e.g., sentiment classification



**positive**

Sentence encoding

How to compute
sentence encoding?

<u>Basic way:</u>
Use final hidden state

*equals*

$h_1$  $h_2$  $h_n$

$\rightarrow$ mean

overall  I  enjoyed  the  movie  a  lot

$x_1$  $x_n$

e.g., sentiment classification

CE loss

# classes

MLP

$V^T$

Sentence encoding

How to compute sentence encoding?

Usually better:
Take element-wise max or mean of all hidden states

element-wise max/mean



overall    I    enjoyed    the    movie    a    lot

Sentence Similarity

$S_1$

He is smart.

RNN

encoding

$S_2$

He is a wise man.

RNN

encoding

Similarity over these encodings

Siamese Architecture

Score $(0,1)$

$W$

$h_a$

$h_b$

He   is   Smart

$S_a$

He is a wise man
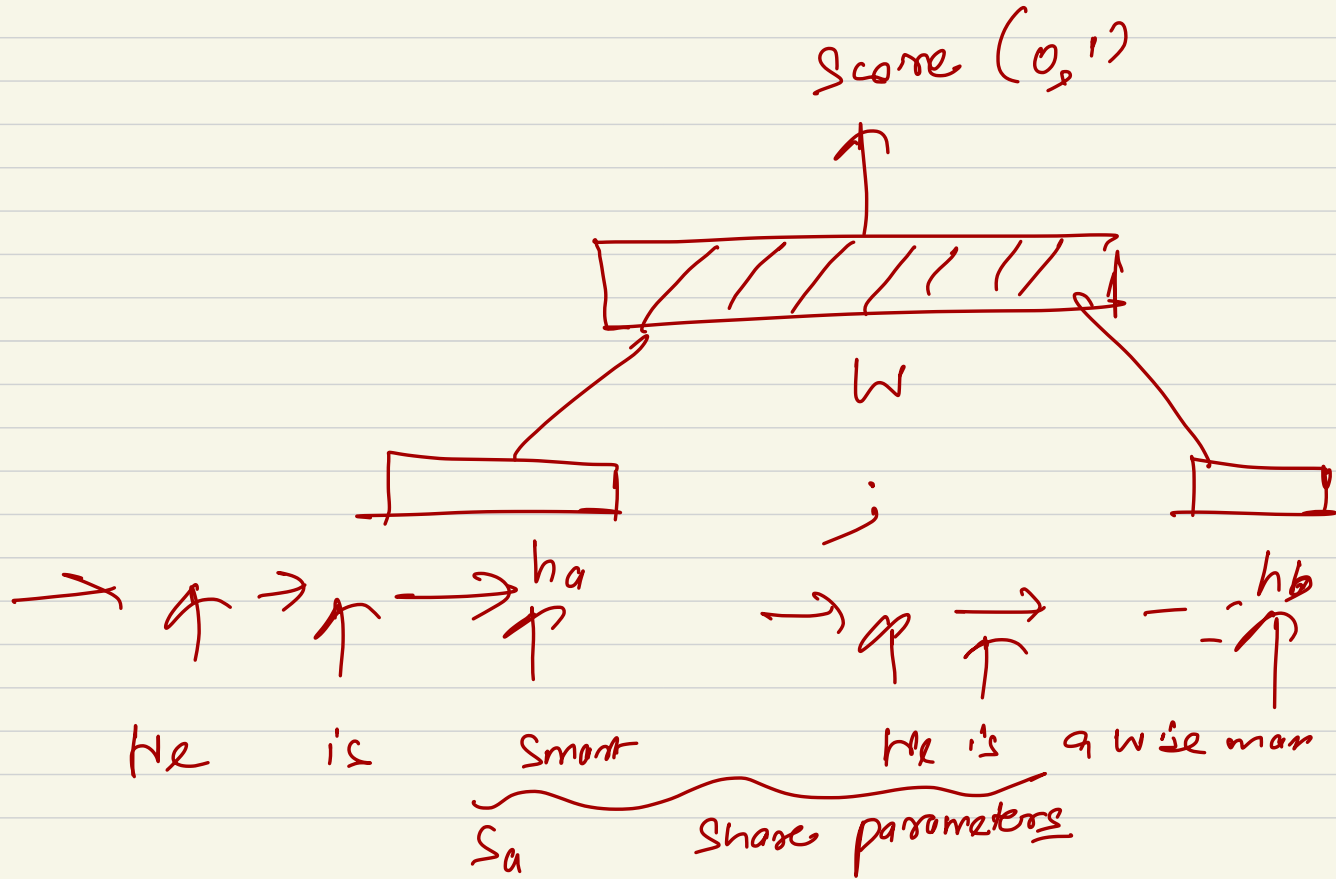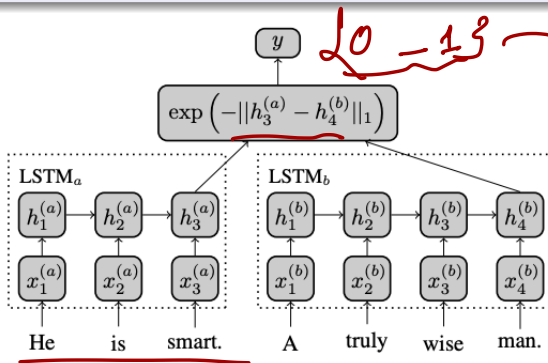
Shared parameters

RNN / Copy

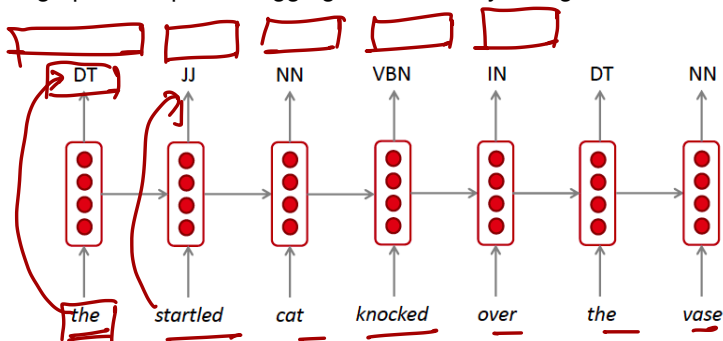# RNNs can be used to find sentence similarity

**Siamese Architecture**

- Pass both sentences through RNNs: the parameters are shared across these
- Use (a function of) the difference between the final hidden states

Sequence labels

e.g., part-of-speech tagging, named entity recognition

output at every step



| DT | JJ | NN | VBN | IN | DT | NN |
|----|----|----|----|----|----|----|

*the* | *startled* | *cat* | *knocked* | *over* | *the* | *vase*

Size of o/p vector ?

| POS tags |

Nikon Coolpix has better image quality than Cannon

adjective

postagees

vector

POS

Embedding Layer (Word + POS)

bidirection

1-hot

100100

LSTM Layer

300-d POS embed.

None
Product 1
Product 2
Aspect
Predicate

initialised embeddings

1-hot vectos

trained word embeddgs

$$x^{(1)} \ - \ - \ x^{(t-1)} \ \overset{ht}{x^{(t)}} \ x^{t+1} \ - \ - \ x^{(n)}$$

**What we have seen till now?**

- The state at time $t$ only captures information from the past $x^{(1)}, \ldots, x^{(t-1)}$, and the present input $x^{(t)}$

- Some models also allow information from past $y$ values to affect the current state when the $y$ values are available

Text classification $\quad$ LM?

$x_1 \quad\quad x_n \quad$ Given $x_1 - - x_{t-1}$ predict $x_t$

POS tagging

# Bidirectional RNNs

## What we have seen till now?
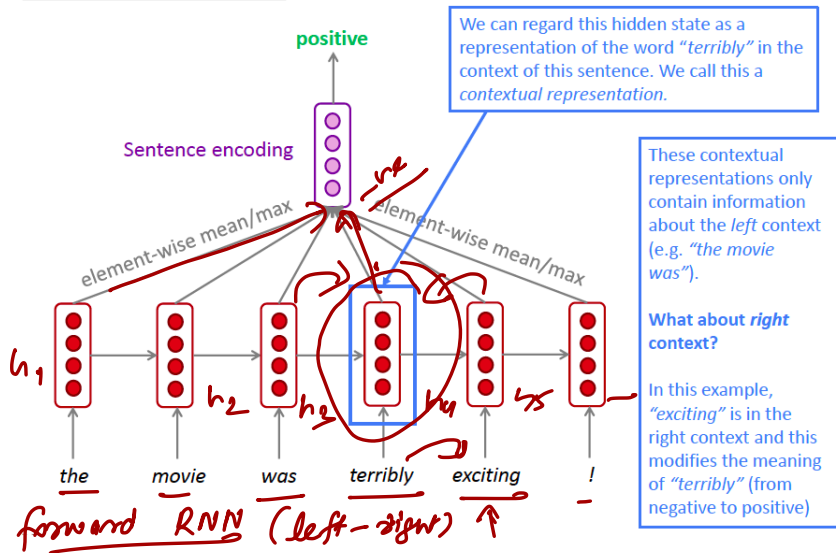
- The state at time $t$ only captures information from the past $x^{(1)}, \ldots, x^{(t-1)}$, and the present input $x^{(t)}$
- Some models also allow information from past $y$ values to affect the current state when the $y$ values are available
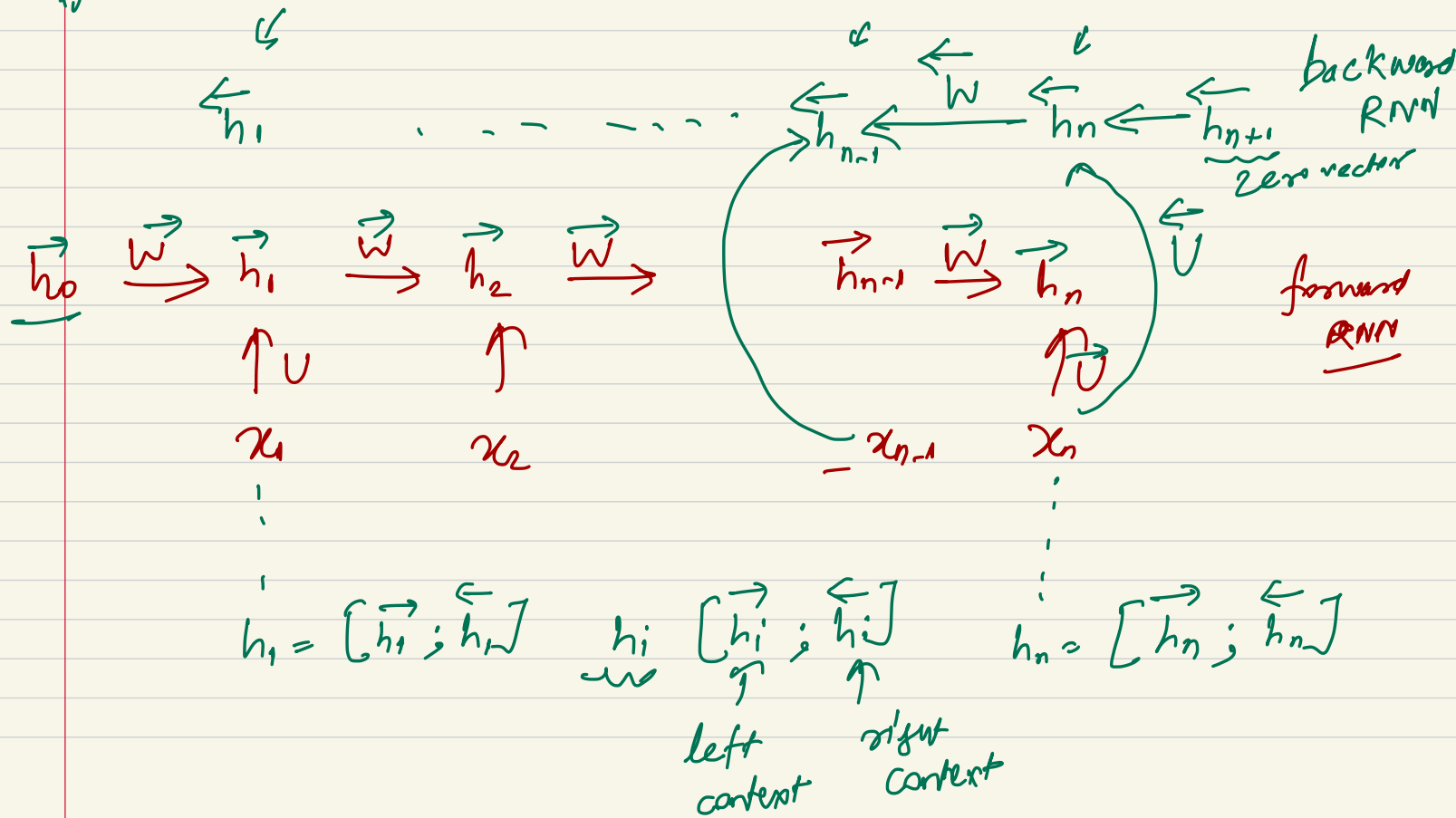
However, in many applications we want to output a prediction of $y^{(t)}$ which may depend on the whole input sequence. e.g., sentiment analysis.

# Bidirectional RNNs: motivation for sentiment classification

**Task: Sentiment Classification**

**positive**

**Sentence encoding**

element-wise mean/max

element-wise mean/max

We can regard this hidden state as a representation of the word *"terribly"* in the context of this sentence. We call this a *contextual representation.*

These contextual representations only contain information about the *left* context (e.g. *"the movie was"*).

**What about *right* context?**

In this example, *"exciting"* is in the right context and this modifies the meaning of *"terribly"* (from negative to positive)

the    movie    was    terribly    exciting    !

forward RNN (left-right)

# Bidirectional RNN



$\overleftarrow{h_1}$ ·········· $\overleftarrow{h_{n-1}}$ $\xleftarrow{W}$ $\overleftarrow{h_n}$ $\leftarrow$ $\overleftarrow{h_{n+1}}$  backward RNN

zero vector

$\overrightarrow{h_0}$ $\xrightarrow{\overrightarrow{W}}$ $\overrightarrow{h_1}$ $\xrightarrow{\overrightarrow{W}}$ $\overrightarrow{h_2}$ $\xrightarrow{\overrightarrow{W}}$ ⋯ $\overrightarrow{h_{n-1}}$ $\xrightarrow{\overrightarrow{W}}$ $\overrightarrow{h_n}$     forward RNN

$\uparrow U$     $\uparrow$     $\uparrow U$

$x_1$     $x_2$     $x_{n-1}$     $x_n$

$$h_1 = [\overrightarrow{h_1} ; \overleftarrow{h_1}] \qquad h_i \quad [\overrightarrow{h_i} ; \overleftarrow{h_i}] \qquad h_n = [\overrightarrow{h_n} ; \overleftarrow{h_n}]$$

left context     right context

# Bidirectional RNNs



This contextual representation of "terribly" has both left and right context!

Not for LM

Concatenated hidden states

Backward RNN

Forward RNN

the    movie    was    terribly    exciting    !

# Bidirectional RNNs

On timestep $t$:

This is a general notation to mean "compute one forward step of the RNN" – it could be a vanilla, LSTM or GRU computation.

Forward RNN $\quad \overrightarrow{\boldsymbol{h}}^{(t)} = \text{RNN}_{\text{FW}}(\overrightarrow{\boldsymbol{h}}^{(t-1)}, \boldsymbol{x}^{(t)})$

Backward RNN $\quad \overleftarrow{\boldsymbol{h}}^{(t)} = \text{RNN}_{\text{BW}}(\overleftarrow{\boldsymbol{h}}^{(t+1)}, \boldsymbol{x}^{(t)})$

Generally, these two RNNs have separate weights

Concatenated hidden states $\quad \boldsymbol{h}^{(t)} = [\overrightarrow{\boldsymbol{h}}^{(t)}; \overleftarrow{\boldsymbol{h}}^{(t)}]$

We regard this as "the hidden state" of a bidirectional RNN. This is what we pass on to the next parts of the network.

# Bidirectional RNNs: simplified diagram



The two-way arrows indicate bidirectionality and the depicted hidden states are assumed to be the concatenated forwards+backwards states.

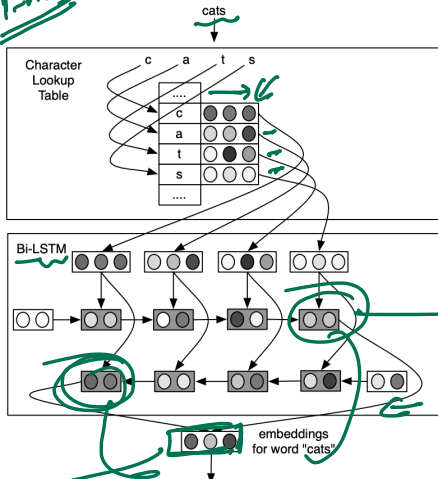Bi-RNNs are only applicable if you have access to the entire input sequence

- thus not applicable to Language Modeling where only the left context is available

embeddings
for word "cats"

final hidden b/w

**Handwritten annotations:**

1-hot

OOV words

- what are character embeddings?

Cats

C a t s

- How do I train my network?

final hidden f/w

y-ŷ MSE Compare with loss "Cats"

Concatenate

## Multi-layer RNNs

The hidden states from RNN layer $i$ are the inputs to RNN layer $i+1$



RNN layer 3

RNN layer 2

RNN layer 1

the    movie    was    terribly    exciting    !

# Vanilla RNN

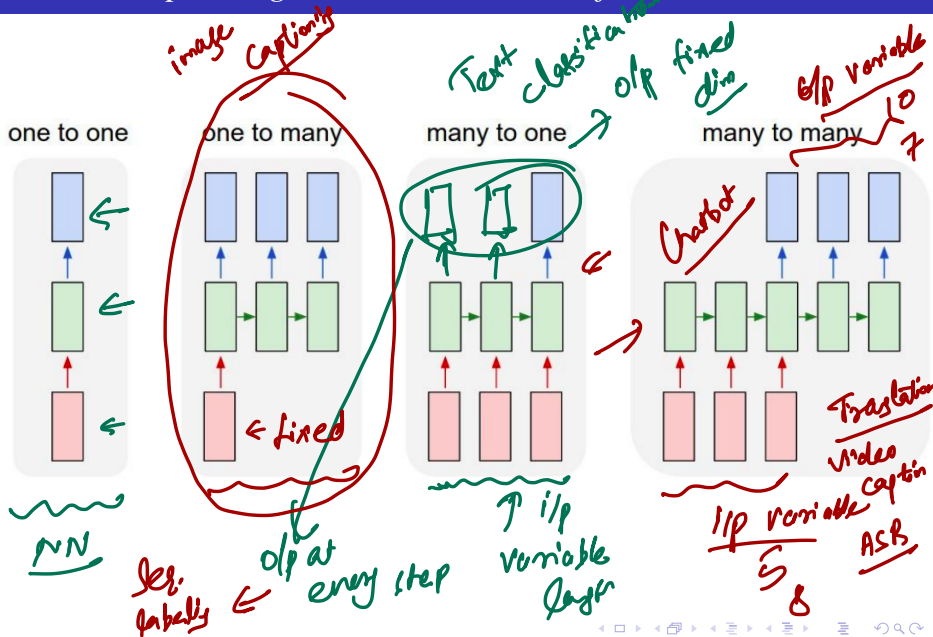$$h_{t-1} \longrightarrow \uparrow \qquad h_t$$

$$x_t$$

GRU, LSTMs   (complicated units)

Bi- RNN        Bi- LSTM

Multi- layer

Multi- layer   Bi-LSTM

# What other paradigm can RNN be used for?

**one to one**

**one to many** — image Captions, output at every step, fixed

**many to one** — Text Classification, o/p fixed dim, i/p variable length

**many to many** — o/p variable, Chatbot, Translation, video captioning, ASR, i/p variable

NN

reg. labels

Seq2Seq p Enc. dec.

Encoder — decoder

# *Sequence to Sequence Models*

Pawan Goyal

CSE, IIT Kharagpur

NLP - Interaction Hour

- Input sequence to a fixed-sized vector
- Input sequence to an output sequence of the same length

- Input sequence to a fixed-sized vector
- Input sequence to an output sequence of the same length

### Any other constraint?

Mapping input sequence to an output sequence, not necessarily of the same length

machine translation, question answering, chatbots, summarization, ...

Machine Translation (MT) is the task of translating a sentence $x$ from one language (the source language) to a sentence $y$ in another language (the target language)

French

| | |
|---|---|
| x: | L'homme est né libre, et partout il est dans les fers |

| | |
|---|---|
| y: | Man is born free, but everywhere he is in chains |

English

- Rousseau

$x$ (for)
Source

$\longrightarrow$

$y$ (En)
target

$\rightarrow$ Choose the word with max. prob.

vocab. of target

$x_1$ — — $x_n$ $\longrightarrow$

$y_1$ $\nearrow$ — · — $y_m$

$h_0 \rightarrow h_1 \rightarrow \rightarrow [h_n] \text{- - - - - - - - - -}$

generate
Sentence

How do I
Training: train params?

Sentence
encoding
for $x$

generate
Sentence

<START>

Run-time

test-phase

$\hookrightarrow$ How do I use my learnt parameters
to generate target sequence?

Generation using
LM

# Sequence-to-sequence architecture

Also known as encoder-decoder architecture

- Input sequence $X = (x^{(1)}, \ldots, x^{(n_x)})$
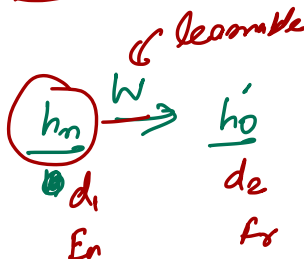- Output sequence $Y = (y^{(1)}, \ldots, y^{(n_y)})$

Also known as encoder-decoder architecture

- Input sequence $X = (x^{(1)}, \ldots, x^{(n_x)})$
- Output sequence $Y = (y^{(1)}, \ldots, y^{(n_y)})$
- Encoder (reader/input) RNN: Emits the context $C$, a vector summarizing the input sequence, usually as a simple function of its final hidden state

*Source encoding*

*Summarization*

*C learnable*

$h_n$ $\xrightarrow{W}$ $h_0'$

$d_1$     $d_2$

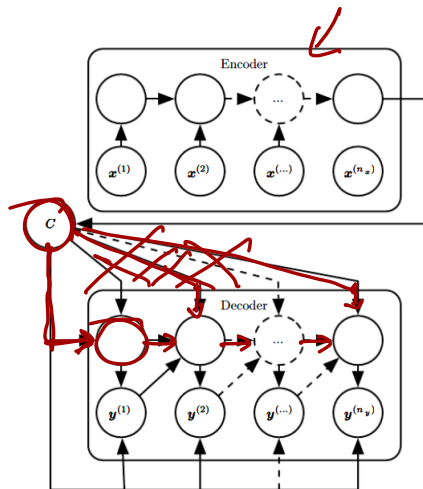$En$     $Fr$

## Sequence-to-sequence architecture

Also known as encoder-decoder architecture

- Input sequence $X = (x^{(1)}, \ldots, x^{(n_x)})$
- Output sequence $Y = (y^{(1)}, \ldots, y^{(n_y)})$
- Encoder (reader/input) RNN: Emits the context $C$, a vector summarizing the input sequence, usually as a simple function of its final hidden state
- Decoder (writer/output) RNN: Is conditioned on the context $C$ to generate the output sequence.

# Sequence-to-sequence architecture

Also known as encoder-decoder architecture

- Input sequence $X = (x^{(1)}, \ldots, x^{(n_x)})$
- Output sequence $Y = (y^{(1)}, \ldots, y^{(n_y)})$
- Encoder (reader/input) RNN: Emits the context $C$, a vector summarizing the input sequence, usually as a simple function of its final hidden state
- Decoder (writer/output) RNN: Is conditioned on the context $C$ to generate the output sequence.

### What is the innovation?

The lengths $n_x$ and $n_y$ can vary from each other
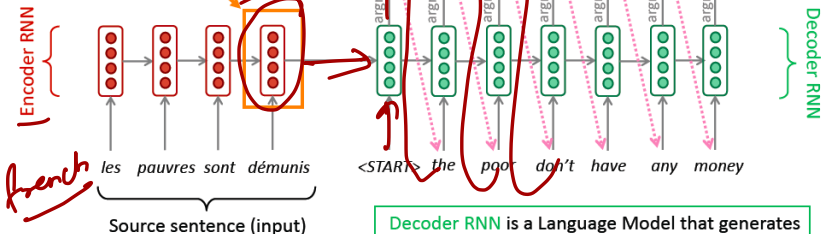
# Encoder and Decoder RNNs

- Last state of the encoder RNN is used as a representation $C$ of the input sequence, provided as input to the decoder RNN
- Decoder RNN is a vector-to-sequence RNN, described earlier.
- Two ways to receive input: as the initial state of the RNN, or can be connected to the hidden units at each time step. These two ways can also be combined.

# Sequence to Sequence Models for Machine Translation



The sequence-to-sequence model

Encoding of the source sentence. Provides initial hidden state for Decoder RNN.

Target sentence (output)

the    poor    don't    have    any    money    <END>

<START>    the    poor    don't    have    any    money

Source sentence (input)

les    pauvres    sont    démunis

Encoder RNN

Decoder RNN

Encoder RNN produces an encoding of the source sentence.

Decoder RNN is a Language Model that generates target sentence conditioned on encoding.

Note: This diagram shows test time behavior: decoder output is fed in ⋯⋯▶ as next step's input

*Many NLP tasks can be phrased as sequence-to-sequence*

- Summarization (long text $\rightarrow$ short text)
- Dialogue (previous utterances $\rightarrow$ next utterance)

*An example of a Conditional Language Model*

- **Language Model** because the decoder is predicting the next word of the target sentence $y$
- **Conditional** because the predictions are also conditioned on the sentence $x$

NMT directly calculates $P(y|x)$ :

$$P(y|x) = P(y_1|x)\, P(y_2|y_1, x)\, P(y_3|y_1, y_2, x) \ldots P(y_T|y_1, \ldots, y_{T-1}, x)$$

Probability of next target word, given target words so far and source sentence $x$

*How to train an NMT system?*

Get a big parallel corpus ...

$$J = \frac{1}{T} \sum_{t=1}^{T} J_t$$

= negative log prob of "he"

= negative log prob of "with"

= negative log prob of <END>

$$= J_1 + J_2 + J_3 + J_4 + J_5 + J_6 + J_7$$

$\hat{y}_1$ $\hat{y}_2$ $\hat{y}_3$ $\hat{y}_4$ $\hat{y}_5$ $\hat{y}_6$ $\hat{y}_7$

teacher forcing

Encoder RNN

Decoder RNN

il    a    m'    entarté    <START>    he    hit    me    with    a    pie

being trained Simultaneously

Source sentence (from corpus)

Target sentence (from corpus)

Seq2seq is optimized as a **single system**.
Backpropagation operates *"end-to-end"*.

One possibility is to generate (decode) the target sentence by taking argmax on each step of the decoder



*This is greedy decoding*

Problems with this method?