

CSCE 638 Natural Language Processing Foundation and Techniques

Lecture 7: Transformers

Kuan-Hao Huang

Spring 2025



(Some slides adapted from Chris Manning, Karthik Narasimhan, Danqi Chen, and Vivian Chen)

Project Sign-Up

- <https://docs.google.com/spreadsheets/d/15Rj4AovtHtIZxILbX1ydrw7IEylamXuV7Dtg7cBD2EU/edit?usp=sharing>
- 3~4 members per team
 - Form teams on your own
 - No solo teams (We have too many students!)

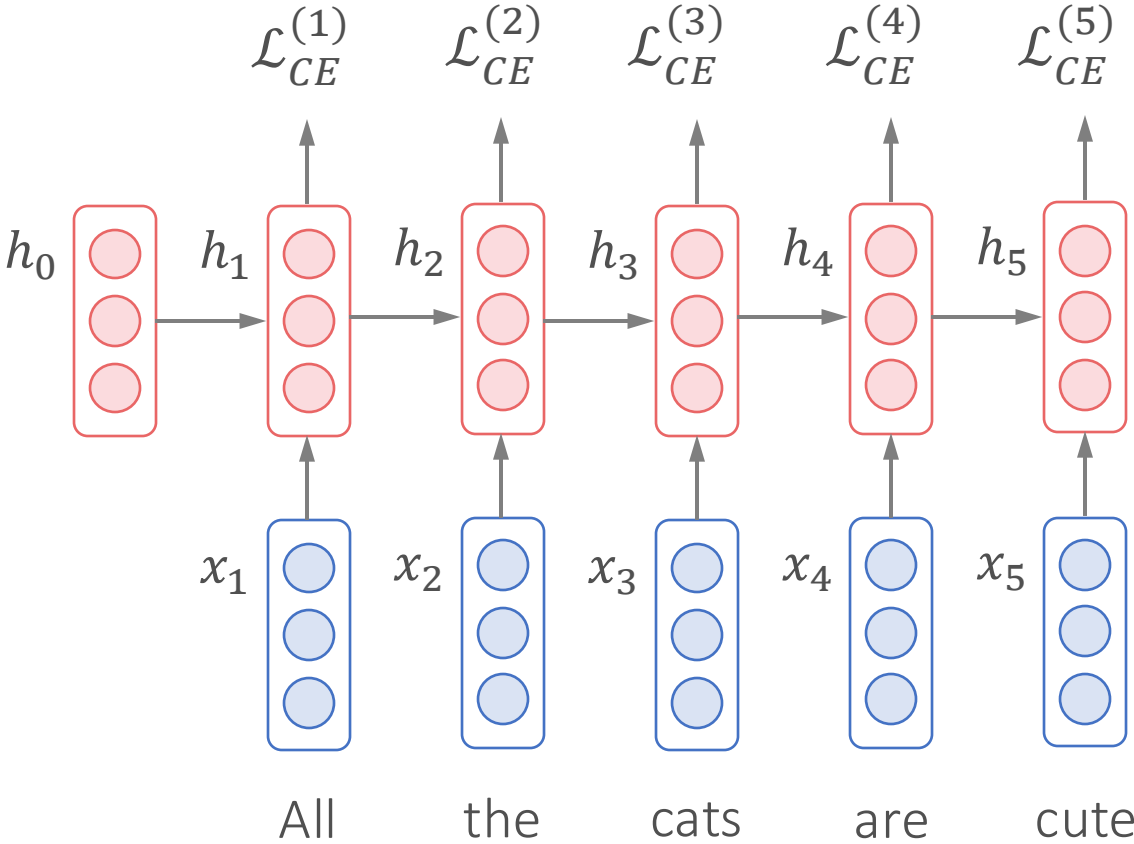
3~4 members per team									
	Project Topic	Member 1 (Name)	Member 1 (E-mail)	Member 2 (Name)	Member 2 (E-mail)	Member 3 (Name)	Member 3 (E-mail)	Member 4 (Name)	Member 4 (E-mail)
Team 1									
Team 2									
Team 3									
Team 4									
Team 5									
Team 6									
Team 7									
Team 8									
Team 9									
Team 10									
Team 11									
Team 12									

Lecture Plan

- Transformers
 - Attention
 - Self-Attention
 - Transformer Encoder
 - Positional Encoding

Recap: RNN as Encoder

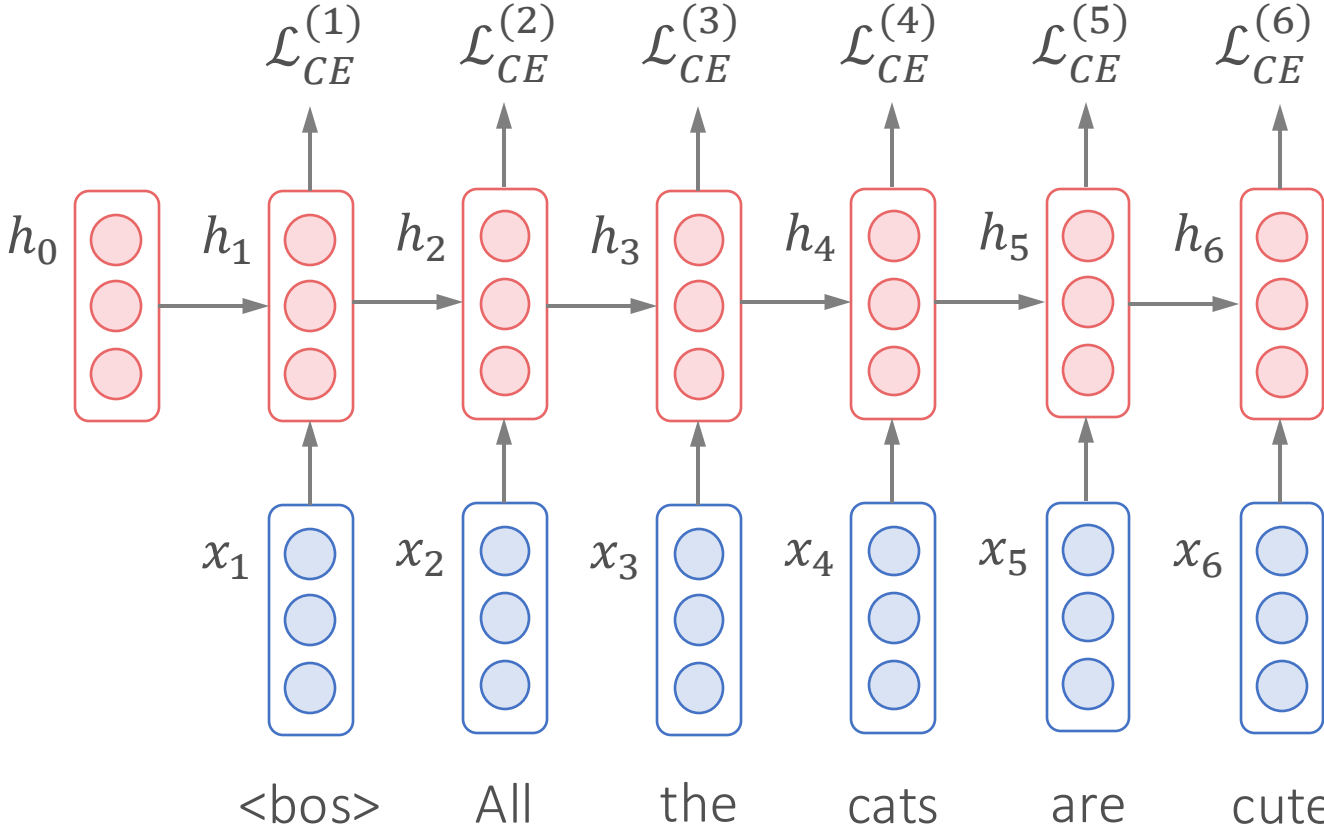
- Sequential labeling: A sequence of **dependent** classification



$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{CE}^{(i)}$$

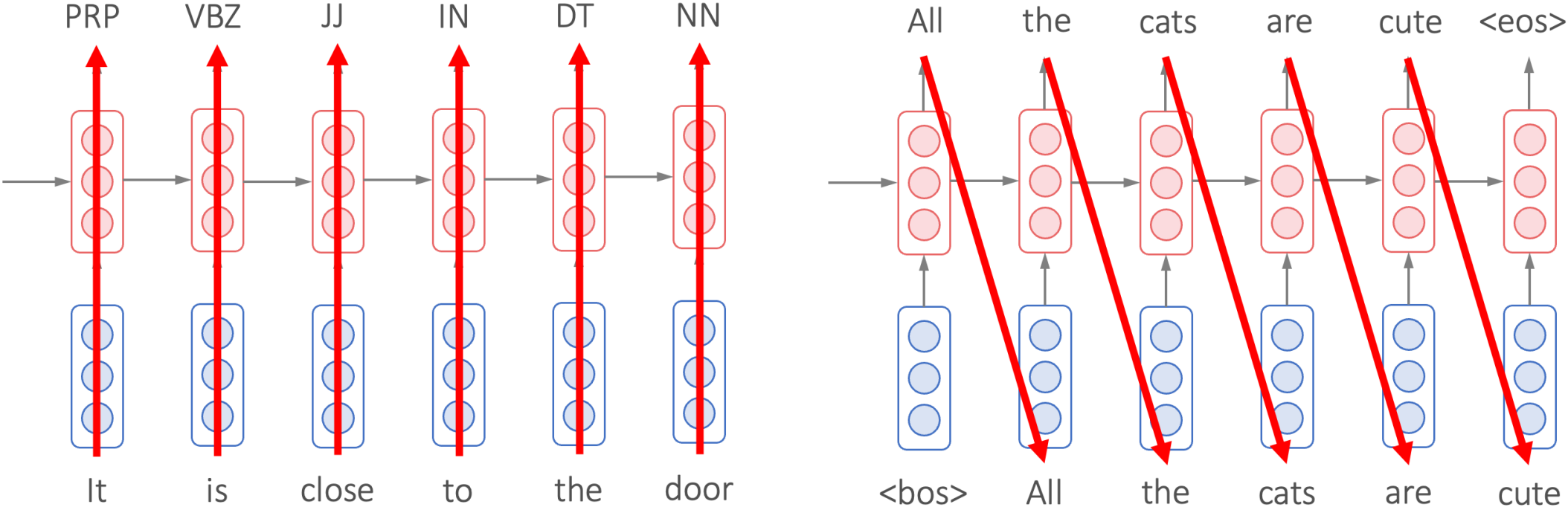
Recap: RNN as Decoder

- RNN Language Modeling
 - Generation is a sequence of word classification



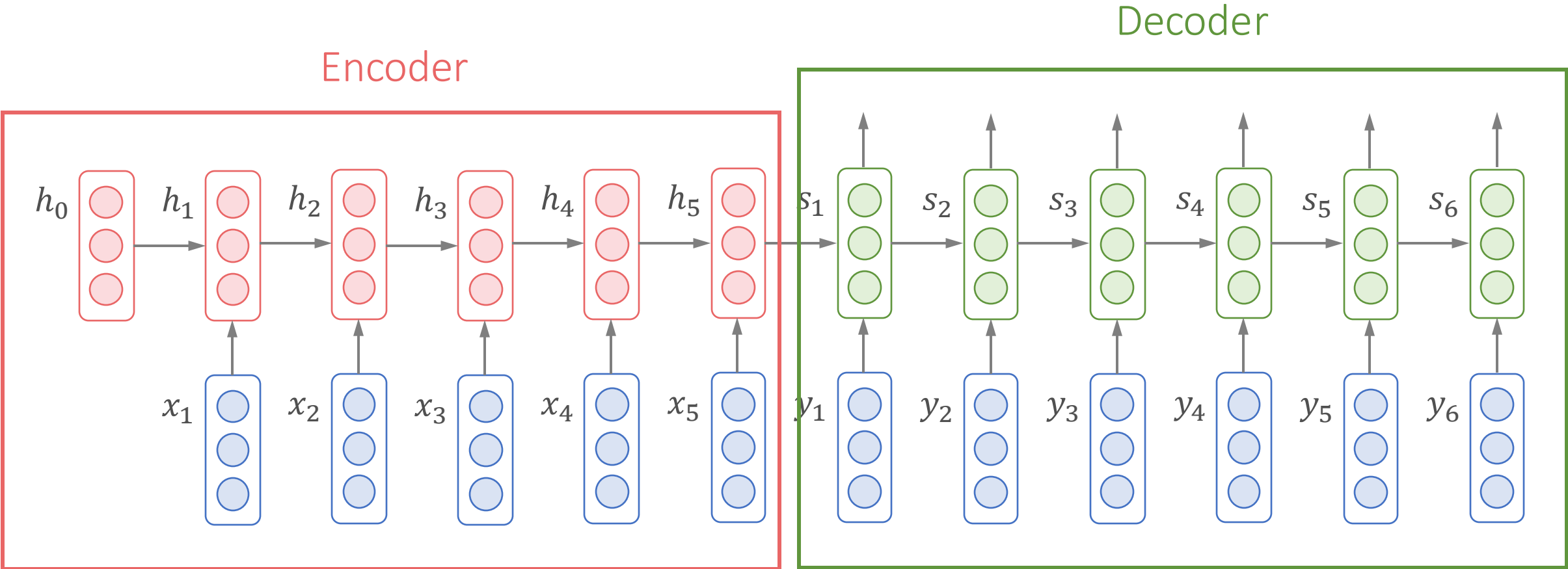
$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{CE}^{(i)}$$

Recap: Encoder vs. Decoder

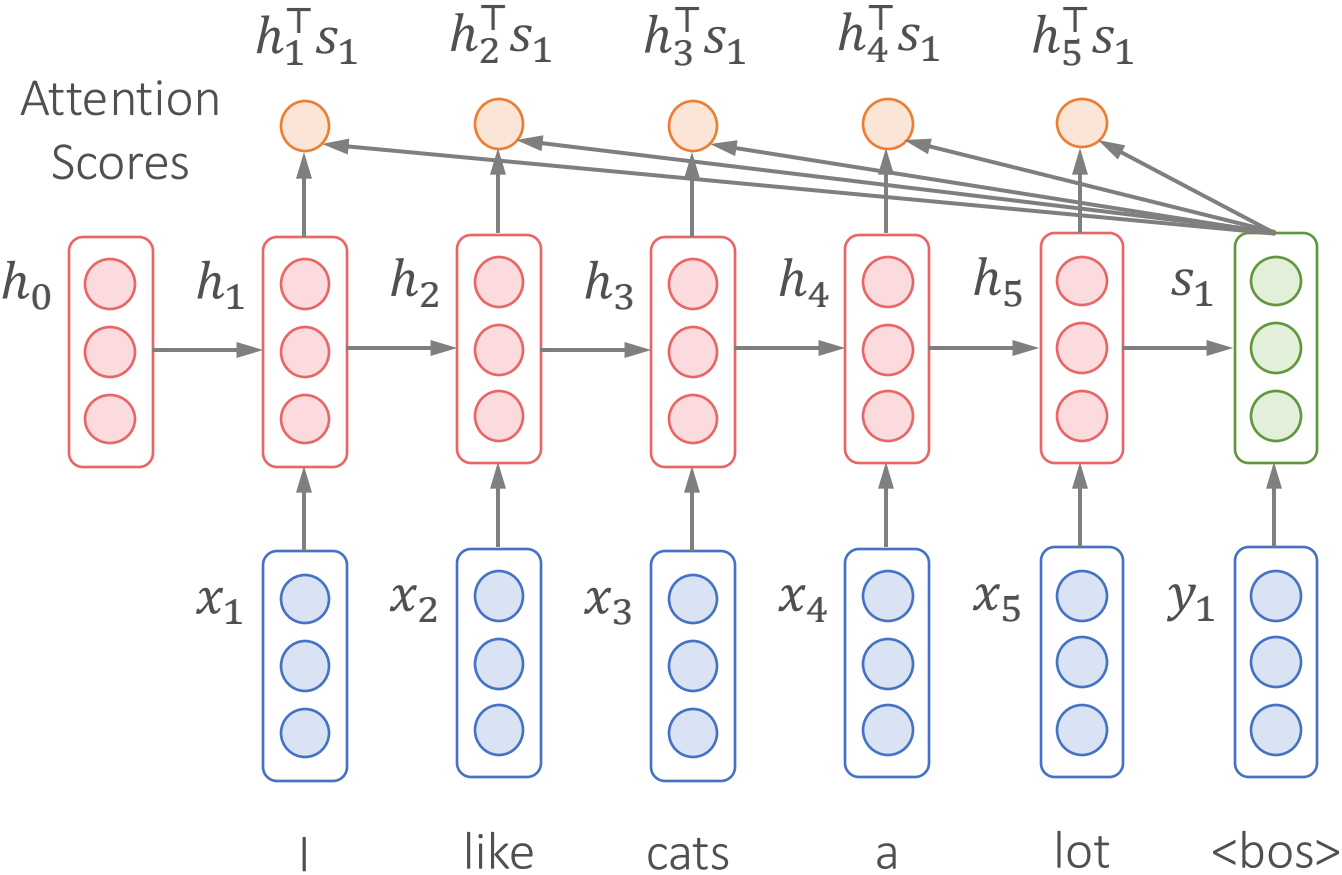


Recap: Sequence-to-Sequence Models (Seq2Seq)

- When we need understanding and generation at the same time

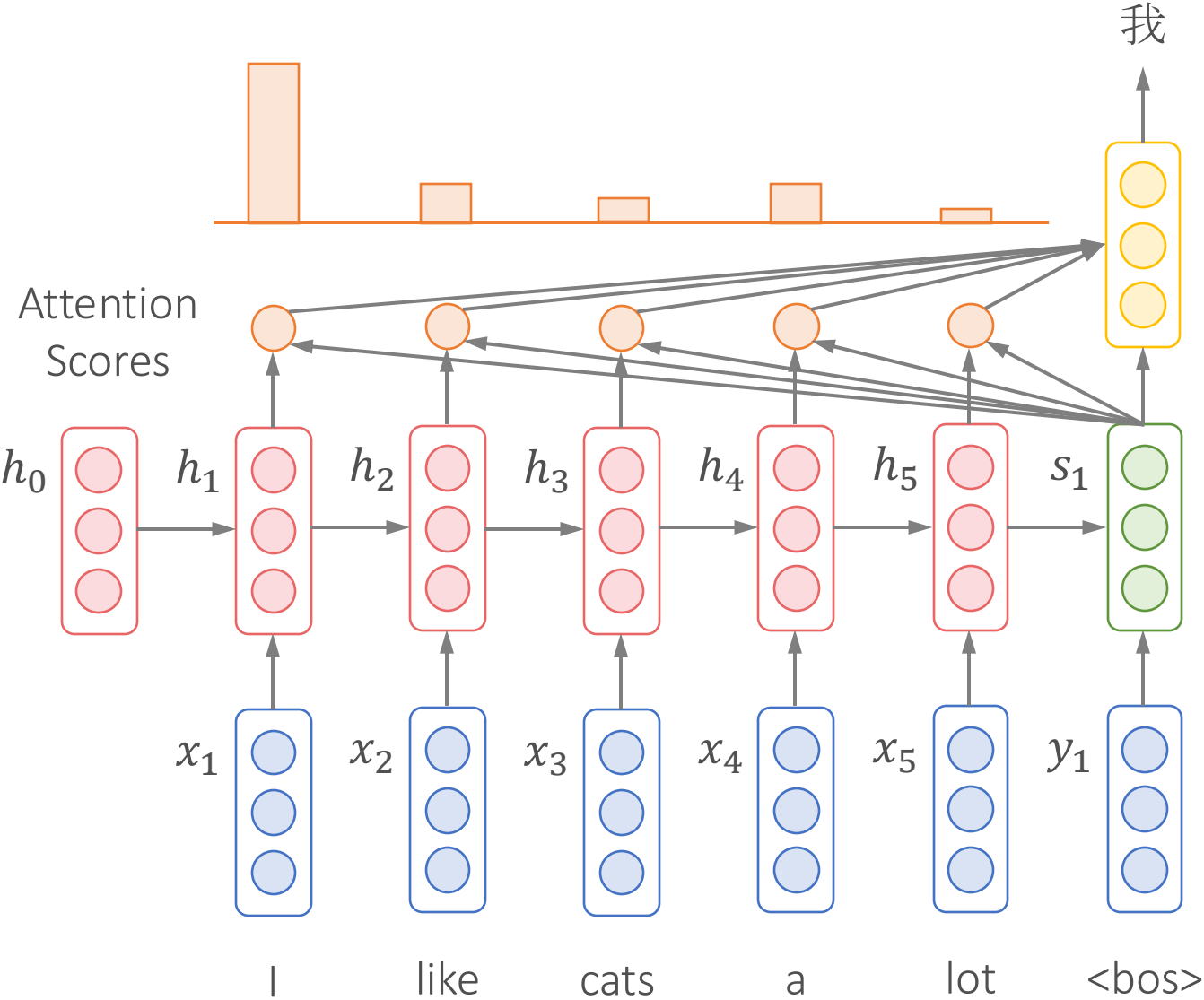


RNN with Attention



Attention Scores $\alpha_i = h_i^T s_1$

RNN with Attention



Attention Scores

$$\alpha_i = h_i^T s_1$$

Normalized Attention Scores

$$\hat{\alpha}_i = \text{softmax}(\alpha_i)$$

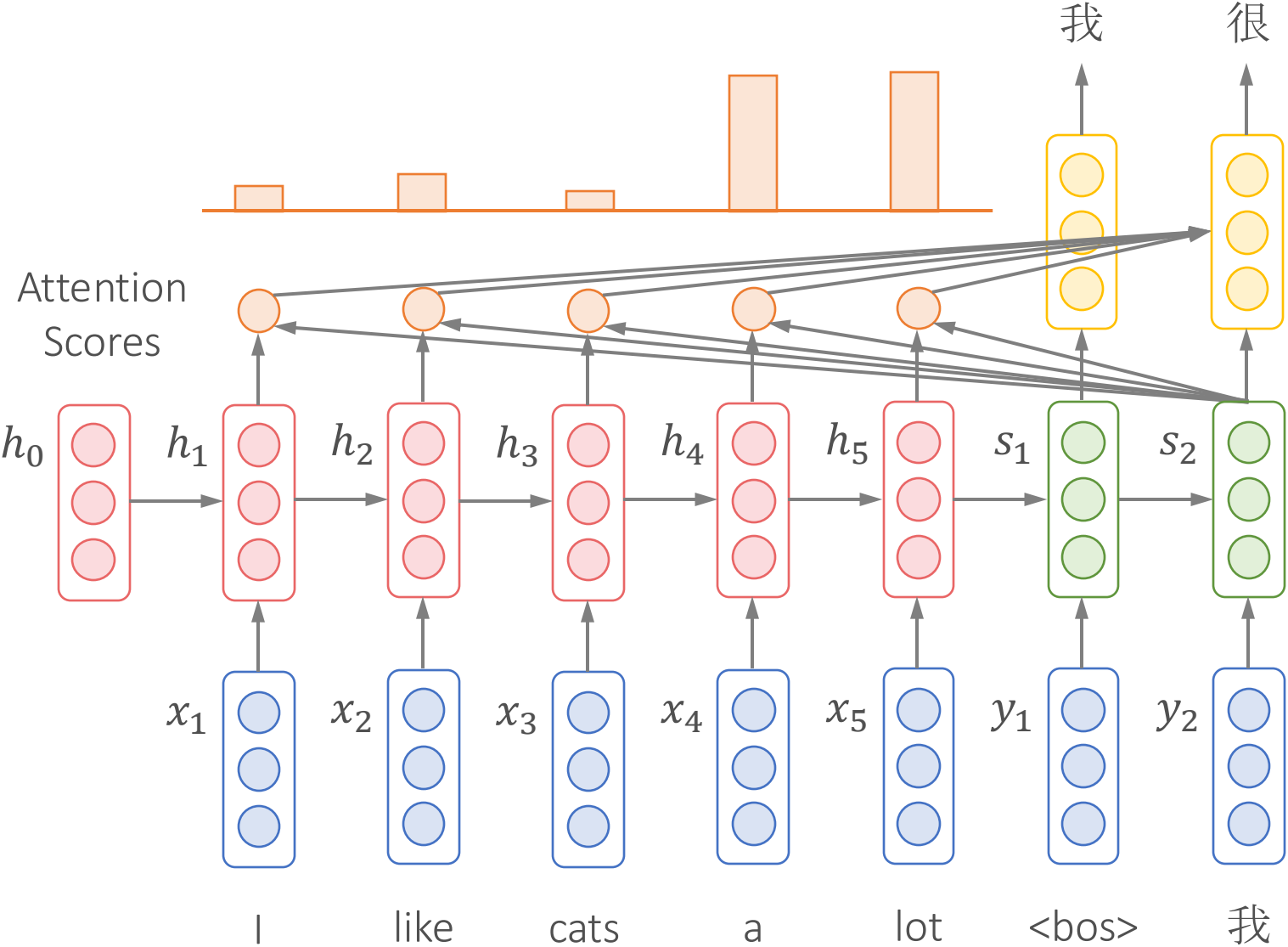
Weighted Sum

$$a = \sum_i \hat{\alpha}_i h_i$$

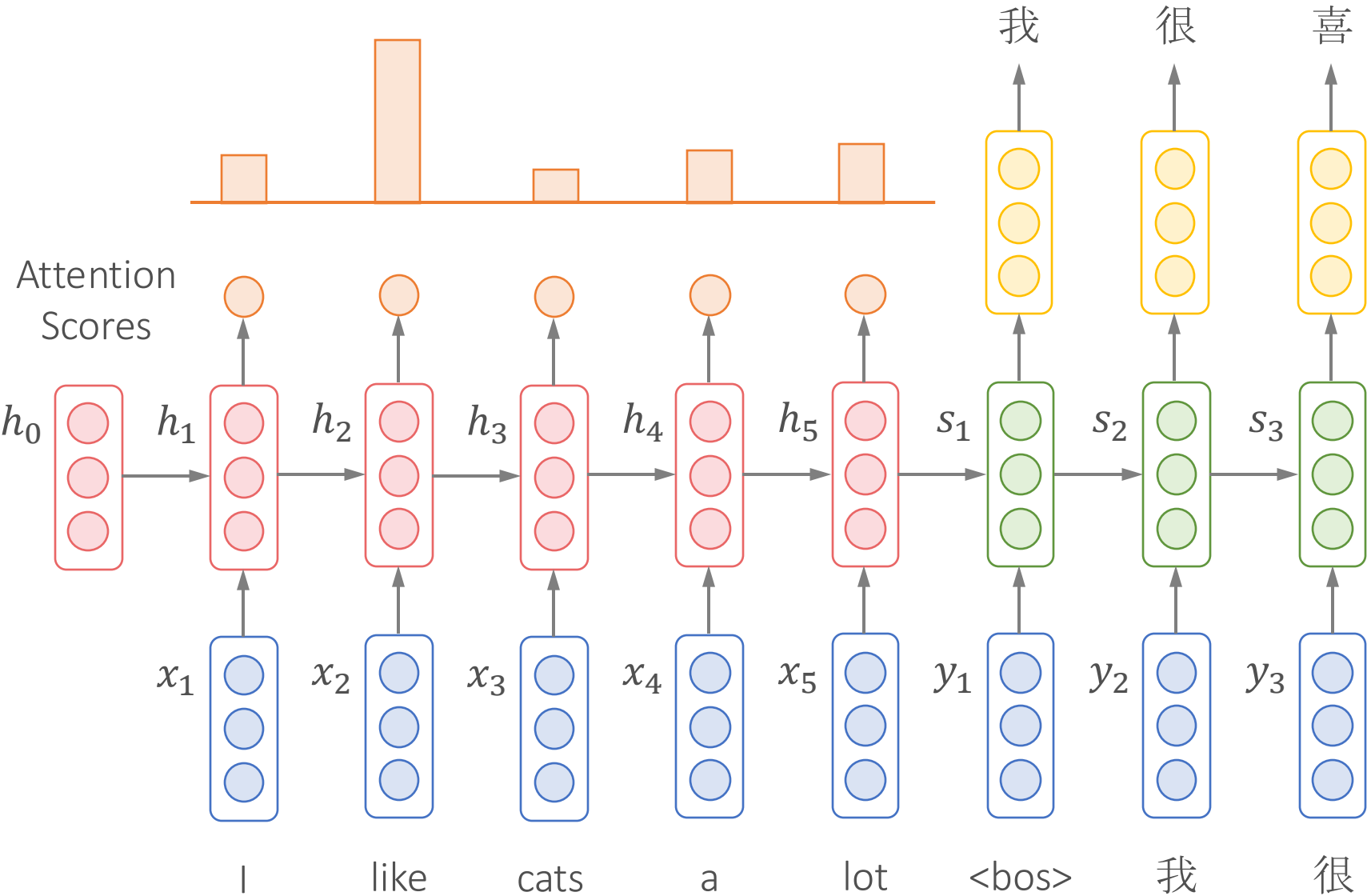
Attention Output

$$\tanh(\mathbf{W}[a; s_1])$$

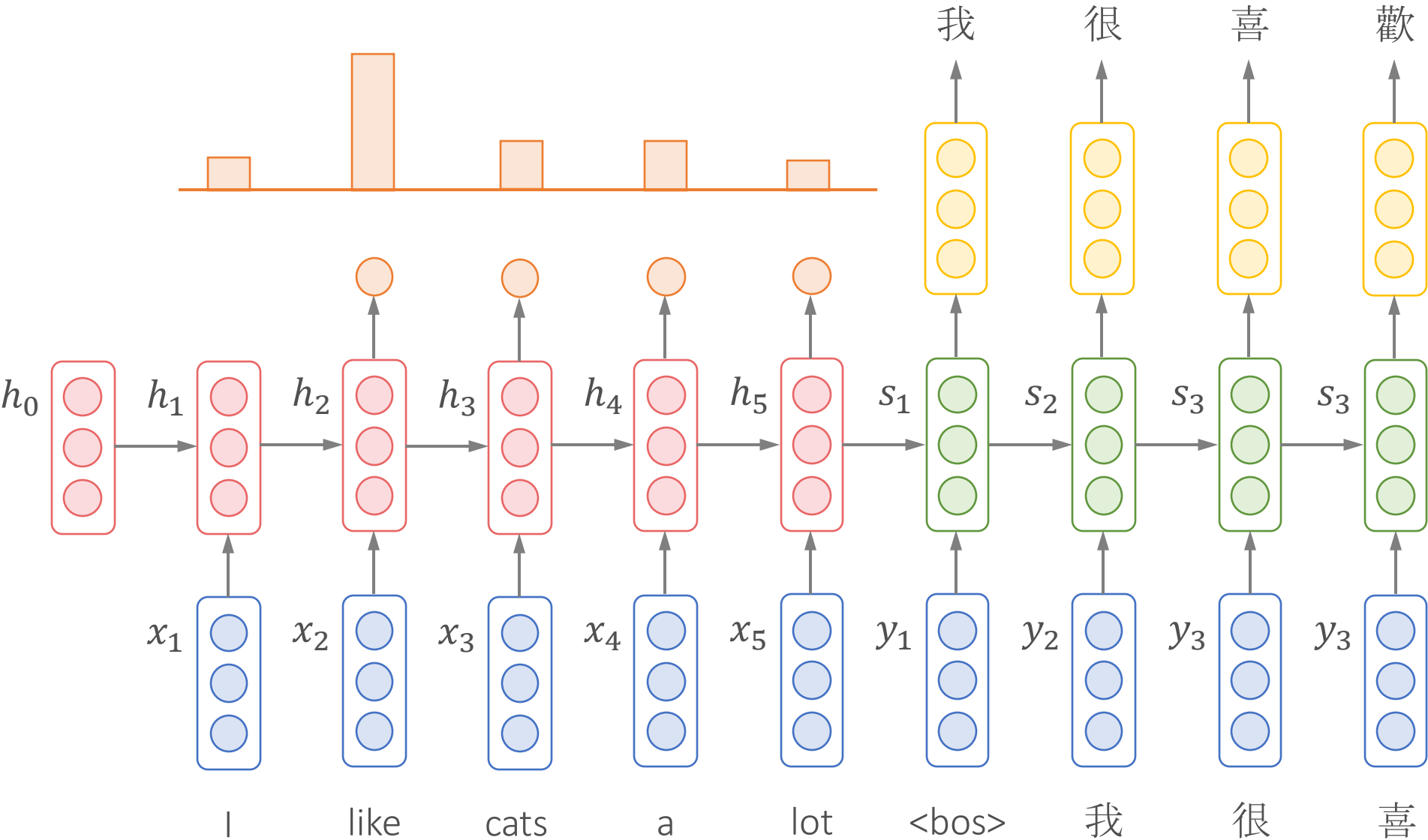
RNN with Attention



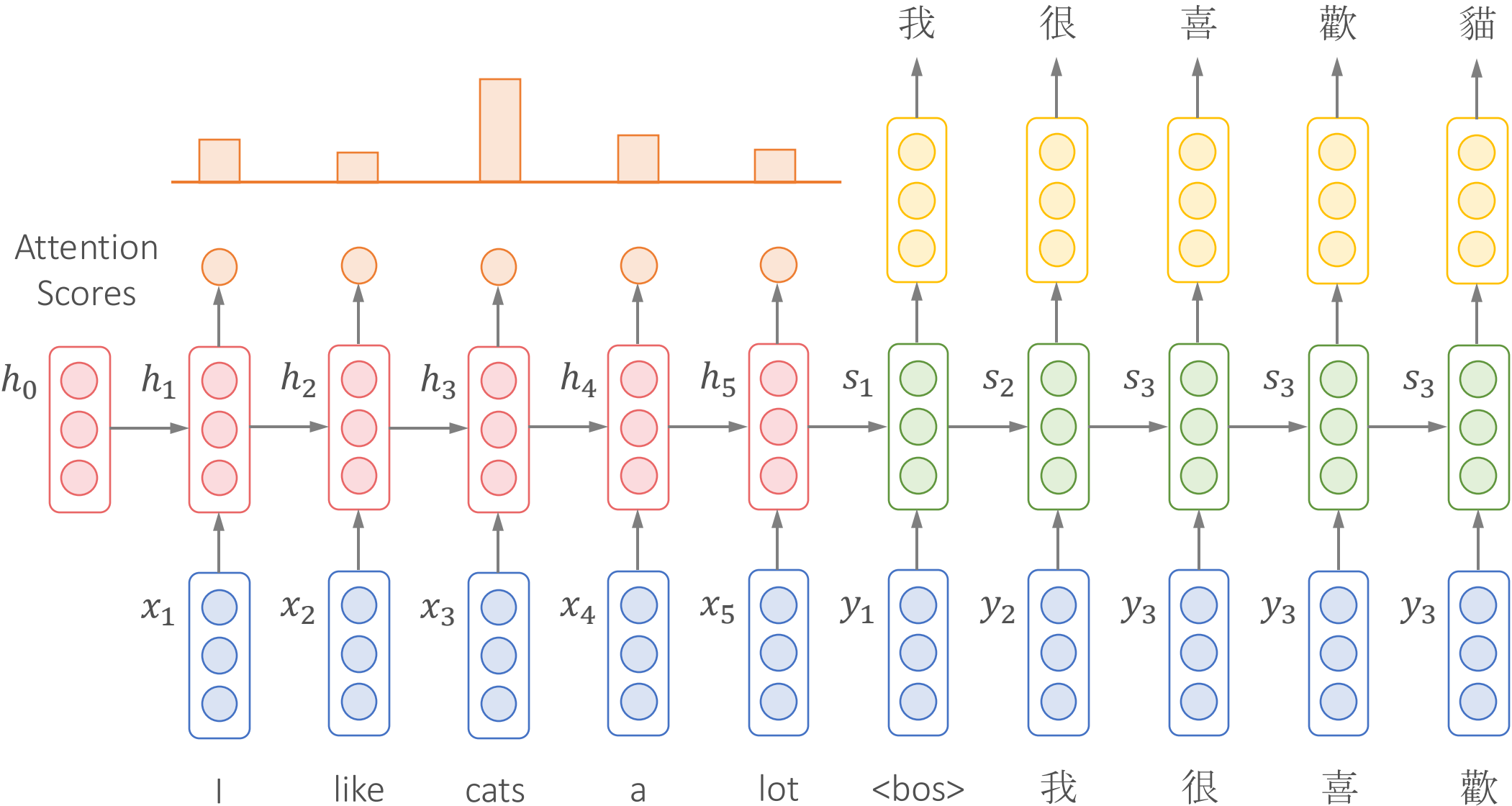
RNN with Attention



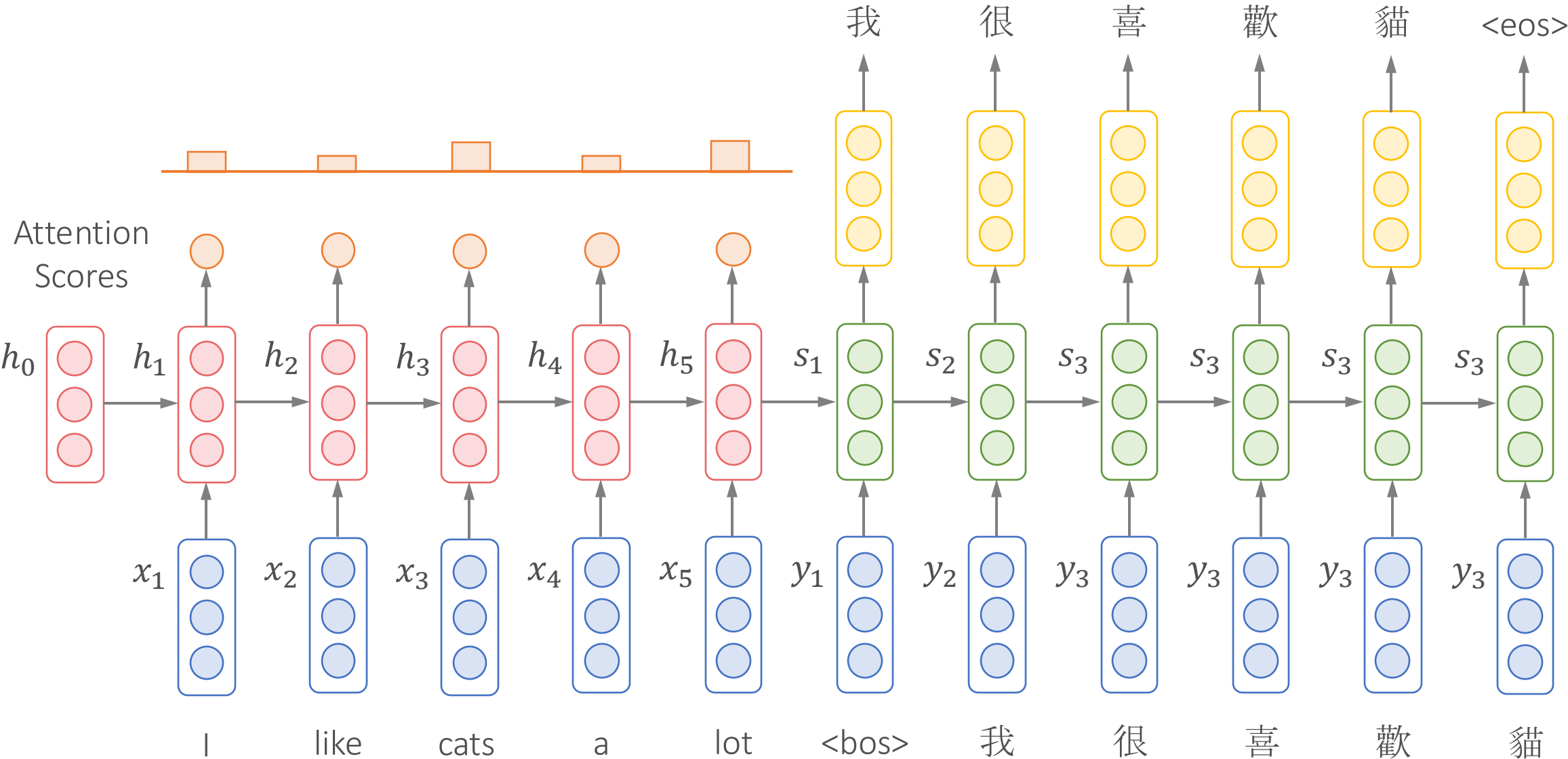
RNN with Attention



RNN with Attention



RNN with Attention



Different Types of Attention

Dot-Product Attention

$$h_i^T s_j$$

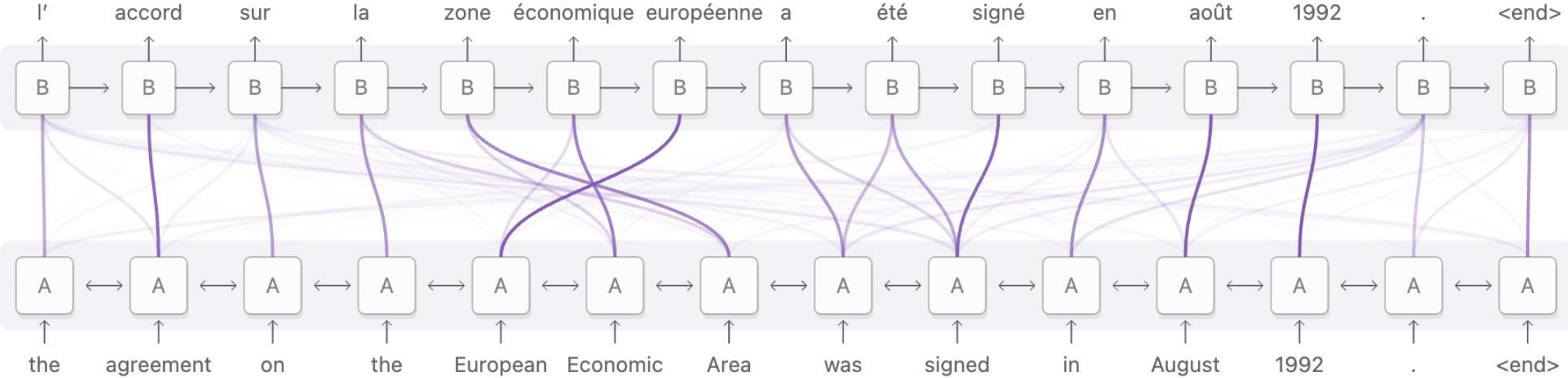
Multiplicative Attention

$$h_i^T W s_j$$

Additive Attention

$$v^T \tanh(W_1 h_i + W_2 s_j)$$

Machine Translation with Attention



Speech Recognition with Attention

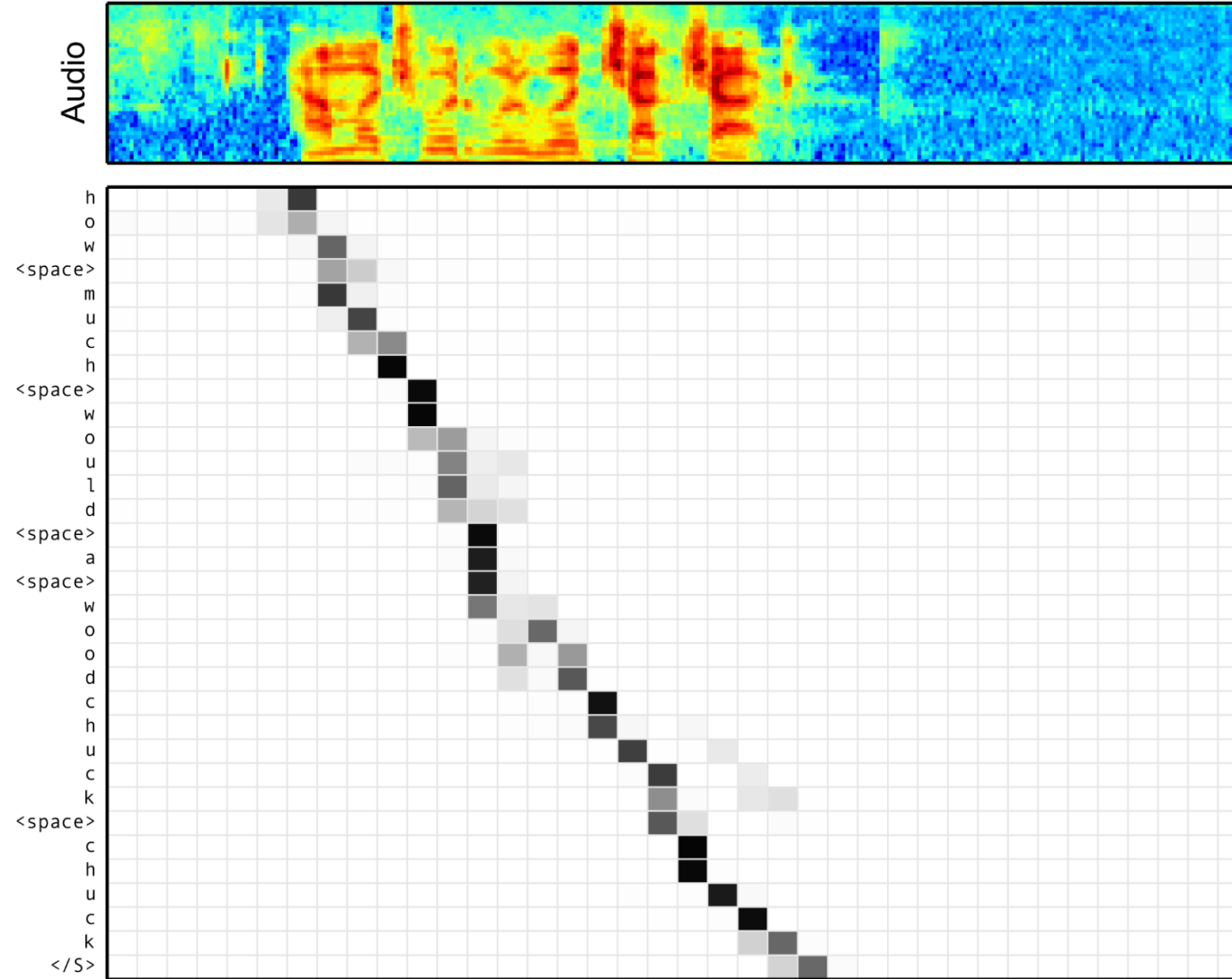
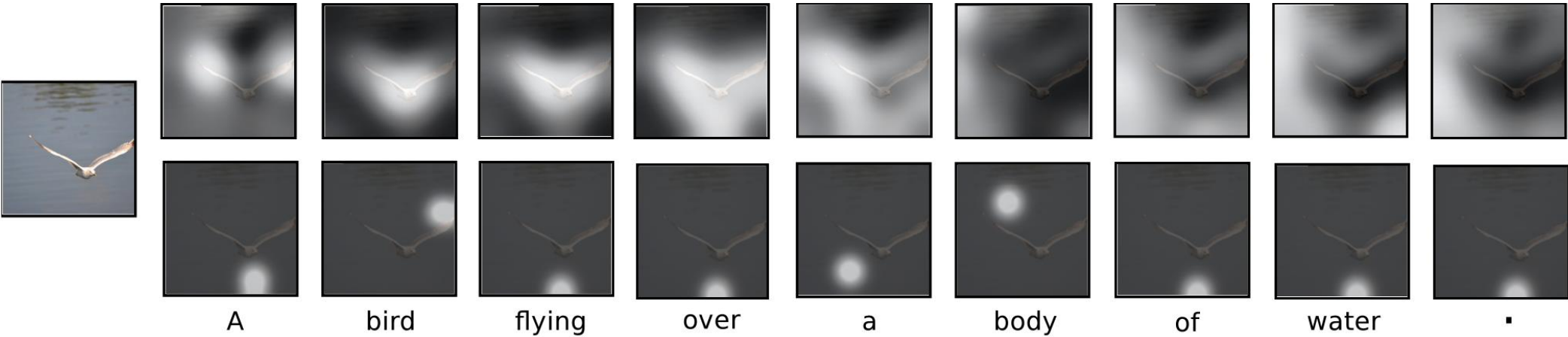


Image Captioning with Attention



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



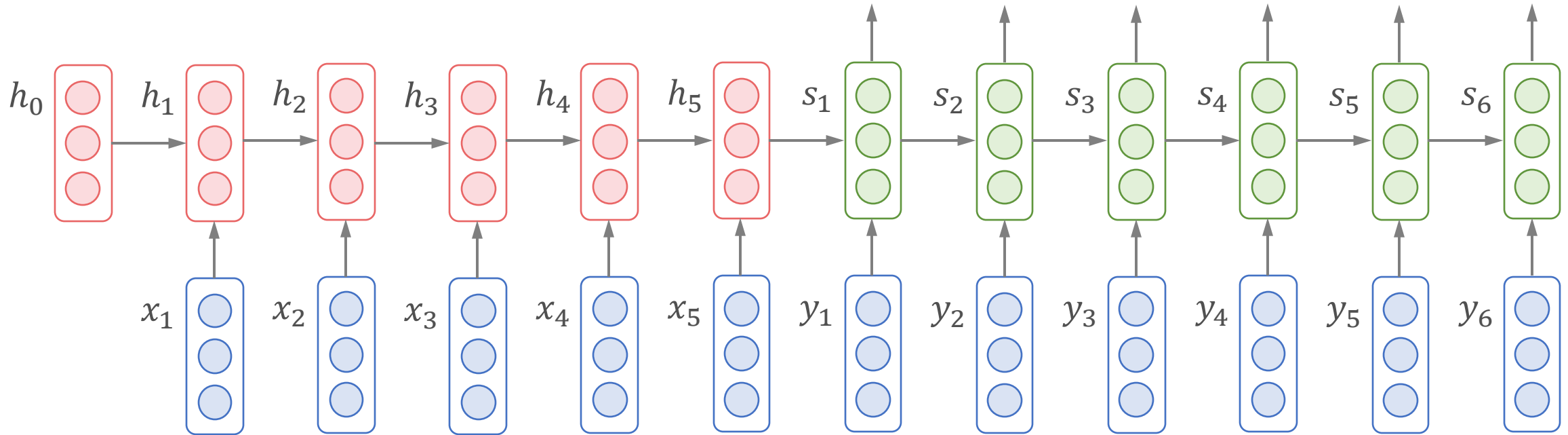
A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

Issues with RNN

- Longer sequences can lead to vanishing gradients → It is hard to capture long-distance information
- Lack parallelizability



Transformers: Attention Is All You Need!

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

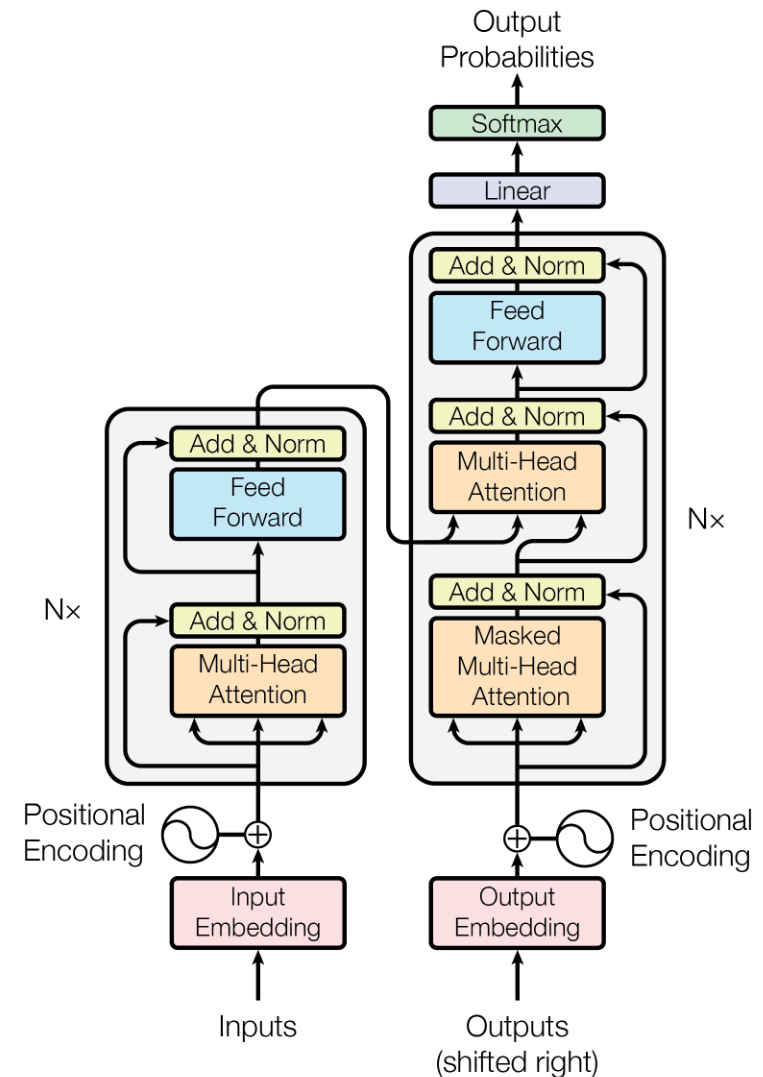
Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

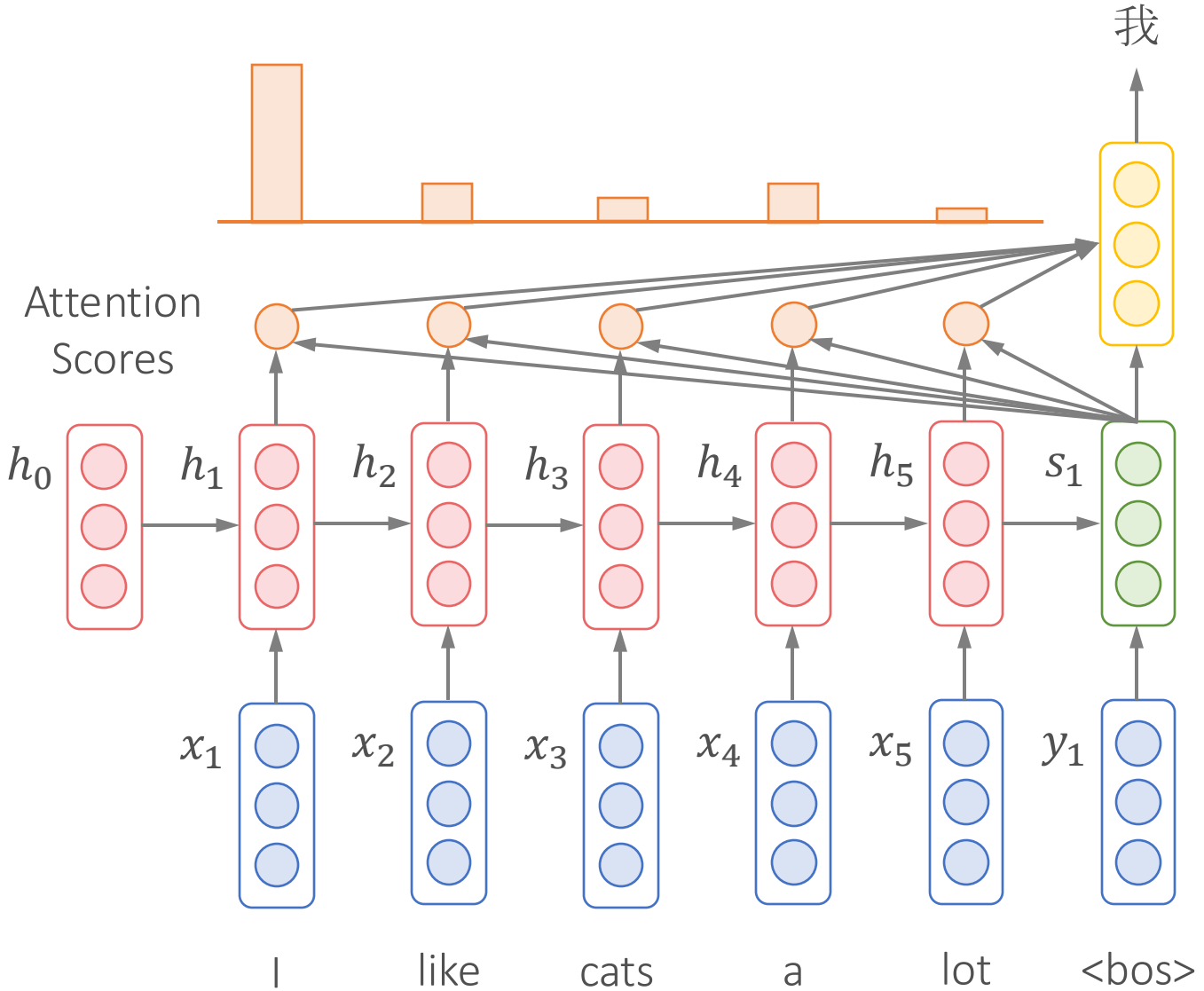
Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com



Look Back at RNN with Attention



Attention Scores

$$\alpha_i = h_i^T s_1$$

Normalized Attention Scores

$$\hat{\alpha}_i = \text{softmax}(\alpha_i)$$

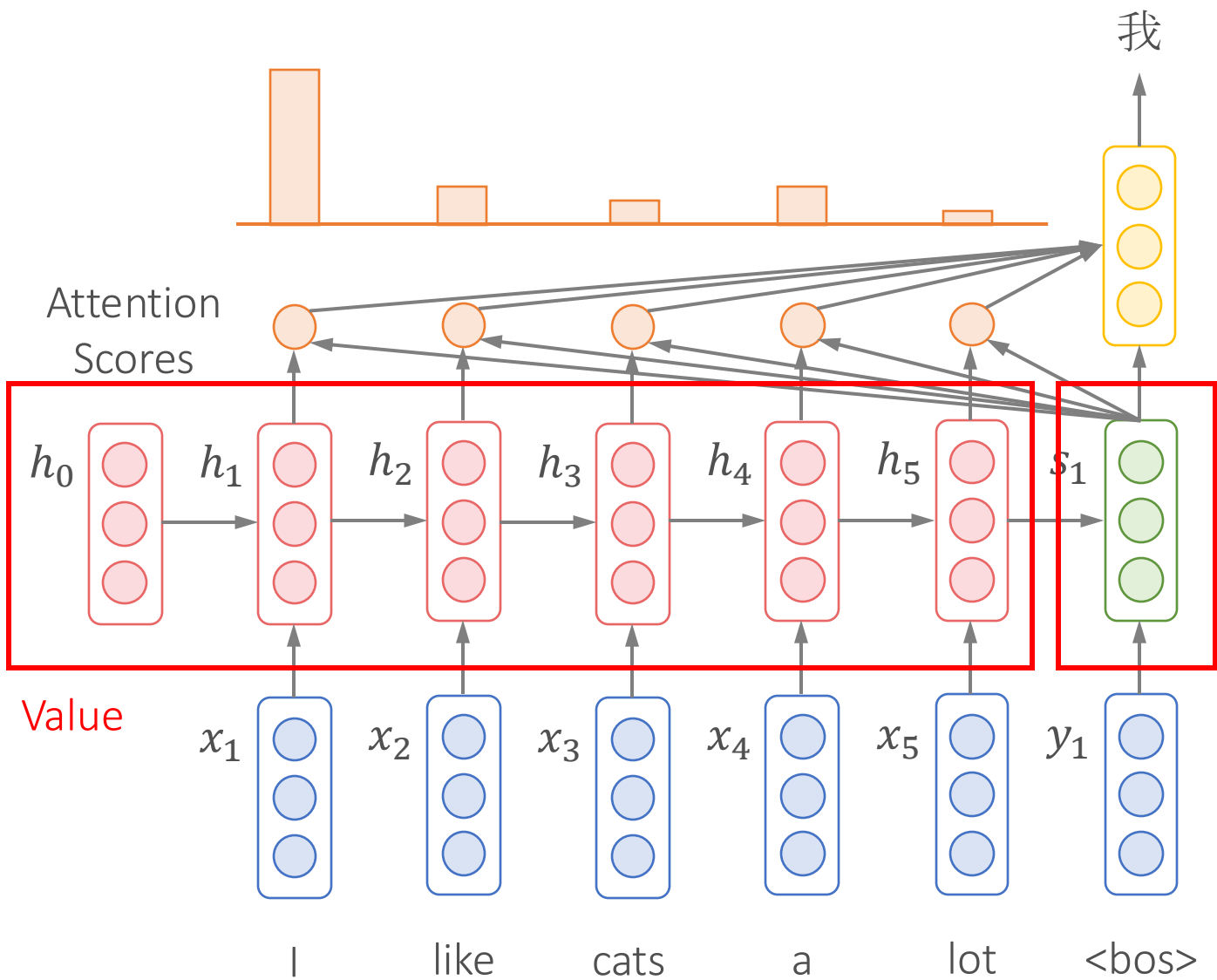
Weighted Sum

$$a = \sum_i \hat{\alpha}_i h_i$$

Attention Output

$$\tanh(\mathbf{W}[a; s_1])$$

Look Back at RNN with Attention



Inner Product of Query and Values

Attention Scores

$$\alpha_i = h_i^T s_1$$

Normalized Attention Scores

$$\hat{\alpha}_i = \text{softmax}(\alpha_i)$$

Weighted Sum

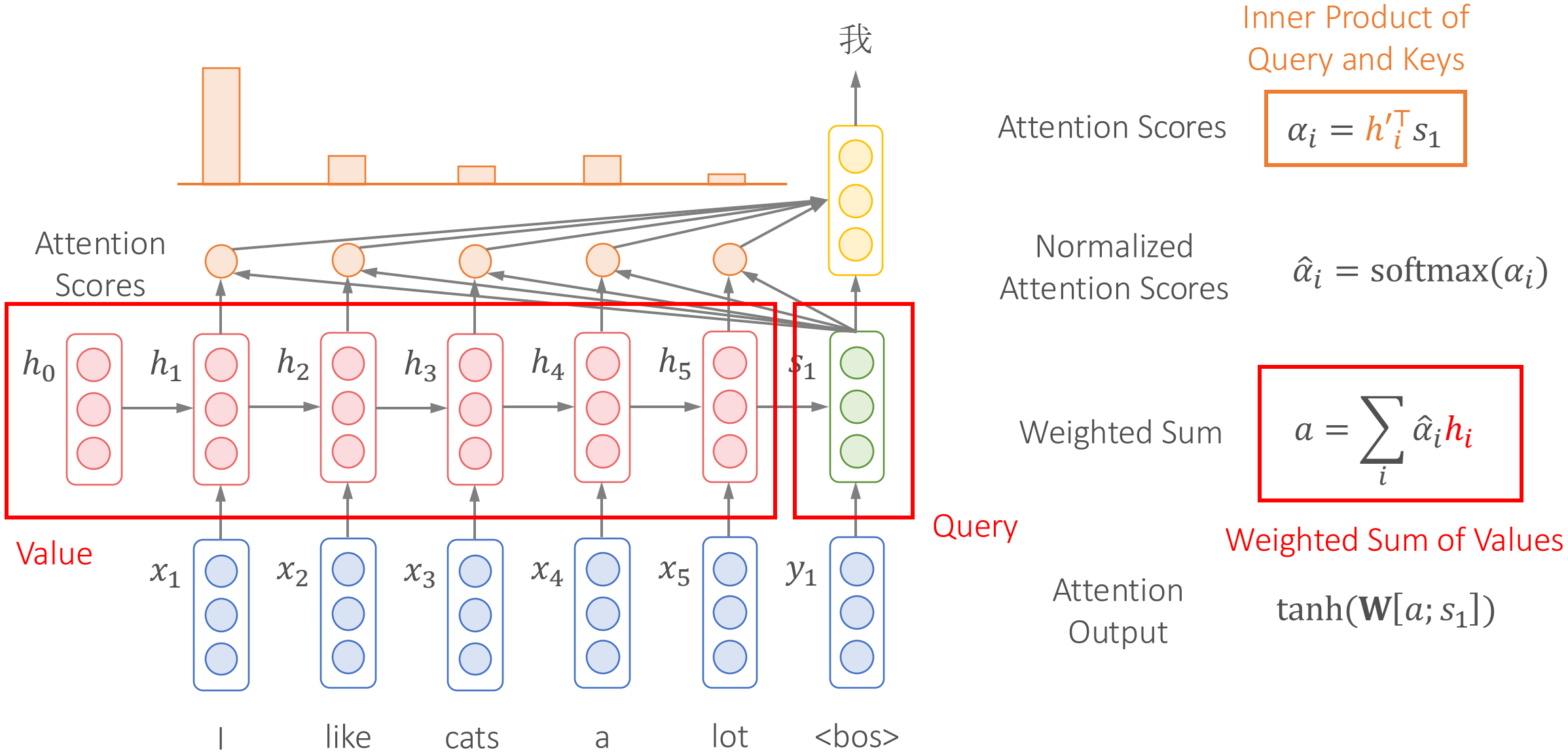
$$a = \sum_i \hat{\alpha}_i h_i$$

Weighted Sum of Values

Attention Output

$$\tanh(\mathbf{W}[a; s_1])$$

Look Back at RNN with Attention – General Version



Inner Product of Query and Keys

$$\alpha_i = h_i^T s_1$$

Attention Scores

Normalized Attention Scores

$$\hat{\alpha}_i = \text{softmax}(\alpha_i)$$

Weighted Sum

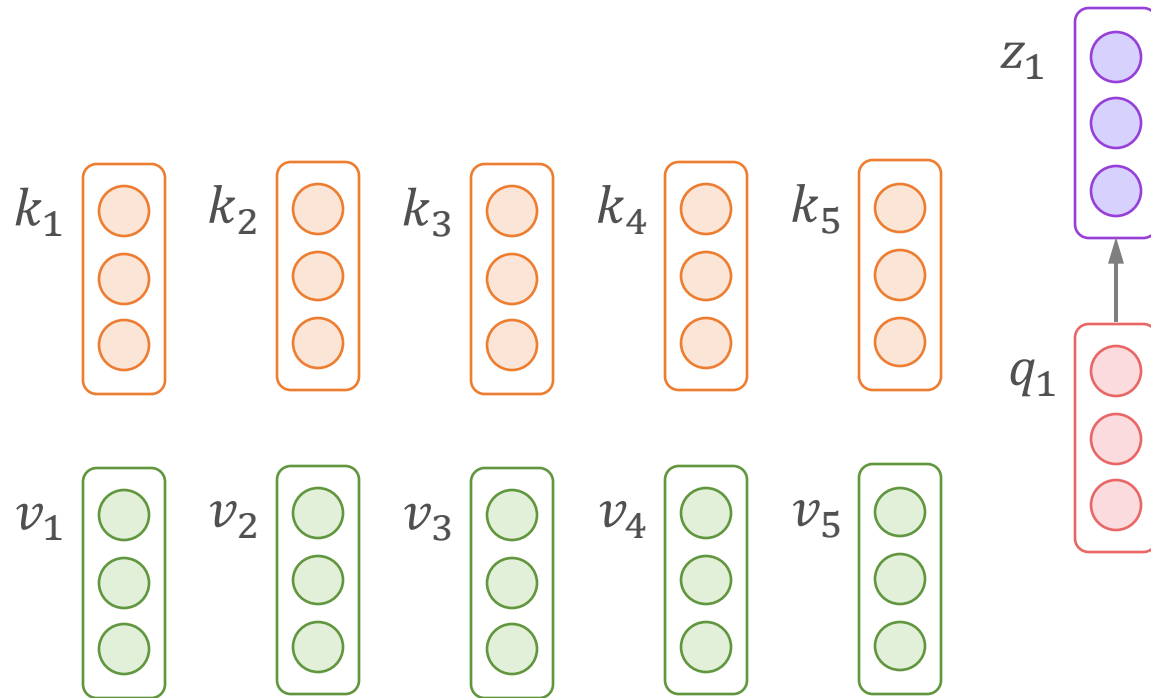
$$a = \sum_i \hat{\alpha}_i h_i$$

Weighted Sum of Values

Attention Output

$$\tanh(\mathbf{W}[a; s_1])$$

Attention – General Version



Attention Scores

$$\alpha_i = k_i^T q_1$$

Normalized
Attention Scores

$$\hat{\alpha}_i = \text{softmax}(\alpha_i)$$

Weighted Sum

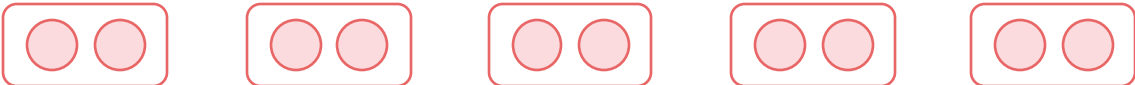
$$z_1 = \sum_i \hat{\alpha}_i v_i$$

From Attention to Self-Attention

- Self-attention = attention from the sequence to itself
 - The queries, keys and values come from the same source
- Any word can be a **query**
- Any word can be a **key**
- Any word can be a **value**

Self-Attention

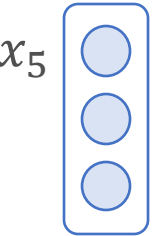
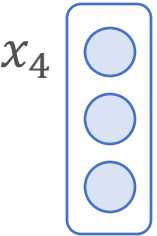
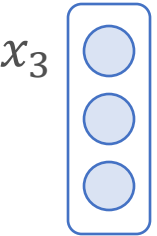
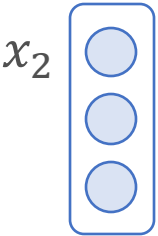
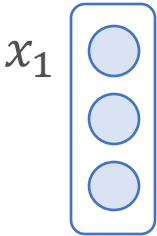
Query $q_i = W^Q x_i$



Key $k_i = W^K x_i$



Value $v_i = W^V x_i$



I

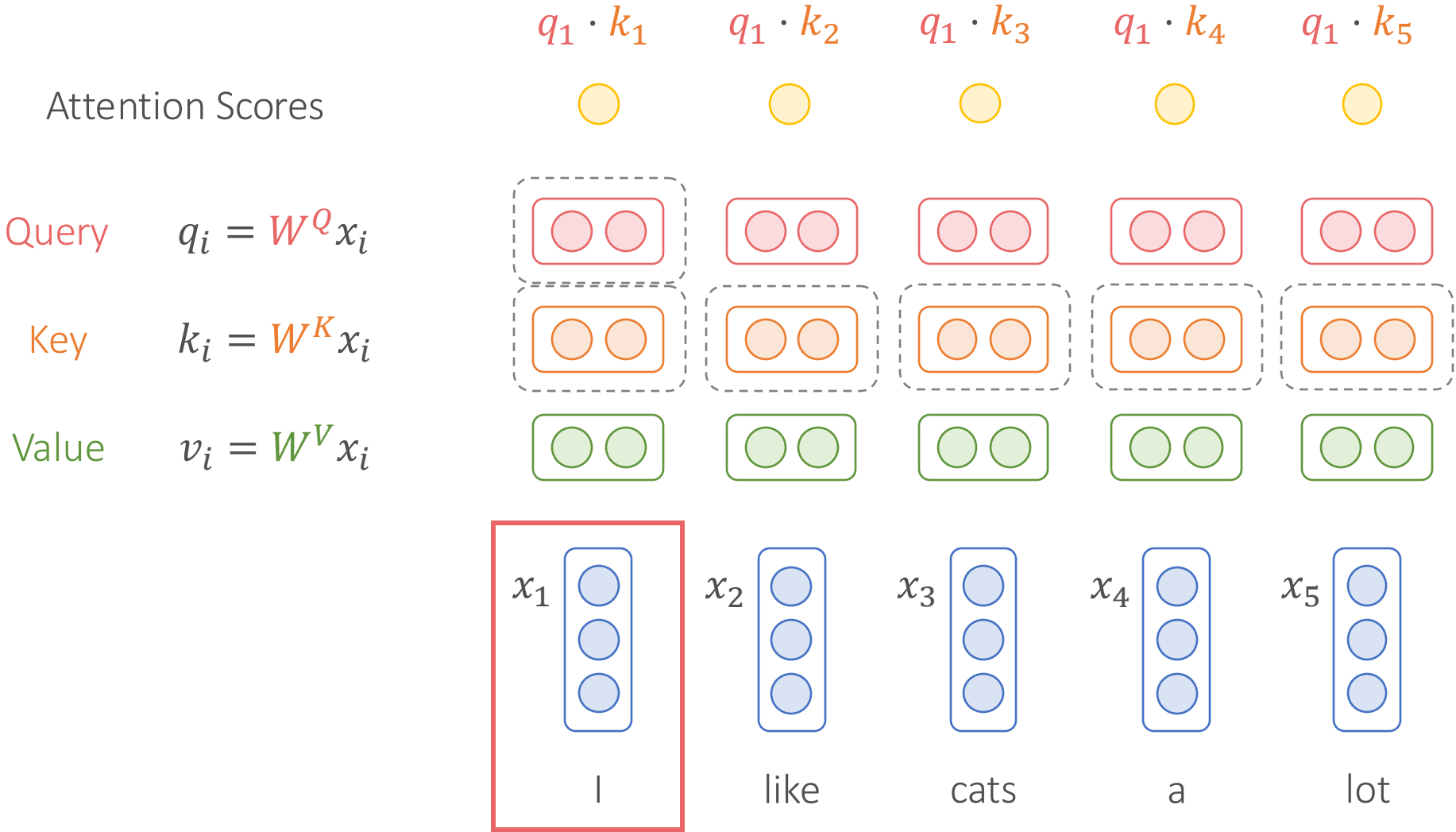
like

cats

a

lot

Self-Attention



Self-Attention

$$\alpha_{1,i} = \text{softmax}\left(\frac{q_1 \cdot k_i}{\sqrt{d}}\right)$$

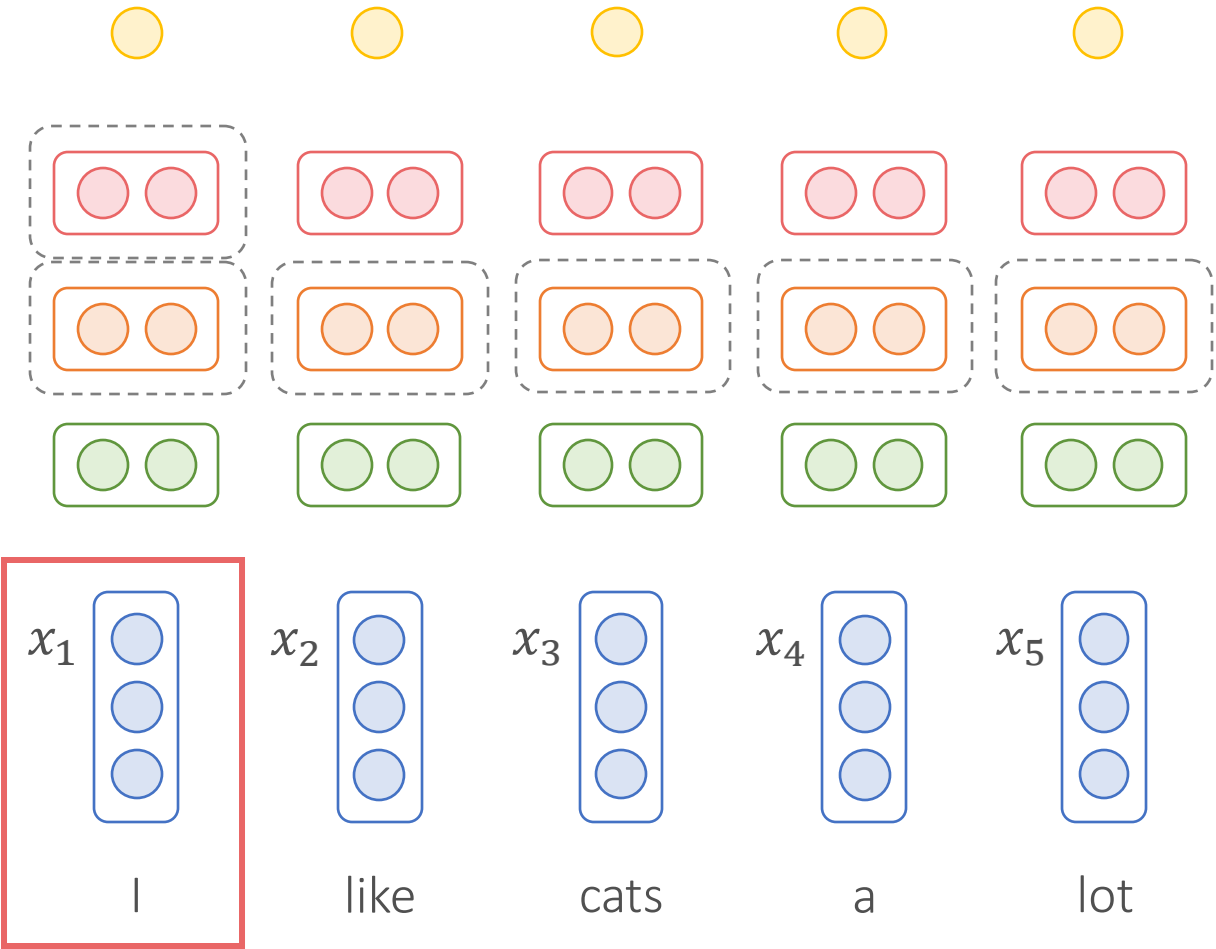
Vector dimension

Normalized
Attention Scores

Query $q_i = W^Q x_i$

Key $k_i = W^K x_i$

Value $v_i = W^V x_i$



Self-Attention

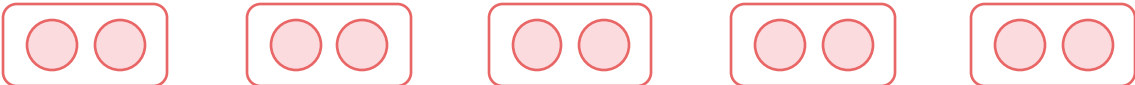
Weighted Sum

$$z_1 = \sum_i \alpha_{1,i} v_i$$

Normalized
Attention Scores



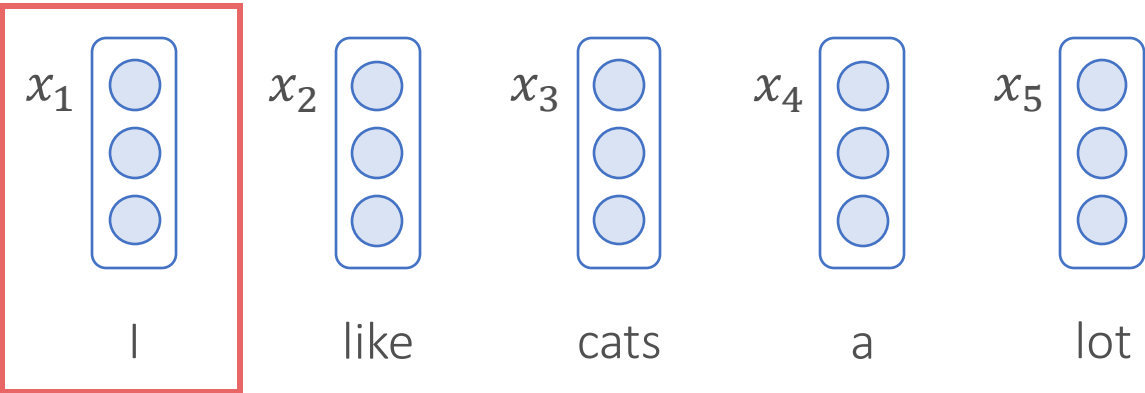
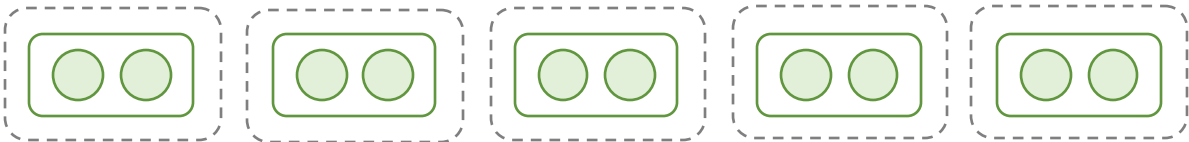
Query $q_i = W^Q x_i$



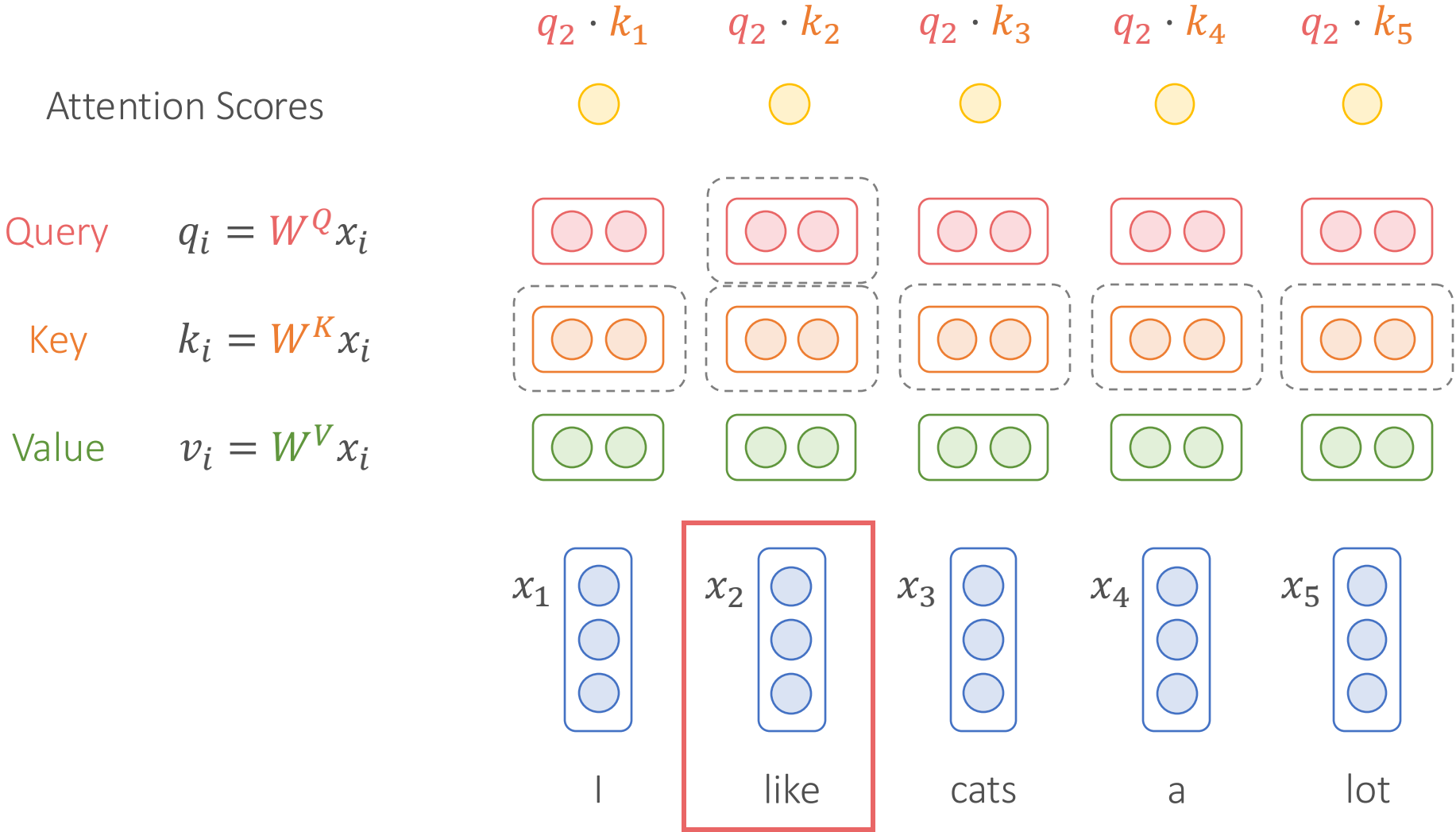
Key $k_i = W^K x_i$



Value $v_i = W^V x_i$



Self-Attention



Self-Attention

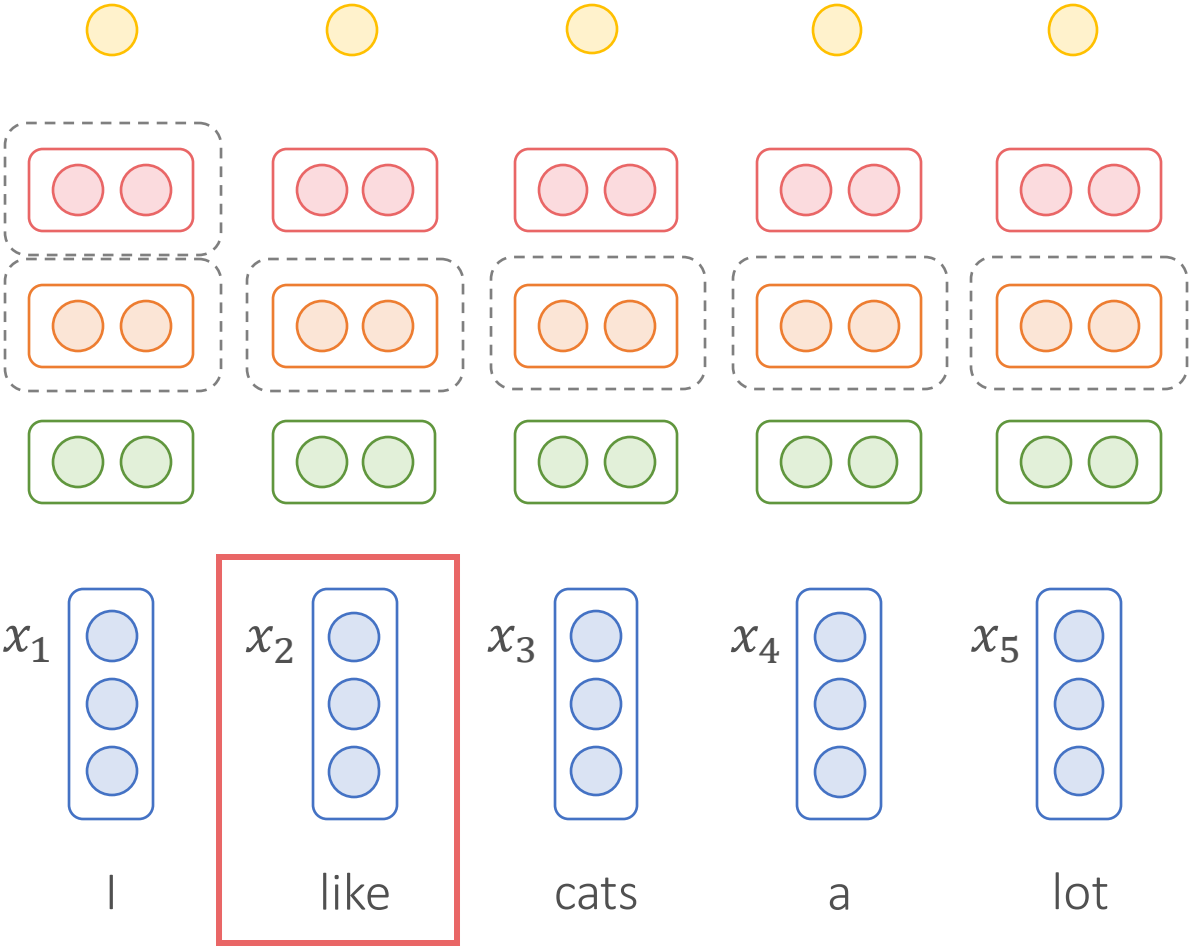
$$\alpha_{2,i} = \text{softmax}\left(\frac{q_2 \cdot k_i}{\sqrt{d}}\right) \quad \text{Vector dimension}$$

Normalized
Attention Scores

Query $q_i = W^Q x_i$

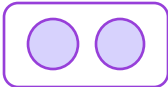
Key $k_i = W^K x_i$

Value $v_i = W^V x_i$



Self-Attention

Weighted Sum

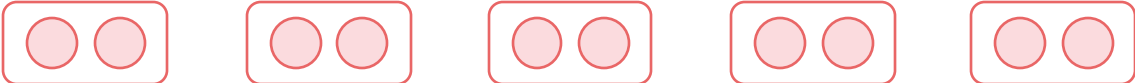


$$z_2 = \sum_i \alpha_{2,i} v_i$$

Normalized Attention Scores



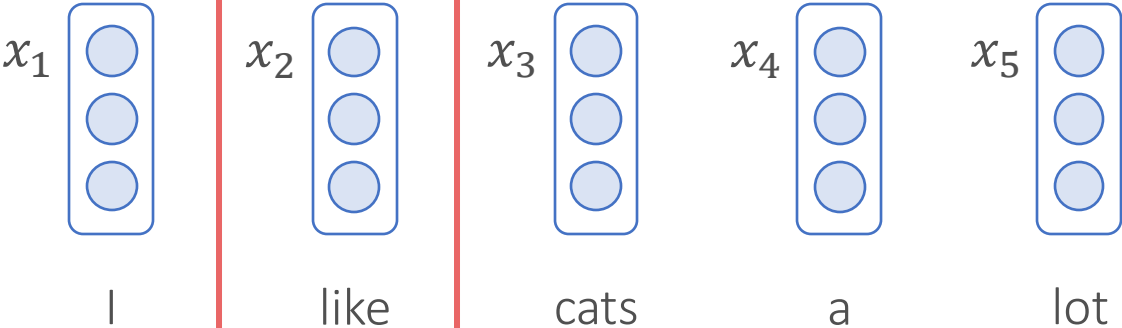
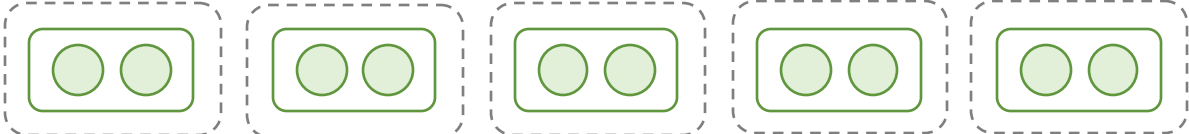
Query $q_i = W^Q x_i$



Key $k_i = W^K x_i$



Value $v_i = W^V x_i$



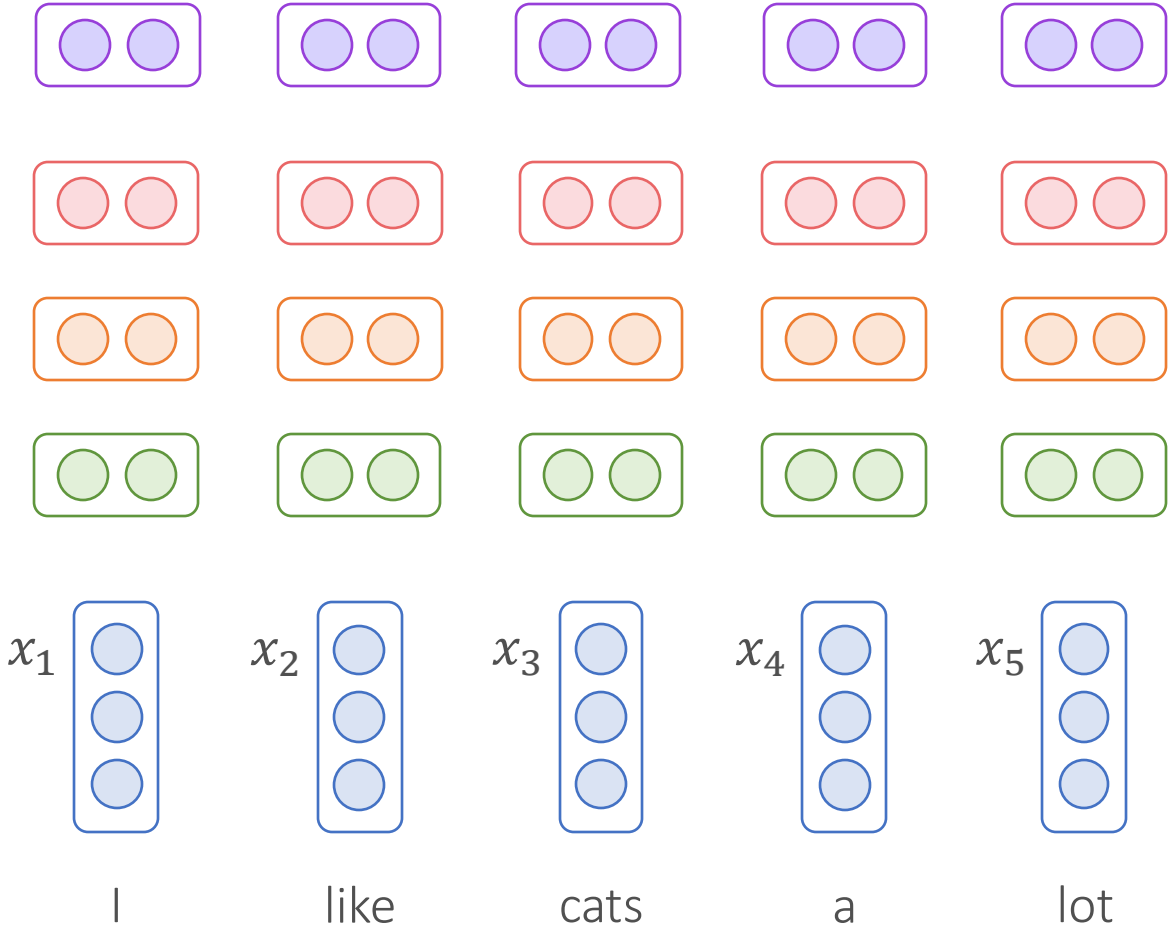
Self-Attention

Self-Attention Output

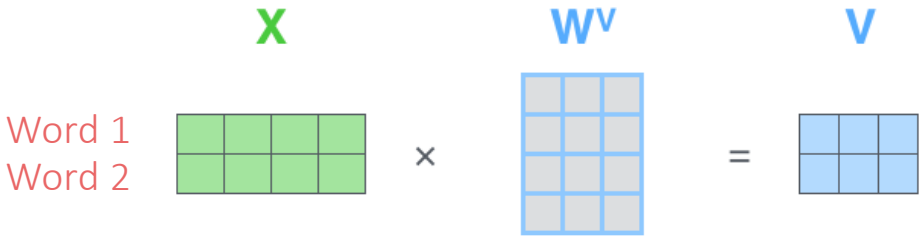
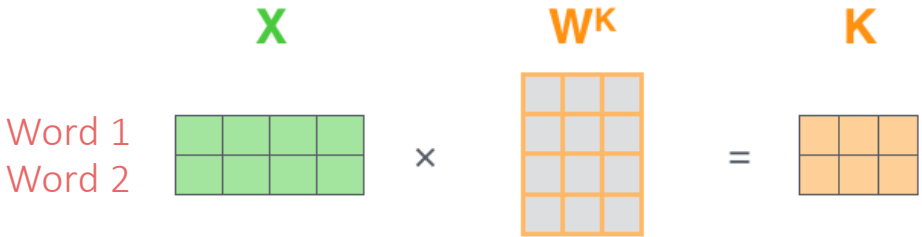
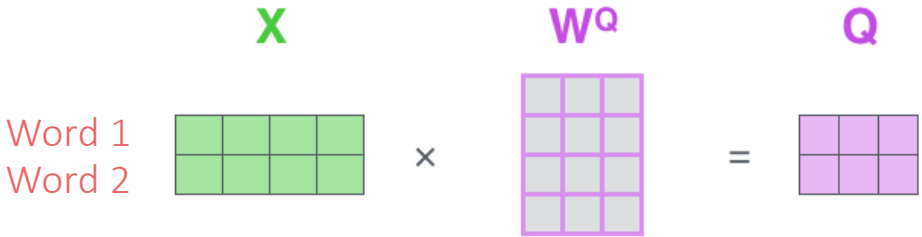
Query $q_i = W^Q x_i$

Key $k_i = W^K x_i$

Value $v_i = W^V x_i$



Self-Attention – Matrix Form



$$\text{softmax} \left(\frac{\begin{matrix} \text{Q} \\ \times \\ \text{K}^T \end{matrix}}{\sqrt{d_k}} \right) \begin{matrix} \text{V} \end{matrix}$$

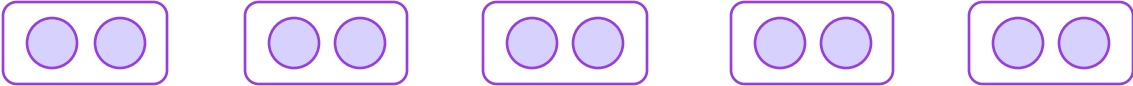
$=$

Z

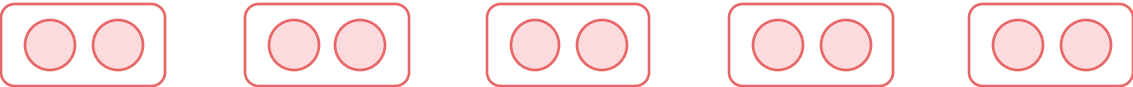
$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Single-Head Attention

Self-Attention Output



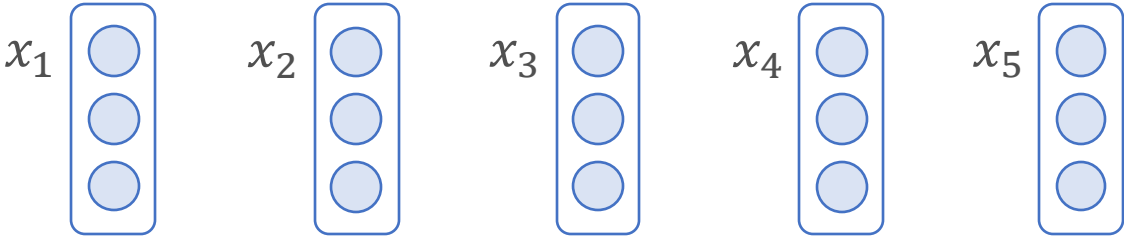
Query $q_i = W^Q x_i$



Key $k_i = W^K x_i$



Value $v_i = W^V x_i$



I like cats a lot

Multi-Head Attention

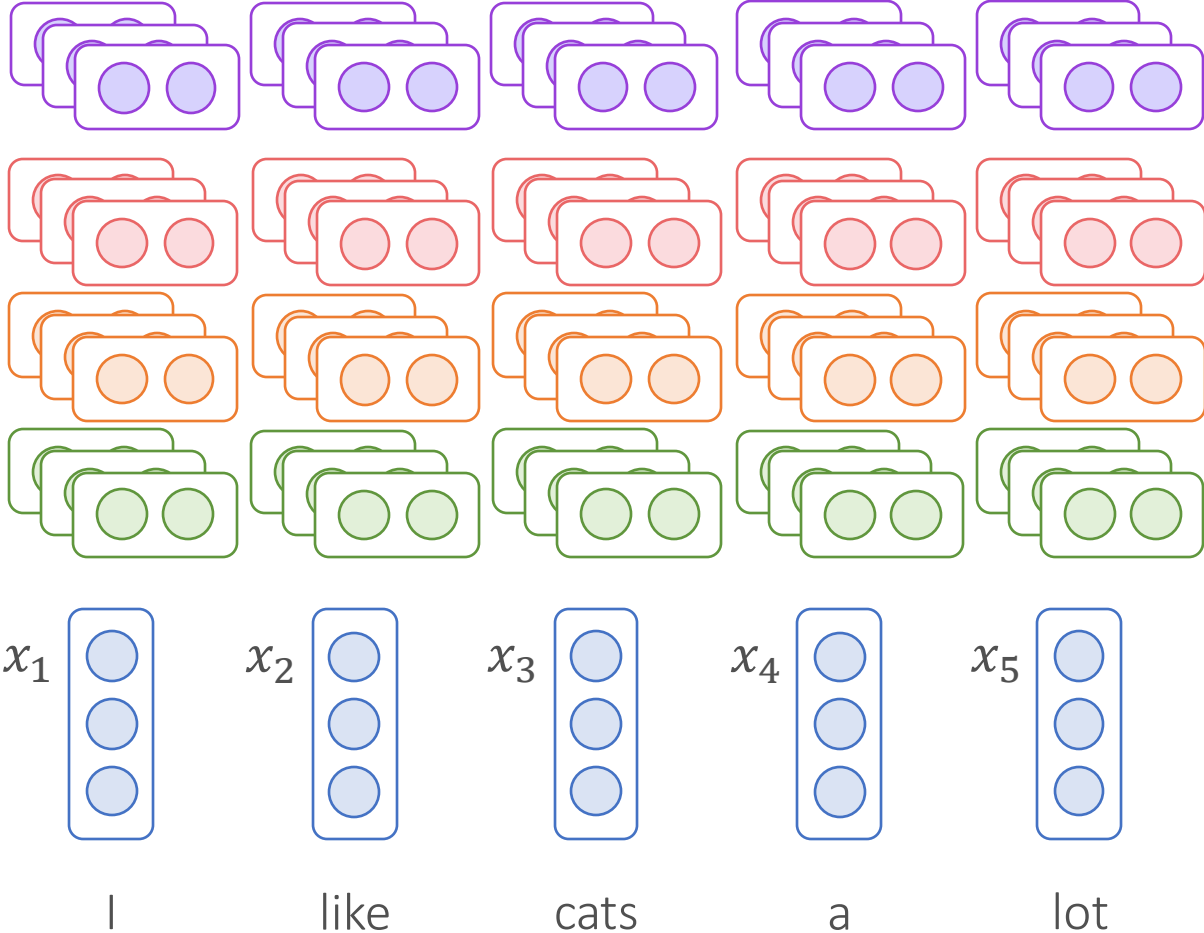
Each attention head focuses on different parts of understanding!

Multi-Attention Output

Query $q_i = W_j^Q x_i$

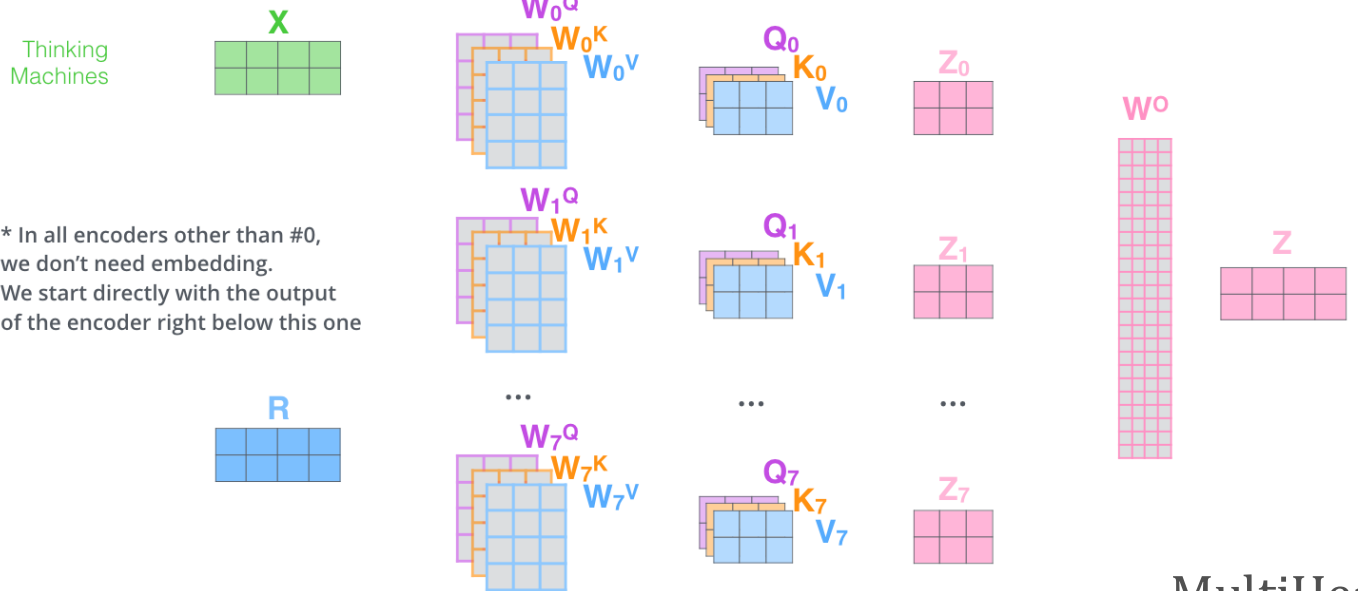
Key $k_i = W_j^K x_i$

Value $v_i = W_j^V x_i$



Multi-Head Attention – Matrix Form

- 1) This is our input sentence*
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting Q/K/V matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^o to produce the output of the layer

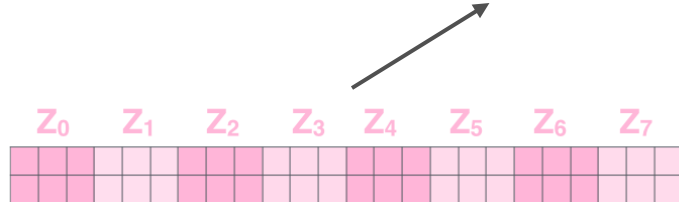


* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

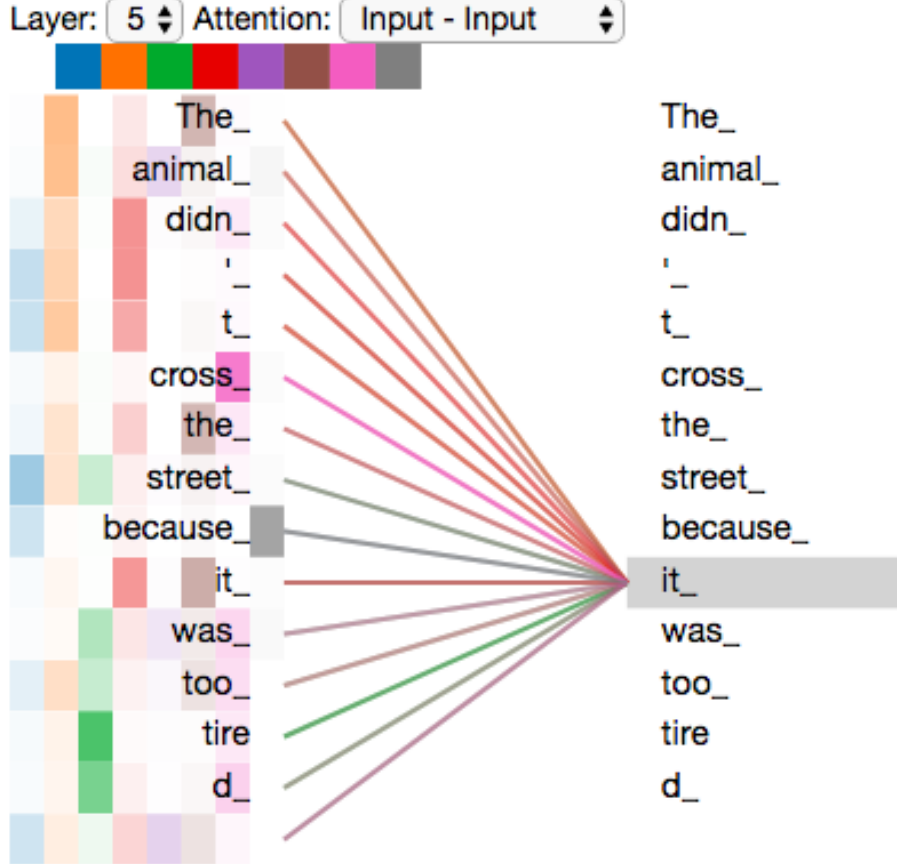
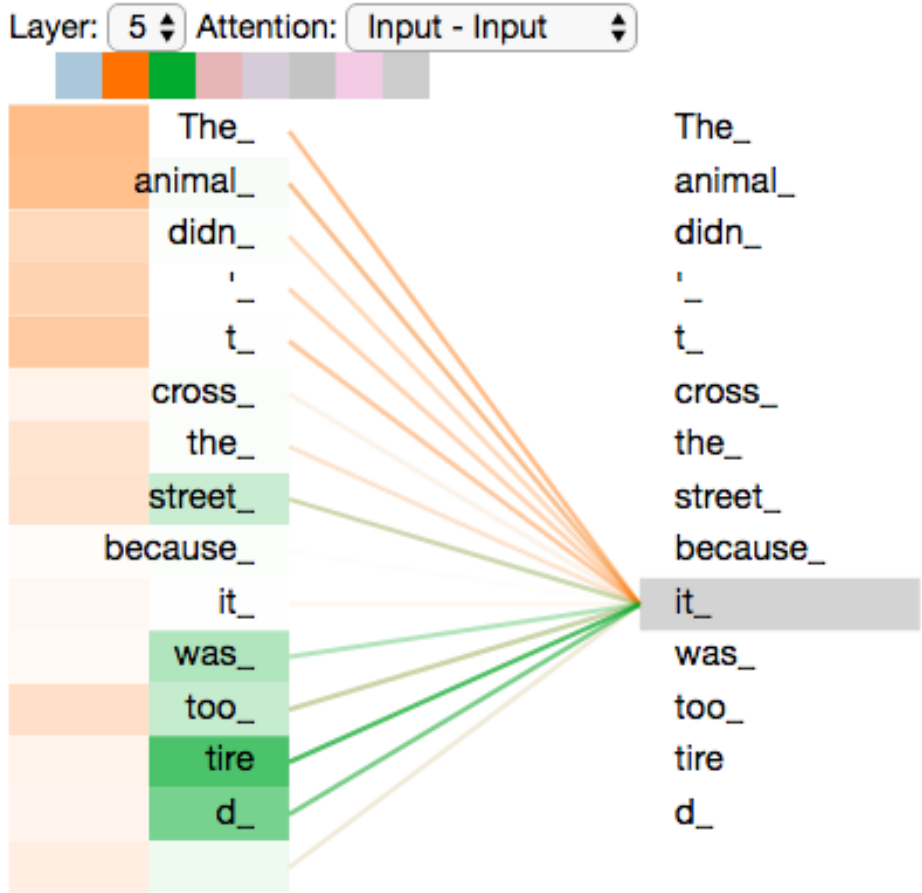
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

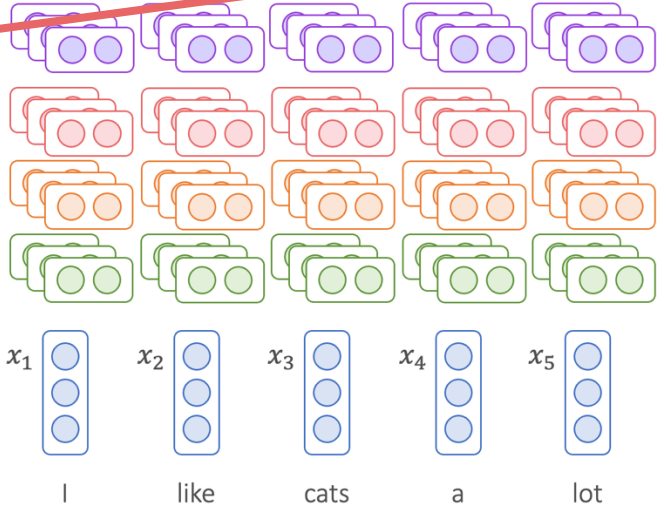
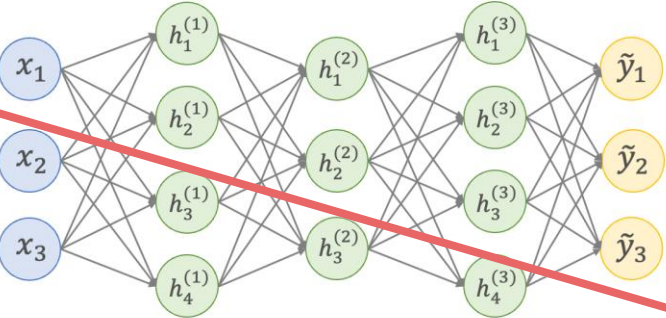
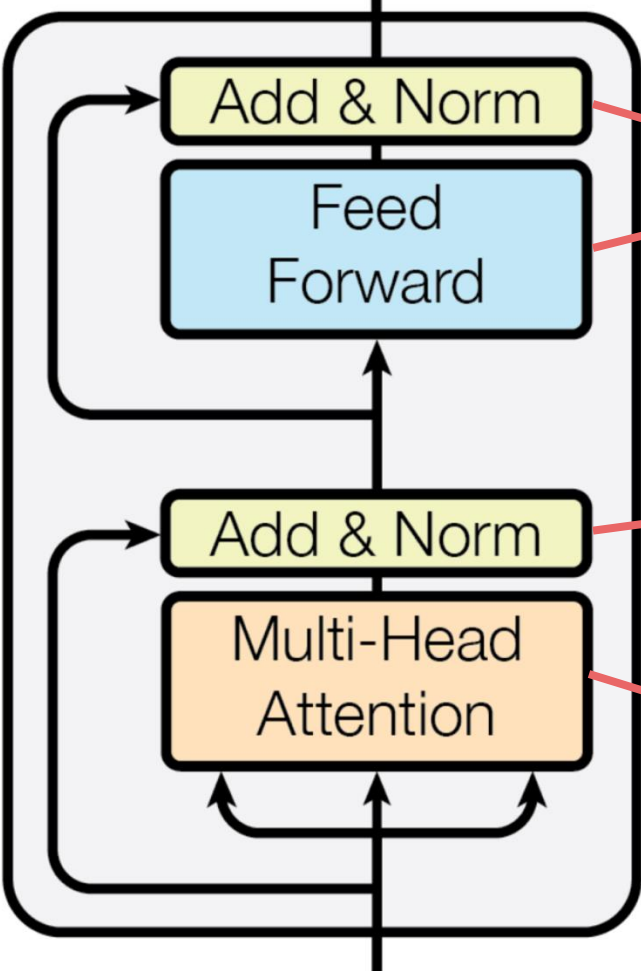
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^o$$



What Does Multi-Head Attention Learn?



Transformer Layer

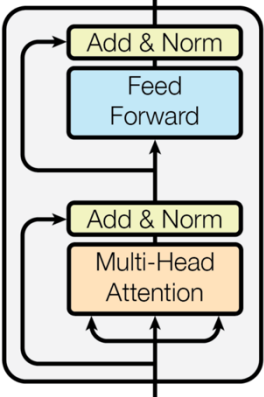
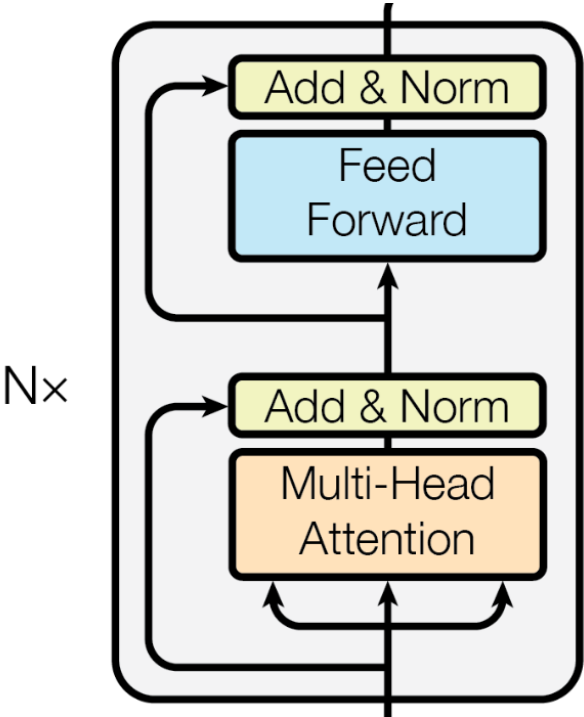


LayerNorm($x + \text{Sublayer}(x)$)

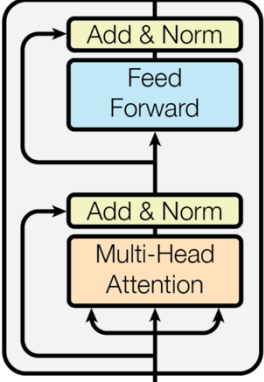
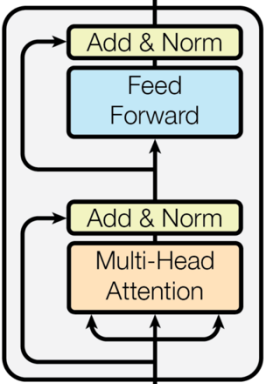
$$y = \frac{x - \mathbf{E}[x]}{\sqrt{\mathbf{Var}[x] + \epsilon}} * \gamma + \beta$$

Residual connection (He et al., 2016)
Layer normalization (Ba et al., 2016)

Transformer Encoder

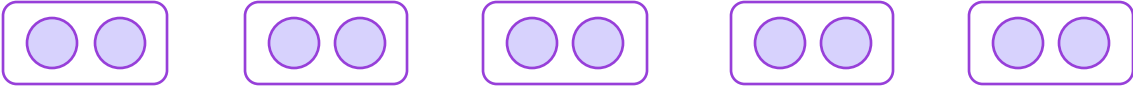


...

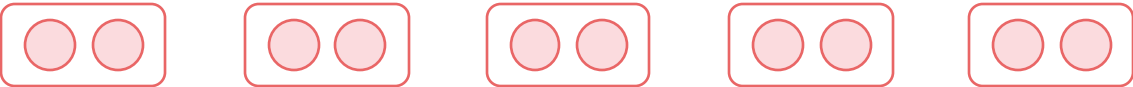


How About Word Order?

Self-Attention Output



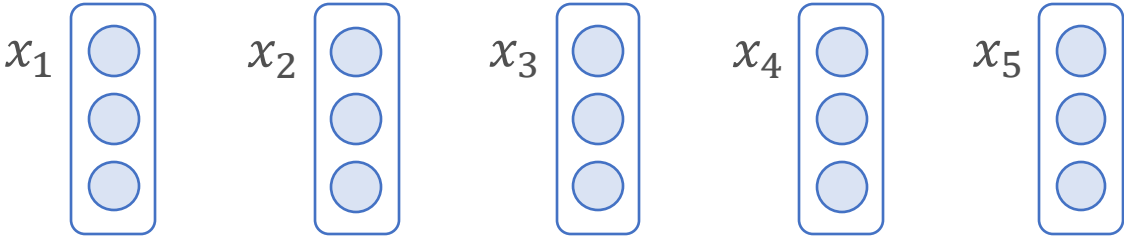
Query $q_i = W^Q x_i$



Key $k_i = W^K x_i$

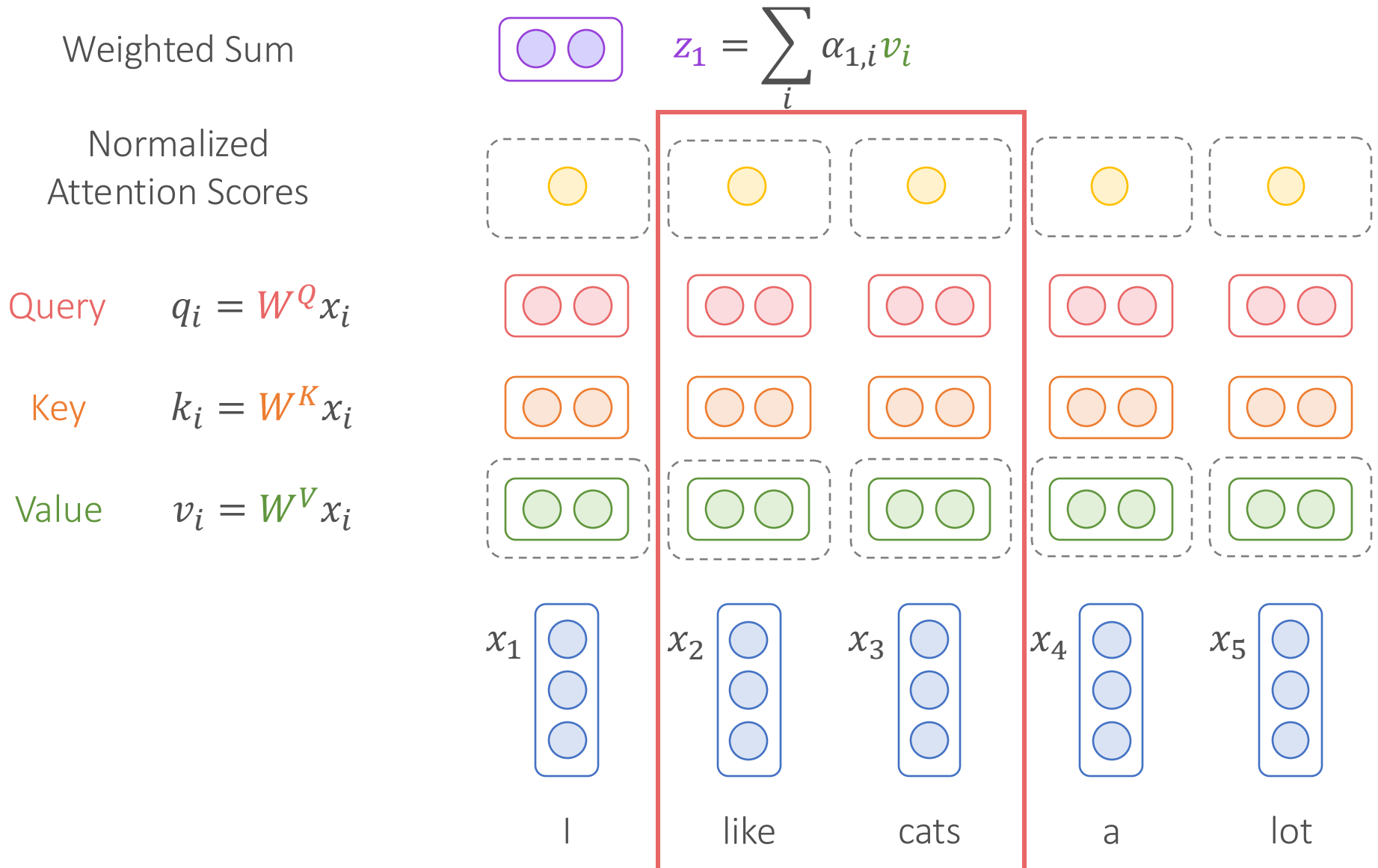


Value $v_i = W^V x_i$

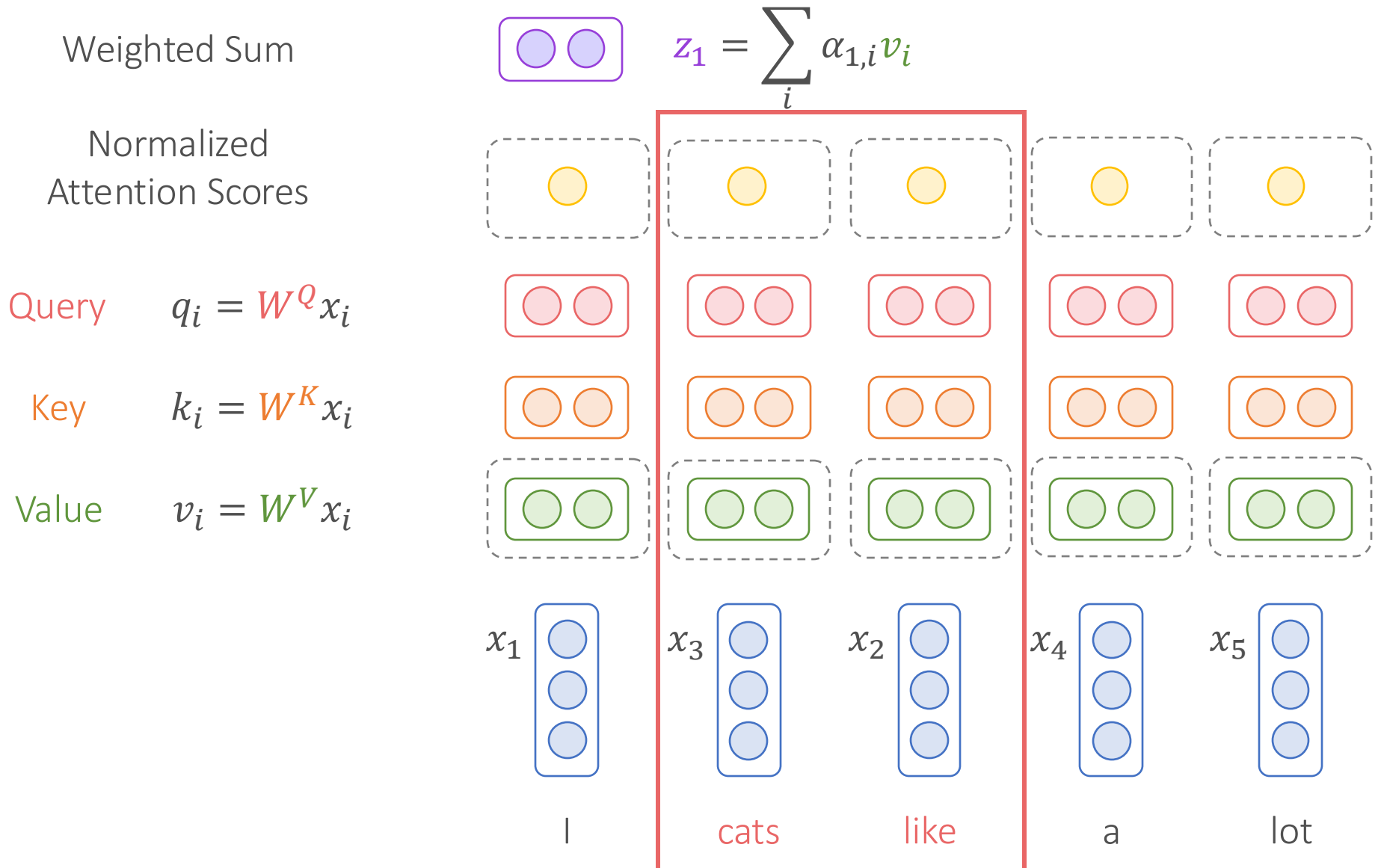


x_1 I x_2 like x_3 cats x_4 a x_5 lot

How About Word Order?

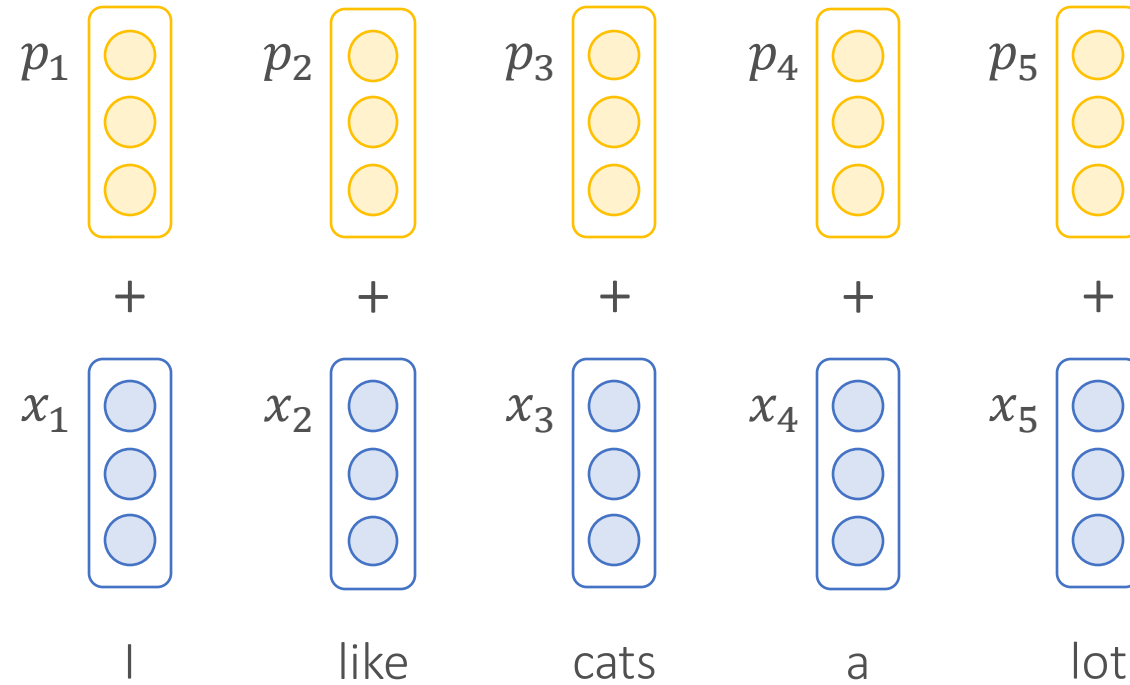


How About Word Order?



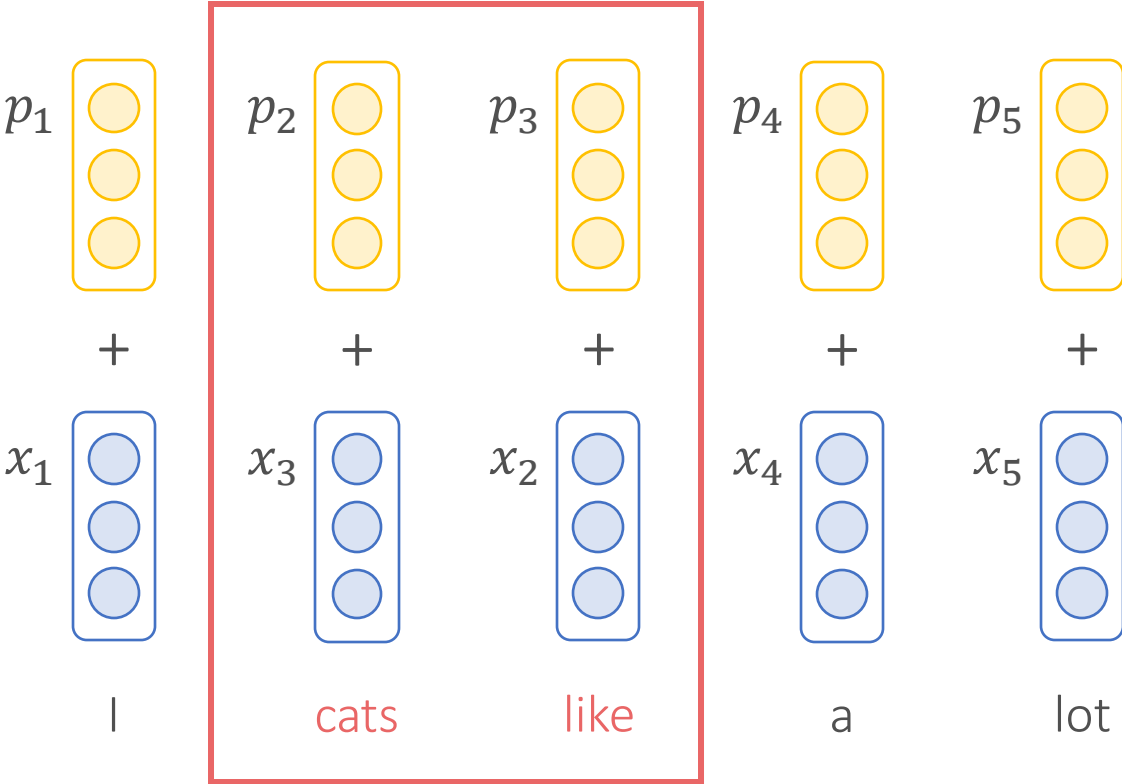
Solution: Positional Encoding

$$x_i \leftarrow x_i + PE_i$$



Solution: Positional Encoding

$$x_i \leftarrow x_i + PE_i$$



Solution: Positional Encoding

- Unique encoding for each position
- Closer positions should have more similar encodings
- Distance between neighboring positions should be the same

Sinusoidal Positional Encoding

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

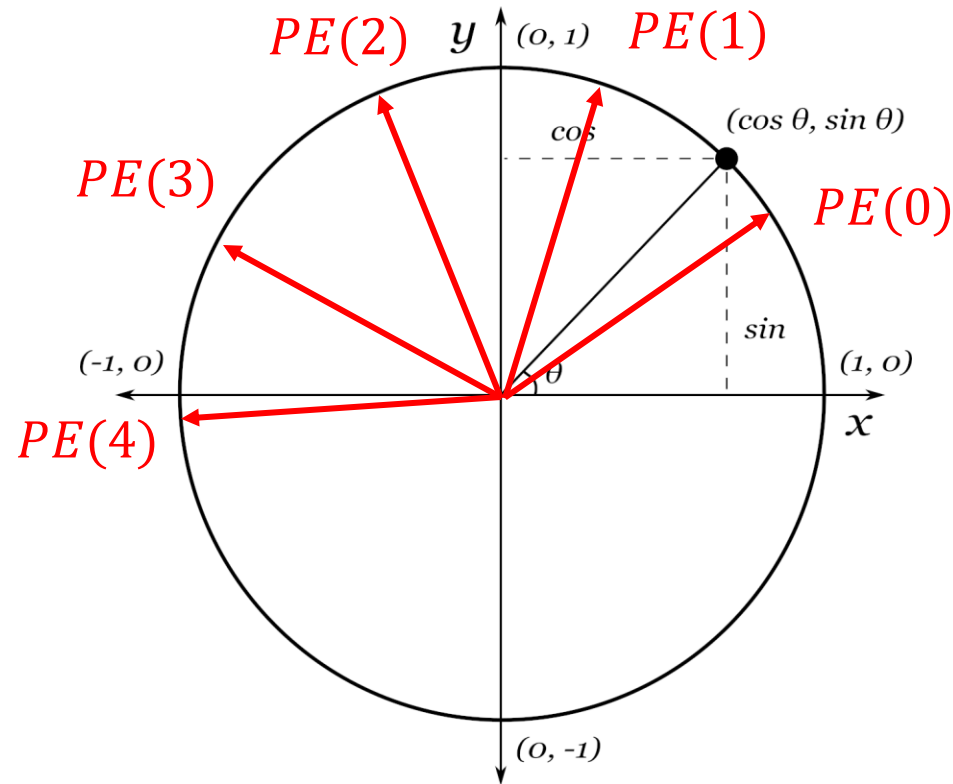
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

Why this?

Sinusoidal Positional Encoding: Intuition

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

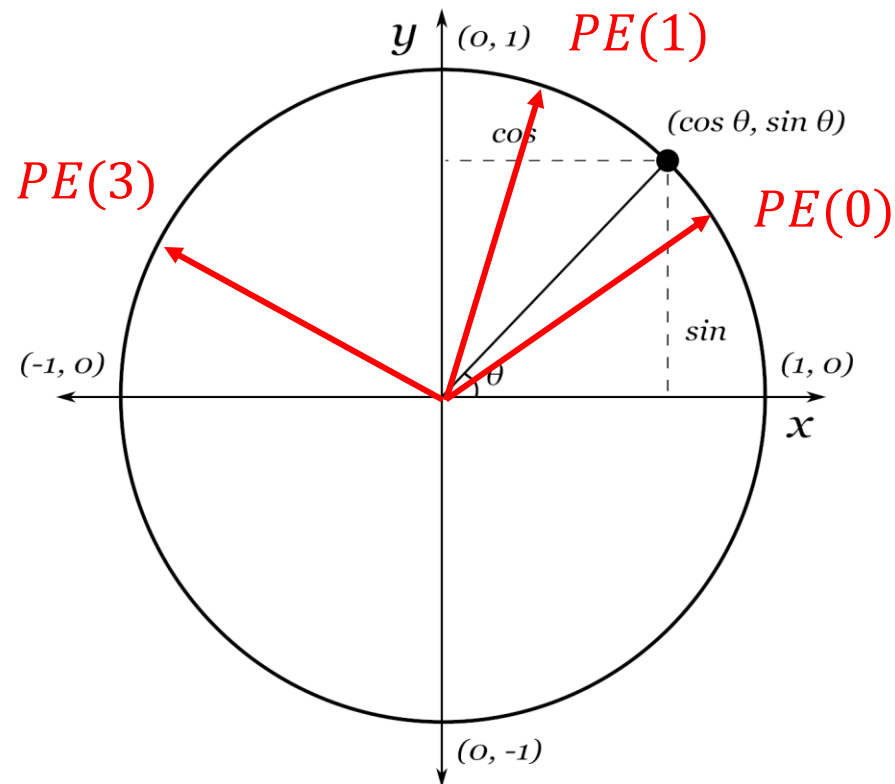
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$



Sinusoidal Positional Encoding: Intuition

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$



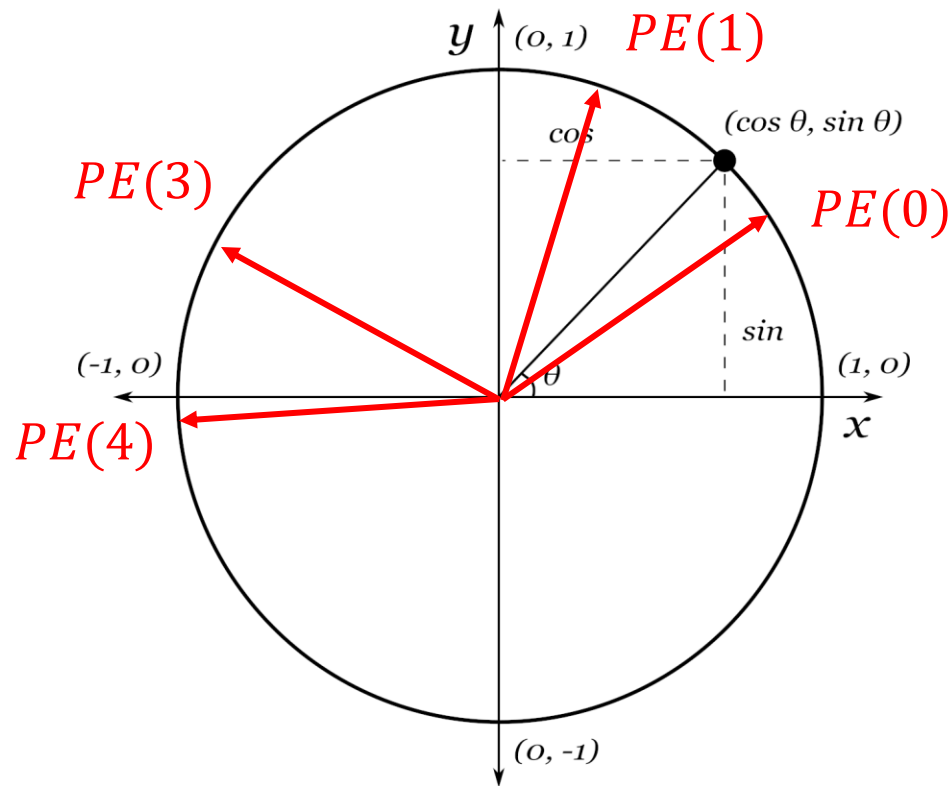
$$\text{Cosine}(PE(0), PE(1)) > \text{Cosine}(PE(0), PE(3))$$

Closer positions should have more similar encodings

Sinusoidal Positional Encoding: Intuition

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$



$$\text{Cosine}(PE(0), PE(1)) = \text{Cosine}(PE(3), PE(4))$$

Distance between neighboring positions should be the same

Sinusoidal Positional Encoding: Intuition

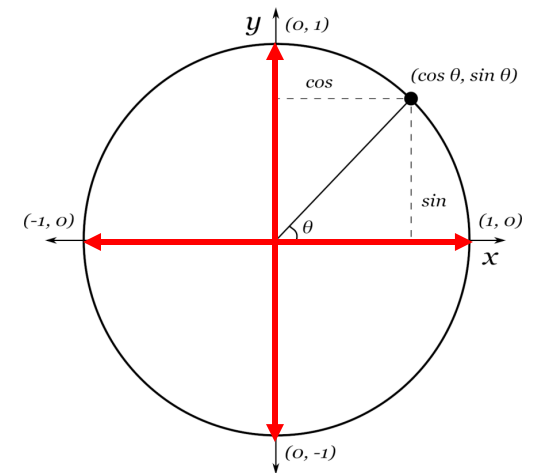
- How to expand to high-dimension?
- Let's consider **binary** positional encoding first
- How to use 4 bits to represent position 0~15?

0 :	0	0	0	0	8 :	1	0	0	0
1 :	0	0	0	1	9 :	1	0	0	1
2 :	0	0	1	0	10 :	1	0	1	0
3 :	0	0	1	1	11 :	1	0	1	1
4 :	0	1	0	0	12 :	1	1	0	0
5 :	0	1	0	1	13 :	1	1	0	1
6 :	0	1	1	0	14 :	1	1	1	0
7 :	0	1	1	1	15 :	1	1	1	1

Sinusoidal Positional Encoding: Intuition

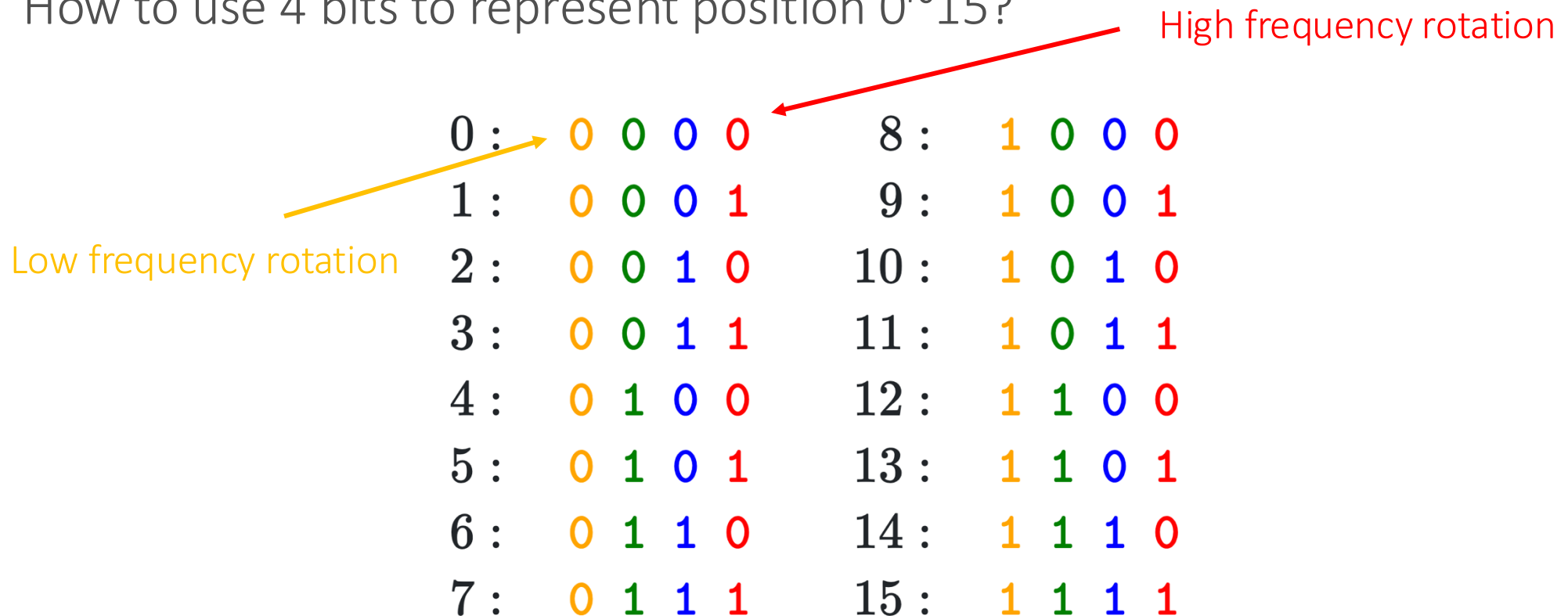
- How to expand to high-dimension?
- Let's consider **binary** positional encoding first
- How to use 4 bits to represent position 0~15?

0 :	0	0	0	0	8 :	1	0	0	0
1 :	0	0	0	1	9 :	1	0	0	1
2 :	0	0	1	0	10 :	1	0	1	0
3 :	0	0	1	1	11 :	1	0	1	1
4 :	0	1	0	0	12 :	1	1	0	0
5 :	0	1	0	1	13 :	1	1	0	1
6 :	0	1	1	0	14 :	1	1	1	0
7 :	0	1	1	1	15 :	1	1	1	1



Sinusoidal Positional Encoding: Intuition

- How to expand to high-dimension?
- Let's consider **binary** positional encoding first
- How to use 4 bits to represent position 0~15?



Sinusoidal Positional Encoding: Intuition

- How to expand to high-dimension?
- Let's consider **binary** positional encoding first
- How to use 4 bits to represent position 0~15?

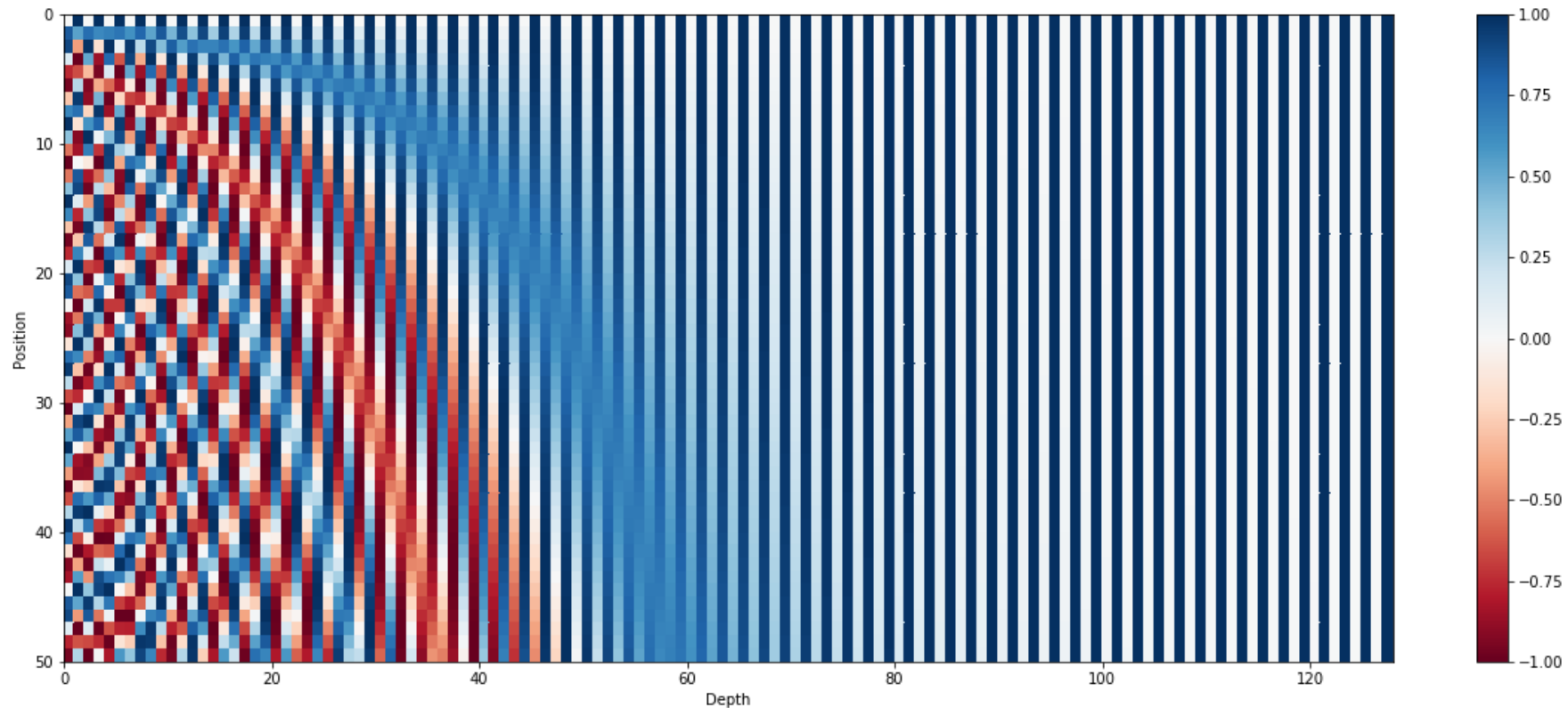
0 :	0	0	0	0	8 :	1	0	0	0
1 :	0	0	0	1	9 :	1	0	0	1
2 :	0	0	1	0	10 :	1	0	1	0
3 :	0	0	1	1	11 :	1	0	1	1
4 :	0	1	0	0	12 :	1	1	0	0
5 :	0	1	0	1	13 :	1	1	0	1
6 :	0	1	1	0	14 :	1	1	1	0
7 :	0	1	1	1	15 :	1	1	1	1

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

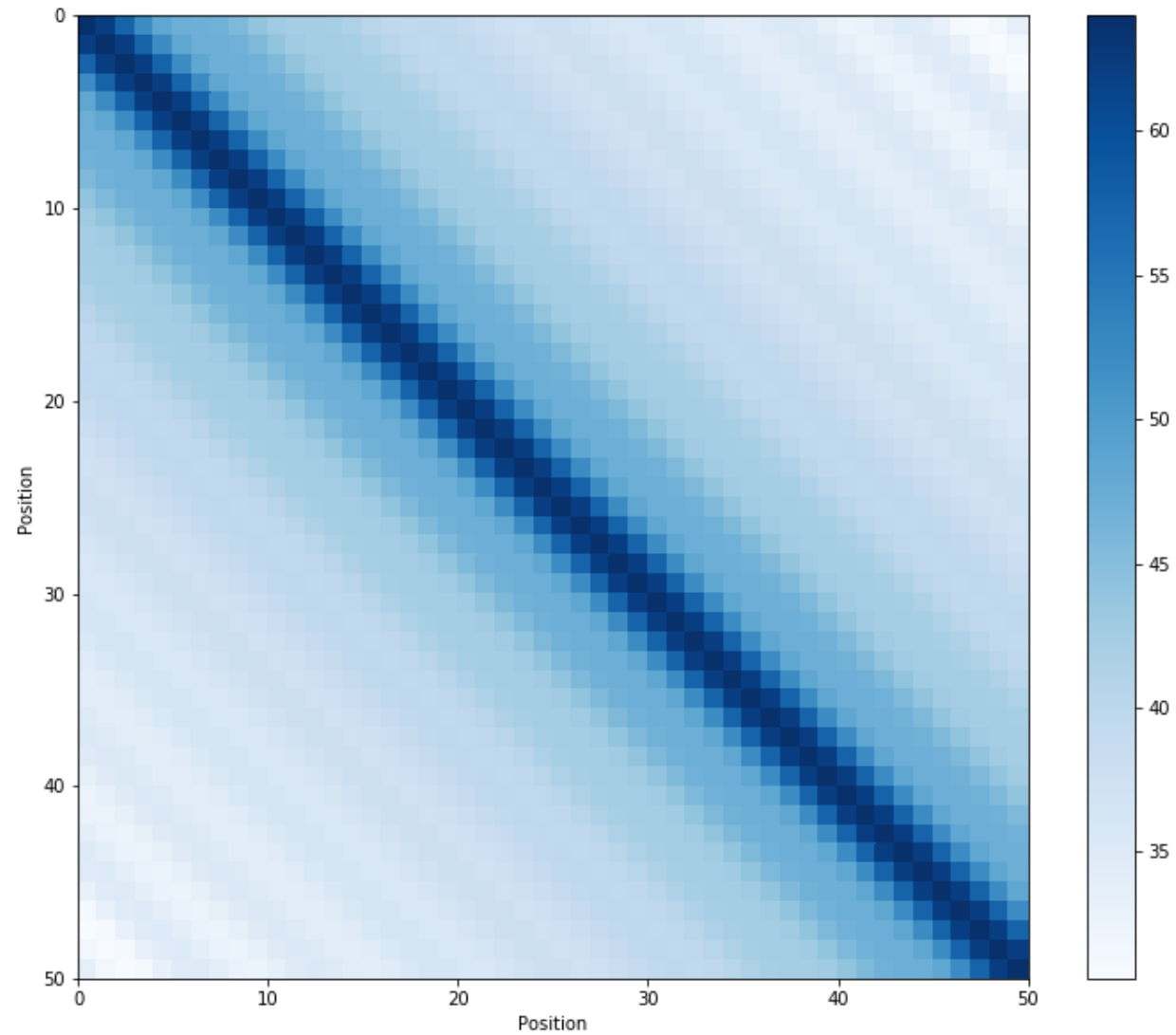
Soft version of alternating bits

Sinusoidal Positional Encoding

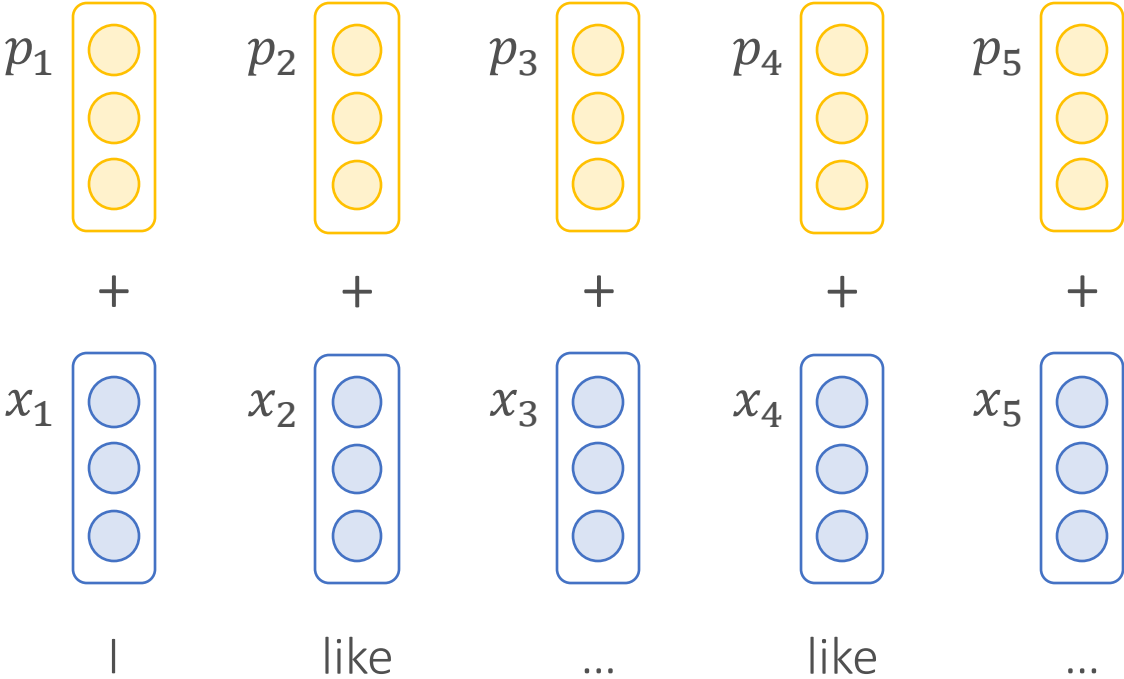
$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$



Sinusoidal Positional Encoding



Positional Encoding



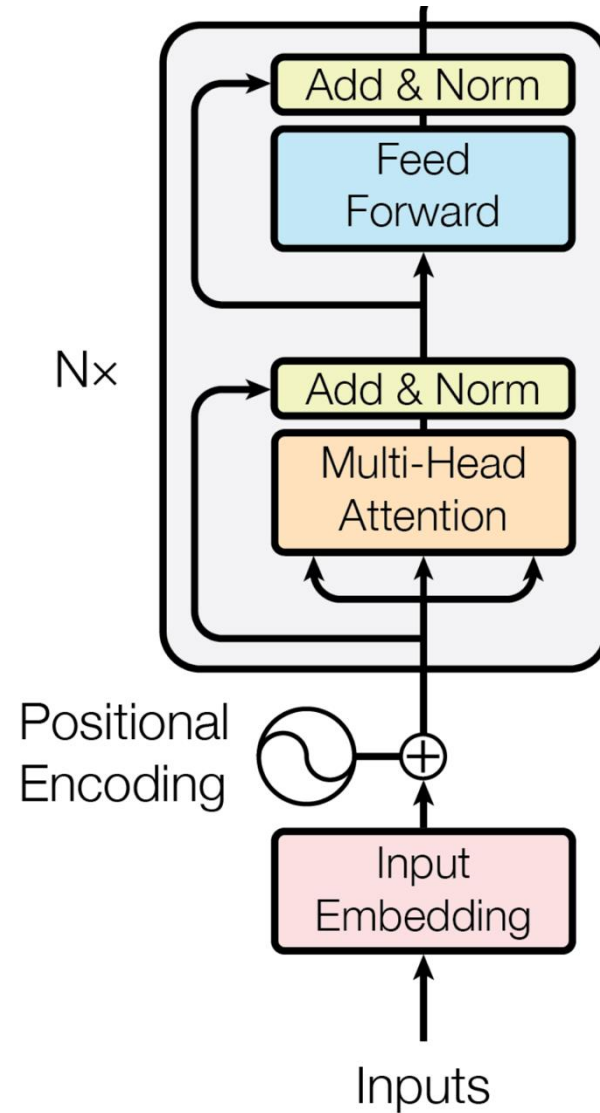
$$(E(I) + PE(1)) (E(\text{like}) + PE(2)) = E(I)E(\text{like}) + E(I)PE(2) + PE(1)E(\text{like}) + PE(1)PE(2)$$

$$(E(I) + PE(1)) (E(\text{like}) + PE(4)) = E(I)E(\text{like}) + E(I)PE(4) + PE(1)E(\text{like}) + PE(1)PE(4)$$

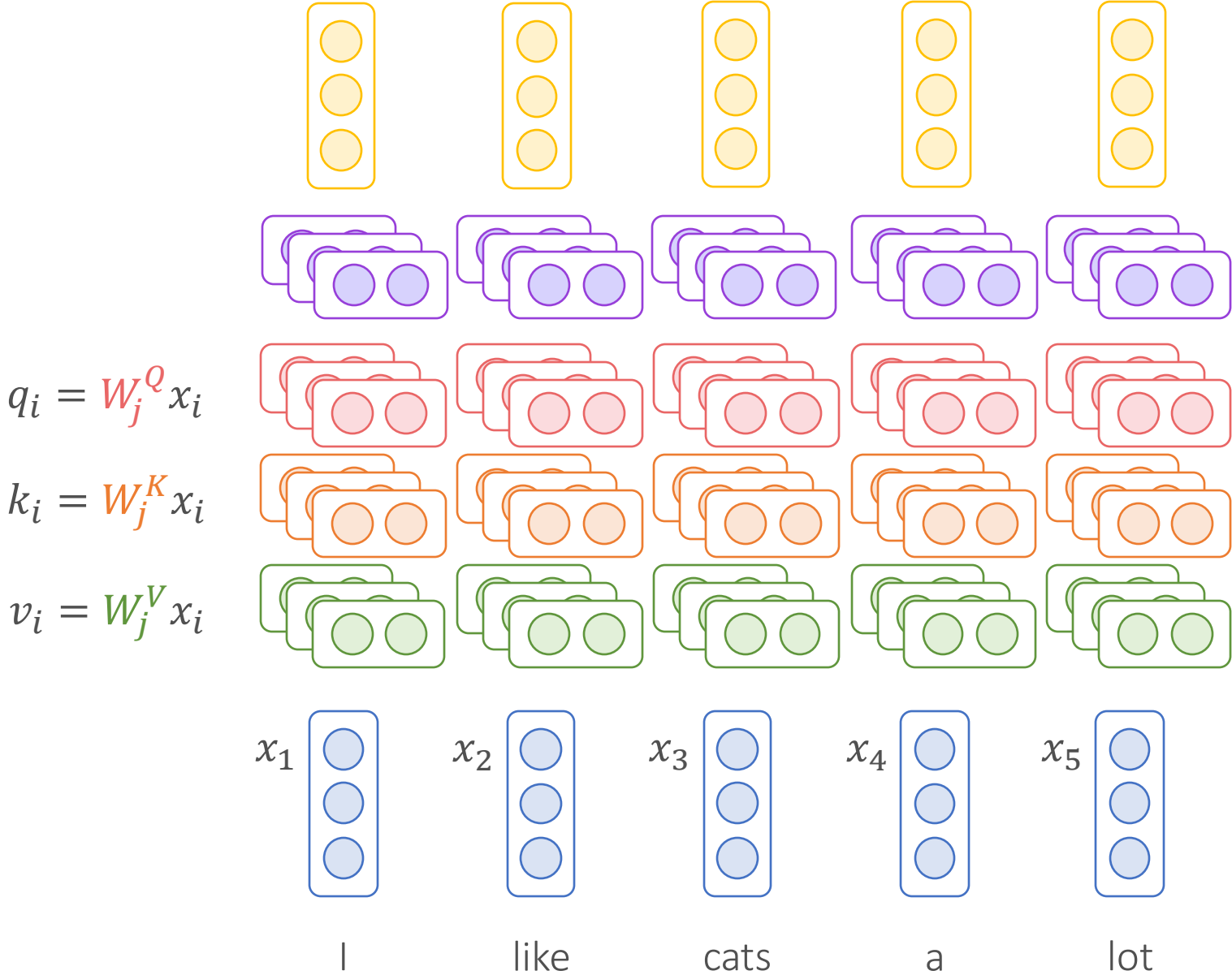
In expectation, they are the same

Position difference

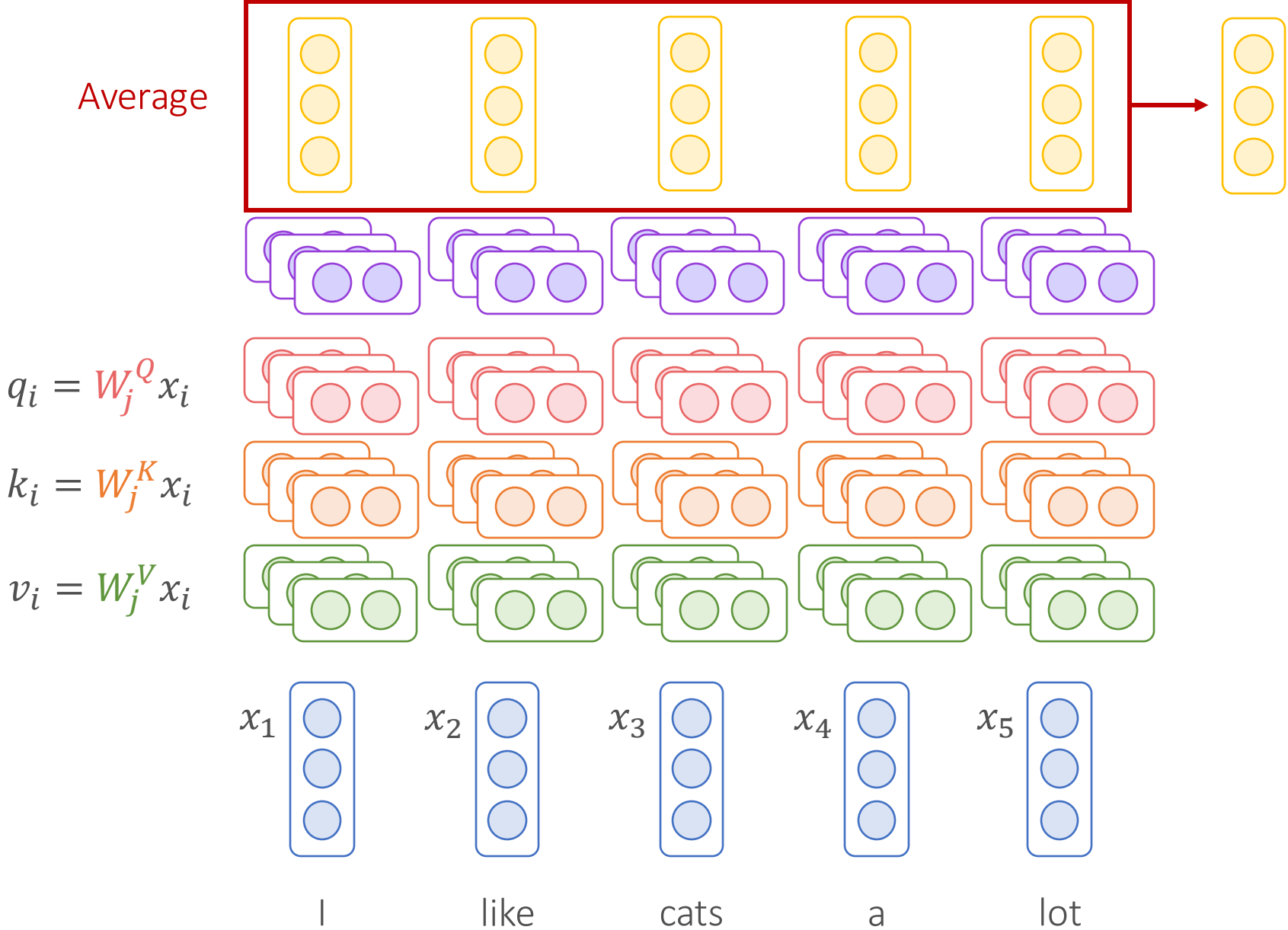
Transformer Encoder with Positional Encoding



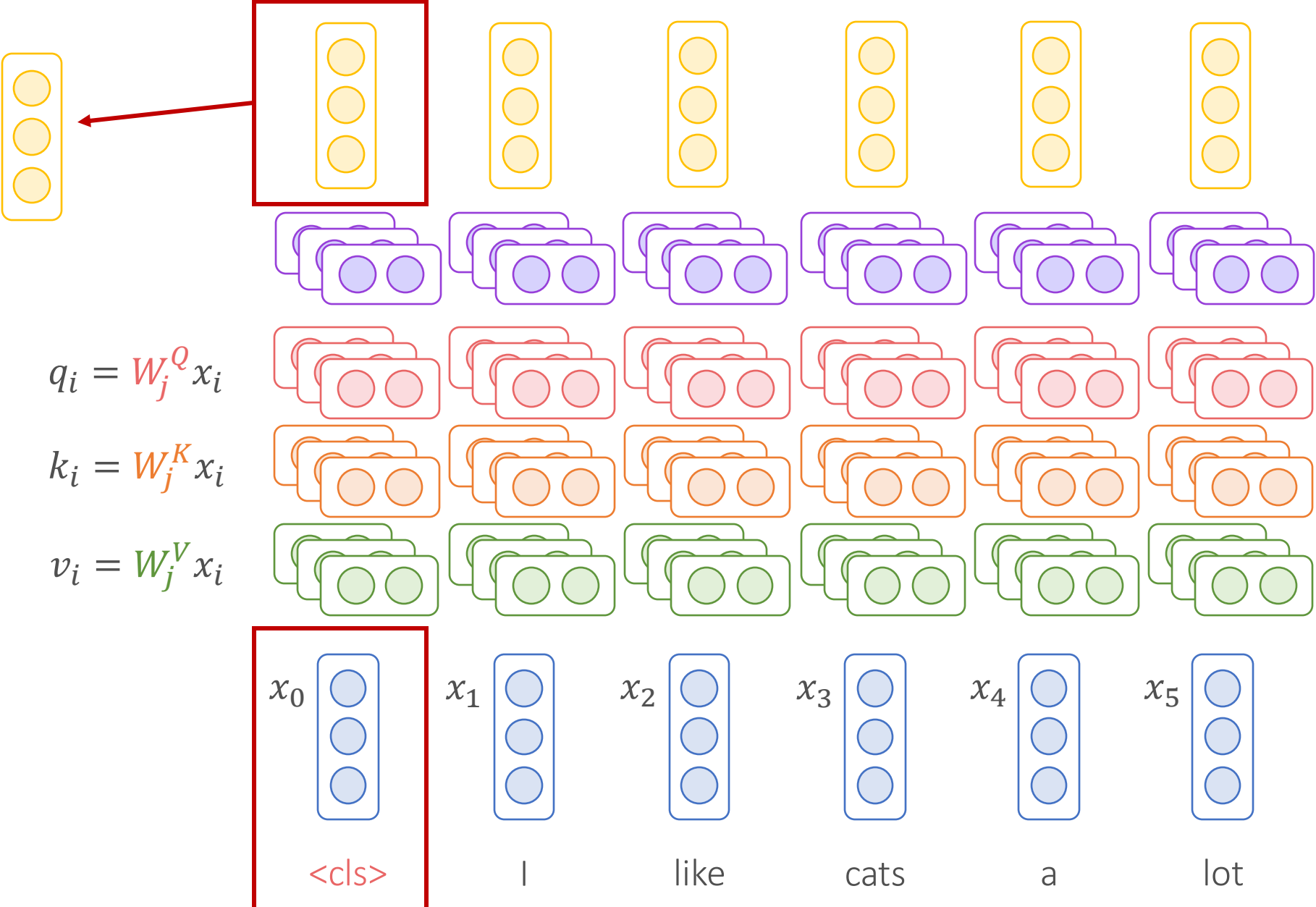
Transformer as Token-Level Encoder



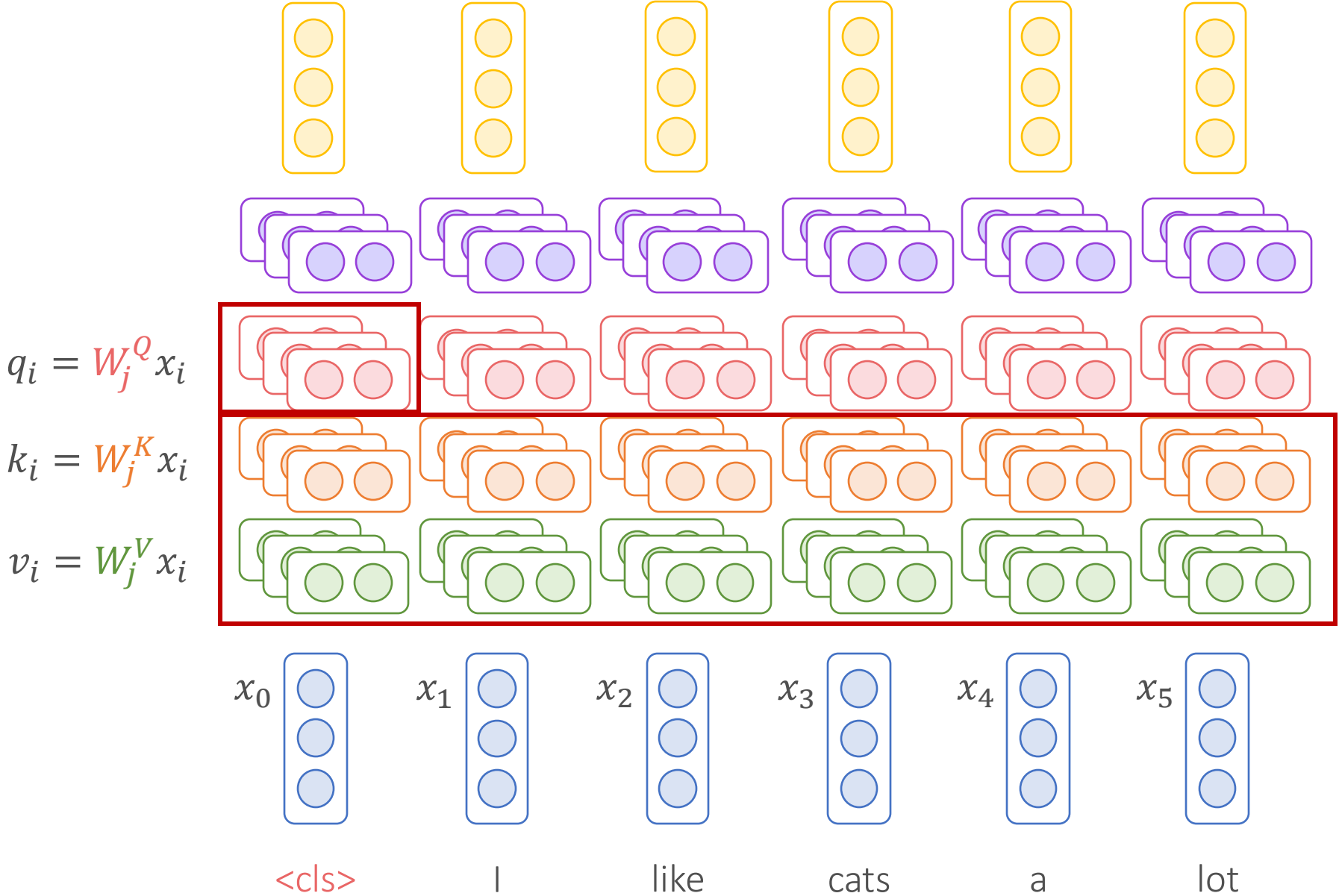
Transformer as Sentence-Level Encoder



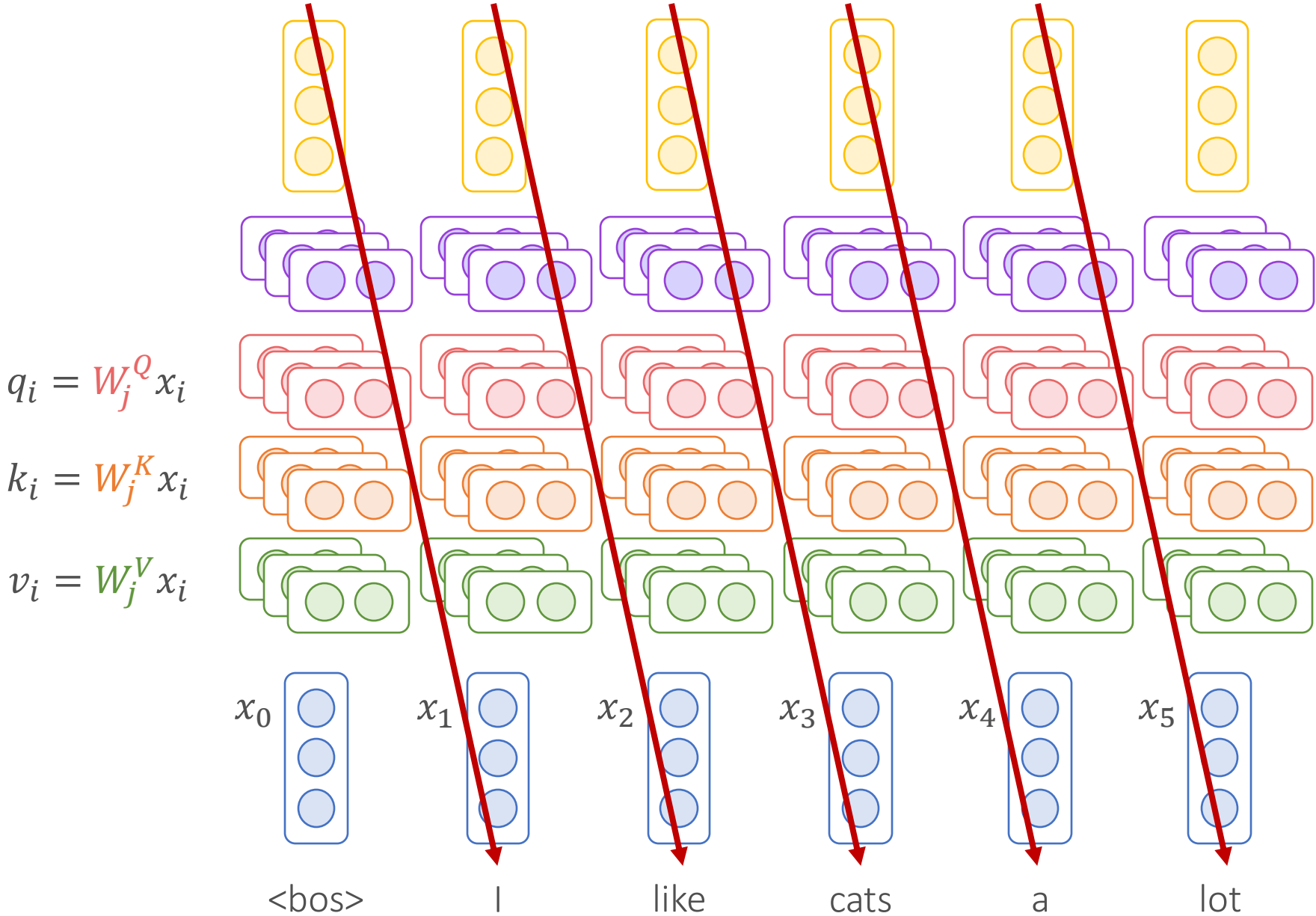
Transformer as Sentence-Level Encoder



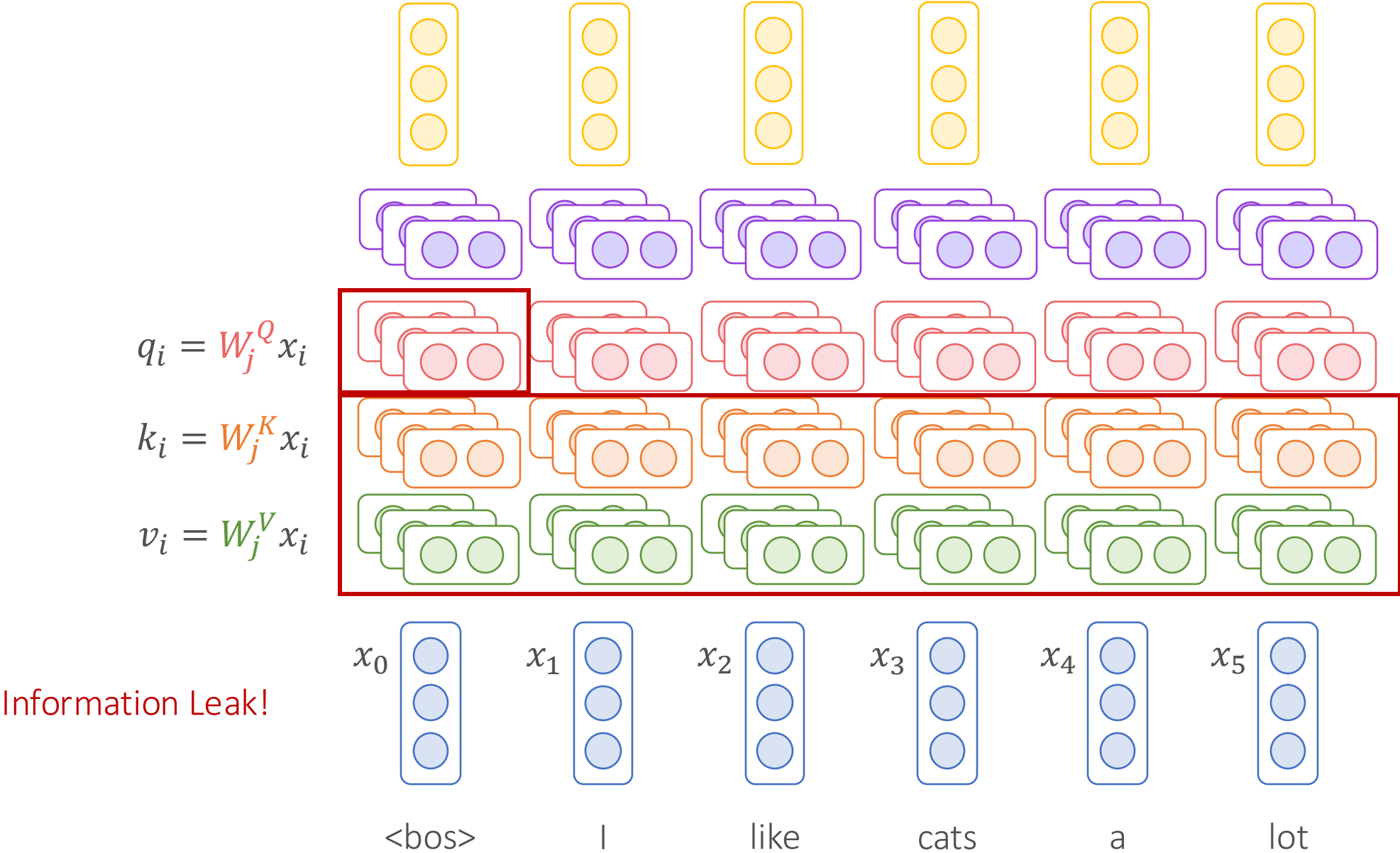
Transformer as Sentence-Level Encoder



Transformer as Decoder?



Transformer as Decoder?



Information Leak!

Lecture Plan

- Transformers
 - Attention
 - Self-Attention
 - Transformer Encoder
 - Positional Encoding