xi[\mathbf{h}_t^1; \dots; \mathbf{h}_t^L]$

- Read keys:  $\mathbf{k}_t^{r,i}$
- Read strengths:  $\beta_t^{r,i}$
- Write key:  $\mathbf{k}_t^w$
- Write strength:  $\beta_t^w$
- Erase vector:  $\mathbf{e}_t$
- Write vector:  $\mathbf{v}_t$
- Free gates:  $\phi_t^i$
- Allocation gate:  $g_t^a$
- Write gate:  $g_t^w$
- Read modes:  $\pi_t^i$

## Memory Addressing: Content-Lookup

R read keys  $\mathbf{k}^{r,i} \in \mathbb{R}^W$ ,  $i = 1 \dots R$

R read strengths  $\beta^{r,i} \in [1, \infty)$ ,  $i = 1 \dots R$

Write key  $\mathbf{k}^w \in \mathbb{R}^W$

Write strength  $\beta^w \in [1, \infty)$



## Memory Addressing: Content-Lookup

R read keys  $\mathbf{k}^{r,i} \in \mathbb{R}^W$ ,  $i = 1 \dots R$

R read strengths  $\beta^{r,i} \in [1, \infty)$ ,  $i = 1 \dots R$

Write key  $\mathbf{k}^w \in \mathbb{R}^W$

Write strength  $\beta^w \in [1, \infty)$

Matching function  $\mathcal{D}$  comparing memory contents

## Memory Addressing: Content-Lookup

R read keys  $\mathbf{k}^{r,i} \in \mathbb{R}^W$ ,  $i = 1 \dots R$

R read strengths  $\beta^{r,i} \in [1, \infty)$ ,  $i = 1 \dots R$

Write key  $\mathbf{k}^w \in \mathbb{R}^W$

Write strength  $\beta^w \in [1, \infty)$

Matching function  $\mathcal{D}$  comparing memory contents

Weighting function  $\mathcal{C}$  normalizing and sharpening matches

$$\mathcal{C}(M, \mathbf{k}, \beta)[i] = \frac{\exp\{\beta \mathcal{D}(\mathbf{k}, M[i, :])\}}{\sum_j \exp\{\beta \mathcal{D}(\mathbf{k}, M[j, :])\}}$$

## Memory Addressing: R/W

Attention dictated by weightings  $\mathbf{w} \in \mathcal{S}^N$

Erase vector  $\mathbf{e}_t \in [0, 1]^W$

Write vector  $\mathbf{v}_t \in \mathbb{R}^W$

Read operations

$$\mathbf{r}_t^i = M_t^\top \mathbf{w}_t^{r,i}$$

# Memory Addressing: R/W

Attention dictated by weightings  $\mathbf{w} \in \mathcal{S}^N$

Erase vector  $\mathbf{e}_t \in [0, 1]^W$

Write vector  $\mathbf{v}_t \in \mathbb{R}^W$

Read operations

$$r_t^i = M_t^\top \mathbf{w}_t^{r,i}$$

Write operations

$$M_t = \underbrace{M_{t-1} \circ (\mathbf{1} - \mathbf{w}_t^w \mathbf{e}_t^\top)}_{\text{erased memory}} + \underbrace{\mathbf{w}_t^w \mathbf{v}_t^\top}_{\text{new write}}$$

# Memory Addressing: Dynamic Allocation

Free gates  $\phi_t^i \in [0, 1]^W$

Allocation gate  $g_t^a \in [0, 1]$

Write gate  $g_t^w \in [0, 1]$

"Free list" scheme

$$\psi_t = \prod_{i=1}^R (1 - \phi_t^i w_{t-1}^{r,i}) \quad (\text{memory retention})$$

$$u_t = (\mathbf{u}_{t-1} + \mathbf{w}_{t-1}^w + \mathbf{u}_{t-1} \circ \mathbf{w}_{t-1}^w) \circ \psi_t \quad (\text{usage tracking})$$

# Memory Addressing: Dynamic Allocation

Free gates  $\phi_t^i \in [0, 1]^W$

Allocation gate  $g_t^a \in [0, 1]$

Write gate  $g_t^w \in [0, 1]$

"Free list" scheme

$$\psi_t = \prod_{i=1}^R (1 - \phi_t^i w_{t-1}^{r,i}) \quad (\text{memory retention})$$

$$u_t = (\mathbf{u}_{t-1} + \mathbf{w}_{t-1}^w + \mathbf{u}_{t-1} \circ \mathbf{w}_{t-1}^w) \circ \psi_t \quad (\text{usage tracking})$$

Attention shift

- Obtain the allocation vector  $\mathbf{a}_t$  by normalizing  $\mathbf{u}_t$
- Shift  $\mathbf{w}_t$  by  $g_t^a \mathbf{a}_t$  and scale by  $g_t^w$

## Memory Addressing: Temporal Linking

Free gates  $\phi_t^i \in [0, 1]^W$

Memory retention

$$\psi_t = \prod_{i=1}^R (1 - \phi_t^i w_{t-1}^{r,i})$$

# Memory Addressing: Temporal Linking

Free gates  $\phi_t^i \in [0, 1]^W$

Memory retention

$$\psi_t = \prod_{i=1}^R (1 - \phi_t^i w_{t-1}^{r,i})$$

Usage tracking

$$u_t = (\mathbf{u}_{t-1} + \mathbf{w}_{t-1}^w + \mathbf{u}_{t-1} \circ \mathbf{w}_{t-1}^w) \circ \psi_t$$



# Memory Addressing: Temporal Linking

Free gates  $\phi_t^i \in [0, 1]^W$

Memory retention

$$\psi_t = \prod_{i=1}^R (1 - \phi_t^i w_{t-1}^{r,i})$$

Usage tracking

$$u_t = (\mathbf{u}_{t-1} + \mathbf{w}_{t-1}^w + \mathbf{u}_{t-1} \circ \mathbf{w}_{t-1}^w) \circ \psi_t$$

Allocation vector  $\mathbf{a}_t$  given by normalizing and sorting  $\mathbf{u}_t$

## Further Reading

- [Neural Turing Machines](#) (Graves, Wayne, Danihelka)
- [Entity Networks](#) (Henaff, Weston, Szlam, Bordes, LeCun)
- [End-to-End Memory Networks](#) (Sukhbaatar, Szlam, Weston, Fergus)
- [Jointly Learning to Align and Translate](#) (Bahdanau, Cho, Bengio)
- [Principles of Probabilistic Programming Languages](#) (Goodman)
- [Backprop as a Functor](#) (Fong, Spivak, Tuyras)
- [Formal Methods for Probabilistic Programming](#) (Selsam, Liang, Dill)

# Conclusion

## Takeaway

We can automatically infer simple functions over complex data structures, in the form of **probability distributions** just by using examples.

# Conclusion

## Takeaway

We can automatically infer simple functions over complex data structures, in the form of **probability distributions** just by using examples.

Thank you!