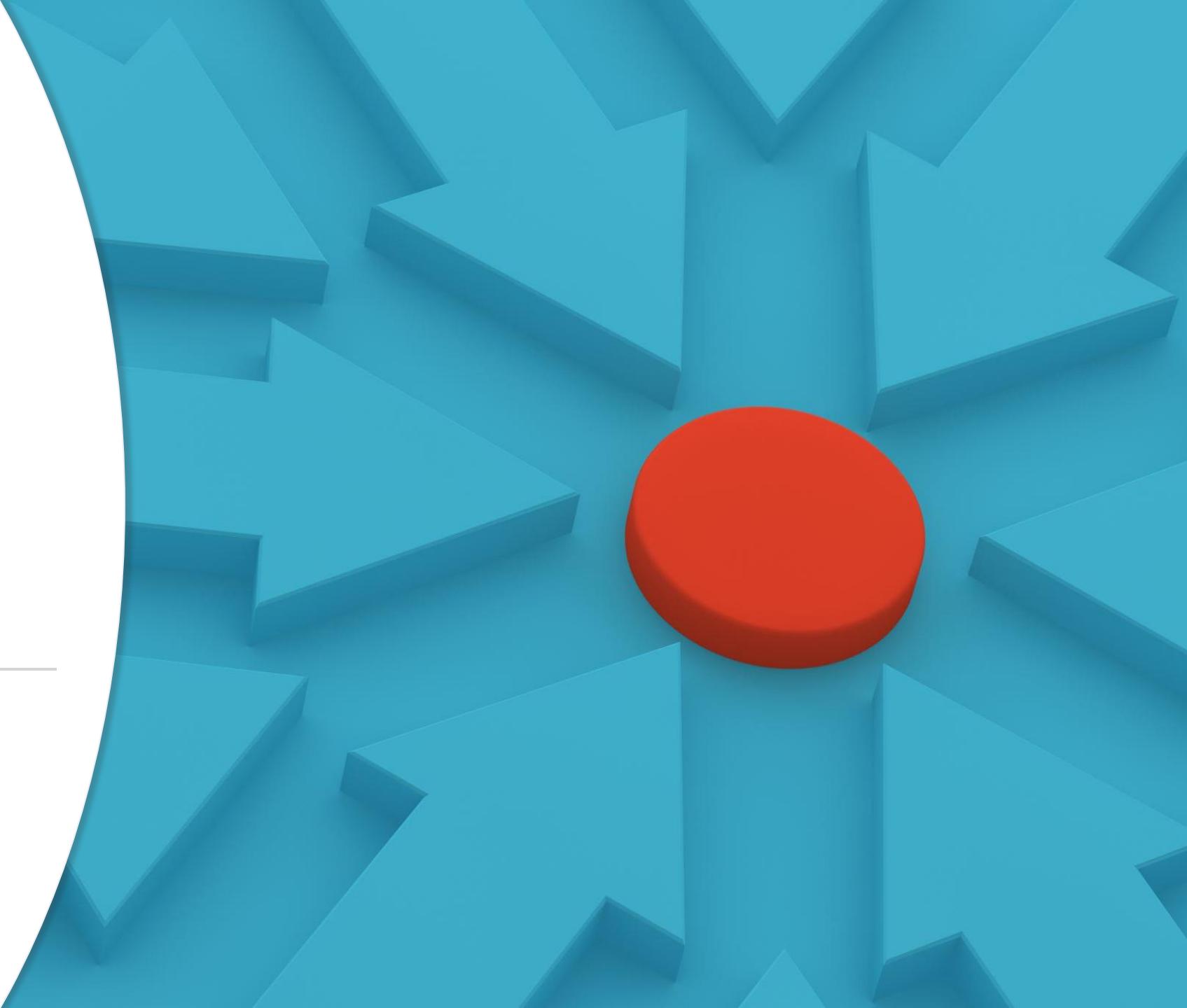




A brief introduction to
GAI Principle

Bo Yuan



Outline

Transformer

LLM Training

DeepSeek

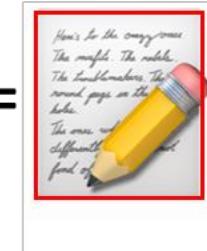
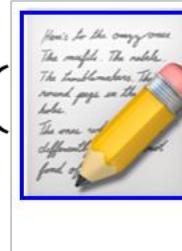
What is GAI?

GAI = Generative Artificial Intelligent

via

LLM = Large Language Model

$$f_{\mathbf{w}}(\quad) =$$



Sequence to Sequence

Hardness of GAI

Q: Write a poem consisting of 100 words.

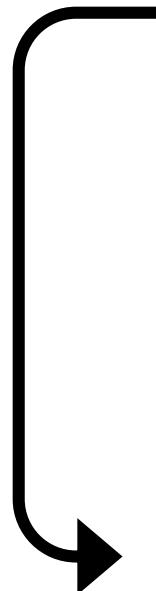
A:

Number of possible answers?

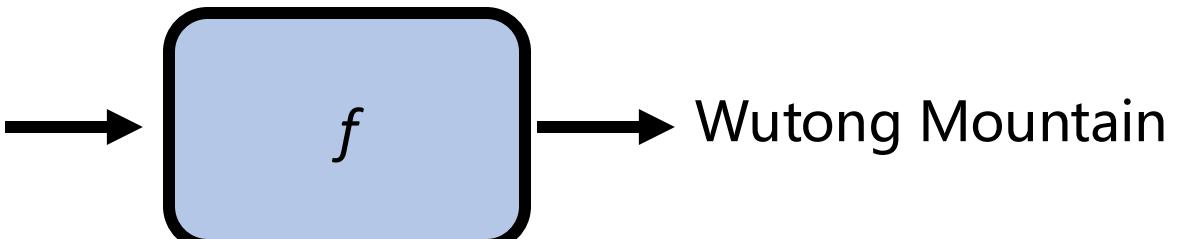


$$1000 \times 1000 \times 1000 \times 1000 \dots = 1000^{100} = 10^{300}$$

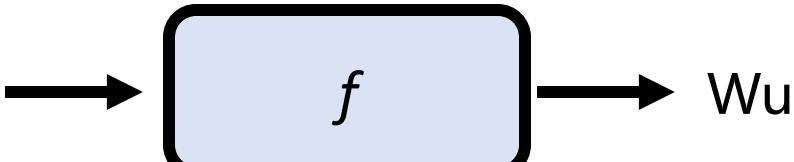
LLM



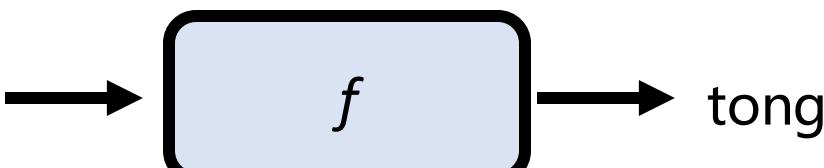
Which is the highest mountain
in Shenzhen?



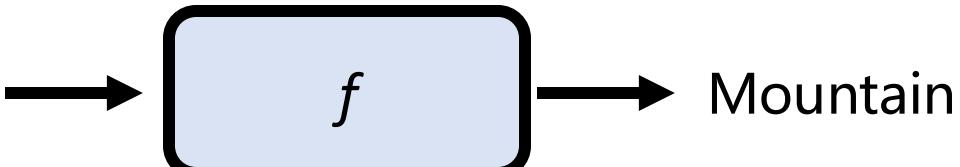
Which is the highest mountain
in Shenzhen?



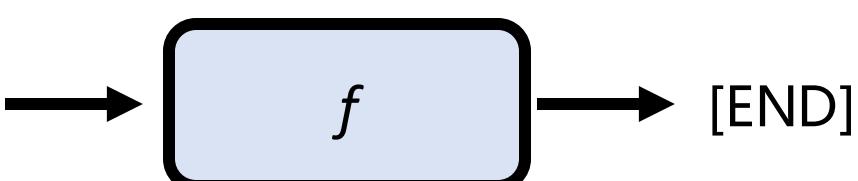
Which is the highest mountain
in Shenzhen? Wu



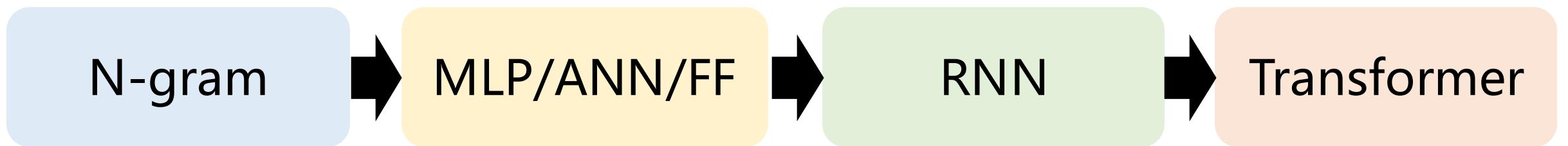
Which is the highest mountain in
Shenzhen? Wutong



Which is the highest mountain in
Shenzhen? Wutong Mountain



Language Model



Transformer

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

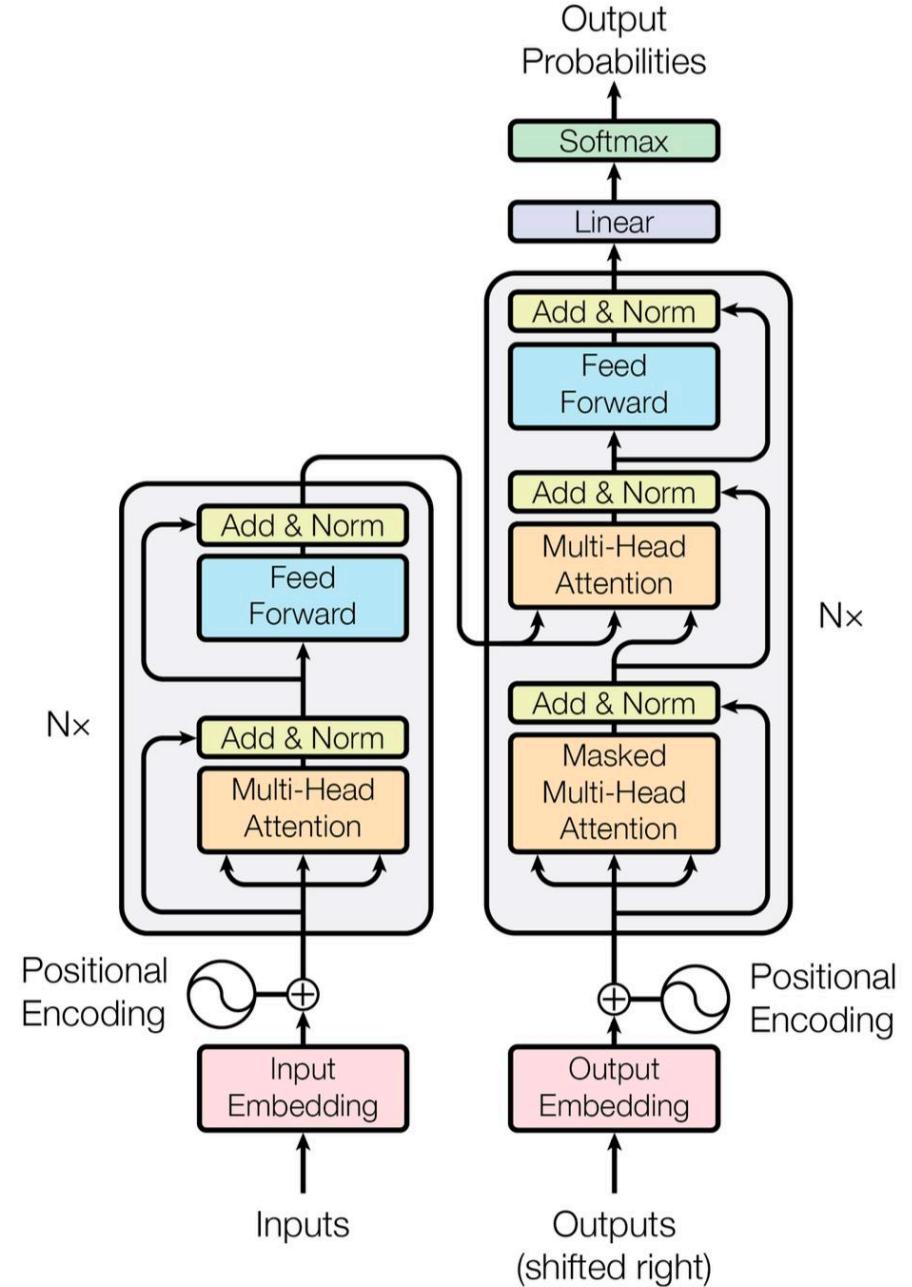
Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

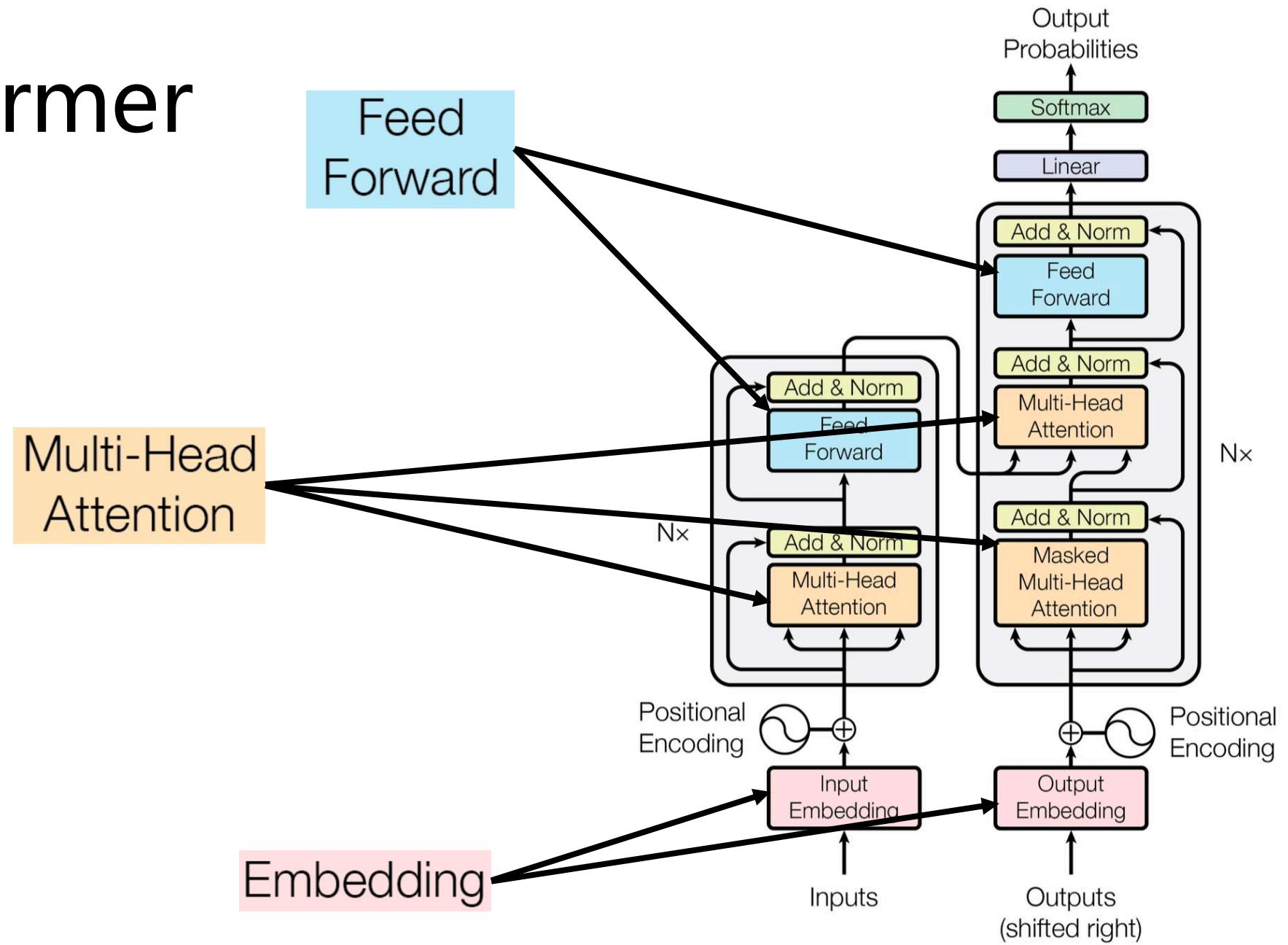
Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

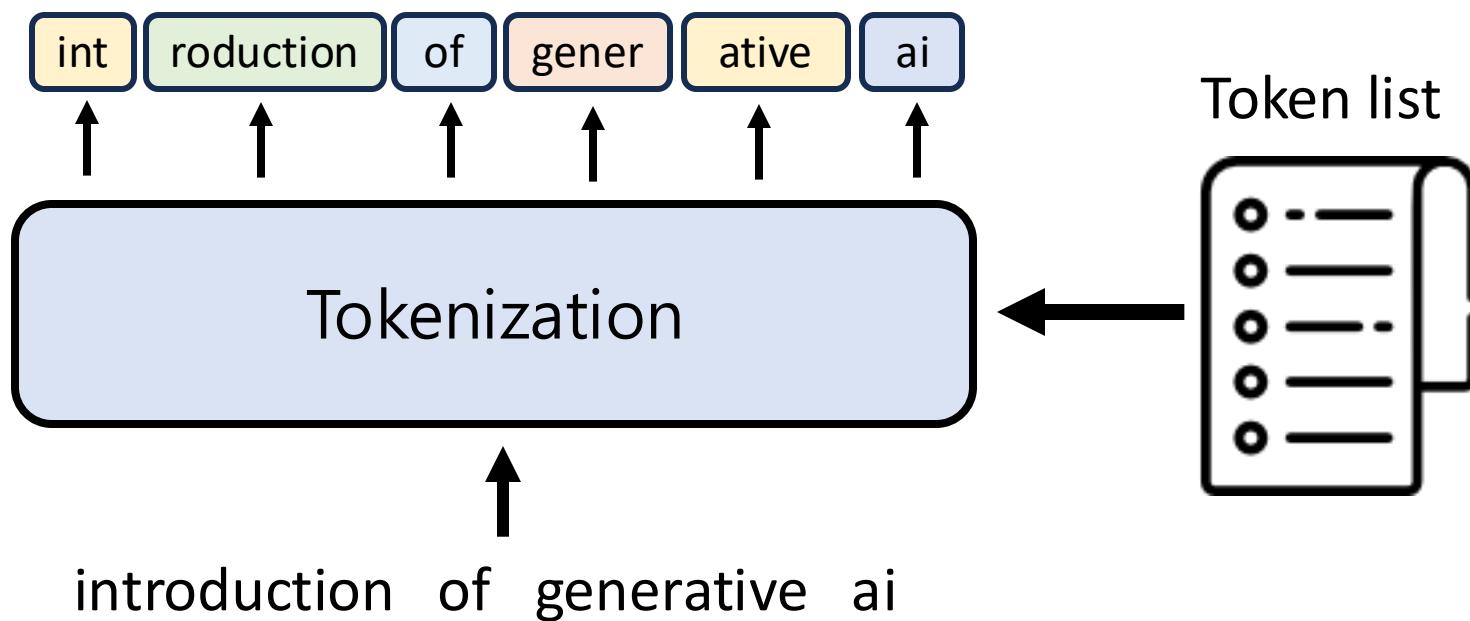
Illia Polosukhin* ‡
illia.polosukhin@gmail.com



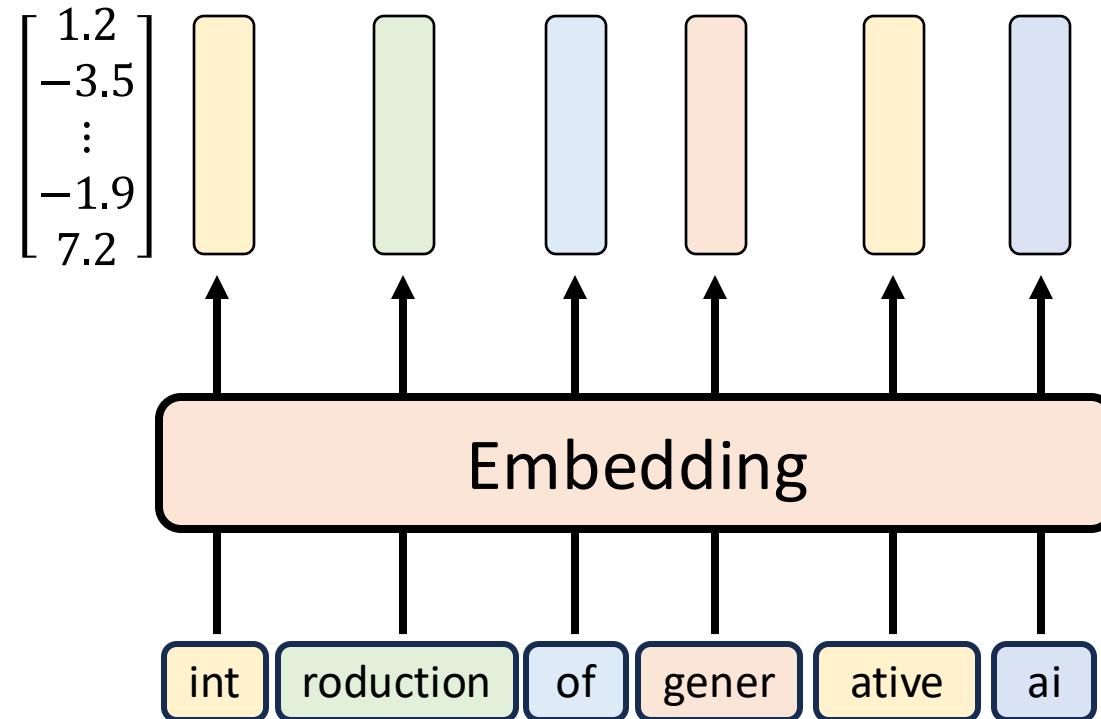
Transformer



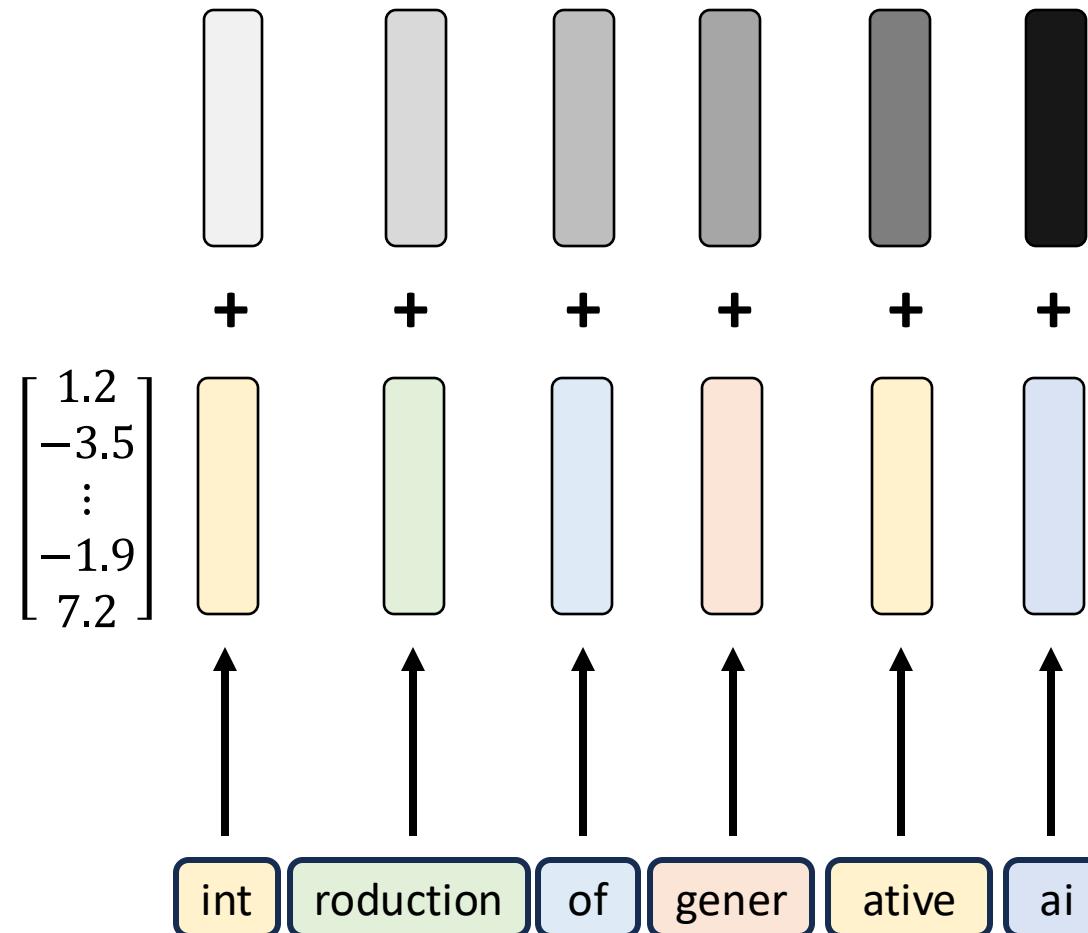
Tokenization (Pre-processing)



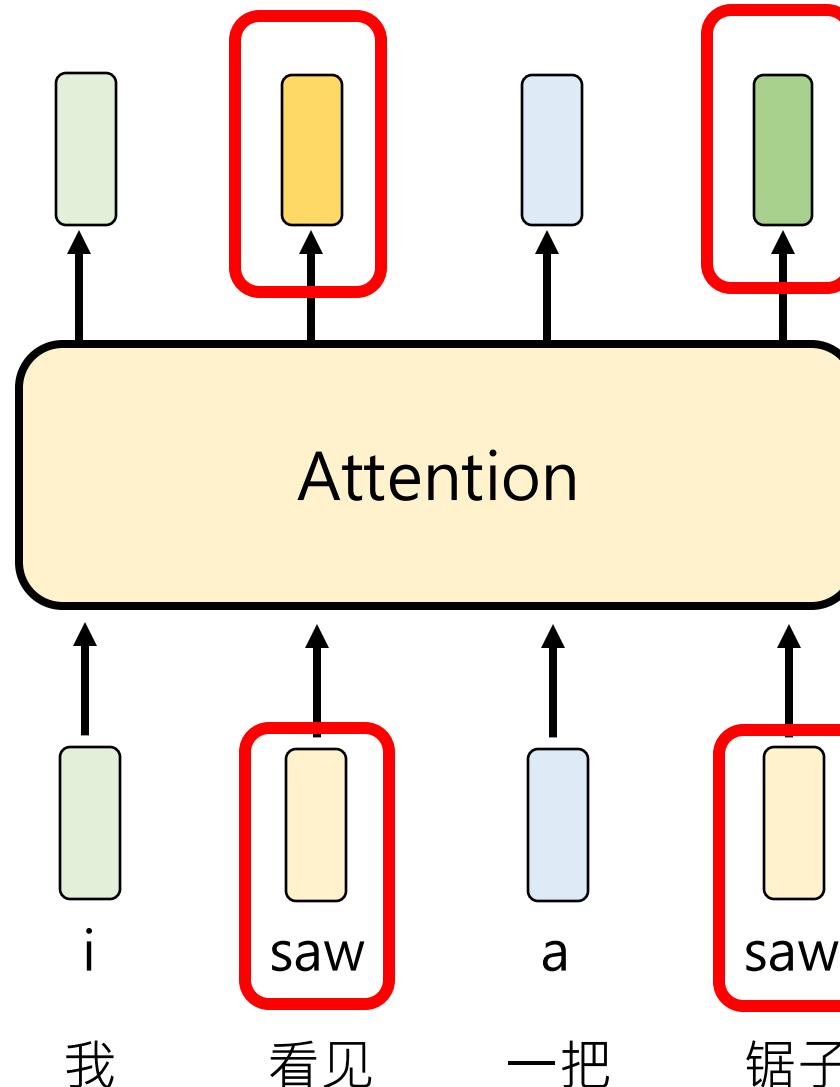
Embedding (Semantic)



Embedding (Positional)



Attention (Context)

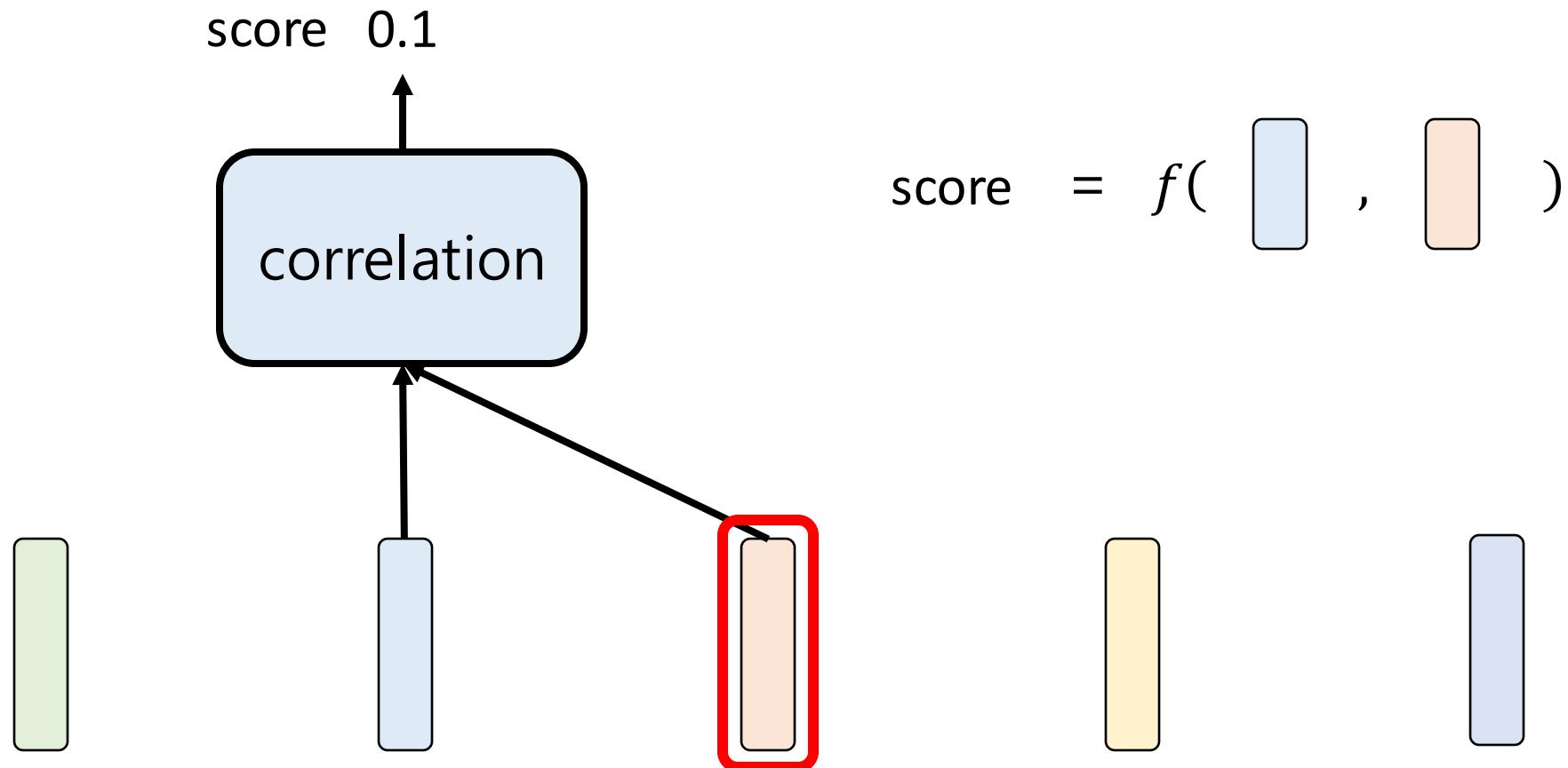


苹果电脑

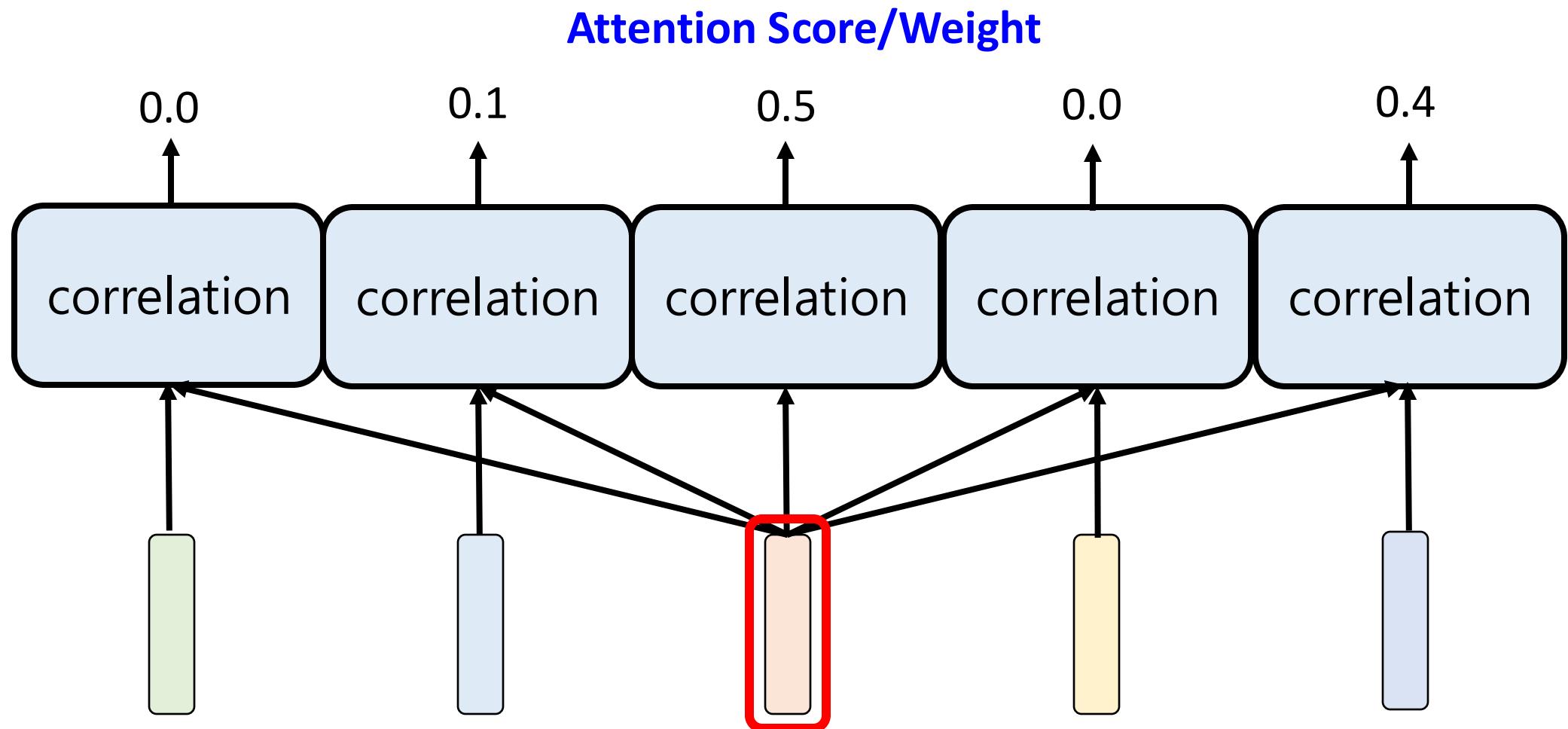
苹果好吃



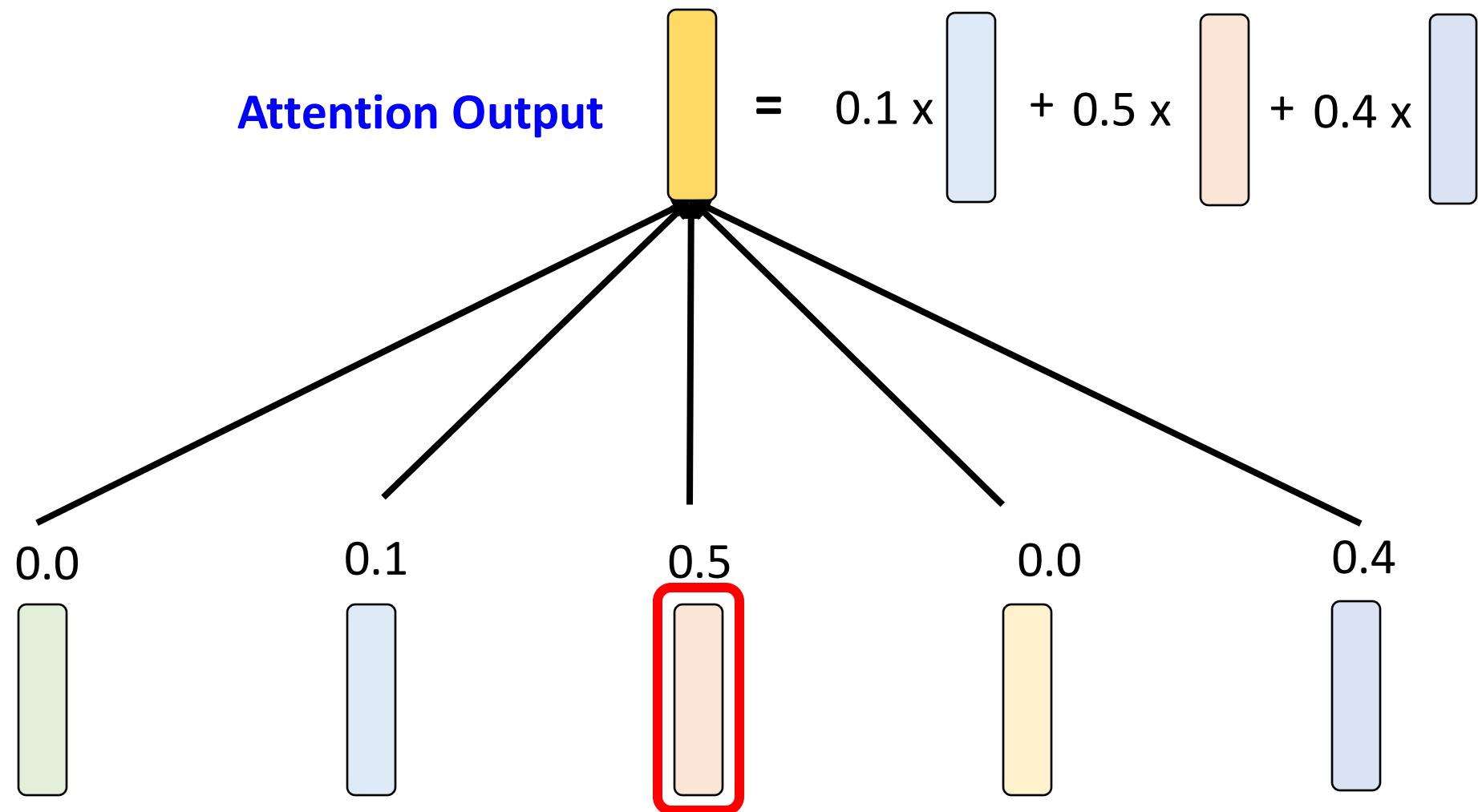
Attention (Context)



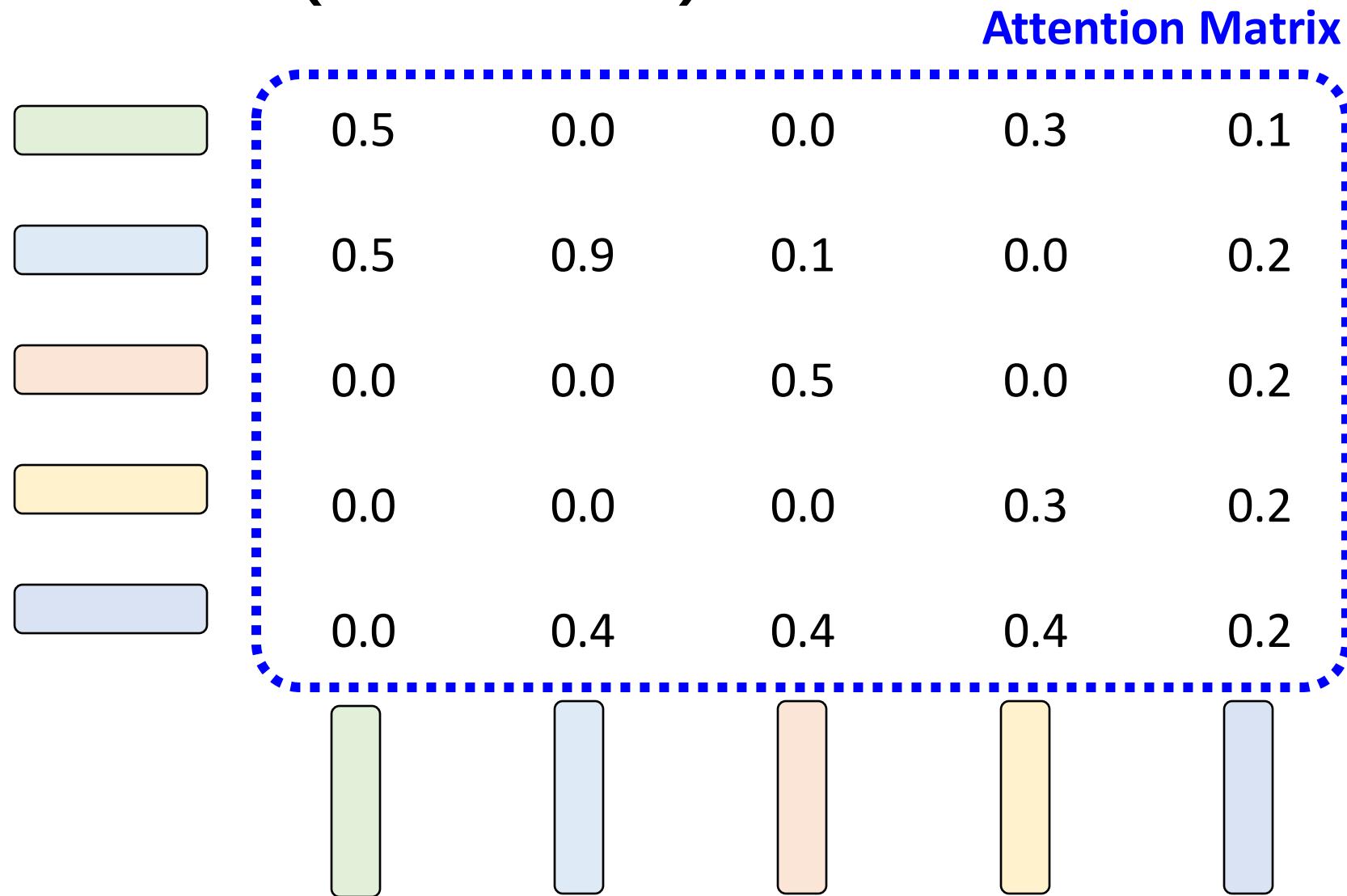
Attention (Context)



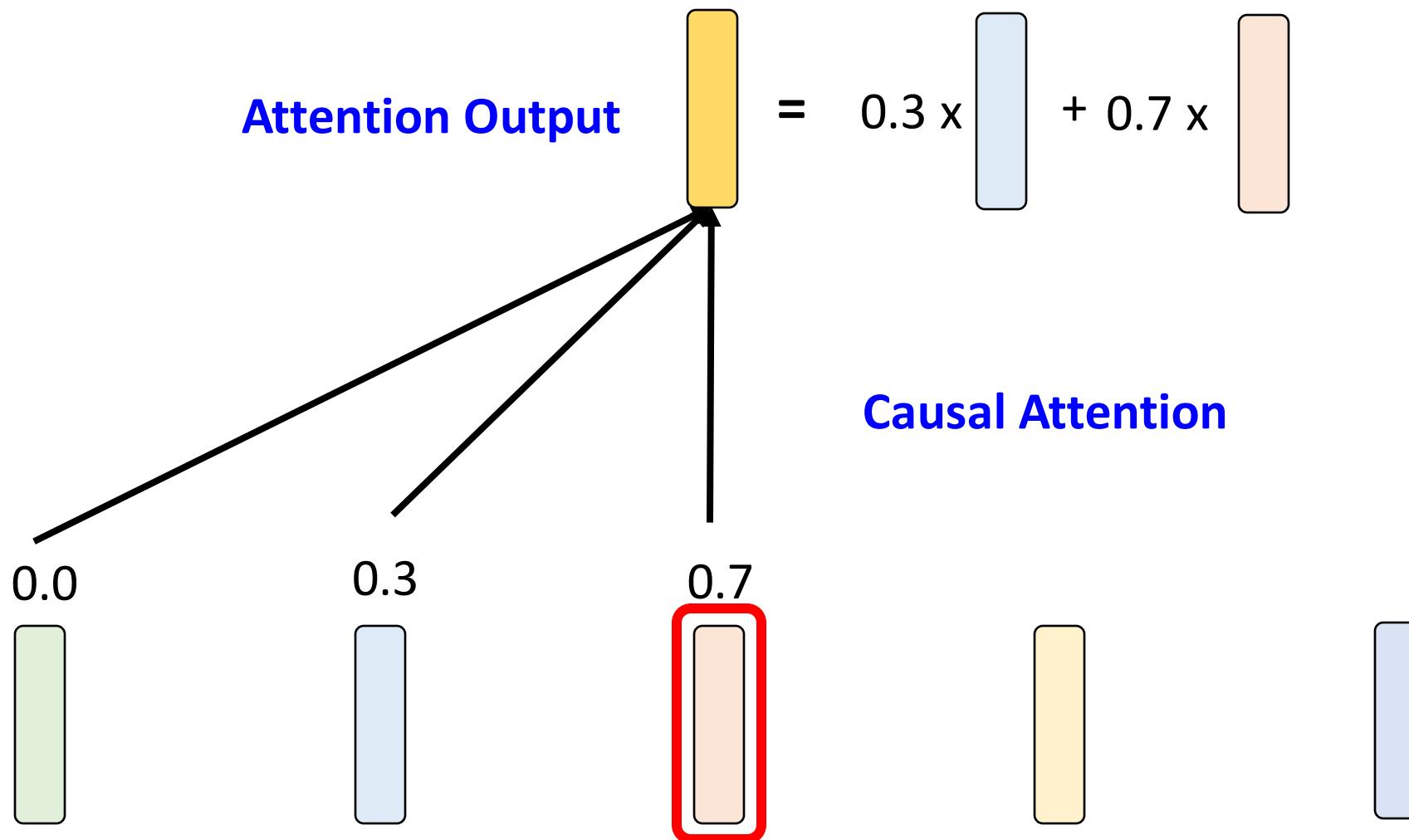
Attention (Context)



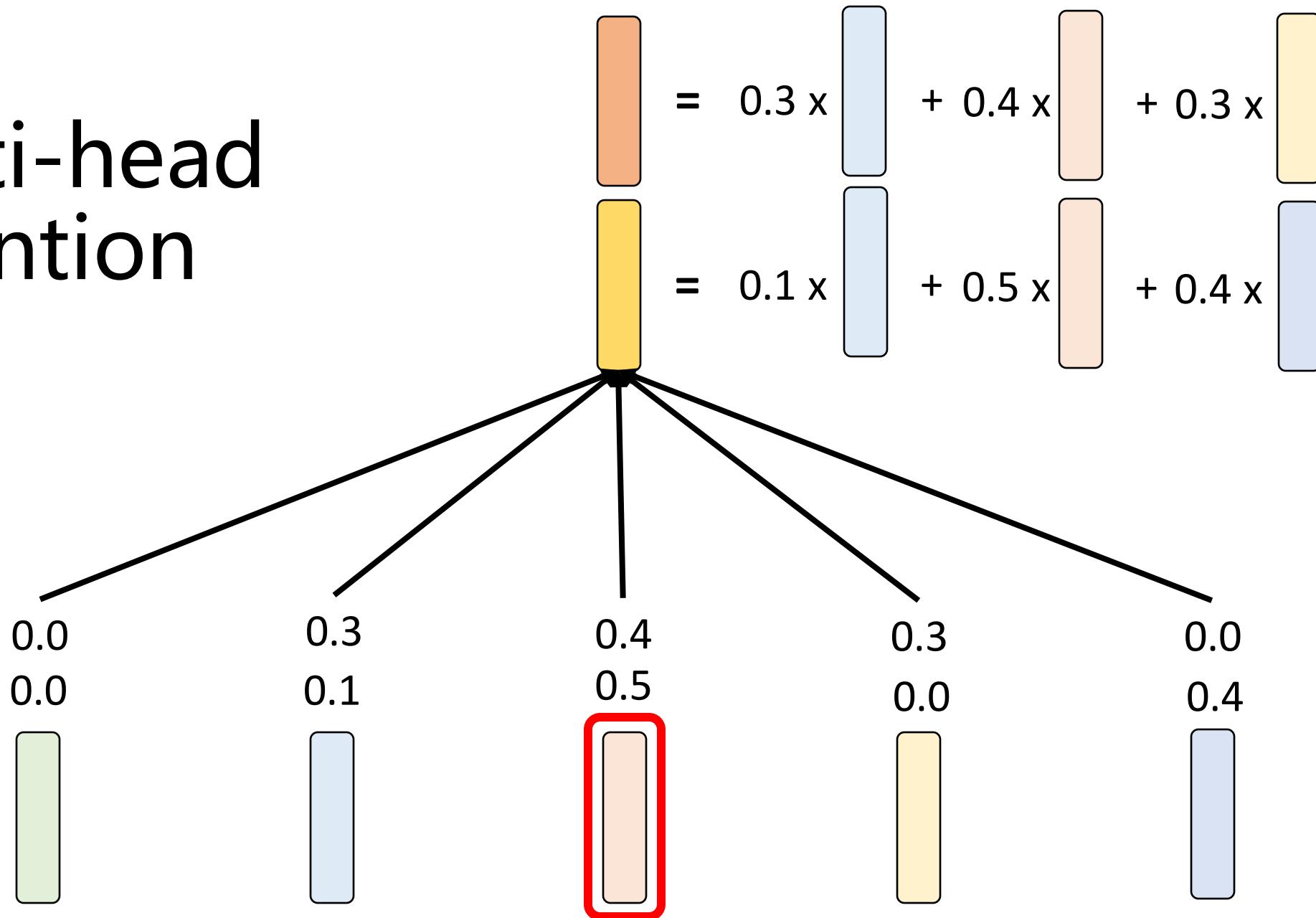
Attention (Context)



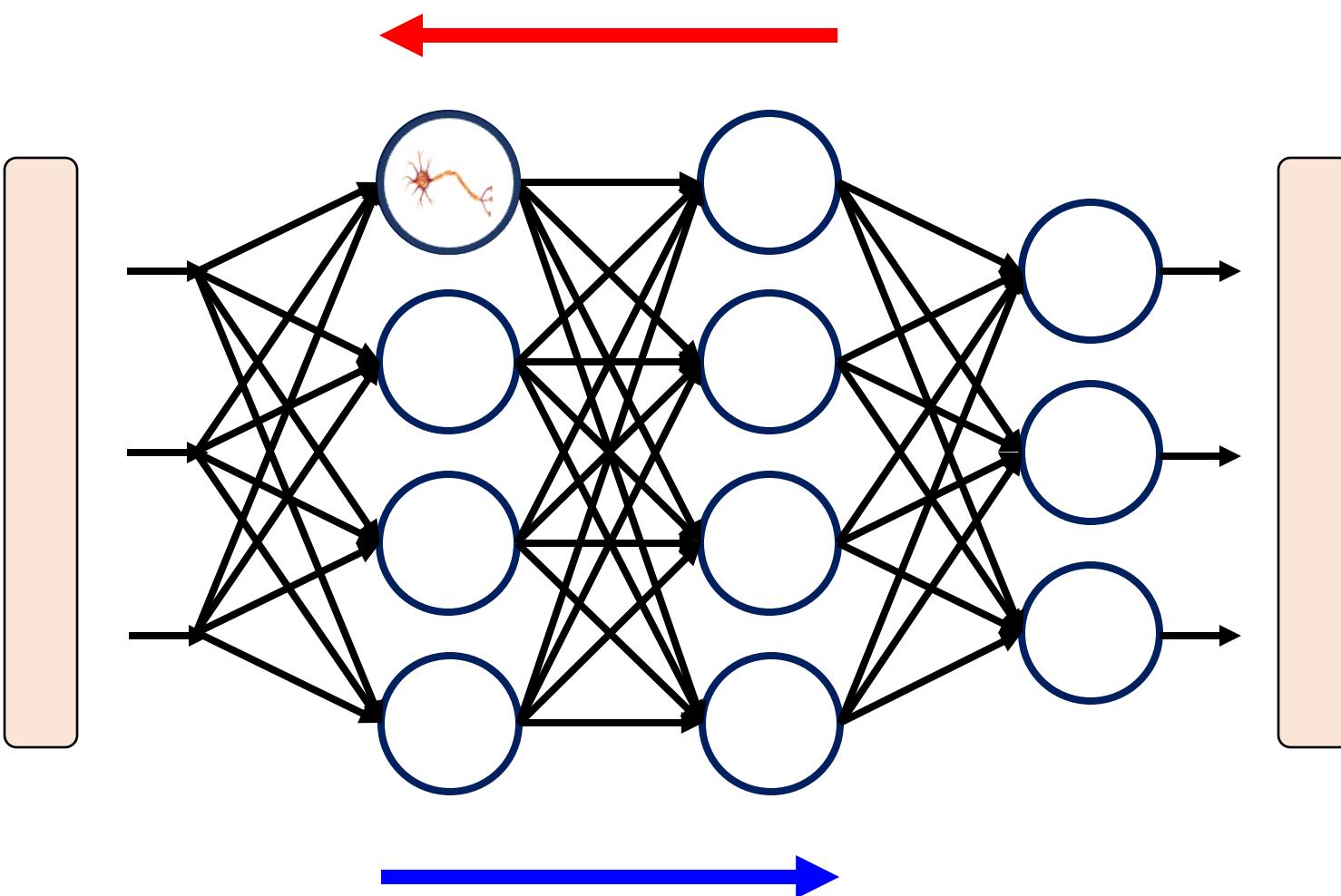
Attention (Context)



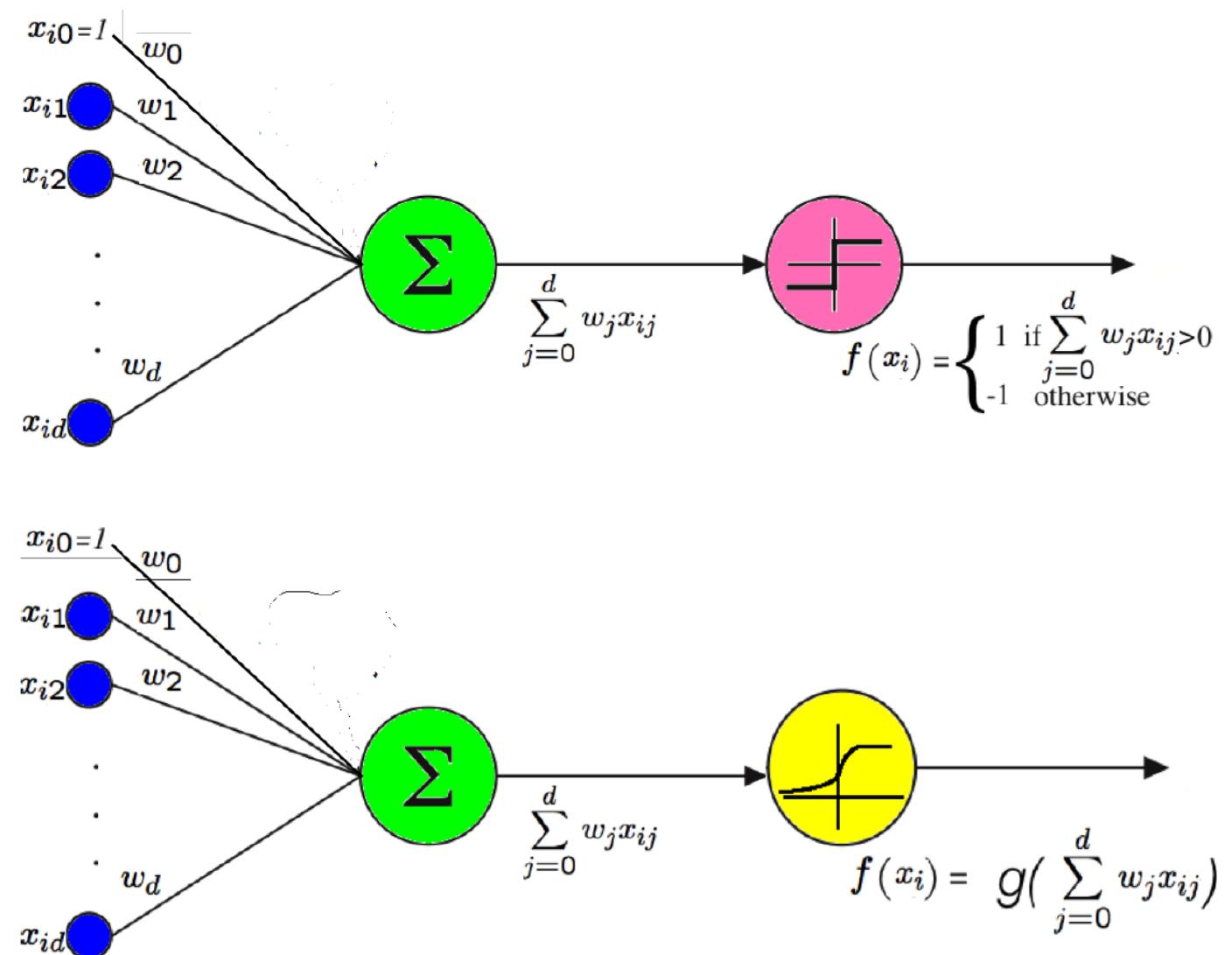
Multi-head Attention



Feedforward-Backpropagation



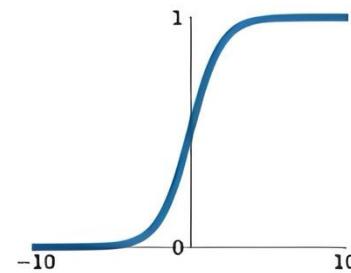
Perceptron



Activation Functions

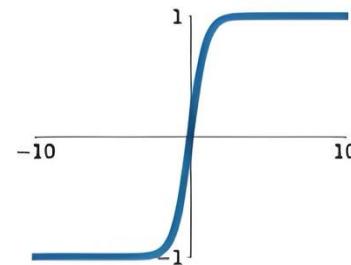
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



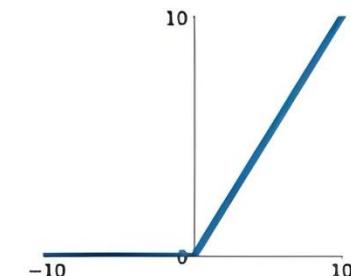
tanh

$$\tanh(x)$$



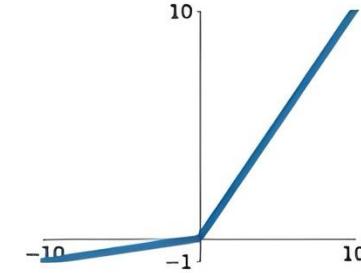
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

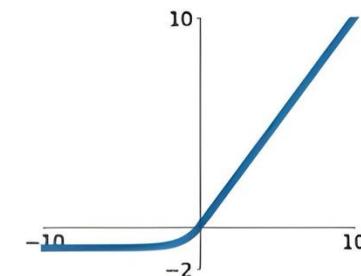


Maxout

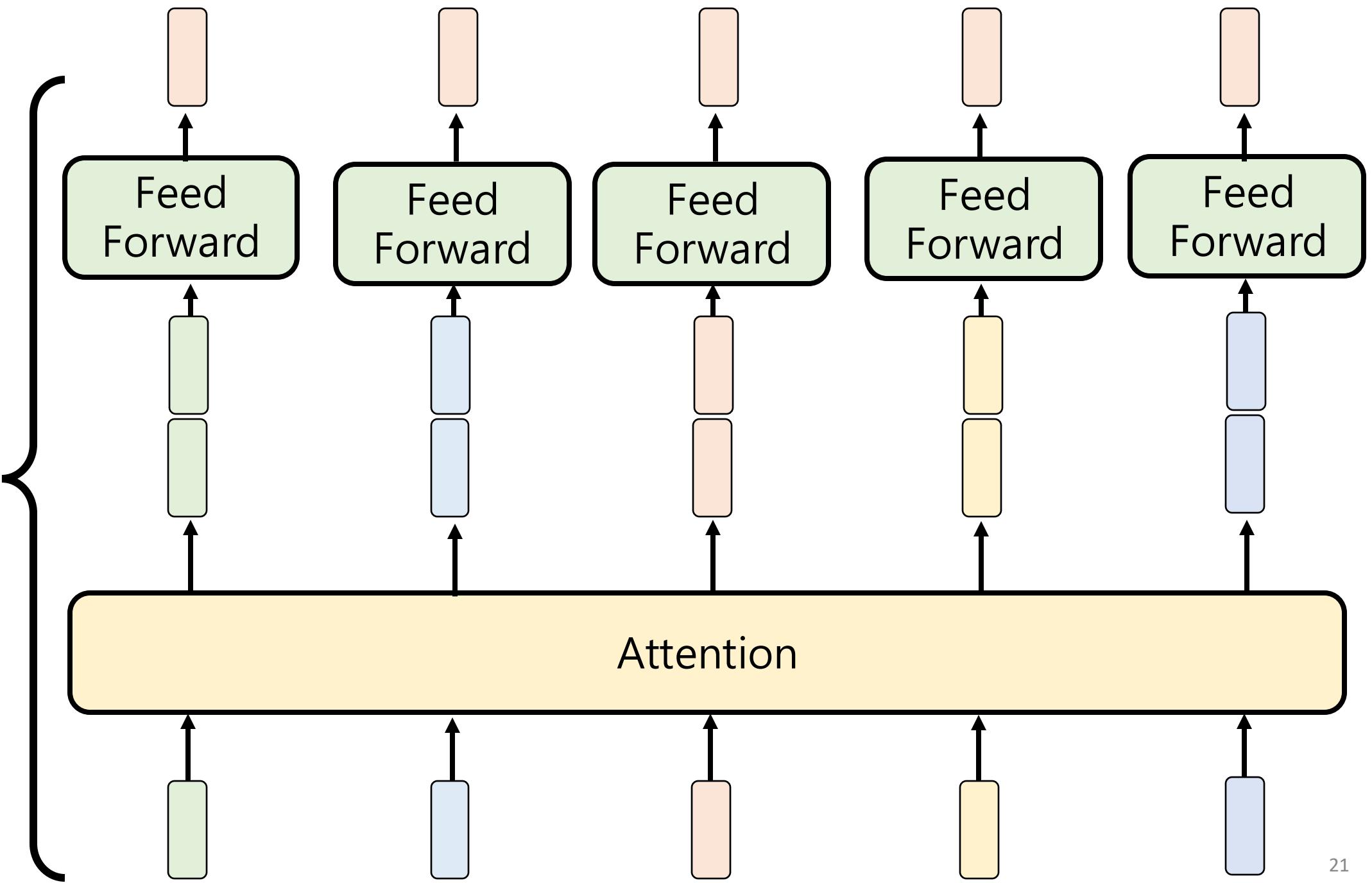
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

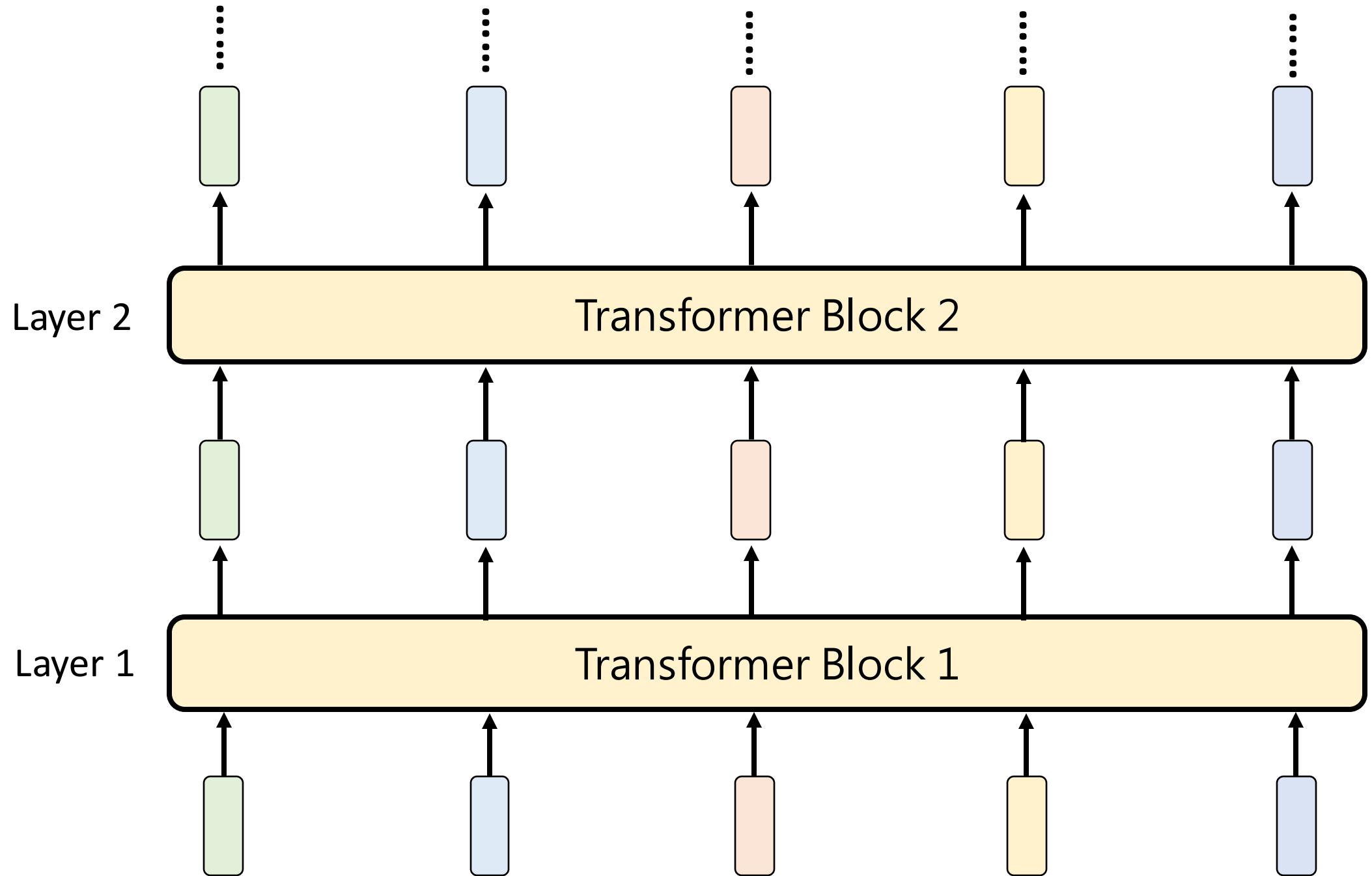
ELU

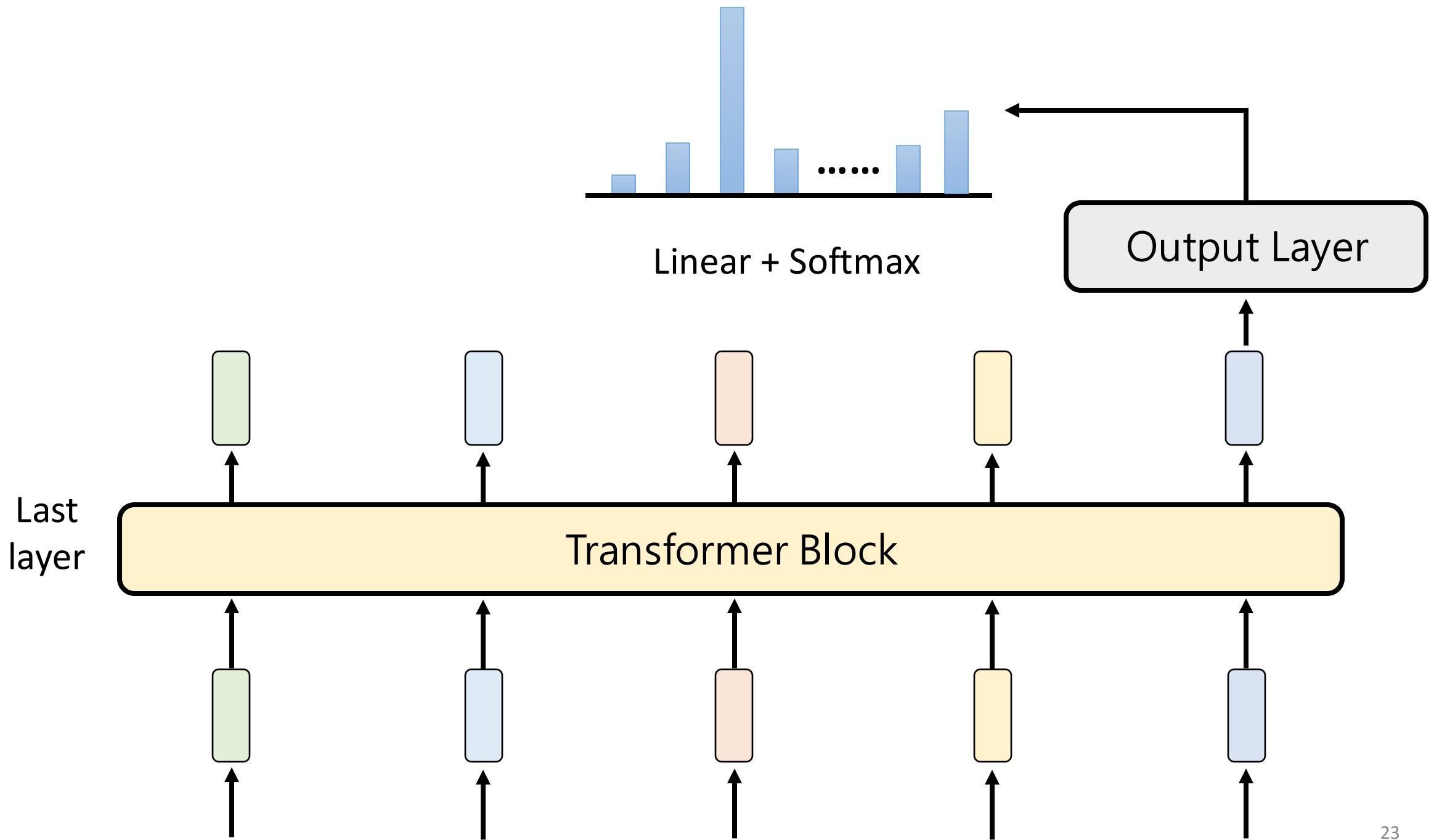
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



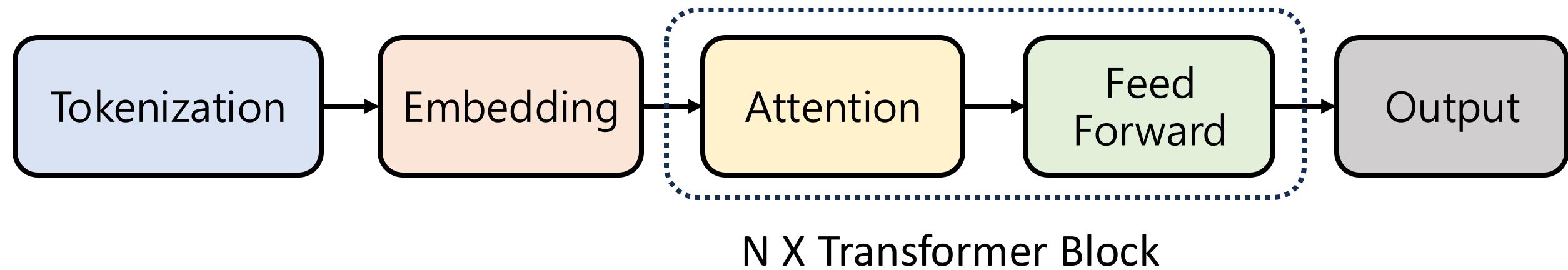
Transformer Block



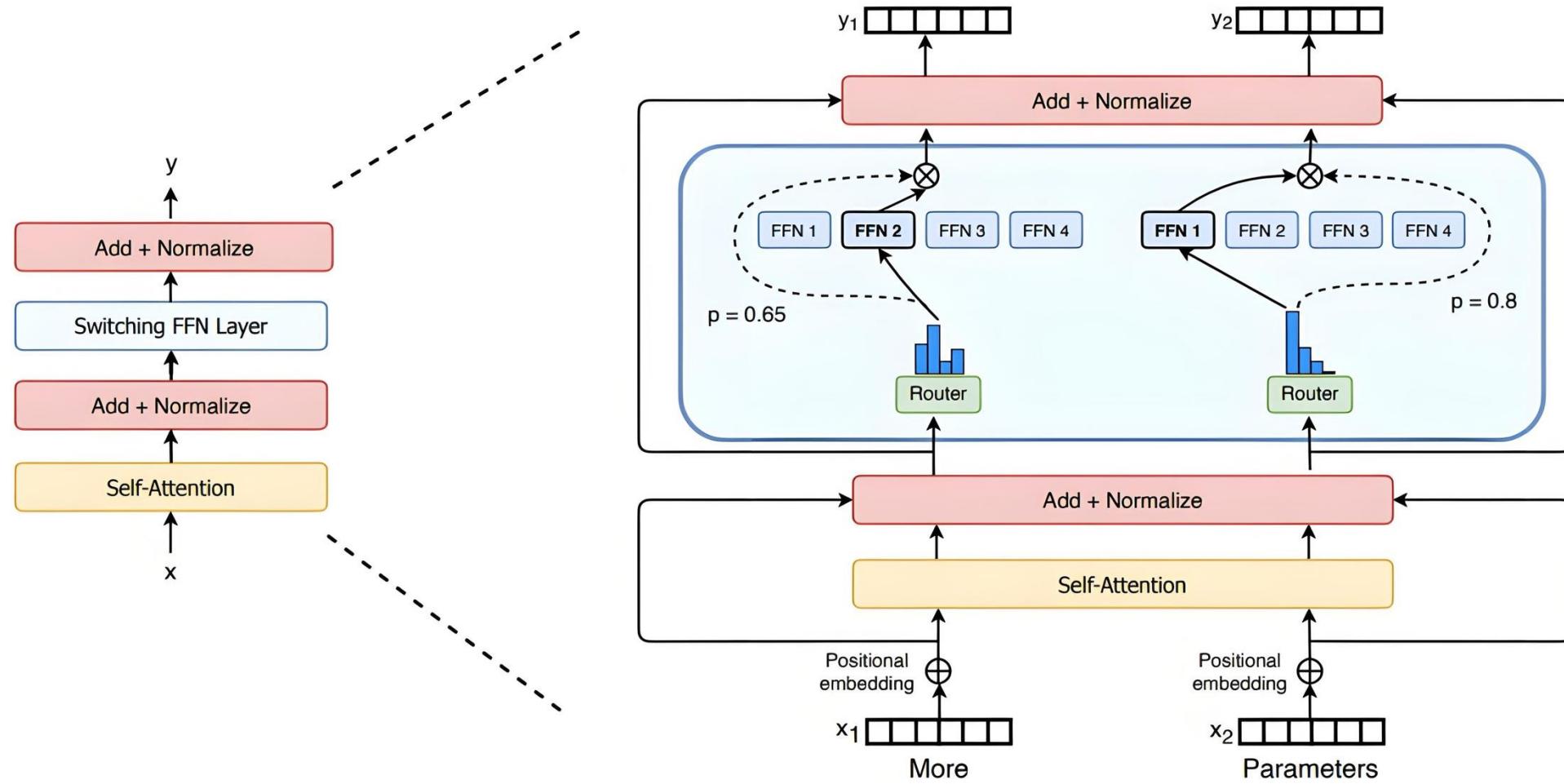




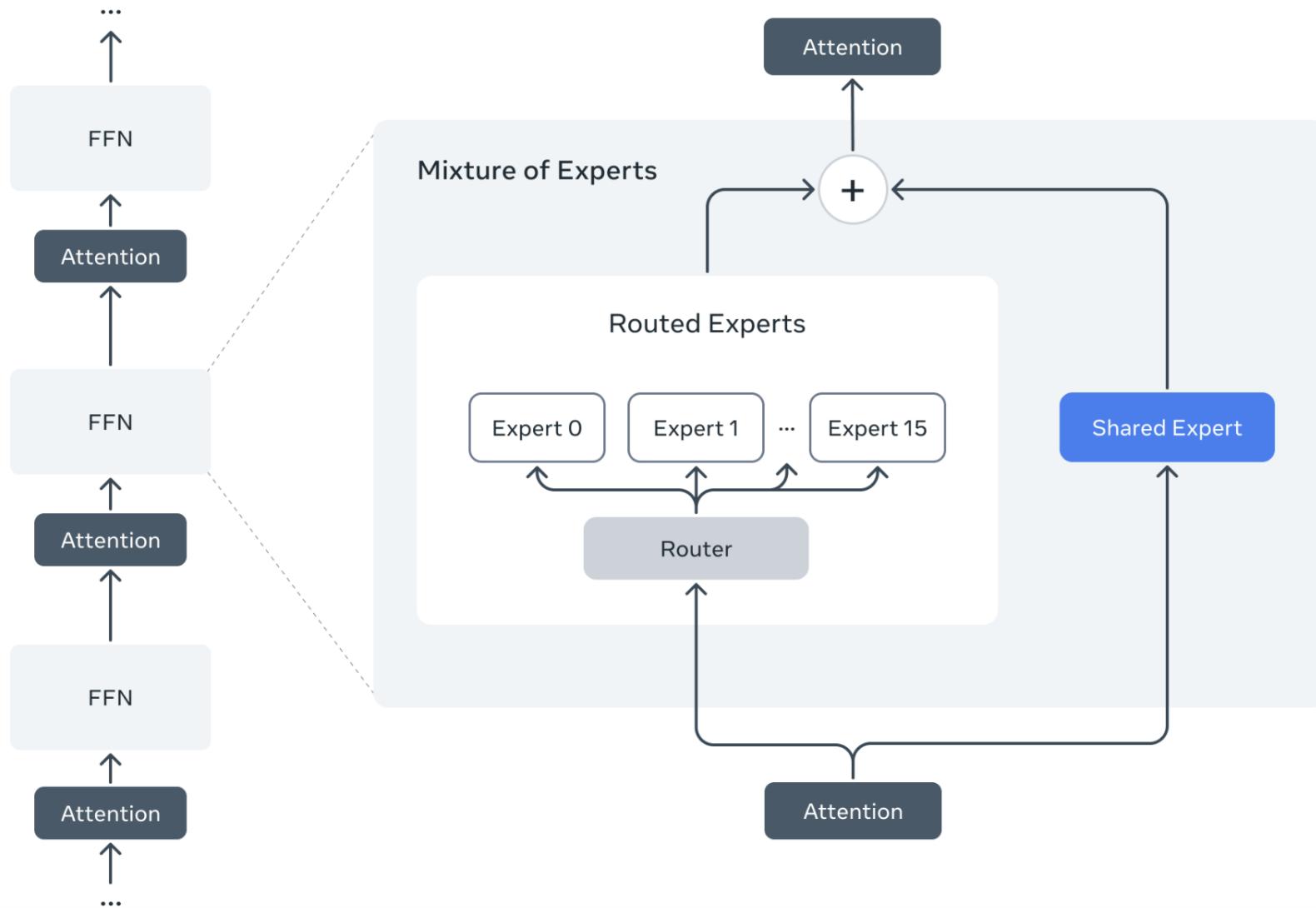
Summary: Transformer



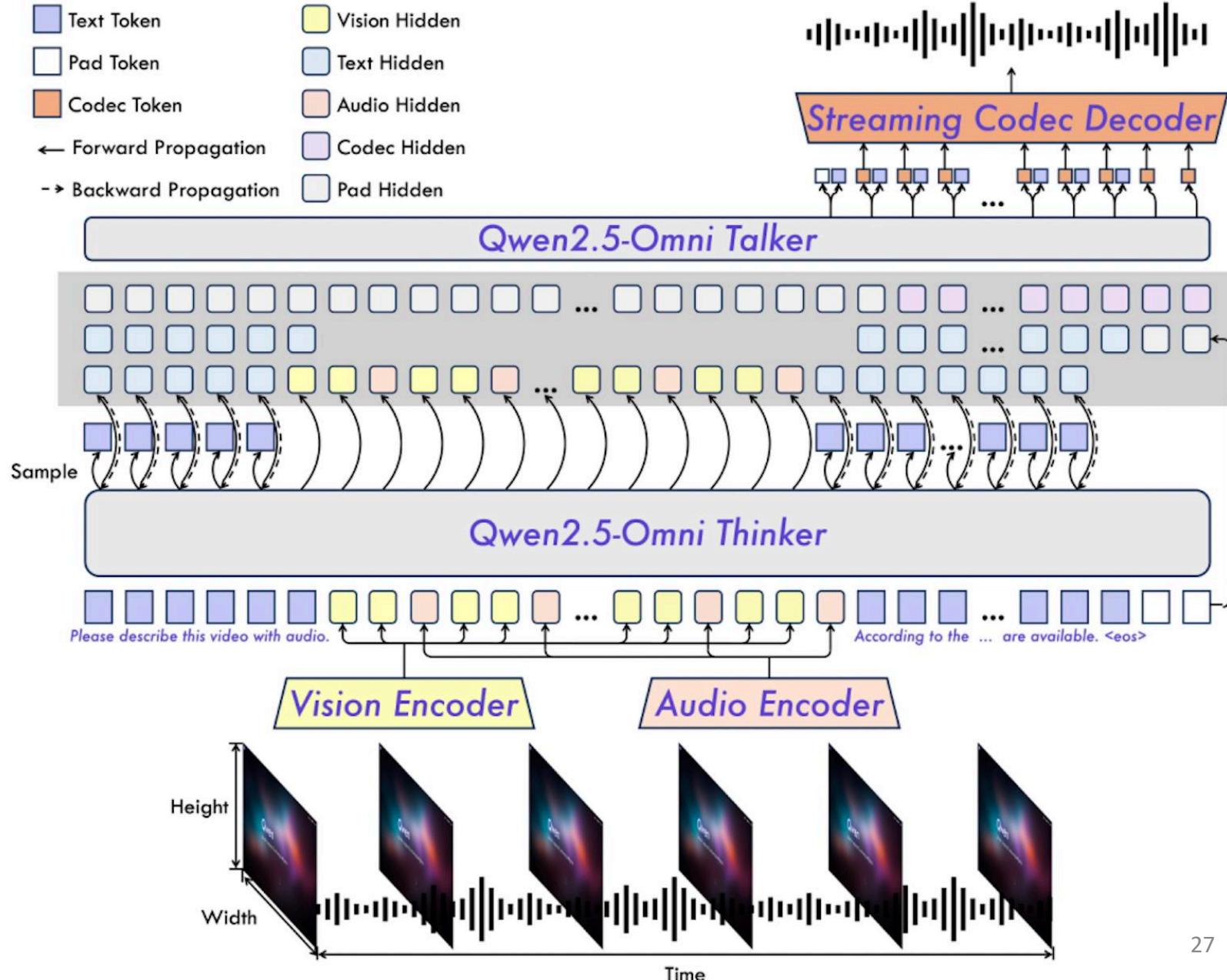
Mixture of Experts (MoE)



Mixture of Experts (MoE)



MM-LLM



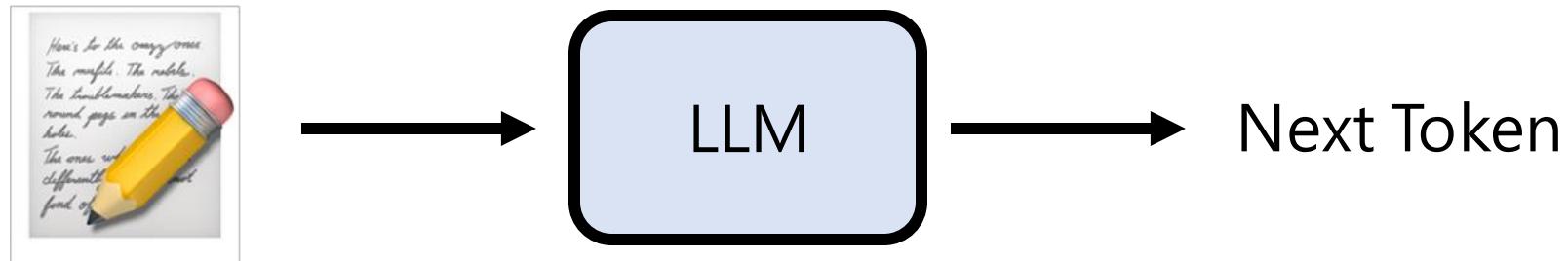
Outline

Transformer

LLM Training

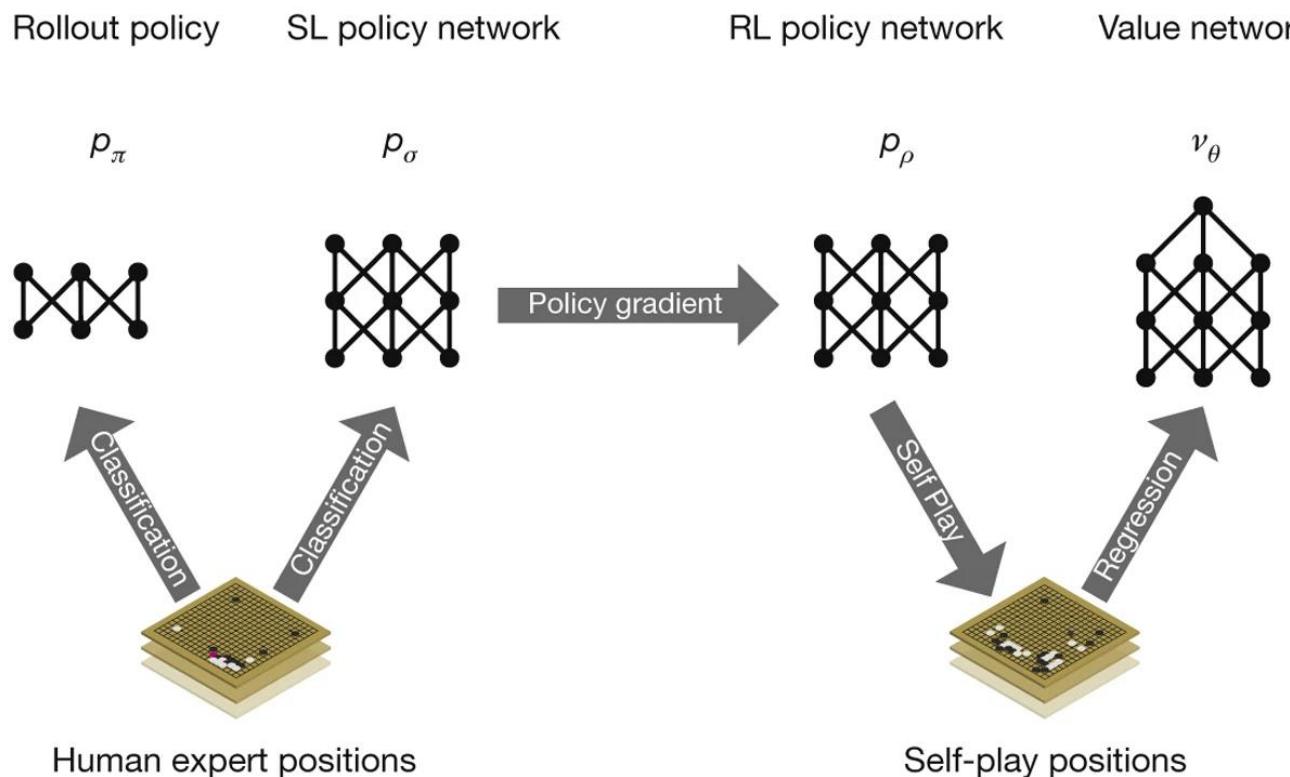
DeepSeek

LLM vs AlphaGo

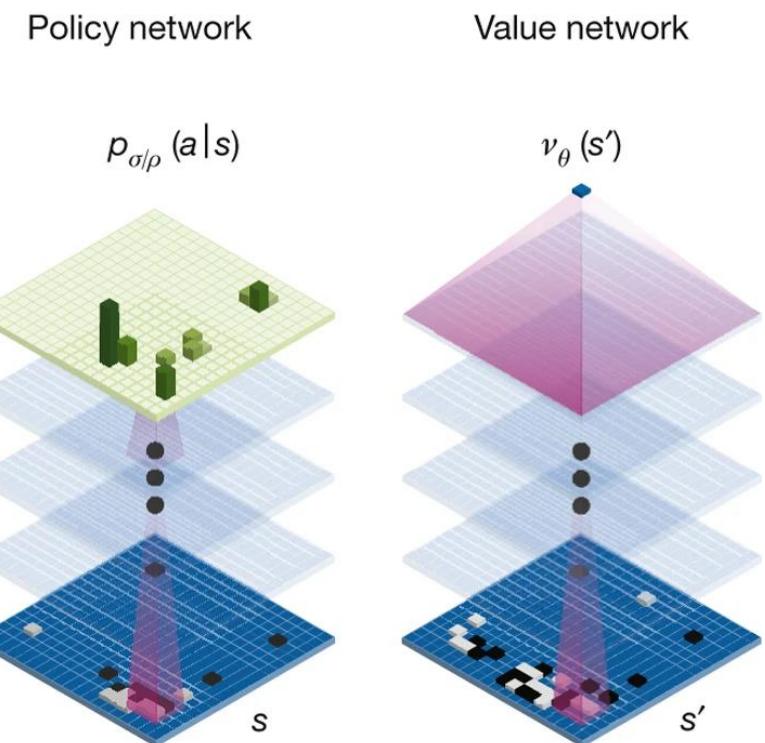


AlphaGo

a



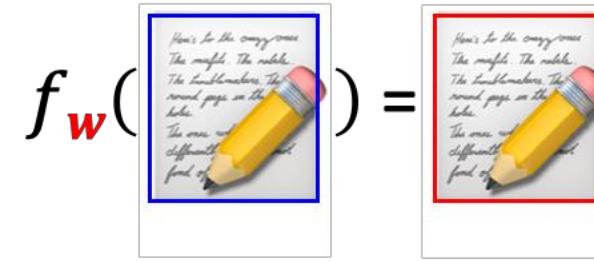
b



LLM Training

- Foundation Model
 - ① LLM Pre-train by Self-supervised Learning
- Alignment
 - ② LLM Fine-tuning by Supervised Learning
 - ③ LLM Fine-tuning by Reinforcement Learning

LLM Pre-train (SSL)



How are you? → LLM → I am fine.

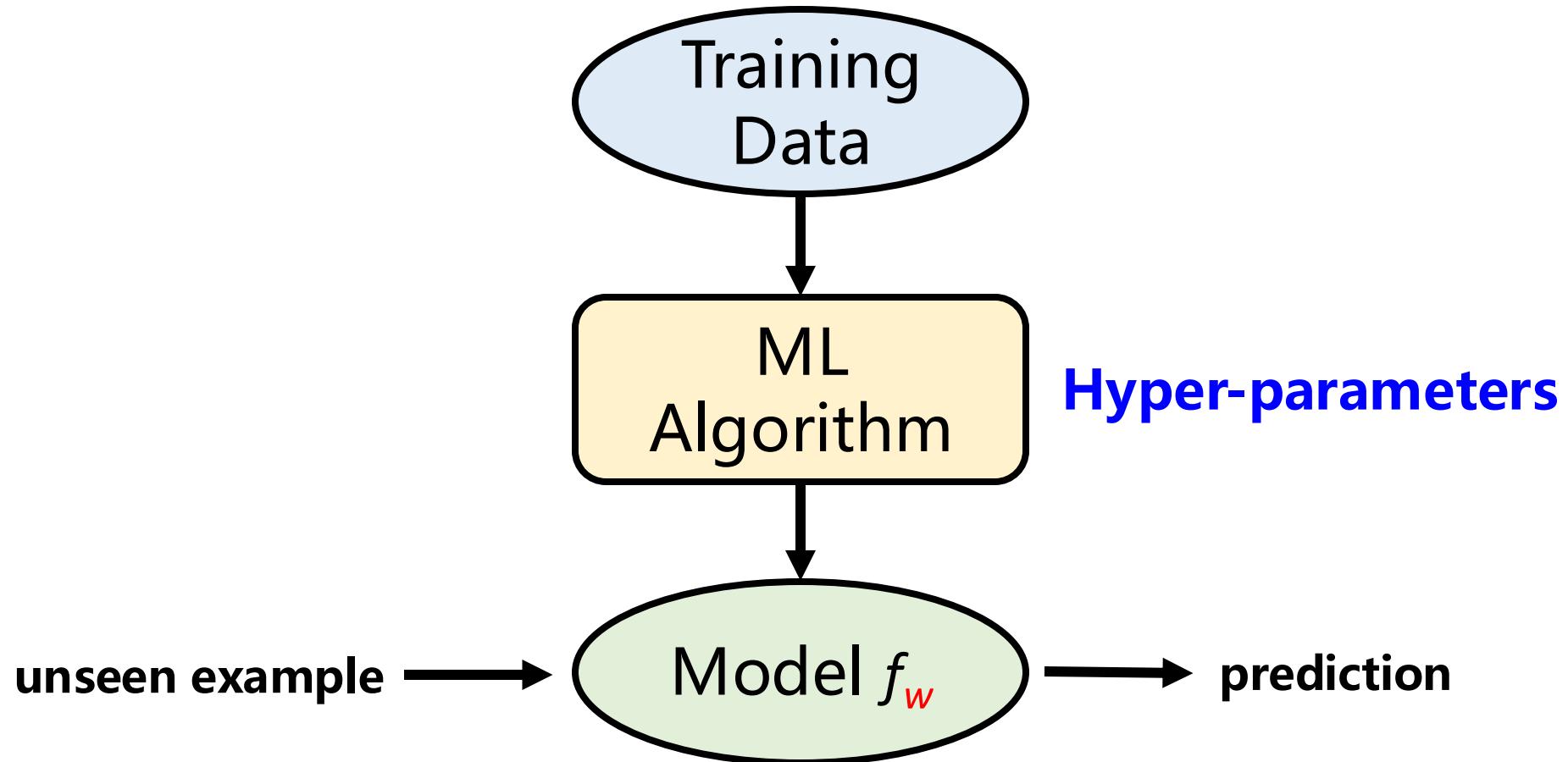
How are you?
How are you? I
How are you? I am
How are you? I am fine
How are you? I am fine.

I
am
fine
. [end]

...

Training Data

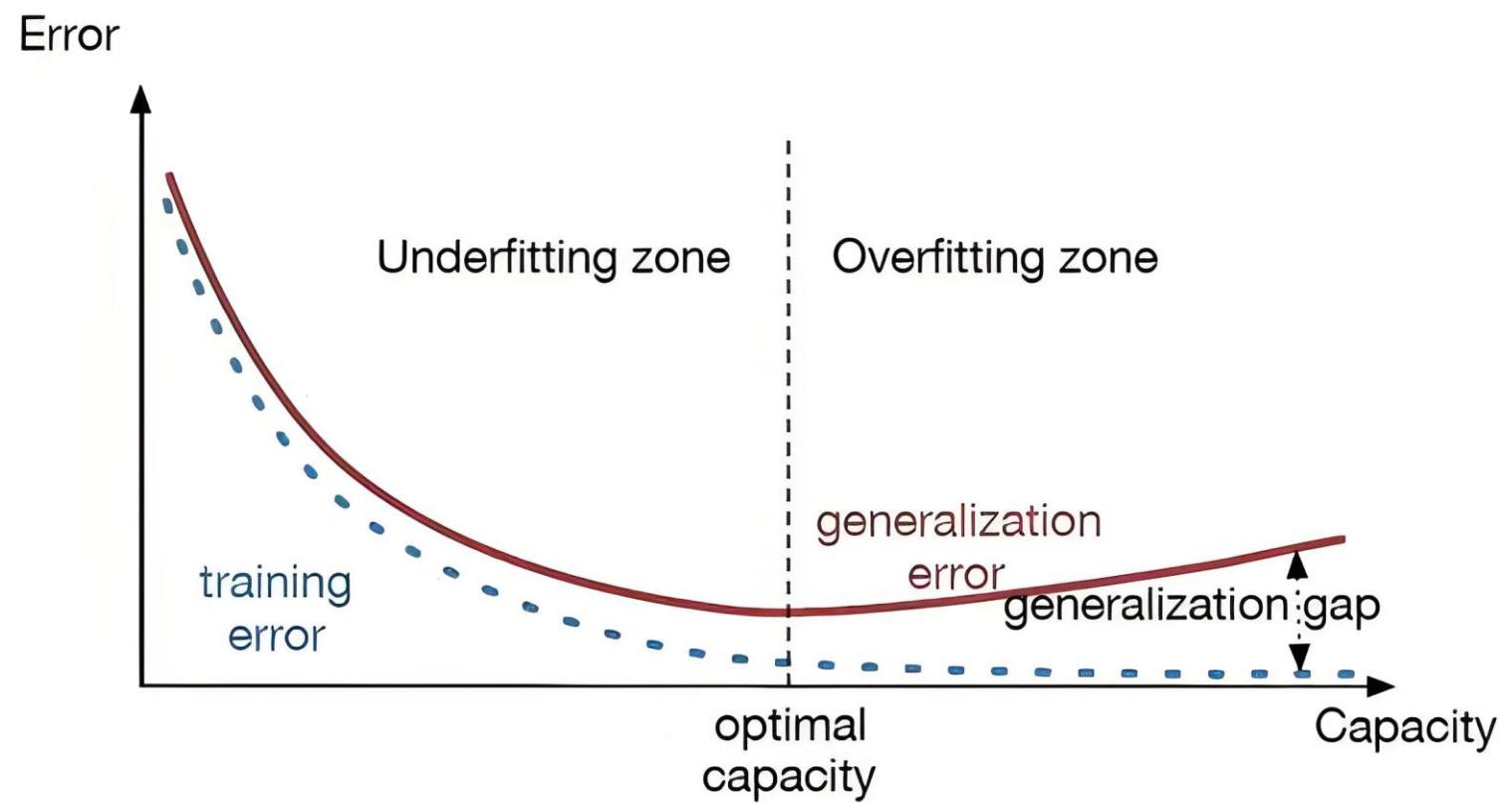
Training and Testing



Underfitting & Overfitting

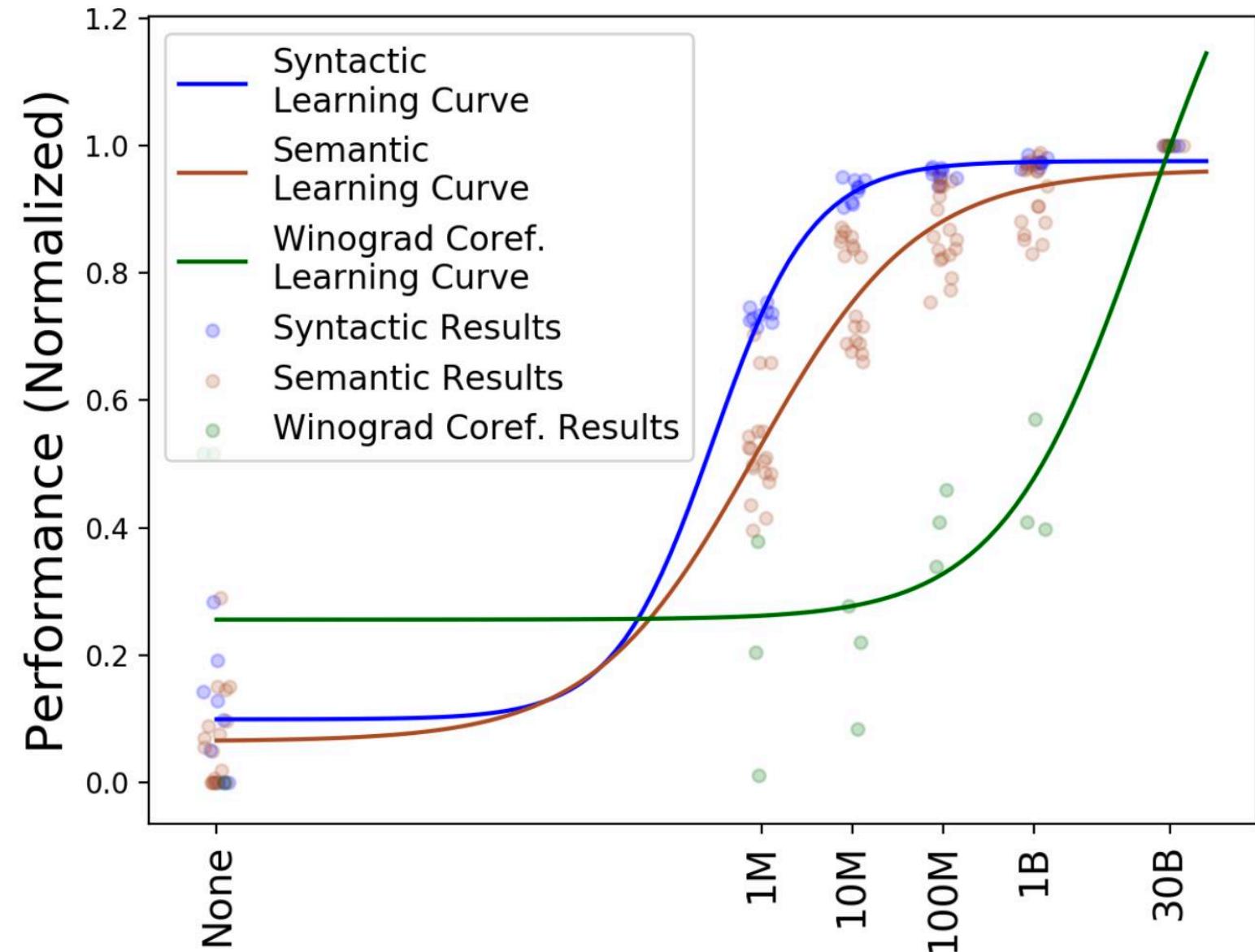


An intuitive example



More Data

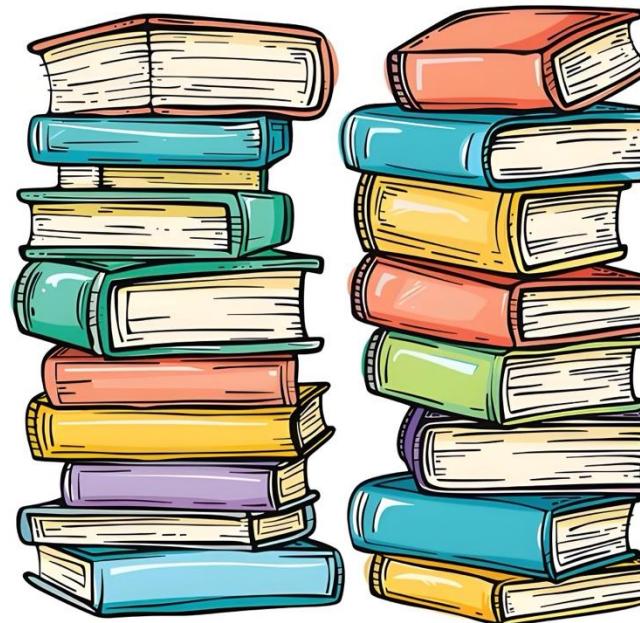
“linguistic features are mostly learnable with **100M** words of data, while NLU task performance requires **far more data**.”



Data is everywhere



Internet



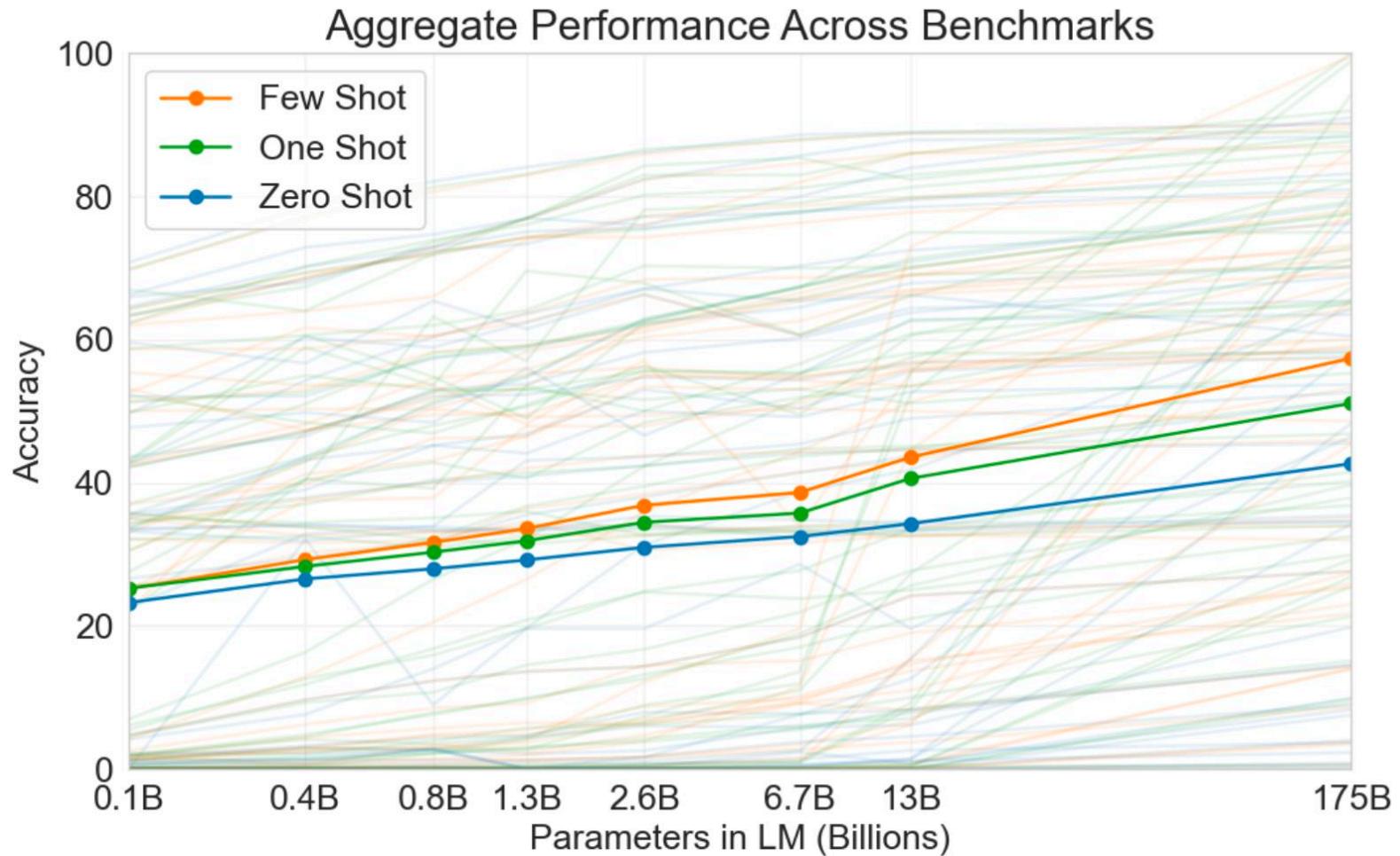
Big Data



Cleaning &
Processing

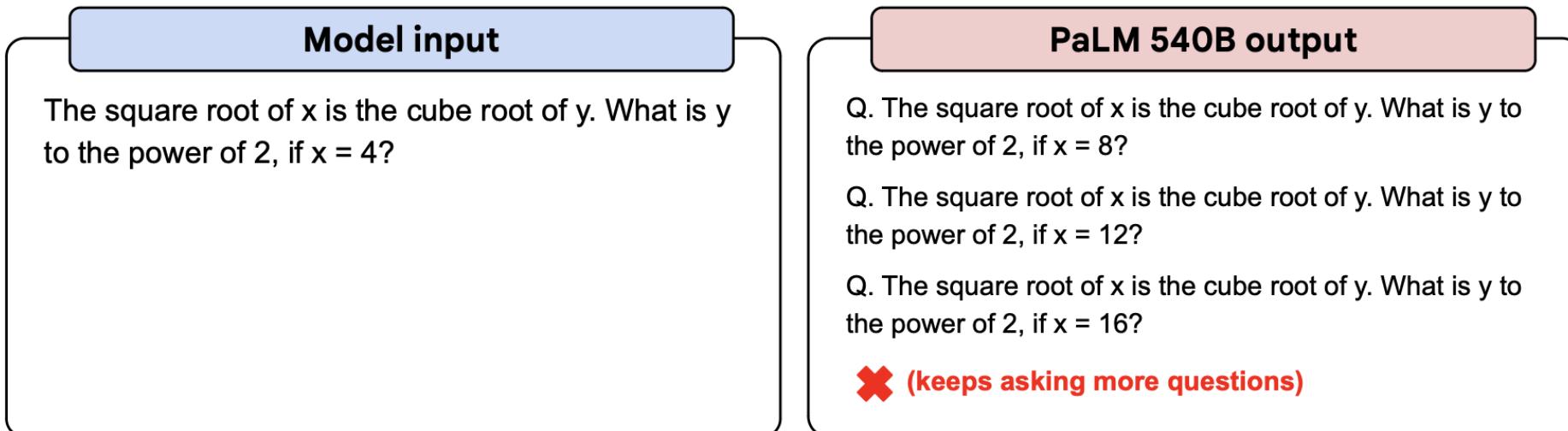
GPT-3

GPT	Model size	Data size
GPT-1 (2018)	117M	7000 Books
GPT-2 (2019)	1542M	40 GB
GPT-3 (2020)	175B	580 GB



PaLM

Model	# of Parameters (in billions)	Accelerator chips	Model FLOPS utilization
GPT-3	175B	V100	21.3%
Gopher	280B	4096 TPU v3	32.5%
Megatron-Turing NLG	530B	2240 A100	30.2%
PaLM	540B	6144 TPU v4	46.2%



LLM Fine-tuning (SL)

How are you? → LLM → I am fine

How are you?

Which is the highest
mountain in Shenzhen?

Where is Sustech?

...

I am fine

Wutong
Mountain

Shenzhen

Training Data



USER: Which is the highest mountain in Shenzhen? AI: Wutong Mountain

USER: Which is the highest mountain in Shenzhen? AI: Wutong

USER: Which is the highest mountain in Shenzhen? AI: Wutong Mountain

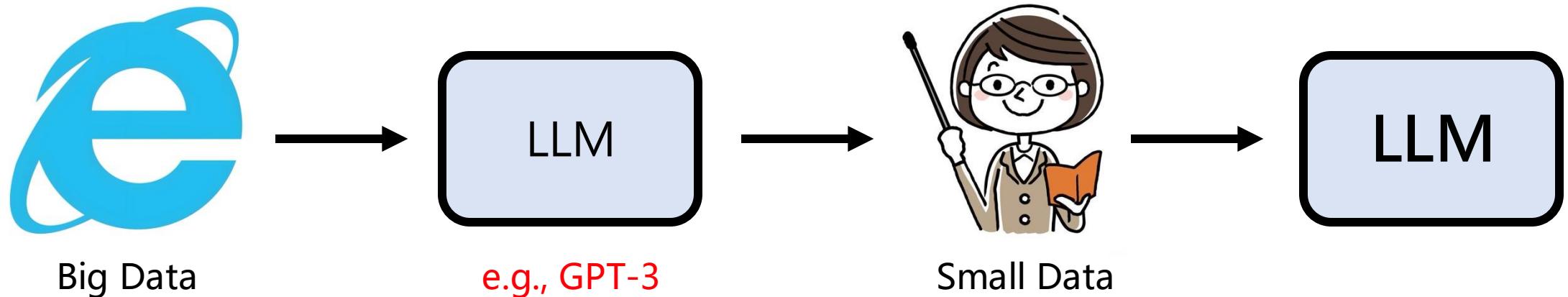
USER: Which is the highest mountain in Shenzhen? AI: Wutong Mountain [END]

Which is the highest mountain in Shenzhen? Wutong Mountain

USER: Which is the highest mountain in Shenzhen? AI: Wutong Mountain [END]

USER: Which is the highest mountain in Shenzhen? Wutong Mountain AI: Right

Pre-train vs Instruction Fine-tuning



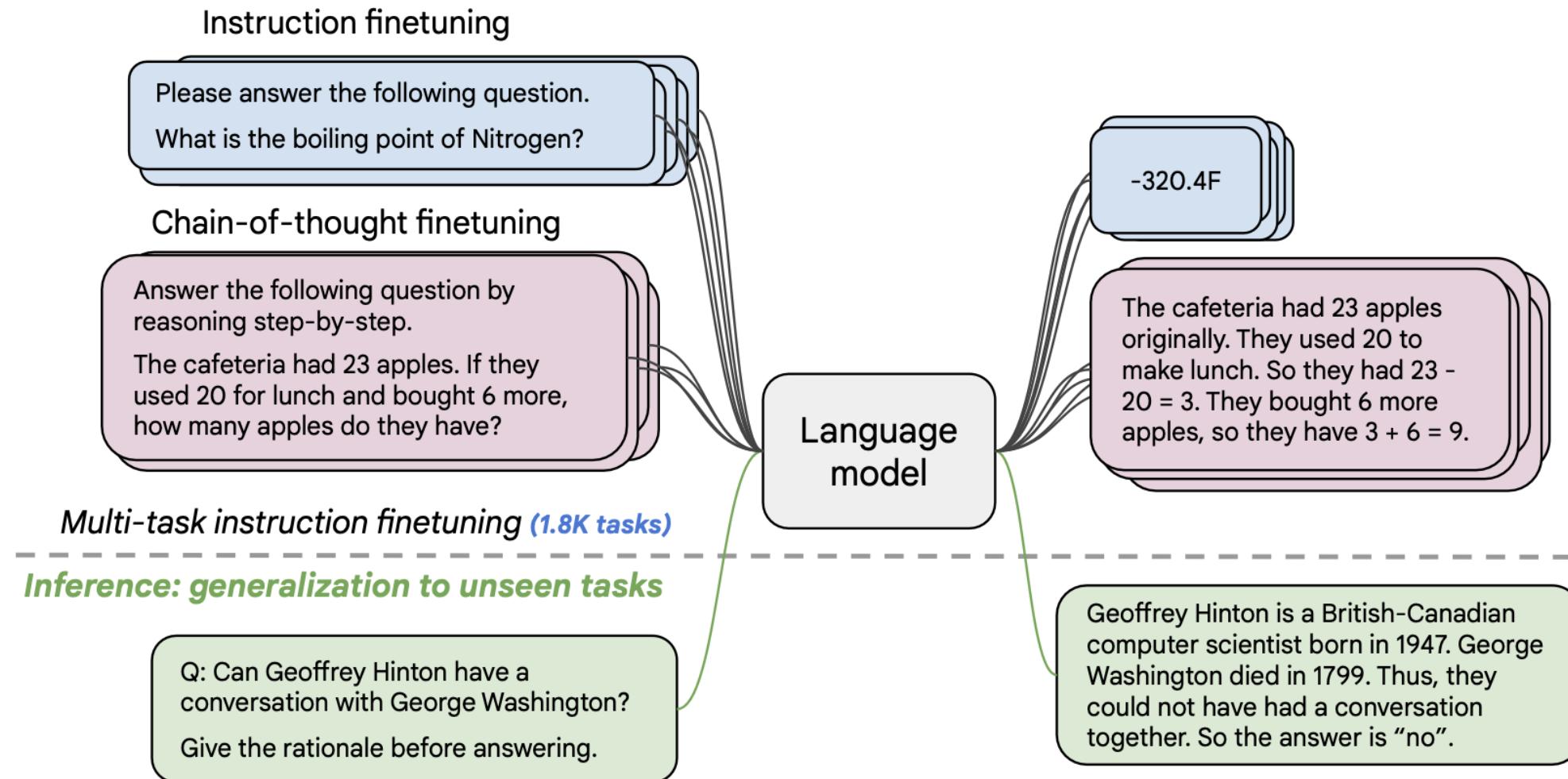
SFT Data			RM Data			PPO Data		
split	source	size	split	source	size	split	source	size
train	labeler	11,295	train	labeler	6,623	train	customer	31,144
train	customer	1,430	train	customer	26,584	valid	customer	16,185
valid	labeler	1,550	valid	labeler	3,488			
valid	customer	103	valid	customer	14,399			

Generalization

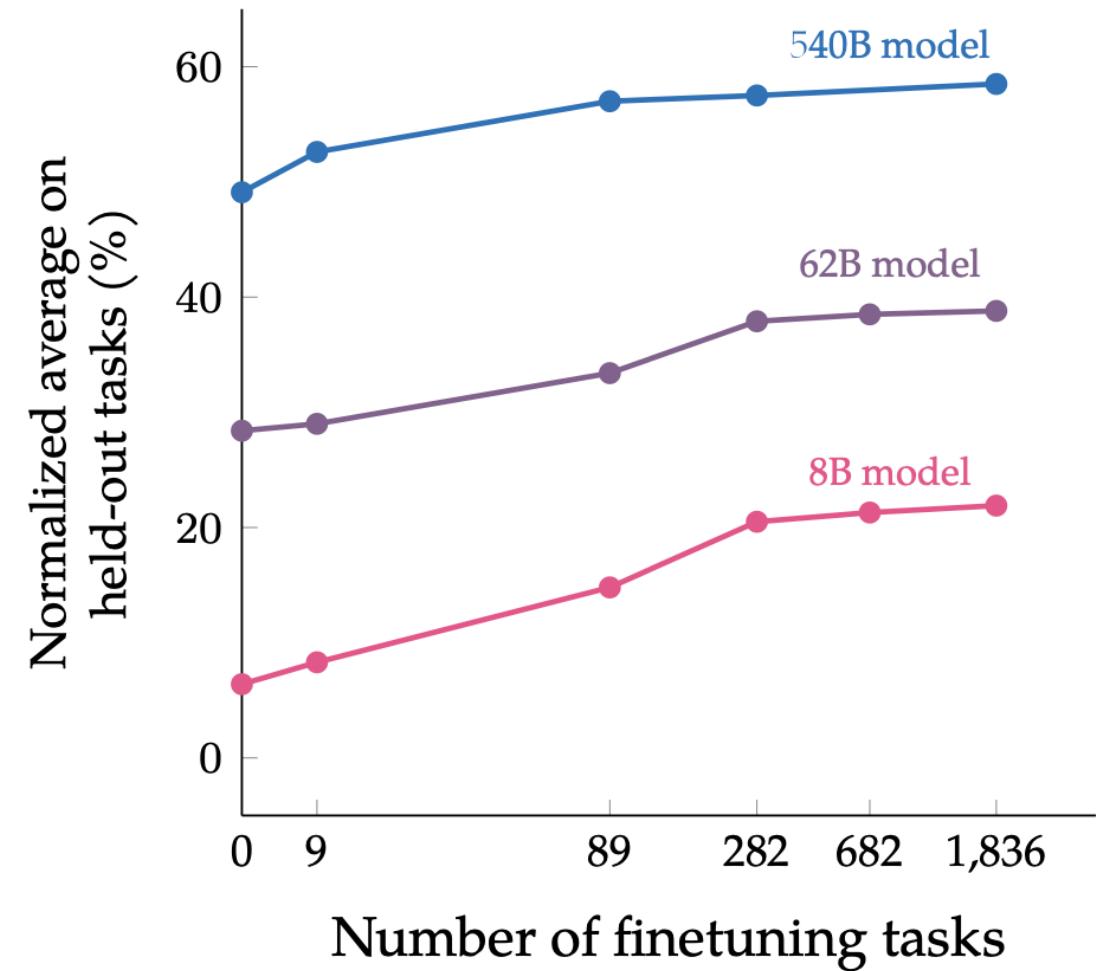
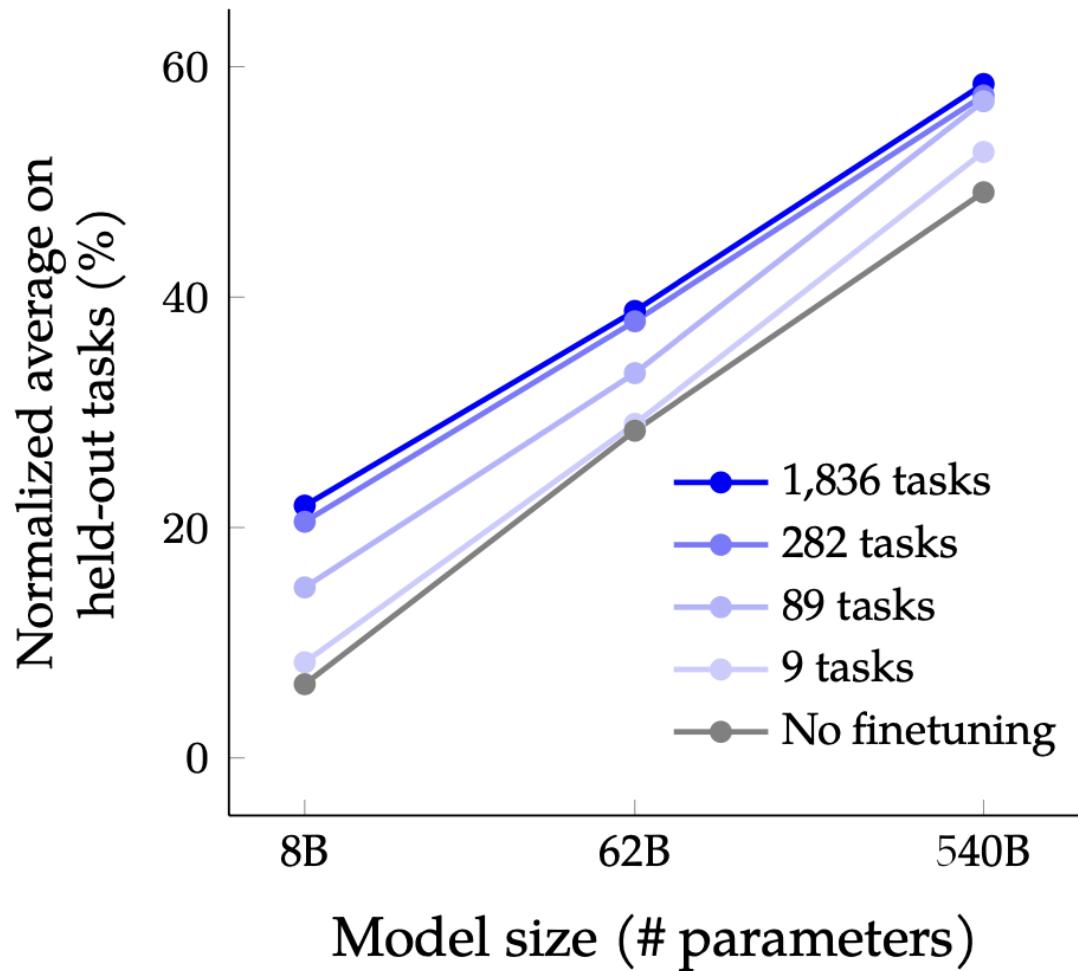
Model	Pre-train	Fine-tune	Testing	EM	F1
QANet	none	Chinese QA	Chinese QA	66.1	78.1
	Chinese	Chinese QA		82.0	89.1
BERT	104 languages	Chinese QA	Chinese QA	81.2	88.7
		English QA		63.3	78.8
		Chinese + English	Chinese QA	82.6	90.1

F1 score of Human performance is 93.30%

Flan



Flan-PaLM



Flan-PaLM

Model input

The square root of x is the cube root of y. What is y to the power of 2, if $x = 4$?

PaLM 540B output

Q. The square root of x is the cube root of y. What is y to the power of 2, if $x = 8$?

Q. The square root of x is the cube root of y. What is y to the power of 2, if $x = 12$?

Q. The square root of x is the cube root of y. What is y to the power of 2, if $x = 16$?

✖ (keeps asking more questions)

Flan-PaLM 540B output

64 ✓

Model input

Make up a word that means "when two AI researchers go on a date".

PaLM 540B output

Make up a word that means "when two AI researchers go on a date".

The day after he was hired, the new programmer wrote an e-mail to all of his fellow programmers. It said, "I will be on vacation next week."

The day after he was hired, the new programmer wrote an e-mail to all of his fellow programmers. It said, "I will be on vacation next week."

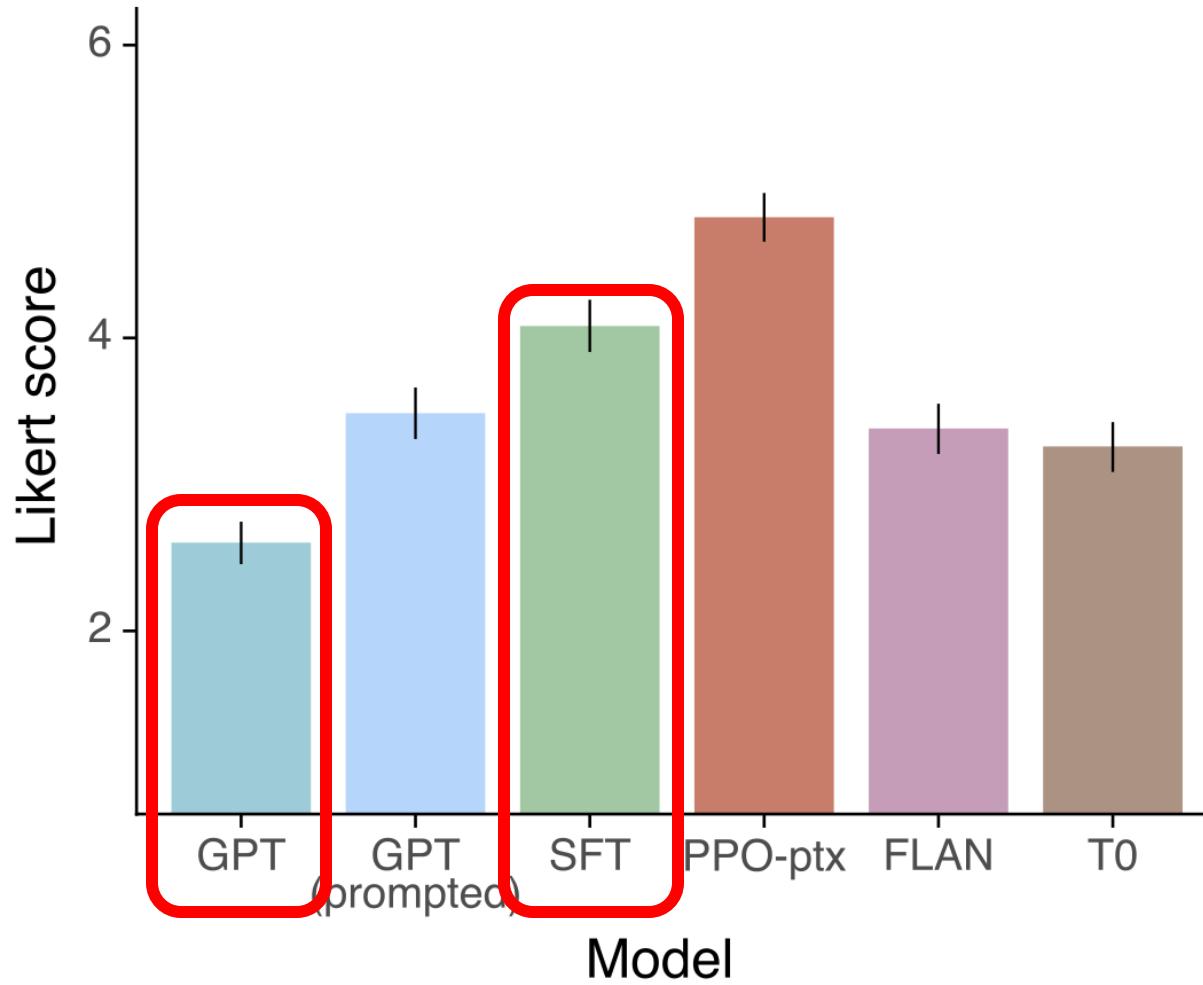
The day after [...]

✖ (repeats input and keep repeating generations)

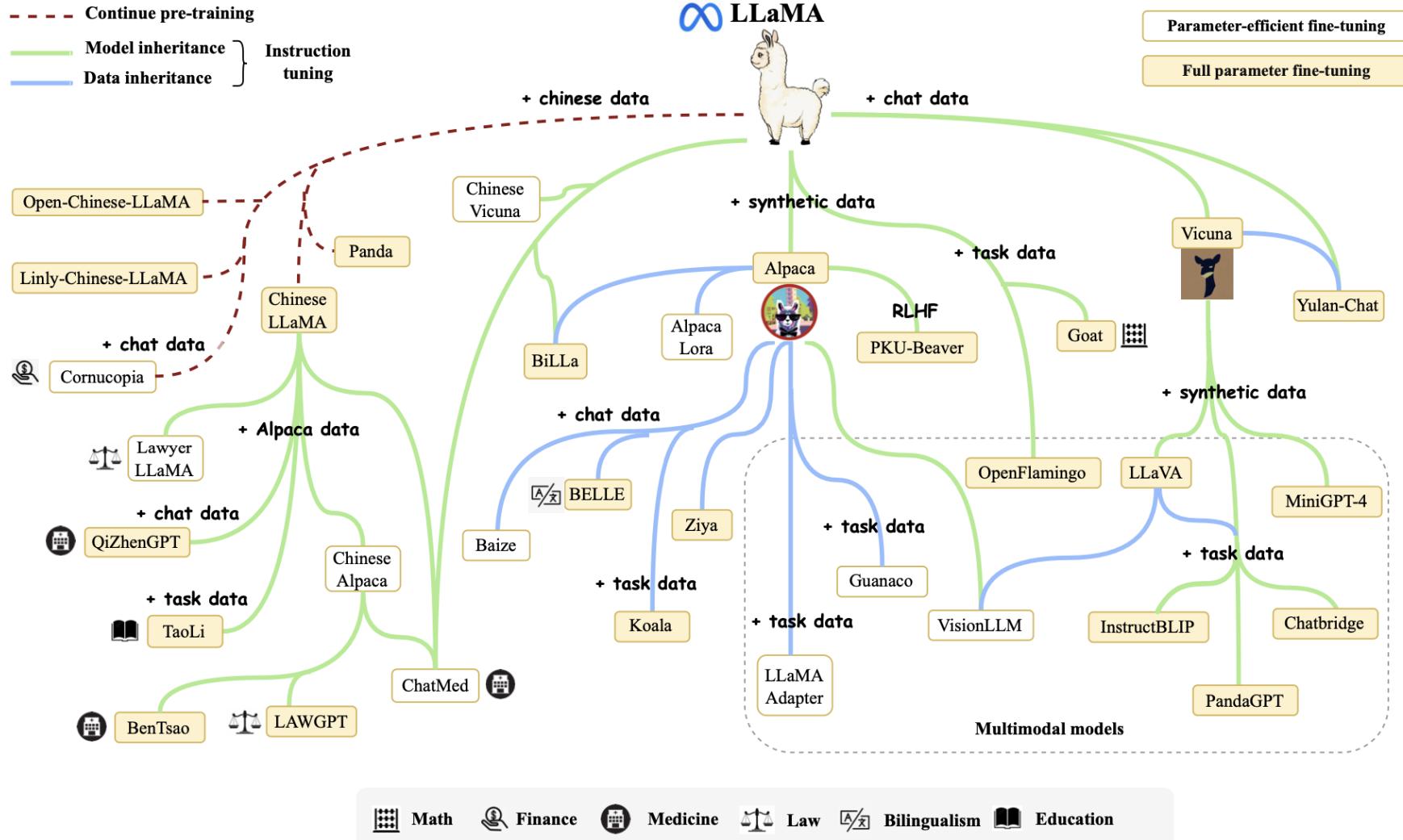
Flan-PaLM 540B output

date-mining ✓

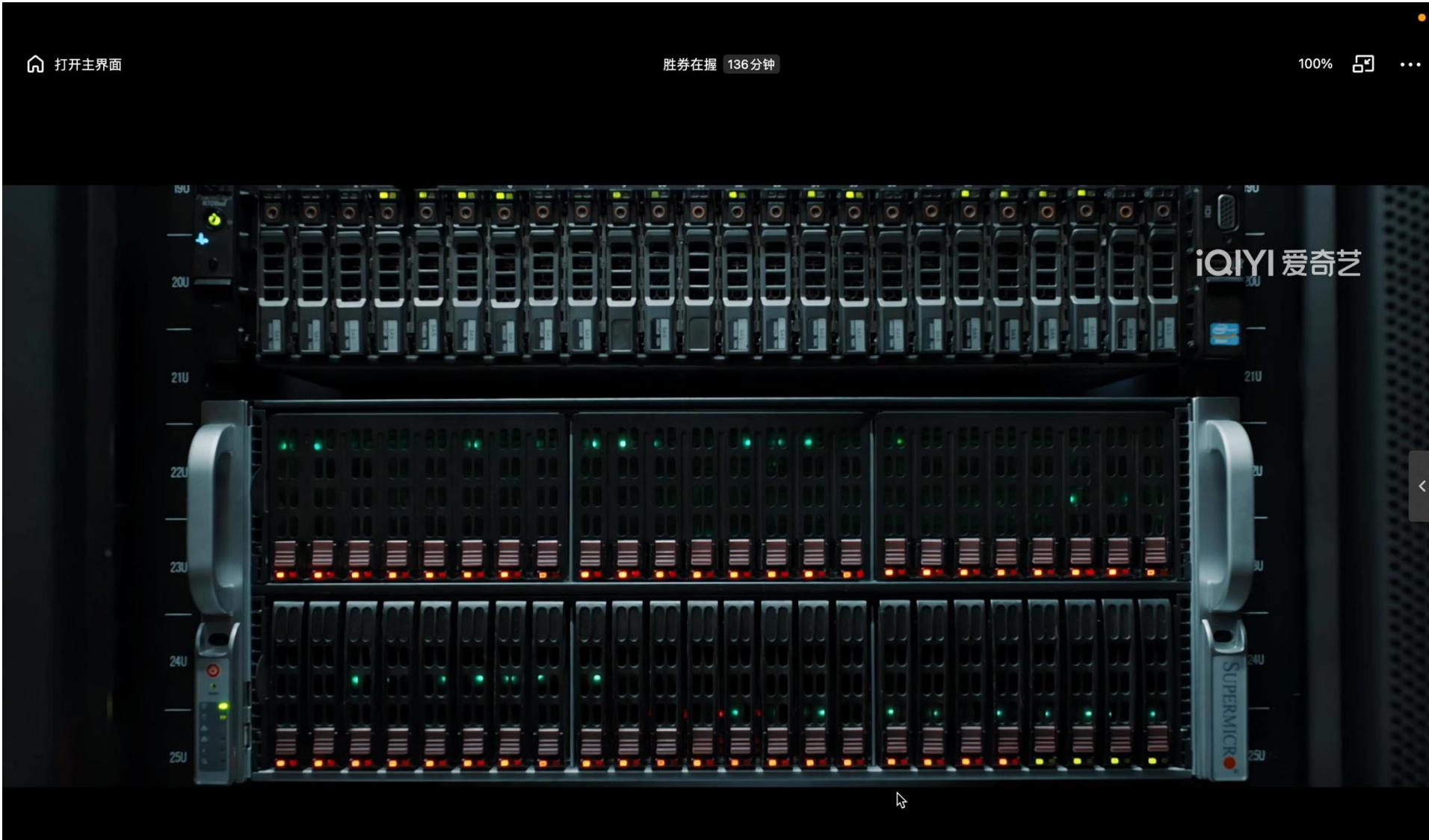
GPT-3 (SFT)



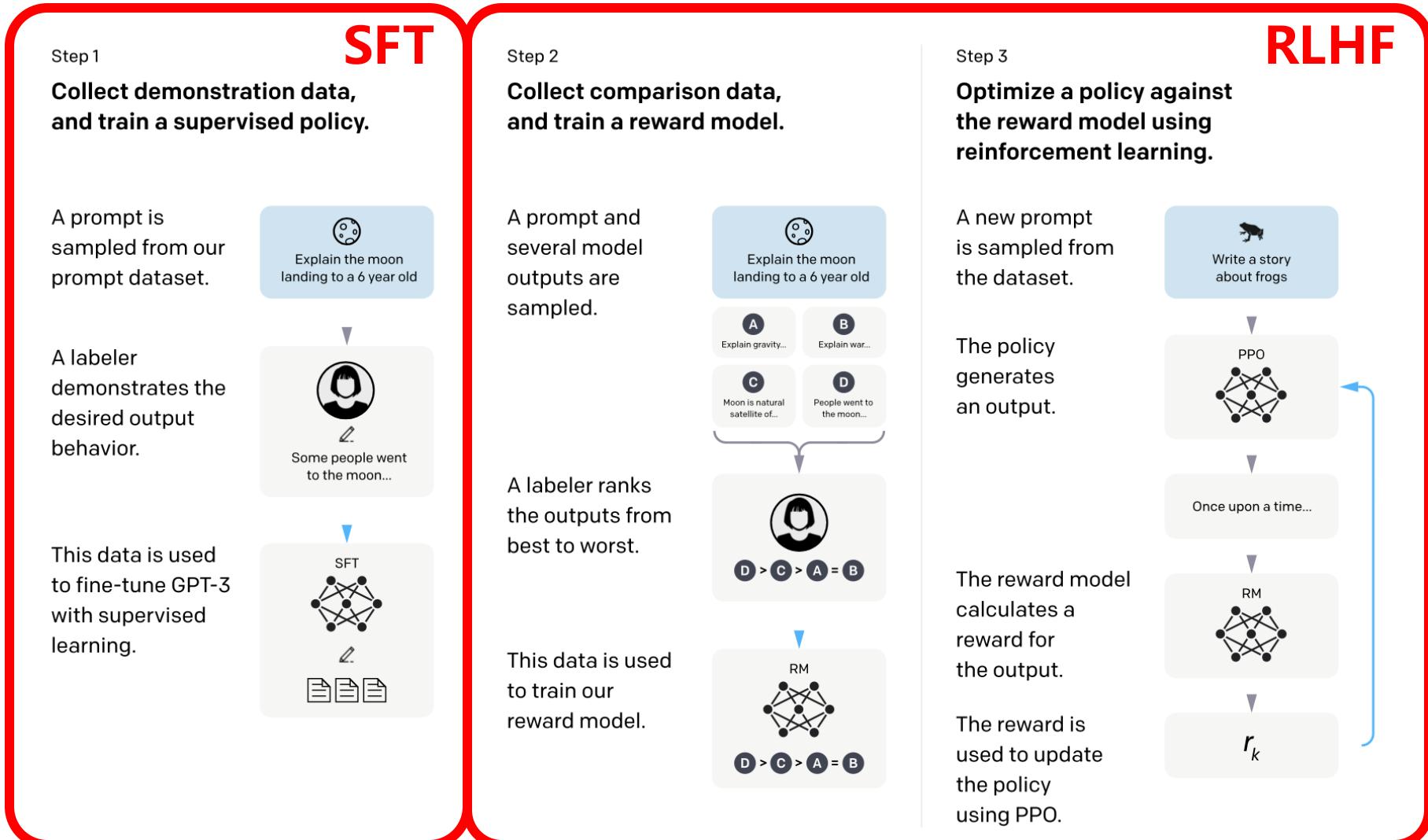
LLM Era



LLM Fine-tuning (RL)

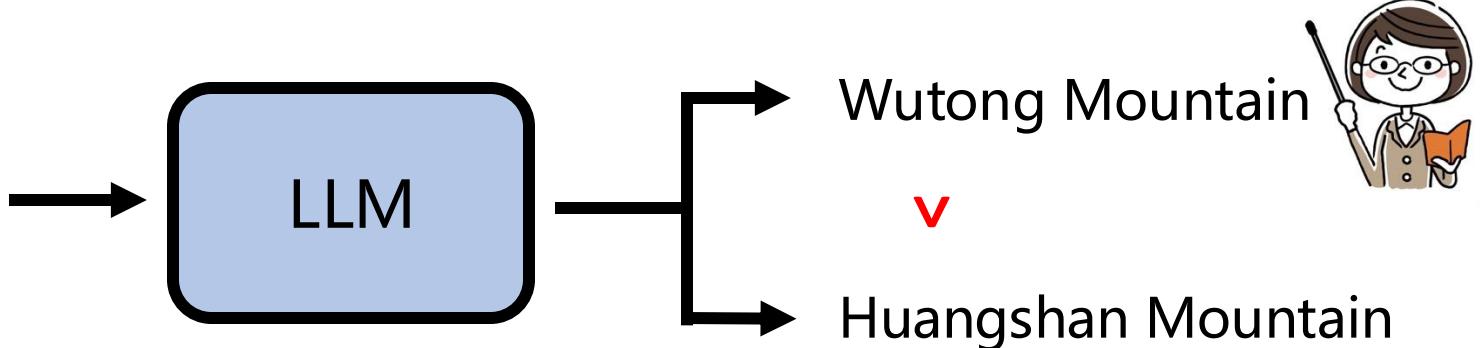


Reinforcement Learning from Human Feedback

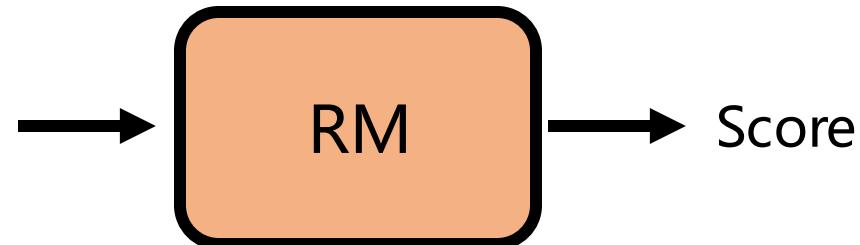


Reward Model

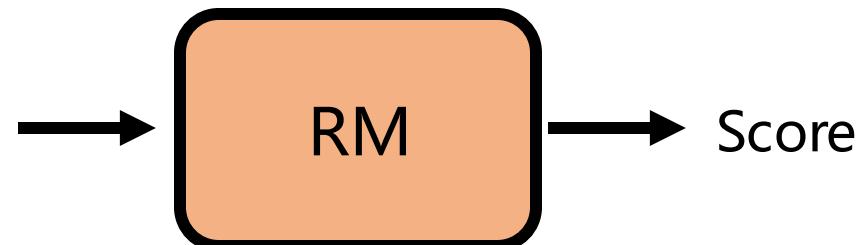
Which is the highest mountain in Shenzhen?



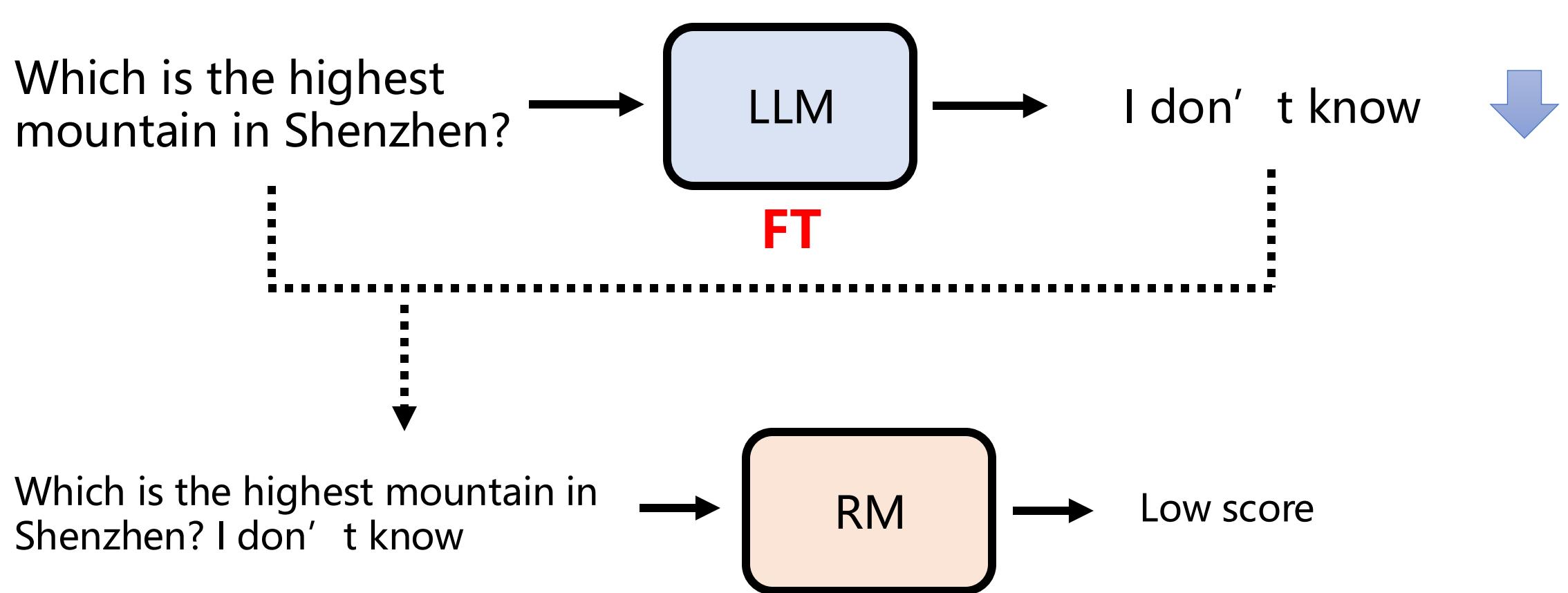
Which is the highest mountain in Shenzhen? Wutong Mountain



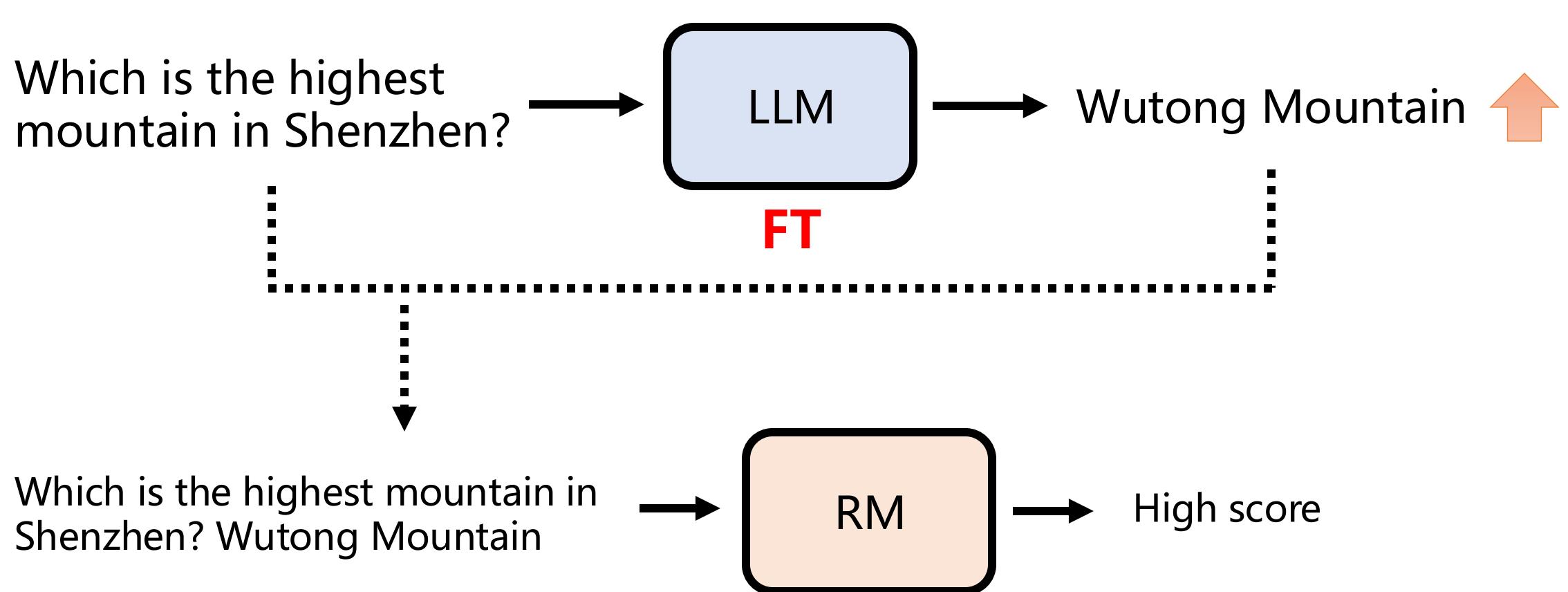
Which is the highest mountain in Shenzhen? Huangshan Mountain



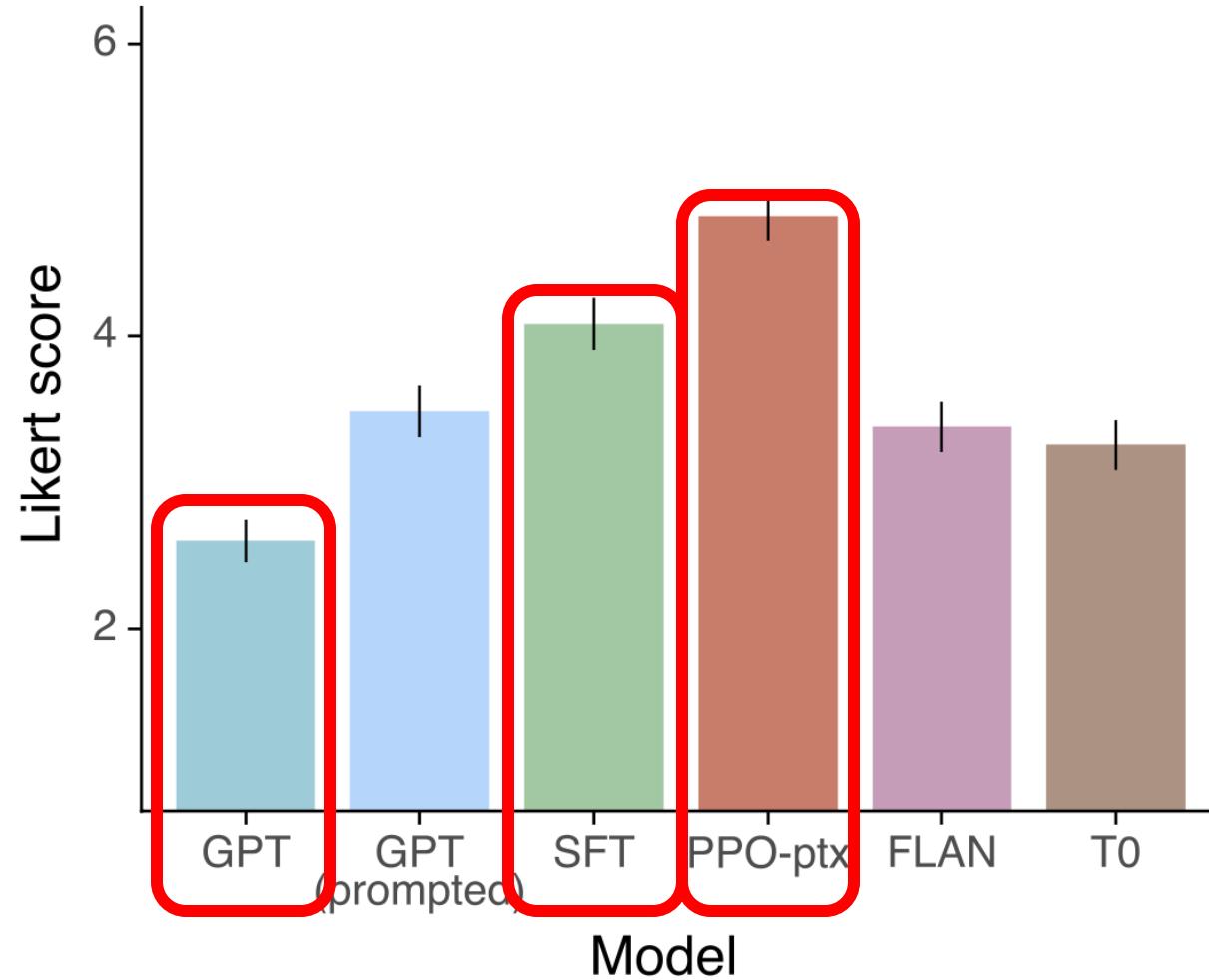
Reward Model



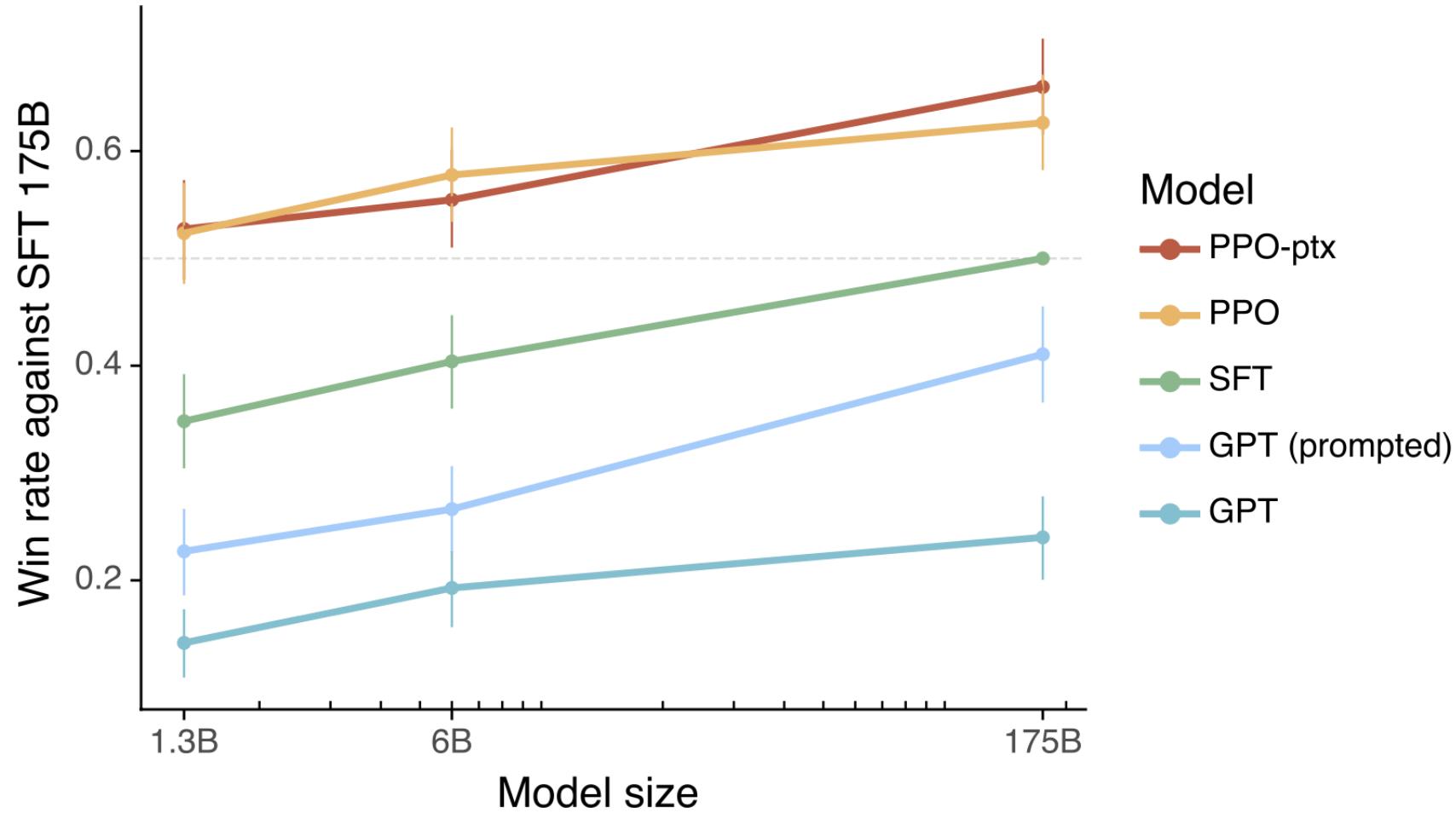
Reward Model



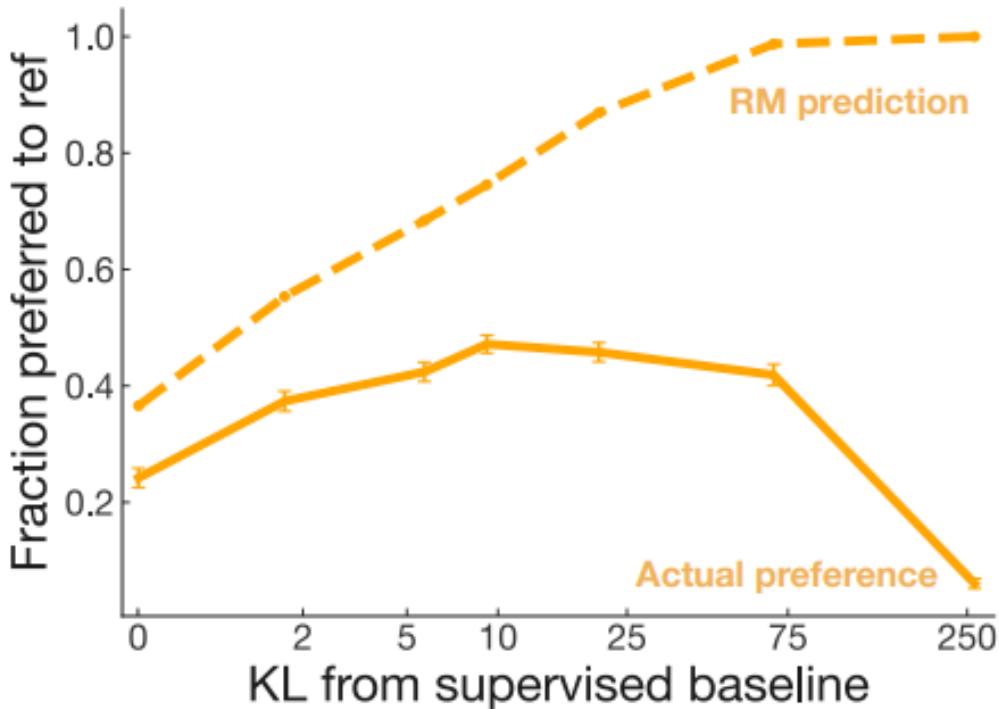
InstructGPT



InstructGPT



Overoptimization



Overoptimization in ChatGPT

- ❖ Some overoptimization symptoms we've seen:
 - ❖ Excessive verbosity (lists of lists of lists)
 - ❖ Excessive apologies, self-doubt
 - ❖ "As an AI language model"
 - ❖ Hedging language, "there's no one-size-fits-all-solution"
 - ❖ Over-refusals

Summary: LLM Training

Training Data

- Phase 1: Self-supervised Learning How are you? I am fine.
- Phase 2: Supervised Learning User: How are you? AI: I am fine.
- Phase 3: Reinforcement Learning User: How are you?
 AI: I am fine. > AI: I love you

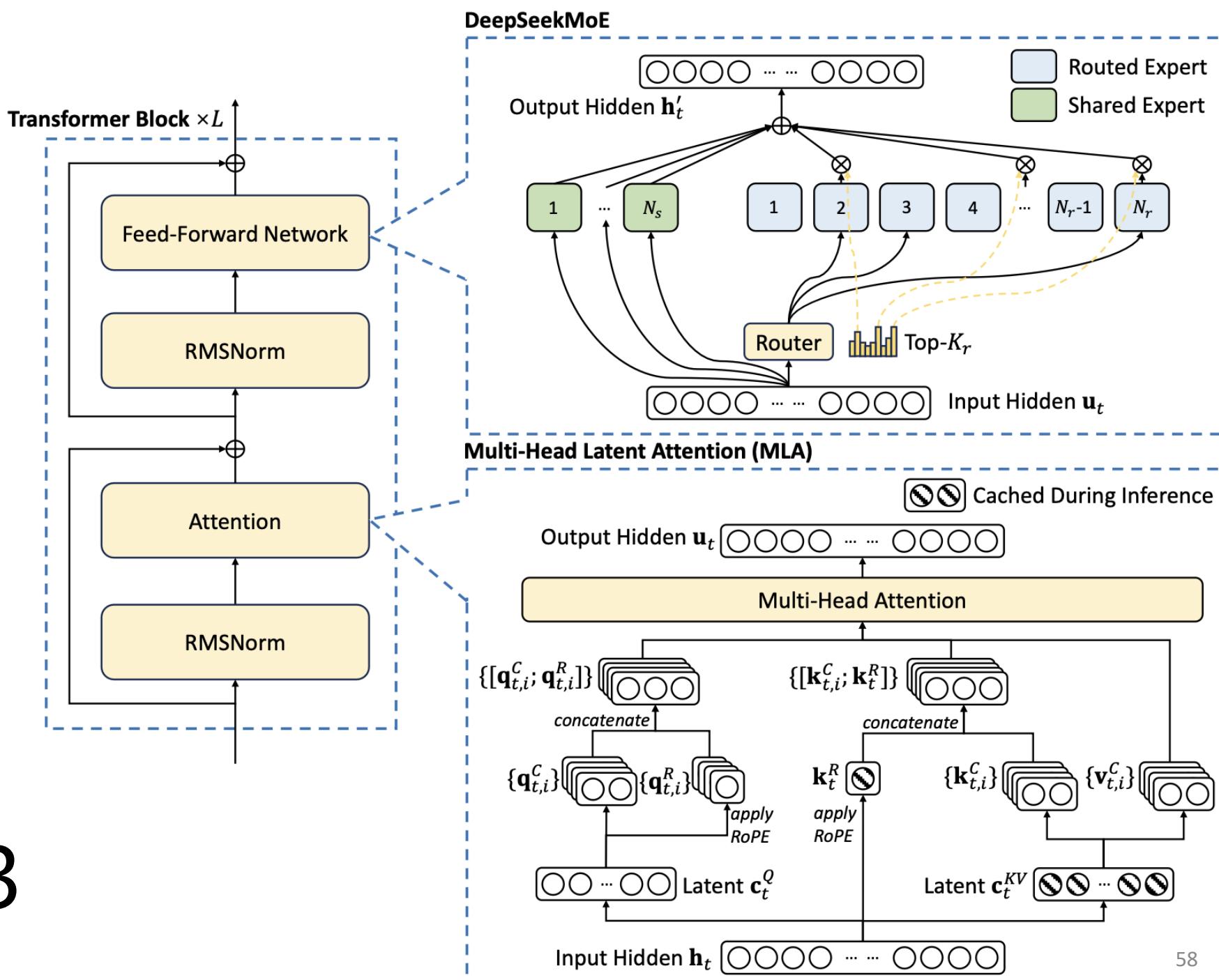
Outline

Transformer

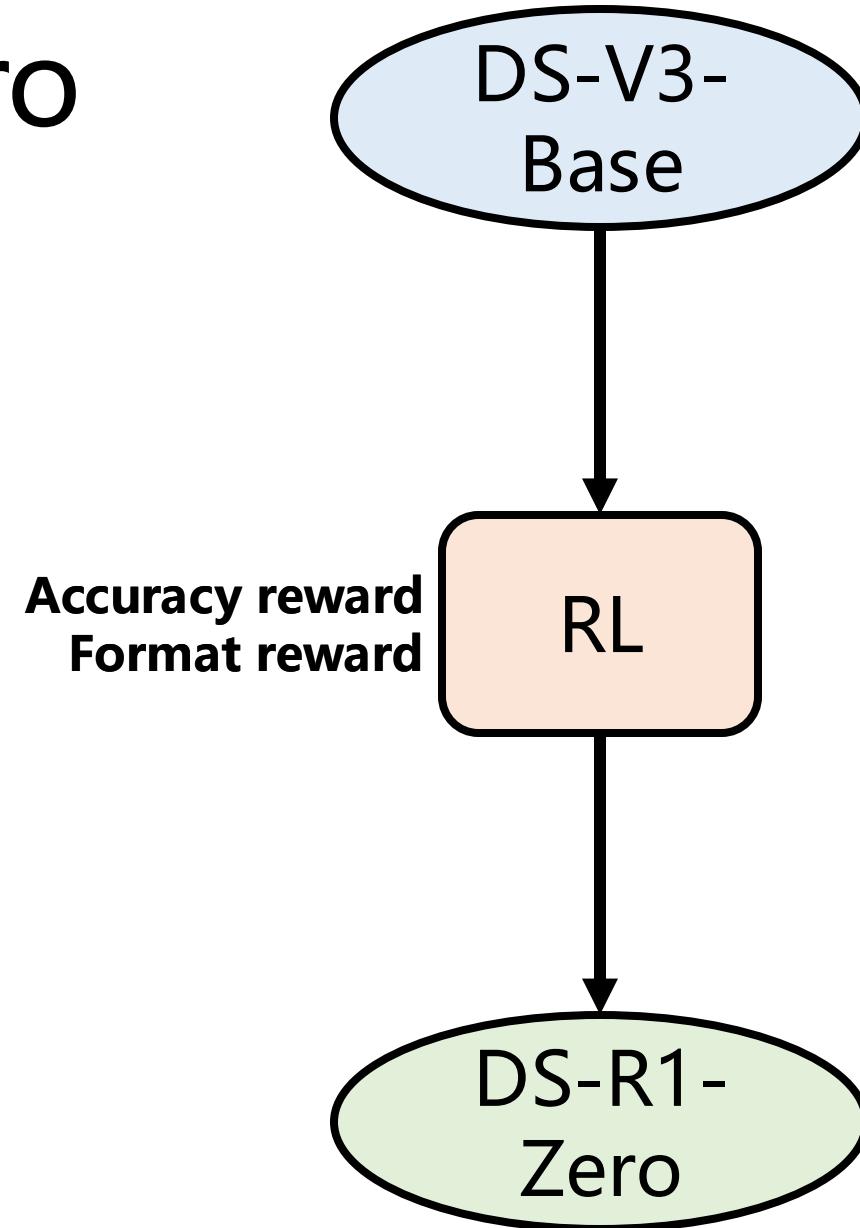
LLM Training

DeepSeek

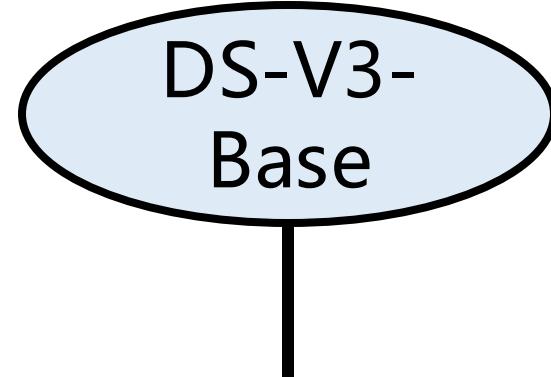
DeepSeek-V3



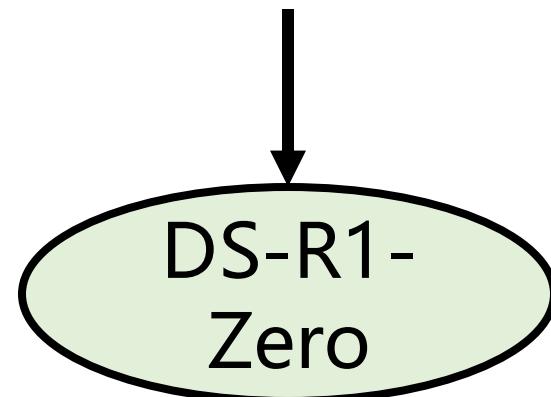
DeepSeek-R1-Zero



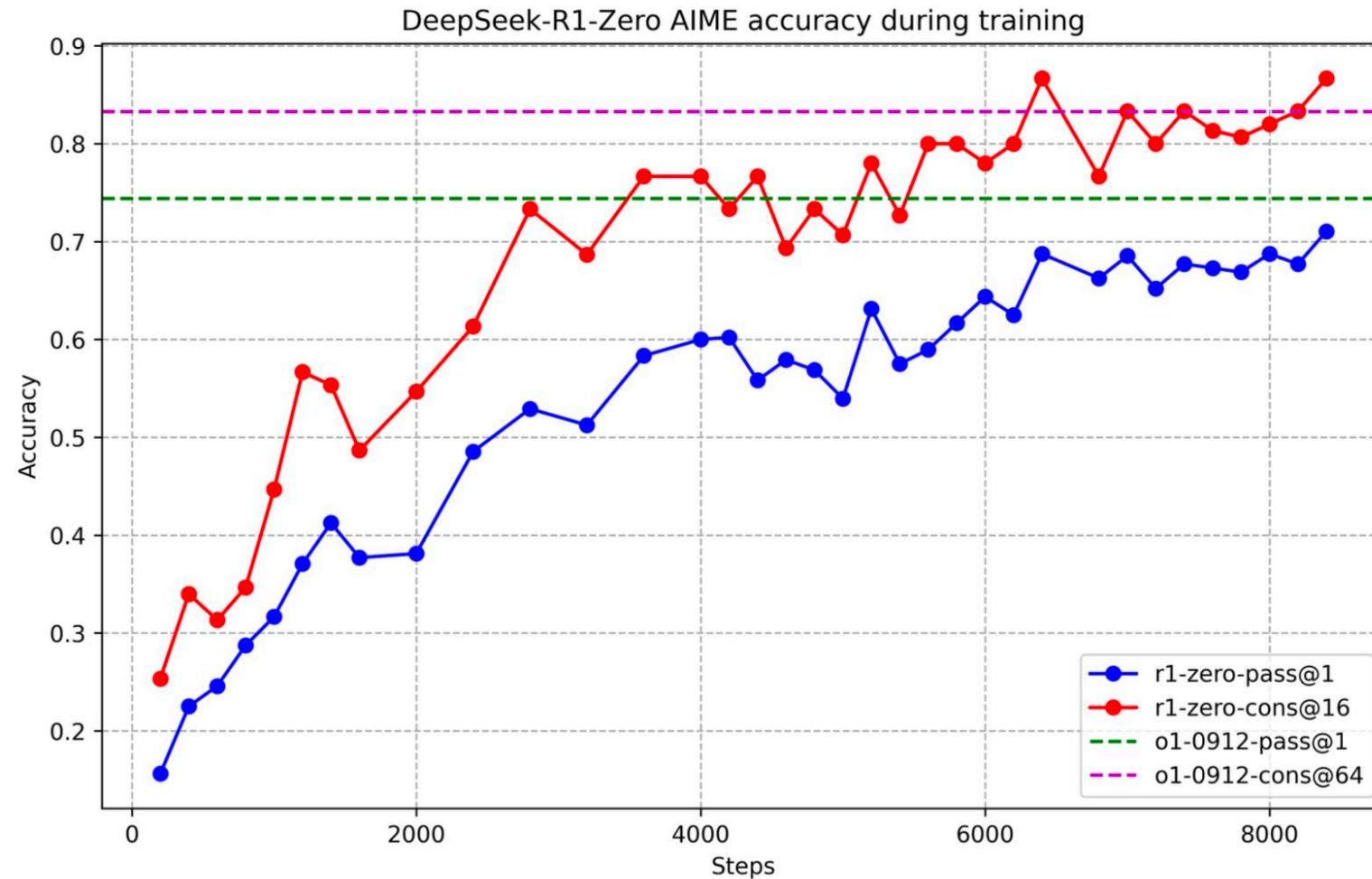
DeepSeek-R1-Zero



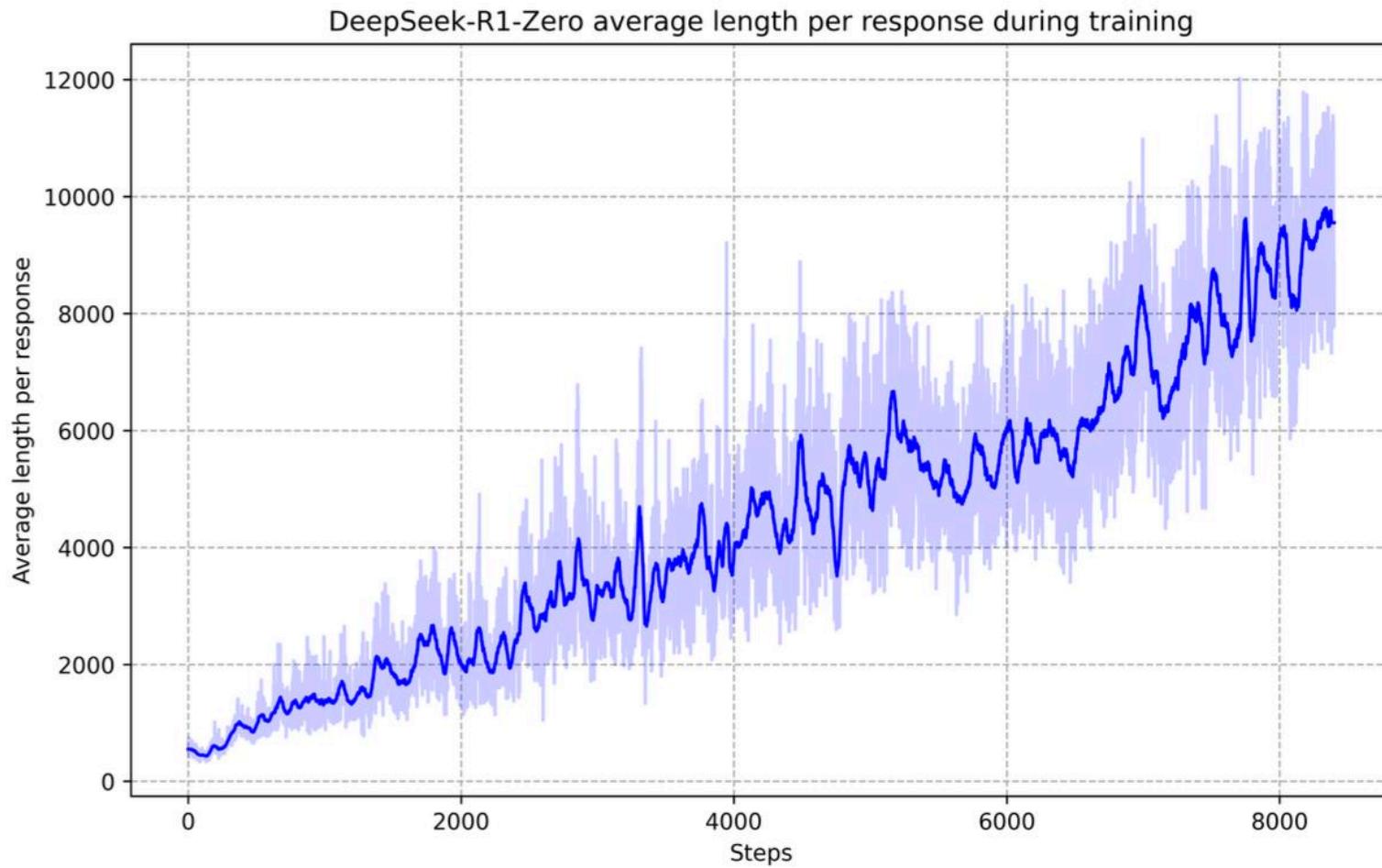
We directly apply RL to the base model without relying on supervised fine-tuning (SFT) as a preliminary step. This approach allows the model to explore chain-of-thought (CoT) for solving complex problems, resulting in the development of DeepSeek-R1-Zero. DeepSeek-R1-Zero demonstrates capabilities such as self-verification, reflection, and generating long CoTs, marking a significant milestone for the research community. Notably, it is the first open research to validate that reasoning capabilities of LLMs can be incentivized purely through RL, without the need for SFT. This breakthrough paves the way for future advancements in this area.



DeepSeek-R1-Zero



DeepSeek-R1-Zero



Aha Moment

Question: If $a > 1$, then the sum of the real solutions of $\sqrt{a - \sqrt{a+x}} = x$ is equal to

Response: <think>

To solve the equation $\sqrt{a - \sqrt{a+x}} = x$, let's start by squaring both \dots

$$(\sqrt{a - \sqrt{a+x}})^2 = x^2 \implies a - \sqrt{a+x} = x^2.$$

Rearrange to isolate the inner square root term:

$$(a - x^2)^2 = a + x \implies a^2 - 2ax^2 + (x^2)^2 = a + x \implies x^4 - 2ax^2 - x + (a^2 - a) = 0$$

\dots

Wait, wait. Wait. That's an aha moment I can flag here.

Let's reevaluate this step-by-step to identify if the correct sum can be \dots

We started with the equation:

$$\sqrt{a - \sqrt{a+x}} = x$$

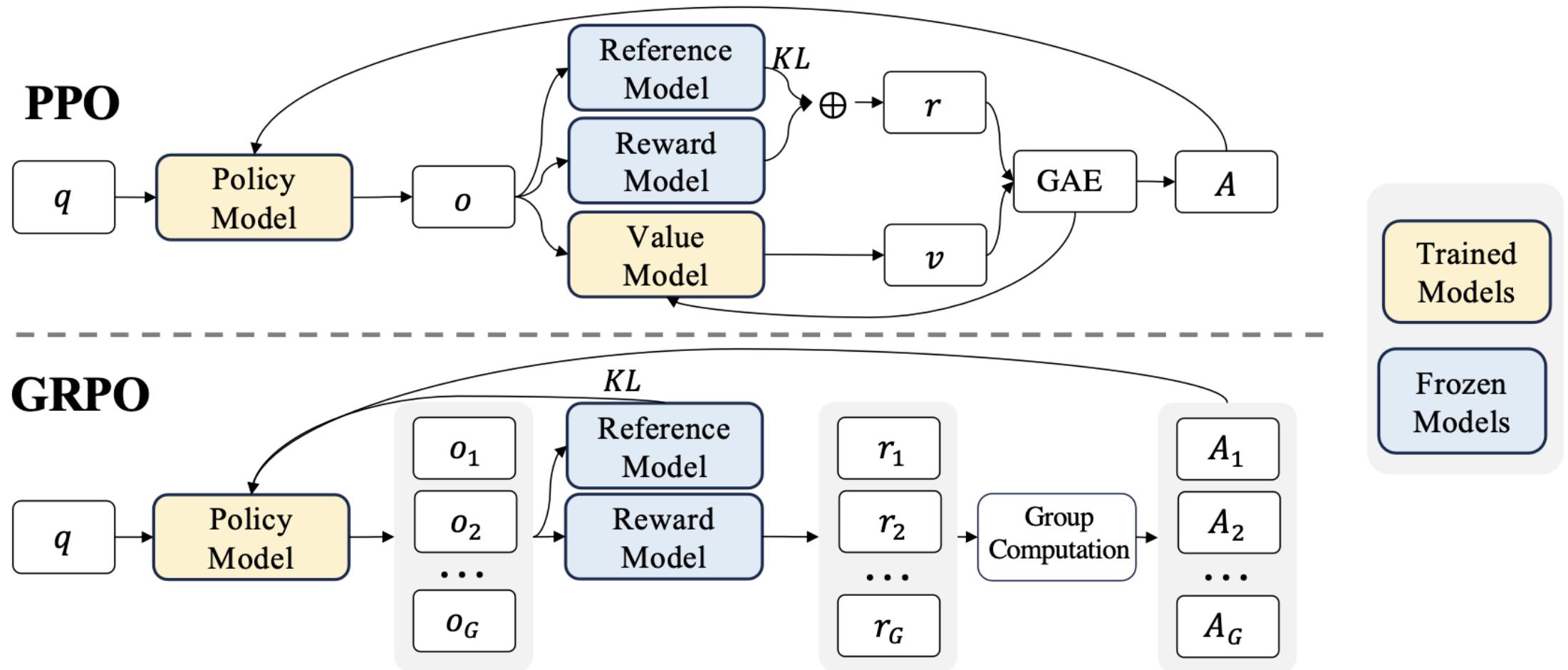
First, let's square both sides:

$$a - \sqrt{a+x} = x^2 \implies \sqrt{a+x} = a - x^2$$

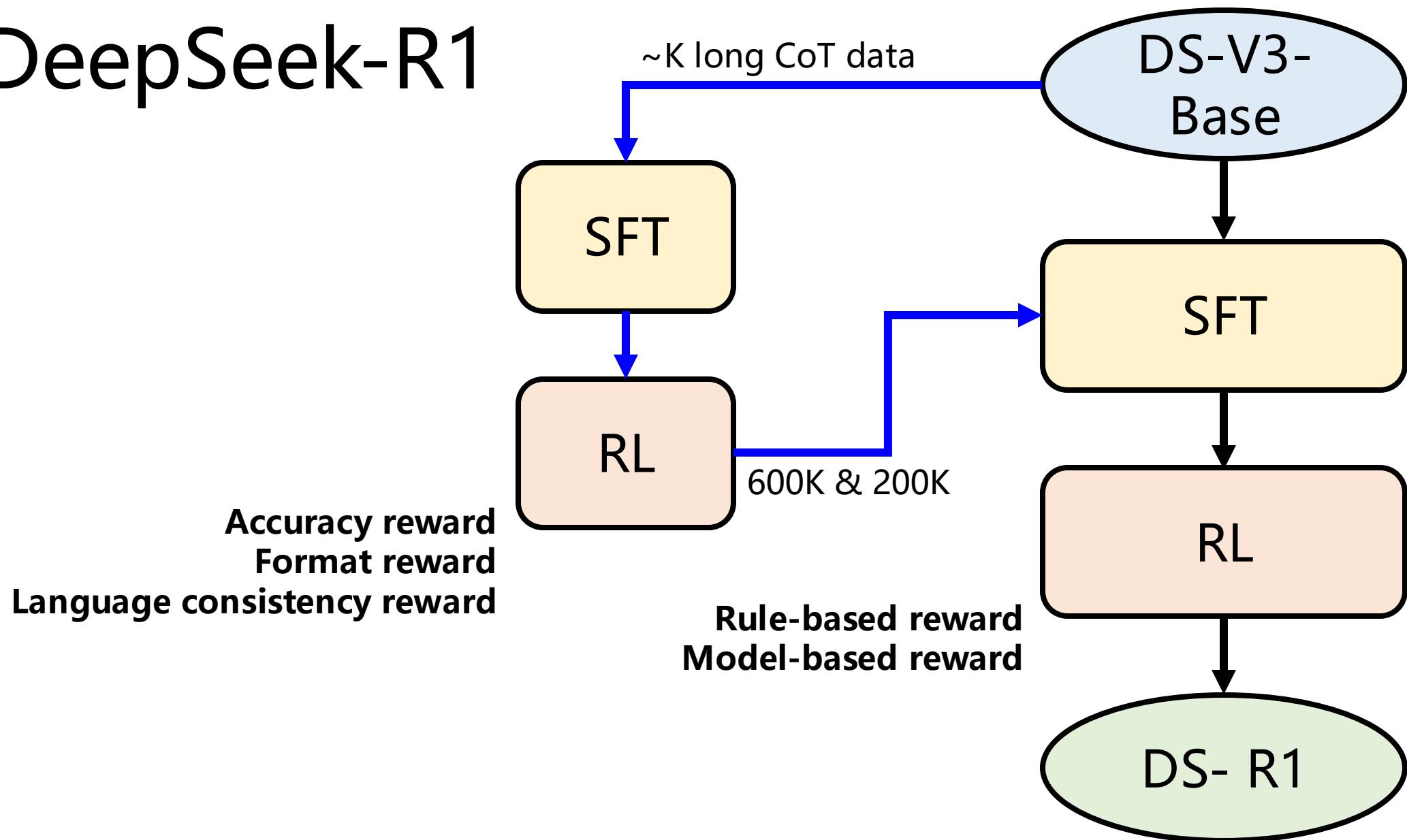
Next, I could square both sides again, treating the equation: \dots

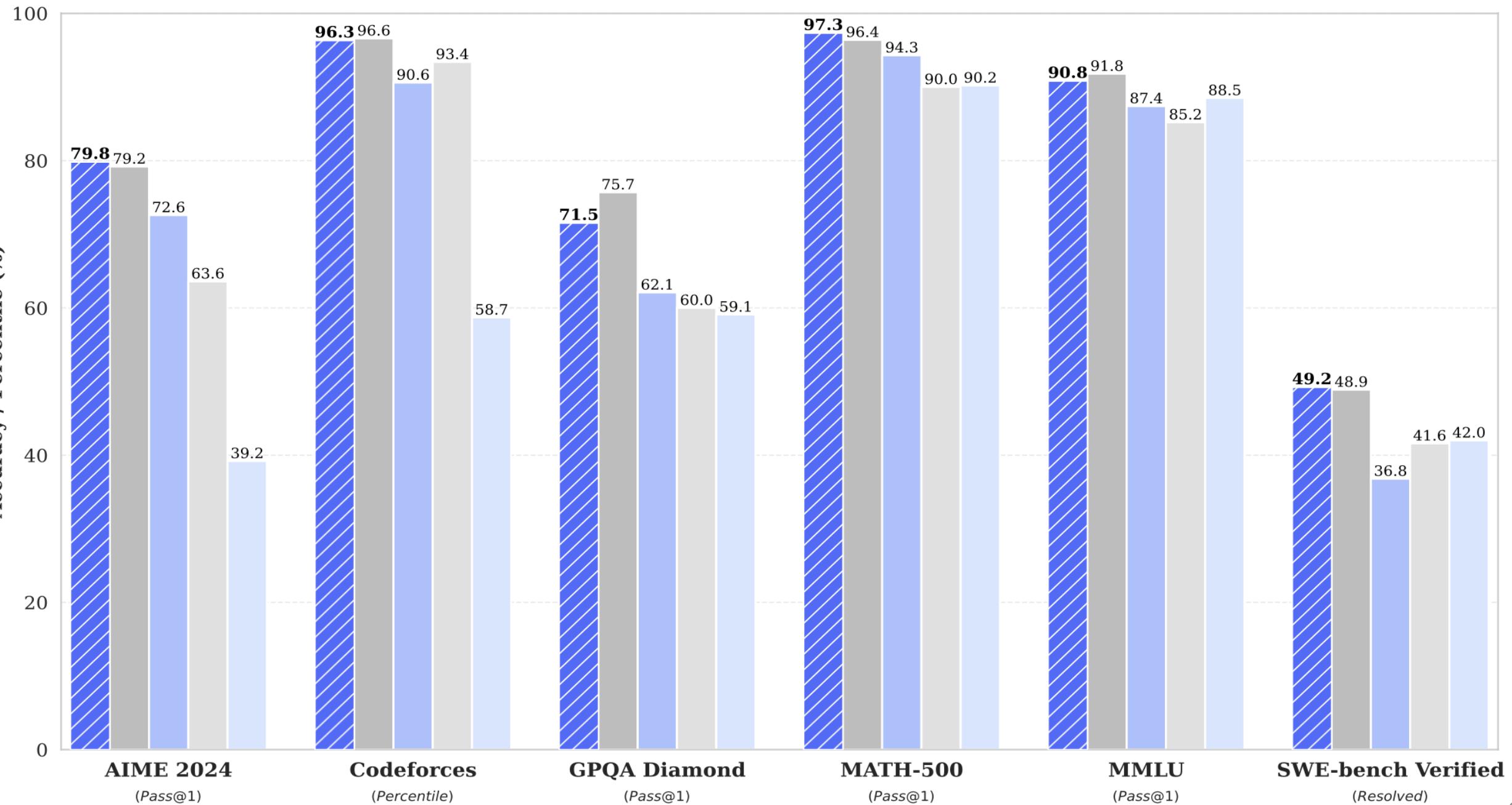
\dots

Group Relative Policy Optimization



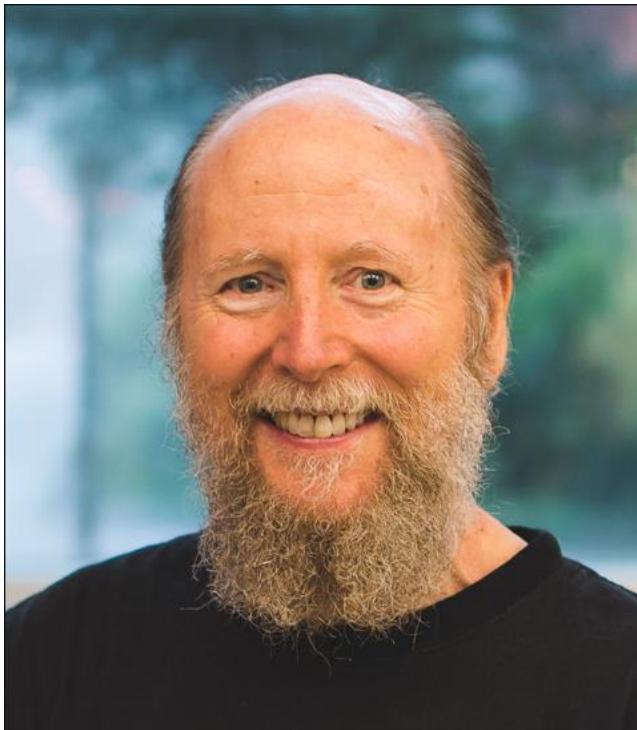
DeepSeek-R1





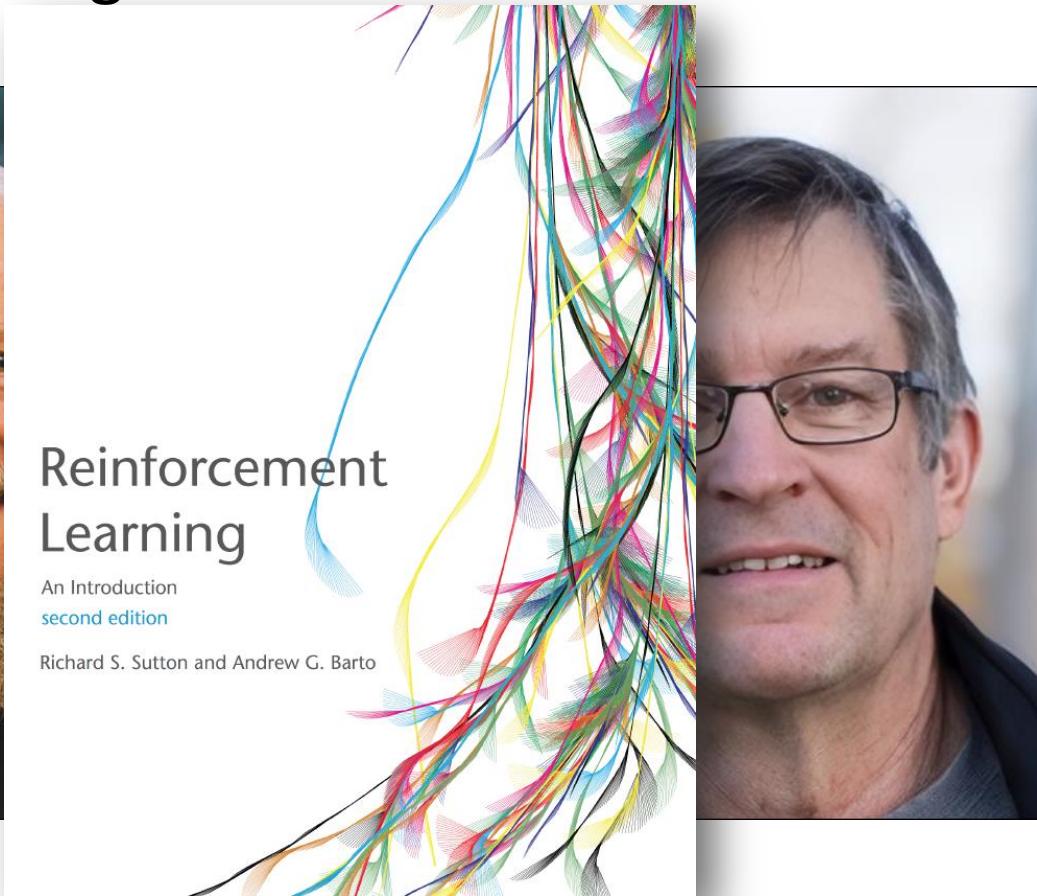
2024 Turing Award

- Andrew Barto and Richard Sutton Recognized as Pioneers of Reinforcement Learning



2024 Turing Award

- Andrew Barto and Richard Sutton Recognized as Pioneers of Reinforcement Learning



THANK YOU