# RNN problems, LSTM, GRU
# and other variants of RNN

Harika Abburi

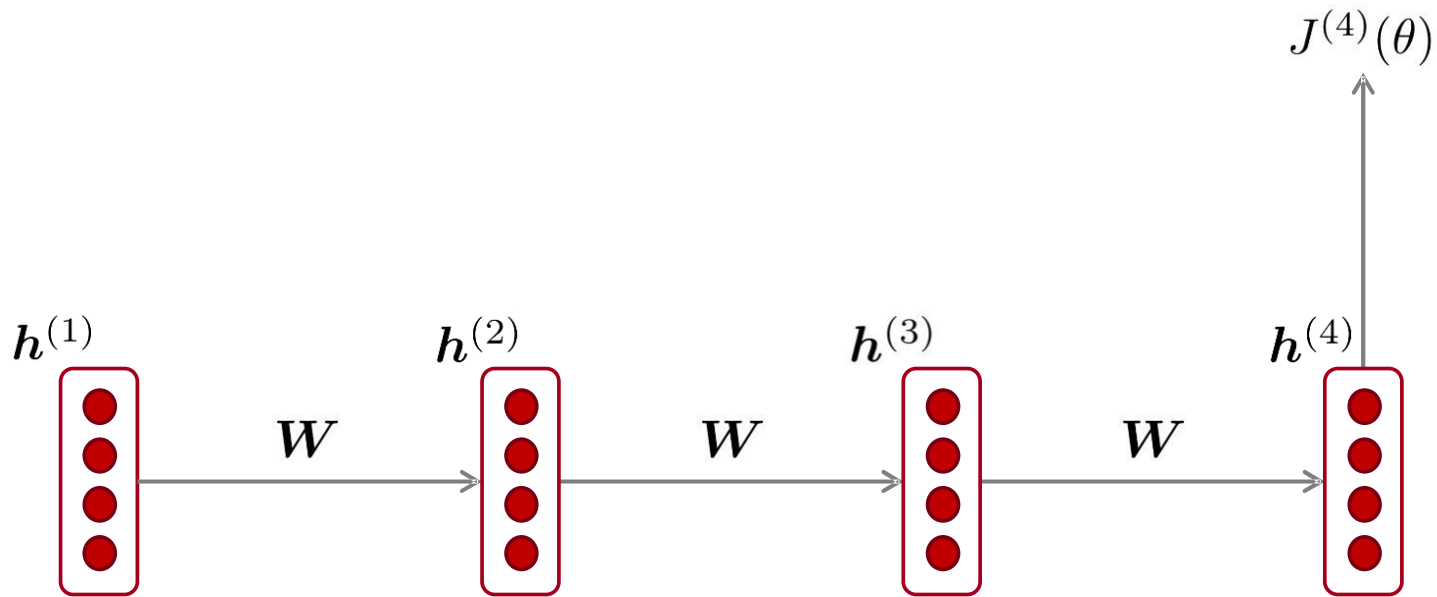Slides adapted from Stanford's NLP with Deep Learning course

# Overview

- Last class we learned:
  - Recurrent Neural Networks (RNNs) and why they're great for Language Modeling (LM).

- Today we'll learn:
  - Problems with RNNs and how to fix them
  - More complex RNN variants
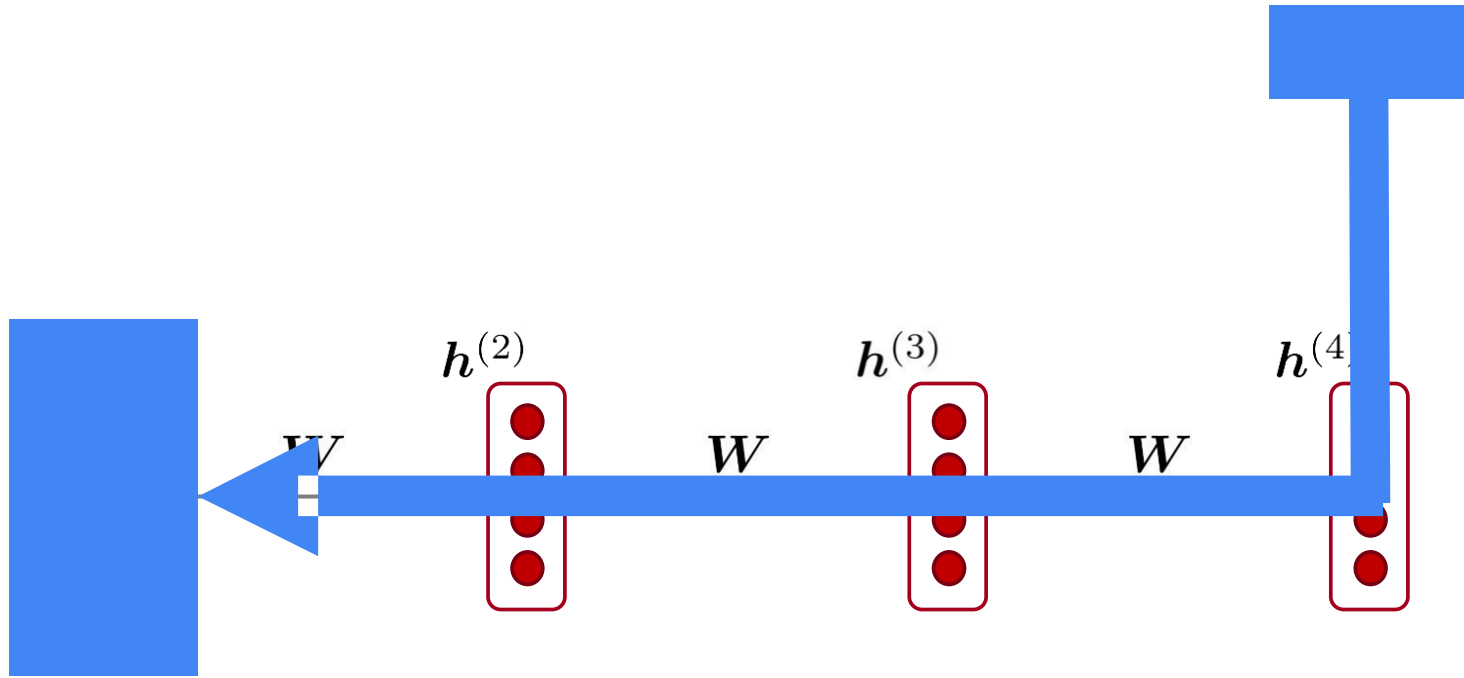
# Today's class

- **Vanishing gradient problem**

  ) **motivates**

- Two new types of RNN: LSTM and GRU

- Other fixes for vanishing (or exploding) gradient:
  - Gradient clipping
  - Skip connections

- More fancy RNN variants:
  - Bidirectional RNNs
  - Multi-layer RNNs

# Vanishing gradient intuition

$$J^{(4)}(\theta)$$

$h^{(1)}$     $h^{(2)}$     $h^{(3)}$     $h^{(4)}$

$W$     $W$     $W$

# Vanishing gradient intuition



$\boldsymbol{h}^{(2)}$     $\boldsymbol{h}^{(3)}$     $\boldsymbol{h}^{(4)}$

$W$     $W$     $W$

$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \ \text{?}$$

# Vanishing gradient intuition



$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \times \frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(2)}}$$

chain rule!

# Vanishing gradient intuition



$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \quad \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \times \qquad \frac{\partial \boldsymbol{h}^{(3)}}{\partial \boldsymbol{h}^{(2)}} \times \quad \frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(3)}}$$

chain rule!

# Vanishing gradient intuition



$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \quad \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \times \qquad\qquad \frac{\partial \boldsymbol{h}^{(3)}}{\partial \boldsymbol{h}^{(2)}} \times \qquad\qquad \frac{\partial \boldsymbol{h}^{(4)}}{\partial \boldsymbol{h}^{(3)}} \times \quad \frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(4)}}$$
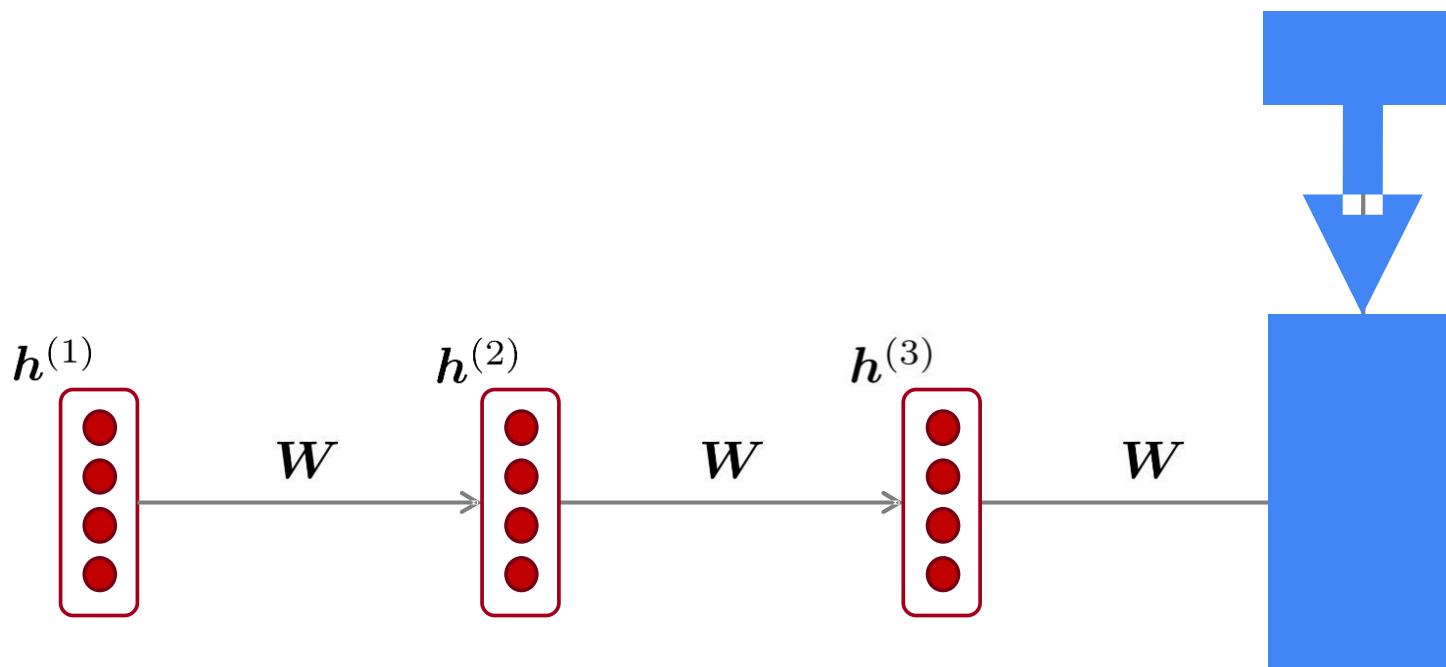
chain rule!

# Vanishing gradient intuition



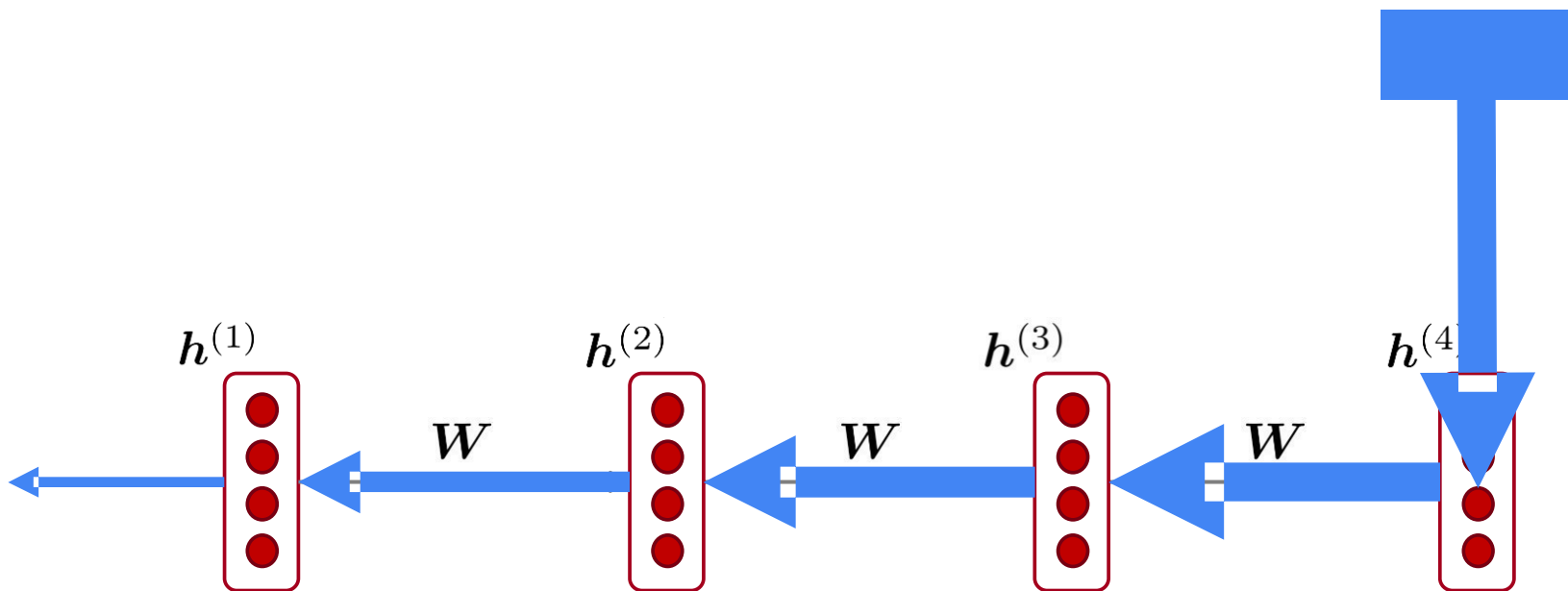$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \boxed{\frac{\partial h^{(2)}}{\partial h^{(1)}}} \times \boxed{\frac{\partial h^{(3)}}{\partial h^{(2)}}} \times \boxed{\frac{\partial h^{(4)}}{\partial h^{(3)}}} \times \frac{\partial J^{(4)}}{\partial h^{(4)}}$$

What happens if these are small?

Vanishing gradient problem: When these are small, the gradient signal gets smaller and smaller as it backpropagates further

# Vanishing gradient proof sketch

- Recall: $$\boldsymbol{h}^{(t)} = \sigma\left(\boldsymbol{W}_h \boldsymbol{h}^{(t-1)} + \boldsymbol{W}_x \boldsymbol{x}^{(t)} + \boldsymbol{b}_1\right)$$

- Therefore: $$\frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} = \operatorname{diag}\left(\sigma'\left(\boldsymbol{W}_h \boldsymbol{h}^{(t-1)} + \boldsymbol{W}_x \boldsymbol{x}^{(t)} + \boldsymbol{b}_1\right)\right) \boldsymbol{W}_h \qquad \text{(chain rule)}$$

- Consider the gradient of the loss $J^{(i)}(\theta)$ on step $i$, with respect to the hidden state $\boldsymbol{h}^{(j)}$ on some previous step $j$.

$$\frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(j)}} = \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}} \prod_{j < t \leq i} \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} \qquad \text{(chain rule)}$$

$$= \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}} \boxed{\boldsymbol{W}_h^{(i-j)}} \prod_{j < t \leq i} \operatorname{diag}\left(\sigma'\left(\boldsymbol{W}_h \boldsymbol{h}^{(t-1)} + \boldsymbol{W}_x \boldsymbol{x}^{(t)} + \boldsymbol{b}_1\right)\right) \qquad \left(\text{value of } \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}}\right)$$

If $W_h$ is small, then this term gets vanishingly small as $i$ and $j$ get further apart
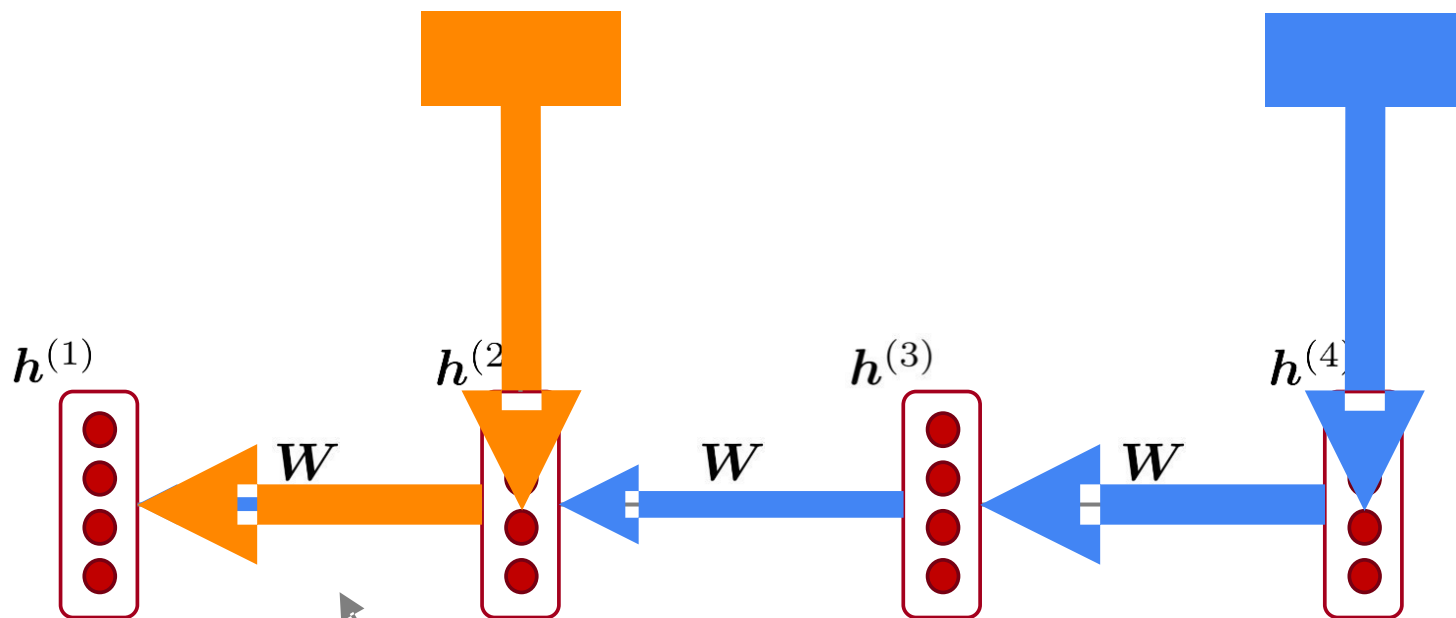
# Vanishing gradient proof sketch

- Consider matrix L2 norms:

$$\left\|\frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(j)}}\right\| \leq \left\|\frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}}\right\| \|\boldsymbol{W}_h\|^{(i-j)} \prod_{j < t \leq i} \left\|\text{diag}\left(\sigma'\left(\boldsymbol{W}_h \boldsymbol{h}^{(t-1)} + \boldsymbol{W}_x \boldsymbol{x}^{(t)} + \boldsymbol{b}_1\right)\right)\right\|$$

- Pascanu et al showed that that if the largest eigenvalue of $W_h$ is less than 1, then the gradient $\left\|\frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(j)}}\right\|$ will shrink exponentially
  - Here the bound is 1 because we have sigmoid nonlinearity

- There's a similar proof relating a largest eigenvalue >1 to exploding gradients

# Why is vanishing gradient a problem?



$h^{(1)}$   $h^{(2)}$   $h^{(3)}$   $h^{(4)}$

$W$   $W$   $W$

Gradient signal from faraway is lost because it's much smaller than gradient signal from close-by.

So model weights are only updated only with respect to near effects, not long-term effects.
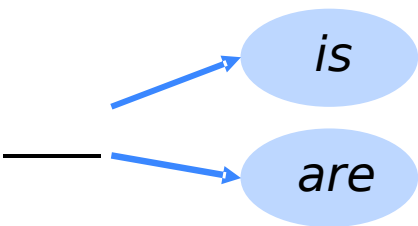
# Why is vanishing gradient a problem?

- Another explanation: Gradient can be viewed as a measure of the *effect of the past on the future*

- If the gradient becomes vanishingly small over longer distances (step $t$ to step $t+n$), then we can't tell whether:
  1. There's no dependency between step $t$ and $t+n$ in the data
  2. We have wrong parameters to capture the true dependency between $t$ and $t+n$

# Effect of vanishing gradient on RNN-LM

- **LM task:** *When she tried to print her tickets, she found that the printer was out of toner. She went to the stationery store to buy more toner. It was very overpriced. After installing the toner into the printer, she finally printed her* _____

- To learn from this training example, the RNN-LM needs to model the dependency between *"tickets"* on the 7 th step and the target word *"tickets"* at the end.

- But if gradient is small, the model can't learn this dependency
  - So the model is unable to predict similar long-distance dependencies at test time

# Effect of vanishing gradient on RNN-LM

- **LM task:** *The writer of the books ___*

  *is*

  *are*

- **Correct answer**: *The writer of the books is planning a sequel*

- **Syntactic** **recency:** *The writer of the books is*     (correct)

- **Sequential** **recency:** *The writer of the books are*     (incorrect)

- Due to vanishing gradient, RNN-LMs are better at learning from sequential recency than syntactic recency, so they make this type of error more often than we'd like [Linzen et al 2016]

"Assessing the Ability of LSTMs to Learn Syntax-Sensitive Dependencies", Linzen et al, 2016. https://arxiv.org/pdf/1611.01368.pdf

# Why is <u>exploding</u> gradient a problem?

- If the gradient becomes too big, then the SGD update step becomes too big:

$$\theta^{new} = \theta^{old} - \overbrace{\alpha}^{\text{learning rate}} \underbrace{\nabla_\theta J(\theta)}_{\text{gradient}}$$

- This can cause bad updates: we take too large a step and reach a bad parameter configuration (with large loss)

- In the worst case, this will result in Inf or NaN in your network (then you have to restart training from an earlier checkpoint)

# Gradient clipping: solution for exploding gradient

- Gradient clipping: if the norm of the gradient is greater than some threshold, scale it down before applying SGD update
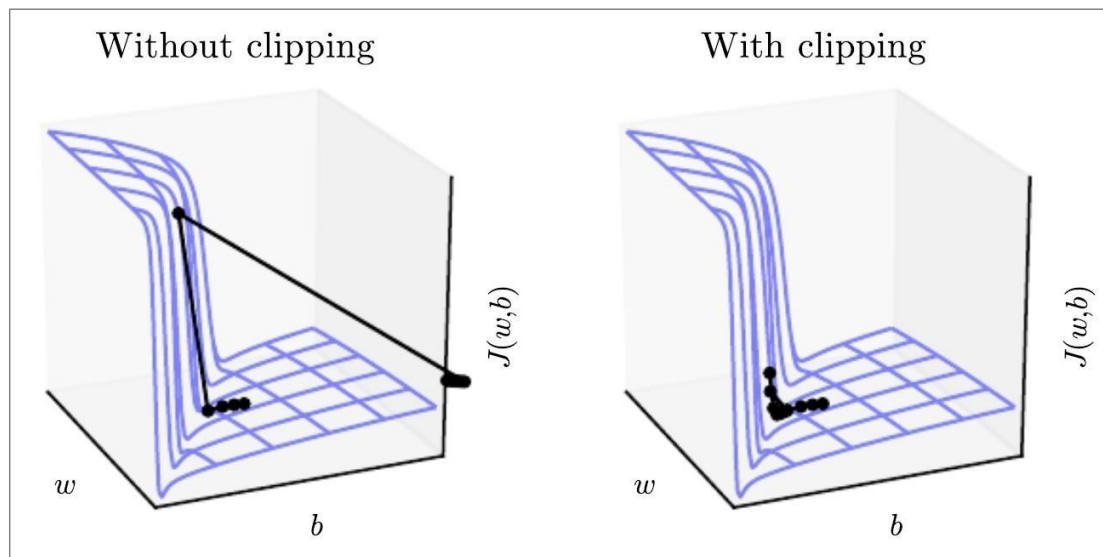
**Algorithm 1** Pseudo-code for norm clipping

$$\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$$
**if** $\|\hat{\mathbf{g}}\| \geq threshold$ **then**
$$\hat{\mathbf{g}} \leftarrow \frac{threshold}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$$
**end if**

- Intuition: take a step in the same direction, but a smaller step

**Source**: "On the difficulty of training recurrent neural networks", Pascanu et al, 2013. http://proceedings.mlr.press/v28/pascanu13.pdf

# Gradient clipping: solution for exploding gradient



- This shows the loss surface of a simple RNN (hidden state is a scalar not a vector)

- The "cliff" is dangerous because it has steep gradient

- On the left, gradient descent takes two very big steps due to steep gradient, resulting in climbing the cliff then shooting off to the right (both bad updates)

- On the right, gradient clipping reduces the size of those steps, so effect is less drastic

# How to fix vanishing gradient problem?

- The main problem is that *it's too difficult for the RNN to learn to preserve information over many timesteps*.

- In a vanilla RNN, the hidden state is constantly being rewritten

$$h^{(t)} = \sigma \left( W_h h^{(t-1)} + W_x x^{(t)} + b \right)$$

- How about a RNN with separate memory?

# Long Short-Term Memory (LSTM)

- A type of RNN proposed by Hochreiter and Schmidhuber in 1997 as asolution to the vanishing gradients problem.

- On step $t$, there is a hidden state $\boldsymbol{h}^{(t)}$ and a cell state $\boldsymbol{c}^{(t)}$
  - Both are vectors length $n$
  - The cell stores long-term information
  - The LSTM can erase, write and read information from the cell

- The selection of which information is erased/written/read is controlled by three corresponding gates
  - The gates are also vectors length $n$
  - On each timestep, each element of the gates can be open (1), closed (0), or somewhere in-between.
  - The gates are dynamic: their value is computed based on the current context

"Long short-term memory", Hochreiter and Schmidhuber, 1997. https://www.bioinf.jku.at/publications/older/2604.pdf

# Long Short-Term Memory (LSTM)

We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep $t$:

**Forget gate:** controls what is kept vs forgotten, from previous cell state

**Input gate:** controls what parts of the new cell content are written to cell

**Output gate:** controls what parts of cell are output to hidden state

**Sigmoid function**: all gate values are between 0 and 1

$$f^{(t)} = \sigma\left(W_f h^{(t-1)} + U_f x^{(t)} + b_f\right)$$

$$i^{(t)} = \sigma\left(W_i h^{(t-1)} + U_i x^{(t)} + b_i\right)$$

$$o^{(t)} = \sigma\left(W_o h^{(t-1)} + U_o x^{(t)} + b_o\right)$$

**New cell content:** this is the new content to be written to the cell

**Cell state**: erase ("forget") some content from last cell state, and write ("input") some new cell content

**Hidden state**: read ("output") some content from the cell

$$\tilde{c}^{(t)} = \tanh\left(W_c h^{(t-1)} + U_c x^{(t)} + b_c\right)$$

$$c^{(t)} = f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)}$$

$$h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$

All these are vectors of same length $n$

Gates are applied using element-wise product

# Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:

# Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:

# How does LSTM solve vanishing gradients?

- The LSTM architecture makes it easier for the RNN to preserve information over many timesteps

    - e.g. if the forget gate is set to remember everything on every timestep, then the info in the cell is preserved indefinitely

    - By contrast, it's harder for vanilla RNN to learn a recurrent weight matrix $W_h$ that preserves info in hidden state

- LSTM doesn't *guarantee* that there is no vanishing/exploding gradient, but it does provide an easier way for the model to learn long-distance dependencies

# LSTMs: real-world success

- In 2013-2015, LSTMs started achieving state-of-the-art results
  - Successful tasks include: handwriting recognition, speech recognition, machine translation, parsing, image captioning
  - LSTM became the dominant approach

- Now (2019), other approaches (e.g. Transformers) have become more dominant for certain tasks.
  - For example in **WMT** (a MT conference + competition):
  - In WMT 2016, the summary report contains "RNN" 44 times
  - In WMT 2018, the report contains "RNN" 9 times and "Transformer" 63 times

**Source:** "Findings of the 2016 Conference on Machine Translation (WMT16)", Bojar et al. 2016, http://www.statmt.org/wmt16/pdf/W16-2301.pdf

**Source:** "Findings of the 2018 Conference on Machine Translation (WMT18)", Bojar et al. 2018, http://www.statmt.org/wmt18/pdf/WMT028.pdf

# Gated Recurrent Units (GRU)

- Proposed by Cho et al. in 2014 as a simpler alternative to the LSTM.
- On each timestep $t$ we have input $\boldsymbol{x}^{(t)}$ and hidden state $\boldsymbol{h}^{(t)}$ (no cell state).

**Update gate:** controls what parts of hidden state are updated vs preserved

**Reset gate:** controls what parts of previous hidden state are used to compute new content

$$\boldsymbol{u}^{(t)} = \sigma\left(\boldsymbol{W}_u \boldsymbol{h}^{(t-1)} + \boldsymbol{U}_u \boldsymbol{x}^{(t)} + \boldsymbol{b}_u\right)$$

$$\boldsymbol{r}^{(t)} = \sigma\left(\boldsymbol{W}_r \boldsymbol{h}^{(t-1)} + \boldsymbol{U}_r \boldsymbol{x}^{(t)} + \boldsymbol{b}_r\right)$$

**New hidden state content:** reset gate selects useful parts of prev hidden state. Use this and current input to compute new hidden content.

$$\tilde{\boldsymbol{h}}^{(t)} = \tanh\left(\boldsymbol{W}_h(\boldsymbol{r}^{(t)} \circ \boldsymbol{h}^{(t-1)}) + \boldsymbol{U}_h \boldsymbol{x}^{(t)} + \boldsymbol{b}_h\right)$$

$$\boldsymbol{h}^{(t)} = (1 - \boldsymbol{u}^{(t)}) \circ \boldsymbol{h}^{(t-1)} + \boldsymbol{u}^{(t)} \circ \tilde{\boldsymbol{h}}^{(t)}$$

**Hidden state:** update gate simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content

**How does this solve vanishing gradient?**
Like LSTM, GRU makes it easier to retain info long-term (e.g. by setting update gate to 0)

"Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation", Cho et al. 2014, https://arxiv.org/pdf/1406.1078v3.pdf

# LSTM vs GRU

- Researchers have proposed many gated RNN variants, but LSTM and GRU are the most widely-used

- The biggest difference is that GRU is quicker to compute and has fewer parameters

- There is no conclusive evidence that one consistently performs better than the other

- LSTM is a good default choice (especially if your data has particularly long dependencies, or you have lots of training data)

- Rule of thumb: start with LSTM, but switch to GRU if you want something more efficient

# Is vanishing/exploding gradient just a RNN problem?

- No! It can be a problem for all neural architectures (including feed-forward and convolutional), especially deep ones.
    - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
    - Thus lower layers are learnt very slowly (hard to train)
    - Solution: lots of new deep feedforward/convolutional architectures that add more direct connections (thus allowing the gradient to flow)

For example:
- Residual connections aka "ResNet"
- Also known as skip-connections
- The identity connection preserves information by default
- This makes deep networks much easier to train



Figure 2. Residual learning: a building block.

"Deep Residual Learning for Image Recognition", He et al, 2015. https://arxiv.org/pdf/1512.03385.pdf
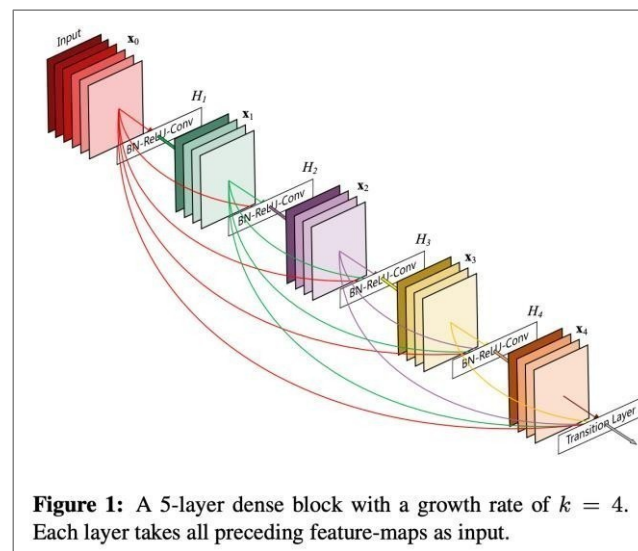
# Is vanishing/exploding gradient just a RNN problem?

- No! It can be a problem for all neural architectures (including feed-forward and convolutional), especially deep ones.
  - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
  - Thus lower layers are learnt very slowly (hard to train)
  - Solution: lots of new deep feedforward/convolutional architectures that add more direct connections (thus allowing the gradient to flow)

For example:
- Dense connections aka "DenseNet"
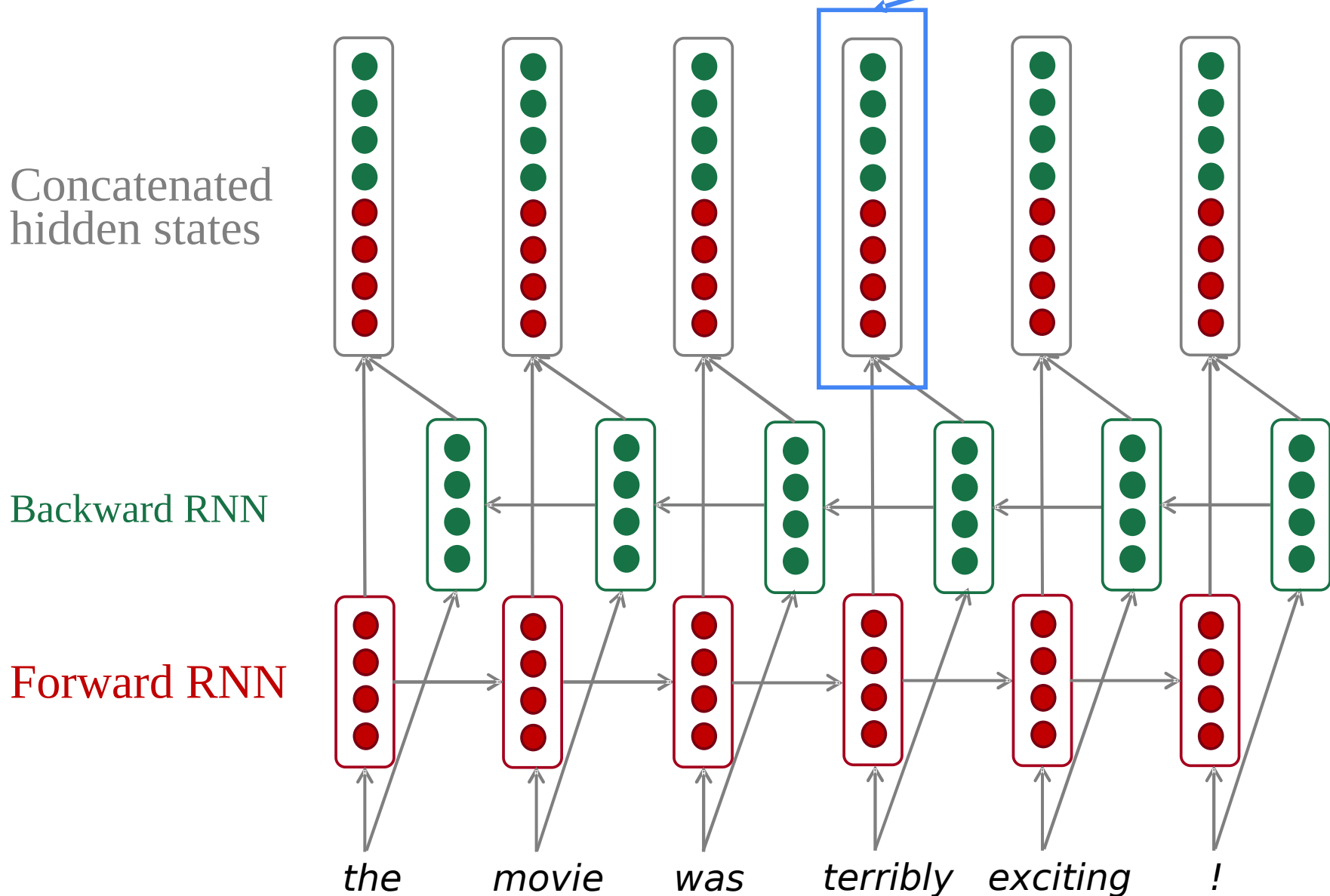- Directly connect everything to everything!



**Figure 1:** A 5-layer dense block with a growth rate of $k = 4$. Each layer takes all preceding feature-maps as input.

"Densely Connected Convolutional Networks", Huang et al, 2017. https://arxiv.org/pdf/1608.06993.pdf

# Bidirectional RNNs: motivation

Task: Sentiment Classification

**positive**

Sentence encoding

element-wise mean/max

element-wise mean/max

*the*   *movie*   *was*   *terribly*   *exciting*   *!*

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)

# **Bidirectional RNNs**

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

Concatenated hidden states

Backward RNN

Forward RNN

*the*        *movie*        *was*        *terribly*    *exciting*        *!*

# Bidirectional RNNs: simplified diagram



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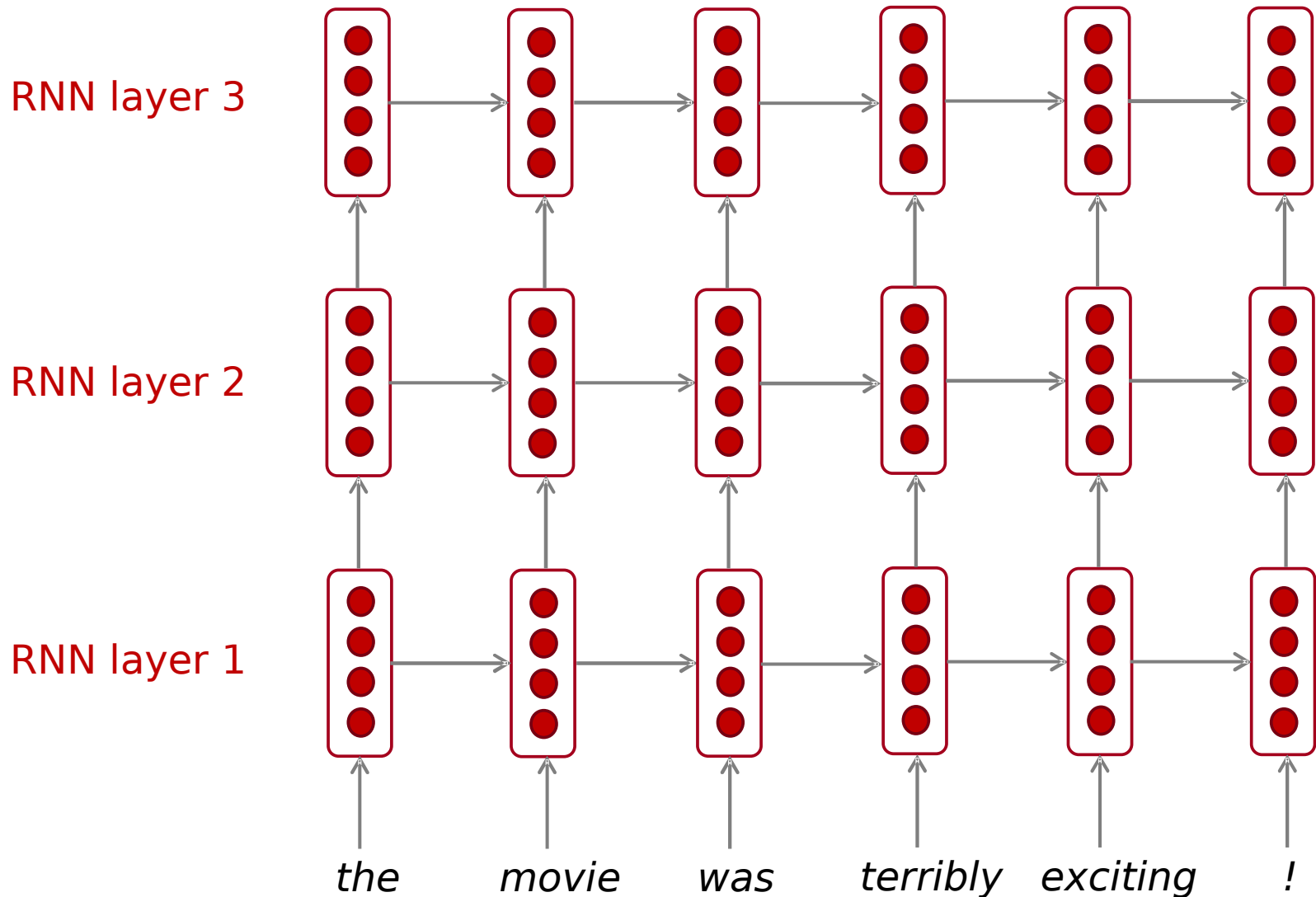
# Multi-layer RNNs

- RNNs are already "deep" on one dimension (they unroll over many timesteps)

- We can also make them "deep" in another dimension by applying multiple RNNs – this is a multi-layer RNN.

- This allows the network to compute more complex representations
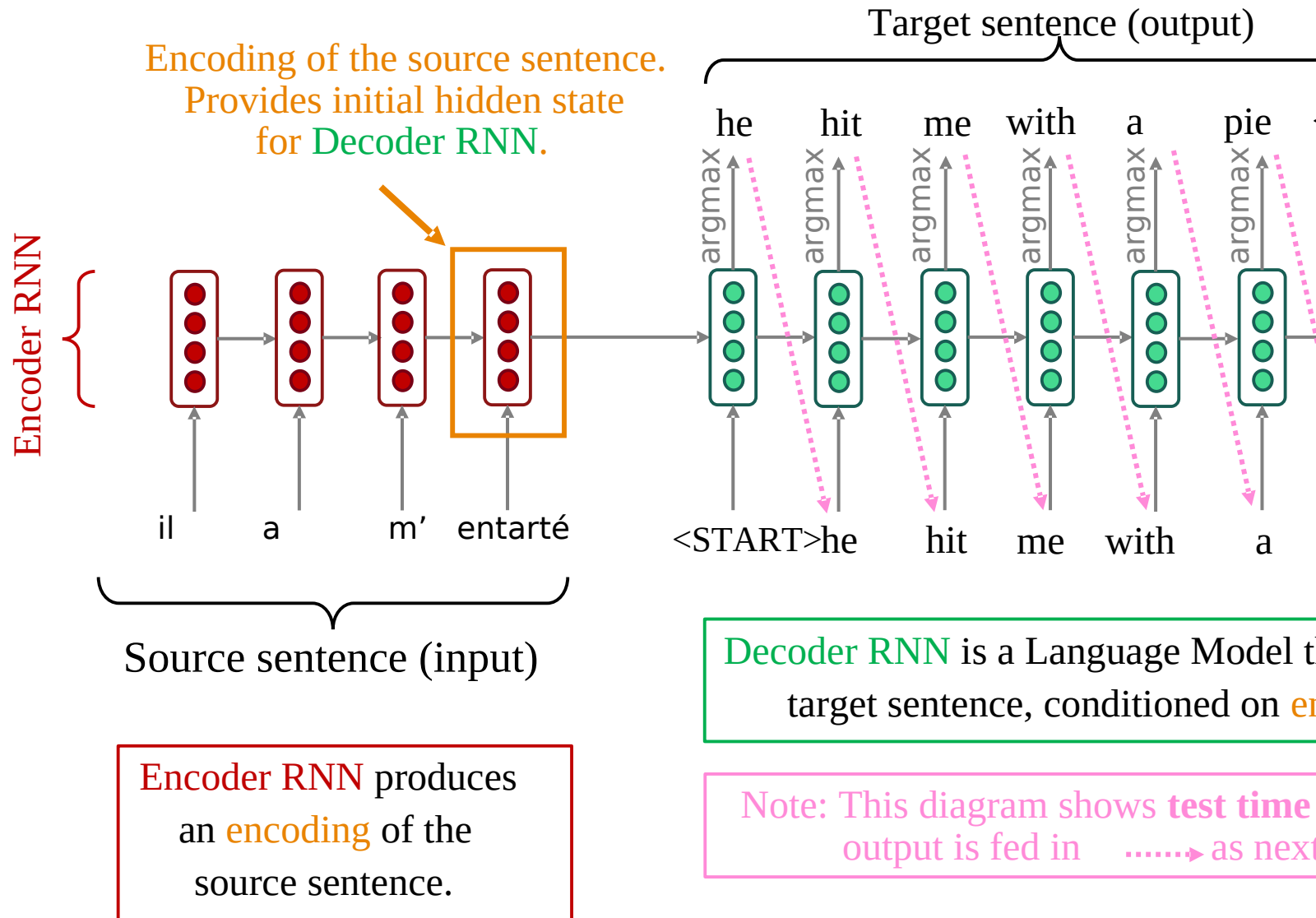
- Multi-layer RNNs are also called *stacked RNNs*.

# **Multi-layer RNNs**

RNN layer 3

RNN layer 2

RNN layer 1

the    movie    was    terribly    exciting    !

# Sequence-to-sequence model

## Neural Machine Translation (NMT)

Target sentence (output)

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

Encoder RNN

he    hit    me    with    a    pie

argmax  argmax  argmax  argmax  argmax  argmax

il    a    m'    entarté

<START>he    hit    me    with    a

Source sentence (input)

Encoder RNN produces
an encoding of the
source sentence.

Decoder RNN is a Language Model t
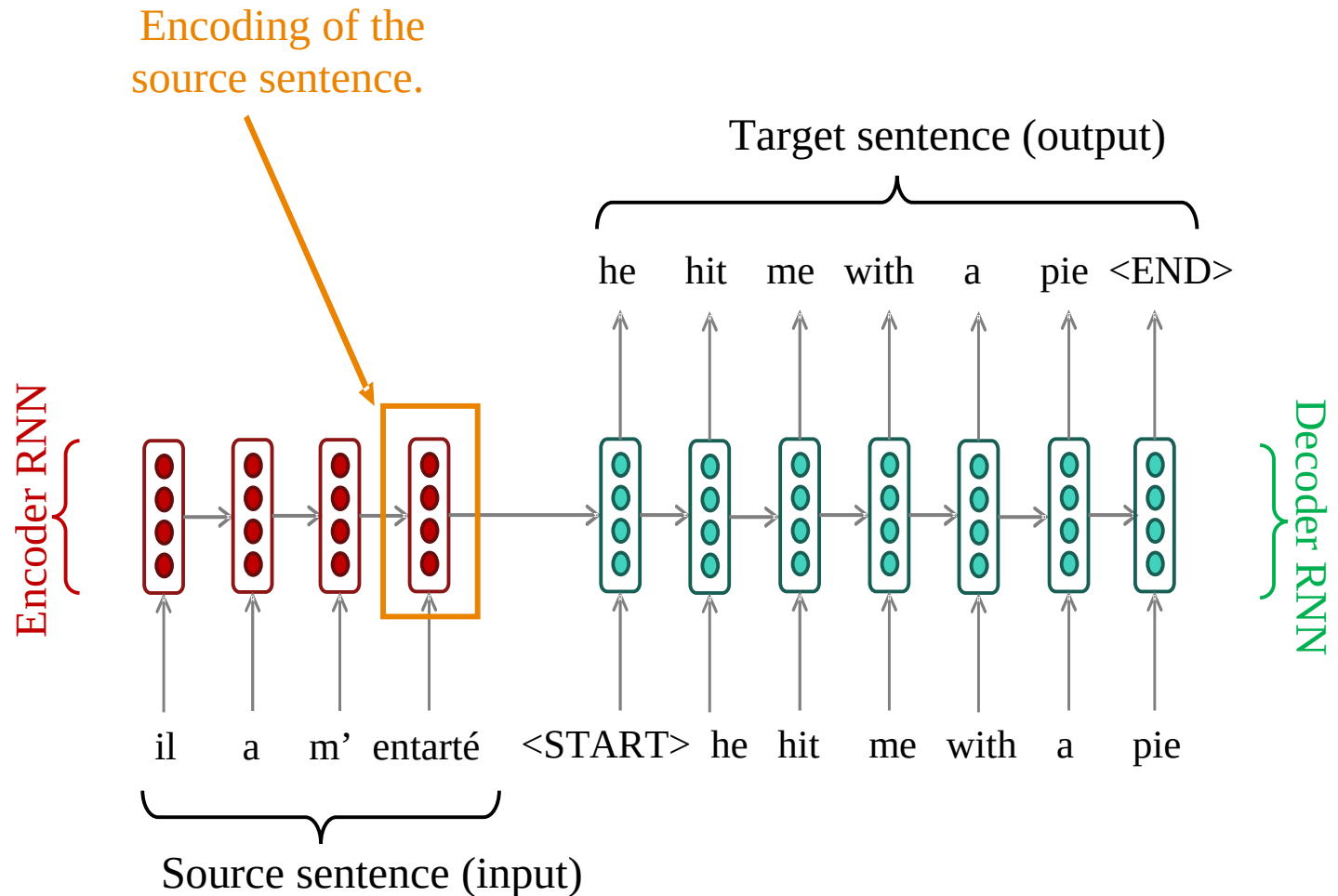target sentence, conditioned on e

Note: This diagram shows **test time**
output is fed in ......▶ as next
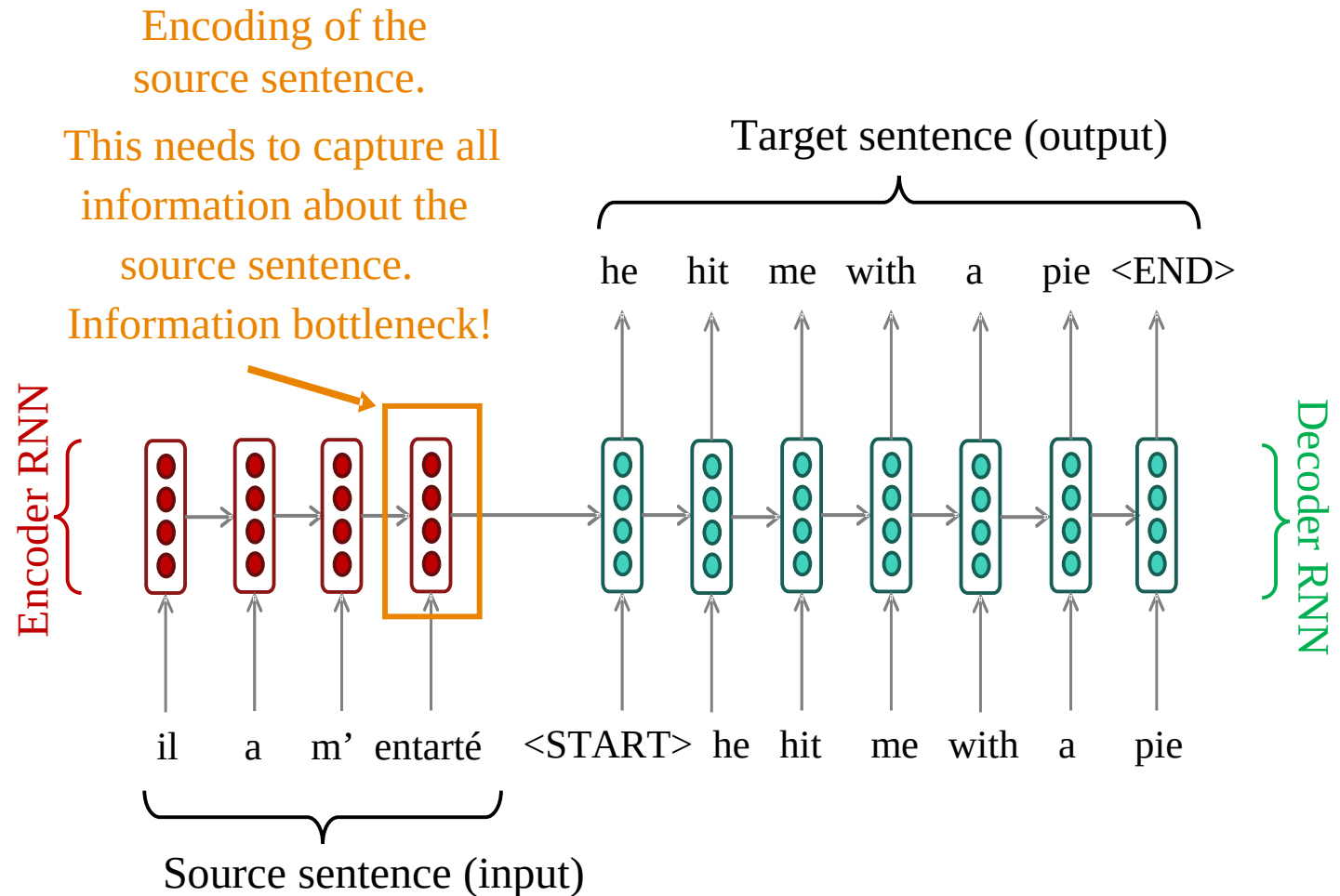
# Sequence-to-sequence is versatile!

- Sequence-to-sequence is useful for more than just MT

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

  - Summarization (long text → short text)

  - Dialogue (previous utterances → next utterance)

  - Parsing (input text → output parse as sequence)

  - Code generation (natural language → Python code)

# Sequence-to-sequence: the bottleneck problem

Encoding of the
source sentence.

Target sentence (output)

he  hit  me  with  a  pie  <END>

Encoder RNN

Decoder RNN

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

Source sentence (input)

**Problems with this architecture?**

# Sequence-to-sequence: the bottleneck problem



Encoding of the source sentence.

This needs to capture all information about the source sentence. Information bottleneck!

Target sentence (output)

he hit me with a pie <END>

Encoder RNN

Decoder RNN

il a m' entarté <START> he hit me with a pie
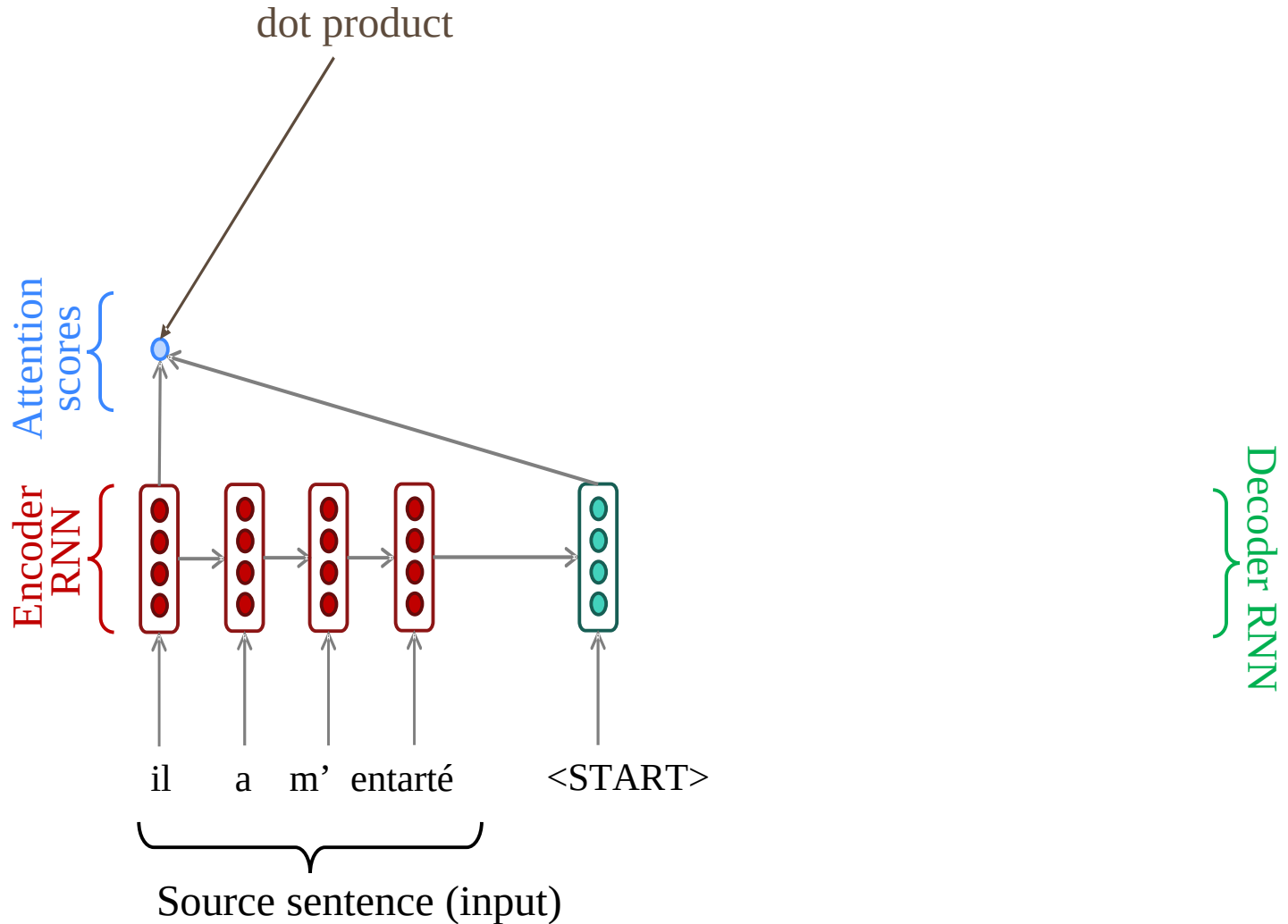
Source sentence (input)

**Problems with this architecture?**
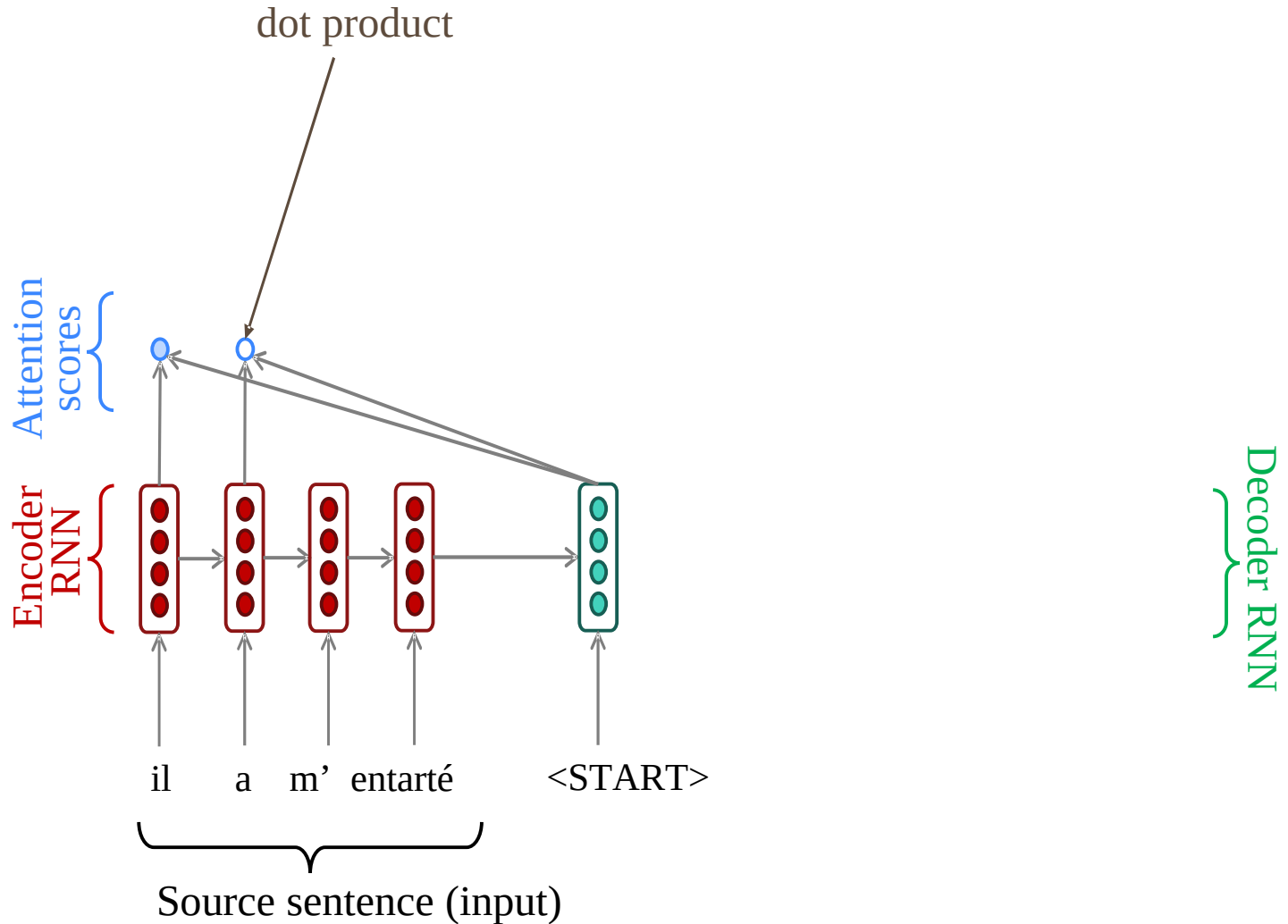
# Attention

- **Attention** provides a solution to the bottleneck problem.

- <u>Core id</u>ea: on each step of the decoder, use direct connection to the encoder to focus on a particular part of the source sequence.

- First, we will show via diagram (no equations), then we will show with equations.

# Sequence-to-sequence with attention

dot product

Attention scores

Encoder RNN

Decoder RNN

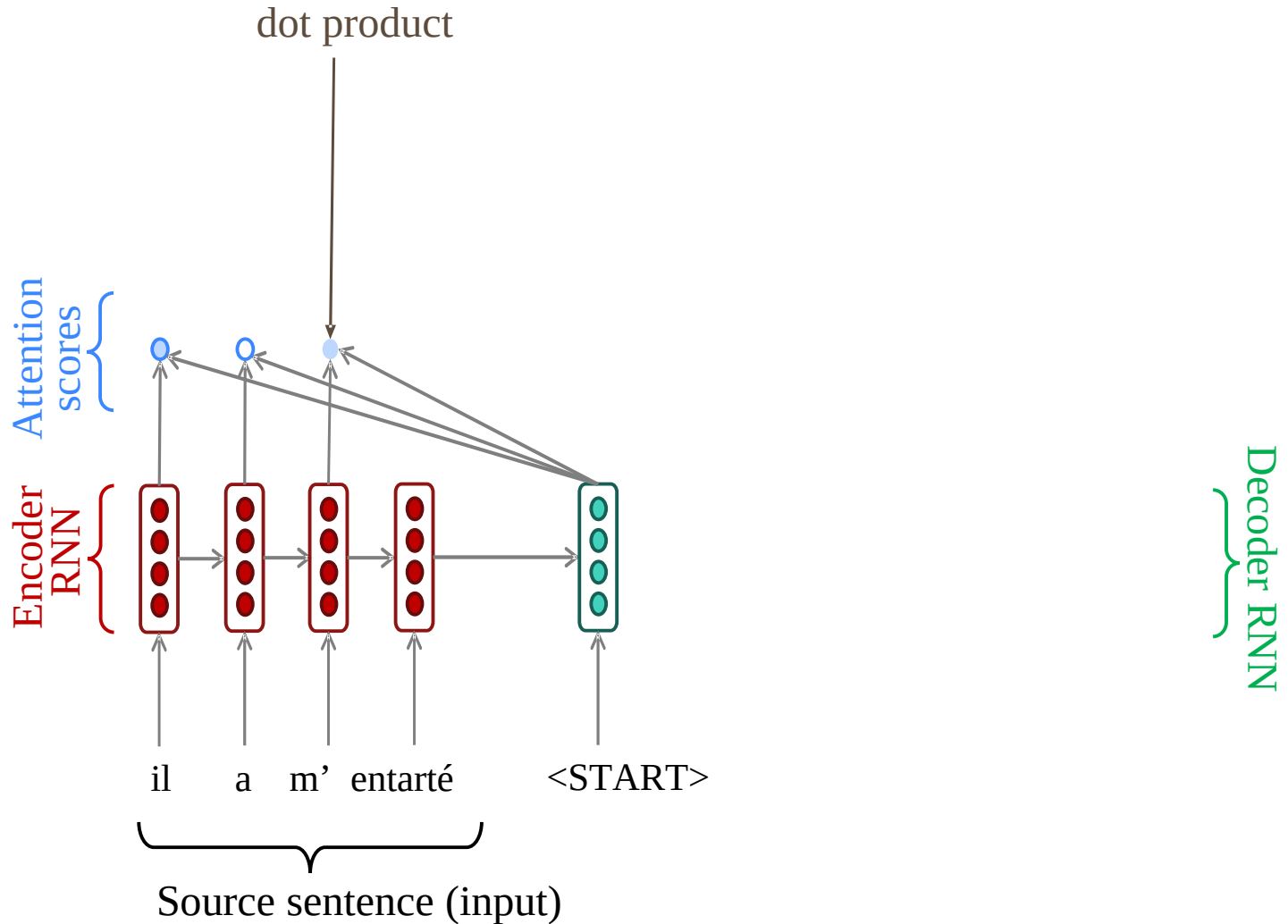il    a    m'   entarté     \<START\>
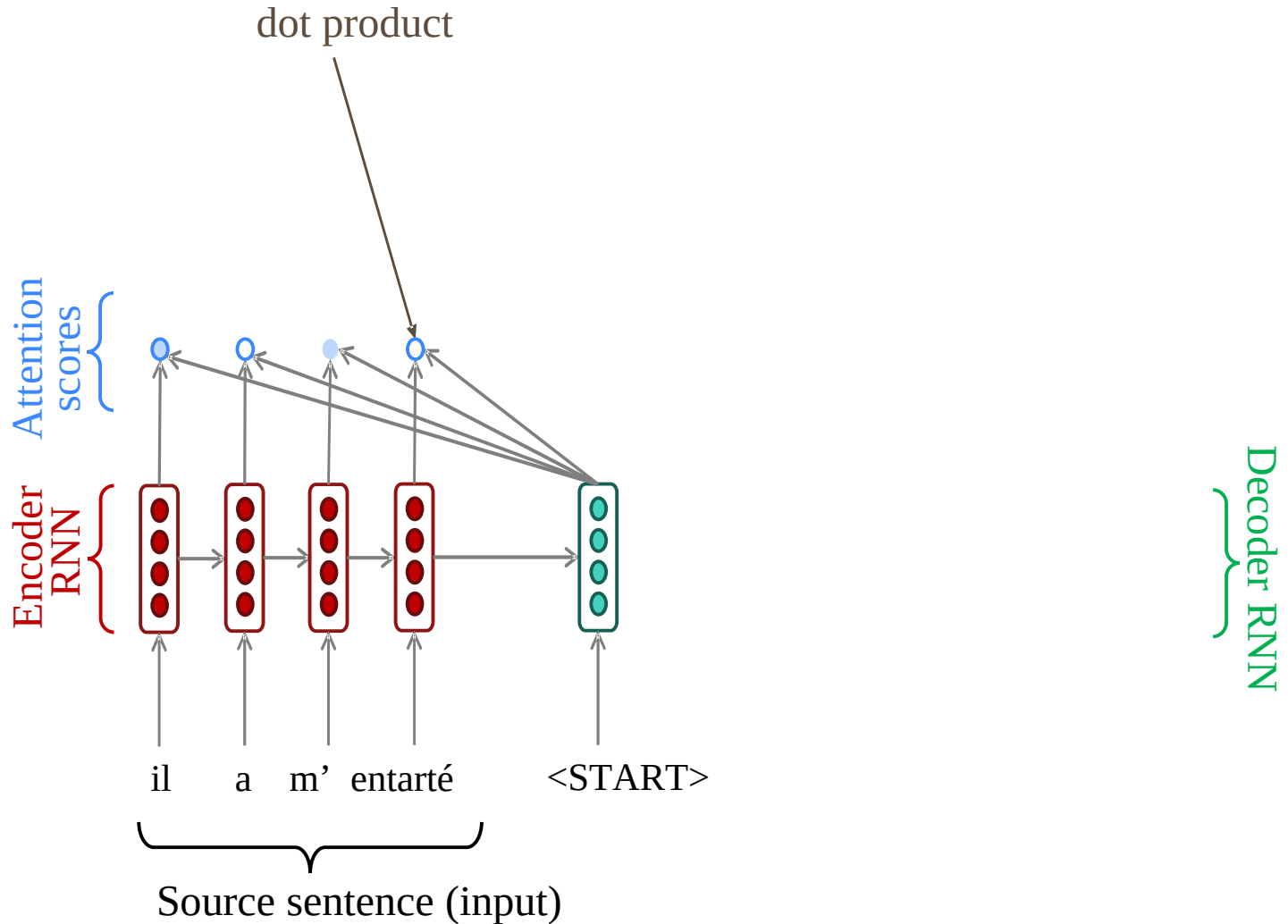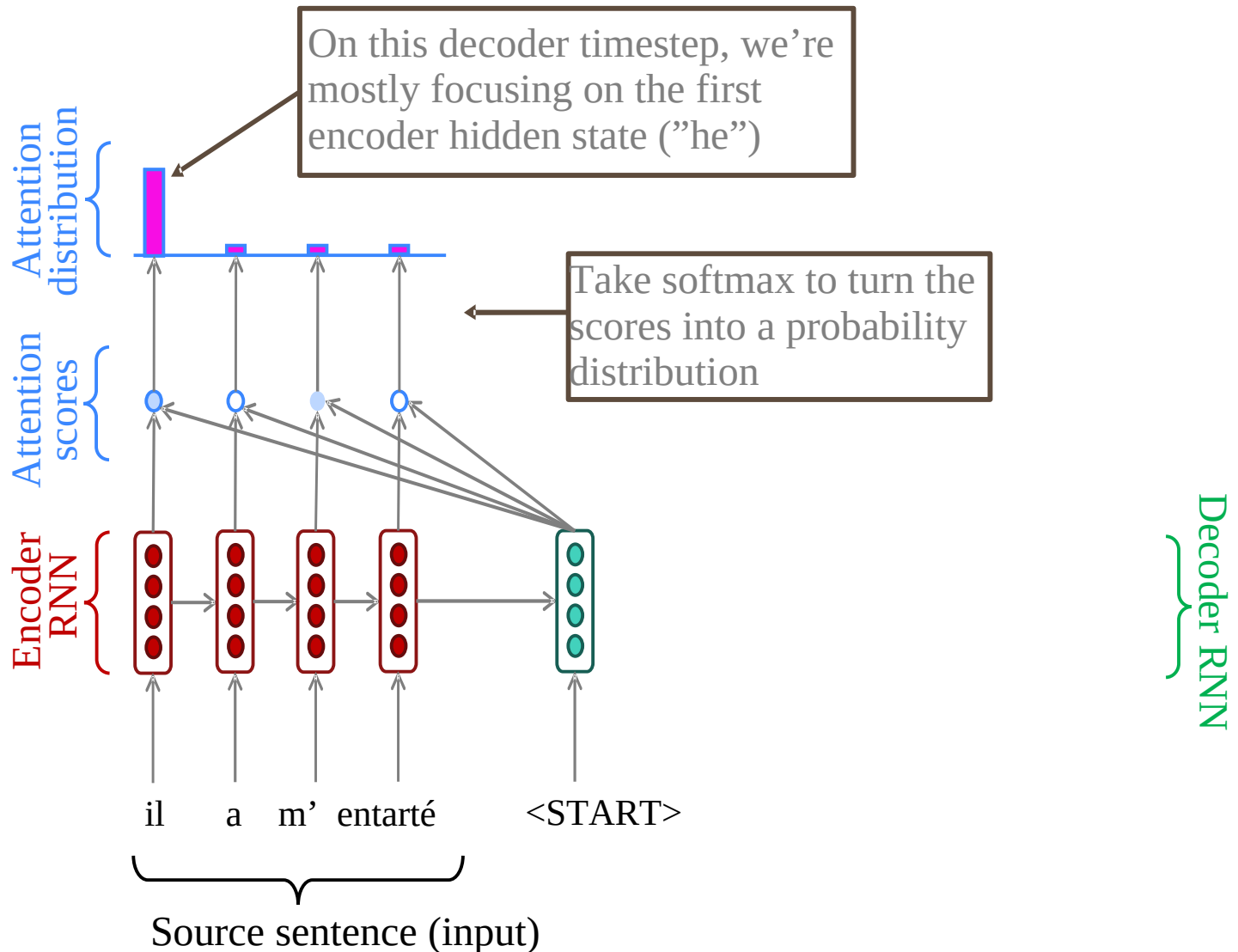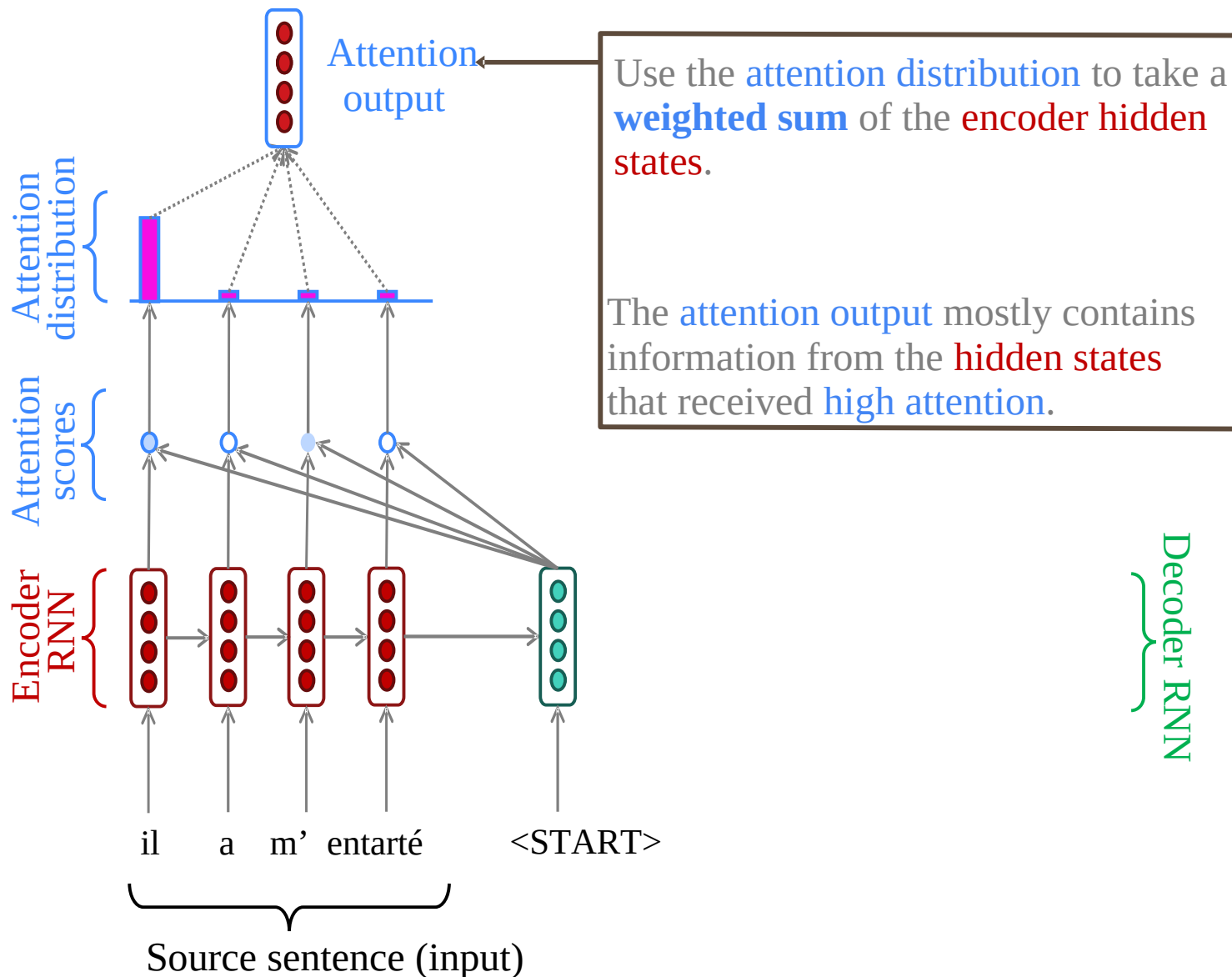
Source sentence (input)

# Sequence-to-sequence with attention

# Sequence-to-sequence with attention

# Sequence-to-sequence with attention



dot product

Attention scores

Encoder RNN

Decoder RNN

il    a    m'  entarté    <START>

Source sentence (input)

# Sequence-to-sequence with attention



On this decoder timestep, we're mostly focusing on the first encoder hidden state ("he")

Take softmax to turn the scores into a probability distribution

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

il    a    m'   entarté    <START>

Source sentence (input)

# Sequence-to-sequence with attention



Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

Use the attention distribution to take a **weighted sum** of the encoder hidden states.

The attention output mostly contains information from the hidden states that received high attention.

il    a    m'   entarté    <START>

Source sentence (input)

# Sequence-to-sequence with attention



**Attention output**

he

$y_!!$

Concatenate attention output with decoder hidden state, then use to compute

**Attention distribution**

**Attention scores**

**Encoder RNN**

**Decoder RNN**

il    a    m'   entarté    <START>

Source sentence (input)

# Sequence-to-sequence with attention



Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

hit

$y"$

Sometimes we take the attention output from the previous step, and also feed it into the decoder (along with the usual decoder input).

il     a     m'     entarté     <START> he

Source sentence (input)

# Sequence-to-sequence with attention

# Sequence-to-sequence with attention



Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

pie

$\hat{y}_6$

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

Source sentence (input)

74

# Attention: in equations

- We have encoder hidden states $h_1, \ldots, h_N \in \mathbb{R}^h$

- On timestep t, we have decoder hidden $s_t \in \mathbb{R}^h$

- We get the attention scores $e^t = [s_t^T h_1, \ldots, s_t^T h_N] \in \mathbb{R}^N$

- We take softmax to get the attention distribution $\alpha^t = \mathrm{softmax}(e^t) \in \mathbb{R}^N$ for this step (this is a probability distribution and sums to 1)

- We use $\alpha^t = \mathrm{softmax}(e^t) \in \mathbb{R}^N$ to take a weighted sum of the encoder hidden states to
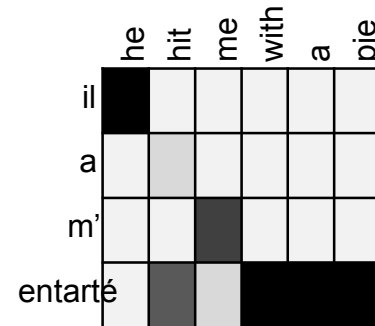
  get the attention output

$$a_t = \sum_{i=1}^N \alpha_i^t h_i \in \mathbb{R}^h$$

- Finally we concatenate the attention output $a_t = \sum_{i=1}^N \alpha_i^t h_i \in \mathbb{R}^h$ with the decoder hidden

  state $s_t \in \mathbb{R}^h$ and proceed as in the non-attention seq2seq model

$$[a_t; s_t] \in \mathbb{R}^{2h}$$

# Attention is great

- Attention significantly improves NMT performance

  - It's very useful to allow decoder to focus on certain parts of the source
- Attention solves the bottleneck problem

  - Attention allows decoder to look directly at source; bypass bottleneck
- Attention helps with vanishing gradient problem

  - Provides shortcut to far away states
- Attention provides some interpretability

  - By inspecting attention distribution, we can see what the decoder was focusing on

# Attention is a general Deep Learning technique

**More general definition of attention**:

Given a set of vector values, and a vector query, **attention** is a technique to compute a weighted sum of the values, dependent on the query.

**Intuition**:

- The weighted sum is a selective summary of the information contained in the values, where the query determines which values to focus on.

- Attention is a way to obtain a fixed-size representation of an arbitrary set of representations (the values), dependent on some other representation (the query).

# There are several attention variants

- We have some values $\boldsymbol{h}_1, \ldots, \boldsymbol{h}_N \in \mathbb{R}^{d_1}$ and a query $\boldsymbol{s} \in \mathbb{R}^{d_2}$

- Attention always involves:

  1. Computing the attention scores $\boldsymbol{e} \in \mathbb{R}^N$ ← There are multiple ways to do this

  2. Taking softmax to get attention distribution $\alpha$:

$$\alpha = \text{softmax}(\boldsymbol{e}) \in \mathbb{R}^N$$

  3. Using attention distribution to take weighted sum of values:

$$\boldsymbol{a} = \sum_{i=1}^{N} \alpha_i \boldsymbol{h}_i \in \mathbb{R}^{d_1}$$

   thus obtaining the attention output **a** (sometimes called the context vector)

# Attention variants

There are several ways you can compute $e \in \mathbb{R}^N$ from $h_1, \ldots, h_N \in \mathbb{R}^{d_1}$ and $s \in \mathbb{R}^{d_2}$
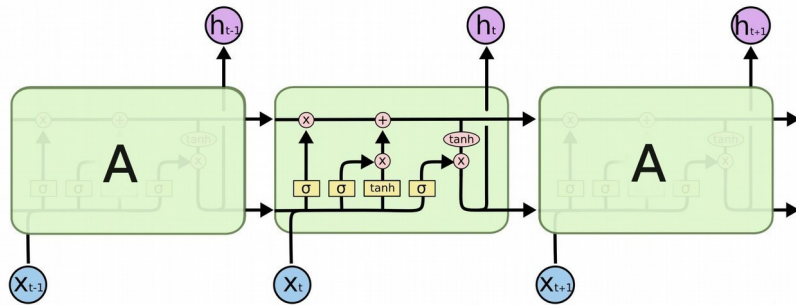
- Basic dot-product attention: $e_i = s^T h_i \in \mathbb{R}$
  - Note: this assumes $d_1 = d_2$
  - This is the version we saw earlier

- Multiplicative attention: $e_i = s^T W h_i \in \mathbb{R}$
  - Where $W \in \mathbb{R}^{d_2 \times d_1}$ is a weight matrix

- Additive attention: $e_i = v^T \tanh(W_1 h_i + W_2 s) \in \mathbb{R}$
  - Where $W_1 \in \mathbb{R}^{d_3 \times d_1}, W_2 \in \mathbb{R}^{d_3 \times d_2}$ are weight matrices and $v \in \mathbb{R}^{d_3}$ is a weight vector.

  - $d_3$ (the attention dimensionality) is a hyperparameter

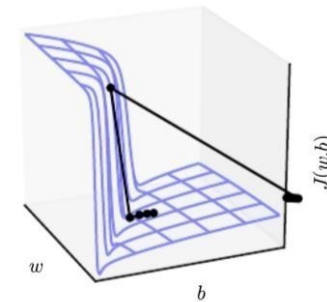**More information:** "Deep Learning for NLP Best Practices", Ruder, 2017. http://ruder.io/deep-learning-nlp-best-practices/index.html#attention

"Massive Exploration of Neural Machine Translation Architectures", Britz et al, 2017, https://arxiv.org/pdf/1703.03906.pdf
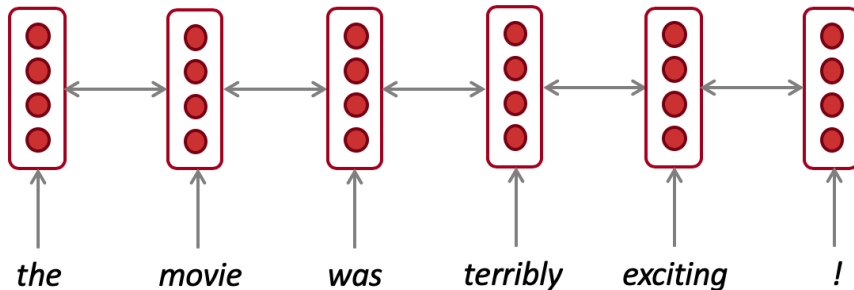
# In summary

Lots of new information today! What are the practical takeaways?
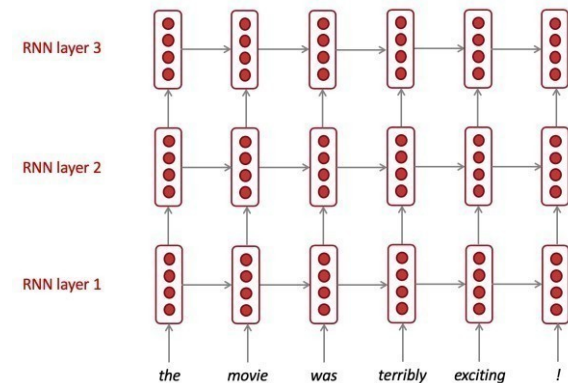


1. LSTMs are powerful but GRUs are faster



2. Clip your gradients



3. Use bidirectionality when possible



4. Multi-layer RNNs are powerful, but you might need skip/dense-connections if it's deep

5. Attention is a way to focus on particular parts of the input

- Improves sequence-to-sequence a lot!

# Thanks!