

# AI · BigData Academy

## NLP Lab. 1

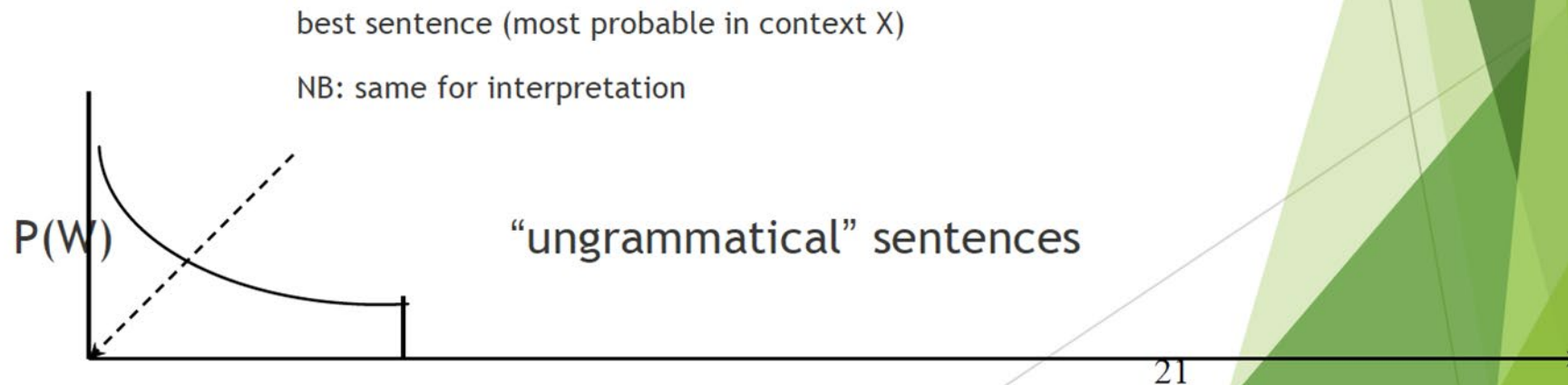
Isoft Lab 백성빈, 이지현, 이채빈

Statistical model, Deep learning model

# Statistical NLP

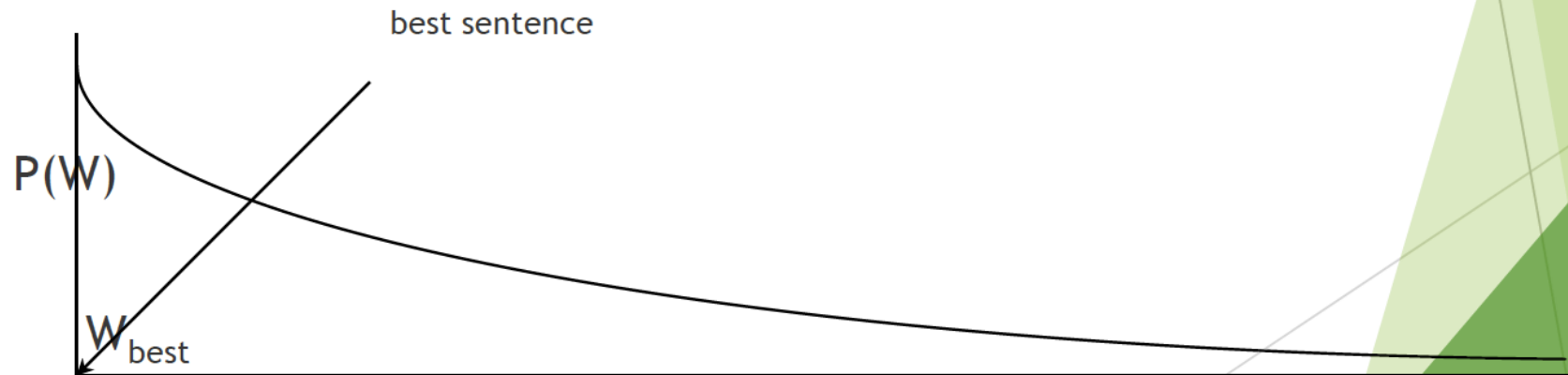
► Imagine:

- Each sentence  $W = \{w_1, w_2, \dots, w_n\}$  gets a probability  $P(W|X)$  in a context  $X$  (think of it in the intuitive sense for now)
- For every possible context  $X$ , sort all the imaginable sentences  $W$  according to  $P(W|X)$ :
- Ideal situation:

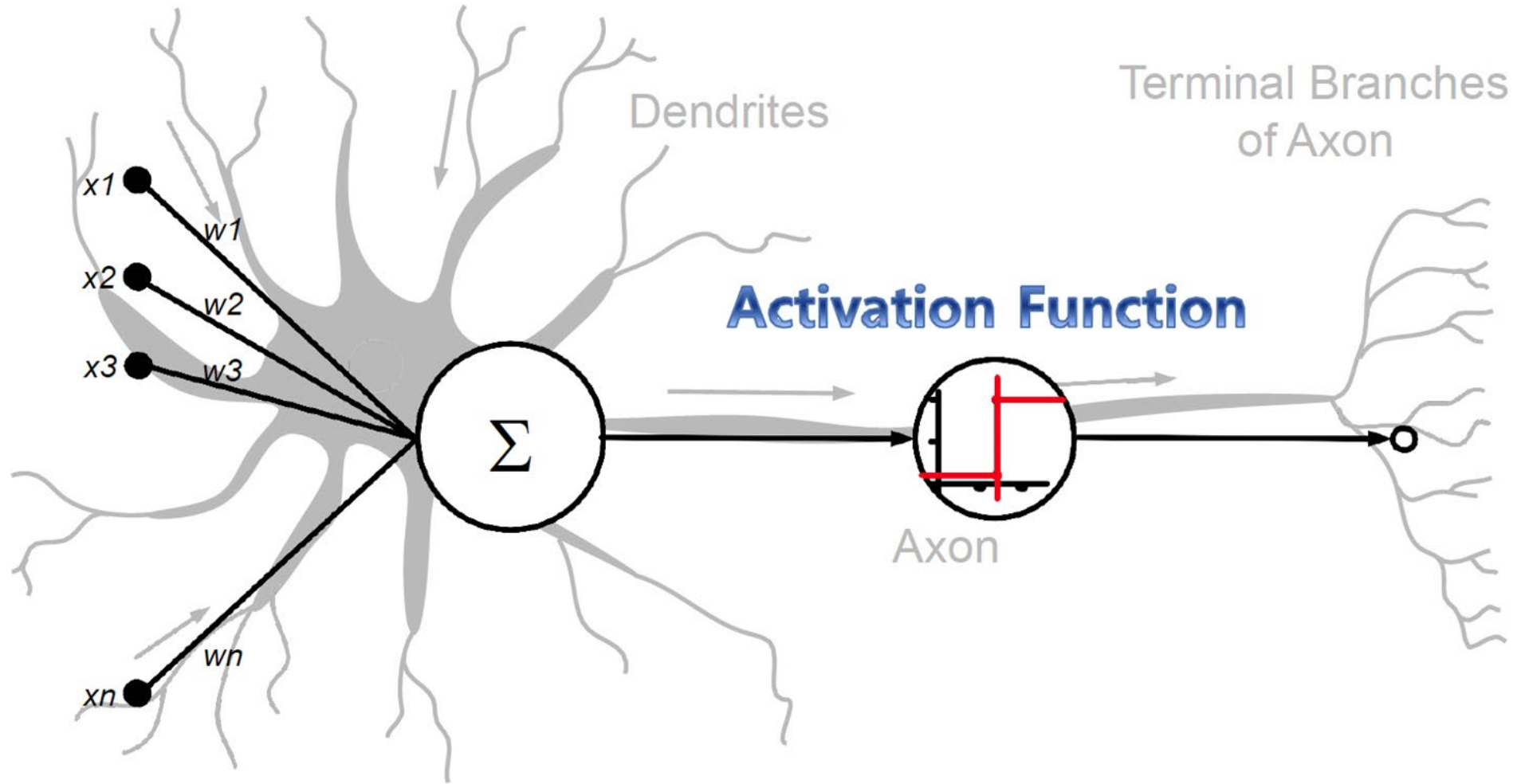


# Real World Situation

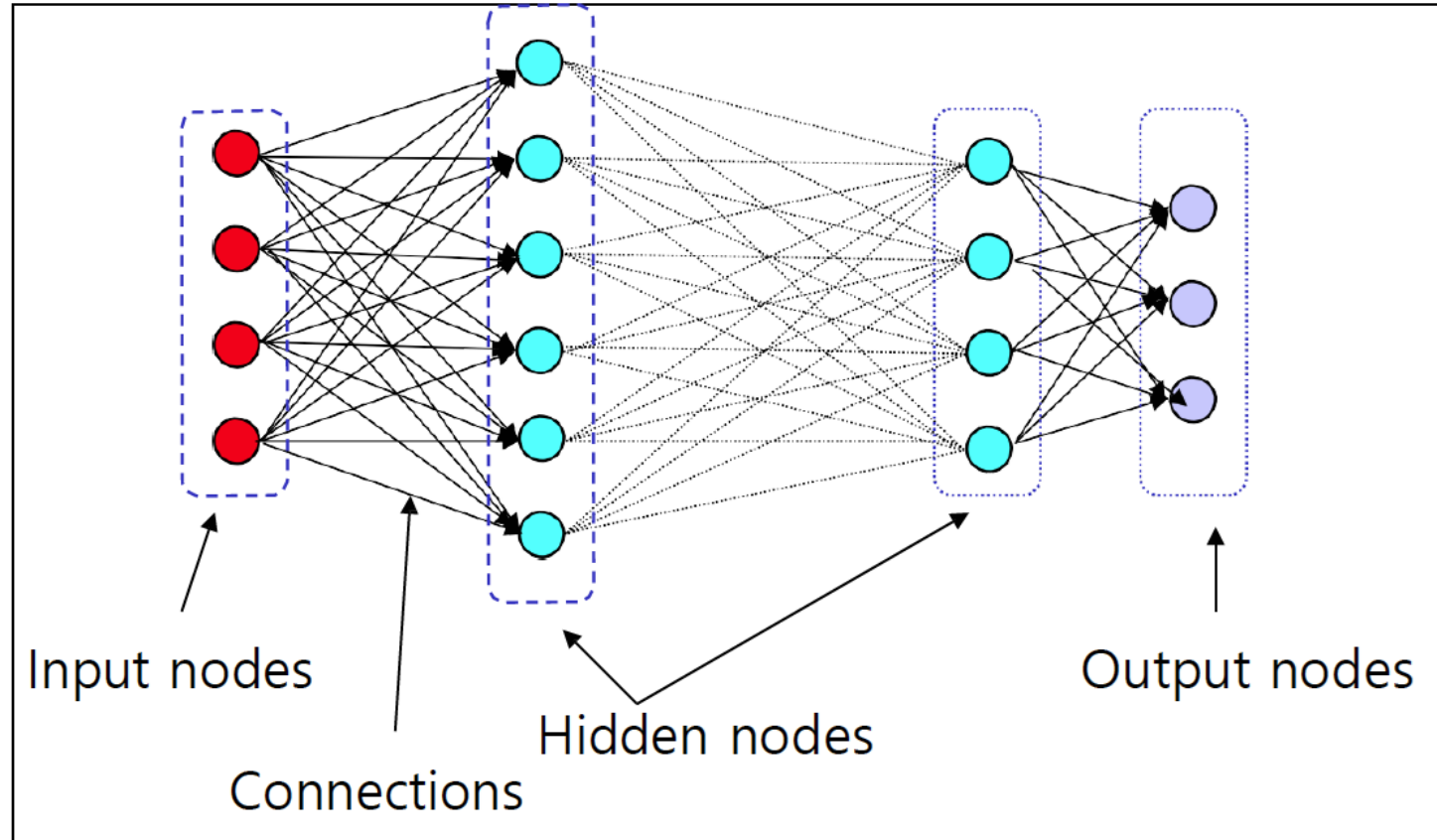
- ▶ Unable to specify set of grammatical sentences today using fixed “categorical” rules (maybe never)
- ▶ Use statistical “model” based on REAL WORLD DATA and care about the best sentence only (disregarding the “grammaticality” issue)



# Artificial Neural Networks (ANN)



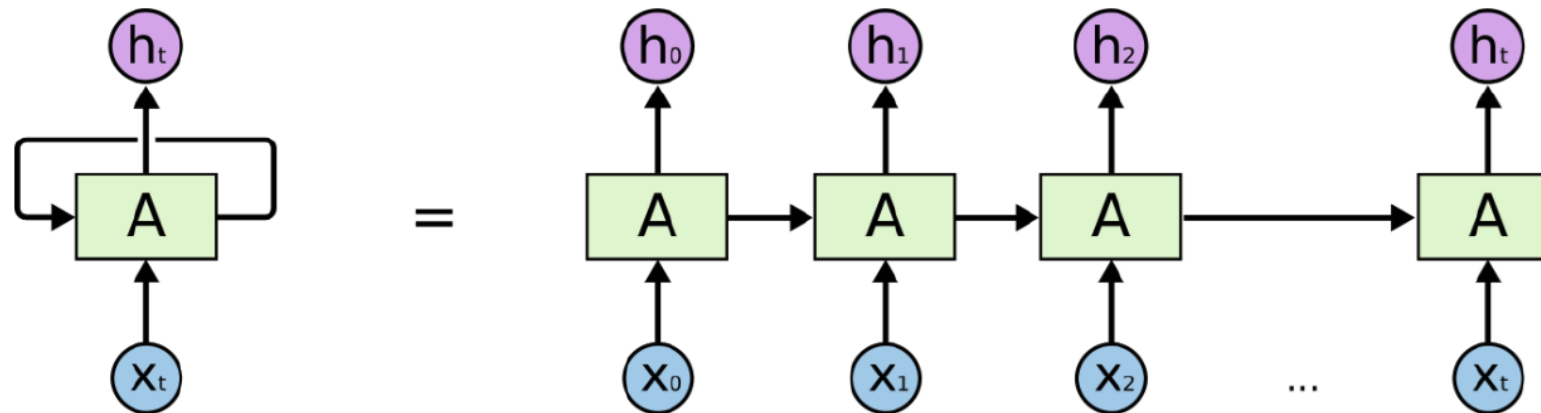
# Layered Networks



$$\begin{aligned} \text{Output : } y &= f(w^1x + w_i^2x + w_i^3x + \cdots + w_i^mx_m) \\ &= f\left(\sum_j w_i^j x_j\right) \end{aligned}$$

# Vanilla recurrent neural networks (RNNs)

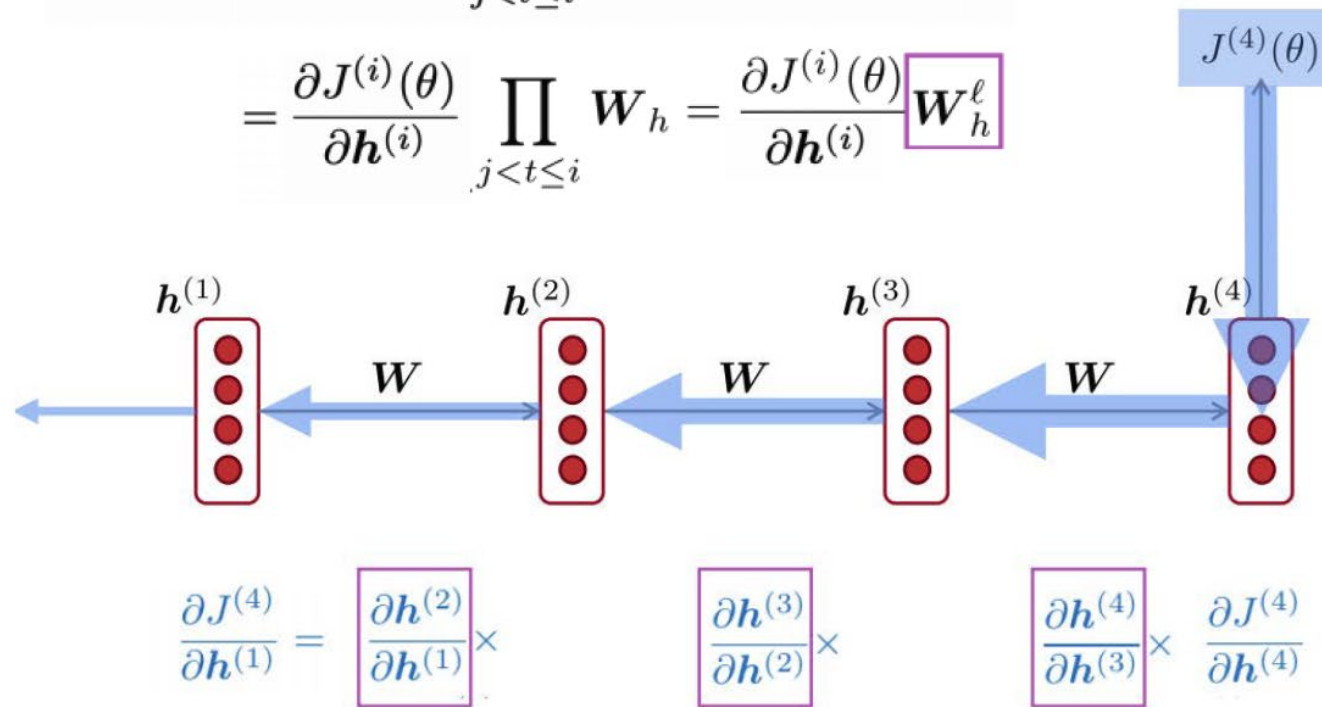
- RNNs have connections from the outputs of previous time steps to inputs of next time steps
- For sequential data, a RNN usually computes hidden state  $h_t$  from the previous hidden state  $h_{t-1}$  and the input  $x_t$ 
  - $h_t = \sigma(W_h h_{t-1} + W_x x_t + b)$



# Vanishing gradient problem

- $h_t = \sigma(W_h h_{t-1} + W_x x_t + b)$
- Let's assume  $\sigma$  is the identity function

$$\begin{aligned}\frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} &= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \\ &= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \mathbf{W}_h = \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \boxed{\mathbf{W}_h^\ell}\end{aligned}$$



$$\text{If all } \frac{\partial \mathbf{h}^t}{\partial \mathbf{h}^{t-1}} < 1 \rightarrow \frac{\partial J^t}{\partial \mathbf{h}^1} \approx 0$$



# Long short-term memory networks (LSTMs)

- LSTMs explicitly keep and update cell memory  $c^{(t)}$  by
  - Removing the previous cell content  $c^{(t-1)}$  by multiplying it with  $f^{(t)}$
  - Adding the new cell content  $\tilde{c}^{(t)}$  multiplied by  $i^{(t)}$
- LSTMs produce output  $h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$

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

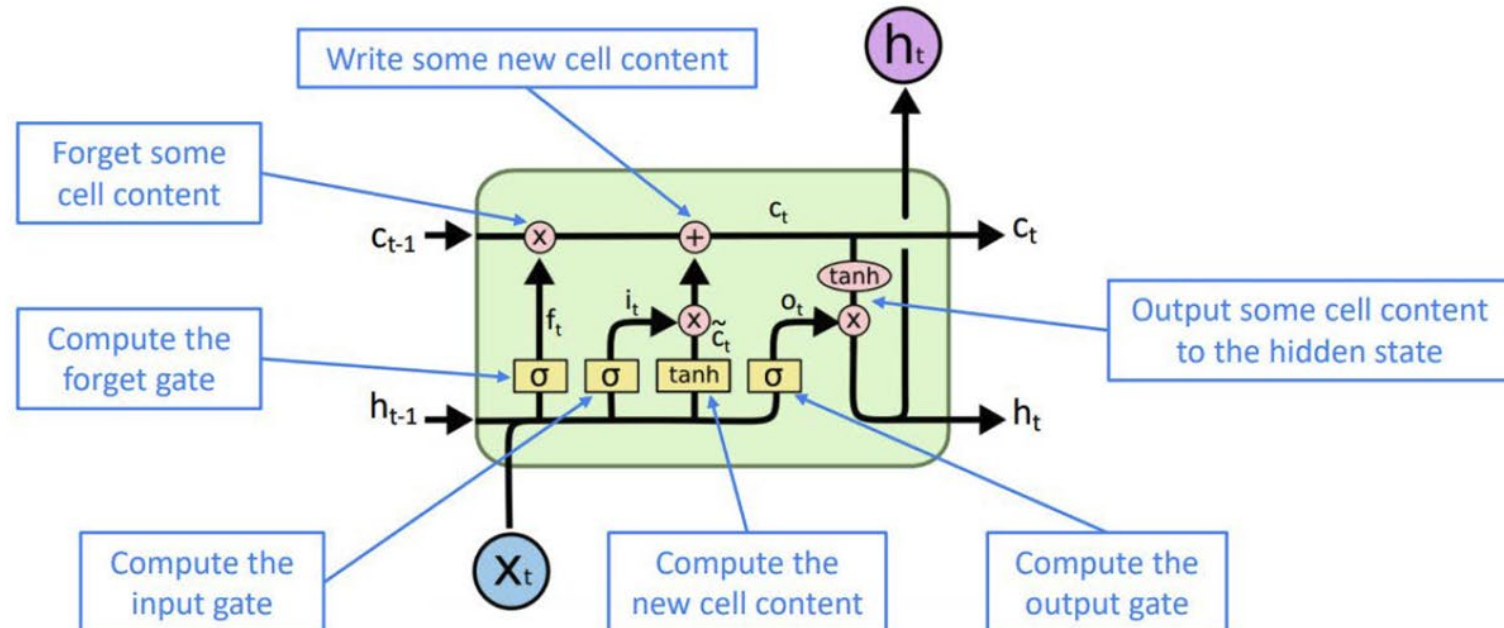
$$i^{(t)} = \sigma \left( W_i h^{(t-1)} + U_i x^{(t)} + b_i \right)$$

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

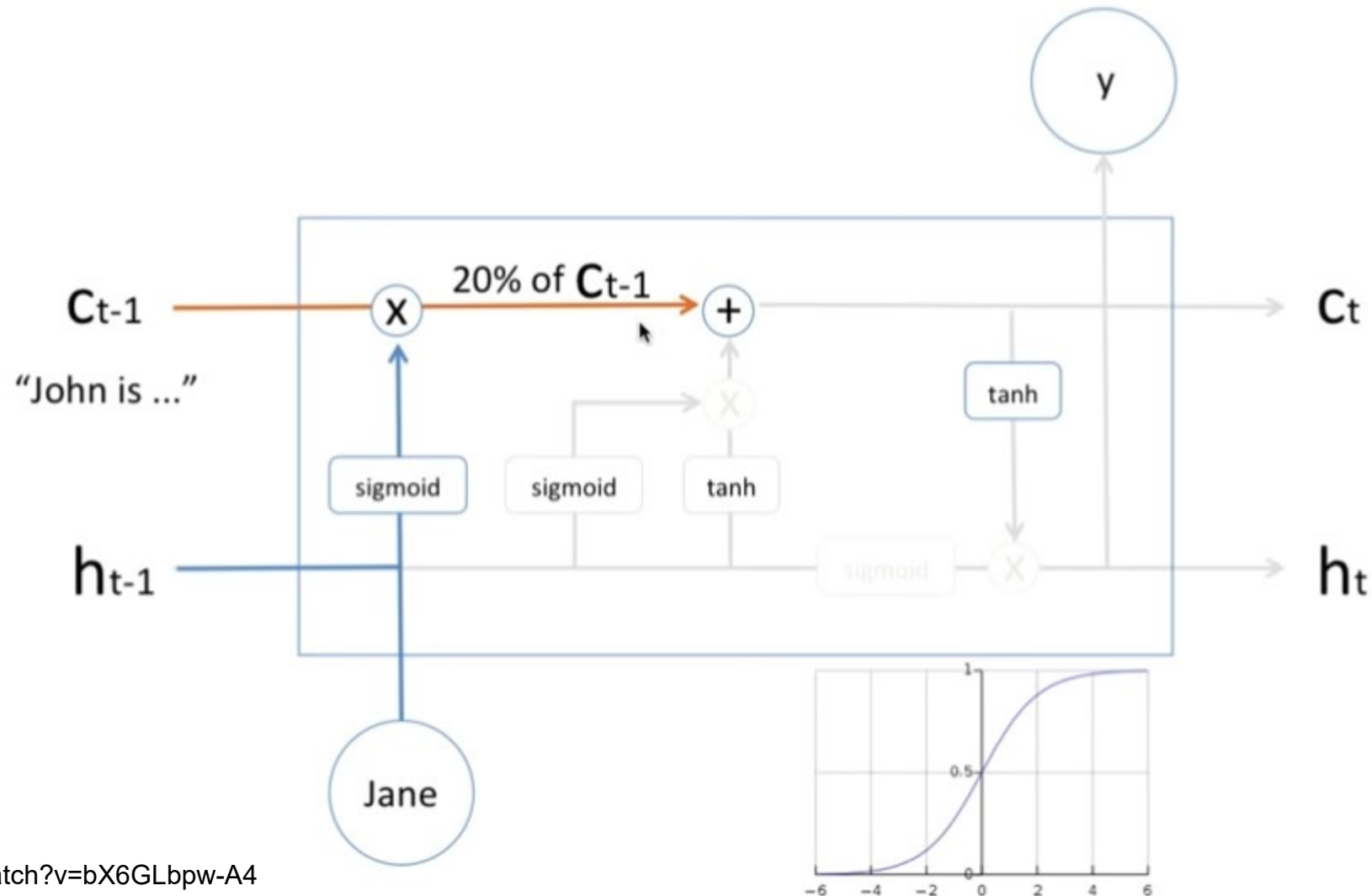
$$\tilde{c}^{(t)} = \tanh \left( W_c h^{(t-1)} + U_c x^{(t)} + b_c \right)$$

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

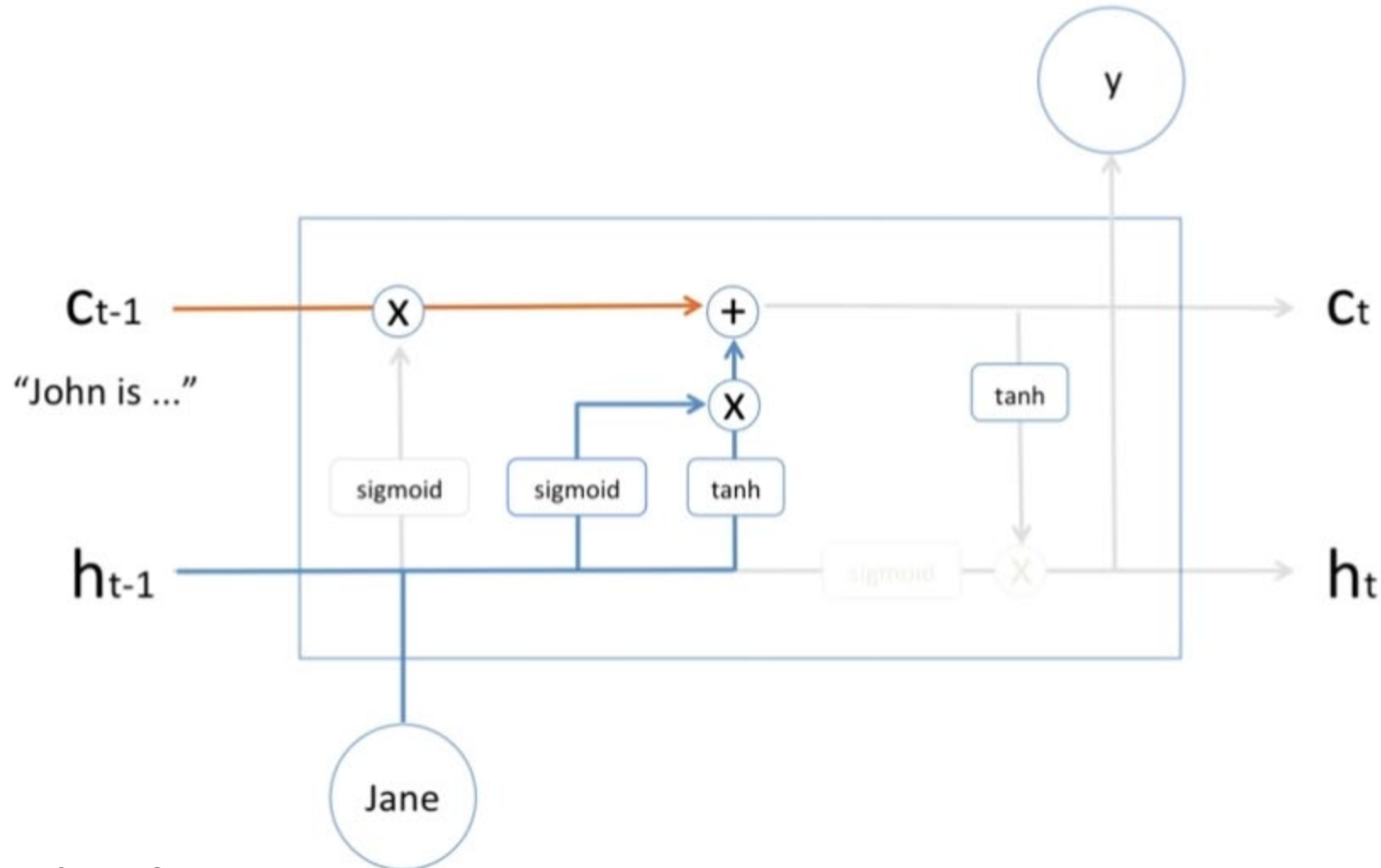
$$h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$



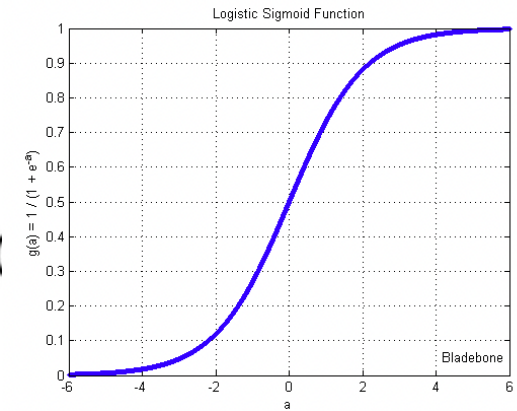
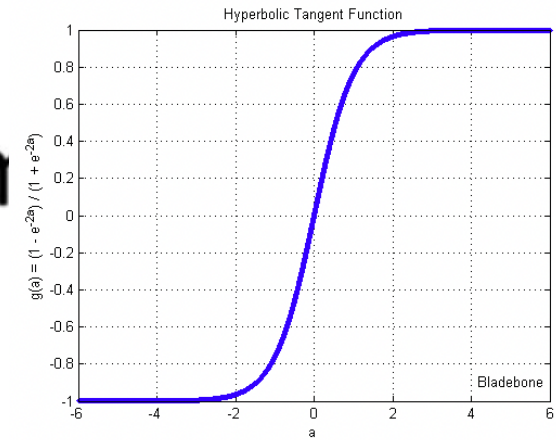
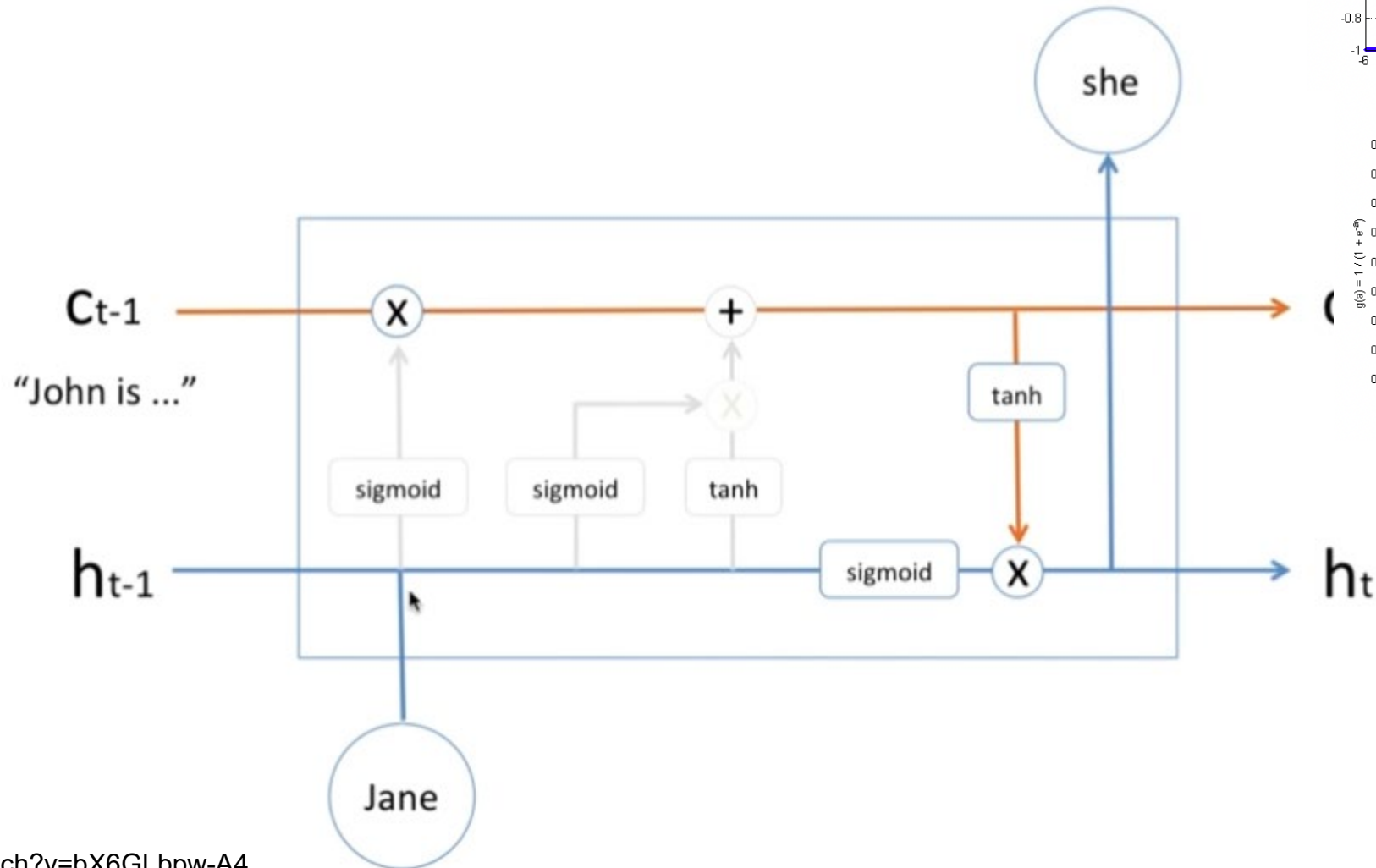
# LSTM cell: forget mechanism



# LSTM cell: input mechanism



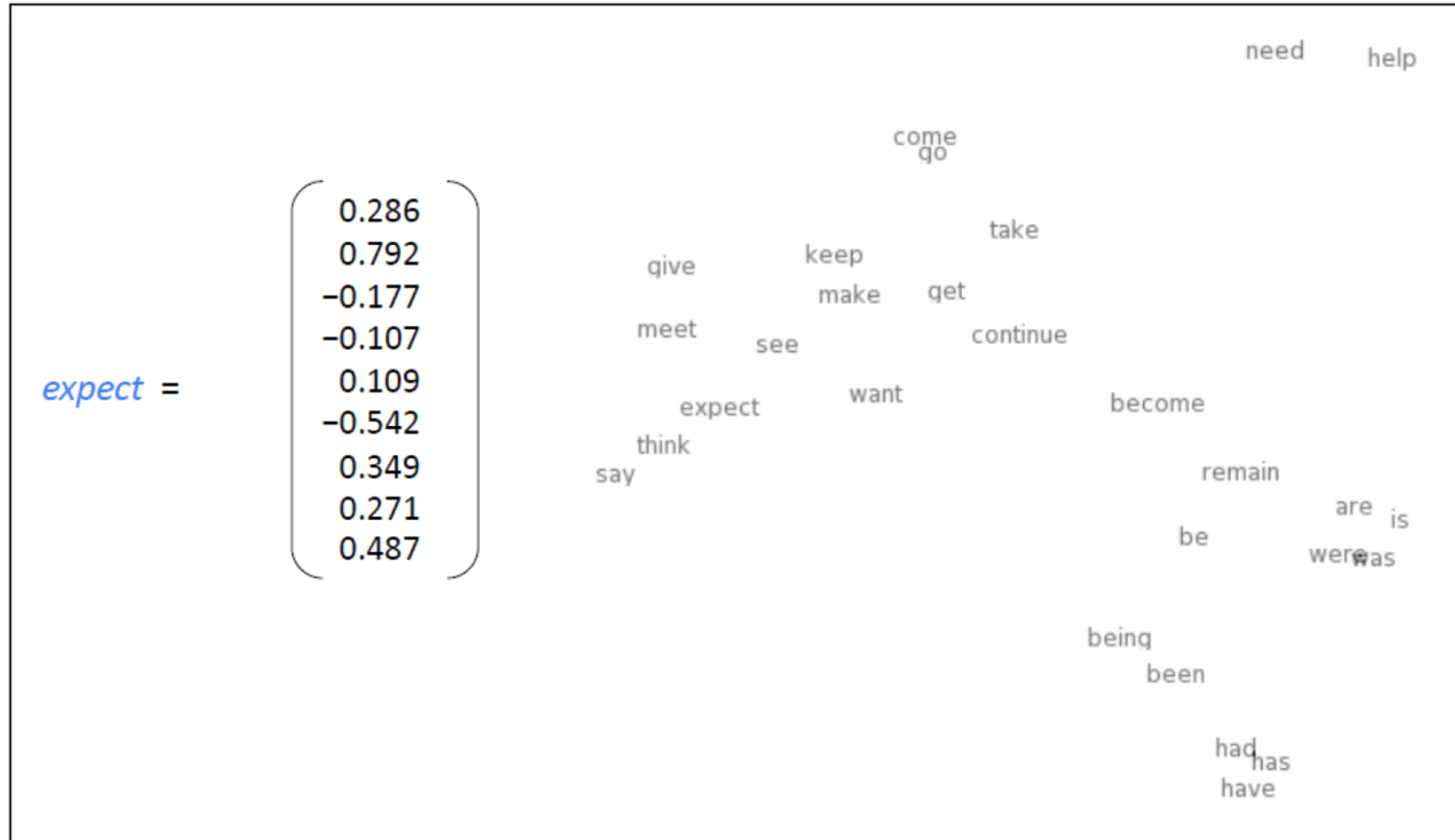
# LSTM cell: output mechanism



# Embeddings

# Word Vector

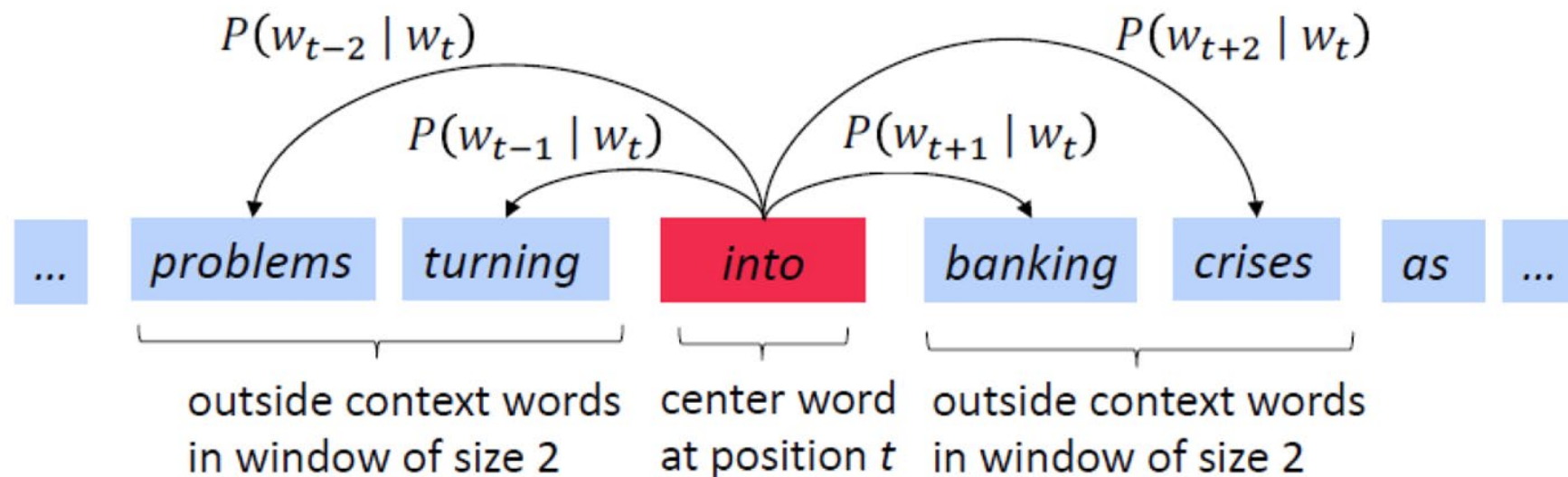
- Represent words as vectors



# Word Vector

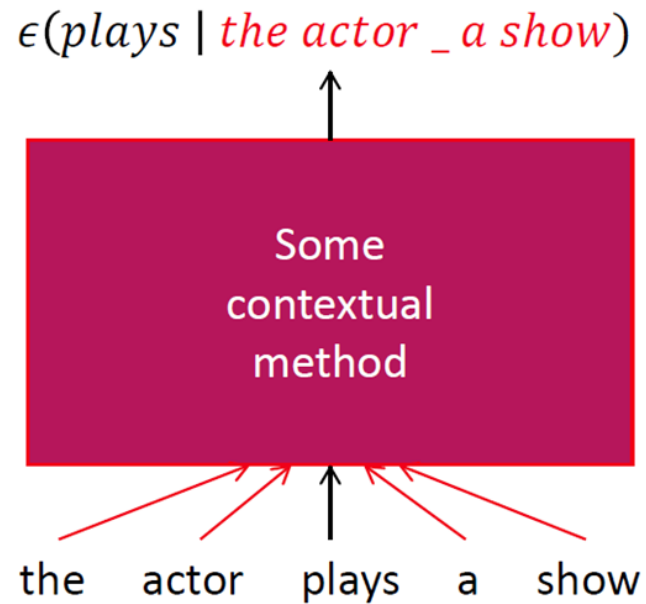
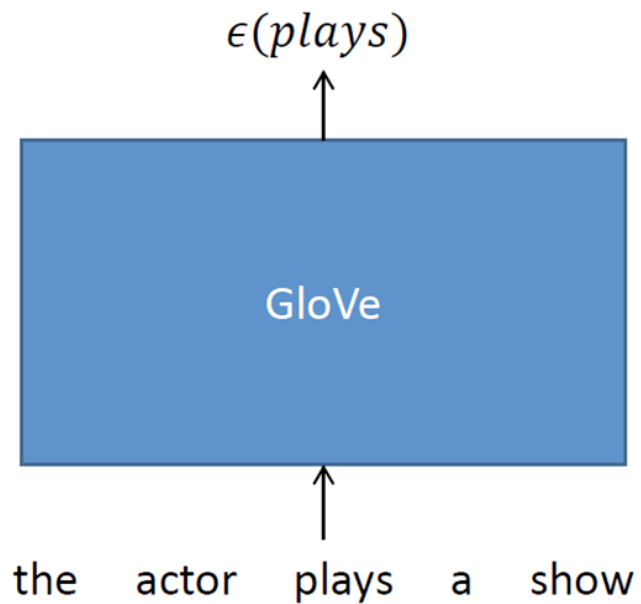
- Distributional semantics : A word's meaning is given by the words that frequently appear close-by
- *"You shall know a word by the company it keeps"*
- Word2vec objective function (skip-grams)

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t; \theta)$$



# Contextual word embedding

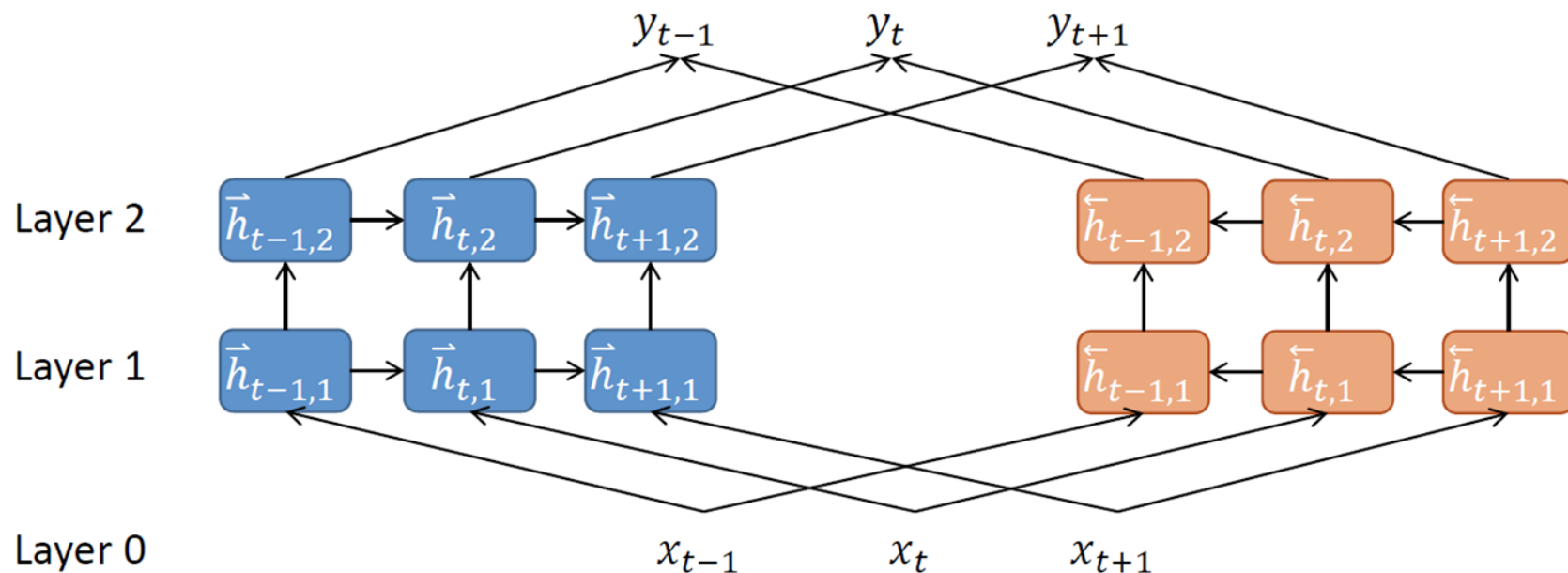
- A word's **contextual embedding** must consider its context





# ELMo: EMBEDDINGS FROM LANGUAGE MODEL

- Multi-layer bidirectional LSTM language model



$$R_k = \{\mathbf{x}_k^{LM}, \vec{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \dots, L\}$$

$$= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \dots, L\},$$

$$\mathbf{h}_{k,0}^{LM} = x_k^{LM} \quad (\text{token representation; GloVe})$$

$$\mathbf{h}_{k,j}^{LM} = \vec{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \quad (\text{LSTM state})$$

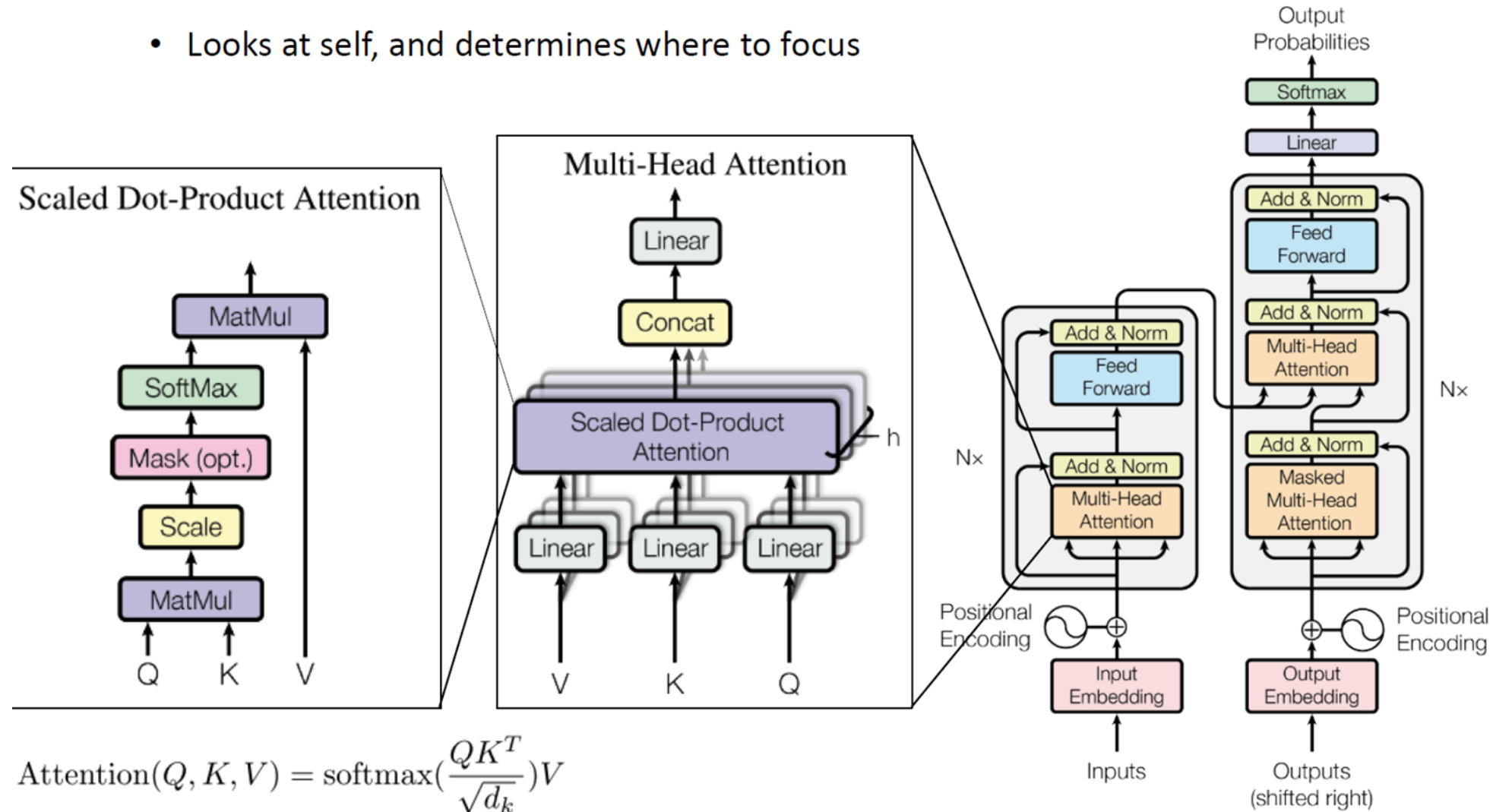
$$\text{ELMo}_k^{\text{task}} = E(R_k; \Theta^{\text{task}}) = \gamma^{\text{task}} \sum_{j=0}^L s_j^{\text{task}} \mathbf{h}_{k,j}^{LM}.$$

$\gamma^{\text{task}}$ : scale (hyper-parameter)

$s_j^{\text{task}}$ : weight (learned)

# Transformer

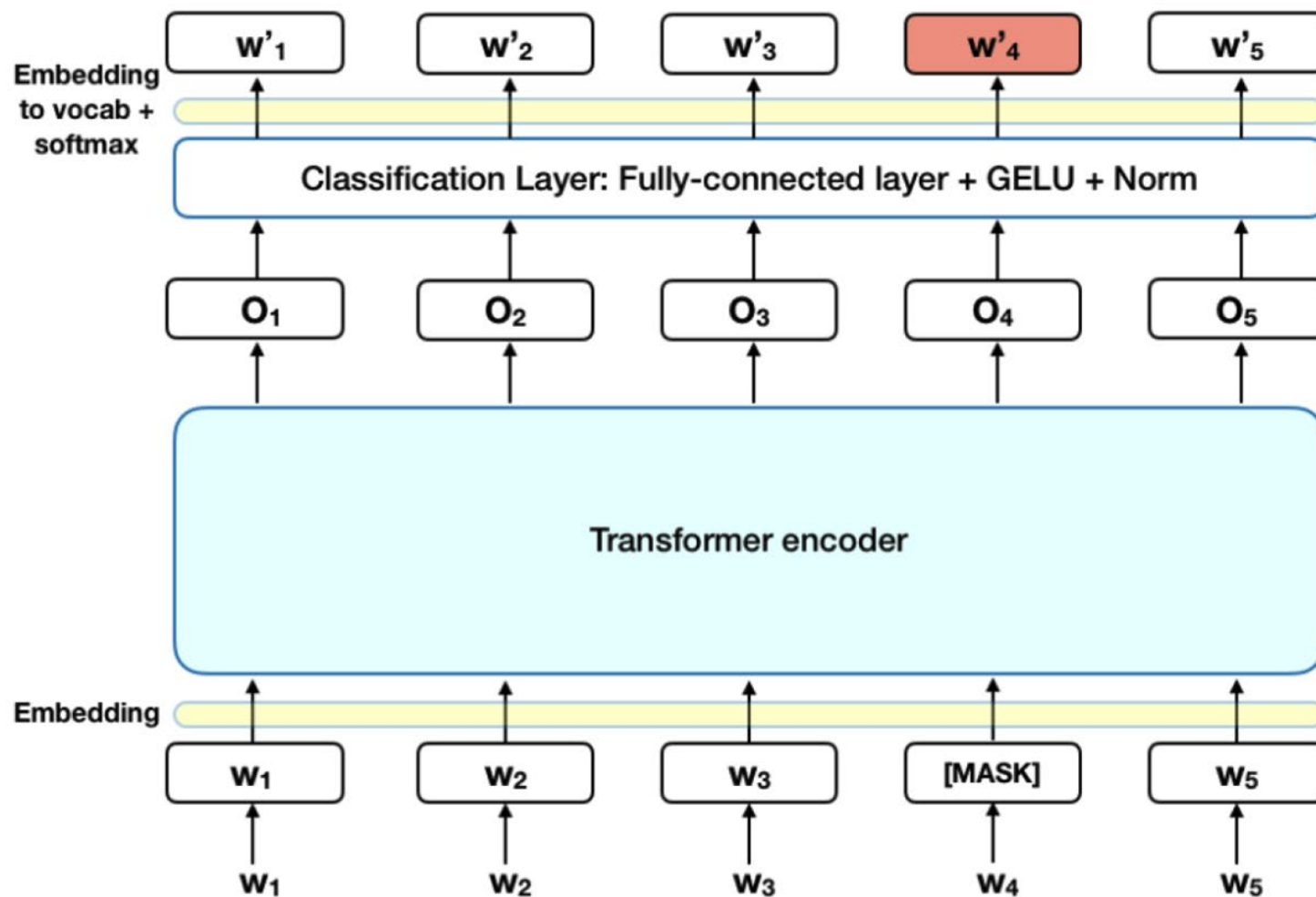
- Parallel self-attention
  - Looks at self, and determines where to focus



# BERT: Bidirectional Encoder Representations from Transformers

- Training 1. Masked words prediction
  - 15% of words are [MASK]ed

\*GELU: Gaussian error linear unit



# BERT: Bidirectional Encoder Representations from Transformers

- Training 2. Next sentence prediction
  - To understand texts more than a sentence

**Input** = [CLS] the man went to [MASK] store [SEP]

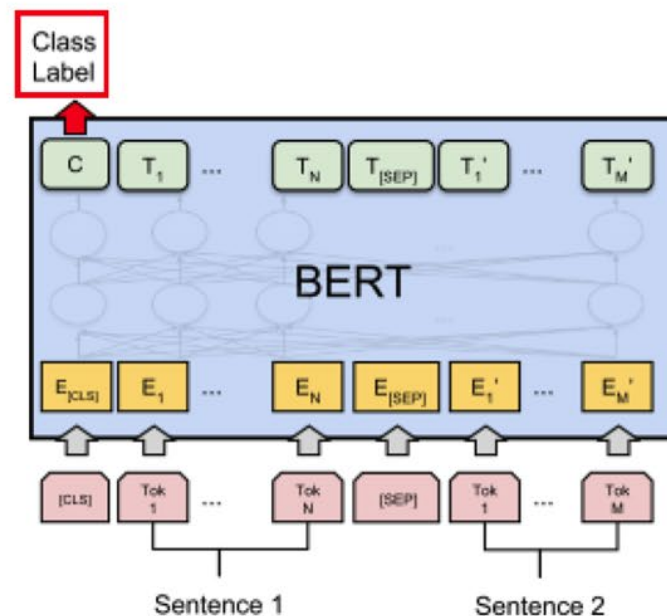
he bought a gallon [MASK] milk [SEP]

**Label** = IsNext

**Input** = [CLS] the man [MASK] to the store [SEP]

penguin [MASK] are flight ##less birds [SEP]

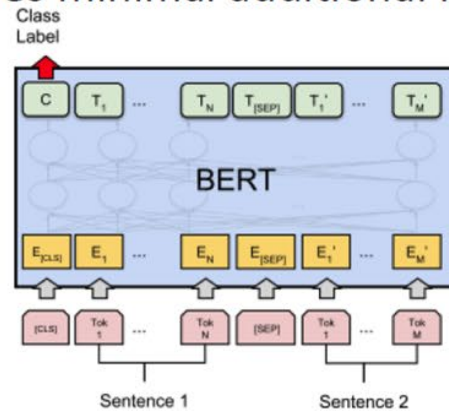
**Label** = NotNext



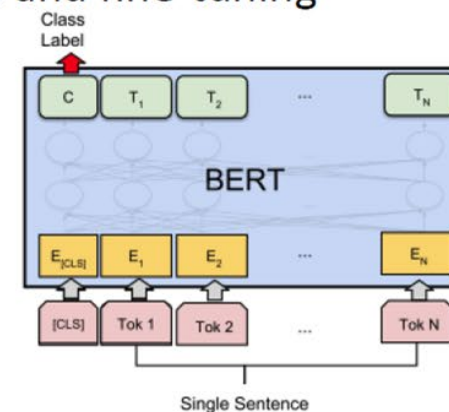
# BERT: Bidirectional Encoder Representations from Transformers

- BERT as universal pre-trained model for NLP
  - BERT requires minimal additional layers and fine-tuning

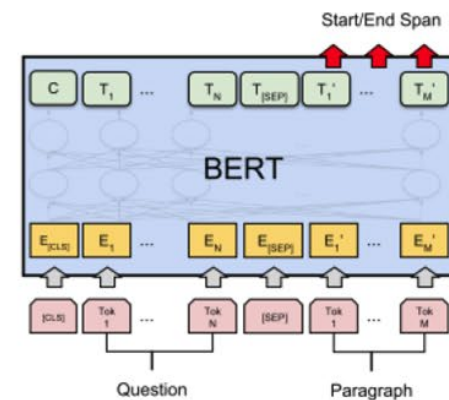
\*GLUE benchmark task



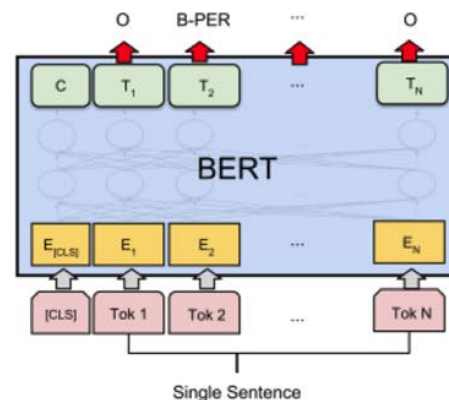
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



(b) Single Sentence Classification Tasks:  
SST-2, CoLA



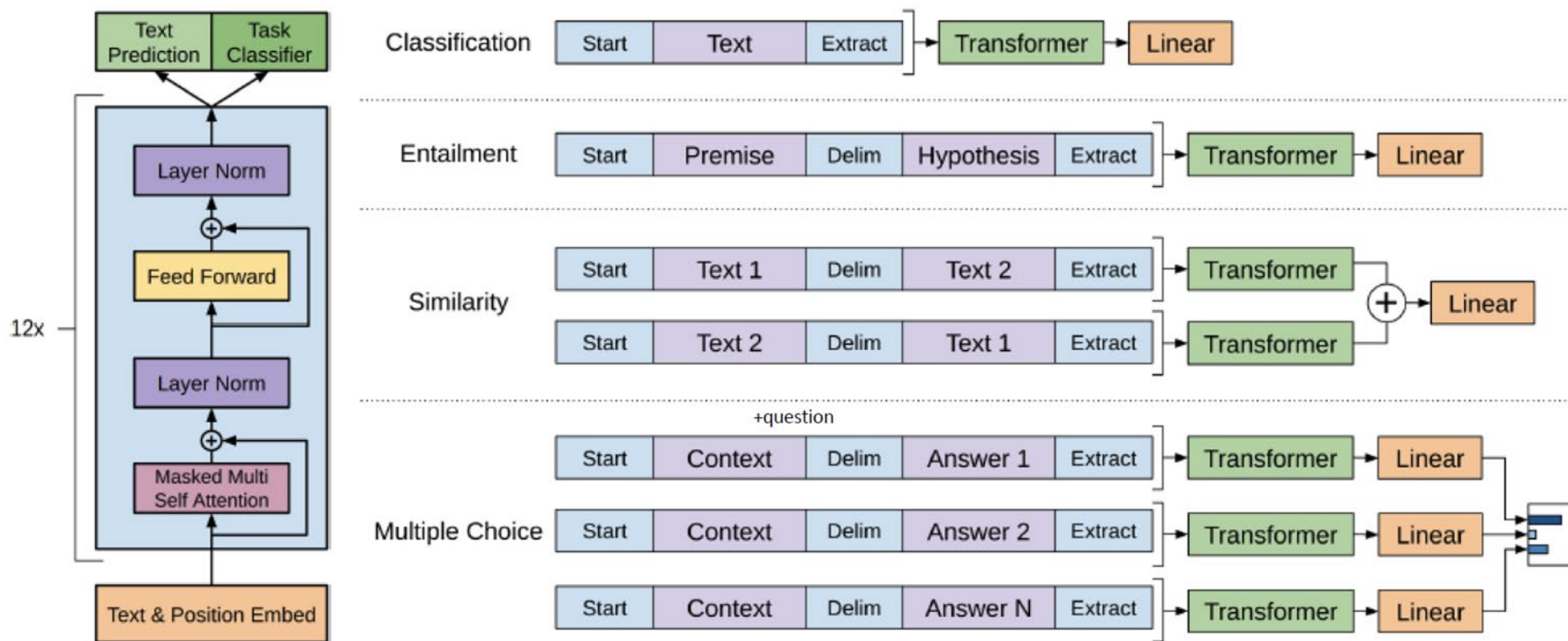
(c) Question Answering Tasks:  
SQuAD v1.1



(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER



# GPT(Generative Pre-Training)



$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta)$$

\*auto-regressive

$$L_2(\mathcal{C}) = \sum_{(x,y)} \log P(y | x^1, \dots, x^m)$$

$$L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C})$$

In pre-training, optimize  $L_1(u)$

$u$ : Unlabeled dataset

$\Theta$ : Model parameters

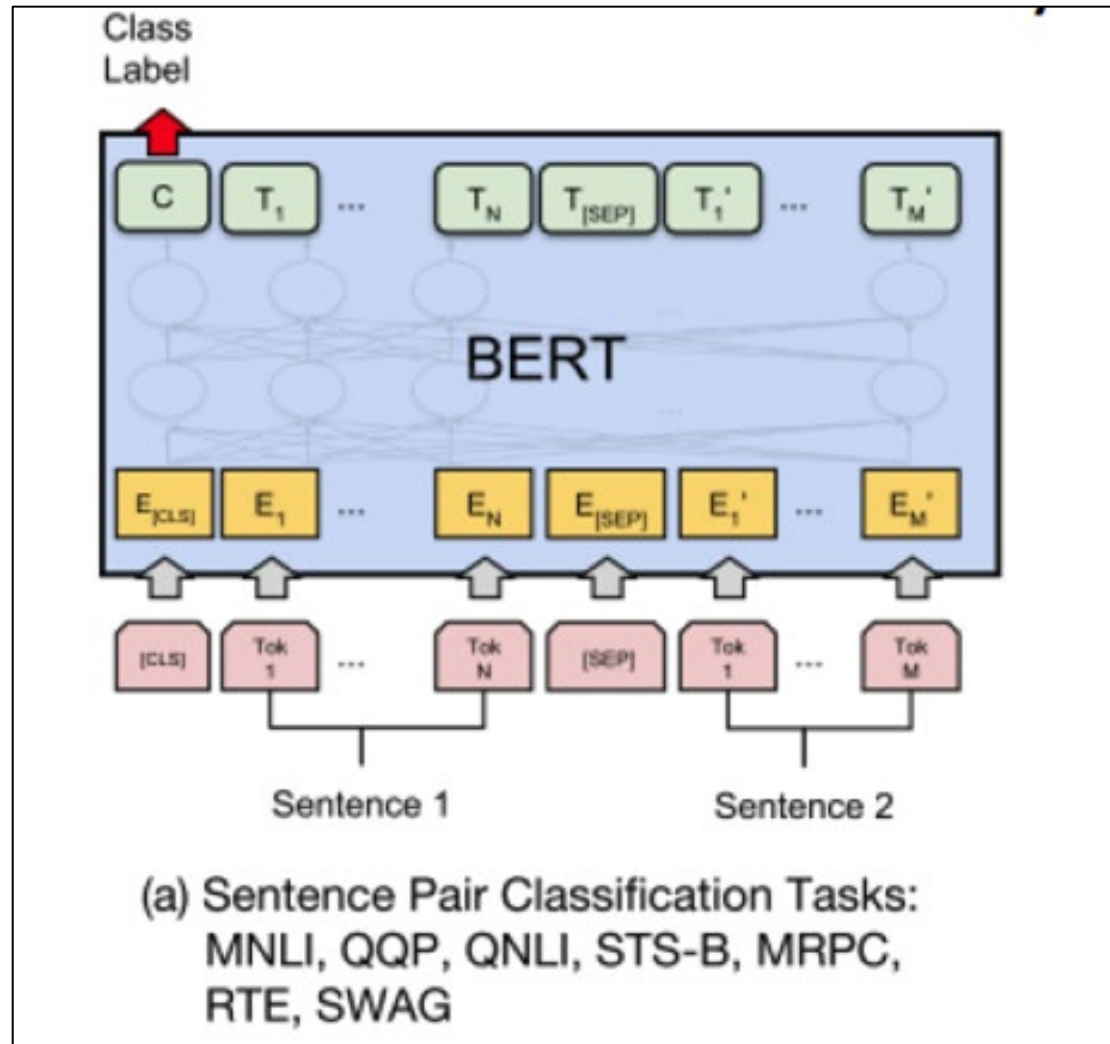
In fine-tuning, optimize  $L_3(c)$

$c$ : Labeled dataset

$\lambda$ : Hyper-parameter weight

# Task 별 아키텍처

# Sentence pair classification





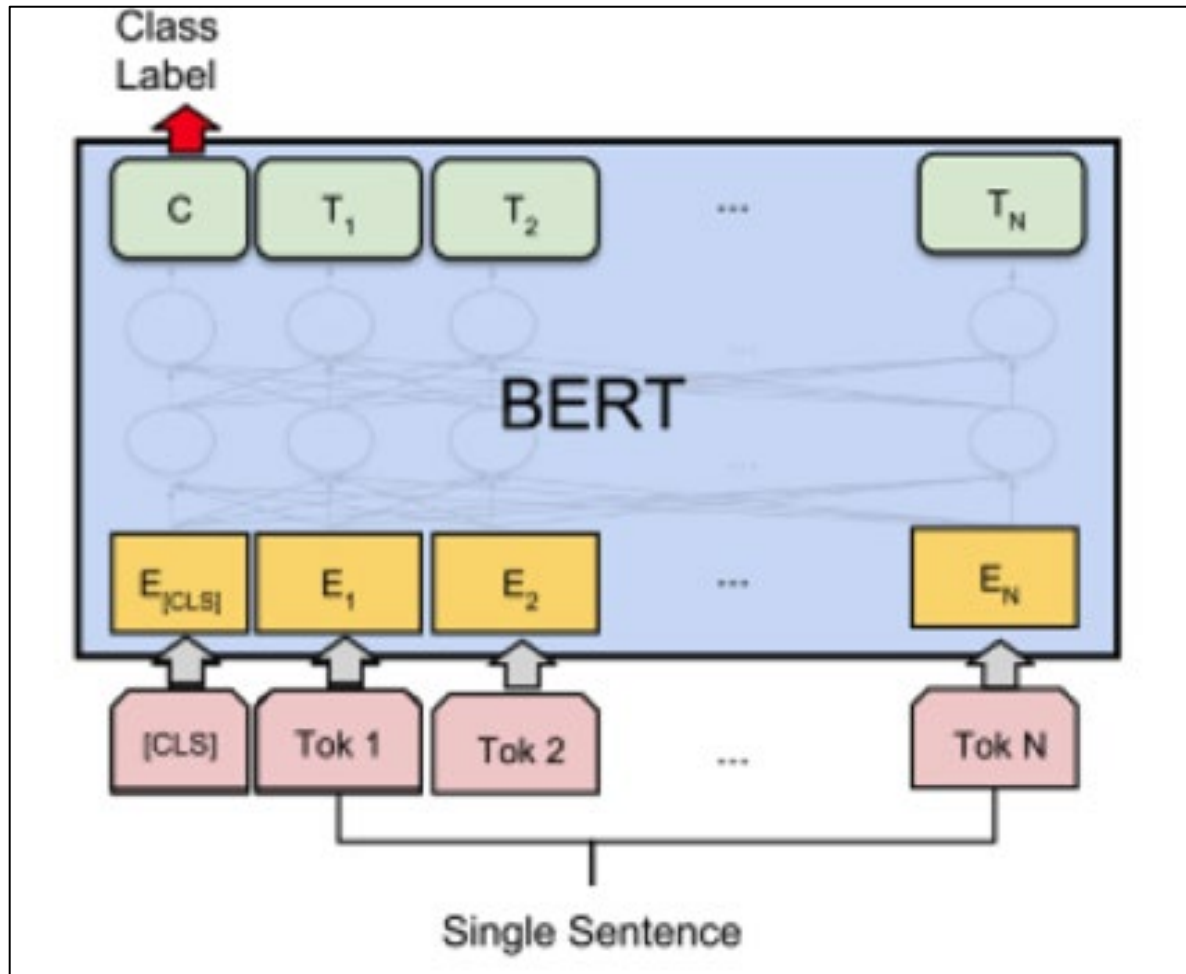
# Sentence pair classification

## SNLI(Stanford Natural Language Interence) Corpus

<https://nlp.stanford.edu/projects/snli/>

Text	Judgments	Hypothesis
A man inspects the uniform of a figure in some East Asian country.	contradiction C C C C C	The man is sleeping.
An older and younger man smiling.	neutral N N E N N	Two men are smiling and laughing at the cats playing on the floor.
A black race car starts up in front of a crowd of people.	contradiction C C C C C	A man is driving down a lonely road.
A soccer game with multiple males playing.	entailment E E E E E	Some men are playing a sport.
A smiling costumed woman is holding an umbrella.	neutral N N E C N	A happy woman in a fairy costume holds an umbrella.

# Single sentence classification task



<Data Sample>

lend some dignity to a dumb story 0

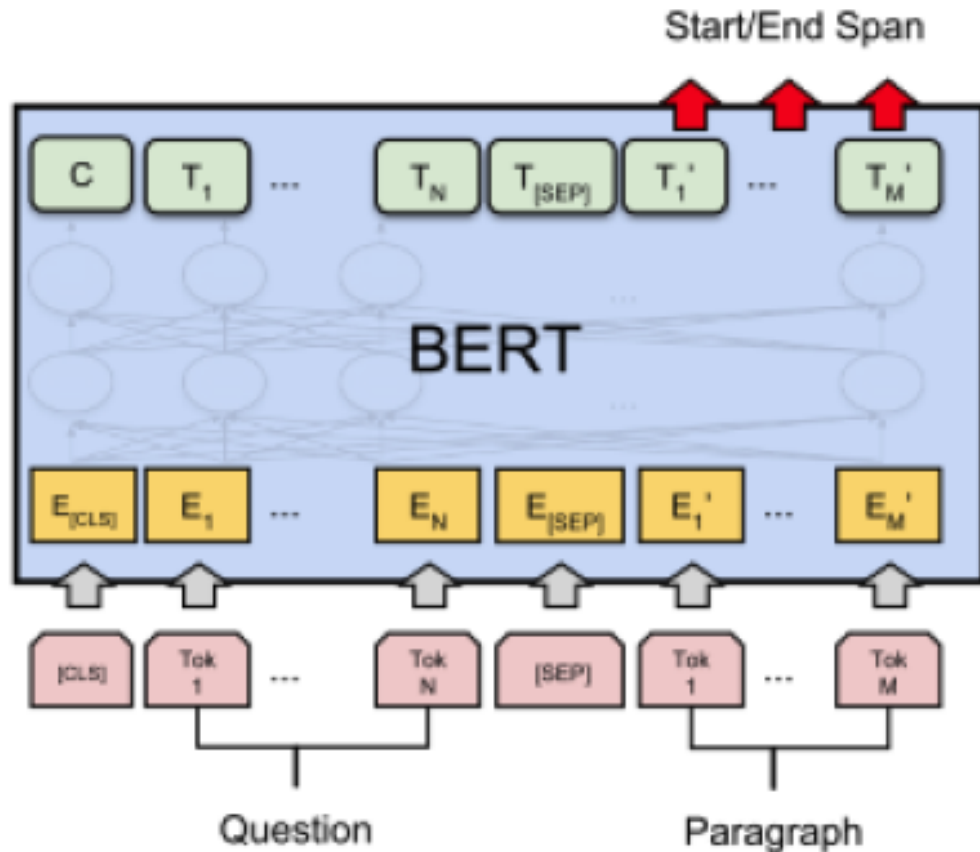
the greatest musicians 1

cold movie 0

with his usual intelligence and subtlety 1

redundant concept 0

# Question Answering



(c) Question Answering Tasks:  
SQuAD v1.1

# Question Answering

Start/End Span

```
In[3]:= sample = ResourceData["SQuAD v1.1", "TrainingData"][[All, 1]]
```

<| Context →

Architecturally, the school has a Catholic character. Atop the Main Building's gold dome is a golden statue of the Virgin Mary. Immediately in front of the Main Building and facing it, is a copper statue of Christ with arms upraised with the legend "Venite Ad Me Omnes". Next to the Main Building is the Basilica of the Sacred Heart. Immediately behind the basilica is the Grotto, a Marian place of prayer and reflection. It is a replica of the grotto at Lourdes, France where the Virgin Mary reputedly appeared to Saint Bernadette Soubirous in 1858. At the end of the main drive (and in a direct line that connects through 3 statues and the Gold Dome), is a simple, modern stone statue of Mary.,

Out[3]=

Question → To whom did the Virgin Mary allegedly appear in 1858 in Lourdes France?,

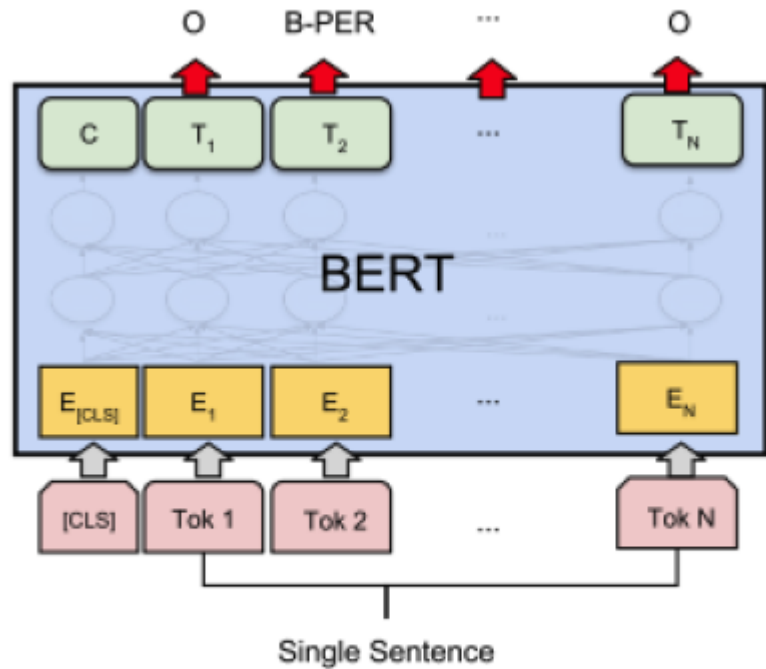
Answer →

Saint Bernadette Soubirous,

AnswerPosition → 516|>

SQuAD v1.1

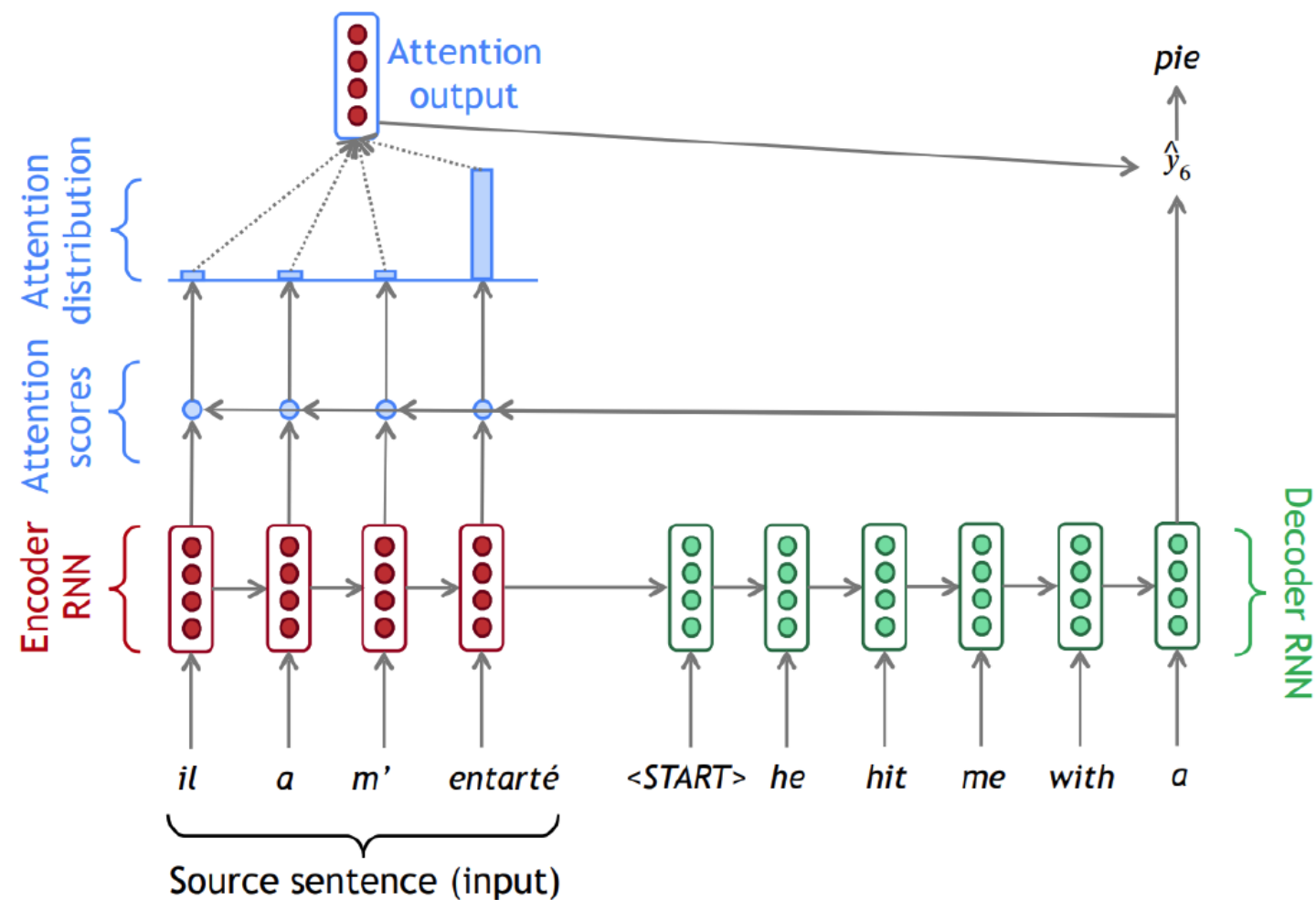
# Single sentence tagging task



(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

U.N.	NNP	I-NP	I-ORG
official	NN	I-NP	O
Ekeus	NNP	I-NP	I-PER
heads	VBZ	I-VP	O
for	IN	I-PP	O
Baghdad	NNP	I-NP	I-LOC
.	.	O	O

# NMT with attention



# HuggingFace Tutorial

# Tokenizer

Models can only process numbers.

so, tokenizers need to convert our text inputs to numerical data.

Split on spaces

Let's	do	tokenization!
-------	----	---------------

Split on punctuation

Let	's	do	tokenization	!
-----	----	----	--------------	---

Word-based

L	e	t	'	s	d	o	t	o	k	e	n	i	z	a	t	i	o	n	!
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Character-based



# Subword tokenization

Let's </w>	do</w>	token	ization</w>	!</w>
------------	--------	-------	-------------	-------

Subword tokenization algorithms rely on the principle that frequently used words should not be split into smaller subwords, but rare words should be decomposed into meaningful subwords.

Byte-level BPE, as used in GPT-2

WordPiece, as used in BERT

SentencePiece or Unigram, as used in several multilingual models

```
from transformers import BertTokenizer
```

```
tokenizer = BertTokenizer.from_pretrained("bert-base-cased")
```

```
tokenizer("Using a Transformer network is simple")
```

```
{'input_ids': [101, 7993, 170, 11303, 1200, 2443, 1110, 3014, 102],  
  'token_type_ids': [0, 0, 0, 0, 0, 0, 0, 0, 0],  
  'attention_mask': [1, 1, 1, 1, 1, 1, 1, 1, 1]}
```

# Encoding

Translating text to numbers is known as encoding.

Encoding is done in a two-step process: the tokenization, followed by the conversion to input IDs

## Step 1 : Tokenization

```
from transformers import AutoTokenizer

tokenizer = AutoTokenizer.from_pretrained("bert-base-cased")

sequence = "Using a Transformer network is simple"
tokens = tokenizer.tokenize(sequence)

print(tokens)
```

The output of this method is a list of strings, or tokens:

```
['Using', 'a', 'transform', '##er', 'network', 'is', 'simple']
```

## Step 2 :From tokens to input IDs

```
ids = tokenizer.convert_tokens_to_ids(tokens)

print(ids)
```

```
[7993, 170, 11303, 1200, 2443, 1110, 3014]
```

# Decoding

```
decoded_string = tokenizer.decode([7993, 170, 11303, 1200, 2443, 1110, 3014])  
print(decoded_string)
```

```
'Using a Transformer network is simple'
```

**Decoding** is going the other way around: from vocabulary indices, we want to get a string.

Let's code