# Machine Translation (Part I)
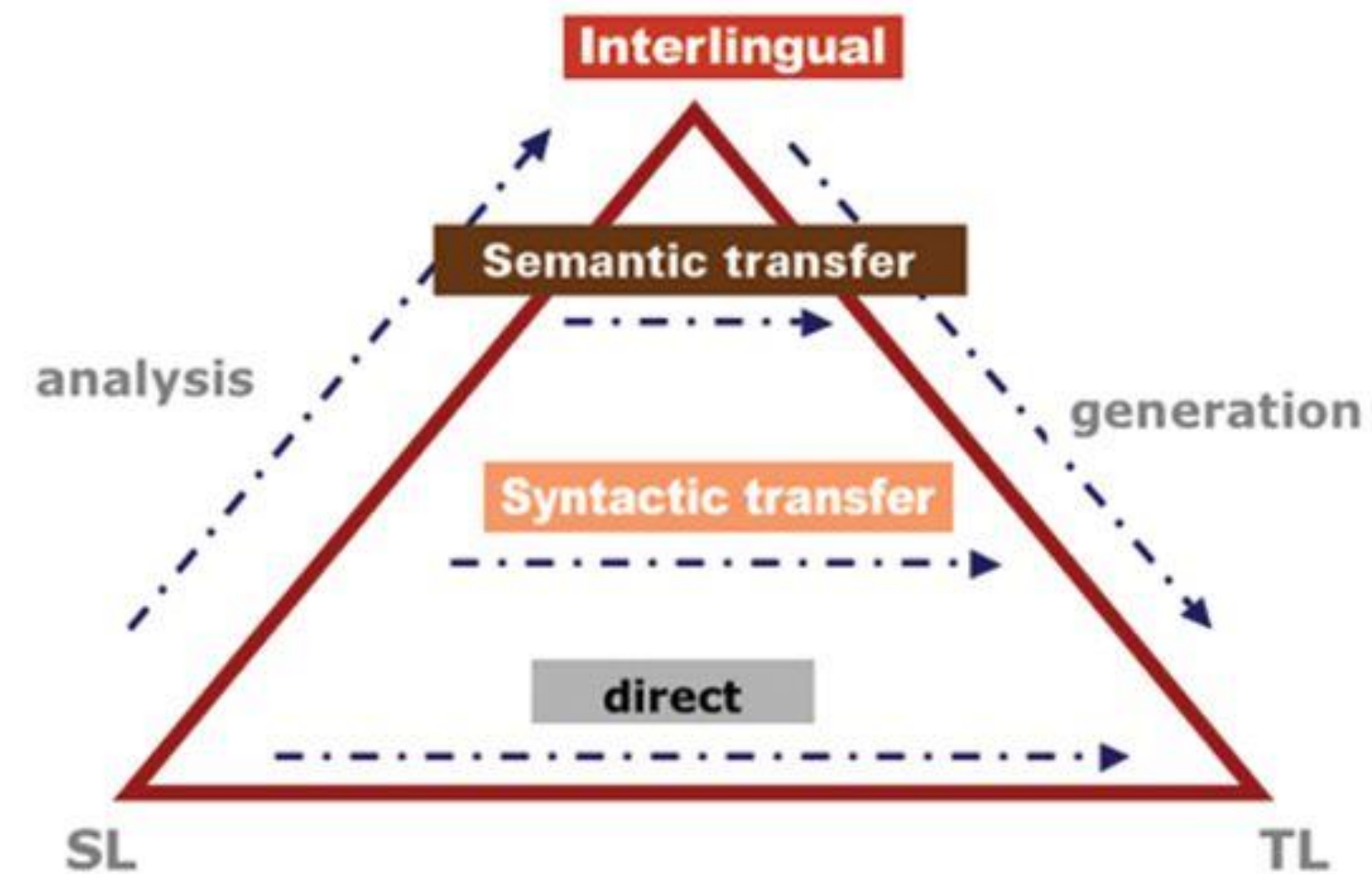
Peerachet Porkaew, NECTEC

# Outline

- Rule-based Machine Translation

- Statistical Machine Translation (SMT)

- Neural Machine Translation (NMT)

  - Sequence-to-Sequence

  - RNN-based NMT + Attention

  - Transformer

# Rule-based Machine Translation

# Rule-based Machine Translation (RBMT)

- Linguistic Knowledge

- Three sub-module :

  - Analysis

  - Transfer

  - Generation

- Template-based Translation


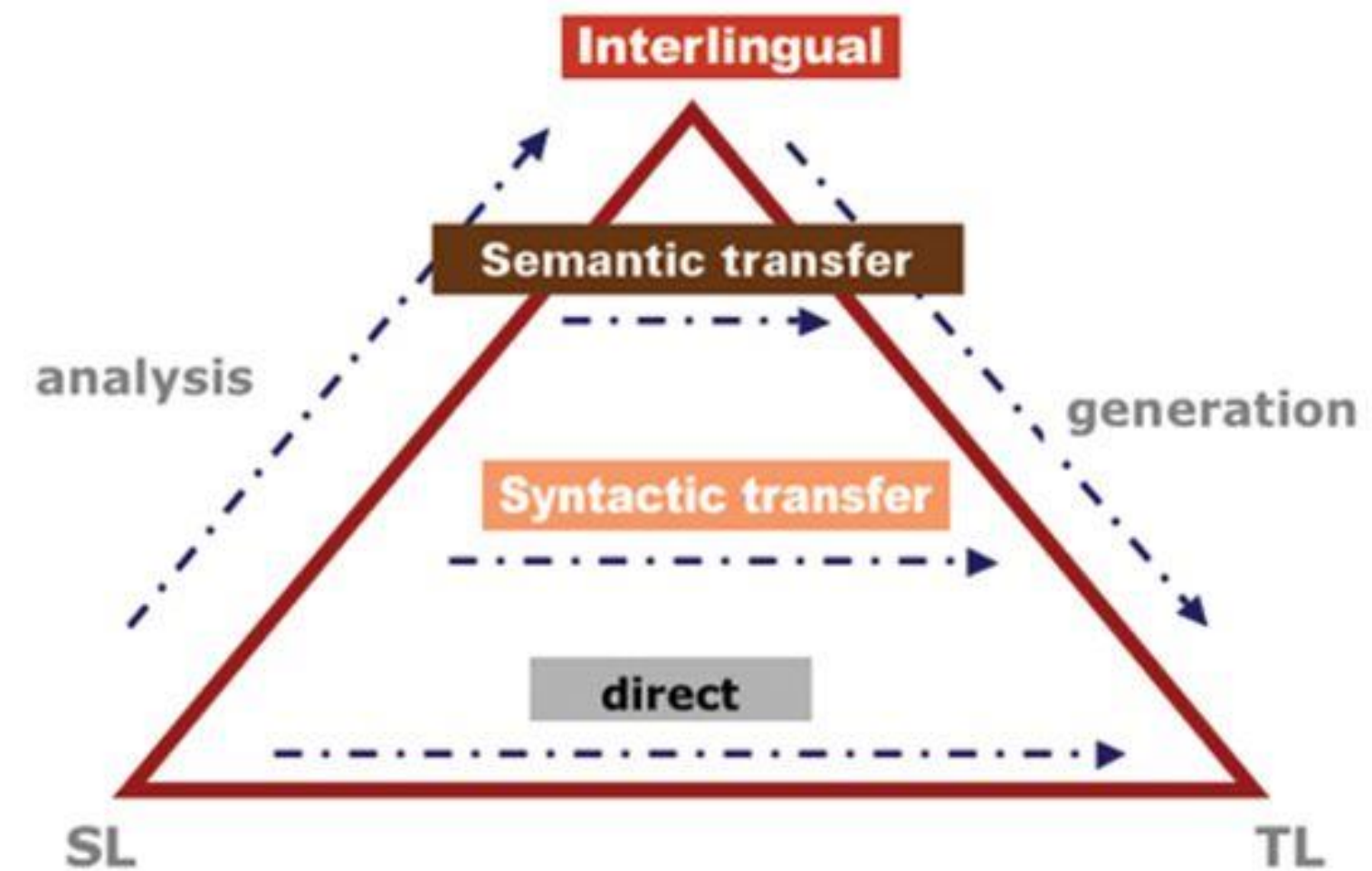
Rule-based MT

The Vauquois Triangle

# Rule-based Machine Translation (RBMT)

- Limitations :

  - Not Automatic

  - Time consuming

  - Conflicts
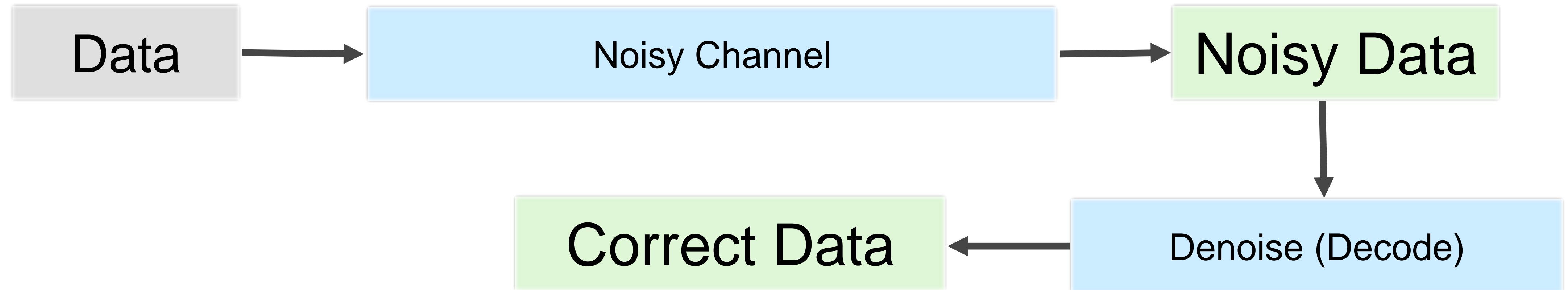
  - Less Clarity

## Rule-based MT



The Vauquois Triangle

# Simple RBMT Demo

# Statistical Machine Translation
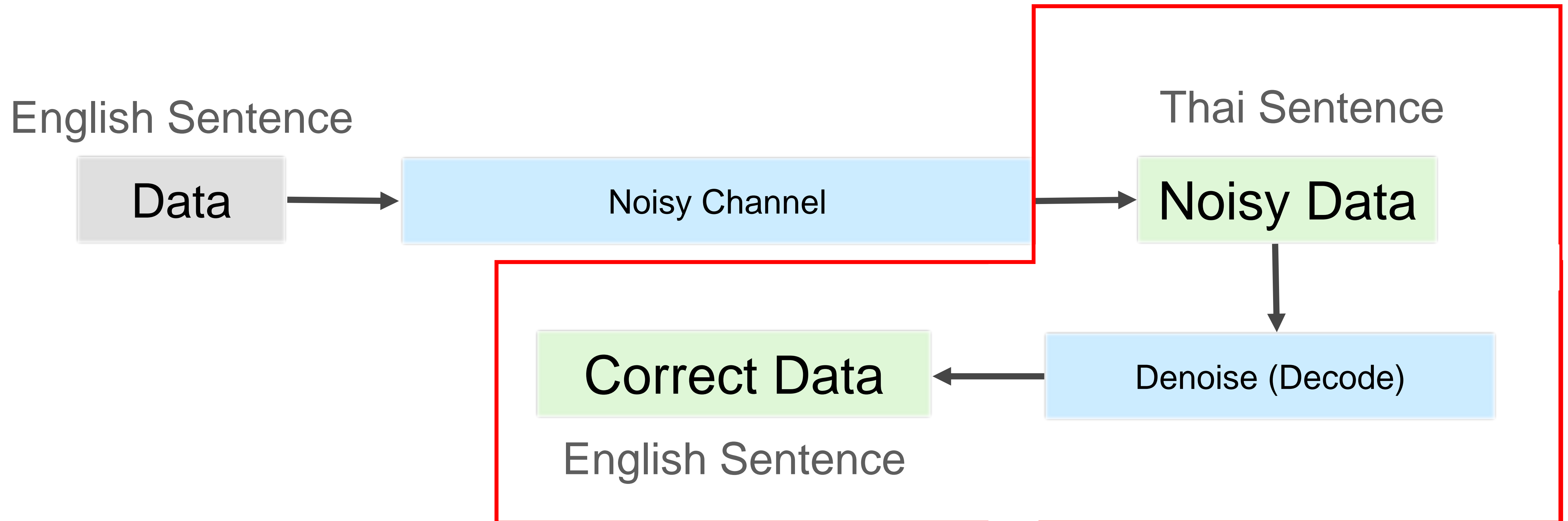
# Statistical Machine Translation

- Noisy Channel Model

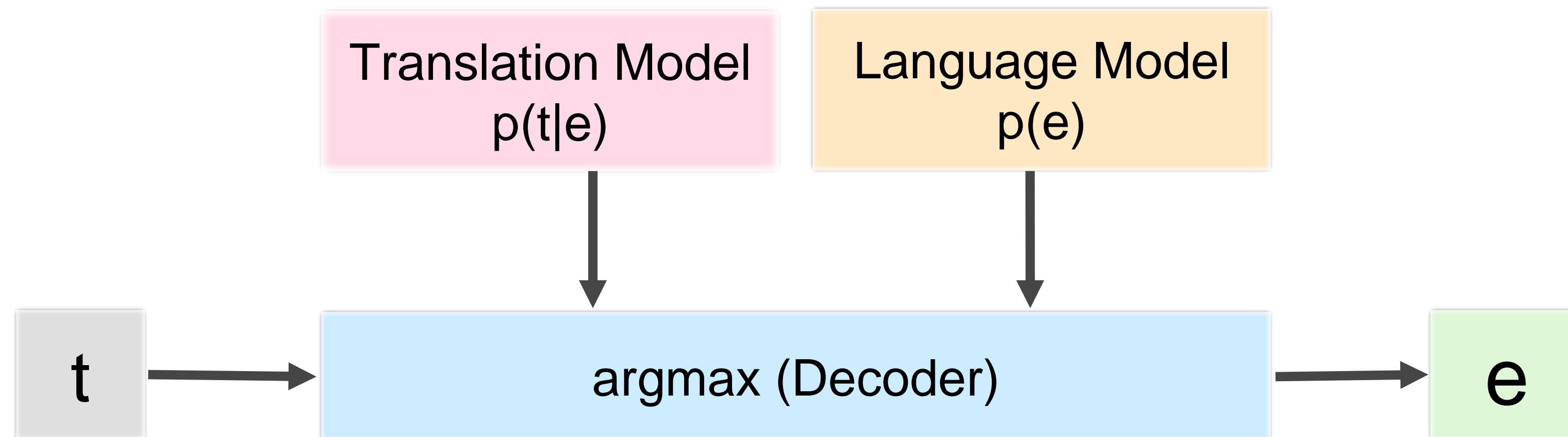# Statistical Machine Translation

- Noisy Channel Model

English Sentence

| Data | → | Noisy Channel | → | Noisy Data |

Thai Sentence

Noisy Data ↓

Correct Data ← Denoise (Decode)

English Sentence

# Statistical Machine Translation

- Source sentence *t*, e.g. Thai

- Target sentence *e*, e.g. English

- Probabilistic formulation using Bayes rule $\quad P(A|B) = \dfrac{P(B|A) \cdot P(A)}{P(B)}$

$$\hat{e} = \operatorname*{argmax}_e \mathrm{p}(e|t)$$
$$\hat{e} = \operatorname*{argmax}_e \mathrm{p}(t|e)p(e)$$

# Statistical Machine Translation (Cont.)

$$\hat{e} = \operatorname{argmax}_e \mathrm{p}(t|e)p(e)$$

Translation Model
p(t|e)

Language Model
p(e)

t

argmax (Decoder)

e

# Statistical Decoder (Simplified version)

| เมื่อวาน | นี้ | ฉัน | ไป | ทะเล | กับ | เพื่อน |
|---|---|---|---|---|---|---|
| Yesterday | this | I | go | the sea | with | friends |
| Yesterday | | I | went to | sea | with | friend |
| Previous day | | me | get to | ocean | with my friend | |

| เมื่อวาน | นี้ | ฉัน | ไป | ทะเล | กับ | เพื่อน |
|---|---|---|---|---|---|---|
| Yesterday | this | I | go | the sea | with | friends |
| Yesterday | | I | went to | sea | with | friend |
| Previous day | | me | get to | ocean | with my friend | |

| เมื่อวาน | นี้ | ฉัน | ไป | ทะเล | กับ | เพื่อน |
|---|---|---|---|---|---|---|
| Yesterday | this | I | go | the sea | with | friends |
| Yesterday | | I | went to | sea | with | friend |
| Previous day | | me | get to | ocean | with my friend | |

# Translation Model *p*(t|e)

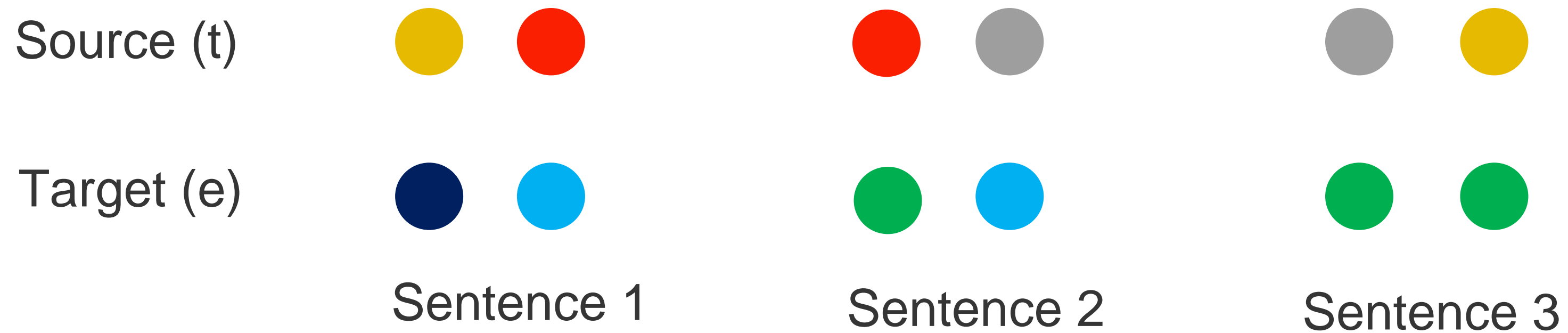- Word-based Translation Model

$$p(t|e) = p\left(\text{ไป}\middle|go\right) = 0.5$$

- Phrase-based Translation Model

$$\mathrm{p}(t_1 t_2, \ldots, t_n | e_1, e_2, \ldots, e_m) = p(\text{เมื่อวาน นี้}| \text{ previous day}) = 0.1$$

# How we get the *p*(t|e) ?

- We do not have alignments. No problem we assume alignment are uniformly paired. ***p(t|e) = c*** (constant)

Source (t)

Target (e)

Sentence 1            Sentence 2            Sentence 3

$p(\,\bullet\,|\,\bullet\,)$            $p(\,\bullet\,|\,\bullet\,)$            $p(\,\bullet\,|\,\bullet\,)$

$p(\,\bullet\,|\,\bullet\,)$            $p(\,\bullet\,|\,\bullet\,)$            $p(\,\bullet\,|\,\bullet\,)$

$p(\,\bullet\,|\,\bullet\,)$            $p(\,\bullet\,|\,\bullet\,)$            $p(\,\bullet\,|\,\bullet\,)$

14

# How we get the $p$(t|e) ?

- But, we do not have alignment. No problem we assume alignment are uniformly paired. **$p(t|e) = c$** (constant)

Source (t)

Target (e)

Sentence 1      Sentence 2      Sentence 3

$p(\; \bullet \;|\; \bullet \;) = 1/2$     $p(\; \bullet \;|\; \bullet \;) = 1/4$     $p(\; \bullet \;|\; \bullet \;) = 2/6$

$p(\; \bullet \;|\; \bullet \;) = 1/2$     $p(\; \bullet \;|\; \bullet \;) = 2/4$     $p(\; \bullet \;|\; \bullet \;) = 1/6$

$p(\; \bullet \;|\; \bullet \;) = 0/2$     $p(\; \bullet \;|\; \bullet \;) = 1/4$     $p(\; \bullet \;|\; \bullet \;) = 3/6$

# How we get the *p*(t|e) ?

- Now, we can get the better alignment from previous knowledge.



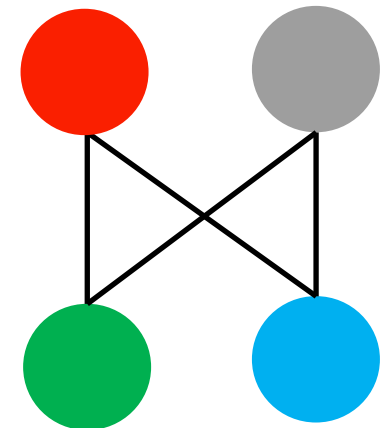Then, we can calculate *p(t|e)* again using these alignment information.
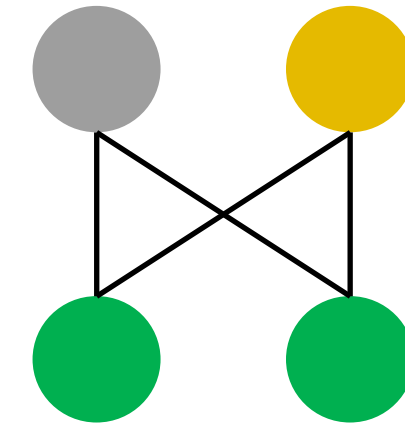
# How we get the *p*(t|e) ?

Source (t)

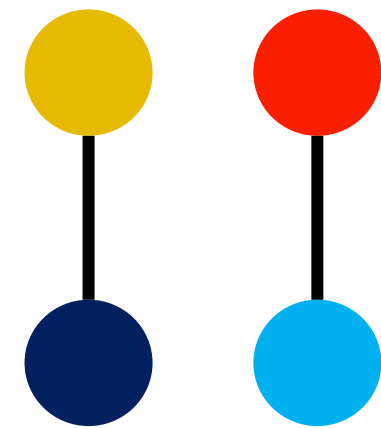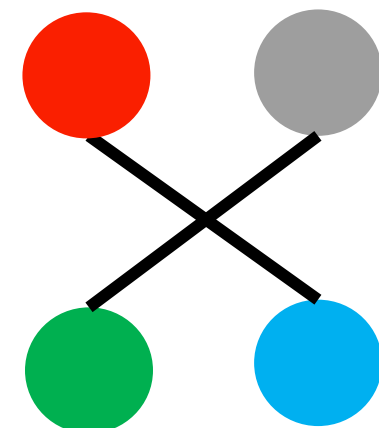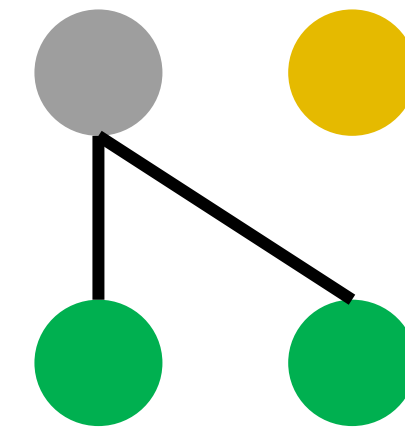Target (e)

Sentence 1

Sentence 2

Sentence 3

Source (t)

Target (e)

Sentence 1

Sentence 2

Sentence 3

Expectation Maximization (EM)

# Translation Model

- **Word-based Translation Model**

$$p(t|e) = p\left(\text{ไป}\middle|go\right) = 0.5$$

- Phrase-based Translation Model

$$p(t_1 t_2, \ldots, t_n | e_1, e_2, \ldots, e_m) = p(\text{เมื่อวาน นี้}| \text{previous day}) = 0.1$$

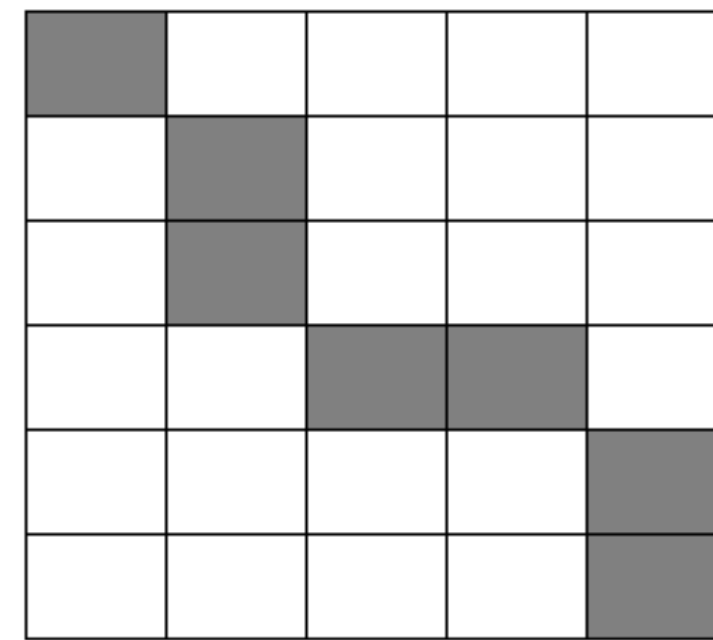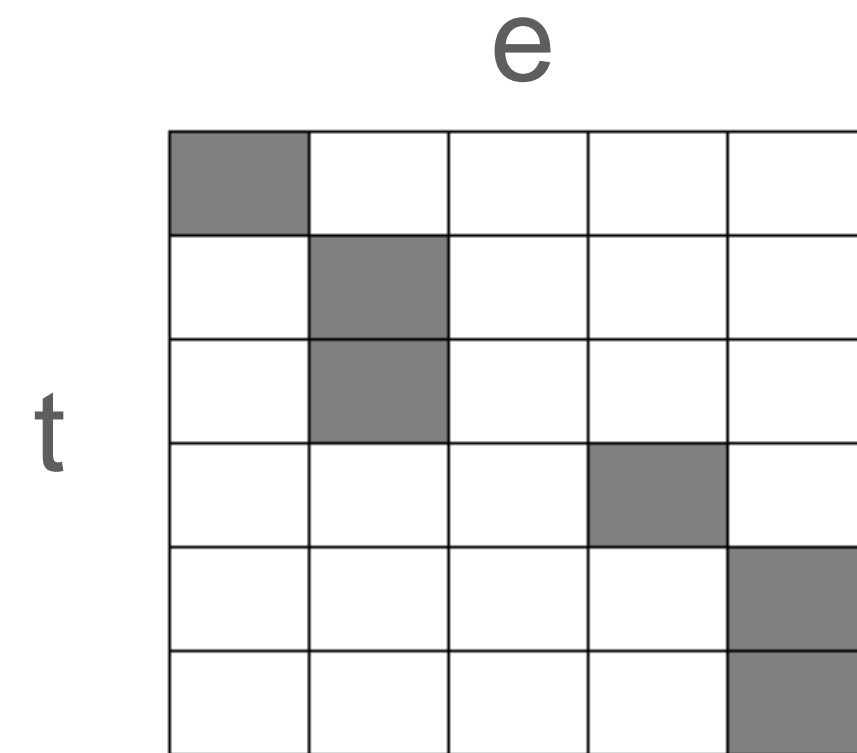# Word-based Translation Model

- **Word-based Translation Model**

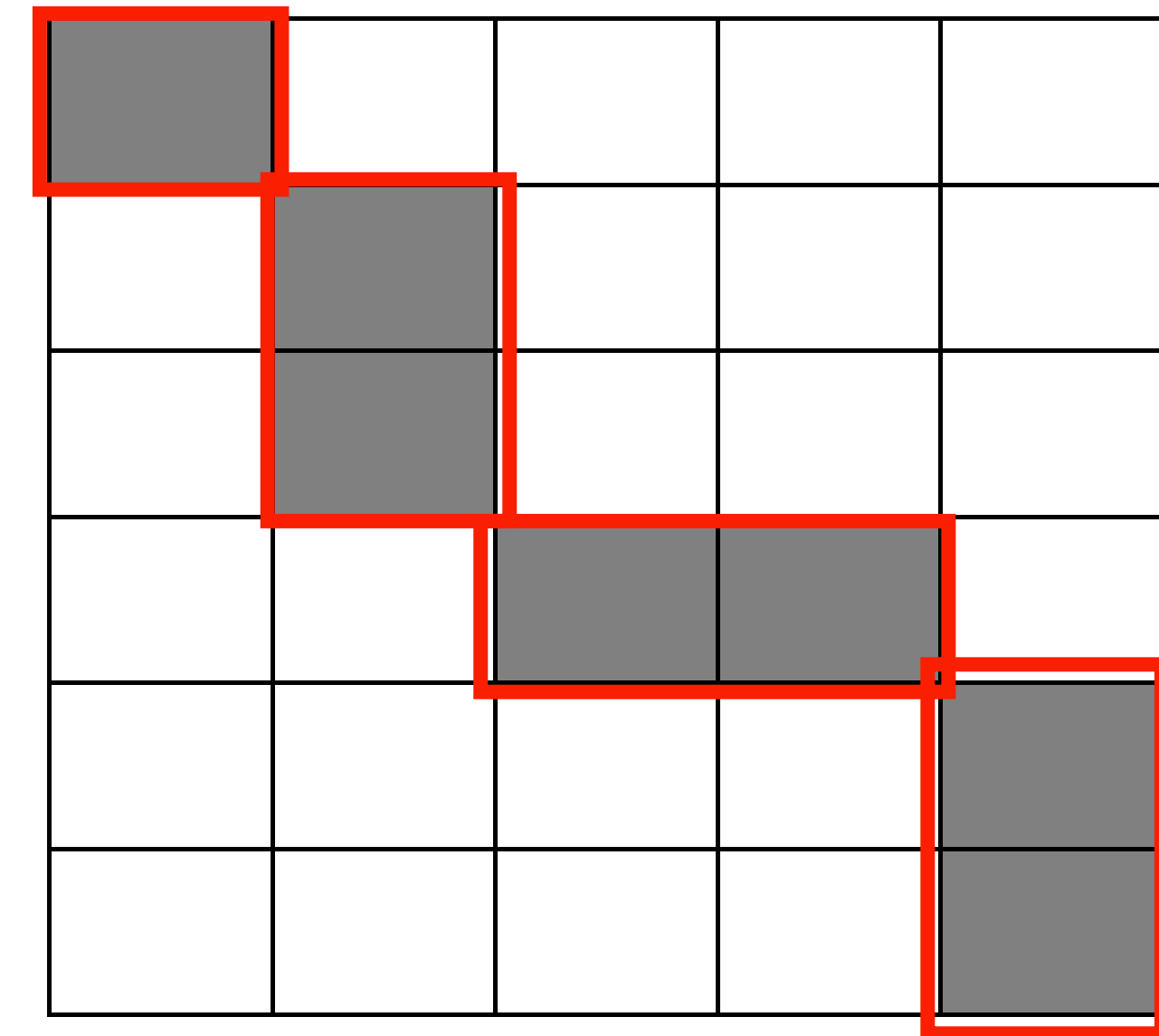$$p(t|e) = p(\text{ไป}|go) = 0.5$$

- **Phrase-based Translation Model**

$$p(t_1 t_2, \ldots, t_n | e_1, e_2, \ldots, e_m) = p(\text{เมื่อวาน นี้}| \text{ previous day}) = 0.1$$

# How we get the P(t|e) for phrases ?

- Phrase Extraction Algorithm

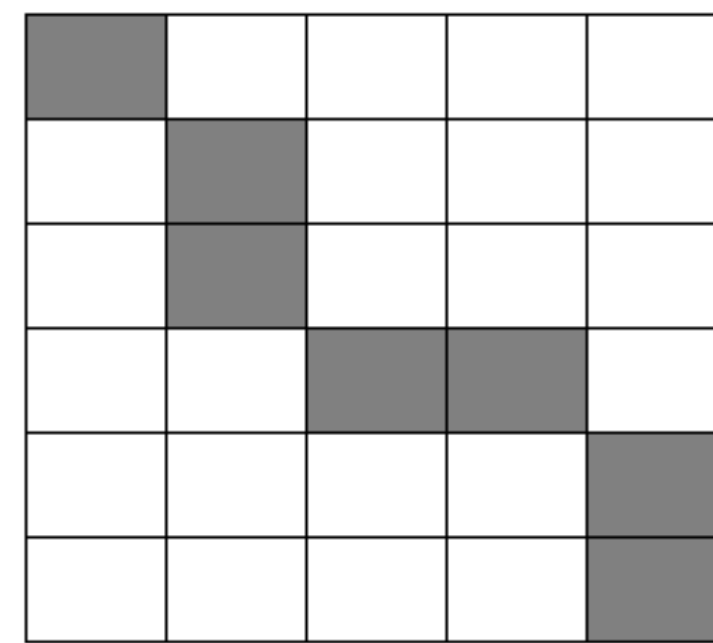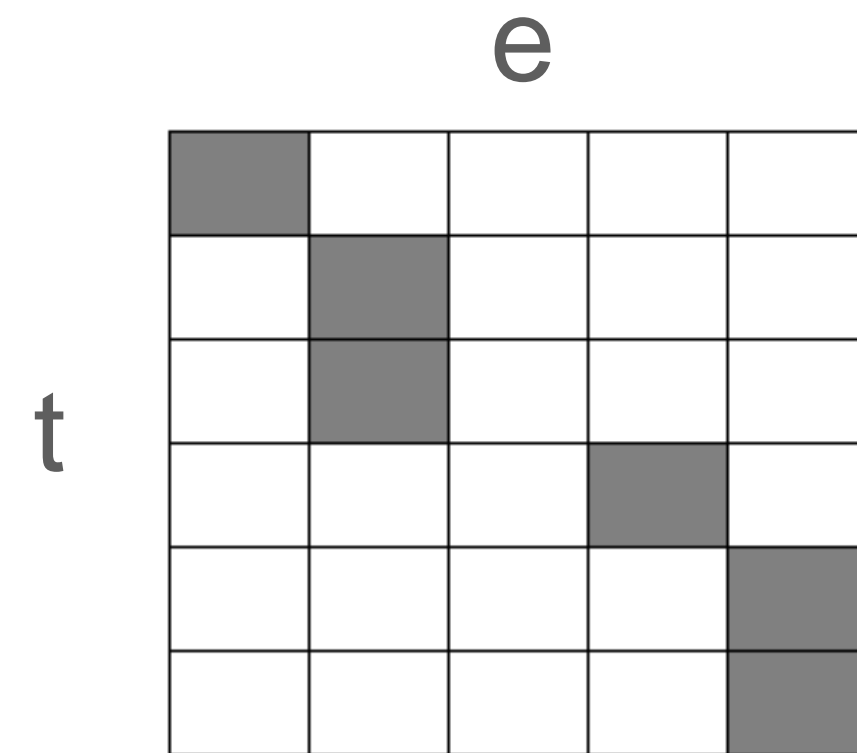  - Expanding single word alignment pairs to multiple-word alignment.
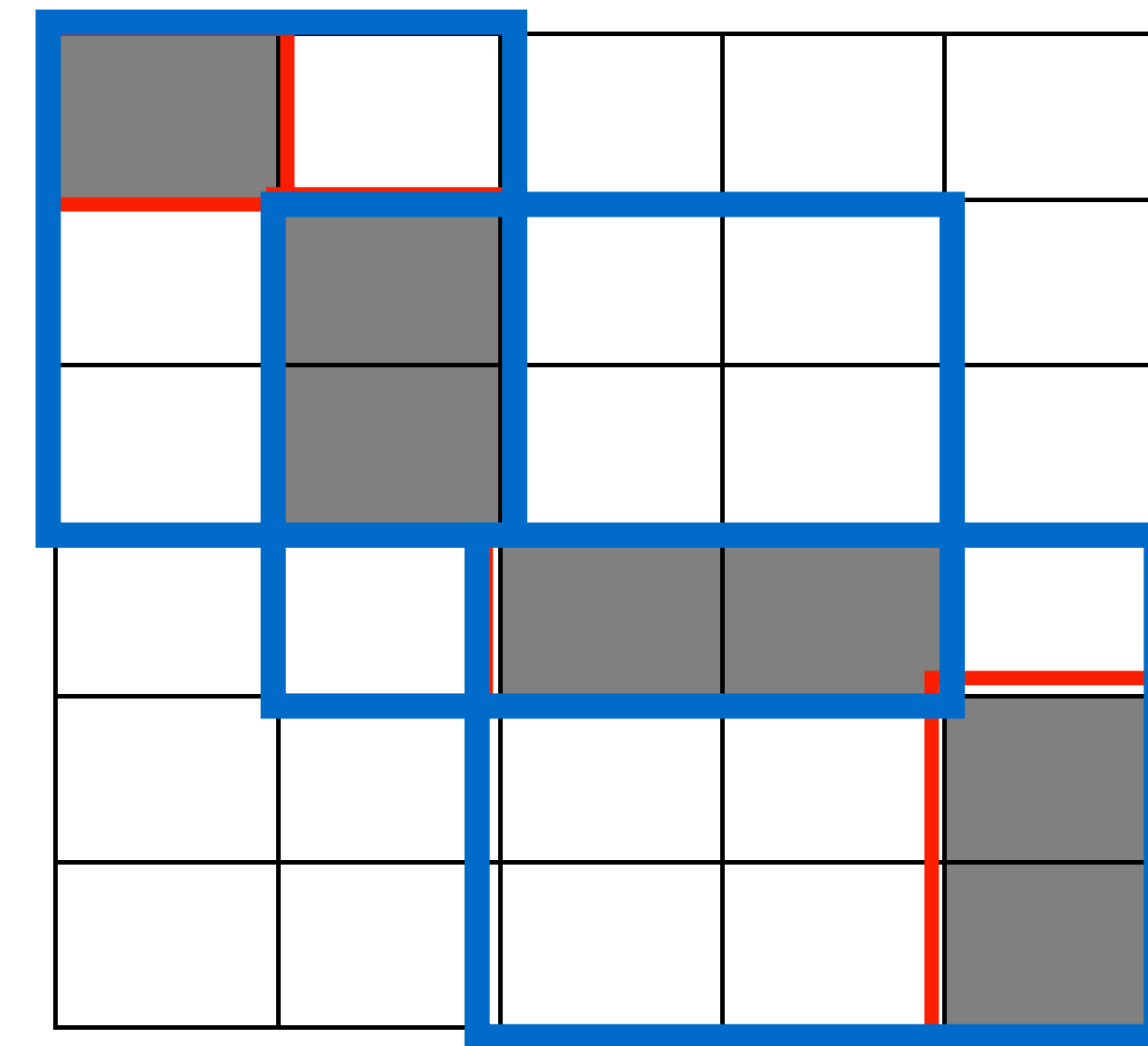
e

t

Fill missing alignment pairs

# How we get the P(t|e) for phrases ?

- Phrase Extraction Algorithm

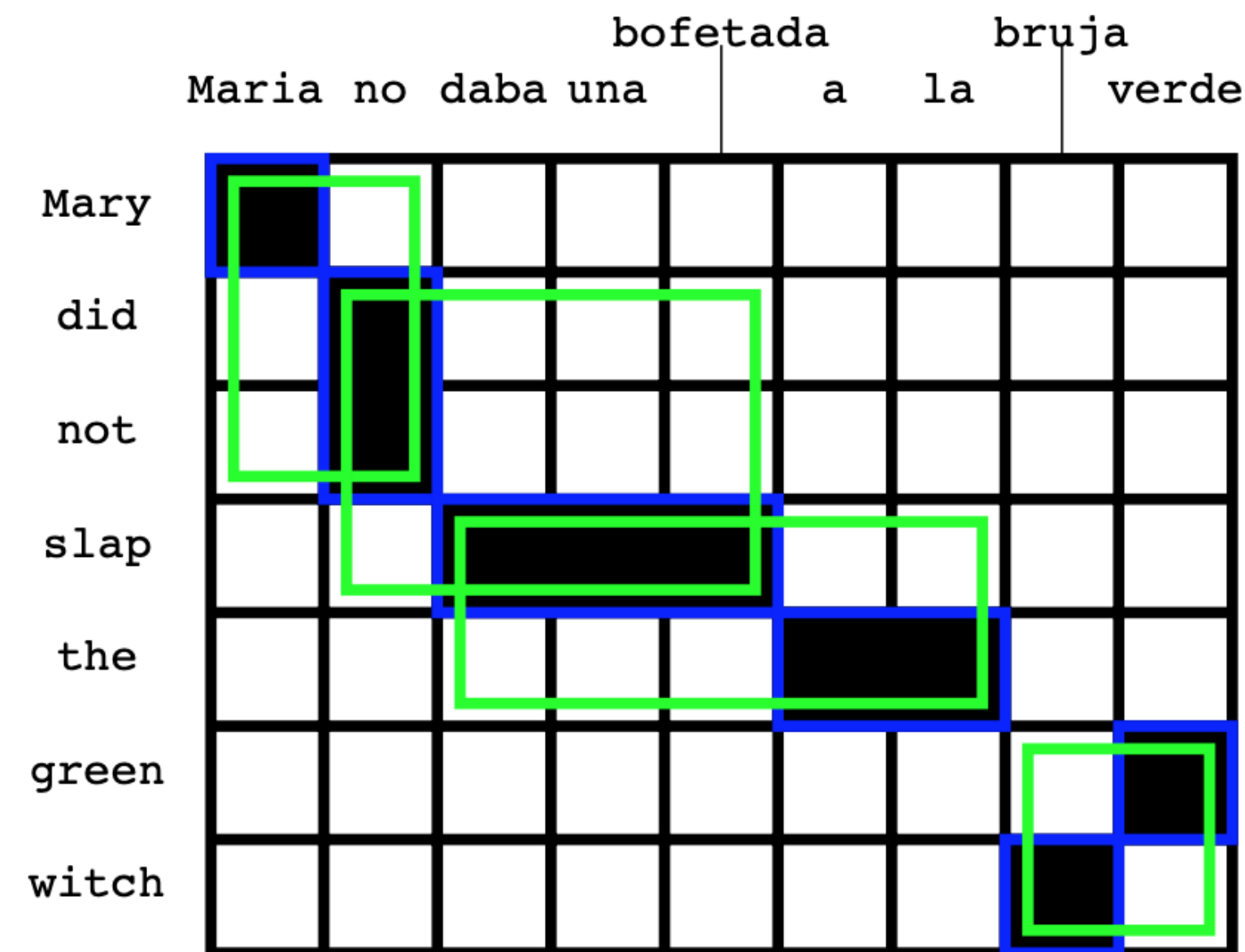  - Expanding single word alignment pairs to multiple-word alignment.



Fill missing alignment pairs

# How we get the P(t|e) for phrases ?

- Phrase Pairs

$$p(t|e) = \frac{count(t,e)}{count(e)}$$



(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch), (verde, green), (Maria no, Mary did not), (no daba una bofetada, did not slap), (daba una bofetada a la, slap the), (bruja verde, green witch)

Image from : www.statmt.org

22

# Statistical Decoder (with Translation model)

| เมื่อวาน | นี้ | ฉัน | ไป | ทะเล | กับ | เพื่อน |
|---|---|---|---|---|---|---|
| Yesterday (0.5) | this (1.0) | I (0.8) | go (0.6) | the sea (0.7) | with (0.8) | friends (0.6) |
| Yesterday (0.7) | | I (0.8) | went to (0.3) | sea (0.2) | with (0.8) | friend (0.4) |
| Previous day (0.1) | | me (0.2) | get to (0.01) | ocean (0.05) | with my friend (0.1) | |

# Language Model $p_{LM}(e)$

- Language Model of the target language.

- Calculate the **"*fluency*"** of sentences

$$p_{LM}(\text{I went to the sea}) > p_{LM}(\text{I went to ocean})$$

# How can we estimate $p_{LM}$ ?

- **N-gram** language model

**1-gram :** $p_{LM}$(I went to the sea) =

$p(I) \times p(went) \times p(to) \times p(the) \times p(sea)$

**2-gram :** $p_{LM}$(I went to the sea) =

$p(I \,|<bos>) \times p(went| \, I) \times p(to|went) \times p(the|to) \times p(sea|the)$

# How can we estimate $P_{LM}$ ?

- **Maximum Likelihood Estimation (MLE)**

- p(I) = count("I") / N

- p(went|I) = count("I went") / count("I")

# Smoothing

- if p(x|y) = 0 ? ➔ P$_{LM}$ = 0

We can use smoothing technique to overcome this situation.

For example, **Add-one smoothing (Laplace Smoothing)**

$$P^*_{\text{Laplace}}(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)+1}{\sum_w (C(w_{n-1}w)+1)} = \frac{C(w_{n-1}w_n)+1}{C(w_{n-1})+V}$$
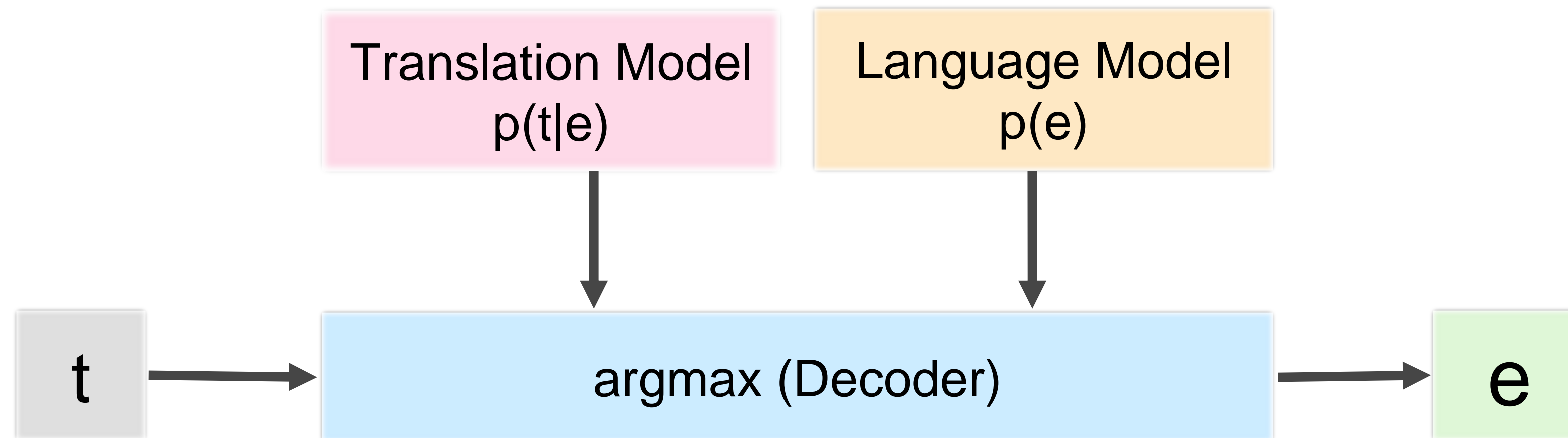
# Smoothing

- **Back-off**

$$\text{if } p(w_i|w_{i-1}) = 0 \text{ then } \hat{p}(w_i|w_{i-1}) = p(w_i)$$

- **Interpolation**

$$\hat{p}(w_i|w_{i-1}) = \lambda_1 p(w_i) + \lambda_2 p(w_i|w_{i-1})$$
$$\lambda_1 + \lambda_2 = 1$$

# Statistical Decoder

$$\hat{e} = \text{argmax}_{e}\, \text{p}(t|e)p(e)$$

# Decoding

- Decoding with Translation Model and Language Model

| เมื่อวาน | นี้ | ฉัน | ไป | ทะเล | กับ | เพื่อน |
|---|---|---|---|---|---|---|
| Yesterday (0.5) | this (1.0) | I (0.8) | go (0.6) | the sea (0.7) | with (0.8) | friends (0.6) |
| Yesterday (0.7) | | I (0.8) | went to (0.3) | sea (0.2) | with (0.8) | friend (0.4) |
| Previous day (0.1) | | me (0.2) | get to (0.01) | ocean (0.05) | with my friend (0.1) | |

Score = 0.7 x 0.8 x 0.6 x 0.7 x 0.8 x 0.6 x
$P_{LM}$("Yesterday I go the sea with friends")

# Decoding

- Decoding with Translation Model and Language Model

| เมื่อวาน | นี้ | ฉัน | ไป | ทะเล | กับ | เพื่อน |
|---|---|---|---|---|---|---|
| Yesterday (0.5) | this (1.0) | I (0.8) | go (0.6) | the sea (0.7) | with (0.8) | friends (0.6) |
| Yesterday (0.7) | | I (0.8) | went to (0.3) | sea (0.2) | with (0.8) | friend (0.4) |
| Previous day (0.1) | | me (0.2) | get to (0.01) | ocean (0.05) | with my friend (0.1) | |

Score = 0.1 x 0.2 x 0.3 x 0.05 x 0.1x
$$P_{LM}(\text{"Previous day me went to ocean with my friend"})$$
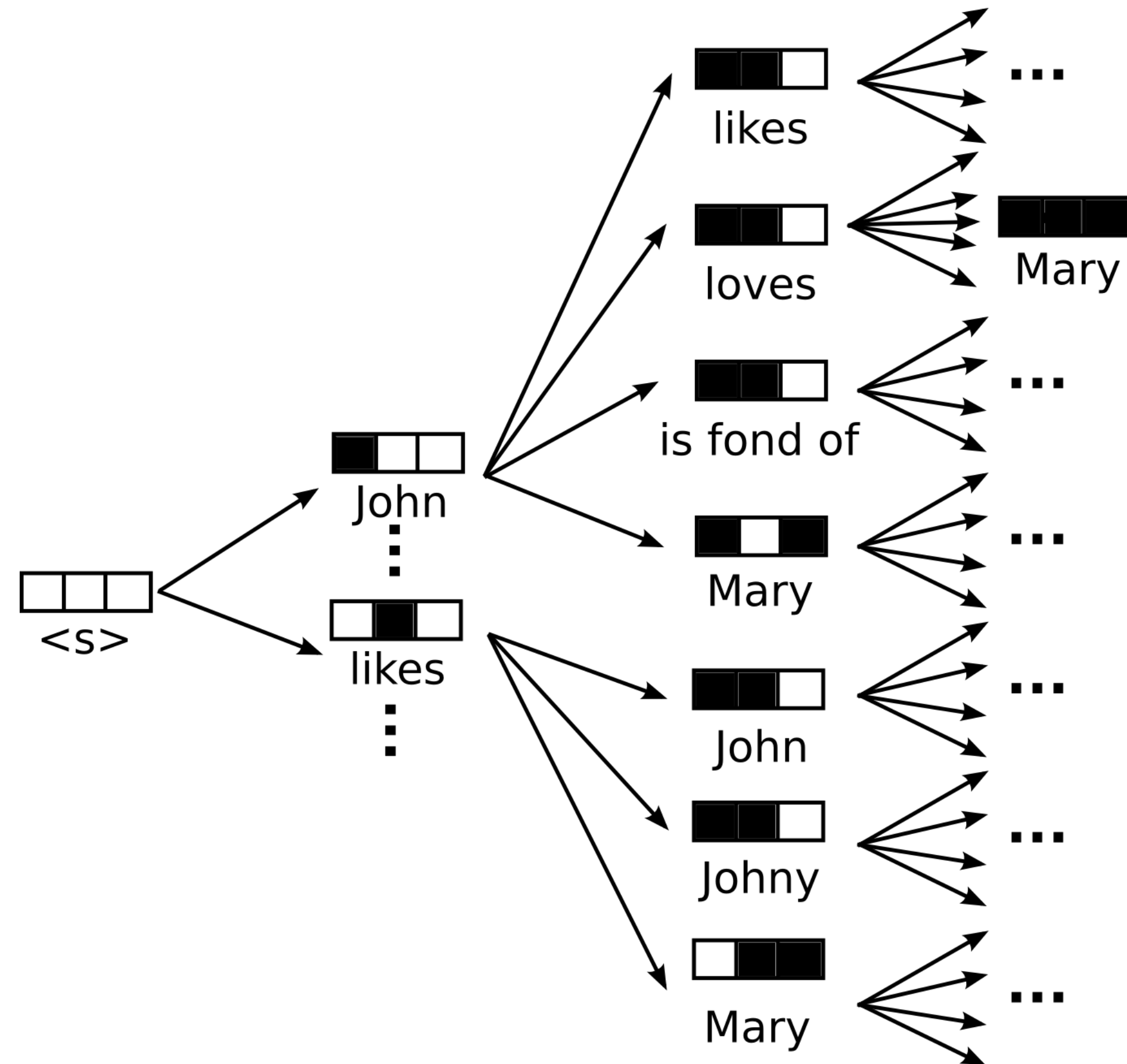
31

# Beam Search

The Beam Search is a tree search algorithm but the data are filtered and sorted using a heuristic function.



Left-to-Right Beam Search

32

# Beam Search



Beam Search in Statistical Decoder
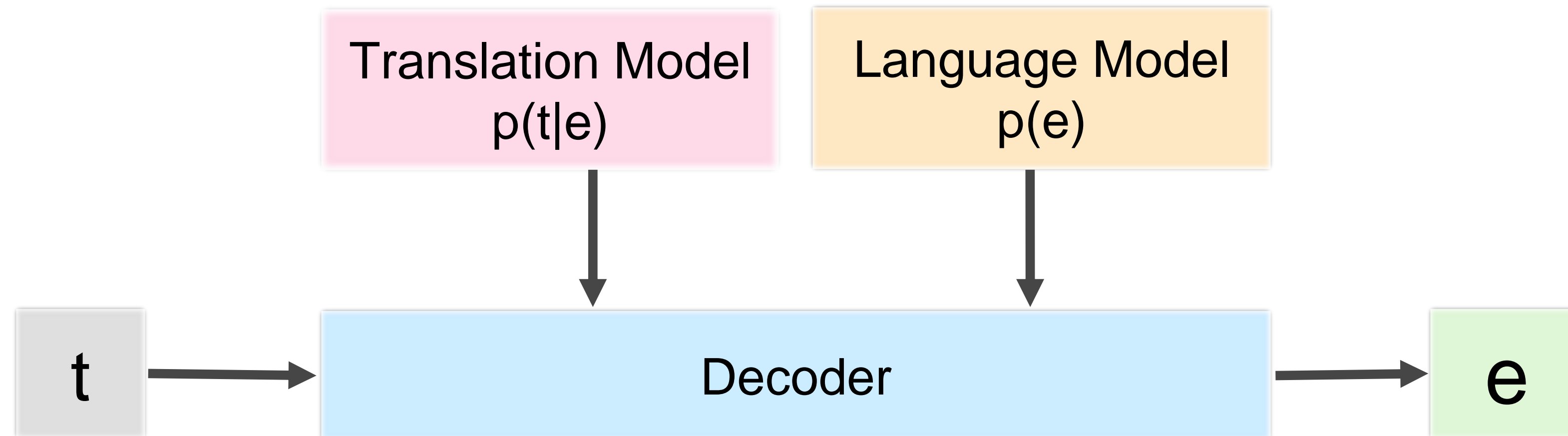
# SMT using

# NLTK Library

```python
from nltk.translate import PhraseTable, StackDecoder
from collections import defaultdict
from math import log

def test_smt():
    phrase_table = PhraseTable()
    phrase_table.add(('book',), ('หนังสือ',), log(0.1))
    phrase_table.add(('this','book',), ('หนังสือ','เล่ม','นี้',), log(0.8))
    phrase_table.add(('this',), ('นี้',), log(0.8))
    phrase_table.add(('costs',), ('ราคา',), log(0.1))
    phrase_table.add(('300',), ('300',), log(0.1))
    phrase_table.add(('300',), ('สาม','ร้อย',), log(0.5))
    phrase_table.add(('baht',), ('บาท',), log(0.8))

    language_prob = defaultdict(lambda: -999.0)
    language_prob[('เล่ม',)] = log(0.5)
    language_prob[('หนังสือ',)] = log(0.4)
    language_prob[('หนังสือ', 'เล่ม', 'นี้')] = log(0.7)
    language_prob[('บาท',)] = log(0.1)
    language_prob[('สาม','ร้อย',)] = log(0.7)
    language_model = type('',(object,),
                        {'probability_change': lambda self, context, phrase: language_prob[ph
                         'probability': lambda self, phrase: language_prob[phrase]})()

    stack_decoder = StackDecoder(phrase_table, language_model)

    out = stack_decoder.translate("this book costs 300 baht".split())
    print(out)
```

```
['หนังสือ', 'เล่ม', 'นี้', 'ราคา', 'สาม', 'ร้อย', 'บาท']
```

# Summary

- Statistical Machine Translation
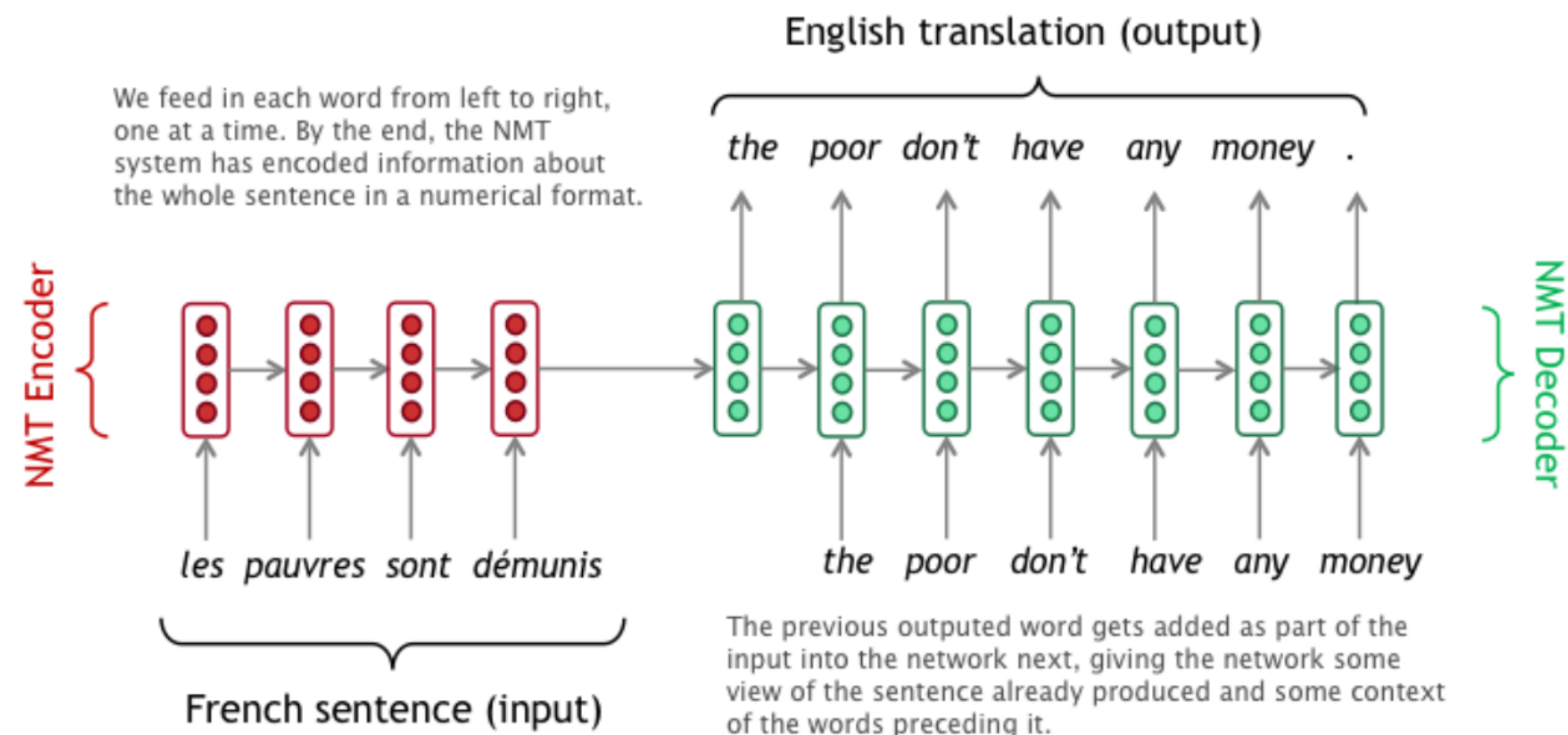
- Translation Model

- Language Model

- Decoder

# SMT Demo

# Machine Translation (Part II)

# Neural Machine Translation

# Neural Machine Translation

- Deep Learning

- End-to-end training

- No word or phrase translation tables required



English translation (output)

We feed in each word from left to right, one at a time. By the end, the NMT system has encoded information about the whole sentence in a numerical format.

the   poor   don't   have   any   money   .

NMT Encoder

NMT Decoder

les   pauvres   sont   démunis

the   poor   don't   have   any   money

French sentence (input)

The previous outputed word gets added as part of the input into the network next, giving the network some view of the sentence already produced and some context of the words preceding it.

# Background

- Recurrent Neural Network (RNN)

- Language model using long short-term memory (LSTM)
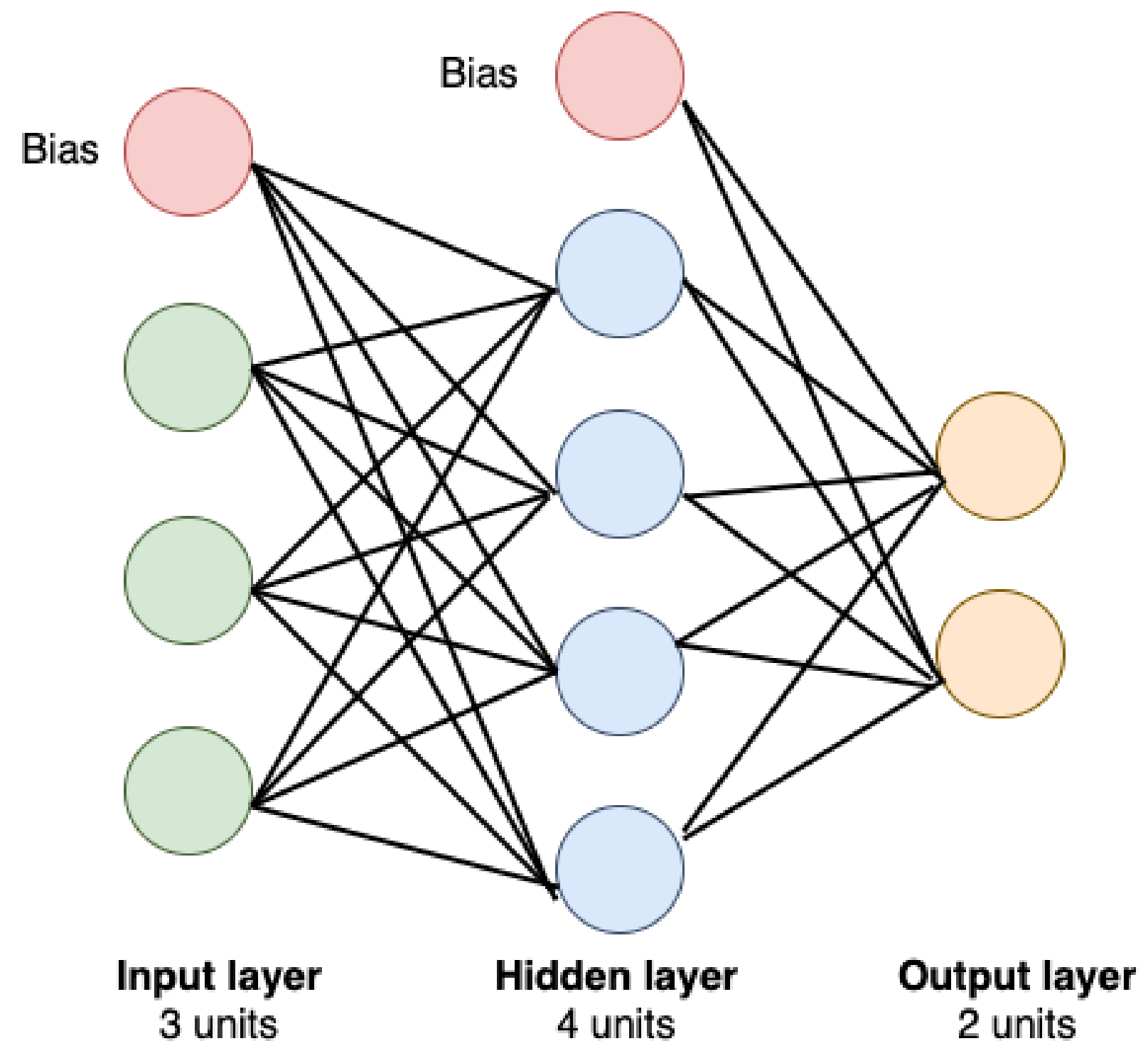
- Negative Log Likelihood (NLL)

- Seq2Seq model



Image courtesy of Chris Olah
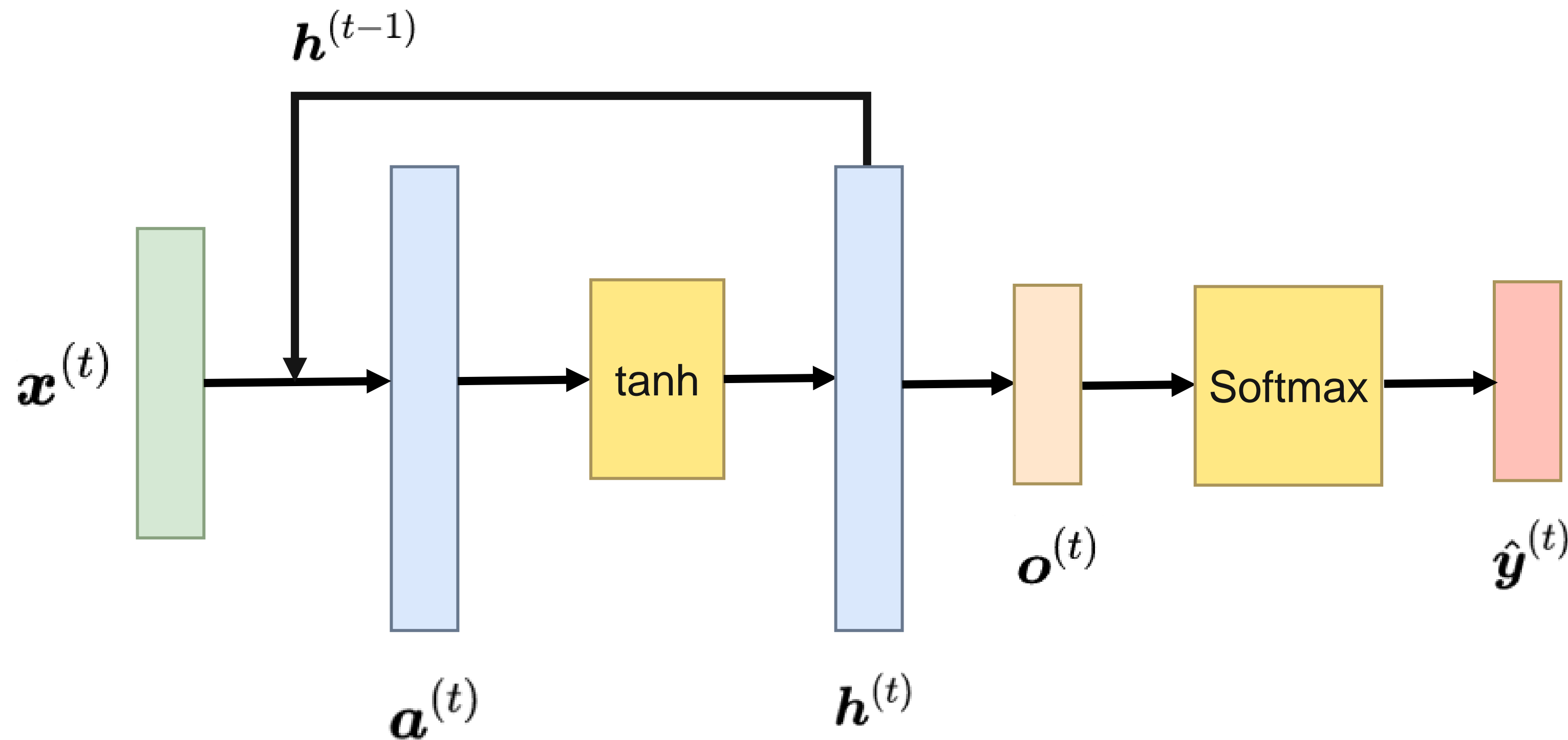
# Neural Network

- Feed Forward Neural Network (FFN)



Simplified version of feed forward neural network

Image taken from https://miro.medium.com/

# Recurrent Neural Network

- Recurrent Neural Network (RNN)

$$\boldsymbol{a}^{(t)} = \boldsymbol{b} + \boldsymbol{W}\boldsymbol{h}^{(t-1)} + \boldsymbol{U}\boldsymbol{x}^{(t)}$$

$$\boldsymbol{h}^{(t)} = \tanh\left(\boldsymbol{a}^{(t)}\right)$$

$$\boldsymbol{o}^{(t)} = \boldsymbol{c} + \boldsymbol{V}\boldsymbol{h}^{(t)}$$

$$\hat{\boldsymbol{y}}^{(t)} = \mathrm{softmax}\left(\boldsymbol{o}^{(t)}\right)$$

$\boldsymbol{h}^{(t-1)}$

$\boldsymbol{x}^{(t)}$

tanh

Softmax

$\boldsymbol{a}^{(t)}$

$\boldsymbol{h}^{(t)}$

$\boldsymbol{o}^{(t)}$

$\hat{\boldsymbol{y}}^{(t)}$

**Softmax Function**

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$$

41

# Recurrent Neural Network

- Recurrent Neural Network (RNN)

# Limitations of Simple Recurrent Unit

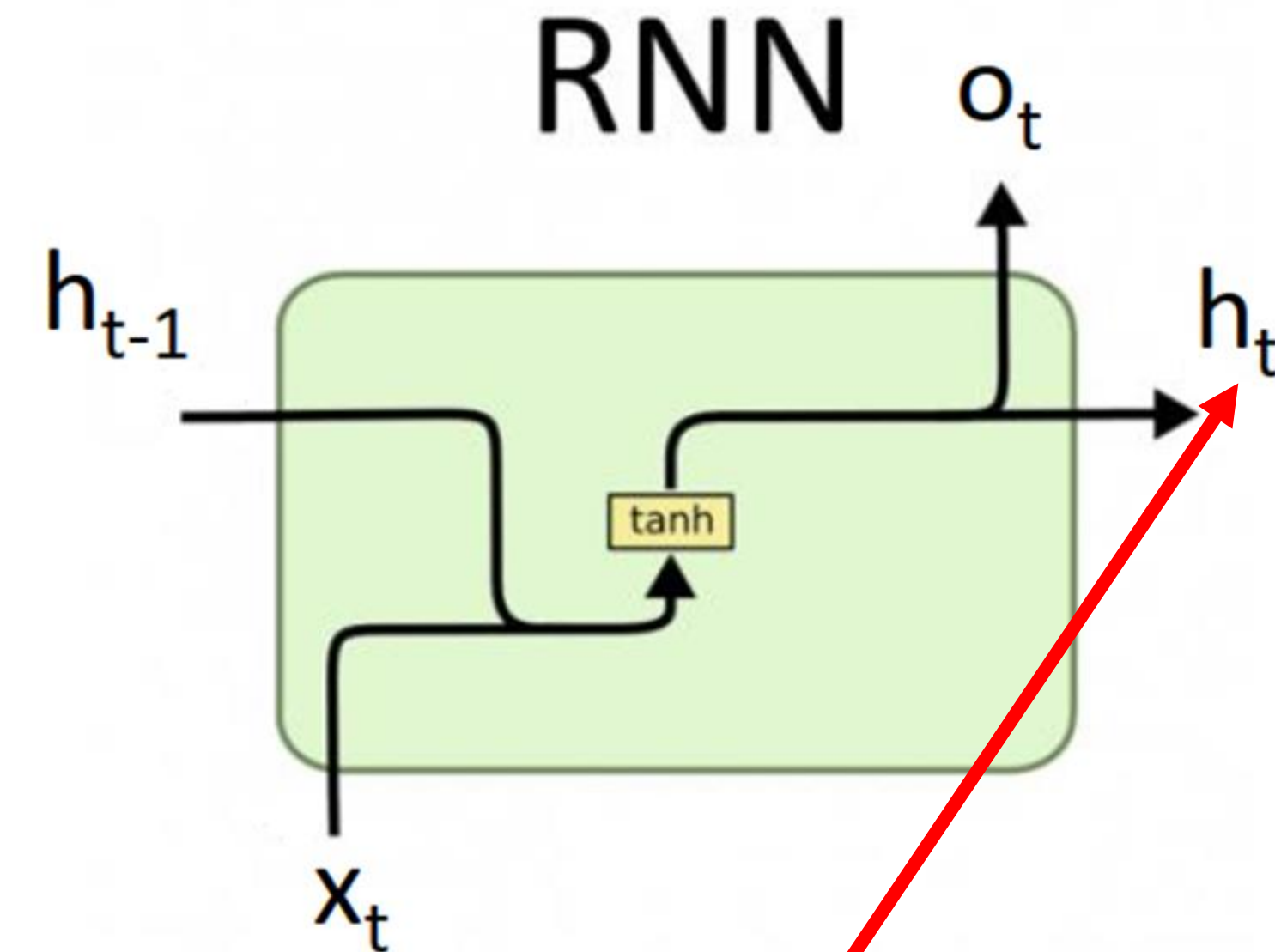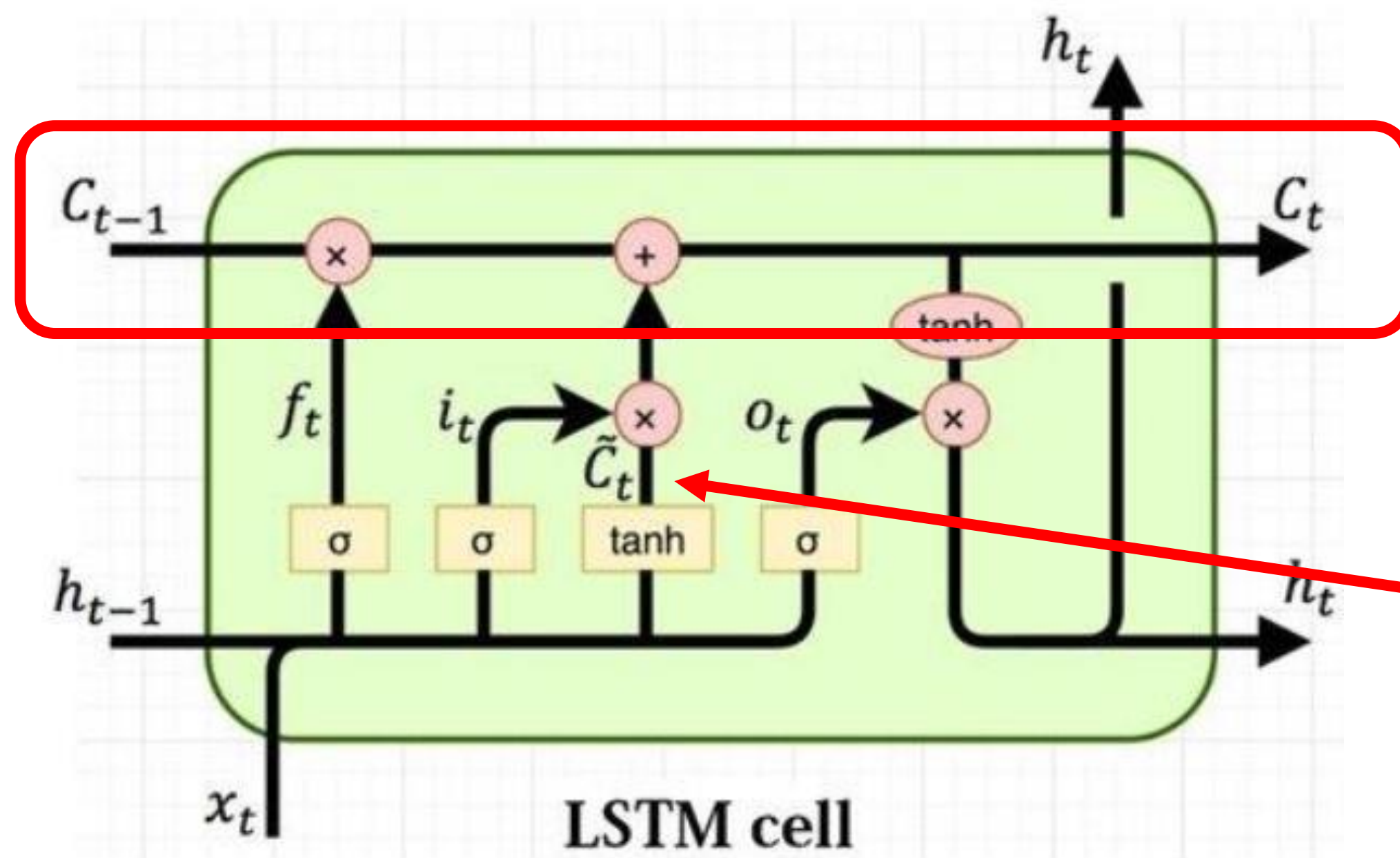- Difficult to train (**gradient vanishing problem**)

- Long distance dependency

เพิ่ง ลง จาก เครื่อง มา หิว มาก เดี๋ยว ช่วย พา ไป หา อะไร _____ หน่อย

{กิน, ทำ, เล่น, ชม, …}

43

# Long short-term memory

- Add "***Cell Memory***" (***C_t***) to RNN units

- Cell memory is controlled by **forget gate** and **input gate**.

- Output is controlled by an **output gate**.

# Long short-term memory



RNN

$$\tilde{C}_t = \tanh(x_t U^g + h_{t-1} W^g)$$

$$C_t = \sigma\Big(f_t * C_{t-1} + i_t * \tilde{C}_t\Big)$$

$$h_t = \tanh(C_t) * o_t$$

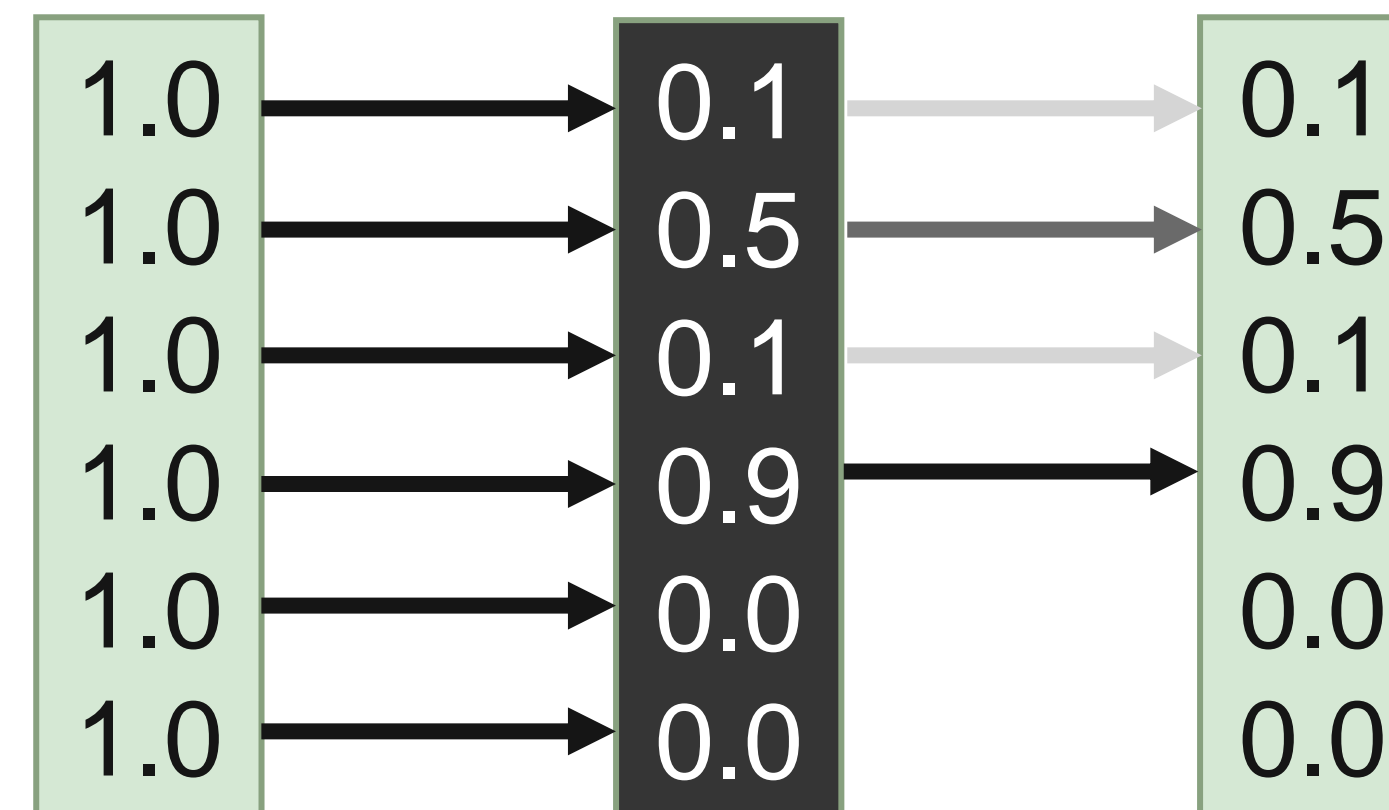LSTM cell

# How gates work ?

$$\tilde{C}_t = \tanh(x_t U^g + h_{t-1} W^g)$$

Forget gate ← $C_t = \sigma\left(\boxed{f_t} * C_{t-1} + \boxed{i_t} * \tilde{C}_t\right)$ → Input gate
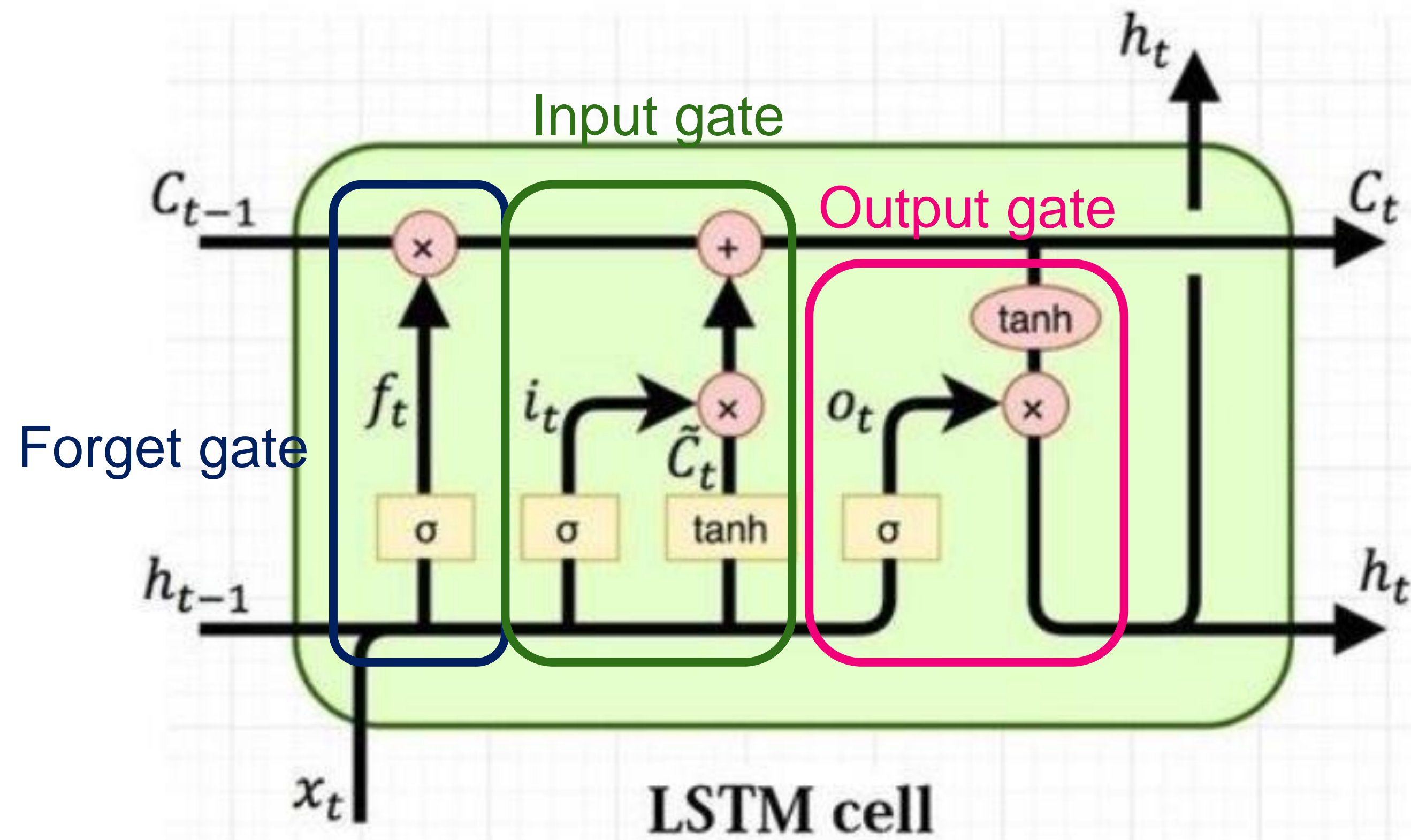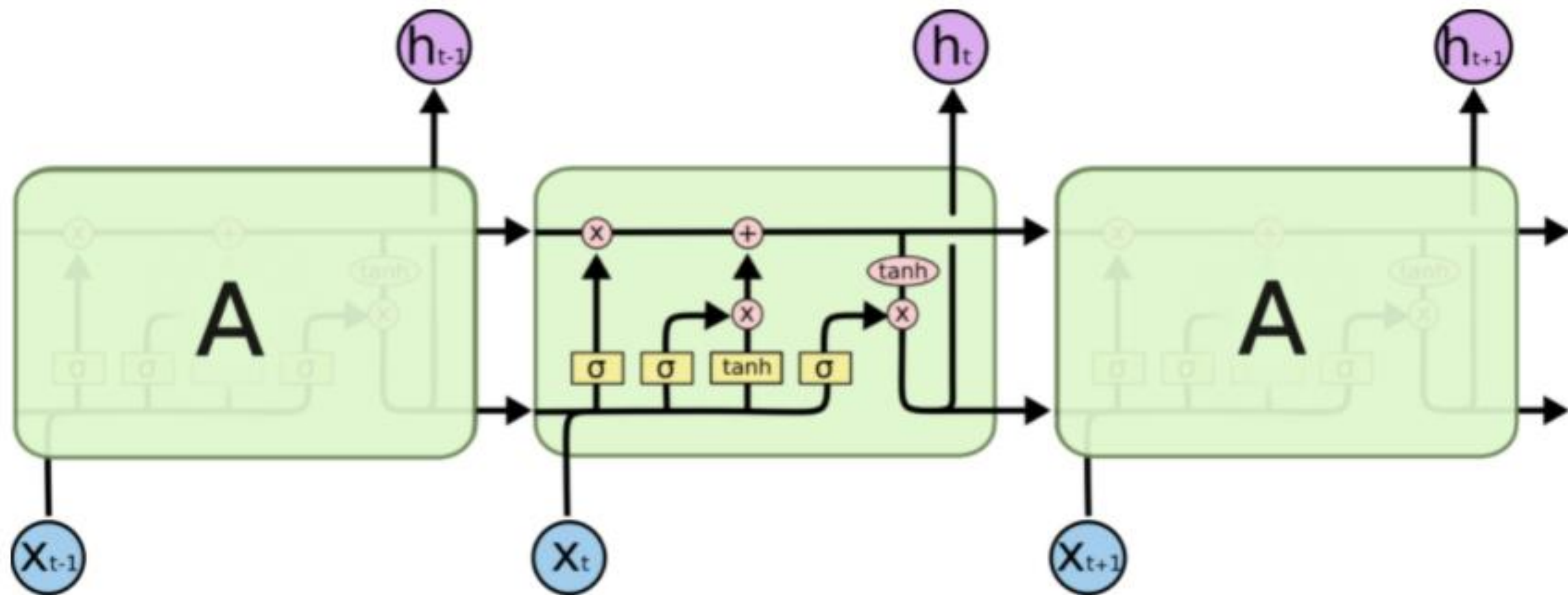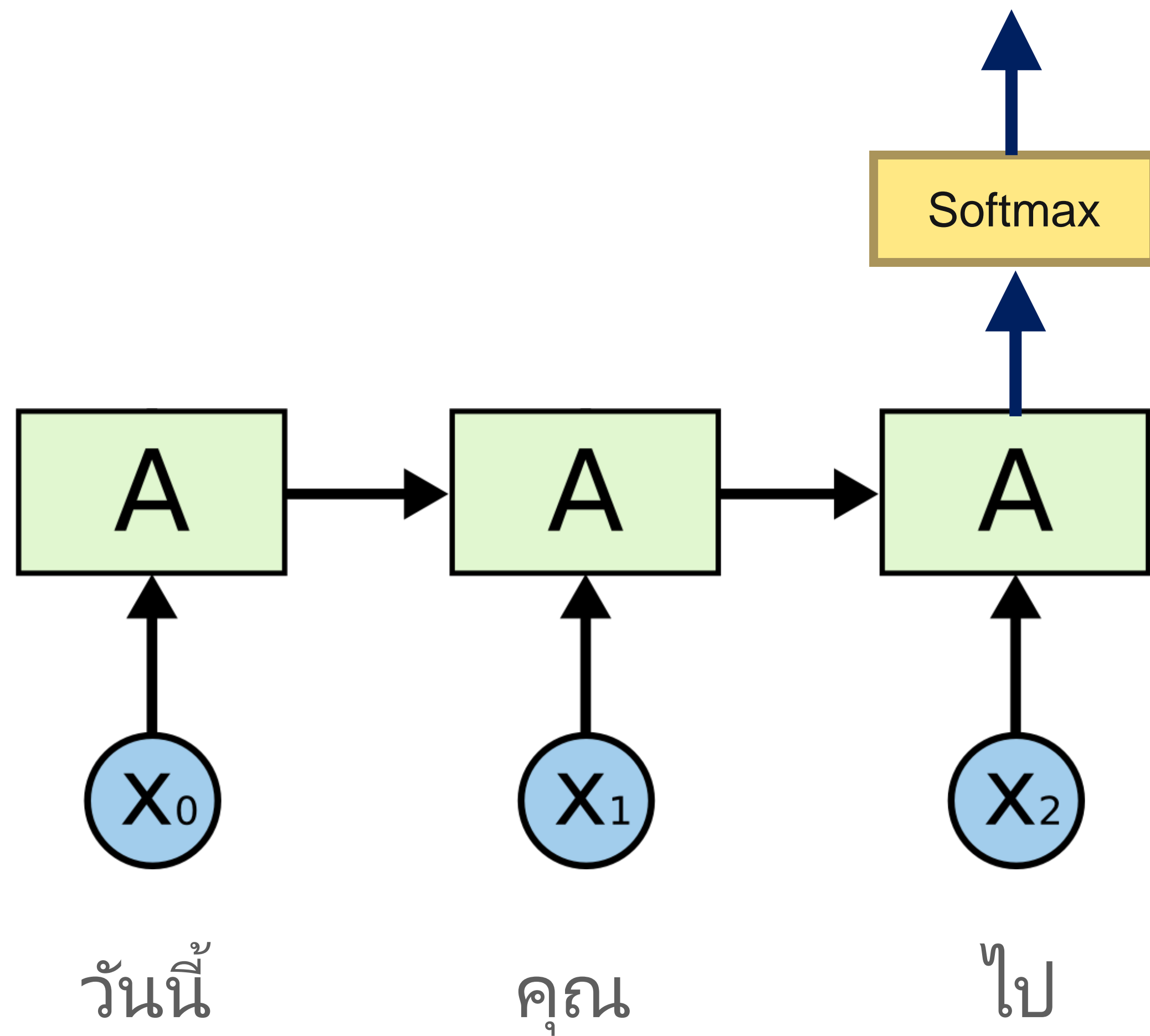
$$h_t = \tanh(C_t) * \boxed{o_t}$$

Output gate

| 1.0 | 0.1 | 0.1 |
| 1.0 | 0.5 | 0.5 |
| 1.0 | 0.1 | 0.1 |
| 1.0 | 0.9 | 0.9 |
| 1.0 | 0.0 | 0.0 |
| 1.0 | 0.0 | 0.0 |

\* Element-wise multiplication

Gate

# How gates work ?



Input gate

Output gate

Forget gate

$C_{t-1}$

$f_t$

$i_t$

$\tilde{C}_t$

$o_t$

tanh

$\sigma$    $\sigma$    tanh    $\sigma$

$h_{t-1}$

$x_t$

$h_t$

$C_t$

$h_t$

**LSTM cell**

$$i_t = \sigma\big(x_t U^i + h_{t-1} W^i\big)$$

$$f_t = \sigma\big(x_t U^f + h_{t-1} W^f\big)$$

$$o_t = \sigma(x_t U^o + h_{t-1} W^o)$$

$\sigma$   = sigmoid function

$$\tilde{C}_t = \tanh(x_t U^g + h_{t-1} W^g)$$

$$C_t = \sigma\big(f_t * C_{t-1} + i_t * \tilde{C}_t\big)$$

$$h_t = \tanh(C_t) * o_t$$

# LSTM

# RNN-based Language Model

$$p(x_t | x_{t-1}, x_{t-2}, x_{t-3}) = \text{softmax}\ (W h_t)$$
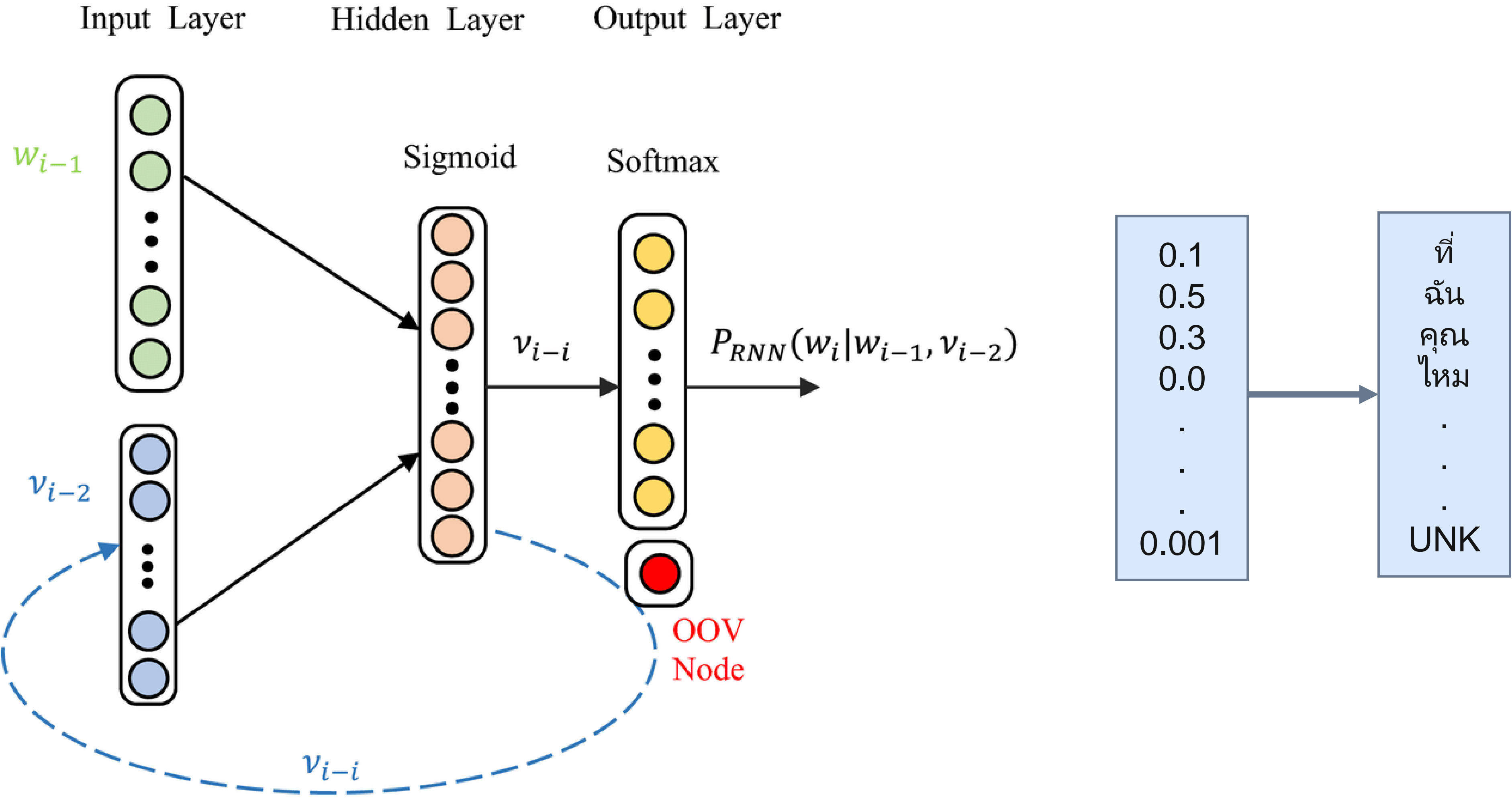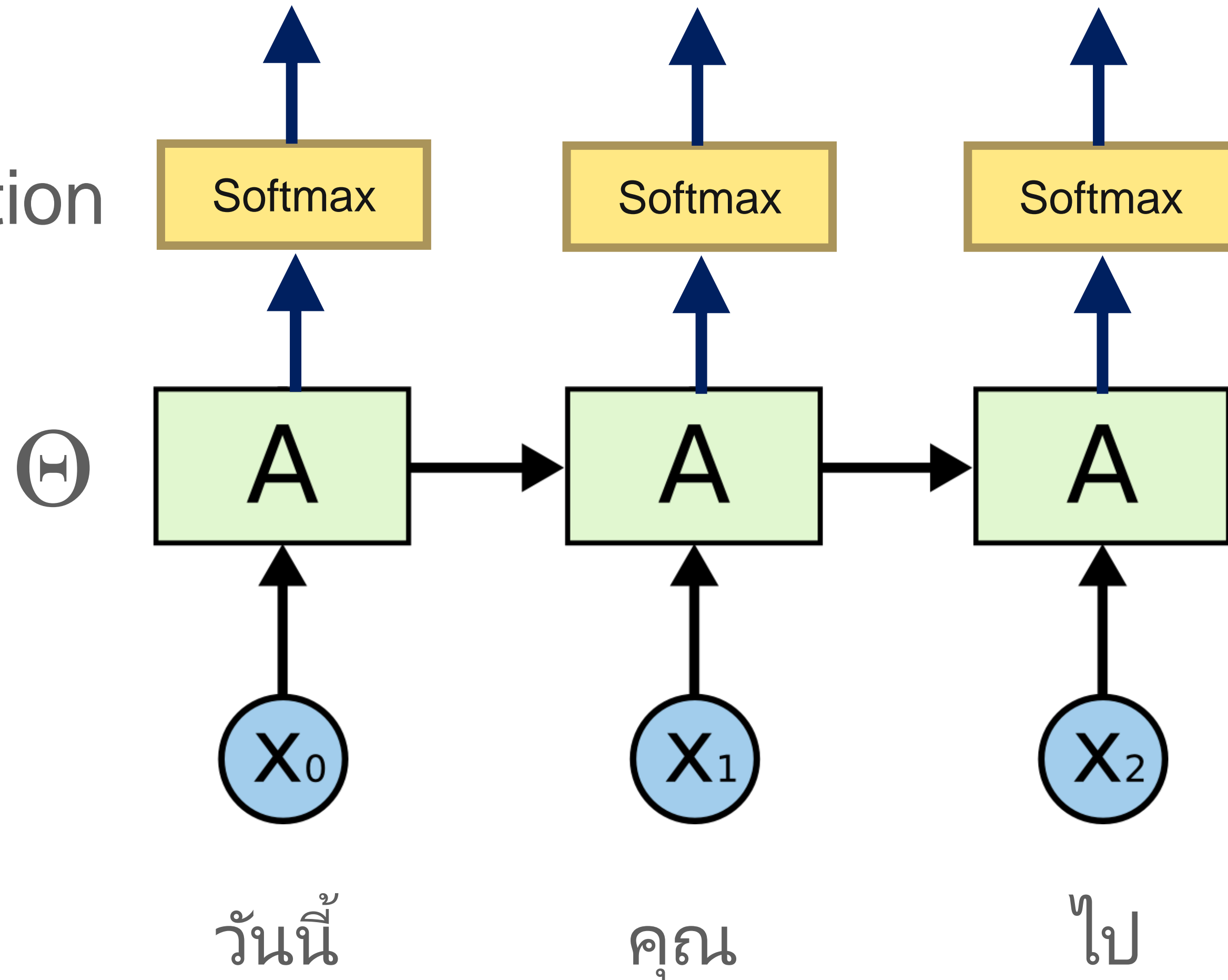
# RNN-based Language Model



Image taken from https://www.researchgate.net/figure/An-architecture-of-recurrent-neural-network-language-model-for-speech-recognition_fig2_333355077

# Training RNN-based Language Model

$$\Theta = \text{argmax}_\Theta \ p(\text{คุณ}|h_0; \Theta) \times \ p(\text{ไป}|h_1; \Theta) \times \ p(\text{ไหน}|h_2; \Theta)$$

Maximum Likelihood Estimation
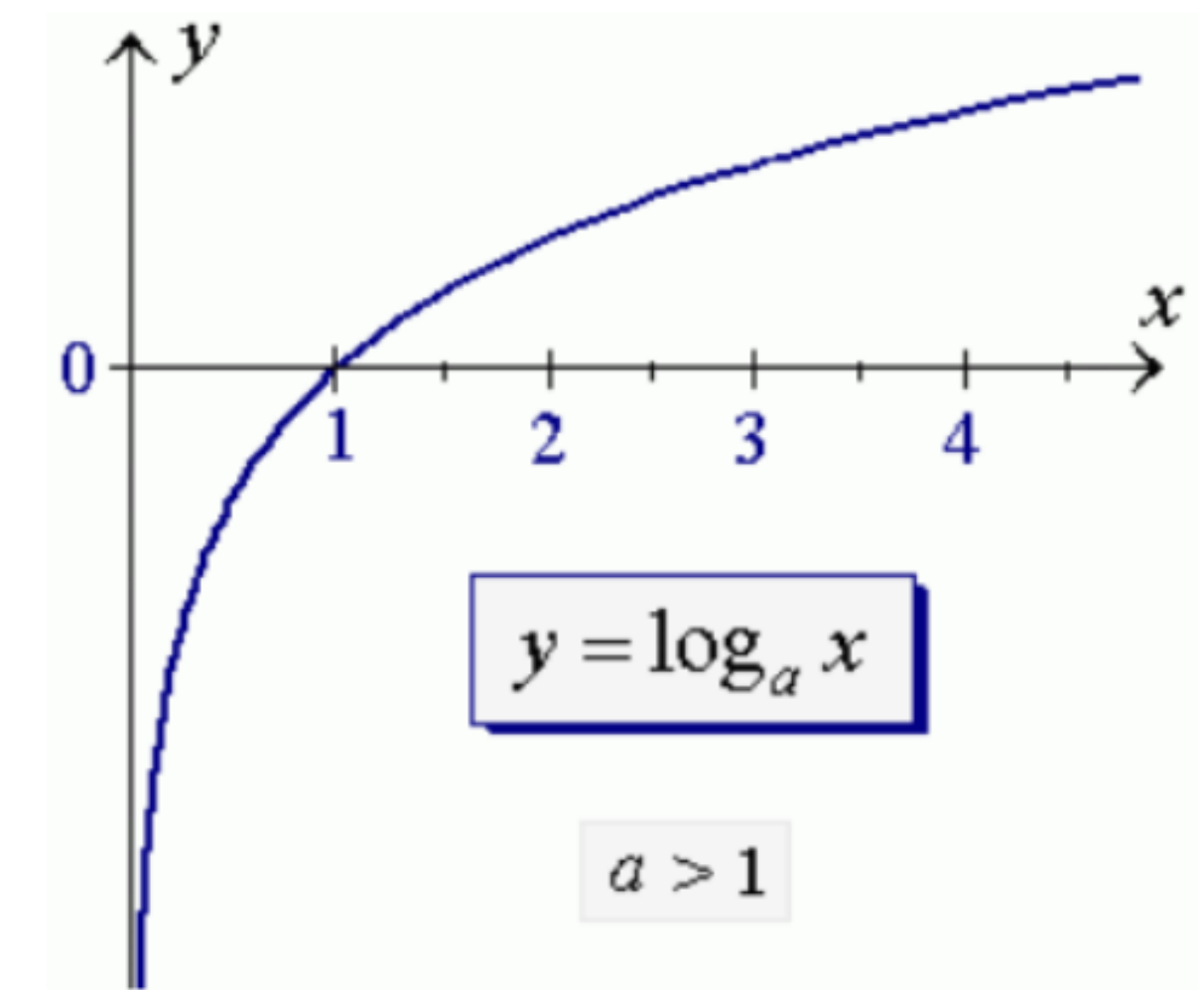
$$L(\theta) = \prod_{i=1}^{n} f(x_i|\theta)$$

$\Theta$

| Softmax | Softmax | Softmax |

| A | A | A |

$x_0$  $x_1$  $x_2$

วันนี้  คุณ  ไป

51

# Negative Log Loss (NLL)

- In many deep learning frameworks, we usually minimize loss functions.

$\Theta = \text{argmax}_\Theta \ p(\text{คุณ}|h_0; \Theta) \times p(\text{ไป}|h_1; \Theta) \times (\text{ไหน}|h_2; \Theta)$

$\Theta = \text{argmax}_\Theta \ \log p(\text{คุณ}|h_0; \Theta) + \log p(\text{ไป}|h_1; \Theta) + \log p(\text{ไหน}|h_2; \Theta)$

$\Theta = \text{argmin}_\Theta \ - \log p(\text{คุณ}|h_0; \Theta) - \log p(\text{ไป}|h_1; \Theta) - \log p(\text{ไหน}|h_2; \Theta)$

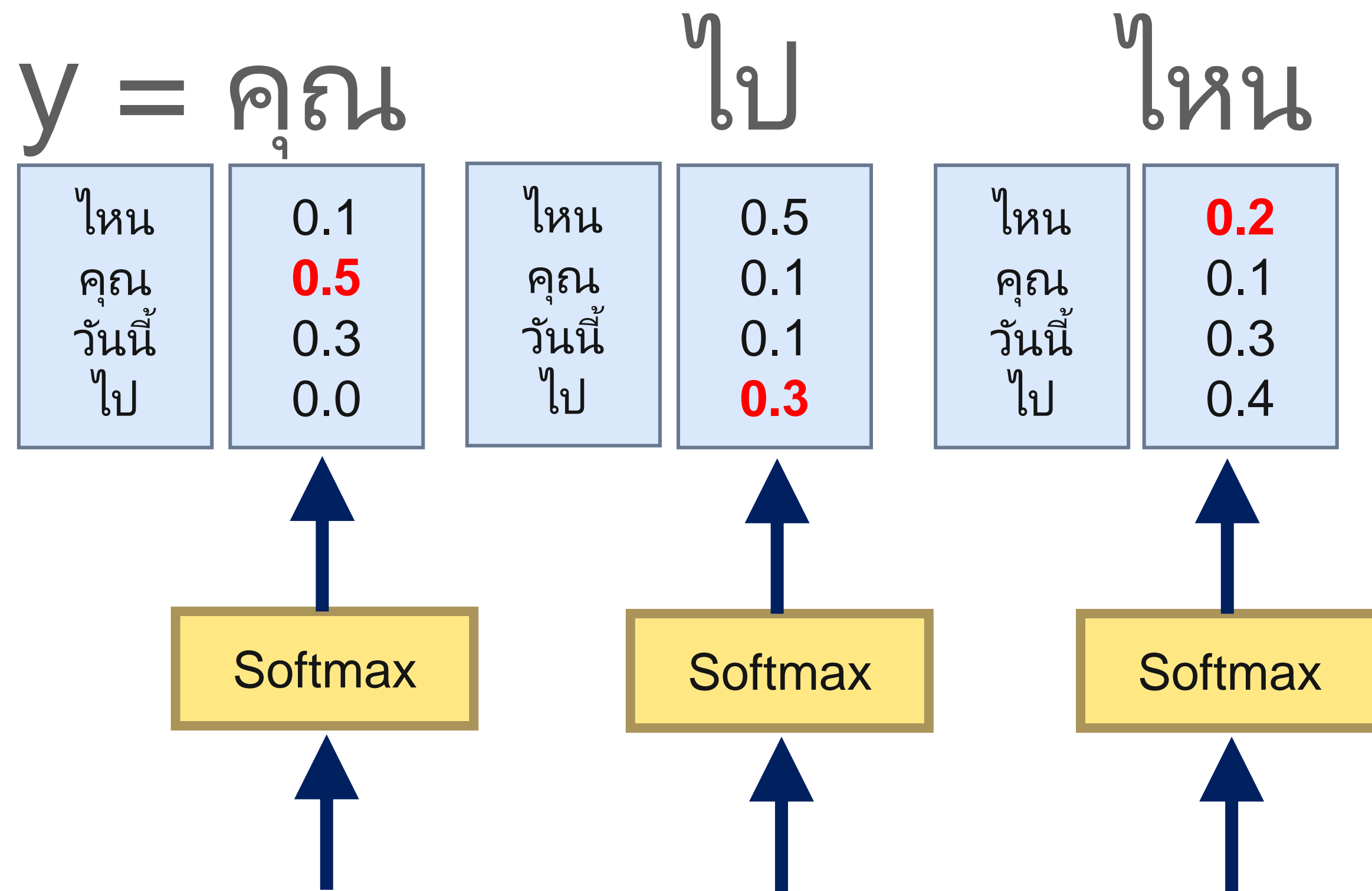$$\text{NLL} = -\sum_{t=1}^{T}\sum_{j=1}^{V} y_j^t \log \hat{y}_j^t$$

$y = \log_a x$

$a > 1$

# Negative Log Loss (NLL)

$$\text{NLL} = -\sum_{t=1}^{T}\sum_{j=1}^{V} \text{y}_j^t \log \hat{y}_j^t$$

- Because $y^t$ are one-hot vectors [ 0 , 0 , …, 1, …, 0, 0 ] , only the position of the ground truth is 1, therefore

$$\text{NLL} = -\sum_{t=1}^{T} \log p(y^t | h_t ; \theta)$$

53

# Negative Log Loss (NLL)

y = คุณ ไป ไหน

| | | | | | |
|---|---|---|---|---|---|
| ไหน | 0.1 | ไหน | 0.5 | ไหน | **0.2** |
| คุณ | **0.5** | คุณ | 0.1 | คุณ | 0.1 |
| วันนี้ | 0.3 | วันนี้ | 0.1 | วันนี้ | 0.3 |
| ไป | 0.0 | ไป | **0.3** | ไป | 0.4 |

Softmax  Softmax  Softmax

$$\text{NLL} = -\sum_{t=1}^{T} \log p(y^t | h_t; \theta)$$

$$= - ( \log 0.5 + \log 0.3 + \log 0.2)$$

54

# Sequence-to-Sequence



English translation (output)

We feed in each word from left to right, one at a time. By the end, the NMT system has encoded information about the whole sentence in a numerical format.

the poor don't have any money .

NMT Encoder

NMT Decoder

les pauvres sont démunis

the poor don't have any money

French sentence (input)

The previous outputed word gets added as part of the input into the network next, giving the network some view of the sentence already produced and some context of the words preceding it.

Image taken from https://medium.com/analytics-vidhya/seq2seq-models-french-to-english-translation-using-encoder-decoder-model-with-attention-9c05b2c09af8

55

# Sequence-to-Sequence

- Encoder : reads source input tokens one-by-one and produce a context vector which represents the source sentence.

- Decoder : generates target tokens from left to right using the context vector and previous output (similar to what language modeling does)



$$p(y_j | s_j, y_{j-1}) = \text{softmax}(W_o t_j)$$

$$t_j = \tanh(W_{t1} s_j + W_{t2} \text{E}[y_{j-1}])$$

56

# Sequence-to-Sequence



Image taken from suriyadeepan.github.io

57

# Training Sequence-to-Sequence

$$\Theta = \text{argmax}_{\Theta} \; p(\text{คุณ}|h_0; \Theta) \; \text{x} \; p(\text{ไป}|h_1; \Theta) \; \text{x} \; p(\text{ไหน}|h_2; \Theta)$$

Maximum Likelihood Estimation

$$L(\theta) = \prod_{i=1}^{n} f(x_i|\theta)$$



$$\Theta = \{\Theta_E, \Theta_D\}$$

58

# Decoding – Greedy

# Decoding – Greedy

# Decoding – Greedy

| | |
|---|---|
| ที่ | 0.2 |
| คุณ | 0.1 |
| วันนี้ | 0.5 |
| ไป | 0.2 |

Score = 0.5 x 0.5

# Decoding – Beam Search

# Summary

- Recurrent Neural Network

- RNN-based Language Model

- Seq2Seq

# Neural Machine Translation with Attention

# Limitations of Seq2Seq Model

- Inefficient for long sentences

- Only good for language pairs with less grammatical variations

# Attention

$$p(y_j \mid s_j, y_{j-1}, c_j) = \text{softmax}(W_o t_j)$$
$$t_j = \tanh(W_{t1} s_j + W_{t2} \mathrm{E}[y_{j-1}] + W_{t3} c_j)$$

# Attention

$$c_j = \text{ATT}(C, s'_j) = \sum_{i=1}^{I} \alpha_{ij} h_i,$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{I} \exp(e_{kj})}$$

$$e_{ij} = v_a^{\text{T}} \tanh(U_a s'_j + W_a h_i)$$
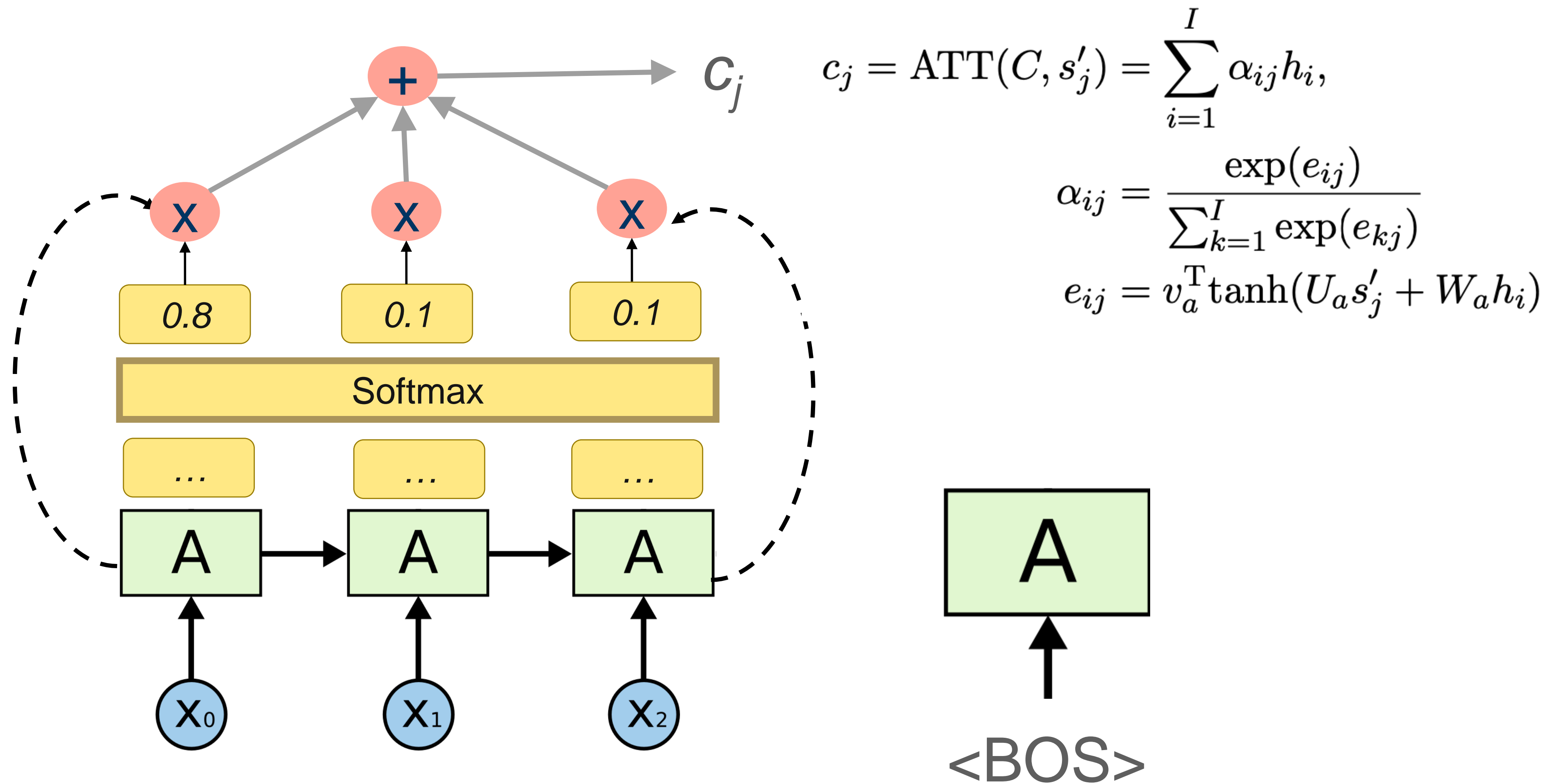
# Attention – Step 1

Intermediate state

$$s'_j$$

NN

A → A → A

A

$x_0$  $x_1$  $x_2$

<BOS>

# Attention – Step 2

$$e_{ij} = v_a^{\mathrm{T}} \tanh(U_a s_j' + W_a h_i)$$



Intermediate state

# Attention – Step 2

$$e_{ij} = v_a^{\mathrm{T}}\tanh(U_a s_j' + W_a h_i)$$

Intermediate state



70

# Attention – Step 3

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{I} \exp(e_{kj})}$$
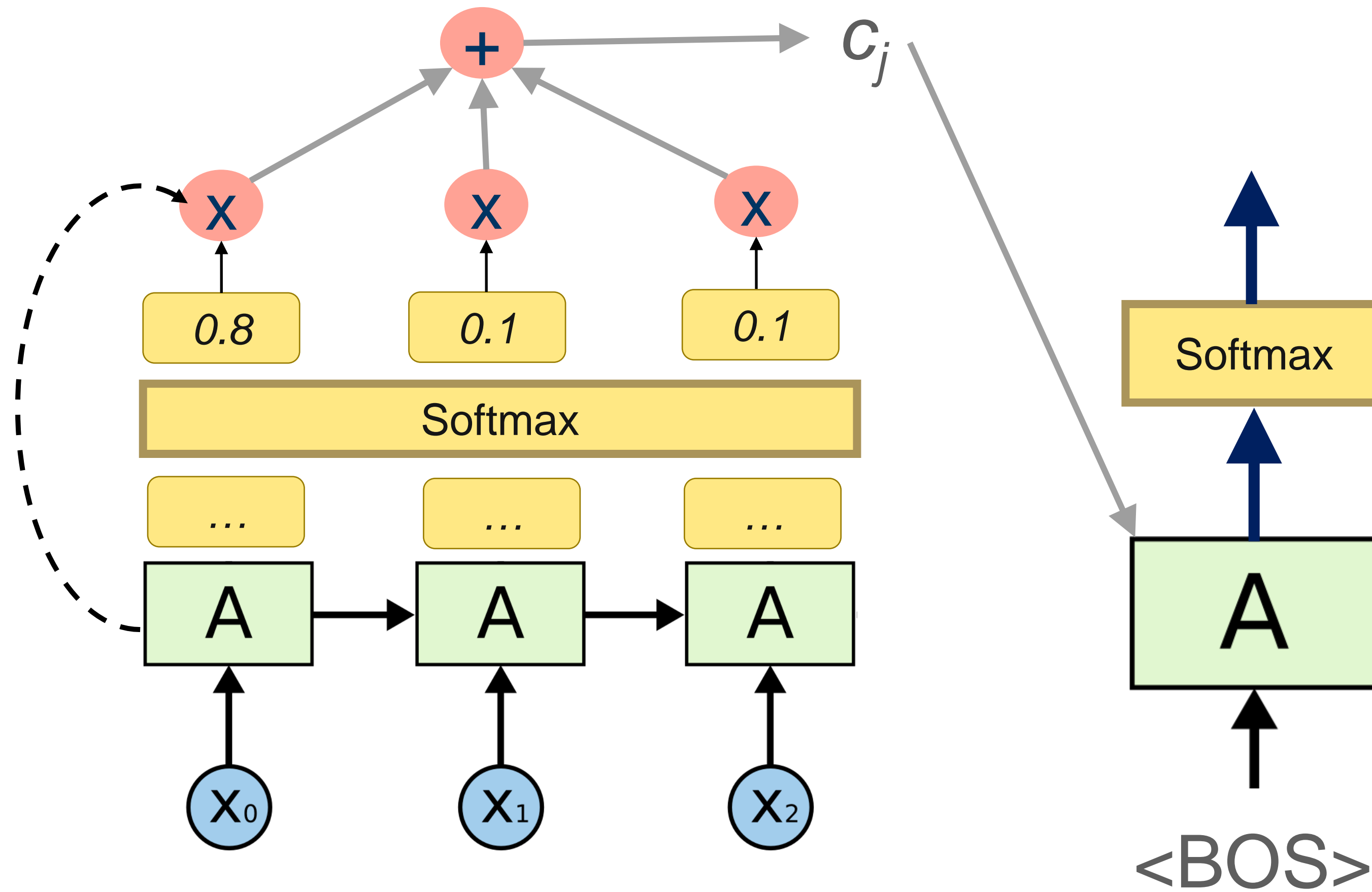
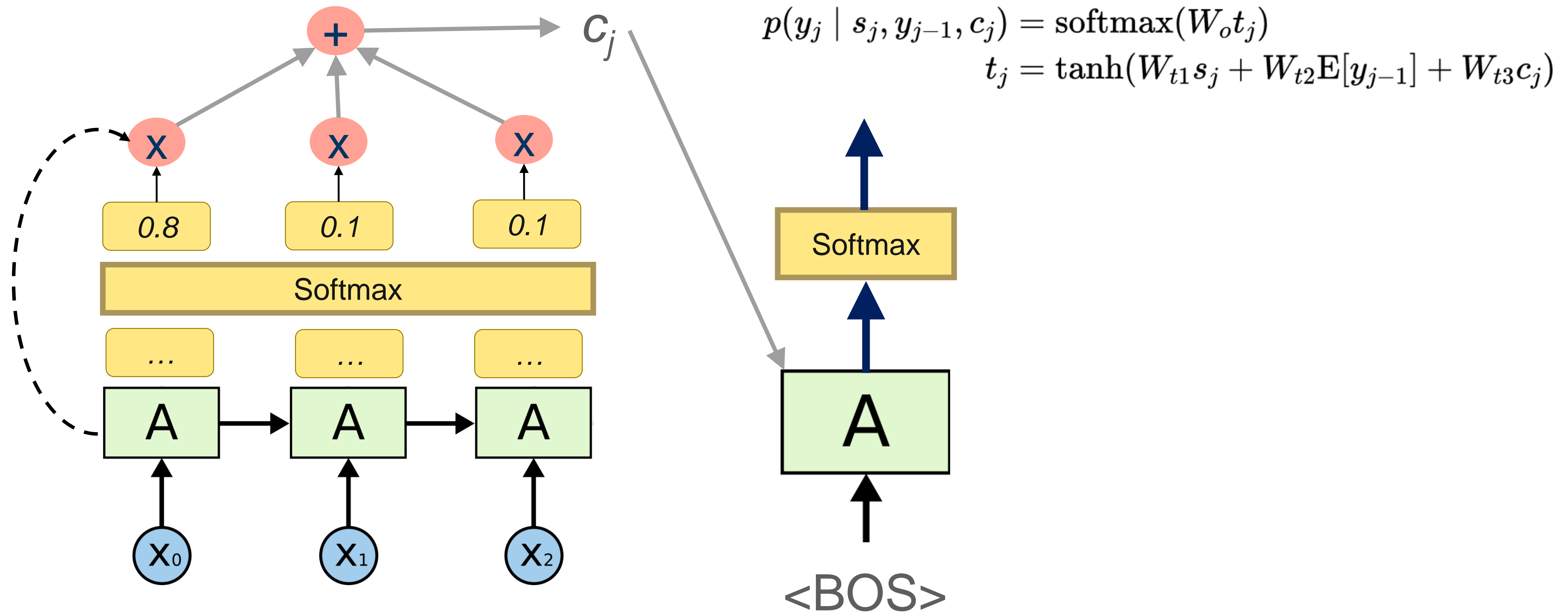$$e_{ij} = v_a^{\mathrm{T}} \tanh(U_a s'_j + W_a h_i)$$

# Attention – Step 4



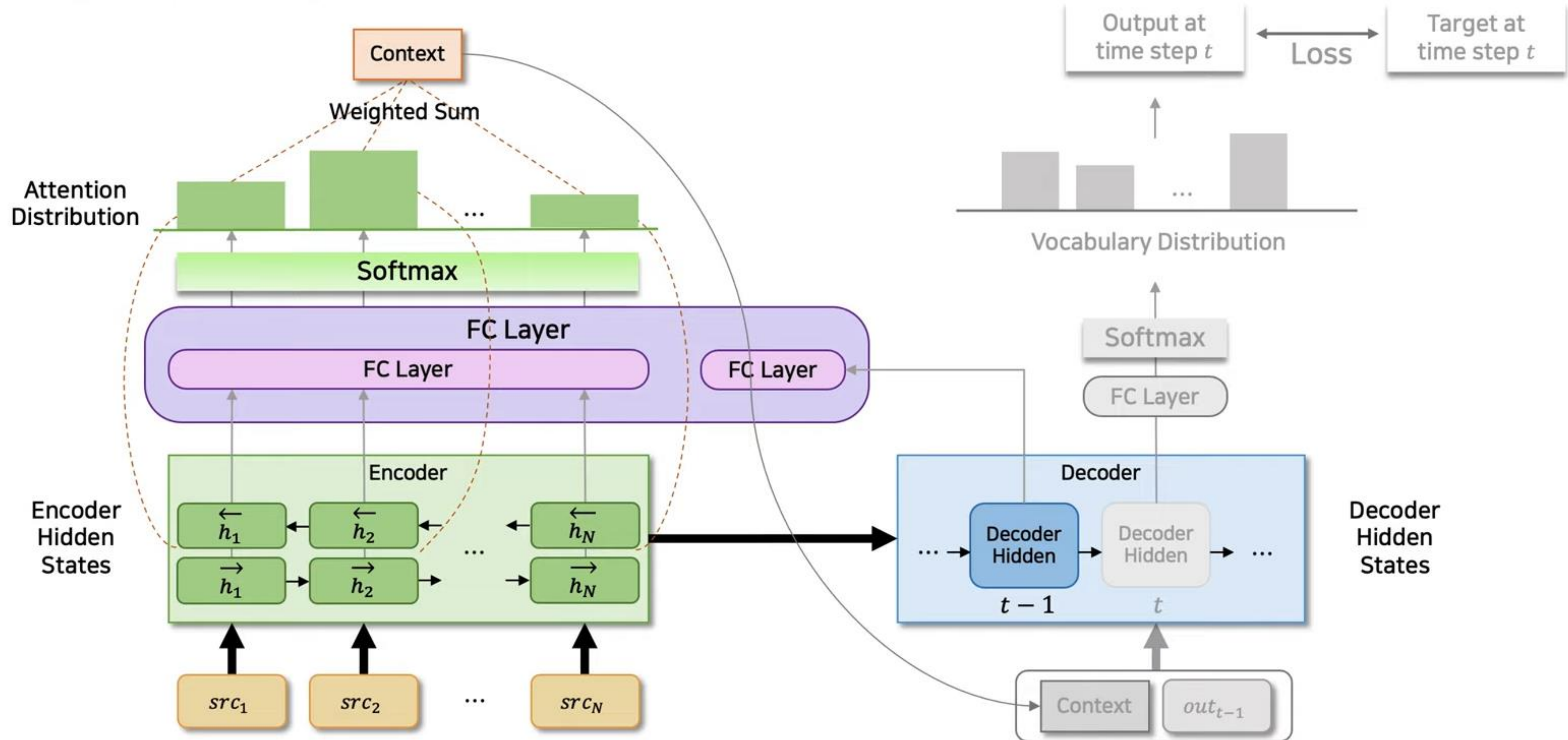$$c_j = \text{ATT}(C, s_j') = \sum_{i=1}^{I} \alpha_{ij} h_i,$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{I} \exp(e_{kj})}$$

$$e_{ij} = v_a^{\mathrm{T}} \tanh(U_a s_j' + W_a h_i)$$

# Attention – Step 5

# Attention



$$p(y_j \mid s_j, y_{j-1}, c_j) = \mathrm{softmax}(W_o t_j)$$
$$t_j = \tanh(W_{t1} s_j + W_{t2} \mathrm{E}[y_{j-1}] + W_{t3} c_j)$$

74

# RNNsearch (Bahdanau et al., 2015)



Image taken from : https://www.youtube.com/watch?v=S2msiG9g7Us

# Results

Neural Machine Translation by Jointly Learning to Align and Translate (Bahdanau et al., 2015)

# Summary

- Limitation of Seq2Seq model

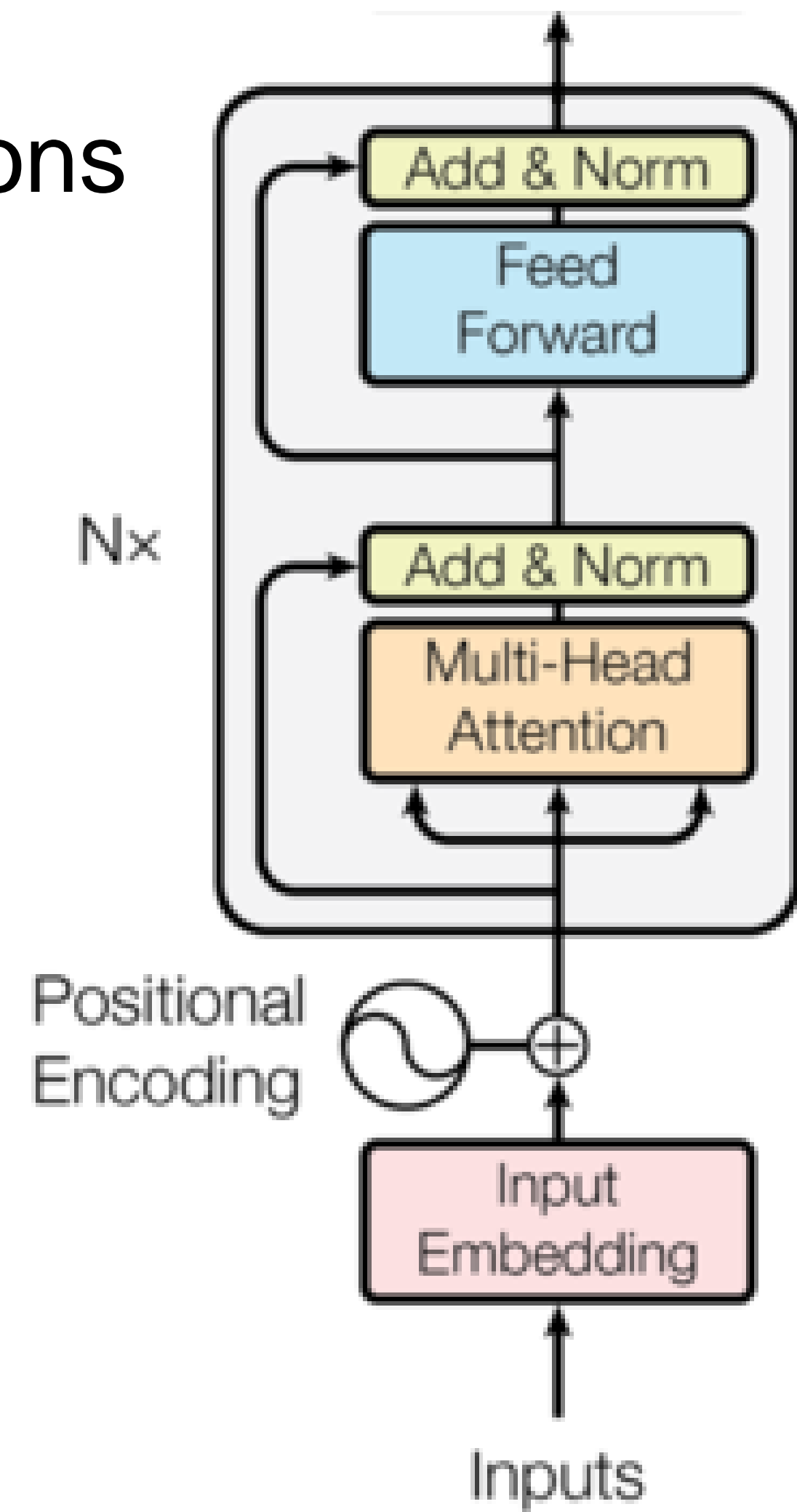- Attention Mechanism

# Transformer

# Transformer

- Attention is all you need (Vaswani et al., 2017)

- An encoder-decoder framework for

  sequence-to-sequence modeling

- No recurrent units



From "Attention is all you need" paper by Vaswani, et al., 2017

# Transformer Encoder

- N layers of Transformer blocks with residual connections

- Parameters in each layer are not shared

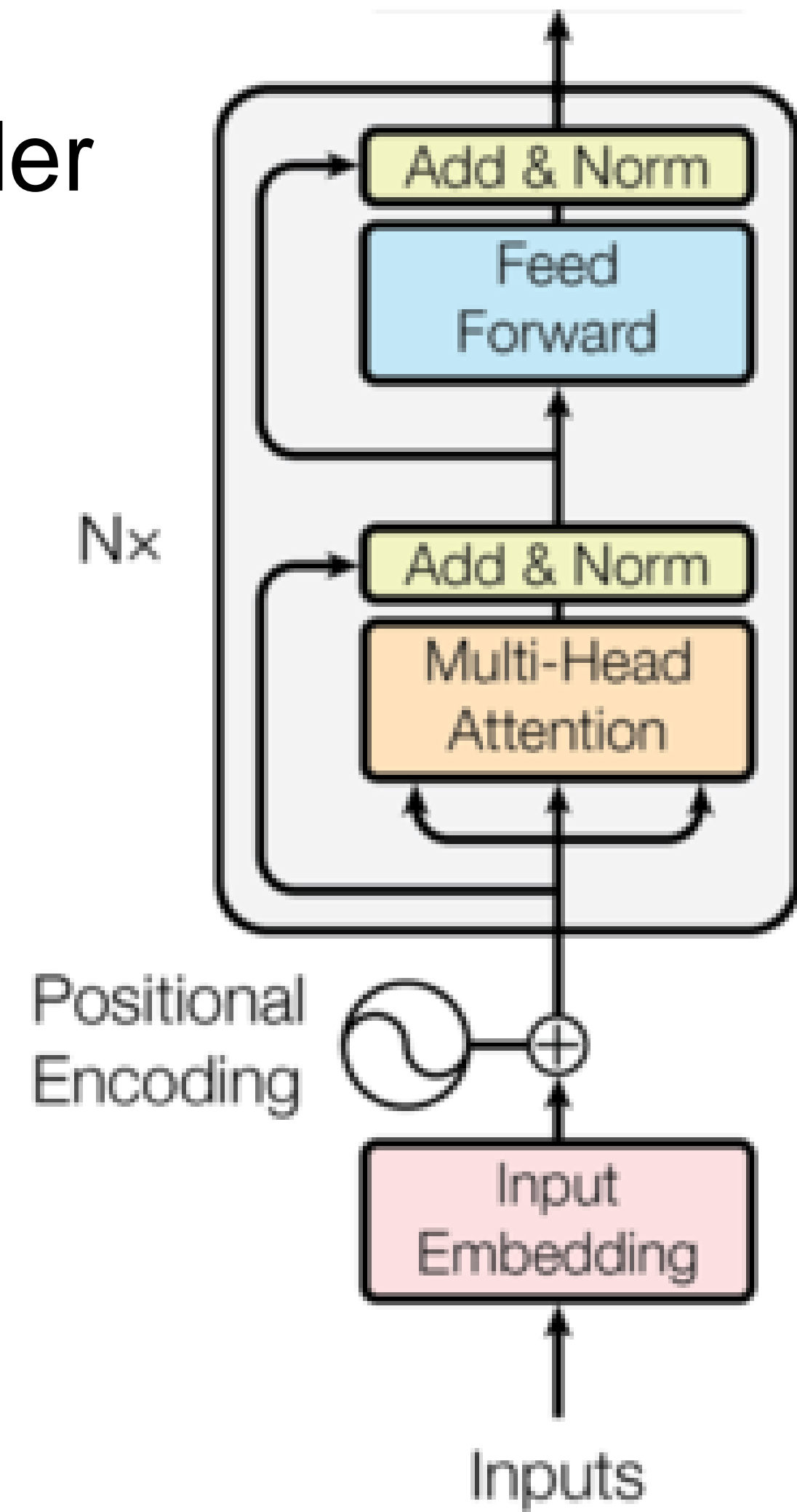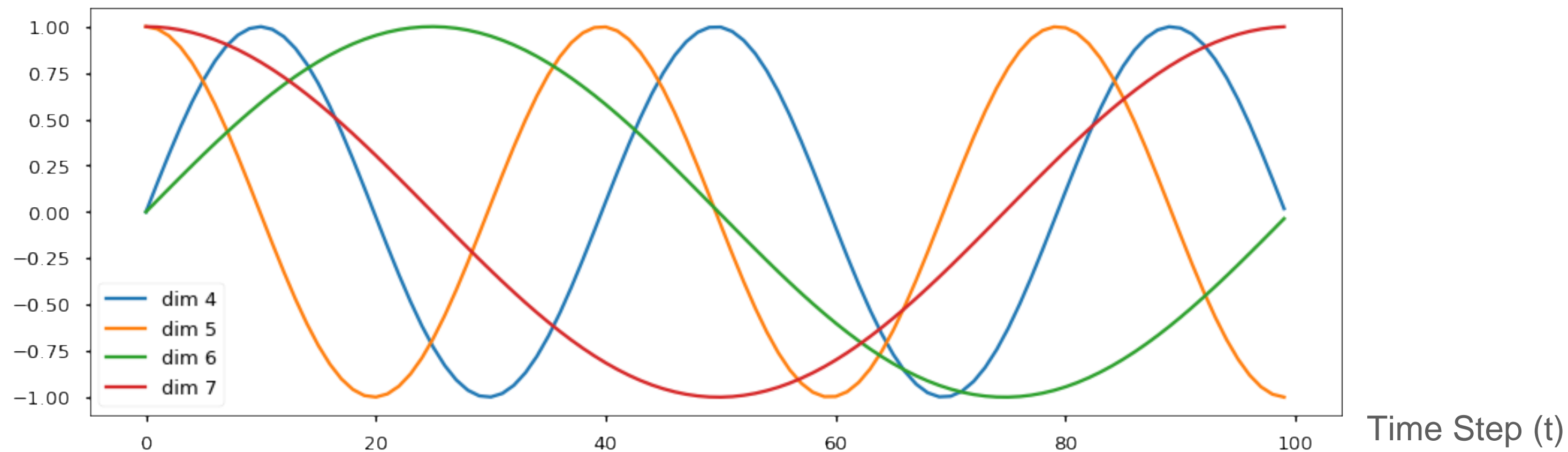- Input tokens are processed in parallel
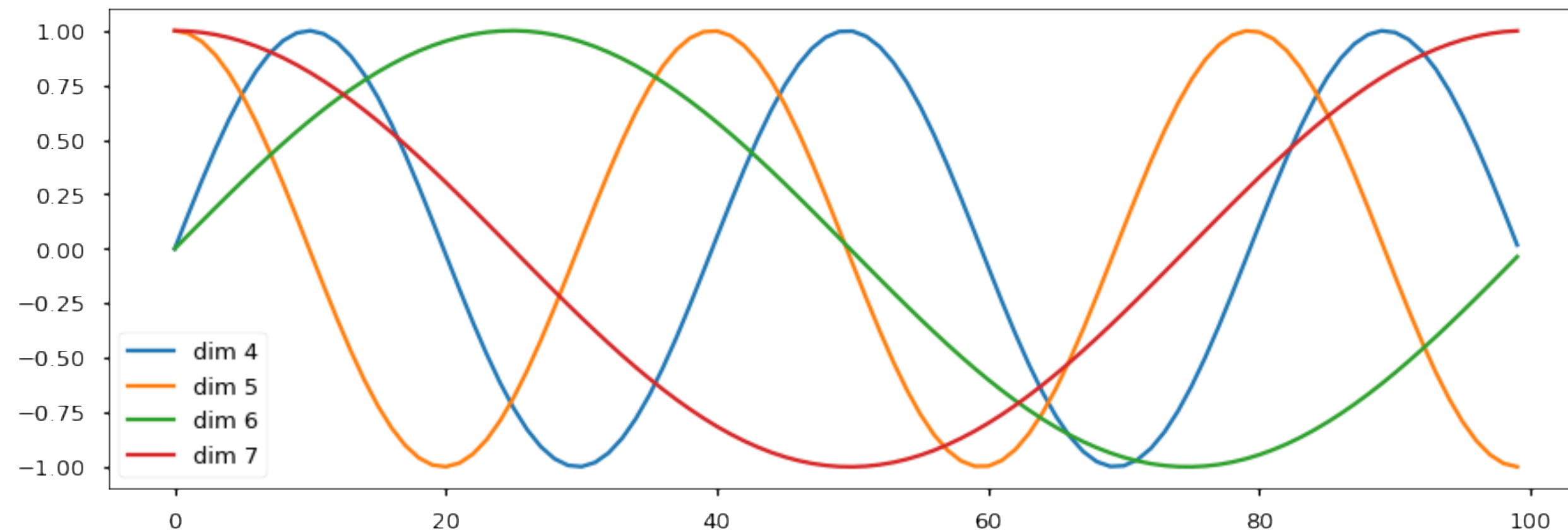
# Positional Encoding

- Positional encoding provide the model about word order

- Sinusoidal position encoding

$$\mathbf{x}'_t = \mathbf{W}_{\text{emb}}\left(\mathbf{x}_t\right) + \vec{p}_t$$



Time Step (t)

# Positional Encoding

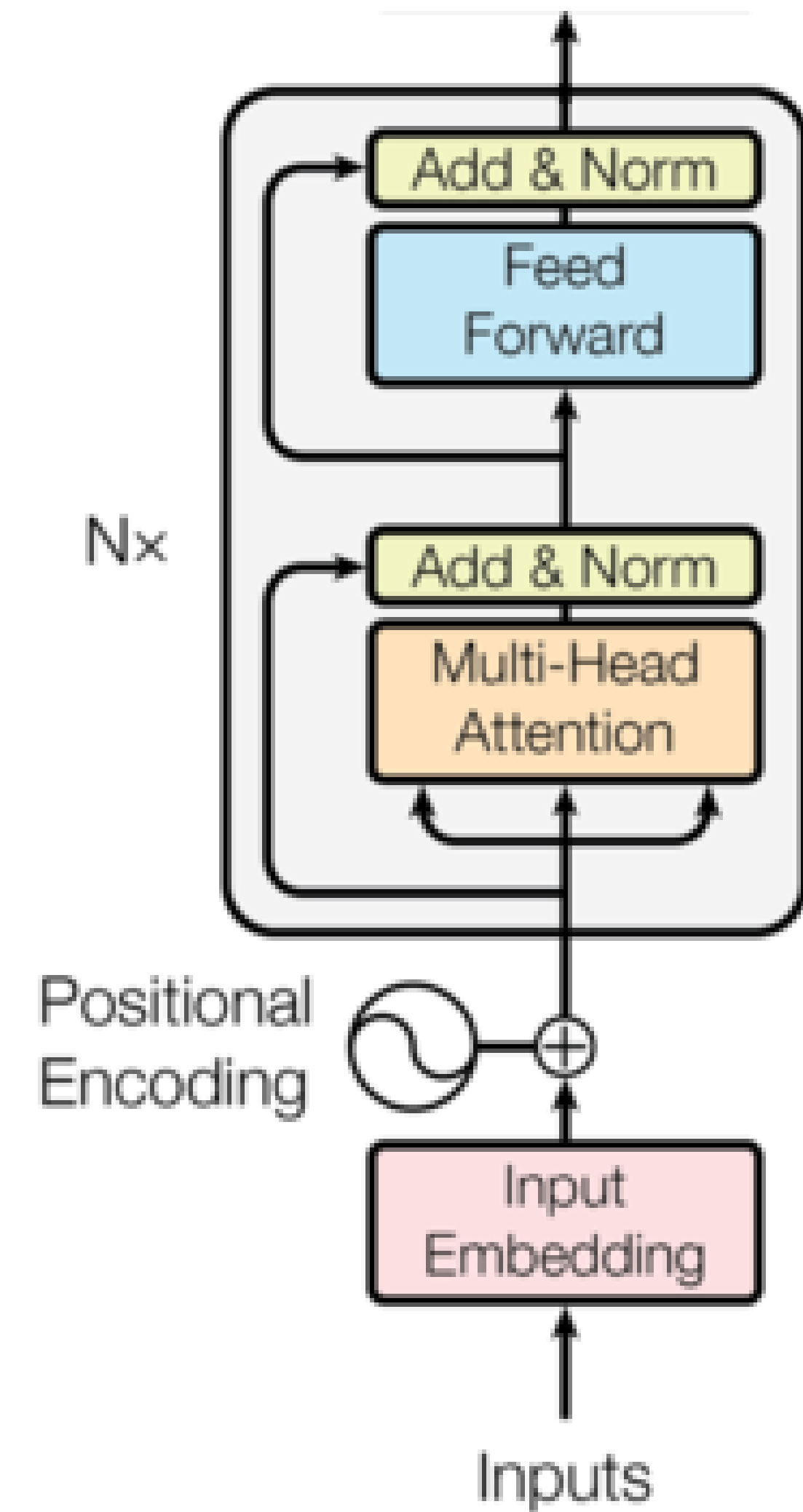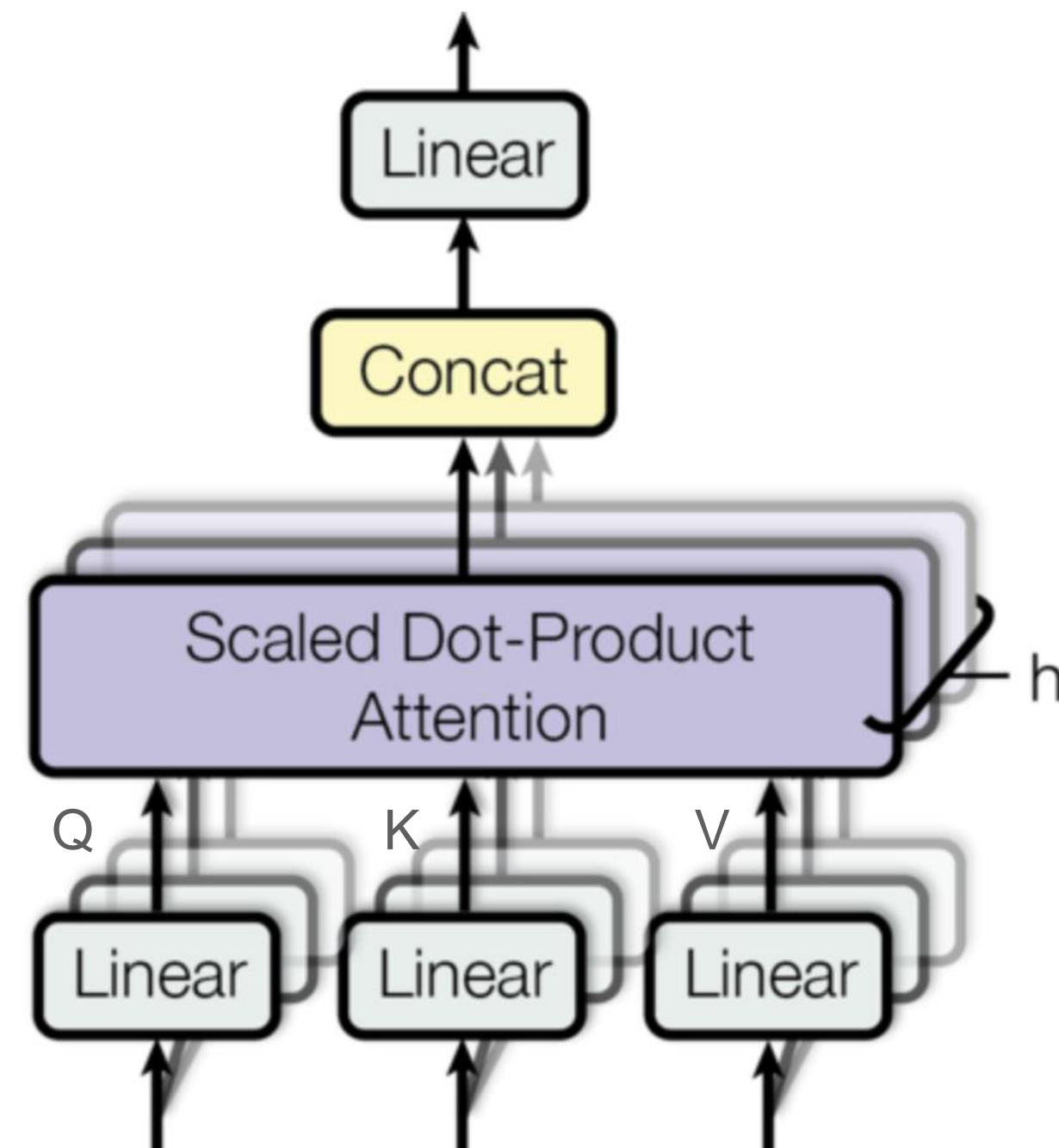$$\mathbf{x}'_t = \mathbf{W}_{\text{emb}}\left(\mathbf{x}_t\right) + \vec{p}_t$$



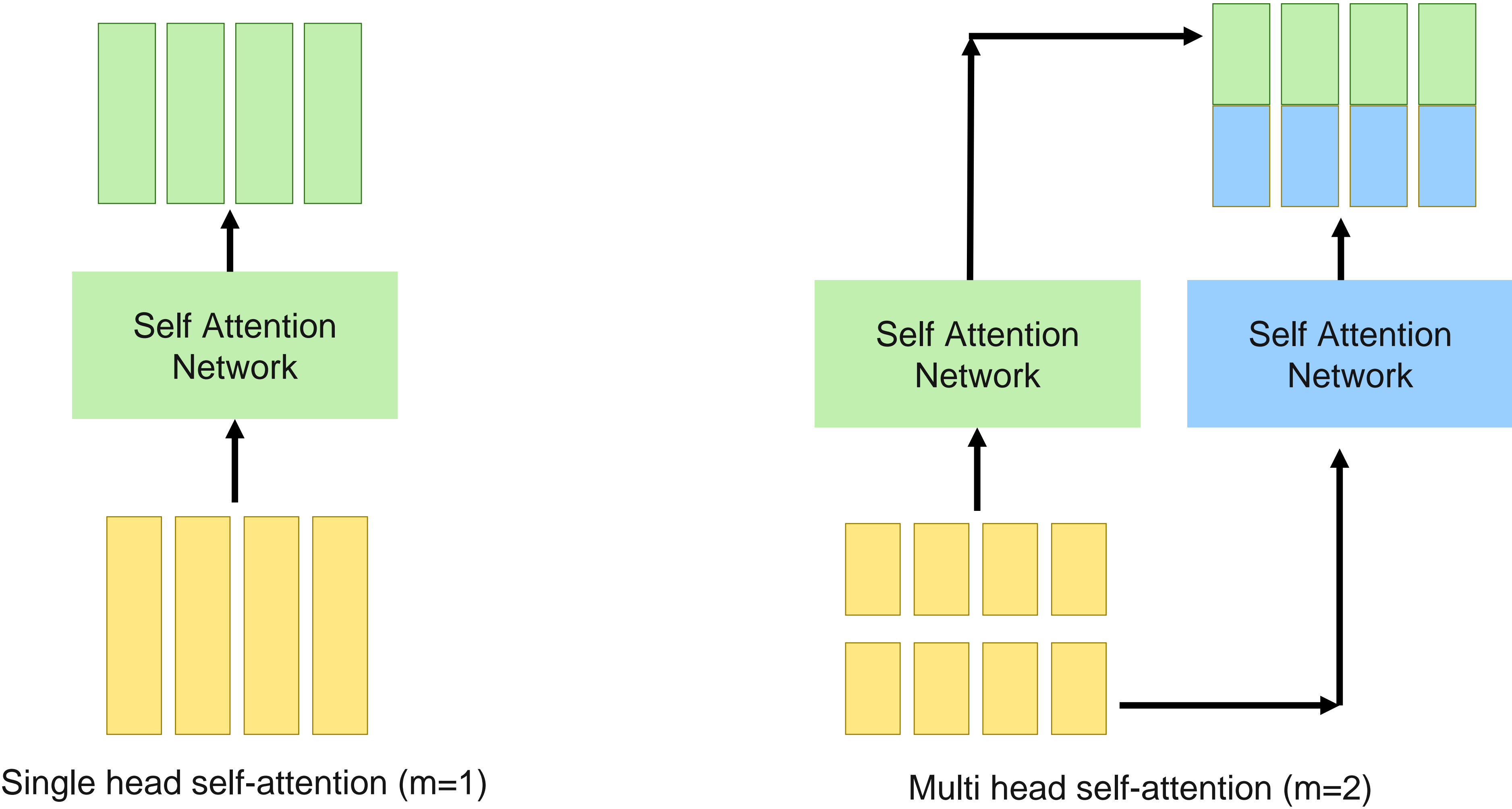Time Step (t)

$$\omega_k = \frac{1}{10000^{2k/d}}$$

$$\overrightarrow{p_t} = \begin{bmatrix} \sin(\omega_1.t) \\ \cos(\omega_1.t) \\ \\ \sin(\omega_2.t) \\ \cos(\omega_2.t) \\ \\ \vdots \\ \\ \sin(\omega_{d/2}.t) \\ \cos(\omega_{d/2}.t) \end{bmatrix}_{d \times 1}$$

82

# Multi-head Attention

- Multi-head Attention is the concatenation of

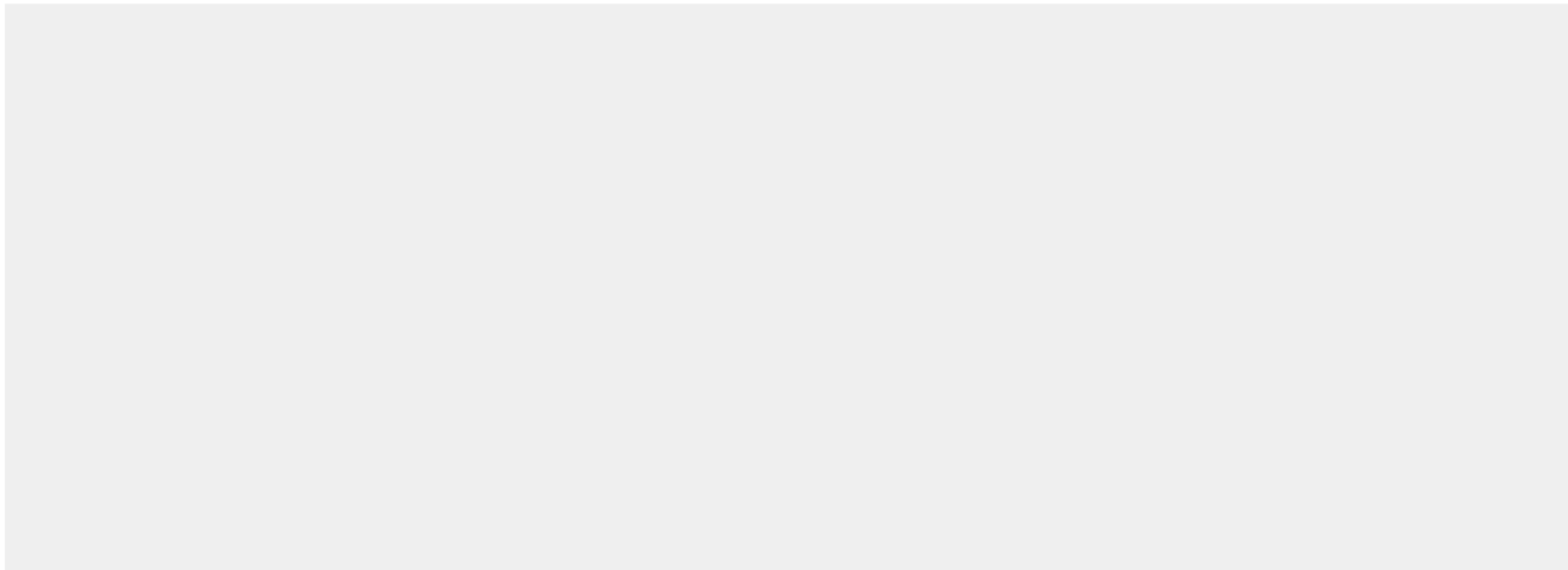the outputs from self-attention network (SAN)

# Multi-head Attention



Single head self-attention (m=1)

Multi head self-attention (m=2)

84

# Self-Attention Network

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

Self-attention

input #1

| 1 | 0 | 1 | 0 |

input #2

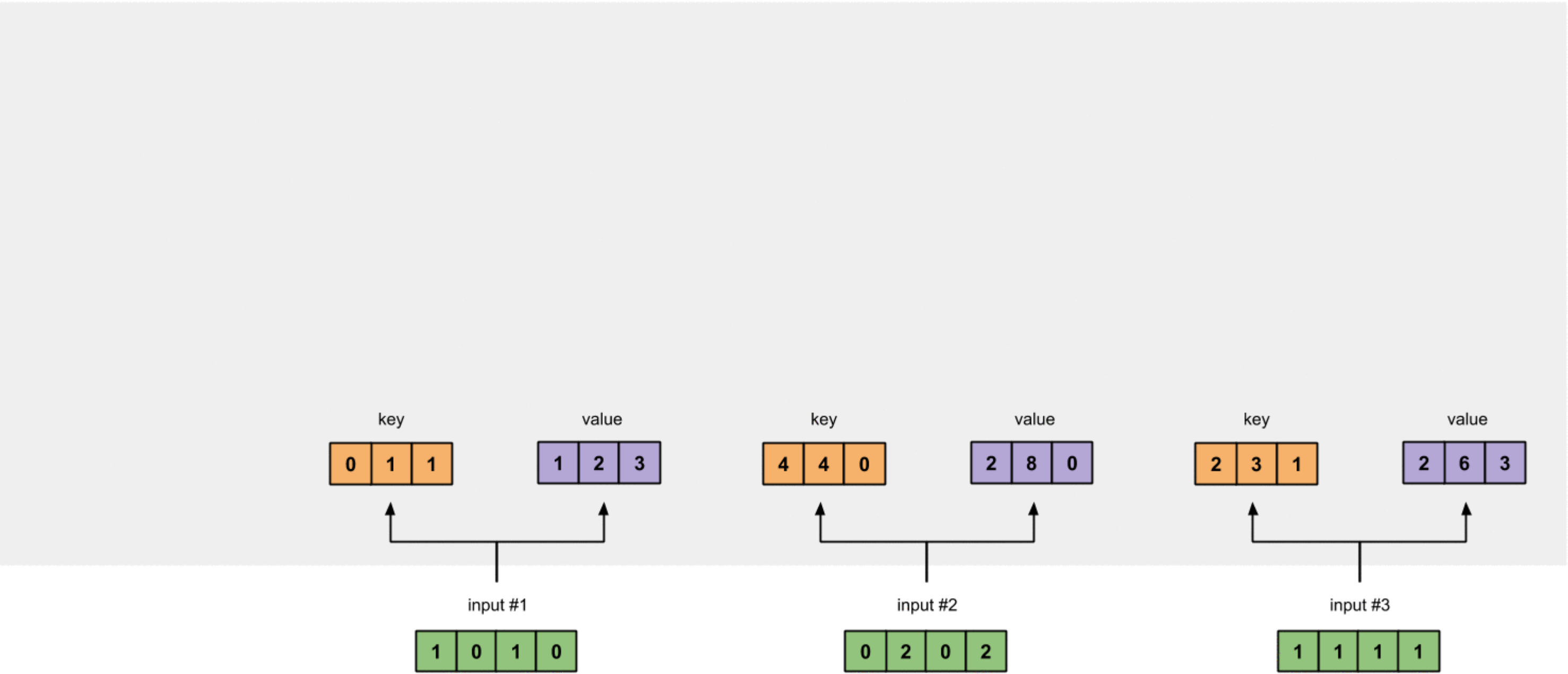| 0 | 2 | 0 | 2 |

input #3

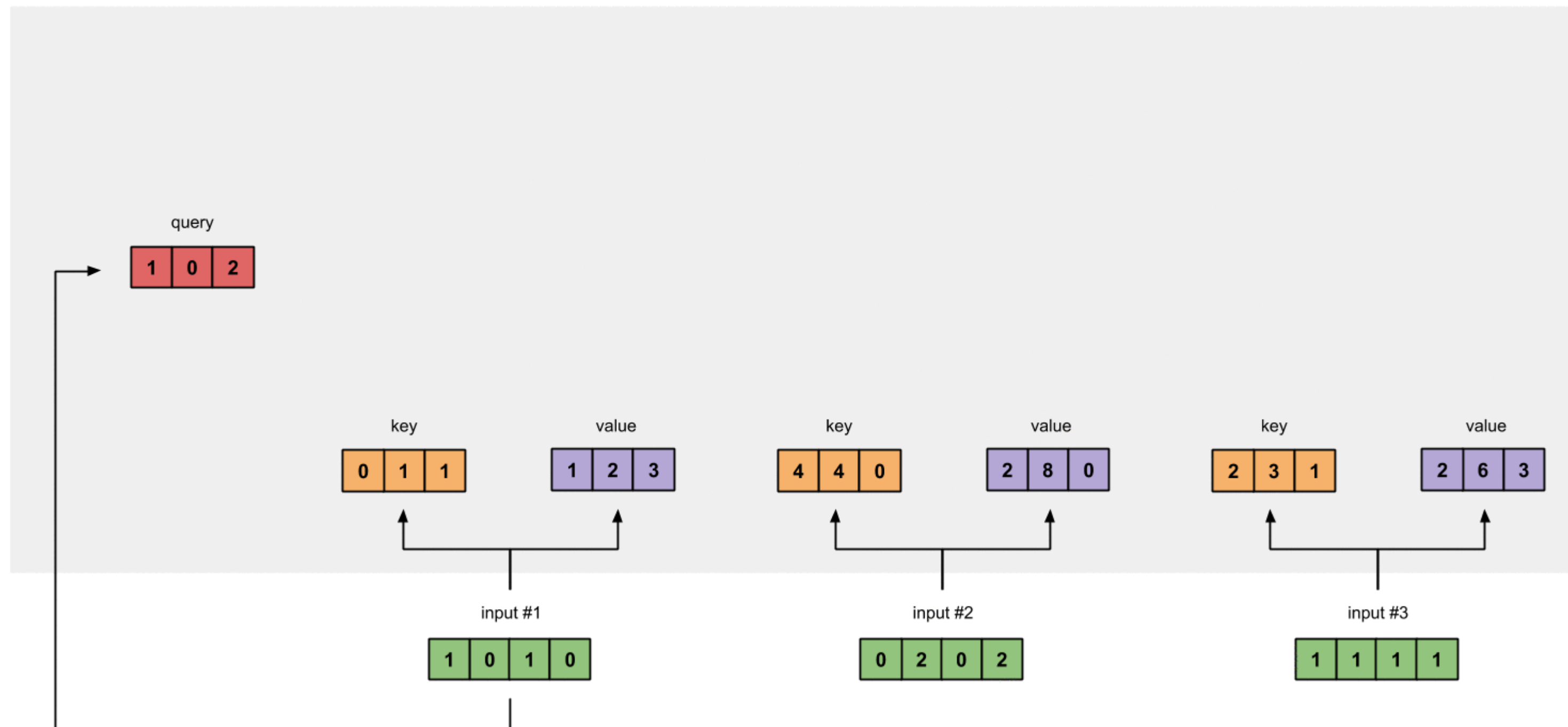| 1 | 1 | 1 | 1 |

# Self-Attention Network

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

Self-attention

# Self-Attention Network

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

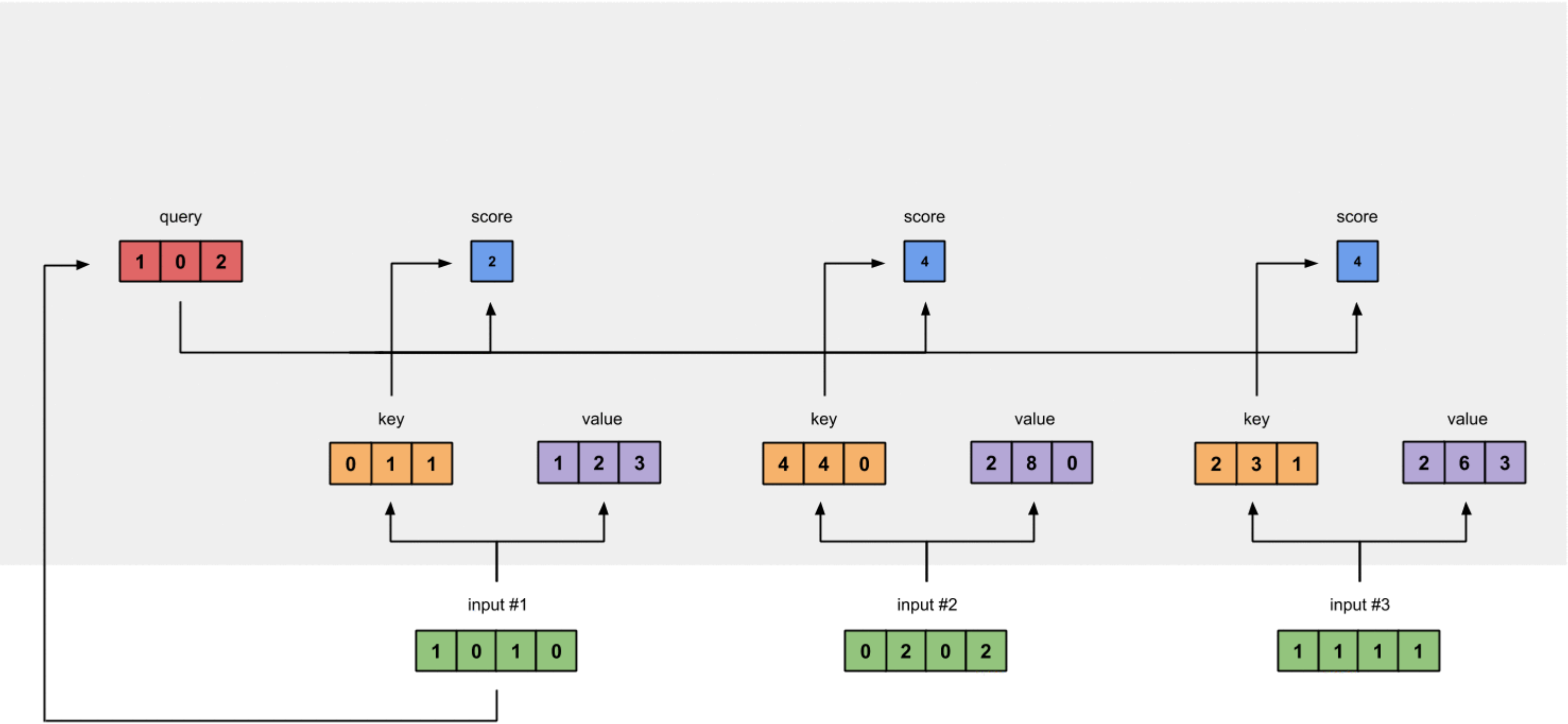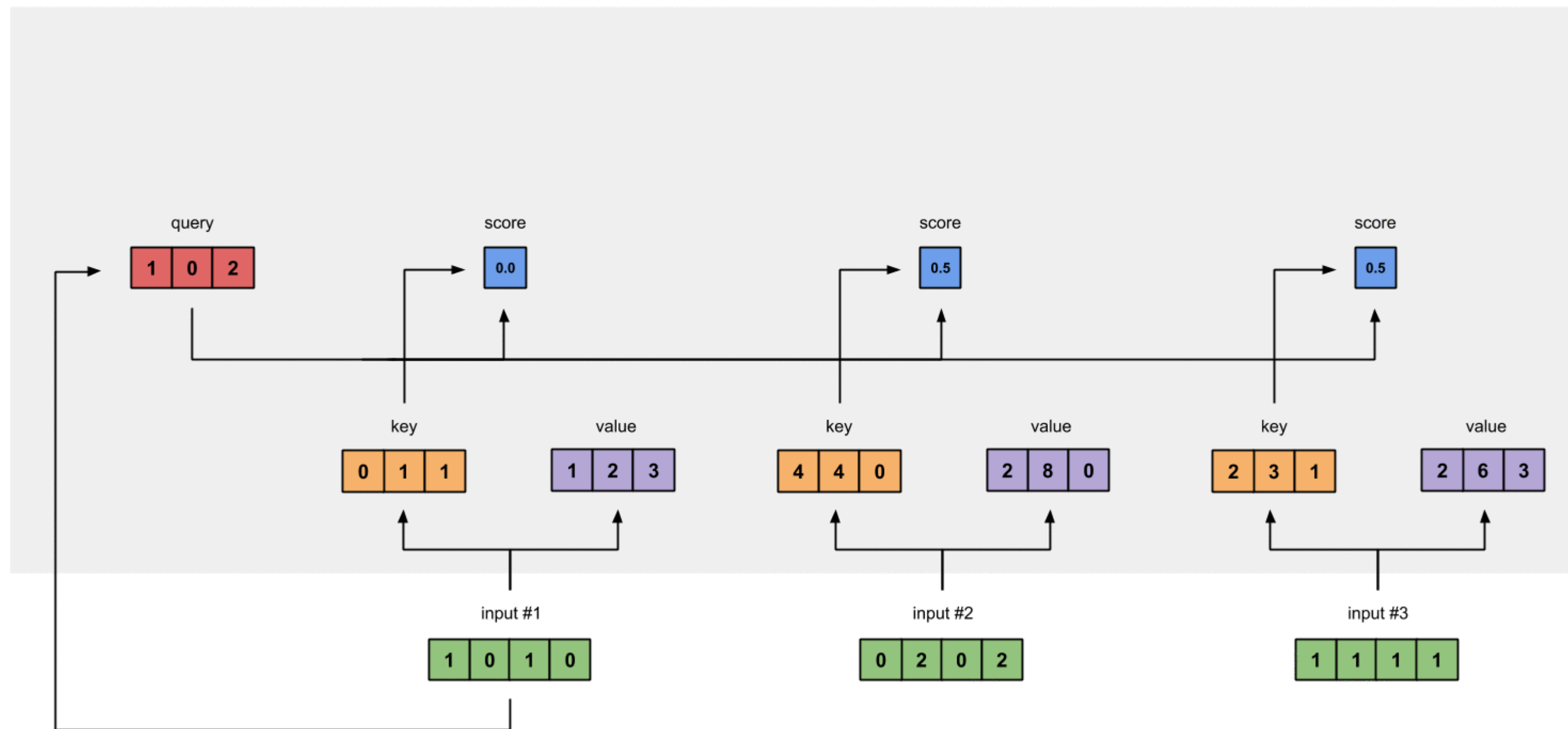Self-attention

query
| 1 | 0 | 2 |

| key | | | value | | | key | | | value | | | key | | | value | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 0 | 1 | 1 | 1 | 2 | 3 | 4 | 4 | 0 | 2 | 8 | 0 | 2 | 3 | 1 | 2 | 6 | 3 |

input #1
| 1 | 0 | 1 | 0 |

input #2
| 0 | 2 | 0 | 2 |

input #3
| 1 | 1 | 1 | 1 |

https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a

# Self-Attention Network

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



Self-attention

https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a

# Self-Attention Network

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



Self-attention

https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a
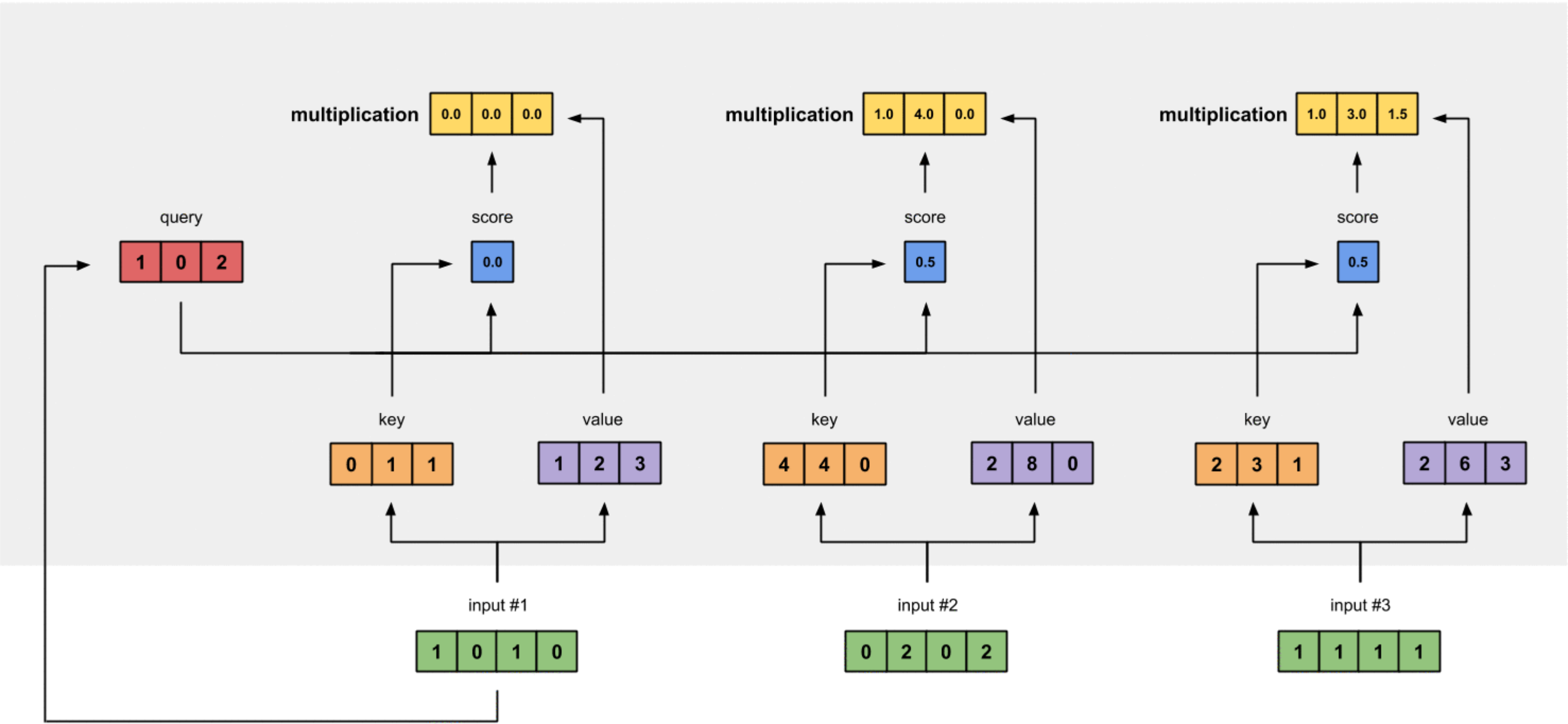
# Self-Attention Network

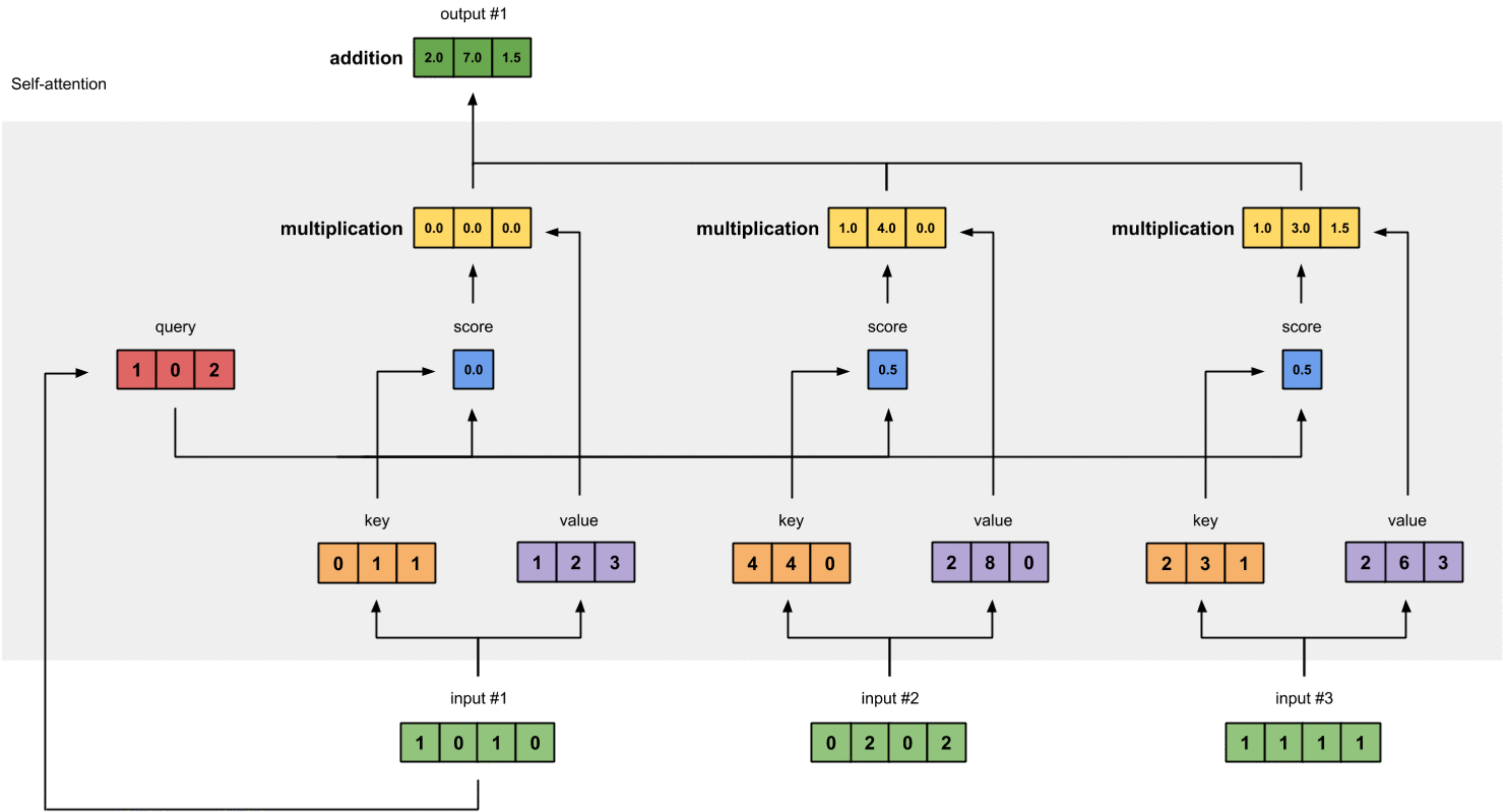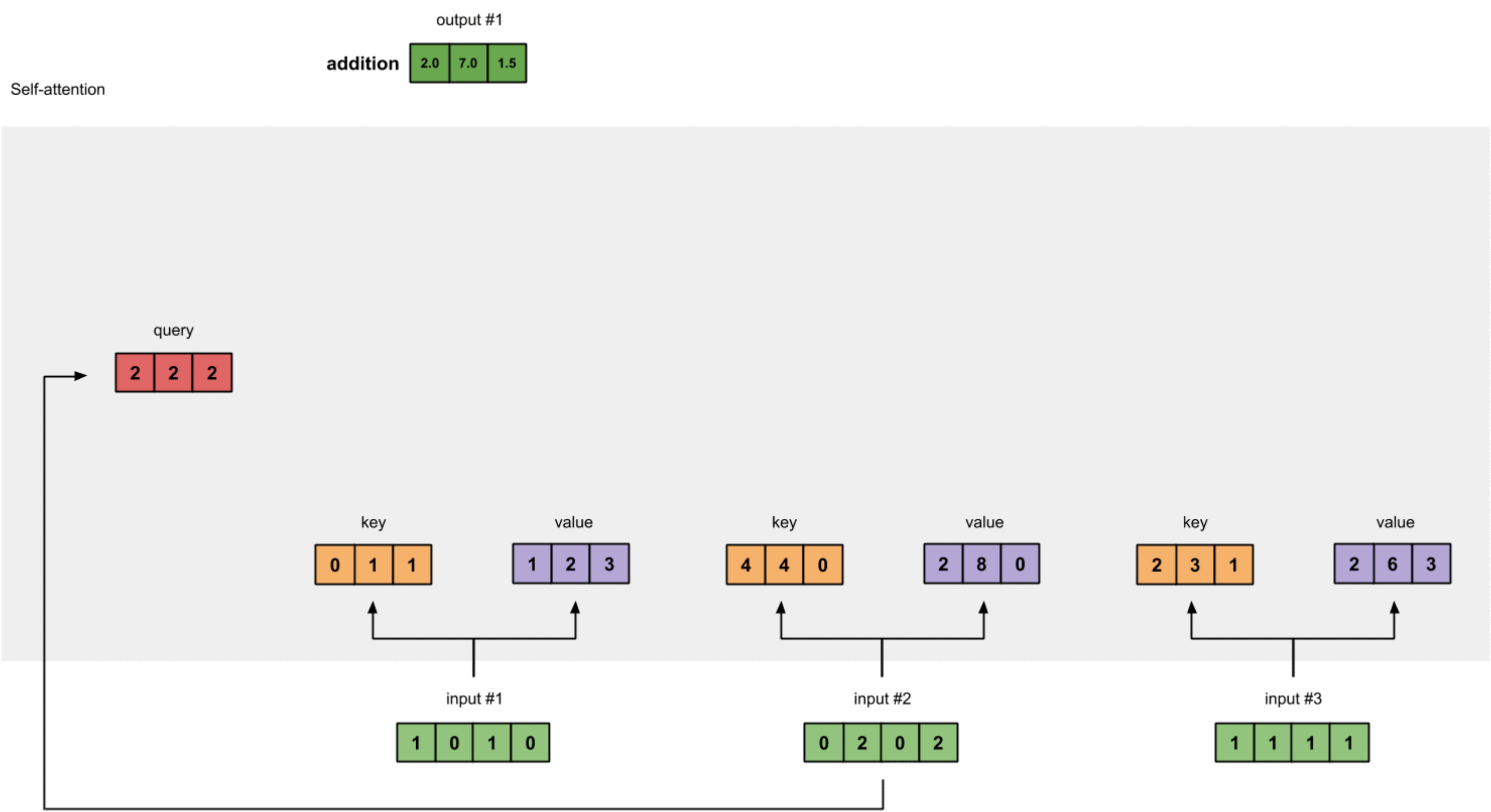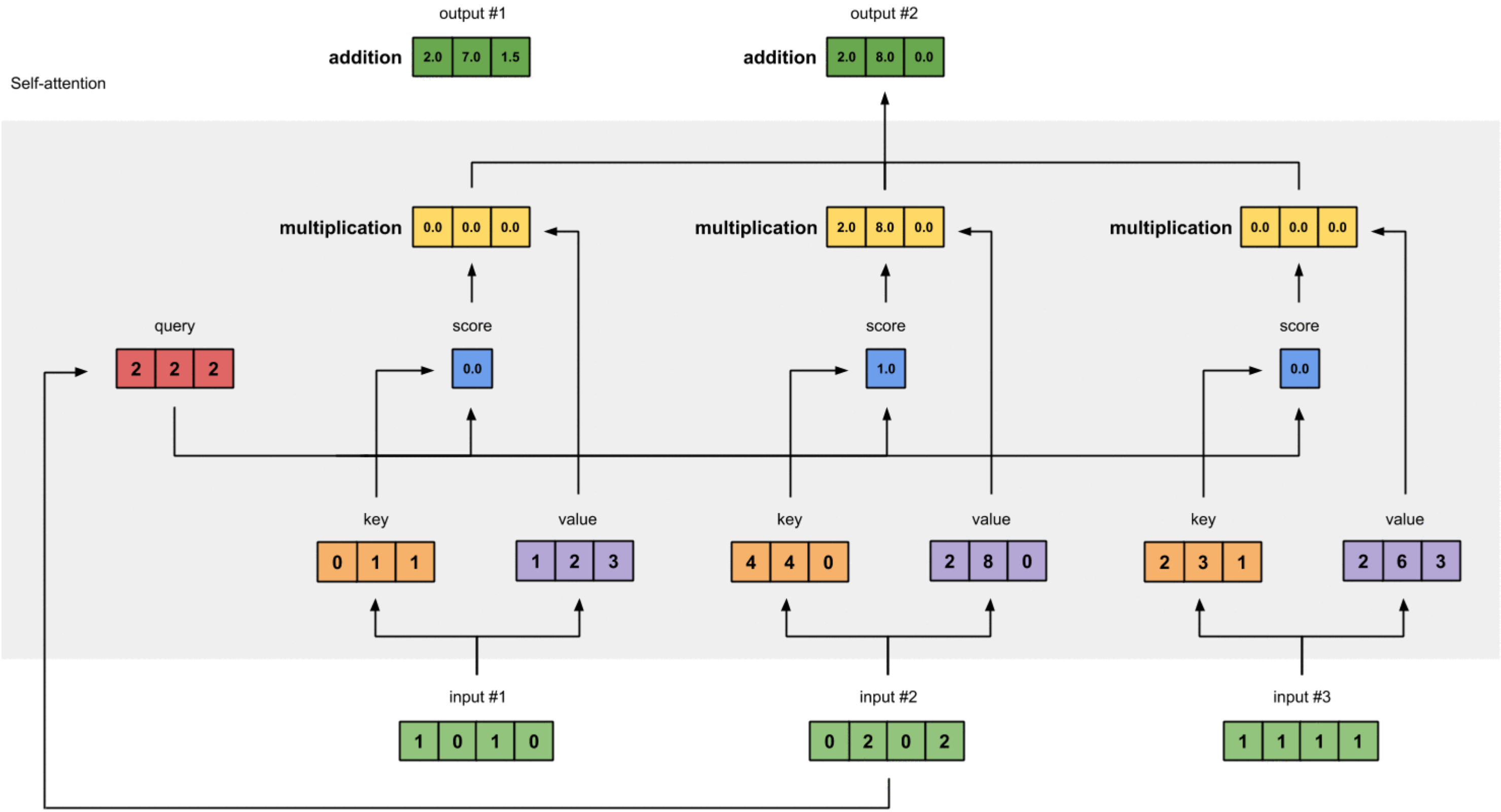$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$



Self-attention

https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a

# Self-Attention Network



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

https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a

# Self-Attention Network

https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a

# Self-Attention Network

https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a

# Self-Attention Network

https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a

# Multi head Self-Attention



Concatenation

# Advantages of self-attention network

- Process in parallel

- Better in modeling long term dependencies, able to freely attend to other position.
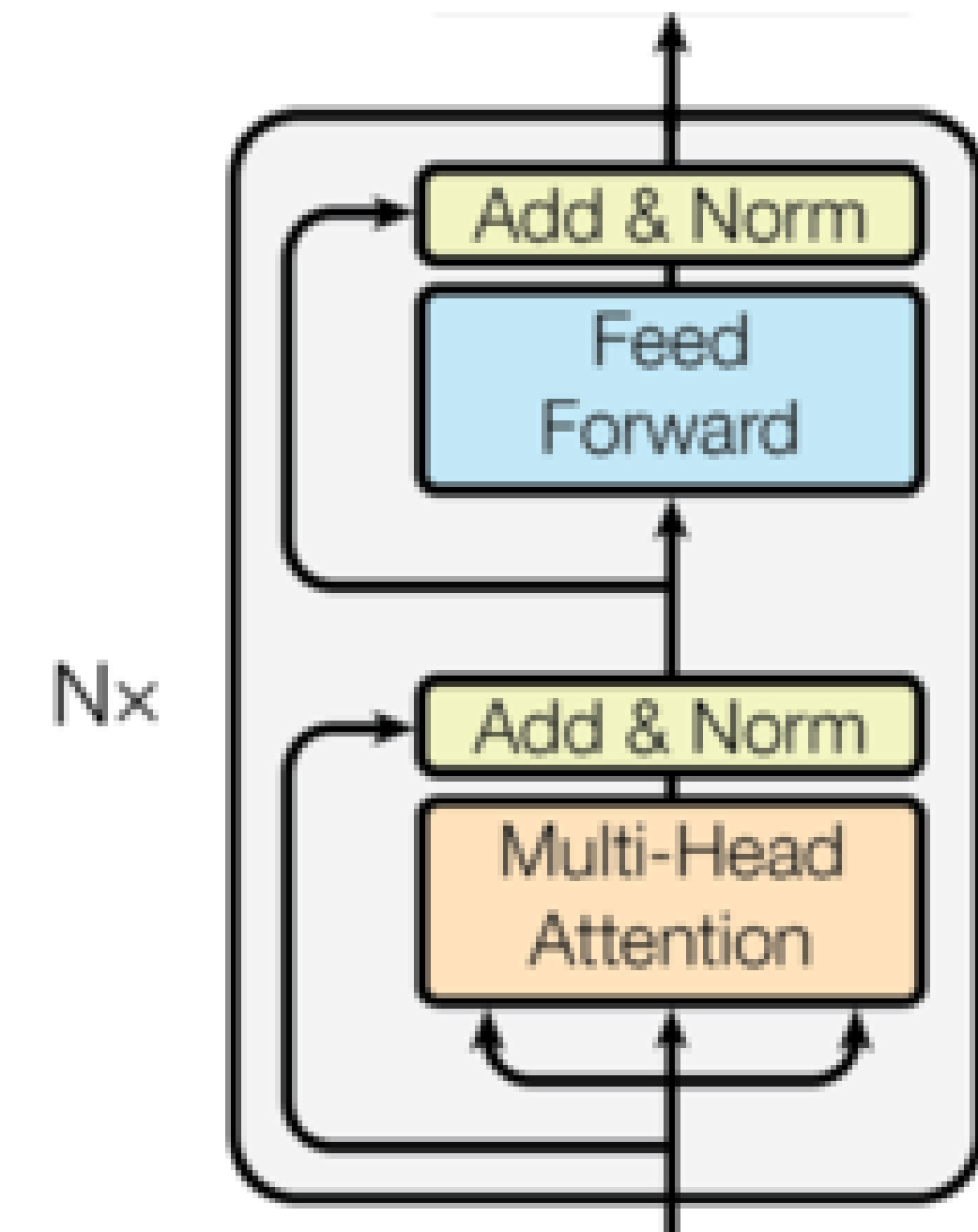
# Encoder : Add & Norm

- Residual Connection

- Layer Normalization

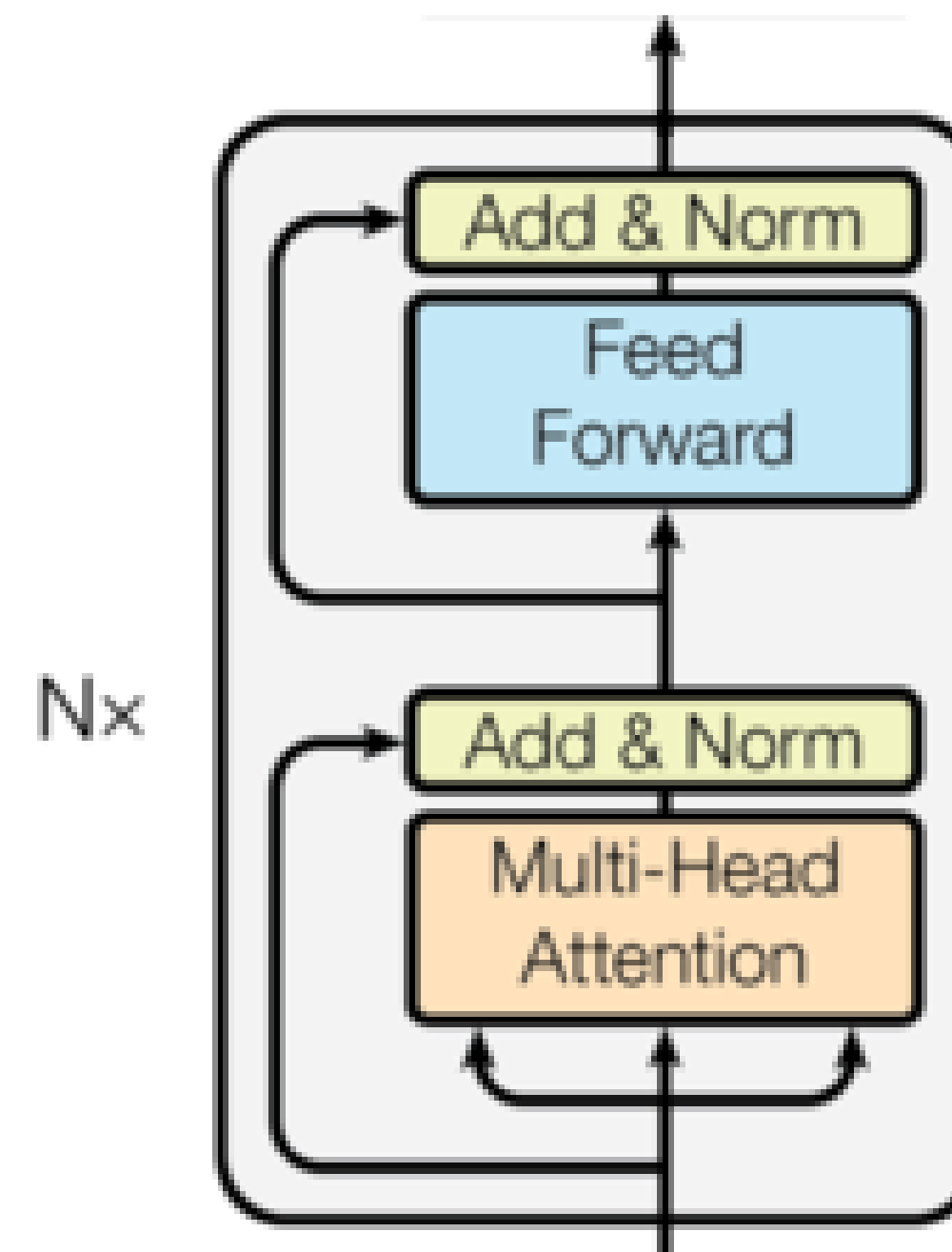$$x' = \mathrm{LayerNorm}\left(\mathrm{SelfATT}(x) + x\right)$$

Layer Normalization – to normalize mean and variance of inputs ( $a^l$ ) for a specific layer ($l$), assume that the layer has $H$ hidden units.

$$\mu^l = \frac{1}{H}\sum_{i=1}^{H} a_i^l \qquad \sigma^l = \sqrt{\frac{1}{H}\sum_{i=1}^{H}\left(a_i^l - \mu^l\right)^2} \qquad \bar{a}_i^l = \frac{g_i^l}{\sigma_i^l}\left(a_i^l - \mu_i^l\right)$$

Nx

Add & Norm

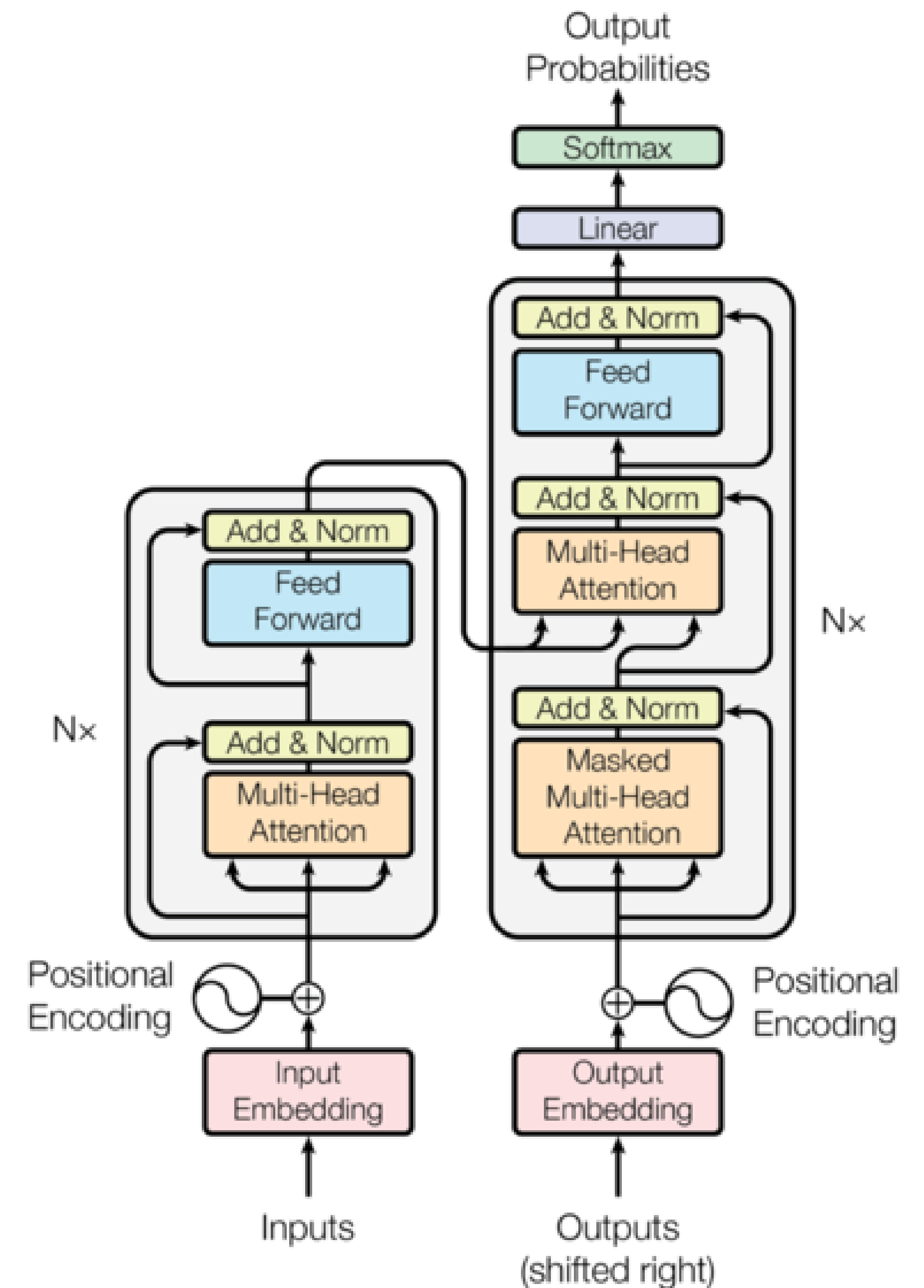Feed Forward

Add & Norm

Multi-Head Attention

97

# Encoder : Feed Forward

- Feed forward network

- Followed by residual connection and layer normalization
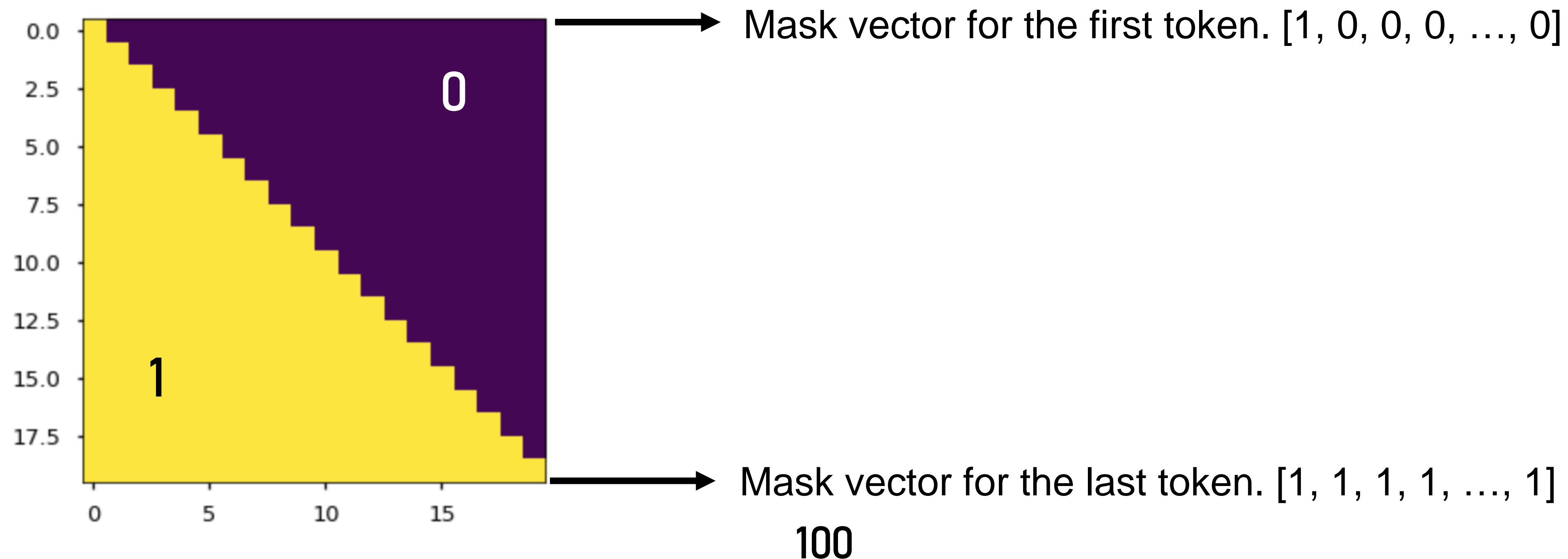
# Transformer : Decoder

- Auto Regressive Decoding

- Based on Transformer layers
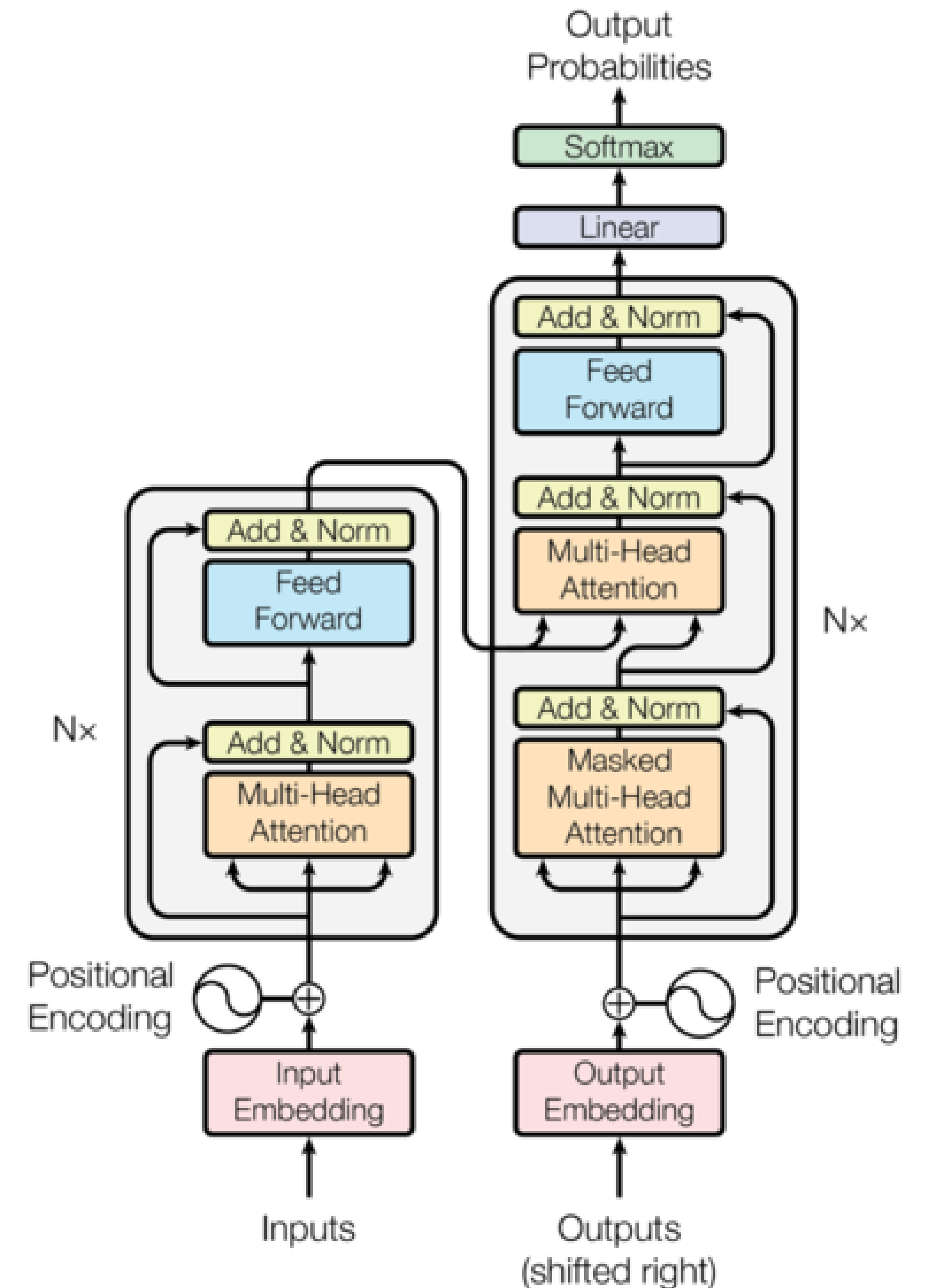
# Decoder : Masked Multi-head Attention

- Auto Regressive Decoding

  Current token can attend only left-side tokens because in decoding step the right-side tokens are not generated. Attention weights are multiplied by this mask matrix.



Mask vector for the first token. [1, 0, 0, 0, …, 0]

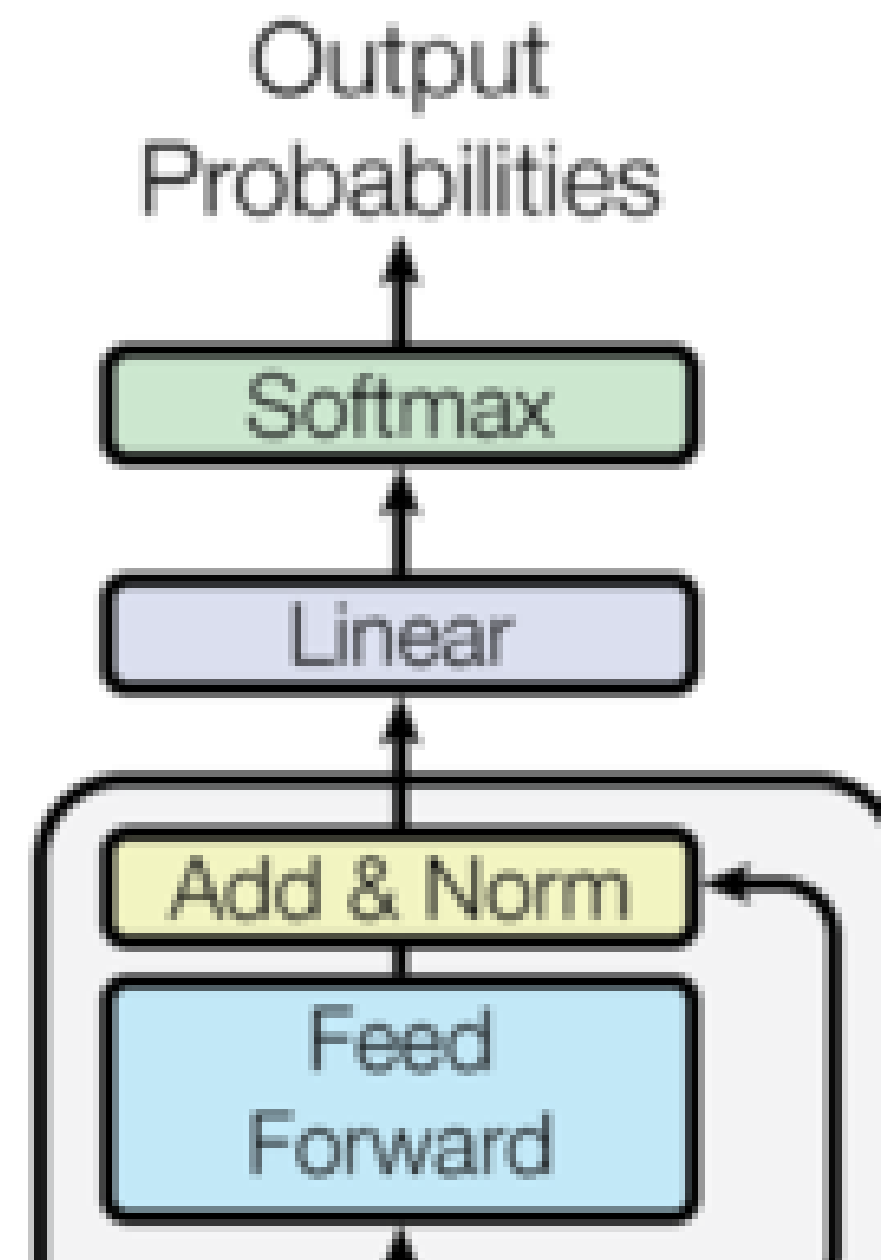Mask vector for the last token. [1, 1, 1, 1, …, 1]

# Decoder-Encoder Attention

- Query data from the encoder outputs.

- Encoder output (K,V) , Decoder State (Q)

# Decoder : Feed Forward and Softmax

- Predict target word distribution

# Optimization

- Label Smoothed Regularization

$$\bar{y}_j^t = (1 - \epsilon)\bar{y}_j + \frac{\epsilon}{V}$$

For example, $V = 3$, $\epsilon = 0.3$

$$\bar{y}_{\text{true(smooth)}} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \\ \epsilon \end{bmatrix} + \frac{1}{3}\begin{bmatrix} \epsilon \\ \epsilon \\ \epsilon \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.8 \end{bmatrix}$$
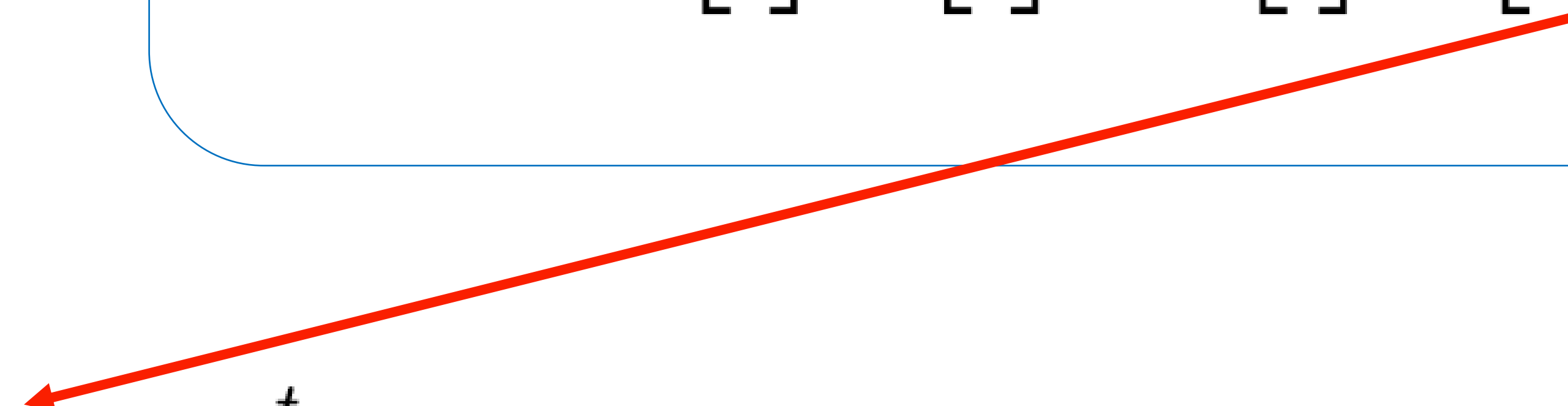
103

# Optimization

- Label Smoothed Regularization

$$\bar{y}_j^t = (1 - \epsilon)\bar{y}_j + \frac{\epsilon}{V}$$

For example, $V = 3$, $\epsilon = 0.3$

$$\bar{y}_{\text{true(smooth)}} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \\ \epsilon \end{bmatrix} + \frac{1}{3}\begin{bmatrix} \epsilon \\ \epsilon \\ \epsilon \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.8 \end{bmatrix}$$

- Label Smoothed NLL

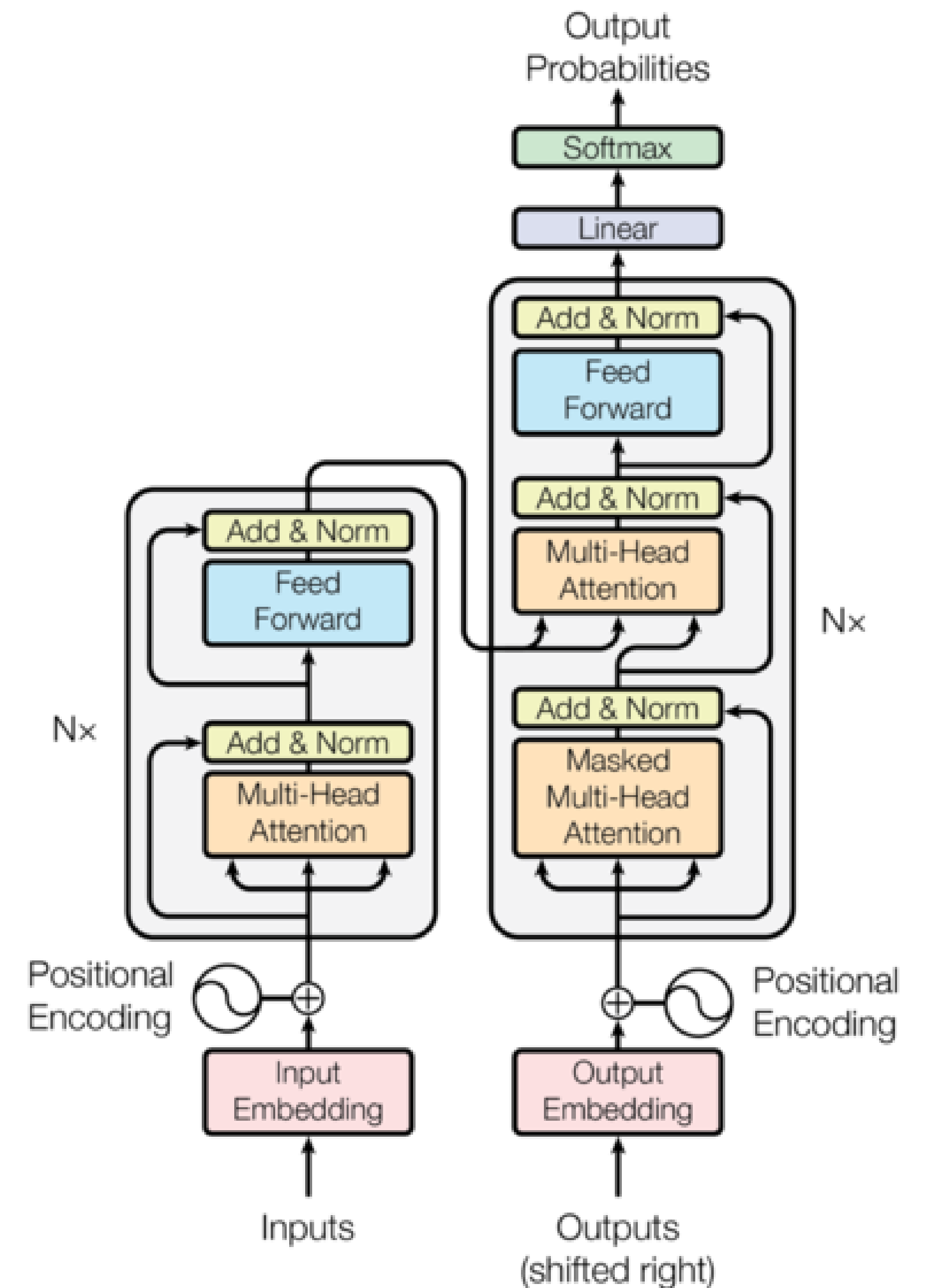$$= -\sum_{t=1}^{T}\sum_{j=1}^{V} \bar{y}_j^t \log \hat{y}_j^t$$

# Decoding : Auto Regressive

- Greedy search

- Beam Search

https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html

# Summary

- Self-Attention Network

- Multi-Head Self Attention

- Positional Encoding

- Encoder

- Decoder

- Optimization

# Data Preparation for MT

# Data Preparation for MT

- **Data Collection**

- **Data Cleansing and Tokenization**

- **Split Train / Validation / Test**

- **Additional Preprocessing Steps for NMT**

  - Subword Preparation

  - Padding and Binarizing

# Data Collection

- **Data Collection**

# Data Cleansing and Tokenization

- Clean empty line

- Align parallel sentences

- Deduplicate

- Tokenize (Word, character segmentation)

- Filter low quality pairs

  - Alignment score

  - Length ratio ( # of source tokens / # of target tokens )

# Split Train / Validation / Test

- Shuffle

- Train / Validation / Test

Train 90%  Valid 5% Test 5%

# Additional Steps for NMT

- Subword Units (Sentence piece, Byte-pair encoding)

ซีรีย์ จีน เรื่อง ดาบมังกรหยก
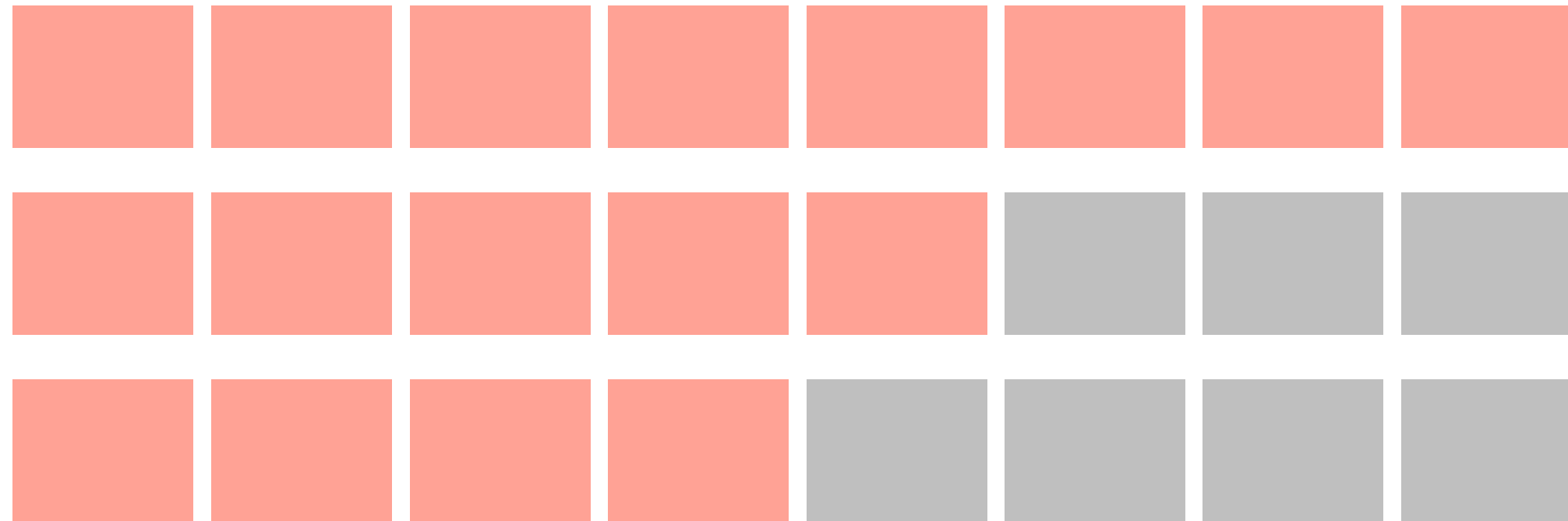
ซีรีย์ จีน เรื่อง ดาบ@@ มังกร@@ หยก

Rico Sennrich, Barry Haddow and Alexandra Birch, Neural Machine Translation of Rare Words with Subword Units, ACL, 2016

# Additional Steps for NMT

- Padding



- Binarizing – Convert strings to tensors

# Summary

- **Data Collection**

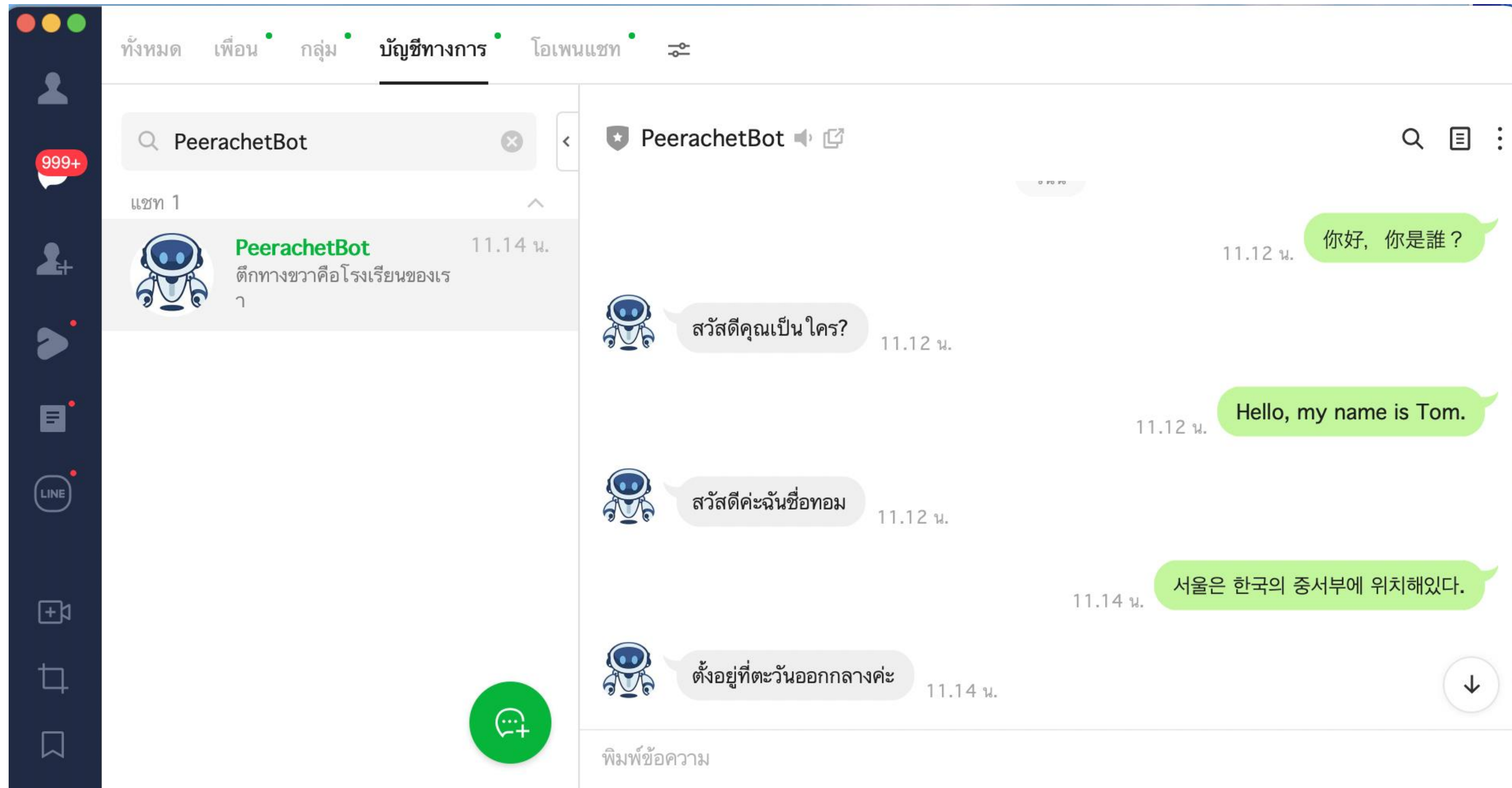- **Data Cleansing and Tokenization**

- **Split Train / Validation / Test**

- **Additional Preprocessing Steps for NMT**

  - Subword Preparation

  - Padding and Binarizing

114

# Application of Seq2Seq Model

# Possible Applications

- **Machine Translation**

- **Text Summarization**

- **Paraphrasing**

- **Question and Answering**

- **Chat Bot**

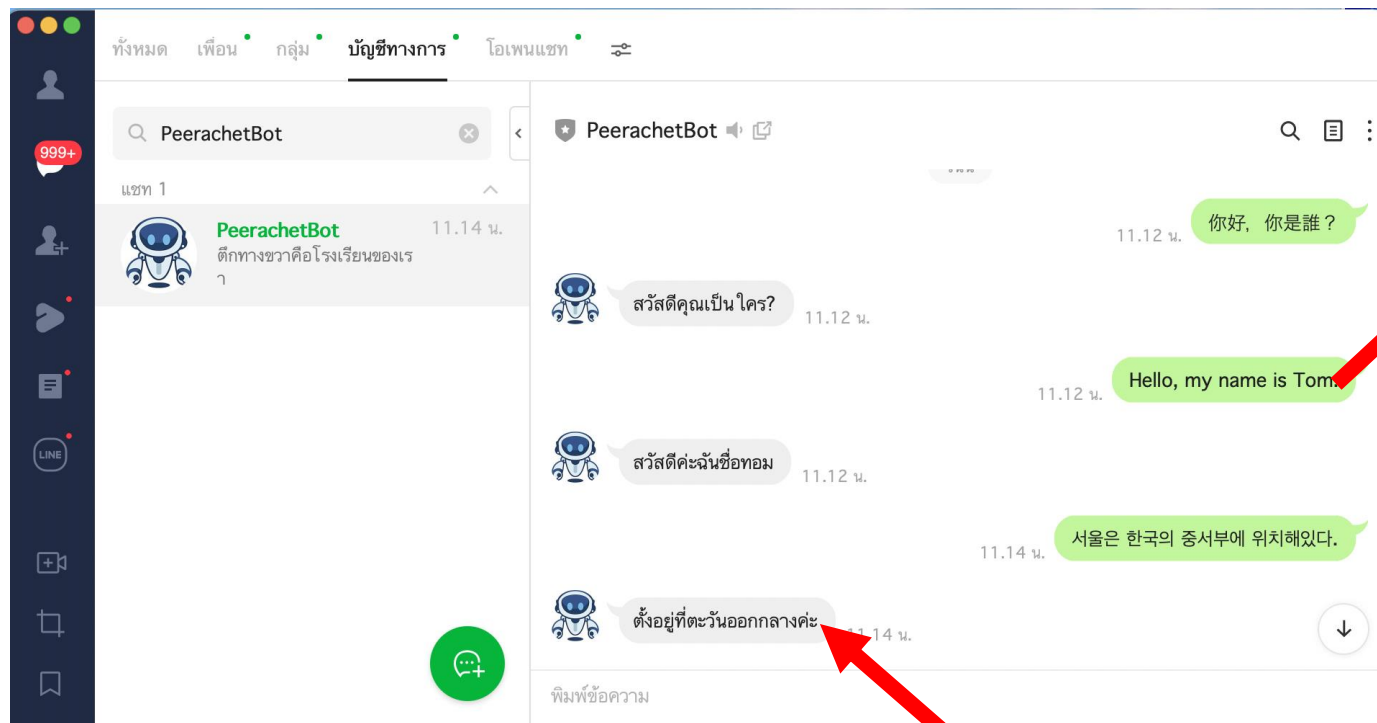# Line Bot Translation Service (Demo only)

# Line Bot Web Hook (Python version)



```python
@handler.add(MessageEvent, message=TextMessage)
def handle_text_message(event):
    text = event.message.text #message from user
    userid = str(event.source)

    print(json.loads(userid)["userId"])

    if detect(text) != "th":
        text = text[0:300]
        print("TRANSLATE ",text)
        now = datetime.now()
        # Format the date and time as a string
        date_time_str = now.strftime("%d/%m/%Y, %H:%M:%S")
        addlog(date_time_str,{"query" : text})



        response = requests.get(f"https://<your translation service url>/api.php?text={text}&lang=th")
        outtext = response.text
        line_bot_api.reply_message(
            event.reply_token,
            TextSendMessage(text=outtext)) #reply the same message from user

    return
```

118

# LINE Developer Console



119

# References

- [StatMT.org](StatMT.org)

- [https://archive.illc.uva.nl/ESSLLI2008/Materials/KoehnCallisonBurch/book.pdf](https://archive.illc.uva.nl/ESSLLI2008/Materials/KoehnCallisonBurch/book.pdf)

- Philipp Koehn, [Statistical Machine Translation](Statistical Machine Translation), 2009

- Christopher Olah, [Understanding LSTM Networks](Understanding LSTM Networks), 2015

- Bahdanau, D., Cho, K. H., & Bengio, Y. [Neural machine translation by jointly learning to align and translate](Neural machine translation by jointly learning to align and translate). ICLR 20151

- Vaswani, Ashish & Shazeer, Noam & Parmar, Niki & Uszkoreit, Jakob & Jones, Llion & Gomez, Aidan & Kaiser, Luksz & Polosukhin, Illia, [Attention is all you need](Attention is all you need) , 2017

- [Transformer: A Novel Neural Network Architecture for Language Understanding](Transformer: A Novel Neural Network Architecture for Language Understanding), 2017

# References

- Jay Alammar, [The Illustrated Transformer](#) , 2018

- Rico Sennrich, Barry Haddow and Alexandra Birch, [Neural Machine Translation of Rare Words with Subword Units](#), ACL, 2016

- [บทที่1 ทำ LINE Bot สามารถโต้ตอบ หรือ Chatbot ด้วย Python (Official) - Saixiii](#)

# Thank you