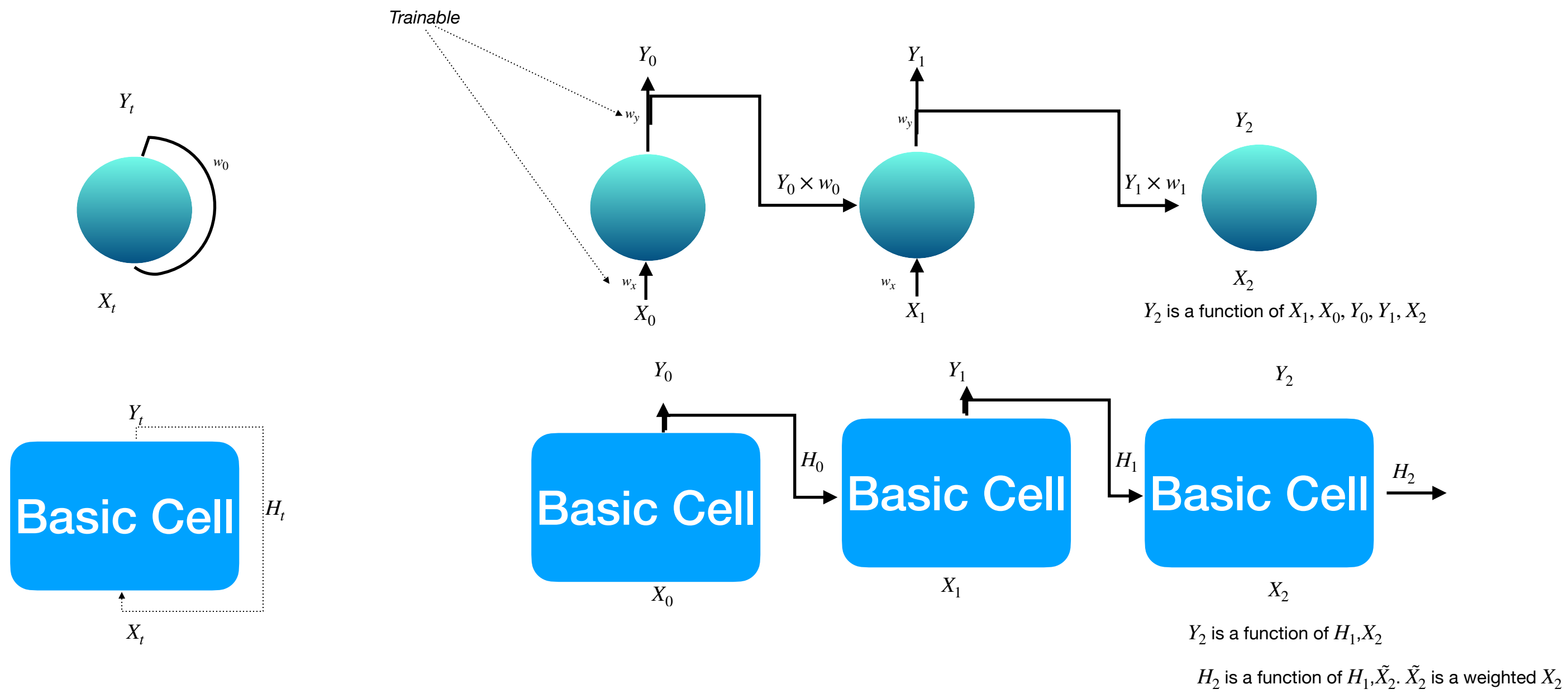# Stateless

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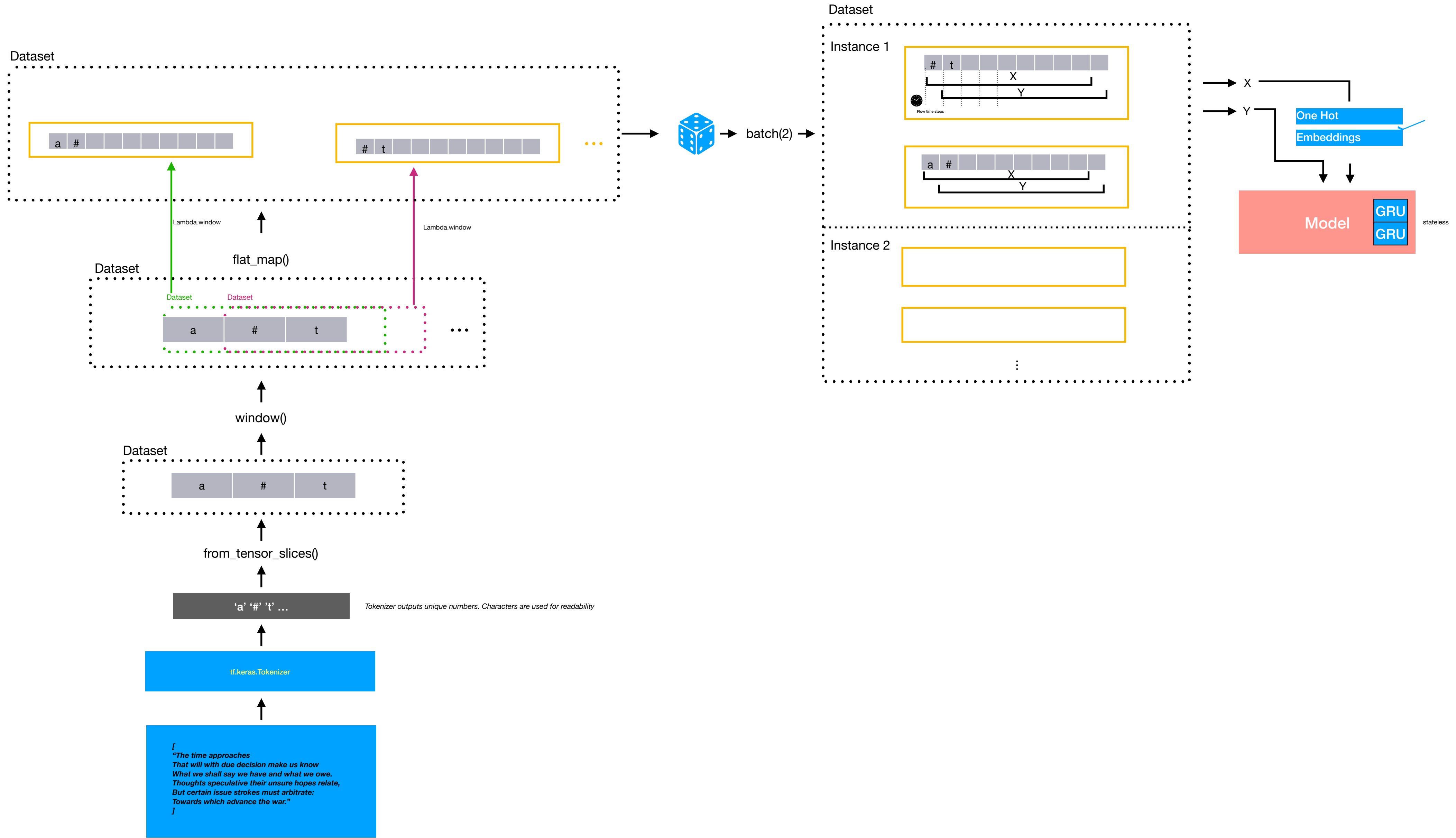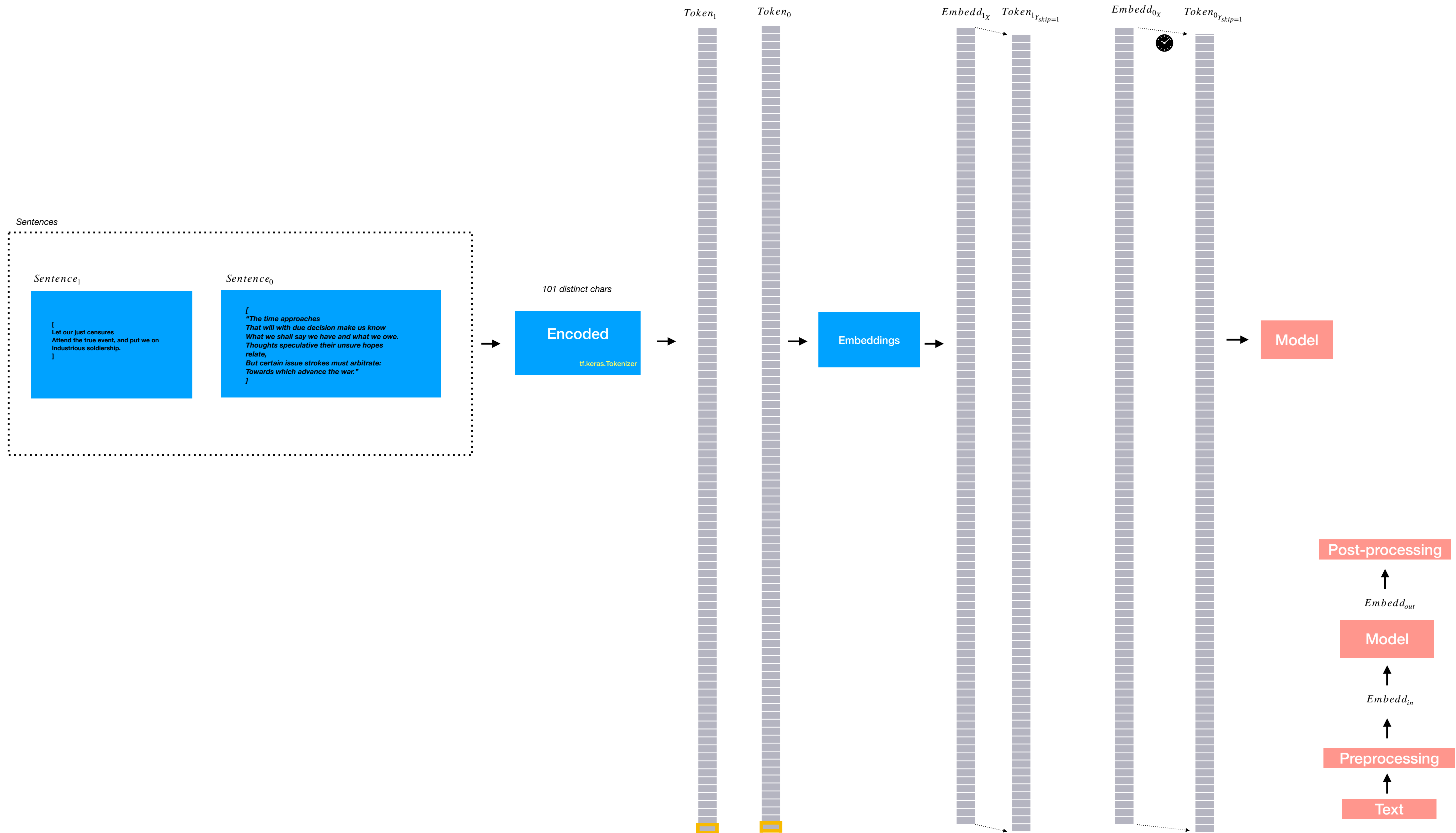


# Stateful

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

*Trainable*

$Y_t$

$w_0$

$X_t$

$Y_0$

$w_y$

$Y_0 \times w_0$

$w_x$

$X_0$

$Y_1$

$w_y$

$Y_1 \times w_1$

$w_x$

$X_1$

$Y_2$

$X_2$

$Y_2$ is a function of $X_1, X_0, Y_0, Y_1, X_2$

$Y_t$

Basic Cell

$H_t$

$X_t$

$Y_0$

Basic Cell

$H_0$

$X_0$

$Y_1$

Basic Cell

$H_1$

$X_1$

$Y_2$

Basic Cell

$H_2$

$X_2$

$Y_2$ is a function of $H_1, X_2$

$H_2$ is a function of $H_1, \tilde{X}_2$. $\tilde{X}_2$ is a weighted $X_2$

Dataset

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

| a | # | | | | | | | | |

| # | t | | | | | | | | |

· · ·

Lambda.window

Lambda.window

flat_map()

Dataset

Dataset          Dataset

| a | # | t |

· · ·

window()

Dataset

| a | # | t |

from_tensor_slices()

'a' '#' 't' ...

*Tokenizer outputs unique numbers. Characters are used for readability*

tf.keras.Tokenizer

[
"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."
]

batch(2)

Dataset

Instance 1

| # | t | | | | | | | |

X

Y

Flow time steps

| a | # | | | | | | | |

X

Y

Instance 2

X

One Hot

Embeddings

Y

Model    GRU
         GRU    stateless

*Sentences*

*Sentence$_1$*

[
Let our just censures
Attend the true event, and put we on
Industrious soldiership.
]

*Sentence$_0$*

[
"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."
]

*101 distinct chars*

**Encoded**

tf.keras.Tokenizer

$Token_1$  $Token_0$

$Embedd_{1_X}$  $Token_{1_{Y_{skip=1}}}$  $Embedd_{0_X}$  $Token_{0_{Y_{skip=1}}}$

**Embeddings**

**Model**

**Post-processing**

↑

$Embedd_{out}$

**Model**

↑

$Embedd_{in}$

↑

**Preprocessing**

↑

**Text**

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

# Stateless

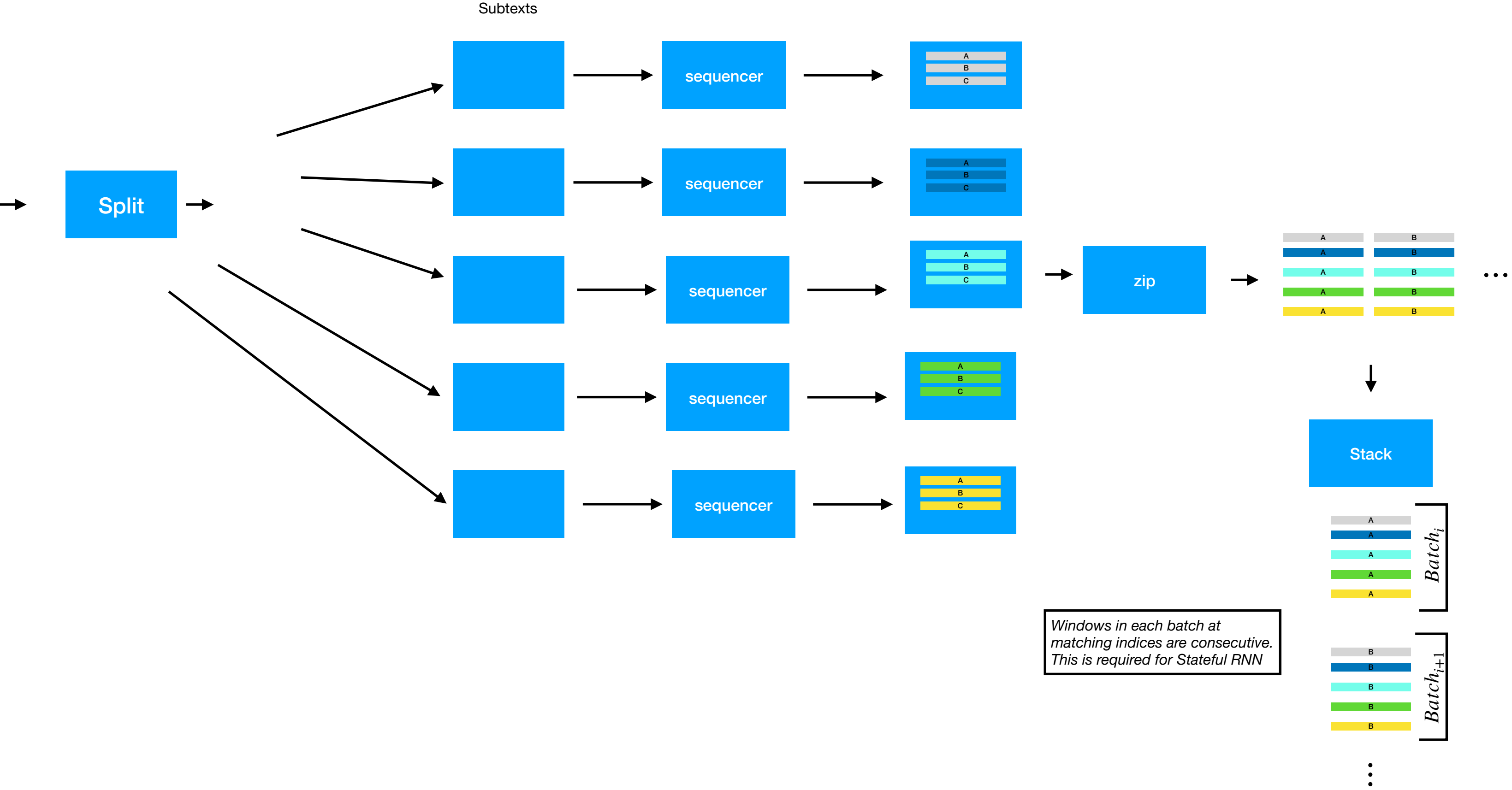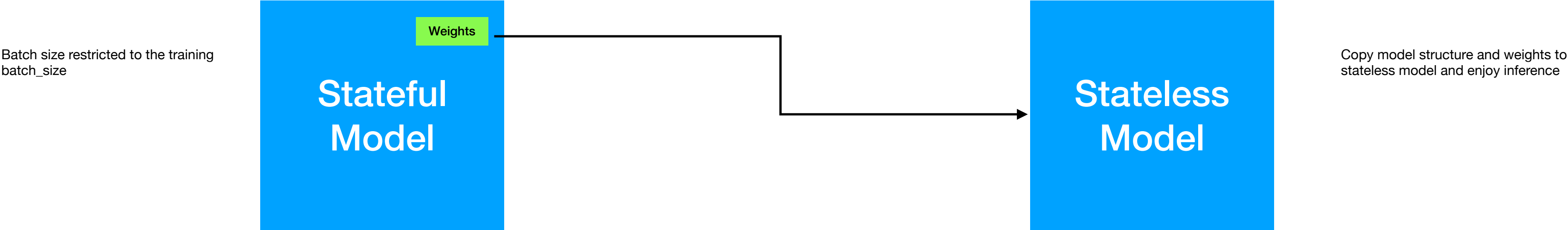**Learns on random portions of text, without any information about the text**

$Batch_i$

RNN

$Batch_{i+1}$

RNN

Hidden state lost after each training iteration.

# Stateful

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

$Batch_i$

RNN

$Batch_{i+1}$

RNN

Preserve hidden state and using to train next batch
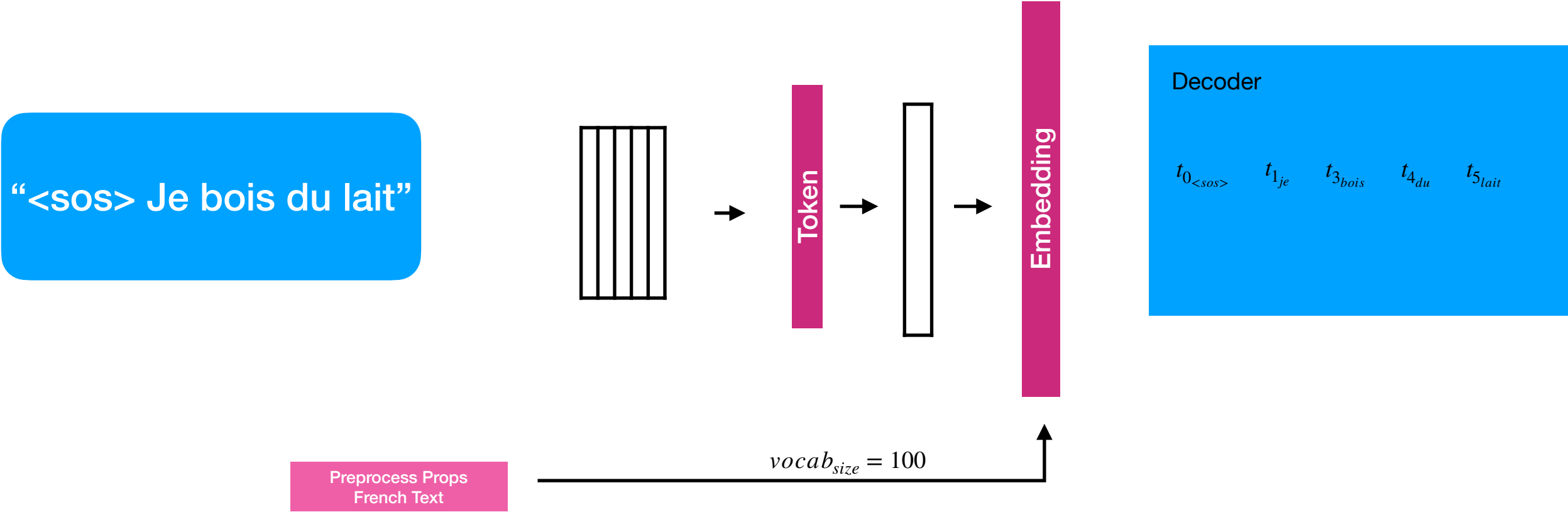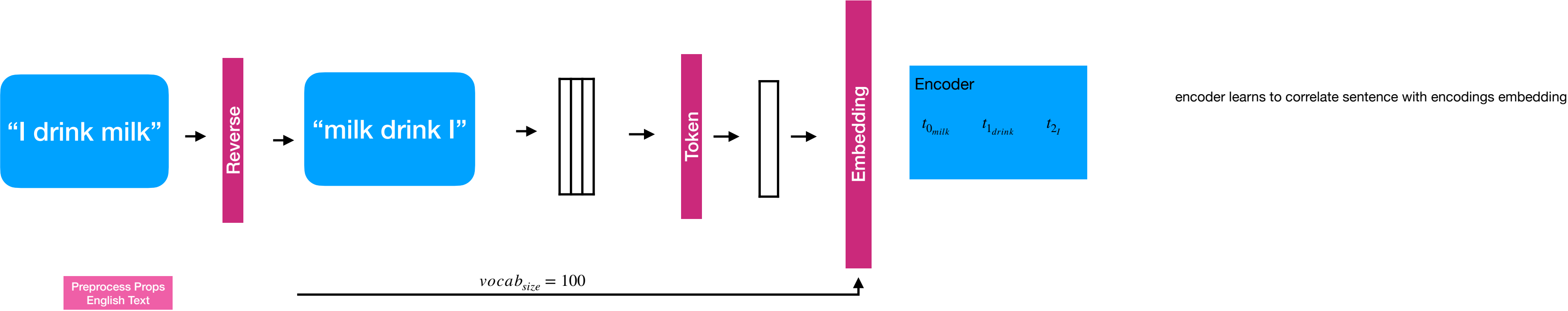
# Batching: Stateful RNN



Subtexts

[
"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:
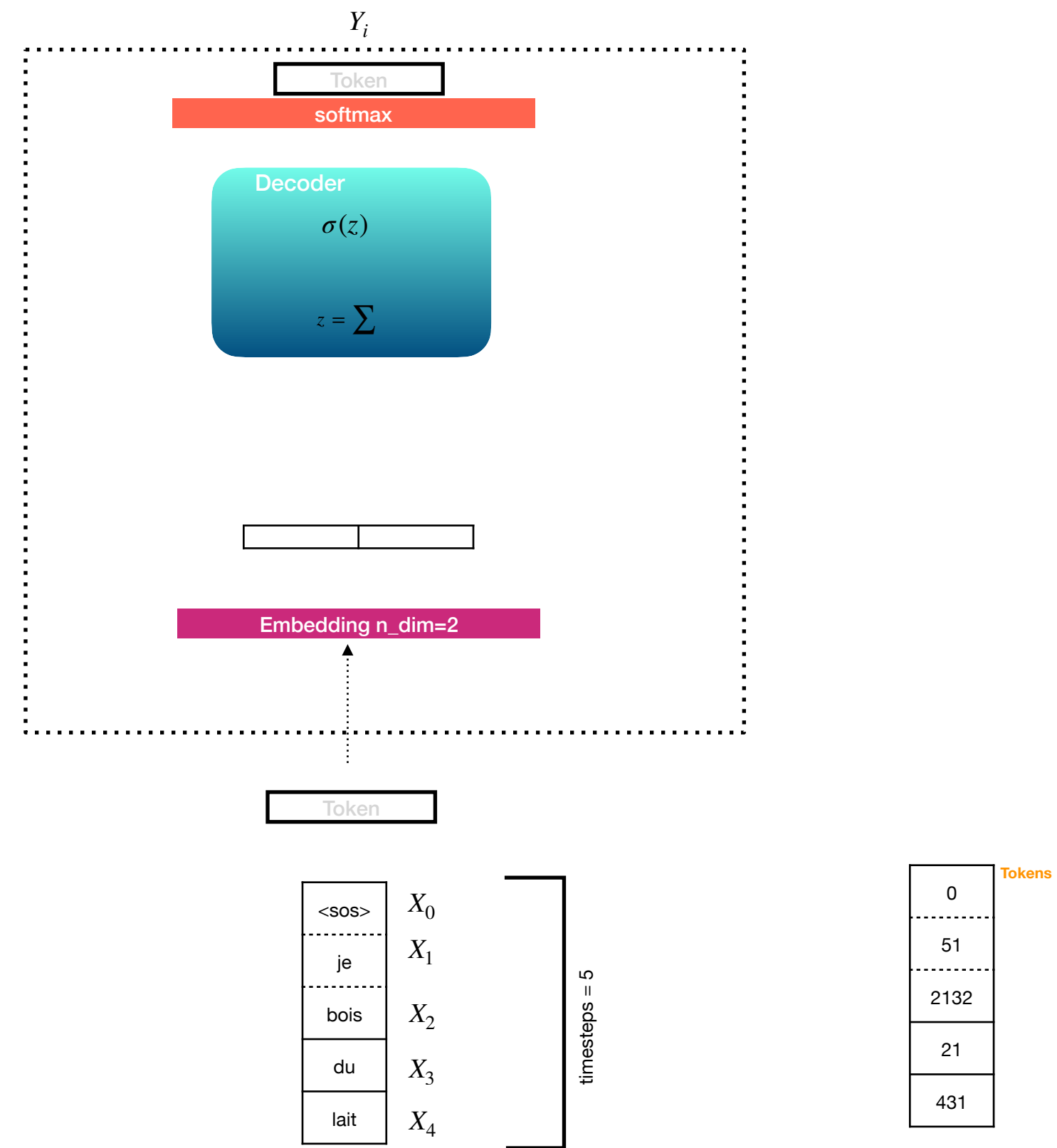Towards which advance the war."
]

Split

sequencer

zip

Stack

Windows in each batch at matching indices are consecutive. This is required for Stateful RNN

$Batch_i$

$Batch_{i+1}$

# Batching: Stateful RNN

Batch size restricted to the training
batch_size

**Stateful
Model**

Weights

**Stateless
Model**

Copy model structure and weights to
stateless model and enjoy inference

# Neural Machine Translation (NMT)

"I drink milk" → **Reverse** → "milk drink I" → → → **Token** → → → **Embedding**

**Encoder**

$t_{0_{milk}}$  $t_{1_{drink}}$  $t_{2_I}$

encoder learns to correlate sentence with encodings embedding

**Preprocess Props English Text**

$vocab_{size} = 100$

"<sos> Je bois du lait" → → **Token** → → **Embedding**

**Decoder**

$t_{0_{<sos>}}$  $t_{1_{je}}$  $t_{3_{bois}}$  $t_{4_{du}}$  $t_{5_{lait}}$

**Preprocess Props French Text**

$vocab_{size} = 100$

$Y_i$

Encoder (RNN)

$\sigma(z)$

$z = \sum$

n_neurons = vocab size

Mini-batch instances = 1

Embedding n_dim=2

vocab_size is preconfigured during preprocessing

*n_dim would actually be around 100-200 to help model generalize better*

Token

| | |
|---|---|
| milk | $X_0$ |
| drink | $X_1$ |
| I | $X_2$ |

timesteps = 3

Tokens

| |
|---|
| 288 |
| 3335 |
| 72 |

$Y_i$

Token

softmax

Decoder

$\sigma(z)$

$z = \sum$

Embedding n_dim=2

Token

| | |
|---|---|
| <sos> | $X_0$ |
| je | $X_1$ |
| bois | $X_2$ |
| du | $X_3$ |
| lait | $X_4$ |

timesteps = 5

Tokens

| |
|---|
| 0 |
| 51 |
| 2132 |
| 21 |
| 431 |

$Y_0$

Encoder (RNN)

$\sigma(z)$

$z = \sum$

$H_{-1}$ $\longrightarrow$ $H_0$

Embedding n_dim=2 $\longleftarrow$ vocab_size is preconfigured during preprocessing

Token

milk $X_0$

$$Y_{-1_{token\_target}} \_or\_ Y_{-1_{token\_out}}$$

Embedding n_dim=2

$Y_0$

Token

softmax

Decoder (RNN)

$\sigma(z)$

$z = \sum$

$H_{2_{encoder}}$ ........ $H_0$

Embedding n_dim=2

Token

<SOS>  $X_0$

$Y_{0_{token\_target}} - or\_Y_{0_{token\_out}}$

Embedding n_dim=2

$Y_1$

Token

softmax

Decoder (RNN)

$\sigma(z)$

$z = \sum$

$H_0$

$H_1$

Embedding n_dim=2

Token

je  $X_1$

$Y_{1_{token\_target}} \_or\_ Y_{1_{token\_out}}$

Embedding n_dim=2

$Y_2$

Token

softmax

Decoder (RNN)

$\sigma(z)$

$z = \sum$

$H_1$

$H_2$

Embedding n_dim=2

Token

bois $X_2$

$Y_{2_{token\_target}} - or\_Y_{2_{token\_out}}$

Embedding n_dim=2

$Y_3$

Token

softmax

Decoder (RNN)

$\sigma(z)$

$z = \sum$

$H_2$

$H_3$

Embedding n_dim=2

Token

du $X_3$

Encoder

$t_{0_{milk}}$   $t_{1_{drink}}$   $t_{2_{I}}$

$Y_0$

RNN

$t_{0_{milk}}$

$Y_1$

RNN

$t_{1_{drink}}$

$Y_2$

RNN

$t_{2_{I}}$

Decoder

$t_{0_{<sos>}}$   $t_{1_{je}}$   $t_{2_{bois}}$   $t_{3_{du}}$   $t_{4_{lait}}$

$Y'_{0_{t=3}}$

RNN

$t_{0_{<sos>}}$

$Y'_{1_{t=4}}$

RNN

$t_{1_{je}}$

$Y'_{2_{t=5}}$

RNN

$t_{2_{bois}}$

$Y'_{3_{t=6}}$

RNN

$t_{3_{du}}$

$Y'_{4_{t=7}}$

RNN

$t_{4_{lait}}$

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

Model learns to translate English sentence into French

Encoder

$t_{0_{milk}}$    $t_{1_{drink}}$    $t_{2_I}$

$Y_0$    $Y_1$    $Y_2$

LTSM    LTSM    LTSM

$t_{0_{milk}}$    $t_{1_{drink}}$    $t_{2_I}$

Decoder

$t_{0_{<sos>}}$    $t_{1_{je}}$    $t_{2_{bois}}$    $t_{3_{du}}$    $t_{4_{lait}}$

$Y'_{0_{t=3}}$    $Y'_{1_{t=4}}$    $Y'_{2_{t=5}}$    $Y'_{3_{t=6}}$    $Y'_{4_{t=7}}$

LTSM    LTSM    LTSM    LTSM    LTSM

$t_{0_{<sos>}}$    $t_{1_{je}}$    $t_{2_{bois}}$    $t_{3_{du}}$    $t_{4_{lait}}$

Model learns to translate English sentence into French

# Batch sentences of varying length

| | | | batch |
|---|---|---|---|
| **Low** | Sentences 1-6 words | Pad = 6 | 6 token  6 token  6 token  6 token |
| **Mid** | Sentences 7-12 words | Pad = 12 | 12 token  12 token  12 token  12 token |
| **Long** | Sentences 13-50 more words | Pad = 50 | 50 token  50 token  50 token  50 token |
| **Superlong** | Sentences 51-1000 more words | Pad = 1000 | 1000 token  1000 token  1000 token  1000 token |

## Batch of equal length Tensors are accepted as inputs to model

$batch_3$  |  $batch_2$  |  $batch_1$  |  $batch_0$

1000 token  1000 token  1000 token  1000 token  |  50 token  50 token  50 token  50 token  |  12 token  12 token  12 token  12 token  |  6 token  6 token  6 token  6 token  → **Model**

"I drink milk" → Model → [vocab vector] → Postprocess → "Je bois du lait"

Output vocab length could be very long

# Handling large output vocab

Text → Model →

Sample Random selection of incorrect words →

Sample Logits output by model for correct word →

Loss based on logits

**Regular RNN Layer
(Causal)**

Predict next steps

Present Input

*Present and past inputs determine next steps*

Past inputs

**Bidirectional
(Non Causal)**

Predict next steps

*Future inputs (lookahead)*

Present Input

*Present and past inputs*

Past inputs

# Bidirectional RNN

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

$Y_{(0)}$

$Y_{(1)}$

$Y_{(2)}$

RNN   RNN   RNN   Cloned RNN in reverse direction

RNN   RNN   RNN

Embedding   Embedding   Embedding

$X_{(0)}$ Queen of the United Kingdon   $X_{(1)}$ the queen of hearts   $X_{(2)}$ the queen bee

Predict next steps

Future inputs (lookahead)

Present Input

Present and past inputs

Past inputs

Predict Output $Y_1$

Past inputs   Present Input $X_1$   Future Output

# Bidirectional RNN

$$Y_{(0)}$$

$$Y_{(1)}$$

$$Y_{(2)}$$

Present Input $X_1$

Future Output

Past inputs

| RNN | RNN | RNN |
| RNN | RNN | RNN |

Embedding

Embedding

Embedding

Phrases (Samples)

| Queen of the United Kingdom | the queen of hearts | the queen bee |

Predict next steps

Future inputs (lookahead)

Present Input

Present and past inputs

Past inputs

Predict Output $Y_1$

Learns next probable sentence by learning from current, present, and future phrases in different forms of English text

# 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 → Estimate next three →

Decoder Clone →

Decoder Clone →

**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**

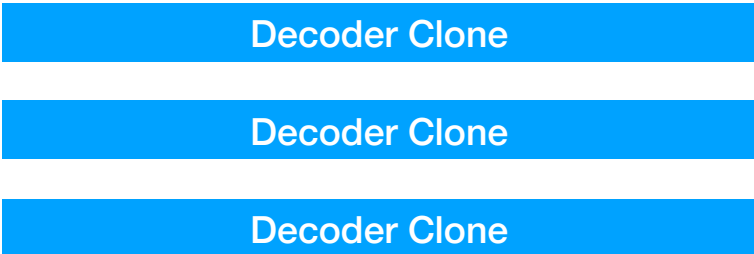Decoder Estimates probability for each word in vocab, only top three make the short list, **given the sentence starts with how**

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 → | **How, Will** | Estimate next three → |
| Decoder Clone → | **How, Are** | Estimate next three → |
| Decoder Clone → | **How, Do** | Estimate next three → |

# Attention Mechanisms

$Y_t$

Encoder (Memory Unit)

$\sigma(z)$

$z = \sum$

$H_{t+1}$

$H_t$

Encoder (Memory Unit)

$\sigma(z)$

$z = \sum$

$H_{t-1}$

$H_t$

| milk | $X_0$ |
| drink | $X_1$ |
| I | $X_2$ |

$Y_t$

Decoder (Memory Unit)

$\sigma(z)$

$z = \sum \alpha_{(t,i)} \times Y_{i_{endcoder}}$

$t - decoder_{time}$

$i - encoder_{time}$

$\alpha_{(t,i)}$

$Y_{i_{encoder}}$

# Attention Mechanisms

$Y_0$

## Encoder (RNN)

$\sigma(z)$

$z = \sum$

$H_2$

$H_3$

## Encoder (RNN)

$\sigma(z)$

$z = \sum$

$H_{-1}$

$H_0$

I  $X_0$

milk  $X_0$

# Attention Mechanisms

# Attention Mechanisms

# Attention Mechanisms

$Y_{(t=-1)_{output\_word}}$

$Y_{(t=-1)_{target\_word}}$

$Y_{(t=0)}$

$H_{2_{enc}}$ ............................➤

**Decoder (RNN)**

$\sigma(z)$

$H_0$ ............................➤

$z = \sum \alpha_{(t,i)} \times Y_i$

$\sum$

$\alpha_{(0,0)}$

$\alpha_{(0,1)}$

$\alpha_{(0,2)}$ ◀.......... *Weight factor on encoder output*

*Note: $i$ and $enc_t$ are equivalent*

$Y_{enc_{t=0}}$

$Y_{enc_{t=1}}$

$Y_{enc_{t=2}}$

# Attention Mechanisms:
## Calculating Weights

$Y_{(t=-1)_{output\_word}}$

$Y_{(t=-1)_{target\_word}}$

$Y_{(t=0)}$

**Decoder (RNN)**

$\sigma(z)$

$H_{2_{enc}}$

$z = \sum \alpha_{(t,i)} \times Y_i$

$H_0$

$\sum$

$\alpha_{(0,0)}$

$\alpha_{(0,1)}$

$\alpha_{(0,2)}$

Weight factor on encoder output

Note: $i$ and $enc_t$ are equivalent

$Y_{enc_{t=0}}$

$Y_{enc_{t=1}}$

$Y_{enc_{t=2}}$

Concat

Concat

Concat

Dense()

Dense()

Dense()

$e_{(0,0)}$

$e_{(0,1)}$

$e_{(0,2)}$

Soft

$\alpha_{(0,0)}$

$\alpha_{(0,1)}$

$\alpha_{(0,2)}$

Time Distributed Dense Layer

Alignment Neural Network

Softmax is not Time Distributed

# Attention Mechanisms

$Y_{(t=0)output\_word}$
$Y_{(t=0)target\_word}$

$Y_{(t=1)}$

Decoder (RNN)

$\sigma(z)$

$H_0$ . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . ►

$z = \sum \alpha_{(1,i)} \times Y_i$

. . . . . . . . . . . . . . . . . . . . ► $H_1$

$$\sum$$

$\alpha_{(1,0)}$     $\alpha_{(1,1)}$     $\alpha_{(1,2)}$ . . . . . . . . *Weight factor on encoder output*

$Y_{enc_{t=0}}$      $Y_{enc_{t=1}}$      $Y_{enc_{t=2}}$      *Note: $i$ and $enc_t$ are equivalent*

$Y_{(t=0)_{output\_word}}$
$Y_{(t=0)_{target\_word}}$

$Y_{(t=1)}$

**Decoder (RNN)**

$\sigma(z)$

$H_0$ ........................➤

$z = \sum \alpha_{(t,i)} \times Y_i$

$H_1$ ........................➤

$\sum$

$\alpha_{(1,0)}$

$\alpha_{(1,1)}$

$\alpha_{(1,2)}$ ............... Weight factor on encoder output

$Y_{enc_{t=0}}$

$Y_{enc_{t=1}}$

$Y_{enc_{t=2}}$

Note: $i$ and $enc_t$ are equivalent

Concat

Dense()

$e_{(1,0)}$

Softmax

$\alpha_{(1,0)}$

Concat

Dense()

$e_{(1,1)}$

$\alpha_{(1,1)}$

Concat

Dense()

$e_{(1,2)}$

$\alpha_{(1,2)}$

Time Distributed Dense Layer

Softmax is not Time Distributed

Alignment Neural Network

Alignment Neural Network finds similarity between encoders and decoders previous hidden state

# Attention Mechanisms
## (Method 2)

$Y_{(t=-1)_{output\_word}}$

$Y_{(t=-1)_{target\_word}}$

Dot product measures similarity

$Y_{(t=0)}$

$v_1$

$v_2$

Decoder (RNN)
$\sigma(\tilde{h})$

$H_{-1}$

$H_0$

$\tilde{h}$

$\Sigma$

$\alpha_{(0,0)}$

$\alpha_{(0,1)}$

$\alpha_{(0,2)}$

$Y_{enc_{t=0}}$

$Y_{enc_{t=1}}$

$Y_{enc_{t=2}}$

Dot

Dot

Dot

Dense(1)

Dense(1)

Dense(1)

$e_{(0,0)}$

$e_{(0,1)}$

$e_{(0,2)}$

Softmax

$\alpha_{(0,0)}$

$\alpha_{(0,1)}$

$\alpha_{(0,2)}$

Time Distributed Dense Layer

Softmax is not Time Distributed

Alignment Neural Network (Minh - Thang method)

Alignment Neural Network finds similarity between encoder and decoder current state

# Positional Encodings



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

| $PE_{0,0}$ |
| $PE_{1,0}$ |
| $PE_{2,0}$ |
| $PE_{3,0}$ |
| $PE_{4,0}$ |
| $PE_{p,i}$ |

| $emb_{0,0}$ |
| $emb_{1,0}$ |
| $emb_{2,0}$ |
| $emb_{3,0}$ |
| $emb_{4,0}$ |
| Embedding (n_dim = 1) |

i - component of the word encoding (dim)
p - position index of word

"Hello my name is Dragon

| $PE_{0,0}$ | $PE_{0,1}$ |
| $PE_{1,0}$ | $PE_{1,1}$ |
| $PE_{2,0}$ | $PE_{2,1}$ |
| $PE_{3,0}$ | $PE_{3,1}$ |
| $PE_{4,0}$ | $PE_{4,1}$ |
| $PE_{p,i}$ | |

| $emb_{0,0}$ | $emb_{0,1}$ |
| $emb_{1,0}$ | $emb_{1,1}$ |
| $emb_{2,0}$ | $emb_{2,1}$ |
| $emb_{3,0}$ | $emb_{3,1}$ |
| $emb_{4,0}$ | $emb_{4,1}$ |
| Embedding (n_dim = 2) | |

"Hello my name is Dragon

# Positional Encodings

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

$freq_i$

Embedding Dimension(Frequency)
$cos_{odd_i}$
$sin_{even_i}$

$p$

Word Position

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

$freq_i$

$p$

---

Embedding Dimension(Frequency)
$cos_{odd_i}$
$sin_{even_i}$

$cos_{odd_{i=101}}$
$sin_{even_{i=100}}$

$P_k$ $P_{k+1}$ $\cdots$ $P_{k}+38$
$P_{k+1+38}$

Word Position

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

$freq_{i=100}$

$freq_{i=101}$

Caution: wave are not representative of cos and sin waves; strictly placeholders

$\Delta$

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

---

Embedding Dimension(Frequency)
$cos_{odd_i}$
$sin_{even_i}$

$cos_{odd_{i=121}}$

$\vdots$

$sin_{even_{i=80}}$

$P_k$ $P_{k+2}$ $\cdots$

Word Position

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

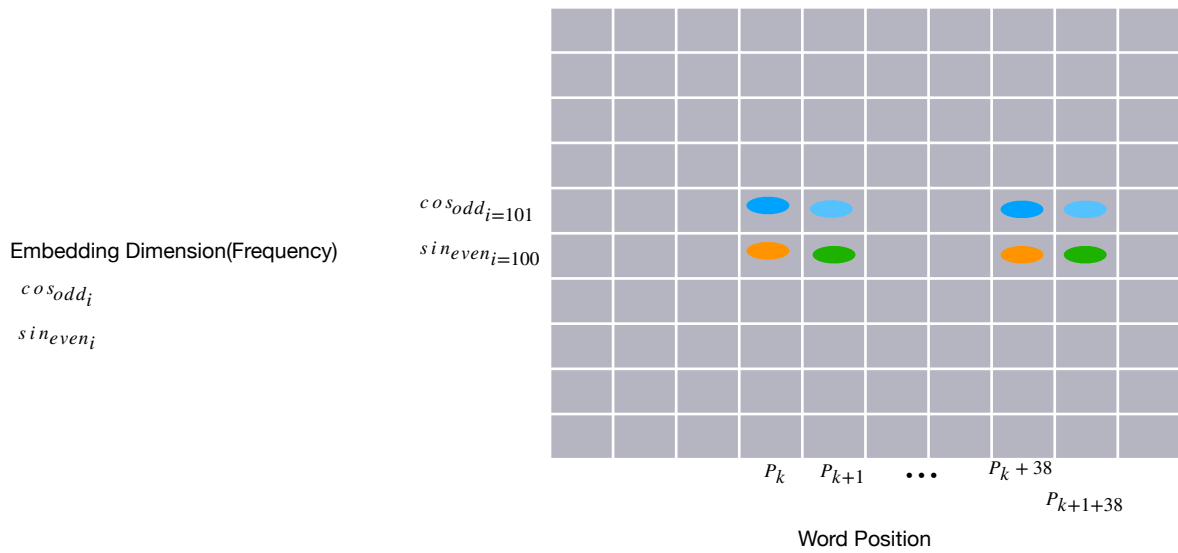# Positional Encodings

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

Embedding Dimension(Frequency)    $sin_{i=100}$

Word Position

$p_i$    $\cdots$    $p_{i+13}$

$freq_i$

$p_i$    $p_{i+13}$    $p_{i+13+\mathbf{25}}$

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

$max\_dim = 4$ (if odd +1)

embedding_dim

$max\_steps = 4$

n_max_words_in_sentence

Meshgrid(max_steps, max_dim/2)

Meshgrid(4, 2)

*Note: Mesh dimensions (2, 4)*

d=max_dim

$$p= \begin{array}{|c|c|c|c|} \hline 0 & 1 & 2 & 3 \\ \hline 0 & 1 & 2 & 3 \\ \hline \end{array}$$

$$i= \begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 1 & 1 & 1 & 1 \\ \hline \end{array}$$ n_dim

n_steps

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

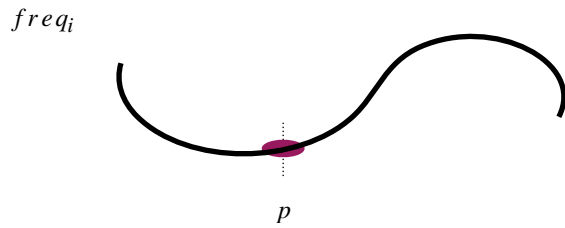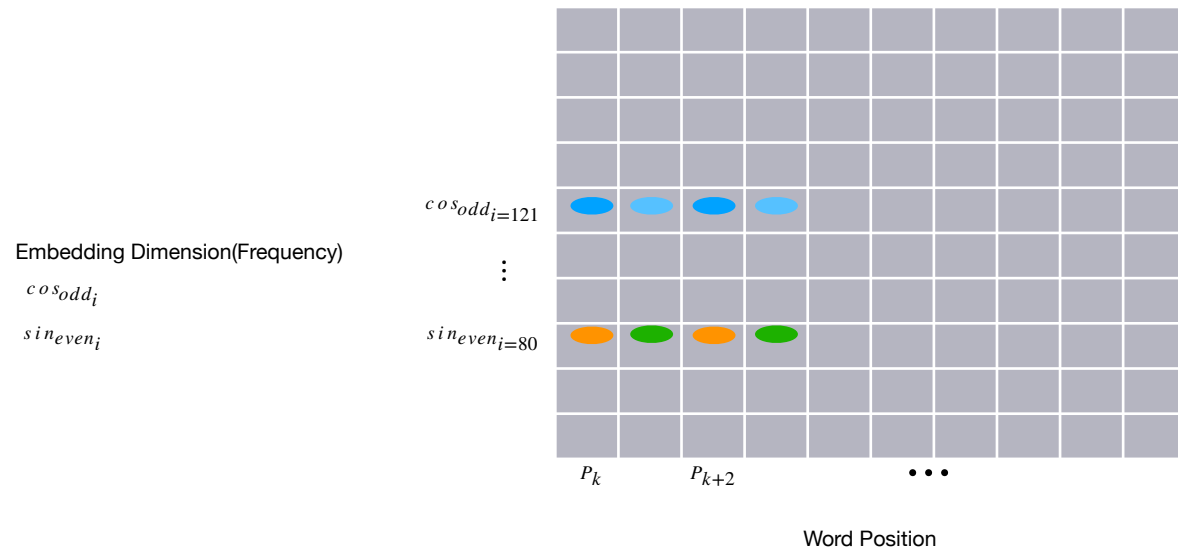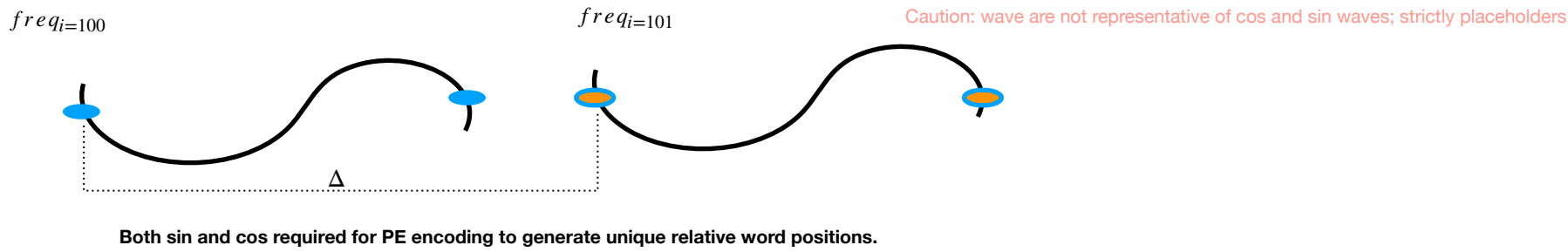| PE_sin(0,0) | PE_sin(0,1) | PE_sin(0,2) | PE_sin(0,3) |
|---|---|---|---|
| PE_sin(0,1) | PE_sin(1,1) | PE_sin(1,2) | PE_sin(1,3) |

| 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

Transpose

*Note: Matrix dimensions (4, 4)*

$PE_{empty} = \left(1, n\_steps, n\_dims\right)$

n_steps

| 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_dim

| 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) |

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

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

# MultiHead Attention Layer

$n\_dimension$

**Q**

$$\begin{bmatrix} Q_{00} & Q_{01} & Q_{02} & Q_{03} \cdots \\ Q_{10} & Q_{11} & Q_{12} & Q_{13} \cdots \\ \vdots & \vdots & \vdots & \vdots \cdots \end{bmatrix}$$

$n\_queries$

$\ast$

$n\_keys$

$n\_dimension$

**K**

$$\begin{bmatrix} K_{00} & K_{01} & K_{02} & K_{03} \cdots \\ K_{10} & K_{11} & K_{12} & K_{13} \cdots \\ \vdots & \vdots & \vdots & \vdots \cdots \end{bmatrix}$$

$n\_dimension$

**Q**

$$\begin{bmatrix} Q_{00} & Q_{01} & Q_{02} & Q_{03} \cdots \\ Q_{10} & Q_{11} & Q_{12} & Q_{13} \cdots \\ \vdots & \vdots & \vdots & \vdots \cdots \end{bmatrix}$$

$n\_queries$

$\ast$

$n\_dimension$

$\mathbf{K^T}$

$n\_keys$

$$\begin{bmatrix} K_{00} & K_{10} & \cdots \\ K_{01} & K_{11} & \cdots \\ K_{02} & K_{12} & \cdots \\ K_{03} & K_{13} & \cdots \\ \vdots & \vdots & \cdots \end{bmatrix}$$

$=$

$\mathbf{QK^T}$ = similarity score

*Score(Query 'subject' , encoded verb )*

$n\_keys$

$n\_queries$

Score for each query/key pair

*Each cell is a score of similarity between Query 'key' and encoded 'key'*

*Query Verb Lookup 'key' with all encoded 'keys'*

$\mathrm{softmax}(\mathbf{QK^T})$

$n\_keys$

$n\_queries$

*Each row of weights sums to 1*

$\ast$

**V**

$n\_dimension$

$n\_keys$

$$\begin{bmatrix} V_{00} & V_{01} & V_{02} & V_{03} \cdots \\ V_{10} & V_{11} & V_{12} & V_{13} \cdots \\ \vdots & \vdots & \vdots & \vdots \cdots \end{bmatrix}$$

*Value 0*
*Value 1*

$=$

$\mathrm{softmax}(\mathbf{QK^T})\mathbf{V}$

$n\_dimension$

$n\_queries$

This is added to Transformer pipeline

*Row(1),Col(1) has a high score. Matrix Multiplication of the softmax and values will maximize value_1.*

*Score(Query 'verb' , encoded verb ) ≈ 0.92*

# MultiHead Attention Layer : Encoder Subsystem

**Multi-Head Attention**

$\mathbf{V} = \mathbf{K} = \mathbf{Q}$

**V**

Words Vectors

Positional Embedding

Input Embedding

Tokenize

Preprocess

Batch of Sentences

Sentence

**K**

n_dimension

$$n\_keys \begin{bmatrix} V_{00} & V_{01} & V_{02} & V_{03} \cdots \\ V_{10} & V_{11} & V_{12} & V_{13} \cdots \\ \vdots & \vdots & \vdots & \vdots \cdots \end{bmatrix}$$

**Q**

n_dimension

$$n\_queries \begin{bmatrix} V_{00} & V_{01} & V_{02} & V_{03} \cdots \\ V_{10} & V_{11} & V_{12} & V_{13} \cdots \\ \vdots & \vdots & \vdots & \vdots \cdots \end{bmatrix}$$

$*$

n_dimension

**K$^T$**

n_keys

$$n\_dimension \begin{bmatrix} V_{00} & V_{10} & \cdots \\ V_{01} & V_{11} & \cdots \\ V_{02} & V_{12} & \cdots \\ V_{03} & V_{13} & \cdots \\ \vdots & \vdots & \cdots \end{bmatrix}$$

$\mathbf{QK^T}$ = similarity score ( Squared Magnitude )

n_keys

n_queries

softmax($\mathbf{QK^T}$)

n_keys

n_queries

$*$

**V**

n_dimension

$$n\_keys \begin{bmatrix} V_{00} & V_{01} & V_{02} & V_{03} \cdots \\ V_{10} & V_{11} & V_{12} & V_{13} \cdots \\ \vdots & \vdots & \vdots & \vdots \cdots \end{bmatrix}$$

softmax($\mathbf{QK^T})\mathbf{V}$

n_dimension

n_keys

*Possible word representative vectors garner more attention*

# 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 )

$K$

n_dimension

$$\begin{bmatrix} V_{00} & V_{01} & V_{02} & V_{03} & \cdots \\ V_{10} & V_{11} & V_{12} & V_{13} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

n_keys

---

?

| Softmax |
| Subsystems |
| Multi-Head Attention    $V = K = Q$ |

$V$

$Q$

n_dimension

$$\begin{bmatrix} V_{00} & V_{01} & V_{02} & V_{03} & \cdots \\ V_{10} & V_{11} & V_{12} & V_{13} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

n_queries

$*$

n_dimension

$K^T$

n_keys

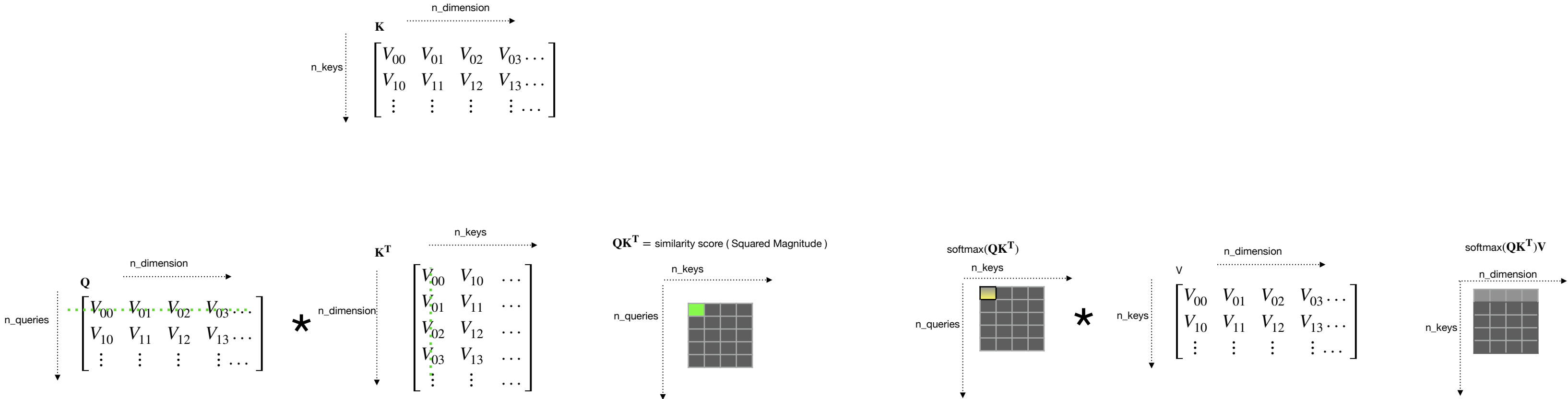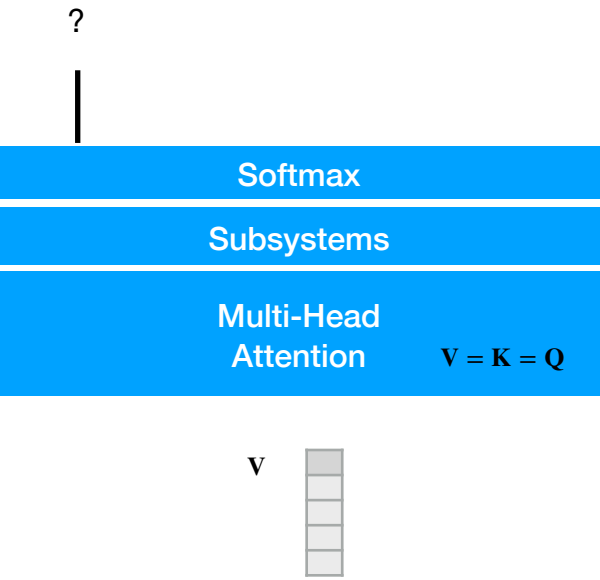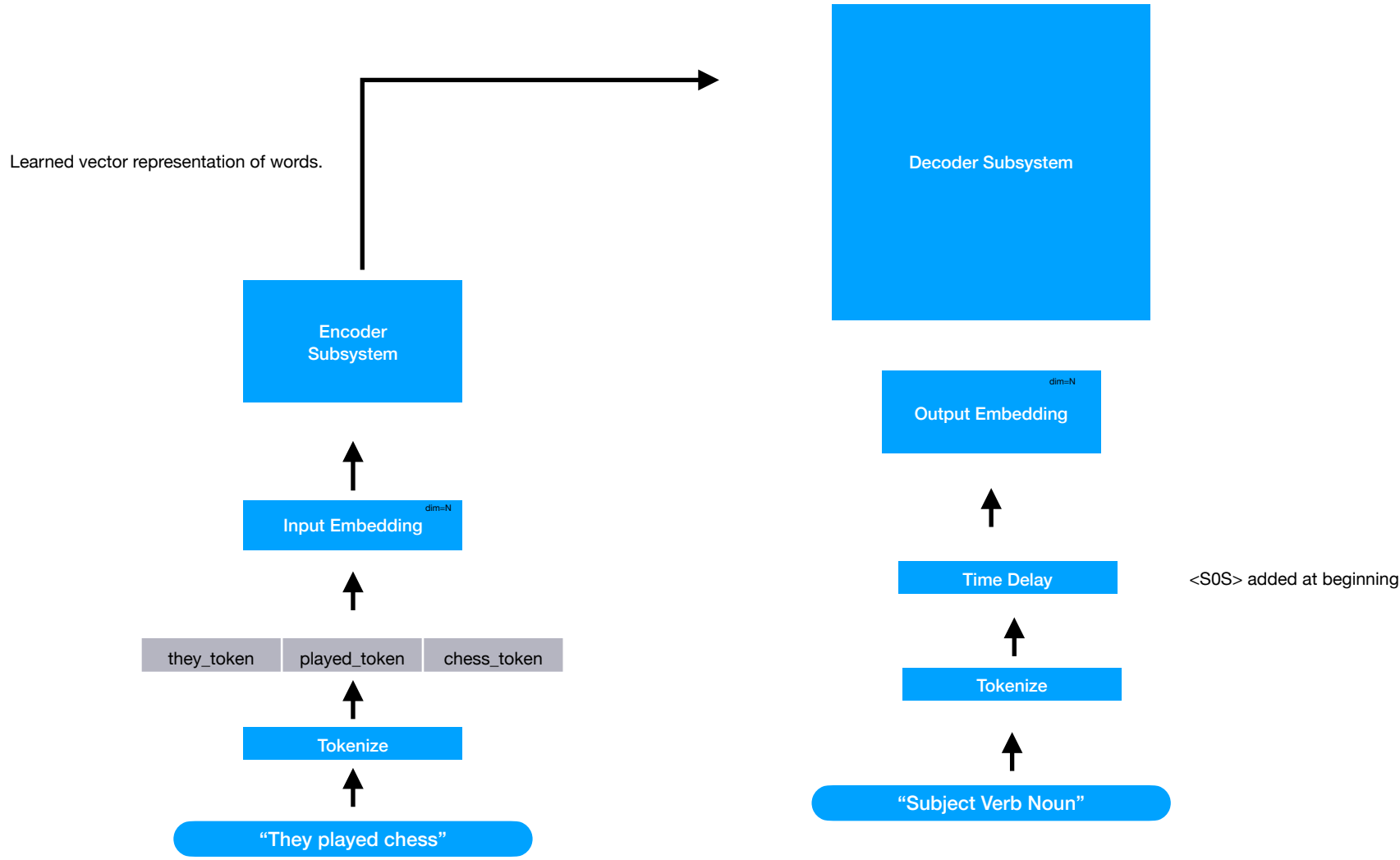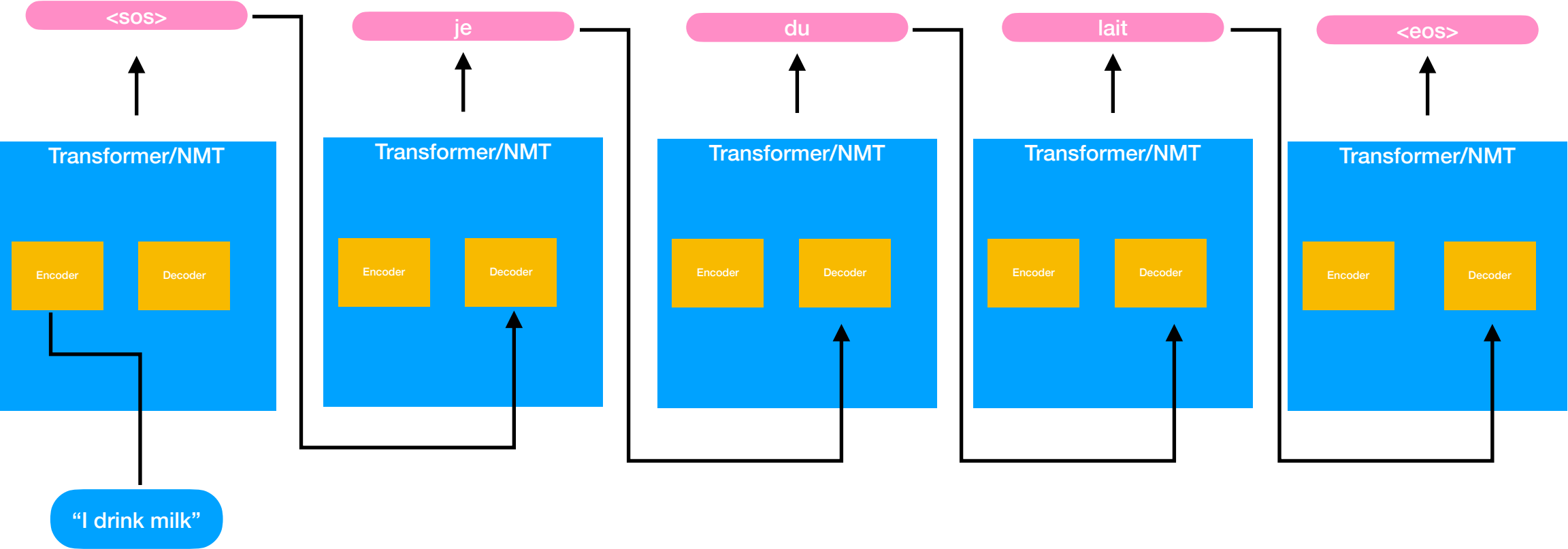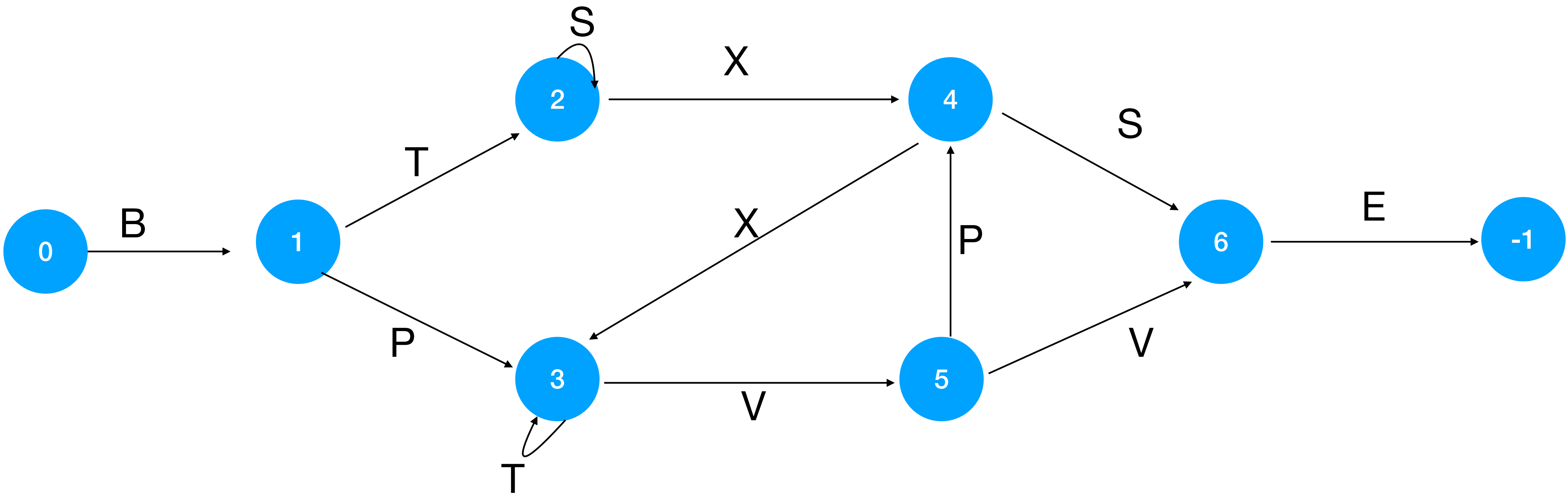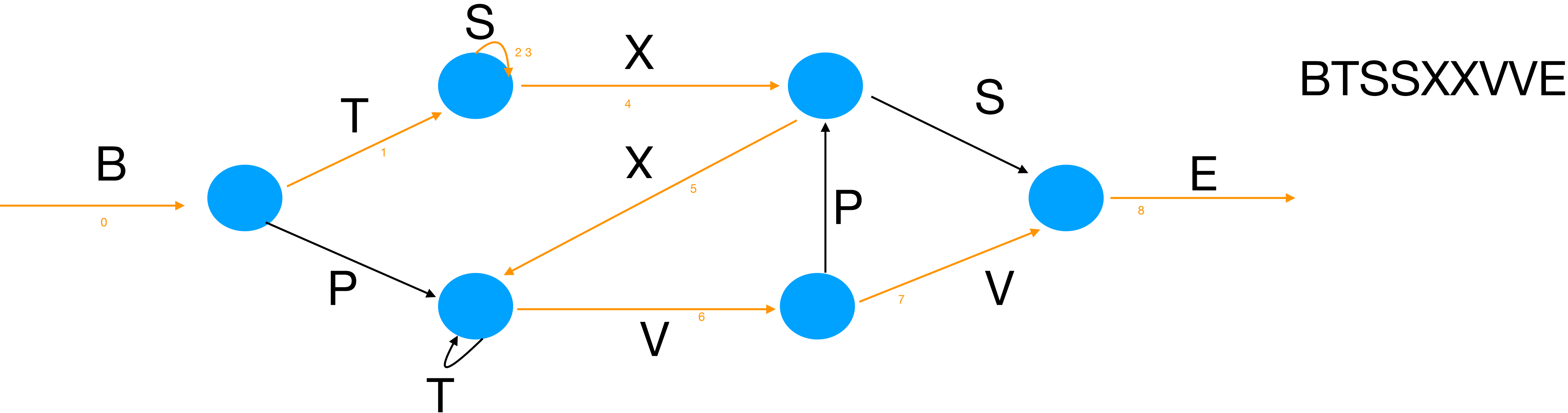$$\begin{bmatrix} V_{00} & V_{10} & \cdots \\ V_{01} & V_{11} & \cdots \\ V_{02} & V_{12} & \cdots \\ V_{03} & V_{13} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

$QK^T$ = similarity score ( Squared Magnitude )

n_keys

n_queries

softmax($QK^T$)

n_keys

n_queries

$*$

n_keys

$V$

n_dimension

$$\begin{bmatrix} V_{00} & V_{01} & V_{02} & V_{03} & \cdots \\ V_{10} & V_{11} & V_{12} & V_{13} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

n_keys

softmax($QK^T$)$V$

n_dimension

---

They    ?

| Softmax |
| Subsystems |
| Multi-Head Attention    $V = K = Q$ |

$V$

$Q$

n_dimension

$$\begin{bmatrix} V_{00} & V_{01} & V_{02} & V_{03} & \cdots \\ V_{10} & V_{11} & V_{12} & V_{13} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

n_queries

$*$

n_dimension

$K^T$

n_keys

$$\begin{bmatrix} V_{00} & V_{10} & \cdots \\ V_{01} & V_{11} & \cdots \\ V_{02} & V_{12} & \cdots \\ V_{03} & V_{13} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

$QK^T$ = similarity score ( Squared Magnitude )

n_keys

n_queries

softmax($QK^T$)

n_keys

n_queries

$*$

n_keys

$V$

n_dimension

$$\begin{bmatrix} V_{00} & V_{01} & V_{02} & V_{03} & \cdots \\ V_{10} & V_{11} & V_{12} & V_{13} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

n_keys

softmax($QK^T$)$V$

n_dimension

---

They    played    ?

| Softmax |
| Subsystems |
| Multi-Head Attention    $V = K = Q$ |

$V$

$Q$

n_dimension

$$\begin{bmatrix} V_{00} & V_{01} & V_{02} & V_{03} & \cdots \\ V_{10} & V_{11} & V_{12} & V_{13} & \cdots \\ V_{20} & V_{21} & V_{22} & V_{23} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

n_queries

$*$

n_dimension

$K^T$

n_keys

$$\begin{bmatrix} V_{00} & V_{10} & V_{20} & \cdots \\ V_{01} & V_{11} & V_{21} & \cdots \\ V_{02} & V_{12} & V_{22} & \cdots \\ V_{03} & V_{13} & V_{23} & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

$QK^T$ = similarity score ( Squared Magnitude )

n_keys

n_queries

softmax($QK^T$)

n_keys

n_queries

$*$

n_keys

$V$

n_dimension

$$\begin{bmatrix} V_{00} & V_{01} & V_{02} & V_{03} & \cdots \\ V_{10} & V_{11} & V_{12} & V_{13} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

n_keys

softmax($QK^T$)$V$

n_dimension

# Train/ Infer

Learned vector representation of words.

Encoder Subsystem

Input Embedding (default)

they_token | played_token | chess_token

Tokenize

"They played chess"

---

Decoder Subsystem

Output Embedding (default)

Time Delay — <S0S> added at beginning

Tokenize

"Subject Verb Noun"

---

<S0S> They

Encoded Vector → Encoder | Decoder → $State_t$

<S0S>

---

<S0S> They Played

Encoder | Decoder → $State_{t+1}$

They

---

<S0S> They Played

$State_{t+1}$ → Encoder | Decoder → $State_{t+2}$

Played

---

<S0S> They Played <E0S>

$State_{t+2}$ → Encoder | Decoder → $State_{t+3}$

Played

# Inference

Reber Grammer

Reber Grammer



BTSSXXVVE

... | | | word | word | **Text Vectorization** | ... | | | 532 | 22 | **From Tensor Slices** | ... | batch = 1 → 532 | batch = 1 → 22 | **GRU** | **GRU**

MDY to ISO

$Y_0$  $Y_1$  $Y_2$  ...

Encoder  Encoder  Encoder

states

$Y'_0$  $Y'_1$  $Y'_2$  $Y'_3$  $Y'_4$  $Y'_5$  ...

Decoder  Decoder  Decoder  Decoder  Decoder  Decoder

Target          Data

| id_2019 | $X_0$ |
| id_- | $X_1$ |
| id_04 | $X_2$ |
| id_- | $X_3$ |
| id_22 | $X_4$ |
| <eos> | $X_5$ |

| id_sos |
| id_2019 |
| id_- |
| id_04 |
| id_- |
| id_22 |

batch

| $X_0$ | id_2019 |
| $X_1$ | id_22 |
| $X_2$ | id_Aptril |

batch

Embedding

Embedding

| "sos" |
| "2019" |
| "-" |
| "04" |
| "-" |
| "22" |

batch

| "2019" |
| "22" |
| "April" |

batch

Vectorize

Vectorize

Preprocess

Preprocess

"April 22 2019"

"2019-04-22"