

# Sequence to Sequence Models with Attention

Mausam

(Slides by Yoav Goldberg, Graham Neubig, Prabhakar Raghavan)

**BACKGROUND**

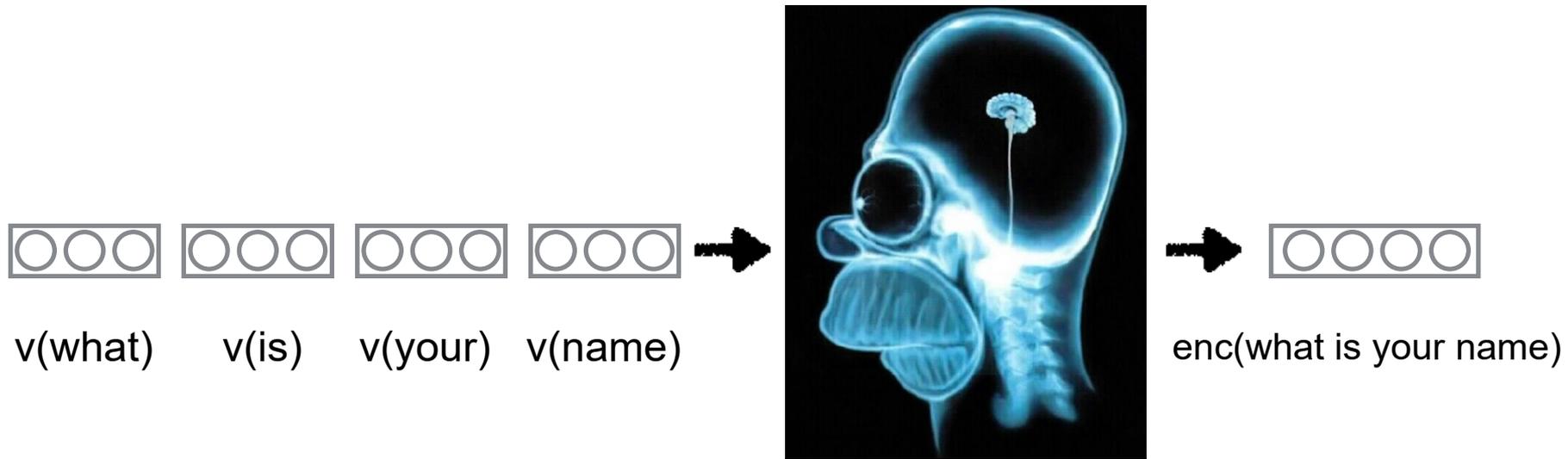
# A Primer on D.L. Building Blocks

- A single vector for an ordered pair of vectors? •  $x;y$
- A single vector for a variable-sized bag of vectors? •  $\sum_i x_i$
- Project a vector to a new space? •  $Wx$
- Are two vectors (from same space) similar? •  $x.y$
- Are two vectors (from different space) similar? •  $xWy$
- A new vector that depends on some vector input? •  $g(Wx+b)$

# A Primer on D.L. Building Blocks

- Output a probability •  $\sigma$
- Output one of two classes •  $\sigma$
- Output one of many classes • softmax
- A feature w/ positive & negative influence • tanh
- A feature w/ positive influence for “deep” nets • ReLu

# Recurrent Neural Networks



- Very strong models of sequential data.
- **Trainable** function from  $n$  vectors to a single vector.

# Recurrent Neural Networks

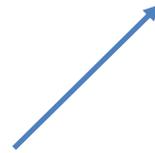
$$RNN(\mathbf{s}_0, \mathbf{x}_{1:n}) = \mathbf{s}_n, \mathbf{y}_n$$

$$\mathbf{x}_i \in \mathbb{R}^{d_{in}}, \mathbf{y}_i \in \mathbb{R}^{d_{out}}, \mathbf{s}_i \in \mathbb{R}^{f(d_{out})}$$

- Very strong models of sequential data.
- **Trainable** function from  $n$  vectors to a single\* vector.

# Recurrent Neural Networks

$$RNN(\mathbf{s}_0, \mathbf{x}_{1:n}) = \mathbf{s}_n, \mathbf{y}_n$$



\*this one is internal. we only care about the  $\mathbf{y}$

$$\mathbf{x}_i \in \mathbb{R}^{d_{in}}, \mathbf{y}_i \in \mathbb{R}^{d_{out}}, \mathbf{s}_i \in \mathbb{R}^{f(d_{out})}$$

- Very strong models of sequential data.
- **Trainable** function from  $n$  vectors to a single\* vector.

# Recurrent Neural Networks

$$RNN(\mathbf{s}_0, \mathbf{x}_{1:n}) = \mathbf{s}_n, \mathbf{y}_n$$

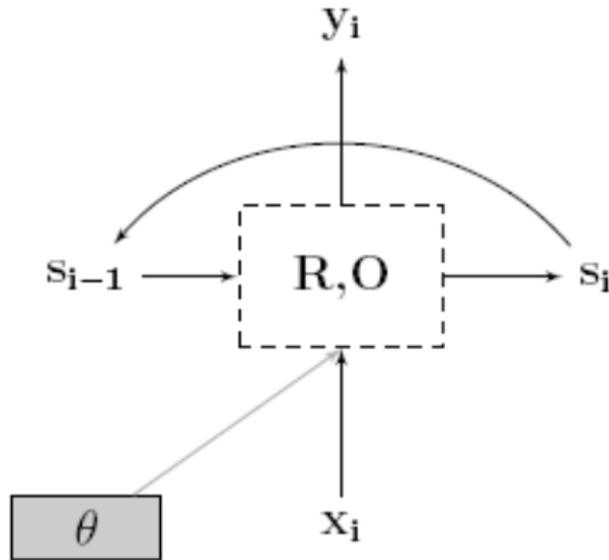
$$\mathbf{s}_i = R(\mathbf{s}_{i-1}, \mathbf{x}_i)$$

$$\mathbf{y}_i = O(\mathbf{s}_i)$$

$$\mathbf{x}_i \in \mathbb{R}^{d_{in}}, \mathbf{y}_i \in \mathbb{R}^{d_{out}}, \mathbf{s}_i \in \mathbb{R}^{f(d_{out})}$$

- **Recursively defined.**
- There's a vector  $\mathbf{y}_i$  for every prefix  $\mathbf{x}_{1:i}$

# Recurrent Neural Networks



$$RNN(\mathbf{s}_0, \mathbf{x}_{1:n}) = \mathbf{s}_n, \mathbf{y}_n$$

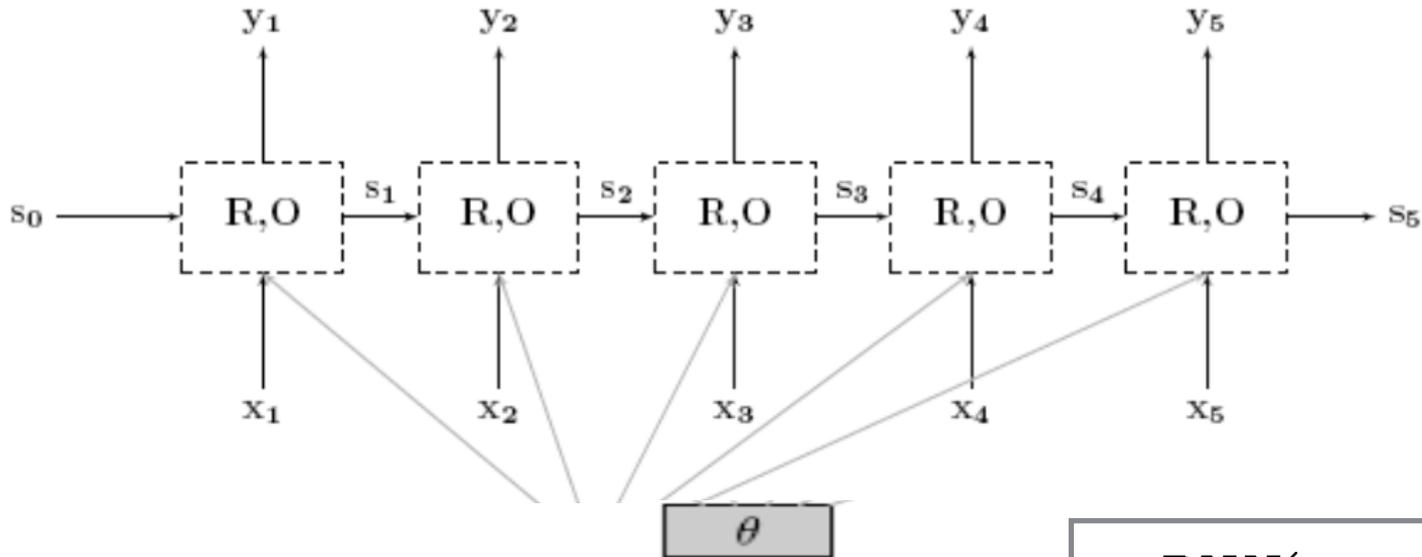
$$\mathbf{s}_i = R(\mathbf{s}_{i-1}, \mathbf{x}_i)$$

$$\mathbf{y}_i = O(\mathbf{s}_i)$$

$$\mathbf{x}_i \in \mathbb{R}^{d_{in}}, \mathbf{y}_i \in \mathbb{R}^{d_{out}}, \mathbf{s}_i \in \mathbb{R}^{f(d_{out})}$$

- **Recursively defined.**
- There's a vector  $\mathbf{y}_i$  for every prefix  $\mathbf{x}_{1:i}$

# Recurrent Neural Networks

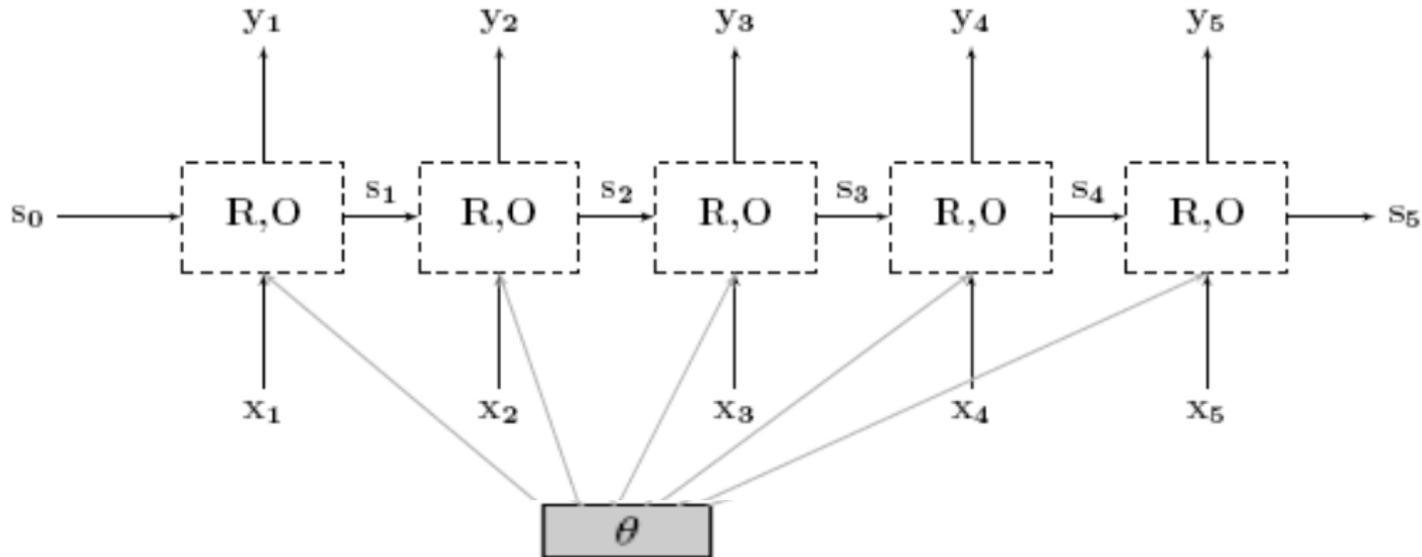


for every finite input sequence,  
can unroll the recursion.

$$RNN(s_0, \mathbf{x}_{1:n}) = \mathbf{s}_n, \mathbf{y}_n$$
$$\mathbf{s}_i = R(\mathbf{s}_{i-1}, \mathbf{x}_i)$$
$$\mathbf{y}_i = O(\mathbf{s}_i)$$
$$\mathbf{x}_i \in \mathbb{R}^{d_{in}}, \mathbf{y}_i \in \mathbb{R}^{d_{out}}, \mathbf{s}_i \in \mathbb{R}^{f(d_{out})}$$

- Recursively defined.
- There's a vector  $\mathbf{y}_i$  for every prefix  $\mathbf{x}_{1:i}$

# Recurrent Neural Networks



for every finite input sequence,  
can unroll the recursion.

State  $s_i$  encodes **history** till this point



# Simple RNN (Elman RNN)

$$R_{SRNN}(\mathbf{s}_{i-1}, \mathbf{x}_i) = \tanh(\mathbf{W}^s \cdot \mathbf{s}_{i-1} + \mathbf{W}^x \cdot \mathbf{x}_i)$$

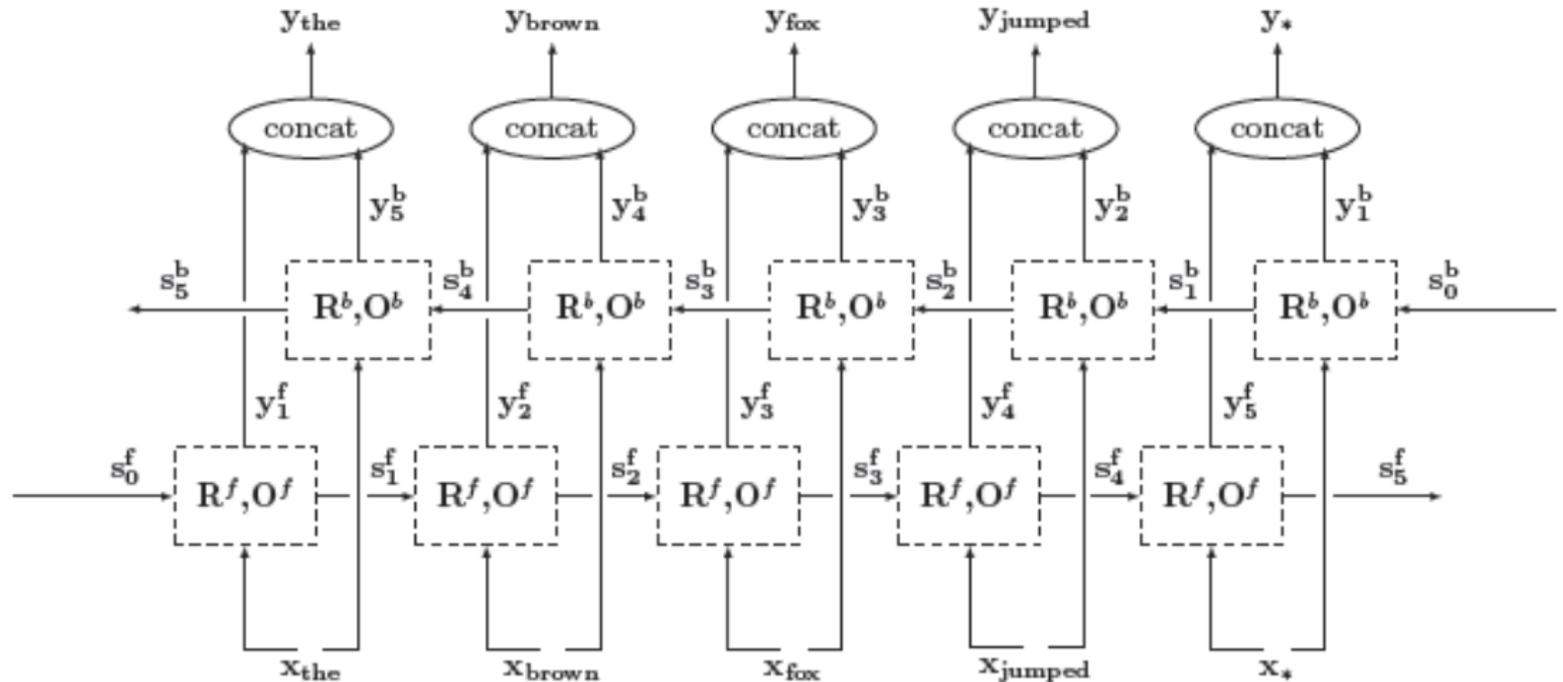
In principle: capture infinite history upto this point

In practice: have issues with long sequences

## RNN → LSTM

Good for backpropagating through long chain sequences

# Bidirectional RNNs

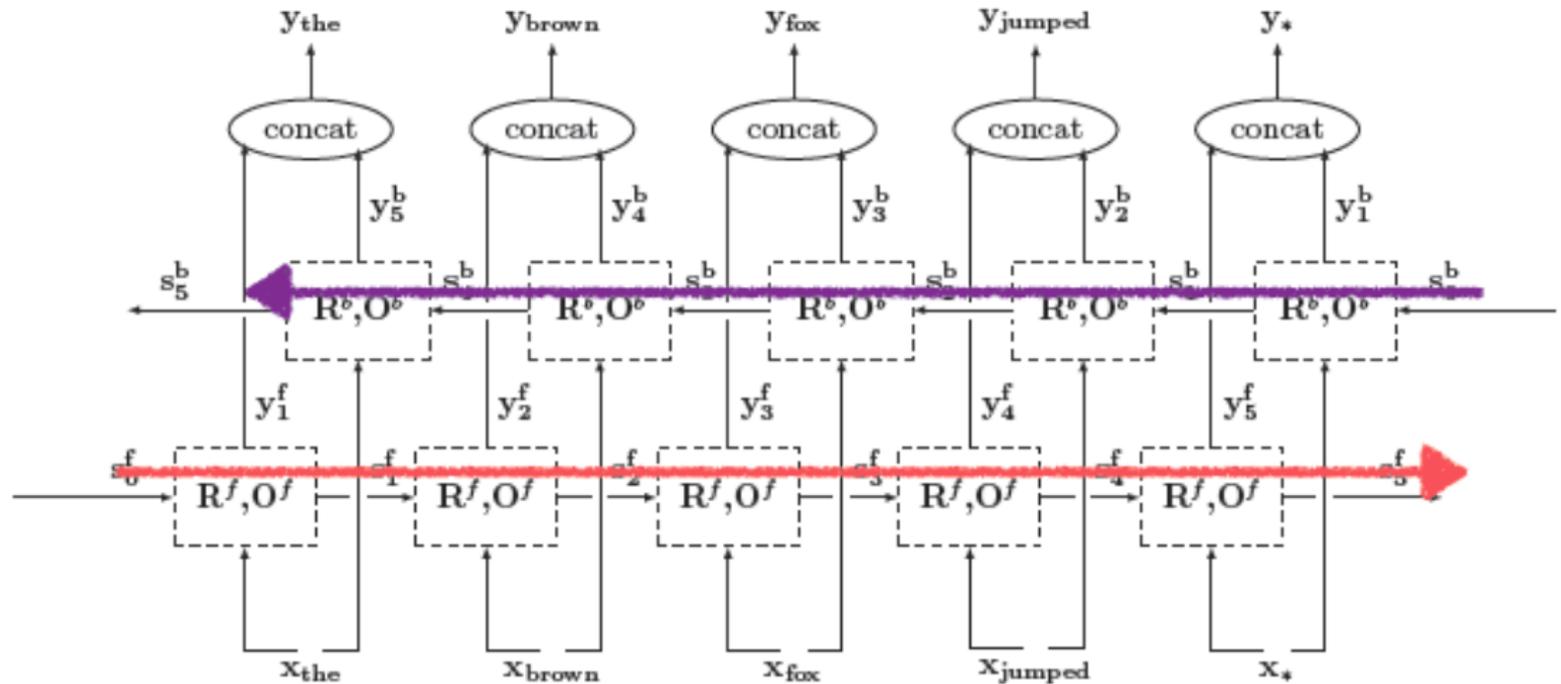


One RNN runs left to right.

Another runs right to left.

Encode **both future and history** of a word.

# Bidirectional RNNs



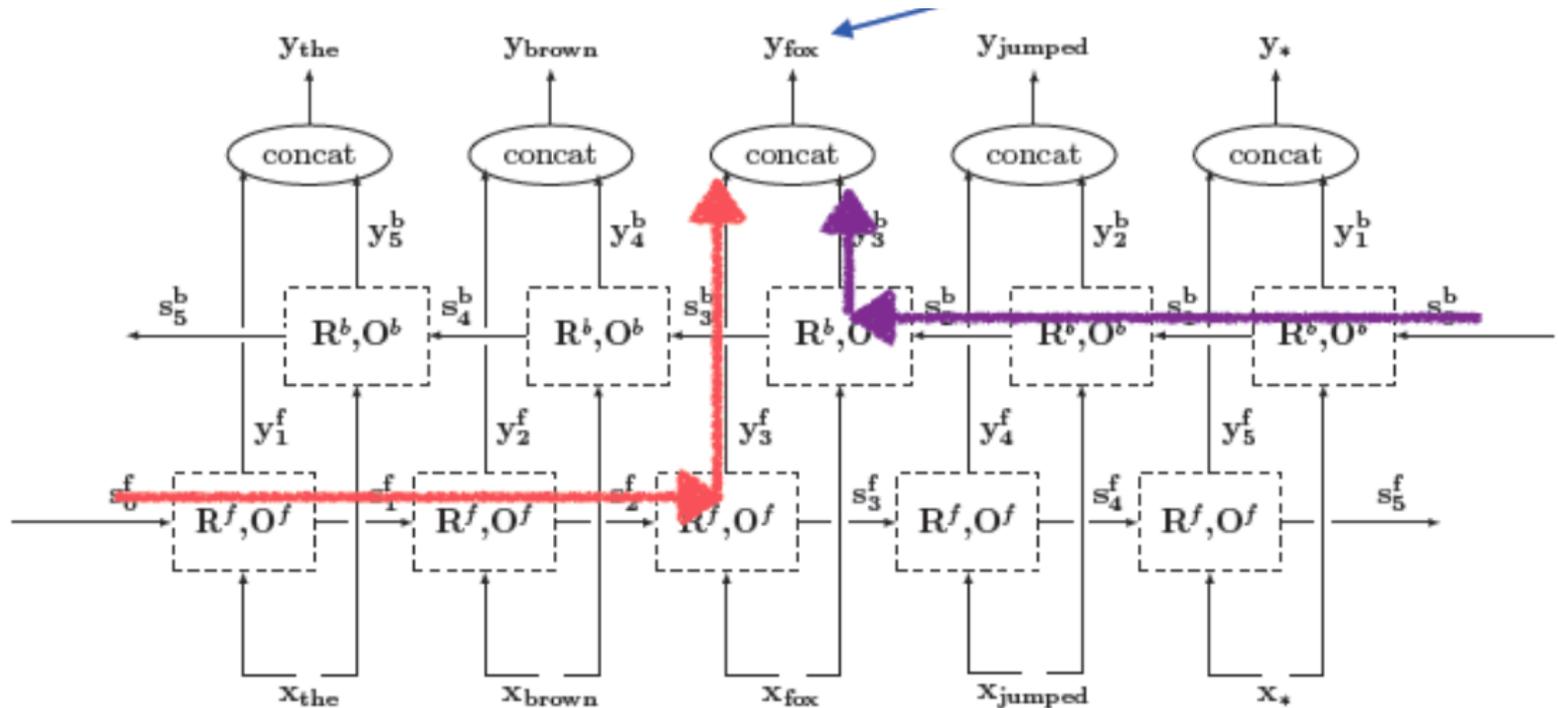
One RNN runs left to right.

Another runs right to left.

Encode **both future and history** of a word.

# Bidirectional RNNs

Infinite window around the word



One RNN runs left to right.

Another runs right to left.

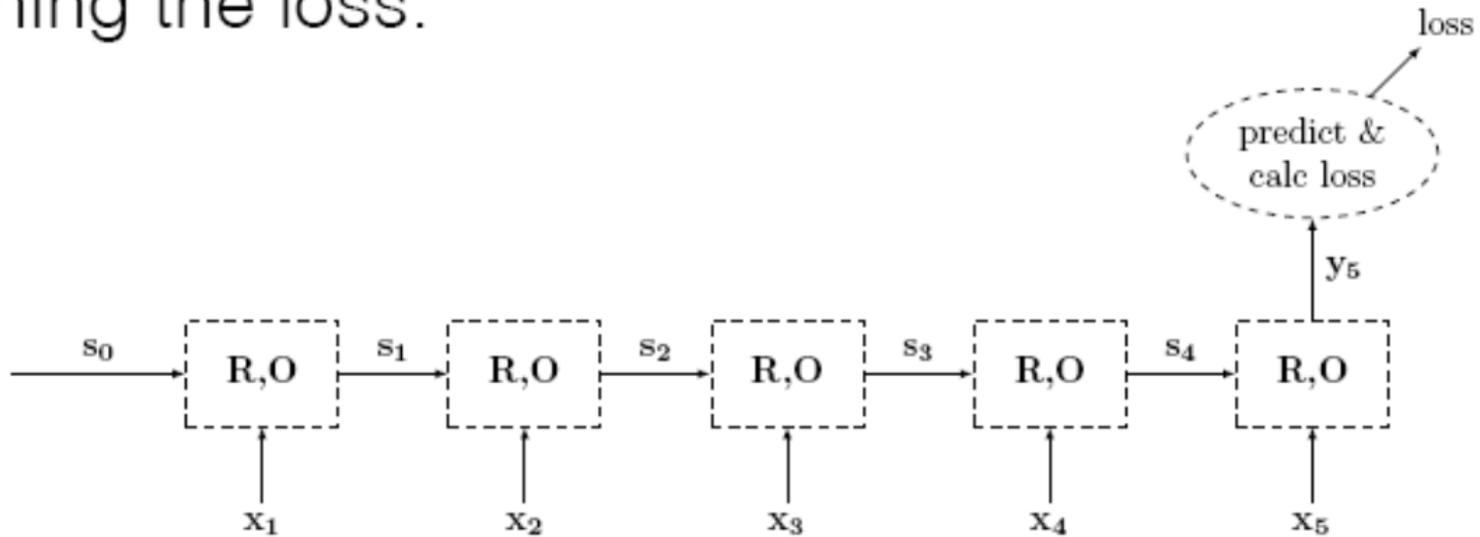
Encode **both future and history** of a word.

# Neural Architectures

- Mapping from a sequence to a single decision.
  - with RNN acceptor.
- Mapping from a sequence to a sequence of same length.
  - with RNN transducer

# RNN Acceptor

Defining the loss.

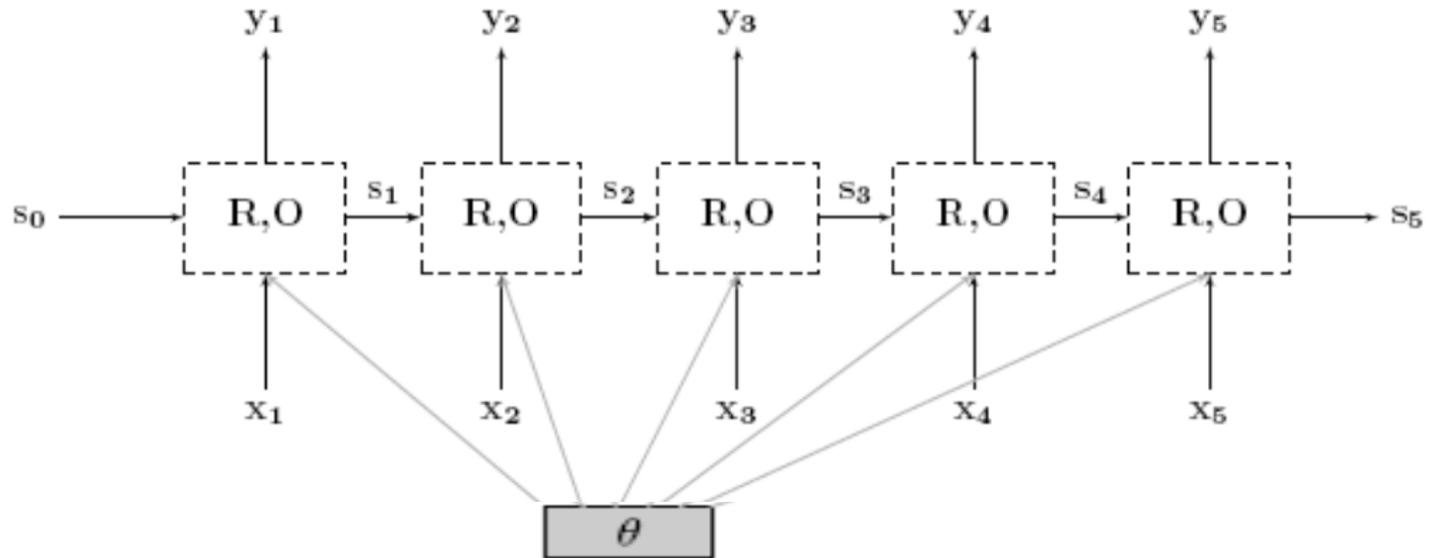


**Acceptor:** predict something from end state.

Backprop the error all the way back.

Train the network to capture meaningful information

# RNN Transducer



what do we do if the input and output sequences are of **different lengths**?

# Sequence Generation

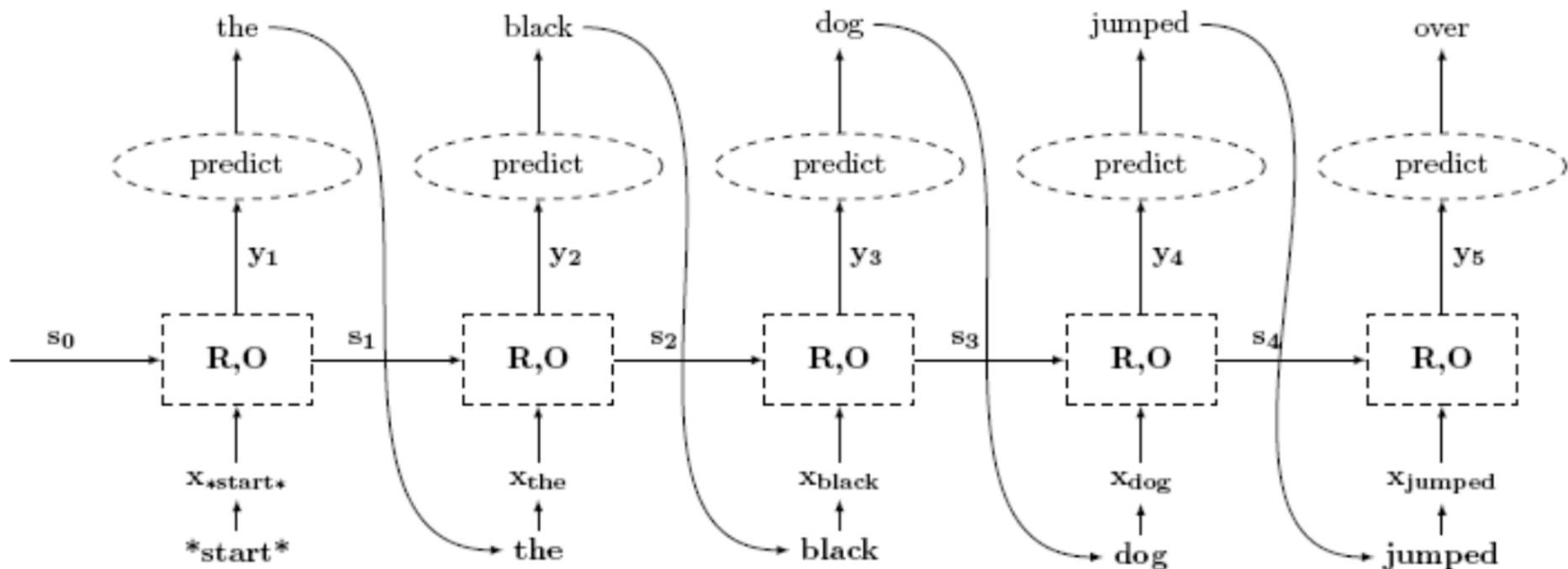
# Aside

How about an architecture for  
**0 to n** mapping.

(Neural Language Model)

# RNN Language Models

- *Training*: similar to an RNN Transducer.
- *Generation*: the output of step  $i$  is input to step  $i+1$ .



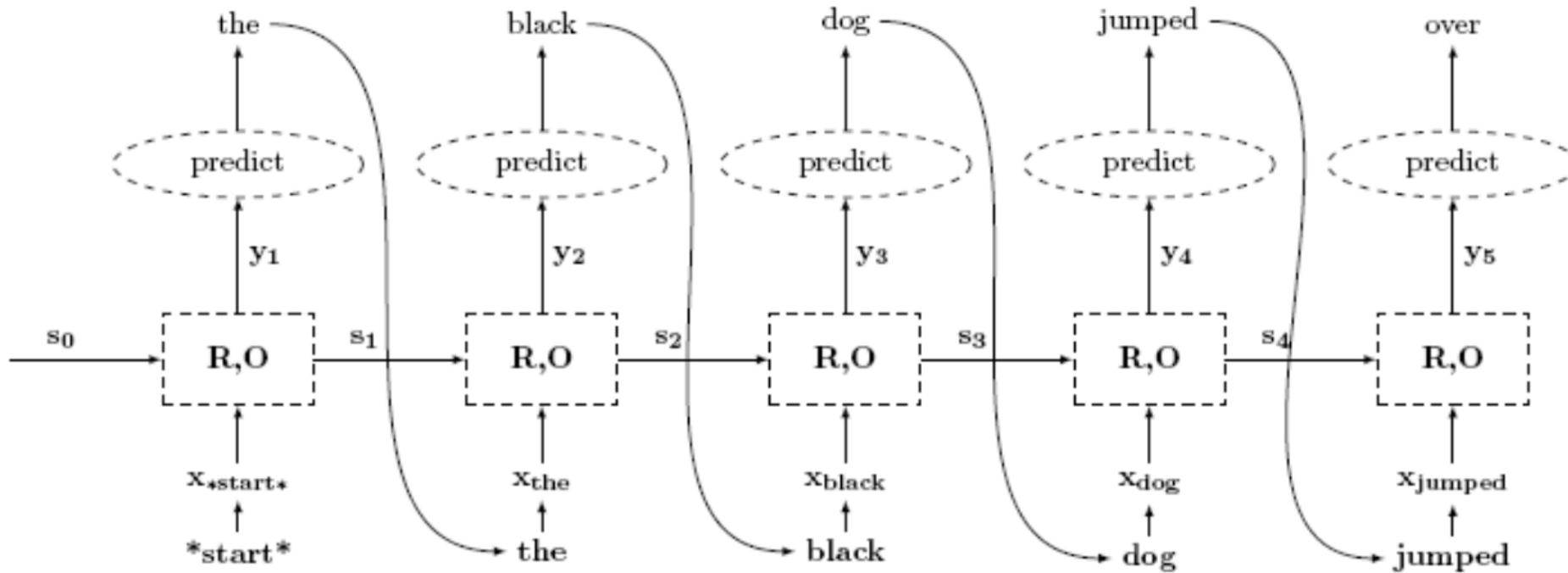
# RNN Language Model for generation

- Define the probability distribution over the next item in a sequence (and hence the probability of a sequence).

$$P(w_{1:n}) = P(w_1)P(w_2 | w_1)P(w_3 | w_{1:2})P(w_4 | w_{1:3}) \dots P(w_n | w_{1:n-1})$$

$$P(w_1, \dots, w_n) = \prod_{i=1}^n P(t_i = w_i | w_1, \dots, w_{i-1})$$

# RNN Language Models



$$p(t_{j+1} = k \mid \hat{t}_{1:j}) = f(\text{RNN}(\hat{\mathbf{t}}_{1:j}))$$

$$\hat{t}_j \sim p(t_j \mid \hat{t}_{1:j-1})$$

$$p(t_{j+1} = k \mid \hat{t}_{1:j}) = f(O(s_{j+1}))$$

$$s_{j+1} = R(\hat{t}_j, s_j)$$

$$\hat{t}_j \sim p(t_j \mid \hat{t}_{1:j-1})$$

# Sequence 2 Sequence

## Part I: No attention

# Back to Original question

How about an architecture for **m to n** mapping.

Generating sentences is nice, but what if we want to add some additional conditioning contexts?

# Conditioned Language Model

- Not just generate text, generate text according to some specification

<u>Input X</u>	<u>Output Y (<b>Text</b>)</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

# RNN Language Model for **Conditioned** generation

Let's add the condition variable to the equation.

$$P(\tau) = \prod_{i=1}^I P(t_i | \underbrace{t_1, \dots, t_{i-1}}_{\text{Context}})$$

Next Word

$$P(\tau | c) = \prod_{j=1}^J P(t_j | \underbrace{c}_{\text{Added Context! (a vector)}}, t_1, \dots, t_{j-1})$$

# How to Pass Context

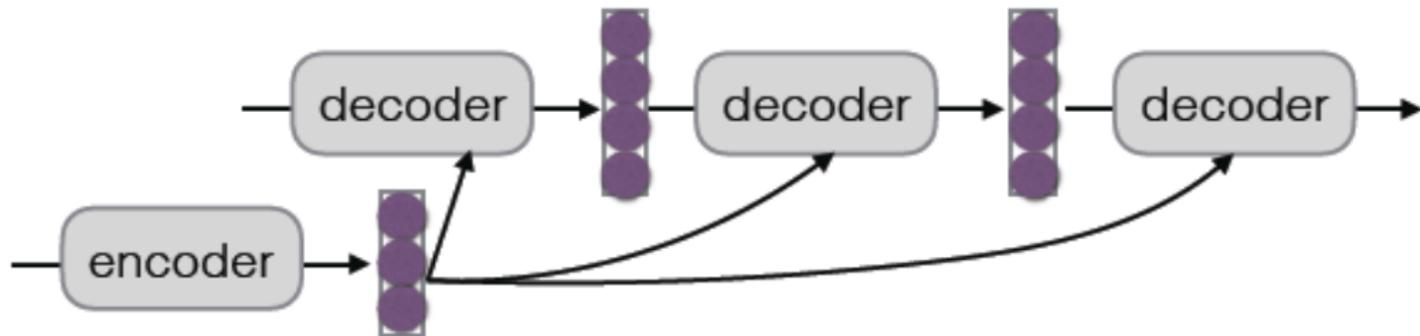
- Initialize decoder w/ encoder (Sutskever et al. 2014)



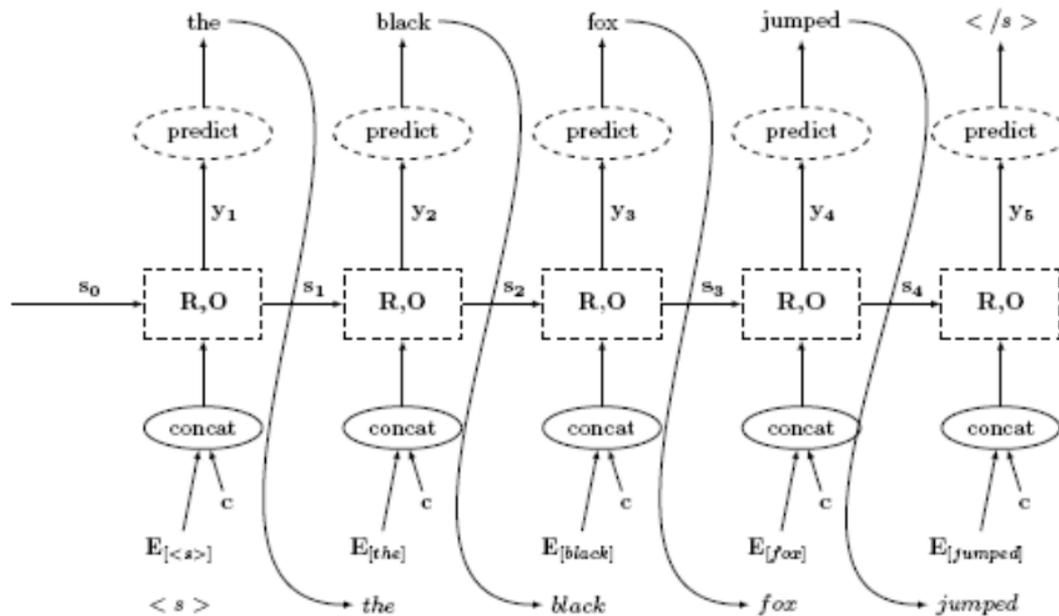
- Transform (can be different dimensions)



- Input at every time step (Kalchbrenner & Blunsom 2013)



# RNN Language Model for **Conditioned** generation



$$p(t_{j+1} = k \mid \hat{t}_{1:j}, c) = f(O(s_{j+1}))$$

$$s_{j+1} = R(s_j, [\hat{t}_j; c])$$

$$\hat{t}_j \sim p(t_i \mid \hat{t}_{1:j-1}, c)$$

# RNN Language Model for **Conditioned generation**

what if we want to condition on an entire sentence?

just encode it as a vector...

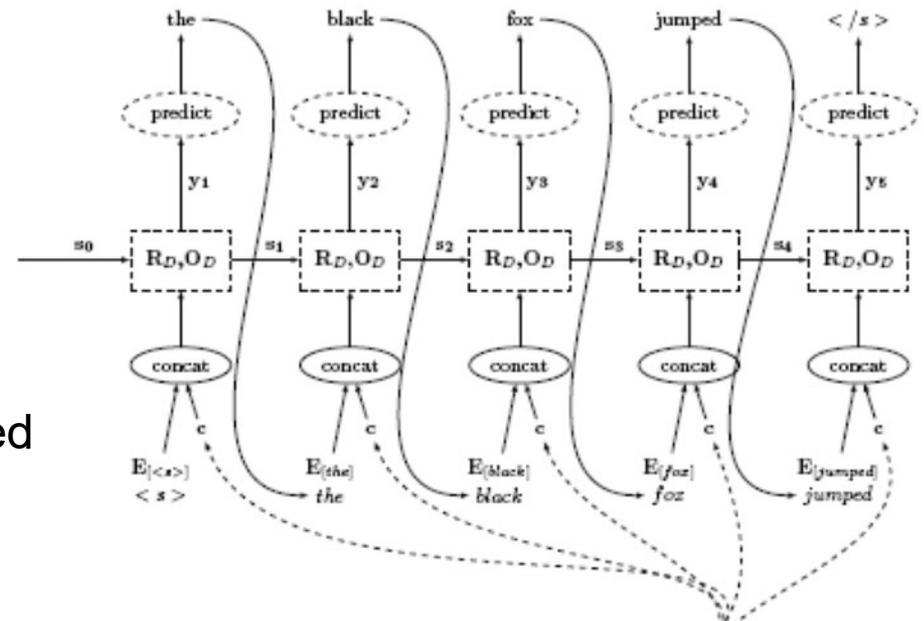
$$\mathbf{c} = \text{RNN}^{\text{enc}}(\mathbf{x}_{1:n})$$

# Sequence to Sequence conditioned generation

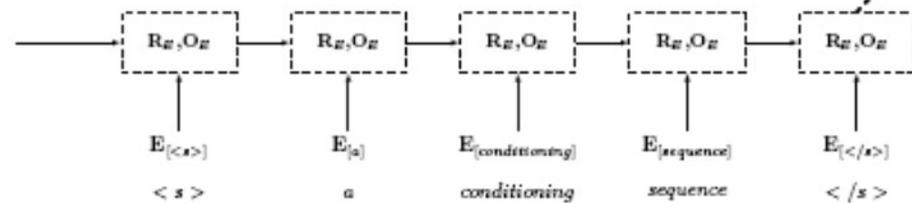
This is also called  
"Encoder Decoder"  
architecture.

Decoder

Decoder is  
just a conditioned  
language model



Encoder



# The Generation Problem

We have a probability model, 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  $y_{j-1} \neq \text{"</s>"}$ :  
   $y_j \sim P(y_j \mid X, y_1, \dots, y_{j-1})$ 
```

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

# Greedy Search

One by one, pick the single highest-probability word

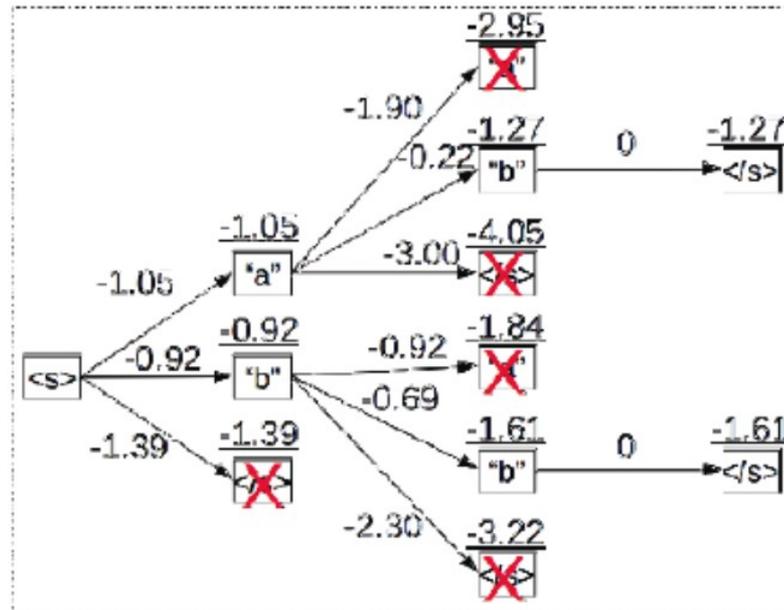
```
while  $y_{j-1} \neq \text{"</s>"}$ :  
   $y_j = \operatorname{argmax} P(y_j \mid X, y_1, \dots, y_{j-1})$ 
```

## Not exact, real problems:

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

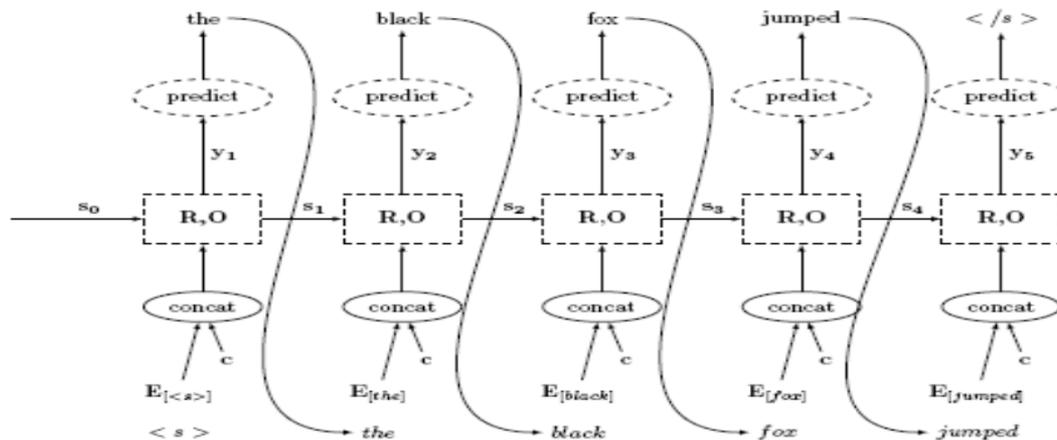
# Beam Search

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



# How to Train this Model?

- Issues with vanilla training
  - Slow convergence. Model instability. Poor skill.
- Simple idea: **Teacher Forcing**
  - Just feed in the *correct* previous word during training
- Drawback: **Exposure bias**
  - Not exposed to mistakes during training



# Solutions to Exposure Bias

- Start with no mistakes, and then
  - gradually introduce them using annealing
- Dropout inputs
  - Helps ensure that the model doesn't rely too heavily on predictions, while still using them
- Corrupt training data

# Evaluation Idea 1

- Steps
  - Use parallel test set
  - Use system to generate translations
  - Compare target translations w/ reference
- We train parameters of our model on a **training set**.
- We test the model's performance on data we haven't seen.
  - A **test set** is an unseen dataset that is different from our training set, totally unused.
  - An **evaluation metric** tells us how well our model does on the test set.

# Metric: Human 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

# Metric: 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

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

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

# Metric: METEOR

- Like BLEU in overall principle, with many other tricks: consider paraphrases, reordering, and function word/content word difference
- **Pros:** Generally significantly better than BLEU, esp. for high-resource languages
- **Cons:** Requires extra resources for new languages (although these can be made automatically), and more complicated

# Evaluation Idea 2

- Don't generate a translation. Instead ask:
- Does our language model prefer good sentences to bad ones?
  - Assign a high probability to “real” sentence

# Intuition of Perplexity

- The Shannon Game:

– How well can we predict the next word?

I always order pizza with cheese and \_\_\_\_\_

The 33<sup>rd</sup> President of the US was \_\_\_\_\_

I saw a \_\_\_\_\_

mushrooms 0.1

pepperoni 0.1

anchovies 0.01

....

fried rice 0.0001

....

and 1e-100

- Unigrams are terrible at this game. (Why?)
- A better model of a text
  - is one which assigns a higher probability to the word that actually occurs

# Perplexity

The best language model is one that best predicts an unseen test set

- Gives the highest  $P(\text{sentence})$

Perplexity is the inverse probability of the test set, normalized by the number of words:

Chain rule:

$$PP(W) = P(w_1 w_2 \dots w_N)^{-\frac{1}{N}}$$

$$= \sqrt[N]{\frac{1}{P(w_1 w_2 \dots w_N)}}$$

$$PP(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1 \dots w_{i-1})}}$$

**Minimizing perplexity is the same as maximizing probability**

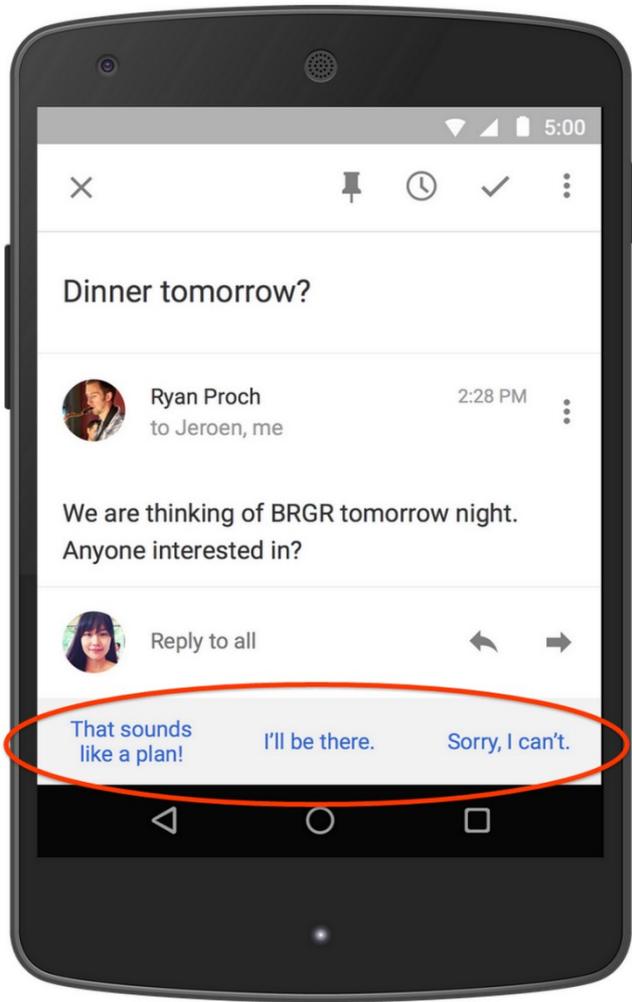
# The Shannon Game intuition for perplexity

- How hard is the task of recognizing digits '0,1,2,3,4,5,6,7,8,9'
  - Perplexity 10
- How hard is recognizing (30,000) names at Microsoft.
  - Perplexity = 30,000
- If a system has to recognize
  - Operator (1 in 4)
  - Sales (1 in 4)
  - Technical Support (1 in 4)
  - 30,000 names (1 in 120,000 each)
  - Perplexity is 53
- Perplexity is weighted equivalent branching factor

# Metric: Perplexity

- Calculate the perplexity of the words in the held-out set *without* doing generation
- **Pros:** Naturally solves multiple-reference problem!
- **Cons:** Doesn't consider decoding or actually generating output.

# Case Study



# Case Study: Smart Reply in Gmail

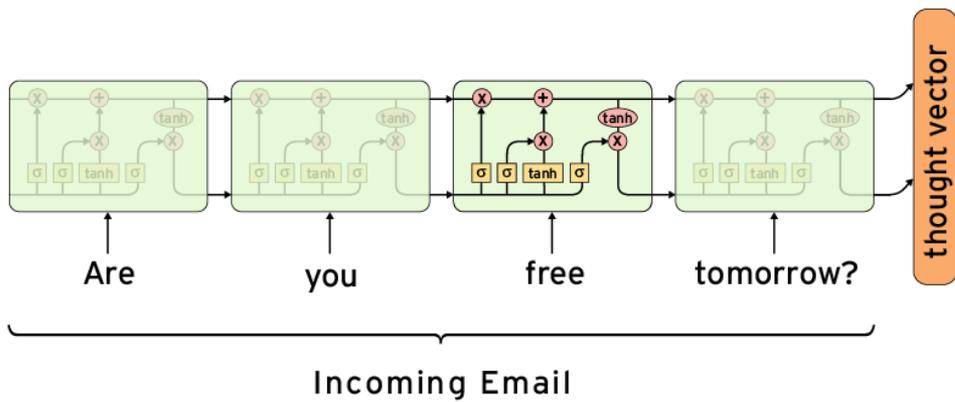


---

# Preprocessing an incoming email

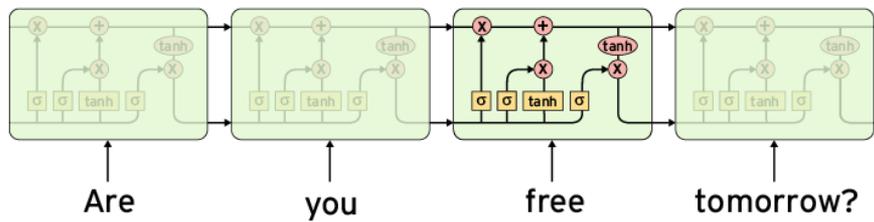
- 
- Language detection
    - Currently handle English, Portuguese, Spanish ... a few more languages are in preparation
  - Tokenization of subject and message body
  - Sentence segmentation
  - Normalization of infrequent words and entities – replaced by special tokens
  - Removal of quoted and forward email portions
  - Removal of greeting and closing phrases (“Hi John”,... “Regards, Mary”)

# ENCODER



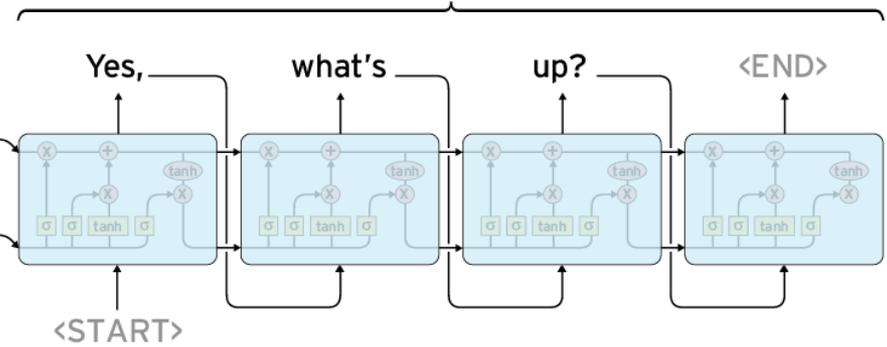
# LSTM translation

ENCODER



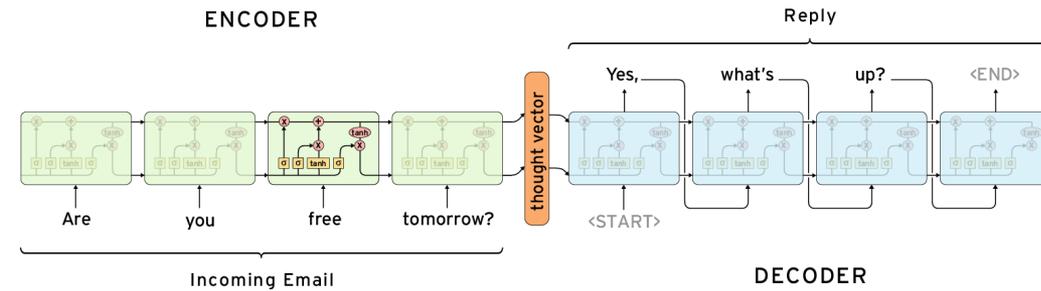
thought vector

Reply



Incoming Email

DECODER



Pick the best suggestions  
(LSTM)



The diagram shows a blue box containing the text "Pick the best suggestions (LSTM)" and the LSTM logo. Two dashed blue lines extend from the top corners of this box towards the encoder and decoder parts of the neural network diagram above, indicating a zoomed-in view of the LSTM components.

---

## Is it worth it?

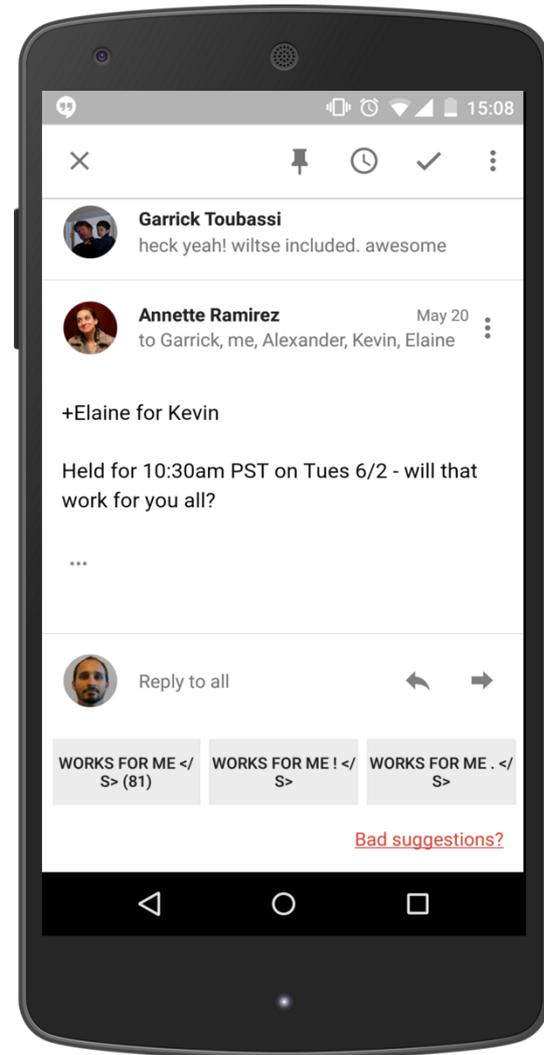
- 
- Precision/accuracy - how well can we guess good replies?
  - Coverage - do most emails have simple, predictable responses?

---

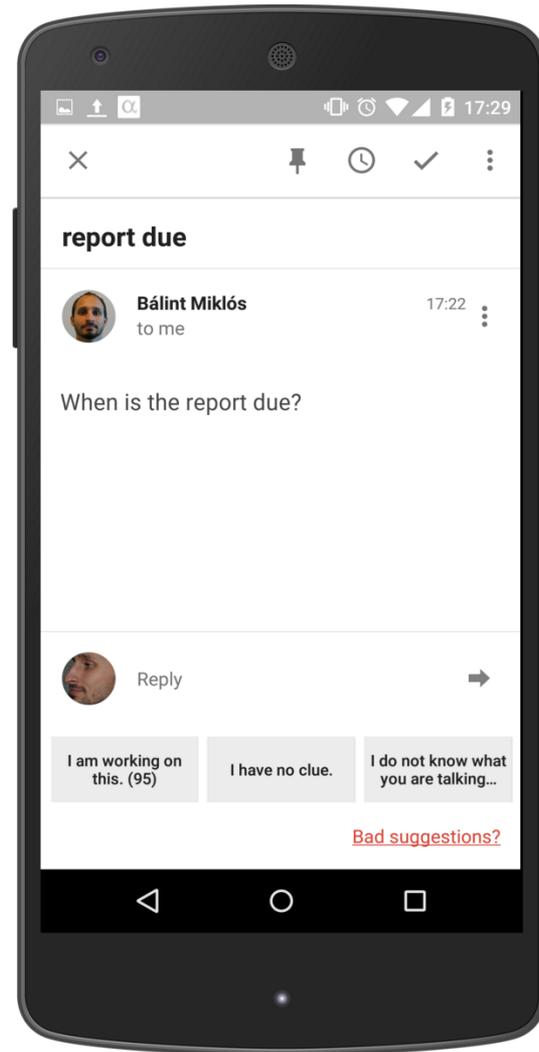
# Metric

- 
- What fraction of the time do users select a suggested reply?
    - How many replies do we suggest? 3
    - Constraint based on user interface, but also users' ability to quickly process choices
  - We get a boost from allowing users to edit responses before sending
    - In early studies, users were nervous that choosing a response would instantly send
    - Careful tuning of this UI gave us bigger gains than a lot of ML tuning

# Some early observations



# Some early observations



# A scoring algorithm != product

- Semantic variation: all 3 suggestions say the same thing ...
  - Can't simply take the 3 highest scoring suggestions
- The “I love you” problem
  - A human can say them, but not a computer ...\*
  - A lot of responses in the training corpus have “I love you”
  - In many cases this isn't appropriate
  - “Family friendliness”
- Sensitivity
  - There are many incoming emails where you don't want the computer to guess replies - Bad news, etc
- \* our expectations of “working” AI are higher than of humans



**Michael Gadberry** @michaelgadberry · 13h

Google Inbox's automated suggested replies are mind-blowingly awesome and accurate. #CheckOutDatStuff #GoogleInbox @inboxbygmail



**Simon Dingle** @SimonDingle · Nov 12

It's like @inboxbygmail has telepathy with its automated responses.



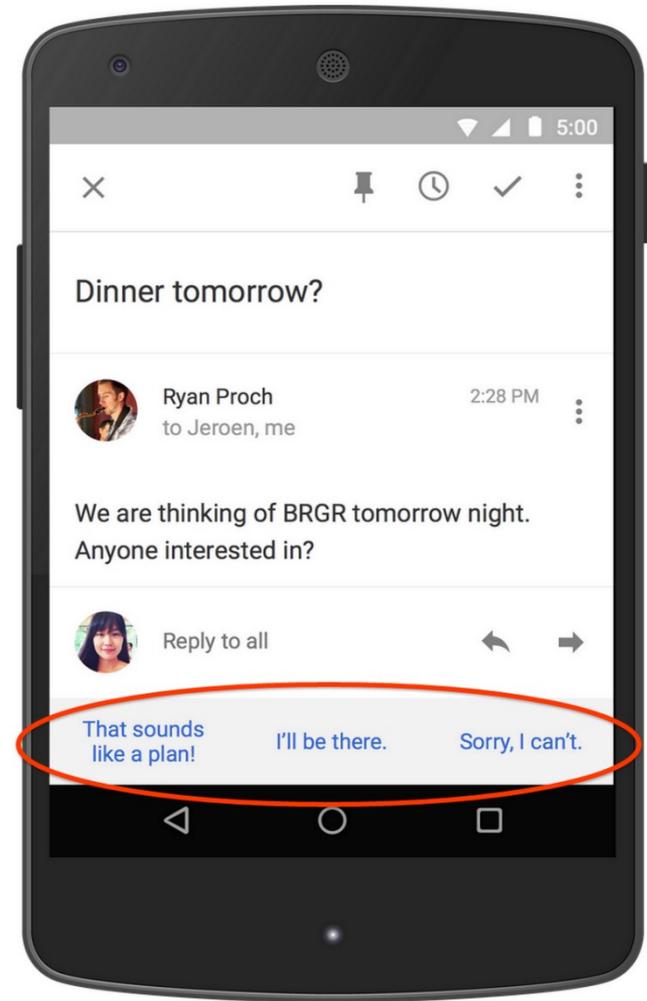
**Tatiana King Jones** @TatianaKing · Nov 12

The new @inboxbygmail auto response choices have been pretty good so far. Have been using them maybe 50% of the time.

# > 10%

of Gmail responses are Smart Replies.

(Users accept computer-generated replies.)



# Encoder-Decoder with different modalities

The encoded conditioning context need not be text, or even a sequence.

# Encoder-Decoder with different modalities

Show and Tell: A Neural Image Caption Generator

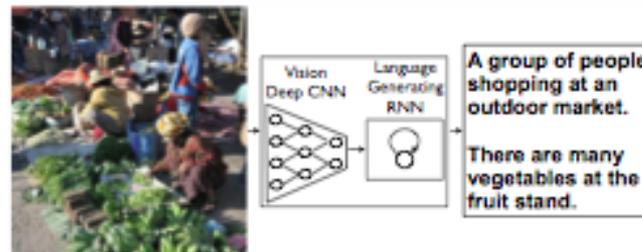
Oriol Vinyals  
Google  
vinyals@google.com

Alexander Toshev  
Google  
toshev@google.com

Samy Bengio  
Google  
bengio@google.com

Dumitru Erhan  
Google  
dumitru@google.com

- Encode: **image** to vector.  
Decode: a sentence describing the image.



This sort-of works.

In my opinion, looks more impressive than really is.

I think it's a man in a business suit standing on a bench.



I am not really confident, but I think it's a man standing on a beach near the water.



I think it's a group of people sitting in front of a crowd.



I am not really confident, but I think it's a close up of a sheep.



Sequence 2 Sequence

Part II: with attention

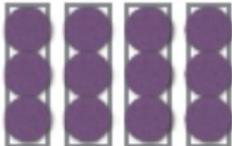
# Sentence Representation

You can't cram the meaning of a whole %&!\$# sentence into a single \$&!#\* vector!



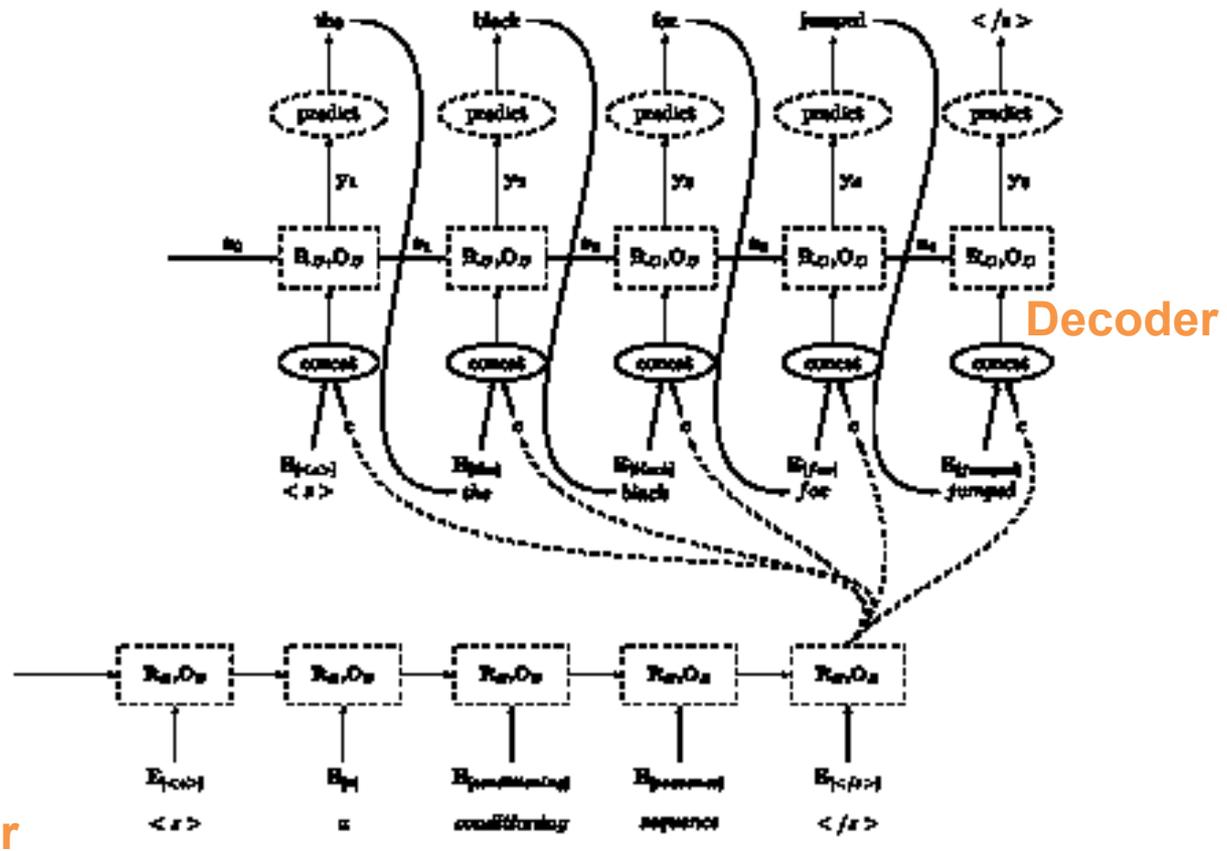
But what if we could use multiple vectors, based on the length of the sentence.

this is an example → 

this is an example → 

# Sequence to Sequence conditioned generation

main idea:  
encoding  
a single vector is  
too restrictive.

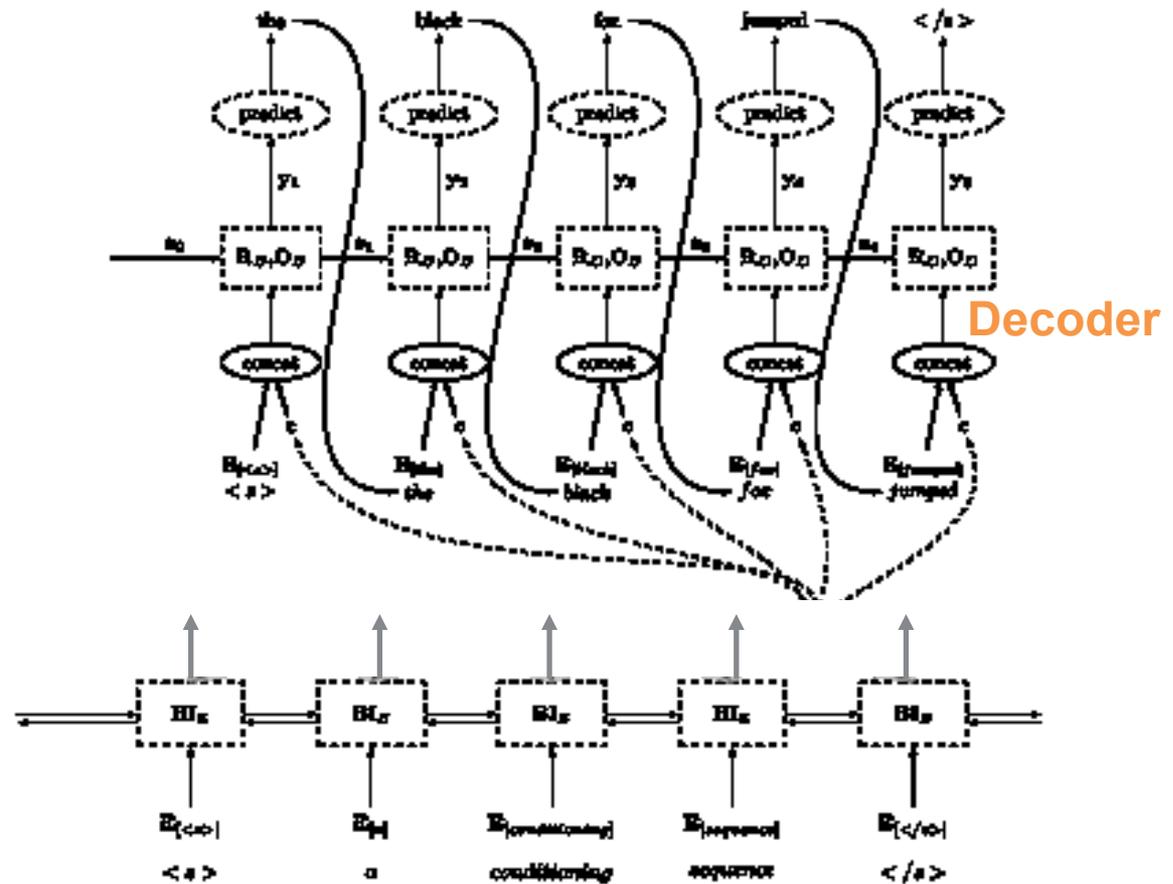


# Attention

- Instead of the encoder producing a single vector for the sentence, it will produce a one vector **for each word**.

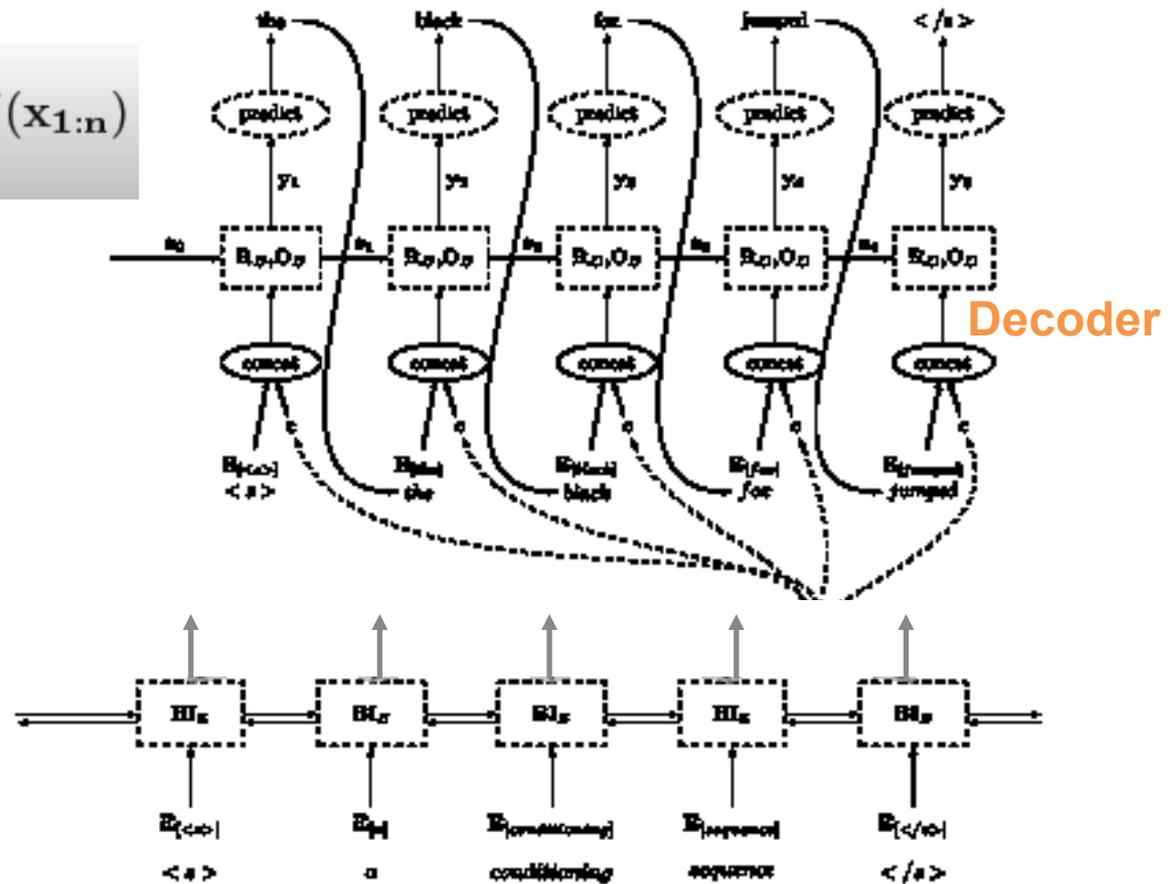


# Sequence to Sequence conditioned generation



# Sequence to Sequence conditioned generation

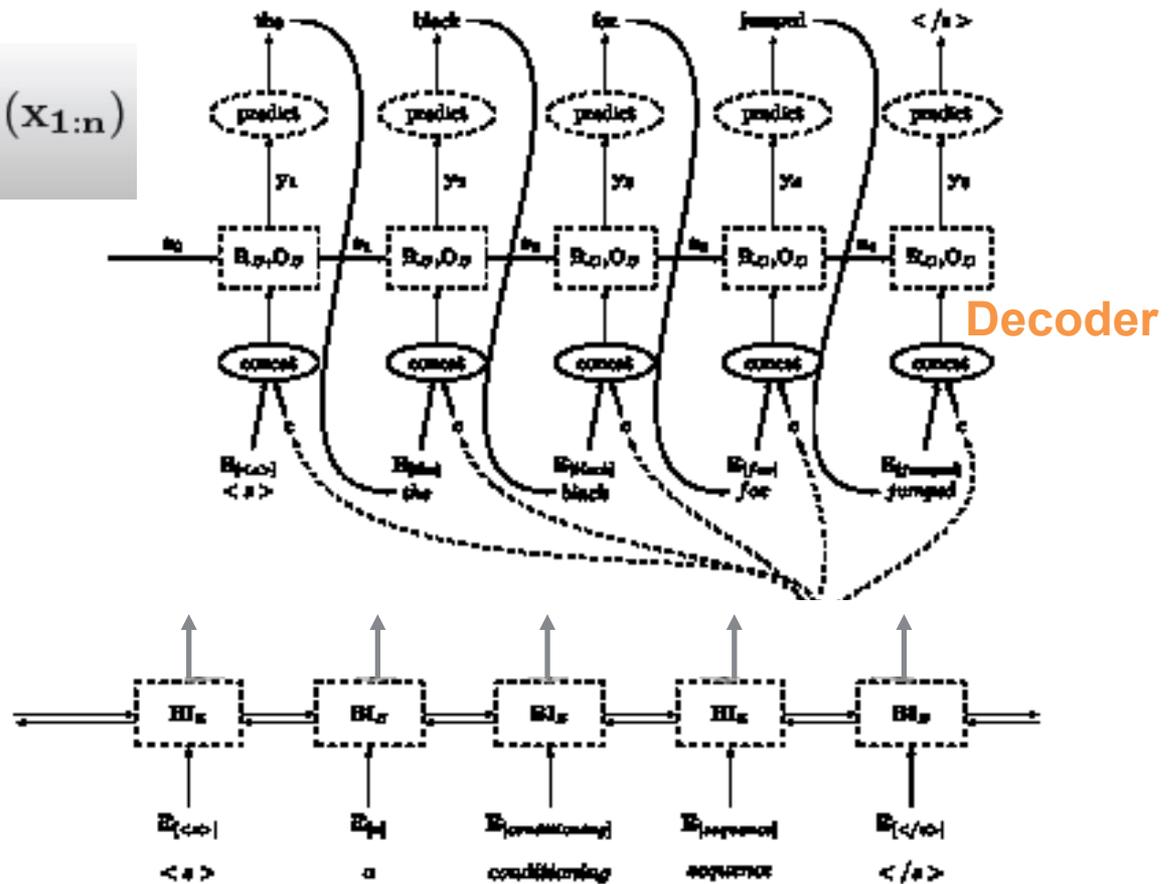
$$c_{1:n} = \text{ENC}(x_{1:n}) = \text{biRNN}^*(x_{1:n})$$



# Sequence to Sequence conditioned generation

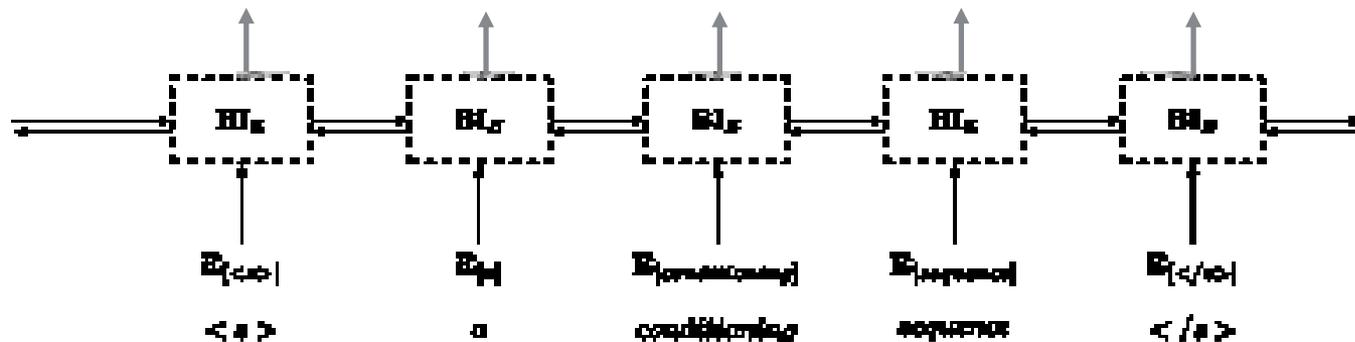
$$c_{1:n} = \text{ENC}(x_{1:n}) = \text{biRNN}^*(x_{1:n})$$

but how do we feed this sequence to the decoder?



# Sequence to Sequence conditioned generation

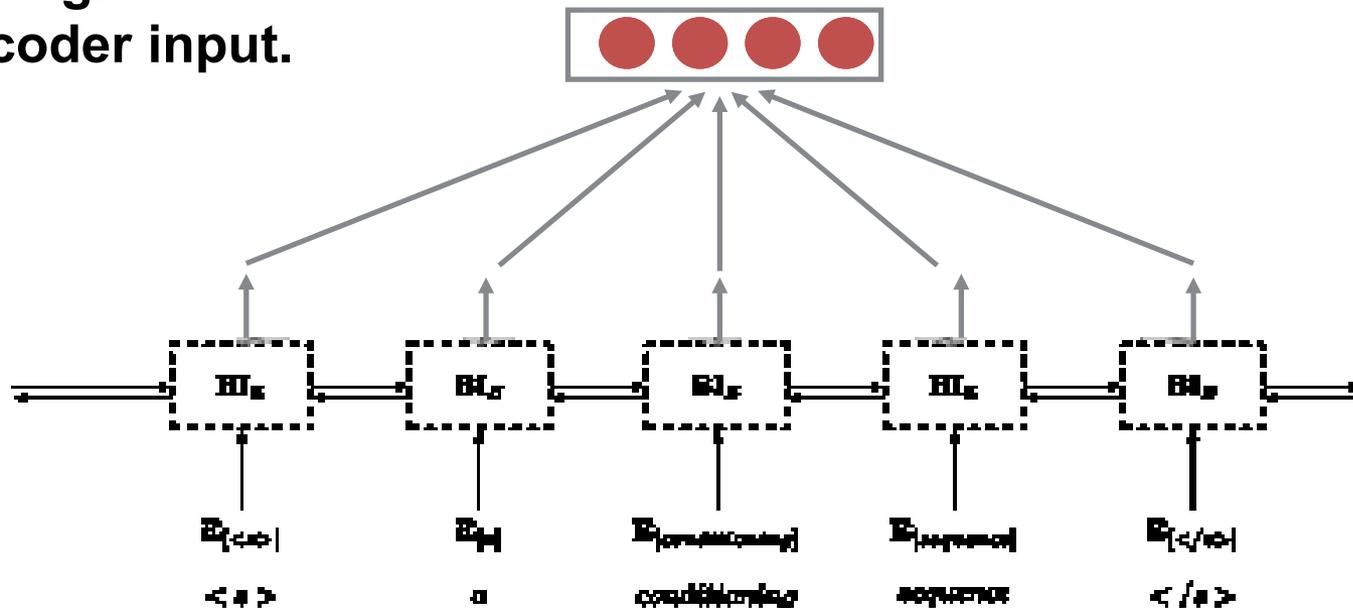
we can combine the different outputs  
into a single vector (attended summary)



# Sequence to Sequence conditioned generation

we can combine the different outputs  
into a single vector (attended summary)

a different single vector  
at each encoder input.



Encoder

$$p(t_{j+1} = k \mid \hat{t}_{1:j}, \mathbf{X}_{1:n}) = f(O(\mathbf{s}_{j+1}))$$

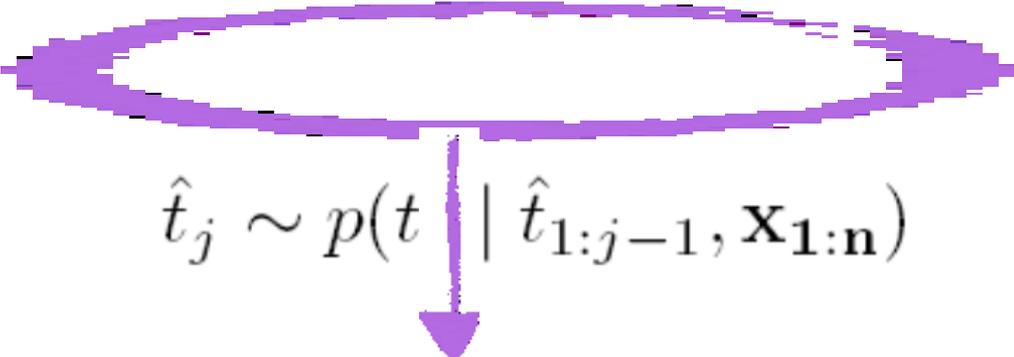
$$\mathbf{s}_{j+1} = R(\mathbf{s}_j, [\hat{\mathbf{t}}_j \circ])$$

$$\circ = \text{attend}(\mathbf{c}_{1:n}, \hat{t}_{1:j})$$

$$\hat{t}_j \sim p(t_j \mid \hat{t}_{1:j-1}, \mathbf{X}_{1:n})$$

$$p(t_{j+1} = k \mid \hat{t}_{1:j}, \mathbf{X}_{1:n}) = f(O(\mathbf{s}_{j+1}))$$

$$\mathbf{s}_{j+1} = R(\mathbf{s}_j, [\hat{\mathbf{t}}_j; \mathbf{c}^j])$$

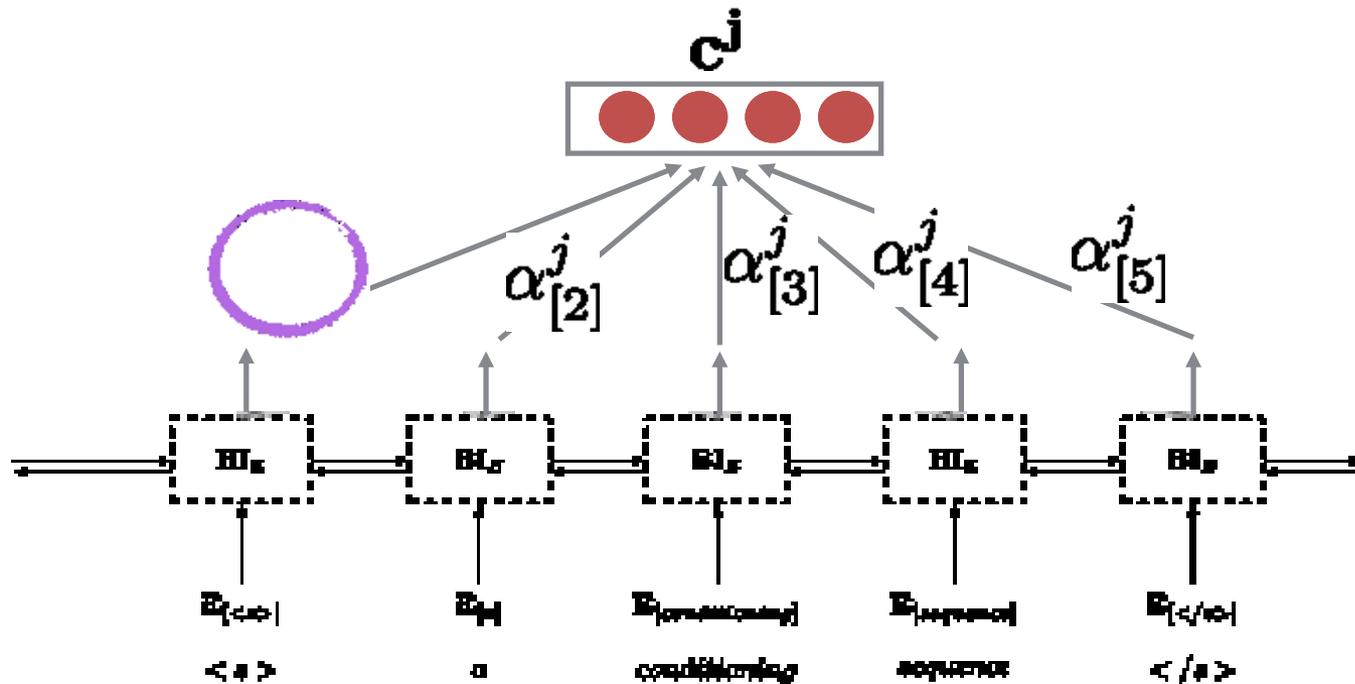

$$\hat{t}_j \sim p(t \mid \hat{t}_{1:j-1}, \mathbf{X}_{1:n})$$

$$\mathbf{c}^j = \sum_{i=1}^n \alpha_{[i]}^j \cdot \mathbf{c}_i$$

# Sequence to Sequence conditioned generation

$$\alpha^j = \text{softmax}(\bar{\alpha}_{[1]}^j, \dots, \bar{\alpha}_{[n]}^j)$$

$$c^j = \sum_{i=1}^n \alpha^j c_i$$



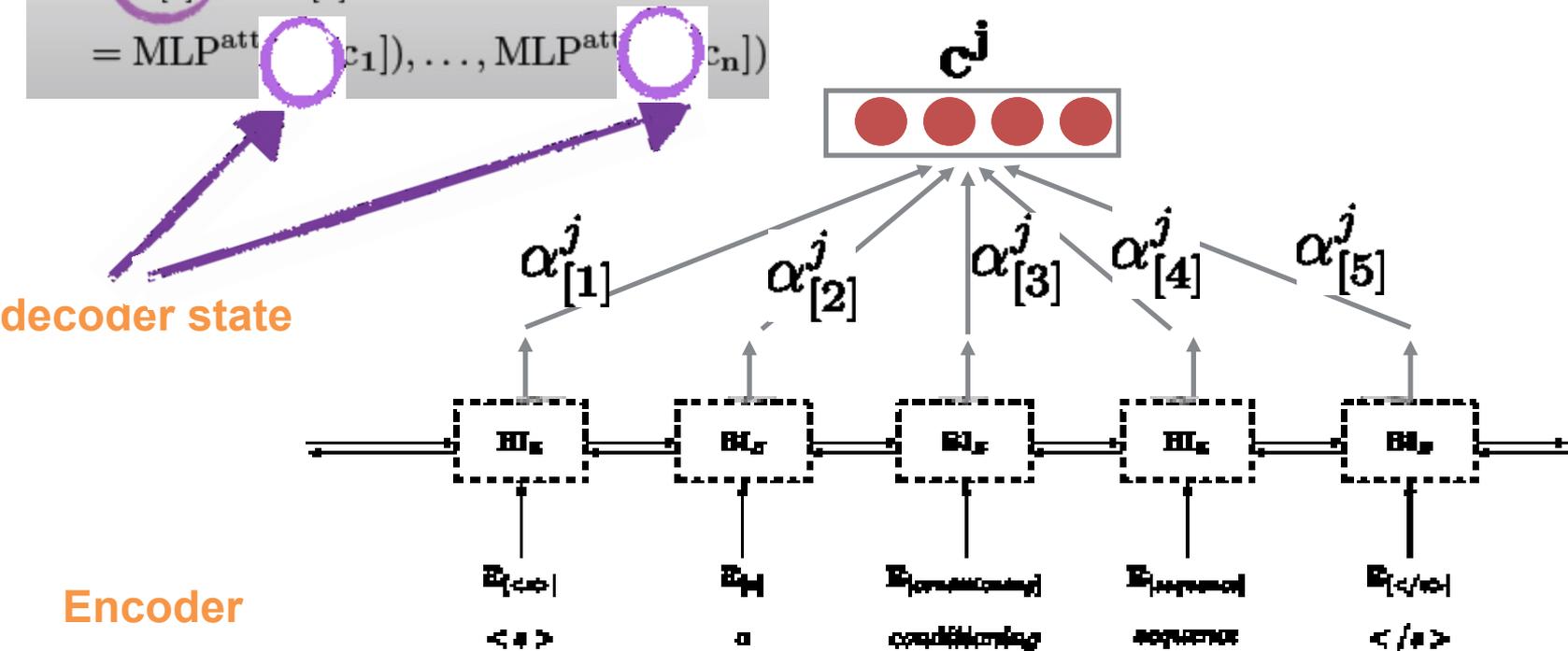
Encoder

# Sequence to Sequence conditioned generation

$$\alpha^j = \text{softmax}(\bar{\alpha}_{[1]}^j, \dots, \bar{\alpha}_{[n]}^j)$$

$$\bar{\alpha}^j = (\bar{\alpha}_{[1]}^j, \dots, \bar{\alpha}_{[n]}^j) = \text{MLP}^{\text{att}}(c_1), \dots, \text{MLP}^{\text{att}}(c_n)$$

$$c^j = \sum_{i=1}^n \alpha_{[i]}^j \cdot c_i$$



# encoder-decoder with attention

$$p(t_{j+1} = k \mid \hat{t}_{1:j}, \mathbf{x}_{1:n}) = f(O_{\text{dec}}(\mathbf{s}_{j+1}))$$

$$\mathbf{s}_{j+1} = R_{\text{dec}}(\mathbf{s}_j, [\hat{t}_j; \mathbf{c}^j])$$

$$\mathbf{c}^j = \sum_{i=1}^n \alpha_{[i]}^j \cdot \mathbf{c}_i$$

$$\mathbf{c}_{1:n} = \text{biRNN}_{\text{enc}}^*(\mathbf{x}_{1:n})$$

$$\alpha^j = \text{softmax}(\bar{\alpha}_{[1]}^j, \dots, \bar{\alpha}_{[n]}^j)$$

$$\bar{\alpha}_{[i]}^j = \text{MLP}^{\text{att}}([\mathbf{s}_j; \mathbf{c}_i])$$

$$\hat{t}_j \sim p(t_j \mid \hat{t}_{1:j-1}, \mathbf{x}_{1:n})$$

$$f(\mathbf{z}) = \text{softmax}(\text{MLP}^{\text{out}}(\mathbf{z}))$$

$$\text{MLP}^{\text{att}}([\mathbf{s}_j; \mathbf{c}_i]) =$$

# encoder-decoder with attention

$$p(t_{j+1} = k \mid \hat{t}_{1:j}, \mathbf{x}_{1:n}) = f(O_{\text{dec}}(\mathbf{s}_{j+1}))$$

$$\mathbf{s}_{j+1} = R_{\text{dec}}(\mathbf{s}_j, [\hat{t}_j; \mathbf{c}^j])$$

$$\mathbf{c}^j = \sum_{i=1}^n \alpha_{[i]}^j \cdot \mathbf{c}_i$$

$$\mathbf{c}_{1:n} = \text{biRNN}_{\text{enc}}^*(\mathbf{x}_{1:n})$$

$$\alpha^j = \text{softmax}(\bar{\alpha}_{[1]}^j, \dots, \bar{\alpha}_{[n]}^j)$$

$$\bar{\alpha}_{[i]}^j = \text{MLP}^{\text{att}}([\mathbf{s}_j; \mathbf{c}_i])$$

$$\hat{t}_j \sim p(t_j \mid \hat{t}_{1:j-1}, \mathbf{x}_{1:n})$$

$$f(\mathbf{z}) = \text{softmax}(\text{MLP}^{\text{out}}(\mathbf{z}))$$

$$\text{MLP}^{\text{att}}([\mathbf{s}_j; \mathbf{c}_i]) =$$

# encoder-decoder with attention

$$p(t_{j+1} = k \mid \hat{t}_{1:j}, \mathbf{x}_{1:n}) = f(O_{\text{dec}}(\mathbf{s}_{j+1}))$$

$$\mathbf{s}_{j+1} = R_{\text{dec}}(\mathbf{s}_j, [\hat{t}_j; \mathbf{c}^j])$$

$$\mathbf{c}^j = \sum_{i=1}^n \alpha_{[i]}^j \cdot \mathbf{c}_i$$

$$\mathbf{c}_{1:n} = \text{biRNN}_{\text{enc}}^*(\mathbf{x}_{1:n})$$

$$\alpha^j = \text{softmax}(\bar{\alpha}_{[1]}^j, \dots, \bar{\alpha}_{[n]}^j)$$

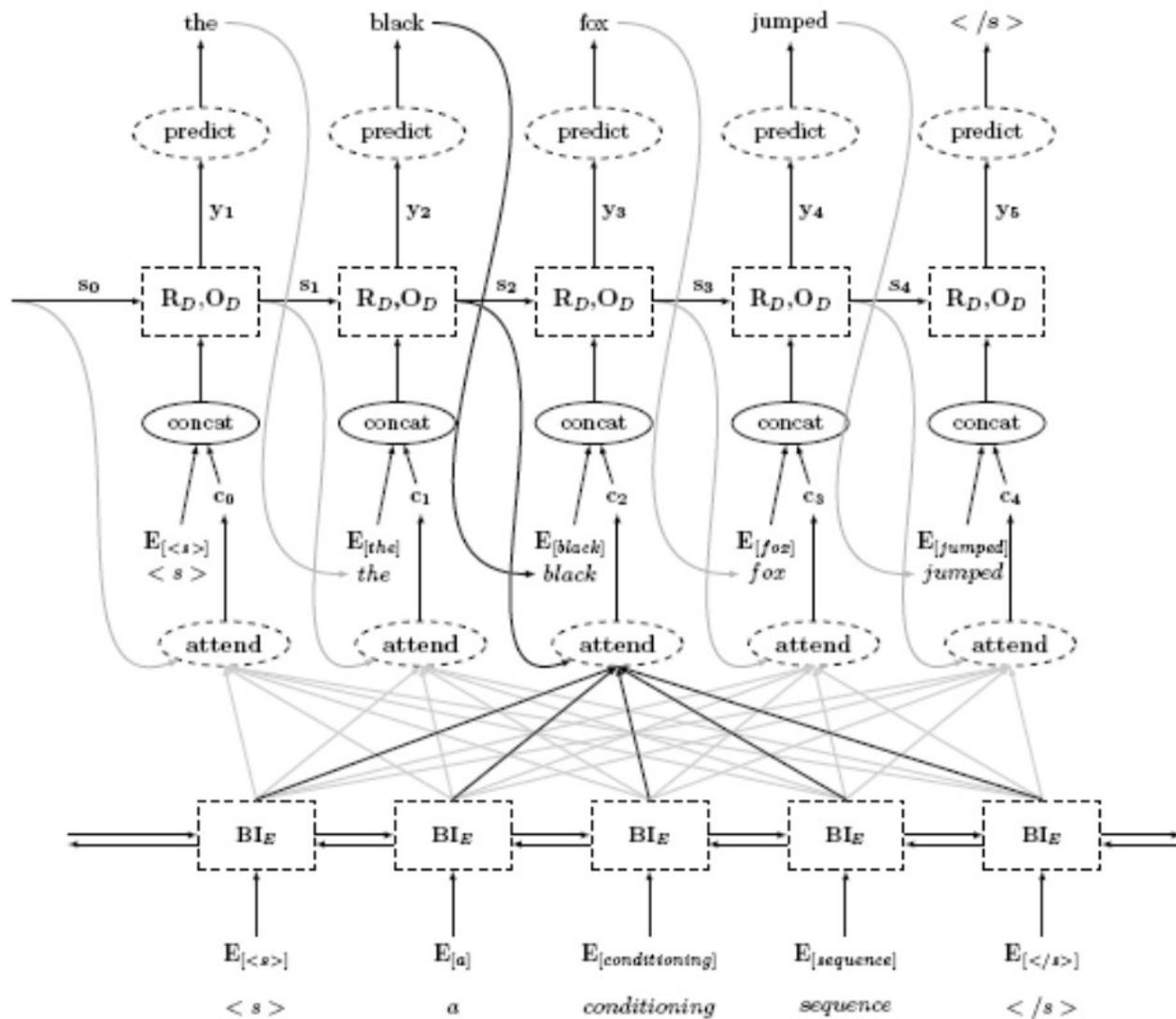
$$\bar{\alpha}_{[i]}^j = \text{MLP}^{\text{att}}([\mathbf{s}_j; \mathbf{c}_i])$$

$$\hat{t}_j \sim p(t_j \mid \hat{t}_{1:j-1}, \mathbf{x}_{1:n})$$

$$f(\mathbf{z}) = \text{softmax}(\text{MLP}^{\text{out}}(\mathbf{z}))$$

$$\text{MLP}^{\text{att}}([\mathbf{s}_j; \mathbf{c}_i]) =$$

# encoder-decoder with attention



# encoder-decoder with attention

- Encoder encodes a sequence of vectors,  $c_1, \dots, c_n$
- At each decoding stage, an MLP assigns a relevance score to each Encoder vector.
- The relevance score is based on  $c_i$  and the state  $s_j$
- Weighted-sum (based on relevance) is used to produce the conditioning context for decoder step  $j$ .

# encoder-decoder with attention

- Decoder "pays attention" to different parts of the encoded sequence at each stage.
- The attention mechanism is "soft" -- it is a mixture of encoder states.
- The encoder acts as a read-only memory for the decoder
- The decoder chooses what to read at each stage

# Attention

- Attention is very effective for sequence-to-sequence tasks.
- Current state-of-the-art systems all use attention.  
(this is basically how Machine Translation works)
- Attention makes models somewhat more ~interpretable.
- (we can see where the model is "looking" at each stage of the prediction process)



# Attention

in the evening until 21:00, there was a further 5mm rain on the town, after 6:00 pm, which had already dropped to Sunday during the night.  
am Abend bis 21 Uhr fielen weitere 5mm Regen auf die Stadt, nach 6:00 mm, die bereits in der Nacht zum Sonntag niedriger gegangen waren.

since then, the island authorities have tried to put an end to the illegal behaviour of non-alcoholic tourists in Magaluf by minimizing the number of participants in the notorious alcohol-free bar.  
die Inselbehörden haben seither versucht, das ordnungswidrige Verhalten alkoholisierter Urlauber in Magaluf zu stoppen, indem die Anzahl der Teilnehmer an den berüchtigten alkoholgetränkten Kneipen minimiert wurde.

## Attention is not Explanation

**Sarthak Jain**  
Northeastern University  
jain.sar@husky.neu.edu

**Byron C. Wallace**  
Northeastern University  
b.wallace@northeastern.edu

# Complexity

- Encoder decoder:
- Encoder-decoder with attention:

# Complexity

- Encoder decoder:  $O(n+m)$
- Encoder-decoder with attention:  $O(nm)$

# Beyond Seq2Seq

- Can think of a general design pattern in neural nets:
  - **Input**: sequence, query
    - **Encode** the input into a sequence of vectors
    - **Attend** to the encoded vectors, based on query (weighted sum, determined by query)
    - **Predict** based on the attended vector

# Attention Functions

**v**: attended vec, **q**: query vec

$\text{MLP}^{\text{att}}(\mathbf{q};\mathbf{v})=$

- Additive Attention:  $\text{ug}(\mathbf{W}^1\mathbf{v} + \mathbf{W}^2\mathbf{q})$
- Dot Product:  $\mathbf{v} \cdot \mathbf{q}$
- Bilinear attention:  $\mathbf{v}^\top \mathbf{W} \mathbf{q}$

# Additive vs Multiplicative

While the two are similar in theoretical complexity, dot-product attention is much faster and more space-efficient in practice, since it can be implemented using highly optimized matrix multiplication code.

While for small values of  $d_k$  the two mechanisms perform similarly, additive attention outperforms dot product attention without scaling for larger values of  $d_k$  [3]. We suspect that for large values of  $d_k$ , the dot products grow large in magnitude, pushing the softmax function into regions where it has extremely small gradients<sup>4</sup>. To counteract this effect, we scale the dot products by  $\frac{1}{\sqrt{d_k}}$ .

$$\frac{\mathbf{v} \cdot \mathbf{q}}{\sqrt{d_k}}$$

$d_k$  is the dimensionality of  $q$  and  $v$

---

Attention Is All You Need

---

Ashish Vaswani\*  
Google Brain  
avaswani@google.com

Noam Shazeer\*  
Google Brain  
noam@google.com

Niki Parmar\*  
Google Research  
nikip@google.com

Jakob Uszkoreit\*  
Google Research  
usz@google.com

Llion Jones\*  
Google Research  
llion@google.com

Aidan N. Gomez\*<sup>†</sup>  
University of Toronto  
aidan@cs.toronto.edu

Lukasz Kaiser\*  
Google Brain  
lukasz.kaiser@google.com

Ilia Polosukhin\*<sup>†</sup>  
ilia.polosukhin@gmail.com

# Key-Value Attention

- Split  $v$  into two vectors  $v=[v_k;v_v]$ 
  - $v_k$ : key vector
  - $v_v$ : value vector
- Use key vector for computing attention  
 $\text{MLP}^{\text{att}}(q;v)= \text{ug}(\mathbf{W}^1v_k + \mathbf{W}^2q)$  //additive
- Use value vector for computing attended summary

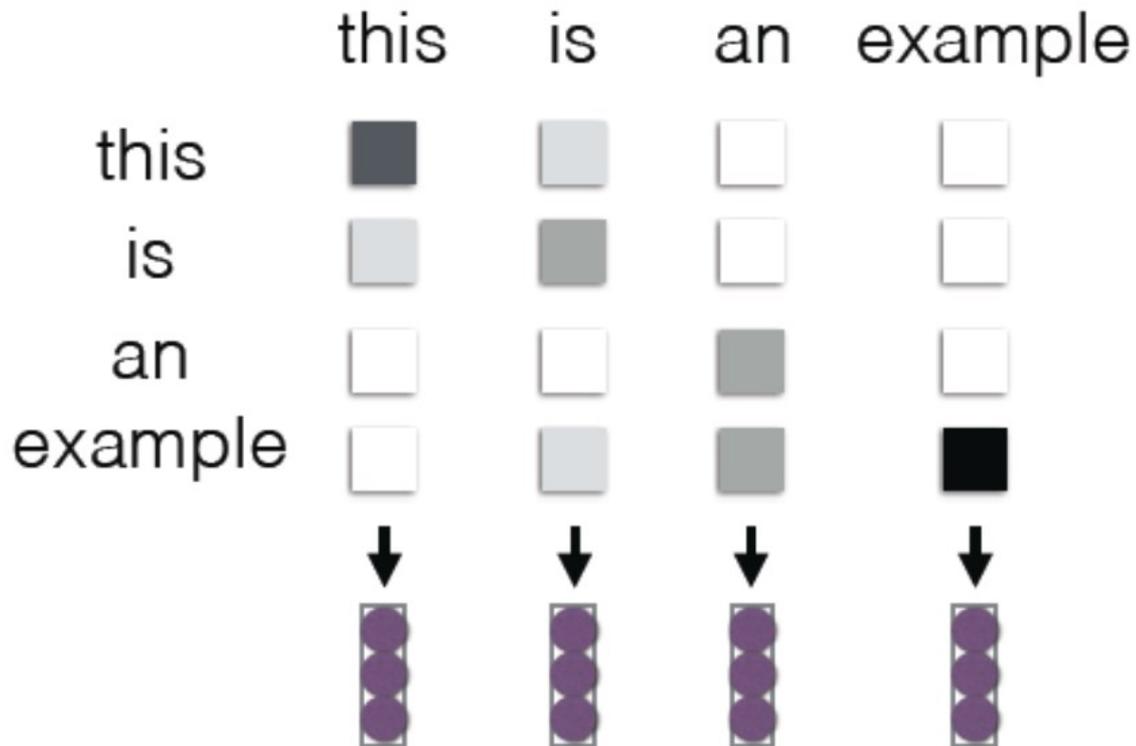
$$v^j = \sum_{i=1}^n \alpha_{[i]}^j \cdot (v_v)_i$$

# Multi-head Key-Value Attention

- For each head
  - Learn different projection matrices  $\mathbf{W}_q$ ,  $\mathbf{W}_k$ ,  $\mathbf{W}_v$
- $\text{MLP}^{\text{att}}(q;v) = [(v_k \mathbf{W}_k) \cdot (q \mathbf{W}_q)] / \text{sqrt}(d_k)$
- For summary use  $v_v \mathbf{W}_v$  (instead of  $v_v$ )
- Train many such heads and
  - use  $\text{aggr}(\text{all such attended summaries})$

# Self-attention/Intra-attention

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



# Do we “need” an LSTM?

- **They are slow**

- Sequential nature of computation makes it tough to optimize operations on GPUs
- Contrast to CNNs: convolutions completely parallelizable

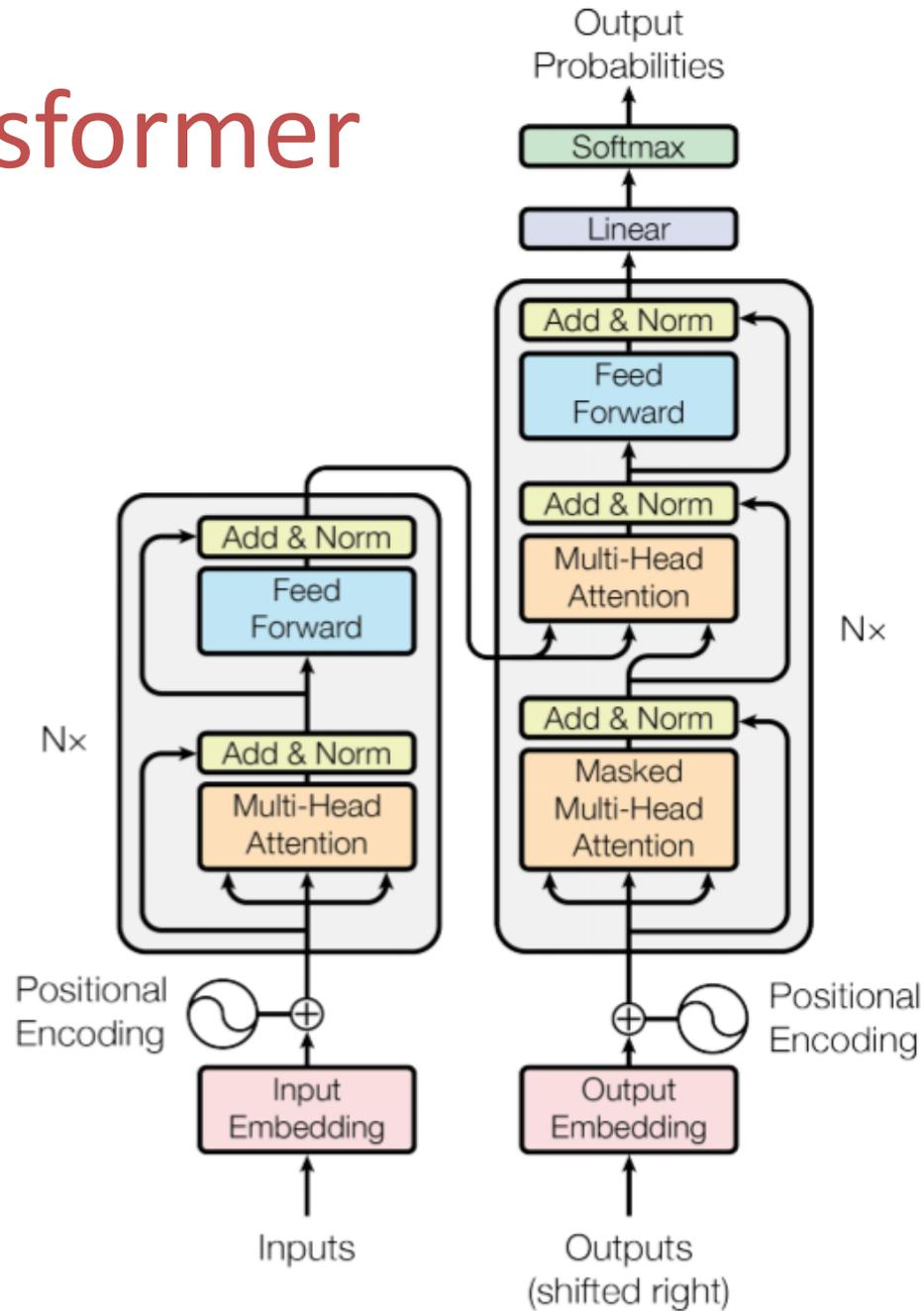
- **They are not deep**

- Vanishing gradient: aggravated for deeper networks
- Less depth → low compositionality power of the network
- Deepest LSTM networks are 8 layered
  - in-contrast to 50-layered Resnets

- **They don't transfer well**

- Networks trained on one task, do not generalize well to even other datasets in the same task, not to speak about other tasks
- ImageNet-trained ResNet fine-tuned on many other datasets

# Transformer



# BERT: Transformer + PreTraining

- In NLP, we are interested in solving a variety of end tasks - Question Answering, Search, etc.
- One approach - train neural models from scratch
- Issue - this involves two things
  - Modelling of Syntax and Semantics of the language
  - Modelling of the end-task
- Pretraining - Learns the modelling of syntax and semantics - through another task
- So the current model can focus exclusively on modelling of end-task

# Pretraining - Masked Language Modelling

- How to pretrain?
- Which base task to choose:
  - Must have abundant data available
  - Must require learning of syntax and semantics
- Language Modelling (Self-supervision)
  - Does not require human annotated labels - abundance of sentences
  - Requires understanding of both syntax and semantics to predict the next word in sentence

# Summary

- RNNs are very capable learners of sequential data.
  - $n \rightarrow 1$ : (bi)RNN acceptor
  - $n \rightarrow n$ : biRNN (transducer)
  - $1 \rightarrow m$ : conditioned generation (conditioned LM)
  - $n \rightarrow m$ : conditioned generation (encoder-decoder)
  - $n \rightarrow m$ : encoder-decoder with attention
- Transformer
  - More scalable than RNN
  - May slowly replace RNNs