

Natural Language Processing

Self Attention and Transformers

Vidhisha Balachandran

vbalacha@cs.cmu.edu



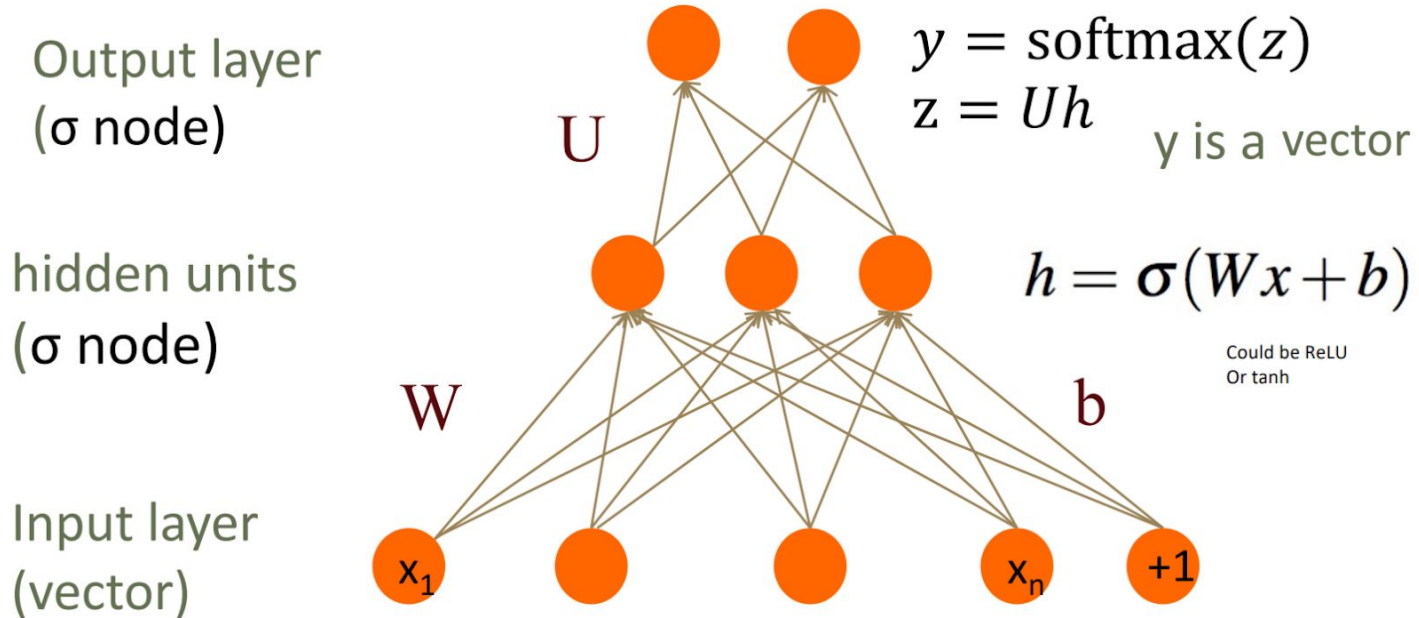
Announcements

- Please fill our survey!
- 3 thing that you like about the course
- 3 things that could be improved

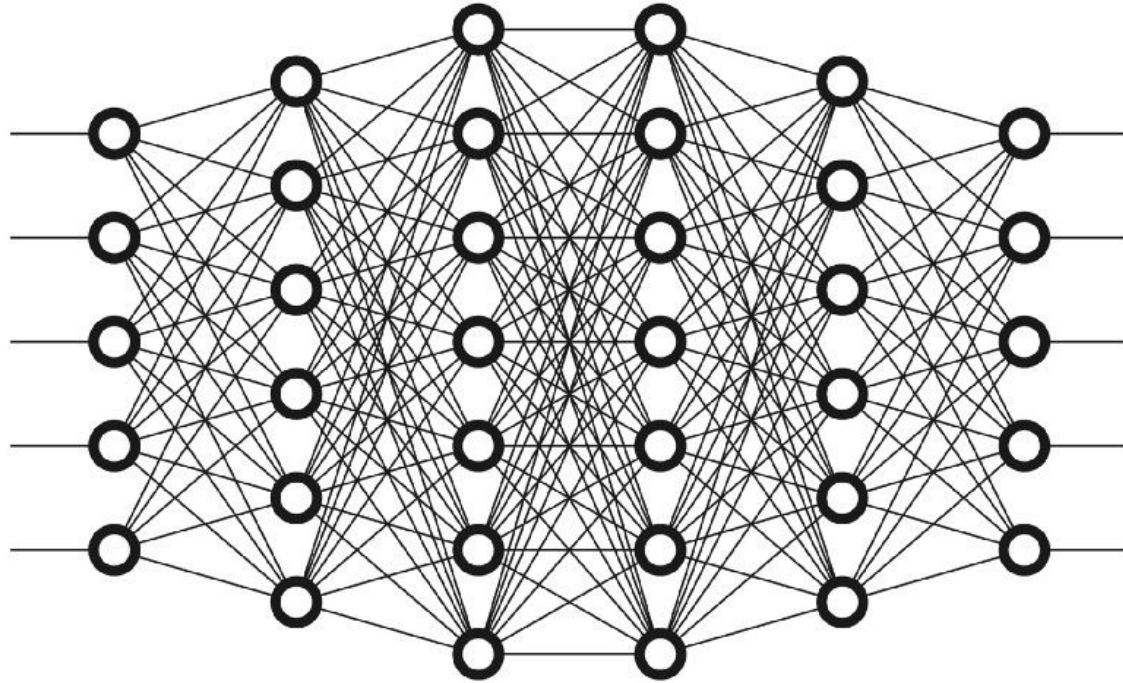
Readings

- [Attention Is All You Need](#)
- [The Illustrated Transformer](#)
- [The Annotated Transformer](#)
- [Language Modeling with Transformers and PyTorch](#)

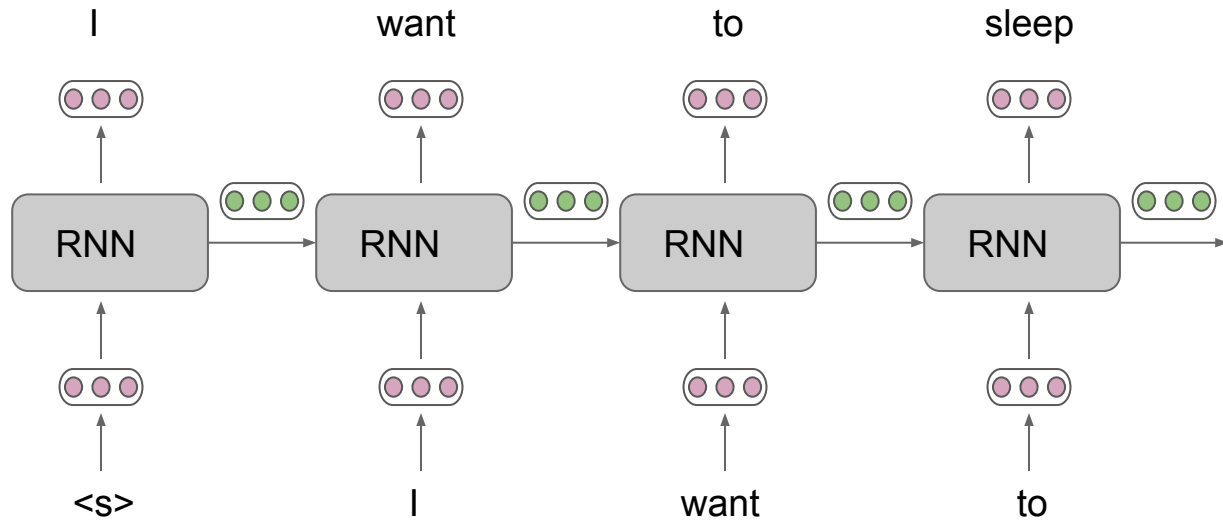
Recap - 2 Layer MLP



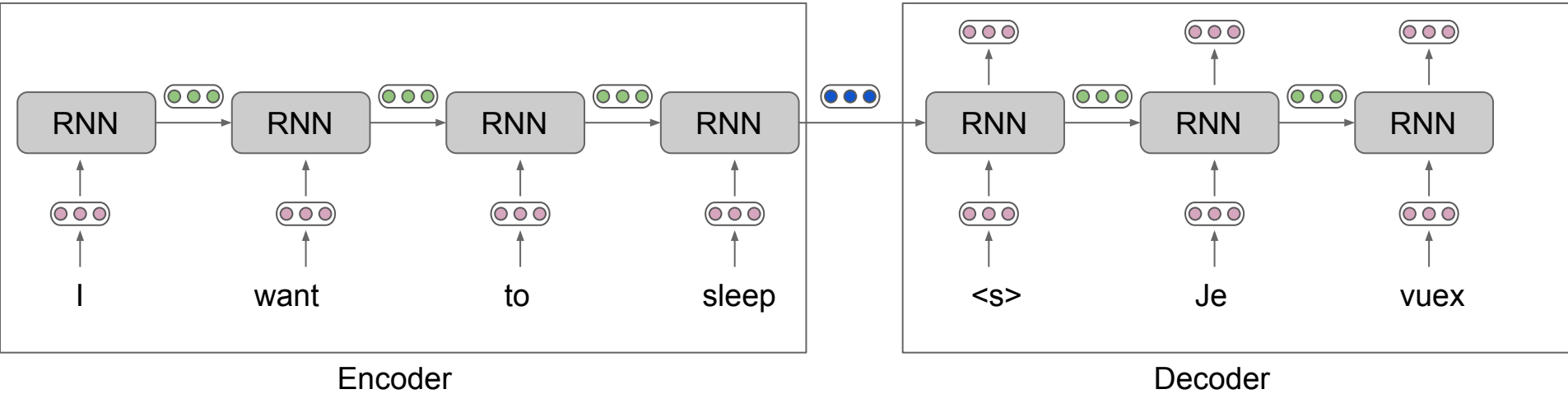
Deep MLP



Recurrent Neural Networks - RNNs



Encoder-Decoder Models



Limitations

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

Long Range Dependencies

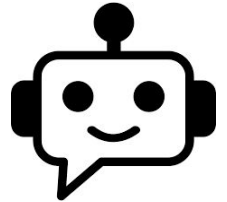


I'm want to watch Wicked! How does the weather in NYC look next week?

It looks sunny with some light rain during the weekend.

Oh! But I don't have a rain jacket :(Is there a store nearby?

There's a marshall's a mile away. They have the navy blue jacket you have been eyeing for a while!



Long Range Dependencies



I'm want to watch Wicked! How does the weather in NYC look next week?

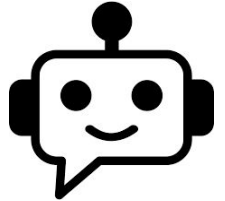
It looks sunny with some light rain during the weekend.

Oh! But I don't have a rain jacket :(Is there a store nearby?

There's a marshall's a mile away. They have the navy blue jacket you have been eyeing for a while!



Ok! Looks like I can actually go! Book **the tickets** for next Wed!



Long Range Dependencies



I'm want to watch **Wicked!** How does the weather in NYC look next week?

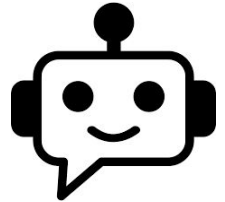
It looks sunny with some light rain during the weekend.

Oh! But I don't have a rain jacket :(Is there a store nearby?

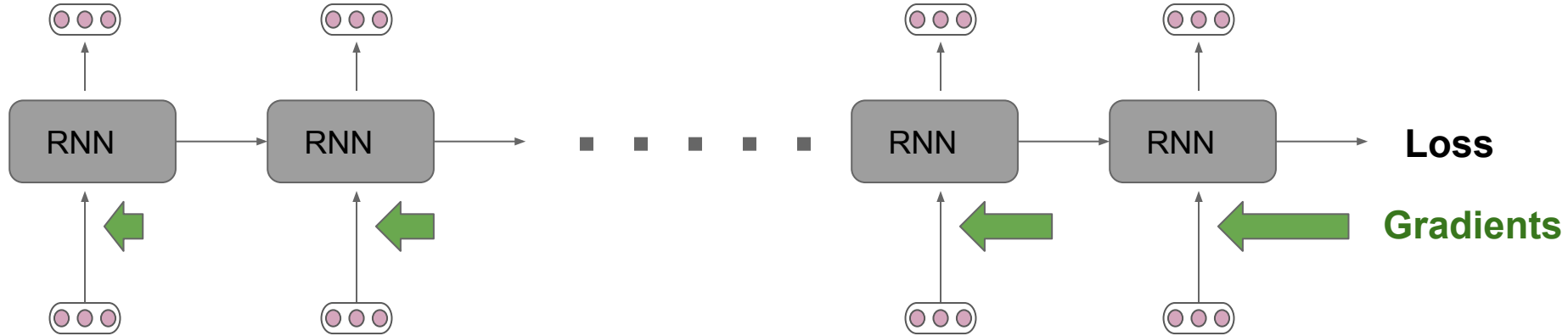
There's a marshall's a mile away. They have the navy blue jacket you have been eyeing for a while!

⋮

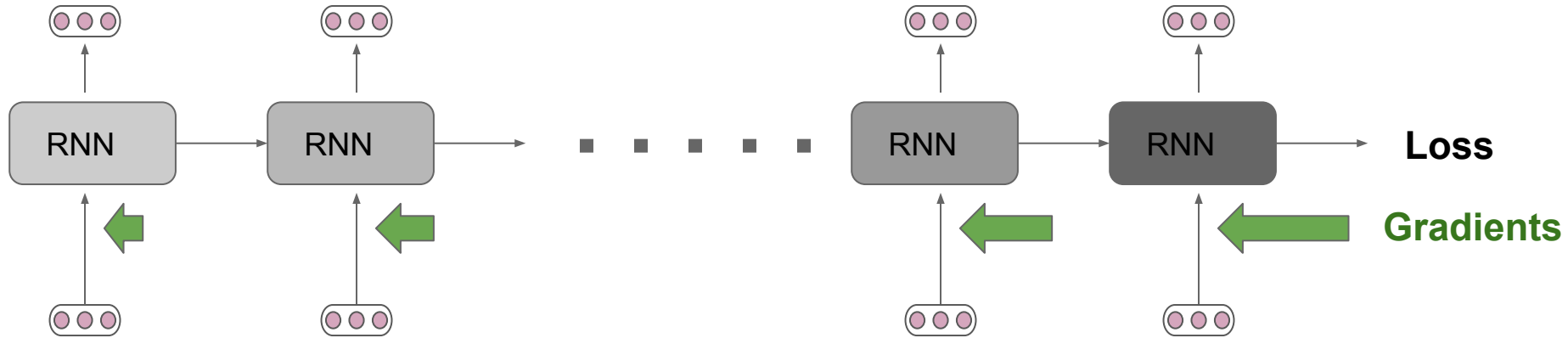
Ok! Looks like I can actually go! Book **the tickets** for next Wed!



Gradient vanishing / explosion



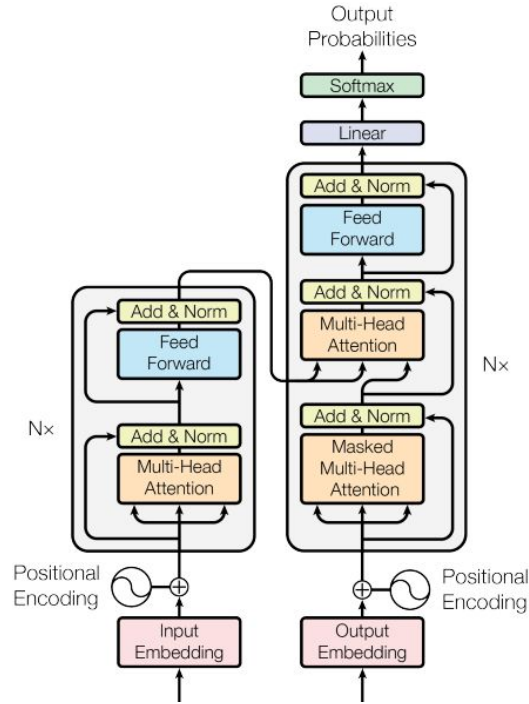
Gradient vanishing / explosion



Limitations

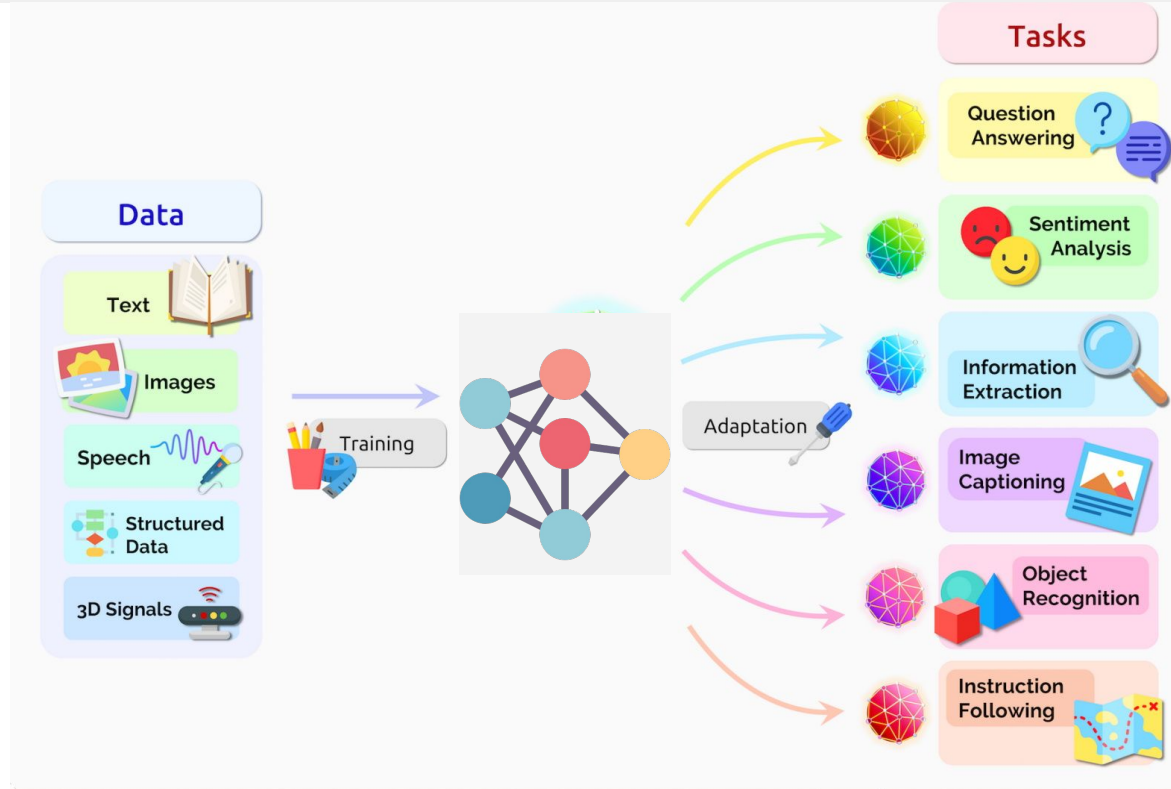
- Long Range Dependencies
- Gradient vanishing / explosion
- Long time to converge
- Expensive computation

Transformer Model



Attention is all you need (Vaswani et.al, 2017)

Wide Applications

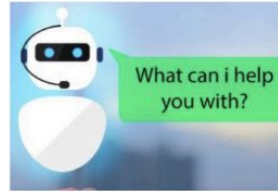


<https://blogs.nvidia.com/blog/2022/03/25/what-is-a-transformer-model/>

Real World Impact



Machine Translation



Smart Assistants



Search Engines



Auto Transcription



Health Record Analysis



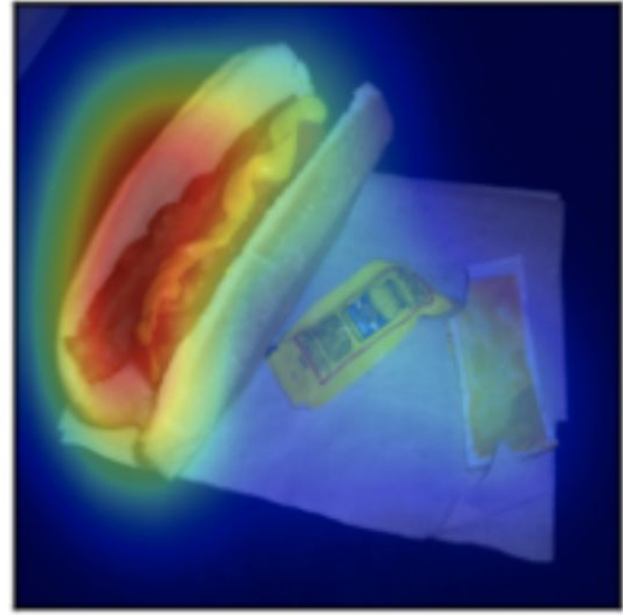
Summarization Engines

and many more

Questions?

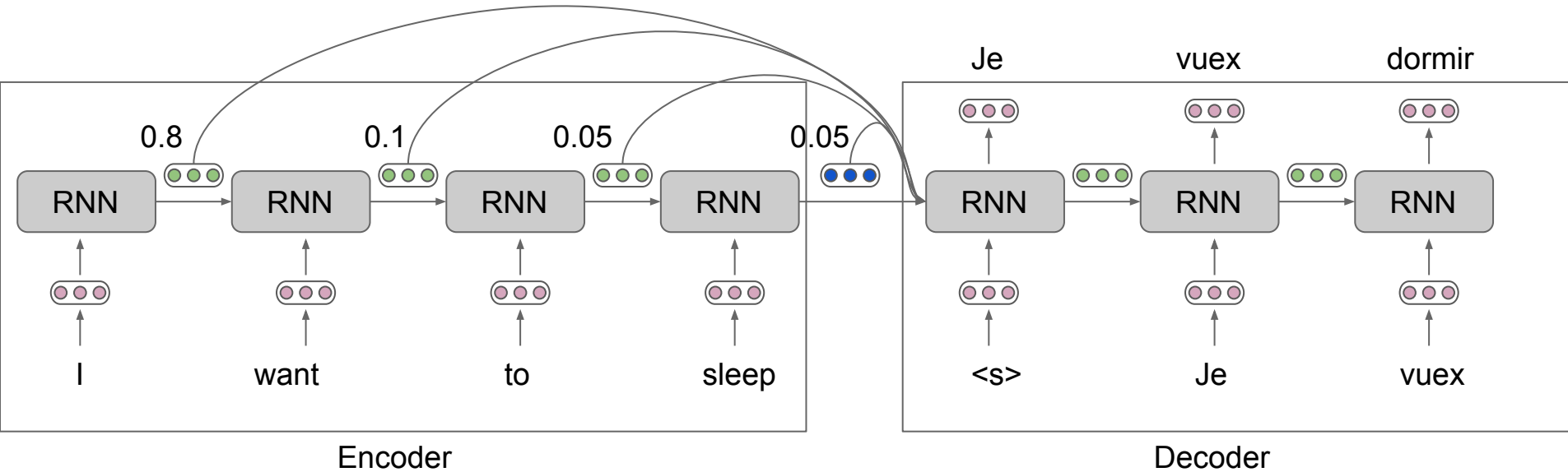
Visual Attention

What toppings are on the hot dog?

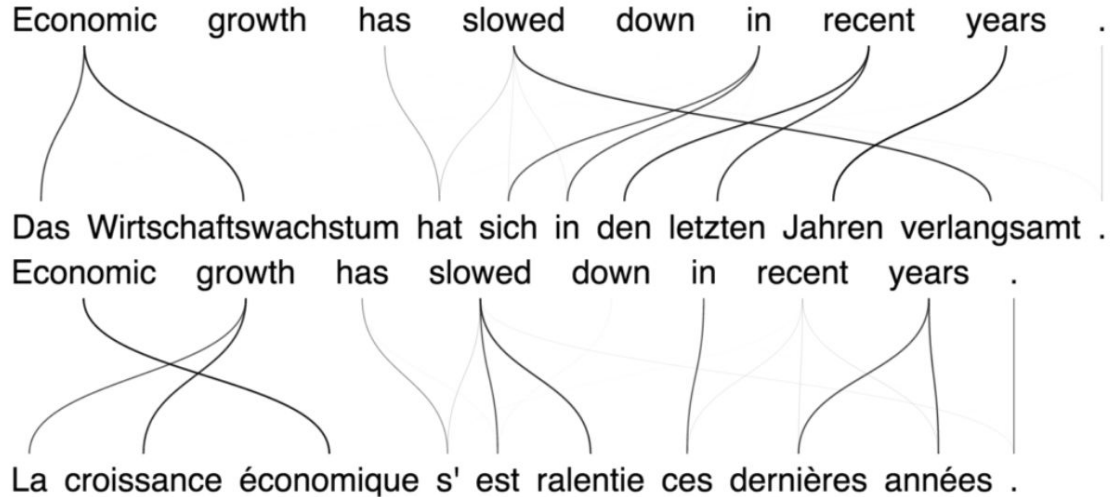
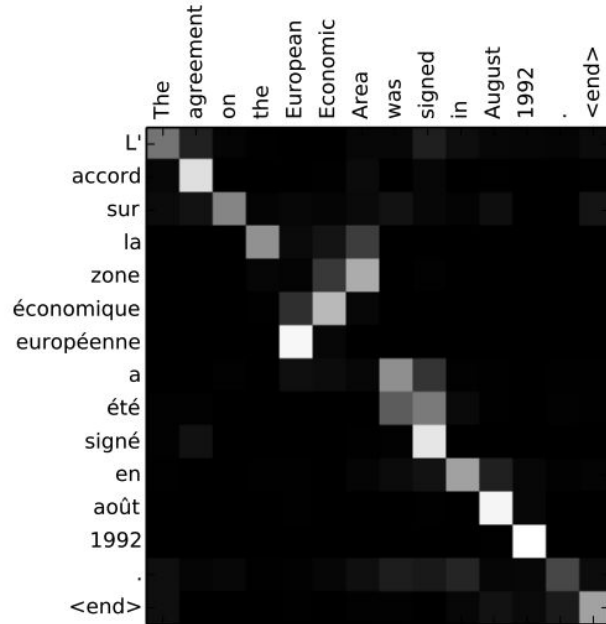


Differential Attention for Visual Question Answering (Patro et.al, 2018)

Cross Attention in NMT

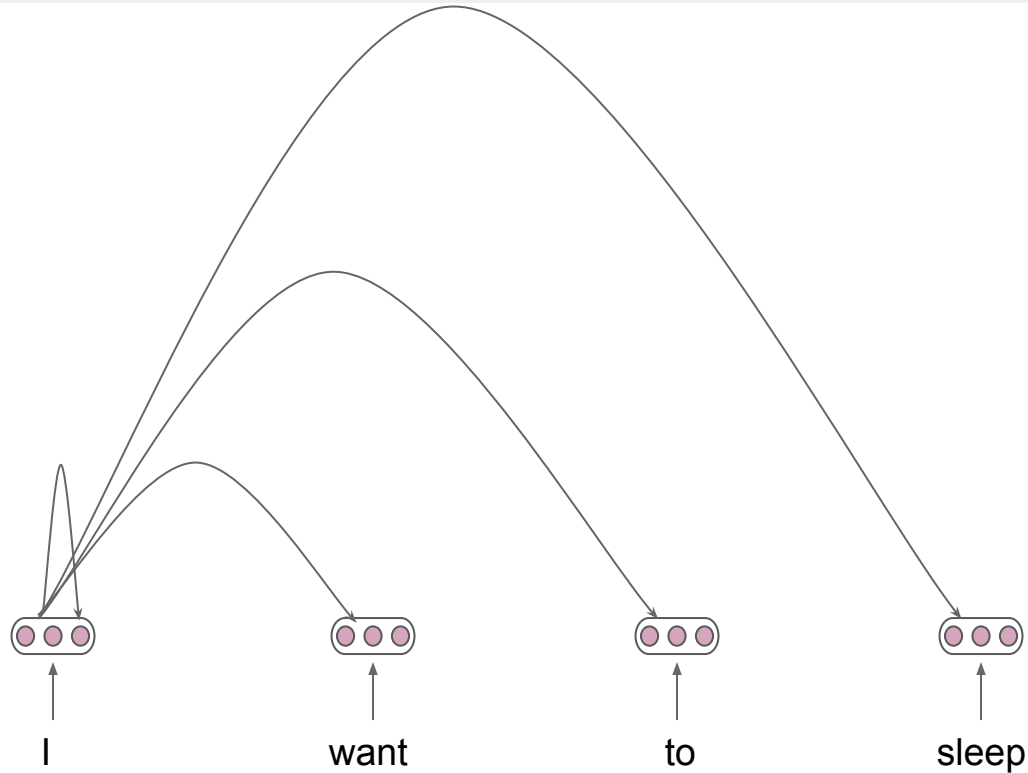


Attention in NMT

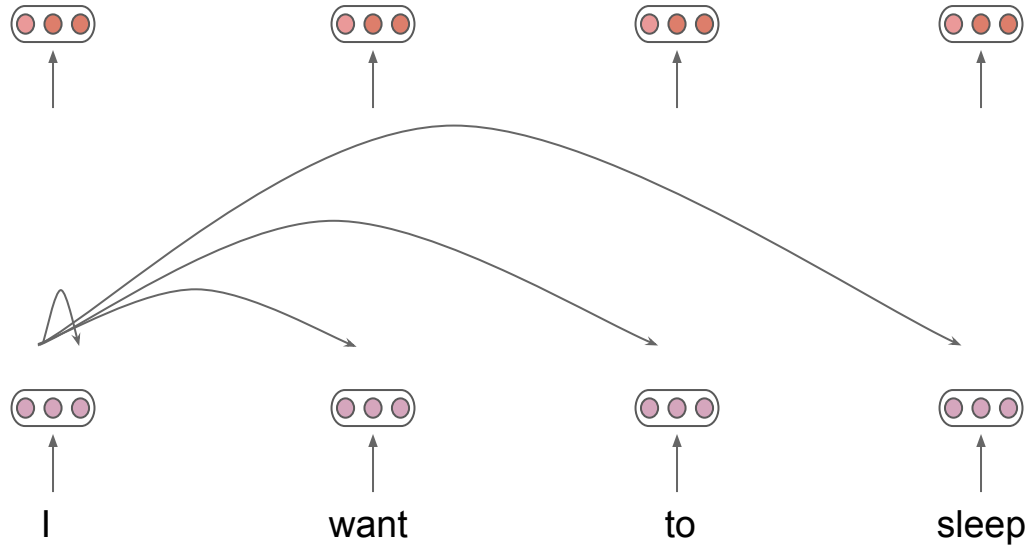


Neural Machine Translation by Jointly Learning to Align and Translate. Bahdanau et al, 2015
<https://developer.nvidia.com/blog/introduction-neural-machine-translation-gpus-part-3/>

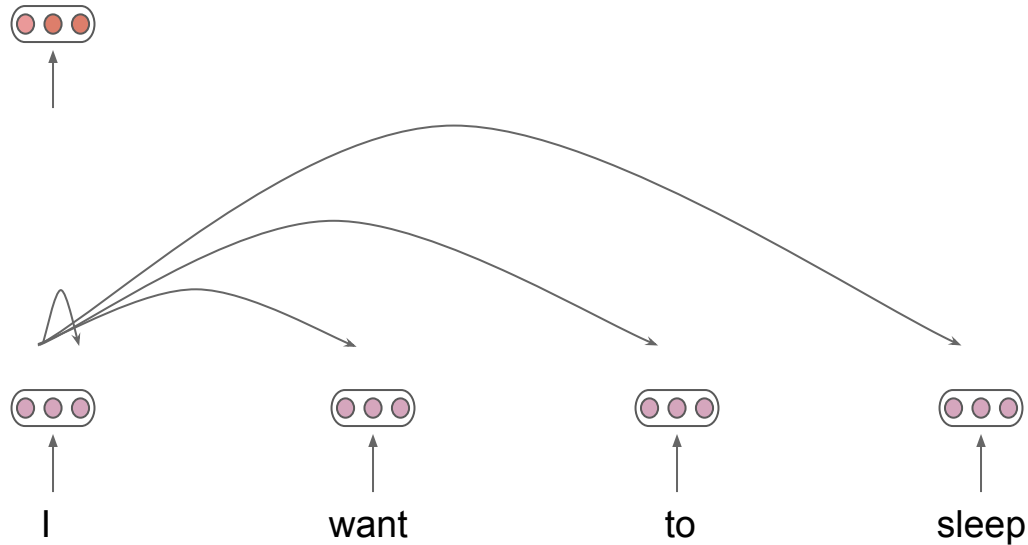
Self Attention



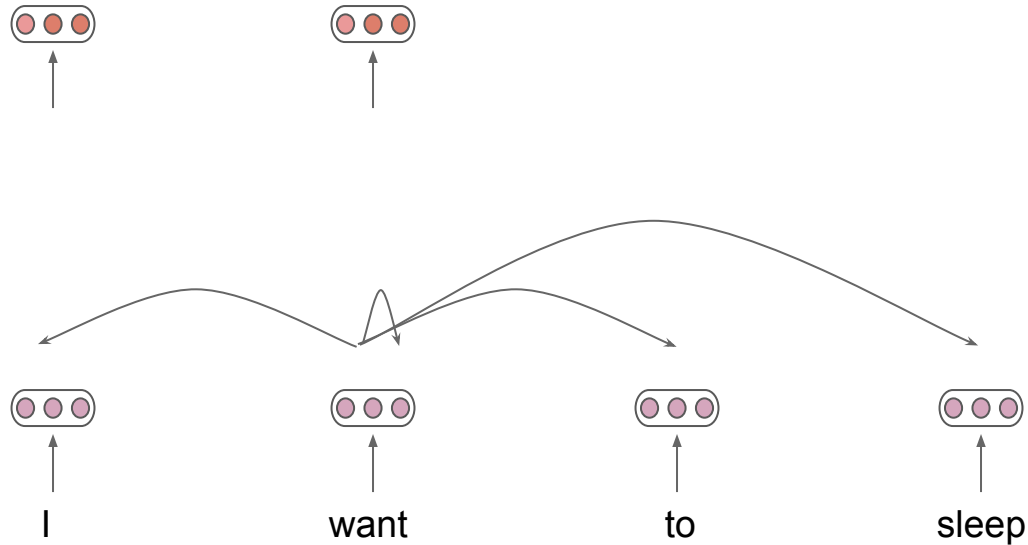
Self Attention - No more recurrence



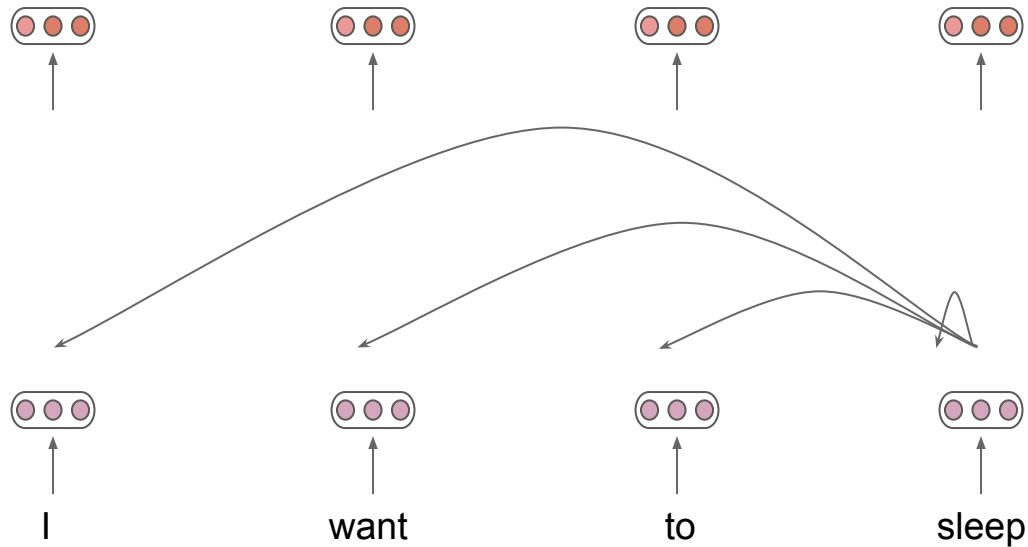
Self Attention - No more recurrence



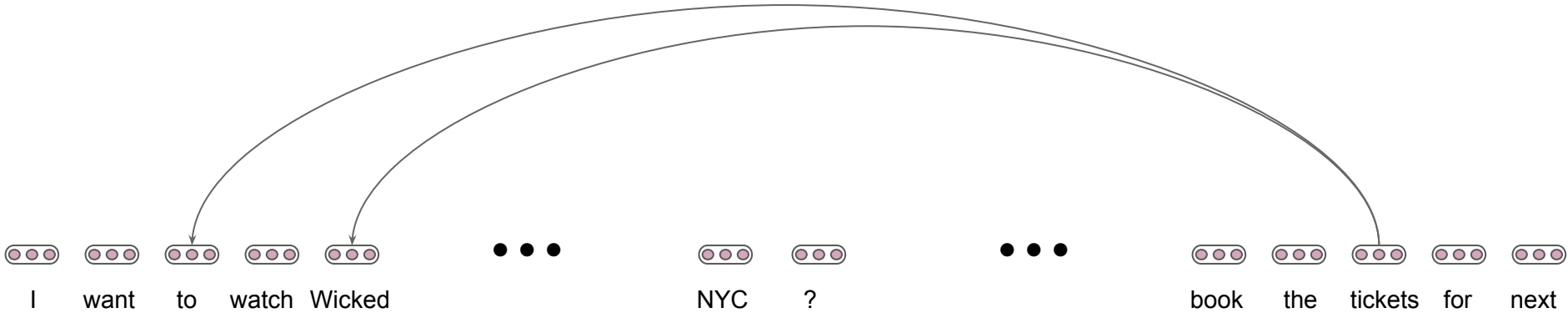
Self Attention - No more recurrence



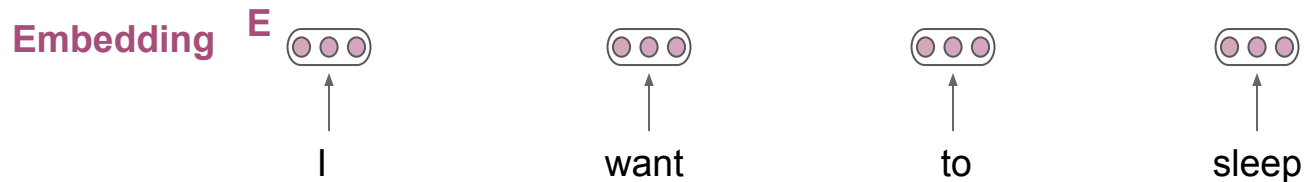
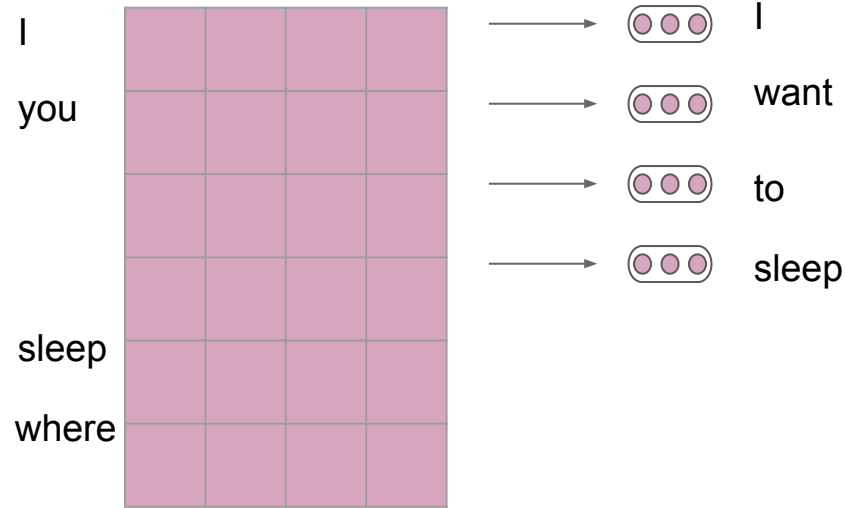
Self Attention - No more recurrence



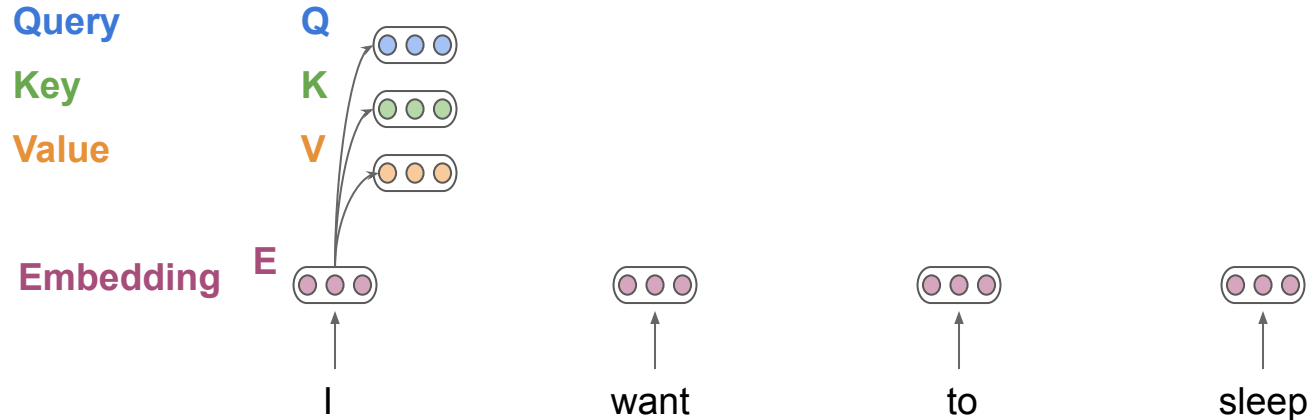
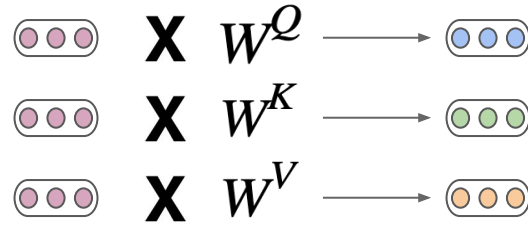
Self Attention for long sequences



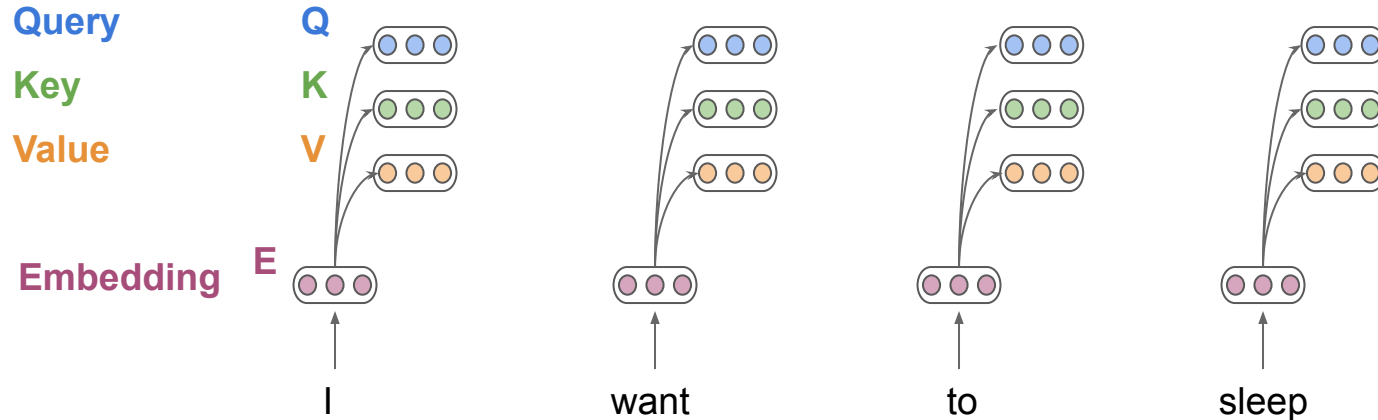
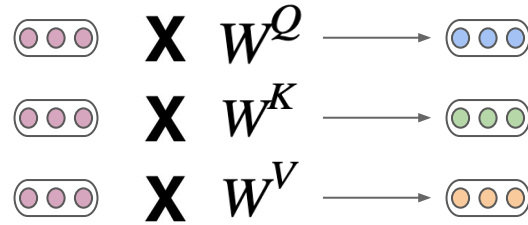
Self Attention - Word Embedding



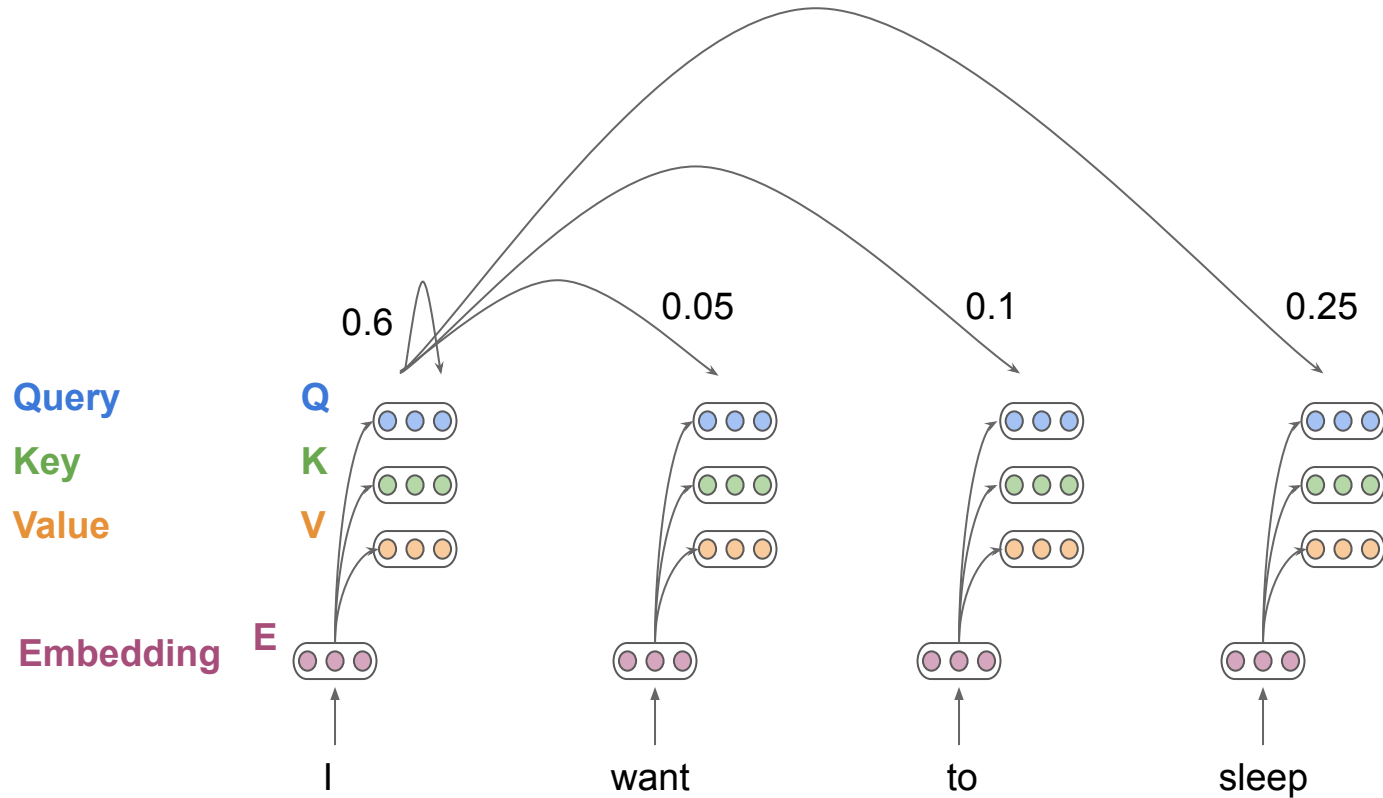
Self Attention - Projection Layer



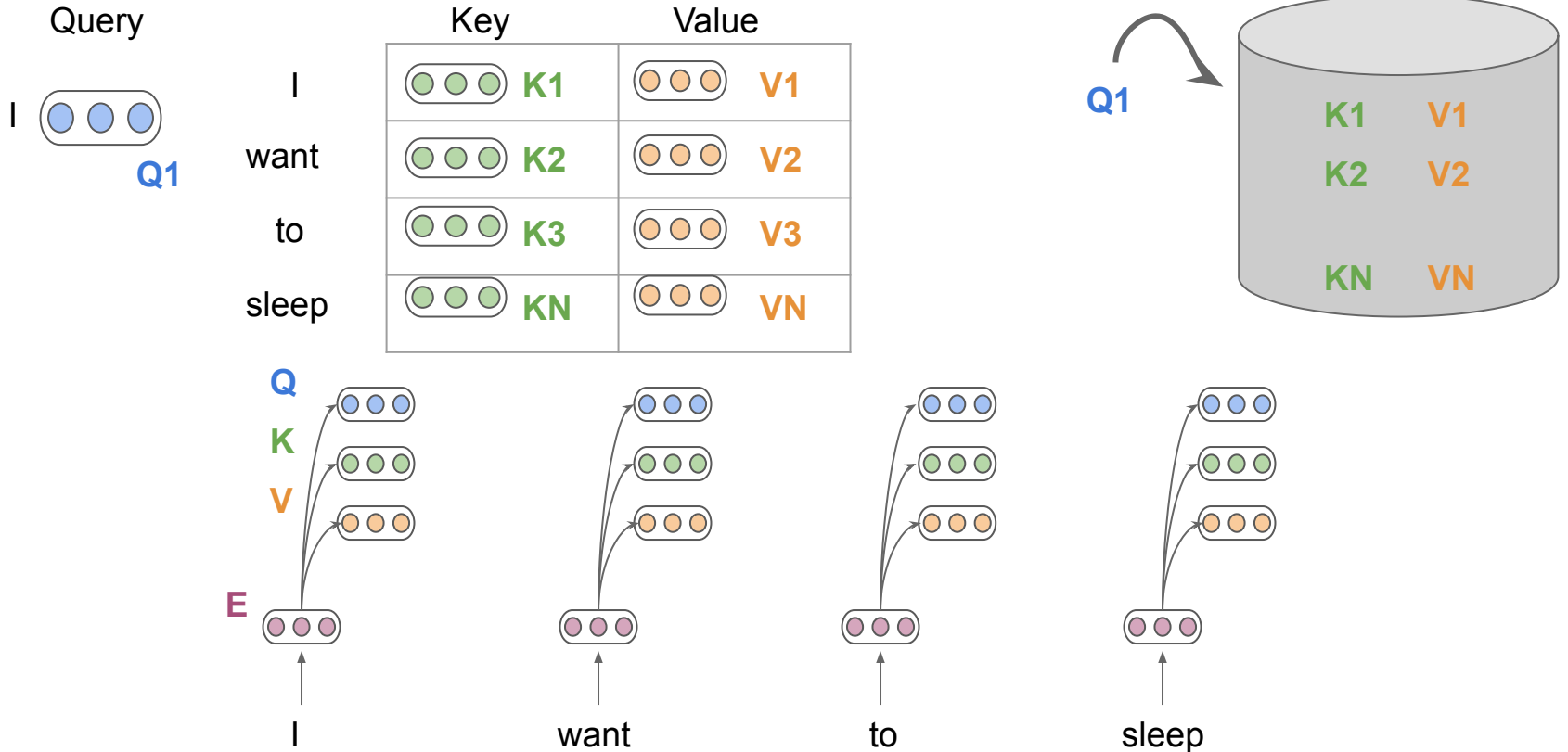
Self Attention - Projection Layer



Self Attention - Attention Scores

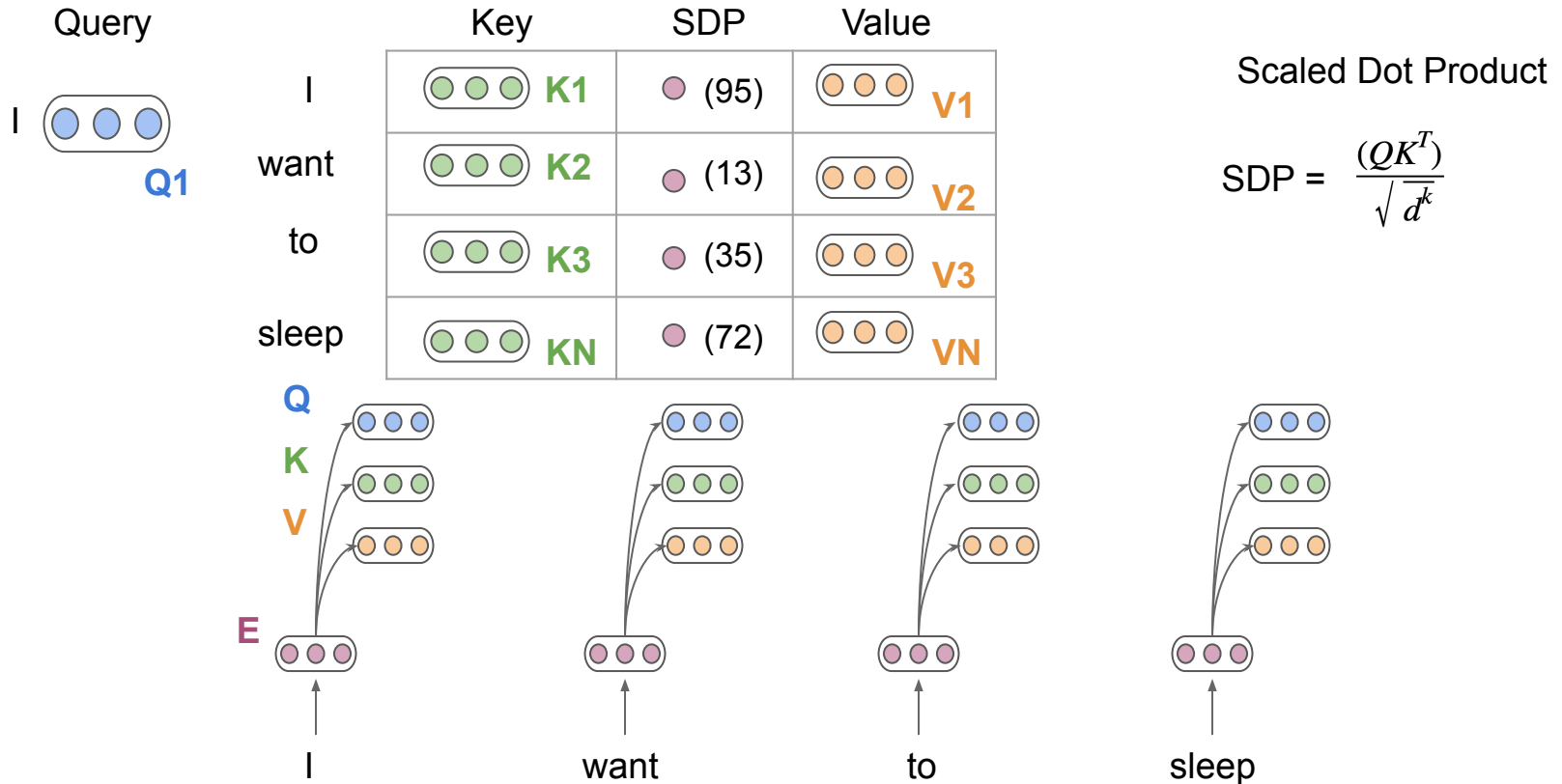


Self Attention



Questions?

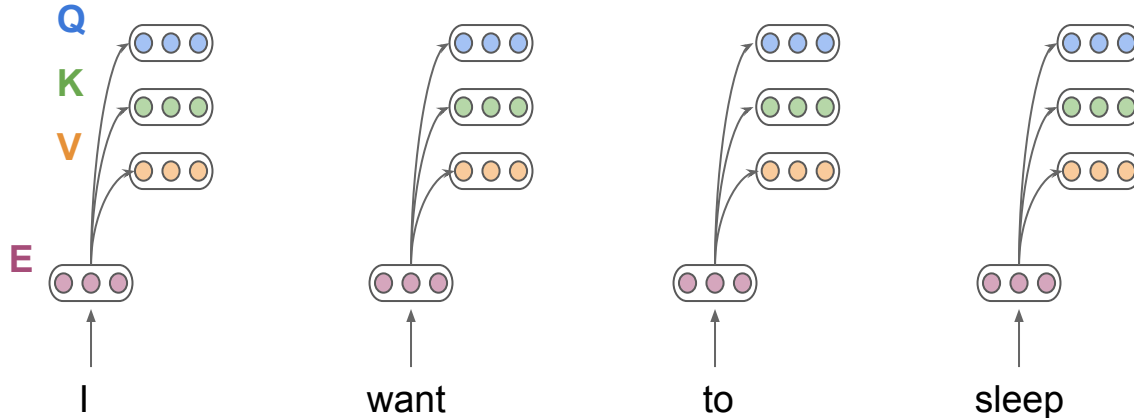
Self Attention - Scaled Dot Product



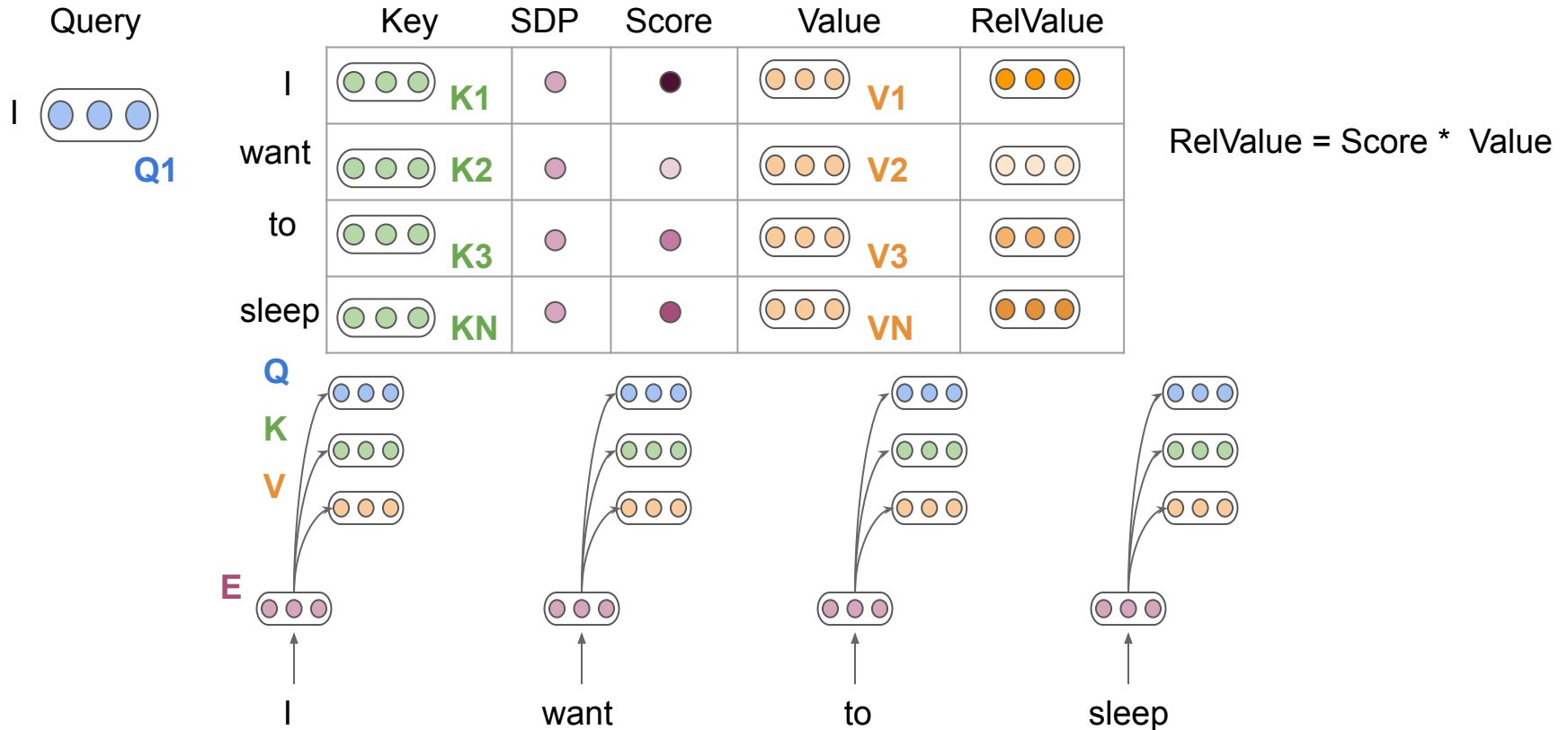
Self Attention - SoftMax

| 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 |

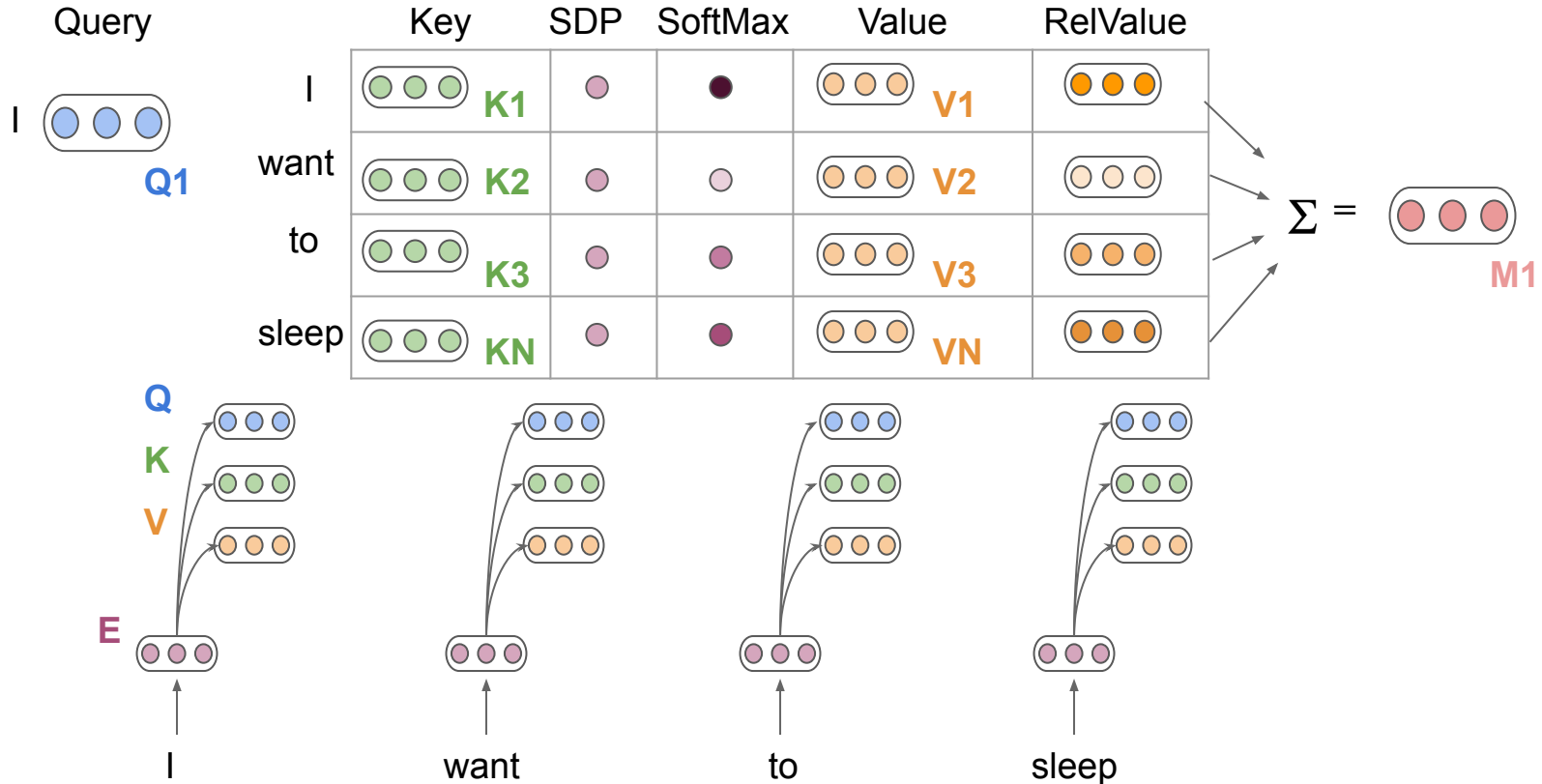
$$\text{score} = \text{softmax} \left(\frac{(QK^T)}{\sqrt{d^k}} \right)$$



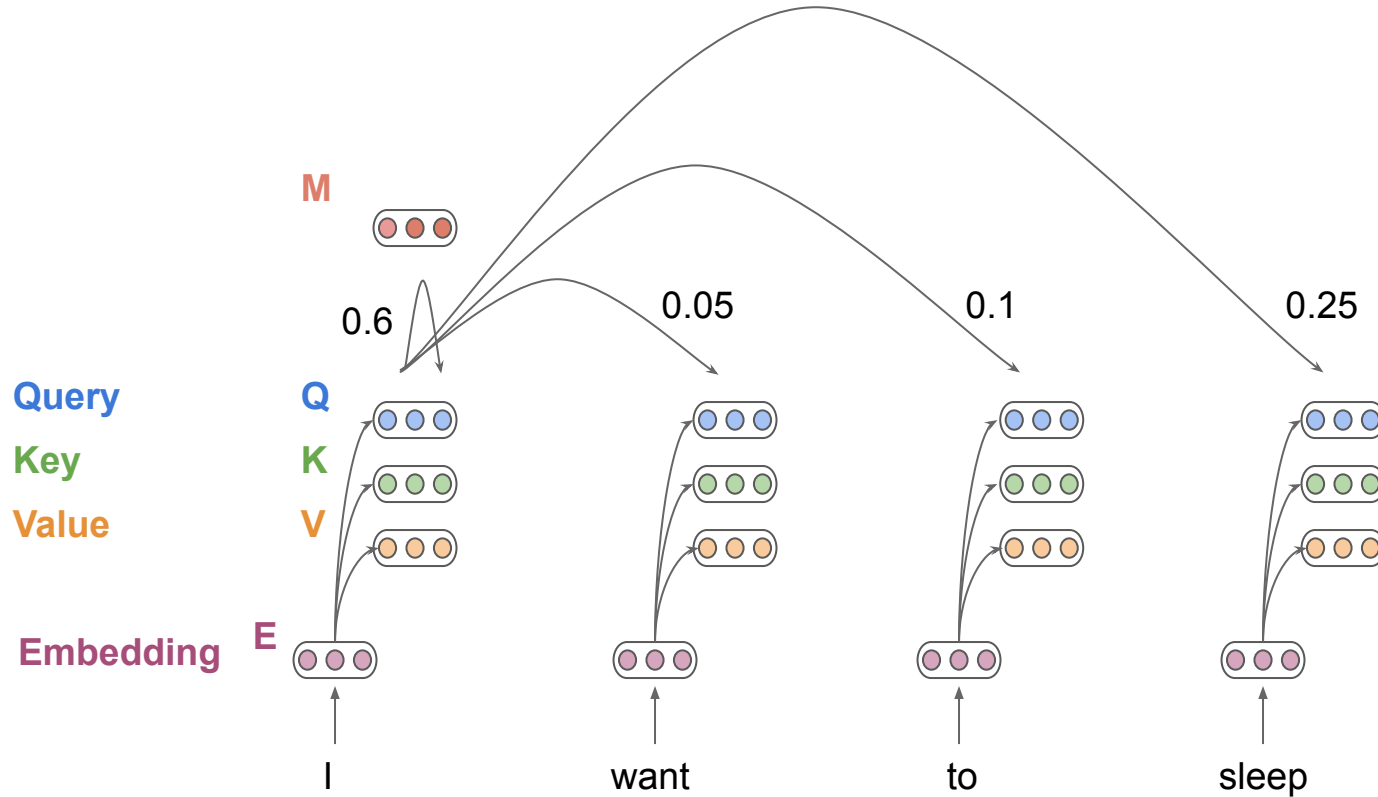
Self Attention - Soft (Relative) Values



Self Attention - Attended Repr

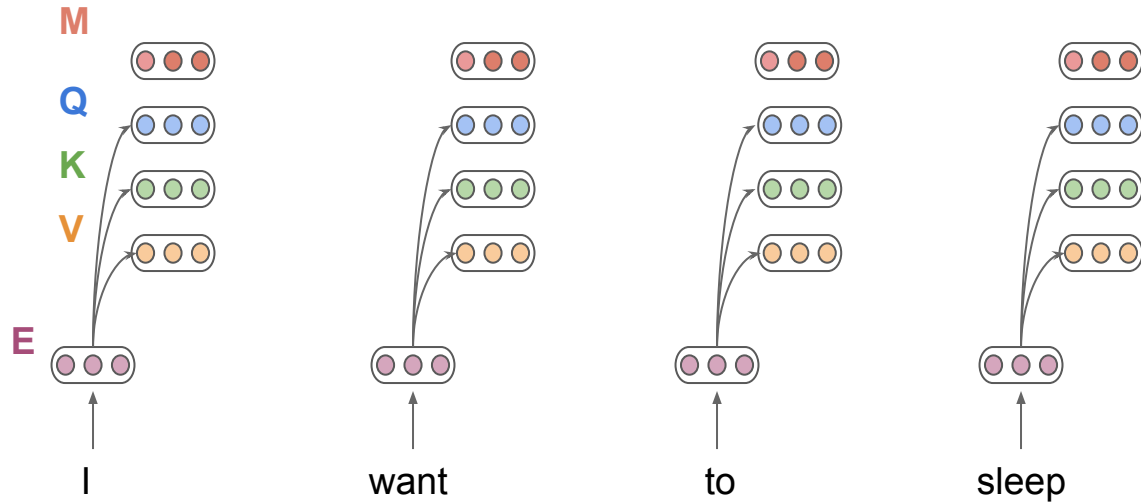


Self Attention - Attended Contextual Rep



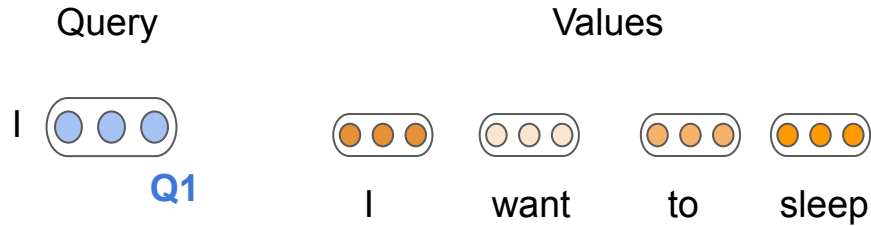
Questions?

Self Attention - Attended Contextual Rep



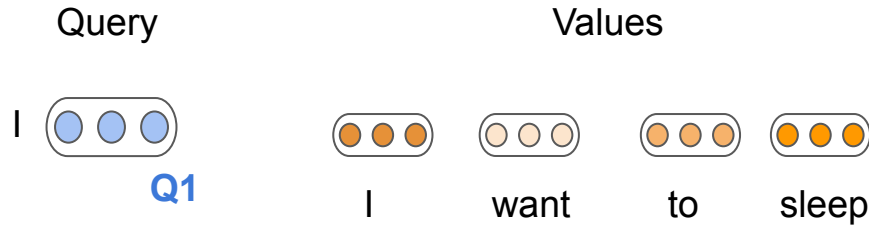
Problem with Self Attention

- Self Attention can focus heavily on the same word!



Problem with Self Attention

- Self Attention can focus heavily on the same word!



- Single representation

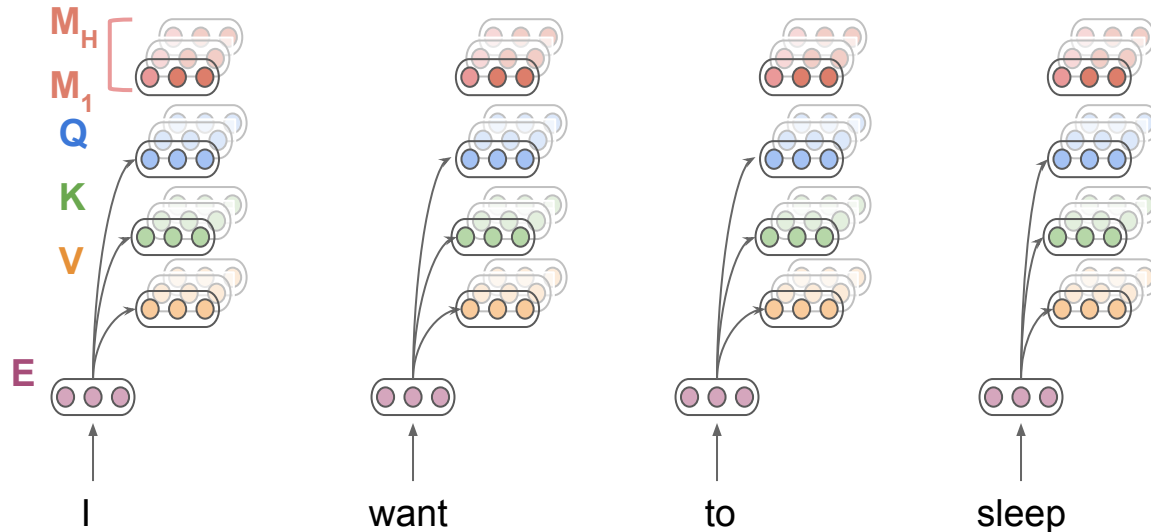
I like **Harry Potter**
(Book)

v/s

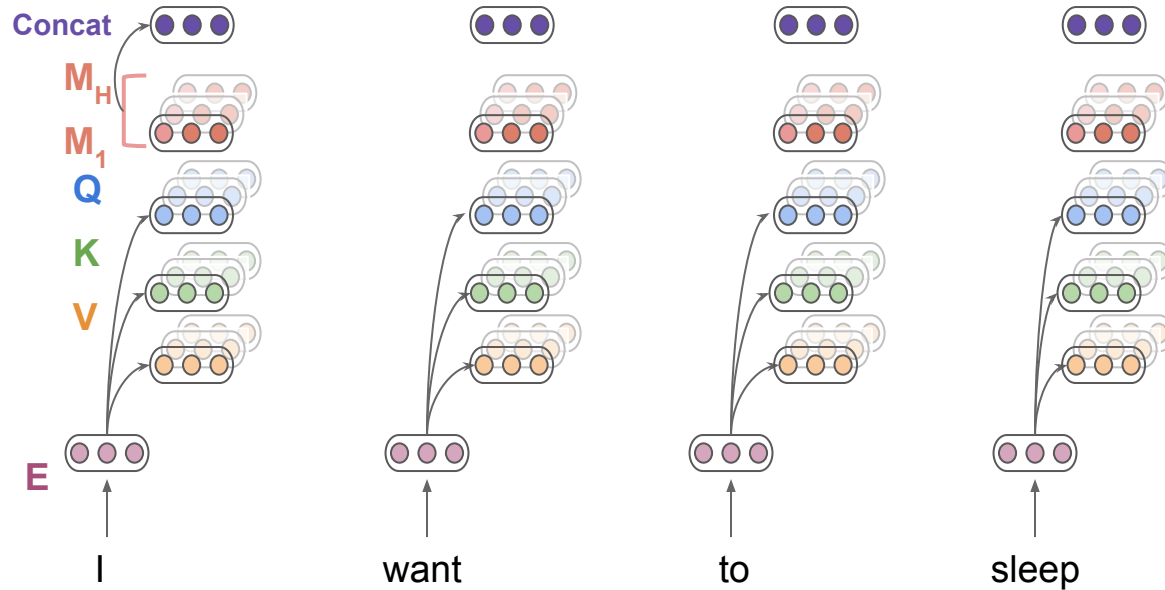
I like **Harry Potter**
(Movie)

Multi-Headed Self Attention

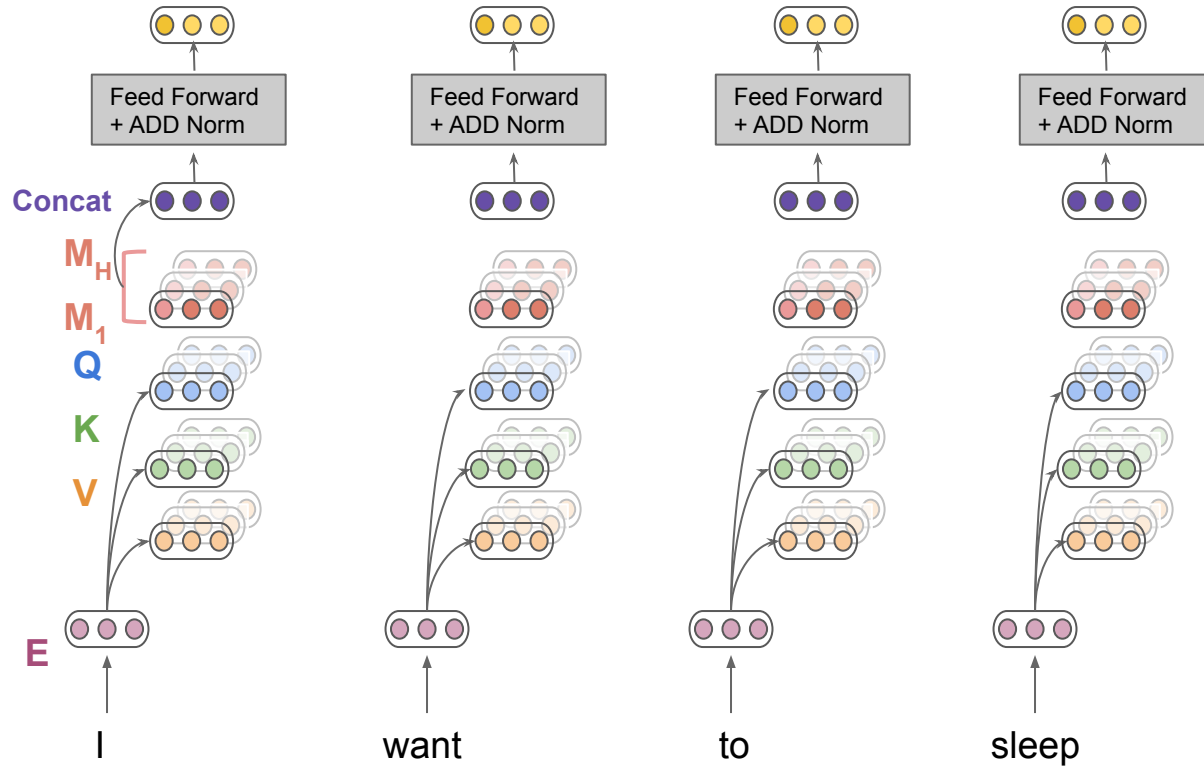
H (no: of heads) Different versions of Q,K,V
Each different repr -> Different attended repr



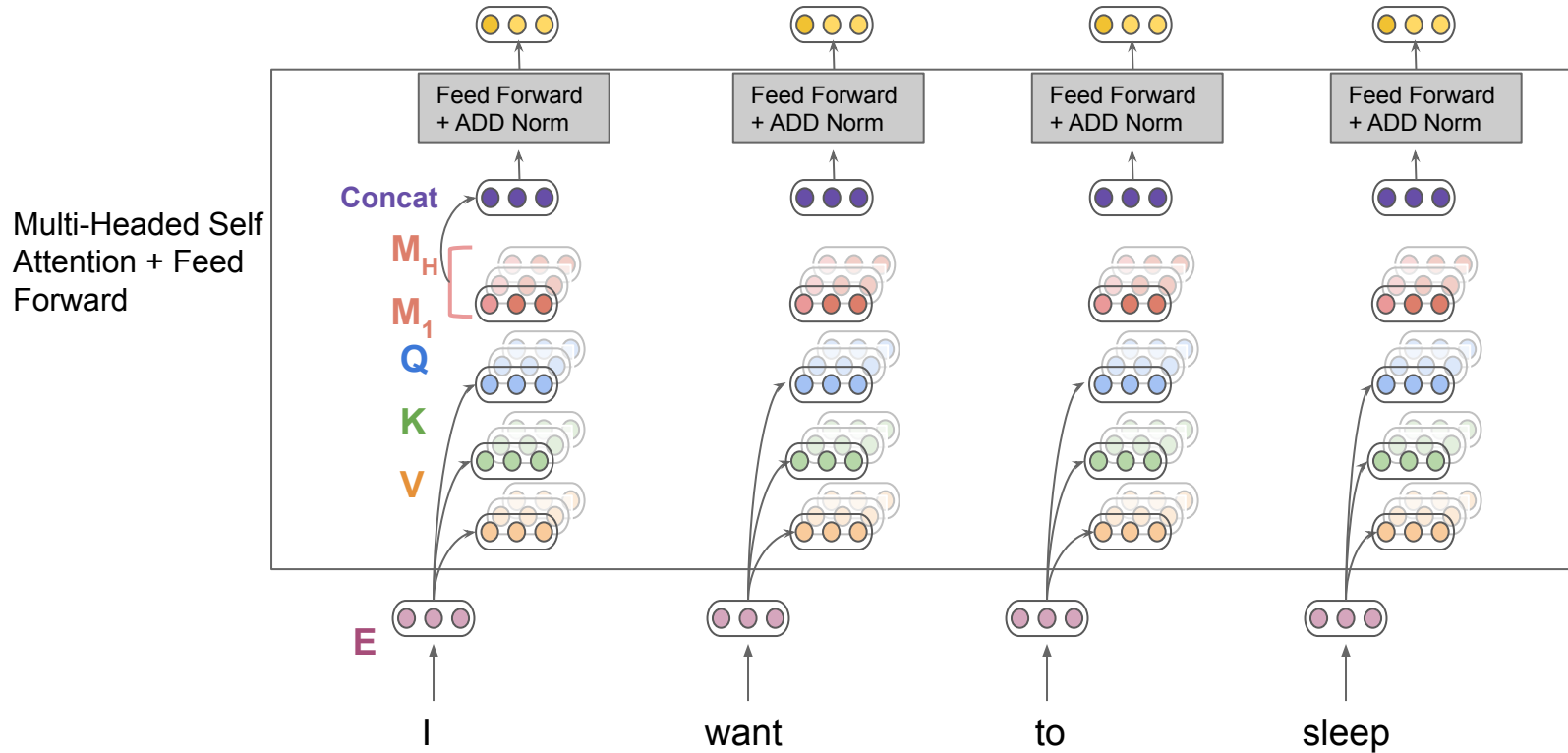
Multi-Headed Self Attention



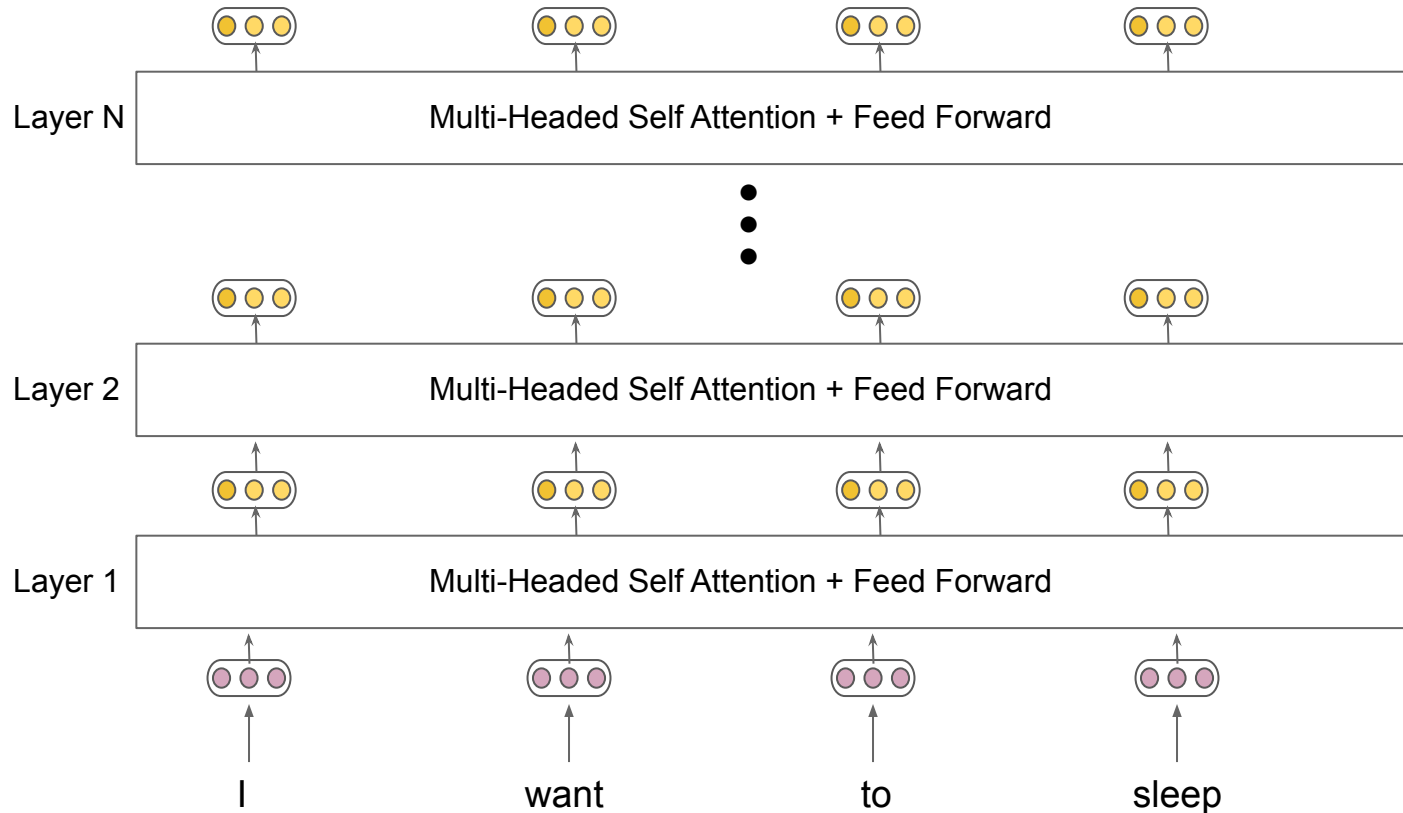
Multi-Headed Self Attention



Multi-Headed Self Attention



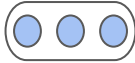
Multi-Headed Self Attention



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 | |

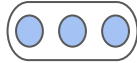
$$\Sigma = \text{capsule with 3 red circles} M$$

I want to sleep

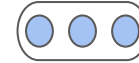
Sleep to I want

Revisiting Self Attention

Query (I)



Query (I)



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

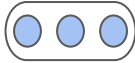
$$\Sigma = \text{capsule with 3 red circles} M$$

I want to sleep

Sleep to I want

Revisiting Self Attention

Query (I)

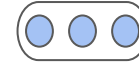


| | Key | SDP | SM | Value | RelValue |
|-------|-----|-----|----|-------|----------|
| I | K1 | | | V1 | |
| want | K2 | | | V2 | |
| to | K3 | | | V3 | |
| sleep | KN | | | VN | |

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

I want to sleep

Query (I)

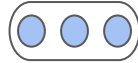


| | Key | SDP | SM | Value | RelValue |
|-------|-----|-----|----|-------|----------|
| Sleep | K1 | | | VN | |
| to | K2 | | | V3 | |
| I | K3 | | | V1 | |
| want | KN | | | V2 | |

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!

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

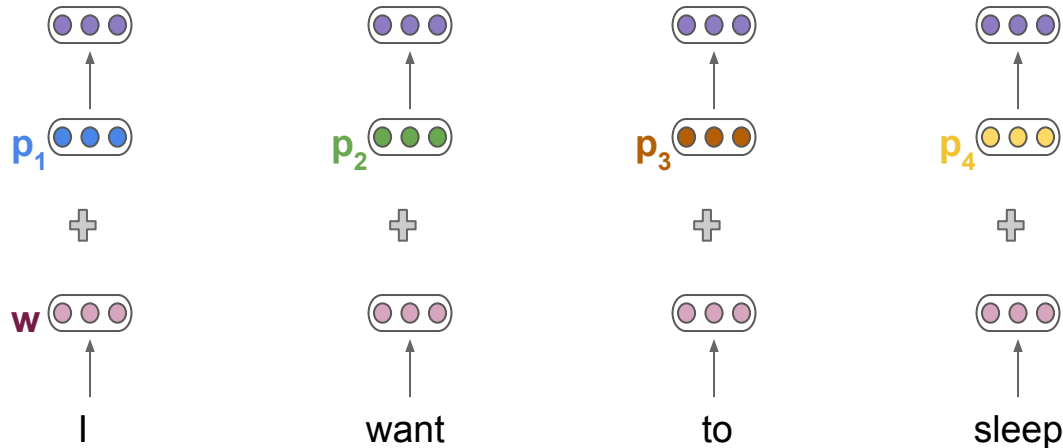
I want to sleep

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

Sleep to I want

Positional Encoding

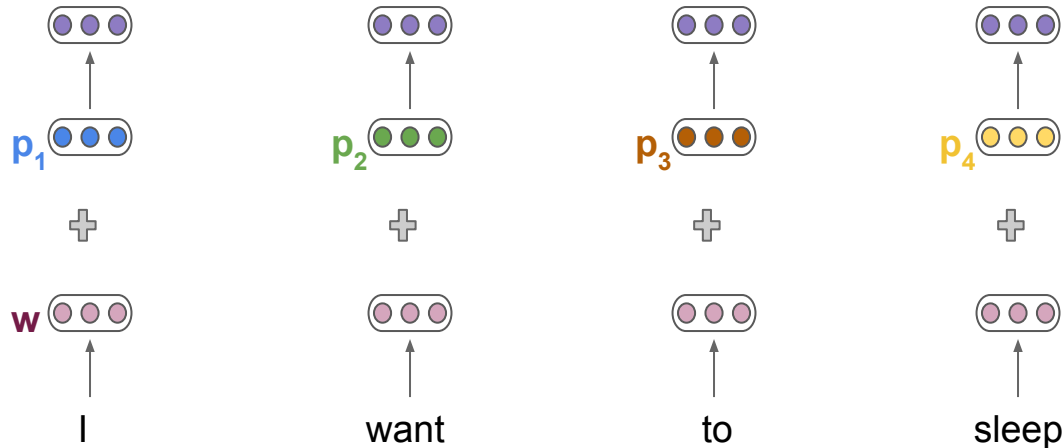
Position embeddings - each position number has an associated embedding



Positional Encoding

Sinusoidal Position embeddings - generalize to any sequence length

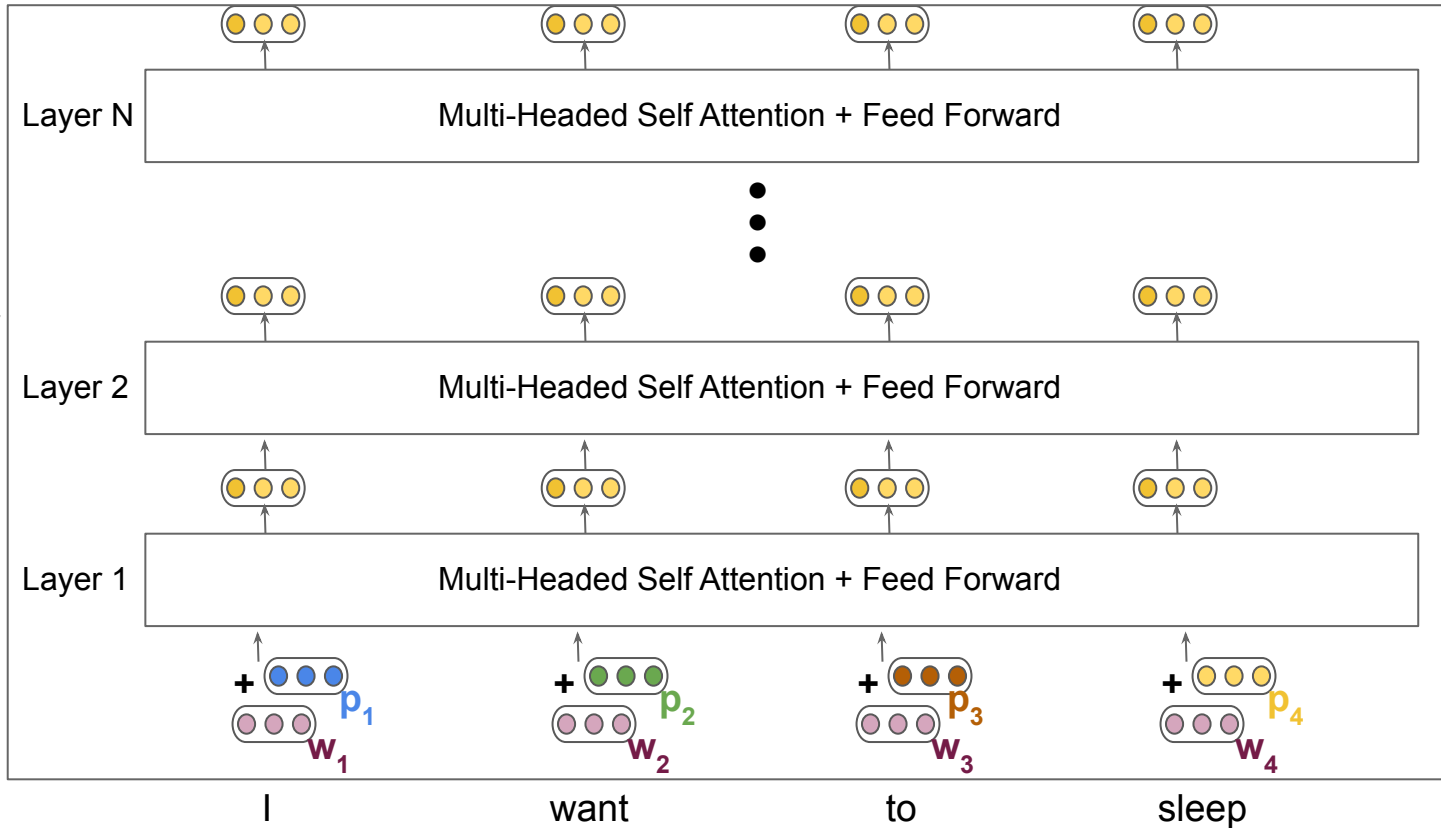
$$p = f(i, t)$$



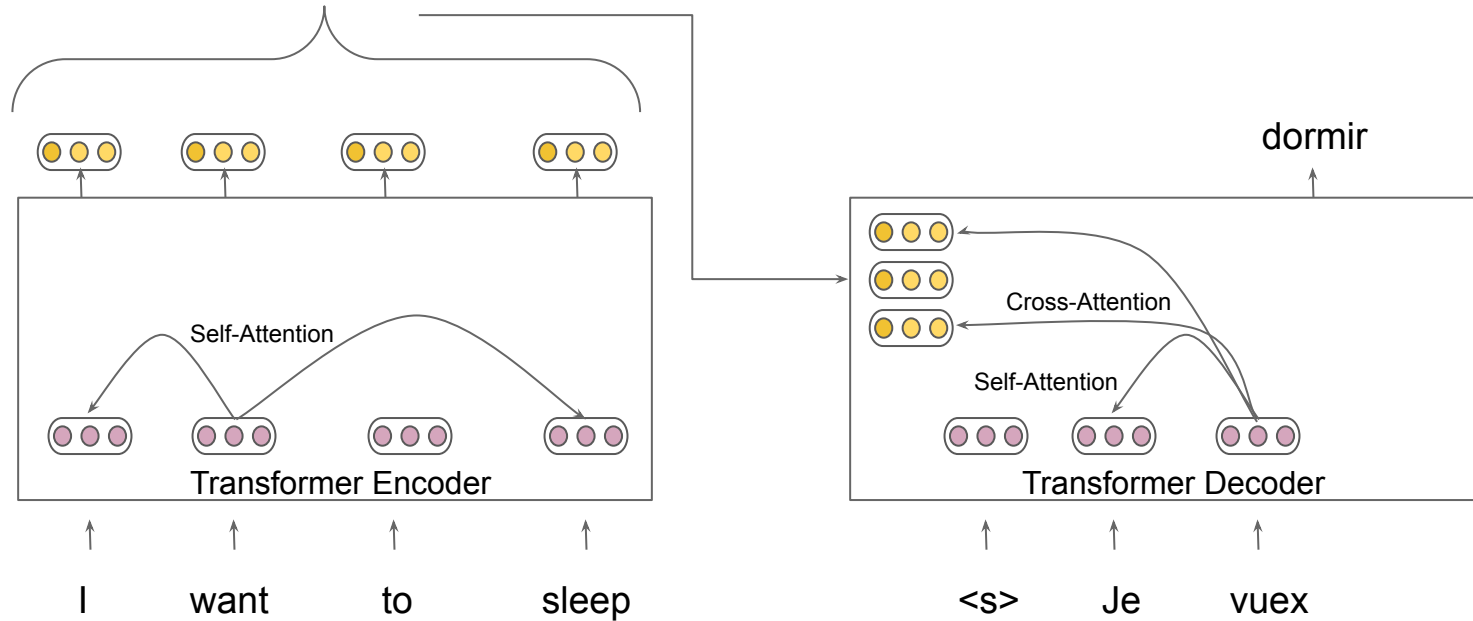
Questions?

Transformer Encoder

N-Layer
Transformer
Encoder



Transformer Encoder - Decoder

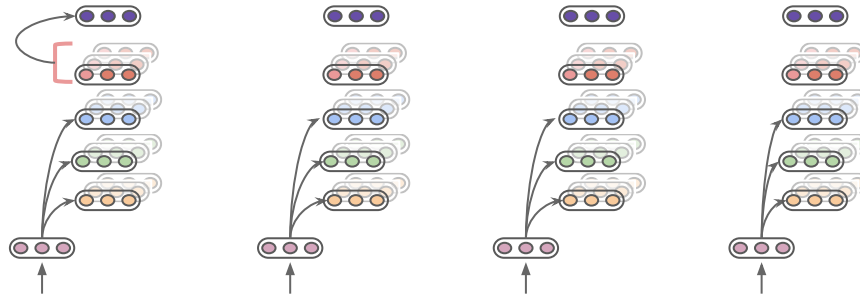


What's so great about Transformers?

- 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

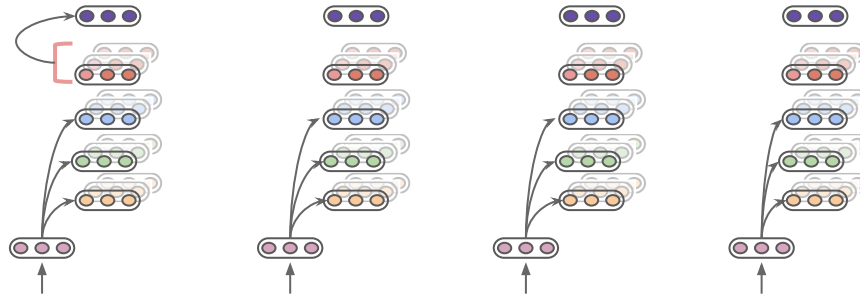
What's so great about Transformers?

- Parallelizable computation - Entire sequence can be processed in parallel

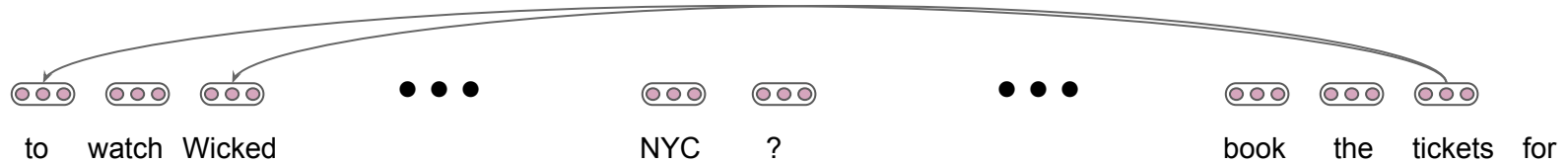


What's so great about Transformers?

- Parallelizable computation - Entire sequence can be processed in parallel

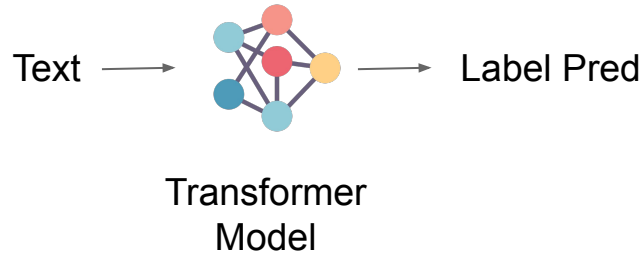


- Rich expressive power - long range dependencies

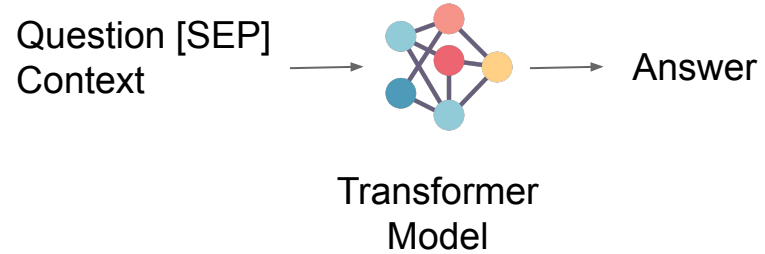


Impact - Wide Applications!

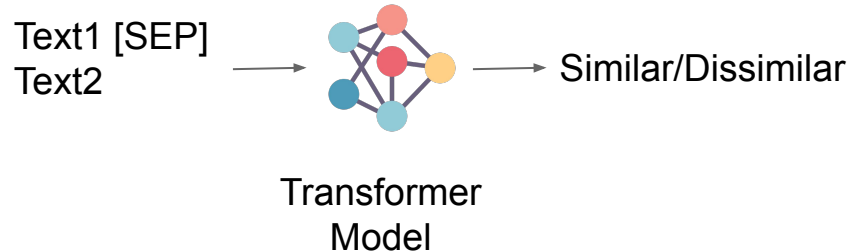
Classification



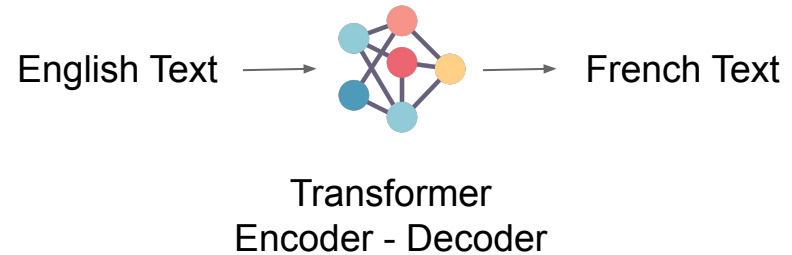
Question Answering



Sentence Similarity



Translation

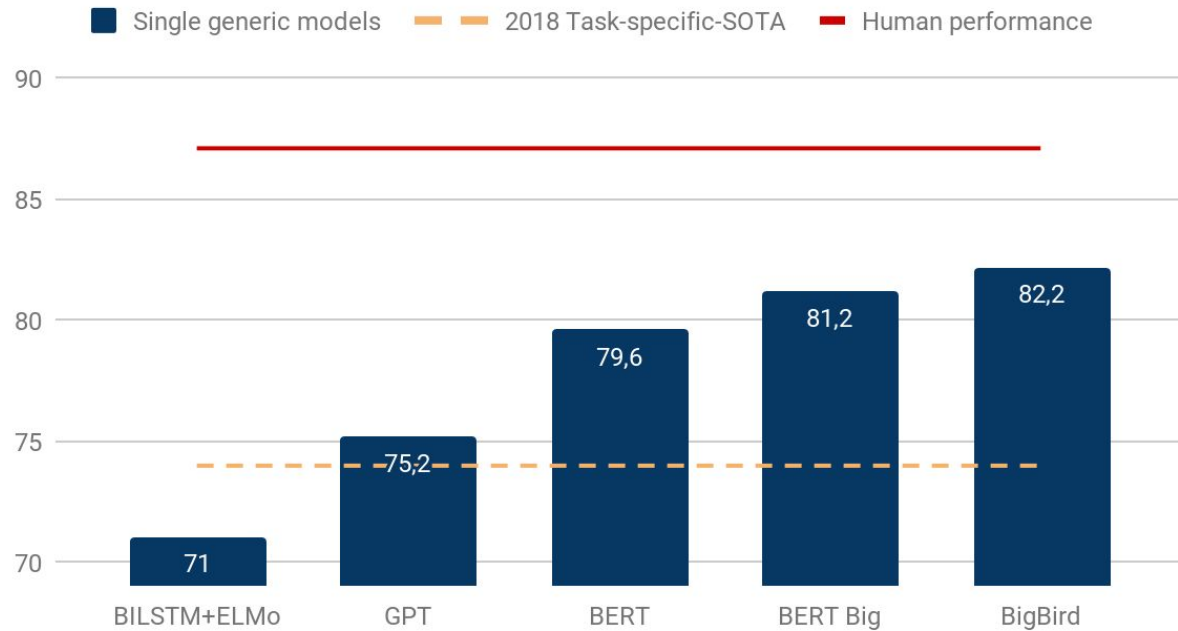


Larger Impact

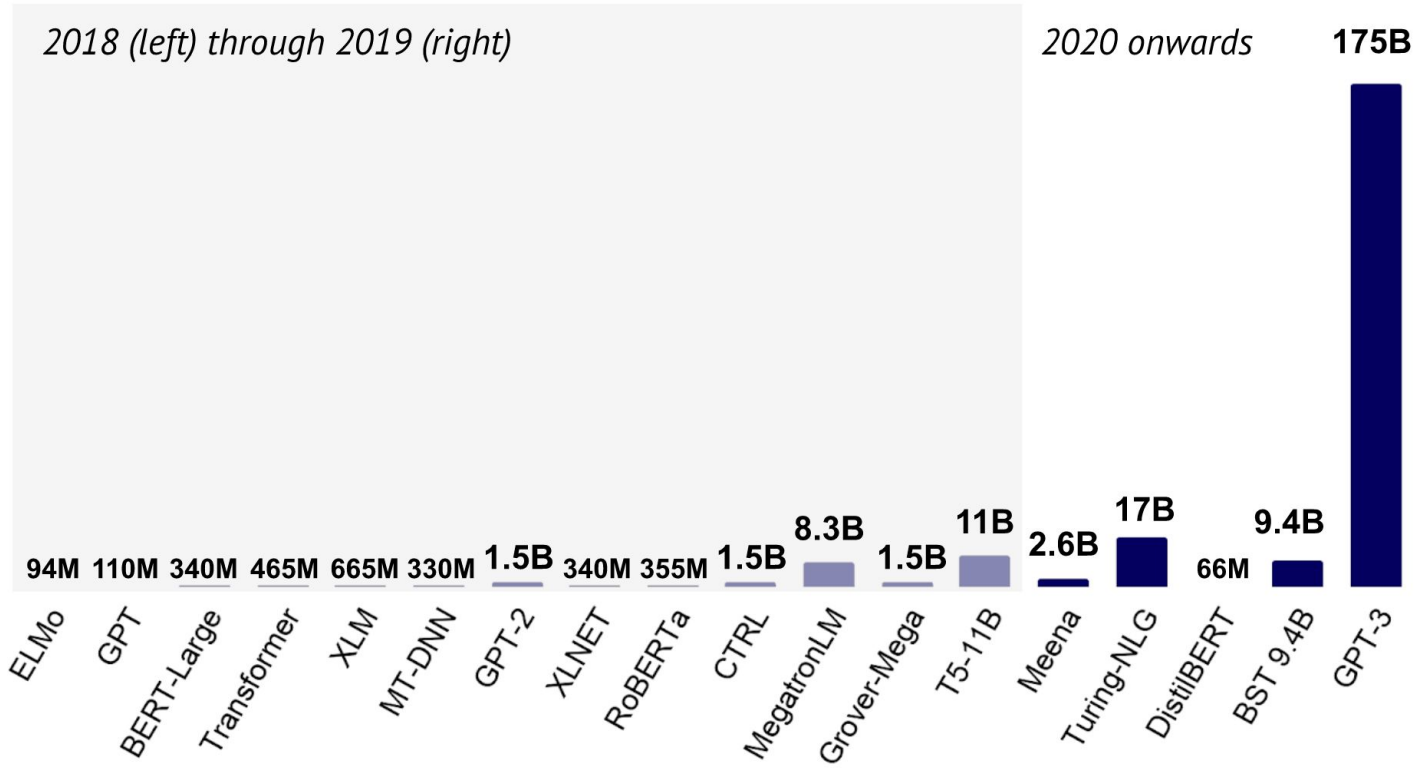


Larger Impact

GLUE scores evolution over 2018-2019



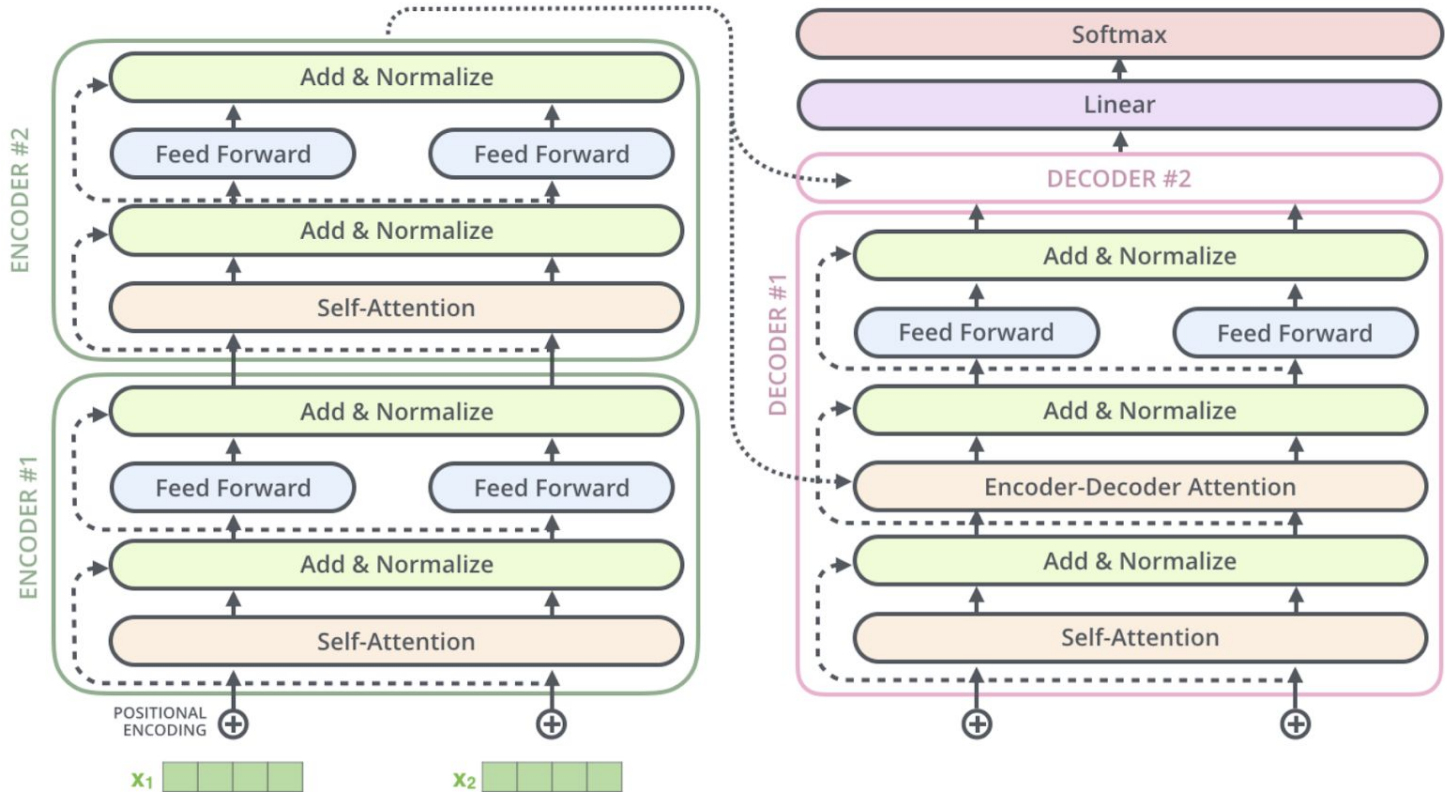
Larger Impact



Thank you!

vbalacha@cs.cmu.edu

Transformer Encoder-Decoder



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