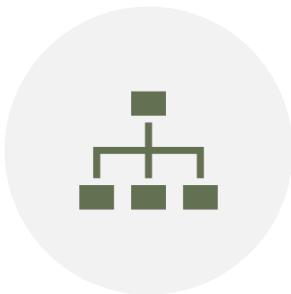


Reminders



HOMEWORK 9 OR PROJECT
MILESTONE 1 ARE DUE TONIGHT
BY 11:59PM



QUIZ ON CHAPTER 28 IS DUE
MONDAY.

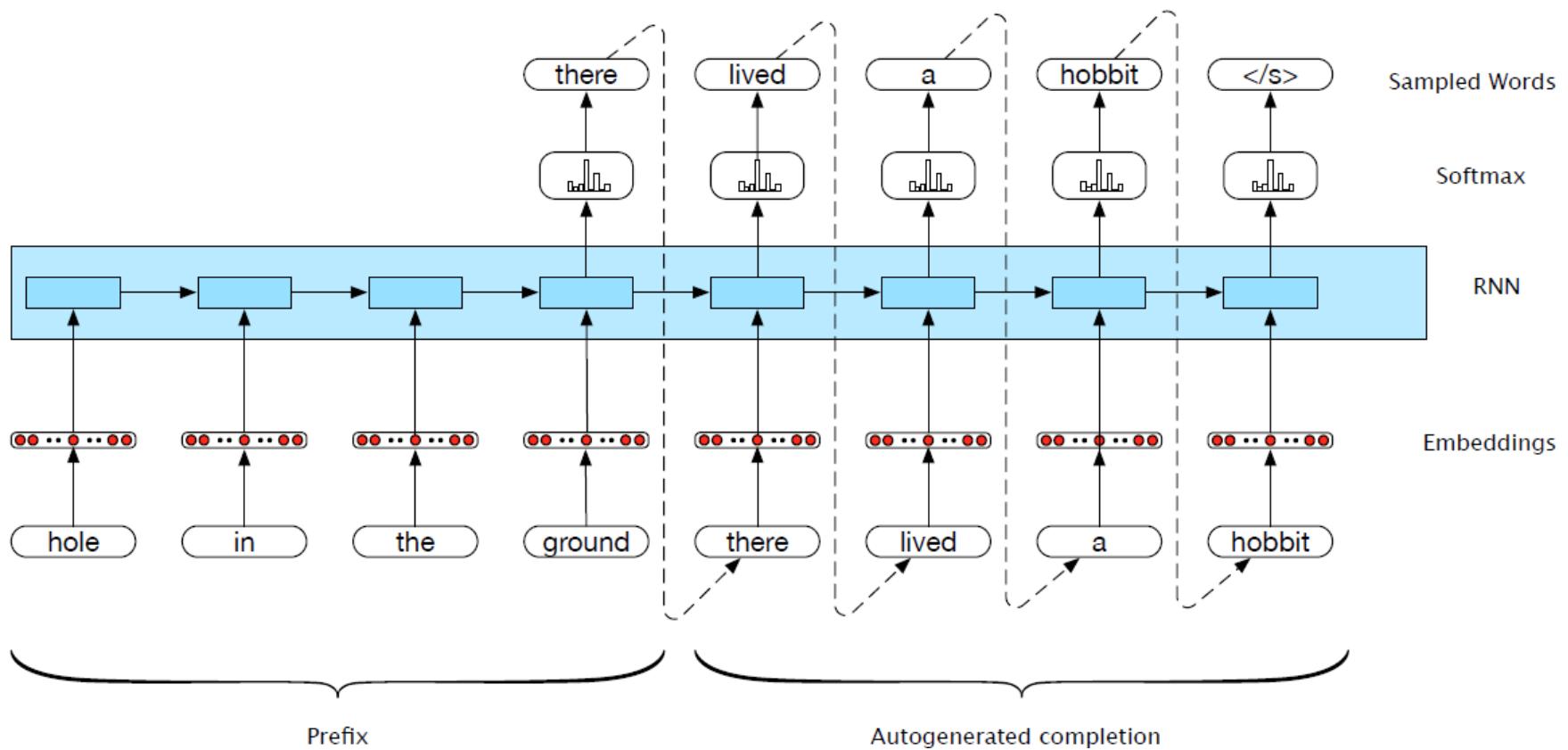


HW10 ON NEURAL MACHINE
TRANSLATION WILL BE RELEASED
SOON. MILESTONE 2 IS READY.

Encoder-Decoder Models

JURAFSKY AND MARTIN CHAPTER 10

Generation with prefix

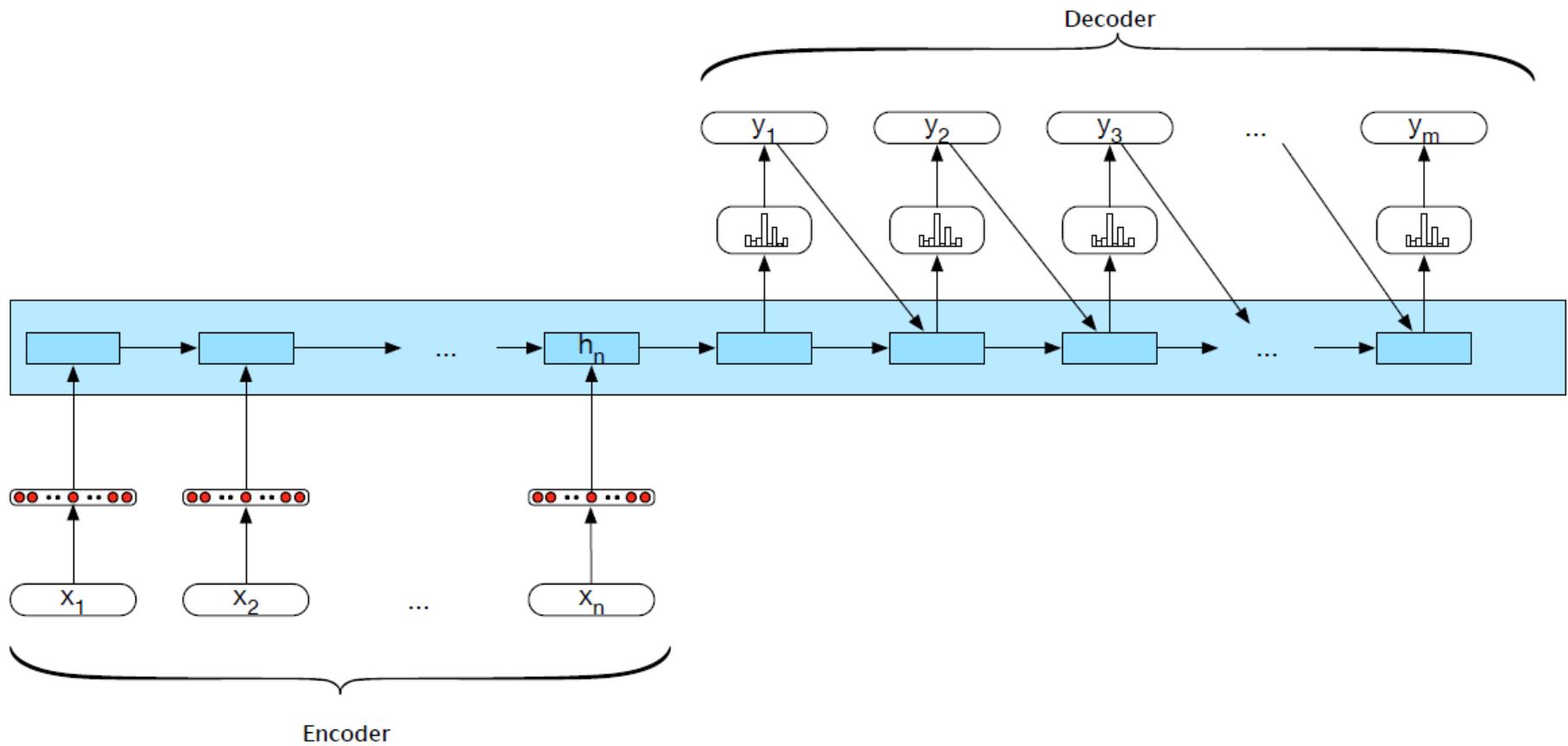


Encoder-Decoder Networks

We can abstract away from the task of MT to talk about the general **encoder-decoder architecture**:

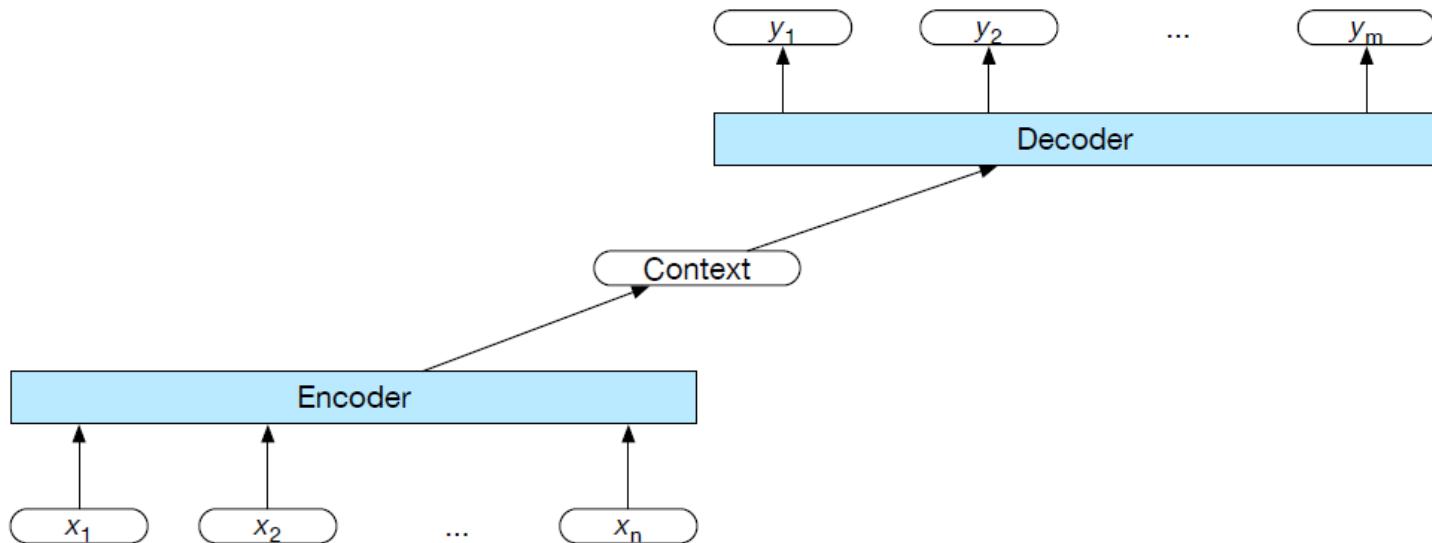
1. An **encoder** takes an input sequence x^n_1 , and generates a corresponding sequence of contextualized representations, h^n_1 .
2. A **context vector**, c , is a function of h^n_1 , and conveys the essence of the input to the decoder.
3. A **decoder** accepts c as input and generates an arbitrary length sequence of hidden states h^m_1 , from which can be used to create a corresponding sequence of output states y^m_1 .

Encoder-decoder networks



Encoder-decoder networks

- An encoder that accepts an input sequence and generates a corresponding sequence of contextualized representations
- A context vector that conveys the essence of the input to the decoder
- A decoder, which accepts context vector as input and generates an arbitrary length sequence of hidden states, from which a corresponding sequence of output states can be obtained

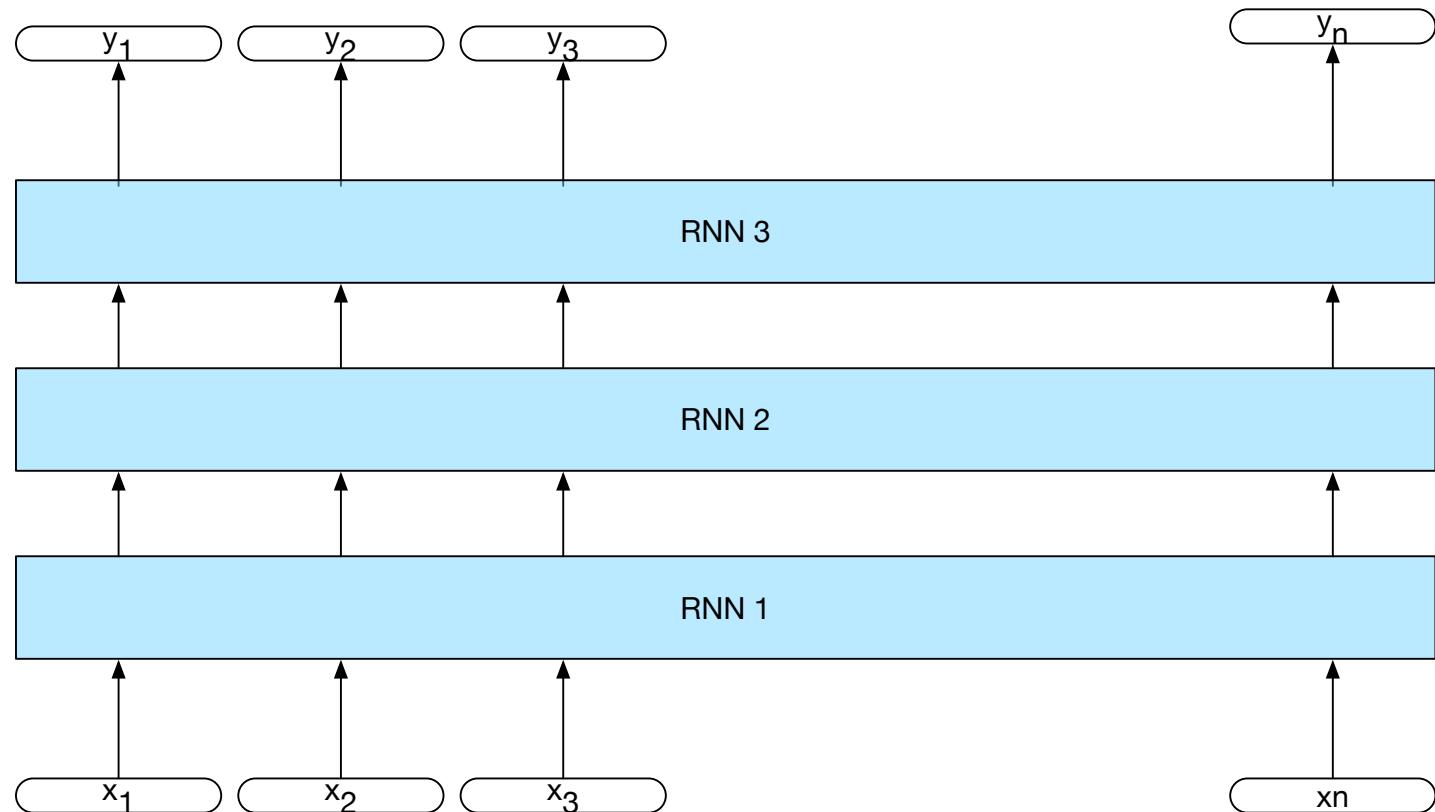


Encoder

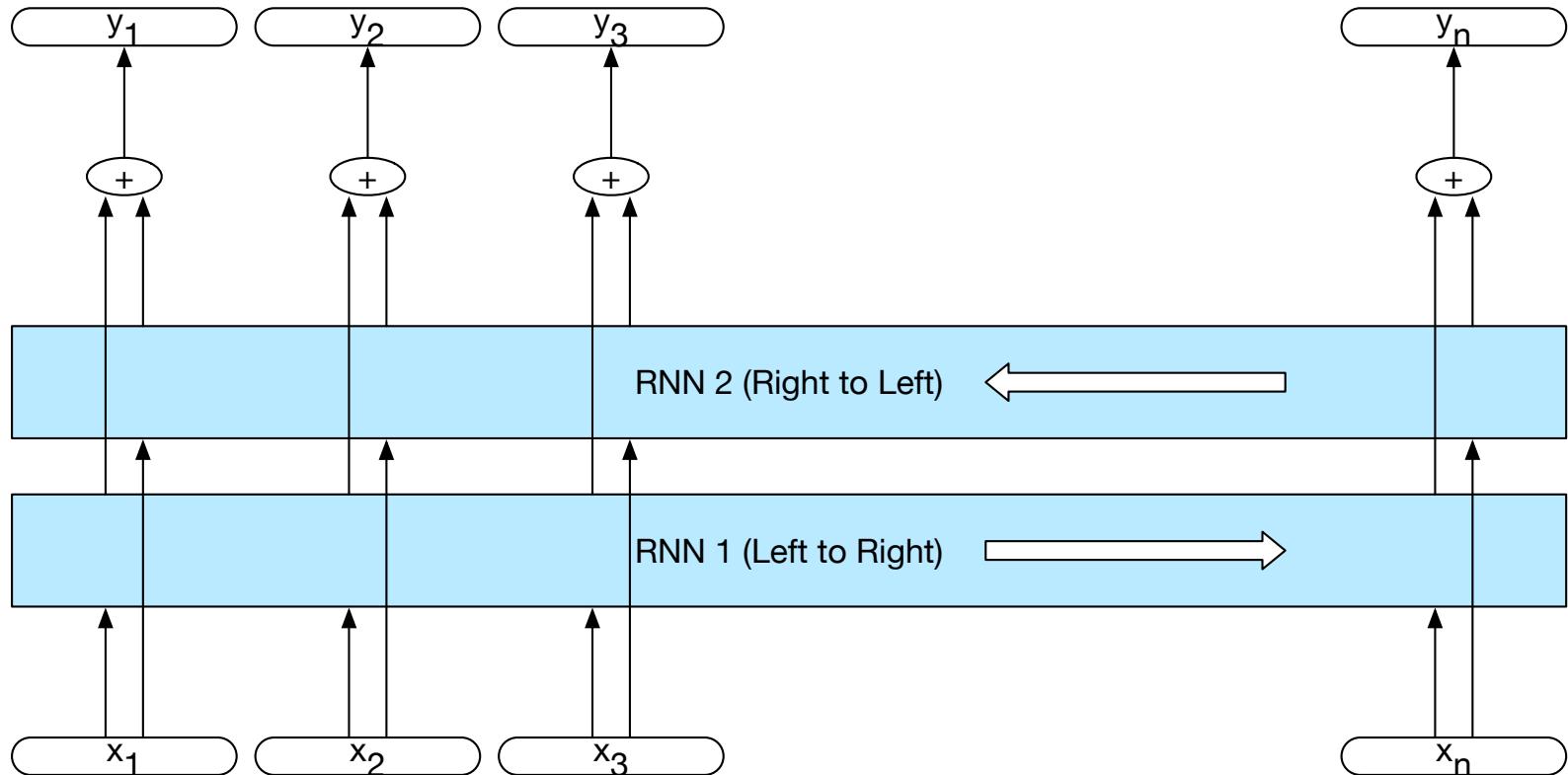
Pretty much any kind of RNN or its variants can be used as an encoder. Researchers have used simple RNNs, LSTMs, GRUs, or even convolutional networks.

A widely used encoder design makes use of stacked Bi-LSTMs where the hidden states from top layers from the forward and backward passes are concatenated

Stacked RNNs



Bidirectional RNNs



Decoder

For the decoder, autoregressive generation is used to produce an output sequence, an element at a time, until an end-of-sequence marker is generated.

This incremental process is guided by the context provided by the encoder as well as any items generated for earlier states by the decoder.

Encoder

$$\begin{aligned} c &= h_n^e \\ h_0^d &= c \end{aligned}$$

Decoder

$$\begin{aligned} h_t^d &= g(\hat{y}_{t-1}, h_{t-1}^d) \\ z_t &= f(h_t^d) \\ y_t &= \text{softmax}(z_t) \end{aligned}$$

Decoder Weaknesses

In early encoder-decoder approaches, the context vector c was only directly available at the beginning of the generation process.

This meant that its influence became less-and-less important as the output sequence was generated.

One solution is to make c available at each step in the decoding process, when generating the hidden states in the decoder

$$h_t^d = g(\hat{y}_{t-1}, h_{t-1}^d, c)$$

and while producing the generated output.

$$y_t = \text{softmax}(\hat{y}_{t-1}, z_t, c)$$

Choosing the best output

For neural generation, where we are trying to generate novel outputs, we can simply sample from the softmax distribution.

In MT where we're looking for a specific output sequence, sampling isn't appropriate and would likely lead to some strange output.

Instead we choose the most likely output at each time step by taking the argmax over the softmax output

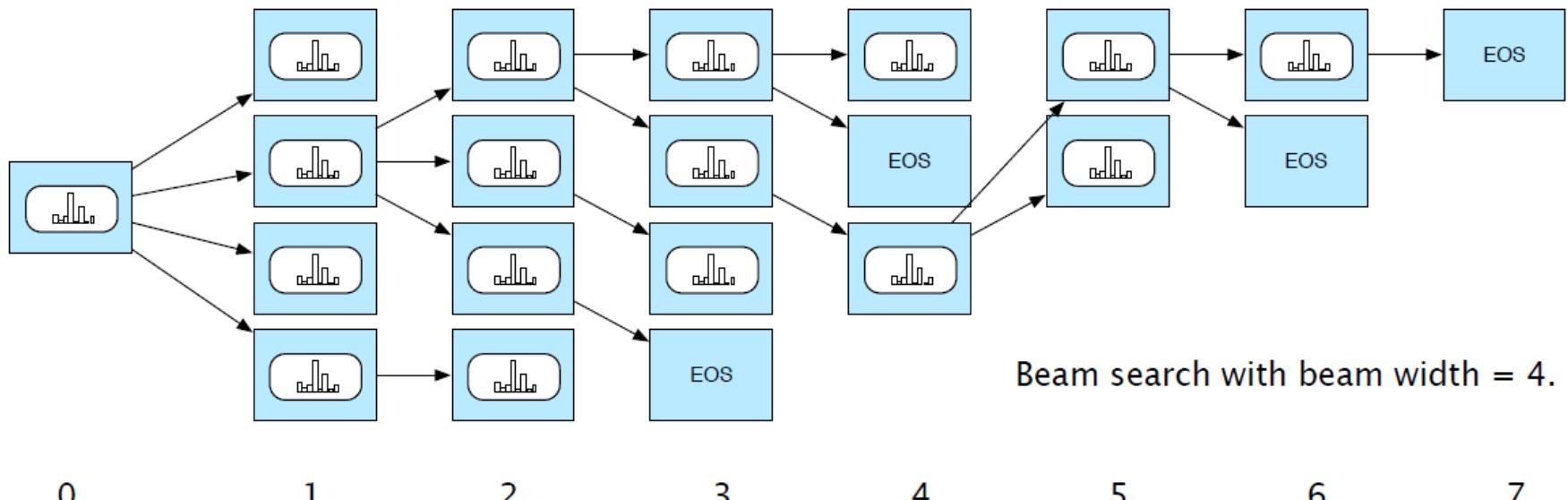
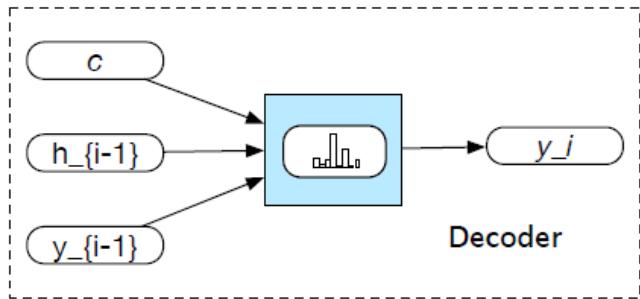
$$\hat{y} = \text{argmax} P(y_i | y_{<i})$$

Beam search

In order to systematically explore the space of possible outputs for applications like MT, we need to control the exponential growth of the search space.

Beam search: combining a breadth-first-search strategy with a heuristic filter that scores each option and prunes the search space to stay within a fixed-size memory footprint, called the beam width

Beam search



0 1 2 3 4 5 6 7

Attention

Weaknesses of the context vector:

- Only directly available at the beginning of the process and its influence will wane as the output sequence is generated
- Context vector is a function (e.g. last, average, max, concatenation) of the hidden states of the encoder. This approach loses useful information about each of the individual encoder states

Potential solution: **attention mechanism**

Attention

- Replace the static context vector with one that is dynamically derived from the encoder hidden states at each point during decoding
- A new context vector is generated at each decoding step and takes all encoder hidden states into derivation
- This context vector is available to decoder hidden state calculations

$$h_i^d = g(\hat{y}_{i-1}, h_{i-1}^d, c_i)$$

Attention

- To calculate c_i , first find relevance of each encoder hidden state to the decoder state. Call it $score(h_{i-1}^d, h_j^e)$ for each encoder state j
- The $score$ can simply be dot product,

$$score(h_{i-1}^d, h_j^e) = h_{i-1}^d \cdot h_j^e$$

Attention

- The score can also be parameterized with weights

$$score(h_{i-1}^d, h_j^e) = h_{i-1}^d W_s h_j^e$$

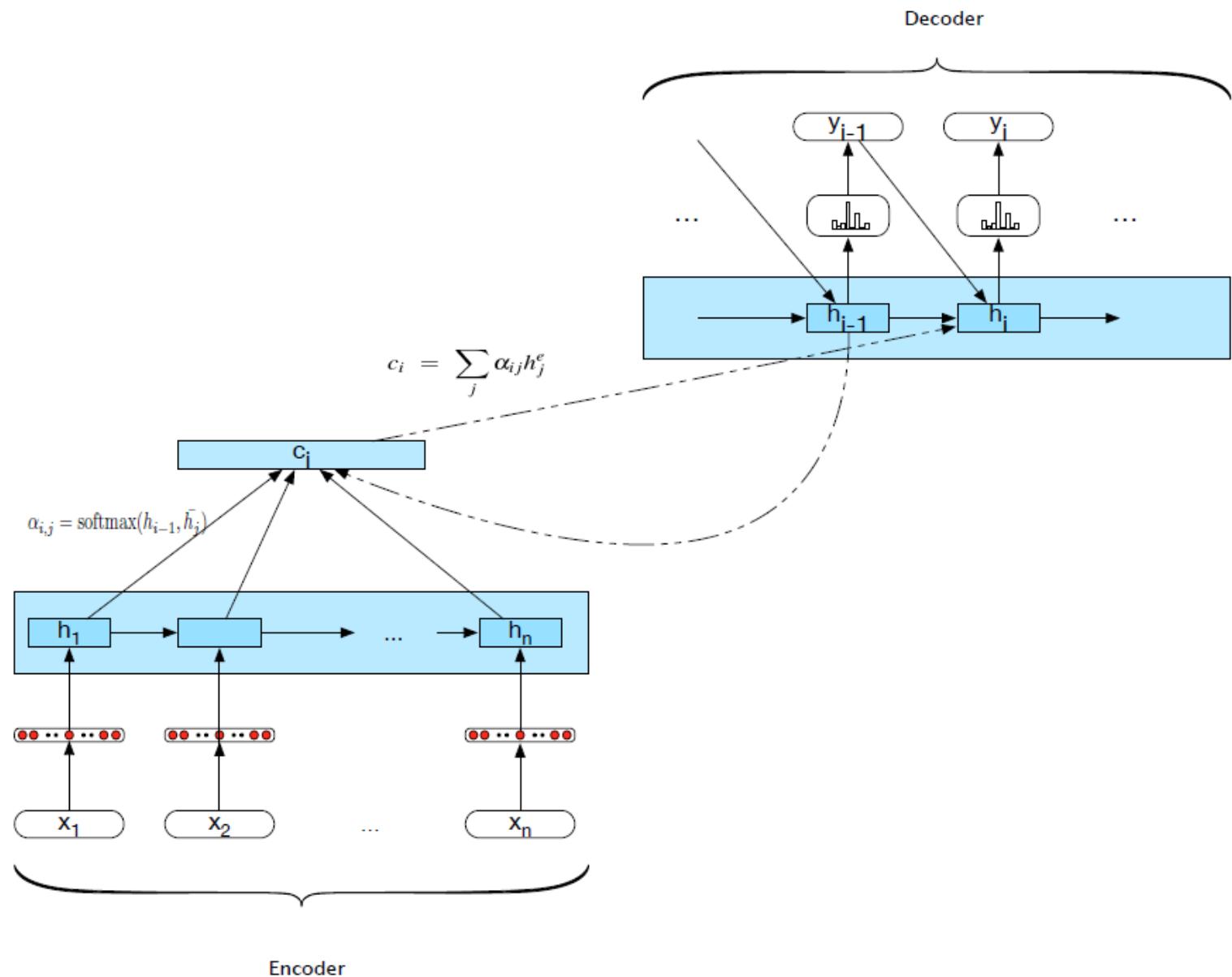
- Normalize them with a softmax to create a vector of weights $\alpha_{i,j}$ that tells us the proportional relevance of each encoder hidden state j to the current decoder state i

$$\alpha_{i,j} = \text{softmax}(score(h_{i-1}^d, h_j^e)) \quad \forall j \in e$$

- Finally, context vector is the weighted average of encoder hidden states

$$c_i = \sum_j \alpha_{i,j} h_j^e$$

Attention mechanism



Applications of Encoder-Decoder Networks

- Text summarization
- Text simplification
- Question answering
- Image captioning
- And more. What do those tasks have in common?

Neural Machine Translation

SLIDES FROM GRAHAM NEUBIG, CMU

NEURAL MACHINE TRANSLATION AND
SEQUENCE-TO-SEQUENCE MODELS: A TUTORIAL

Machine Translation

Translation from one language to another

I'm giving a talk at University of Pennsylvania



ペンシルベニア大学で講演をしています。

Long-distance Dependencies

Agreement in number, gender, etc.

1. **He** does not have very much confidence in **himself**.
2. **She** does not have very much confidence in **herself**.

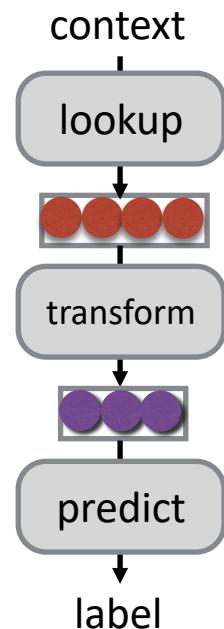
Selectional preference:

1. The **reign** has lasted as long as the life of the **queen**.
2. The **rain** has lasted as long as the life of the **clouds**.

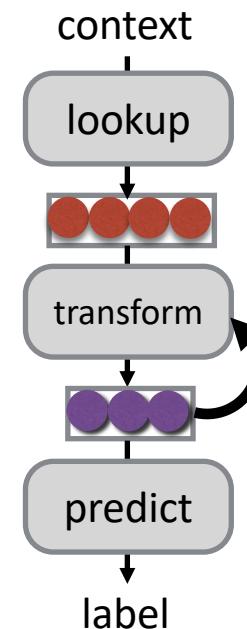
Recurrent Neural Networks

Tools to “remember” information

Feed-forward NN

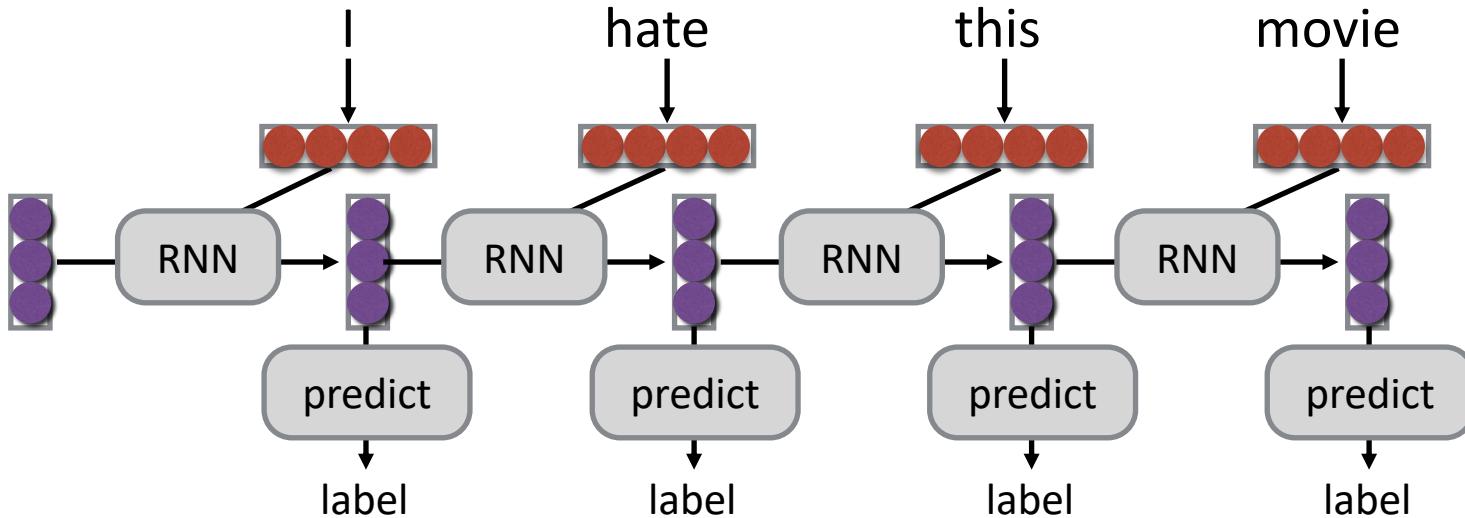


Recurrent NN

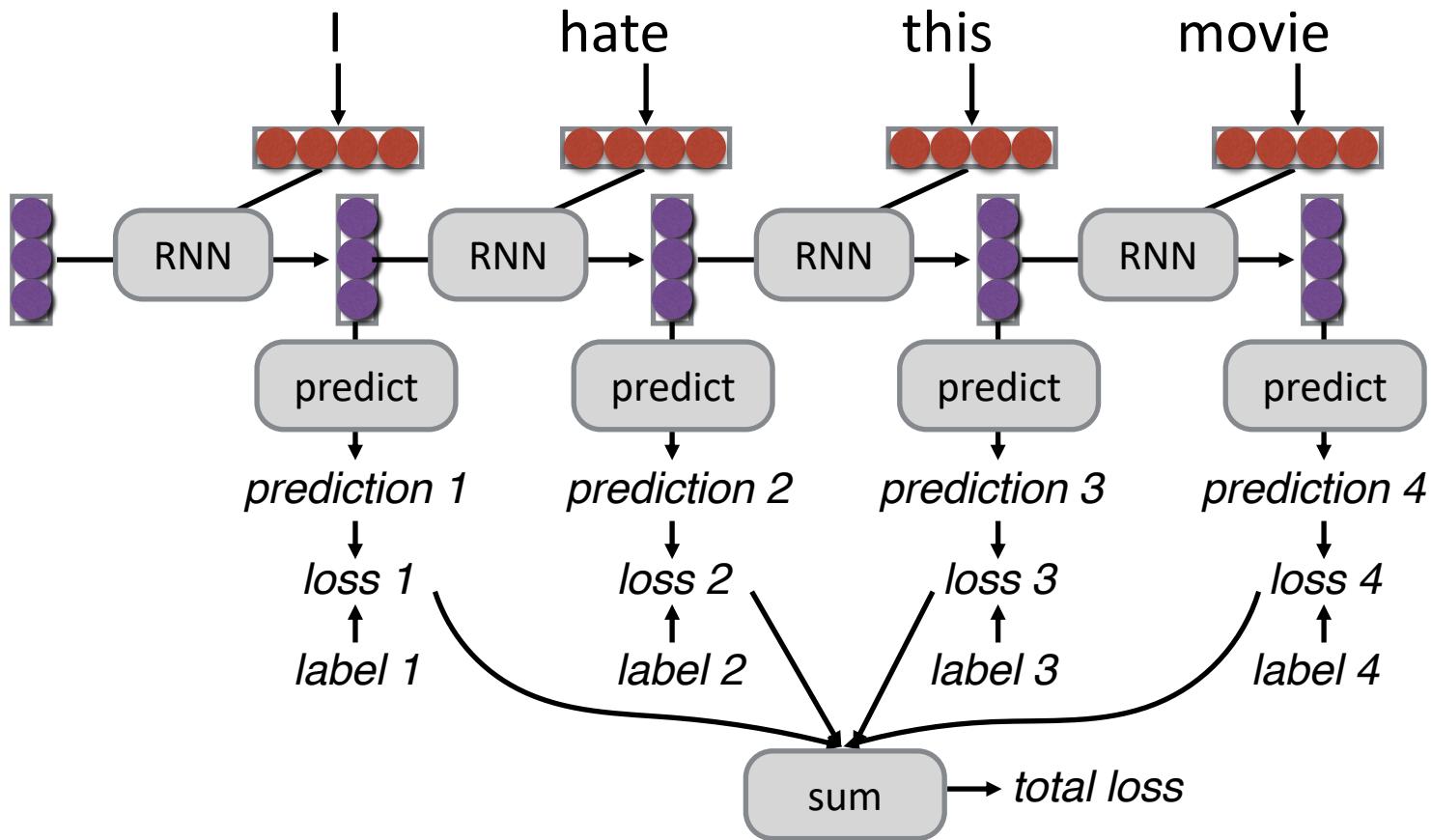


Unrolling in Time

What does processing a sequence look like?

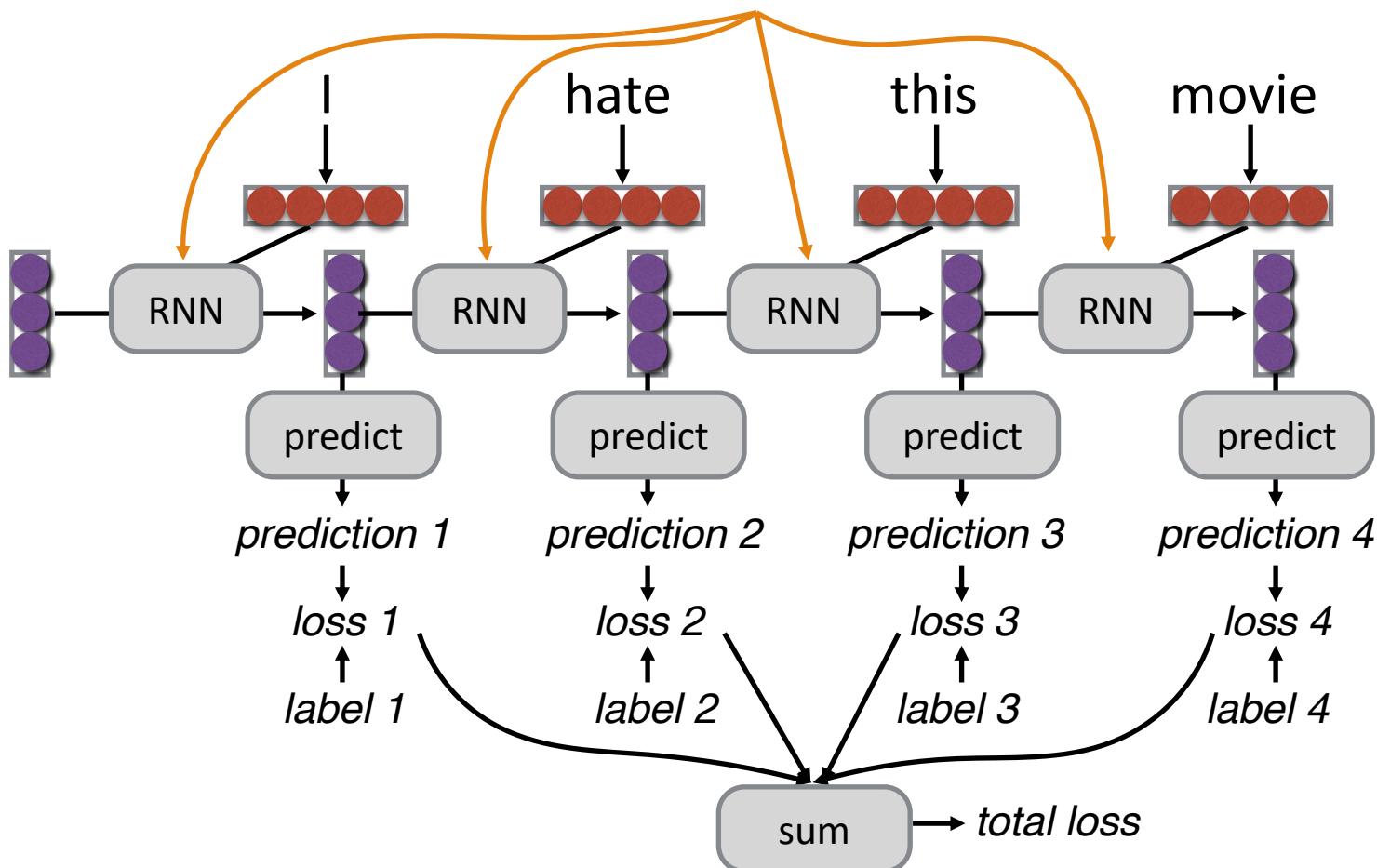


Training RNNs



Parameter Tying

Parameters are shared! Derivatives are accumulated.



What Can RNNs Do?

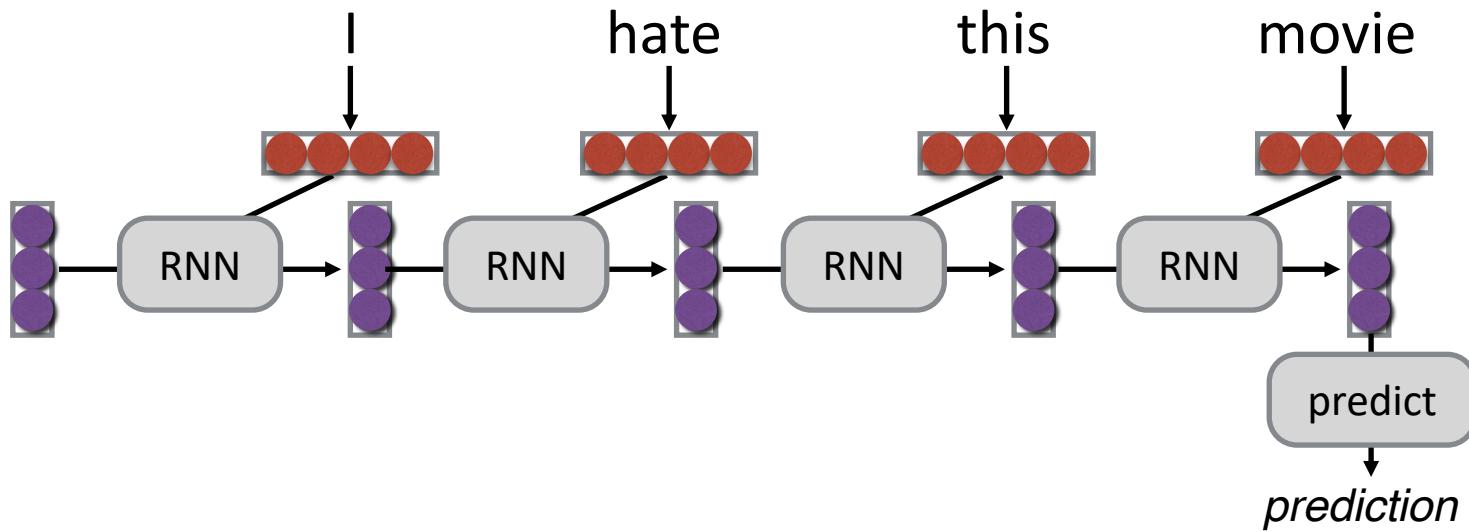
Represent a sentence

- Read whole sentence, make a prediction

Represent a context within a sentence

- Read context up until that point

Representing Sentences

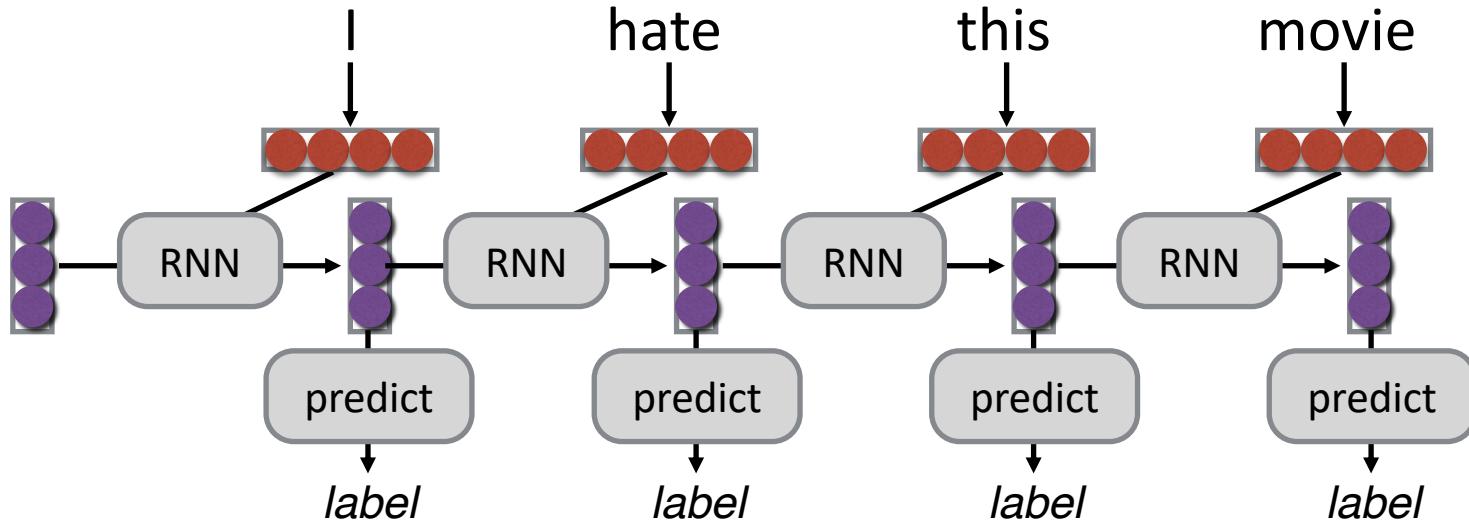


Sentence classification

Conditioned generation

Retrieval

Representing Contexts



Tagging

Language Modeling

Calculating Representations for Parsing, etc.

Language Models

Language models are generative models of text

$$s \sim P(x)$$



“The Malfoys!” said Hermione.

Harry was watching him. He looked like Madame Maxime. When she strode up the wrong staircase to visit himself.

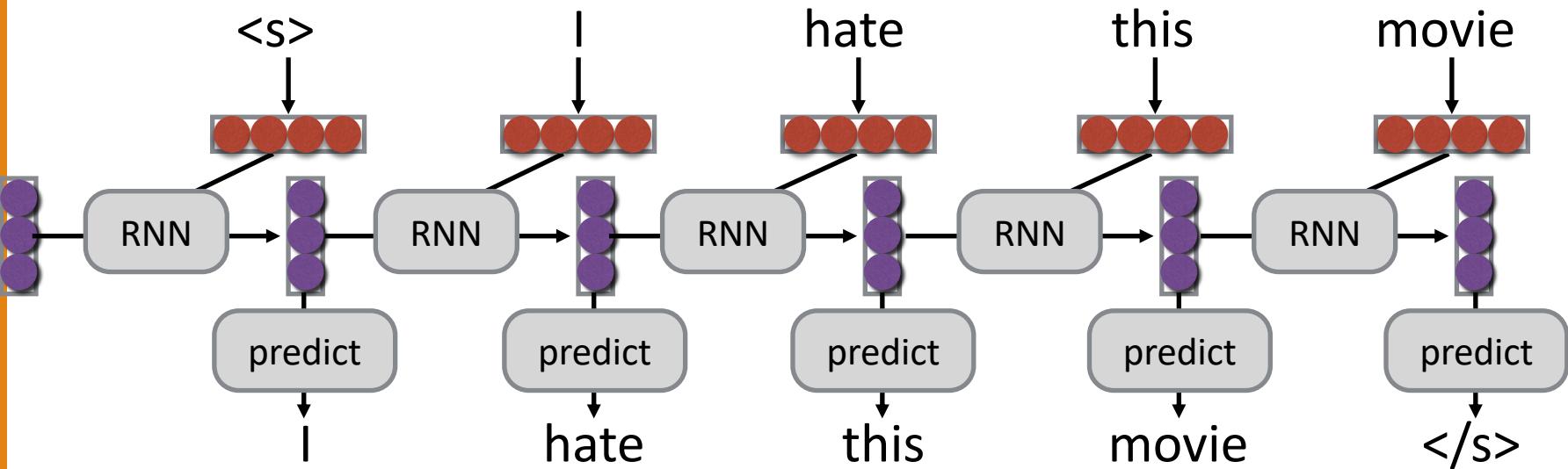
“I’m afraid I’ve definitely been suspended from power, no chance—indeed?” said Snape. He put his head back behind them and read groups as they crossed a corner and fluttered down onto their ink lamp, and picked up his spoon. The doorbell rang. It was a lot cleaner down in London.

Calculating the Probability of a Sentence

$$P(X) = \prod_{i=1}^I P(x_i \mid x_1, \dots, x_{i-1})$$


Next Word Context

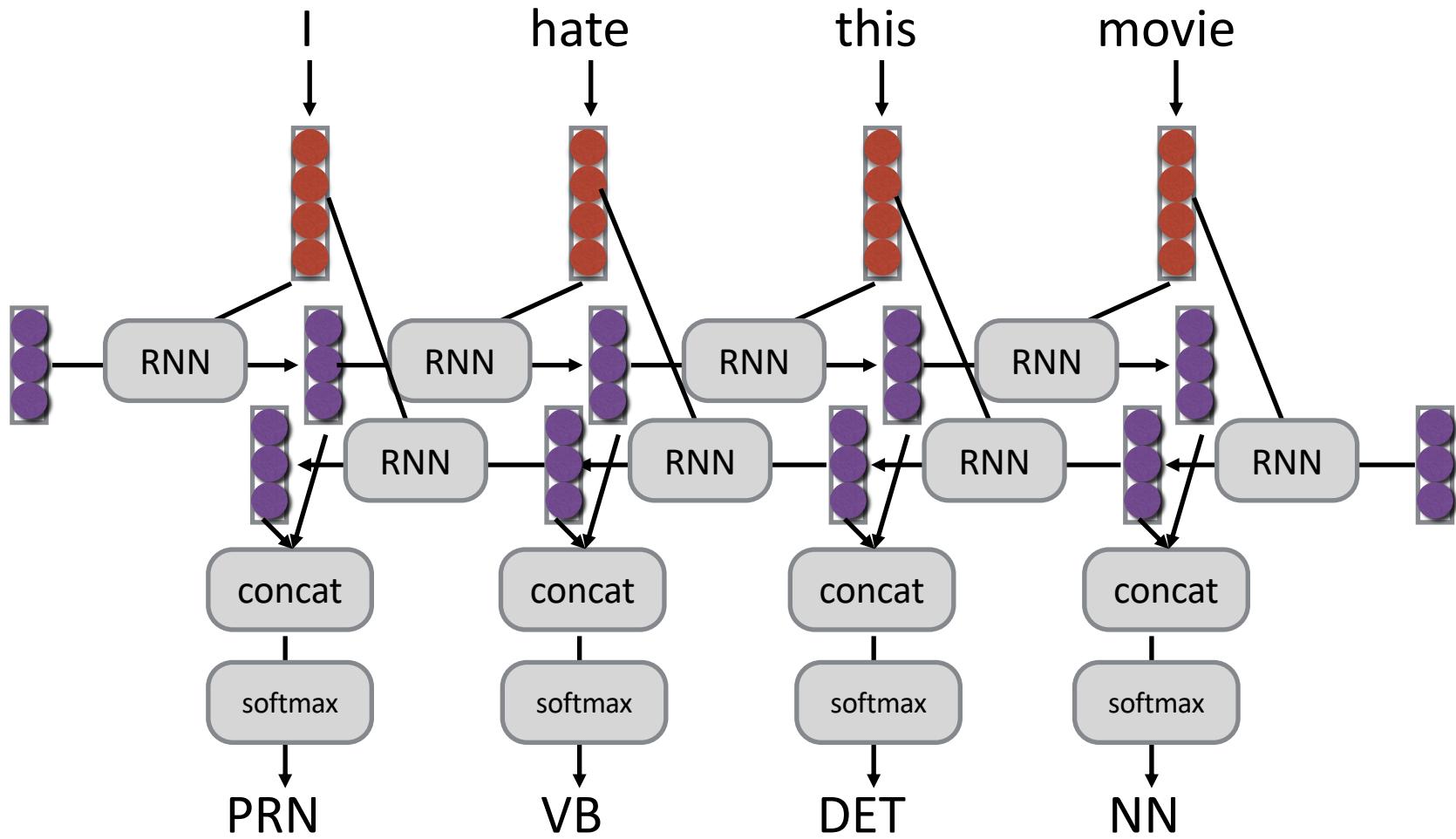
Language Models with RNNs



At each step, calculate probability of next word

Bi-directional RNNs

A simple extension, run the RNN in both directions



Conditional Language Modeling for Machine Translation

Conditional Language Models

Not just generate text, generate text according to some specification

<u>Input X</u>	<u>Output Y (Text)</u>	<u>Task</u>
Structured Data	NL Description	NL Generation
English	Japanese	Translation
Document	Short Description	Summarization
Utterance	Response	Response Generation
Image	Text	Image Captioning
Speech	Transcript	Speech Recognition

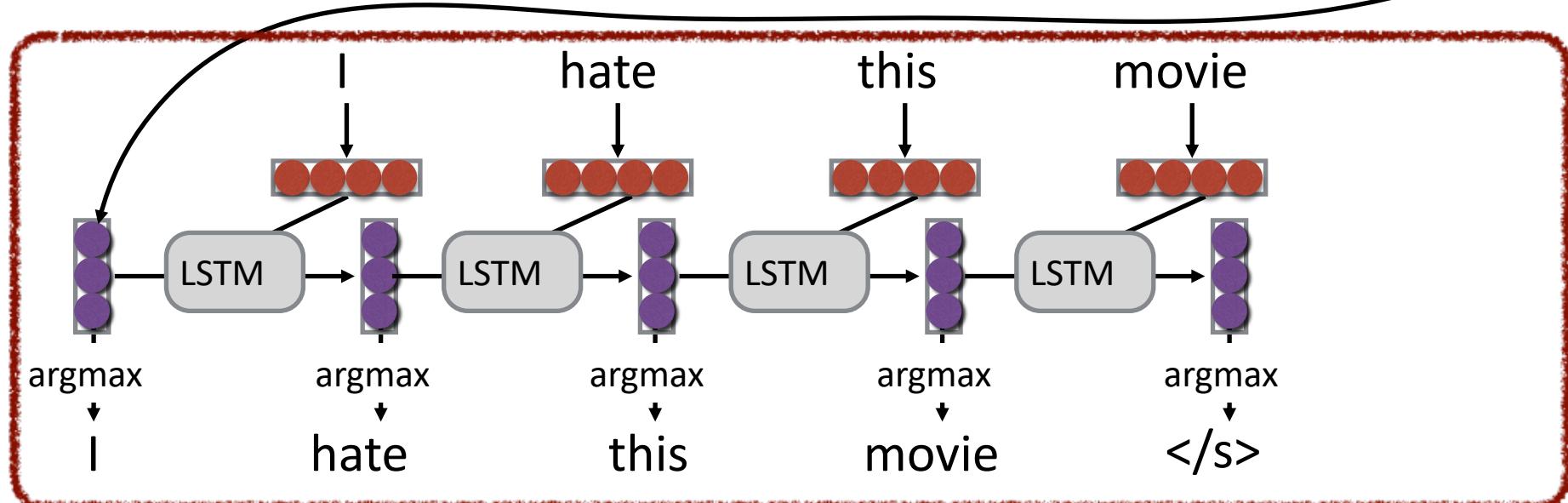
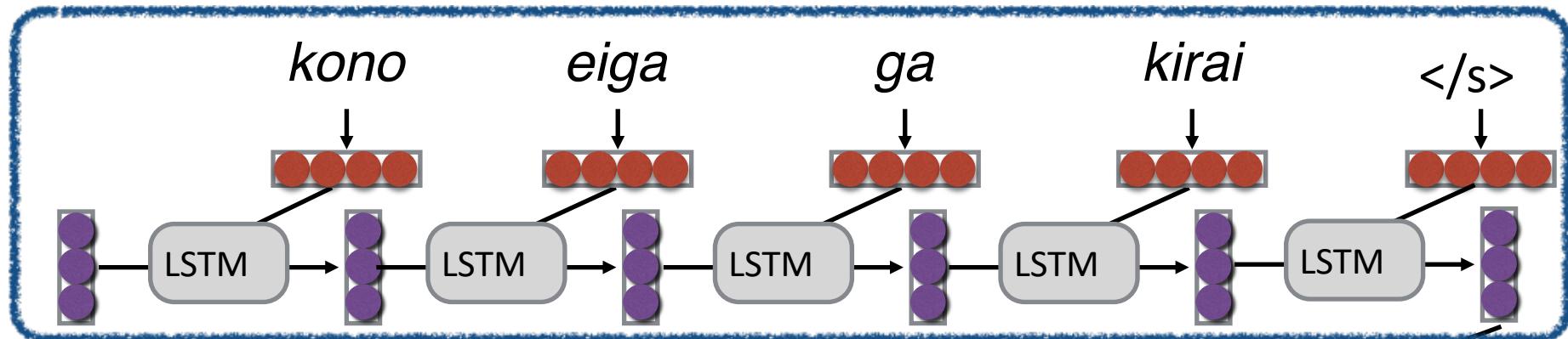
Conditional Language Models

$$P(Y|X) = \prod_{j=1}^J P(y_j | X, y_1, \dots, y_{j-1})$$

Added Context!

One type of Conditional LM

Encoder

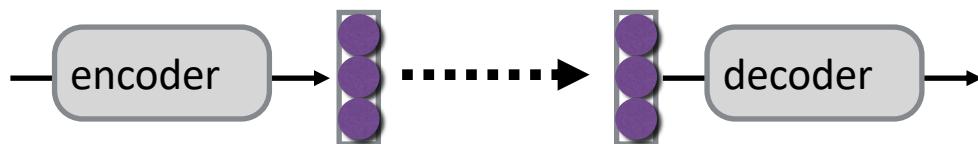


Decoder

Sutskever et al. 2014

How to pass the hidden state?

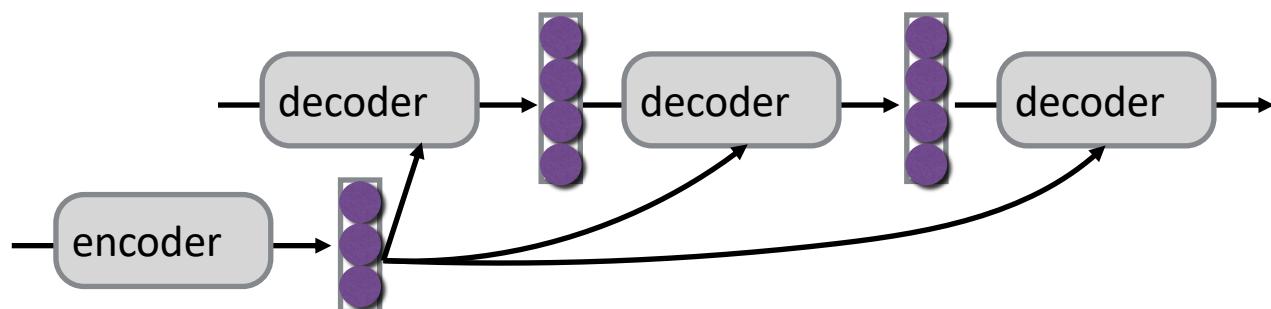
Initialize decoder w/ encoder (Sutskever et al. 2014)



Transform (can be different dimensions)



Input at every time step (Kalchbrenner & Blunsom 2013)



Training Conditional LMs

Get parallel corpus of inputs and outputs

Maximize likelihood

Standard corpora for MT:

- WMT Conference on Machine Translation runs an evaluation every year with large-scale (e.g. 10M sentence) datasets
- Smaller datasets, e.g. 200k sentence TED talks from IWSLT, can be more conducive to experimentation

The Generation Problem

We have a model of $P(Y|X)$, how do we use it to generate a sentence?

Two methods:

- **Sampling:** Try to generate a *random* sentence according to the probability distribution.
- **Argmax:** Try to generate the sentence with the *highest* probability.

Ancestral Sampling

Randomly generate words one-by-one.

```
while yj-1 != "</s>":  
    yj ~ P(yj | X, y1, ..., yj-1)
```

An **exact method** for sampling from P(X), no further work needed.

Greedy Search

One by one, pick the single highest-probability word

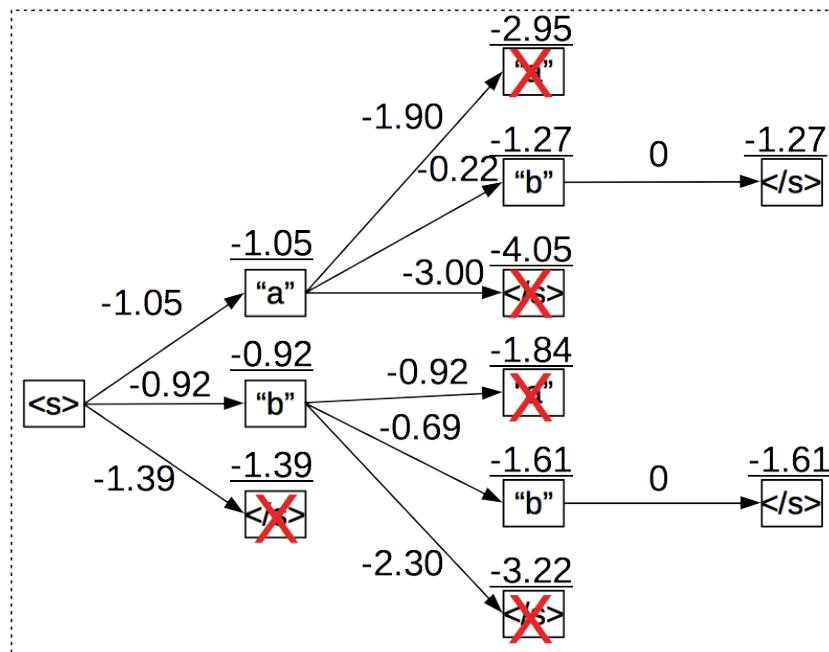
```
while yj-1 != "</s>":  
    yj = argmax P(yj | X, y1, ..., yj-1)
```

Not exact, which causes real problems:

1. Will often generate the “easy” words first
2. Will prefer multiple common words to one rare word

Beam Search

Instead of picking one high-probability word, maintain several paths



How do we Evaluate?

Basic Evaluation Paradigm

Use parallel test set

Use system to generate translations

Compare target translations w/ reference

Human Evaluation

Ask a human to do evaluation

太郎が花子を訪れた

Taro visited Hanako the Taro visited the Hanako Hanako visited Taro

Adequate?	Yes	Yes	No
Fluent?	Yes	No	Yes
Better?	1	2	3

- Final goal, but slow, expensive, and sometimes inconsistent

BLEU

Works by comparing n-gram overlap w/ reference

Reference: Taro visited Hanako

System: the Taro visited the Hanako

1-gram: 3/5

2-gram: 1/4

Brevity: $\min(1, |\text{System}|/|\text{Reference}|) = \min(1, 5/3)$ brevity penalty = 1.0

$$\begin{aligned}\text{BLEU-2} &= (3/5 * 1/4)^{1/2} * 1.0 \\ &= 0.387\end{aligned}$$

- **Pros:** Easy to use, good for measuring system improvement
- **Cons:** Often doesn't match human eval, bad for comparing very different systems

Other Options

METEOR: Considers synonyms

Translation Edit Rate (TER): Considers number of edits to be turned into a good translation

chrF: Considers score on the character level

RIBES: Considers reordering

etc. etc.

Attention

Sentence Representations



You can't cram the
meaning of a whole
sentence into a
single vector!

Sentence Representations

But what if we could use multiple vectors, based on the length of the sentence?

this is an example



this is an example



Basic Idea

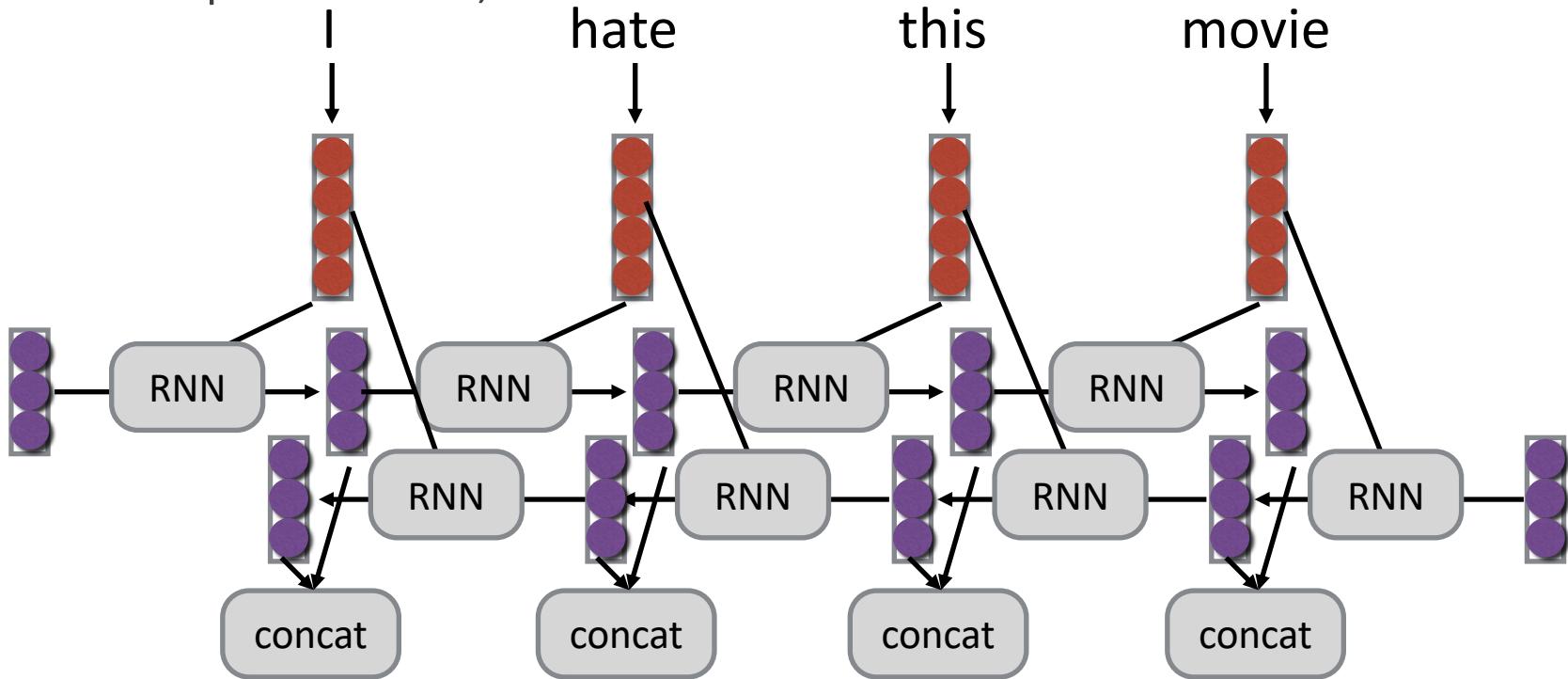
Encode each word in the sentence into a vector

When decoding, perform a linear combination of these vectors,
weighted by “attention weights”

Use this combination in picking the next word

Encoder Bi-RNNs

A simple extension, run the RNN in both directions

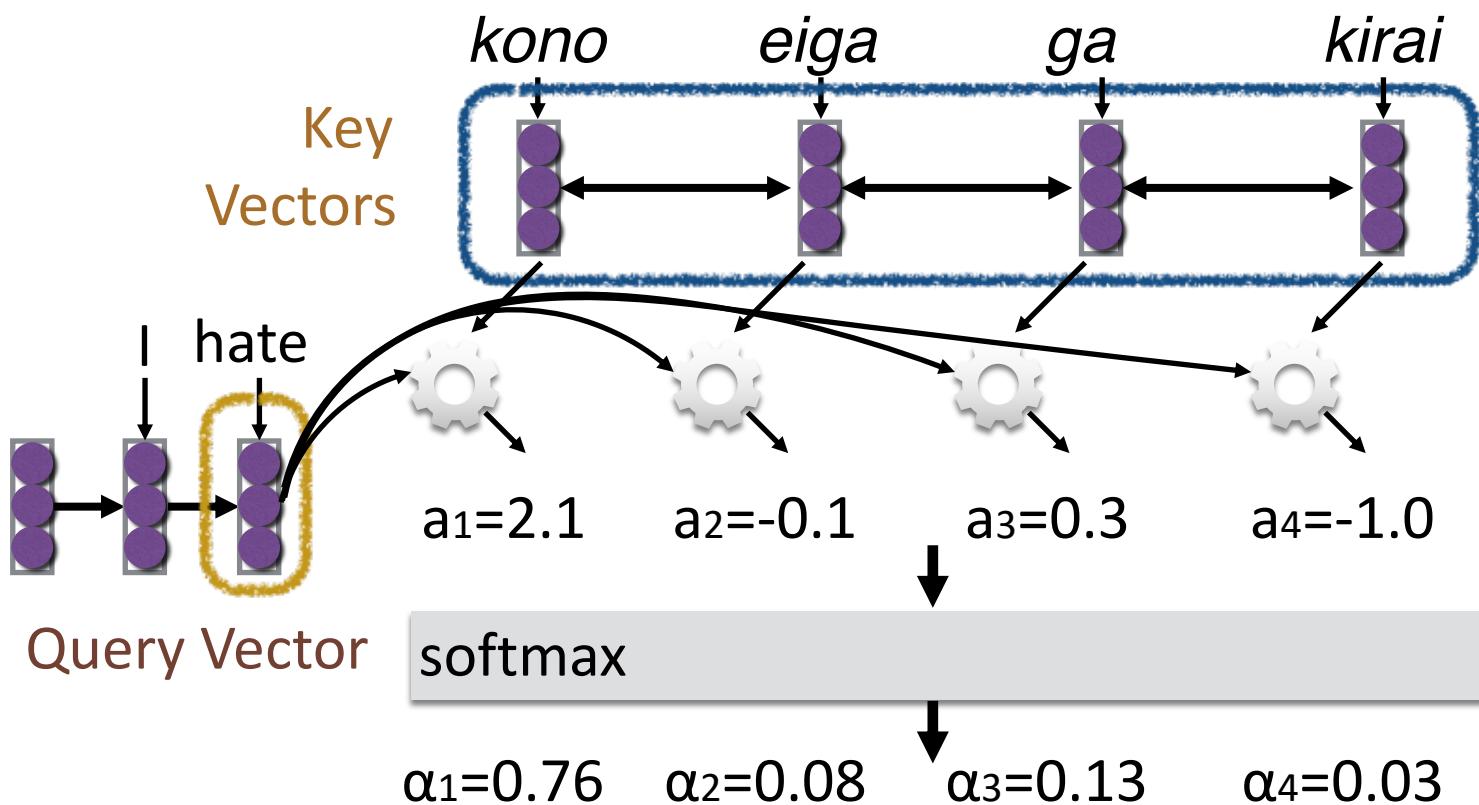


Calculating Attention (1)

Use “query” vector (decoder state) and “key” vectors (all encoder states)

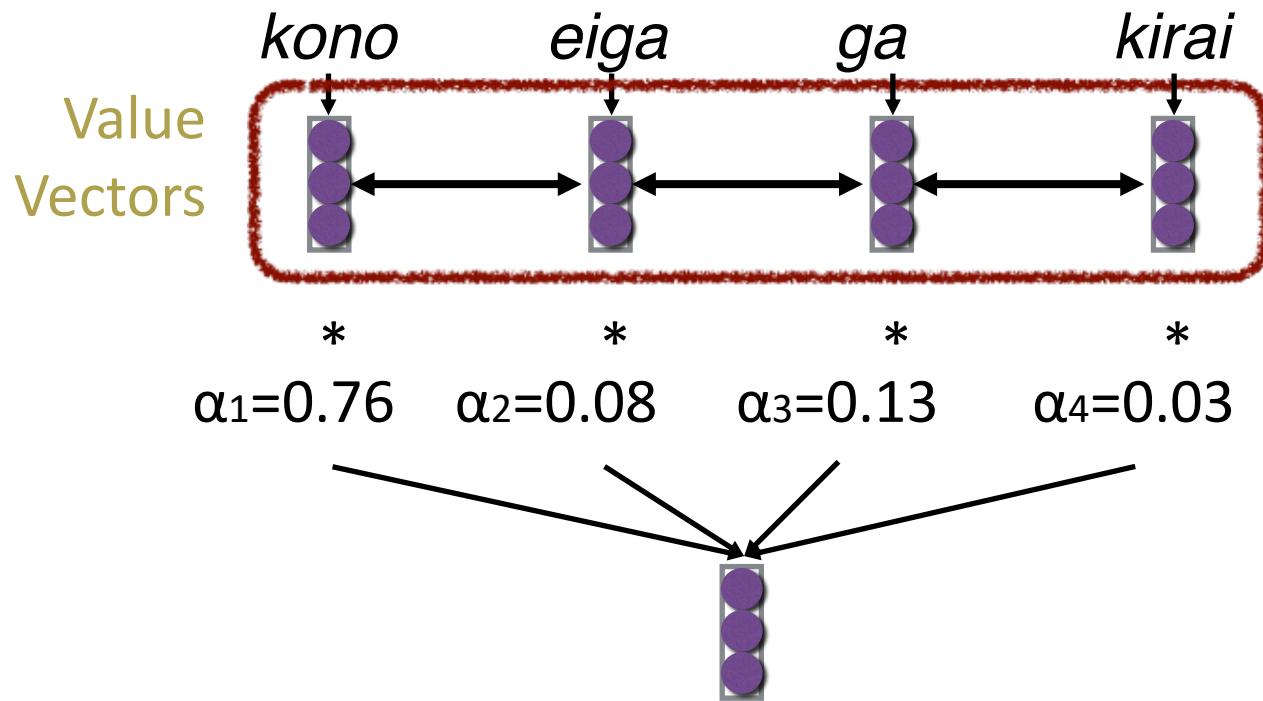
For each query-key pair, calculate weight

Normalize to add to one using softmax



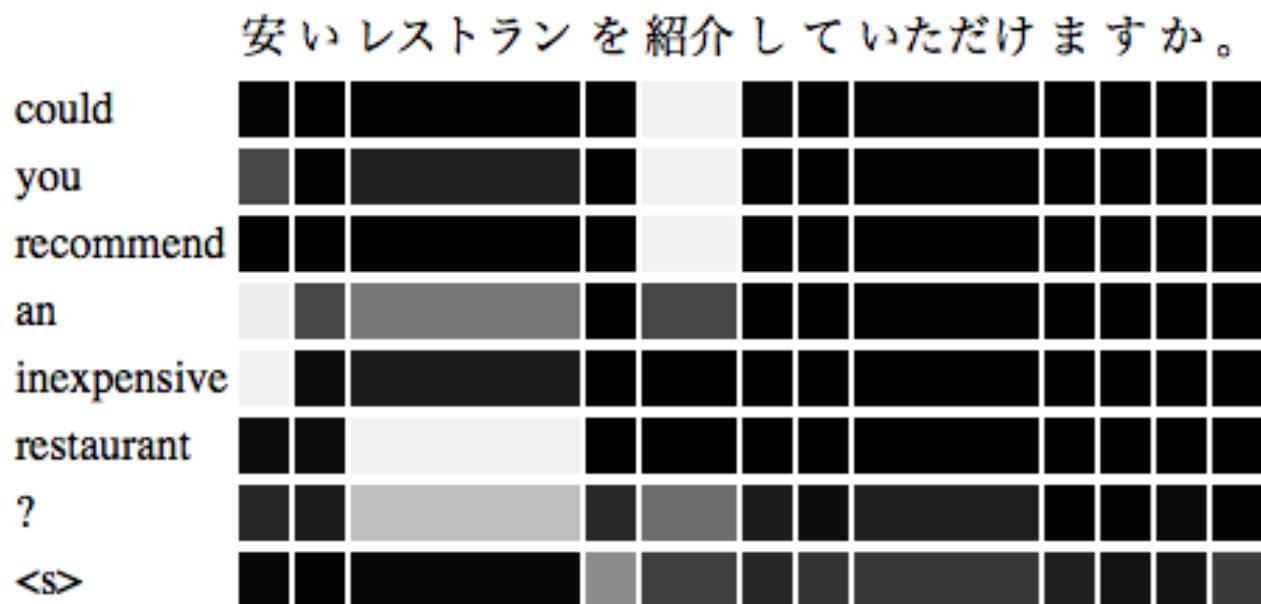
Calculating Attention (2)

Combine together value vectors (usually encoder states, like key vectors) by taking the weighted sum



Use this in any part of the model you like

A Graphical Example



Attention Score Functions (1)

\mathbf{q} is the query and \mathbf{k} is the key

Multi-layer Perceptron (Bahdanau et al. 2015)

$$a(\mathbf{q}, \mathbf{k}) = \mathbf{w}_2^\top \tanh(W_1[\mathbf{q}; \mathbf{k}])$$

- Flexible, often very good with large data

Bilinear (Luong et al. 2015)

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

Attention Score Functions (2)

Dot Product (Luong et al. 2015)

$$a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^\top \mathbf{k}$$

- No parameters! But requires sizes to be the same.

Scaled Dot Product (Vaswani et al. 2017)

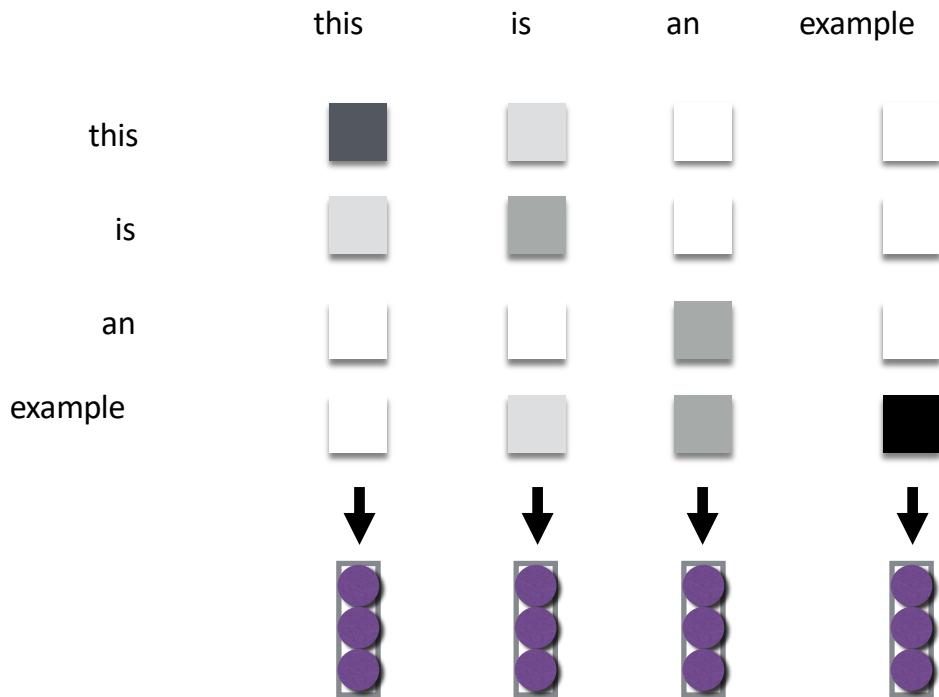
- Problem: scale of dot product increases as dimensions get larger
- Fix: scale by size of the vector

$$a(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q}^\top \mathbf{k}}{\sqrt{|\mathbf{k}|}}$$

Extensions to Attention

Intra-Attention / Self-Attention (Cheng et al. 2016)

Each element in the sentence attends to other elements → context sensitive encodings!



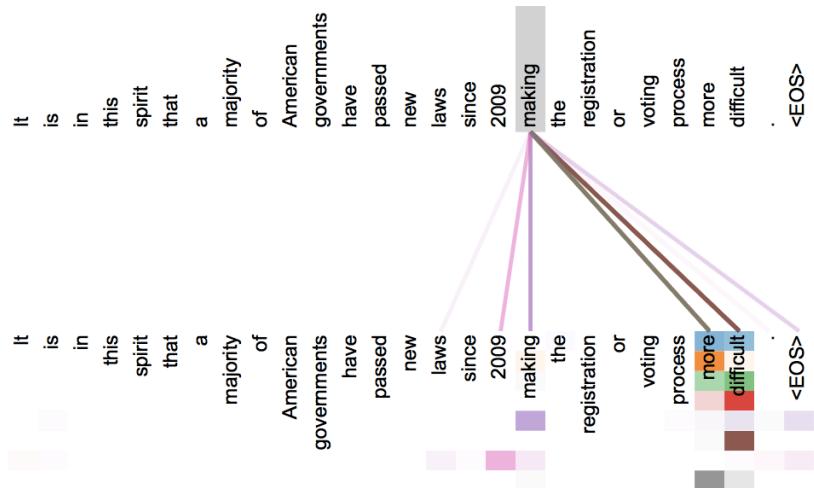
Multi-headed Attention

Idea: multiple attention “heads” focus on different parts of the sentence

- e.g. Different heads for “copy” vs regular (Allamanis et al. 2016)

	Target	Attention Vectors	λ
m_1	set	$\alpha = \langle s \rangle \{ \text{this} . \text{use Browser Cache} = \text{use Browser Cache} ; \} \langle /s \rangle$ $\kappa = \langle s \rangle \{ \text{this} . \text{use Browser Cache} = \text{use Browser Cache} ; \} \langle /s \rangle$	0.012
m_2	use	$\alpha = \langle s \rangle \{ \text{this} . \text{use Browser Cache} = \text{use Browser Cache} ; \} \langle /s \rangle$ $\kappa = \langle s \rangle \{ \text{this} . \text{use Browser Cache} = \text{use Browser Cache} ; \} \langle /s \rangle$	0.974
m_3	browser	$\alpha = \langle s \rangle \{ \text{this} . \text{use Browser Cache} = \text{use Browser Cache} ; \} \langle /s \rangle$ $\kappa = \langle s \rangle \{ \text{this} . \text{use Browser Cache} = \text{use Browser Cache} ; \} \langle /s \rangle$	0.969
m_4	cache	$\alpha = \langle s \rangle \{ \text{this} . \text{use Browser Cache} = \text{use Browser Cache} ; \} \langle /s \rangle$ $\kappa = \langle s \rangle \{ \text{this} . \text{use Browser Cache} = \text{use Browser Cache} ; \} \langle /s \rangle$	0.583
m_5	END	$\alpha = \langle s \rangle \{ \text{this} . \text{use Browser Cache} = \text{use Browser Cache} ; \} \langle /s \rangle$ $\kappa = \langle s \rangle \{ \text{this} . \text{use Browser Cache} = \text{use Browser Cache} ; \} \langle /s \rangle$	0.066

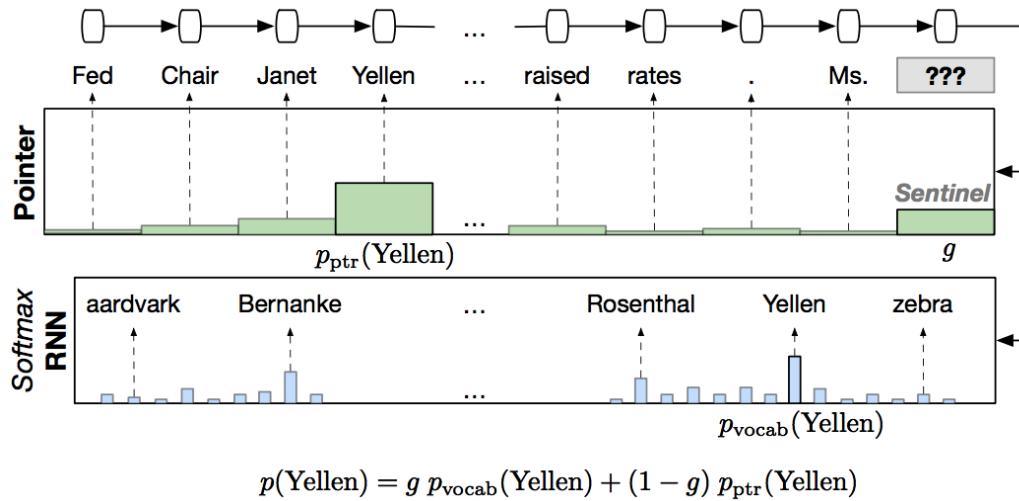
- Or multiple independently learned heads (Vaswani et al. 2017)



- Or one head for every hidden node! (Choi et al. 2018)

Attending to Previously Generated Things

In language modeling, attend to the previous words
(Merity et al. 2016)



In translation, attend to either input or previous output (Vaswani et al. 2017)

An Interesting Case
Study:

“Attention is All You
Need”

(Vaswani et al. 2017)

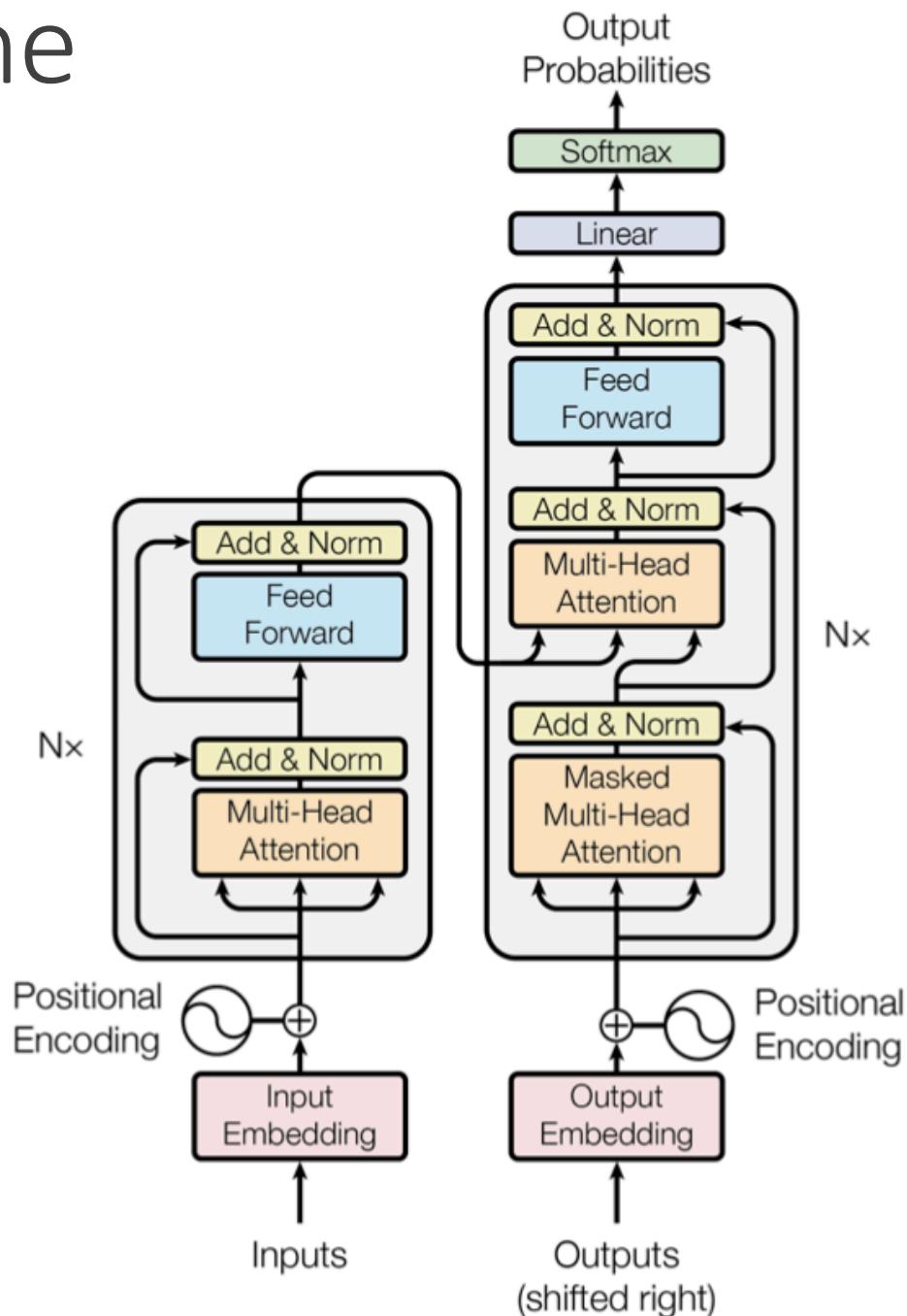
Summary of the “Transformer”

(Vaswani et al. 2017)

A sequence-to-sequence model based entirely on attention

Also have attention on the output side!
Calculate probability of next word by attention over previous words.

Fast: only matrix multiplications



Attention Tricks

Self Attention: Each layer combines words with others

Multi-headed Attention: 8 attention heads function independently

Normalized Dot-product Attention: Remove bias in dot product when using large networks

Positional Encodings: Make sure that even if we don't have RNN, can still distinguish positions

Training Tricks

Layer Normalization: Help ensure that layers remain in reasonable range

Specialized Training Schedule: Adjust default learning rate of the Adam optimizer

Label Smoothing: Insert some uncertainty in the training process

Masking for Efficient Training

Masking for Training

We want to perform training in as few operations as possible using big matrix multiplies

We can do so by “masking” the results for the output

How to Get Started?

Getting Started

Find training data, (e.g. TED talks from IWSLT), in your favorite language

Download a toolkit (e.g. OpenNMT, fairseq, Sockeye, xnmt) and run it on the data

Calculate the BLEU score and look at the results

Think of what's going right, what's going wrong!

Questions?

To Learn More:
"Neural Machine Translation and
Sequence-to-sequence Models: A Tutorial"