# 人工智能技术及应用
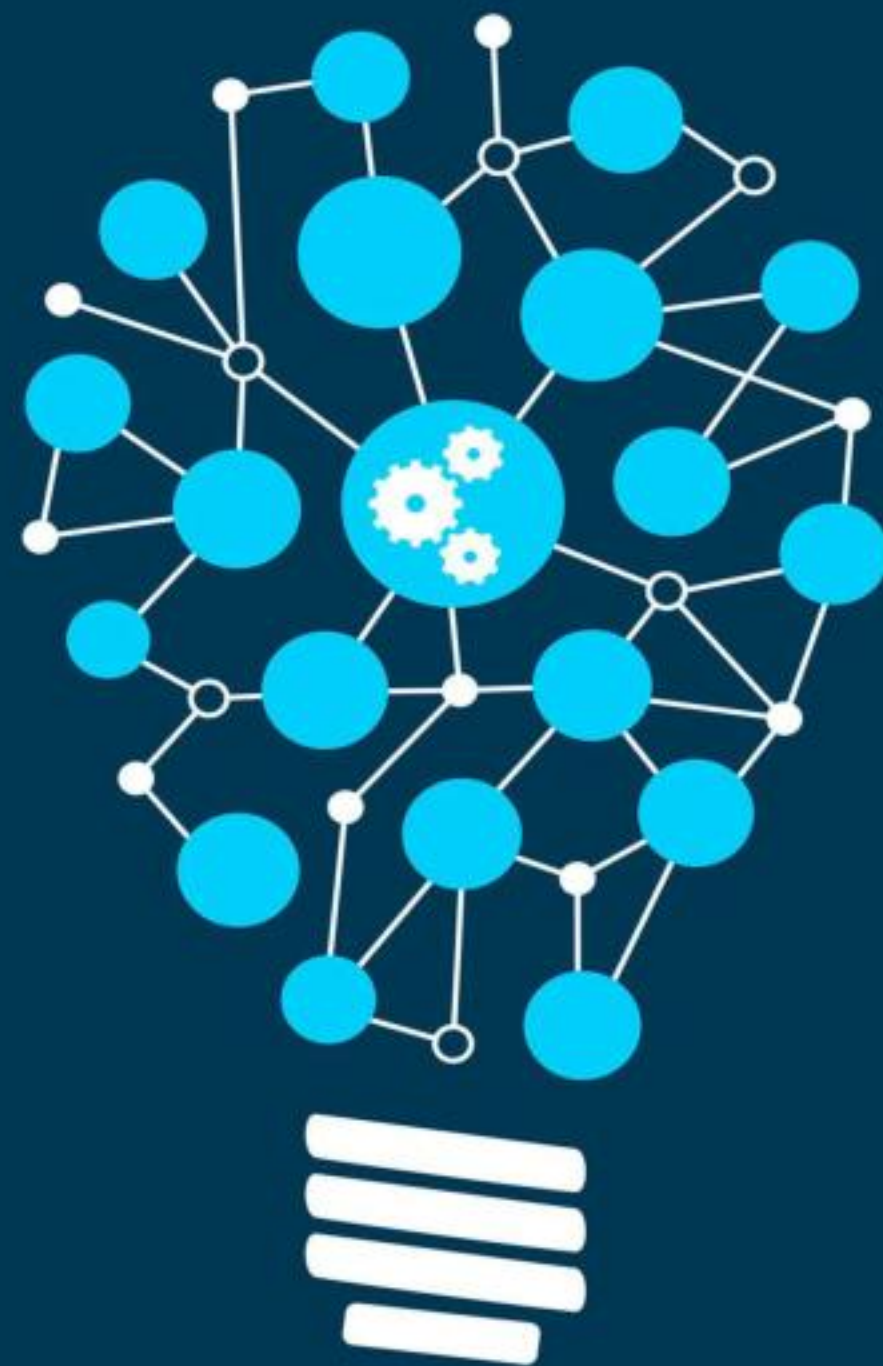
## Artificial Intelligence and Application
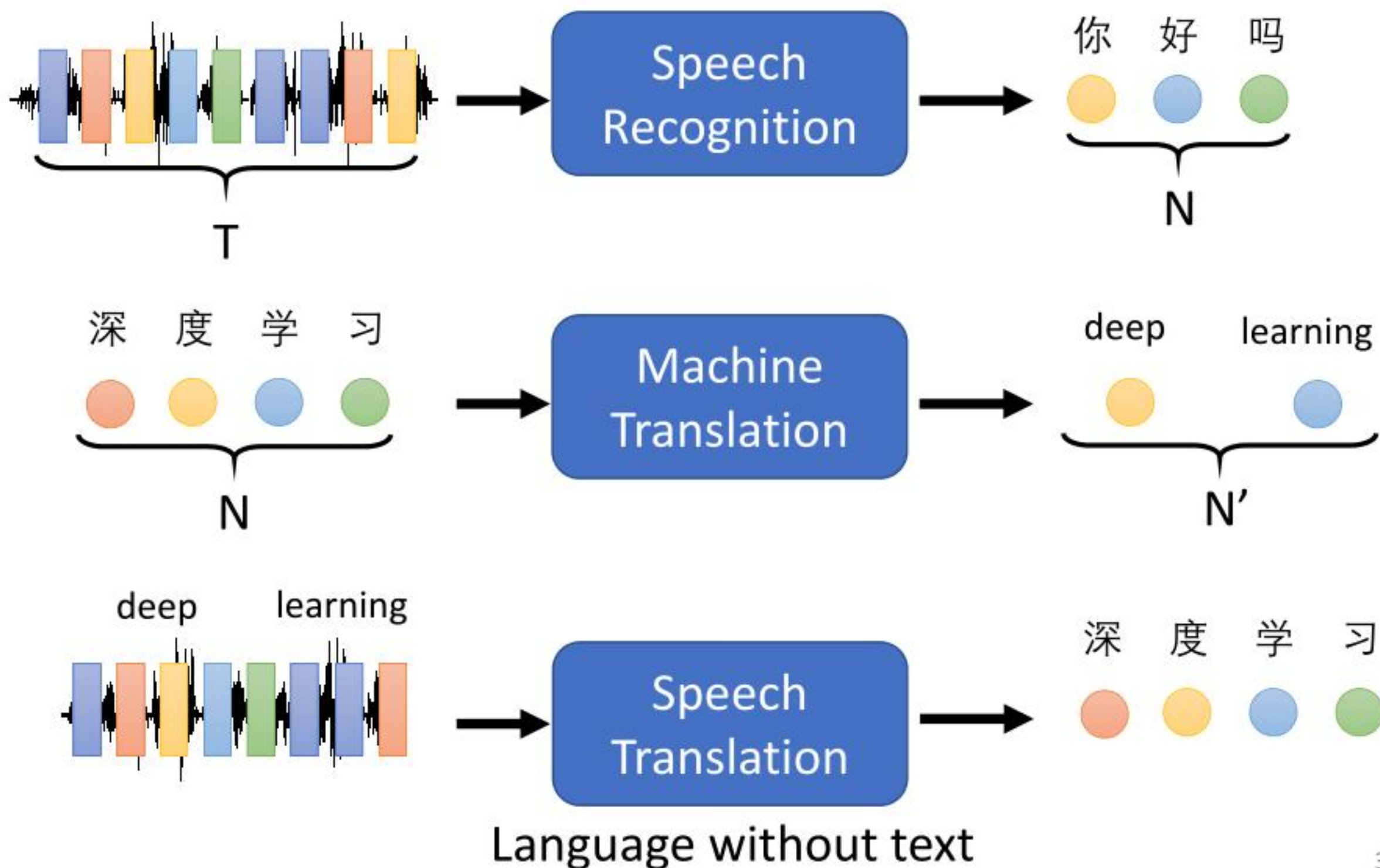
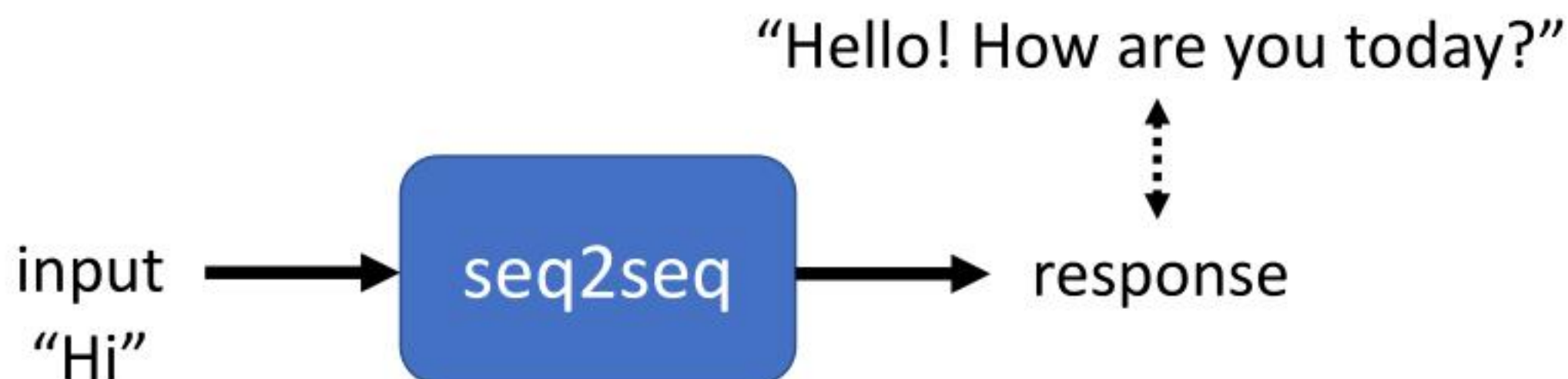# *Sequence-to-sequence (Seq2seq)*

Input a sequence, output a sequence

The output length is determined by model.

# Seq2seq for Chatbot

"Hello! How are you today?"

↕

input ⟶ **seq2seq** ⟶ response

"Hi"

Training data:

> [PERSON 1:] Hi
> [PERSON 2:] Hello ! How are you today ?
> [PERSON 1:] I am good thank you , how are you.
> [PERSON 2:] Great, thanks ! My children and I were just about to watch Game of Thrones.
> [PERSON 1:] Nice ! How old are your children?
> [PERSON 2:] I have four that range in age from 10 to 21. You?
> [PERSON 1:] I do not have children at the moment.
> [PERSON 2:] That just means you get to keep all the popcorn for yourself.
> [PERSON 1:] And Cheetos at the moment!
> [PERSON 2:] Good choice. Do you watch Game of Thrones?
> [PERSON 1:] No, I do not have much time for TV.
> [PERSON 2:] I usually spend my time painting: but, I love the show.

# *Most Natural Language Processing applications ...*

## Question Answering (QA)

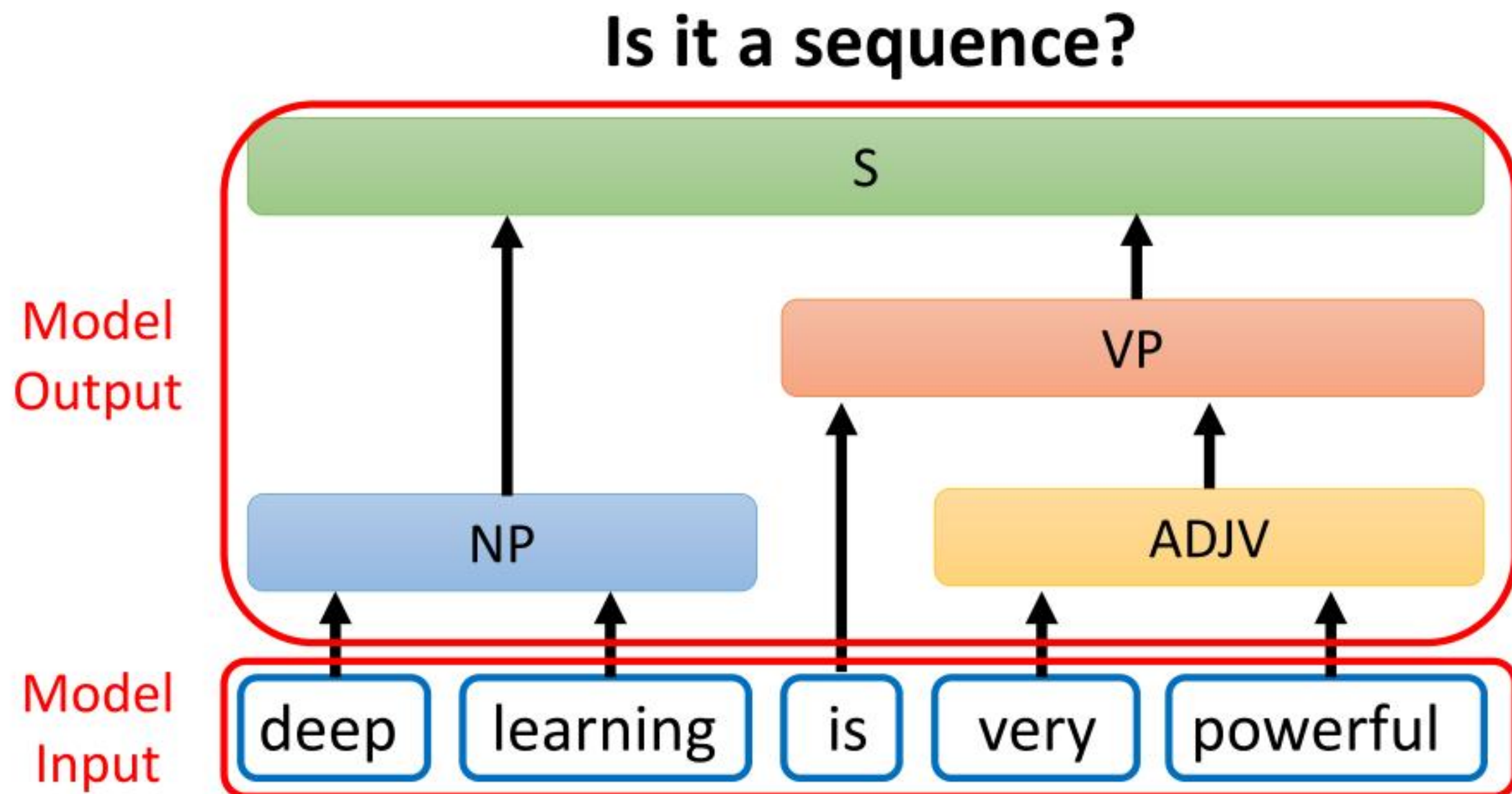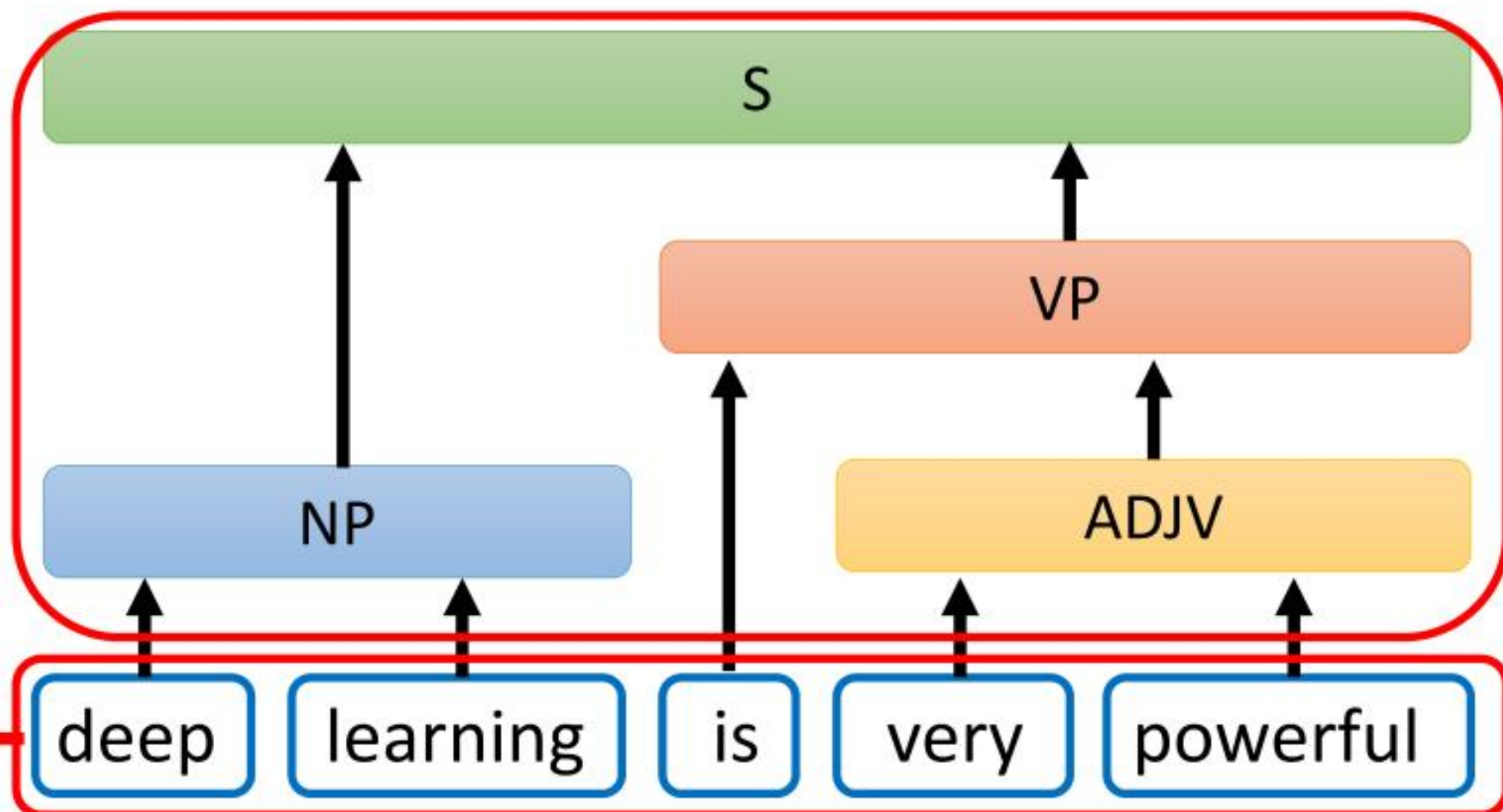| Question | Context | Answer |
|---|---|---|
| What is a major importance of Southern California in relation to California and the US? | ...Southern California is a major economic center for the state of California and the US.... | major economic center |
| What is the translation from English to German? | Most of the planet is ocean water. | Der Großteil der Erde ist Meerwasser |
| What is the summary? | Harry Potter star Daniel Radcliffe gains access to a reported £320 million fortune... | Harry Potter star Daniel Radcliffe gets £320M fortune... |
| Hypothesis: Product and geography are what make cream skimming work. Entailment, neutral, or contradiction? | Premise: Conceptually cream skimming has two basic dimensions – product and geography. | Entailment |
| Is this sentence positive or negative? (sentiment analysis) | A stirring, funny and finally transporting re-imagining of Beauty and the Beast and 1930s horror film. | positive |

decaNLP

**QA** can be done by **seq2seq**

question, context → Seq2seq → answer

https://arxiv.org/abs/1806.08730
https://arxiv.org/abs/1909.03329

5

# *Seq2seq for Syntactic Parsing*

## Is it a sequence?

# Seq2seq for Syntactic Parsing

(S    (NP    deep    learning    )    (VP    is

(ADJV    very    powerful    )    )    )

Seq2seq!



7

# Seq2seq for Syntactic Parsing

(S    (NP    deep    learning    )    (VP    is

(ADJV    very    powerful    )    )    )

**Grammar as a Foreign Language**

Oriol Vinyals*
Google
vinyals@google.com

Lukasz Kaiser*
Google
lukaszkaiser@google.com

Terry Koo
Google
terrykoo@google.com

Slav Petrov
Google
slav@google.com

Ilya Sutskever
Google
ilyasu@google.com

Geoffrey Hinton
Google
geoffhinton@google.com

https://arxiv.org
/abs/1412.7449

deep    learning    is    very    powerful

# Seq2seq for Multi-label Classification

An object can belong to multiple classes.

Class 1
Class 3

Class 1

Class 3
Class 9
Class 17

Class 10

Seq2seq

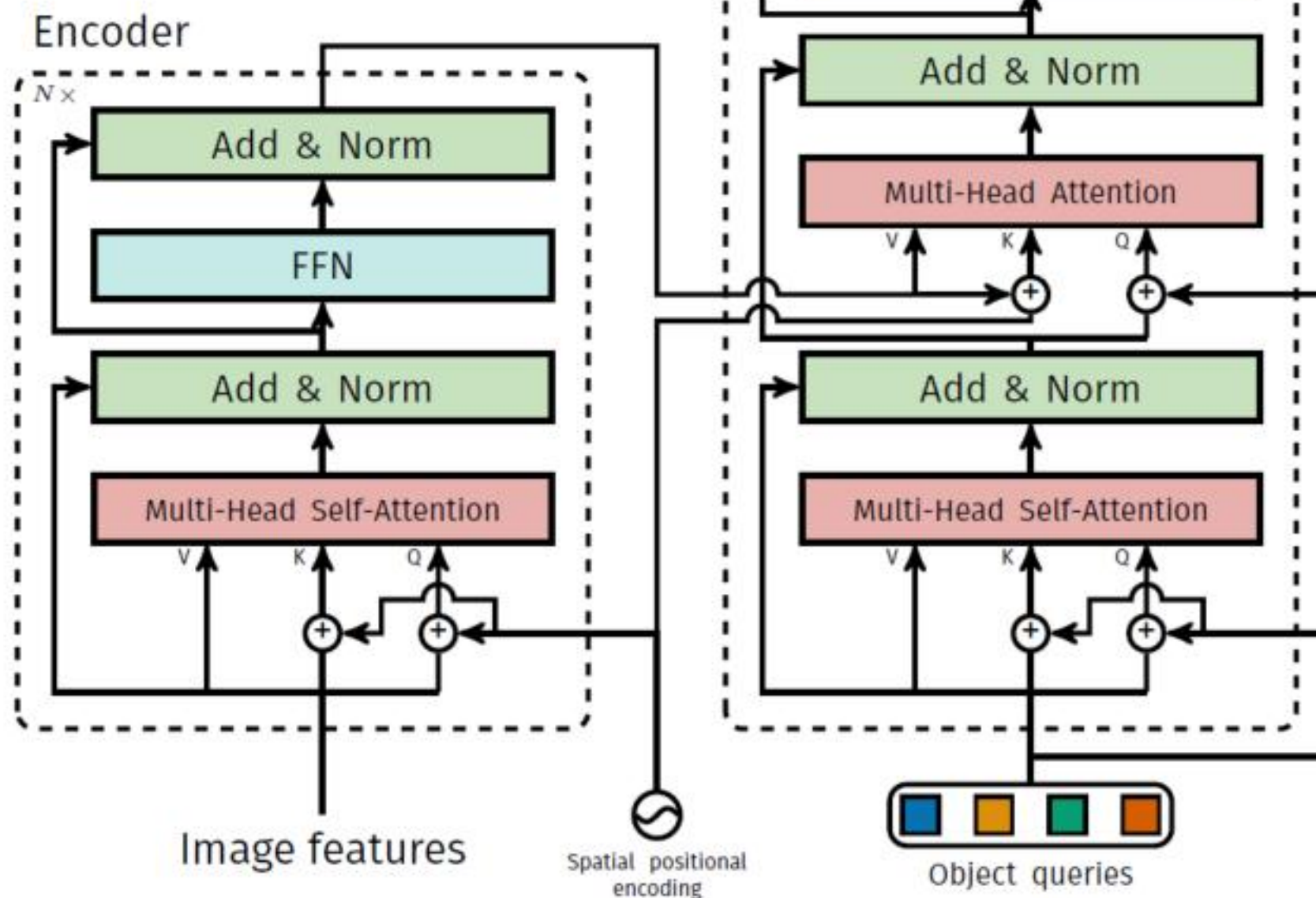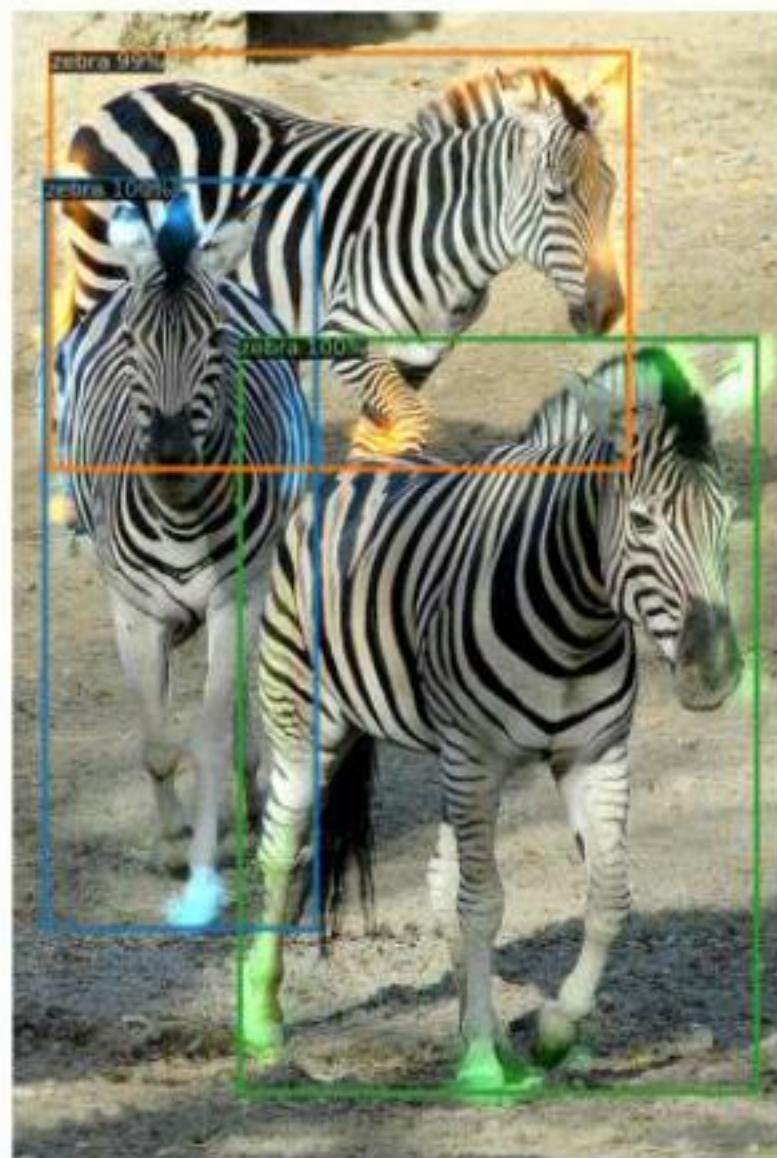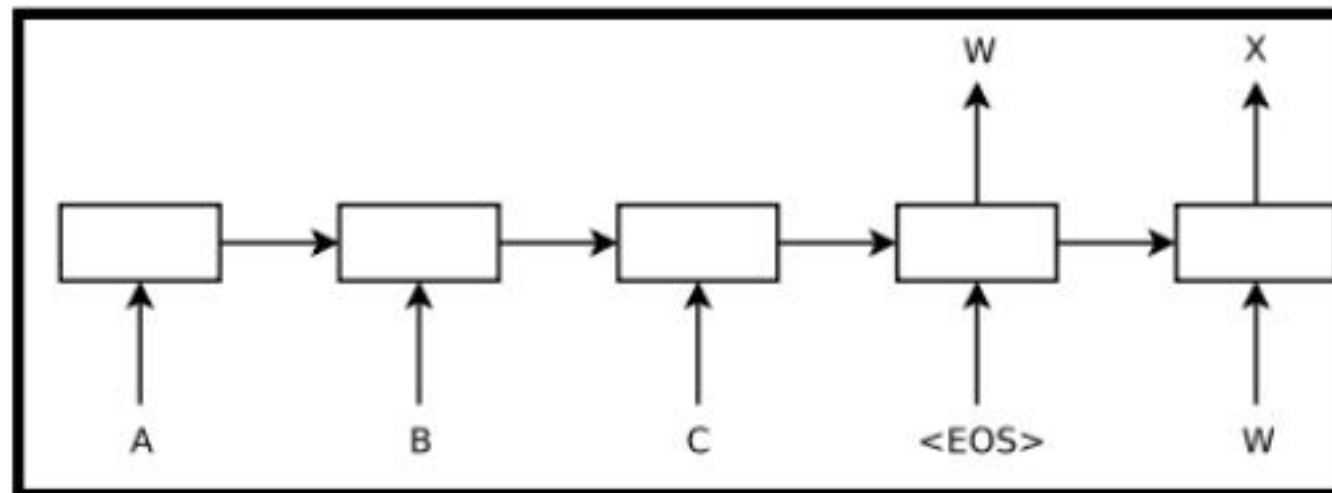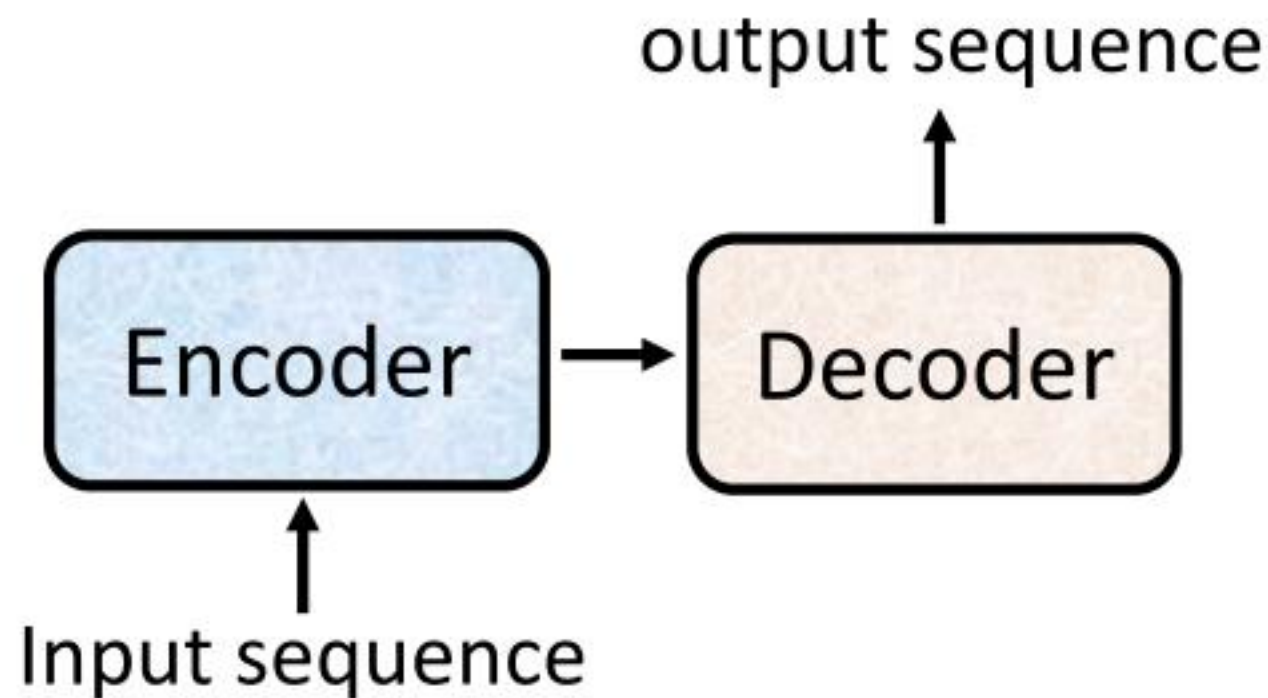Class 9    Class 7    Class 13

https://arxiv.org/abs/1909.03434
https://arxiv.org/abs/1707.05495

# Seq2seq for Object Detection

https://arxiv.org/abs/2005.12872

# Seq2seq

output sequence

Encoder → Decoder

Input sequence

W          X

A          B          C          <EOS>          W

Sequence to Sequence Learning with
Neural Networks
https://arxiv.org/abs/1409.3215

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm                  Add & Norm

Feed                        Multi-Head
Forward                     Attention
                                              Nx

Add & Norm                  Add & Norm

Nx    Multi-Head               Masked
      Attention              Multi-Head
                             Attention

Positional                              Positional
Encoding                                Encoding

Input                       Output
Embedding                   Embedding

Inputs                      Outputs
                            (shifted right)

Transformer
https://arxiv.org/abs/1706.03762
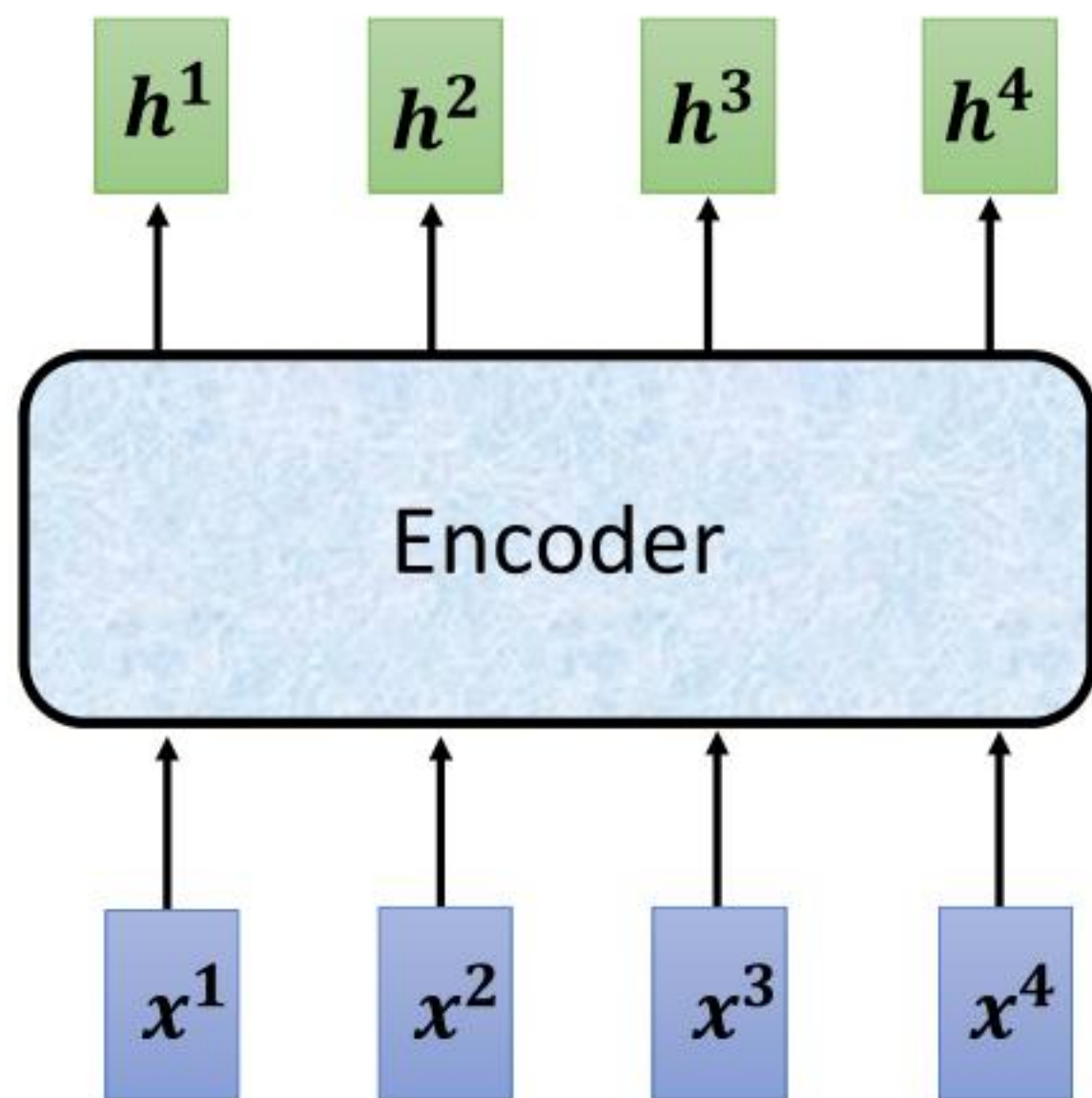
# Encoder

output sequence

Encoder → Decoder

Input sequence

# Encoder

## Transformer's Encoder

You can use **RNN** or **CNN**.

$h^1$ $h^2$ $h^3$ $h^4$

Encoder

$x^1$ $x^2$ $x^3$ $x^4$

Add & Norm

Feed Forward

Nx

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

residual $\boldsymbol{a} + \boldsymbol{b}$

$\boldsymbol{b}$ $+$ $\boldsymbol{a}$

$$\begin{bmatrix} x'_1 \\ x'_2 \\ \vdots \\ x'_K \end{bmatrix} \qquad x'_i = \frac{x_i - m}{\sigma}$$

Layer Norm
https://arxiv.org/
abs/1607.06450

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_K \end{bmatrix}$$

mean $m$

standard

deviation $\sigma$

norm

FC

Self-attention

15

# To learn more ......

- On Layer Normalization in the Transformer Architecture
- https://arxiv.org/abs/2002.047 45

- PowerNorm: Rethinking Batch Normalization in Transformers
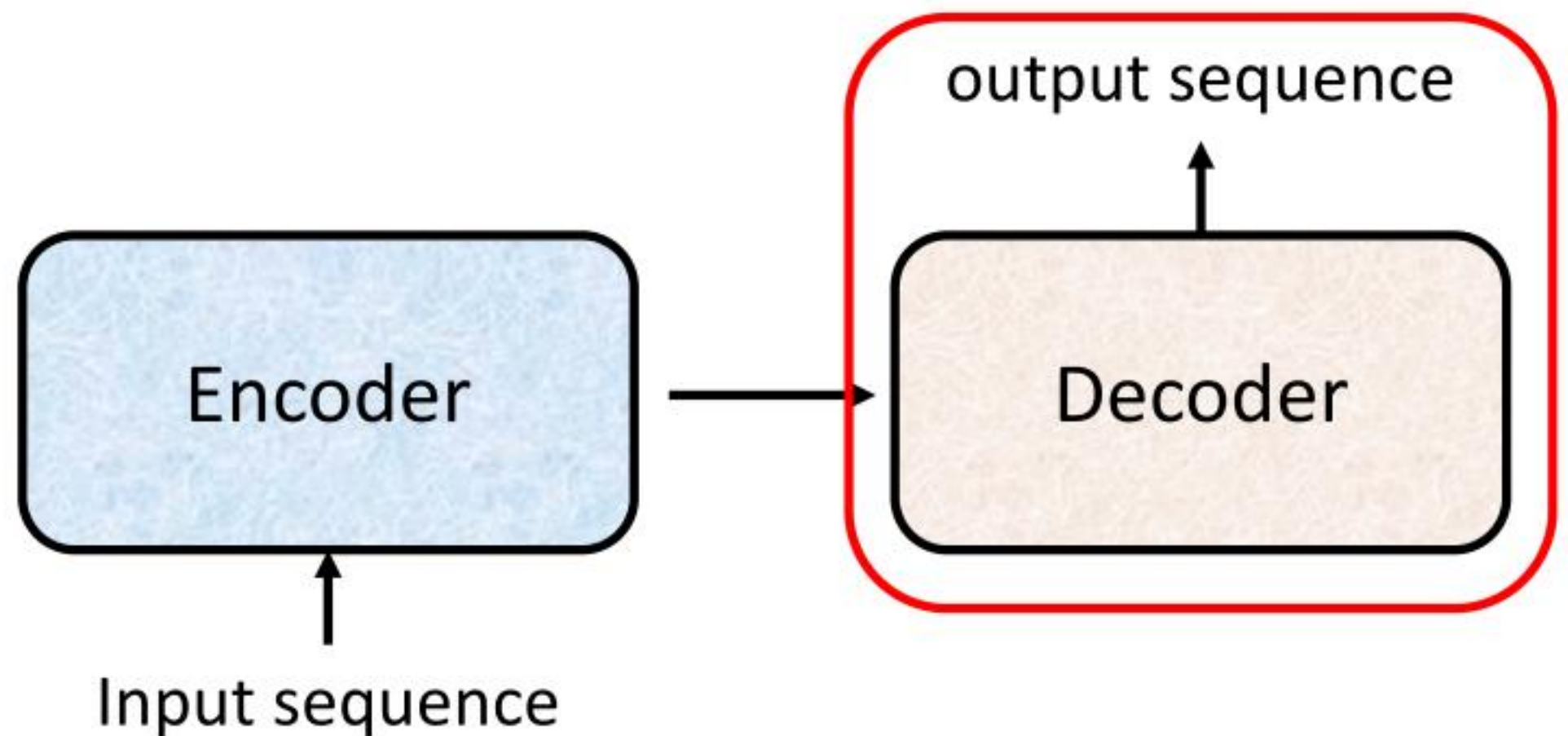- https://arxiv.org/abs/2003.078 45

# Decoder

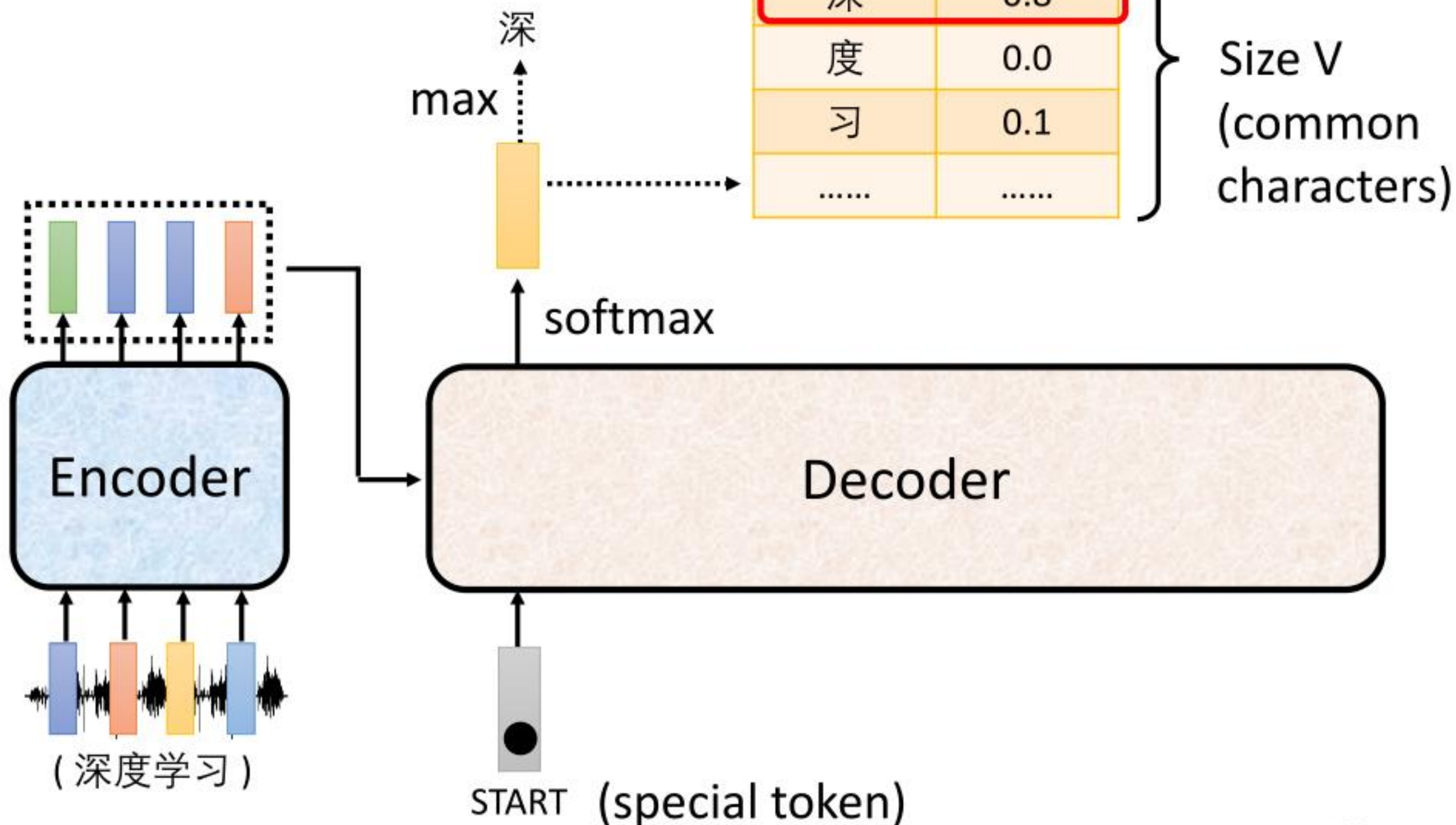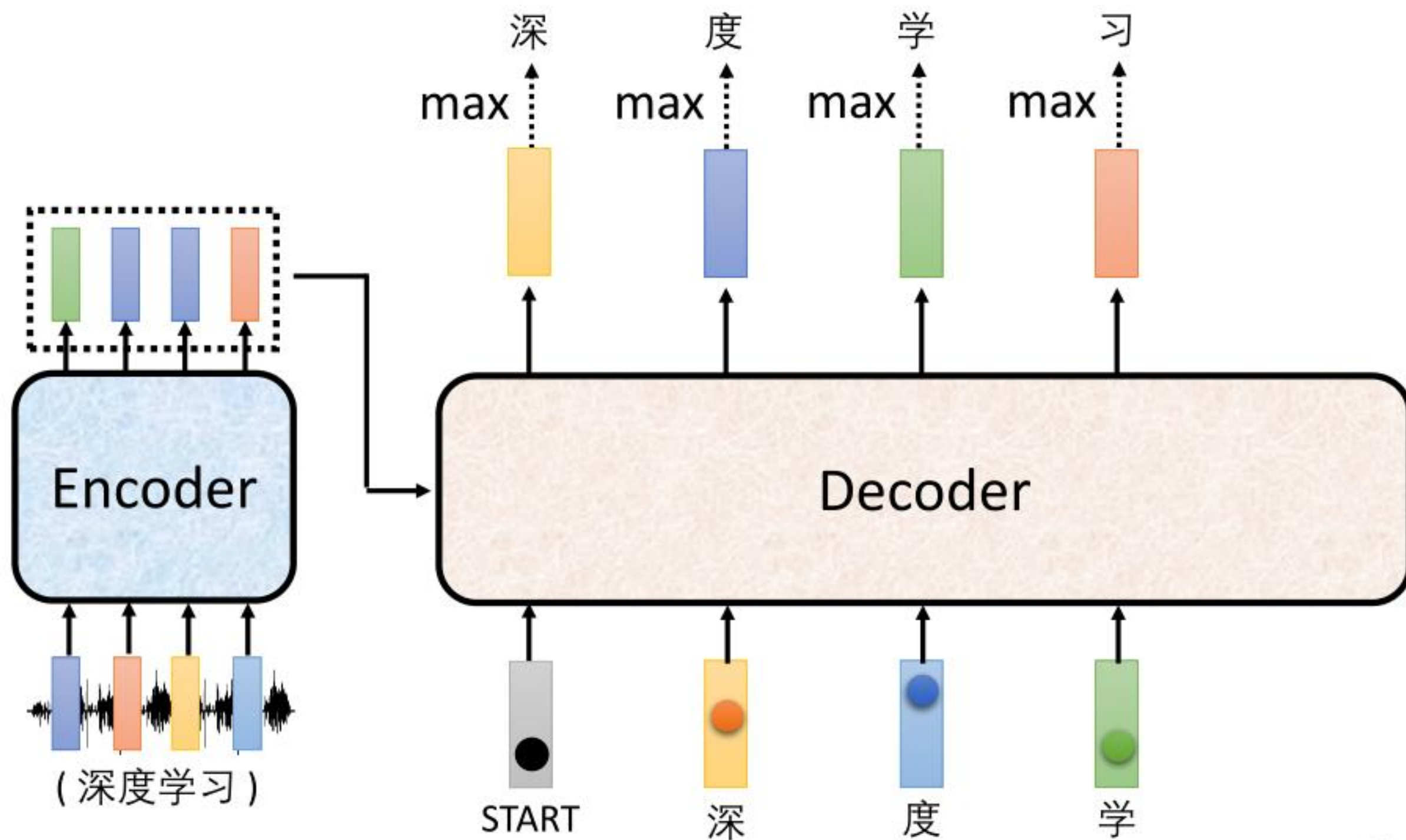

Encoder

Input sequence

Decoder

output sequence

# Decoder – Autoregressive (AT)

# Autoregressive
(Speech Recognition as example)

distribution

| 学 | 0.0 |
|---|---|
| 深 | 0.8 |
| 度 | 0.0 |
| 习 | 0.1 |
| ...... | ...... |

Size V (common characters)
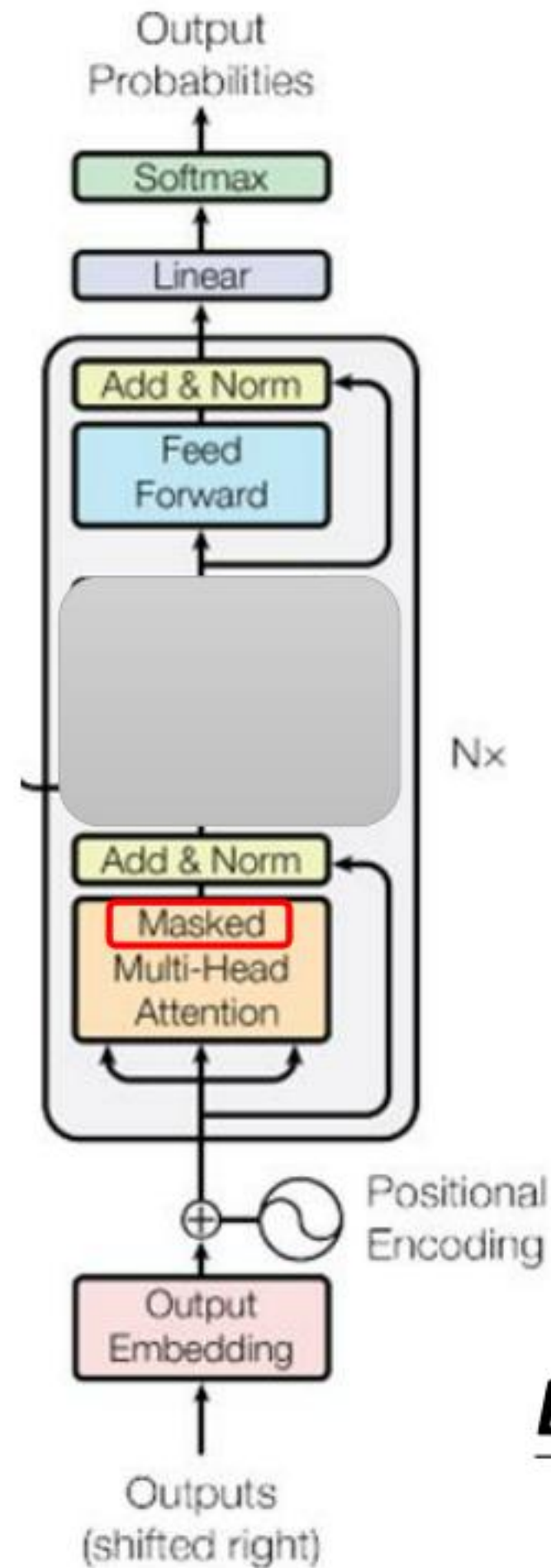
深

max

softmax

Encoder

Decoder

( 深度学习 )

START (special token)

# Autoregressive

# Autoregressive

ignore the input from the encoder here ☺

深　度　学　习

max　max　max　max

Decoder

START　深　度　学

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Nx

Positional Encoding
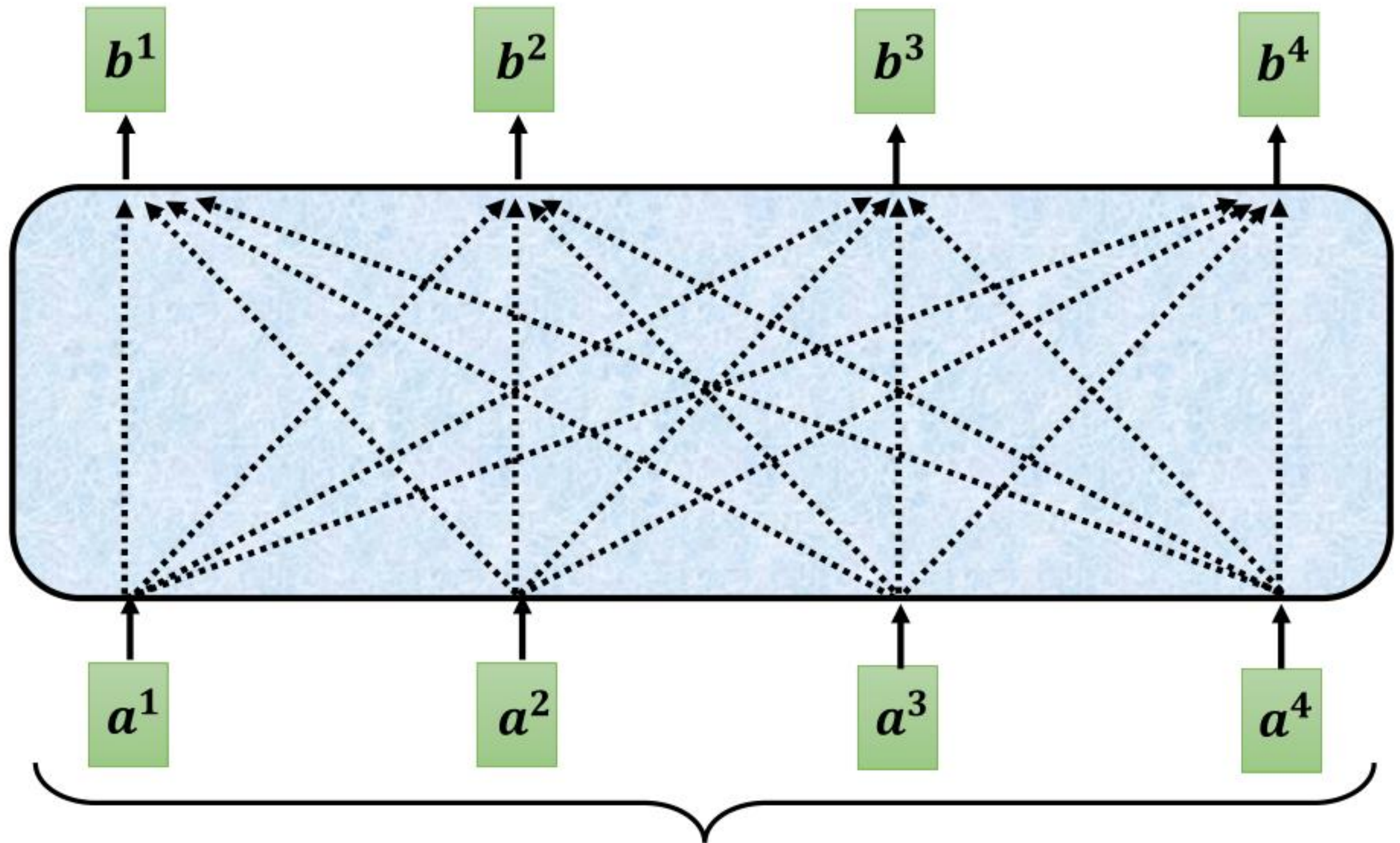
Output Embedding
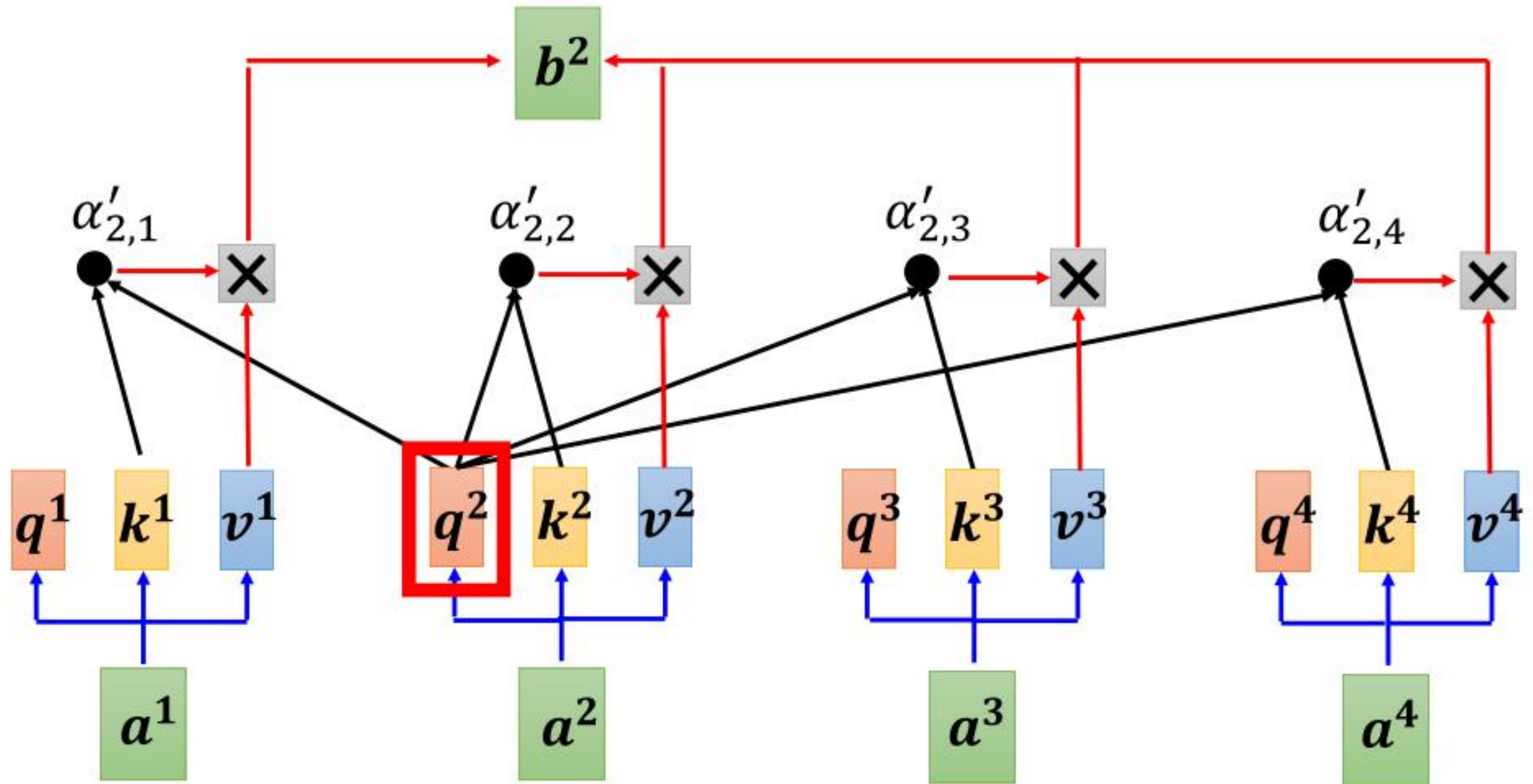
Outputs (shifted right)

**Encoder**

**Decoder**

24

# Self-attention ➡ Masked Self-attention
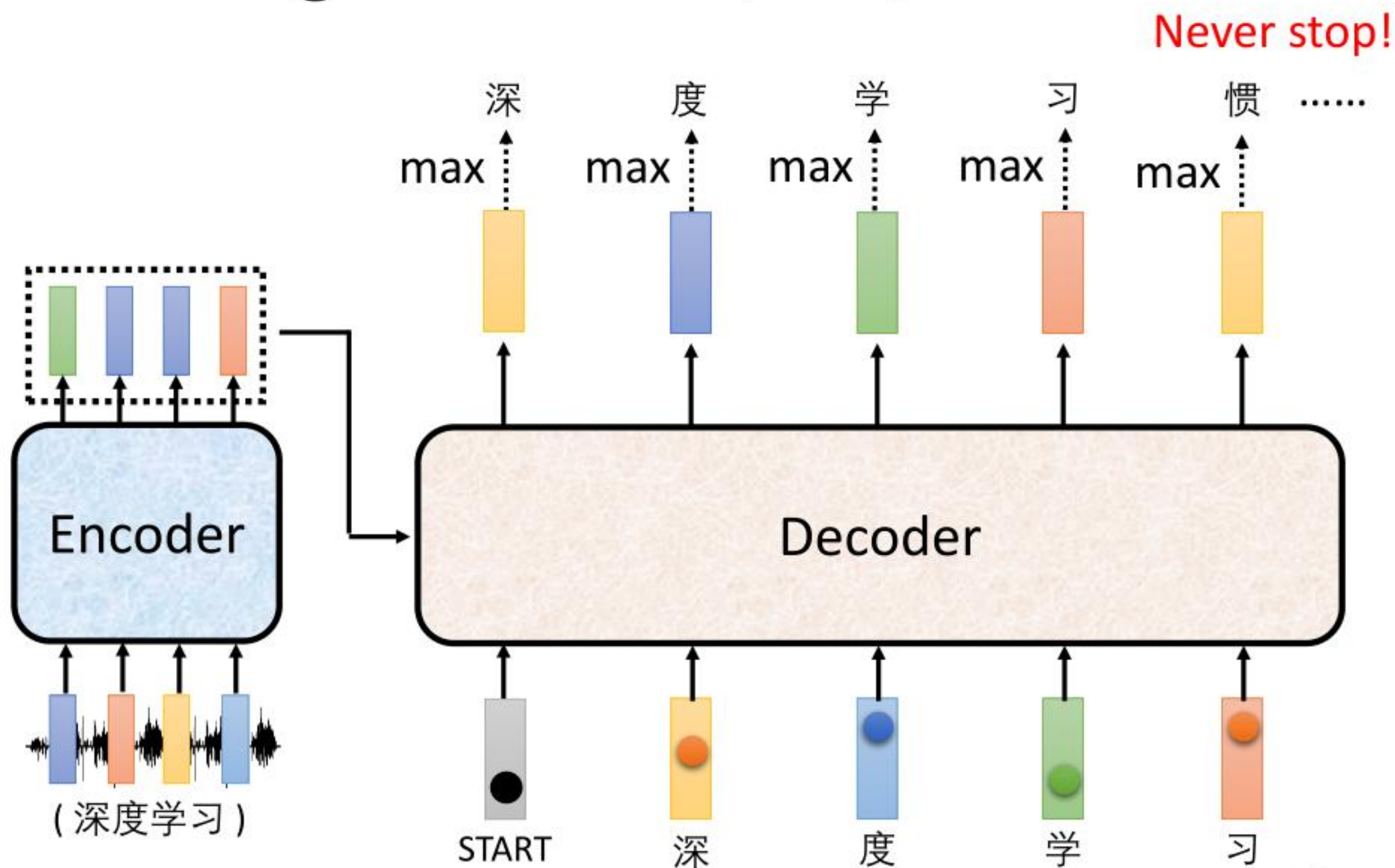


Can be either **input** or **a hidden layer**

# *Self-attention* → *Masked Self-attention*
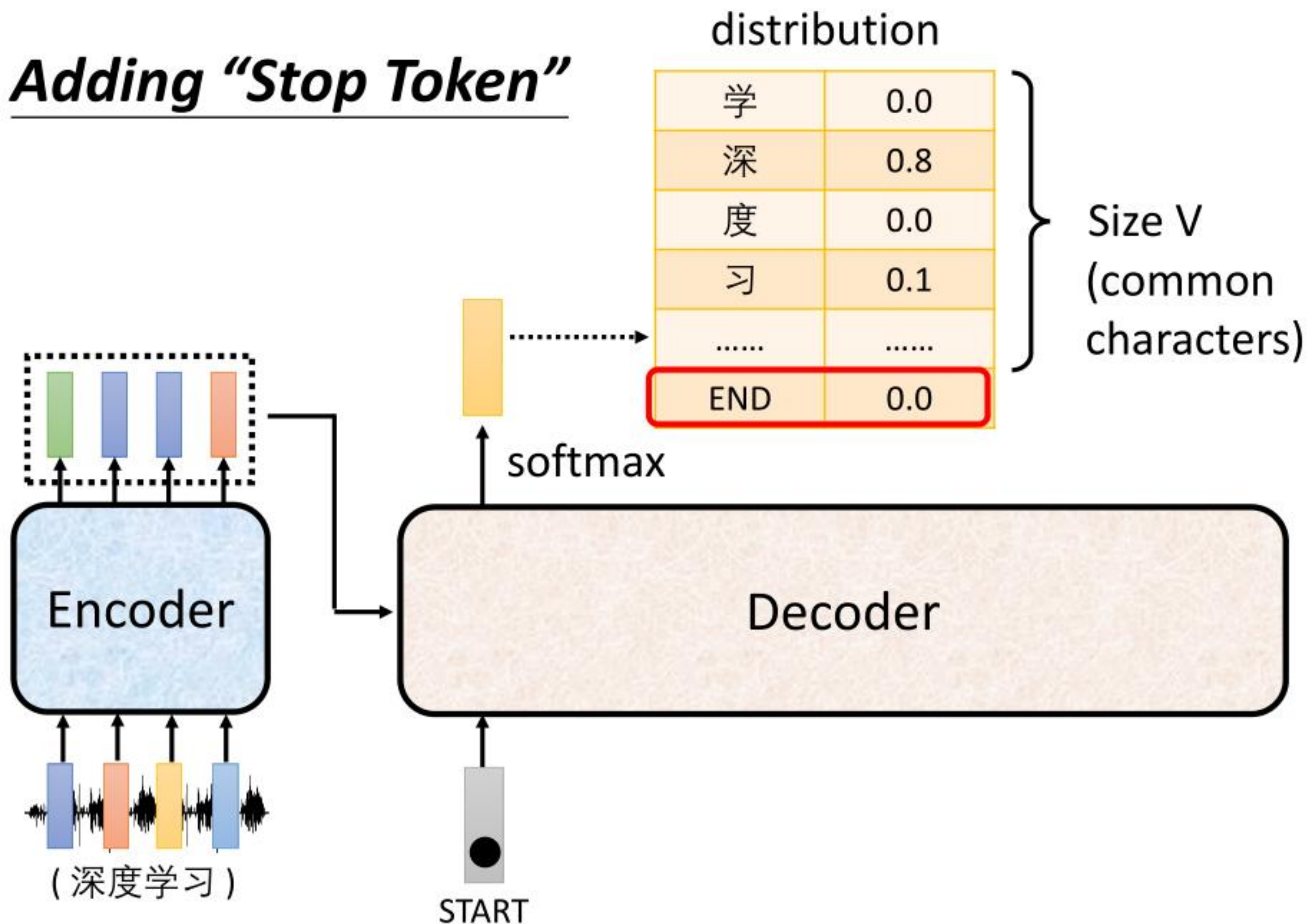


Why masked? Consider how does decoder work

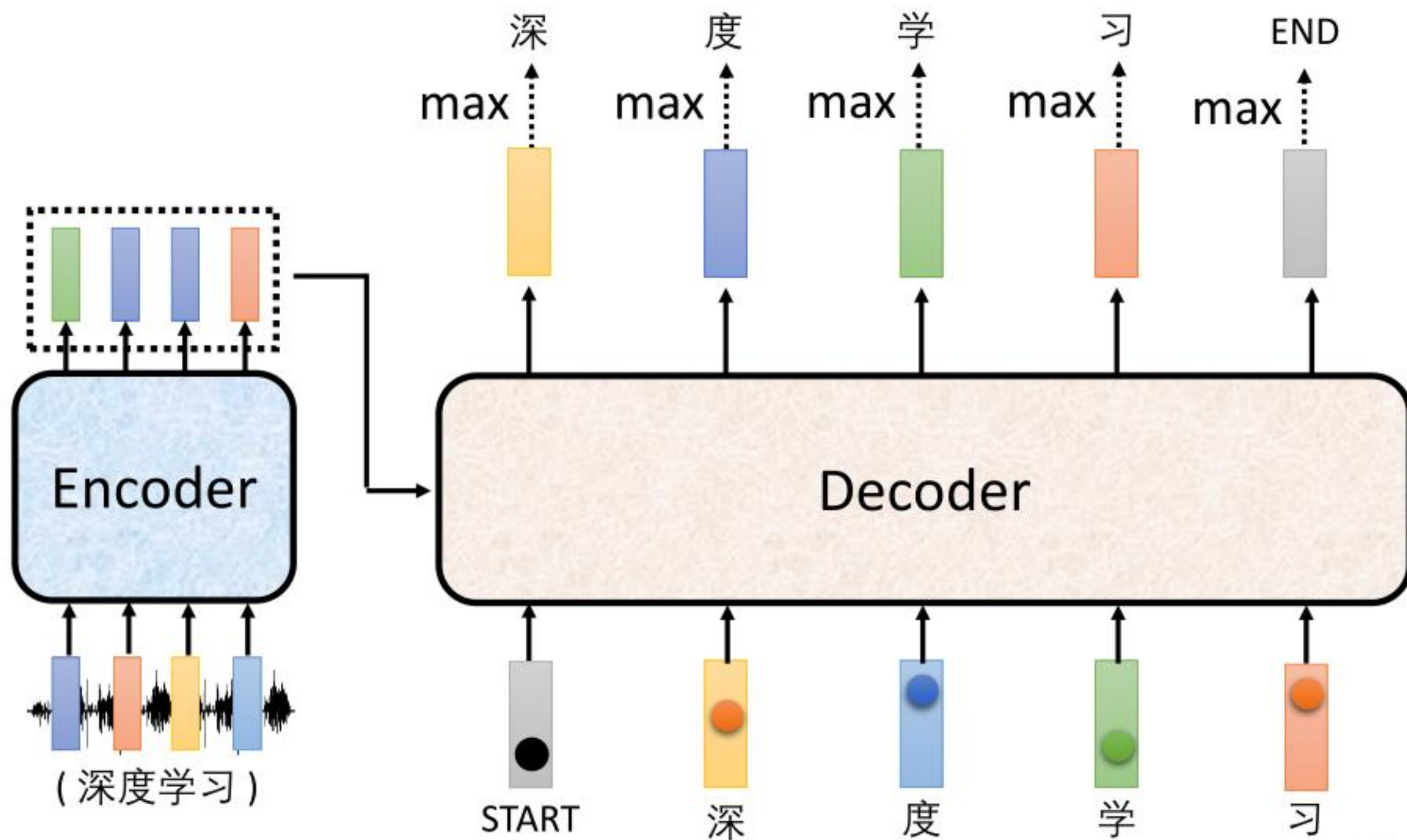# Autoregressive

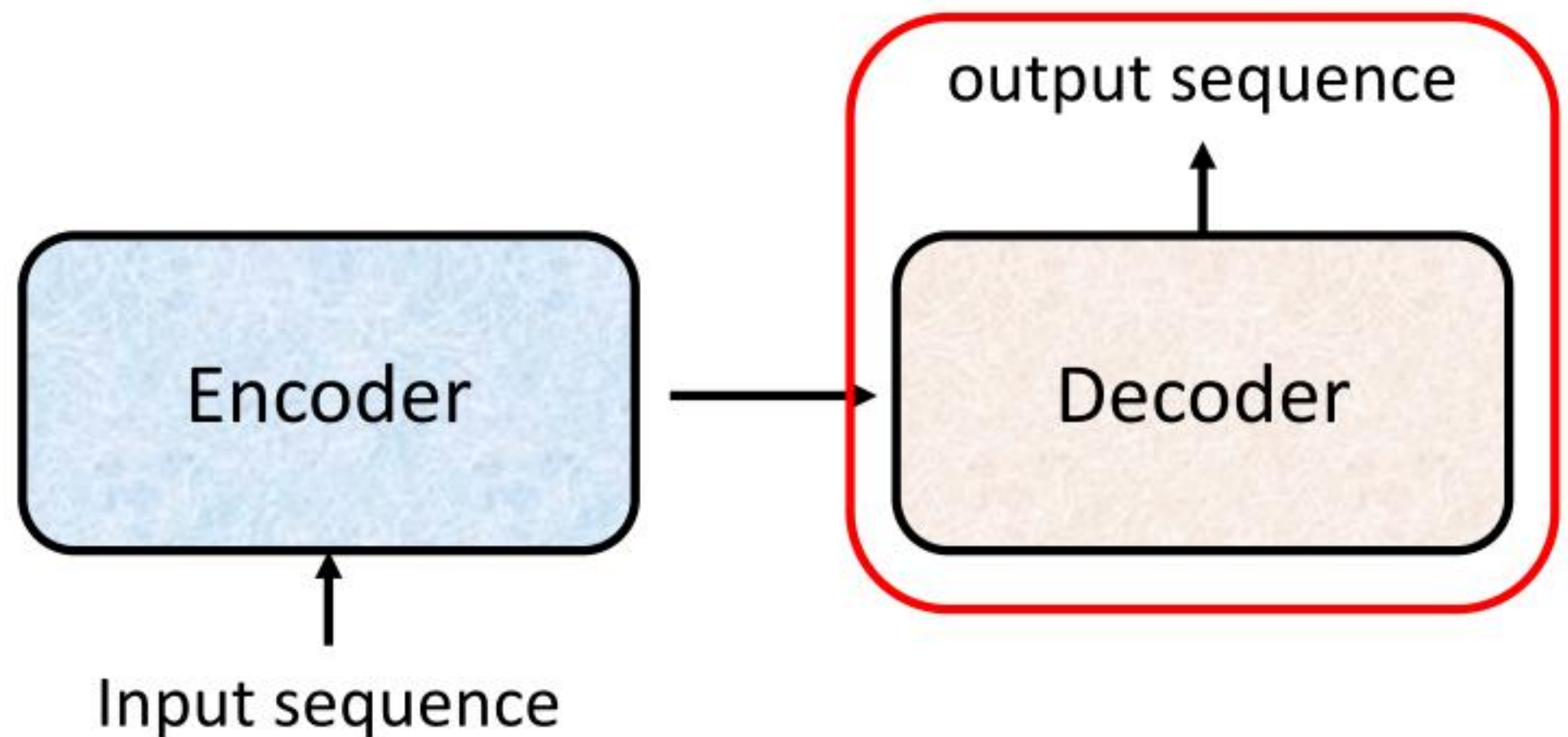We do not know the correct output length.

# *Adding "Stop Token"*

distribution

| | |
|---|---|
| 学 | 0.0 |
| 深 | 0.8 |
| 度 | 0.0 |
| 习 | 0.1 |
| ...... | ...... |
| END | 0.0 |

Size V (common characters)

softmax

Encoder

Decoder

( 深度学习 )

START

# Autoregressive

# Decoder – Non-autoregressive (NAT)



output sequence

Encoder

Decoder

Input sequence

# AT v.s. NAT



$w_1$  $w_2$  $w_3$  END

AT Decoder

START  $w_1$  $w_2$  $w_3$
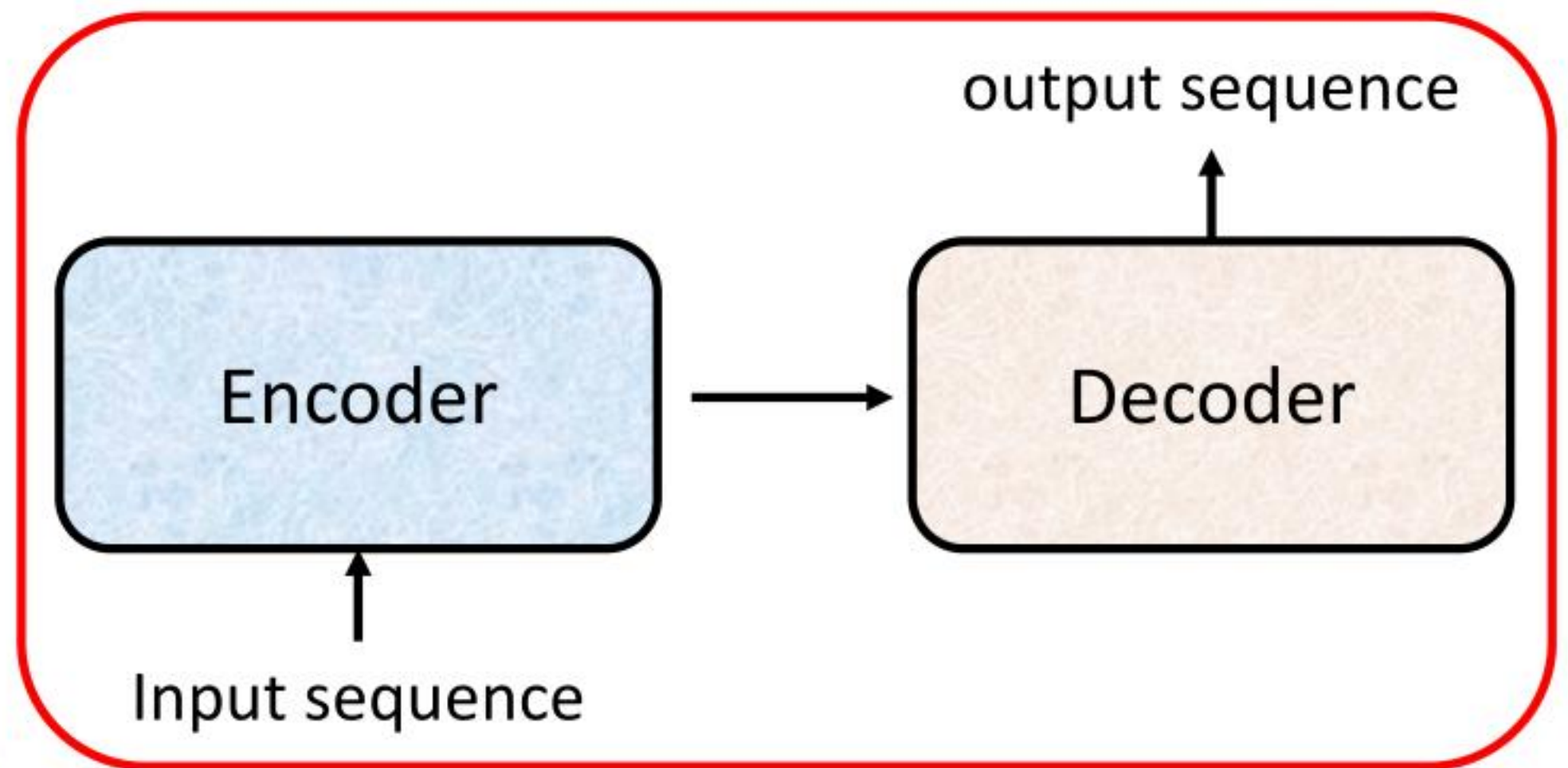
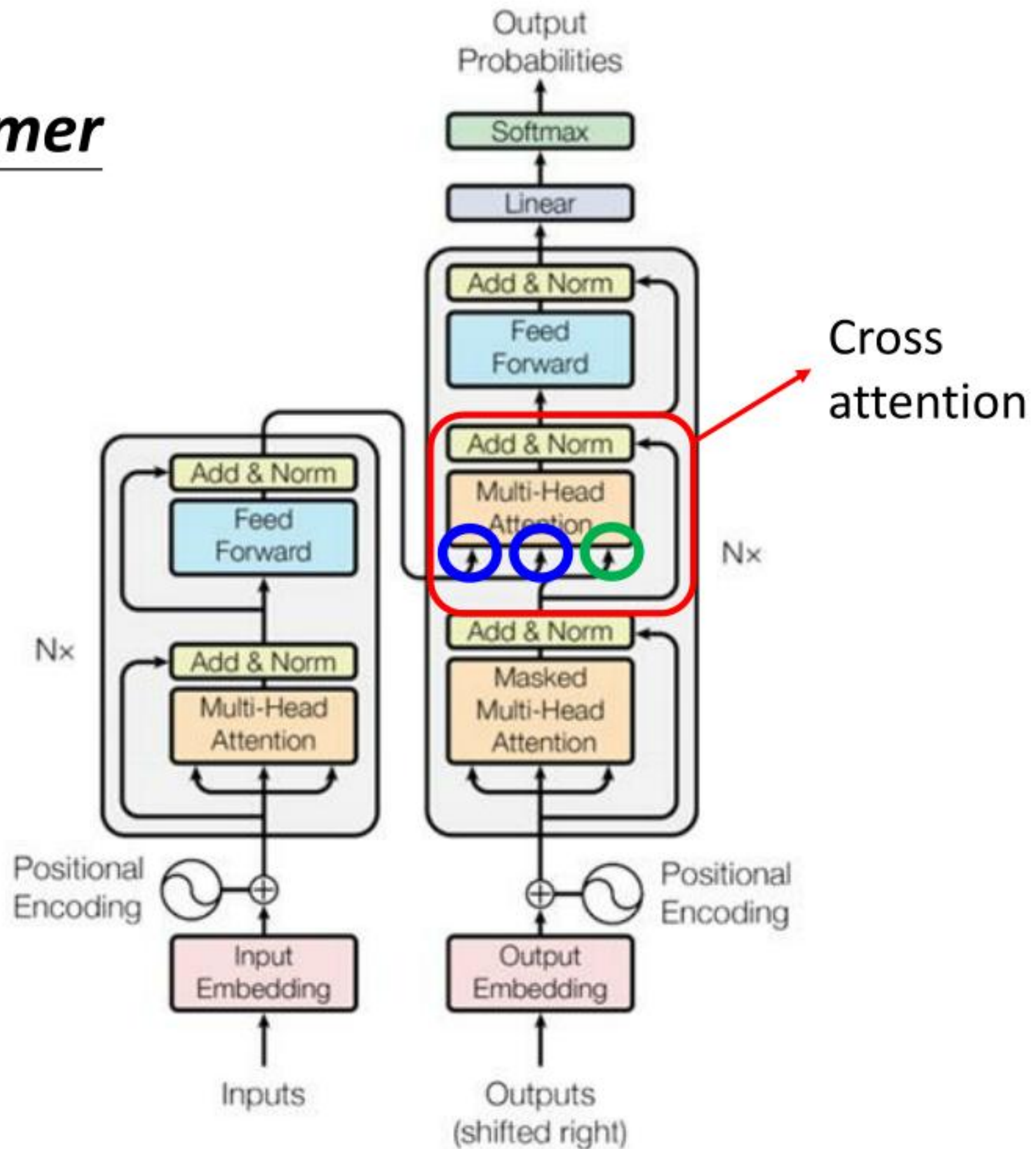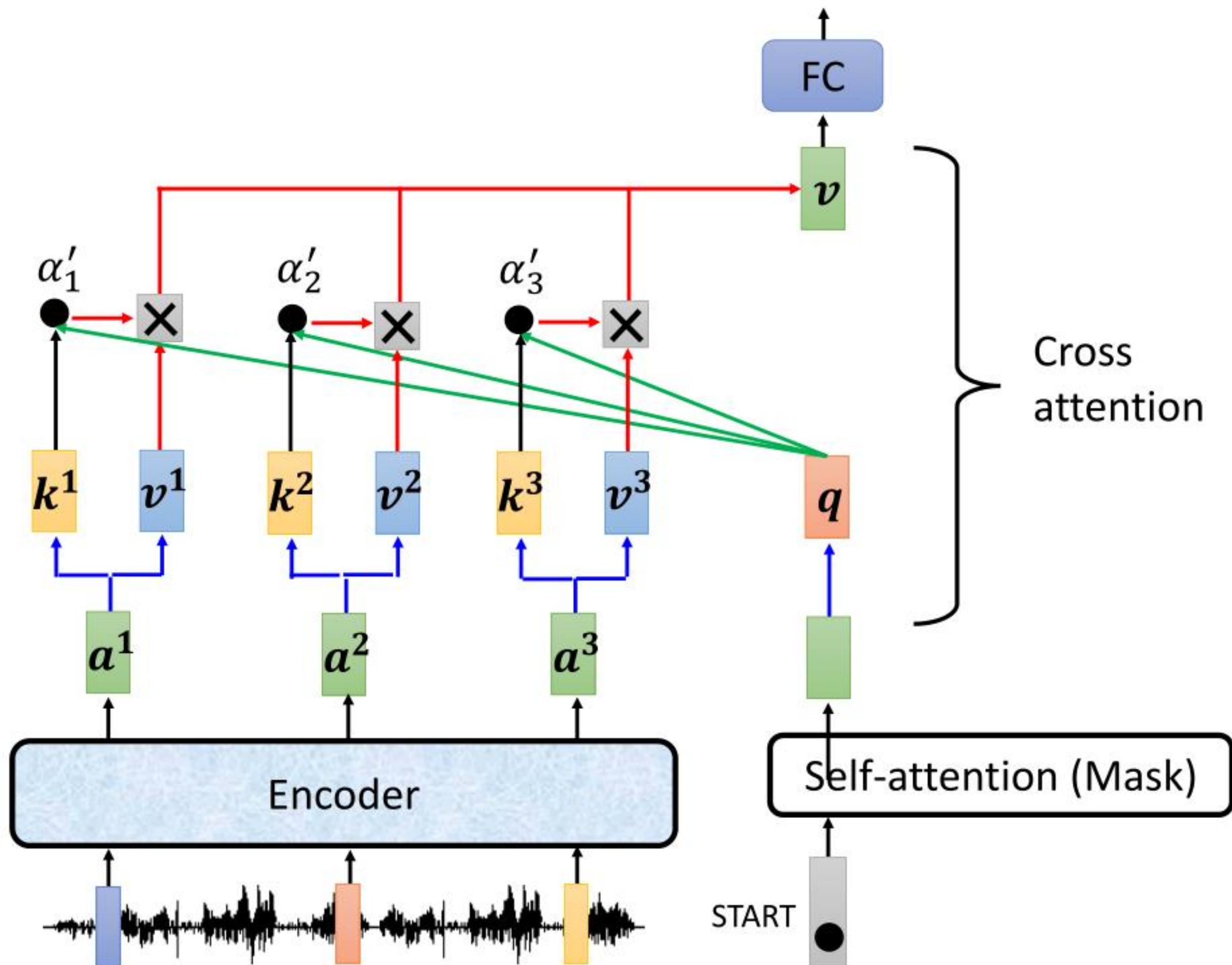$w_1$  $w_2$  END  $w_4$

NAT Decoder

START  START  START  START

➤ How to decide the output length for NAT decoder?

- Another predictor for output length

- Output a very long sequence, ignore tokens after END

➤ Advantage: parallel, more stable generation (e.g., TTS)

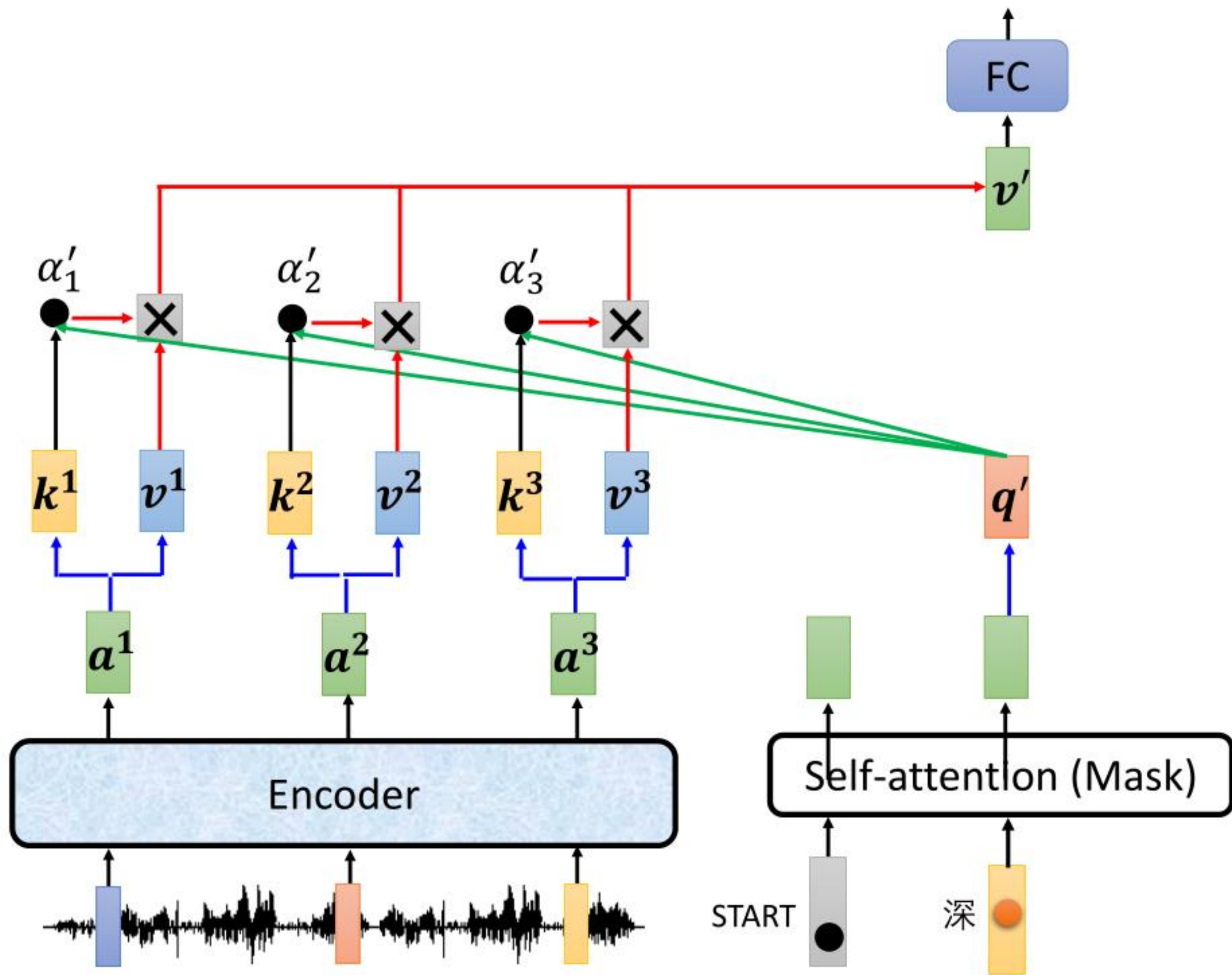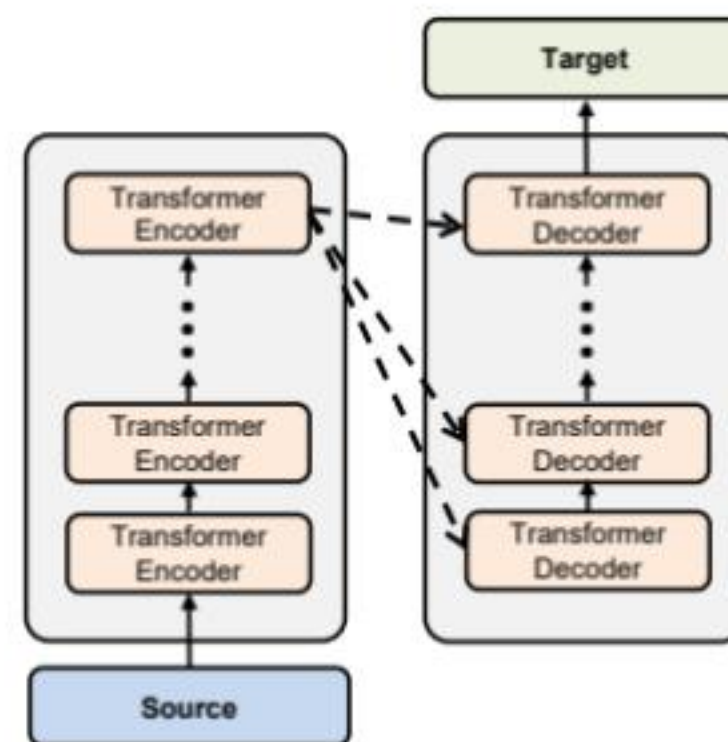➤ NAT is usually worse than AT  (why? **Multi-modality**)

# Encoder-Decoder



output sequence

Encoder → Decoder

Input sequence

# *Transformer*



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Cross attention

Nx

Add & Norm

Masked Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Positional Encoding

Input Embedding

Inputs

Positional Encoding

Output Embedding

Outputs (shifted right)
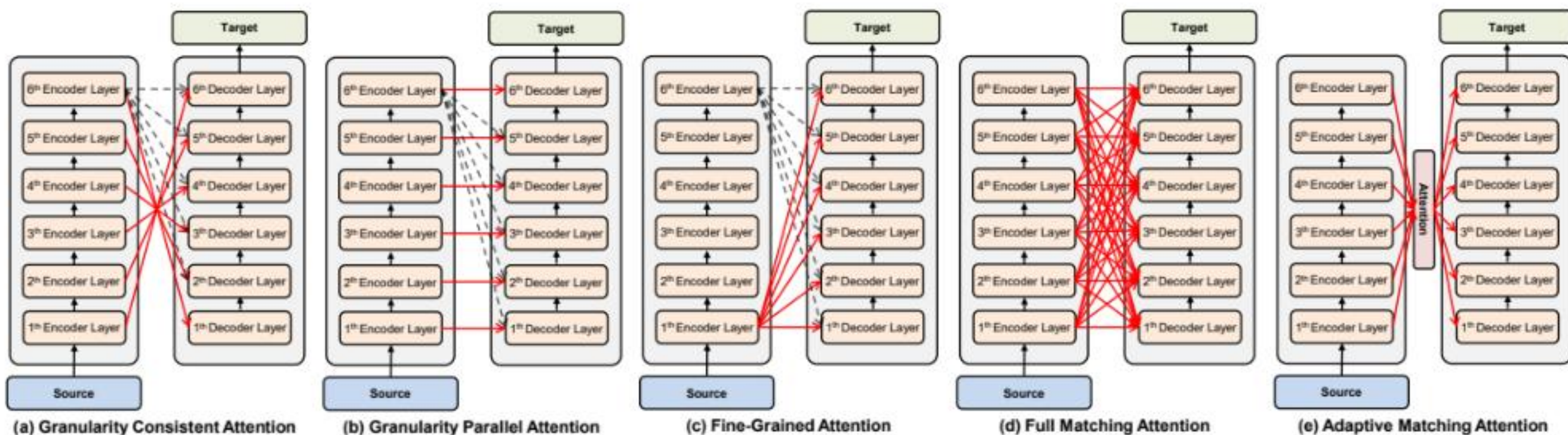
33

# Cross Attention

Source of image:
https://arxiv.org/abs/2005.08081



(a) Conventional Transformer
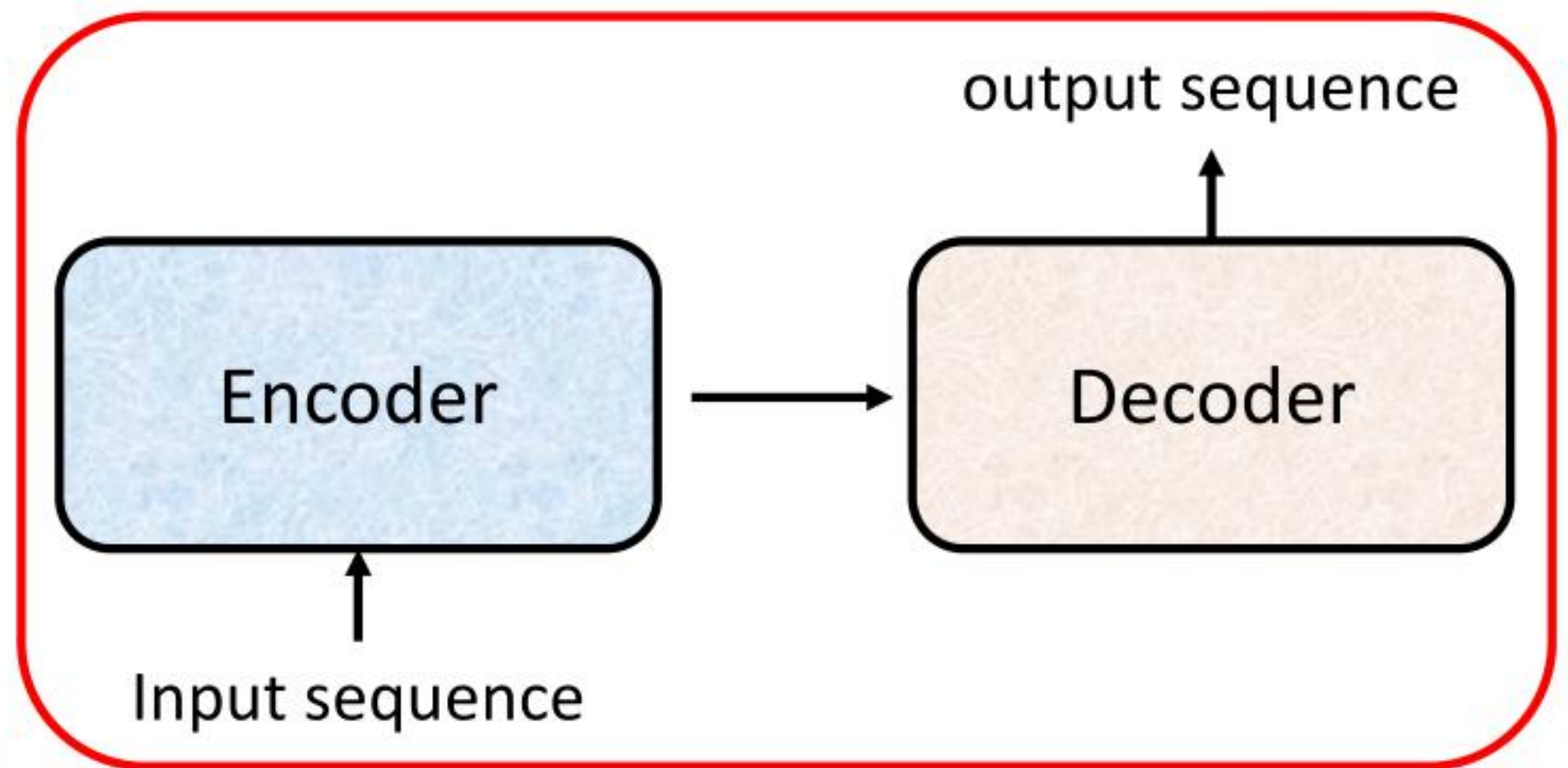


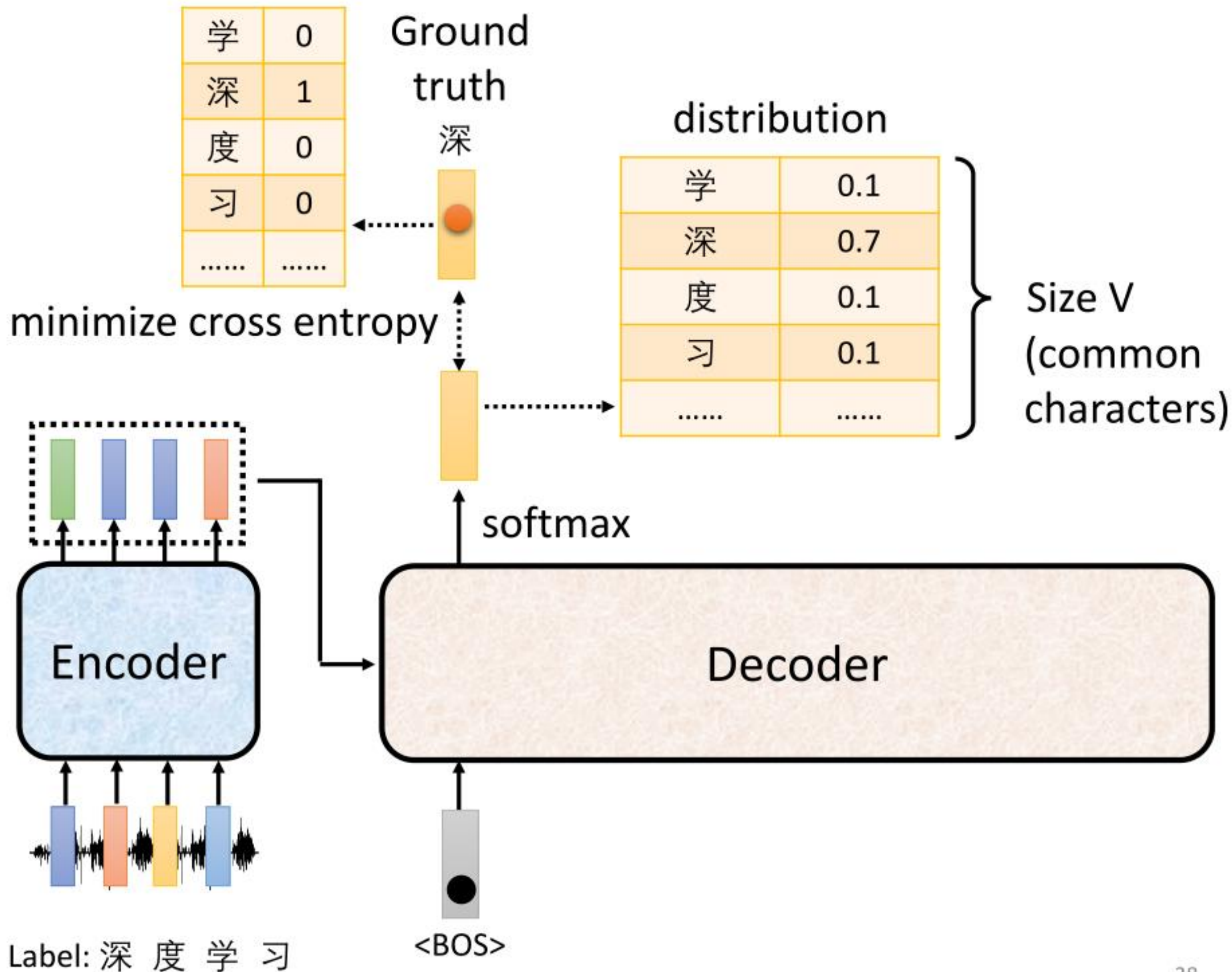(a) Granularity Consistent Attention    (b) Granularity Parallel Attention    (c) Fine-Grained Attention    (d) Full Matching Attention    (e) Adaptive Matching Attention

# Training

| 学 | 0 |
|---|---|
| 深 | 1 |
| 度 | 0 |
| 习 | 0 |
| …… | …… |

Ground truth

深

minimize cross entropy

distribution

| 学 | 0.1 |
|---|---|
| 深 | 0.7 |
| 度 | 0.1 |
| 习 | 0.1 |
| …… | …… |

Size V (common characters)

softmax

Encoder

Decoder

Label: 深 度 学 习

<BOS>

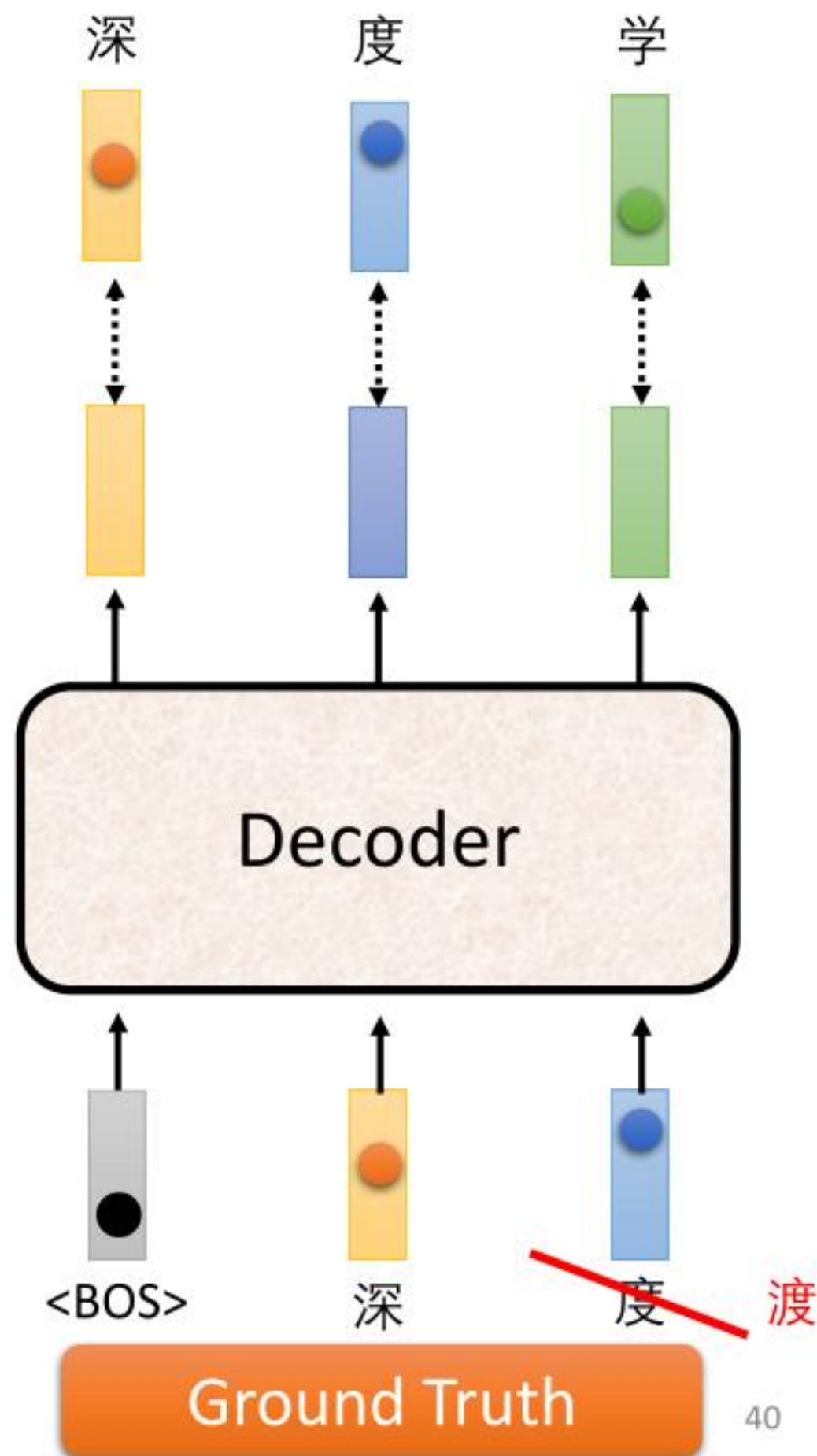# **Teacher Forcing**: using the ground truth as input.



minimize cross entropy

深　度　学　习　<EOS>

Encoder

Decoder

Label:深　度　学　习

<BOS>

Ground Truth

深　　　　度　　　　学　　　　习

There is a mismatch! ☹
**exposure bias**

Decoder

深 渡 学
深 渡

START 深 渡

Decoder

深 度 学

<BOS> 深 度 渡

Ground Truth

40

# Scheduled Sampling

- **Original Scheduled Sampling**

  https://arxiv.org/abs/1506.03099

- **Scheduled Sampling for Transformer**

  https://arxiv.org/abs/1906.07651

- **Parallel Scheduled Sampling**

  https://arxiv.org/abs/1906.04331



41

# Tips

output sequence

Encoder → Decoder

Input sequence

# Copy Mechanism

## *Machine Translation*

French: Guillaume et Cesar ont une voiture bleue a Lausanne.

Copy  Copy  Copy

English: Guillaume and Cesar have a blue car in Lausanne.
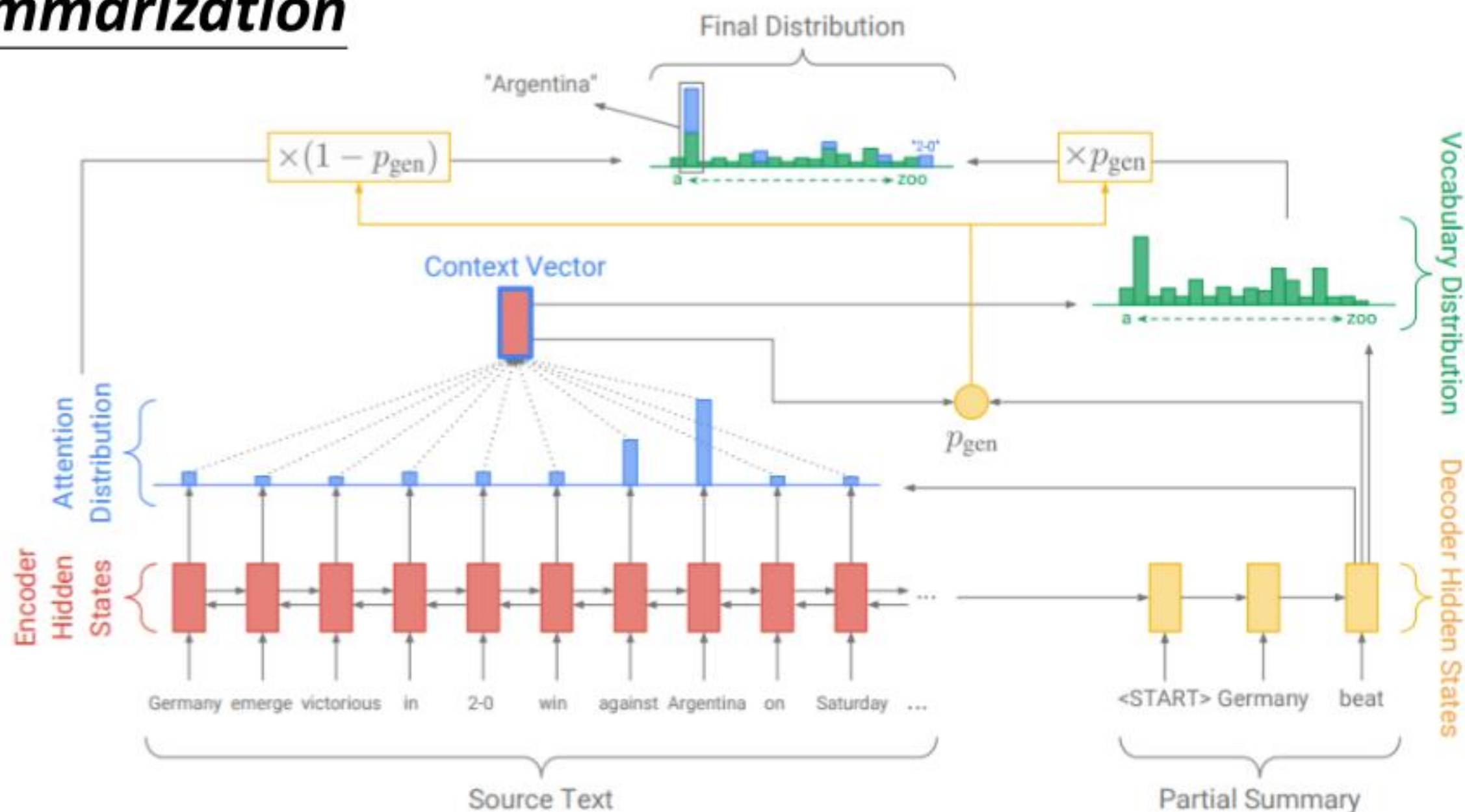
## *Chat-bot*

User: 你好，我是甘道夫

Machine:甘道夫你好，很高兴认识你

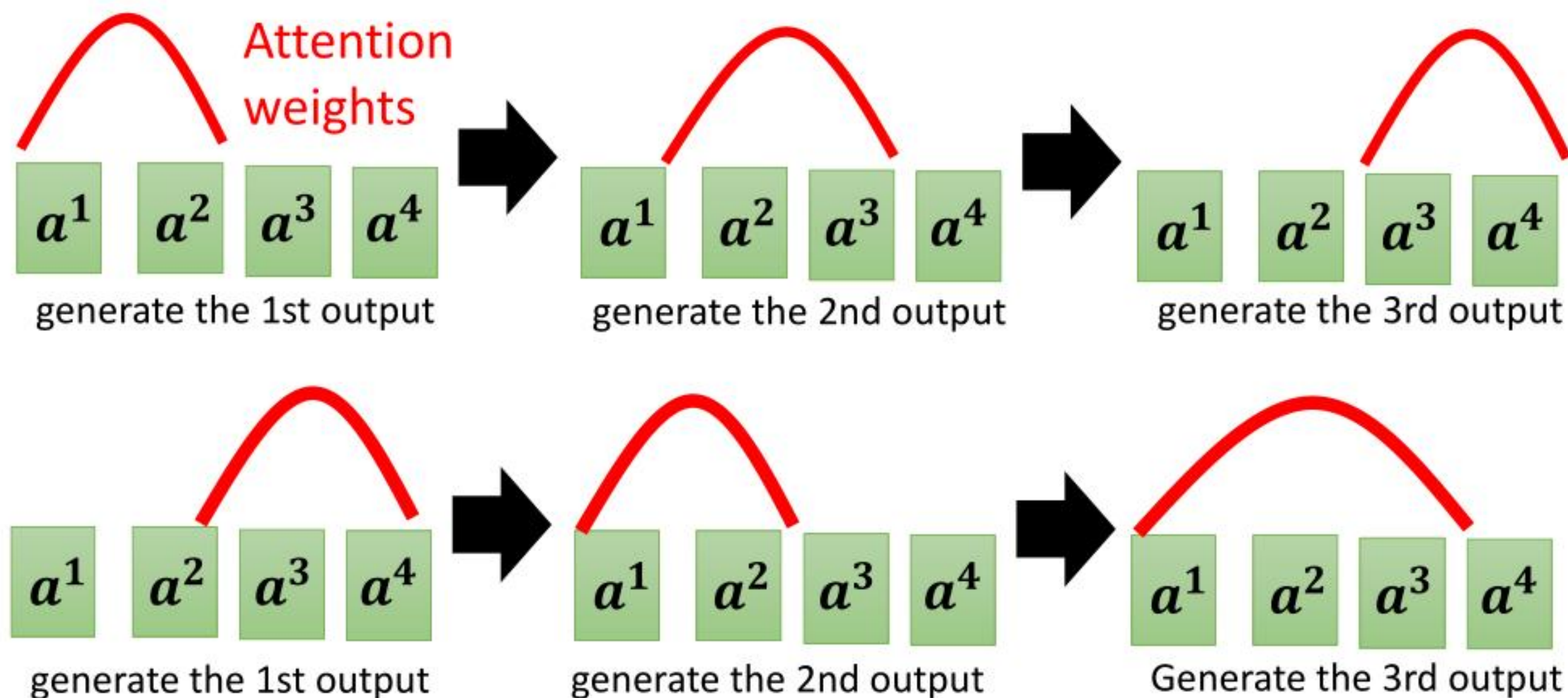# Copy Mechanism

## *Summarization*

# Guided Attention

Monotonic Attention
Location-aware attention

In some tasks, input and output are monotonically aligned.
For example, speech recognition, TTS, etc.



Attention weights

generate the 1st output → generate the 2nd output → generate the 3rd output

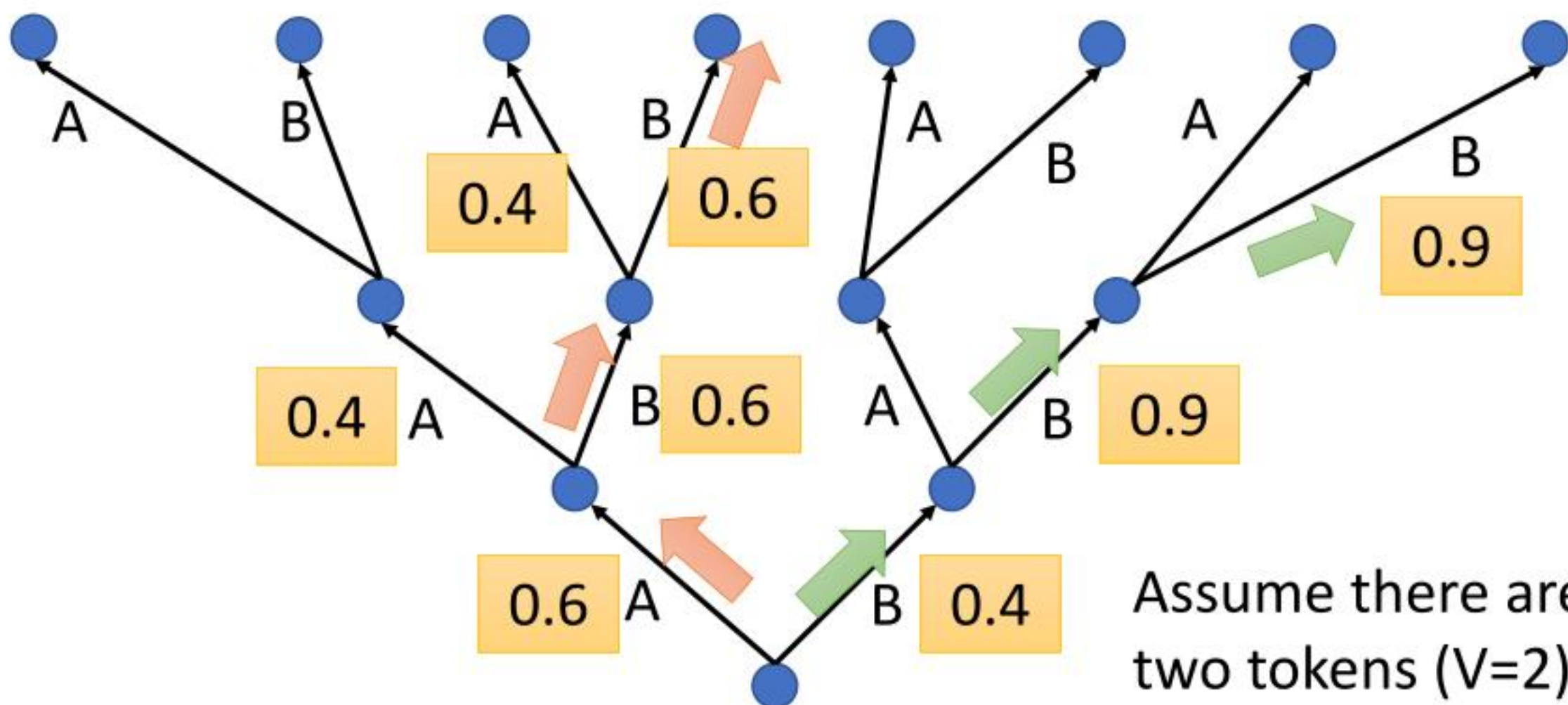generate the 1st output → generate the 2nd output → Generate the 3rd output

**Something wrong!**

# Beam Search

The red path is **Greedy Decoding**.

The green path is the best one.

Not possible to check all the paths …    → Beam Search



Assume there are only two tokens (V=2).

# Sampling

## The Curious Case of Neural Text Degeneration

**Context**: In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

**Beam Search, *b*=32**:
"The study, published in the Proceedings of the National Academy of Sciences of the United States of America (PNAS), was conducted by researchers from the Universidad Nacional Autónoma de México (UNAM) and the Universidad Nacional Autónoma de México (UNAM/Universidad Nacional Autónoma de México/Universidad Nacional Autónoma de México/Universidad Nacional Autónoma de México/Universidad Nacional Autónoma de ..."
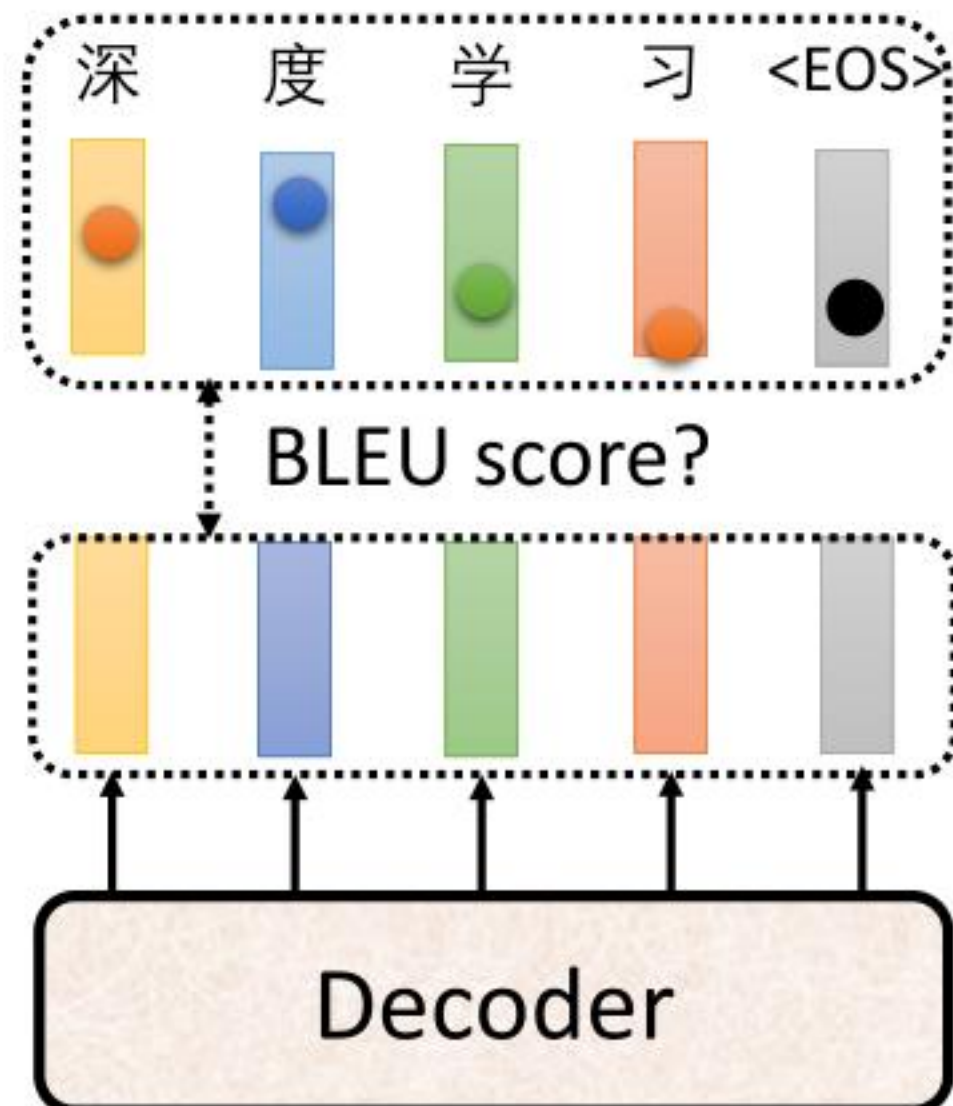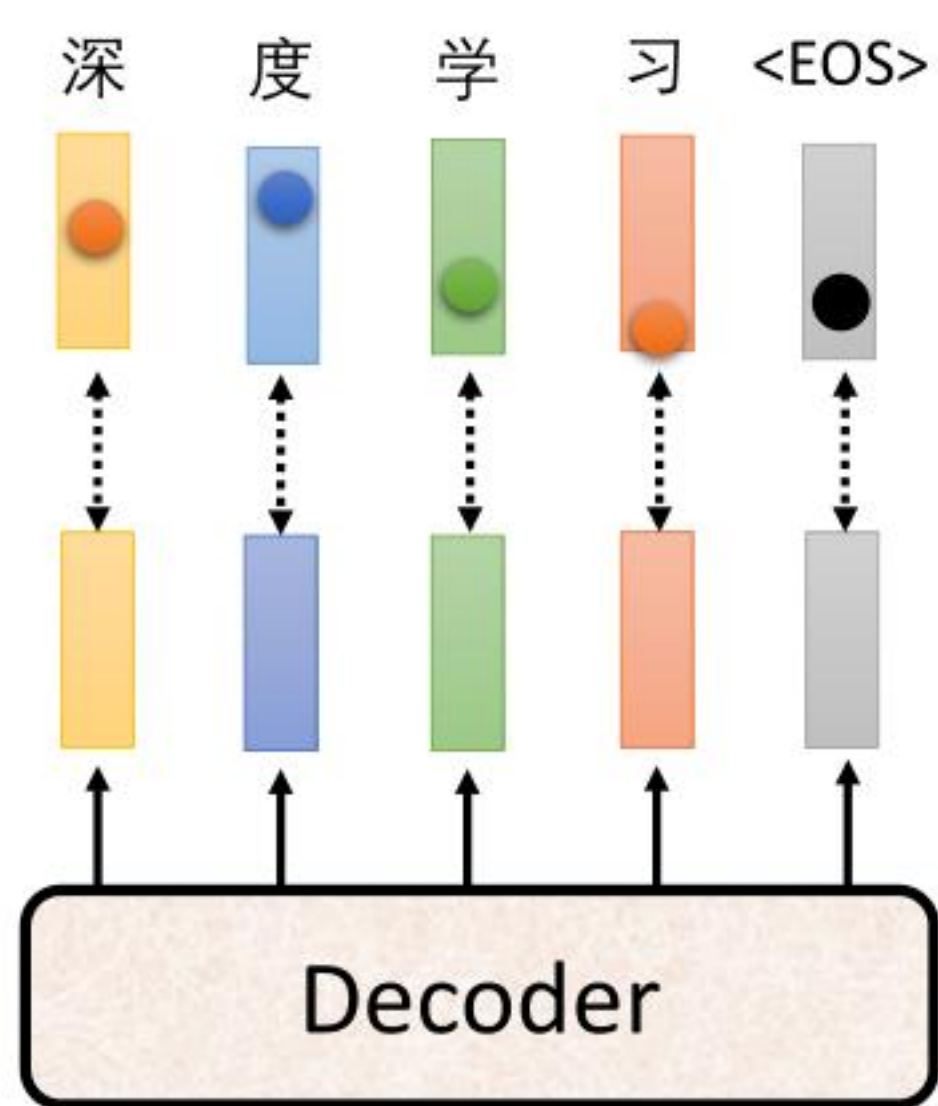
**Pure Sampling**:
They were cattle called Bolivian Cavalleros; they live in a remote desert uninterrupted by town, and they speak huge, beautiful, paradisiacal Bolivian linguistic thing. They say, 'Lunch, marge.' They don't tell what the lunch is," director Professor Chuperas Omwell told Sky News. "They've only been talking to scientists, like we're being interviewed by TV reporters. We don't even stick around to be interviewed by TV reporters. Maybe that's how they figured out that they're cosplaying as the Bolivian Cavalleros."

Randomness is needed for decoder when generating sequence in some tasks.

Accept that nothing is perfect. True beauty lies in the cracks of imperfection. ☺

# Optimizing Evaluation Metrics?



How to do the optimization?

When you don't know how to optimize, just use
reinforcement learning (RL)!   https://arxiv.org/abs/1511.06732

# Concluding Remarks: Transformer