

Self Attention and Transformers

Vidhisha Balachandran

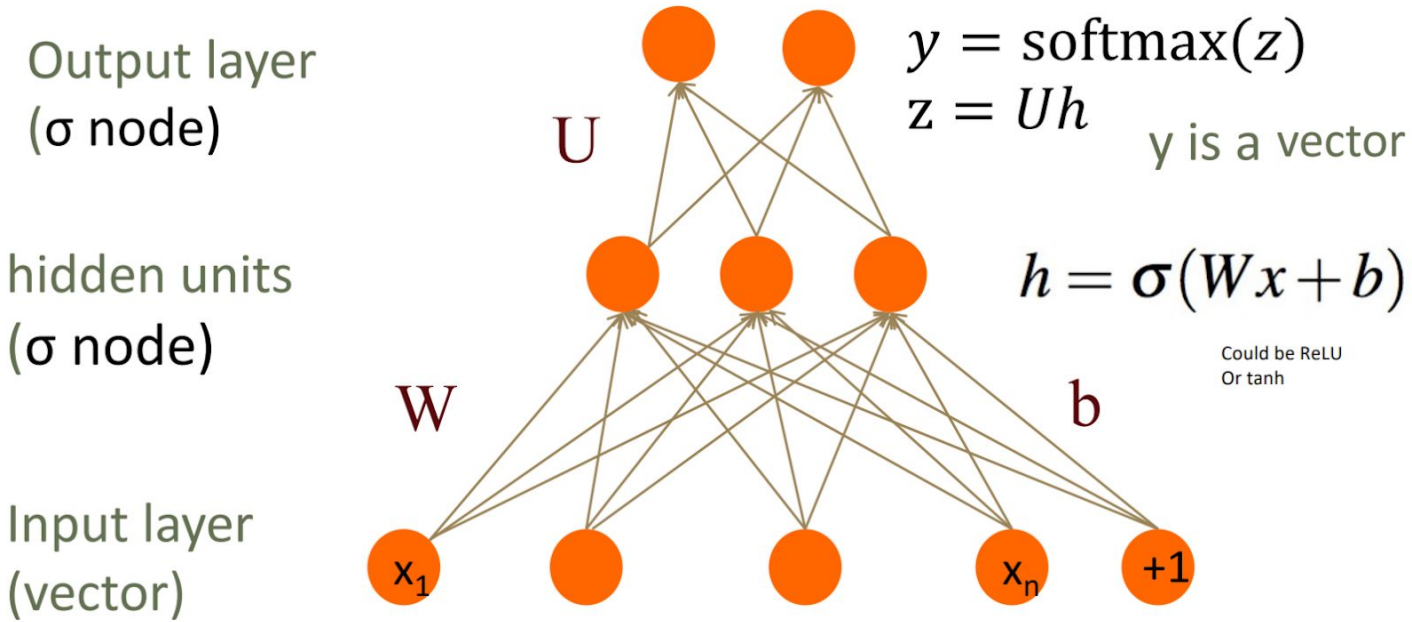
vidhishab@microsoft.com

Inspired by slides from Emma Strubell, Graham Neubig

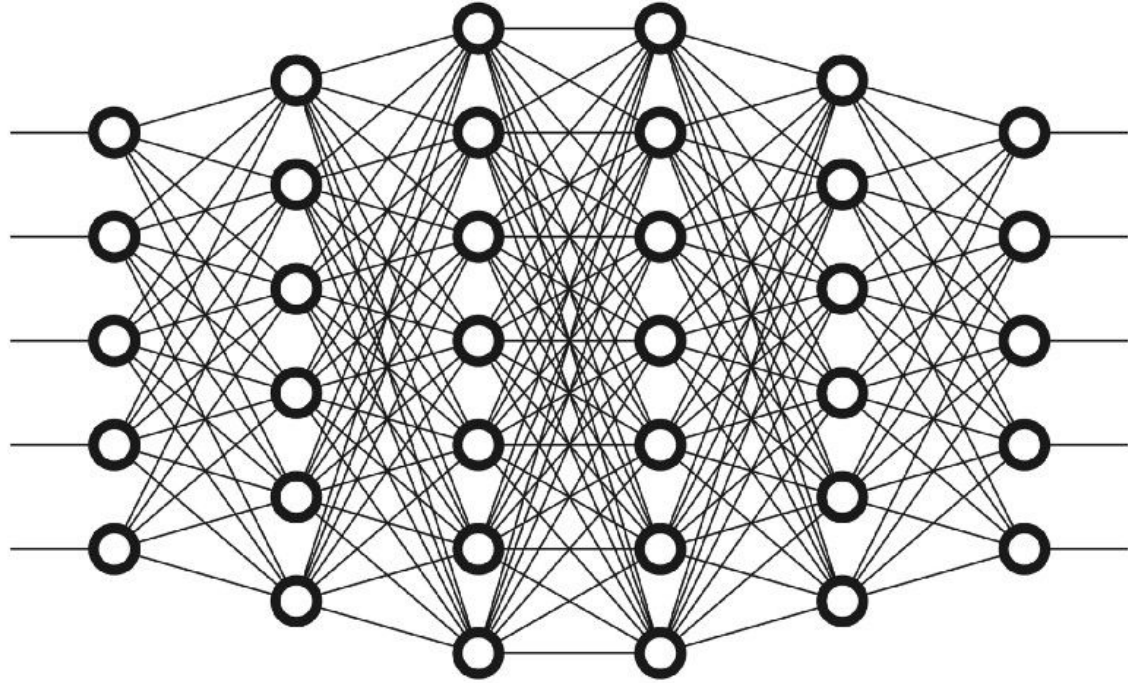
Readings

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

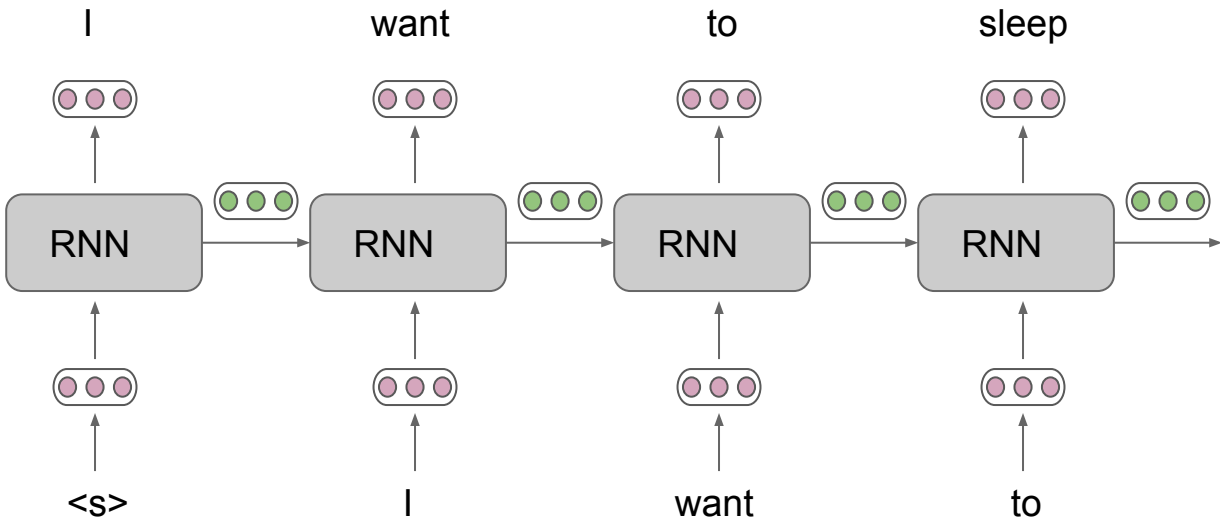
Recap - 2 Layer MLP



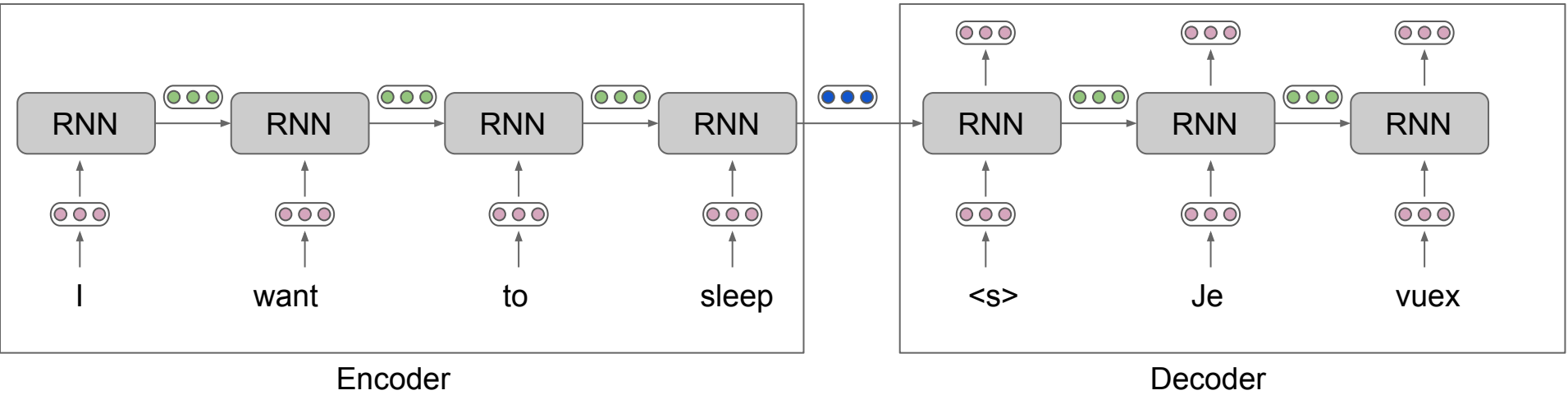
Deep MLP



Recurrent Neural Networks - RNNs



Encoder-Decoder Models



Limitations

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

Long Range Dependencies

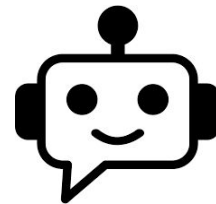


I'm want to watch Wicked! How does the weather in NYC look next week?

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?

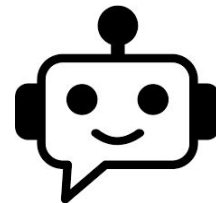
It looks sunny with some light rain during the weekend.

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?

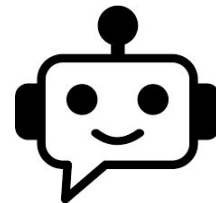
It looks sunny with some light rain during the weekend.

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

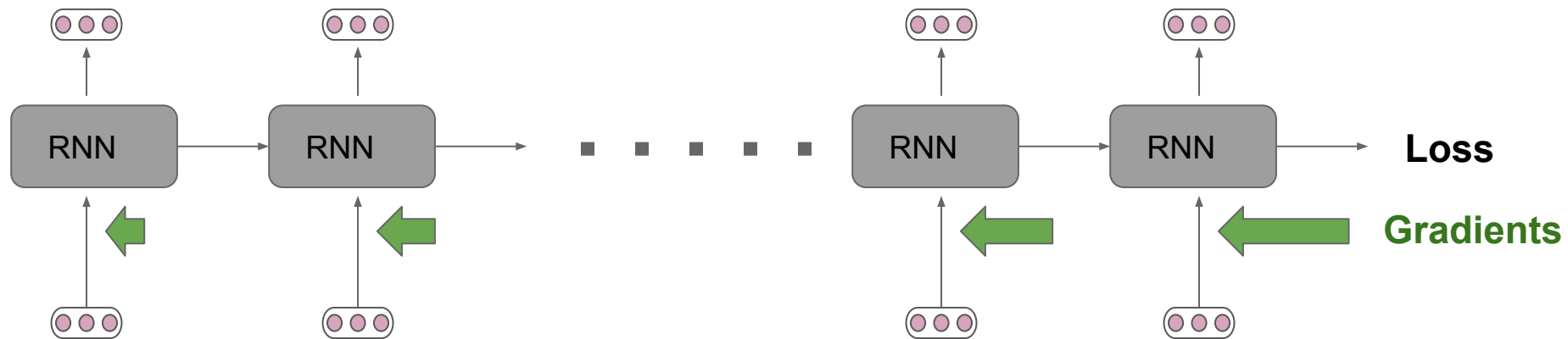
There's a marshall's a mile away. They have the navy blue jacket you have been eyeing for a while!

⋮

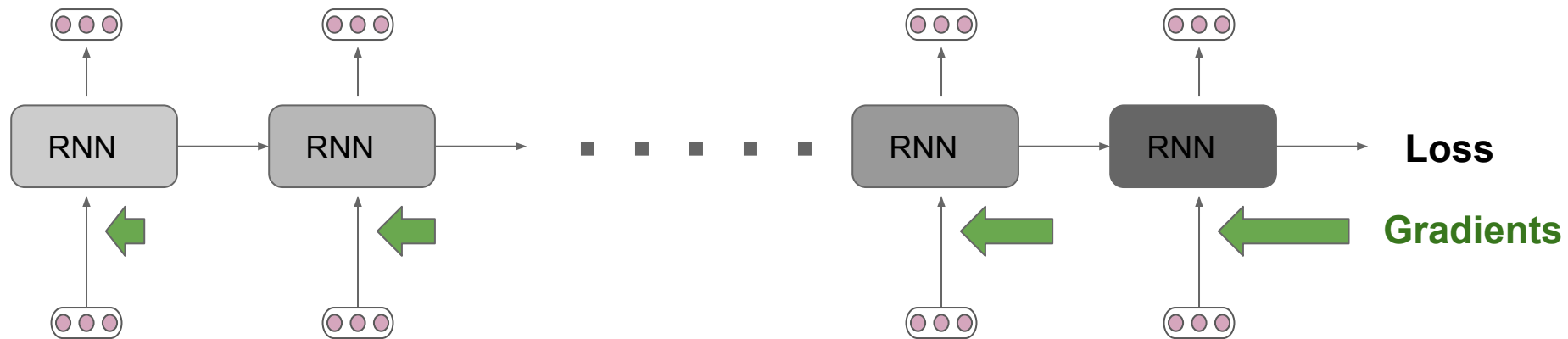
Ok! Looks like I can actually go! Book **the tickets** for next Wed!



Gradient vanishing / explosion



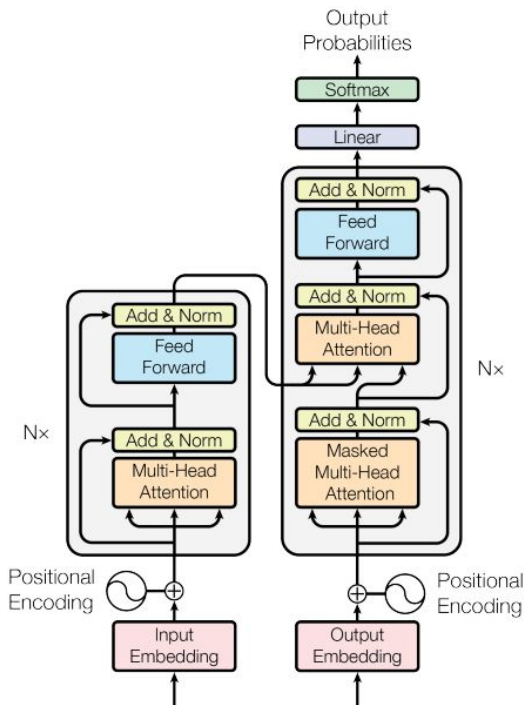
Gradient vanishing / explosion



Limitations

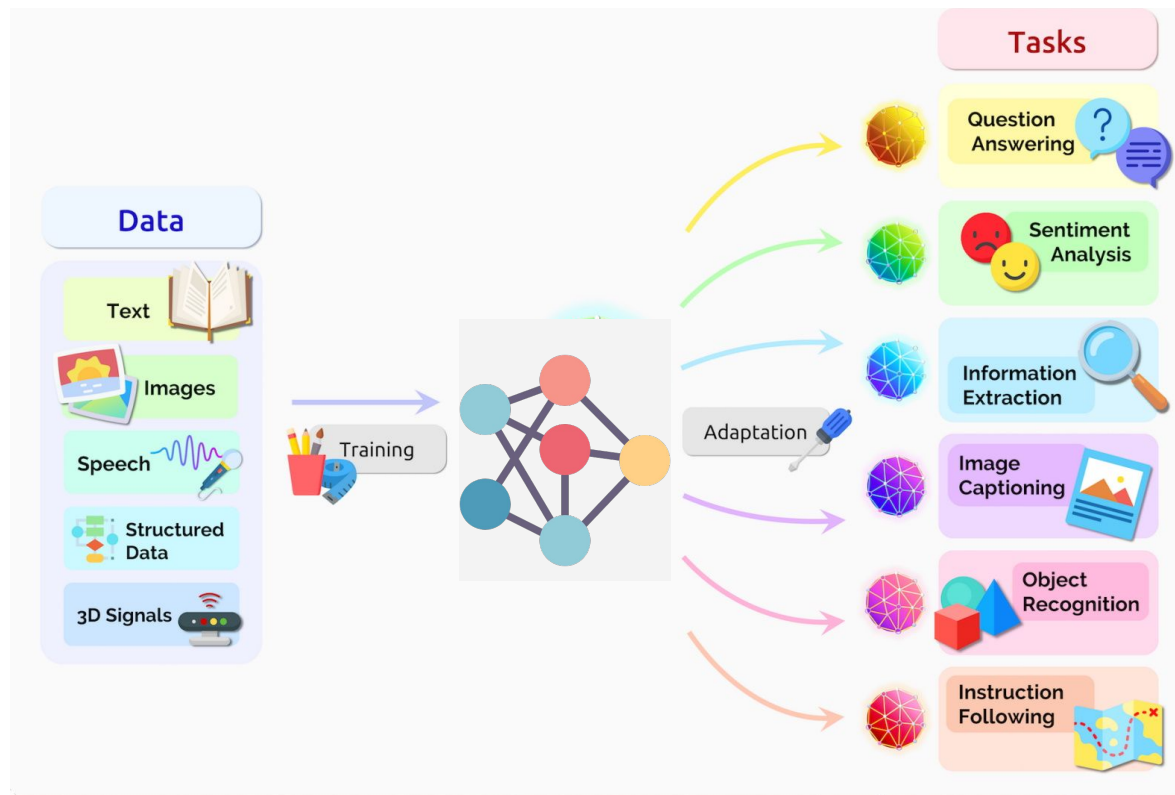
- Long Range Dependencies
- Gradient vanishing / explosion
- Long time to converge
- Expensive computation

Transformer Model



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

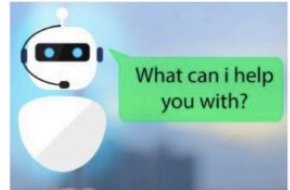
Wide Applications



Real World Impact



Machine Translation



Smart Assistants



Search Engines



Auto Transcription



Health Record Analysis



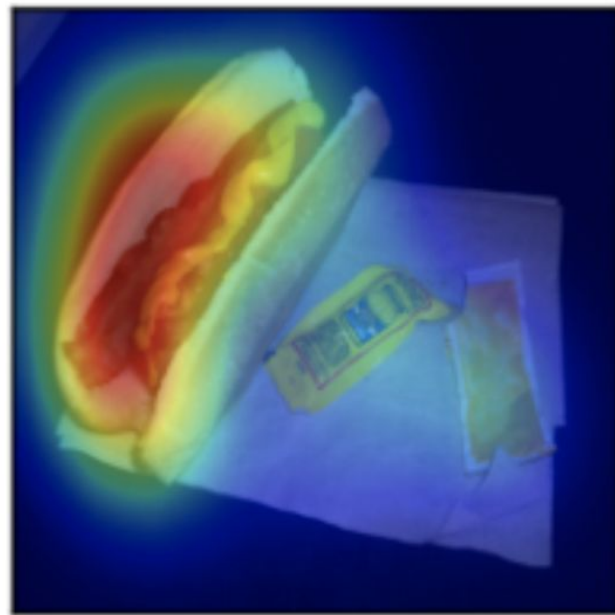
Summarization Engines

and many more

Questions?

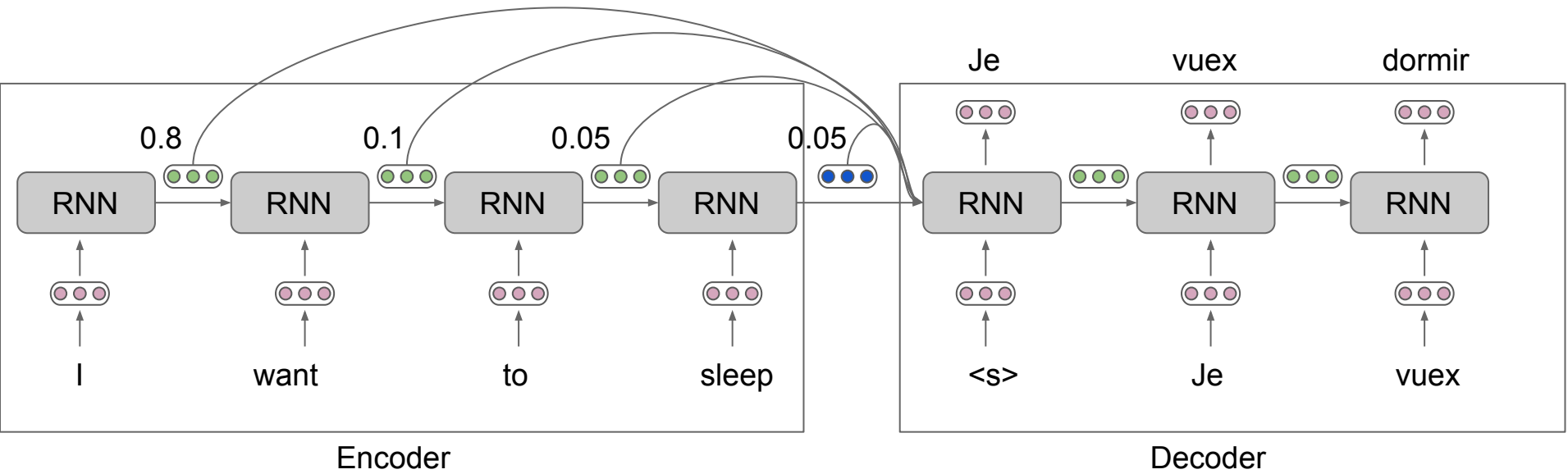
Visual Attention

What toppings are on the hot dog?

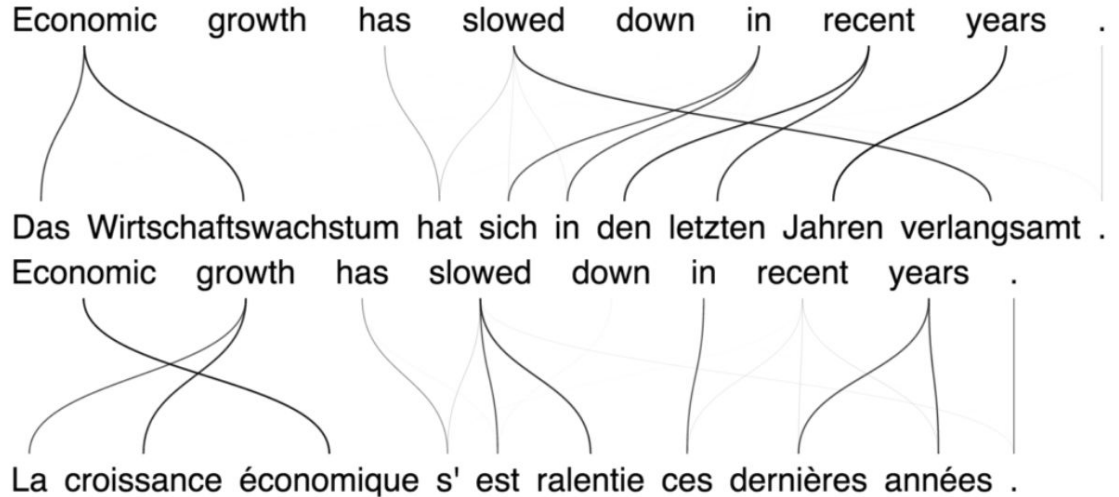
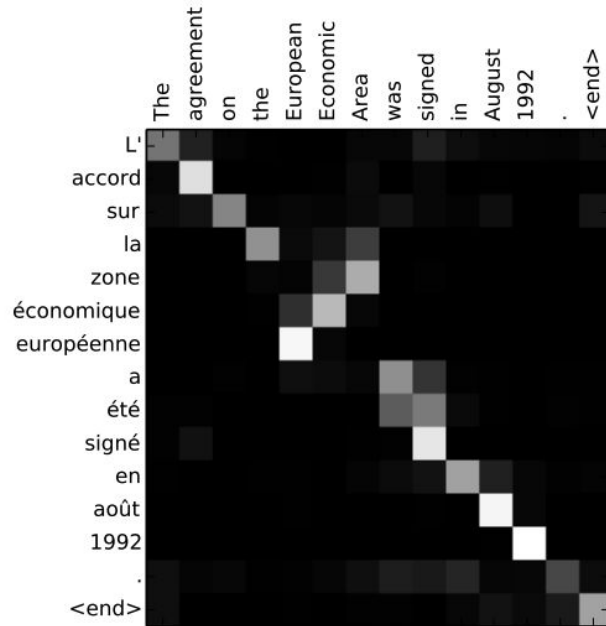


Differential Attention for Visual Question Answering (Patro et.al, 2018)

Cross Attention in NMT

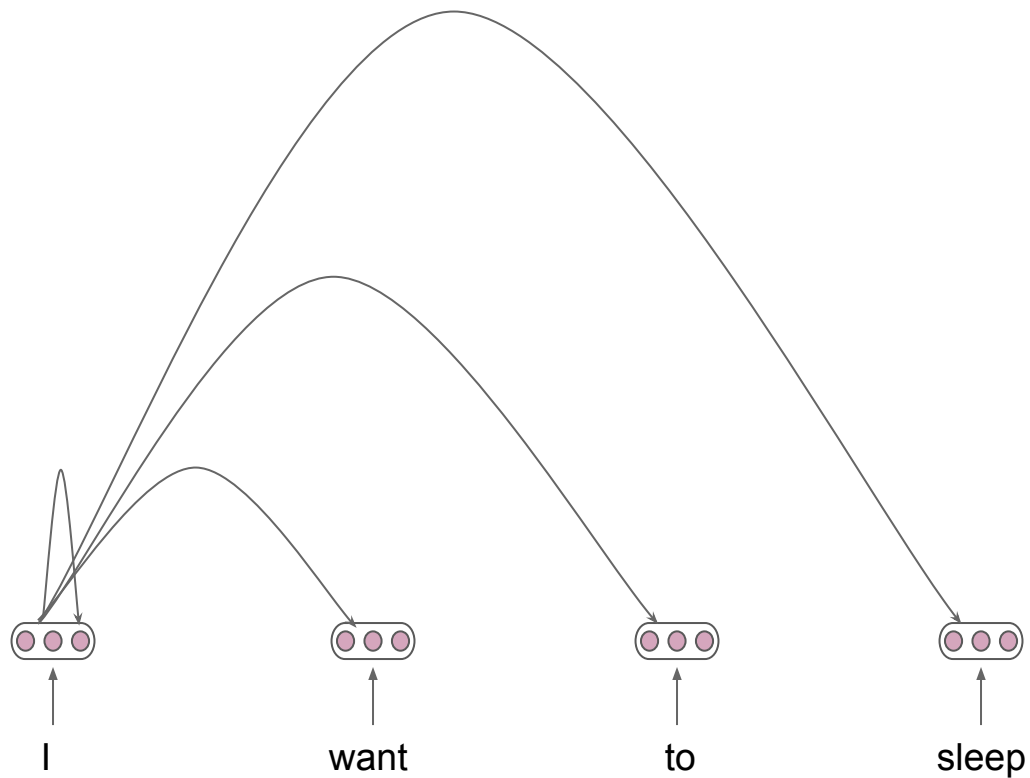


Attention in NMT

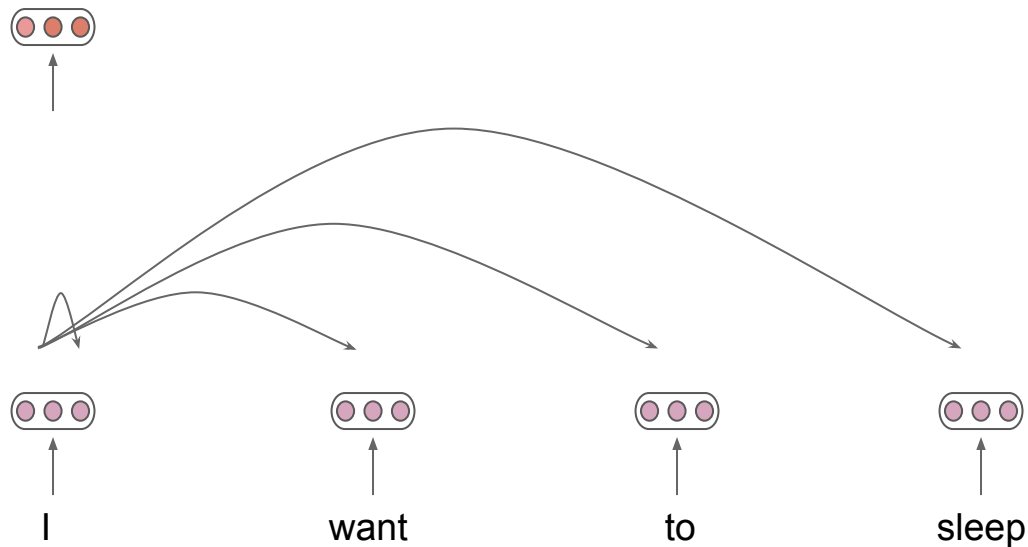


Neural Machine Translation by Jointly Learning to Align and Translate. Bahdanau et al, 2015
<https://developer.nvidia.com/blog/introduction-neural-machine-translation-gpus-part-3/>

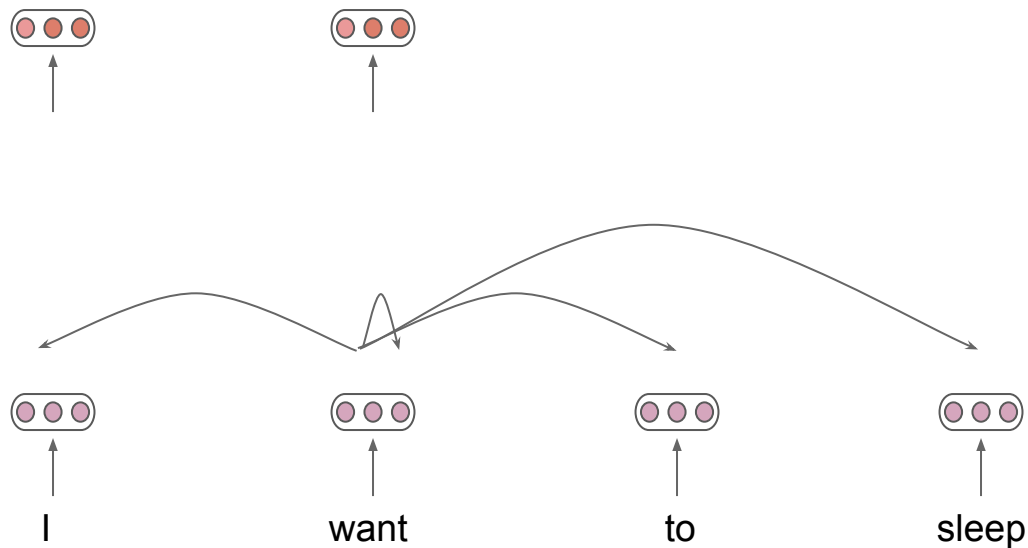
Self Attention



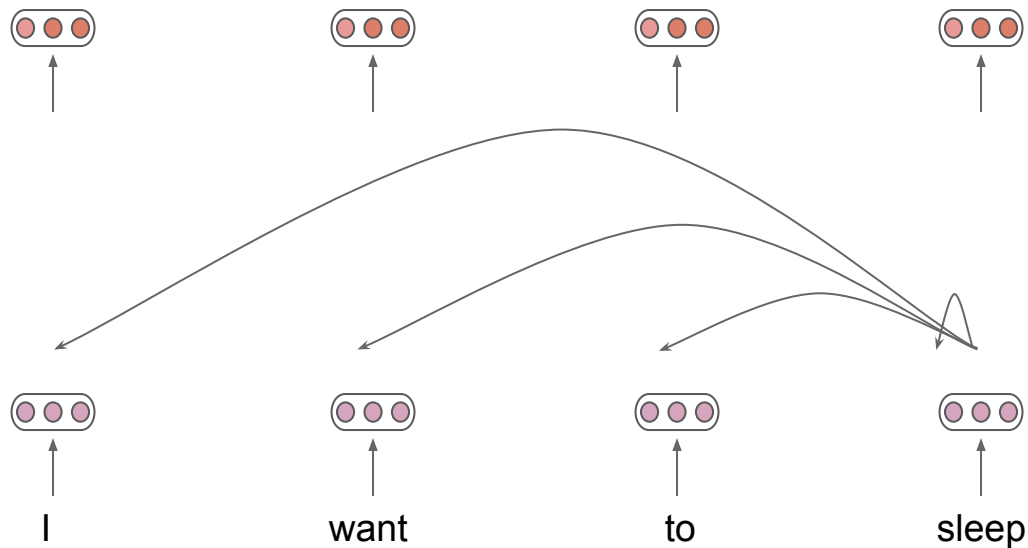
Self Attention - No more recurrence



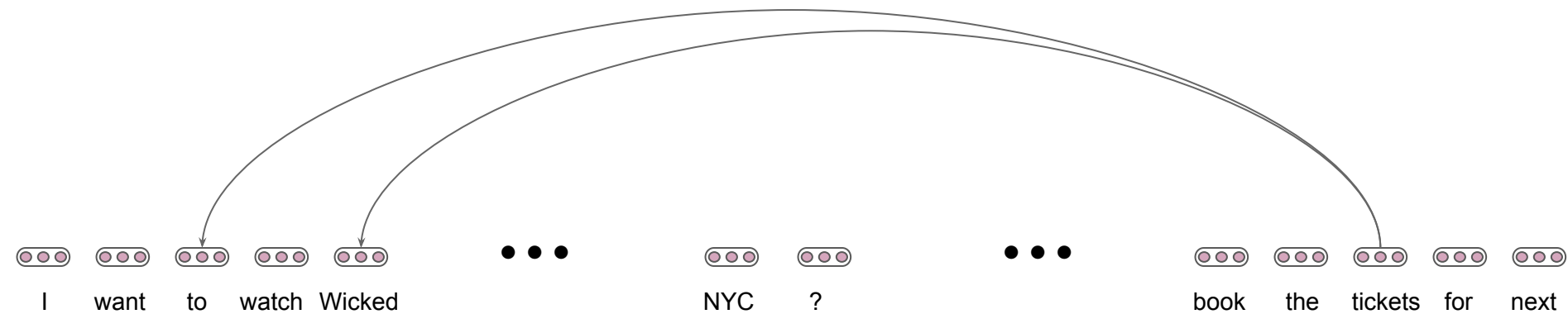
Self Attention - No more recurrence



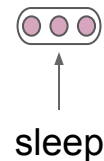
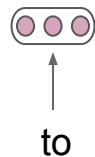
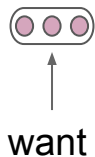
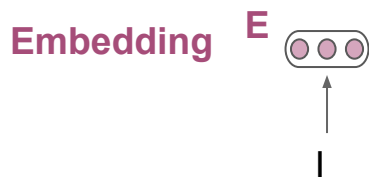
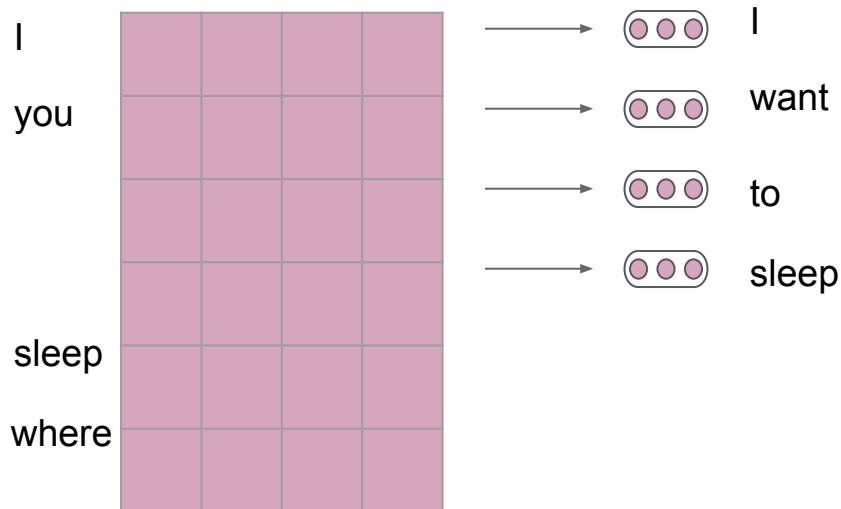
Self Attention - No more recurrence



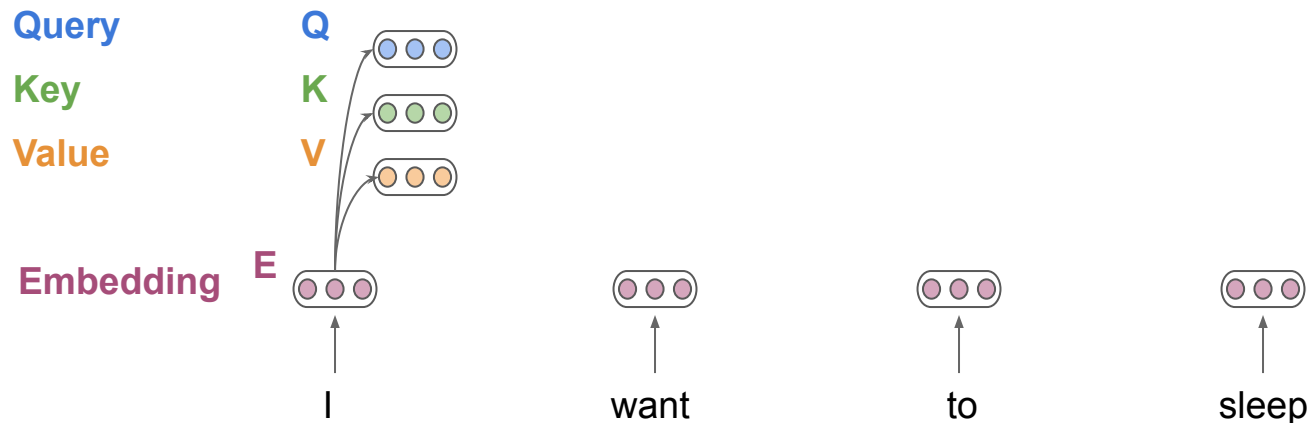
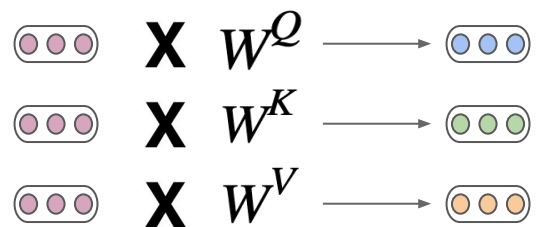
Self Attention for long sequences



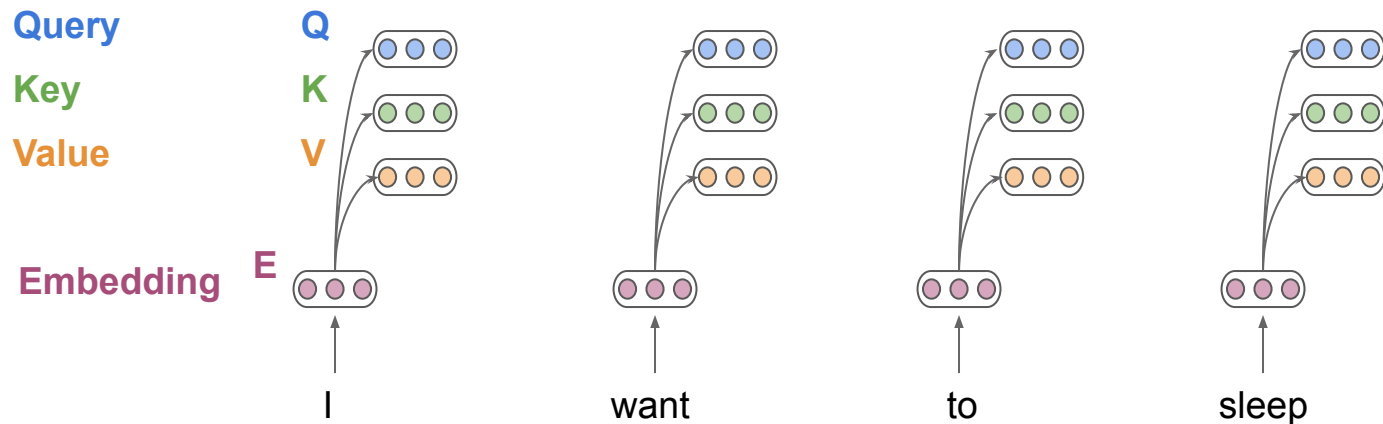
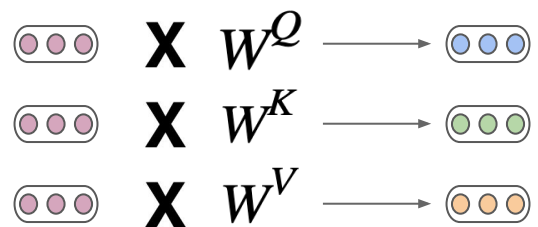
Self Attention - Word Embedding



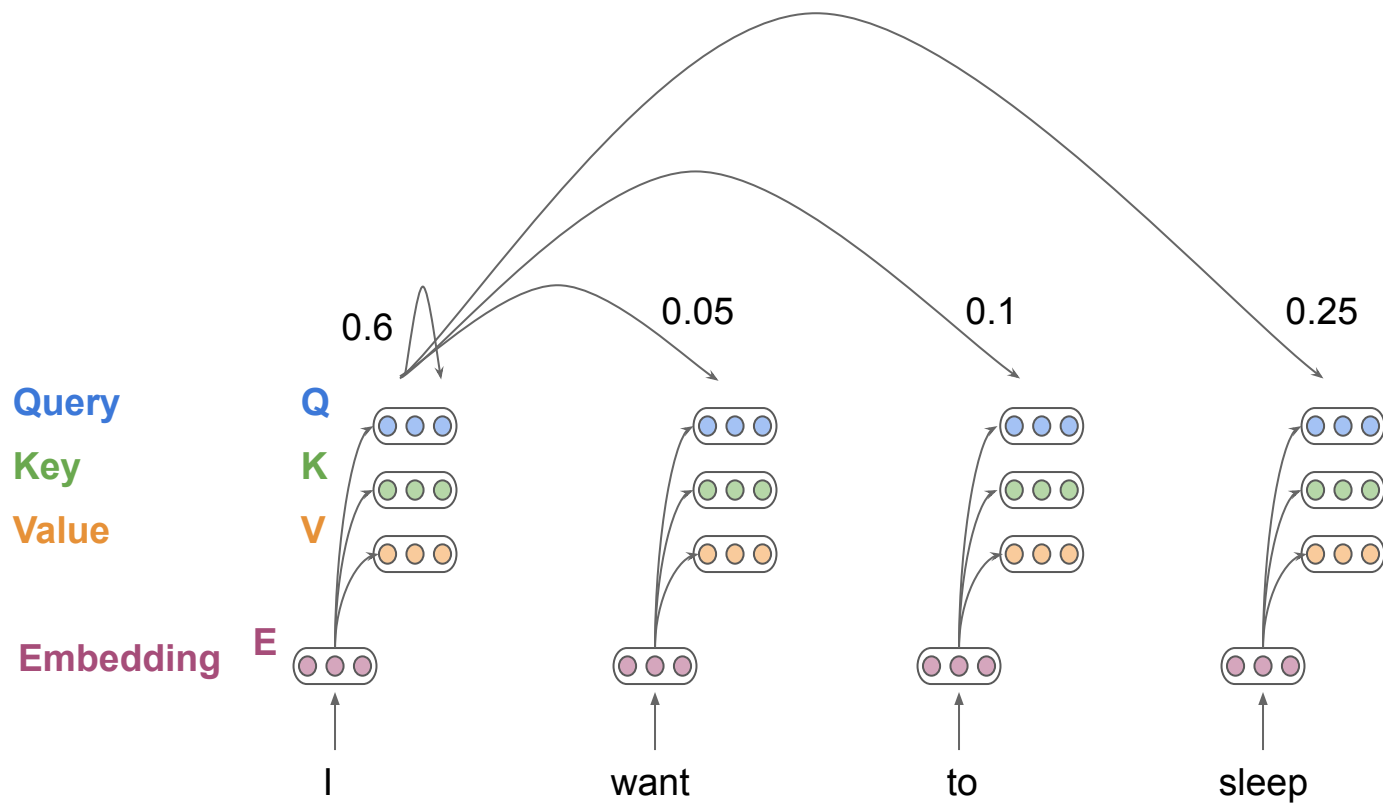
Self Attention - Projection Layer



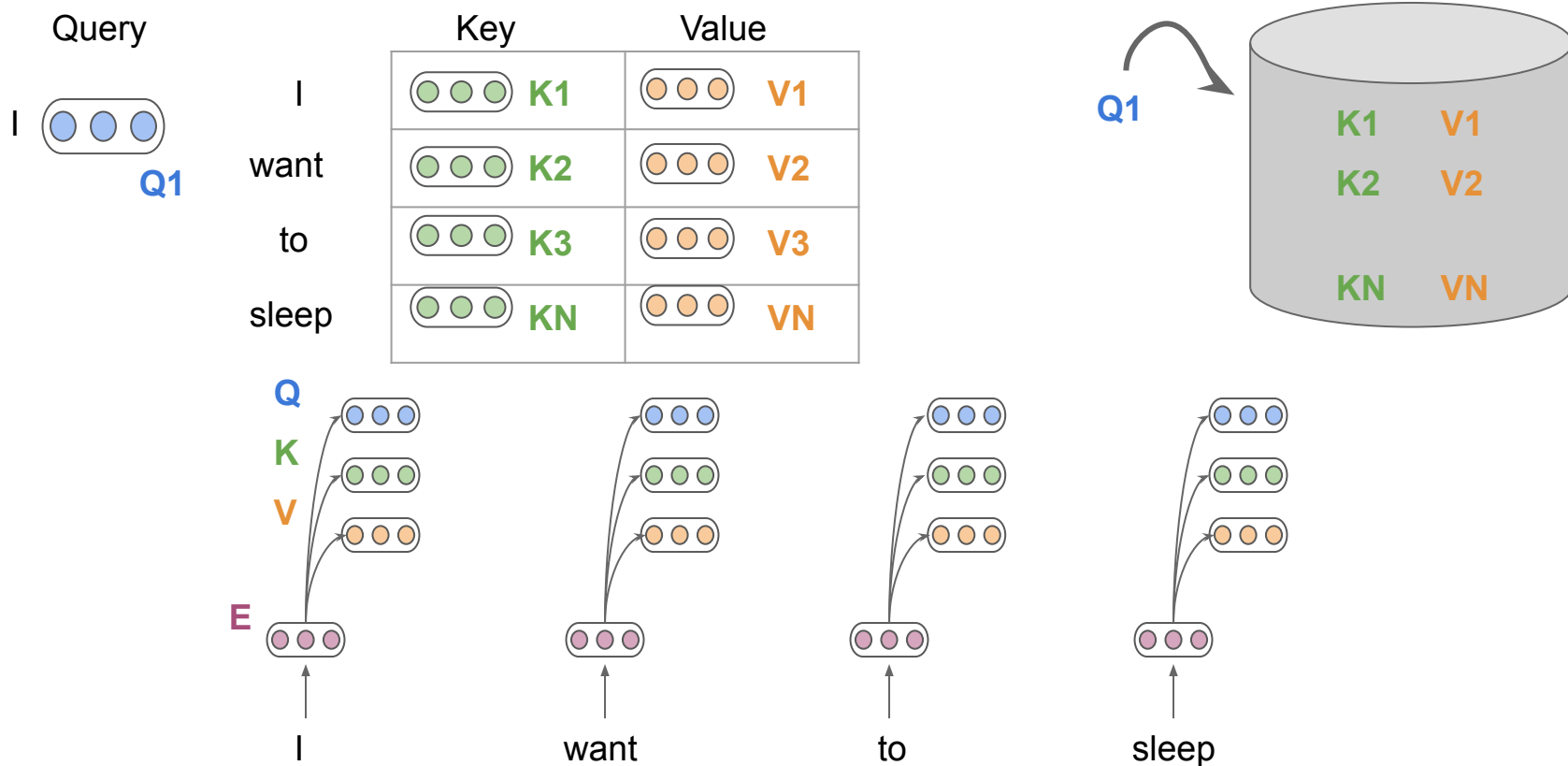
Self Attention - Projection Layer



Self Attention - Attention Scores

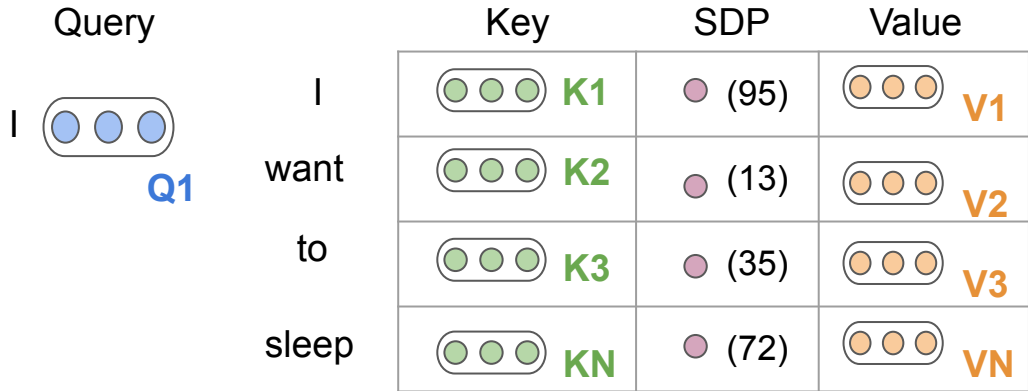


Self Attention



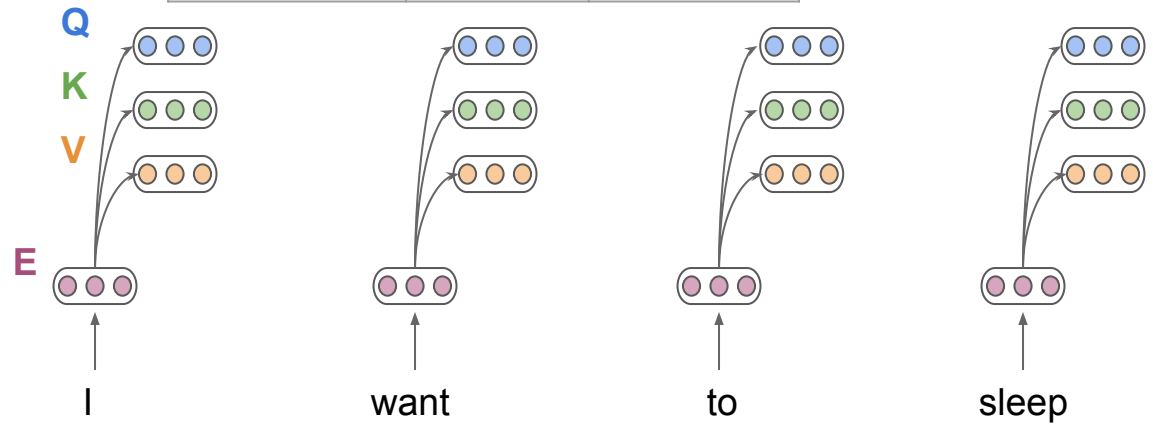
Questions?

Self Attention - Scaled Dot Product



Scaled Dot Product

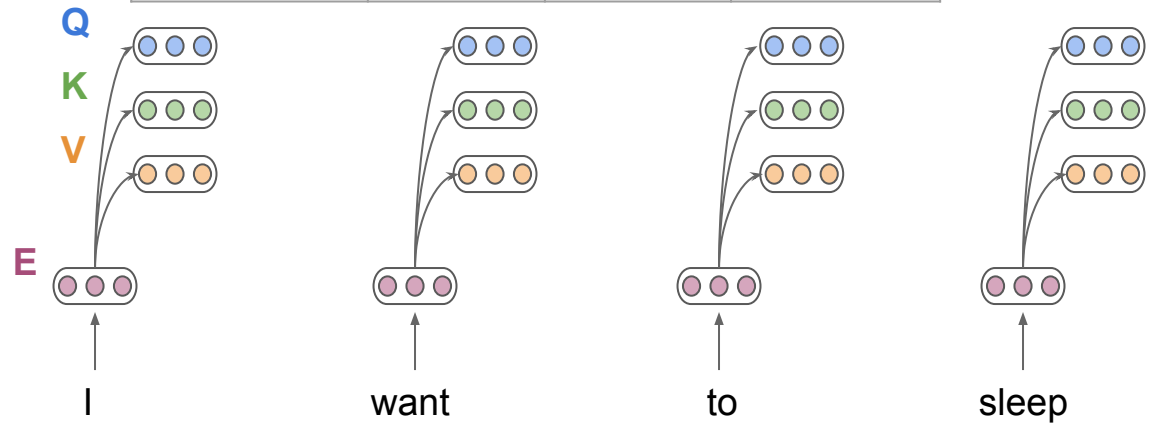
$$SDP = \frac{(QK^T)}{\sqrt{d^k}}$$



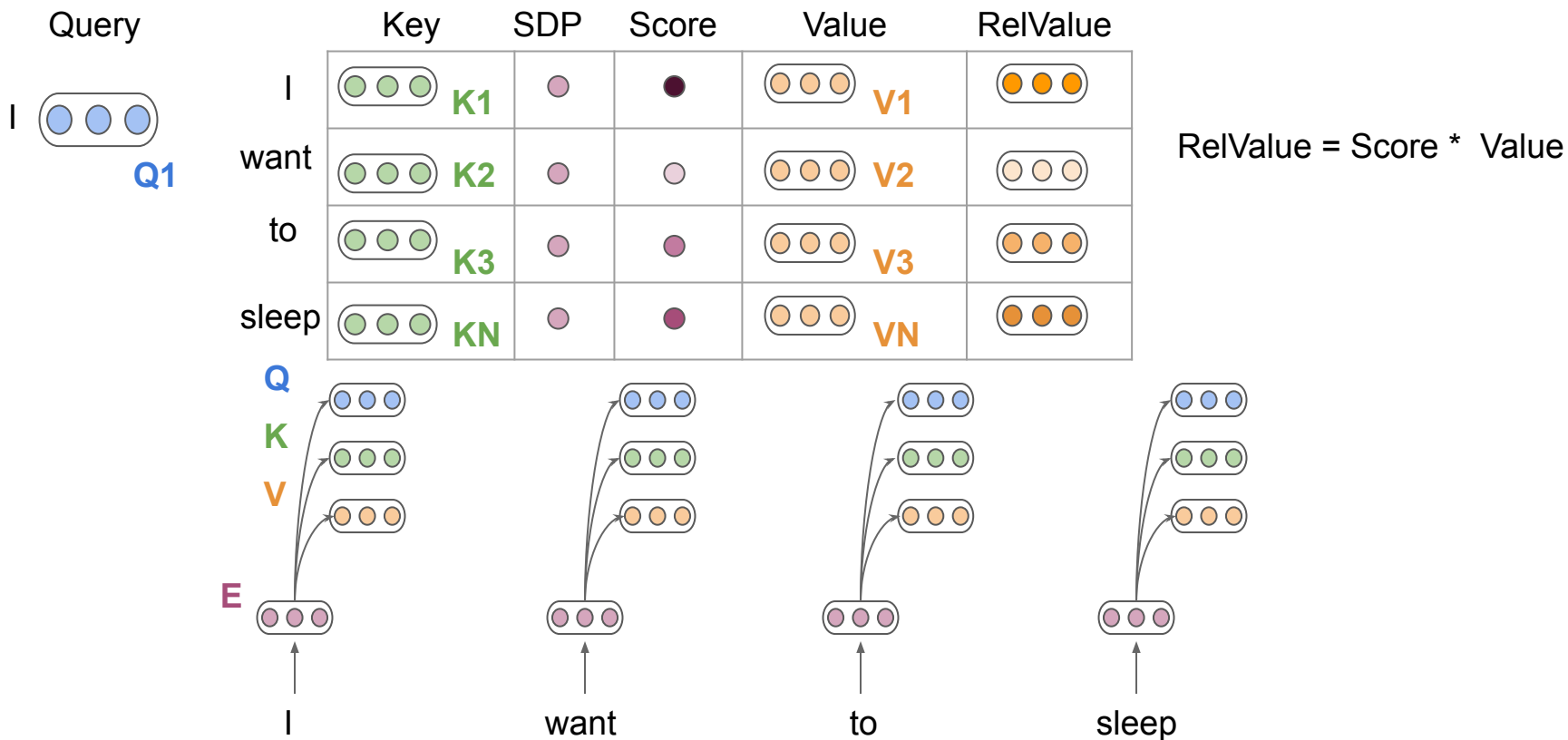
Self Attention - SoftMax

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

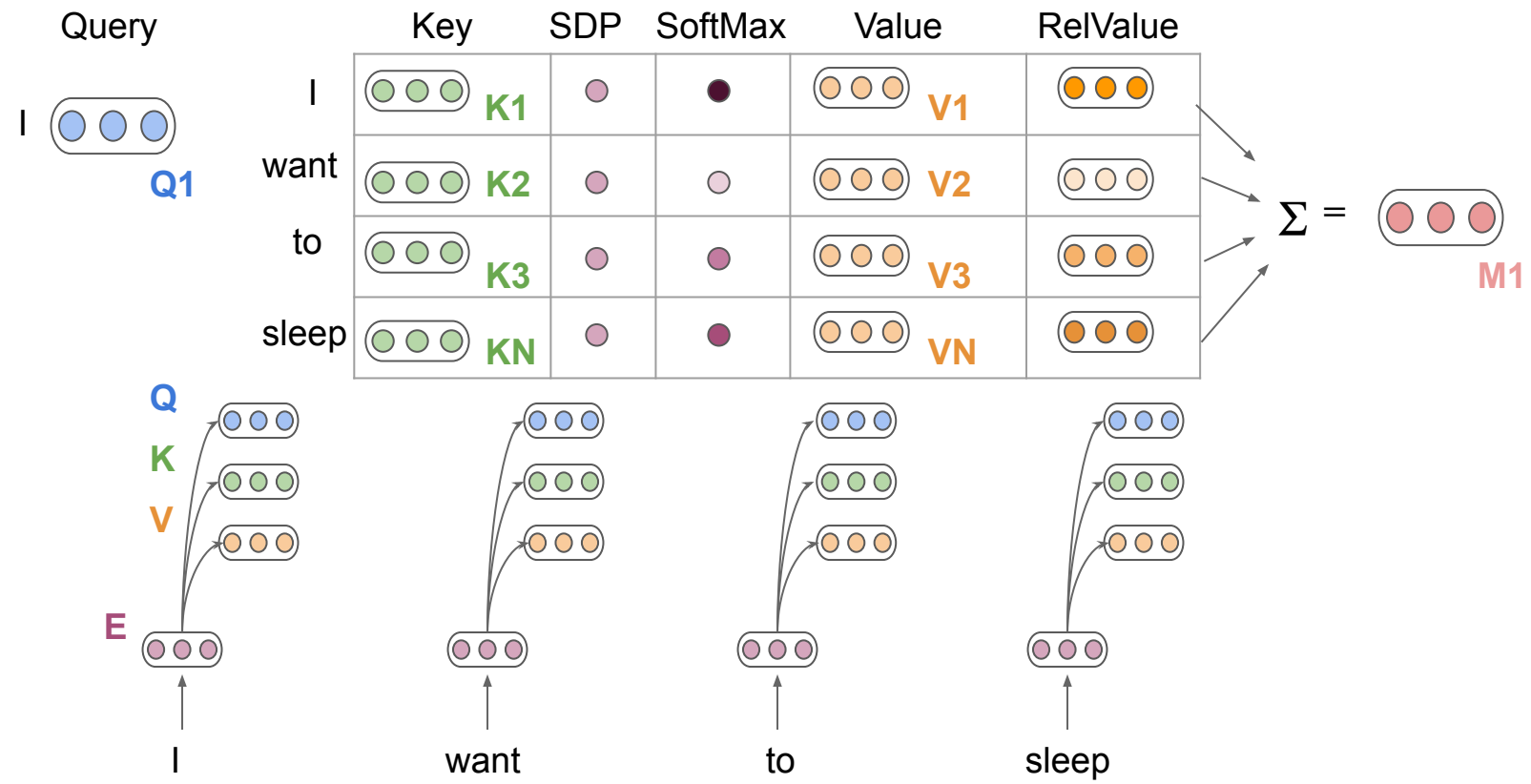
$$\text{score} = \text{softmax} \left(\frac{(QK^T)}{\sqrt{d^k}} \right)$$



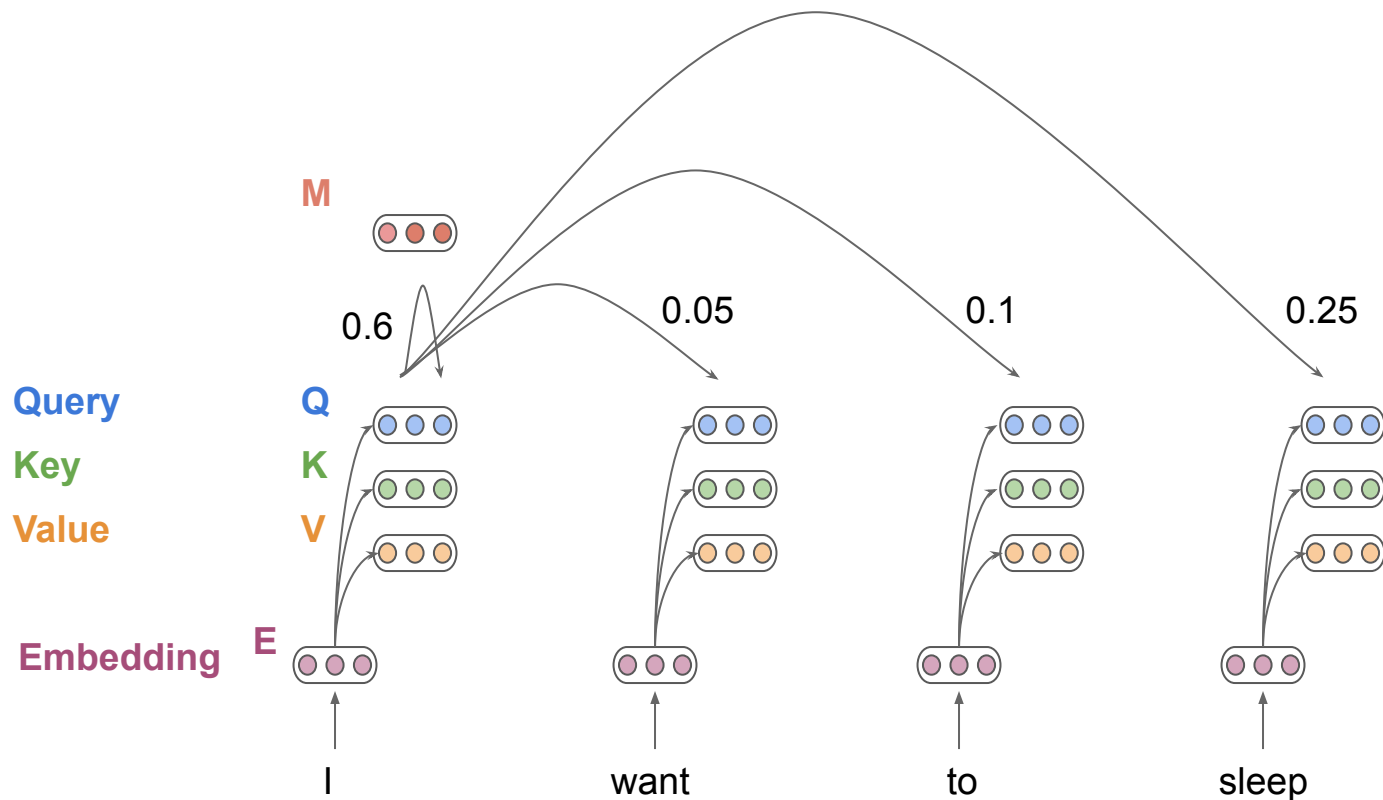
Self Attention - Soft (Relative) Values



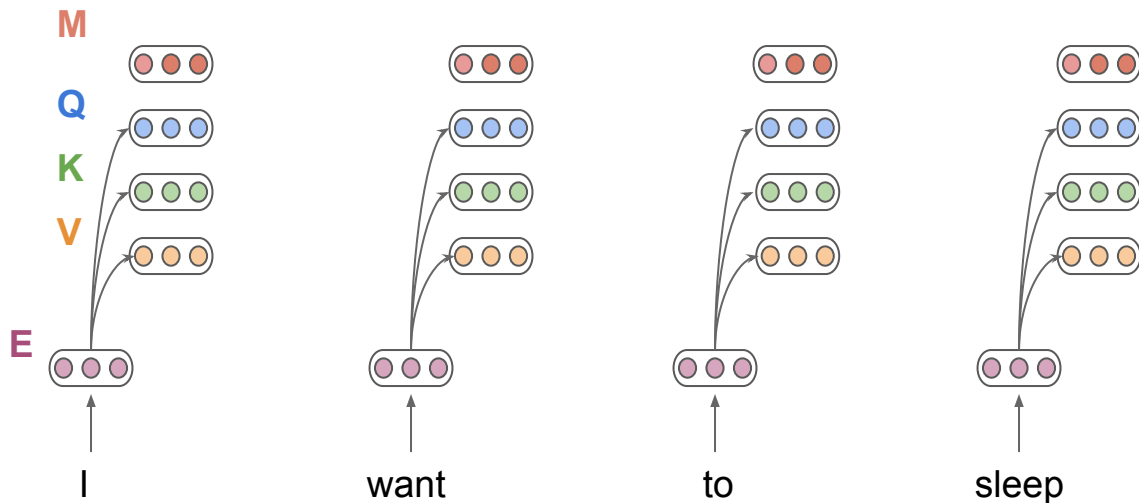
Self Attention - Attended Repr



Self Attention - Attended Contextual Rep



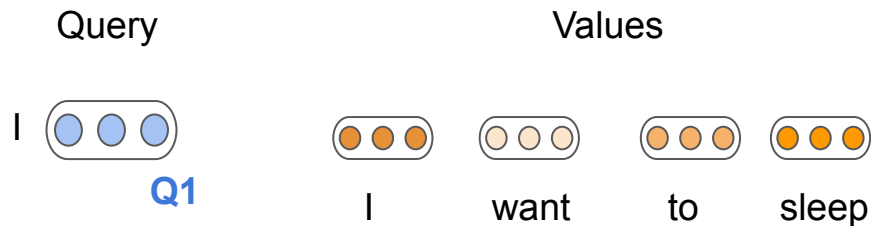
Self Attention - Attended Contextual Rep



Questions?

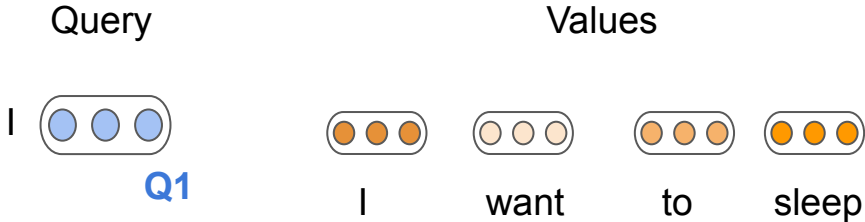
Problem with Self Attention

- Self Attention can focus heavily on the same word!



Problem with Self Attention

- Self Attention can focus heavily on the same word!



- Single representation

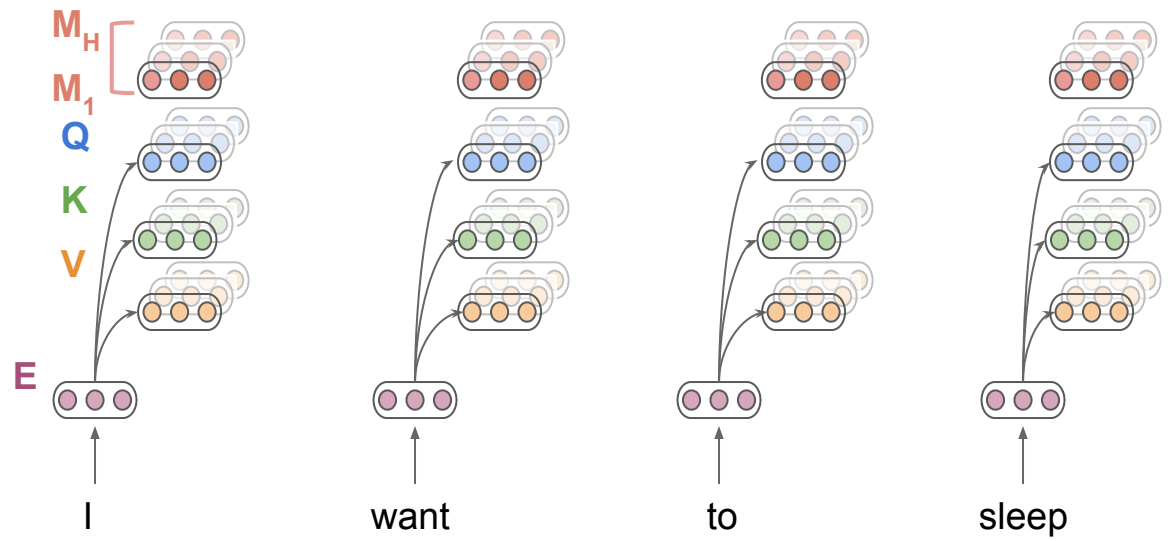
I like **Harry Potter**
(Book)

v/s

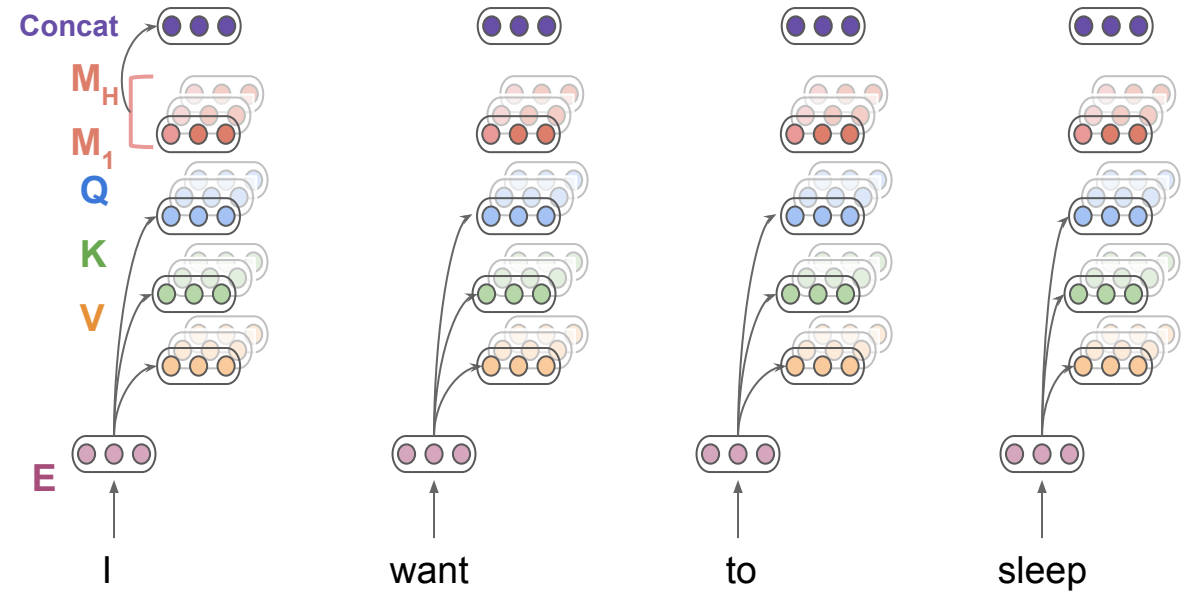
I like **Harry Potter**
(Movie)

Multi-Headed Self Attention

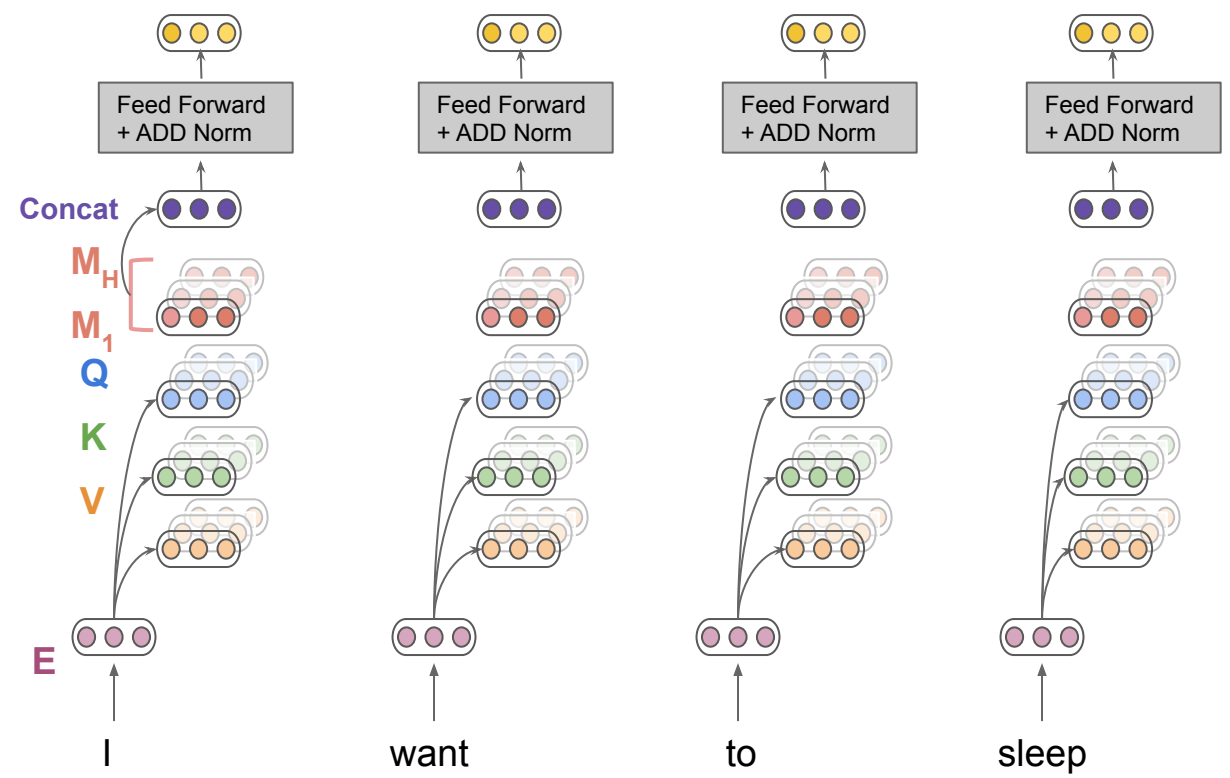
H (no: of heads) Different versions of Q,K,V
Each different repr -> Different attended repr



Multi-Headed Self Attention

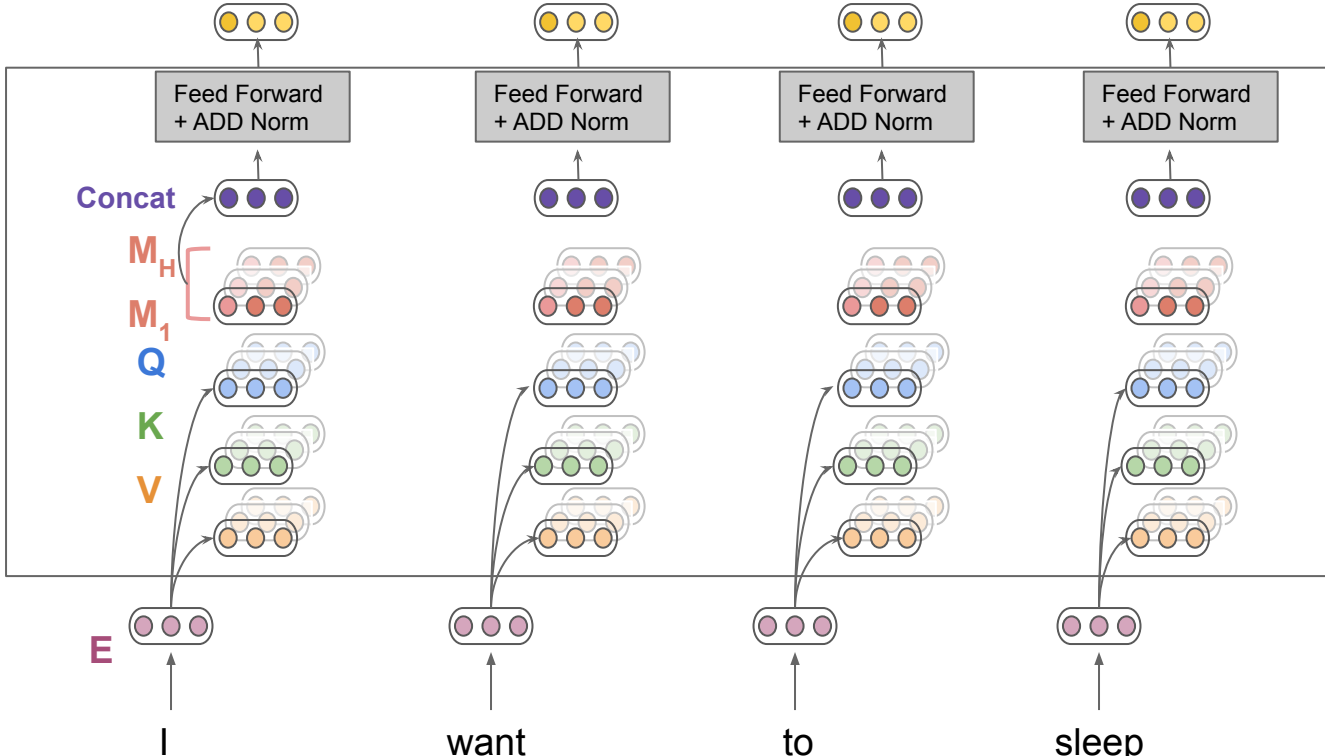


Multi-Headed Self Attention

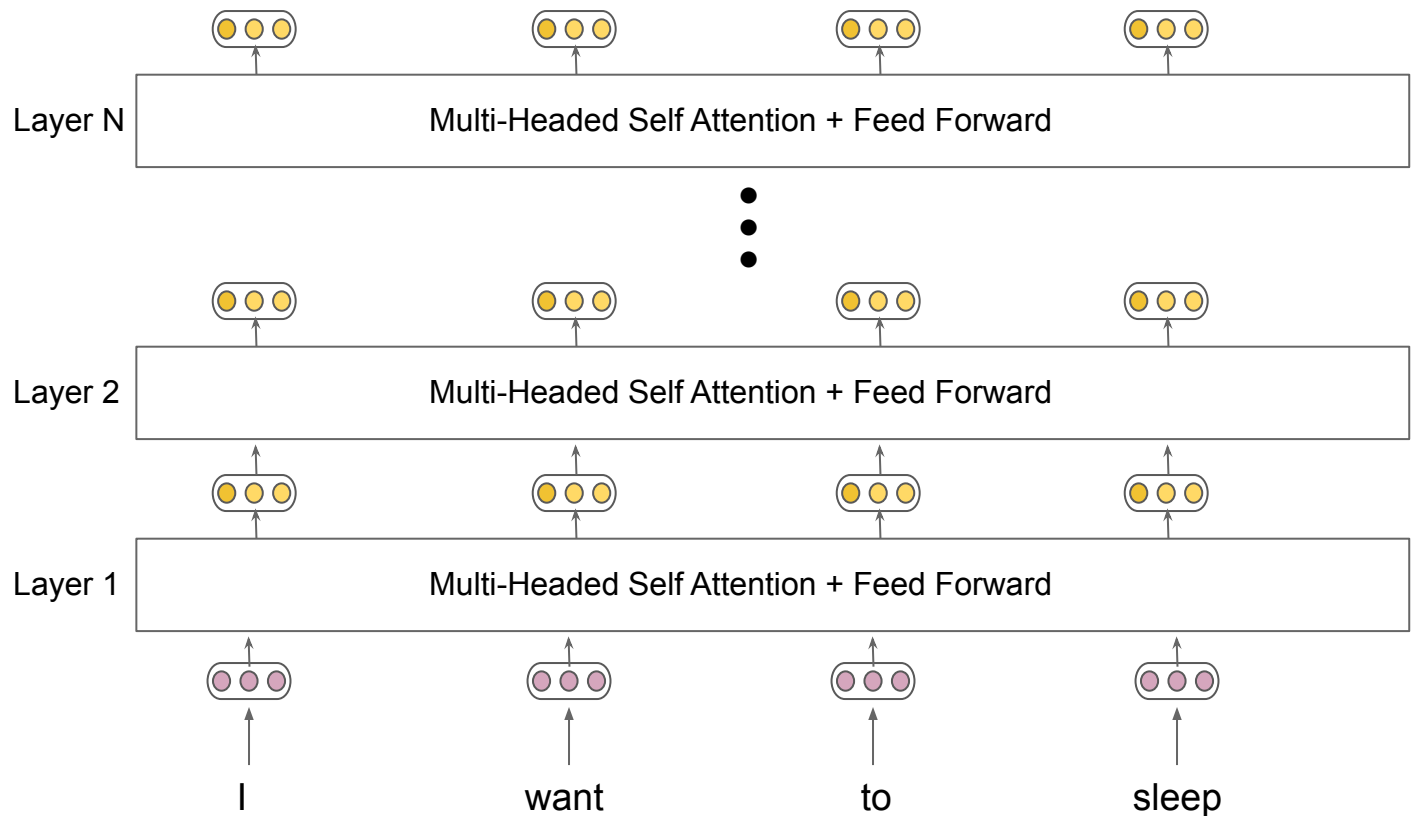


Multi-Headed Self Attention

Multi-Headed Self Attention + Feed Forward



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	

$\Sigma =$ 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	

$\Sigma =$ 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	

$\Sigma =$ 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	

$\Sigma =$ 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	

$\Sigma =$ M

I want to sleep

Query (I)



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

$\Sigma =$ M

Sleep to I want

Revisiting Self Attention



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

Same representation for both sentences - But positions matter!



I want to sleep



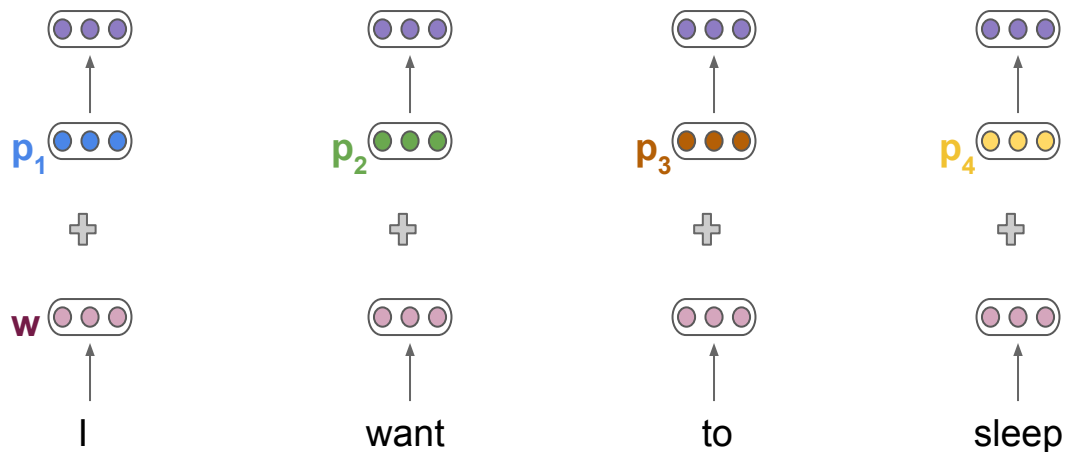
	Key	SDP	SM	Value	RelValue
Sleep	K1			VN	
to	K2			V3	
I	K3			V1	
want	KN			V2	



Sleep to I want

Positional Encoding

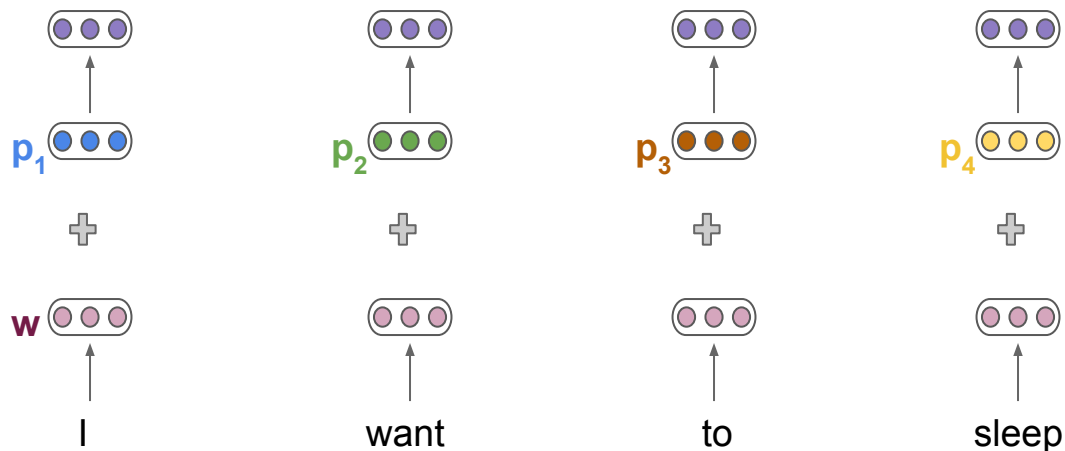
Position embeddings - each position number has an associated embedding



Positional Encoding

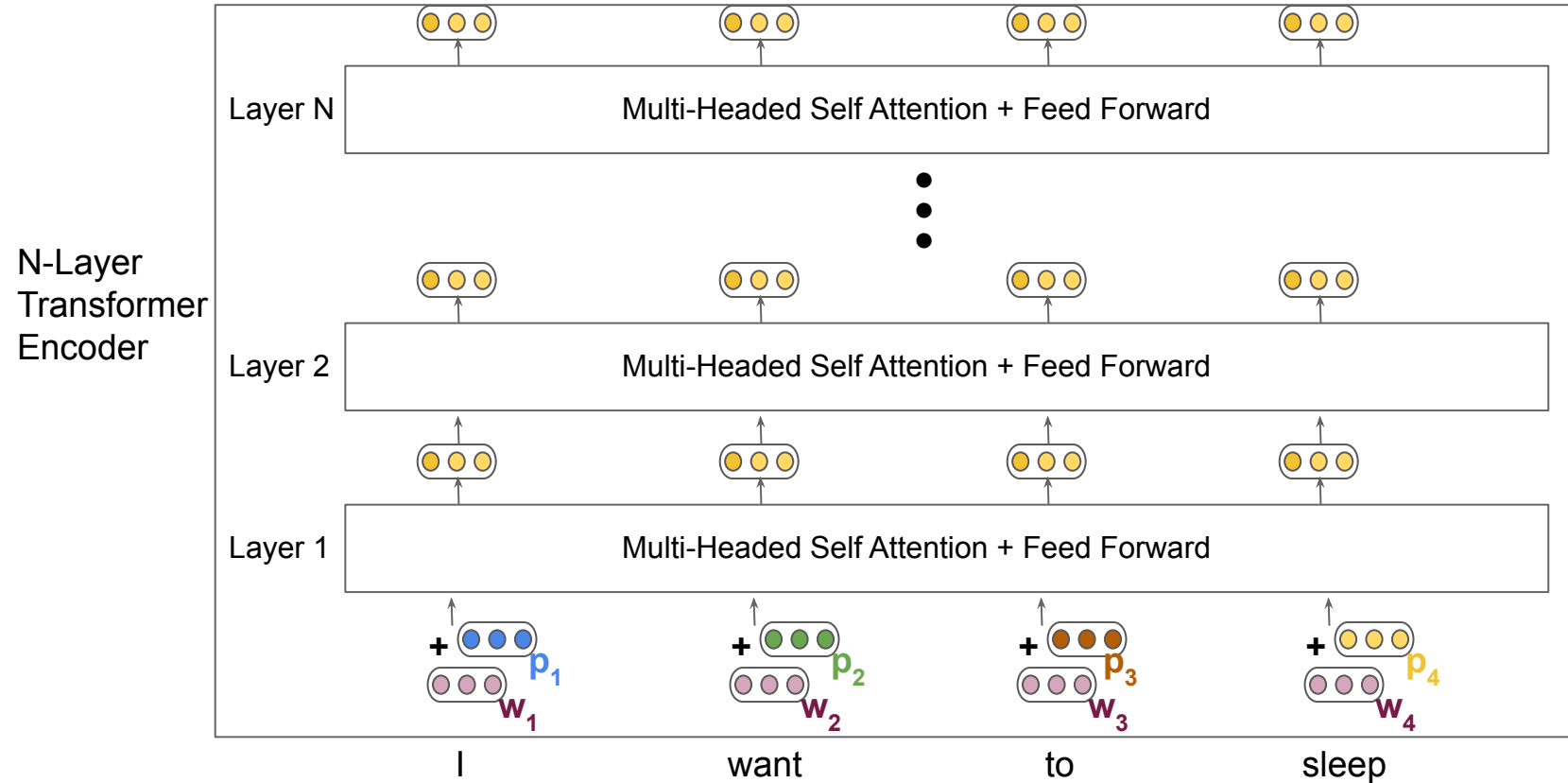
Sinusoidal Position embeddings - generalize to any sequence length

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

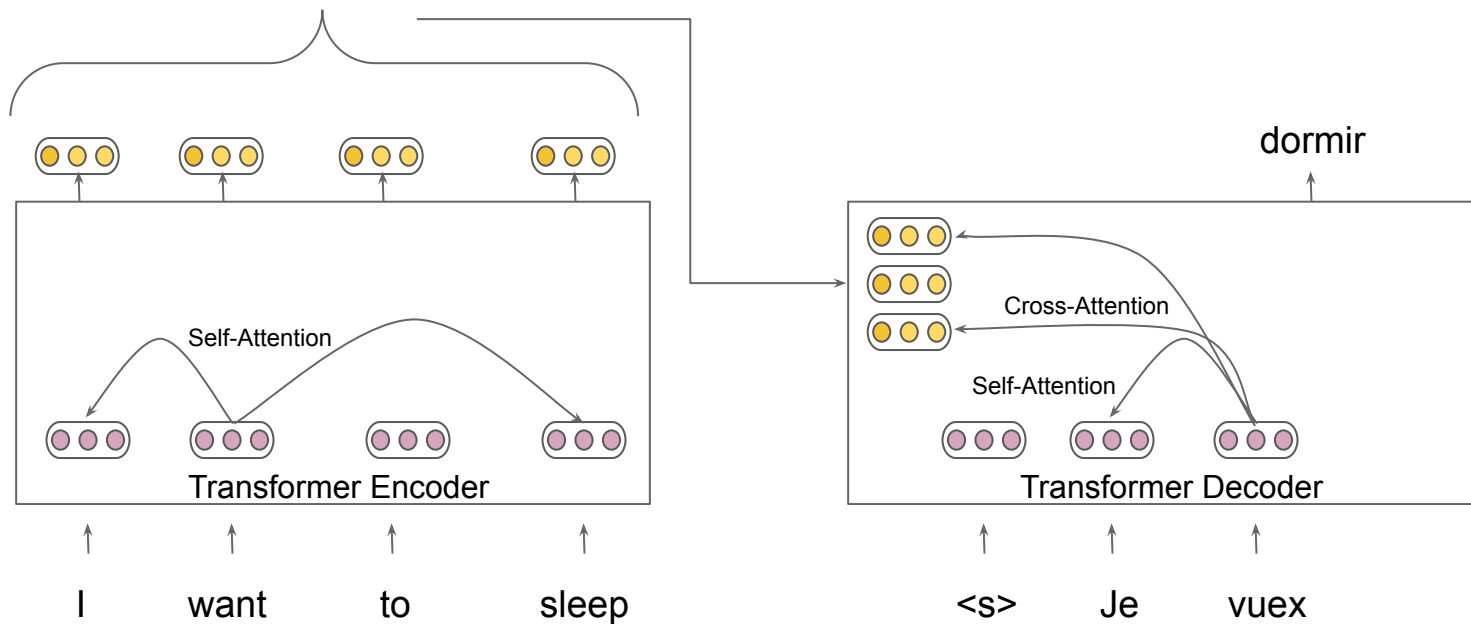


Questions?

Transformer Encoder



Transformer Encoder - Decoder

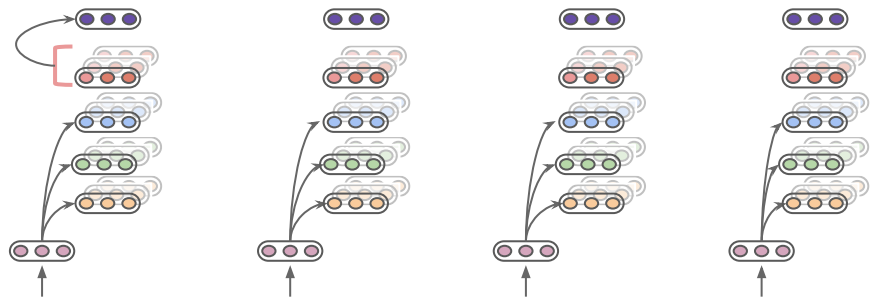


What's so great about Transformers?

- 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

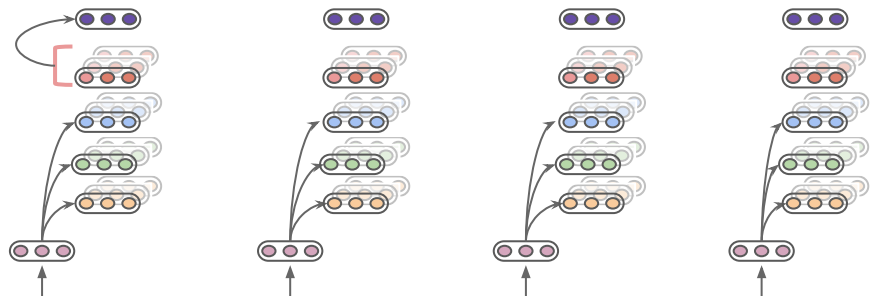
What's so great about Transformers?

- Parallelizable computation - Entire sequence can be processed in parallel

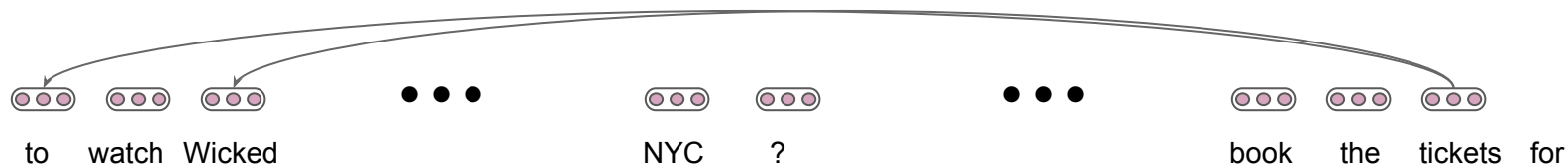


What's so great about Transformers?

- Parallelizable computation - Entire sequence can be processed in parallel

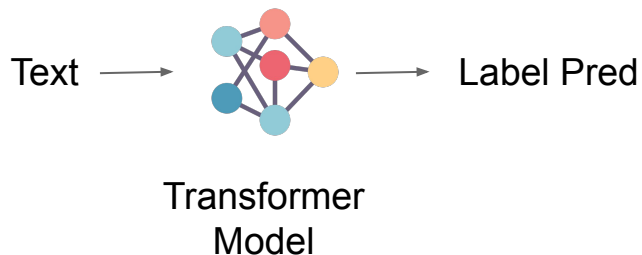


- Rich expressive power - long range dependencies

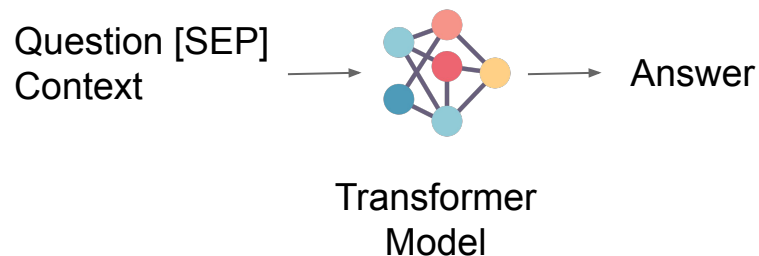


Impact - Wide Applications!

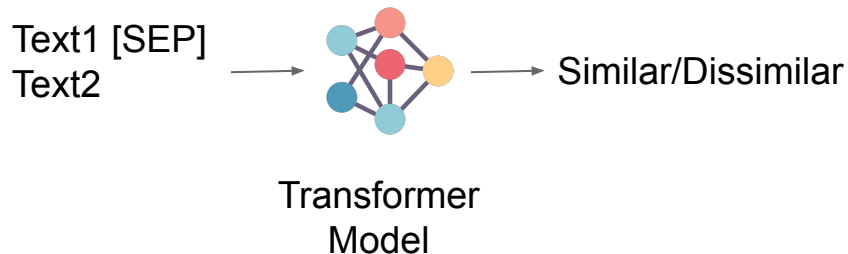
Classification



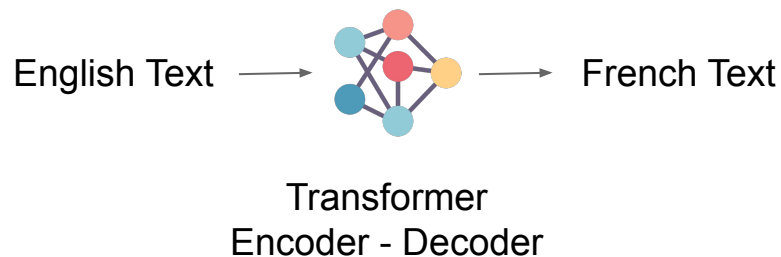
Question Answering



Sentence Similarity

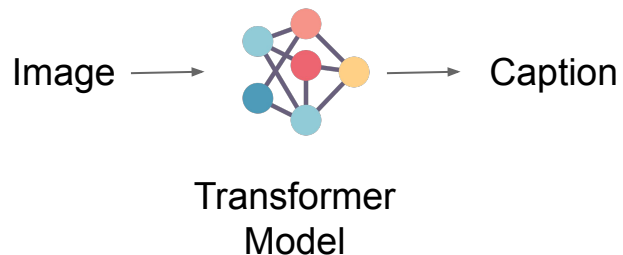


Translation



Impact - Wide Applications!

Captioning



Visual Question Answering

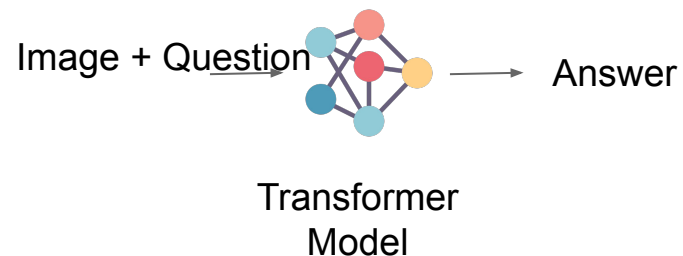
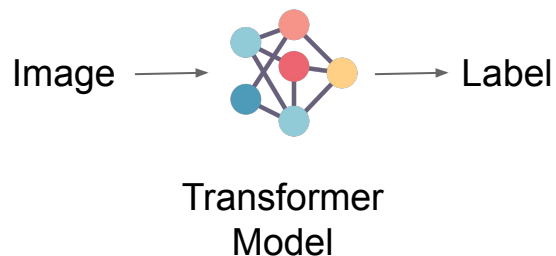
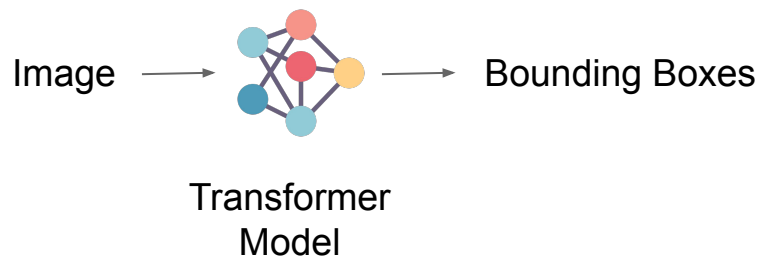


Image Classification



Object Detection

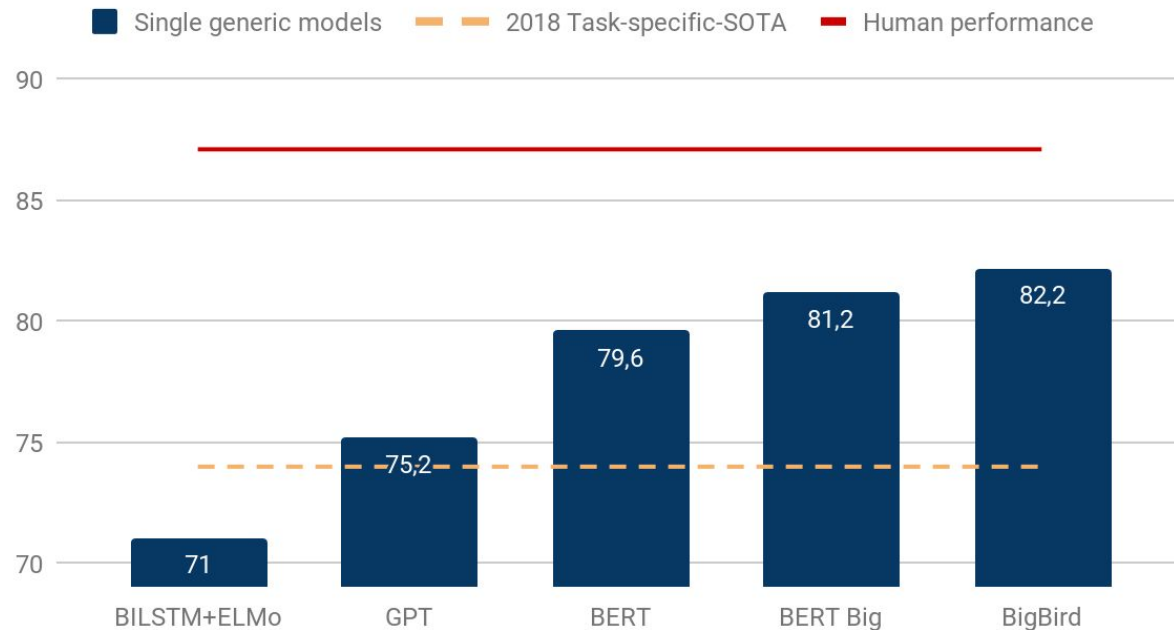


Larger Impact

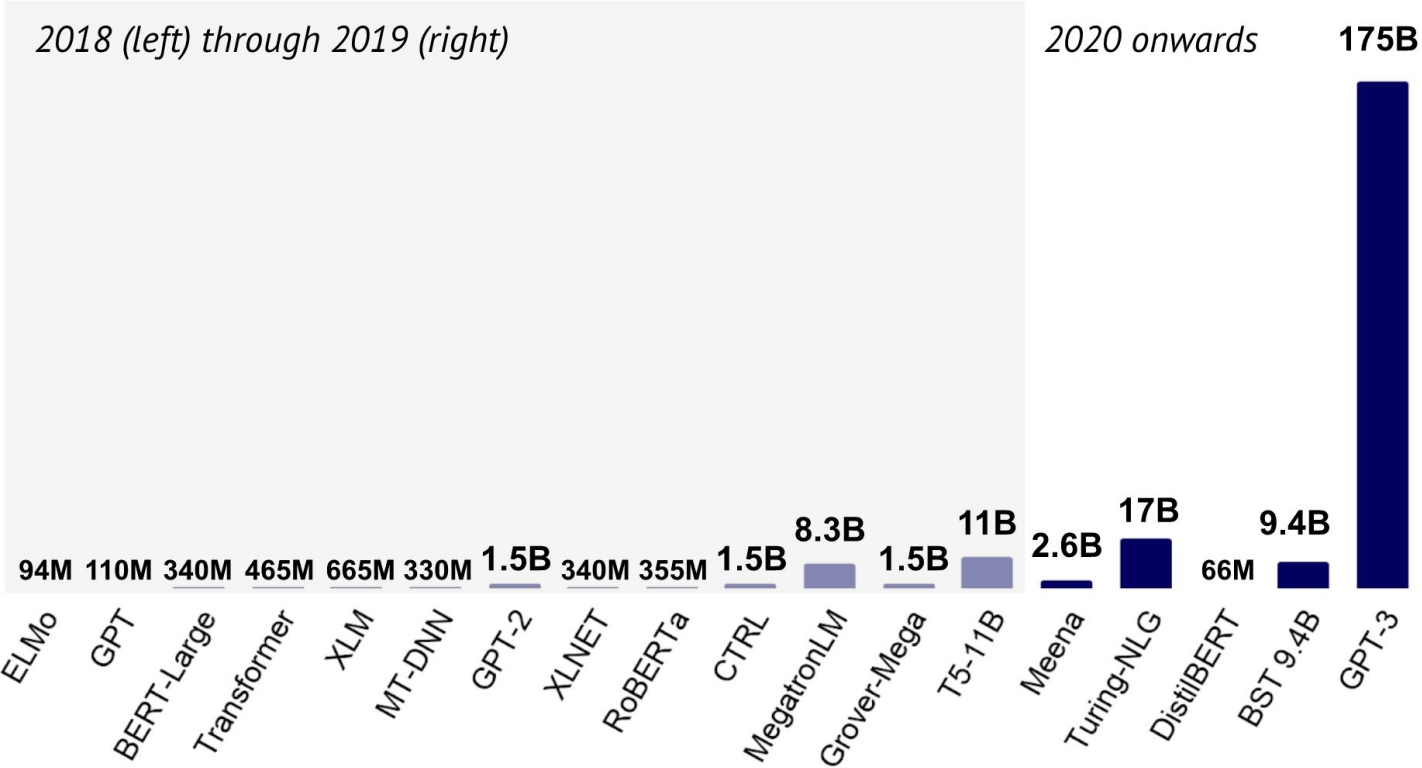


Larger Impact

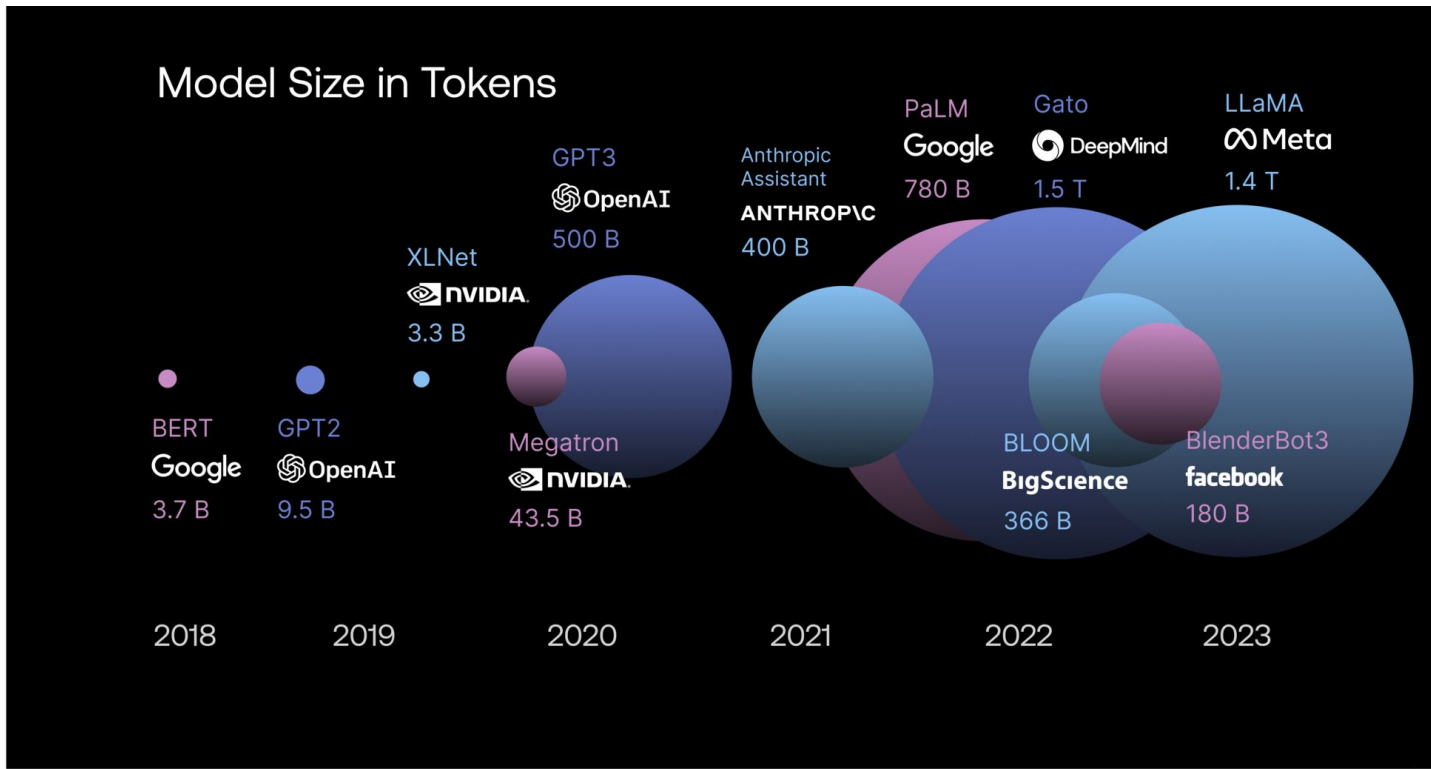
GLUE scores evolution over 2018-2019



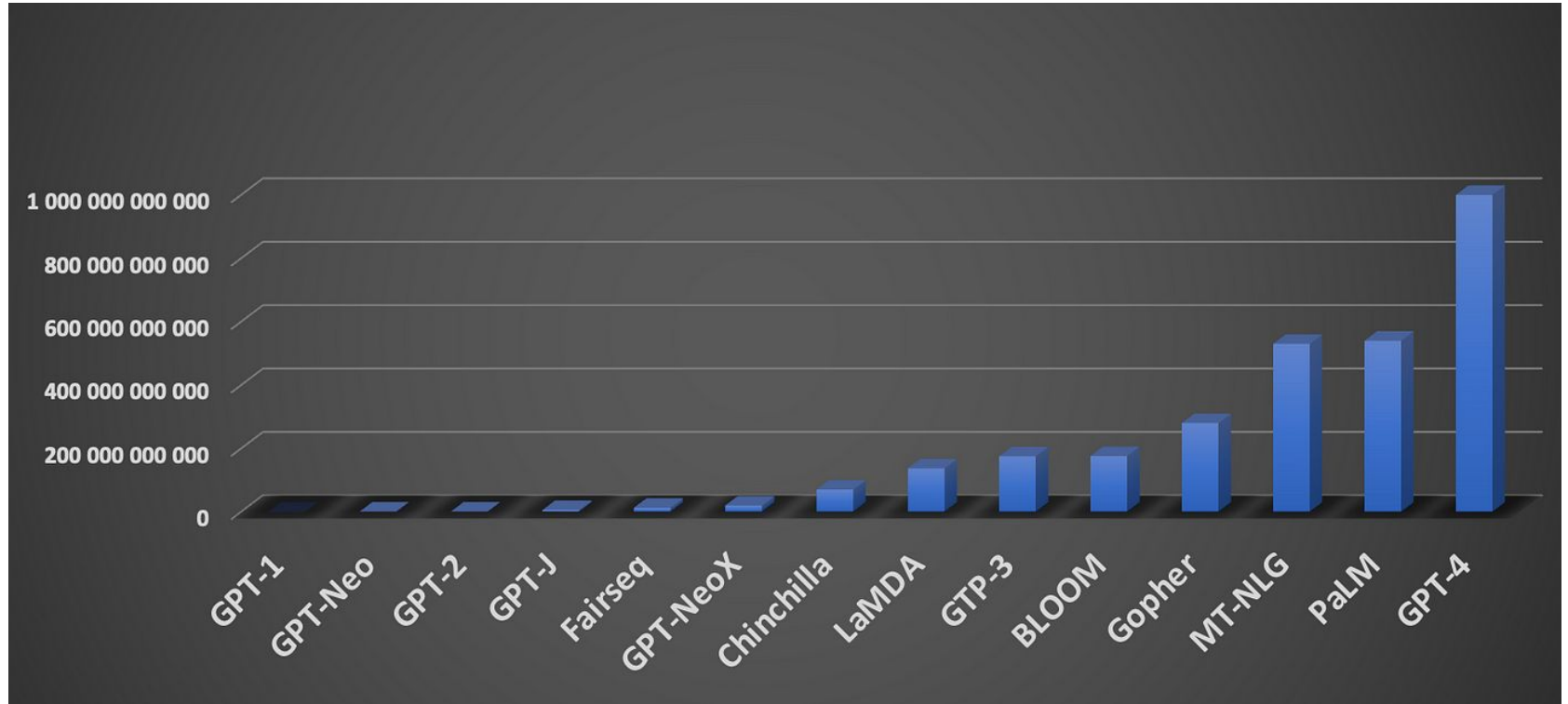
Larger Impact



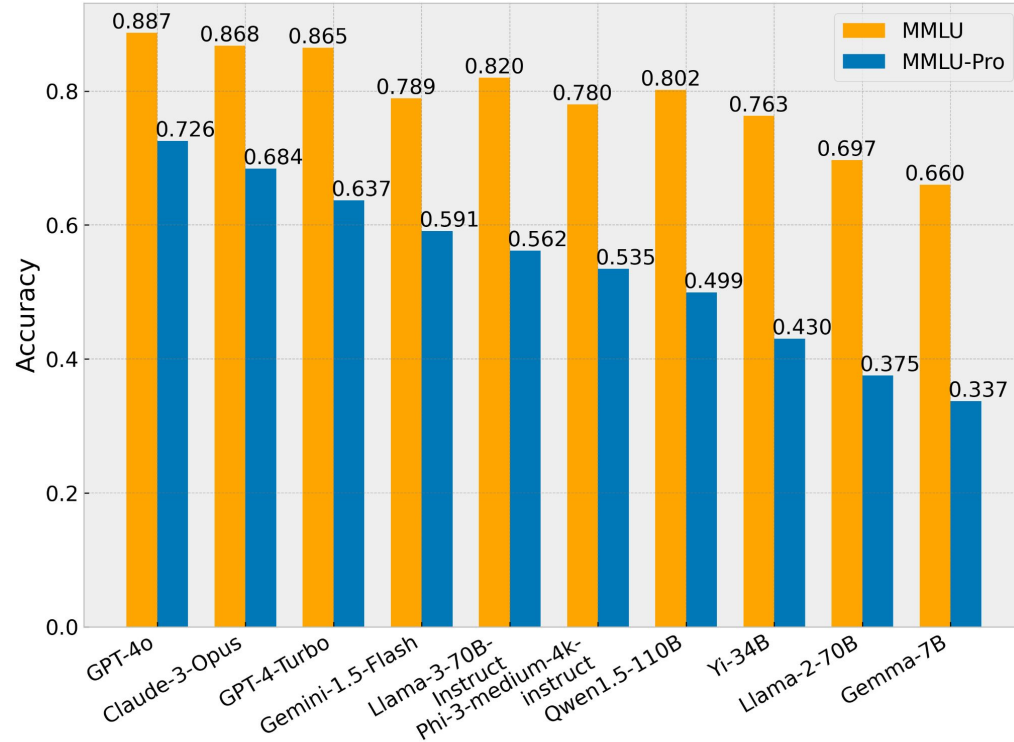
Larger Impact



Larger Impact



Larger Impact



Larger Impact

The
Economist

AI will revolutionise research. But could it transform science altogether?

REUTERS

How will leveraging AI change the future of legal services?

Microsoft Research Blog

GPT-4's potential in shaping the future of radiology

Microsoft

Announcing Microsoft Copilot, your everyday AI companion

POLITICO

More schools want your kids to use ChatGPT. Really.

Education leaders are embracing technology that set off a plagiarism panic just months ago.

THE FIFTY

Healthcare IT News

NYU Langone Health LLM can predict hospital readmissions

The Verge

Bing, Bard, and ChatGPT: How AI is rewriting the internet

Thank you!

vidhishab@microsoft.com

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