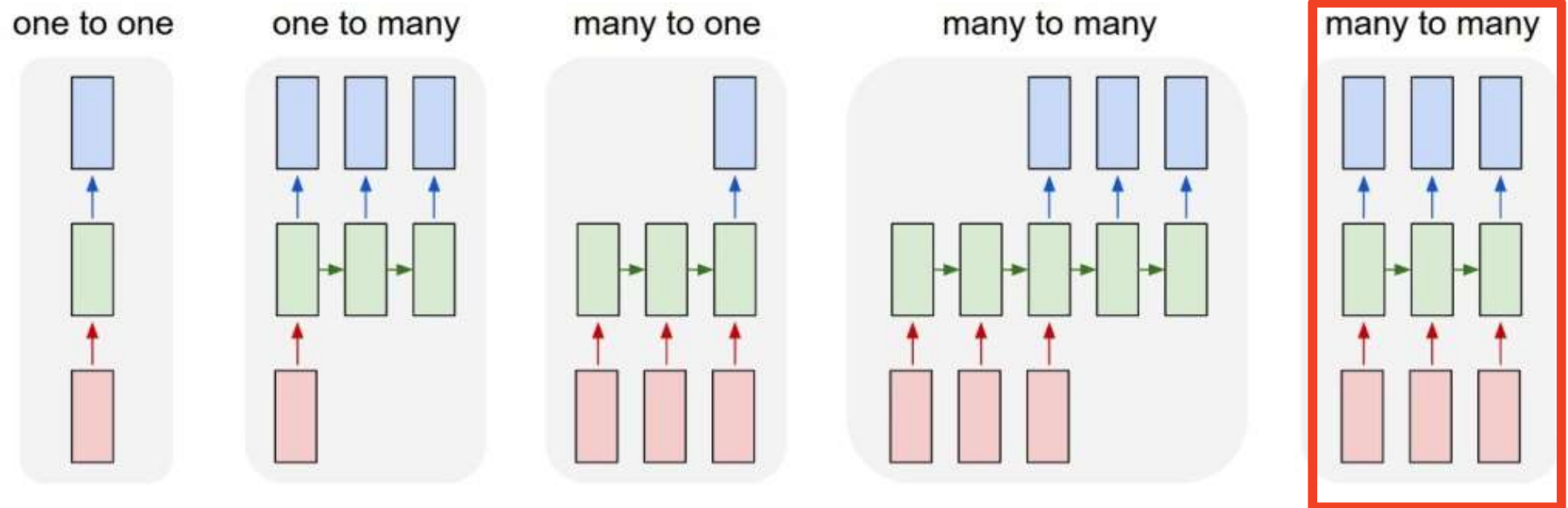


Lecture 11: Sequence-to-sequence architectures

Learning settings

slide credit: A. Karpathy



One-to-one: image to class label

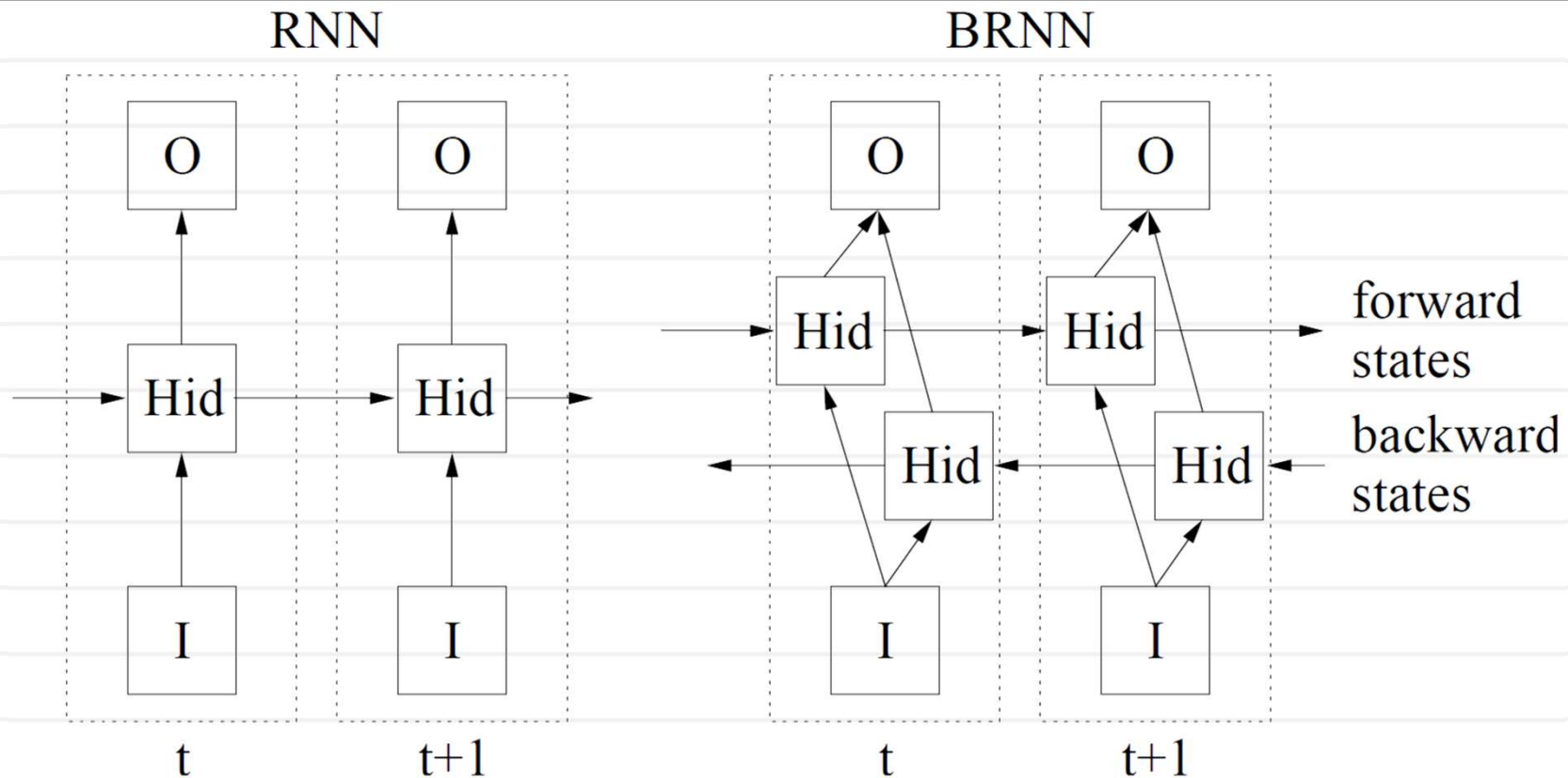
One-to-many: text generation/image captioning

Many-to-one: sentiment analysis

Many-to-many 1: machine translation

Many-to-many 2: online classification (e.g. POS tagging)

Bi-directional RNN



```
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]

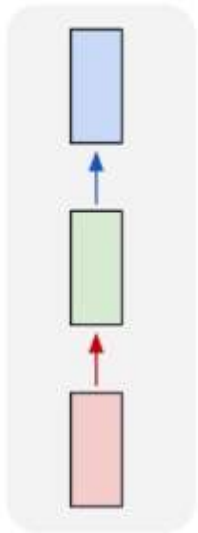
Uni-directional vs bi-directional

- 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

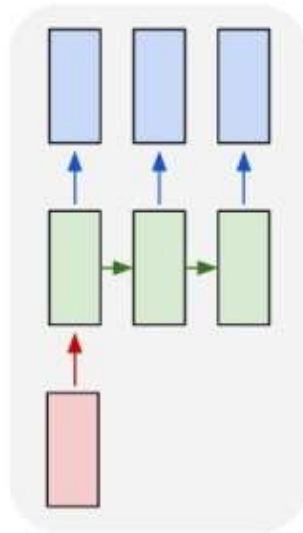
Learning settings

slide credit: A. Karpathy

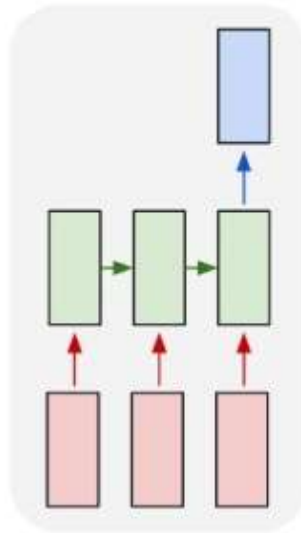
one to one



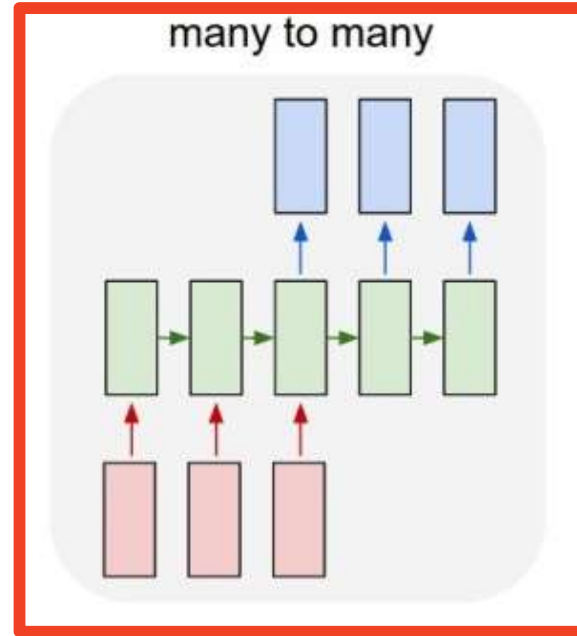
one to many



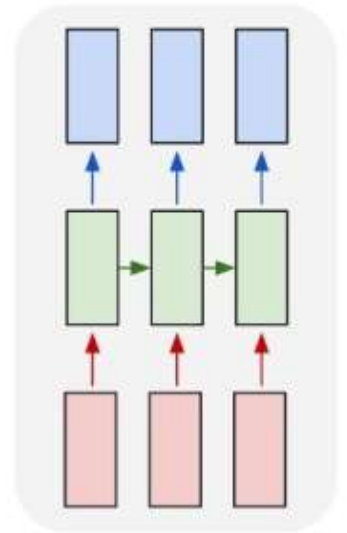
many to one



many to many

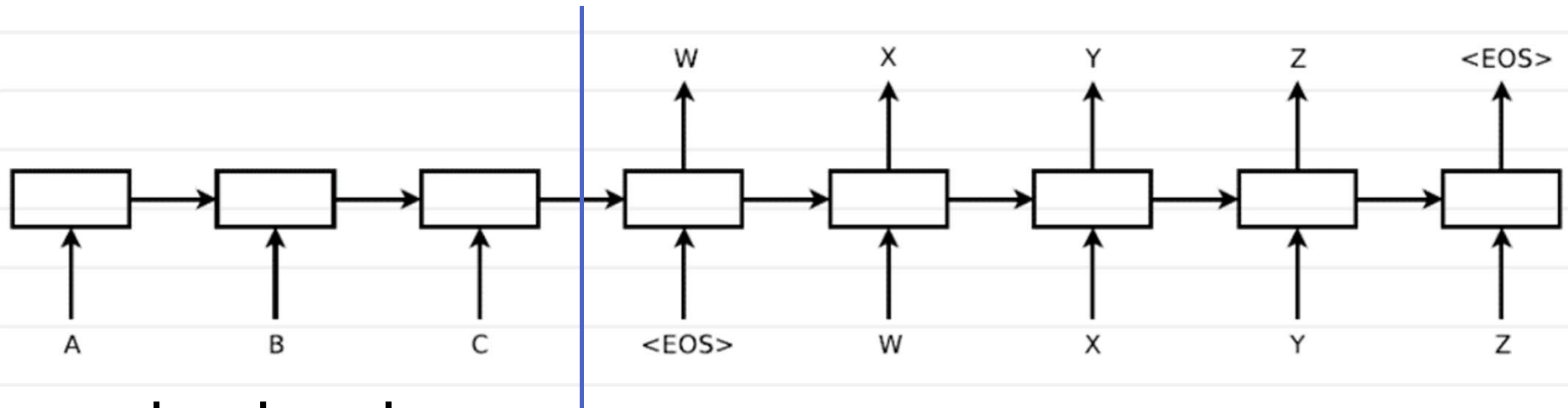


many to many



aka "seq2seq"

Sequence-to-sequence machine translation



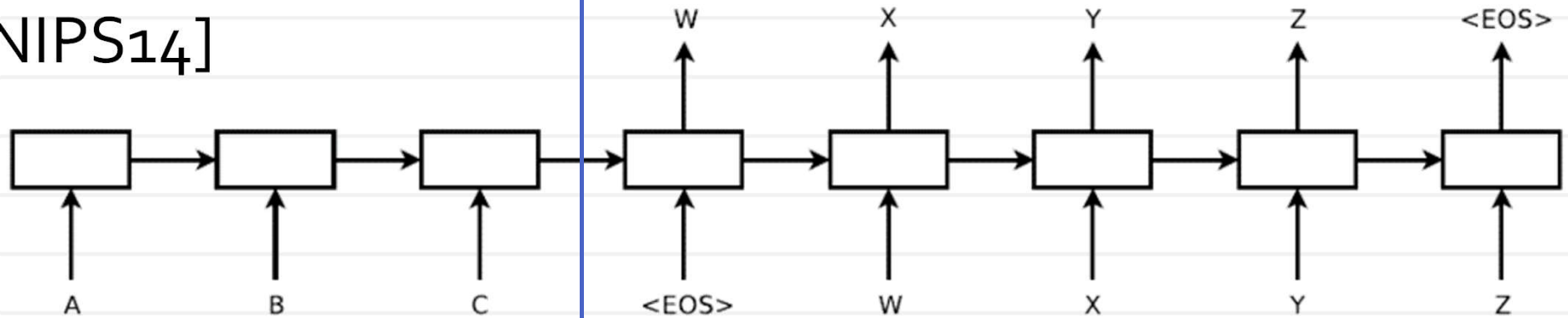
Important notes:

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]



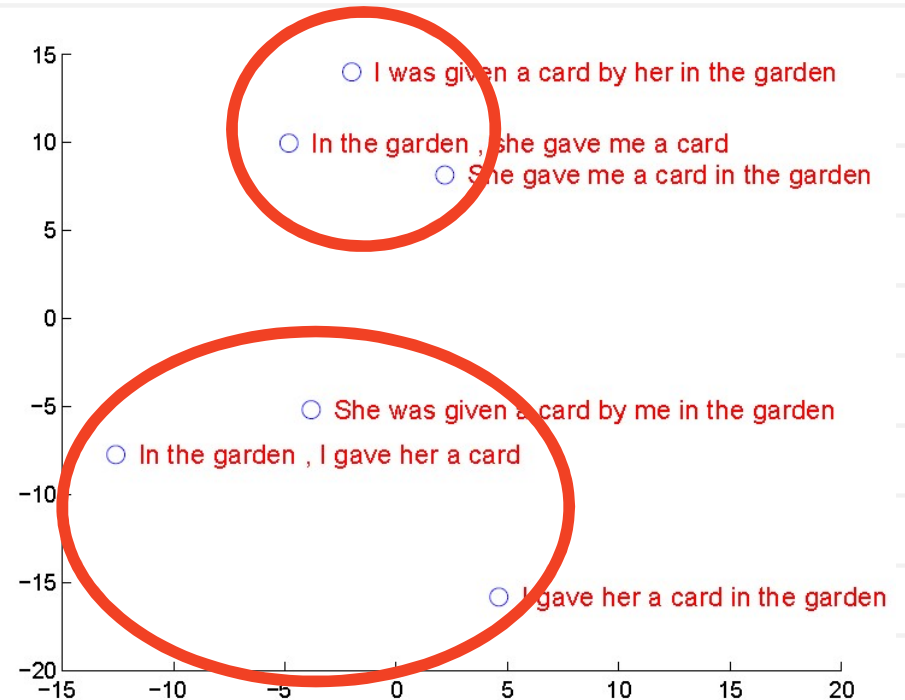
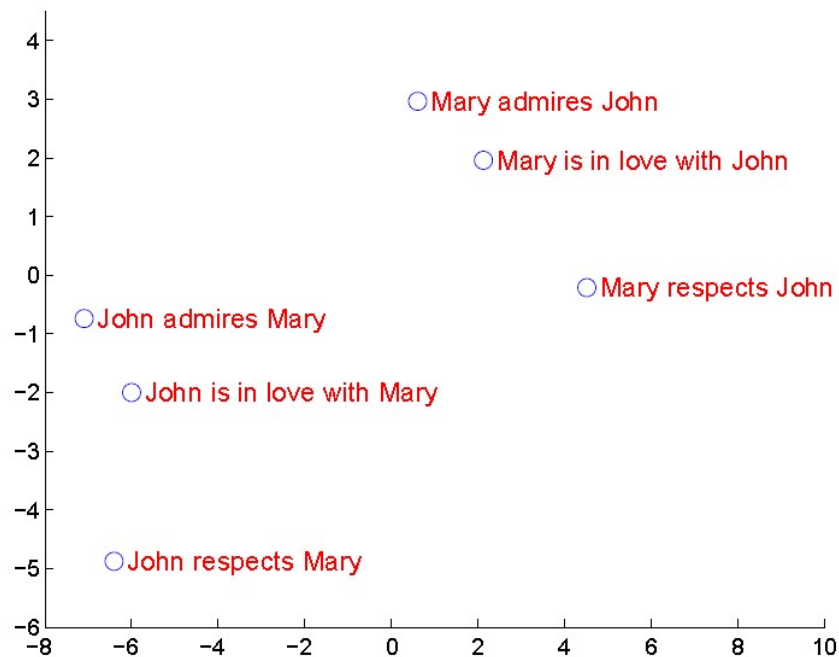
Decoding proceeds by *beam search*:

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- \rightarrow 2

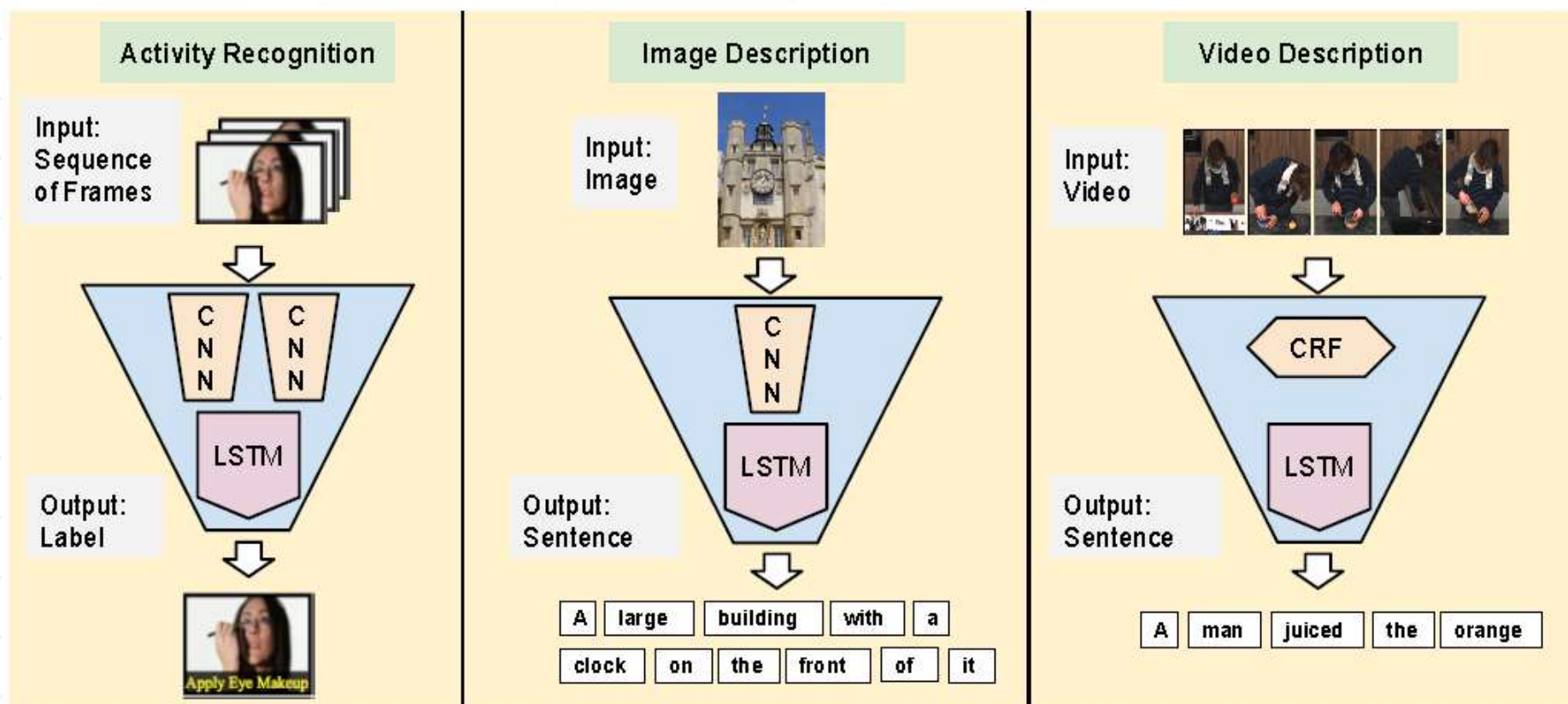
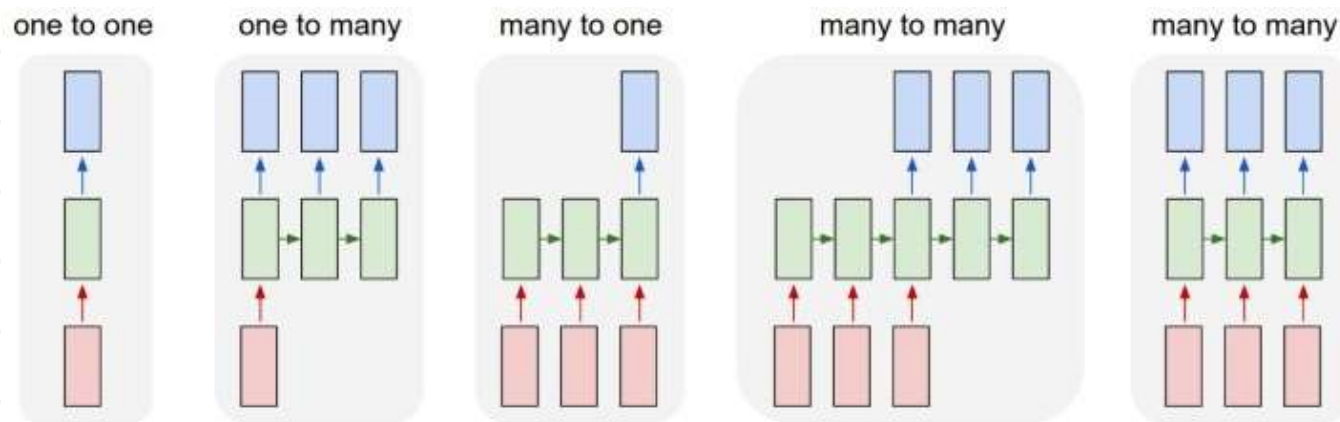
[Sutskever et al. NIPS14]

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 ” .

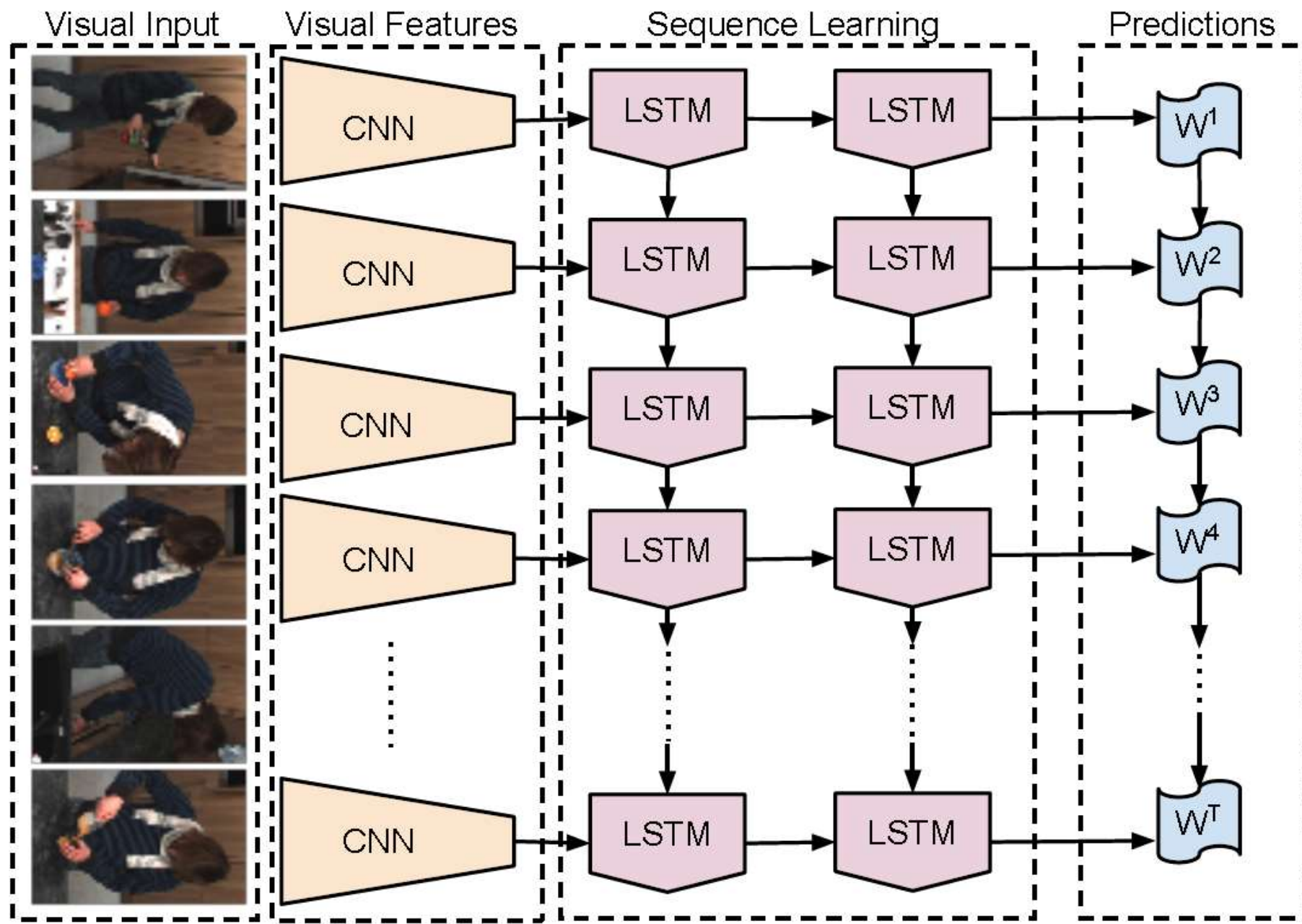
[Sutskever et al. NIPS14]

Image/video captioning



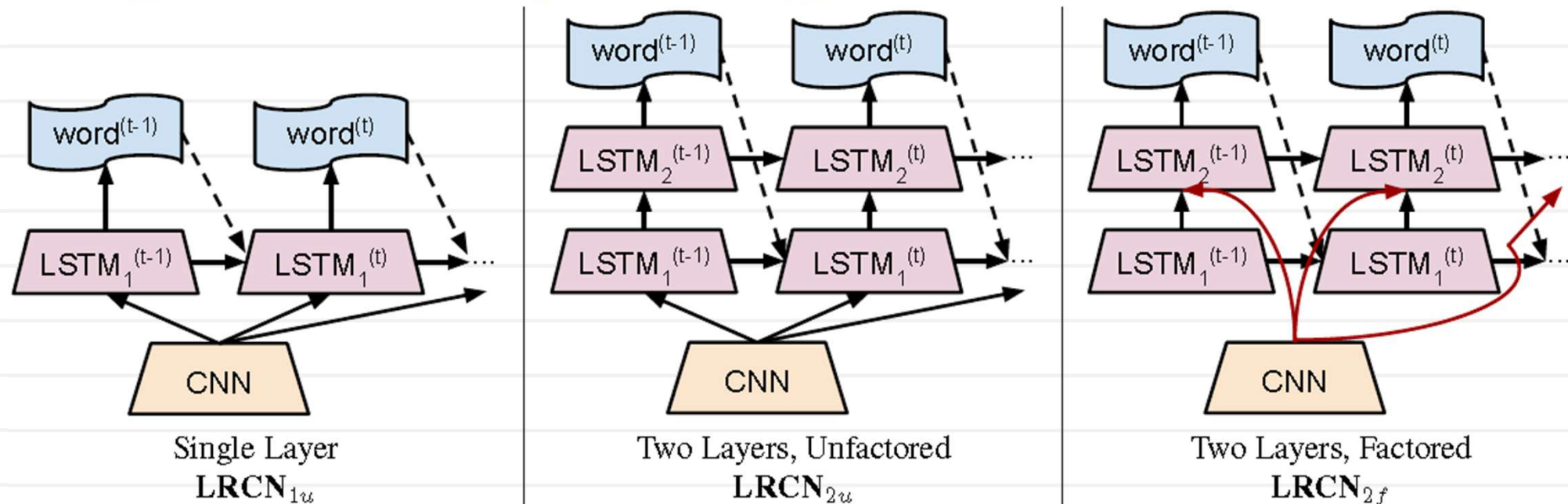
[Donahue et al. 2015]

Image/video captioning



[Donahue et al. 2015]*

Image/video captioning



- 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	Medr
LRCN _{1u}	14.1	31.3	39.7	24
LRCN _{2u}	3.8	12.0	17.9	80
LRCN _{2f}	17.5	40.3	50.8	9
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 surfboard.



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



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

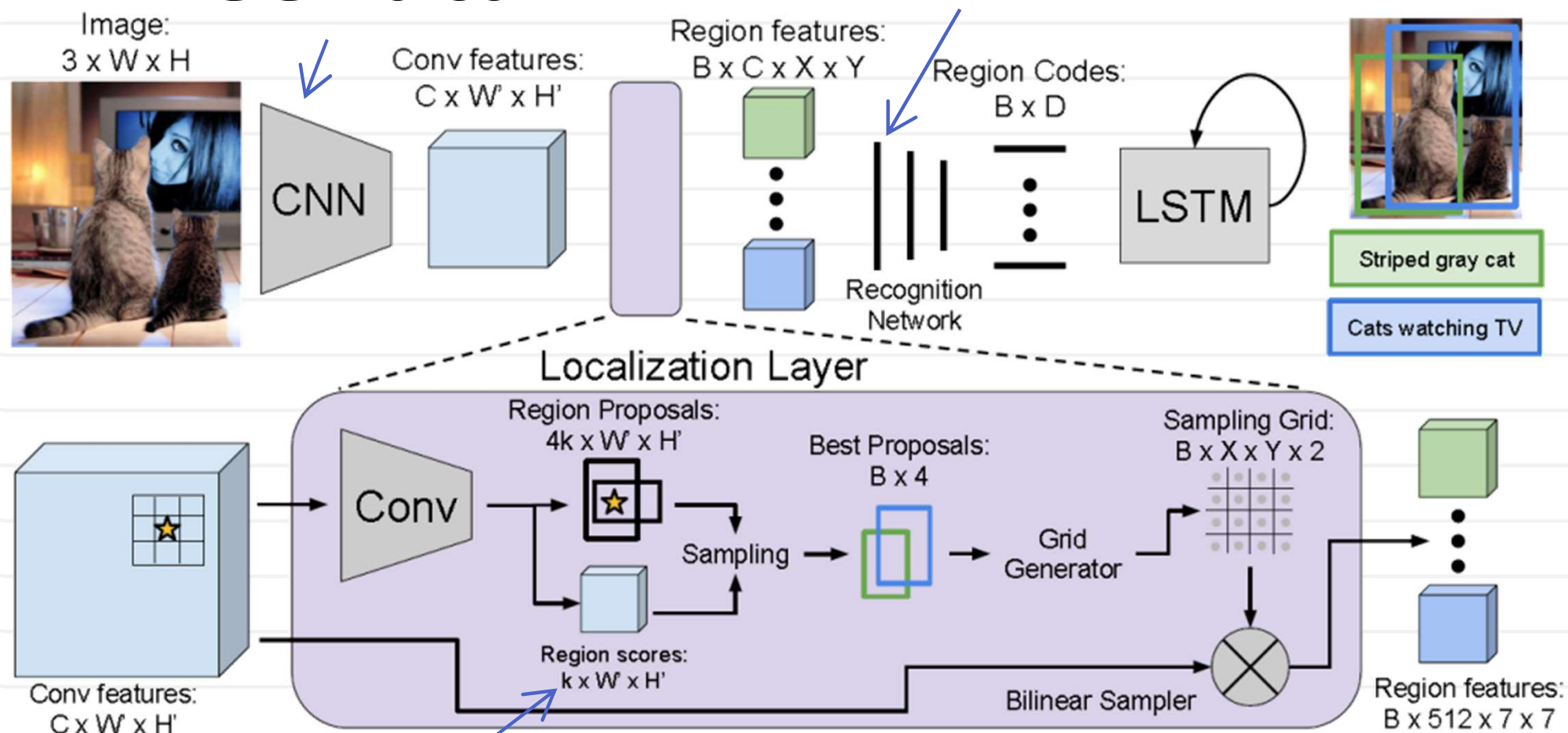


A person holding a cell phone in their hand.

[Donahue et al. 2015]

End-to-end dense image captioning

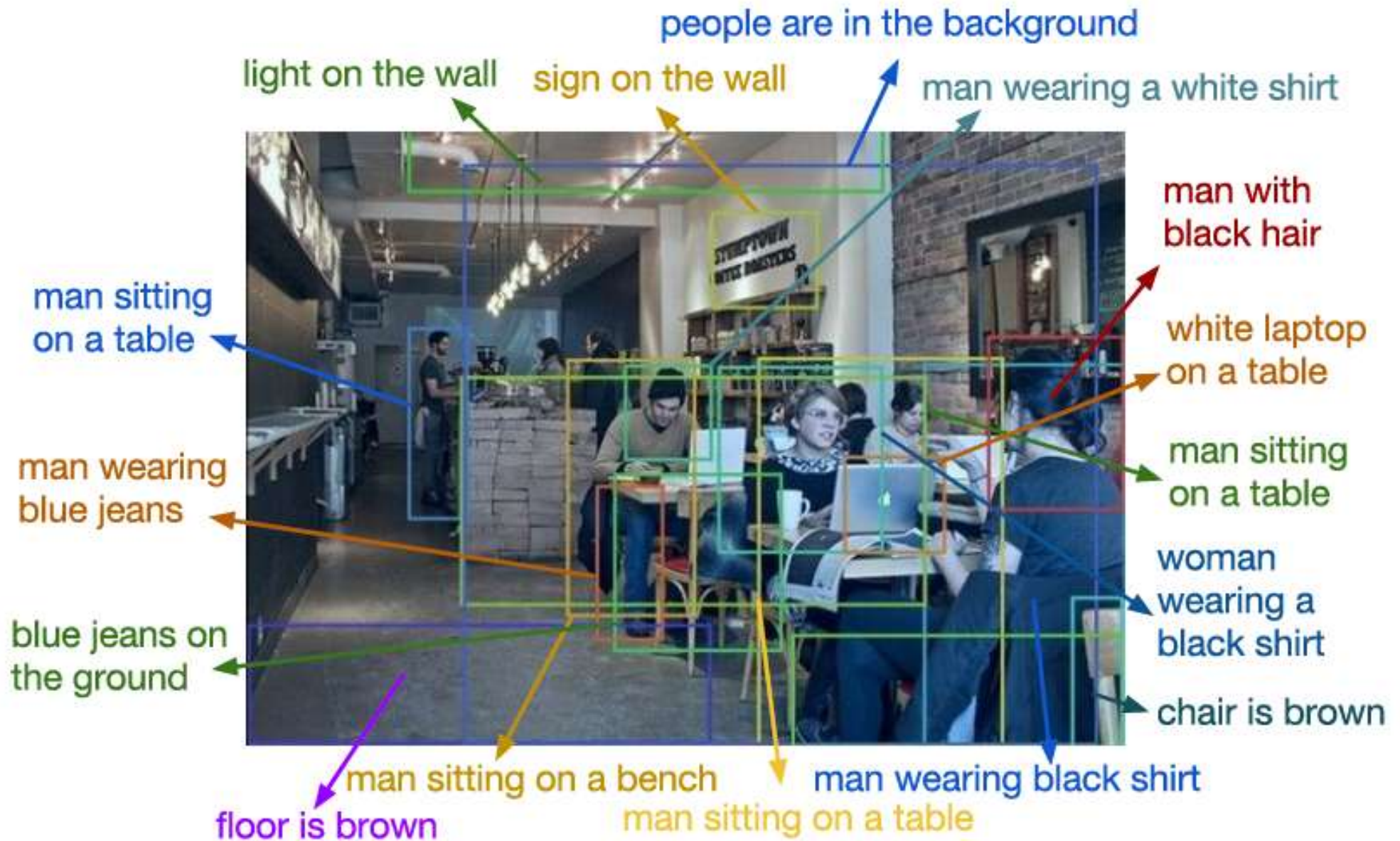
VGG-16 conv fully-connected, 2 layers+dropout



k-anchors at $W' \times H'$ positions

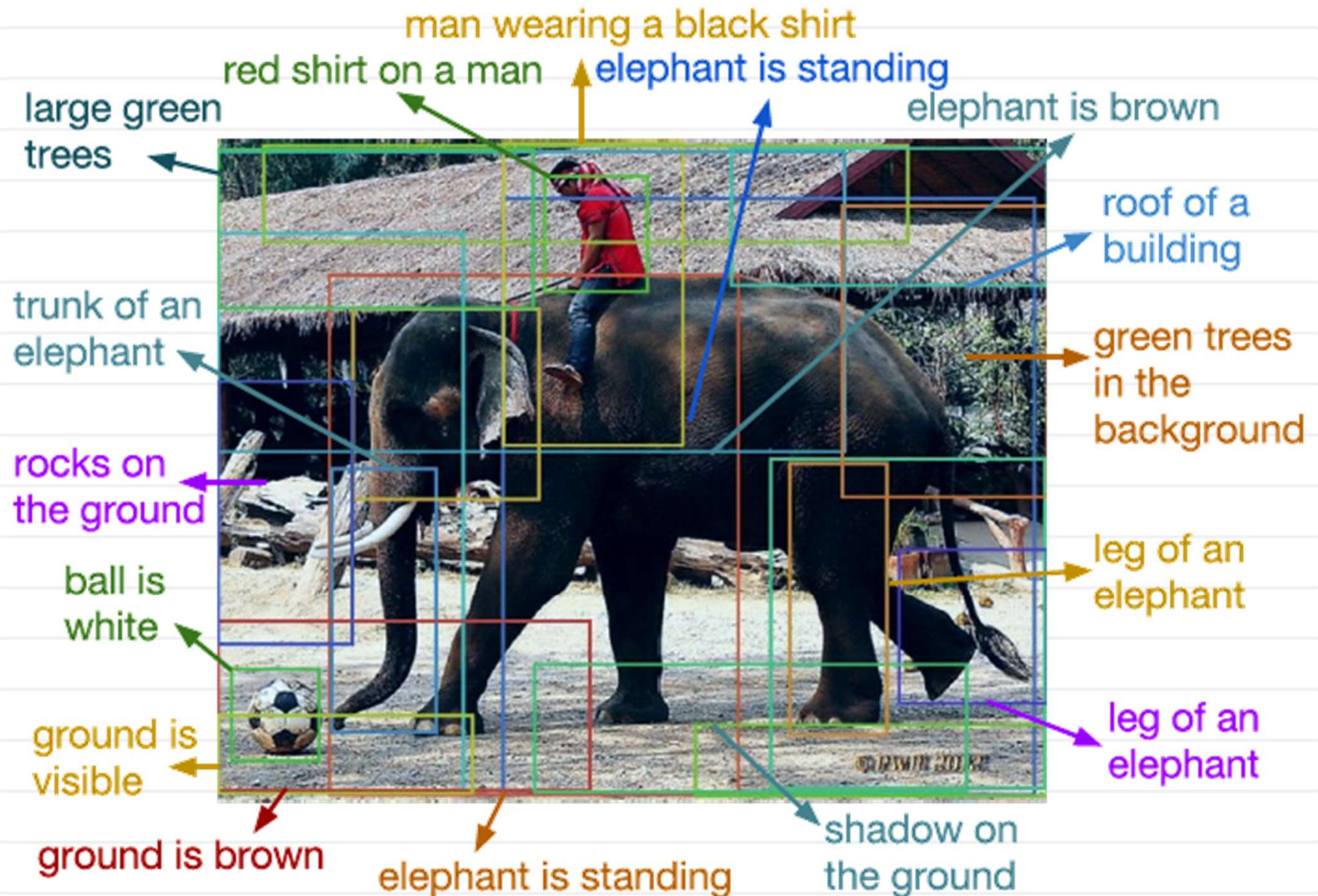
[Johnson et al, CVPR16]

End-to-end dense image captioning



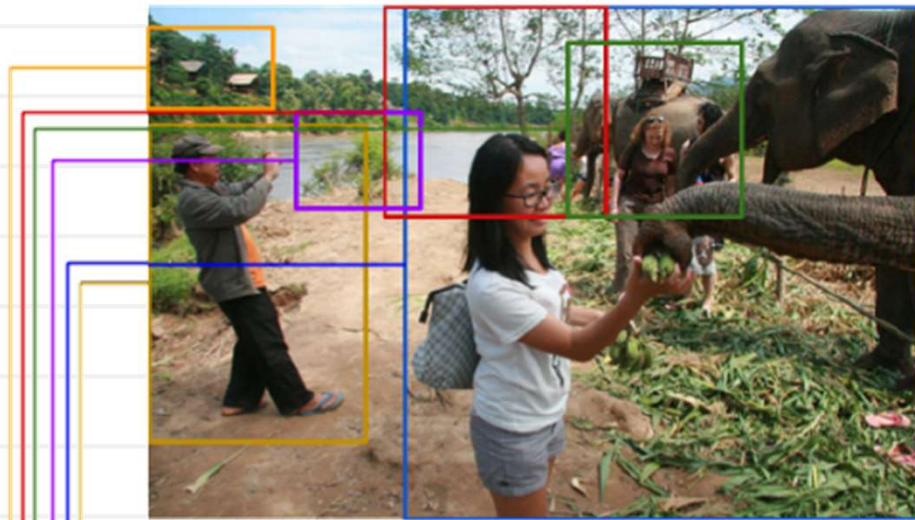
[Johnson et al, CVPR16]

End-to-end dense image captioning



[Johnson et al, CVPR16]

Training set: “visual genome”



Girl feeding elephant
Man taking picture
Huts on a hillside

→ **A man taking a picture.**

Flip flops on the ground
Hillside with water below
Elephants interacting with people
Young girl in glasses with backpack
Elephant that could carry people

→ **An elephant trunk taking two bananas.**

→ **A bush next to a river.**

People watching elephants eating
A woman wearing glasses.
A bag
Glasses on the hair.

→ **The elephant with a seat on top**

A woman with a purple dress.
A pair of pink flip flops.
A handle of bananas.

→ **Tree near the water**

A blue short.

→ **Small houses on the hillside**

A woman feeding an elephant
A woman wearing a white shirt and shorts
A man taking a picture

A man wearing an orange shirt
An elephant taking food from a woman
A woman wearing a brown shirt
A woman wearing purple clothes
A man wearing blue flip flops
Man taking a photo of the elephants
Blue flip flop sandals
The girl's white and black handbag
The girl is feeding the elephant
The nearby river
A woman wearing a brown t shirt
Elephant's trunk grabbing the food
The lady wearing a purple outfit
A young Asian woman wearing glasses
Elephants trunk being touched by a hand
A man taking a picture holding a camera
Elephant with carrier on it's back
Woman with sunglasses on her head
A body of water
Small buildings surrounded by trees
Woman wearing a purple dress
Two people near elephants
A man wearing a hat
A woman wearing glasses
Leaves on the ground

- “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

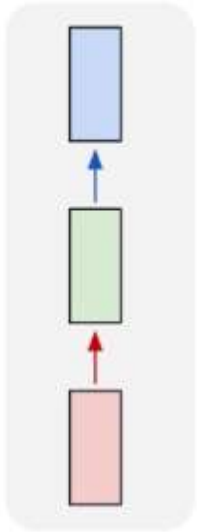
Everything Mapped to

Wordnet Synsets

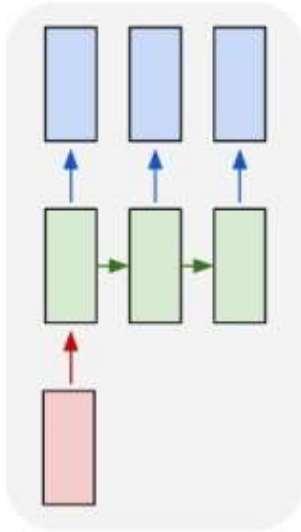
[Krishna et al. 2016]

Sequence-to-sequence machine translation

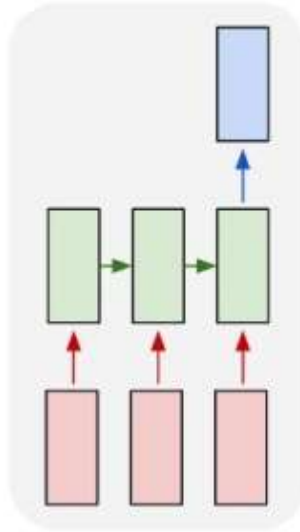
one to one



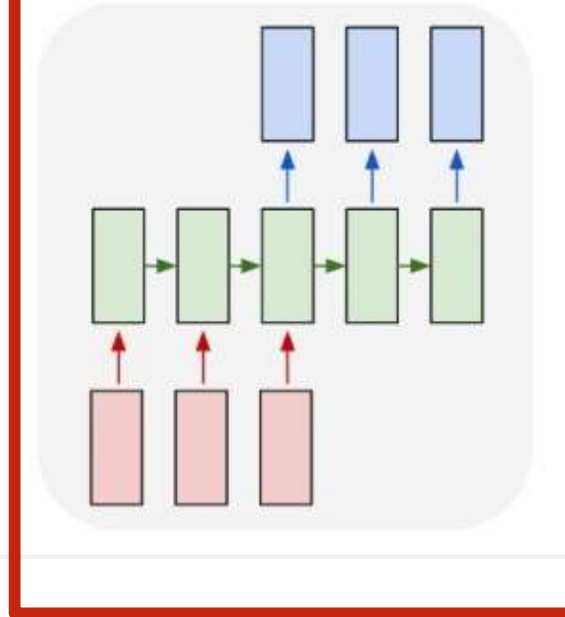
one to many



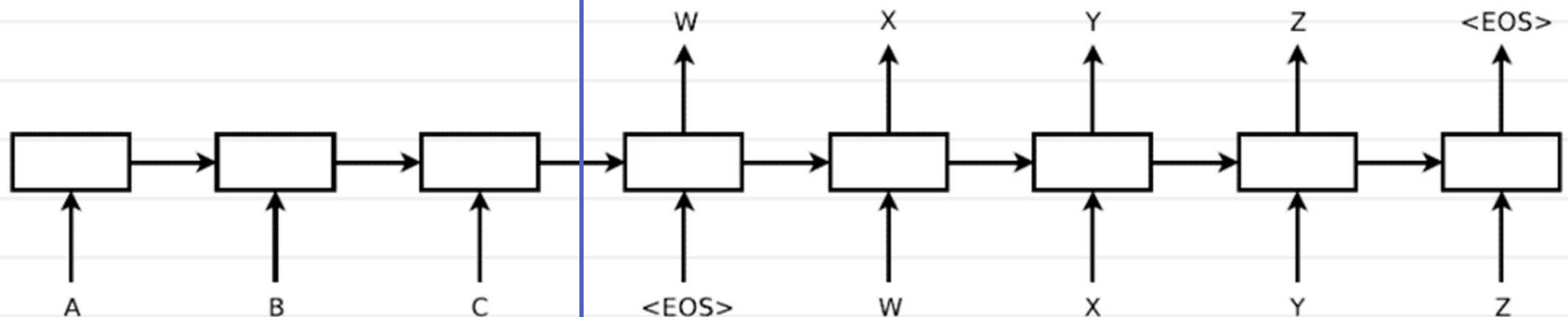
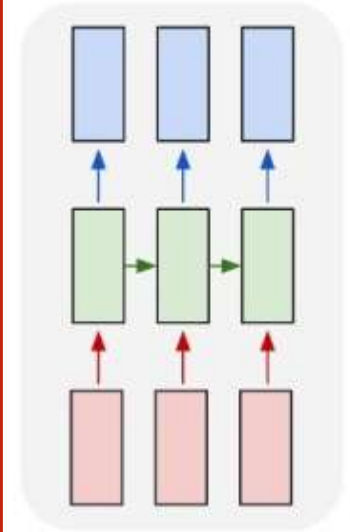
many to one



many to many

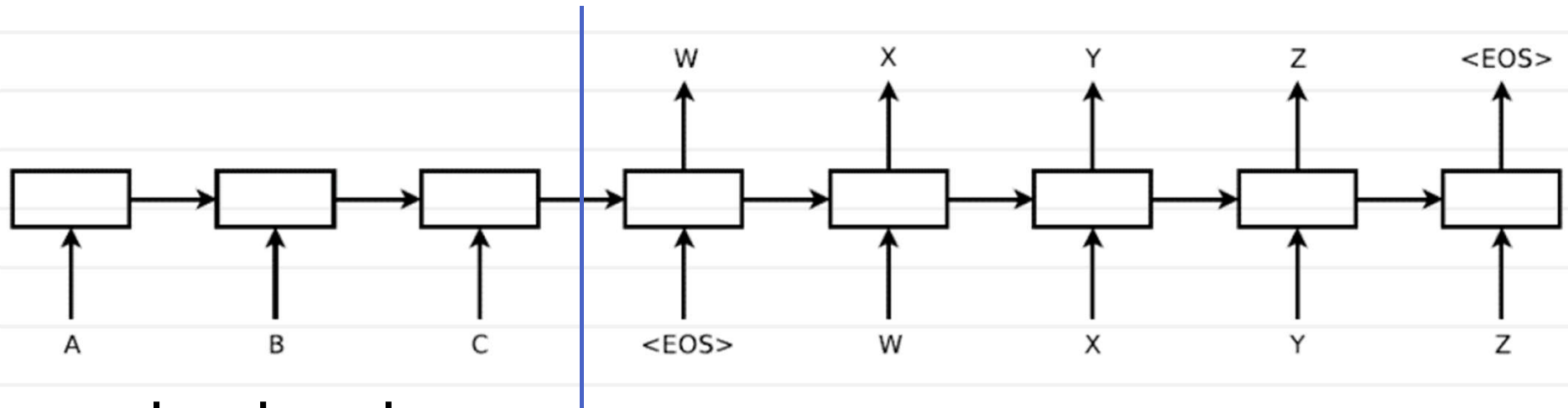


many to many



[Sutskever et al. NIPS14]

Sequence-to-sequence machine translation

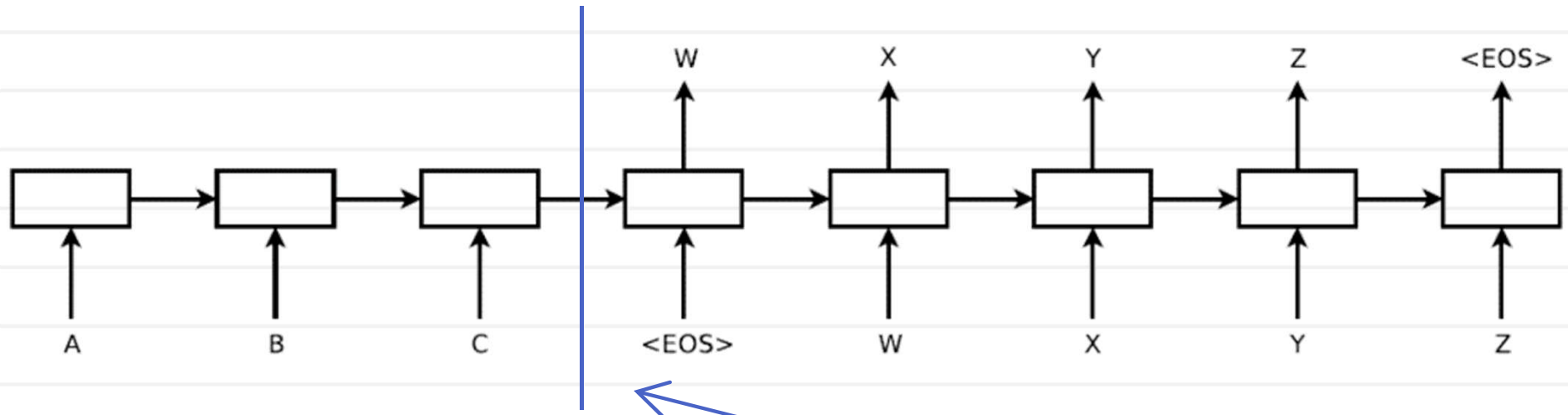


Important notes:

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



Problem:

all the meaning has to be carried from here

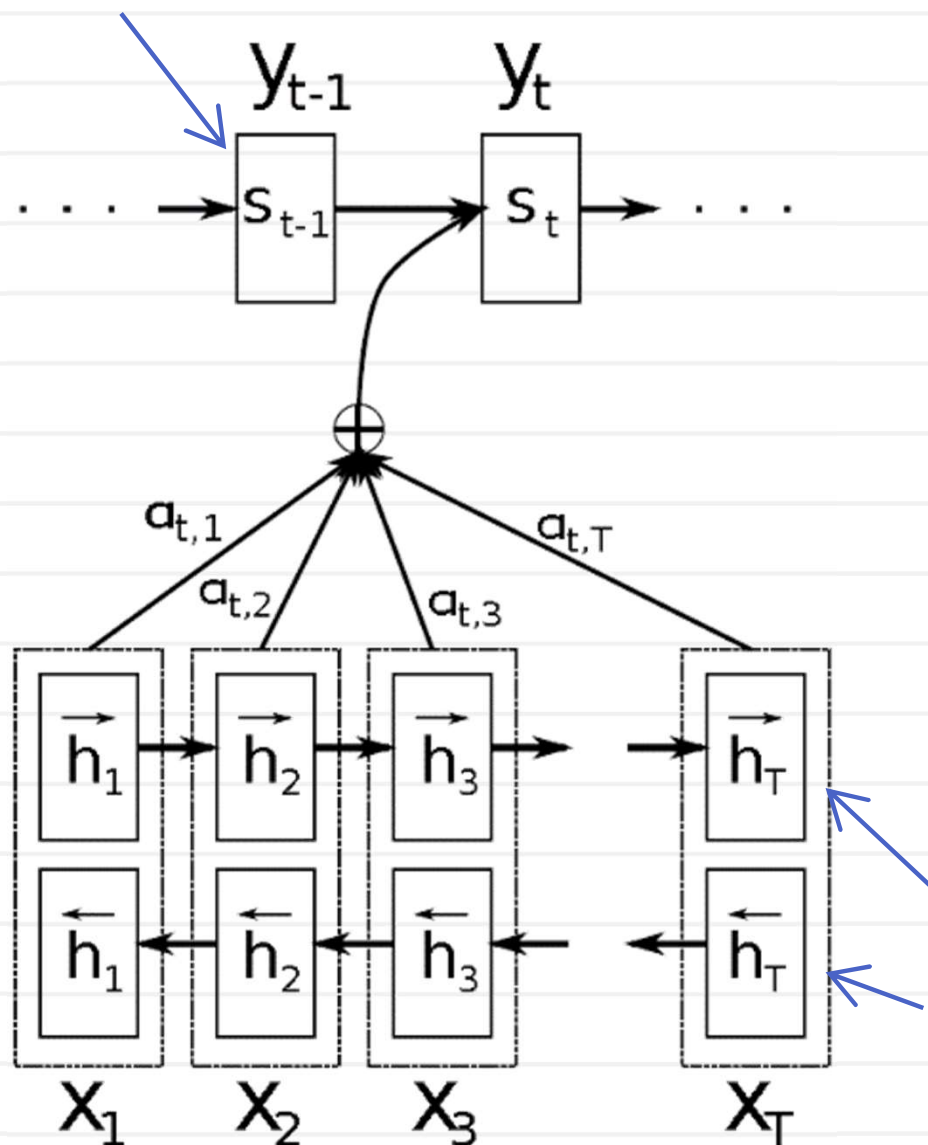
- Large memory needed
- Information has to survive for a very long time



[Sutskever et al. NIPS14]

Translation with attention

decoder RNN



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

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

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

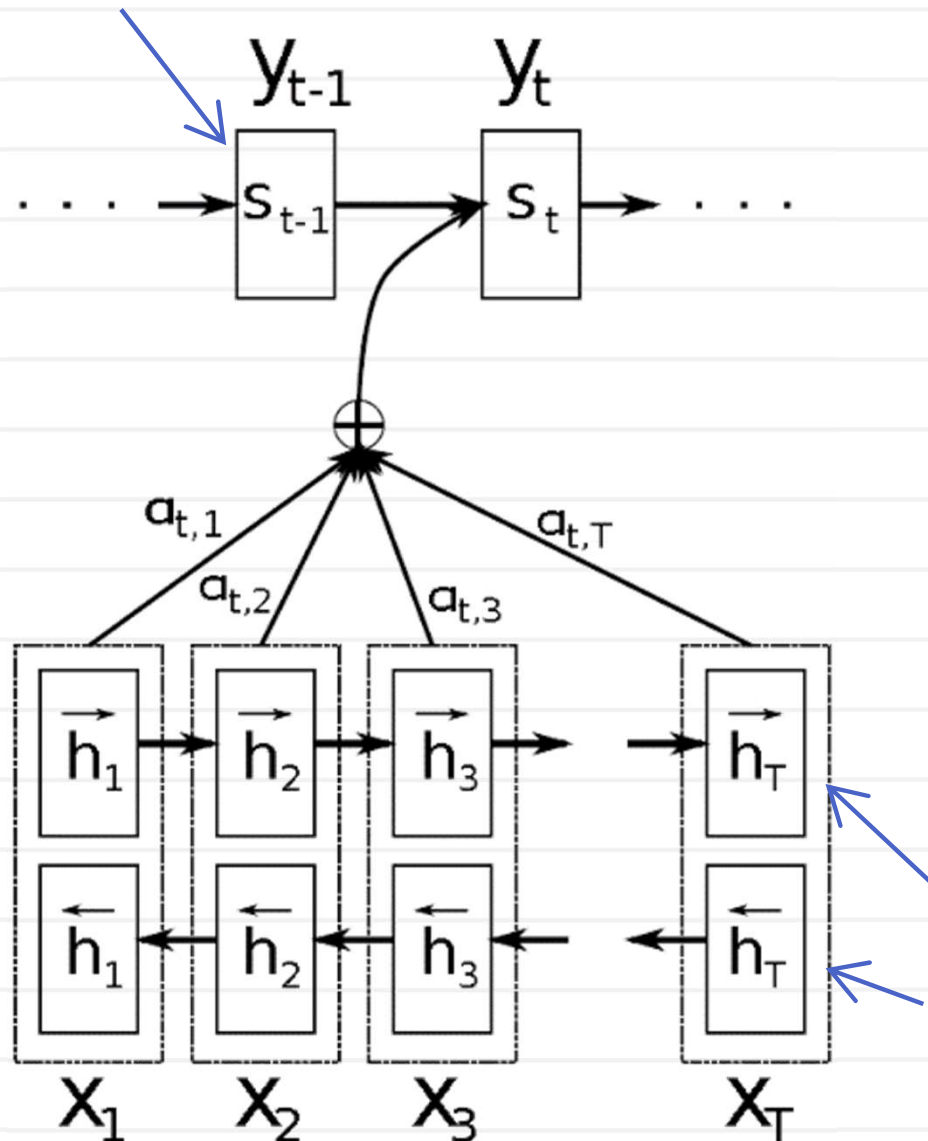
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

encoder RNN 1
encoder RNN 2

[Bahdanau et al. 2015]

Translation with attention

decoder RNN



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

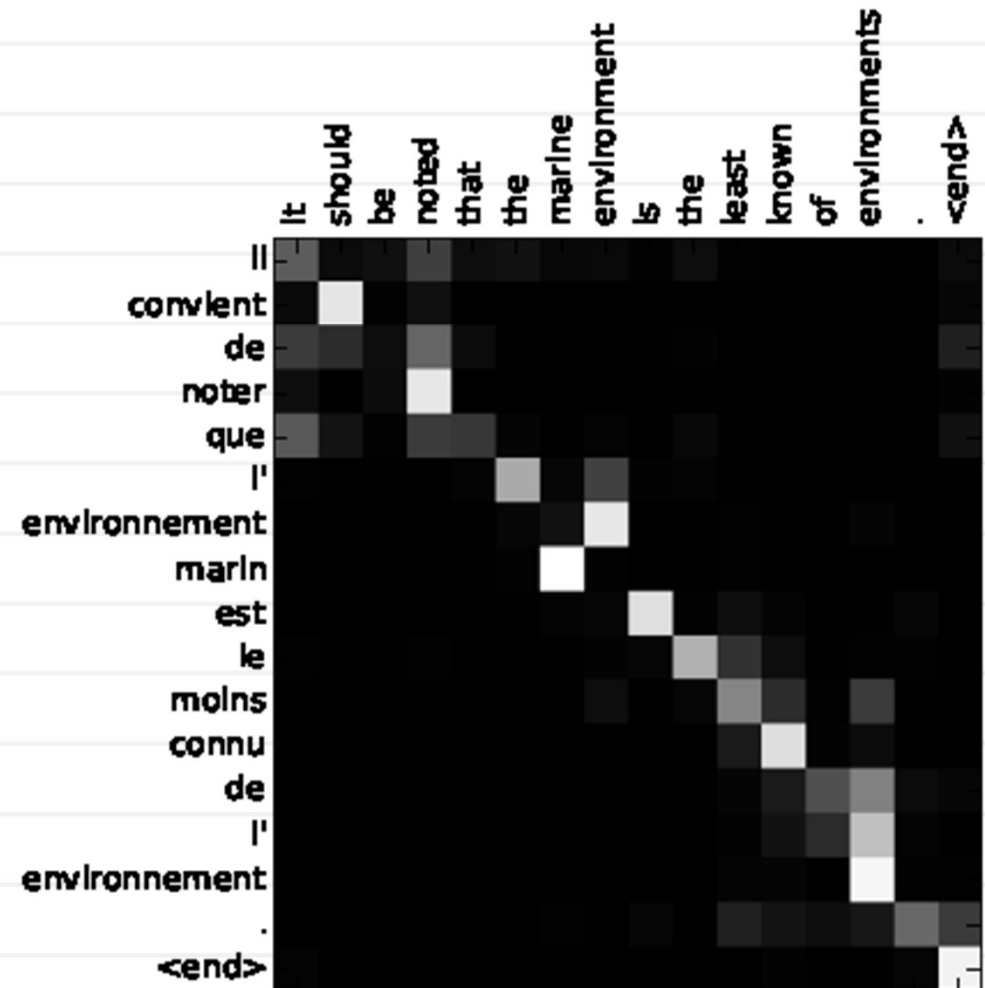
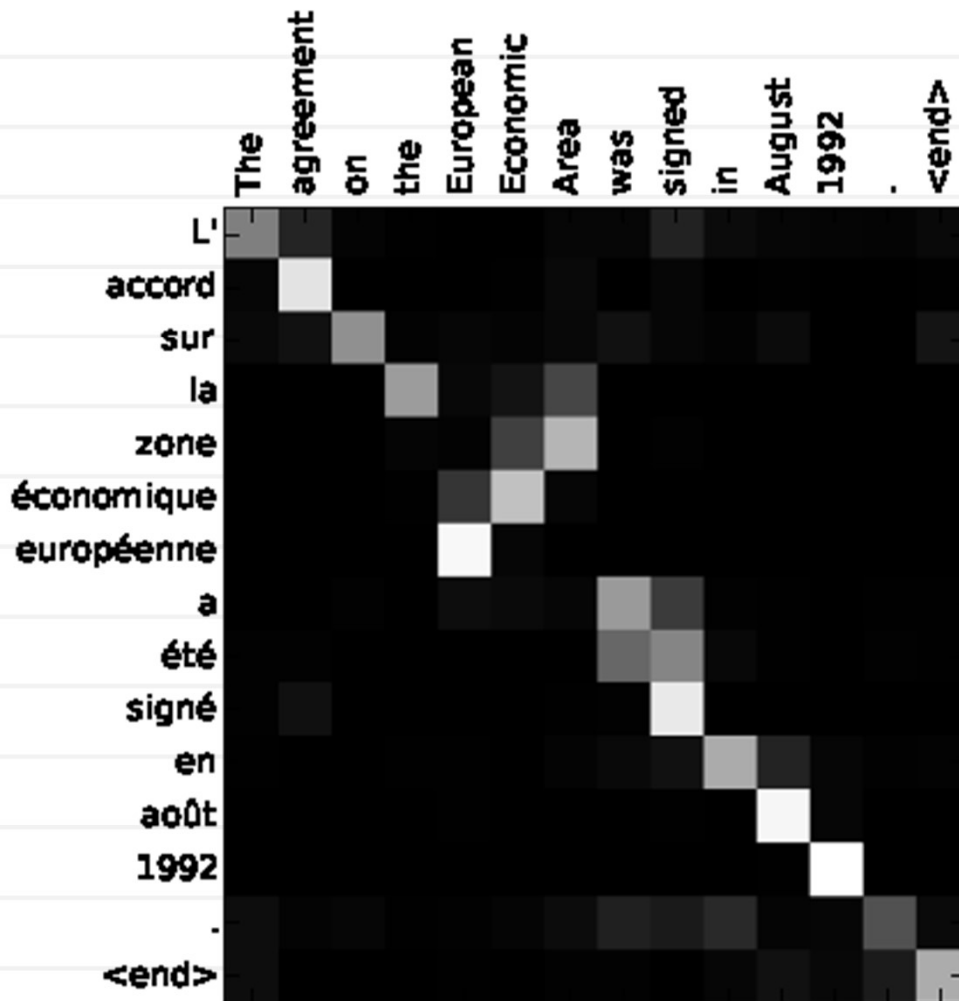
- *Attention model:* feed-forward neural network
- All components are trained end-to-end

encoder RNN 1

encoder RNN 2

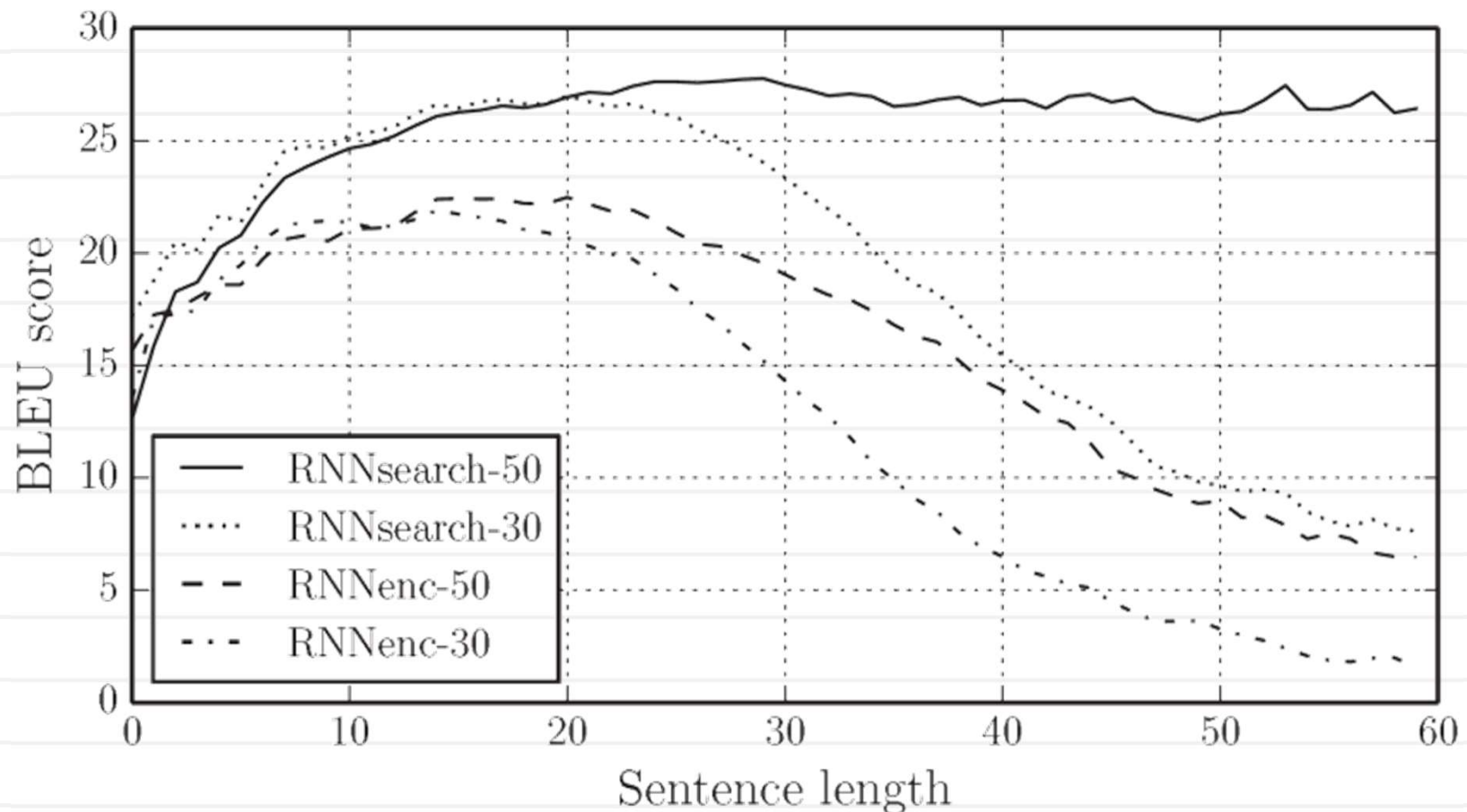
[Bahdanau et al. 2015]

Translation with attention



[Bahdanau et al. 2015]

Translation with attention



- BLEU-score \approx 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 to carry out a diagnosis or a procedure, based on his status as a health care worker at a hospital.

LSTM system:

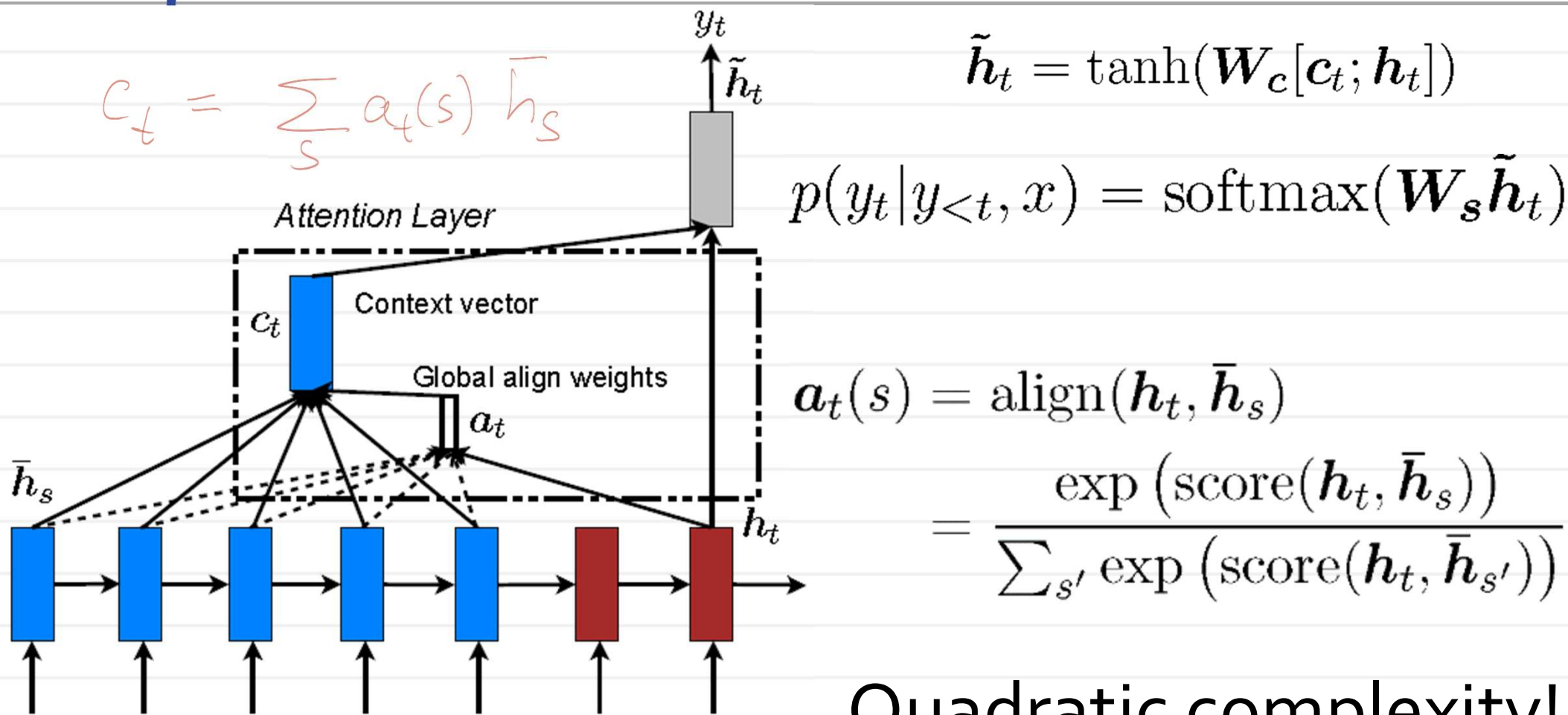
Un privilège d'admission est le droit d'un médecin de reconnaître un patient à l'hôpital ou un centre médical d'un diagnostic ou de prendre un diagnostic en fonction de son état de santé.

Attention-based system:

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

[Bahdanau et al. 2015]

Simpler translation with attention



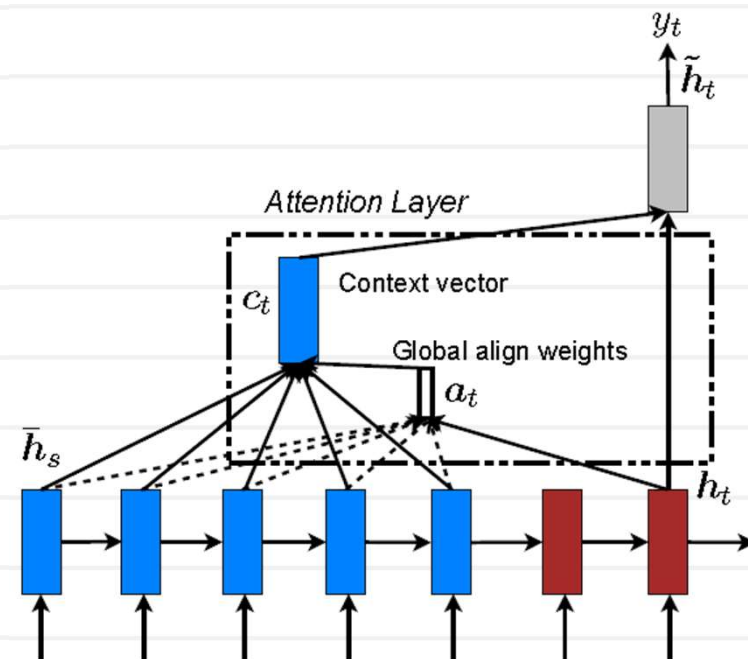
Quadratic complexity!

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^\top \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^\top \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{v}_a^\top \tanh(\mathbf{W}_a[\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat} \end{cases}$$

[Luong et al. 2015]

Recap

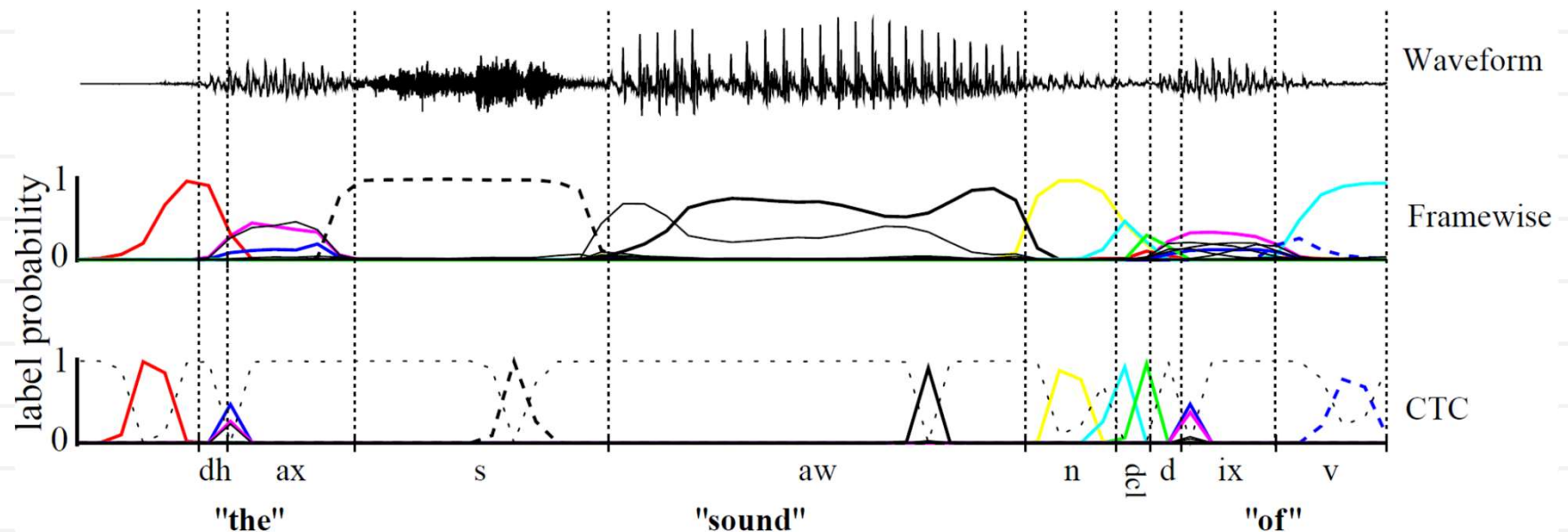
- Attention solved the limited memory problem
- Complexity is quadratic (in the length of sequence)



Online seq2seq with monotonic alignment

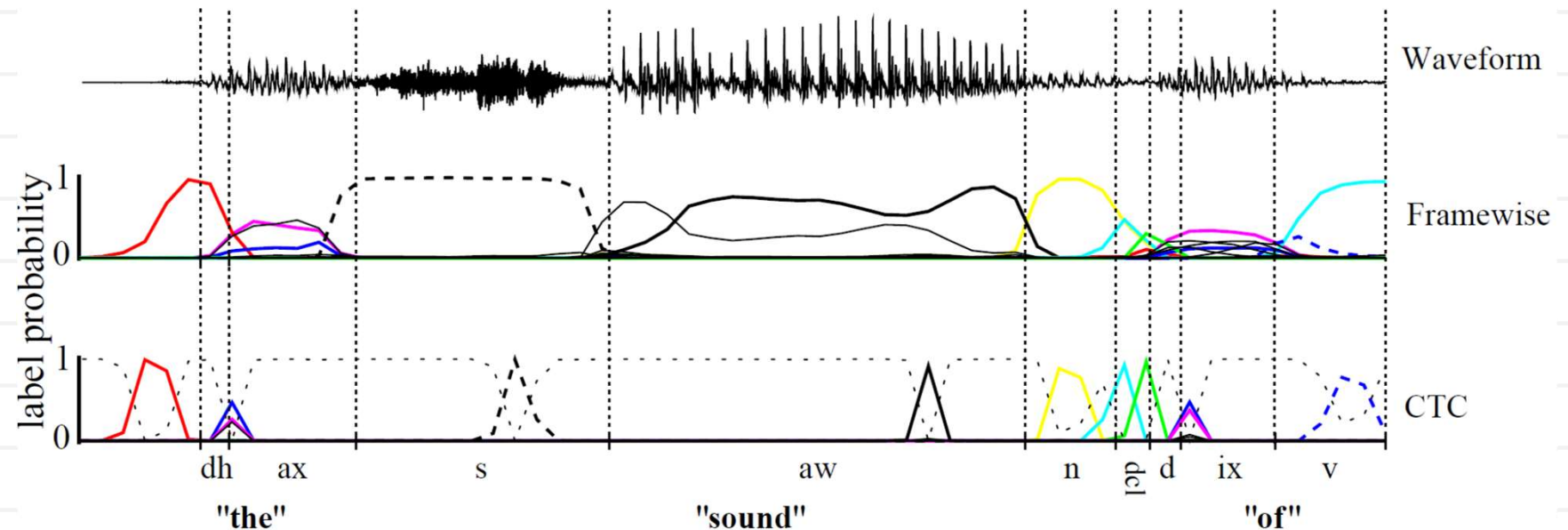
Many problems are sequence 2 sequence with monotonic alignment:

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



[Graves et al. 2006]

Online seq2seq with monotonic alignment



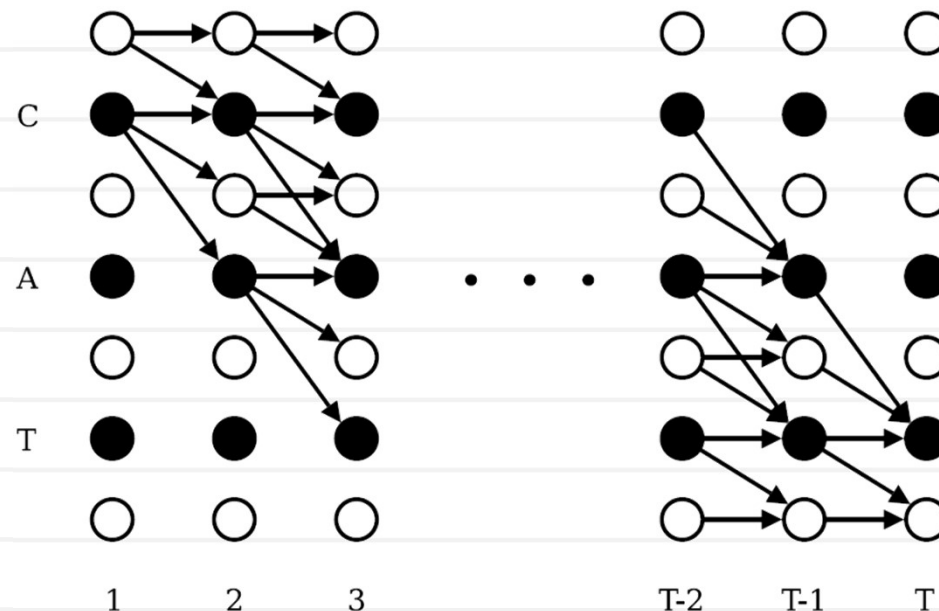
Decoding: 'aaa__bb_c____ddaa' \rightarrow abcd

What should be the loss that encourage correct parsing?

Answer: connectionist temporal classification (CTC) loss

[Graves et al. 2006]

CTC-loss

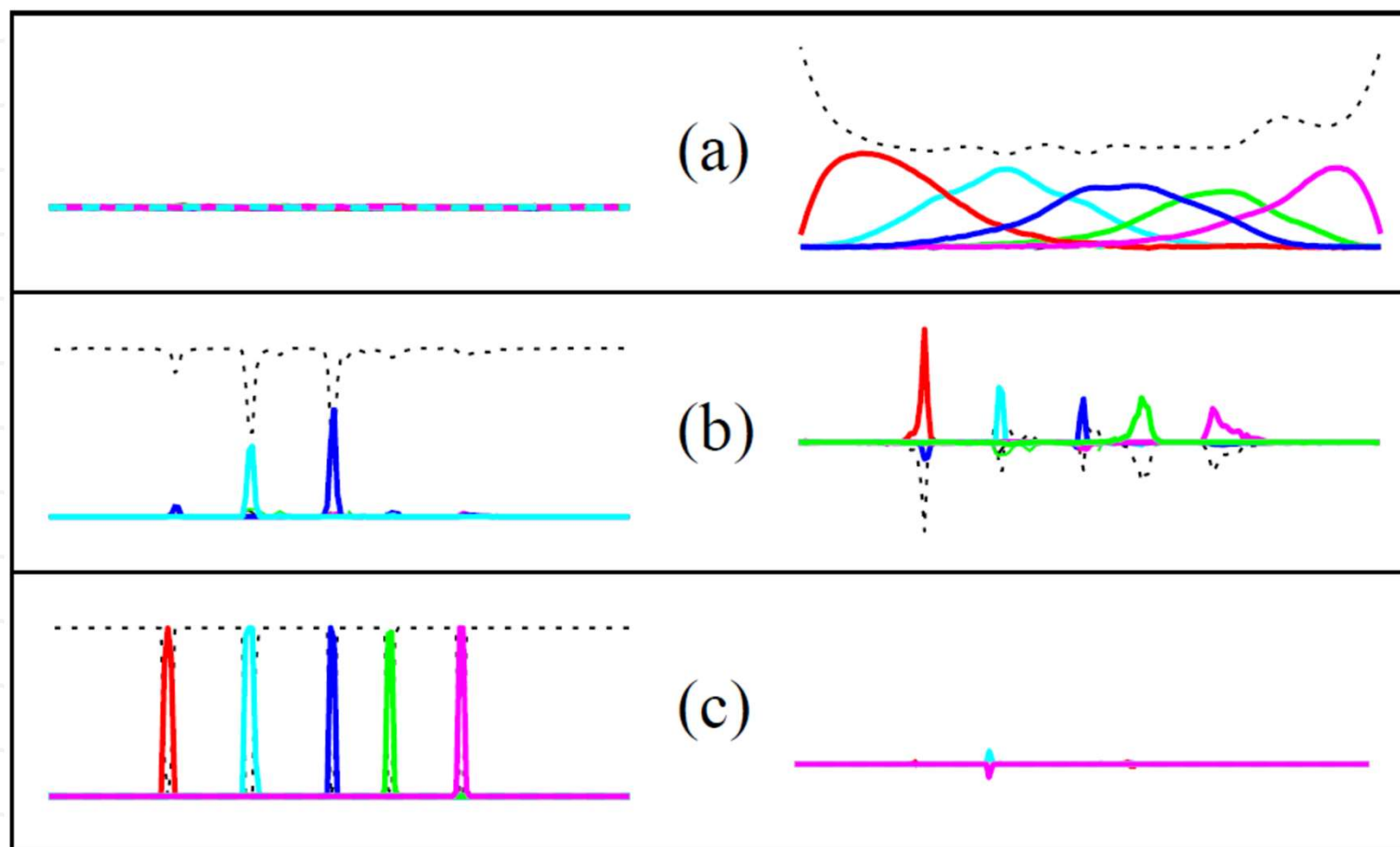


- Augment the output state with *blank*
- Predict probabilities of each symbol (inc. blank) at each time moment
- Compute the probability of each lattice vertex under correct paths using forward-backward
- Push log-probabilities up (*ML training*) proportionally to the current probability

[Graves et al. 2006]

Evolution of the CNC 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=mLxsbWAYlpw>

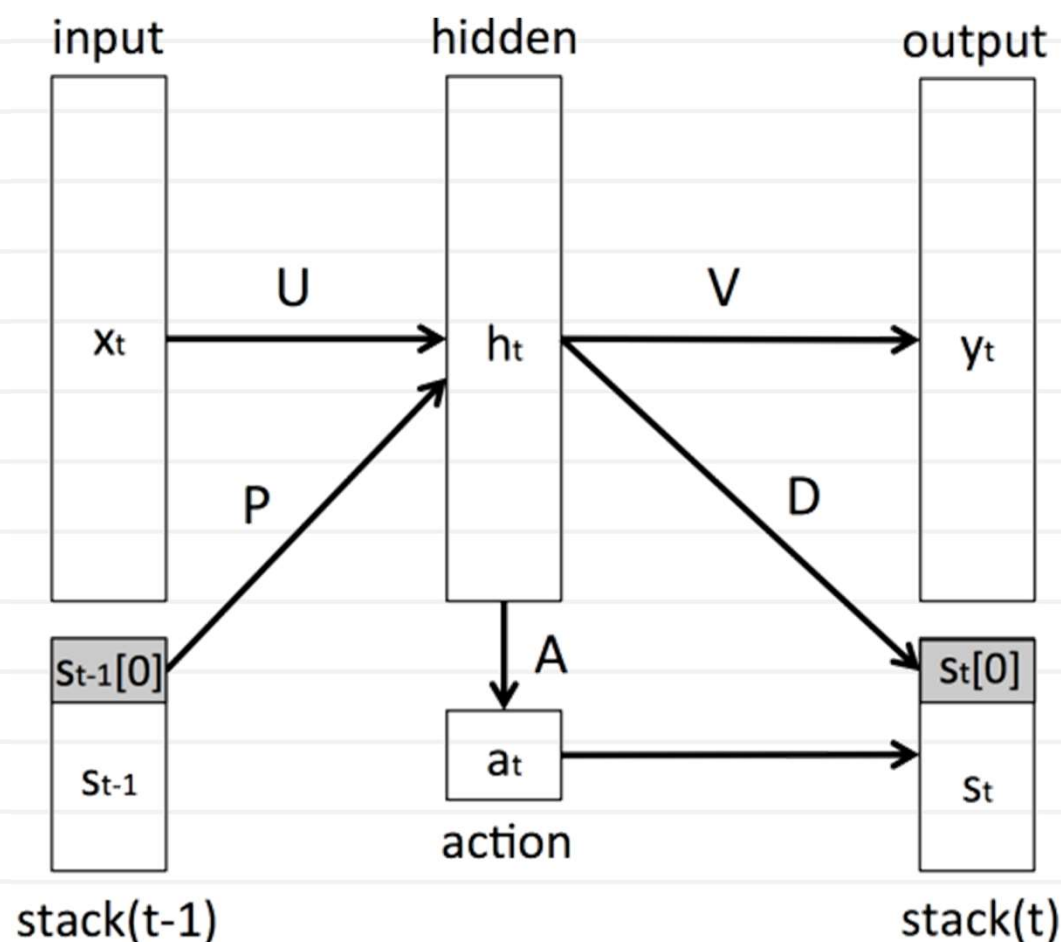
Stack augmented RNNs

- 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 (think RAM and CPU!)

[Joulin and Mikolov NIPS15]

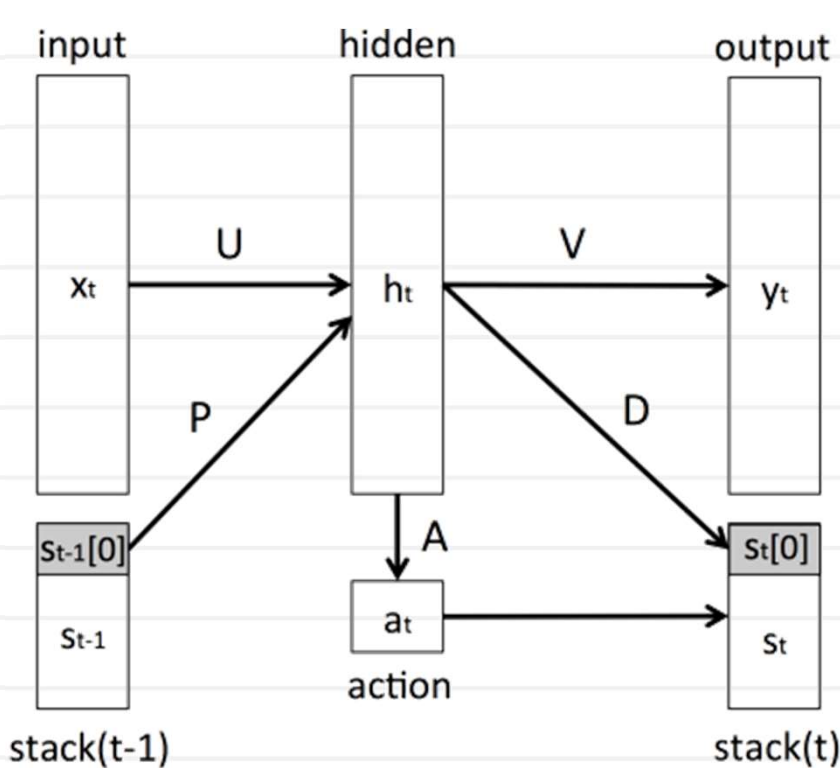
Stack augmented RNNs

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



[Joulin and Mikolov NIPS15]

Stack augmented RNNs



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

$$y_t = SM(Vh_t)$$

$$a_t = SM(Ah_t)$$

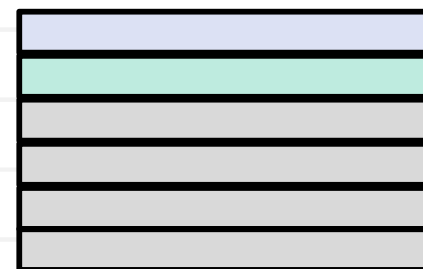
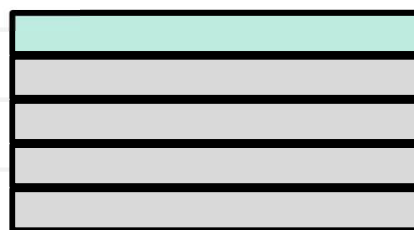
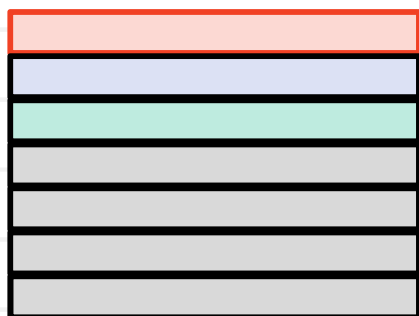
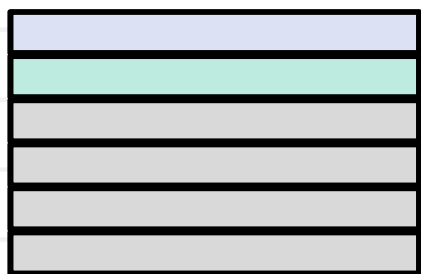
$$G(Dh_t)$$

s_{t-1}^k

PUSH

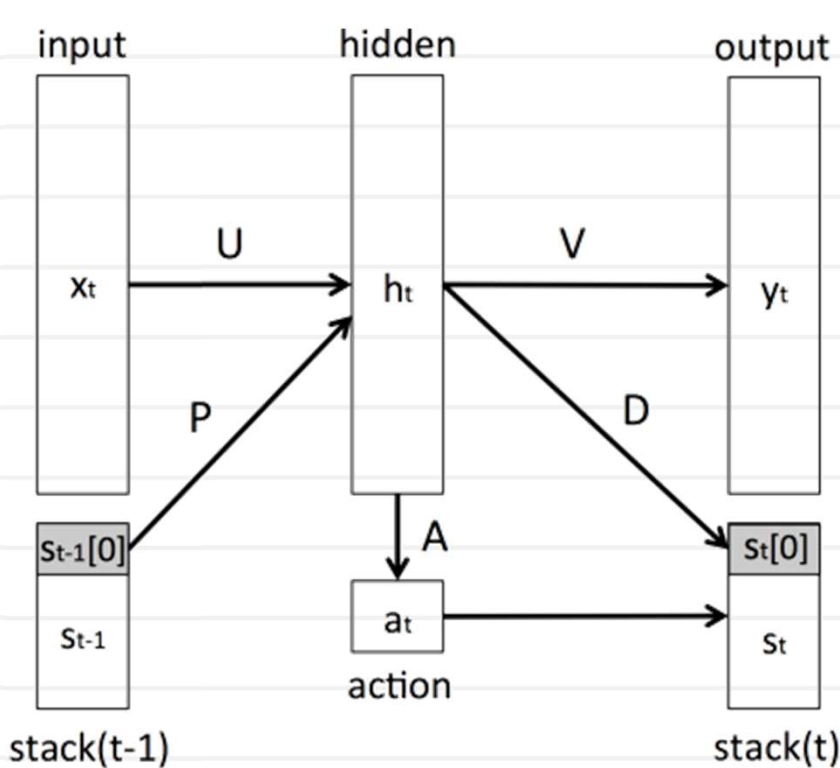
POP

NO



[Joulin and Mikolov NIPS15]

Stack augmented RNNs



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

$$y_t = SM(Vh_t)$$

$$a_t = SM(Ah_t)$$

$$S_t^0 = a_t[Push] G(Dh_t) + a_t[Pop] S_{t-1}^1 + a_t[No] S_{t-1}^0$$

$$S_t^i = a_t[Push] S_{t-1}^{i-1} + a_t[Pop] S_{t-1}^{i+1} + a_t[No] S_{t-1}^i$$

[Joulin and Mikolov NIPS15]

Binary addition with stack-RNN

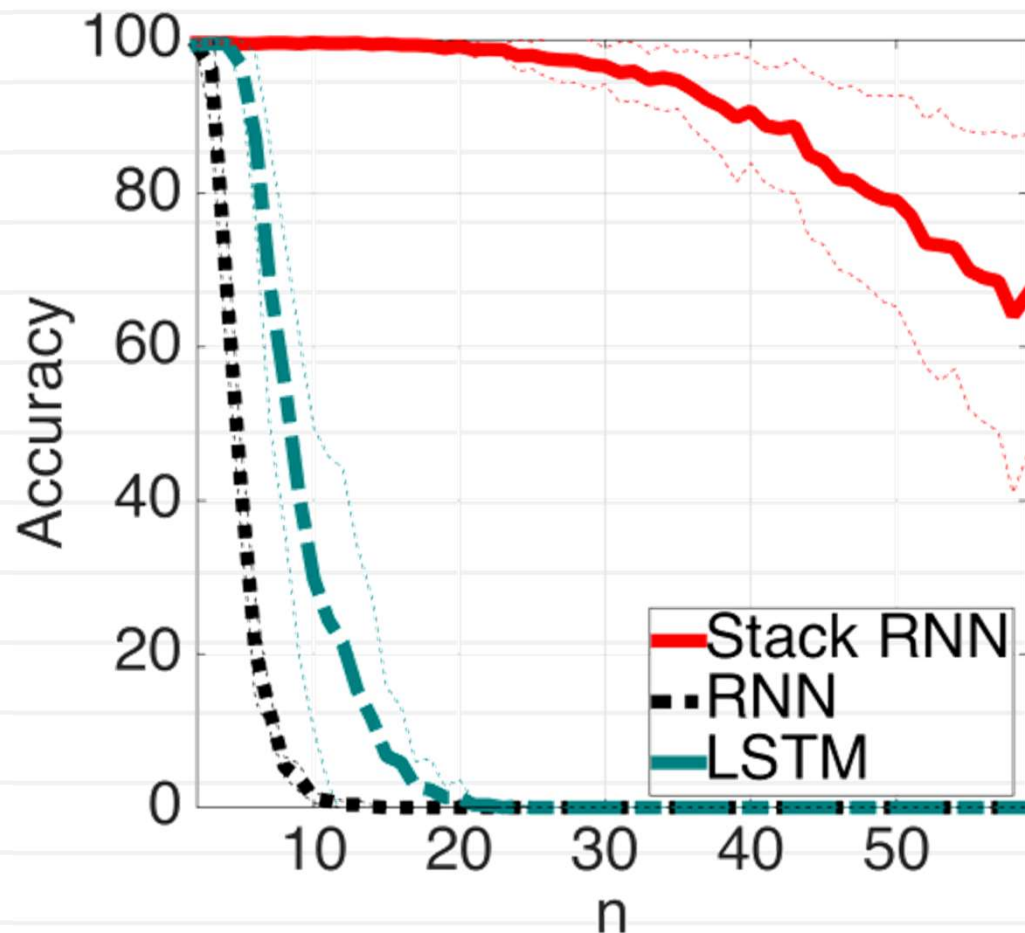
Goal: train a network that can add binary numbers.

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	1	0	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	0	0	0	1	0	Number 2
Stack 4:								1	0	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	0	0	1	0	0	0	0	1	1	0		Carry
Stack 6:			1	0	0	0	1	1						0	1	0	0	0	0	1	-1		Number 1
Stack 7:																							Junk
Stack 8:																							Junk
Stack 9:																							Junk
Stack 10:																							Junk

PUSH POP

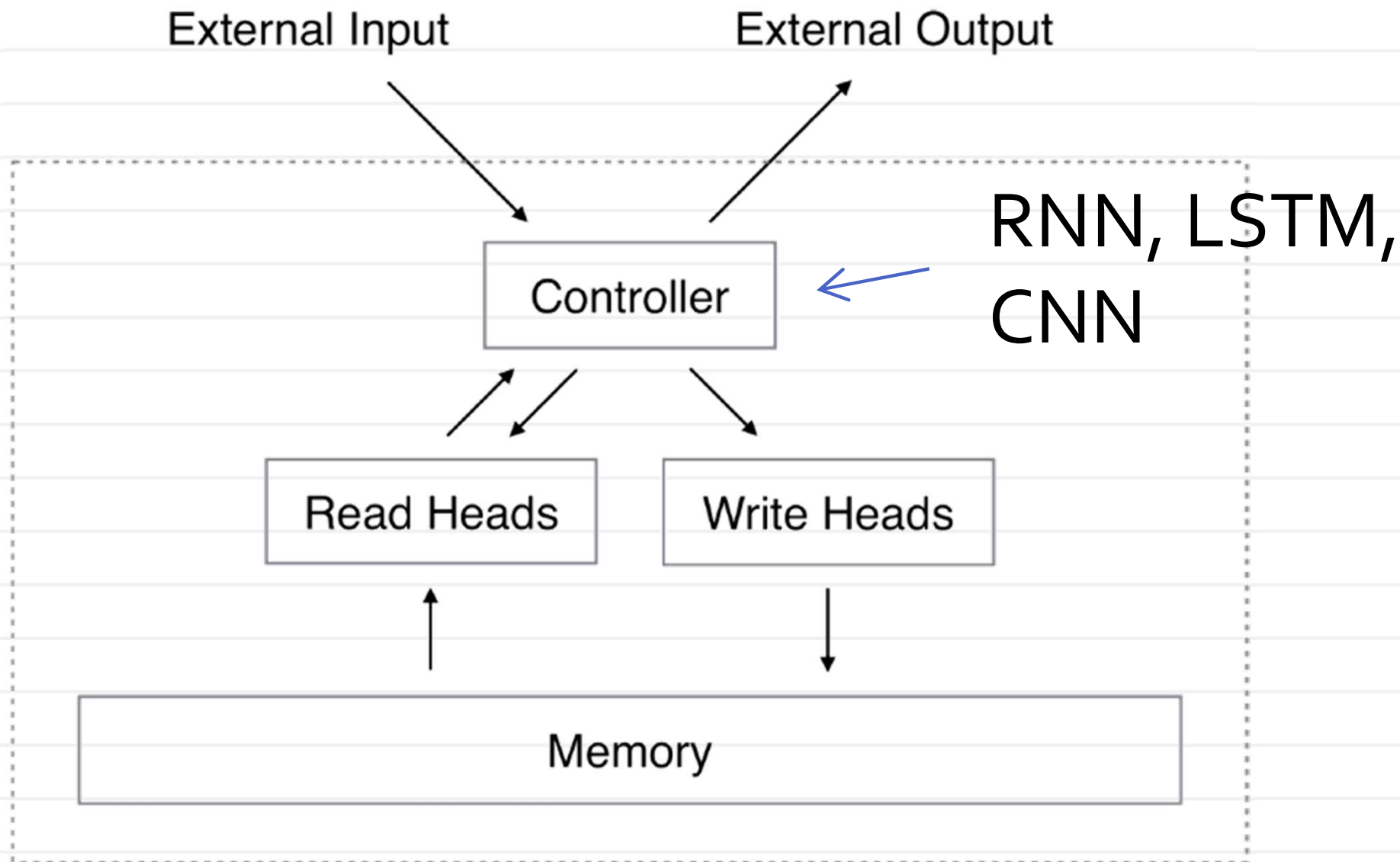
NB: the answer is reversed, i.e. $101+11 = 0001$

Binary addition with stack-RNN



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

Neural Turing Machine



[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

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)