**Lecture 11: Sequence-to-sequence architectures. Neural attention and memory.**
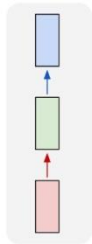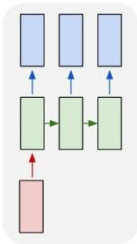
# Learning settings

One-to-one: image to class label
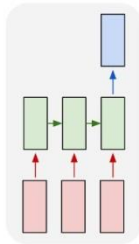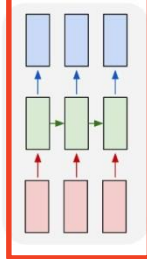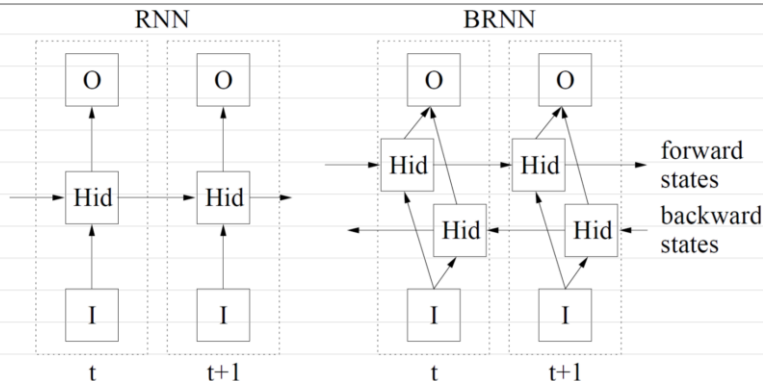One-to-many: text generation/image captioning
Many-to-one: sentiment analysis
Many-to-many 1: machine translation
Many-to-many 2: online classification (e.g. POS tagging)

# Bi-directional RNN



RNN           BRNN

forward states

backward states

**for** $t = 1$ to $T$ **do**
   Do forward pass for the forward hidden layer, storing activations at
   each timestep
**for** $t = T$ to $1$ **do**
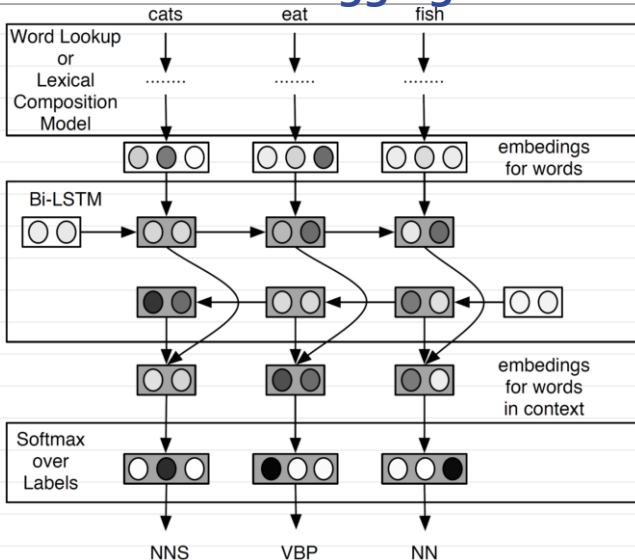   Do forward pass for the backward hidden layer, storing activations at
   each timestep
**for** $t = 1$ to $T$ **do**
   Do forward pass for the output layer, using the stored activations from
   both hidden layers

[A Graves, PhD thesis]

# Bi-LSTM POS tagging



$$\mathbf{l}_i = \tanh(\mathbf{L}^f \mathbf{s}_i^f + \mathbf{L}^b \mathbf{s}_i^b + \mathbf{b}_l)$$

[Ling et al. EMNLP15]

# Bi-LSTM POS tagging

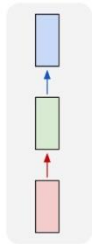| | acc | parameters | words/sec |
|---|---|---|---|
| Word Lookup | 96.97 | 2000k | 6K |
| Convolutional (S&Z) | 96.80 | 42.5k | 4K |
| Forward RNN | 95.66 | 17.5k | 4K |
| Backward RNN | 95.52 | 17.5k | 4K |
| Bi-RNN | 95.93 | 40k | 3K |
| Forward LSTM | 97.12 | 80k | 3K |
| Backward LSTM | 97.08 | 80k | 3K |
| Bi-LSTM $d_{CS} = 50$ | 97.22 | 70k | 3K |
| Bi-LSTM | **97.36** | 150k | 2K |

[Ling et al. EMNLP15]

## Uni-directional vs bi-directional

- Bi-directional is not applicable when "future" is unavailable
- When future is available bi-directional is almost always better
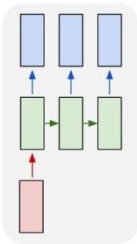- E.g. NLP (batch mode), bioinformatics

# Learning settings

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

One-to-one: image to class label
One-to-many: text generation/image captioning
Many-to-one: sentiment analysis
Many-to-many 1: machine translation
Many-to-many 2: online classification (e.g. POS tagging)

# Online seq2seq with monotonic alignment

Many problems are sequence 2 sequence with monotonic alignment:

- Not one-to-one as sequence prediction or POS tagging
- More constrained than general seq2seq



[Graves et al. 2006]

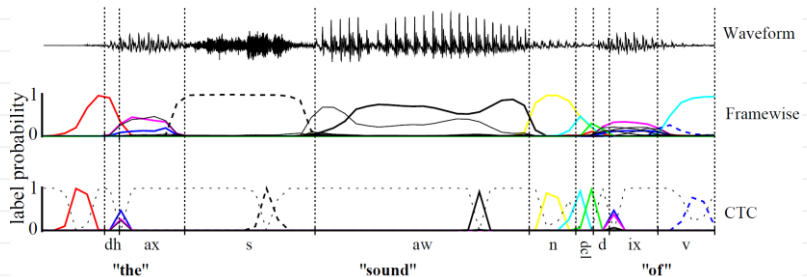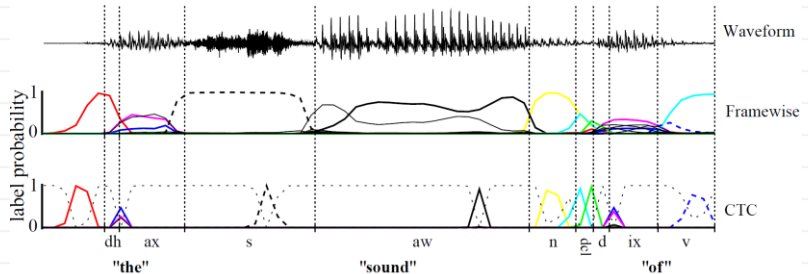# Online seq2seq with monotonic alignment



Decoding: 'aaa__bb_c____ddaa' → abcda

What should be the loss that encourage correct parsing?

Answer: connectionist temporal classification (CTC) loss

[Graves et al. 2006]

# CTC-loss



- Augment the output state with *blank*
- Predict probabilities of each symbol (inc. blank) at each time moment
- Compute the probability of each lattice vertex under correct paths using forward-backward
- Push log-probabilities up (*ML training*) proportionally to the current probability

[Graves et al. 2006]

# Evolution of the CTC signal

GT sequence:



Prediction      Gradient w.r.t. prediction

[Graves et al. 2006]

# LSTM demo: handwriting recognition



LSTM RNN Demo by Nikhil Buduma:

https://www.youtube.com/watch?v=mLxsbWAYIpw

# Image/video captioning



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

## Activity Recognition

Input: Sequence of Frames

CNN  CNN

LSTM

Output: Label

Apply Eye Makeup

## Image Description

Input: Image

CNN

LSTM

Output: Sentence

A large building with a clock on the front of it

## Video Description

Input: Video

CRF

LSTM

Output: Sentence

A man juiced the orange

[Donahue et al. 2015]

# Image/video captioning



[Donahue et al. 2015]*

# Image/video captioning



Single Layer
$\mathbf{LRCN}_{1u}$

Two Layers, Unfactored
$\mathbf{LRCN}_{2u}$

Two Layers, Factored
$\mathbf{LRCN}_{2f}$

- Train on 108,000 images with descriptions
- Test on 1000 images (5 descr per image)
- For each image score 5000 descriptions
- See if top-k has a correct description:

|  | R@1 | R@5 | R@10 | Med$r$ |
|---|---|---|---|---|
| $\mathbf{LRCN}_{1u}$ | 14.1 | 31.3 | 39.7 | 24 |
| $\mathbf{LRCN}_{2u}$ | 3.8 | 12.0 | 17.9 | 80 |
| $\mathbf{LRCN}_{2f}$ | **17.5** | **40.3** | **50.8** | **9** |
| $\mathbf{LRCN}_{4f}$ | 15.8 | 37.1 | 49.5 | 10 |

[Donahue et al. 2015]

# Image/video captioning

## Best results:



A female tennis player in action on the court.

A group of young men playing a game of soccer

A man riding a wave on top of a surf-board.

A baseball game in progress with the batter up to plate.

A brown bear standing on top of a lush green field.

A person holding a cell phone in their hand.

[Donahue et al. 2015]

# End-to-end dense image captioning



VGG-16 conv    fully-connected, 2 layers+dropout

Image:
3 x W x H

Conv features:
C x W' x H'

Region features:
B x C x X x Y

Region Codes:
B x D

CNN

Recognition Network

LSTM

Striped gray cat

Cats watching TV

Localization Layer

Conv features:
C x W' x H'

Conv

Region Proposals:
4k x W' x H'

Region scores:
k x W' x H'

Sampling

Best Proposals:
B x 4

Grid Generator

Sampling Grid:
B x X x Y x 2

Bilinear Sampler

Region features:
B x 512 x 7 x 7

k-anchors at W'xH' positions

[Johnson et al, CVPR16]

# End-to-end dense image captioning



[Johnson et al, CVPR16]

# End-to-end dense image captioning



[Johnson et al, CVPR16]

# Training set: "visual genome"



- "New Image-net"

108,249 Images
4.2 Million Region Descriptions
1.7 Million Visual Question Answers
2.1 Million Object Instances
1.8 Million Attributes
1.8 Million Relationships
Everything Mapped to Wordnet Synsets

[Krishna et al. 2016]

# Learning settings

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

*aka* "seq2seq"

# Sequence-to-sequence machine translation



Important notes:

1. Fixed lexicon (160,000 English, 80,000 French) + 'UNK' word
2. Deep (four layers, 1000 cells in each)
3. Reversing input sequence helps a lot
4. Using two different LSTMs
5. Decoding proceeds by *beam search*

[Sutskever et al. NIPS14]

# Sequence-to-sequence machine translation

[Sutskever et al. NIPS14]



Decoding proceeds by *beam search:*

1. At the first step generate top-K words
2. At each step, expand each of the K in top-L ways (gives KL results)
3. Pick the best K out of KL results

NB: needs some mechanism to compare sequences of different lengths

# Sequence-to-sequence machine translation

Learned embeddings:



PCA 1000-> 2

[Sutskever et al. NIPS14]

# Sequence-to-sequence machine translation

| Type | Sentence |
|---|---|
| Our model | Ulrich UNK , membre du conseil d' administration du constructeur automobile Audi , affirme qu' il s' agit d' une pratique courante depuis des années pour que les téléphones portables puissent être collectés avant les réunions du conseil d' administration afin qu' ils ne soient pas utilisés comme appareils d' écoute à distance . |
| Truth | Ulrich Hackenberg , membre du conseil d' administration du constructeur automobile Audi , déclare que la collecte des téléphones portables avant les réunions du conseil , afin qu' ils ne puissent pas être utilisés comme appareils d' écoute à distance , est une pratique courante depuis des années . |
| Our model | " Les téléphones cellulaires , qui sont vraiment une question , non seulement parce qu' ils pourraient potentiellement causer des interférences avec les appareils de navigation , mais nous savons , selon la FCC , qu' ils pourraient interférer avec les tours de téléphone cellulaire lorsqu' ils sont dans l' air " , dit UNK . |
| Truth | " Les téléphones portables sont véritablement un problème , non seulement parce qu' ils pourraient éventuellement créer des interférences avec les instruments de navigation , mais parce que nous savons , d' après la FCC , qu' ils pourraient perturber les antennes-relais de téléphonie mobile s' ils sont utilisés à bord " , a déclaré Rosenker . |
| Our model | Avec la crémation , il y a un " sentiment de violence contre le corps d' un être cher " , qui sera " réduit à une pile de cendres " en très peu de temps au lieu d' un processus de décomposition " qui accompagnera les étapes du deuil " . |
| Truth | Il y a , avec la crémation , " une violence faite au corps aimé " , qui va être " réduit à un tas de cendres " en très peu de temps , et non après un processus de décomposition , qui " accompagnerait les phases du deuil " . |

[Sutskever et al. NIPS14]
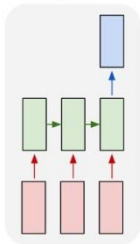
# Sequence-to-sequence machine translation



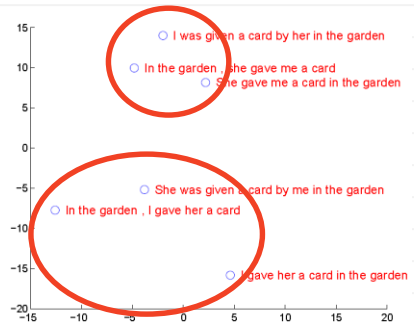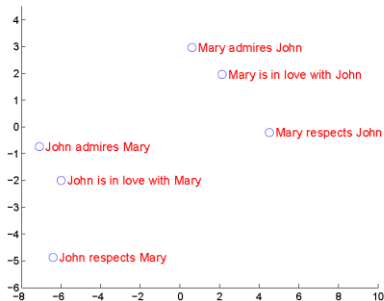| one to one | one to many | many to one | many to many | many to many |

[Sutskever et al. NIPS14]

# Sequence-to-sequence machine translation



Problem:
all the meaning has to be carried from here

- Large memory needed
- Information has to survive
  for a very long time

[Sutskever et al. NIPS14]

## Translation with attention



decoder RNN

$$s_i = f(s_{i-1}, y_{i-1}, c_i)$$

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

$$e_{ij} = \boxed{a}(s_{i-1}, h_j)$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

encoder RNN 1
encoder RNN 2

[Bahdanau et al. 2015]

# Translation with attention



decoder RNN

$$e_{ij} = a(s_{i-1}, h_j)$$

- *Attention model:* feed-forward neural network
- All components are trained end-to-end

encoder RNN 1
encoder RNN 2

[Bahdanau et al. 2015]

# Translation with attention



[Bahdanau et al. 2015]

# Translation with attention



- BLEU-score ≈ precision over n-grams
- Trained either with <30 word phrases or with <50 word phrases

[Bahdanau et al. 2015]

# Translation with attention

*An admitting privilege is the right of a doctor to admit a patient to a hospital or a medical centre <u>to carry out a diagnosis or a procedure, based on his status as a health care worker at a hospital.</u>*

## LSTM system:

*Un privilège d'admission est le droit d'un médecin de reconnaître un patient à l'hôpital ou un centre médical <u>d'un diagnostic ou de prendre un diagnostic en fonction de son état de santé.</u>*

## Attention-based system:

*Un privilège d'admission est le droit d'un médecin d'admettre un patient à un hôpital ou un centre médical <u>pour effectuer un diagnostic ou une procédure, selon son statut de travailleur des soins de santé à l'hôpital.</u>*

[Bahdanau et al. 2015]

# Simpler translation with attention

$$c_t = \sum_s a_t(s)\bar{h}_s$$



*Attention Layer*

$c_t$ — Context vector

Global align weights

$a_t$

$\bar{h}_s$

$h_t$

$y_t$

$\tilde{h}_t$

$$\tilde{h}_t = \tanh(\boldsymbol{W_c}[\boldsymbol{c_t}; \boldsymbol{h_t}])$$

$$p(y_t | y_{<t}, x) = \text{softmax}(\boldsymbol{W_s}\tilde{\boldsymbol{h}}_t)$$

$$\boldsymbol{a_t}(s) = \text{align}(\boldsymbol{h_t}, \bar{\boldsymbol{h}}_s)$$

$$= \frac{\exp\left(\text{score}(\boldsymbol{h_t}, \bar{\boldsymbol{h}}_s)\right)}{\sum_{s'} \exp\left(\text{score}(\boldsymbol{h_t}, \bar{\boldsymbol{h}}_{s'})\right)}$$

## Quadratic complexity!

$$\text{score}(\boldsymbol{h_t}, \bar{\boldsymbol{h}}_s) = \begin{cases} \boldsymbol{h}_t^\top \bar{\boldsymbol{h}}_s & \textit{dot} \\ \boldsymbol{h}_t^\top \boldsymbol{W_a} \bar{\boldsymbol{h}}_s & \textit{general} \\ \boldsymbol{v}_a^\top \tanh\left(\boldsymbol{W_a}[\boldsymbol{h_t}; \bar{\boldsymbol{h}}_s]\right) & \textit{concat} \end{cases}$$

[Luong et al. 2015]

# Recap

- Attention solved the limited memory problem
- Complexity is quadratic (in the length of sequence)

# "Attention is all you need": single head

General purpose "single-head" attention:



$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

- Recombining previous layer (*V*) in a long-range way using few parameters

[Vaswani et al. NIPS17]

# "Attention is all you need": multiple heads

[Vaswani et al. NIPS17]



General purpose "multi-head" attention:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

# Transformer architecture



- Applying multi-head attention several times first for input, then for output
- Each unit is residual
- Emitting output one word at a time
- Positional encoding adds position-dependent features

[Vaswani et al. NIPS17]

# "Attention is all you need"



[Vaswani et al. NIPS17]

# "Attention is all you need"



[Vaswani et al. NIPS17]

# "Attention is all you need"



[Vaswani et al. NIPS17]

# "Attention is all you need"



[Vaswani et al. NIPS17]

# "Attention is all you need"



[Vaswani et al. NIPS17]

## Stack augmented RNNs

- Inherent limitation of RNNs: memory capacity

- Increasing memory by *n* gives the increase of parameters by $n^2$

- **Conclusion:** we need to decouple memory and operations (thnik RAM and CPU!)

[Joulin and Mikolov NIPS15]

# Stack augmented RNNs

Conclusion: we need to decouple memory
and operations (think RAM and CPU!)



[Joulin and Mikolov NIPS15]

# Stack augmented RNNs

$$h_t = \sigma(Ux_t + Wh_{t-1} + Ps_{t-1}^k)$$



$$y_t = \text{SoftMax}(Vh_t)$$

$$a_t = \text{SoftMax}(Ah_t)$$

$$\sigma(Dh_t)$$

Push    Pop    No

[Joulin and Mikolov NIPS15]

# Stack augmented RNNs

$$a_t = \text{SoftMax}(Ah_t)$$

Actions: Push, Pop, No

$$s_t^0 = a_t^{\text{Push}}\sigma(Dh_t) + a_t^{\text{Pop}}s_{t-1}^1 + a_t^{\text{No}}s_{t-1}^0$$

$$s_t^i = a_t^{\text{Push}}s_{t-1}^{i-1} + a_t^{\text{Pop}}s_{t-1}^{i+1} + a_t^{\text{No}}s_{t-1}^i$$



$\sigma(Dh_t)$

Push    Pop    No

[Joulin and Mikolov NIPS15]

# Binary addition with stack-RNN

Goal: train a network that can add binary numbers.



|  | Inputs: | . | 1 | 0 | 0 | 0 | 1 | 1 | + | 1 | 1 | 1 | 0 | = | 1 | 0 | 0 | 0 | 1 | 1 | . |  |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | Predictions: | 0 | 0 | . | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | . | 0 |  |
| Stack 1: | 0 |  |  |  |  | 1 |  |  | 1 |  |  |  |  |  | 0 | Counter |
| Stack 2: | 1 | -1 |  |  |  |  |  |  | 1 |  |  |  | 0 |  | 1 | End of number 2 |
| Stack 3: | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | Number 2 |
| Stack 4: |  |  |  |  |  | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |  | Length of number 2 |
| Stack 5: | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | Carry |
| Stack 6: |  | 1 | 0 | 0 | 0 | 1 | 1 |  | 0 | 1 | 0 | 0 | 0 | 1 | -1 | Number 1 |
| Stack 7: |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  | Junk |
| Stack 8: |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  | Junk |
| Stack 9: |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  | Junk |
| Stack 10: |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  | Junk |

PUSH POP

NB: the answer is reversed, i.e. 101+11 = 0001

# Binary addition with stack-RNN



- Training with total lengths upto 20
- 100 hidden units and 10 1-dim stacks

# Neural Turing Machine



External Input    External Output

Controller ← RNN, LSTM, CNN

Read Heads    Write Heads

Memory

[Graves et al. 2014]

## Outlook

- RNNs allow to solve many problems with sequences (as inputs or outputs)
- CTC-loss is useful for monotonically aligned input-output tasks
- The *attention* idea is working and is used across different domains (e.g. computer vision)
- Learning a computer to "program" is ambitious and promising
- Currently works only for simplistic algorithms
- Differentiability requires real-valued (soft) values
- Learning systems that make discrete choices is harder (but possible)

# Bibliography

A. Graves. Supervised Sequence Labelling with Recurrent Neural Networks. Textbook, Studies in Computational Intelligence, Springer, 2012

Sepp Hochreiter, Jürgen Schmidhuber: Long Short-Term Memory. Neural Computation 9(8): 1735-1780 (1997)

Ilya Sutskever, Oriol Vinyals, Quoc V. Le:
Sequence to Sequence Learning with Neural Networks. NIPS 2014: 3104-3112

Jeff Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Trevor Darrell, Kate Saenko:
Long-term recurrent convolutional networks for visual recognition and description. CVPR 2015: 2625-2634

Justin Johnson, Andrej Karpathy, Li Fei-Fei, DenseCap: Fully Convolutional Localization Networks for Dense Captioning. CVPR 2016

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, Fei-Fei Li: Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. CoRR abs/1602.07332 (2016)

D. Bahdanau, K. Cho, and Y. Bengio: Neural machine translation by jointly learning to align and translate. In ICLR 2015.

# Bibliography

Wang Ling, Chris Dyer, Alan W. Black, Isabel Trancoso, Ramon Fermandez, Silvio Amir, Luís Marujo, Tiago Luís: Finding Function in Form: Compositional Character Models for Open Vocabulary Word Representation. EMNLP 2015: 1520-1530

Minh-Thang Luong, Hieu Pham, Christopher D. Manning:
Effective Approaches to Attention-based Neural Machine Translation. CoRR abs/1508.04025 (2015)

Alex Graves, Santiago Fernández, Faustino J. Gomez, Jürgen Schmidhuber:
Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. ICML 2006: 369-376

Armand Joulin, Tomas Mikolov:
Inferring Algorithmic Patterns with Stack-Augmented Recurrent Nets. NIPS 2015: 190-198

Alex Graves, Greg Wayne, Ivo Danihelka:
Neural Turing Machines. CoRR abs/1410.5401 (2014)

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin:
Attention is All you Need. NIPS 2017: 6000-6010