# EE782 Advanced Topics in Machine Learning
# 06 Neural Networks for NLP

Amit Sethi

Electrical Engineering, IIT Bombay

# Module objectives

- Design LSTM based networks for various NLP tasks

- Articulate the information bottleneck in encoder-decoder architectures

- Explain how attention alleviates the information bottleneck

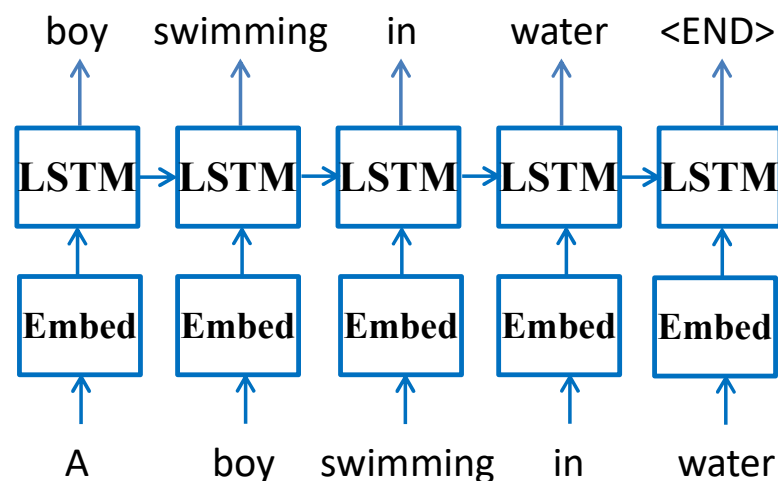- Understand how attention can do almost all tasks done by the other layers

# Pre-processing for NLP

- The most basic pre-processing is to convert words into an embedding using Word2Vec or GloVe

- Otherwise, a one-hot-bit input vector can be too long and sparse, and require lots on input weights
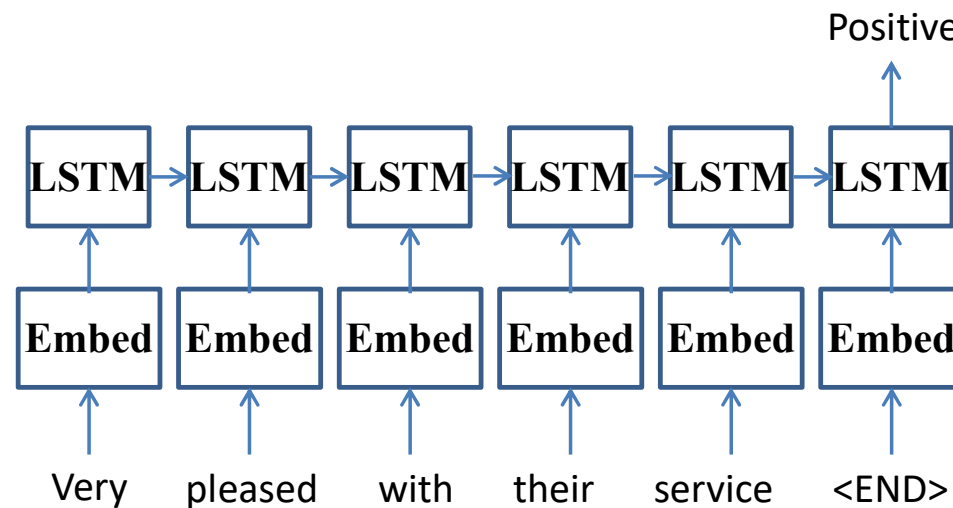
# Pre-training LSTMs

- Learning to predict the next word can imprint powerful language models in LSTMs

- This captures the grammar and syntax

- Usually, LSTMs are pre-trained on corpora
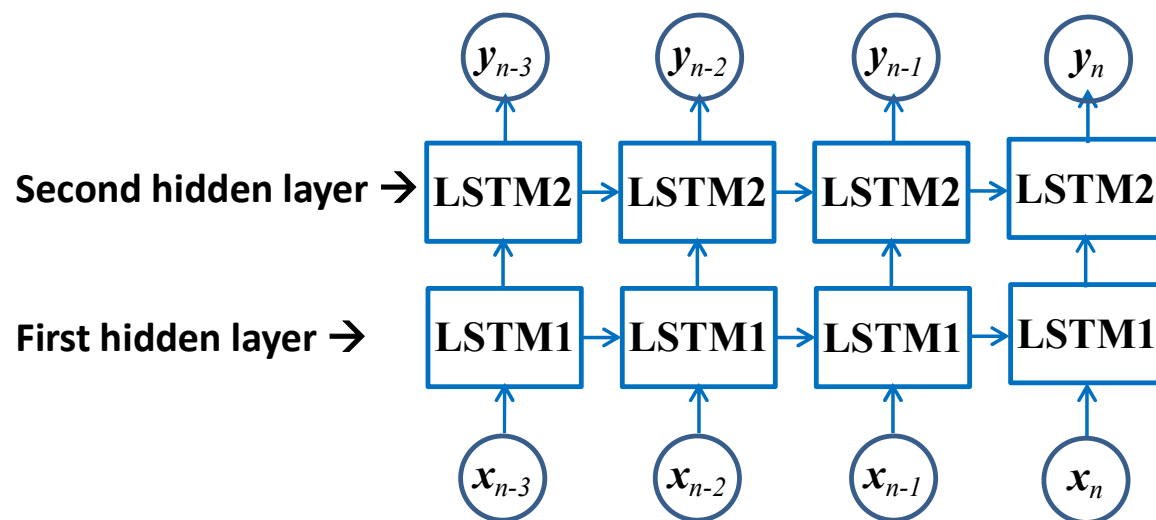
# Sentiment analysis

- Very common for customer review or new article analysis
- Output before the end can be discarded (not used for backpropagation)
- This is a many-to-one task
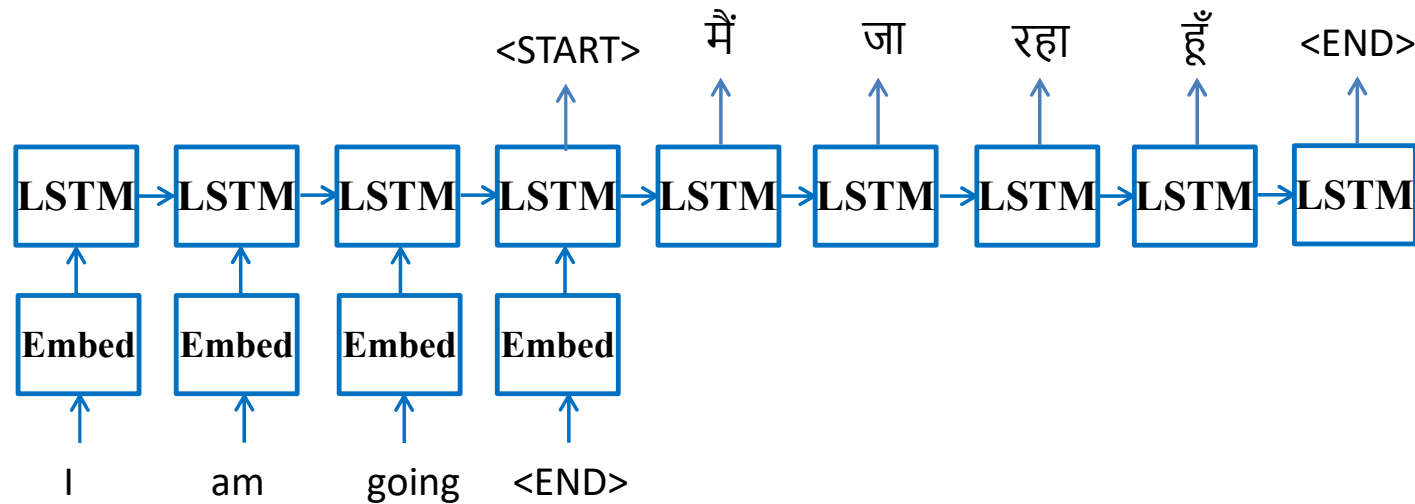
# Multi-layer LSTM

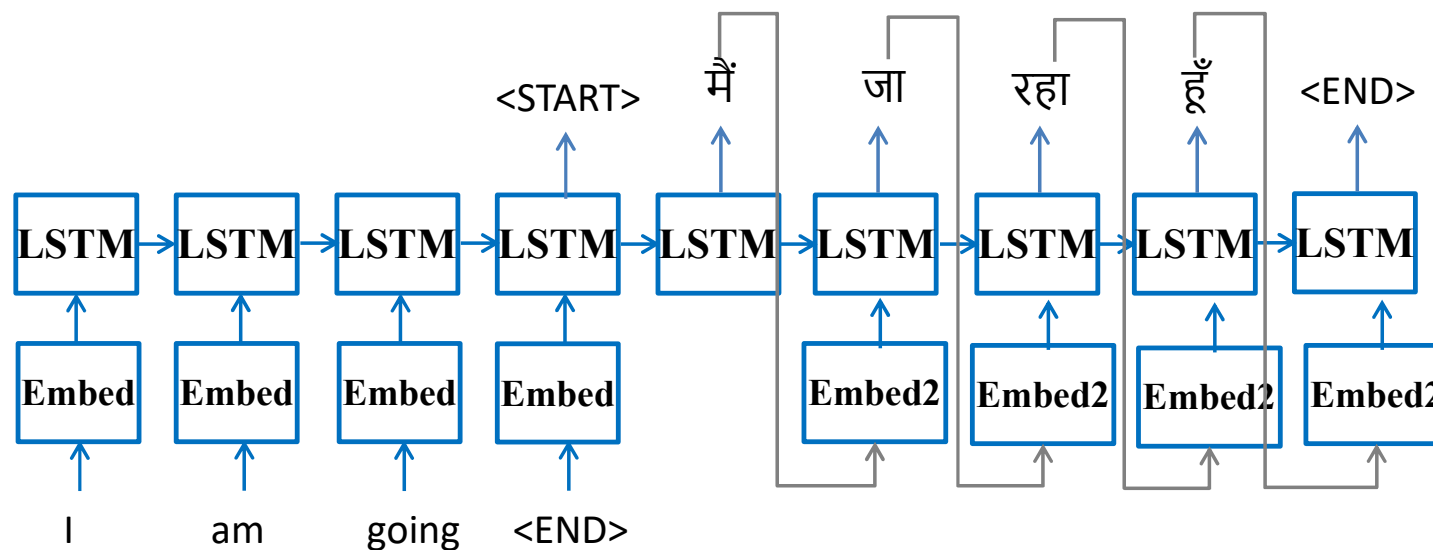- More than one hidden layer can be used

# Machine translation

- A naïve model would be to use a many-to-many network and directly train it

# Machine translation

- One could also feed in the output to the next instance input to predict a coherent structure
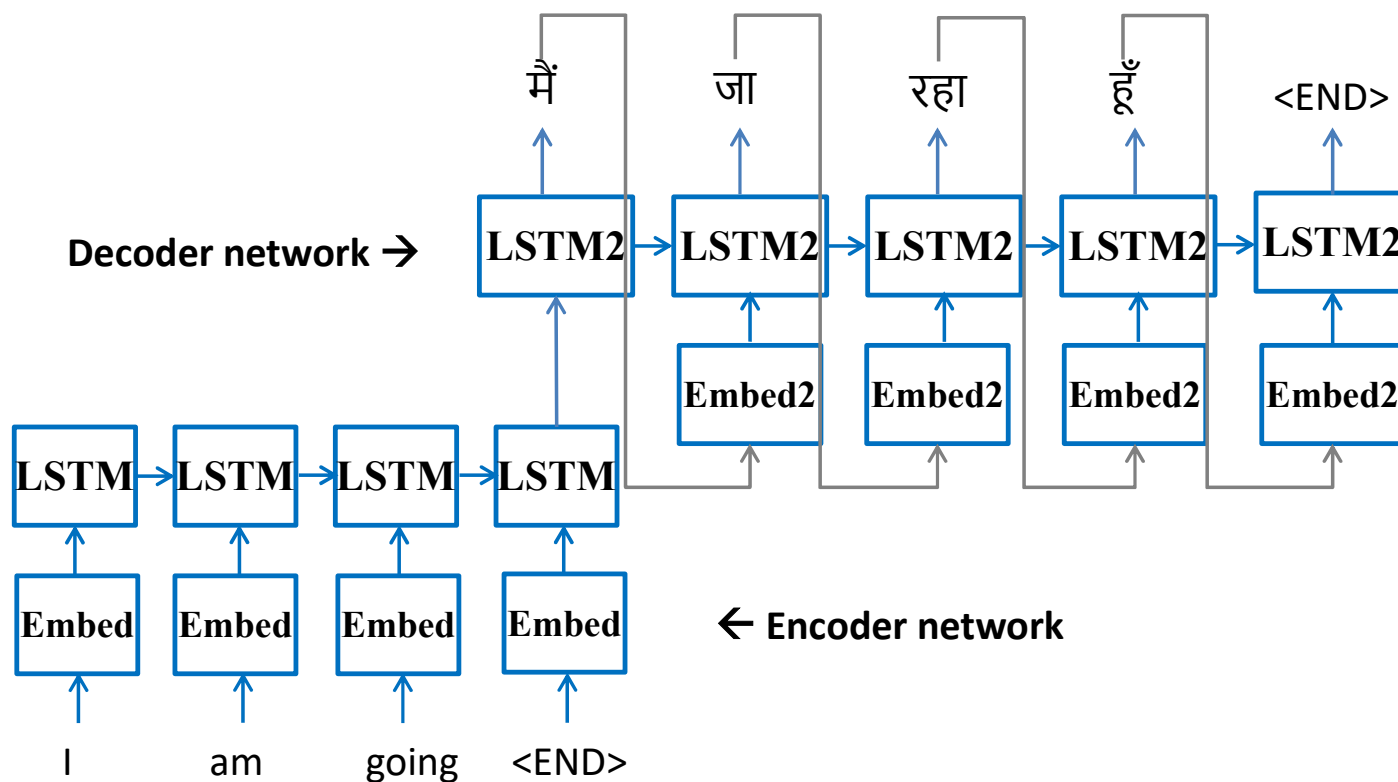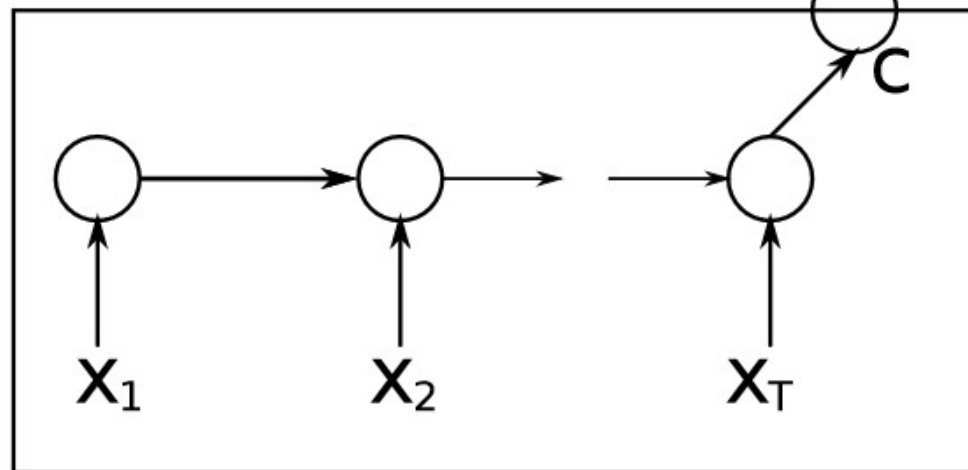
# Machine translation
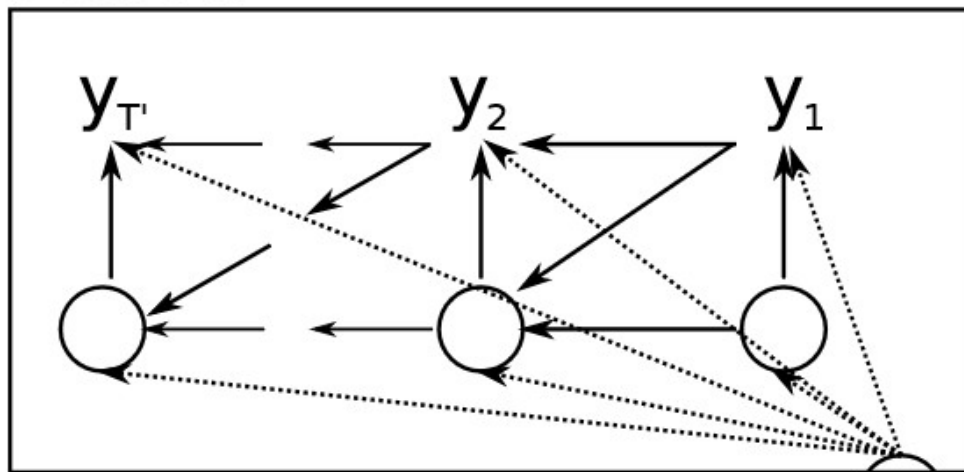
- In actuality, one would use separate LSTMs pre-trained on two different languages

# Machine translation using encoder-decoder

# Bi-directional LSTM

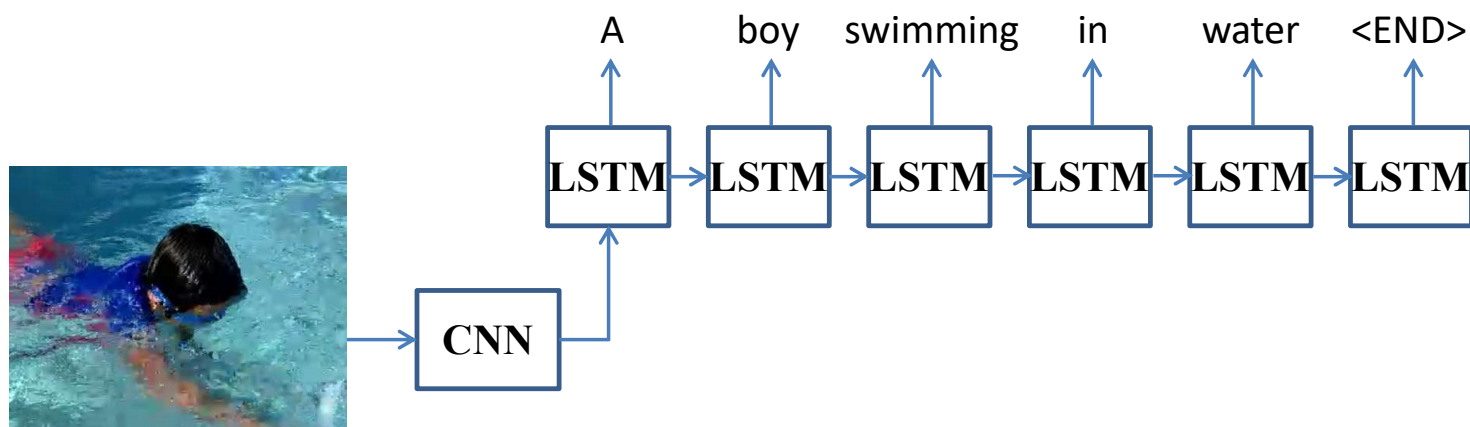- Many problems require a reverse flow of information as well
- For example, POS tagging may require context from future words

# Sentence generation

- Very common for image captioning

- Input is given only in the beginning

- This is a one-to-many task

# Sentence generation

- Very common for image captioning

- Input is given only in the beginning

- This is a one-to-many task

# Video Caption Generation

# Information bottlenech in the context vector

The entire info about the previous sentence has to flow through this context vector *c*

**Decoder network →**

| मैं | जा | रहा | हूँ | <END> |

LSTM2 → LSTM2 → LSTM2 → LSTM2 → LSTM2

Embed2 | Embed2 | Embed2 | Embed2

**← Encoder network**

LSTM → LSTM → LSTM → LSTM

Embed | Embed | Embed | Embed

I | am | going | <END>

# Pooling summarizes feature vectors over several locations

- Average pooling: $1/N \sum_i x_i$

- Max pooling: $\max_i (x_i)$

- Attention: $\sum_i a_i x_i$
    - Often, $a_i = \text{softmax}(e_i)$
    - And $e_i = f(q, k_i)$

# Attention is like intelligent pooling

- E.g. kernel methods:   $y = \sum_i k(x,x_i)\, y_i$
  - Query is $x$
  - Key is $x_i$
  - Value is $y_i$
  - Permutation-invariant in $x_i$
- In softmax attention, the sum of the attentions is *1*

# Side-branches in neural networks



Next output

Next output

Next output

C

Feature Extraction 1

Feature Extraction 2

Feature Extraction

Feature Extraction

Attention

Previous output

Previous output

Previous output

Concatenation

Residual

Basic Attention

# LSTM with Attention Mechanism

# LSTM with Attention Mechanism

$p(y_t \mid \{y_1, \cdots, y_{t-1}\}, c) = g(y_{t-1}, s_t, c)$

**Context**

$c = q(\{h_1, \cdots, h_{Tx}\})$

$h_t = f(x_t, h_{t-1})$

मैं  जा  रहा  हूँ  <END>

I    am   going   <END>

$x_t$

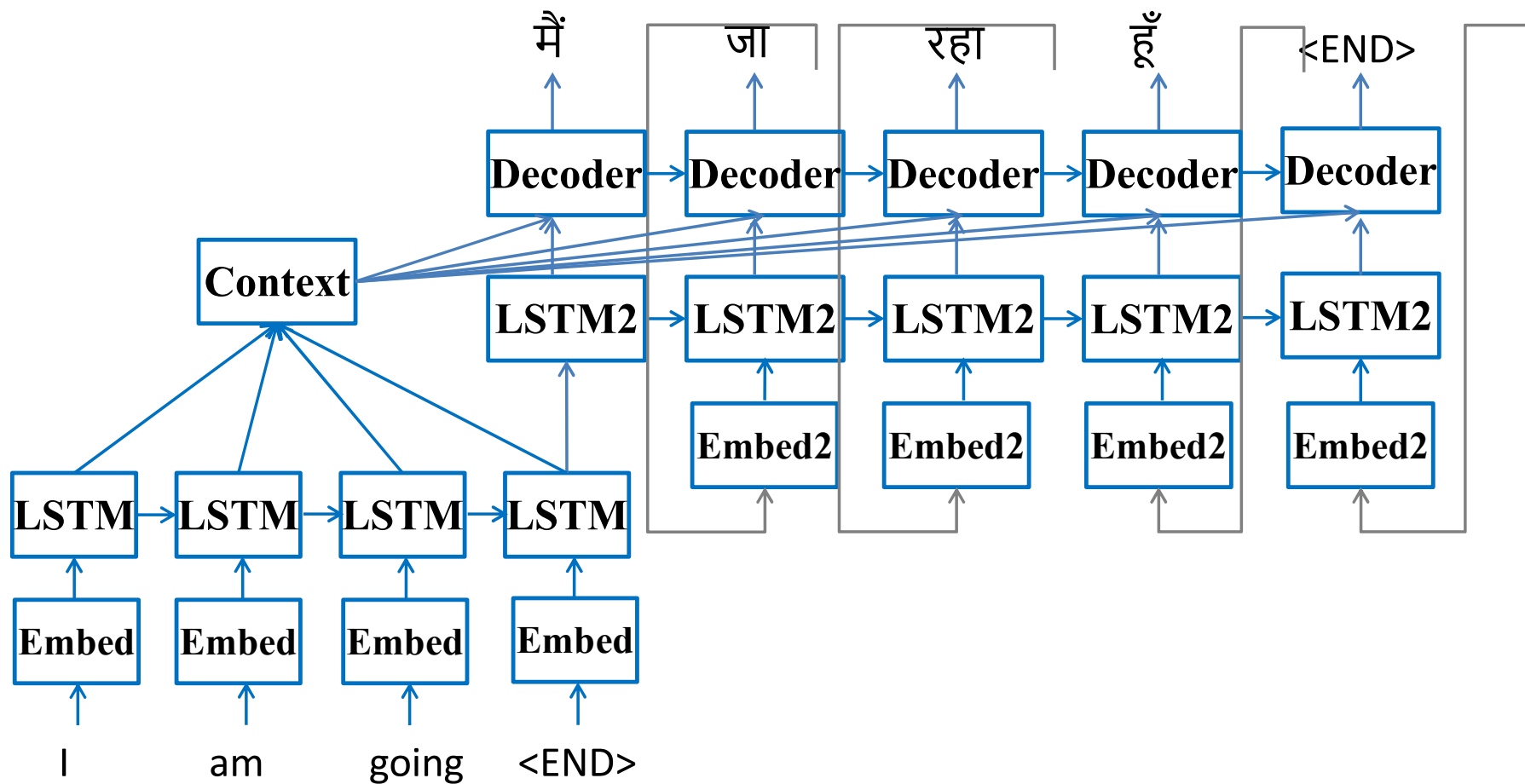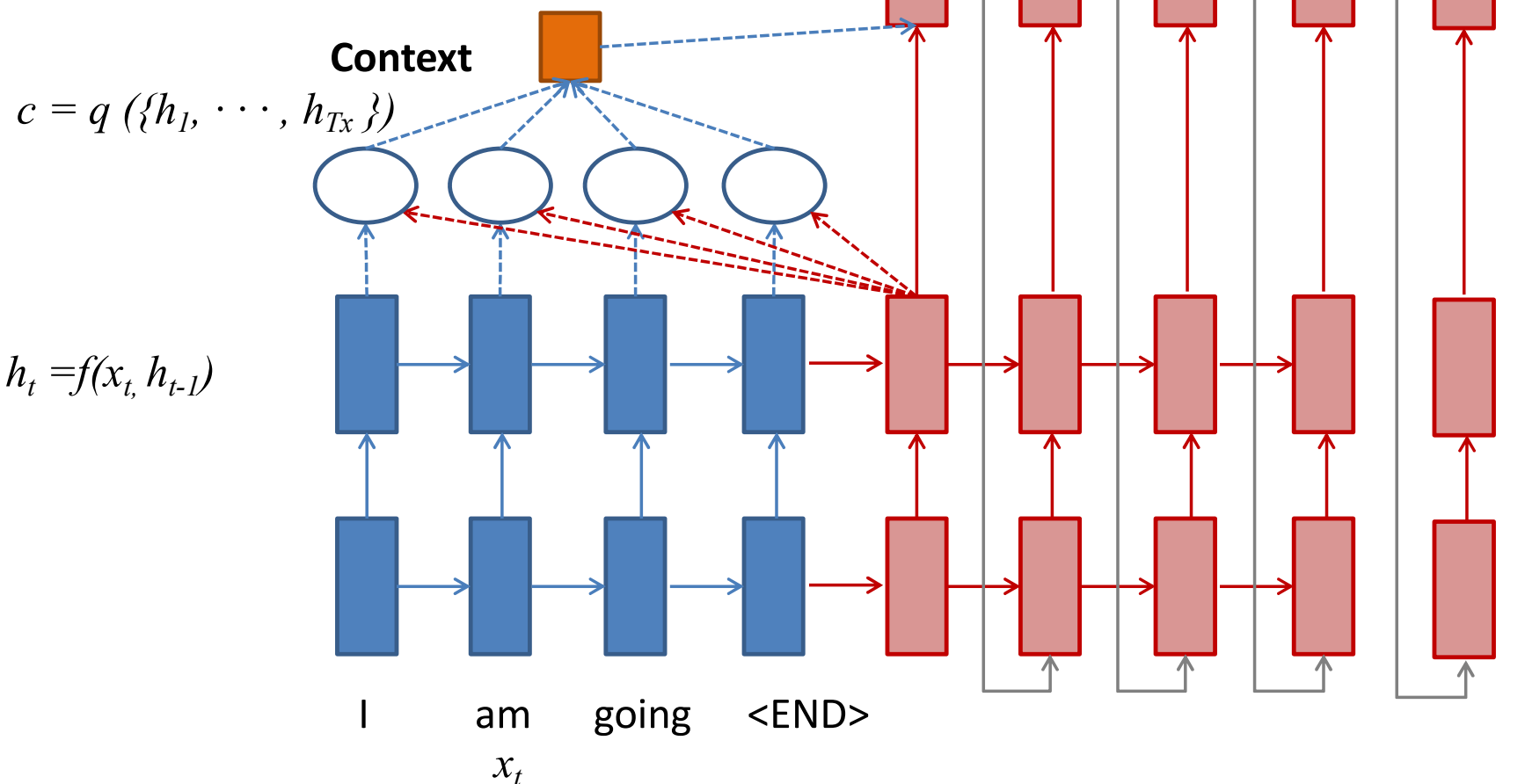"Neural Machine Translation by Jointly Learning to Align and Translate" by Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio, ICLR 2015

# LSTM with Attention Mechanism

$$p(y_i \,|\, y_1, \ldots, y_{i-1}, \boldsymbol{x}) = g(y_{i-1}, s_i, c_i)$$

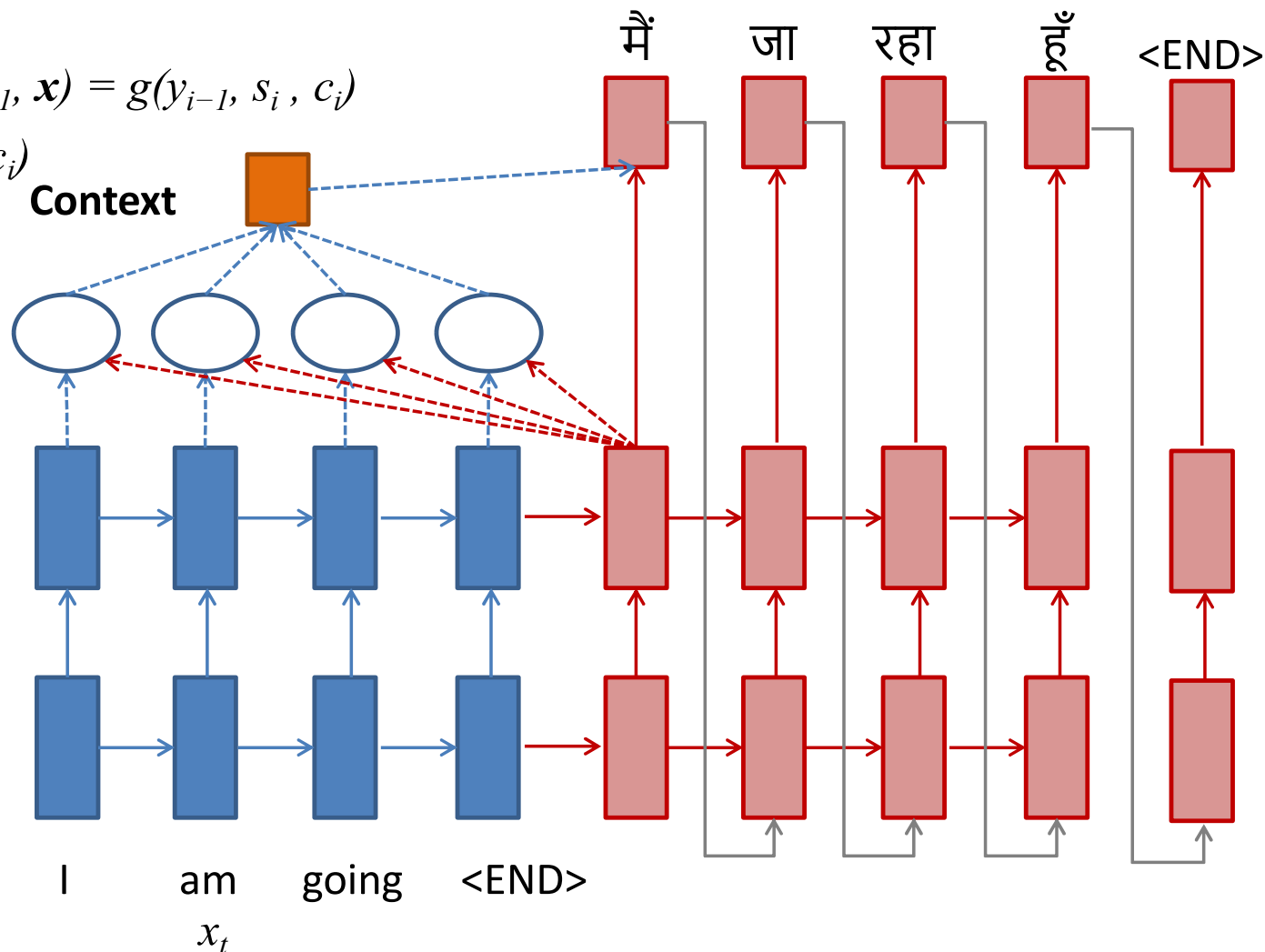$$s_i = f(s_{i-1}, y_{i-1}, c_i)$$

**Context**

$$c_i = \sum_j \alpha_{ij} h_j$$
$$\alpha_{ij} = softmax(e_{ij})$$
$$e_{ij} = a(s_{i-1}, h_j)$$

$$hj = f(xj, h_{j-1})$$

मैं    जा    रहा    हूँ    <END>

I    am    going    <END>
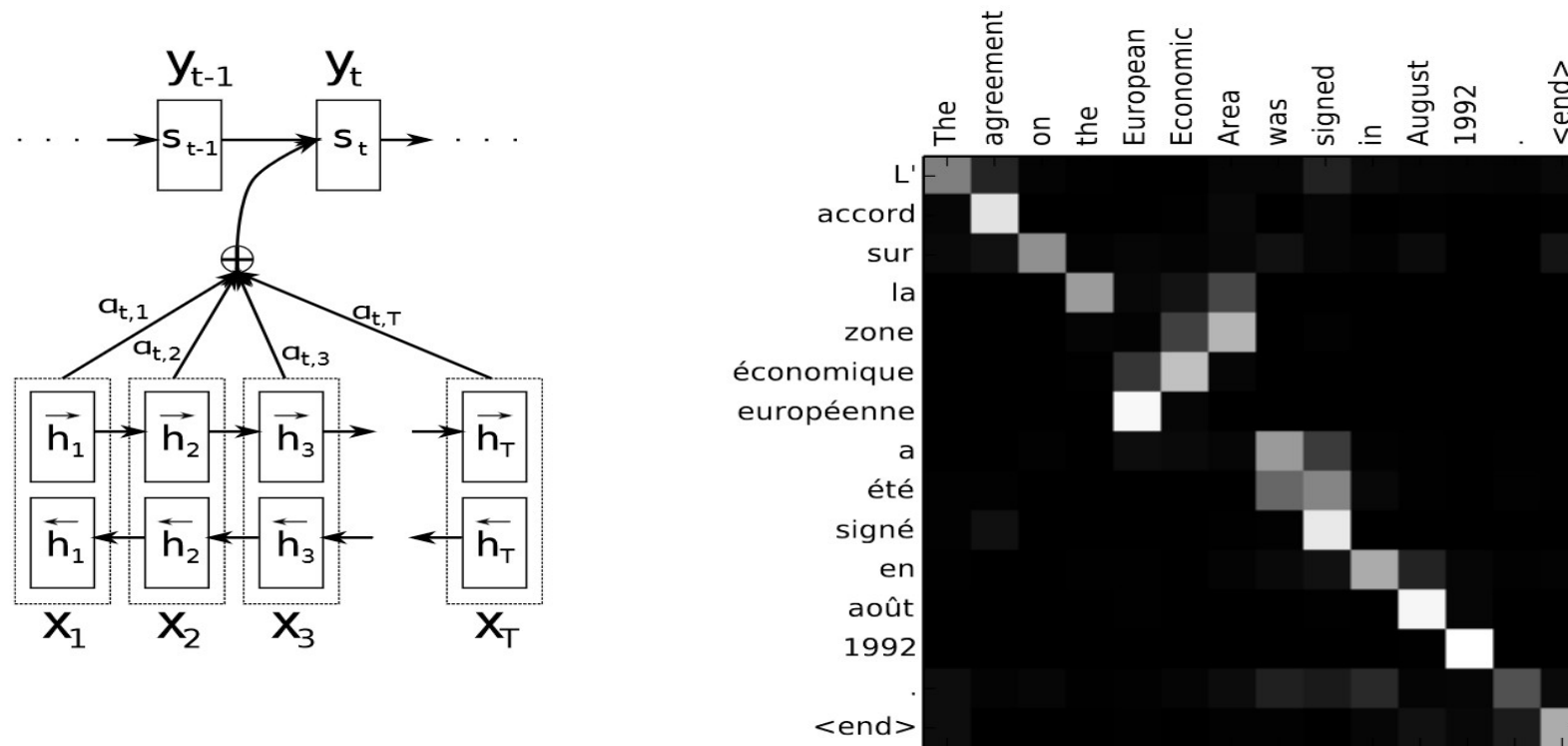
$$x_t$$

"Neural Machine Translation by Jointly Learning to Align and Translate" by Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio, ICLR 2015

# Attention between encoder and decoder

# How attention changes information flow

- Previously

$$h_t = f\left(x_t, h_{t-1}\right)$$

$$c = q\left(\{h_1, \cdots, h_{T_x}\}\right)$$

$$q\left(\{h_1, \cdots, h_T\}\right) = h_T$$

$$p(\mathbf{y}) = \prod_{t=1}^{T} p(y_t \mid \{y_1, \cdots, y_{t-1}\}, c)$$

$$p(y_t \mid \{y_1, \cdots, y_{t-1}\}, c) = g(y_{t-1}, s_t, c)$$



- With attention

$$p(y_i|y_1, \ldots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i)$$

$$s_i = f(s_{i-1}, y_{i-1}, c_i)$$

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j \qquad \alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

$$e_{ij} = a(s_{i-1}, h_j)$$

# No vs. global vs. local attention



Source: *"Effective Approaches to Attention-based Neural Machine Translation," by Luong, Pham, Manning, 2015*
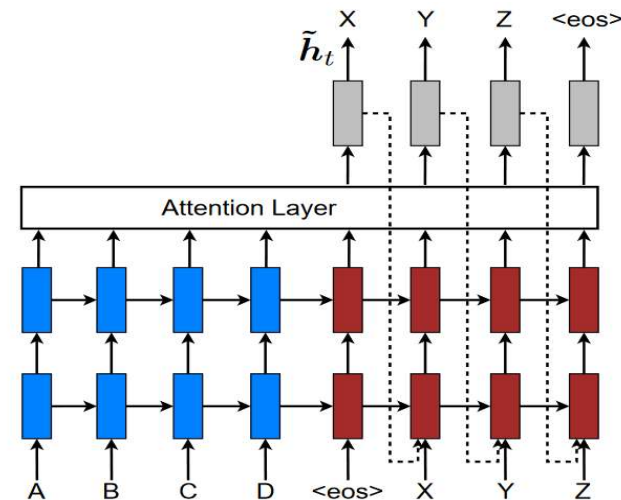
# Transformer networks

# Attention in Transformer networks



Scaled Dot-Product Attention

Multi-Head Attention

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

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

*Source: "Attention Is All You Need," by Vaswani et al., 2017*

# Details of transfomer by Vaswani et al. (2017)

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

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

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

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

*Source: "Attention Is All You Need," by Vaswani et al., 2017*

# Layer norm visualized



"Batch Normalization" by Ioffe, Szegedy, 2015
"Layer Normalization" by Ba, Kiros, Hinton, 2016
"Group Normalization" by Wu, He, 2017

# BERT



Source: "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," by Devli et al., 2018

# BERT

# BERT

| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

*Source: "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," by Devli et al., 2018*

# Vision Transformer



**Vision Transformer (ViT)**

**Transformer Encoder**

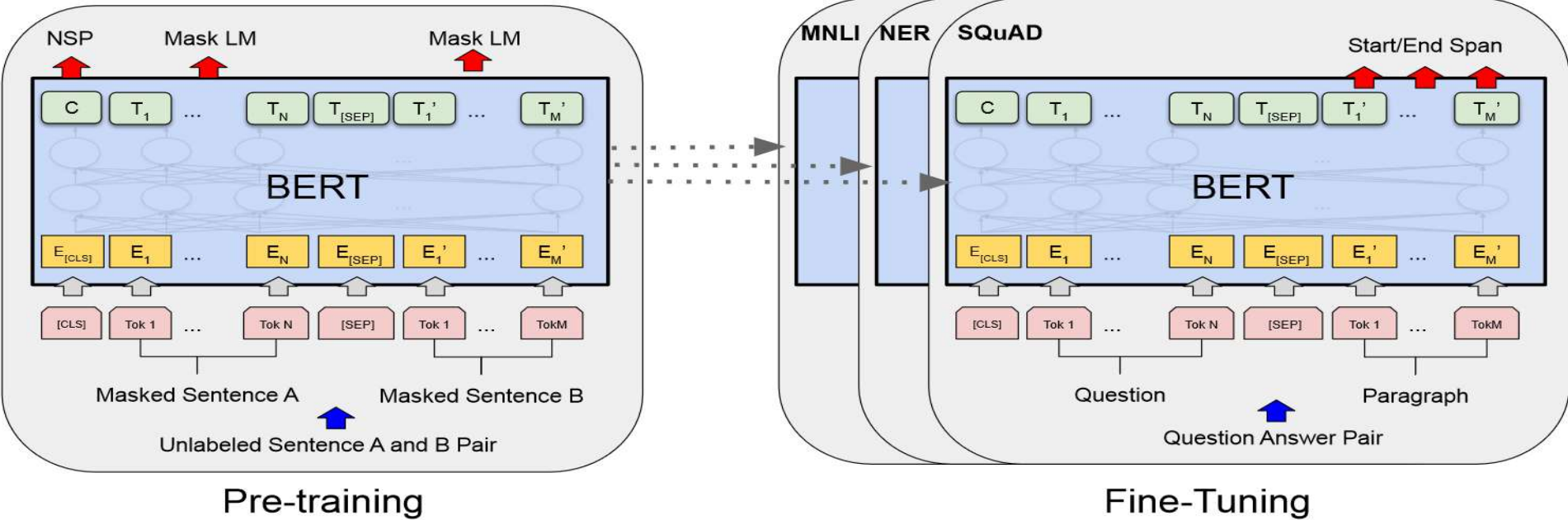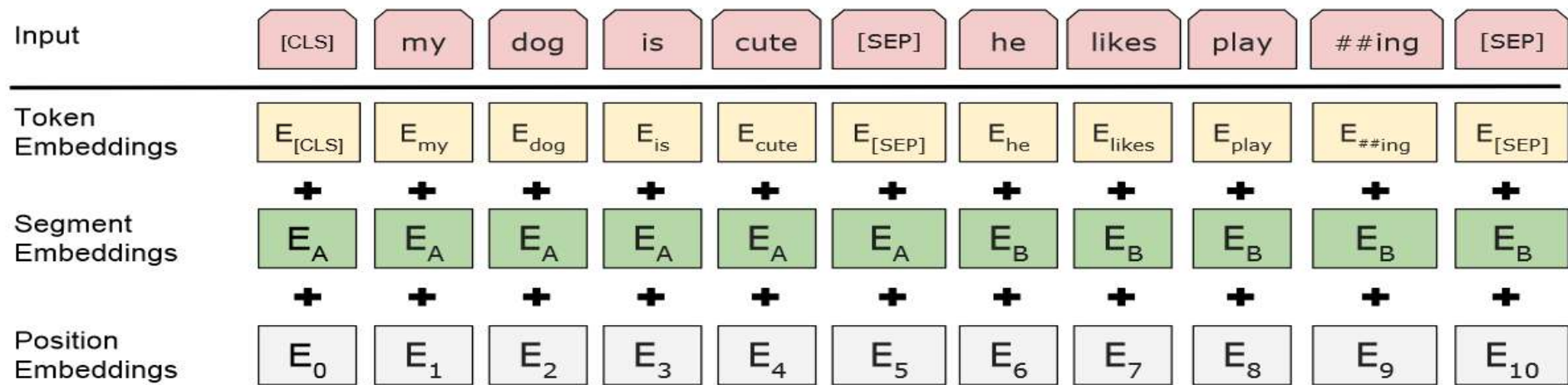Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable "classification token" to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).

"An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," Dosovitskiy et al. ICLR 2021 https://arxiv.org/pdf/2010.11929.pdf

# Vision Transformer

$$\mathbf{z}_0 = [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \cdots ; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{pos}, \qquad \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{pos} \in \mathbb{R}^{(N+1) \times D}$$
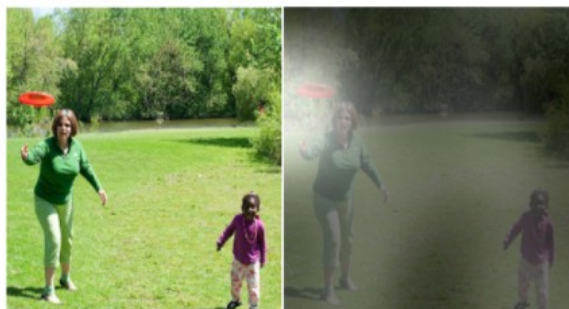
$$\mathbf{z'}_\ell = \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, \qquad \ell = 1 \ldots L$$

$$\mathbf{z}_\ell = \text{MLP}(\text{LN}(\mathbf{z'}_\ell)) + \mathbf{z'}_\ell, \qquad \ell = 1 \ldots L$$

$$\mathbf{y} = \text{LN}(\mathbf{z}_L^0)$$

$$[\mathbf{q}, \mathbf{k}, \mathbf{v}] = \mathbf{z} \mathbf{U}_{qkv} \qquad \mathbf{U}_{qkv} \in \mathbb{R}^{D \times 3D_h},$$

$$A = \text{softmax}\left(\mathbf{q}\mathbf{k}^\top / \sqrt{D_h}\right) \qquad A \in \mathbb{R}^{N \times N},$$
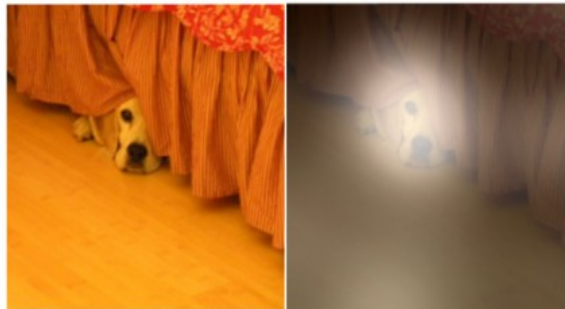
$$\text{SA}(\mathbf{z}) = A\mathbf{v}.$$

$$\text{MSA}(\mathbf{z}) = [\text{SA}_1(z); \text{SA}_2(z); \cdots ; \text{SA}_k(z)] \mathbf{U}_{msa} \qquad \mathbf{U}_{msa} \in \mathbb{R}^{k \cdot D_h \times D}$$

"An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," Dosovitskiy et al. ICLR 2021 https://arxiv.org/pdf/2010.11929.pdf

# Association of words to image regions
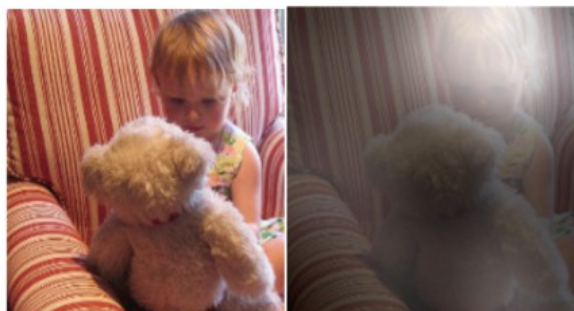


A woman is throwing a <u>frisbee</u> in a park.

A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

A little <u>girl</u> sitting on a bed with a teddy bear.

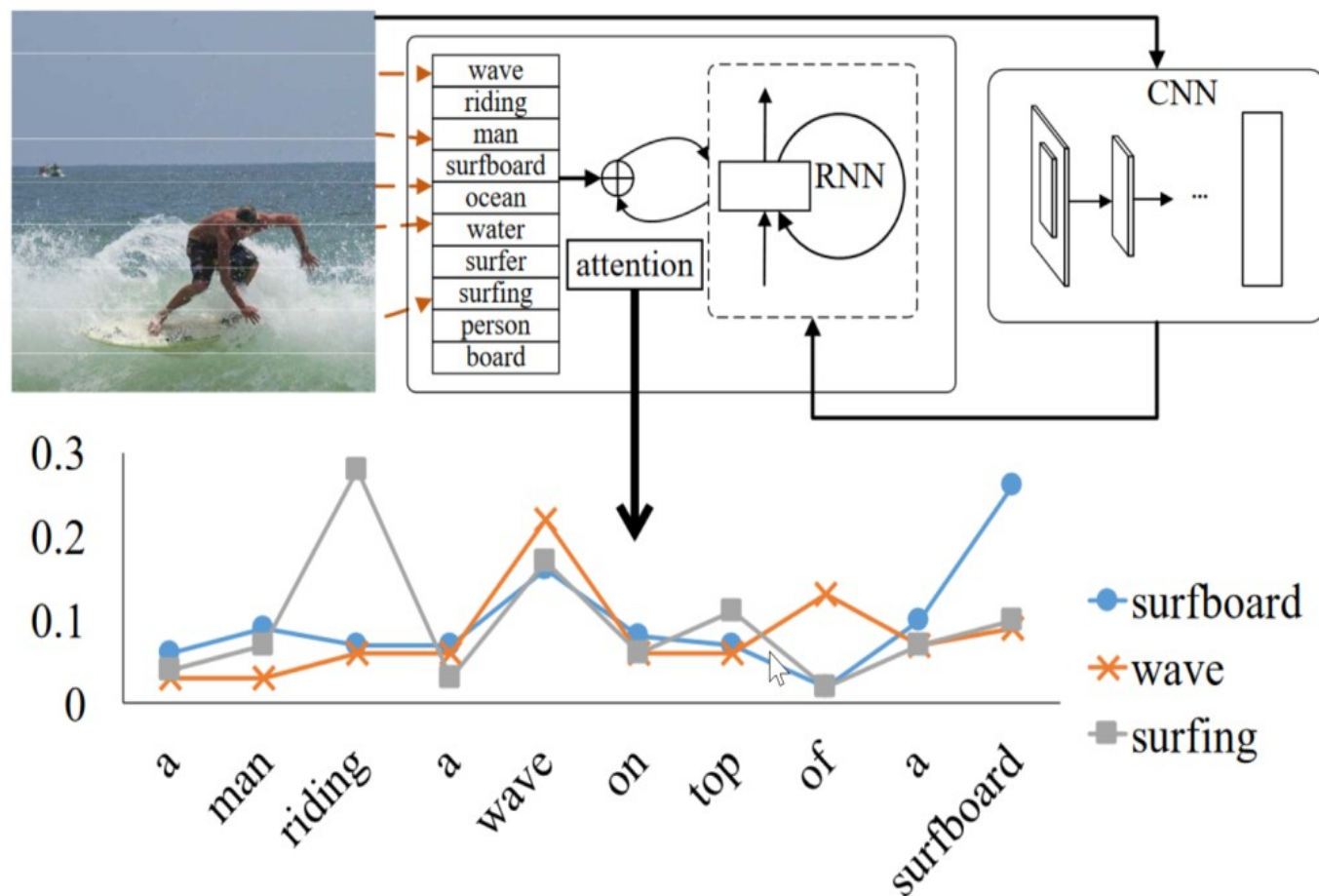A group of <u>people</u> sitting on a boat in the water.

A giraffe standing in a forest with <u>trees</u> in the background.

Source: Xu, Kelvin, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. "Show, attend and tell: Neural image caption generation with visual attention." In *International conference on machine learning*, pp. 2048-2057. 2015.

# Injecting attribute extraction into the attention process



Source: You Q, Jin H, Wang Z, Fang C, Luo J. Image captioning with semantic attention. InProceedings of the IEEE conference on computer vision and pattern recognition 2016 (pp. 4651-4659).
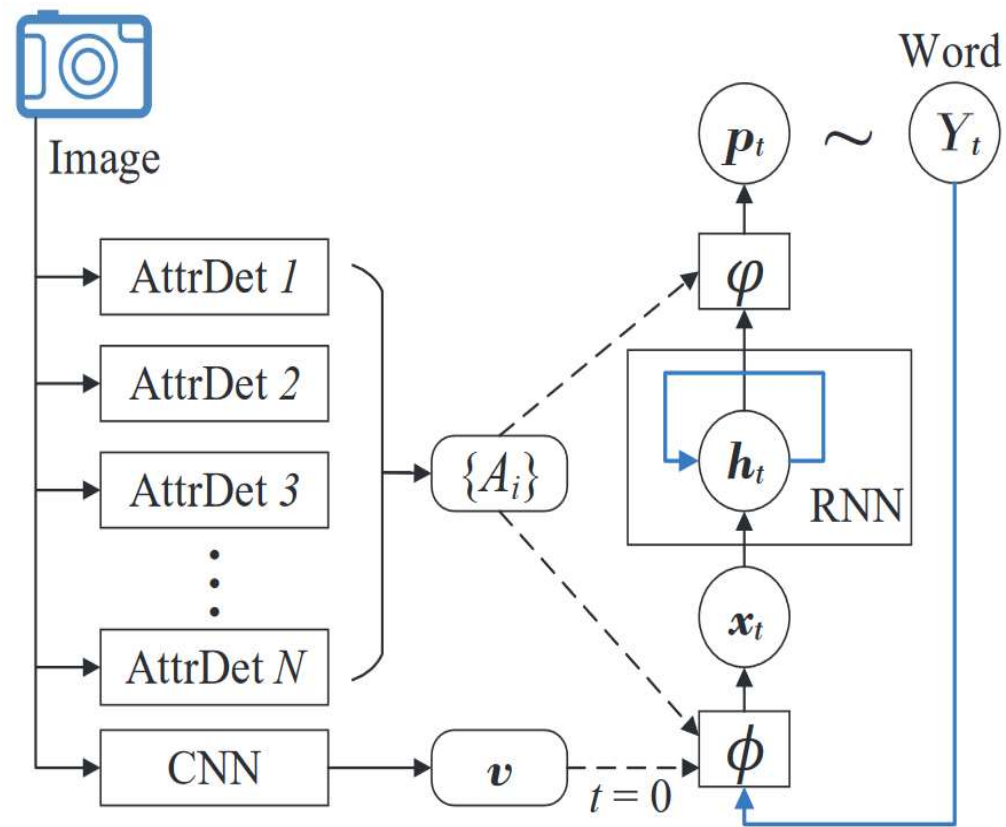
# Injecting attribute extraction into the attention process



Source: You Q, Jin H, Wang Z, Fang C, Luo J. Image captioning with semantic attention. InProceedings of the IEEE conference on computer vision and pattern recognition 2016 (pp. 4651-4659).
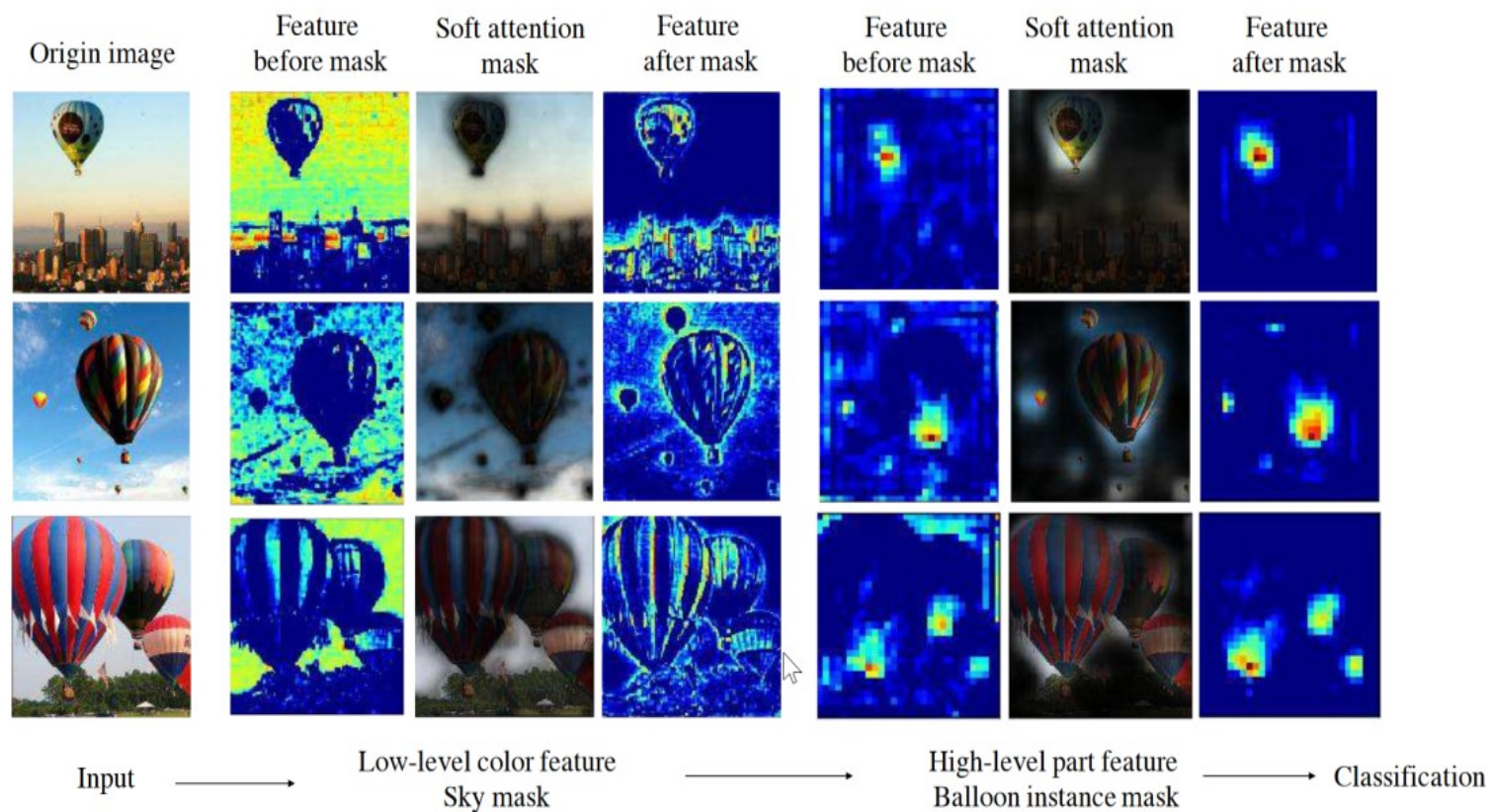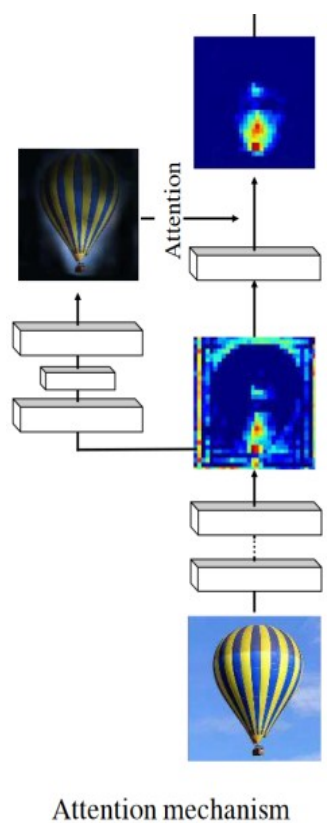
# Outline

- Attention in NLP
- **Attention in vision**
  - Attention for image captioning
  - **Attention for image recognition**
  - Attention for segmentation and detection
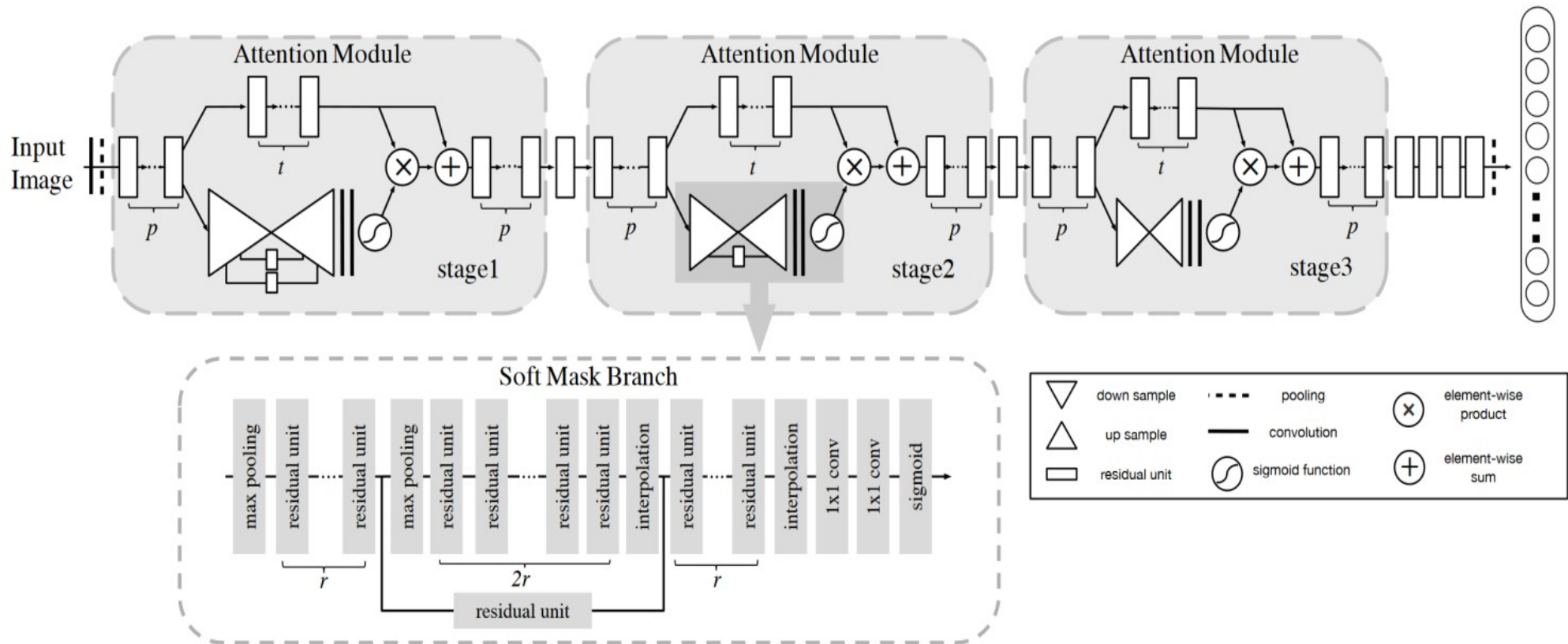  - Attention in other vision applications

# Residual attention network



Source: Wang F, Jiang M, Qian C, Yang S, Li C, Zhang H, Wang X, Tang X. Residual attention network for image classification. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition 2017 (pp. 3156-3164).

# Two branches of residual attention network



**Code snippet (GitHub: koichiro11 / residual-attention-network):**

with tf.variable_scope("attention"):
  output = (1 + output_soft_mask) * output_trunk
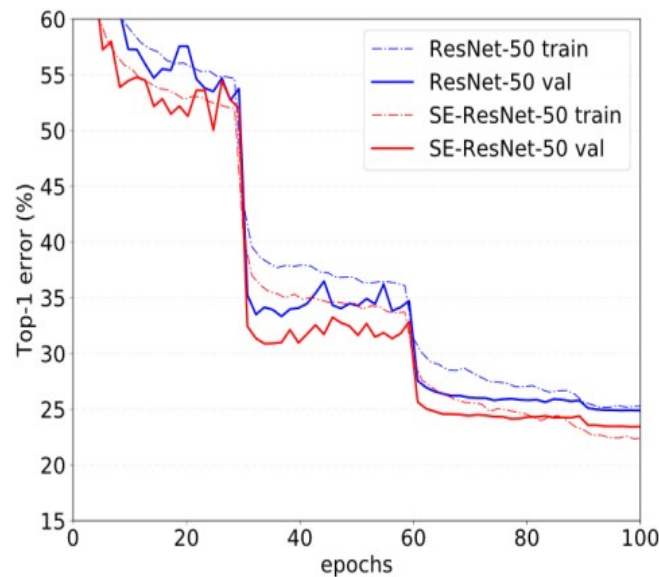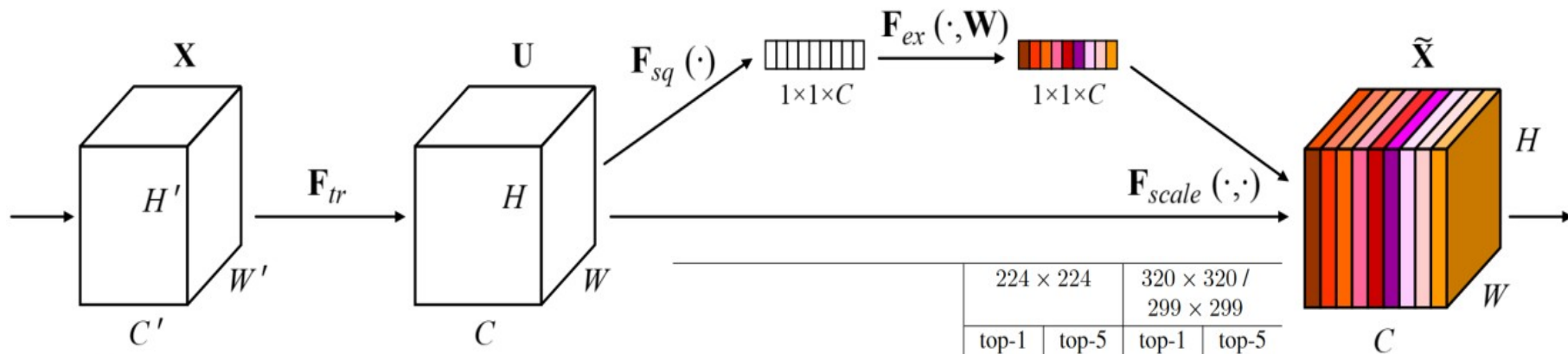
with tf.variable_scope("output"):
  output_soft_mask = tf.layers.conv2d…
  output_soft_mask = tf.nn.sigmoid(output_soft_mask)

Paper: Wang F, Jiang M, Qian C, Yang S, Li C, Zhang H, Wang X, Tang X. Residual attention network for image classification. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition 2017 (pp. 3156-3164).
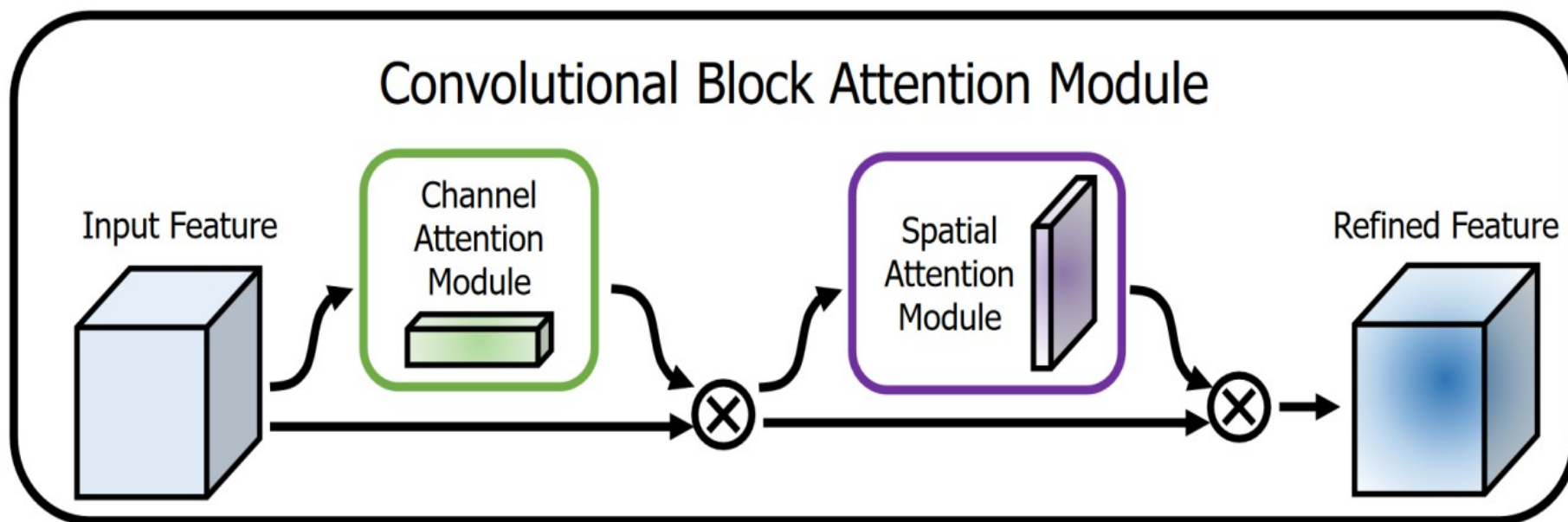
# Squeeze and excitation networks



| | | 224 × 224 | | 320 × 320 / 299 × 299 | |
|---|---|---|---|---|---|
| | | top-1 err. | top-5 err. | top-1 err. | top-5 err. |
| ResNet-152 [10] | | 23.0 | 6.7 | 21.3 | 5.5 |
| ResNet-200 [11] | | 21.7 | 5.8 | 20.1 | 4.8 |
| Inception-v3 [44] | | - | - | 21.2 | 5.6 |
| Inception-v4 [42] | | - | - | 20.0 | 5.0 |
| Inception-ResNet-v2 [42] | | - | - | 19.9 | 4.9 |
| ResNeXt-101 (64 × 4d) [47] | | 20.4 | 5.3 | 19.1 | 4.4 |
| DenseNet-264 [14] | | 22.15 | 6.12 | - | - |
| Attention-92 [46] | | - | - | 19.5 | 4.8 |
| Very Deep PolyNet [51] † | | - | - | 18.71 | 4.25 |
| PyramidNet-200 [8] | | 20.1 | 5.4 | 19.2 | 4.7 |
| DPN-131 [5] | | 19.93 | 5.12 | 18.55 | 4.16 |
| **SENet-154** | | **18.68** | **4.47** | **17.28** | **3.79** |
| NASNet-A (6@4032) [55] † | | - | - | 17.3‡ | 3.8‡ |
| **SENet-154 (post-challenge)** | | - | - | **16.88‡** | **3.58‡** |

Source: Hu J, Shen L, Sun G. Squeeze-and-excitation networks. InProceedings of the IEEE conference on computer vision and pattern recognition 2018 (pp. 7132-7141).

# Convolutional block attention module



Convolutional Block Attention Module

Source: Woo S, Park J, Lee JY, So Kweon I. Cbam: Convolutional block attention module. InProceedings of the European Conference on Computer Vision (ECCV) 2018 (pp. 3-19).
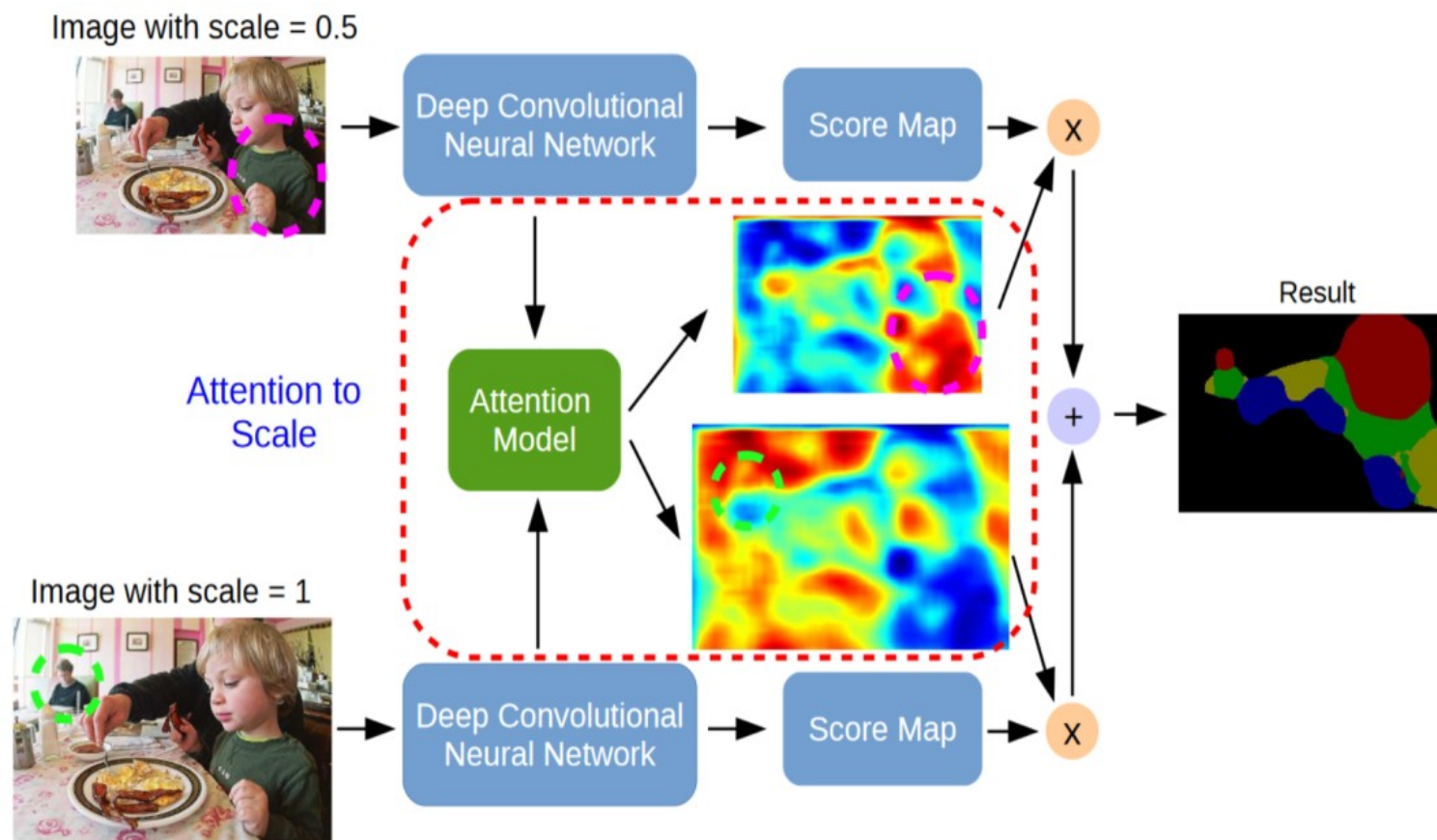
# Outline

- Attention in NLP
- **Attention in vision**
  - Attention for image captioning
  - Attention for image recognition
  - **Attention for segmentation and detection**
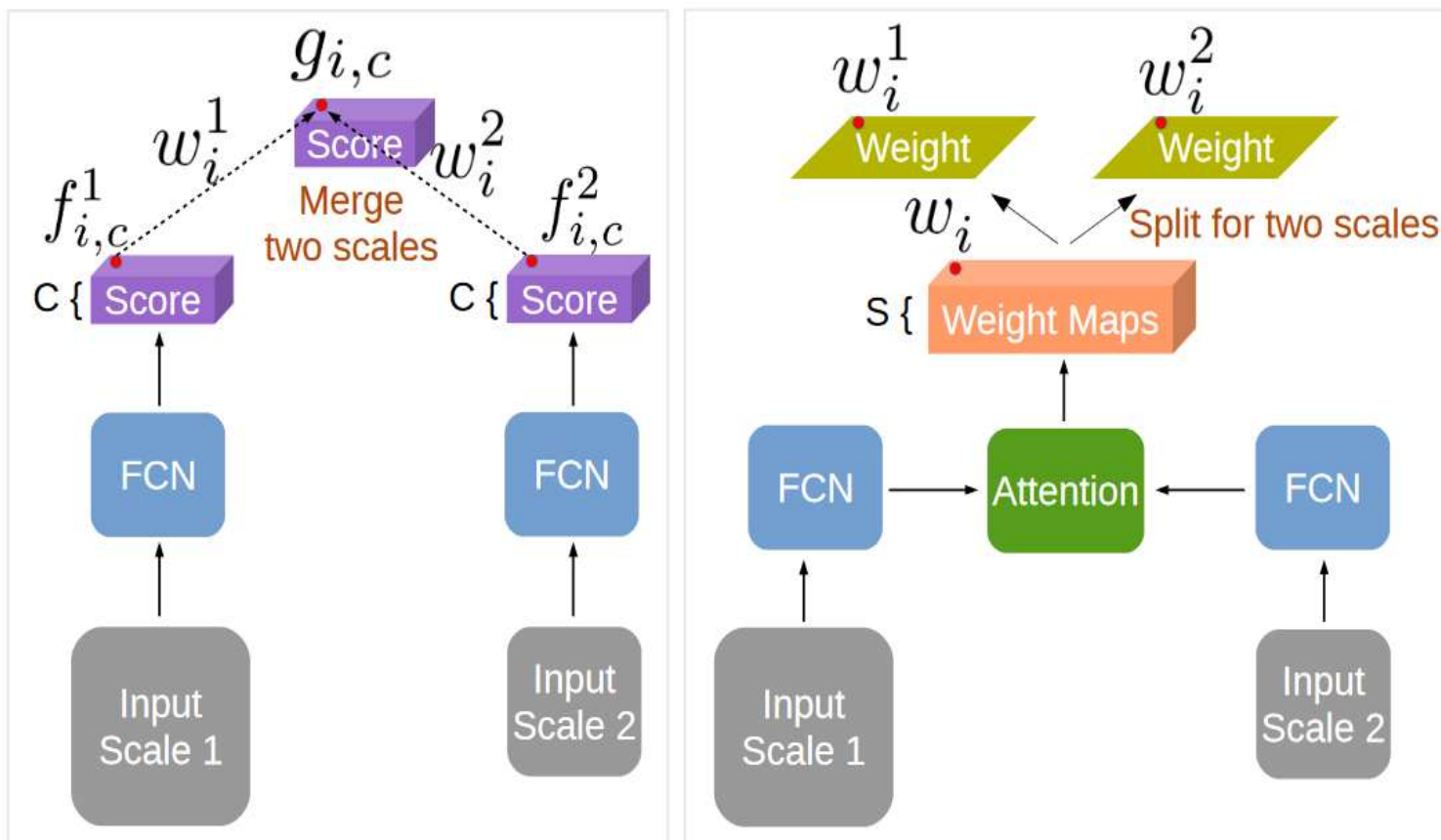  - Attention in other vision applications

# Multi-scale attention for segmentation

Source: Chen LC, Yang Y, Wang J, Xu W, Yuille AL. Attention to scale: Scale-aware semantic image segmentation. InProceedings of the IEEE conference on computer vision and pattern recognition 2016 (pp. 3640-3649).
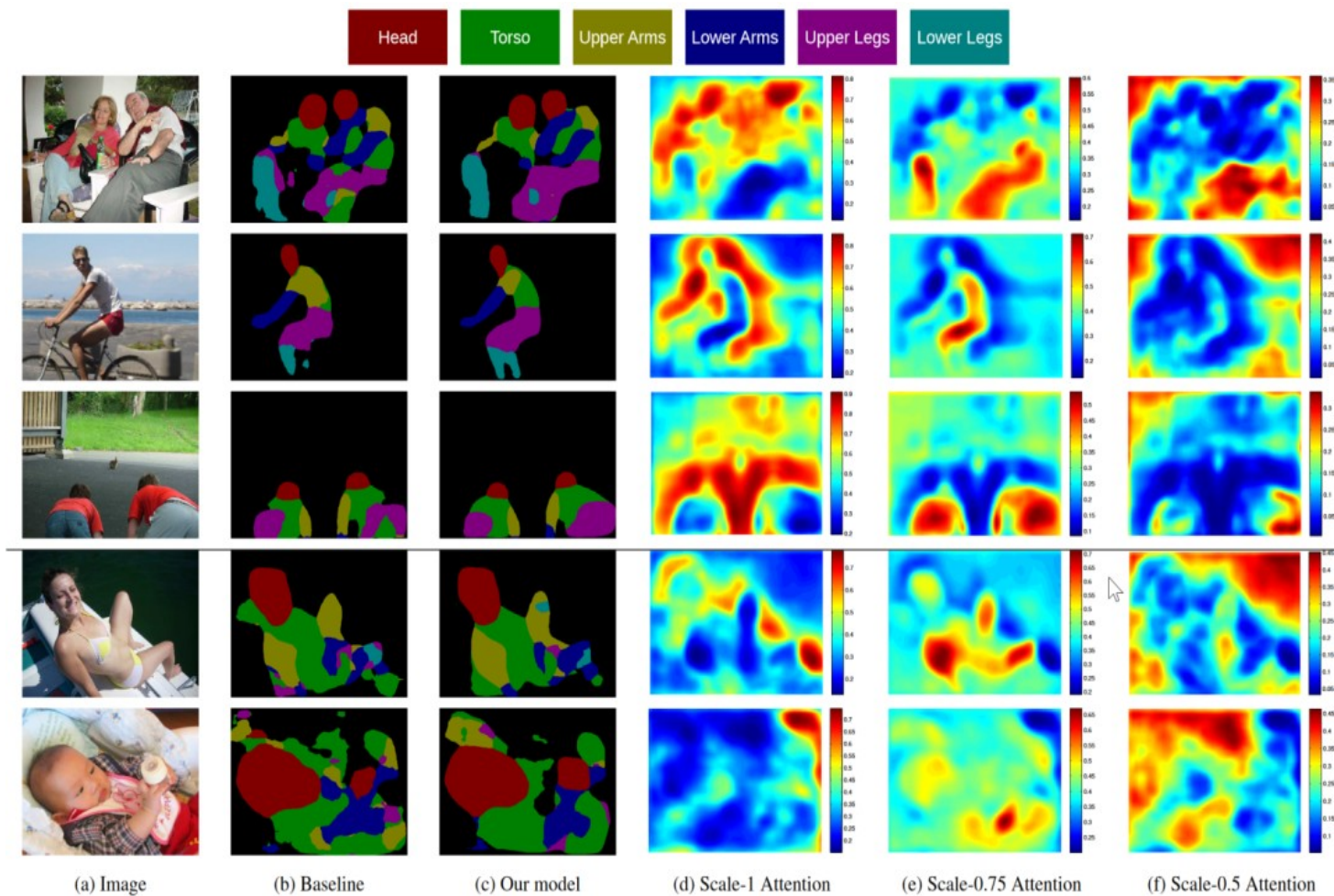
# Multi-scale attention for segmentation



Source: Chen LC, Yang Y, Wang J, Xu W, Yuille AL. Attention to scale: Scale-aware semantic image segmentation. InProceedings of the IEEE conference on computer vision and pattern recognition 2016 (pp. 3640-3649).

# Multi-scale attention for segmentation



| Head | Torso | Upper Arms | Lower Arms | Upper Legs | Lower Legs |

(a) Image   (b) Baseline   (c) Our model   (d) Scale-1 Attention   (e) Scale-0.75 Attention   (f) Scale-0.5 Attention
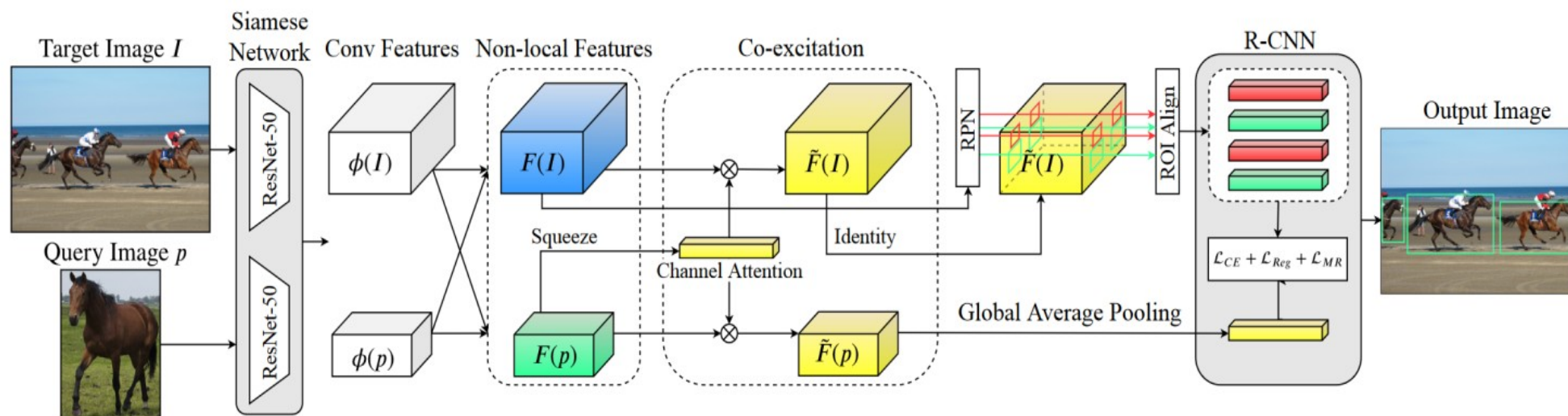
Source: Chen LC, Yang Y, Wang J, Xu W, Yuille AL. Attention to scale: Scale-aware semantic image segmentation. InProceedings of the IEEE conference on computer vision and pattern recognition 2016 (pp. 3640-3649).

# One-shot object detection



Source: Hsieh TI, Lo YC, Chen HT, Liu TL. One-Shot Object Detection with Co-Attention and Co-Excitation. InAdvances in Neural Information Processing Systems 2019 (pp. 2721-2730).
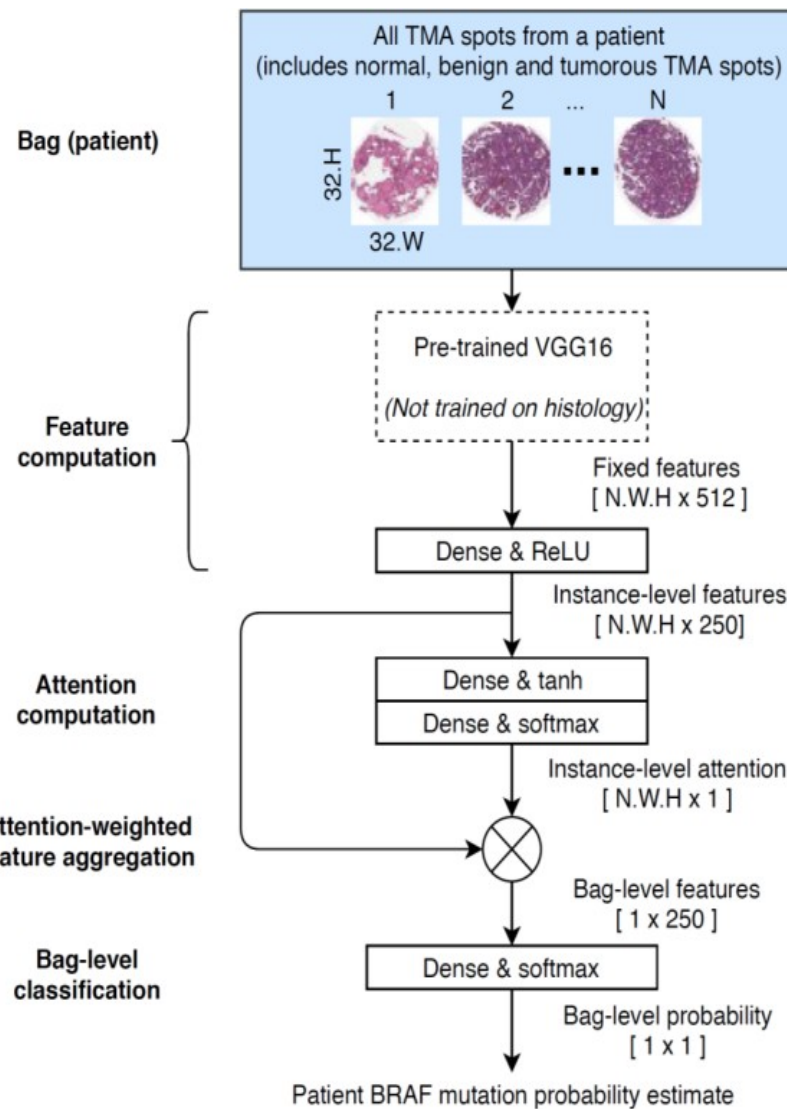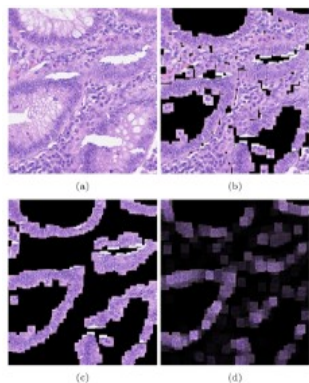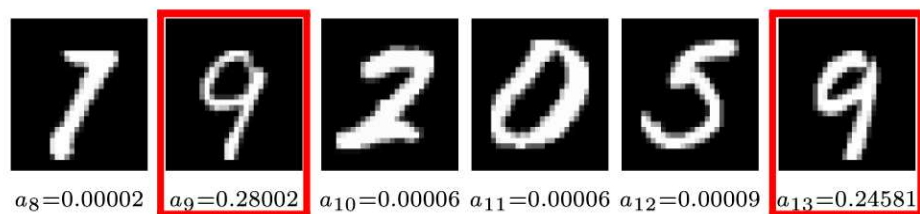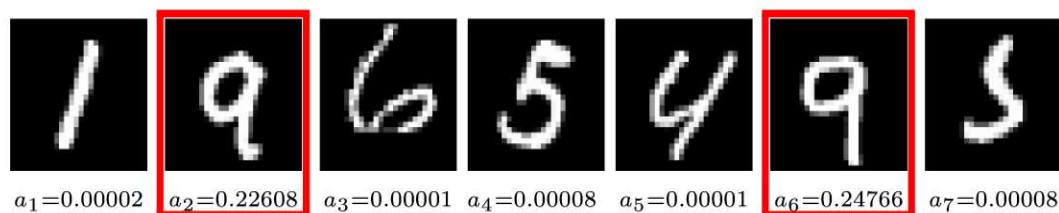
# Outline

- Attention in NLP
- **Attention in vision**
  - Attention for image captioning
  - Attention for image recognition
  - Attention for segmentation and detection
  - **Attention in other vision applications**

# Attention-based Multiple Instance Learning

Source: Ilse M, Tomczak JM, Welling M. Attention-based deep multiple instance learning. In 35th International Conference on Machine Learning, ICML 2018 2018 Jan 1 (pp. 3376-3391).
And, adapted in unpublished work by Deepak Anand, Kumar Yashashwi, Neeraj Kumar, Swapnil Rane, Peter Gann, and Amit Sethi