



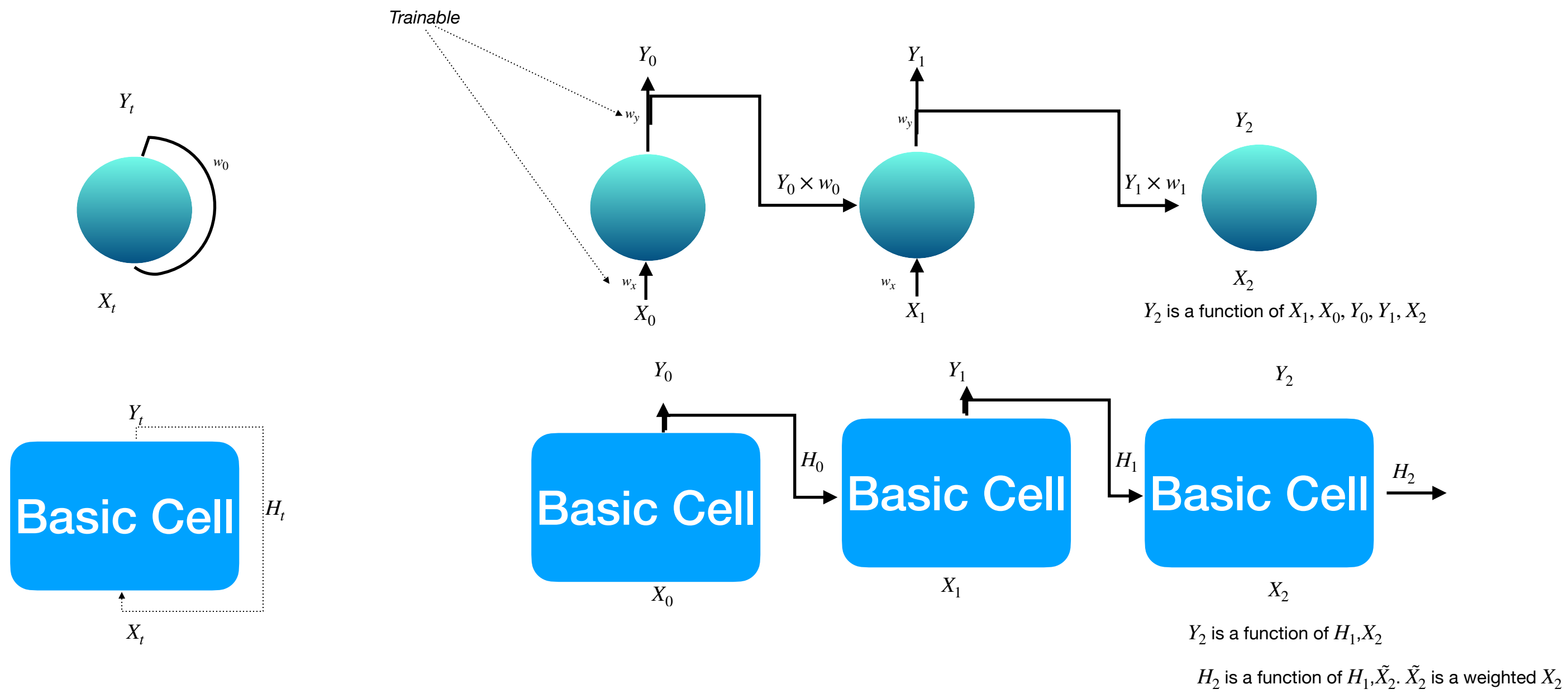
Stateless

Learns on random portions of text, without any information about the text

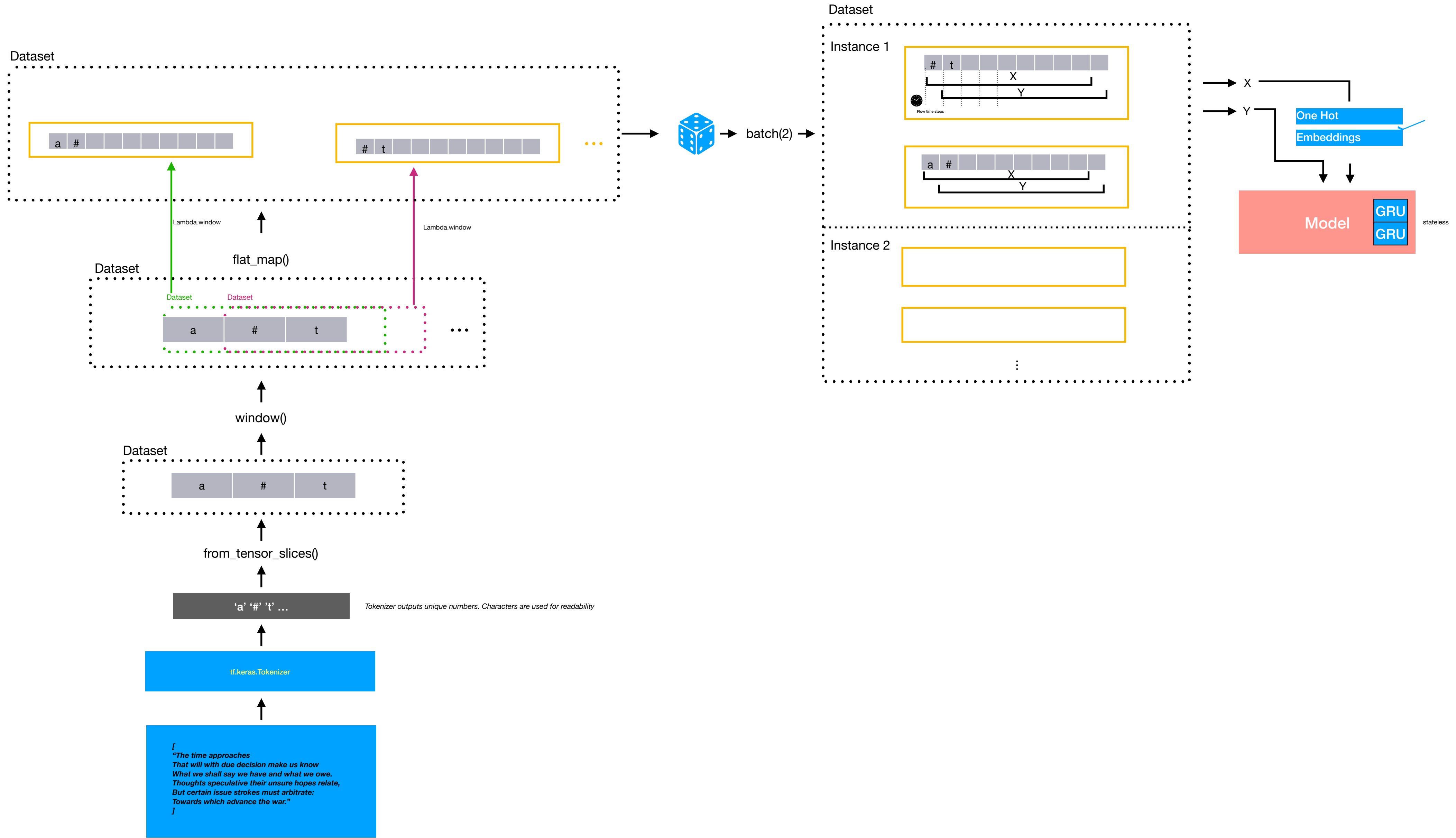


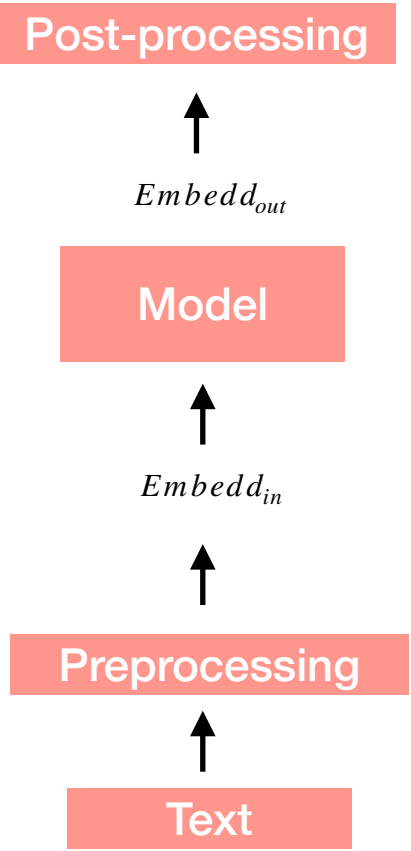
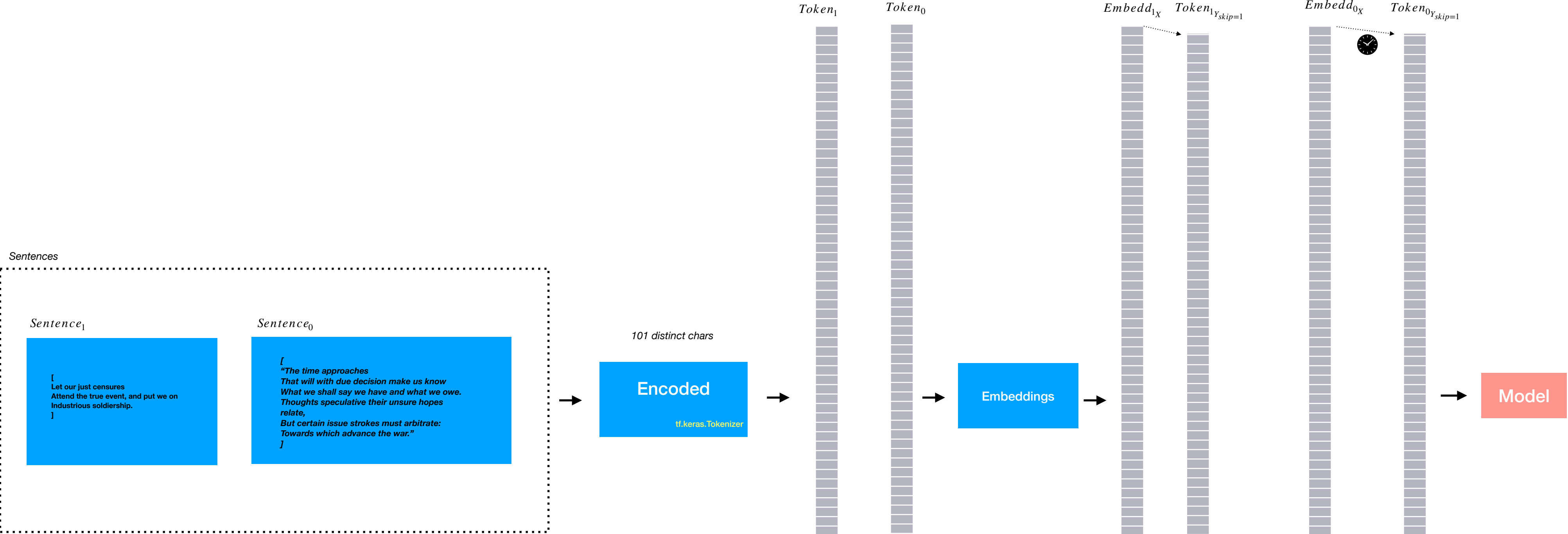
Stateful

Learned hidden state preserved , allowing the model to learn longer patterns



Dataset of **tensors** of size N . Batch size equal to 1

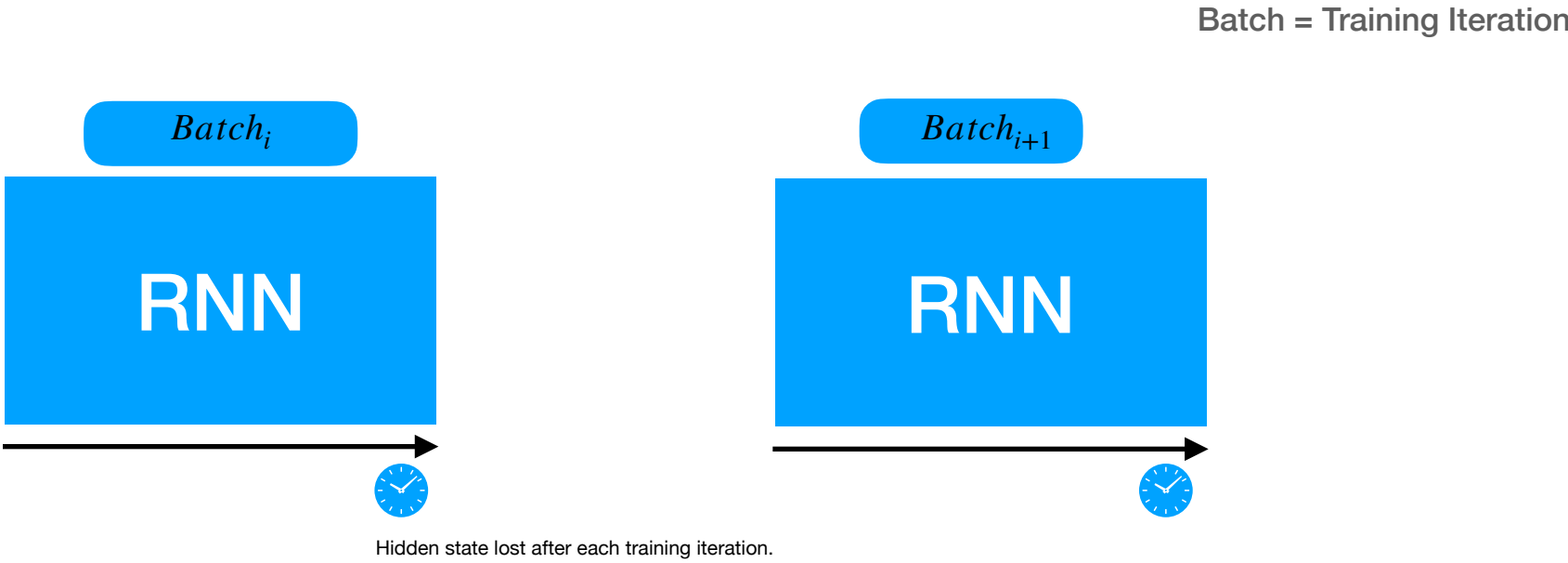




Model train gray data to predict single sample.
Models goal is to train sentence and predict next character

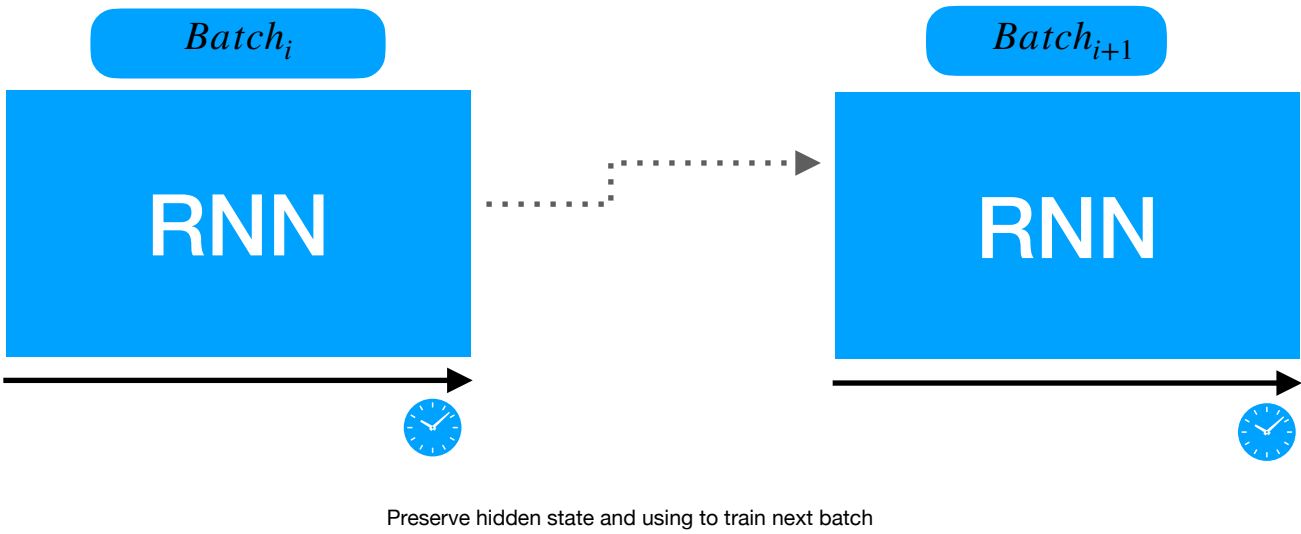
Stateless

Learns on random portions of text, without any information about the text



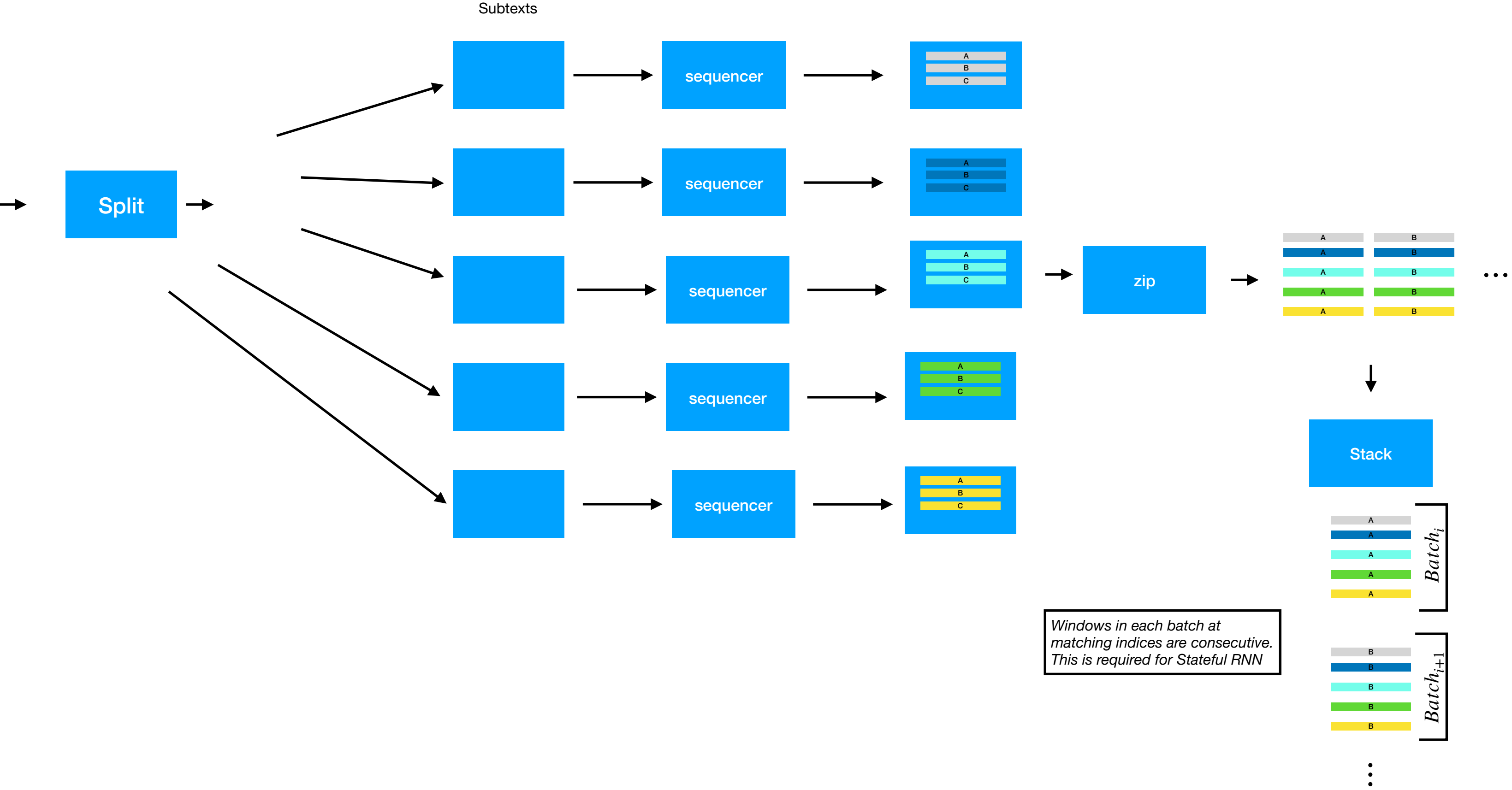
Stateful

Learned hidden state preserved , allowing the model to learn longer patterns

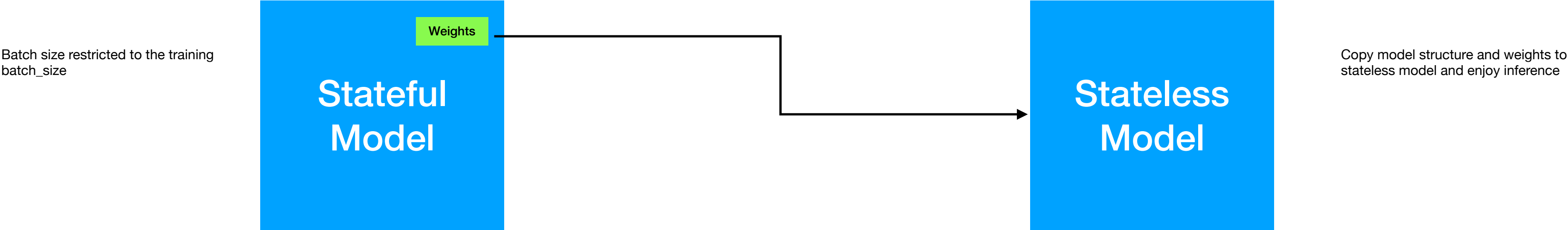


Batching: Stateful RNN

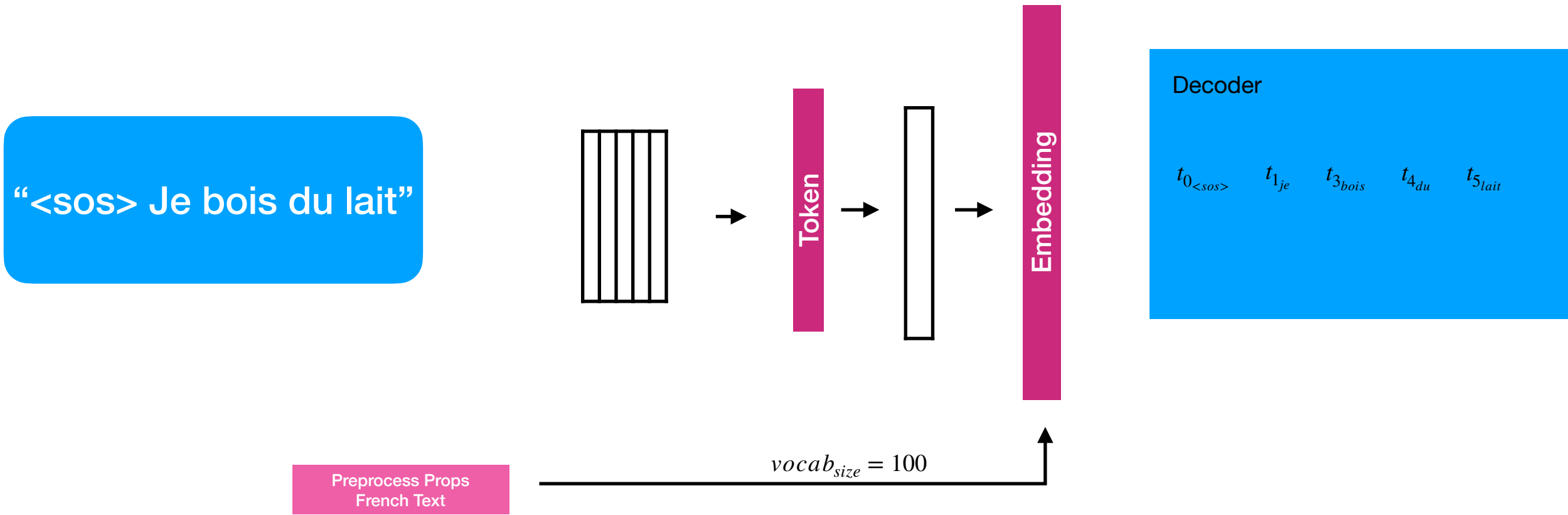
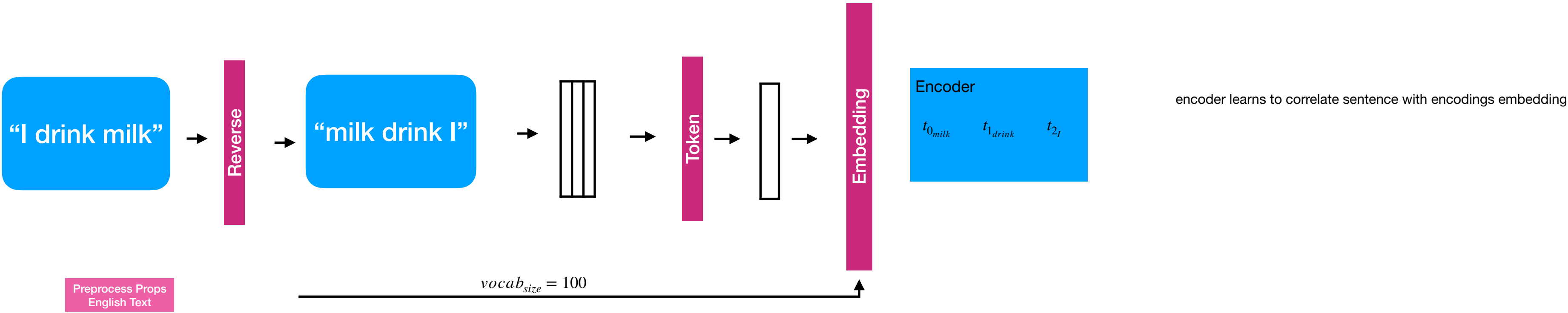
[
 “The time approaches
 That will with due decision make us know
 What we shall say we have and what we owe.
 Thoughts speculative their unsure hopes relate,
 But certain issue strokes must arbitrate:
 Towards which advance the war.”
]

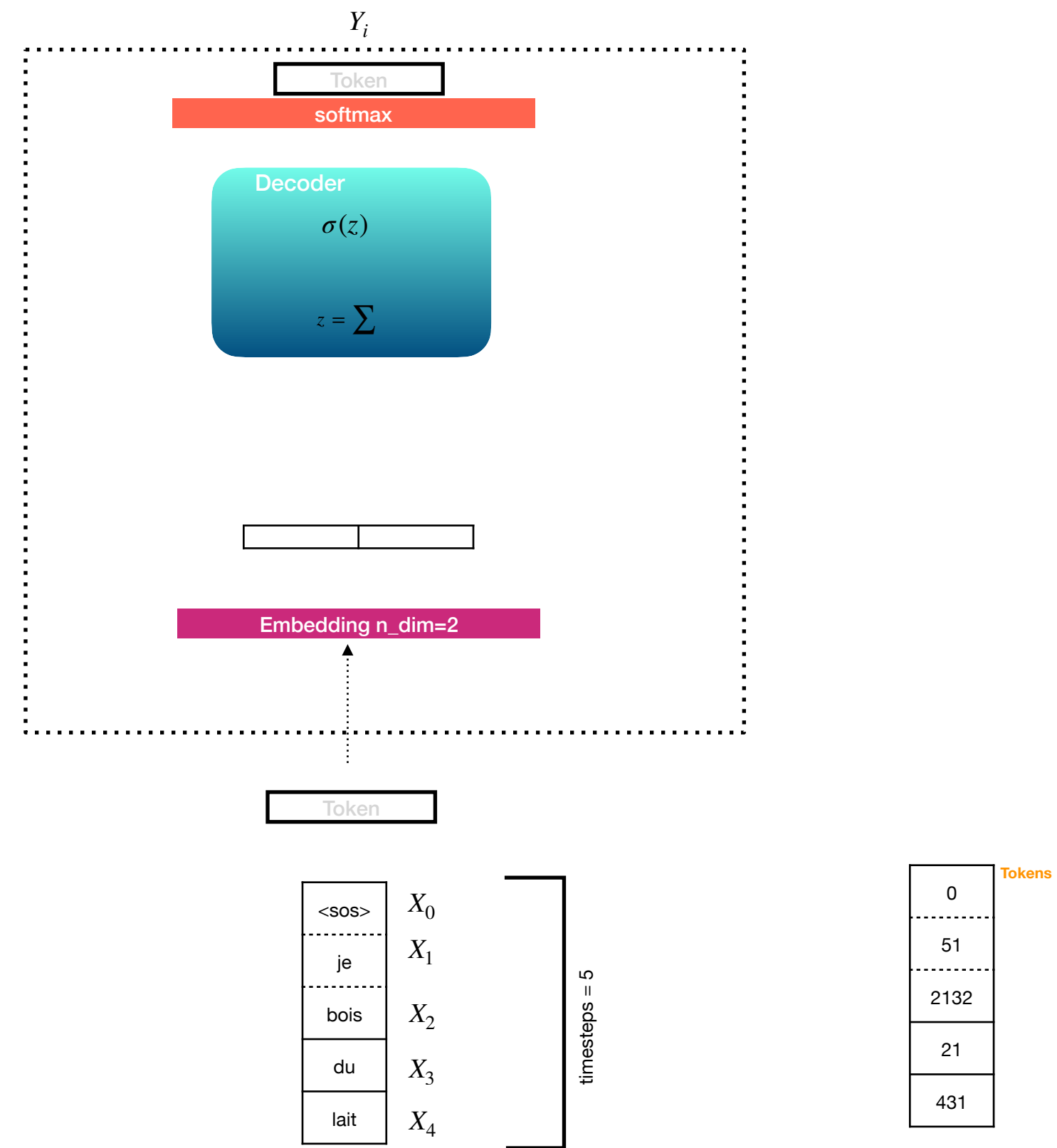
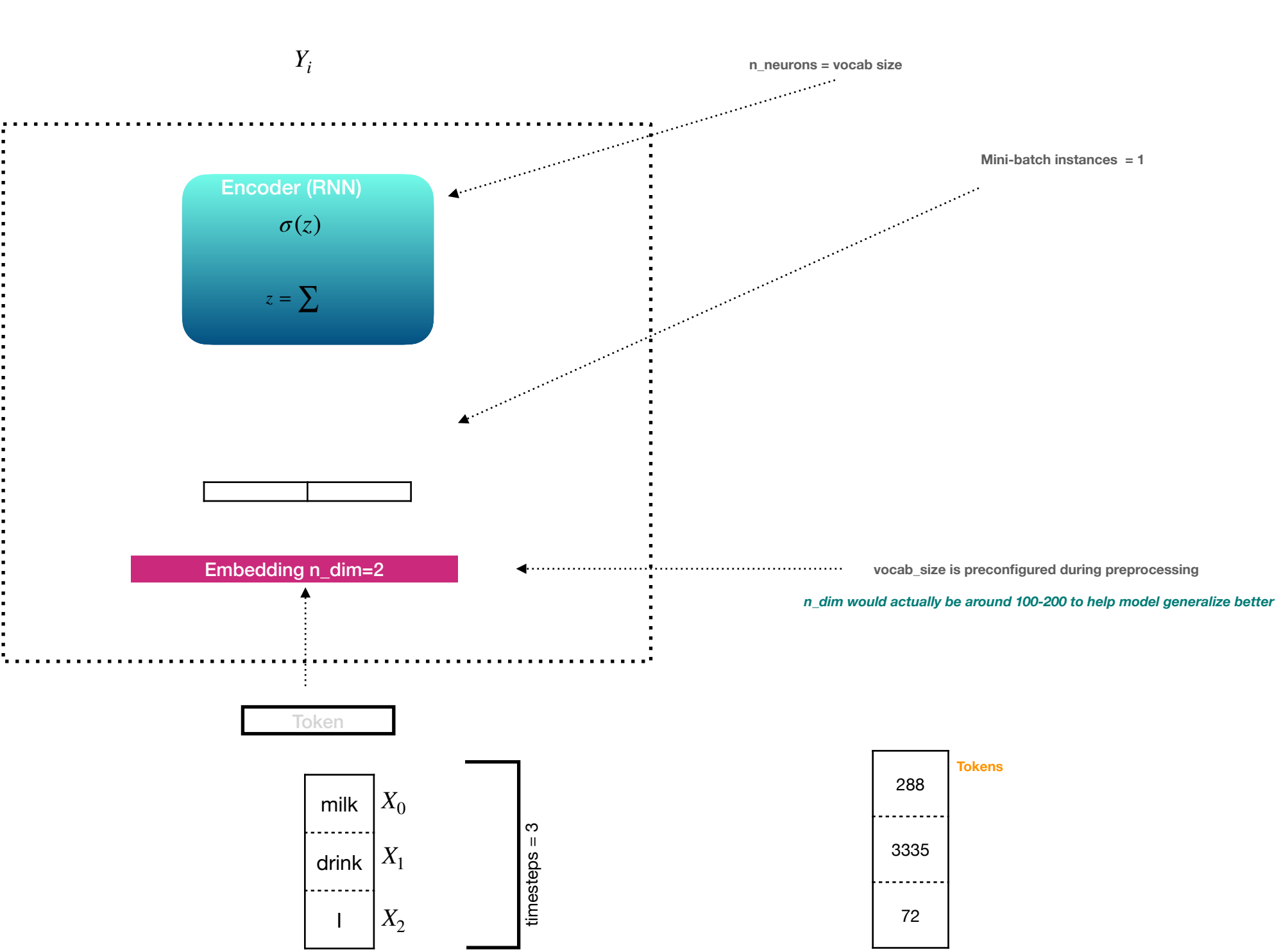


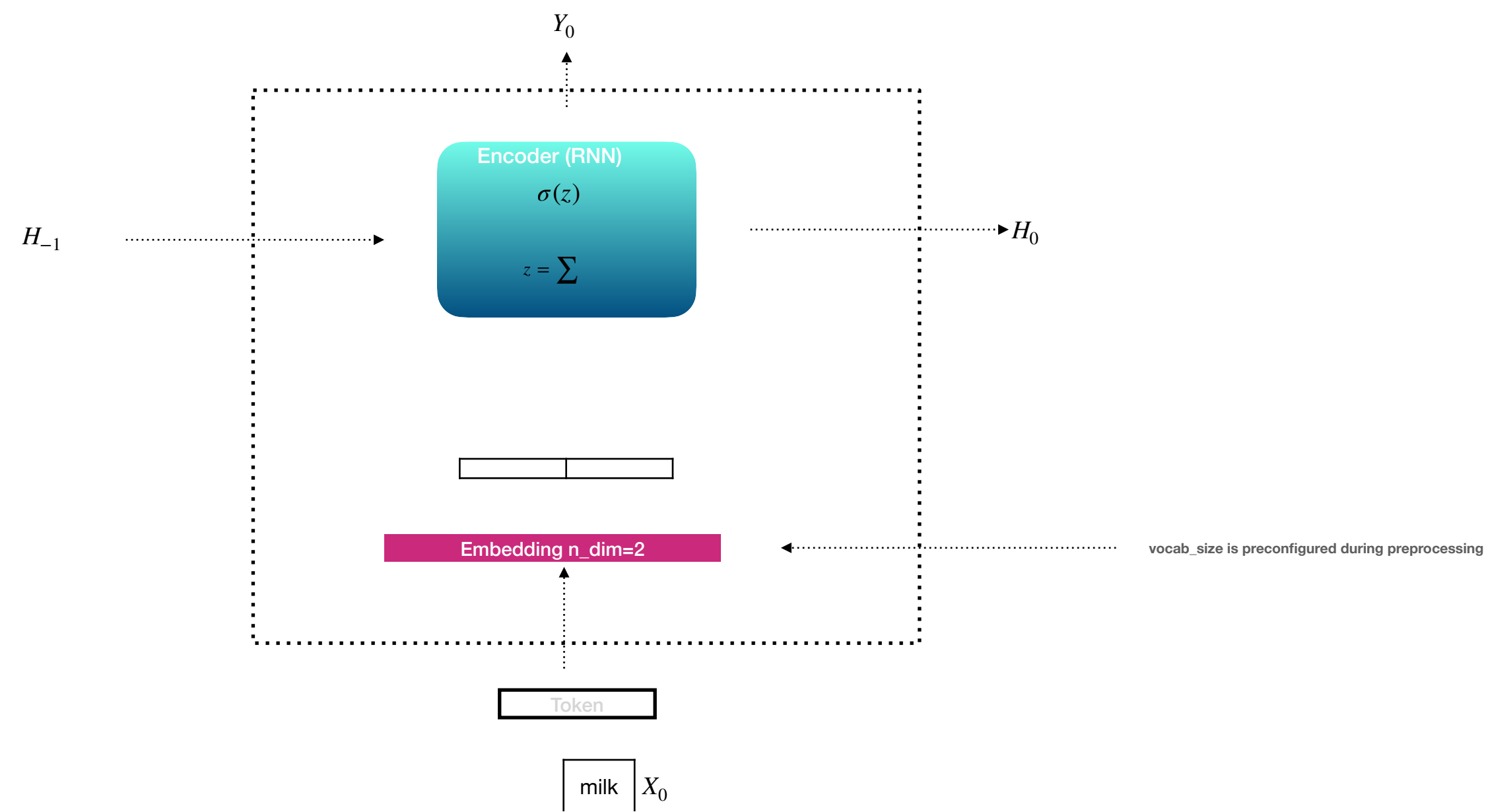
Batching: Stateful RNN

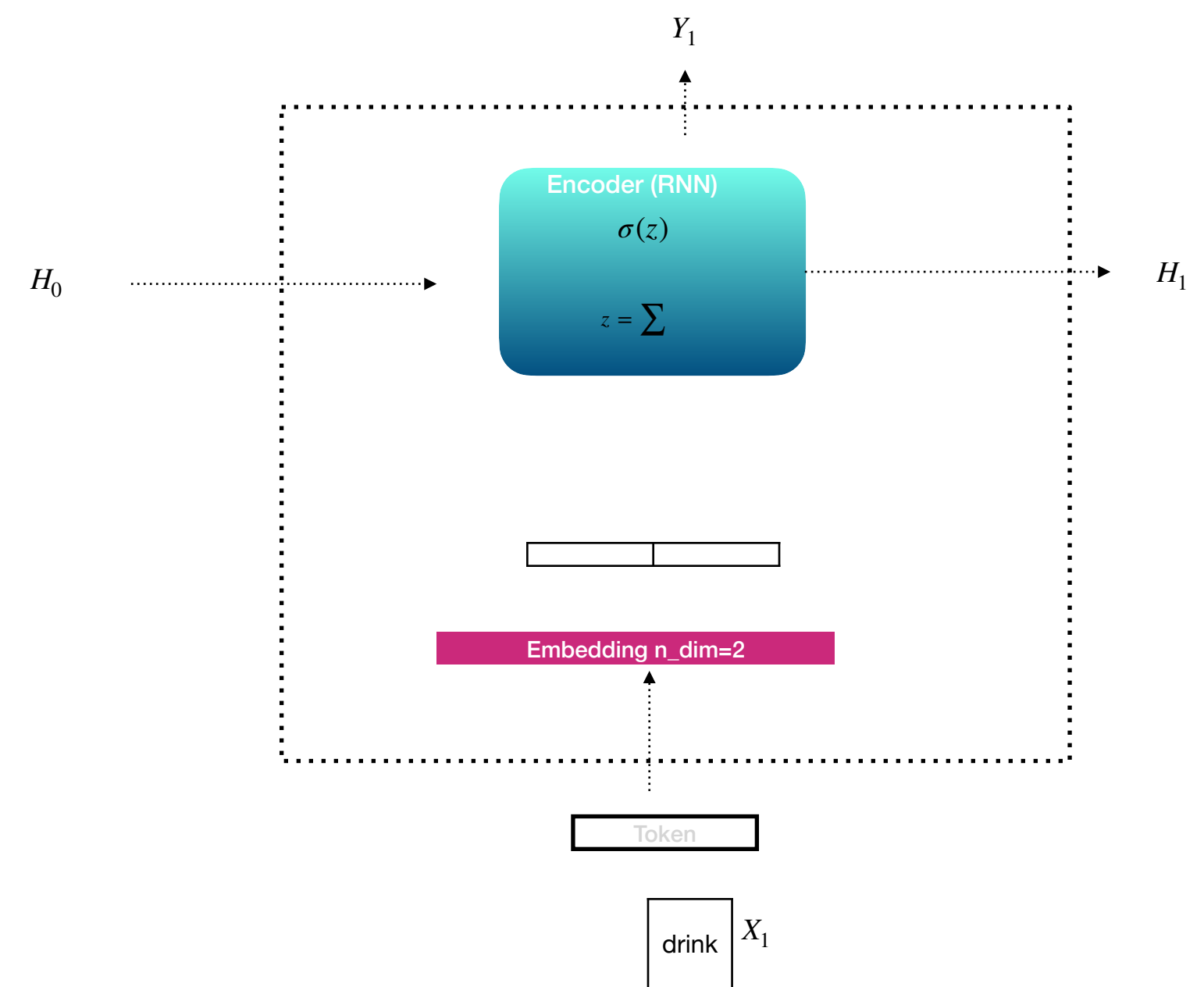


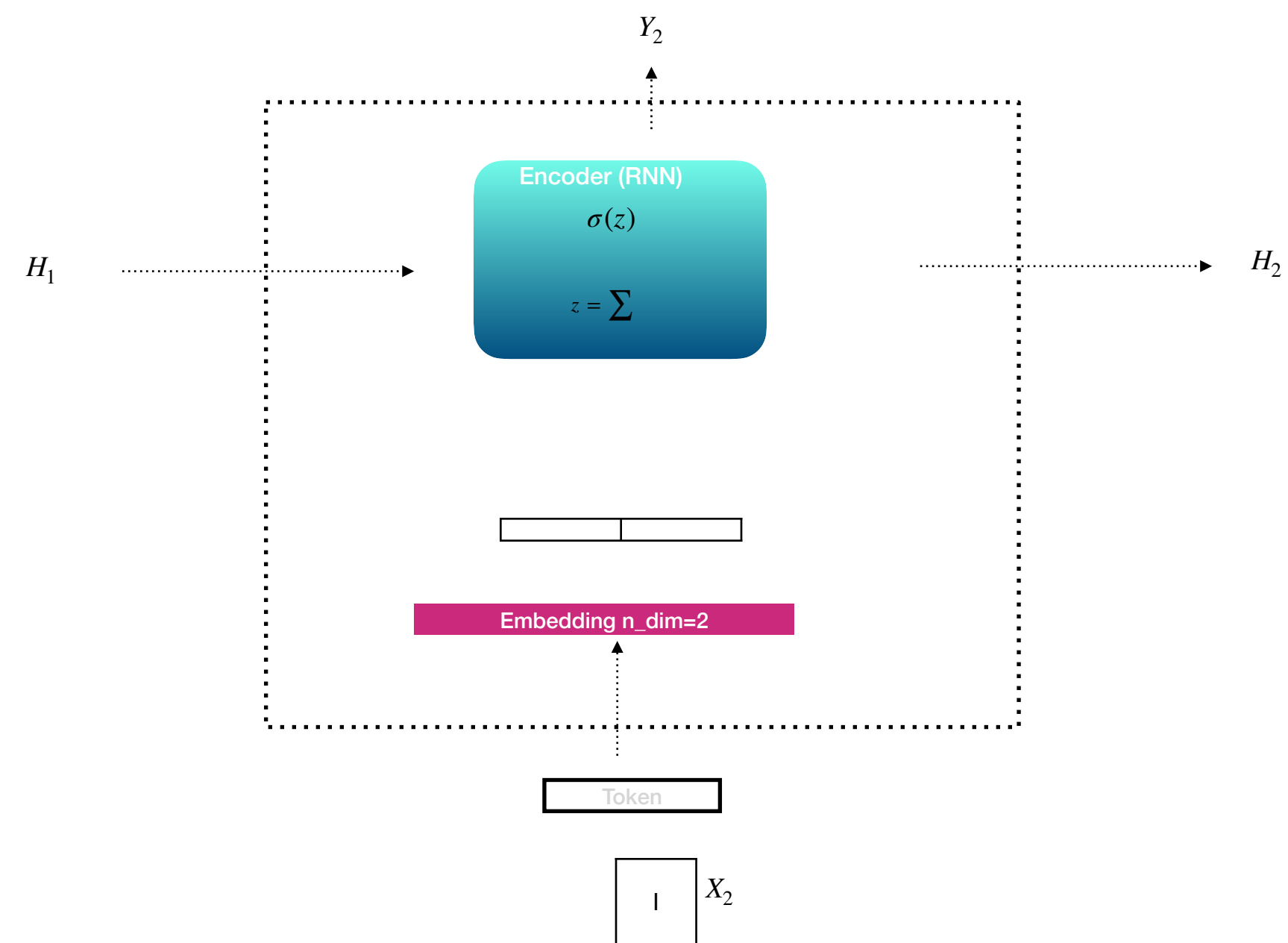
Neural Machine Translation (NMT)



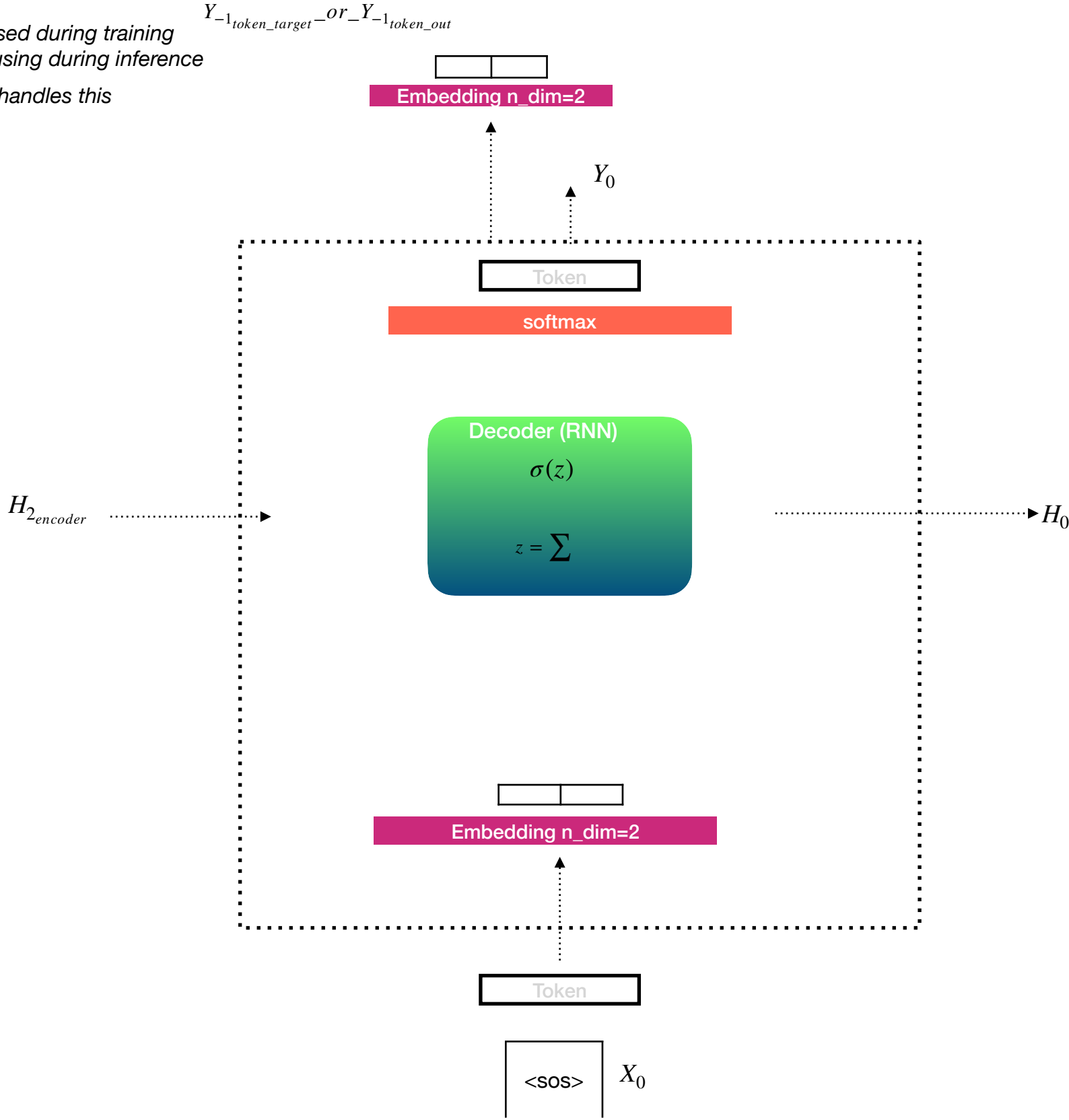


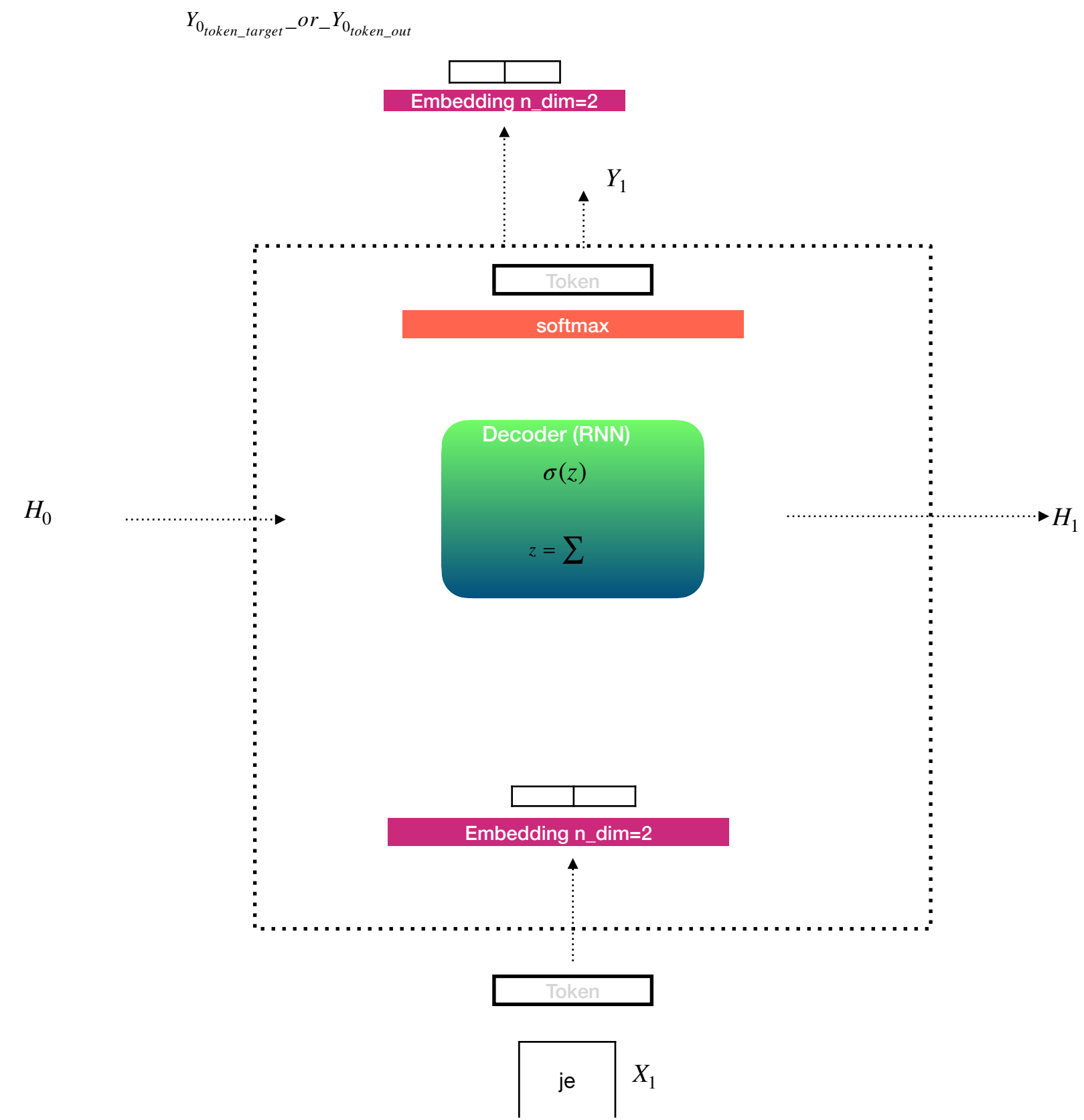


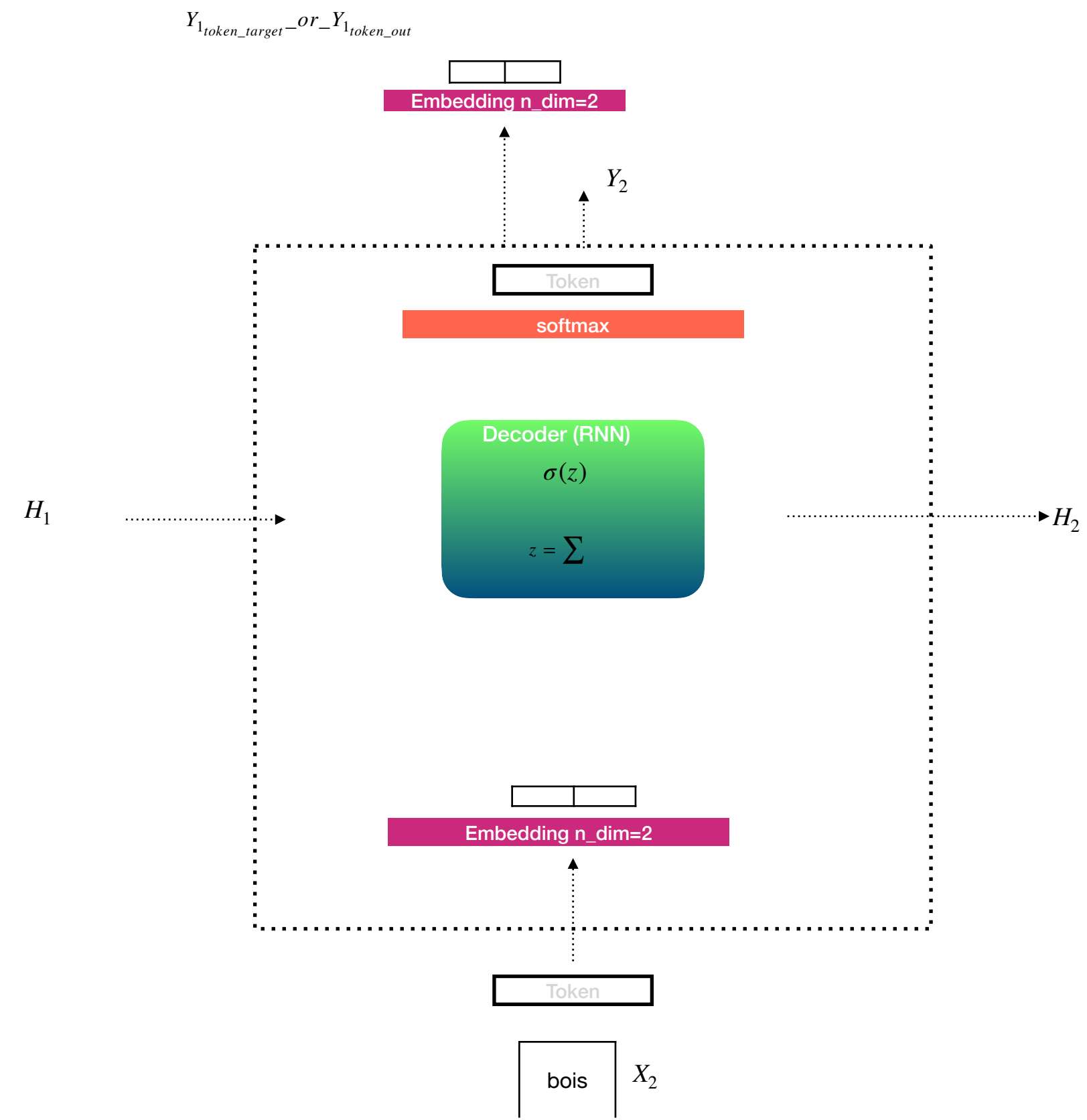


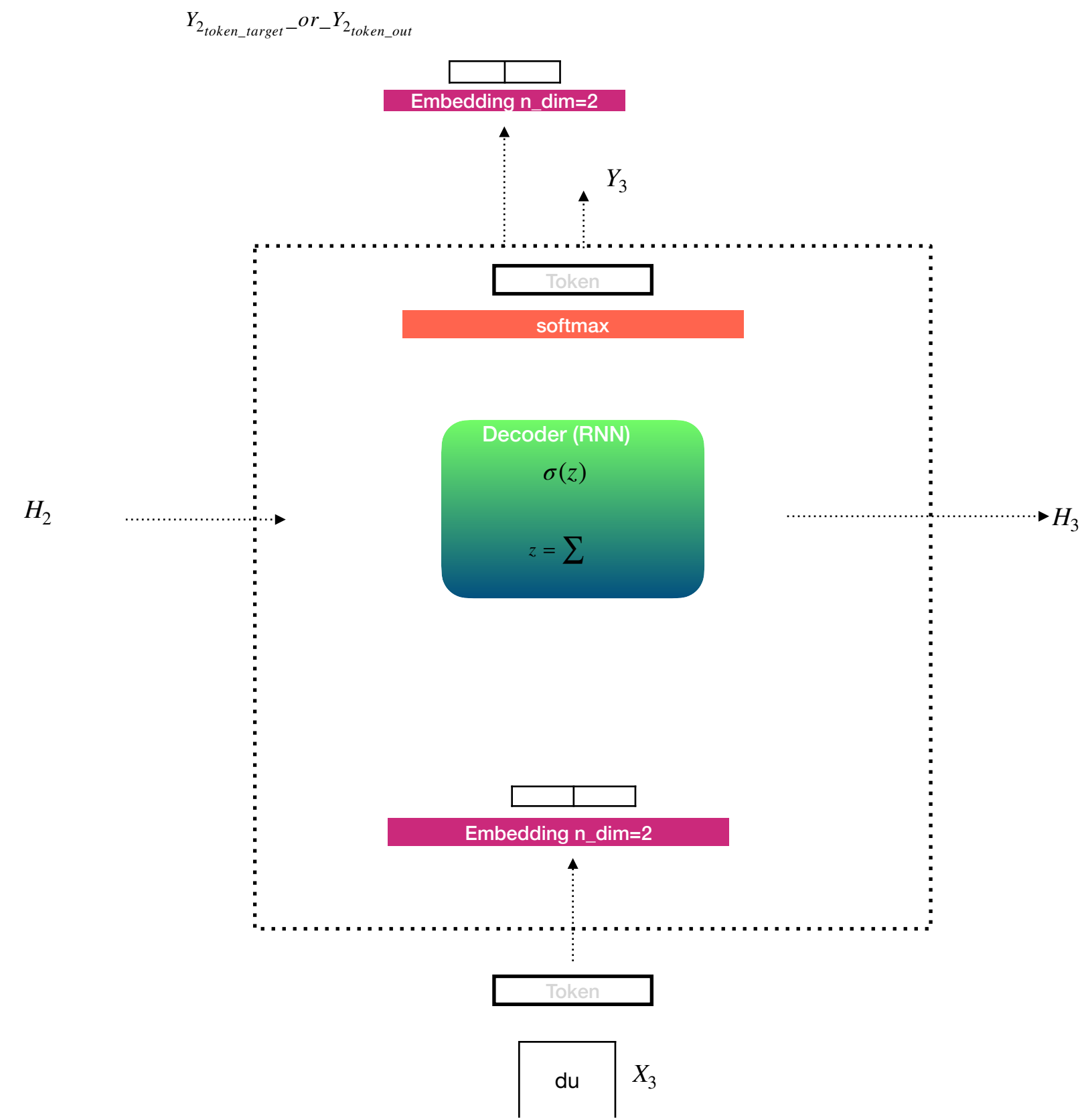


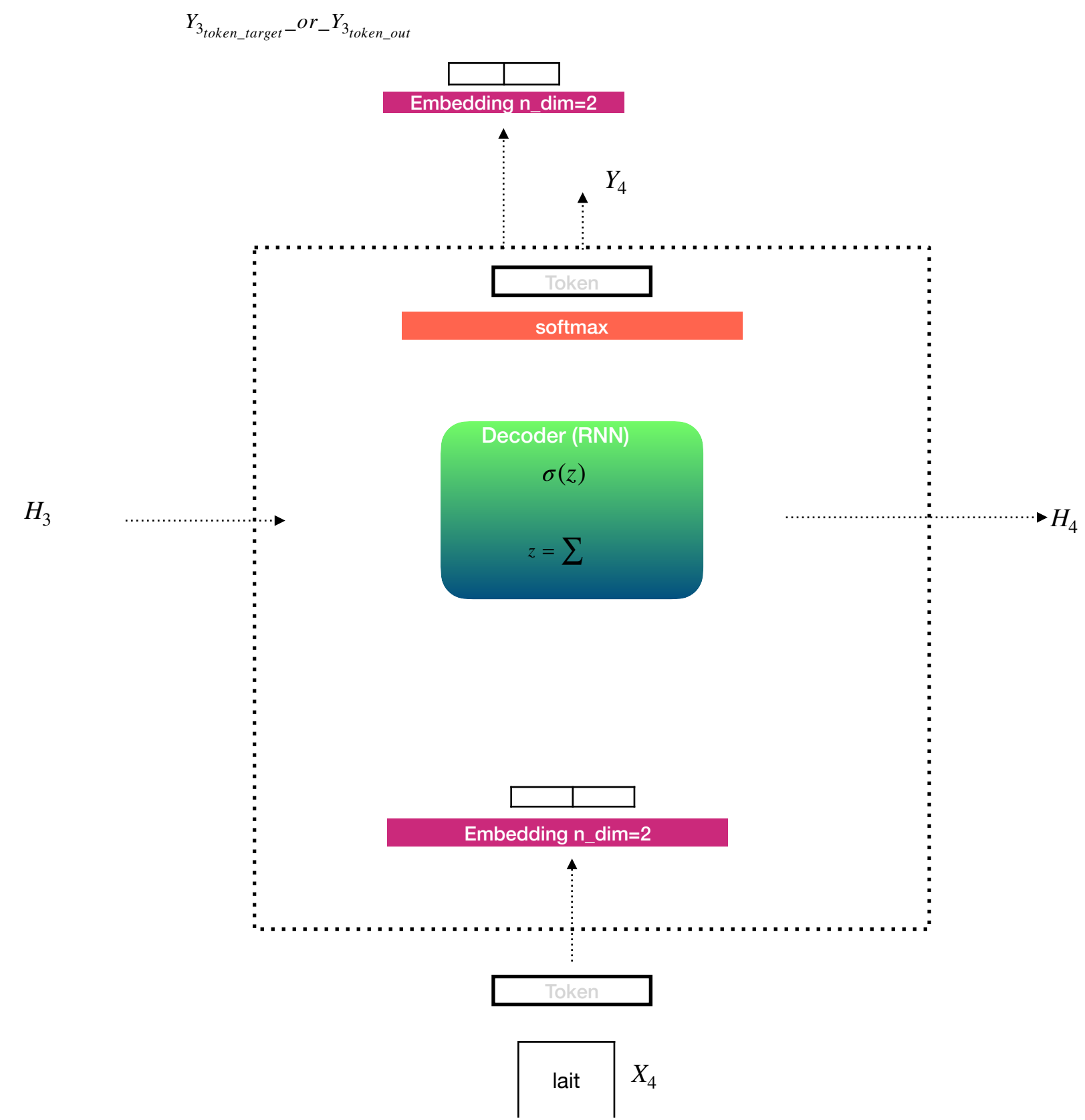
- Prev target token used during training
 - Prev output token using during inference
- Training sampler handles this



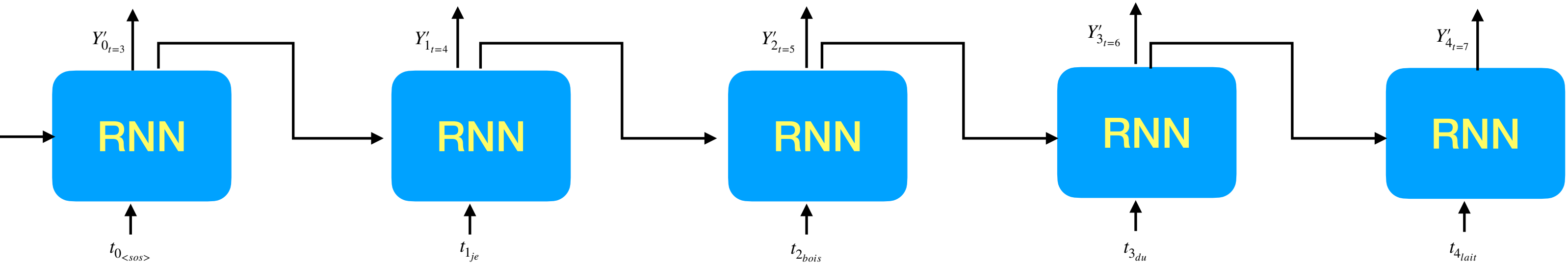
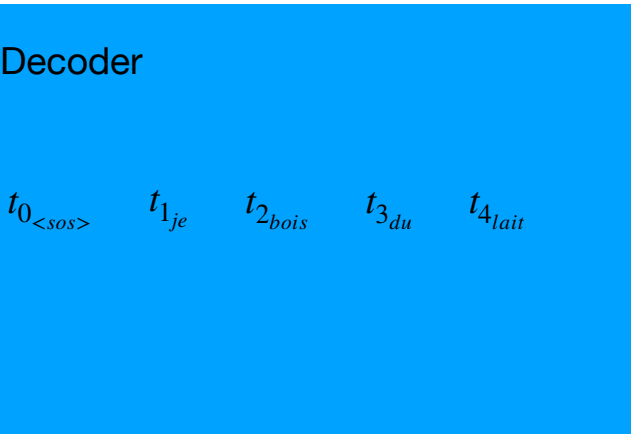
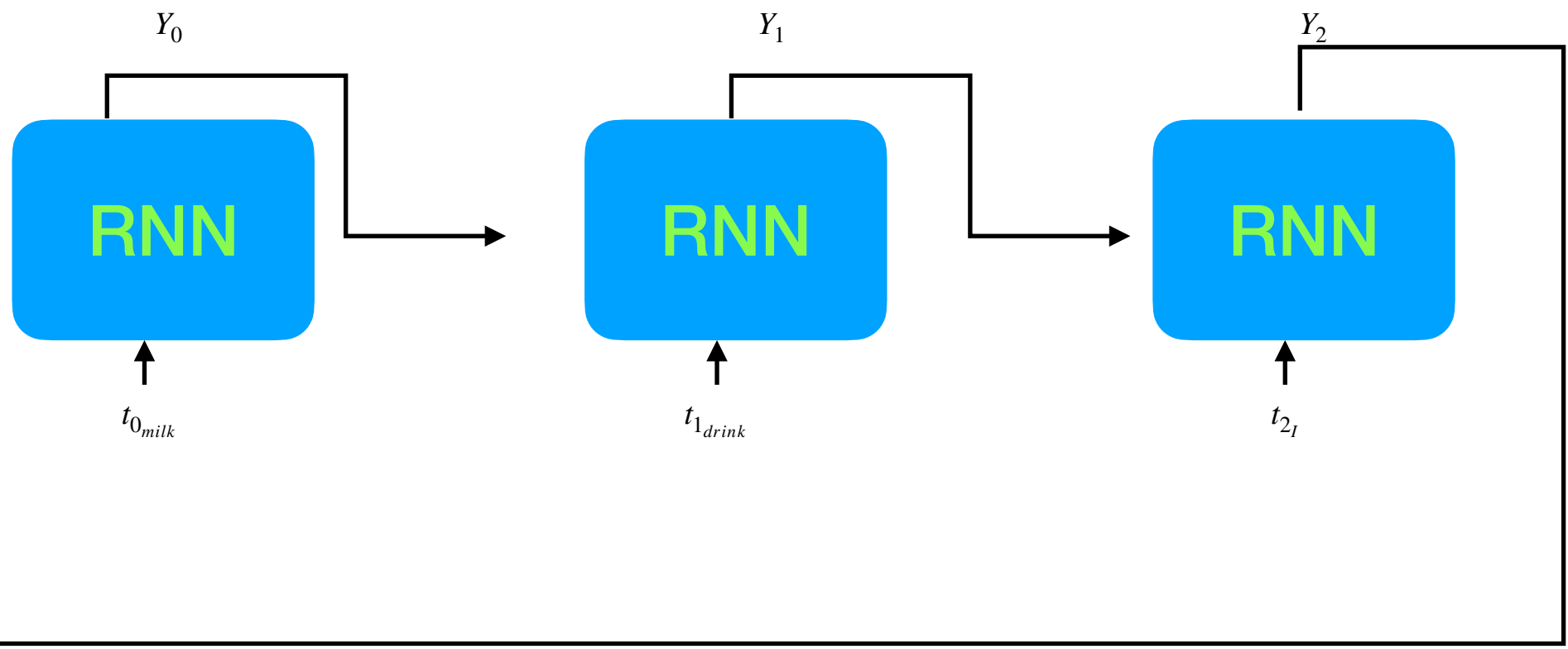
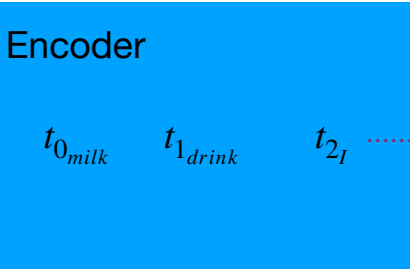






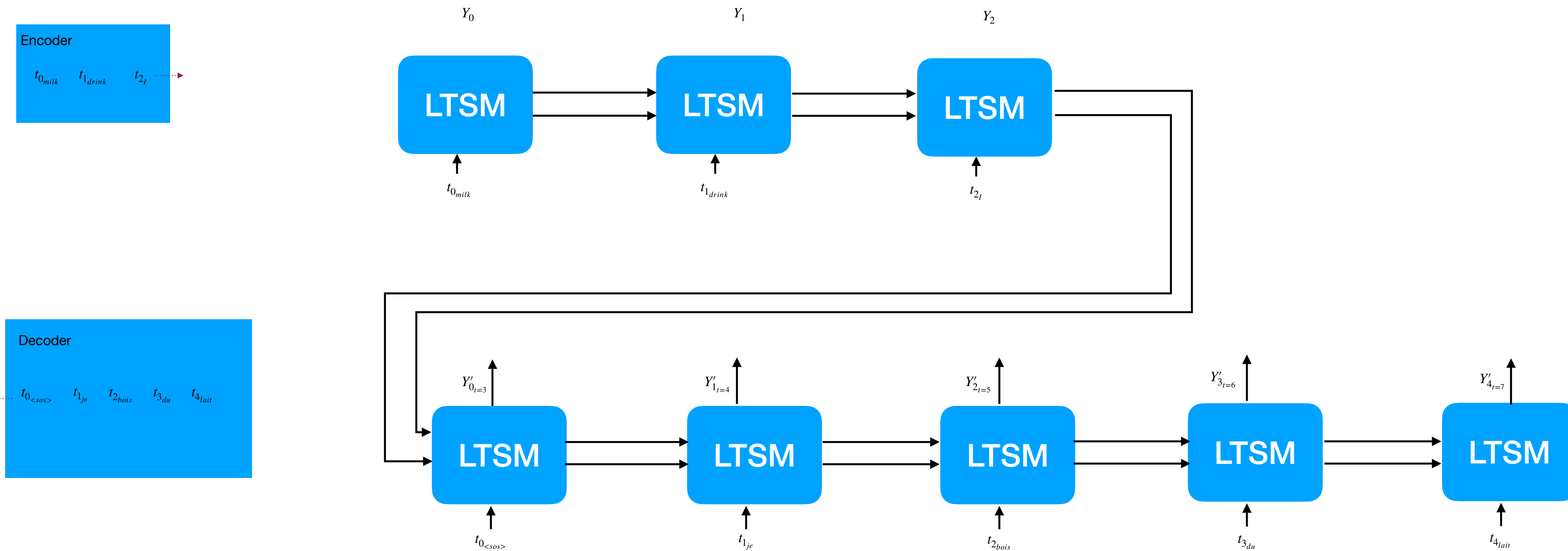


Note: Single RNN being fed samples every time step



Note: Single RNN being fed samples every time step. Decoder waits for Encoder to unravel

Model learns to translate English sentence into French

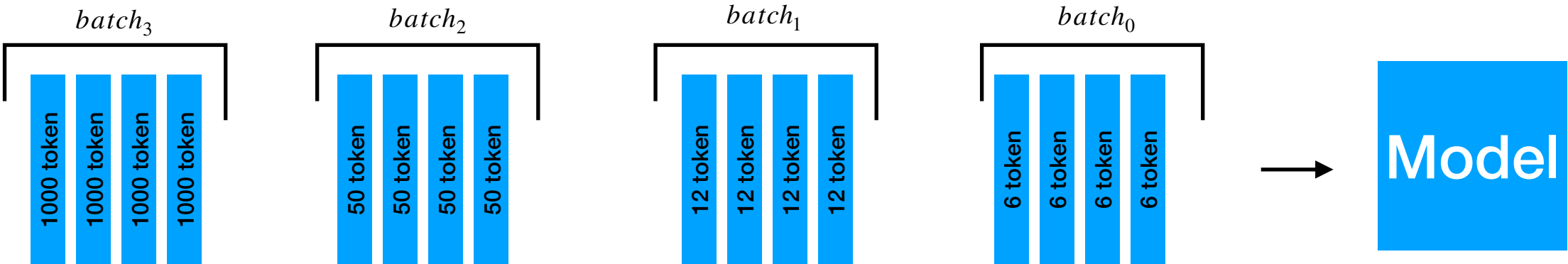


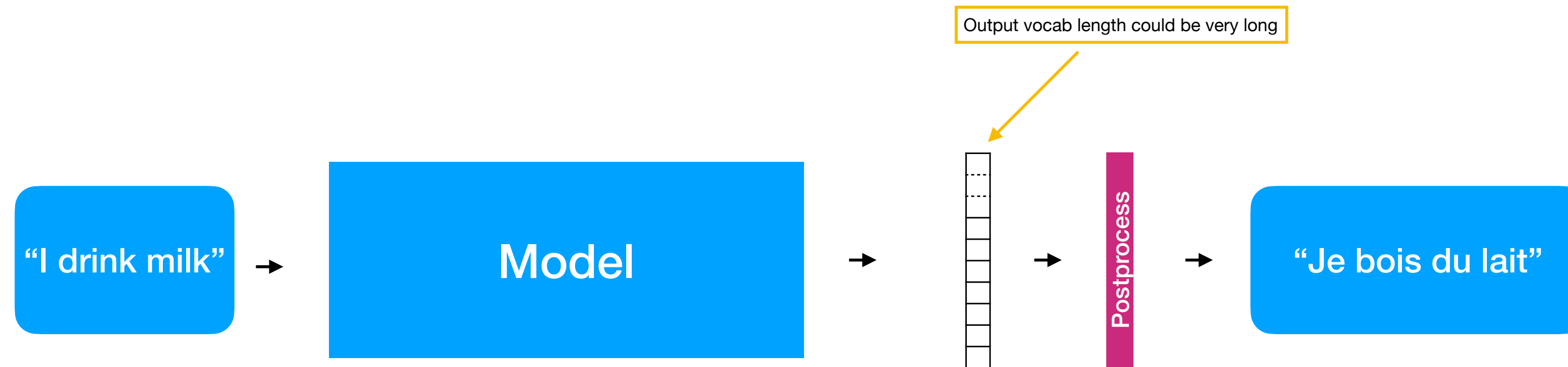
Model learns to translate English sentence into French

Batch sentences of varying length

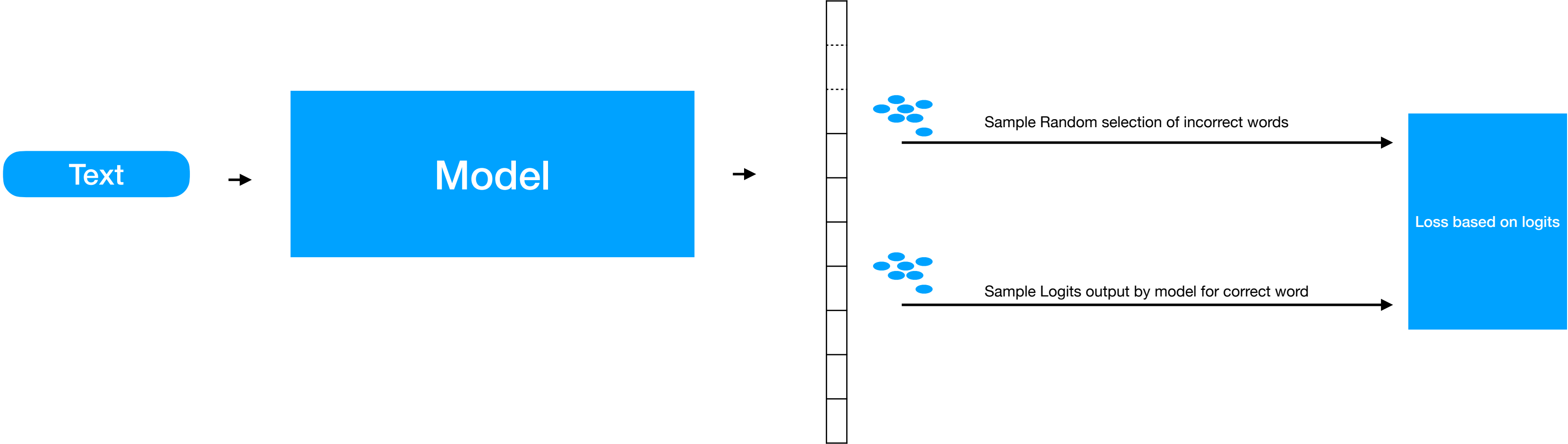


Batch of equal length Tensors are accepted as inputs to model

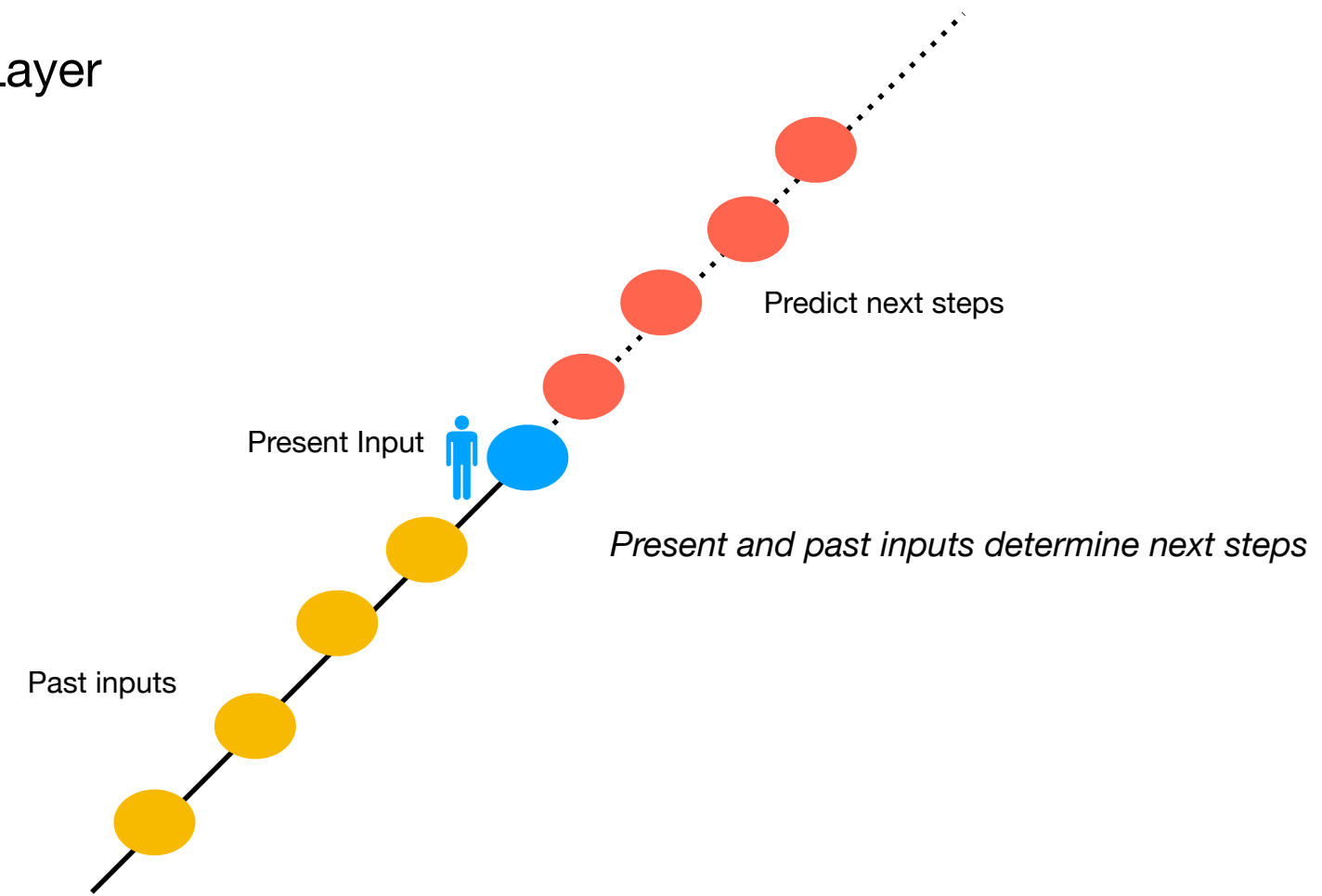




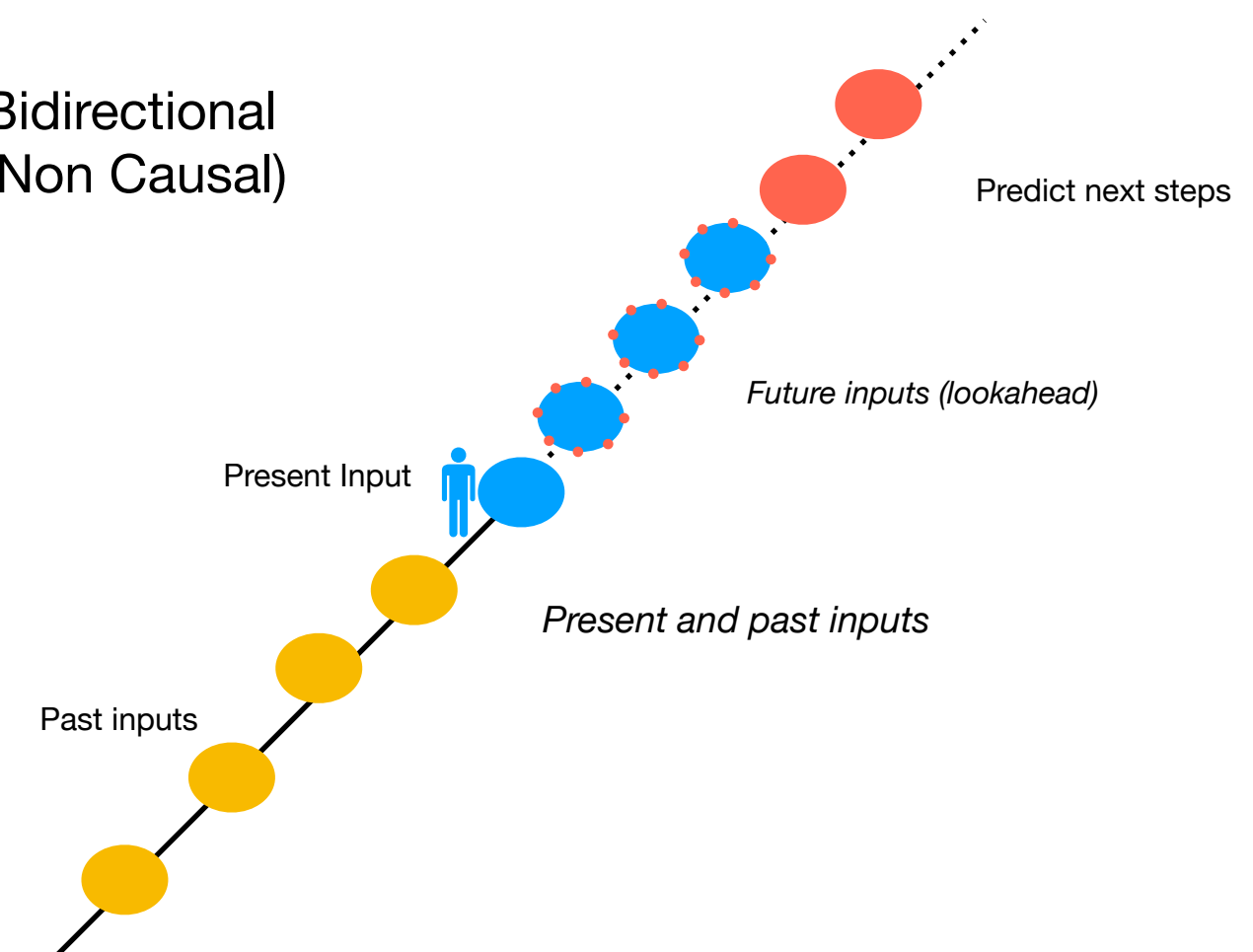
Handling large output vocab



Regular RNN Layer
(Causal)

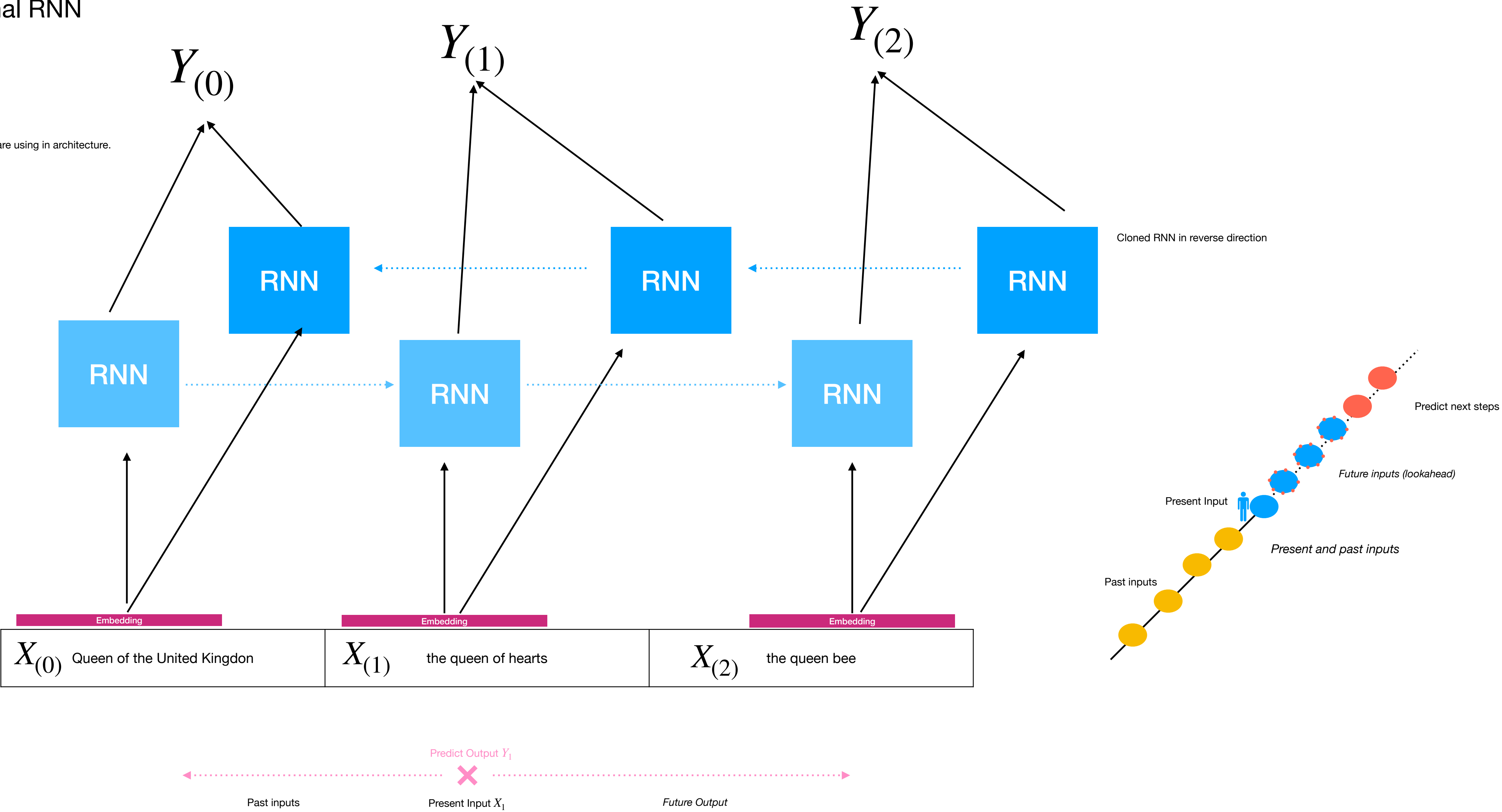


Bidirectional
(Non Causal)

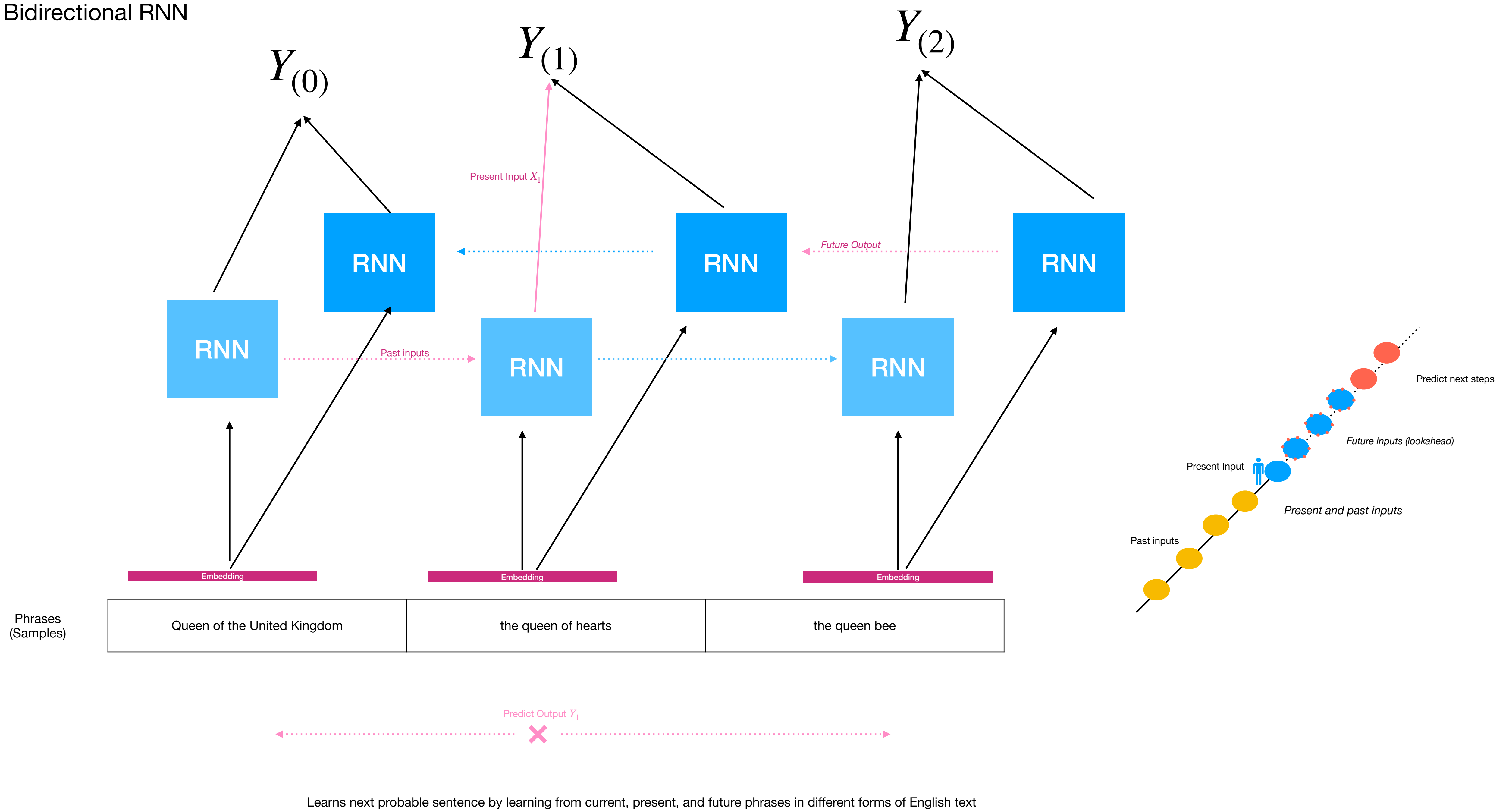


Bidirectional RNN

Note: Only two RNN are using in architecture.
Recursive unit!



Bidirectional RNN



Beam Search (n=3)

“Comment vas-tu”



Encoder-Decoder

T0

How (75 %)

What (3 %)

You (1 %)

10,000 probabilities of each word

Decoder Estimates probability for each word in vocab, only top three make the short list

Decoder Clone

Decoder Clone

Decoder Clone

Estimate next three

T1

How, Will (36 %)

How, Are (32 %)

How, Do (16 %)

10,000 probabilities of next word

Probability of sentence it completes (i.e. two words)

(Conditional probability of word) * (probability of sentence it completes)

How, will -> 36 * 75 = **0.27**

How, are -> 32 * 75 = **0.24**

How, do -> 16 * 75 = **0.12**

Decoder Estimates probability for each word in vocab, only top three make the short list, **given the sentence starts with how**

Estimate next three

What, are (50%)

?

?

10,000 probabilities of next word

Probability of sentence it completes

What, are -> 3 * 50 = 0.015

?, ? -> -

?, ? -> -

Decoder Estimates probability for each word in vocab, only top three make the short list, **given the sentence starts with what**

Estimate next three

?

?

?

10,000 probabilities of next word

Probability of sentence it completes

?, ? -> -

?, ? -> -

?, ? -> -

Decoder Estimates probability for each word in vocab, only top three make the short list, **given the sentence starts with you**

Decoder Clone

Decoder Clone

Decoder Clone

How, Will

How, Are

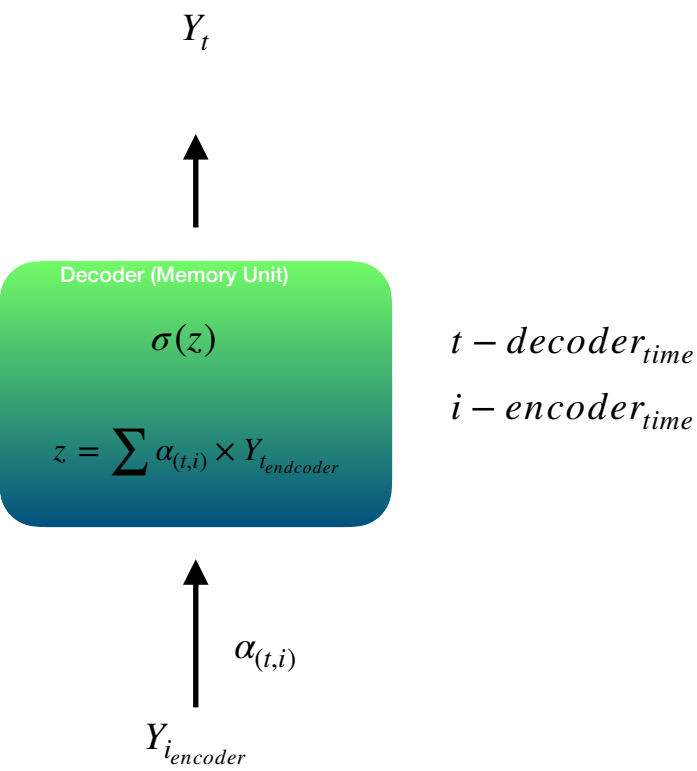
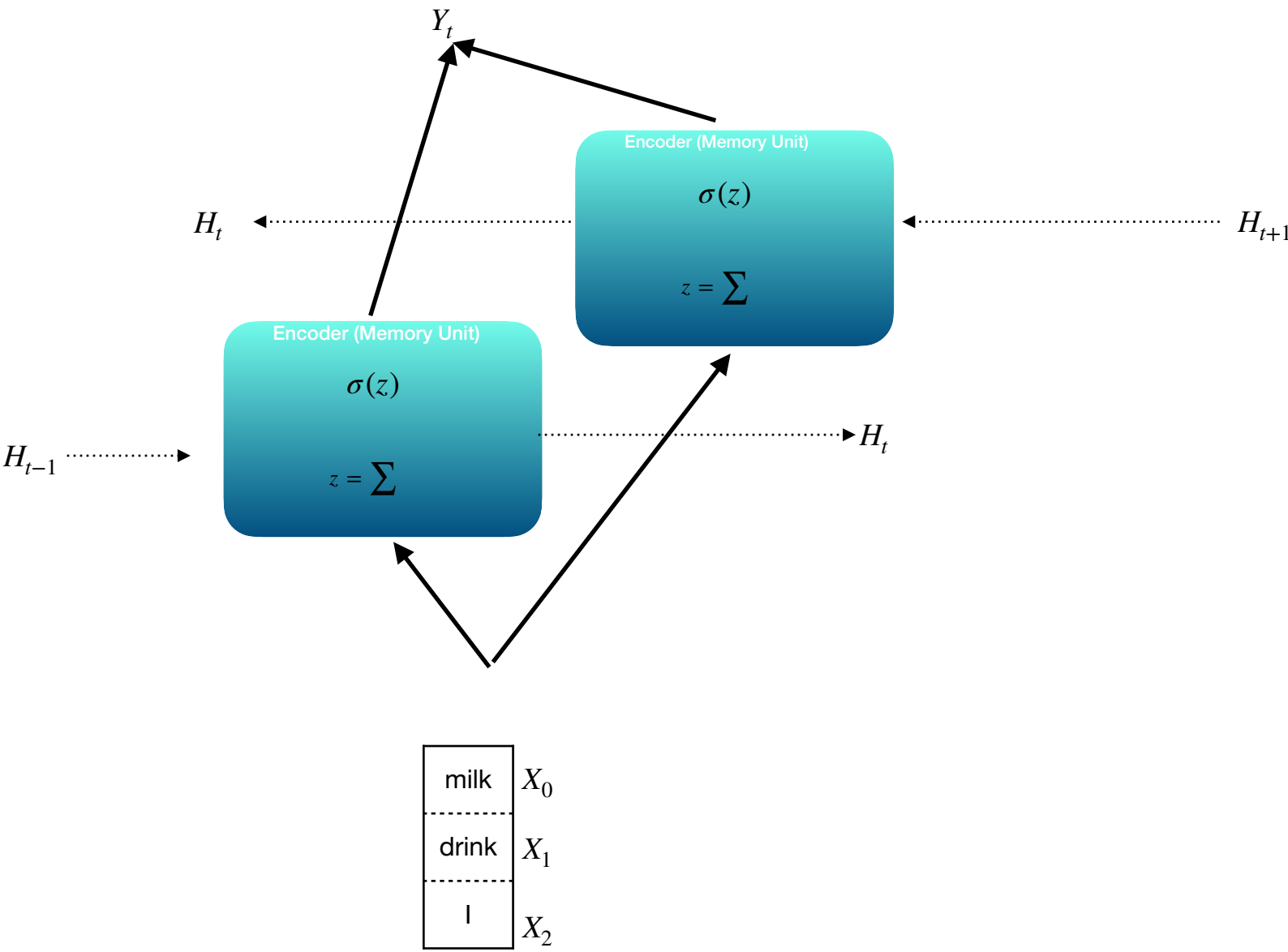
How, Do

Estimate next three

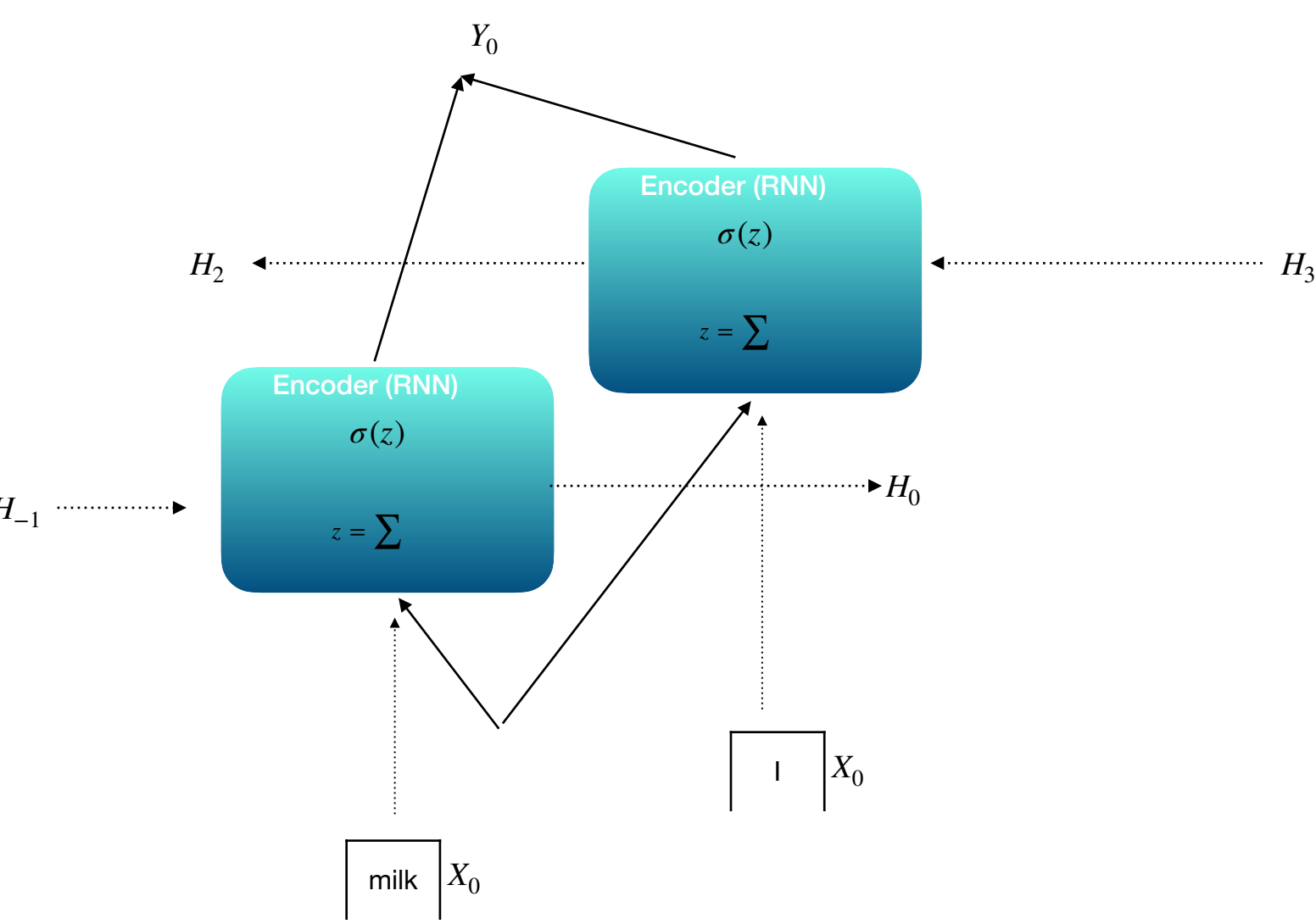
Estimate next three

Estimate next three

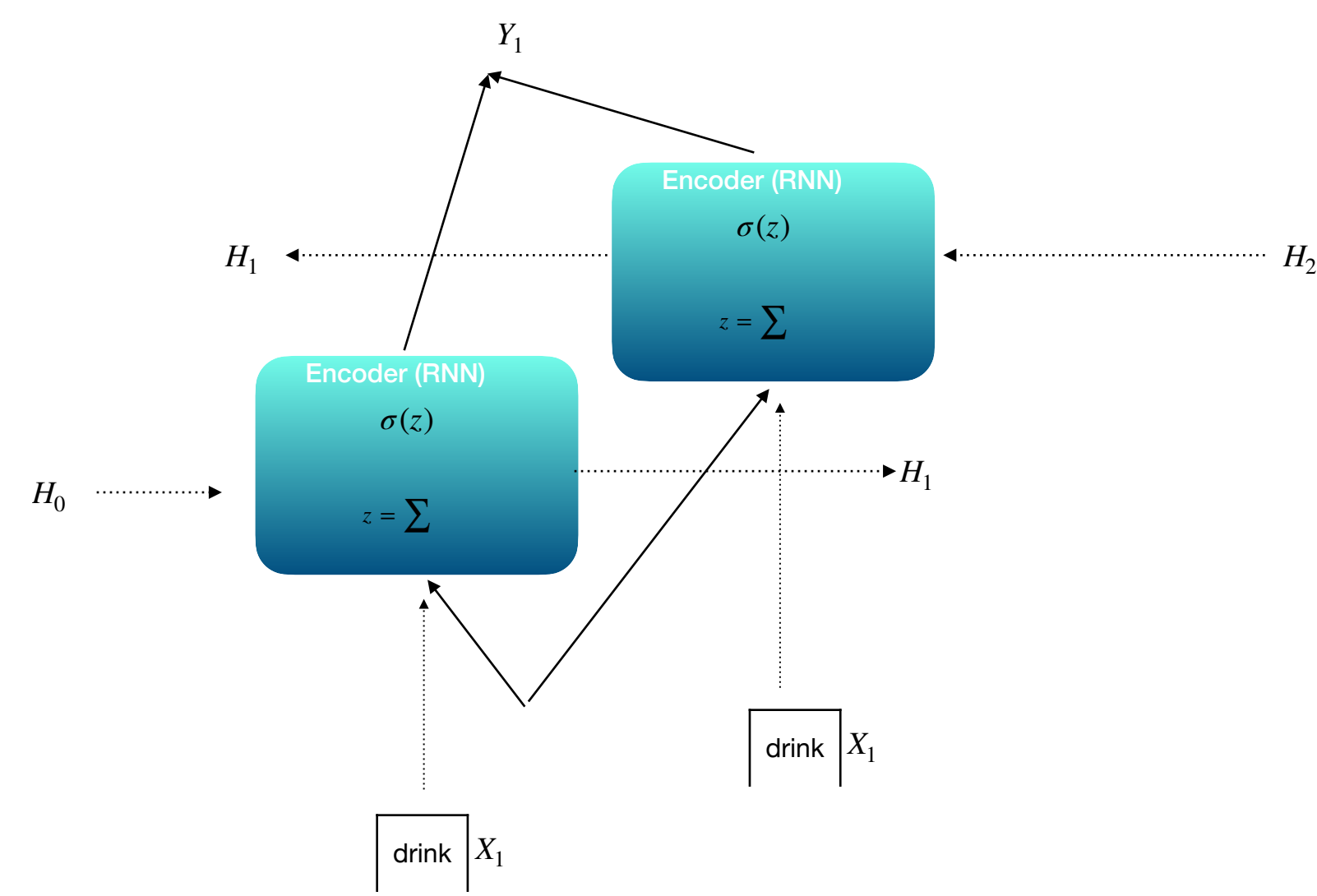
Attention Mechanisms



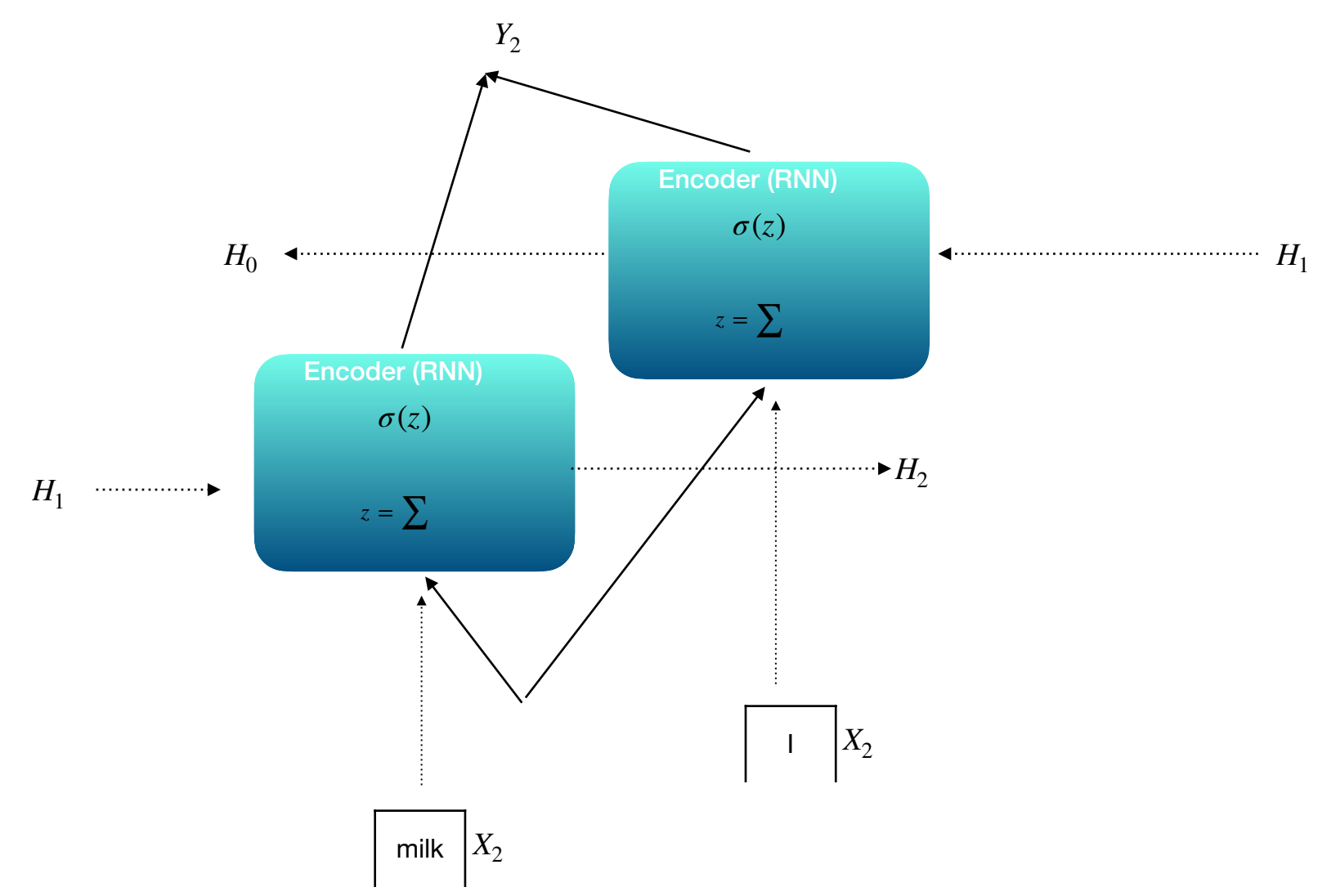
Attention Mechanisms



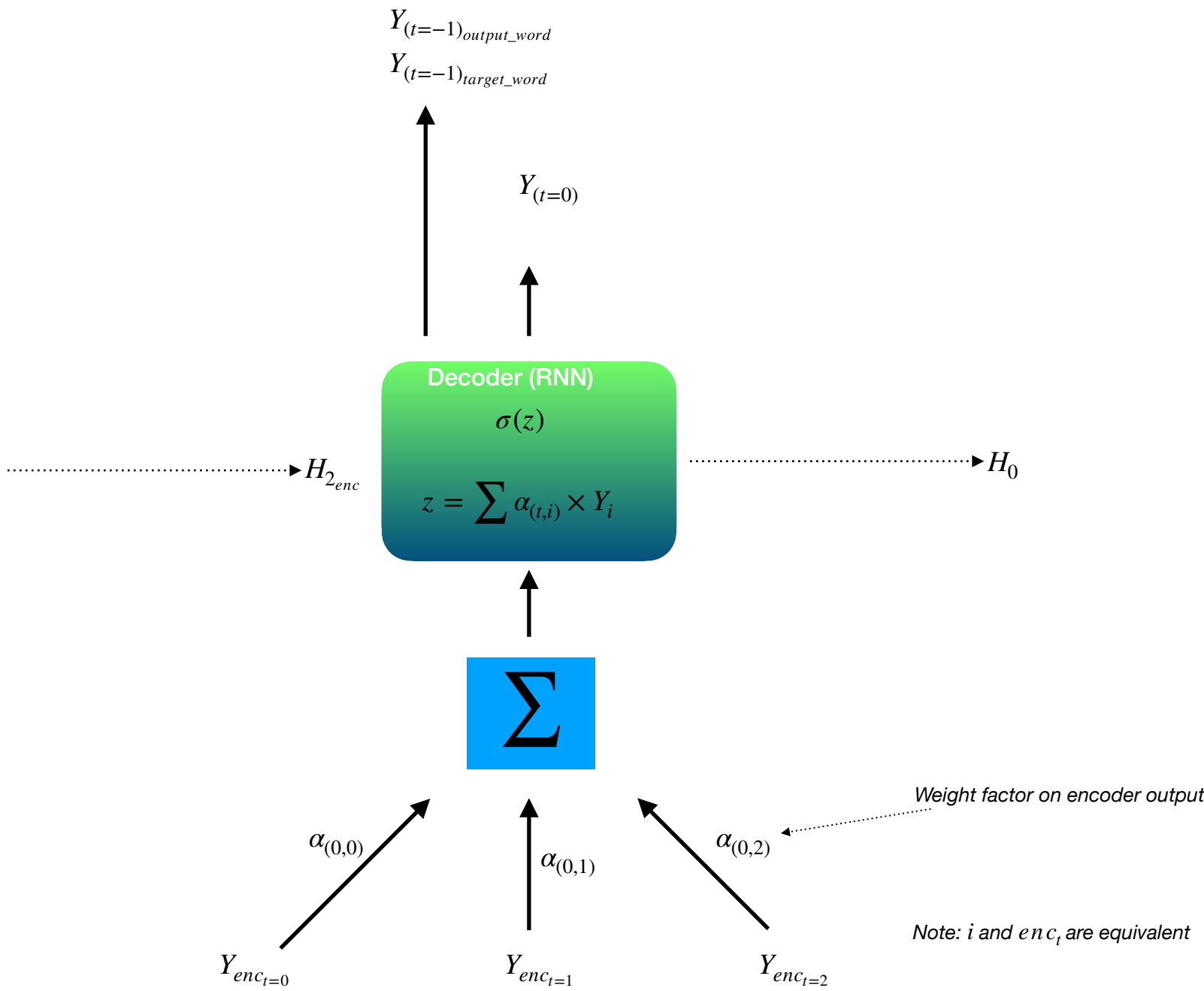
Attention Mechanisms



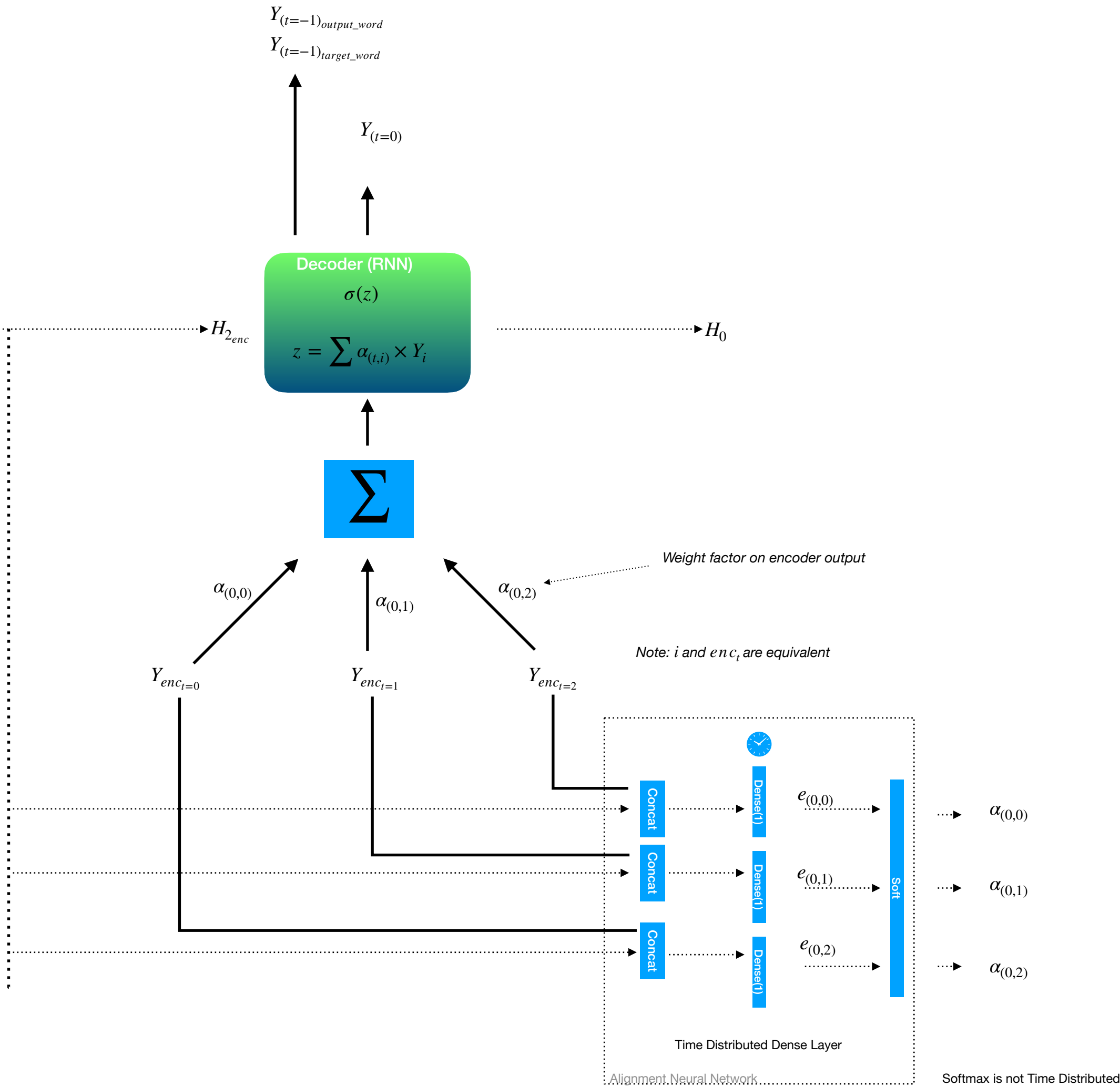
Attention Mechanisms



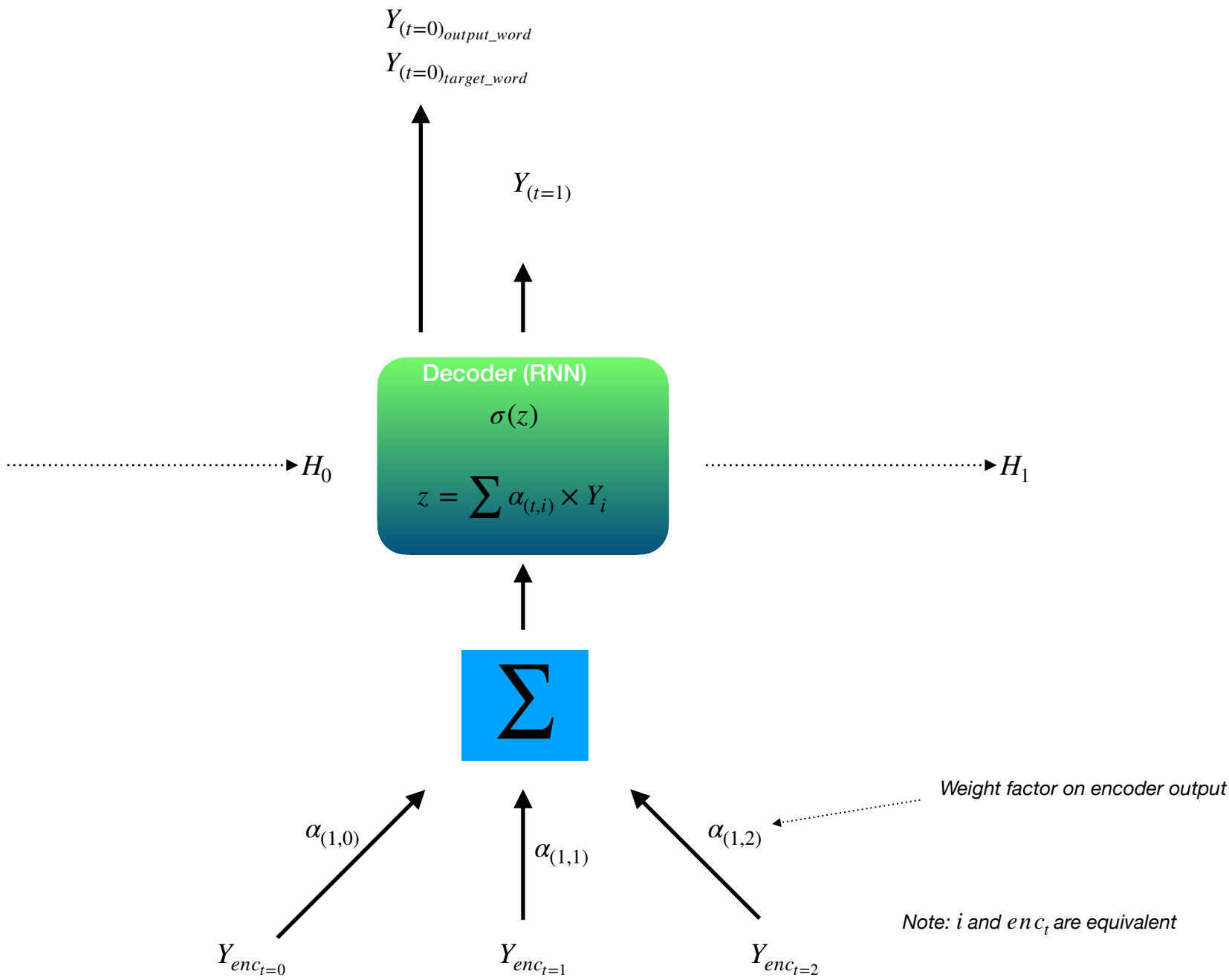
Attention Mechanisms



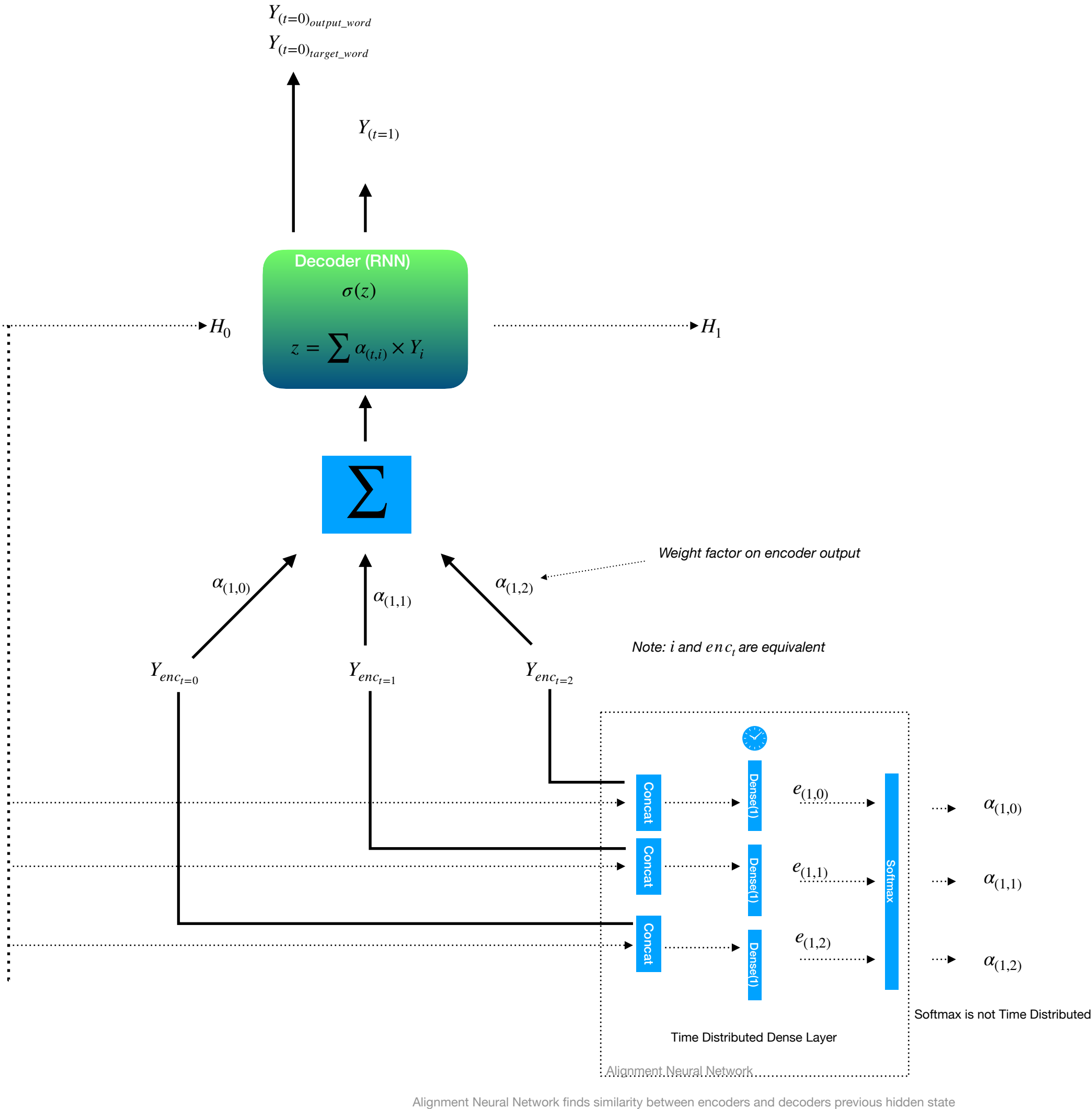
Attention Mechanisms:
Calculating Weights



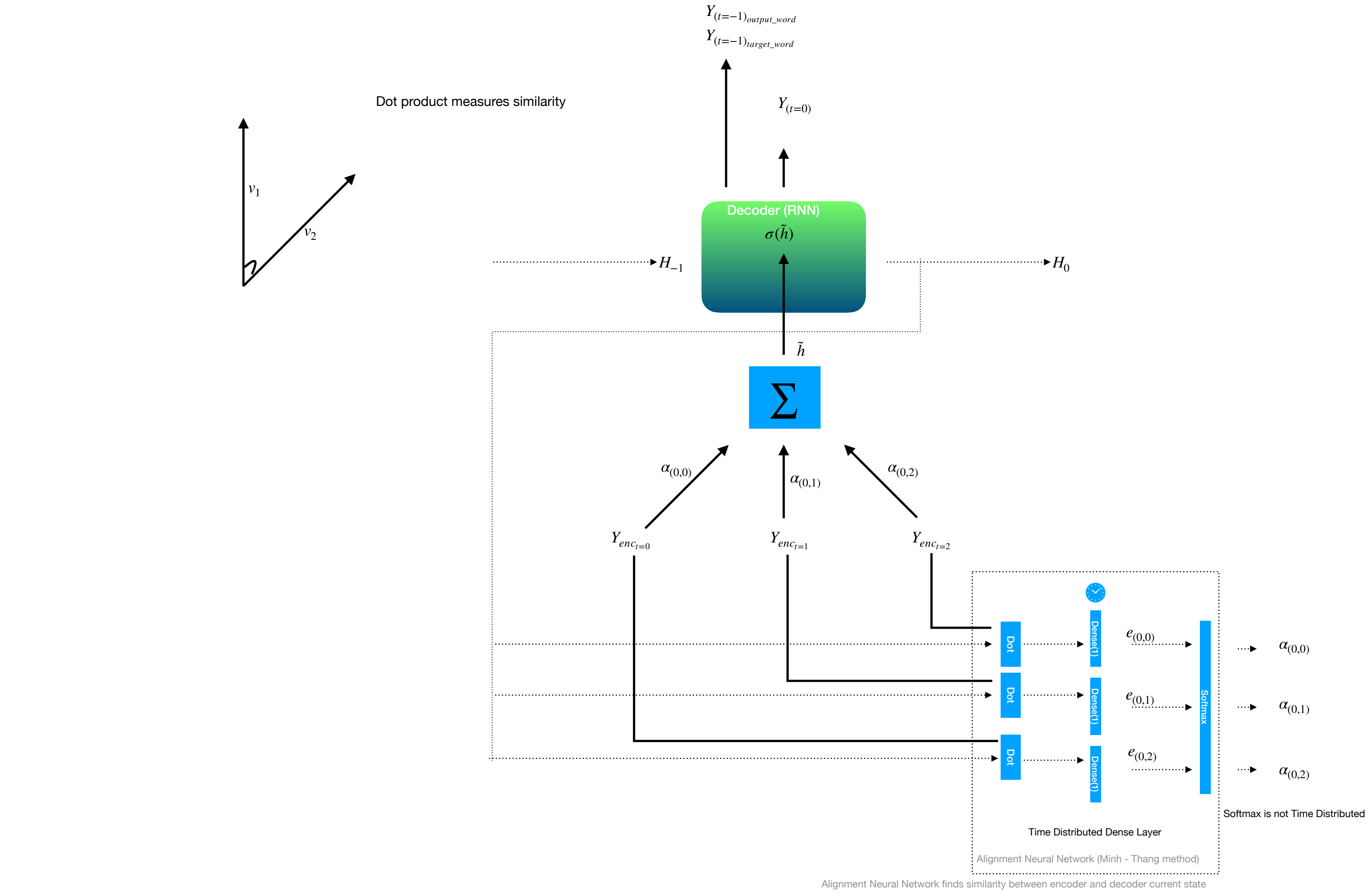
Attention Mechanisms



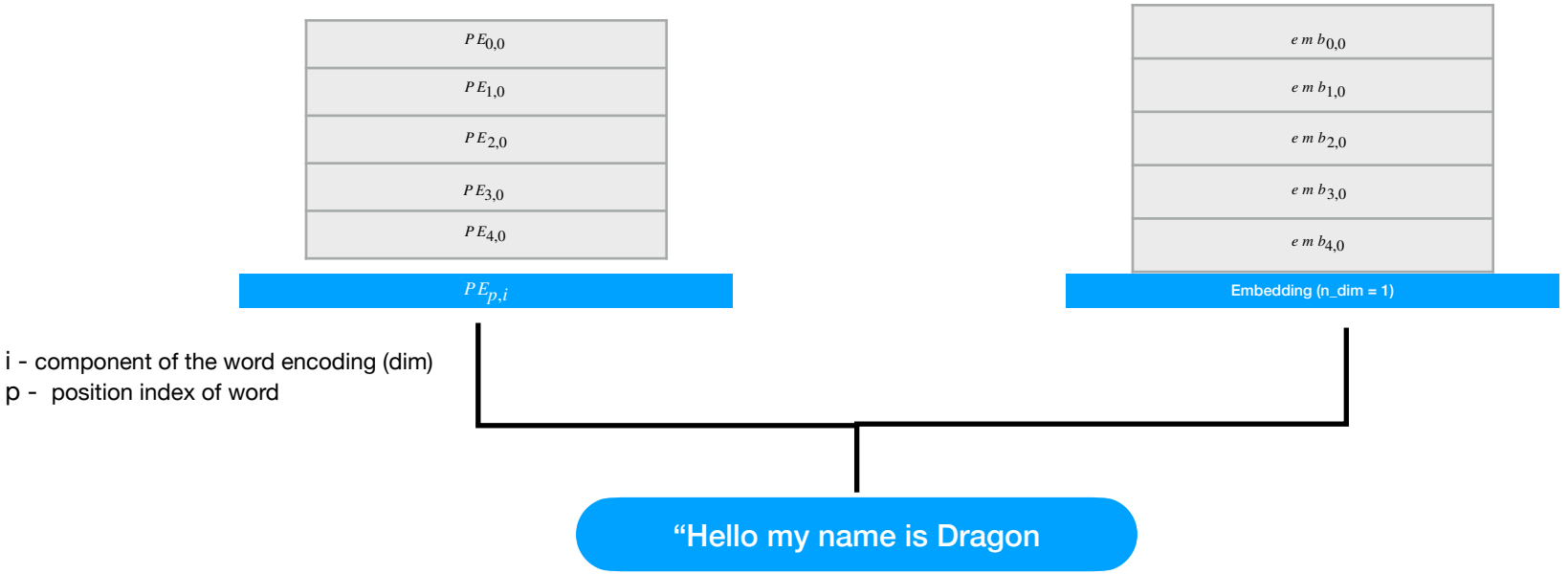
Attention Mechanisms:
Calculating Weights



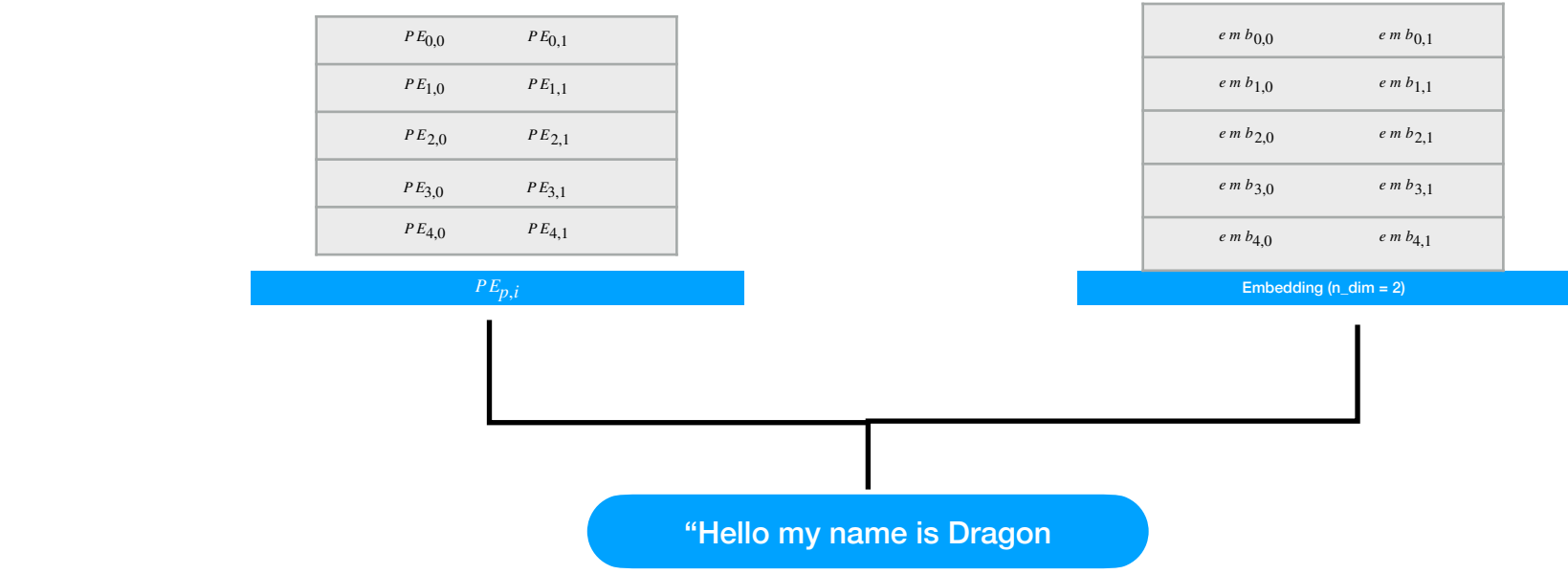
Attention Mechanisms
(Method 2)



Positional Encodings

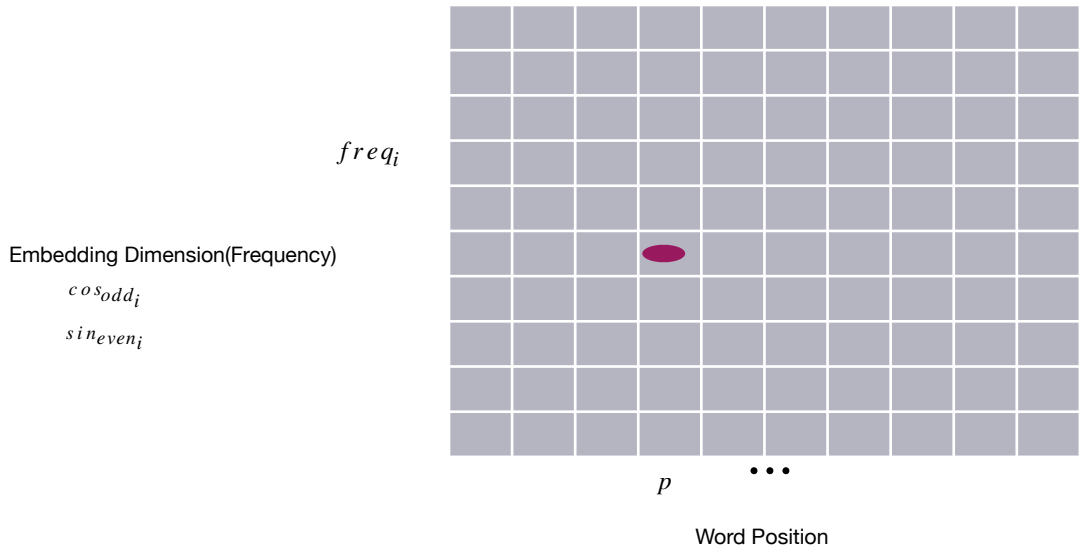


$$PE_{p,i} = \begin{cases} \sin(p/10000^{i/d}) & \text{if } i \text{ is even} \\ \cos(p/10000^{(i-1)/d}) & \text{if } i \text{ is odd} \end{cases}$$

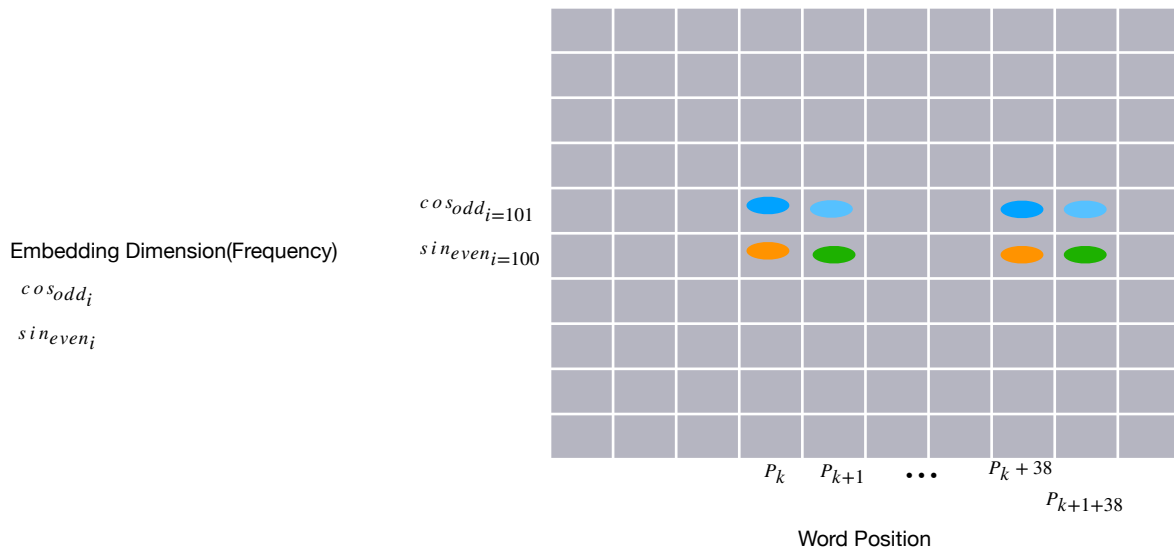
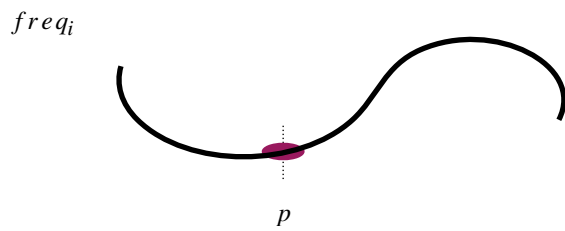


Positional Encodings

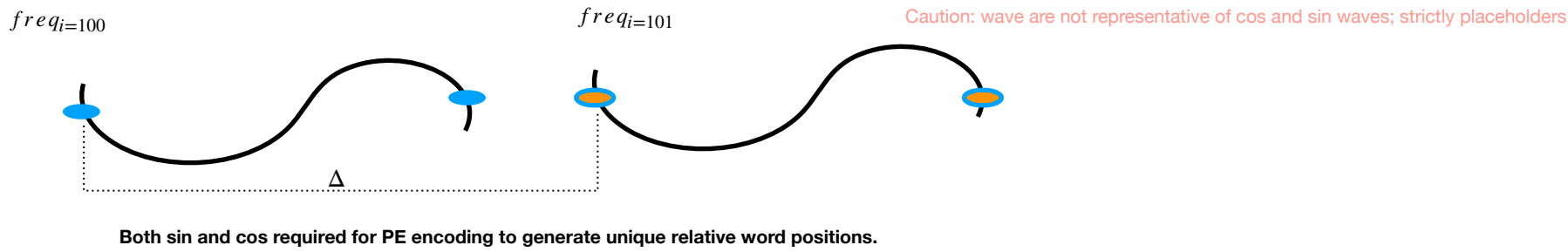
$$PE_{p,i} = \begin{cases} \sin(p/10000^{i/d}) & \text{if } i \text{ is even} \\ \cos(p/10000^{(i-1)/d}) & \text{if } i \text{ is odd} \end{cases}$$



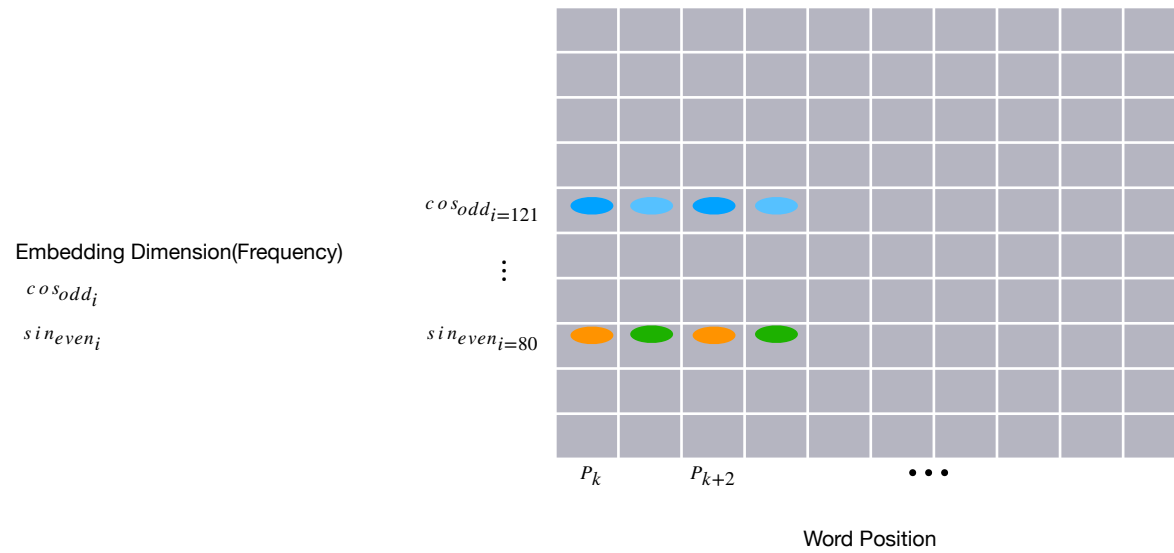
Note: Unique positional encoding at different frequencies (i) and positions(p)



Words located 38 word distances apart have a positional encoding in dimensions 100 and 101



Both sin and cos required for PE encoding to generate unique relative word positions.



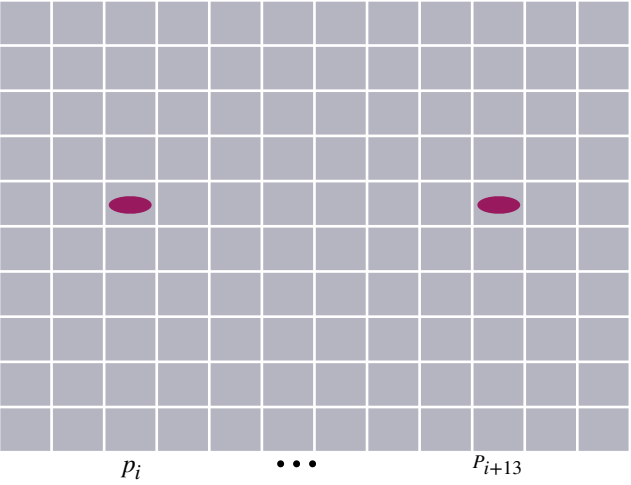
Words located 2 word distances apart have a positional encoding in dimensions 80 and 121

Positional Encodings

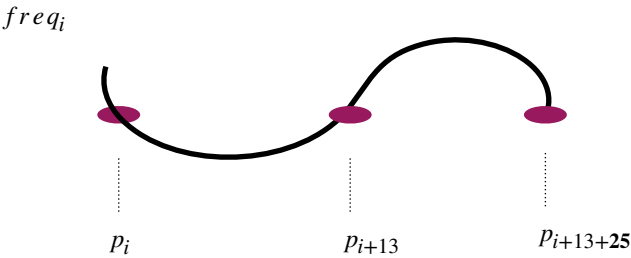
$$PE_{p,i} = \begin{cases} \sin(p/10000^{i/d}) & \text{for all } i \text{ even} \\ \cos(p/10000^{(i-1)/d}) & \text{for all } i \text{ odd} \end{cases}$$

Embedding Dimension(Frequency)

$\sin_{i=100}$

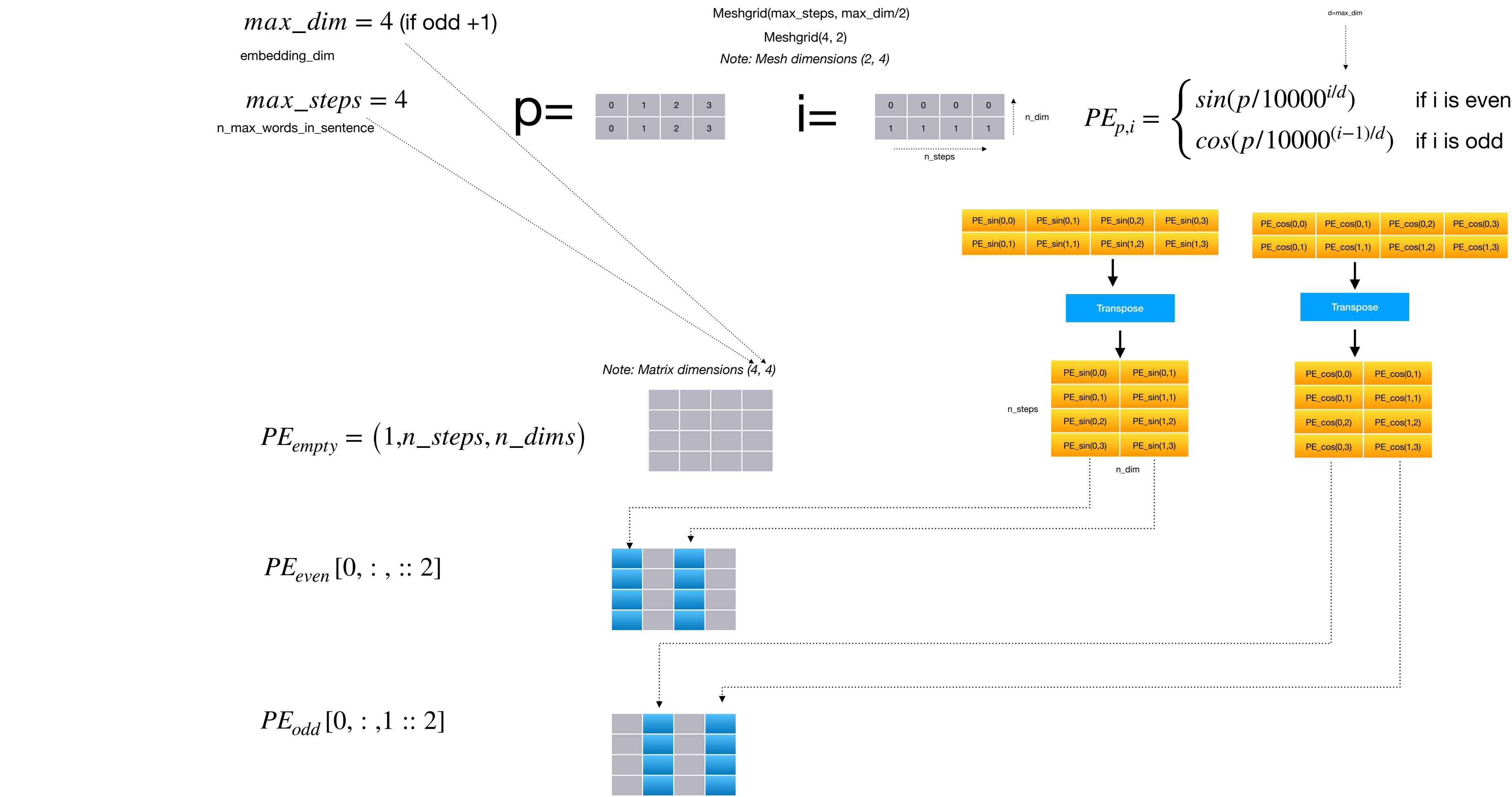


Word Position



Deltas are not unique. This is an example of aliasing. Model cannot learn from ambiguous encoding.

Positional Encodings : Simple Example



Note: Matrix dimensions (4, 4)

$PE_{empty} = (1, n_steps, n_dims)$

$PE_{even}[0, :, :, 2]$

$PE_{odd}[0, :, 1 :: 2]$

PE_sin(0,0)	PE_sin(0,1)	PE_sin(0,2)	PE_sin(0,3)
PE_sin(0,1)	PE_sin(0,1)	PE_sin(1,2)	PE_sin(1,3)

↓

Transpose

↓

PE_sin(0,0)	PE_sin(0,1)
PE_sin(0,1)	PE_sin(1,1)
PE_sin(0,2)	PE_sin(1,2)
PE_sin(0,3)	PE_sin(1,3)

n_steps

n_dim

PE_cos(0,0)	PE_cos(0,1)	PE_cos(0,2)	PE_cos(0,3)
PE_cos(0,1)	PE_cos(1,1)	PE_cos(1,2)	PE_cos(1,3)

↓

Transpose

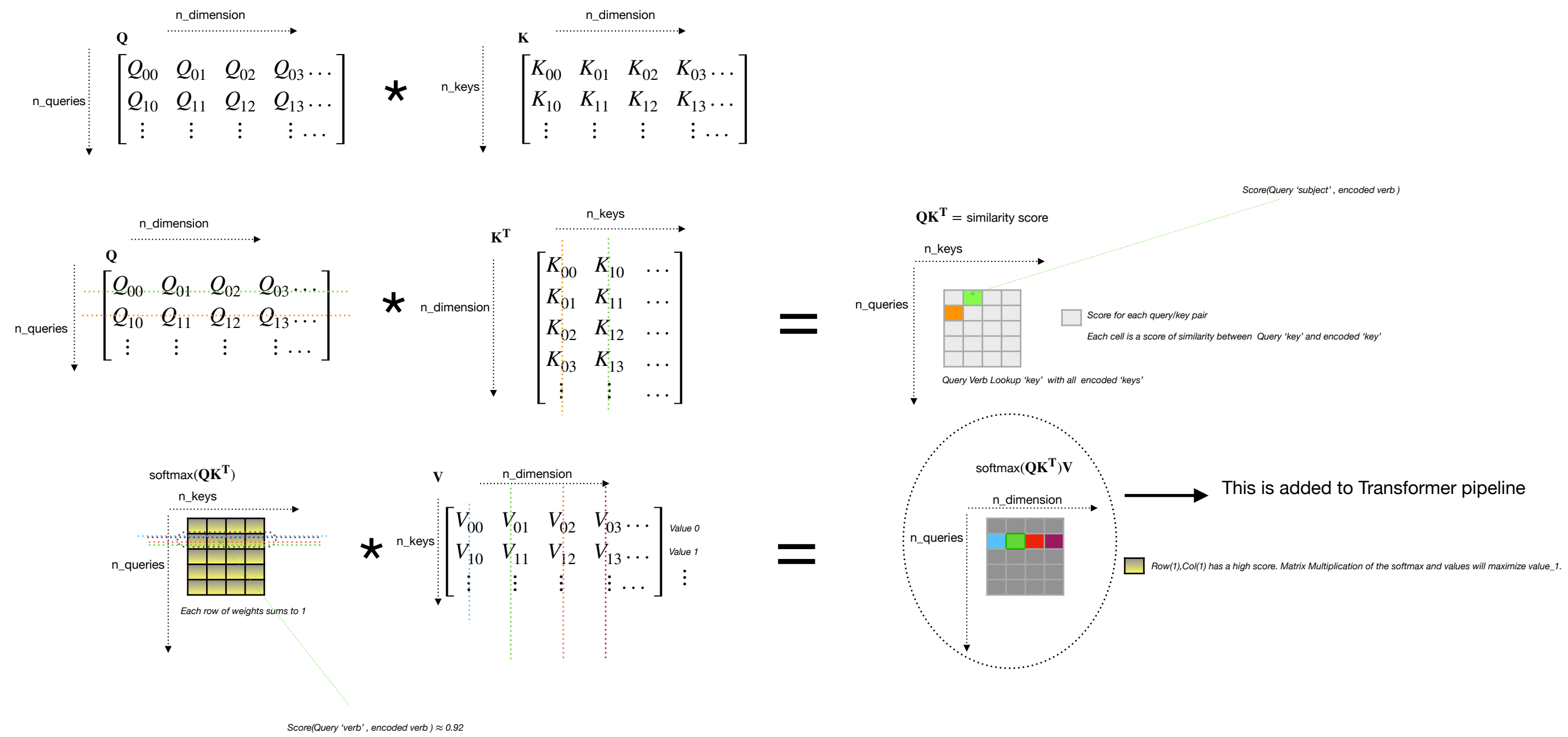
↓

PE_cos(0,0)	PE_cos(0,1)
PE_cos(0,1)	PE_cos(1,1)
PE_cos(0,2)	PE_cos(1,2)
PE_cos(0,3)	PE_cos(1,3)

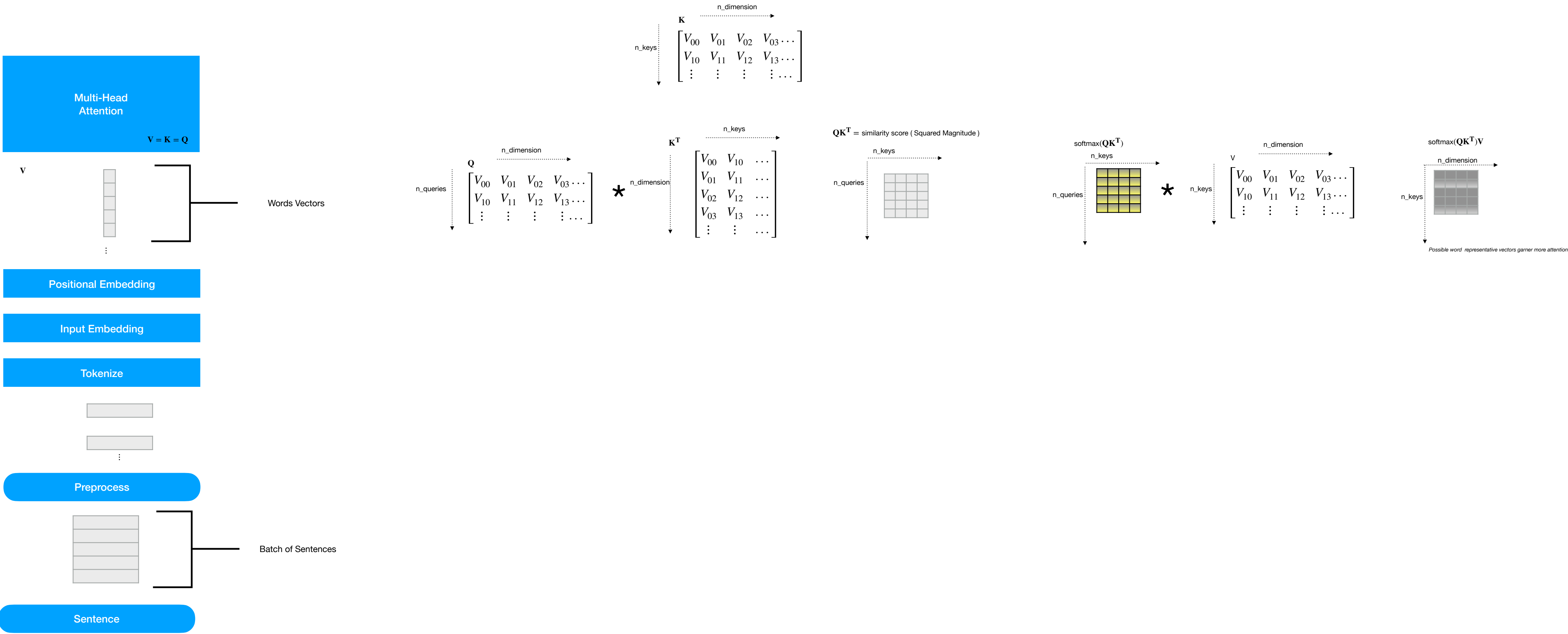
n_steps

n_dim

MultiHead Attention Layer



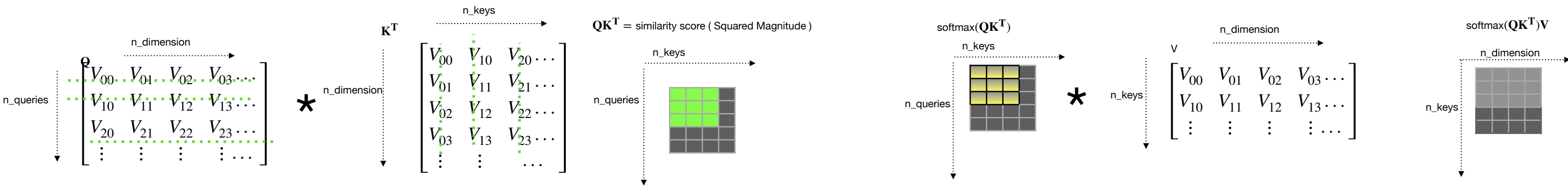
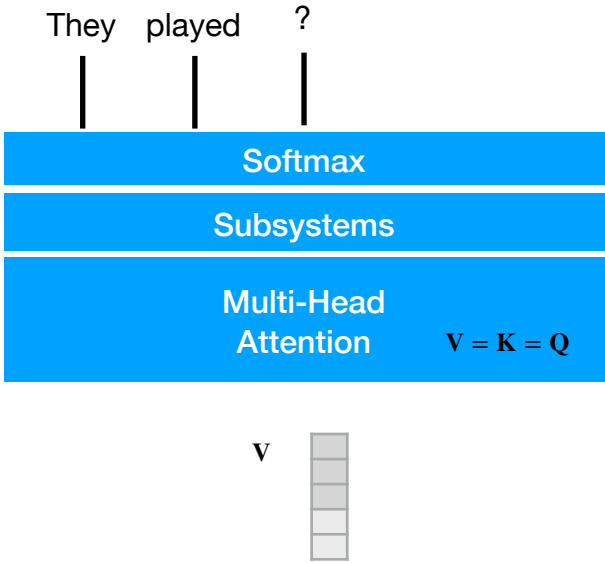
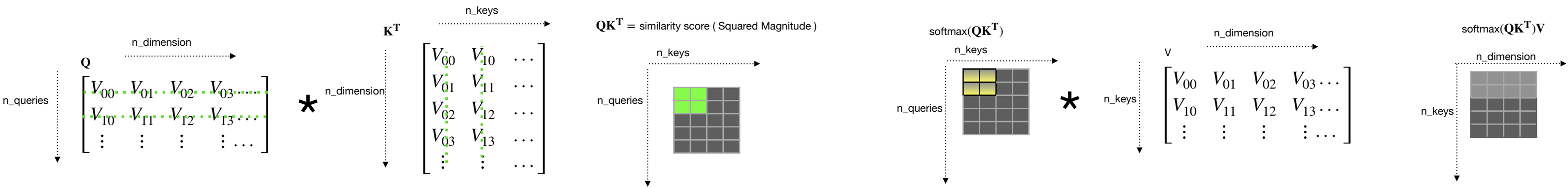
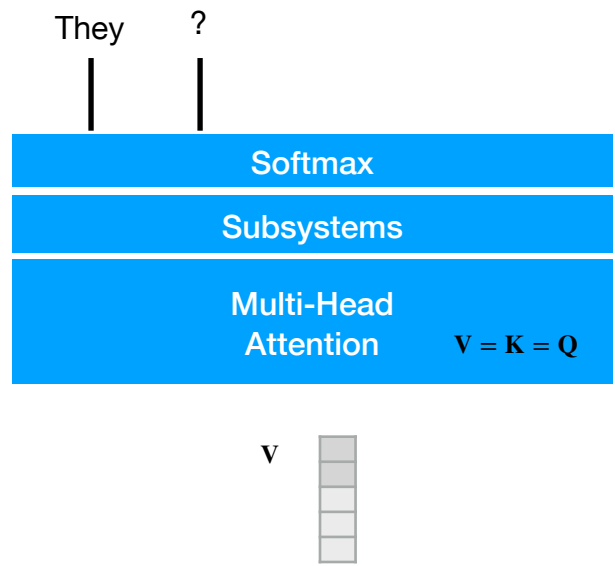
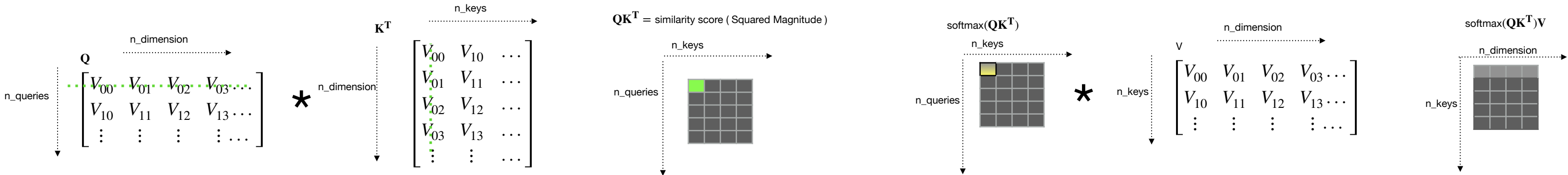
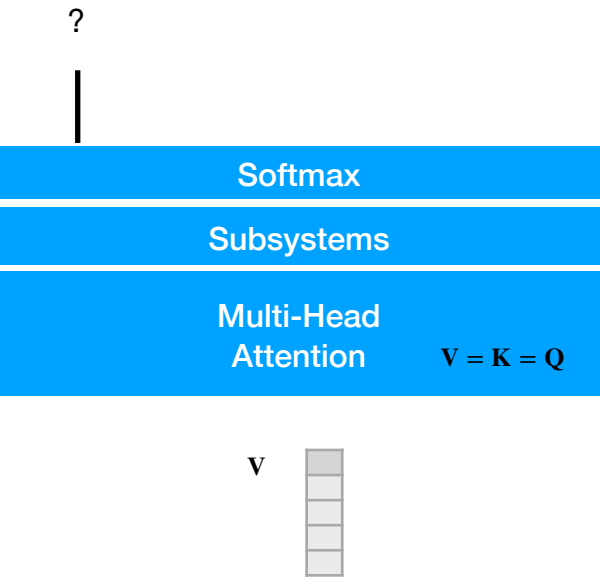
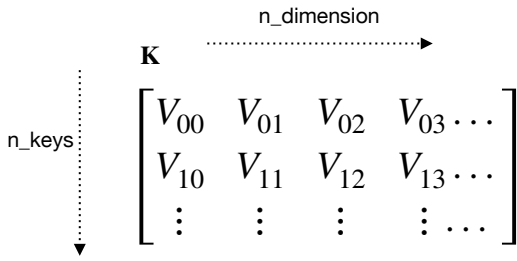
MultiHead Attention Layer : Encoder Subsystem

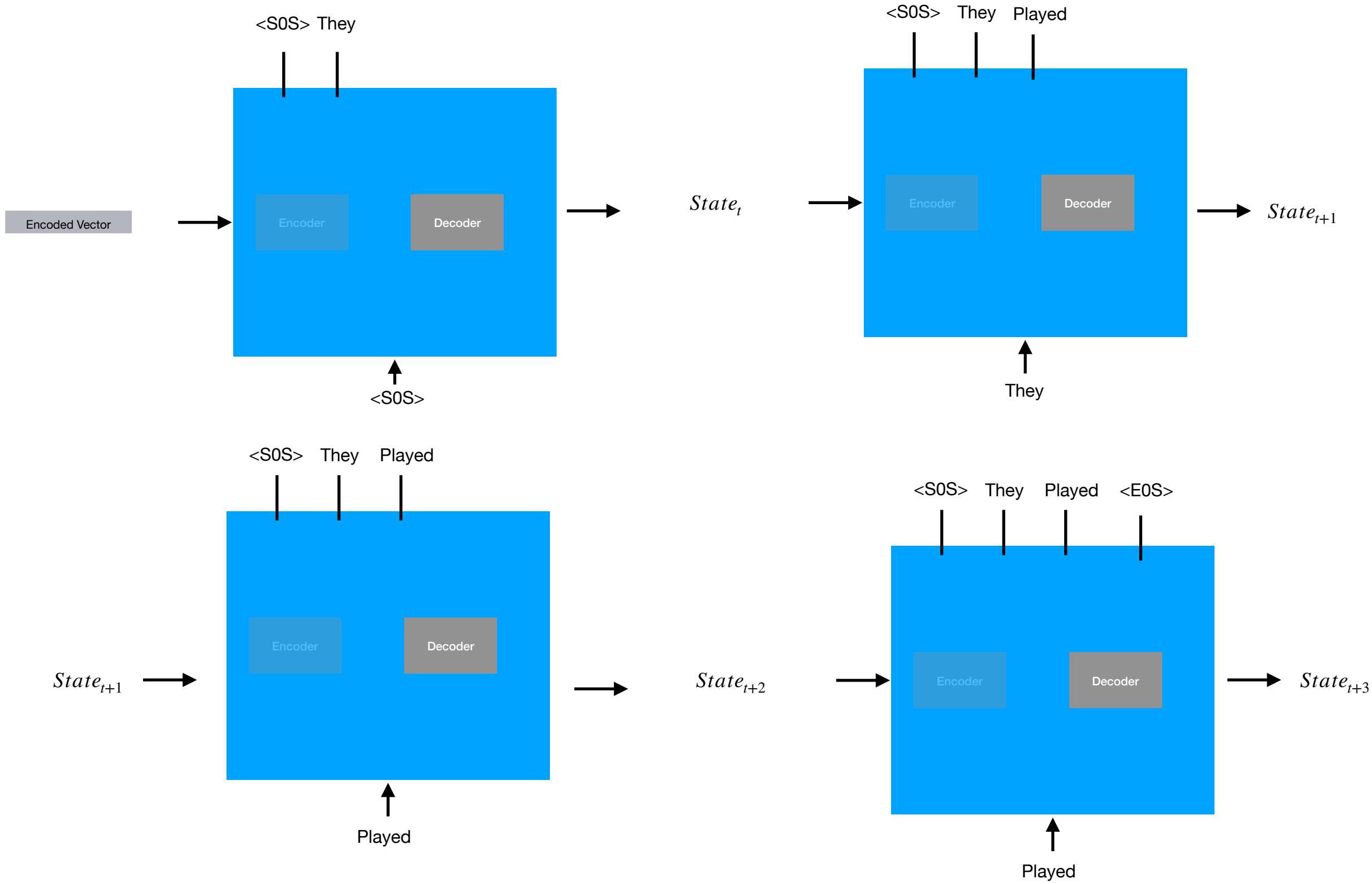
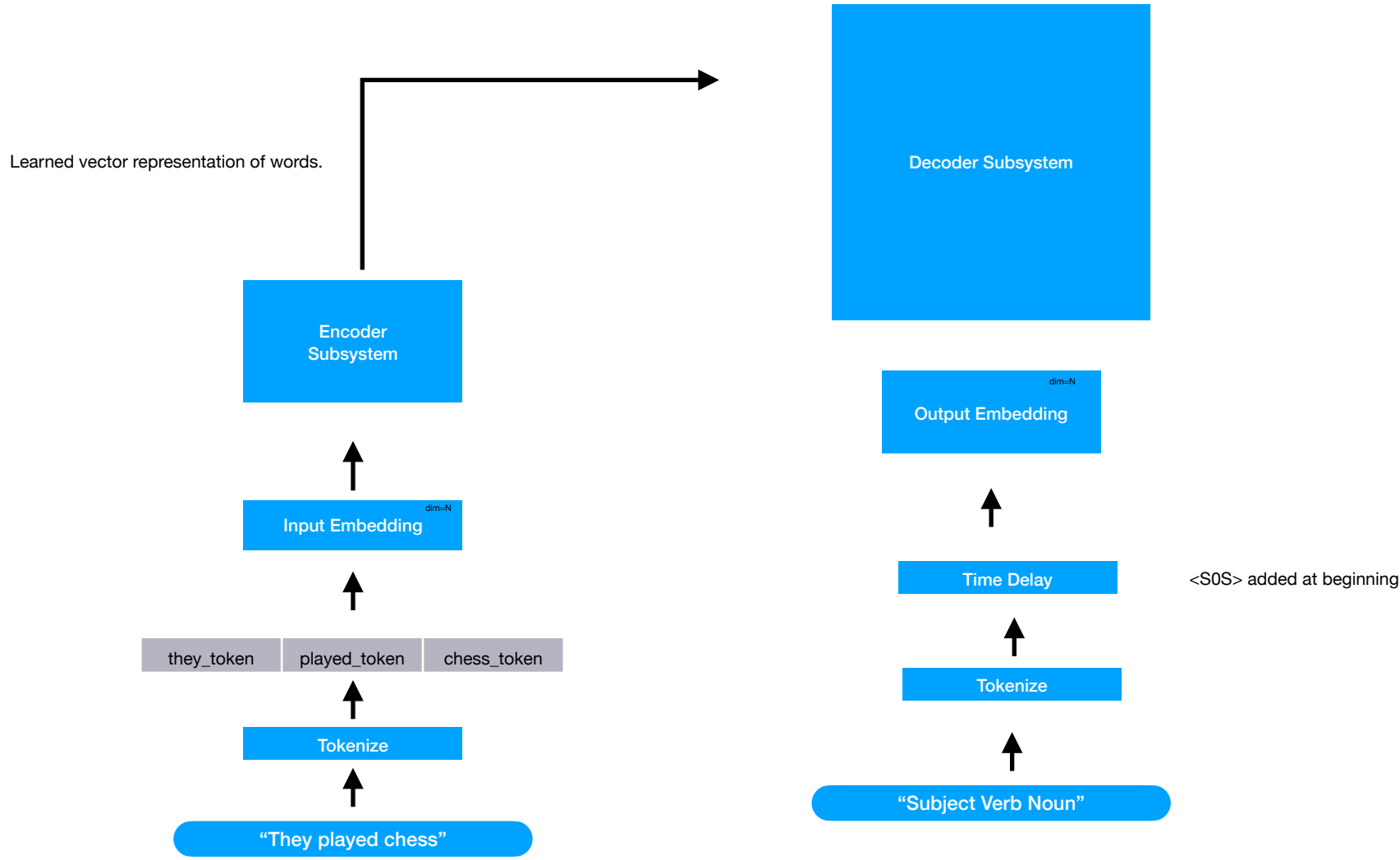


MaskedMultiHead Attention Layer : Decoder Subsystem

Mask MH Attention. Words cannot compare to words in the future.
Prevent a word from comparing itself to words located after it.
Masking can be done by adding large negative number to

QK^T = similarity score (Squared Magnitude)





Inference

