

Large Language Models

demystified

Prachya Boonkwan (P'/N' Arm)

Language and Semantic Technology Lab (LST)

National Electronics and Computer Technology Center (NECTEC)

prachya.boonkwan@nectec.or.th, kaamanita@gmail.com

URL ⇒ <https://tinyurl.com/2pvf3kcr>



Who? Me?

- Nickname: **Arm** (P'/N' Arm, etc.)
- Born: Aug 1981
- Work: researcher at NECTEC since 2005
- Education
 - B.Eng & M.Eng, CPE Kasetsart University
 - Obtained Ministry of Science Scholarship in early 2008
 - Did a PhD in Informatics (AI & Computational Linguistics) at University of Edinburgh, UK



Large Language Models Demystified

- Introduction to AI
- Large language models
- Transformer Model
- Training ChatGPT
- OpenThaiGPT
- Future direction

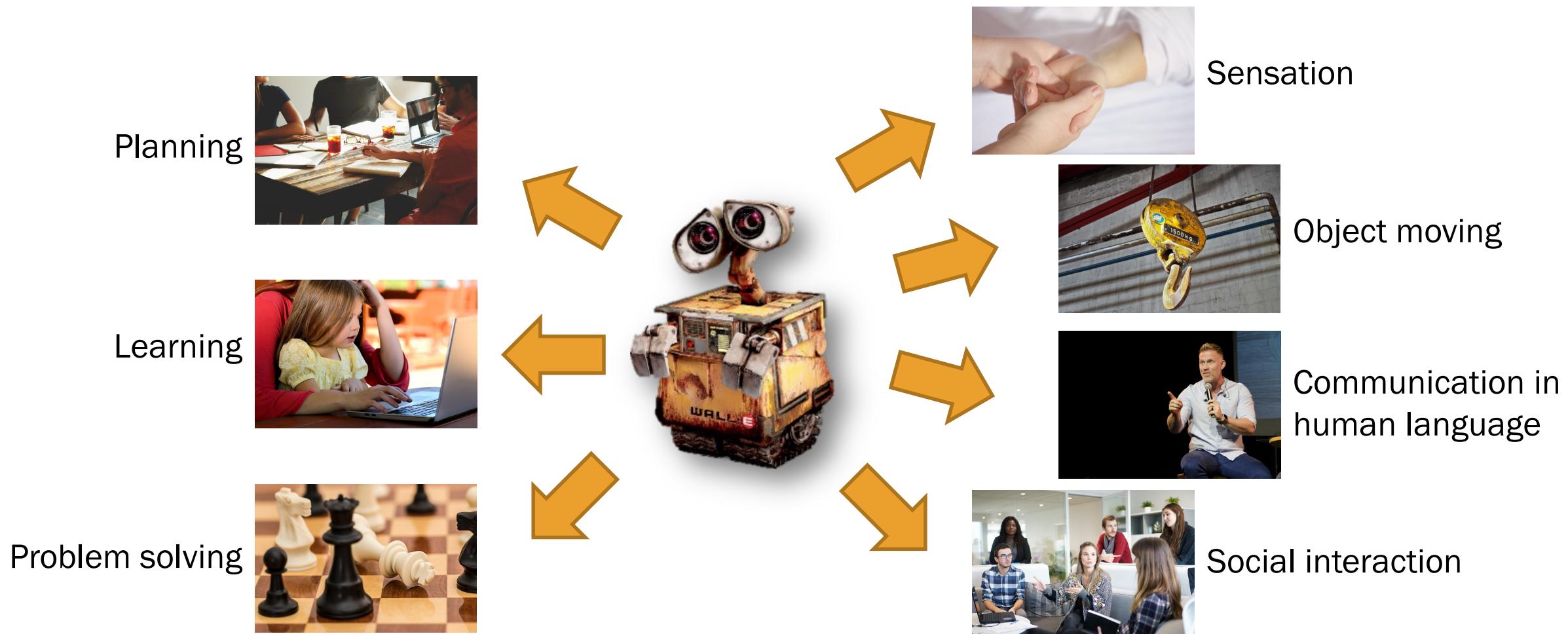


<https://tinyurl.com/2pvf3kcr>

1. What exactly is AI?

What is Artificial Intelligence?

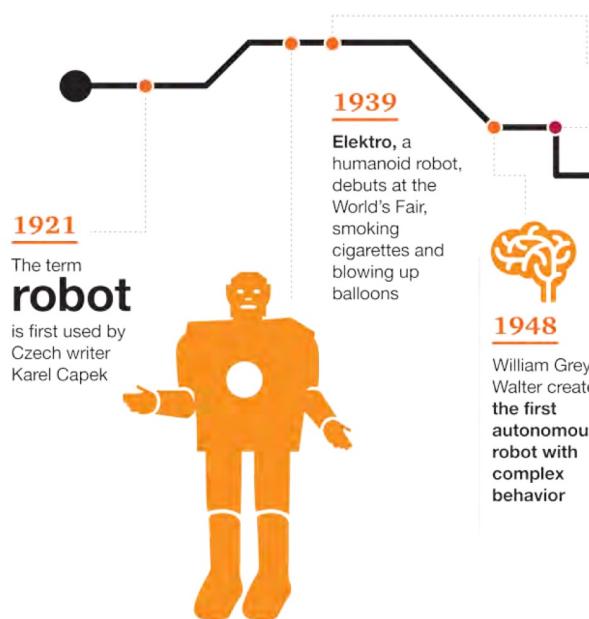
- Imitation of natural intelligence with machine





The rise of Robotics and AI

Fueled by advances in computing power and connectivity, the fields of robotics and artificial intelligence have grown rapidly



1941 Isaac Asimov formulates the **Three Laws of Robotics:**

- A robot may not injure a human being or, through inaction, allow a human being to be harmed.
- A robot must obey orders given it by human beings except where such orders would conflict with the First Law.
- A robot must protect its own existence as long as such protection does not conflict with the First or Second Law.

1956

Field of AI research founded at a conference at Dartmouth

1960

Frank Rosenblatt constructs **Mark I Perceptron**, a computer that learned new skills by trial and error

1972

Stanford researcher develops **PARRY**, designed to simulate a paranoid schizophrenic.

1974

Intel produces its second-generation 8080 general-purpose chips

1968

Mobile robot "Shakey" is introduced. It's controlled by a computer the size of a room

1979

SCARA, an articulated robot arm, is developed for assembly lines

1984

Doug Lenat and his team start **Cyc**, to codify millions of pieces of knowledge that compose human common sense

1985

Jaron Lanier's VPL Research, Inc., sells first VR glasses and gloves; Lanier coins the phrase

1986

Honda creates the E0, the first of a series of humanoid robots that walk on two feet

1988

The first **HelpMate** service robot begins work at Danbury Hospital

Turing's Test.

It tests a machine's ability to "think" by answering a series of questions. In essence, the tester must think the machine's answers are coming from a human

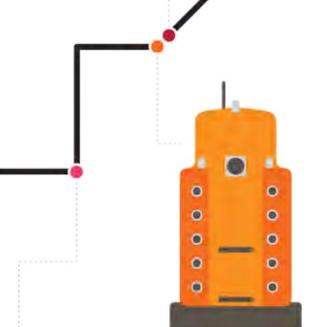
Minimize and maximize

Shrinking disk sizes and exponentially growing capacity help fuel robotics and AI

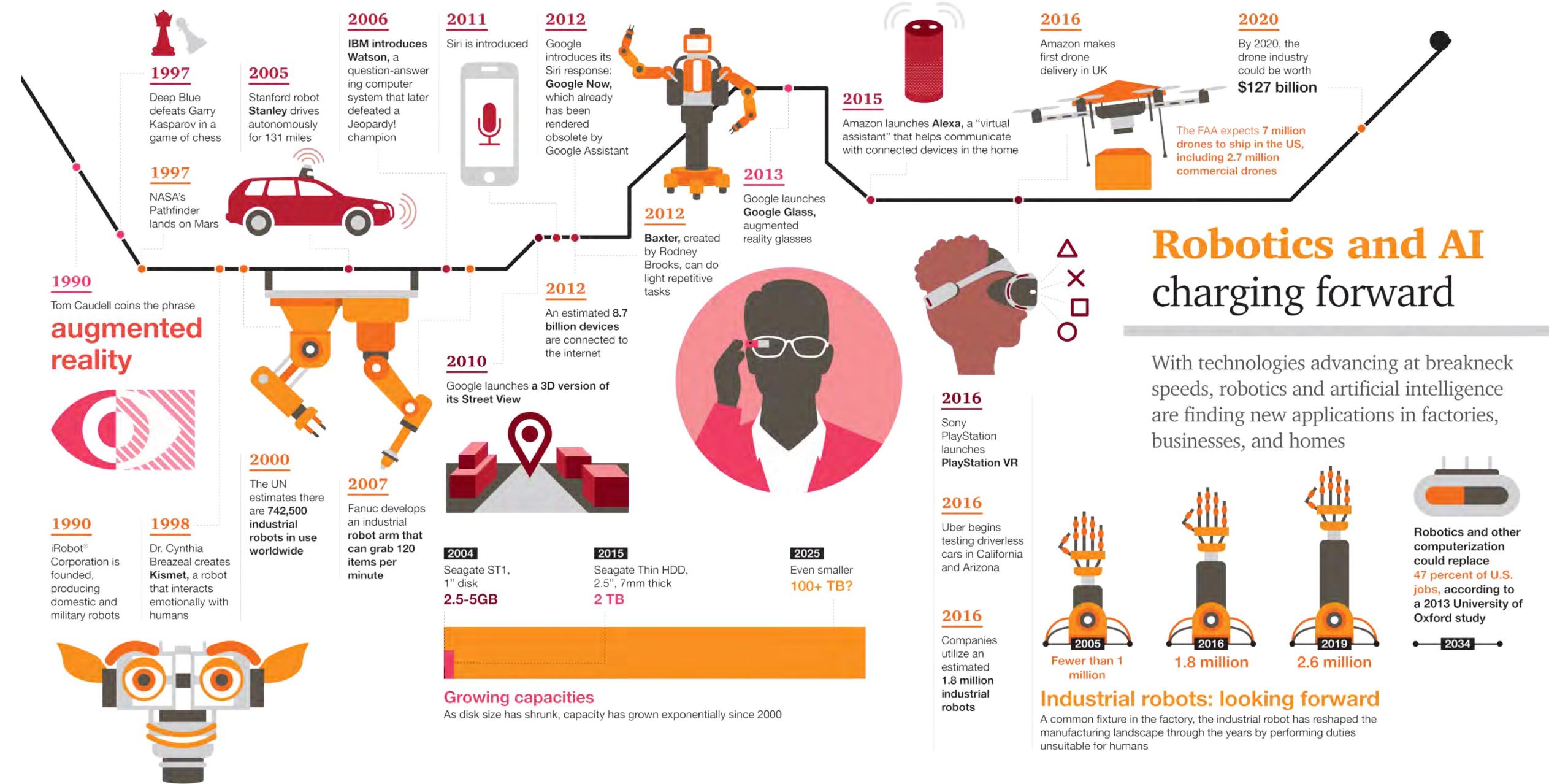


1988 **Nope, I'm human.**

Researchers launch **Jabberwocky**, an AI chatbot designed to learn through conversation

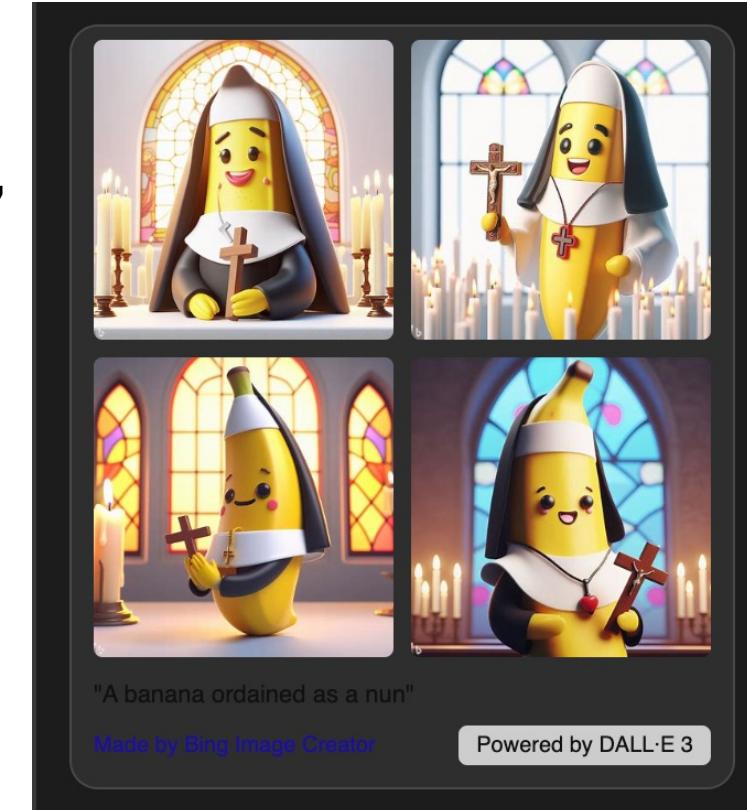


virtual reality



Generative AI

- Algorithms (such as ChatGPT) that can be used to create new content, including audio, code, images, text, simulations, and videos
- Popular generative AI models
 - ChatGPT for text (OpenAI)
 - DALL-E and Stable Diffusion (MidJourney) for image
 - VALL-E for voice (Microsoft)
 - Deepfake for video



นี่คือคลิปวิดีโอ พระพยอม กัลยาณ
ถ่ายไว้เมื่อวันที่ 28 สิงหาคม 2562

https://www.youtube.com/watch?v=wWjnXa_k2NE

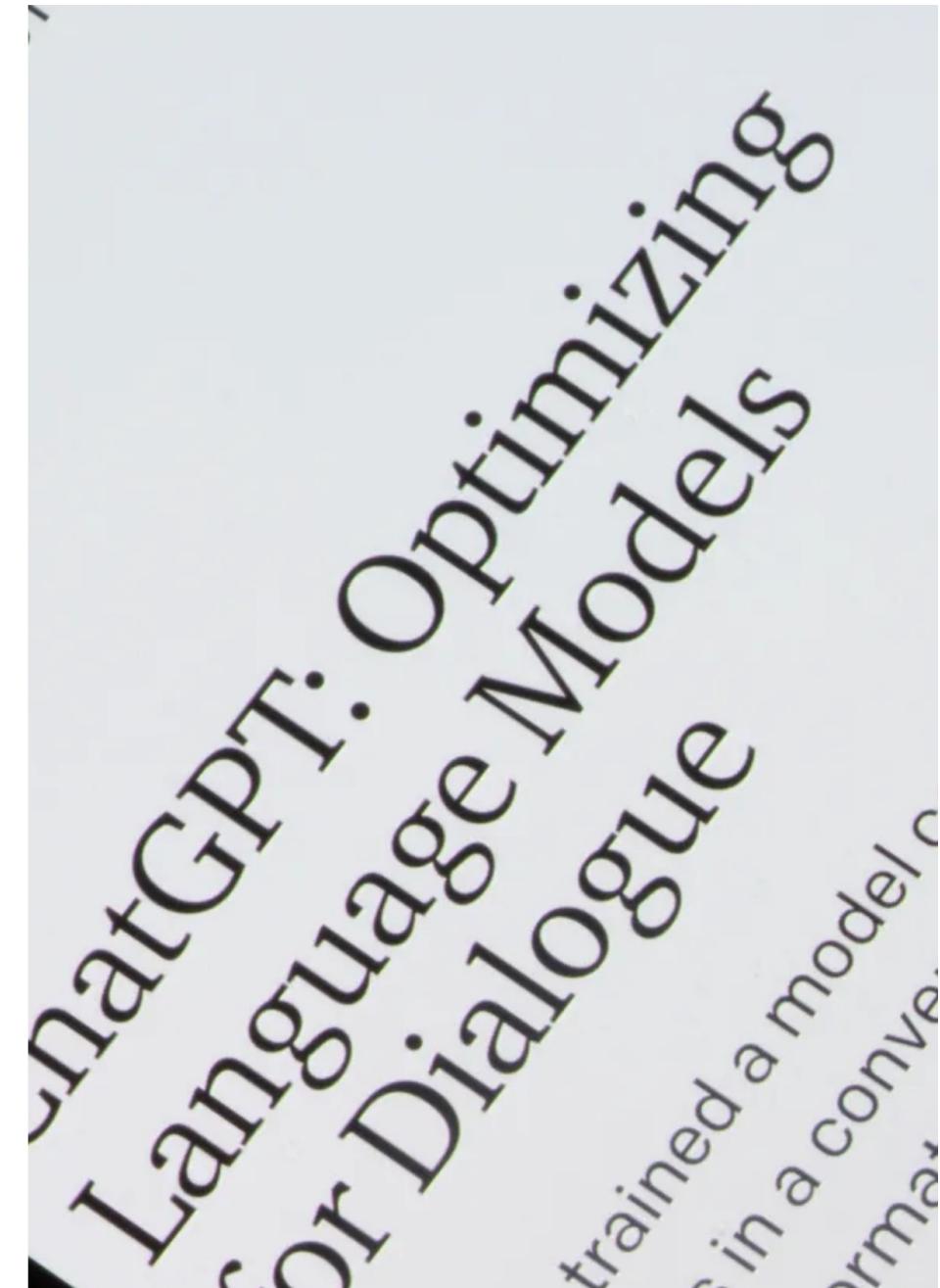
2. Large Language Models

ChatGPT

ChatGPT is an advanced AI chatbot trained by [OpenAI](#) which interacts in a conversational way. The dialogue format makes it possible for ChatGPT to answer followup questions, admit its mistakes, challenge incorrect premises, and reject inappropriate requests.

ChatGPT relies on the powerful [GPT-3.5](#) technology. GPT stands for Generative Pre-Trained Transformer, a complex neural network based on the revolutionary [Attention concept](#).

<https://chat.openai.net>

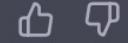




I'm having a hair loss problem. Can taking a daily zinc pill solve the problem?



Zinc is an essential mineral that is important for a number of processes in the body, including hair growth. Zinc deficiency can lead to hair loss, and supplementing with zinc may help to promote hair growth in some people. However, hair loss can also be caused by a variety of other factors, such as genetics, hormonal imbalances, medical conditions, and certain medications.



If you are experiencing hair loss, it's best to consult a doctor or a dermatologist for an accurate diagnosis and appropriate treatment recommendations. They may recommend taking a zinc supplement in addition to other treatments, such as topical products, lifestyle changes, or medications. However, it's important to keep in mind that taking too much zinc can also have negative side effects, so it's important to take only the recommended amount.

In conclusion, while zinc supplementation may help with hair loss in some cases, it's not a cure-all solution and it's important to work with a healthcare professional to determine the best course of action for your specific situation.



ChatGPT (OpenAI)



Copilot (Microsoft)



perplexity

Perplexity.AI



Gemini (Google)

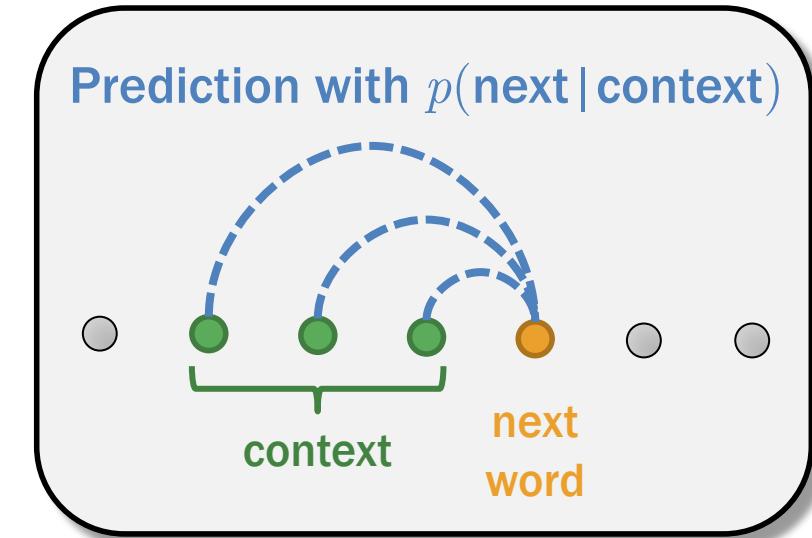


Claude.AI (Anthropic)

Large Language Model (LLM)

- Statistical prediction for how strings are produced in a language
- Interpreted as a generative model
 1. Generate the first word w_1
 2. Keep generating the **next word** w_k based on the previous words (a.k.a. **context**) $w_1 \dots w_{k-1}$ until the whole sentence of length N is produced

$$P(w_1 \dots w_N) = p(w_1) \prod_{k=2}^N p(\text{next word} | \underbrace{w_1 \dots w_{k-1}}_{\text{context}})$$



We need at least 1 billion words of text data to train a stable LLM, which learns **word collocations** and **phrase structures**

What is GPT?



อรอนา สิงห์ศรี

บ้อง เป็น สาว ----- --- --- --- -- ----- ----- ----- ----- ----- ----- ----- -----

(This is a very popular Thai country song in the 1980s)

What is GPT?



อรอนา สิงห์ศรี

บ้อง เป็น สาว ขอนแก่น --- --- --- --- --- --- --- --- --- --- ---

(This is a very popular Thai country song in the 1980s)

What is GPT?



อรอนา สิงห์ศรี

บ้อง เป็น สาว ขอนแก่น ยัง --- --- -- ----- ----- ----- ----- ----- ----- ----- -----

(This is a very popular Thai country song in the 1980s)

What is GPT?



อรอนา สิงห์ศรี

บ้าน เป็น สาว บ้านแก่น ยัง บ ---

(This is a very popular Thai country song in the 1980s)

What is GPT?



อรอนา สิงห์ศรี

บ้อง เป็น สาว ขอนแก่น ยัง บ่ เคย -- ----- ----- ----- ----- ----- ----- -----

(This is a very popular Thai country song in the 1980s)

GPT: Generative Pretrained Transformer

เพลง: สาวอีสานรอรัก



อรอนุมา สิงห์ศิริ

น้อง เป็น สาว ขอนแก่น ยัง บ่ เคย มี ແພນ บ้าน อยู่ ແດນ อีสาน

<https://www.youtube.com/watch?v=stLQVIau1ns>

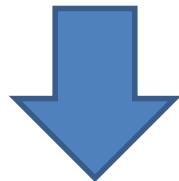
Transformer Model (Vaswani et al., 2016)

- Sequence-to-sequence generation

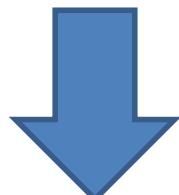
- Translation:** It learns how to produce a target sequence from a source sequence, given a very large dataset of sequence pairs
- Pros:** It learns **word collocations** and **phrase structures** on the input and output sequences, and associates them cross-lingually in the table of **translation alignments**
- Cons:** It consists of an expansive amount of neuron cells, and the training process can be quite time-consuming

Who | is | the | current | president | of | the | US

Source: sequence of words (prompt)



TRANSFORMER



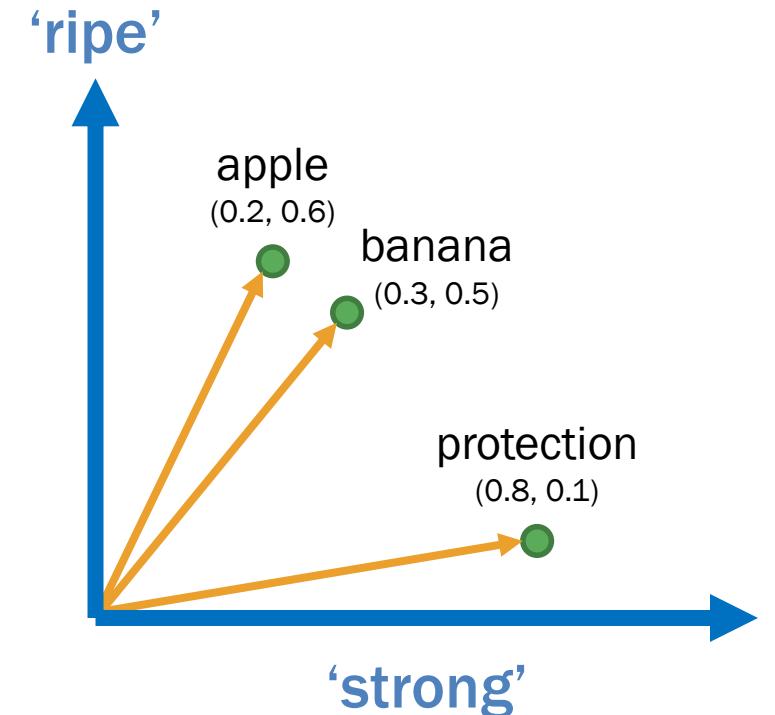
The | president | of | the | US | is | Joe | Biden

Target: sequence of words (response)

3. Transformer Model

Recap: Word Vector

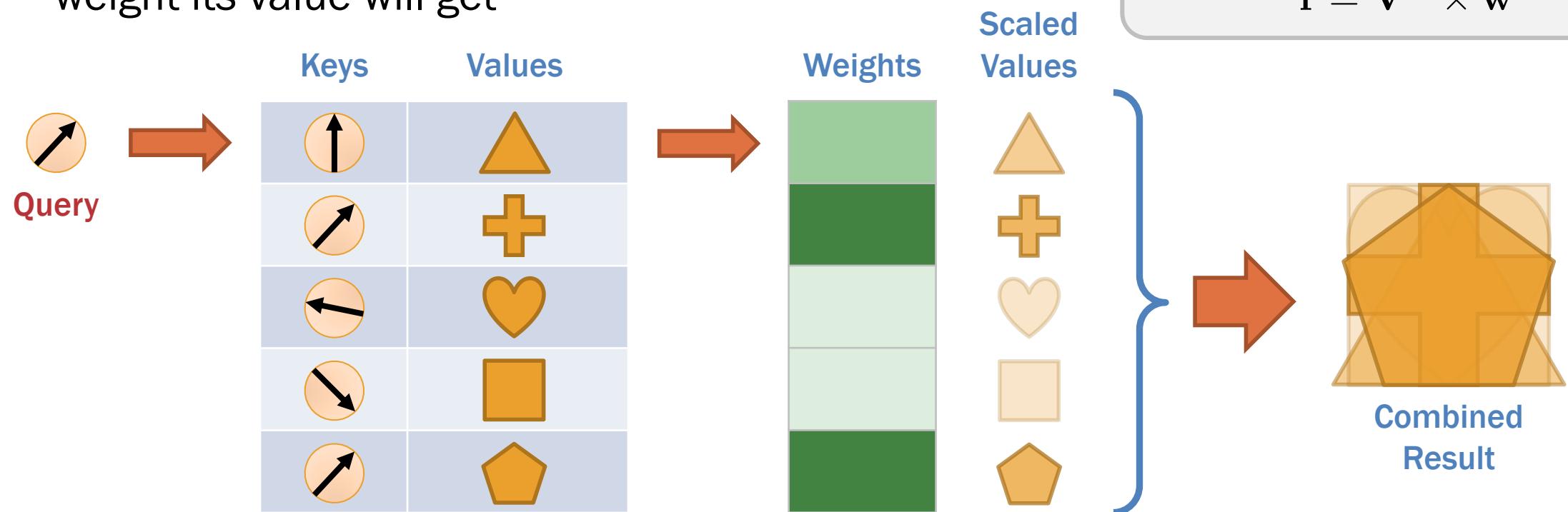
- Distributional semantics
 - Measured by co-occurrence of words and their contexts (i.e. **context distribution**)
 - Various types of similarity metrics and context are employed (Dagan+, 2008)
 - Backoff method (Katz, 1987) with interpolation (Jelinek+, 1980) are required for smoothing the zero counts
 - Bigram co-occurrence is commonly used as the context representation (Brown+, 1992)
 - Well-known techniques include Latent Semantic Analysis (Dumais, 2005) and PMI (Church+, 1990)



In this over-simplified example,
each numeric element is:
 $p(\text{word} | \text{context})$

Scaled Dot-Product Attention

- Semantic similarity \Rightarrow search engine
 - Query is compared against each key with dot product
 - The more similar the key is to the query, the more weight its value will get



$$w_i \propto k_i \cdot q$$

Simple
Form

$$\mathbf{r} = \sum_{i=1}^N w_i \mathbf{v}_i$$

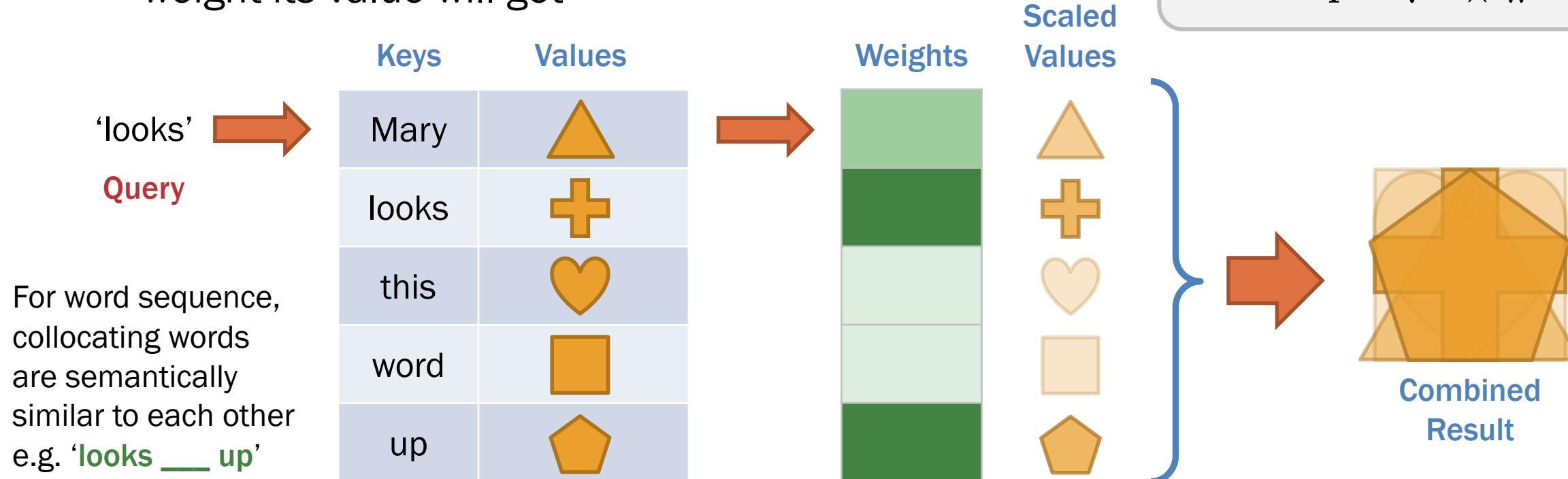
Matrix
Form

$$\mathbf{w} = \text{Softmax}(\mathbf{K} \times \mathbf{q})$$

$$\mathbf{r} = \mathbf{V}^\top \times \mathbf{w}$$

Scaled Dot-Product Attention

- Semantic similarity \Rightarrow search engine
 - Query is compared against each key with dot product
 - The more similar the key is to the query, the more weight its value will get



$$w_i \propto k_i \cdot q$$

Simple
Form

$$\mathbf{r} = \sum_{i=1}^N w_i \mathbf{v}_i$$

Matrix
Form

$$\mathbf{w} = \text{Softmax}(\mathbf{K} \times \mathbf{q})$$

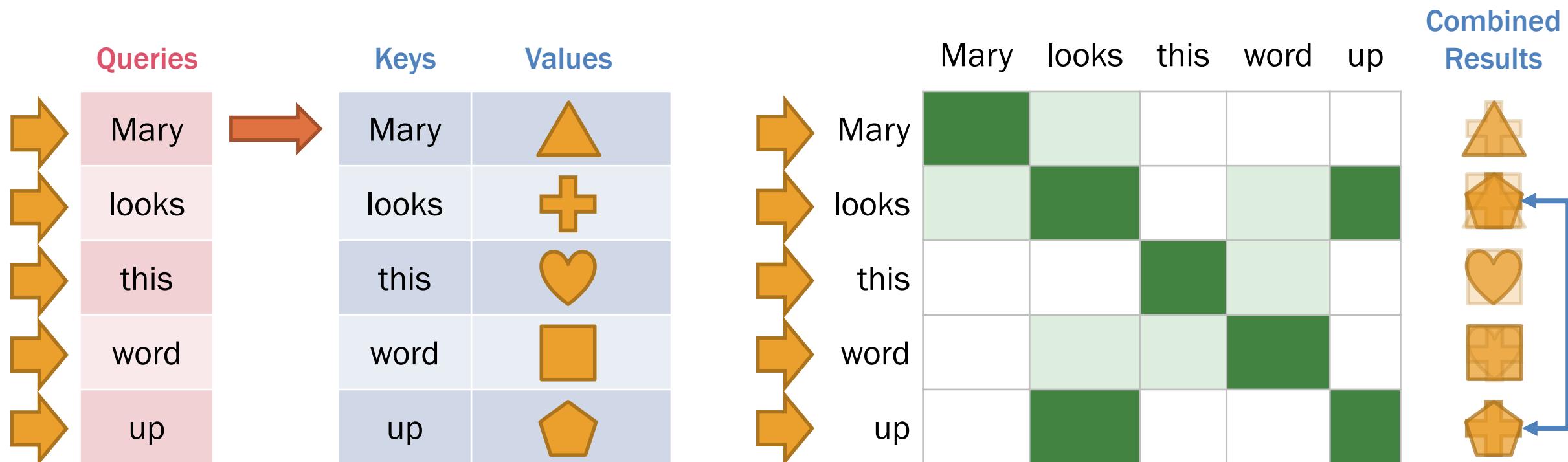
$$\mathbf{r} = \mathbf{V}^\top \times \mathbf{w}$$

Combined
Result

Self-Attention

- Scaled dot-product attention whose queries and keys are the same
- Collocations will have almost similar results

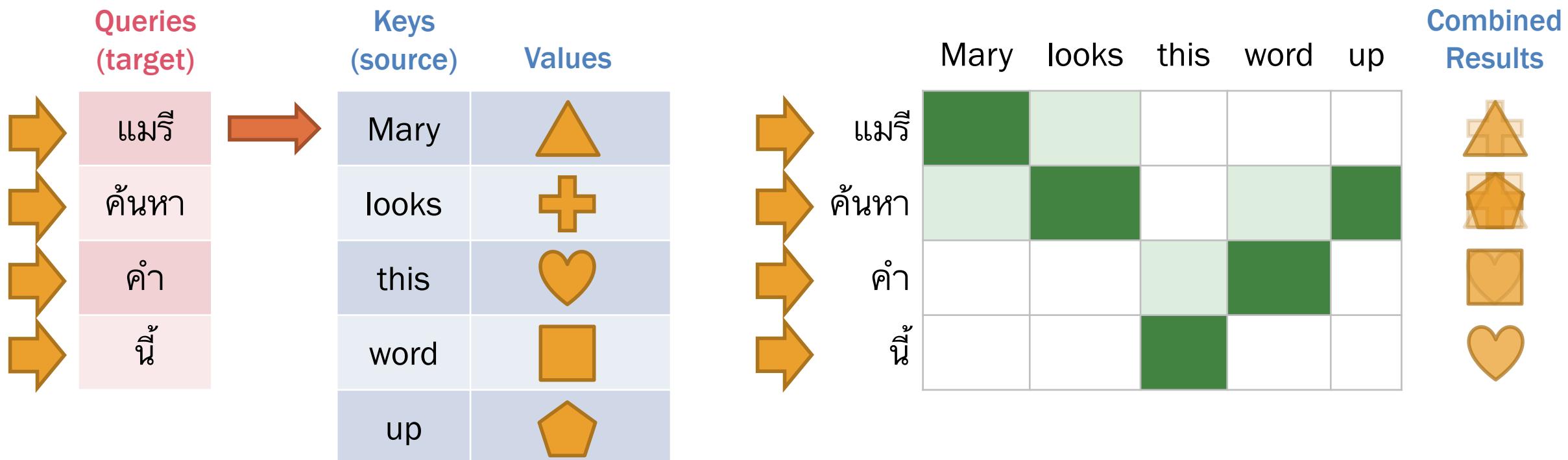
$$\begin{aligned} \text{Matrix Form } \mathbf{W} &= \text{Softmax}(\mathbf{K} \times \mathbf{K}^T) \\ \text{Form } \mathbf{R} &= \mathbf{W} \times \mathbf{V} \end{aligned}$$



Cross-Attention

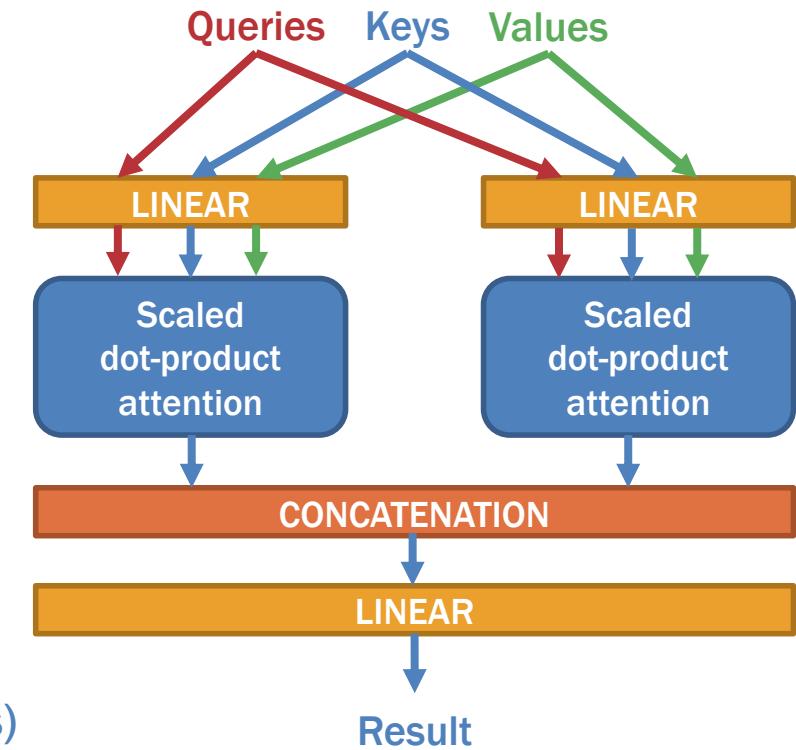
- Scaled dot-product attention whose queries are the target and whose keys are the source
- Collocation alignment via semantic similarity

$$\begin{aligned} \text{Matrix Form } \mathbf{W} &= \text{Softmax}(\mathbf{Q} \times \mathbf{K}^\top) \\ \mathbf{R} &= \mathbf{W} \times \mathbf{V} \end{aligned}$$



Multihead Attention

- Scaled dot-product attention has a drawback
 - It recognizes **only one** type of word collocation
 - If we assume more than one type of word collocation per sequence, then we have to combine multiple attention heads [default = 8 heads]



HEAD 1 (looks __ up)

Mary	Poppins	looks	this	word	up
Mary					
Poppins					
looks					
this					
word					
up					

HEAD 2 (Mary Poppins)

Mary	Poppins	looks	this	word	up
Mary					
Poppins					
looks					
this					
word					
up					

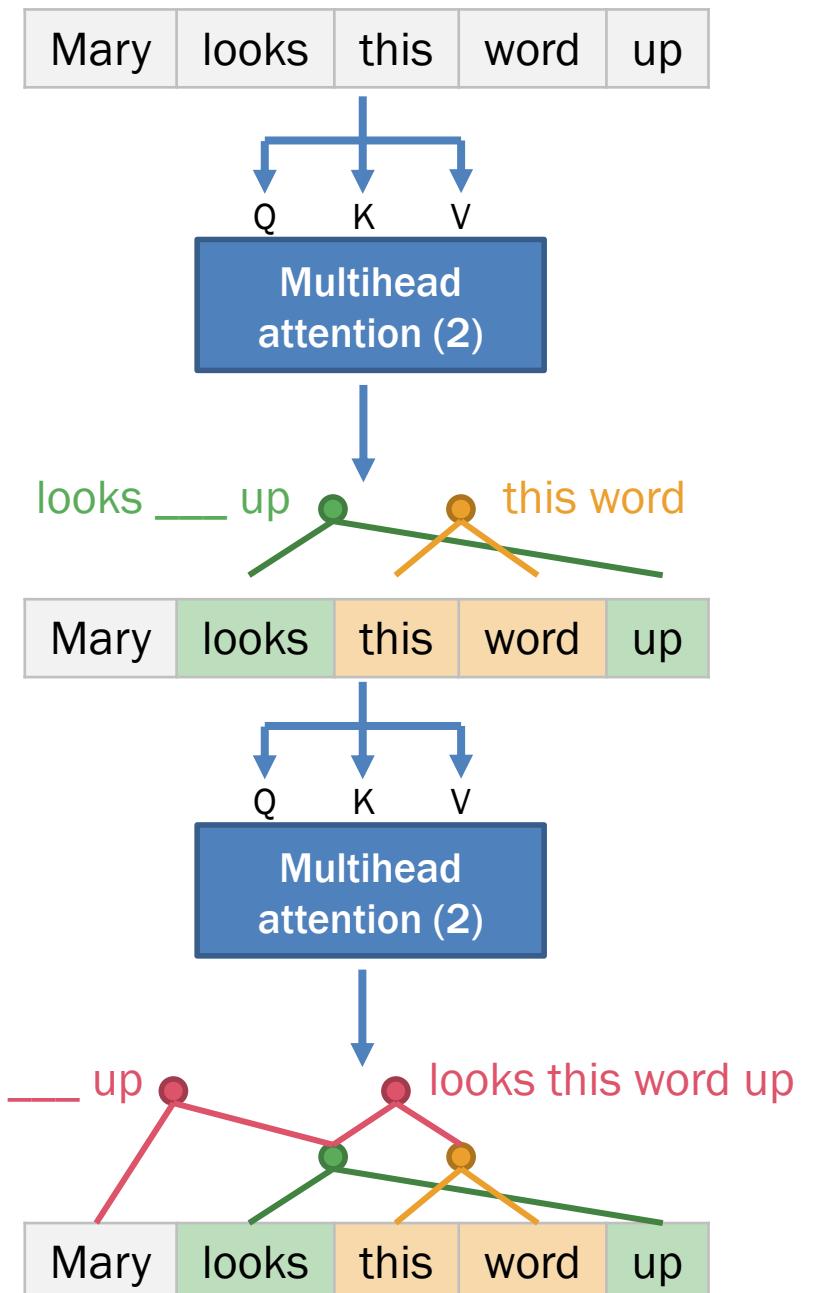
Notation

↓ ↓ ↓
Q K V

**Multihead
attention (n)**

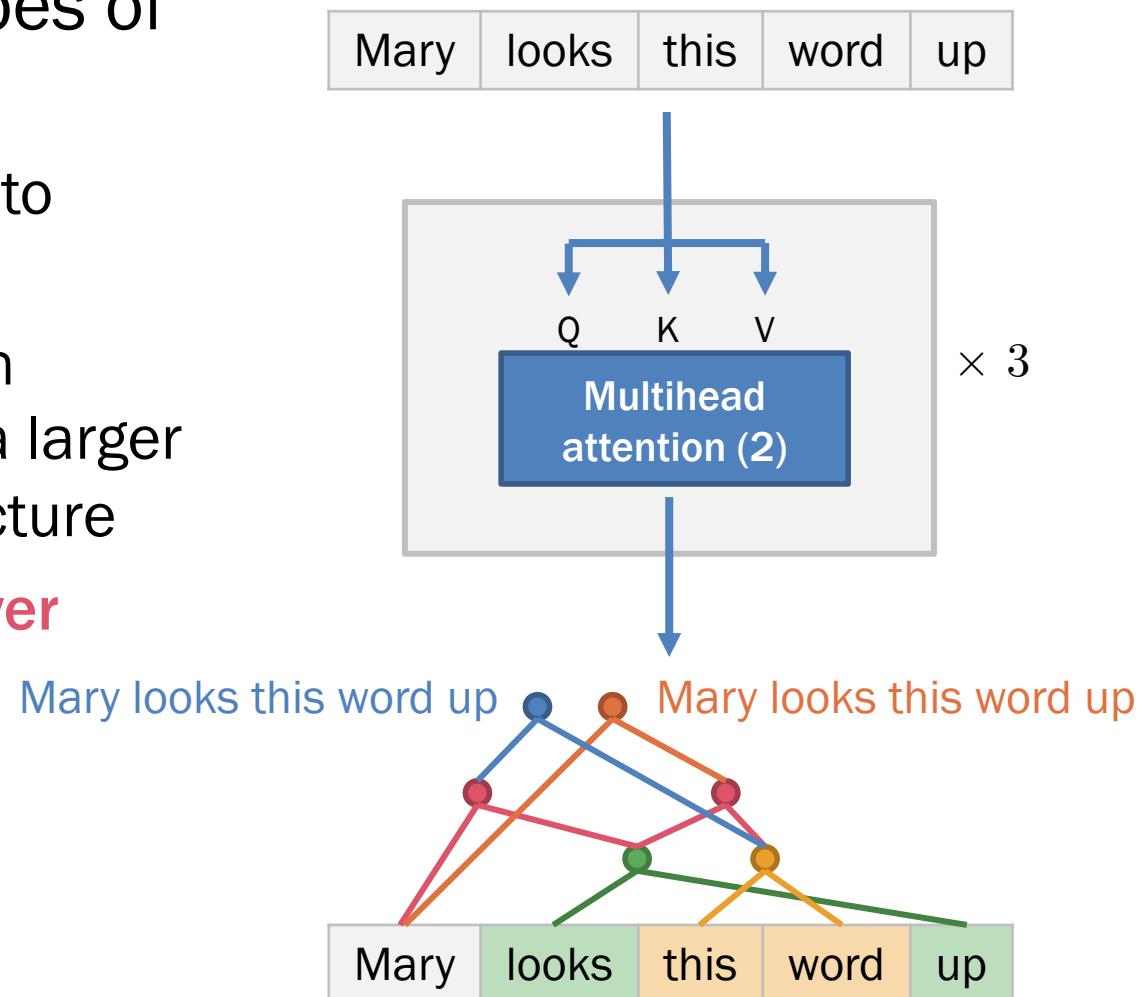
Phrase Structure

- H -head self-attention recognizes H types of word collocation per sequence
 - One layer can combine consecutive words to become a phrase
 - More layers of multihead self-attention can combine consecutive phrases to become a larger phrase or even a sentence \Rightarrow phrase structure
 - Each layer is simply called an **encoding layer**



Phrase Structure

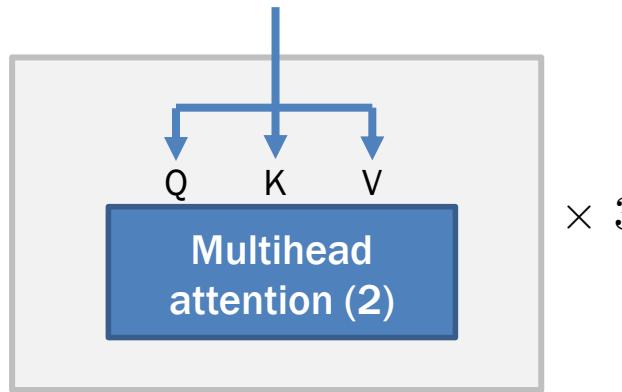
- H -head self-attention recognizes H types of word collocation per sequence
 - One layer can combine consecutive words to become a phrase
 - More layers of multihead self-attention can combine consecutive phrases to become a larger phrase or even a sentence \Rightarrow phrase structure
 - Each layer is simply called an **encoding layer**



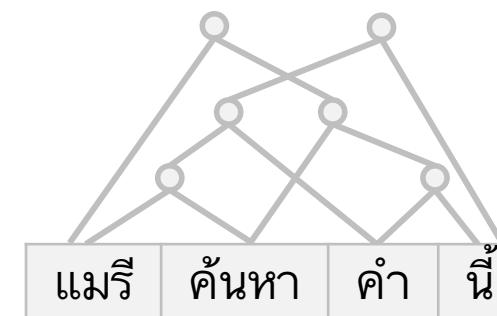
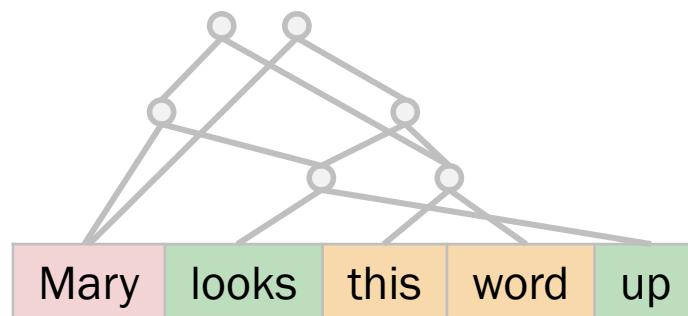
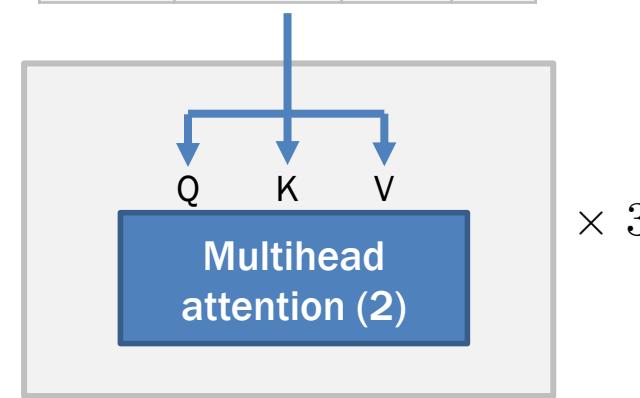
Alignment of Phrase Structures

- H -head alignment attention recognizes H pairs of phrase structures

Mary looks this word up

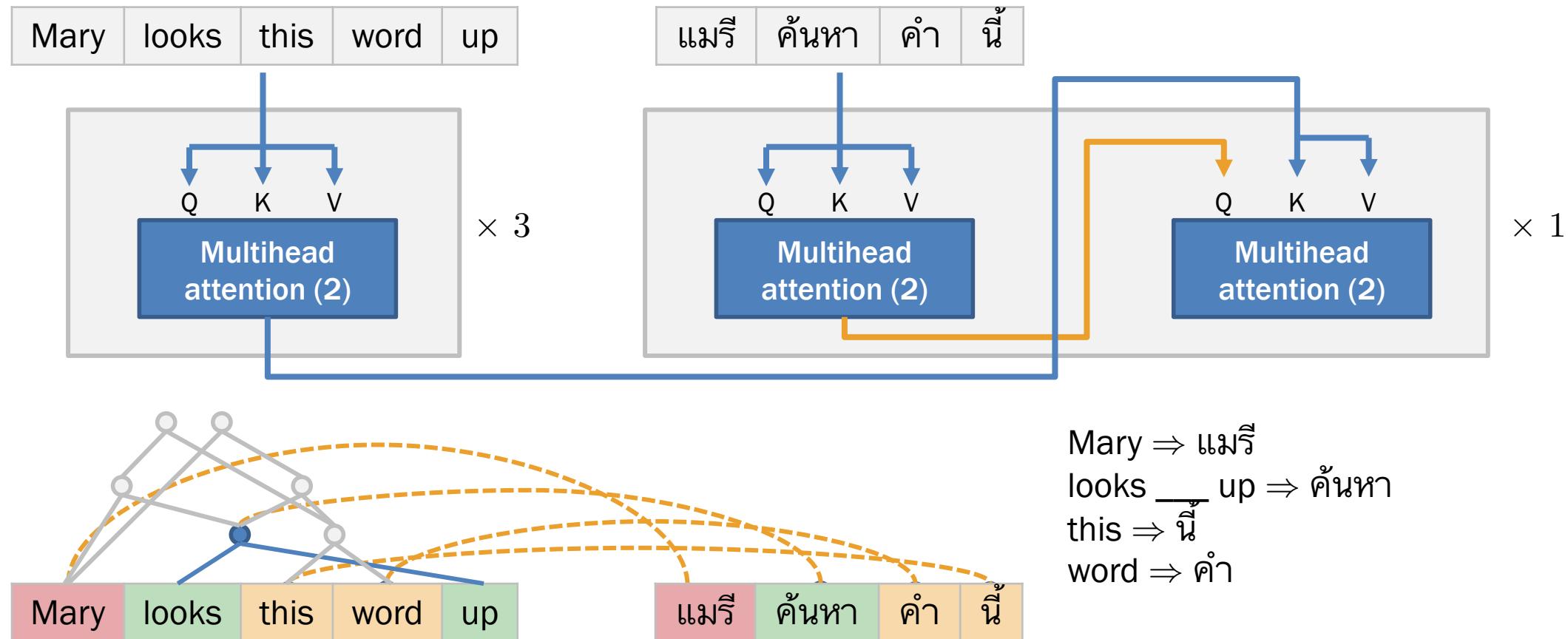


แมรี ค้นหา คำ นี่



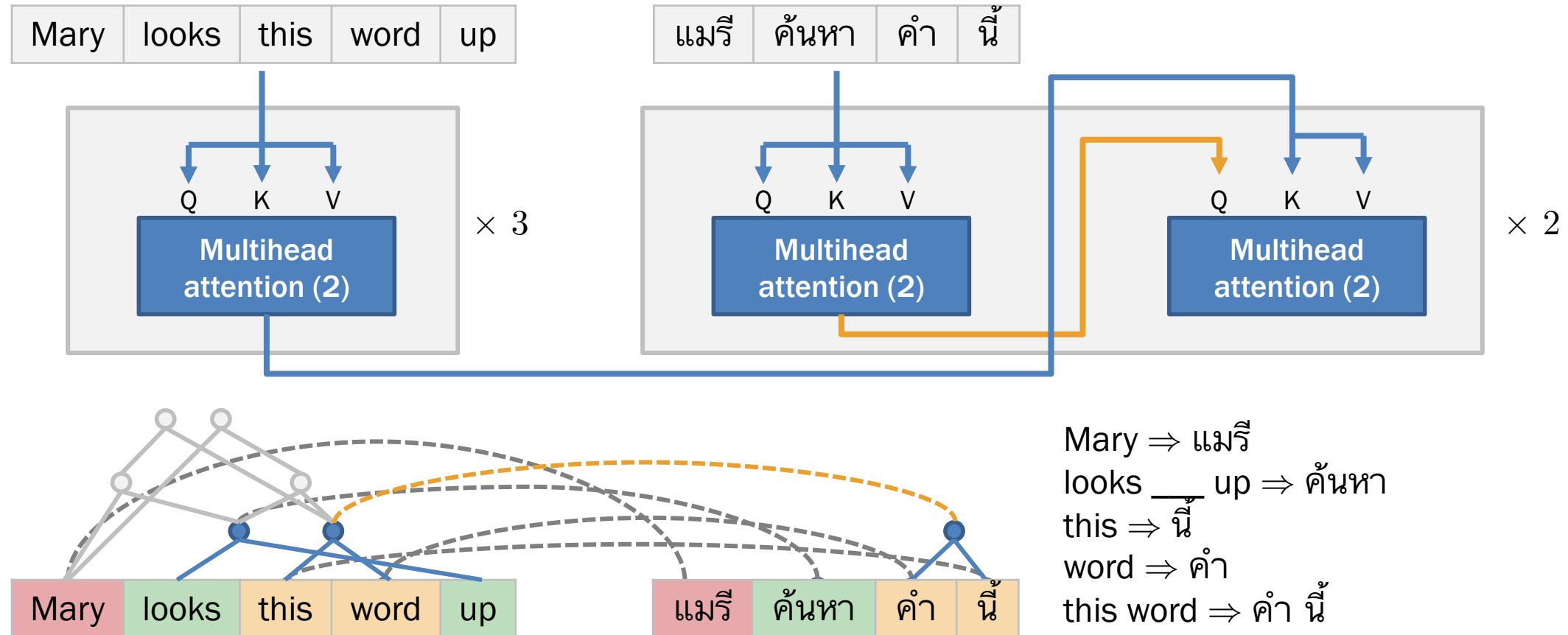
Alignment of Phrase Structures

- H -head alignment attention recognizes H pairs of phrase structures \Rightarrow **decoding layer**



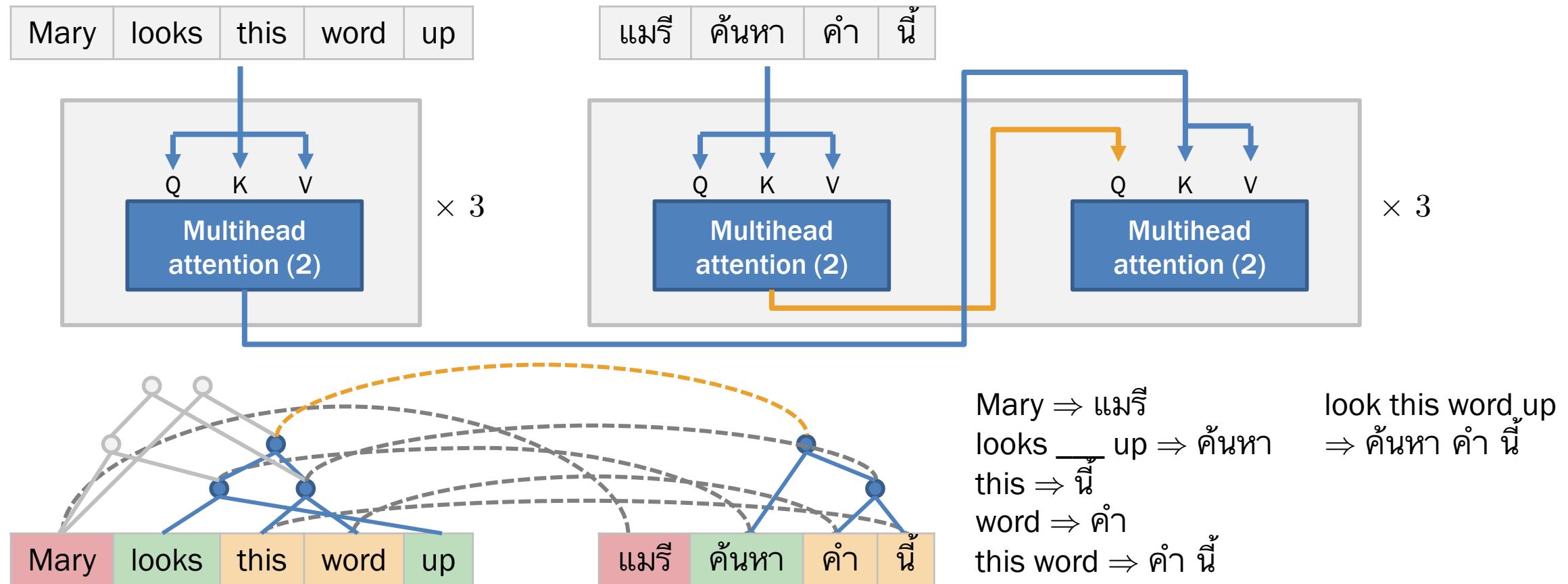
Alignment of Phrase Structures

- H -head alignment attention recognizes H pairs of phrase structures \Rightarrow **decoding layer**



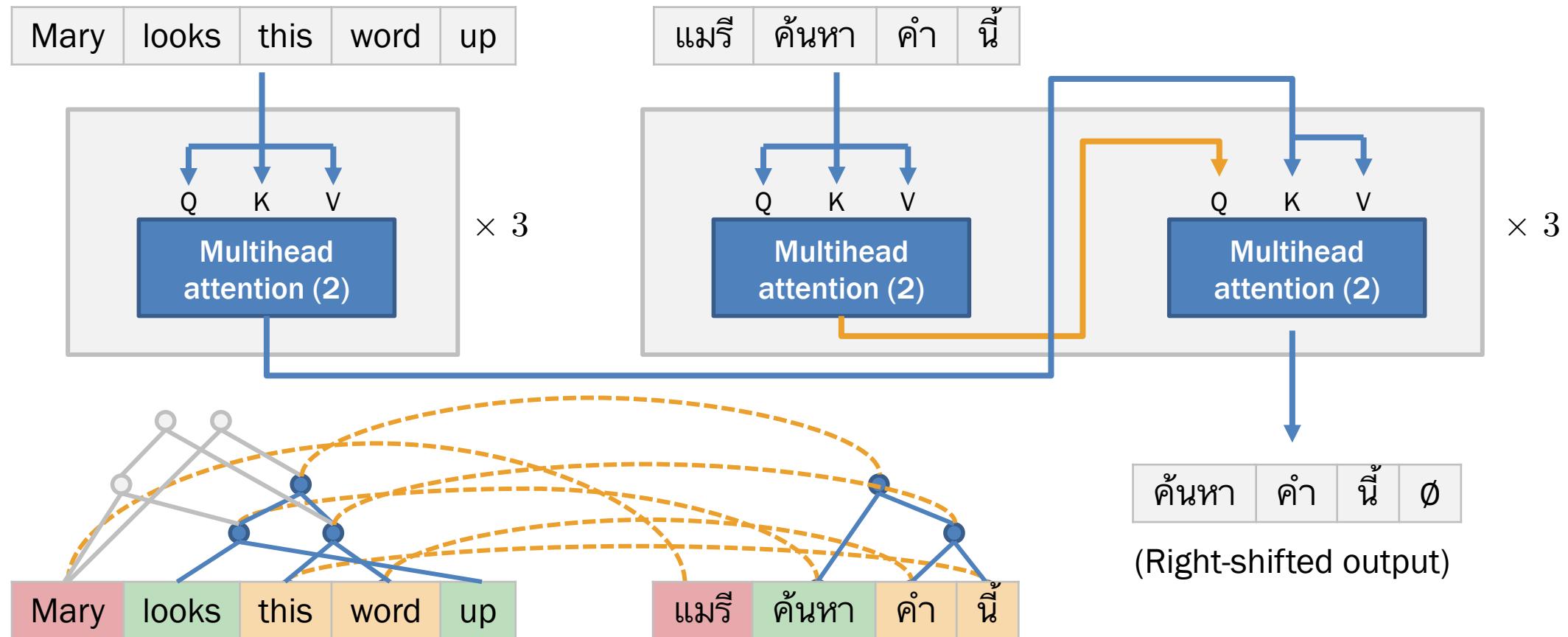
Alignment of Phrase Structures

- H -head alignment attention recognizes H pairs of phrase structures \Rightarrow **decoding layer**

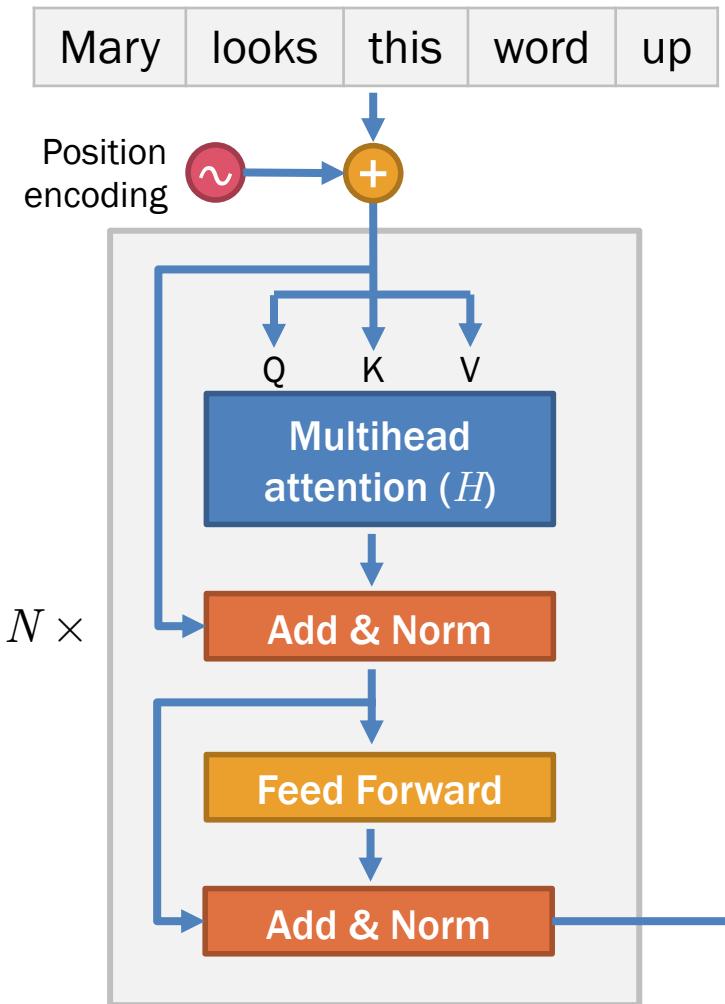


Alignment of Phrase Structures

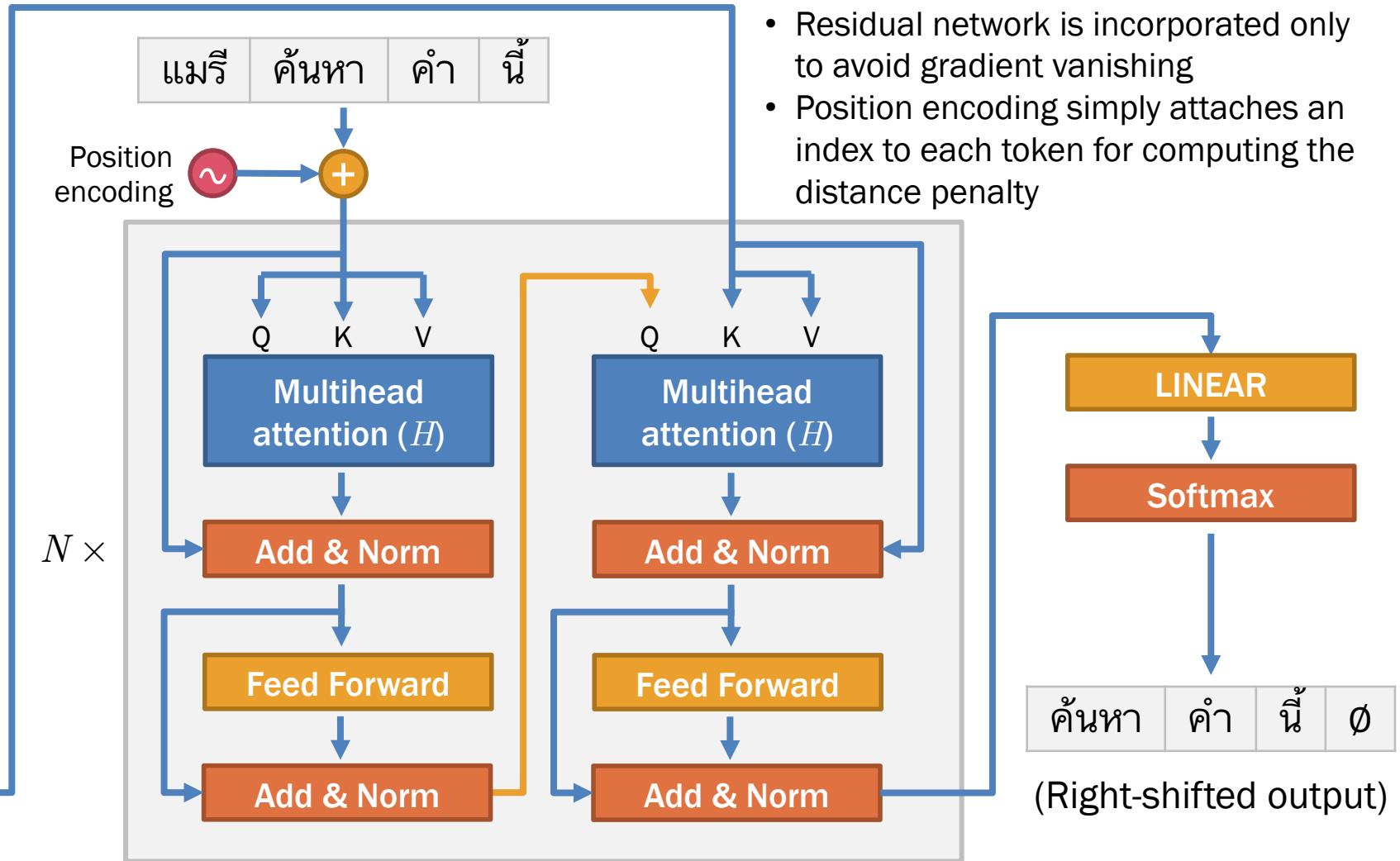
- H -head alignment attention recognizes H pairs of phrase structures \Rightarrow **decoding layer**



Overview of the Transformer Model



1. Extract phrase structures in the source

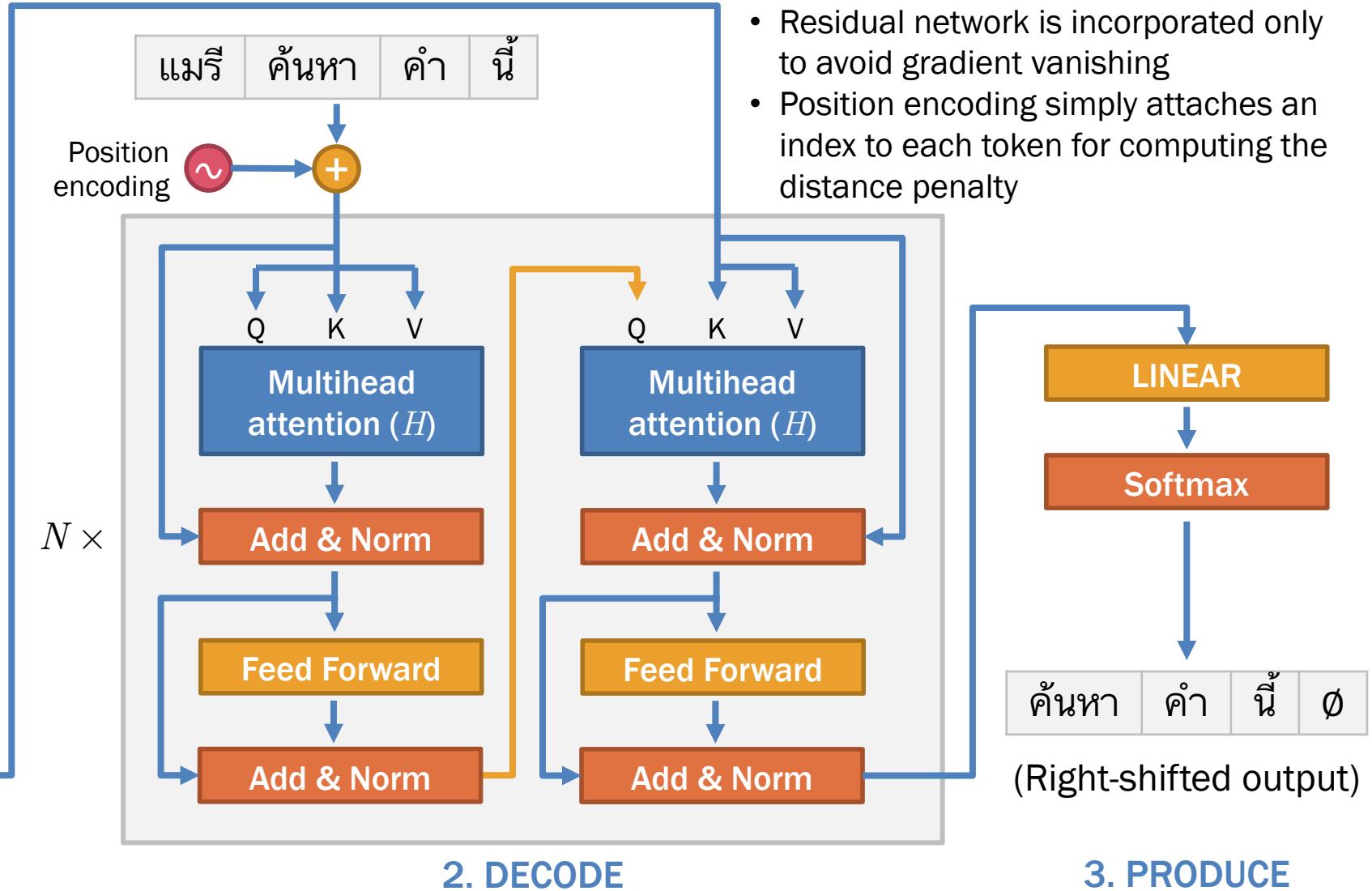
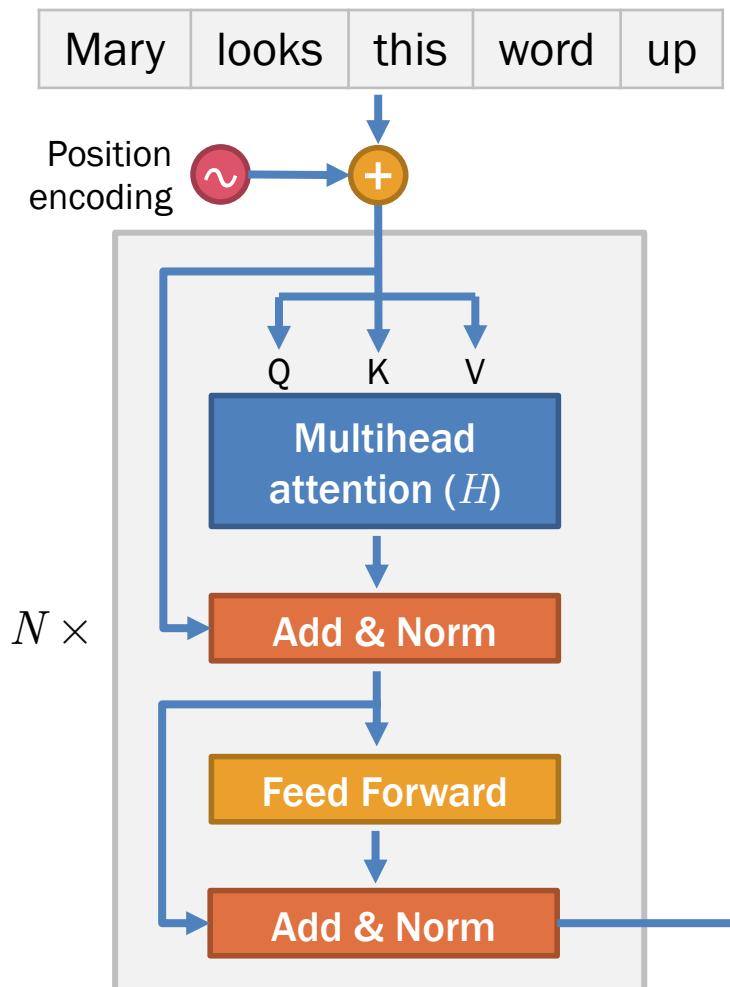


2. Align phrase structures to the target

3. Produce translation

- Residual network is incorporated only to avoid gradient vanishing
- Position encoding simply attaches an index to each token for computing the distance penalty

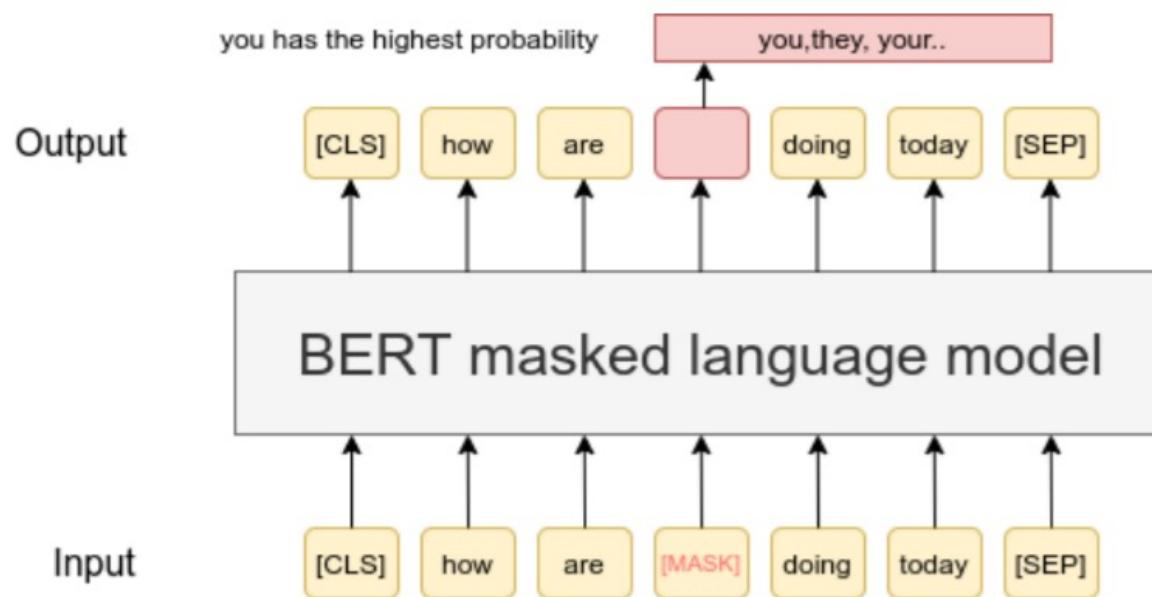
Overview of the Transformer Model



4. Training an LLM

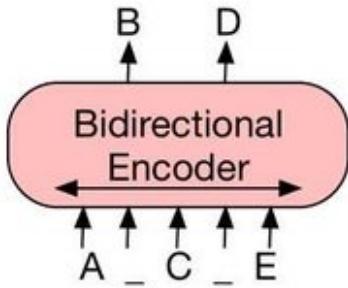
Masked Language Model (MLM)

- We replace some words in the input with blanks and compute the loss of word prediction on these blanks in the output

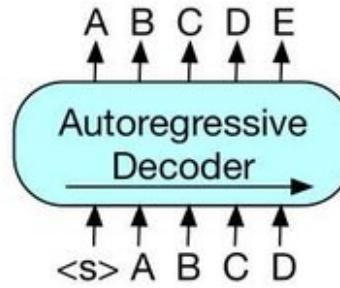


- Each input text is marked at some words by [MASK]
- Masking percentage = 15%
- Once marked, the masks will not be changed
- Special tokens
 - [CLS] = classifier token
 - [SEP] = separator token
 - [MASK] = mask

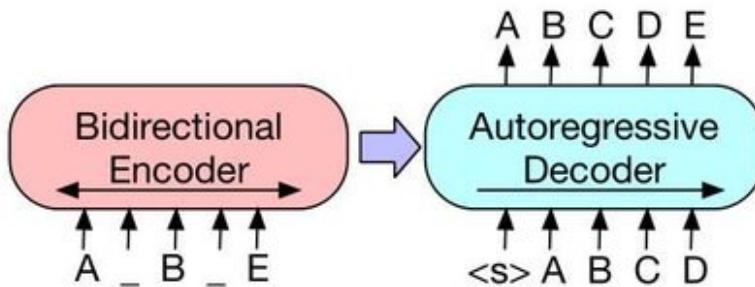
Approaches of Training LLMs (Lewis et al., 2019)



(a) BERT: Random tokens are replaced with masks, and the document is encoded bidirectionally. Missing tokens are predicted independently, so BERT cannot easily be used for generation.



(b) GPT: Tokens are predicted auto-regressively, meaning GPT can be used for generation. However words can only condition on leftward context, so it cannot learn bidirectional interactions.

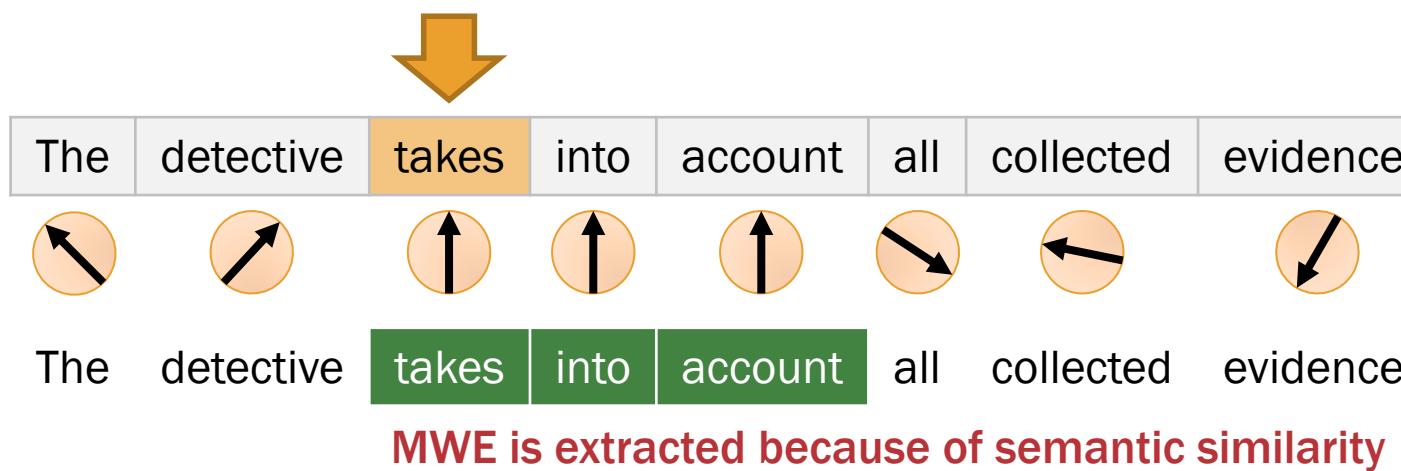


(c) BART: Inputs to the encoder need not be aligned with decoder outputs, allowing arbitrary noise transformations. Here, a document has been corrupted by replacing spans of text with a mask symbol. The corrupted document (left) is encoded with a bidirectional model, and then the likelihood of the original document (right) is calculated with an autoregressive decoder. For fine-tuning, an uncorrupted document is input to both the encoder and decoder, and we use representations from the final hidden state of the decoder.

- **BERT:**
 - Bidirectional encoder
- **GPT:**
 - Generative Pretrained Transformer
 - Autoregressive (unidirectional) decoder
- **BART:**
 - Bidirectional encoder + autoregressive decoder

Pros: Multiword Expression (MWE)

- It recognizes the idiosyncratic collocations of at least 2 words
 - E.g. ‘peanut butter’, ‘car park’, ‘kick the bucket’, ‘take into account’, ‘break up’
 - It learns MWEs by comparing each word with the remaining to reveal semantic similarity



Pros: Moderate-Distance Dependency

- It recognizes word collocation that is separate within a moderate distance
 - E.g. ‘look ____ up’, ‘ask ____ out’, ‘pay ____ for’
 - It learns moderate-distance dependency with semantic similarity and distance penalty



My	husband	paid	his	money	for	a	course	for	car	mechanics
----	---------	------	-----	-------	-----	---	--------	-----	-----	-----------



My	husband	paid	his	money	for	a	course	for	car	mechanics
----	---------	------	-----	-------	-----	---	--------	-----	-----	-----------

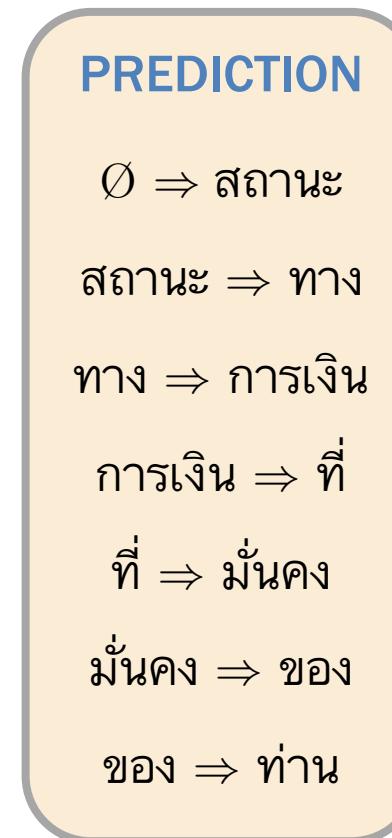
The similarity is penalized by the distance



Pros: Moderate Reordering

- It learns to reorder words with next-word prediction (language model), cross-lingual semantic similarity, and distance penalty

your	stable	financial	status
your	stable	financial	status
your	stable	financial	status
your	stable	financial	status
your	stable	financial	status
your	stable	financial	status
your	stable	financial	status
your	stable	financial	status



Next-word prediction takes into account an entire input sequence

สถานะ
สถานะ ทาง
สถานะ ทาง การเงิน
สถานะ ทาง การเงิน ที่
สถานะ ทาง การเงิน ที่ มั่นคง
สถานะ ทาง การเงิน ที่ มั่นคง ของ
สถานะ ทาง การเงิน ที่ มั่นคง ของ ท่าน

Pros: Conceptualization

- It learns to conceptualize a long subsequence into a shorter one with semantic similarity and distance penalty
 - E.g. ‘initiated a scheme for building ____’ is conceptualized into ‘invented’ and consequently translated into ‘ประดิษฐ์’

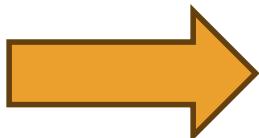
Stevenson	initiated	a	scheme	for	building	the	first	locomotive	
Stevenson	initiated	a	scheme	for	building	the	first	locomotive	PREDICTION
Stevenson	initiated	a	scheme	for	building	the	first	locomotive	$\emptyset \Rightarrow \text{สตีเวนสัน}$
Stevenson	initiated	a	scheme	for	building	the	first	locomotive	$\text{สตีเวนสัน} \Rightarrow \text{ประดิษฐ์}$
Stevenson	initiated	a	scheme	for	building	the	first	locomotive	$\text{ประดิษฐ์} \Rightarrow \text{รถจักรไอน้ำ}$
Stevenson	initiated	a	scheme	for	building	the	first	locomotive	$\text{รถจักรไอน้ำ} \Rightarrow \text{คัน}$
Stevenson	initiated	a	scheme	for	building	the	first	locomotive	$\text{คัน} \Rightarrow \text{แรก}$

4. Training ChatGPT

Constructing ChatGPT

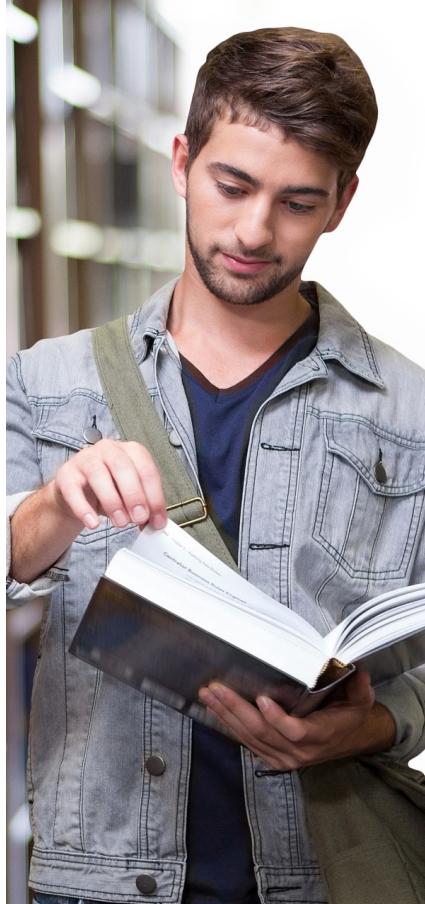


Toddler
babbling a
gibberish



TRAIN

Let it read a
voluminous
amount of texts



High-school
student with
adequate
language
proficiency



FINE-TUNE

Train it to do a
set of specific
tasks



Intelligent
chatbot

1) Model:
Transformer Model

2) Language Model:
GPT, BERT, LLaMa, etc.

3) Dataset:
Instructions and Conversations

Instruction Dataset (Ouyang et al., 2022)

Table 1: Distribution of use case categories from our API prompt dataset.

Use-case	(%)
Generation	45.6%
Open QA	12.4%
Brainstorming	11.2%
Chat	8.4%
Rewrite	6.6%
Summarization	4.2%
Classification	3.5%
Other	3.5%
Closed QA	2.6%
Extract	1.9%

Table 2: Illustrative prompts from our API prompt dataset. These are fictional examples inspired by real usage—see more examples in Appendix A.2.1.

Use-case	Prompt
Brainstorming	List five ideas for how to regain enthusiasm for my career
Generation	Write a short story where a bear goes to the beach, makes friends with a seal, and then returns home.
Rewrite	This is the summary of a Broadway play: """ {summary} """ This is the outline of the commercial for that play: """

Use Case	Example
brainstorming	List five ideas for how to regain enthusiasm for my career
brainstorming	What are some key points I should know when studying Ancient Greece?
brainstorming	What are 4 questions a user might have after reading the instruction manual for a trash compactor? {user manual}
	1.

- **Prompt: “instruction”**
- Cleverly designed set of instructions and responses for a chatbot
- Covering frequently asked questions and their answers

User Prompts and Chats (Ouyang et al., 2022)

closed qa	<p>Answer the following question: What shape is the earth?</p> <p>A) A circle B) A sphere C) An ellipse D) A plane</p>
closed qa	<p>Tell me how hydrogen and helium are different, using the following facts:</p> <p>{list of facts}</p>
open qa	<p>I am a highly intelligent question answering bot. If you ask me a question that is rooted in truth, I will give you the answer. If you ask me a question that is nonsense, trickery, or has no clear answer, I will respond with "Unknown".</p> <p>Q: What is human life expectancy in the United States? A: Human life expectancy in the United States is 78 years.</p> <p>Q: Who was president of the United States in 1955? A:</p>
open qa	<p>Who built the statue of liberty?</p>
open qa	<p>How do you take the derivative of the sin function?</p>
open qa	<p>who are the indigenous people of New Zealand?</p>

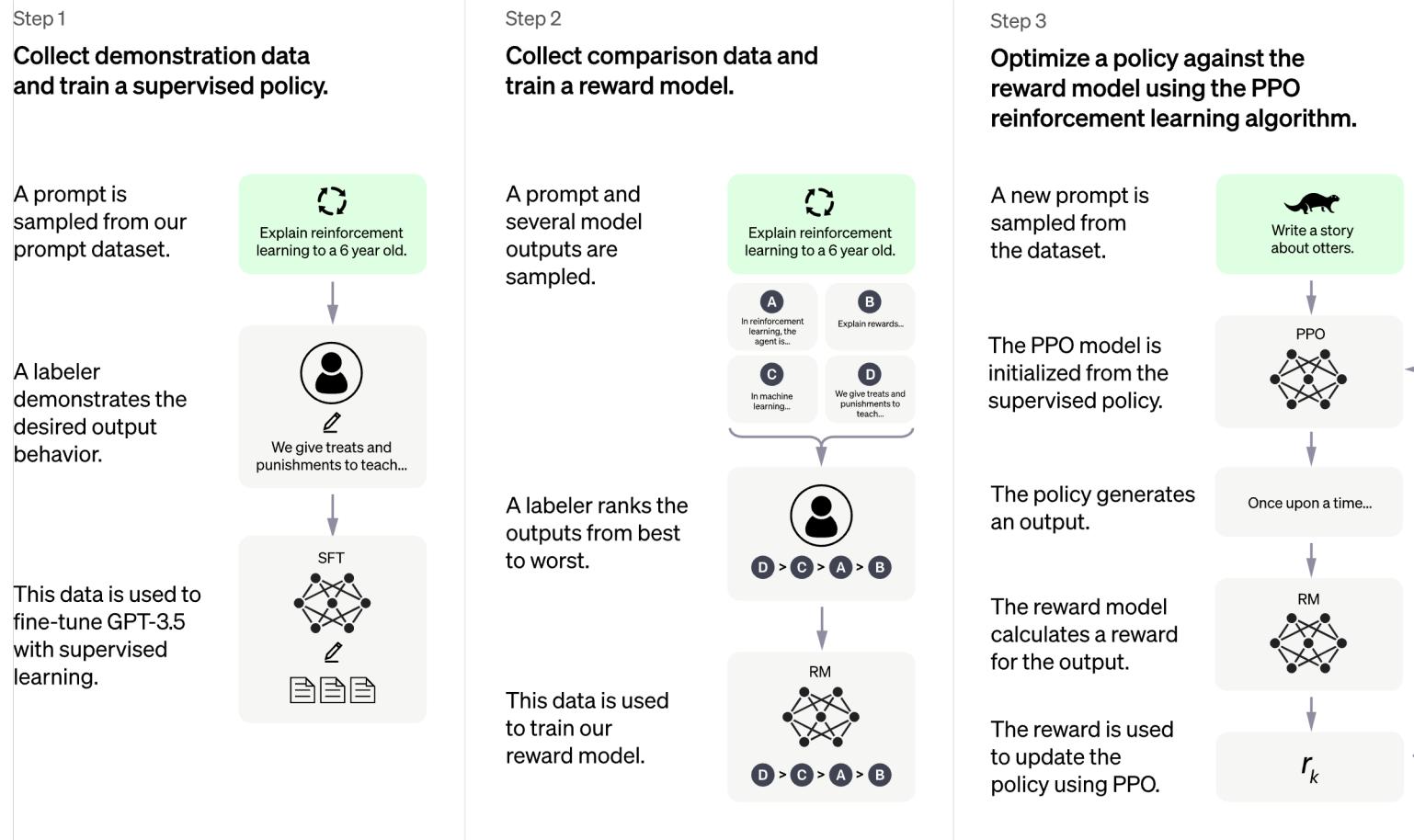
User Prompts and Chats (Ouyang et al., 2022)

Use Case	Example
chat	<p>The following is a conversation with an AI assistant. The assistant is helpful, creative, clever, and very friendly.</p> <p>Human: Hello, who are you? AI: I am an AI created by OpenAI. How can I help you today? Human: I'd like to cancel my subscription. AI:</p>
chat	<p>Marv is a chatbot that reluctantly answers questions with sarcastic responses:</p> <p>You: How many pounds are in a kilogram? Marv: This again? There are 2.2 pounds in a kilogram. Please make a note of this. You: What does HTML stand for? Marv: Was Google too busy? Hypertext Markup Language. The T is for try to ask better questions in the future. You: When did the first airplane fly? Marv:</p>
chat	<p>This is a conversation with an enlightened Buddha. Every response is full of wisdom and love.</p> <p>Me: How can I achieve greater peace and equanimity? Buddha:</p>

User Prompts and Chats (Ouyang et al., 2022)

Use Case	Example
classification	<p>The following is a list of companies and the categories they fall into:</p> <p>Apple, Facebook, Fedex</p> <p>Apple Category: Technology</p> <p>Facebook Category: Social Media</p> <p>Fedex Category:</p>
extract	<p>Text: {text}</p> <p>Keywords:</p>
generation	"Hey, what are you doing there?" Casey was startled. He hadn't even begun to
generation	The name of the next Star Wars movie is
generation	This is the research for an essay: ==== {description of research} ==== Write a high school essay on these topics: ====

Prompts + Human Ranking + RL



- 3 steps
 - Fine-tune the language model with the instruction dataset
 - Retrain the reward model for chat response with human ranking
 - Optimize the policy model w.r.t. the reward model with the PPO Algorithm (proximal policy optimization)

5. OpenThaiGPT

Training Datasets

Pretraining datasets

Aa Name	Source Type	Access Type	token size (GPT Thai)	size (GB)
mC4	Web Crawl	Public	16B	56 GB
OSCAR 2019	Web Crawl	Public	3.4B	16GB
OSCAR 2021	Web Crawl	Public	3.8B	16GB
OSCAR 2022	Web Crawl	Public	10.8B	60GB
OSCAR 2023	Web Crawl	Public Access needed	17B	100GB
CC100	Web Crawl	Public	12B	80GB
LST20	News	Public Access needed	3M (Compound Word)	
Twitter	Social	P Conan Private	560M	
Linetoday	News	P Conan Private	105M	
TraffyFondue	Complaint	P Conan Private	5M	
Wikipedia	Wiki	Public	10M	
prachathai67k	News	Public	160M	
ThaiPBS	News	Public		
Scb-th-en(extract th)	MISC	Public	50.3M	
wisesight_sentiment	Social	Public	0.978M	
Wongnai_reviews	Social	Public	11.1M	
ThaiRath		Public		
Best		Public Access needed		

- OpenThaiGPT was trained on the dataset consisting of over **2 trillion tokens** (sub-words): CC100, OSCAR, and mC4
- After cleansing, decontamination, and anonymization, we obtained **37.3 billion tokens**
- Pantip.com supplied additional **20.0 billion tokens** of social media datasets
- In total, we obtained **60 billion tokens** for training

LLM Project @ NECTEC

- It simulates how humans string words to become a meaningful sentence
- Trained models:
 - **LLaMa (GPT)** and **Mistral** for language generation
 - **ELECTRA** for document classification
 - **BART** for machine translation and document summarization
 - **DeBERTa** for classification and zero-shot learning
 - **Multimodality:** text, speech, images, videos
- Making Thai one of the pivot languages for LLM via knowledge distillation and model grafting



OpenThaiGPT 1.0.0

Available on <https://openthaigpt.openservice.in.th> and Hugging Face

The screenshot shows the OpenThaiGPT 1.0.0-alpha web interface. At the top, there's a header with the project name and a sub-header about the first Thai implementation of a 7B-parameter LLaMA v2 Chat model. Below the header, there's a section titled "Examples" with six buttons: "ลดความอ้วนต้องทำอย่างไร", "วางแผนที่ยวainภูเก็ต", "เขียนบทความ", "เขียนโค้ด", "คำนวณคณิตศาสตร์", and "แปลภาษา". Underneath this is a "Instruction" field containing the text "อธิบายเรื่องเคมีควบคุมตั้มให้หน่อย". Below it is an "Input" field with the text "คำダメ (ไม่จำเป็น)". A checked checkbox labeled "Stream output" is present. A large orange "Generate" button is centered below the input fields. At the bottom, there are "Stop / Cancel" and "Clear" buttons. The "Output" section at the bottom contains a large block of generated text in Thai, which is a continuation of the instruction about chemistry.

- **Base:** LLaMa v.2
- **Sizes:** 7 and 13 billion parameters (neurons)
- 70B parameters will be released in Feb 2024
- Natural word-level generation
- BPE tokenizer with 20,000 Thai frequent tokens
- Launched on AI4Thai

🌟 Open LLM Leaderboard

💡 The 🌟 Open LLM Leaderboard aims to track, rank and evaluate open LLMs and chatbots.

🌟 Submit a model for automated evaluation on the 🌟 GPU cluster on the "Submit" page! The leaderboard's backend runs the great [Eleuther AI Language Model Evaluation Harness](#) - read more details in the "About" page!

Archived on 14 SEP 2023

[LLM Benchmark](#)
 [About](#)
 [Submit here!](#)

Select columns to show

Average
 ARC
 HellaSwag
 MMLU
 TruthfulQA
 Type
 Precision
 Hub License
 #Params (B)
 Hub

 Model sha

Show gated/private/deleted models

thai

Model types

pretrained
 fine-tuned
 instruction-tuned
 RL-tuned

Model sizes

Unknown
 < 1.5B
 ~3B
 ~7B
 ~13B
 ~35B
 60B+

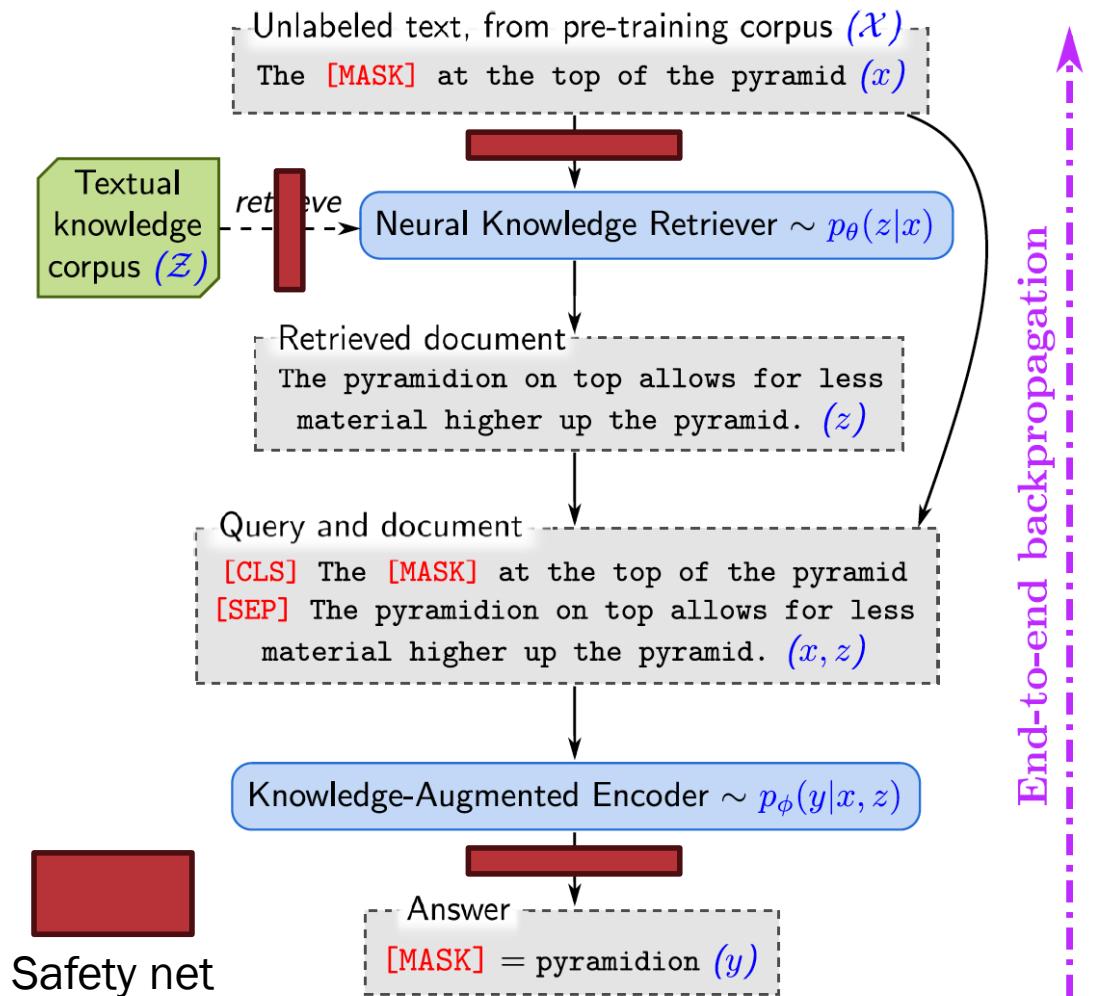
T	Model	Average	ARC	HellaSwag	MMLU	TruthfulQA
	openthaigpt/openthaigpt-1.0.0-alpha-7b-chat-ckpt-hf	53.25	50.85	74.89	40.02	47.23
	wannaphong/openthaigpt-0.1.0-beta-full-model_for_open_llm_leaderboard	51.3	51.28	77.46	33.18	43.28
	pythainlp/wangchanglm-7.5B-sft-en-sharded	38.7	34.47	59.81	26.37	34.15
	nvthainlp/wangchanglm-7.5B-sft-enth	38.05	33.79	58.99	24.52	34.9

Capabilities

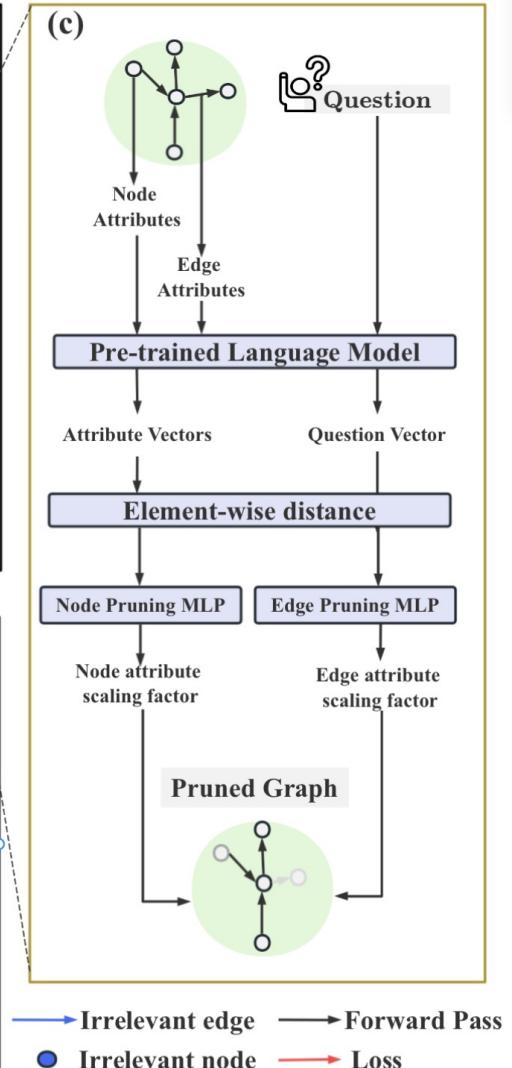
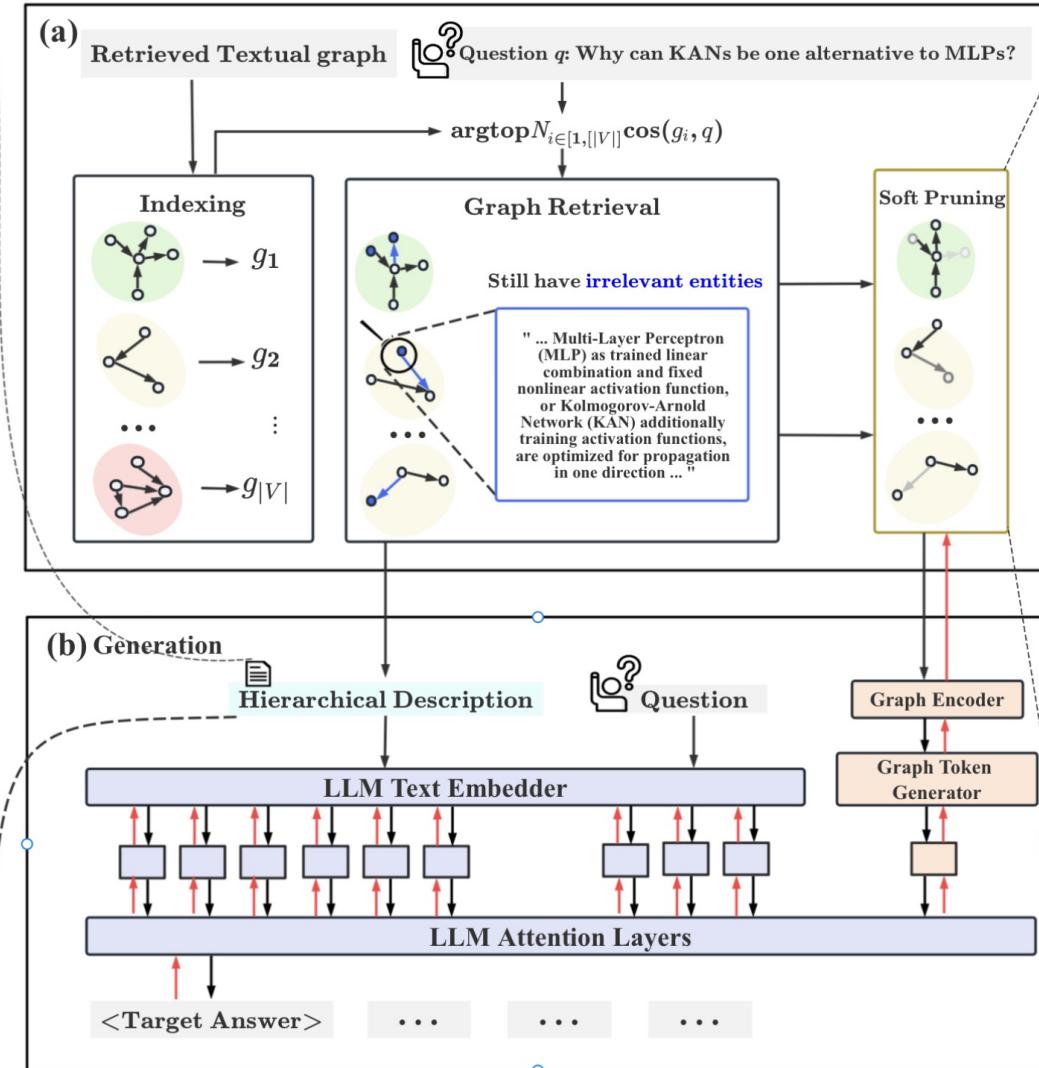
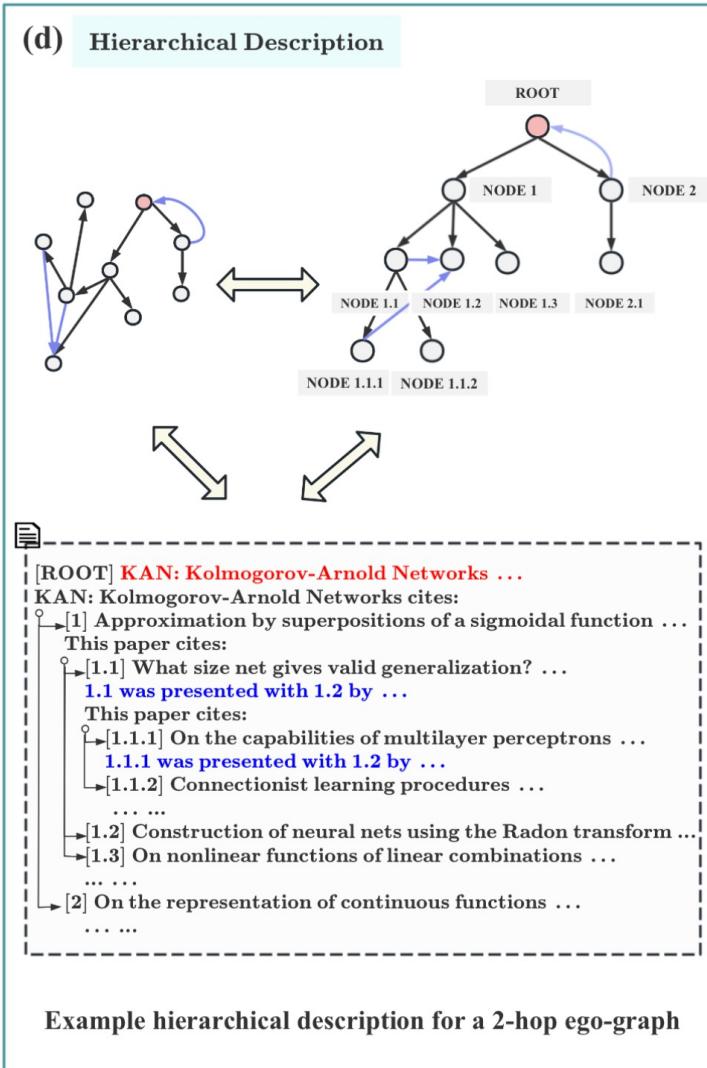
- Question answering from Wikipedia
- Trip planning
- Explaining steps and formulas
- Writing short descriptions
- Writing short pieces of code
- Solving math problems described in natural language
- Translation from/to Thai

Retrieval-Augmented Generation (RAG)

- Enhancing question answering by enriching the query with relevant info
 - Mitigating the hallucination issue in traditional LLMs
 - In-context learning:** simply attaching the query with relevant documents retrieved with information retrieval
 - Vector databases (e.g. FAISS and VectorDB) are used to store internal documents on on-premise servers
 - Safety net:** content filtering



RAG with Knowledge Graph



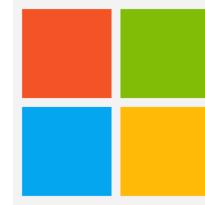
World Arena of LLMs

Multilingual (English)



Google

BERT (340M)
GLaM (1.2T)
LaMDA (137B)
PaLM (540B)
Minerva (540B)
PaLM-2 (340B)



Microsoft

Megatron-Turing (530B)



OpenAI

GPT-2 (1.5B)
GPT-3.5 (175B)
GPT-4 (1T)



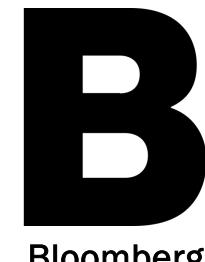
DeepMind

Gopher (280B)
Chinchilla (70B)



Meta

OPT (175B)
Galactica (120B)
LLaMa (65B)
LLaMa-2 (70B)



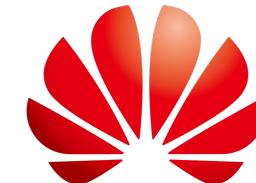
Bloomberg

BloombergGPT (50B)

Chinese



ERNIE (260B)



HUAWEI

PanGu- Σ (1.085T)

Thai



ARTIFICIAL INTELLIGENCE ASSOCIATION OF THAILAND



a member of NSTDA

OpenThaiGPT (7B)
equivalent to GPT-3
OpenThaiGPT (13B)
equivalent to GPT-3.5



WangchanBERTa
WangchanGLM

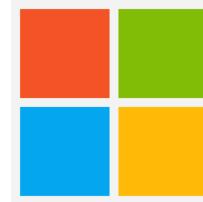
Typhoon (7B)

World Arena of LLMs

Multilingual (English)



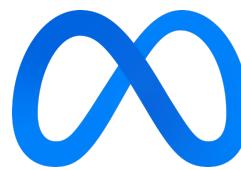
BERT (340M)
GLaM (1.2T)
LaMDA (137B)
PaLM (540B)
Minerva (540B)
PaLM-2 (340B)



Megatron-Turing (530B)



OpenAI



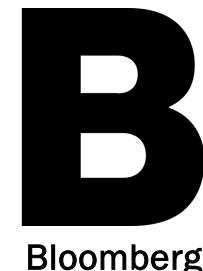
Meta

GPT-2 (1.5B)
GPT-3.5 (175B)
GPT-4 (1T)



DeepMind

Gopher (280B)
Chinchilla (70B)



Bloomberg

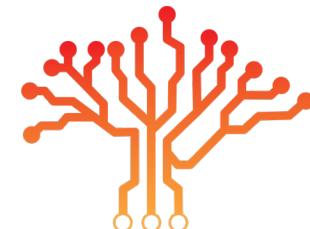
BloombergGPT (50B)

Southeast Asian Languages



Alibaba Group
阿里巴巴集团

SEA-LLM (7B, 13B)



AI SINGAPORE®
SEA-LION (7B)

Thai



OpenThaiGPT (7B)
equivalent to GPT-3
OpenThaiGPT (13B)
equivalent to GPT-3.5



WangchanBERTa
WangchanGLM

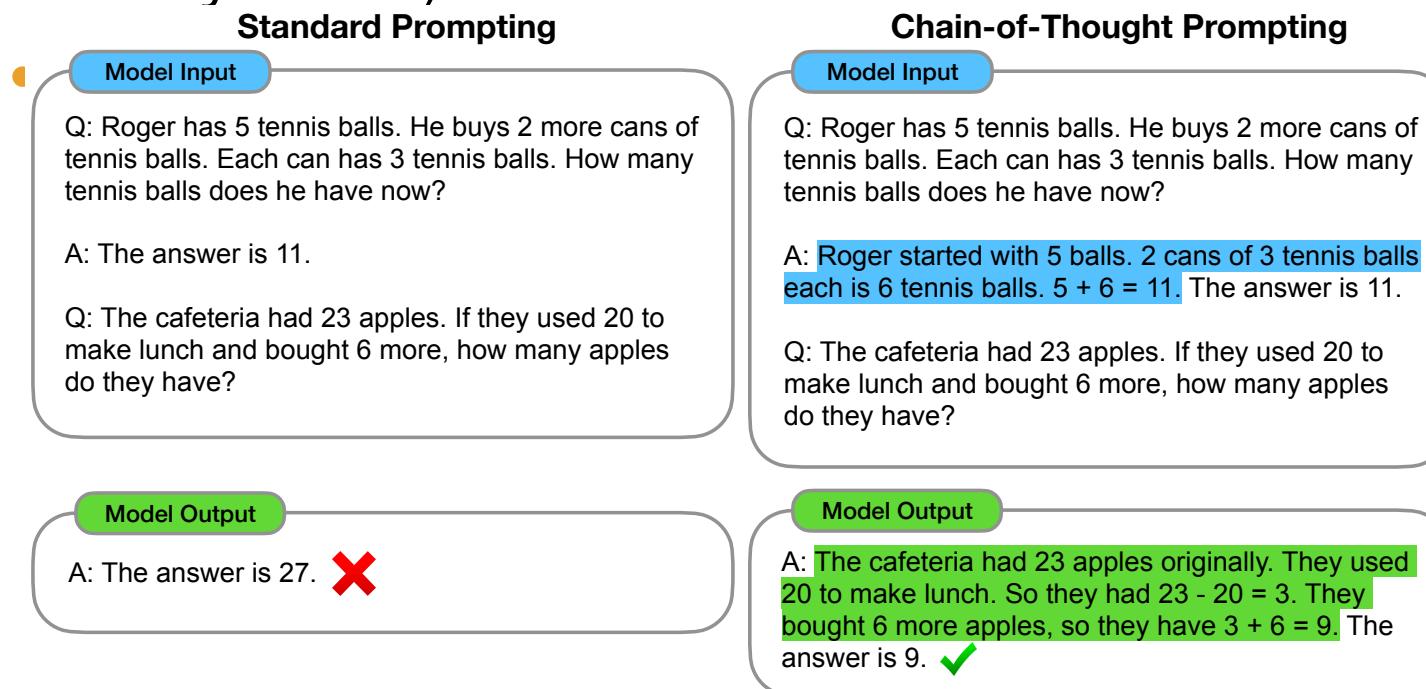


Typhoon (7B)

6. Future Direction

Future: PaLM (Wei et al., 2022)

- Pathways Language Model
 - 540B parameters trained with Pathways Transformer model (single model capable of doing many tasks)



- Chain of thoughts is added to the target sequence
- Reasoning and answering are learned in tandem
- Interpretability is introduced to Transformer

Future: PaLM (Wei et al., 2022)

- Chain of thoughts

<p>Math Word Problems (free response)</p> <p>Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?</p> <p>A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.</p>	<p>Math Word Problems (multiple choice)</p> <p>Q: How many keystrokes are needed to type the numbers from 1 to 500? Answer Choices: (a) 1156 (b) 1392 (c) 1480 (d) 1562 (e) 1788</p> <p>A: There are 9 one-digit numbers from 1 to 9. There are 90 two-digit numbers from 10 to 99. There are 401 three-digit numbers from 100 to 500. $9 + 90(2) + 401(3) = 1392$. The answer is (b).</p>	<p>CSQA (commonsense)</p> <p>Q: Sammy wanted to go to where the people were. Where might he go? Options: (a) race track (b) populated areas (c) desert (d) apartment (e) roadblock</p> <p>A: The answer must be a place with a lot of people. Race tracks, desert, apartments, and roadblocks don't have a lot of people, but populated areas do. So the answer is (b).</p>
<p>StrategyQA</p> <p>Q: Yes or no: Would a pear sink in water?</p> <p>A: The density of a pear is about 0.6 g/cm³, which is less than water. Thus, a pear would float. So the answer is no.</p>	<p>Date Understanding</p> <p>Q: The concert was scheduled to be on 06/01/1943, but was delayed by one day to today. What is the date 10 days ago in MM/DD/YYYY?</p> <p>A: One day after 06/01/1943 is 06/02/1943, so today is 06/02/1943. 10 days before today is 05/23/1943. So the answer is 05/23/1943.</p>	<p>Sports Understanding</p> <p>Q: Is the following sentence plausible? "Joao Moutinho caught the screen pass in the NFC championship."</p> <p>A: Joao Moutinho is a soccer player. The NFC championship is part of American football, not soccer. So the answer is no.</p>
<p>SayCan (Instructing a robot)</p> <p>Human: How would you bring me something that isn't a fruit?</p> <p>Explanation: the user wants something to eat that isn't a fruit. An energy bar is not a fruit, so I will bring the user an energy bar.</p> <p>Plan: 1. find(energy bar) 2. pick(energy bar) 3. find(user) 4. put(energy bar) 5. done().</p>	<p>Last Letter Concatenation</p> <p>Q: Take the last letters of the words in "Lady Gaga" and concatenate them.</p> <p>A: The last letter of "Lady" is "y". The last letter of "Gaga" is "a". Concatenating them is "ya". So the answer is ya.</p>	<p>Coin Flip (state tracking)</p> <p>Q: A coin is heads up. Maybelle flips the coin. Shalonda does not flip the coin. Is the coin still heads up?</p> <p>A: The coin was flipped by Maybelle. So the coin was flipped 1 time, which is an odd number. The coin started heads up, so after an odd number of flips, it will be tails up. So the answer is no.</p>

Future: PaLM-E (Driess et al., 2023)

- Embodied multimodal PaLM with 562B parameters

Mobile Manipulation



Human: Bring me the rice chips from the drawer. Robot: 1. Go to the drawers, 2. Open top drawer. I see 3. Pick the green rice chip bag from the drawer and place it on the counter.

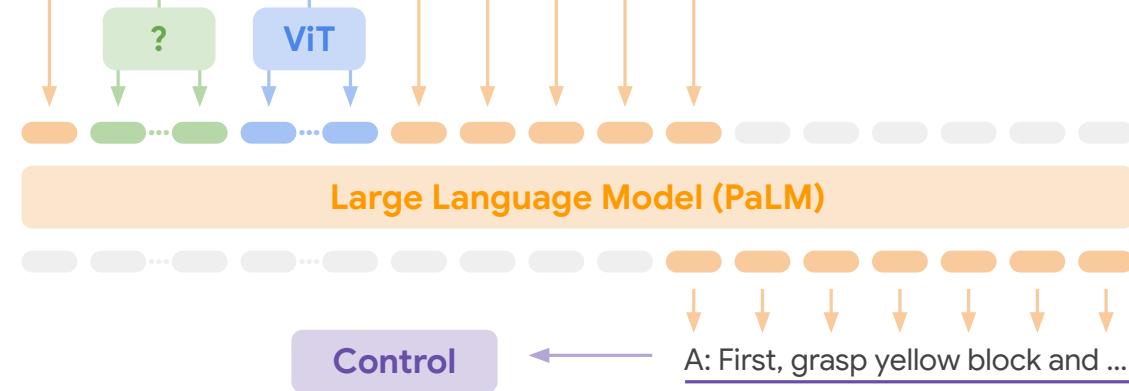
Visual Q&A, Captioning ...



Given Q: What's in the image? Answer in emojis.
A: .

PaLM-E: An Embodied Multimodal Language Model

Given ... Q: How to grasp blue block? A: First, grasp yellow block

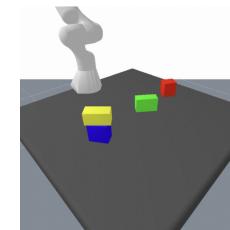


Describe the following :
A dog jumping over a hurdle at a dog show.

Language Only Tasks

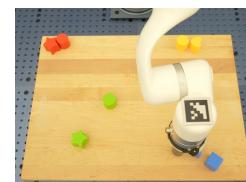
Here is a Haiku about embodied language models:
Embodied language models are the future of natural language

Task and Motion Planning



Given Q: How to grasp blue block?
A: First grasp yellow block and place it on the table, then grasp the blue block.

Tabletop Manipulation



Given Task: Sort colors into corners.
Step 1. Push the green star to the bottom left.
Step 2. Push the green circle to the green star.

Future: PaLM-E (Driess et al., 2023)

Visually-conditioned Jokes + Few-shot prompting



I see (image omitted). Description: A cow sitting on a rug. Joke: Why is the cow sitting on a rug? Because it wants to be a cow-ch! ... I see . Description: a donkey is carrying a dog, cat, and rooster. Joke: what do you call a donkey with a rooster on his back. A rooster booster.

Zero-shot Multimodal Chain-of-Thought



Given . Q: Can I go down this street on a bicycle, yes or no? A: Let's think step by step. 1. do not enter. 2. except bicycles. 3. do not entry except bicycles. 4. yes.



Given . Q: How many championship rings did the player on the left win in his career? A: Let's think step by step. 1. kobe bryant is on the left. 2. kobe bryant has won 5 championship rings.



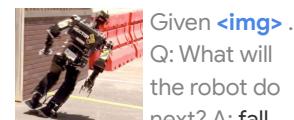
Given . Q: Who are the two teams playing in this photo, and which was the last to win a championship, and which year did they win, and who was their star player that year? A: Let's think step by step. 1. team in white is new york knick. 2. team in green is boston celtic. 3. last to win championship was boston celtic. 4. year was 2008. 5. star player was paul pierce. 6. i am not sure if this is correct. i am not basketball fan. i just google it.

Robot Visual Perception, Dialogue, and Planning



Robot: I am a robot operating in a kitchen. Given , what do I see? a bowl of eggs, a bowl of flour, and a bowl of sugar.. and a spatula.. and a spoon.. and a bag of flour.. and a bag of sugar. When a human asks me to do a task, I will respond with the sequence of actions I would do to accomplish the task with only the items I see. Human: Use all of the ingredients you see to make a cake batter. 1. crack egg. 2. put egg in bowl. 3. put flour in bowl. 4. put sugar in bowl. 5. mix. 6. put in pan. 7. bake. 8. eat. 9. clean up.

Physical Prediction



Given . Q: If a robot wanted to be useful here, what steps should it take? A: clean table..pick up trash.. pick up chairs.. wipe chairs.. put chairs down.



Spatial Grounding



Given . Q: What flavor is the doughnut on the right? Strawberry.. Q: What flavor is the doughnut on the left? A blueberry.

Photo 1: <img1>. Photo 2: <img2>. Q: What is in Photo 1 but NOT in Photo 2? A: sunglasses.

Photo 1: <img1>. Photo 2: <img2>. Q: What is in Photo 1 but NOT in Photo 2? A: Let's think step by step and explain. Photo 1 has sunglasses on top of folded clothes. Photo 2 does not have sunglasses on top of folded clothes. Therefore, sunglasses are in Photo 1 but not in Photo 2.

End-to-End Egocentric Q&A, Dialogue



Photo 1, at 10:30 am: <img1>. Photo 2, at 12:45 pm: <img2>. Photo 3, at 3:45 pm: <img3>. Q: I forget, what did I have for lunch, and what time was it? A: Let's think step by step. 1. you had a sandwich for lunch. 2. it was 12:45 pm.

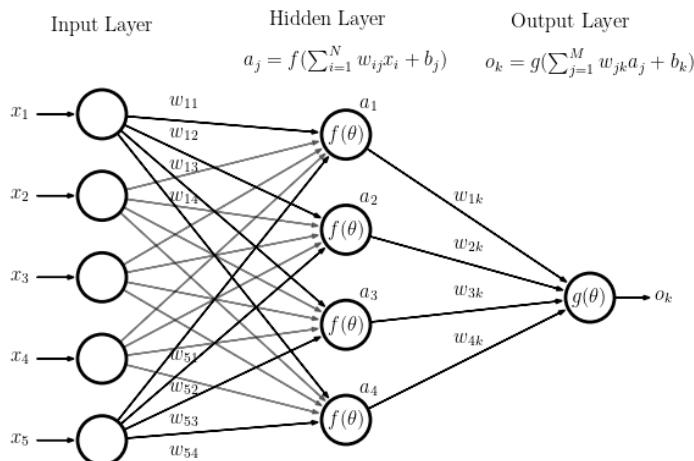
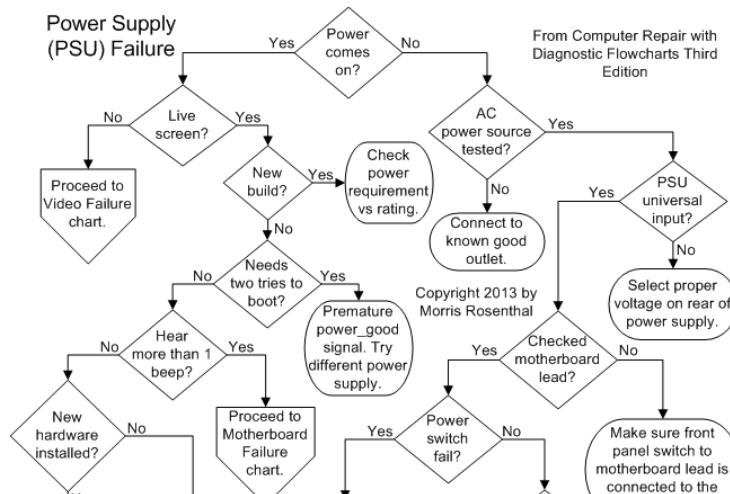
OCR-free Math Reasoning



Given . Q: I am getting just two custom pizza for me and my friend. How much should I pay in total? A: Let's think step by step. 1st step: 2 custom pizzas. 2nd step: 2 x \$9.99. 3rd step: \$19.98.

Quō Vādis, AI? (Where are you going to, AI?)

- Three waves of AI



Wave 1: DESCRIBE

Handcrafted knowledge
and expert systems

<https://stackoverflow.com/questions/14092782/how-to-save-a-decision-flow-chart-as-xml>

Wave 2: CATEGORIZE

Blackbox AI,
statistical learning,
and neural networks

http://www.astroml.org/book_figures/appendix/fig_neural_network.html

Wave 3: EXPLAIN

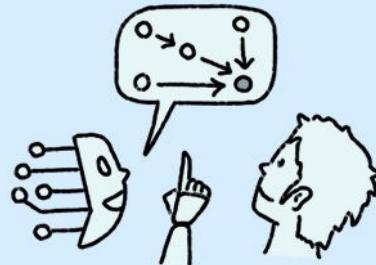
Whitebox AI,
context-aware adaptation,
and knowledge graphs

<https://www.linkedin.com/pulse/explainable-ai-worthwhile-george-baras/>

Contextual AI

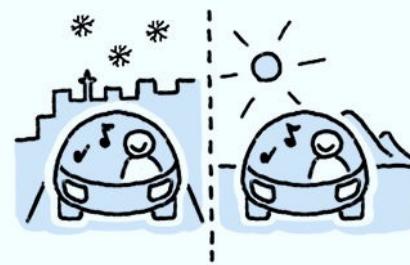
PILLARS OF CONTEXTUAL A.I.

INTELLIGIBLE



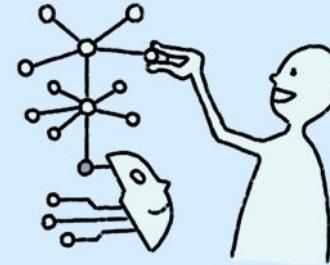
Able to explain what it knows, how it knows, and what it's doing.

ADAPTIVE



Able to meet user's expectations in different environments.

CUSTOMIZABLE



Able to be fully controlled by the user.

CONTEXT-AWARE



Able to perceive at the same level as a human does.

Challenges

- AI hallucination
 - Generated contents may be incorrect due to confabulation (connecting unrelated stories)
 - Distinction between **hallucination** (recognized mistake) vs. **creativity**
- Unethical uses of Generative AI
 - Plagiarism and abuses in learning
 - Generating fake news and misinformation
 - Leaks of personal information

Types of Hallucination

- Out-of-context generation
- Contradiction to previously generated contents
- Misinformation generation

Thank You

prachya.boonkwan@nectec.or.th

kaamanita@gmail.com