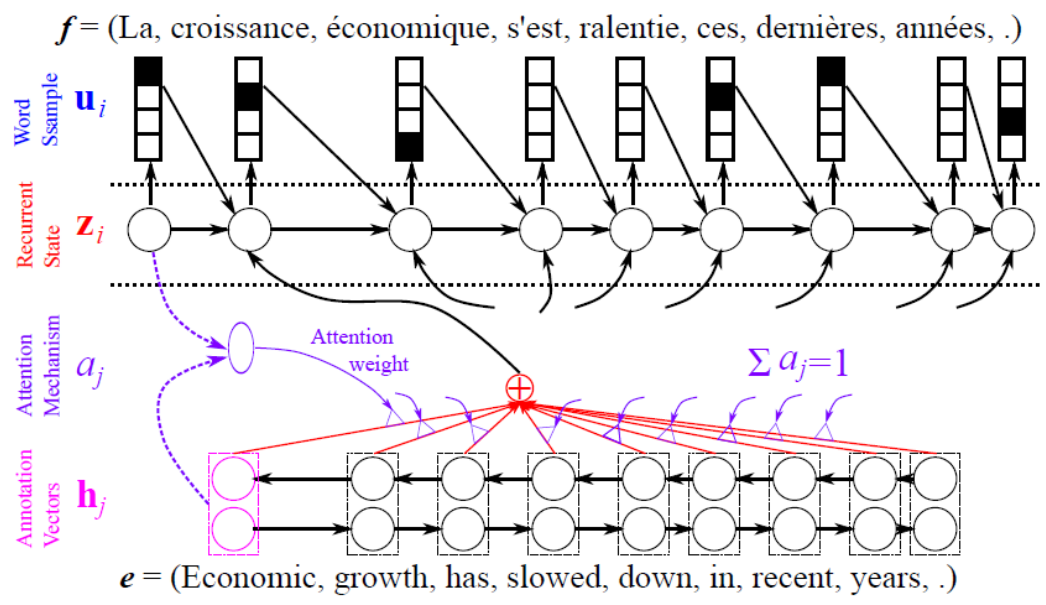
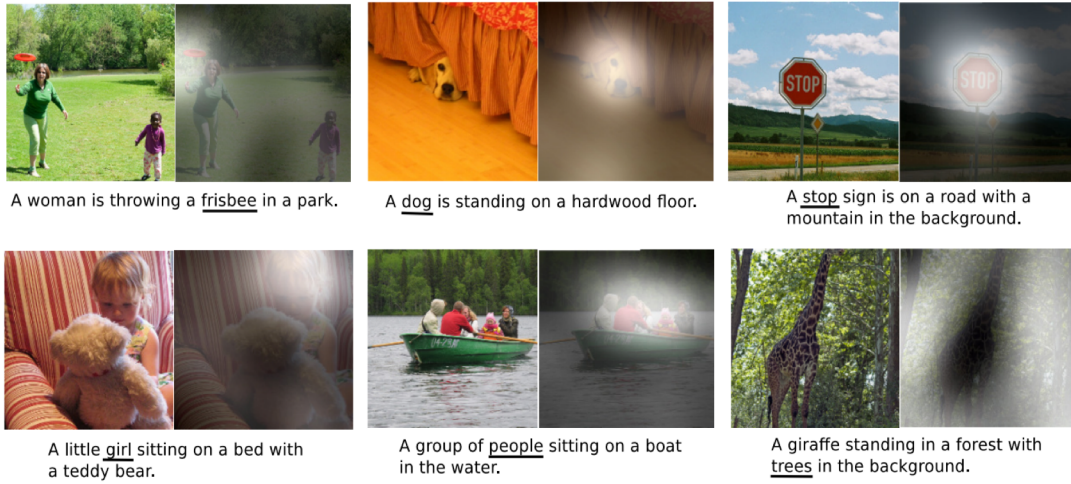


# Sequence-to-sequence models with attention



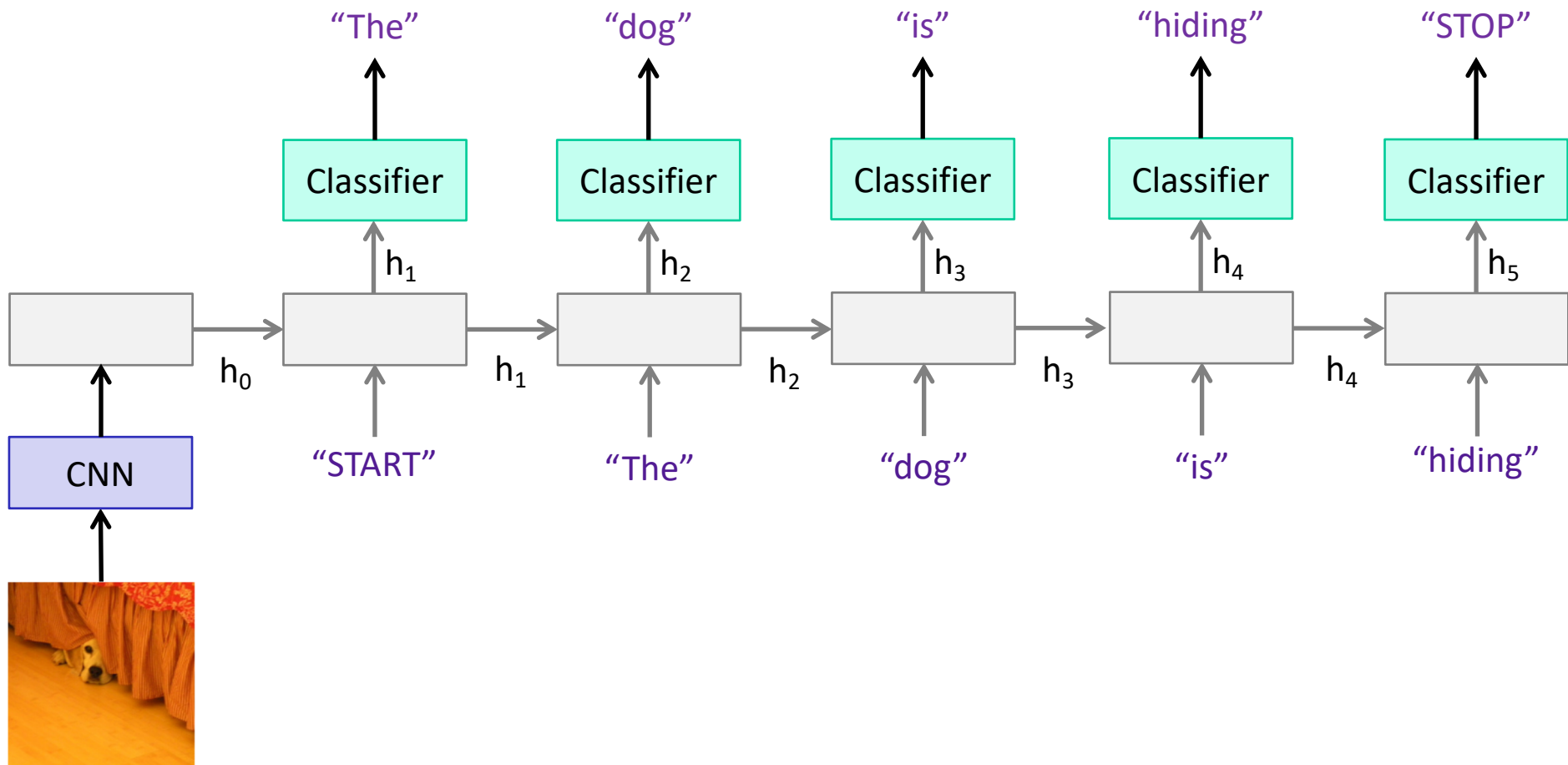
# Overview

---

- Image captioning with attention
- Neural machine translation with attention
  - Recurrent models with global and local attention
  - Google Neural Machine Translation
  - Convolutional sequence to sequence models
  - Attention without recurrence or convolutions

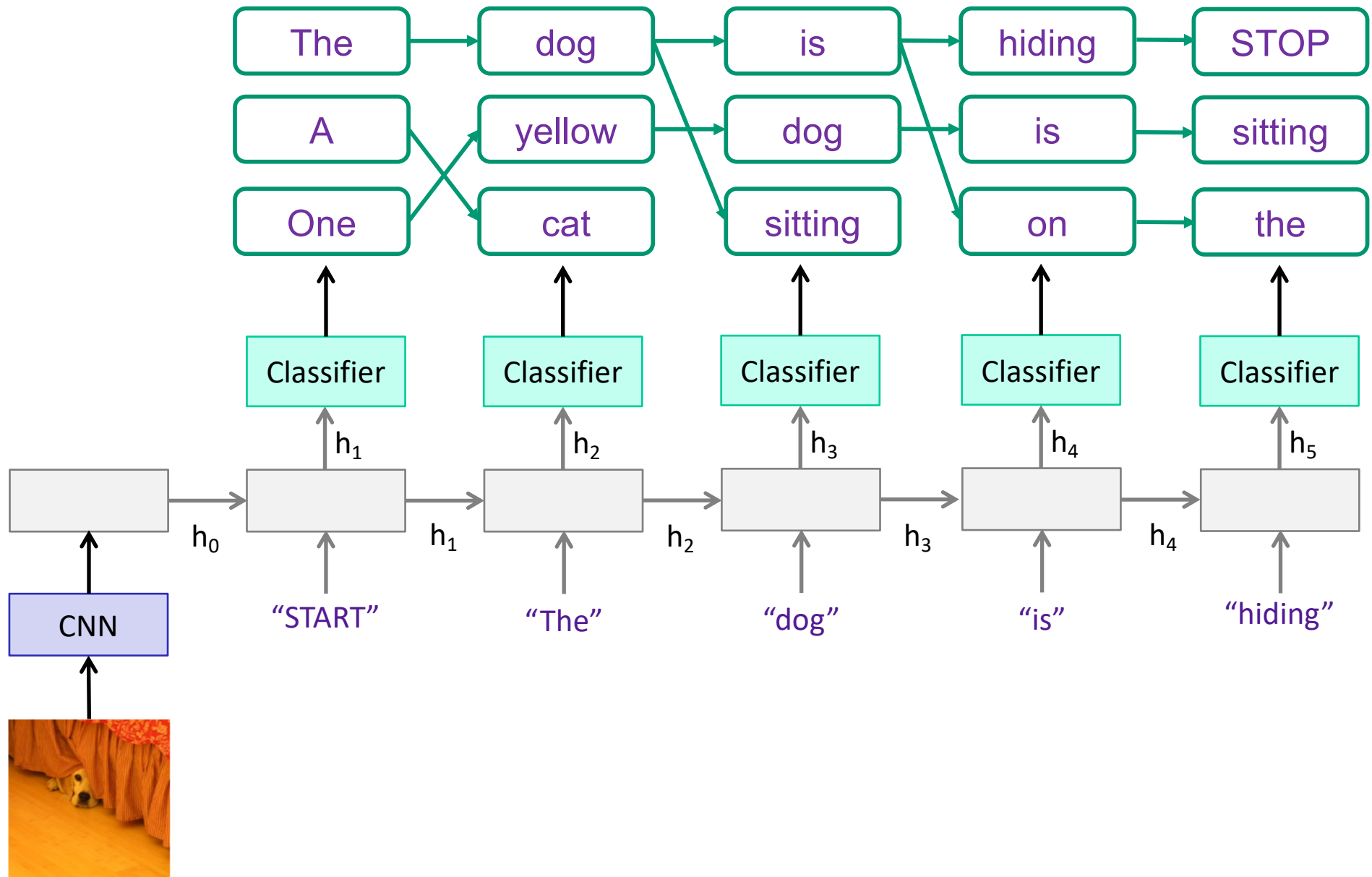
# Review: Image captioning

---



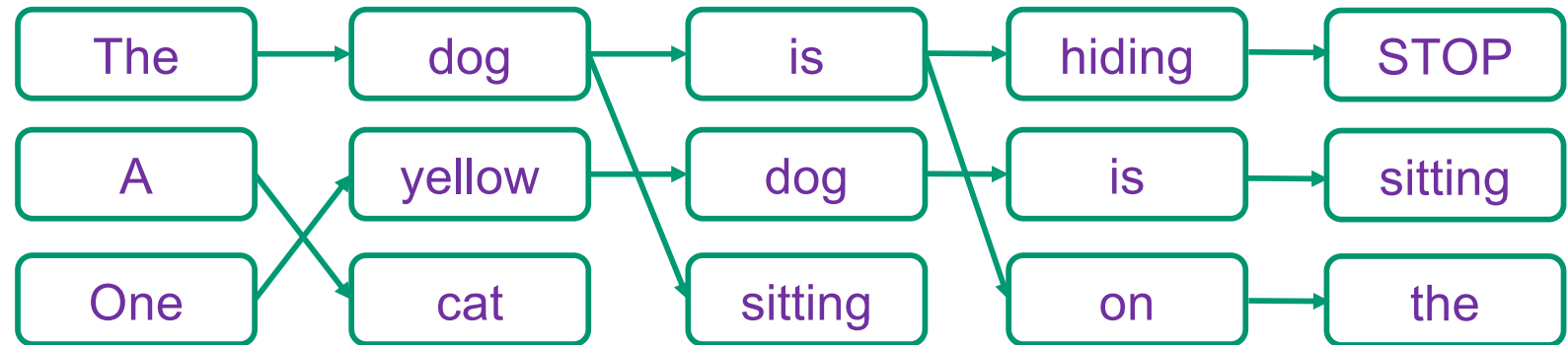
# Review: Image captioning

---



# Beam search

---



- Maintain  $k$  top-scoring candidate sentences (according to sum of per-word log-likelihoods)
  - At each step, generate all their successors and reduce to  $k$  (*beam width*)

# How to evaluate image captioning?

---

Reference sentences (written by human annotators):



- “A dog hides underneath a bed with its face peeking out of the bed skirt”
- “The small white dog is peeking out from under the bed”
- “A dog is peeking its head out from underneath a bed skirt”
- “A dog peeking out from under a bed”
- “A dog that is under a bed on the floor”

Generated sentence:

- “A dog is hiding”

# BLEU: Bilingual Evaluation Understudy

---

- **N-gram precision:** count the number of n-gram matches between candidate and reference translation, divide by total number of n-grams in candidate translation
  - Clip counts by the maximum number of times an n-gram occurs in any reference translation
  - Multiply by *brevity penalty* to penalize short translations
- Most commonly used measure despite well-known shortcomings



Overview

Challenges

Download

Evaluate

Leaderboard

Table-C5

Table-C40

2015 Captioning Challenge

Last update: June 8, 2015. Visit [CodaLab](#) for the latest results.

	CIDEr-D	Meteor	ROUGE-L	BLEU-1	BLEU-2	BLEU-3	BLEU-4	
m-RNN (Baidu/ UCLA) <sup>[16]</sup>	0.886	0.238	0.524	0.72	0.553	0.41	0.302	
m-RNN <sup>[15]</sup>	0.817	0.219	0.504	0.719	0.515	0.404	0.299	
MSR Captiva							0.308	
Google <sup>[4]</sup>	CIDEr-D	CIDEr: Consensus-based Image Description Evaluation						0.309
Berkeley LR	METEOR	Meteor Universal: Language Specific Translation Evaluation for Any Target Language						0.277
Nearest Neig	Rouge-L	ROUGE: A Package for Automatic Evaluation of Summaries						0.28
MSR <sup>[8]</sup>	BLEU	BLEU: a Method for Automatic Evaluation of Machine Translation						0.291
Montreal/Toronto <sup>[10]</sup>	0.85	0.243	0.513	0.689	0.515	0.372	0.268	
PicSOM <sup>[13]</sup>	0.833	0.231	0.505	0.683	0.51	0.377	0.281	
Tsinghua Bigeye <sup>[14]</sup>	0.673	0.207	0.49	0.671	0.494	0.35	0.241	
MLBL <sup>[7]</sup>	0.74	0.219	0.499	0.666	0.498	0.362	0.26	
Human <sup>[5]</sup>	0.854	0.252	0.484	0.663	0.469	0.321	0.217	

Metrics

CIDEr-D

CIDEr: Consensus-based Image Description Evaluation

METEOR

Meteor Universal: Language Specific Translation Evaluation for Any Target Language

Rouge-L

ROUGE: A Package for Automatic Evaluation of Summaries

BLEU

BLEU: a Method for Automatic Evaluation of Machine Translation





Overview

Challenges

Download

Evaluate

Leaderboard

Table-C5

Table-C40

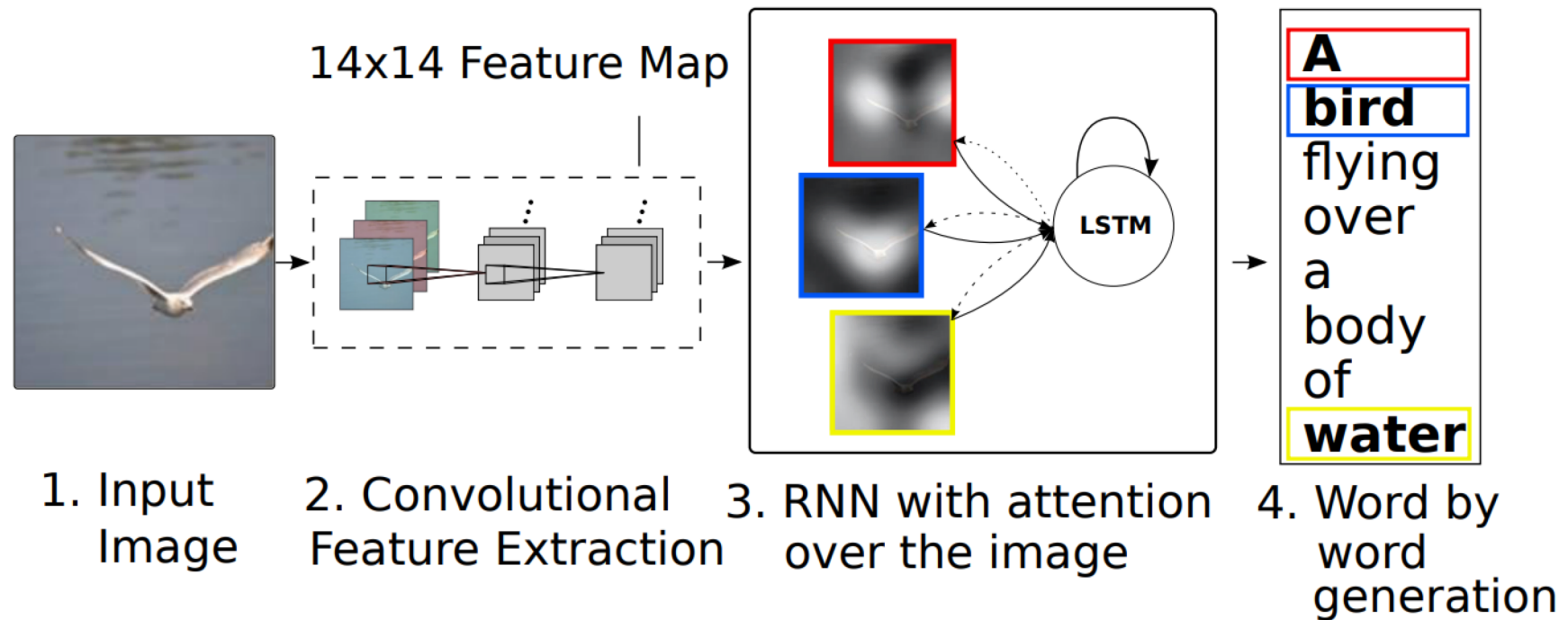
2015 Captioning Challenge

Last update: June 8, 2015. Visit [CodaLab](#) for the latest results.

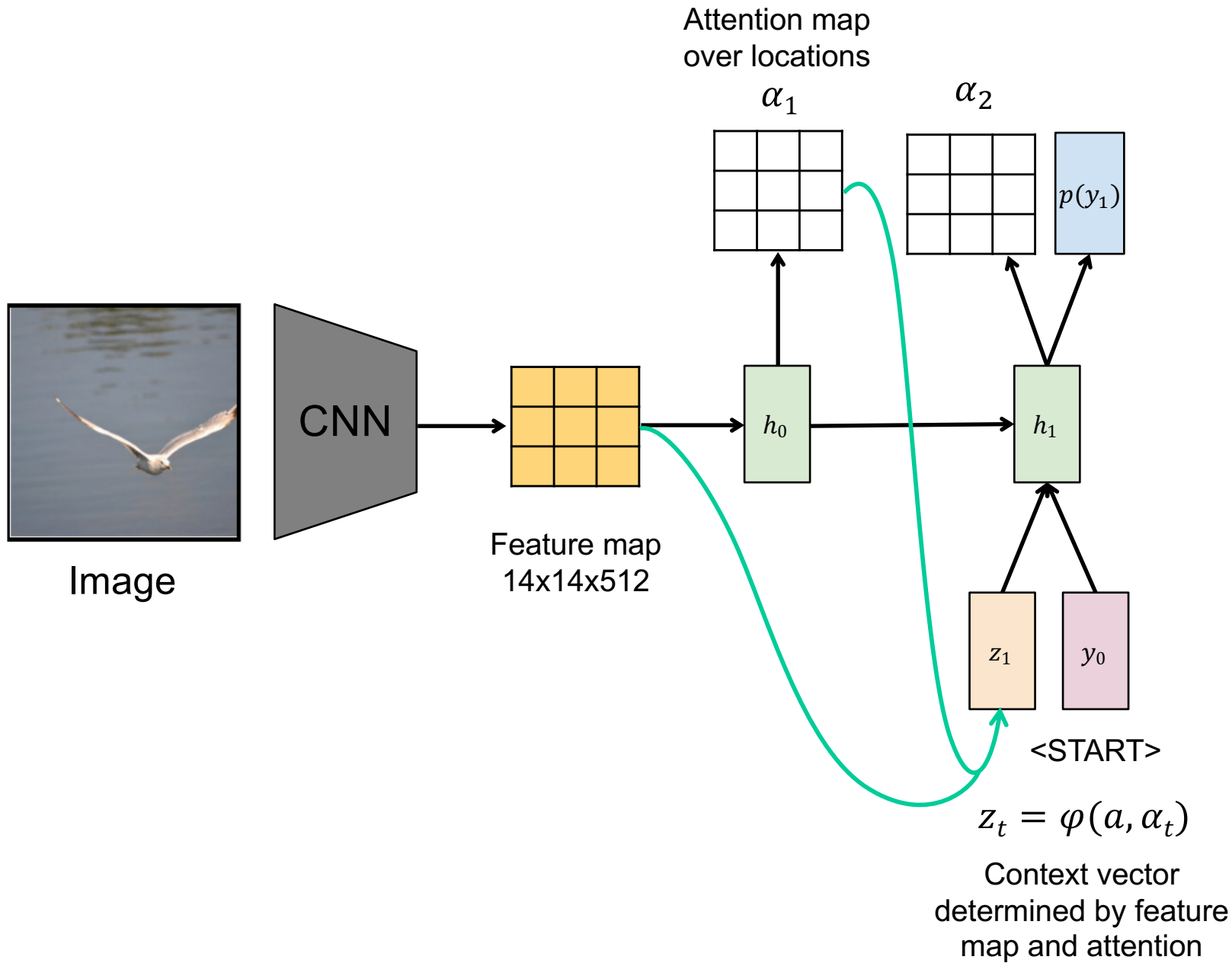
	M1	M2	M3	M4	M5
Human <sup>[5]</sup>	0.638	0.675	4.836	3.428	0.352
Google <sup>[4]</sup>	0.272	0.247	4.107	2.740	0.222
MSR <sup>[8]</sup>	M1	Percentage of captions that are evaluated as better or equal to human caption.			
Montreal	M2	Percentage of captions that pass the Turing Test.			
MSR Ca	M3	Average correctness of the captions on a scale 1-5 (incorrect - correct).			
Berkeley	M4	Average amount of detail of the captions on a scale 1-5 (lack of details - very detailed).			
m-RNN <sup>[1]</sup>	M5	Percentage of captions that are similar to human description.			
Nearest Neighbor <sup>[11]</sup>	0.216	0.255	3.801	2.716	0.196
PicSOM <sup>[13]</sup>	0.202	0.250	3.965	2.552	0.182
Brno University <sup>[3]</sup>	0.194	0.213	3.079	3.482	0.154
m-RNN (Baidu/ UCLA) <sup>[16]</sup>	0.190	0.241	3.831	2.548	0.195
MIL <sup>[6]</sup>	0.168	0.197	3.349	2.915	0.159
MLBL <sup>[7]</sup>	0.167	0.196	3.659	2.420	0.156

# Captioning with attention

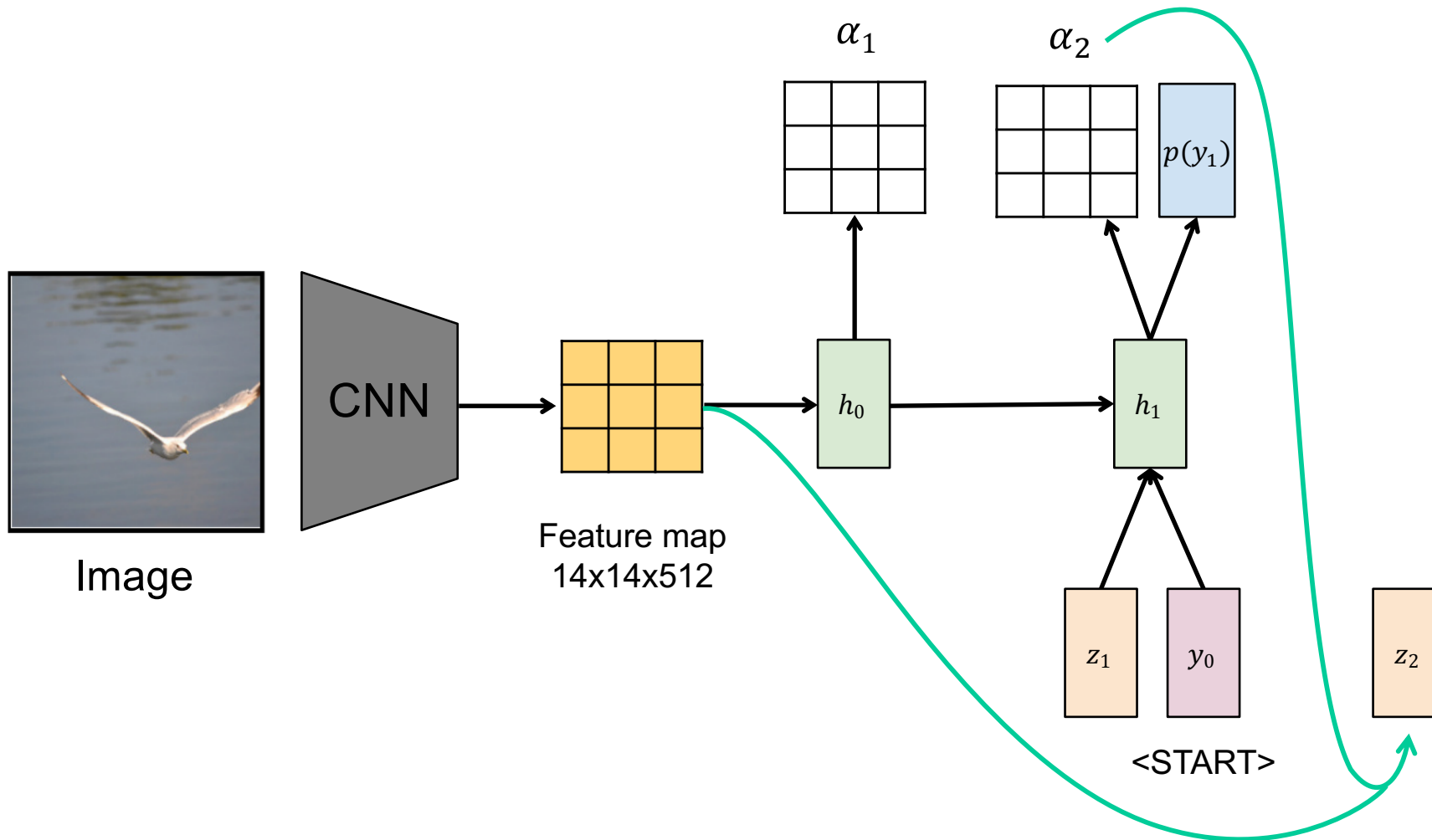
---



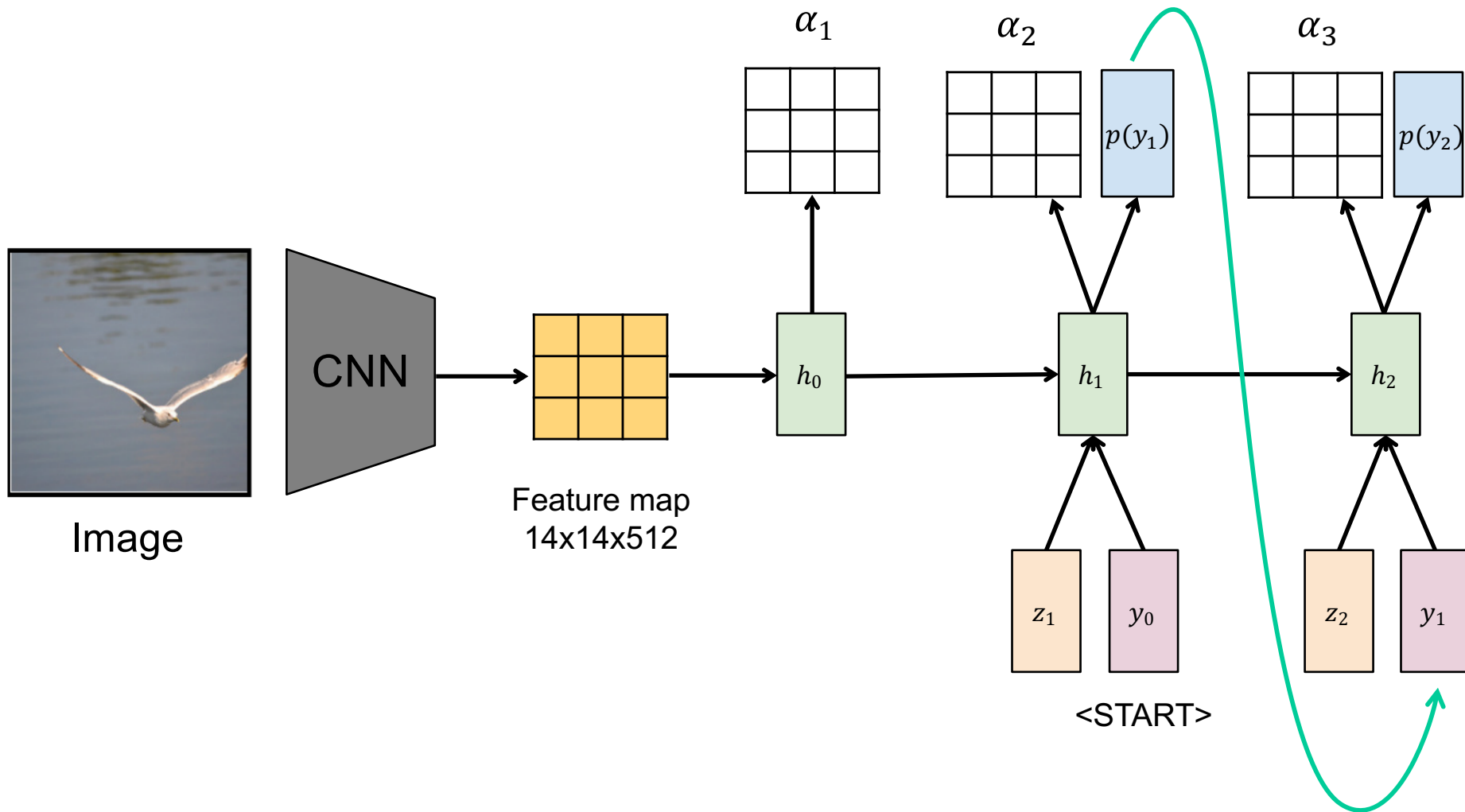
# Captioning with attention



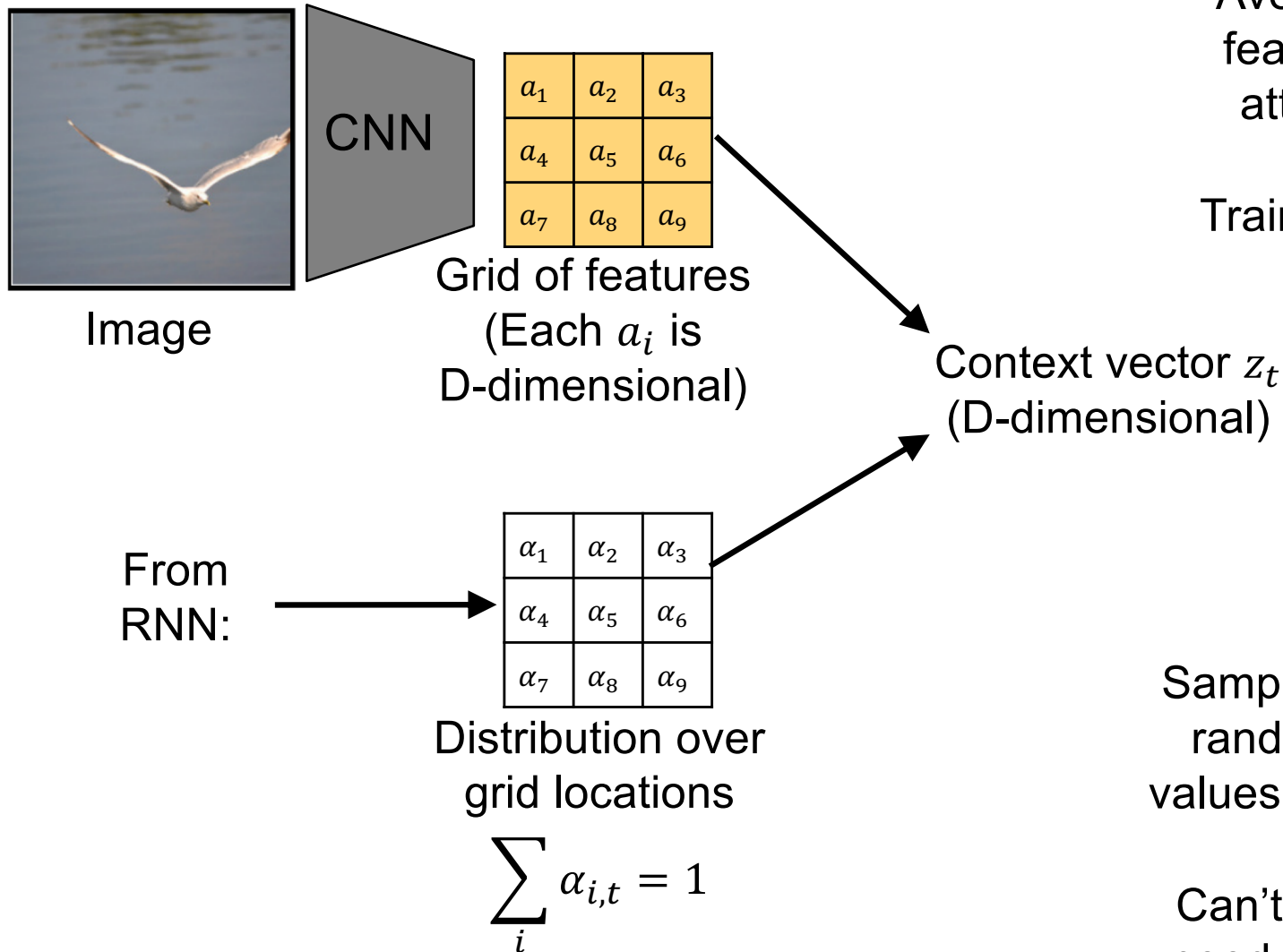
# Captioning with attention



# Captioning with attention



# “Soft” and “hard” attention



## Soft attention:

Average over locations of feature map weighted by attention:  $z_t = \sum_i \alpha_{i,t} a_i$

Train with gradient descent

## Hard attention:

Sample ONE location:  $z_t$  is a random variable taking on values  $a_i$  with probabilities  $\alpha_{i,t}$

Can't use gradient descent; need reinforcement learning

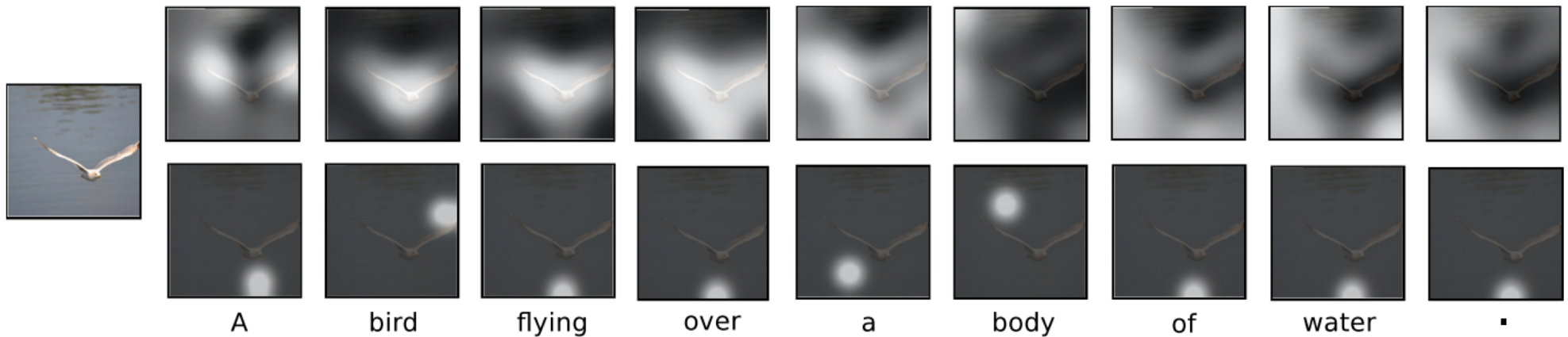
# “Soft” and “hard” attention

---

## Soft attention:

Average over locations of feature map weighted by attention:

$$z_t = \sum_i \alpha_{i,t} a_i$$



## Hard attention:

Sample ONE location:  $z_t$  is a random variable taking on values  $a_i$  with probabilities  $\alpha_{i,t}$

# Results

---

Dataset	Model	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR
Flickr8k	Google NIC	63	41	27	-	-
	Soft-Attention	<b>67</b>	44.8	29.9	19.5	18.93
	Hard-Attention	<b>67</b>	<b>45.7</b>	<b>31.4</b>	<b>21.3</b>	<b>20.30</b>
Flickr30k	Google NIC	66.3	42.3	27.7	18.3	-
	Soft-Attention	66.7	43.4	28.8	19.1	<b>18.49</b>
	Hard-Attention	<b>66.9</b>	<b>43.9</b>	<b>29.6</b>	<b>19.9</b>	18.46
COCO	Google NIC	66.6	46.1	32.9	24.6	-
	Soft-Attention	70.7	49.2	34.4	24.3	<b>23.90</b>
	Hard-Attention	<b>71.8</b>	<b>50.4</b>	<b>35.7</b>	<b>25.0</b>	23.04



# Example Results

---

- Good captions



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.

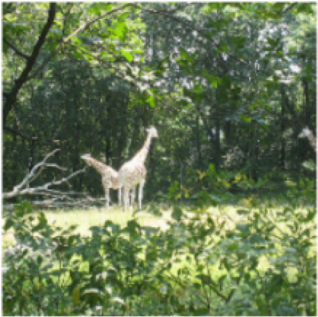


A giraffe standing in a forest with trees in the background.

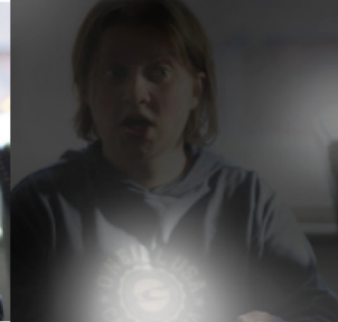
# Example Results

---

- Mistakes



A large white bird standing in a forest.



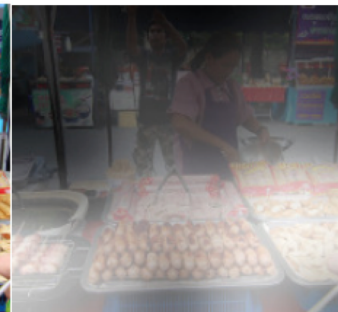
A woman holding a clock in her hand.



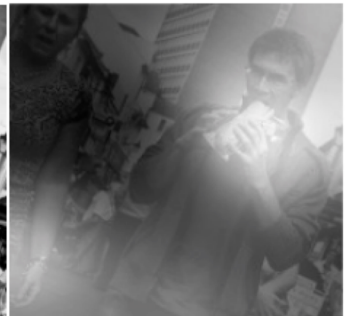
A man wearing a hat and a hat on a skateboard.



A person is standing on a beach with a surfboard.



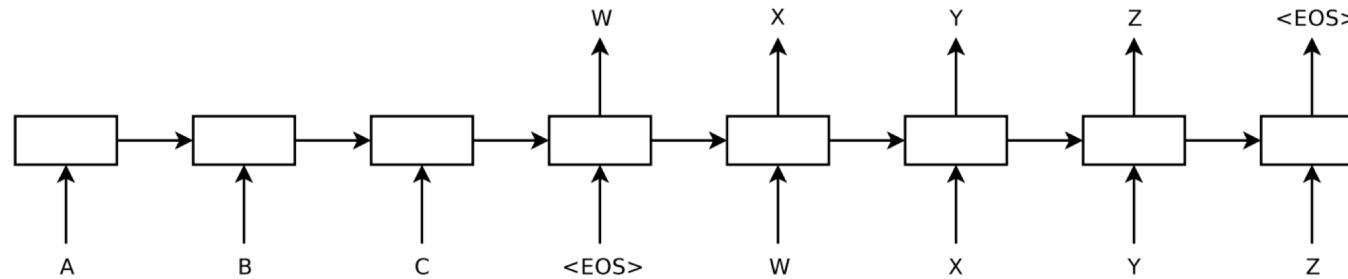
A woman is sitting at a table with a large pizza.



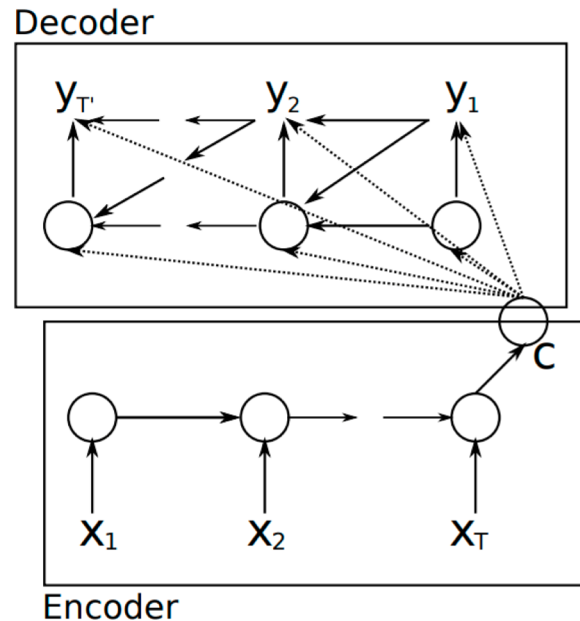
A man is talking on his cell phone while another man watches.

# Machine translation: Vanilla Seq2Seq

---



I. Sutskever, O. Vinyals, Q. Le, [Sequence to Sequence Learning with Neural Networks](#), NIPS 2014

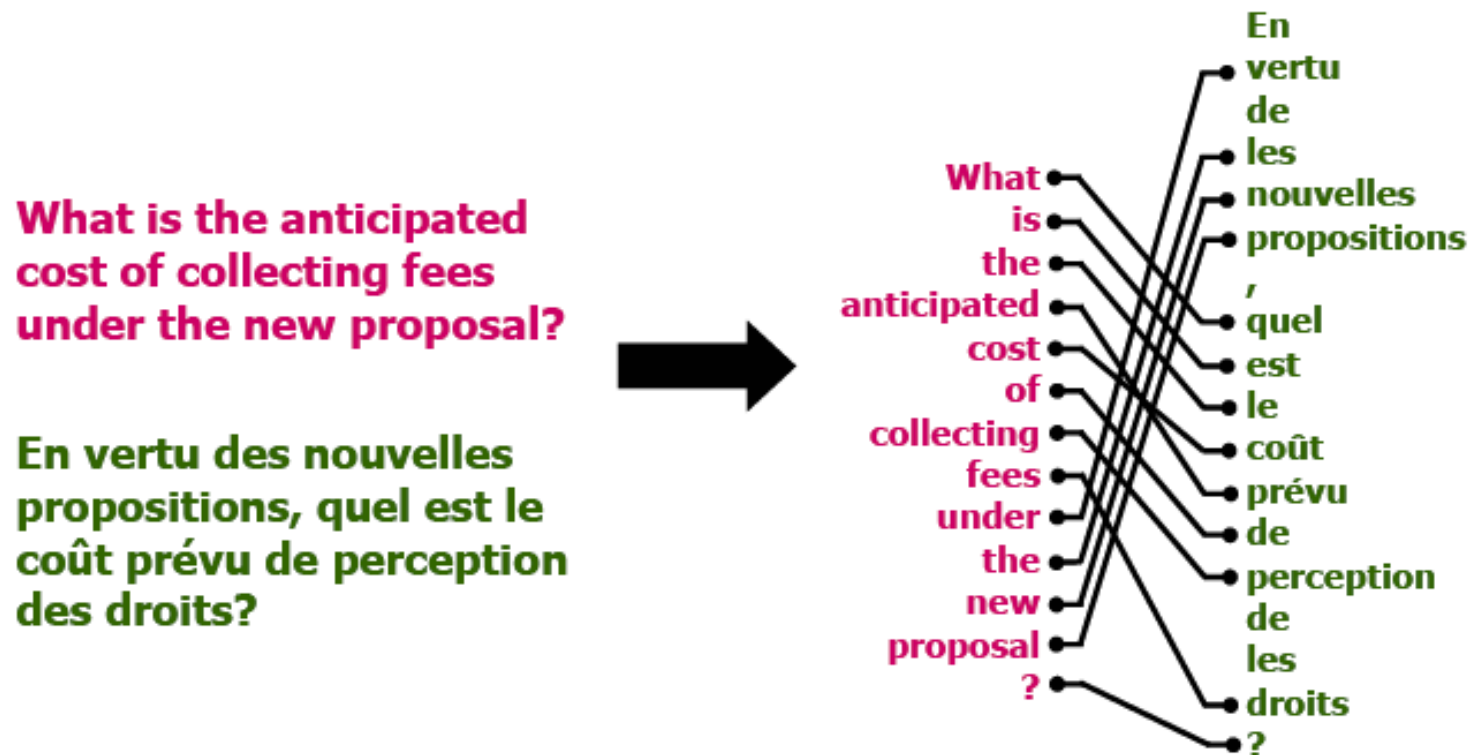


K. Cho, B. Merriënboer, C. Gulcehre, F. Bougares, H. Schwenk, and Y. Bengio, [Learning phrase representations using RNN encoder-decoder for statistical machine translation](#), ACL 2014

# Machine translation with attention

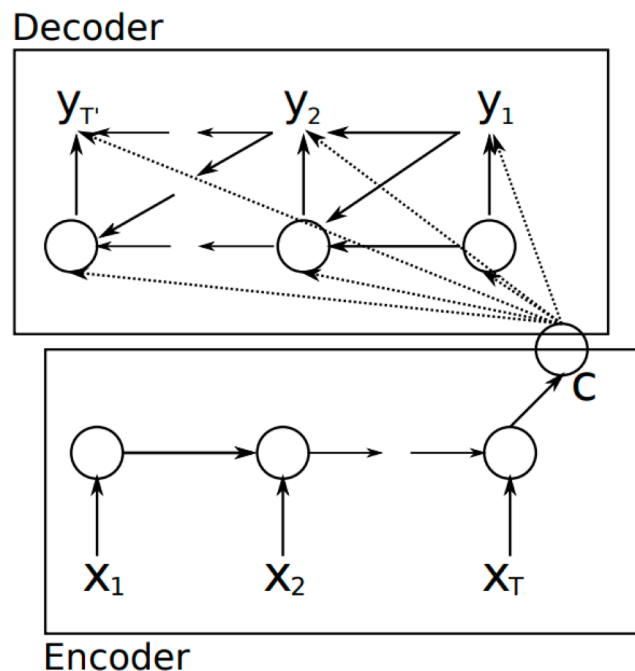
---

- Key idea: translation requires *alignment*



# Machine translation with attention

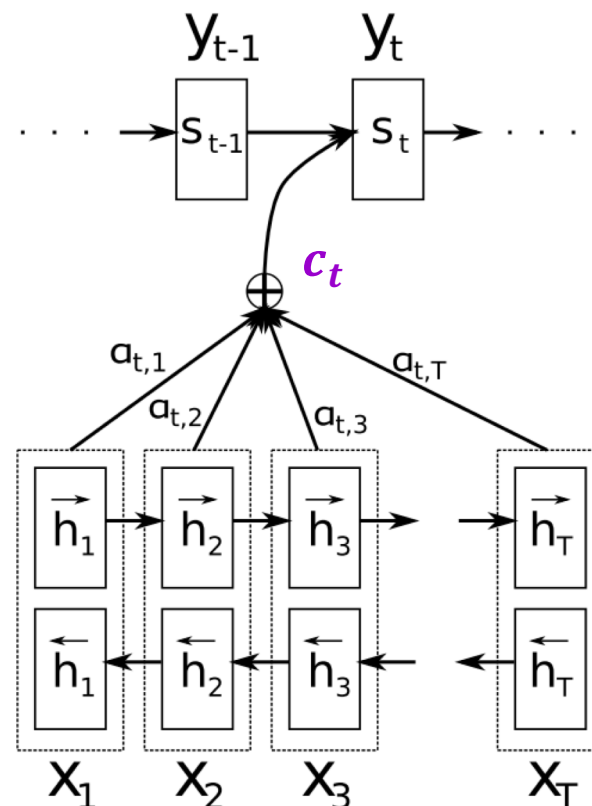
## Standard encoder-decoder



A fixed context vector  $c = h_T$  is used for decoding each word.

$$h_t = f(x_t, h_{t-1})$$

## Attention-based model

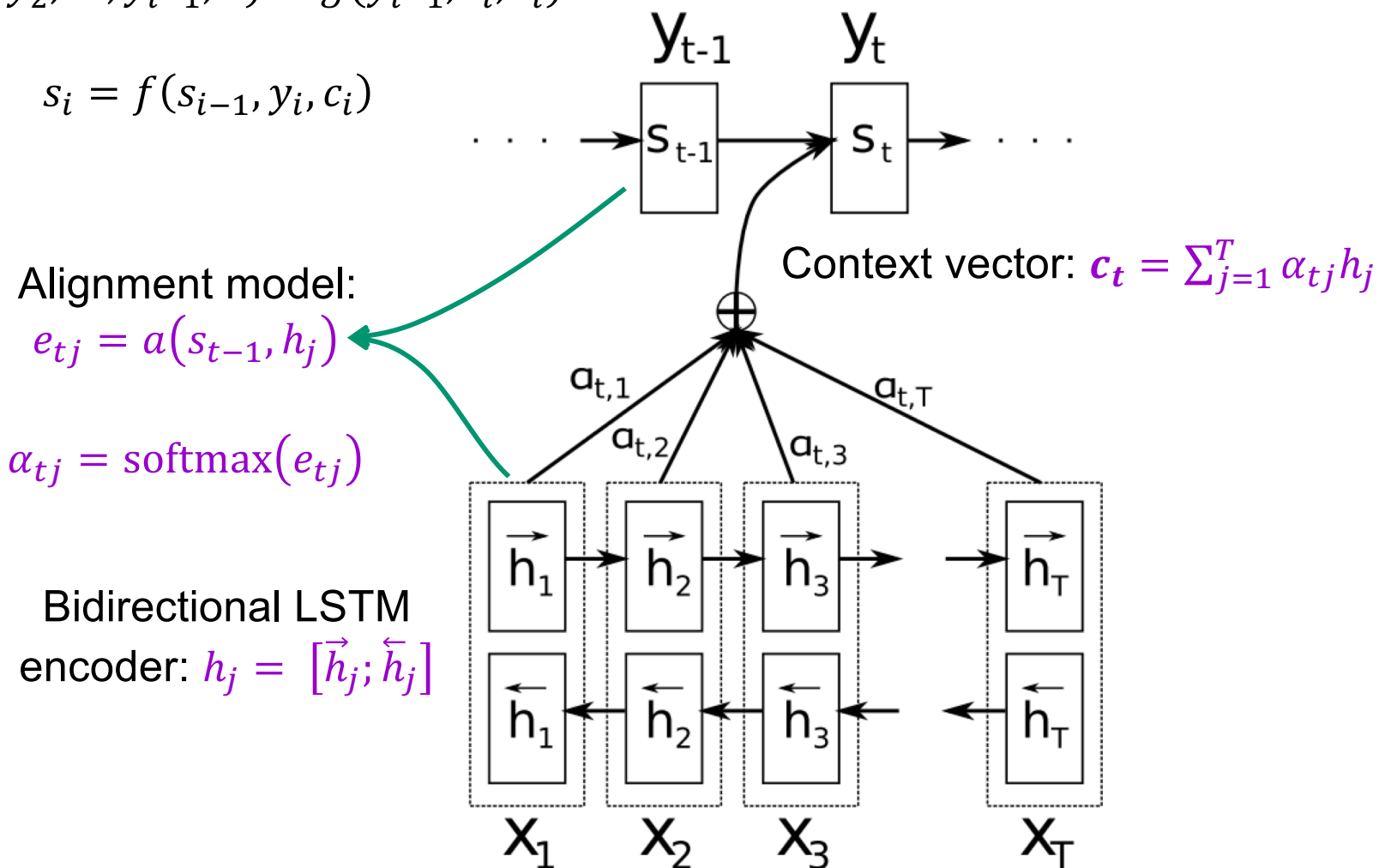


Context vector  $c_t$  pays attention to different phrases in the source when generating each word

# Global attentional model

$$p(y_i | y_1, y_2, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i)$$

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



# Attention model

$$\alpha_{tj}$$

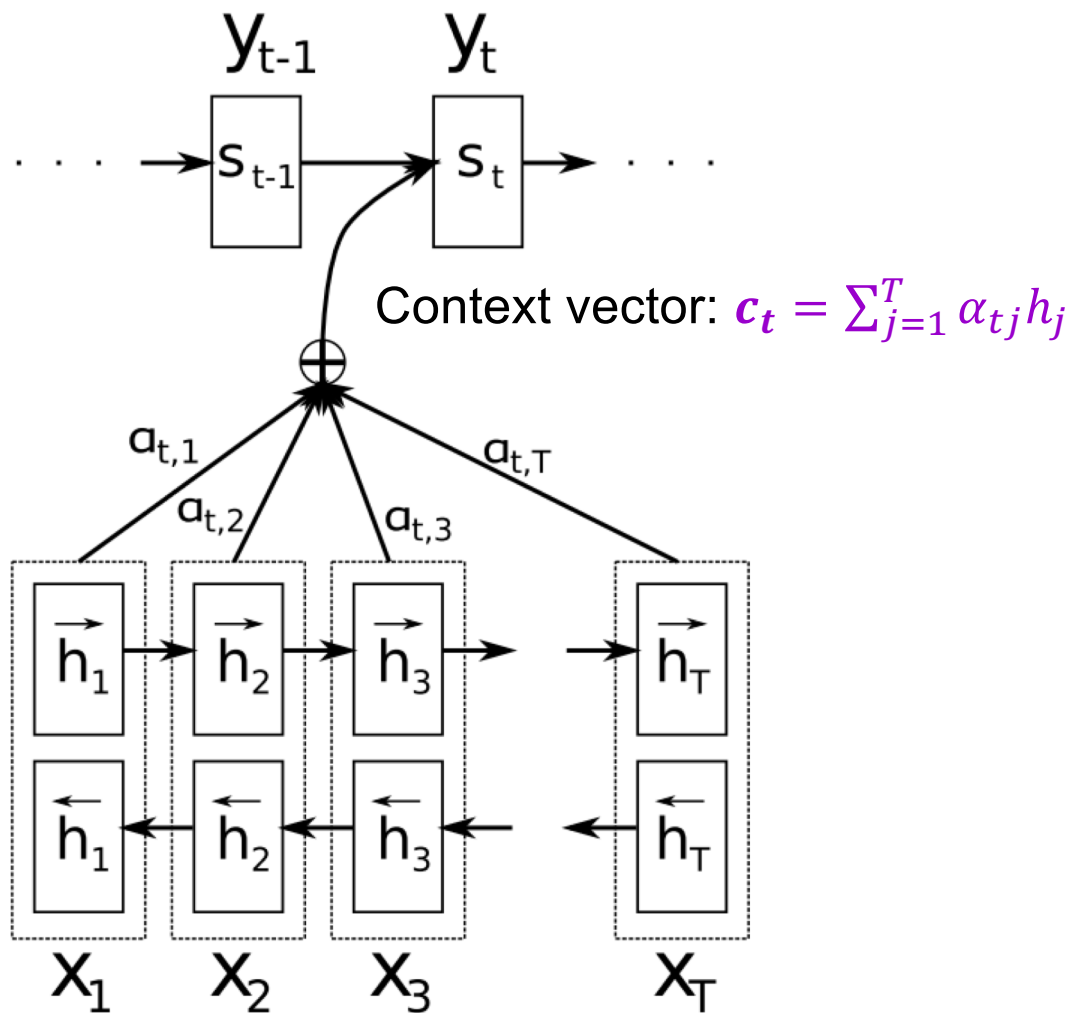
How much attention should output  $y_t$  pay to input  $x_j$

$$\alpha_{tj} = \text{softmax}(e_{tj})$$

How to compute  $e_{tj}$ ?

$$e_{tj} = a(s_{t-1}, h_j)$$

Train this using small NN  
This model is effectively trying to align encoder hidden states with decoder hidden states

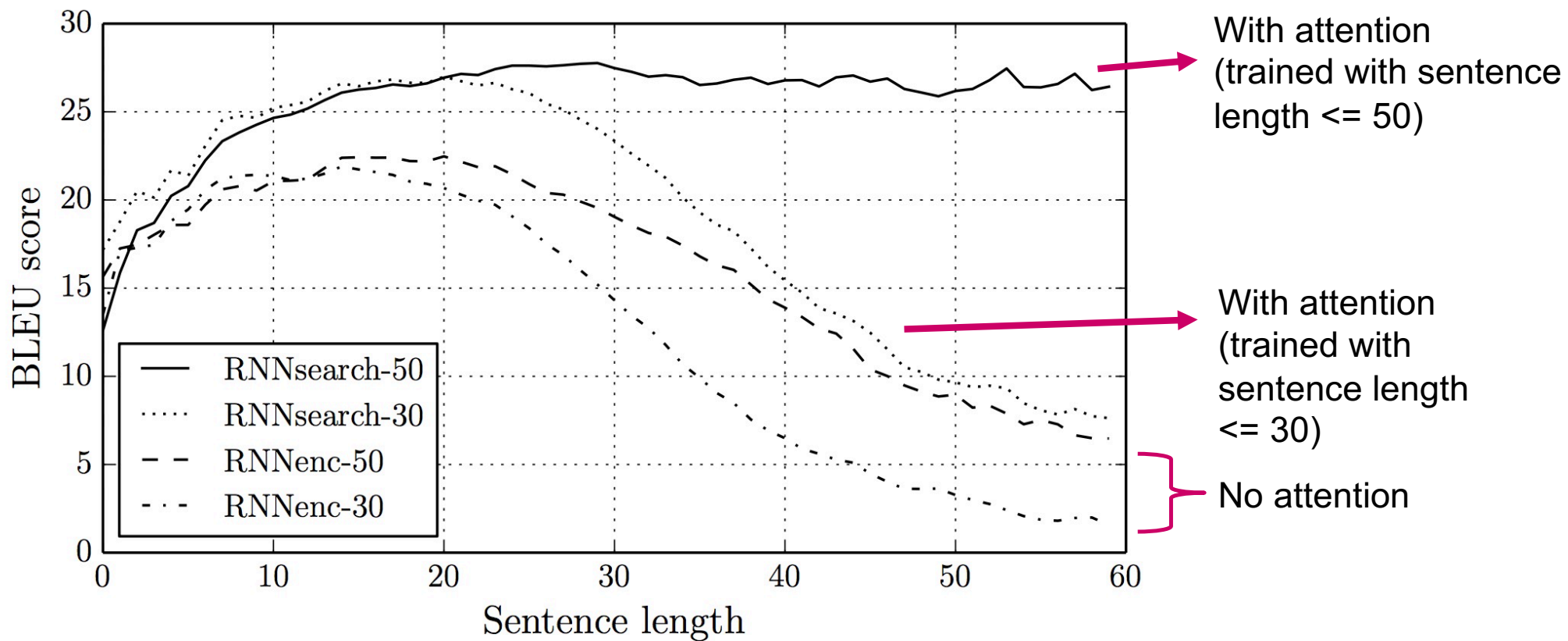






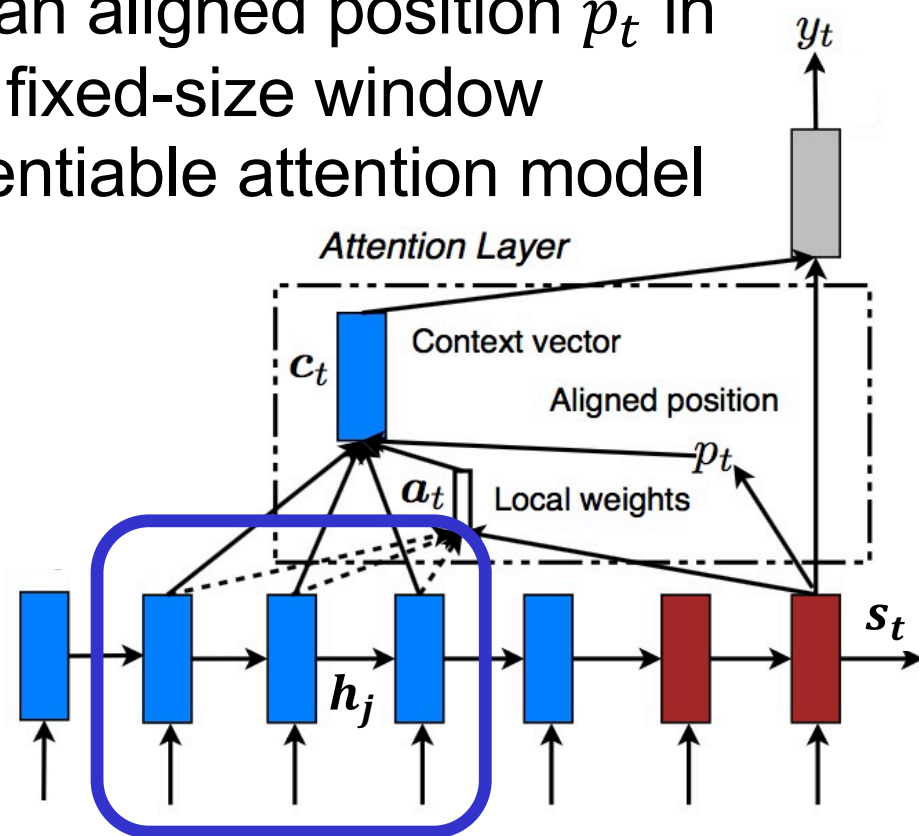
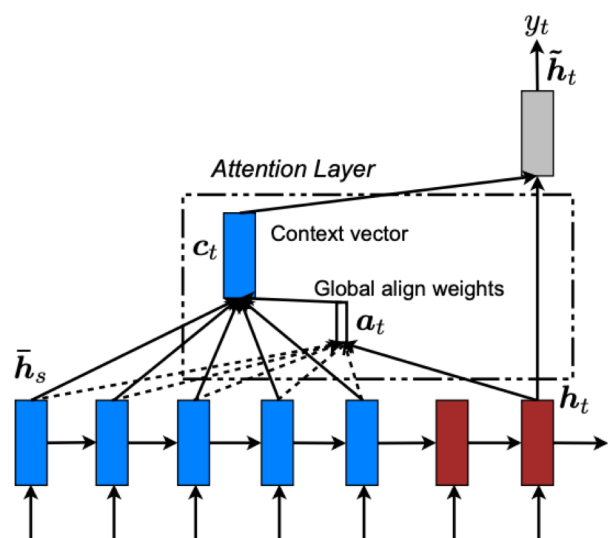
# Quantitative evaluation

---

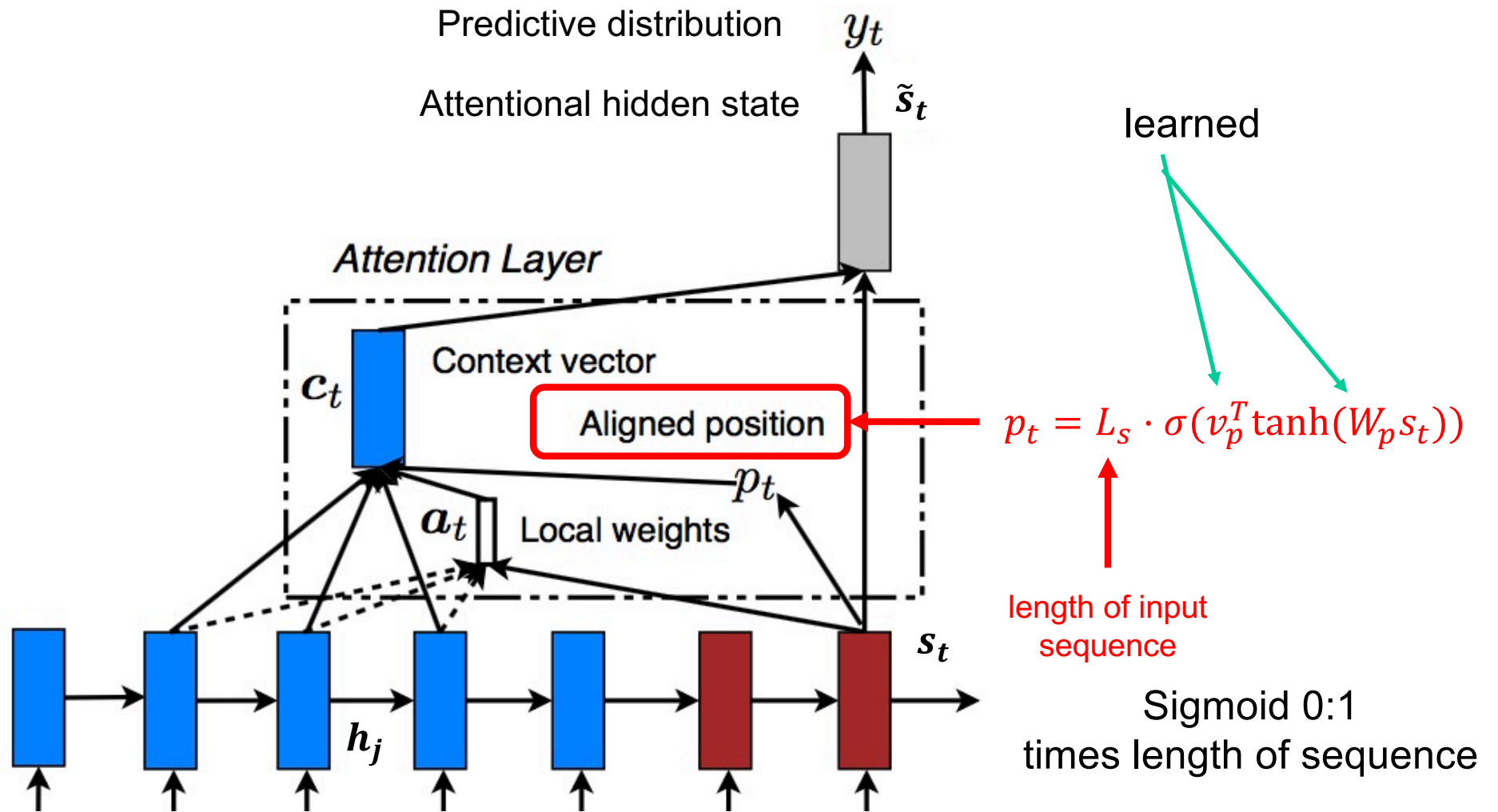


# Local attention

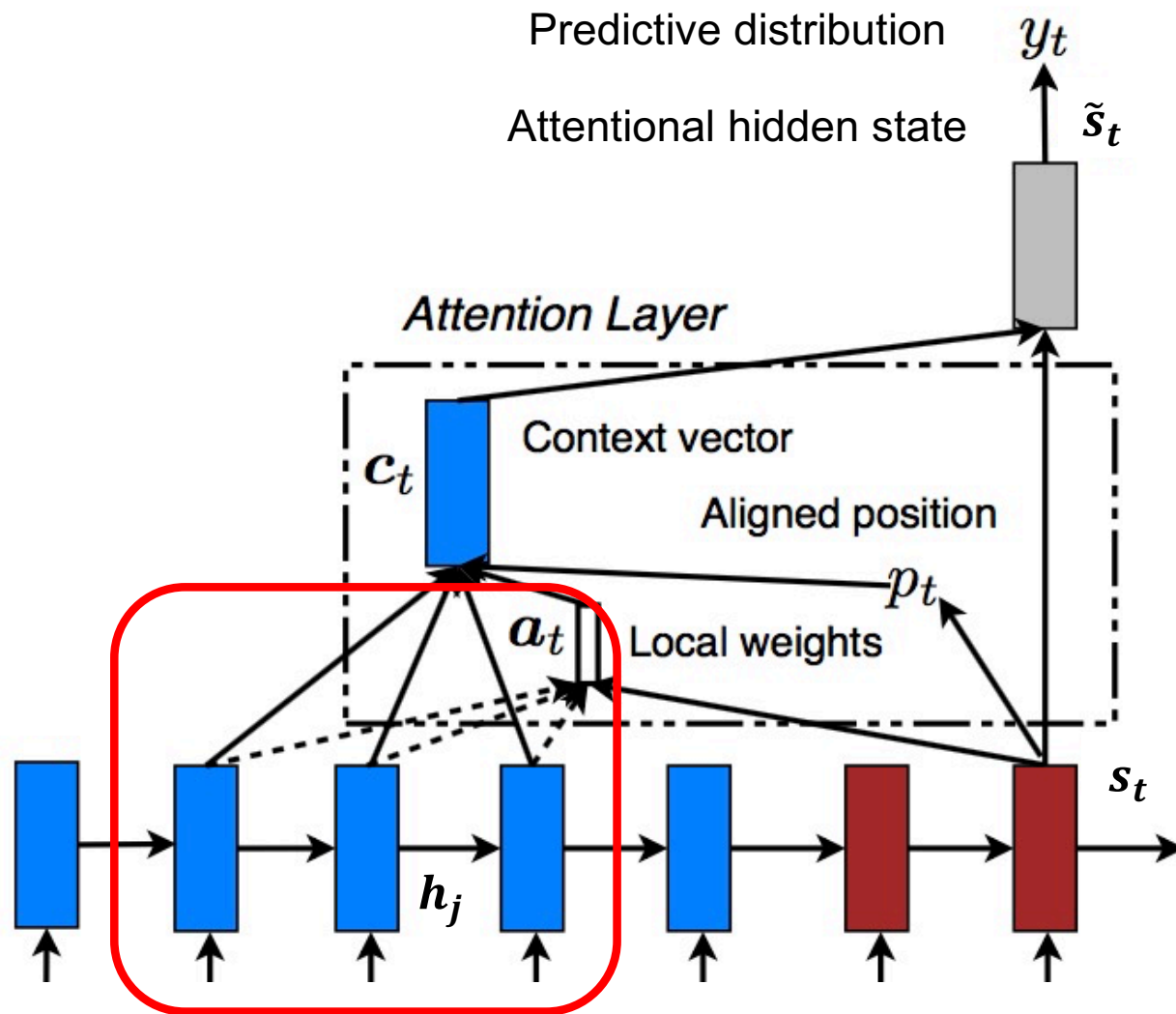
- Key idea: design mechanism similar to “hard attention” but differentiable
- Global attention - Attend over whole input
- For each target word, predict an aligned position  $p_t$  in the source; form context from fixed-size window around  $p_t$  - here simple differentiable attention model



# Local attention



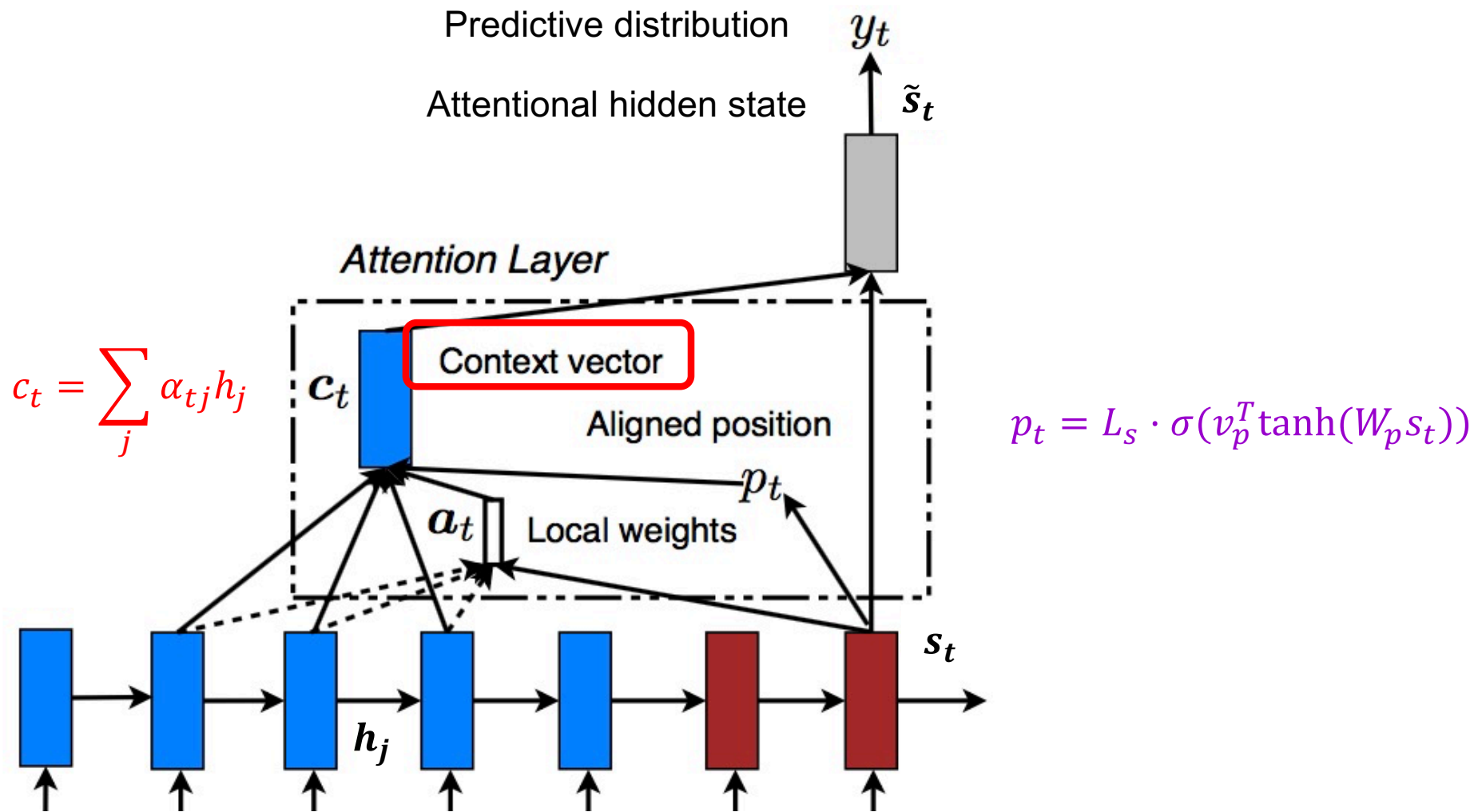
# Local attention



$$p_t = L_s \cdot \sigma(v_p^T \tanh(W_p s_t))$$

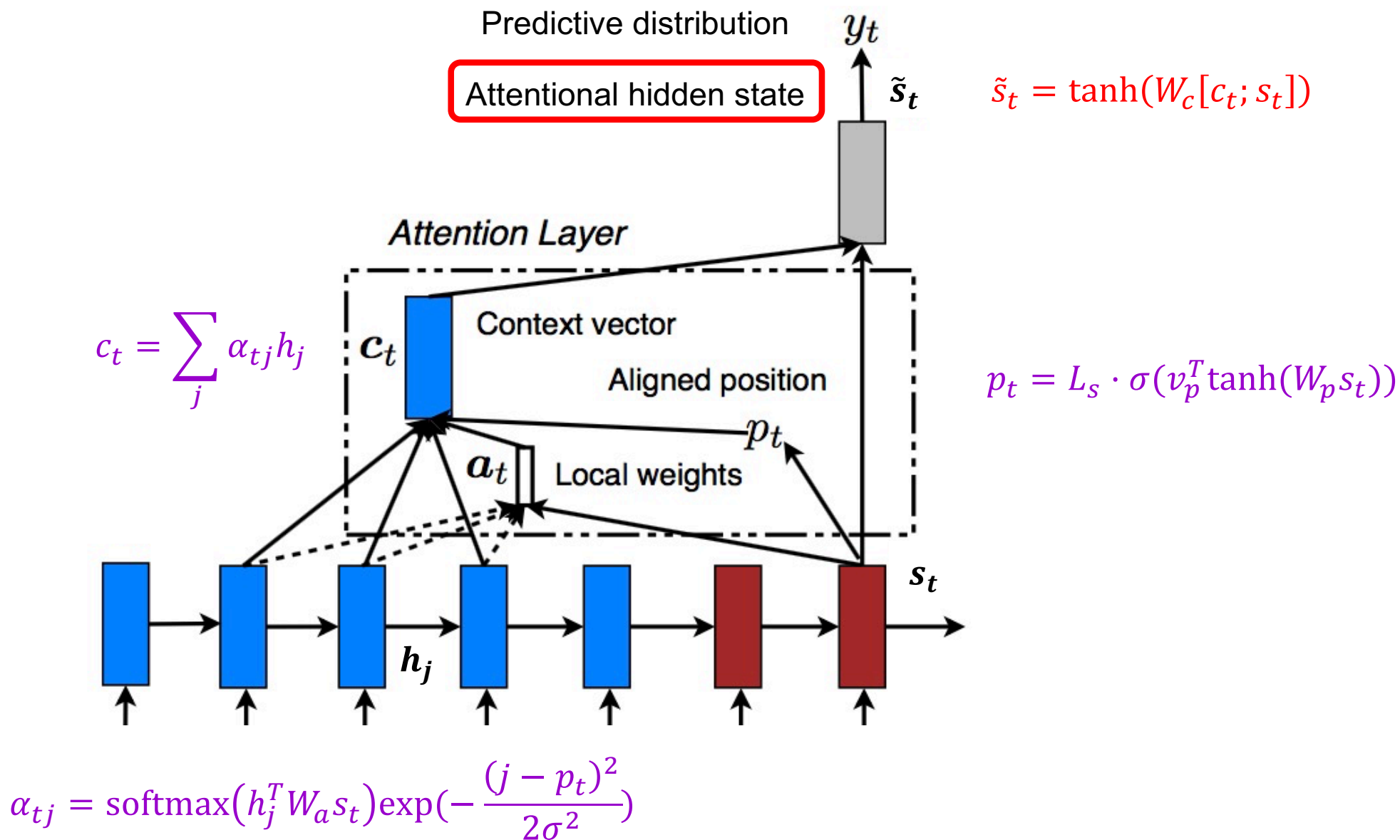
$$\alpha_{tj} = \text{softmax}(h_j^T W_a s_t) \exp\left(-\frac{(j - p_t)^2}{2\sigma^2}\right)$$

# Local attention

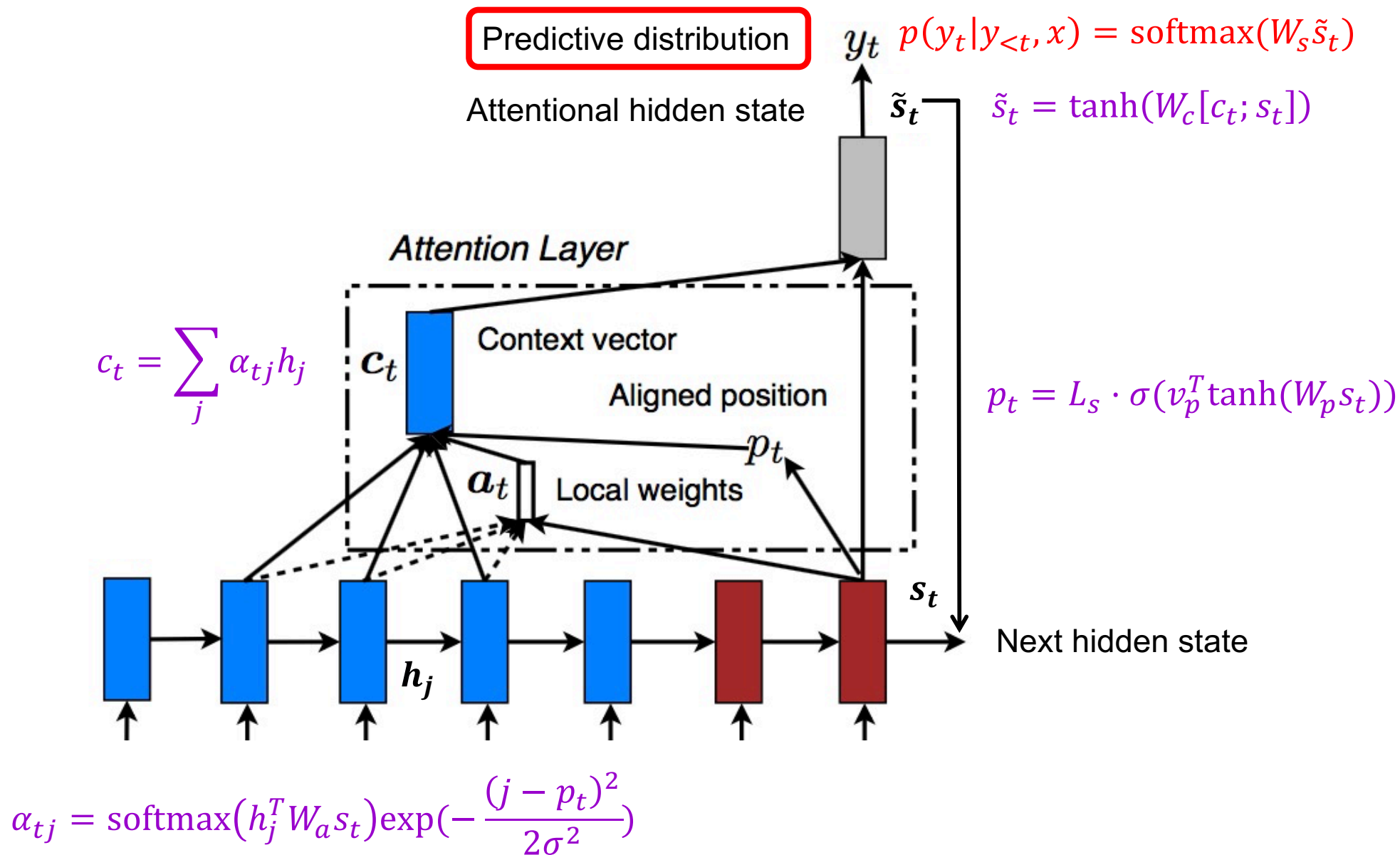


$$\alpha_{tj} = \text{softmax}(h_j^T W_a s_t) \exp\left(-\frac{(j - p_t)^2}{2\sigma^2}\right)$$

# Local attention

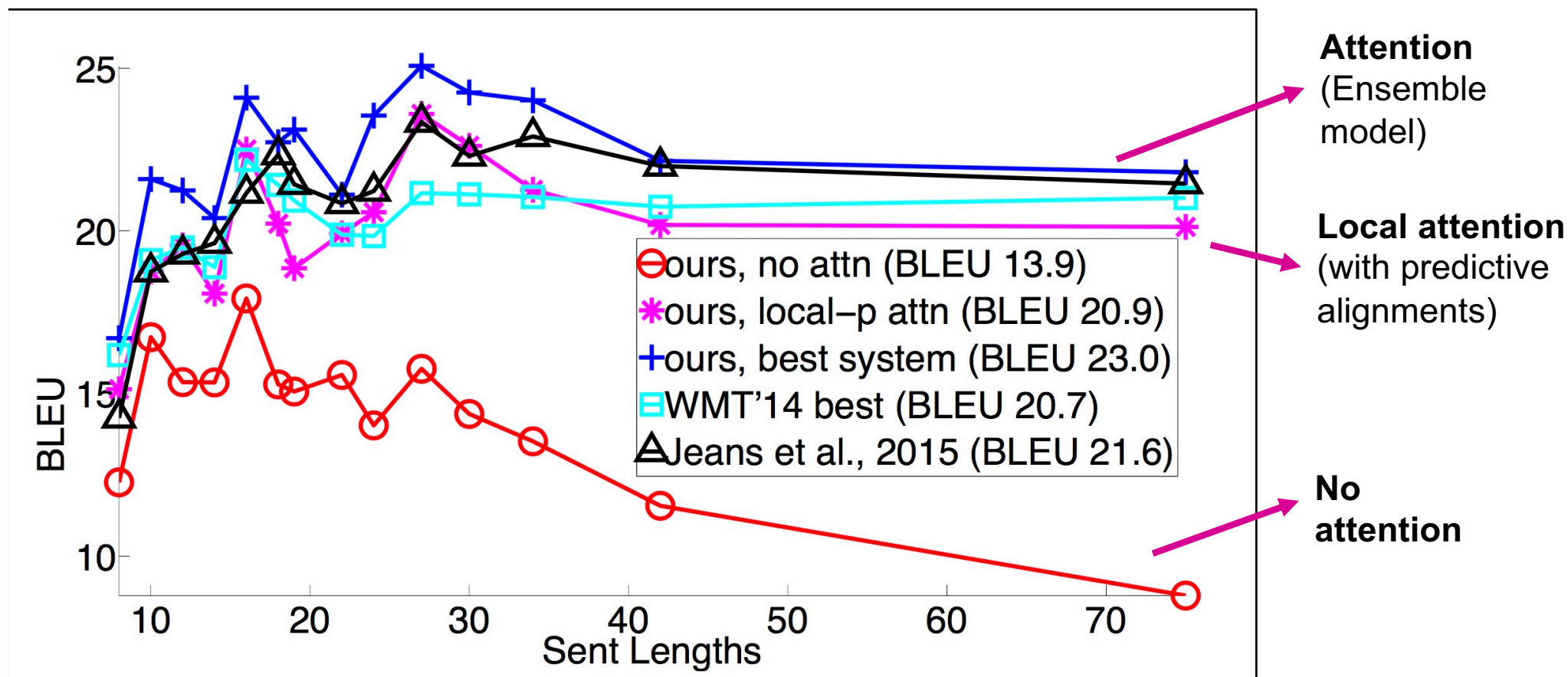


# Local attention



# Results

- English-German translation





# Google Neural Machine Translation (GNMT)

## Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation

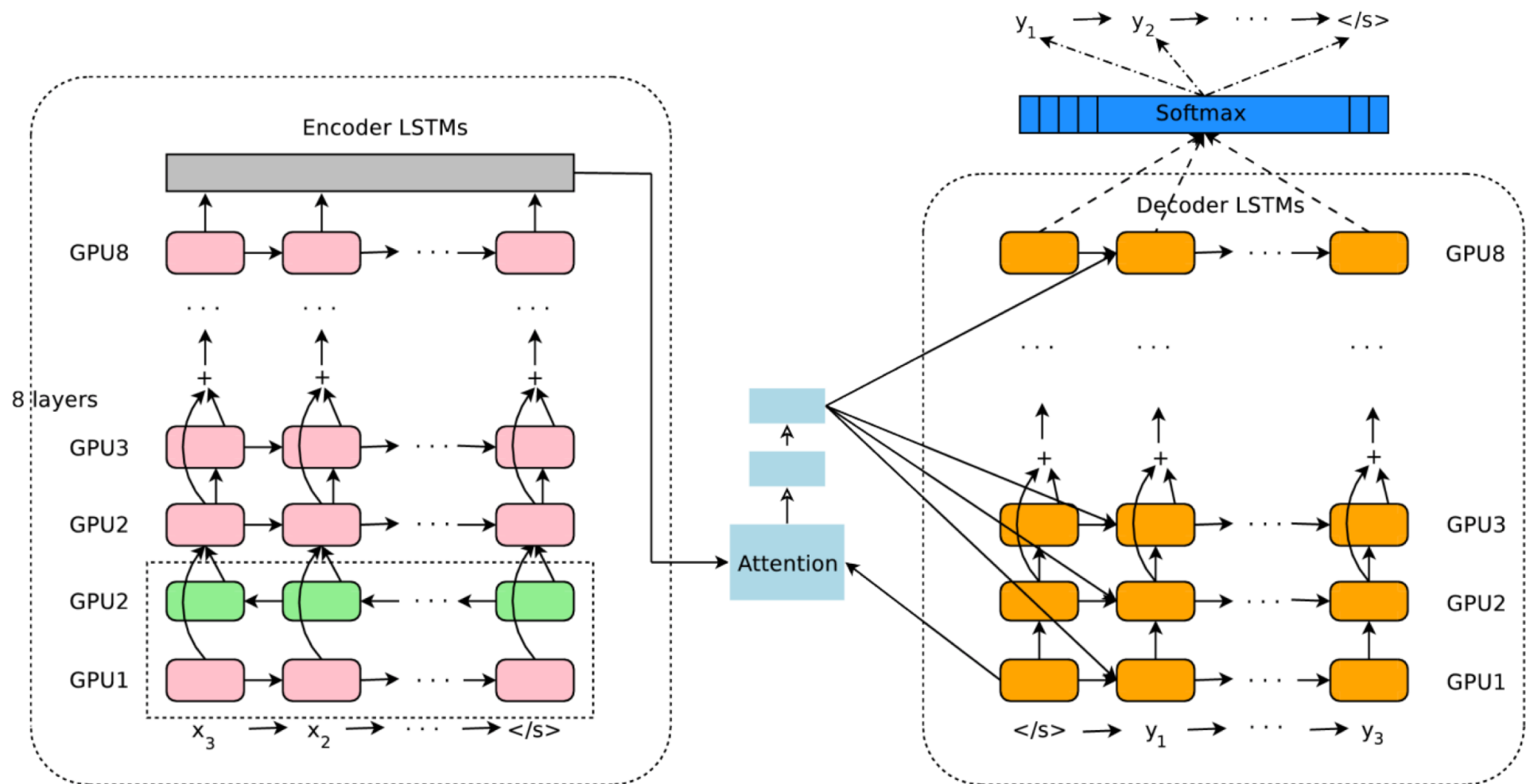
Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi  
yonghui,schuster,zhifengc,qvl,mnorouzi@google.com

Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, Jeffrey Dean

Y. Wu et al., [Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation](#), arXiv 2016

<https://www.nytimes.com/2016/12/14/magazine/the-great-ai-awakening.html>

# Google Neural Machine Translation (GNMT)



Y. Wu et al., [Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation](#), arXiv 2016

# Google Neural Machine Translation (GNMT)

---

- **Standard training objective:** maximize log-likelihood of ground truth output given input:

$$\sum_i \log P_W(Y_i^* | X_i)$$

- Not related to task-specific reward function (e.g., BLEU score)
- Does not encourage “better” predicted sentences to get better likelihood
- **GMNT objective:** expectation of rewards over possible predicted sentences  $Y$ :

$$\sum_i \sum_Y P_W(Y | X_i) r(Y, Y_i^*)$$

- Use variant of BLEU score to compute reward
- Reward is not differentiable -- need reinforcement learning to train (initialize with ML-trained model)

# Google Neural Machine Translation (GNMT)

---

- Results on production data (500 randomly sampled sentences from Wikipedia and news websites)

Table 10: Mean of side-by-side scores on production data

	PBMT	GNMT	Human	Relative Improvement
English → Spanish	4.885	5.428	5.550	87%
English → French	4.932	5.295	5.496	64%
English → Chinese	4.035	4.594	4.987	58%
Spanish → English	4.872	5.187	5.372	63%
French → English	5.046	5.343	5.404	83%
Chinese → English	3.694	4.263	4.636	60%

**Side-by-side scores:** range from 0 (“completely nonsense translation”) to 6 (“perfect translation”), produced by human raters fluent in both languages

**PBMT:** Translation by phrase-based statistical translation system used by Google

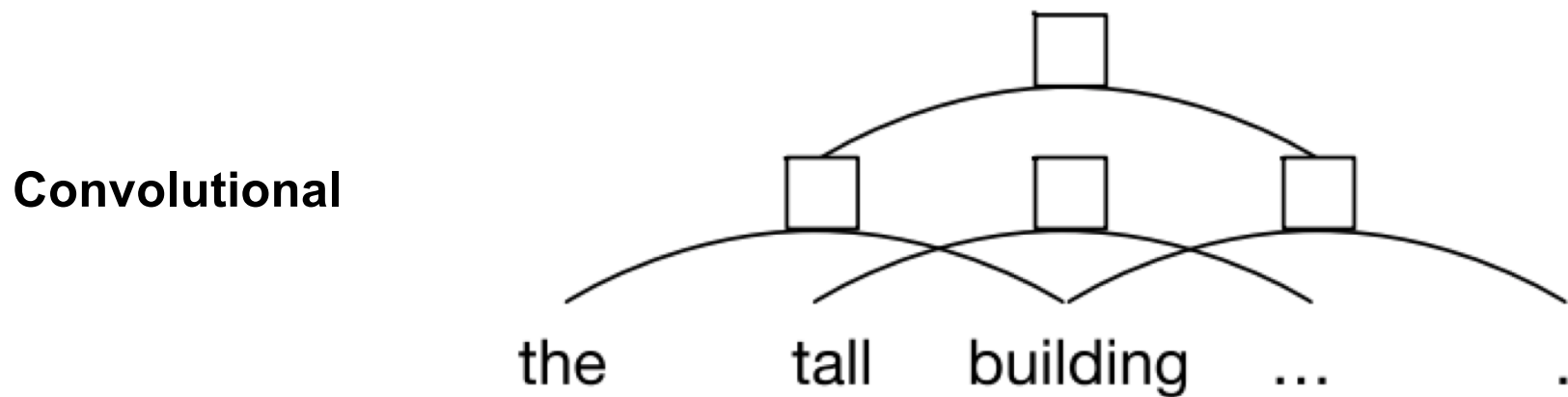
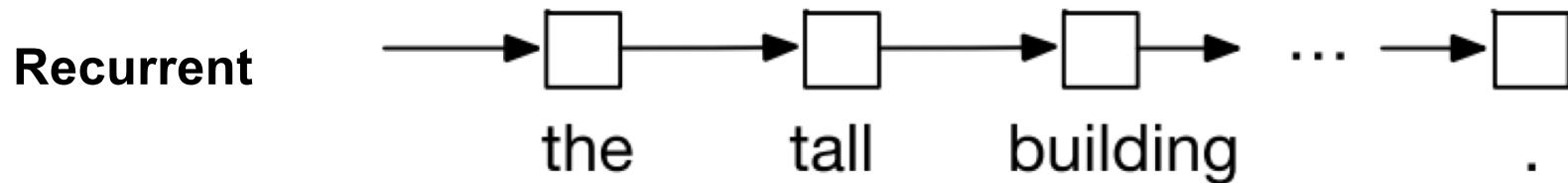
**GNMT:** Translation by our GNMT system

**Human:** Translation by humans fluent in both languages

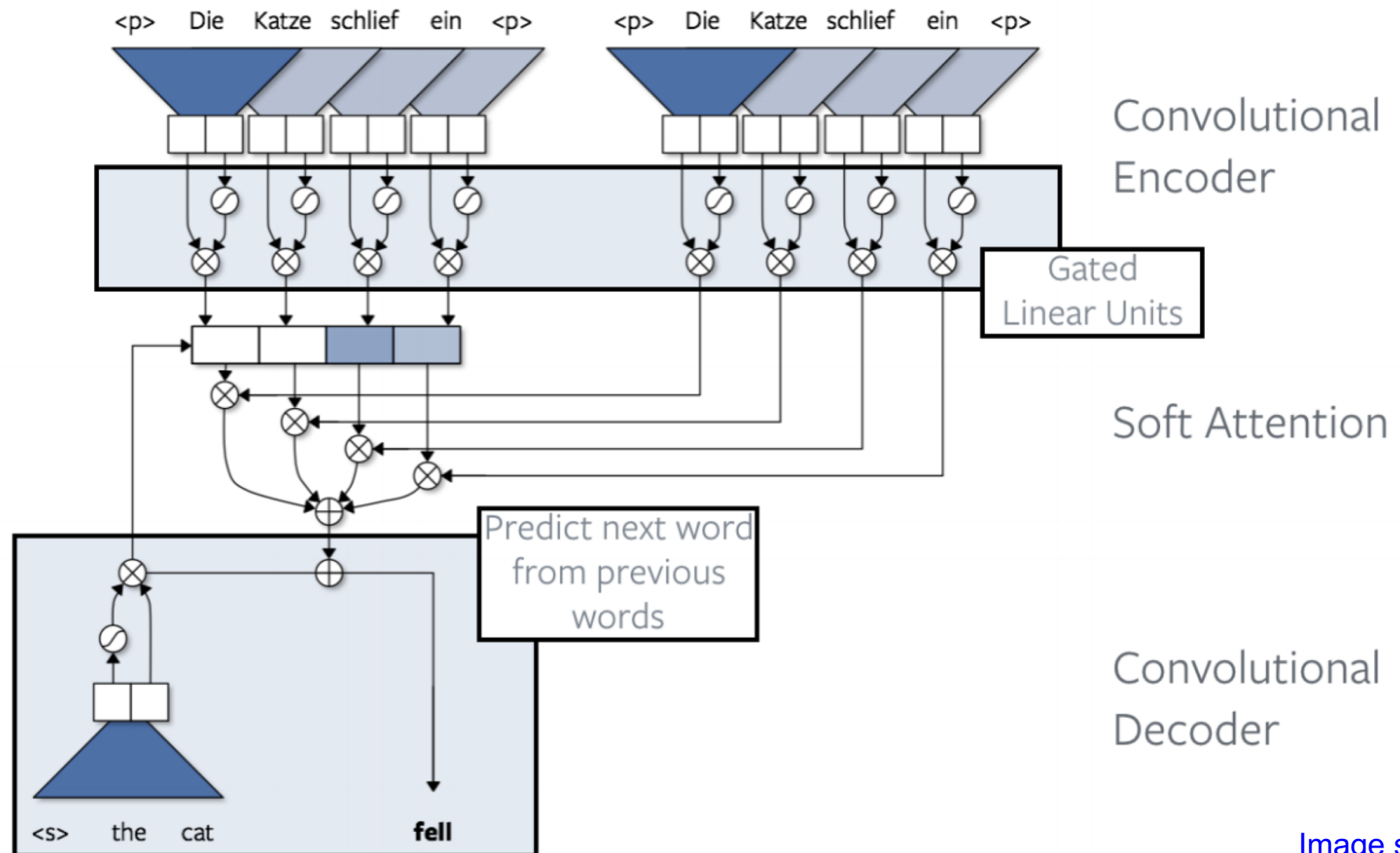
# Convolutional sequence models

---

- Instead of recurrent networks, use 1D convolutional networks



# Convolutional sequence to sequence learning



[Image source](#)

J. Gehring, M. Auli, D. Grangier, D. Yarats, Y. Dauphin, [Convolutional sequence to sequence learning](#), ICML 2017

# Convolutional sequence to sequence learning

---

- Results

<b>WMT'14 English-German</b>	<b>BLEU</b>
Luong et al. (2015) LSTM (Word 50K)	20.9
Kalchbrenner et al. (2016) ByteNet (Char)	23.75
Wu et al. (2016) GNMT (Word 80K)	23.12
Wu et al. (2016) GNMT (Word pieces)	24.61
ConvS2S (BPE 40K)	25.16

<b>WMT'14 English-French</b>	<b>BLEU</b>
Wu et al. (2016) GNMT (Word 80K)	37.90
Wu et al. (2016) GNMT (Word pieces)	38.95
Wu et al. (2016) GNMT (Word pieces) + RL	39.92
ConvS2S (BPE 40K)	40.51

# Convolutional sequence models

---

- From the conclusion:

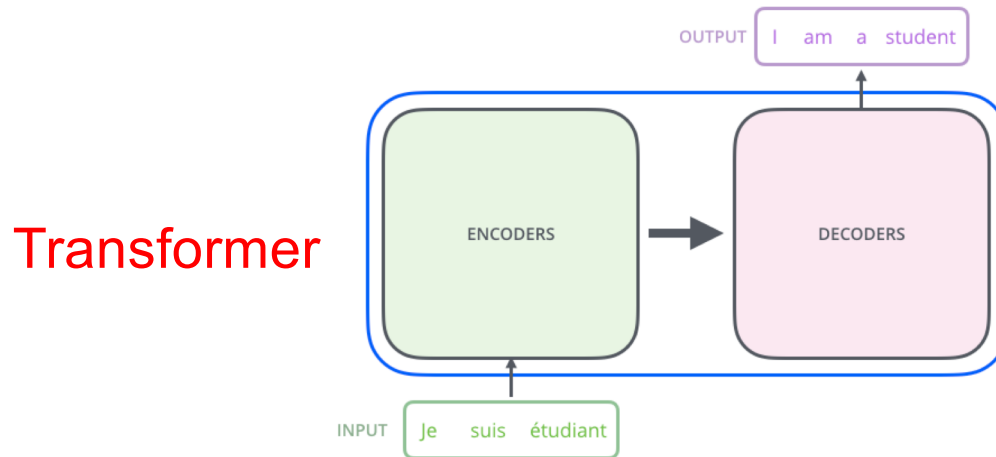
The preeminence enjoyed by recurrent networks in sequence modeling may be largely a vestige of history. Until recently, before the introduction of architectural elements such as dilated convolutions and residual connections, convolutional architectures were indeed weaker. Our results indicate that with these elements, a simple convolutional architecture is more effective across diverse sequence modeling tasks than recurrent architectures such as LSTMs. Due to the comparable clarity and simplicity of TCNs, we conclude that convolutional networks should be regarded as a natural starting point and a powerful toolkit for sequence modeling.



# Attention is all you need

---

- NMT architecture using only FC layers and attention
- More efficient and parallelizable than recurrent or convolutional architectures, faster to train, better accuracy



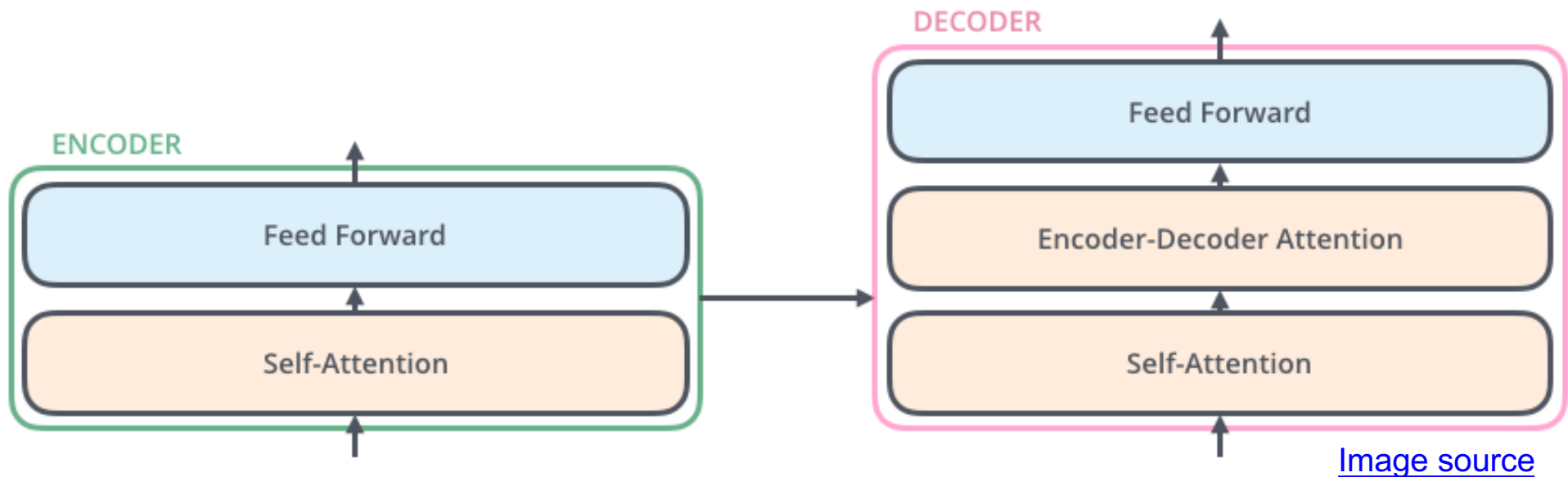
# Attention is all you need

---

- NMT architecture using only FC layers and attention

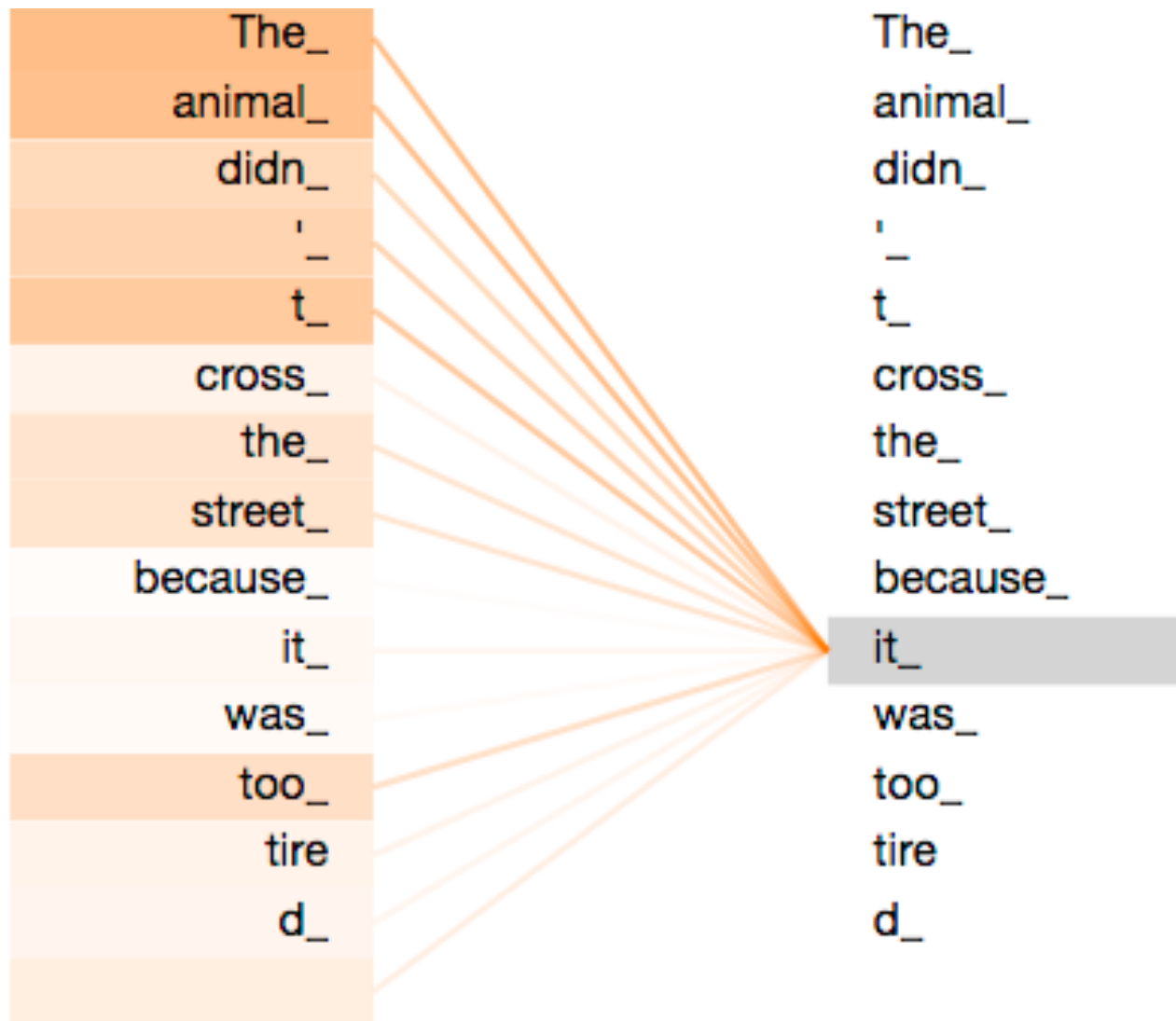
**Encoder:** receives entire input sequence and outputs encoded sequence of the same length

**Decoder:** predicts one word at a time, conditioned on encoder output and previously predicted words

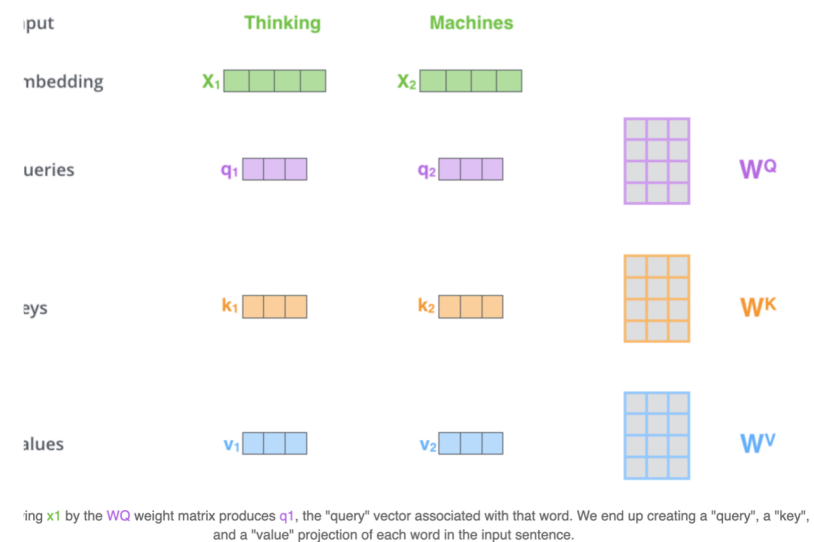
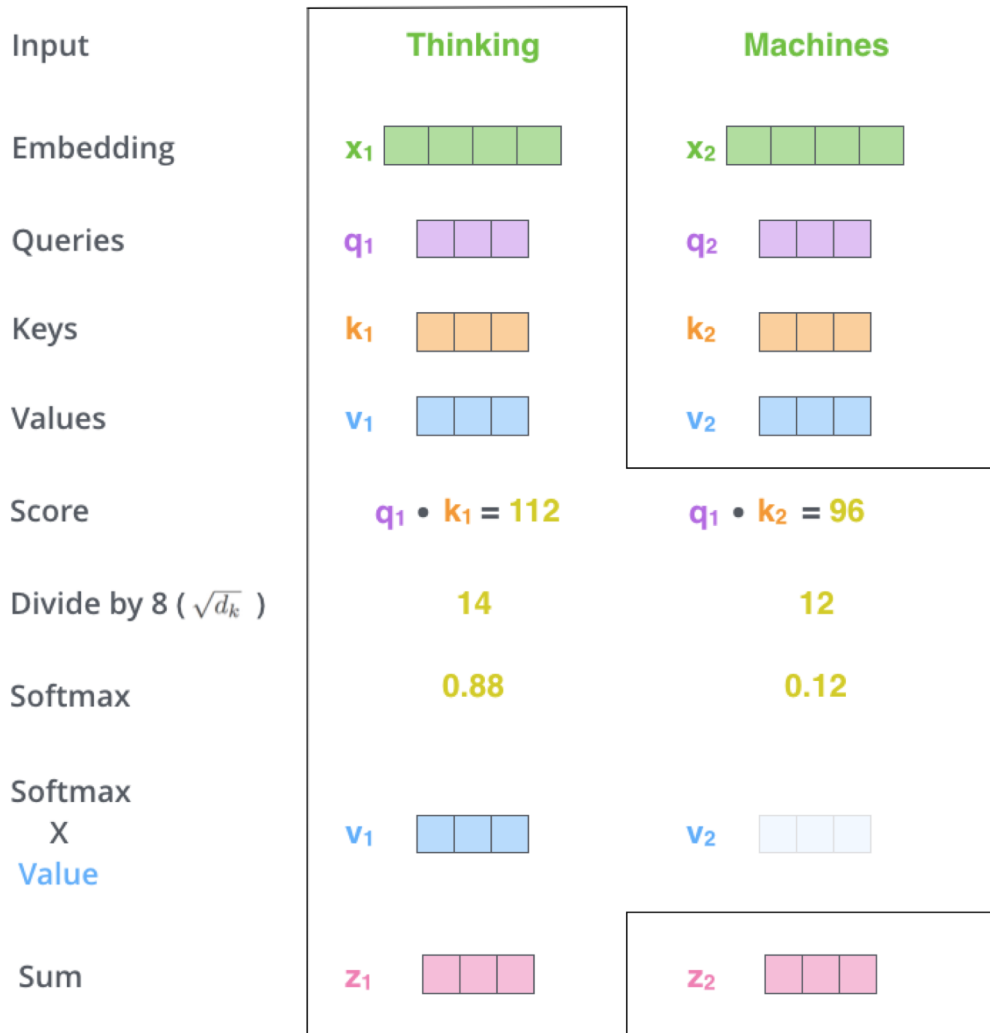


# Self-attention

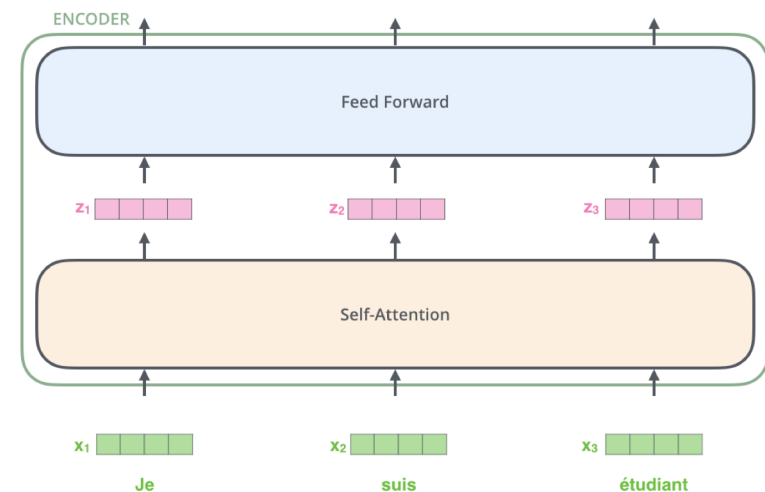
---



# Self-attention details



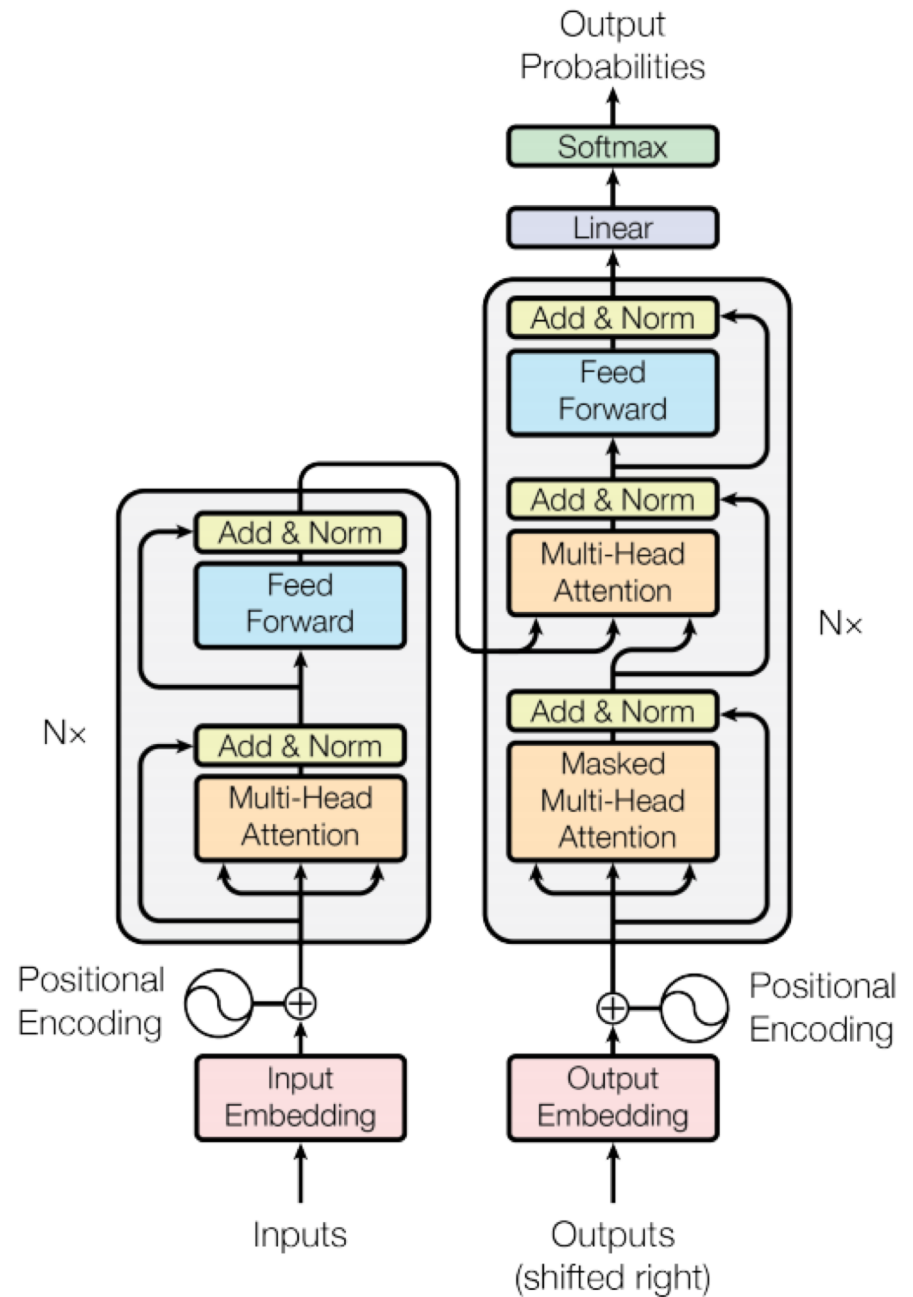
Generate better word embeddings depending on context



# Transformer architecture in detail

Additional bells and whistles

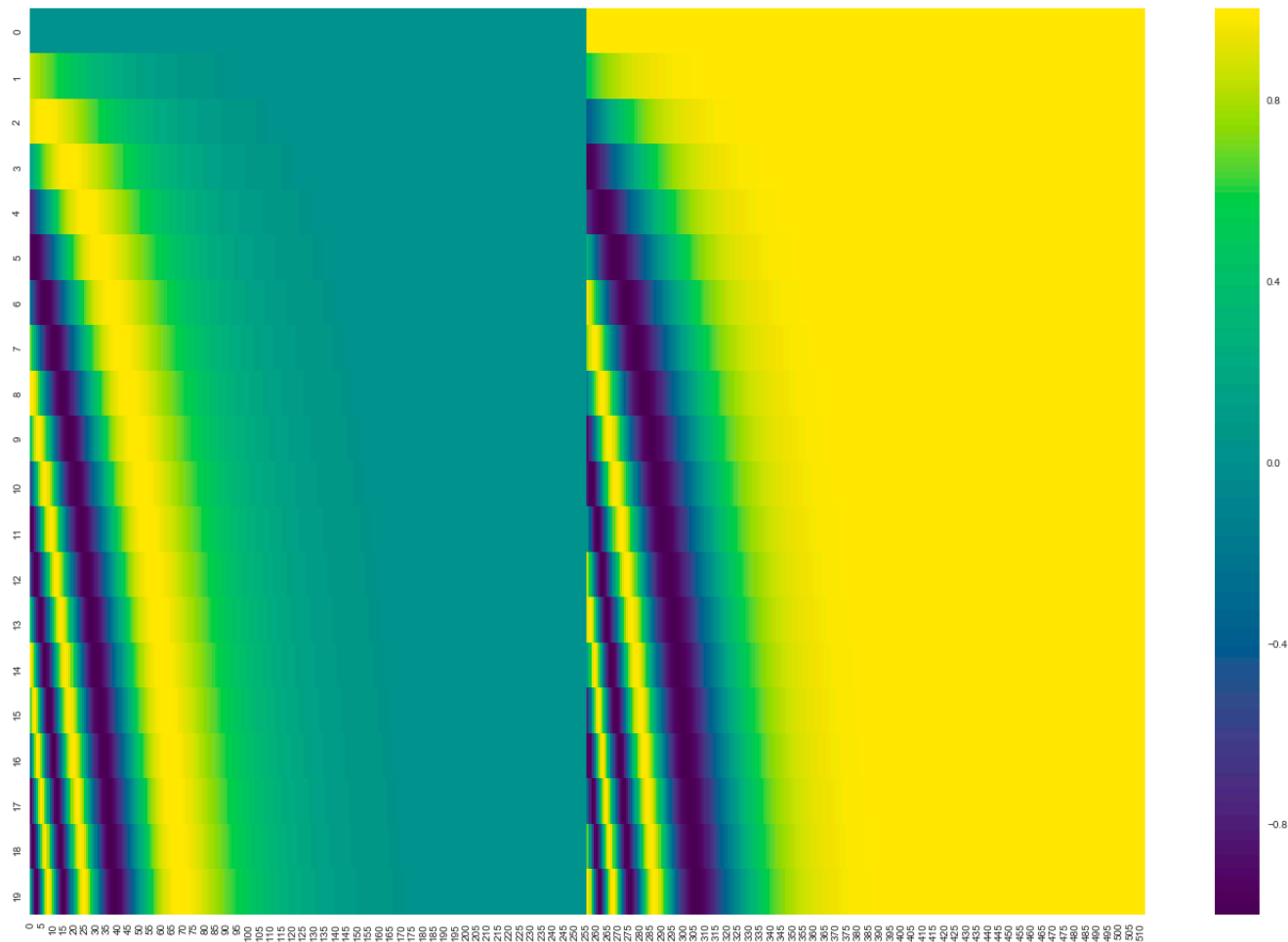
Multiple attention heads



# Positional encoding

---

- Hand-crafted encoding (using sines and cosines) is added to every input vector



# Attention mechanism

---

- *Scaled dot product attention:*

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

- $Q, K, V$  are matrices with rows corresponding to queries, keys, and values,  $d_k$  is the dim. of the keys

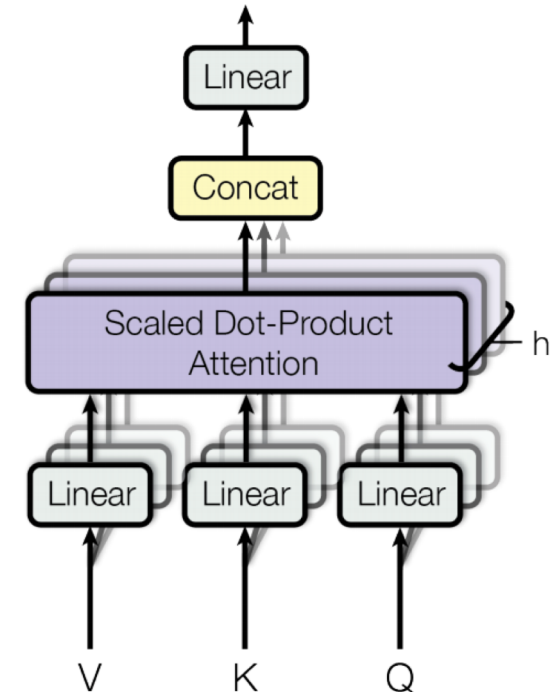
# Attention mechanism

---

- *Scaled dot product attention:*

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

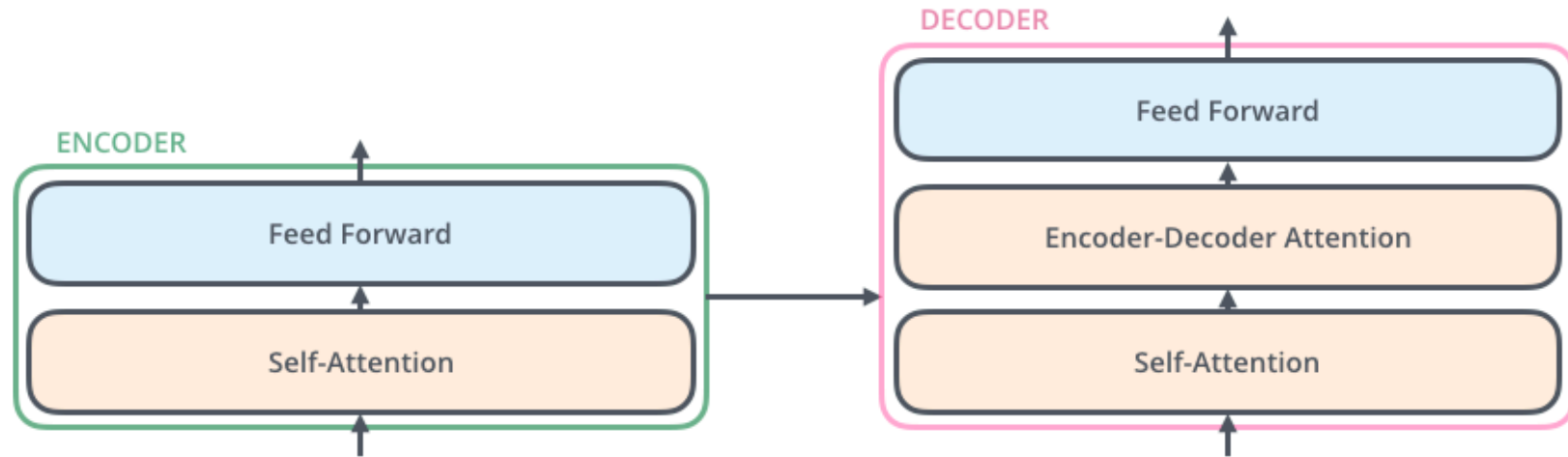
- $Q, K, V$  are matrices with rows corresponding to queries, keys, and values,  $d_k$  is the dim. of the keys
- *Multi-head attention:* run  $h$  attention models in parallel on top of different linearly projected versions of  $Q, K, V$ ; concatenate and linearly project the results





# Attention mechanism

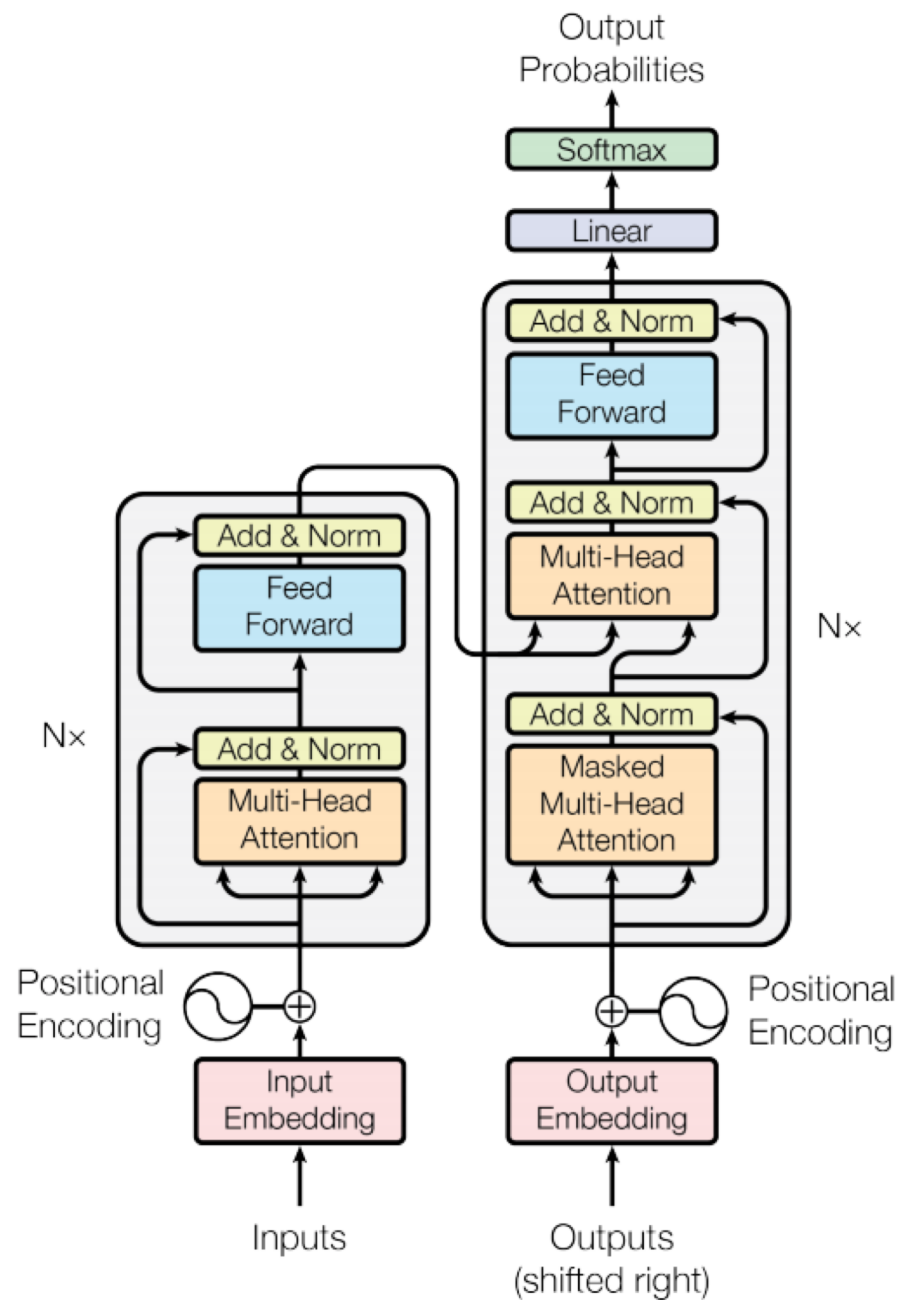
---



- *Encoder-decoder attention:* queries come from previous decoder layer, keys and values come from output of encoder
- *Encoder self-attention:* queries, keys, and values come from previous layer of encoder
- *Decoder self-attention:* values corresponding to future outputs are masked out

# Transformer architecture in detail

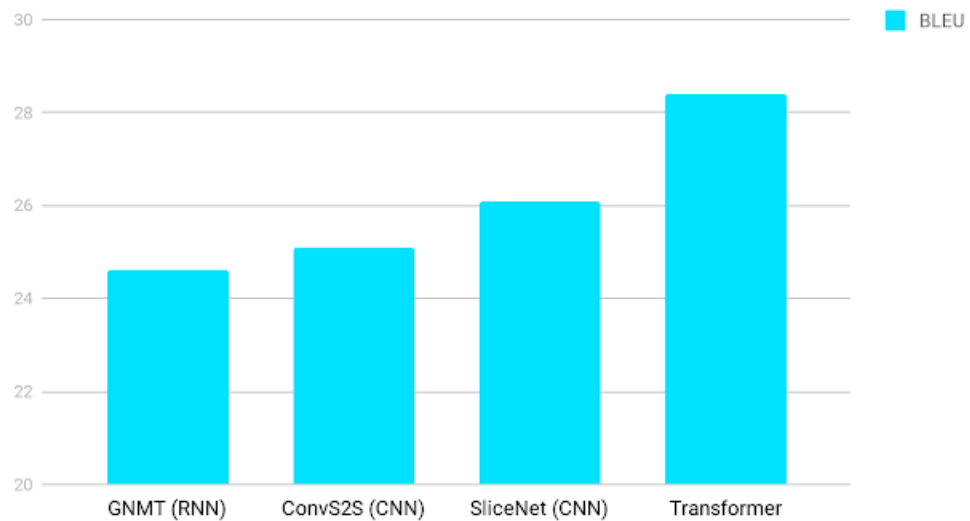
---



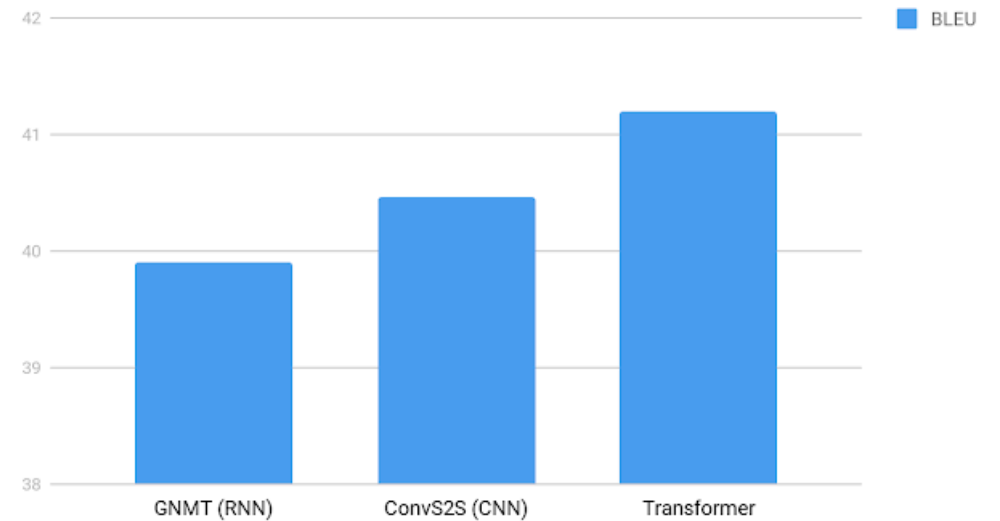
# Results

---

English German Translation quality



English French Translation Quality



<https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>

# Other ideas

---

Training NLP – requires aligned datasets (of ImageNet variety)

***Universal Language Model Fine-tuning for Text Classification*** Jeremy Howard  
, Sebastian Ruder\*

How to finetune the language layers and classifier layers  
(architectures, loss functions, dropout)

Better word embeddings trained directly from language models

***Deep contextualized word representations*** Matthew E. Peters† , Mark  
Neumann† , Mohit Iyyer† , Matt Gardner

Better ways to represent vectors ELMo vector assigned to a token  
or word is actually a function of the entire sentence containing  
that word. Therefore, the same word can have different word  
vectors under different contexts.

ELMo word top of a two-layer bidirectional language model (biLM).

This biLM model has two layers stacked together.

Trained in unsupervised way

# ELMO

---



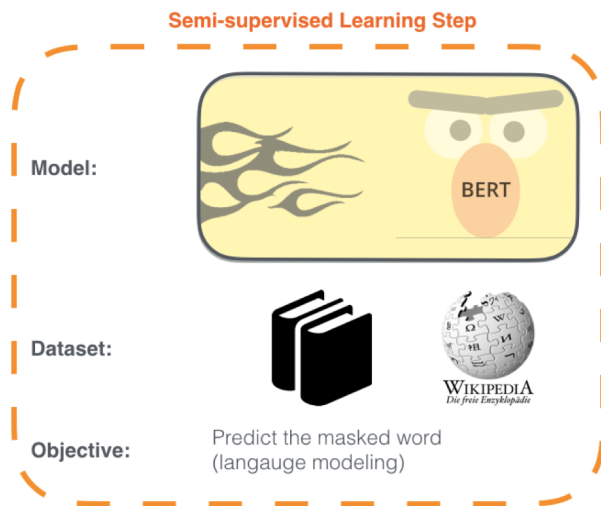
# BERT

## Bidirectional Encoder Representations from Transformers

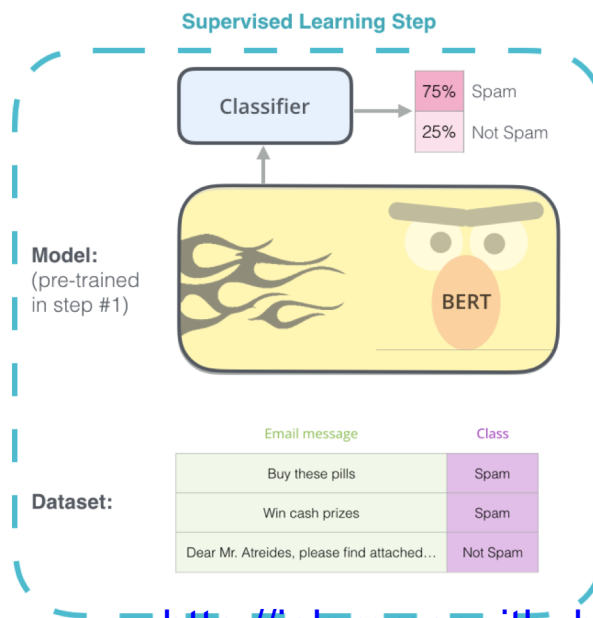
- Pretrain NLP representations
- Universal language models which can be adapted to many language tasks
- Seq2Seq models + attention – good to machine translation
- What about other language tasks ?

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



2 - **Supervised** training on a specific task with a labeled dataset.

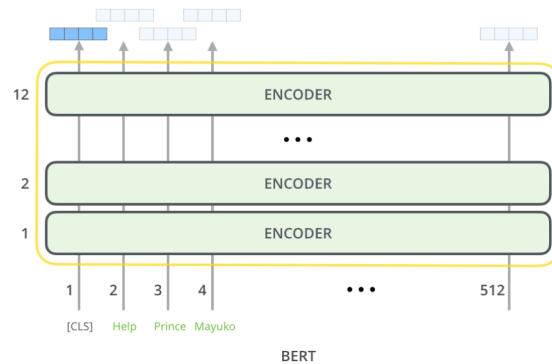
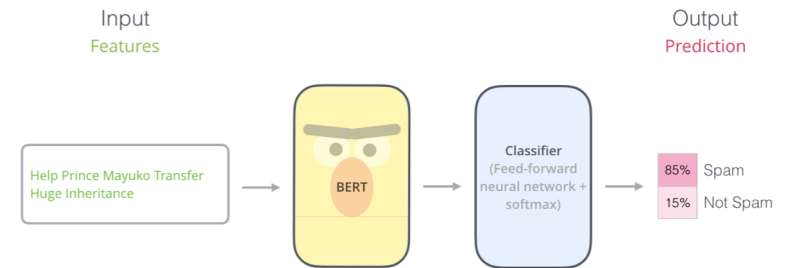


<http://jalammar.github.io/illustrated-bert/>

---

It can be used for different tasks

- Spam/not spam
- Fact checking fact/no fact
- Sentiment analysis positive/not
- Visual question answering
  
- Bert pretrained encoder of the transformer
- For classification – focus only on output in the first token – feed to feedforward NN



# Parting thoughts

---

- Methodology for text generation problems
  - Evaluation is tricky
  - Maximum likelihood training is not the most appropriate (but alternatives involve optimizing non-differentiable objectives)
- Attention appears to be a game-changer for NMT (for image captioning, not as much)
  - But there is much more to MT than attention (dealing with unknown words, etc.)
- Recurrent architectures are not the only option for sequence modeling
  - Convolutional and feedforward alternatives should be considered