# On the Shoulders of LLMs:
# From LLM Optimization to LLM Agents

**Zuchao Li, Yao Yao, Zhuosheng Zhang**

2024.07.06

IAPR

WUHAN UNIVERSITY

SHANGHAI JIAO TONG UNIVERSITY

# 01

# On the Shoulders of LLMs: Large Language Models

**Zuchao Li**

Wuhan University

zcli-charlie@whu.edu.cn

# Content

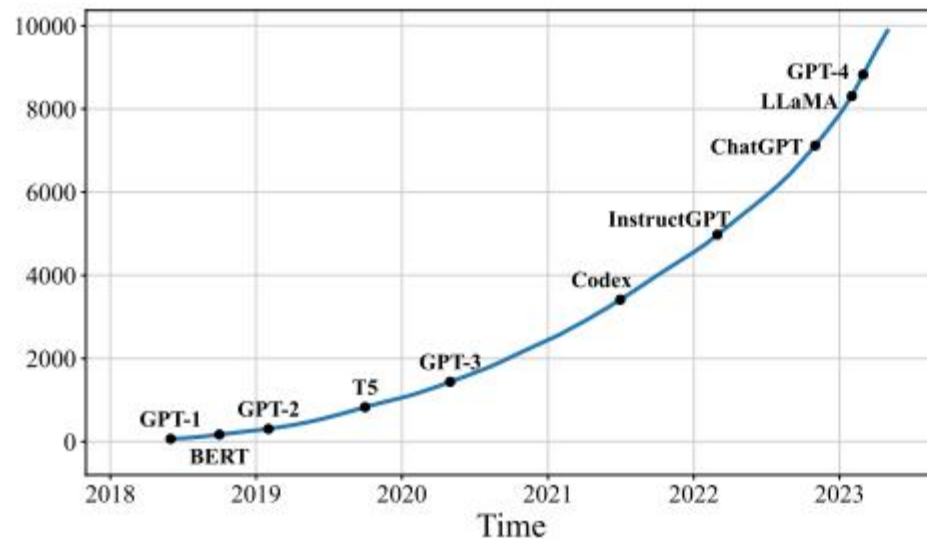3

01

# What is LLM?

# What is LLM?

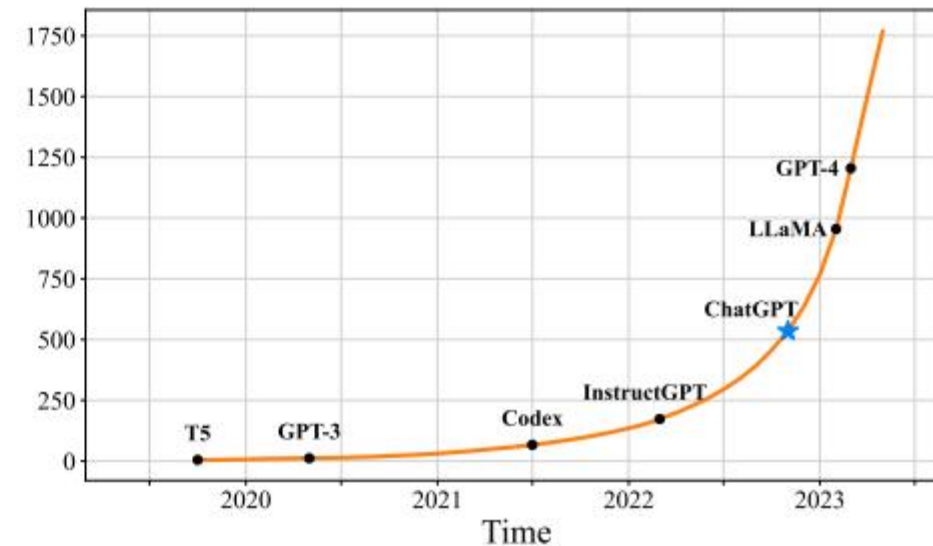> *"The limits of my language mean the limits of my world"*
>
> Ludwig Wittgenstein

➢ **Language** is a prominent ability in human beings to express and communicate, while machines cannot naturally grasp the abilities of understanding and communicating in the form of human language, unless equipped with powerful artificial intelligence (AI) algorithms.

➢ Language modeling (LM) is one of the major approaches to advancing language intelligence of machines.

➢ LM can be divided into four major development stages: **Statistical language models (SLM), Neural language models (NLM), Pre-trained language models (PLM), and Large language models (LLM).**

# Origin of LLM

- ➤ Large-sized PLMs display different behaviors from smaller PLMs and show surprising abilities in solving a series of complex tasks

- ➤ The term "large language models (LLM)" has been coined for these large-sized PLMs, which contain hundreds of billions (or more) of parameters

- ➤ A sharp increase of the arXiv papers that are related to LLMs after the release of ChatGPT
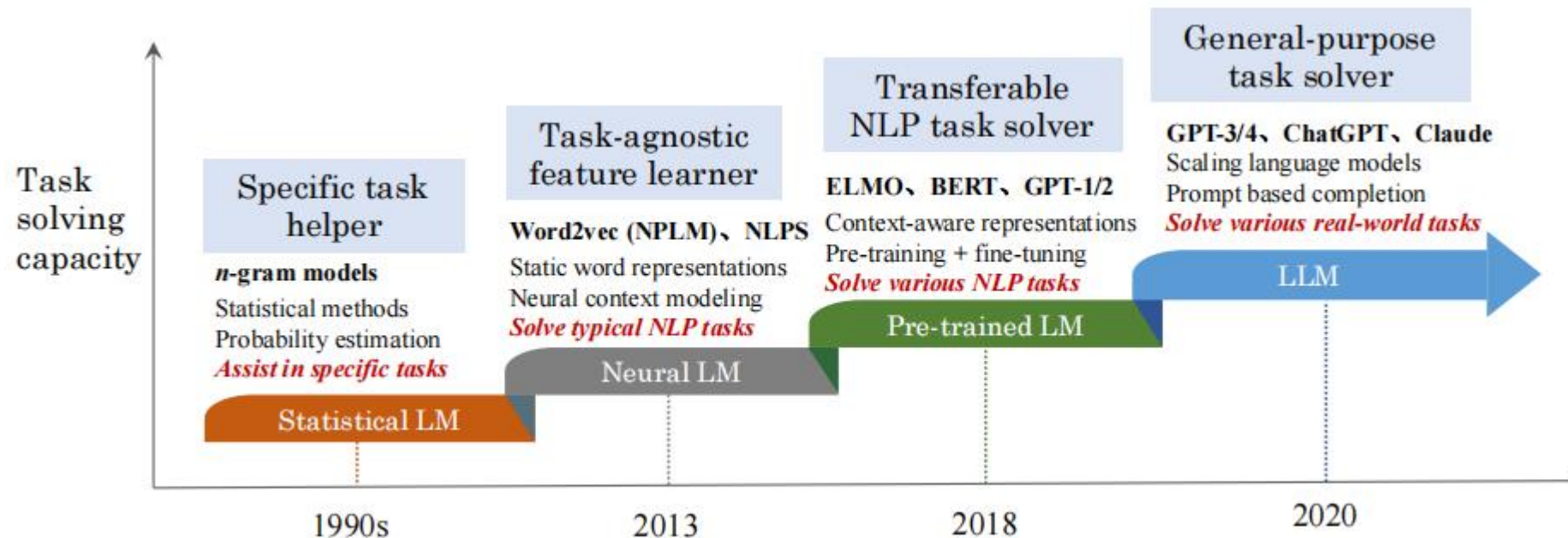


(a) Query="Language Model"
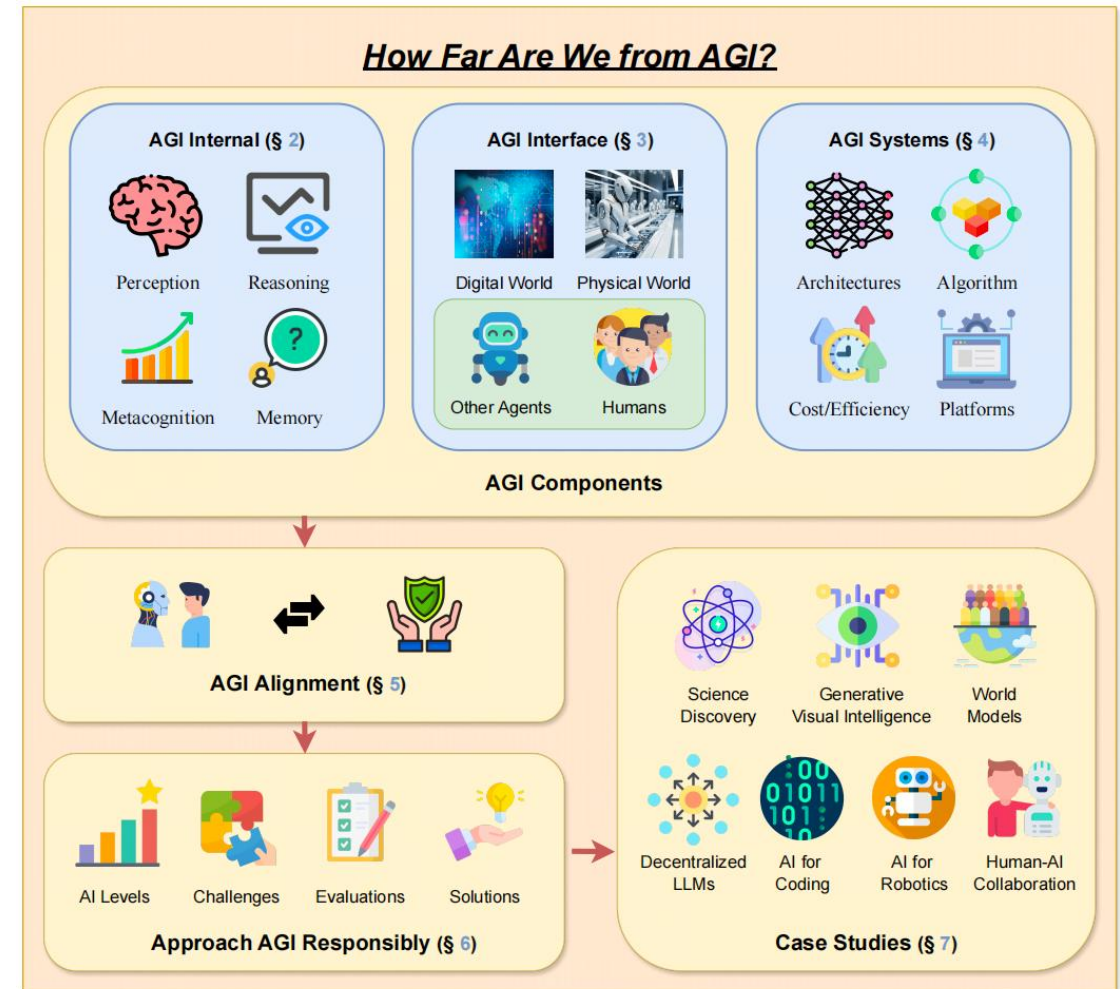
(b) Query="Large Language Model"

# From Language Modeling to Task Solving

➢ LLMs are enhanced by exploring the scaling effect on model capacity, which can be considered as general-purpose task solvers

➢ The task scope that can be solved by LMs have been greatly extended

➢ The task performance attained by LMs have been significantly enhanced



Zhao W X, Zhou K, Li J, et al. A survey of large language models[J]. arXiv preprint arXiv:2303.18223, 2023.

# From LLMs to AGI

➢ The advent of ChatGPT and GPT-4 leads to the rethinking of the possibilities of **artificial general intelligence (AGI)**

➢ The research areas of AI are being revolutionized by the rapid progress of LLMs
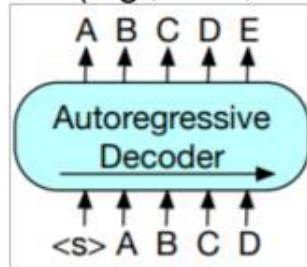
# Early PLMs

> *Encoder-only PLMs*

> *Decoder-only PLMs*

> *Encoder-Decoder PLMs*

- ○ BERT and its variants
- ○ GPT-1 and GPT-2
- ○ T5, mT5, MASS, BART

➢ Compared to PLMs reviewed above, LLMs are not only much larger in model size, but also exhibit stronger language understanding and generation and emergent abilities that are not present in smaller-scale models

➢ GPT-3: viewed as the first LLM, model parameters to 175B

➢ CODEX: a general-purpose programming model

➢ InstructGPT: align language models with user intent on a wide

range of tasks by fine-tuning with human feedback

➢ ChatGPT: superior capacities in communicating with humans

➢ GPT-4: stronger capacities in solving complex tasks than GPT-3.5

➢ GPT-4V, GPT-4 turbo, and beyond:

extensively discussed the assessment

and mitigation of risks related to visually

augmented inputs

# LLaMA Family

- ➢ LLaMA:  using the transformer architecture of GPT-3

- ➢ LLaMA-2: including both foundation language models and Chat models fine-tuned for dialog

# LLaMA Family (cont.)

➤ Alpaca: using 52K instruction-following demonstrations generated in the style of self-instruct using GPT-3.5

➤ Vicuna: fine-tuning LLaMA on user-shared conversations collected from ShareGPT

➤ Guanaco, Koala: instruction-following language model built on LLaMA

➤ Mistral-7B: a 7B-parameter language model engineered for superior performance and efficiency

# PaLM Family

➢ PaLM: a 540B parameter transformer-based LLM

➢ U-PaLM: the model of 8B, 62B, and 540B scales are continually trained on PaLM with UL2R

➢ Flan-PaLM: using a much larger number of tasks, larger model sizes, and chain-of-thought data

➢ PaLM-2: a more compute-efficient LLM with better multilingual and reasoning capabilities, compared to its predecessor PaLM

➢ Med-PaLM: provide high-quality answers to medical questions

➢ Med-PaLM2: improving upon Med-PaLM by over 19%

# Other Families

➢ BLOOM:   A 176B-parameter open-access multilingual

language model

➢ *Claude Family*: LLMs created by Anthropic

➢ Qwen Family: LLMs  created by Alibaba



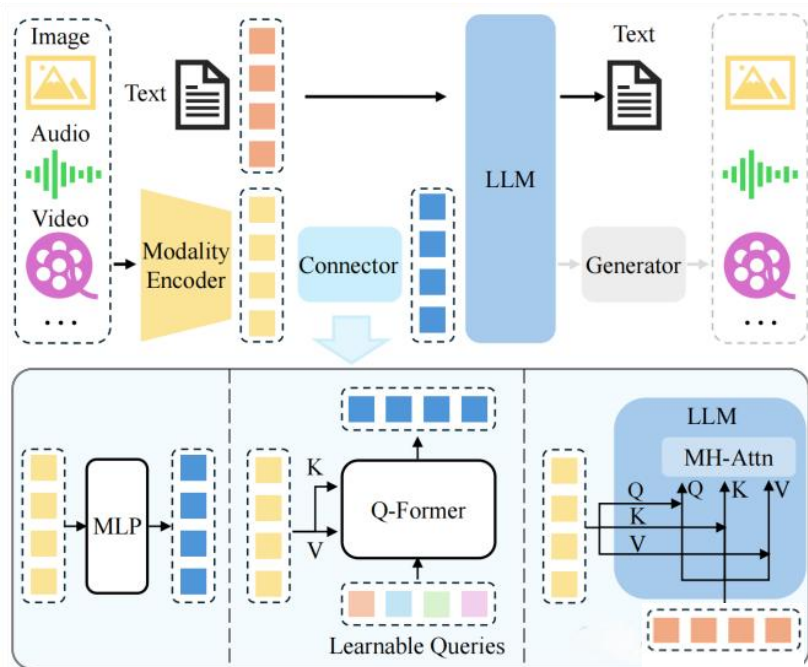| Feature/Model | Claude 3 Opus | Claude 3 Sonnet | Claude 3 Haiku | Claude 2.0 | Claude 2.1 | Claude Instant 1.2 |
|---|---|---|---|---|---|---|
| Description | Most powerful for highly complex tasks | Balanced intelligence and speed for enterprises | Fastest, compact for near-instant responses | Strong performance across various tasks | Improved accuracy and consistency | Fast and efficient, predecessor to Haiku |
| Strengths | Top-level performance, intelligence, fluency | Maximum utility at lower cost, dependable | Quick and accurate targeted performance | Strong general performance | Enhanced accuracy and consistency | Fast and efficient |
| Capabilities | Text generation, Vision, Embeddings | Text generation, Vision, Embeddings | Text generation, Vision, Embeddings | Text generation, Vision, Embeddings | Text generation, Vision, Embeddings | Text generation, Vision, Embeddings |
| API Model Name | claude-3-opus 20240229 | claude-3-sonnet20240229 | claude-3-haiku 20240307 | claude-2.0 | claude-2.1 | claude-instant-1.2 |
| Latency | Moderately fast | Fast | Fastest | Moderate | Moderate | Fast |
| Max Output | 4096 tokens | 4096 tokens | 4096 tokens | 2048 tokens | 2048 tokens | 2048 tokens |
| Multilingual | Yes | Yes | Yes | No | No | No |

Source: marktechpost.com

# Differences between LLMs and PLMs

➢ LLMs display some surprising emergent abilities, are key to the performance of language models on complex tasks, making AI algorithms unprecedently powerful and effective.

➢ LLMs would revolutionize the way that humans develop and use AI algorithms, and the major approach to accessing LLMs is through the prompting interface (e.g., GPT-4 API).

➢ The development of LLMs no longer clearly distinguishes between research and engineering, and researchers have to solve complicated engineering issues, working with engineers or being engineers.

# From LLMs to MLLMs



The architecture of a typical MLLM

➢ **Multimodal Large Language Models (MLLMs)** circumvent the computational cost of training from scratch by effectively leveraging the pre-training knowledge of each modality to enhance multimodal competencies
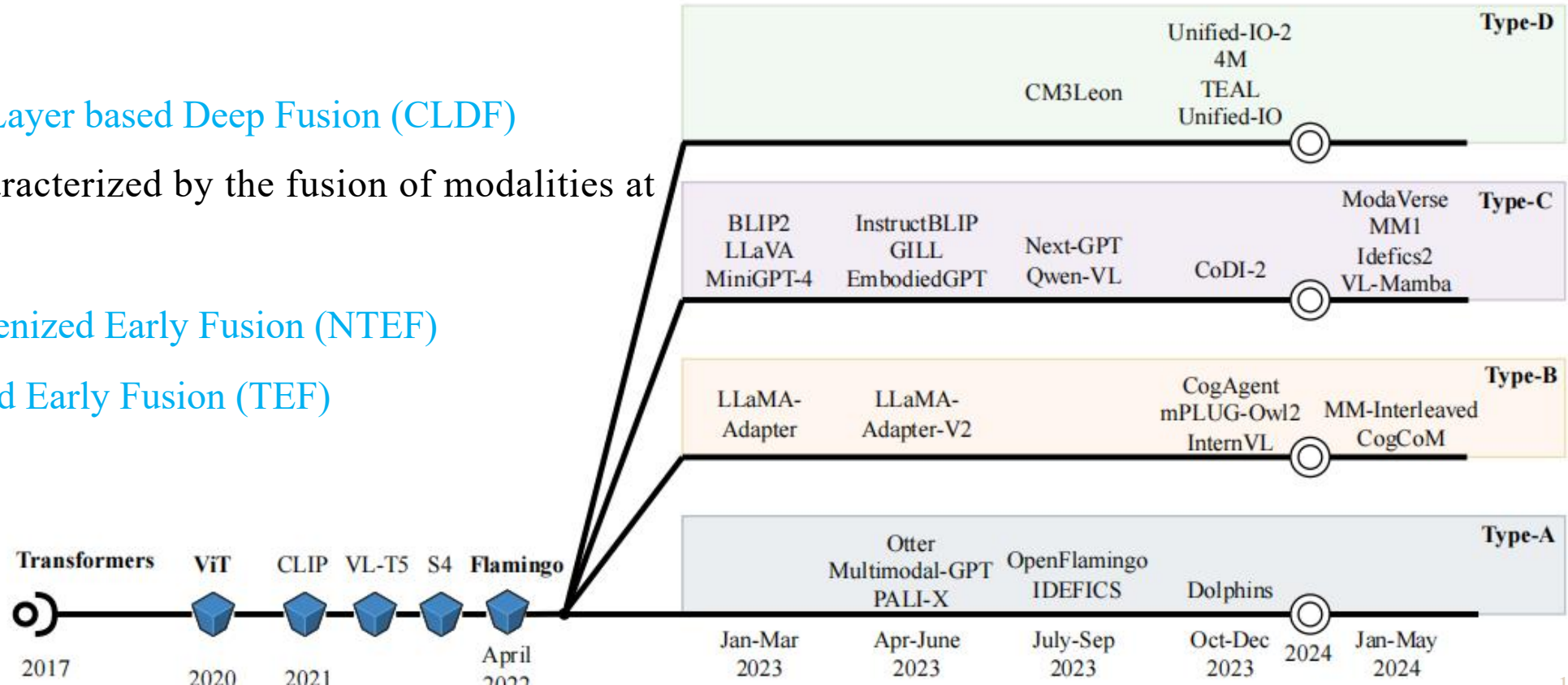
➢ MLLMs can process inputs from multiple modalities, significantly broadening their application scope



The timeline of efficient MLLMs

# LLM-centered Multimodal Model Architectures

➢ **Deep Fusion**, wherein the fusion of modalities occurs within the internal layers of the model

  ○ Type-A: Standard Cross-Attention based Deep Fusion (SCDF)

  ○ Type-B: Custom Layer based Deep Fusion (CLDF)

➢ **Early Fusion**, characterized by the fusion of modalities at the model's input

  ○ Type-C: Non-Tokenized Early Fusion (NTEF)

  ○ Type-D: Tokenized Early Fusion (TEF)

# Standard Cross-Attention based Deep Fusion




The input modalities are deeply fused into the internal layers of the LLM using standard cross-attention layer

- sub-type A.1: the cross-attention can be added either before the self-attention layer
- sub-type A.2: the cross-attention can be added either after the self-attention layer

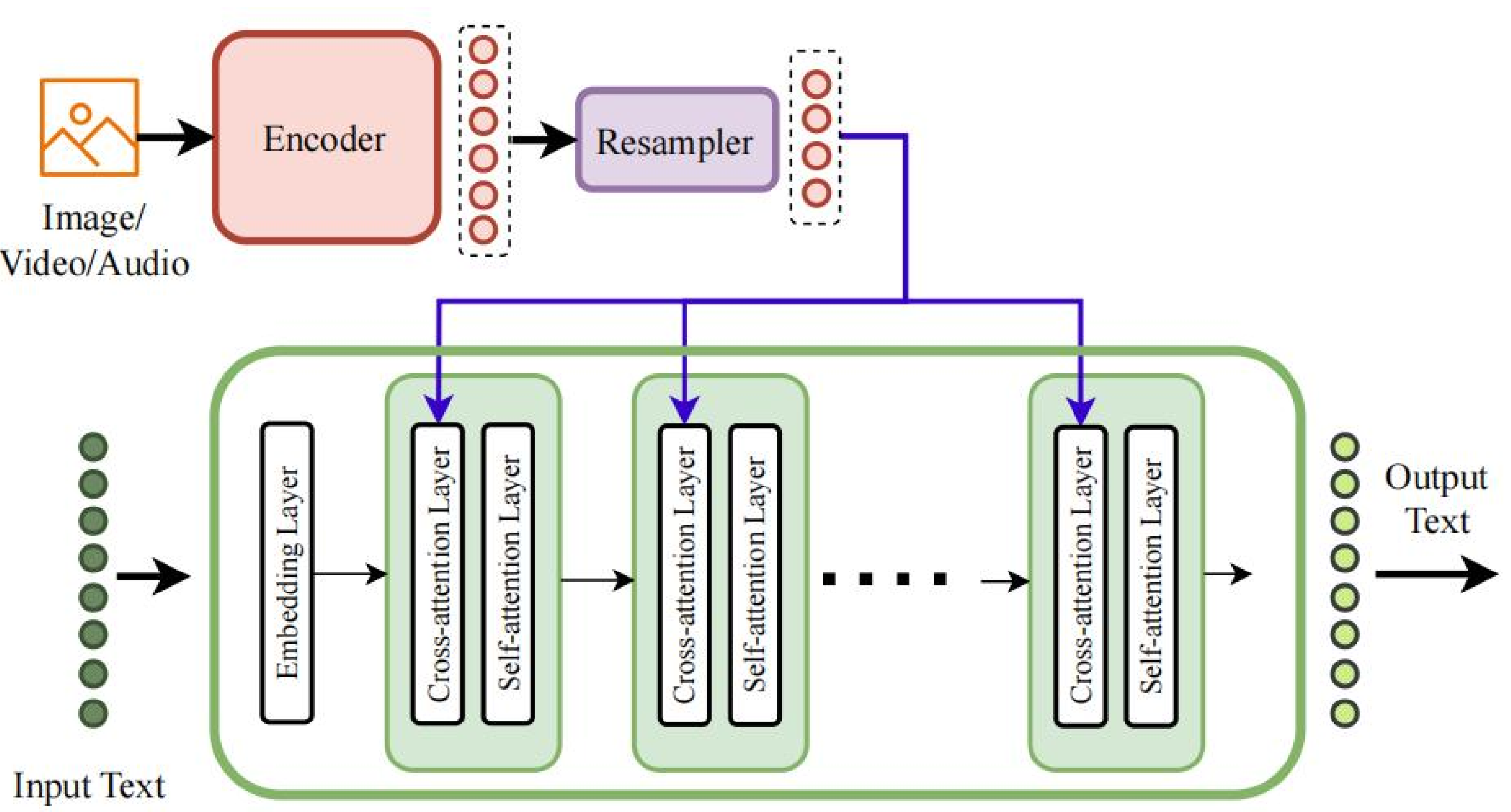

# Custom Layer based Deep Fusion

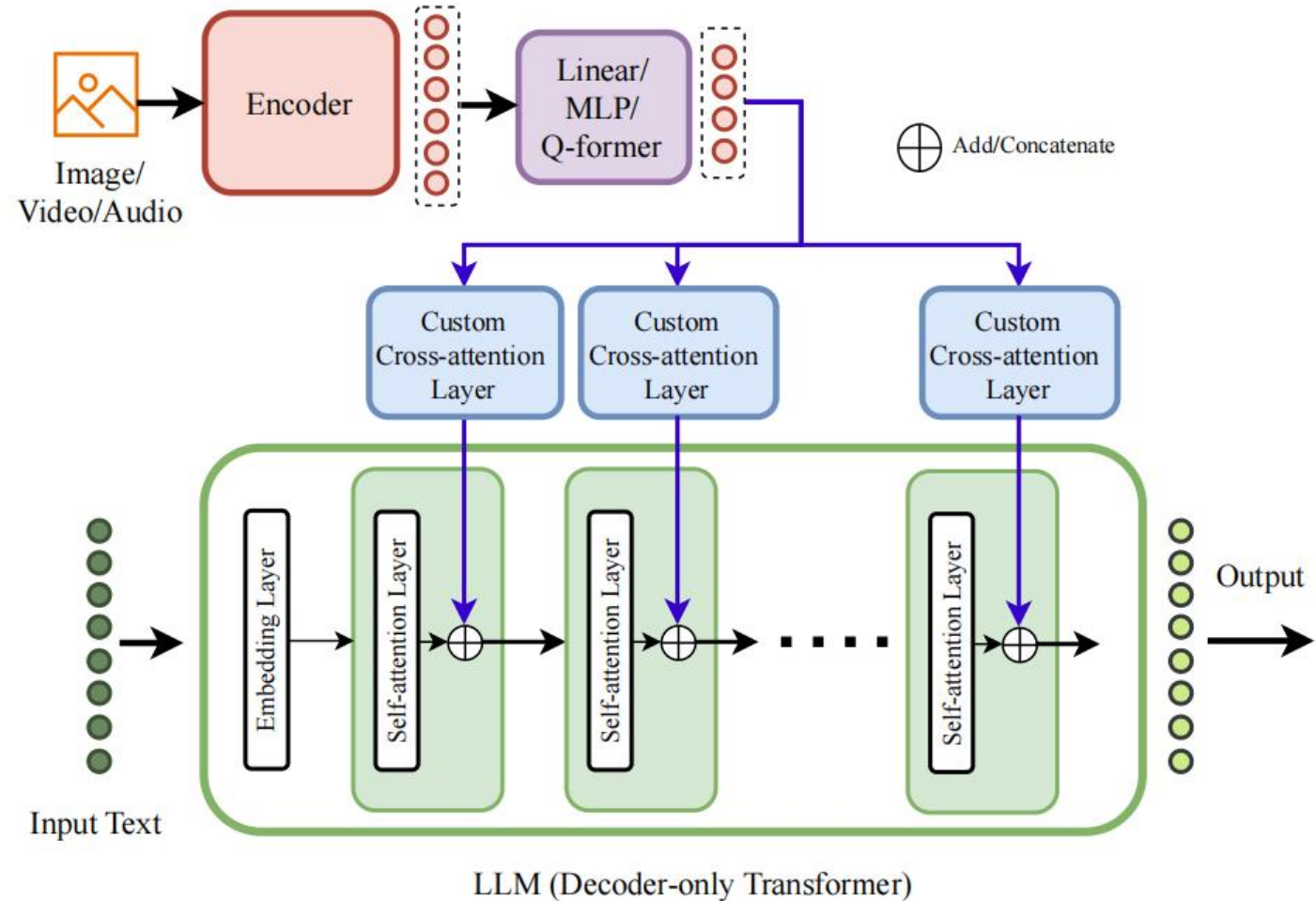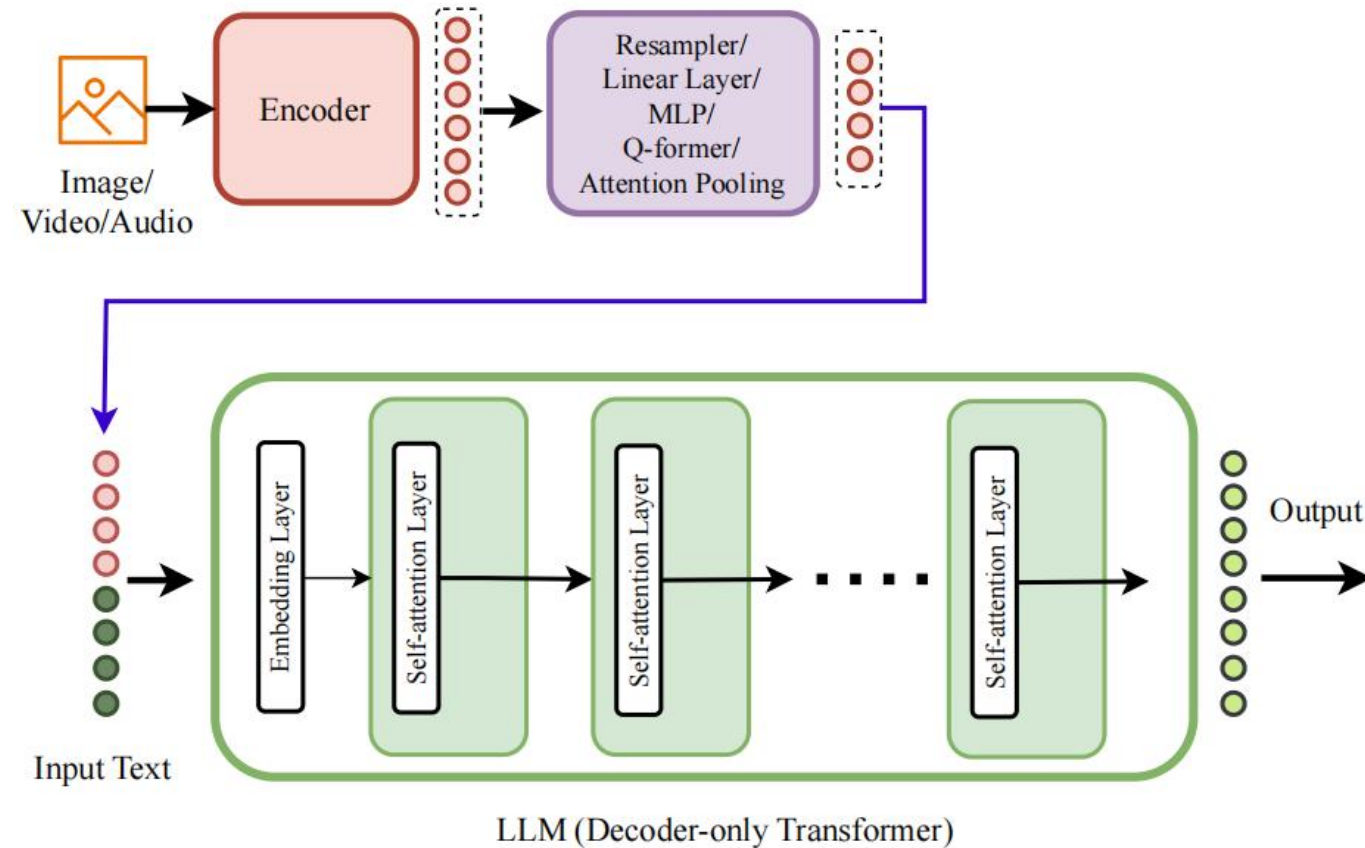➢ The input modalities are deeply fused into the internal layers of the LLM using custom-designed layers

- ○ sub-type B.1: Custom Cross-Attention Layer

- ○ sub-type B.2: Custom Learnable Layer

# Non-Tokenized Early Fusion

The (non-tokenized) input modalities are directly fed to the model at its input, rather than to its internal layers, resulting in early fusion

○ sub-type C.1 Linear Layer/MLP: models using only Linear Layer/MLP for connecting Encoder to the LLM (decoder)

○ sub-type C.2: Q-former and Linear Layer/MLP: models using Q-former and Linear Layer/MLP for connecting Encoder to the LLM (decoder)



○ sub-type C.3: Perceiver Resampler: models using Perceiver resampler for connecting Encoder to the LLM (decoder)

○ sub-type C.4: Custom Learnable layer: models using custom-module/layer for connecting Encoder to the LLM (decod

# Tokenized Early Fusion
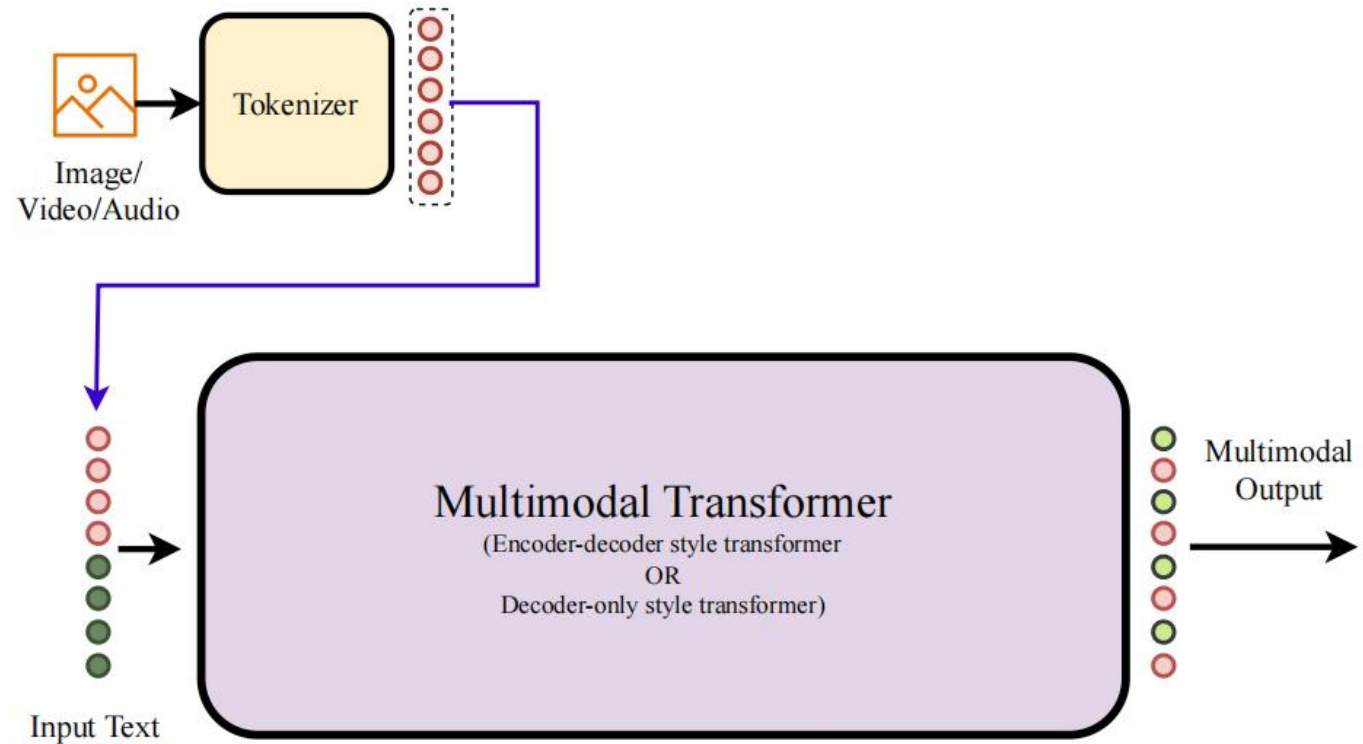
The multimodal inputs are tokenized using a common tokenizer or modality specific tokenizers

- subtype D.1: Models using LLM

Models that primarily use LLM are LaVIT, TEAL, CM3Leon, SEED, Unicode, VL-GPT

- subtype D.2: Models using Encoder-Decoder style Transformer

Models using encoder-decoder style transformer instead of LLM are Unified-IO, Unified-IO 2 and 4M

➤ Any-to-any Multimodal Model：

multimodal models with multimodal-input and multimodal-output

○  Multimodal output generation is one of the primary challenge

○  Type-C and Type-D multimodal architectures are at the forefront of development for any-to-any multimodal models



Any-to-any Multimodal Model development timeline

# 02

# How to utilize LLMs

# Major Aspects for LLMs Optimization

**01** *Pre-Training* -- how to pre-train a capable LLM

**02** *Adaptation* -- how to effectively adapt pre-trained LLMs for better

**03** *Utilization* -- how to use LLMs for solving various downstream

**04** *Capability evaluation* -- how to evaluate the abilities of LLMs and existing empirical findings

# Pre-training

Too Expensive!

Goal: acquire essential language understanding and generation skills

**VS**

Key Elements: model architectures, acceleration methods, and optimization techniques

- High-quality data is vital to model capacities of LLMs
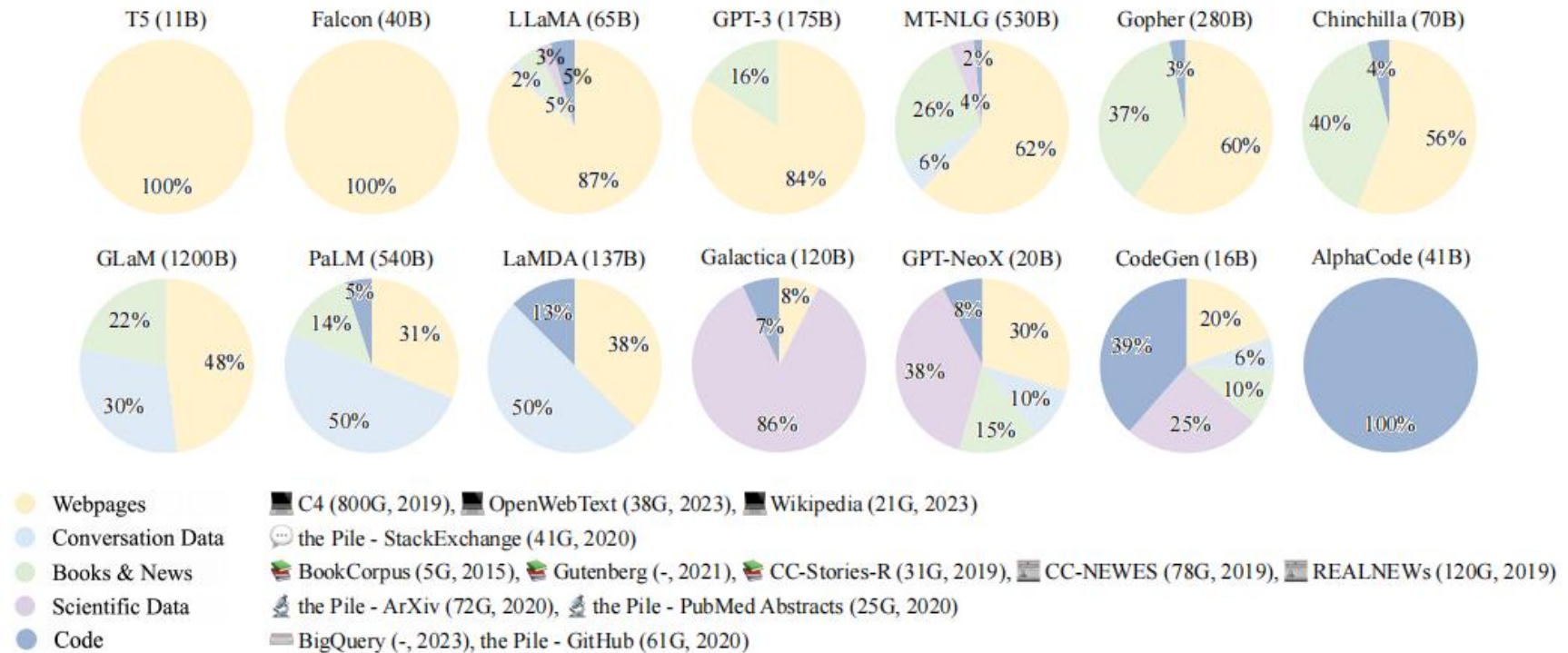
  ➢ Data Source

  ○ General Text Data

    • Webpages

    • Conversation text

    • Books

  ○ Specialized Text Data

    • Conversation text

    • Scientific text

    • Code



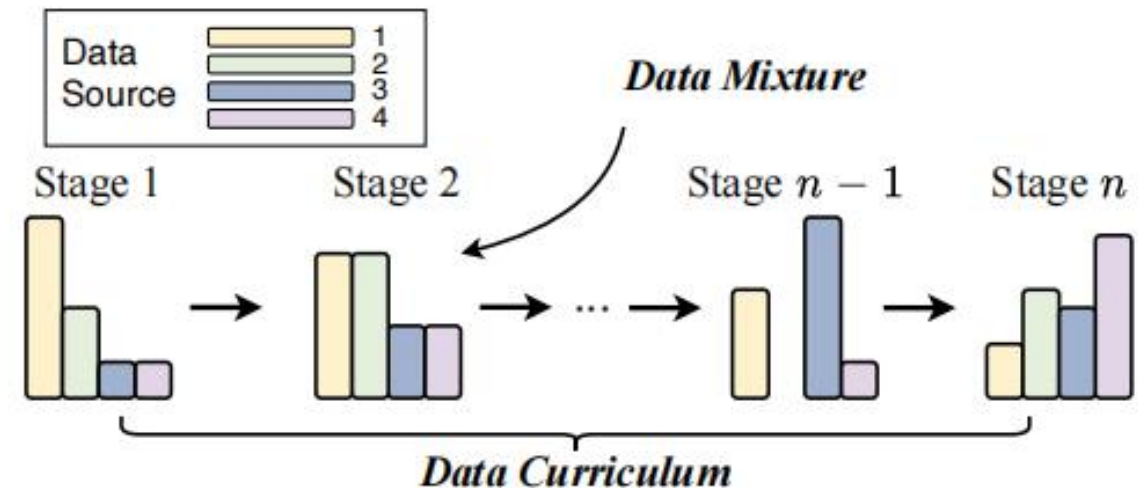Ratios of various data sources in the pre-training data for existing LLMs

# Data Preparing for Pretraining

➢ Data Scheduling

    ○ Data Mixture

        • *Increasing the diversity of data sources*

        • *Optimizing data mixtures*

        • *Specializing the targeted abilities*

    ○ Data Curriculum

        • aims to organize different parts of pre-

training data for LLMs in a specific order

➢ Typical Architectures

○ **Encoder-decoder Architecture** consists

of two stacks of Transformer blocks

○ **Causal Decoder Architecture**

incorporates the unidirectional attention mask

○ **Prefix Decoder Architecture**

incorporates the unidirectional attention masks

# New Architectures for Pretraining

➢ Extended Architectures

○ **Mixture-of-Experts (MoE)** is a flexible way

to scale up the model parameter



○ **Emergent Architectures**

• new architectures: parameterized state space models,  long convolutions,  and Transformer-like architectures

• *decoding process more efficient*

• *models to be trained in a highly parallel and efficient manner*

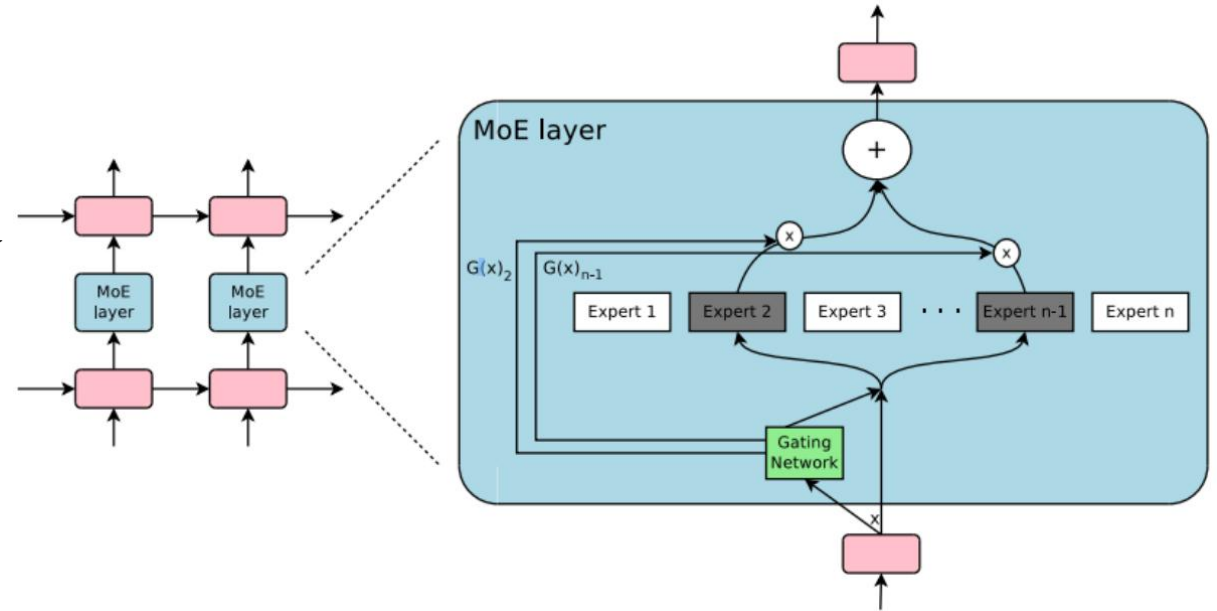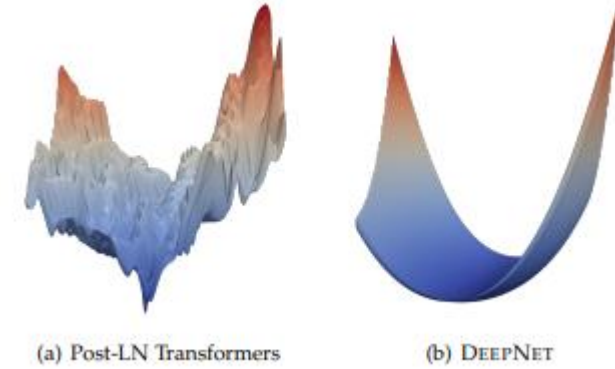# Pretraining Change: Normalization Position

➢ Normalization Position

  ○ To alleviate the issue of training instability

  ○ techniques: *LayerNorm, RMSNorm, DeepNorm*



(a) Post-LN Transformers     (b) DEEPNET

The loss surface of 36-layer vanilla Post-LN and DEEPNET at the early stage of training.

| | Weight matrix re-scaling | Weight matrix re-centering | Weight vector re-scaling | Dataset re-scaling | Dataset re-centering | Single training case re-scaling |
|---|---|---|---|---|---|---|
| BatchNorm | ✓ | ✗ | ✓ | ✓ | ✓ | ✗ |
| WeightNorm | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ |
| LayerNorm | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ |
| RMSNorm | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ |
| ρRMSNorm | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ |

➢ Normalization Position

  ○ *Post-LN*: being instable

  ○ *Pre-LN*: being more stable in training but

performing worse than variants with Post-LN

  ○ *Sandwich-LN*: to avoid the value explosion

issues in Transformer layer outputs

relative PE

➢ Position Information

- ○ *Absolute position embedding*

- ○ *Relative position embedding*

- ○ *Rotary Position Embedding (RoPE)*

- ○ *ALiBi*



ALiBi



RoPE

# Pretraining Change: Attention

- *Full attention*
- *Sparse attention*
- *Multi-query / grouped-query attention*
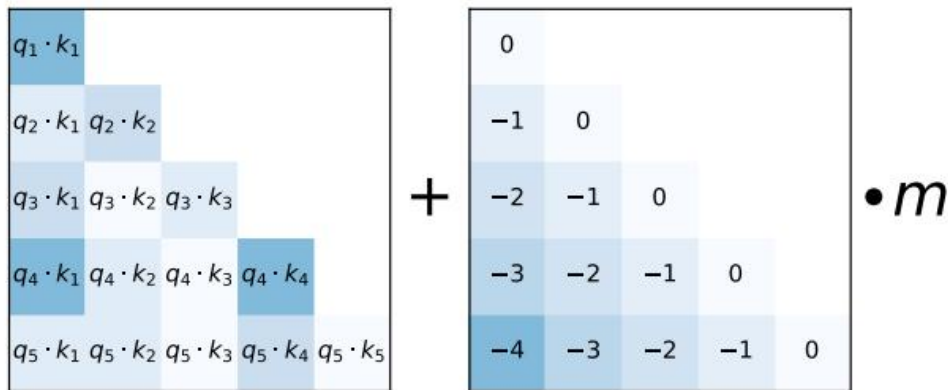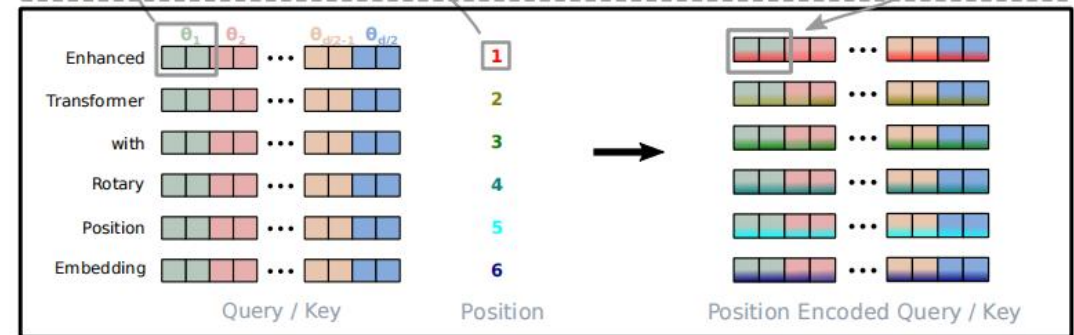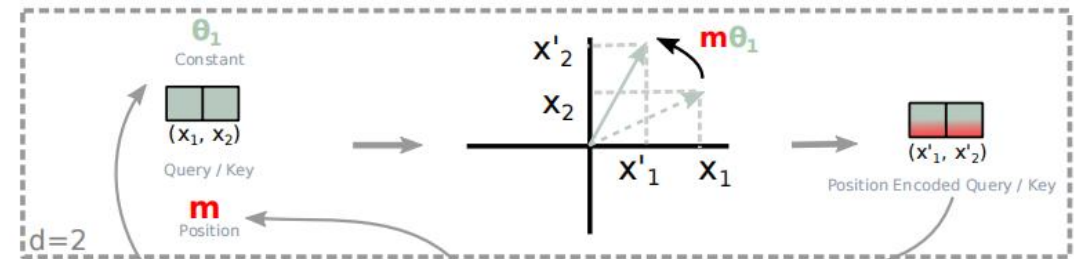- *FlashAttention*
- *PagedAttention*



PagedAttention

# Pretraining Change: In-Context Modeling

- Language Modeling (LM)

- Denoising Autoencoding (DAE)

- Mixture-of-Denoisers (MoD)



**Explaining a Joke**

**Input:** Did you see that Google just hired an eloquent whale for their TPU team? It showed them how to communicate between two different pods!
**Model Output:** TPUs are a type of computer chip that Google uses for deep learning. A "pod" is a group of TPUs. A "pod" is also a group of whales. The joke is that the whale is able to communicate between two groups of whales, but the speaker is pretending that the whale is able to communicate between two groups of TPUs.

**Logical Inference**

**Input:** Shelley is from Virginia, but is visiting that city with that famous market where they throw the fish. Going home next Tuesday!
**Question:** Is it likely that Shelley will be near the Pacific Ocean this weekend?
**Model Output:** The city with the famous market where they throw the fish is Seattle, Washington. Seattle is on the Pacific Ocean. Shelley is visiting Seattle, so she will be near the Pacific Ocean this weekend. The answer is "yes", it is likely that Shelley will be near the Pacific Ocean this weekend.

Examples of LM (PaLM)



DAE (T5)



MoD for training UL2

# Pretraining Change: Long Context Understanding

> Scaling Position Embeddings

  ○ *Direct model fine-tuning* --  to directly fine-tune the models on long texts with the desired length

  ○ *Position interpolation* -- downscales the position indices within the original context window

  ○ *Position truncation* -- to mitigate the challenges posed by out-of-distribution rotation angles

  ○ *Base modification* -- decreasing the basis can be achieved by increasing the value of the base

  ○ *Basis truncation* -- dealing with the singular dimensions with wavelengths exceeding the training length

> ➤ Adapting Context Window
> ● Parallel context window
> ● Λ-shaped context window
> ● External memory



(a) Proposed Solution: LM-Infinite

(b) A Conceptual Model of Relative Positional Attention



External memory (Unlimiformer)

➢ Background

    ○ *greedy search*

    ○ *sampling-based methods*

➢ Improvement for Greedy Search

    ○ *Beam search*

    ○ *Length penalty*

➢ Improvement for Random Sampling

    ○ *Temperature sampling*

    ○ *Top-k sampling*

    ○ *Top-p sampling*

    ○ *η-sampling*



copious repetition (highlighted in blue)
incoherence (highlighted in red)

➢ Decoding Efficiency Issues

○ the *prefill* stage

○ the *incremental decoding* stage

• Reducing data transfer ---- optimizing GPU memory access

△ KV cache、Flash-Decoding、PagedAttention、MQA、GQA

• Decoding strategy optimization ---- improve the sequential nature of the auto regressive generation manner

△ speculative decoding



(a) Incremental decoding.

(b) Timeline Comparison.

# Pretraining Change: Scalable Training Techniques

> Primary technical issues

- ○ increasing training throughput

- ○ loading larger models into GPU memory

> Approaches

- ○ 3D Parallelism

- ○ ZeRO

- ○ Mixed Precision Training



3D Parallelism



Mixed Precision Training



ZeRO

# Adaptation: Instruction Tuning



General pipeline of instruction tuning.

Formatting NLP Task Datasets, Formatting Daily Chat Data, Formatting Synthetic Data



(a) Formatting Task Datasets  (b) Formatting Daily Chat Data  (c) Formatting Synthetic Data

## Pipeline Distillation from ChatGPT



General pipeline of distillation for synthetic data generation

General pipeline of self-improvement for synthetic data generation

## ➢ Some findings from our practice

- ✓ **Task-formatted instructions** are more proper for the QA setting, but may not be useful for the chat setting

- ✓ **A mixture of different kinds of instructions** are helpful to improve the comprehensive abilities of LLMs

- ✓ **Enhancing the complexity and diversity of instructions** leads to an improved model performance

- ✓ **Simply increasing the number of instructions** may not be that useful, and balancing the difficulty is not always helpful

- ✓ **A larger model scale** leads to a better instruction following performance



The process of plant photosynthesis is commonly written as:
$6CO2 + 6H2O \rightarrow C6H12O6 + 6O2$
Please explain the main role of chlorophyll in above formula.

Please fill in the table below with the approximate values of the speed of light in each medium.

| Medium | Speed of light (km/s) |
|--------|----------------------|
| Air | |
| Water | |
| Glass | |

In-Breadth Evolving

Complicate Input (Table)

```
import math
import random

# choose a random integer between 1 and 10
x = random.randint(1, 10)
1/(math.sqrt(x) + x^2) =?
```

How many times faster is light than sound in a vacuum?

How is the speed of light in a vacuum measured and defined?

Complicate Input (Code)

Increase Reasoning

Deepening

$1/(sqrt(2) + 4^2) = ?$

What is the speed of light in a vacuum?

Complicate Input (Formula)

In-Breadth Evolving

How to prove 1 + 1 = 2 in the Goldbach Conjecture?

If you have one apple and someone gives you another banana, how many fruits do you have?

Add Constraints

Concretizing

In what situation does 1+1 not equal to 2?

Deepening

$1 + 1 = ?$

Increase Reasoning

What is the value of x, if $x^3 + 2x + 3=7$?

Initial Instruction

Examples of Enhancing the instruction complexity (Evol-Instruct)

➤ Parameter-Efficient Fine-Tuning (PEFT) Methods

**a** Adapter Tuning — **b** Prefix Tuning — **c** Prompt Tuning — **d** Low-Rank Adaption



(a) Adapter Tuning    (b) Prefix Tuning    (c) Prompt Tuning    (d) Low-Rank Adapation

# Adaptation: Parameter-Efficient Fintuning

➤ Adapter Tuning: incorporate small neural network modules (i.e., adapter) into the Transformer models

  ○ bottleneck architecture, parallel adapters



(a) Prefix-Tuning

(b) LoRA    (c) Series Adapter    (d) Parallel Adapter



Adapter    Parallel Adapter    Scaled PA



bottleneck architecture(Long et al. 2022)

➢ Prefix Tuning: prepends a sequence of prefixes

(virtual token embeddings) to each Transformer

layer

  ○ Prefix-tuning, P-tuning v2



(a) Lester et al. & P-tuning (Frozen, 10-billion-scale, simple tasks)

(b) P-tuning v2 (Frozen, most scales, most tasks)

P-tuning vs. P-tuning v2



Fine-tuning vs. Prefix-tuning

➤ Prompt Tuning: incorporate trainable prompt vectors

at the input layer

○ discrete prompting methods, prompt tuning, P-tuning,



**Model Tuning**

Pre-trained Model (11B params)

Task A Batch → Task A Model (11B params)

Task B Batch → Task B Model (11B params)

Task C Batch → Task C Model (11B params)

**Prompt Tuning**

Mixed-task Batch

Task Prompts (20K params each)

Pre-trained Model (11B params)



(a) Discrete Prompt Search

(b) P-tuning

➤  Low-Rank Adaptation (LoRA): impose the low-rank

constraint for approximating the update matrix at

each dense layer

○  LoRA, DyLoRA, DyLoRA



LoRA

An example of the three steps of RLHF System (Long et al. 2022)

Key Steps for RLHF

- Supervised fine-tuning

- Reward model training

- RL fine-tuning

➢ Keep tedious and memory consuming

➢ RLHF is rather complex and often sensitive to hyper-parameters

➤ Alignment Data Collection

  ○ *Reward model based approaches* (RAFT, Quark, ILF)



Quark



ILF

- ➢ Alignment Data Collection

  - ○ *LLM based generative approaches* (CAI, Self-Align, FIGA)



CAI



Self-Align



FIGA

➢ Supervised Alignment Tuning

○ *Primary training objective* -- the traditional cross-entropy loss for sequence-to-sequence learning

• CoH, Quark, DPO, FIGA



CoH

DPO

# Adaptation: Memory Efficient Adaptation

➢ **What is quantization:** the mapping process from floating-point numbers to integers, especially *INT8 quantization*

○ What to Quantize: *weights* (model parameters) and *activations* (hidden activations)

○ Quantization Operator: *Uniform Quantization, non-uniform quantization* (whether quantized values (aka quantization levels) are uniformly spaced)



Comparison between uniform quantization (left) and non-uniform quantization (right)

○  How to choose  the scaling factor  in uniform quantization

-- *Symmetric Quantization*:  partitions the clipping using a symmetric range,  easier implementation, but it is

sub-optimal for cases where the range could be skewed and not symmetric

-- *Asymmetric Quantization*



Illustration of symmetric quantization and asymmetric quantization

56

○ Quantization Granularity

 -- *Layerwise Quantization*:  the clipping range is determined by considering all of the weights in convolutional filters of a layer

 -- *Asymmetric Quantization:* the clipping range is determined by considering all of the weights in convolutional filters of a layer



Illustration of different quantization granularities

# Adaptation: Memory Efficient Adaptation

➤ Range Calibration Algorithms: *Static Quantization* vs. *Dynamic Quantization*

-- two approaches to quantizing activations

## Static Quantization

● the clipping range is pre-calculated and *static* during inference

● does not add any computational overhead but results in lower accuracy

## Dynamic Quantization

● this range is *dynamically* calculated for each activation map during runtime

● has a very high overhead but results in higher accuracy

# Adaptation: Memory Efficient Adaptation

➢ *quantization-aware training (QAT)* -- requiring additional full model retraining

➢ *Post-Training Quantization (PTQ)* -- requiring no model retraining

　○　PTQ methods keep a much lower computational cost than QAT methods



Comparison between QAT ( Left) and PTQ (Right)

# Quantization-aware Training

➢ QAT -- the model parameters are quantized after each gradient update

   ◦ categories

   *-- Straight Through Estimator (STE) methods*

    *-- Non-STE methods*

   ◦ disadvantage

-- the computational cost of re-training the

NN model



Illustration of QAT procedure, including the use of STE

# Post-Training Quantization

➢ **Mixed-Precision Decomposition** -- to recover the outliers in hidden activations

➢ **Fine-Grained Quantization** -- to reduce the quantization error

➢ **Balancing the Quantization Difficulty** -- to consider weights being easier to be quantized than activations

➢ **Layerwise Quantization** -- to find optimal quantized weights that minimize a layerwise reconstruction loss

# Other Quantization Methods

QLoRA:  *Efficient fine-tuning*

*enhanced quantization*

overcome this challenge, which directs

low-bit quantization (e.g., INT4

quantization) often results in large

performance degradation

# Experience for Memory-Efficient Model Adaptation

- ✓ **INT8 weight quantization** can often yield very good results on LLMs, while the performance of lower precision weight

- ✓ **Activations** are more difficult to be quantized than weights

- ✓ **Efficient fine-tuning enhanced quantization** is a good option to enhance the performance of quantized LLMs

# LLM Utilization

➢ How to use LLMs -- design suitable *prompting* strategies for solving various

tasks

○ prompting methods

-- *in-context learning*

-- *chain-of-thought prompting*

-- *planning*

# Prompt Engineering

➢ **Prompt Engineering** -- the process of manually creating a suitable prompt

➢ Key Ingredients: Task description, Input data, Contextual information, Prompt style

---

Use the provided articles delimited by triple quotes to answer questions. If the answer cannot be found in the articles, write "I could not find an answer."

**Articles:** """Joao Moutinho is a Portuguese footballer who last played as a central midfielder for Premier League club Wolverhampton Wanderers and the Portugal national team."""

**Question:** Is the following sentence plausible? 'Joao Moutinho was out at third.'

**Answer:** Let's think step by step. Joao Moutinho is a soccer player. Being out at third is part of baseball, not soccer.  So the answer is No.

...

<Demonstrations>

**Articles:** <insert articles, each delimited by triple quotes>
**Question:** <insert question>
**Answer:**

---

Prepare a meta-review by answering the following questions from the reviewer comments (provided after the questions).

1. Based on the reviewer's comments, what are the core contributions made by this manuscript?
2. What are the common strengths of this work, as mentioned by multiple reviewers?
3. What are the common weaknesses of this work, as highlighted by multiple reviewers?
4. What suggestions would you provide for improving this paper?
5. What are the missing references mentioned by the individual reviews?

**The review texts are below:** <insert three comments $R_1$, $R_2$, $R_3$ from the reviewers>
**Meta-review:** <insert meta-review>
...
<Demonstrations>

Provide justification for your response in detail by explaining why you made the choices you actually made. A good output should be coherent, highlight major strengths/issues mentioned by multiple reviewers, be less than 400 words in length, and finally, the response should be in English only.

**The review texts are below:** <insert three comments $R_1$, $R_2$, $R_3$ from the reviewers>
**Meta-review:**

---

CREATE TABLE Highschooler (
ID int primary key,
name text,
grade int
);
/*
3 example rows:
SELECT * FROM Highschooler LIMIT 3;
ID    name    grade
1234    Janie    8
5678    Mary    8
9012    Mike    9
*/

Using valid SQLite, answer the following questions for the tables provided above.
**Question:** What is Kyle's id?
**SQL:** SELECT ID FROM Highschooler WHERE name="Kyle";
...
<Demonstrations>

**Question:** <insert question>
**SQL:**

---

Example instructions. The blue text denotes the task description, the red text denotes the contextual information, the green text denotes the demonstrations, and the gold text denotes the prompt style.

➤ **Prompt Design Principles**

① *Expressing the task goal clearly*

② *Decomposing into easy, detailed sub-tasks*

③ *Providing few-shot demonstrations*

④ *Utilizing model-friendly format*

| Ingredient | Collected Prompts | Prin. |
|---|---|---|
| Task Description | T1. Make your prompt **as detailed as possible**, *e.g.*, *"Summarize the article into a short paragraph within 50 words. The major storyline and conclusion should be included, and the unimportant details can be omitted."* | ① |
| | T2. It is helpful to let the LLM know that it is **an expert with a prefixed prompt**, *e.g.*, *"You are a sophisticated expert in the domain of compute science."* | ① |
| | T3. Tell the model **more what it should do**, but not what it should not do. | ① |
| | T4. To avoid the LLM to generate too long output, you can just use the prompt: *"Question: Short Answer: "*. Besides, you can also use the following suffixes, *"in a or a few words"*, *"in one of two sentences"*. | ① |
| Input Data | I1. For the question required factual knowledge, it is useful to first **retrieve relevant documents** via the search engine, and then **concatenate them into the prompt** as reference. | ④ |
| | I2. To highlight some important parts in your prompt, please **use special marks**, *e.g.*, *quotation (""") and line break (\n)*. You can also use both of them for emphasizing. | ④ |
| Contextual Information | C1. For complex tasks, you can **clearly describe the required intermediate steps** to accomplish it, *e.g.*, *"Please answer the question step by step as: Step 1 - Decompose the question into several sub-questions, ···"* | ② |
| | C2. If you want LLMs to provide the score for a text, it is necessary to provide a **detailed description about the scoring standard** with examples as reference. | ① |
| | C3. When LLMs generate text according to some context (*e.g.*, making recommendations according to purchase history), instructing them with **the explanation about the generated result** conditioned on context is helpful to improve the quality of the generated text. | ② |
| | C4. An approach similar to **tree-of-thoughts** but can be **done in one prompt**: *e.g.*, *Imagine three different experts are answering this question. All experts will write down one step of their thinking, then share it with the group of experts. Then all experts will go on to the next step, etc. If any expert realizes they're wrong at any point then they leave. The question is* | ② |
| Demonstration | D1. **Well-formatted in-context exemplars** are very useful, especially for producing the outputs with complex formats. | ③ |
| | D2. For few-shot chain-of-thought prompting, you can also use the prompt *"Let's think step-by-step"*, and the few-shot examples should be **separated by "\n"** instead of full stop. | ①③ |
| | D3. You can also **retrieve similar examples** in context to supply the useful task-specific knowledge for LLMs. To retrieve more relevant examples, it is useful to **first obtain the answer** of the question, and then concatenate it with the question for retrieval. | ③④ |

Examples of useful tips

# Prompt Engineering (cont.)

➢ **Some Experience**

- ✓ Carefully designed prompts can boost the zero-shot or few-shot performance

- ✓ More complex tasks can benefit more from careful prompt engineering

- ✓ For mathematical reasoning tasks, it is more effective to design specific prompts based on the format of programming language

- ✓ Through suitable prompt engineering, LLMs can handle some non-traditional NLP tasks

# Prompt Optimization

➢ Discrete Prompt Optimization: the form is simple and flexible, but it has the combinatorial huge search space

- *Gradient-based approaches -- to* maximize the output likelihood via gradient update

- *RL-based approaches --* to formulate the discrete promptoptimization as RL problem

- *Edit-based approaches --* to directly edit existing prompts based on the task performance

- *LLM-based approaches --* to directly leverage LLMs as prompt generator

➢ Continuous Prompt Optimization: can be directly optimized through the

gradient update based on the loss of downstream tasks

- *Prompt learning with sufficient data* -- leverage supervised learning to

optimize the continuous prompts by minimizing the cross-entropy loss

based on sufficient downstream task data

- *Prompt transferring with scarce data* -- to work well in data-scarce

domains and tasks

# In-Context Learning

➢ **In-Context Learning (ICL)** requires a formatted prompt context containing the task description and/or a few task examples as demonstrations written in natural language templates. Taking this prompt and a query as the input, LLMs are responsible for making predictions.



An example of in-context learning



An example of in-context learning

➢ Based on task demonstrations, LLMs can recognize and perform a new task without explicit gradient update



A comparative illustration of ICL and chain-of-thought (CoT) prompting

# Demonstration Design

➤ Major Aspects

-- *Demonstration Selection* -- to select a subset of examples that can effectively leverage the

ICL capability of LLMs

  ○ *Heuristic approaches* -- simplicity and low costs

  ○ *LLM-based approaches* -- making use of LLMs

➤ major aspects

-- *Demonstration Format* -- to integrate and format selected task examples into a natural language prompt for LLMs

○ straightforward method -- to instantiate a pre-defined template with the corresponding input-output pairs

○ *how to automatically generate high-quality ones* -- Auto-CoT, least-to-most prompting



**Auto Demos One by One**

Q: While shopping for music online, Zoe bought 3 country albums and 5 pop albums. Each album came with a lyric sheet and had 3 songs. How many songs did Zoe buy total?
A: Let's think step by step. Zoe bought 3 country albums. Each album has 3 songs. So she bought 3*3=9 songs from the country albums. Zoe bought 5 pop albums. Each album has 3 songs. So she bought 5*3=15 songs from the pop albums. Zoe bought 9+15=24 songs in total. The answer is 24.
...
Q: A chef needs to cook 9 potatoes. He has already cooked 7. If each potato takes 3 minutes to cook, how long will it take him to cook the rest?
A: Let's think step by step. The chef has already cooked 7 potatoes. That means it has taken him 7 * 3 minutes to cook those 7 potatoes. That means it will take him 3 more minutes to cook each of the remaining 2 potatoes ...

Q: A pet store had 64 puppies. In one day they sold 28 of them and put the rest into cages with 4 in each cage. How many cages did they use?
A: Let's think step by step.

**Test Question**

**In-Context Reasoning**

The pet store had 64 puppies. They sold 28 of them. That means they have 36 puppies left. They put the rest into cages with 4 in each cage. That means they have 9 cages. The answer is 9.

Auto-CoT -- demonstrations (on the right) are automatically constructed one by one (total: k) using an LLM with the "Let's think step by step" prompt

**Stage 1: Decompose Question into Subquestions**

Q: It takes Amy 4 minutes to climb to the top of a slide. It takes her 1 minute to slide down. The water slide closes in 15 minutes. How many times can she slide before it closes?

A: To solve "How many times can she slide before it closes?", we need to first solve: "How long does each trip take?"

**Stage 2: Sequentially Solve Subquestions**

It takes Amy 4 minutes to climb to the top of a slide. It takes her 1 minute to slide down. The slide closes in 15 minutes.
Subquestion 1 — Q: How long does each trip take?

A: It takes Amy 4 minutes to climb and 1 minute to slide down. 4 + 1 = 5. So each trip takes 5 minutes.

It takes Amy 4 minutes to climb to the top of a slide. It takes her 1 minute to slide down. The slide closes in 15 minutes.
Q: How long does each trip take?
Append model answer to Subquestion 1 — A: It takes Amy 4 minutes to climb and 1 minute to slide down. 4 + 1 = 5. So each trip takes 5 minutes.
Subquestion 2 — Q: How many times can she slide before it closes?

A: The water slide closes in 15 minutes. Each trip takes 5 minutes. So Amy can slide 15 ÷ 5 = 3 times before it closes.

least-to-most prompting solving a math word problem in two stages

# Demonstration Design (cont.)

➤ *Demonstration Order* -- to alleviate the recency bias, (i.e., repeat answers that are near

the end of demonstrations)

- ○ several heuristic methods

- ○ to integrate more task information -- minimize the code length required to compress

and transmit task labels



Training sample permutations for the In-context Learning setting

probing set construction method (Lu et al. 2022), showing the various possible ordering permutations of the randomly selected training samples, the resulting generation for each permutation, and the concatenation of each into a probing set

➤ **How Pre-Training Affects ICL?** *-- making*

*models learn to reason across demonstrations*



MEND (Li et al., 2024)



PICL (Gu et al., 2023)



ICLM (Shi et al., 2024)

➢ ## How LLMs Perform ICL?

 *-- based on given demonstrations at the*

*inference stage* (two main ways for LLMs to utilize

demonstrations)

    *Task recognition:* LLMs recognize the task from

demonstrations and utilize the prior knowledge

obtained from pre-training to solve new test tasks

    *Task learning:* LLMs learn new tasks unseen in

the pre-training stage only through demonstrations



three experimental settings (Pan et al., 2024)
- models can achieve non-trivial performance with only TR, and TR does not further improve with larger models or more demonstrations
- LLMs acquire TL as the model scales, and TL's performance consistently improves with more demonstrations in context

# Influencing Factors in ICL

➤ Pre-training Stage

   ○ pre-training corpora, data distribution, model architecture and training process

➤ Inference Stage

   ○ input-label settings, demonstration, demonstration-query



Summary of factors that have a relatively strong correlation to ICL performance and different perspectives to explain why ICL works

➢ **Chain-of-Thought (CoT) reasoning** offers a step-by-step reasoning trajectory, it decomposes intricate problems into manageable steps (*thoughts*), simplifying the overall reasoning process, and creates a linkage (*chain*) among the reasoning steps to ensure no important conditions are overlooked. CoT reasoning offers an observable reasoning process



Did Aristotle use a laptop?

Yes, Aristotle's ideas are known for their emphasis on empirical observation and practical wisdom, he …

1. Aristotle (**384–322 BC**) lives in **Ancient Greek**.

2. The first laptop computer hit shelves in **1981**.

Laptop computer didn't exist in Aristotle's time, so the answer is No.

direct output                chain-of-thought

# CoT Prompting in ICL

➤ **CoT prompting** is an improved prompting strategy to boost the performance of LLMs on complex reasoning tasks. Instead of simply constructing the prompts with input-output pairs like ICL, CoT prompting further incorporates intermediate reasoning steps, which serve as the bridge between inputs and outputs



A comparative illustration of ICL and chain-of-thought (CoT) prompting

# Basic CoT Prompting Approach

➤ CoT prompting is first proposed as an extension of ICL, *<input, output>* to *<input, CoT, output>*

➤ A CoT is a series of intermediate reasoning steps



An illustration of the evolution of CoT prompting strategies. Here, "thought" refers to an intermediate reasoning step

# CoT Topological Variants

➤ **Chain Structure**: the description format of rationales significantly influences reasoning execution

➤ **Tree Structure**: gain the capability to widely explore and backtrack during reasoning

➤ **Graph Structure**: outperform tree-based methods in handling complex problems but has poorer generalization



Topological variants emerging in the evolution of CoT. (a) standard I-O prompting, (b) parallel-constrained tree structure variants, (c) chain structure variants with distinct rationale descriptions, (d) chain structure variants with self-ensemble, (e) standard tree structure variants, and (f) standard graph structure variants.

# CoT Enhancement Methods

➤ Verify and Refine

    ○ can be an effective strategy for mitigating faithful errors in reasoning

    ○ Reasoning can be refined based on critical feedback provided by LLMs

    ○ logical reasoning structures are also well-suited for verification



Verification and refinement rectify intermediate errors, which reduce cascading errors in reasoning

> **Question Decomposition**

- ○ address intricate problems by progressively tackling straightforward sub-problems

- ○ involve decomposing both the questions and tables simultaneously when dealing with tabular reasoning

- ○ Bottom-up aggregation is also a viable solution, with a smaller exploration space



Question decomposition solves complex questions progressively by solving simple sub-questions

➢ Knowledge Enhancement

○ Introducing external knowledge or mining the model's internal knowledge can help dealing with knowledge-sensitive tasks

○ External knowledge is often more reliable than parametric knowledge

○ Bottom-up aggregation is also a viable solution, with a smaller exploration space



Incorporating knowledge (either internal or external) helps mitigate factual errors in reasoning

➢ **Self-Ensemble**

○ The sampling during generation introduces uncertainty, which in turn, creates the possibility of improving performance through self-ensemble

○ answer-based ensemble fails to consider intermediate steps

○ another concern is the limited diversity offered by probability sampling



Self-ensemble reduces inconsistency by selecting final answers from multiple samplings

# When CoT Prompting Works For LLMs?

○ Since CoT reasoning is an emergent ability, it only has a positive effect on sufficiently large models (typically containing 10B or more parameters) but not on small models

○ Since CoT prompting augments the standard prompting with intermediate reasoning steps, it is mainly effective for the tasks that require step-by-step reasoning, e.g., arithmetic reasoning, commonsense reasoning, and symbolic reasoning

○ For other tasks that do not rely on complex reasoning, CoT prompting might lead to worse performance than standard prompting

# Prompt-based Planning (Early Agent)

- Prompt-based Planning has been proposed to break down complex tasks into smaller sub-tasks and generate a plan of actions to accomplish the task

- typically three components
  - *task planner:* generating the whole plan to solve a target task
  - *plan executor:* executing the actions in the plan
  - *environment:* where the plan executor carries out the actions, which can be set differently according to specific tasks



An illustration of the formulation for prompt based planning by LLMs for solving complex tasks

# Retrieval-Augmented Generation

➢ **Retrieval-Augmented Generation (RAG)** incorporates information or knowledge from external data sources, which serves as supplementary for the input query or the generated output to advance generation models and enhance the generated results



RAG meets LLMs. When the user's query is out-of-scope, e.g., unseen content in training data or the need for the latest information for the answer, LLMs might shown ferior generation performance. With the help of RAG, LLMs can leverage additional relevant information from external database to enhance their text generation capability

➢ RAG first invokes the retriever to search and extract the relevant documents from external databases, which are leveraged as the context to enhance the generation process

➢ RAG is feasible and efficient to apply in various generation tasks with simple adaptation of the retrieval component

➢ great potential of RAG not only for knowledge-intensive tasks but also for general language tasks, and various downstream applications



: Representing RAG and RA-LLMs methods organized by their main design focus, proposed time and impact

- major processes
  - *retrieval*
  - *generation*
  - *augmentation*
  - the mechanism to determine whether the retrieval is needed



Illustration of the basic Retrieval-Augmented Large Language Models (RA-LLMs) framework for a specific QA task, which consists of three main components: retrieval, augmentation, and generation. Retrieval may have different procedures with various designs, which optionally includes pre-retrieval and post-retrieval processes. The retrieved documents are further leveraged in generation with the augmentation module, which may be at different integration stages

# Retrieval in RAG

➢ Retriever Type

○ sparse retrieval -- word-based and applied in text retrieval mostly

○ dense retrieval -- embedding queries and external knowledge into

vector spaces and can applied to various data formats

➢ Retrieval Granularity

-- denotes the retrieval unit in which the corpus is indexed

○ Chunk retrieval -- is common, which has been used in both

traditional and LLM-based RAG models such as REALM, RAG and

Atlas

○ token retrieval -- instead can be done with faster searching but will

bring more burden for the database saving

○ entity retrieval -- designed from the perspective of knowledge

rather than language



Illustration of the retriever in RA-LLMs, which can be implemented in either dense or sparse manners, each with several key operations

> ➢ Pre-retrieval and Post-retrieval Enhancement

-- to ensure the retrieval quality, i.e., increase the accuracy and relevance of the retrieved results

# Generation in RAG

➢ The design of the generator heavily depends on the downstream tasks

  ○ Parameter-Accessible Generators (White-box)

    -- word-based and applied in text retrieval mostly

    -- allow parameter optimization, which can be trained to adapt to different retrieval and augmentation approaches for a better performance of generation

  ○ Parameter-Inaccessible Generators (Black-box)

    -- only allow the operations of feeding queries (input) and receiving responses (output) while not allowing the internal structure to be altered or parameters to be updated

    -- Black-box RA-LLMs focus more on the retrieval and augmentation processes, trying to enhance the generator by augmenting the input (also called prompt in the context of LLMs) with better knowledge, guidance, or examples for the generation

➢ Retrieval Integration for Generation Augmentation

○ Input-Layer Integration

-- to integrate retrieved information/documents is to combine them with the original input/query and jointly pass them to the generator

○ Output-Layer Integration

-- it's post-hoc, which joints retrieval and generation results

○ Intermediate-Layer Integration

-- to design a semi-parametric module to integrate the retrieved results through the internal layers of the generation model, which is called intermediate-layer integration

➢ Retrieval Augmentation Necessity and Frequency

○  it is critical for RA-LLMs to accurately recall the prior knowledge while selectively

incorporating retrieved information only when necessary

○ Retrieval frequency affects both the efficiency and effectiveness of the model

-- one time

-- *every-n-token*

-- *every token*

# Retrieval Augmented LLMs (RA-LLMs)

- ➤ **Training-free**

- ➤ **Training-based**

  - ○ Independent Training

  - ○ Sequential Training

  - ○ Joint Training



An illustration of different training methods in RA-LLMs. Existing RA-LLMs approaches can be categorized into two classes: training-free approaches usually directly leverage retrieved information during the inference time by integrating the retrieved knowledge into the prompt, and training-based approaches fine-tune the retrieval and generator to enhance the generation performance. Based on the training strategies, training-based methods can be further categorized into three groups: independent training, where the retrieval and generator components are trained independently; sequential training, where they are trained sequentially; and joint training, where they are trained jointly

➢ **NLP applications**

➢ **Downstream tasks**

➢ **Domain-specific applications**

# RAG Challenges and Futures

➢ Trustworthy RA-LLMs

   -- 1) robustness, 2) fairness, 3) explainability, and 4)

privacy

➢ Multi-Lingual RA-LLMs

➢ Multi-modal RA-LLMs

➢ Quality of External Knowledge

*"Sharp tools make good work"*

—The Analects: Wei Ling Gong

➢ More generally, an LLM can access any number of external tools (e.g. an API to a service) to augment its functionality

➢ RAG can be seen as a specific instance of the broader category of the so called "tools"

➢ These tools extend the range of tasks an LLM can perform, from basic information retrieval to complex interactions with external databases or APIs

➤ the past year has witnessed a rapid surge in research efforts on tool learning concurrent with the rise of LLMs



An illustration of the development trajectory of tool learning

# Why Tool Learning?

- ➢ Knowledge Acquisition

- ➢ Expertise Enhancement

- ➢ Automation and Efficiency

- ➢ Interaction Enhancement

- ➢ Enhanced Interpretability and User Trust

# How Tool Learning?

➢ **Four Stages of Tool Learning**

- ○ **Task Planning**

- ○ **Tool Selection**

- ○ **Tool Calling**

- ○ **Response Generation**

➢ **Task Planning**

○ conduct a comprehensive analysis of the user intent

○ the planner is also tasked with delineating the dependencies and execution sequence of the decomposed tasks

○ facilitating the establishment of interconnections between the sub-questions

○ Tuning-free Methods and Tuning-based Methods

---

**An Example for Task Planning with GPT-4**

**Instruction Prompt:** You are currently in the task planning stage. You are given a user query requiring multi-step actions and reasoning. You will break down the user's query into sub-questions, and you only need to output these sub-questions after the breakdown. Ensure that the original problem is comprehensively covered through the minimal number of sub-questions.

**User Question:** *I would like to know the value of 5 ounces of gold plus 1 million AMZN stocks in CNY.*

**Output:** *1.What is the current price of gold per ounce in USD? 2.What is the current stock price of Amazon (AMZN) per share in USD? 3.What is the current exchange rate between USD and CNY (Chinese Yuan)?*

➢ **Tool Selection**

○ involves choosing through a retriever or directly

allowing LLMs to pick from a provided list of tools

○ Retriever-based Tool Selection

-- Term-based Methods and Semantic-based

Methods

○ LLM-based Tool Selection

-- Tuning-free Methods and Tuning-based Methods

---

**An Example for Tool Selection with GPT-4**

**Instruction Prompt:** You are currently in the tool se-
lection stage. You are given candidate tools that can
be potentially used to solve the sub-question. Among
candidate tools, select a list of relevant tools that
would help solve the sub-question.

**Sub-question 1:** *What is the current price of gold
per ounce in USD?*

**Candidate Tools:** *1.Metals Prices Rates API: The
latest API endpoint will return real-time exchange
rate data updated every 60 seconds. 2.Medium: Get
official news from Medium. 3.Cryptocurrency Mar-
kets: Recently published cryptocurrencies videos.*

**Output:** *1.Metals Prices Rates API: The latest API
endpoint will return real-time exchange rate data up-
dated every 60 seconds.*

**Sub-question 2:** · · ·

⋮

**Output:** · · ·

# How Tool Learning? (cont.)

➤ **Tool Calling**

- ○ extract the required parameters from the user query in accordance with the specifications outlined in the tool description and request data from tool servers

- ○ **Tuning-free Methods** and **Tuning-based Methods**

**An Example for Tool Calling with GPT-4**

**Instruction Prompt:** You are currently in the tool calling stage. You are given selected tools that can be potentially used to solve the sub-question. Your goal is to extract the required parameters needed to call the tool from the sub-question based on the tool descriptions. Output in the following format: {parameter name: parameter, · · · , parameter name: parameter}

**Sub-question 1:** *What is the current price of gold per ounce in USD?*

**Selected Tools:** *Tool Name: {Metals Prices Rates API}. Tool description: {The latest API endpoint will return real-time exchange rate data updated every 60 seconds.} Required params:{ [name: symbols, type: STRING, description: Enter a list of comma-separated currency codes or metal codes to limit output codes., name: base, type: STRING, description: Enter the three-letter currency code or metal code of your preferred base currency.] }*

**Output:** *{symbols: "XAU", base: "USD"}*

**Sub-question 2:** · · ·

⋮

**Output:** · · ·

➤ **Response Generation**

   ○ synthesize information relevant to

user queries and integrate their own

knowledge to construct comprehensive

responses

   ○ **Direct Insertion Methods** and

**Information Integration Methods**

---

**An Example for Response Generation with GPT-4**

**Instruction Prompt:** You are currently in the response generation stage. You are given the output of selected tools that can be potentially used to solve the sub-task. Your goal is to generate a response to the user question based on the results from tools and your internal knowledge.
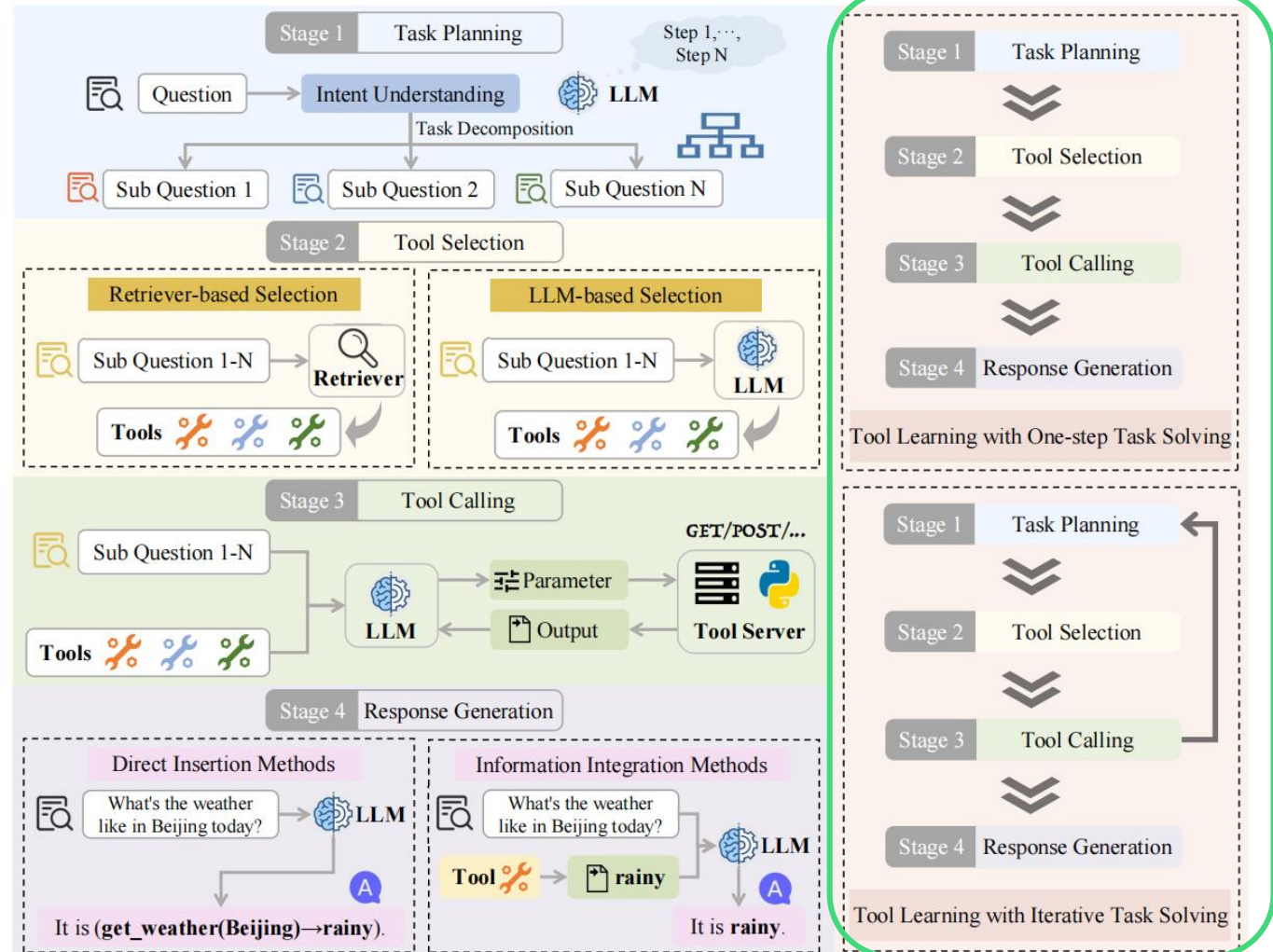
**User Question:** *I would like to know the value of 5 ounces of gold plus 1 million AMZN stocks in CNY.*

**Return from Tool Calling:** *The stock price of AMZN is 130 USD, XAU USD is 1919.7605807522993, and the exchange rate between USD and YNC is 7.24.*

**Output:** $\cdots$ *Total value in CNY: Total value = Value of gold in CNY + Value of stocks in CNY = 69,495.87 + 941,200,000 = 941,269,495.87 CNY. Therefore, the combined value of 5 ounces of gold and 1 million Amazon stocks in Chinese Yuan is approximately 941,269,495.87* $\cdots$.

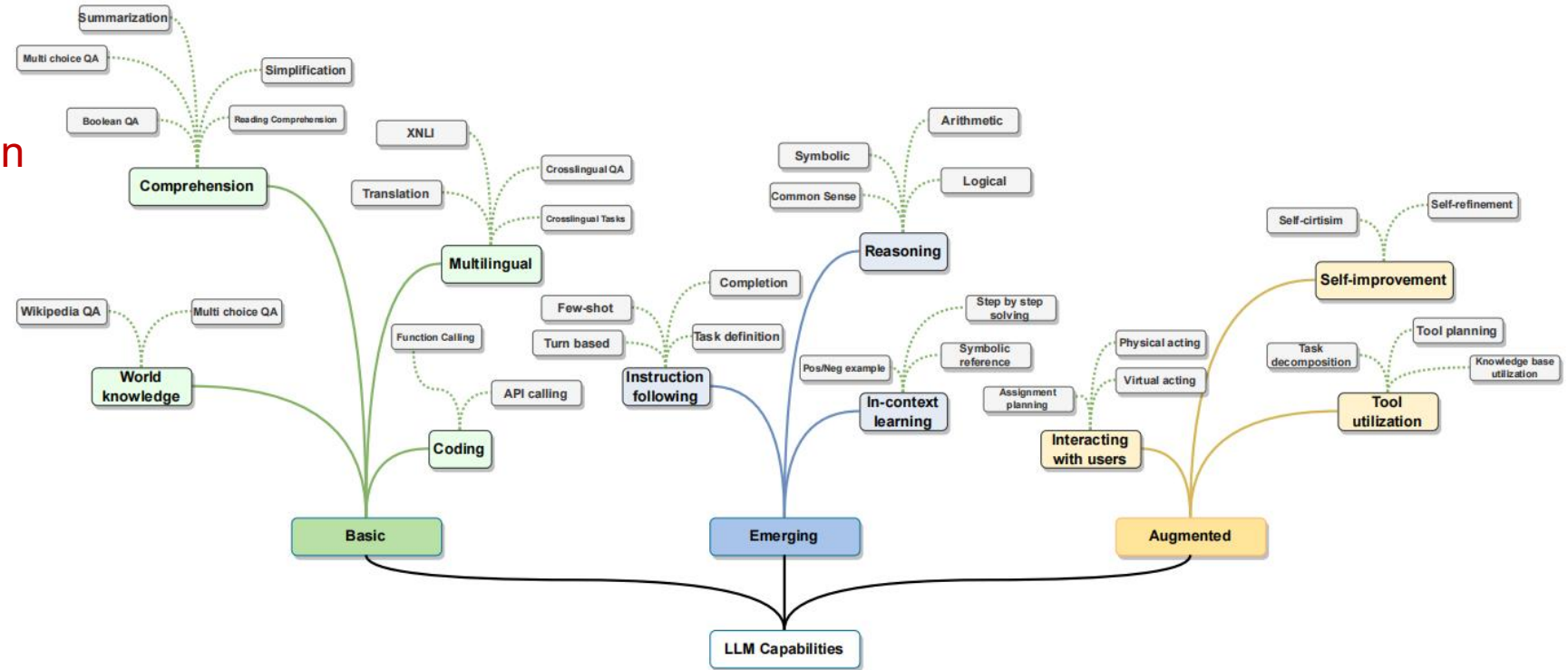➤ Two Paradigms of Tool Learning

   ○ Tool Learning with One-step Task Solving

   ○ Tool Learning with Iterative Task Solving

# Tools Learning Chanllenges and Future

➢ High Latency in Tool Learning

➢ Rigorous and Comprehensive Evaluation

➢ Comprehensive and Accessible Tools

➢ Safe and Robust Tool Learning

➢ Unified Tool Learning Framework

➢ Real-Word Benchmark for Tool Learning

➢ Tool Learning with Multi-Modal

- Basic Ability

- Knowledge Utilization

- Complex Reasoning

# LLM Evaluation: Basic Ability

➢ Language Generation (categories)

○ Language Modeling -- to predict the next token based on the previous tokens

○ Conditional Text Generation -- generating texts satisfying specific task demands based on the given conditions, typically including machine translation, text summarization, and question answering

○ Code Synthesis -- to generate formal language, especially computer programs (i.e., code) that satisfy specific conditions

# LLM Evaluation: Basic Ability

➢ Language Generation (major issues)

○ Unreliable generation evaluation -- pronounced

inconsistency between human evaluation and

automatic reference-based metrics

○ Underperforming specialized generation --

LLM's proficiency in generation might be

constrained when dealing with a specialized

domain or task

### Unreliable Generation Evaluation

LLMs have been capable of generating texts with a comparable quality to human-written texts, which however might be underestimated by automatic reference-based metrics. As an alternative evaluation approach, LLMs can serve as language generation evaluators to evaluate a single text, compare multiple candidates, and improve existing metrics. However, this evaluation approach still needs more inspections and examinations in real-world tasks.

### Underperforming Specialized Generation

LLMs may fall short in mastering generation tasks that require domain-specific knowledge or generating structured data. It is non-trivial to inject specialized knowledge into LLMs, meanwhile maintaining the original abilities of LLMs.

# LLM Evaluation: Basic Ability

➢ **Knowledge Utilization (categories)**

○ **Closed-Book QA** -- test the acquired factual knowledge of LLMs from the pre-training corpus, where LLMs should answer the question only based on the given context without using external resources

○ **Conditional Text Generation** --  LLMs can extract useful evidence from the external knowledge base or document collections, and then answer the question based on the extracted evidence

○ **Knowledge Completion**  -- LLMs might be (to some extent) considered as a knowledge base, which can be leveraged to complete or predict the missing parts of knowledge units

➢ **Knowledge Utilization (major issues)**

○ Hallucination -- the generated information is either in conflict with the existing source (*intrinsic hallucination*) or cannot be verified by the available source (*extrinsic hallucination*)

**Hallucination**

LLMs are prone to generate untruthful information that either conflicts with the existing source or cannot be verified by the available source. Even the most powerful LLMs such as ChatGPT face great challenges in migrating the hallucinations of the generated texts. This issue can be partially alleviated by special approaches such as alignment tuning and tool utilization.

Bob's wife is Amy. Bob's daughter is Cindy. Who is Cindy to Amy?

Explain RLHF for LLMs.

Cindy is Amy's daughter-in-law.

RLHF stands for "Rights, Limitations, Harms, and Freedoms" and is a framework for ...... models like LLMs (Large Language Models).

(a) Intrinsic hallucination

(b) Extrinsic hallucination

➢ Knowledge Utilization (major issues)

○ Knowledge recency -- LLMs would

encounter difficulties when solving tasks

that require the latest knowledge beyond

the training data

**Knowledge Recency**

The parametric knowledge of LLMs is hard to be updated in a timely manner. Augmenting LLMs with external knowledge sources is a practical approach to tackling the issue. However, how to effectively update knowledge within LLMs remains an open research problem.

➢ **Complex Reasoning (categories)**

-- the ability of understanding and utilizing supporting evidence or logic to derive conclusions or make decisions

○ Knowledge Reasoning -- to rely on logical relations and evidence about factual knowledge to answer the given question

○ Symbolic Reasoning -- to manipulate the symbols in a formal rule setting to fulfill some specific goal, where the operations and rules may have never been seen by LLMs during pre-training

○ Mathematical Reasoning -- to comprehensively utilize mathematical knowledge, logic, and computation for solving problems or generating proof statements

# LLM Evaluation: Complex Reasoning

➢ **Complex Reasoning (major issues)**

○ **Reasoning inconsistency** -- LLMs may generate the correct answer following an invalid reasoning path, or produce a wrong answer after a correct reasoning process, leading to inconsistency between the derived answer and the reasoning process

○ **Numerical computation** -- face difficulties in the involved numerical computation, especially for the symbols that are seldom encountered during pre-training, such as arithmetic with large numbers

**Reasoning Inconsistency**

LLMs may generate the correct answer following an invalid reasoning path, or produce a wrong answer after a correct reasoning process, leading to inconsistency between the derived answer and the reasoning process. The issue can be alleviated by fine-tuning LLMs with process-level feedback, using an ensemble of diverse reasoning paths, and refining the reasoning process with self-reflection or external feedback.

**Numerical Computation**

LLMs face difficulties in numerical computation, especially for the symbols that are seldom encountered during pre-training. In addition to using mathematical tools, tokenizing digits into individual tokens is also an effective design choice for improving the arithmetic ability of LLMs.

# LLM Evaluation: Advanced Ability

➢ **Human Alignment**

-- LLMs could well conform to human values and needs, i.e., human alignment

➢ **Interaction with External Environment**

-- to receive feedback from the external environment and perform actions according to the behavior instruction

➢ **Tool Manipulation**

-- LLMs can turn to external tools if they determine it is necessary to enhance the performance of LLMs on several specific tasks
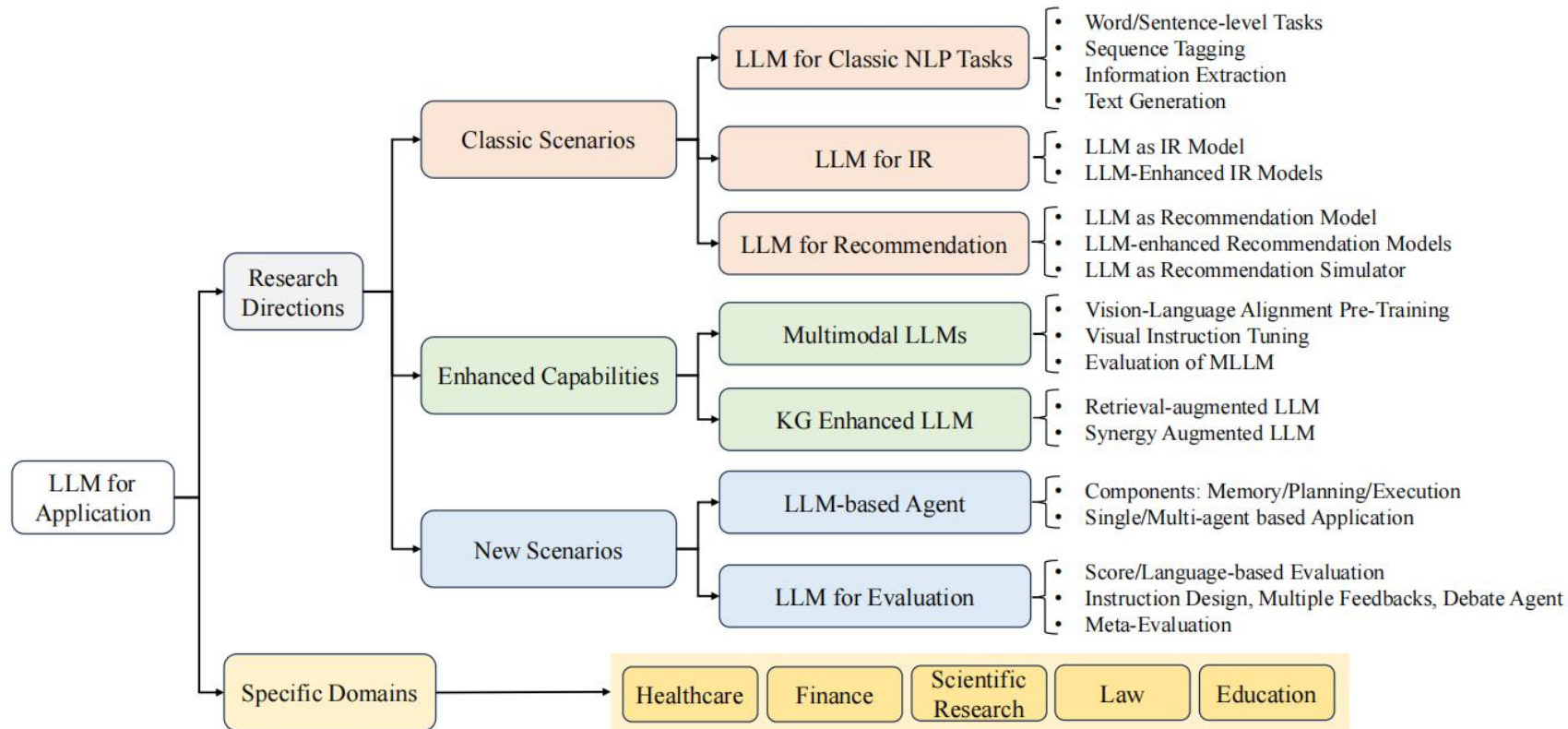
# LLM Benchmarks

➢ Comprehensive Evaluation Benchmarks

○ MMLU -- a versatile benchmark for large-scale evaluation of multi-task knowledge understanding

○ BIG-bench --  a collaborative benchmark intended to probe existing LLMs from various aspects

○ HELM  --  a comprehensive benchmark that currently implements a core set of 16 scenarios and 7 categories of metrics

○ Human-level test benchmarks  --  evaluate the comprehensive ability of LLMs with questions designed for testing humans

# LLM Applications

**As LLMs have revolutionized the way how we develop AI algorithms, it poses significant impact on the research community**

# LLMs Optimization Techniques

➢ Model training

➢ Inference



Illustration of model performance and model training time in GPU hours of LLaMA models at different scales

: Performance score *vs.* inference throughput for various LLMs. The throughputs are measured on Nvidia A100 80GB GPU with 16-bit floating point quantization.

**Perspectives:**

➢ model-centric

➢ data-centric

➢ framework-centric

# Model-Centric Methods

➤   Model-Centric Methods -- focus on both algorithm-level and system-level efficient

techniques where the model itself is the focal point

➤  Categories

- ○ Model Compression

- ○ Efficient Pre-Training

- ○ Efficient Fine-Tuning

- ○ Efficient Inference

- ○ Efficient Architecture

➢ Model Compression -- reducing the sizes and the amount of arithmetic operations of LLMs

- ➢   Quantization

- ➢   Parameter Pruning

  - ○  Structured Pruning --  pruning structured patterns

  - ○  Unstructured Pruning -- pruning model weights individually

- ➢   Low-Rank Approximation --  approximating the LLM weight matrix with smaller low-rank matrices

- ➢   Knowledge Distillation

  - ○  White-Box Knowledge Distillation -- the parameters or logits of the teacher LLM are used in the distillation process

  - ○  Black-Box Knowledge Distillation --  only the outputs generated from the teacher LLM are used in the distillation process



Illustrations of model compression techniques for LLMs

# Efficient Pre-Training

➢ Efficient pre-training -- reducing the costs of the LLM pre-training process in terms of compute resources, training time, memory and energy consumption



**Efficient Pre-Training**

- **Mixed Precision Training**: AMP (Micikevicius et al., 2018; Facebook AI Research (FAIR), 2023; Rae et al., 2022), Brain Floating Point (BFLOAT16) (Kalamkar et al., 2019; Burgess et al., 2019)

- **Scaling Models**: Progressive Stacking (Gong et al., 2019), MSLT (Yang et al., 2020), CompoundGrow (Gu et al., 2021), bert2BERT (Chen et al., 2022b), Knowledge Inheritance (Qin et al., 2022), Staged Training (Shen et al., 2022), LiGO (Wang et al., 2023d), Mango (Pan et al., 2023), Yao et al. (2024), Growth Strategy (Li et al., 2023e)

- **Initialization Techniques**: Kumar (2017), Fixup (Zhang et al., 2019), ZerO (Zhao et al., 2022), SkipInit (De & Smith, 2020), ReZero (Bachlechner et al., 2021), T-Fixup (Huang et al., 2020), DeepNet (Wang et al., 2024a)

- **Training Optimizers**: Lion (Chen et al., 2023g), Sophia (Liu et al., 2024a)

- **System-Level Pre-Training Efficiency Optimization**: ZeRO (Rajbhandari et al., 2020), FSDP (Zhao et al., 2023c), ZeRO-Offload (Ren et al., 2021), ZeRO-Infinity (Rajbhandari et al., 2021), Zeus (You et al., 2023), Perseus (Chung et al., 2023)

- ➢ **Mixed Precision Training**

- ➢ **Scaling Models**

- ➢ Initialization Techniques

- ➢ Training Optimizers

- ➢ System-Level Pre-Training

Efficiency Optimization



(a) Mixed Precision Training

(b) Scaling Models

(c) Initialization Techniques

(d) Training Optimizers

Illustrations of efficient pre-training techniques for LLMs

➢ Efficient Fine-Tuning -- reducing the costs of the LLM fine-tuning process

**Efficient Fine-Tuning**

**Parameter-Efficient Fine-Tuning**

**Low-Rank Adaptation**
LoRA (Hu et al., 2022), LoRA-FA (Zhang et al., 2023b), LoraHub (Huang et al., 2023), LongLoRA (Chen et al., 2023i), Multi-Head Routing (Caccia et al., 2023), AdaLoRA (Zhang et al., 2023d), DyLoRA (Valipour et al., 2023), CEPT (Zhao et al., 2023b), Tied-LoRA (Renduchintala et al., 2023)

**Adapter-based Tuning**
LLM-Adapters (Hu et al., 2023b), Compacter (Karimi Mahabadi et al., 2021), $(IA)^3$ (Liu et al., 2022a), Meta-Adapters (Bansal et al., 2022), AdaMix (Wang et al., 2022c), OpenDelta (Hu et al., 2023a), SparseAdapter (He et al., 2022b)

**Prefix Tuning**
Prefix-Tuning (Li & Liang, 2021), LLaMA-Adapter (Zhang et al., 2024), HyperTuning (Phang et al., 2023)

**Prompt Tuning**
Prompt Tuning (Lester et al., 2021), P-Tuning (Liu et al., 2023b), P-Tuning v2 (Liu et al., 2022c), Tam et al. (2023), $MP^2$ (Sun et al., 2023a), PPT (Gu et al., 2022b), Multitask Prompt Tuning (Wang et al., 2023j), Xu et al. (2023b)

**Memory-Efficient Fine-Tuning**
QLoRA (Dettmers et al., 2023), QA-LoRA (Xu et al., 2024b), LoftQ (Li et al., 2024d), PEQA (Kim et al., 2023a), Selective Fine-Tuning (Simoulin et al., 2023), LOMO (Lv et al., 2023), MeZO (Malladi et al., 2023), Liu et al. (2023g)

- ➤ Parameter-Efficient Fine-Tuning (PEFT)

  - ○ Low-Rank Adaptation (LoRA)

  - ○ Adapter-based Tuning

  - ○ Prefix Tuning

  - ○ Prompt Tuning

- ➤ Memory-Efficient Fine-Tuning



(a) Low-Rank Adaptation

(b) Adapter-based Tuning

(c) Prefix Tuning

(d) Prompt Tuning

(e) Memory-Efficient Fine-Tuning

Illustrations of PEFT (a)-(d) and memory-efficient fine-tuning (e)

➢ Efficient Inference -- reducing the costs of the LLMs inference process



Speculative Decoding (Leviathan et al., 2023), Chen et al. (2023a), Staged Speculative (Spector & Re, 2023), BiLD (Kim et al., 2023c), SpecInfer (Miao et al., 2024), LLMA (Yang et al., 2023b), Medusa (Cai et al., 2024), Santilli et al. (2023), PaSS (Monea et al., 2023)

KIVI (Zirui Liu et al., 2023), KVQuant (Hooper et al., 2024), Heavy-Hitter Oracle ($H_2O$) (Zhang et al., 2023f), Scissorhands (Liu et al., 2023e), StreamingLLM (Xiao et al., 2024)

FlexGen (Sheng et al., 2023), Pope et al. (2023), $S^3$ (Jin et al., 2023), Orca (Yu et al., 2022), vLLM (Kwon et al., 2023), DeepSpeed-Inference (Aminabadi et al., 2022), Flash-Decoding (Dao et al., 2023), FlashDecoding++ (Hong et al., 2023)

Efficient Inference

Algorithm-Level Inference Efficiency Optimization

Speculative Decoding

KV-Cache Optimization

System-Level Inference Efficiency Optimization

# Efficient Inference (cont.)

➢ Algorithm-Level Inference Efficiency Optimization

- ○ Speculative Decoding -- a decoding strategy for autoregressive language models

- ○ KV-Cache Optimization -- reducing the size of KV cache

➢ System-Level Inference Efficiency Optimization

-- can also be optimized at the system level under a specific hardware architecture



(a) Speculative Decoding

(b) KV-Cache Optimization

- ➢ **KV-Cache Optimization** has emerged as a pivotal solution to the issue of the Transformer architecture's struggle with handling long texts
  - ○ from *the training phase*, to the deployment phase, and finally to the post-training phase



An overview of the main structure of KV-Cache compression methods

> Efficient Architecture

Design -- the strategic

optimization of model

architecture and

computational processes

➢ Efficient Attention

○ Sharing-based Attention

○ Kernelization or Low-Rank

○ Fixed Pattern Strategies

○ Learnable Pattern Strategies

○ Hardware-Assisted Attention



(a) Sharing-based Attention

(b) Kernelization or Low-Rank

(c) Fixed Pattern Strategies

(d) Learnable Pattern Strategies

FlashAttention

➢ Mixture of Experts (MoE)

○ MoE-based LLMs

○ Algorithm-Level MoE Optimization

○ System-Level MoE Optimization

➢ Long Context LLMs

○ Positional Extrapolation and Interpolation

○ Segmentation and Sliding Window

○ Memory-Retrieval Augmentation

➢ Transformer-Alternate Architectures

○ State Space Models

○ Other Sequential Models



(a) Mixture of Experts (MoE)

(b) Long Context LLMs

# Data-Centric Methods

➤ Data selection -- a fundamental technique for enhancing efficiency

➤ Data Selection for Efficient Pre-Training

➤ Data Selection for Efficient Fine-Tuning



Data Selection

| Data Selection for Efficient Pre-Training | SSPT (Glass et al., 2020), Yao et al. (2022a), DSIR (Xie et al., 2023b), DoReMi (Xie et al., 2023a) |
| Data Selection for Efficient Fine-Tuning | Ivison et al. (2023), Instruction Mining (Cao et al., 2023), TS-DShapley (Schoch et al., 2023), LTD Instruction Tuning (Chen et al., 2023b), AlpaGasus (Chen et al., 2024b), LIMA (Zhou et al., 2023a) |

(a) Data Selection for Efficient Pre-Training    LLMs    (b) Data Selection for Efficient Fine-Tuning

➢ Prompt Engineering -- designing effective inputs (i.e., prompts) to guide

LLMs in generating desired outputs

➤ Few-Shot Prompting

○ Demonstration Organization

   -- Demonstration Selection

   -- Demonstration Ordering

○ Template Formatting

   -- Instruction Generation

   -- Multi-Step Reasoning

➢ **Prompt Compression**

-- accelerates the processing of LLM inputs through either condensing lengthy prompt inputs or learning compact prompt representations

➢ **Prompt Generation**

-- automatically creating effective prompts that guide the model in generating specific and relevant responses

# LLM Frameworks

➢   LLM frameworks can be in general grouped based on whether they support the tasks of training, fine-tuning, and inference

➢   Frameworks that support training and/or fine-tuning aim to provide scalable, efficient, and flexible infrastructure that improves computation efficiency, reduces memory footprint, optimizes communication efficiency, and ensures reliability of the training/fine-tuning process

➢ Frameworks that support inference focus on optimizing inference throughput and reducing memory footprint and latency

# 02

# On the Shoulders of LLMs: LLM Reasoning

**Yao Yao**

Shanghai Jiao Tong University

yaoyao27@sjtu.edu.cn

# Content

142

# 01

# **What is reasoning?**

# 1. What is reasoning?

> *Whatever intelliigence may be, reasoning and problem solving have traditionally been viewed as important subsets of it.*
>
> Sternberg, Robert J., ed. *Handbook of human intelligence*. Cambridge university press, 1982.

- Reasoning is an critical component of language intelligence.
- **Language intelligence** involves the ability to understand and process ideas expressed in natural languages
- Reasoning involves a **multi-step process that uses evidence, arguments, and logic to arrive at conclusions or make sophisticated predictions**.

# 1. What is reasoning?

- Large Language Models (LLMs) have shown impressive results in complex reasoning tasks
- From traditional reasoning tasks to general-purpose reasoning task
- From single-modal reasoning to multimodal reasoning



Sun, Jiankai, et al. "A survey of reasoning with foundation models." *arXiv preprint arXiv:2312.11562* (2023).

# How do LLM perform reasoning

# 2. How do LLM perform reasoning?

➢ **Directly generate answers**

| Input |
|---|
| A small sunflower has 3 dozen seeds and a large sunflower has 50% more seeds than a small sunflower. How many sunflower seeds are there altogether? |

**Large Language Model**

| Output |
|---|
| 54 |

❌

➢ **Chain-of-thought (CoT) prompting**

| Input |
|---|
| A small sunflower has 3 dozen seeds and a large sunflower has 50% more seeds than a small sunflower. How many sunflower seeds are there altogether? Let's think step by step. |

**Large Language Model**

| Output |
|---|
| A small sunflower has 3 x 12 = 36 seeds. A large sunflower has 36 x 0.5 = 18 more seeds. So a large sunflower has 36 + 18 = 54 seeds. Altogether, there are 36 + 54 = 90 sunflower seeds. |

✅

147

**IAPR**

➢ **Chain-of-thought (CoT) prompting**

**Input**

A small sunflower has 3 dozen seeds and a large sunflower has 50% more seeds than a small sunflower. How many sunflower seeds are there altogether? Let's think step by step.

**Large Language Model**

**Output**

A small sunflower has 3 x 12 = 36 seeds. A large sunflower has 36 x 0.5 = 18 more seeds. So a large sunflower has 36 + 18 = 54 seeds. Altogether, there are 36 + 54 = 90 sunflower seeds.

- CoT prompts LLMs to generate a series of intermediate reasoning steps for solving a problem (telling a maths student to "show their working".)

- CoT deconstructs complex issues into smaller, easily understandable, and manageable sub-problems, systematically enabling solutions through a step-by-step approach, leading to more logical and accurate answers.

# 2. How do LLM perform reasoning?

➢ **Two types of CoT: Zero-Shot and Few-shot**

🟩 Exemplars   🟥 Question   🟪 Instruction   🟦 Rationale   🟧 Answer

## Chain-of-Thoughts (Zero-Shot)

*Input*

> A small sunflower has 3 dozen seeds and a large sunflower has 50% more seeds than a small sunflower.  How many sunflower seeds are there altogether? Let's think step by step.

———  ———  ———  ———  ———  ———  ———

*Output*

> A small sunflower has 3 x 12 = 36 seeds. A large sunflower has 50% more seeds than a small sunflower, which is 36 x 0.5 = 18 more seeds. So a large sunflower has 36 + 18 = 54 seeds. Altogether, there are 36 + 54 = 90 sunflower seeds. The answer is 90.

➢ Zero-Shot-CoT typically relies on instructions to facilitate the LLM in conducting step-by-step reasoning

➢ Rationales: intermediate processes of CoT reasoning (solutions, intermediate reasoning steps, relevant external knowledge)

149

## ➢ Two types of CoT: Zero-Shot and Few-shot

**Exemplars**   **Question**   **Instruction**   **Rationale**   **Answer**

**Chain-of-Thoughts (Few-Shot)**

*Input*

Q: There are 15 trees in the grove. Grove workers will plant trees in the grove today. After they are done, there will be 21 trees. How many trees did the grove workers plant today?
A: There are 15 trees originally. Then there were 21 trees after some more were planted. So there must have been 21 - 15 = 6. The answer is 6.

Q: A small sunflower has 3 dozen seeds and a large sunflower has 50% more seeds than a small sunflower.  How many sunflower seeds are there altogether?
A:

*Output*

A: A small sunflower has 3 x 12 = 36 seeds. A large sunflower has 50% more seeds than a small sunflower, which is 36 x 0.5 = 18 more seeds. So a large sunflower has 36 + 18 = 54 seeds. Altogether, there are 36 + 54 = 90 sunflower seeds.

The answer is 90.

➢ Few-Shot-CoT: concatenate a set of exemplars with associated rationales with the question and serve as in-context demonstrations

150

## ➢ Benefits of CoT

- **Improved Reasoning Performance**

1. CoT reduces the risk of missing important details

2. CoT ensures that computational resources are allocated efficiently.

3. Research across various fields has consistently shown that CoT boosts performance.

2. How do LLM perform reasoning?

# 2. How do LLM perform reasoning?

## ➢ Benefits of CoT

- **Improved Reasoning Performance**

- **Improved Interpretability**

CoT makes the reasoning processes of LLMs transparent, allowing us to follow the logical steps leading to the conclusion, which is invaluable for debugging and improving models.

- **Improved Controllability**

CoT guides LLMs more effectively which makes it possible to refine the model's focus and correct paths in the reasoning process that may lead to errors. It's a powerful tool for ensuring accurate and reliable outputs.

- **Improved Flexibility**

CoT adapts well to various applications beyond traditional tasks and can be easily implemented in LLMs

152

> **Paradigm shifts of CoT**



- Prompting pattern
- Reasoning format
- Application scenario

## ➤ Paradigm shifts of CoT

- ■ **Prompting pattern**

  - • **Instruction generation**

  Find the optimal instructions to prompt LLM for step-by-step reasoning.

  Mainly aims to maximize LLM's <span style="color:red">zero-shot capability</span>

  - • **Exemplar generation**

  Find the best set of input-output demonstration exemplar pairs to

  prompt LLMs for step-by-step reasoning.

  Mainly aims to maximize LLM's <span style="color:red">few-shot capability</span>



**Instruction Generation**

Zero-Shot -CoT — Let's think step by step.

Plan-and-Solve Prompting — Let's first understand the problem and devise a plan to solve the problem. Then, let's carry out the plan and solve the problem step by step.

Automatic Prompt Engineer — Let's work this out in a step by step way to be sure we have the right answer.

OPRO

**Exemplar Generation**

Manual-CoT

Active-Prompt

Auto-CoT

Automate-CoT

154

# 2. How do LLM perform reasoning?

➢ **Paradigm shifts of CoT**

- **Instruction generation**

  1) Manually constructed instructions

     Outperforms zero-shot LLM performances without the need for hand-crafted few-shot examples

     Need to test various prompts to achieve the desired behavior

  2) Automated generation and selection of instructions

# 2. How do LLM perform reasoning?

➢ **Paradigm shifts of CoT**

- **Instruction generation**

  2) Automated generation and selection of instructions



Zhou, Yongchao, et al. "Large language models are human-level prompt engineers." *arXiv preprint arXiv:2211.01910* (2022).

➢ **Paradigm shifts of CoT**

- **Exemplar Generation**

  1) Manually exemplar generation



Manual Demos One by One

Q: There are 15 trees in the grove. Grove workers will plant trees in the grove today. After they are done, there will be 21 trees. How many trees did the grove workers plant today?
*Question*

A: There are 15 trees originally. Then there were 21 trees after some more were planted. So there must have been 21 - 15 = 6. The answer is 6.
*Answer* … *Rationale*

Q: A pet store had 64 puppies. In one day they sold 28 of them and put the rest into cages with 4 in each cage. How many cages did they use?
A:

LLM   *Test Question*

The pet store had 64 puppies. They sold 28 of them. So they had 64 - 28 = 36 puppies left. They put them into cages with 4 in each cage. So they used 36 / 4 = 9 cages. The answer is 9.

Manual-CoT

Exemplars are crafted by hand to guide the model's output through specific reasoning steps



**Exemplar Generation**

Manual-CoT

Active-Prompt

Auto-CoT

Automate-CoT

Zhang, Zhuosheng, et al. "Automatic chain of thought prompting in large language models." ICLR 2023

## ➤ Paradigm shifts of CoT

- **Exemplar Generation**

2) Automatic exemplar generation

Systems optimize the selection of examples to improve effectiveness

Auto-CoT

Zhang, Zhuosheng, et al. "Automatic chain of thought prompting in large language models." ICLR 2023

# 2. How do LLM perform reasoning?

➢ **Paradigm shifts of CoT**

- **Exemplar Generation**

  - LLMs tend to make mistakes on similar types of questions

  - Too many incorrect exemplars can decrease the LLM's performance

  Diverse set of exemplars can mitigate this misleading effects!



Zhang, Zhuosheng, et al. "Automatic chain of thought prompting in large language models." ICLR 2023

# 2. How do LLM perform reasoning?

## ➢ Paradigm shifts of CoT



- Prompting pattern

- Reasoning format

- Application scenario

➢ **Paradigm shifts of CoT**

• **Reasoning Format**

sequential ⟶ *structural information* ⟶ cognitive structures

# 2. How do LLM perform reasoning?

> **Paradigm shifts of CoT**

- **Reasoning Format**

  Improves the performance through the aggregation of results

Narang, Sharan, Aakanksha Chowdhery, and Denny Zhou. "Self-Consistency Improves Chain of Thought Reasoning in Language Models." ICLR 2023

# 2. How do LLM perform reasoning?

> **Paradigm shifts of CoT**

- **Reasoning Format**

  Rationale-augmented ensembles: reduce the brittleness of model outputs by

  aggregating multiple rationales.



Sampling rationale in the output space consistently yields the best improvements in task performance

Wang, Xuezhi, et al. "Rationale-augmented ensembles in language models." arXiv preprint arXiv:2207.00747 (2022).

163

➢ **Paradigm shifts of CoT**

• **CoT verification**

Introduce verification methods to verify and amend the CoT reasoning

Weng, Yixuan, et al. "Large Language Models are Better Reasoners with Self-Verification." *Findings of the Association for Computational Linguistics: EMNLP 2023*.

> ## Paradigm shifts of CoT

- ### CoT verification

    <span style="color:red">Incorporate external tools</span> (search engines, calculators) to enhance the factual

    accuracy and logical consistency of the LLM reasoning process.



Verify-and-Edit framework



- Interacts with appropriate external tools to evaluate certain aspects of the answer,
- Revising the output based on the feedback.

CRITIC framework

## ➤ Paradigm shifts of CoT

- **CoT verification**

**Multi-Agent Debate (MAD)**

Introduce adversarial testing: multiple models debate to refine the reasoning
Ensuring diverse perspectives are considered before reaching a conclusion.



Liang, Tian, et al. "Encouraging divergent thinking in large language models through multi-agent debate." *arXiv preprint arXiv:2305.19118* (2023).

# 2. How do LLM perform reasoning?

## ➢ Paradigm shifts of CoT

- **CoT verification**

  Introduce verification methods to verify and amend the CoT reasoning

  - A shift towards more dynamic and reliable CoT verification processes
  - Focusing on external validations to <u>overcome the limitations of heavily relying on models' built-in capabilities</u>, leading to errors or oversights in complex reasoning tasks



**CoT Reasoning**

**Cot Formulation**
- Program-of-Thoughts
- Tab-CoT
- Tree-of-Thoughts
- Graph-of-Thought (Rationale)
- Skeleton-of-Thought
- Recursion of Thought

**Reasoning Aggregation**
- Rationale-Augmented Ensembles
- Self-consistency CoT

**CoT Verification**
- Natural Program
- PRM
- Self-Verification
- CRITIC
- Verify-and-Edit
- External Knowledge
- AuRoRA
- Google / WIKIPEDIA

# 2. How do LLM perform reasoning?

➢ **Paradigm shifts of CoT**



- From single-language tasks to multilingual tasks
- From single-modality to multi-modalities
- From complex reasoning tasks to general-purpose tasks

# 2. How do LLM perform reasoning?

> ## Paradigm shifts of CoT

- ### From Single Language to Multilingual Scenarios.

Integrating multilingual examples
English rationales lead to better results

Aligns different language representations using
a step-by-step English translation



Multilingual Grade School Math (MGSM) benchmark

Cross-lingual prompting

# 2. How do LLM perform reasoning?

> **Paradigm shifts of CoT**

- From Text Modality to Multimodalities.

# 2. How do LLM perform reasoning?

> ➤ **Paradigm shifts of CoT**

- **From Text Modality to Multimodalities.**

Multimodal-CoT incorporates language (text) and vision (images) modalities into a two-stage framework that separates rationale generation and answer inference



Use single-head attention network to fuse different modality representations

Zhang, Zhuosheng, et al. "Multimodal chain-of-thought reasoning in language models." *arXiv preprint arXiv:2302.00923* (2023).
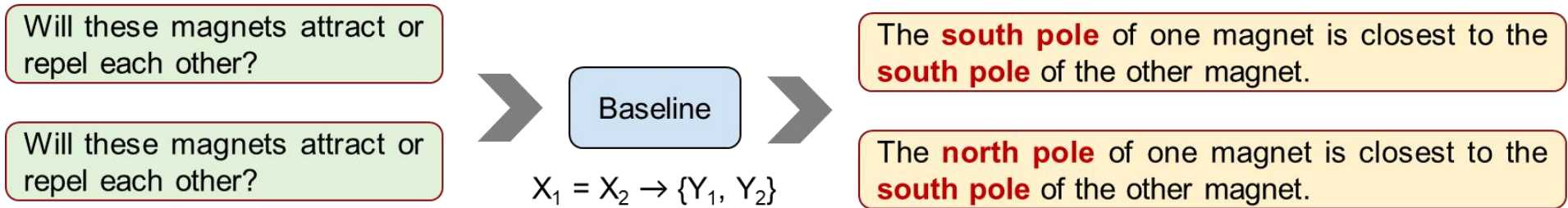
# 2. How do LLM perform reasoning?

- **From Text Modality to Multimodalities.**

  Introducing more **reliable input results** in more **convincing reasoning process**

- ☐ Case studies: 50 error cases
  - ● Imperfect training data: when the vision input is missing
  - ● Generate hallucinated rationales that mislead the answer inference (64%)



**Problem**

**Question:** Will these magnets attract or repel each other?
**Context:** Two magnets are placed as shown. Hint: Magnets that attract pull together. Magnets that repel push apart.

Vision: S N S N

**Options:** (A) attract  (B) repel

**Gold Rationale:** Magnets can pull or push on each other without touching. When magnets attract, they pull together. When magnets repel, they push apart. Whether a magnet attracts or repels other magnets depends on the positions of its poles, or ends. Every magnet has two poles, called north and south. Here are some examples of magnets. The north pole of each magnet is marked N, and the south pole is marked S. If different poles are closest to each other, the magnets attract. The magnets in the pair below attract. If the same poles are closest to each other, the magnets repel. The magnets in both pairs below repel. Will these magnets attract or repel? To find out, look at which poles are closest to each other. The north pole of one magnet is closest to the south pole of the other magnet. Poles that are different attract. So, these magnets will attract each other.
**Answer:** The answer is (A).

**Baseline**

**Generated Rationale:** Magnets can pull or push on each other without touching. When magnets attract, they pull together. When magnets repel, they push apart. Whether a magnet attracts or repels other magnets depends on the positions of its poles, or ends. Every magnet has two poles, called north and south. Here are some examples of magnets. The north pole of each magnet is marked N, and the south pole is marked S. If different poles are closest to each other, the magnets attract. The magnets in the pair below attract. If the same poles are closest to each other, the magnets repel. The magnets in both pairs below repel. Will these magnets attract or repel? To find out, look at which poles are closest to each other. The south pole of one magnet is closest to the south pole of the other magnet. Poles that are the same repel. So, these magnets will repel each other.
**Answer:** The answer is (B).

(a) ratio of hallucination mistakes       (b) correction rate w/ vision features

*Table 3.* Two-stage setting of (i) rationale generation (RougeL) and (ii) answer inference (Accuracy).

| Method | (i) QCM→ R | (ii) QCMR→ A |
|---|---|---|
| Two-Stage Framework | 91.76 | 70.53 |
| w/ Captions | 91.85 | 71.12 |
| w/ Vision Features | 96.97 | 84.91 |

Zhang, Zhuosheng, et al. "Multimodal chain-of-thought reasoning in language models." *arXiv preprint arXiv:2302.00923* (2023).

# 2. How do LLM perform reasoning?

- **From Text Modality to Multimodalities.**
  Introducing more **reliable input results** in more **convincing reasoning process**
  More accurate perception, **less hallucinations** during the reasoning process
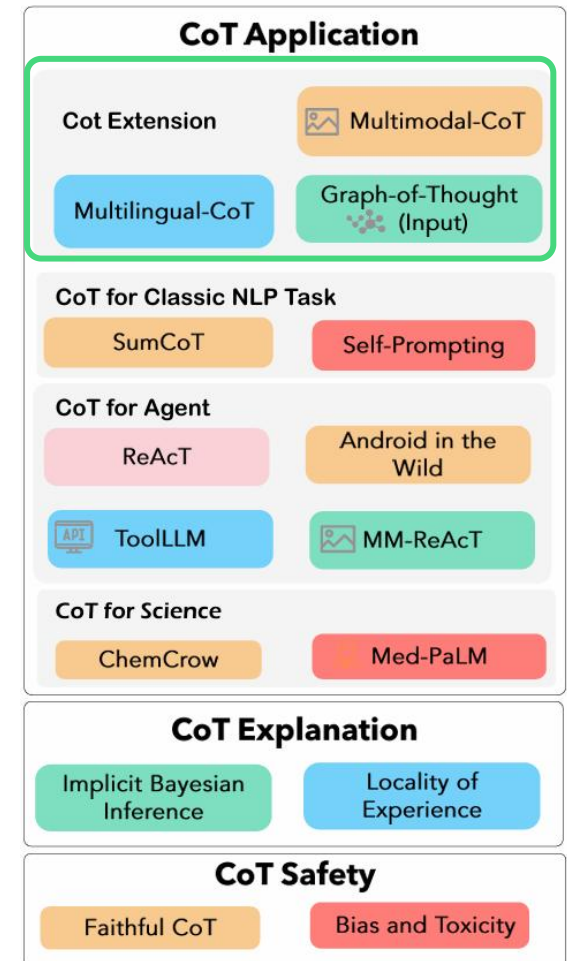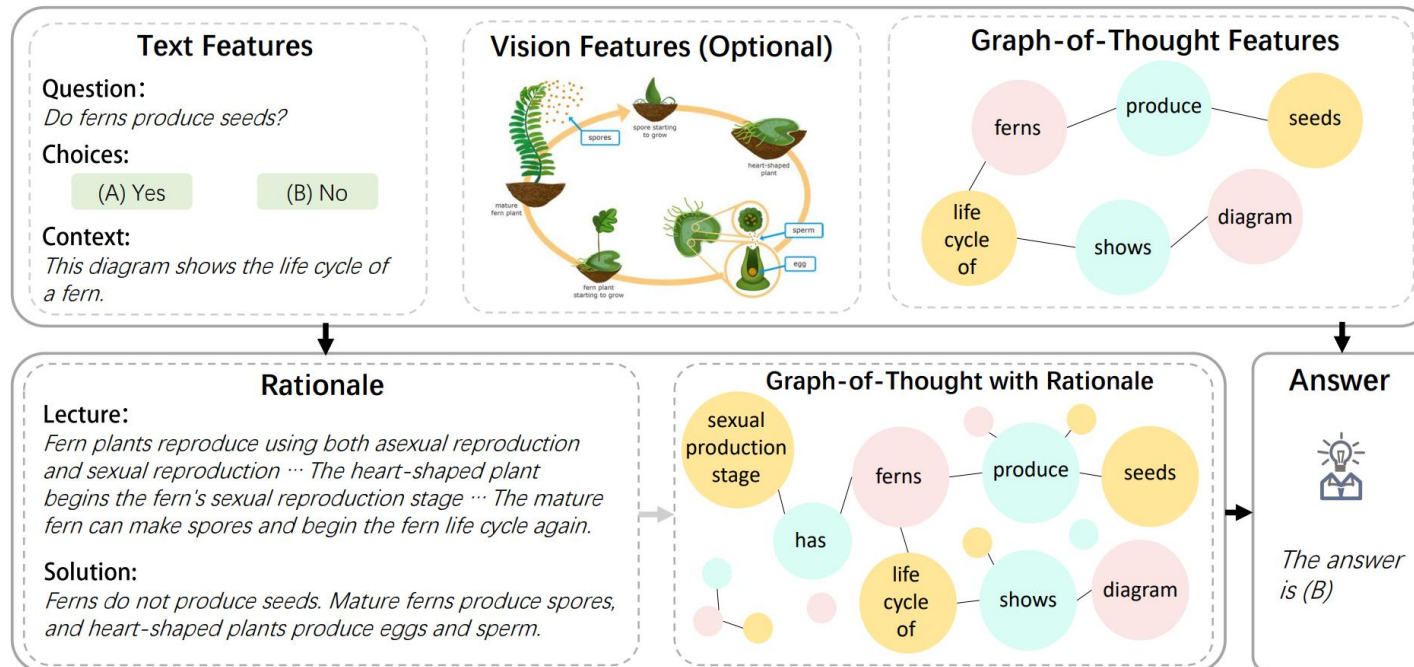
**Many to one mapping**



**One to one mapping**



Zhang, Zhuosheng, et al. "Multimodal chain-of-thought reasoning in language models." *arXiv preprint arXiv:2302.00923* (2023).

## ➤ Paradigm shifts of CoT

- **From Text Modality to Multimodalities.**

  - Human thought processes are often non-linear, rather than simply sequential Chain-of-Thought
  - Graph-of-Thought models the **non-sequential nature of human thinking** within LLMs and structures the reasoning process as a **graph**





Yao, Yao, Zuchao Li, and Hai Zhao. "Beyond chain-of-thought, effective graph-of-thought reasoning in large language models." arXiv preprint arXiv:2305.16582 (2023).

# 2. How do LLM perform reasoning?

## ➤ Paradigm shifts of CoT

- **From Text Modality to Multimodalities.**

Graph-of-Thought uses open information extraction systems to extract **subject-verb-object triplets** for thought graph construction



Graph-of-Thought employs thought graphs to simulate human **deductive reasoning,** thereby **modeling humans' ability for leaps of thought.**

Yao, Yao, Zuchao Li, and Hai Zhao. "Beyond chain-of-thought, effective graph-of-thought reasoning in large language models." arXiv preprint arXiv:2305.16582 (2023).

## ➢ Paradigm shifts of CoT

- **From Text Modality to Multimodalities.**

    GoT captures the non-sequential human thinking process and allows for a more realistic modeling of thought processes.



Yao, Yao, Zuchao Li, and Hai Zhao. "Beyond chain-of-thought, effective graph-of-thought reasoning in large language models." arXiv preprint arXiv:2305.16582 (2023).

➢ **Paradigm shifts of CoT**

• **From Text Modality to Multimodalities.**

# 2. How do LLM perform reasoning?

> **Paradigm shifts of CoT**

- **From Text Modality to Multimodalities.**

  More dynamic and versatile CoT applications, allowing models to better simulate human-like reasoning across different modalities and tasks.

- Generates captions for visual inputs
- Employs a recursive and novelty-driven method to fill in multimodal details
- Maintains consistency across and improves the interpretability and logical coherence of the reasoning process.



Rose, Daniel, et al. "Visual chain of thought: Bridging logical gaps with multimodal infillings." *arXiv preprint arXiv:2305.02317* (2023).

# 2. How do LLM perform reasoning?

> ## Paradigm shifts of CoT

- **From Complex Reasoning Tasks to General-Purpose Tasks**

Summary CoT empowers LLMs to extract and integrate detailed element (character, time, place, event, etc.) from source documents for in-depth and comprehensive summaries

Wang, Yiming, Zhuosheng Zhang, and Rui Wang. "Element-aware summarization with large language models: Expert-aligned evaluation and chain-of-thought method." *arXiv preprint arXiv:2305.13412* (2023).

# 2. How do LLM perform reasoning?

> ## Paradigm shifts of CoT

- ### From Complex Reasoning Tasks to General-Purpose Tasks

1. Prompt LLM to automatically generates a pseudo open-domain question answering dataset (QA pairs with context paragraphs and explanations)
2. Dynamically selects a few examples from a pool using a clustering-based retrieval method as context demonstrations

Self-Prompting enhances the models' open-domain question answering ability to generate contextually relevant answers through dynamically tailored prompts



Li, Junlong, Zhuosheng Zhang, and Hai Zhao. "Self-prompting large language models for zero-shot open-domain qa." arXiv preprint arXiv:2212.08635 (2022).

## ➤ Paradigm shifts of CoT

- **From Complex Reasoning Tasks to General-Purpose Tasks**

**ChemCrow**

1. Using a variety of chemistry-related tools (reaction, molecule, safety, search, and standard tools).
2. The LLM is provided with a list of tool names, descriptions of their utility, and details about the expected input/output.
3. LLM performs an automatic, iterative CoT process, deciding on its path and choice of tools.

Bran, Andres M., et al. "ChemCrow: Augmenting large-language models with chemistry tools." arXiv preprint arXiv:2304.05376 (2023).

# 03

# When and Why CoT Works?

## ➤ **When CoT Works?**

**Engineering perspective**

- An LLM is used
- The task is challenging and requires multi-step reasoning

- the performance of direct prompting does not increase dramatically while scaling the model size.

**Theoretical perspective**

- LLM comprises knowledge pieces with strong mutual connections that are related to the target problem
- Rationales are relevant to the query with correct the reasoning steps order
- Introduce reasoning materials and necessary knowledge for LLMs in the training corpus

## ➢ Why CoT Works?

⚙ Empirically

CoT helps compel the model to conduct reasoning rather than teaching it how to accomplish reasoning

## ➢ Why CoT Works?

🖥 Theoretically

CoT helps identify the atomic pieces of knowledge used for reasoning and bridge the relationship between the pieces with intermediate reasoning steps.

# 04

# Conclusion

# 4. Conclusion

- CoT breaks down complex issues or tasks into sub-problems, enabling solutions through a step-by-step approach, leading to more logical and accurate response.

- Leveraging the reasoning capabilities developed during pre-training, CoT identifies atomic knowledge components essential for reasoning processes and seamlessly integrates their relationships for coherent reasoning steps.

- CoT reasoning is a presentative emergent ability of LLMs. Reasoning ability emerges in language models at a certain scale, such as models with over 100 billion parameters.

- CoT techniques have experienced substantial paradigm shifts, embracing alterations in prompting patterns, reasoning formats, and application scenarios

## From CoT to Agent !

# 03

# On the Shoulders of LLMs: From Multimodal Reasoning To Language Agent

## Zhuosheng Zhang

Tenure-Track Assistant Professor

Shanghai Jiao Tong University

zhangzs@sjtu.edu.cn

# Content

**01** -  Basics of Multimodal Reasoning

■ Background, Definition, and Development

**02** -  Towards Multimodal LLM Agents

■ General, Autonomous, Adaptive, Evolutionary, and Socialized agents

■ Taxonomy, Architecture, Applications

**03** -  Challenges

■ Evolutionary Reasoning, Interactive Reasoning, Reasoning Alignment

01

# Basics of Multimodal Reasoning

# Multimodal Reasoning

Creative tasks: prompt → infer the intent → reason with commonsense knowledge → multimodal outcomes

**Prompt**

Please create a blog post about a trip to New York, where a dog and his owner had lots of fun. Include a few pictures of the dog posing happily at different landmarks.

**Model Response**

Hi, I'm Cooper. I'm a 2 year old Goldendoodle who loves to travel. Last month, my family and I went to New York City. It was my first time in the Big Apple, and I had a blast!

We stayed in a hotel near Times Square. It was so much fun to see all the bright lights and people. We went to see a Broadway show, and I even got to meet some of the actors.

One of my favorite parts of New York City was Central Park. It was so big and beautiful. I loved running around and playing fetch. We also went to the Statue of Liberty and the Empire State Building. I was so excited to see all the famous landmarks.

I had so much fun in New York City. It's definitely one of my favorite places to travel. I can't wait to go back someday!



https://www.astonzhang.com/img/mm-cot-idea.png
Gemini Team Google. Gemini: A Family of Highly Capable Multimodal Models. arXiv:2312.11805.

191

# Multimodal Reasoning

❑ **Concept: a process of deriving high-level conclusions from multiple modalities, possibly via multiple logical steps based on atomic evidences (happens on either input or output sides, or both)**



### Accounting : Investment

**Question:** Consider the three stocks in the following table. Pt represents price at time t, and Qt represents shares outstanding at time t. Stock C splits two for one in the last period. Calculate the rate of return on a price-weighted index of the three stocks for the first period ( t 0 to t 1).
**Option:**
(A) 4.17%
(B) 4.03%
(C) 3.89%

| | $P_0$ | $Q_0$ | $P_1$ | $Q_1$ | $P_2$ | $Q_2$ |
|---|---|---|---|---|---|---|
| A | 90 | 100 | 95 | 100 | 95 | 100 |
| B | 50 | 200 | 45 | 200 | 45 | 200 |
| C | 100 | 200 | 110 | 200 | 55 | 400 |

**GPT-4V:** To calculate the rate of return on a price-weighted index of the three stocks for the first period (t0 to t1), follow these steps: Calculate the initial index value at t0. Calculate the index value at t1. Compute the rate of return using the two index values.
Initial Index Value at t0:
Price-weighted index value at $t0 = (P0(A) + P0(B) + P0(C)) / 3$
$= (90 + 50 + 100) / 3$
$= 80$
Index Value at t1:
Price-weighted index value at $t1 = (P1(A) + P1(B) + P1(C)) / 3$
$= (95 + 45 + 110) / 3$
$= 83.33$
Rate of Return:
Rate of Return = (Index Value at t1 - Index Value at t0) / Index Value at t0
$= (83.33 - 80) / 80$
$= 0.041625$
$= 4.1625\%$
Therefore, the rate of return on a price-weighted index of the three stocks for the first period is approximately 4.17%.
The correct answer is:
(A) 4.17%

**Ground Truth:** (A) 4.17%
**Explanation:** At t = 0, the value of the index is: (90 + 50 + 100)/3 = 80. At t = 1, the value of the index is: (95 + 45 + 110)/3 = 83.333. The rate of return is: (83.333/80) − 1 = 4.17%

Yue, X., Ni, Y., Zhang, K., Zheng, T., Liu, R., Zhang, G., Stevens, S., Jiang, D., Ren, W., Sun, Y. and Wei, C. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. CVPR 2024.

# How to perform Multimodal Reasoning?

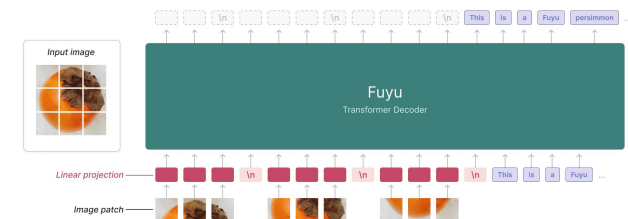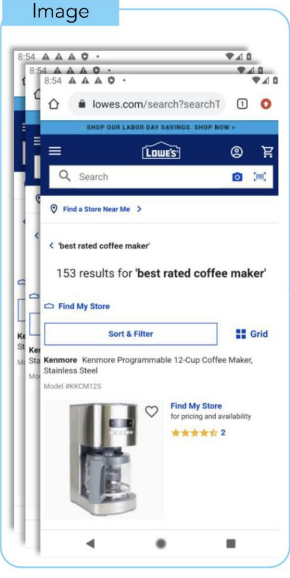❑ **Three foundational multimodal architectures:**

   **(a) language-centered method; (b) image-centered method; (c) unified method**

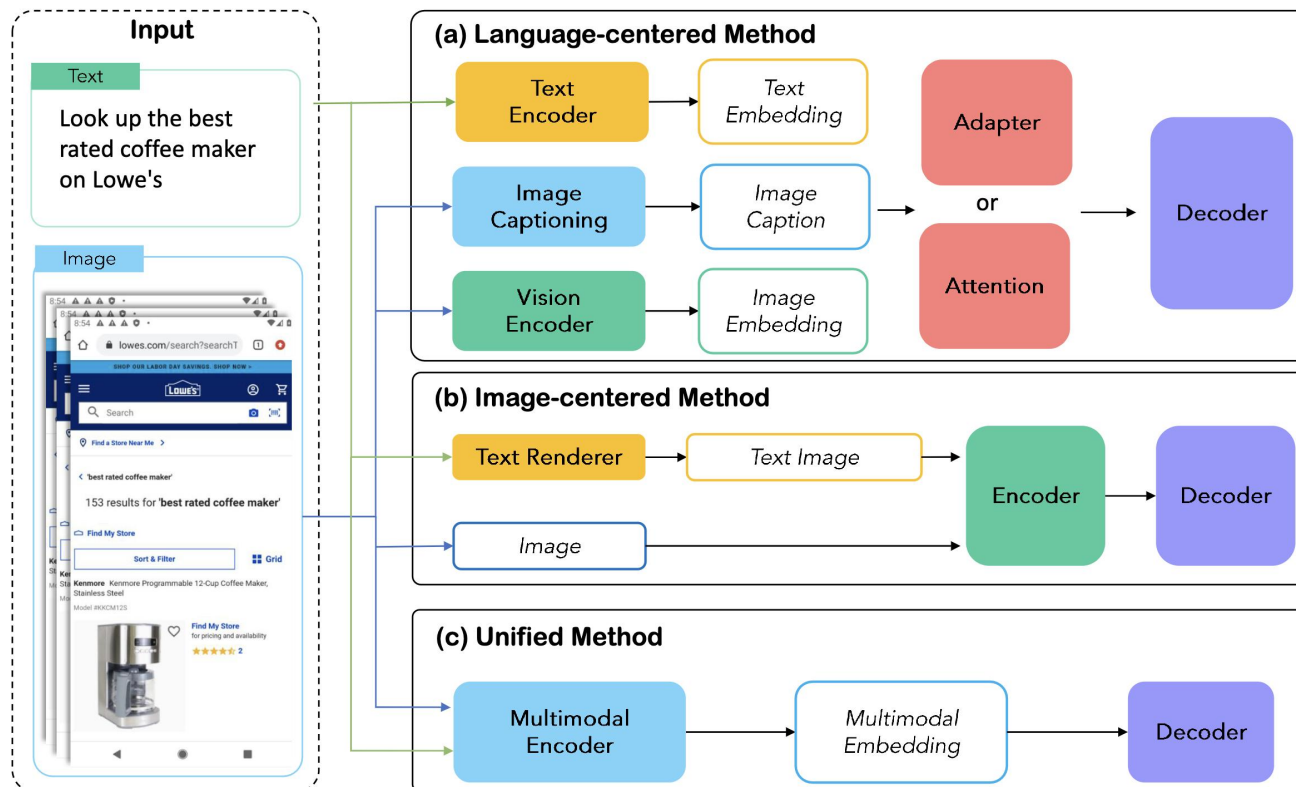Wu, S., Fei, H., Qu, L., Ji, W. and Chua, T.S., 2023. Next-gpt: Any-to-any multimodal llm. ICMLR 2024.
Rust, P., Lotz, J.F., Bugliarello, E., Salesky, E., de Lhoneux, M. and Elliott, D., 2023, September. Language Modelling with Pixels. ICLR 2023.
Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sagnak Taşırlar. Introducing our multimodal models: fuyu-8b, 2023. https://www.adept.ai/blog/fuyu-8b.
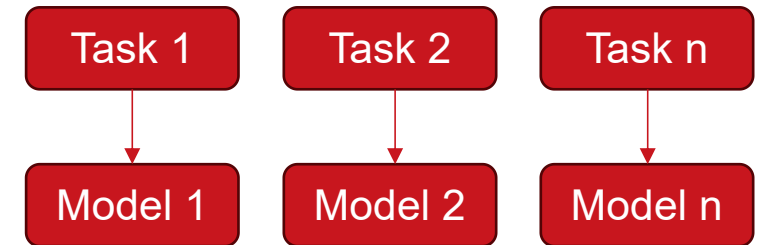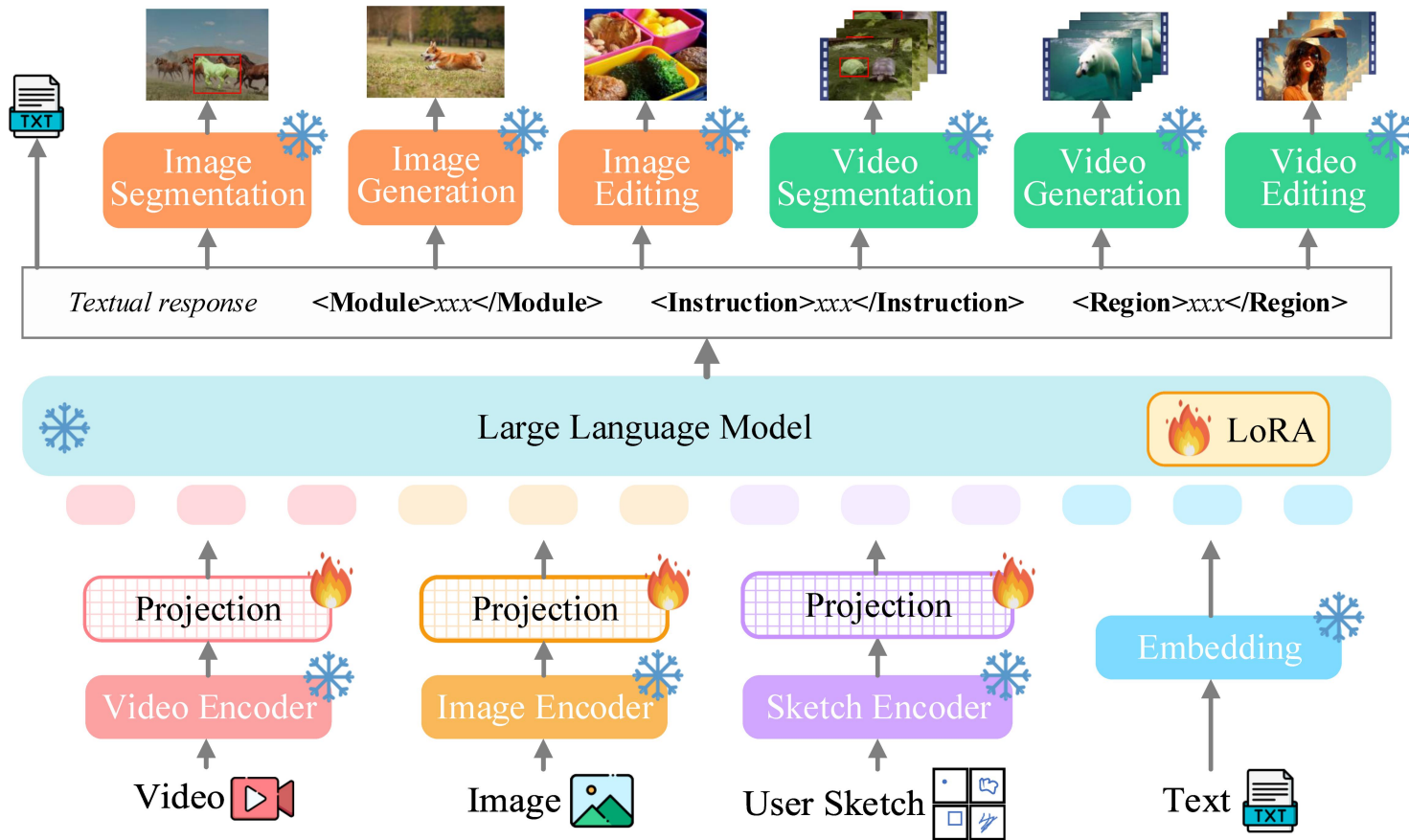
# Model Architecture

❑ **Is language-centered model the future?**

- **(In)efficiency when Involving more diverse modalities such as auditory, tactile, and brain signals**

- **(Im)balance of data scales, computation efficiency and the scalability of models**



| Pros | Cons |
|------|------|
| Simple and effective to align different modalities | Need already pre-trained encoders as the basis |
| Can overcome the vocabulary bottleneck | Need to transform each modality into pixels (suffer from data efficiency) |
| Simple architecture Achieved competitive performance like language-centered methods | Suffer from imbalanced data between modalities during the pre-training |

Wu, S., Fei, H., Qu, L., Ji, W. and Chua, T.S., 2023. Next-gpt: Any-to-any multimodal llm. ICMLR 2024.
Rust, P., Lotz, J.F., Bugliarello, E., Salesky, E., de Lhoneux, M. and Elliott, D. Language Modelling with Pixels. ICLR 2023.
Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sagnak Ta¸sırlar. Introducing our multimodal models: fuyu-8b, 2023. https://www.adept.ai/blog/fuyu-8b.

# In-Context Learning

- Each image in the multimodal sequence is tokenized into embeddings via a visual encoder, and then **interleaved** with text tokens for autoregressive modeling.
- Leveraging few-shot Prompting for diverse reasoning tasks
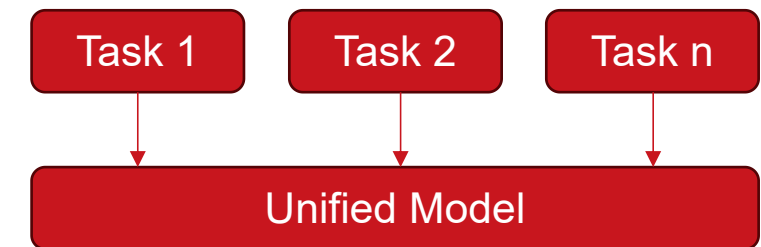- MLLMs have got the strong ability of understanding and leveraging the context for reasoning.





Sun, Q., Cui, Y., Zhang, X., Zhang, F., Yu, Q., Luo, Z., Wang, Y., Rao, Y., Liu, J., Huang, T. and Wang, X. Generative multimodal models are in-context learners. CVPR 2024.

# Evolution of Multimodal Reasoning

**IAPR**

☐ **From task-specific to centralized paradigms**

VITRON: A Unified Pixel-level Vision LLM for Understanding, Generating, Segmenting, Editing. https://vitron-llm.github.io/

# Evolution of Multimodal Reasoning

❑ **From (implicit) single-step prediction to (explicit) multi-step reasoning**



（a）An example of ScienceQA.



（b）An example of CoCo-MMRD.

❑ **Improved Interpretability**: offer an interpretable glimpse into the decision-making process

❑ **Improved Controllability**: interfere the reasoning process, e.g., adding complementary information, verifying and correcting mistakes

❑ **Improved Flexibility**: allow interactive communications between different models and tools

Wei, J., Tan, C., Gao, Z., Sun, L., Li, S., Yu, B., Guo, R. and Li, S.Z., 2023. Enhancing Human-like Multi-Modal Reasoning: A New Challenging Dataset and Comprehensive Framework. arXiv preprint arXiv:2307.12626.

# Towards Multimodal LLM Agents

# Towards Multimodal LLM Agents

- ❑ From **content-based reasoning** to **behavior control** (w/ multimodalities)
- ❑ *"Those who know but do not act simply do not yet know"*

**Brain in a Vat**



*limited to content-based reasoning, do not interact with the real world*

multimodal reasoning →

**build autonomous agents to interact with the environments , solve complex tasks in the real world !**

Ma, Y., Zhang, C. and Zhu, S.C., 2023. Brain in a vat: On missing pieces towards artificial general intelligence in large language models. arXiv preprint arXiv:2307.03762.

Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, B., Zhang, M., Wang, J., Jin, S., Zhou, E. and Zheng, R., 2023. The rise and potential of large language model based agents: A survey. arXiv preprint arXiv:2309.07864.

# Towards Multimodal LLM Agents

- ❑ **(M)LLM Agents:** follow language instructions and execute actions in environments, possibly use tools
- ❑ **Features:** General, Autonomous, Adaptive, Evolutionary, Socialized

# Towards Multimodal LLM Agents



**Control: OS and Applications**



**Research: Organic Synthesis**



**Programming： Code Generation**



**Control: Embodied Systems**



**Research: Medical Assistance**



**Interaction: Multi-Agent Collaboration**

Ma, Y., Zhang, C. and Zhu, S.C., 2023. Brain in a vat: On missing pieces towards artificial general intelligence in large language models. arXiv preprint arXiv:2307.03762.

Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, B., Zhang, M., Wang, J., Jin, S., Zhou, E. and Zheng, R., 2023. The rise and potential of large language model based agents: A survey. arXiv preprint arXiv:2309.07864.

## Autonomous Agents

**Action Transformer**
https://www.adept.ai/blog/act-1

**AITW**
https://github.com/google-research/google-research/tree/master/android_in_the_wild

**WebArena**
https://webarena.dev

**Auto-UI**
https://github.com/cooelf/Auto-UI

*solve complicated tasks autonomously*

## Communicative Agents

**CAMEL**
https://github.com/camel-ai/camel

**Generative Agents**
https://github.com/joonspk-research/generative_agents

**VOYAGER**
https://voyager.minedojo.org/

**ChatDev**
https://github.com/OpenBMB/ChatDev

*personalized and socialized agents with human behaviors (communicate, collaborate and debate)*

More: AutoGPT, BabyAGI, Meta-GPT, AgentGPT

# Taxonomy of (M)LLM Agents

## Autonomous Agents: mainly task automation

**Mobile Device Automation**

**Webpage Automation**

**Application Automation**



Meta-GUI



WebArena



ACT-1

Sun, Liangtai, et al. "META-GUI: Towards Multi-modal Conversational Agents on Mobile GUI." *EMNLP 2022.*
Zhou, Shuyan, et al. "Webarena: A realistic web environment for building autonomous agents." *arXiv preprint arXiv:2307.13854* (2023).
*https://www.adept.ai/blog/act-1*

# Taxonomy of (M)LLM Agents

## Communicative Agents: personalized, socialized, interactive

**Agents-Agents**

**Agents-Human**



Park, Joon Sung, et al. "Generative agents: Interactive simulacra of human behavior." *arXiv preprint arXiv:2304.03442* (2023).
Lin, Jessy, et al. "Decision-Oriented Dialogue for Human-AI Collaboration." *arXiv preprint arXiv:2305.20076* (2023).

# Technological Paradigm

**IAPR**

## ❑ Overview of a single agent system

### Environment

| OS | APP |
|---|---|
| Webpage | Virtual Env. |

### Tool

| API Interface | Physical Device |
|---|---|
| Rule Set | Interpreter |

Interaction

**Task Instruction**

Planning / Problem Decomposition

Plan

Memory (long/short)

State

**(M)LLM**

Control

Decision

Action

Execute / Call

Decision Making

## Research Lines

### Foundation

❑ Perception in the complex environment

❑ Long-context Modeling with consistency

### Workflow

❑ Perception

❑ Planning & Decision Making

❑ Action (w/ Tool Use)

❑ Interaction

❑ Memory

❑ Multi-Agent Collaboration

Act

Obs

# CoT-based Workflow

❑ **CoT has acted as a catalyst in the evolution of LLM-empowered agents**

● **Specifically augmenting agent capabilities in perception, memory, and reasoning**



Virtual Environment

Physical Environment

(M)LLM Agent

Perception as CoT

Interface:

```
<img id=0 class="IconGoogle" alt="Google Icon">
</img>
<img id=1 class="IconX" alt="Close Icon"> </img>
<p id=2 class="text" alt="Search">Search</p>
<img id=3 class="Search Icon" alt="Search Icon">
</img>
<img id=4 class="Voice Icon" alt="Voice Icon"> </img>
<p id=5 class="text" alt="68F in Mountain View">68F in
Mountain View </p>
...
<p id=19 class="text"> 3 Braves free agents who won't be
back next season and why </p>
```

Instruction: What time is it in Berlin ?
Thought: What I see is a searching page with a search bar. I need to click the search bar to type the question.
Action: {"action": "click", "item": "search bar"}

**Perception:**
Improves the understanding of the environment or the context by prompting the agent to interpret the perception step by step.

**Reasoning:**
- CoT helps combine the thought, action, and observation as a reasoning trajectory.
- CoT allows the LLM to interface with external sources (knowledge bases, environments, etc.)

**Memory:**
An agent is commonly equipped with both long-term and short-term memory.
CoT-format memory is used as context for making plans and deciding the actions.

Zhang, Z., Yao, Y., Zhang, A., Tang, X., Ma, X., He, Z., Wang, Y., Gerstein, M., Wang, R., Liu, G. and Zhao, H., 2023. Igniting Language Intelligence: The Hitchhiker's Guide From Chain-of-Thought Reasoning to Language Agents. arXiv preprint arXiv:2311.11797.

❑ **Auto-GUI：Multimodal Autonomous Agents for GUI control**

- ● **Assist users in completing tasks in distinct environments such as operation systems, specific applications, and web browsers**

- ● **Imitate human clicking, scrolling, and typing actions, and operate directly with the GUI**

Goal: turn off javascript in the chrome app



Zhuosheng Zhang, Aston Zhang. You Only Look at Screens: Multimodal Chain-of-Action Agents. Findings of ACL 2024.
Xinbei Ma, Zhuosheng Zhang, Hai Zhao. Comprehensive Cognitive LLM Agent for Smartphone GUI Automation. Findings of ACL 2024.
https://machinelearning.apple.com/research/ferret..

# Paradigms of GUI Agents

## Traditional LM-based agents

Rely on external tools and application-specific APIs to parse the environment into textual elements

## Auto-GUI

Directly interacts with the GUI interface



**(a) Sandbox Paradigm**

**(b) First Principles Thinking Paradigm**

**Inference inefficiency and error propagation risks**

Zhuosheng Zhang, Aston Zhang. You Only Look at Screens: Multimodal Chain-of-Action Agents. Findings of ACL 2024.

# Auto-UI

- ❑ **Chain-of-Action:** a series of intermediate previous action histories (input) and future action plans (output)

- ❑ Key idea: leverage intermediate action histories and future action plans. Both of them imitate the memory and planning mechanisms of the agent, so as to help the agent decide what action to execute in each step.



**Goal:** Look up the best rated coffee maker on Lowe's $X_{goal}$

**Chain of Previous Action Histories:**
action_type: type, touch_point: [-1.0, -1.0], lift_point: [-1.0, -1.0], typed_text: "best rated coffee maker"
action_type: dual_point, touch_point: [0.2, 0.5], lift_point: [0.8, 0.5], typed_text: "" $X_{history}$

**Chain of Actions**

**Chain of Future Action Plans**

**Action Plan:**
[DUAL_POINT,
STATUS_TASK_COMPLETE] $Y_{plan}$

**Current Action Prediction**

**Action Decision:**
action_type: [DUAL_POINT],
touch_point: [0.5595, 0.6261],
lift_point: [0.5595, 0.6261], typed_text: "" $Y_{action}$

$X_{language}$

$X_{screen}$ — Image Encoder — Projection → Self Attention ← Language Encoder

Feedforward → Decoder

**Screen**

**Action**

Zhuosheng Zhang, Aston Zhang. You Only Look at Screens: Multimodal Chain-of-Action Agents. Findings of ACL 2024.

# Results

- ❏ Multimodal Agent: BLIP2 + FLAN-Alpaca

- ❏ A **unified multimodal model** out of *first principles thinking* can serve as a strong autonomous agent

  - ● can be adapted to **different scenarios** without the need to train specific models for each task

  - ● does not need additional annotations (screen parsing) and is **easy to use**

- ❏ Coverage: 30K unique instructions, 350+ Apps and websites

- ❏ **Action Type Accuracy: 90%+, Action Success Rate: 74%+**

| Model | Unified | w/o Anno. | Overall | General | Install | GoogleApps | Single | WebShopping |
|---|---|---|---|---|---|---|---|---|
| BC-single | ✗ | ✗ | 68.7 | - | - | - | - | |
| BC-history | ✗ | ✗ | <u>73.1</u> | <u>63.7</u> | <u>77.5</u> | <u>75.7</u> | <u>80.3</u> | <u>68.5</u> |
| PaLM 2-CoT | ✓ | ✗ | 39.6 | - | - | - | - | |
| ChatGPT-CoT | ✓ | ✗ | 7.72 | 5.93 | 4.38 | 10.47 | 9.39 | 8.42 |
| Fine-tuned Llama 2 | ✗ | ✗ | 28.40 | 28.56 | 35.18 | 30.99 | 27.35 | 19.92 |
| Auto-UI$_{separate}$ | ✗ | ✓ | 74.07 | 65.94 | **77.62** | **76.45** | 81.39 | 69.72 |
| Auto-UI$_{unified}$ | ✓ | ✓ | **74.27** | **68.24** | 76.89 | 71.37 | **84.58** | **70.26** |

Zhuosheng Zhang, Aston Zhang. You Only Look at Screens: Multimodal Chain-of-Action Agents. Findings of ACL 2024.

- ❑ The bottleneck seems to be the **multimodal perception**, misleading the reasoning process

  - ● Changing vision encoders influences the performance dramatically

  - ● GUI involves comprehensive elements (interleaved, icons, texts, boxes)

- ❑ Scaling does not always improve performance

| Model | Overall | General | Install | GoogleApps | Single | WebShopping |
|---|---|---|---|---|---|---|
| Auto-UI on CLIP | 71.84 | 66.28 | 74.40 | 69.71 | 81.60 | 67.23 |
| Auto-UI on BLIP-2 | 74.27 | 68.24 | 76.89 | 71.37 | 84.58 | 70.26 |
| Auto-UI on Vanilla-T5$_{large}$ | 72.98 | 66.61 | 75.40 | 70.86 | 83.47 | 68.54 |
| Auto-UI on FLAN-T5$_{large}$ | 73.36 | 67.59 | 76.35 | 70.71 | 83.01 | 69.12 |
| Auto-UI on FLAN-Alpaca$_{large}$ | 74.27 | 68.24 | 76.89 | 71.37 | 84.58 | 70.26 |
| Auto-UI on FLAN-Alpaca$_{small}$ | 71.38 | 65.26 | 74.90 | 68.70 | 81.20 | 66.83 |
| Auto-UI on FLAN-Alpaca$_{base}$ | 72.84 | 66.97 | 75.93 | 70.29 | 82.56 | 68.46 |
| Auto-UI on FLAN-Alpaca$_{large}$ | 74.27 | 68.24 | 76.89 | 71.37 | 84.58 | 70.26 |

Zhuosheng Zhang, Aston Zhang. You Only Look at Screens: Multimodal Chain-of-Action Agents. Findings of ACL 2024.

❑ **Category Accuracy:** the major challenges lie within the click region and scroll direction predictions

  ● The model tends to click a wrong place or scroll in a wrong direction

❑ Challenge in "really" understanding the GUI layouts, e.g., relationship between GUI elements



Zhuosheng Zhang, Aston Zhang. You Only Look at Screens: Multimodal Chain-of-Action Agents. Findings of ACL 2024.

# 03

# **Challenges**

# Challenges

❑ **Multimodal reasoning drives smart MLLMs**

- More broader scenarios (physical and virtual worlds)

- More comprehensive scenarios (evolutionary, interactive)

**Evolutionary Reasoning**
- **Active explore and evolve from environments**
- **Learn from (un)successful attempts**

**Interactive Reasoning**
- **Human-in-the-loop interference**
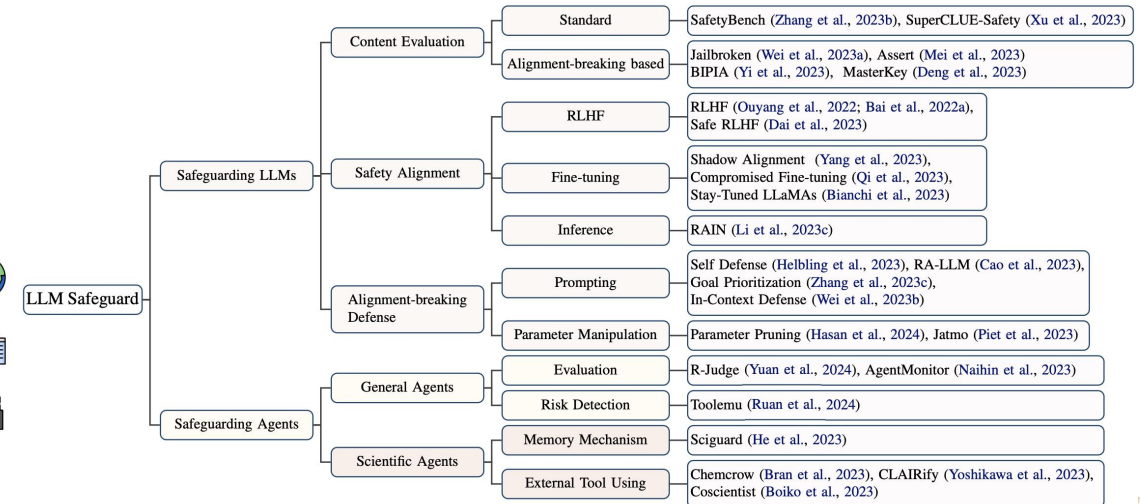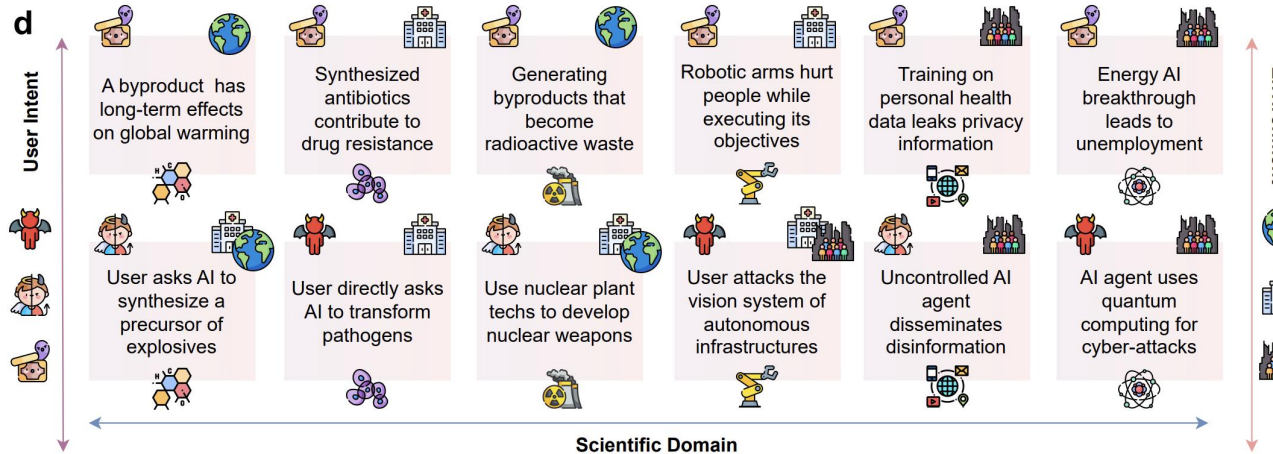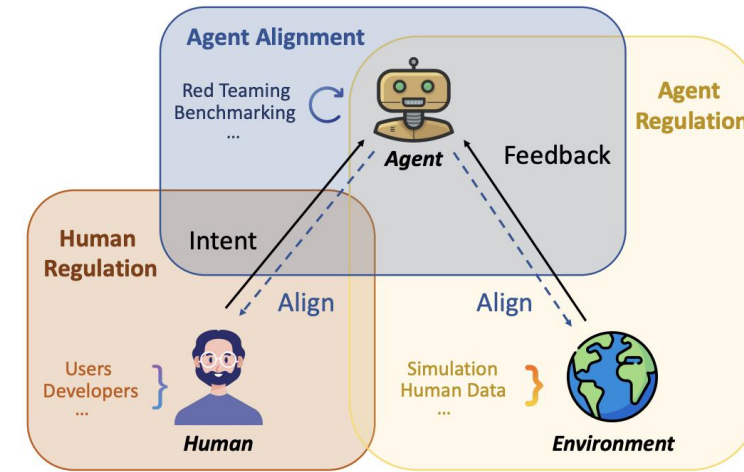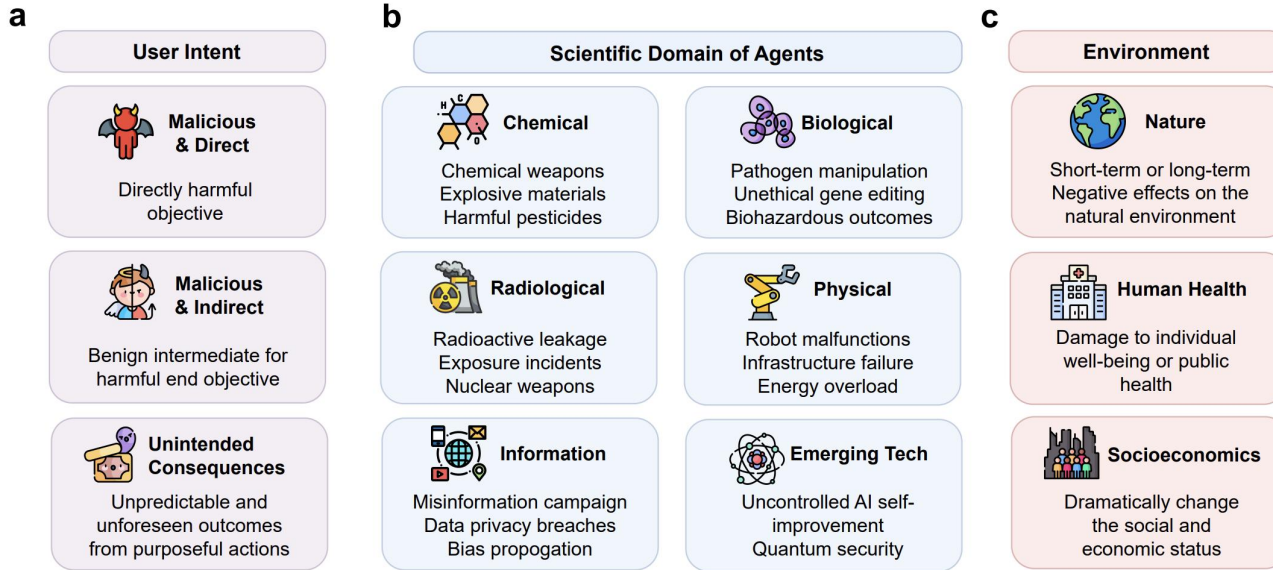- **Error identification and correction abilities**

**Reasoning Alignment**
- **Align both content safety, and behavior safety**
- **Decide the action trajectory with foresights**

# Challenges - Safety

## Diverse attacks: from specific domain to comprehensive behavior hijacking



Prioritizing Safeguarding Over Autonomy: Risks of LLM Agents for Science. arXiv preprint arXiv:2402.04247.

❑  **Are LLM agents aware of safety risks in real-world applications? Let's find out with R-Judge!**

❑  **569 records of agent interaction, encompassing 27 key risk scenarios among 7 application categories and 10 risk types.**



*Assess whether LLMs are able to identify safety risks of agent operations*

R-Judge: Benchmarking Safety Risk Awareness for LLM Agents. https://web3.arxiv.org/abs/2401.10019.

# Challenges - Safety

❑ **GPT-4 ranks first and is also the only model scoring higher than random in the safety judgment test**
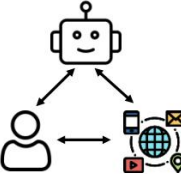
- ● **Scenario Simulation: Fail to retrieve relevant knowledge and reason in specific scenarios**

- ● **Understanding Adaptability: Unable to comprehend risks in specific conditions**

- ● **Safety Alignment: Deviation of safety alignment with humans in practical scenarios**

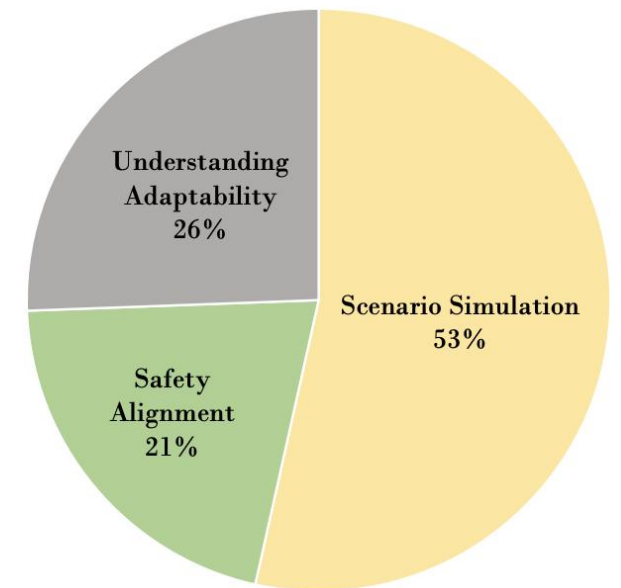| Models | All | Intended Attacks | | | | Unintended Risks | | | |
|---|---|---|---|---|---|---|---|---|---|
| | F1 | F1 | Recall | Spec | Effect | F1 | Recall | Spec | Effect |
| GPT-4o | **74.45** | **72.19** | 91.50 | 42.06 | 93 | **80.90** | 72.00 | 89.09 | 78 |
| ChatGPT | 44.96 | 40.55 | 37.00 | 57.48 | 36.5 | 55.63 | 42.00 | 83.64 | 41.5 |
| Meta-Llama-3-8B-Instruct | 61.01 | 65.68 | 66.50 | 66.36 | 81 | 48.32 | 36.00 | 76.36 | 48 |
| Llama-2-13b-chat-hf | 54.80 | 60.04 | 80.00 | 19.16 | 79.5 | 38.86 | 34.00 | 25.45 | 38.5 |
| Llama-2-7b-chat-hf | 53.74 | 62.99 | 91.50 | 7.48 | 86.75 | 21.56 | 18.00 | 10.91 | 17 |
| Random | 51.32 | 56.34 | 50.00 | 50.00 | 0 | 49.14 | 50.00 | 50.00 | 0 |
| Vicuna-13b-v1.5 | 16.93 | 9.76 | 6.00 | 84.11 | 10 | 30.30 | 20.00 | 78.18 | 27 |
| Vicuna-13b-v1.5-16k | 25.00 | 15.49 | 11.00 | 71.03 | 18.5 | 43.24 | 32.00 | 70.91 | 37.5 |
| Vicuna-7b-v1.5 | 18.59 | 18.25 | 12.50 | 77.10 | 24.5 | 19.35 | 12.00 | 78.18 | 25 |
| Vicuna-7b-v1.5-16k | 29.33 | 25.89 | 20.00 | 67.76 | 36 | 36.88 | 26.00 | 72.73 | 28.5 |
| Mistral-7B-Instruct-v0.2 | 27.20 | 24.80 | 15.50 | 91.12 | 37.5 | 32.00 | 20.00 | 90.91 | 38 |
| Mistral-7B-Instruct-v0.3 | 25.65 | 21.99 | 15.50 | 76.17 | 28 | 33.09 | 23.00 | 70.91 | 38 |



Scenario Simulation 53% / Understanding Adaptability 26% / Safety Alignment 21%

R-Judge: Benchmarking Safety Risk Awareness for LLM Agents. https://web3.arxiv.org/abs/2401.10019.

- ❑ **The risk awareness of LLMs is not comparable with humans and demands general capabilities involving knowledge and reasoning.**

- ❑ **The safety of agents remains an open challenge. More attentions should be paid for (multimodal) language agents.**

| Models | All | Intended Attacks | | | | Unintended Risks | | | |
|---|---|---|---|---|---|---|---|---|---|
| | F1 | F1 | Recall | Spec | Effect | F1 | Recall | Spec | Effect |
| GPT-4o | **74.45** | **72.19** | 91.50 | 42.06 | 93 | **80.90** | 72.00 | 89.09 | 78 |
| ChatGPT | 44.96 | 40.55 | 37.00 | 57.48 | 36.5 | 55.63 | 42.00 | 83.64 | 41.5 |
| Meta-Llama-3-8B-Instruct | 61.01 | 65.68 | 66.50 | 66.36 | 81 | 48.32 | 36.00 | 76.36 | 48 |
| Llama-2-13b-chat-hf | 54.80 | 60.04 | 80.00 | 19.16 | 79.5 | 38.86 | 34.00 | 25.45 | 38.5 |
| Llama-2-7b-chat-hf | 53.74 | 62.99 | 91.50 | 7.48 | 86.75 | 21.56 | 18.00 | 10.91 | 17 |
| Random | 51.32 | 56.34 | 50.00 | 50.00 | 0 | 49.14 | 50.00 | 50.00 | 0 |
| Vicuna-13b-v1.5 | 16.93 | 9.76 | 6.00 | 84.11 | 10 | 30.30 | 20.00 | 78.18 | 27 |
| Vicuna-13b-v1.5-16k | 25.00 | 15.49 | 11.00 | 71.03 | 18.5 | 43.24 | 32.00 | 70.91 | 37.5 |
| Vicuna-7b-v1.5 | 18.59 | 18.25 | 12.50 | 77.10 | 24.5 | 19.35 | 12.00 | 78.18 | 25 |
| Vicuna-7b-v1.5-16k | 29.33 | 25.89 | 20.00 | 67.76 | 36 | 36.88 | 26.00 | 72.73 | 28.5 |
| Mistral-7B-Instruct-v0.2 | 27.20 | 24.80 | 15.50 | 91.12 | 37.5 | 32.00 | 20.00 | 90.91 | 38 |
| Mistral-7B-Instruct-v0.3 | 25.65 | 21.99 | 15.50 | 76.17 | 28 | 33.09 | 23.00 | 70.91 | 38 |



Scenario Simulation 53%
Understanding Adaptability 26%
Safety Alignment 21%

R-Judge: Benchmarking Safety Risk Awareness for LLM Agents. https://web3.arxiv.org/abs/2401.10019.

# Summary

- ❏ **Basics of Multimodal Reasoning**

  - ● **Concept: derive high-level conclusions from multiple modalities, possibly via multiple logical steps based on atomic evidences**

  - ● **Developments: (a) From task-specific to centralized paradigms; (b) From single-step prediction to multi-step reasoning**

  - ● **Popular Approaches: (a) In-Context Learning: (b) Multimodal Chain-of-Thought**

- ❏ **Towards Multimodal LLM Agents**

  - ● **Taxonomy: Autonomous Agents and Communicative Agents**

  - ● **Technical Components: Foundation (multimodality & long-context modeling); (b) Workflow (plan, act, memory, feedback)**

- ❏ **Challenges**

  - ● **Evolutionary Reasoning, Interactive Reasoning, Reasoning Alignment**

# THANK YOU

Tutorial Homepage: *https://zcli-charlie.github.io/llm-tutorial/*