

CMPT 983

Grounded Natural Language Understanding

January 17, 2022

Review of deep learning models

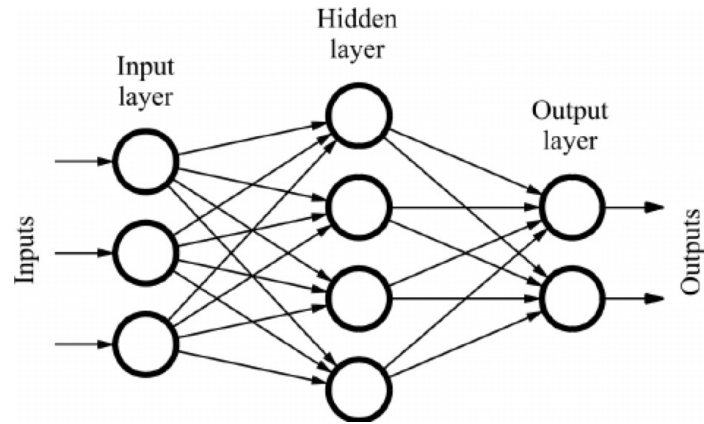
Today

- Review of basic deep learning building blocks
 - CNNs
 - RNNs
 - Attention
 - Transformers

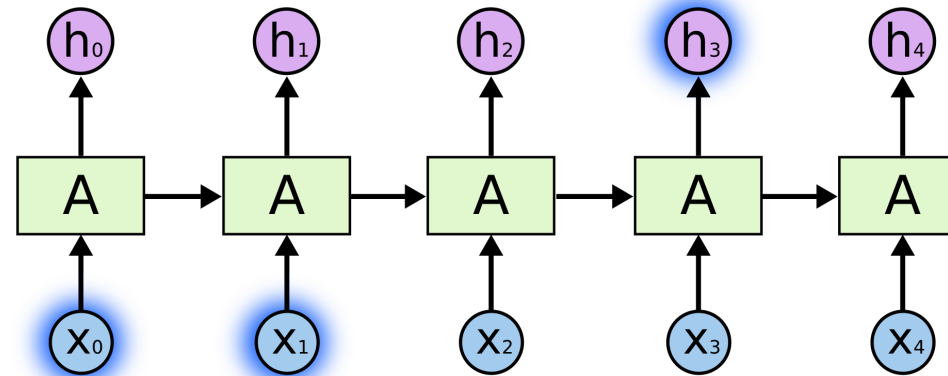
Deep learning models

Neural network architectures

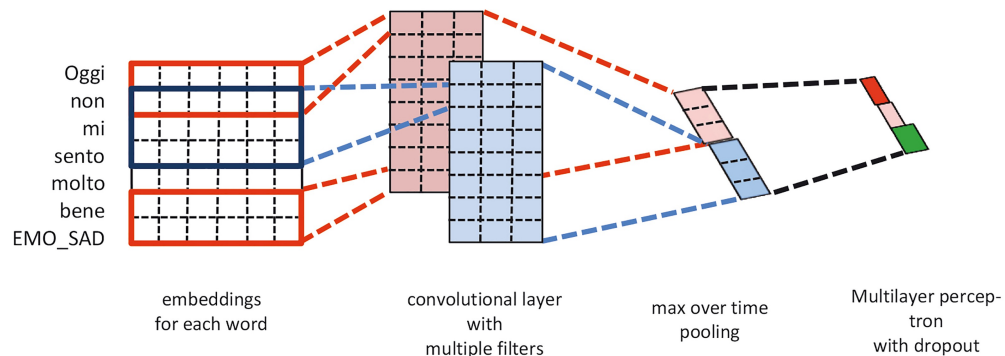
Feed-forward NNs



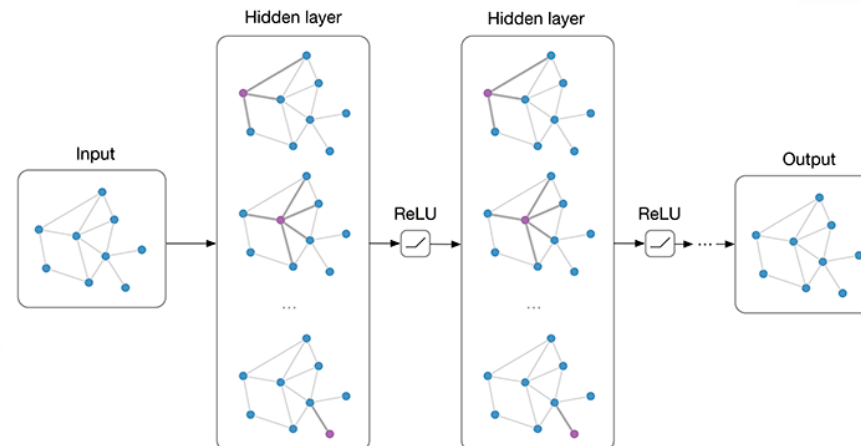
Recurrent networks (RNNs)



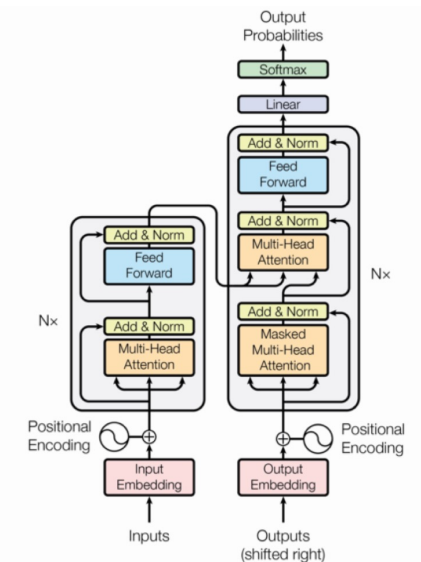
Convolution networks (CNNs)



Graph NNs



Transformers



- All network architectures can be used to model **images, text, 3D representations**, etc.
- Traditionally:
 - CNNs for images – scale/translation invariance
 - RNNs for sequences (text)
 - Transformers were introduced for machine translation
 - Now used for images and 3D shapes as well
 - Currently SOTA vision+language models are all using transformers!

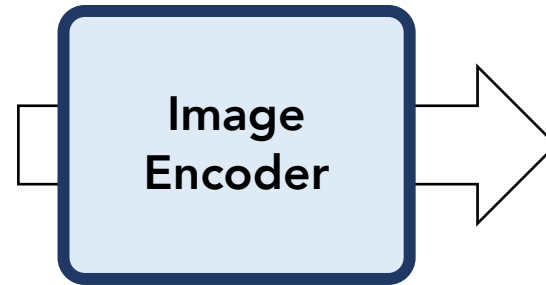
Modelling Images

Modelling Images

I



Image



V

Useful Visual Feature

Modelling Images

Convolutional Neural Networks

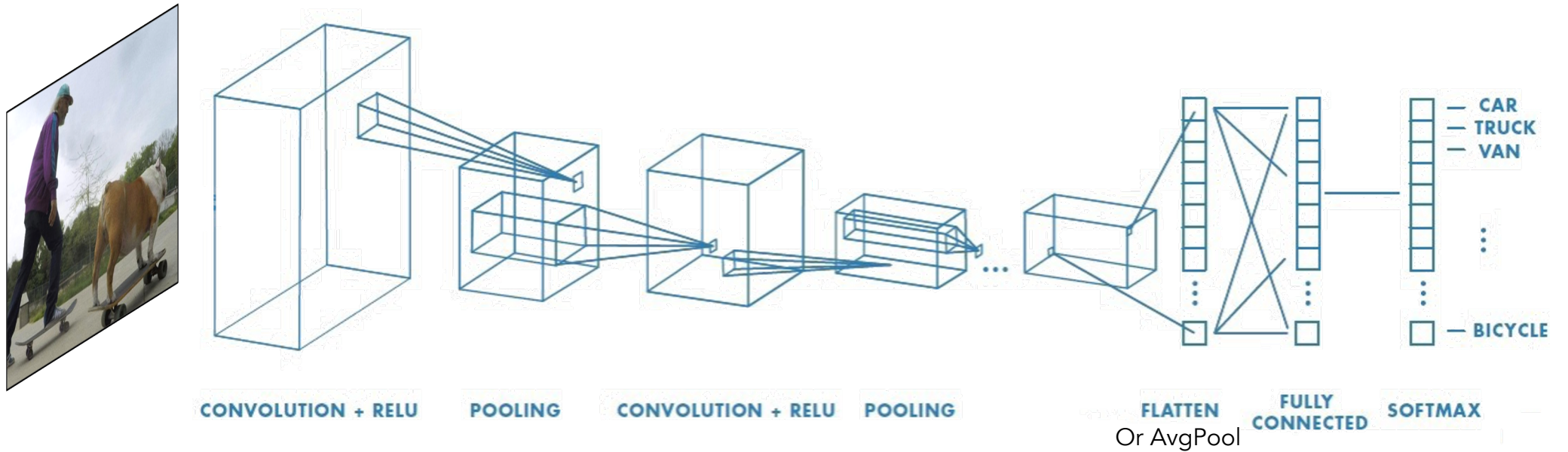
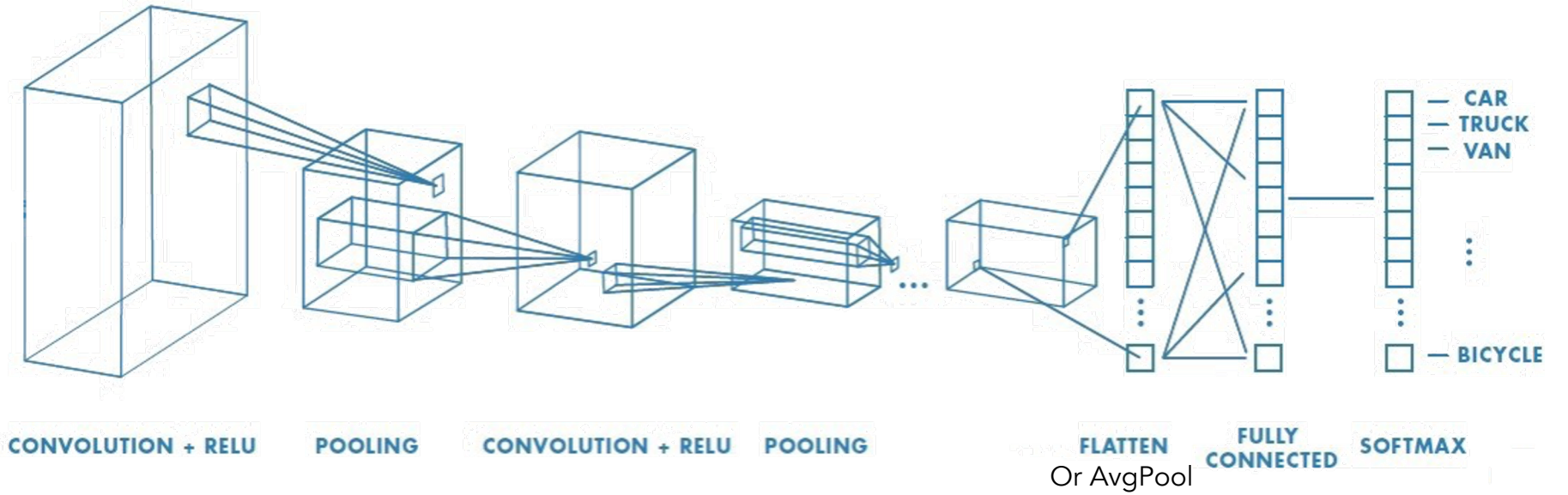
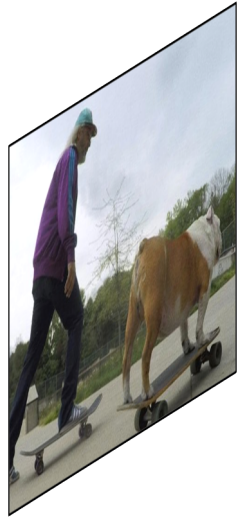


Image Credit: MathWorks

Modelling Images



Modelling Images

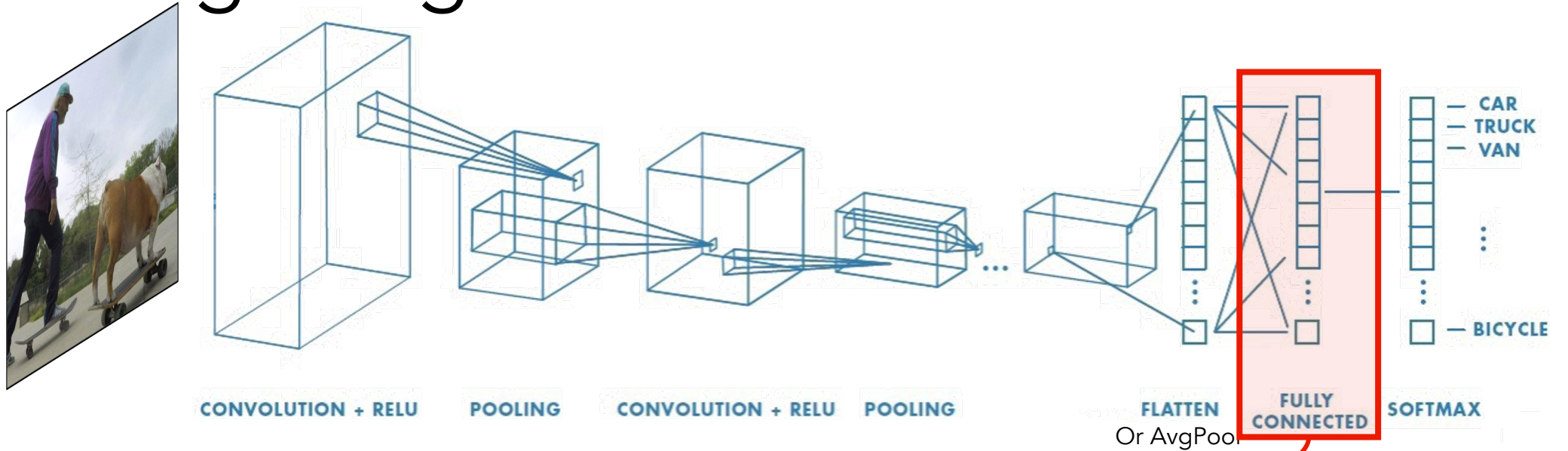
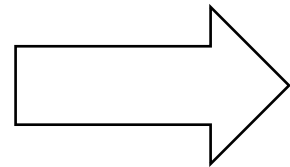


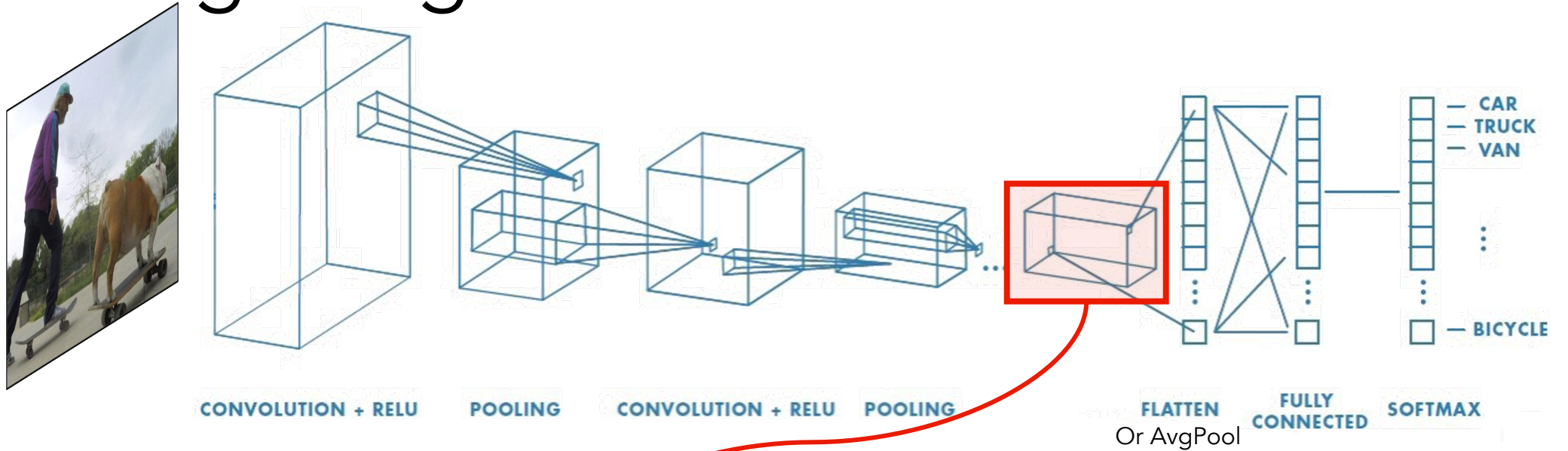
Image Level Feature



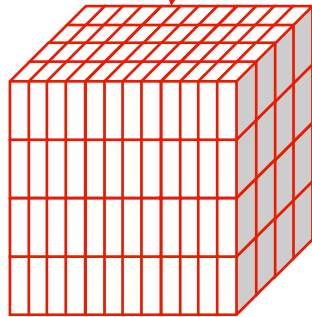
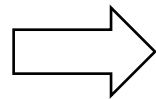
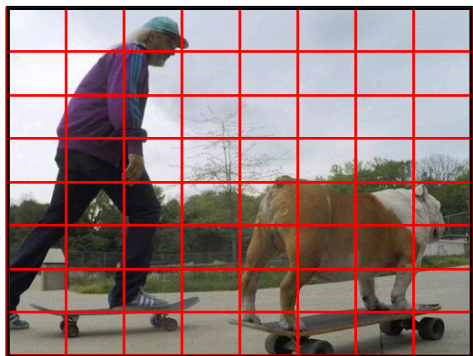
$$V \in \mathbb{R}^{1 \times d}$$

- No spatial information
- Highly compressed

Modelling Images



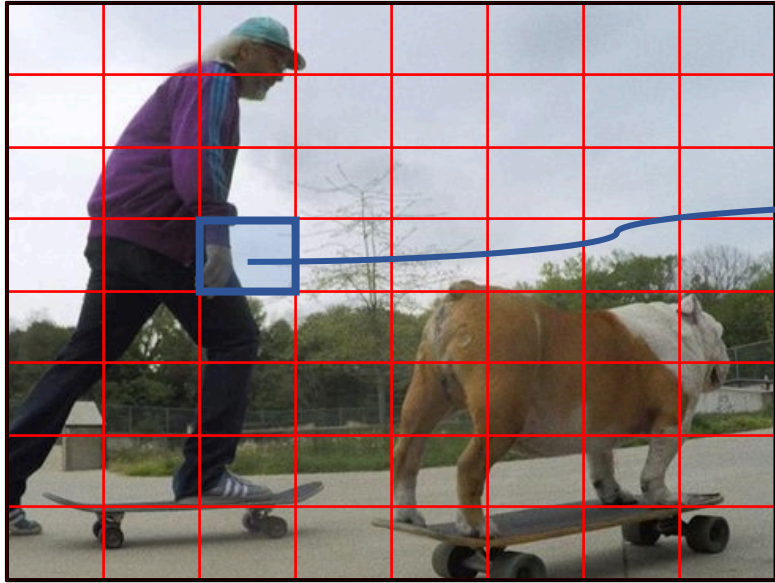
Spatial Image Features



$$V \in \mathbb{R}^{w \times h \times d}$$

- Feature vector per grid cell
- Captures some spatial info
- Uniform grid...

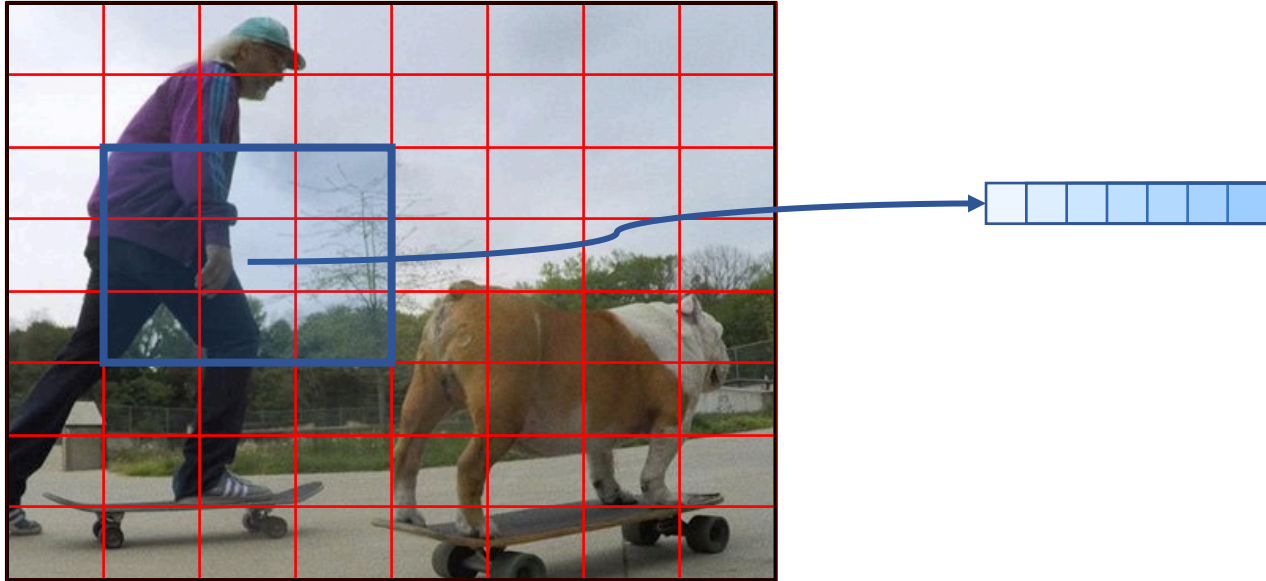
Modelling Images



Grid-based features

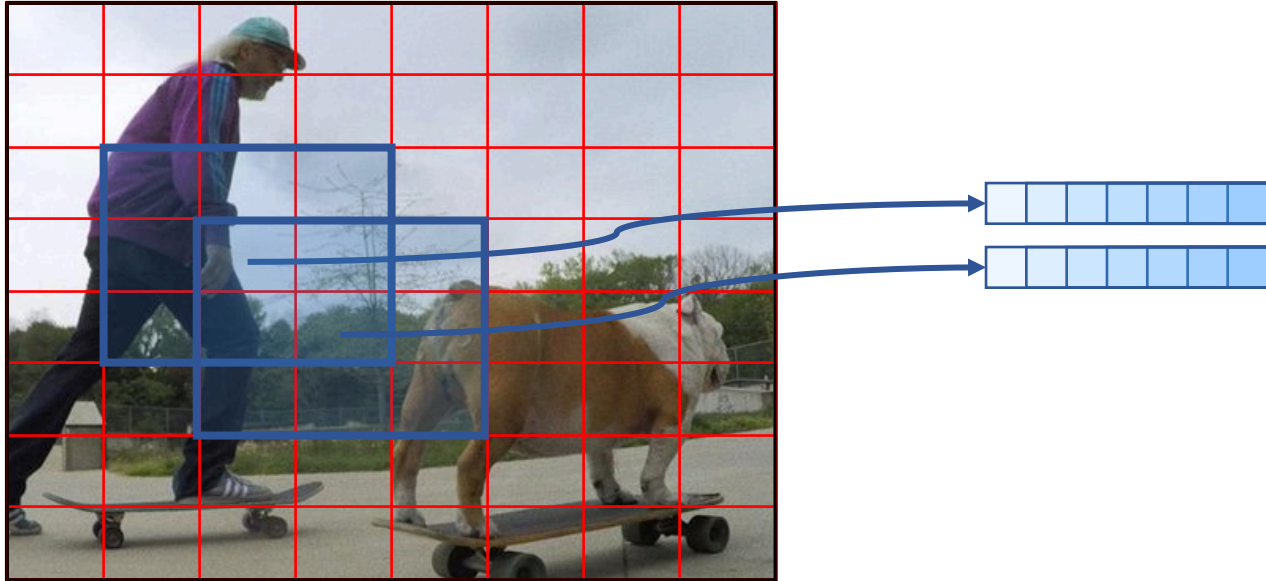


Modelling Images



*Considering receptive field it is actually much more like

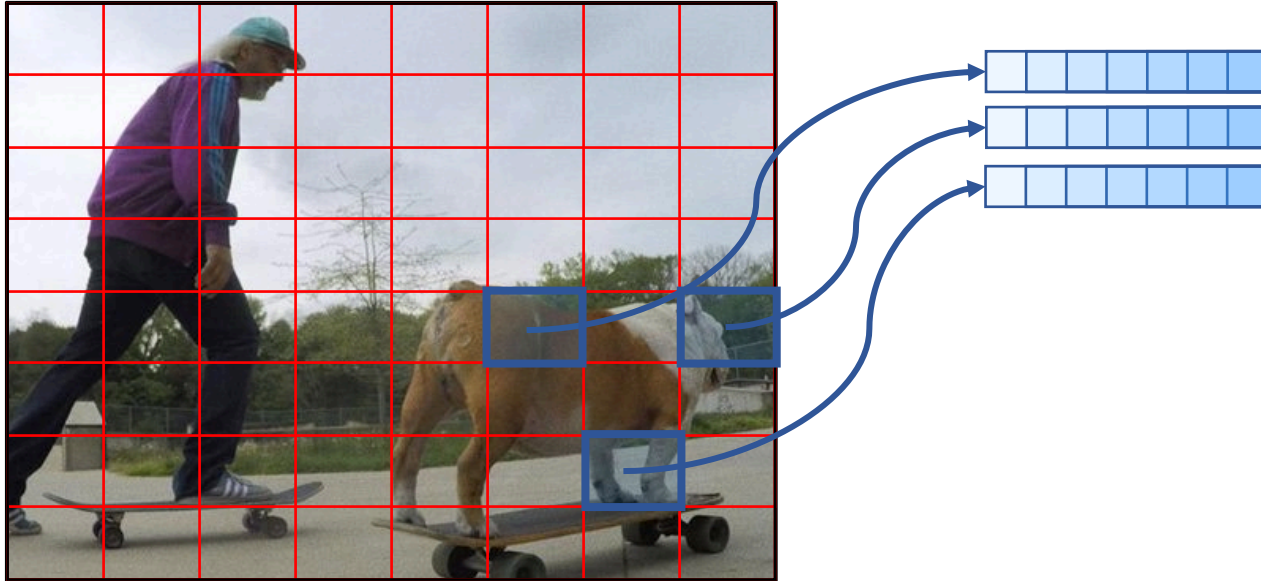
Modelling Images



*Considering receptive field it is actually much more like

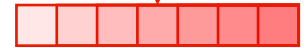
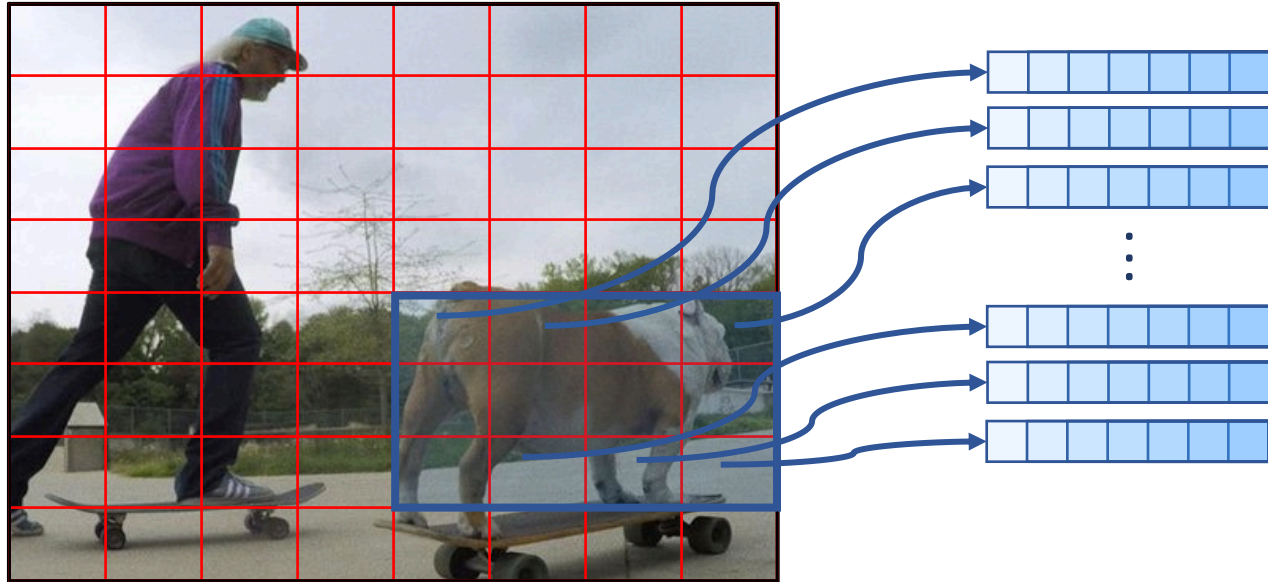
Modelling Images

"dog"



Modelling Images

"dog"



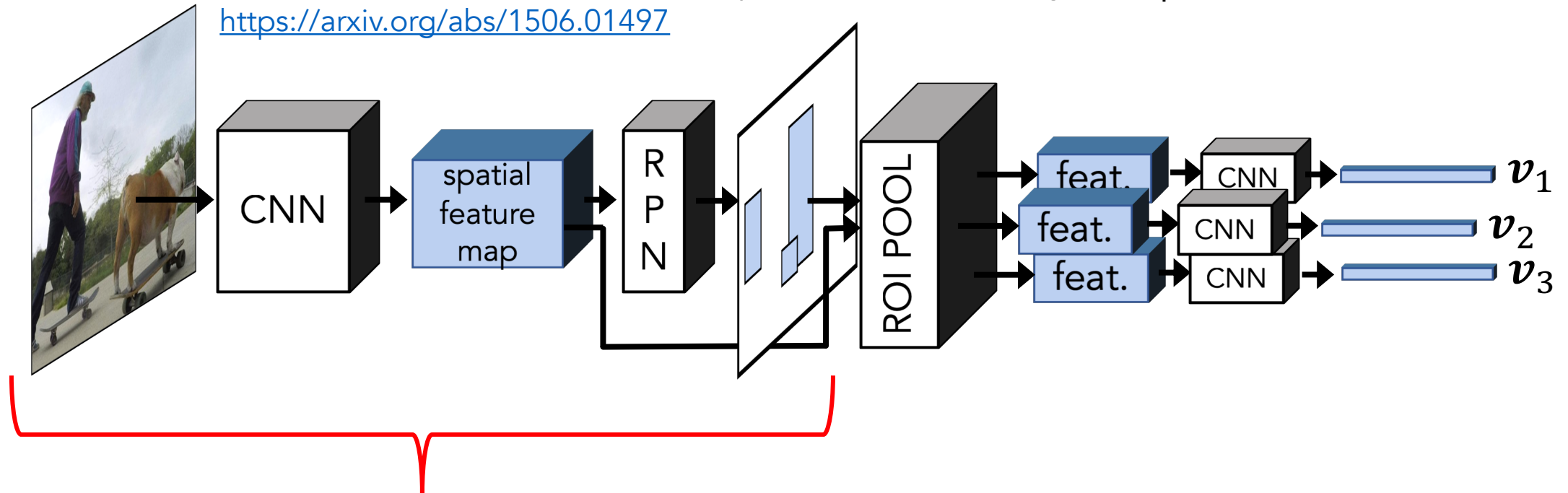
Modelling Images

Idea: Switch to **object detection** models as the backbone for image representation

- Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering arxiv.org/abs/1707.07998

Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks

<https://arxiv.org/abs/1506.01497>

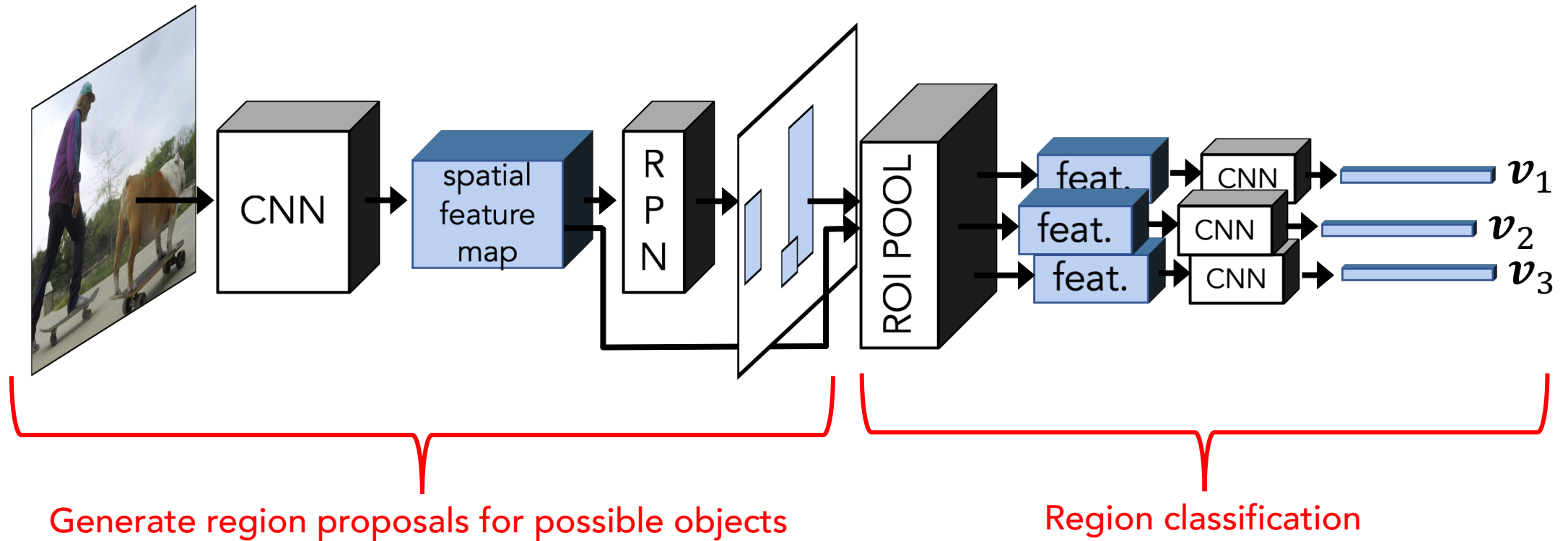


Generate region proposals for possible objects

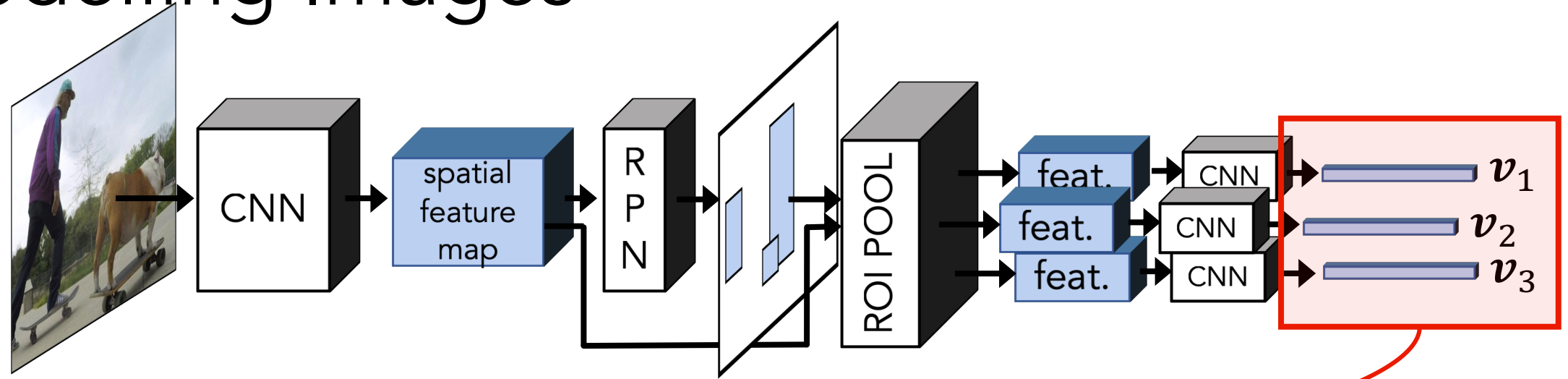
Modelling Images

Idea: Switch to **object detection** models as the backbone for image representation

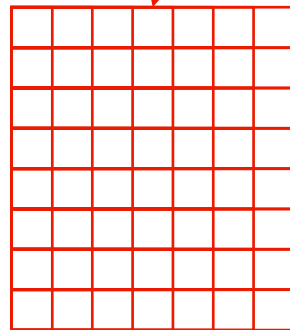
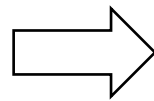
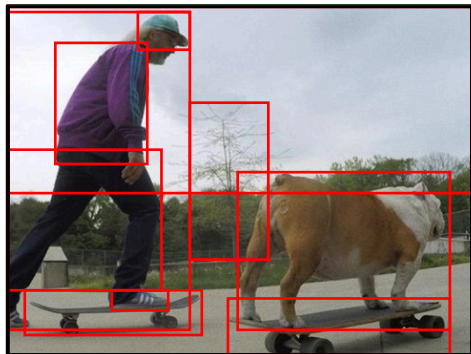
- Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering arxiv.org/abs/1707.07998



Modelling Images



Object-Centric Image Features



$$V \in \mathbb{R}^{k \times d}$$

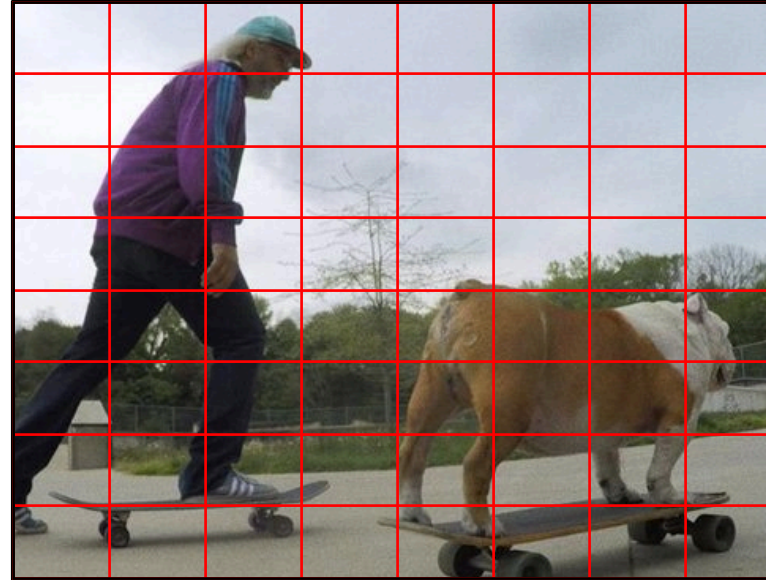
- Feature vector per bounding box
- Spatial features can be added
- Object-centric

Modelling Images

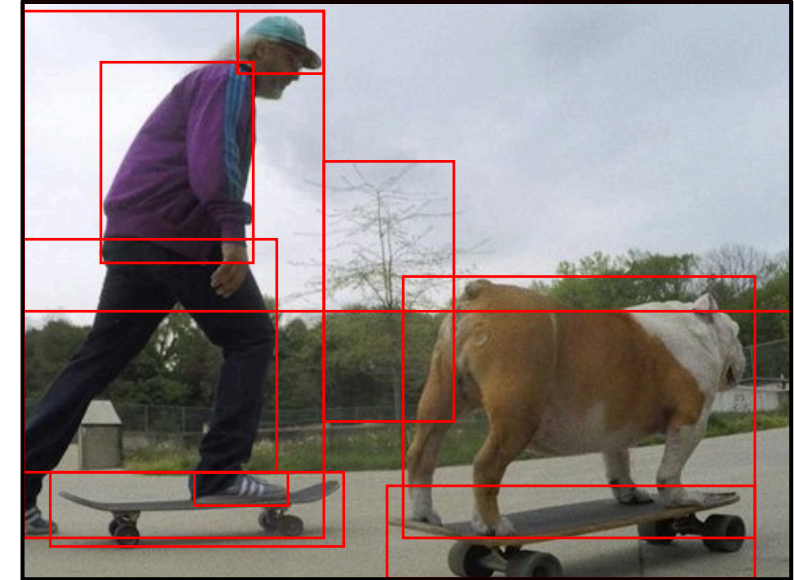
Image Level Features



Spatial / Conv Features



Detection Features



ResNet 101
Trained on ImageNet

FasterRCNN - ResNet 101
Trained on Visual Genome

These are almost never fine-tuned for downstream tasks in vision-and-language.

Modelling Images: Pretraining

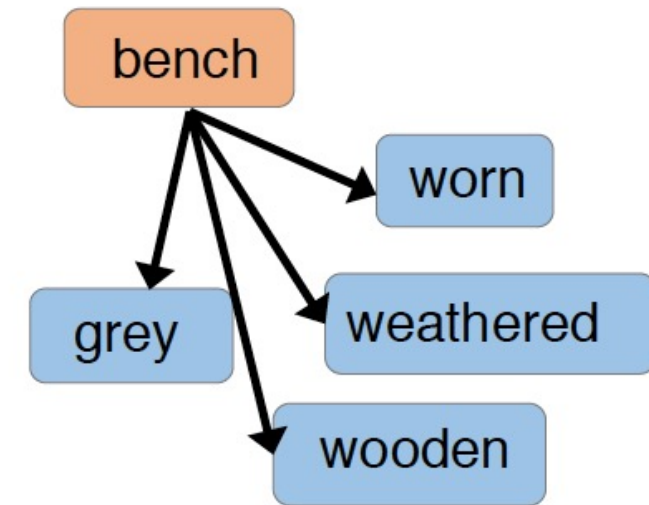
ResNet 101 Pre-training on ImageNet

- 1000 object classes (many fine-grained)



Faster R-CNN Pre-training on Visual Genome

- 1600 object classes
- 400 attribute classes



Modelling Sequences

Modelling Sequences

Recurrent Neural Networks

- Ideal for processing sequential data containing possibly long-term dependencies.
- Various implementations (e.g. simple RNN, LSTM, GRU) expose the same API

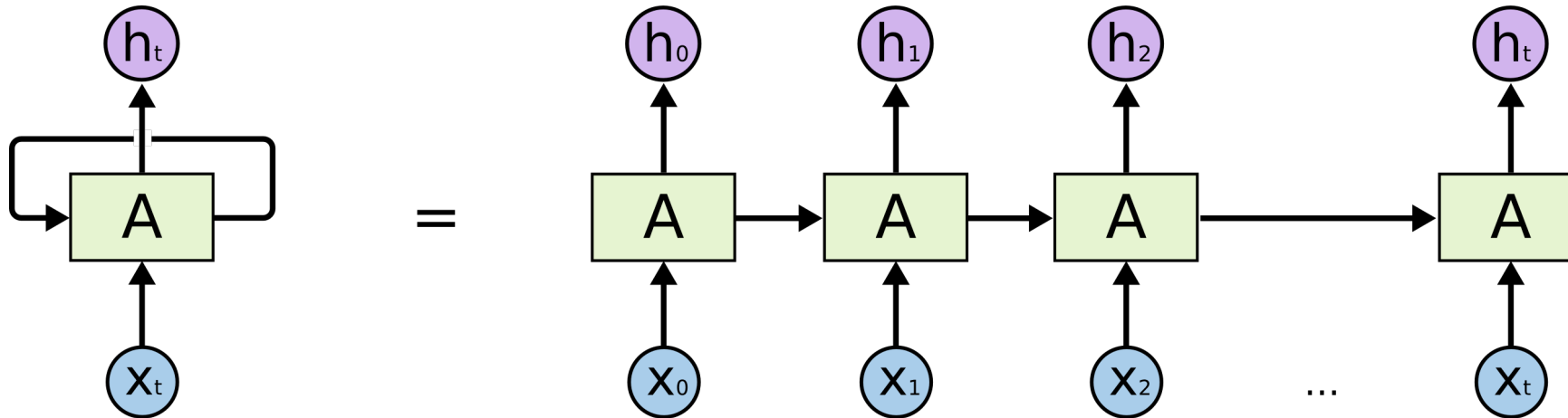


Image Credit: Christopher Olah

Modelling Sequences

Recurrent Neural Networks

- Ideal for processing sequential data containing possibly long-term dependencies.
- Various implementations (e.g. simple RNN, LSTM, GRU) expose the same API

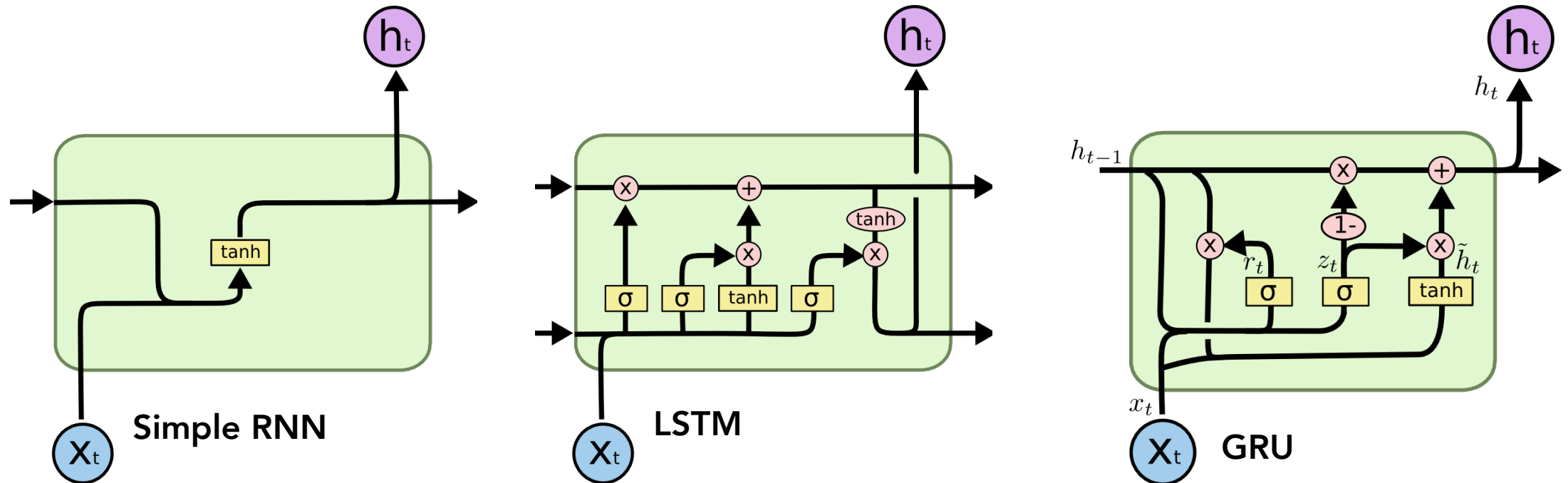
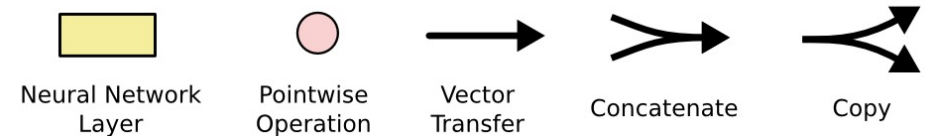


Image Credit: Christopher Olah

Slide credit: Stefan Lee



Modelling Sequences

Recurrent Neural Networks

- Ideal for processing sequential data containing possibly long-term dependencies.
- Various implementations (e.g. simple RNN, LSTM, GRU) expose the same API

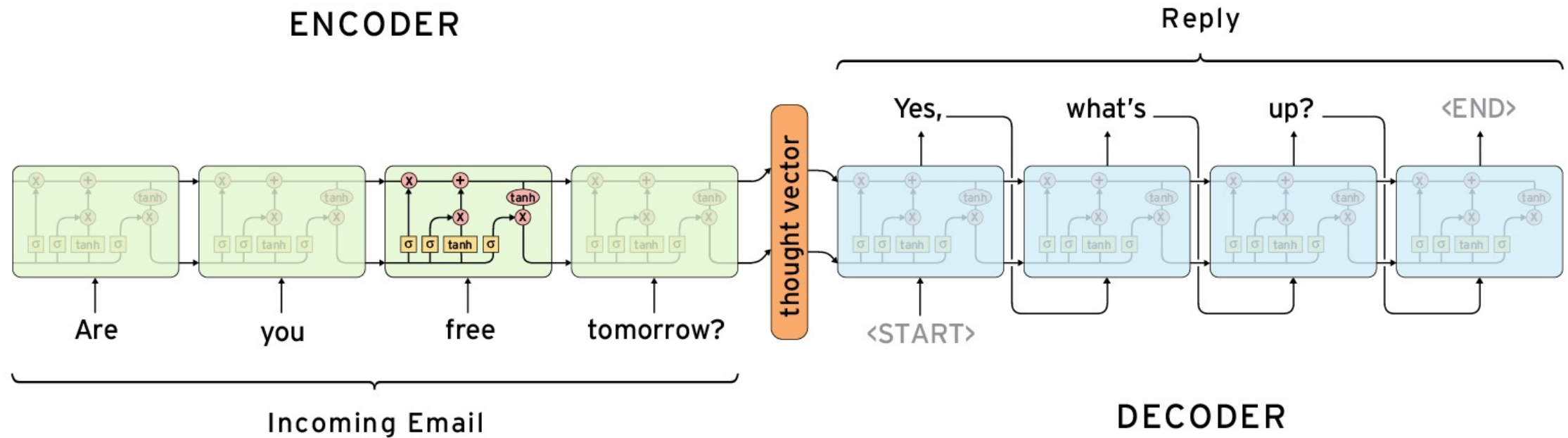
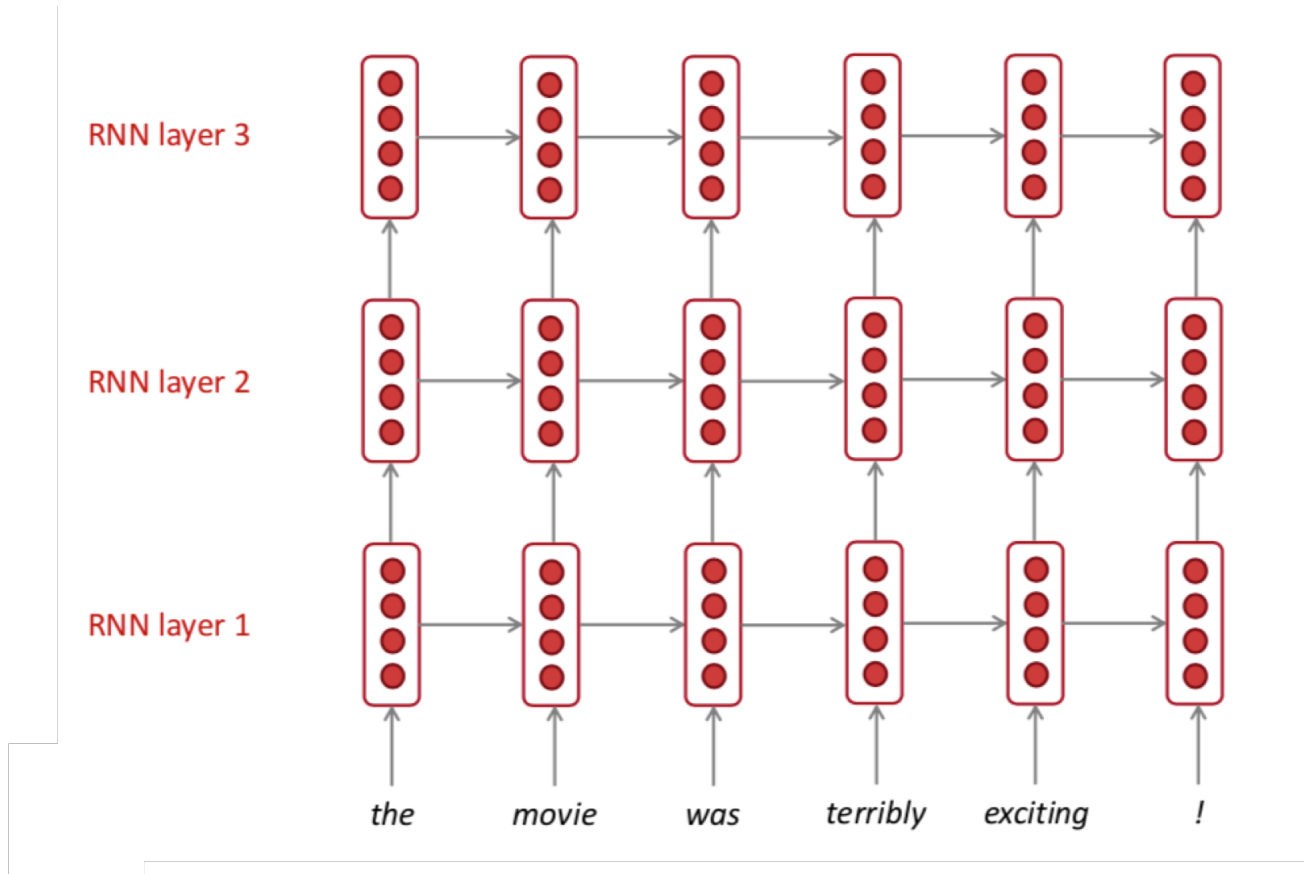


Image Credit: Christopher Olah

Multi-layer (stacked) RNNs

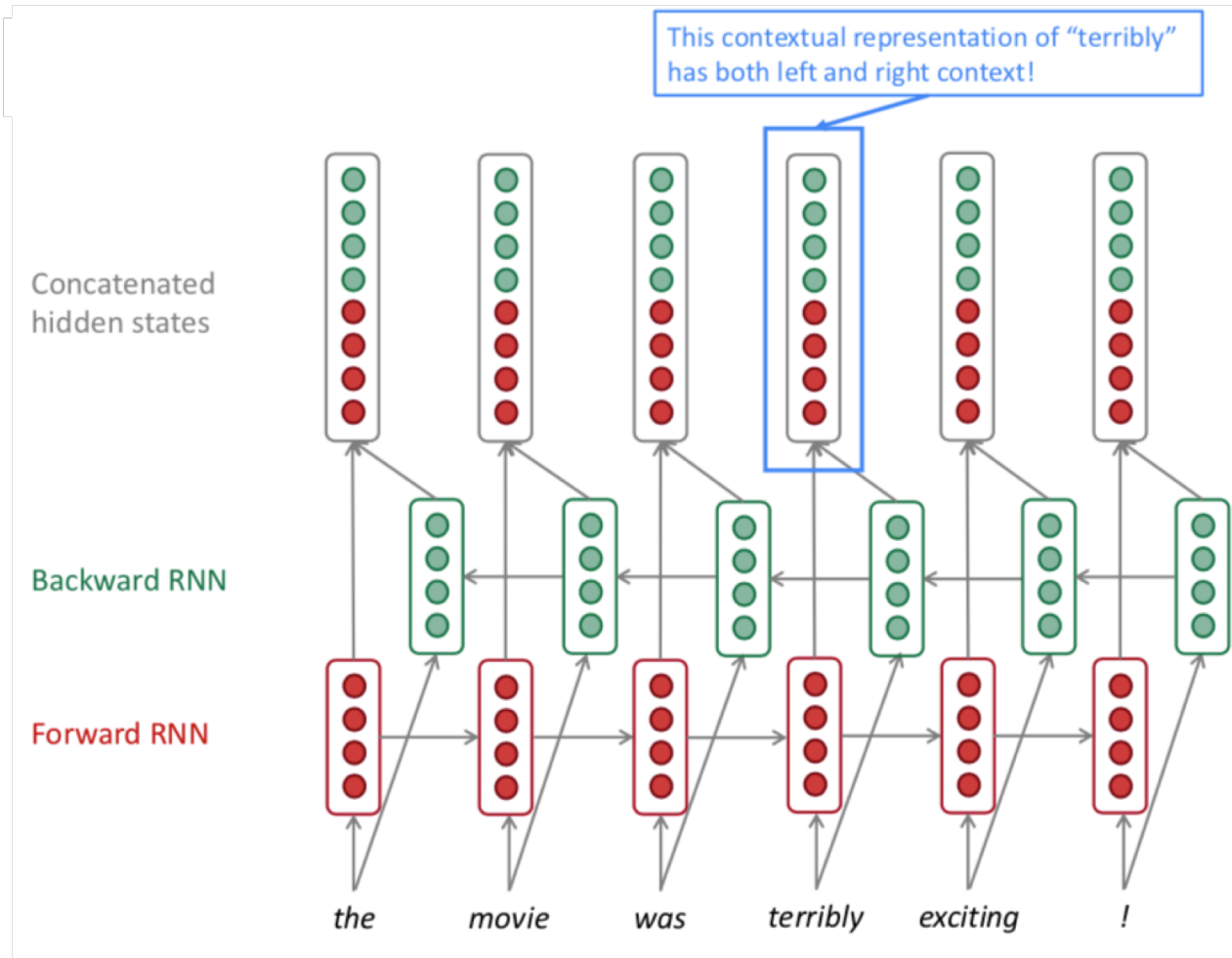


The hidden states from RNN layer i are the inputs to RNN layer $i + 1$

In practice, using 2 to 4 layers is common (usually better than 1 layer)

Transformer-based networks can be up to 24 layers with lots of skip-connections.

Bidirectional RNNs



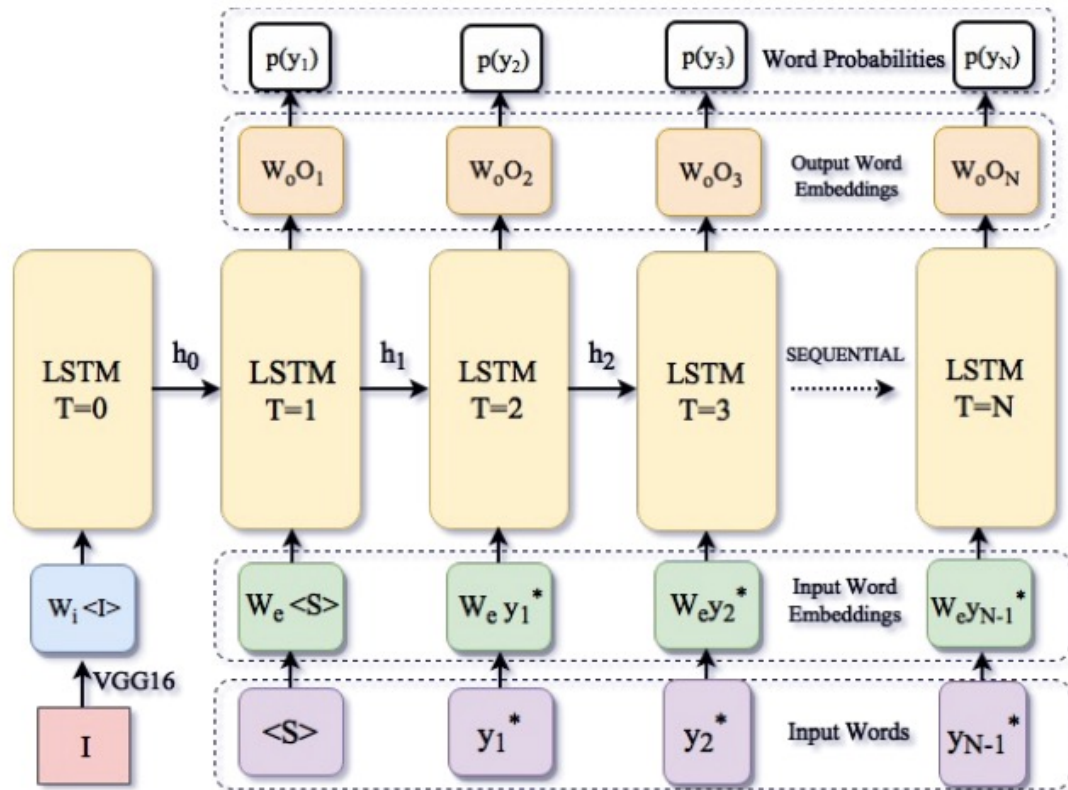
Incorporate information from both directions

Useful in encoder

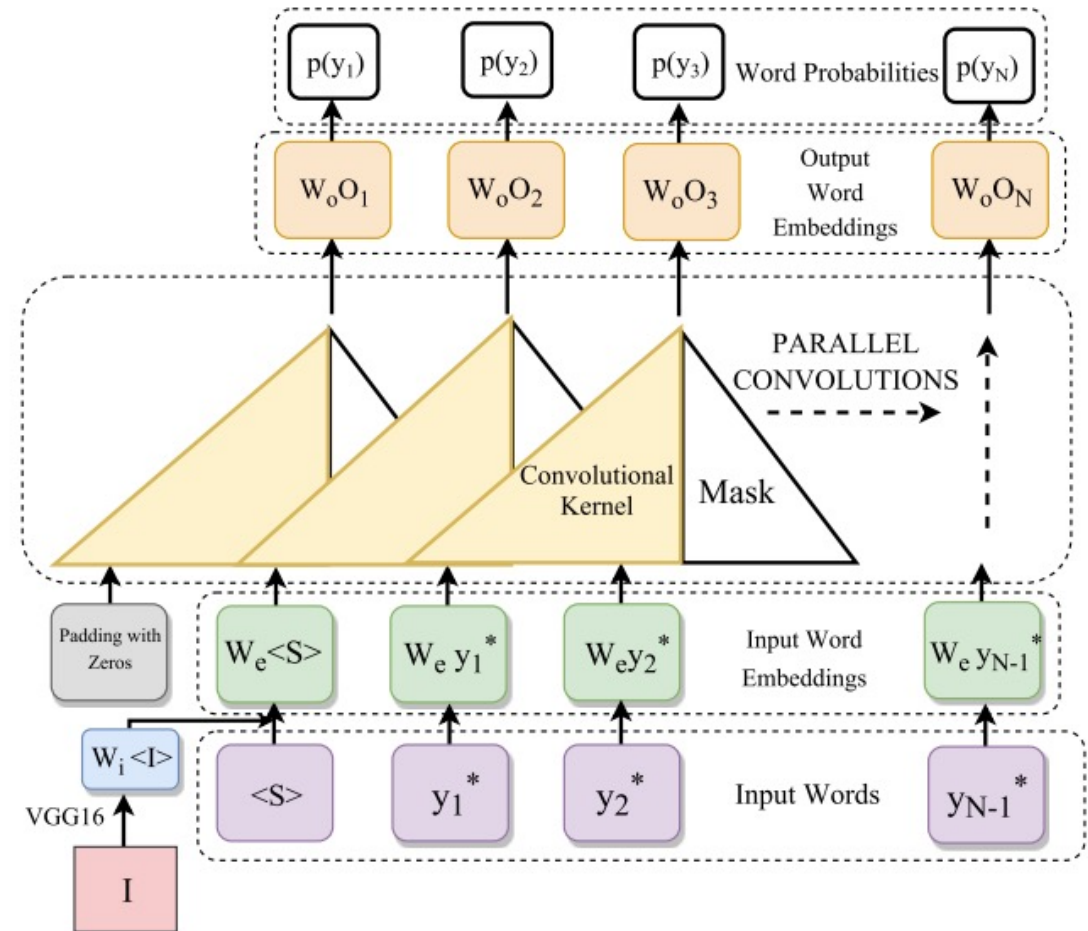
Modelling Sequences

CNNs as a fixed-time horizon alternative:

- Parallel computation!
- Tricky encoding.

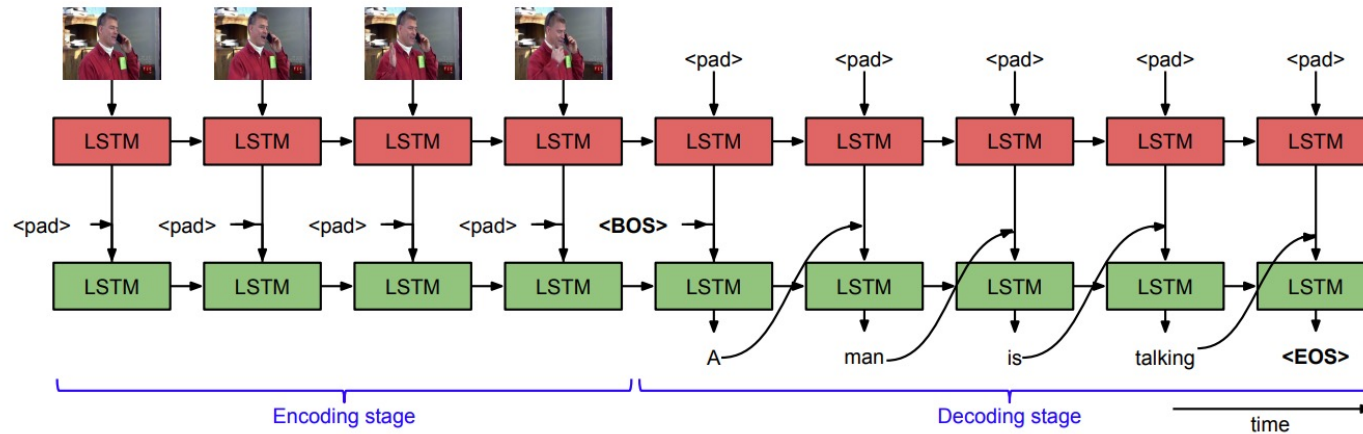


Aneja et al. CVPR 2018



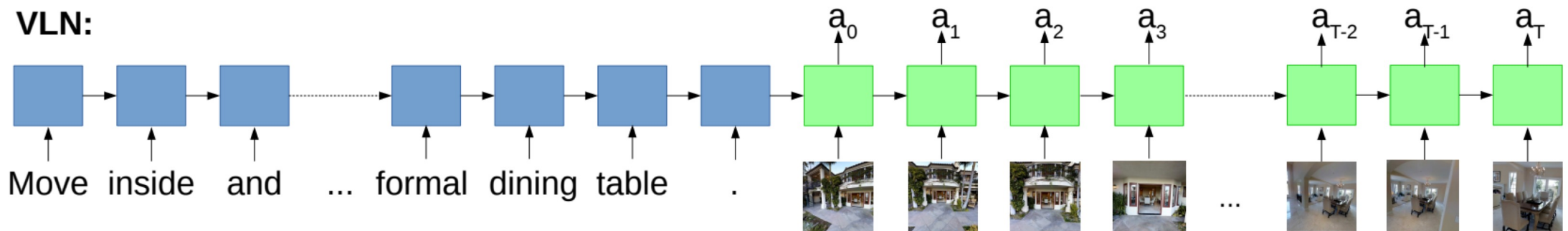
Multimodal seq2seq models

- Video captioning (video frames to text)



- Embodied AI (text + frames to actions)

VLN:

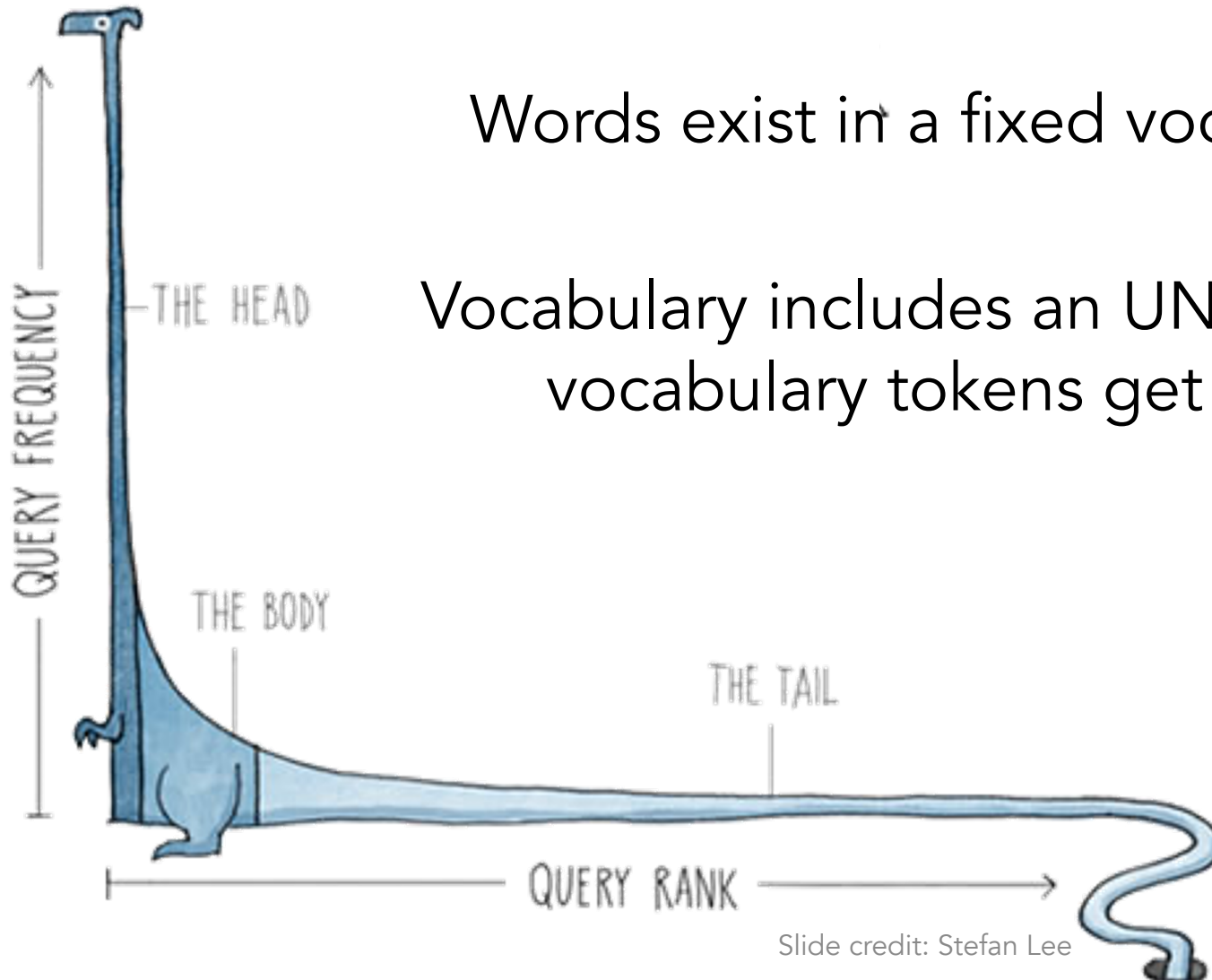


Some Notes on Representing Text

Words and Vocabularies

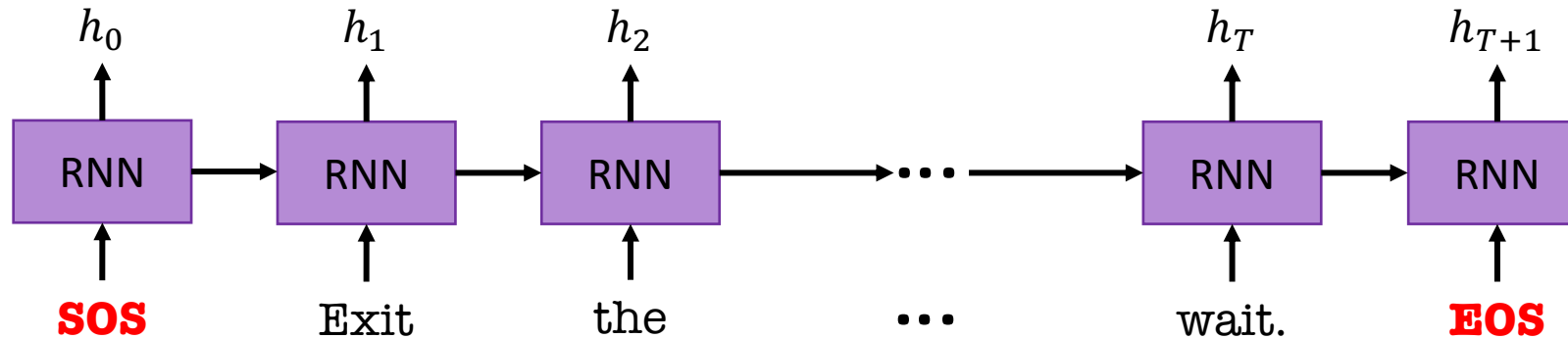
Words exist in a fixed vocabulary, i.e. $w \in V$

Vocabulary includes an UNK token – any out of vocabulary tokens get mapped to this.



Some Notes on Representing Text

Quirks of Common Practice



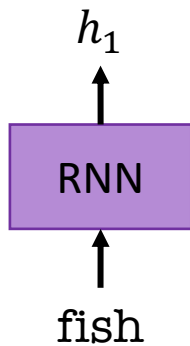
Start of Sequence

End of Sequence

Some Notes on Representing Text

What is actually input to represent tokens?

- One-hot vector \rightarrow learned representation
 - For $V = \{cat, dog, fish\}$, $w_{fish} = [0 \ 0 \ 1]$.



$$\begin{array}{c} \text{fish} \\ [0 \ 0 \ 1] \end{array} \begin{array}{c} \mathbf{W} \\ \begin{bmatrix} 0.2 & 1 & 1.5 & 0.8 & -0.2 & 1.2 \\ -1.3 & 2 & -2 & 1.2 & 0.56 & 0.1 \\ 0.13 & 0.2 & 0.95 & 0.2 & -1.3 & 0.5 \end{bmatrix} \end{array}$$

$$w_{fish} * W = [0.13 \ 0.2 \ 0.95 \ 0.2 \ -1.3 \ 0.5]$$

Initialize to random vectors and learn the embeddings during training

Some Notes on Representing Text

What is actually input to represent tokens?

cat dog

- Use pretrained word embeddings
 - Word2Vec
 - GloVE

fish

$$w_{fish} = GloVE("fish")$$

- Can do a mix of these
 - initialize learned embeddings with pretrained values

Attention

Need for "attention"

- Uses encoding of entire input when generating each output token
- Maybe would be useful to **focus** on a part of the input as the output tokens are generated

Translation

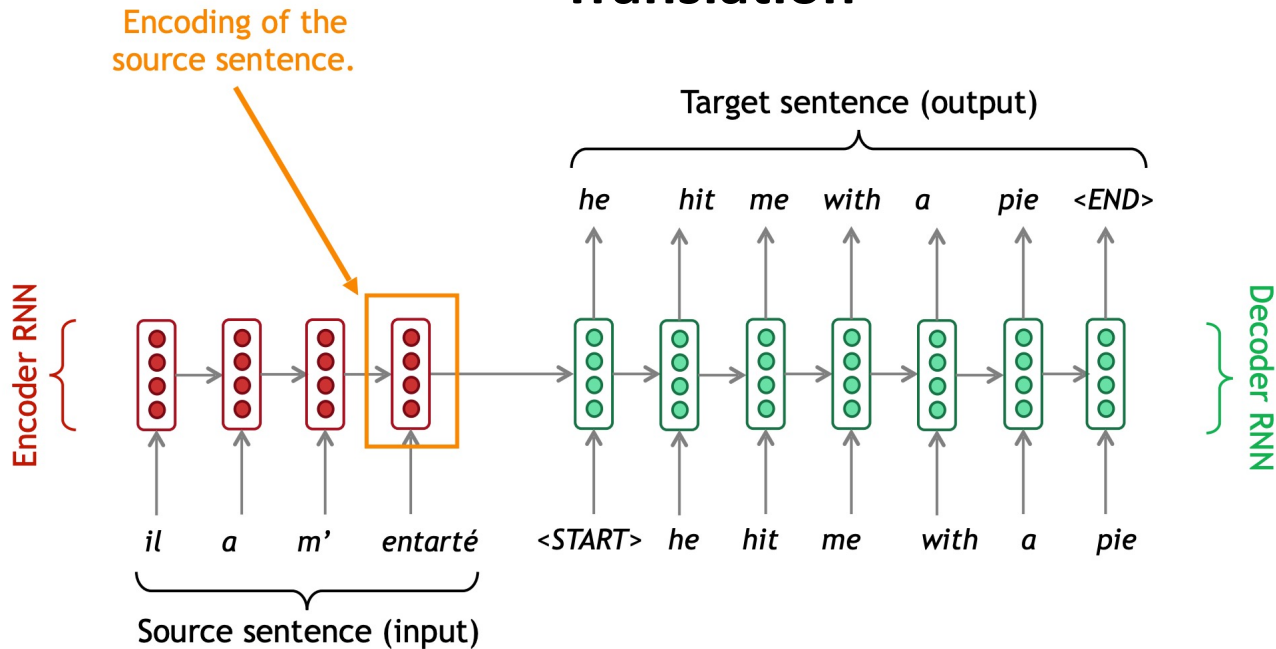


Image credit: Abigail See

Captioning

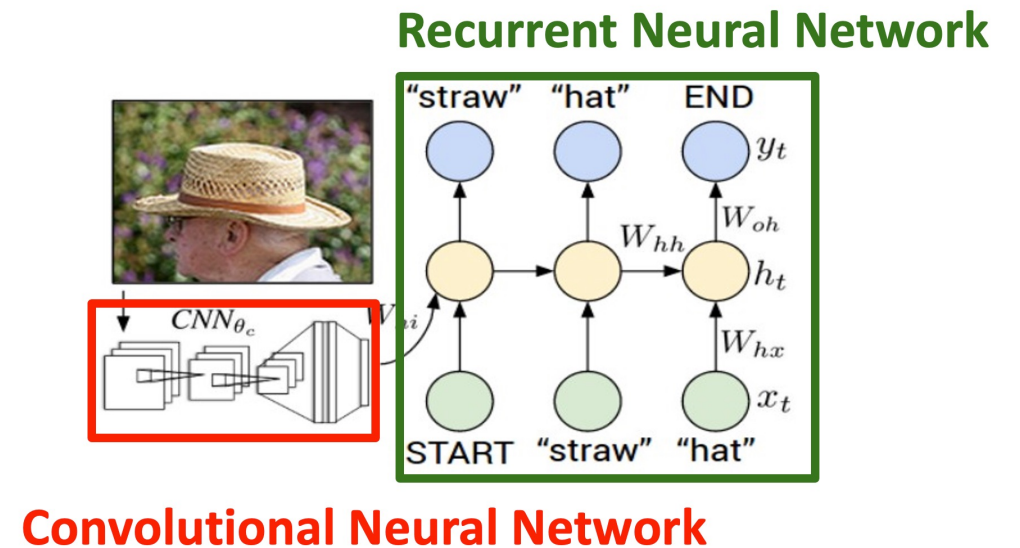
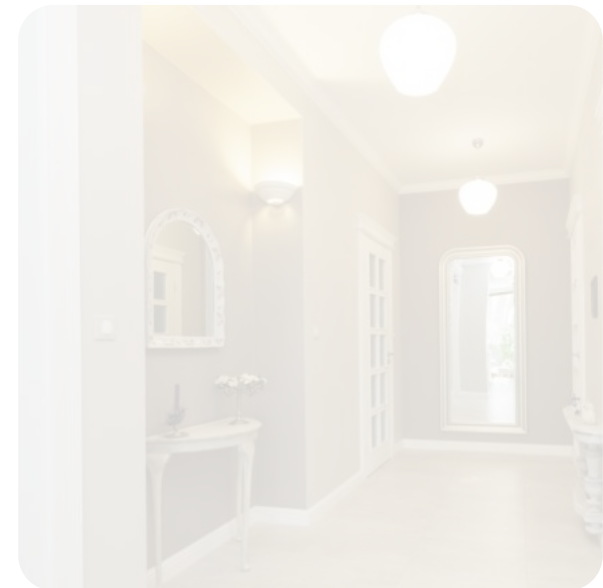


Image credit: Andrej Karpathy

Attention for VLN

Not every part of an input is important given the task context

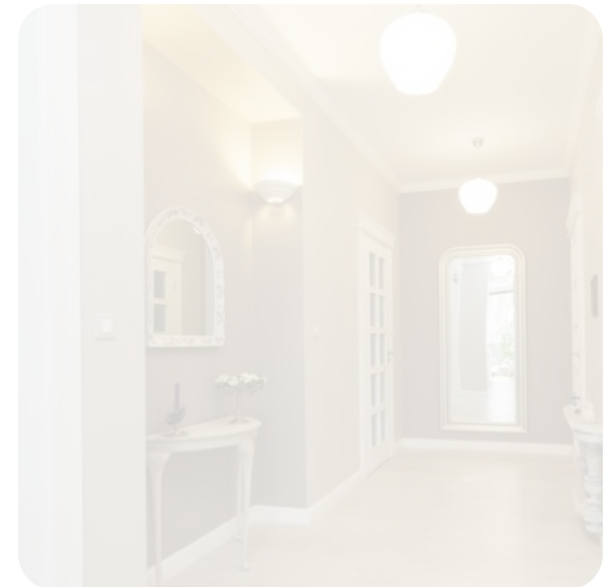
Exit the bathroom. Turn left and exit the room using the door on the left. Wait there.



Attention for VLN

Not every part of an input is important given the task context

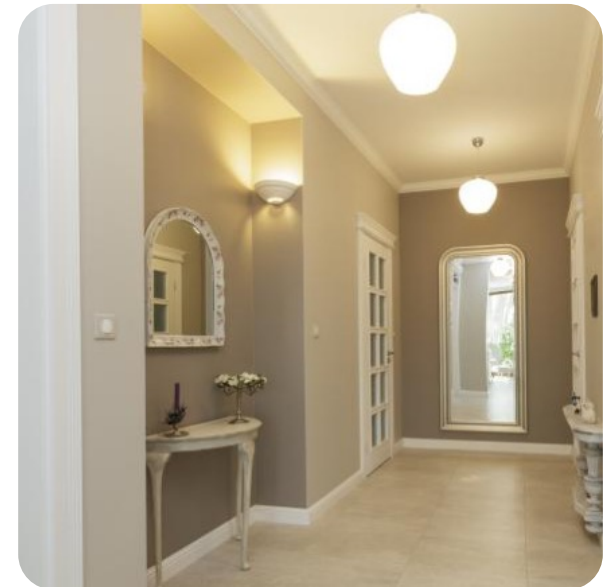
Exit the bathroom. Turn left and exit the room using the door on the left. Wait there.



Attention for VLN

Not every part of an input is important given the task context

Exit the bathroom. Turn left and exit the room using the door on the left. **Wait there.**



Attention for VLN

Not every part of an input is important given the task context

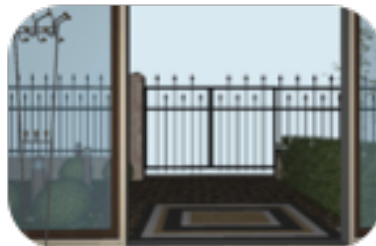
Take a right when you see the mirrored wardrobe.



Attention for VQA

Not every part of an input is important given the task context

What shape is the doormat?



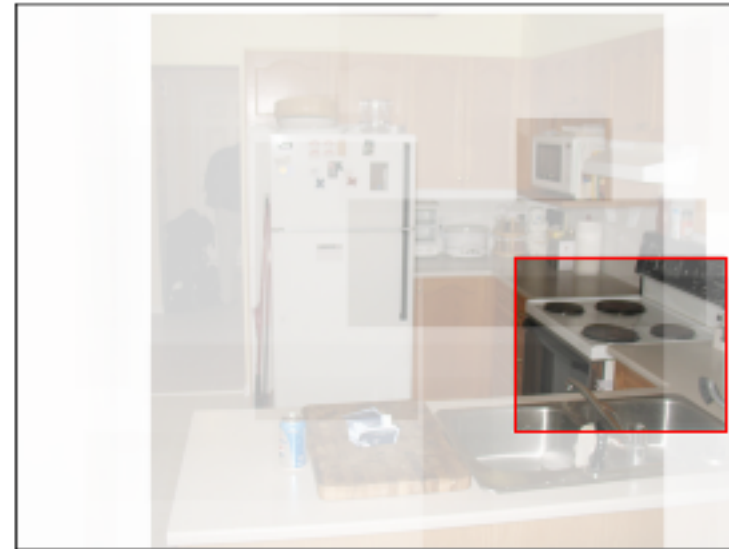
Attention

- The concept of 'attention' has seen widespread use...
 - In many language and / or vision tasks, attention works extremely well!
 - Attention improves interpretability of neural networks

Q: What room are they in?



A: kitchen



Focus on part of input

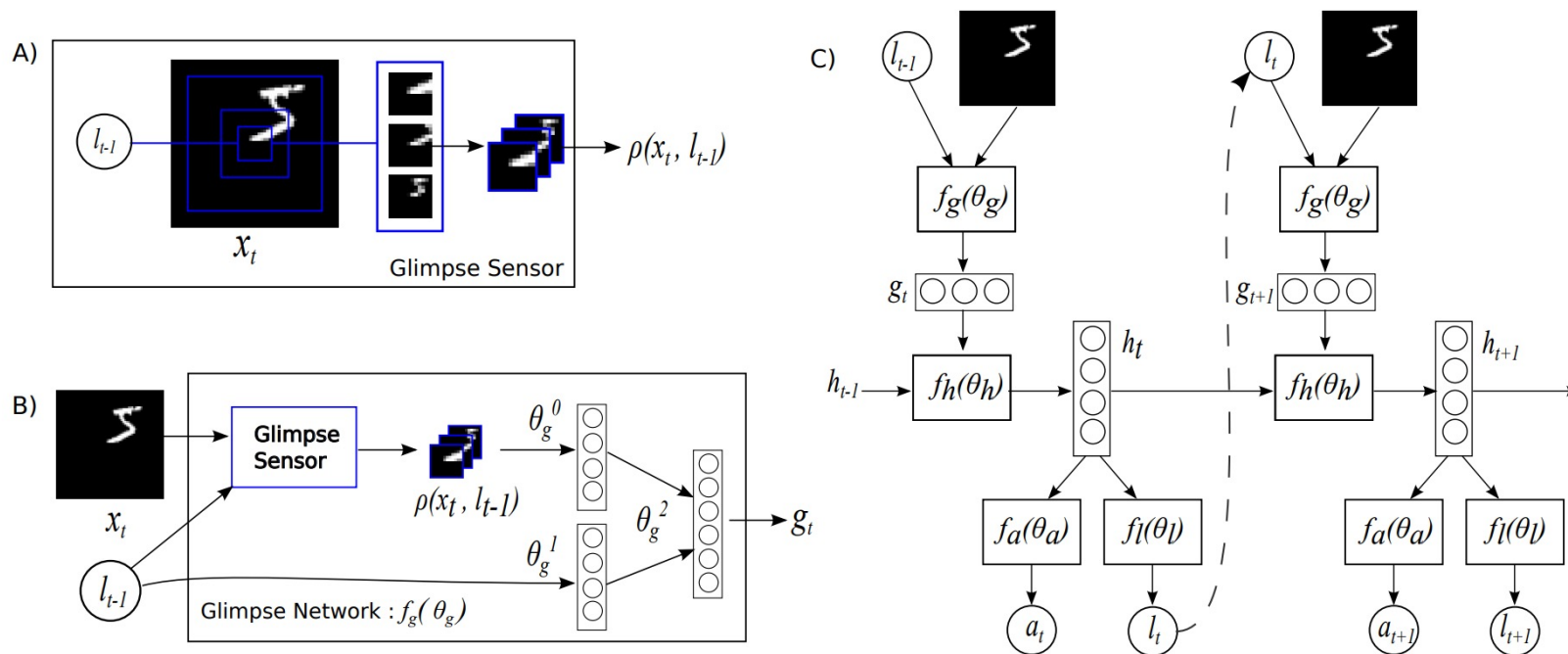
Attention mechanisms - summary

- Attention is probably one of the simplest and most effective ideas in deep learning – proven across many different domains
- Practically: focus on part of input by taking a weighted sum of different input parts
- With sufficient data, attention mechanisms can learn to ground language in visual content from 'distant' supervision
- Given the complexity of biological attention systems, assume there is still much to explore... particularly temporal aspects in context of LV&A

Attention

- Major impact on computer vision and NLP from 2014/5

Recurrent Models of Visual Attention

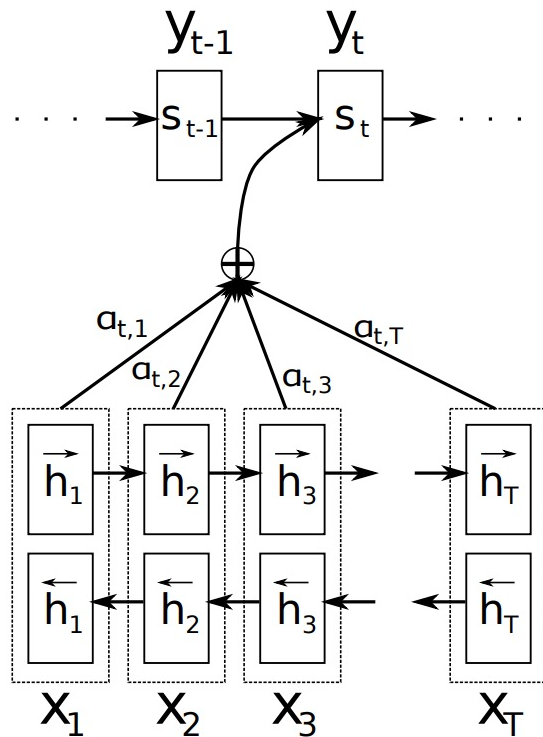


Mnih et al. NIPS 2014

Attention

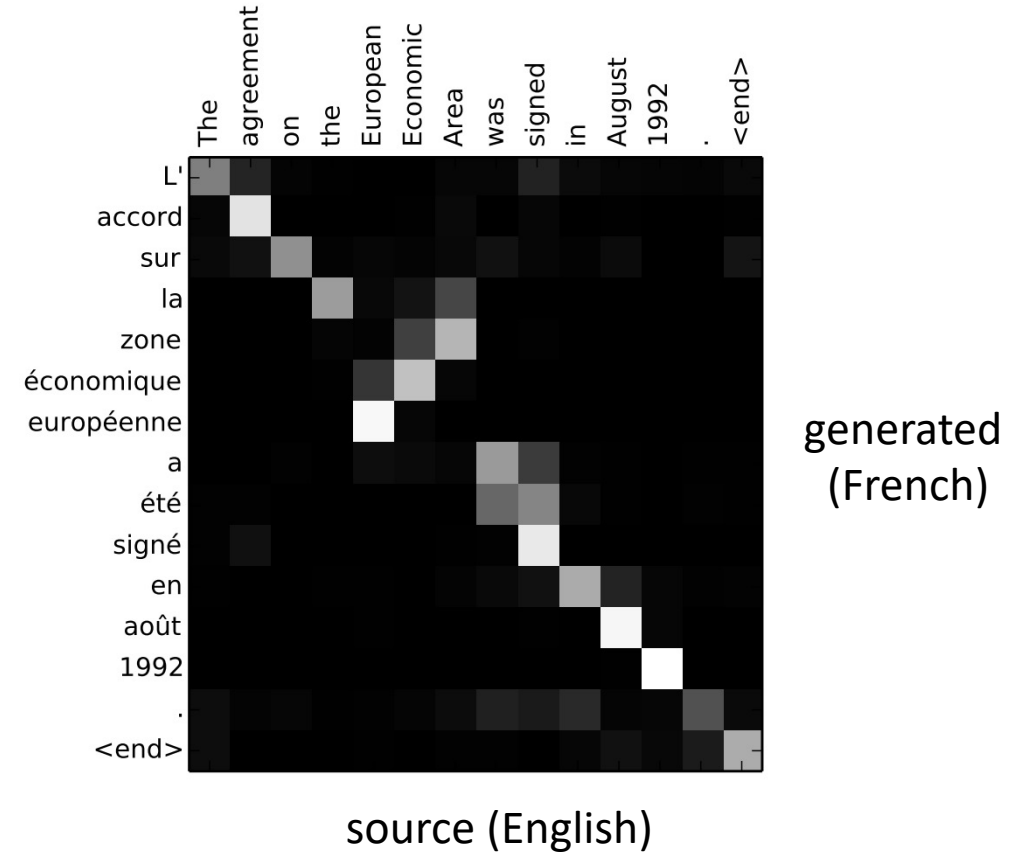
- Major impact on computer vision and NLP from 2014/5

Neural Machine Translation by Jointly Learning to Align and Translate



Bahdanau et al. ICLR 2015

Slides Credit: Peter Anderson

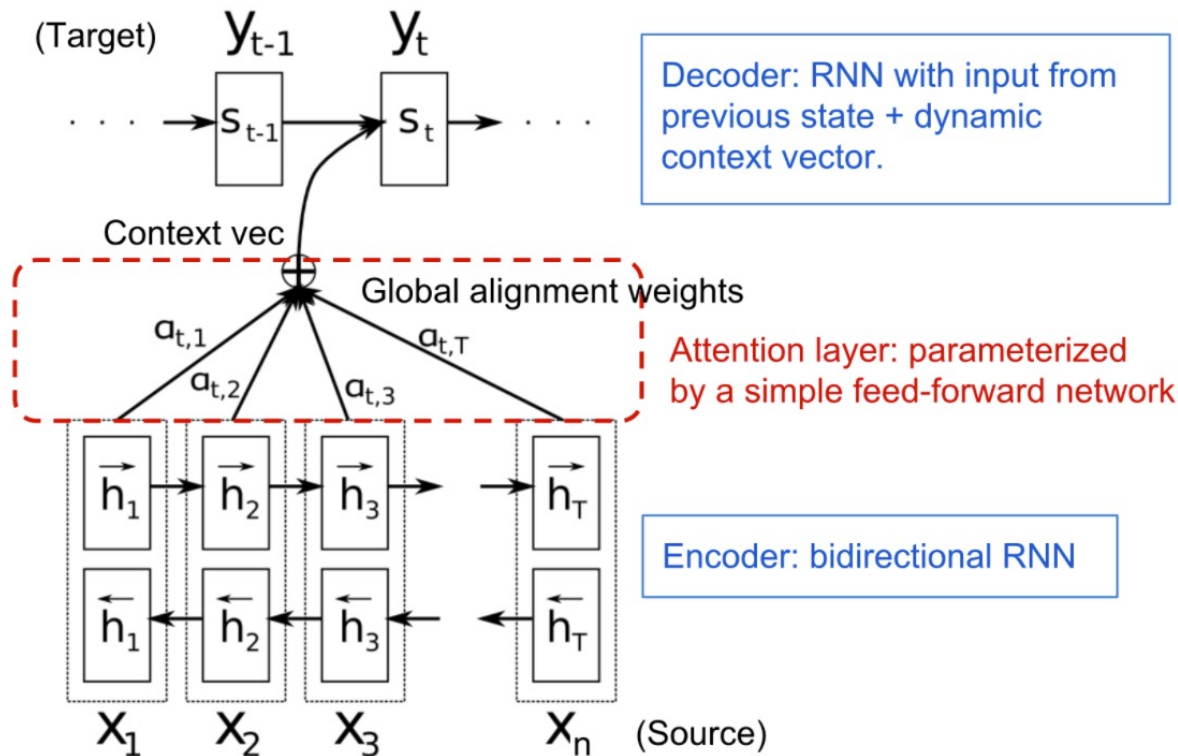


Attention weights α_i (0 = black, 1 = white)

Attention

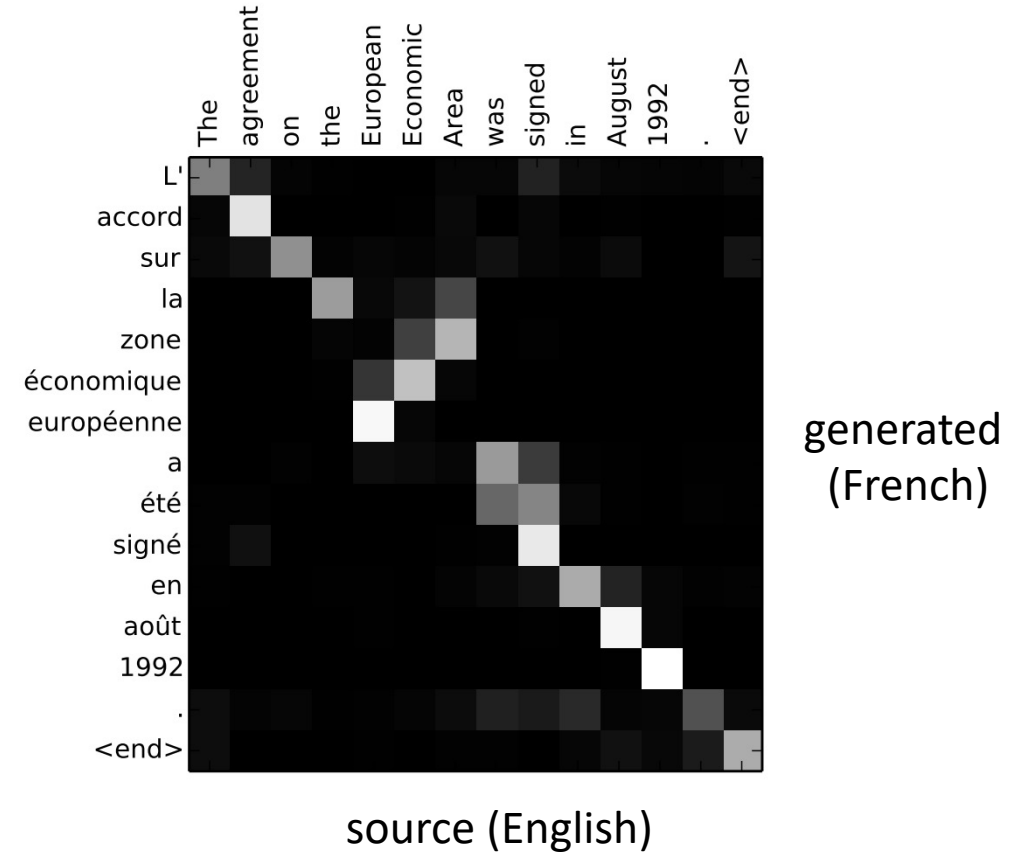
- Major impact on computer vision and NLP from 2014/5

Neural Machine Translation by Jointly Learning to Align and Translate



Bahdanau et al. ICLR 2015

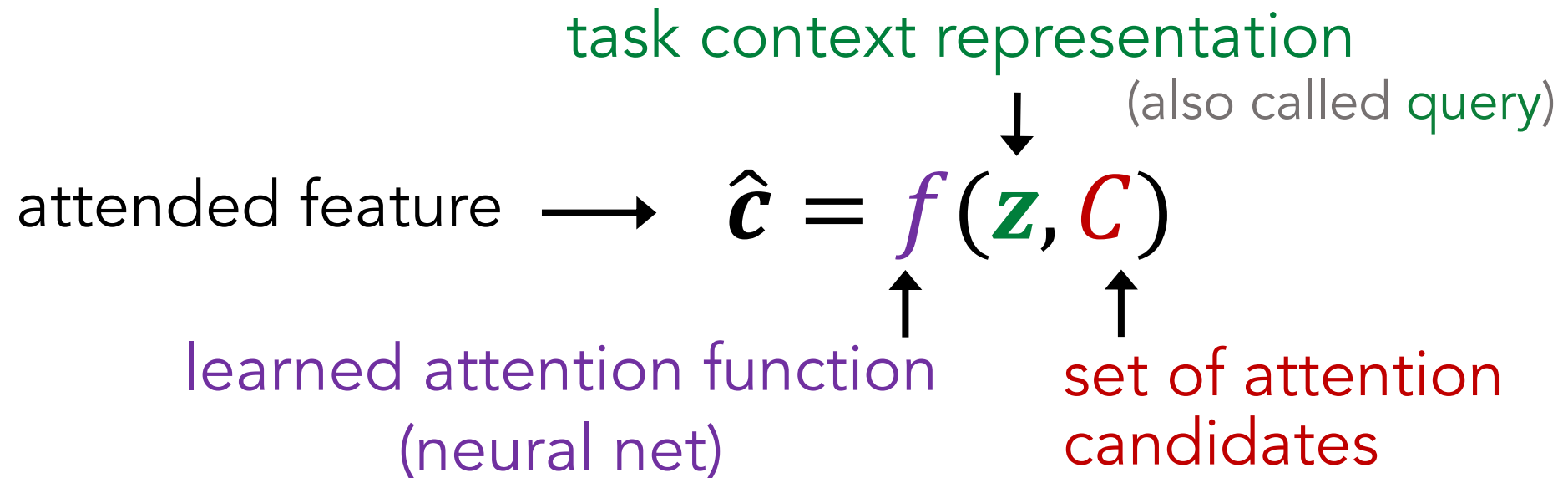
Slides Credit: Peter Anderson



Attention

Attention in Neural Networks:

A learned mechanism that **learns to focus** on a subset of the **input** that is most **relevant to the current task**.



Computing attention

Attention function, f

$$a_i = g(\mathbf{c}_i, \mathbf{z})$$

$$\boldsymbol{\alpha} = \text{softmax}(\mathbf{a})$$

$$\hat{\mathbf{c}} = \sum_{i=1}^k \alpha_i \mathbf{c}_i$$

Attention scores: \mathbf{a} (unnormalized)

Attention weights: $\boldsymbol{\alpha}$ (normalized)

Final attention output

Weighted sum of context features

Attention score $a_i = g(\mathbf{c}_i, \mathbf{z})$

how well does the attention candidate \mathbf{c}_i match the query \mathbf{z}

- Dot-product attention:

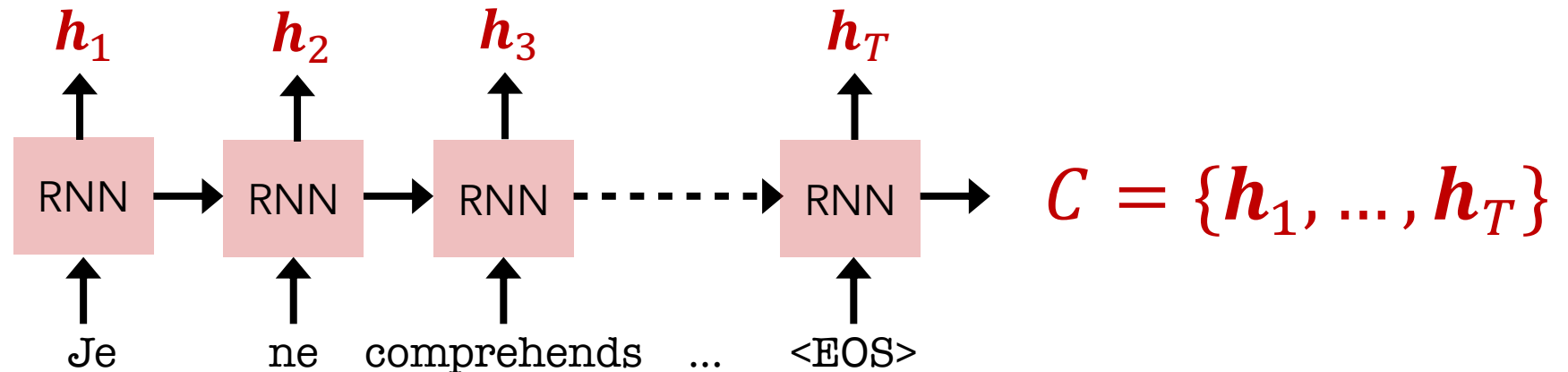
$$g(\mathbf{c}_i, \mathbf{z}) = \mathbf{z}^\top \mathbf{c}_i$$

- Neural network

$$g(\mathbf{c}_i, \mathbf{z}) = v^\top \tanh(W_1 \mathbf{c}_i + W_2 \mathbf{z})$$

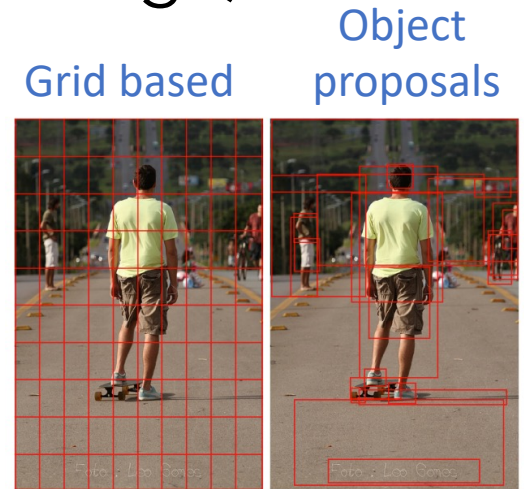
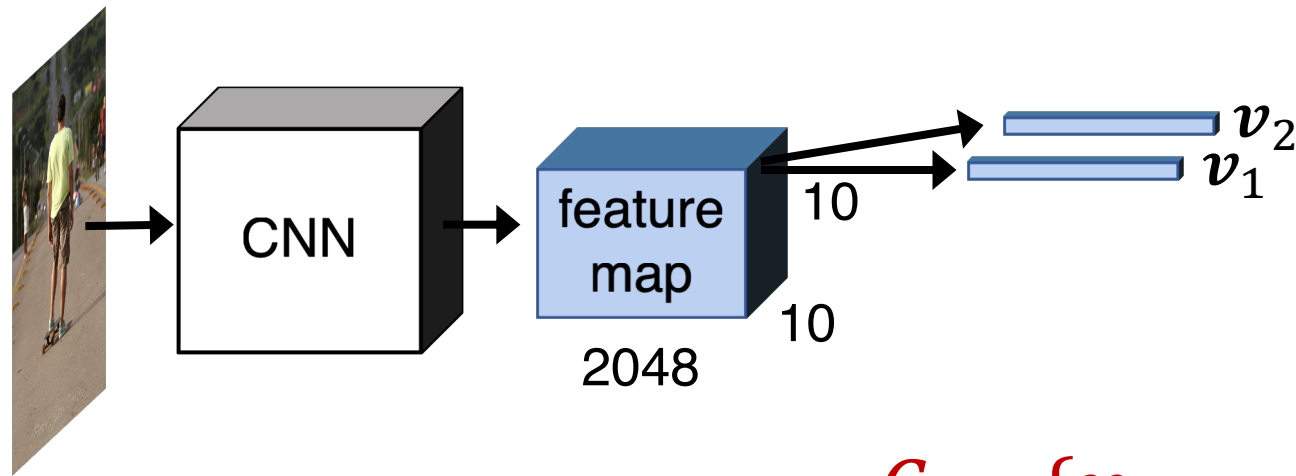
Attention over Text

Attention candidates, \mathcal{C} typically defined by the hidden states of an encoder (e.g. one feature vector for each word in the input text)



Attention over Visual Features

- Attention candidates, \mathcal{C} typically defined by the spatial output of a CNN (feature vectors for different parts of the image)

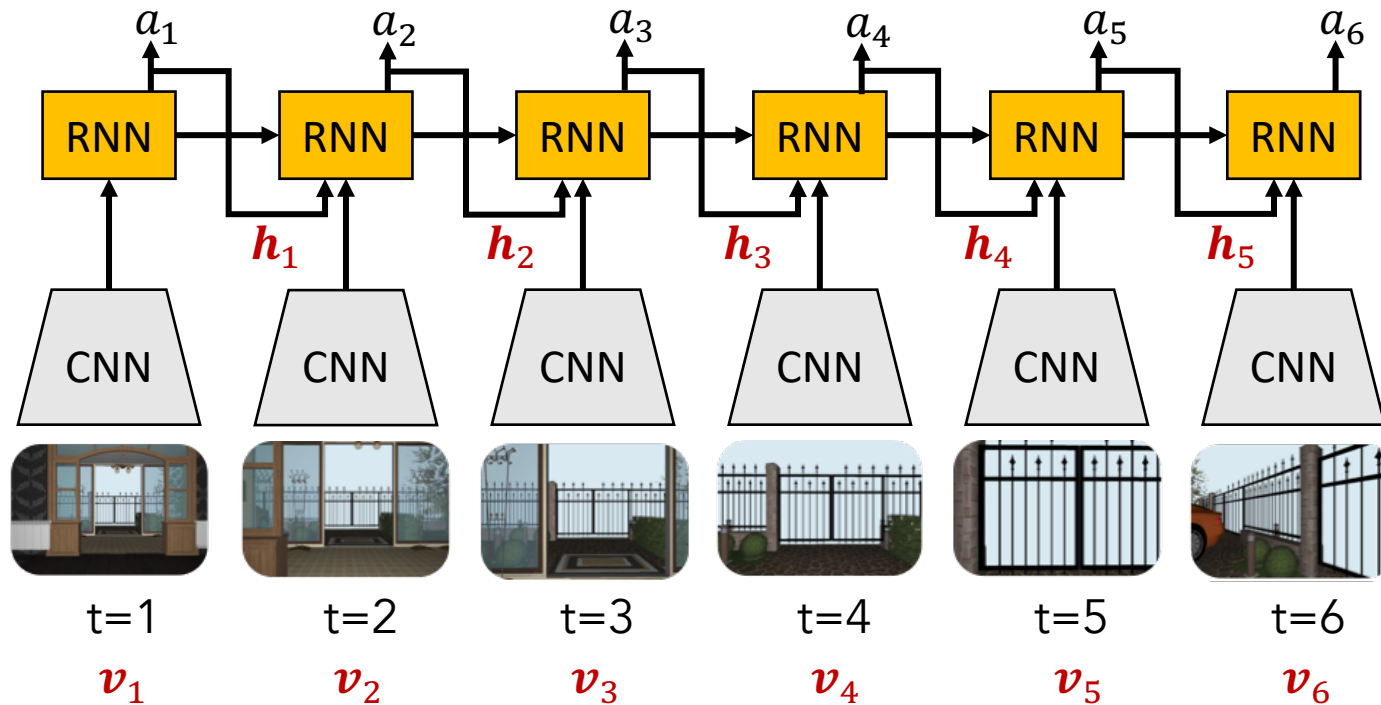


$$\mathcal{C} = \{v_1, \dots, v_{100}\}$$

Attention over Agent Experience

Embodied AI (visual language navigation)

- Attention candidates, \mathcal{C} as agent hidden state or visual vectors



$$\mathcal{C} = \{h_1, \dots, h_5\}$$

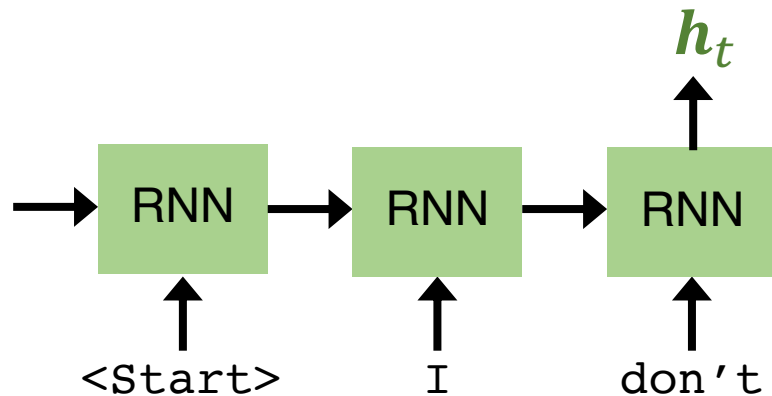
or

$$\mathcal{C} = \{v_1, \dots, v_6\}$$

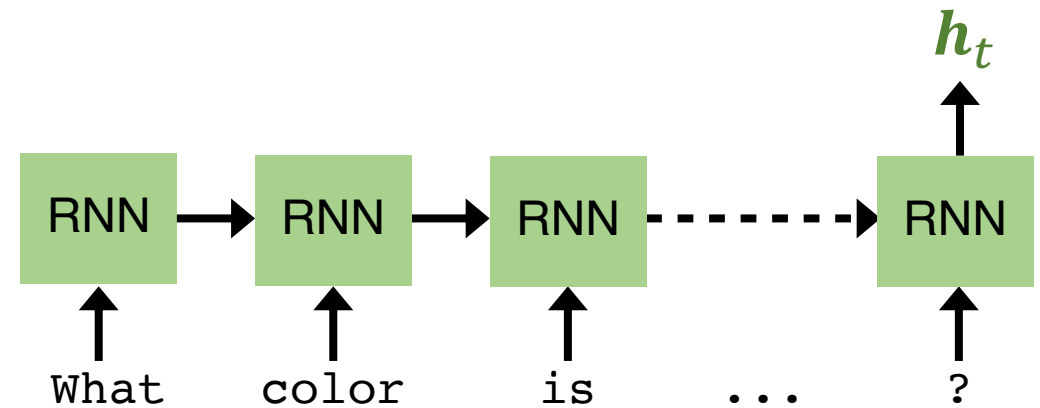
Task context for Attention

- Task context representation, z , is often an RNN encoding

Machine translation / image captioning:
Decoder hidden state

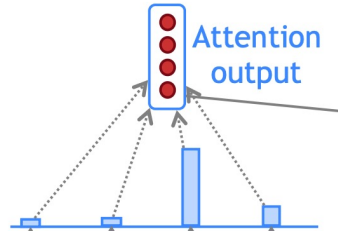


VQA: Question encoding
(final encoder hidden state)



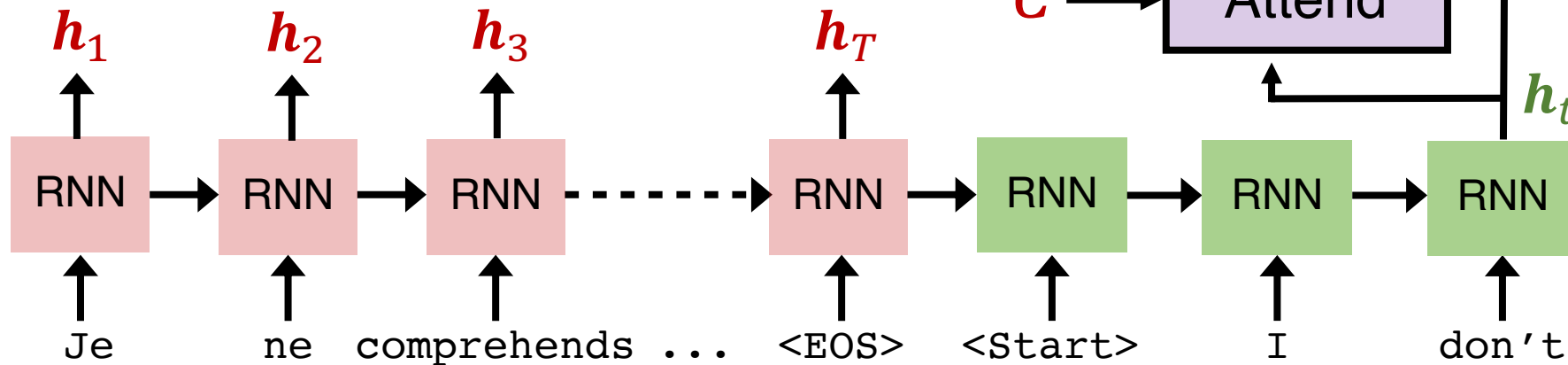
Putting it all together

- Attentive machine translation model



\hat{c}_t is the weighted sum of encoder hidden states

$$C = \{c_1, \dots, c_T\}$$



understand

Predict next word to generate

Softmax

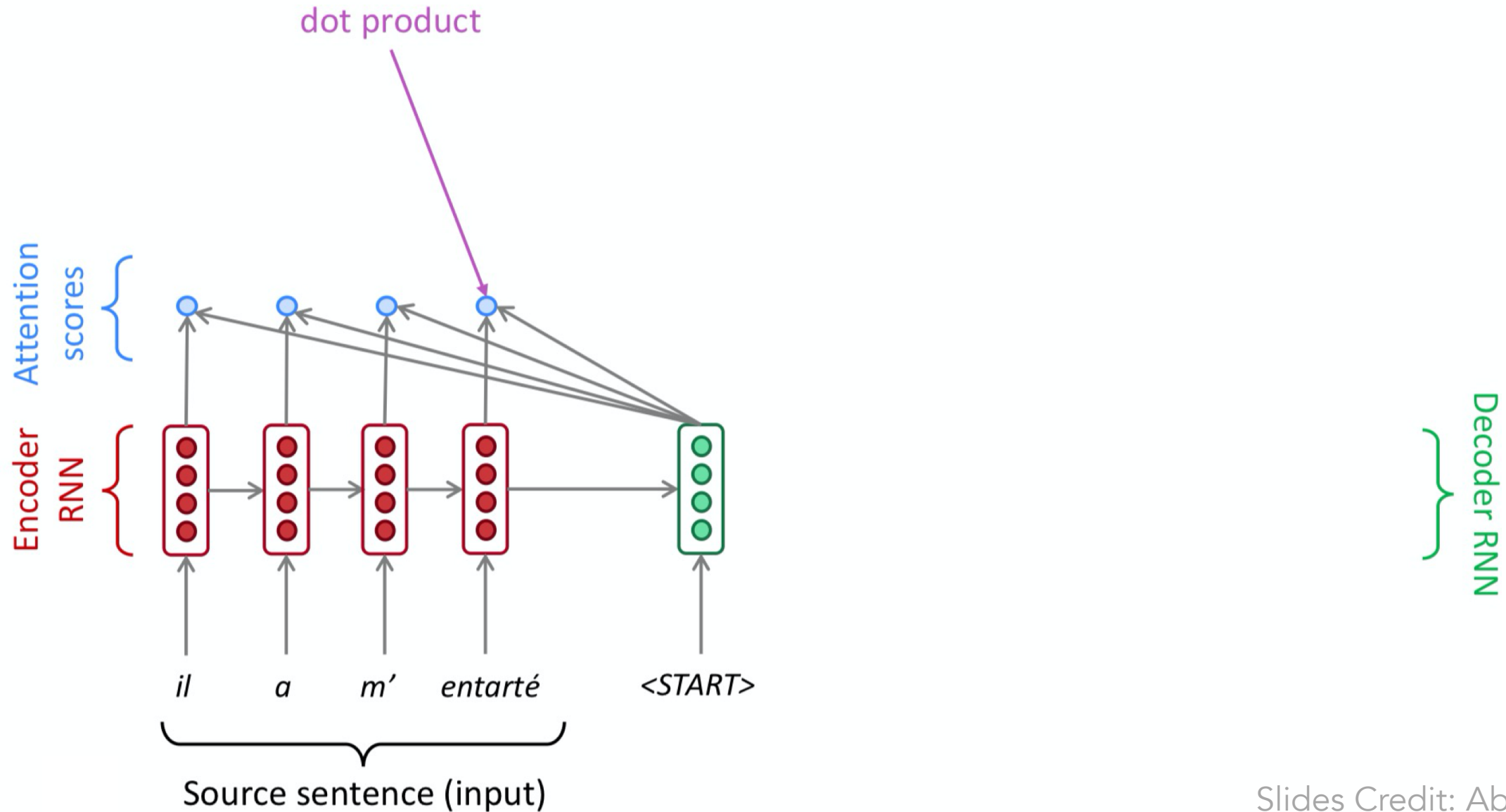
Feedforward Net

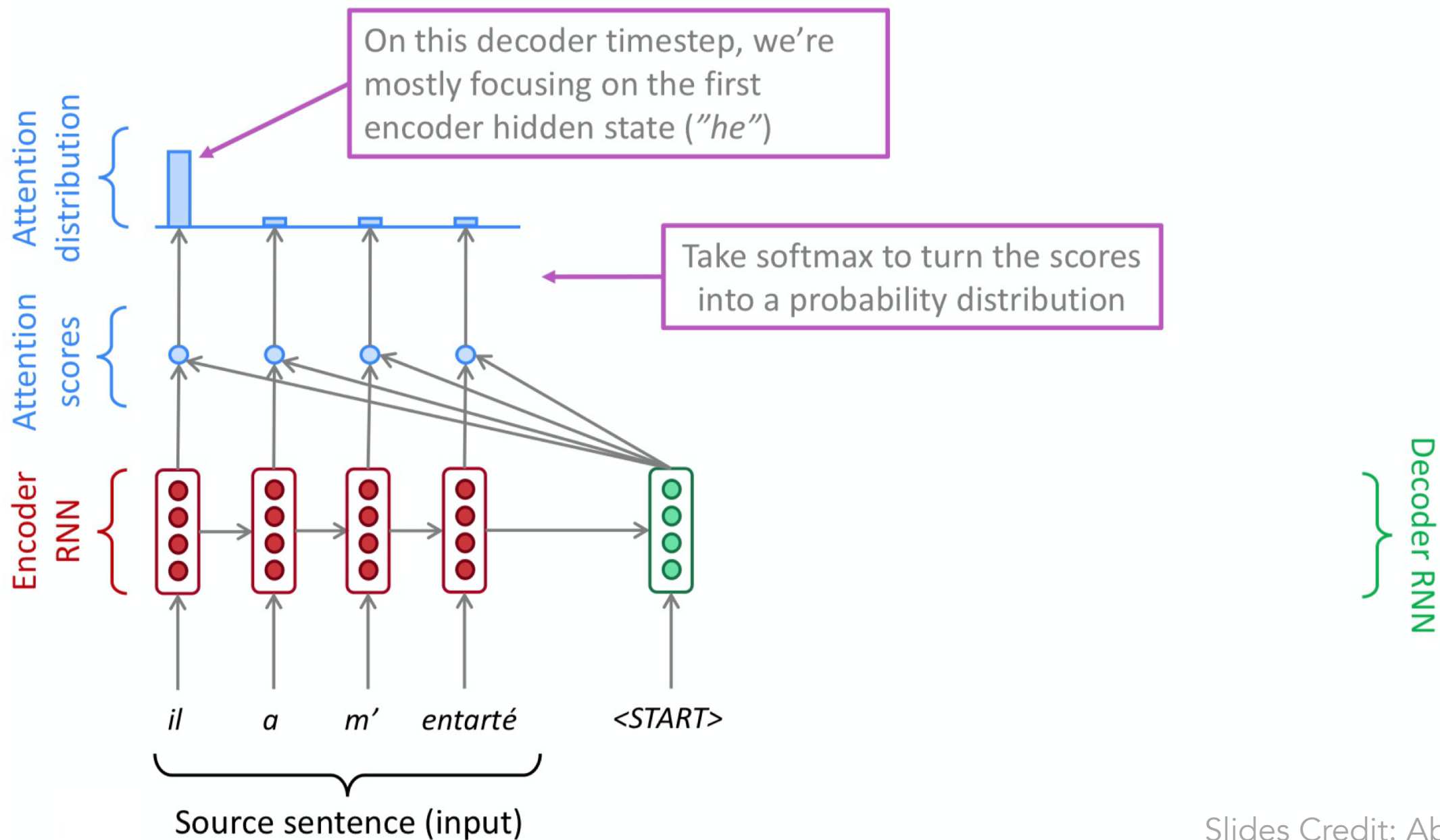
Concatenation

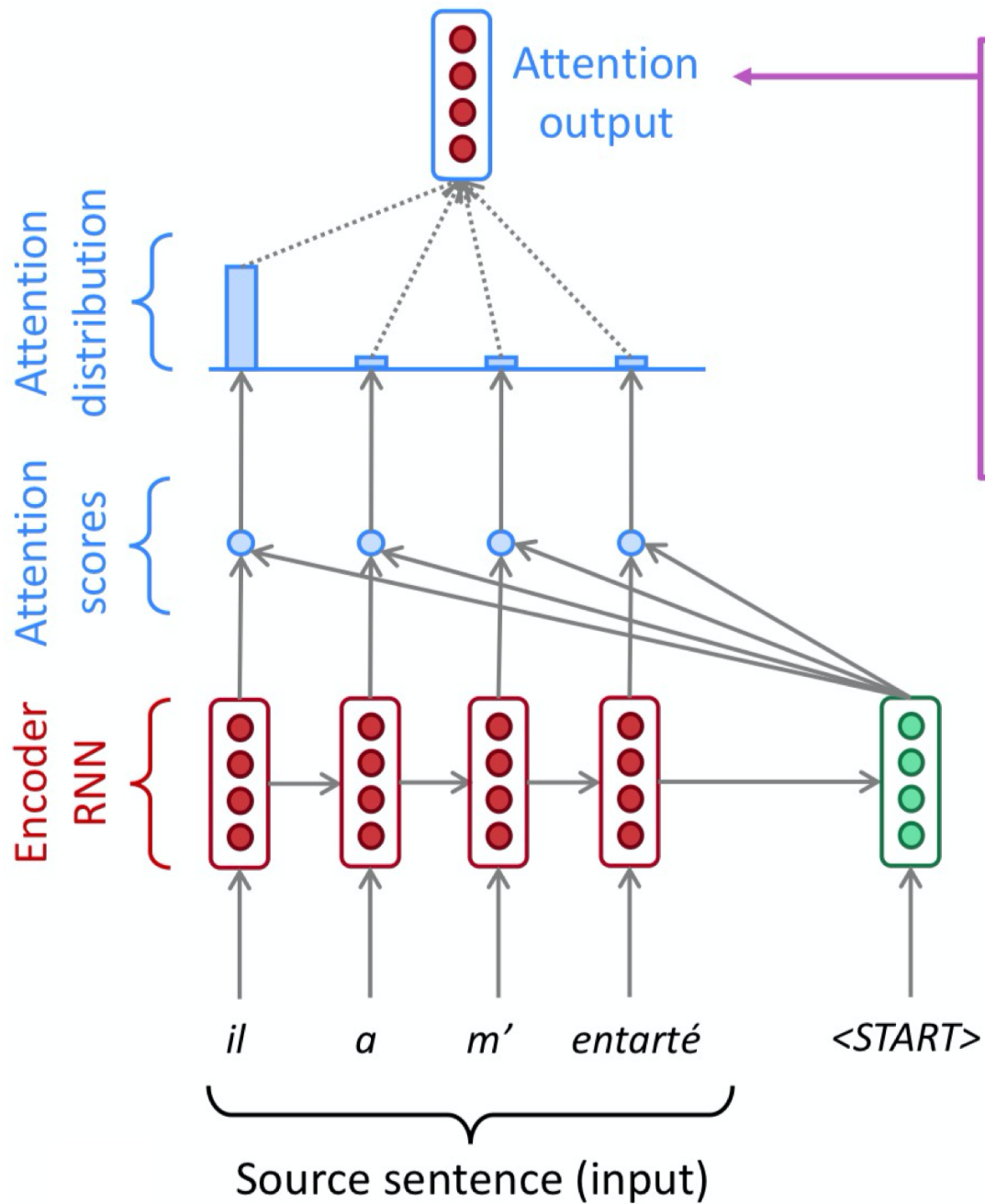
Attend

Interlude: Attention for machine translation

Attention for machine translation

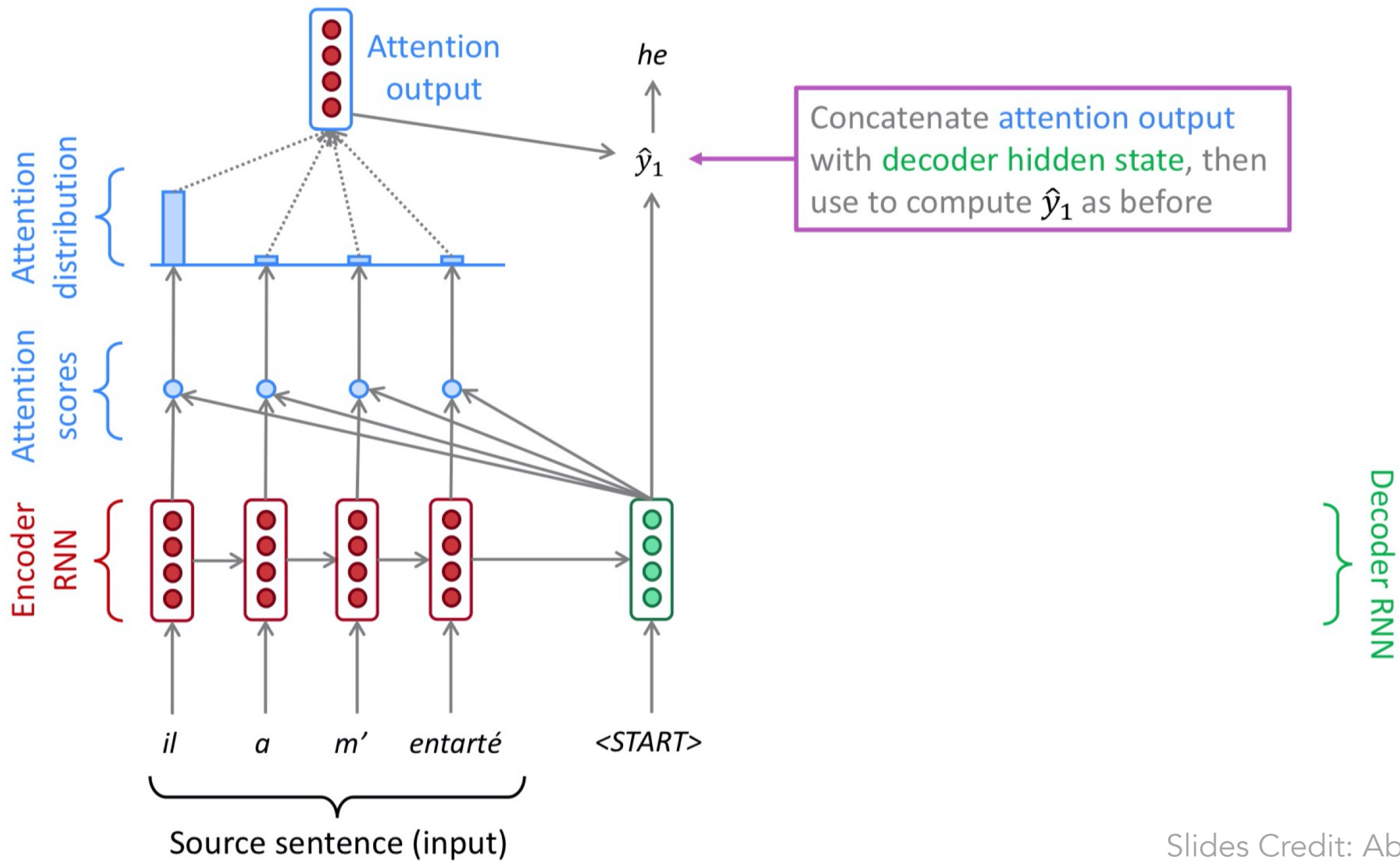


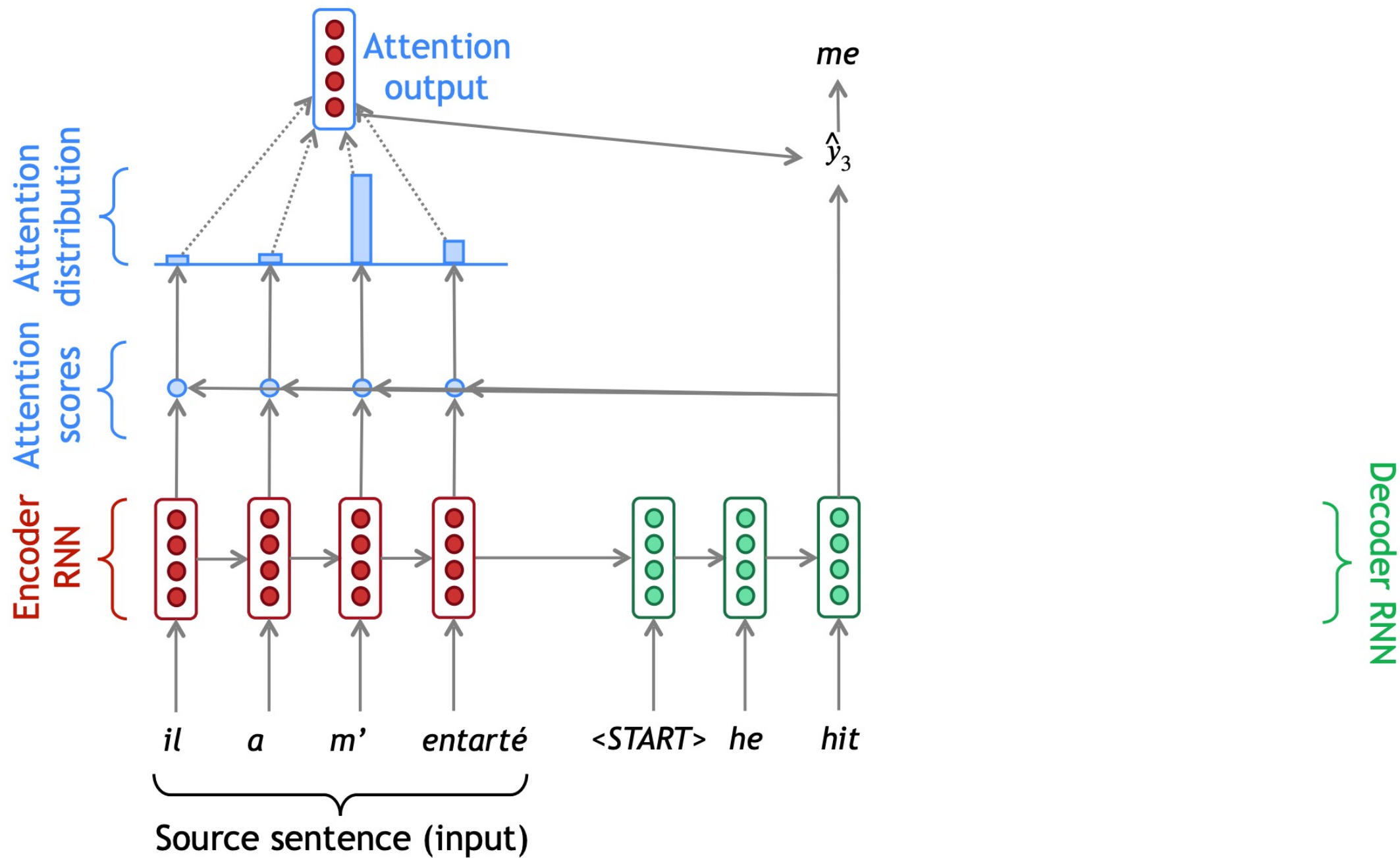


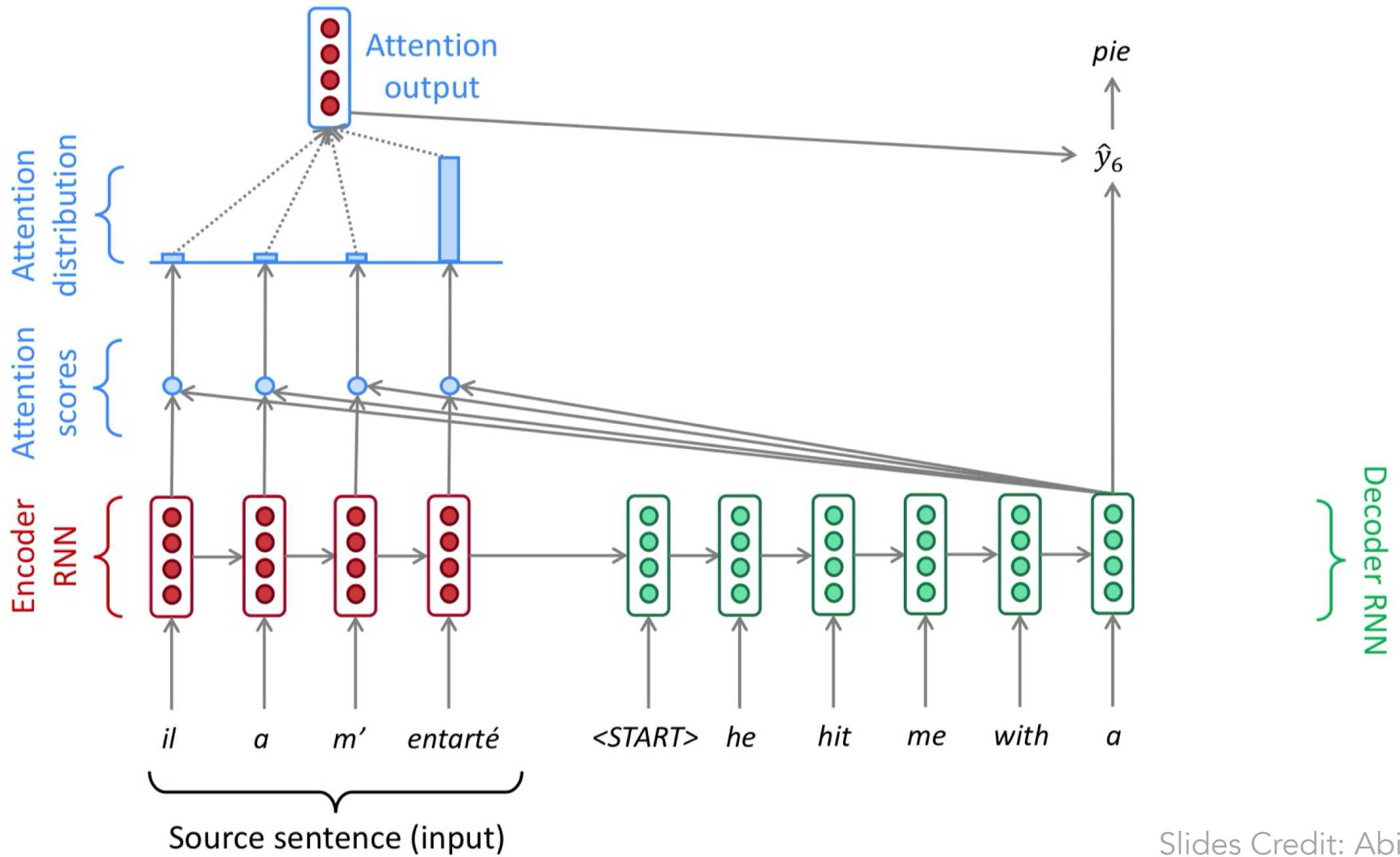


Use the attention distribution to take a **weighted sum** of the **encoder hidden states**.

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

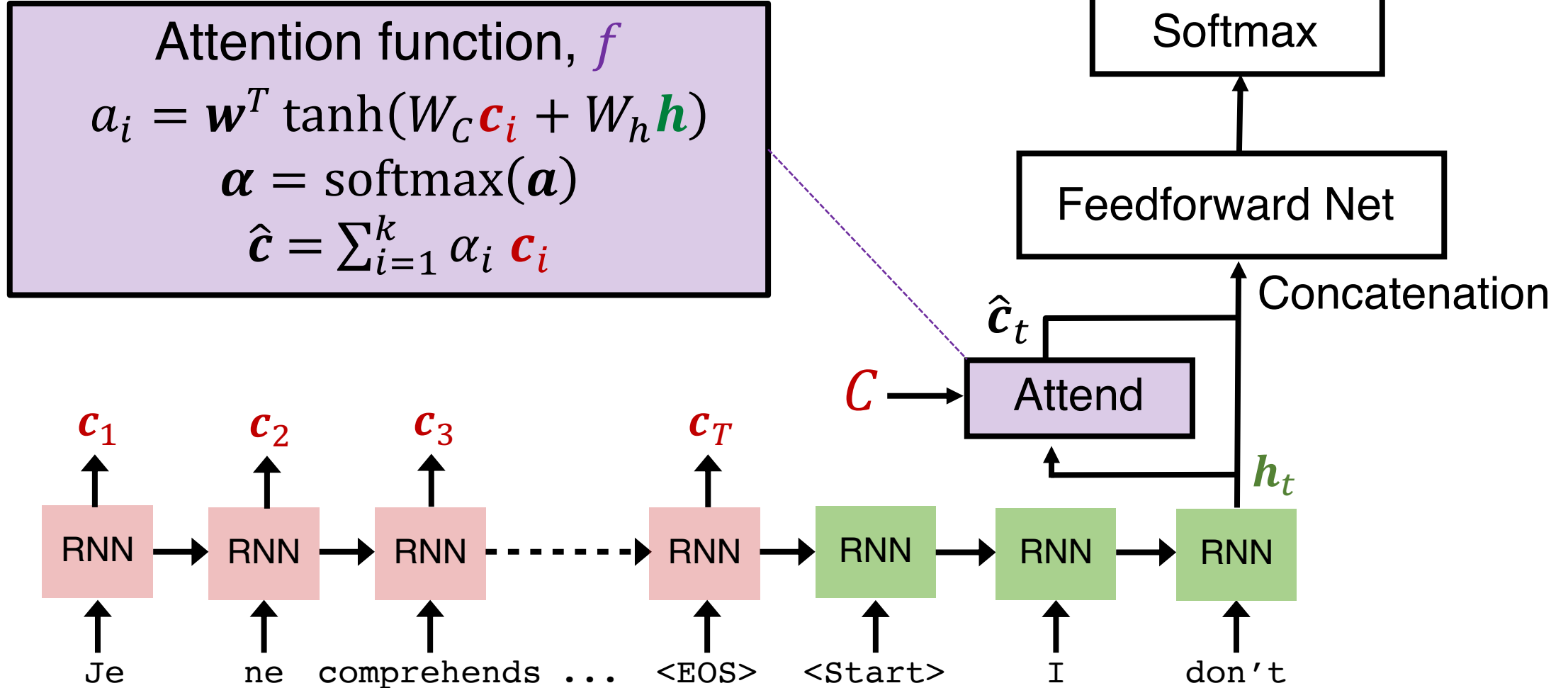






Putting it all together

- Attentive machine translation model



Putting it all together

- Attentive **Visual Question Answering** model

Attention over visual regions!

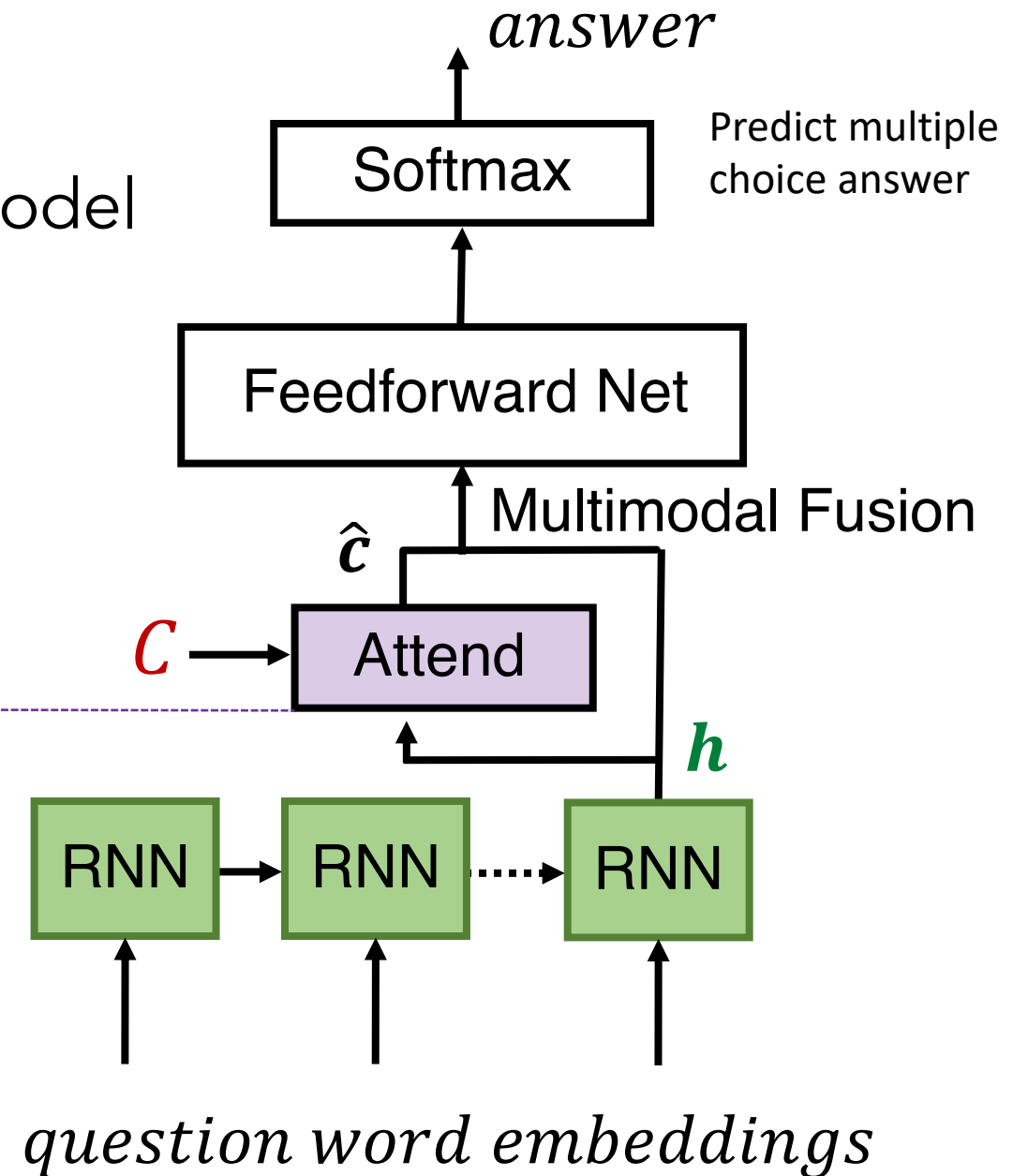
$$C = \{v_1, \dots, v_{100}\}$$

Attention function, f

$$a_i = \mathbf{w}^T \tanh(W_C \mathbf{c}_i + W_h \mathbf{h})$$

$$\boldsymbol{\alpha} = \text{softmax}(\mathbf{a})$$

$$\hat{\mathbf{c}} = \sum_{i=1}^k \alpha_i \mathbf{c}_i$$



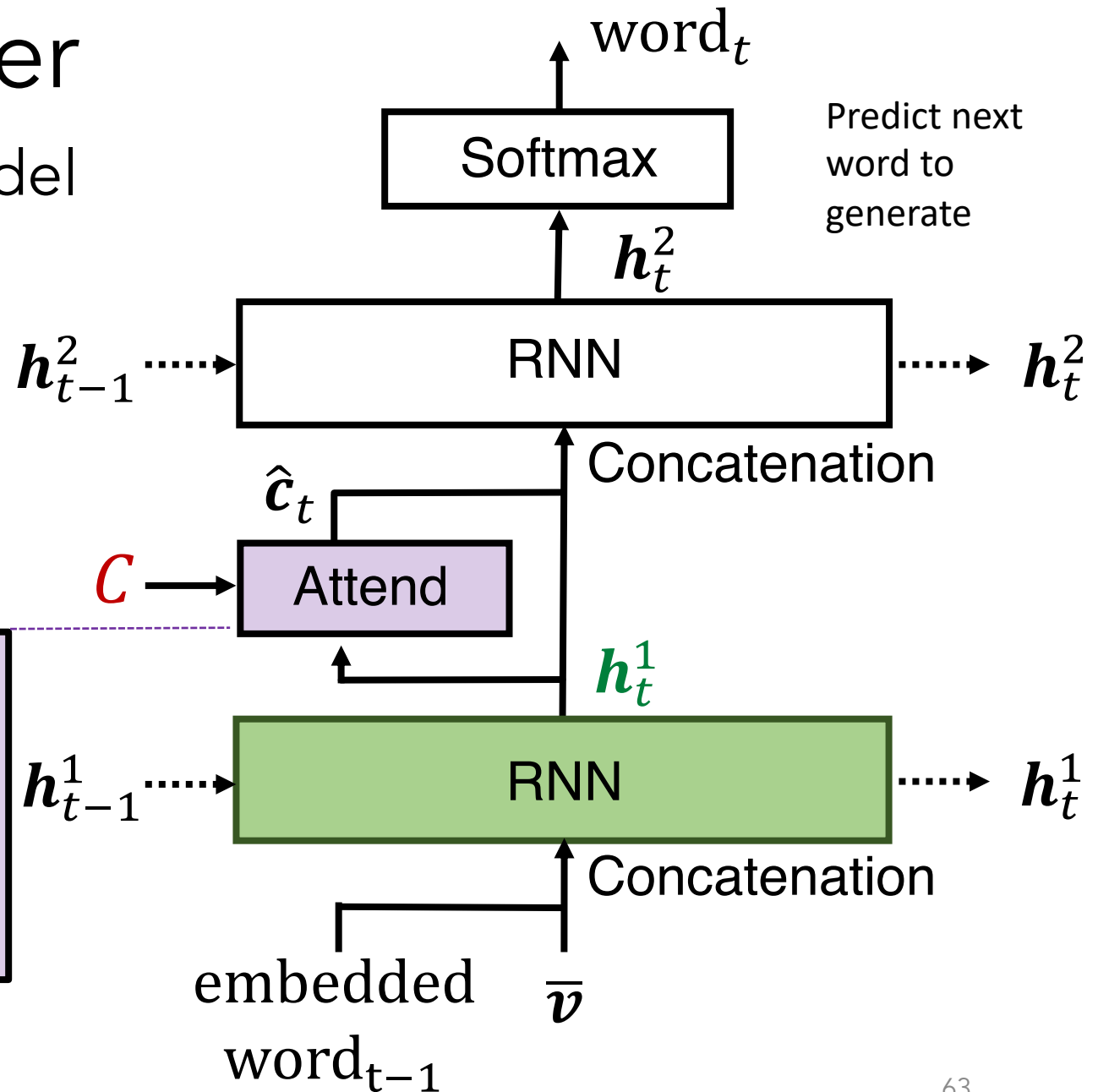
Putting it all together

- Attentive Image Captioning model

Attention over visual regions!
 $C = \{v_1, \dots, v_{100}\}$

Attention function, f

$$a_i = \mathbf{w}^T \tanh(W_C \mathbf{c}_i + W_h \mathbf{h})$$
$$\boldsymbol{\alpha} = \text{softmax}(\mathbf{a})$$
$$\hat{\mathbf{c}} = \sum_{i=1}^k \alpha_i \mathbf{c}_i$$



Types of attention scores

Attention function, f

$$a_i = g(\mathbf{c}_i, \mathbf{z})$$

$$\boldsymbol{\alpha} = \text{softmax}(\mathbf{a})$$

$$\hat{\mathbf{c}} = \sum_{i=1}^k \alpha_i \mathbf{c}_i$$

- Dot-product attention:

$$g(\mathbf{c}_i, \mathbf{z}) = \mathbf{z}^\top \mathbf{c}_i$$

- Scaled dot-product attention:

$$g(\mathbf{c}_i, \mathbf{z}) = \mathbf{z}^\top \mathbf{c}_i / \sqrt{d}$$

- Bilinear / multiplicative attention:

$$g(\mathbf{c}_i, \mathbf{z}) = \mathbf{z}^\top \mathbf{W} \mathbf{c}_i \in \mathbb{R}$$

where \mathbf{W} is a weight matrix

- Additive attention (essentially MLP):

$$g(\mathbf{c}_i, \mathbf{z}) = \mathbf{v}^\top \tanh(\mathbf{W}_1 \mathbf{c}_i + \mathbf{W}_2 \mathbf{z})$$

where $\mathbf{W}_1, \mathbf{W}_2$ are weight matrices and \mathbf{v} is a weight vector

Query-key-value view of attention

Attention function, f

$$a_i = g(\mathbf{c}_i, \mathbf{z})$$

$$\boldsymbol{\alpha} = \text{softmax}(\mathbf{a})$$

$$\hat{\mathbf{c}} = \sum_{i=1}^k \alpha_i \mathbf{c}_i$$



Attention function, f

$$a_i = g(\mathbf{k}_i, \mathbf{q})$$

$$\boldsymbol{\alpha} = \text{softmax}(\mathbf{a})$$

$$\hat{\mathbf{c}} = \sum_{i=1}^k \alpha_i \mathbf{v}_i$$

Projected query, key, value



$$\mathbf{q} = W_Q \mathbf{z}$$

$$\mathbf{k}_i = W_K \mathbf{c}_i$$

$$\mathbf{v}_i = W_V \mathbf{c}_i$$



Matrix form

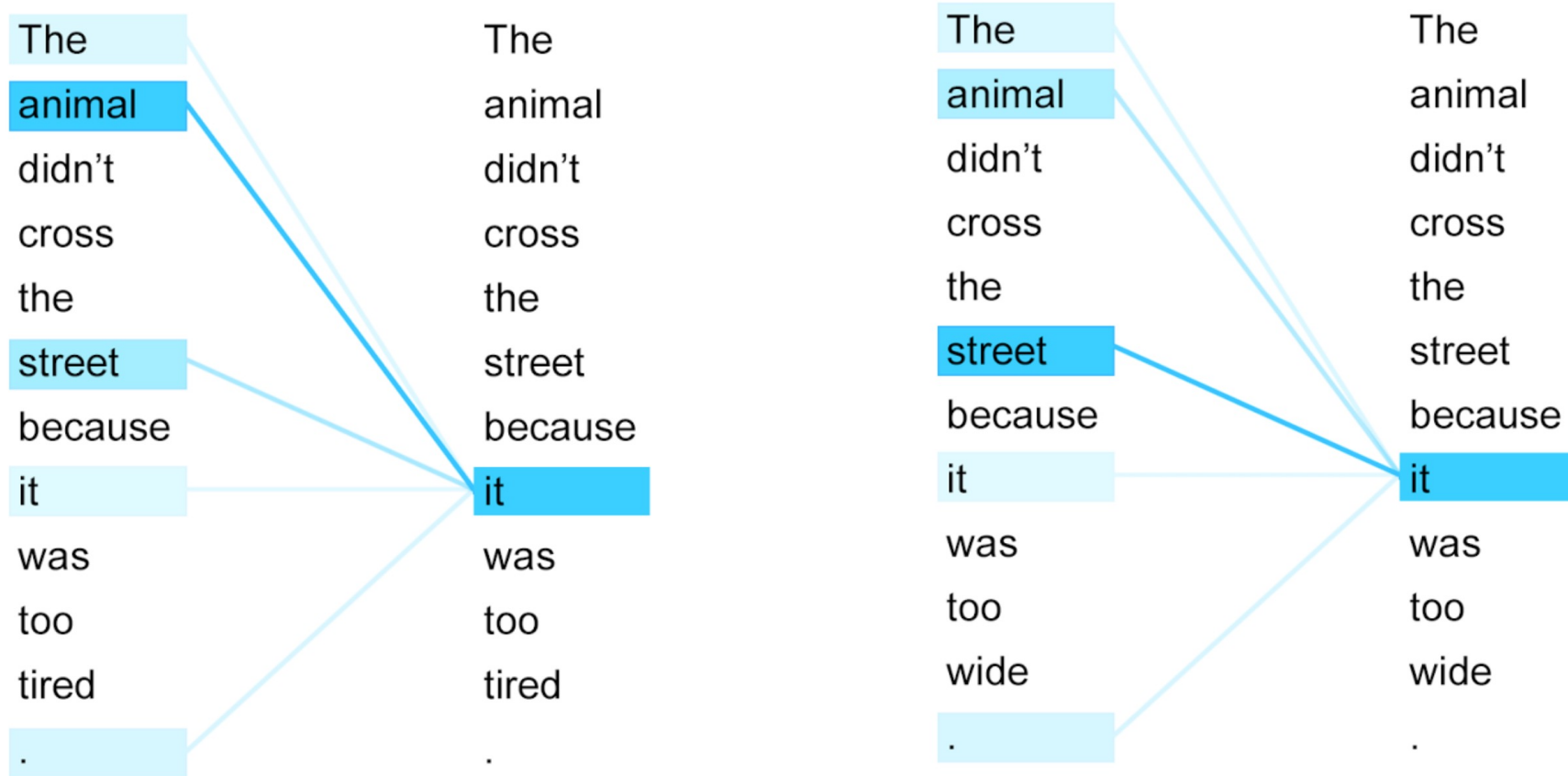
$$\mathbf{q} = W_Q \mathbf{z}$$

$$K = W_K C^T$$

$$V = W_V C^T$$

Self-attention

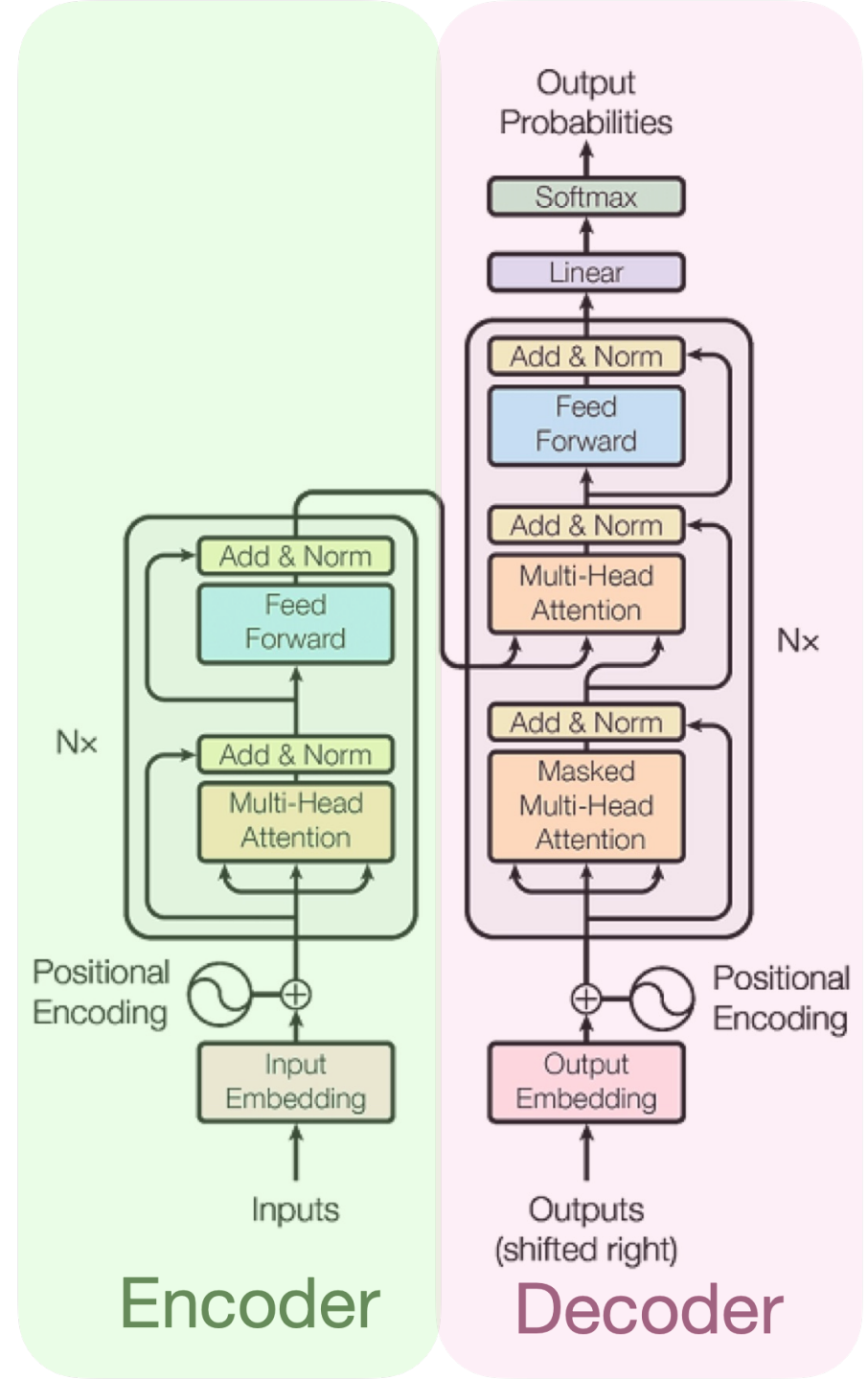
- Attention (correlation) with different parts of itself



- Transformers: modules with scaled dot-product self-attention

Transformers

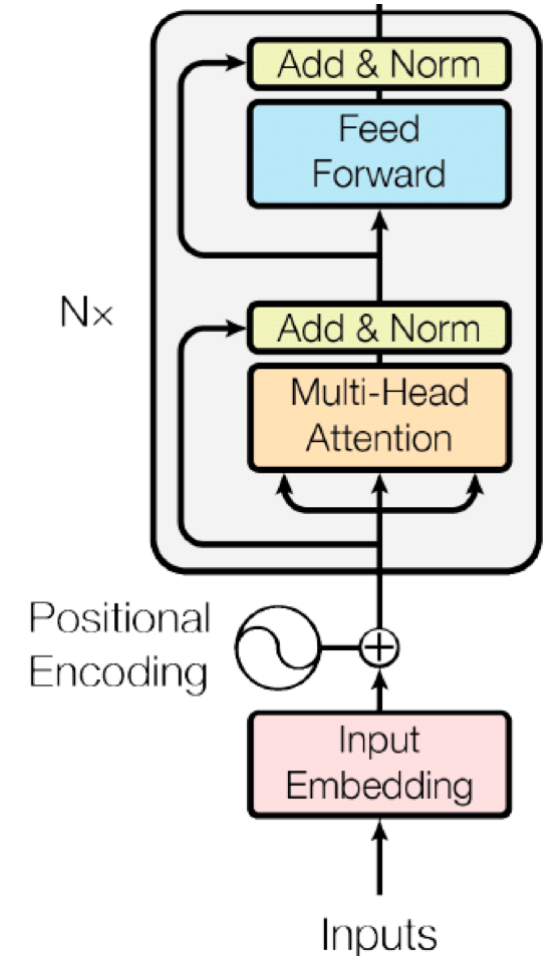
- NIPS'17: Attention is All You Need
- Originally proposed for NMT (encoder-decoder framework)
- Key idea: **Multi-head self-attention**
- No recurrence structure so training can be parallelized



Modelling Sequences -- Transformers

- Each Transformer block has two sub-layers
 - Multi-head attention
 - 2-layer feedforward NN (with ReLU)
- Each sublayer has a residual connection and a layer normalization
 - $\text{LayerNorm}(x + \text{SubLayer}(x))$
- Input layer has a positional encoding

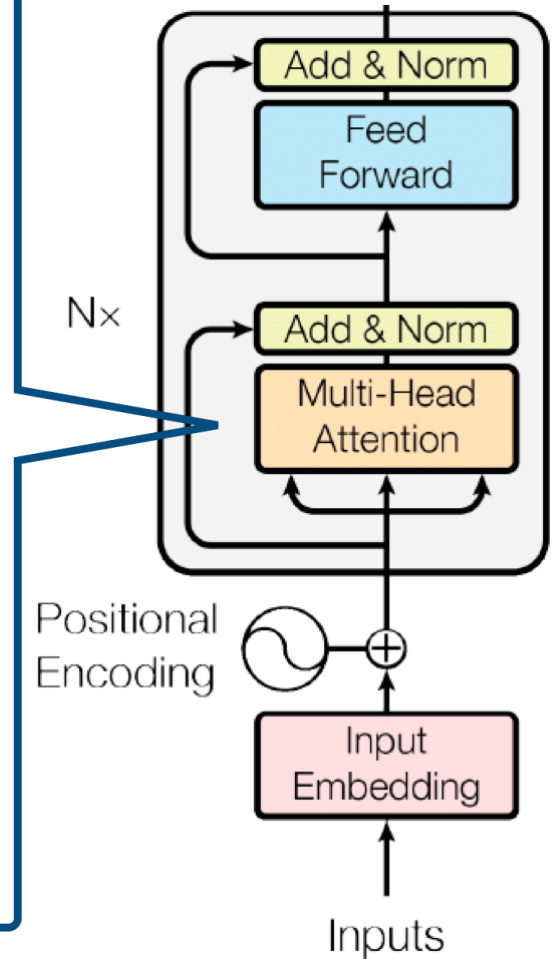
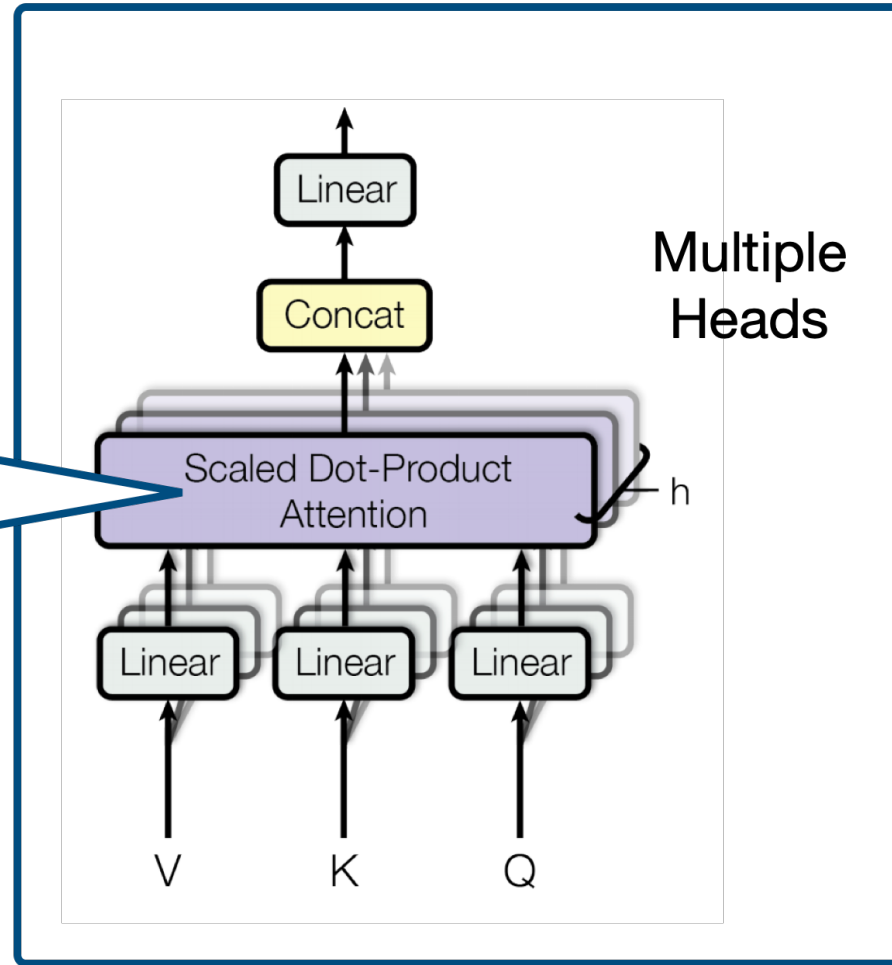
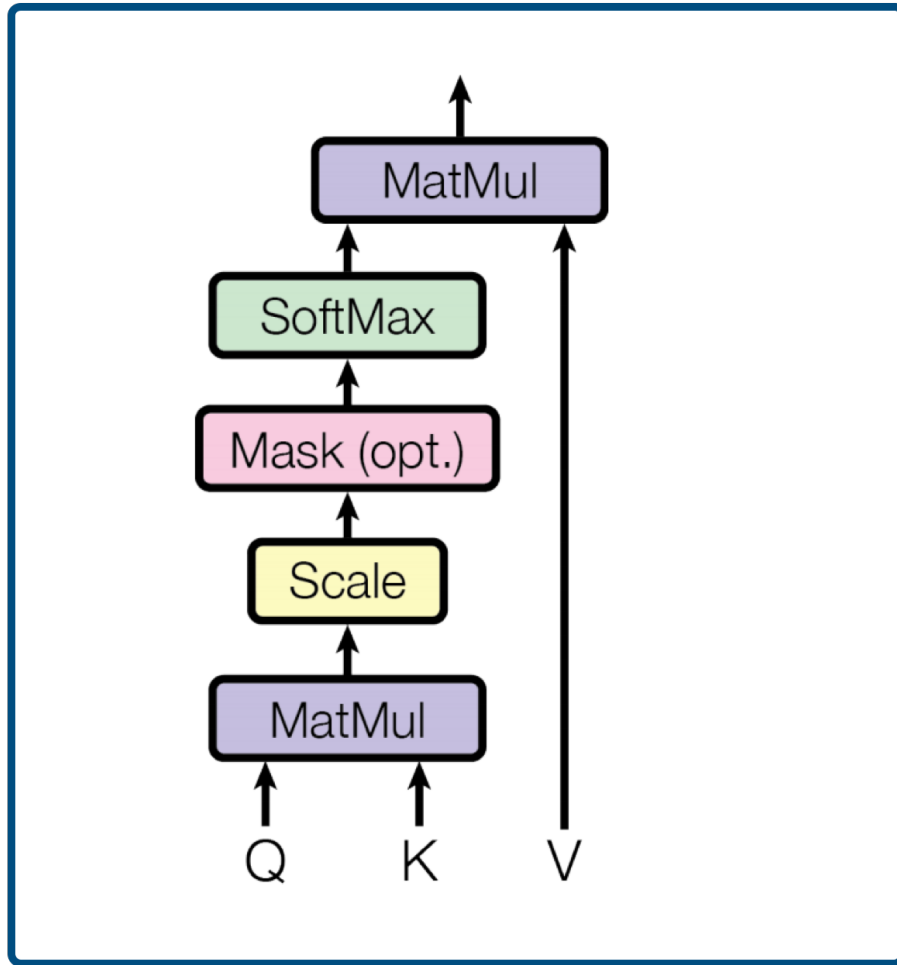
Helps the training process!



Modelling Sequences -- Transformers

Scaled Dot-Product Attention

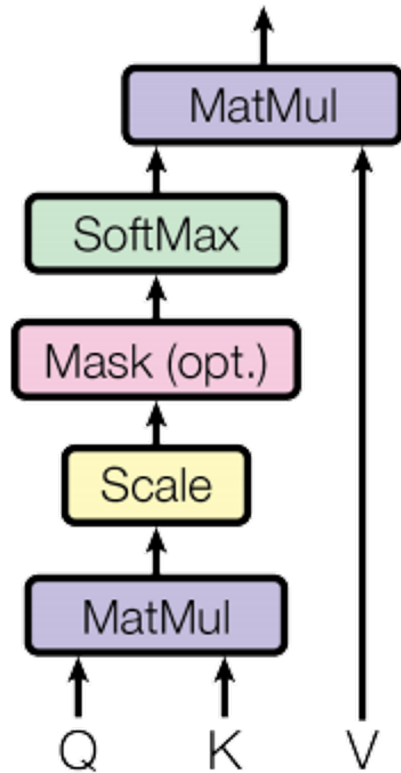
self-attention



Scaled Dot Product Attention

Efficient, stable training

Scaled Dot-Product Attention



Let $X \in \mathbb{R}^{M \times d_x}$ be a matrix of task context vectors to attend to
Let $C \in \mathbb{R}^{N \times d_c}$ be a matrix of input vectors to attend over

SDPAttention(X, C):

$$Q = W_Q X^T \quad W_Q \in \mathbb{R}^{d_h \times d_x}$$

$$K = W_K C^T \quad W_K \in \mathbb{R}^{d_h \times d_c}$$

$$V = W_V C^T \quad W_V \in \mathbb{R}^{d_v \times d_c}$$

$$\text{Return } \hat{V} = \text{softmax} \left(\frac{Q^T K}{\sqrt{d_h}} \right) V$$

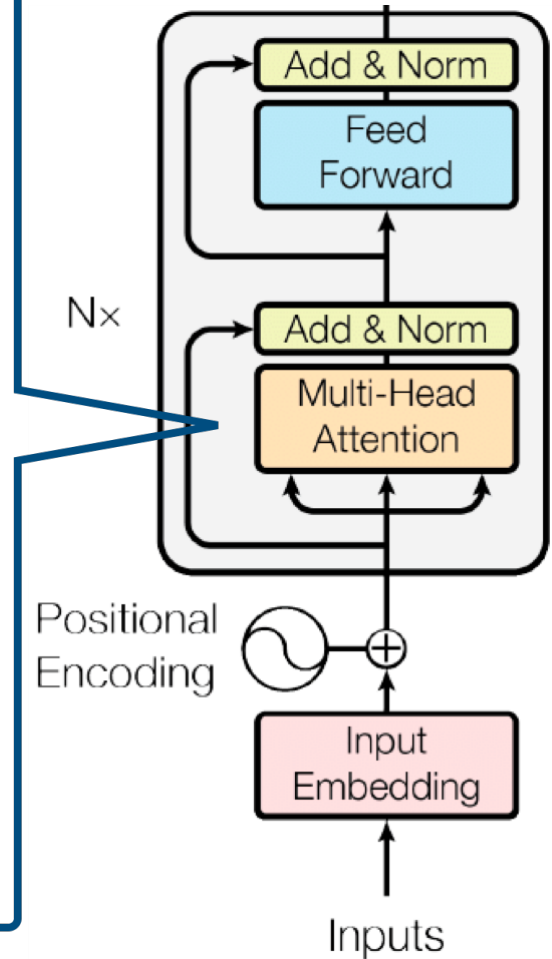
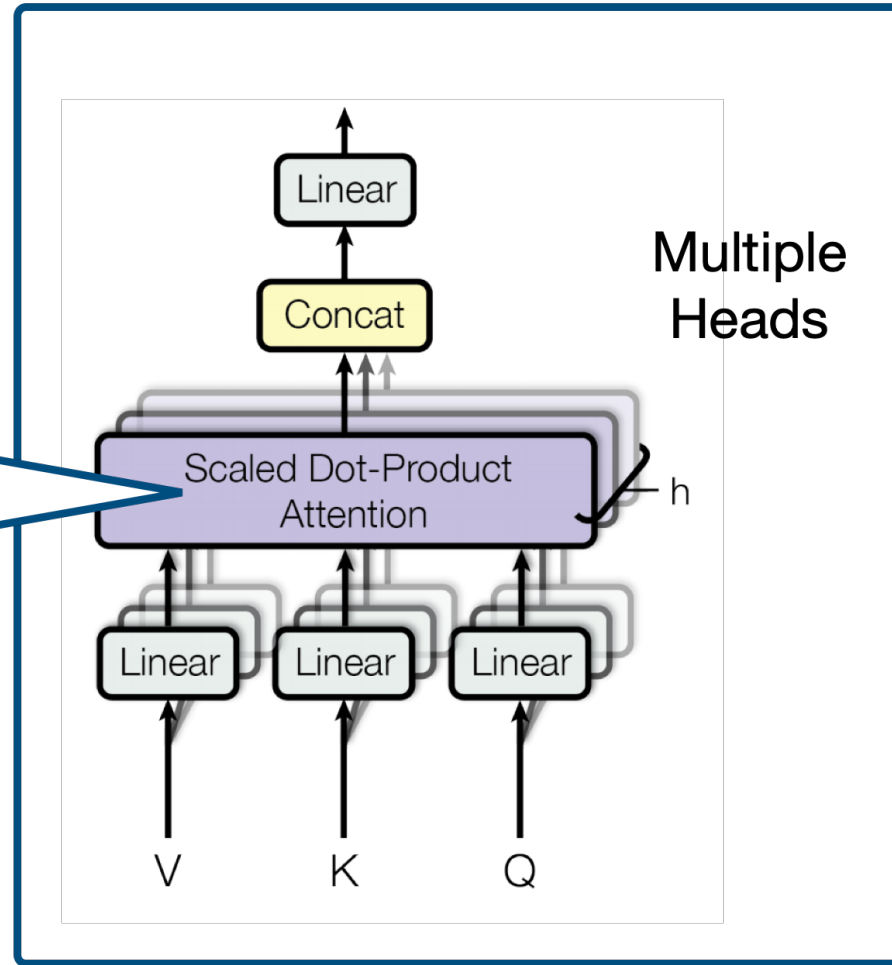
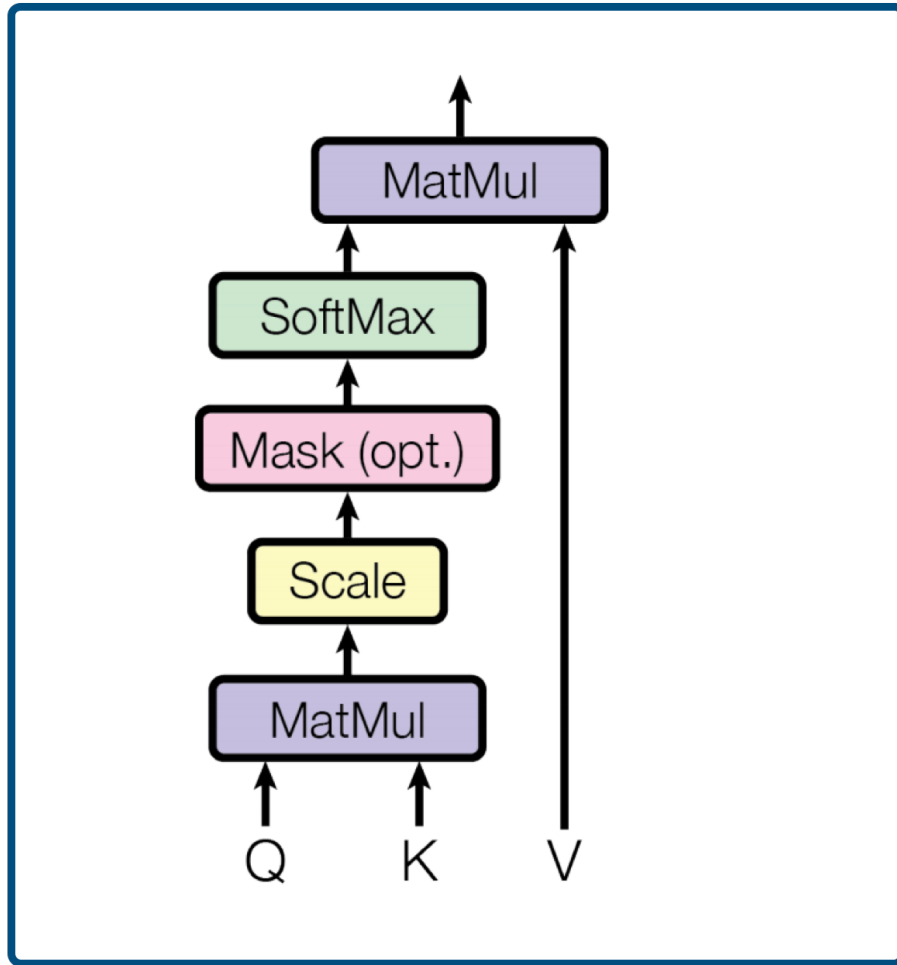
$\hat{V} \in \mathbb{R}^{M \times d_v}$ be a matrix of attended values

Attention Is All You Need <https://arxiv.org/pdf/1706.03762.pdf>

Modelling Sequences -- Transformers

Scaled Dot-Product Attention

self-attention $SDPAttention(C, C)$:



Multi-head attention

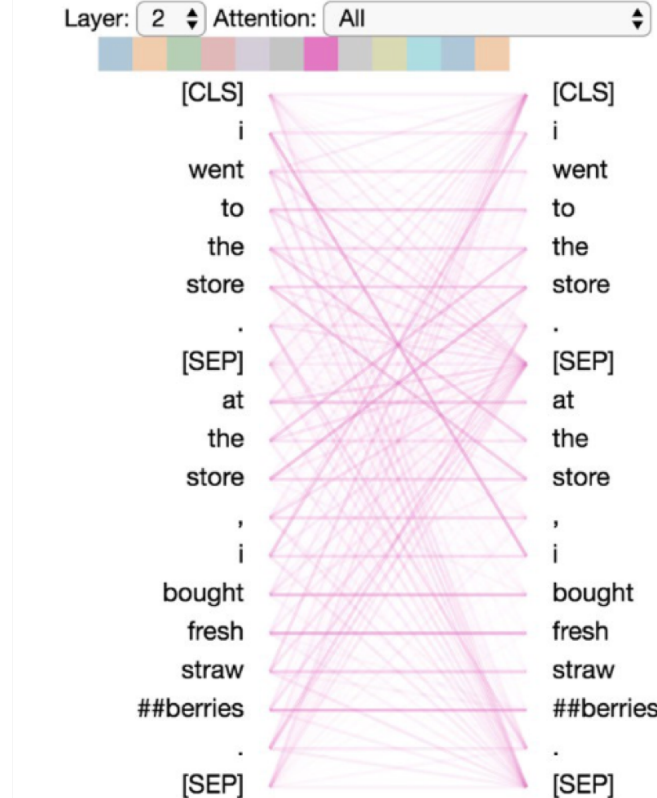
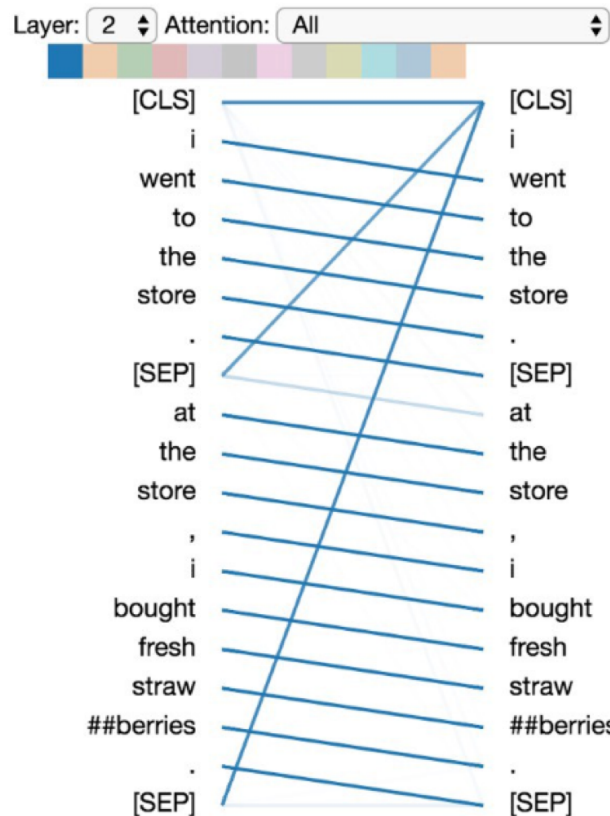
One head is not expressive enough. Let's have multiple heads!

$$A(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W_O$$

$$\text{head}_i = A(W_{Q_i} X^T, W_{K_i} X^T, W_{V_i} X^T)$$

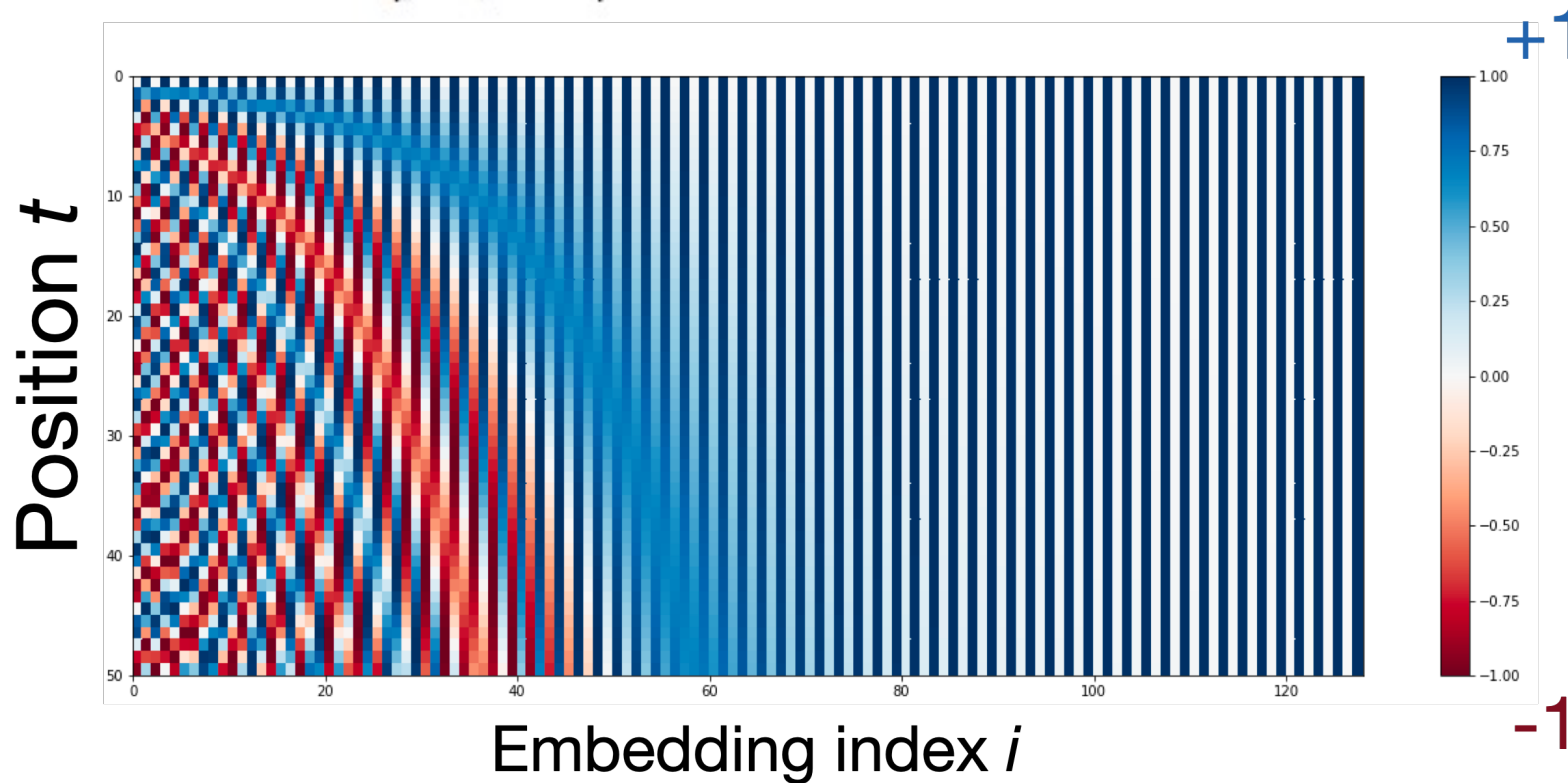
In practice, $h = 8$,

$$d = d_{out}/h, W_O \in \mathbb{R}^{d_{out} \times d_{out}}$$

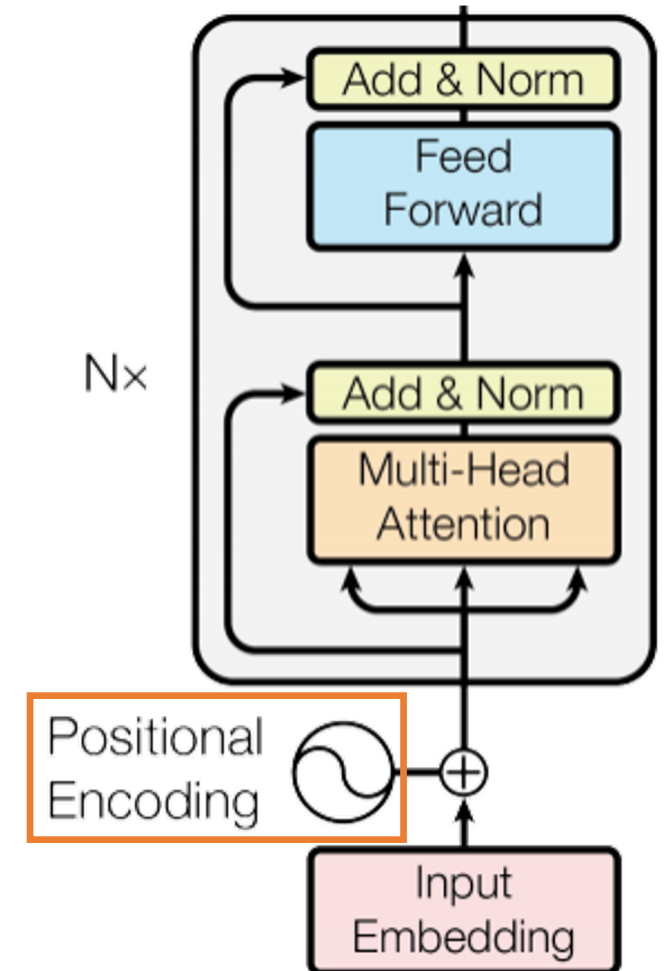


Transformers: Encoding position

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$



***SDPAttention*(Y, Y):**



Attention Is All You Need <https://arxiv.org/pdf/1706.03762.pdf>

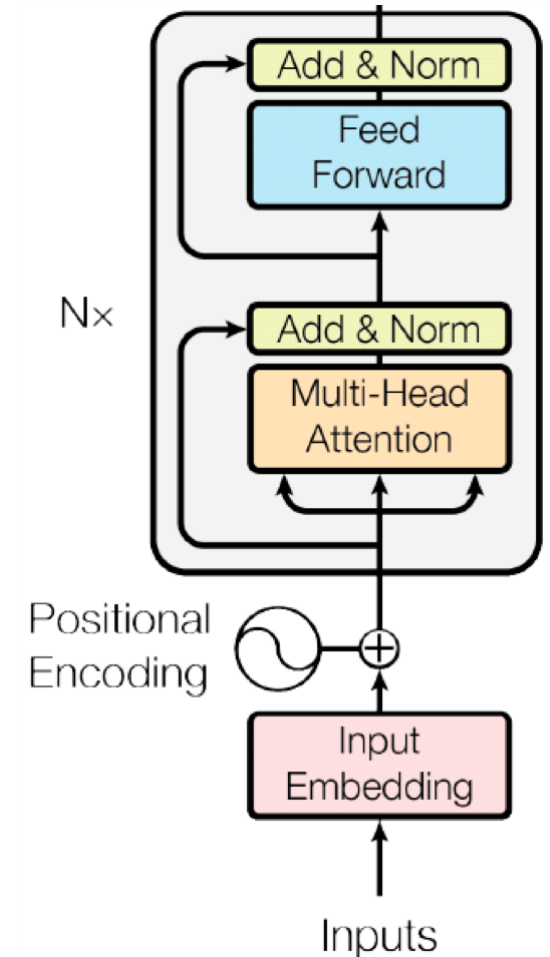
Modelling Sequences -- Transformers

- Each Transformer block has two sub-layers
 - Multi-head attention
 - 2-layer feedforward NN (with ReLU)
Provides non-linearity

- Each sublayer has a residual connection and a layer normalization
 $\text{LayerNorm}(x + \text{SubLayer}(x))$

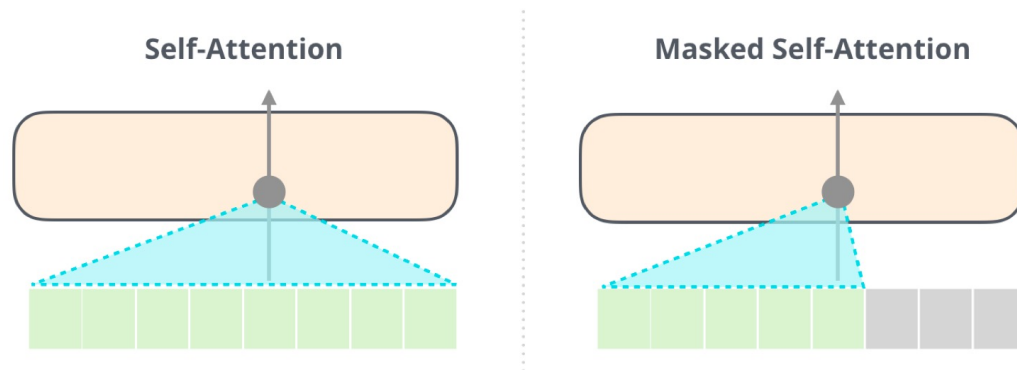
Helps the training process!

- Input layer has a positional encoding

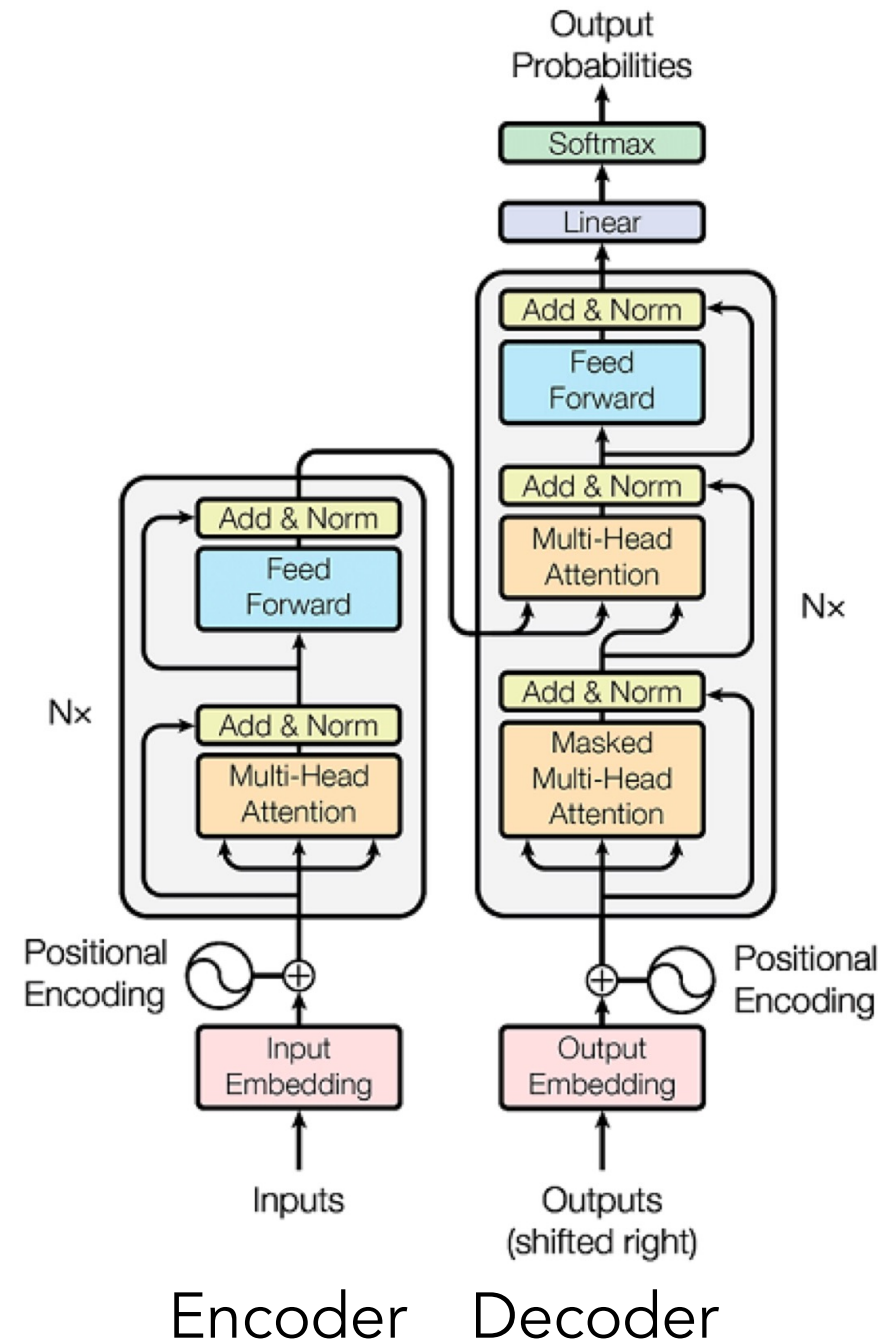


Transformers

- Encoder: Multi-headed self-attention
- Decoder
 - Masked self-attention

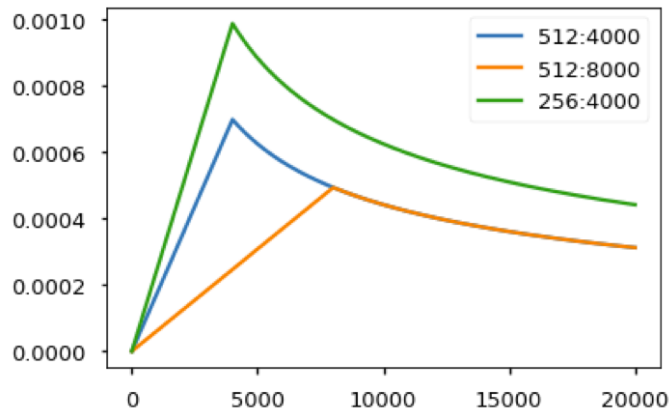


- Cross attention
 - queries: previous decoder layer
 - keys/values: output of encoder
- Autoregressive decoding



Transformers

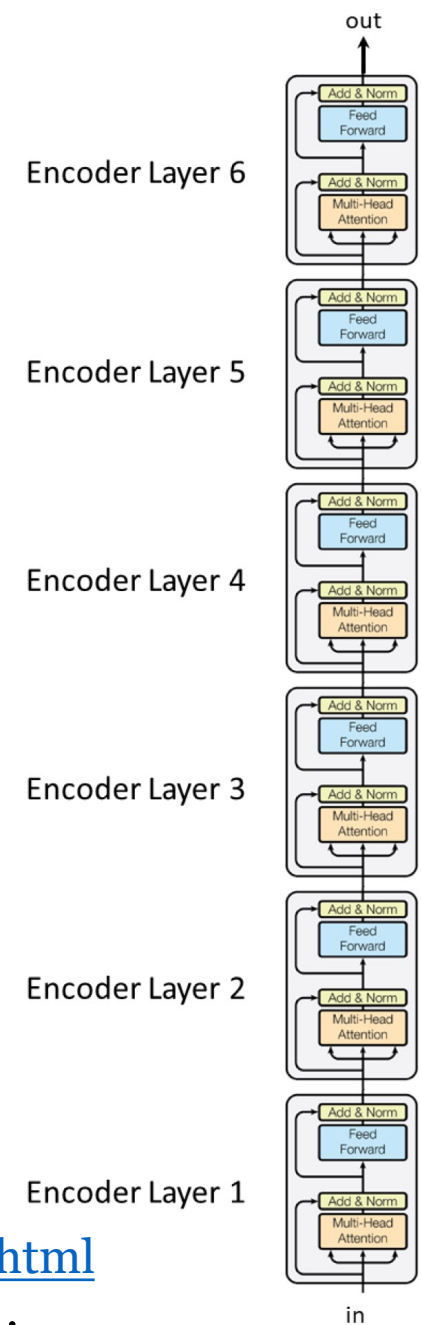
- Stacked into multi-layers
- For language, input embedding is subwords
 - Byte-pair encoding (BPE) / Word pieces
- Other training details:
 - Learning rate with warmup and decay



- Label smoothing: one-hot vector + noise

The Annotated Transformer <http://nlp.seas.harvard.edu/2018/04/03/attention.html>











A Jupyter notebook which explains how Transformer works line by line in PyTorch!



Transformers are used for everything!

10 Novel Applications using Transformers [DL]

Transformers have had a lot of success in training neural language models. In the past few weeks, we've seen several trending papers with code applying Transformers to new types of task:

-  **Transformer for Image Synthesis** - [Esser et al. \(2020\)](#)
-  **Transformer for Multi-Object Tracking** - [Sun et al. \(2020\)](#)
-  **Transformer for Music Generation** - [Hsiao et al. \(2021\)](#)
-  **Transformer for Dance Generation with Music** - [Huang et al. \(2021\)](#)
-  **Transformer for 3D Object Detection** - [Bhattacharyya et al. \(2021\)](#)
-  **Transformer for Point-Cloud Processing** - [Guo et al. \(2020\)](#)
-  **Transformer for Time-Series Forecasting** - [Lim et al. \(2020\)](#)
-  **Transformer for Vision-Language Modeling** - [Zhang et al. \(2021\)](#)
-  **Transformer for Lane Shape Prediction** - [Liu et al. \(2020\)](#)
-  **Transformer for End-to-End Object Detection** - [Zhu et al. \(2021\)](#)

Next time

- Multimodal representations