



Large Models – COMP4423 Computer Vision

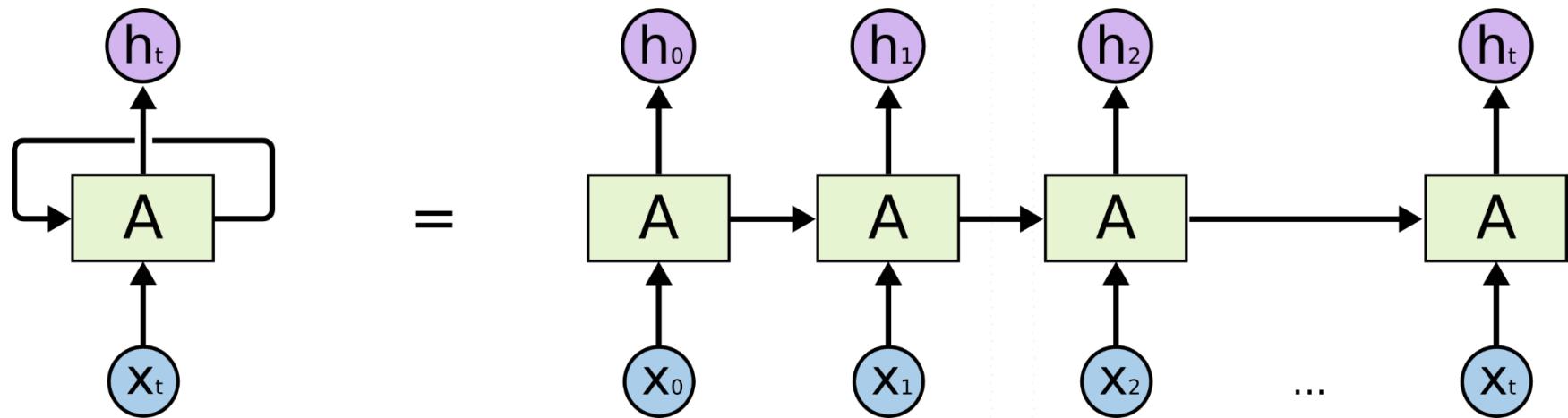
Xiaoyong Wei (魏驍勇)

x1wei@polyu.edu.hk

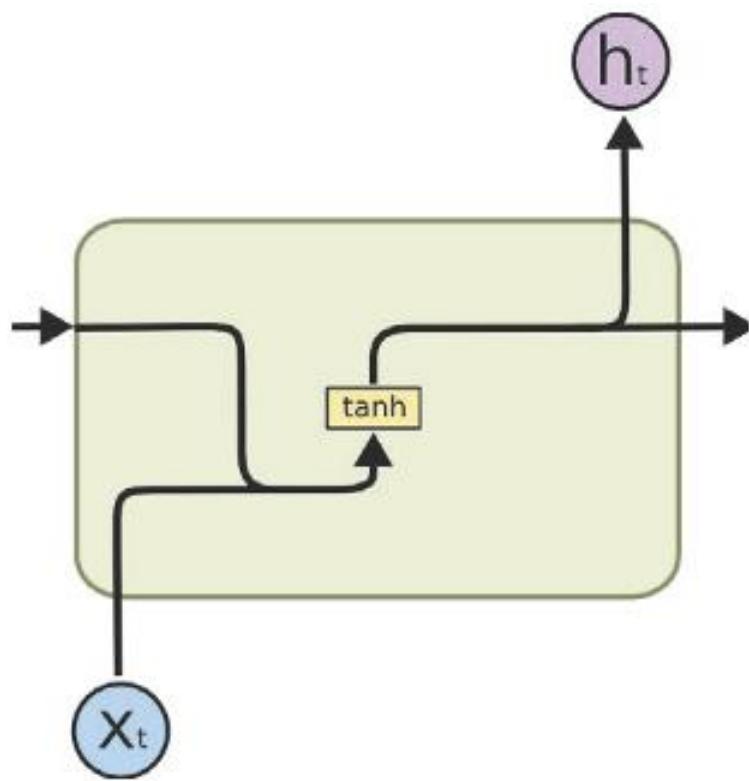


RNN and Image Captioning

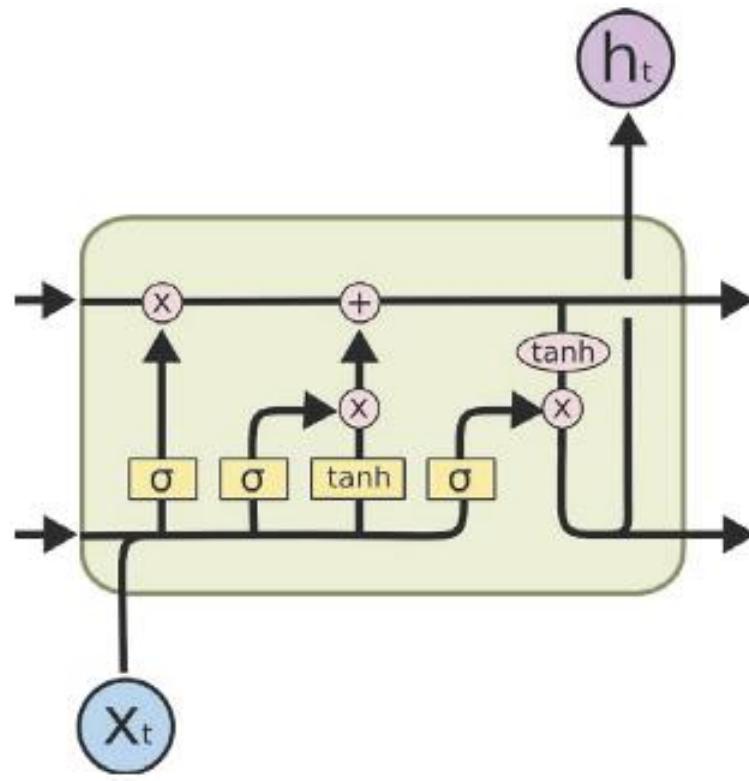
Recurrent Neural Networks



RNN vs. LSTM



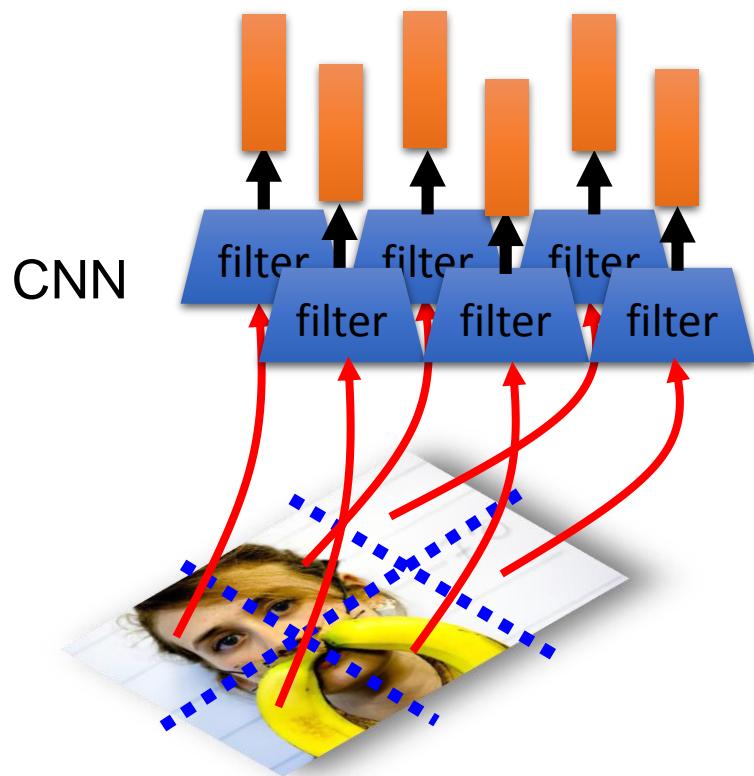
(a) RNN



(b) LSTM

Image Captioning

A vector for each region



z^0 is initial parameter, it is also learned

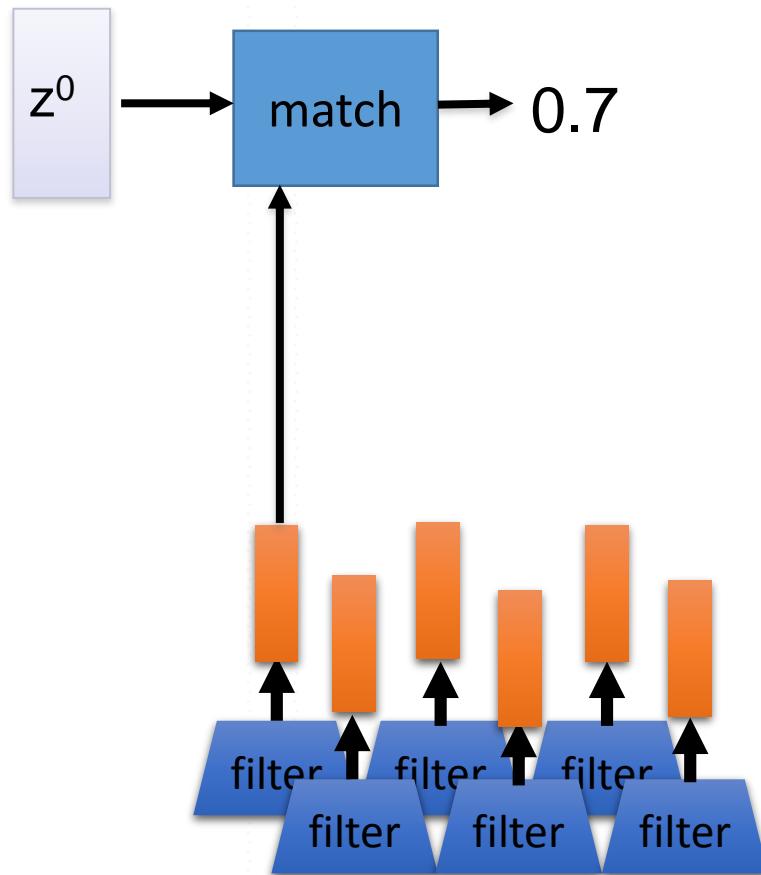


Image Captioning

A vector for each region

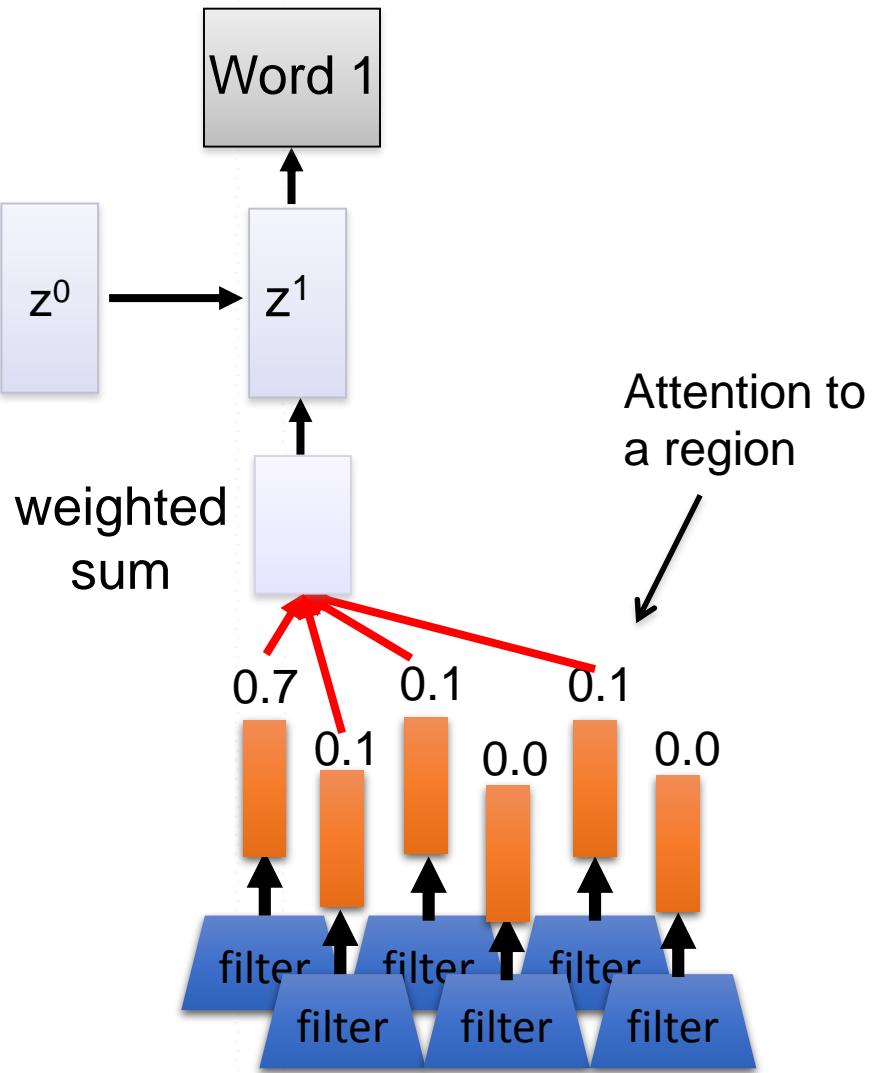
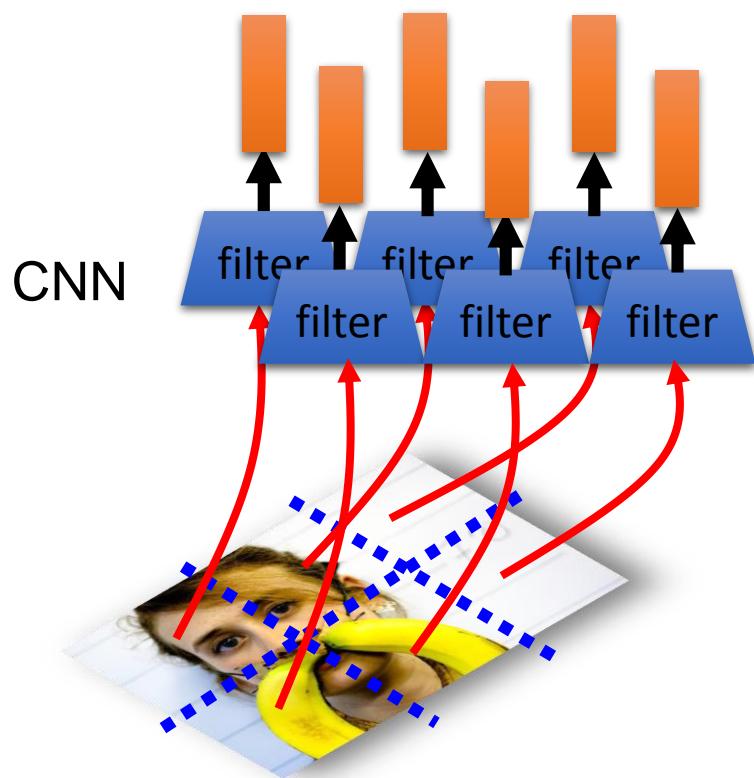


Image Captioning

A vector for each region

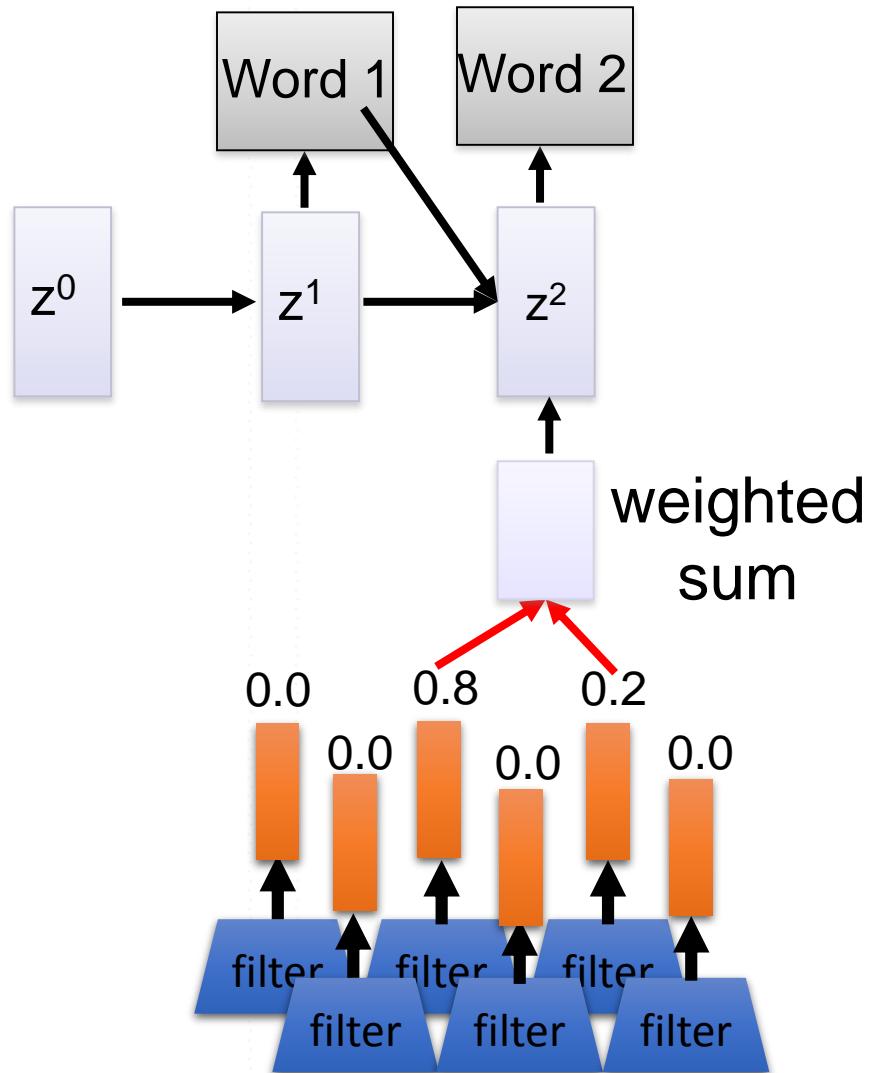
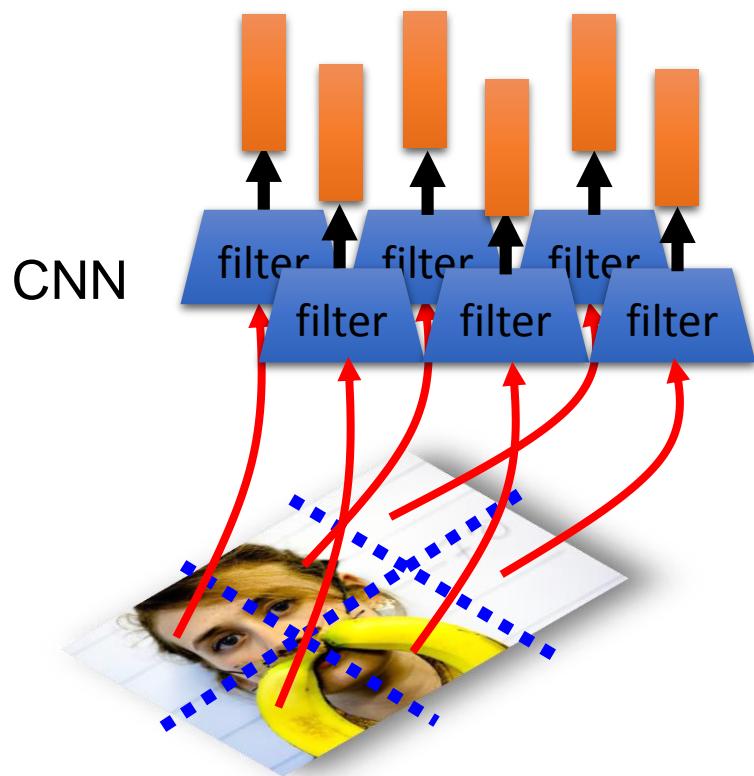


Image Captioning



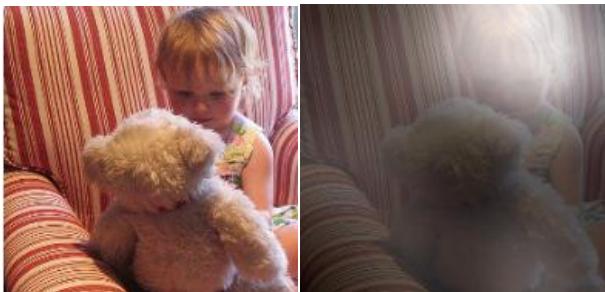
A woman is throwing a frisbee in a park.



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.

Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, Yoshua Bengio, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML, 2015

Image Captioning



A large white bird standing in a forest.



A woman holding a clock in her hand.



A man wearing a hat and a hat on a skateboard.



A person is standing on a beach with a surfboard.



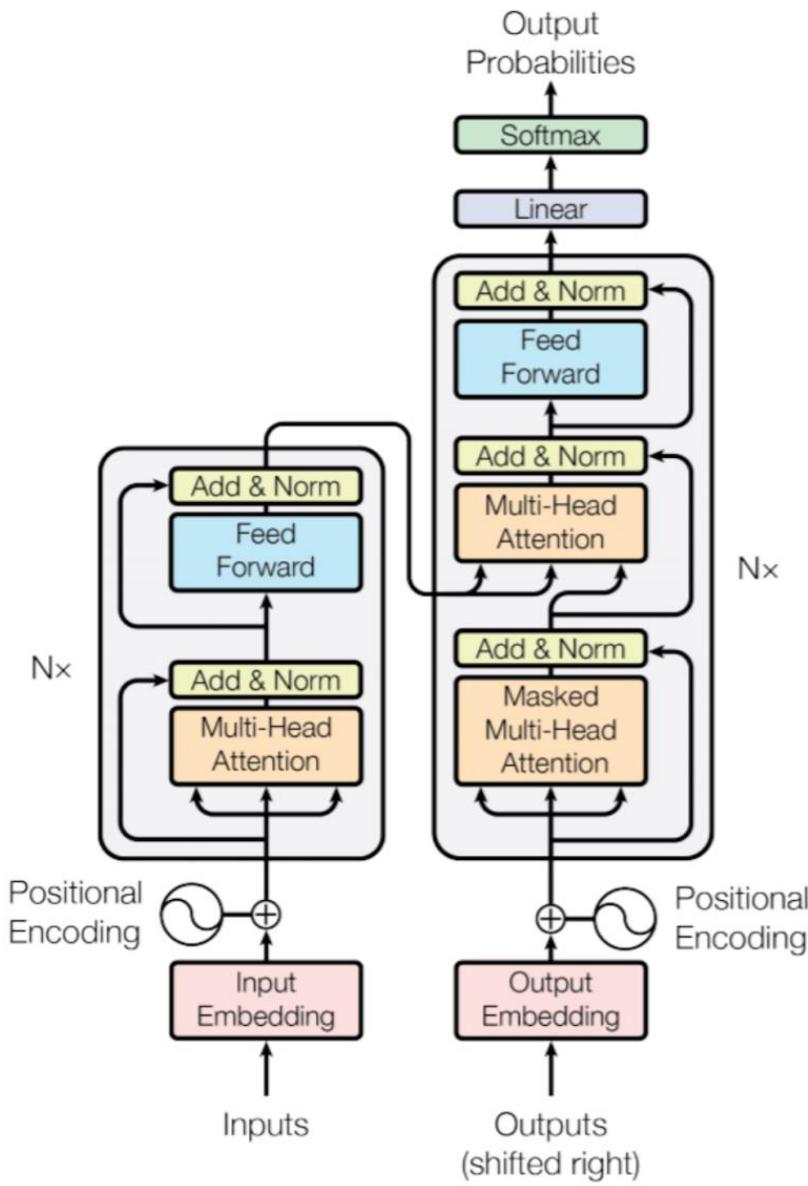
A woman is sitting at a table with a large pizza.



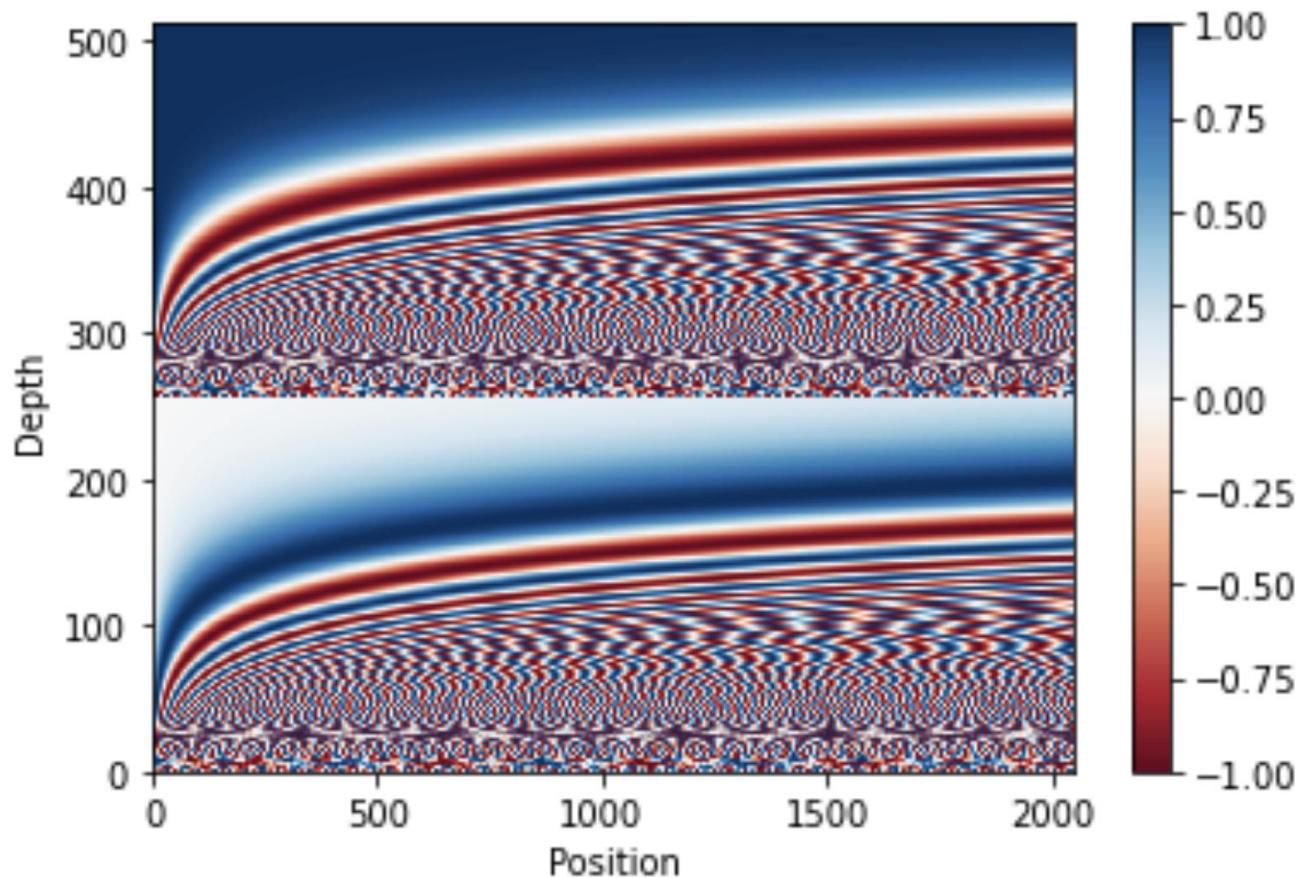
A man is talking on his cell phone while another man watches.

Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, Yoshua Bengio, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML, 2015

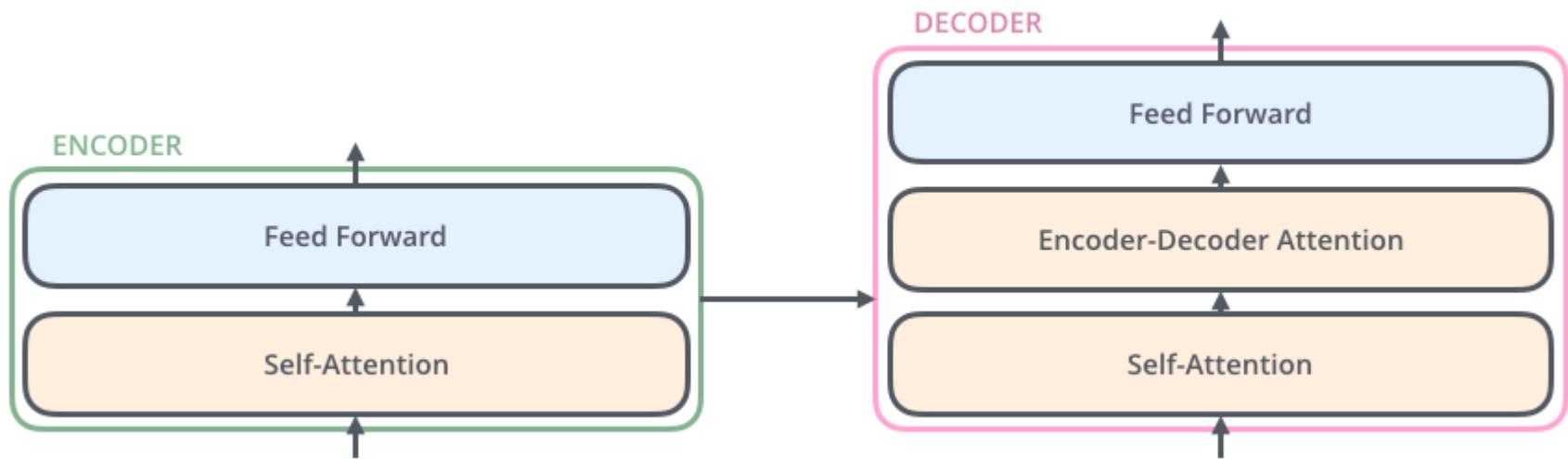
Transformers



Positional Encoding



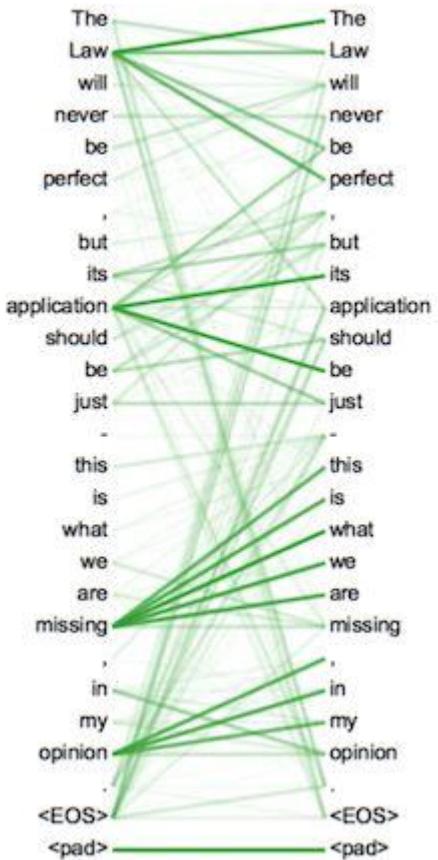
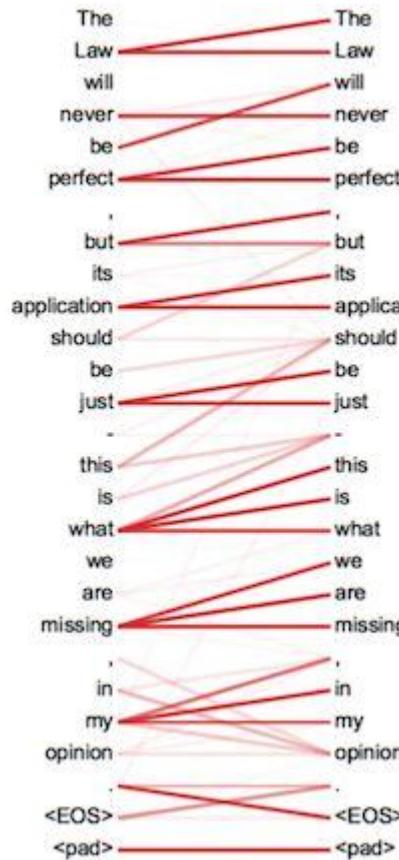
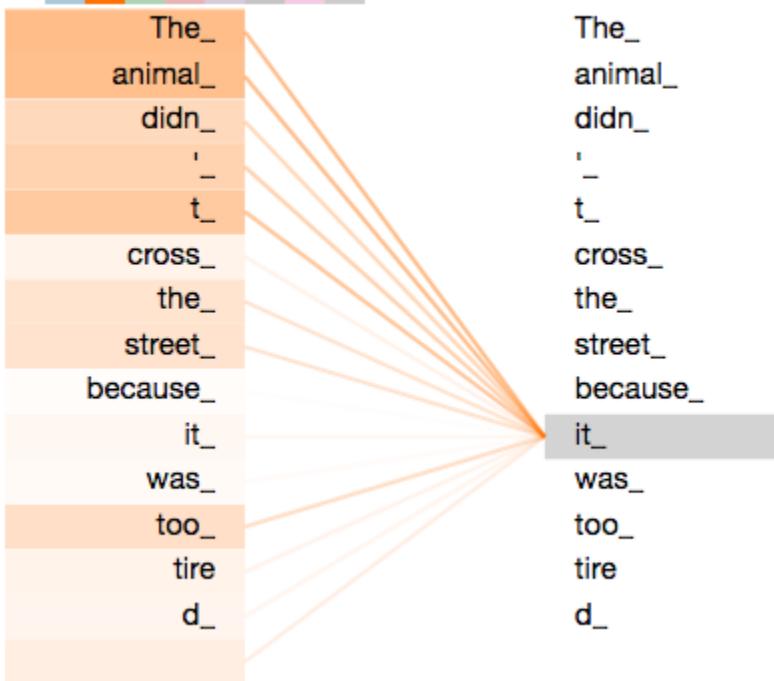
Simplified Version



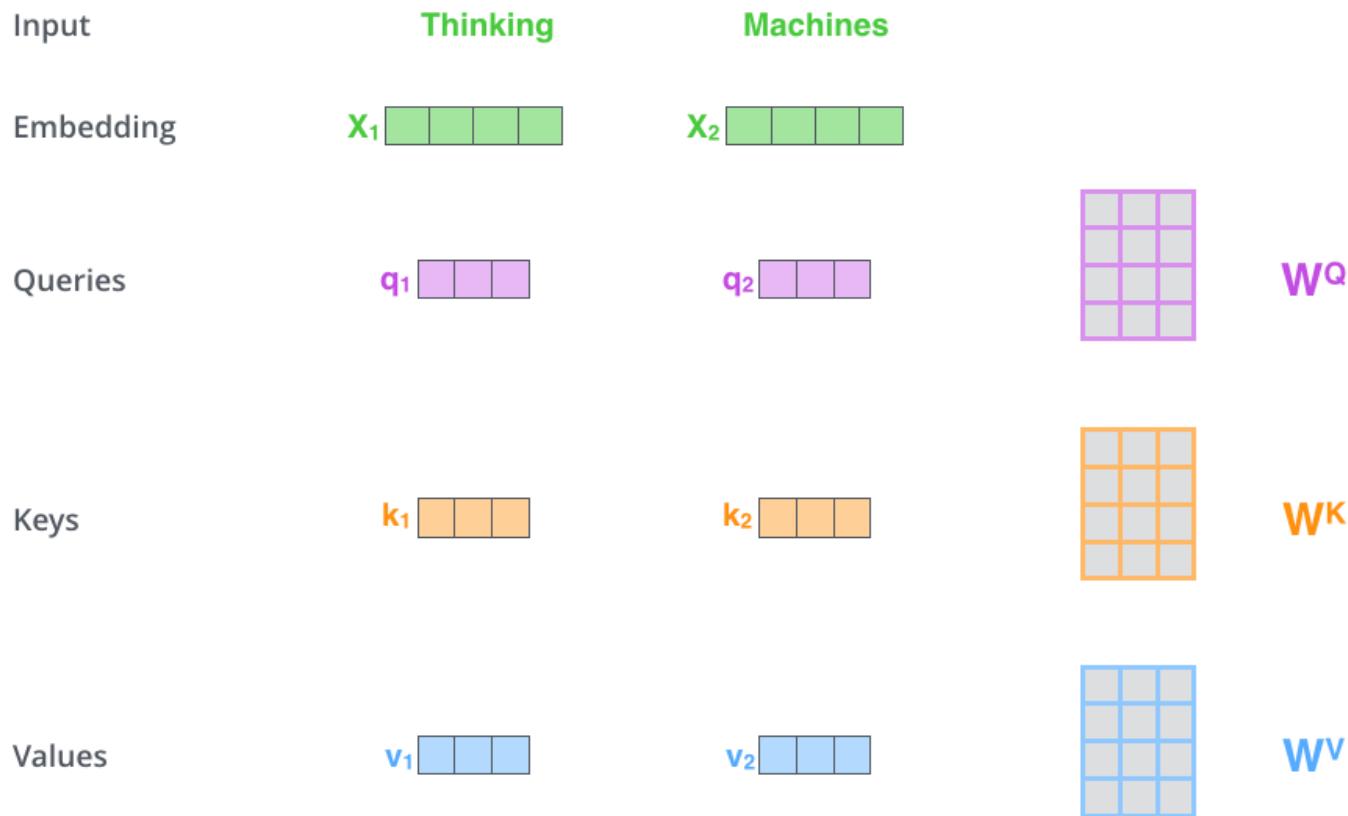
<http://jamalmar.github.io/illustrated-transformer/>

Self Attention

Layer: 5 Attention: Input - Input



Self Attention

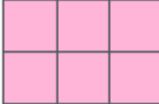


<http://jalammar.github.io/illustrated-transformer/>

Self Attention

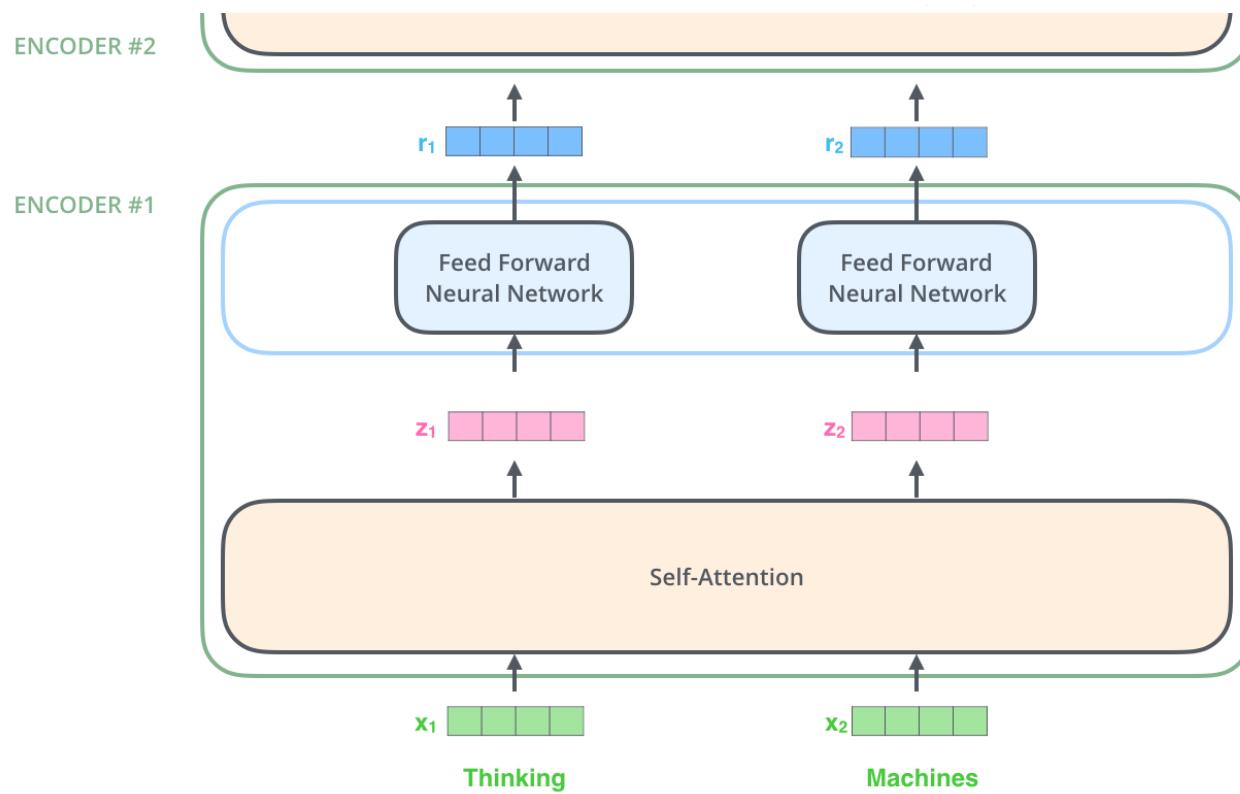
$$\text{softmax} \left(\frac{\begin{matrix} \mathbf{Q} & \times & \mathbf{K}^T \\ \begin{matrix} \text{purple} & \end{matrix} & \times & \begin{matrix} \text{orange} & \end{matrix} \\ \hline \end{matrix}}{\sqrt{d_k}} \right) \mathbf{V}$$

\mathbf{Z}

= 

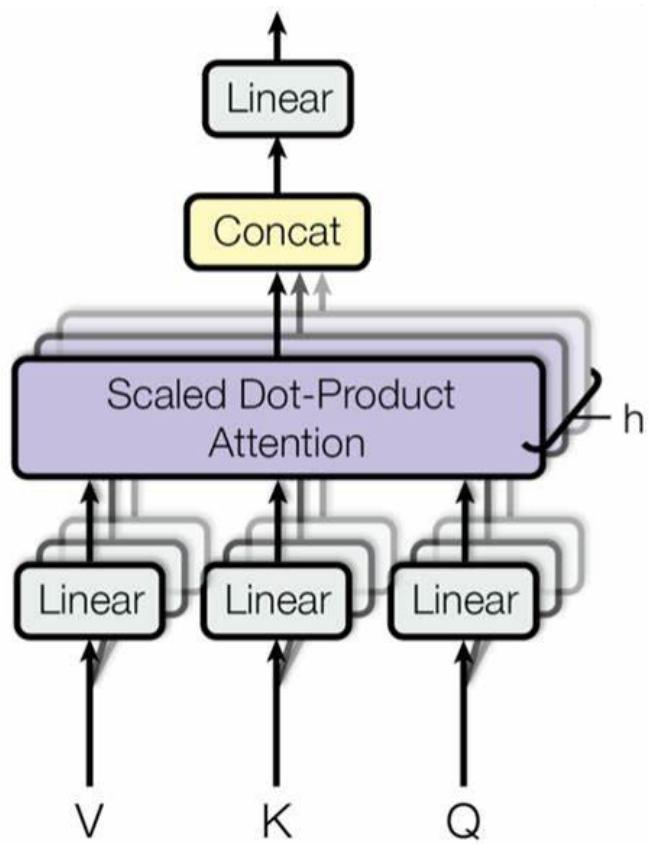
<http://jalammar.github.io/illustrated-transformer/>

Encoder

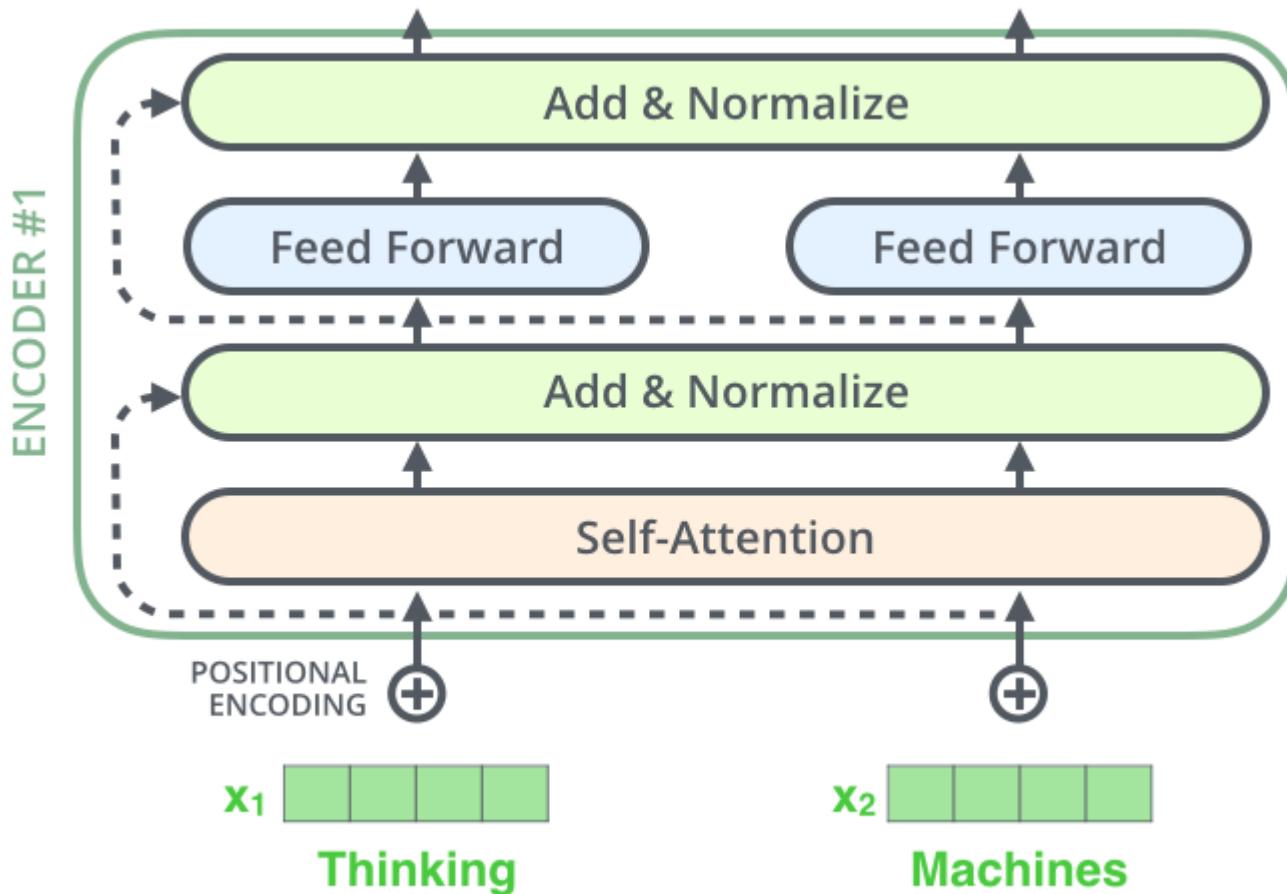


<http://jalamar.github.io/illustrated-transformer/>

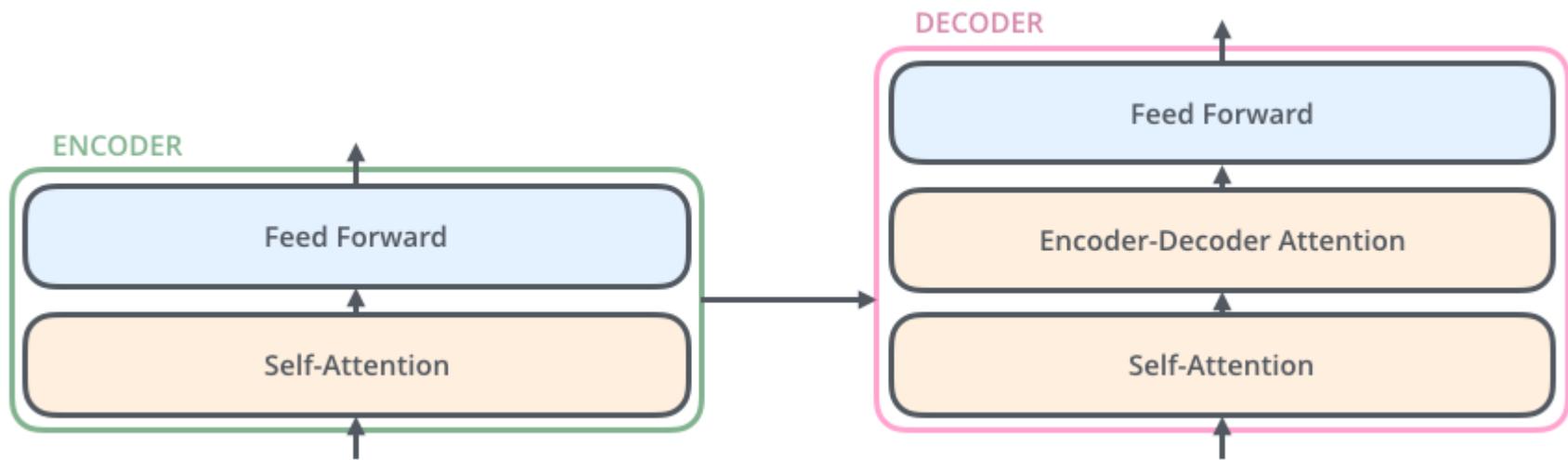
Multi-Head Attention



Residuals



Simplified Version

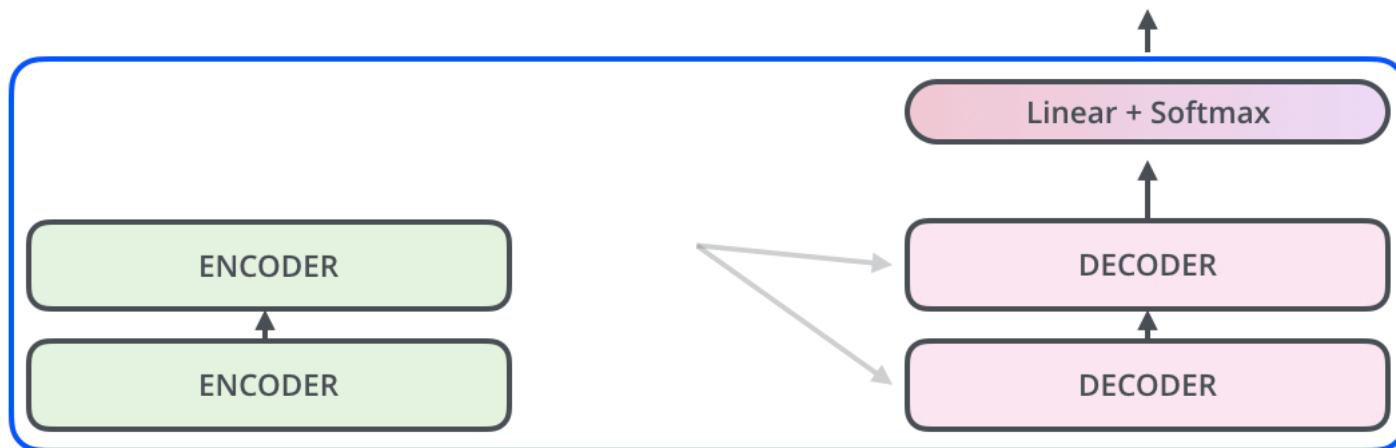


<http://jamalmar.github.io/illustrated-transformer/>

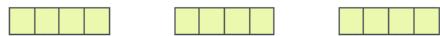
Decoding

Decoding time step: 1 2 3 4 5 6

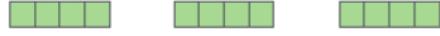
OUTPUT



EMBEDDING
WITH TIME
SIGNAL



EMBEDDINGS

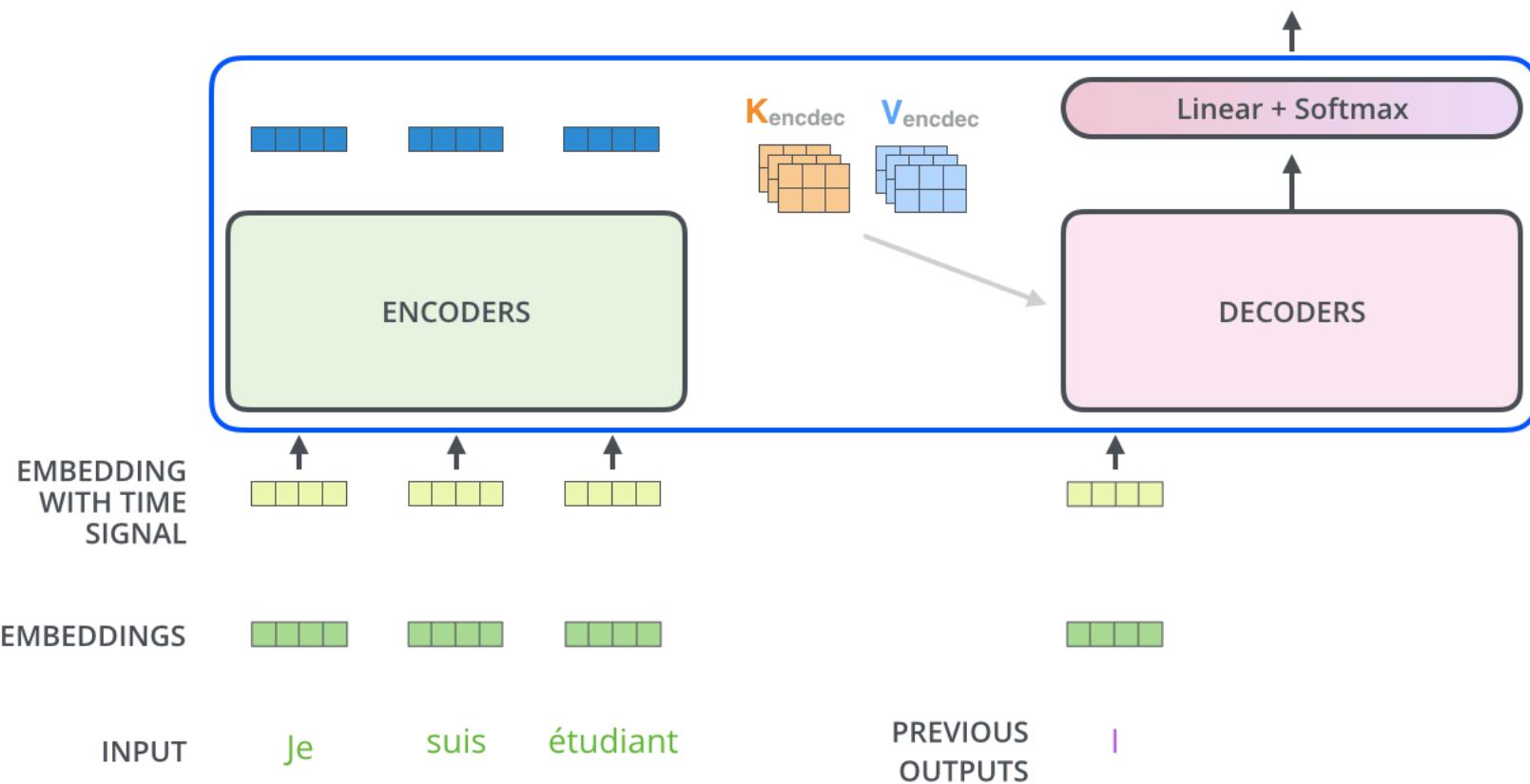


INPUT Je suis étudiant

Decoding

Decoding time step: 1 2 3 4 5 6

OUTPUT |



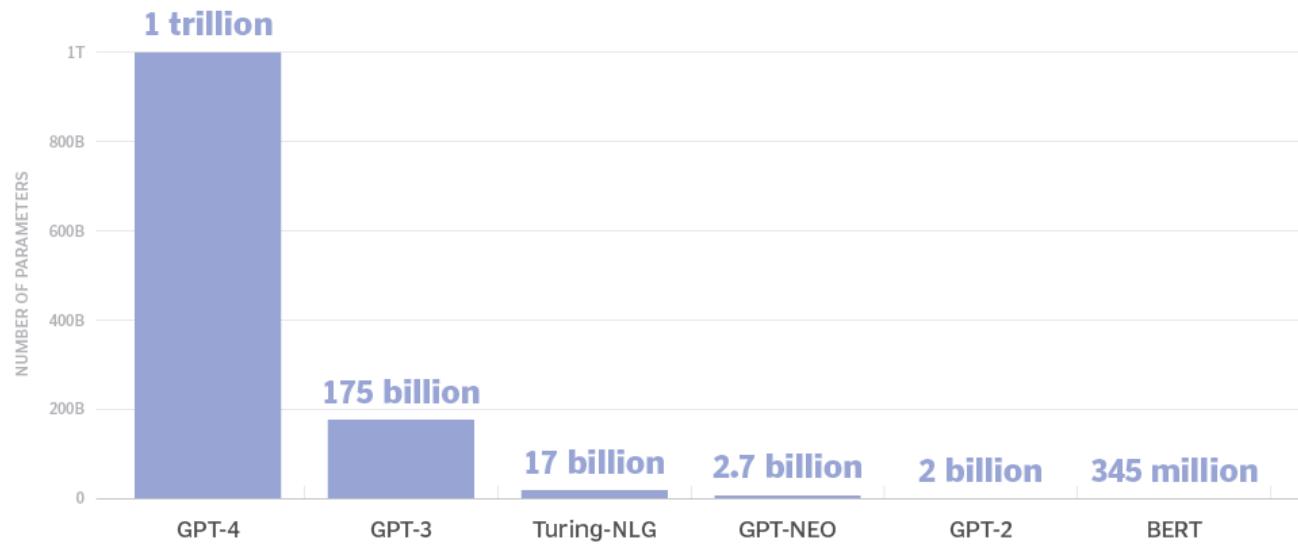
Transformers Models

- >BERT
- >GPT2, GPT3, GPT3.5(ChatGPT), GPT4
- >ViT

Large Language Models

They are LARGE

Parameters of transformer-based language models

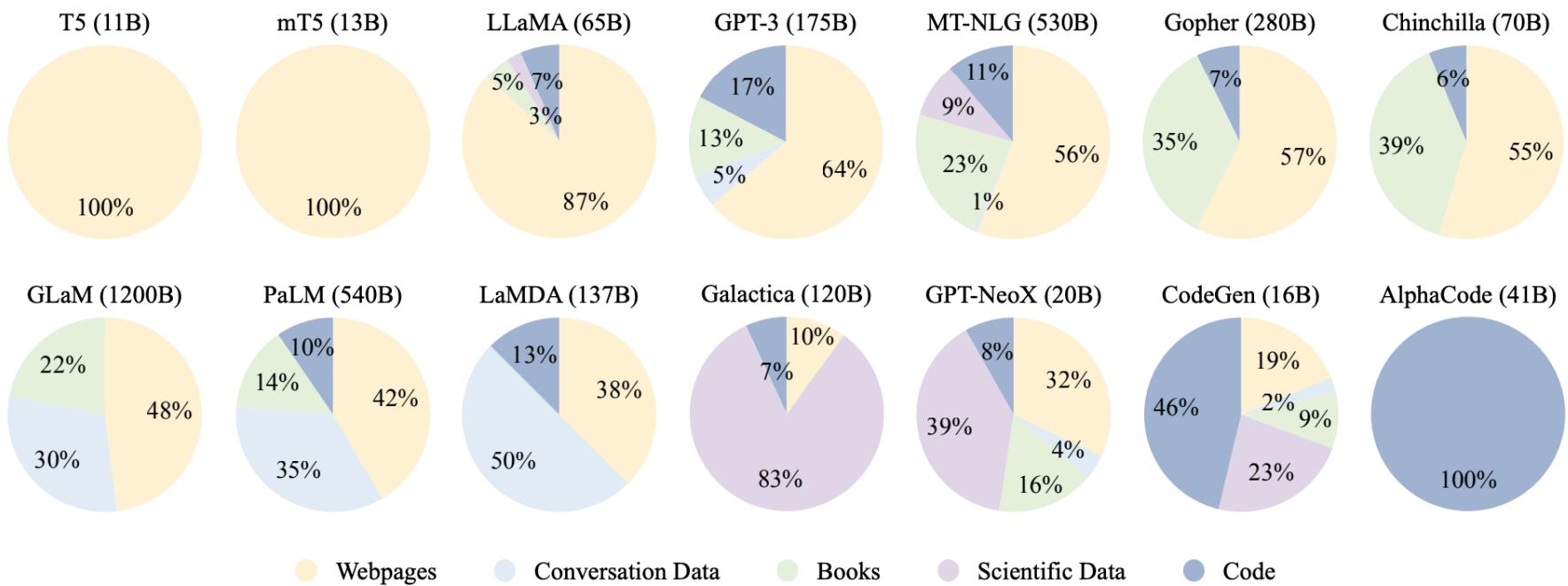


Why are they so LARGE?

>**Emergent abilities** of LLMs are formally defined as “the abilities that are not present in small models but arise in large models”

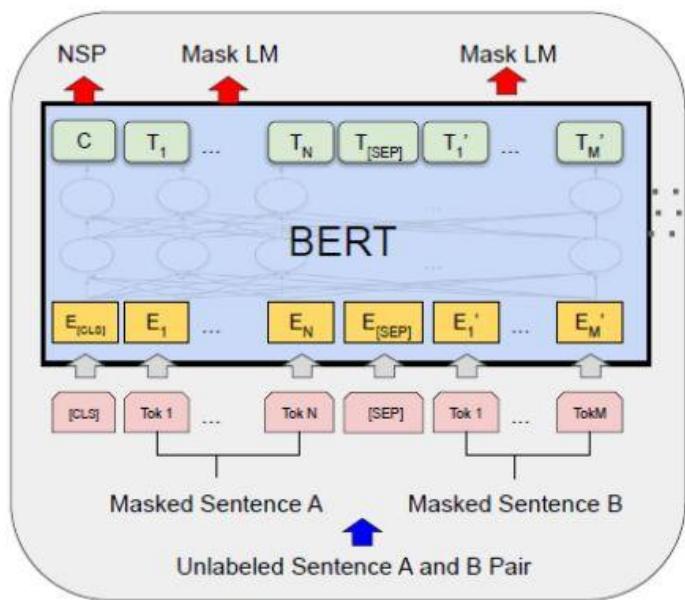
How can they be so LARGE?

>Large corpora: Books, CommonCrawl, Reddit links, Wikipedia, Code, Conversation, ...

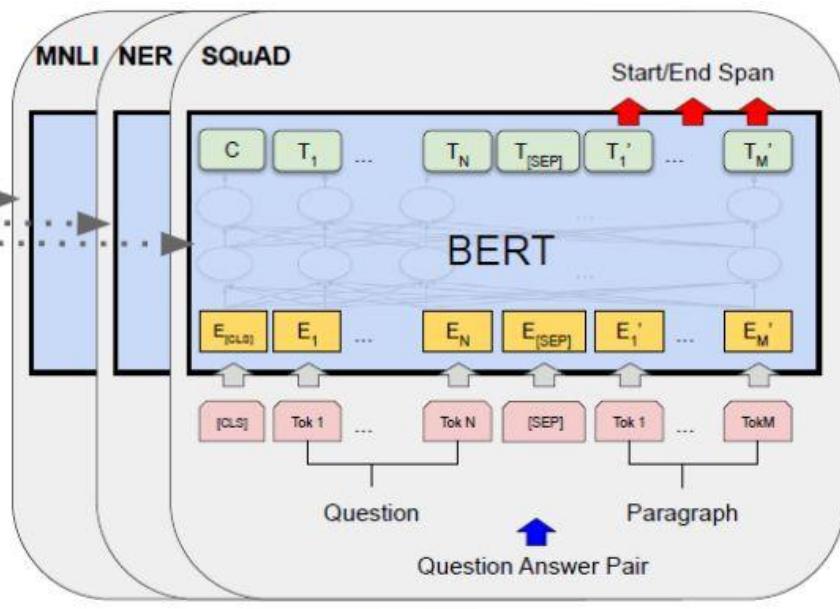


How can they be so LARGE?

> Pre-training: Language modeling, Denoising autoencoding, ...



Pre-training



Fine-Tuning

How can they chat like us?

In-Context Learning

Answer the following mathematical reasoning questions:

N x

Q: If you have 12 candies and you give 4 candies to your friend, how many candies do you have left?

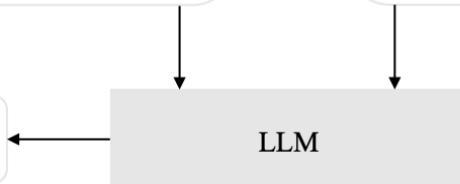
A: The answer is 8.

Q: If a rectangle has a length of 6 cm and a width of 3 cm, what is the perimeter of the rectangle?

A: The answer is 18 cm.

Q: Sam has 12 marbles. He gives 1/4 of them to his sister. How many marbles does Sam have left?

A: The answer is 9.



Chain-of-Thought Prompting

Answer the following mathematical reasoning questions:

N x

Q: If a rectangle has a length of 6 cm and a width of 3 cm, what is the perimeter of the rectangle?

A: For a rectangle, add up the length and width and double it. So, the perimeter of this rectangle is $(6 + 3) \times 2 = 18$ cm.
The answer is 18 cm.

Q: Sam has 12 marbles. He gives 1/4 of them to his sister. How many marbles does Sam have left?

A: He gives $(1 / 4) \times 12 = 3$ marbles. So Sam is left with $12 - 3 = 9$ marbles.
The answer is 9.

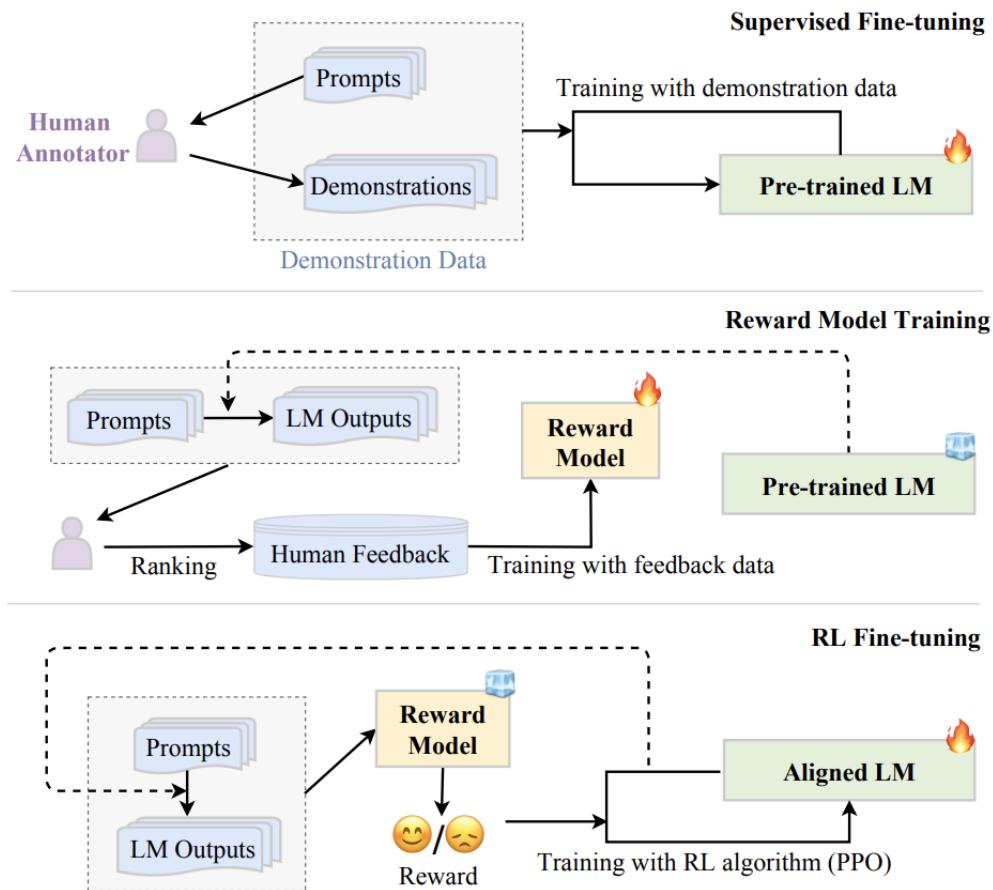
: Task description

: Demonstration

: Chain-of-Thought

: Query

How can they chat like us?

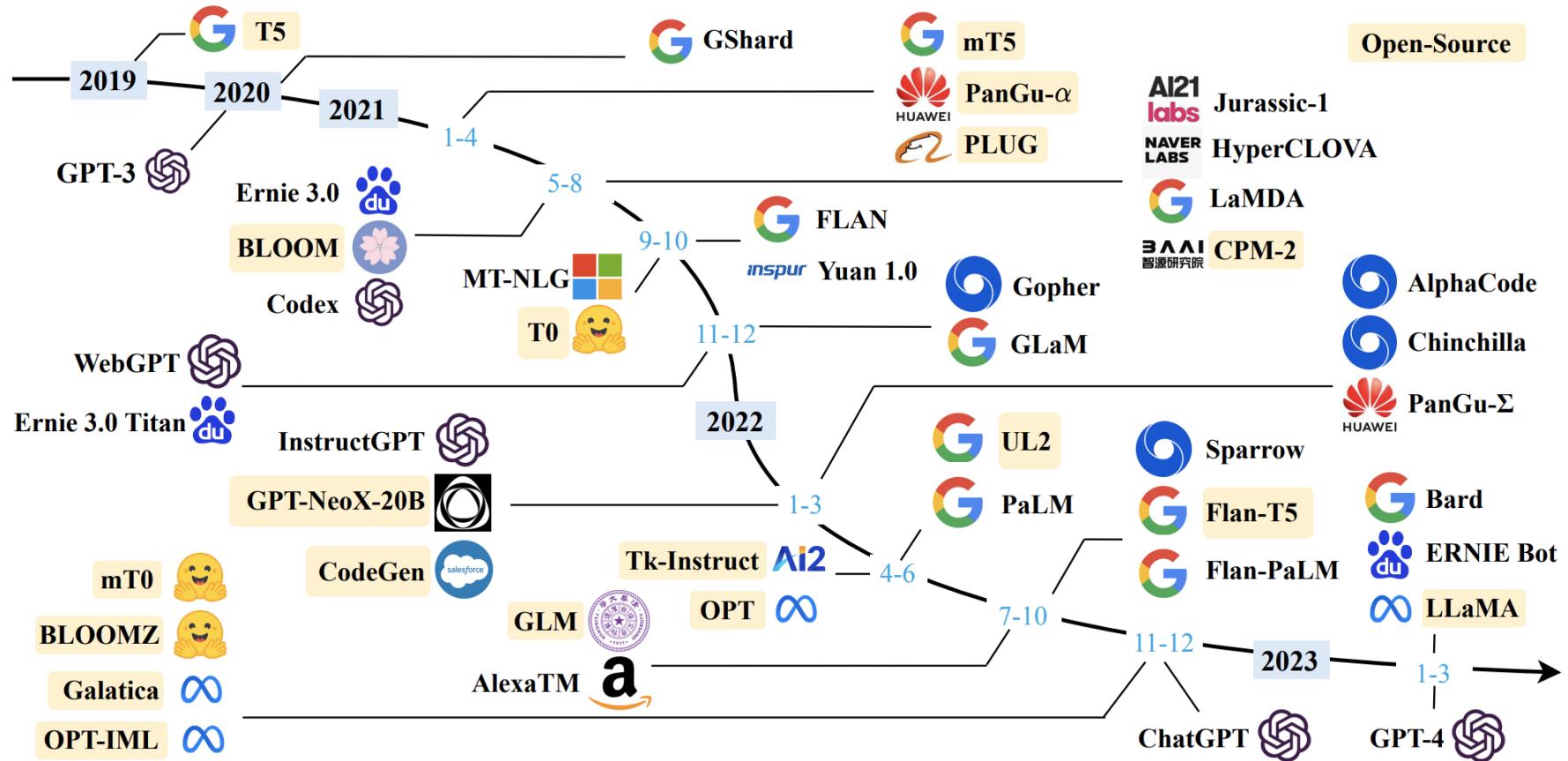


Are they just chatbots?

- > **Text generation.** The ability to generate text on any topic that the LLM has been trained on is a primary use case.
- > **Translation.** For LLMs trained on multiple languages, the ability to translate from one language to another is a common feature.
- > **Content summary.** Summarizing blocks or multiple pages of text is a useful function of LLMs.

Are they just chatbots?

- >**Rewriting content.** Rewriting a section of text is another capability.
- >**Classification and categorization.** An LLM is able to classify and categorize content.
- >**Sentiment analysis.** Most LLMs can be used for sentiment analysis to help users to better understand the intent of a piece of content or a particular response.
- >**Conversational AI and chatbots.** LLMs can enable a conversation with a user in a way that is typically more natural than older generations of AI technologies.



What can WE do with them?

GPT-4

GPT-4 can accept images as inputs and generate captions, classifications, and analyses.

Input

What can I make with these ingredients?



Output

There are many options for what you can make with these ingredients! Some possibilities include:

- Pancakes or waffles
- Crepes
- French toast
- Omelette or frittata
- Quiche
- Custard or pudding
- Cake or cupcakes
- Muffins or bread- Cookies or biscuits

These are just a few examples, but the possibilities are endless!

Visual ChatGPT

Enter text and press enter, or upload an image

I

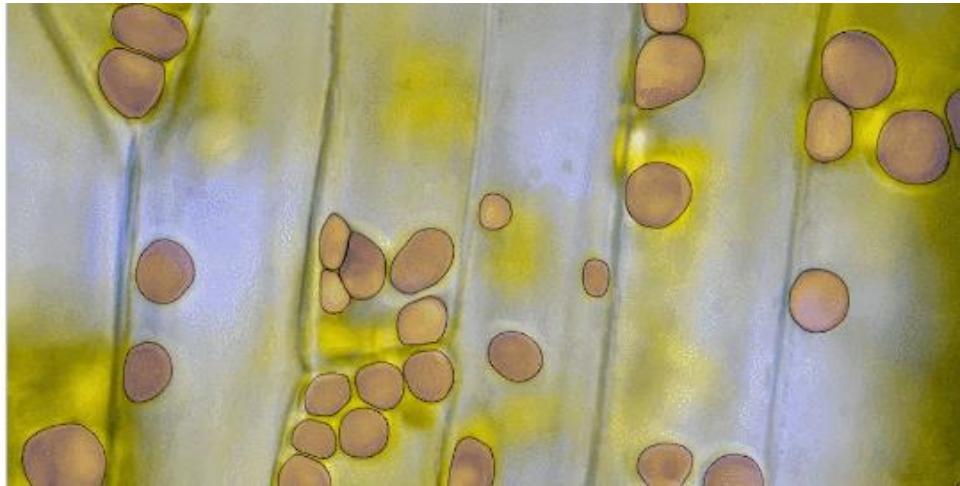
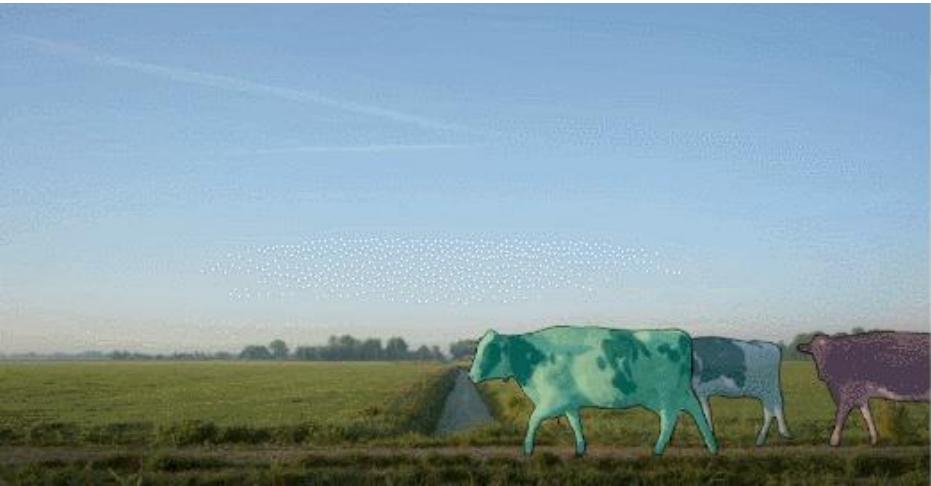
Clear

Upload

Meta - Segment Anything



Meta - Segment Anything

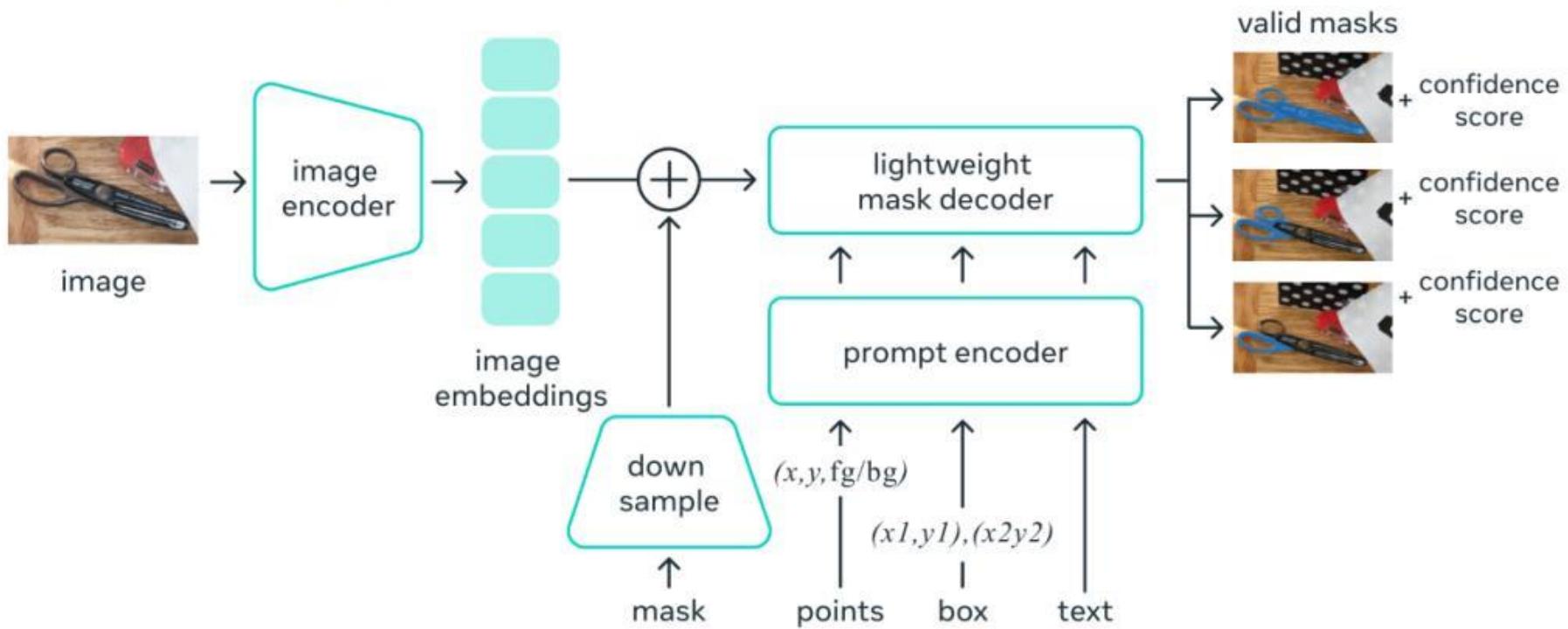


Meta - Segment Anything

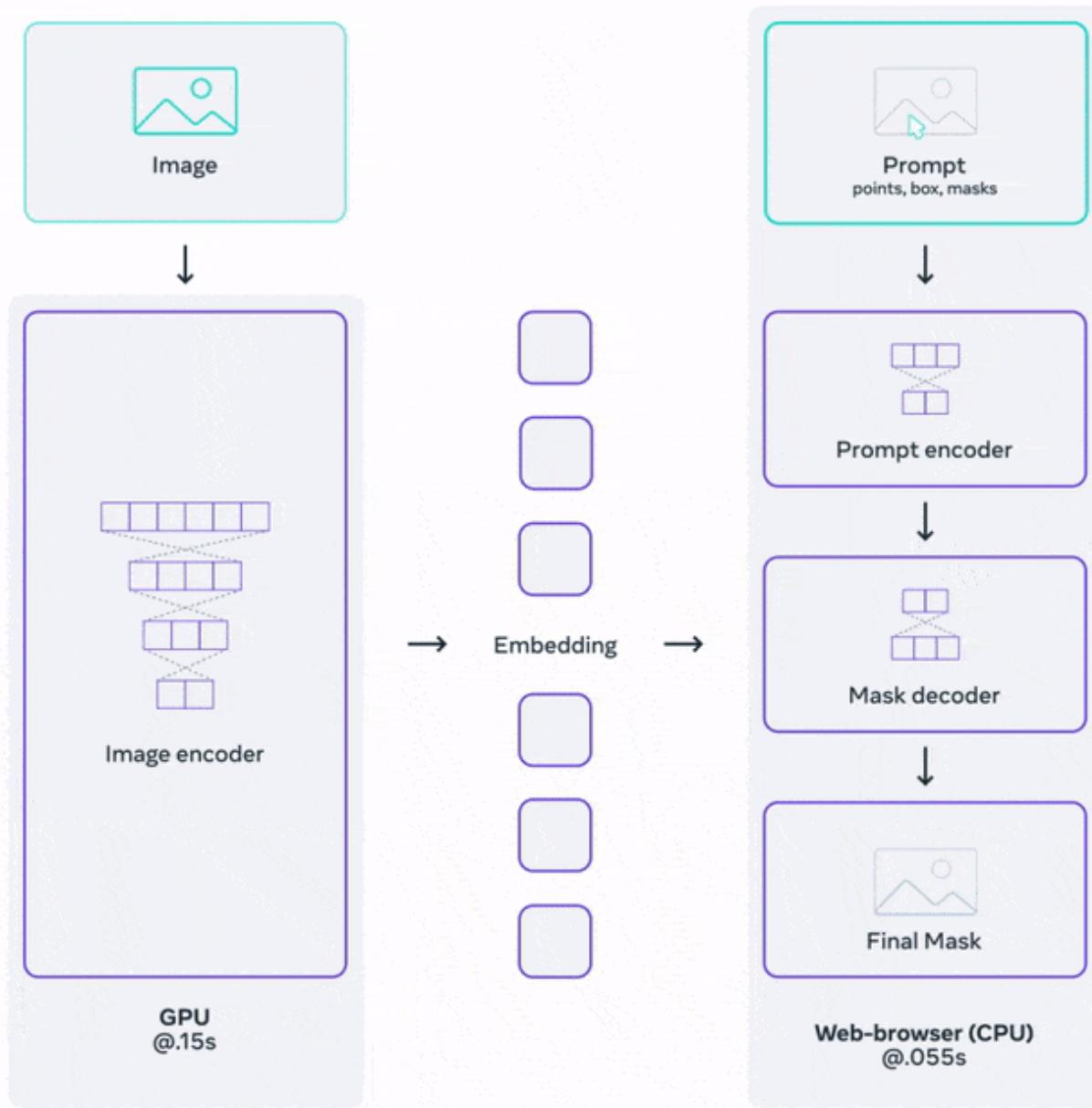


Meta - Segment Anything

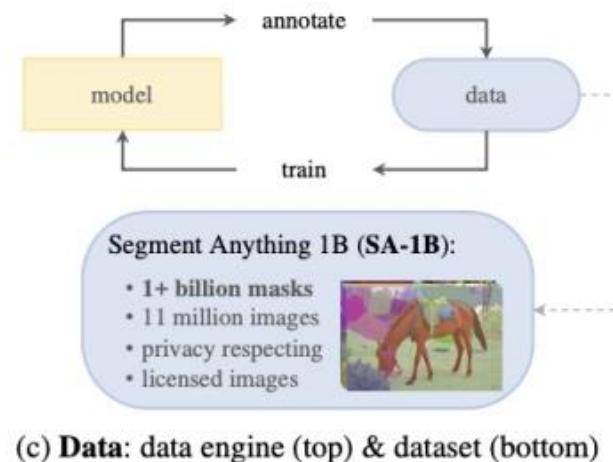
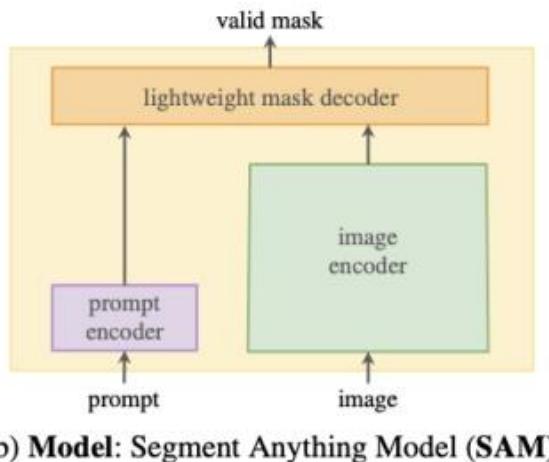
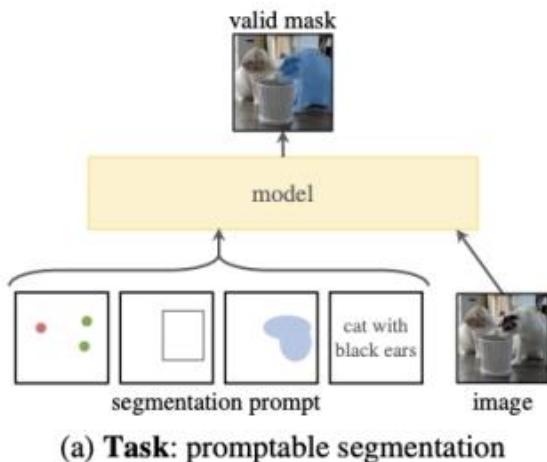
Universal segmentation model



Meta SAM



Meta - Segment Anything





Thank you!

