

# Introduction to Large Language Models

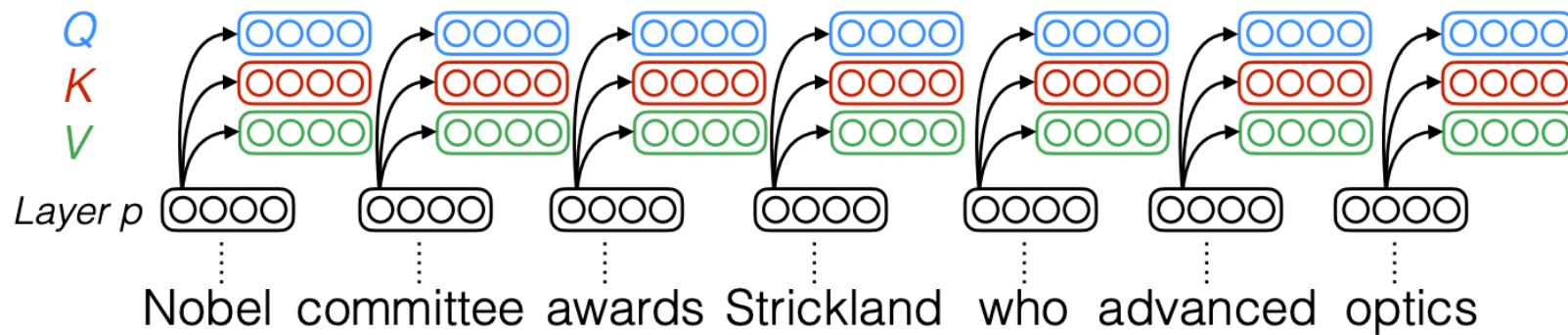
## **CONTENT**

- Attention mechanism
  - Model structure
- Pre-training & IT & RLHF
  - Scaling
  - Multi Modal LLM

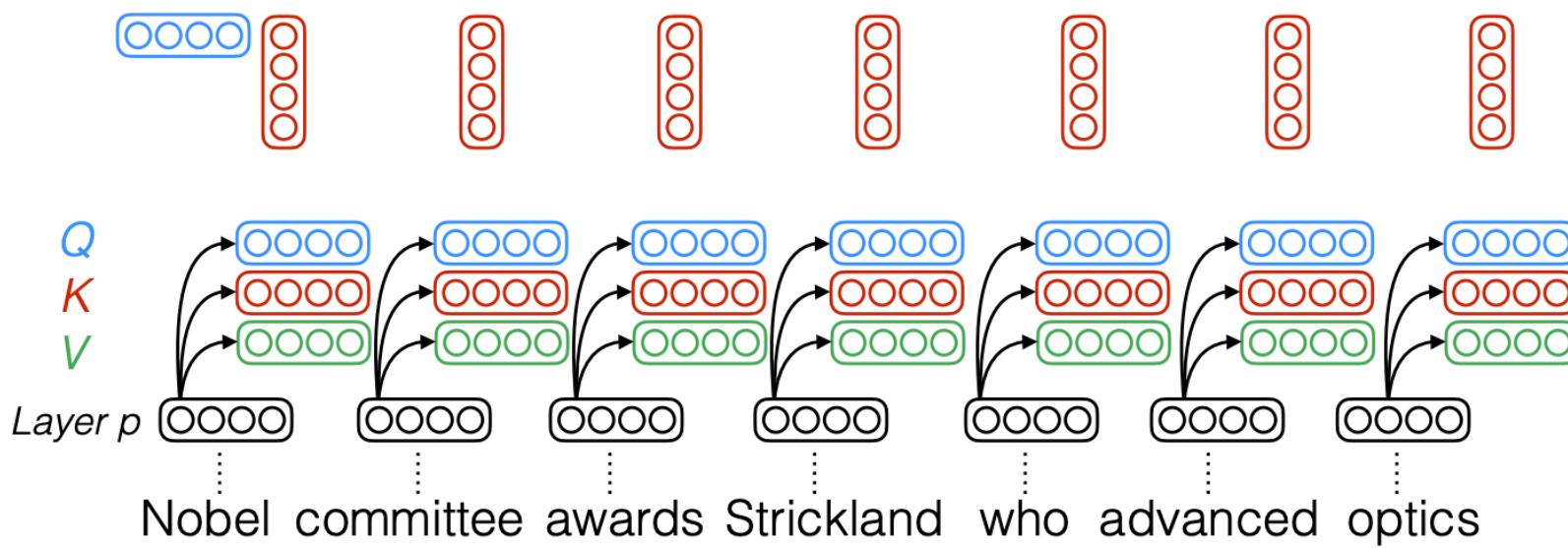


# Attention mechanism

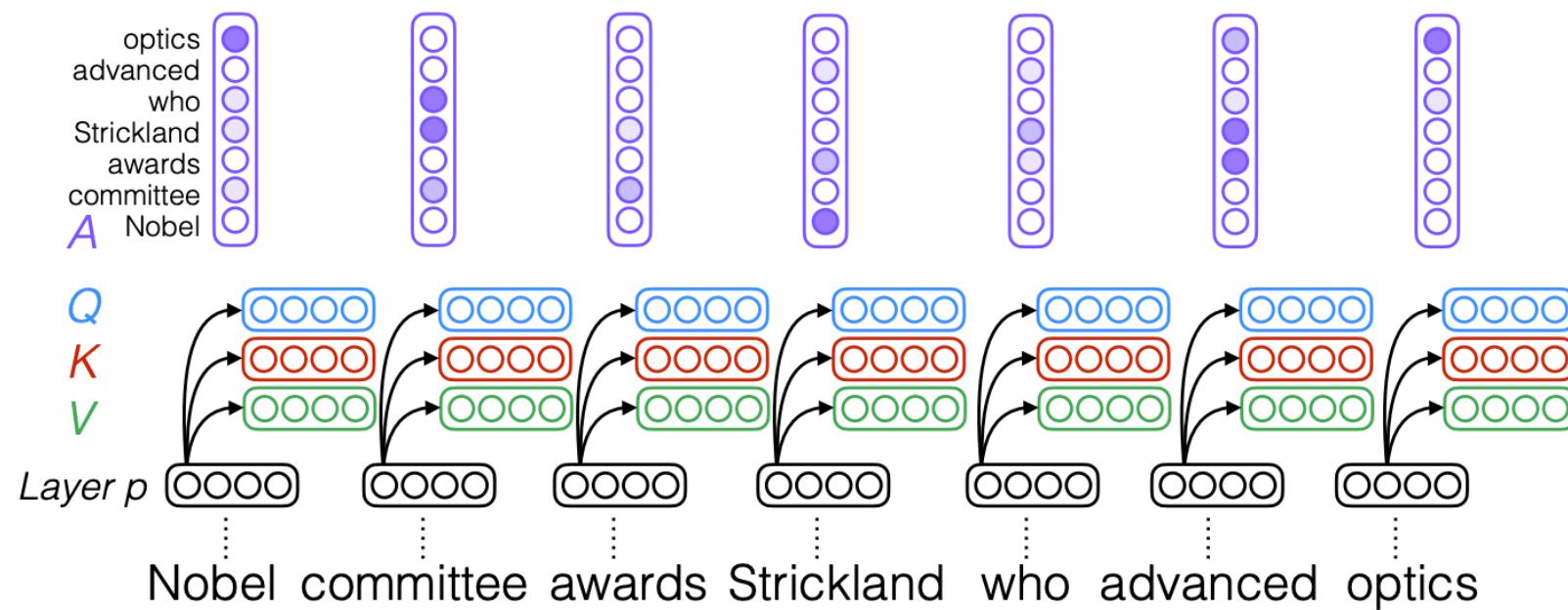
## Attention mechanism



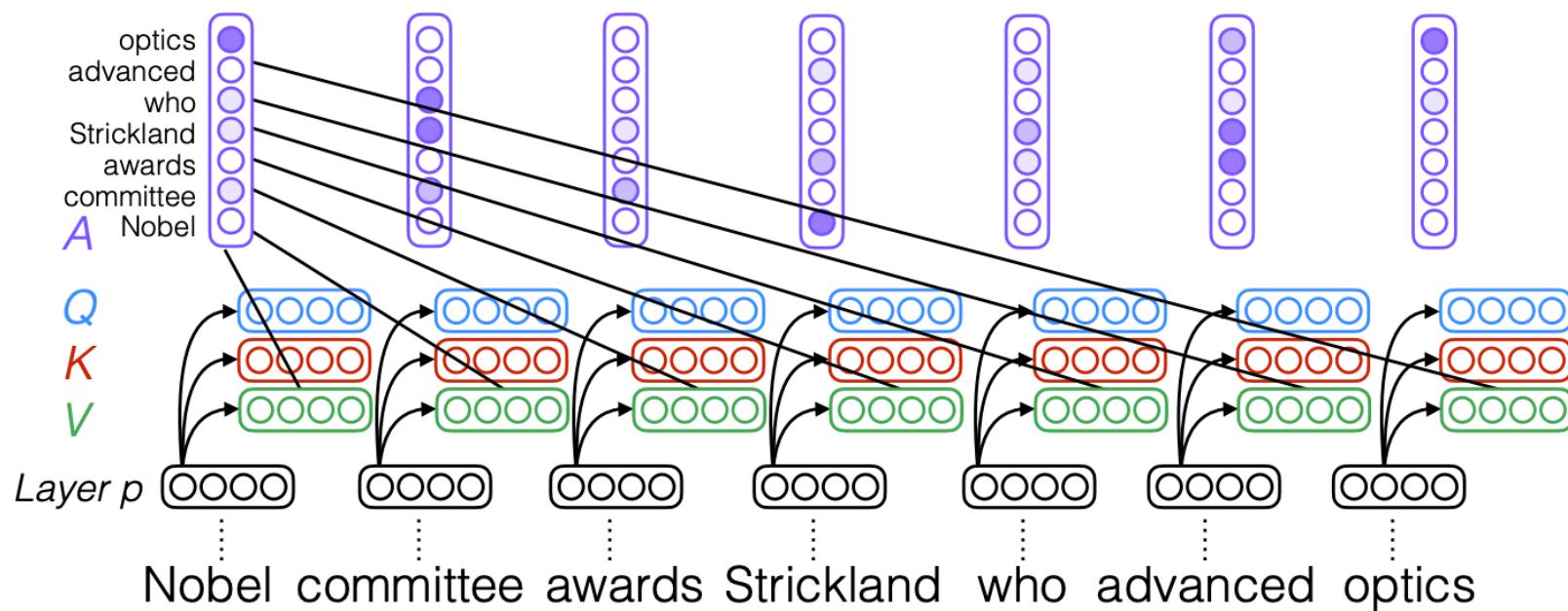
## Attention mechanism



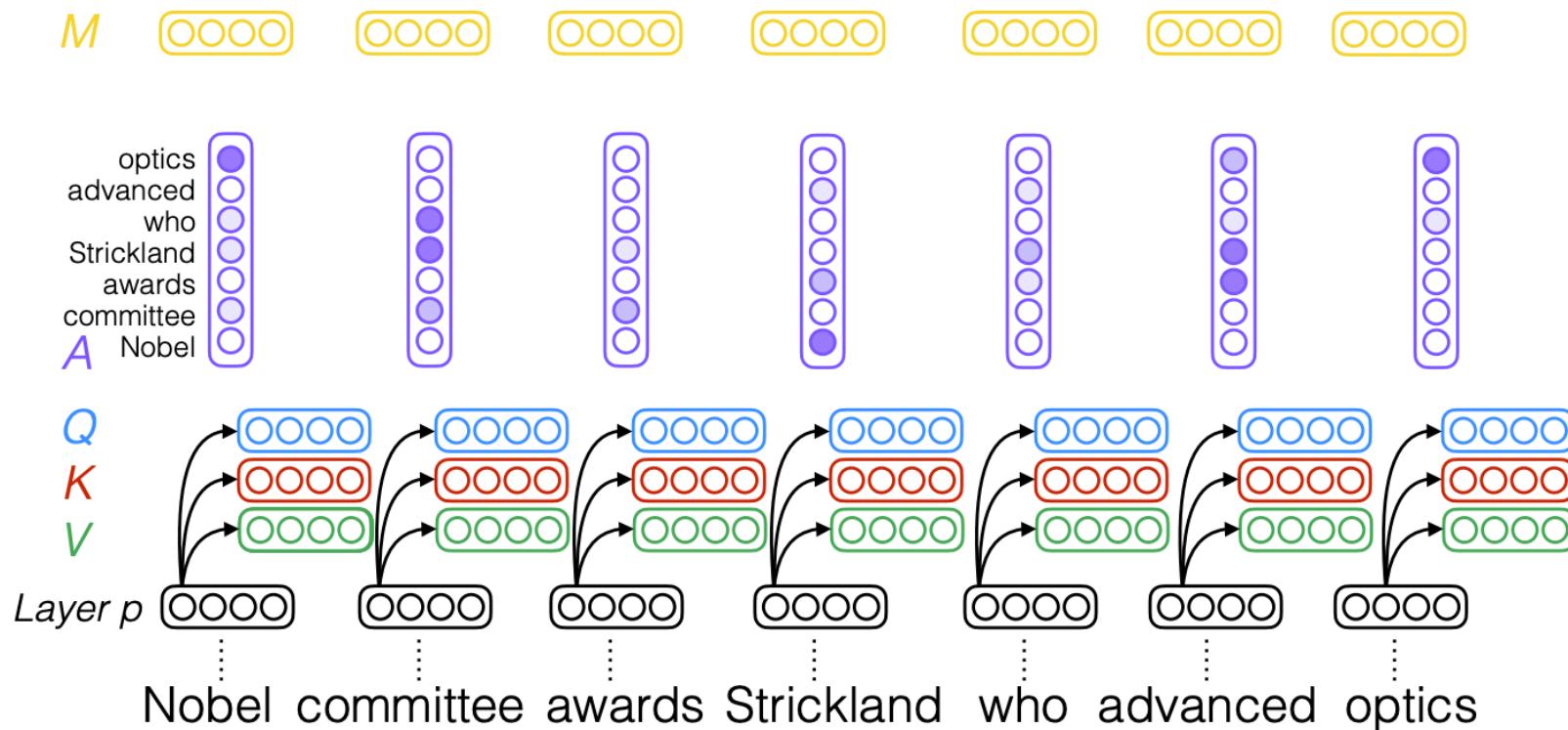
## Attention mechanism



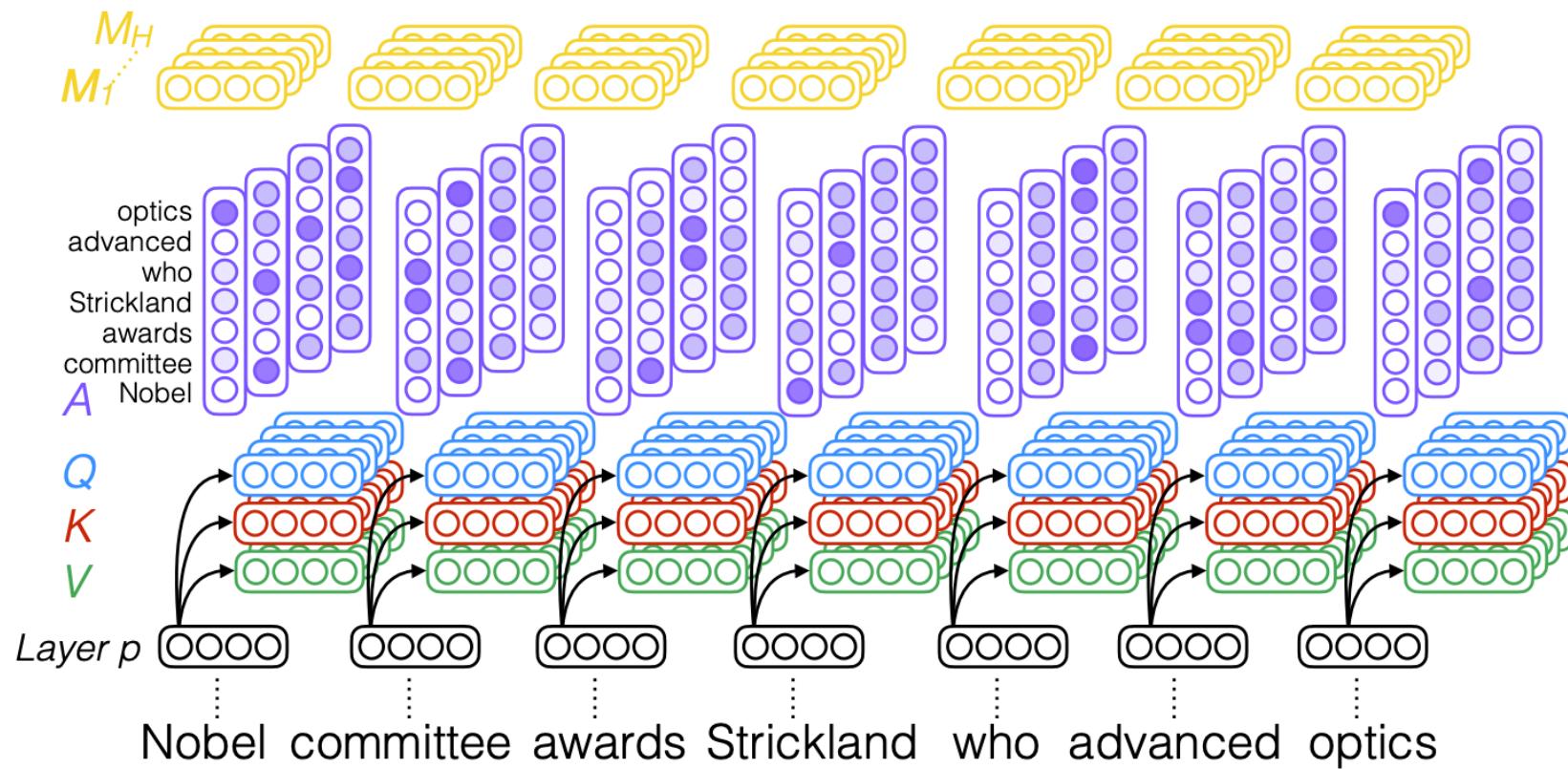
## Attention mechanism



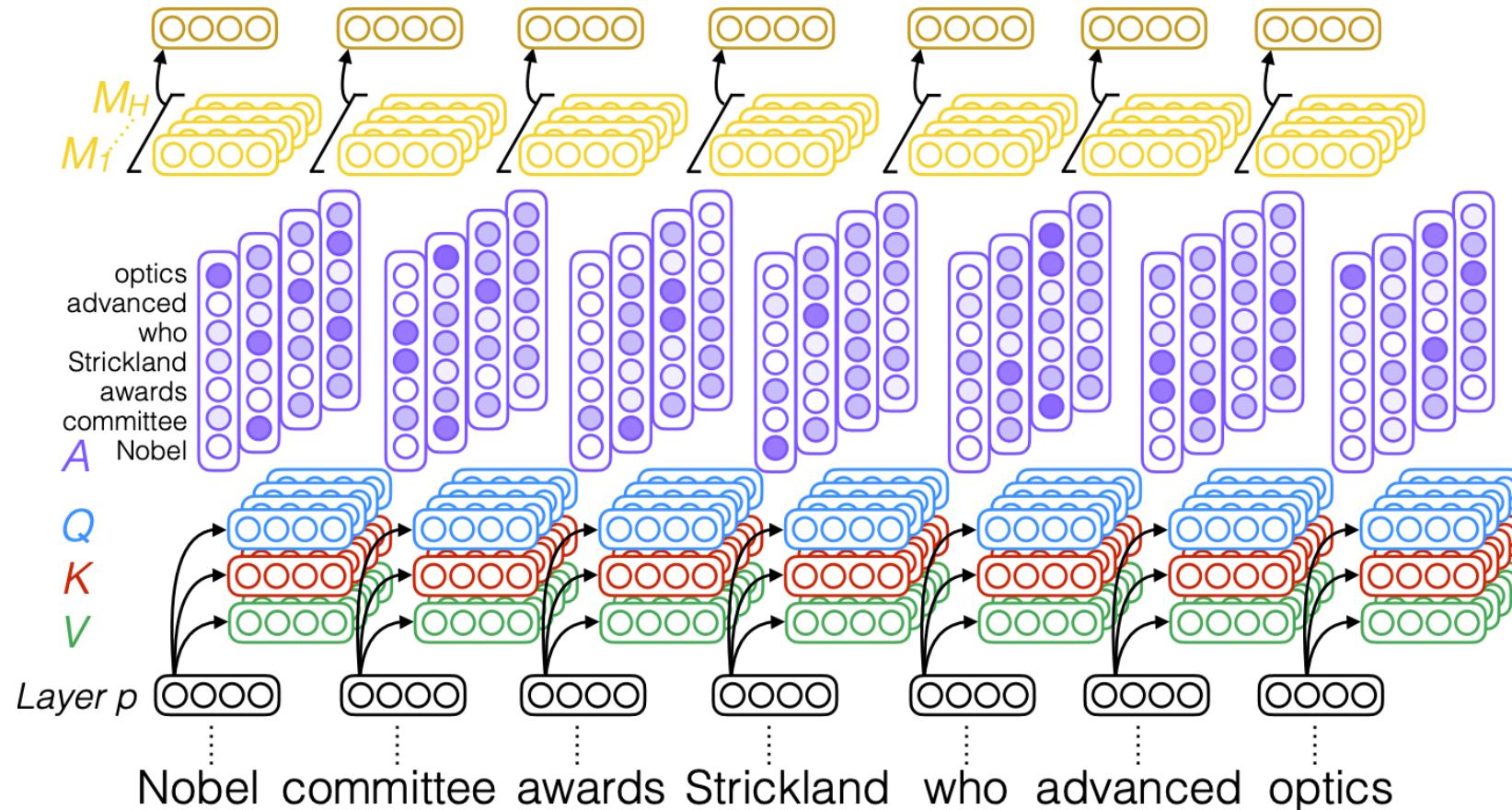
## Attention mechanism



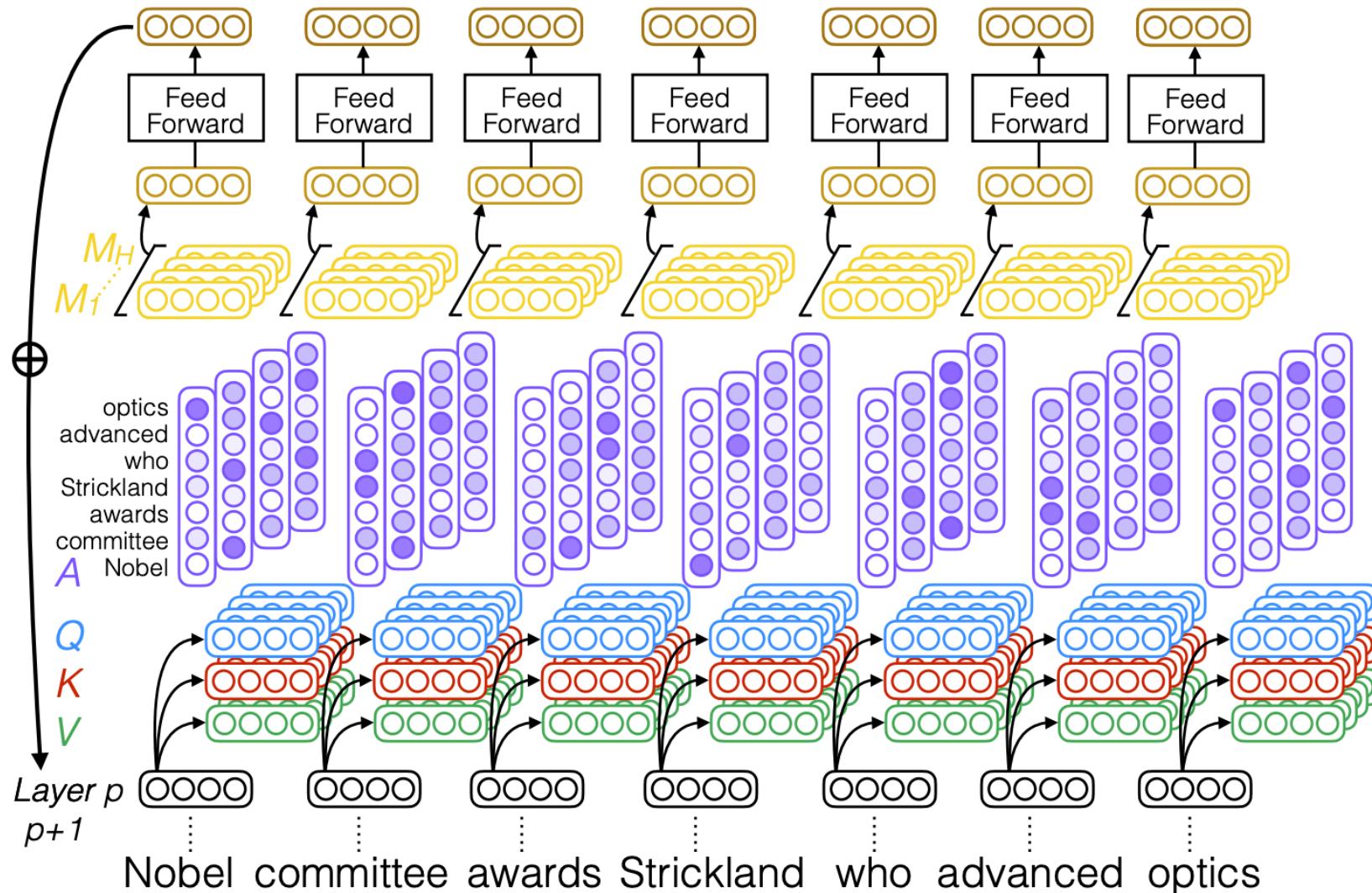
## Attention mechanism



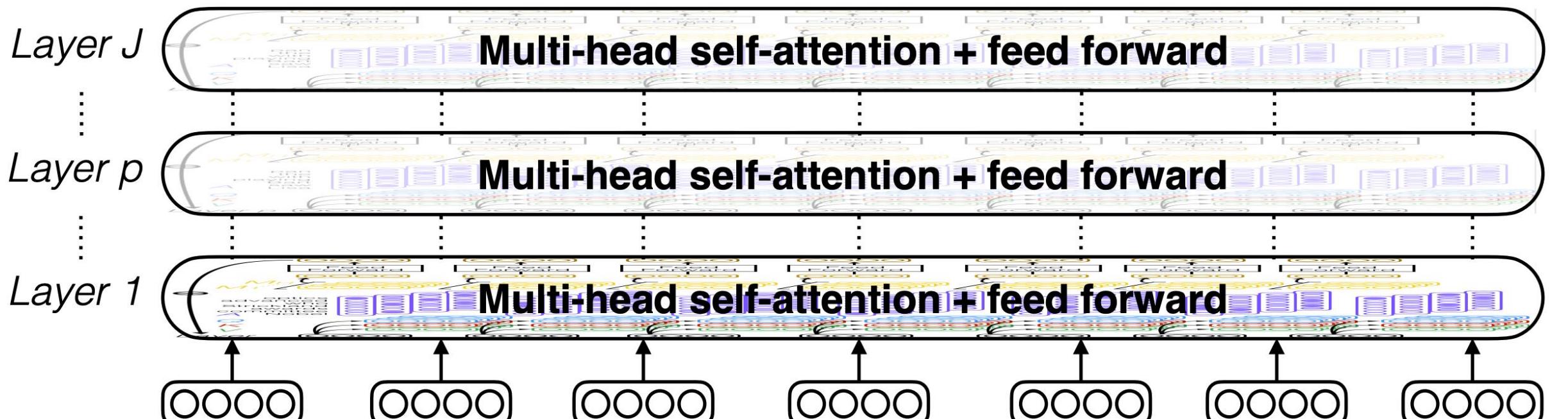
## Attention mechanism



## Attention mechanism



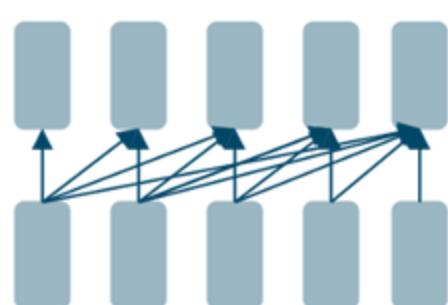
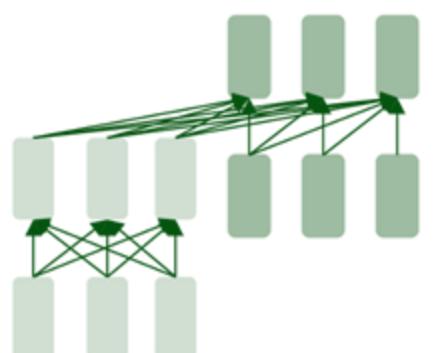
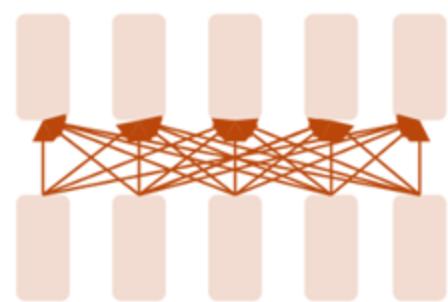
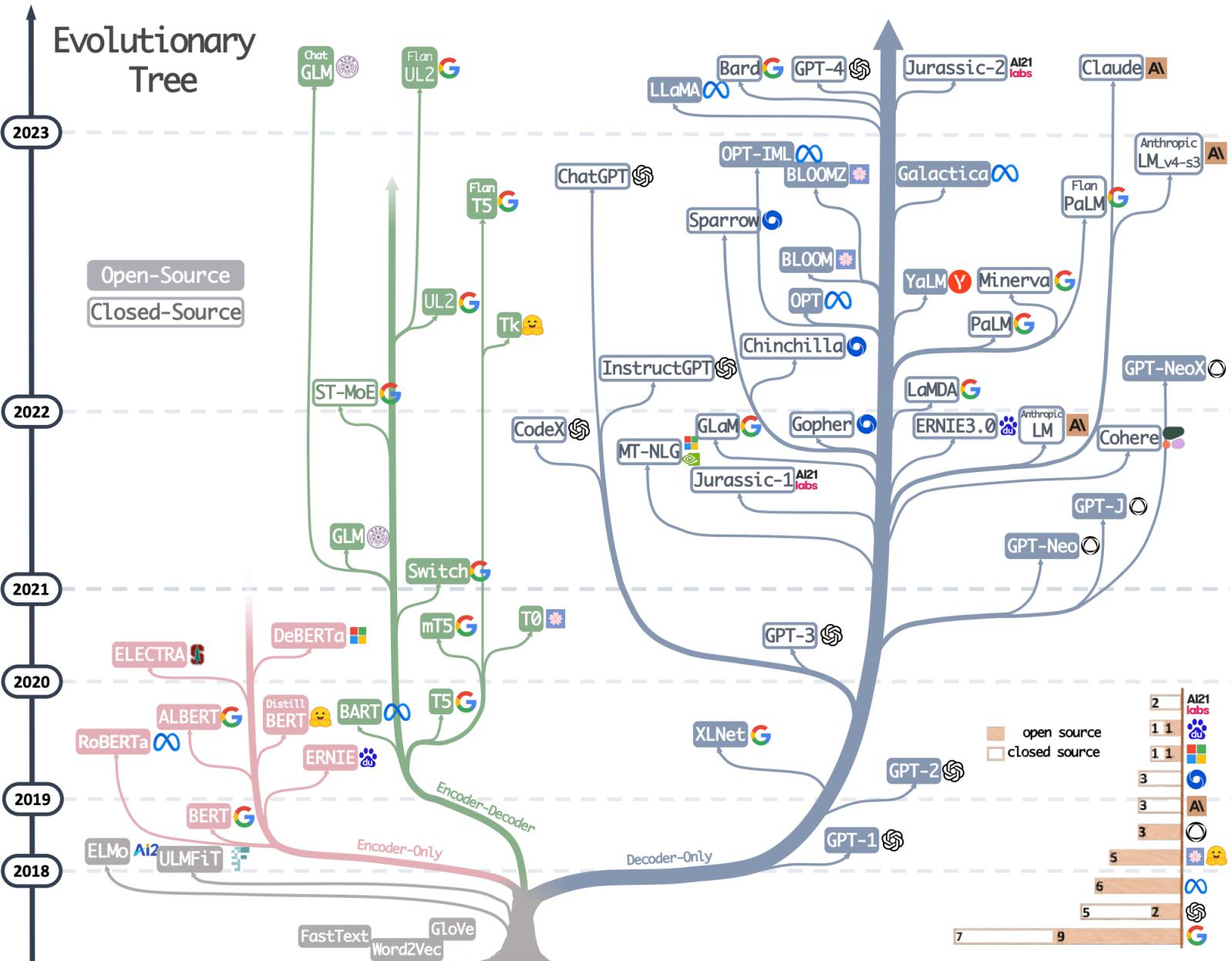
## Attention mechanism



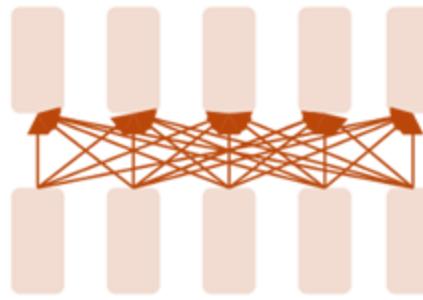
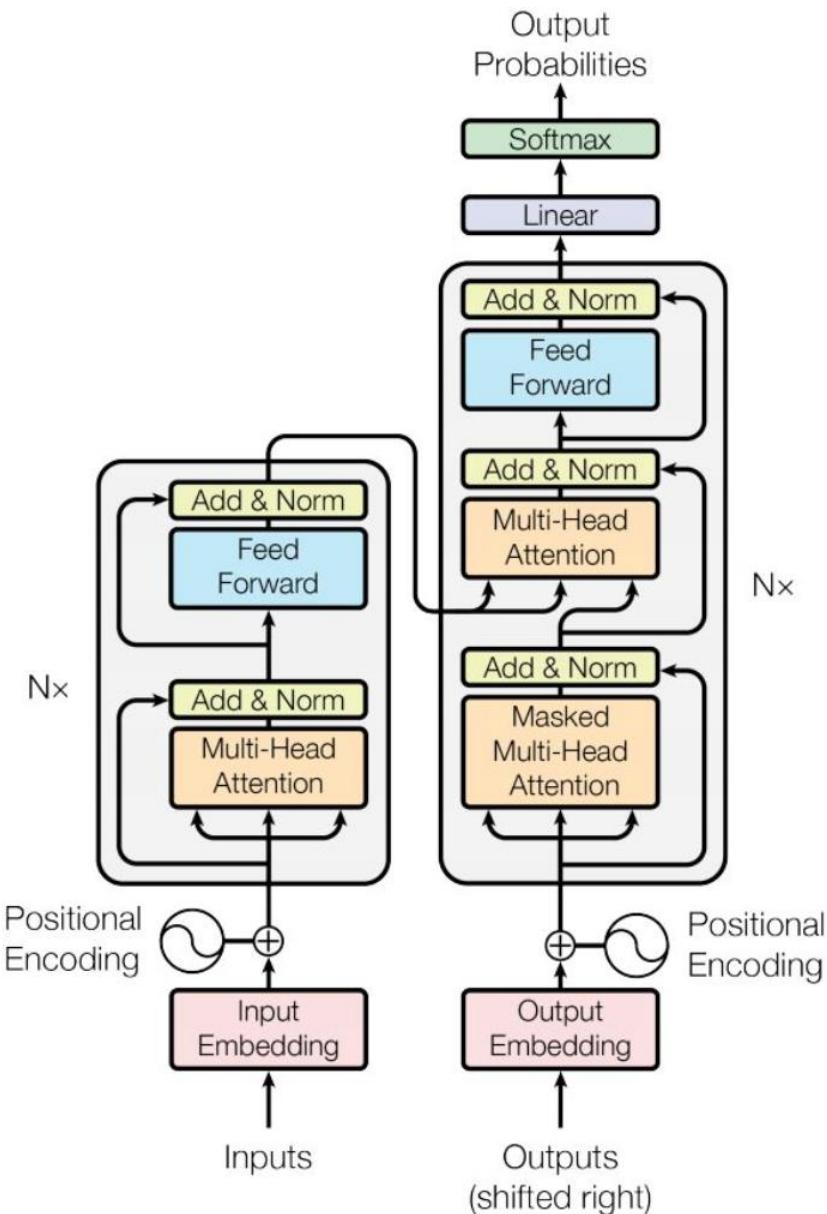
Nobel committee awards Strickland who advanced optics



# Model structure



## Encoder & encoder-decoder



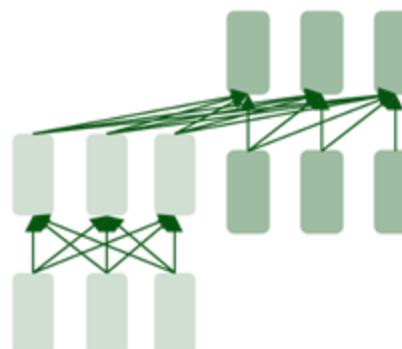
encoder

Trained by predicting words from surrounding words on both sides.

**good:** Strong comprehension ability.

**bad:** Limited generation ability.

Application: Discrimination task



encoder-decoder

Trained to map from one sequence to another

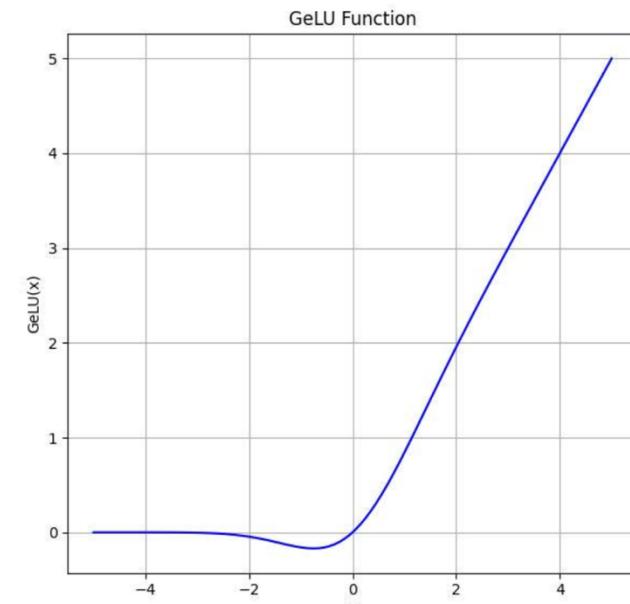
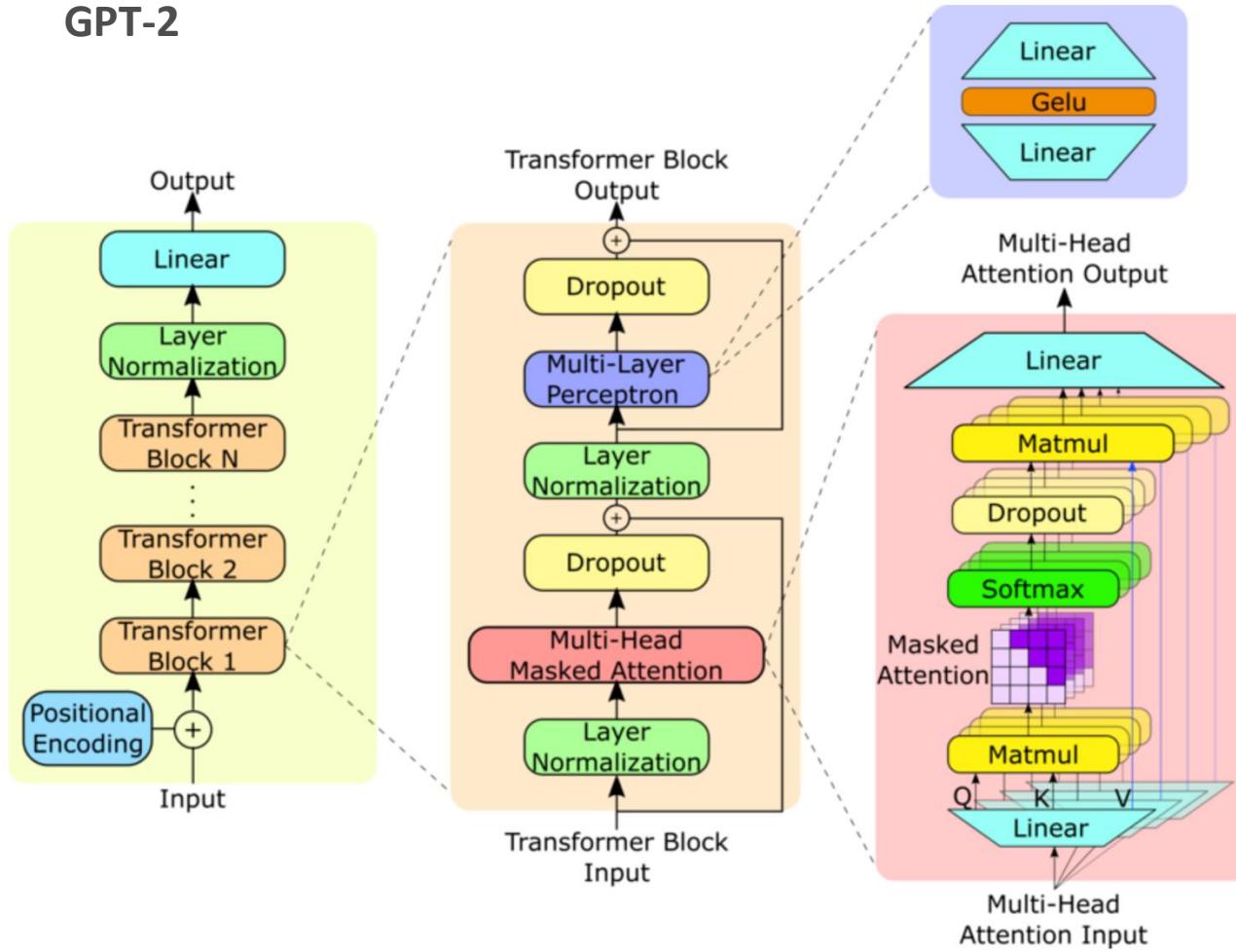
**good:** good at both comprehension and generation

**bad:** Hard and expensive to train

Application: Translation task

# Decoder

## GPT-2

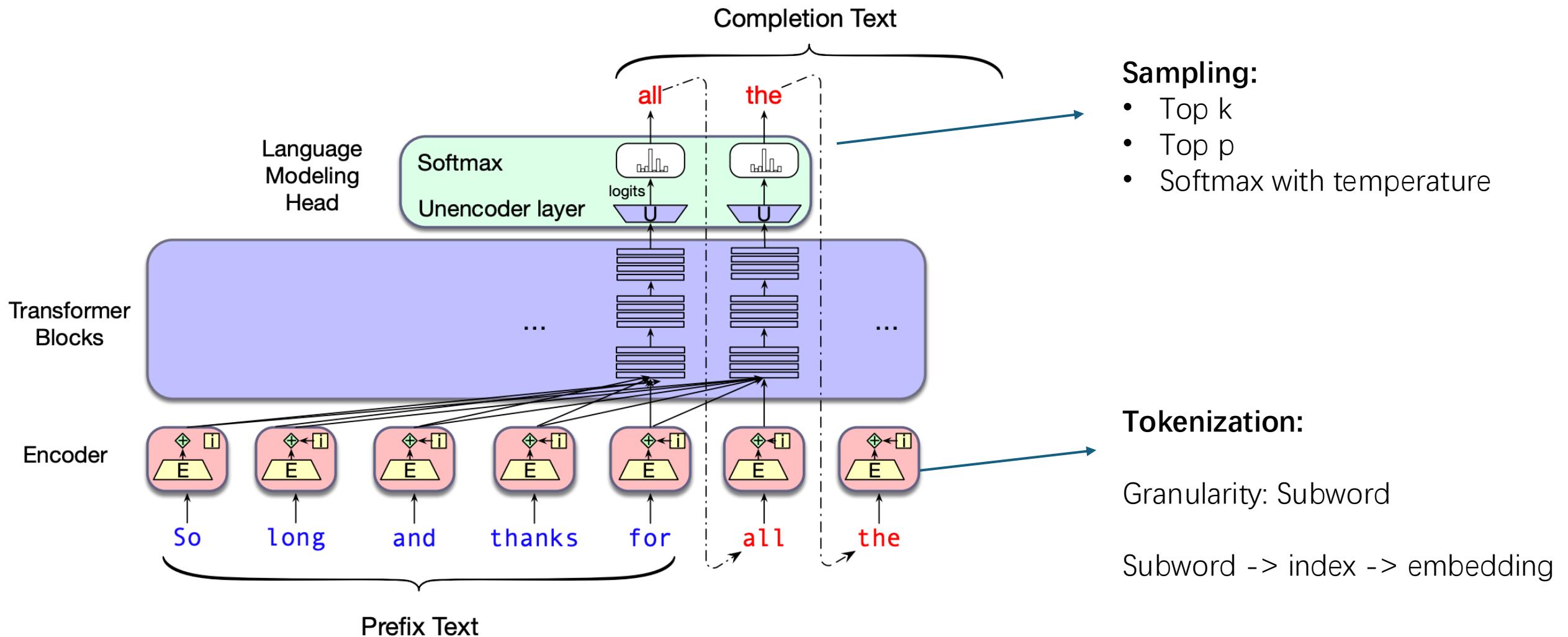


Predict from left to right

Good at generation

Easy for scaling up !

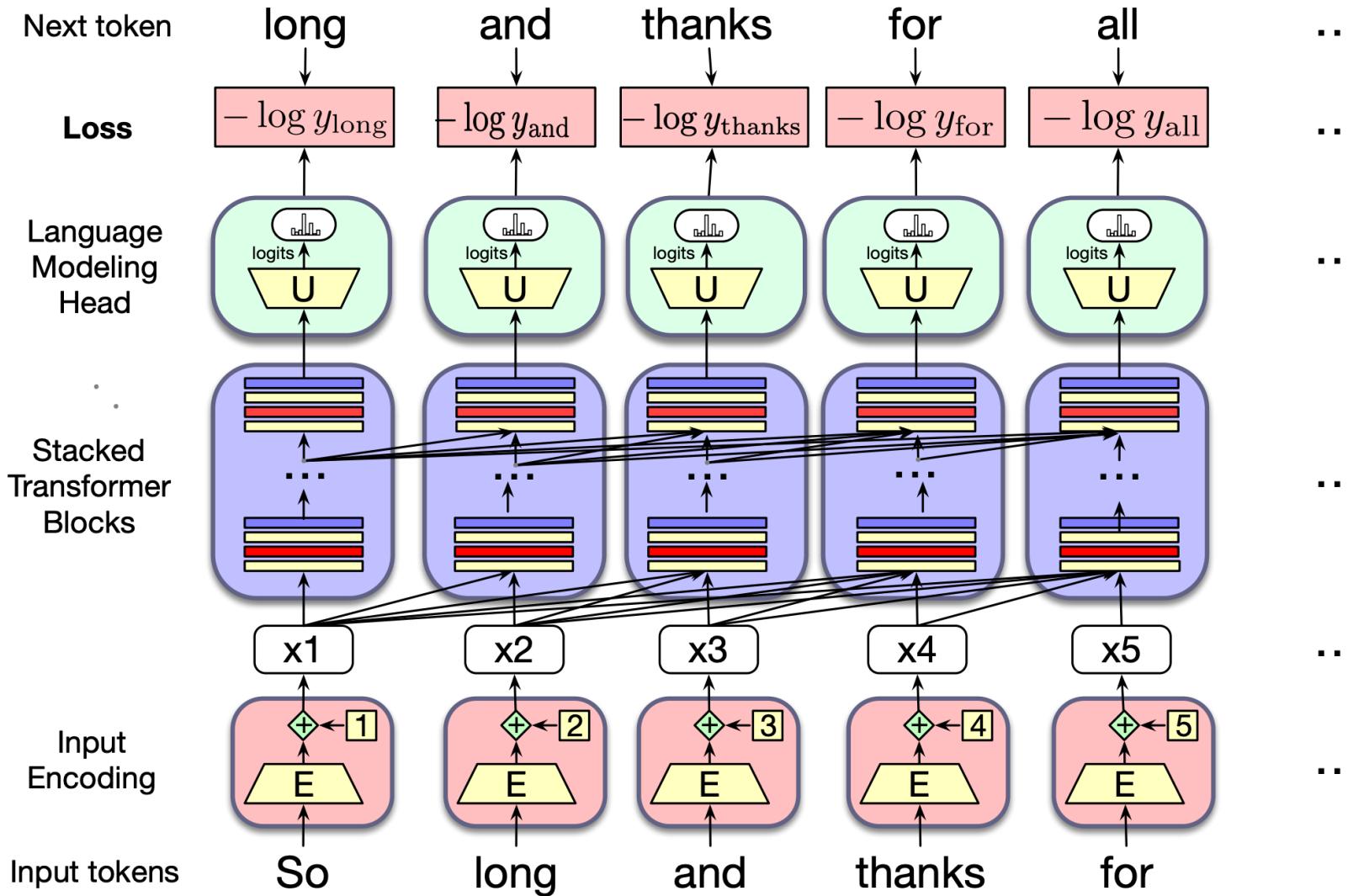
## Decoder





Pre-training & IT & RLHF

## Pre-training

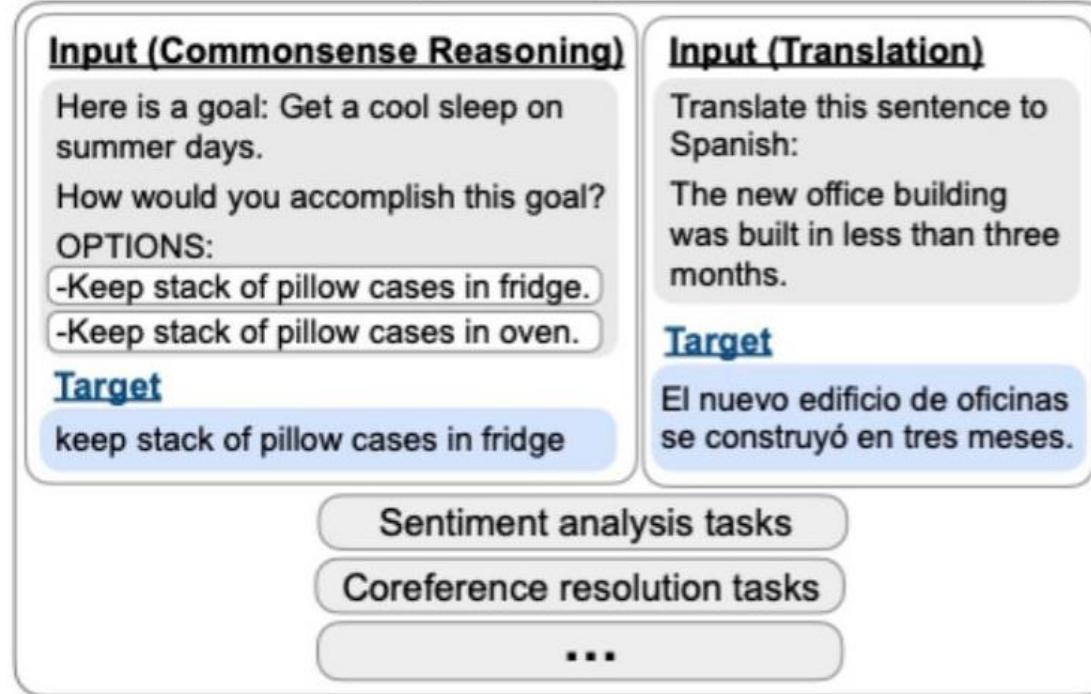


- Train to predict the next token (self-supervised training)
- Large-scale corpora from the internet
- Cross-entropy loss
- Good generalist auto-completes

Task specifical ?  
In context learning ~

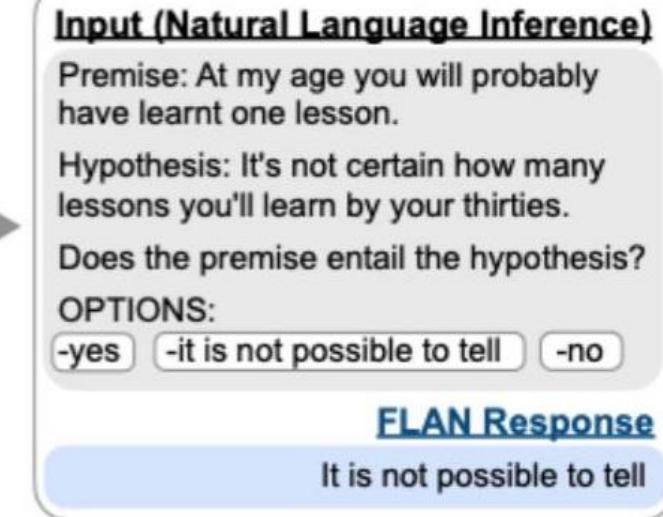
# Instruction-tuning

## Finetune on many tasks (“instruction-tuning”)



Only supervised the blue part (answer)

## Inference on unseen task type



- Input and Target : Instruction + input as input with the target in SFT
- Objective function : Loss computed only for target tokens in SFT
- Purpose : good SFT builds models that can do many unseen tasks

Expensive data labeling...

Not optimal answer for human ...

# Reinforce learning with human feedback

LLMs may produce text that can cause direct harm – **allowing easy access to dangerous information.** Therefore, LLMs should be trained to produce outputs that align with human preferences and values.

## Collect comparison data, and train a reward model.

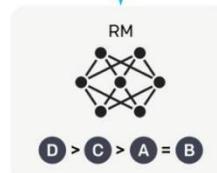
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.

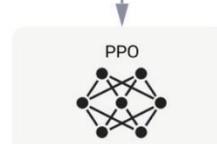


## Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.



The policy generates an output.



Once upon a time...

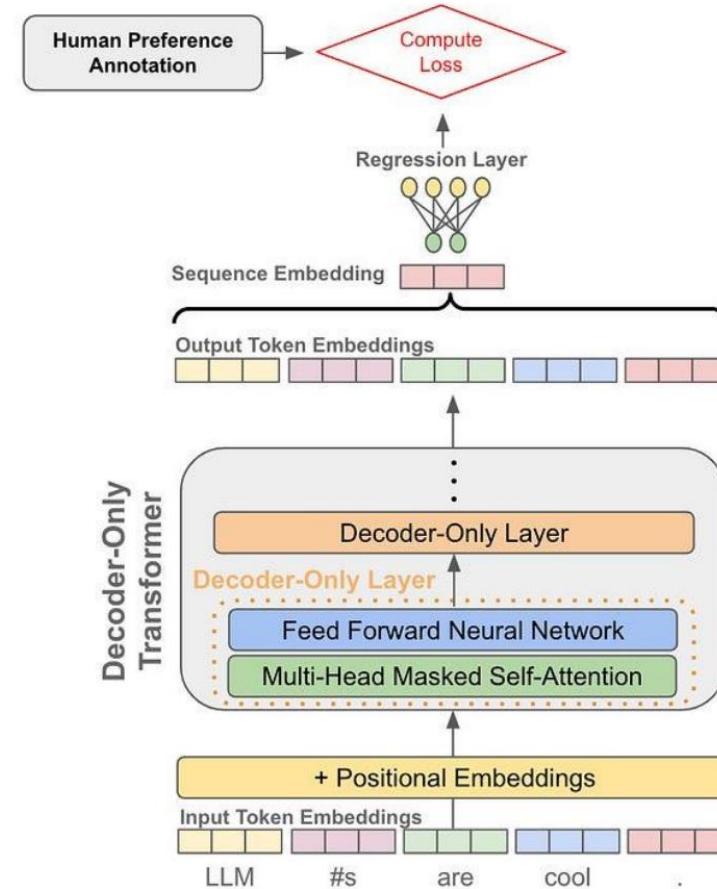


$r_k$

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

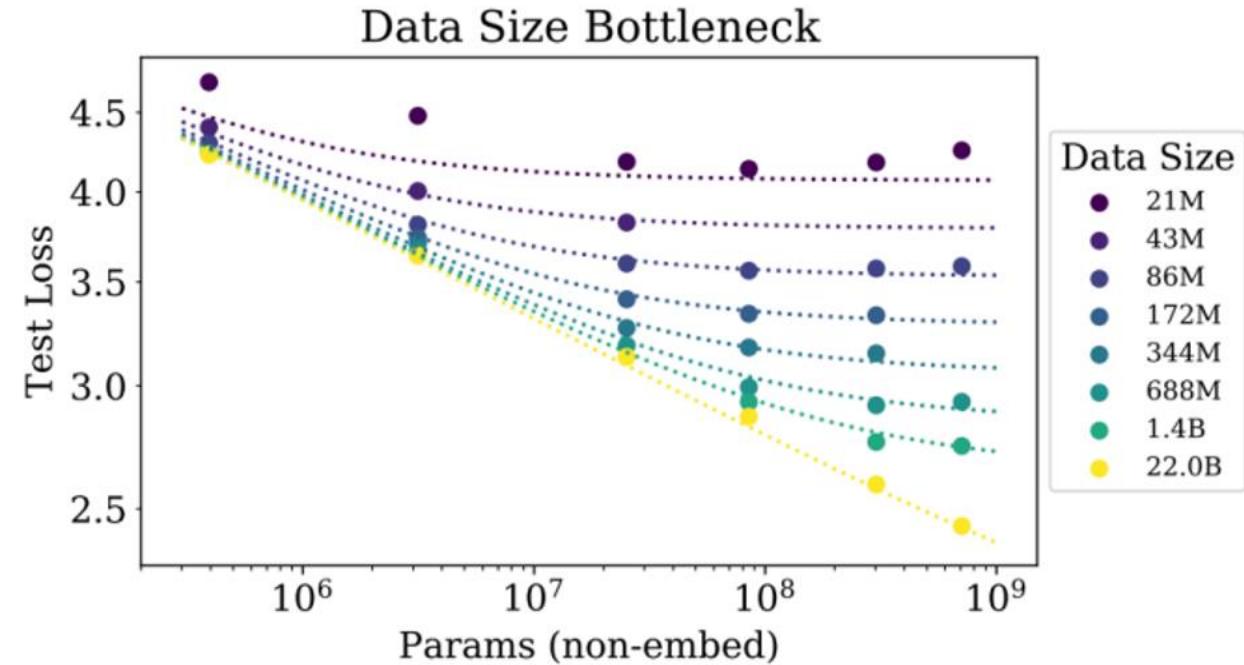
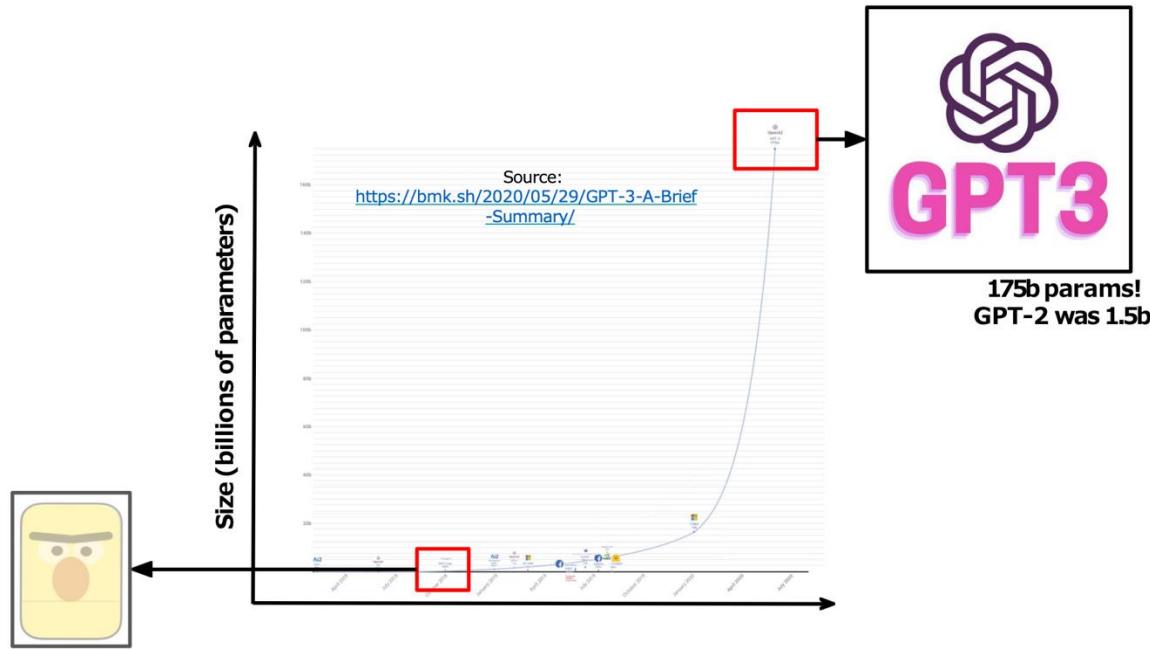
## Reward Model Structure





# Scaling

## Scaling up

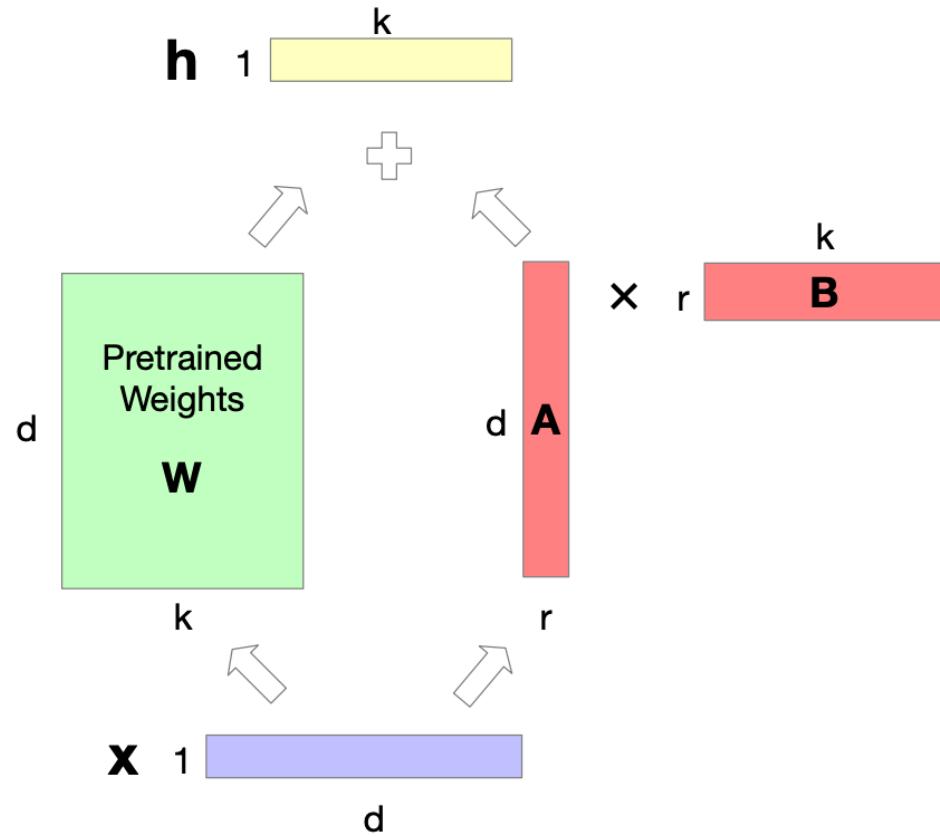


- GPT-3 trained on text can do arithmetic problems like addition and subtraction
- Different abilities “emerge” at different scales

Larger means stronger

Efficiency ...

## Parameter-Efficient Finetuning (PEFT)



Keep dense projection matrix frozen,  
update the low rank matrix

## KV Cache

$$\begin{array}{c}
 \left[ \begin{array}{c} Q \\ \vdots \\ q_1 \\ q_2 \\ q_3 \\ q_4 \end{array} \right]_{N \times d_K} \times \left[ \begin{array}{c} K^T \\ \vdots \\ \mathbb{1} \\ \mathbb{1} \\ \mathbb{1} \\ \mathbb{1} \end{array} \right]_{d_K \times N} = \left[ \begin{array}{c} QK^T \\ \vdots \\ q_1 \cdot k_1 & q_1 \cdot k_2 & q_1 \cdot k_3 & q_1 \cdot k_4 \\ q_2 \cdot k_1 & q_2 \cdot k_2 & q_2 \cdot k_3 & q_2 \cdot k_4 \\ q_3 \cdot k_1 & q_3 \cdot k_2 & q_3 \cdot k_3 & q_3 \cdot k_4 \\ q_4 \cdot k_1 & q_4 \cdot k_2 & q_4 \cdot k_3 & q_4 \cdot k_4 \end{array} \right]_{N \times N} \\
 \left[ \begin{array}{c} Q \\ q_4 \end{array} \right]_{1 \times d_K} \times \left[ \begin{array}{c} K^T \\ \vdots \\ \mathbb{1} \\ \mathbb{1} \\ \mathbb{1} \\ \mathbb{1} \end{array} \right]_{d_K \times N} = \left[ \begin{array}{c} QK^T \\ \vdots \\ q_4 \cdot k_1 & q_4 \cdot k_2 & q_4 \cdot k_3 & q_4 \cdot k_4 \end{array} \right]_{1 \times N}
 \end{array}$$

$$\begin{array}{c}
 V \\
 \times \\
 V
 \end{array}
 = \begin{array}{c}
 A \\
 \times \\
 A
 \end{array}$$

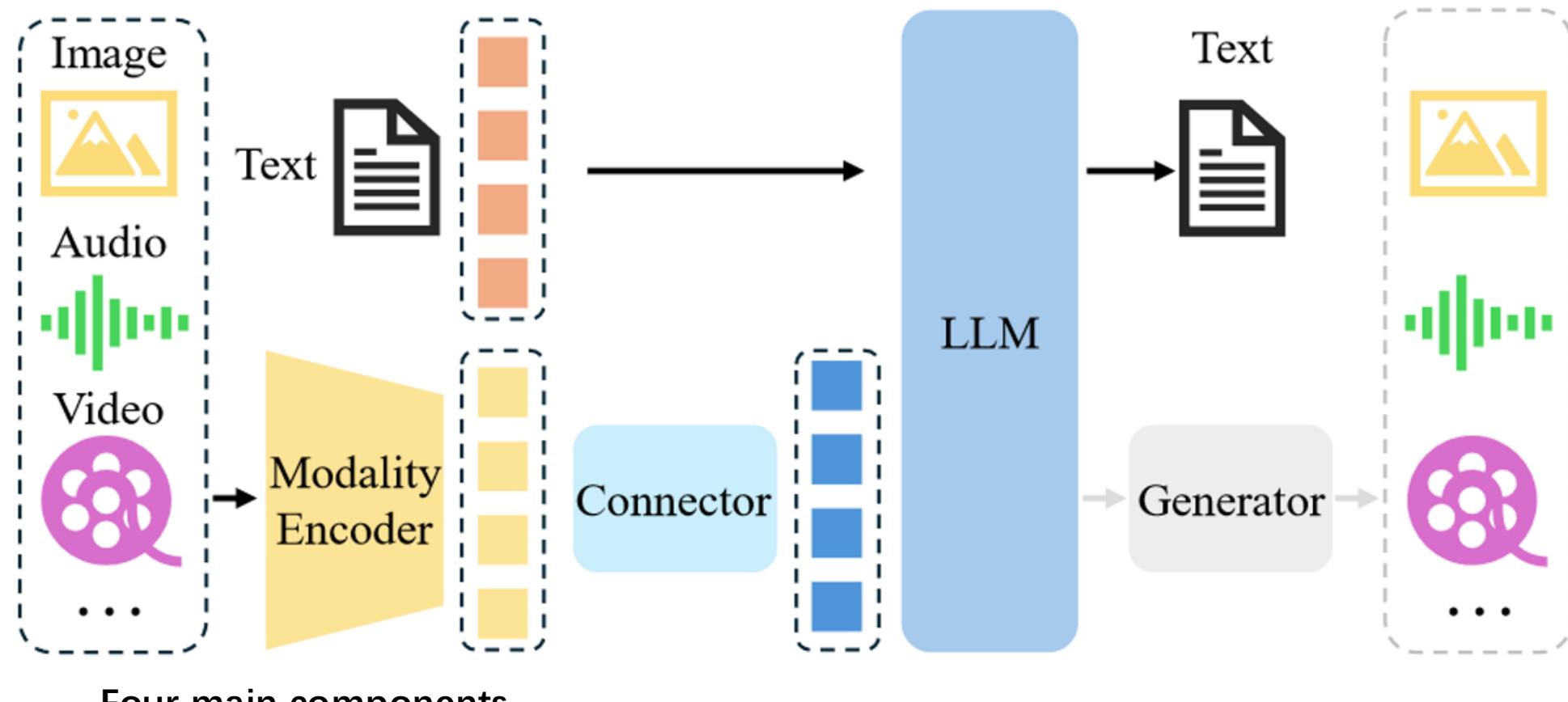
$$\begin{array}{c}
 V \\
 \times \\
 V
 \end{array}
 = \begin{array}{c}
 a \\
 \times \\
 a
 \end{array}$$

Avoid recompute the past keys and values during inference



# Multi Modal LLM

# Multi Modal LLM



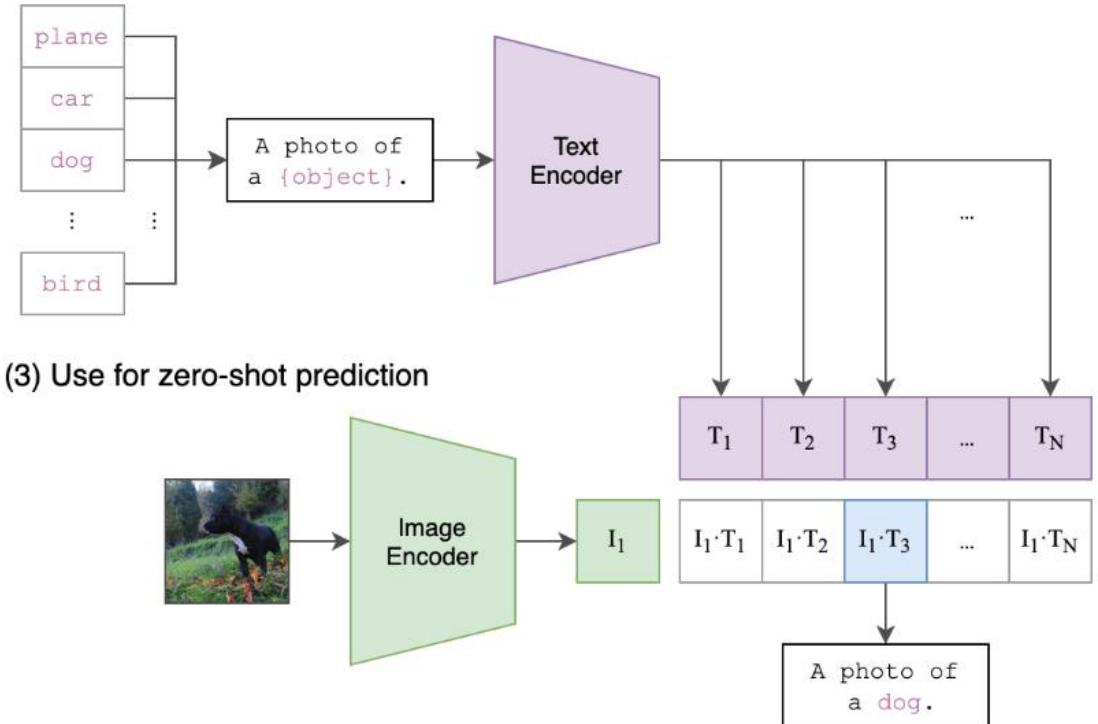
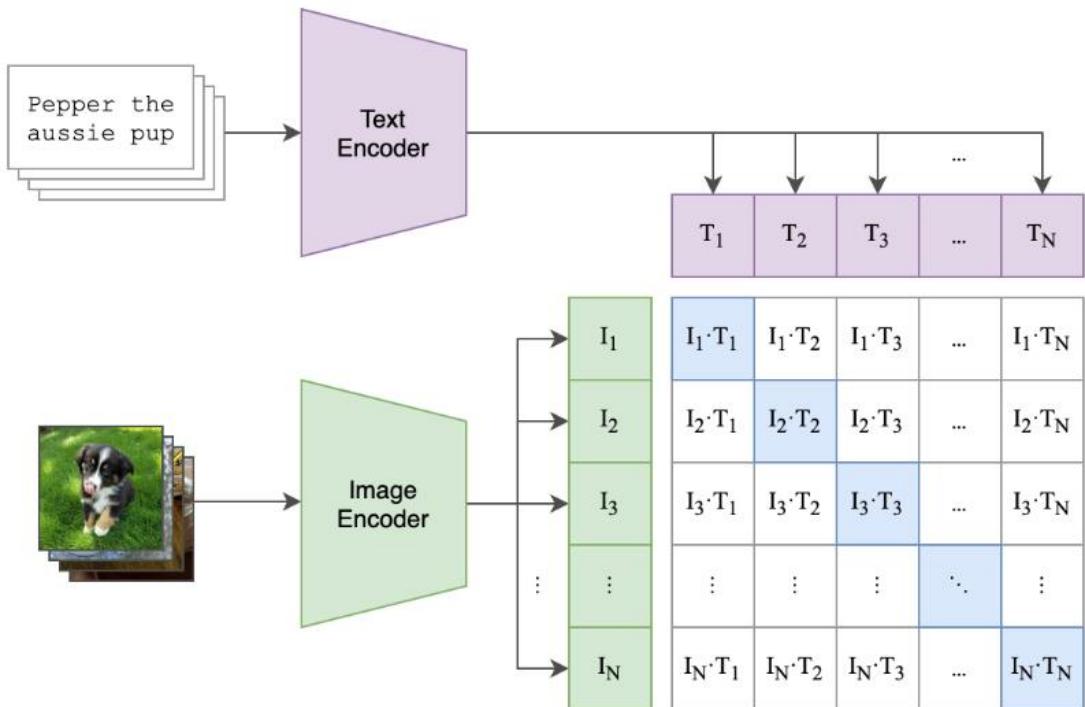
(1) multimodal encoder

(2) connector

(3) large language model

(4) multimodal generator

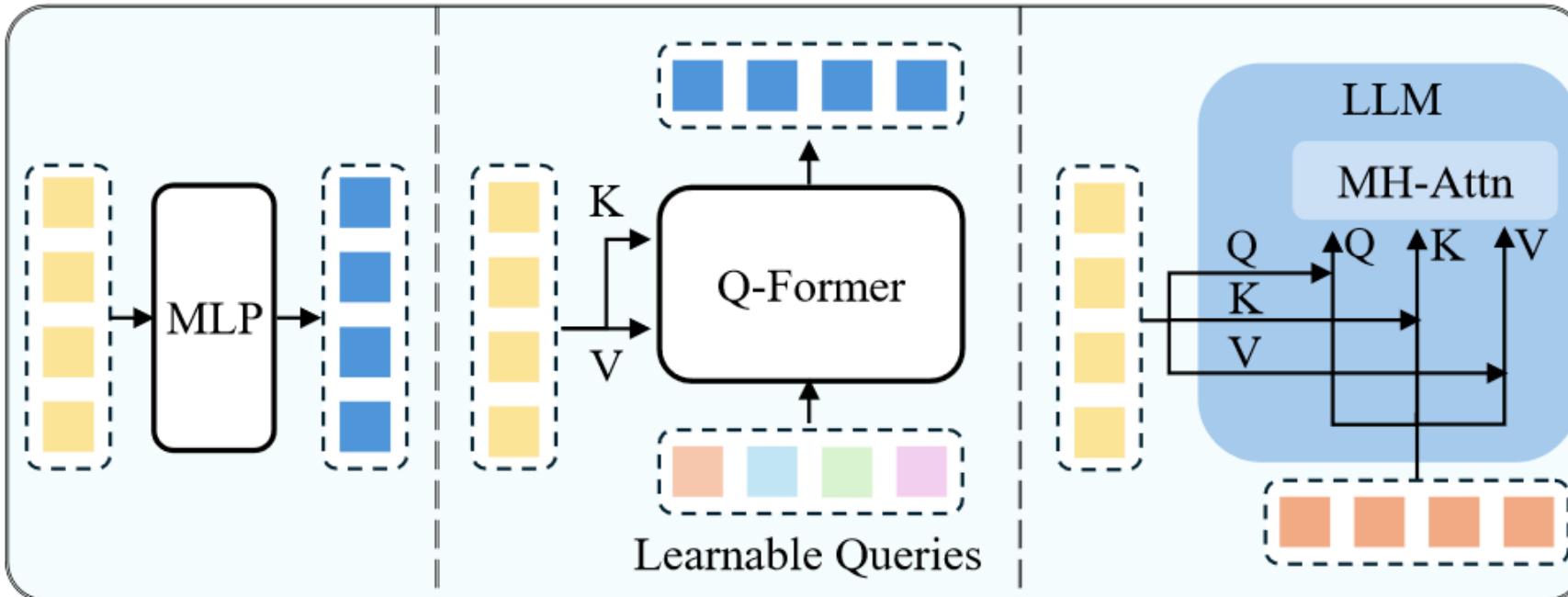
# Multimodal Encoder



## Key factor:

1. Encoder parameter number
2. Pretrained dataset size
3. Resolution ratio

## Connector

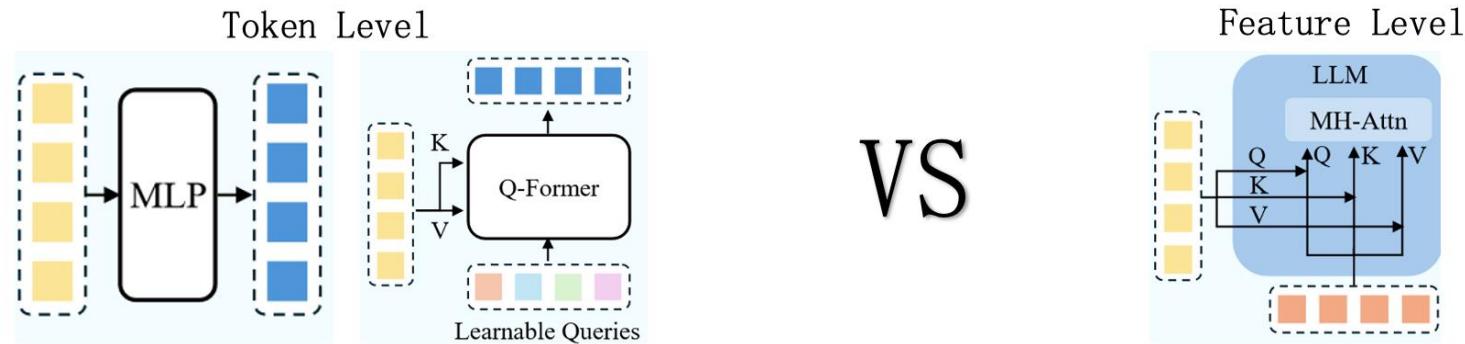


The role of connectors is to integrate multimodal information, which can be divided into token-level and feature-level based on the fusion hierarchy.

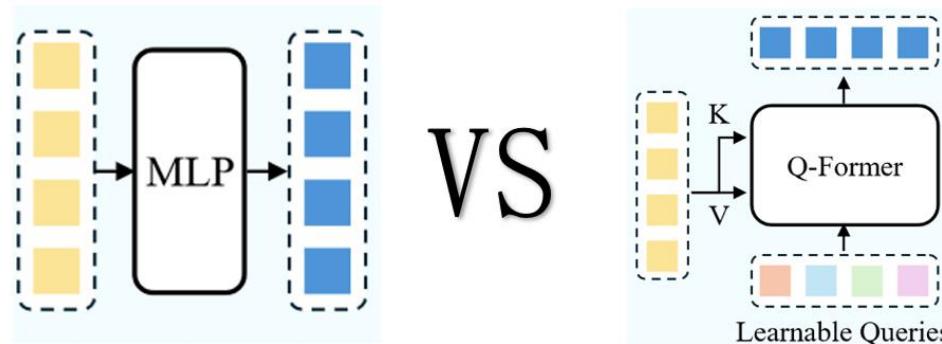
Three mainstream connectors are:

MLP (token-level), Q-Former (token-level), and multi-head attention (feature-level).

## Connector



Token-level performs better in VQA benchmark tests (VQA: Visual Question Answering)

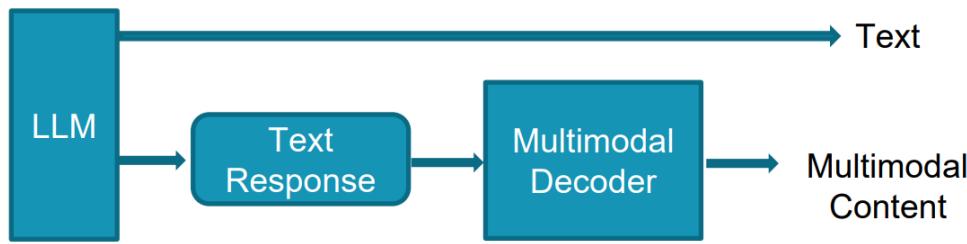


The type of the connector is far less important than the number of visual tokens and the input resolution.

# Multimodal Generator

There are two implementation methods for the multimodal decoder:

- (1) Using the text output as the input.
  - (2) Using the embeddings corresponding to specific tokens as the input.
- 

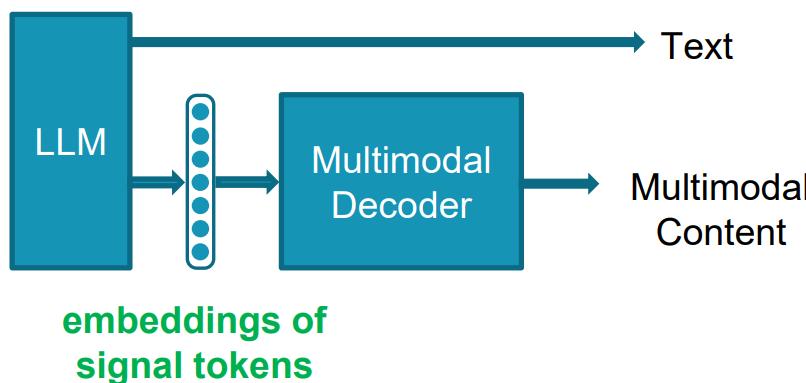


Advantages:

- Efficient, no need to fine-tune the Large Language Model (LLM).
- High lower limit of performance.

Disadvantages:

- Lack the ability of end-to-end fine-tuning.
  - Low upper limit of performance, and some multimodal tasks cannot be translated into text.
- 



Characteristics:

- It can be fine-tuned in an end-to-end manner.
- It has a high upper limit of performance and can convey information that cannot be carried by text, such as: visual spatial relationships.



Further...

LLM generate image?

