

Sequence-to-Sequence Architectures

Kapil Thadani
kapil@cs.columbia.edu



Outline

- Machine translation
 - Phrase-based MT
 - Encoder-decoder architecture
- Attention mechanism
 - Bahdanau et al (2015)
 - Visualizations
 - Variants
 - Transformers
- Decoding large vocabularies
 - Alternative approaches
 - Copying mechanism
- Autoencoders
 - Denoising autoencoders
 - Variational autoencoders (VAEs)

Previously: processing text with RNNs

Inputs

- One-hot vectors for words/characters/previous output
- Embeddings for words/sentences/context
- CNN over characters/words/sentences

⋮

Recurrent layers

- Forward, backward, bidirectional, deep
- Activations: σ , tanh, gated (LSTM, GRU), ReLU initialized with identity

⋮

Outputs

- Softmax over words/characters/labels
- Absent (i.e., pure encoders)

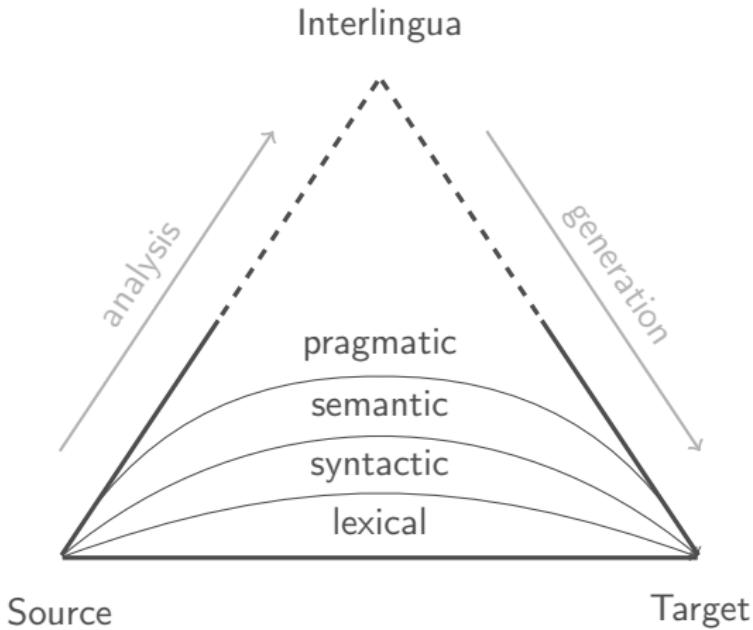
⋮

Machine Translation

"One naturally wonders if the problem of translation could conceivably be treated as a problem in cryptography. When I look at an article in Russian, I say: 'This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.'

— Warren Weaver
Translation (1955)

The MT Pyramid

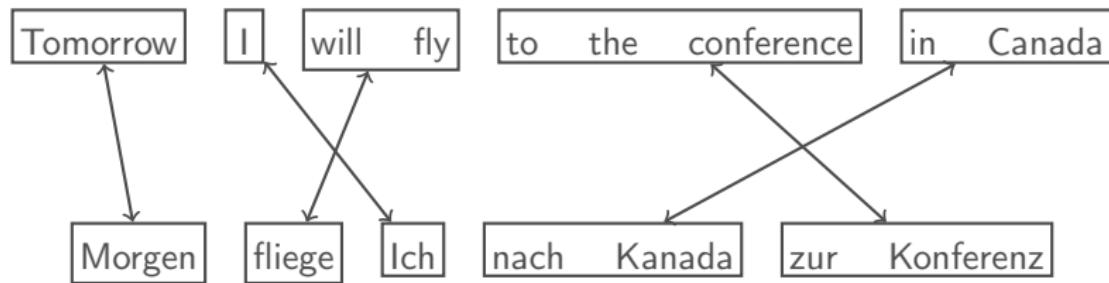


Phrase-based MT

Tomorrow I will fly to the conference in Canada

Morgen fliege Ich nach Kanada zur Konferenz

Phrase-based MT



Phrase-based MT

1. Collect bilingual dataset $\langle S_i, T_i \rangle \in \mathcal{D}$
 2. Unsupervised phrase-based alignment
 - ▶ phrase table π
 3. Unsupervised n-gram language modeling
 - ▶ language model ψ
 4. Supervised decoder
 - ▶ parameters θ
- $$\begin{aligned}\hat{T} &= \arg \max_T p(T|S) \\ &= \arg \max_T p(S|T, \pi, \theta) \cdot p(T|\psi)\end{aligned}$$

	kdybys	tam	byl	.	ted'	bys	to	věděl
if	█							
you			█					
were		█						
there								
you					█			
would						█		
know							█	
it					█			
now						█		

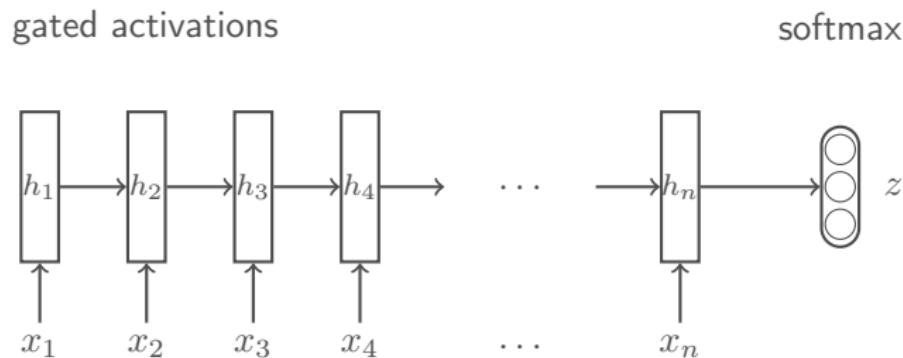
Neural MT

1. Collect bilingual dataset $\langle S_i, T_i \rangle \in \mathcal{D}$
2. Unsupervised phrase-based alignment
 - ▶ phrase tables
3. Unsupervised n-gram language modeling
 - ▶ language models
4. Supervised encoder-decoder framework
 - ▶ parameters θ

RNN

Input words x_1, \dots, x_n

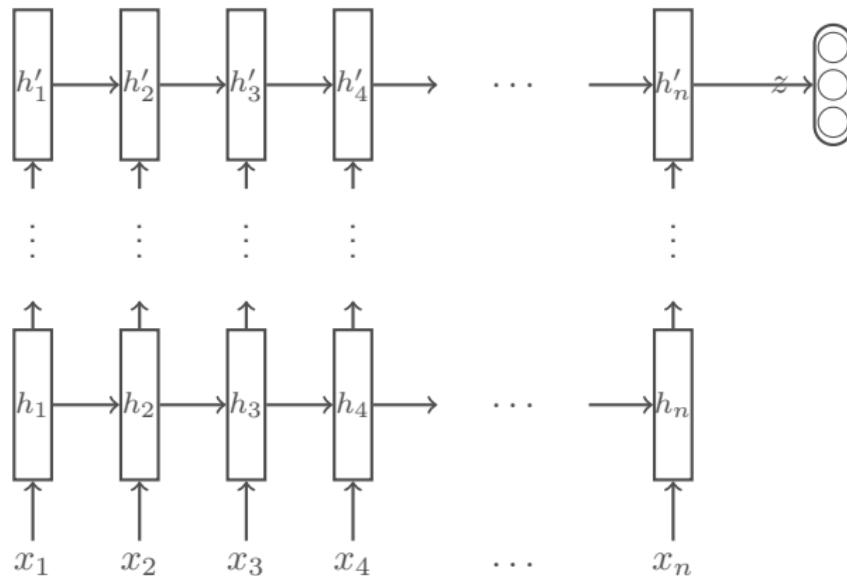
Output label z



Deep RNN

Input words x_1, \dots, x_n

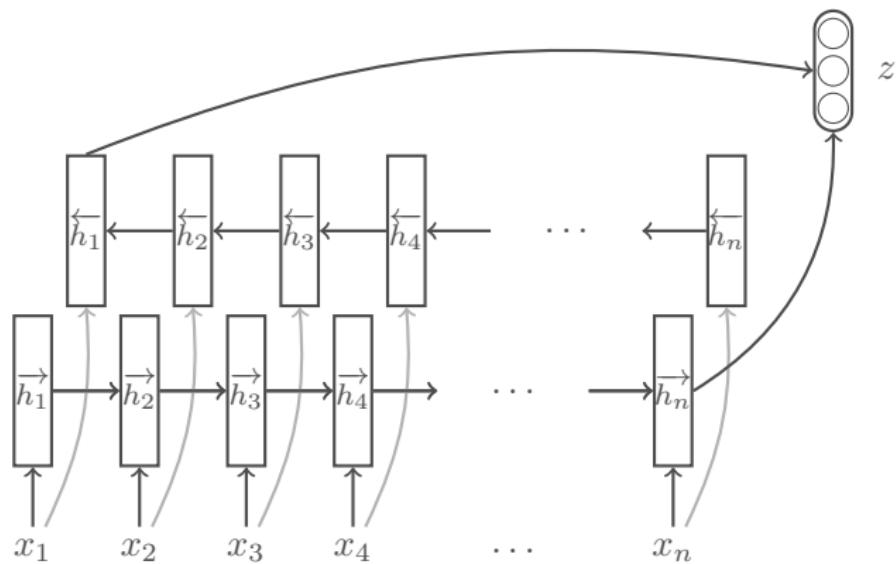
Output label z



Bidirectional RNN

Input words x_1, \dots, x_n

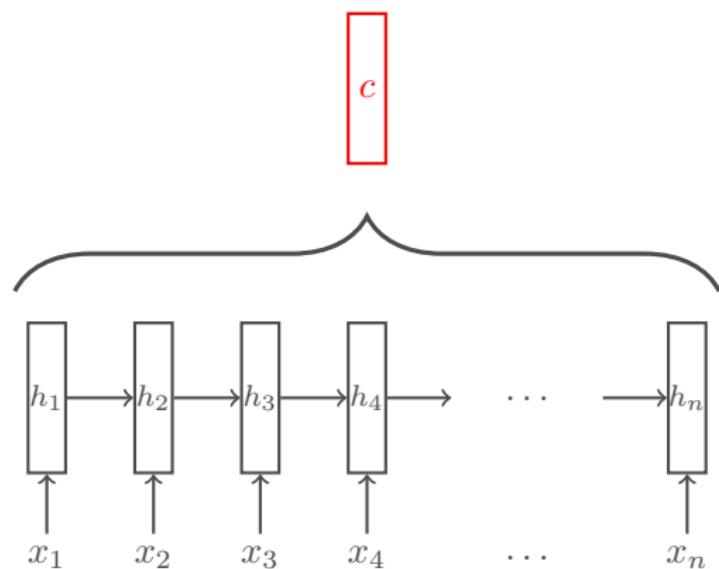
Output label z



RNN encoder

Input words x_1, \dots, x_n

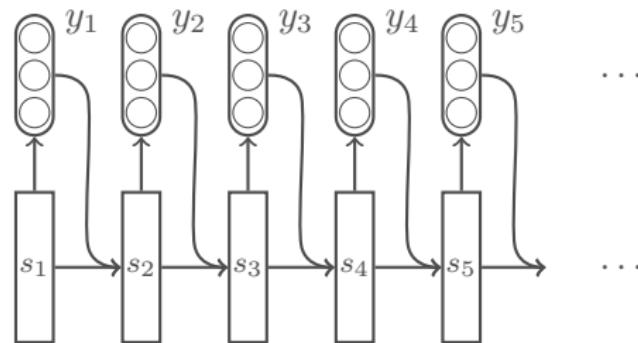
Output encoding c



RNN language model

Input words y_1, \dots, y_k

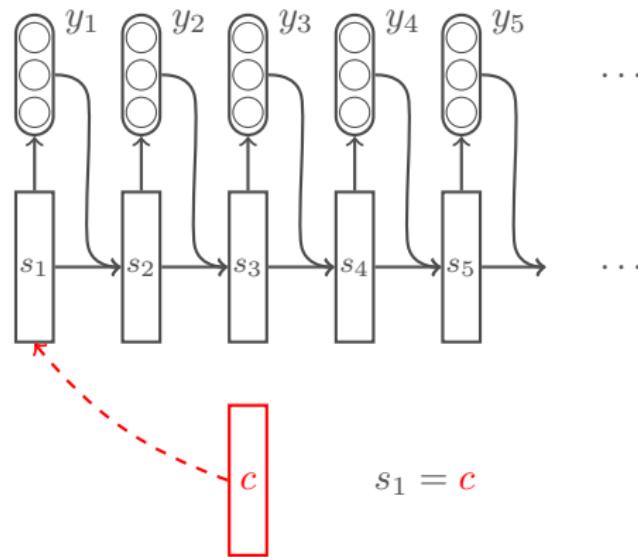
Output following words y_k, \dots, y_m



RNN decoder

Input context c

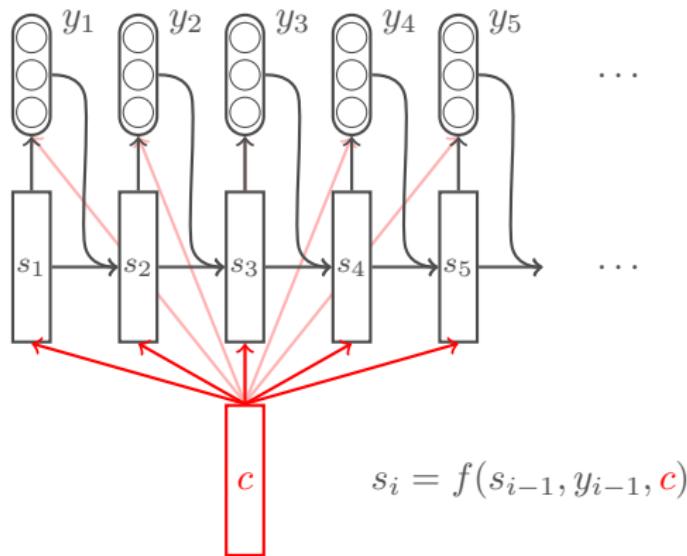
Output words y_1, \dots, y_m



RNN decoder

Input context c

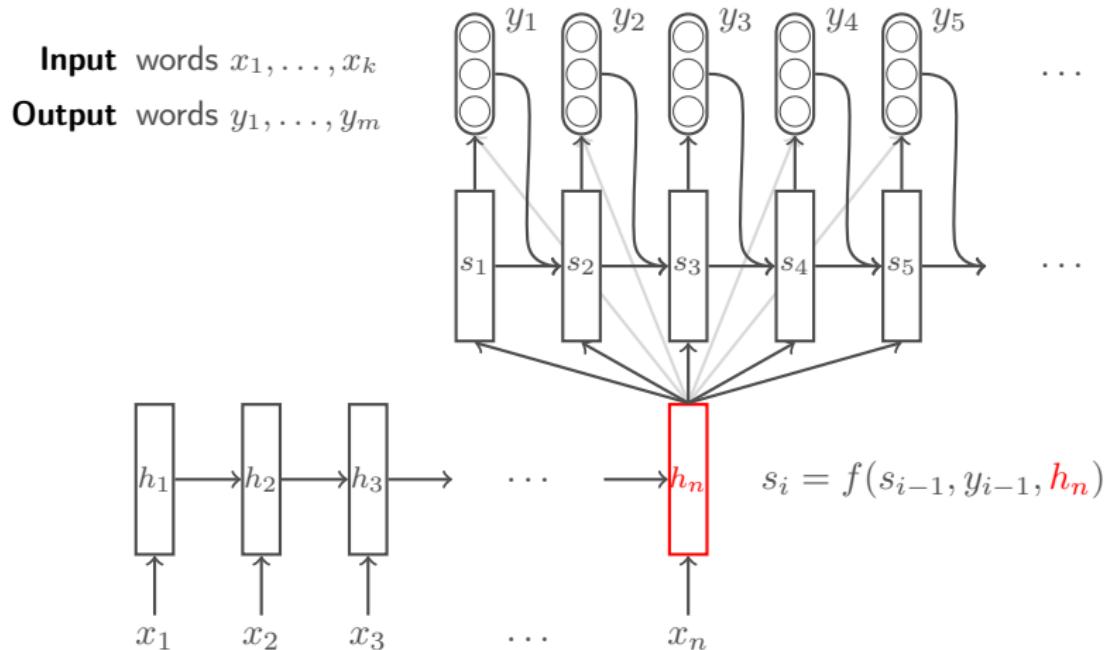
Output words y_1, \dots, y_m



Sequence-to-sequence learning

Sutskever, Vinyals & Le (2014)

Sequence to Sequence Learning with Neural Networks



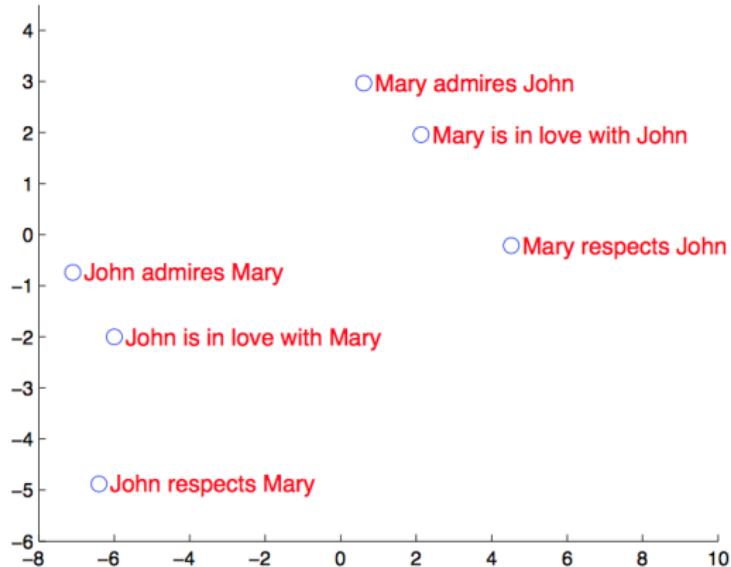
Sequence-to-sequence learning

Sutskever, Vinyals & Le (2014)

Sequence to Sequence Learning with Neural Networks

Produces a fixed length representation of input

- “sentence embedding” or “thought vector”



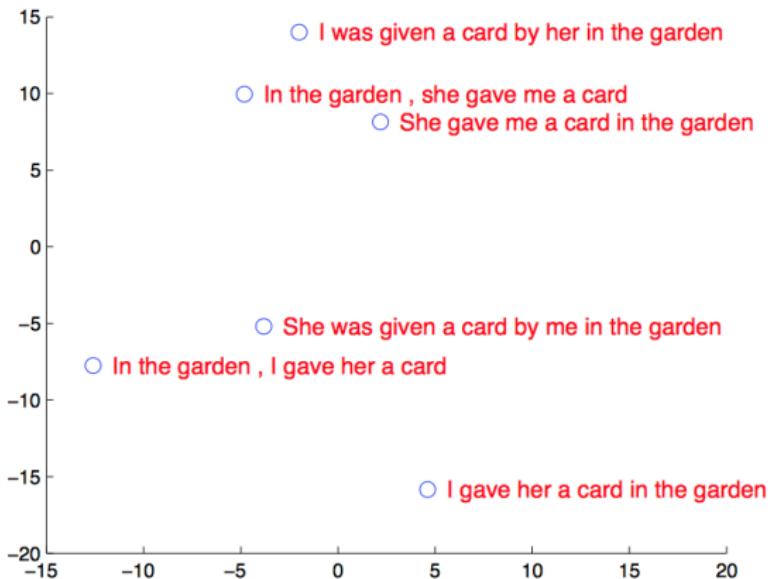
Sequence-to-sequence learning

Sutskever, Vinyals & Le (2014)

Sequence to Sequence Learning with Neural Networks

Produces a fixed length representation of input

- “sentence embedding” or “thought vector”



Sequence-to-sequence learning

Sutskever, Vinyals & Le (2014)

Sequence to Sequence Learning with Neural Networks

LSTM units do not solve vanishing gradients

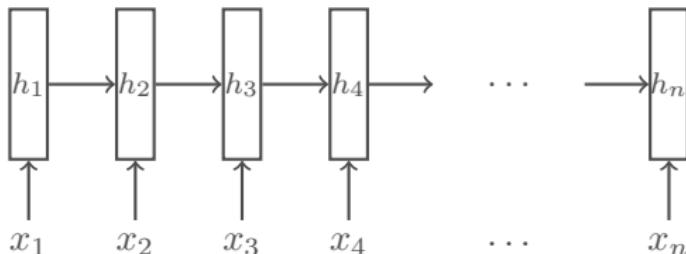
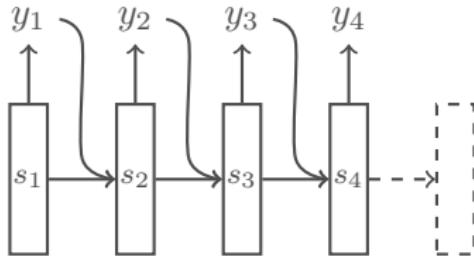
- Poor performance on long sentences
- Need to reverse the input

Method	test BLEU score (ntst14)
Bahdanau et al. [2]	28.45
Baseline System [29]	33.30
Single forward LSTM, beam size 12	26.17
Single reversed LSTM, beam size 12	30.59

Attention-based translation

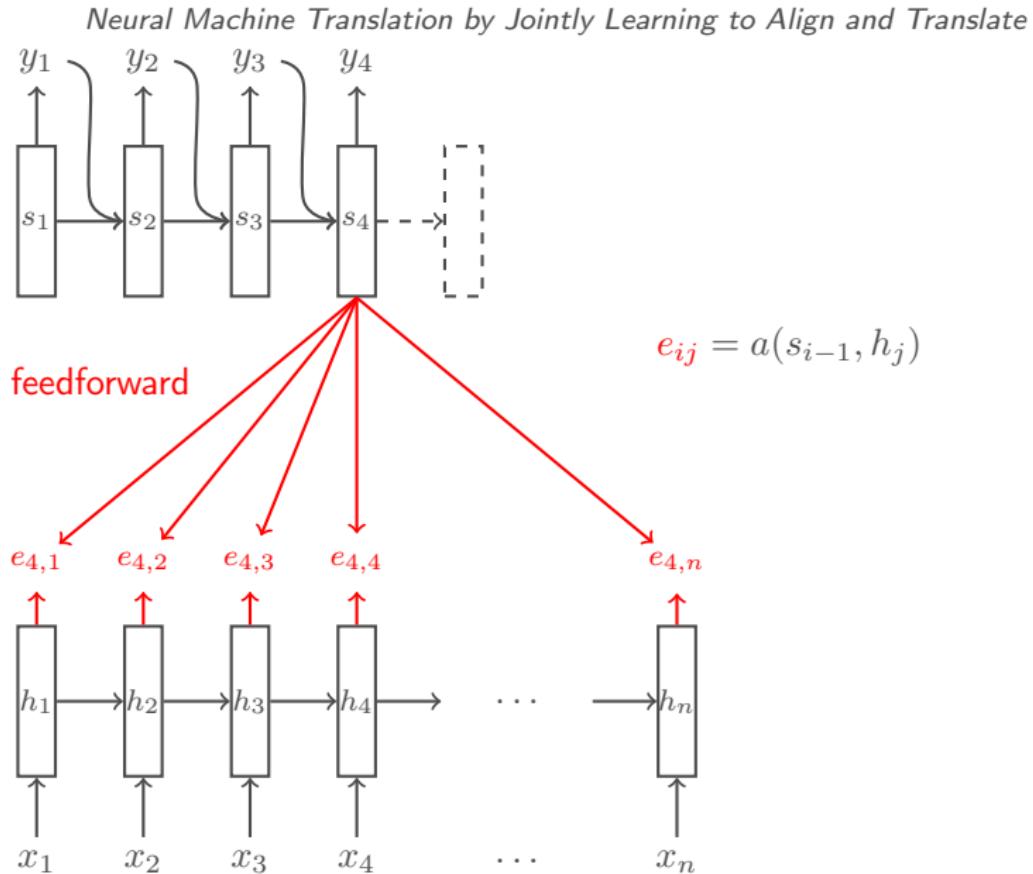
Bahdanau et al (2015)

Neural Machine Translation by Jointly Learning to Align and Translate



Attention-based translation

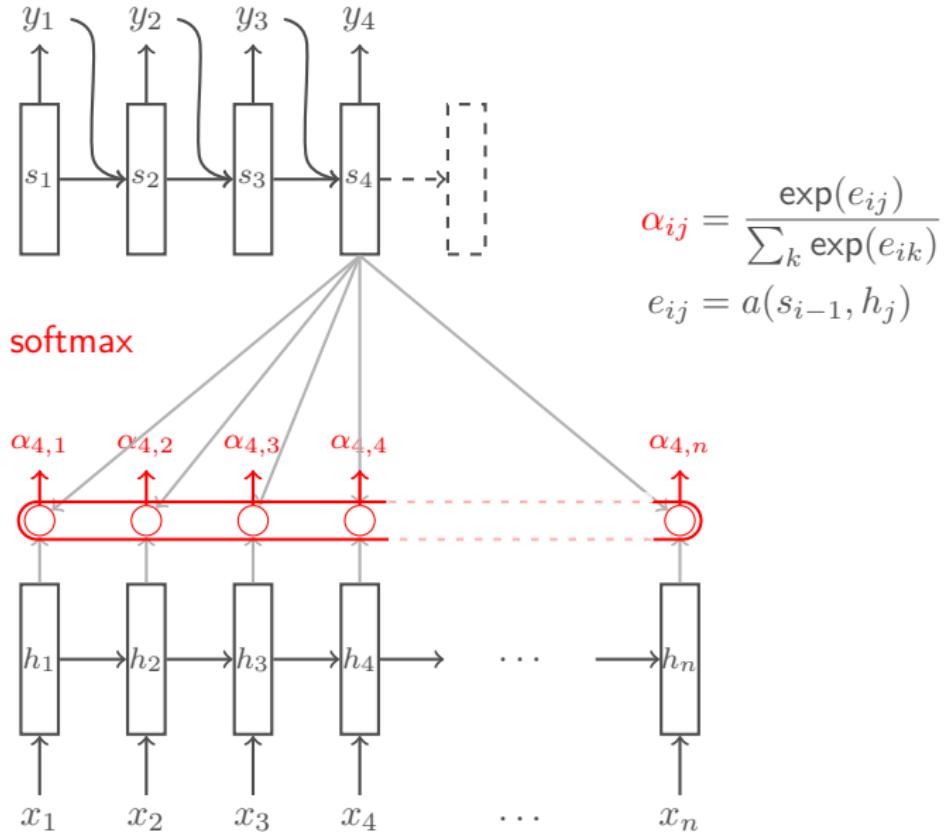
Bahdanau et al (2015)



Attention-based translation

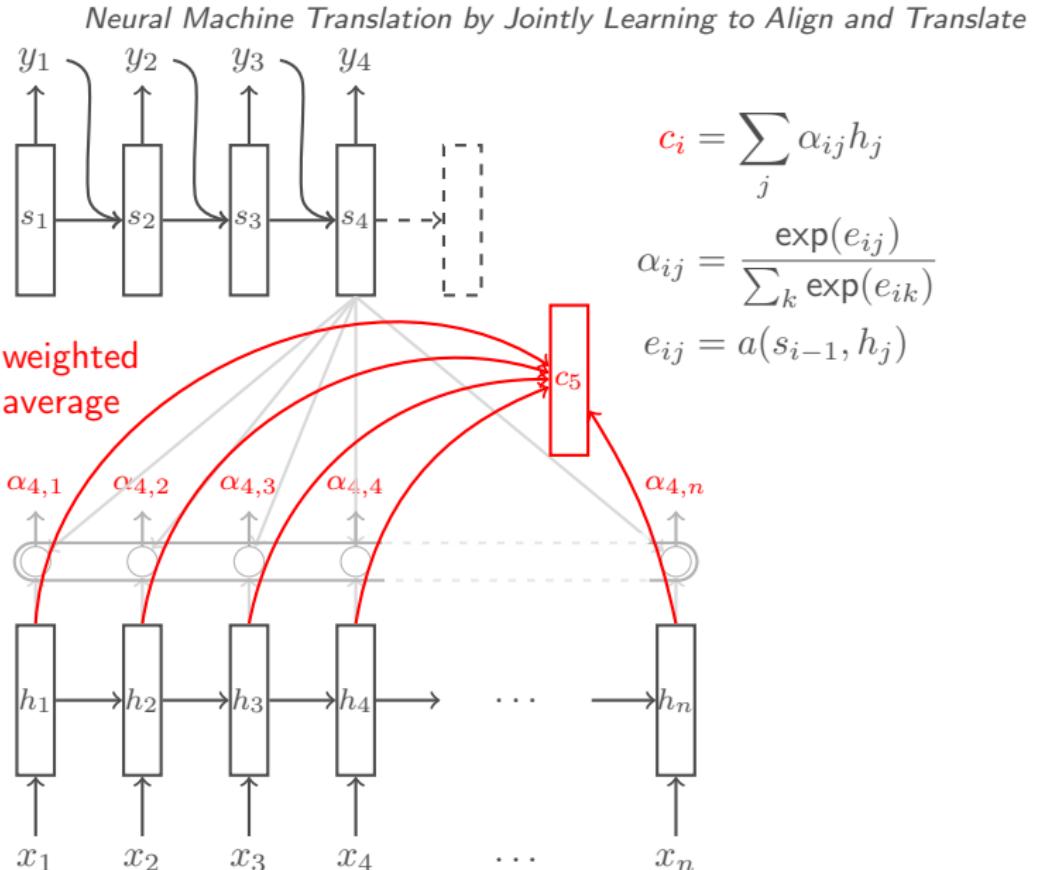
Bahdanau et al (2015)

Neural Machine Translation by Jointly Learning to Align and Translate



Attention-based translation

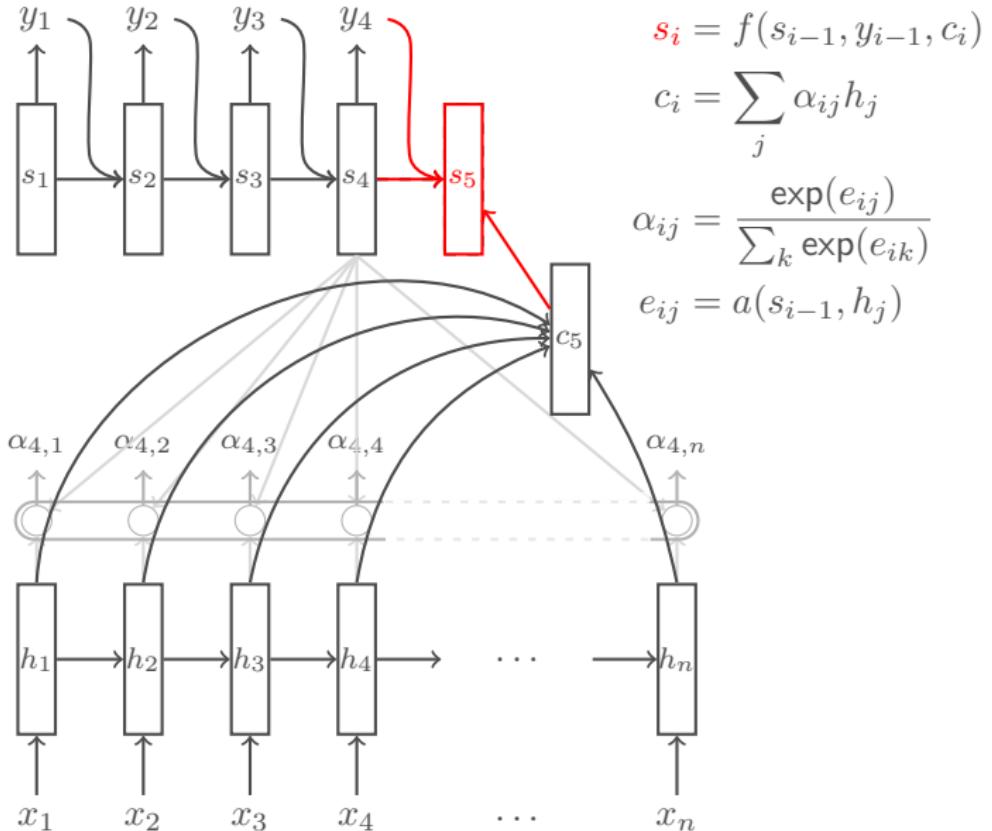
Bahdanau et al (2015)



Attention-based translation

Bahdanau et al (2015)

Neural Machine Translation by Jointly Learning to Align and Translate



Attention-based translation

Bahdanau et al (2015)

Neural Machine Translation by Jointly Learning to Align and Translate

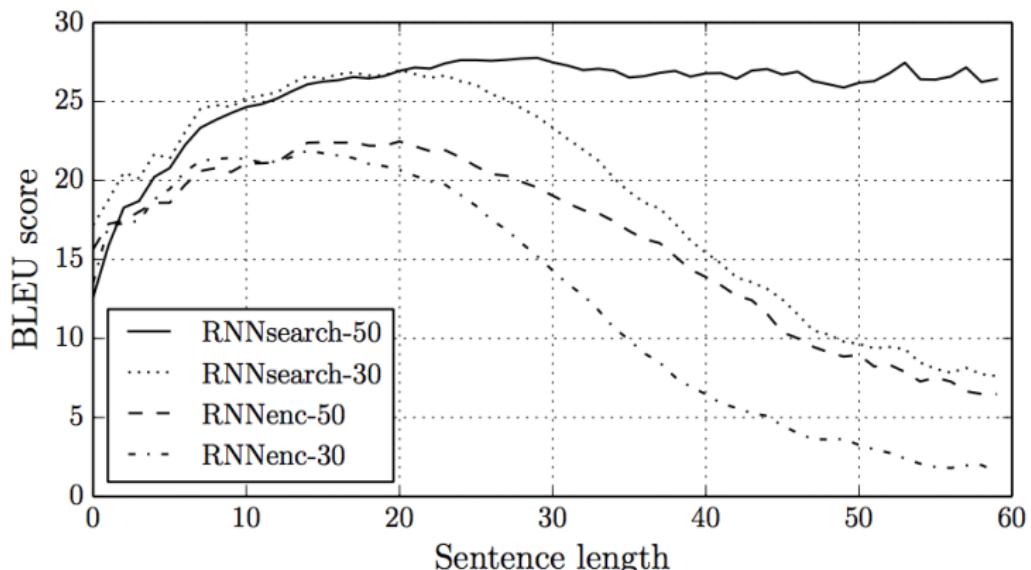
- Bidirectional encoder, GRU activations
- Softmax for y_i depends on y_{i-1} and an additional hidden layer

- + Backprop directly to attended regions, avoiding vanishing gradients
- + Can visualize attention weights α_{ij} to interpret prediction
- Inference is $\mathcal{O}(mn)$ instead of $\mathcal{O}(m)$ for seq-to-seq

Attention-based translation

Bahdanau et al (2015)

Neural Machine Translation by Jointly Learning to Align and Translate

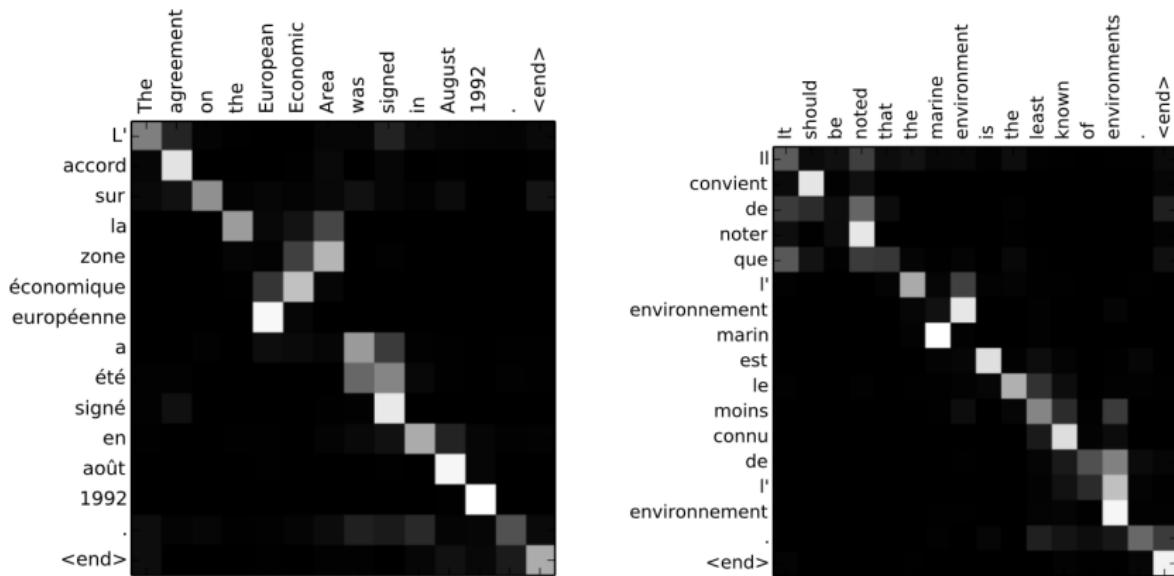


Improved results on long sentences

Attention-based translation

Bahdanau et al (2015)

Neural Machine Translation by Jointly Learning to Align and Translate



Sensible induced alignments

Attention-based translation

16

Bahdanau et al (2015)

Neural Machine Translation by Jointly Learning to Align and Translate

La destruction de l'équipement signifie que la Syrie ne peut plus produire de nouvelles armes chimiques.

<end>

Destruction
of
the
equipment
means
that
Syria
can
no
longer
produce
new
chemical
weapons

<end>

Sensible induced alignments

Natural language inference

Given a premise, e.g.,

The purchase of Houston-based LexCorp by BMI for \$2Bn prompted widespread sell-offs by traders as they sought to minimize exposure. LexCorp had been an employee-owned concern since 2008.

and a hypothesis, e.g.,

BMI acquired an American company. (1)

predict whether the premise

- o entails the hypothesis
- o contradicts the hypothesis
- o or remains neutral

Natural language inference

Given a premise, e.g.,

The purchase of Houston-based LexCorp by BMI for \$2Bn prompted widespread sell-offs by traders as they sought to minimize exposure. LexCorp had been an employee-owned concern since 2008.

and a hypothesis, e.g.,

BMI bought employee-owned LexCorp for \$3.4Bn. (2)

predict whether the premise

- o entails the hypothesis
- o contradicts the hypothesis
- o or remains neutral

Natural language inference

Given a premise, e.g.,

The purchase of Houston-based LexCorp by BMI for \$2Bn prompted widespread sell-offs by traders as they sought to minimize exposure. LexCorp had been an employee-owned concern since 2008.

and a hypothesis, e.g.,

BMI is an employee-owned concern. (3)

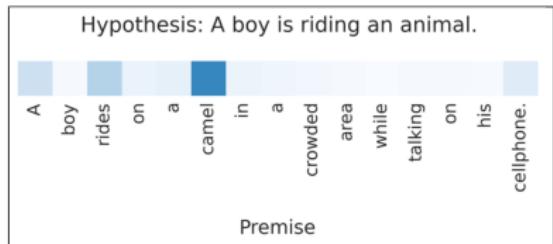
predict whether the premise

- o entails the hypothesis
- o contradicts the hypothesis
- o or remains neutral

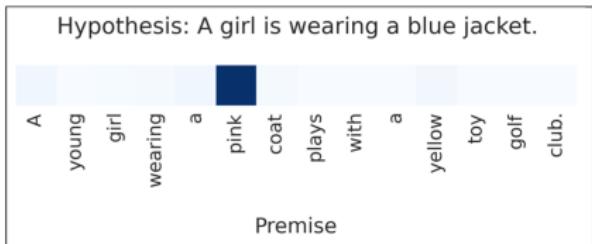
Natural language inference

Rocktäschel et al (2016)

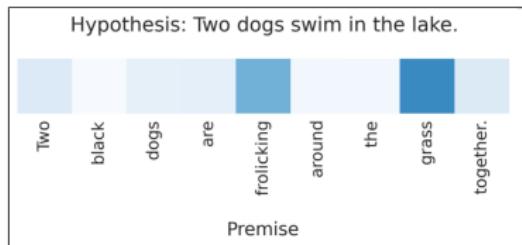
Reasoning about Entailment with Neural Attention



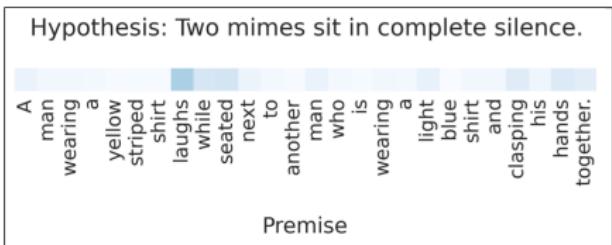
(a)



(b)



(c)



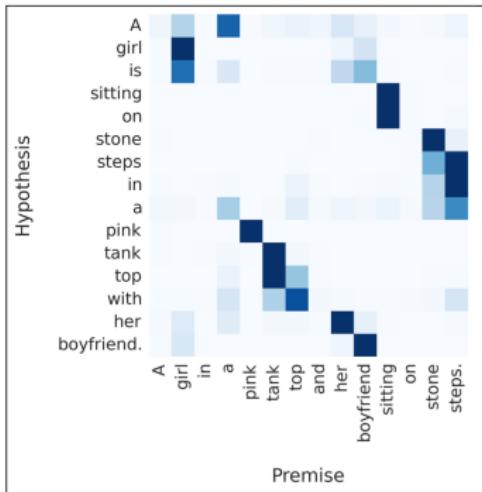
(d)

Attention conditioned on h_T

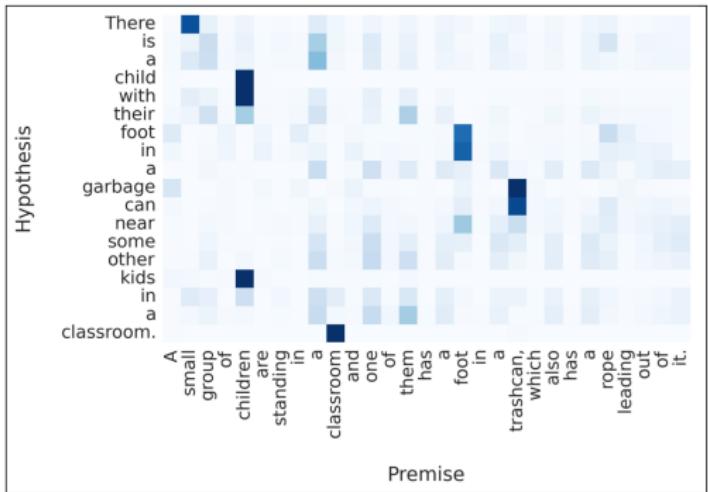
Natural language inference

Rocktäschel et al (2016)

Reasoning about Entailment with Neural Attention



(a)



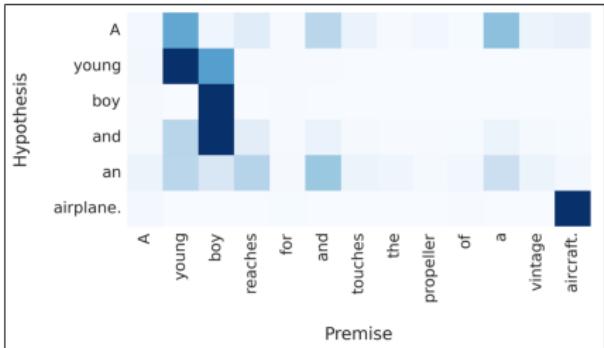
(b)

Attention conditioned on h_1, \dots, h_T : Synonymy, importance

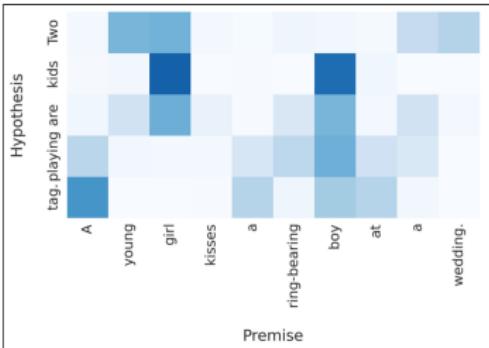
Natural language inference

Rocktäschel et al (2016)

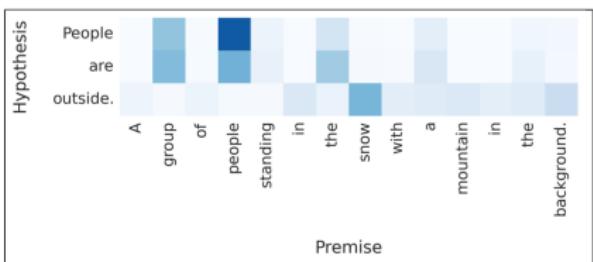
Reasoning about Entailment with Neural Attention



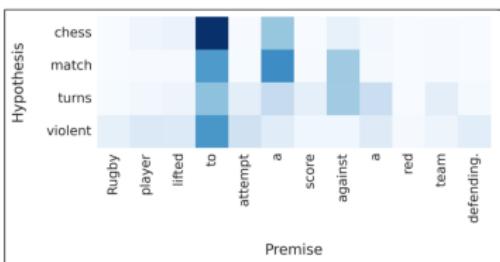
(c)



(d)



(e)



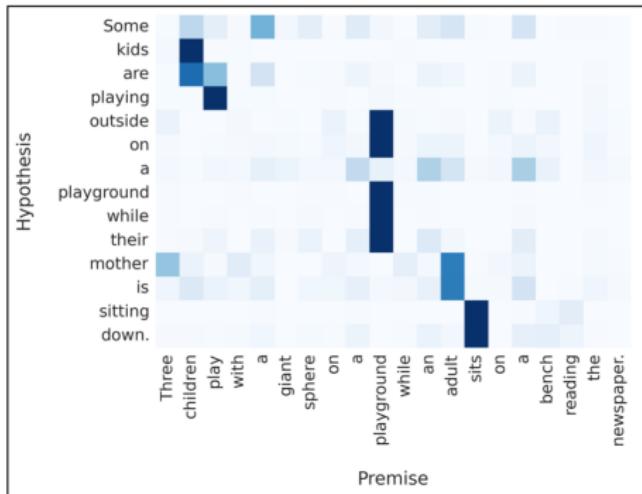
(f)

Attention conditioned on h_1, \dots, h_T : Relatedness

Natural language inference

Rocktäschel et al (2016)

Reasoning about Entailment with Neural Attention



(g)

Attention conditioned on h_1, \dots, h_T : Many:one

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

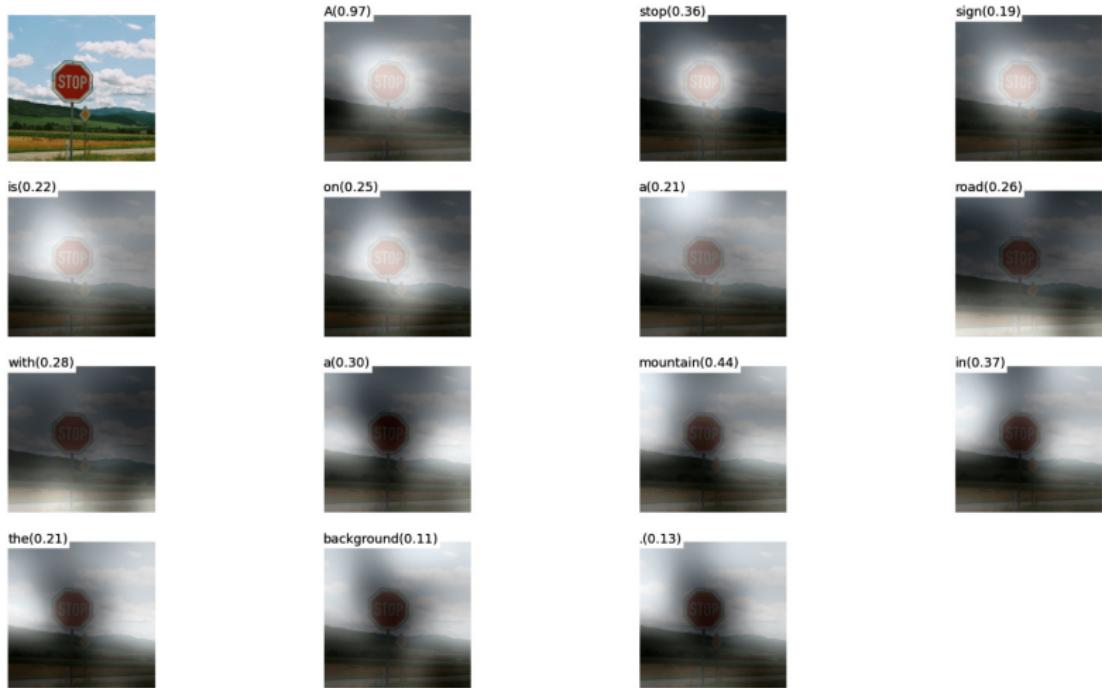
The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

Attention over images

Xu et al (2015)

Show, Attend & Tell: Neural Image Caption Generation with Visual Attention

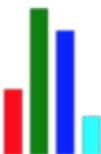
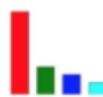


(b) A stop sign is on a road with a mountain in the background.

Attention over videos

Yao et al (2015)

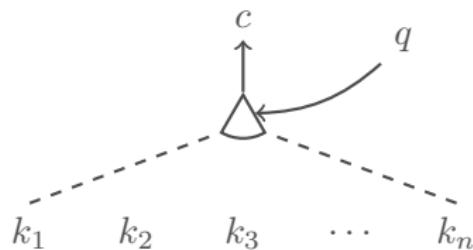
Describing Videos by Exploiting Temporal Structure



+Local+Global: Someone is frying a fish in a pot

Attention variants

$c = \text{ATTENTION}(\text{query } q, \text{ keys } k_1 \dots k_n)$



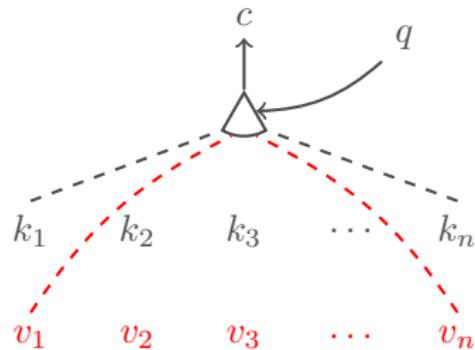
$$\alpha_i = \text{softmax}(\text{score}(q, k_i))$$

$$c = \sum_i \alpha_i, k_i$$

Attention variants

$c = \text{ATTENTION}(\text{query } q, \text{ keys } k_1 \dots k_n, \text{ values } v_1 \dots v_n)$

e.g., memory networks (Weston et al, 2015; Sukhbataar et al, 2015)

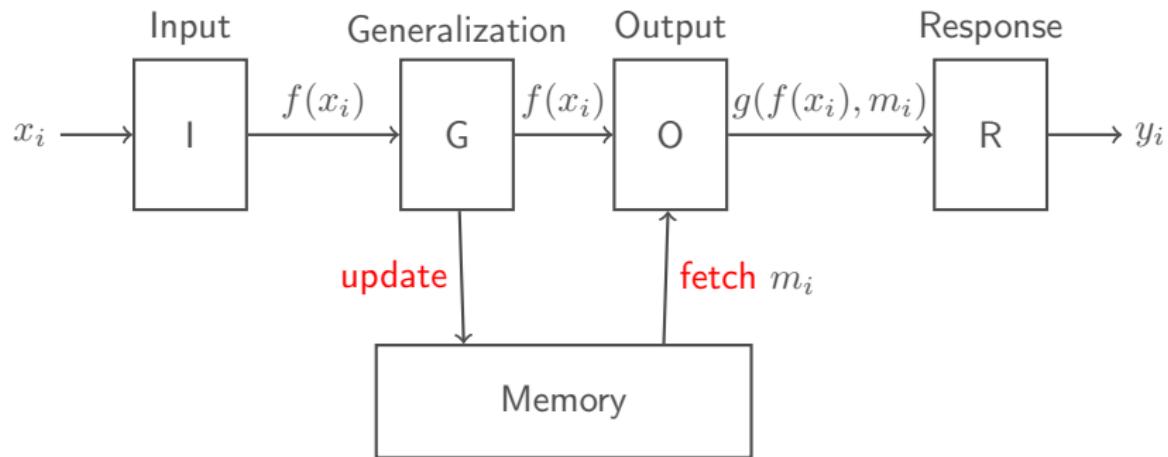


$$\alpha_i = \text{softmax}(\text{score}(q, k_i))$$

$$c = \sum_i \alpha_i, v_i$$

Attention variants

Weston et al (2015)
Memory Networks



MemN2N (Sukhbataar et al, 2015)

- + Soft attention over memories
- + Multiple memory lookups (hops)
- + End-to-end training

Attention scoring functions

- Additive (Bahdanau et al, 2015)

$$\text{score}(q, k) = \mathbf{u}^\top \tanh(\mathbf{W}[q; k])$$

- Multiplicative (Luong et al, 2015)

$$\text{score}(q, k) = q^\top \mathbf{W}k$$

- Scaled dot-product (Vaswani et al, 2017)

$$\text{score}(q, k) = \frac{q^\top k}{\sqrt{d_k}}$$

Attention variants

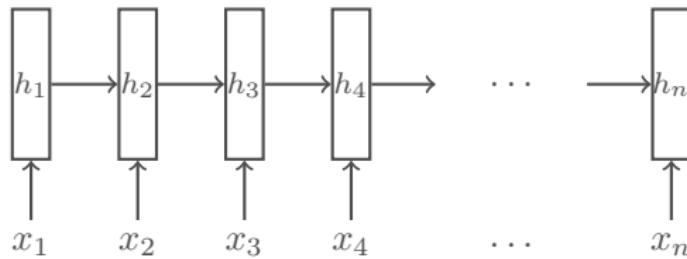
- Stochastic hard attention (Xu et al, 2015)
- Local attention (Luong et al, 2015)
- Monotonic attention (Yu et al, 2016; Raffel et al, 2017)
- Self attention (Cheng et al, 2016; Vaswani et al, 2017)
- Convolutional attention (Allamanis et al, 2016)
- Structured attention (Kim et al, 2017)
- Multi-headed attention (Vaswani et al, 2017)

Transformer

Vaswani et al (2017)

Attention is All You Need

RNN encoder

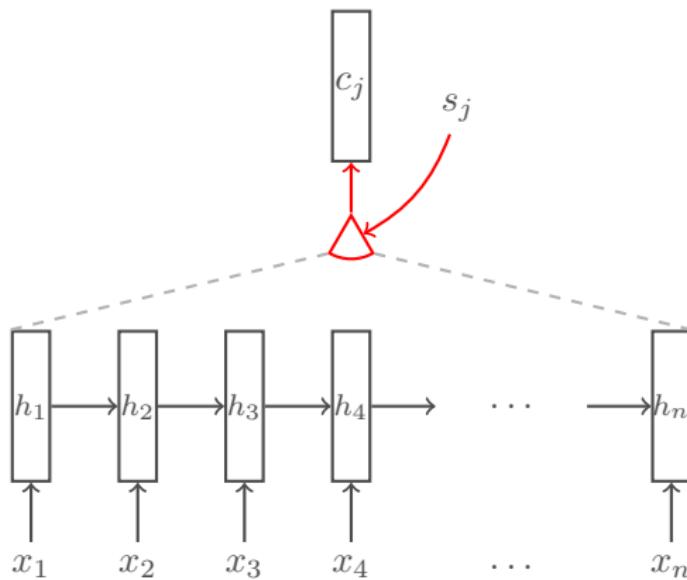


Transformer

Vaswani et al (2017)

Attention is All You Need

RNN encoder with **attention**

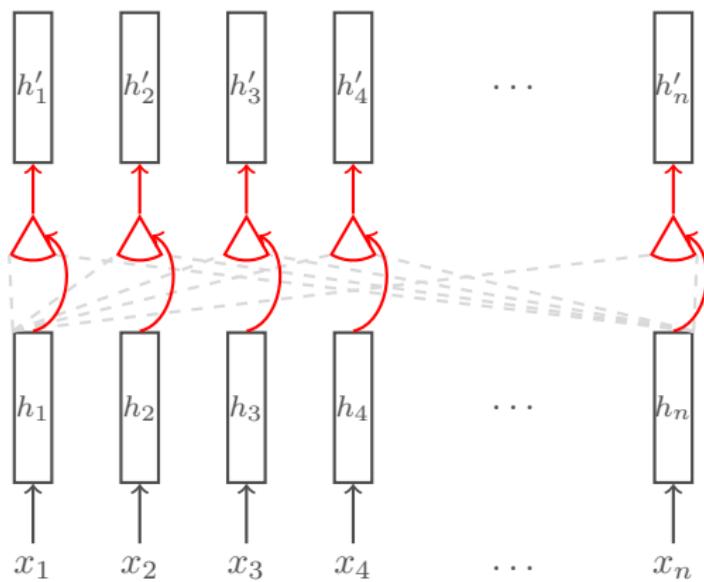


Transformer

Vaswani et al (2017)

Attention is All You Need

Deep encoder with self-attention

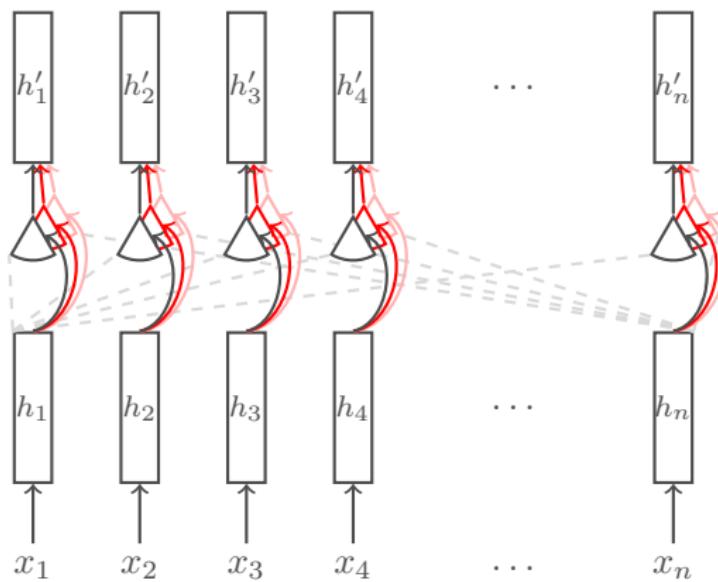


Transformer

Vaswani et al (2017)

Attention is All You Need

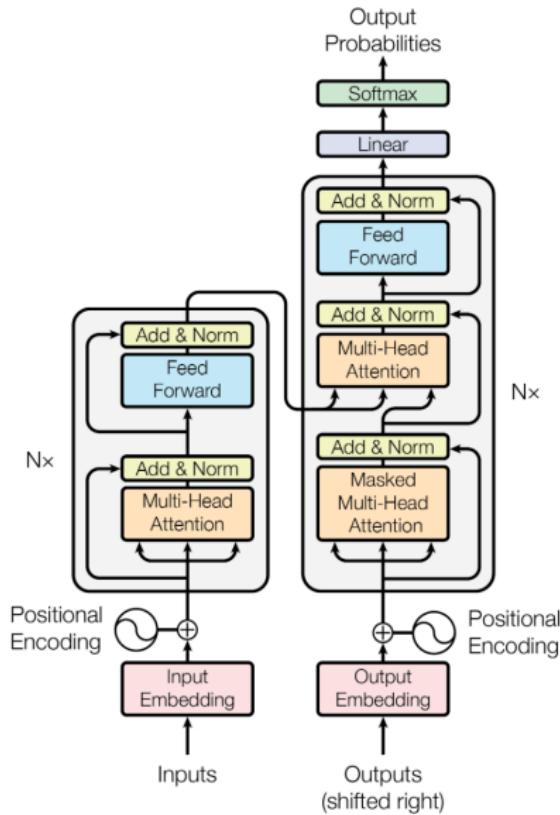
Deep encoder with **multi-headed** self-attention



Transformer

Vaswani et al (2017)

Attention is All You Need



Transformer

Vaswani et al (2017)

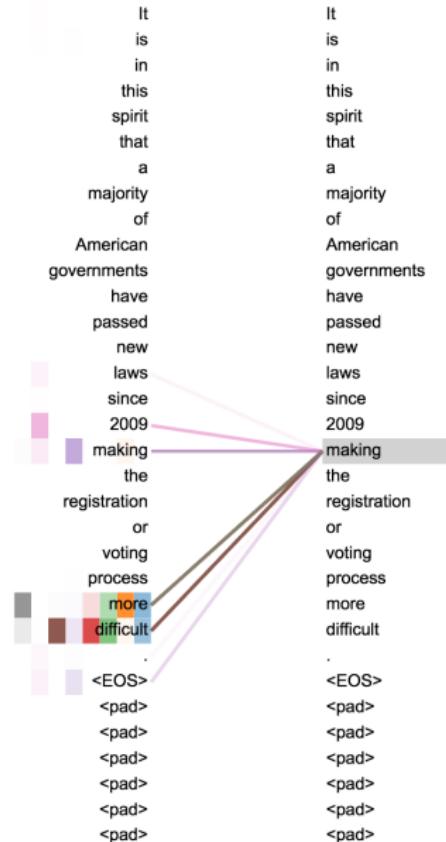
Attention is All You Need

- Self-attention at every layer instead of recurrence
 - + Inference can be parallelized
- No sensitivity to input position
 - Positional embeddings required
 - + Can apply to sets
- Deep architecture (6 layers) with multi-headed attention
 - + Higher layers appear to learn linguistic structure
- Scaled dot-product attention with masking
 - + Avoids bias in simple dot-product attention
 - + Fewer parameters needed for rich model
- Improved runtime and performance on translation, parsing, etc

Transformer

Vaswani et al (2017)

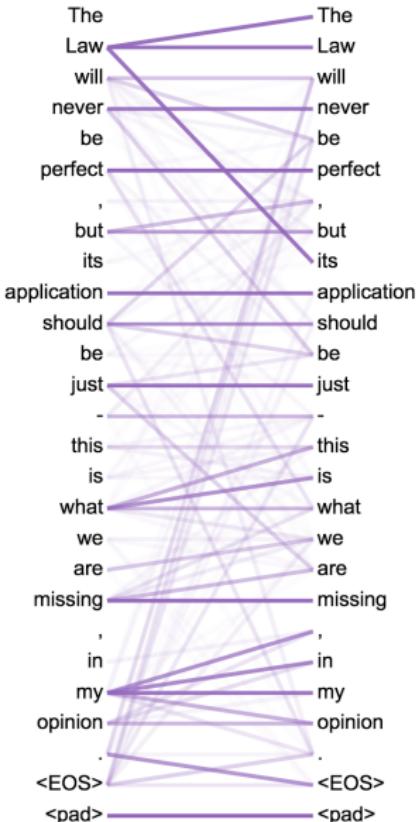
Attention is All You Need



Transformer

Vaswani et al (2017)

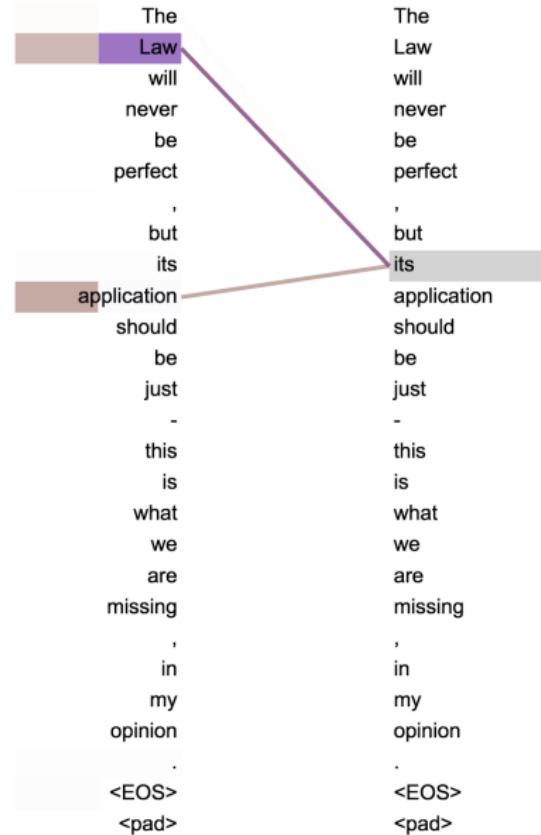
Attention is All You Need



Transformer

Vaswani et al (2017)

Attention is All You Need



Transformer

Vaswani et al (2017)

Attention is All You Need



Transformer

Vaswani et al (2017)

Attention is All You Need



Large vocabularies

Sequence-to-sequence models can typically scale to 30K-50K words

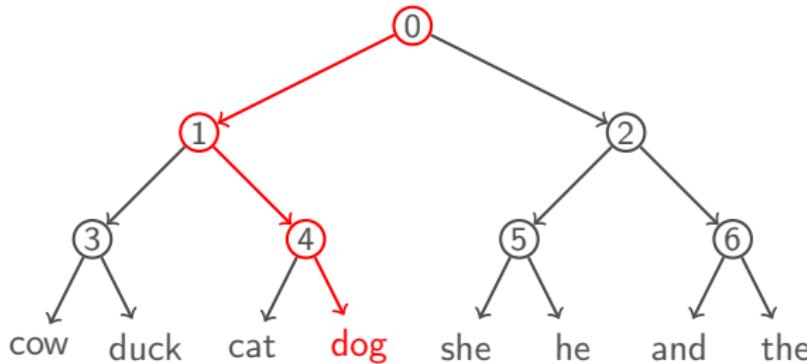
But real-world applications need at least 500K-1M words

Large vocabularies

Alternative 1: Hierarchical softmax

- Predict path in binary tree representation of output layer
- Reduces to $\log_2(V)$ binary decisions

$$p(w_t = \text{"dog"} | \dots) = (1 - \sigma(U_0 h_t)) \times \sigma(U_1 h_t) \times \sigma(U_4 h_t)$$



Large vocabularies

Jean et al (2015)

On Using Very Large Target Vocabulary for Neural Machine Translation

Alternative 2: Importance sampling

- Expensive to compute the softmax normalization term over V

$$p(y_i = w_j | y_{<i}, x) = \frac{\exp(W_j^\top f(s_i, y_{i-1}, c_i))}{\sum_{k=1}^{|V|} \exp(W_k^\top f(s_i, y_{i-1}, c_i))}$$

- Use a small subset of the target vocabulary for each update
- Approximate expectation over gradient of loss with fewer samples
- Partition the training corpus and maintain local vocabularies in each partition to use GPUs efficiently

Large vocabularies

Sennrich et al (2016)

Neural machine translation of rare words with subword units

Alternative 3: Subword units

- Reduce vocabulary by replacing infrequent words with sub-words

Jet makers feud over seat width with big orders at stake



_ J et _ makers _ fe ud _ over _ seat _ width _ with _ big _ orders _ at _ stake

- Code for byte-pair encoding:

<https://github.com/rsennrich/subword-nmt>

Copying mechanism

Gu et al (2016)

Incorporating Copying Mechanism in Sequence-to-Sequence Learning

In monolingual tasks, copy rare words directly from the input

- Generation via standard attention-based decoder

$$\psi_g(y_i = w_j) = W_j^\top f(s_i, y_{i-1}, c_i) \quad w_j \in V$$

- Copying via a non-linear projection of input hidden states

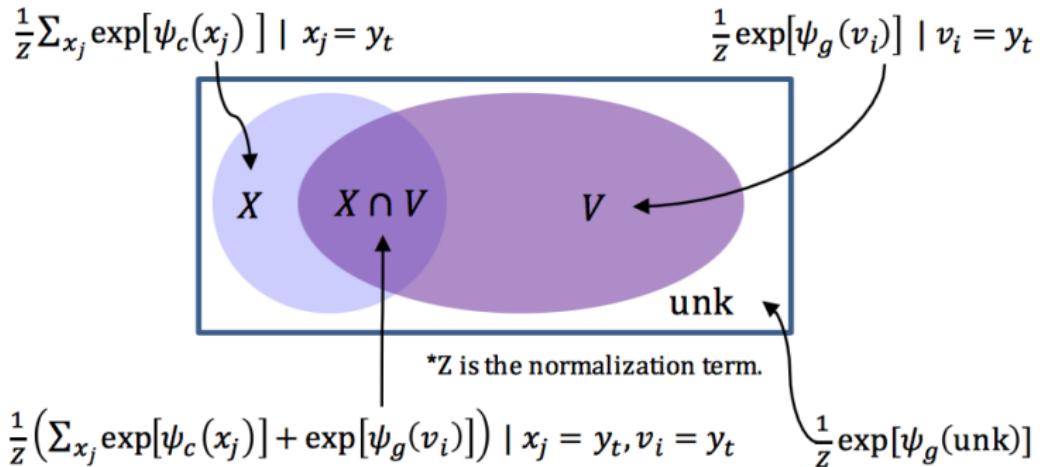
$$\psi_c(y_i = x_j) = \tanh(h_j^\top U) f(s_i, y_{i-1}, c_i) \quad x_j \in X$$

- Both modes compete via the softmax

$$p(y_i = w_j | y_{<i}, x) = \frac{1}{Z} \left(\exp(\psi_g(w_j)) + \sum_{k: x_k = w_j} \exp(\psi_c(x_k)) \right)$$

Copying mechanism

Incorporating Copying Mechanism in Sequence-to-Sequence Learning

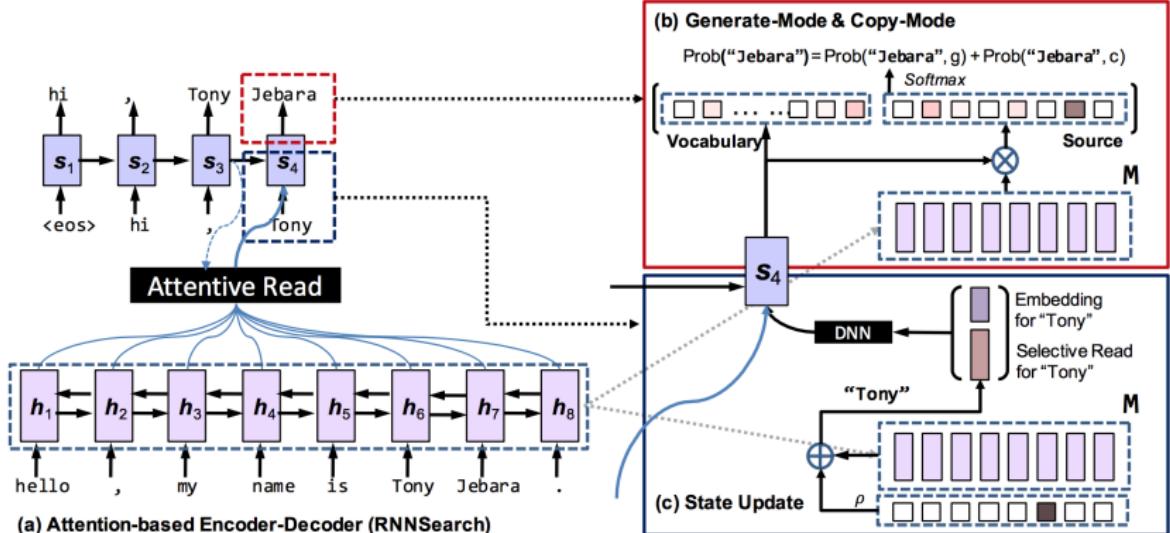


Decoding probability $p(y_t | \dots)$

Copying mechanism

Gu et al (2016)

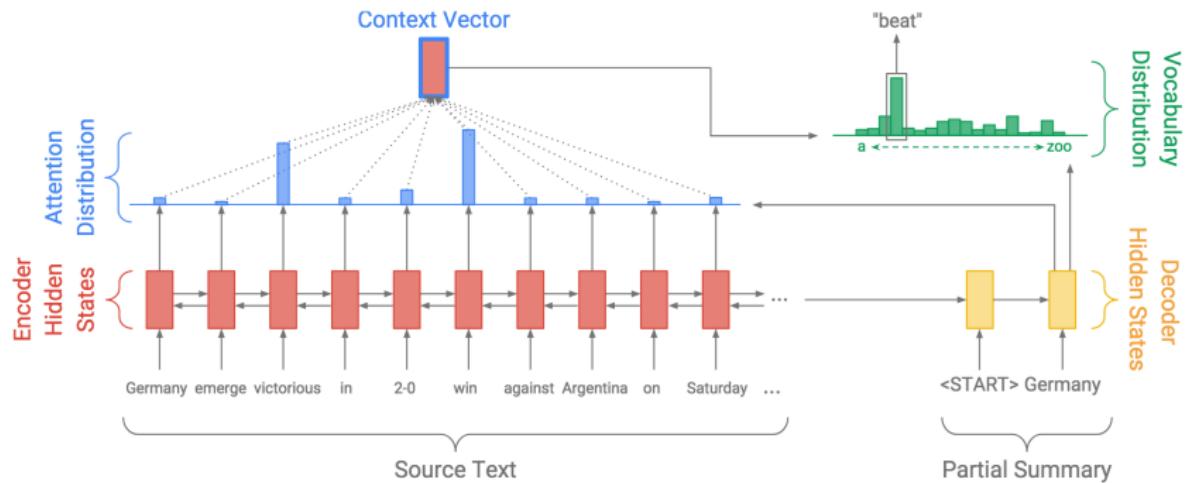
Incorporating Copying Mechanism in Sequence-to-Sequence Learning



Copying mechanism

See et al (2017)

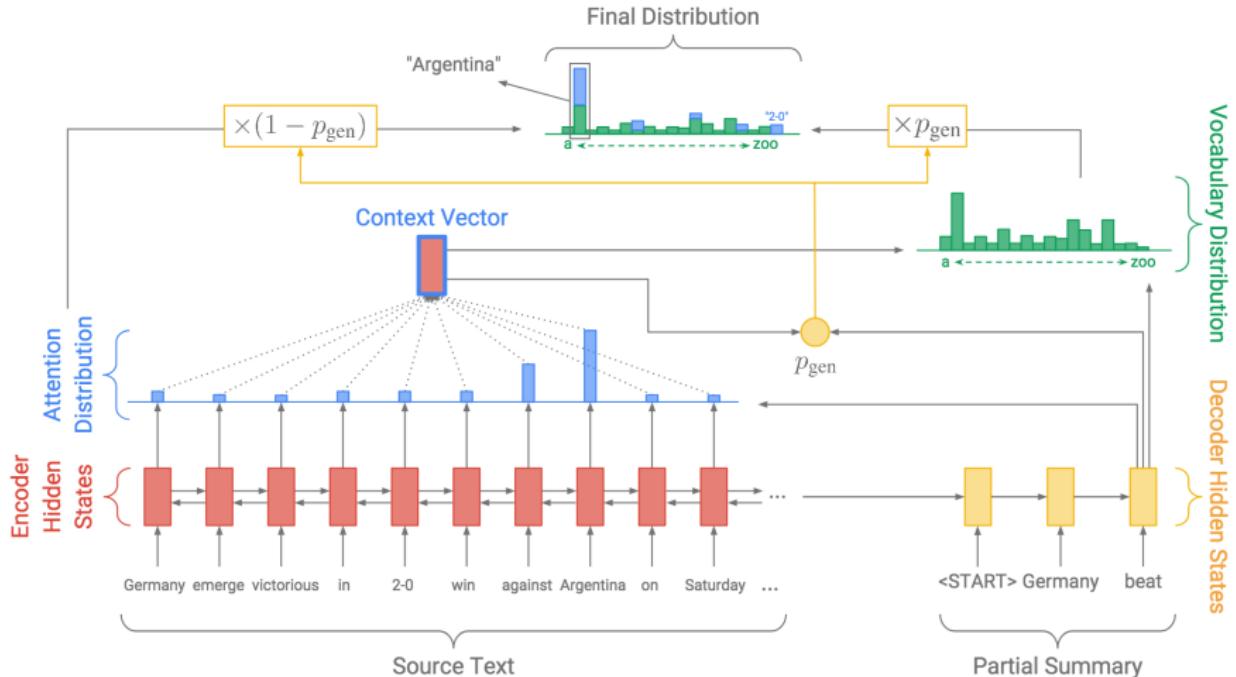
Get to the Point: Summarization with Pointer Generator Networks



Attention for common words

Copying mechanism

Get to the Point: Summarization with Pointer Generator Networks



Copying from input for rarer words

Autoencoders

Given input x , learn an encoding z that can be decoded to reconstruct x

For sequence input x_1, \dots, x_n , can use standard MT models

- Is attention viable?
- + Useful for pre-training text classifiers (Dai et al, 2015)

Denoising autoencoders

Hill et al (2016)

Learning Distributed Representations of Sentences from Unlabeled Data

Given **noisy** input \tilde{x} , learn an encoding z that can be decoded to reconstruct x

Noise functions for sequences

- Drop words with probability p_{drop}
- Swap words with probability p_{swap}

- + Useful as features for a linear classifier
- + Can learn sentence representations without ordered docs

Dataset	Sentence 1	Sentence 2	/5
STS 2014	News The problem is simpler than that.	<i>Mexico wishes to avoid more violence.</i> <i>The problem is simple.</i>	4 3.8
	Forum A social set or clique of friends.	<i>An unofficial association of people or groups.</i>	3.6
	WordNet Taking Aim #Stopgunviolence #Congress #NRA	<i>Obama, Gun Policy and the N.R.A.</i>	1.6
	Twitter A woman riding a brown horse.	<i>A young girl riding a brown horse.</i>	4.4
	Images Iranians Vote in Presidential Election.	<i>Keita Wins Mali Presidential Election.</i>	0.4
	Headlines SICK (test+train) A lone biker is jumping in the air.	<i>A man is jumping into a full pool.</i>	1.7

Variational autoencoders (VAEs)

Kingma & Welling (2015)

Auto-encoding Variational Bayes

Autoencoders often don't generalize well to new data, noisy representations

Approximate the posterior $p(z|x)$ with variational inference

- Encoder: induce $q(z|x)$ with parameters θ
- Decoder: sample z and reconstruct x with parameters ϕ
- Loss:

$$\ell_i = -\mathbb{E}_{z \sim q_\theta(z|x_i)} \log p_\phi(x_i|z) + \text{KL}(q_\theta(z|x_i)||p(z))$$

Estimate gradients using *reparameterization trick* for Gaussians

$$\mathcal{N}(\mu, \sigma^2) = \mathcal{N}(0, 1)\sigma + \mu$$

- + Can interpolate smoothly in representation space, e.g.,
<https://giphy.com/gifs/26ufgj5LH3YK01Zlu>