

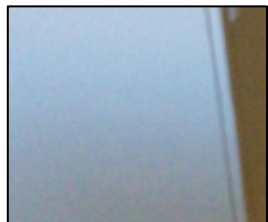
# Attention and Transformer

Some slides from James Hays, Justin Johnson

# Scene Understanding



# Context is important



# Language understanding

... serve ...

# Language understanding

... great **serve** from Djokovic ...



# Language understanding

... be right back after I **serve** these salads ...





**Brendan Dolan-Gavitt**

@moyix

The latest generation of adversarial image attacks is, uh, somewhat simpler to carry out [openai.com/blog/multimoda...](https://openai.com/blog/multimodal-adversarial-attacks)

### Attacks in the wild

We refer to these attacks as *typographic attacks*. We believe attacks such as those described above are far from simply an academic concern. By exploiting the model's ability to read text robustly, we find that even *photographs of hand-written text* can often fool the model. Like the Adversarial Patch,<sup>22</sup> this attack works in the wild; but unlike such attacks, it requires no more technology than pen and paper.

Attack text label iPod ▾



|              |       |
|--------------|-------|
| Granny Smith | 85.6% |
| iPod         | 0.4%  |
| library      | 0.0%  |
| pizza        | 0.0%  |
| toaster      | 0.0%  |
| dough        | 0.1%  |



|              |       |
|--------------|-------|
| Granny Smith | 0.1%  |
| iPod         | 99.7% |
| library      | 0.0%  |
| pizza        | 0.0%  |
| toaster      | 0.0%  |
| dough        | 0.0%  |

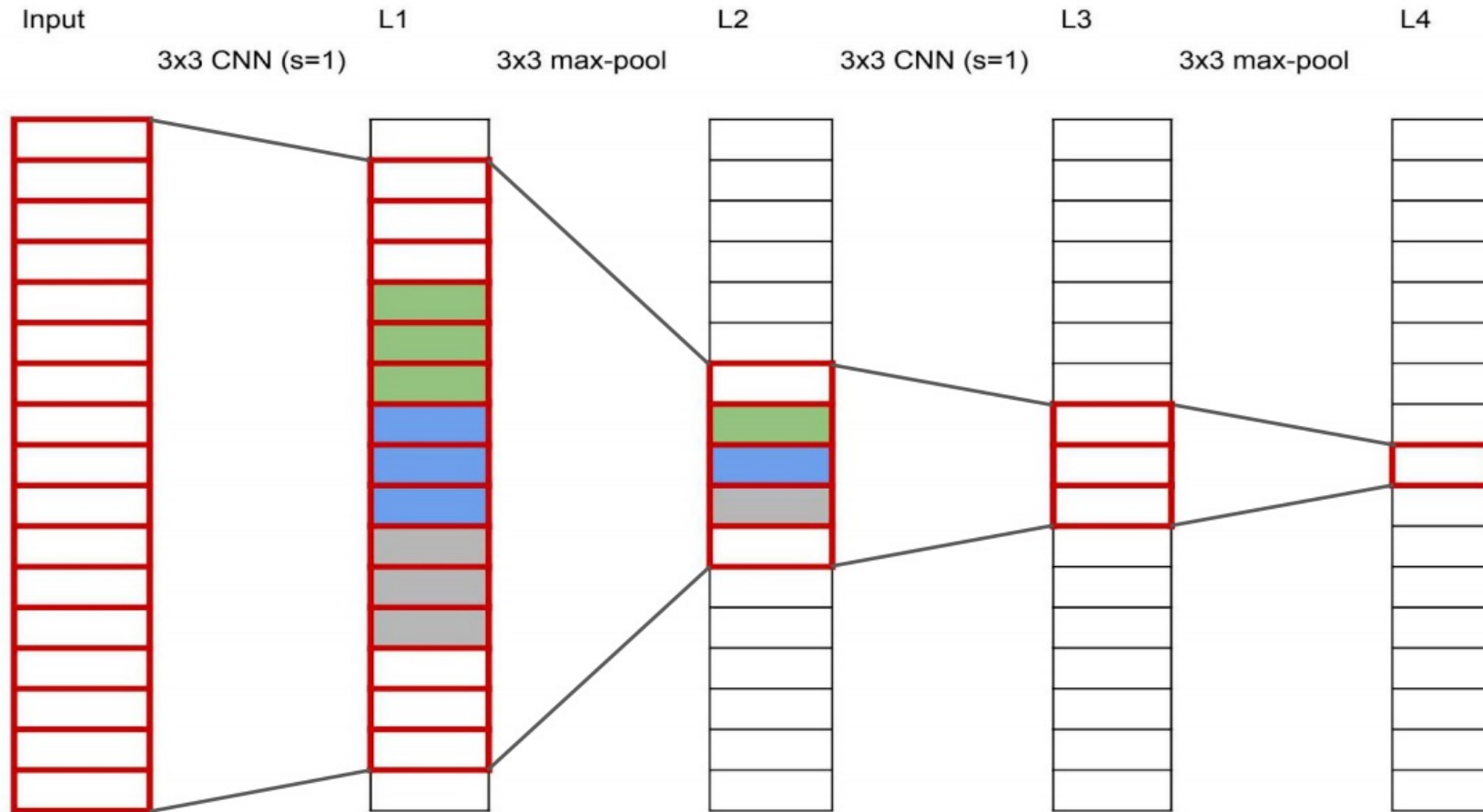
When we put a label saying "iPod" on this Granny Smith apple, the model erroneously classifies it as an iPod in the zero-shot setting.

So how do we fix these problems?

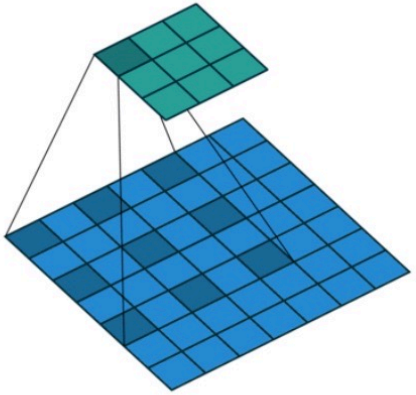




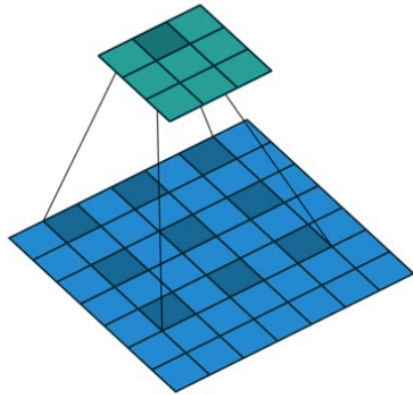
# Receptive field



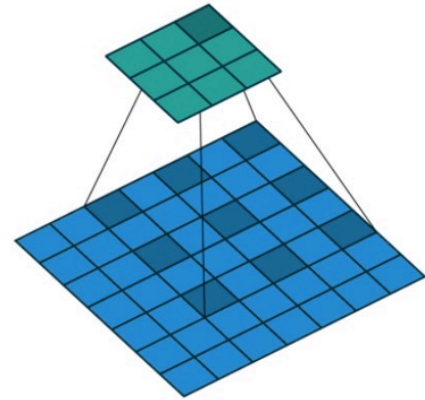
# Dilated Convolution



No padding, no stride, dilation



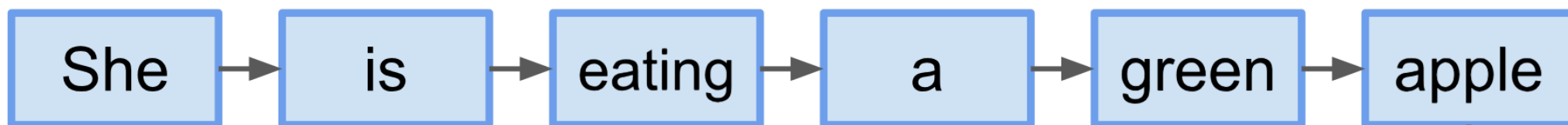
No padding, no stride, dilation



No padding, no stride, dilation

## Sequence 2 Sequence models in language

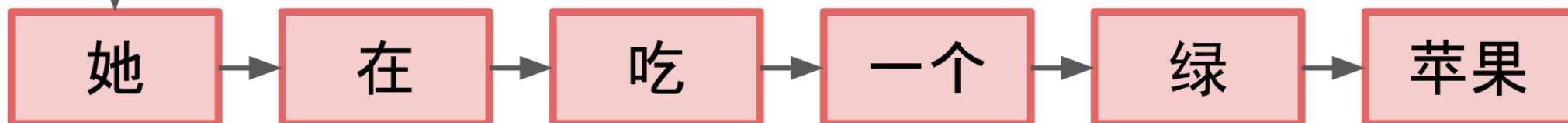
**Encoder**



Context vector (length: 5)

[0.1, -0.2, 0.8, 1.5, -0.3]

**Decoder**



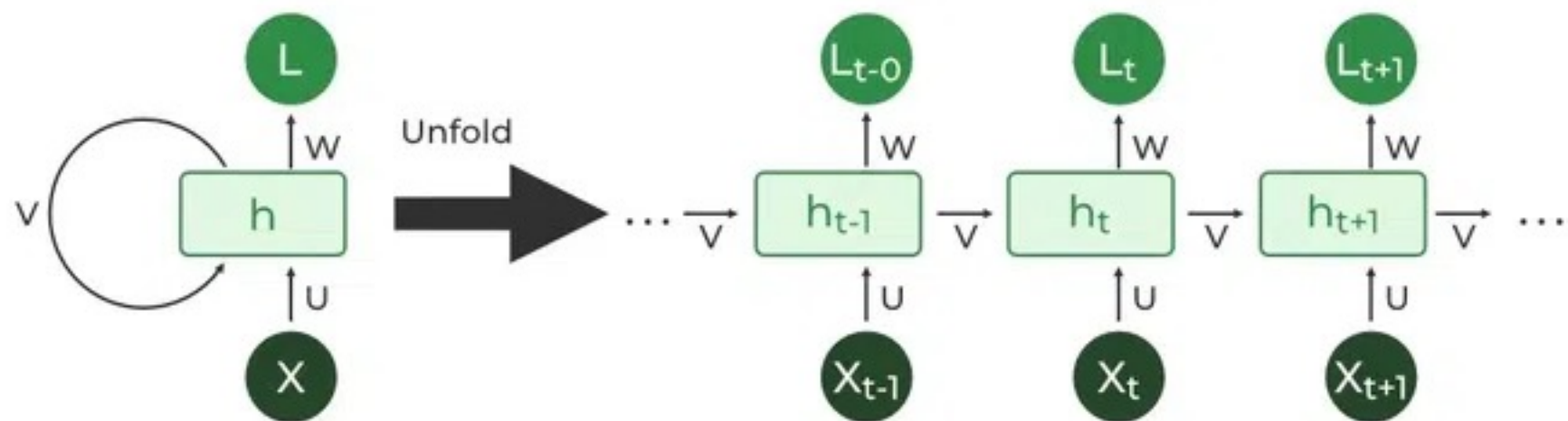
**Problem: Input sequence bottlenecked through fixed sized context vector.**

How to model global context and large receptive field/sequence?

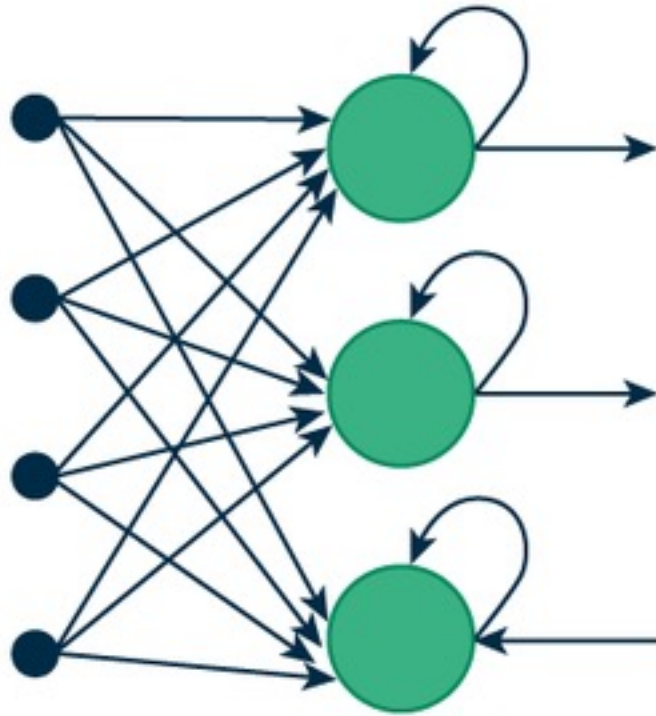
# Outline

- Recurrent Neural Network (RNN)
- Attention and Transformer
- Vision Transformer for Image Classification

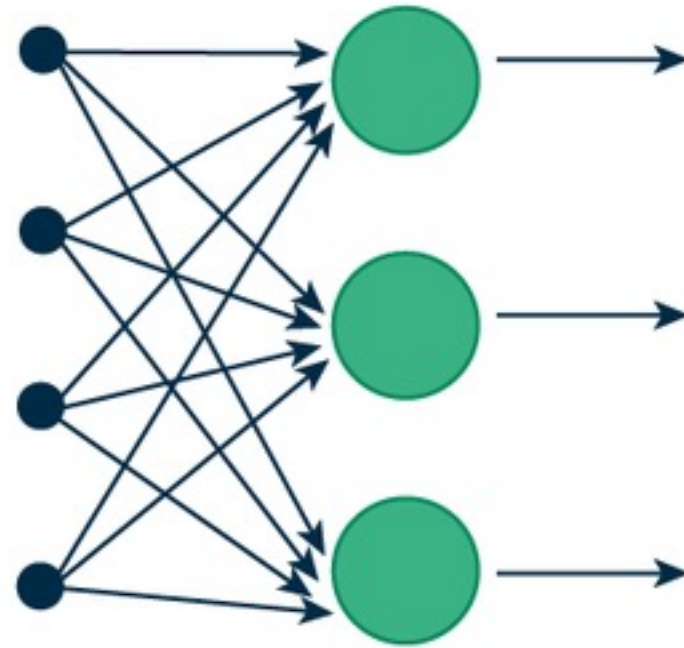
# Recurrent Neural Network



# Recurrent Neural Network



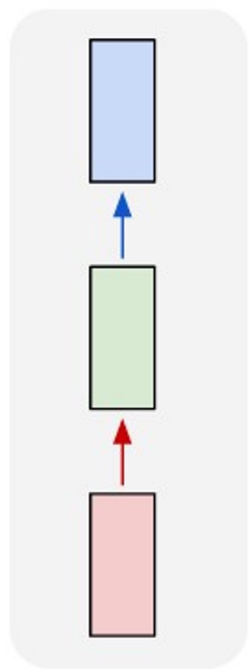
(a) Recurrent Neural Network



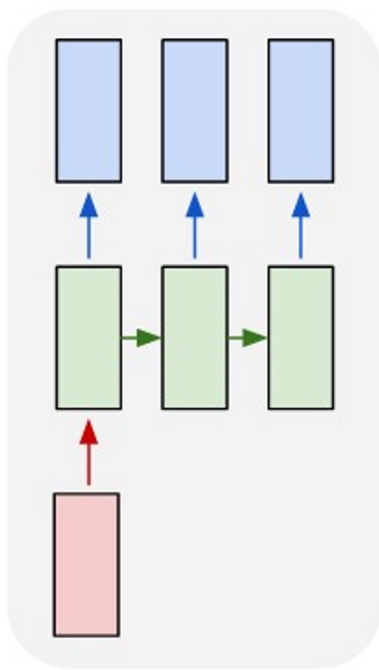
(b) Feed-Forward Neural Network

# Recurrent Neural Networks

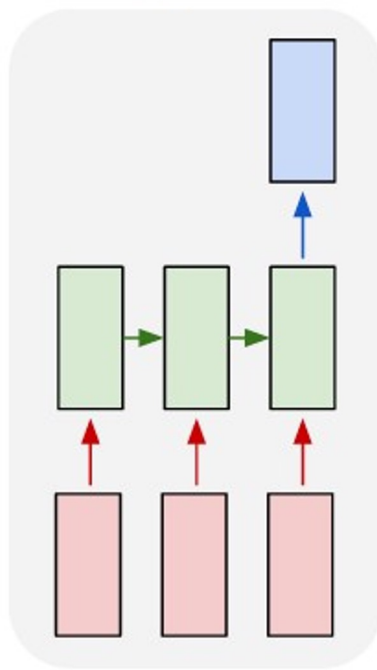
one to one



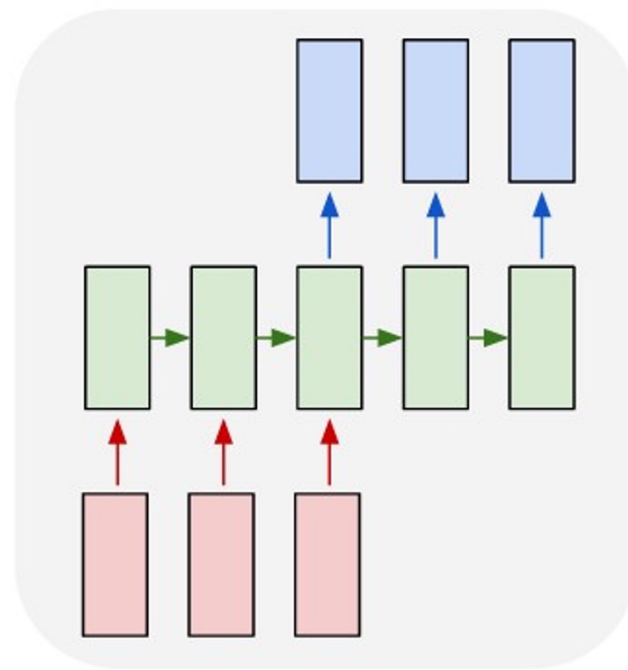
one to many



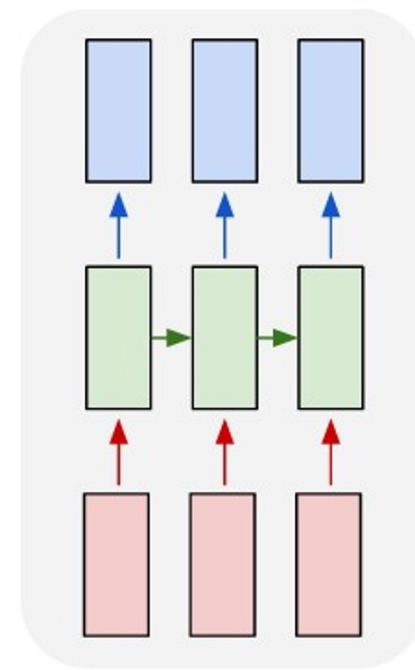
many to one



many to many



many to many



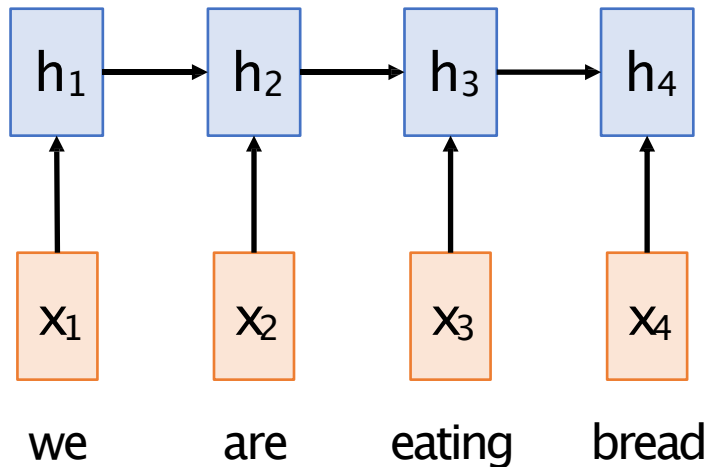


# Sequence-to-Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$



# Sequence-to-Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

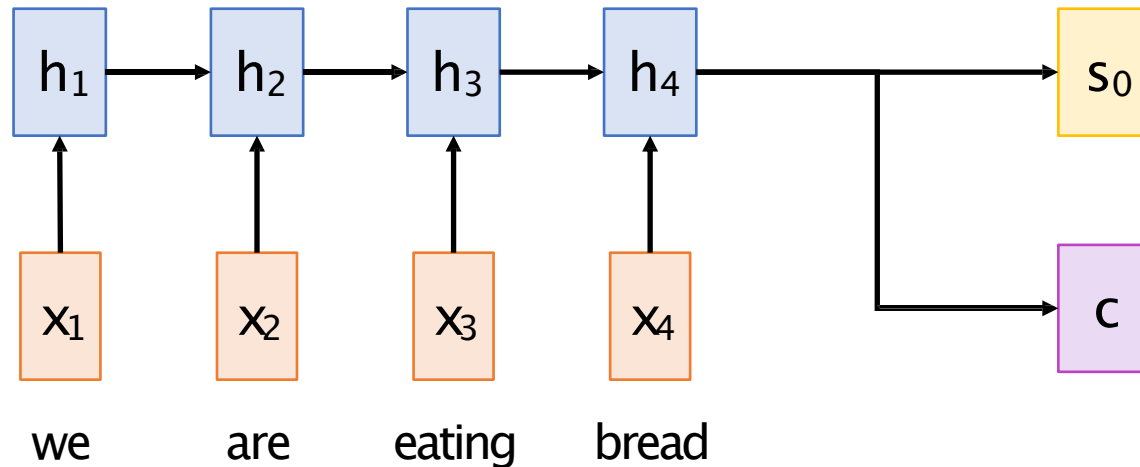
**Output:** Sequence  $y_1, \dots, y_T$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c = h_T$ )



# Sequence-to-Sequence with RNNs

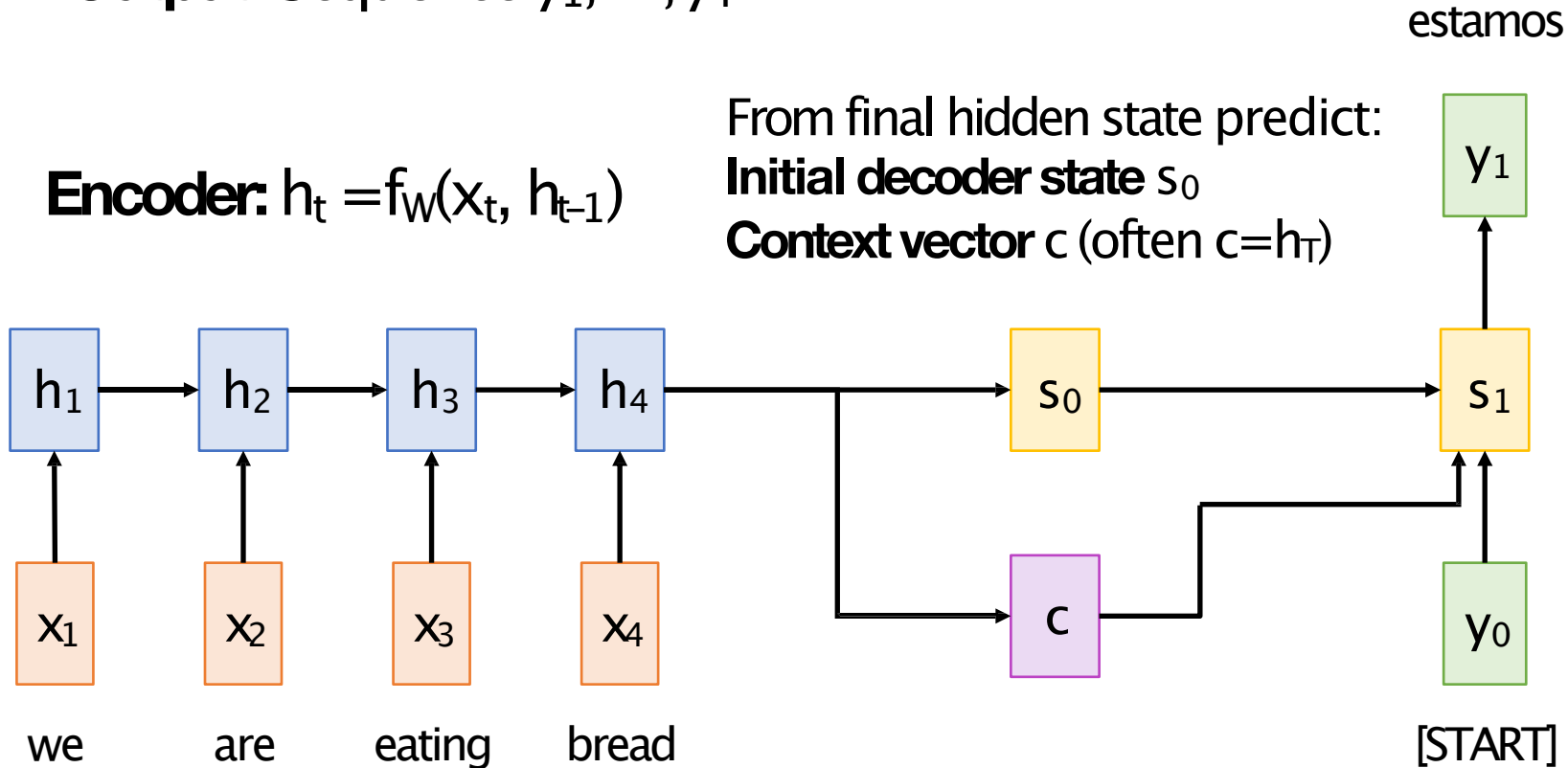
**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$

From final hidden state predict:  
**Initial decoder state**  $s_0$   
**Context vector**  $c$  (often  $c = h_T$ )



# Sequence-to-Sequence with RNNs

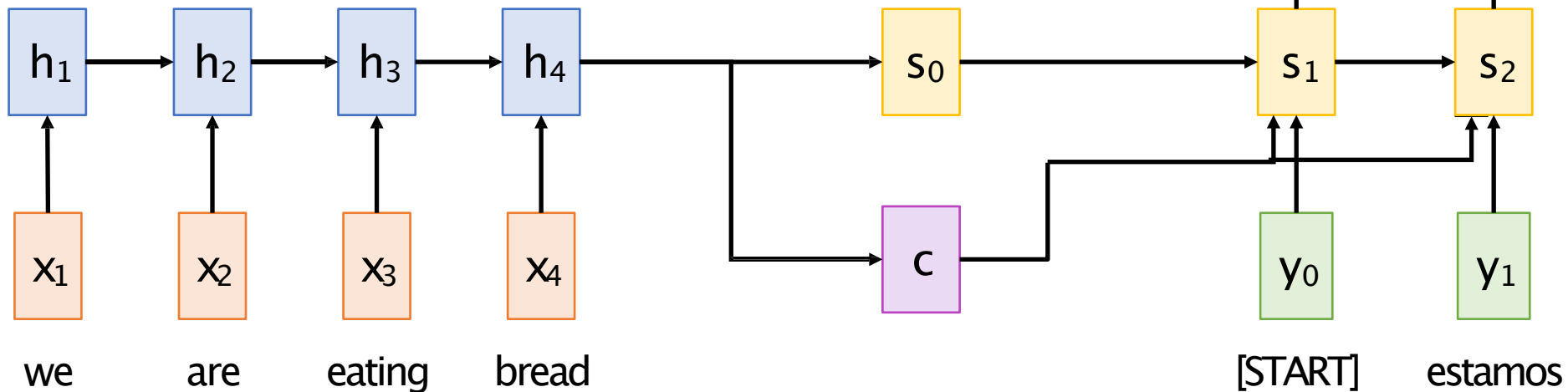
**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$

From final hidden state predict:  
**Initial decoder state**  $s_0$   
**Context vector**  $c$  (often  $c = h_T$ )



# Sequence-to-Sequence with RNNs

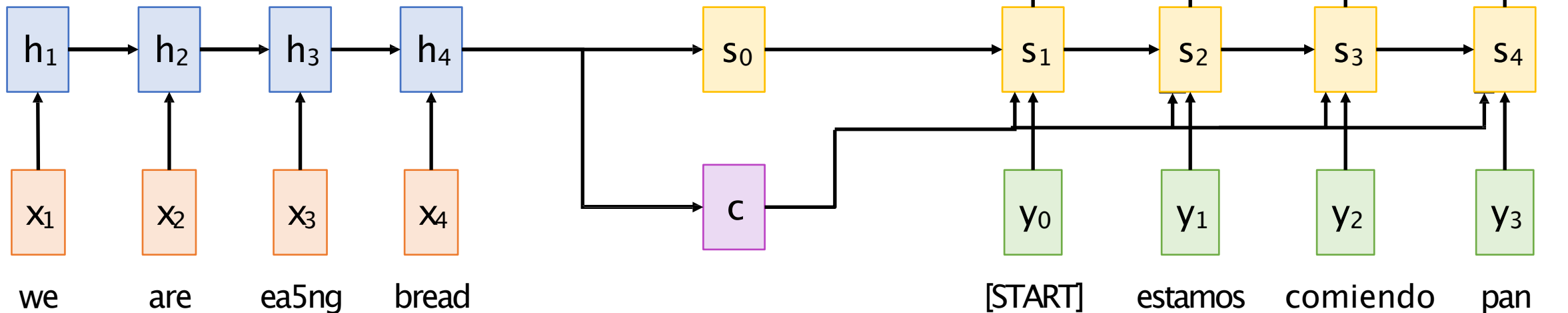
**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$

From final hidden state predict:  
**Initial decoder state**  $s_0$   
**Context vector**  $c$  (often  $c = h_T$ )



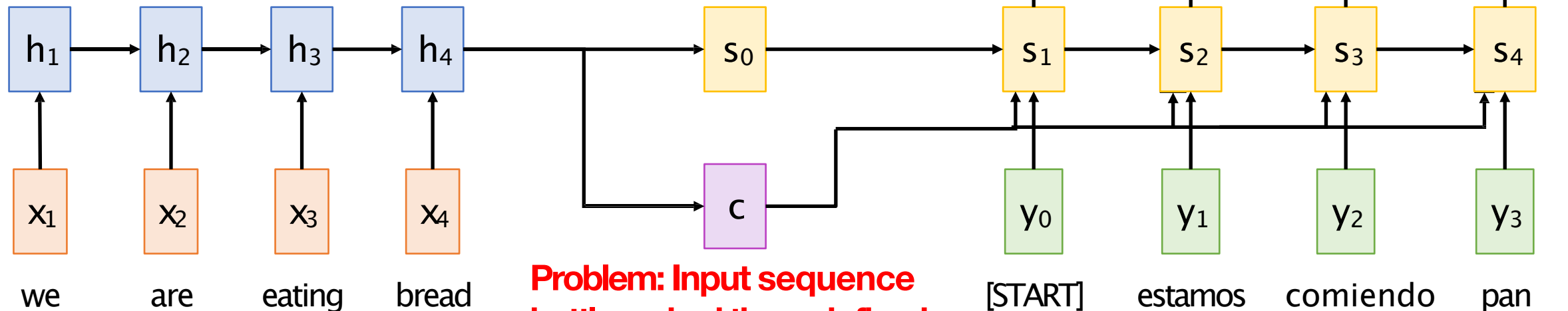
# Sequence-to-Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$

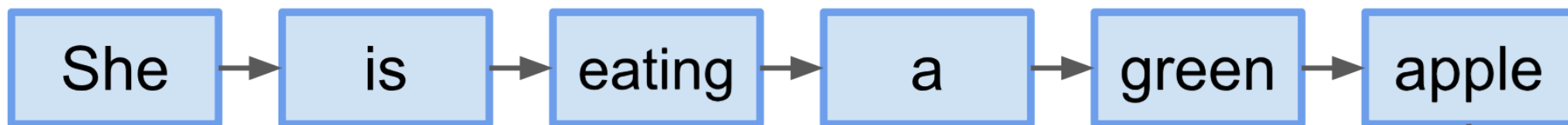


From final hidden state predict:  
**Initial decoder state**  $s_0$   
**Context vector**  $c$  (often  $c=h_T$ )

**Problem: Input sequence bottlenecked through fixed-sized vector. What if  $T=1000$ ?**

# Sequence 2 Sequence models in language

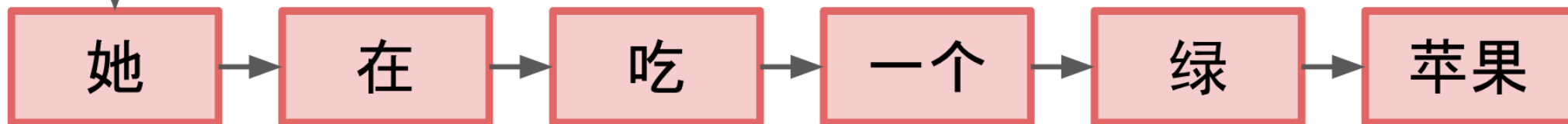
**Encoder**



Context vector (length: 5)

$[0.1, -0.2, 0.8, 1.5, -0.3]$

**Decoder**



---

# 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\* †**  
University of Toronto  
aidan@cs.toronto.edu

**Łukasz Kaiser\***  
Google Brain  
lukaszkaizer@google.com

**Illia Polosukhin\* †**  
illia.polosukhin@gmail.com

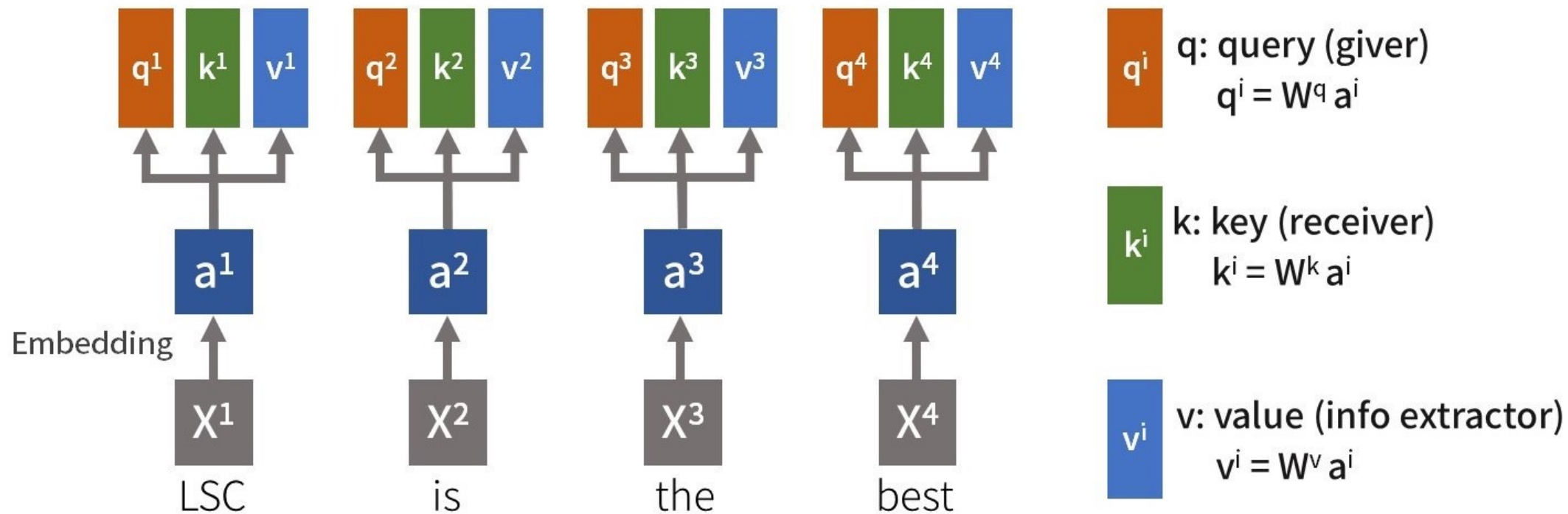
## Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based on a sequence-to-sequence model with multi-head attention without an encoder-decoder architecture.

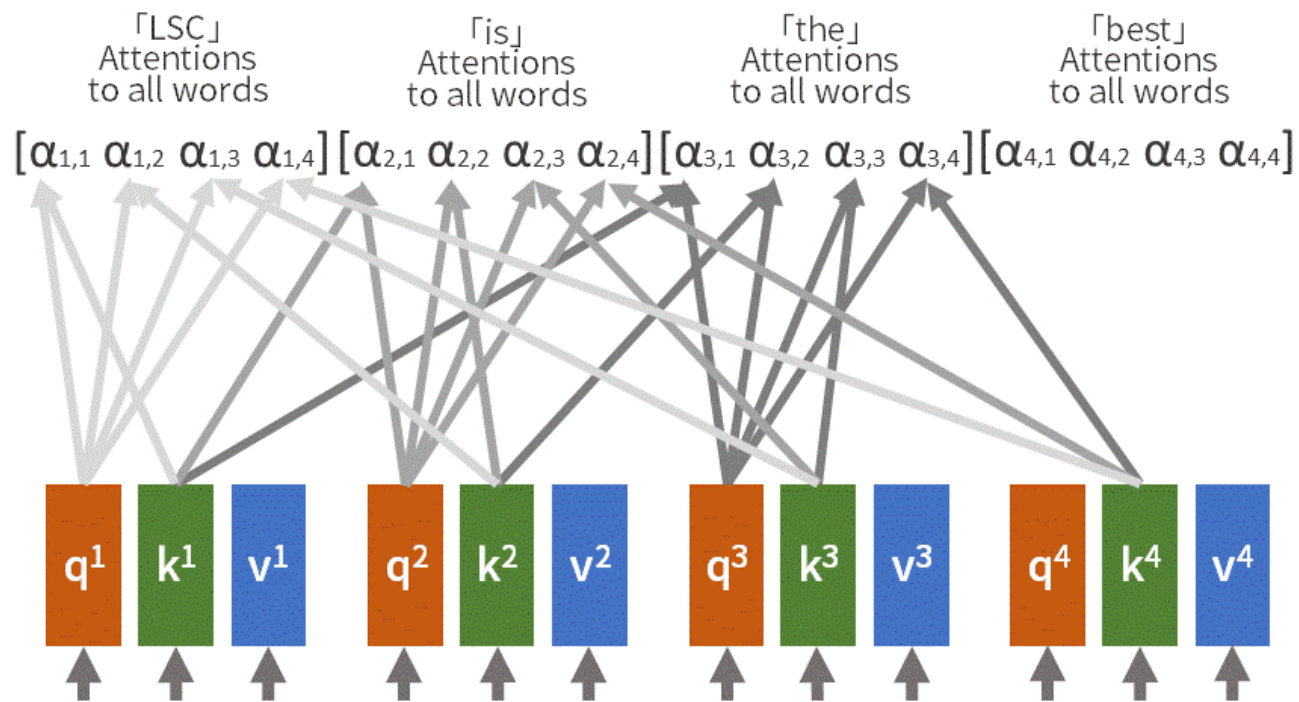


# Attention Operation

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



**Input: LSC is the best!**



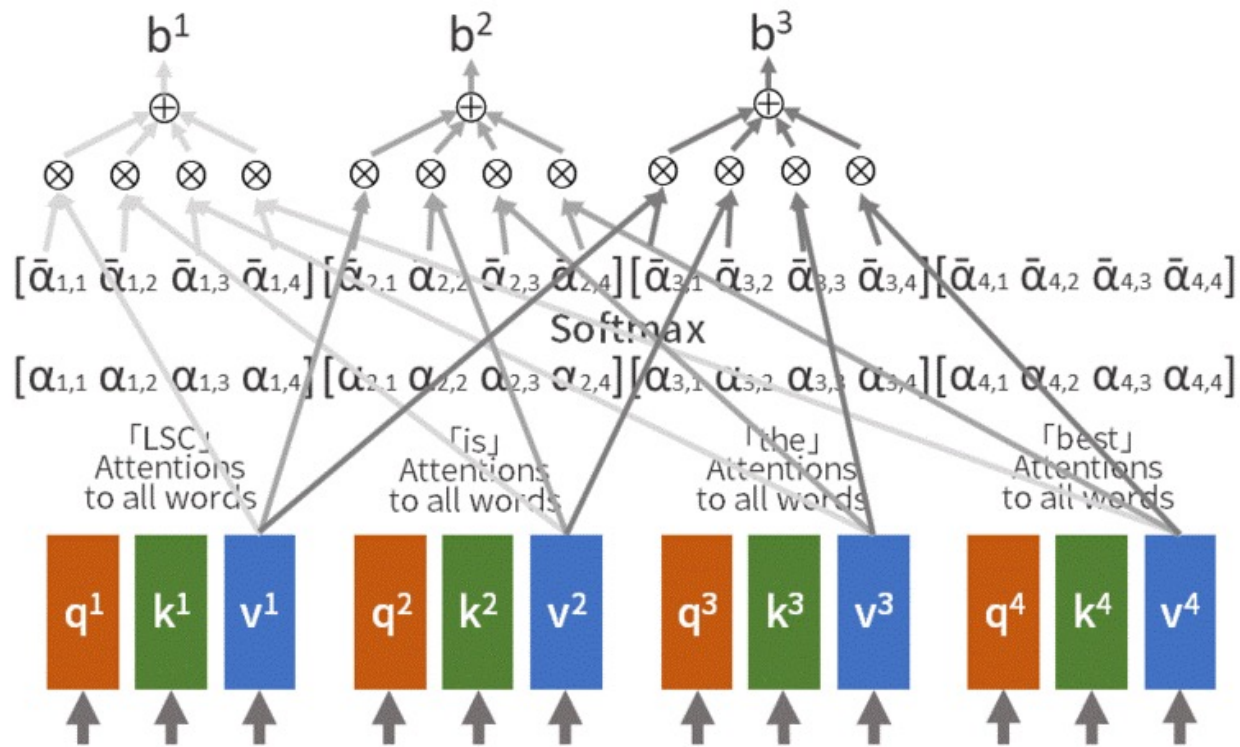
$$\alpha_{i,j} = \frac{q^i \cdot k^j}{\sqrt{d}}$$

d: dimension of q, k

A =

$$\begin{bmatrix} \alpha_{1,1} & \alpha_{1,2} & \alpha_{1,3} & \alpha_{1,4} \\ \alpha_{2,1} & \alpha_{2,2} & \alpha_{2,3} & \alpha_{2,4} \\ \alpha_{3,1} & \alpha_{3,2} & \alpha_{3,3} & \alpha_{3,4} \\ \alpha_{4,1} & \alpha_{4,2} & \alpha_{4,3} & \alpha_{4,4} \end{bmatrix}$$

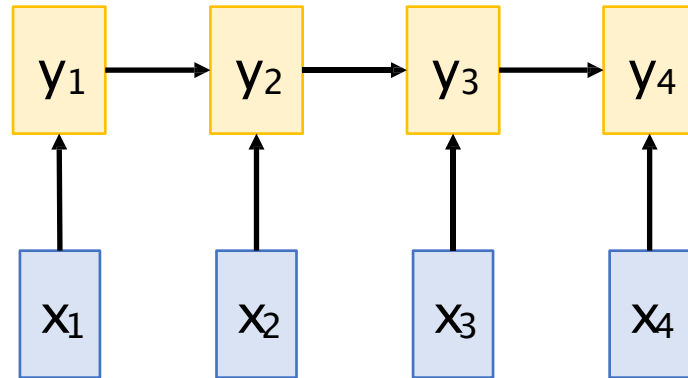
Attention Matrix



$$b^i = \sum_j \bar{\alpha}_{i,j} v^j$$

# Ways of Processing Sequences

## Recurrent Neural Network

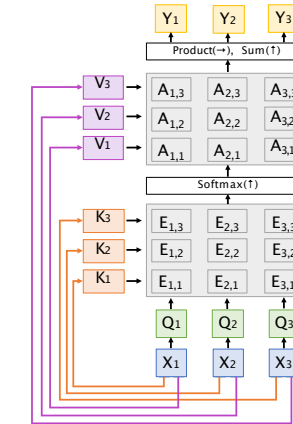


Works on **Ordered Sequences**

(+) **Good at long sequences:** After one RNN layer,  $h_T$  "sees" the whole sequence

(-) **Not parallelizable:** need to compute hidden states sequentially

## Self-Attention



Works on **Sets of Vectors**

(-) **Good at long sequences:** after one self-attention layer, each output "sees" all inputs!

(+) **Highly parallel:** Each output can be computed in parallel

(-) **Very memory intensive**

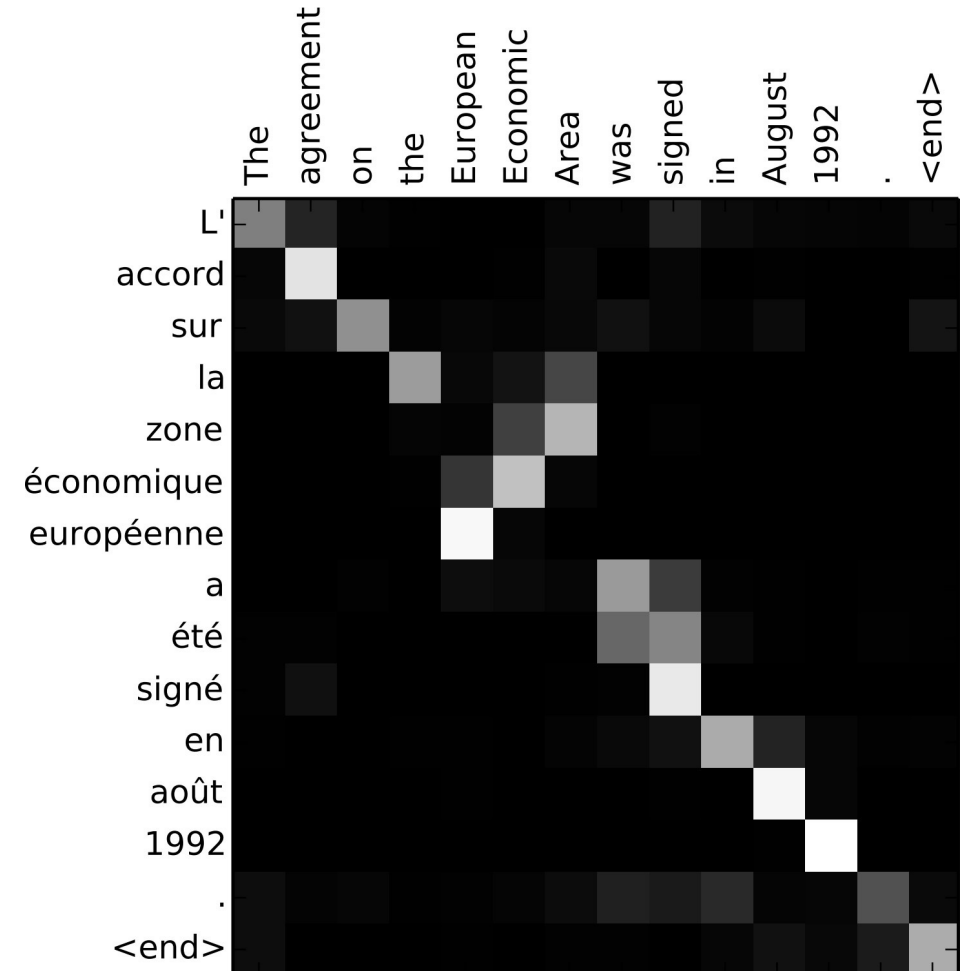
# Sequence-to-Sequence with Attention

**Example:** English to French translation

**Input:** “The agreement on the European Economic Area was signed in August 1992.”

**Output:** “L’accord sur la zone économique européenne a été signé en août 1992.”

Visualize attention weights  $a_{t,i}$



# Sequence-to-Sequence with RNNs and Attention

**Example:** English to French translation

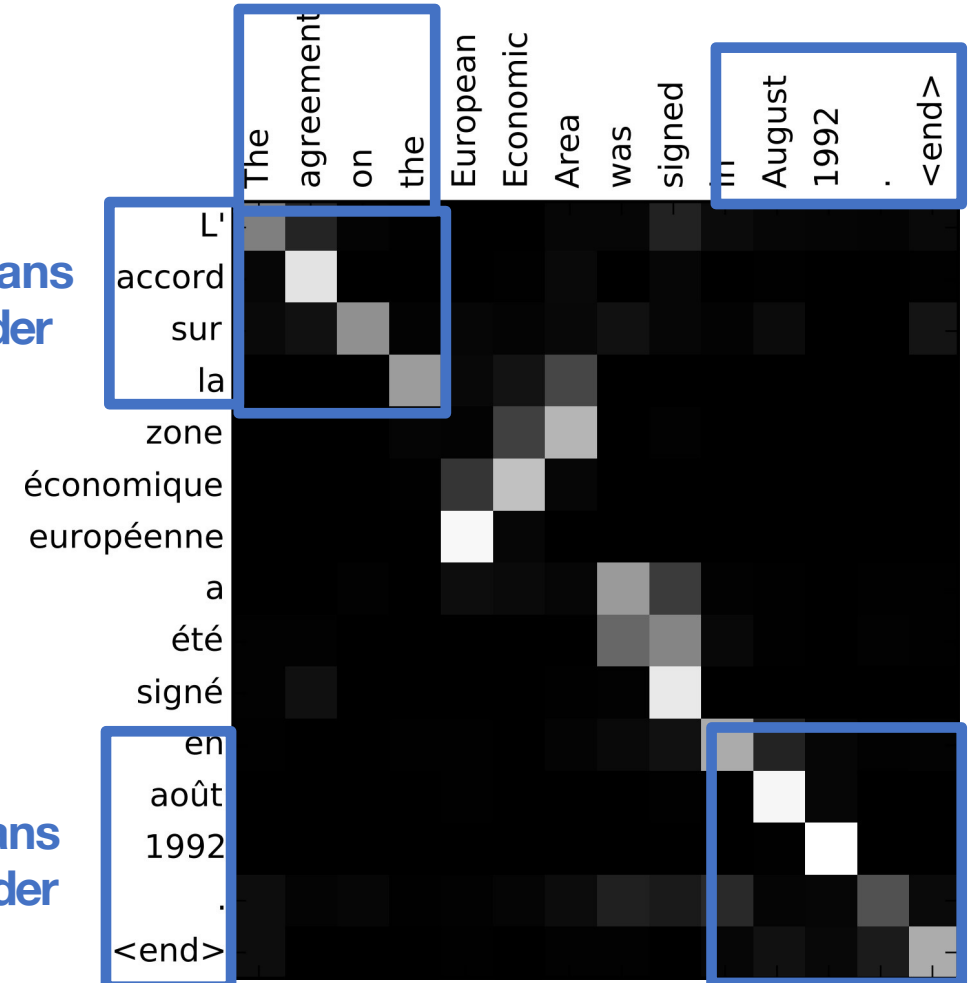
**Input:** “**The agreement on the European Economic Area was signed in August 1992.**”

**Output:** “**L'accord sur la zone économique européenne a été signé en août 1992.**”

Visualize attention weights  $a_{t,i}$

Diagonal attention means words correspond in order

Diagonal attention means words correspond in order



# Sequence-to-Sequence with RNNs and Attention

**Example:** English to French translation

**Input:** “The agreement on the European Economic Area was signed in August 1992.”

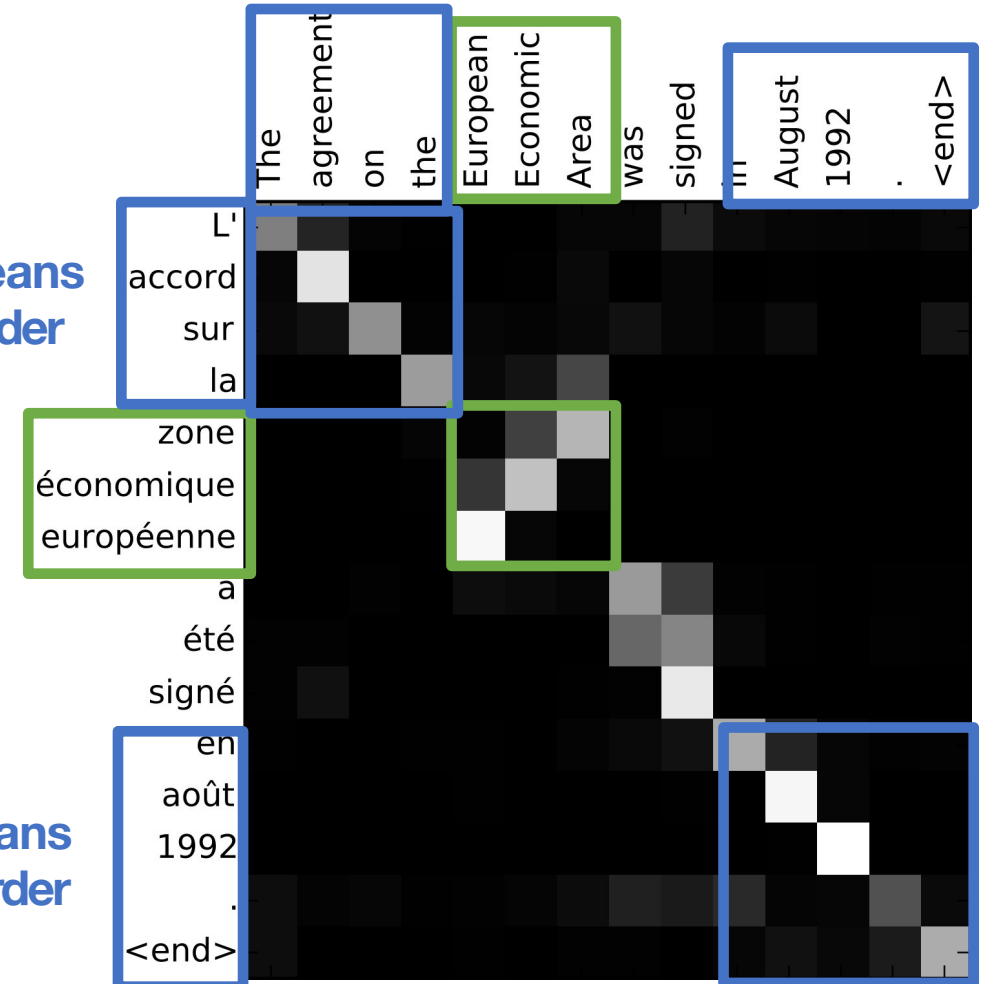
**Output:** “L'accord sur la zone économique européenne a été signé en août 1992.”

Visualize attention weights  $a_{t,i}$

Diagonal attention means words correspond in order

Attention figures out different word orders

Diagonal attention means words correspond in order





# Sequence-to-Sequence with RNNs and Attention

**Example:** English to French translation

**Input:** “The agreement on the European Economic Area was signed in August 1992.”

**Output:** “L'accord sur la zone économique européenne a été signé en août 1992.”

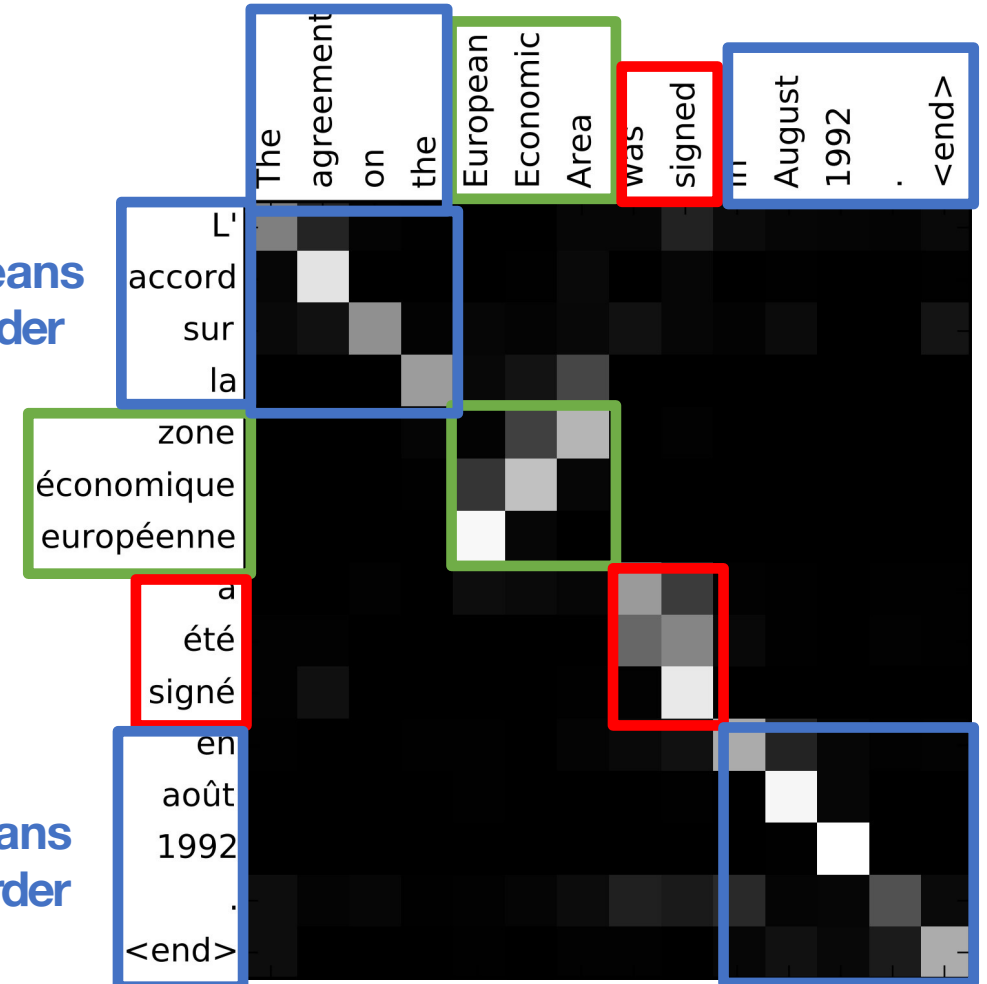
Visualize attention weights  $a_{t,i}$

Diagonal attention means words correspond in order

A)en+on figures out different word orders

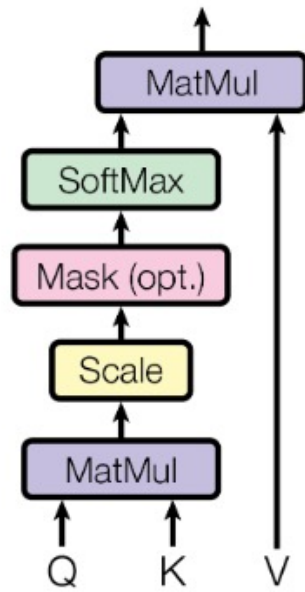
Verb conjugation

Diagonal a)en+on means words correspond in order

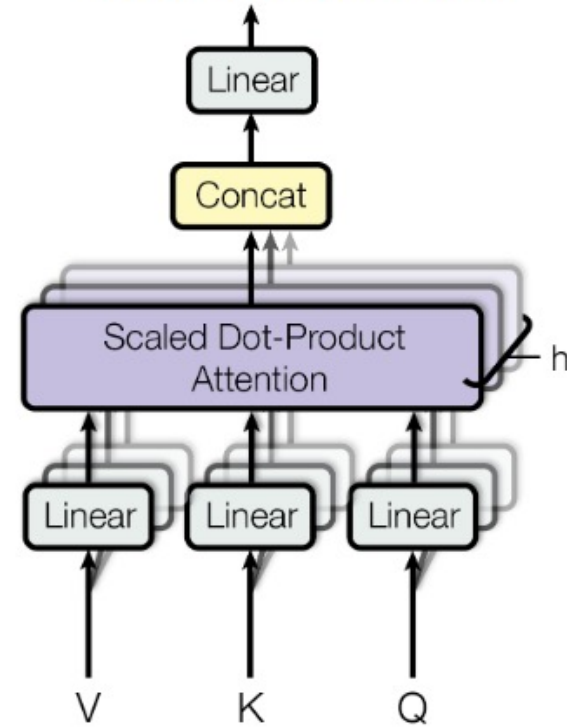


# Multi-head attention

Scaled Dot-Product Attention



Multi-Head Attention



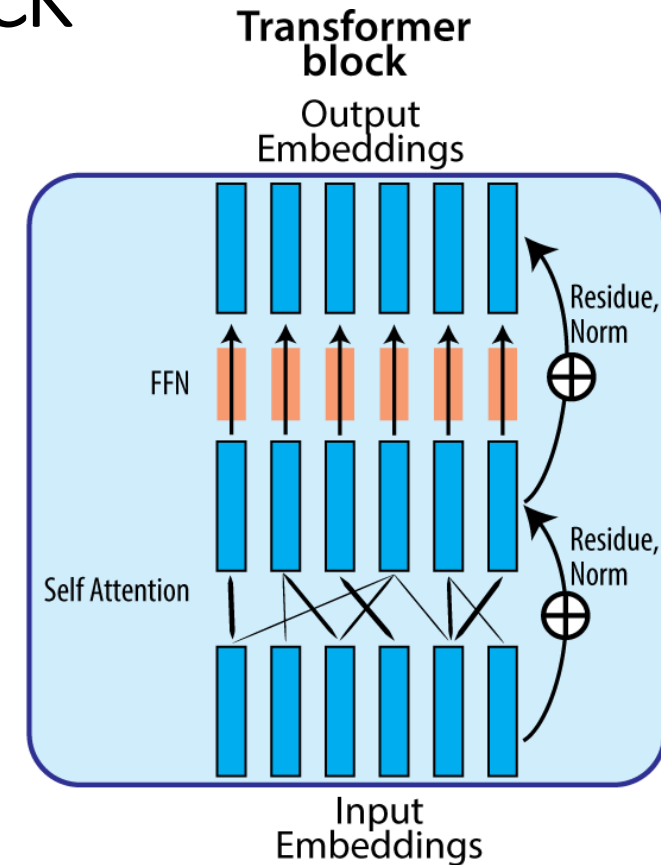
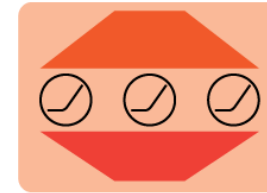
# Cross vs Self Attention

- Cross Attention
  - Key, and Value from one set of tokens
  - Query from another set of tokens
  - E.g. words in one language pay attention to words in **another**.
- Self Attention
  - Key, Value, and Query from the same set of tokens

# From Attention to TransformerBlock

- A **transformer block** has
  - Self Attention
    - information exchange *between tokens*
  - Feed forward network
    - Information transform *within tokens*
    - E.g. a multi-layer perceptron with 1 hidden layer
    - GeLU activation is commonly used.
  - Normalization (Layer normalization)
- A transformer model contains N x **transformer block**

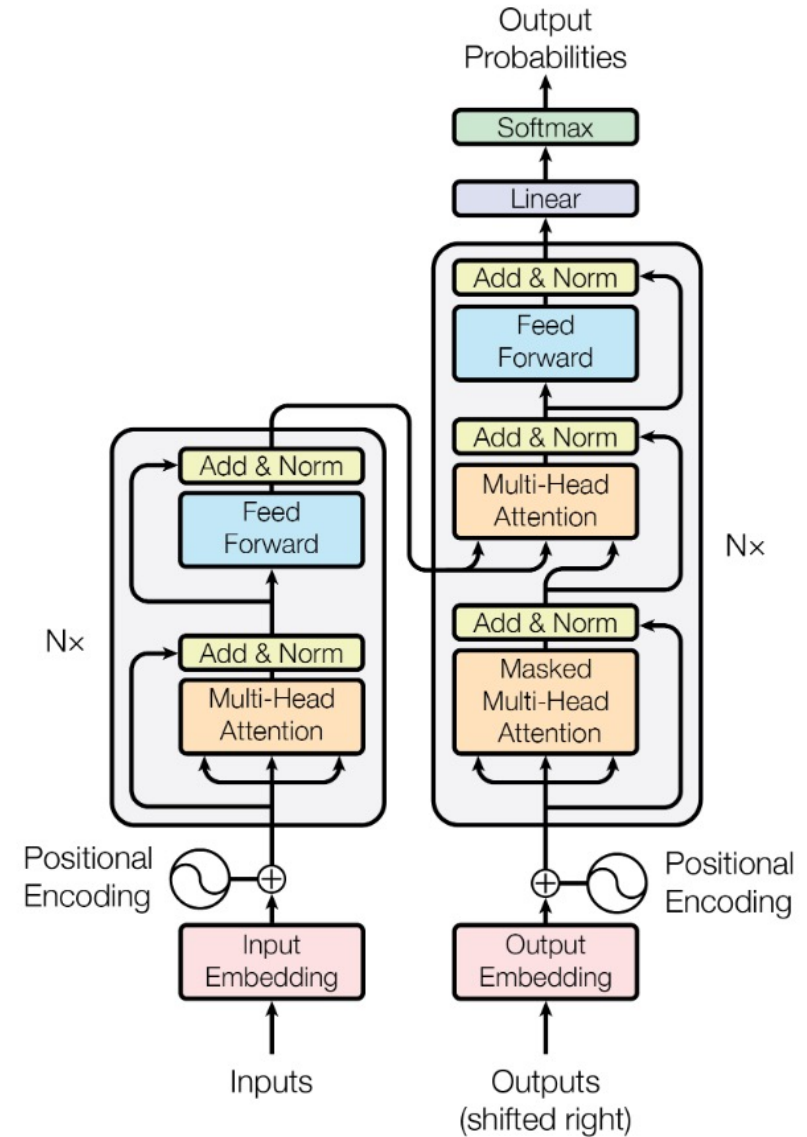
FFN  
(2 layer MLP)



# Transformer Architecture

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

| Model                           | BLEU        |              | Training Cost (FLOPs)                 |                     |
|---------------------------------|-------------|--------------|---------------------------------------|---------------------|
|                                 | EN-DE       | EN-FR        | EN-DE                                 | EN-FR               |
| ByteNet [18]                    | 23.75       |              |                                       |                     |
| Deep-Att + PosUnk [39]          |             | 39.2         |                                       | $1.0 \cdot 10^{20}$ |
| GNMT + RL [38]                  | 24.6        | 39.92        | $2.3 \cdot 10^{19}$                   | $1.4 \cdot 10^{20}$ |
| ConvS2S [9]                     | 25.16       | 40.46        | $9.6 \cdot 10^{18}$                   | $1.5 \cdot 10^{20}$ |
| MoE [32]                        | 26.03       | 40.56        | $2.0 \cdot 10^{19}$                   | $1.2 \cdot 10^{20}$ |
| Deep-Att + PosUnk Ensemble [39] |             | 40.4         |                                       | $8.0 \cdot 10^{20}$ |
| GNMT + RL Ensemble [38]         | 26.30       | 41.16        | $1.8 \cdot 10^{20}$                   | $1.1 \cdot 10^{21}$ |
| ConvS2S Ensemble [9]            | 26.36       | <b>41.29</b> | $7.7 \cdot 10^{19}$                   | $1.2 \cdot 10^{21}$ |
| Transformer (base model)        | 27.3        | 38.1         | <b><math>3.3 \cdot 10^{18}</math></b> |                     |
| Transformer (big)               | <b>28.4</b> | <b>41.8</b>  | $2.3 \cdot 10^{19}$                   |                     |



# AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE


Alexey Dosovitskiy<sup>\*,†</sup>, Lucas Beyer<sup>\*</sup>, Alexander Kolesnikov<sup>\*</sup>, Dirk Weissenborn<sup>\*</sup>,  
Xiaohua Zhai<sup>\*</sup>, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,  
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby<sup>\*,†</sup>

<sup>\*</sup>equal technical contribution, <sup>†</sup>equal advising

Google Research, Brain Team

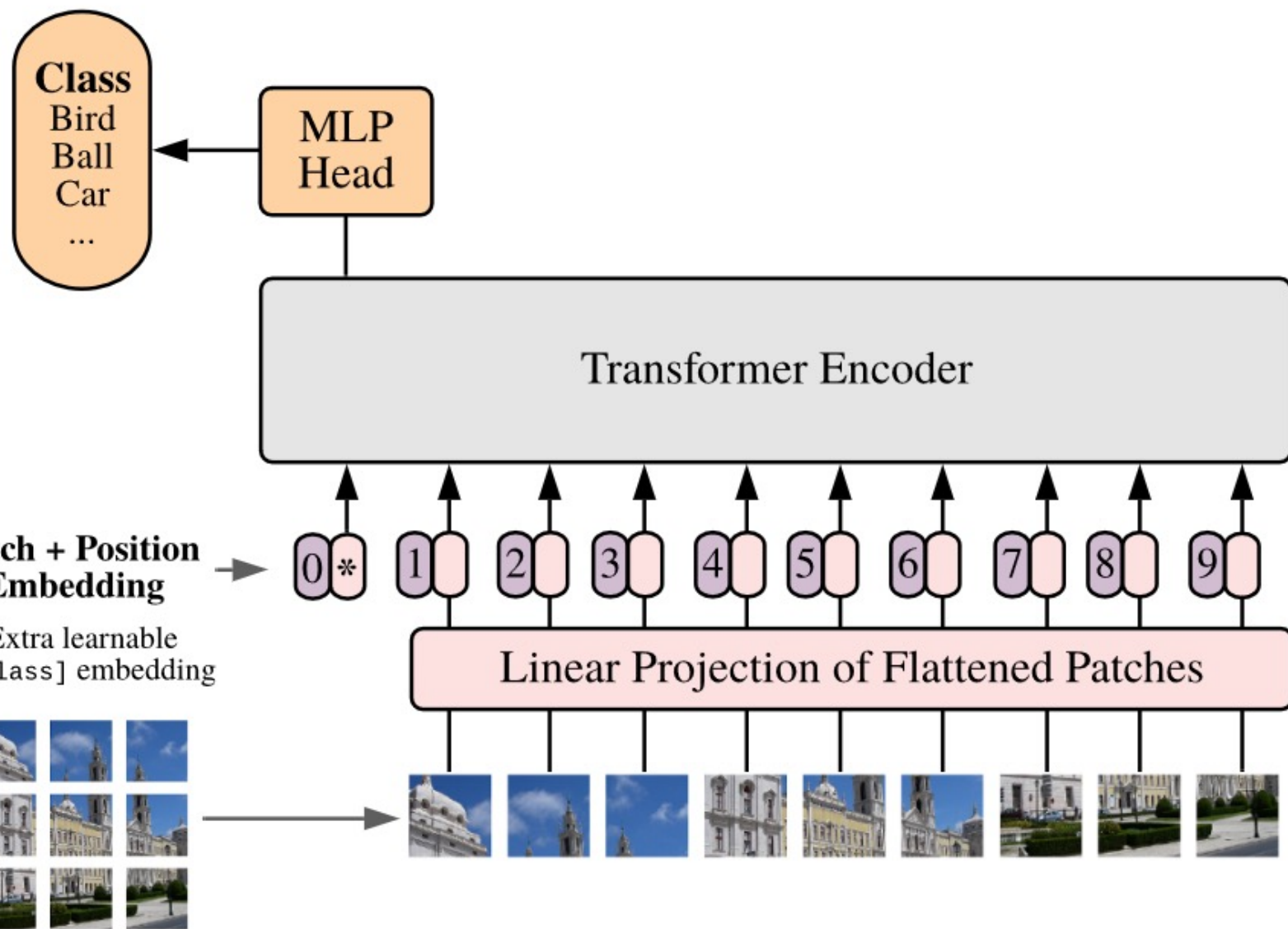
{adosovitskiy, neilhoulby}@google.com

## ABSTRACT

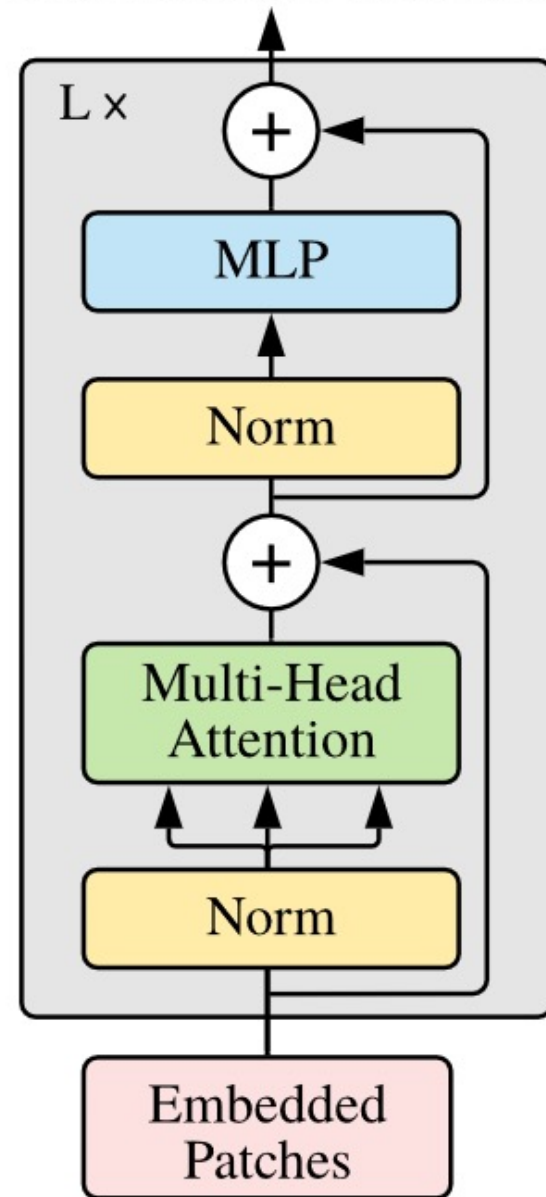
While the Transformer architecture has become the de-facto standard for natural language processing tasks, its applications to computer vision remain limited. In vision, attention is either applied in conjunction with convolutional networks, or used to replace certain components of convolutional networks while keeping their overall structure in place. We show that this reliance on CNNs is not necessary and a pure transformer applied directly to sequences of image patches can perform very well on image classification tasks. When pre-trained on large amounts of data and transferred to multiple mid-sized or small image recognition benchmarks (ImageNet, CIFAR-100, VTAB, etc.), Vision Transformer (ViT) attains excellent results compared to state-of-the-art convolutional networks while requiring substantially fewer computational resources to train. 



# Vision Transformer (ViT)



# Transformer Encoder



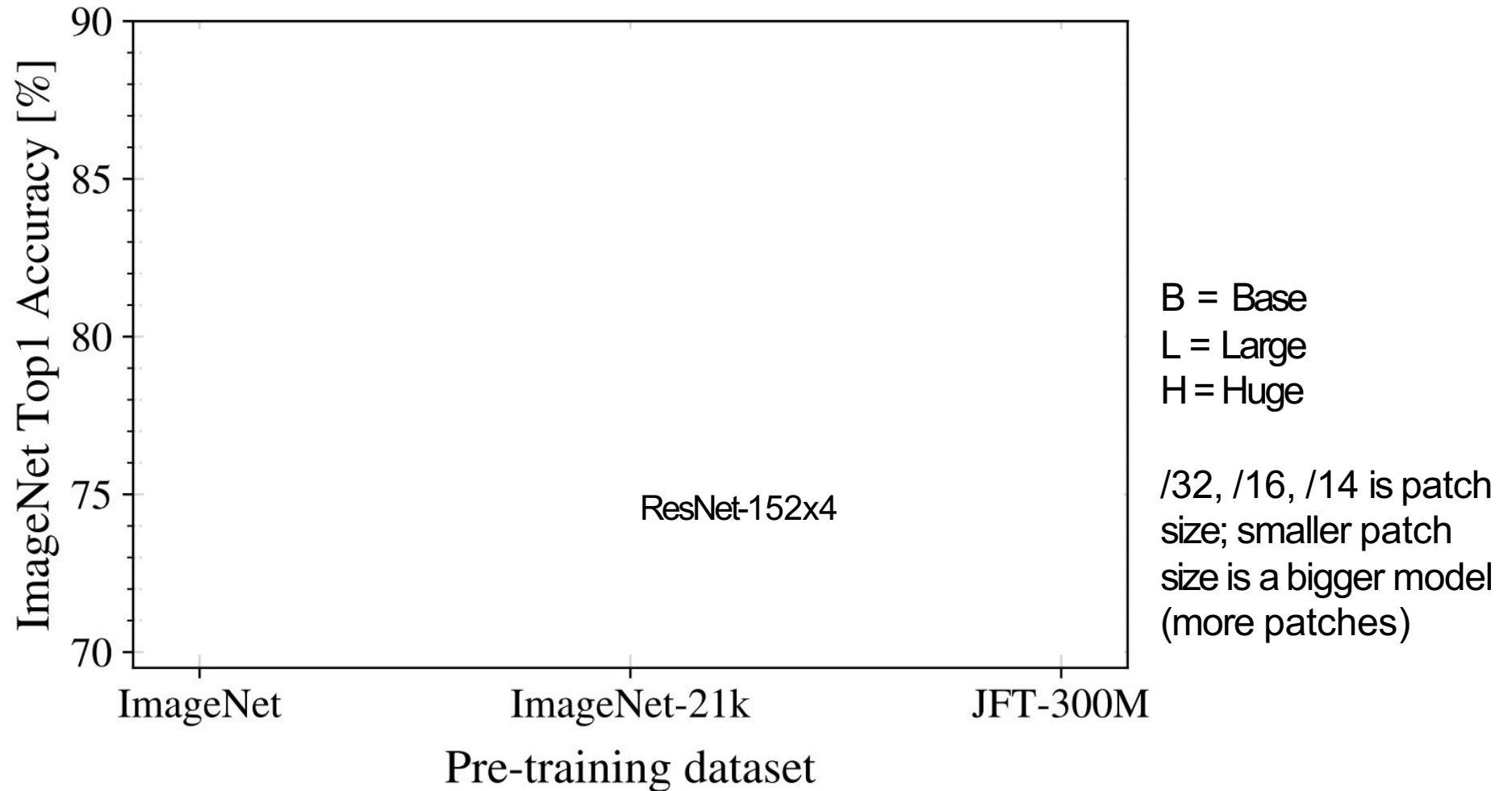


| Model     | Layers | Hidden size $D$ | MLP size | Heads | Params |
|-----------|--------|-----------------|----------|-------|--------|
| ViT-Base  | 12     | 768             | 3072     | 12    | 86M    |
| ViT-Large | 24     | 1024            | 4096     | 16    | 307M   |
| ViT-Huge  | 32     | 1280            | 5120     | 16    | 632M   |

Table 1: Details of Vision Transformer model variants.

|                    | Ours-JFT<br>(ViT-H/14)  | Ours-JFT<br>(ViT-L/16)  | Ours-I21K<br>(ViT-L/16) | BiT-L<br>(ResNet152x4) | Noisy Student<br>(EfficientNet-L2) |
|--------------------|-------------------------|-------------------------|-------------------------|------------------------|------------------------------------|
| ImageNet           | <b>88.55</b> $\pm 0.04$ | 87.76 $\pm 0.03$        | 85.30 $\pm 0.02$        | 87.54 $\pm 0.02$       | 88.4/88.5*                         |
| ImageNet ReaL      | <b>90.72</b> $\pm 0.05$ | 90.54 $\pm 0.03$        | 88.62 $\pm 0.05$        | 90.54                  | 90.55                              |
| CIFAR-10           | <b>99.50</b> $\pm 0.06$ | 99.42 $\pm 0.03$        | 99.15 $\pm 0.03$        | 99.37 $\pm 0.06$       | —                                  |
| CIFAR-100          | <b>94.55</b> $\pm 0.04$ | 93.90 $\pm 0.05$        | 93.25 $\pm 0.05$        | 93.51 $\pm 0.08$       | —                                  |
| Oxford-IIIT Pets   | <b>97.56</b> $\pm 0.03$ | 97.32 $\pm 0.11$        | 94.67 $\pm 0.15$        | 96.62 $\pm 0.23$       | —                                  |
| Oxford Flowers-102 | 99.68 $\pm 0.02$        | <b>99.74</b> $\pm 0.00$ | 99.61 $\pm 0.02$        | 99.63 $\pm 0.03$       | —                                  |
| VTAB (19 tasks)    | <b>77.63</b> $\pm 0.23$ | 76.28 $\pm 0.46$        | 72.72 $\pm 0.21$        | 76.29 $\pm 1.70$       | —                                  |
| TPUv3-core-days    | 2.5k                    | 0.68k                   | 0.23k                   | 9.9k                   | 12.3k                              |

# Vision Transformer (ViT) vs ResNets

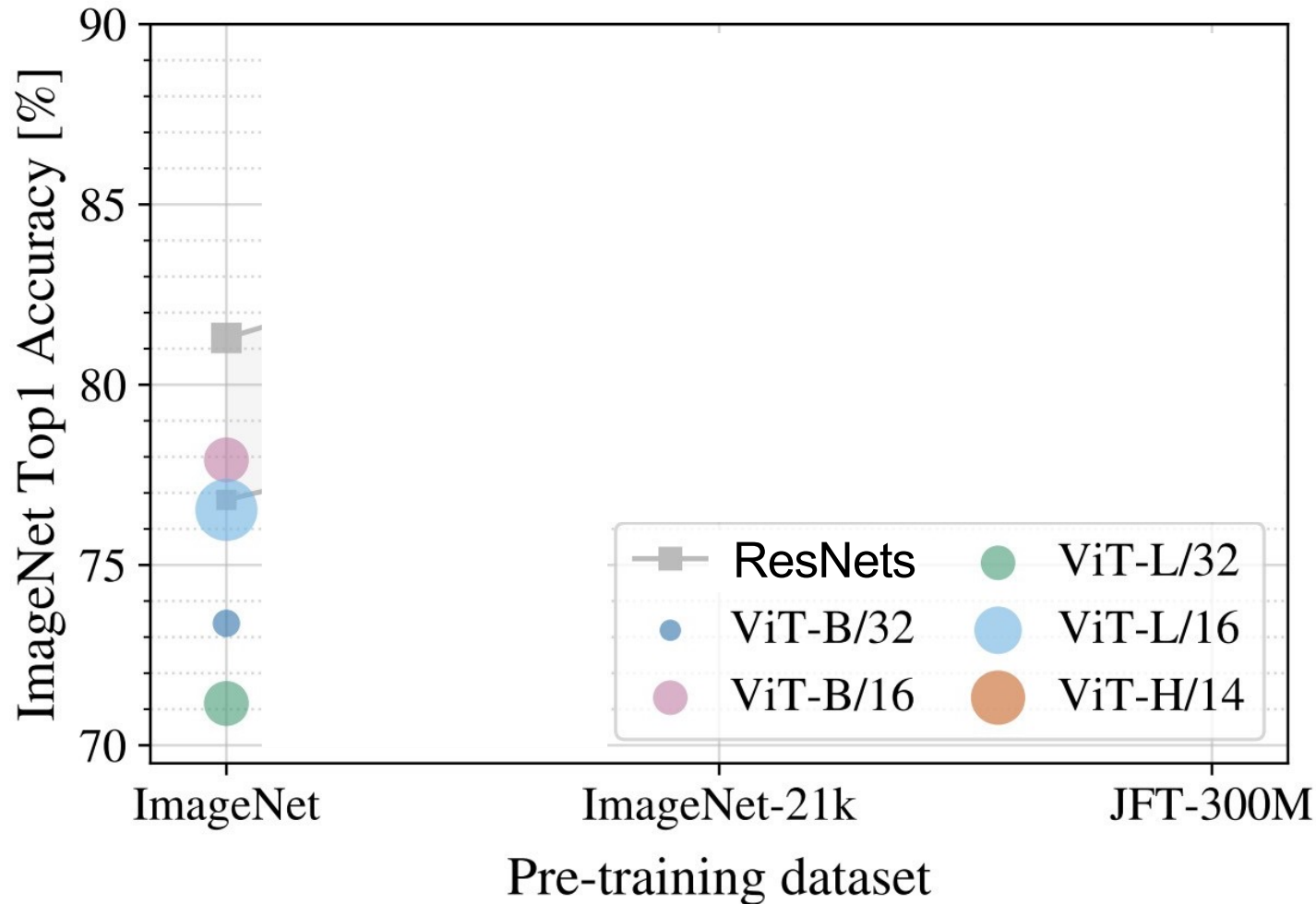


Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

# Vision Transformer (ViT) vs ResNets

Recall: ImageNet dataset has 1k categories, 1.2M images

When trained on ImageNet, ViT models perform worse than ResNets



B = Base  
L = Large  
H = Huge

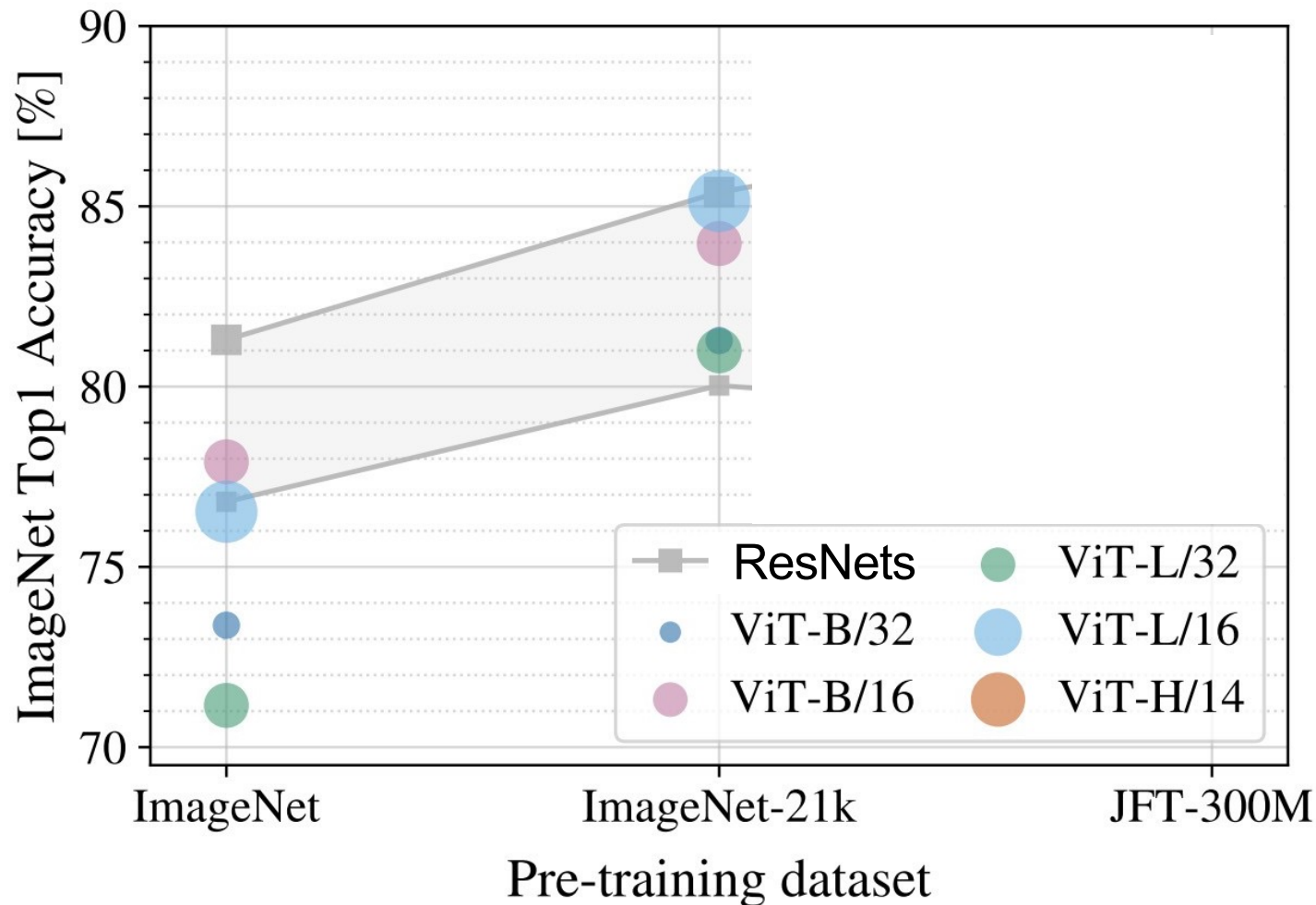
/32, /16, /14 is patch size; smaller patch size is a bigger model (more patches)

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

# Vision Transformer (ViT) vs ResNets

ImageNet-21k has 14M images with 21k categories

If you pretrain on ImageNet-21k and fine-tune on ImageNet, ViT does better: big ViTs match big ResNets



B = Base  
L = Large  
H = Huge

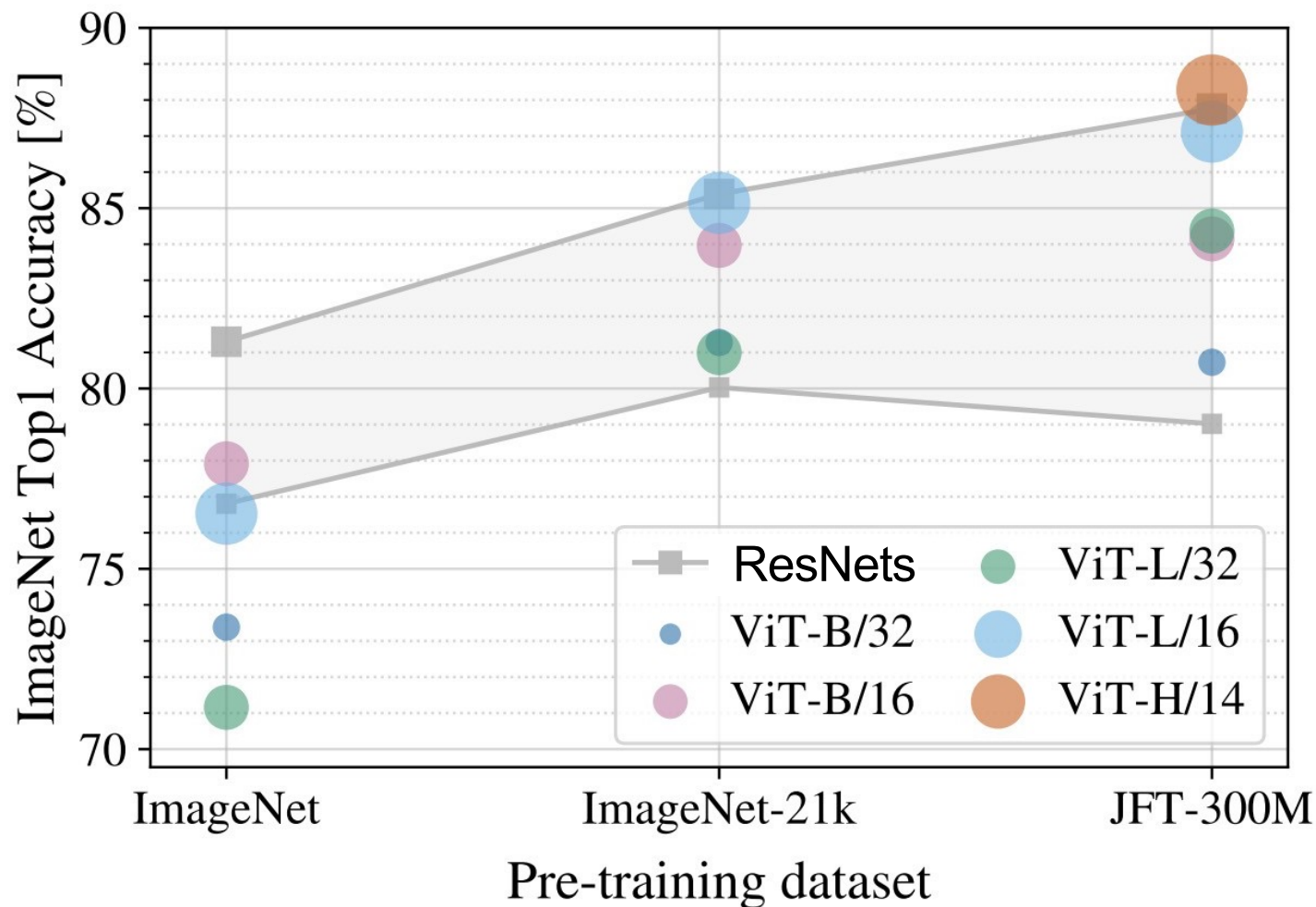
/32, /16, /14 is patch size; smaller patch size is a bigger model (more patches)

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

# Vision Transformer (ViT) vs ResNets

JFT-300M is an internal Google dataset with 300M labeled images

If you pretrain on JFT and finetune on ImageNet, large ViTs outperform large ResNets



B = Base  
L = Large  
H = Huge

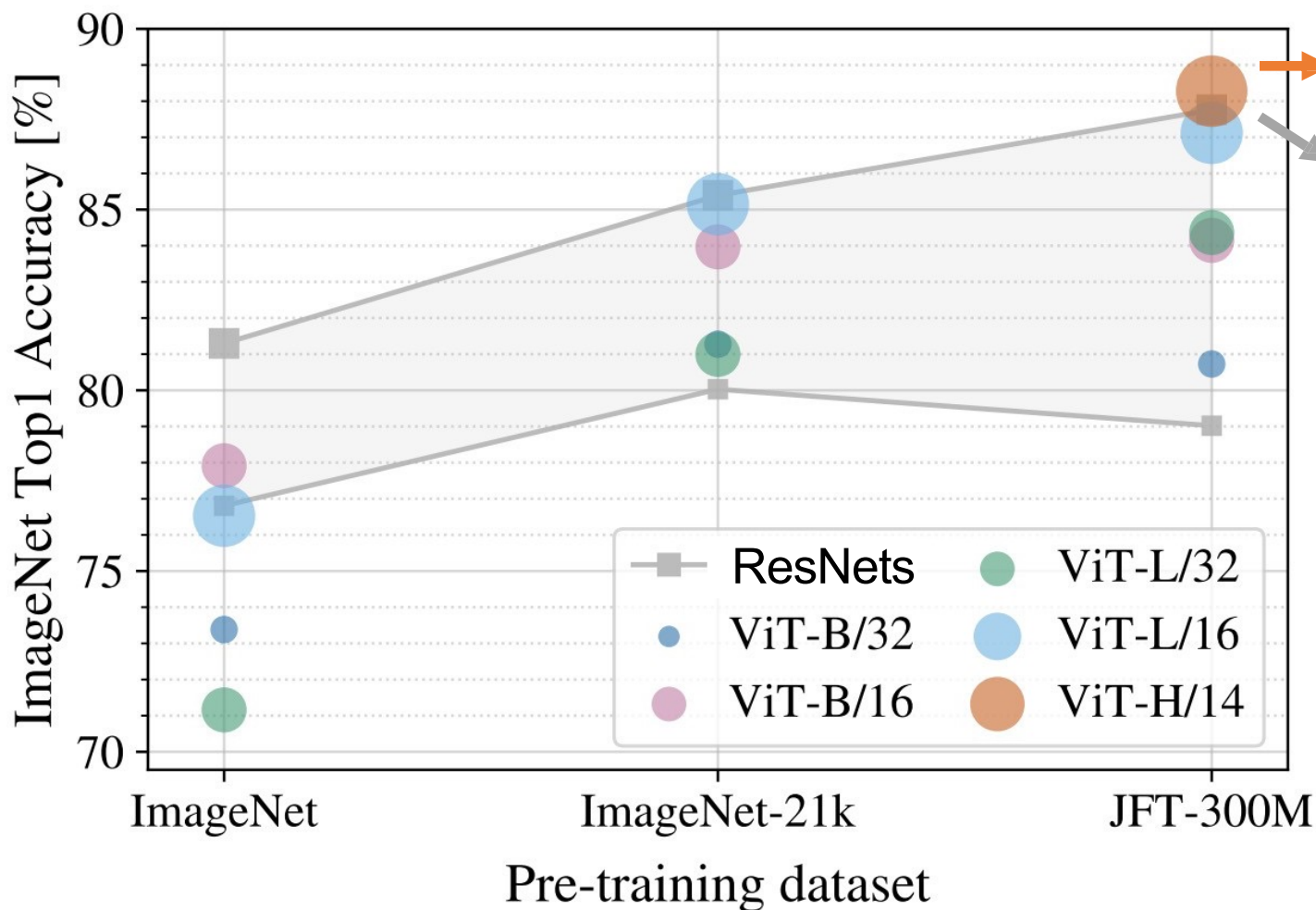
/32, /16, /14 is patch size; smaller patch size is a bigger model (more patches)

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

# Vision Transformer (ViT) vs ResNets

JFT-300M is an internal Google dataset with 300M labeled images

If you pretrain on JFT and finetune on ImageNet, large ViTs outperform large ResNets



ViT: 2.5k TPU-v3 core days of training

ResNet: 9.9k TPU-v3 core days of training

ViTs make more efficient use of GPU / TPU hardware (matrix multiply is more hardware-friendly than conv)

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

# Attention and Transformer

- Sequence to sequence network with RNN
- Attention module and Transformer network
- Vision Transformer vs CNN for Image Classification