

# Natural Language Processing

## Self Attention and Transformers

Kabir Ahuja

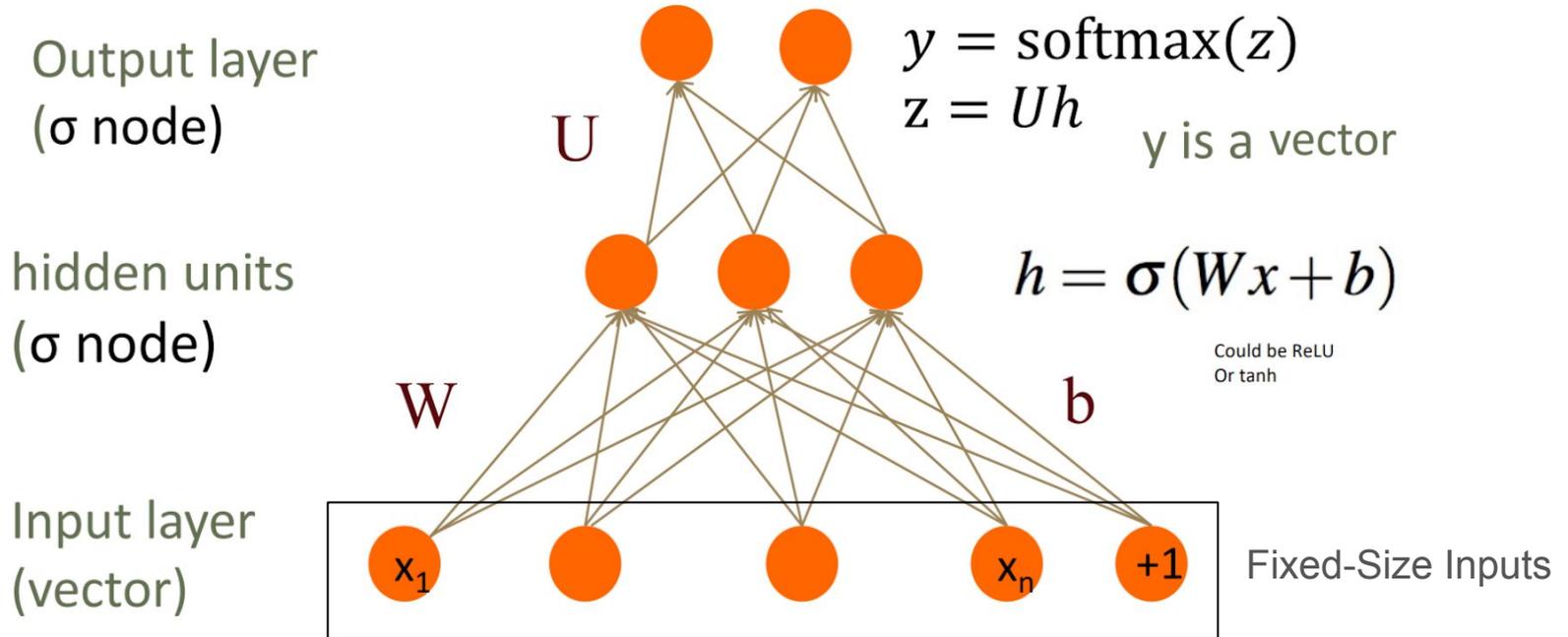
[kahuja@cs.washington.edu](mailto:kahuja@cs.washington.edu)

Adapted from slides from by Vidhisha Balachandran, Emma Strubell, Graham Neubig

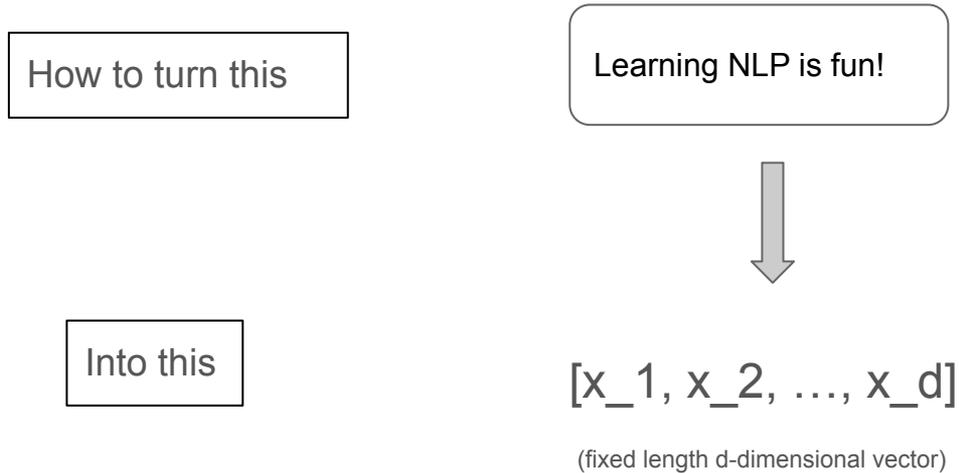
# Readings

- [Attention Is All You Need](#)
- [The Illustrated Transformer](#)
- [The Annotated Transformer](#)
- [Language Modeling with Transformers and PyTorch](#)
- [RoFormer: Enhanced Transformer with Rotary Position Embedding](#)
- [How Rotary Position Embedding Supercharges Modern LLMs \[RoPE\]](#)

# Recap - 2 Layer MLP

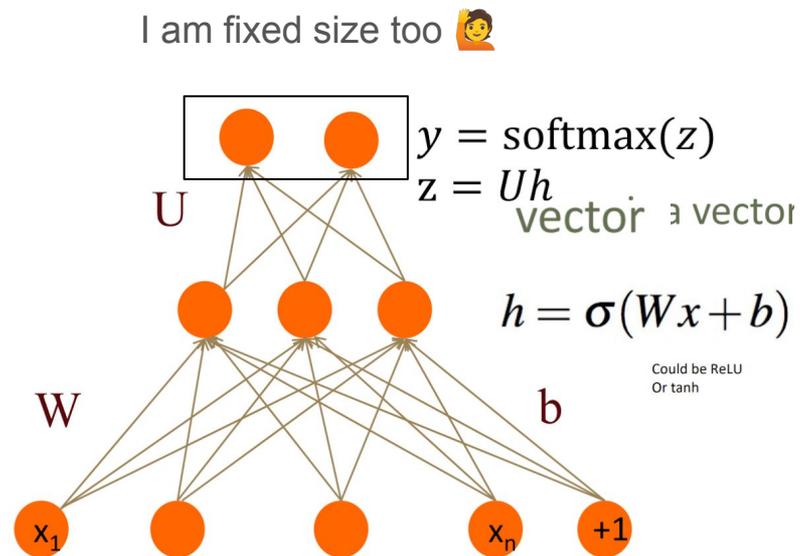


# Fixed size representations for Natural Language



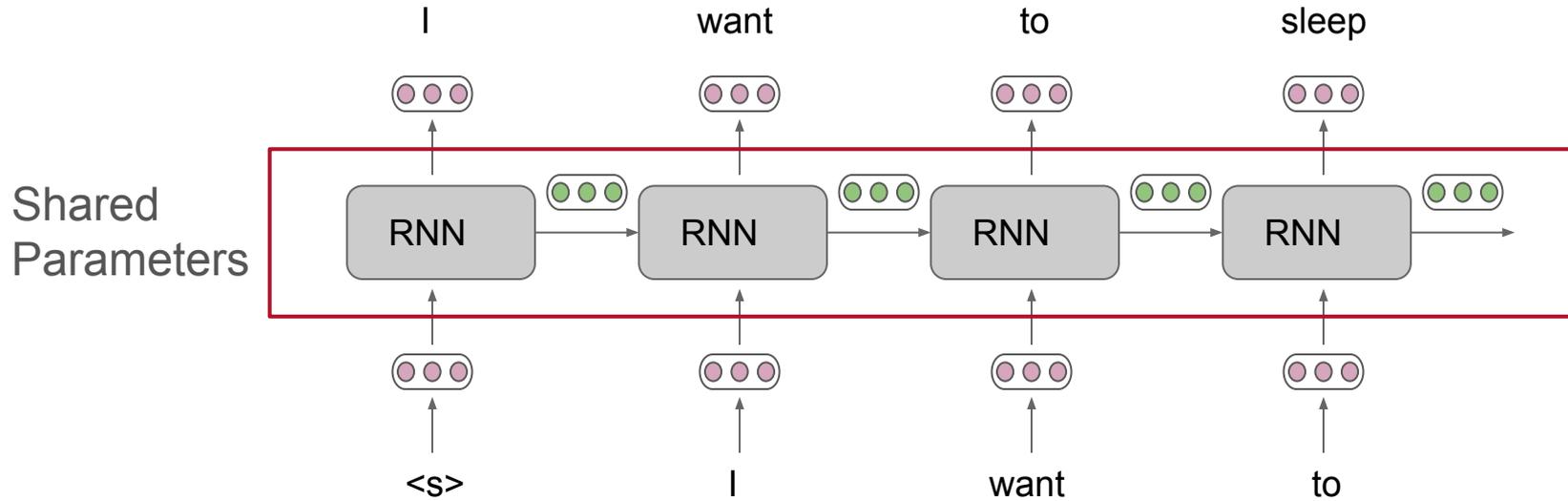
# Problems with Fixed-Size Representations

1. Loss of important contextual information
2. Order invariant
3. Often Poor Generalization (Especially in case of Linguistic Features)
4. What if the output is a sequence too? Think about tasks like machine translation, good luck getting the sentence back from the fixed representation :)

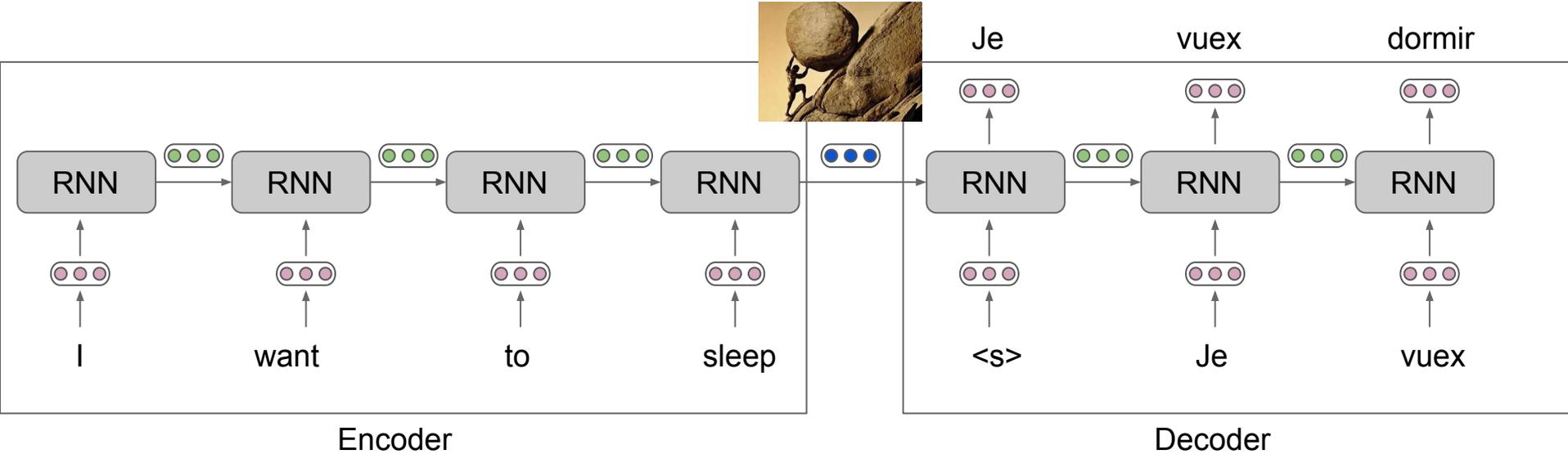


# Sequence Models To Rescue

# Recurrent Neural Networks - RNNs

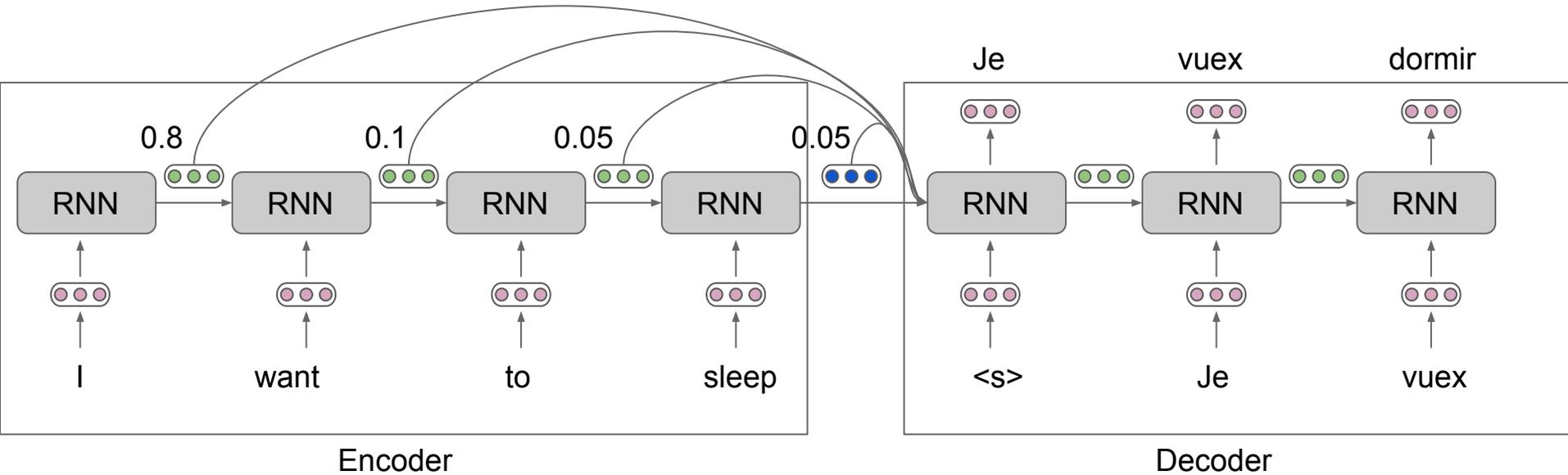


# Encoder-Decoder Models



Typically Encoder and Decoder would be separate networks i.e. have their own separate weights

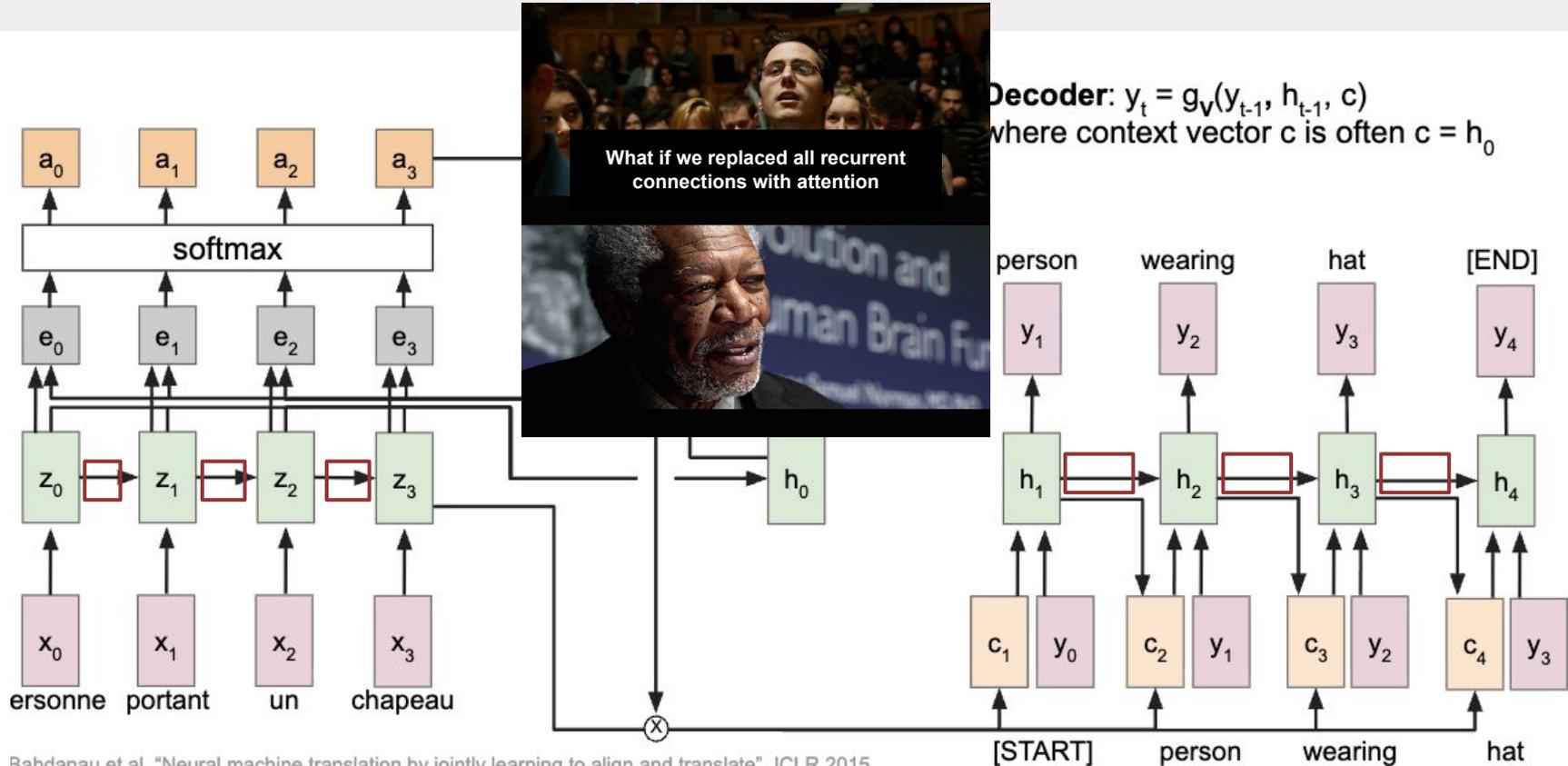
# Encoder-Decoder Models With Attention



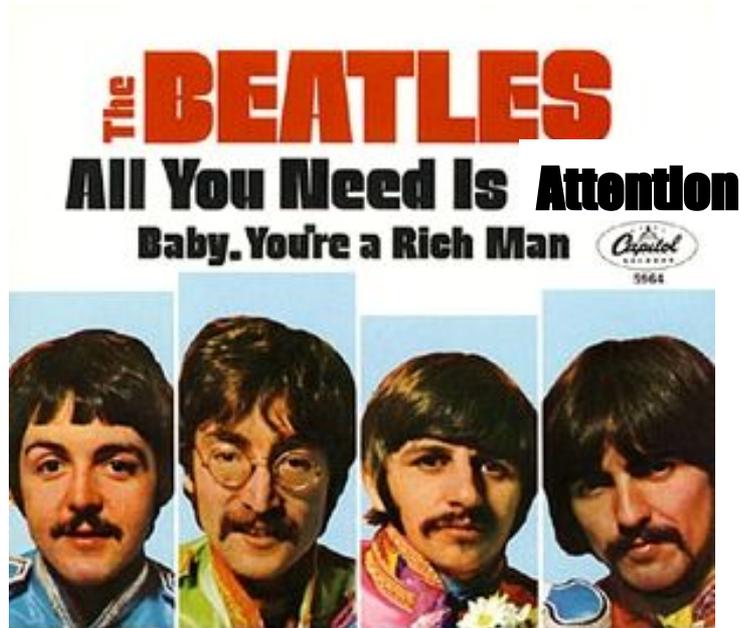
# What's wrong with RNNs?

- Long Range Dependencies
- Gradient vanishing / explosion
- Long time to converge
- Expensive computation

# Encoder-Decoder Models With Attention



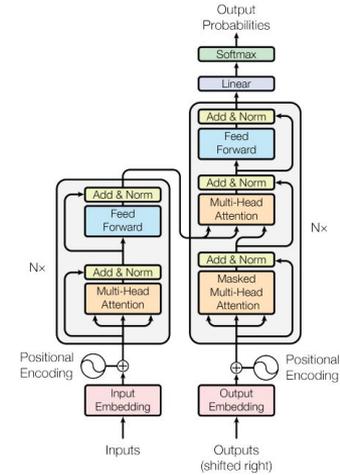
Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015



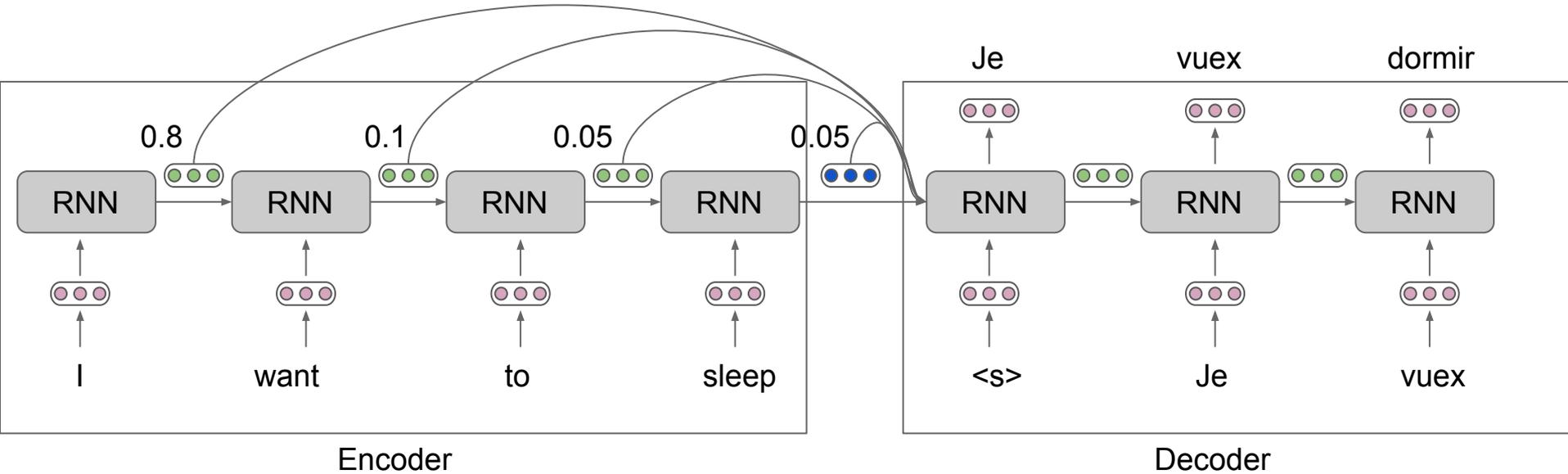
**Attention Is All You Need**

---

<p>Ashish Vaswani* Google Brain avaswani@google.com</p>	<p>Noam Shazeer* Google Brain noam@google.com</p>	<p>Niki Parmar* Google Research niki@google.com</p>	<p>Jakob Uszkoreit* Google Research usz@google.com</p>
<p>Llion Jones* Google Research llion@google.com</p>	<p>Aidan N. Gomez*<sup>1</sup> University of Toronto aidan@cs.toronto.edu</p>	<p>Lukas Kaiser* Google Brain lukasz.kaiser@google.com</p>	
<p>Illia Polosukhin*<sup>1</sup> illia.polosukhin@gmail.com</p>			

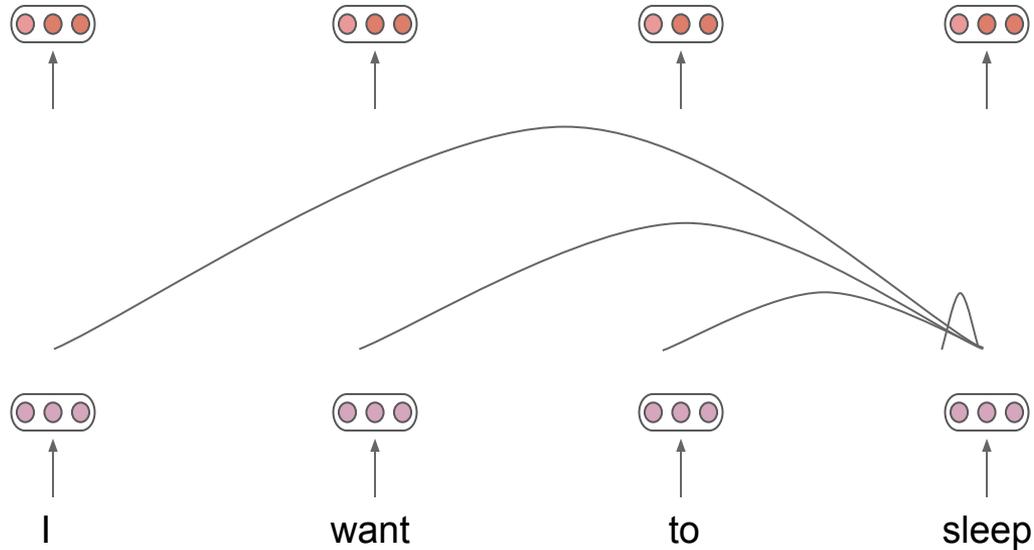


# Last look at RNNs

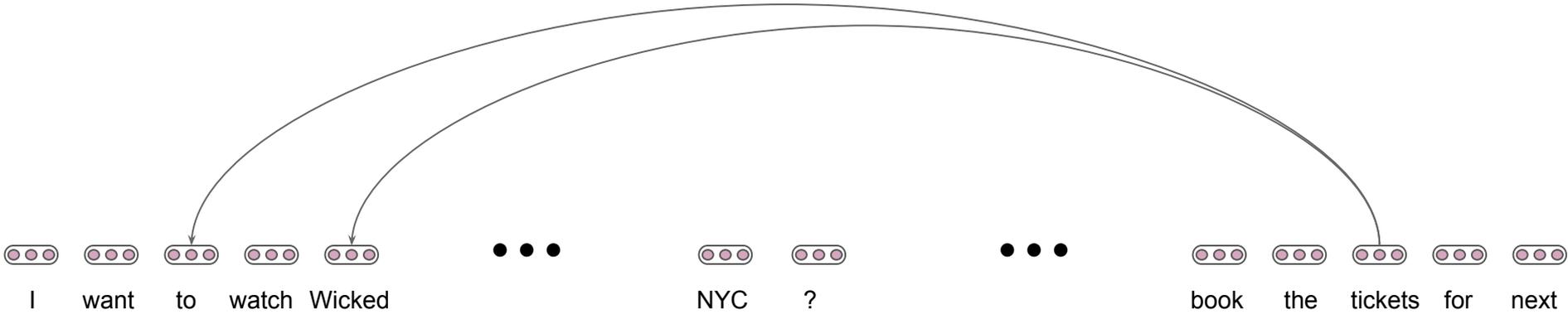


# Self Attention - No more recurrence

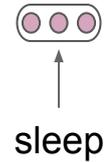
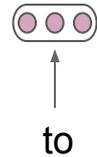
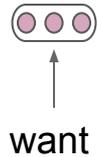
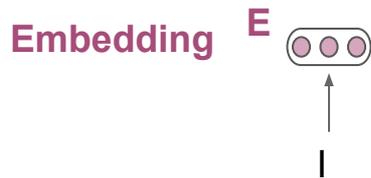
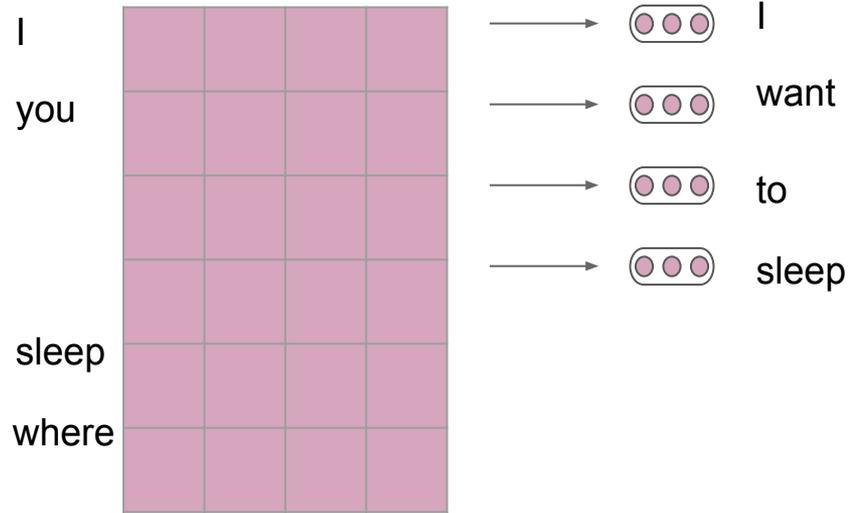
## Contextual Representations



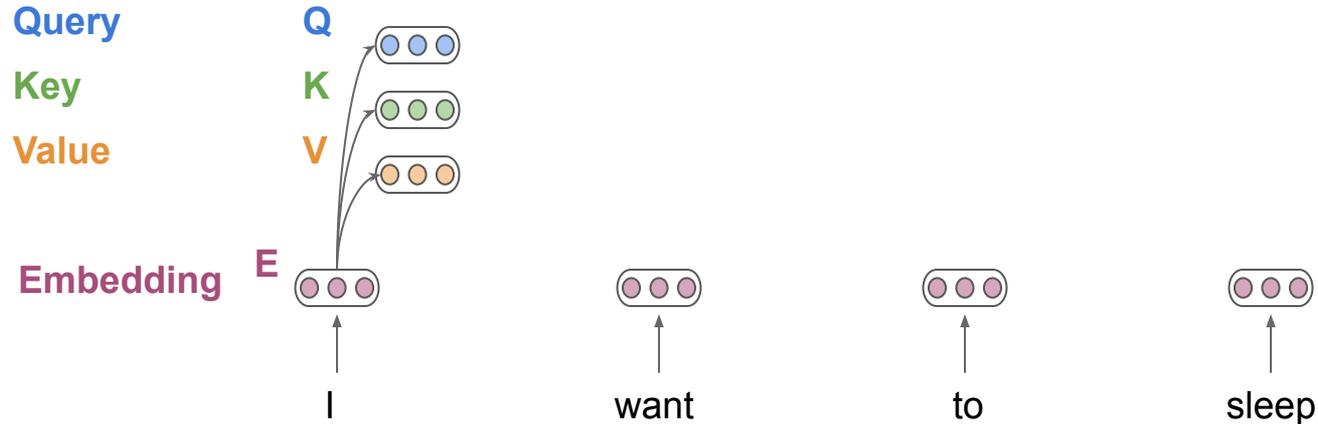
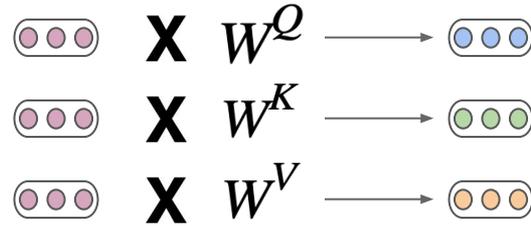
# Self Attention for long sequences



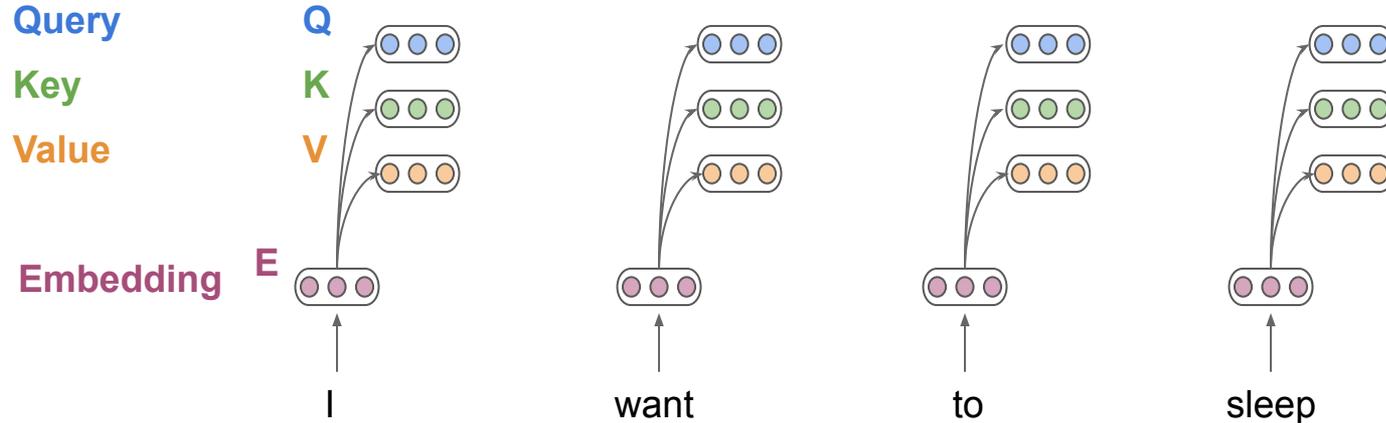
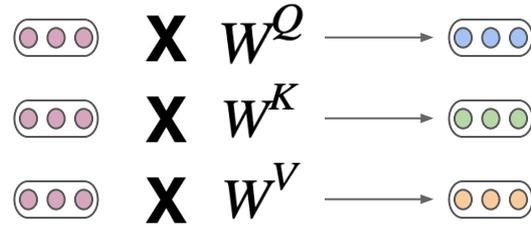
# Self Attention - Word Embedding



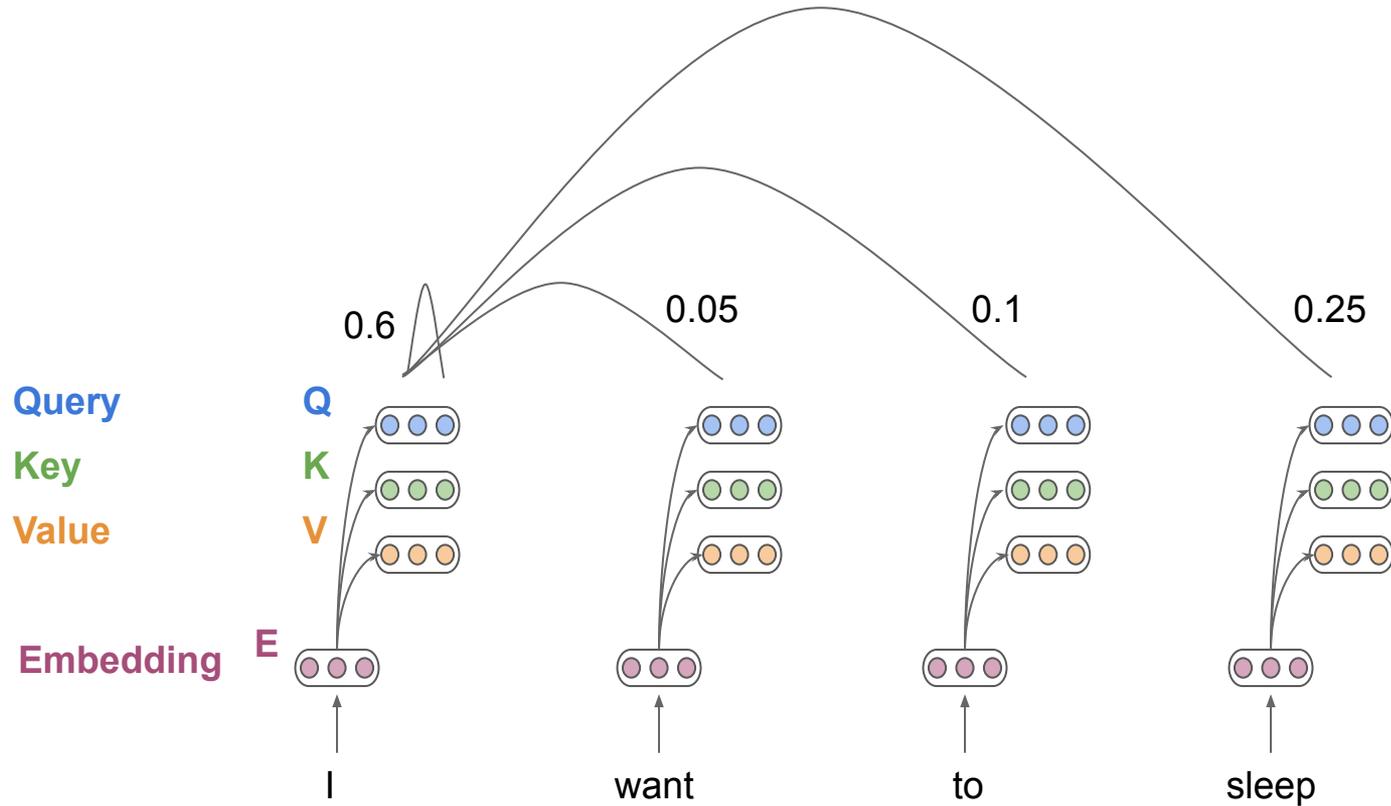
# Self Attention - Projection Layer



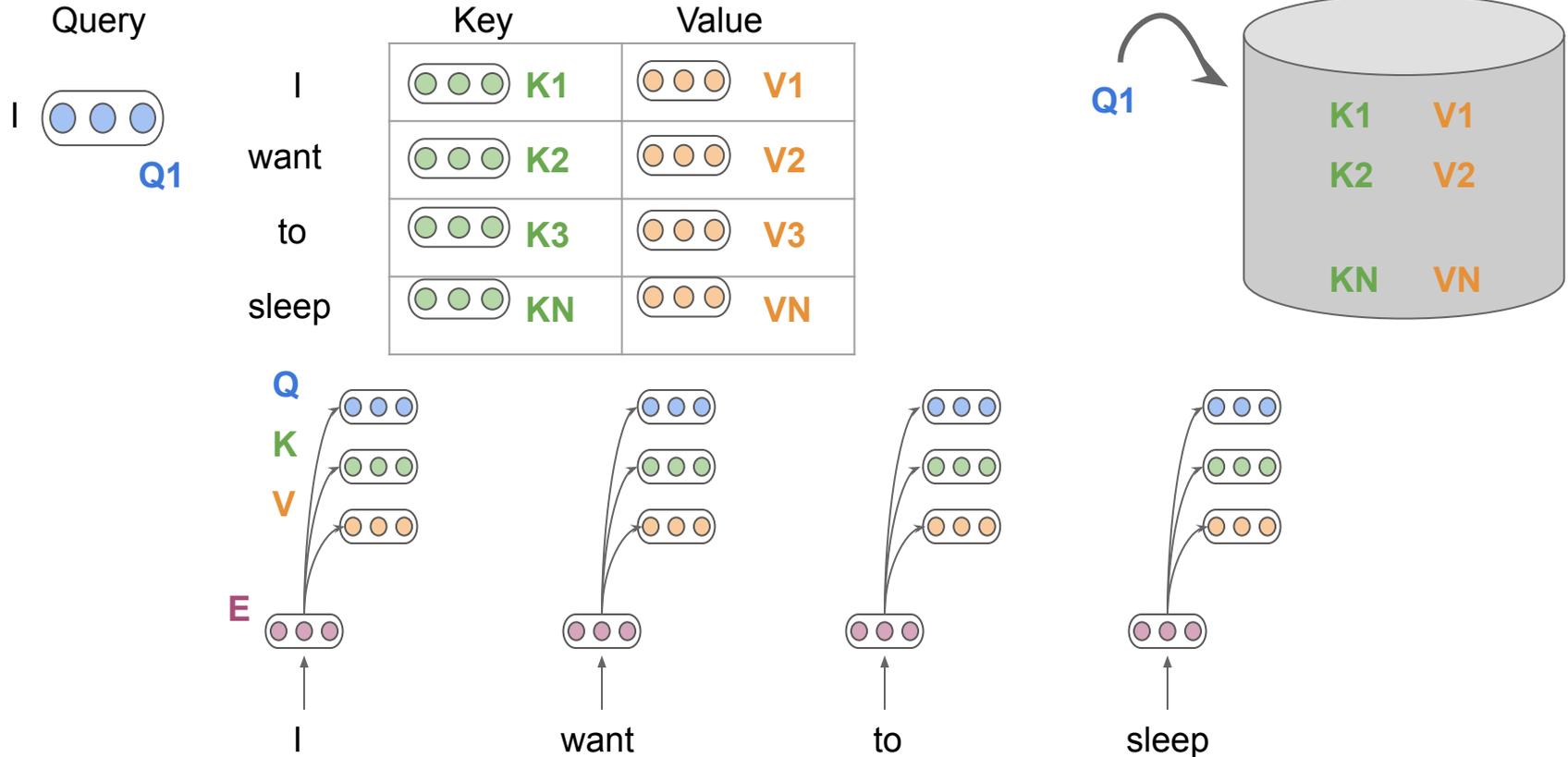
# Self Attention - Projection Layer



# Self Attention - Attention Scores



# Self Attention



# Questions?

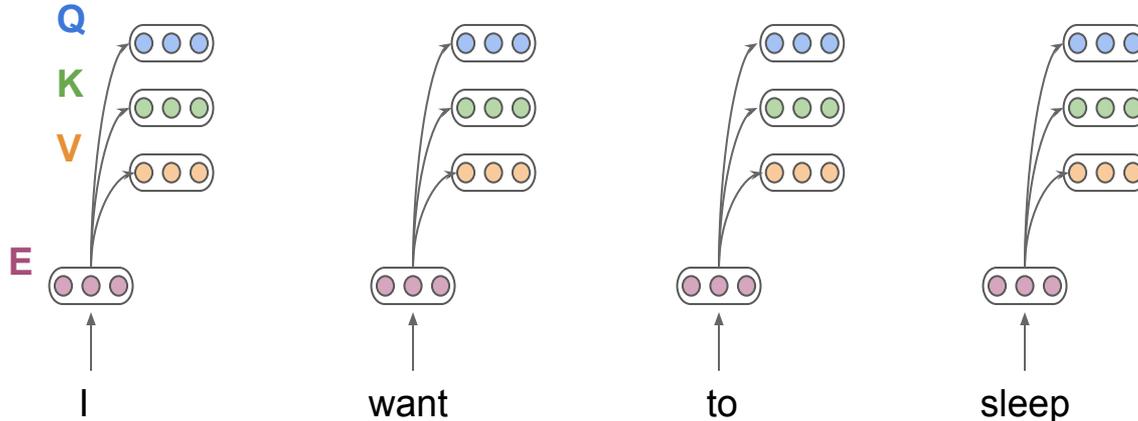
# Self Attention - Scaled Dot Product



Scaled Dot Product

$$SDP = \frac{(QK^T)}{\sqrt{d^k}}$$

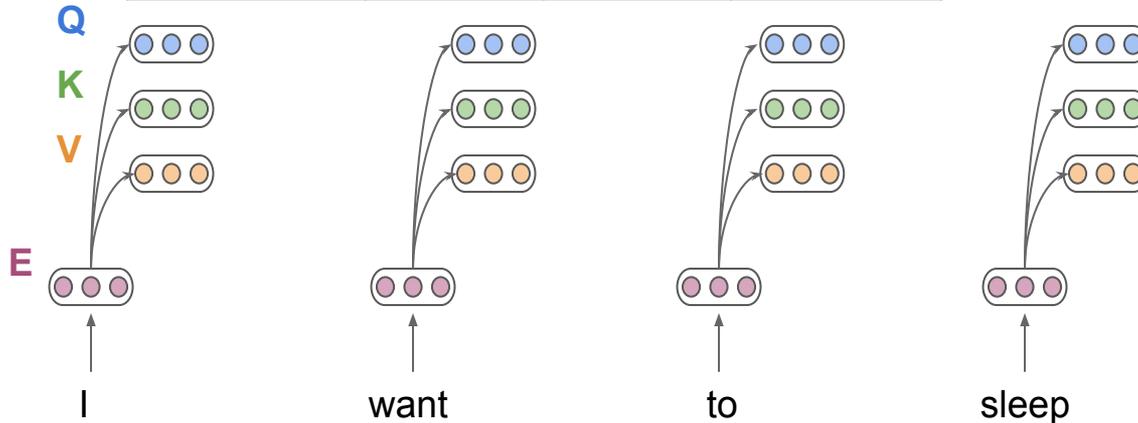
Added to ensure that the dot-product has unit variance



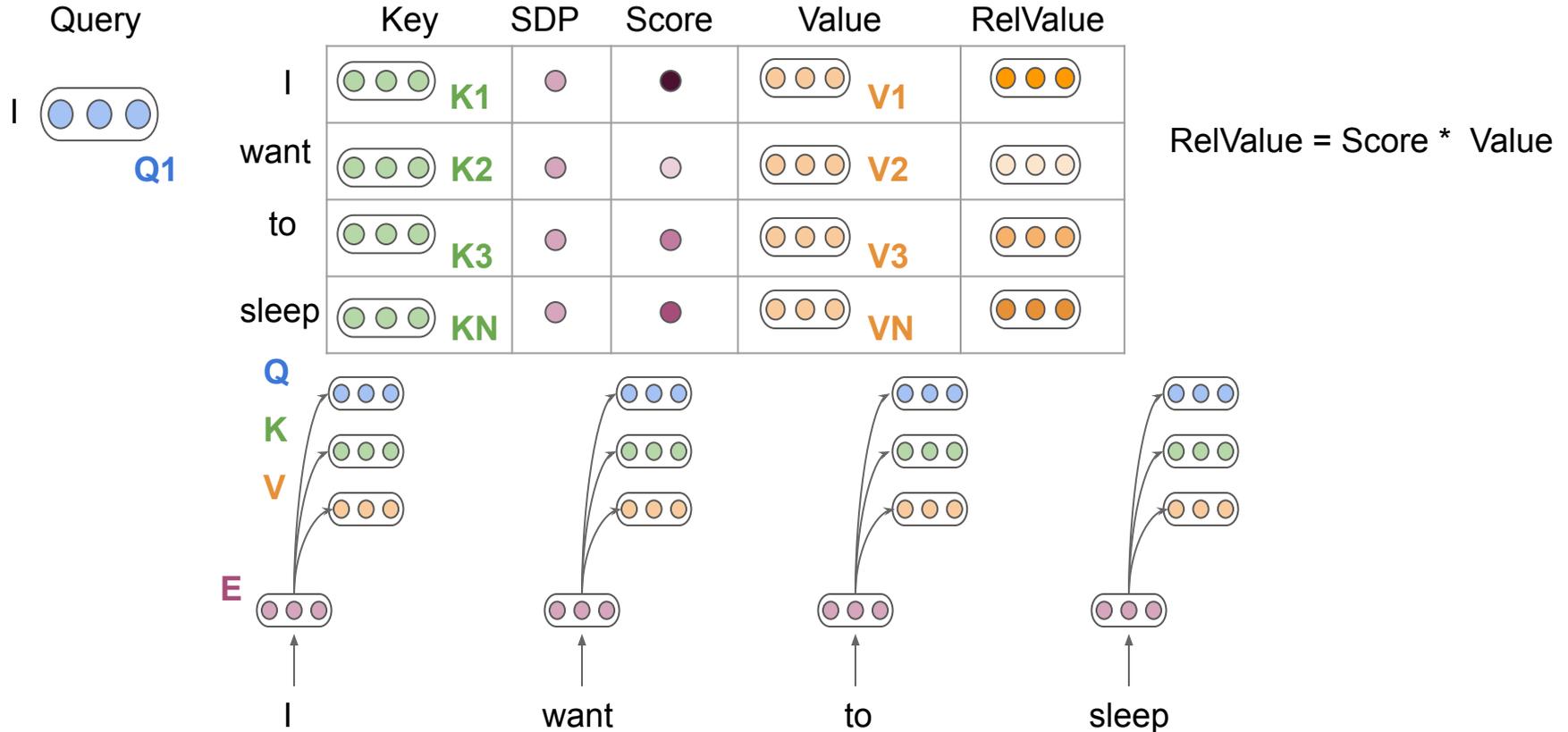
# Self Attention - SoftMax

Query	Key	SDP	Score	Value
I	K1	(95)	(0.6)	V1
want	K2	(13)	(0.05)	V2
to	K3	(35)	(0.1)	V3
sleep	KN	(72)	(0.25)	VN

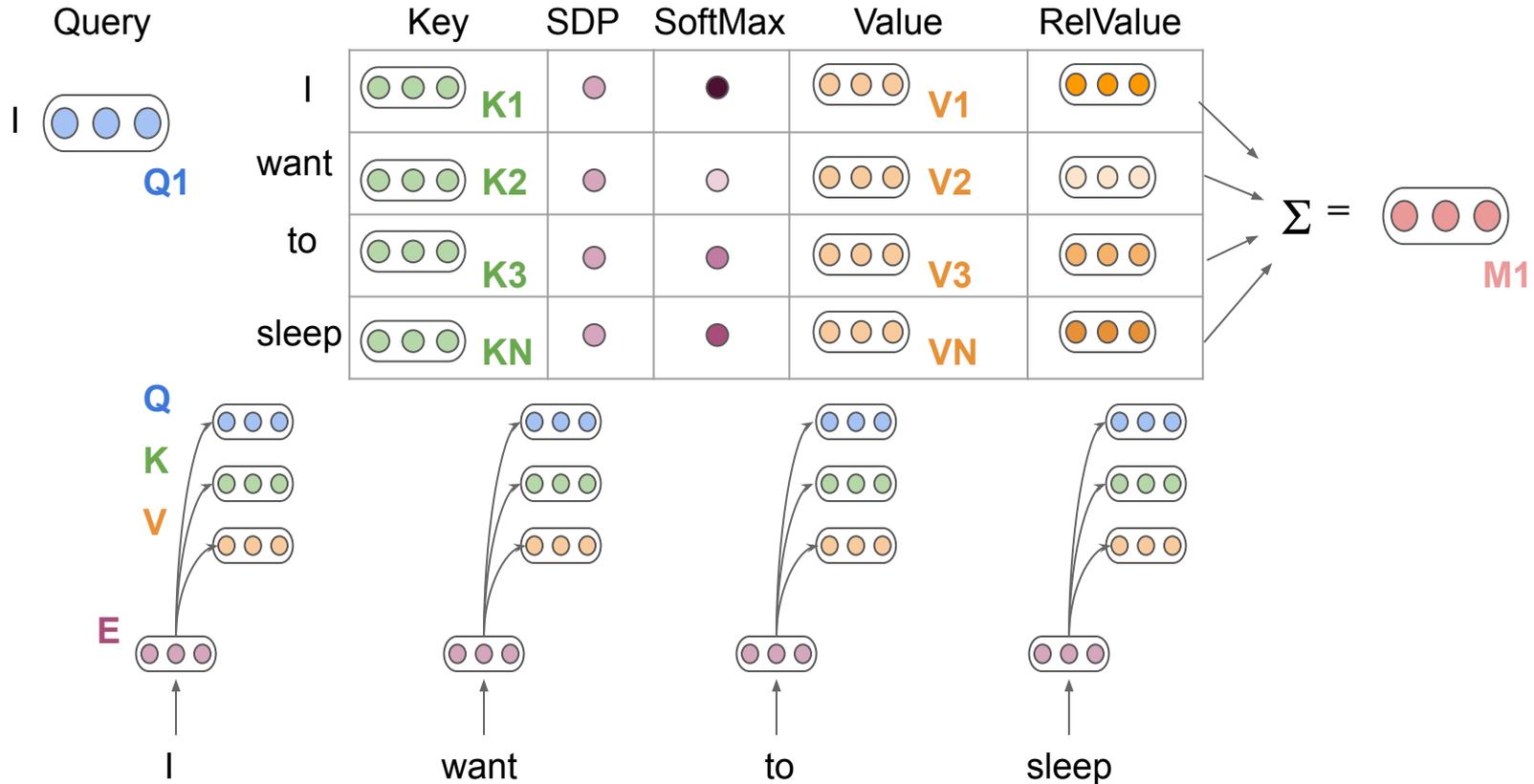
$$\text{score} = \text{softmax} \left( \frac{(QK^T)}{\sqrt{d^k}} \right)$$



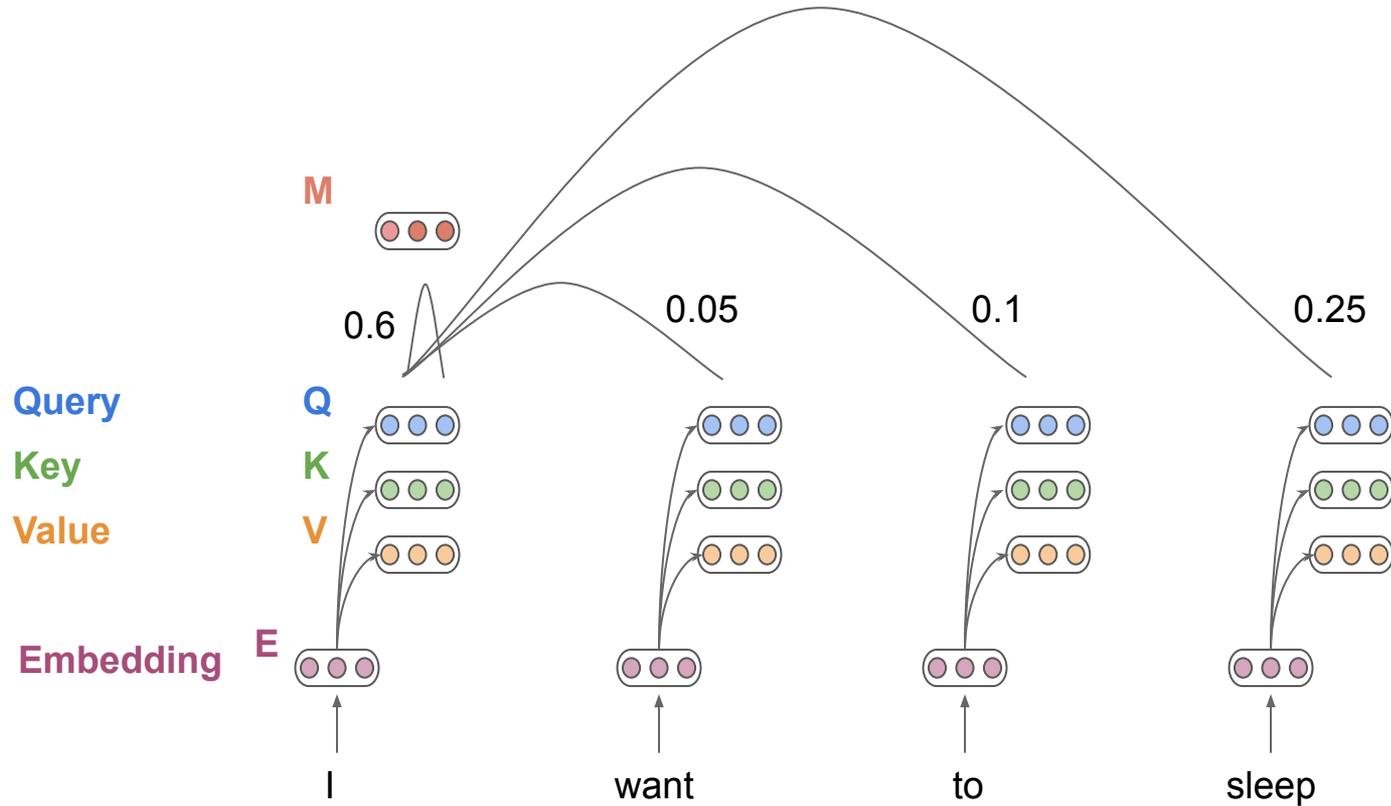
# Self Attention - Soft (Relative) Values



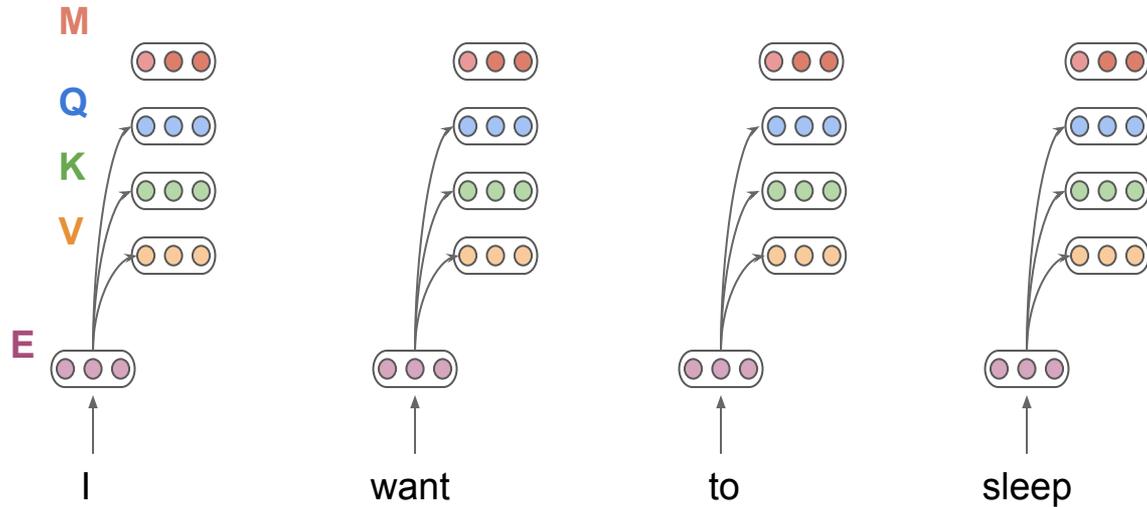
# Self Attention - Attended Repr



# Self Attention - Attended Contextual Rep



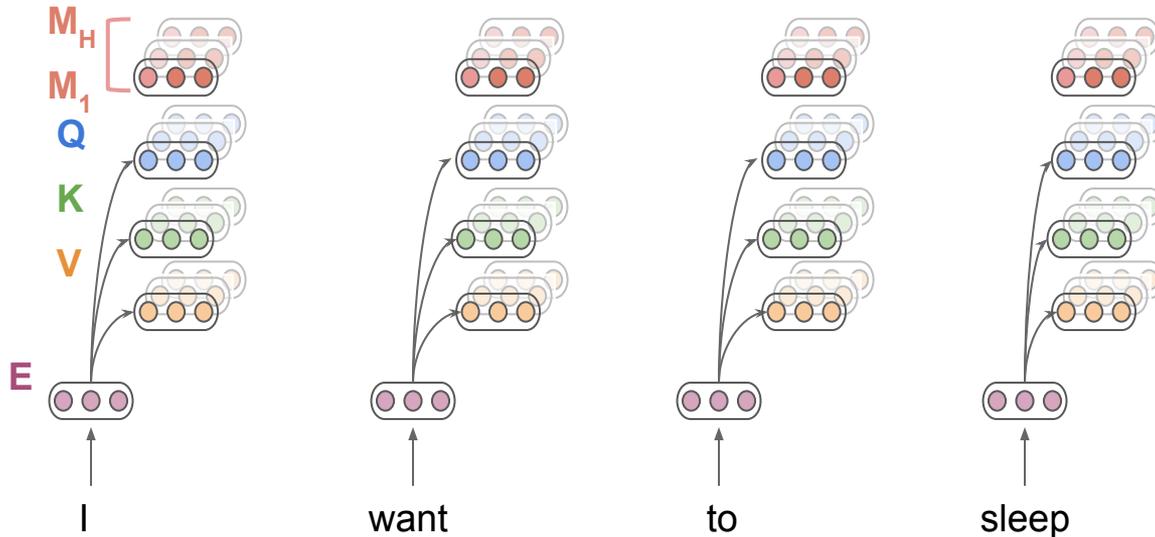
# Self Attention - Attended Contextual Rep



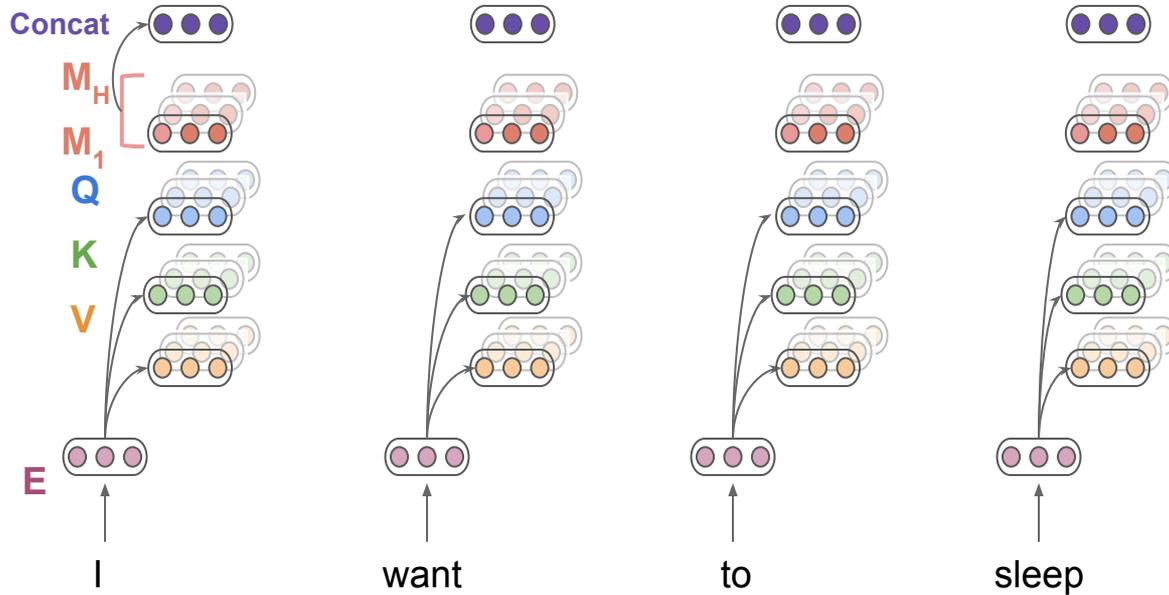
# Multi-Headed Self Attention

Love me some representation power!

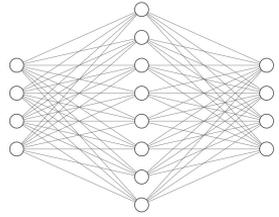
H (no: of heads) Different versions of Q,K,V  
Each different repr -> Different attended repr



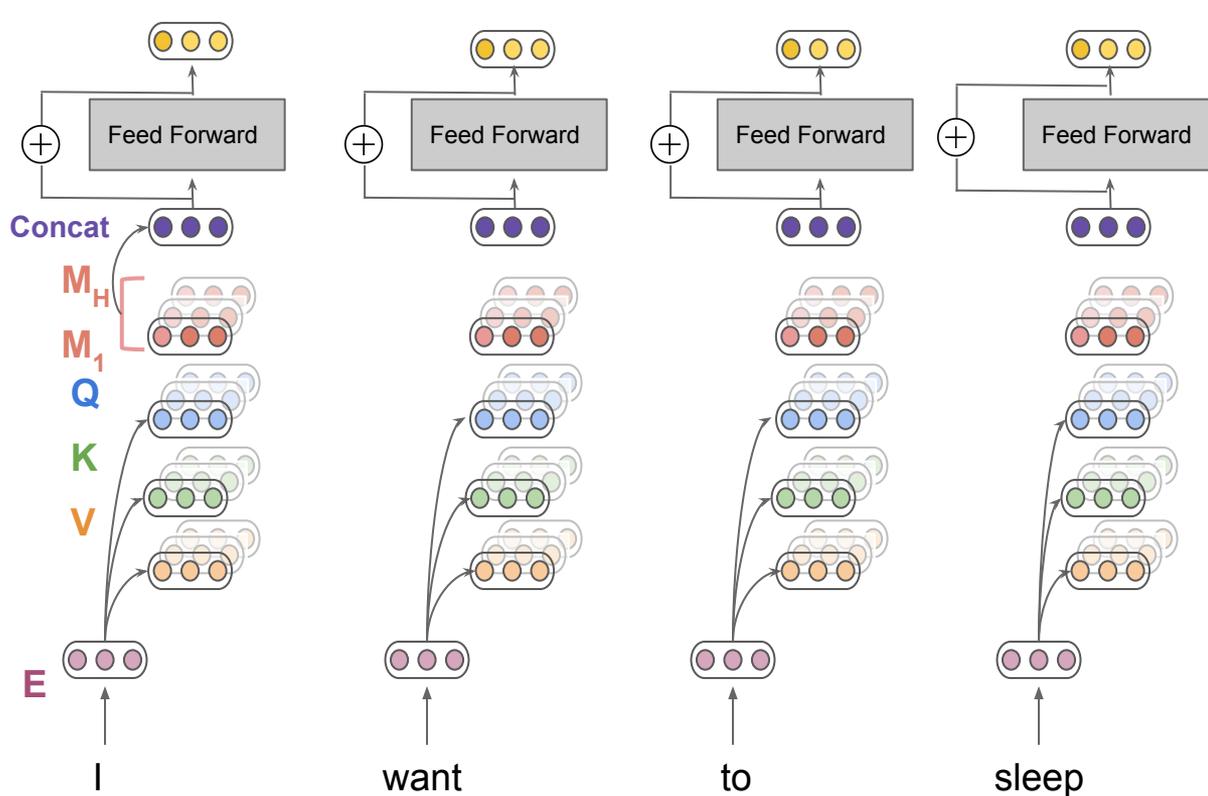
# Multi-Headed Self Attention



# Multi-Headed Self Attention

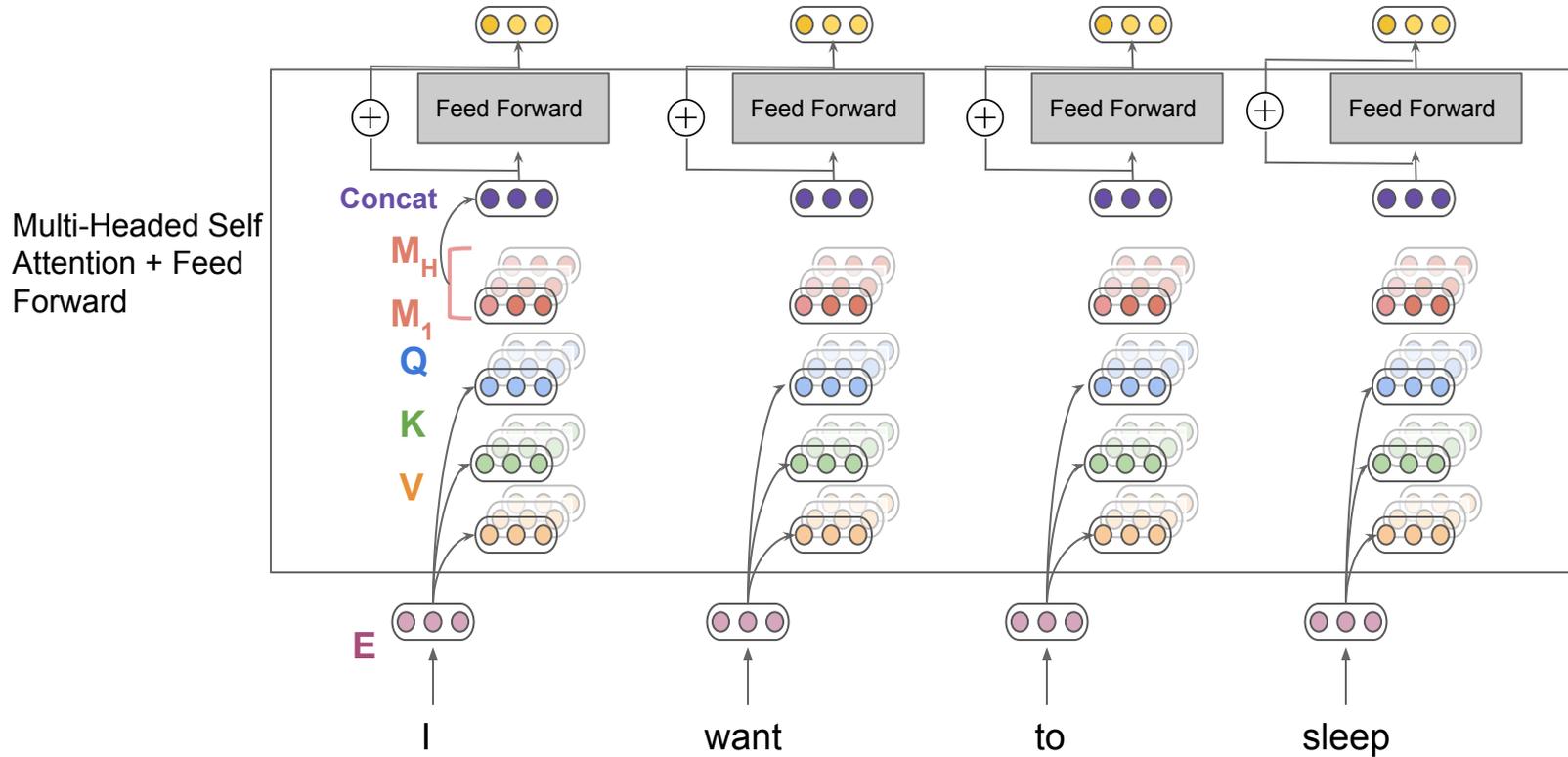


It's Me MLP

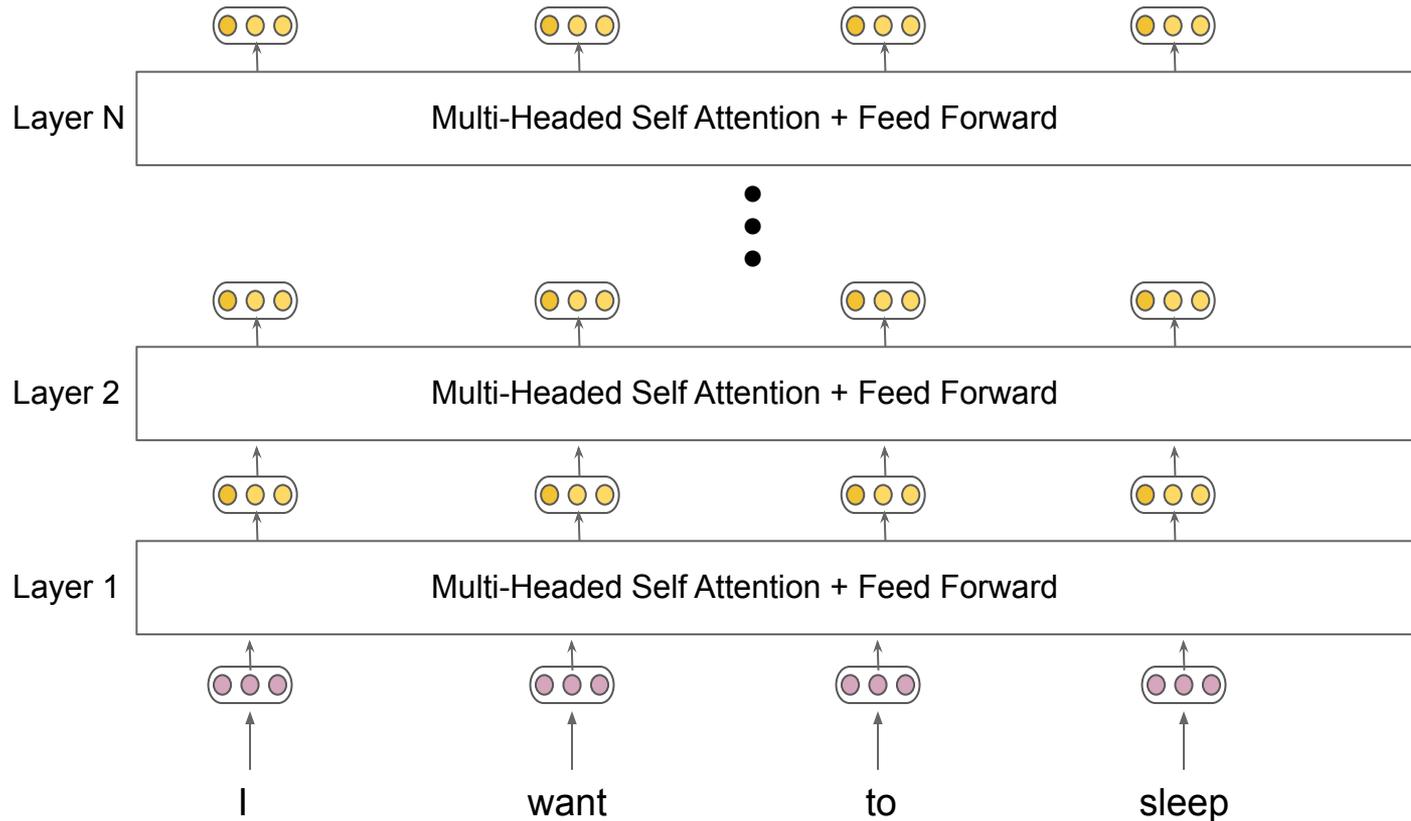


Residual Connections :  
Help with cleaner gradient flows during back-prop

# Multi-Headed Self Attention



# Multi-Headed Self Attention



# Questions?

# Revisiting Self Attention

Query (I)



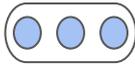
	Key	SDP	SM	Value	RelValue
I	K1			V1	
want	K2			V2	
to	K3			V3	
sleep	KN			VN	

$$\Sigma = \text{capsule with 3 red circles} M$$

I      want      to      sleep

# Revisiting Self Attention

Query (I)



	Key	SDP	SM	Value	RelValue
I	K1			V1	
want	K2			V2	
to	K3			V3	
sleep	KN			VN	

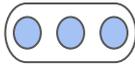
$$\Sigma = \text{capsule with 3 red circles} M$$

I    want    to    sleep

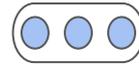
Sleep    to    I    want

# Revisiting Self Attention

Query (I)



Query (I)



	Key	SDP	SM	Value	RelValue
I	K1			V1	
want	K2			V2	
to	K3			V3	
sleep	KN			VN	

$$\Sigma = \text{capsule with 3 red circles} M$$

I    want    to    sleep

Sleep    to    I    want

# Revisiting Self Attention

Query (I)

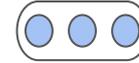


	Key	SDP	SM	Value	RelValue
I	K1			V1	
want	K2			V2	
to	K3			V3	
sleep	KN			VN	

$$\Sigma = \text{ M}$$

I want to sleep

Query (I)



	Key	SDP	SM	Value	RelValue
Sleep	K1			VN	
to	K2			V3	
I	K3			V1	
want	KN			V2	

Sleep to I want

# Revisiting Self Attention

Query (I)



	Key	SDP	SM	Value	RelValue
I	K1			V1	
want	K2			V2	
to	K3			V3	
sleep	KN			VN	

$$\Sigma = \text{ M}$$

I want to sleep

Query (I)



	Key	SDP	SM	Value	RelValue
Sleep	K1			VN	
to	K2			V3	
I	K3			V1	
want	KN			V2	

$$\Sigma = \text{ M}$$

Sleep to I want

# Revisiting Self Attention

Query (I)



Query (I)



	Key	SDP	SM	Value	RelValue
I	K1			V1	
want	K2			V2	
to	K3			V3	
sleep	KN			VN	

	Key	SDP	SM	Value	RelValue
Sleep	K1			VN	
to	K2			V3	
I	K3			V1	
want	KN			V2	

**Same representation for both sentences - But positions matter!**

$$\Sigma = \text{ M}$$

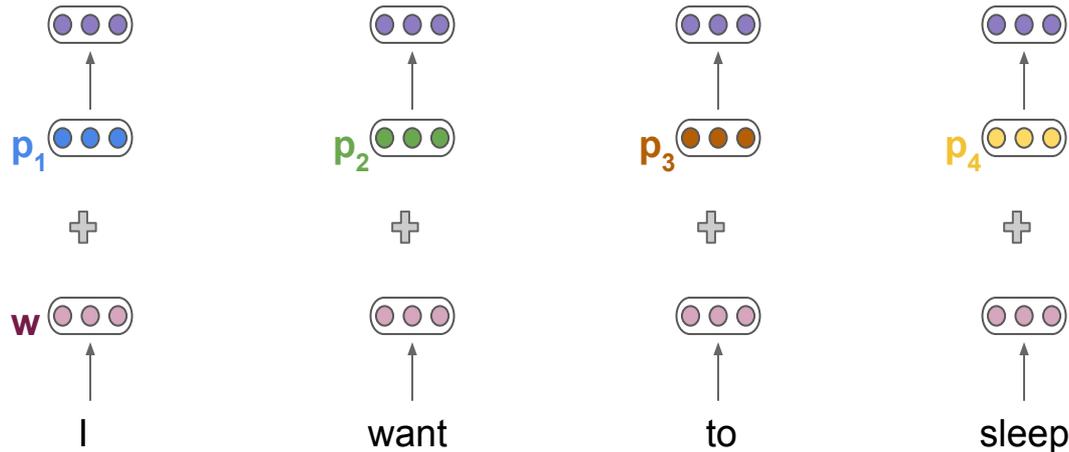
I want to sleep

$$\Sigma = \text{ M}$$

Sleep to I want

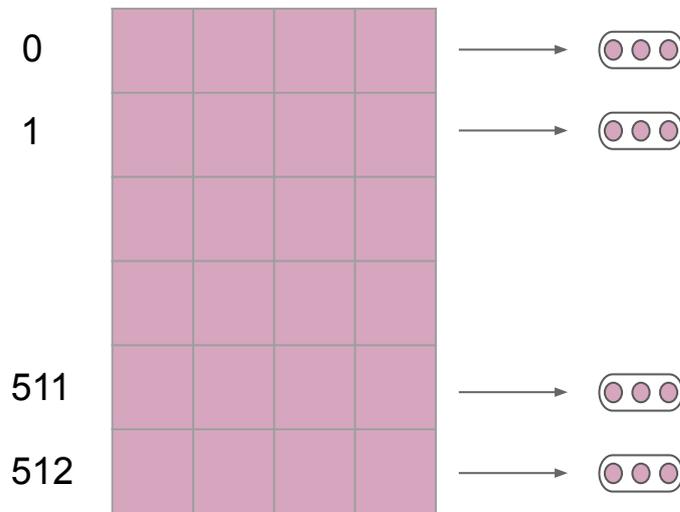
# Positional Encoding

Position embeddings - each position number has an associated embedding



# Learnable Positional Encoding

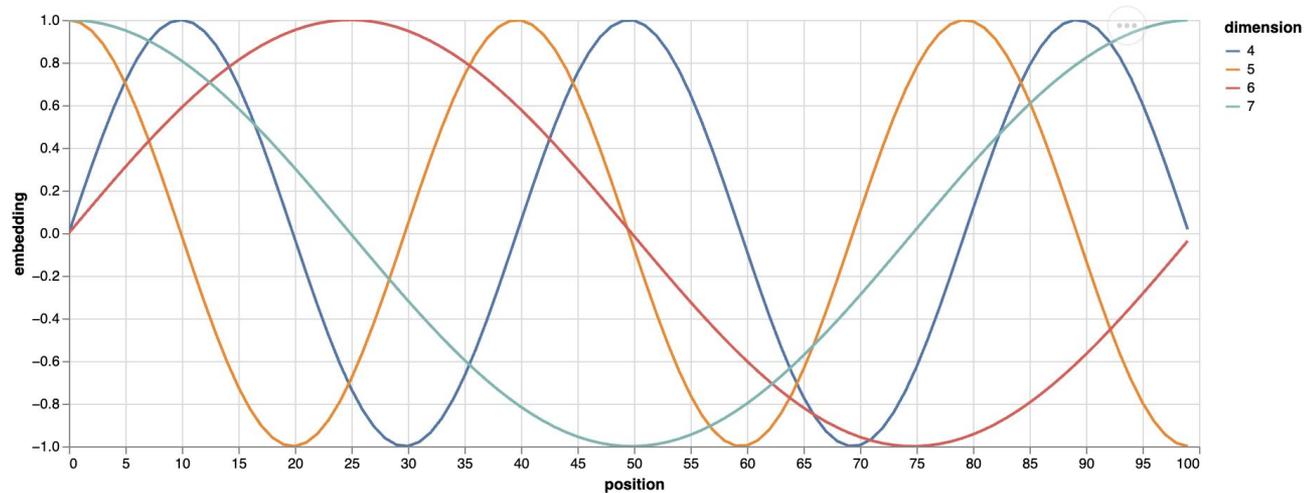
- Simplest type of positional encodings
- Initialize position encodings as random vectors for each position from 0 to MAX LENGTH
- Used in GPT-1, GPT-2, GPT-3
- CONS: Doesn't generalize to sequences of length greater than MAX LENGTH



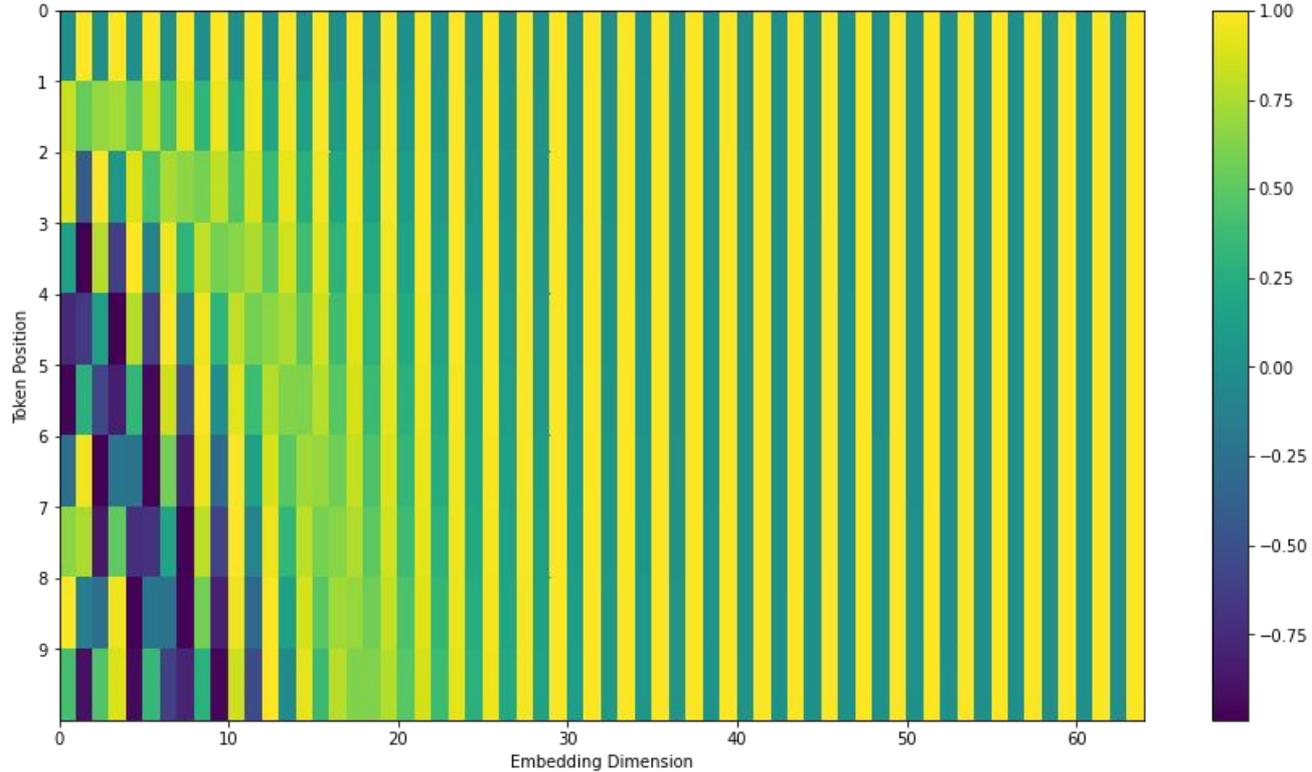
# Sinusoidal Positional Encoding

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

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

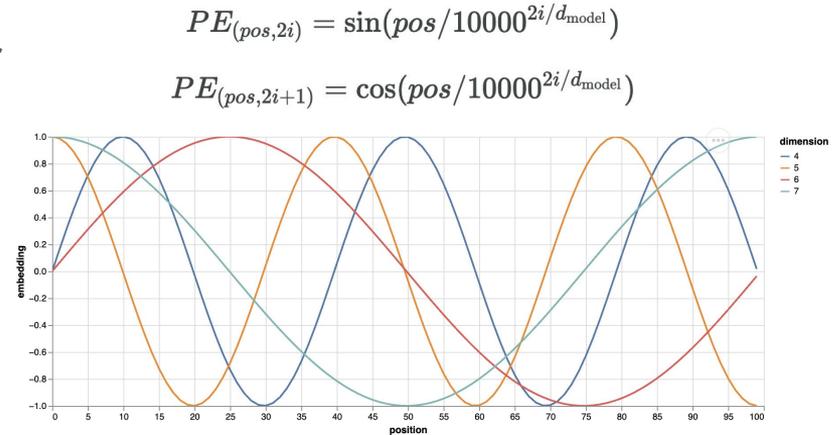


# Sinusoidal Positional Encoding



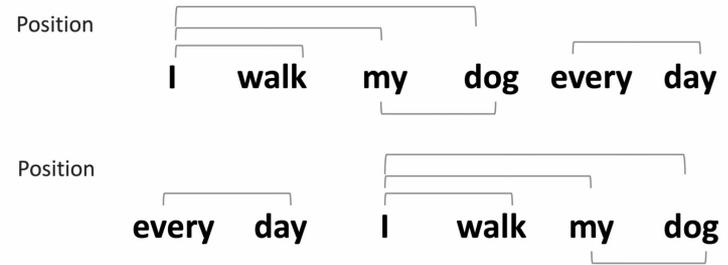
# Sinusoidal Positional Encoding

- Introduced in the original Transformers paper and was used in models like BERT
- Major promise of this embeddings over using learnable encodings was that these might lead to length generalization beyond maximum length seen during training
- Alas, that's not the case and sinusoidal embeddings are usually as bad as learnable embeddings when it comes to length generalization



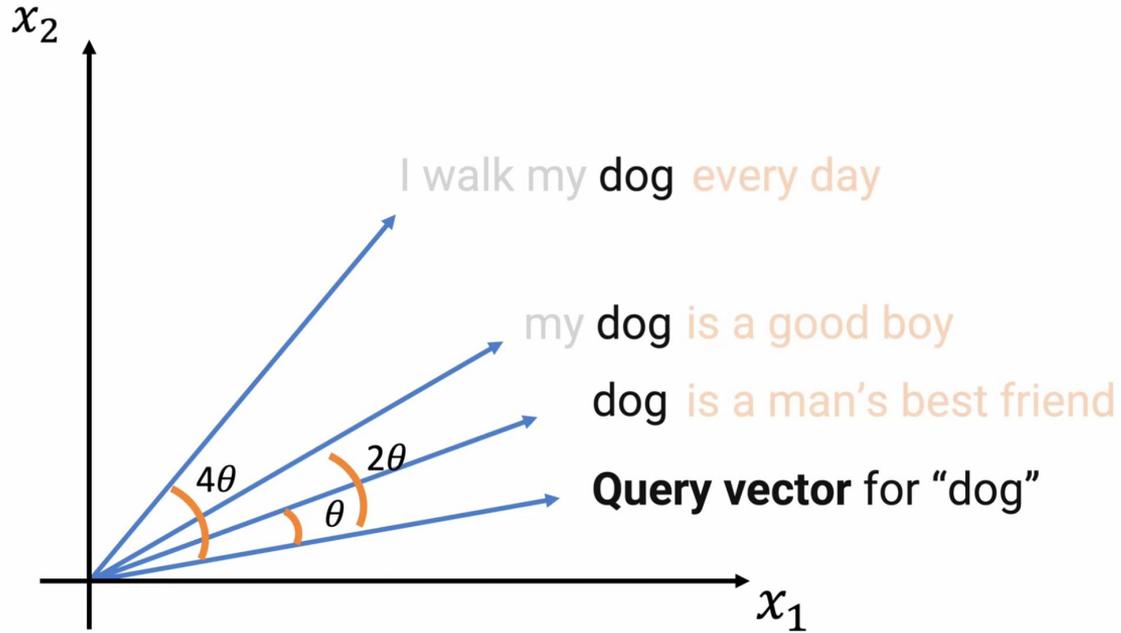
# Relative Position Encodings

- Both learnable and sinusoidal position encodings were examples of Absolute Position Encodings
- i.e. positional information depends on the absolute position of the token
- In relative position encodings we provide the positional information in form of the relative difference between token position
- This information is added while taking the dot product between the query and key vectors



From Jia-Bin Huan's Youtube Video: <https://www.youtube.com/watch?v=SMBklmDWOyQ&t=683s>

# Rotary Positional Encodings (RoPE)



Rotary Positional Embeddings, 2021

From Jia-Bin Huan's Youtube Video: <https://www.youtube.com/watch?v=SMBkImDWOyQ&t=683s>

# Rotary Positional Encodings (RoPE)

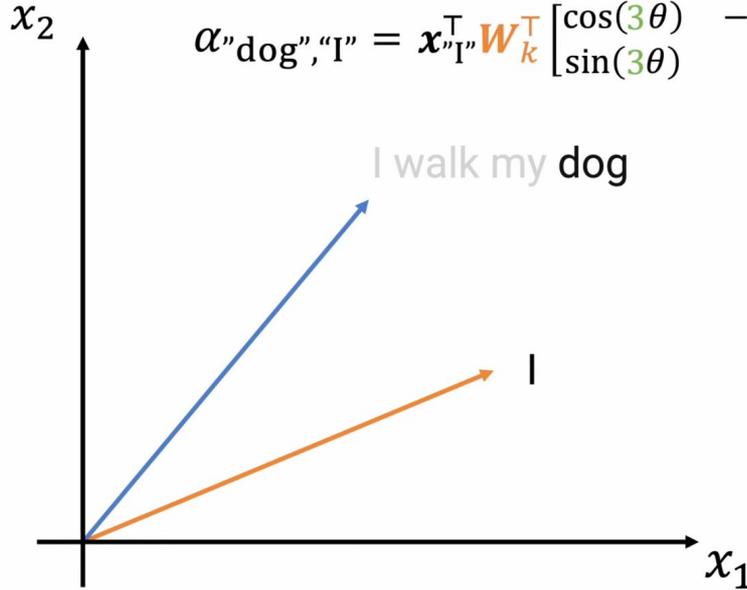
$$q_{\text{"dog"}} = \begin{bmatrix} \cos(4\theta) & -\sin(4\theta) \\ \sin(4\theta) & \cos(4\theta) \end{bmatrix} W_q \mathbf{x}_{\text{"dog"}}$$

Query vector for "dog"

$$k_{\text{"I"}} = \begin{bmatrix} \cos(1\theta) & -\sin(1\theta) \\ \sin(1\theta) & \cos(1\theta) \end{bmatrix} W_k \mathbf{x}_{\text{"I"}}$$

Key vector for "I"

$$\alpha_{\text{"dog"}, \text{"I"}} = \mathbf{x}_{\text{"I"}}^T W_k^T \begin{bmatrix} \cos(3\theta) & -\sin(3\theta) \\ \sin(3\theta) & \cos(3\theta) \end{bmatrix} W_q \mathbf{x}_{\text{"dog"}}$$



From Jia-Bin Huan's Youtube Video: <https://www.youtube.com/watch?v=SMBkImDWOyQ&t=683s>

# Rotary Positional Encodings (RoPE)

$$q_{\text{"dog"}} = \begin{bmatrix} \cos(6\theta) & -\sin(6\theta) \\ \sin(6\theta) & \cos(6\theta) \end{bmatrix} W_q x_{\text{"dog"}}$$

Query vector for "dog"

$$k_{\text{"I"}} = \begin{bmatrix} \cos(3\theta) & -\sin(3\theta) \\ \sin(3\theta) & \cos(3\theta) \end{bmatrix} W_k x_{\text{"I"}}$$

Key vector for "I"

$$\alpha_{\text{"dog", "I"}} = x_{\text{"I"}}^T W_k^T \begin{bmatrix} \cos(3\theta) & -\sin(3\theta) \\ \sin(3\theta) & \cos(3\theta) \end{bmatrix} W_q x_{\text{"dog"}}$$



From Jia-Bin Huan's Youtube Video: <https://www.youtube.com/watch?v=SMBkImDWOyQ&t=683s>

# Rotary Positional Encodings (RoPE)

$$R_{\Theta}^m W_q \mathbf{x}_m$$

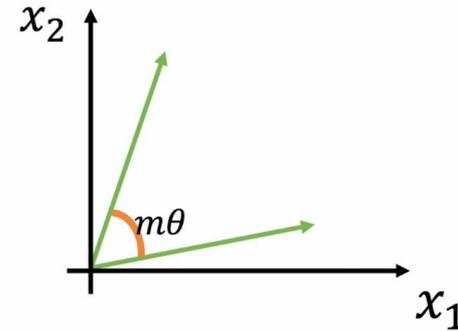
query vector for the token at position  $m$

$$R_{\Theta}^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$

$$R_{\Theta}^n W_k \mathbf{x}_n$$

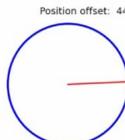
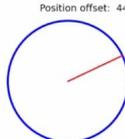
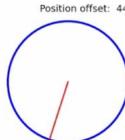
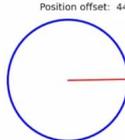
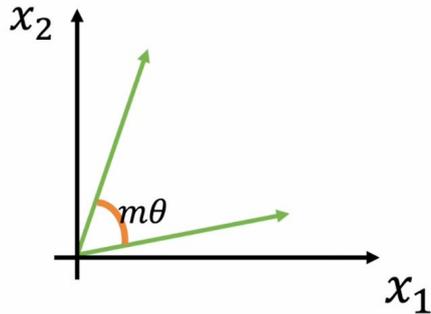
key vector for the token at position  $n$

$$\alpha_{m,n} = \mathbf{x}_n^T W_k^T R_{\Theta}^{m-n} W_q \mathbf{x}_m$$



# Rotary Positional Encodings (RoPE)

$$R_{\Theta}^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$



$$\begin{bmatrix} \square \\ \square \end{bmatrix} = \begin{bmatrix} \cos(m\theta_1) & -\sin(m\theta_1) \\ \sin(m\theta_1) & \cos(m\theta_1) \end{bmatrix} \begin{bmatrix} \square \\ \square \end{bmatrix}$$

$$\begin{bmatrix} \square \\ \square \end{bmatrix} = \begin{bmatrix} \cos(m\theta_2) & -\sin(m\theta_2) \\ \sin(m\theta_2) & \cos(m\theta_2) \end{bmatrix} \begin{bmatrix} \square \\ \square \end{bmatrix}$$

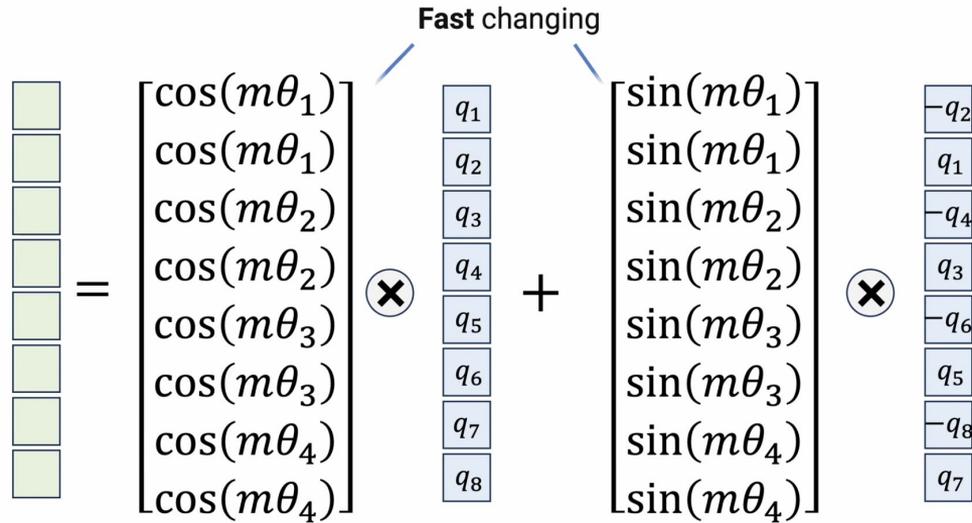
$$\begin{bmatrix} \square \\ \square \end{bmatrix} = \begin{bmatrix} \cos(m\theta_3) & -\sin(m\theta_3) \\ \sin(m\theta_3) & \cos(m\theta_3) \end{bmatrix} \begin{bmatrix} \square \\ \square \end{bmatrix}$$

$$\begin{bmatrix} \square \\ \square \end{bmatrix} = \begin{bmatrix} \cos(m\theta_4) & -\sin(m\theta_4) \\ \sin(m\theta_4) & \cos(m\theta_4) \end{bmatrix} \begin{bmatrix} \square \\ \square \end{bmatrix}$$

$W_q \mathbf{x}_m$

From Jia-Bin Huan's Youtube Video: <https://www.youtube.com/watch?v=SMBkImDWOyQ&t=683s>

# Rotary Positional Encodings (RoPE)



$$f_q(\mathbf{x}_m, m)$$

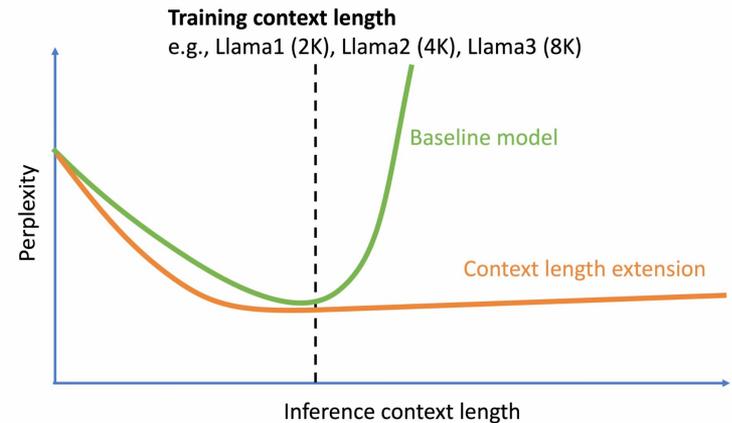
$$\theta_i = N^{-2i/d}$$

$$N = 10,000$$

From Jia-Bin Huan's Youtube Video: <https://www.youtube.com/watch?v=SMBkImDWoyQ&t=683s>

# Rotary Positional Encodings (RoPE)

- Generally show faster convergence than Absolute Position Encoding methods
- The length generalization is still not guaranteed!
- Unseen relative distances i.e. the ones that exceed the maximum sequence length during training remain an issue
- Neural-Tangent-Kernel (NTK) RoPE is an extension to RoPE that enables generalization to long sequences

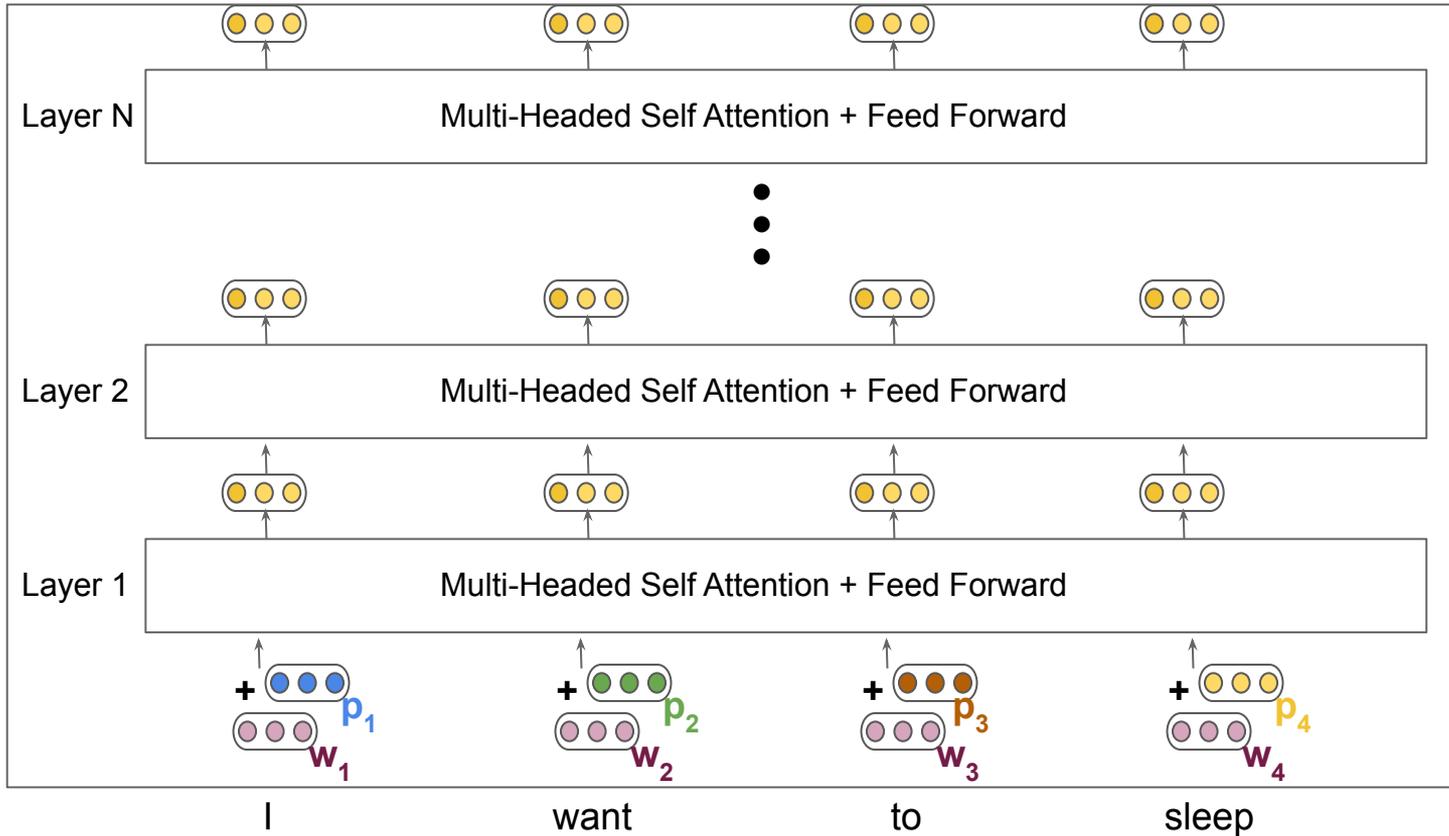


From Jia-Bin Huan's Youtube Video: <https://www.youtube.com/watch?v=SMBkImDWOyQ&t=683s>

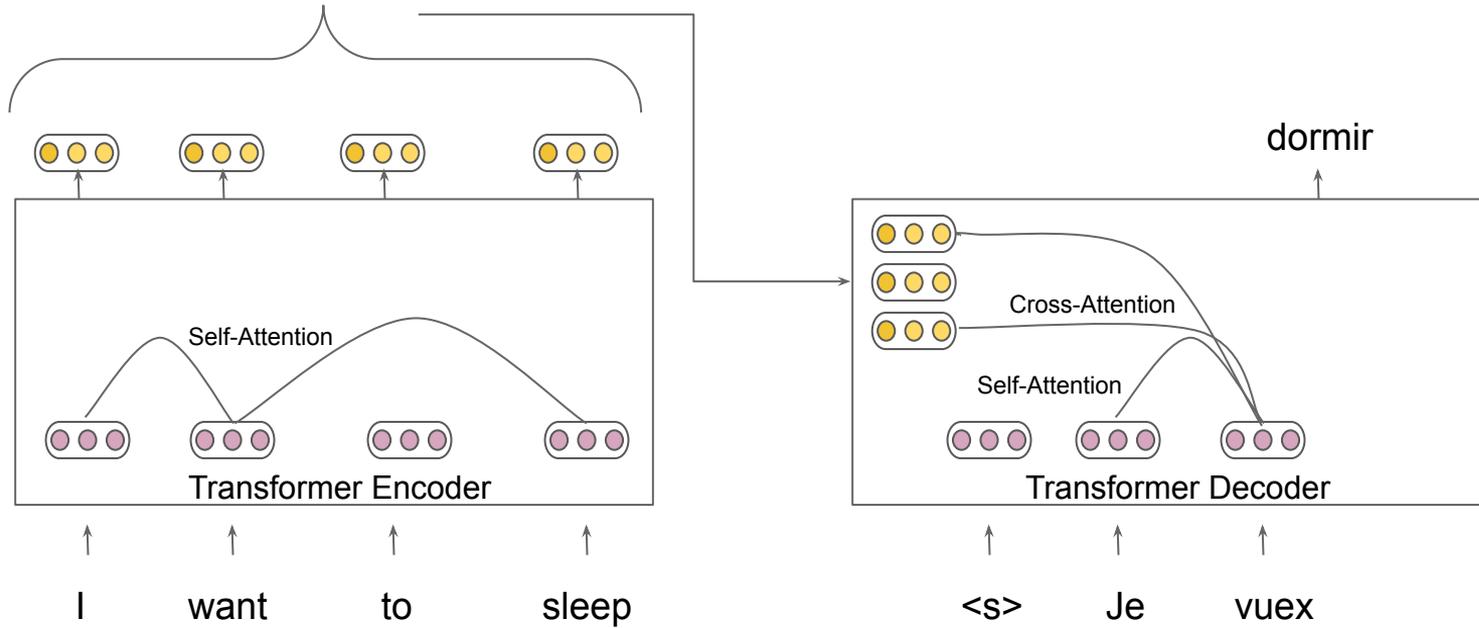
# Questions?

# Transformer Encoder

N-Layer  
Transformer  
Encoder



# Transformer Encoder - Decoder

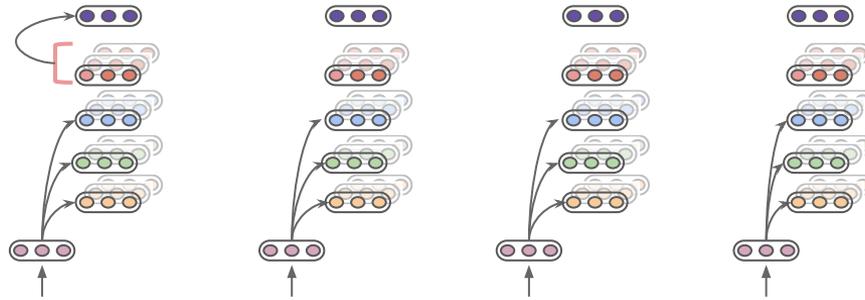


# What's so great about Transformers?

- Parallelizable computation
  - Entire sequence, All queries, all attention heads computed in parallel
  - Benefits from fast matrix multiplication on GPUs
- Rich expressive power
  - Every token connected to every other token
  - Can form long range dependencies
- Depth not proportional to seq length
  - Reduces exploding/vanishing gradient problem
  - Converges faster

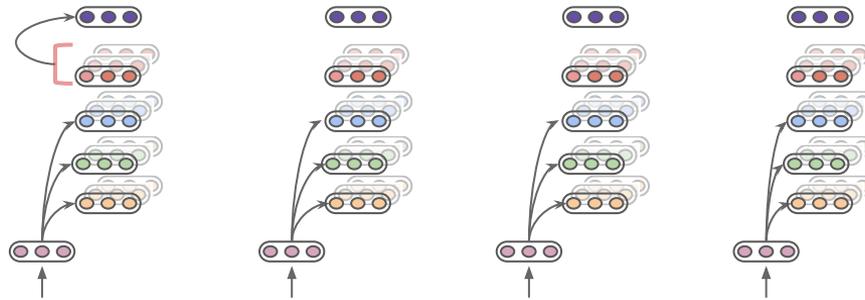
# What's so great about Transformers?

- Parallelizable computation - Entire sequence can be processed in parallel

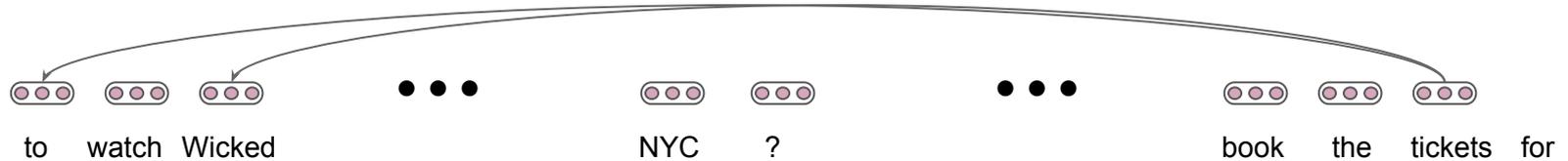


# What's so great about Transformers?

- Parallelizable computation - Entire sequence can be processed in parallel



- Rich expressive power - long range dependencies



# Thank you!

kahuja@cs.washington.edu

# Results/Impact

- Improves results, Establishes SOTA in various tasks!
  - Machine Translation
  - Constituency Parsing
  - Language Modeling
  - and more!
- Computationally faster!
  - No sequential computation - Entire sequence processed in parallel