



Transformer

李宏毅

Hung-yi Lee



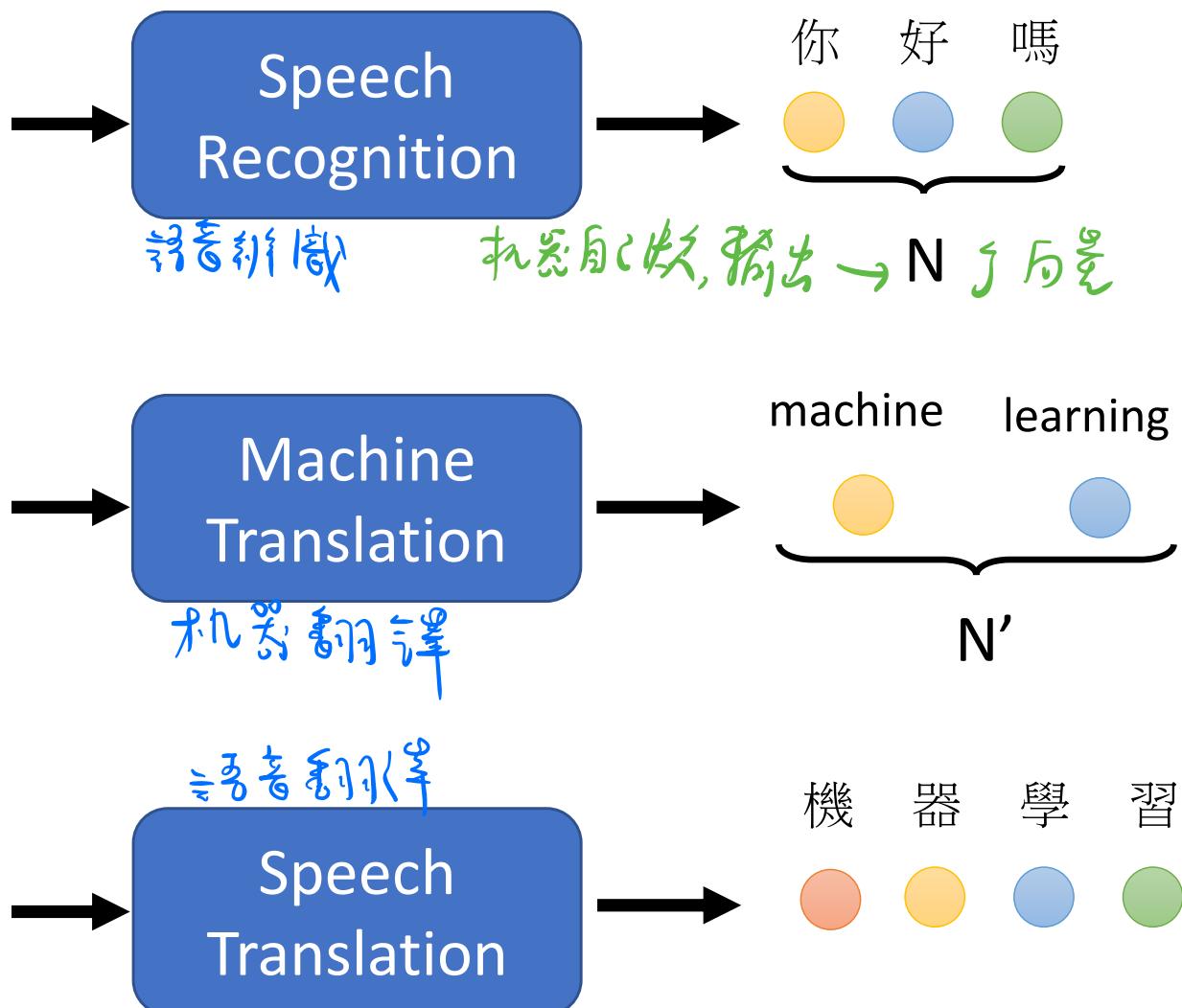
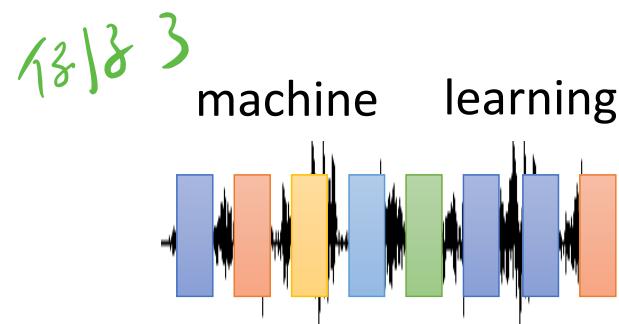
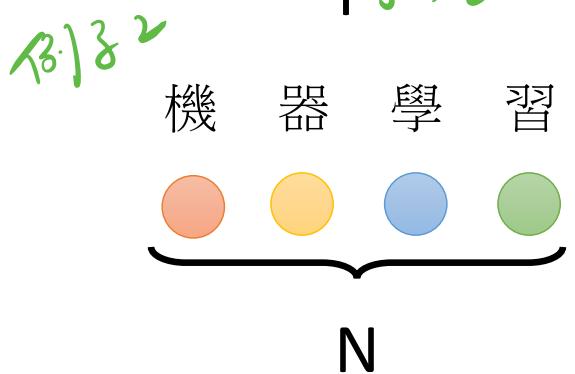
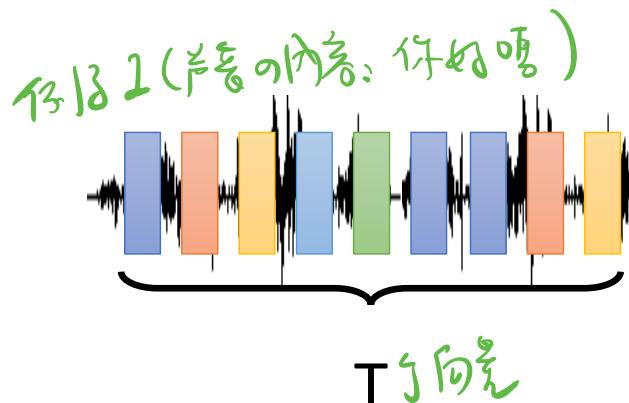
BERT

Sequence-to-sequence (Seq2seq)

transformer
Seq2Seq model

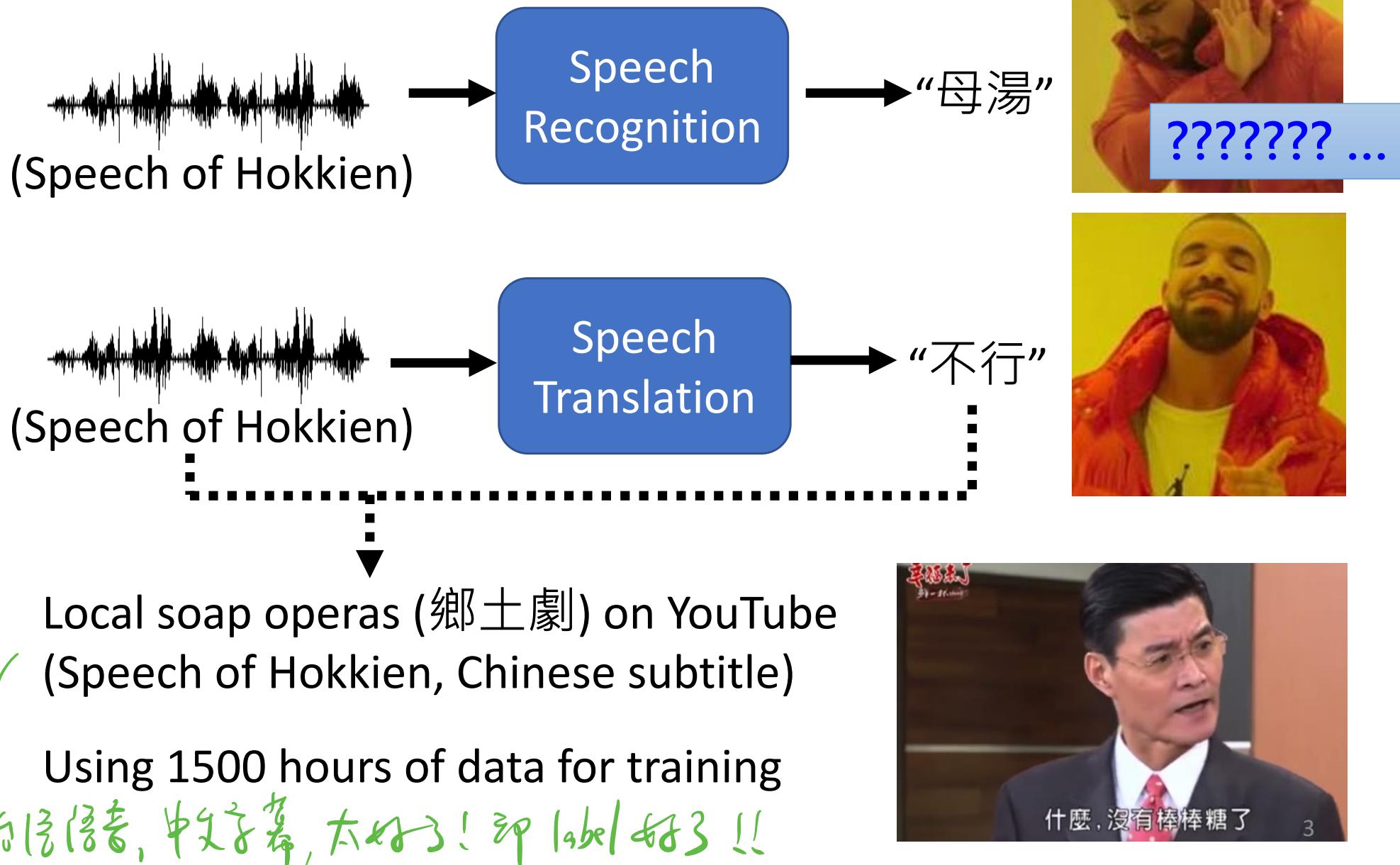
Input a sequence, output a sequence

The output length is determined by model.



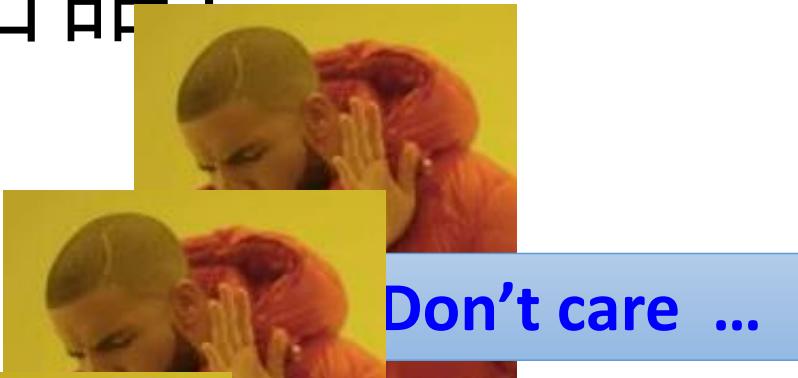
Language without text (e.g. 台语)

Hokkien (閩南語、台語)



Hokkien (閩南語、台語)

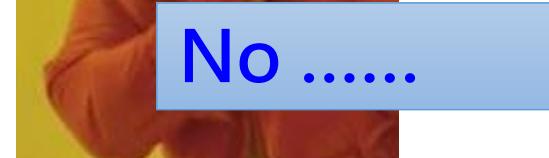
- Background music & noises?



- Noisy transcriptions?



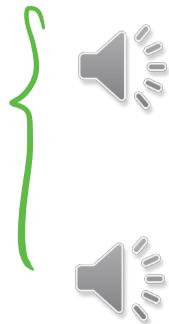
- Phonemes of Hokkien?



“硬train—發”
(Ying Train Yi Fa)

Hokkien (閩南語、台語)

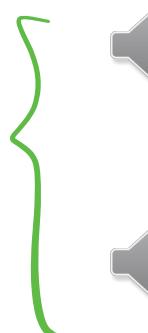
成功
case



你的身體撐不住

沒事你為什麼要請假

失敗
case



要生了嗎 Answer:不會膩嗎

我有幫廠長拜託

Answer: 我拜託廠長了

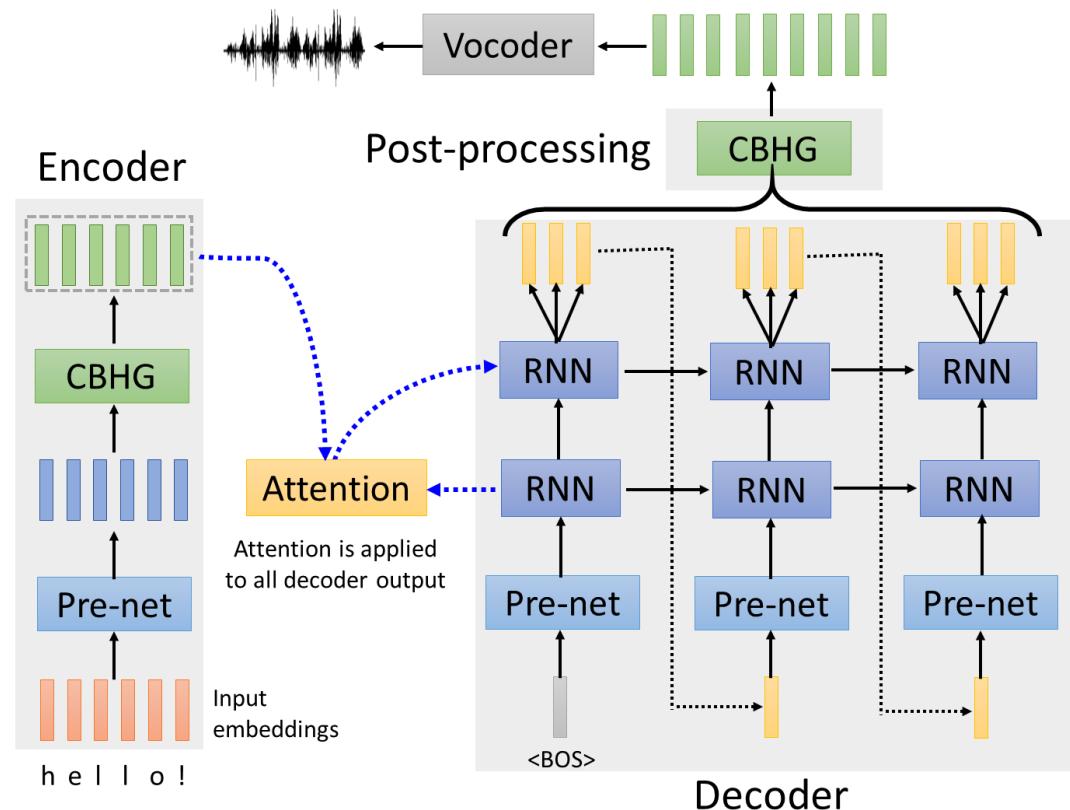
Text-to-Speech (TTS) Synthesis

(輸入文字，輸出音波圖
⇒ 語音合成)

Taiwanese Speech Synthesis

Source of data: 台灣婧聲2.0

感謝張凱為同學提供實驗結果



歡迎來到台大語言處理實驗室

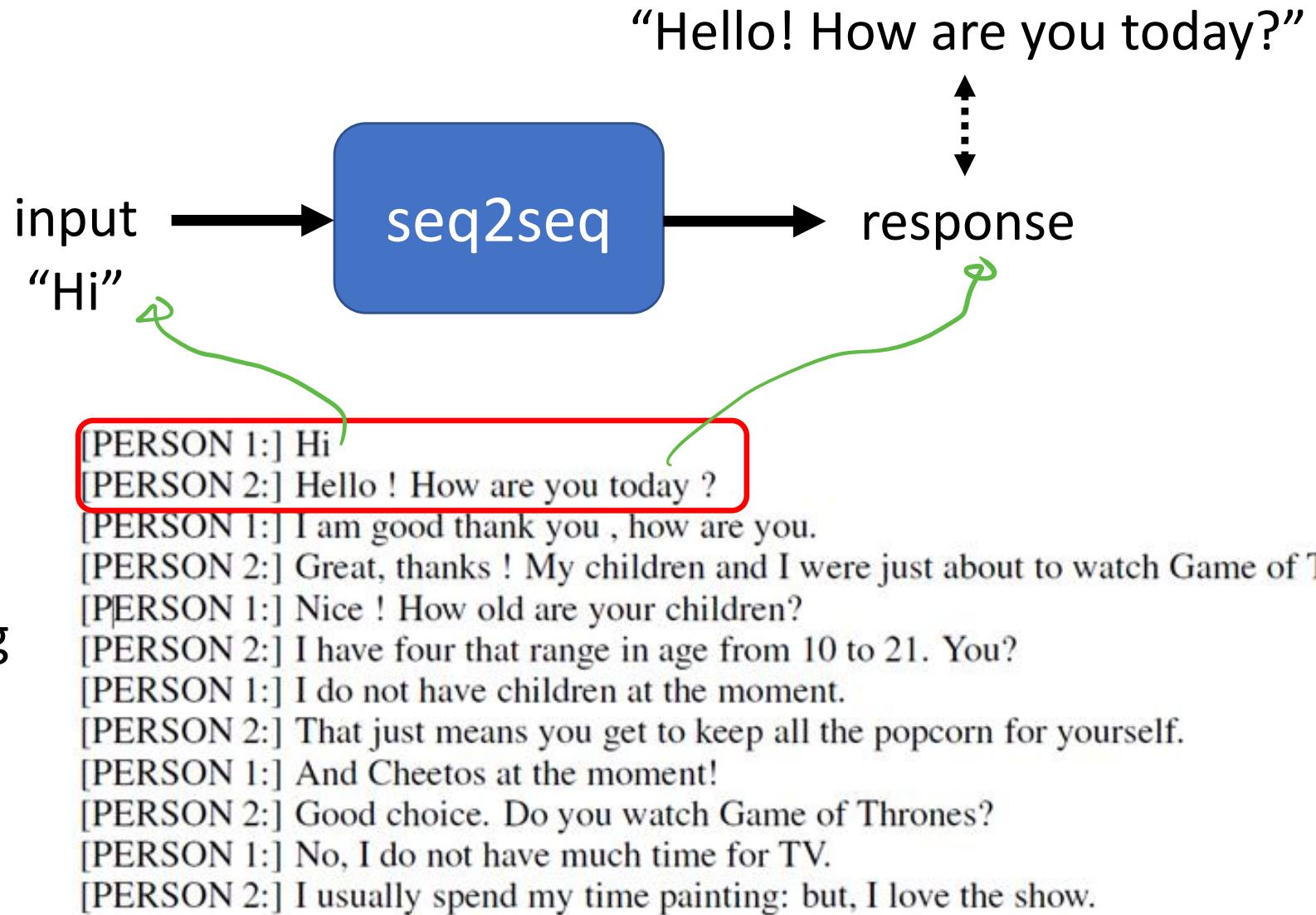


最近肺炎真嚴重，要記得戴口罩、
勤洗手，有病就要看醫生



文字上的例子

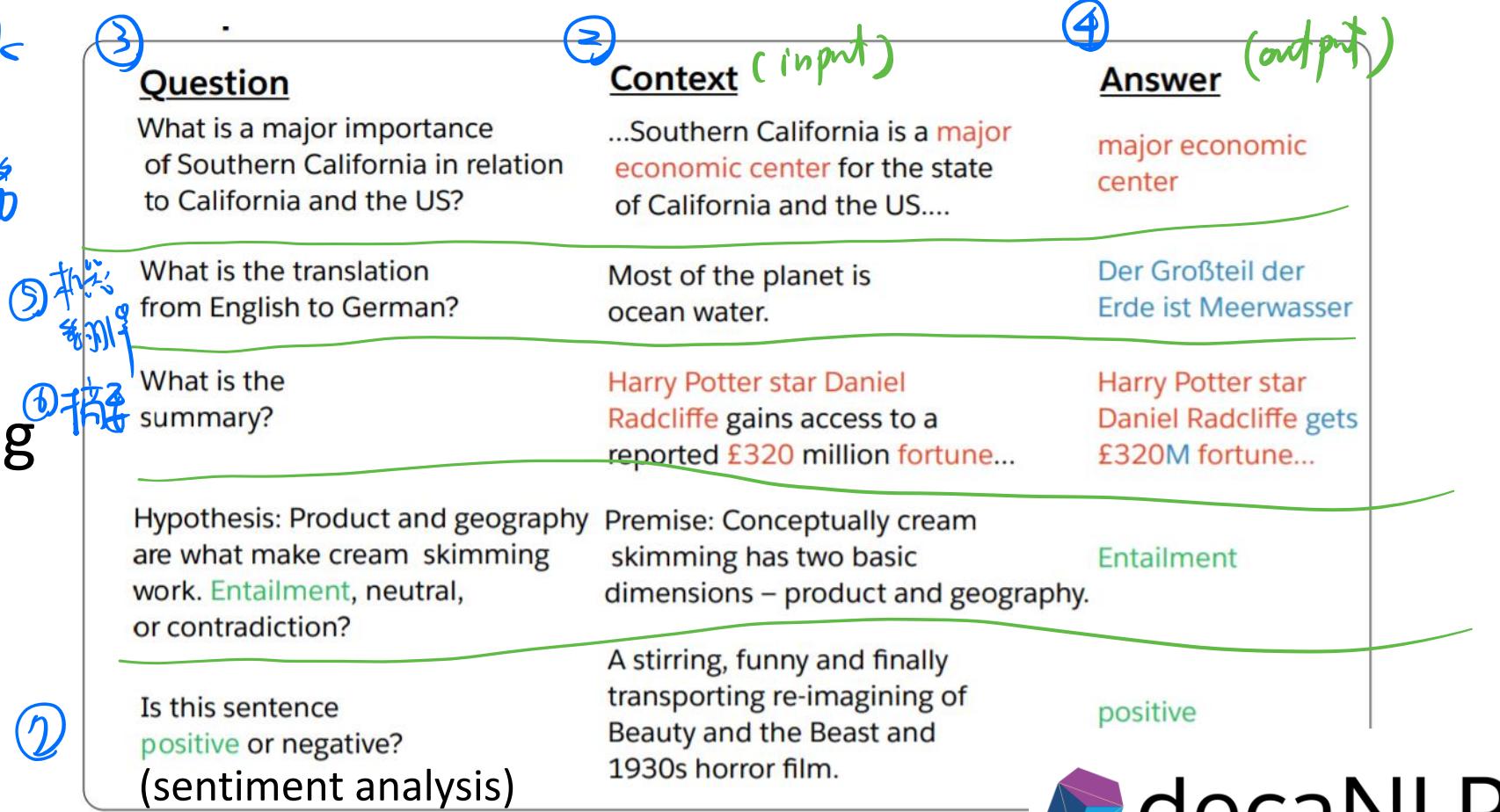
Seq2seq for Chatbot



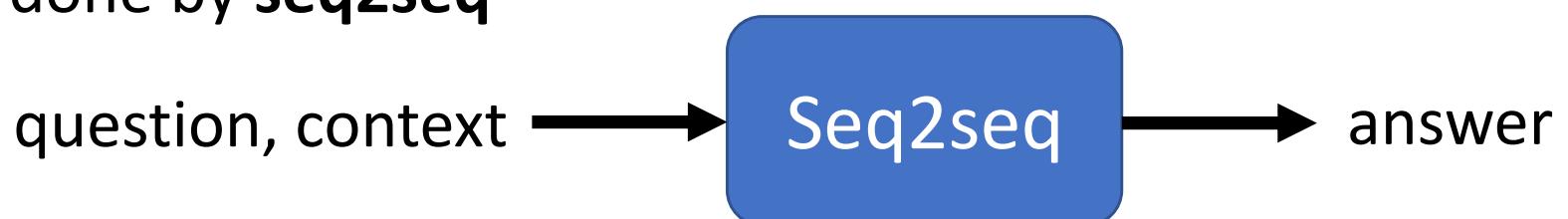
Most Natural Language Processing applications ...

① 很多 NLP 的 task
即 \rightarrow formulate
成 \rightarrow 的任務

Question Answering (QA)



QA can be done by seq2seq



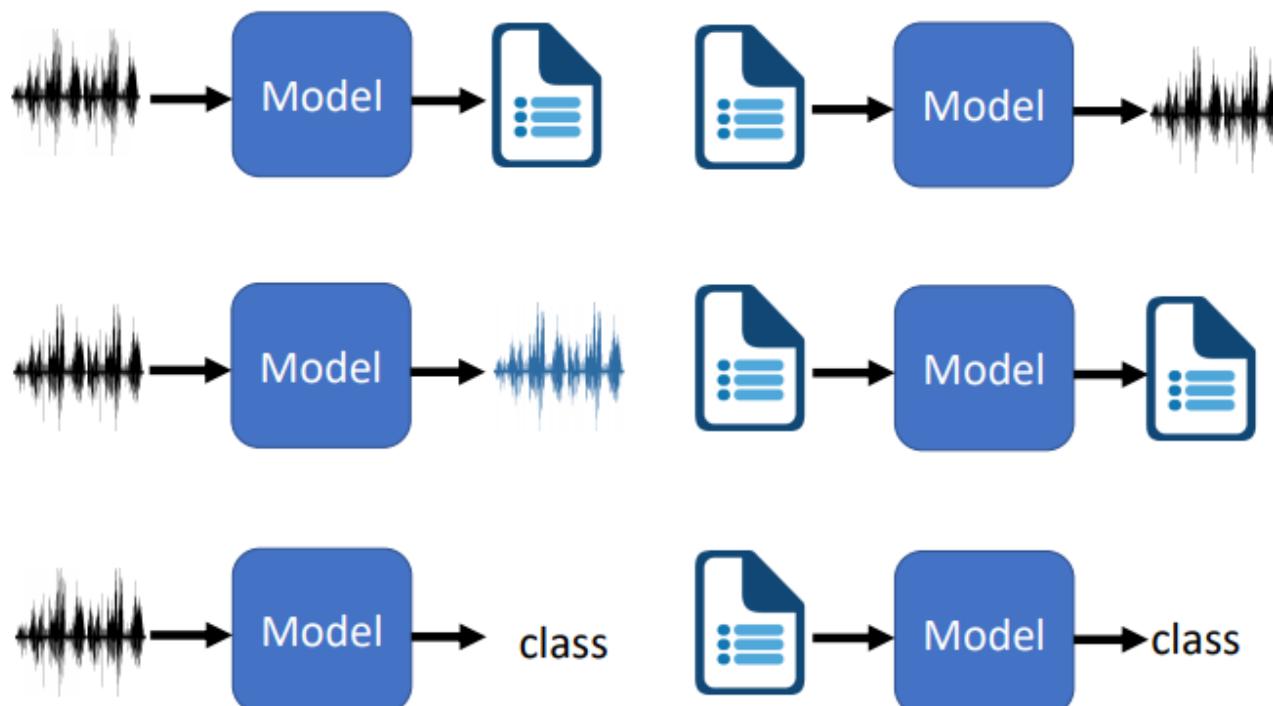
<https://arxiv.org/abs/1806.08730>
<https://arxiv.org/abs/1909.03329>

Deep Learning for Human Language Processing

深度學習與人類語言處理

↙
seg² seg
像端切
但對應付
很多 task
但如是想
要往做好
某種 task,
通常就不用
seg² seg.
而在這門課會詳談另外 task 下, 哪些的 model 是會

One slide for this course

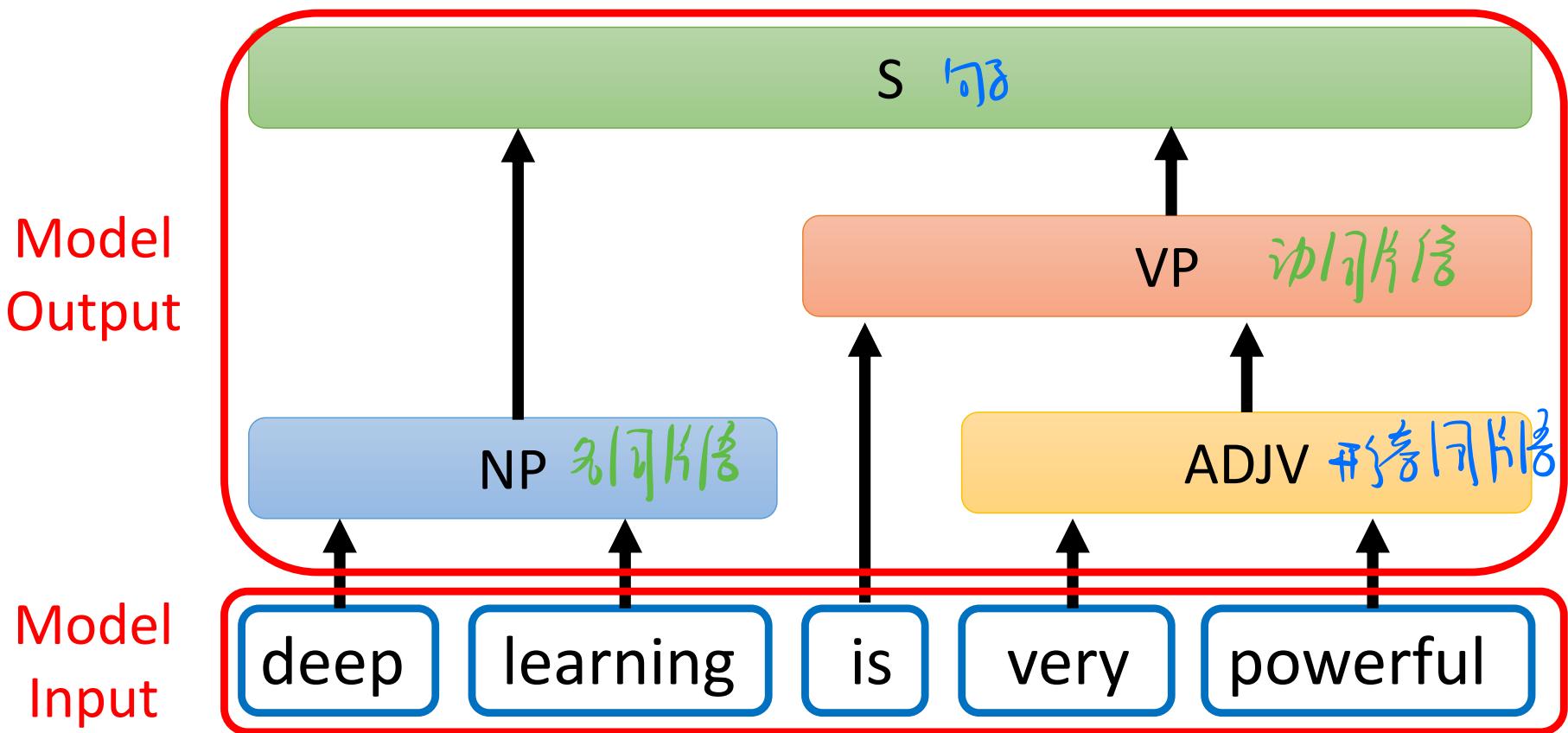


Source webpage: <https://speech.ee.ntu.edu.tw/~hylee/dlhlp/2020-spring.html>

Seq2seq for Syntactic Parsing

文法解析 \Rightarrow Seq2seq 解析

Is it a sequence?

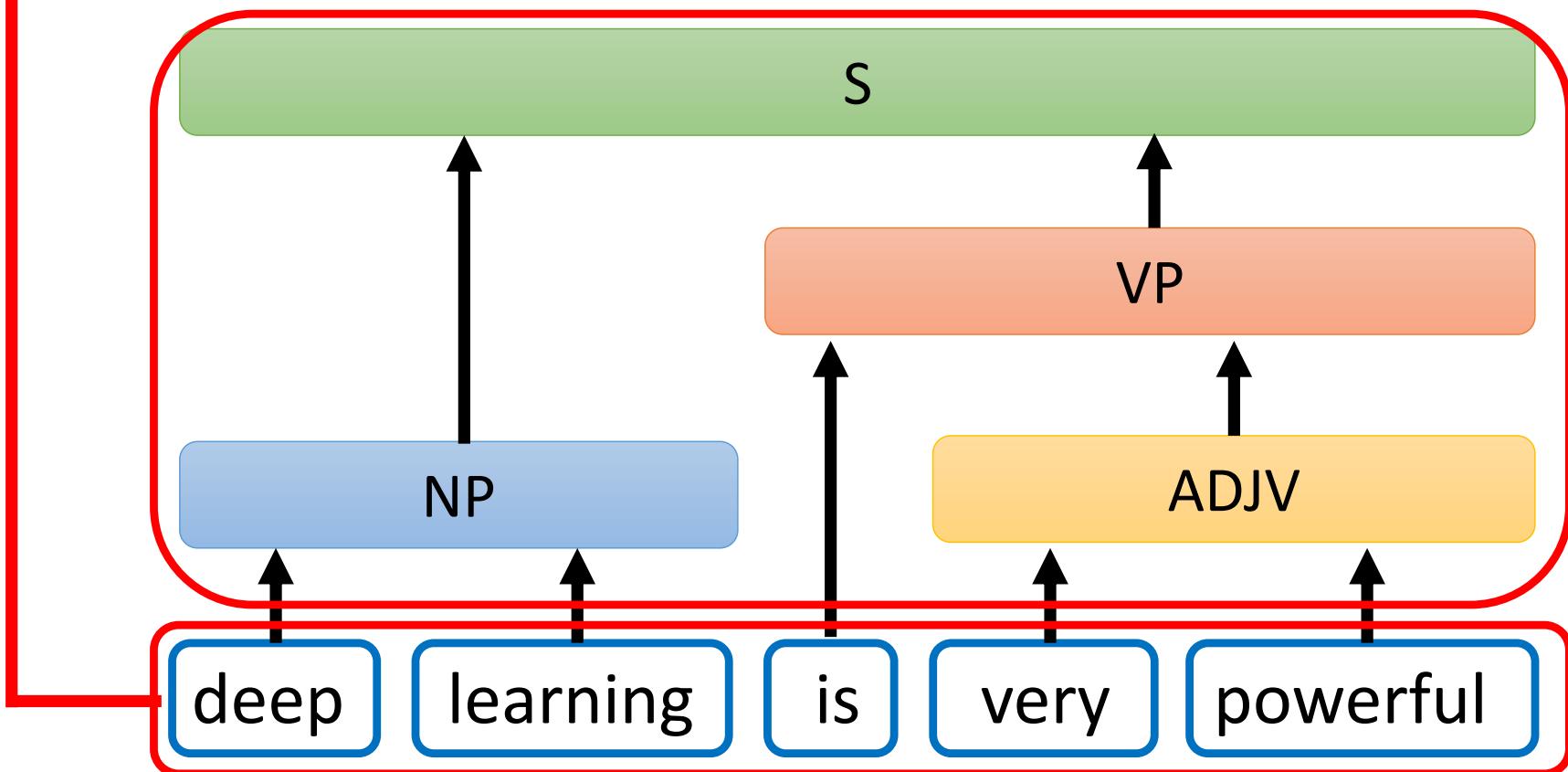


Seq2seq for Syntactic Parsing

↓ parsing for tree

(S (NP deep learning)) (VP is
(ADJV very powerful)))

Seq2seq!



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

c.f. Multi-class Classification

An object can belong to multiple classes.



Class 1
Class 3



Class 1



Class 3
Class 9



Class 10

Class 17

输入一篇文本



Seq2seq



Class 9



Class 7



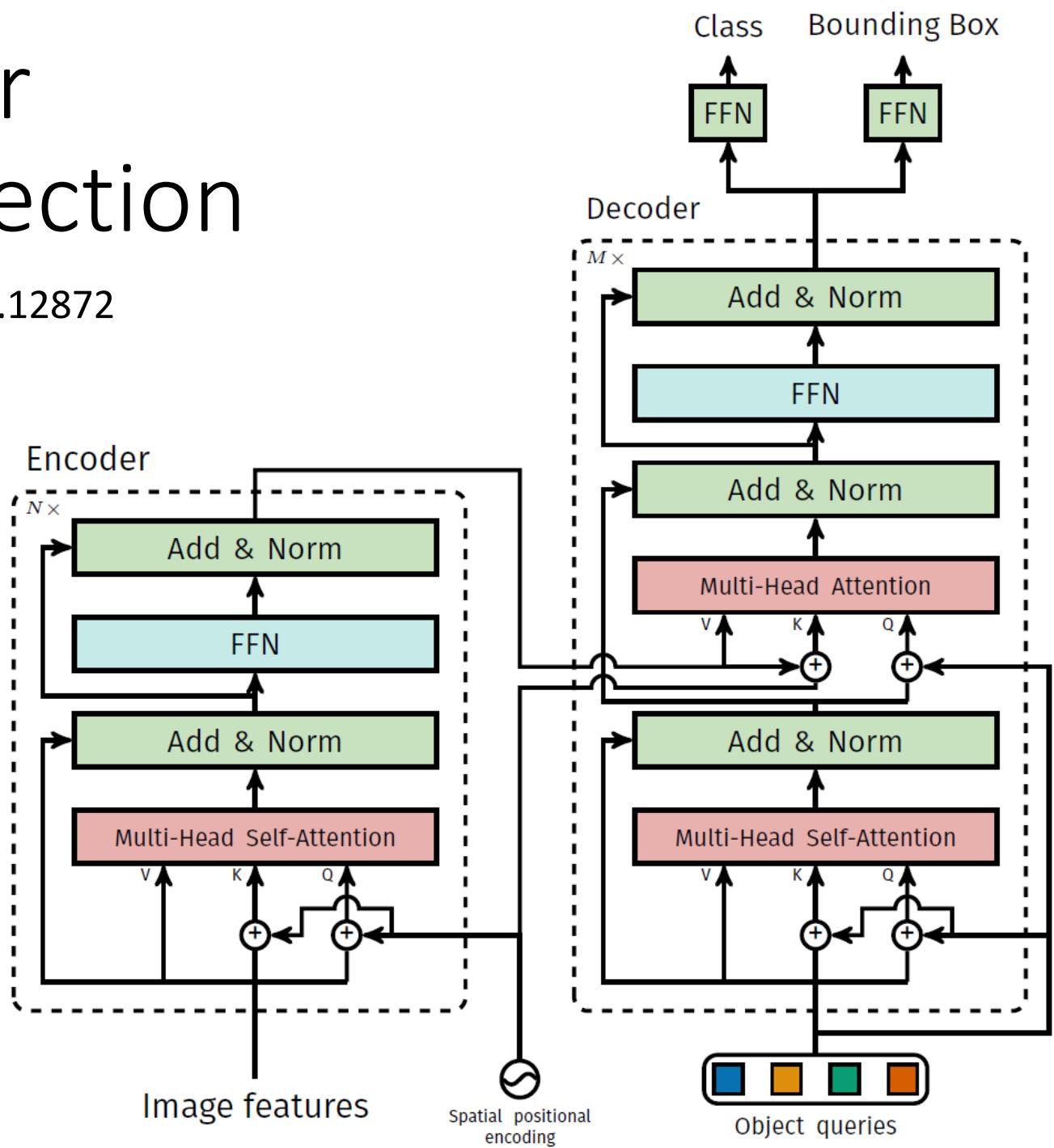
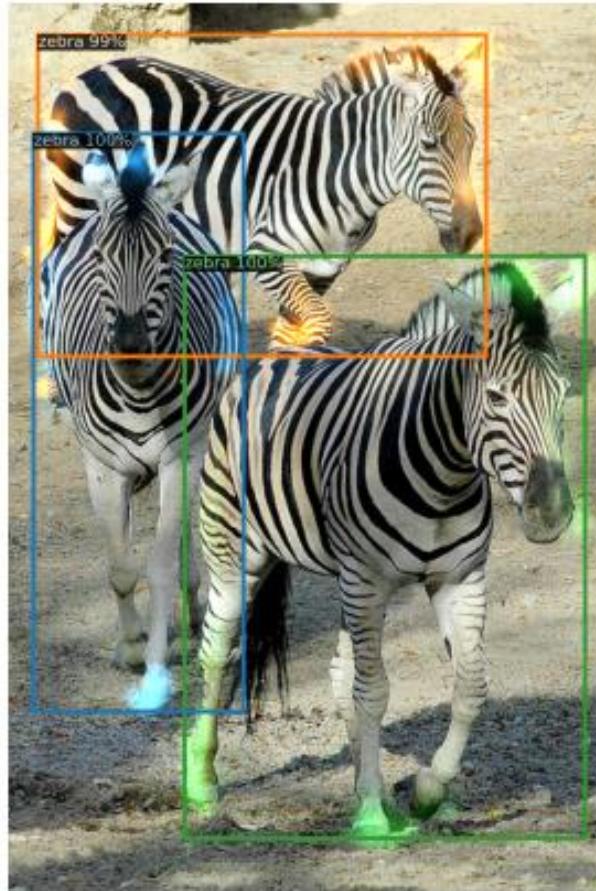
Class 13

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

model 自己决定, 然后输出哪几类

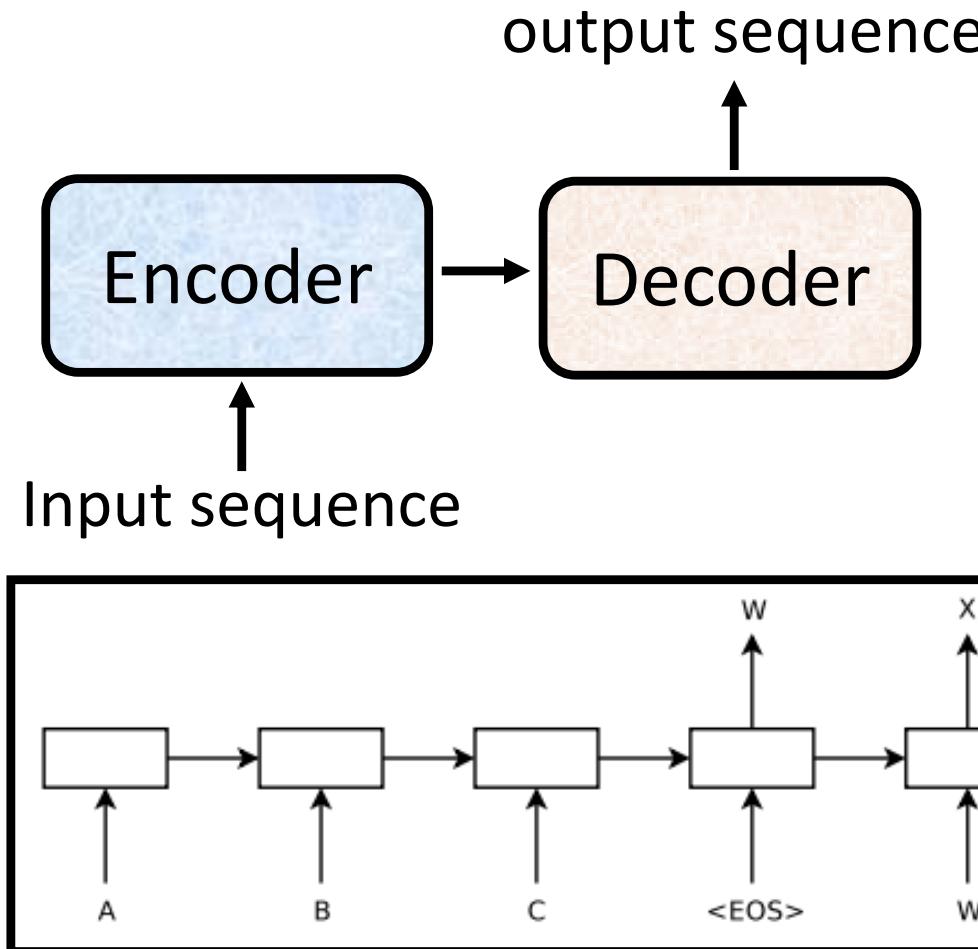
Seq2seq for Object Detection

<https://arxiv.org/abs/2005.12872>



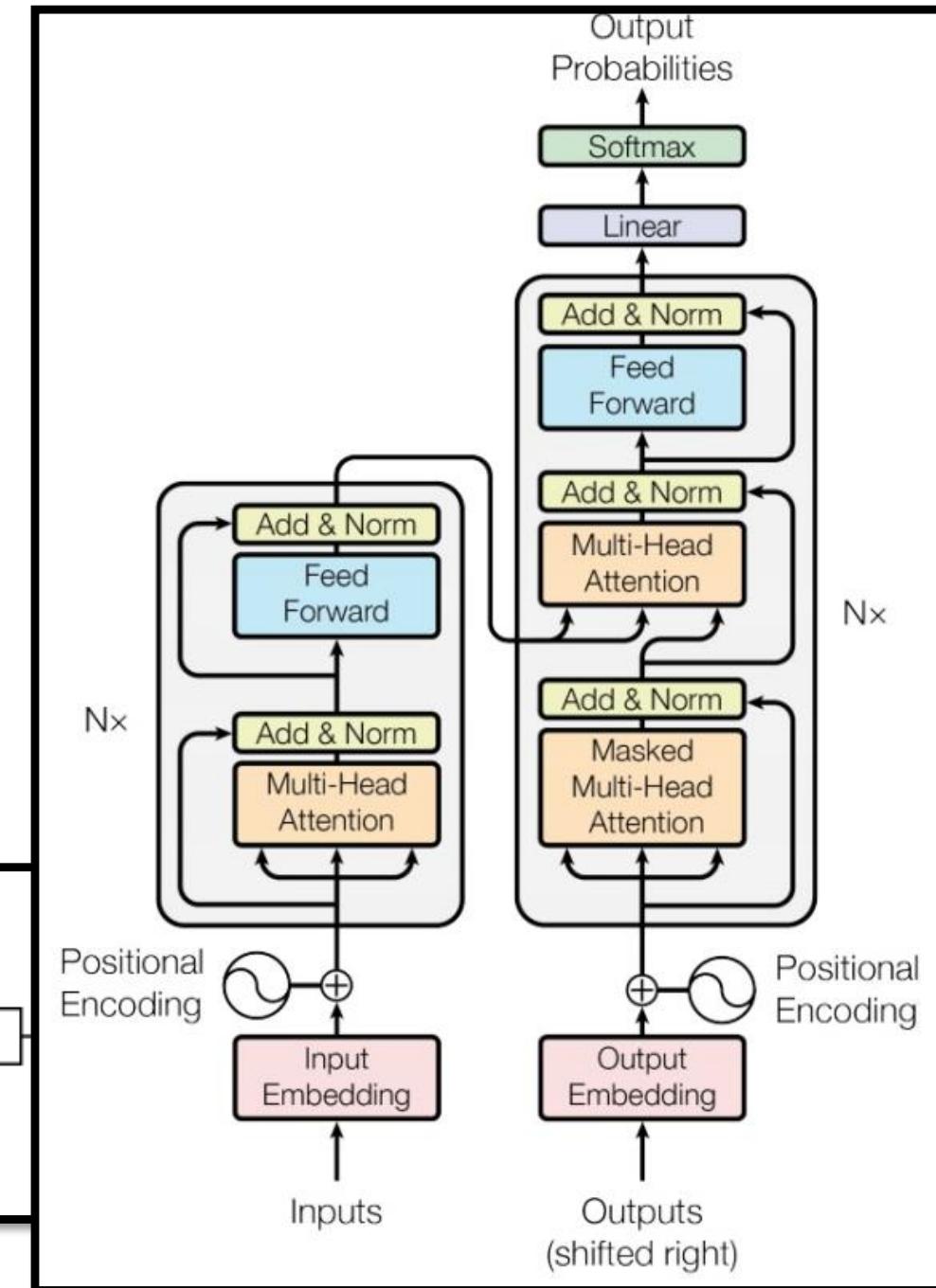
開始講序列模型 seq²seq

Seq2seq



Sequence to Sequence Learning with
Neural Networks

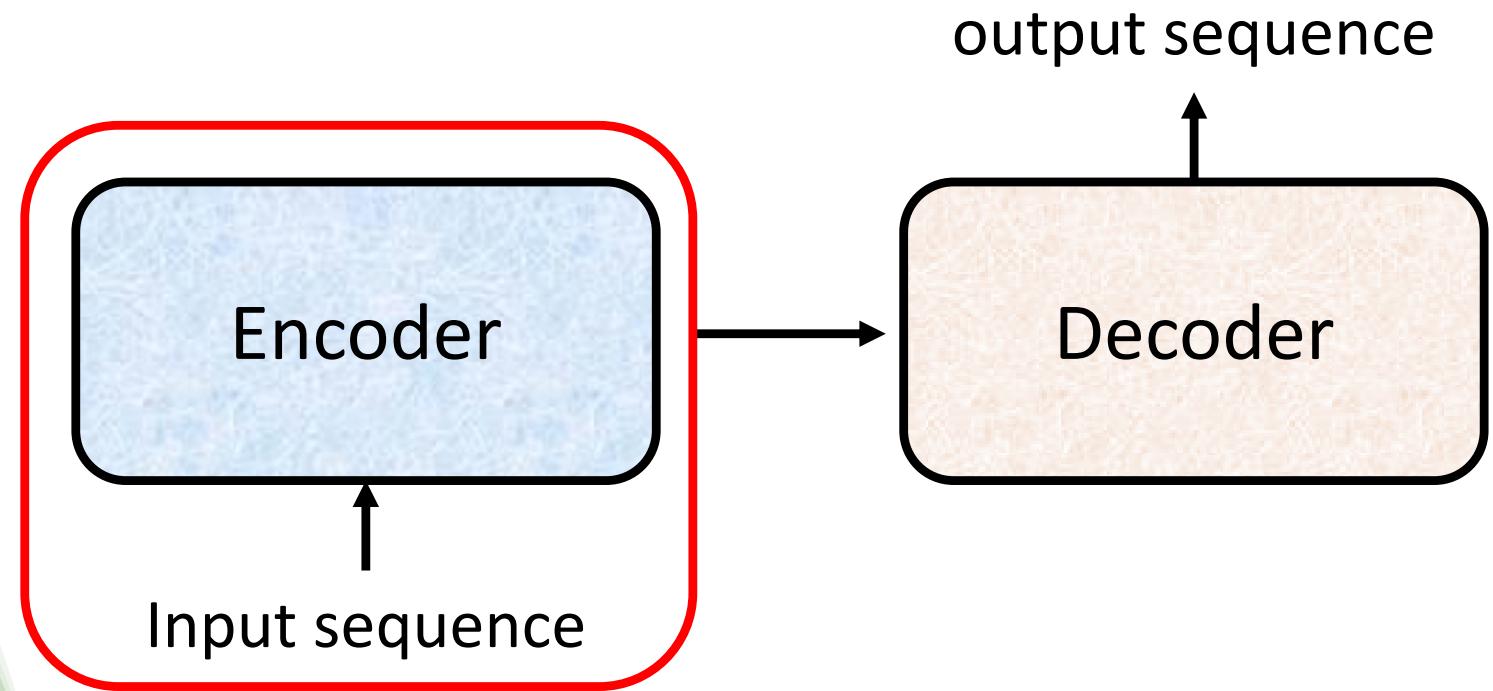
<https://arxiv.org/abs/1409.3215>



Transformer

<https://arxiv.org/abs/1706.03762>

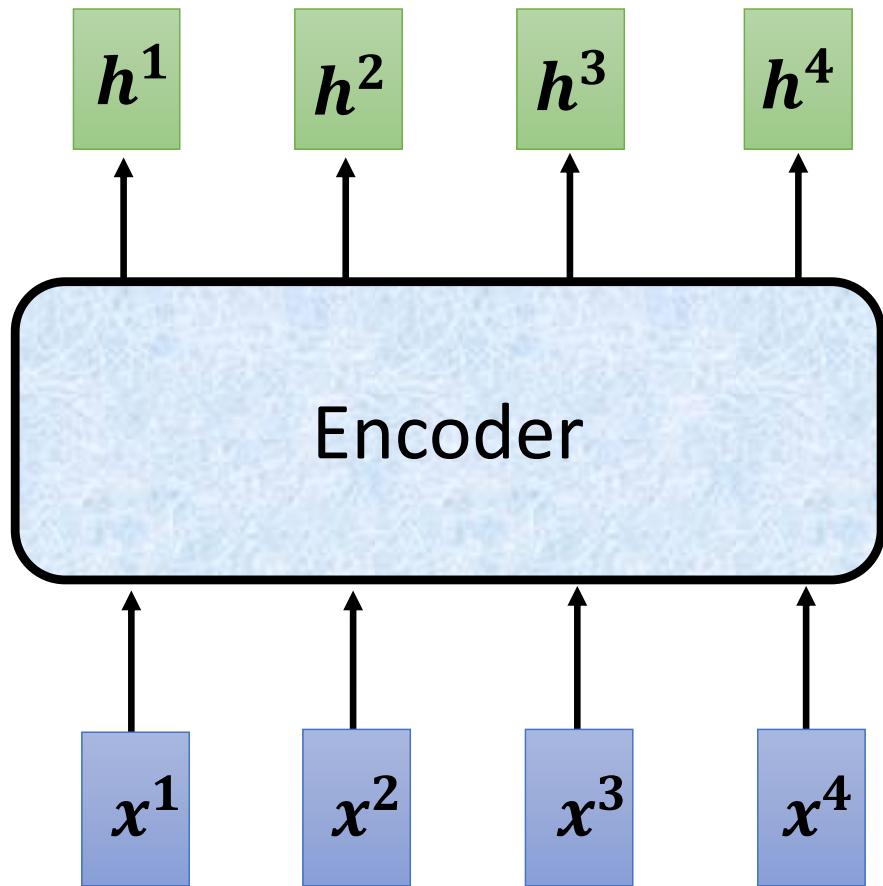
Encoder



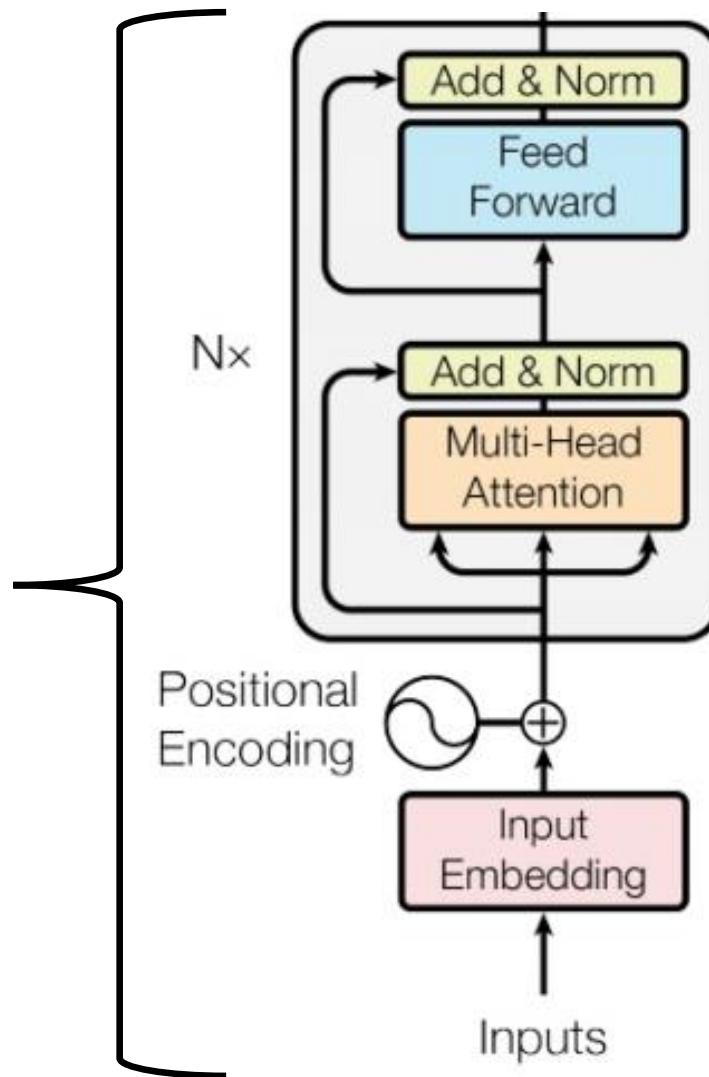
Encoder

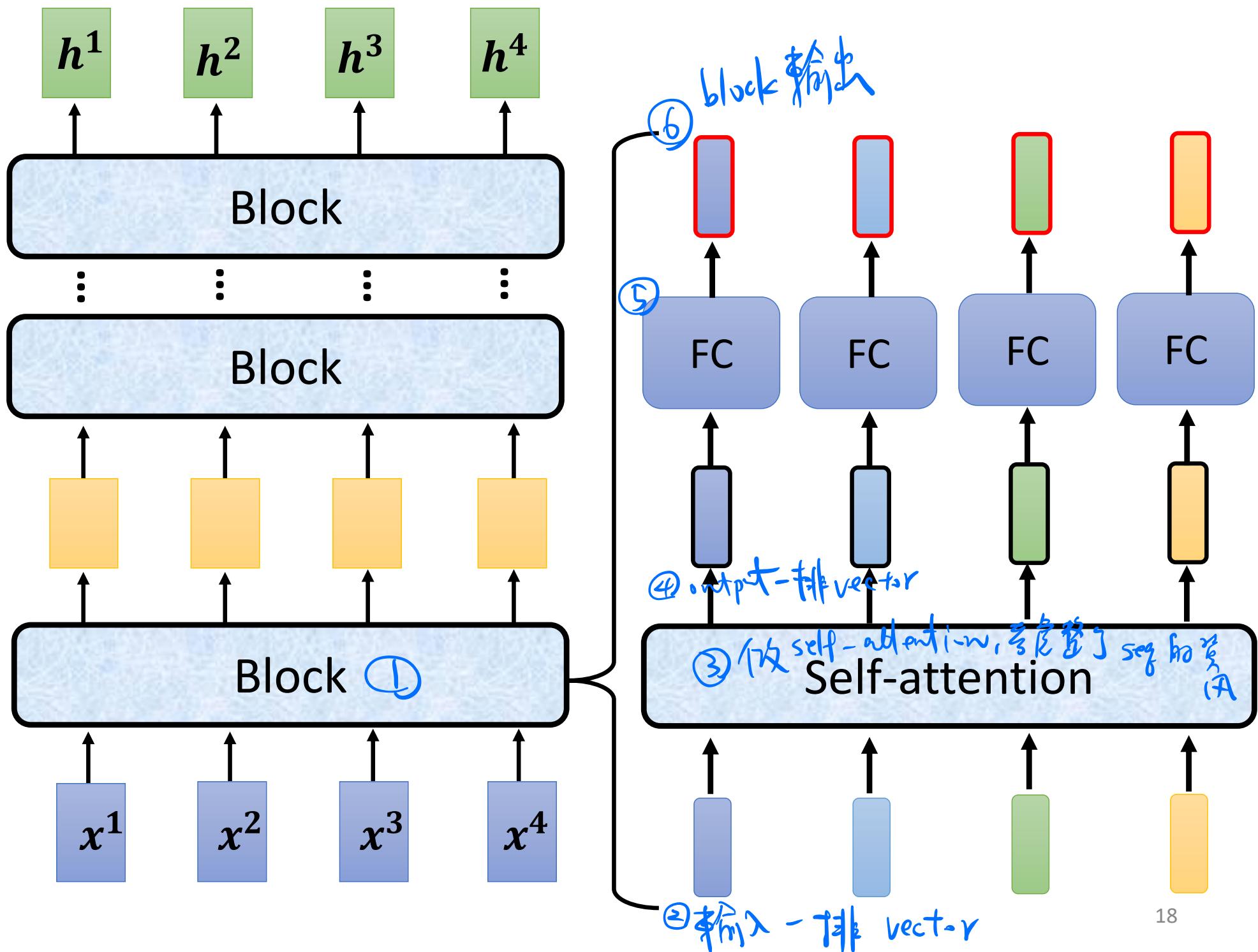
要做啥 $\Rightarrow \begin{cases} \text{input} \\ \text{output} \end{cases}$ N个向量 \Rightarrow 要做到这了，
N个向量 \Rightarrow 可用 RNN or CNN or
今日之角

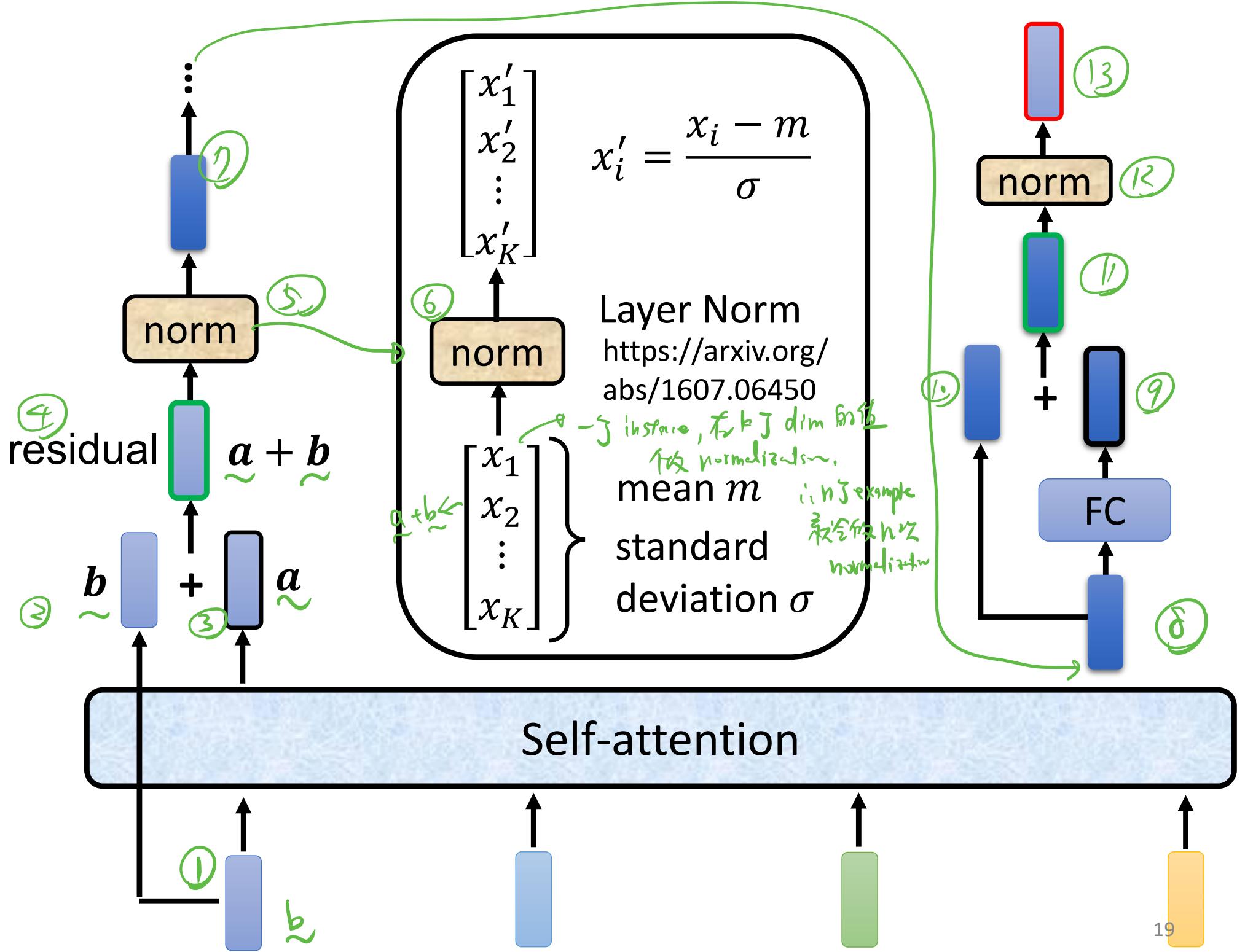
You can use **RNN** or **CNN**.



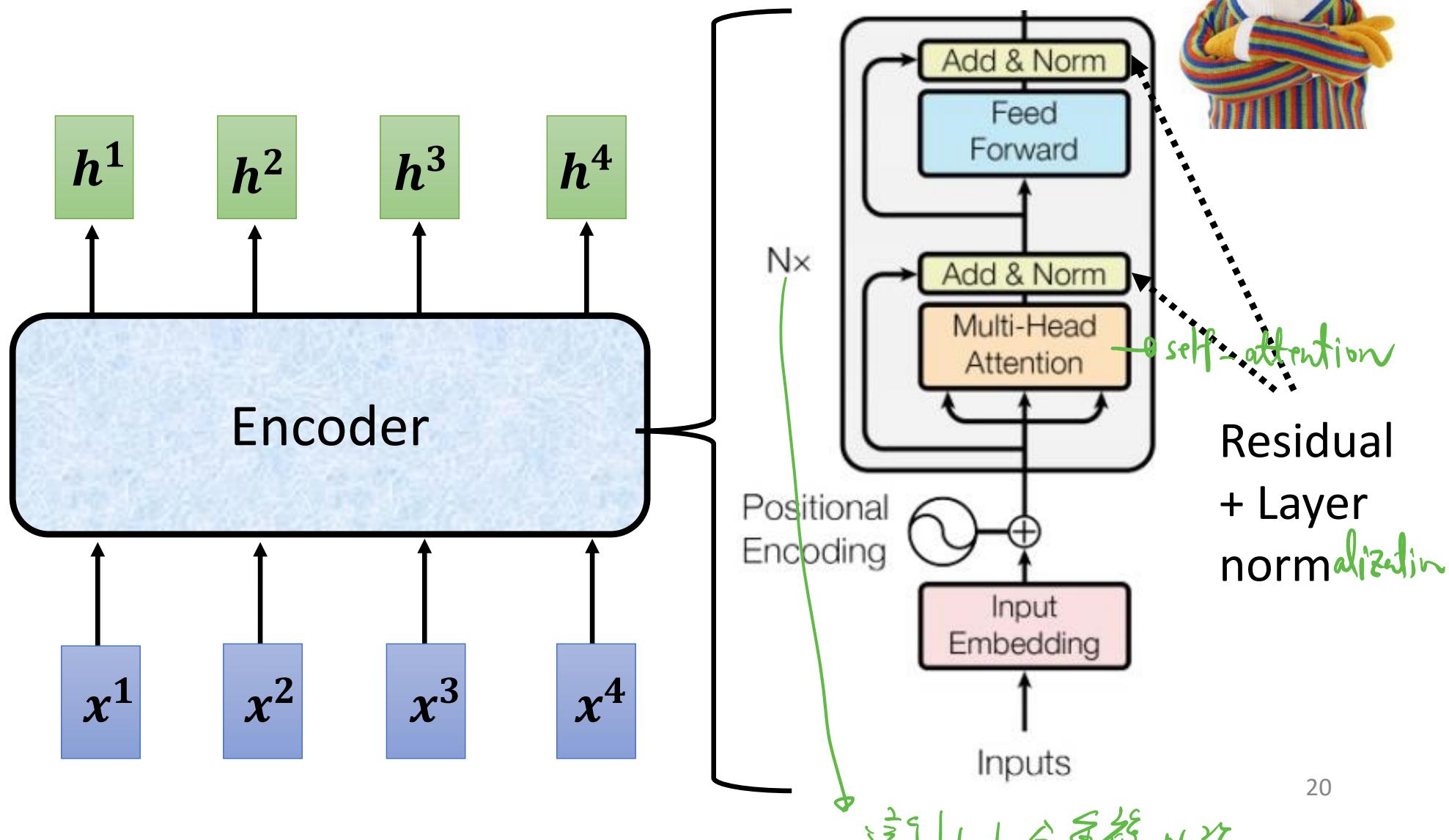
Transformer's Encoder







I use the **same** network architecture as transformer encoder.



20 block 有 273 N 人

To learn more

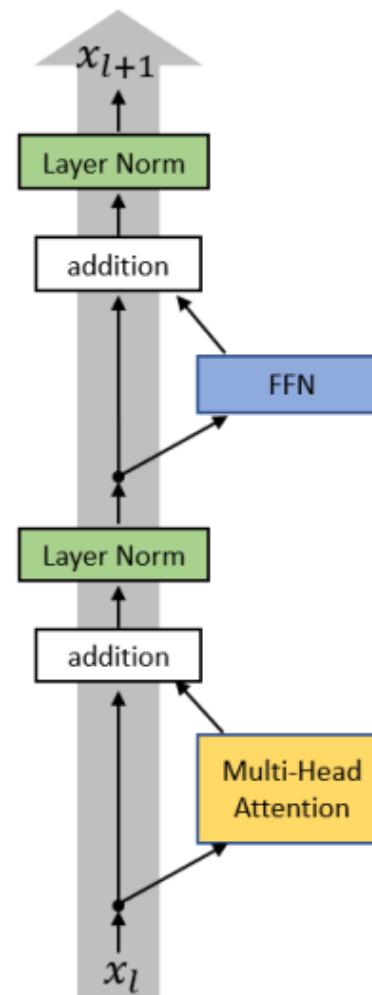
① $\Rightarrow Q = \frac{1}{\sqrt{d}}$ 像 transformer 的 encoder 是这样设计的？

- ②
- On Layer Normalization in the Transformer Architecture
 - <https://arxiv.org/abs/2002.04745>

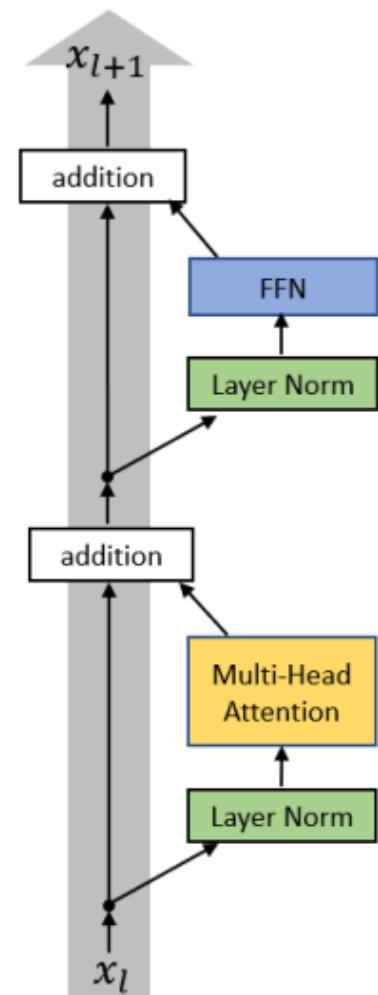
layer normalization 和直接连，有意义？

- ④
- PowerNorm: Rethinking Batch Normalization in Transformers
 - <https://arxiv.org/abs/2003.07845>

這篇 paper 說什麼，
在 transformer 中，layer norm
和 batch norm 重疊。



(a)



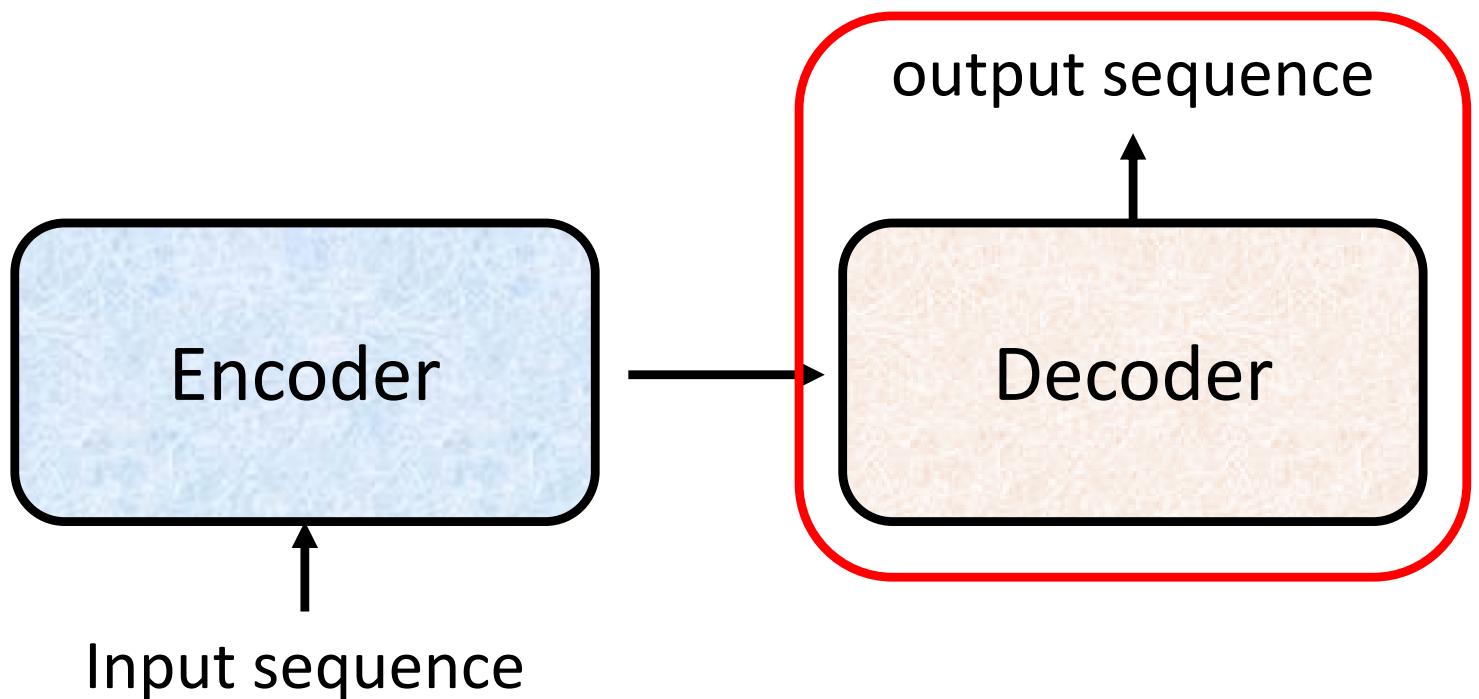
(b)

This transformer

③

修改，更直观
效果更好，
不必纠结于
序的设计

Decoder

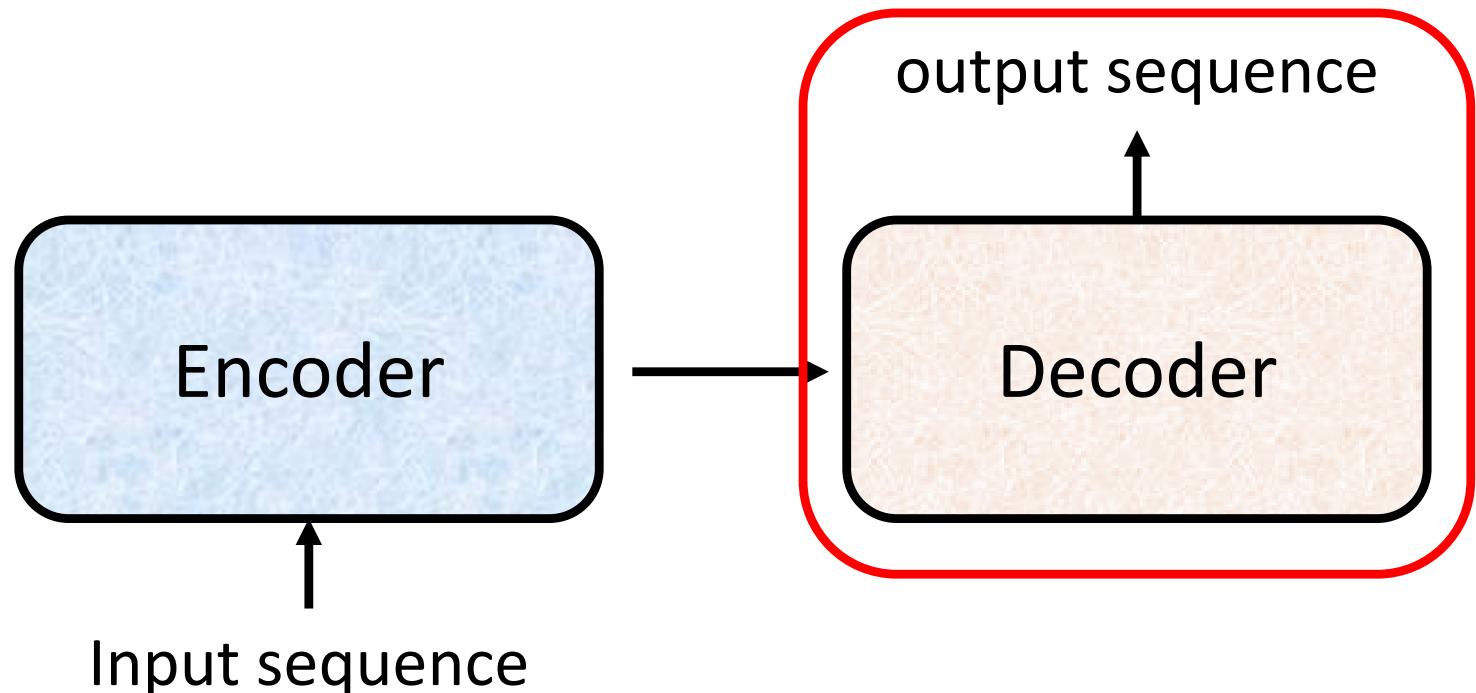


Decoder有2種

- └ Auto regressive
- └ Non-autoregressive

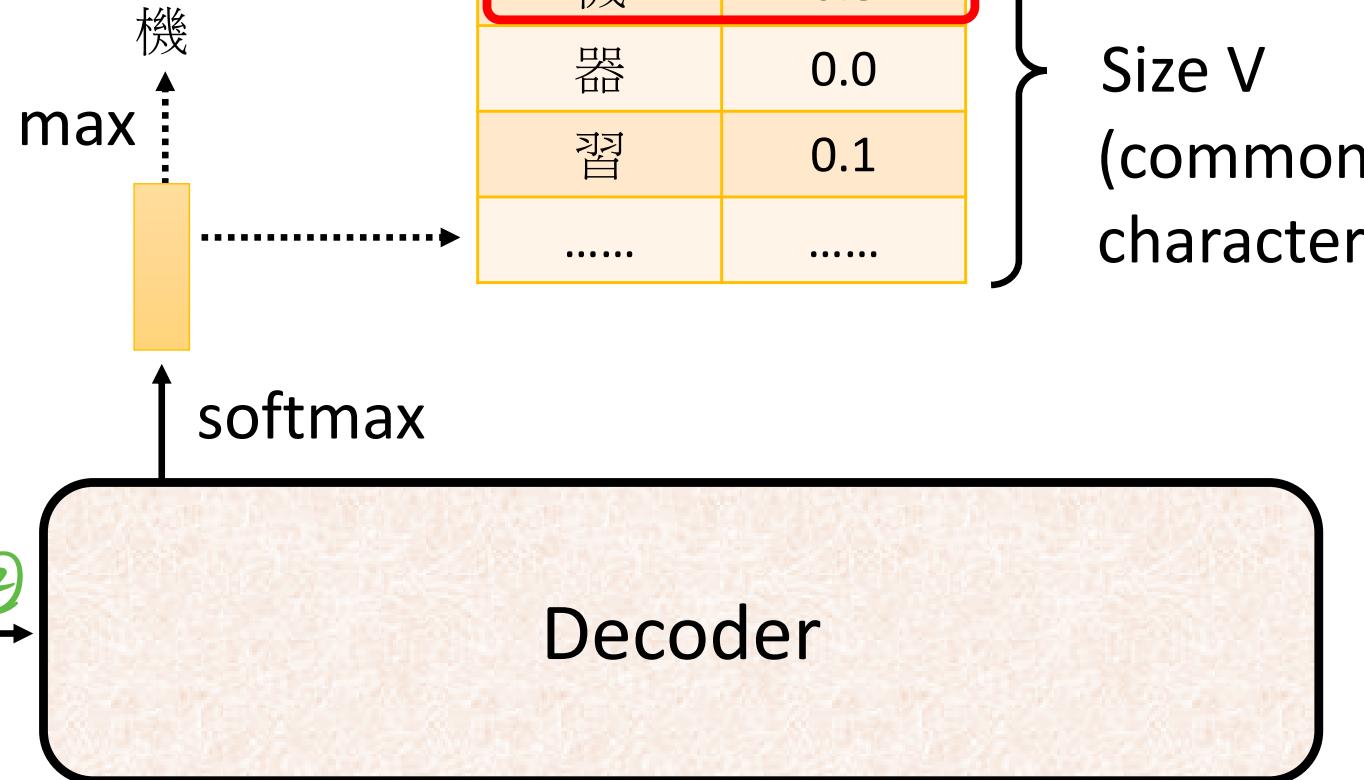
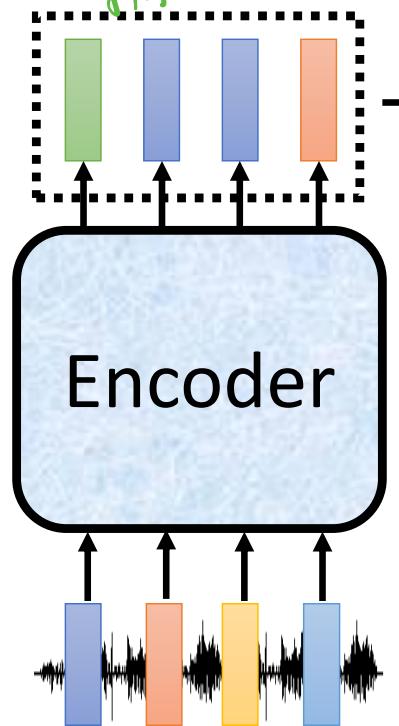
Decoder

– Autoregressive (AT)



Autoregressive (Speech Recognition as example)

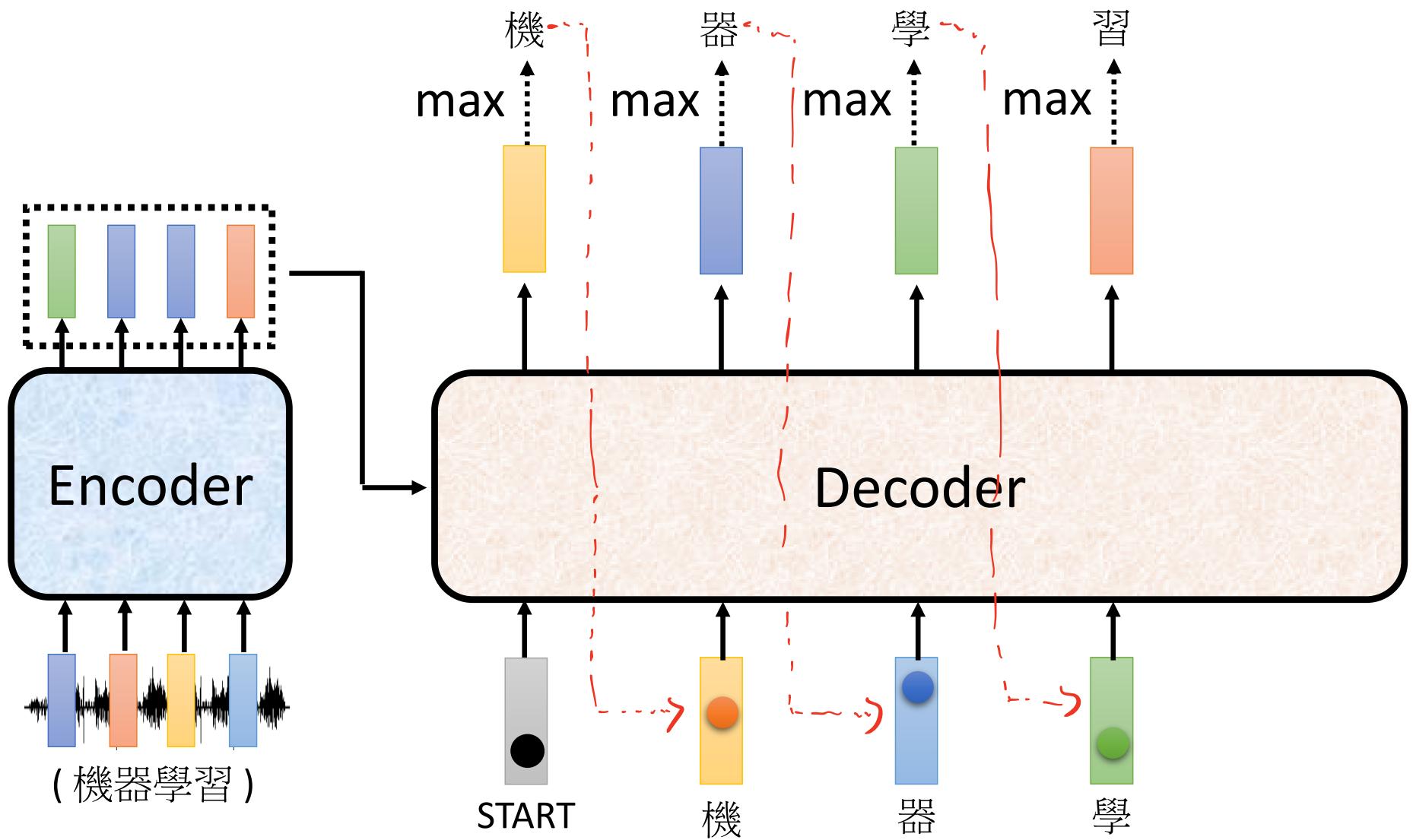
④ encoder 何の事：
 → 入力-非 vector
 → 演出-非 vector

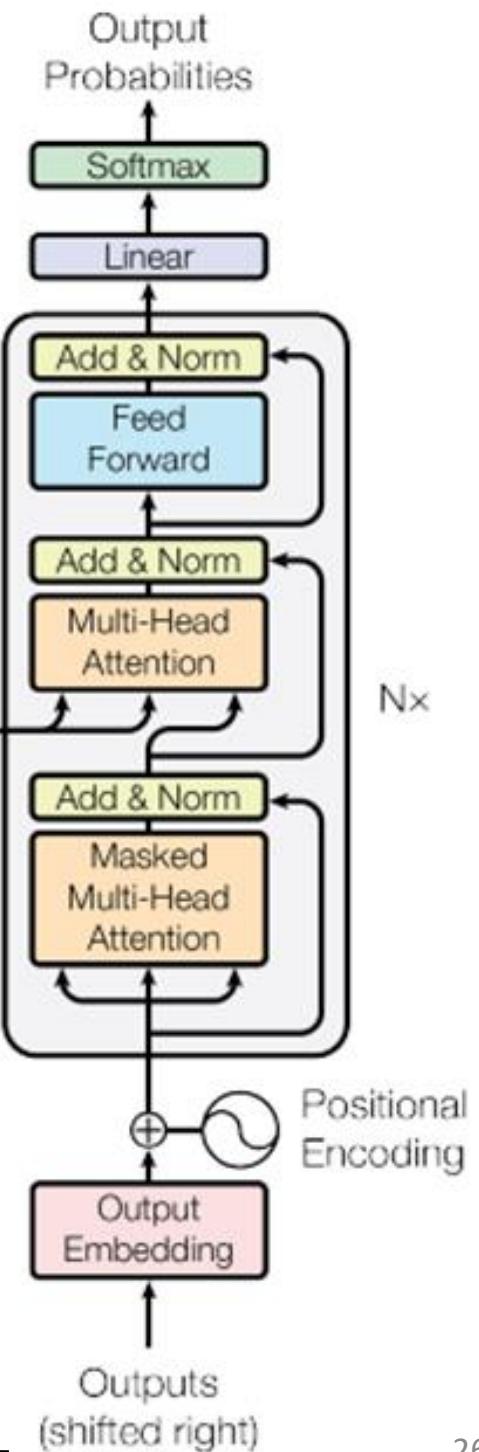
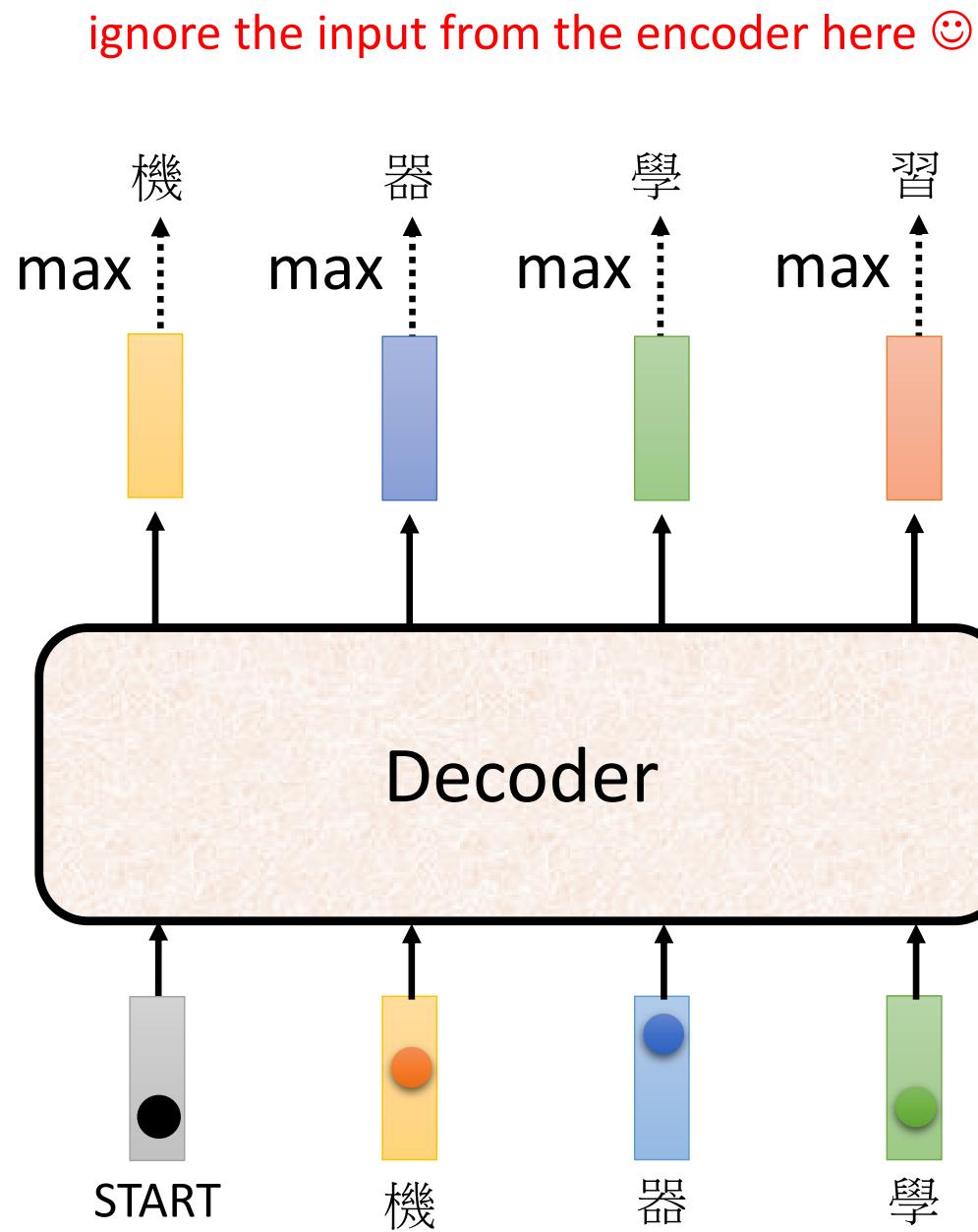


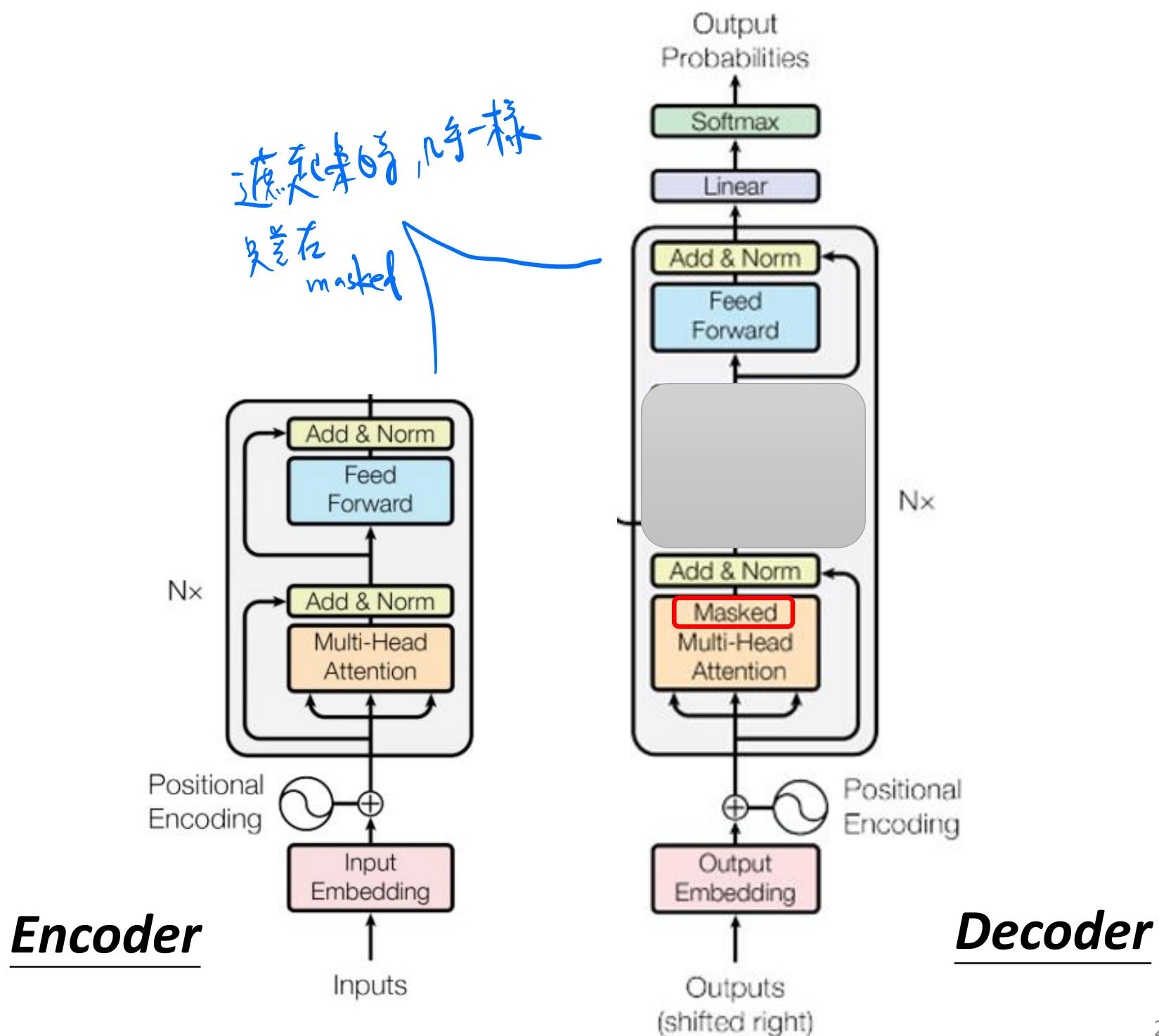
(機器學習)

BOS: Begin Of Sentence

Autoregressive



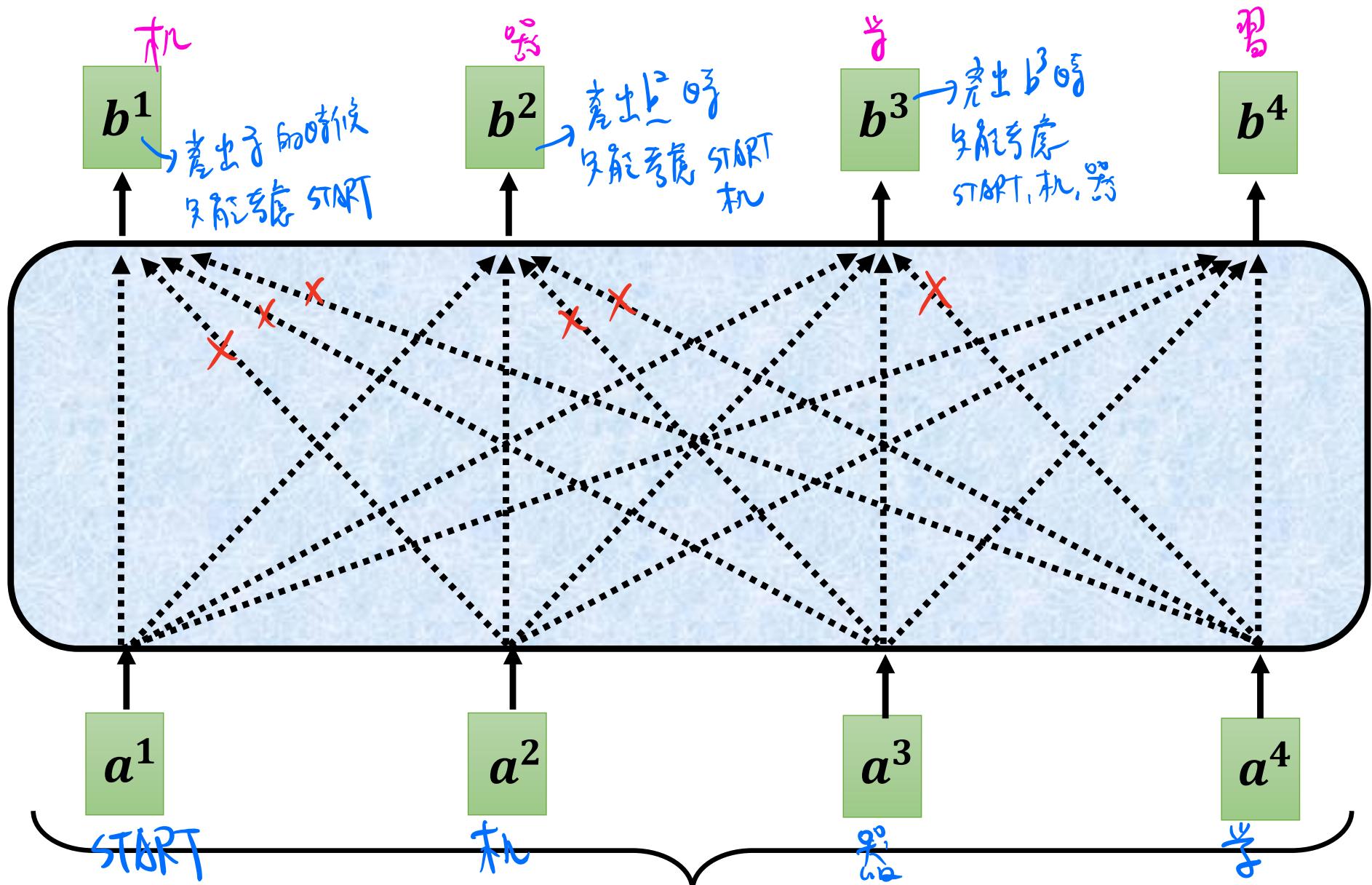




Encoder

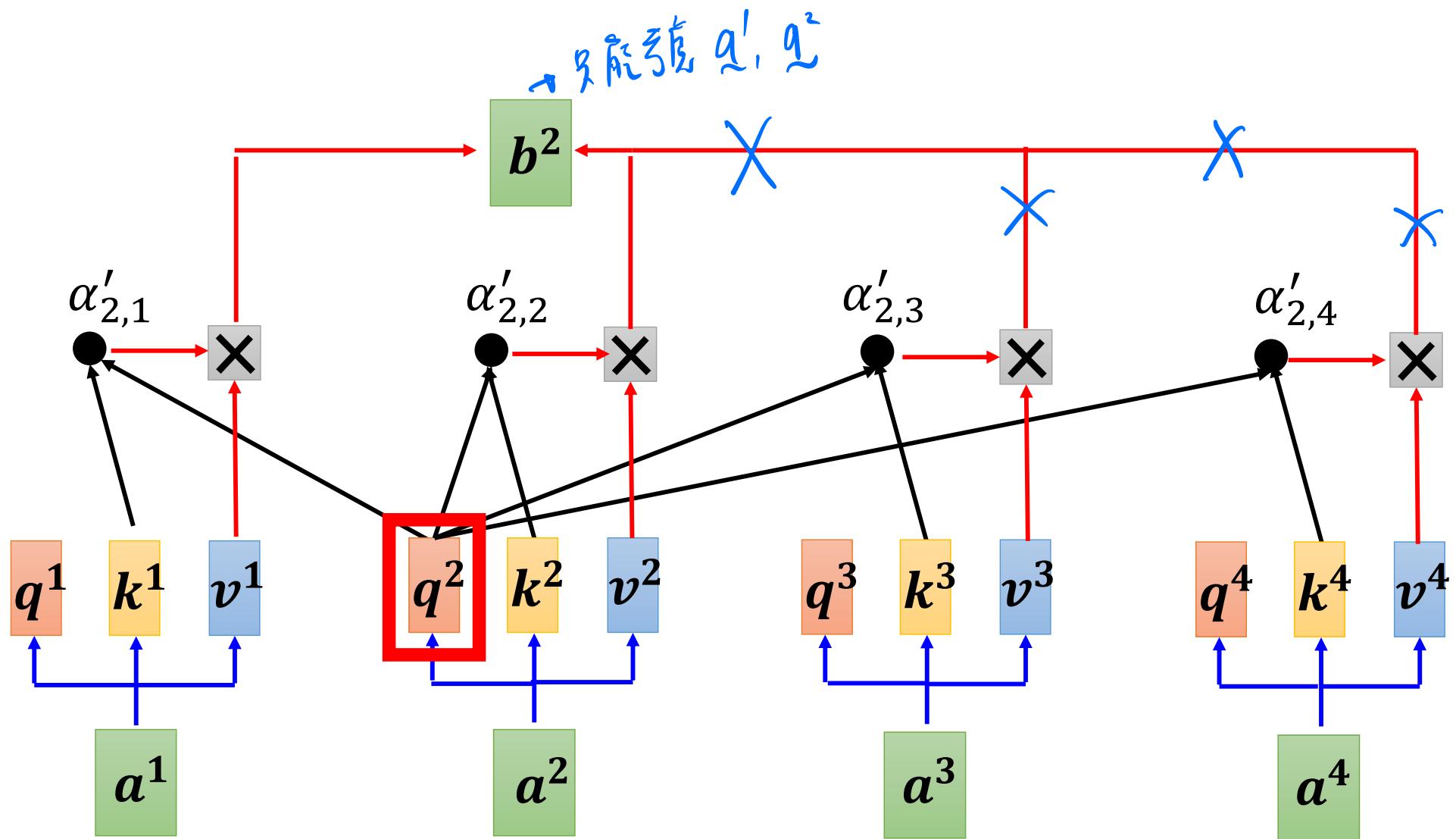
Decoder

Self-attention → *Masked Self-attention*



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

Self-attention → Masked Self-attention



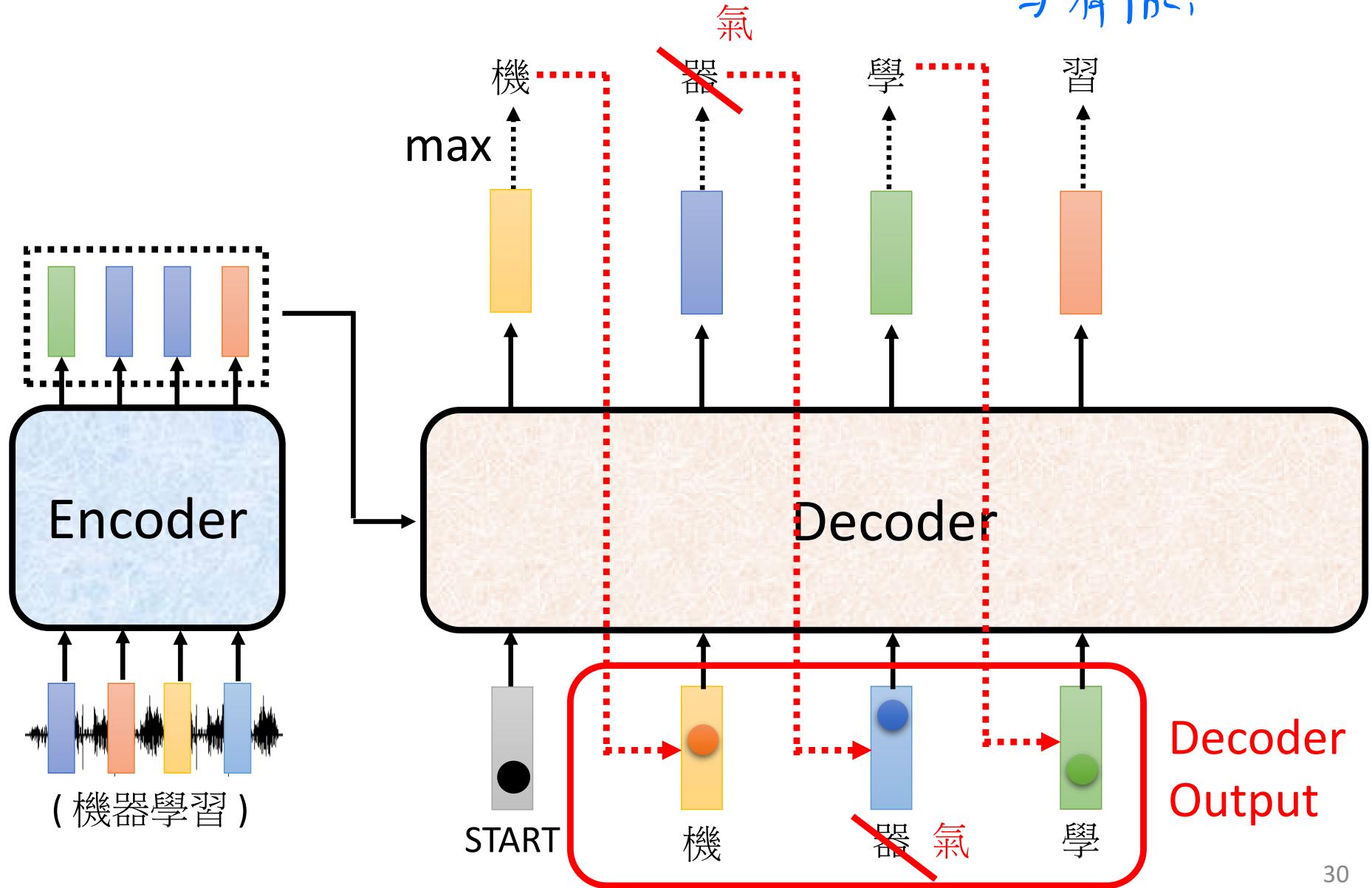
Why masked? Consider how does decoder work

⇒ 當 b^2 時，根本還沒有到 q^3 和 q^4

所以 a^4 會直接跳過

Autoregressive

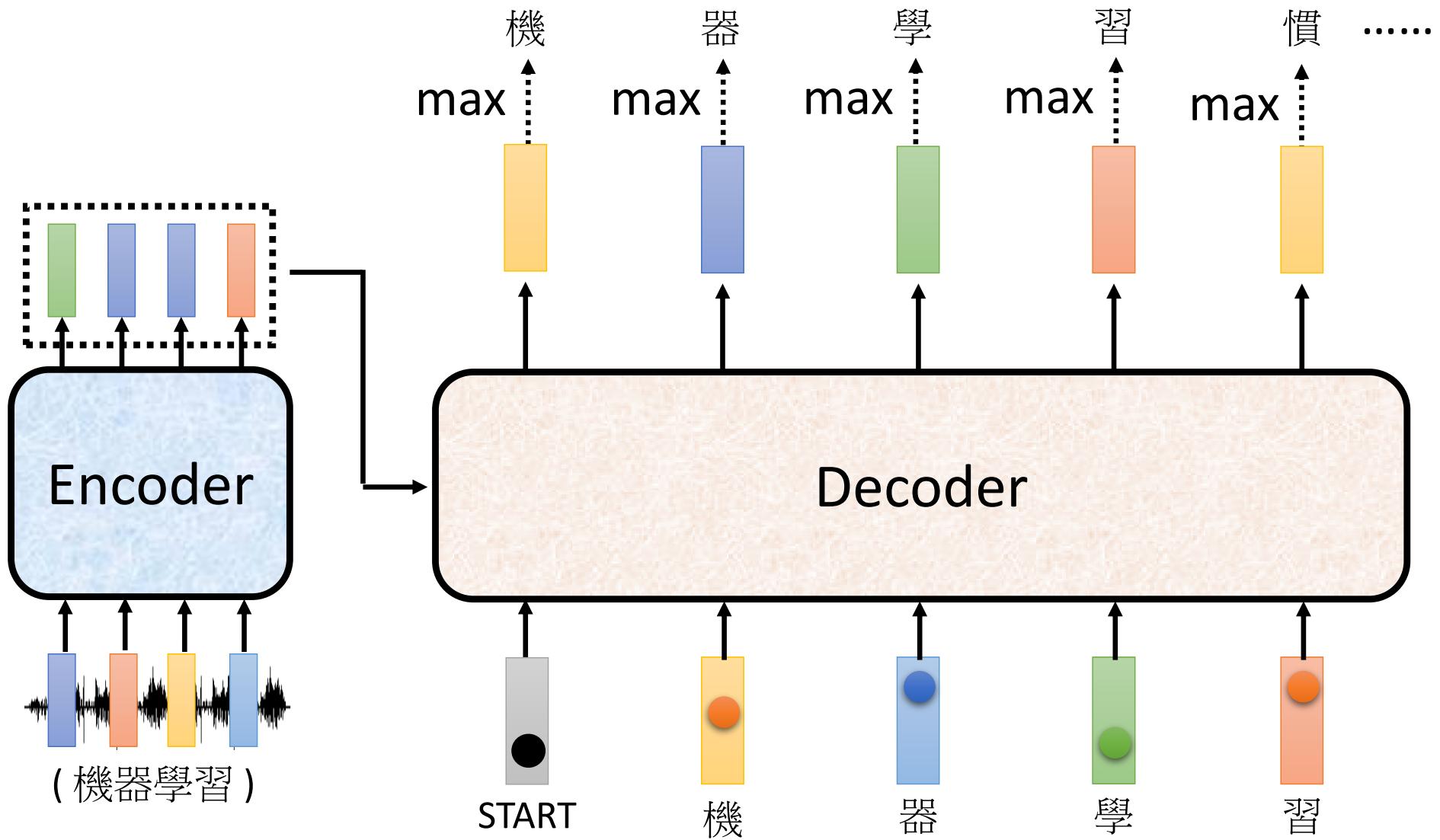
Q: 會不會一步錯，影響到後面步驟？
⇒ 有可能，之後會錯。



Autoregressive

We do not know the correct output length.

Never stop!

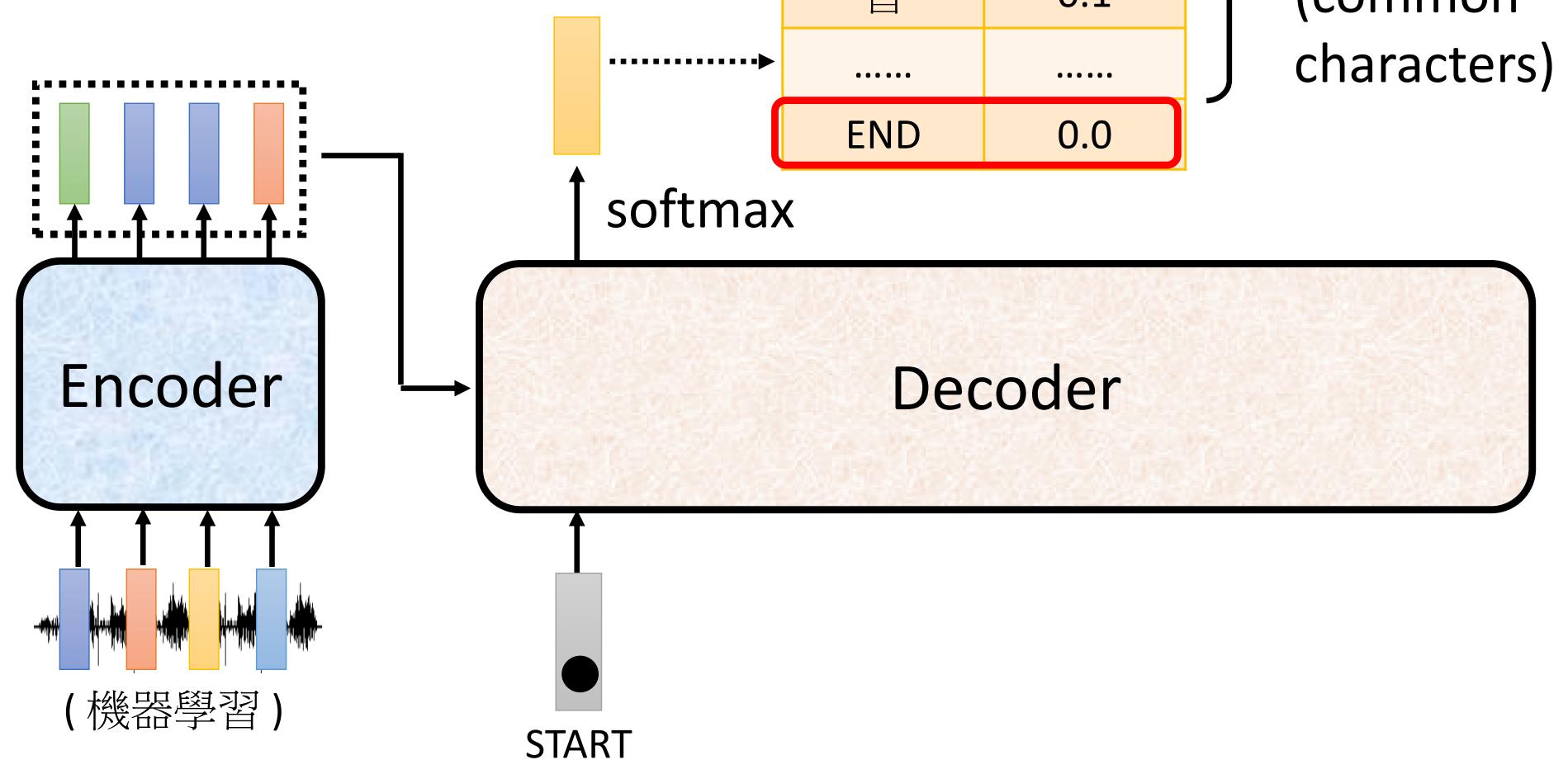


推文接龍 (Tweet Solitaire)

推	: 1:	超	06/12 10:39
推	: n:	人	06/12 10:40
推	: tation:	正	06/12 10:41
→	: host:	大	06/12 10:47
推	: :	中	06/12 10:59
推	: 403:	天	06/12 11:11
推	: :	外	06/12 11:13
推	: 527:	飛	06/12 11:17
→	: 990b:	仙	06/12 11:32
→	: 512:	草	06/12 12:15

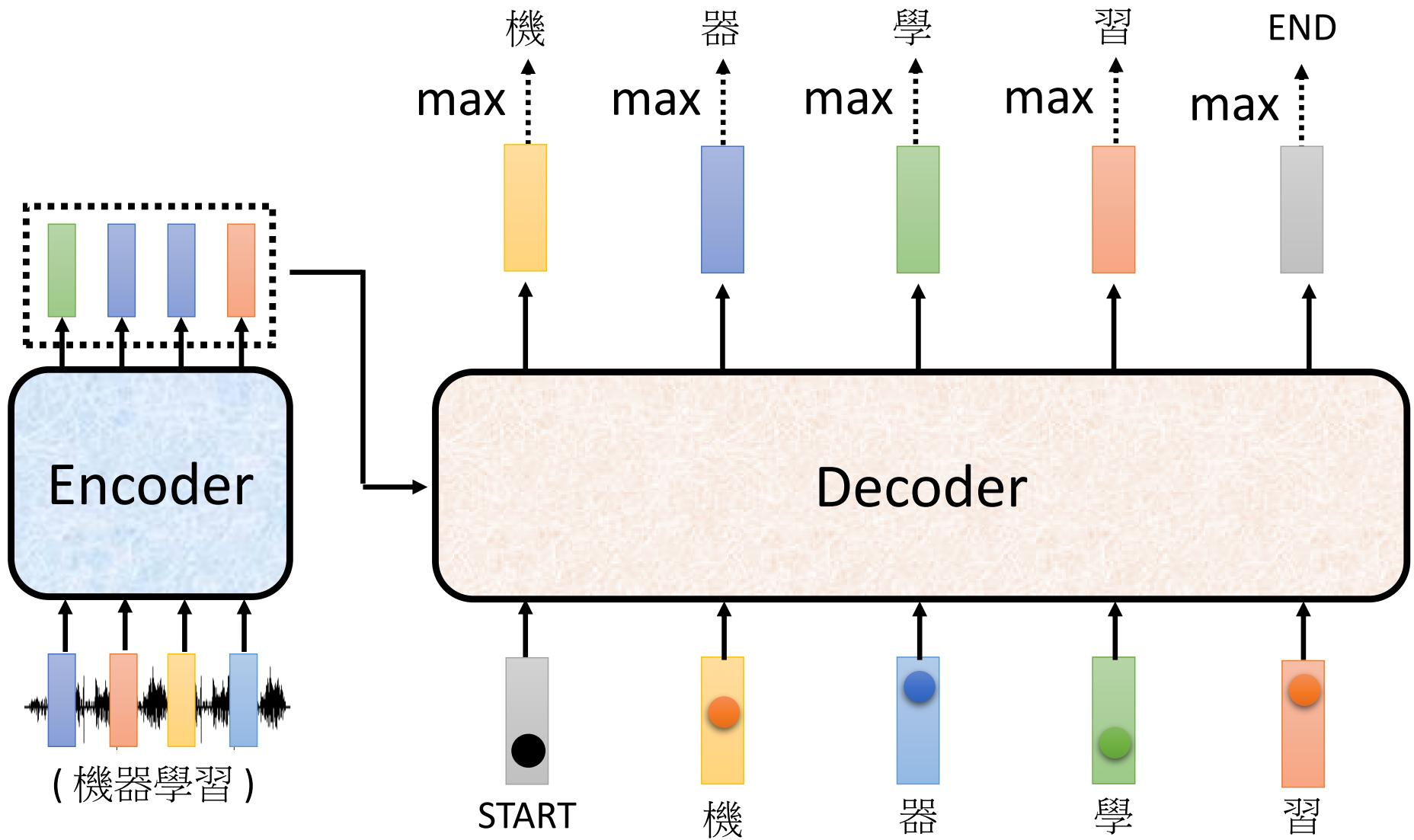
推 tlkagk: =====斷=====

Adding “Stop Token”



Autoregressive

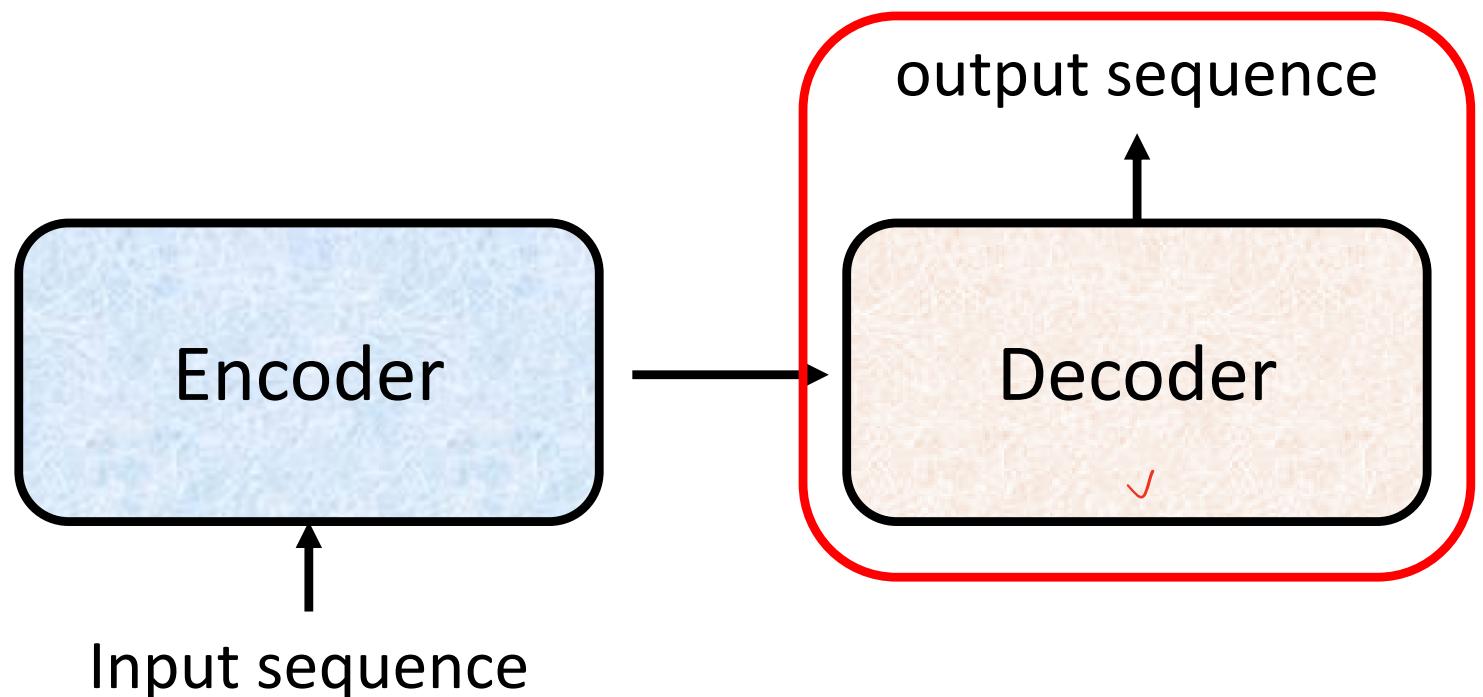
Stop at here!



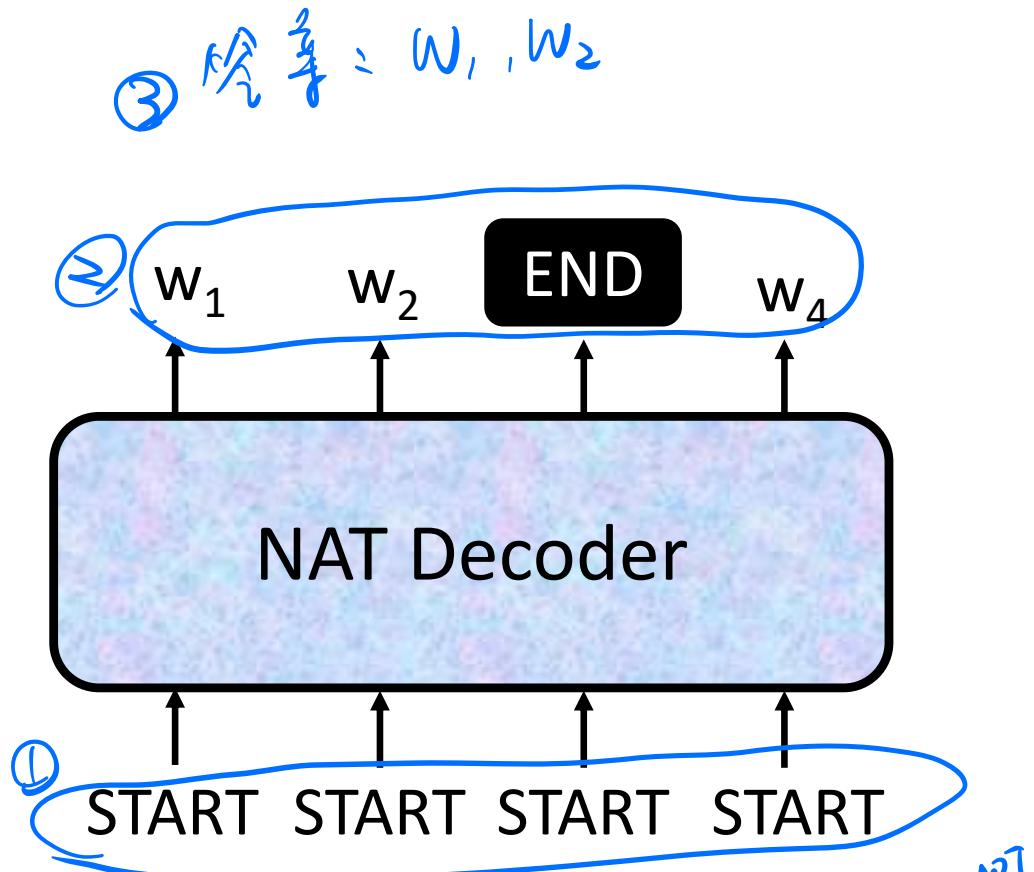
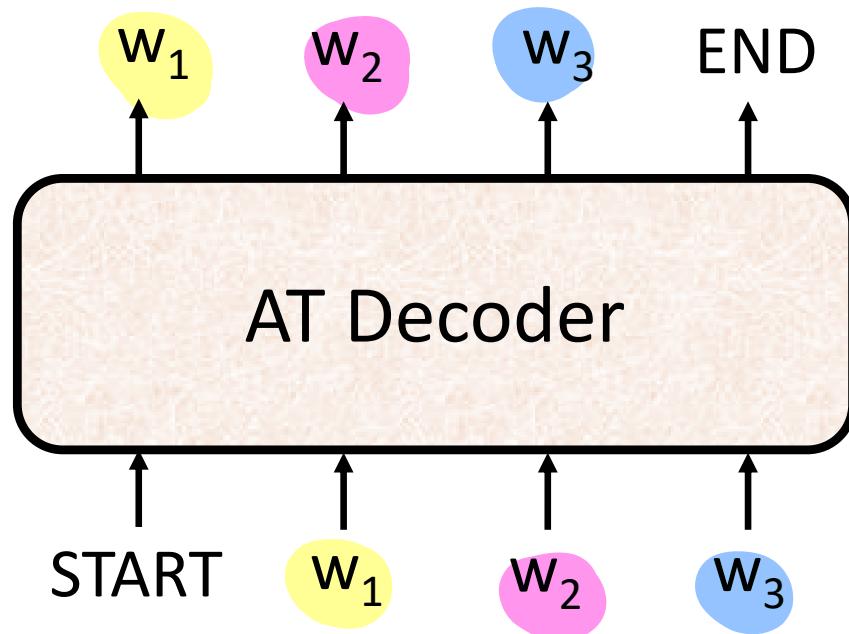
2頁搞些什，介紹一下 NAT

Decoder

– Non-autoregressive (NAT)



AT v.s. NAT



- How to decide the output length for NAT decoder? ④ の要数 START
 - Another predictor for output length ⑤ *# train-j predictor, input encoder, output n_j START e.g. 4*
 - Output a very long sequence, ignore tokens after END
⑥ 各種作成直前の他長のSTART, e.g. 300J, 然後 output 取り到END以降
- Advantage: parallel, more stable generation (e.g., TTS)
- NAT is usually worse than AT (why? Multi-modality)

To learn more

NATの障壁 (大坑)

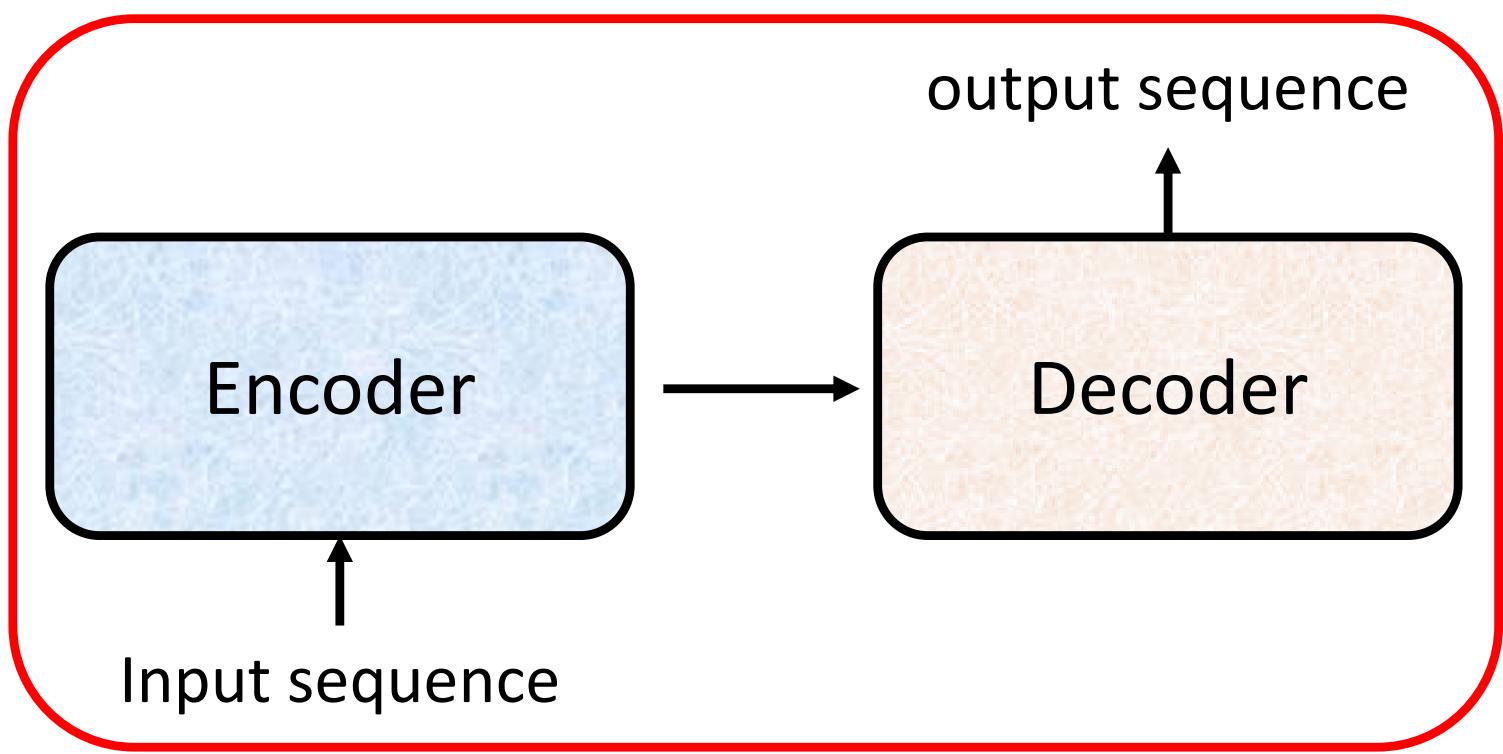


<https://youtu.be/jvyKmU4OM3c>

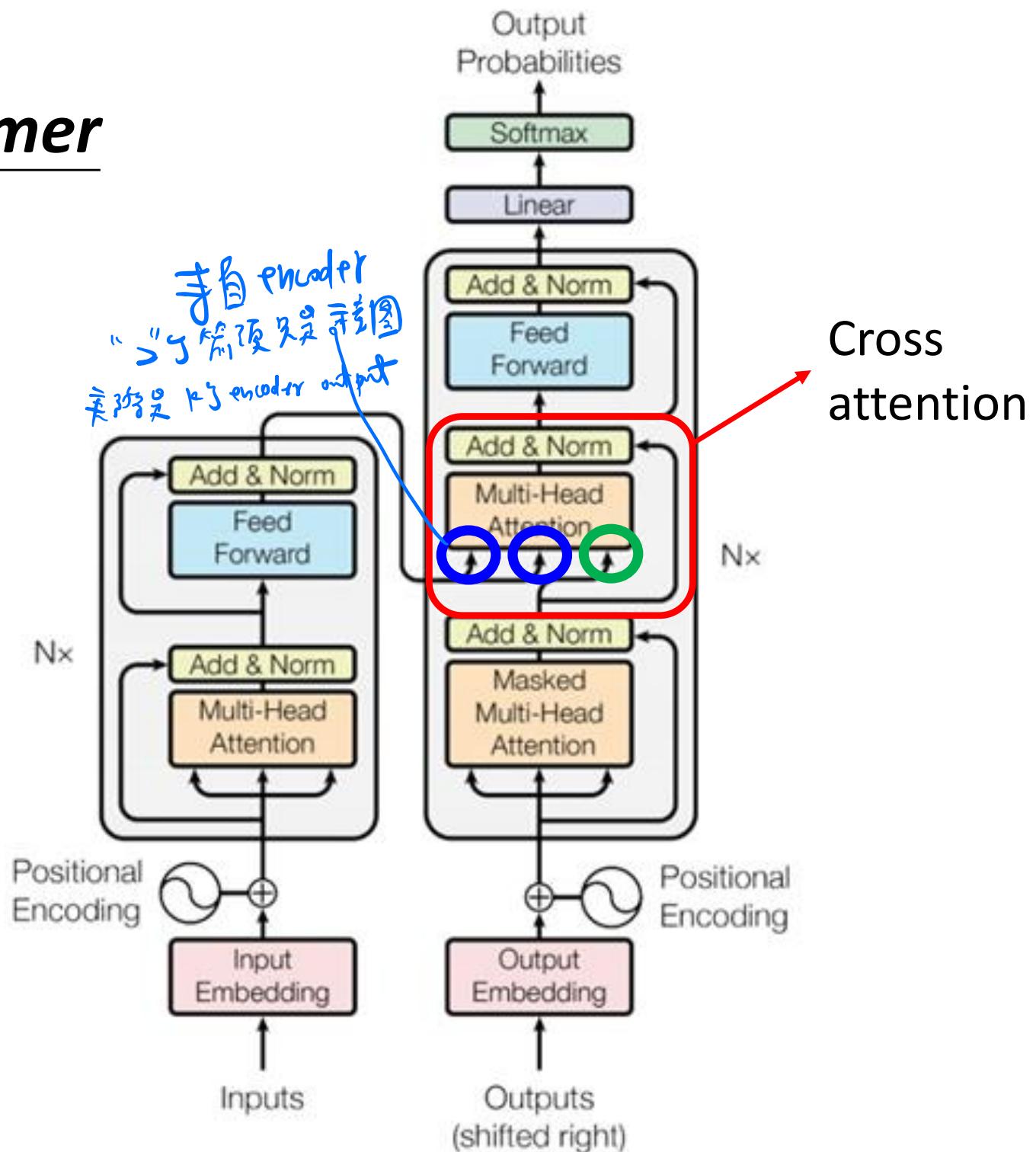
(in Mandarin)

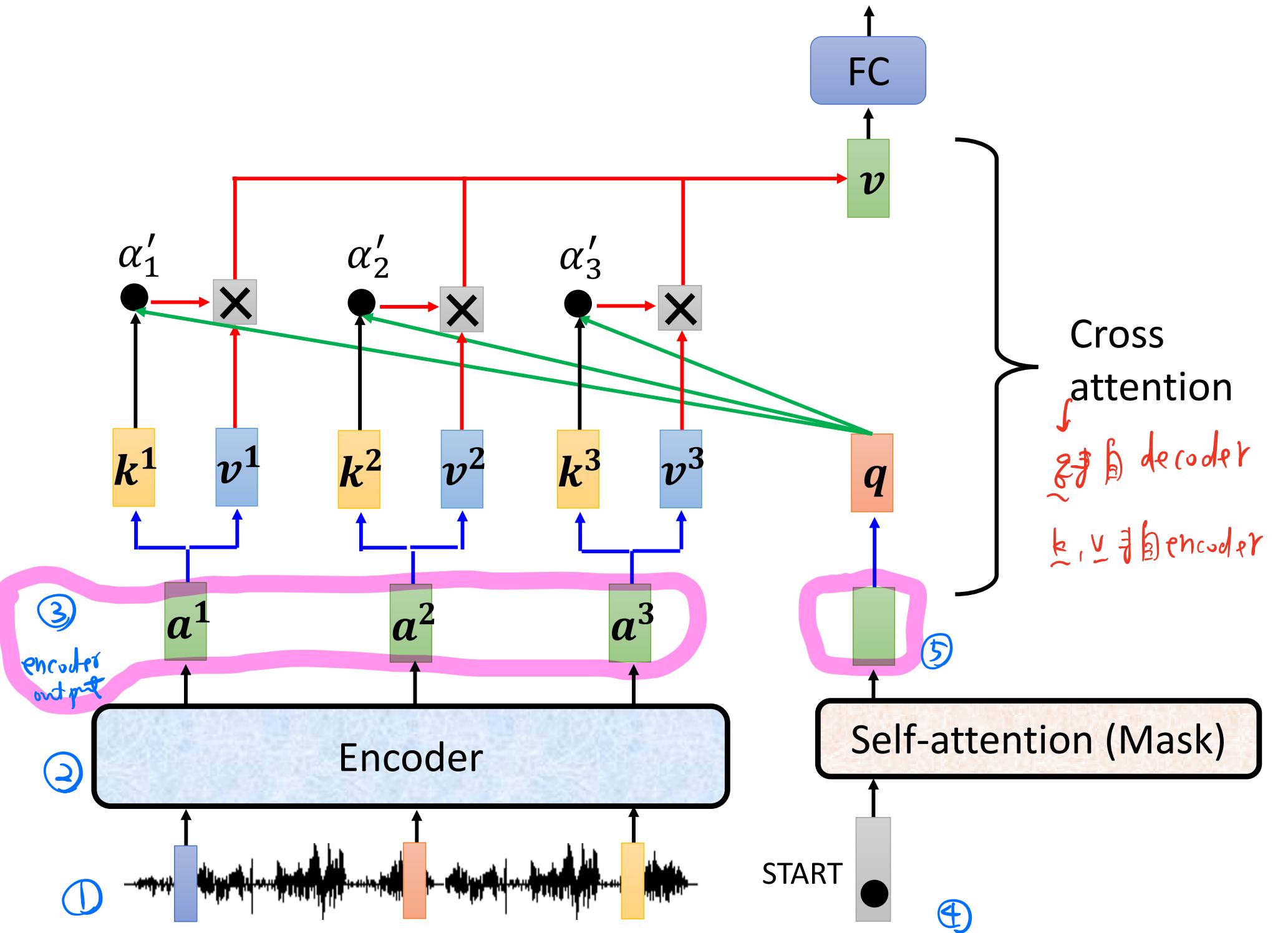
回頭者 decoder, 同じ被説明する地図 ⇒ cross-attention

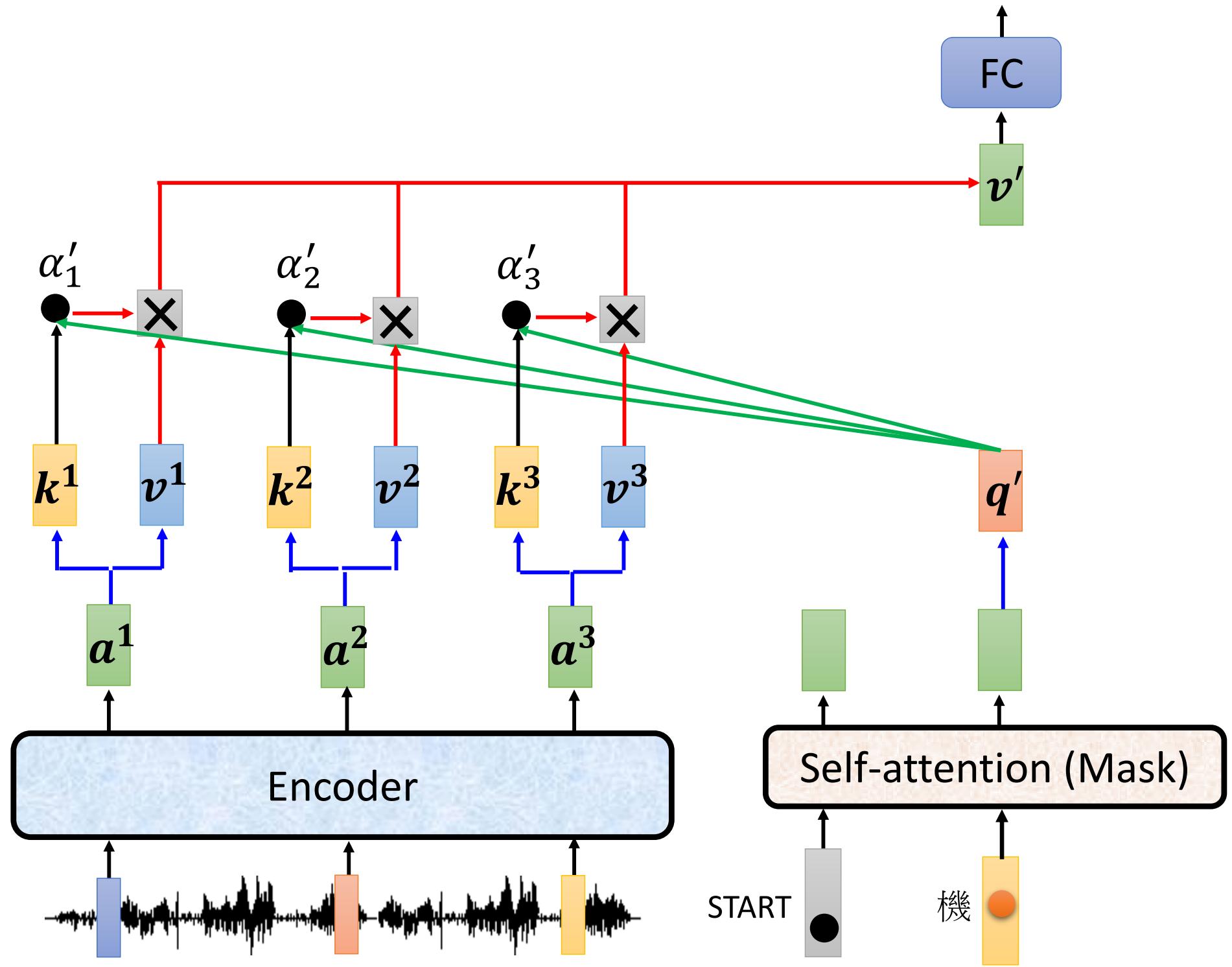
Encoder-Decoder



Transformer



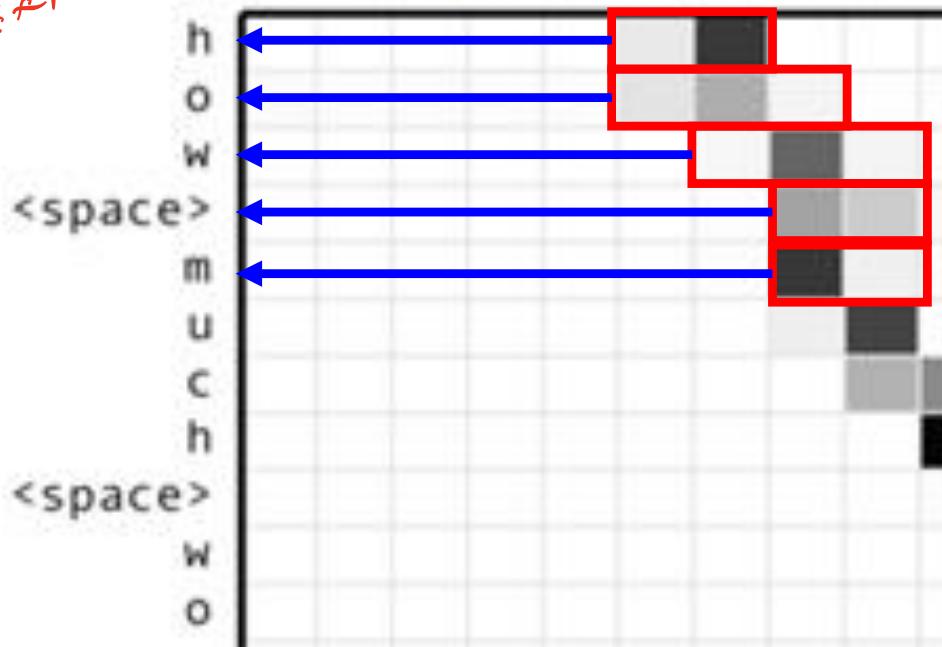
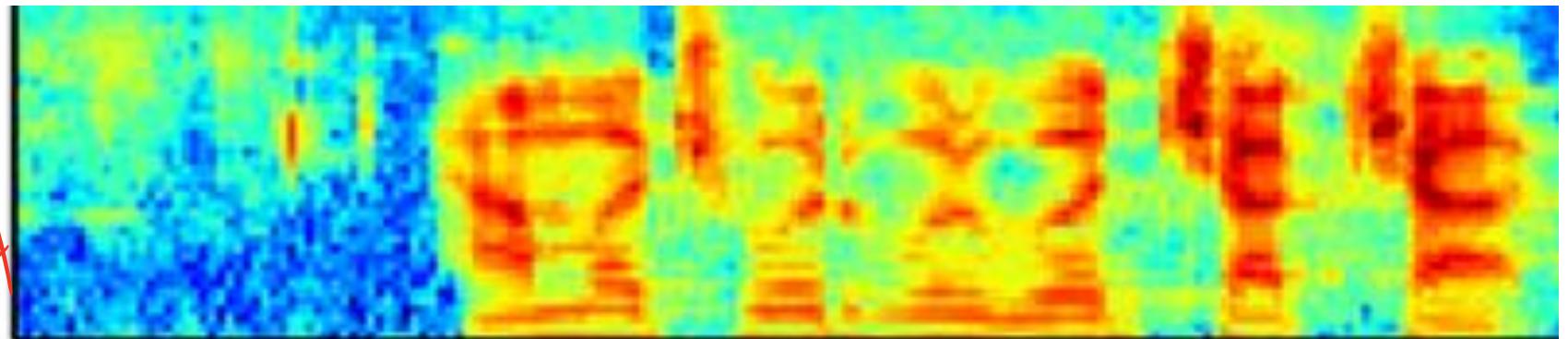




Cross Attention

Listen, attend and spell: A neural network for large vocabulary conversational speech recognition
<https://ieeexplore.ieee.org/document/7472621>

X $t_1 \dots t_n$ ↗
y $y_1 \dots y_n$ ↗
z $z_1 \dots z_n$ ↗
w $w_1 \dots w_n$ ↗
h $h_1 \dots h_n$ ↗
→ $\hat{w}_1 \dots \hat{w}_n$ ↗



(This is not transformer.)

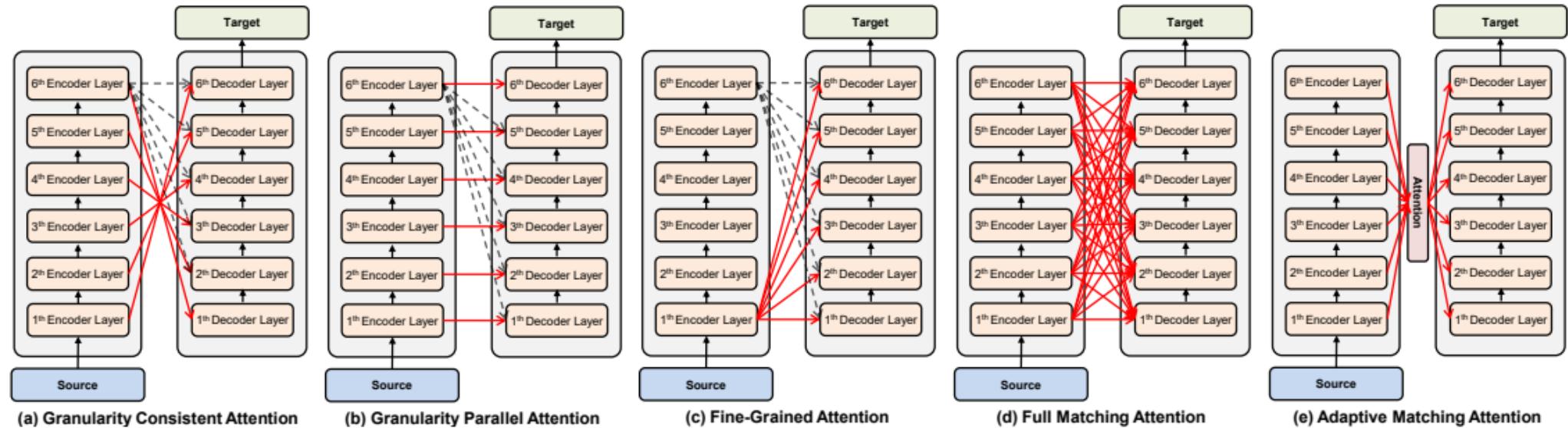
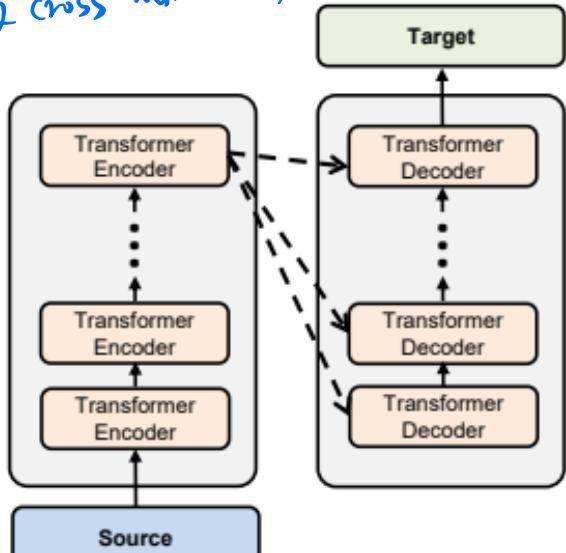
Cross Attention

Source of image:

② <https://arxiv.org/abs/2005.08081>

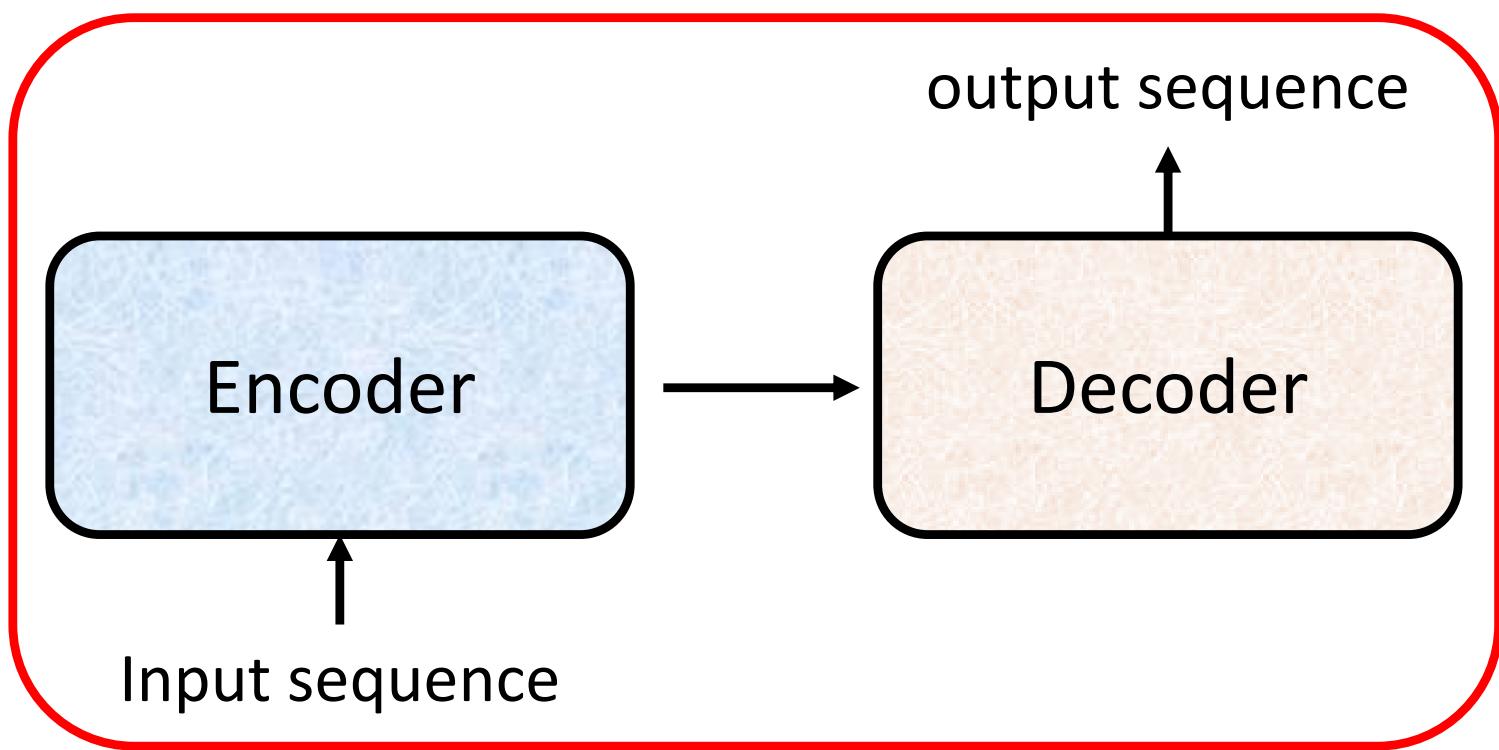
這篇 paper 裡試用の式、做筆記

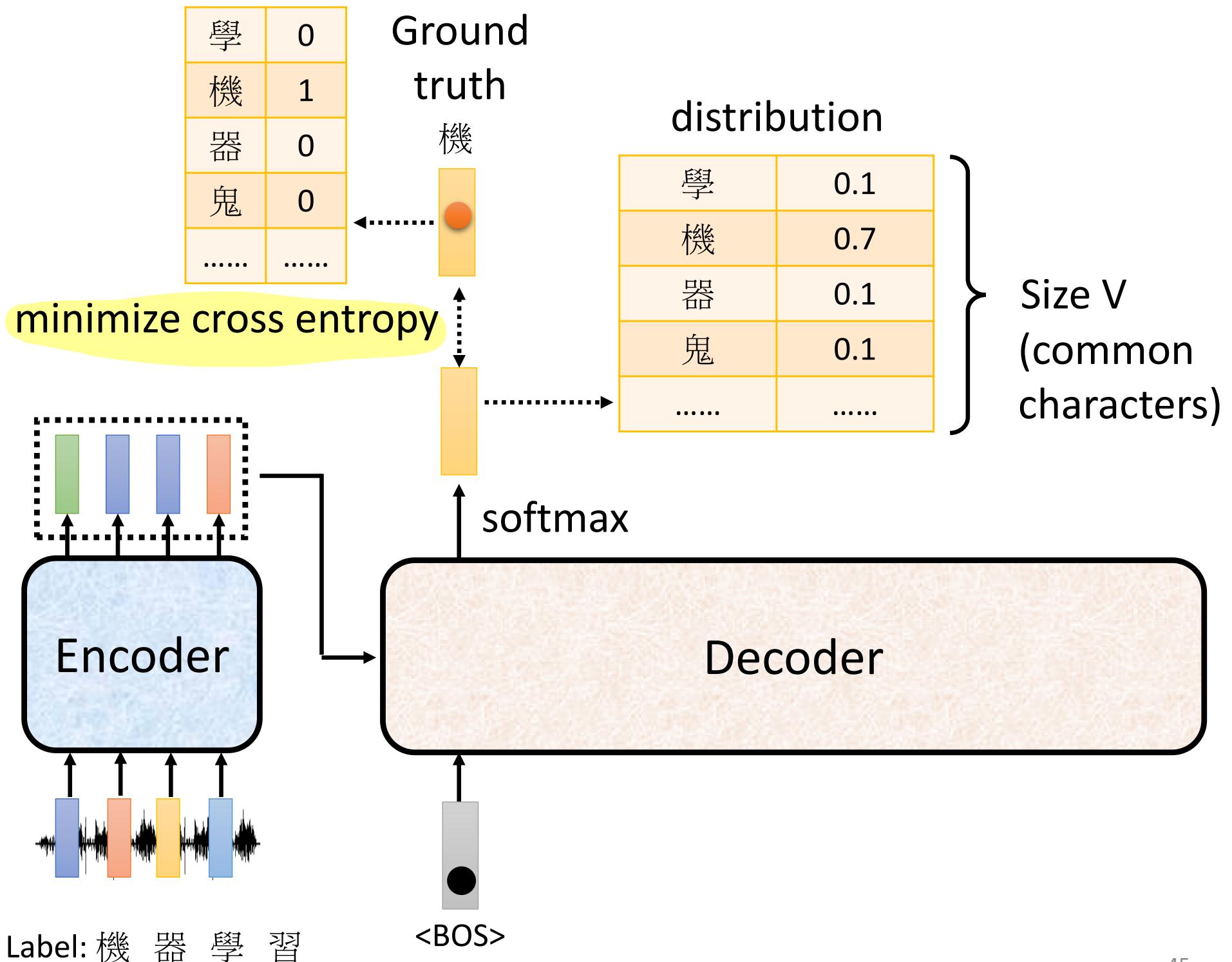
Q This paper is using encoder from -16 to 5, L, P
decoder fit cross attention,



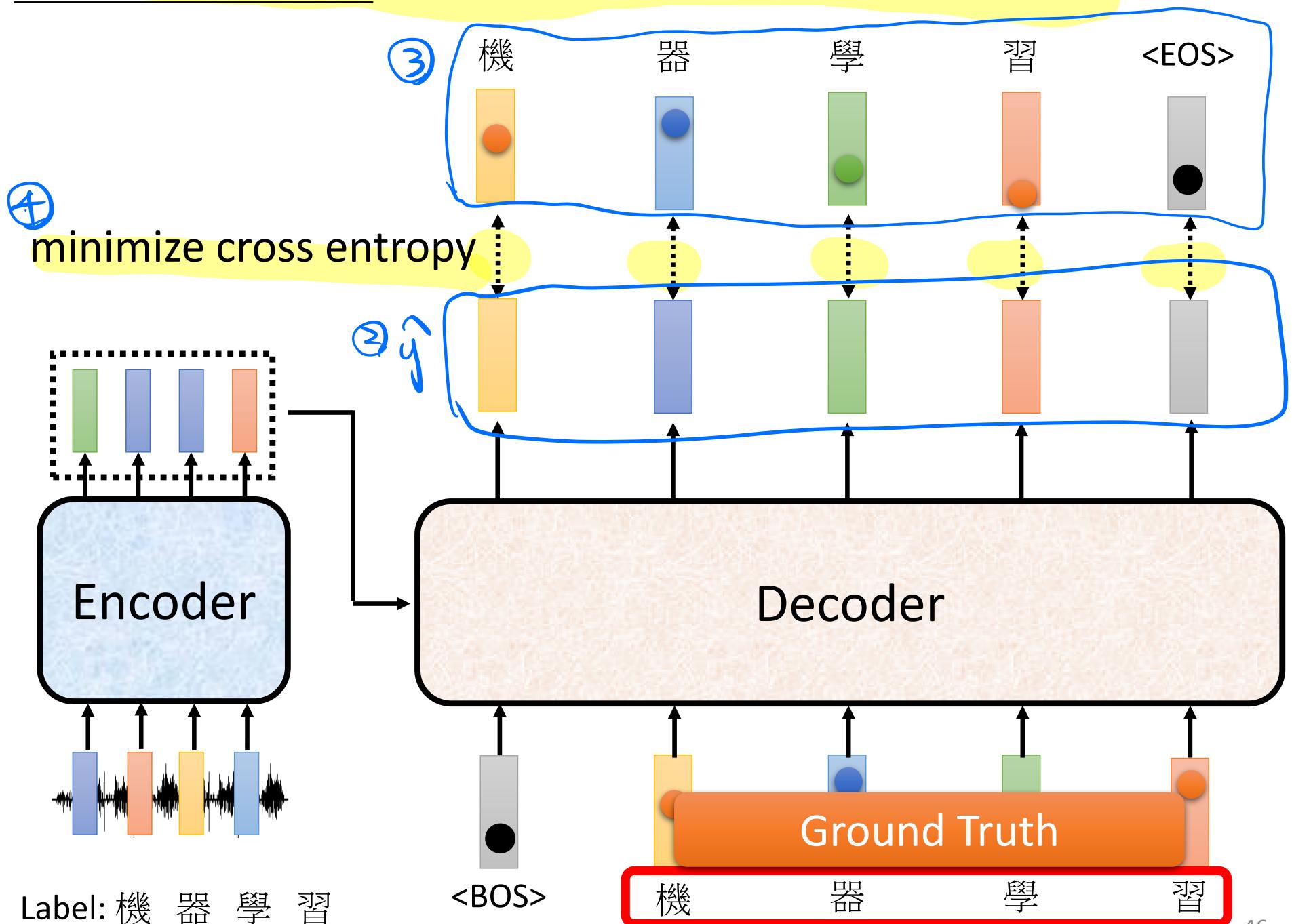
剛才在講 inference,
現在要講怎麼 train

Training



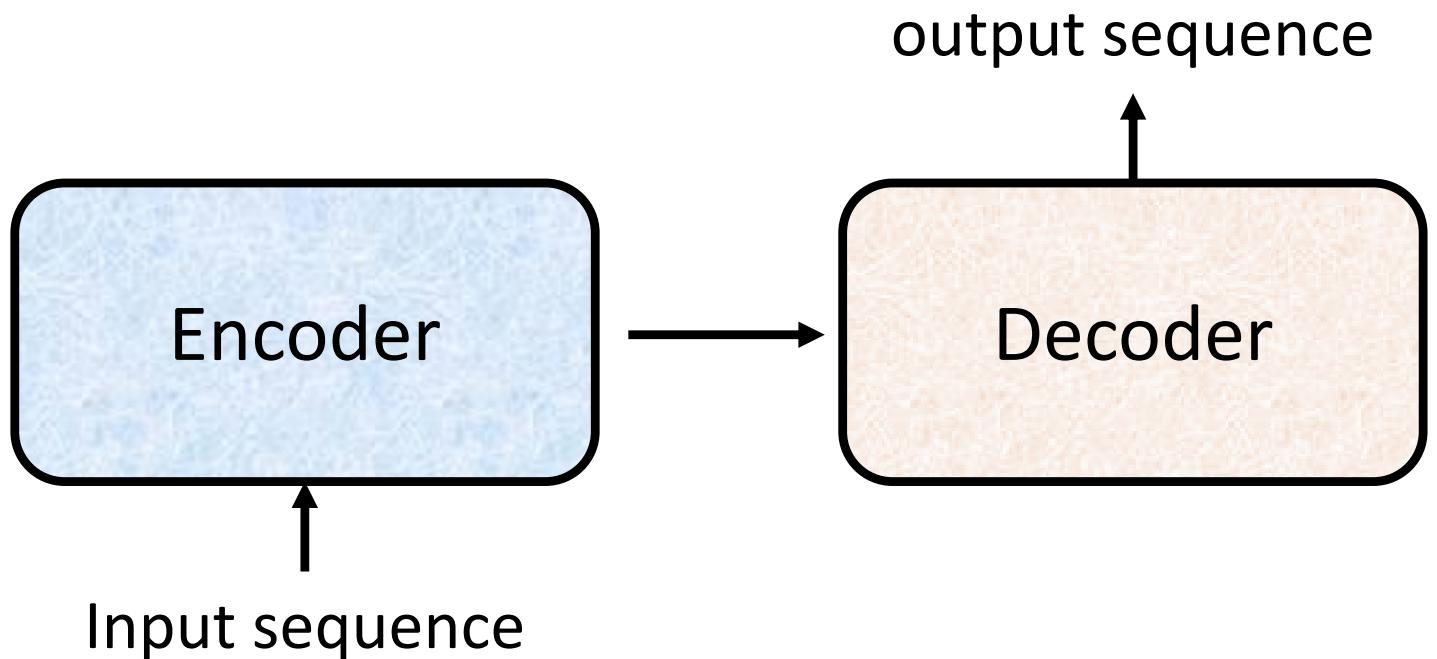


Teacher Forcing: using the ground truth as input.



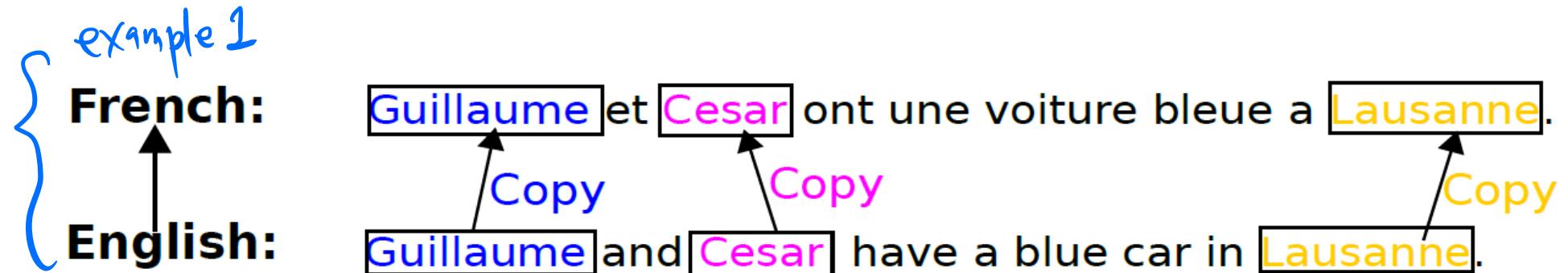
① 令 encoder 的 input vector 为 ground truth

Tips

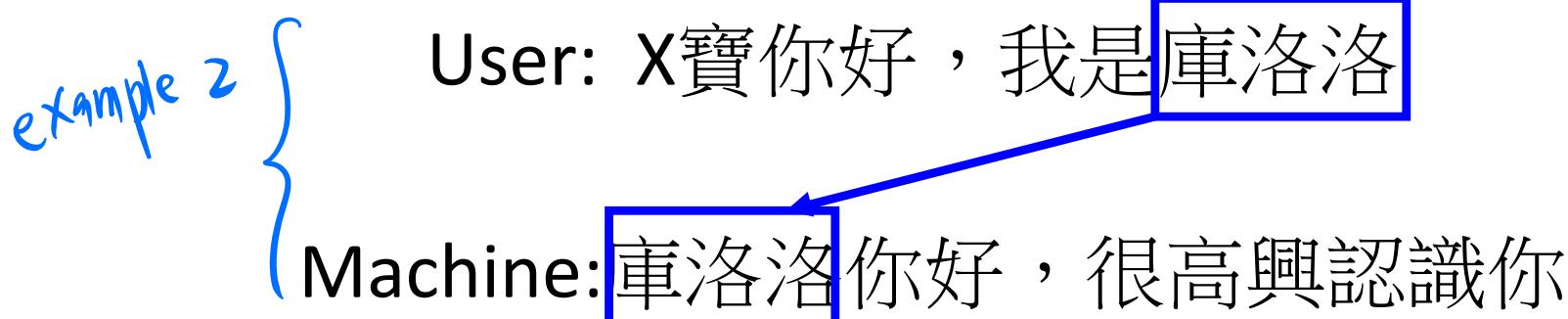


Copy Mechanism

Machine Translation



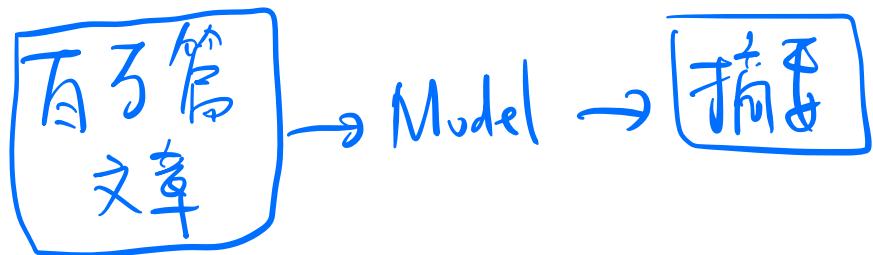
Chat-bot



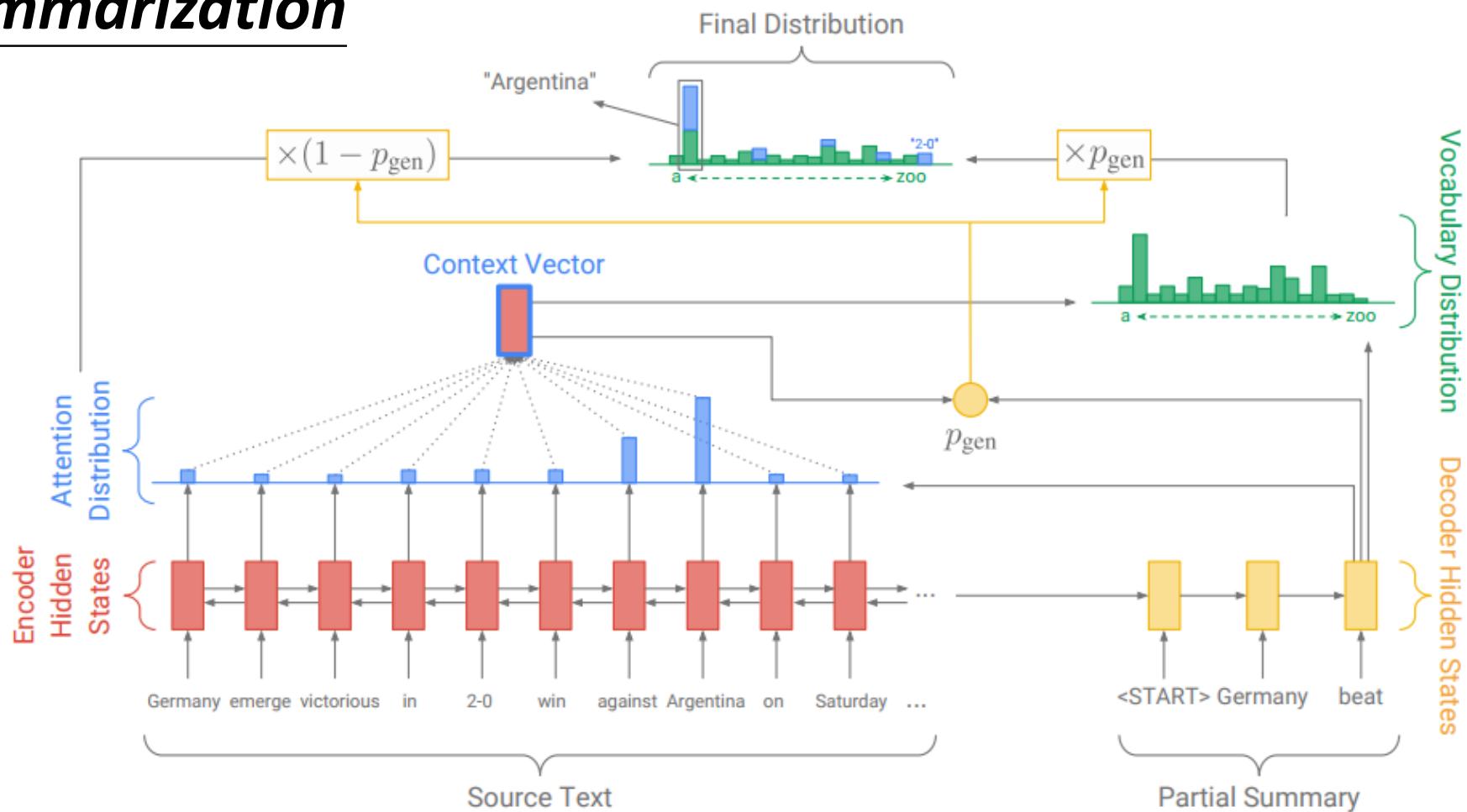
Copy Mechanism

example 3 = 摘要 (需要去 copy 15 个词)

Summarization



<https://arxiv.org/abs/1704.04368>



Copy Mechanism

Pointer Network



<https://youtu.be/VdOyqNQ9aww>

Incorporating Copying Mechanism in Sequence-to-Sequence Learning

<https://arxiv.org/abs/1603.06393>

Guided Attention



高雄發大財我現在要出征



發財發財發財發財



發財發財發財



發財發財



發財 (Missing an input character!)



這樣的 e- 語音合成任閒的，文字是由左公列右

① 総合 attention 有り因式の状況、
② output の attention は後で取り扱う

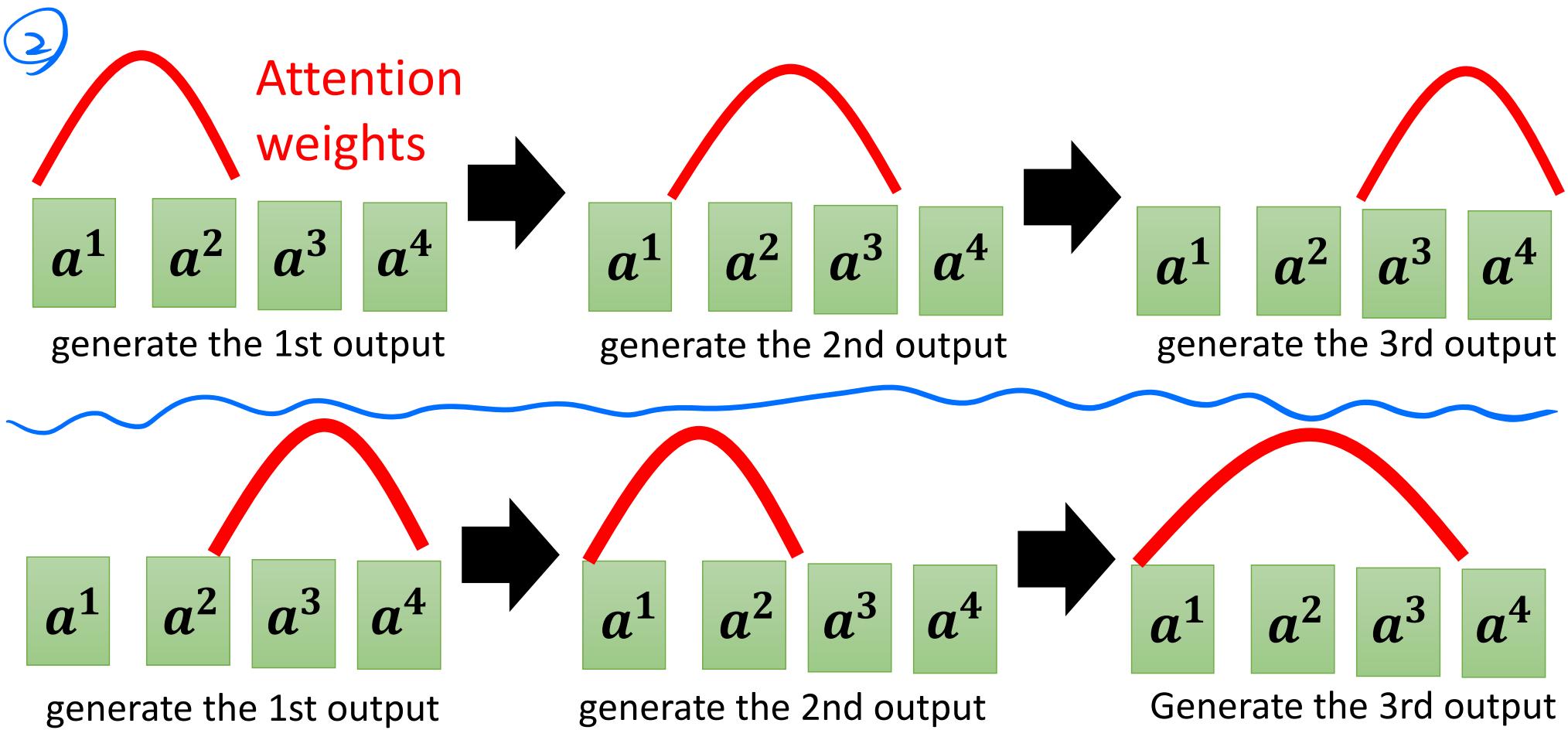
Guided Attention

Monotonic Attention
Location-aware attention

大抵
自己
言語

In some tasks, input and output are monotonically aligned.

For example, speech recognition, TTS, etc.



③ 総合 attention は上手い

④ Something wrong!

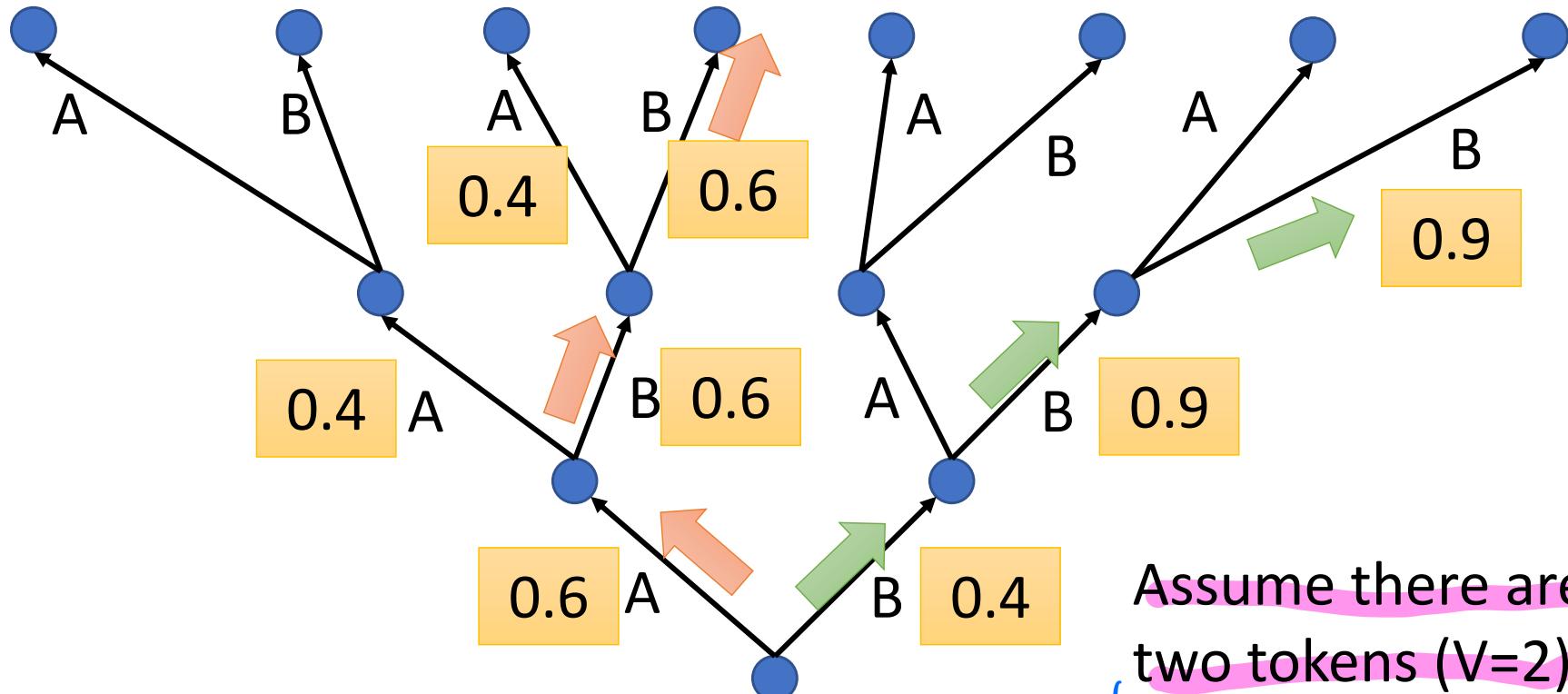
Beam Search

- ② The red path is **Greedy Decoding**.
- ③ The green path is the best one.
- ④ Not possible to check all the paths ... → Beam Search 回溯

→ 那 output 是 A, B, B

→ output 是 B, B, B

⑤ → Beam Search



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

→ 很多世上只有 2 个字, A, B

Beam search 有時有用, 有時沒用, 這篇文章

Sampling

告訴你 beam search 是用

The Curious Case of Neural Text Degeneration

<https://arxiv.org/abs/1904.09751>

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 México/Universidad Nacional Autónoma de ...")

Pure Sampling: ⑧

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."

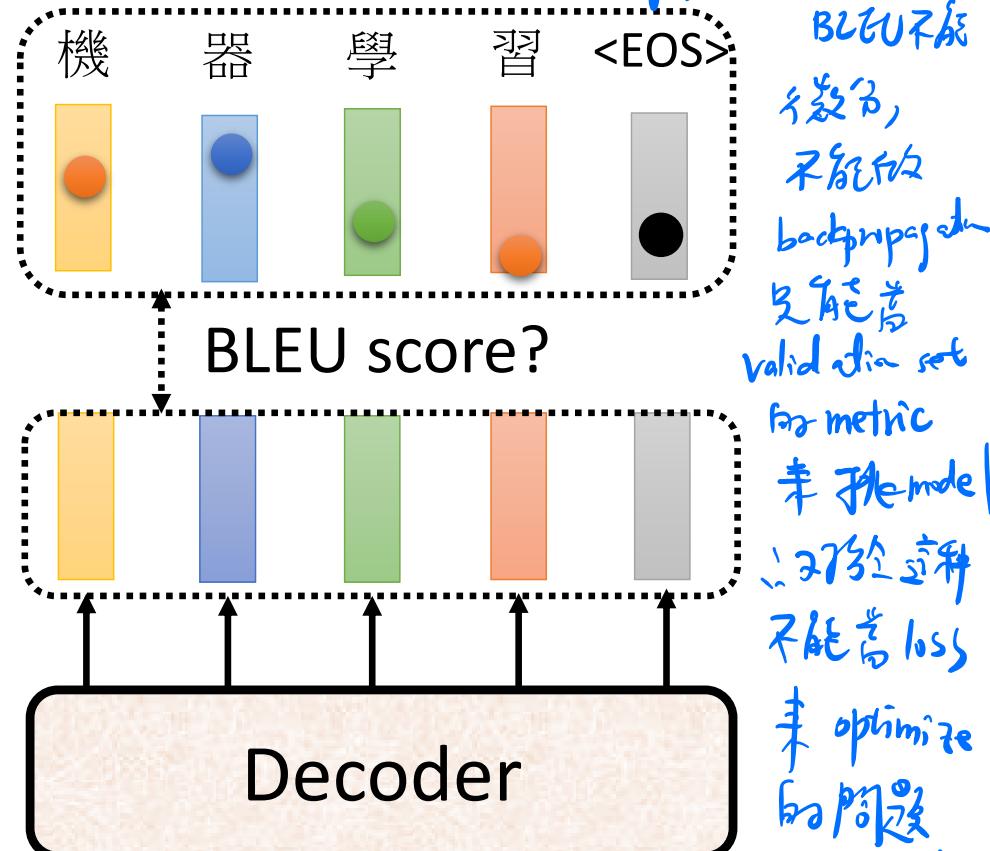
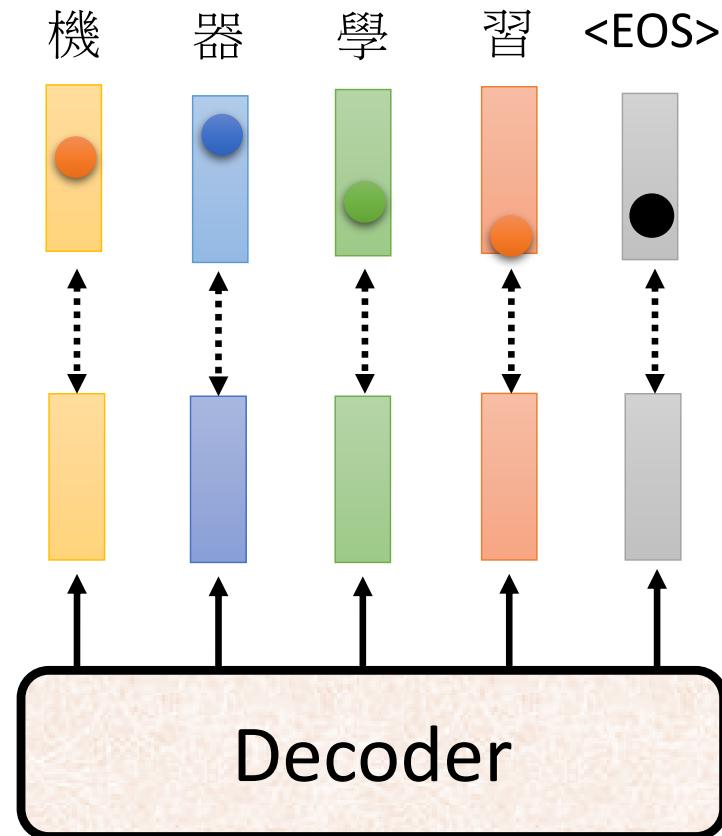
Randomness is needed for decoder when generating sequence in some tasks.

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

在追求完美之中

Optimizing Evaluation Metrics?

← 作業中，最優化的模型好像 metric 是 BLEU，但 training 却用 cross-entropy，why? => ''



How to do the optimization?

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

把 BLEU 作為 reward, decoder 為 agent, 不重 train ->

There is a mismatch! 😞

exposure bias

training obj, input to decoder 即是 ground truth

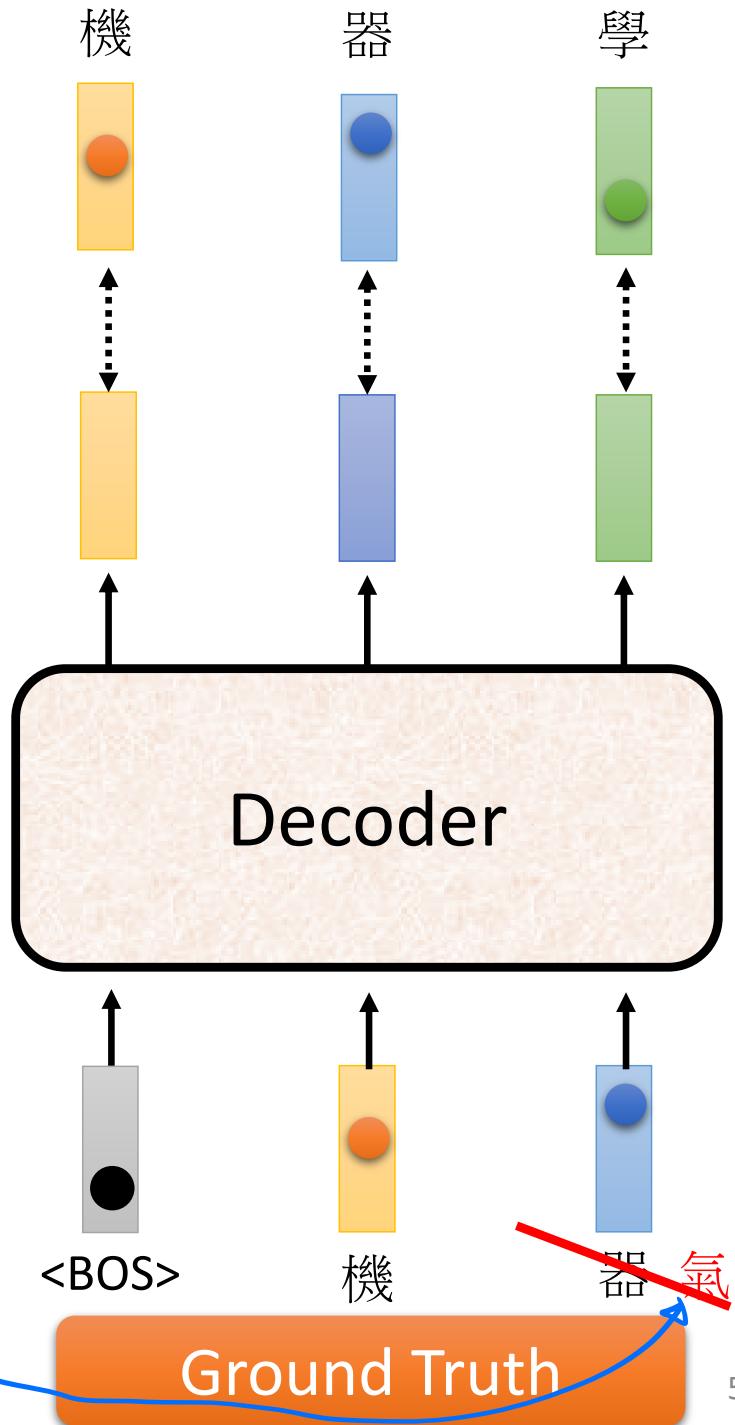
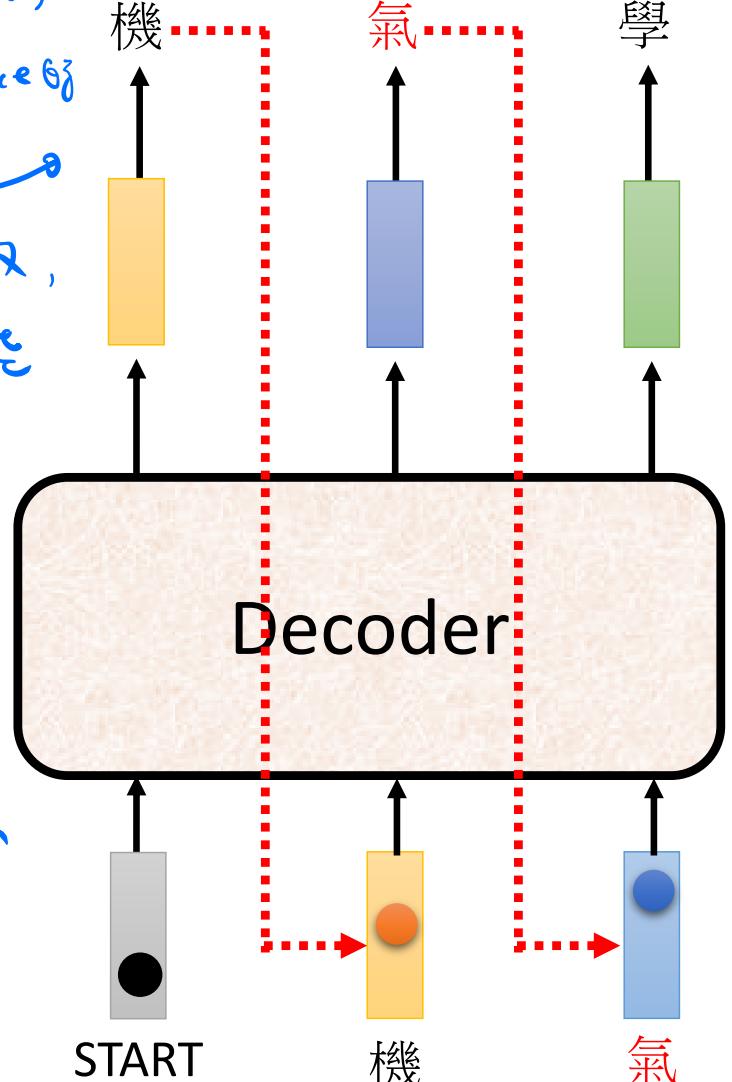
但 inference obj

是這樣做，
那就不能
一步錯
步步錯

那，
解決方法

training obj

故意加一些錯誤的 noise



立場叫

Scheduled Sampling

如下圖

- Original Scheduled Sampling

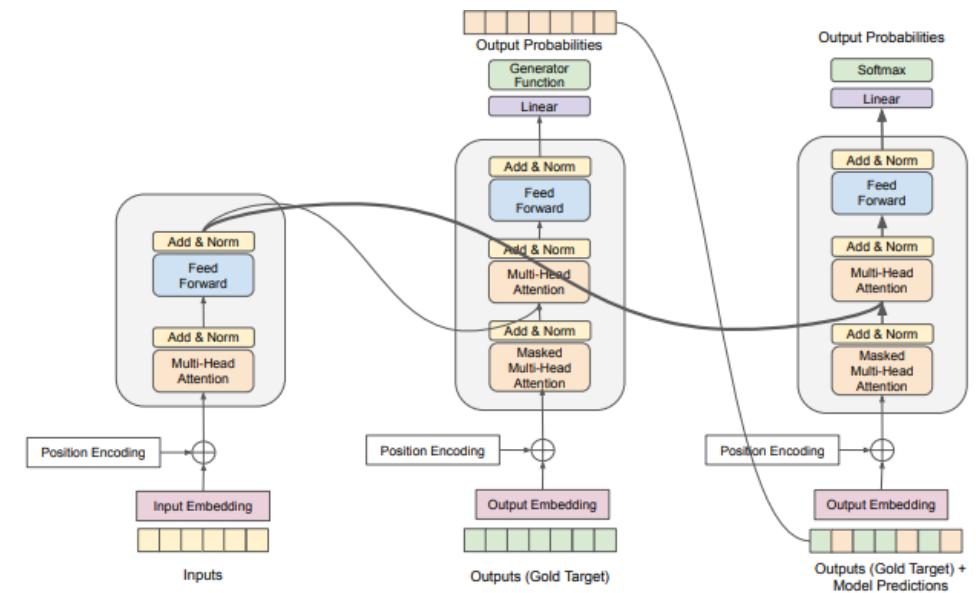
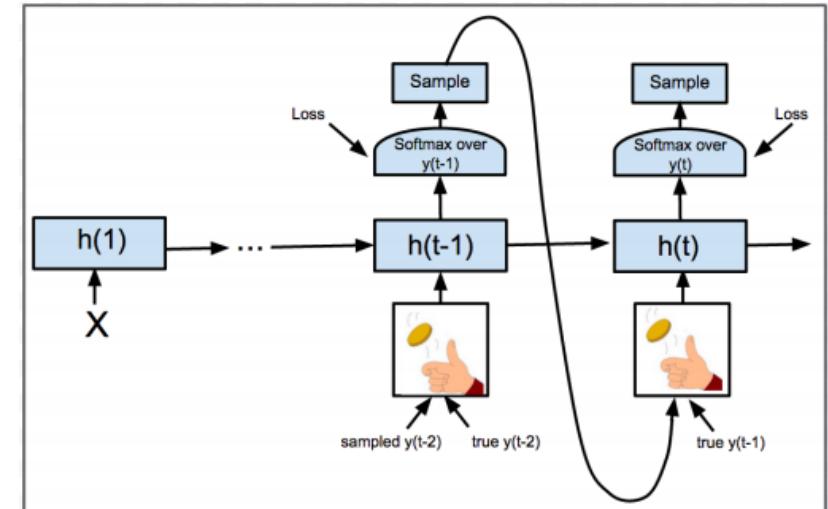
<https://arxiv.org/abs/1506.03099>

- Scheduled Sampling for Transformer

<https://arxiv.org/abs/1906.07651>

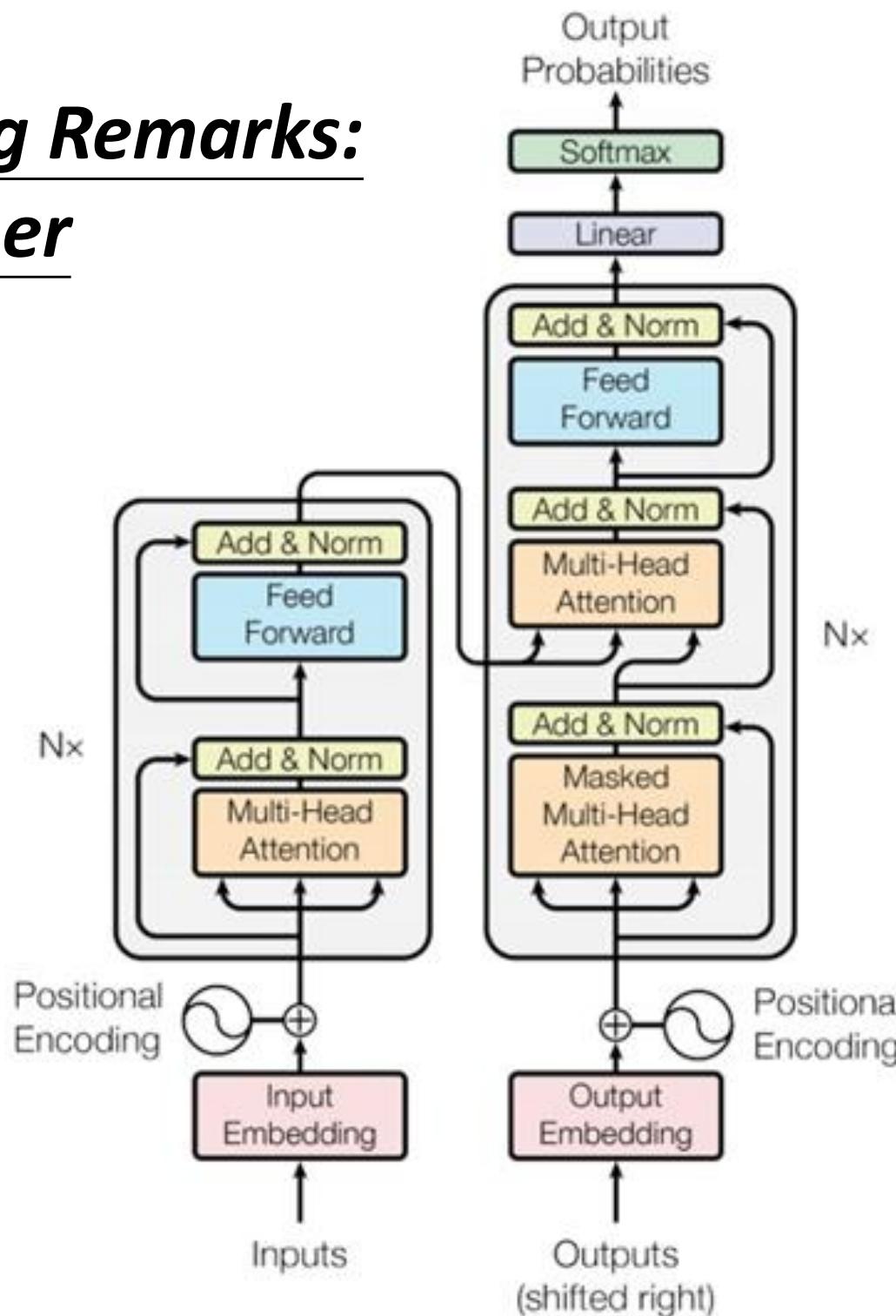
- Parallel Scheduled Sampling

<https://arxiv.org/abs/1906.04331>



Schedule Sampling

Concluding Remarks: Transformer



Q&A