

Introduction To Seq2Seq

10/06/2018

(adapted from ICML 2017 tutorial)

Outline

- An Overview of seq2seq (~25 minutes)
 - Applications
- Attention (~20 minutes)
 - Applications
- Break (15 mins)
- Loss functions (~15 minutes)
- Advanced Topics
 - Autoregressive Models (~5 mins)
 - Online Seq2Seq (~10 mins)
- New Architectures (~10 mins)
- Questions, discussion, etc. (15 minutes)

Seq2Seq Overview

Sequences!

- Words, Letters

50 years ago, the fathers of artificial intelligence convinced everybody that logic was the key to intelligence. Somehow we had to get computers to do logical reasoning. The alternative approach, which they thought was crazy, was to forget logic and try and understand how networks of brain cells learn things. Curiously, two people who rejected the logic based approach to AI were Turing and Von Neumann. If either of them had lived I think things would have turned out differently... now neural networks are everywhere and the crazy approach is winning.

Geoff Hinton

- Speech



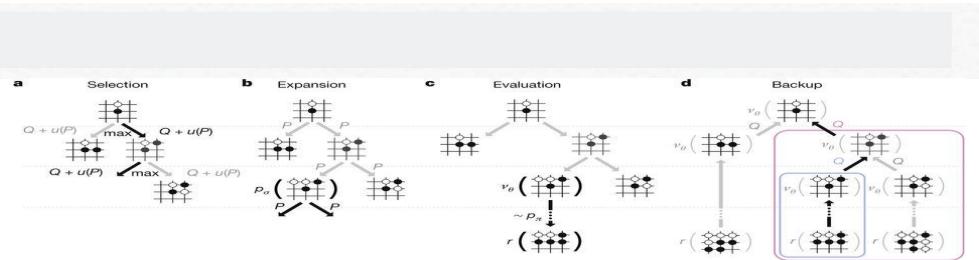
- Images, Videos



- Programs

```
while (*d++ = *s++);
```

- Sequential Decision Making (RL)



Classical Models for Sequence Prediction

- Sequence prediction was classically handled as a structured prediction tasks
 - Most were built on conditional independence assumptions
 - Other such as DAGGER were based on supervisory signals and auxiliary information

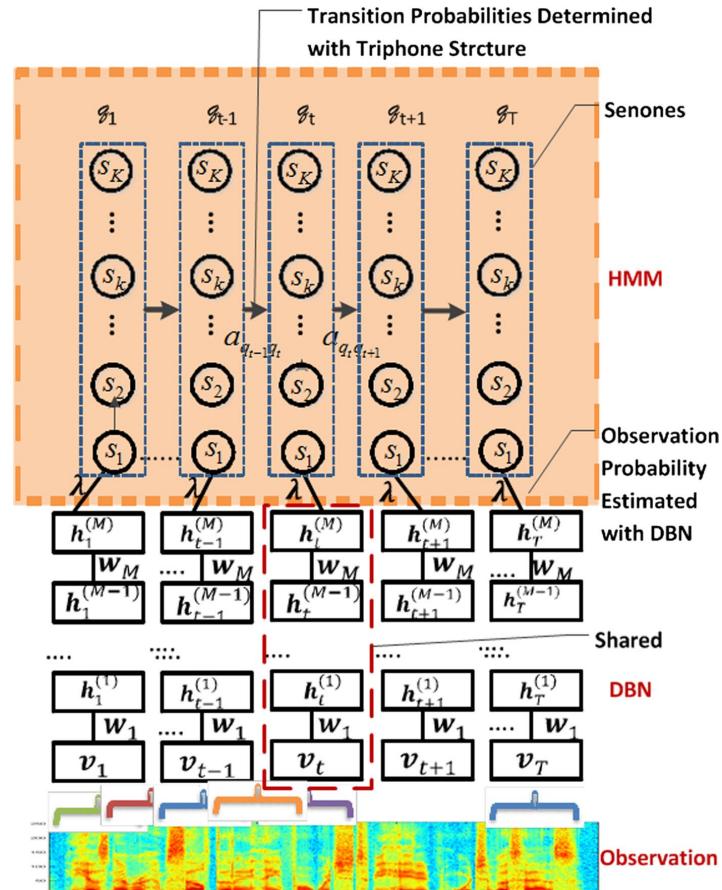


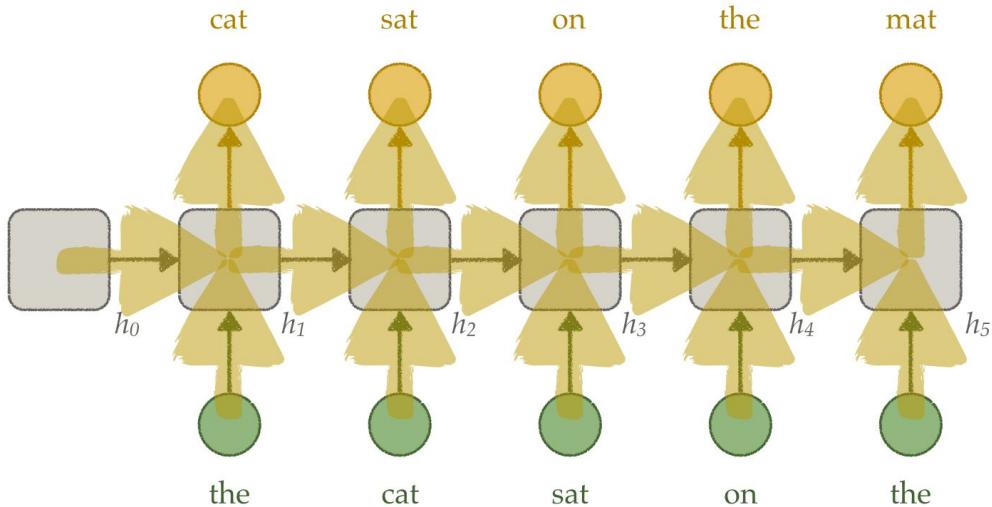
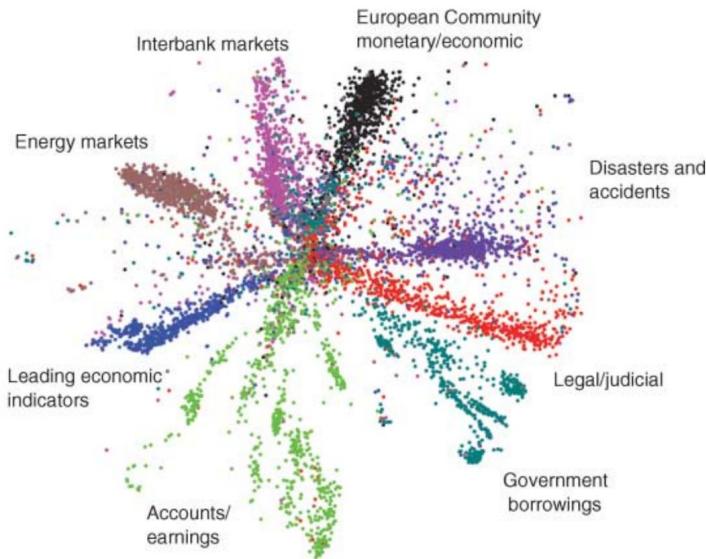
Figure credit: Li Deng

Two Key Ingredients

Neural Embeddings



Recurrent Language Models

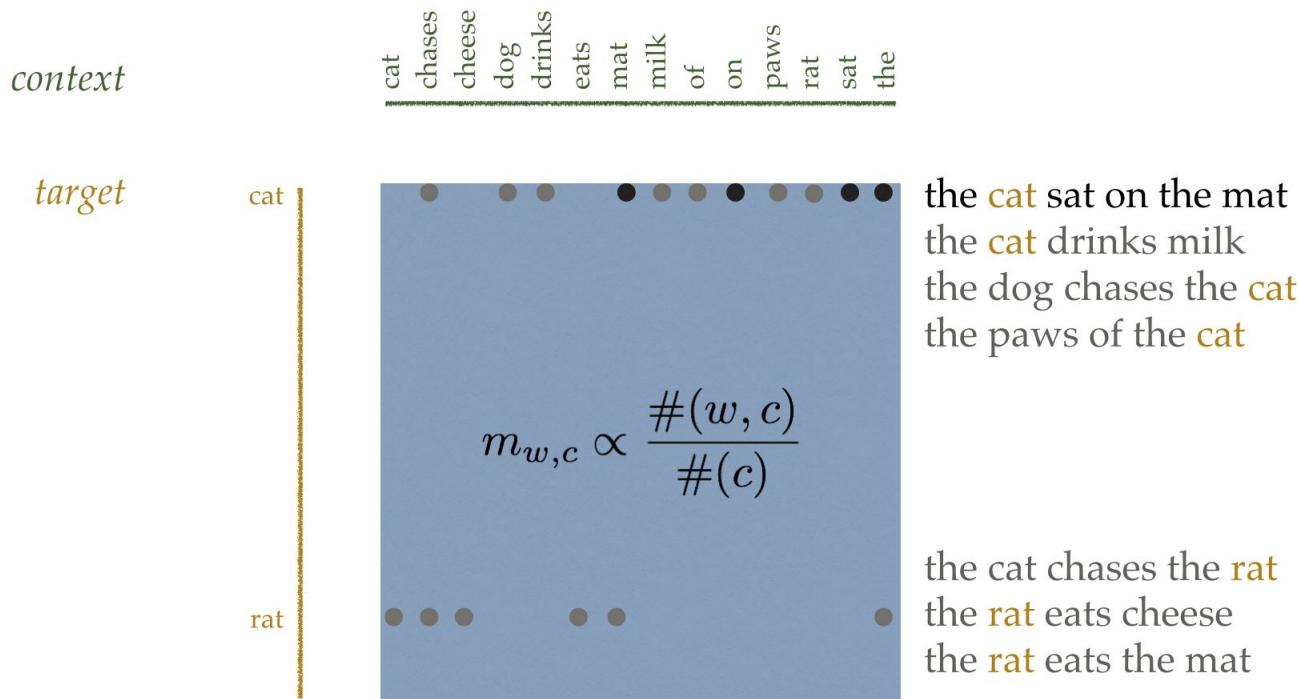


1. Hinton, G., Salakhutdinov, R. "Reducing the Dimensionality of Data with Neural Networks." *Science* (2006)
2. Mikolov, T., et al. "Recurrent neural network based language model." *Interspeech* (2010)

Language Models

<i>context</i>						<i>target</i>	$P(w_t w_{t-1}, w_{t-2}, \dots w_{t-5})$
the	cat	sat	on	the	mat		0.15
w_{t-5}	w_{t-4}	w_{t-3}	w_{t-2}	w_{t-1}	w_t		
the	cat	sat	on	the	rug		0.12
the	cat	sat	on	the	hat		0.09
the	cat	sat	on	the	dog		0.01
the	cat	sat	on	the	the		0
the	cat	sat	on	the	sat		0
the	cat	sat	on	the	robot		?
the	cat	sat	on	the	printer		?

n-grams



n-grams

$$P(w_1, w_2, \dots, w_{T-1}, w_T) \approx \prod_{t=1}^T P(w_t | w_{t-1}, \dots, w_{t-n+1})$$

the	cat	sat	on	the	mat	$P(w_1)$
the	cat	sat	on	the	mat	$P(w_2 w_1)$
the	cat	sat	on	the	mat	$P(w_3 w_2, w_1)$
the	cat	sat	on	the	mat	$P(w_4 w_3, w_2)$
the	cat	sat	on	the	mat	$P(w_5 w_4, w_3)$
the	cat	sat	on	the	mat	$P(w_6 w_5, w_4)$

The Chain Rule

$$P(w_1, w_2, \dots, w_{T-1}, w_T) = \prod_{t=1}^T P(w_t | w_{t-1}, w_{t-2}, \dots, w_1)$$

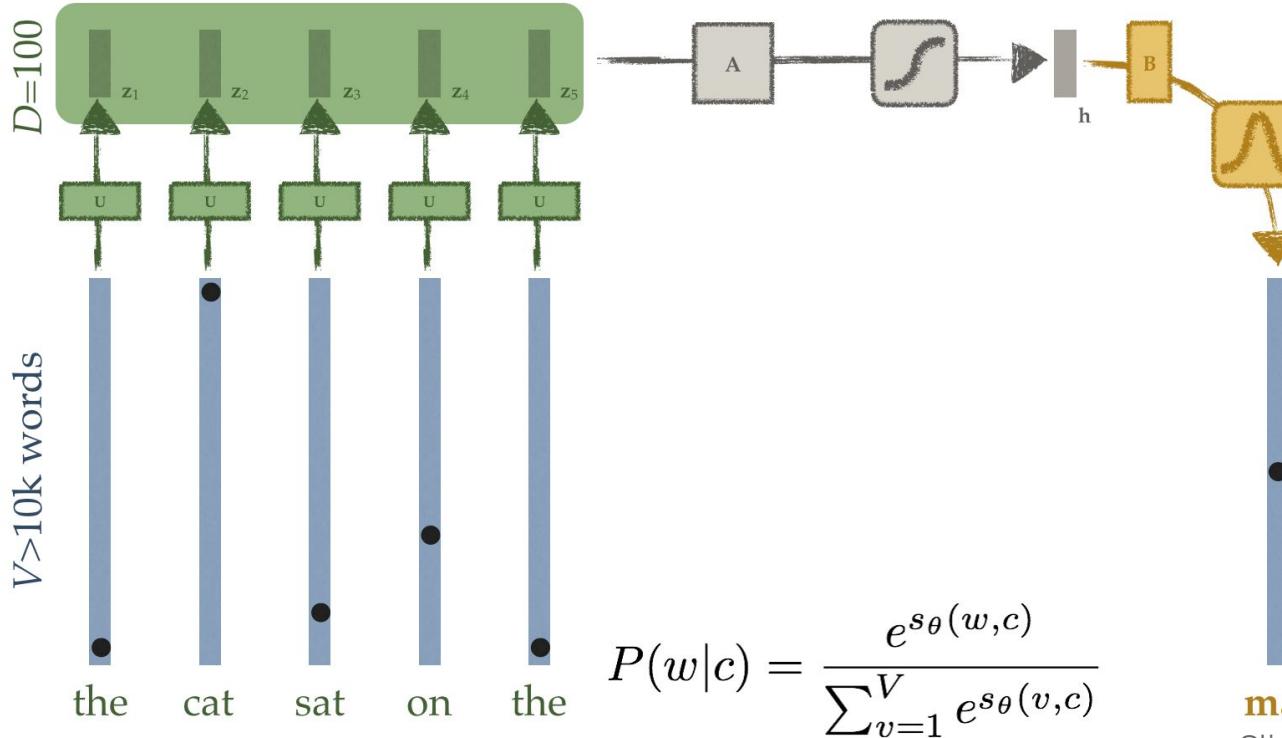
the	cat	sat	on	the	mat	$P(w_1)$
the	cat	sat	on	the	mat	$P(w_2 w_1)$
the	cat	sat	on	the	mat	$P(w_3 w_2, w_1)$
the	cat	sat	on	the	mat	$P(w_4 w_3, w_2, w_1)$
the	cat	sat	on	the	mat	$P(w_5 w_4, w_3, w_2, w_1)$
the	cat	sat	on	the	mat	$P(w_6 w_5, w_4, w_3, w_2, w_1)$

A Key Insight: vectorizing context

Bengio, Y. et al., "A Neural Probabilistic Language Model", *JMLR* (2001, 2003)

Mnih, A., Hinton, G., "Three new graphical models for statistical language modeling", *ICML* 2007

$$p(w_t | w_1, \dots, w_{t-1}) = p_\theta(w_t | f_\theta(w_1, \dots, w_{t-1}))$$

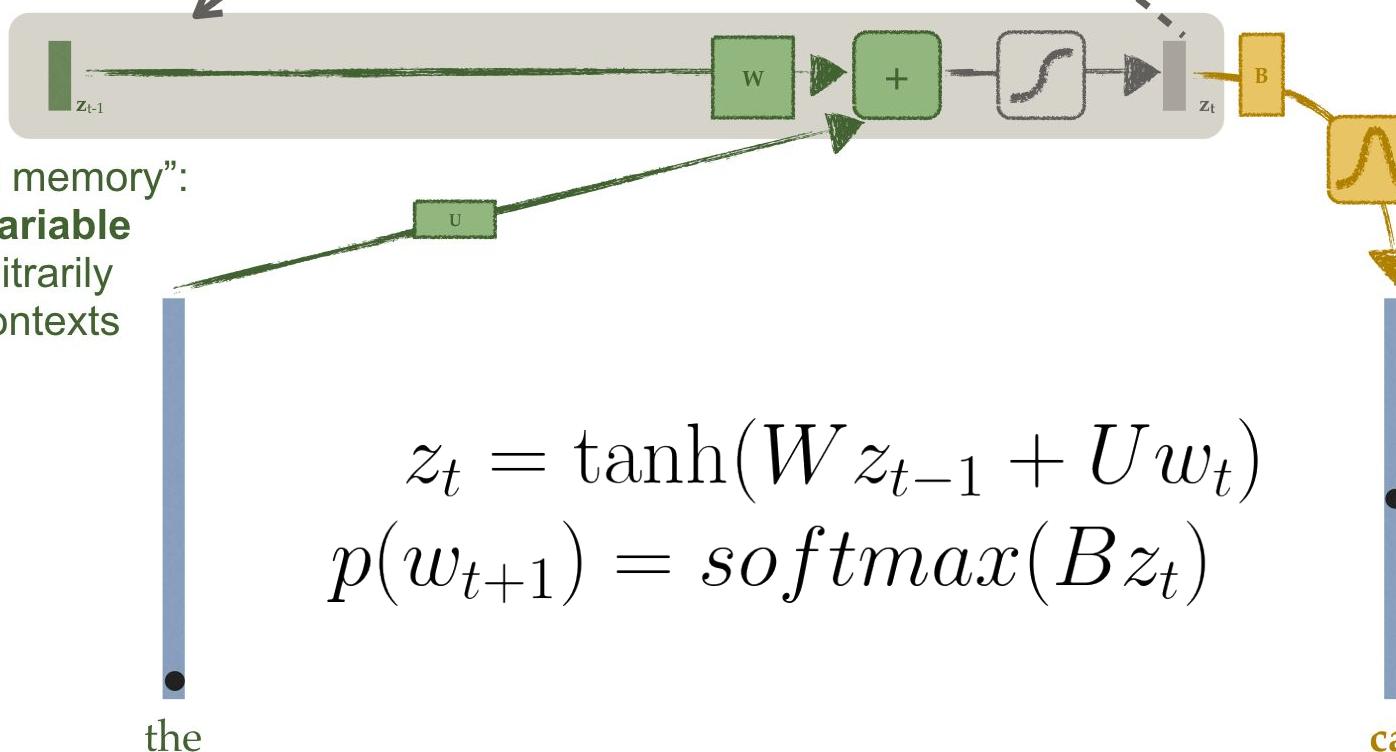


mat

Slide Credit: Piotr Mirowski

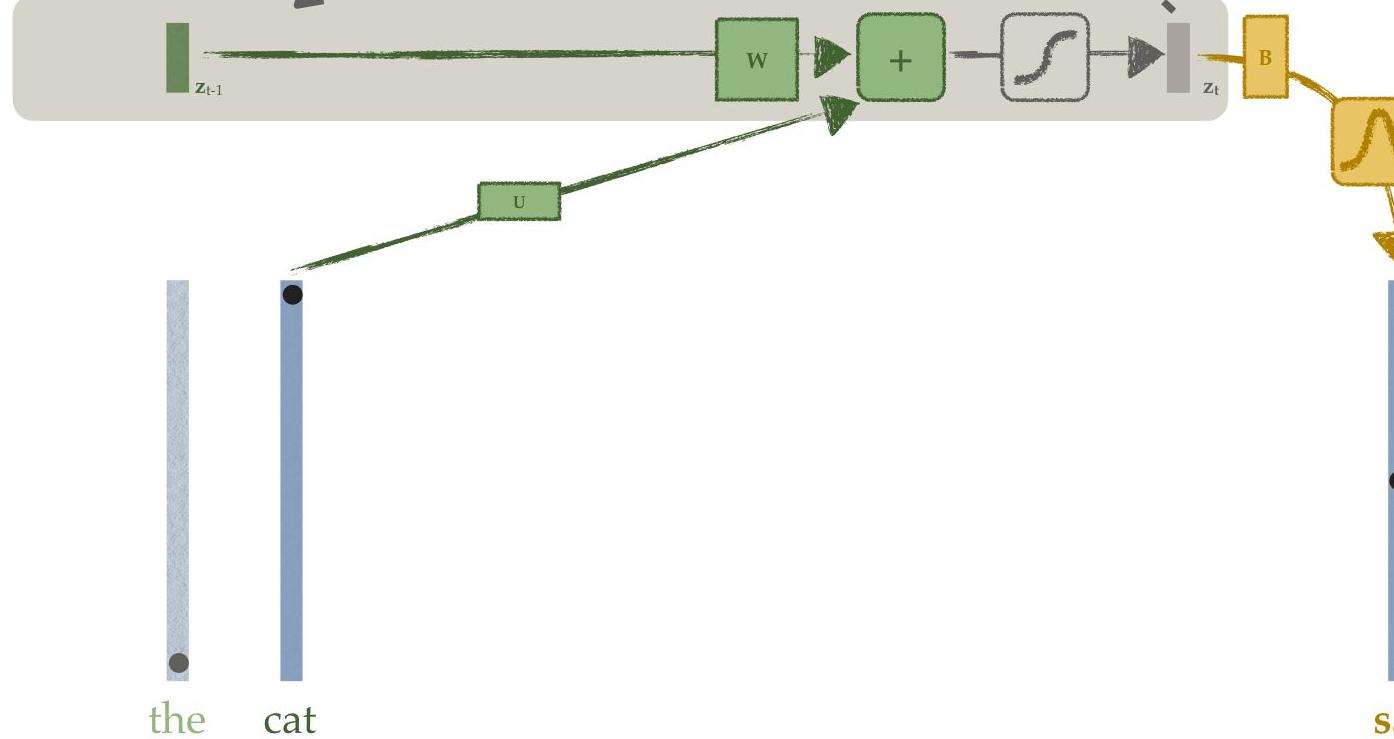
Recurrent Neural Network Language Models

[Jeffrey L Elman (1991) "Distributed representations, simple recurrent networks and grammatical structure", *Machine Learning*; Tomas Mikolov et al. (2010) "Recurrent neural network based language model", *INTERSPEECH*]

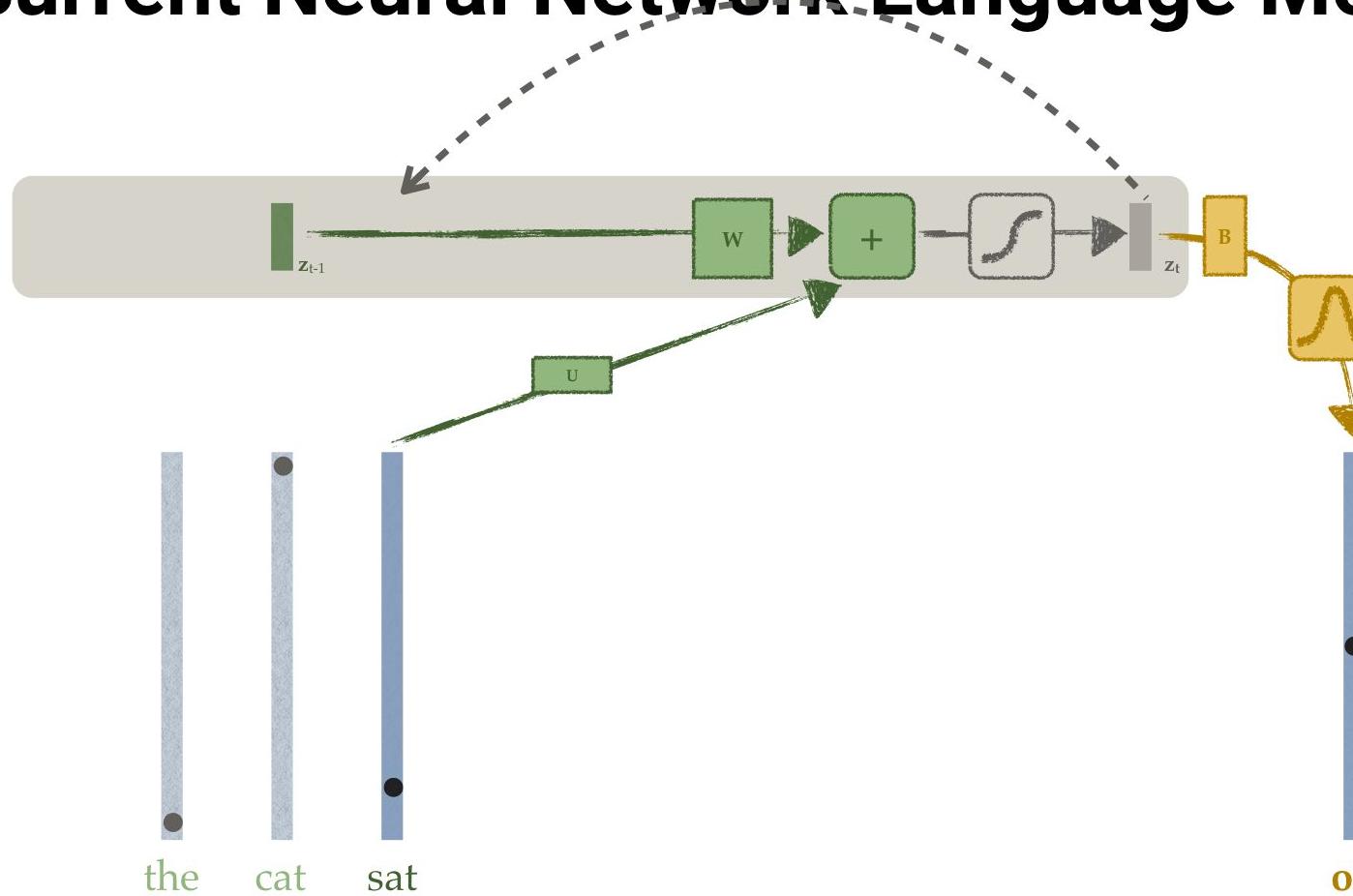


$$\begin{aligned} z_t &= \tanh(Wz_{t-1} + Uw_t) \\ p(w_{t+1}) &= softmax(Bz_t) \end{aligned}$$

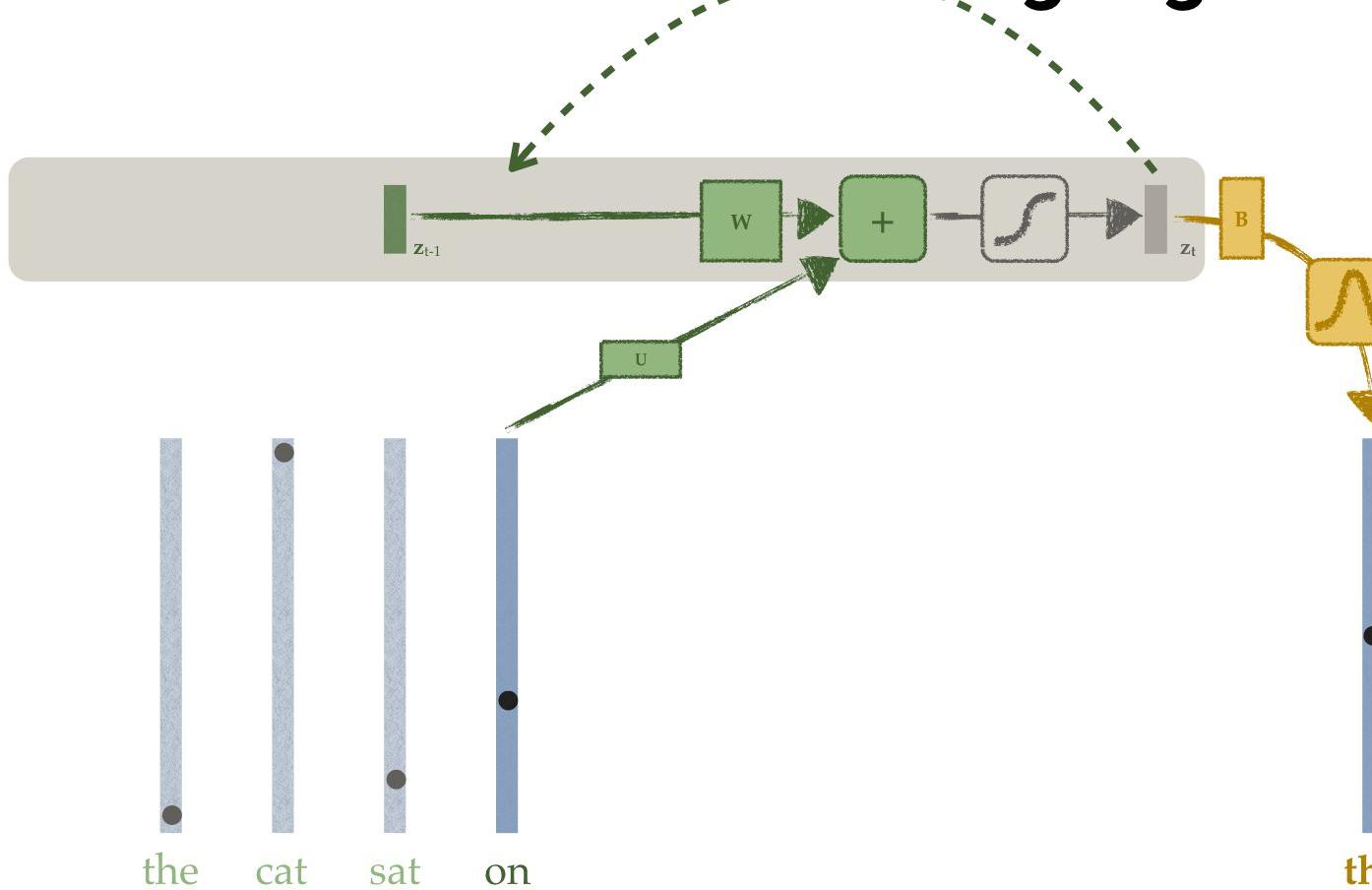
Recurrent Neural Network Language Models



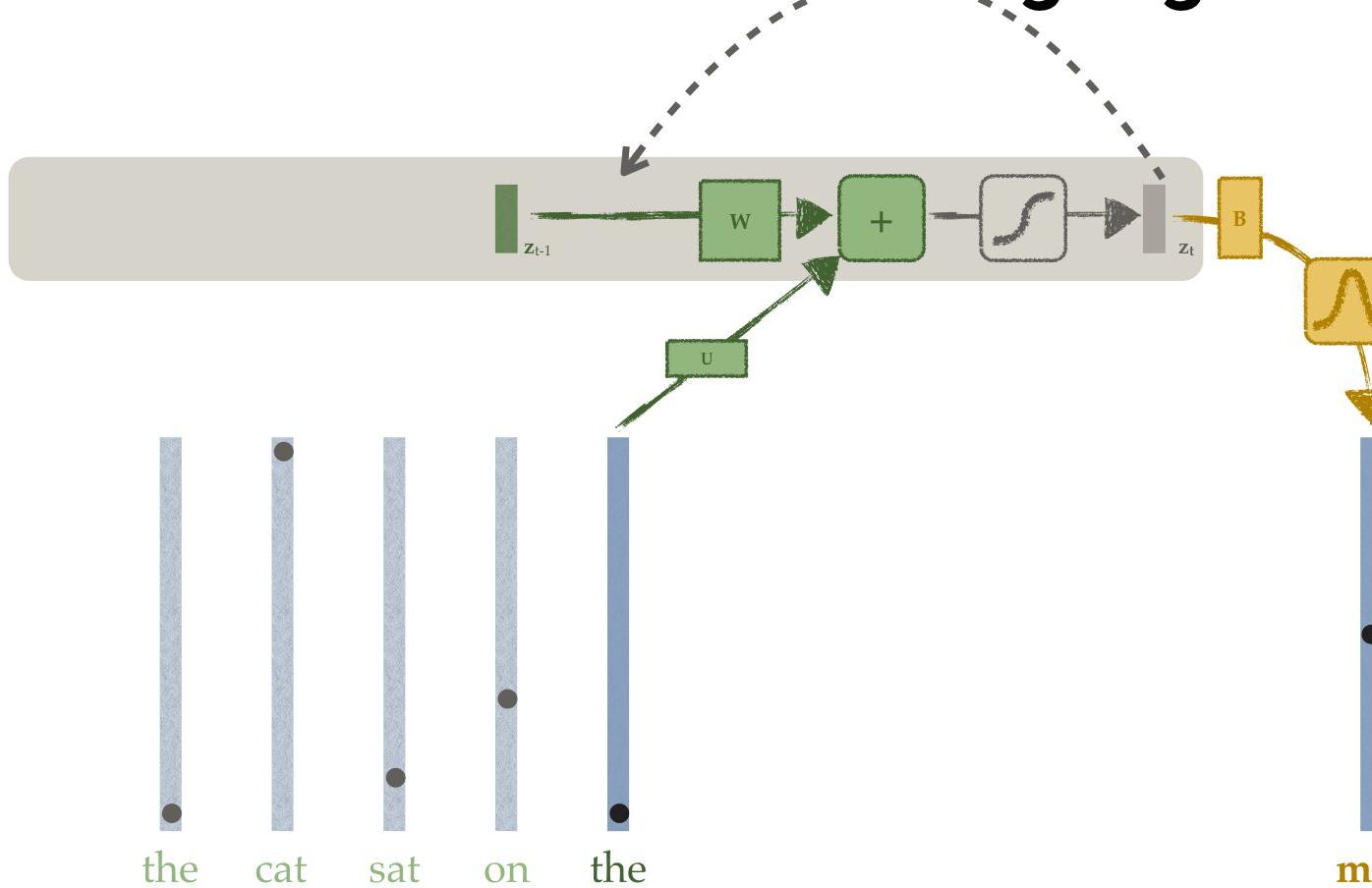
Recurrent Neural Network Language Models



Recurrent Neural Network Language Models



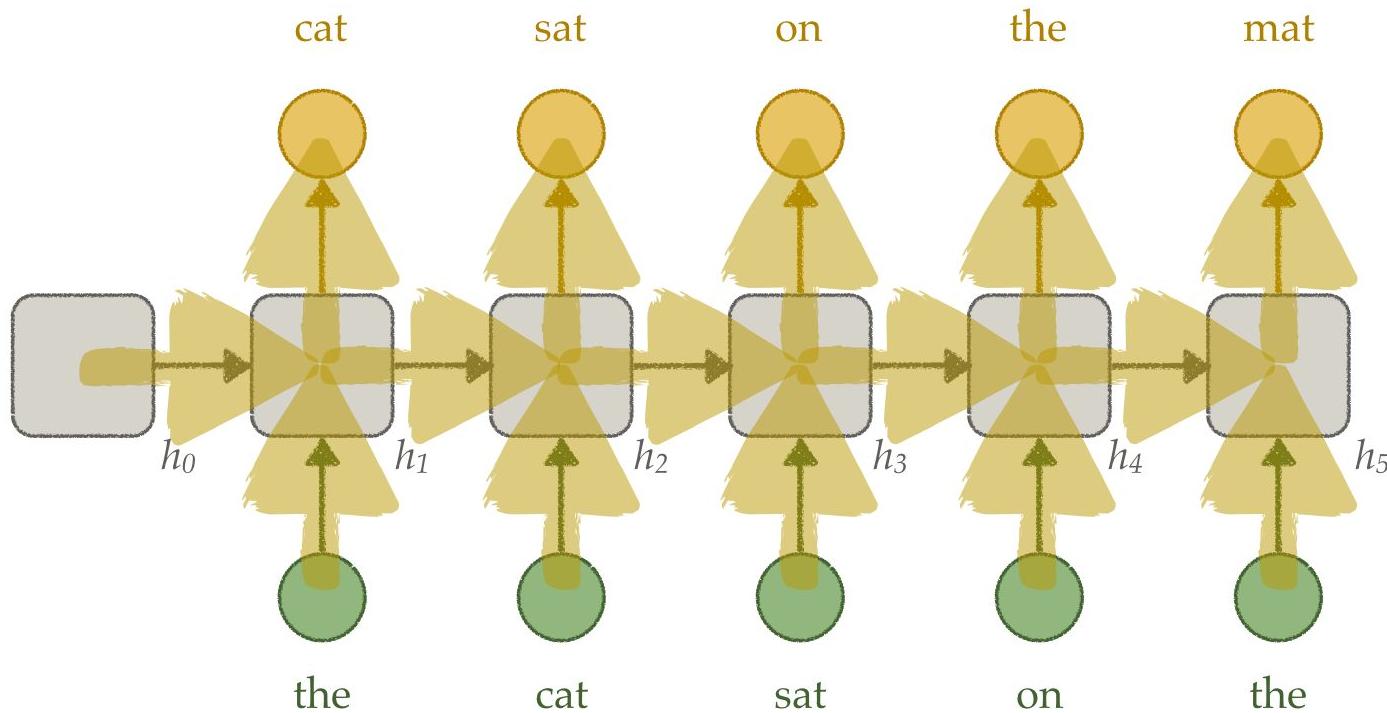
Recurrent Neural Network Language Models



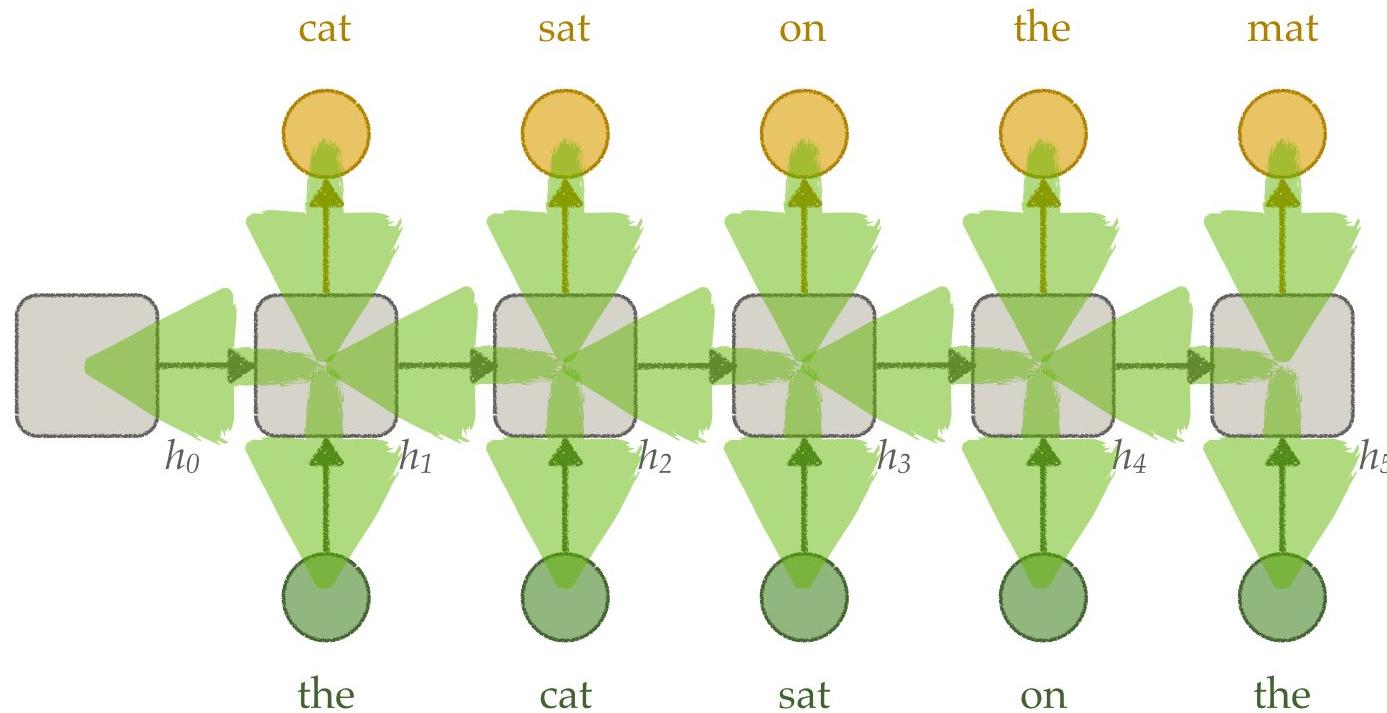
What do we Optimize?

$$\theta^* = \arg \max_{\theta} E_{w \sim data} \log P_{\theta}(w_1, \dots, w_T)$$

Recurrent Neural Network Language Models



Recurrent Neural Network Language Models



Mini-Turing Test

Some of the obese people lived five to eight years longer than others. (from language model)

Abu Dhabi is going ahead to build solar city and no pollution city. (from data distribution)

Or someone who exposes exactly the truth while lying.

VIERA , FLA . -- Sometimes, Rick Eckstein dreams about baseball swings.

For decades, the quintessentially New York city has elevated its streets to the status of an icon.

The lawsuit was captioned as United States ex rel.

Seq2Seq

Joint Language and Translation Modeling with Recurrent Neural Networks

Michael Auli, Michel Galley, Chris Quirk, Geoffrey Zweig
Microsoft Research
Redmond, WA, USA
{michael.auli,mgalley,chrisq,gzweig}@microsoft.com

Abstract

We present a joint language and translation model based on a recurrent neural network which predicts target words based on an unbounded history of both source and target words. The weaker independence assumption of this model results in a more larger space of models compared to a feed-forward-based language or translation models. We tackle this issue with a new lattice-based training algorithm that estimate the dependencies between words. Our joint model builds on a well known recurrent neural language model (Mikolov, 2012) augmented by a layer of additional inputs for the source language. We show comparable accuracy compared to the traditional channeled model features. Our best result improve the output of a system trained on WMT 2012 French-English dataset by up to 1.5 BLEU, and by 1.1 BLEU on average across several test sets.

1 Introduction

Recently, several feed-forward neural network-based language and translation models have achieved impressive accuracy improvements on statistical machine translation tasks (Almeida et al., 2011; Le et al., 2012b; Schwenk et al., 2012). In this paper we focus on recurrent neural network architectures, which have recently advanced the state of the art in language modeling (Mikolov et al., 2010; Mikolov et al., 2012) and machine translation (MT) using multi-task feed-forward based networks in both perplexity and word error rate in speech recognition (Ariyoshi et al., 2012; Sundermeyer et al., 2012). The major attraction of recurrent architectures is their potential to capture long-span dependencies since

predictions are based on an *unbounded history* of previous words. This is in contrast to feed-forward networks as well as conventional n-gram models, both of which are limited to fixed-length contexts. Based on the success of this model architecture, we base our joint language and translation model on an extension of the recurrent neural network language model (Mikolov and Zweig, 2012) that introduces a layer of additional inputs (S2).

Most previous work on neural networks for speech recognition or machine translation used a re-scoring setup based on n-best lists (Ariyoshi et al., 2012; Molko et al., 2012) for evaluation, though some explore the algorithmic engineering challenges of direct decoder-degeneration.¹ Instead, we exploit *lattices*, which offer a much richer representation of the decoder output, since they compactly encode an exponential number of translation hypotheses in polynomial space. In contrast, n-best lists are typically very redundant, representing only a few combinations of words in the sequence. This is a major challenge in lattice rescoring with a recurrent neural network model is the effect of the unbounded history on search since the usual dynamic programming assumptions which are exploited for efficiency do not hold up anymore. We apply a novel algorithm to the task of rescoring with an unbounded language model as empirically demonstrated its effectiveness (S3).

The main contribution is learning to jointly improve improvements with the recurrent neural network language model over a competitive n-gram baseline across several language pairs. We even observe consistent gains when pairing the model with a large n-gram model trained on up to 575 times more training data. The joint model is trained on a state-of-the-art system when rescoring a n-best lists of translations.

1 Introduction

In most statistical approaches to machine translation the basic units of translation are phrases that are composed of one or more words. A crucial component of translation systems are models that estimate translation probabilities for pairs of phrases, one phrase being from the source language and the other from the target language. Some models compare phrase pairs and others compare as distinct forms the surface forms of the phrases are distinct. Although distinct phrase pairs often share significant simili-

¹One notable exception is Li et al. (2012) who propose reading lattices with a feed-forward network-based model.

Recurrent Continuous Translation Models

Nal Kalchbrenner
Department of Computer Science
University of Oxford
nal.kalchbrenner,phil.blunsom@cs.ox.ac.uk

Abstract

We introduce a class of probabilistic translation models called Recurrent Continuous Translation Models that are purely based on continuous representations for words, phrases and sentences and do not rely on discrete units of phrases and sentences. The models have an encoder and a decoding aspect. The generation of the translation is modelled with a target Recurrent Language Model, while the translation on the source side is modelled with a Continuous Sentence Model. Through various experiments, we show that if models obtain a perplexity with respect to gold translations that is comparable to state-of-the-art alignment-based translation models. Secondly, we show that they are remarkably sensitive to the word order, syntax, and meaning of the input sentence despite lacking a alignment. Finally we show that they match a state-of-the-art system when rescoring a n-best lists of translations.

1 Introduction

The main contribution is learning to jointly improve improvements with the recurrent neural network language model over a competitive n-gram baseline across several language pairs. We even observe consistent gains when pairing the model with a large n-gram model trained on up to 575 times more training data. The joint model is trained on a state-of-the-art system when rescoring a n-best lists of translations.

arXiv:1406.1078v3 [cs.CL] 3 Sep 2014

Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation

Kyunghyun Cho
Bart van Rijsbergen
Université de Montréal
firstname.lastname@umontreal.ca

Caglar Gulcehre
Jacobs University, Germany
d.bahdanau@jacobs-university.de

Dmitry Bahdanau
Université de Montréal
first.name.last.name@umontreal.ca

Fethi Bougares
Holger Schwenk
Université du Maine, France
firstname.lastname@lumii.univ-lleman.fr

Yoshua Bengio
Université de Montréal, CIFAR Senior Fellow
find.melon@the.web

Abstract

In this paper, we propose a novel neural network model based on an RNN Encoder-Decoder that consists of two recurrent neural networks (RNN). One RNN encodes a sequence of symbols into a fixed-length vector representation, and the other decodes the representation into another sequence of symbols. The encoder and the decoder of the proposed model are jointly trained to predict the probability of a target sequence given a source sequence. The performance of a statistical machine translation system is empirically found to improve by using the conditional probabilities of phrase pairs produced by the RNN Encoder-Decoder as an addition to the existing log-linear model. Qualitatively, we show that the proposed model learns a semantically and syntactically meaningful representation of linguistic phrases.

1 Introduction

Deep neural networks have shown great success in various applications such as object recognition (see, e.g., Krizhevsky et al., 2012) and speech recognition (see, e.g., Dahl et al., 2012; Povey et al., 2011; Chen et al., 2012). Furthermore, recent research has shown that neural networks can be successfully used in a number of tasks in natural language processing (NLP). These include, but are not limited to, language modeling (Bengio et al., 2003), paraphrase detection (Socher et al., 2011) and word embedding extraction (Mikolov et al., 2013). In the field of statistical machine translation (SMT), deep neural networks have begun to show promising results. (Schwenk, 2012) summarizes a successful usage of feed-forward neural networks in the framework of phrase-based SMT systems.

We qualitatively analyze the trained RNN Encoder-Decoder by comparing its phrase scores with those given by the existing translation model.

The qualitative analysis shows that the RNN Encoder-Decoder is better at capturing linguistic regularities in the phrase table, indirectly explaining the observed improvements in the overall machine translation performance. The further analysis of the model reveals that the RNN Encoder-Decoder learns a continuous space representation of a phrase that preserves both the semantic and syntactic structure of the phrase.

Sequence to Sequence Learning with Neural Networks

Ilya Sutskever
Google
ilya.sutskever@google.com

Oriol Vinyals
Google
vinyals@google.com

Quoc V. Le
Google
quoc.v.le@google.com

Abstract

Deep Neural Networks (DNNs) are powerful models that have achieved excellent performance on difficult learning tasks. Although DNNs work well whenever large labeled training sets are available, they cannot be used to map sequences to sequences. In this paper, we present a general end-to-end approach to sequence learning that makes minimal assumptions on the sequence structure. Our method uses a recurrent neural network that takes a sequence of inputs and maps it to a vector of a fixed dimension. We then train a deep LSTM to decode the target sequence from this vector. Our main result is that an English to French translation task from the WMT'14 dataset, the translations produced by the LSTM achieve a BLEU score of 34.8 on the entire test set, where the LSTM's BLEU score was previously 31.5. Additionally, we show that the LSTM's performance is not significantly different on long sentences. For comparing a phrase-based SMT system achieves a BLEU score of 33.3 on the same dataset. When we used the LSTM to rerank the 1000 hypotheses produced by the aforementioned SMT system, its BLEU score increased to 36.5, which is close to the previous best result on this task. The LSTM also handles semantic dependencies between words in a sentence, as well as word order and are relatively invariant to the active and the passive voice. Finally, we found that reversing the order of the words in all source sentences (but not target sentences) improved the LSTM's performance markedly, because doing so introduced many short term dependencies between the source and the target sentence which made the optimization problem easier.

1 Introduction

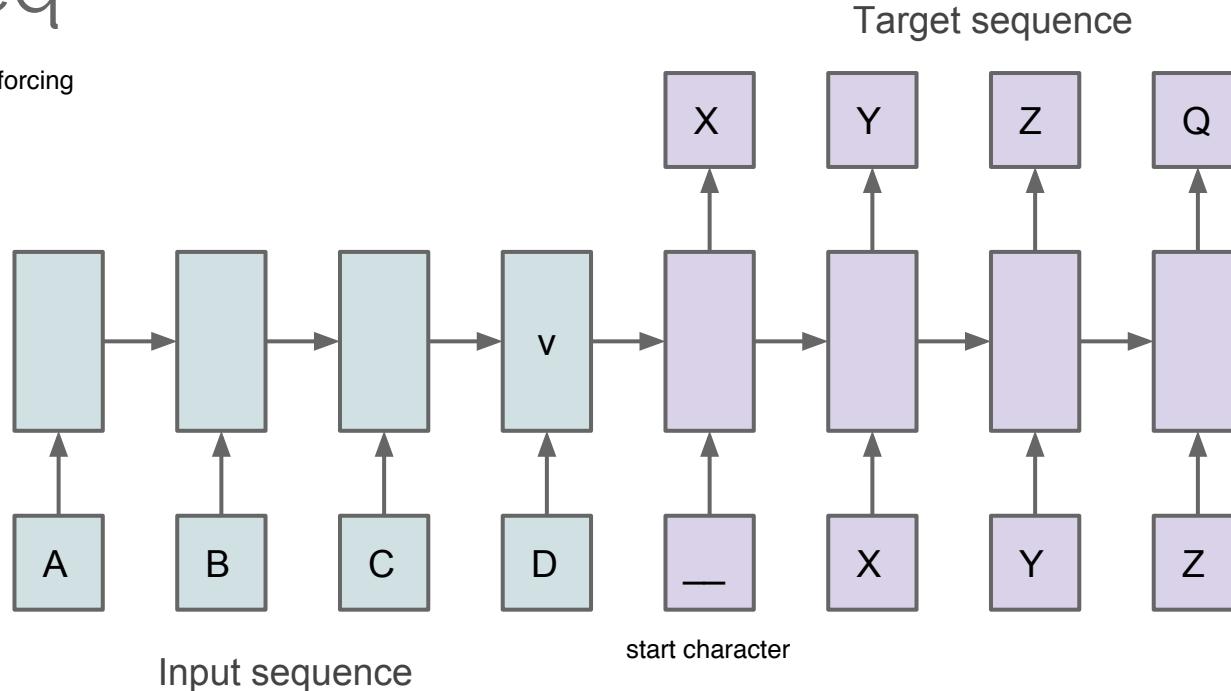
Deep Neural Networks (DNNs) are extremely powerful machine learning models that achieve excellent performance on difficult problems such as speech recognition [13,7] and visual object recognition [19, 6, 21, 20]. DNNs are powerful because they can perform arbitrary parallel computation for a modest number of steps. A simple example of the power of DNNs is the ability to learn an N-gram language model using a 2-layer fully connected network [27], while linear networks are related to conventional statistical models, they learn an intricate computation. Furthermore, large DNNs can be trained with supervised backpropagation whenever the labeled training set has enough information to specify the network's parameters. Thus, if there exists a parameter setting of a large DNN that achieves good results (for example, because humans can solve the task very rapidly), supervised backpropagation will find it and solve the task.

Despite their flexibility and power, DNNs can only be applied to problems whose inputs and targets can be sensibly encoded with vectors of fixed dimensionality. It is a significant limitation, since many important problems are best expressed with sequences whose lengths are not known *a-priori*. For example, speech recognition and machine translation are sequential problems. Likewise, question answering can also be seen as mapping a sequence of words representing the question to a

1. Auli, M., et al. "Joint Language and Translation Modeling with Recurrent Neural Networks." *EMNLP (2013)*
2. Kalchbrenner, N., et al. "Recurrent Continuous Translation Models." *EMNLP (2013)*
3. Cho, K., et al. "Learning Phrase Representations using RNN Encoder-Decoder for Statistical MT." *EMNLP (2014)*
4. Sutskever, I., et al. "Sequence to Sequence Learning with Neural Networks." *NIPS (2014)*

Seq2Seq

training with teacher forcing

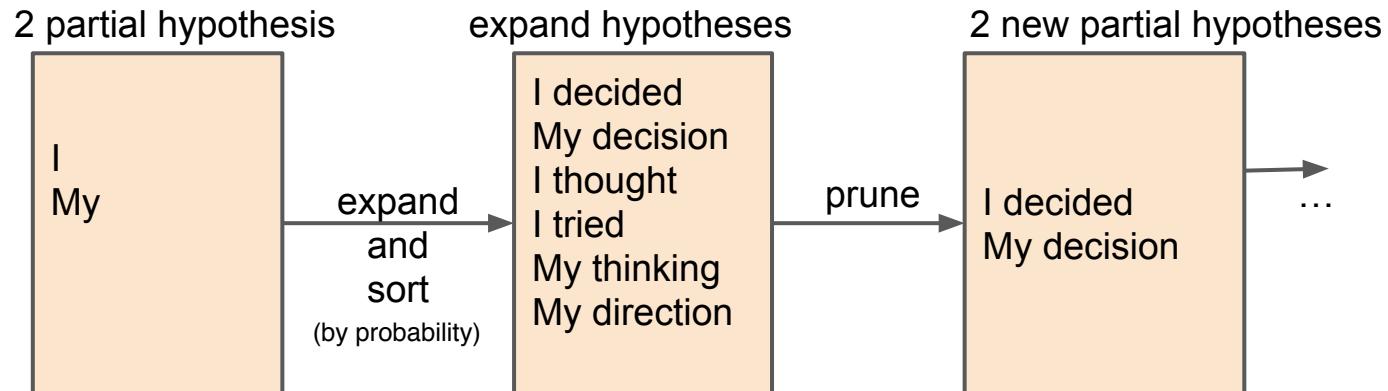


$$P(y_1, \dots, y_{T'} | x_1, \dots, x_T) = \prod_{t=1}^{T'} p(y_t | v, y_1, \dots, y_{t-1})$$

Decoding in a Nutshell (Beam Size 2)

beam search in inference

$$y^* = \arg \max_{y_1, \dots, y_{T'}} P(y_1, \dots, y_{T'} | x_1, \dots, x_T)$$

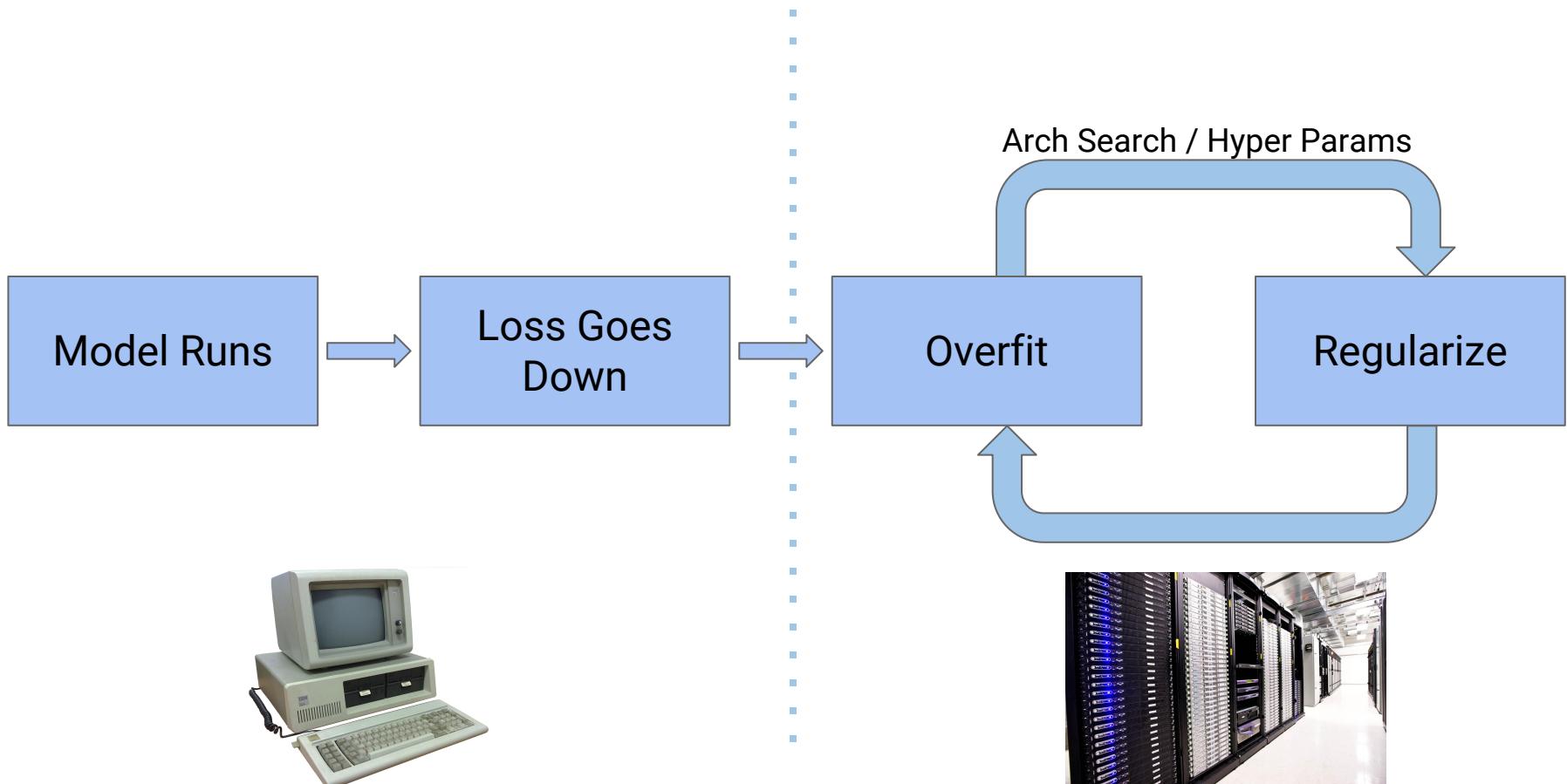


Code

Source: <https://github.com/keveman/tensorflow-tutorial/blob/master/PTB%20Word%20Language%20Modeling.ipynb>

```
class LSTMCell(object):
    def __init__(self, state_size):
        self.state_size = state_size
        self.W_f = tf.Variable(self.initializer())
        self.W_i = tf.Variable(self.initializer())
        self.W_o = tf.Variable(self.initializer())
        self.W_C = tf.Variable(self.initializer())
        self.b_f = tf.Variable(tf.zeros([state_size]))
        self.b_i = tf.Variable(tf.zeros([state_size]))
        self.b_o = tf.Variable(tf.zeros([state_size]))
        self.b_C = tf.Variable(tf.zeros([state_size]))
    def __call__(self, x_t, h_t1, C_t1):
        X = tf.concat(1, [h_t1, x_t])
        f_t = tf.sigmoid(tf.matmul(X, self.W_f) + self.b_f)
        i_t = tf.sigmoid(tf.matmul(X, self.W_i) + self.b_i)
        o_t = tf.sigmoid(tf.matmul(X, self.W_o) + self.b_o)
        Ctilde_t = tf.tanh(tf.matmul(X, self.W_C) + self.b_C)
        C_t = f_t * C_t1 + i_t * Ctilde_t
        h_t = o_t * tf.tanh(C_t)
        return h_t, C_t
    def initializer(self):
        return tf.random_uniform([2*self.state_size, self.state_size],
                               -0.1, 0.1)
```

Deep Learning Vicious Cycle



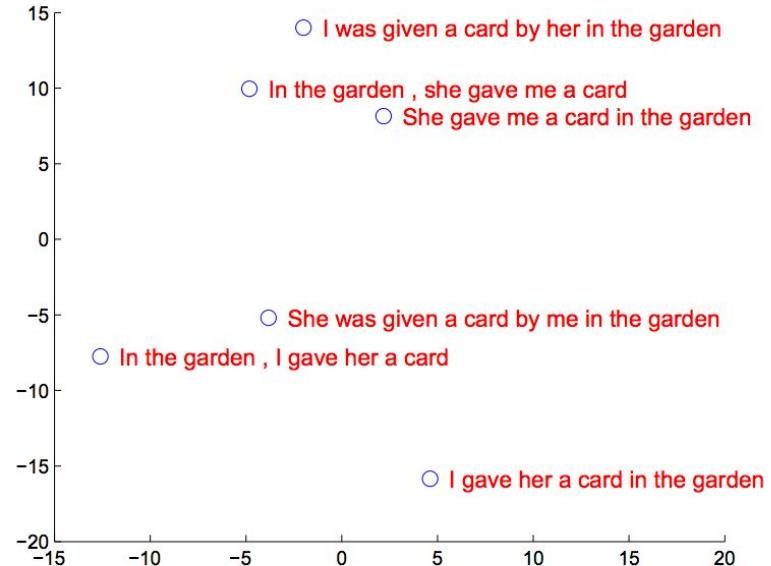
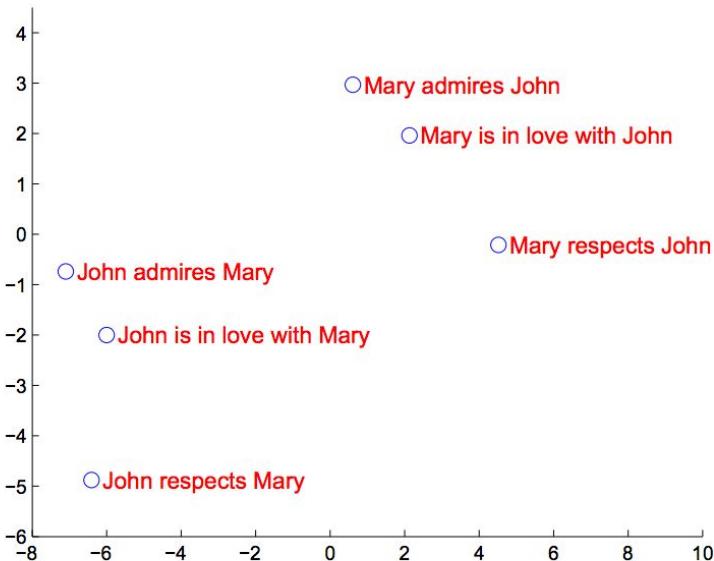
(Some) Tricks of the Trade

- Long sequences?
 - Attention
 - Bigger state
- Can't overfit?
 - Bigger hidden state
 - Deep LSTM + Skip Connections
- Overfit?
 - Dropout + Ensembles
- Tuning
 - Keep calm and decrease your learning rate
 - Initialization of parameters is critical (in seq2seq we used $U(-0.05, 0.05)$)
 - Clip the gradients!
 - E.g. if $\|grad\| > 5$: $grad = grad / \|grad\| * 5$

Applications

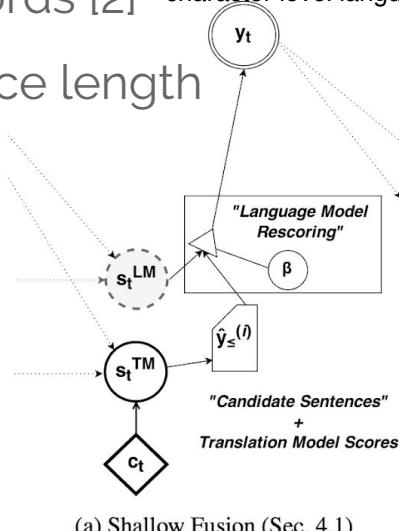
Machine Translation

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
Ensemble of 5 reversed LSTMs, beam size 1	33.00
Ensemble of 2 reversed LSTMs, beam size 12	33.27
Ensemble of 5 reversed LSTMs, beam size 2	34.50
Ensemble of 5 reversed LSTMs, beam size 12	34.81

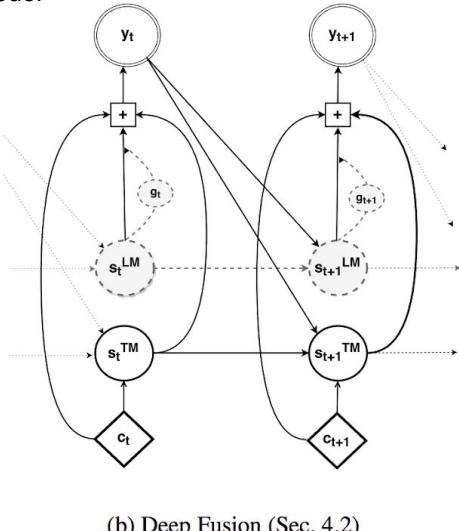


Machine Translation - other concerns

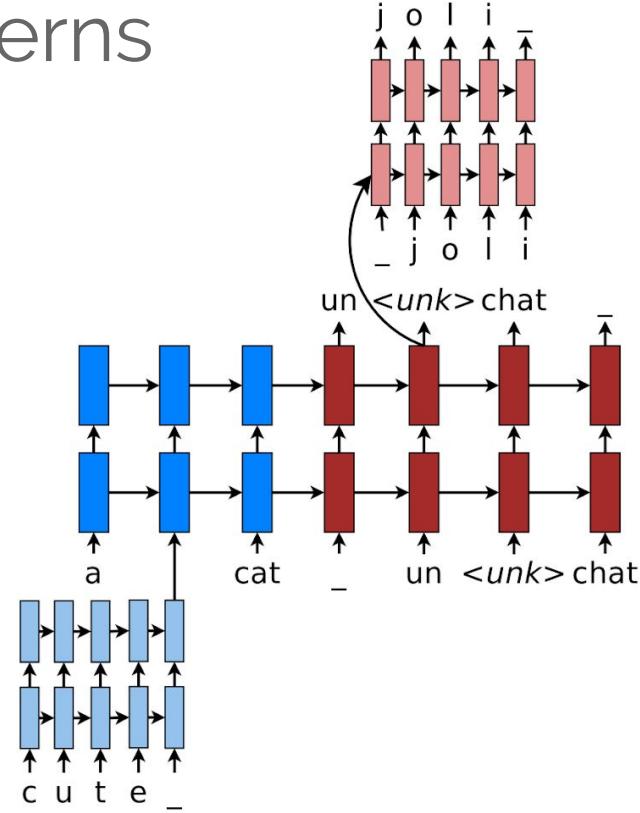
- Using Language Models [1] lack of training data
- OOV words [2] character level language model
- Sequence length



(a) Shallow Fusion (Sec. 4.1)



(b) Deep Fusion (Sec. 4.2)



1. Gulcehre, C., et al. "On using monolingual corpora in neural machine translation." *arXiv* (2015).
2. Luong, T., and Manning, C. "Achieving open vocabulary neural MT with hybrid word-character models." *arXiv* (2016).

Image Captioning

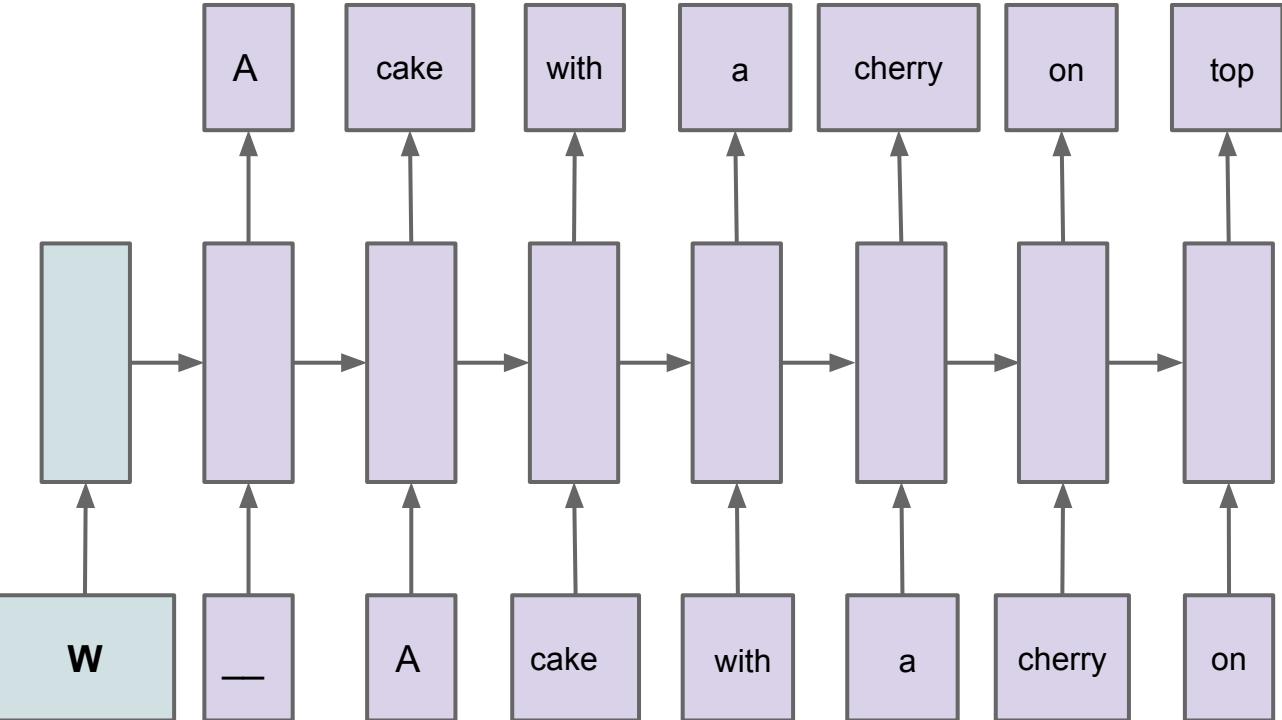
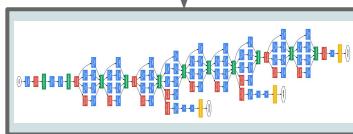
$p(\text{English} \mid \text{French})$



$p(\text{English} \mid \text{Image})$

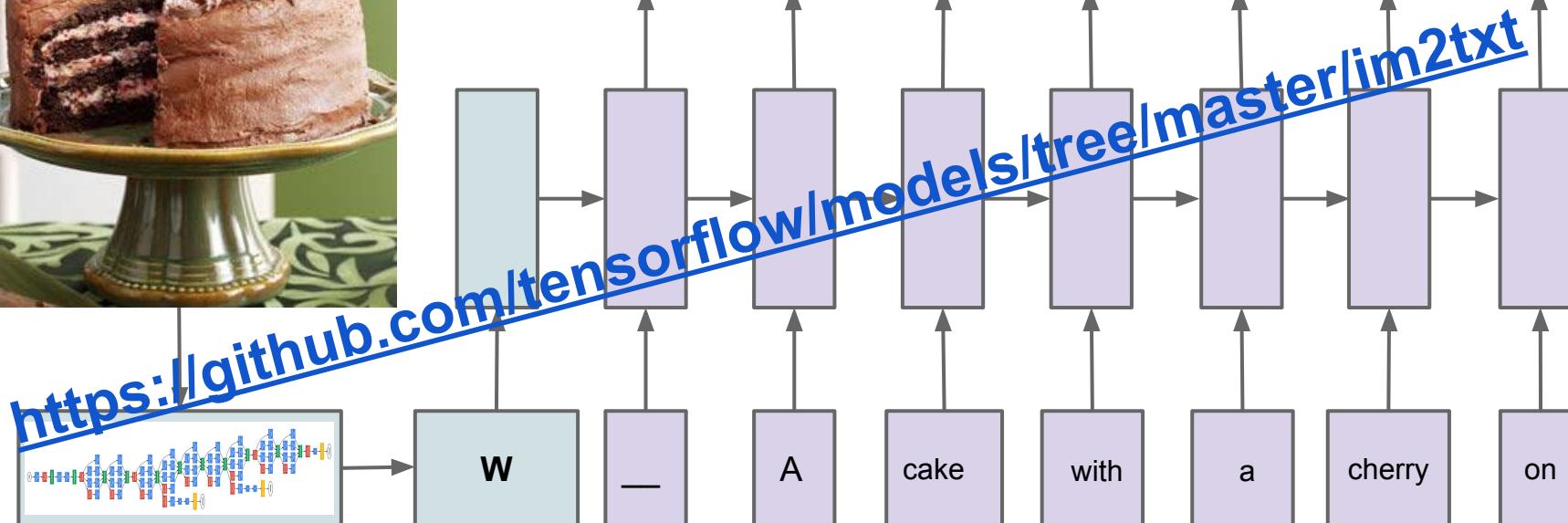
1. Vinyals, O., et al. "Show and Tell: A Neural Image Caption Generator." *CVPR* (2015).
2. Mao, J., et al. "Deep captioning with multimodal recurrent neural networks (m-rnn)." *ICLR* (2015).
3. Karpathy, A., Li, F., "Deep visual-semantic alignments for generating image descriptions." *CVPR* (2015).

Image Captioning



$$\theta^* = \arg \max_{\theta} p(S|I)$$

Image Captioning



$$\theta^* = \arg \max_{\theta} p(S|I)$$

Image Captioning



Human: A close up of two bananas with bottles in the background.

BestModel: A bunch of bananas and a bottle of wine.

InitialModel: A close up of a plate of food on a table.

Image Captioning



Human: A woman holding up a yellow banana to her face.

BestModel: A woman holding a banana up to her face.

InitialModel: A close up of a person eating a hot dog.

Image Captioning



Human: A man outside cooking with a sub in his hand.

BestModel: A man is holding a sandwich in his hand.

InitialModel: A man cutting a cake with a knife.

Image Captioning



Human: Someone is using a small grill to melt his sandwich.

BestModel: A person is cooking some food on a grill.

InitialModel: A pizza sitting on top of a white plate.

Image Captioning



Human: A blue , yellow and red train travels across the tracks near a depot.

BestModel: A blue and yellow train traveling down train tracks.

InitialModel: A train that is sitting on the tracks.

Learning To Execute [Zaremba & Sutskever, 2014]

Learning To Execute

- One of the first (modern) examples of learning algorithms
- 2014--??? “era of discovery” → Apply seq2seq to *everything*

Input:

```
j=8584  
for x in range(8):  
    j+=920  
    b=(1500+j)  
    print((b+7567))
```

Target: 25011.

Input:

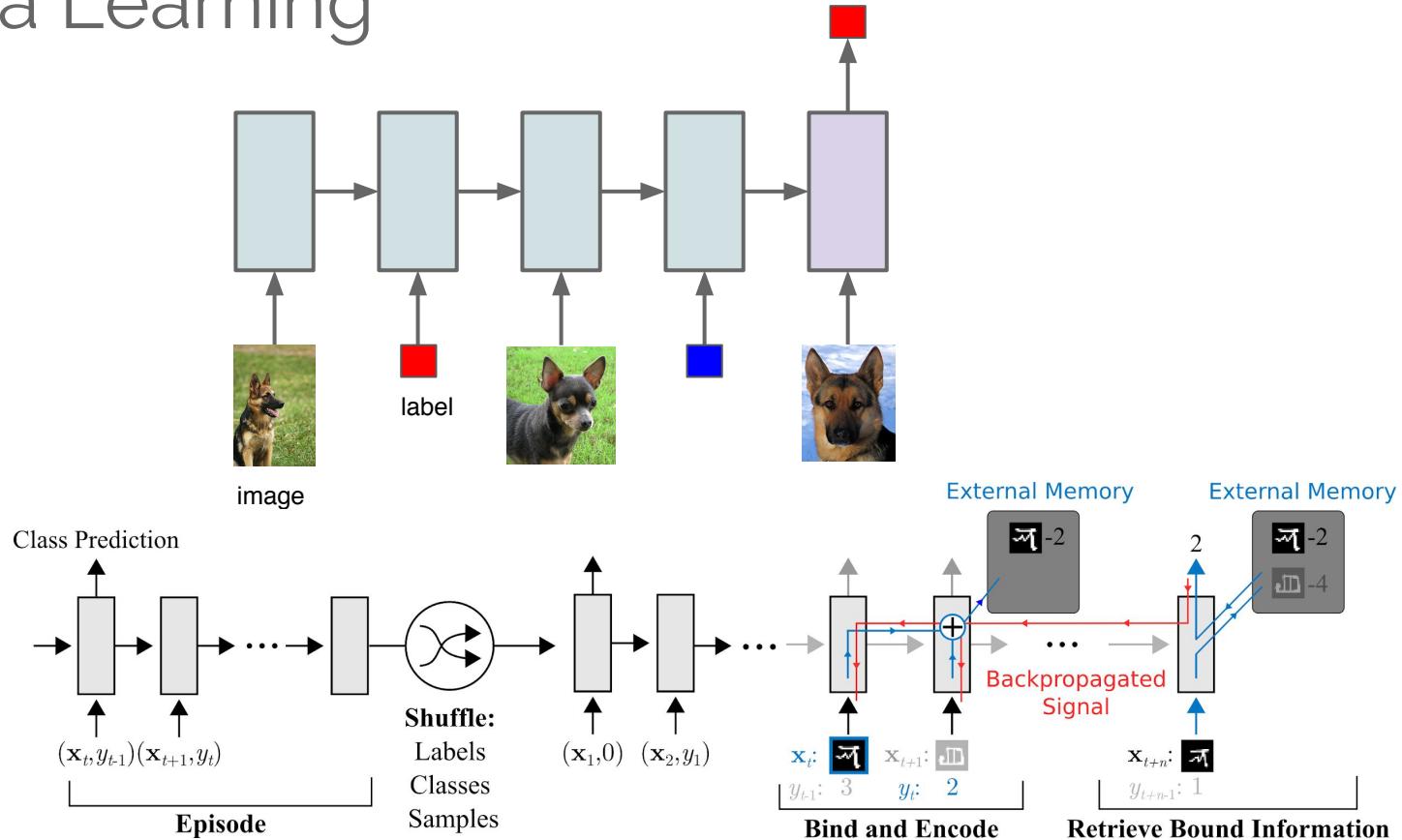
```
i=8827  
c=(i-5347)  
print((c+8704) if 2641<8500 else 5308)
```

Target: 12184.

Input:

```
vqppkn  
sqdvfljmnc  
y2vxddsepnimcbvubkomhrpliibtwztbljipcc  
Target: hkhpg
```

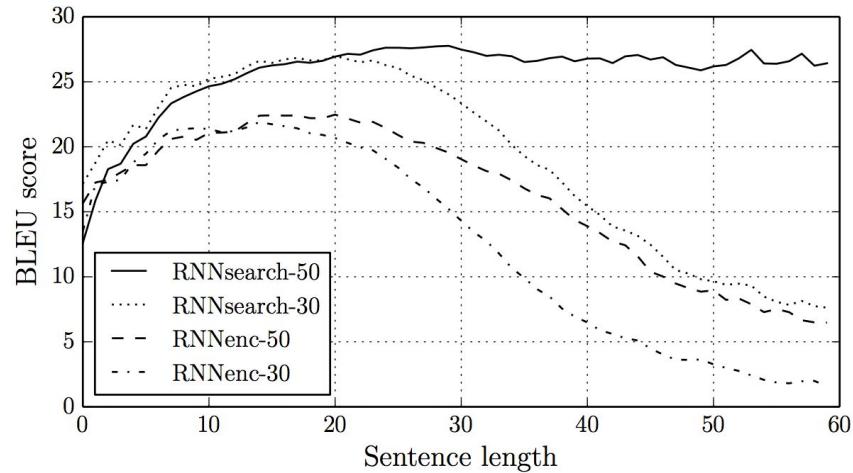
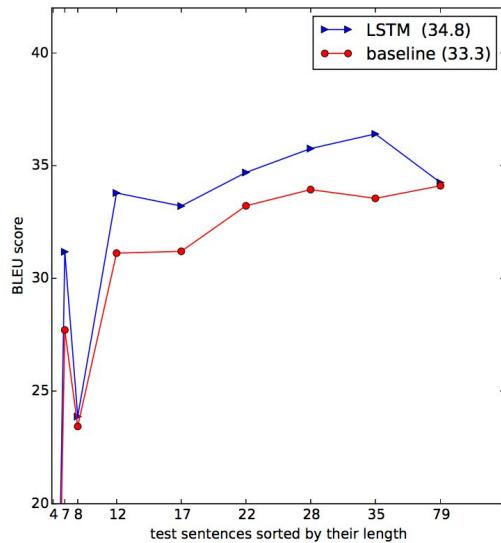
Meta Learning



1. Young, A., et al. "Meta-learning with backpropagation." *Neural Nets* (2001).
2. Santoro, A., et al. "Meta-Learning with Memory-Augmented Neural Networks." *ICML* (2016).
3. Vinyals, O., et al. "Matching Networks for One Shot Learning." *NIPS* (2016).

Seq2Seq - Limitations

- Fixed Size Embeddings are easily overwhelmed by long inputs or long outputs



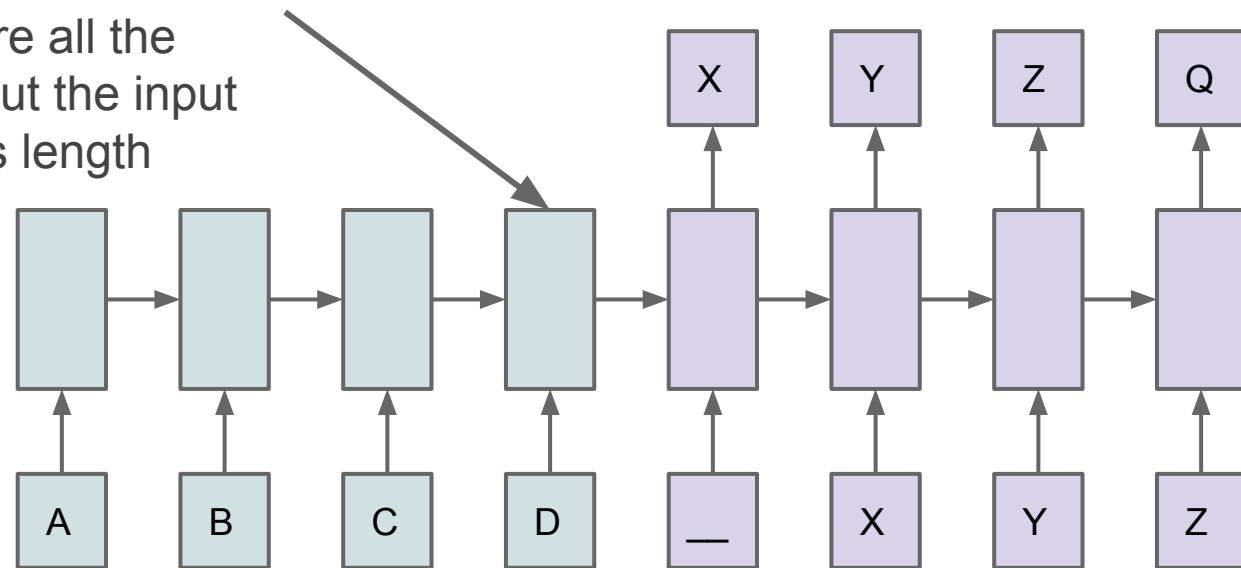
1. Sutskever, I., et al. "Sequence to Sequence Learning with Neural Networks." *NIPS* (2014)
2. Bahdanau, D., et al. "Neural Machine Translation by Jointly Learning to Align and Translate." *ICLR* (2015)

Attention

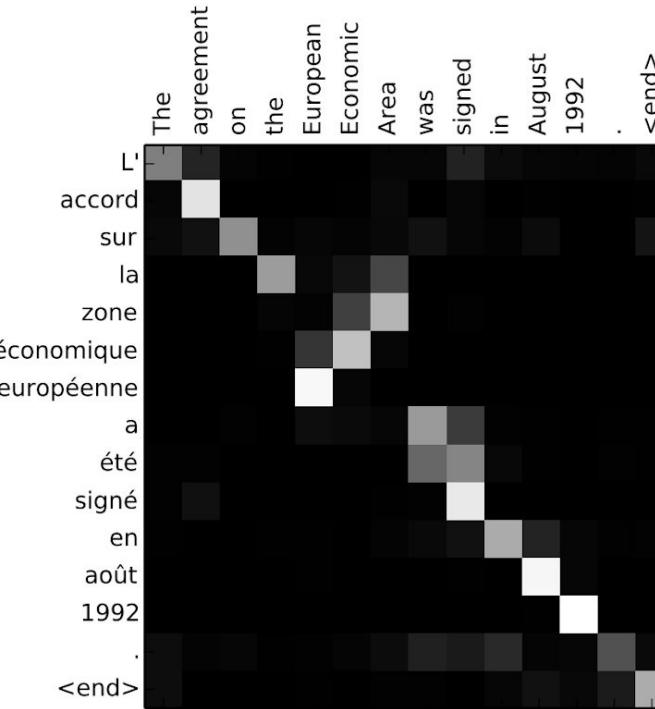
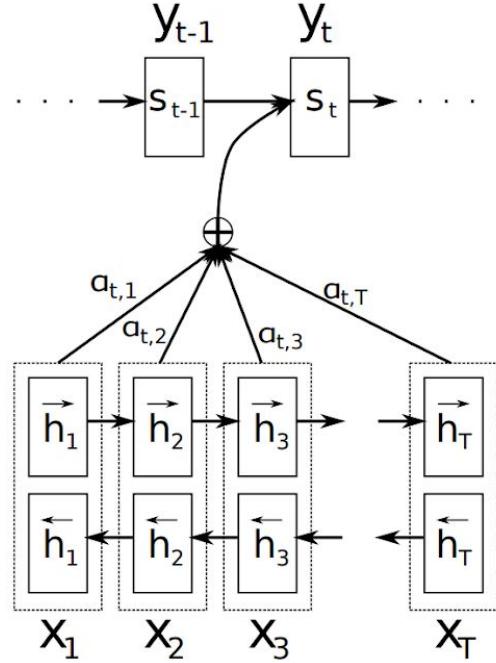
Seq2Seq -- The issue with long inputs

- Same embedding informs the entire output
- Needs to capture all the information about the input regardless of its length

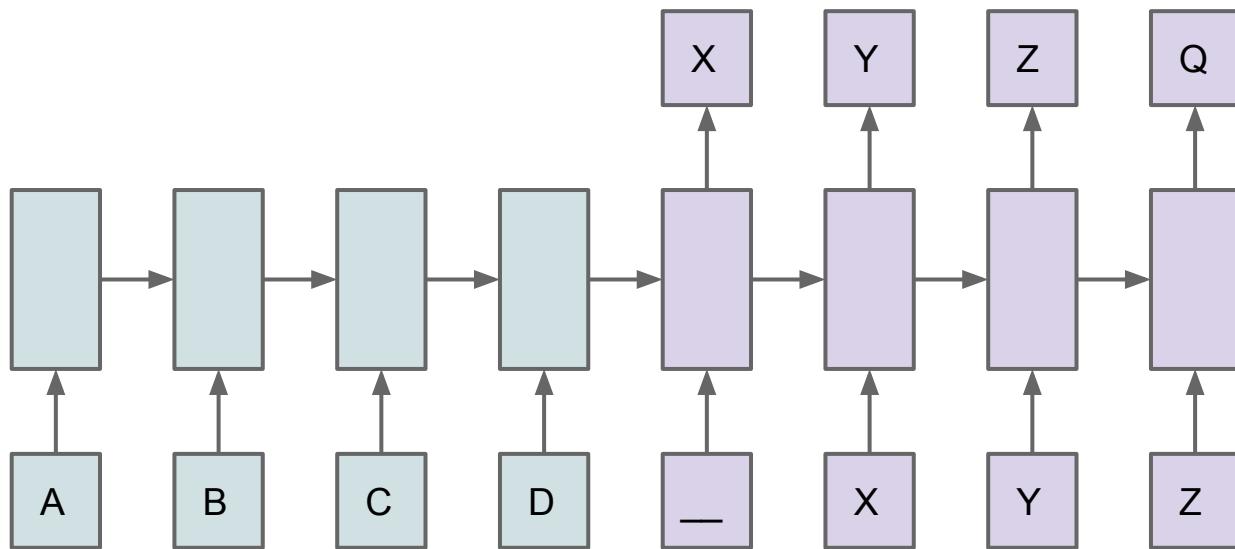
Is there a better way to pass the information from encoder to the decoder ?



Seq2Seq with Attention

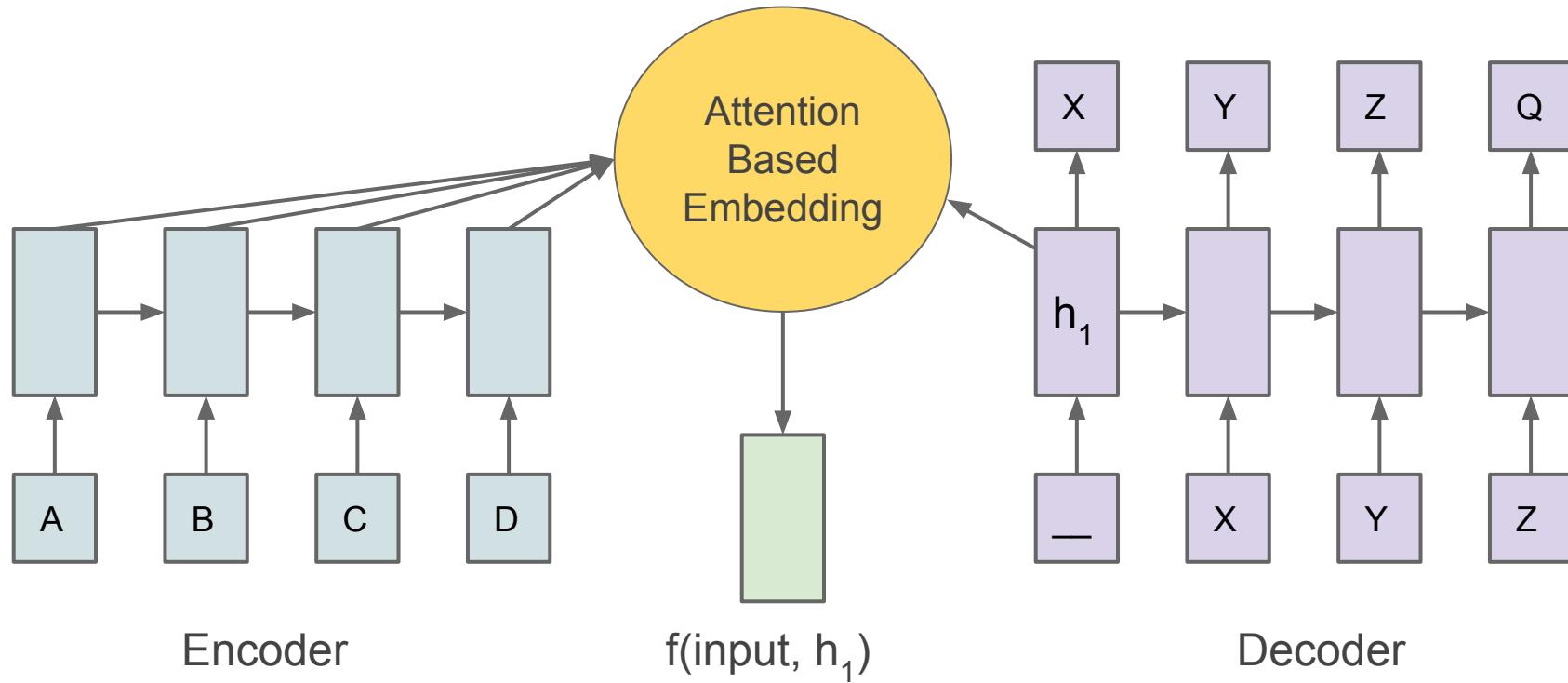


Seq2Seq



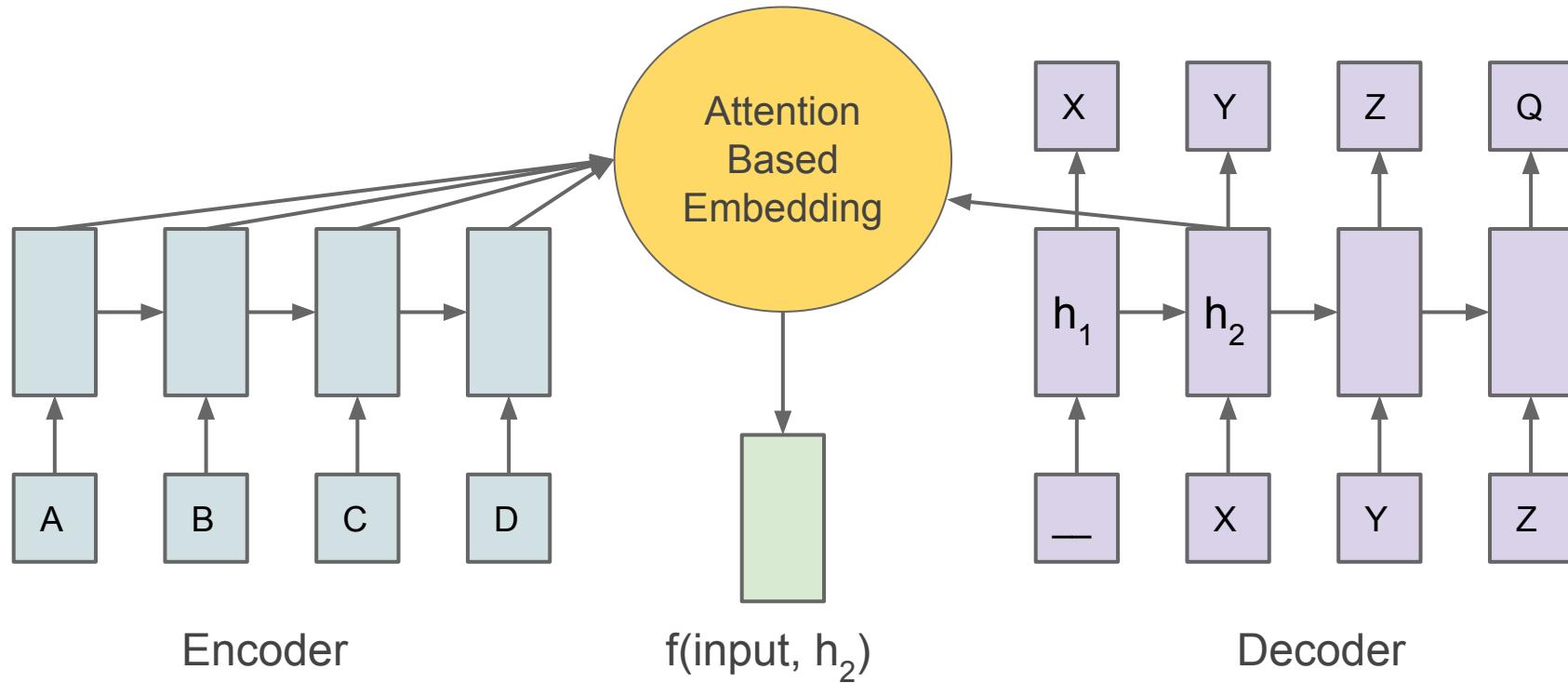
Seq2Seq with Attention

- A different embedding computed for every output step



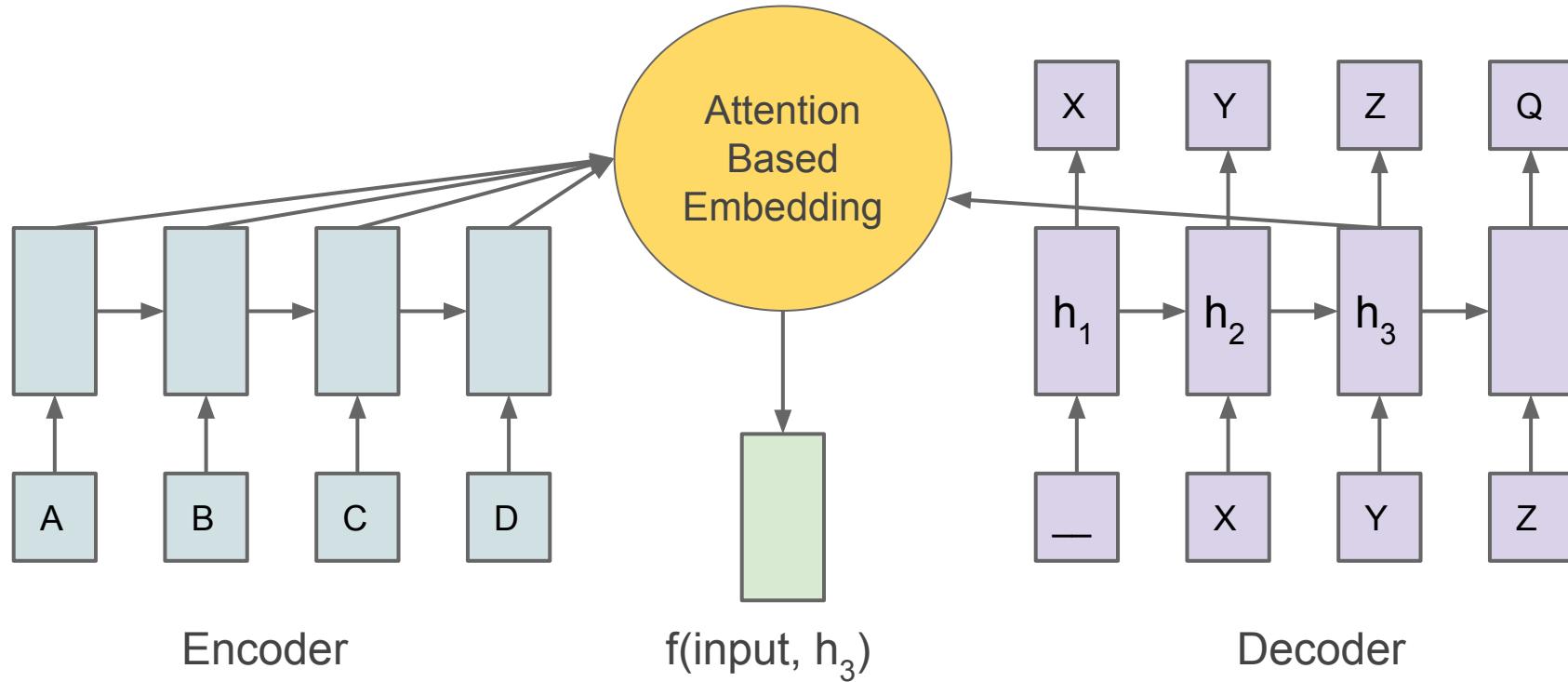
Seq2Seq with Attention

- A different embedding computed for every output step



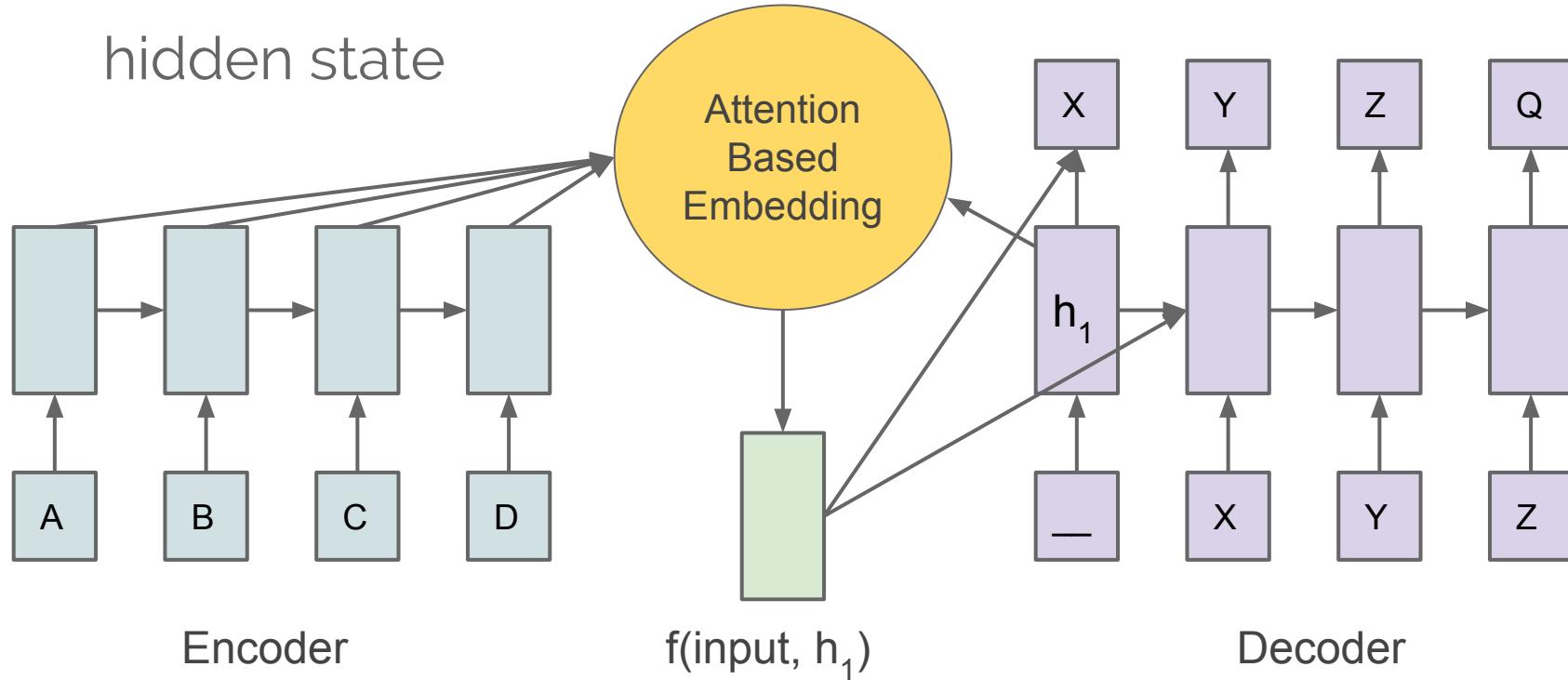
Seq2Seq with Attention

- A different embedding computed for every output step



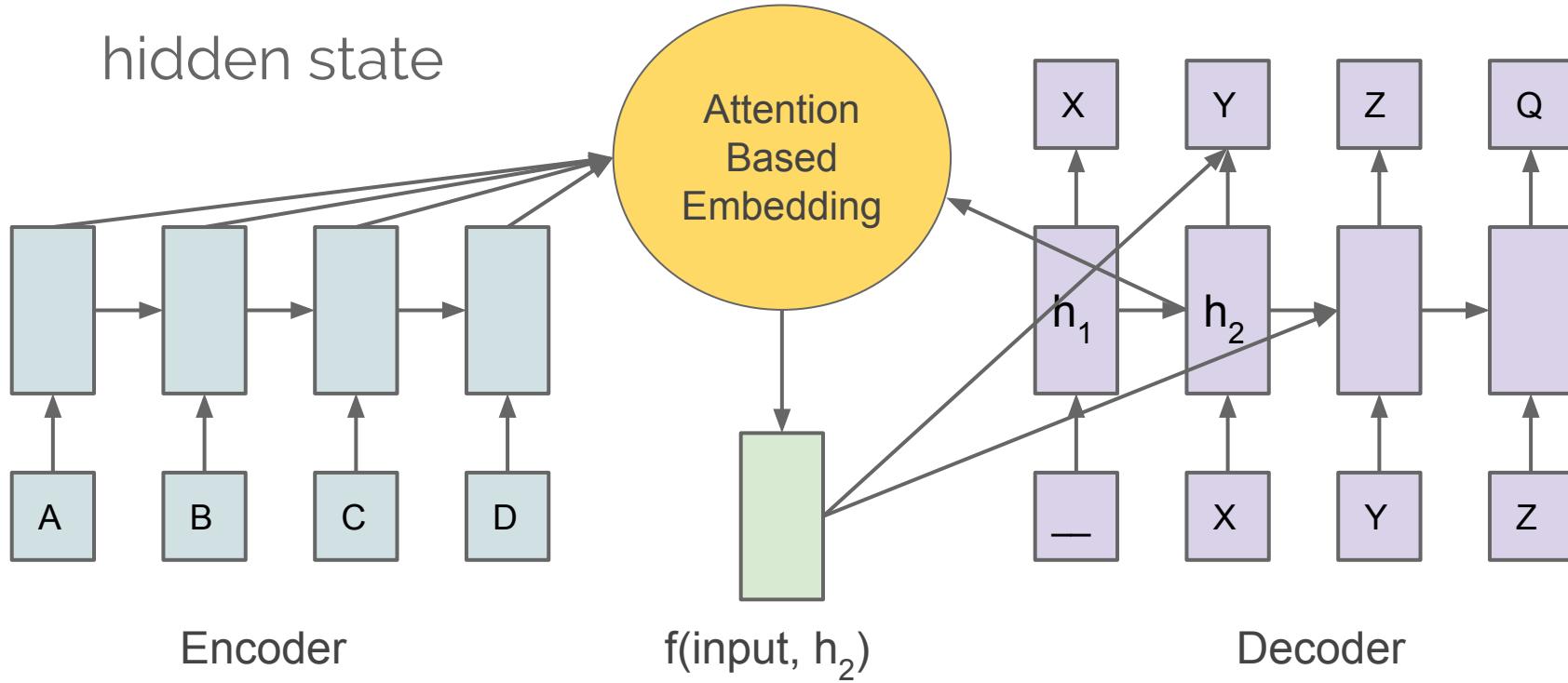
Seq2Seq with Attention

- Embedding used to predict output, and compute next hidden state



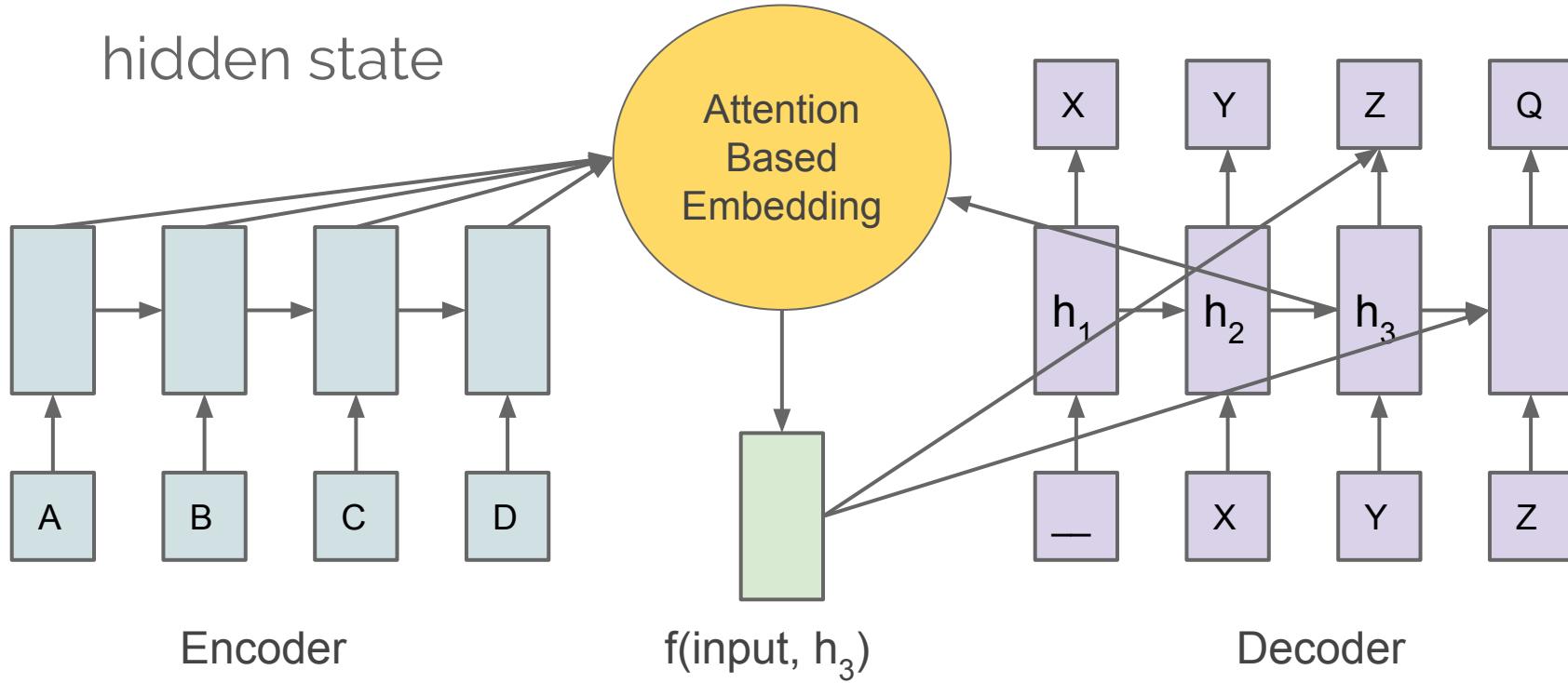
Seq2Seq with Attention

- Embedding used to predict output, and compute next hidden state



Seq2Seq with Attention

- Embedding used to predict output, and compute next hidden state



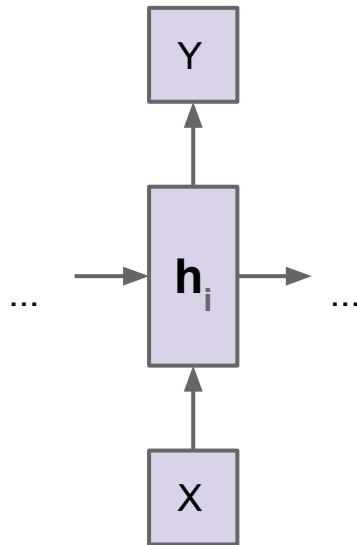
Attention Based Embedding

- Linear blending of embedding RNN states $e_1 e_2 e_3 e_4$ is a natural choice
- How to produce the coefficients (attention vector) for blending ?
 - Content based coefficients based on query state h_i and embedding RNN states $e_1 e_2 e_3 e_4$

Dot product Attention

example 1

- Inputs: “I am a cat.”
- Input RNN states: $e_1 e_2 e_3 e_4$
- Decoder RNN state at step i (query): h_i
- Compute scalars $h_i^T e_1, h_i^T e_2, h_i^T e_3, h_i^T e_4$ representing similarity / relevance between encoder steps and query.
- Normalize $[h_i^T e_1, h_i^T e_2, h_i^T e_3, h_i^T e_4]$ with softmax to produce attention weights, e.g. [0.0 0.05 0.9 0.05]



Content Based Attention

example 2

Attention [Bahdanau, Cho and Bengio, 2014]

$$u_j = v^T \tanh(W_1 e_j + W_2 d) \quad j \in (1, \dots, n)$$

$$a_j = \text{softmax}(u_j) \quad j \in (1, \dots, n)$$

$$d' = \sum_{j=1}^n a_j e_j$$

1-layer attention network

Neural Turing Machine [1] & Memory Networks [2]:

Given a key vector \mathbf{k} and a strength scalar s emitted by the controller, get a weighting (w_1, \dots, w_N) over memory locations where

$$w_i = \frac{e^{\hat{w}_i}}{\sum_j e^{\hat{w}_j}} \quad ; \quad \hat{w}_i = \log(1 + e^s) C(\mathbf{k}, \mathbf{m}_i) \quad ; \quad C(\mathbf{k}, \mathbf{m}_i) = \frac{\mathbf{k} \cdot \mathbf{m}_i}{\|\mathbf{k}\| \|\mathbf{m}_i\|}$$

and $C(\mathbf{k}, \mathbf{m}_i)$ is the cosine similarity between the key and the word \mathbf{m}_i at memory location i

1. Graves, A., et al. "Neural Turing Machines." *arxiv* (2014)
2. Weston, J., et al. "Memory Networks." *arxiv* (2014)

Other strategies for attention models

- Tensored attention
 - Minh-Thang Luong, Hieu Pham, and Christopher D. Manning. "Effective Approaches to Attention-based Neural Machine Translation." EMNLP'15.
- Multiple heads
- Pyramidal encoders
 - William Chan, Navdeep Jaitly, Quoc Le, Oriol Vinyals. "Listen Attend and Spell". ICASSP 2015.
- Hierarchical Attention
 - Andrychowicz, Marcin, and Karol Kurach. "Learning efficient algorithms with hierarchical attentive memory." *arXiv preprint arXiv:1602.03218* (2016).
- Hard Attention
 - Xu, Kelvin, et al. "Show, attend and tell: Neural image caption generation with visual attention." ICML 2015

Applications

Sentence to Constituency Parse Tree

1. Read a sentence
 2. Flatten the tree into a sequence (adding (,))
 3. “Translate” from sentence to parse tree

John has a dog . →

```

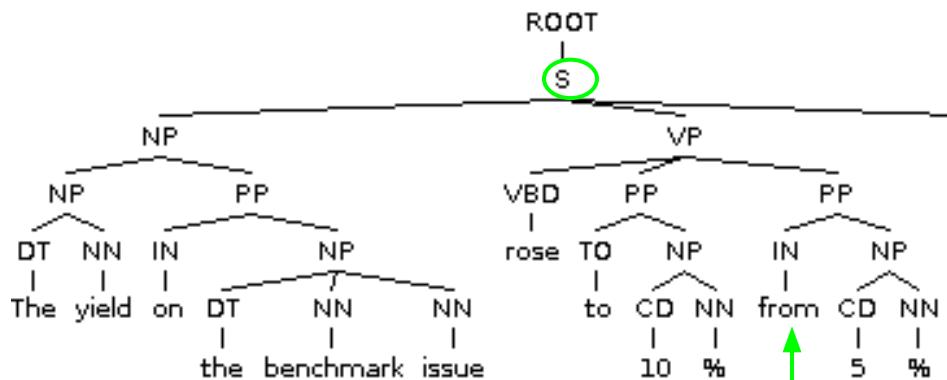
graph TD
    S --- NP1[NP]
    S --- VP[VP]
    NP1 --- NNP[NNP]
    VP --- VBZ[VBZ]
    VP --- NP2[NP]
    NP2 --- DT[DT]
    NP2 --- NN[NN]
    
```

John has a dog . → (S (NP NNP)_{NP} (VP VBZ (NP DT NN)_{NP})_{VP} .)_S

Parsing Example

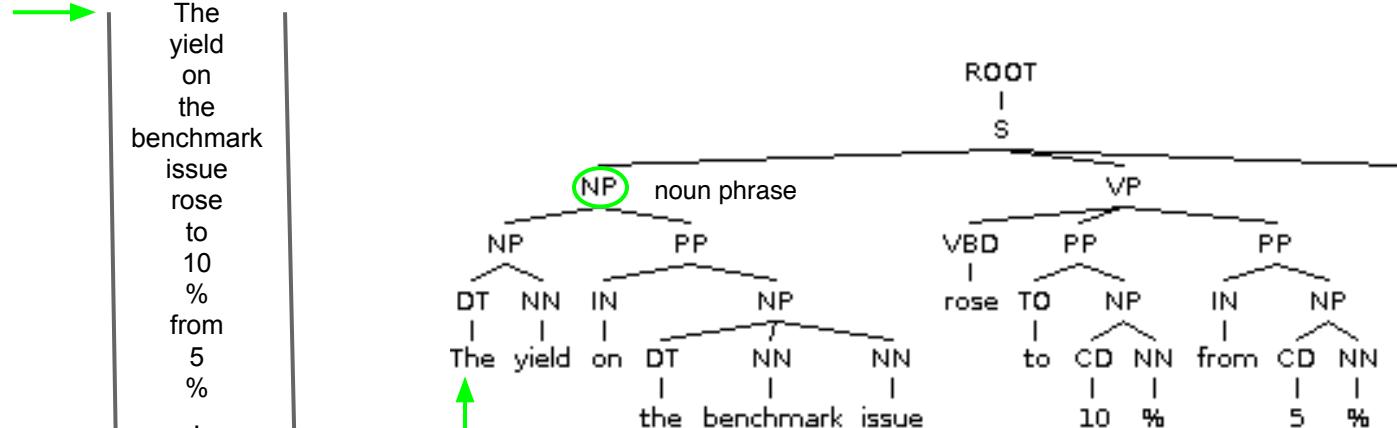
attention follows the parsing structure

The
yield
on
the
benchmark
issue
rose
to
10
%
from
5
%
.



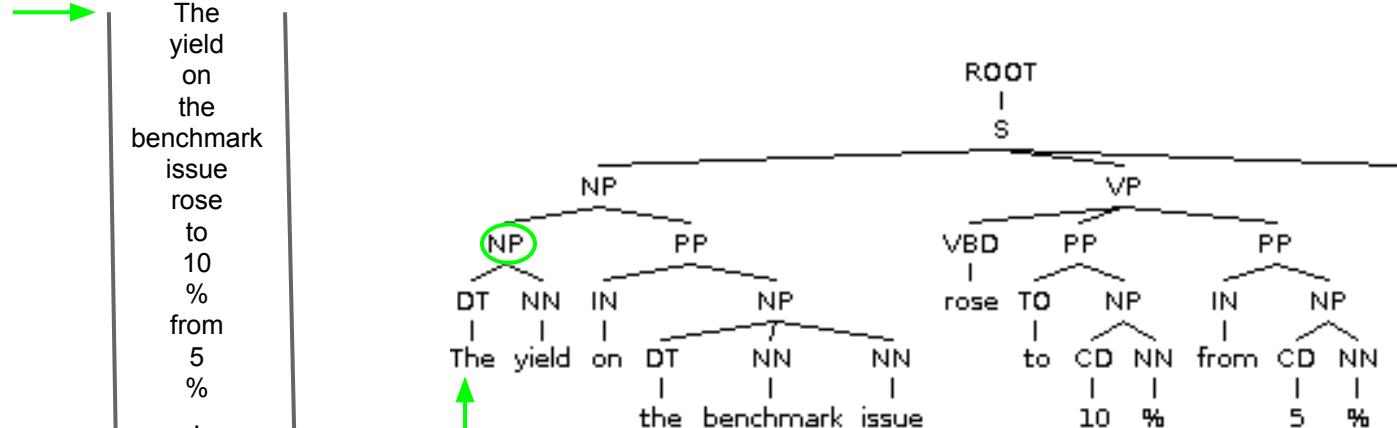
Encoder LSTM
(reverse input)

Parsing Example

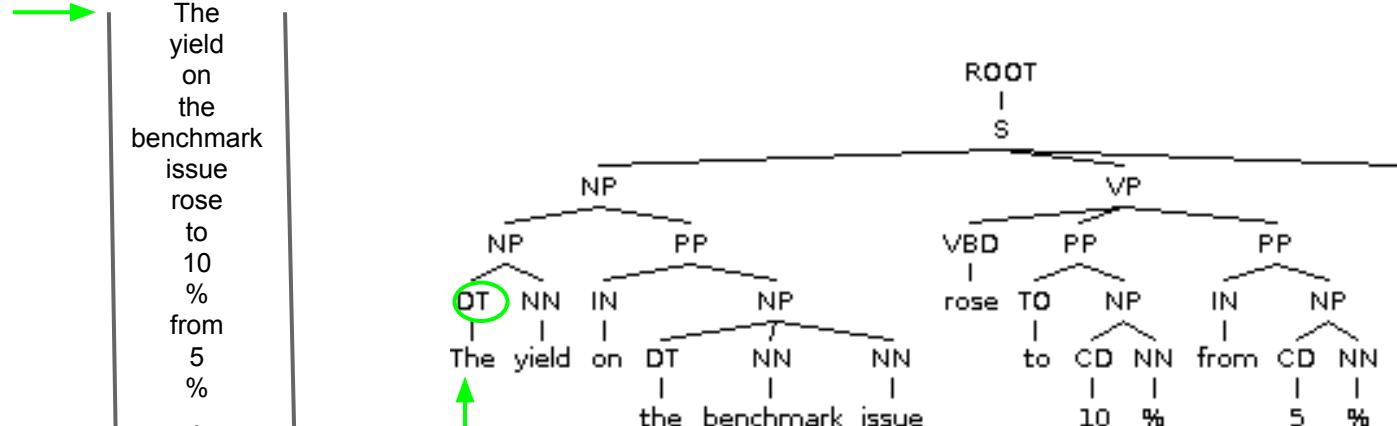


Encoder LSTM
(reverse input)

Parsing Example



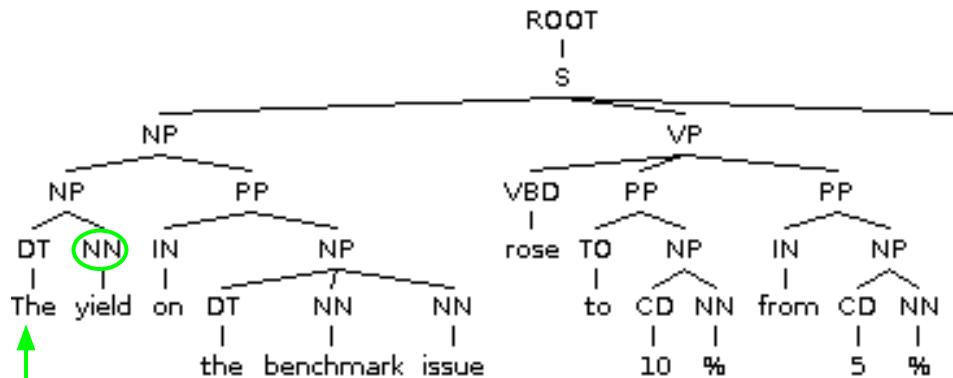
Parsing Example



Encoder LSTM
(reverse input)

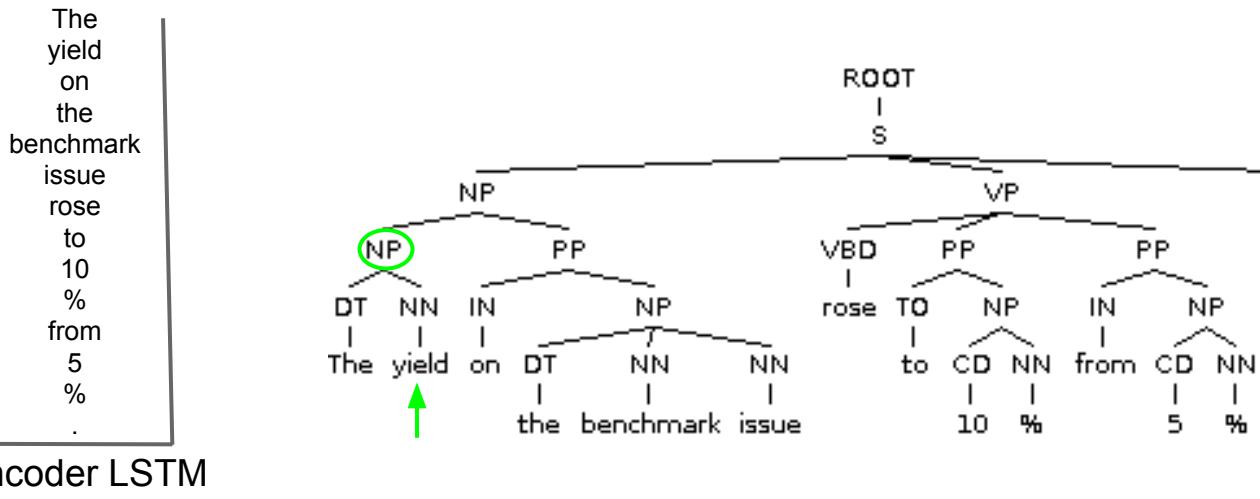
Parsing Example

The yield on the benchmark issue rose to 10 % from 5 %

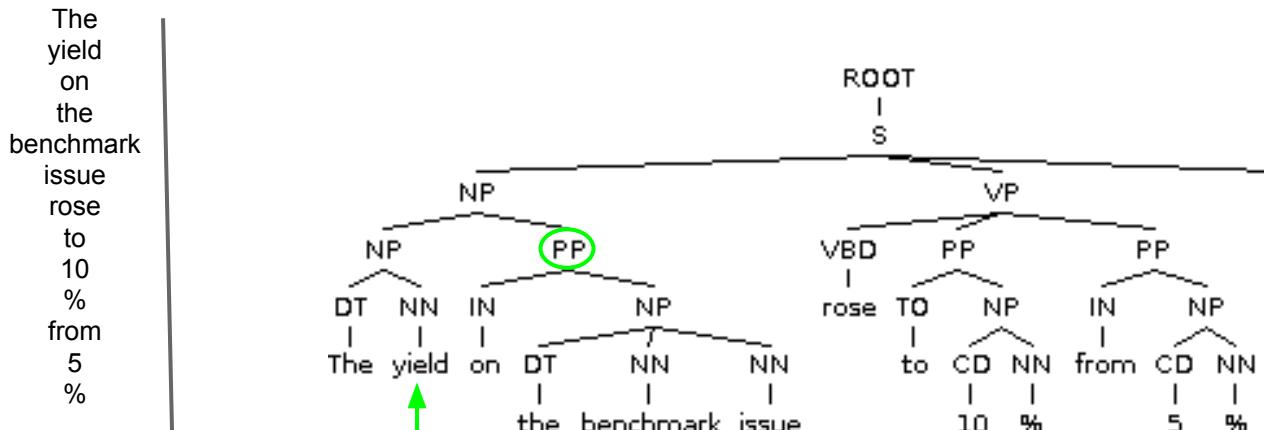


Encoder LSTM (reverse input)

Parsing Example

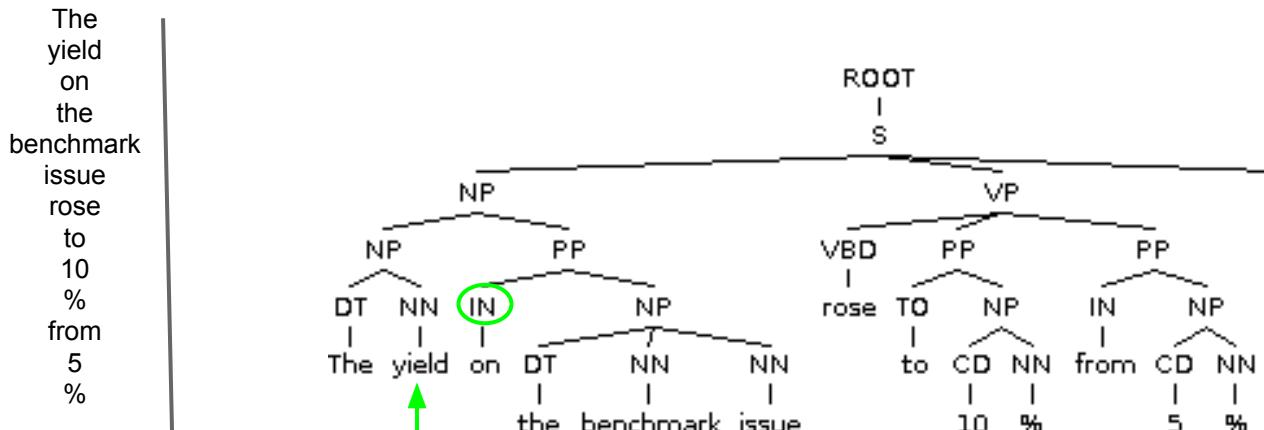


Parsing Example



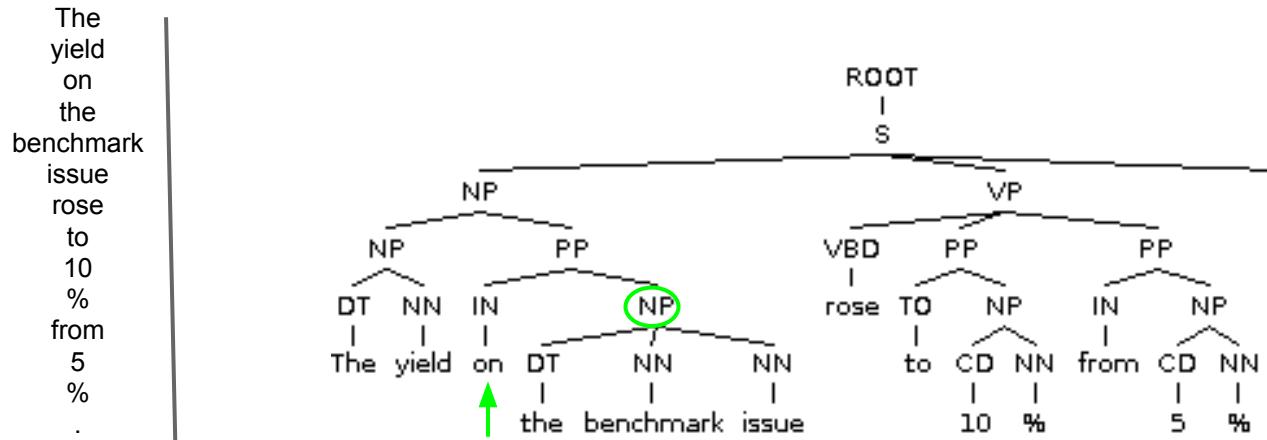
Encoder LSTM
(reverse input)

Parsing Example



Encoder LSTM
(reverse input)

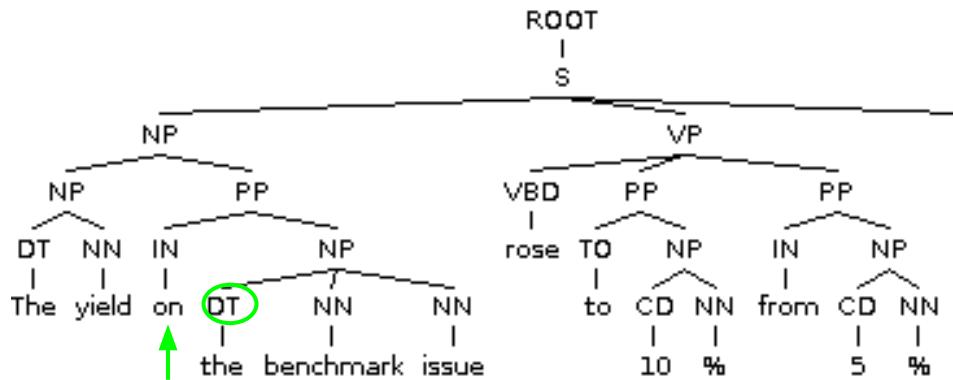
Parsing Example



Encoder LSTM
(reverse input)

Parsing Example

The
yield
on
the
benchmark
issue
rose
to
10
%
from
5
%
.

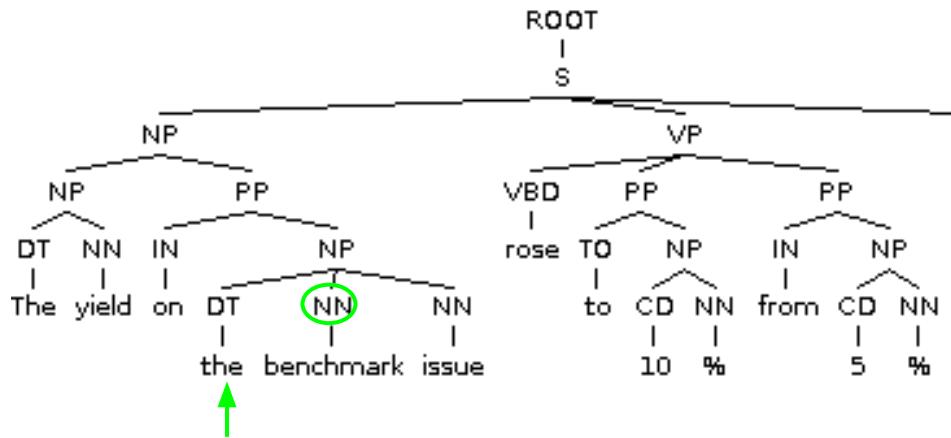


Encoder LSTM
(reverse input)

Parsing Example

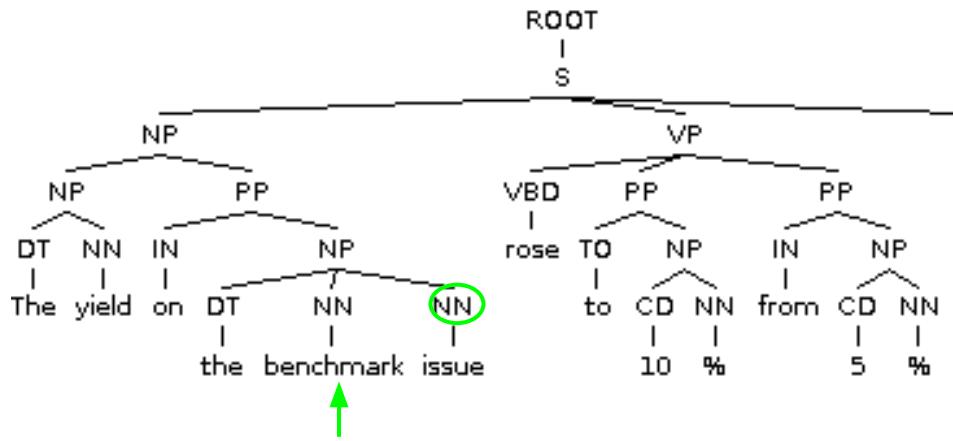
The
yield
on
the
benchmark
issue
rose
to
10
%
from
5
%
.

Encoder LSTM
(reverse input)



Parsing Example

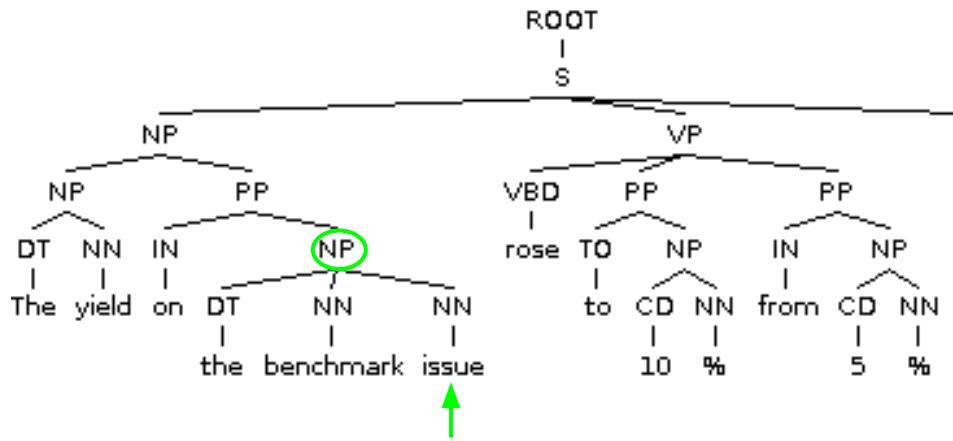
The
yield
on
the
benchmark
issue
rose
to
10
%
from
5
%
.



Encoder LSTM
(reverse input)

Parsing Example

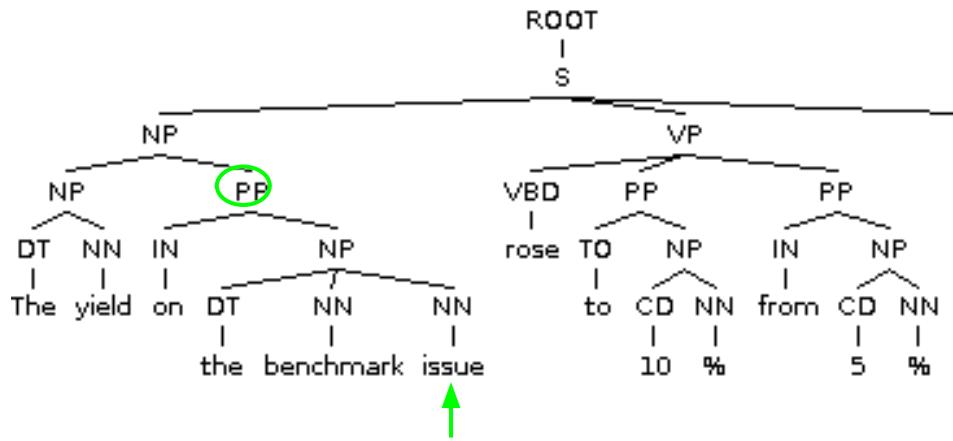
The yield on the benchmark issue rose to 10 % from 5 %



Encoder LSTM (reverse input)

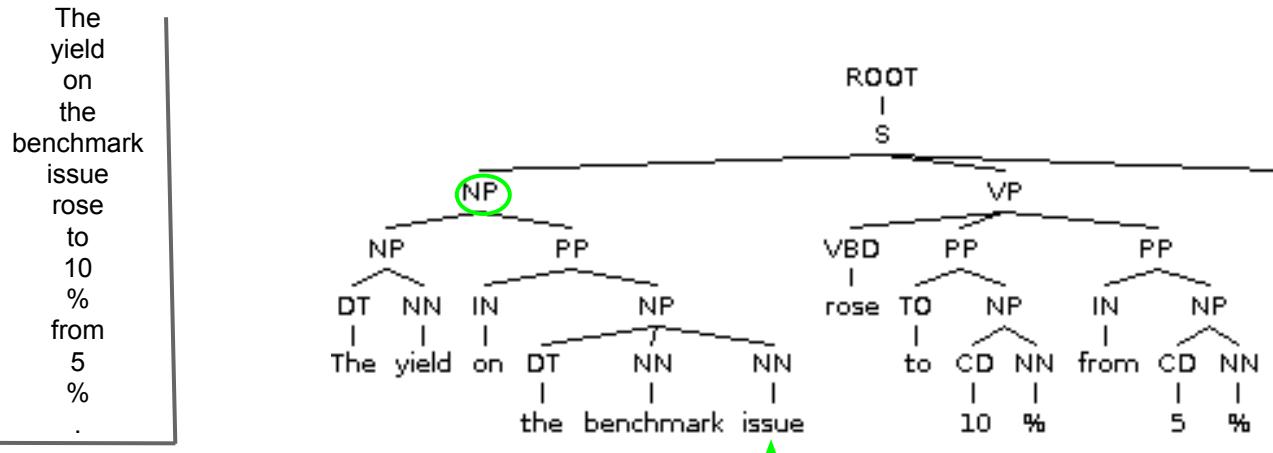
Parsing Example

The
yield
on
the
benchmark
issue
rose
to
10
%
from
5
%
.



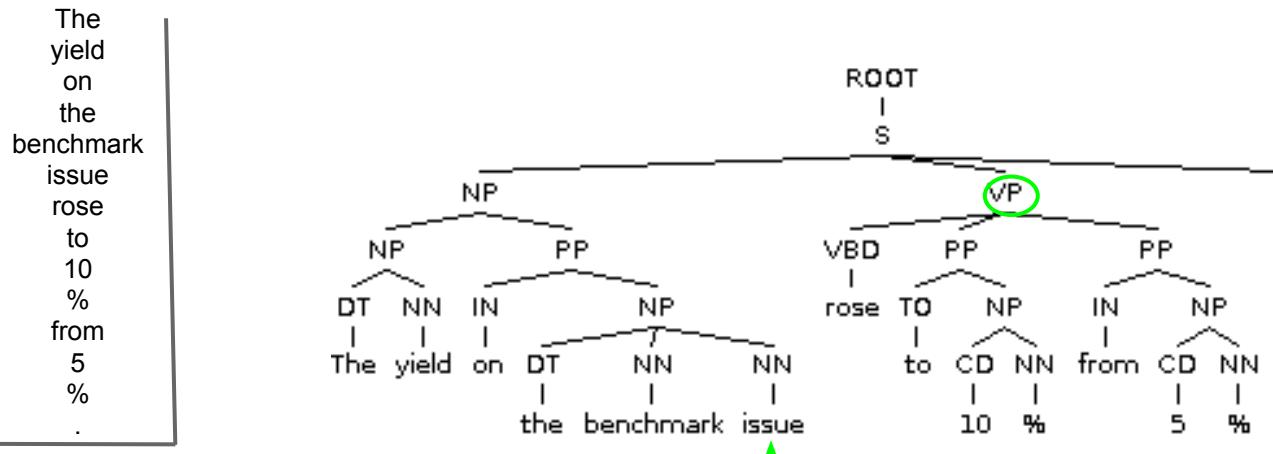
Encoder LSTM
(reverse input)

Parsing Example



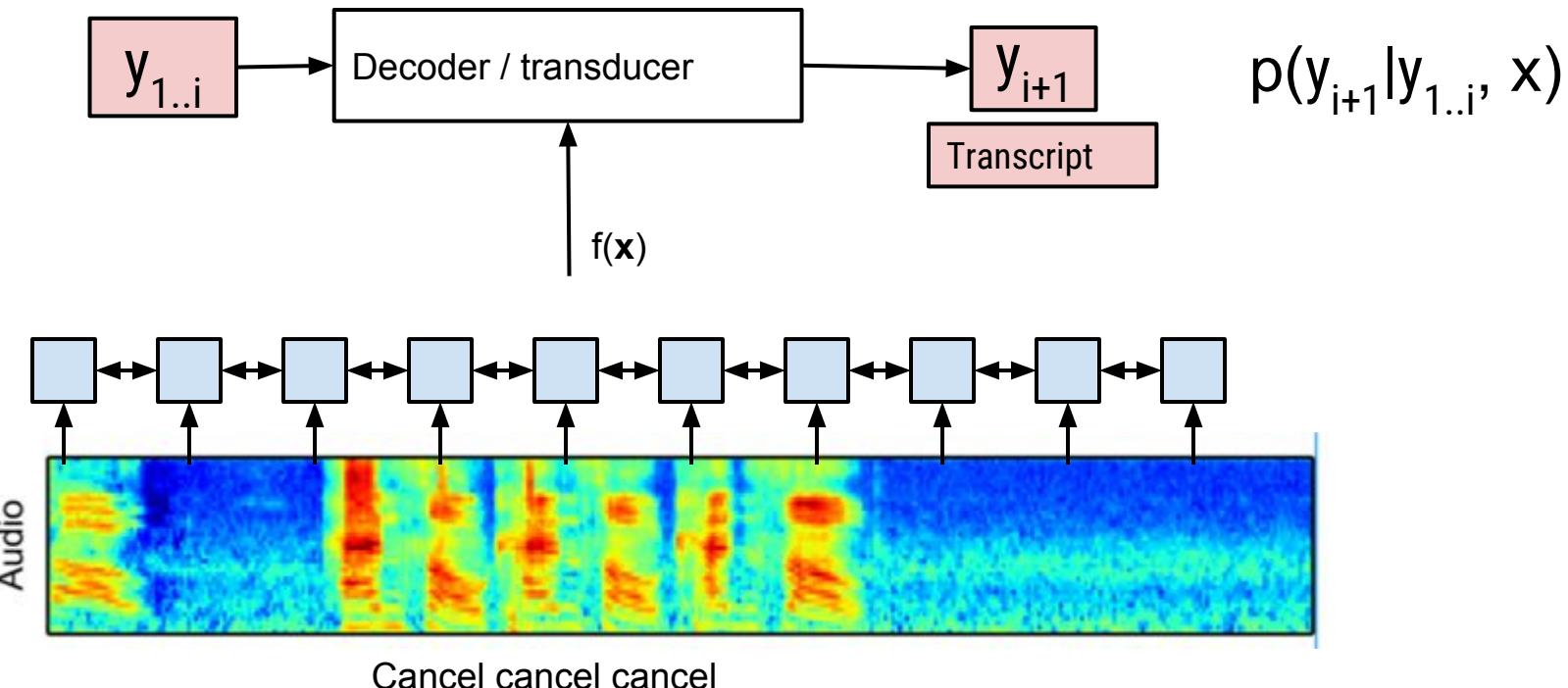
Encoder LSTM
(reverse input)

Parsing Example

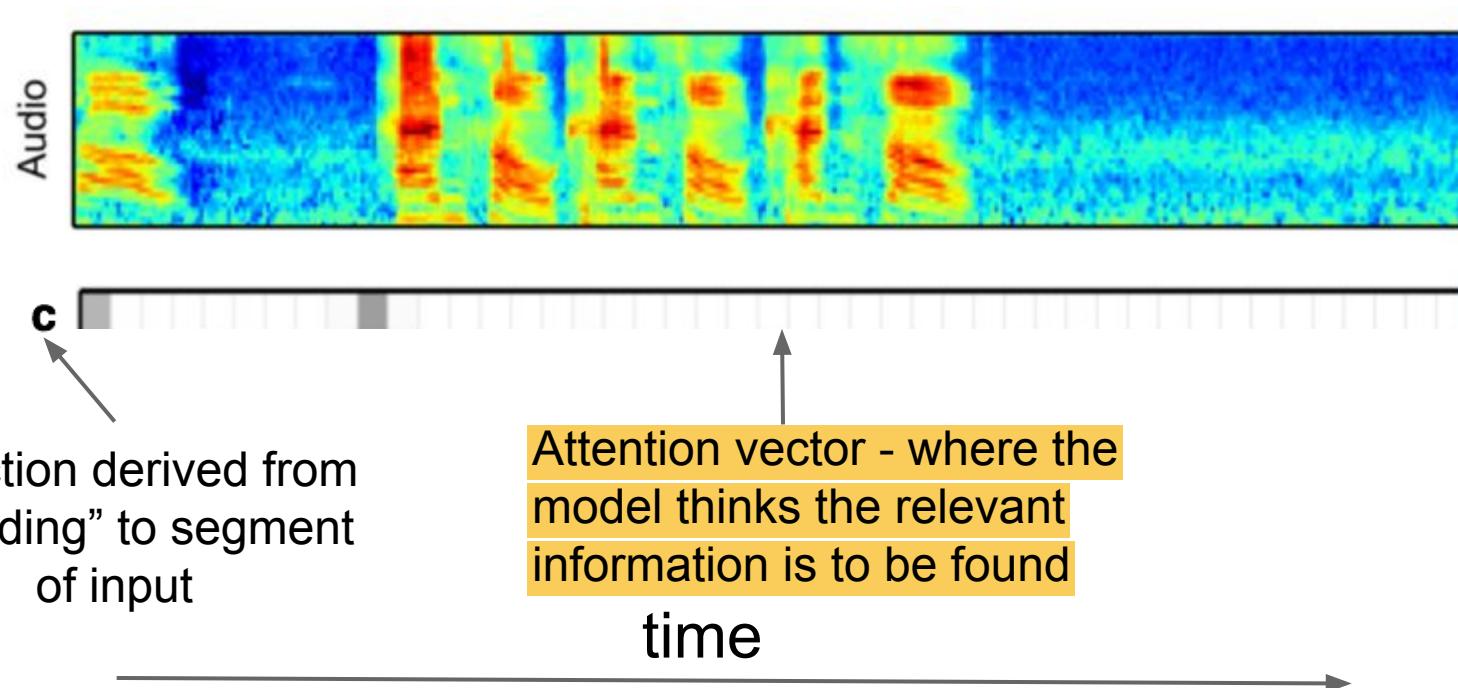


Speech Recognition

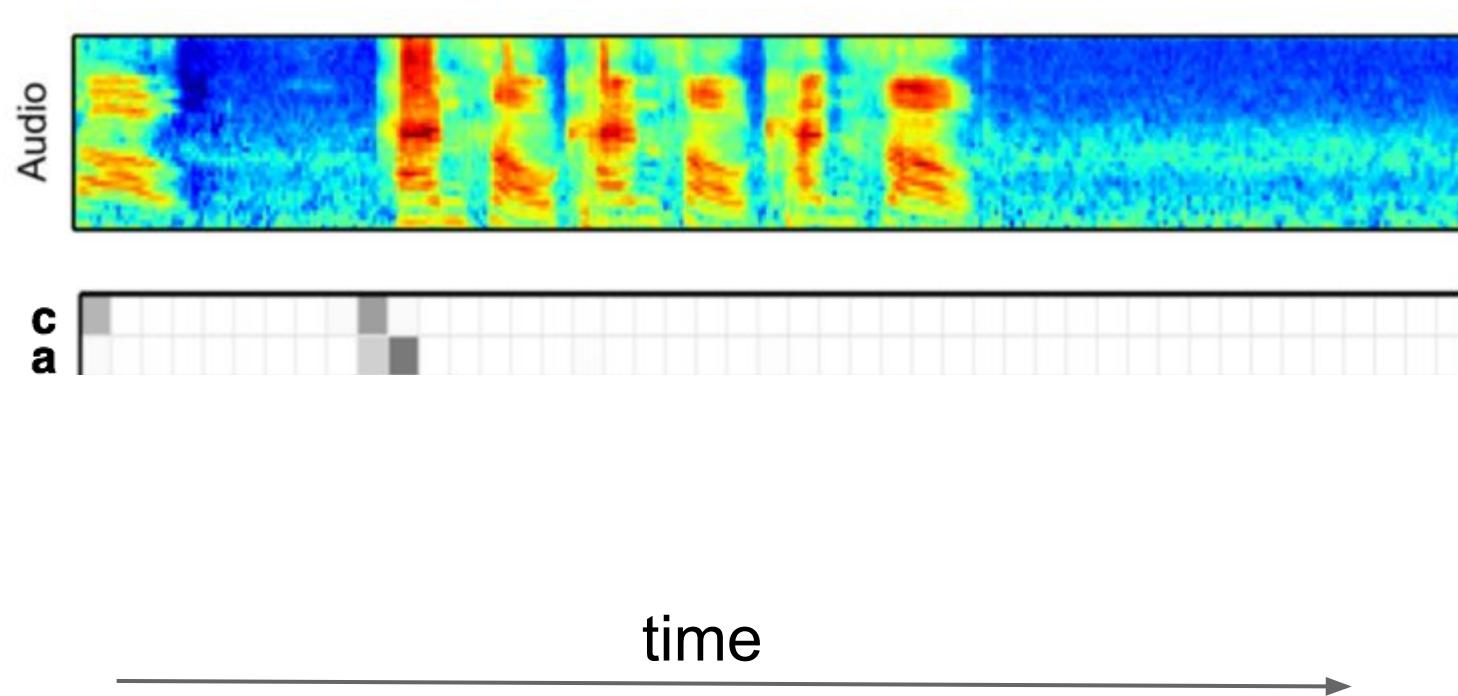
spectrum → characters



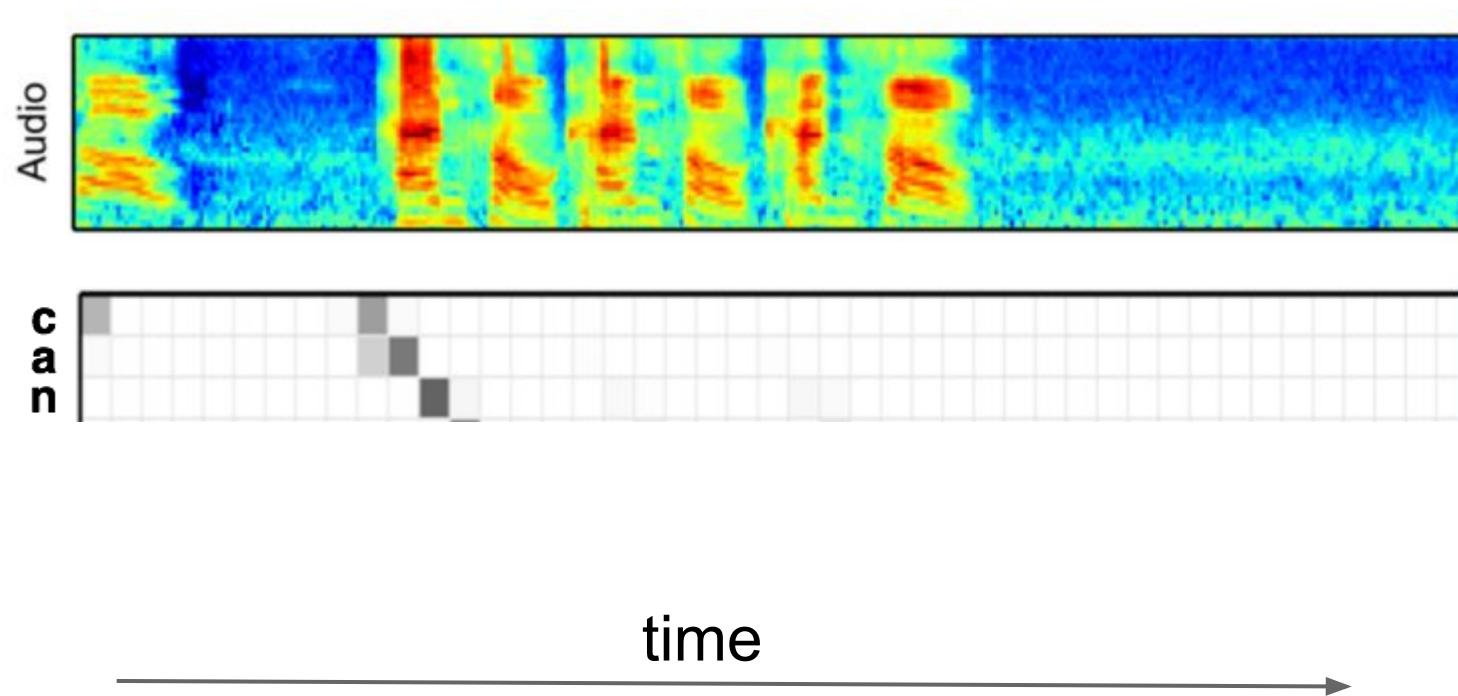
Attention Example



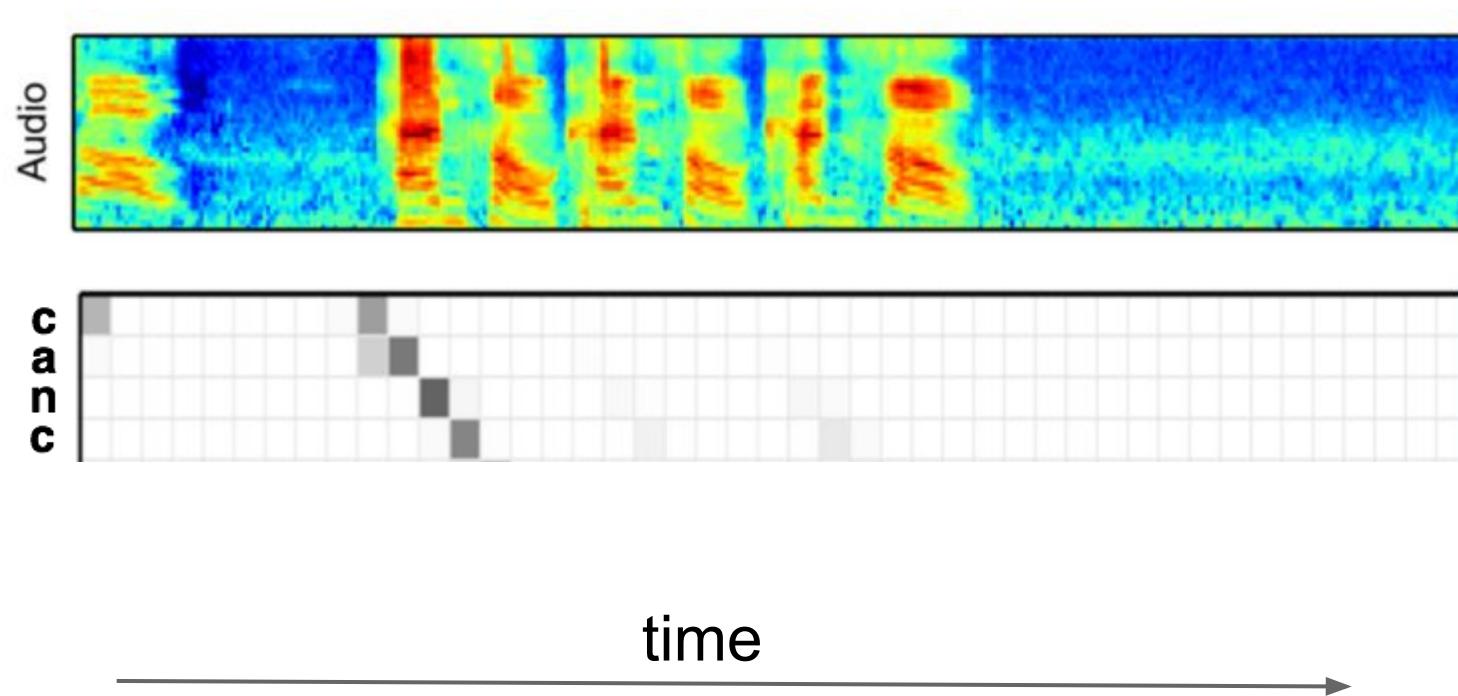
Attention Example



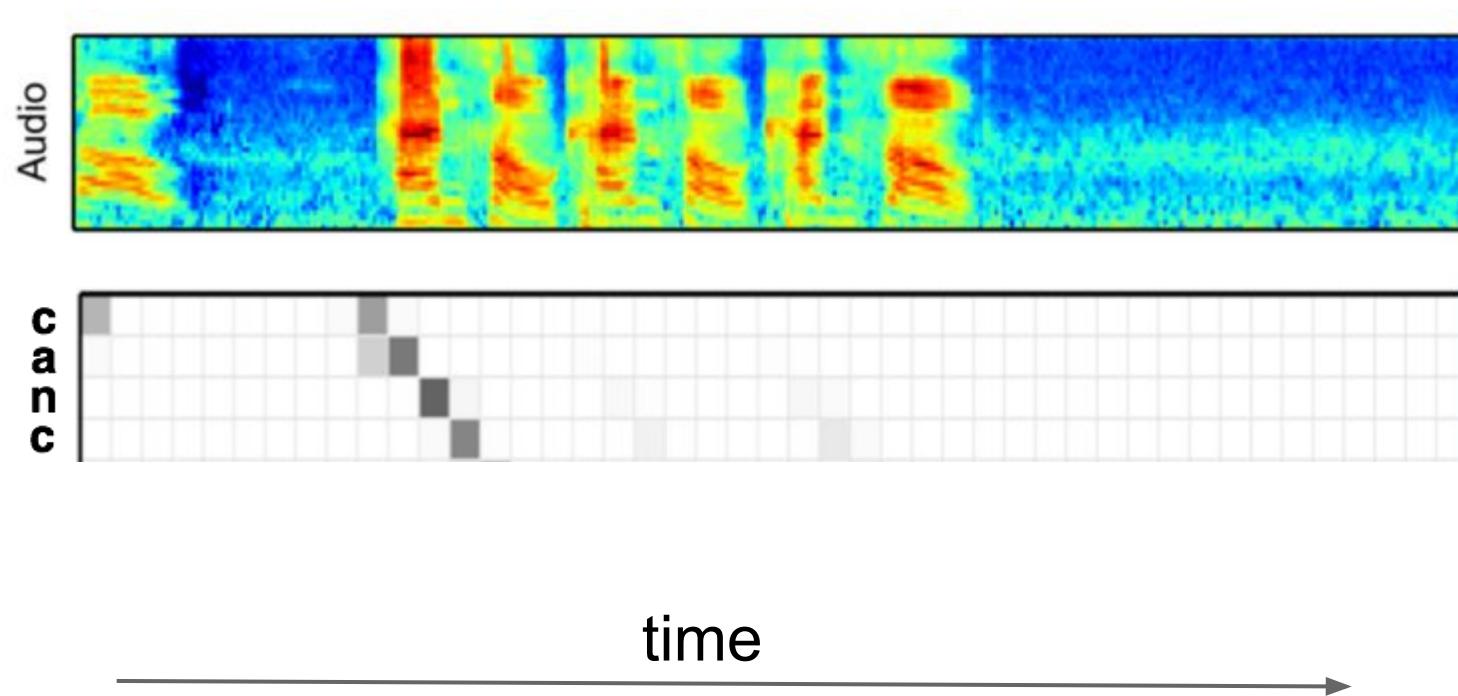
Attention Example



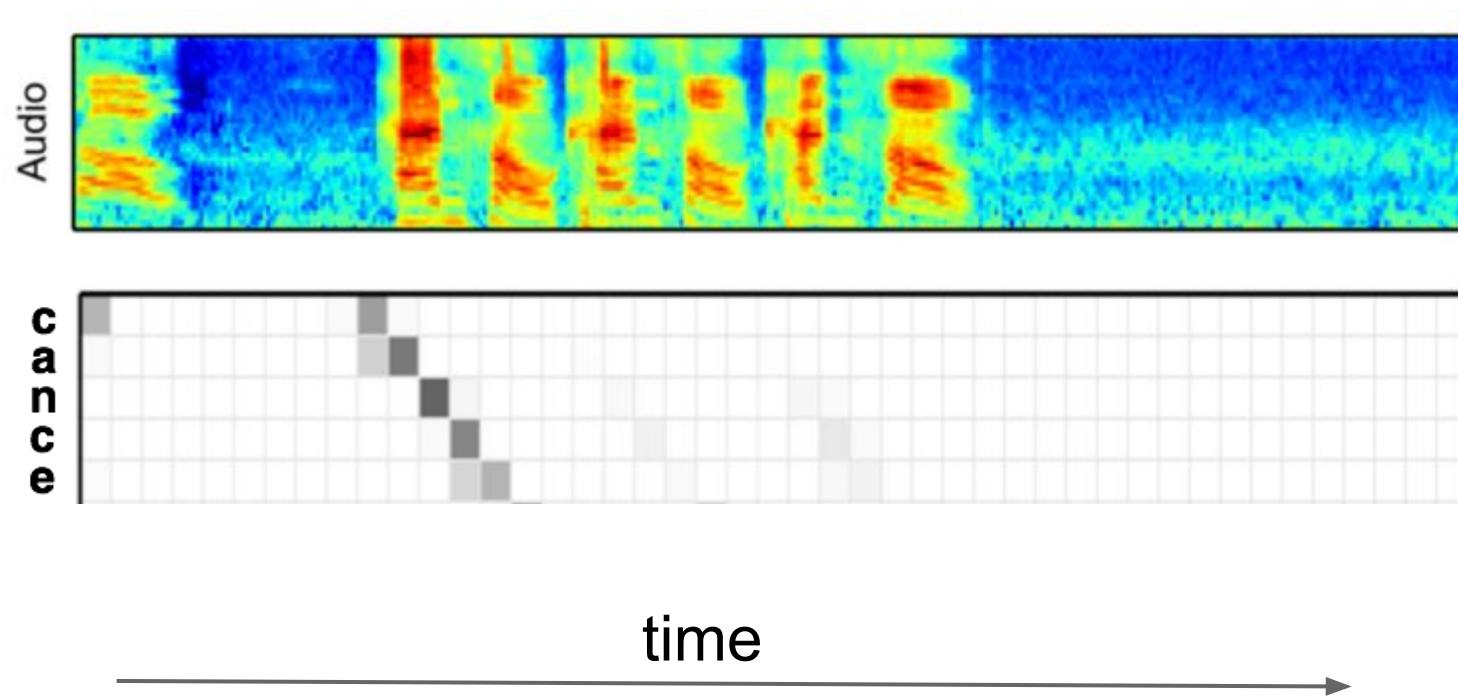
Attention Example



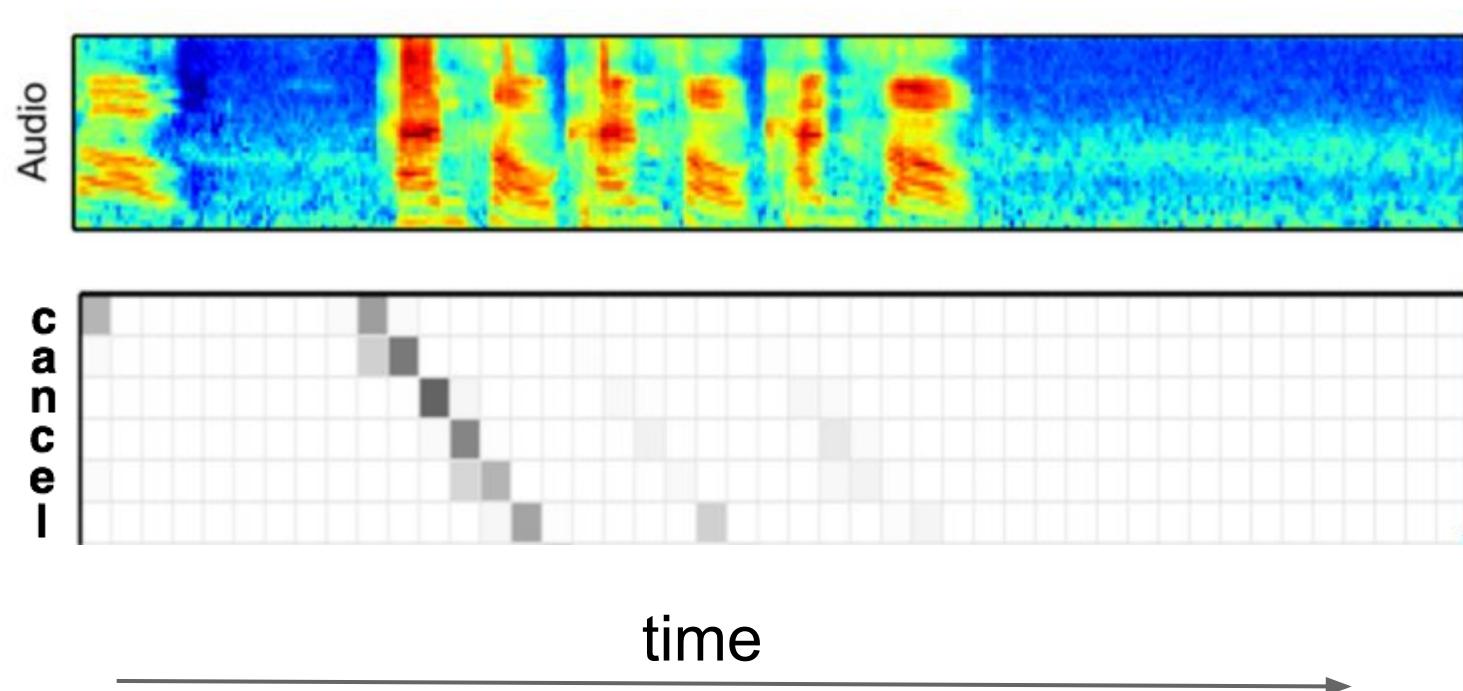
Attention Example



Attention Example

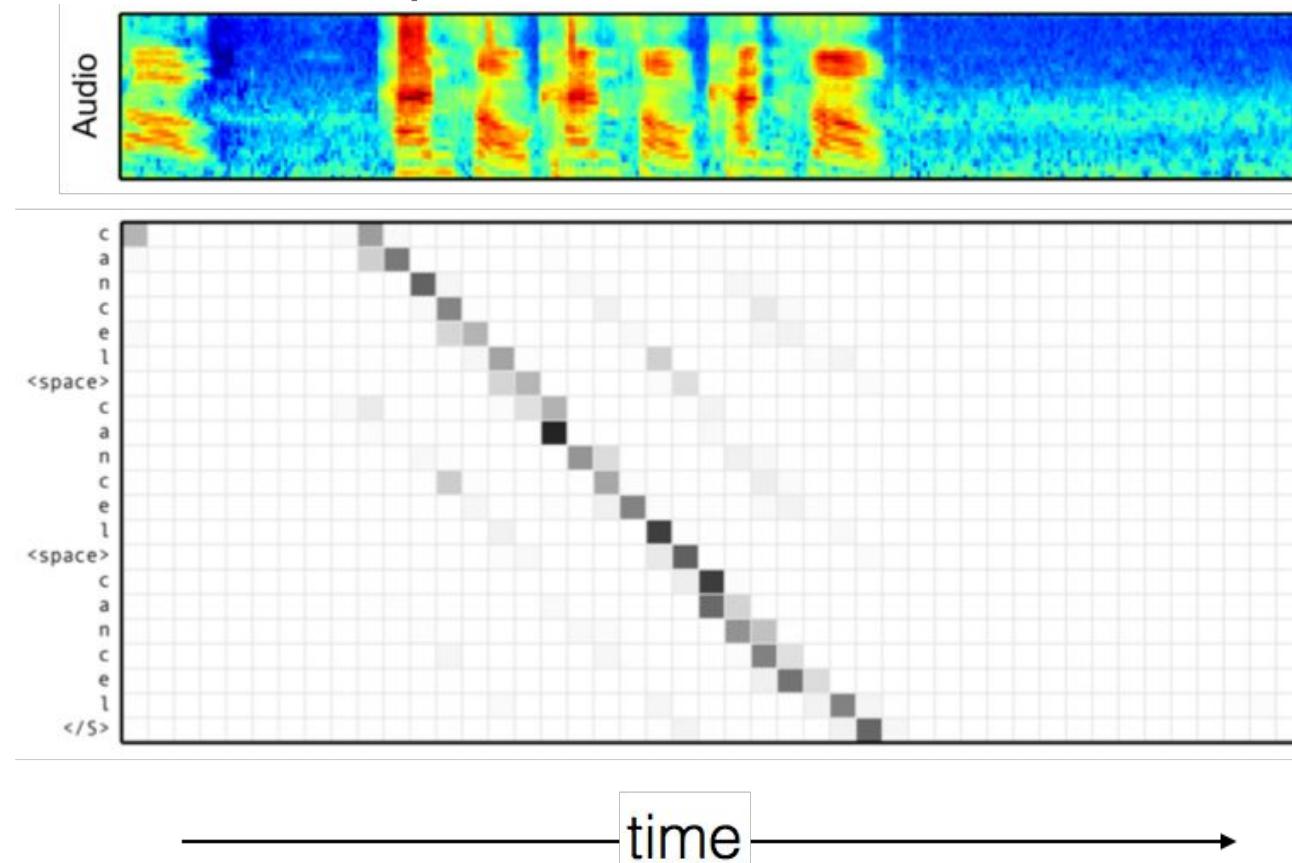


Attention Example



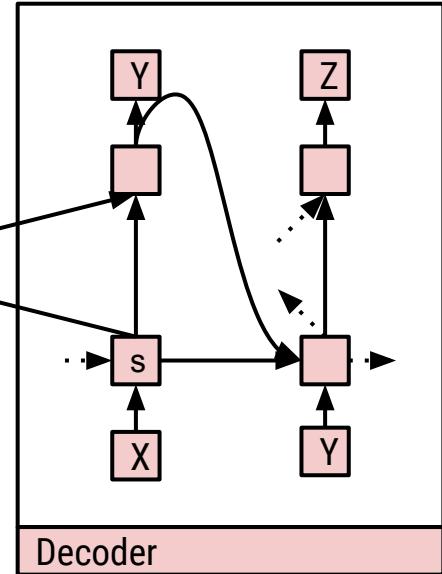
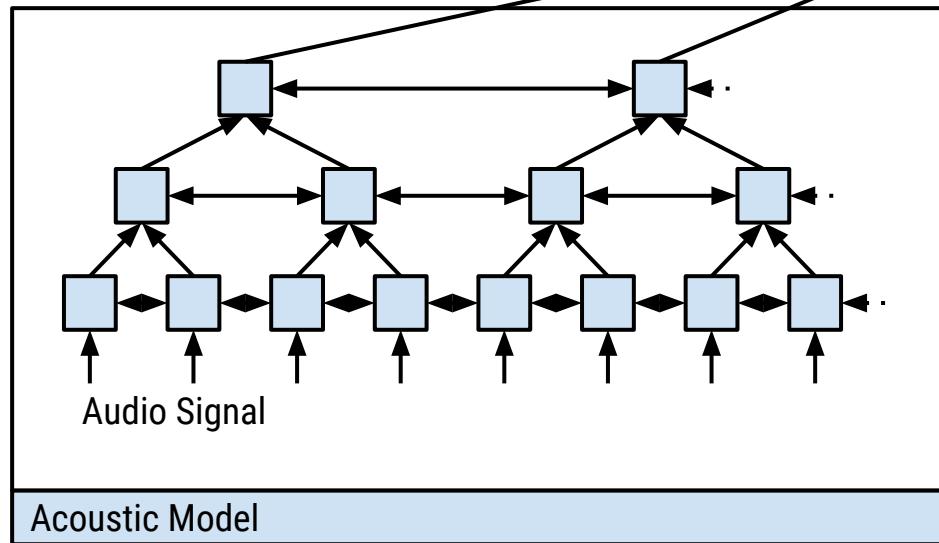
Attention Example

overcomes some difficulties of HMM, still not good enough for longer sequences



Listen Attend and Spell (LAS)

- Reducing time resolution with a pyramidal encoder



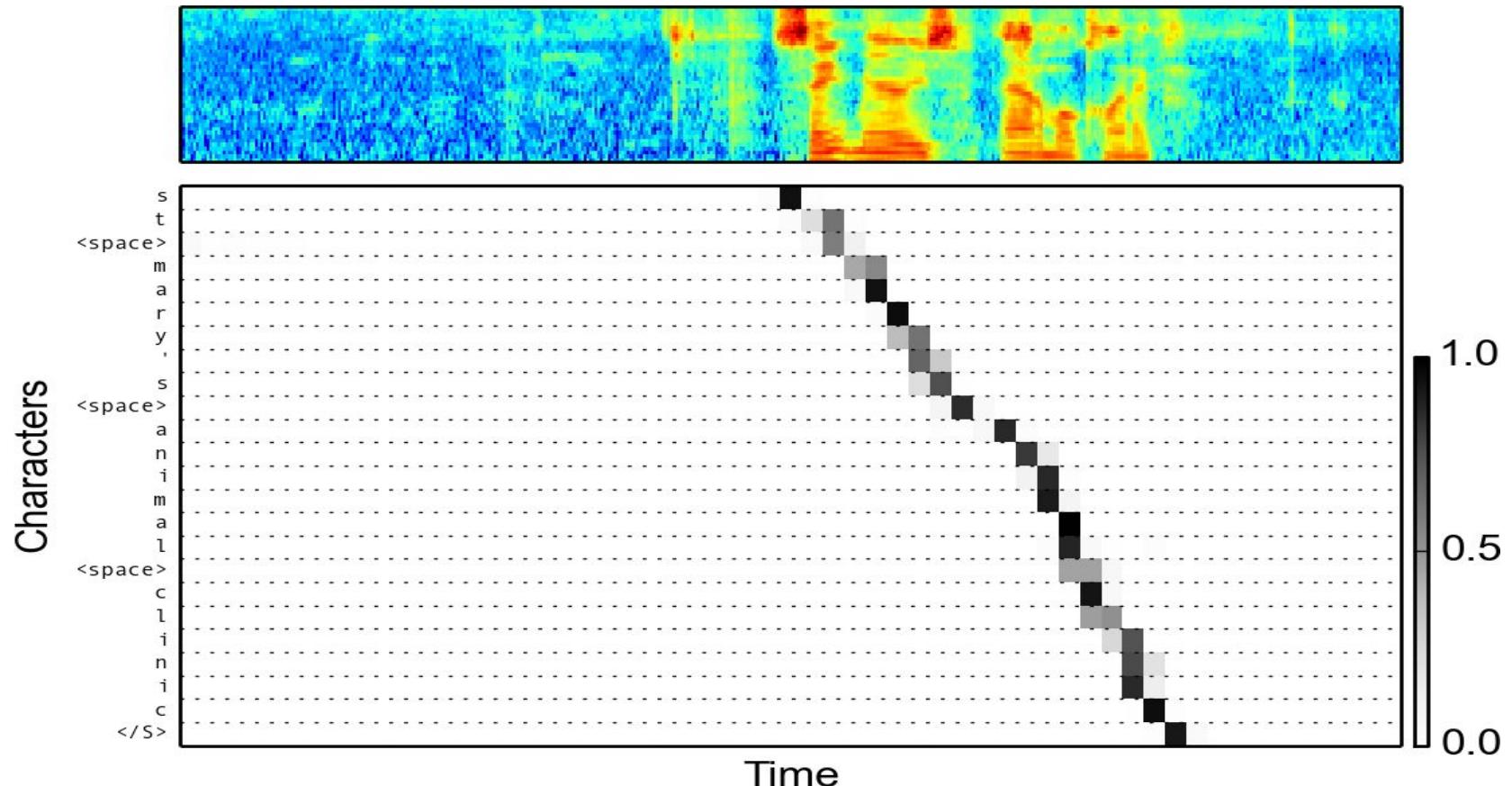
LAS - Multimodal outputs

word error rate

Beam	Text	LogProb	WER
Truth	call aaa roadside assistance	-	-
1	call aaa roadside assistance	-0.5740	0.00
2	call triple a roadside assistance	-1.5399	50.0
3	call trip way roadside assistance	-3.5012	50.0
4	call xxx roadside assistance	-4.4375	25.0

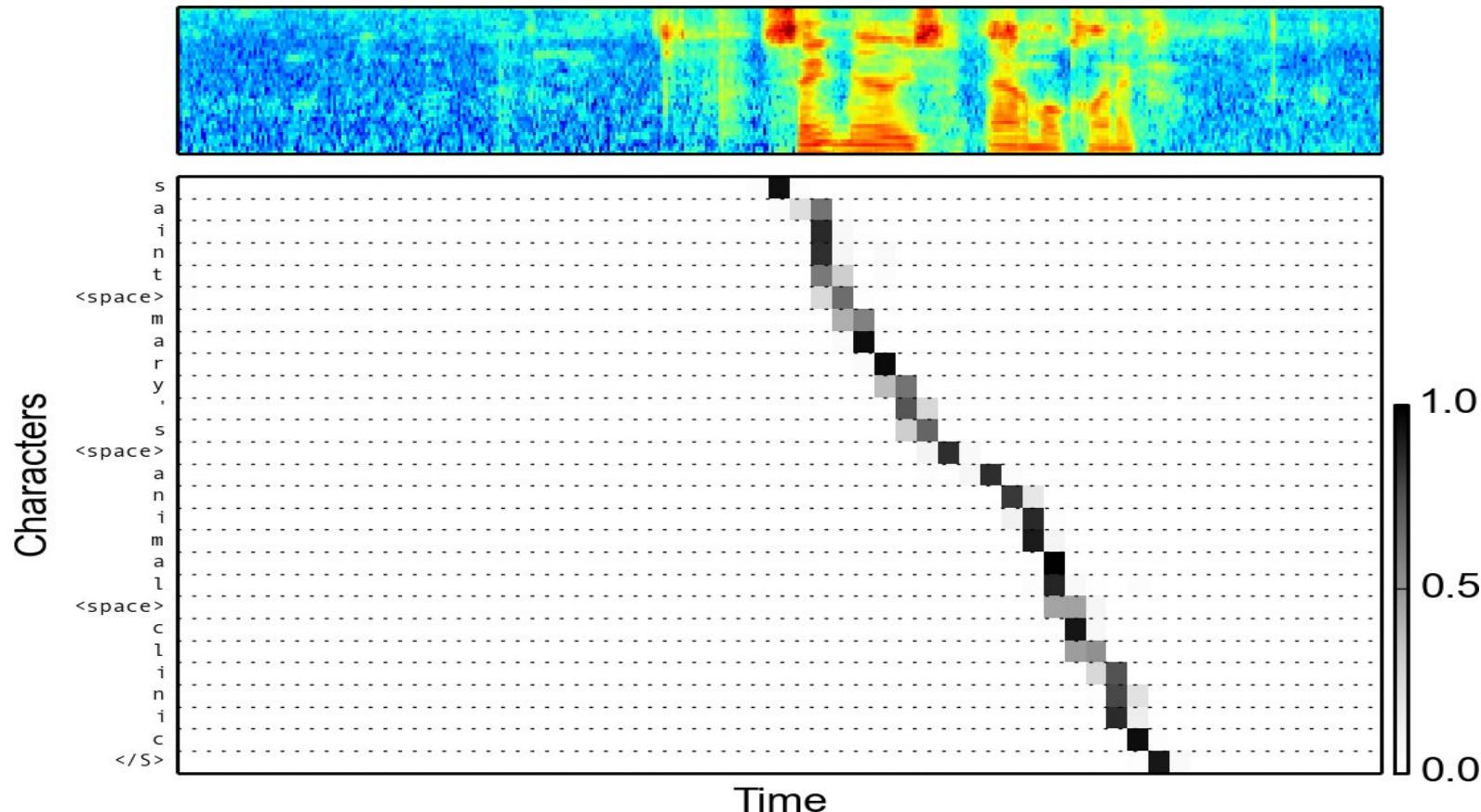
Very different outputs for same input!

LAS - Impact of Autoregressive Inputs on Attention



LAS - Impact of Autoregressive Inputs on Attention

model outputs differently at different time



LAS - Results

Model	Clean WER	Noisy WER
CLDNN-HMM (baseline)	8	8.9
Listen Attend and Spell (LAS)	14	16.5
LAS + external language model	10.3	12.0

Comparable to some of our best models without extensive engineering !