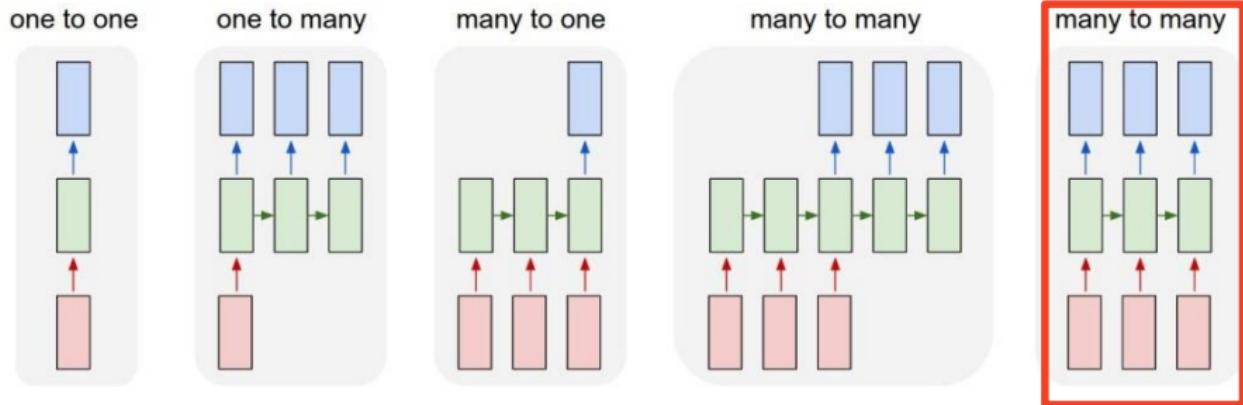


# Lecture 11: Sequence-to-sequence architectures. Neural attention and memory.

# Learning settings

slide credit: A. Karpathy



One-to-one: image to class label

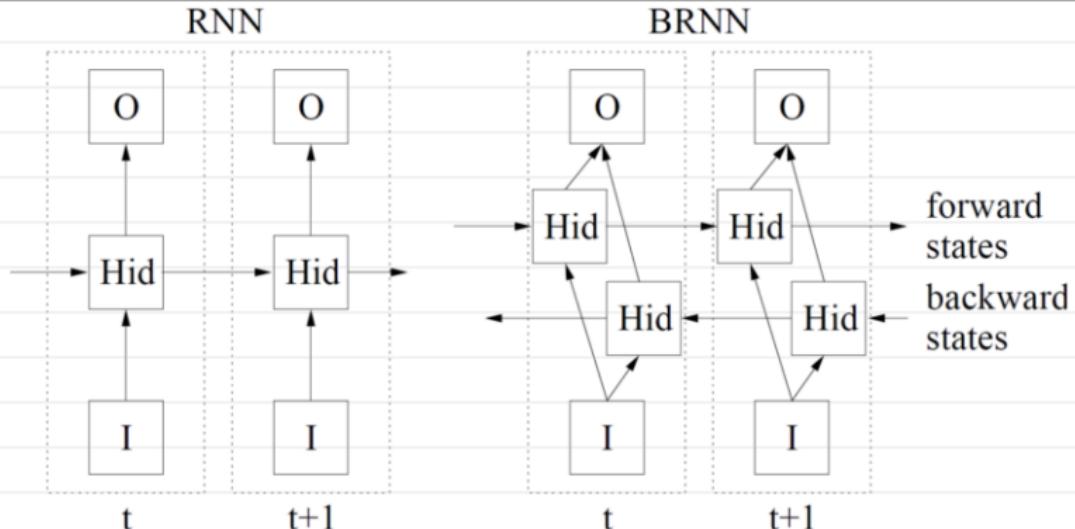
One-to-many: text generation/image captioning

Many-to-one: sentiment analysis

Many-to-many 1: machine translation

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

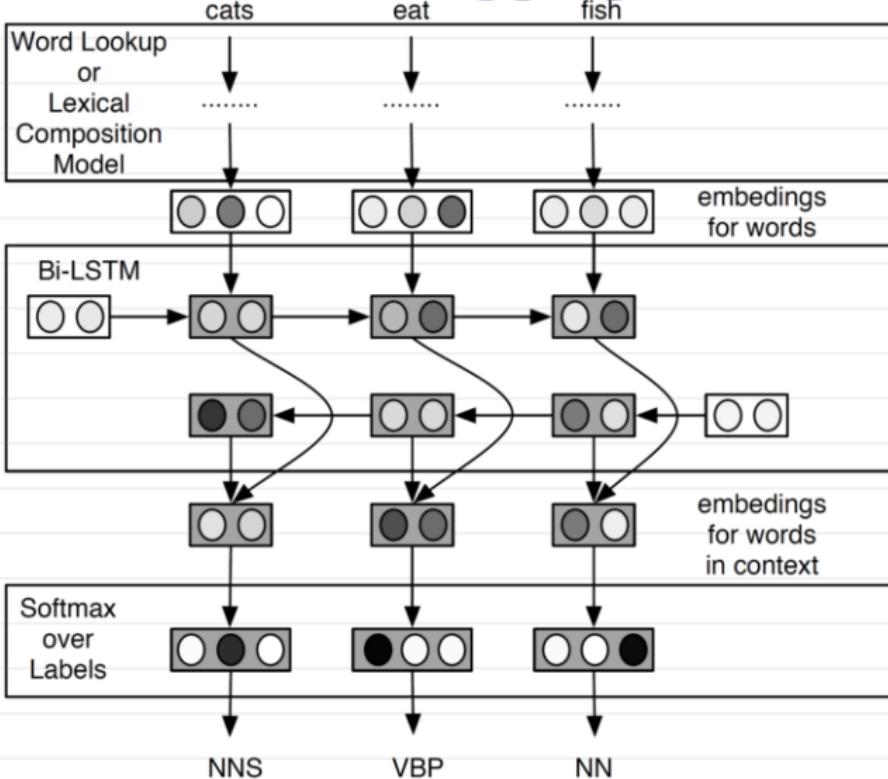
# Bi-directional RNN



```
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]

# Bi-LSTM POS tagging



$$l_i = \tanh(\mathbf{L}^f s_i^f + \mathbf{L}^b s_i^b + \mathbf{b}_l)$$

[Ling et al.  
EMNLP15]

# Bi-LSTM POS tagging

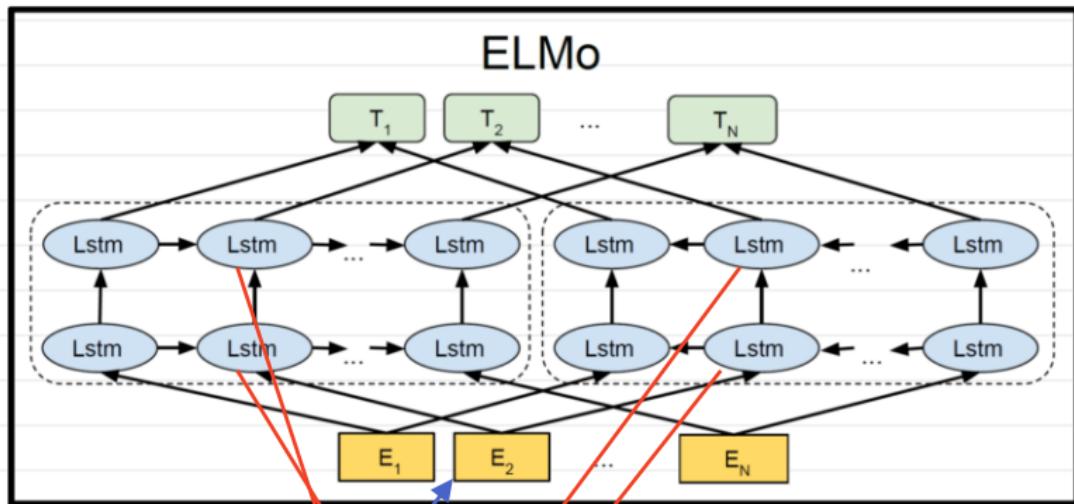
	acc	parameters	words/sec
Word Lookup	96.97	2000k	6K
Convolutional (S&Z)	96.80	42.5k	4K
Forward RNN	95.66	17.5k	4K
Backward RNN	95.52	17.5k	4K
Bi-RNN	95.93	40k	3K
Forward LSTM	97.12	80k	3K
Backward LSTM	97.08	80k	3K
Bi-LSTM $d_{CS} = 50$	97.22	70k	3K
Bi-LSTM	<b>97.36</b>	150k	2K

[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

# BiLM/ELMo



- BiLM is learned for language modeling (next word prediction)
- ELMo coefficients are fine-tuned for the new task:

$$\text{ELMo}_k^{\text{task}} = E(R_k; \Theta^{\text{task}}) = \gamma^{\text{task}} \sum_{j=0}^L s_j^{\text{task}} \mathbf{h}_{k,j}^{LM}$$

BiLM word descriptor

[Peters et al. ACL18]

# BiLM/ELMO

Source	Nearest Neighbors
GloVe play	playing, game, games, played, players, plays, player, Play, football, multiplayer
biLM	Chico Ruiz made a spectacular <u>play</u> on Alusik 's grounder {...}
	Olivia De Havilland signed to do a Broadway <u>play</u> for Garson {...}

Table 4: Nearest neighbors to “play” using GloVe and the context embeddings from a biLM.

## SNLI task [Bowman et al. EMNLP15]:

Text	Judgments	Hypothesis
A man inspects the uniform of a figure in some East Asian country.	contradiction C C C C	The man is sleeping
An older and younger man smiling.	neutral N N N N N	Two men are smiling and laughing at the cats playing on the floor.
A black race car starts up in front of a crowd of people.	contradiction C C C C	A man is driving down a lonely road.
A soccer game with multiple males playing.	entailment E E E E E	Some men are playing a sport.
A smiling costumed woman is holding an umbrella.	neutral N N E C N	A happy woman in a fairy costume holds an umbrella.

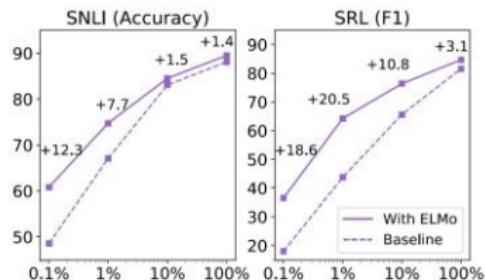
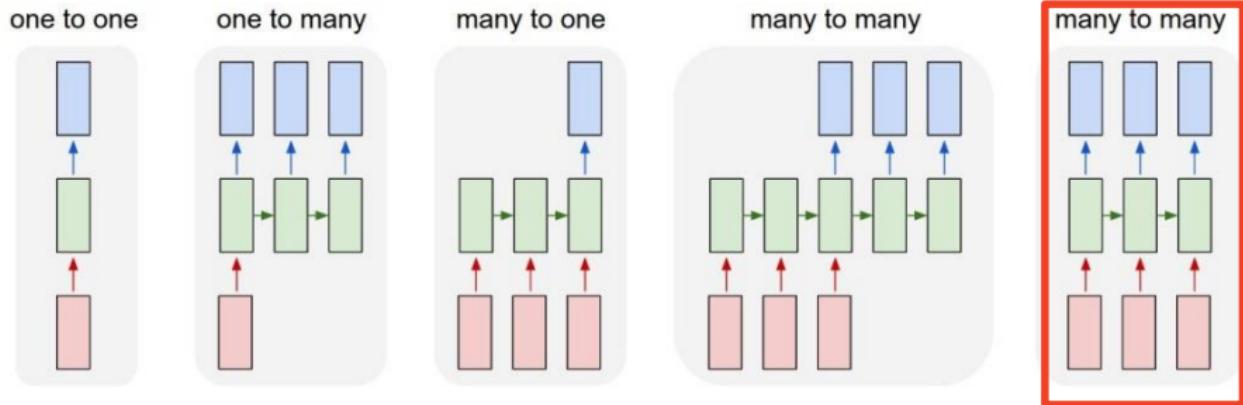


Figure 1: Comparison of baseline vs. ELMo performance for SNLI and SRL as the training set size is varied from 0.1% to 100%.

## [Peters et al. ACL18]

# 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

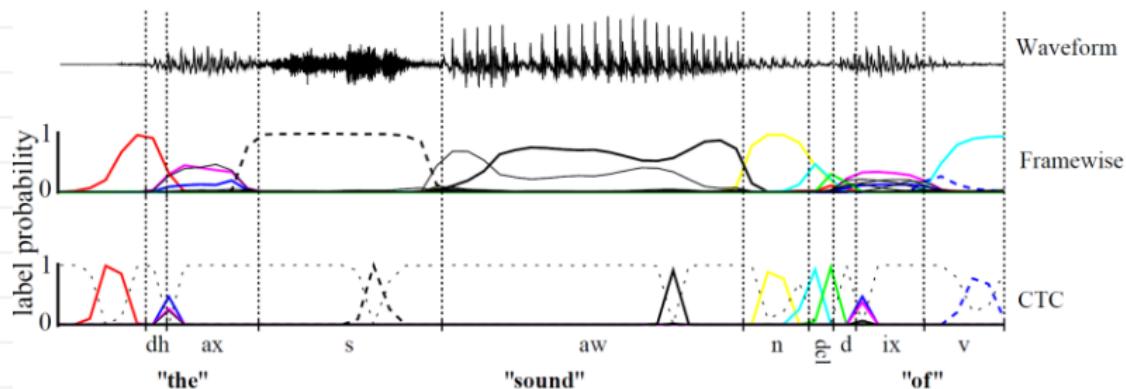
Many-to-many 1: machine translation

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

# Online seq2seq with monotonic alignment

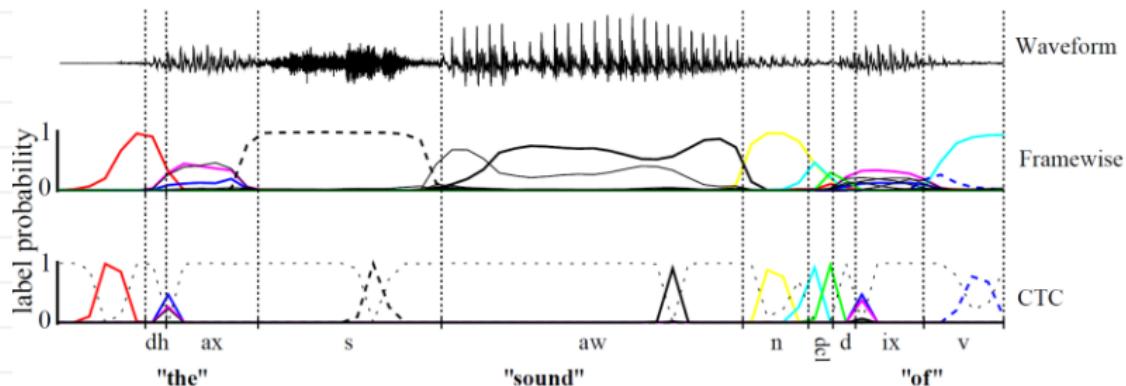
Many problems are sequence 2 sequence with monotonic alignment:

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



[Graves et al. 2006]

# Online seq2seq with monotonic alignment



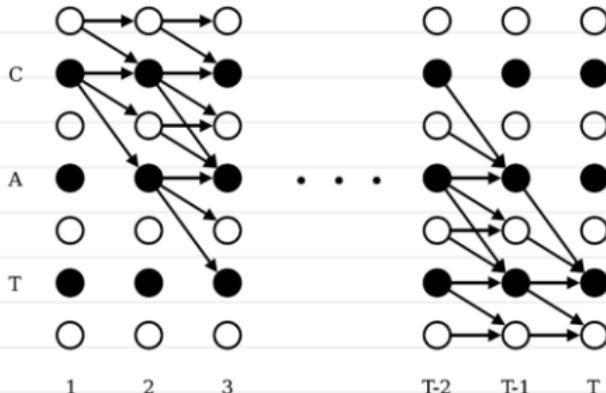
Decoding: 'aaa\_\_bb\_c\_\_\_ddaa' → abcda

What should be the loss that encourage correct parsing?

Answer: connectionist temporal classification (CTC) loss

[Graves et al. 2006]

# CTC-loss

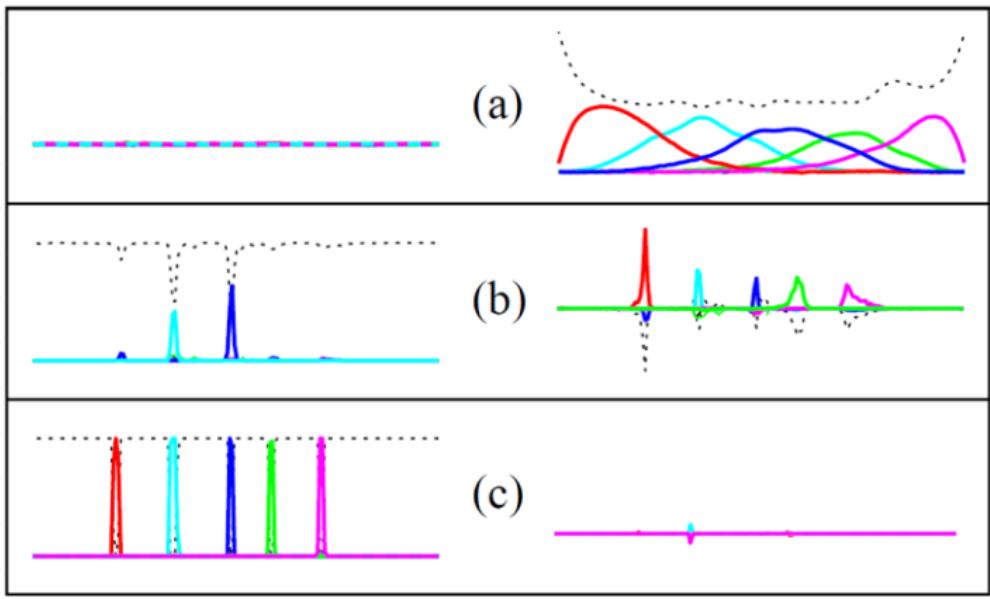


- Augment the output state with *blank*
- Predict probabilities of each symbol (inc. blank) at each time moment
- Compute the probability of each lattice vertex under correct paths using forward-backward
- Push log-probabilities up (*ML training*) proportionally to the current probability

[Graves et al. 2006]

# Evolution of the CTC signal

GT sequence:



Prediction

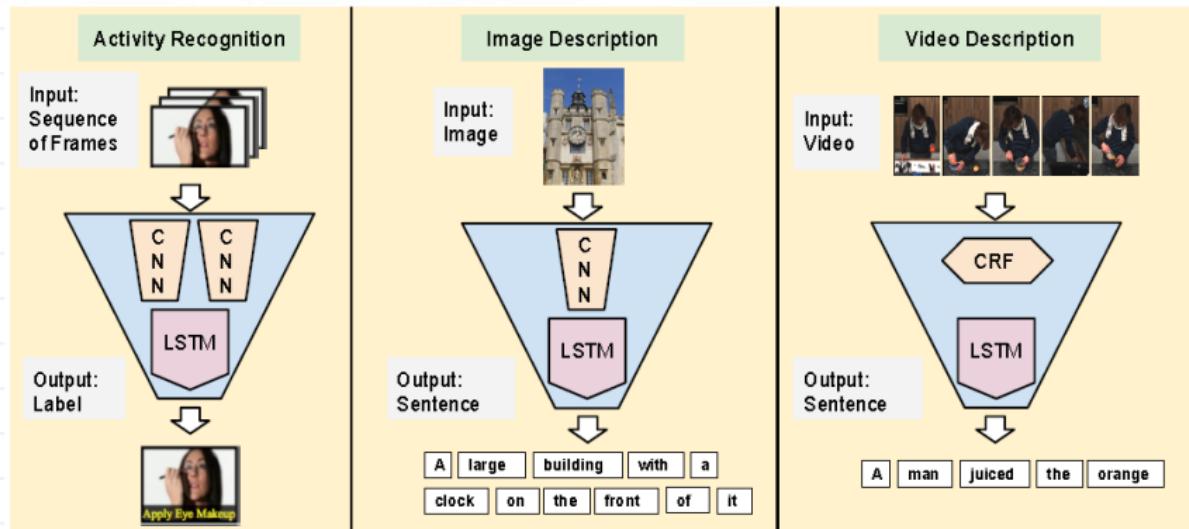
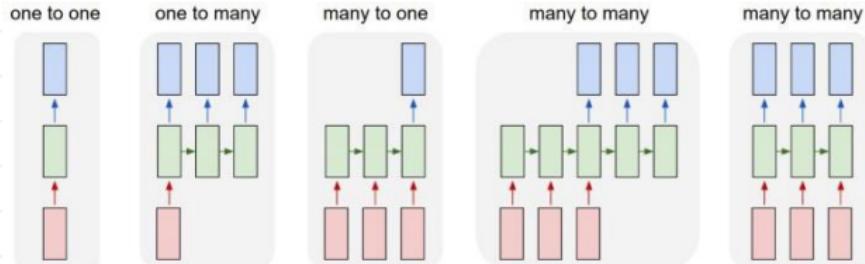
Gradient w.r.t. prediction

[Graves et al. 2006]

# LSTM demo: handwriting recognition

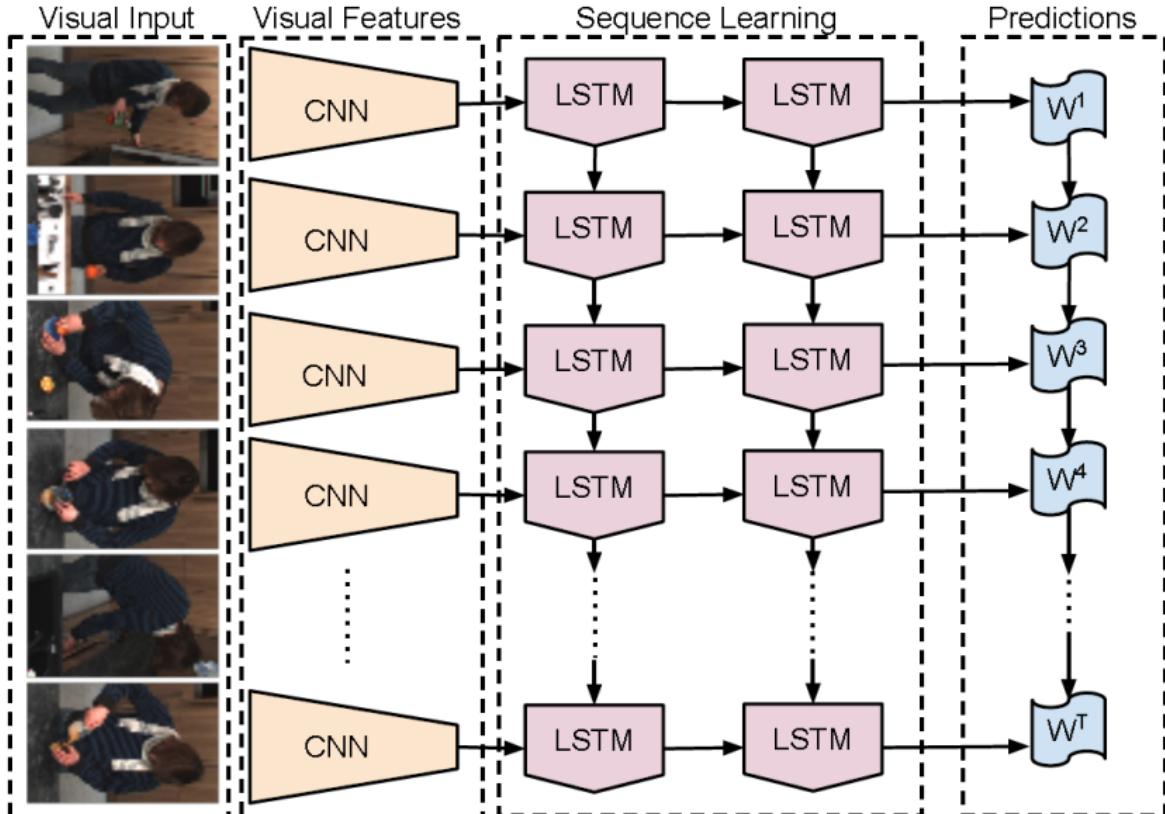
LSTM RNN Demo by Nikhil Buduma:  
<https://www.youtube.com/watch?v=mLxsbWAYIpw>

# Image/video captioning



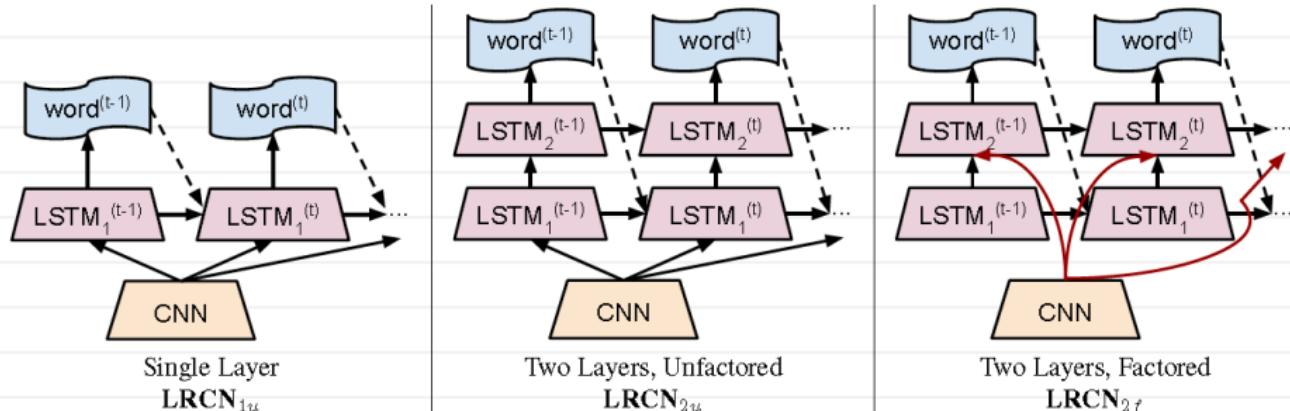
[Donahue et al. 2015]

# 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 <sub>1u</sub>	14.1	31.3	39.7	24
LRCN <sub>2u</sub>	3.8	12.0	17.9	80
LRCN <sub>2f</sub>	<b>17.5</b>	<b>40.3</b>	<b>50.8</b>	<b>9</b>
LRCN <sub>4f</sub>	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.



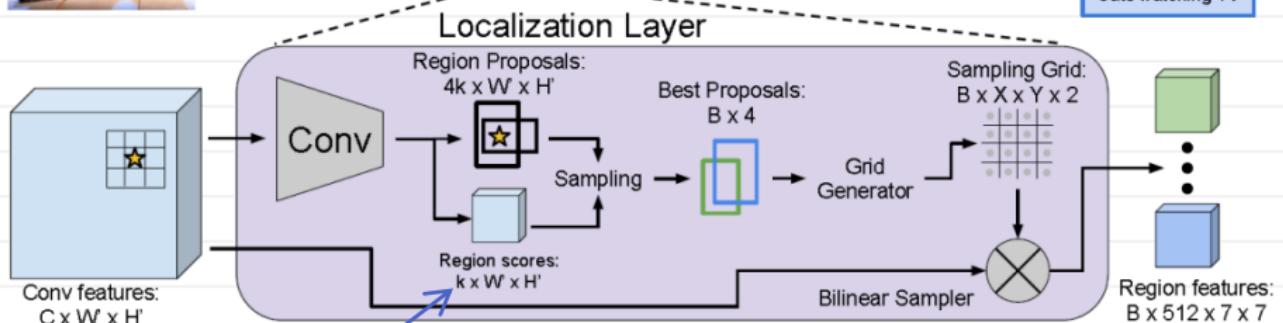
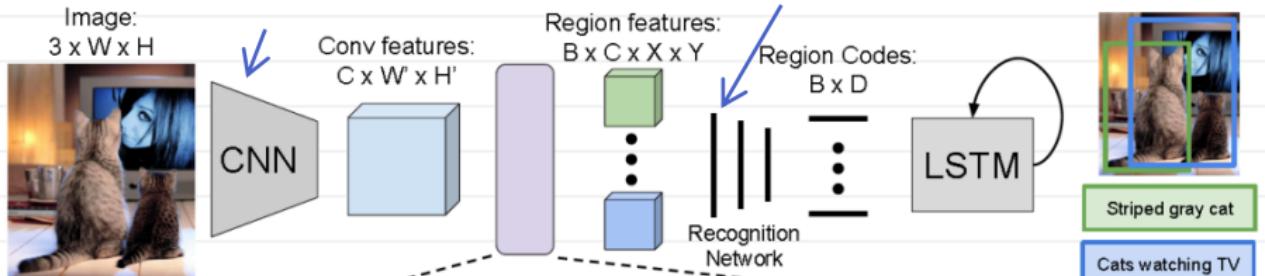
A person holding a cell phone in their hand.

[Donahue et al. 2015]

# End-to-end dense image captioning

VGG-16 conv

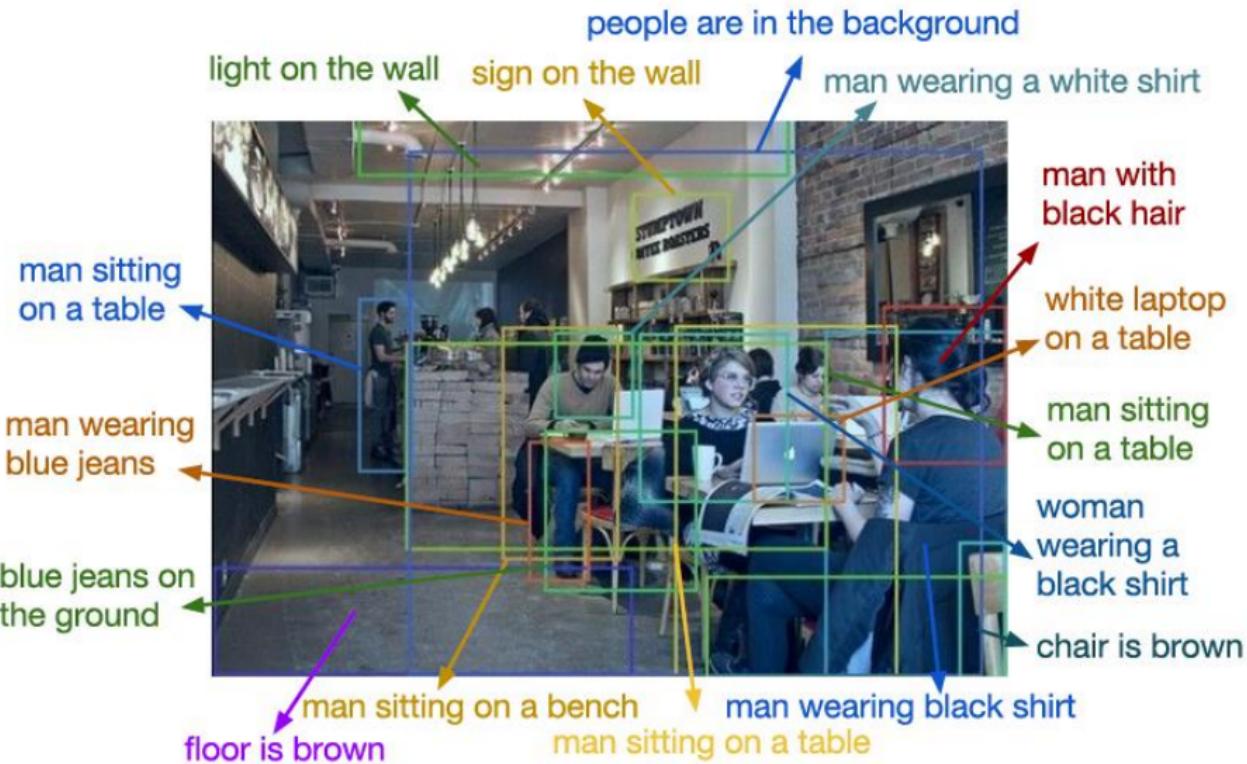
fully-connected, 2 layers+dropout



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

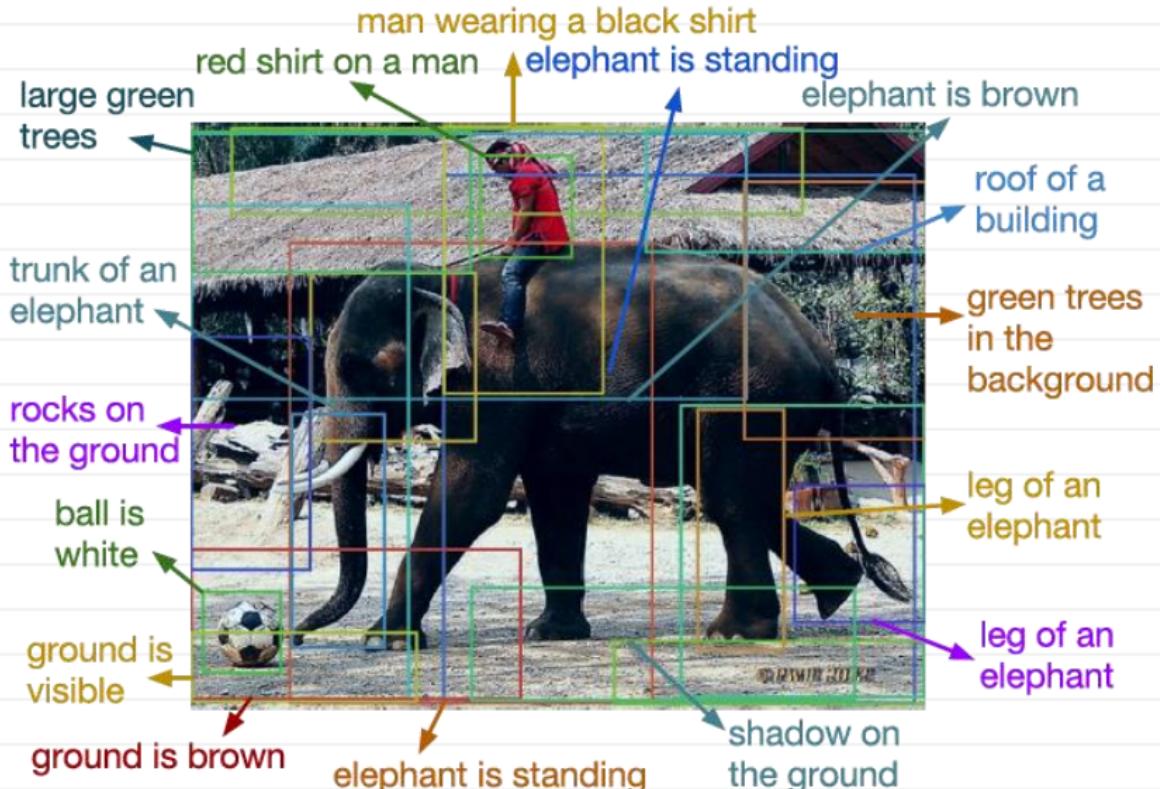
[Johnson et al, CVPR16]

# End-to-end dense image captioning



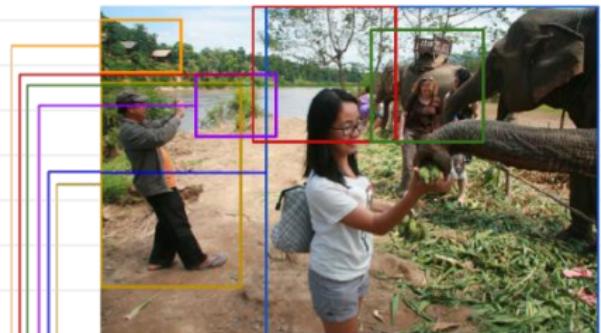
[Johnson et al, CVPR16]

# 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 shirt.  
▶ Small houses on the hillside  
A woman feeding an elephant  
A woman wearing a white shirt and shorts  
A man taking a picture

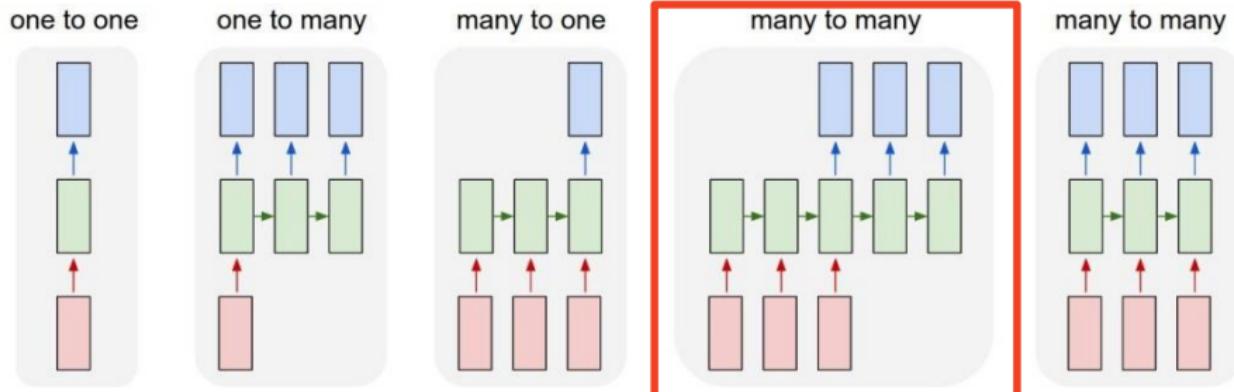
A man wearing an orange shirt  
An elephant taking food from a woman  
A woman wearing a brown shirt  
A woman wearing purple clothes  
A man wearing blue flip flops  
Man taking a photo of the elephants  
Blue flip flop sandals  
The girl's white and black handbag  
The girl is feeding the elephant  
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  
Elephant's trunk being touched by a hand  
A man taking a picture holding a camera  
Elephant with carrier on its 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]

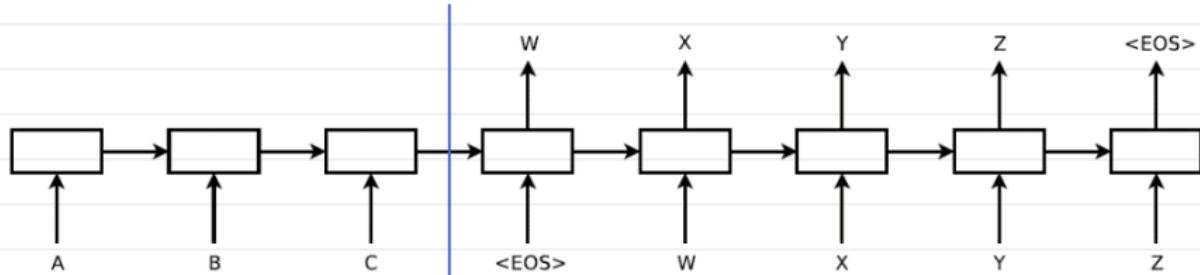
# Learning settings

slide credit: A. Karpathy



aka "seq2seq"

# Sequence-to-sequence machine translation



Important notes:

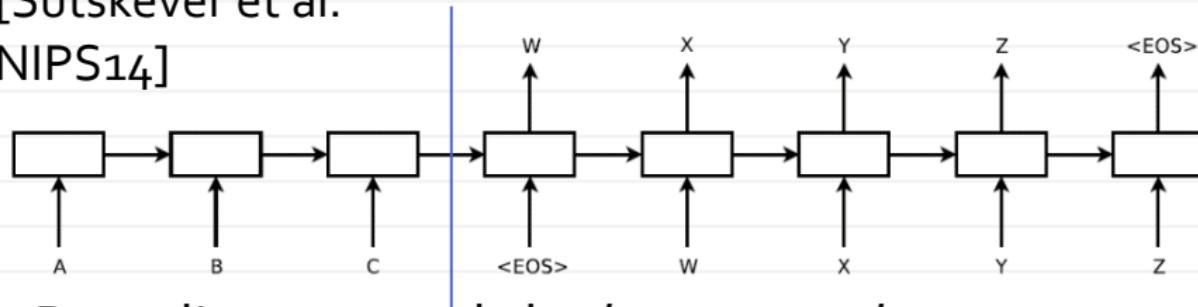
1. Fixed lexicon (160,000 English, 80,000 French) + 'UNK' word
2. Deep (four layers, 1000 cells in each)
3. Reversing input sequence helps a lot
4. Using two different LSTMs
5. Decoding proceeds by *beam search*

[Sutskever et al. NIPS14]

# Sequence-to-sequence machine translation

[Sutskever et al.

NIPS14]



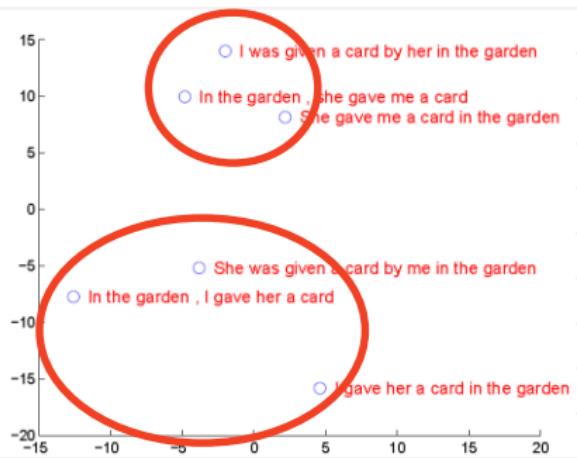
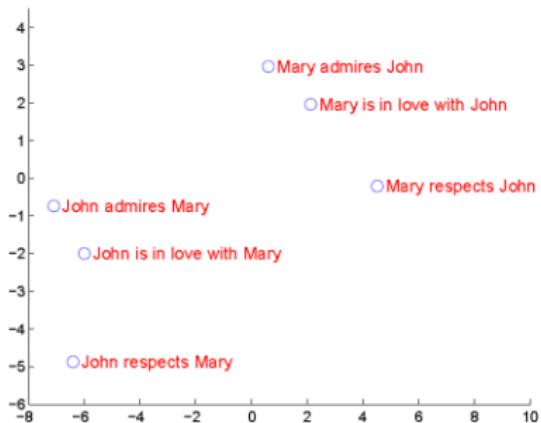
Decoding proceeds by *beam search*:

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]

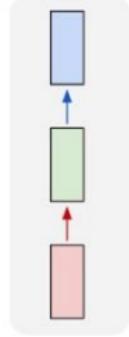
# Sequence-to-sequence machine translation

Type	Sentence
<b>Our model</b>	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 .
<b>Truth</b>	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 .
<b>Our model</b>	" 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 .
<b>Truth</b>	" 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 .
<b>Our model</b>	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 " .
<b>Truth</b>	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 " .

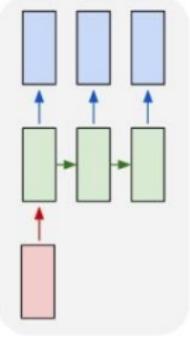
[Sutskever et al. NIPS14]

# Sequence-to-sequence machine translation

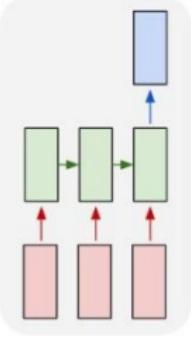
one to one



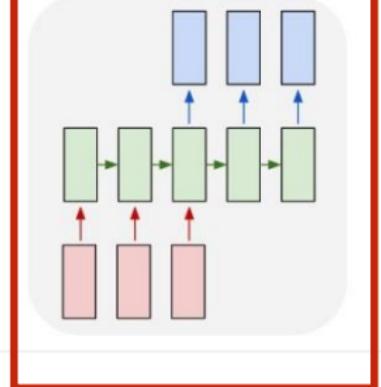
one to many



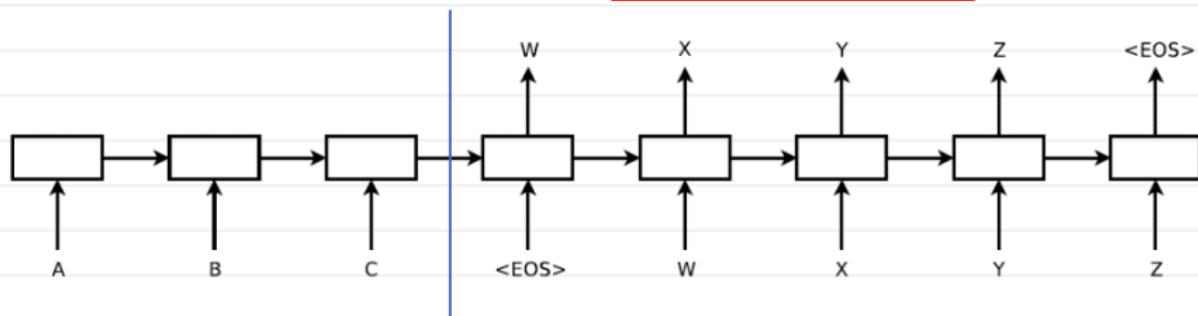
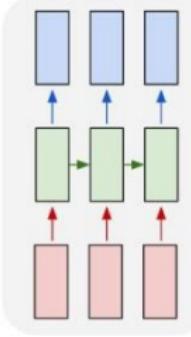
many to one



many to many

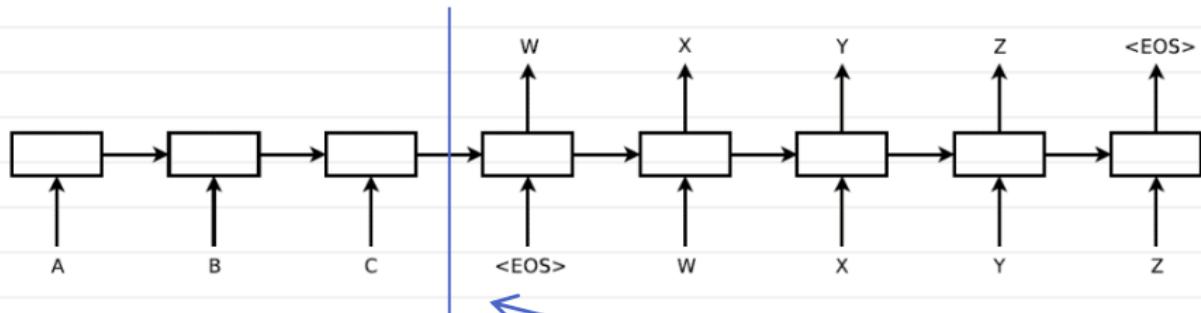


many to many



[Sutskever et al. NIPS14]

# Sequence-to-sequence machine translation



Problem:  
all the meaning has to be carried from here

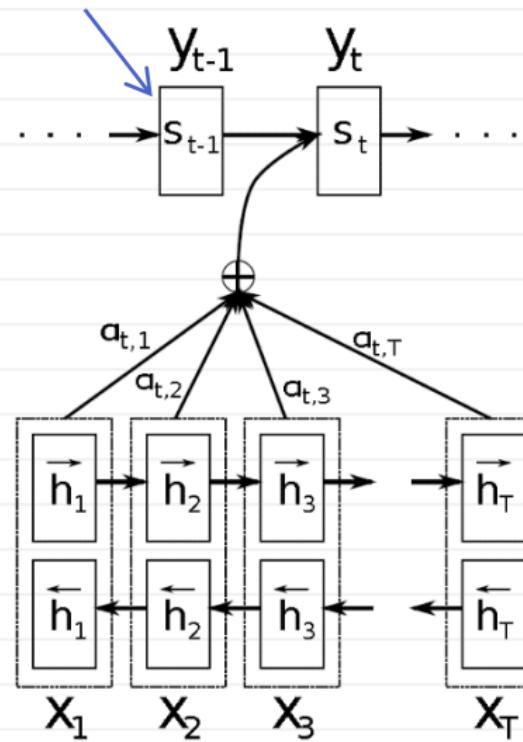
- Large memory needed
- Information has to survive for a very long time



[Sutskever et al. NIPS14]

# Translation with attention

decoder RNN



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

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

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

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

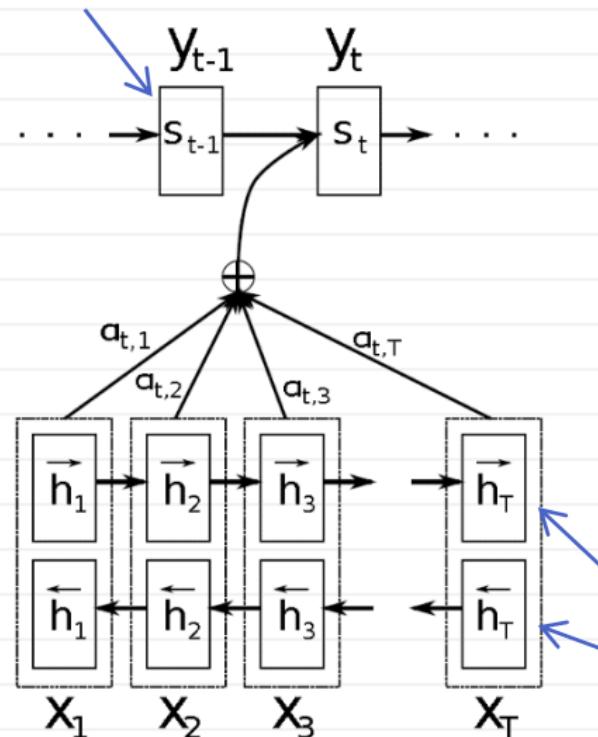
encoder RNN 1

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

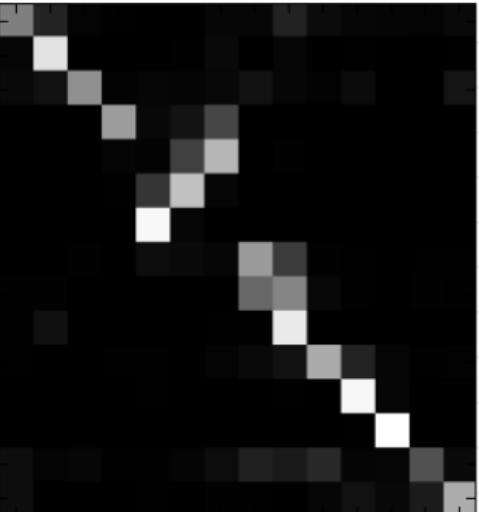
encoder RNN 1  
encoder RNN 2

[Bahdanau et al. 2015]

# Translation with attention

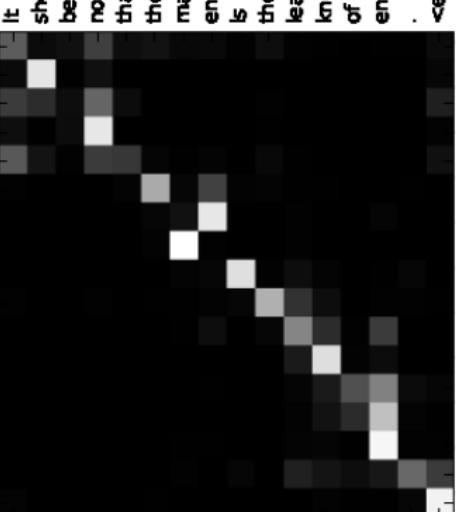
The  
agreement  
on  
the  
European  
Economic  
Area  
was  
signed  
in  
August  
1992  
. <end>

L'  
accord  
sur  
la  
zone  
économique  
européenne  
a  
été  
signé  
en  
août  
1992  
<end>

An attention matrix visualization for the first sentence. It is a 10x10 grid where each cell's color represents the attention weight between the corresponding source word and target word. A prominent diagonal band of high attention (lighter shades) runs from the top-left ('The') to the bottom-right ('environments'). Other words like 'European' and 'Economic' also have some attention weights shown.

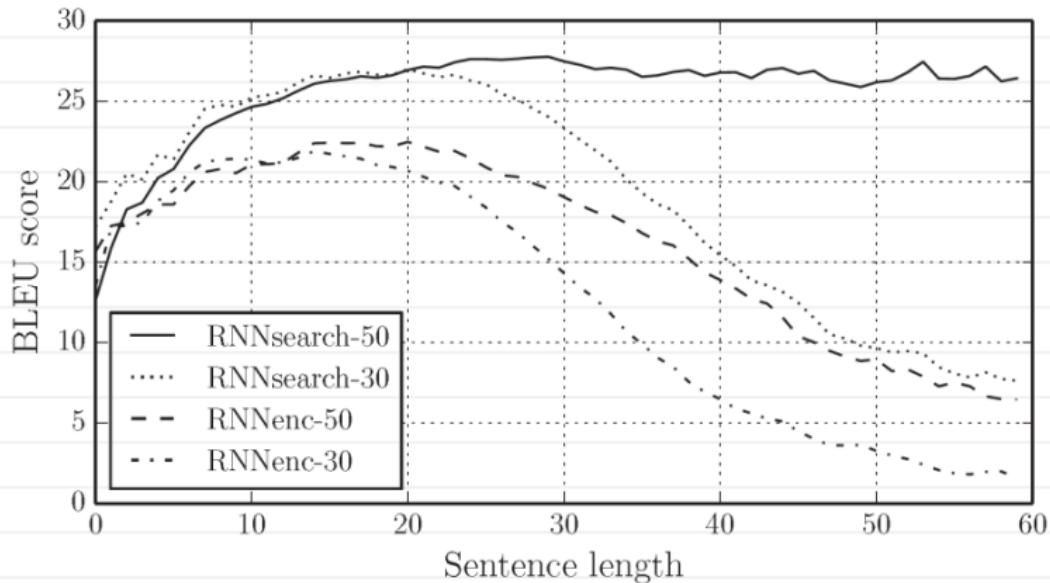
It  
should  
be  
noted  
that  
the  
marine  
environment  
is  
the  
least  
known  
of  
environments  
. <end>

Il  
convient  
de  
noter  
que  
l'  
environnement  
marin  
est  
le  
moins  
connu  
de  
l'  
environnement  
. <end>

An attention matrix visualization for the second sentence. Similar to the first, it shows a strong diagonal band of high attention weights. The words 'should', 'be', 'noted', 'that', 'the', 'marine', 'environment', 'is', 'the', 'least', 'known', 'of', and 'environments' all receive significant attention from their corresponding source words.

[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:

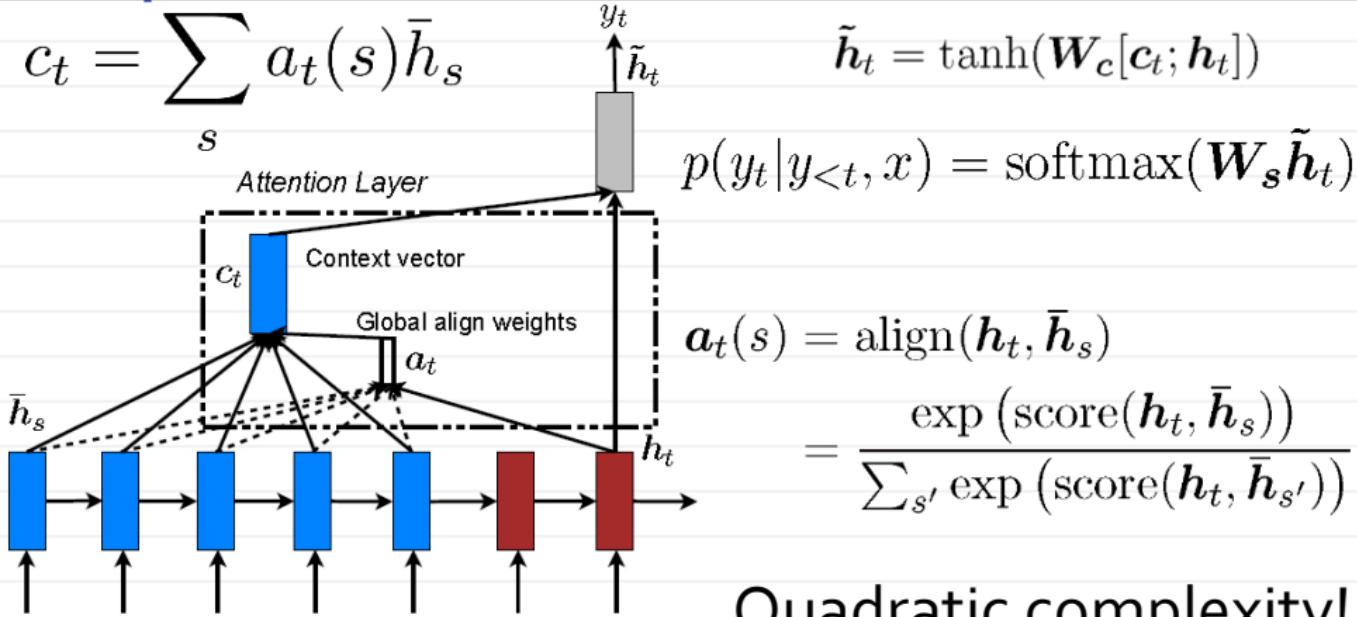
Un privilège d'admission est le droit d'un médecin de reconnaître un patient à l'hôpital ou un centre médical d'un diagnostic ou de prendre un diagnostic en fonction de son état de santé.

## 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



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

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

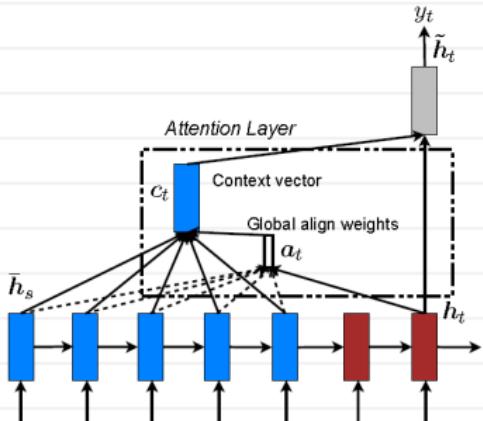
Quadratic complexity!

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

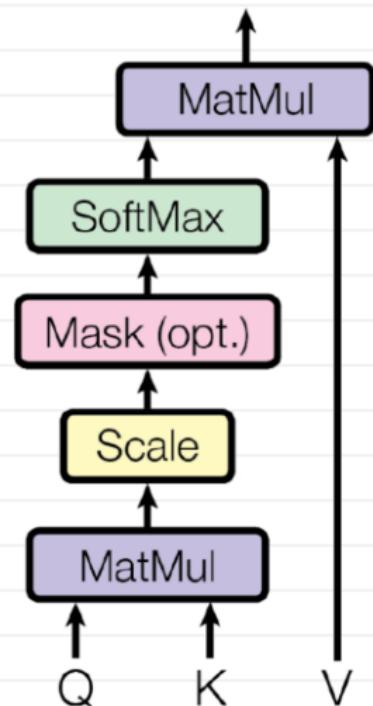
[Luong et al. 2015]

# Recap

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



# “Attention is all you need”: single head



General purpose “single-head” attention:

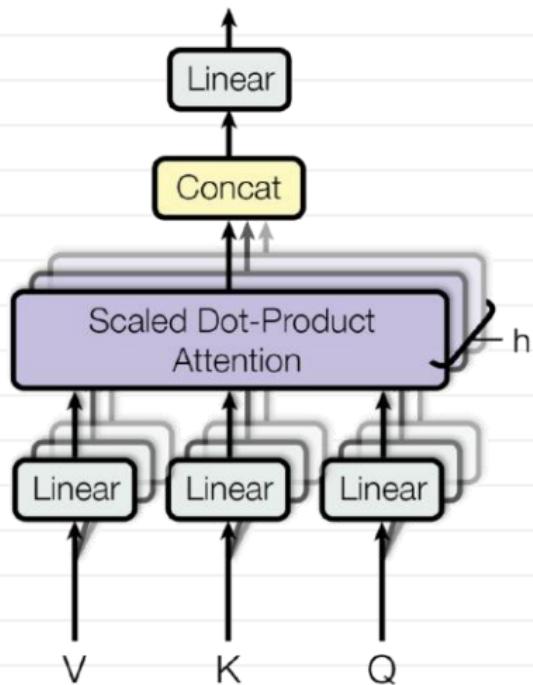
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

- 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]

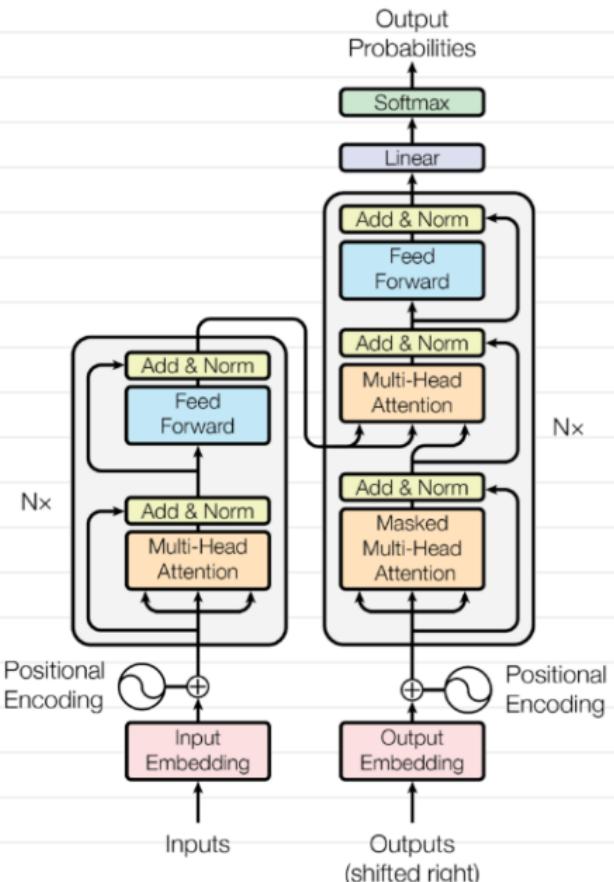


General purpose  
“multi-head” attention:

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

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

# Transformer architecture



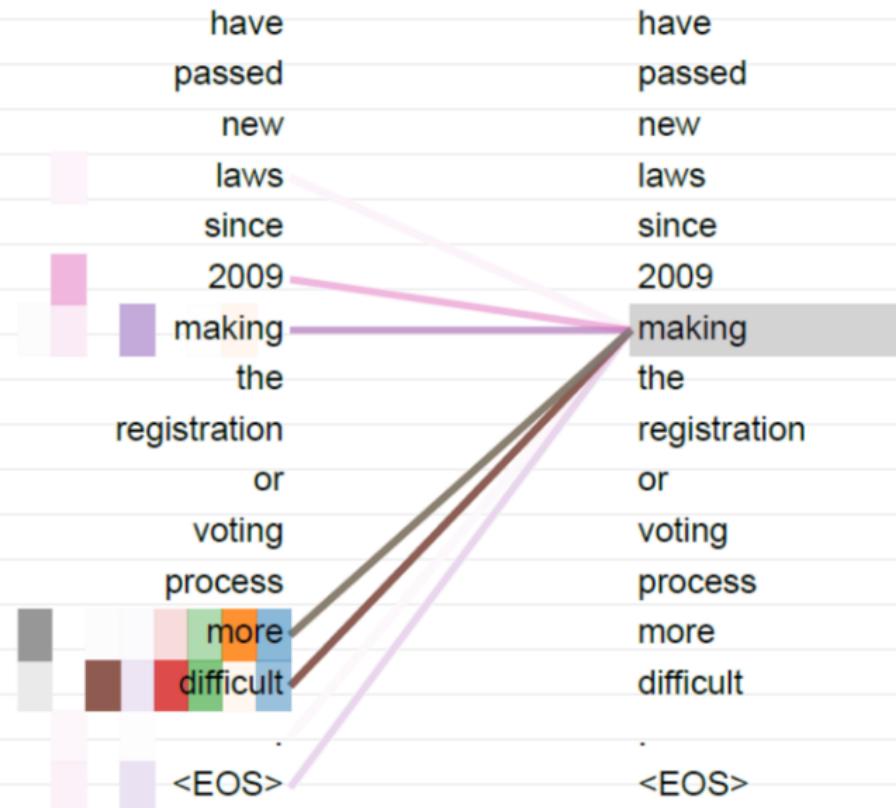
- Applying multi-head attention several times first for input, then for output
- Each unit is residual
- Emitting output one word at a time
- Positional encoding adds position-dependent features:

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

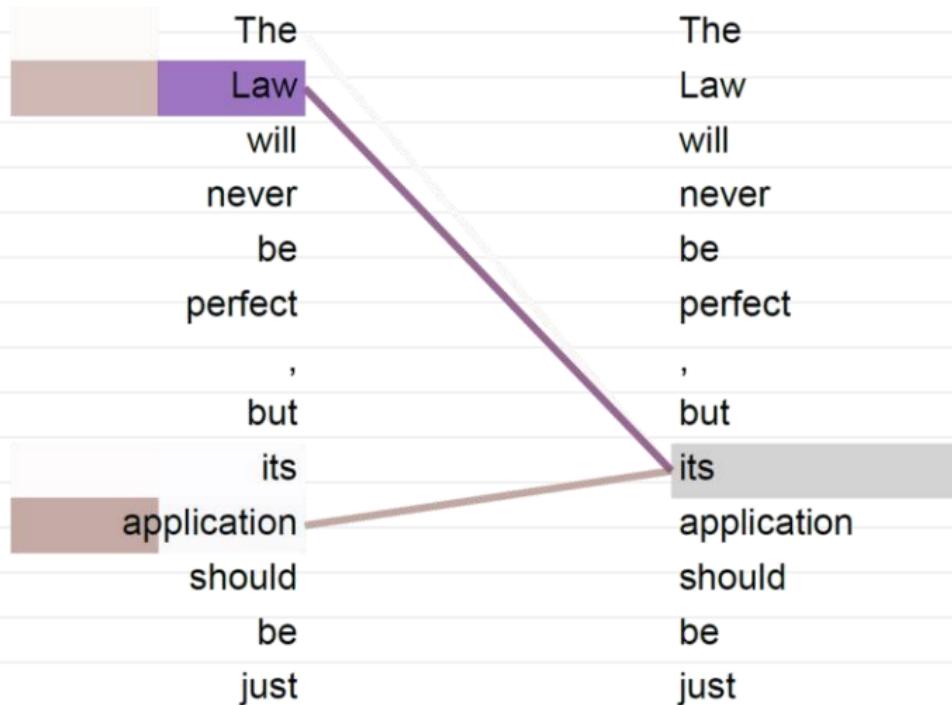
[Vaswani et al. NIPS17]

# “Attention is all you need”



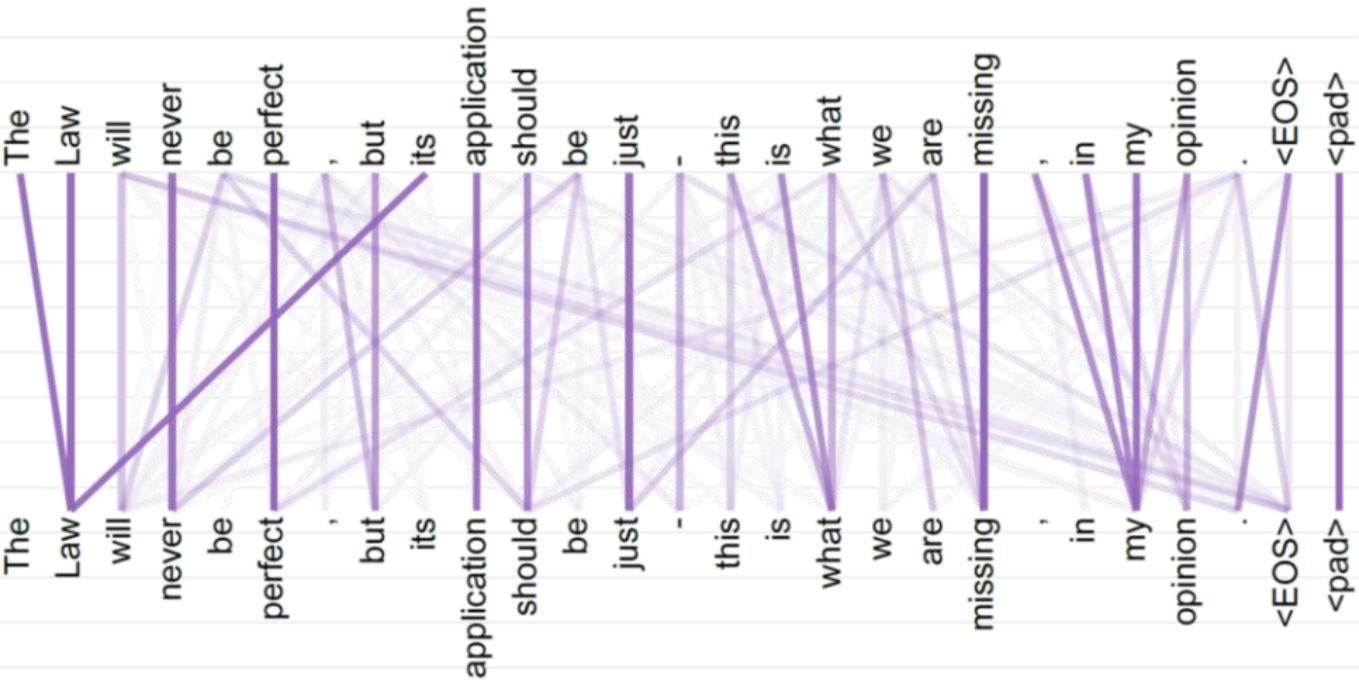
[Vaswani et al. NIPS17]

# “Attention is all you need”



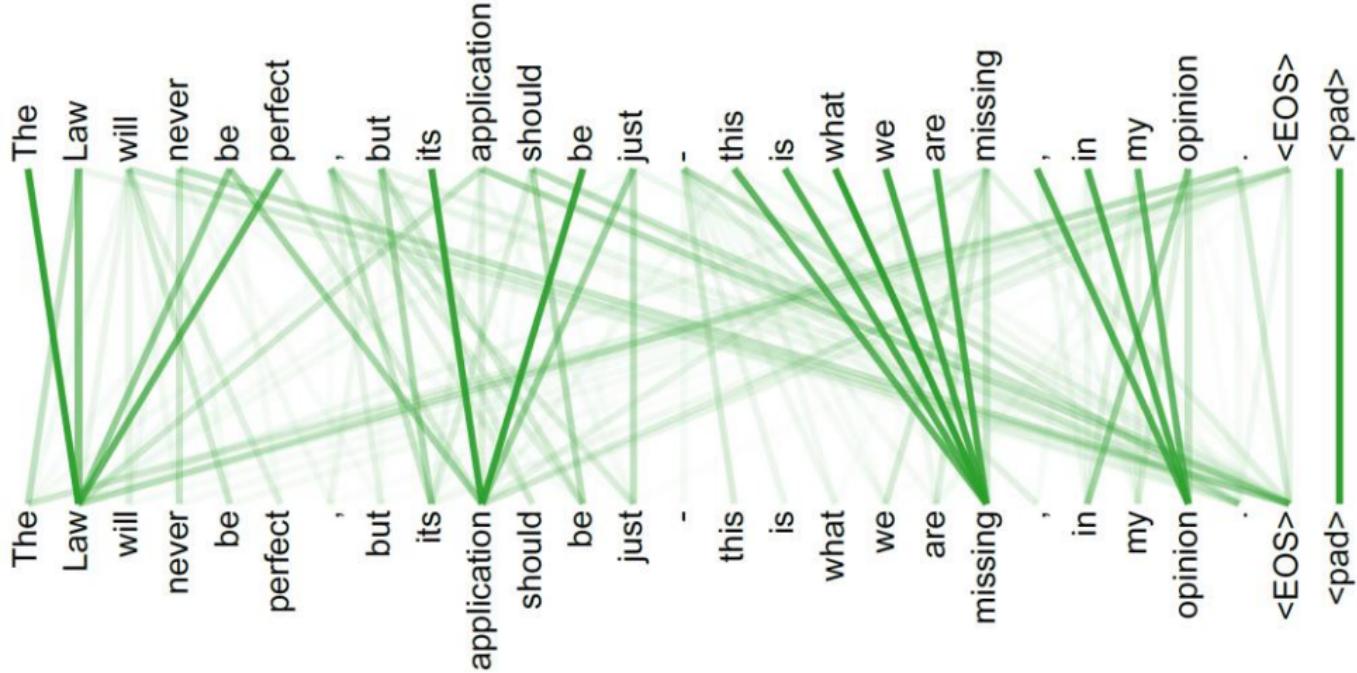
[Vaswani et al. NIPS17]

# “Attention is all you need”



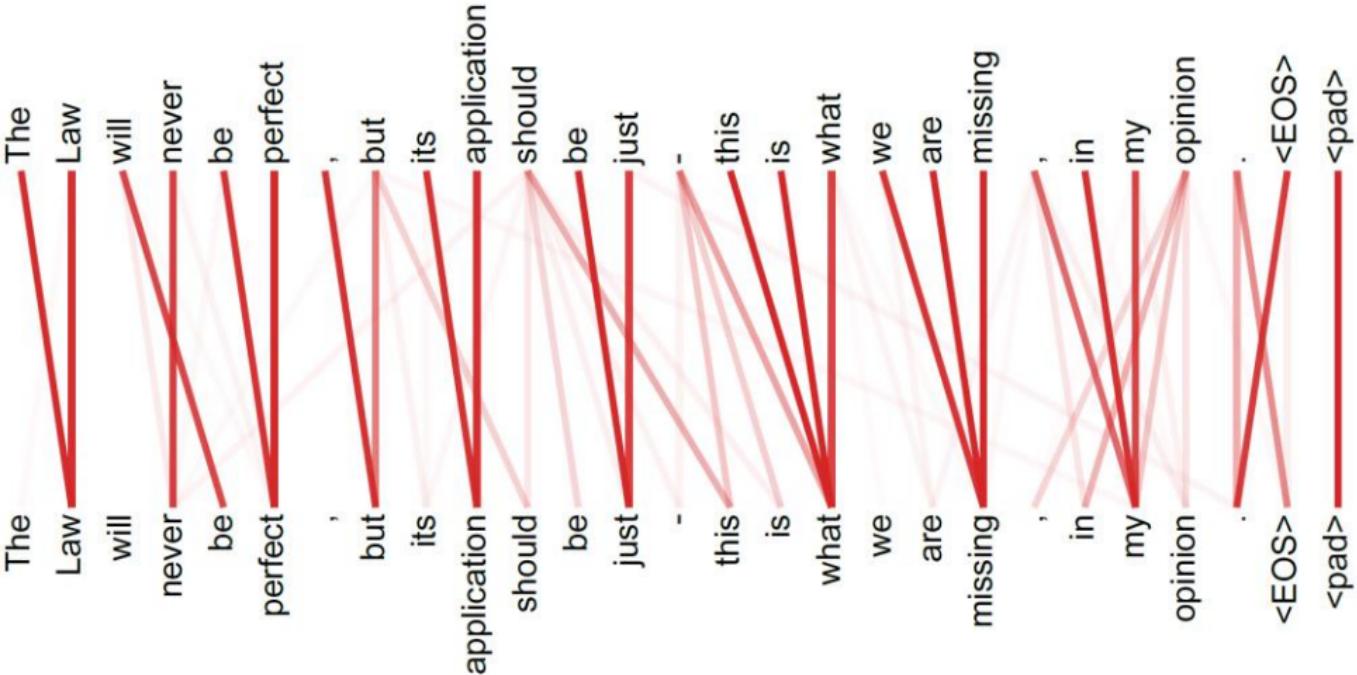
[Vaswani et al. NIPS17]

# "Attention is all you need"



[Vaswani et al. NIPS<sub>17</sub>]

# "Attention is all you need"



[Vaswani et al. NIPS17]

# “Attention is all you need”

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types.  $n$  is the sequence length,  $d$  is the representation dimension,  $k$  is the kernel size of convolutions and  $r$  the size of the neighborhood in restricted self-attention.

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

seq. of dilated conv

- Key property: attention layer complexity scales linearly in  $d$  (in attention layers) – same as depthwise separable convolution

[Vaswani et al. NIPS17]

# “Attention is all you need”

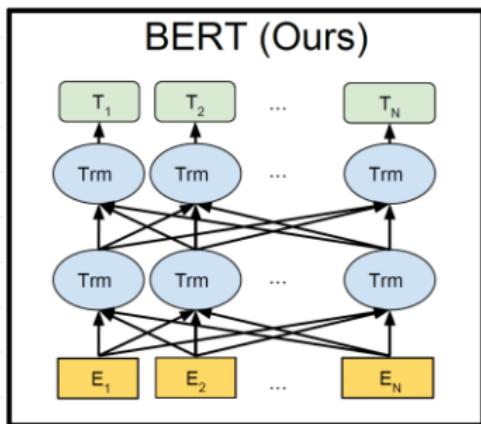
Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations  
Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [15]	23.75			
Deep-Att + PosUnk [32]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [31]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [8]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [26]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [32]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [31]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [8]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1		<b><math>3.3 \cdot 10^{18}</math></b>
Transformer (big)	<b>28.4</b>	<b>41.0</b>		$2.3 \cdot 10^{19}$

[Vaswani et al. NIPS17]

# BERT: pretraining by language modeling

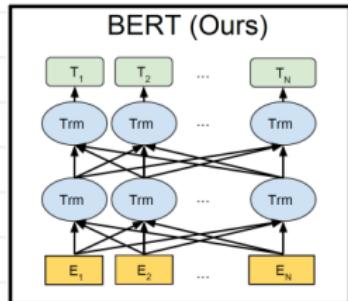
[Devlin et al. 2018]



- Pretraining six-level attention on language prediction tasks
- During training mask some words and replace others (~denoising autoenc idea)
- For sentence-level tasks, use the final layer representation of the first word
- **Fine-tune everything on the target task**

# BERT: pretraining by language modeling

[Devlin et al. 2018]



System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>91.1</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>81.9</b>

Table 1: GLUE Test results, scored by the GLUE evaluation server. The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set. OpenAI GPT = (L=12, H=768, A=12); BERT<sub>BASE</sub> = (L=12, H=768, A=12); BERT<sub>LARGE</sub> = (L=24, H=1024, A=16). BERT and OpenAI GPT are single-model, single task. All results obtained from <https://gluebenchmark.com/leaderboard> and <https://blog.openai.com/language-unsupervised/>.

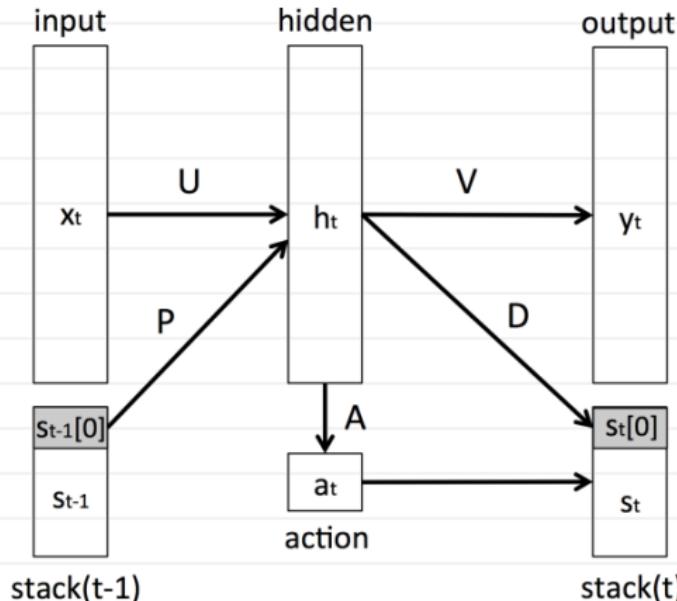
## Stack augmented RNNs

- 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

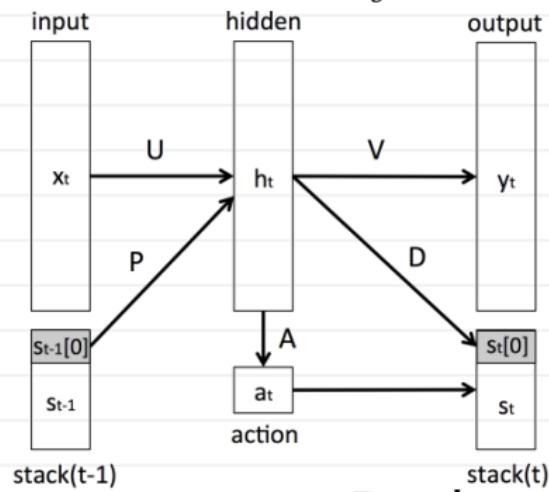
Conclusion: we need to decouple memory and operations (think RAM and CPU!)



[Joulin and Mikolov NIPS15]

# Stack augmented RNNs

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



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

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

$$\sigma(Dh_t)$$



[Joulin and Mikolov NIPS15]

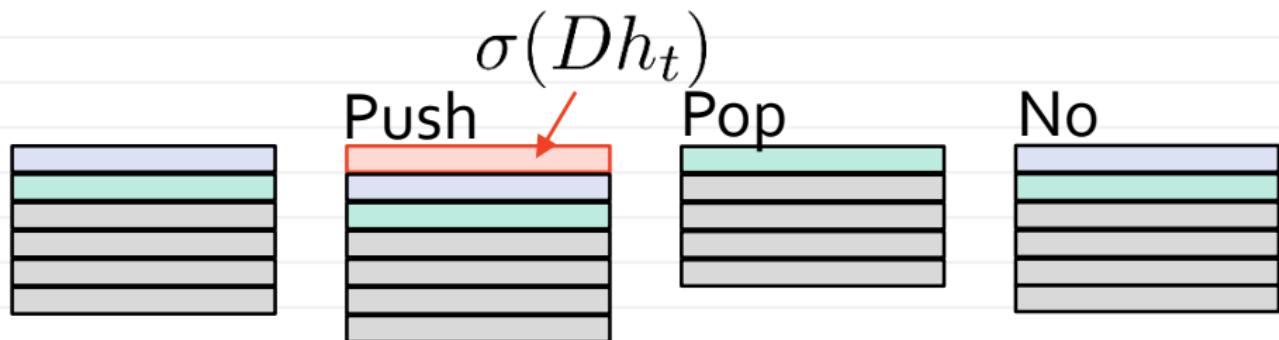
# Stack augmented RNNs

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

Actions: Push, Pop, No

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

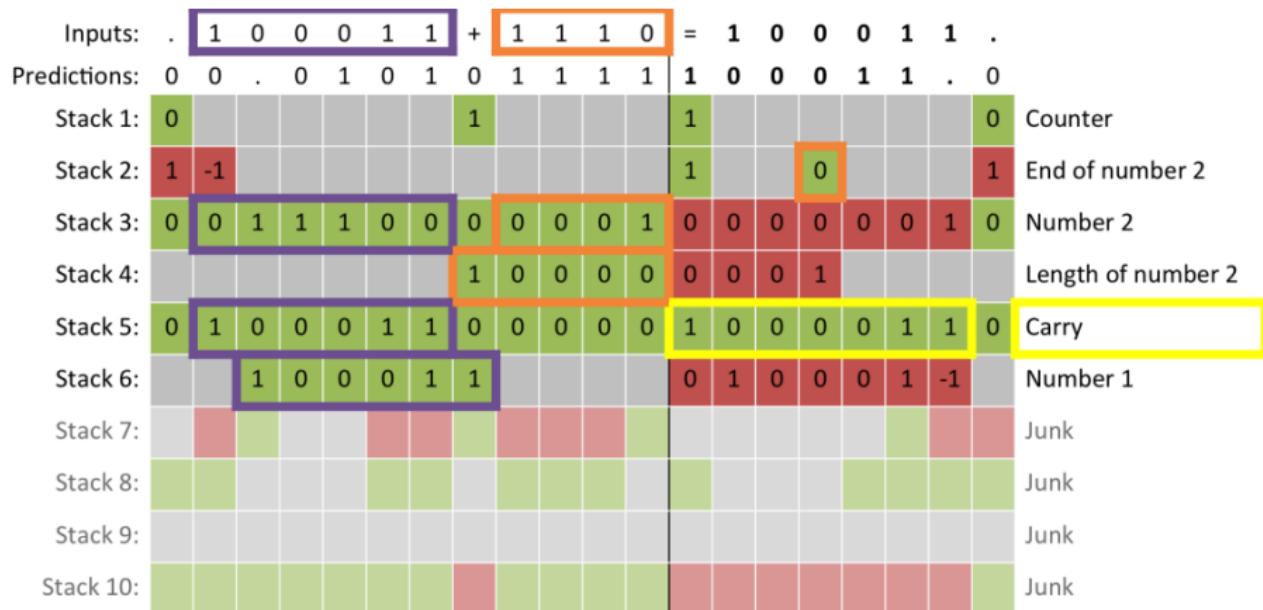
$$s_t^i = a_t^{\text{Push}} s_{t-1}^{i-1} + a_t^{\text{Pop}} s_{t-1}^{i+1} + a_t^{\text{No}} s_{t-1}^i$$



[Joulin and Mikolov NIPS15]

# Binary addition with stack-RNN

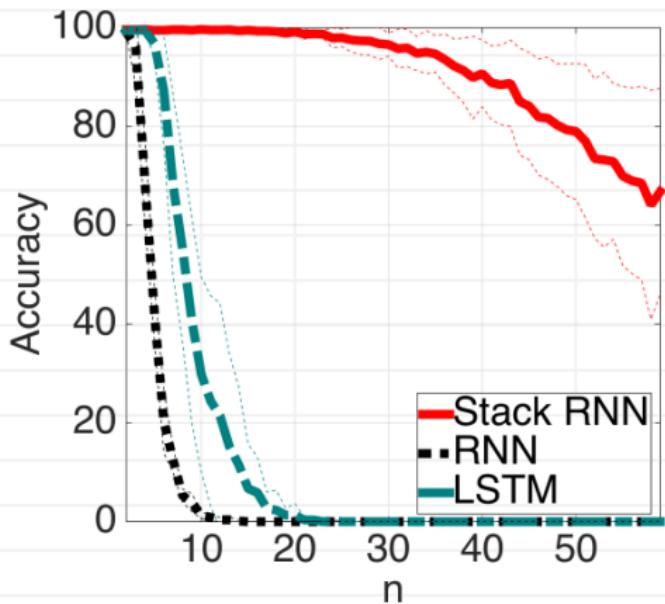
Goal: train a network that can add binary numbers.



PUSH POP

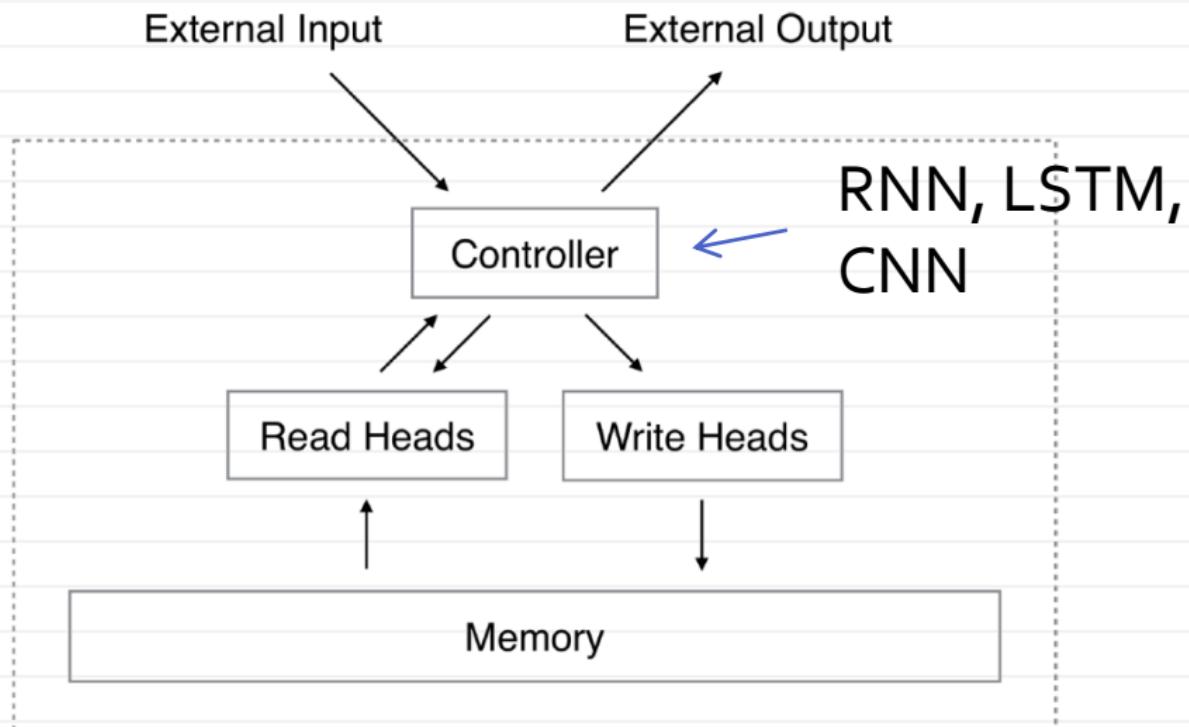
NB: the answer is reversed, i.e.  $101+11 = 0001$

# 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

Wang Ling, Chris Dyer, Alan W. Black, Isabel Trancoso, Ramon Fernandez, 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

# Bibliography

Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, Luke Zettlemoyer:

Deep Contextualized Word Representations. NAACL-HLT 2018: 2227-2237

Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova:

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. Arxiv 18