

AI 驱动 软件研发 全面进入数字化时代

中国·北京 08.18-19

AI+
software
Development
Digital
summit



ALM关键技术与应用范式

王昊奋 同济大学

科技生态圈峰会 + 深度研习 —— 1000+ 技术团队的选择



2023K+
全球软件研发行业创新峰会
上海站

会议时间 | 06.09-10



2023K+
全球软件研发行业创新峰会
北京站

会议时间 | 07.21-22



2024K+
全球软件研发行业创新峰会
深圳站

会议时间 | 05.17-18



K+峰会详情



会议时间 | 08.18-19

AiDD AI+软件研发数字峰会
北京站



会议时间 | 11.17-18

AiDD AI+软件研发数字峰会
深圳站



AiDD峰会详情

▶ 演讲嘉宾



王昊奋

OpenKG发起人/CCF知识图谱SIG主席

同济大学百人计划特聘研究员、OpenKG发起人之一，业界知名知识图谱与自然语言处理专家，腾讯云最具价值专家TVP，CCF术语工委副主任、知识图谱SIG主席、上海秘书长，中国中文信息学会理事，畅销书《知识图谱方法、实践与应用》的作者，曾作为2家AI独角兽企业的CTO；具有超过16年的知识图谱研发和技术管理经验。受邀在世界人工智能大会等诸多国际与国内智能峰会上担任讲者，并在自然语言处理国内顶级会议NLPCC多次担任知识图谱方向主席，长期作为ISWC, WWW, AAAI等人工智能国际顶级会议程序委员会委员。

目录

CONTENTS

1. LLM的崛起
2. ALM关键技术
3. 应用落地范式
4. 未来展望

PART 01 **LLM的崛起**

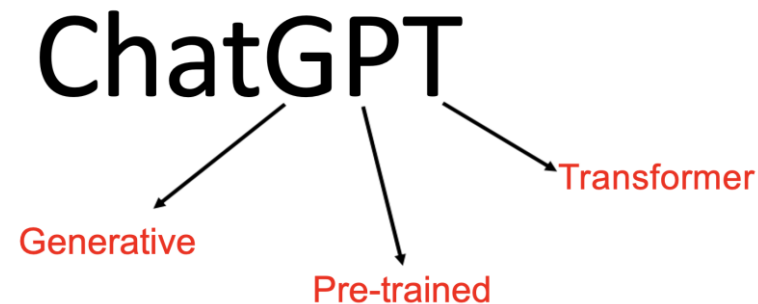


▶ 大语言模型掀起迈向AGI的浪潮

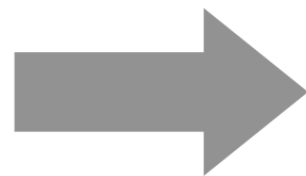
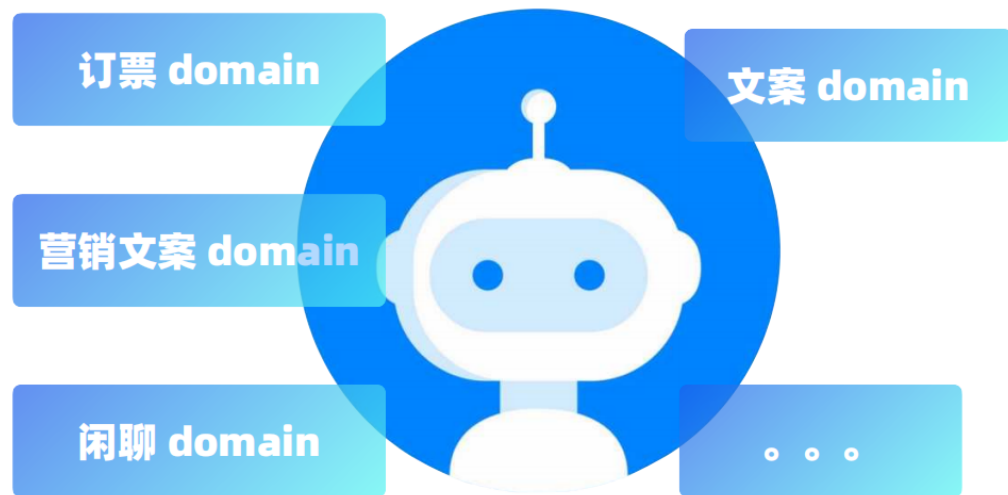
ChatGPT		
Examples	Capabilities	Limitations
"Explain quantum computing in simple terms" →	Remembers what user said earlier in the conversation	May occasionally generate incorrect information
"Got any creative ideas for a 10 year old's birthday?" →	Allows user to provide follow-up corrections	May occasionally produce harmful instructions or biased content
"How do I make an HTTP request in Javascript?" →	Trained to decline inappropriate requests	Limited knowledge of world and events after 2021


ChatGPT是由美国人工智能公司OpenAI在2022年11月推出的**生成式对话预训练大语言模型**。

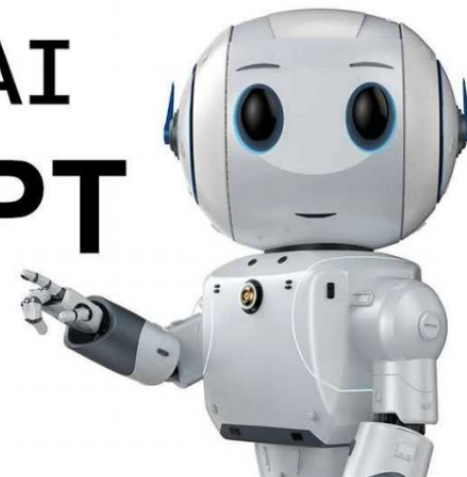
它以对话的方式进行交互。对话形式使得其**能够回答后续问题，承认错误，质疑不正确的前提，并拒绝不适当的请求**



▶ 2016年的Chatbot VS ChatGPT下的Chatbot



 OpenAI
ChatGPT



域(domain)、意图(intent)、词槽(slot)定义智能的年代

大语言模型 + Prompting
通用型人工智能



siri



小爱同学



小度



▶ 大语言模型掀起迈向AGI的浪潮

NLP任务:

- 文本分类
- 信息抽取
- 文本摘要
- 智能问答
- 阅读理解
- 机器翻译
- 文本生成
- 语法纠正
-

应用场景:

- 信息分类
- 文本结构化
- 摘要说明
- 对话问答
- 复杂文本理解
- 多种语言翻译
- 内容创作
- 信息纠错
-

句子简化 TL;DR summarization Summarize text by adding a 'tl;dr:' to the en...	修复代码Bug Python bug fixer Find and fix bugs in source code.	文字转表情符号 Movie to Emoji Convert movie titles into emoji.
表格填充数据 Spreadsheet creator Create spreadsheets of various kinds of dat...	语言聊天机器人 JavaScript helper chatbot Message-style bot that answers JavaScript ...	程序代码翻译 Translate programming languages Translate from one programming language ...
机器学习机器人 ML/AI language model tutor	清单制作 Science fiction book list maker	代码解释 Explain code
问&答 Q&A Answer questions based on existing knowle...	语法纠正 Grammar correction Corrects sentences into standard English.	Python代码解释 Python to natural language Explain a piece of Python code in human un...
内容概况 Summarize for a 2nd grader Translates difficult text into simpler concep...	生成OpenAi的代码 Natural language to OpenAI API Create code to call to the OpenAI API usin...	时间复杂度计算 Calculate Time Complexity Find the time complexity of a function.
程序命令生成 Text to command Translate text into programmatic commands.	语言翻译 English to other languages Translates English text into French, Spanish...	高级情绪评分 Advanced tweet classifier Advanced sentiment detection for a piece o...
Stripe国际API生成 Natural language to Stripe API	SQL语句生成 SQL translate	关键字提取 Keywords
程序语言转换 JavaScript to Python Convert simple JavaScript expressions into ...	好友聊天 Friend chat Emulate a text message conversation.	ESRB文本分类 ESRB rating Categorize text based upon ESRB ratings.
颜色生成 Mood to color Turn a text description into a color.	程序文档生成 Write a Python docstring An example of how to create a docstring for ...	美食制作 (后果自负) Recipe creator (eat at your own risk) Create a recipe from a list of ingredients.
段落创作 Analogy maker Create analogies. Modified from a communi...	代码压缩 JavaScript one line function Turn a JavaScript function into a one liner.	摆烂聊天 Marv the sarcastic chat bot Marv is a factual chatbot that is also sarcas...
故事创作 Micro horror story creator Creates two to three sentence short horror ...	人称转换 Third-person converter Converts first-person POV to the third-pers...	点评生成 Restaurant review creator Turn a few words into a restaurant review.
摘要说明 Notes to summary Turn meeting notes into a summary.	头脑风暴 VR fitness idea generator Create ideas for fitness and virtual reality g...	面试 Interview questions Create interview questions.

▶ 大语言模型掀起迈向AGI的浪潮

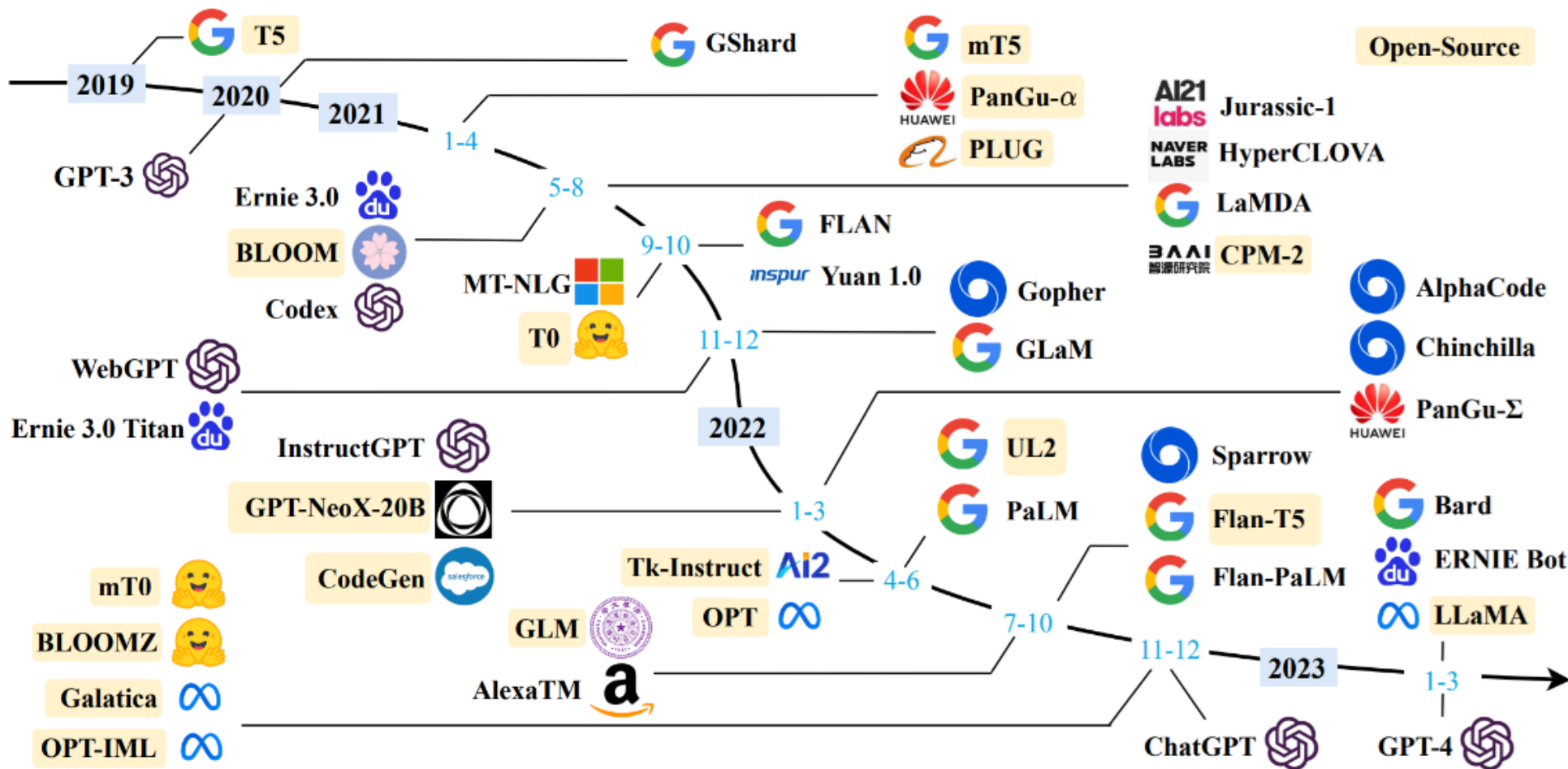
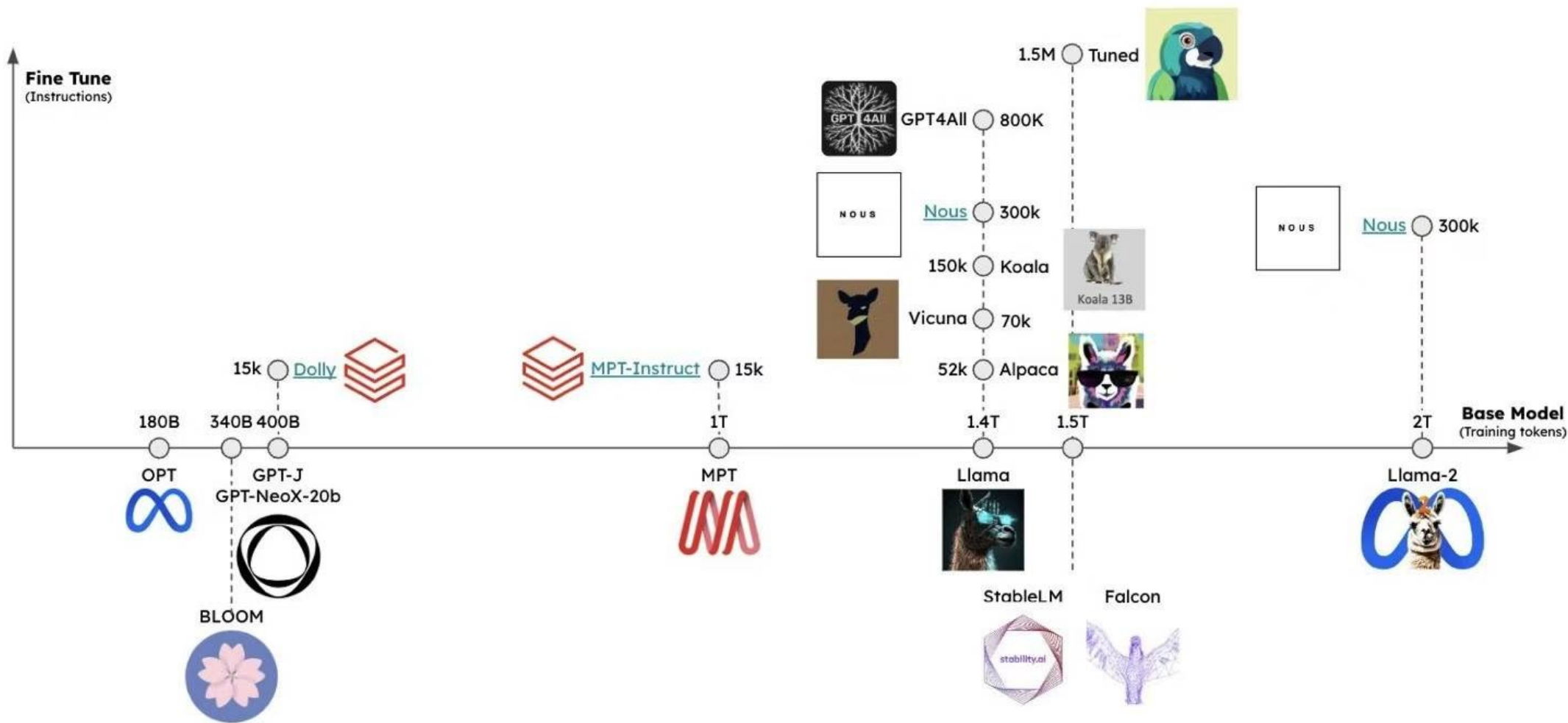


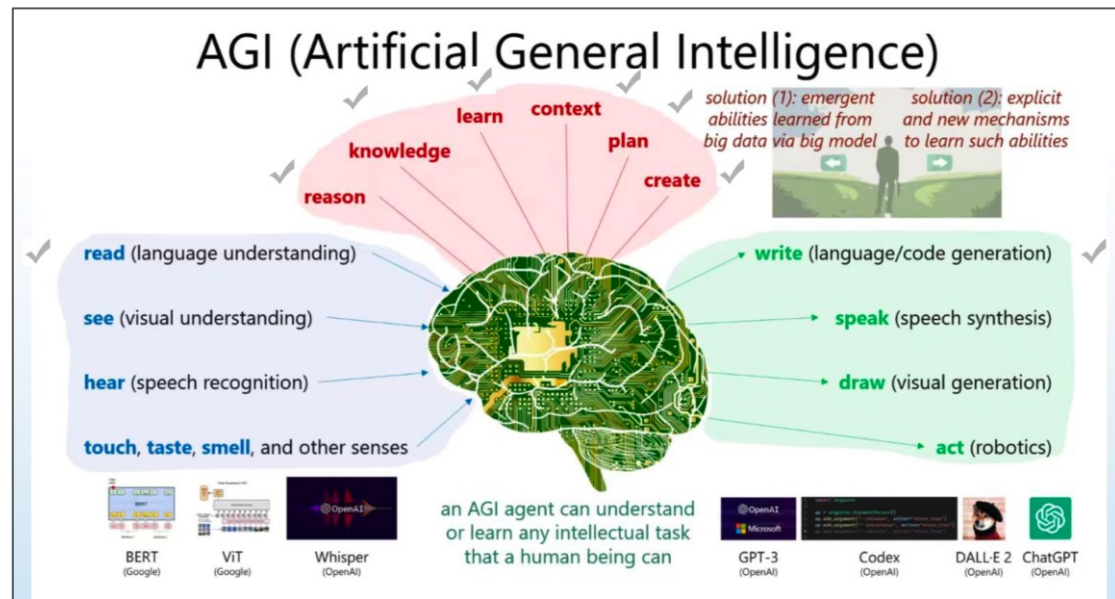
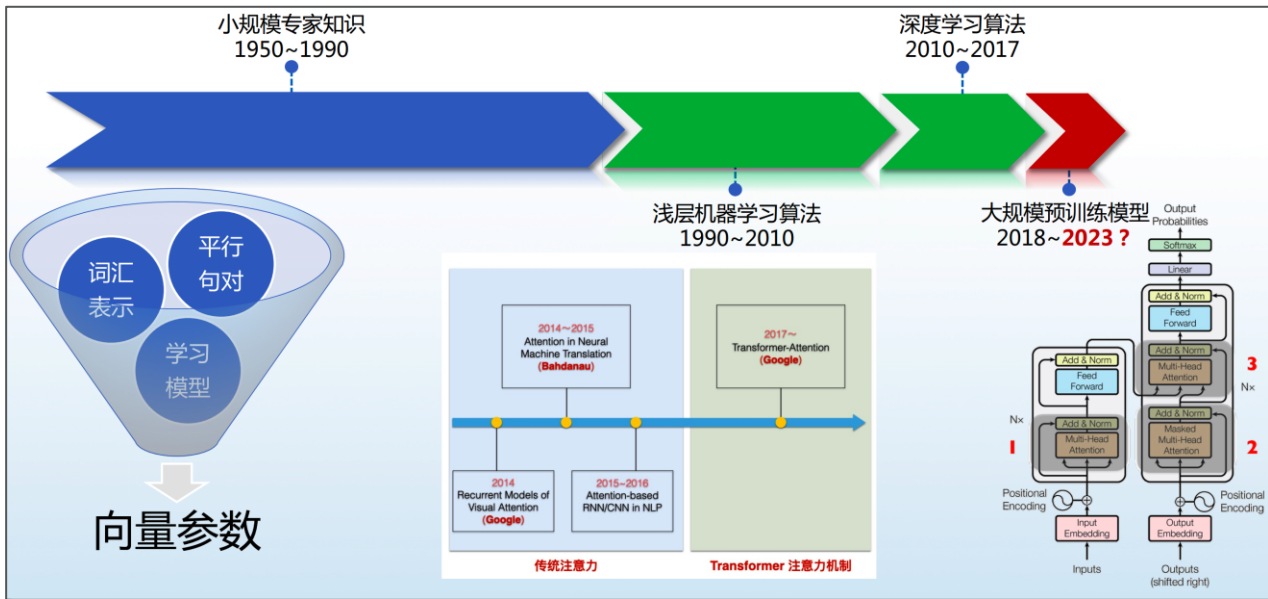
Fig. 1. A timeline of existing large language models (having a size larger than 10B) in recent years. We mark the open-source LLMs in yellow color.

A Survey of Large Language Models, 2023

▶ 开源基础模型+精调促进生态繁荣



大语言模型掀起迈向AGI的浪潮



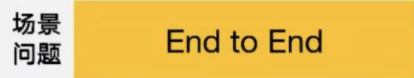
传统机器学习

把一个场景问题拆解成若干子问题，各个击破



深度学习

端到端地解决单一场景问题

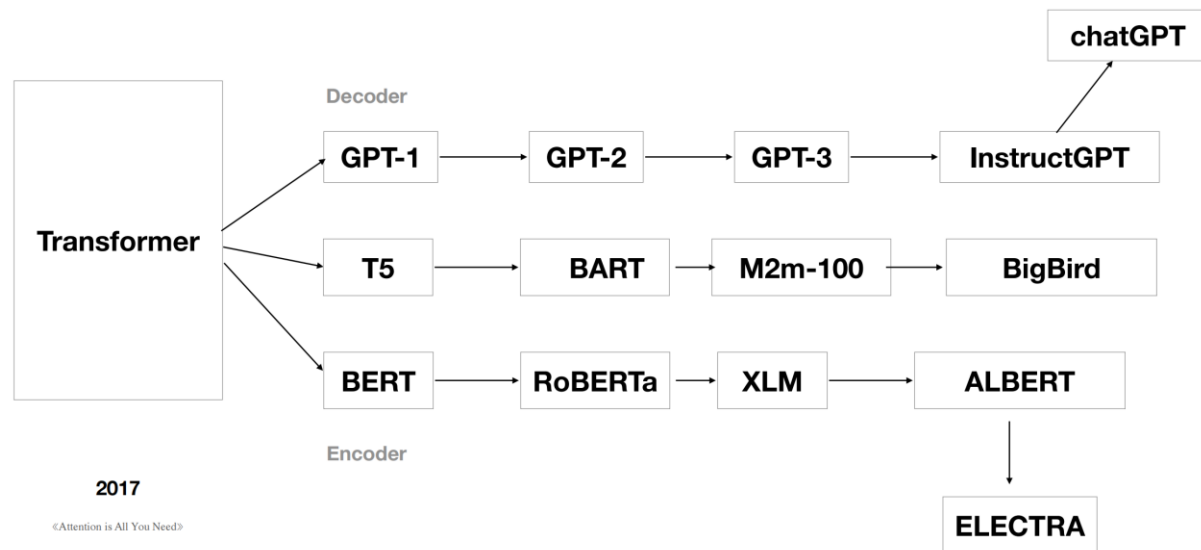
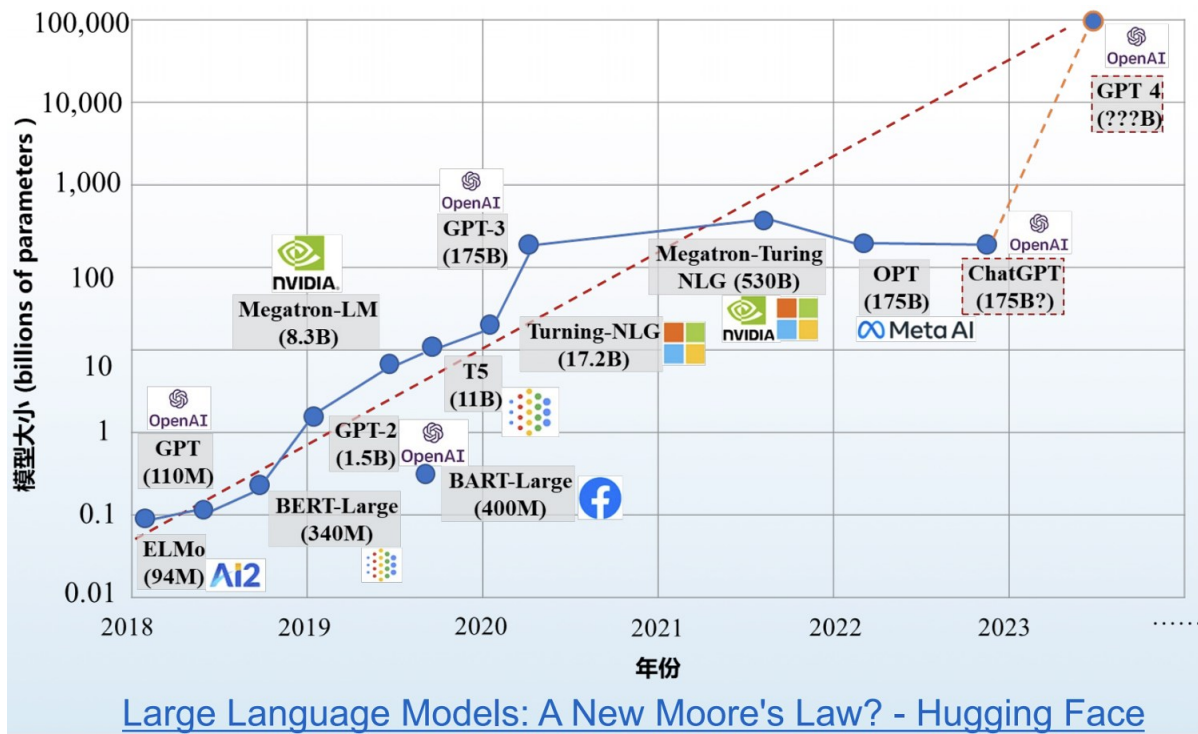


基于大模型的学习

多场景端到端问题一次解决



▶ 预训练语言模型的跃进式发展



AE语言模型：能够AR语言模型：只能获取单向信息，不能同时获取双向信息，代表为 GPT系列

获取双向信息进行预测，代表为 BERT系列

混合模型： Encoder-Decoder架构，Encoder采取自编码语言模型， Decoder采取自回归语言模型，代表为T5和BART等

▶ BERT 和 GPT 都是文字的模型，都在“猜概率”



GPT 和 BERT 的区别是什么



GPT (Generative Pre-training Transformer) 和 BERT (Bidirectional Encoder Representations from Transformers) 都是基于 Transformer 模型的预训练语言模型，可以用于各种自然语言处理任务，例如文本分类、命名实体识别、机器翻译等等。两者的主要区别在于以下几个方面：

1. 预训练方式不同：GPT 是一种单向的语言模型，即通过左侧的单词序列预测右侧的单词序列，而 BERT 则是一种双向的语言模型，即通过同时考虑左右两侧的单词序列来预测中心单词的表示。
2. 目标不同：GPT 的目标是生成下一个单词，即通过给定一段文本，预测下一个单词的概率分布；而 BERT 的目标是预测中心单词，即通过给定一段文本，预测每个单词的隐含表示，其中中心单词的表示可以用于其他任务中。
3. 预训练数据集不同：GPT 主要使用了互联网上的大规模文本数据集，例如维基百科、新闻、小说等等；而 BERT 则主要使用了 BooksCorpus 和 Wikipedia 数据集，其中 BooksCorpus 包含了 800 万本图书的文本，Wikipedia 数据集包含了维基百科的文本。
4. 模型结构不同：GPT 主要由多个 Transformer 解码器组成，而 BERT 则由多个 Transformer 编码器组成，其中 BERT 的最后一层会输出整个输入序列的表示，而 GPT 则只输出最后一个单词的表示。

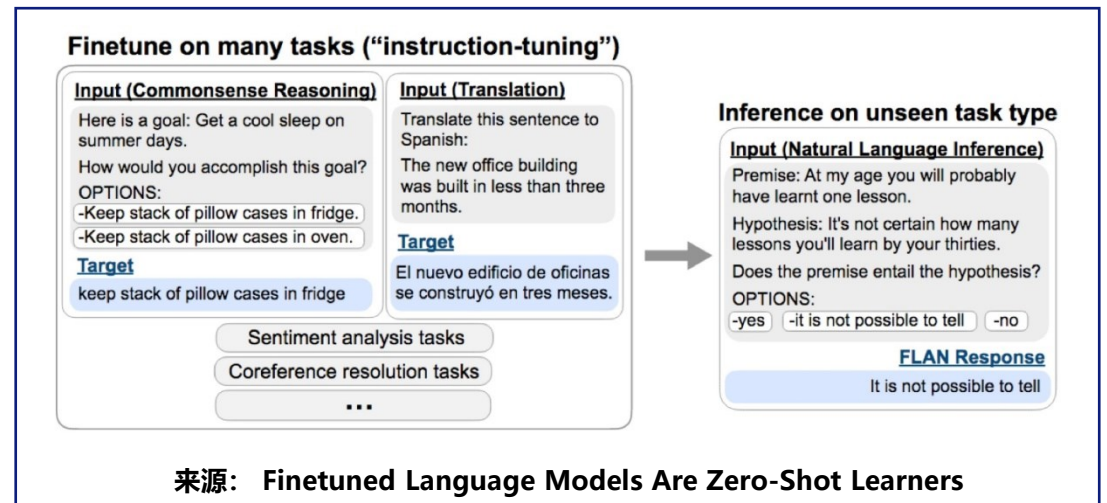
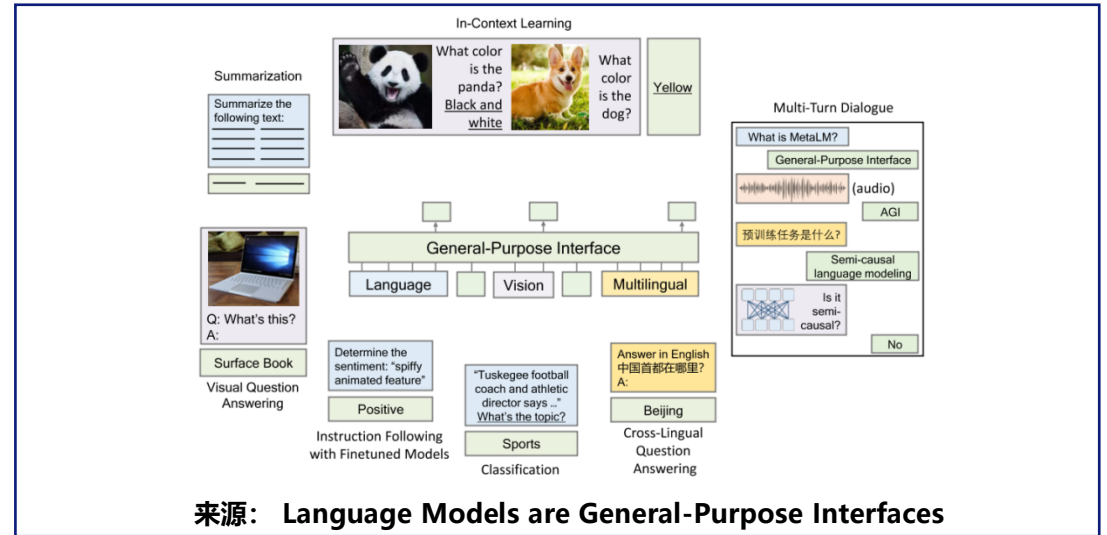
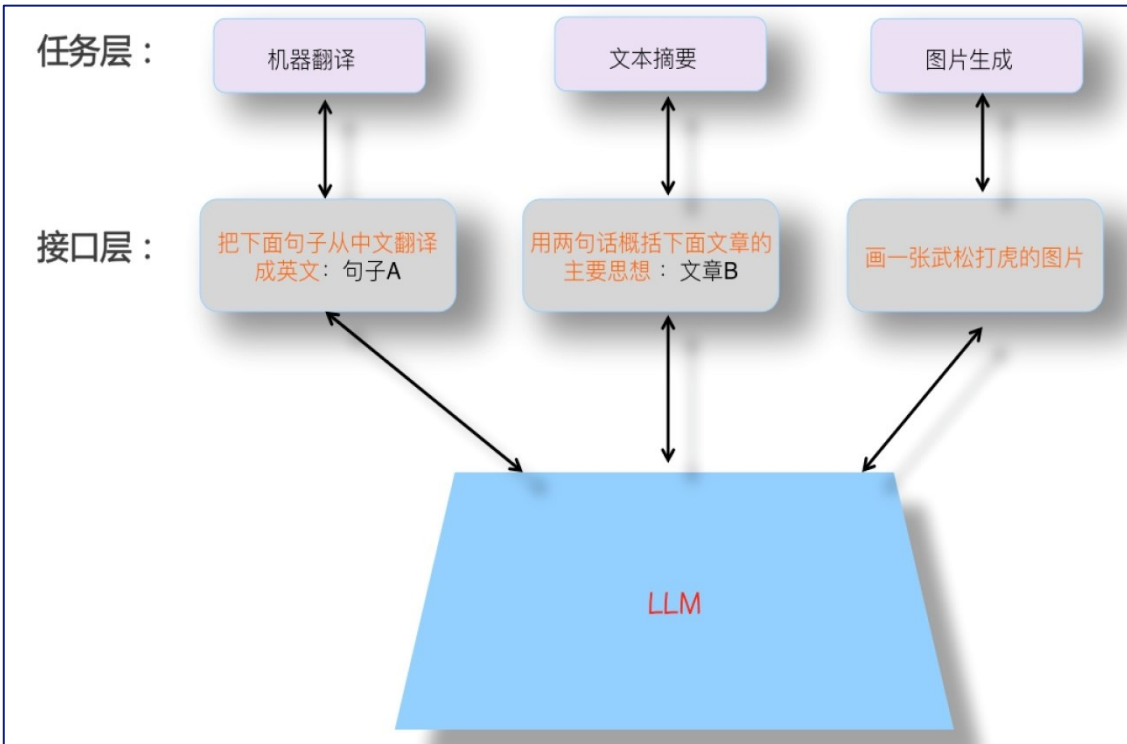
总之，GPT 和 BERT 在预训练方式、目标、预训练数据集和模型结构等方面有一定的区别。两者都是目前自然语言处理领域的前沿研究方向，可以应用于各种自然语言处理任务中，并取得了非常好的效果。

BERT	GPT
双向语言模型	单向语言模型
猜中间的字，就像完形填空	猜下一个字，就像写作文
我 __ 20 号回家	我打算20号__

▶ 大语言模型作为基础提供接口 – 用提示表达

基础模型/大模型: 指通过在大规模的数据上训练后能适应一系列下游任务的模型

- Zero Shot Prompting
- Few Shot Prompting
- In Context Learning
- Instruction



▶ 提示工程 - 什么是提示

任务的描述



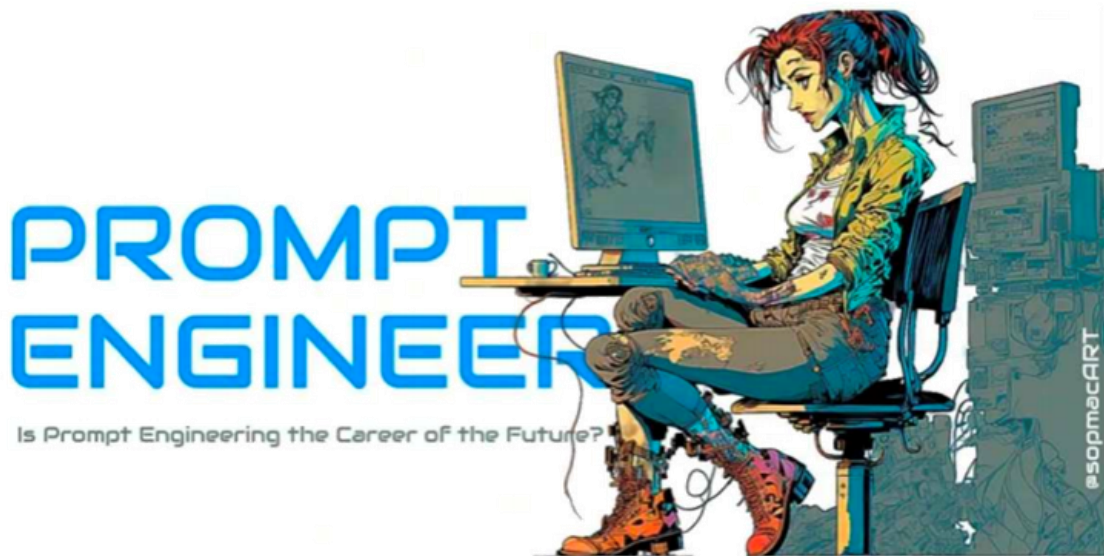
样例

▶ 提示工程师年薪达到\$300K

Prompt Engineer成为最火职业

2022年12月,Riley Goodside加入ScaleAI成为了全世界第一个Prompt Engineer(提示工程师), 据说年薪达到百万。

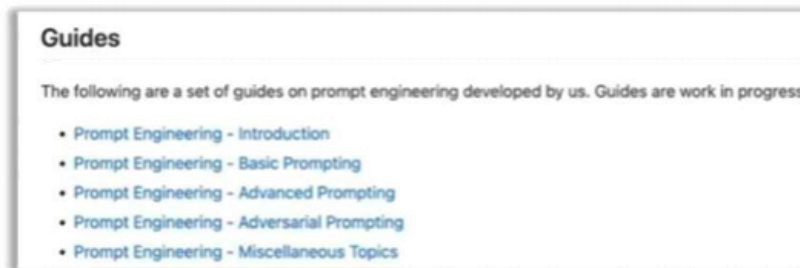
随着ChatGPT的走红, 各厂纷纷开始招募PromptEngineer来更好的"驾驶"自己的AI模型。各类大佬也发现Prompt之于模型的重要。Prompt Engineer自然而然的成为了硅谷最"紧俏"的职位之一。



OpenAI创始人: “为聊天机器人写一个非常棒的提示是一个非常高的技能, 这也是使用自然语言编程的早期例子。”



在GitHub上, 一个指导人们使用Prompt的指引在短时间内达到了5k+ 🌟。



▶ 大语言模型的“特征工程”围绕提示展开：提示工程

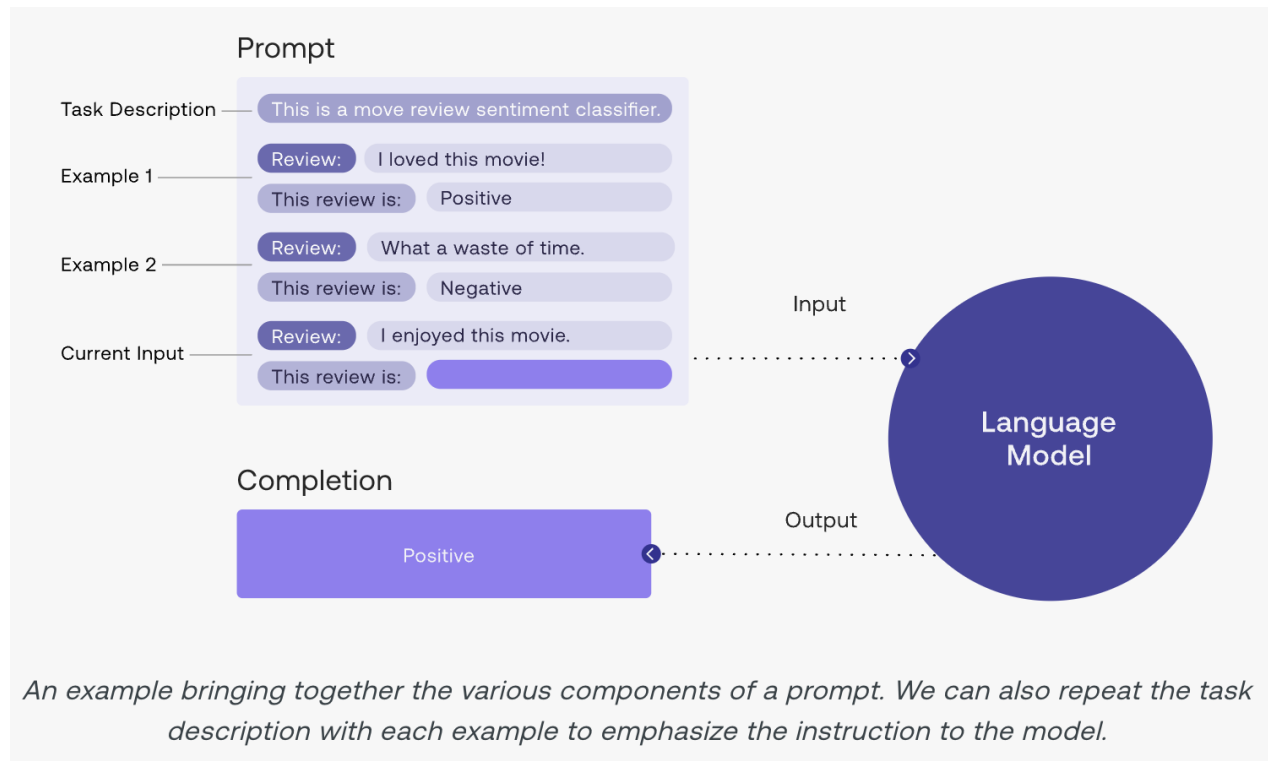
• A prompt is composed with the following components:

- Instructions
- Context
- Input data
- Output indicator

Classify the text into neutral, negative or positive

Text: I think the food was okay.

Sentiment:



▶ 大语言模型的“特征工程”围绕提示展开：提示工程示例

Text Summarization

Context

Antibiotics are a type of medication used to treat bacterial infections. They work by either killing the bacteria or preventing them from reproducing, allowing the body's immune system to fight off the infection.

Antibiotics are usually taken orally in the form of pills, capsules, or liquid solutions, or sometimes administered intravenously. They are not effective against viral infections, and using them inappropriately can lead to antibiotic resistance.

Instruction

Explain the above in one sentence:

Antibiotics are medications used to treat bacterial infections by either killing the bacteria or stopping them from reproducing, but they are not effective against viruses and overuse can lead to antibiotic resistance.

Question Answering

Answer the question based on the context below. Keep the answer short and concise. Respond "Unsure about answer" if not sure about the answer.

Context: Teplizumab traces its roots to a New Jersey drug company called Ortho Pharmaceutical. There, scientists generated an early version of the antibody, dubbed OKT3. Originally sourced from mice, the molecule was able to bind to the surface of T cells and limit their cell-killing potential. In 1986, it was approved to help prevent organ rejection after kidney transplants, making it the first therapeutic antibody allowed for human use.

Question: What was OKT3 originally sourced from?

Answer: Mice.

▶ 大语言模型的“特征工程”围绕提示展开：提示工程示例

Text Classification

Classify the text into neutral, negative or positive.

Text: I think the food was okay.

Sentiment: **Neutral**

Role Playing

The following is a conversation with an AI research assistant. The assistant tone is technical and scientific.

Human: Hello, who are you?

AI: Greeting! I am an AI research assistant. How can I help you today?

Human: Can you tell me about the creation of blackholes?

AI: Sure! Black holes are regions of spacetime where the gravitational force is so strong that nothing, not even light, can escape from it. They are created when a very massive star dies and its core collapses in on itself, forming a singularity of infinite density. The intense gravity of the singularity pulls in all the matter and radiation around it, creating the black hole.

▶ 大语言模型的“特征工程”围绕提示展开：提示工程示例

Code Generation

```
"""
Table departments, columns = [DepartmentId,
DepartmentName]
Table students, columns = [DepartmentId, StudentId,
StudentName]
Create a MySQL query for all students in the Computer
Science Department
"""

SELECT StudentId, StudentName
FROM students
WHERE DepartmentId IN (SELECT DepartmentId FROM
departments WHERE DepartmentName = 'Computer Science');
```

Reasoning

The odd numbers in this group add up to an even number: 15, 32, 5, 13, 82, 7, 1.

Solve by breaking the problem into steps. First, identify the odd numbers, add them, and indicate whether the result is odd or even.

```
Odd numbers: 15, 5, 13, 7, 1
Sum: 41
41 is an odd number.
```

▶ 大语言模型的“特征工程”围绕提示展开：提示工程示例

Zero-Shot

Chain-of-Thought

- Involves adding "Let's think step by step" to the original prompt

I went to the market and bought 10 apples. I gave 2 apples to the neighbor and 2 to the repairman. I then went and bought 5 more apples and ate 1. How many apples did I remain with?

11 apples

I went to the market and bought 10 apples. I gave 2 apples to the neighbor and 2 to the repairman. I then went and bought 5 more apples and ate 1. How many apples did I remain with?

Let's think step by step.

First, you started with 10 apples. You gave away 2 apples to the neighbor and 2 to the repairman, so you had 6 apples left. Then you bought 5 more apples, so now you had 11 apples. Finally, you ate 1 apple, so you would remain with 10 apples.

Source: [Large Language Models are Zero-Shot Reasoners](#)

Self-Consistency

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: We start with 15 trees. Later we have 21 trees. The difference must be the number of trees they planted. So, they must have planted $21 - 15 = 6$ trees. The answer is 6.

Q: If there are 3 cars in the parking lot and 2 more cars arrive, how many cars are in the parking lot?
A: There are 3 cars in the parking lot already. 2 more arrive. Now there are $3 + 2 = 5$ cars. The answer is 5.

...

Q: Olivia has \$23. She bought five bagels for \$3 each. How much money does she have left?
A: She bought 5 bagels for \$3 each. This means she spent 5

Q: When I was 6 my sister was half my age. Now I'm 70 how old is my sister?
A:

Output 1
Output 2
Output 3

When I was 6 my sister was half my age, so she was 3. Now I am 70, so she is $70 - 3 = 67$. The answer is 67.

When the narrator was 6, his sister was half his age, which is 3. Now that the narrator is 70, his sister would be $70 - 3 = 67$ years old. The answer is 67.

When the narrator was 6, his sister was half his age, which is 3. Now that the narrator is 70, his sister would be $70 - 3 = 67$ years old. The answer is 67.

▶ 经常被问到的问题：ChatGPT 为什么不能查询天气？

ChatGPT 本身不联网

ChatGPT 知识更新困难。一方面整个模型的重新训练成本很大，另一方面知识更新也会带来知识遗忘的隐忧，即你不知道他这次更新是不是在学会什么的同时，也忘记了什么。也就是说ChatGPT在解决这个问题之前，他的知识将始终落后一段时间。目前的 ChatGPT 还远远取代不了 Google。



大规模搜索

相比之下，Google 既有庞大的索引网站数据库，也提供图像、地图等功能和工具，让你快速访问其他网站更是它的「基本素养」。事实上，Google 最受欢迎的搜索就是「Facebook」、「[亚马逊](#)」等关键词。



生成式搜索

在 ChatGPT 的基础设定里，信息的时效性已经落后了。ChatGPT 基于互联网的数十亿个文本示例训练，学习生涯停留在了 2021 年，它也没有联网，不实时调用外部网络资源，无法像 Google 地图般根据位置为你推荐餐馆，不能帮你打开 Twitter、Facebook。

▶ GPT = Giant Parrot Talking? (统计学习)



▶ 大语言模型存在的不足

OL 河南省洛阳市的旅游公司有哪些

作为AI助手，我无法查询到具体的旅游公司信息，建议您可以通过以下方式查询：

1. 在百度、携程、马蜂窝等旅游网站上搜索“洛阳旅游公司”，可以看到相关公司的信息和评价。
2. 在当地的商业指南或黄页中查找相关公司的联系方式和地址。
3. 咨询当地的旅游局或旅游协会，了解当地的旅游公司情况。

OL 今天天气怎么样 **实时问题**

抱歉，我作为一个人工智能语言模型，无法实时获取当天的天气信息。请您查询当地的天气预报或者查看当地天气应用以获取最准确的信息。

OL 你如何看待2022年的世界杯比赛结果

作为一个人工智能语言模型，我没有情感和预测能力，也不能预测未来的事件。因此，我不能对2022年世界杯的比赛结果进行任何预测或偏见。不过，世界杯是一个全球性的运动盛会，它吸引了数百万人的关注和热情，我们可以期待在比赛中看到精彩的比赛和顶尖球员的出色表现。

OL 建筑工地有一批砖，最上层两块砖，第2层6块砖，第3层10块砖.....，依次每层比其上一层多4块，已知最下层有2106块砖，这堆砖共有多少块？

我们可以先计算出这堆砖共有多少层。设共有n层，则最下面一层有：

$$2106 = (2 + (2 + 4) + (2 + 2 \times 4) + \dots + (2 + (n - 2) \times 4) + (2 + (n - 1) \times 4))$$

可以对上式进行化简，得到：

$$2106 = [2 \times n + 4 \times (1 + 2 + \dots + (n - 1))] / 2$$

化简后得到：

$$2106 = (2n + 4n^2 - 4n) / 2$$

即：

$$2n^2 + n - 1053 = 0$$

解得 $n \approx 32.6$ ，即共有32层。

然后，我们可以利用等差数列求和公式来计算这堆砖的总数。设第1层有a1块砖，公差为d，则第n层有：

$$a_n = a_1 + (n - 1)d$$

利用等差数列求和公式，可以得到：

$$\text{砖的总数} = (a_1 + a_n) \times n / 2 = (2a_1 + (n - 1)d) \times n / 2$$

代入 $a_1 = 2$ ， $d = 4$ ， $n = 32$ ，可得砖的总数为：

$$(2 \times 2 + (32 - 1) \times 4) \times 32 / 2 = 2112$$

因此，这堆砖共有2112块。

正确答案：
555458块

- 存在“幻觉”，容易一本正经胡说八道
- 提示直接决定回复质量，经常容易“说正确的废话”
- 无法理解与更新个性化知识，如无法针对个性需求提供信息
- 无法处理动态、实时问题，如股票价格、天气变化等
- 知识无法更新，不能提供最新信息
- 不擅长数学计算和逻辑推理等

▶ 大语言模型落地的限制

知识层面

识

ChatGPT 缺少知识驱动，对领域问题依然有理解瓶颈。

1. 本质上，ChatGPT是一个高级的纯数据驱动的统计模型——连四则运算都搞不定。
2. OpenAI对于知识的态度一向是：模型参数即知识。

数据层面

数

ChatGPT 本身在数据适配性、合规性和安全性上，针对中国应用依然具有限制。

1. 墙外语料与监管要求冲突
2. 大量现有领域知识用LLM再训练可能浪费

可信层面

黑

ChatGPT 仍然是一个黑盒模型，不仅计算过程不可解释，产出也不可信。

1. ChatGPT具有初步的逻辑分析和规划能力，但是模型仍然缺乏可解释和溯源能力，无法支持复杂决策任务的可信计算。
2. ChatGPT的结果，无法证明，也无法证伪，你不一定敢信

成本层面

贵

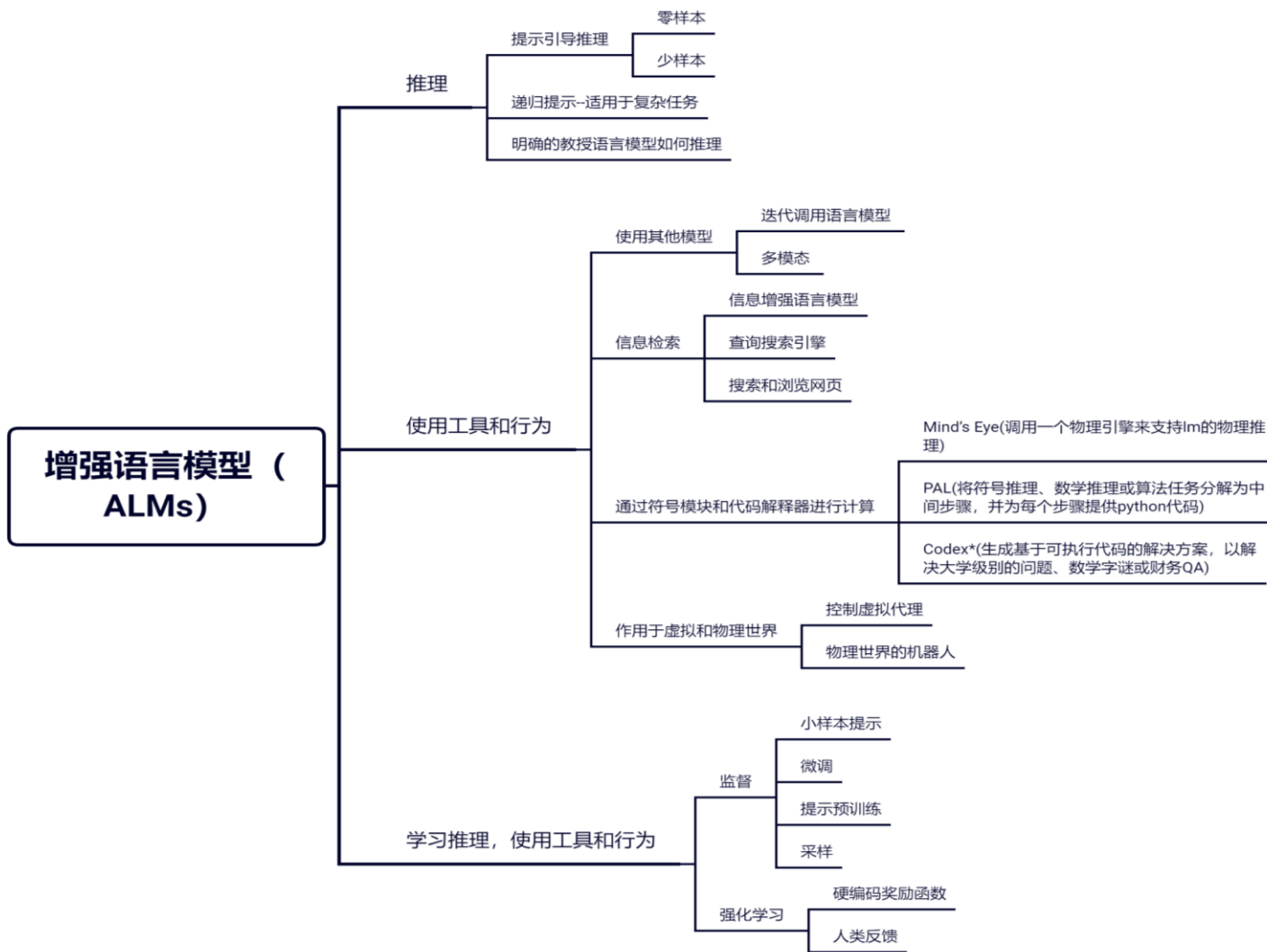
大模型的稳定训练并实现优异性能需要极高的计算成本和工程实现能力。

1. 初代GPT-3的学习耗费1200 万美元
2. 除OpenAI和DeepMind，几乎没有完整掌握技术栈的机构

增强语言大模型ALM的必要性和可能方向

❖ 针对当前LLM的不足，研究者们提出了一些改进措施，例如使LM利用外部工具，用LM的权重中不包含的重要缺失信息来增强上下文理解，形成更强大的智能体；这些模型统称为**增强语言模型（ALMs）**。

- ❖ **推理（Reasoning）**：将复杂任务分解成更简单的子任务，LM可以自己 或使用工具更容易地解决。
- ❖ **工具（Tool）**：收集外部信息，或者对ALM感知的虚拟或物理世界产生影响。
- ❖ **行为（Act）**：调用一个对虚拟或物理世界有影响的工具并观察其结果，将其纳入ALM的当前上下文。
- ❖ **结合使用**：推理和工具可以放在同一个模块里，二者都是通过增强LM的上下文来更好地预测缺失；收集额外信息的工具和对虚拟或物理世界产生影响的工具可以被LM以同样的方式调用。



来源: Augmented Language Models: a Survey (Yann Lecun et al.)

PART 02

ALM关键技术



增强关键技术1：高级提示工程

Chain-of-Thought (Wei et al., 2022)

Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 tennis balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: The bakers at the Beverly Hills Bakery baked 200 loaves of bread on Monday morning. They sold 93 loaves in the morning and 39 loaves in the afternoon. A grocery store returned 6 unsold loaves. How many loaves of bread did they have left?

Model Output

A: The bakers started with 200 loaves. They sold 93 in the morning and 39 in the afternoon. So they sold $93 + 39 = 132$ loaves. The grocery store returned 6 loaves. So they had $200 - 132 - 6 = 62$ loaves left. The answer is 62.



Program-aided Language models (this work)

Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 tennis balls.

`tennis_balls = 5`

2 cans of 3 tennis balls each is

`bought_balls = 2 * 3`

tennis balls. The answer is

`answer = tennis_balls + bought_balls`

Q: The bakers at the Beverly Hills Bakery baked 200 loaves of bread on Monday morning. They sold 93 loaves in the morning and 39 loaves in the afternoon. A grocery store returned 6 unsold loaves. How many loaves of bread did they have left?

Model Output

A: The bakers started with 200 loaves

`loaves_baked = 200`

They sold 93 in the morning and 39 in the afternoon

`loaves_sold_morning = 93`

`loaves_sold_afternoon = 39`

The grocery store returned 6 loaves.

`loaves_returned = 6`

The answer is

`answer = loaves_baked - loaves_sold_morning - loaves_sold_afternoon + loaves_returned`

`>>> print(answer)`

74



Source: PAL: Program-aided Language Models

增强关键技术1：高级提示工程

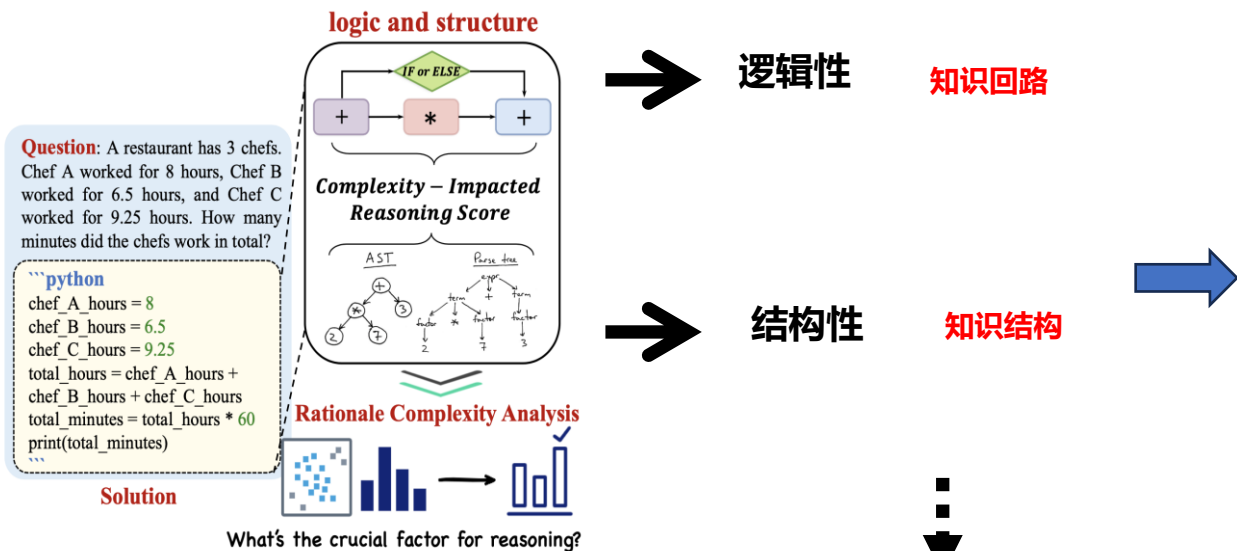
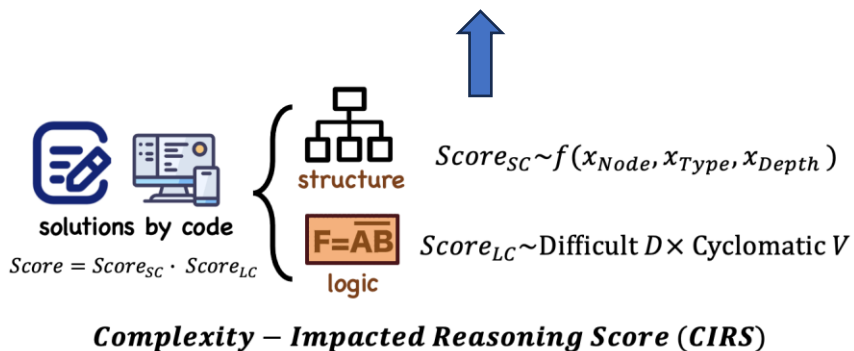


Figure 1: We leverage code data to analyze what kind of data is crucial for reasoning abilities of LLMs models.

如何去衡量对推理的作用?



When Do Program-of-Thoughts Work for Reasoning? (to appear)

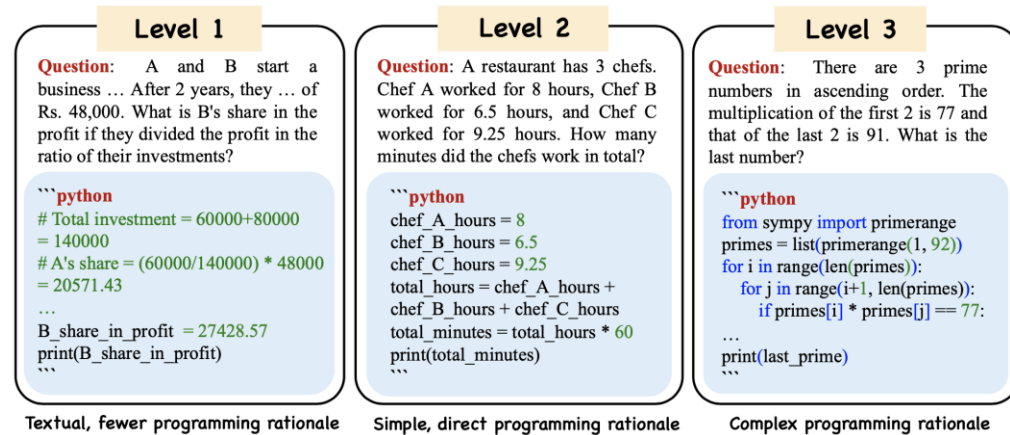


Figure 4: Examples of different of complexities and levels.

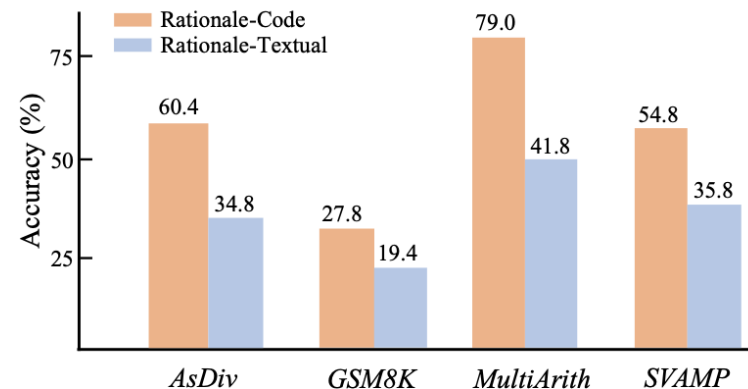
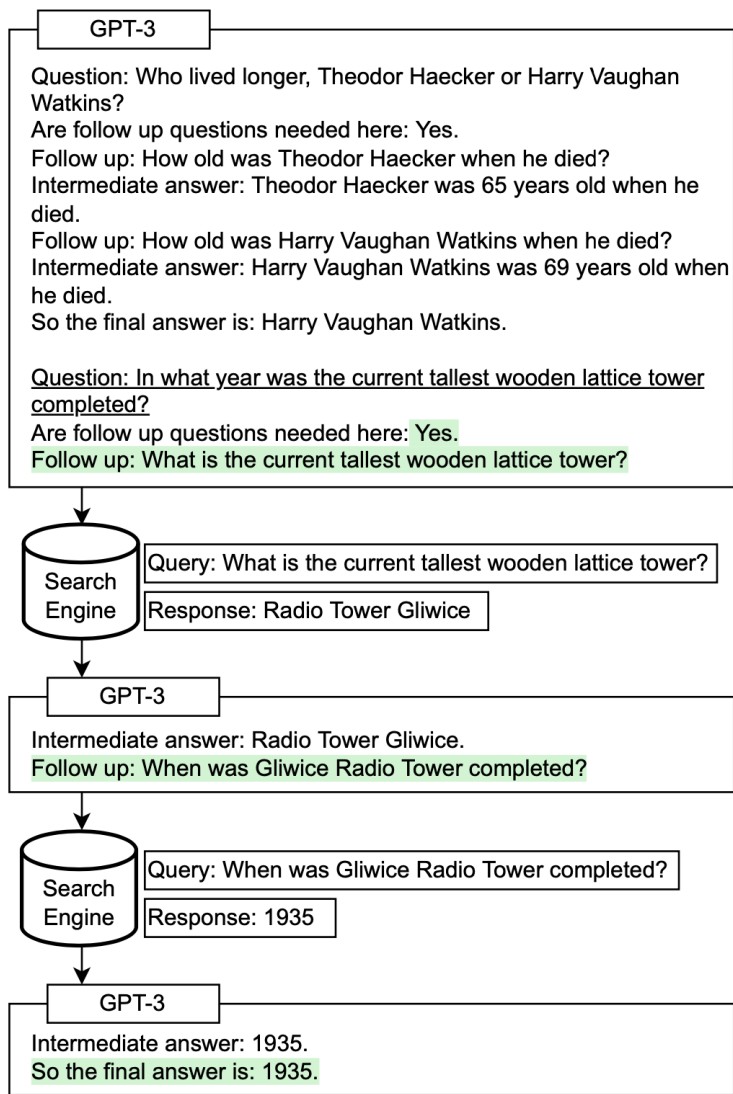


Figure 6: Comparison for textual and code rationales.

增强关键技术1：高级提示工程



Self-Ask ([Press et al. 2022](#)) is a method to repeatedly prompt the model to ask *following-up questions* to construct the thought process iteratively.

Measuring and Narrowing the Compositionality Gap in Language Models, 2022

增强关键技术1：高级提示工程

(1) Hotspot QA

Question: Aside from the Apple Remote, what other device can control the program Apple Remote was originally designed to interact with?

(1a) Standard

Answer: iPod



(1b) CoT (Reason Only)

Thought: Let's think step by step. Apple Remote was originally designed to interact with Apple TV. Apple TV can be controlled by iPhone, iPad, and iPod Touch. So the answer is iPhone, iPad, and iPod Touch.

Answer: iPhone, iPad, iPod Touch



(1c) Act-Only

Act 1: Search[Apple Remote]

Obs 1: The Apple Remote is a remote control ...

Act 2: Search[Front Row]

Obs 2: Could not find [Front Row]. Similar: ...

Act 3: Search[Front Row (software)]

Obs 3: Front Row is a discontinued media center software ...

Act 4: Finish[yes]



(1d) ReAct (Reason + Act)

Thought 1: I need to search Apple Remote and find the program it was originally designed to interact with.

Act 1: Search[Apple Remote]

Obs 1: The Apple Remote is a remote control introduced in October 2005 by Apple ... originally designed to control the Front Row media center program ...

Thought 2: Apple Remote was originally designed to control the Front Row media center program. I need to search Front Row next and find what other device can control it.

Act 2: Search[Front Row]

Obs 2: Could not find [Front Row]. Similar: ['Front Row Seat to Earth', 'Front Row Motorsports', 'Front Row (software)', ...]

Thought 3: Front Row is not found. I need to search Front Row (software) .

Act 3: Search[Front Row (software)]

Obs 3: Front Row is a discontinued media center software ...

Thought 4: Front Row (software) is controlled by an Apple Remote or the keyboard function keys. So the answer is keyboard function keys.

Act 4: Finish[keyboard function keys]



Source: ReAct: Synergizing Reasoning and Acting in Language Models

增强关键技术1：高级提示工程

单阶段方法

Single-Stage

Q: There are 3 cars in the parking lot and 2 more cars arrive. How many cars are in the parking lot?

C: There are 3 cars in the parking lot already. 2 more arrive. Now there are $3 + 2 = 5$ cars.

A: The answer is 5.

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

Language Model

C: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$.

A: The answer is 11.

- Sensitivity of in-context learning
- complexity、diversity、explicitly

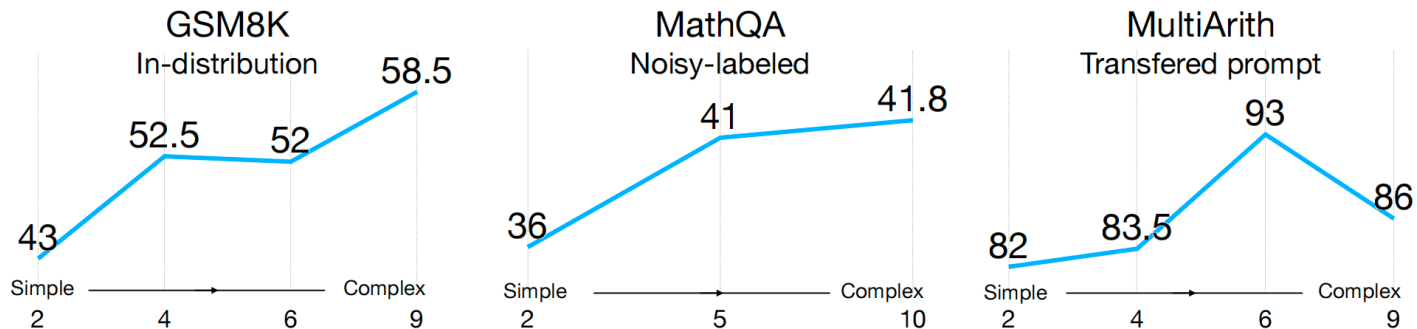


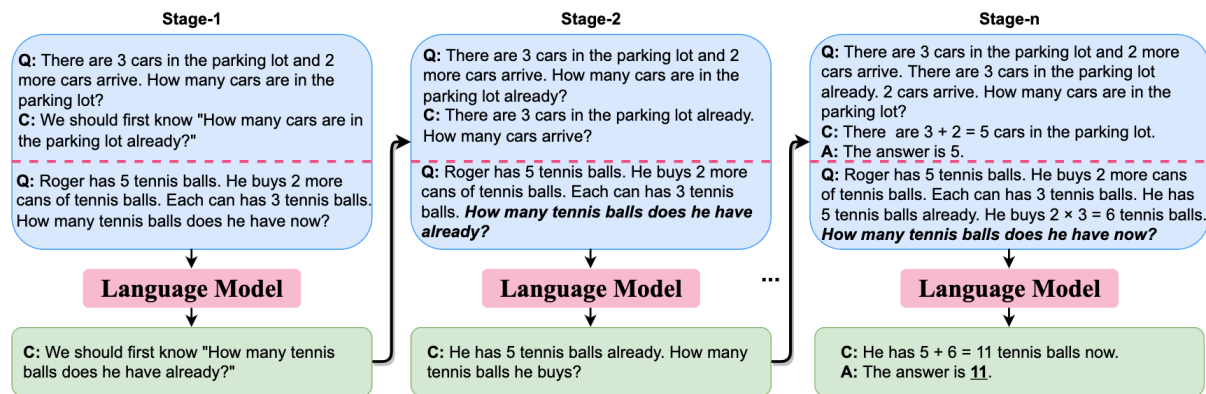
Figure 2: Validation set performance. X-axis means reasoning steps and y-axis means accuracy. More reasoning steps in prompts overall achieve higher accuracy when prompts are in-distribution (left), noisily labeled (middle), and out of distribution (right).

在CoT方法上简单直接的改进，现有工作表明in-context learning对于其中包含的例子的顺序、质量等非常敏感，一个非常小的改变可能会引起模型效果上很大的下降，想要进一步的优化CoT，一个直观的手段是优化其中的例子，现有方法主要是提高例子的复杂度、多样性、明确性等

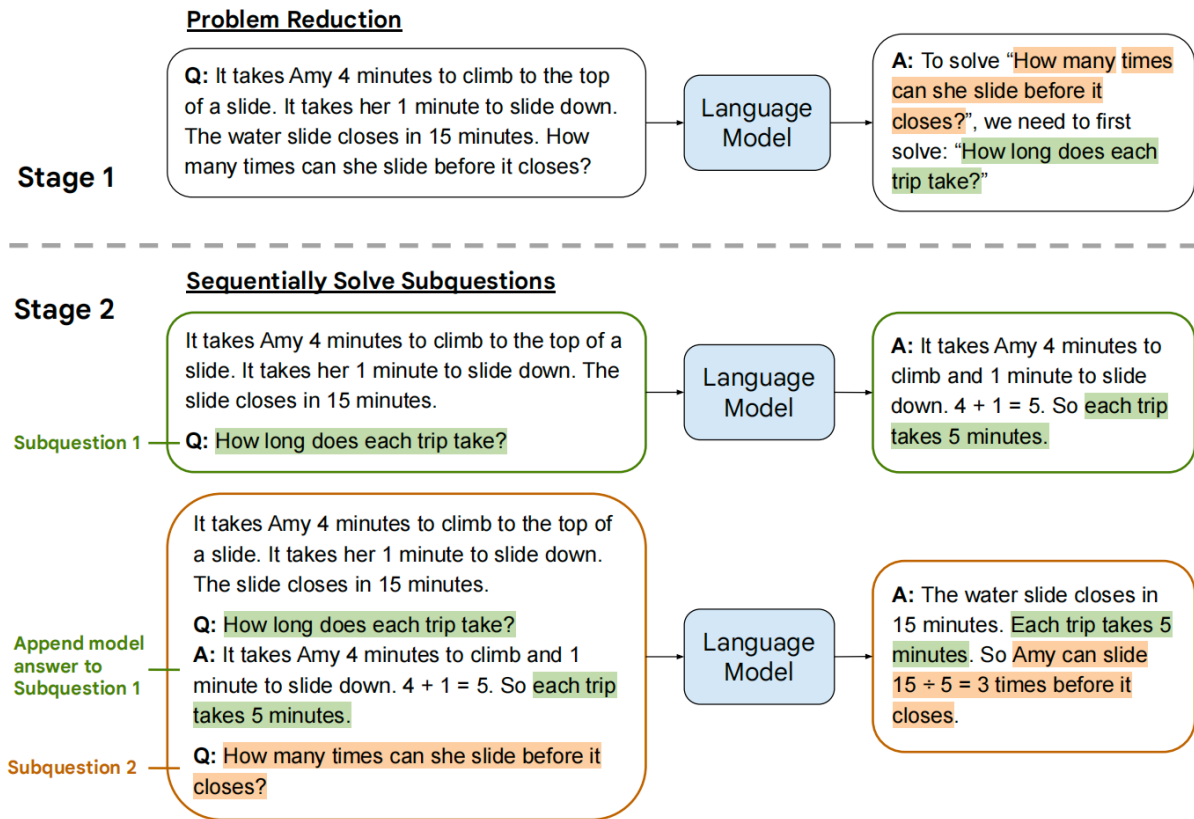
Complexity-Based Prompting for Multi-step Reasoning, ICLR 2023

增强关键技术1：高级提示工程

多阶段方法



人类在推理问题时往往不是一次性的，而是反复思考的多次输入输出，受此启发便产生了相对于单阶段的多阶段类提示方法，该类方法将一个复杂问题拆成多个子问题按多阶段进行推理



Least-to-Most Prompting Enables Complex Reasoning in Large Language Models 2022

增强关键技术1：高级提示工程

提示分解法

QC: Concatenate the first letter of every word in "Jack Ryan" using spaces
Q1: [split] What are the words in "Jack Ryan"?
#1: ["Jack", "Ryan"]
Q2: (foreach) [str_pos] What is the first letter of #1?
#2: ["J", "R"]
Q3: [merge] Concatenate #2 with spaces
#3: "J R"
Q4: [EOQ]
...

decomp

Q: What are the words in "Elon Musk Tesla"?
A: ["Elon", "Musk", "Tesla"]

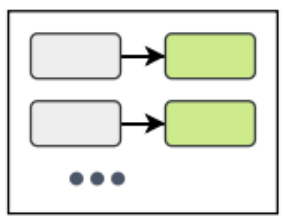
Q: What are the letters in "C++"?
A: ["C", "+", "+"]
...

split

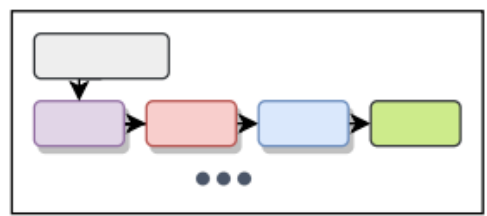
Q: Concatenate ["n", "i", "e"]
A: "nie"

Q: Concatenate ["n", "i", "c", "e"] using spaces
A: "n i c e"
...

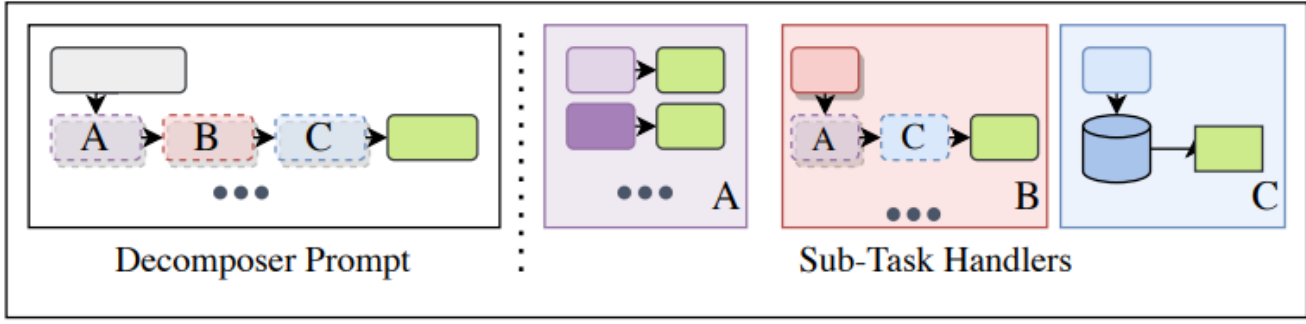
merge



Standard Prompting



Chain-of-Thought Prompting



Decomposed Prompting

Decomposed Prompting: A Modular Approach for Solving Complex Tasks 2022

增强关键技术1：高级提示工程

思维链结构化

Tree of Thoughts

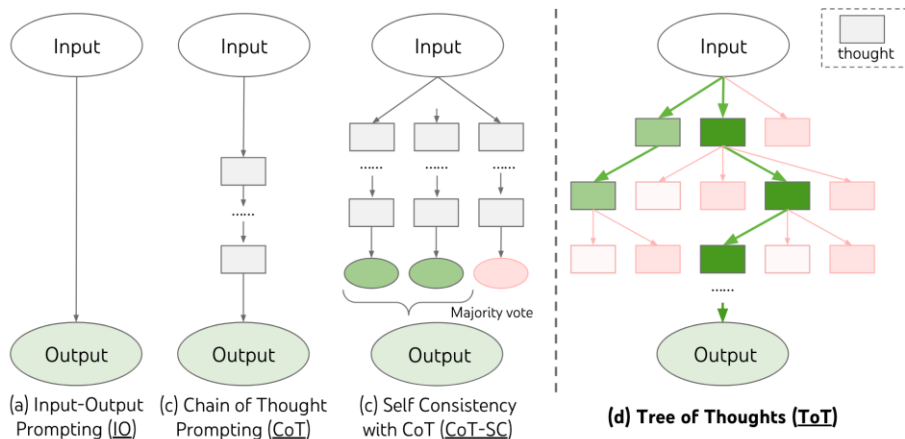


Figure 1: Schematic illustrating various approaches to problem solving with LLMs. Each rectangle box represents a *thought*, which is a coherent language sequence that serves as an intermediate step toward problem solving. See concrete examples of how thoughts are generated, evaluated, and searched in Figures 046.

Graph of Thoughts

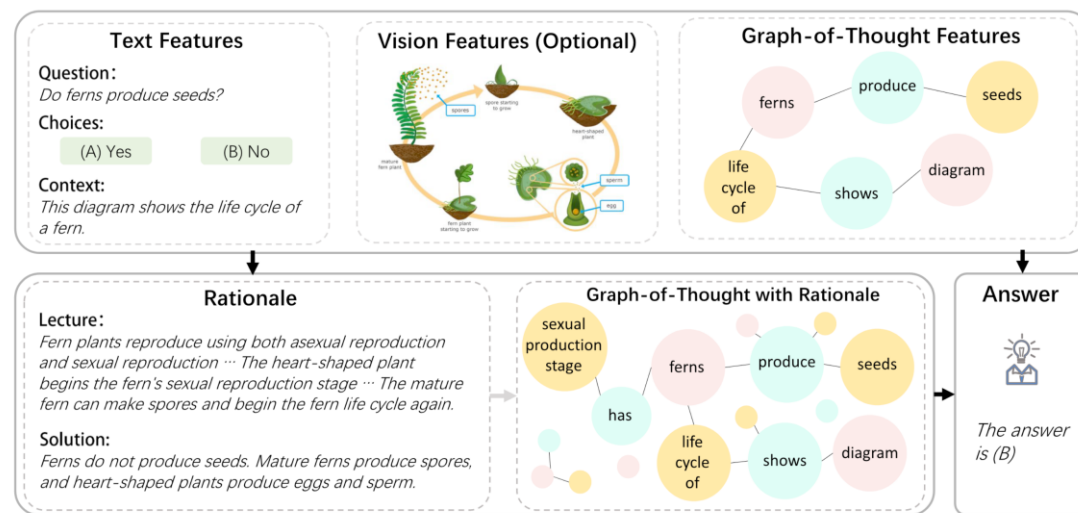
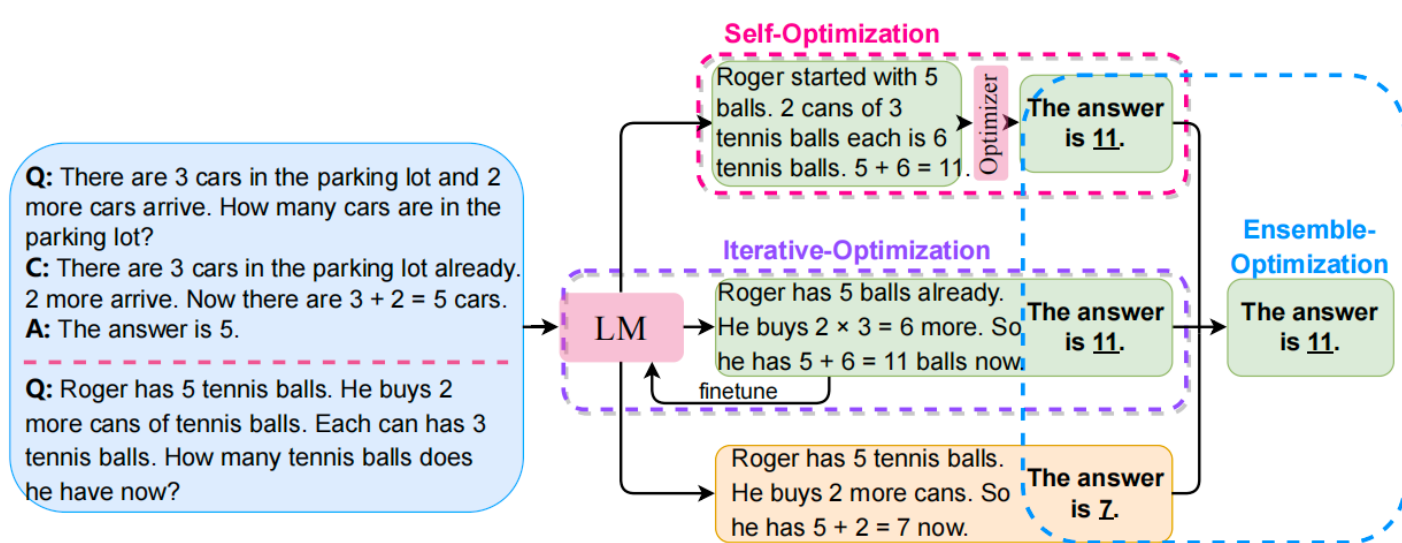


Figure 1: An example of GoT reasoning. Vision features are optional and are only required in multimodal reasoning task.

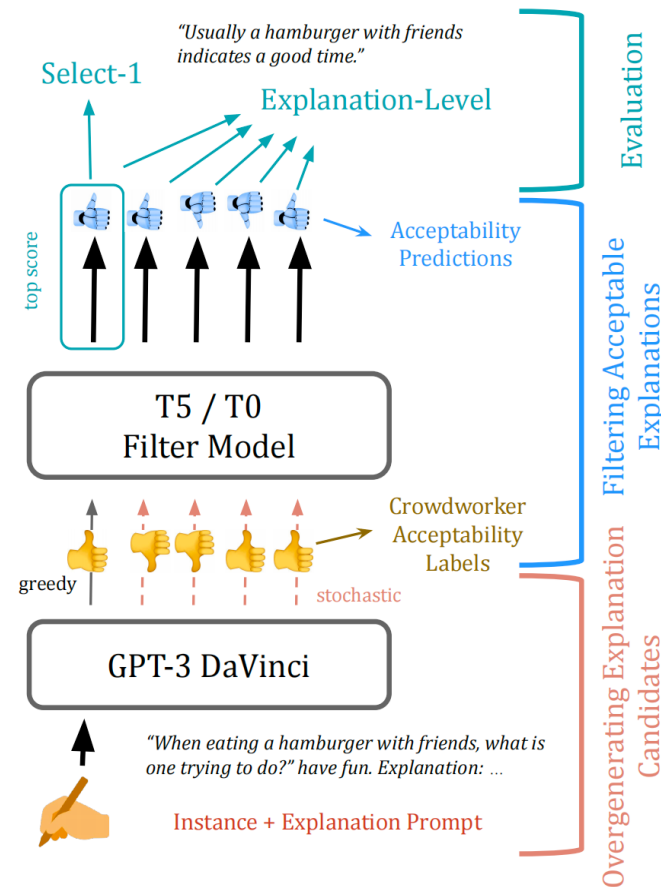
结构化更好或逻辑性更强的提示性知识起到对复杂推理产生约束作用，产生更好Step-by-step的思维过程

增强关键技术1：高级提示工程

Self-optimization



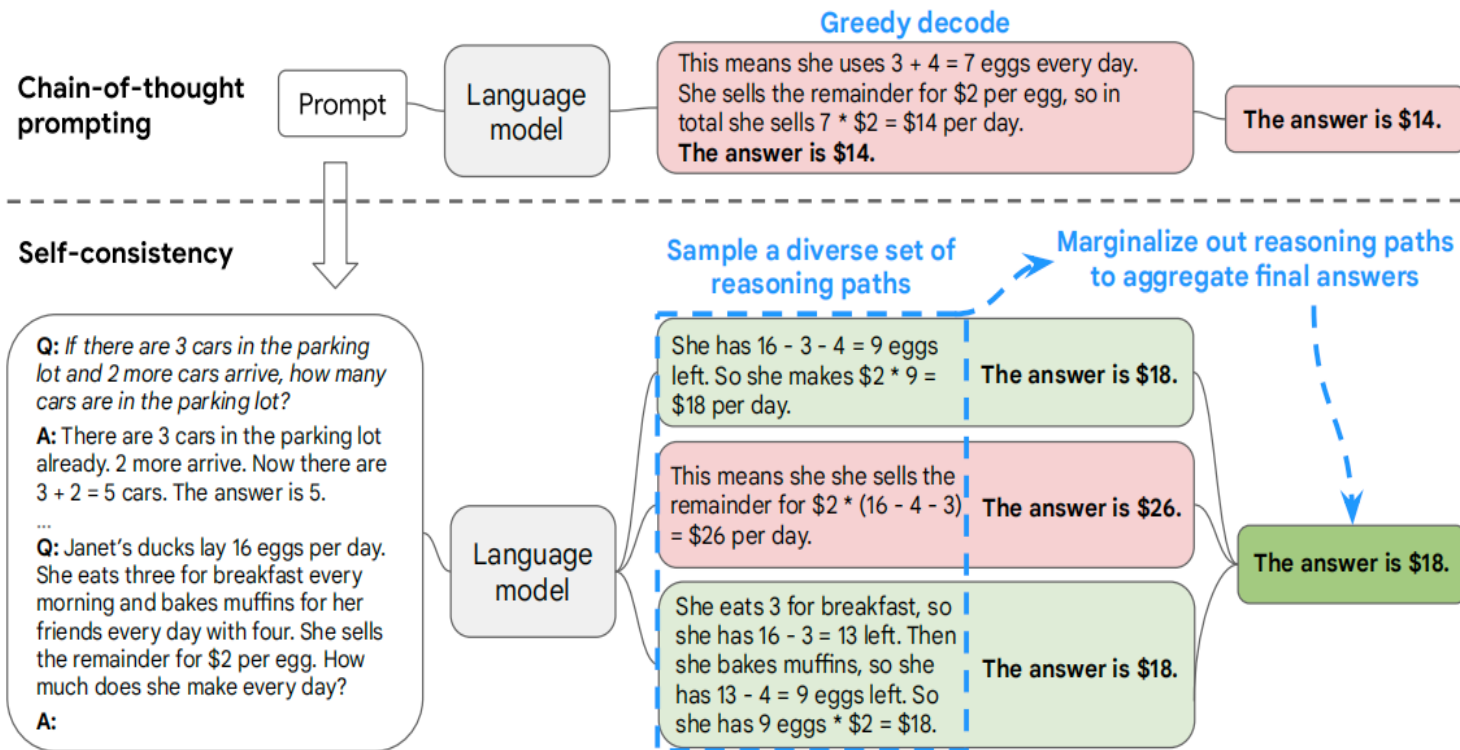
单阶段和多阶段代表的prompt工程方法是在模型的输入上进行策略增强的，而过程优化方法还可以在模型输出上进行策略增强。推理过程是CoT prompting的关键，self-optimization方法是引入相应的矫正模块或过滤模块等对推理路径进行优化



Reframing Human-AI Collaboration for Generating Free-Text Explanations, NAACL 2022

增强关键技术1：高级提示工程

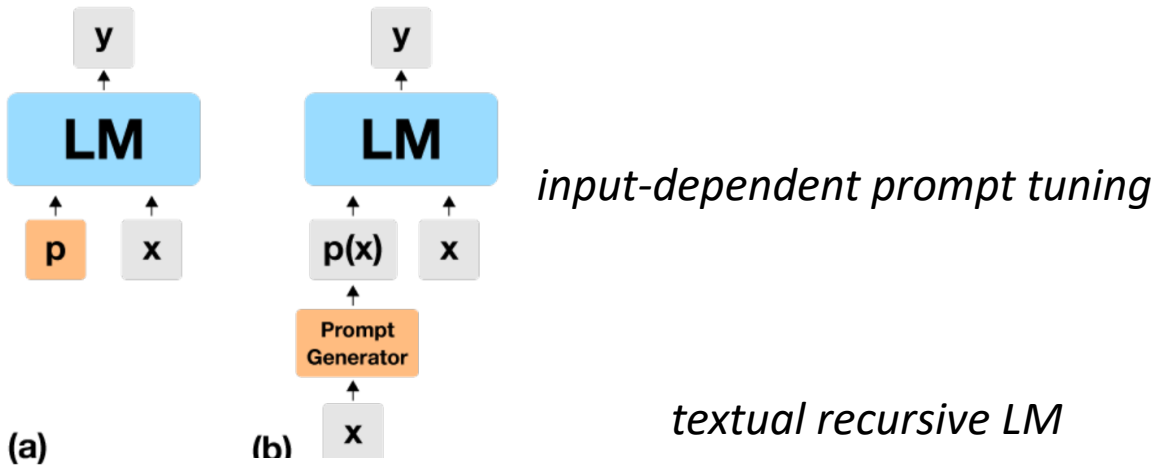
Ensemble-Optimization



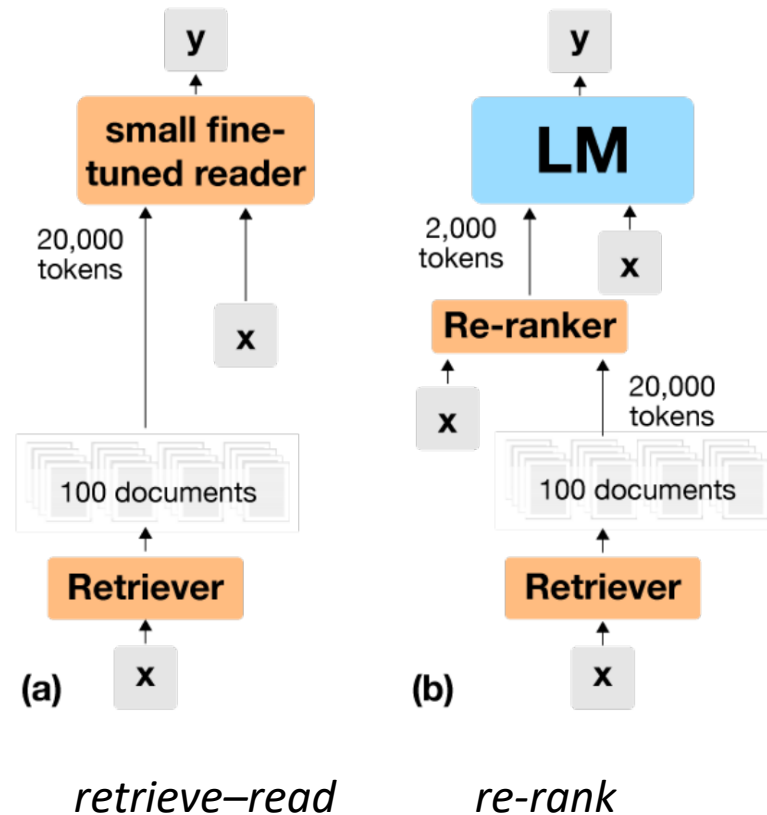
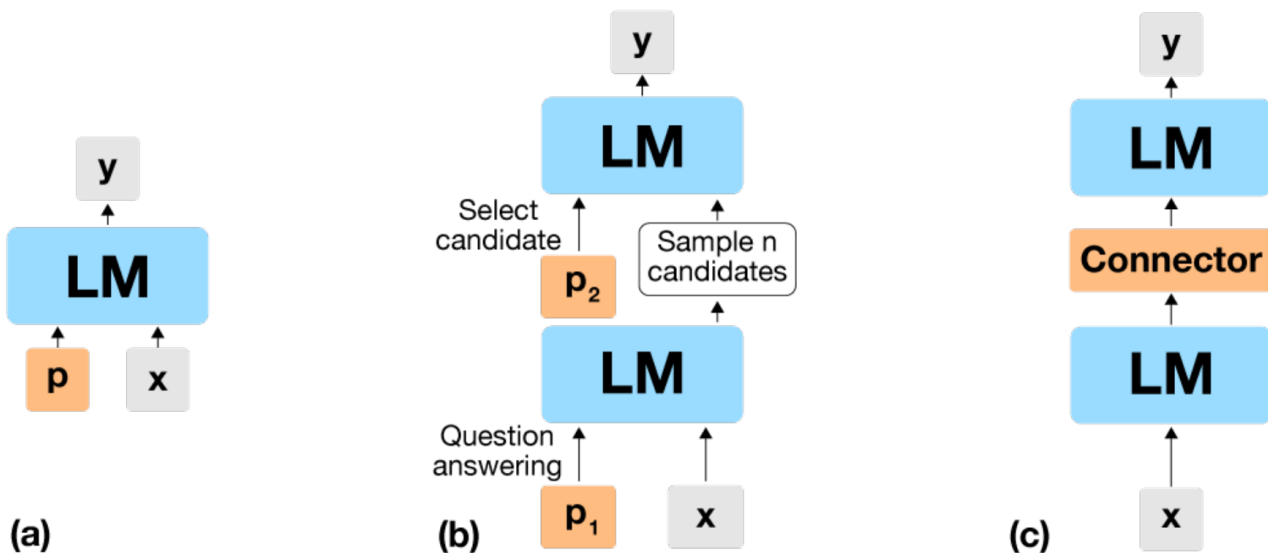
条条大路通罗马，推理问题往往不只有一个路径可以到达正确答案，往往一个问题有多种解法，由此便产生了集成优化的方法，该类方法以Google的self-consistency为代表，使用在生成式解码过程中常用的采样手段，比如temperature、top-k等，让语言模型产生多条路径，将所有路径的答案汇总生成最终答案

Self-Consistency Improves Chain of Thought Reasoning in Language Models 2022

增强关键技术1：提示工程的“设计模式”



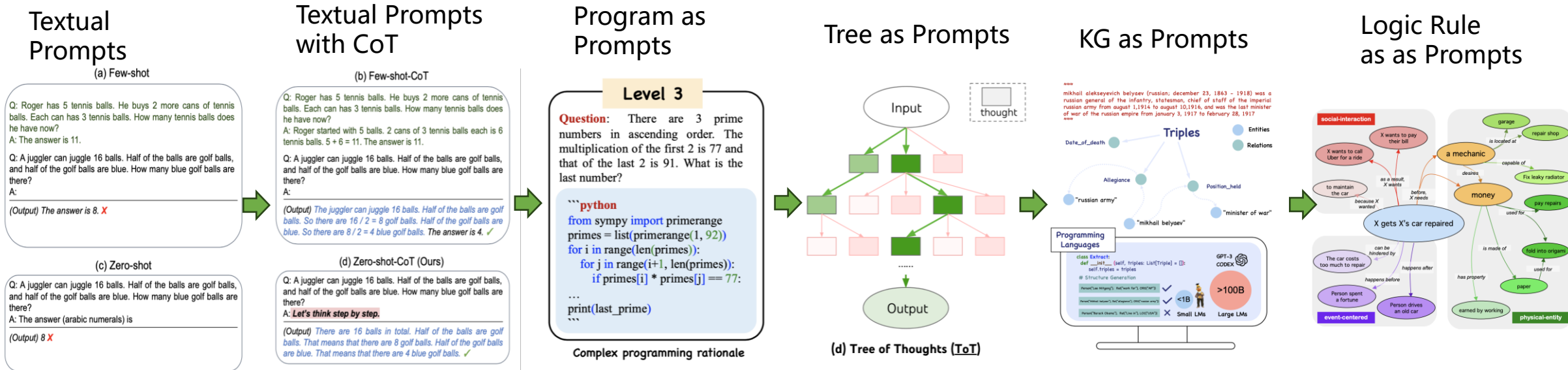
textual recursive LM



STANDING ON THE SHOULDERS OF GIANT FROZEN LANGUAGE MODELS 2022

增强关键技术1：提示工程即知识工程

大量实践表明提示设计的好坏和详细程度，以及提示表示的结构化和逻辑性水平对模型输出结果影响巨大，**提示工程本质就是知识工程**，目的也是从人获取先验知识，来指导模型训练或激活模型推理能力。



提示知识的结构化水平和逻辑性逐步增加，模型推能力更强，但提示知识获取的难度逐步增加，规模涌现越不容易实现

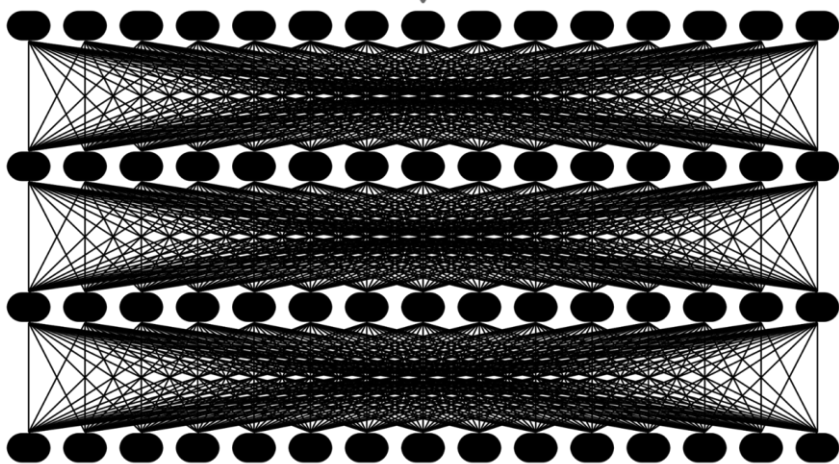
表示 vs 推理

表示水平与推理能力的正比关系仍然存在
表示复杂度与规模化的矛盾关系仍然存在

表示 vs 规模

▶ 增强关键技术2：检索/知识增强

🔍 When was Stanford University founded?



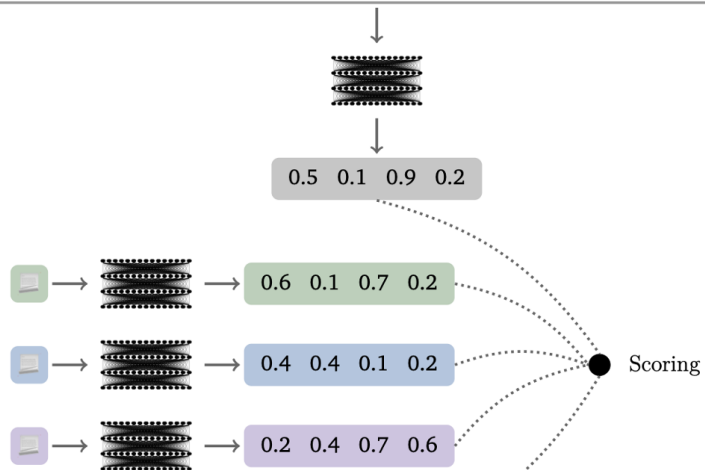
Stanford University was founded in 1891.

- 效率
- 时效性更新
- 溯源
- 理解能力
- 合成表达

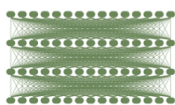


增强关键技术2：检索/知识增强

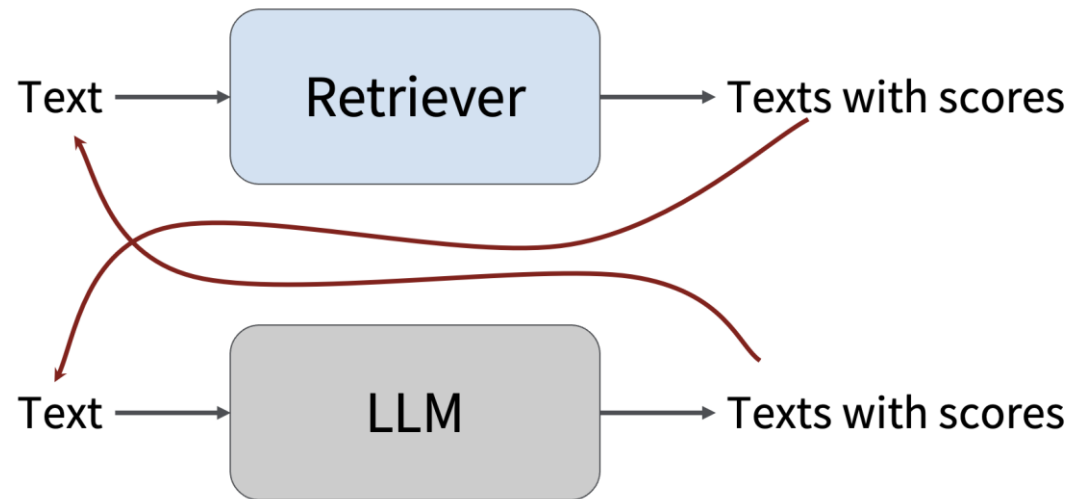
Which MVP of a game Red Flaherty umpired was elected to the Baseball Hall of Fame?



“Red Flaherty umpired in World Series 1955, 1958, 1965, and 1970.” [Red Flaherty](#)
“The 1965 World Series is remembered for MVP Sandy Koufax.” [1965 World Series](#)
“Sandy Koufax was elected to the Baseball Hall of Fame.” [Sandy Koufax](#)



Sandy Koufax, elected to the Hall of Fame in 1972 [\[link\]](#), was the MVP of the 1965 World Series [\[link\]](#), where Red Flaherty was an umpire [\[link\]](#).

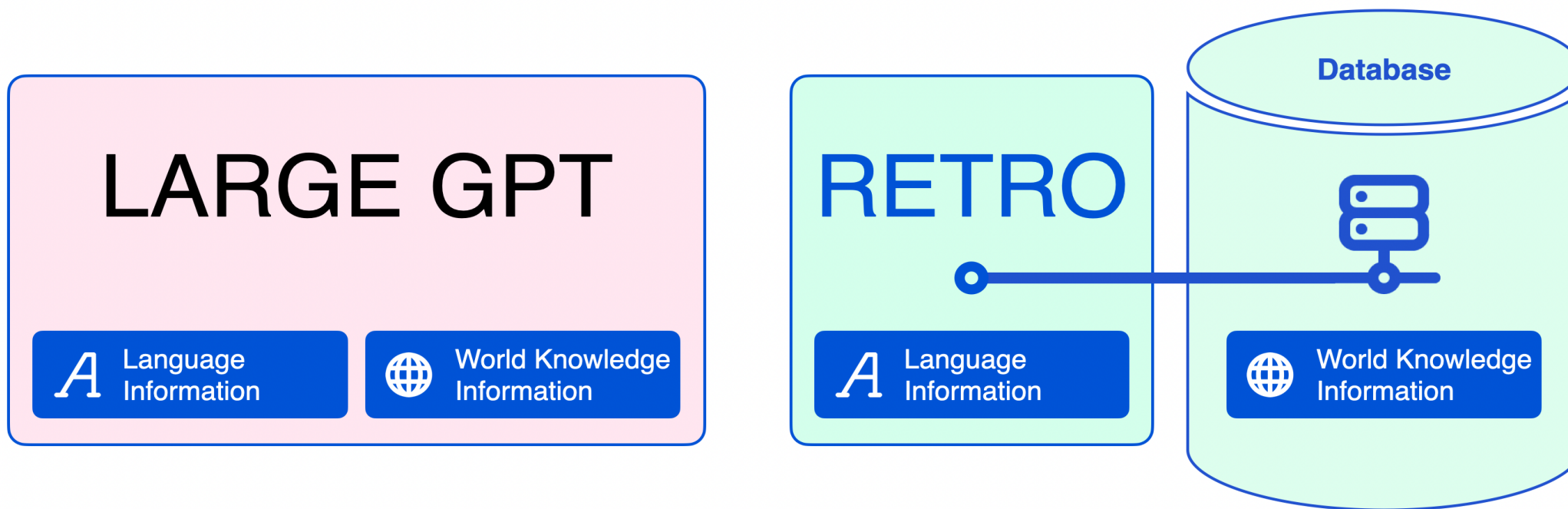


- 效率
- 时效性更新
- 溯源
- 理解能力
- 合成表达



▶ 增强关键技术2：检索/知识增强

- 改进语言模型的另一种途径：通过检索包括网页、书籍、新闻和代码在内的文本段落数据库来增强检索，生成了一种新的语言模型RETRO
- RETRO模型利用从大型语料库中检索到的文档块，基于与前面标记的局部相似性来增强自回归语言模型。该模型可以从零开始训练，也可以快速改装带检索的预训练Transformer，仍然取得良好的性能。



DeepMind's RETRO (Retrieval-Enhanced TRansfOrmer)

增强关键技术2：检索/知识增强

知识增强与结构增强：Structure-inducing Pre-training

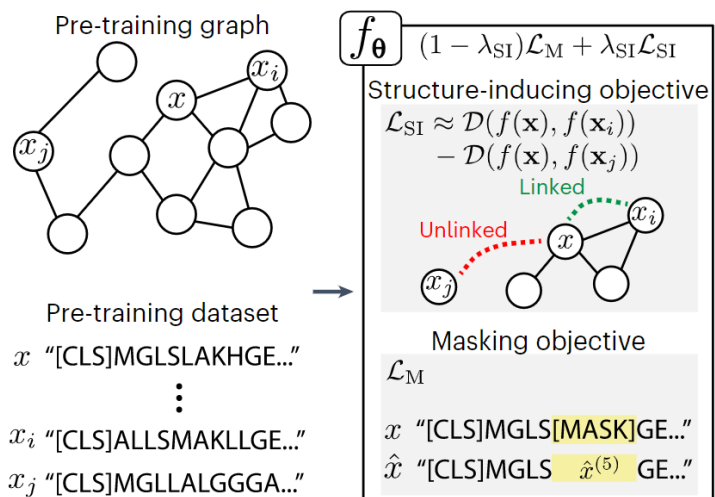
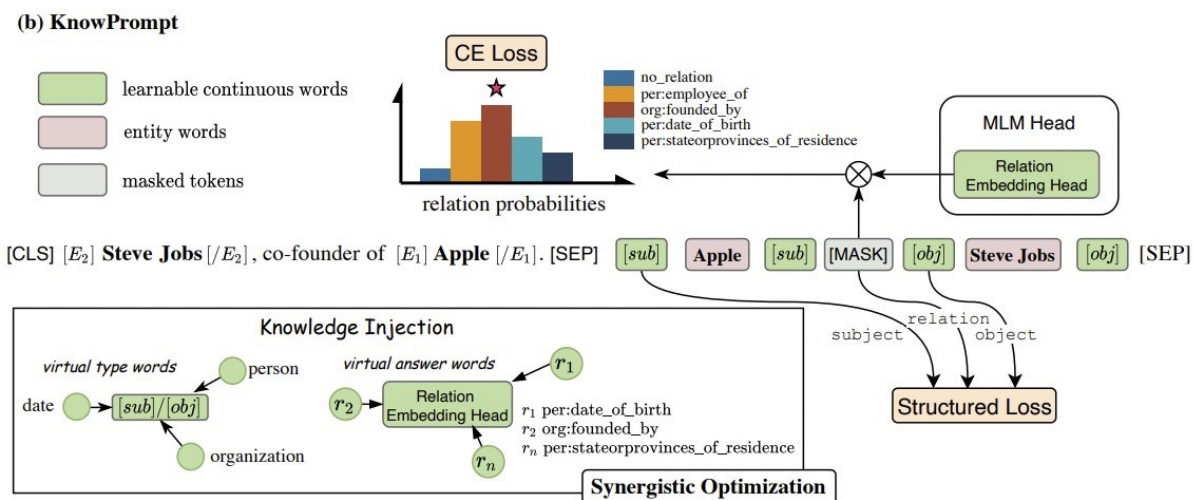


Fig. 1 | Our PT framework. We re-cast the PT formulation by taking a PT graph G_{PT} as an auxiliary input. G_{PT} is used to define a structure-inducing objective \mathcal{L}_{SI} , which pushes a PT encoder f_{θ} to embed samples such that samples are close in the latent space if and only if they are linked in G_{PT} .



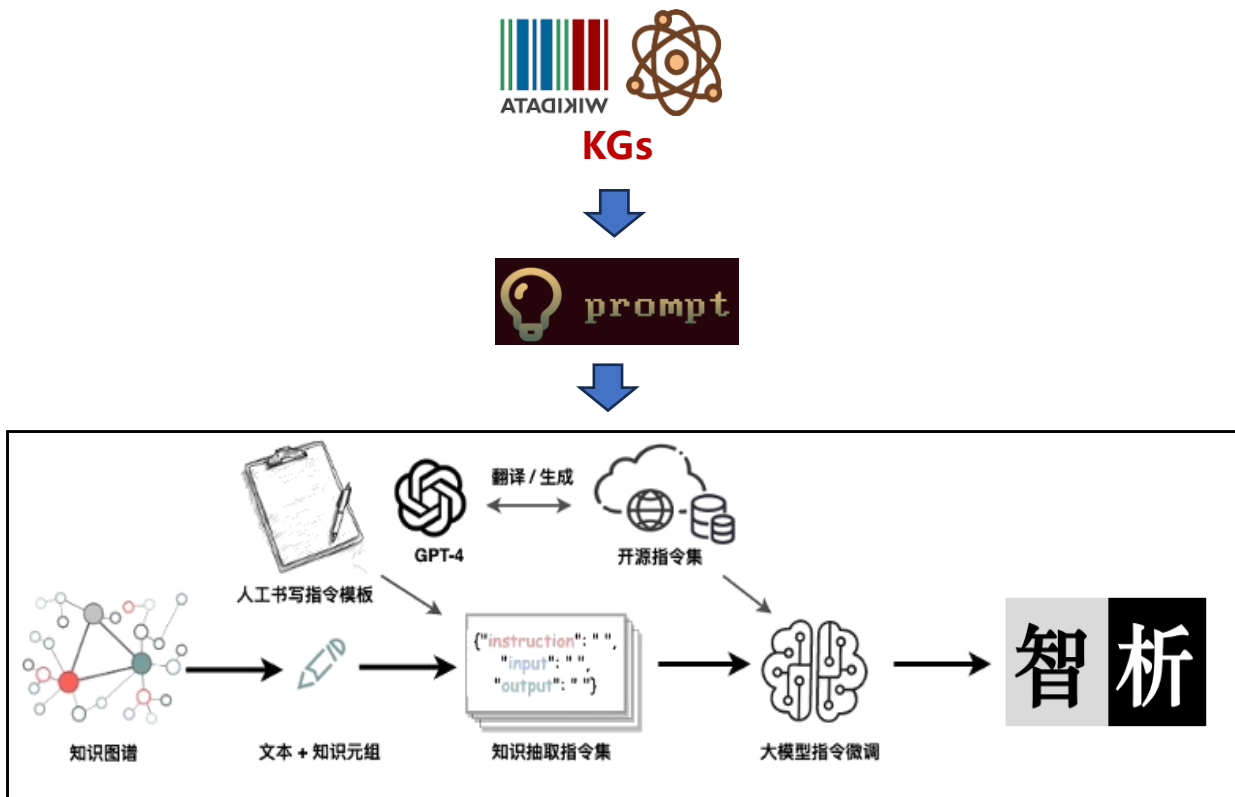
Inter-sample structure

Intra-sample structure

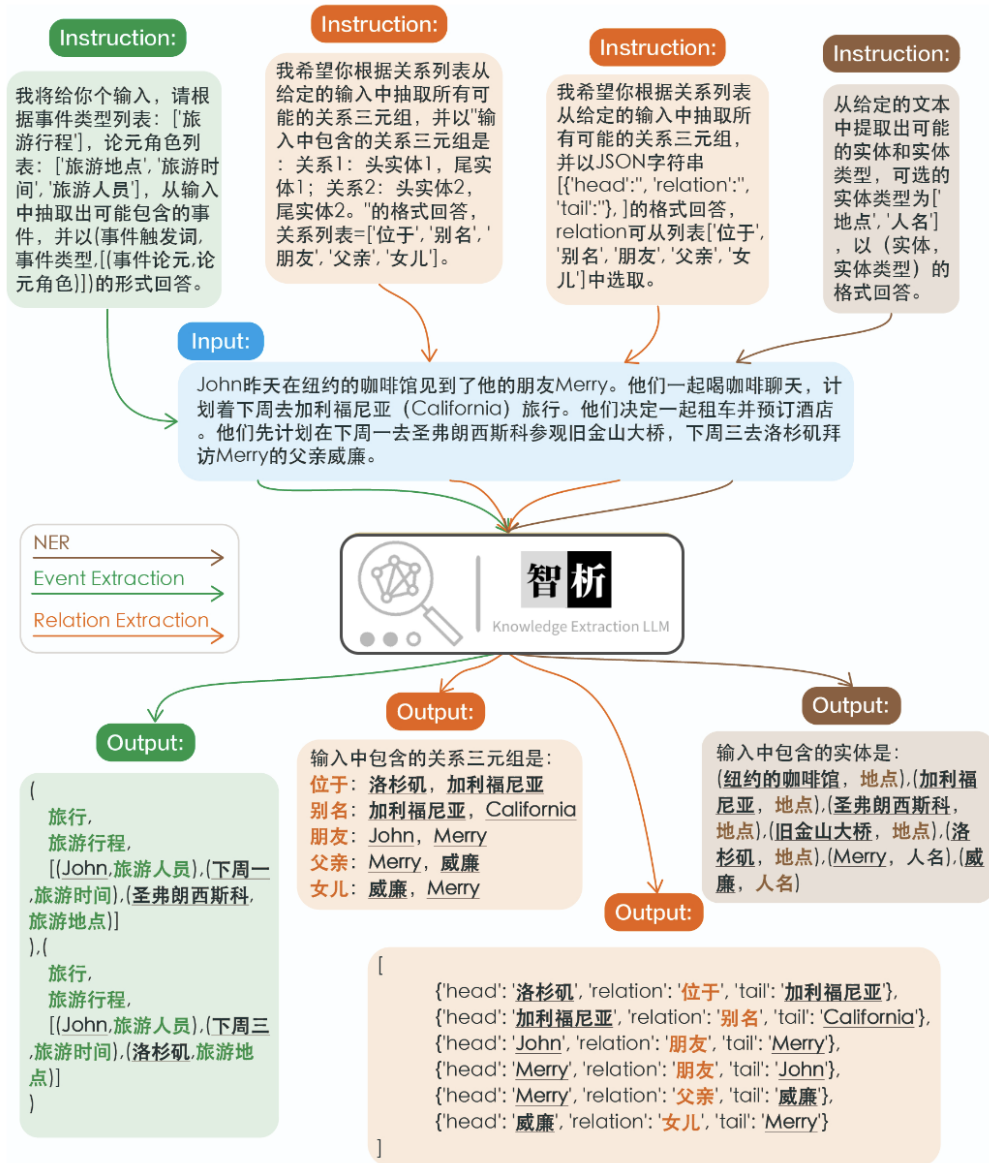
Structure-inducing pre-training. Nature Machine Intelligence 2023

增强关键技术2：检索/知识增强

知识图谱与指令精调：KG2Instruct



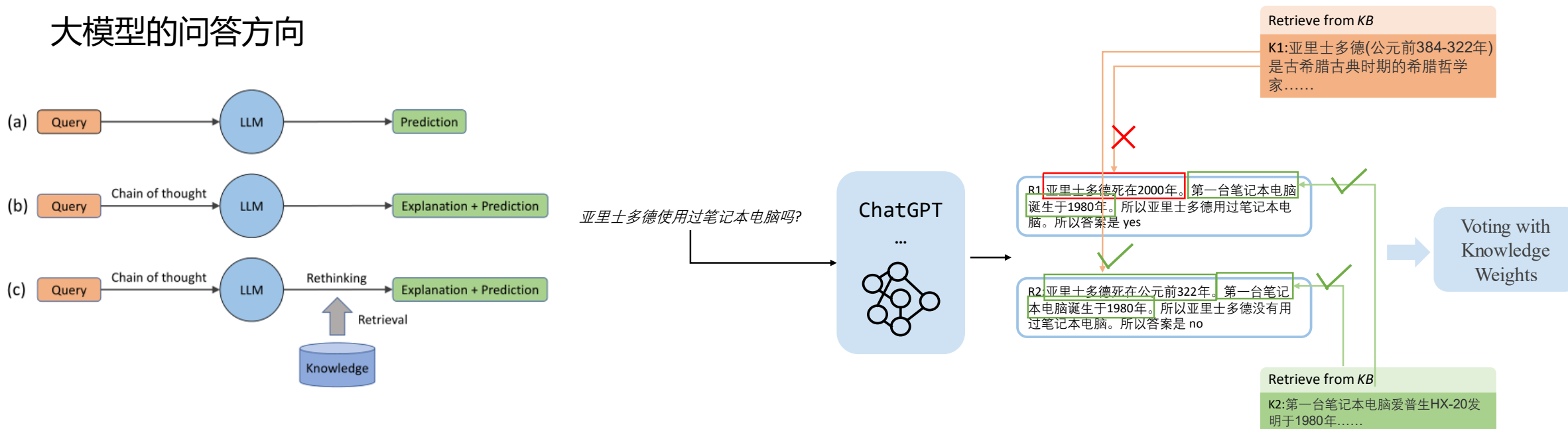
• <http://knowlm.openkg.cn/>



增强关键技术2：检索/知识增强

- 将知识图谱作为预训练模型推理时的外挂知识库，类似检索增强
- 具体的例子： Rethinking

大模型的问答方向

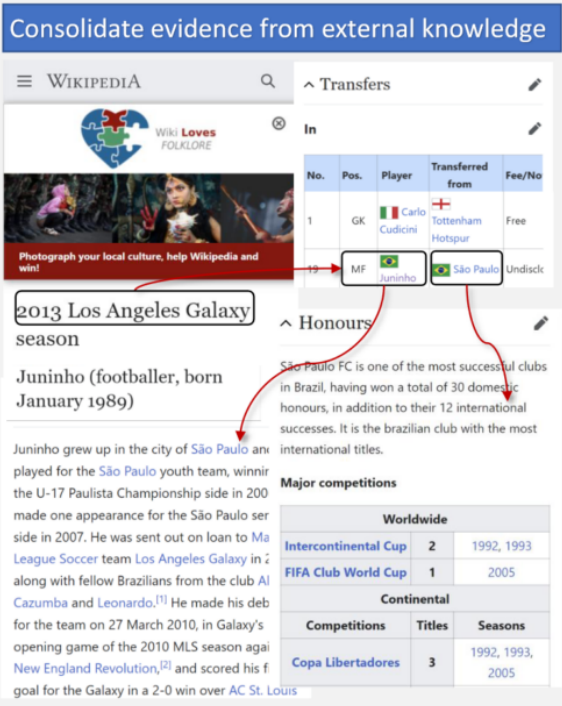


增强关键技术2：检索/知识增强



Which 2013 Los Angeles Galaxy player transferred in from the team with 12 international titles ?

Consolidate evidence from external knowledge



2013 Los Angeles Galaxy season

Juninho (footballer, born January 1989)

Juninho grew up in the city of São Paulo and played for the São Paulo youth team, winning the U-17 Paulista Championship side in 2007. He was sent out on loan to Major League Soccer team Los Angeles Galaxy in 2008 along with fellow Brazilians from the club Al Cazumba and Leonardo.^[1] He made his debut for the team on 27 March 2010, in Galaxy's opening game of the 2010 MLS season against New England Revolution,^[2] and scored his first goal for the Galaxy in a 2-0 win over AC St. Louis.

Worldwide		
Intercontinental Cup	2	1992, 1993
FIFA Club World Cup	1	2005
Continental		
Competitions	Titles	Seasons
Copa Libertadores	3	1992, 1993, 2005

Revise response via automatic feedback

Candidate response:
Jaime Penedo is transferred in from C.S.D. Municipal, a team with 12 international titles.

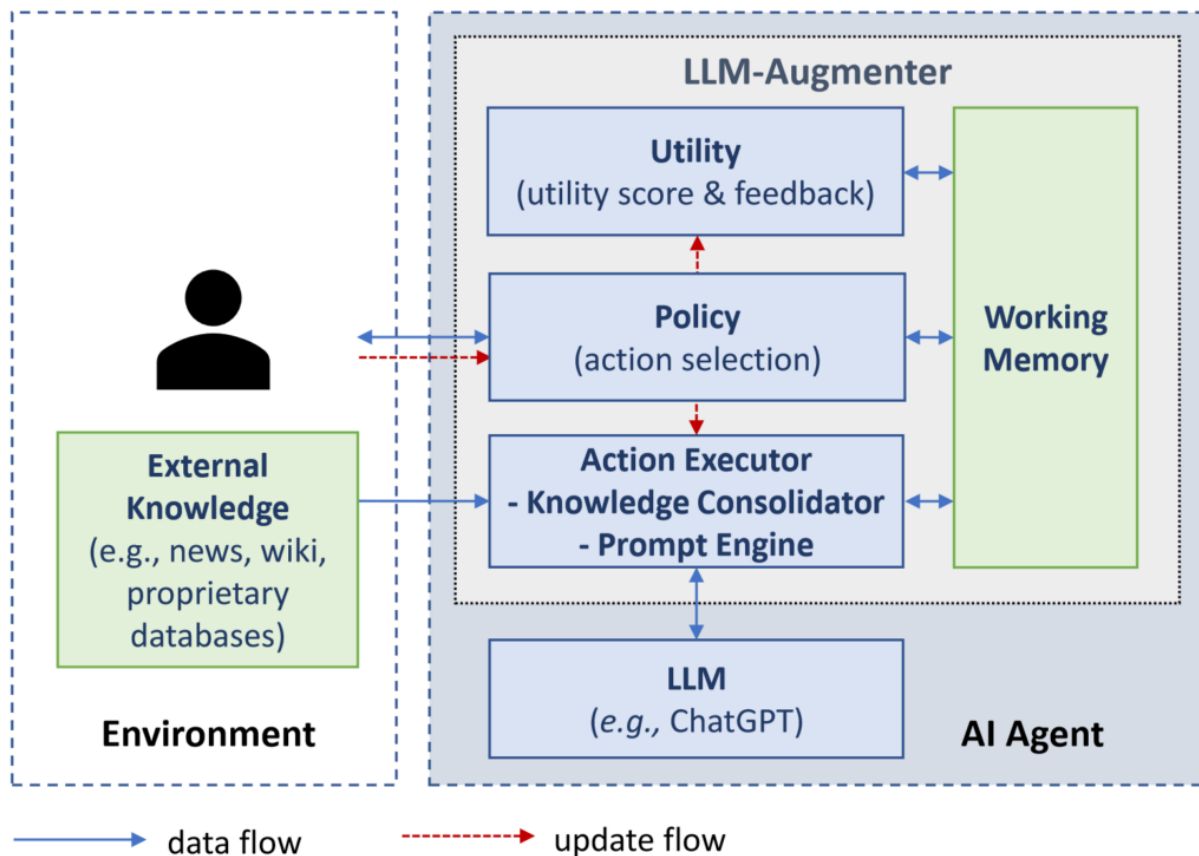
Feedback:
The player Jaime Penedo is transferred in from C.S.D. Municipal, but there is no information about the number of international titles of this team.

Revised candidate response:
Juninho is transferred in from São Paulo, a team with 12 international titles.

AI Agent (LLM-Augmenter + LLM)



Juninho is transferred in from São Paulo, a team with 12 international titles.



Check Your Facts and Try Again: Improving Large Language Models with External Knowledge and Automated Feedback 2023

▶ 增强关键技术2：检索/知识增强

LLM 推理能力

- 代码参与预训练：从底层增强LLM推理能力
- Prompt工程：在使用时激发LLM推理能力

KG 推理能力

- 图关联分析
- 专家推理
- 复杂过程推理
- 时空联合推理
- 案例推理

LLM推理能力

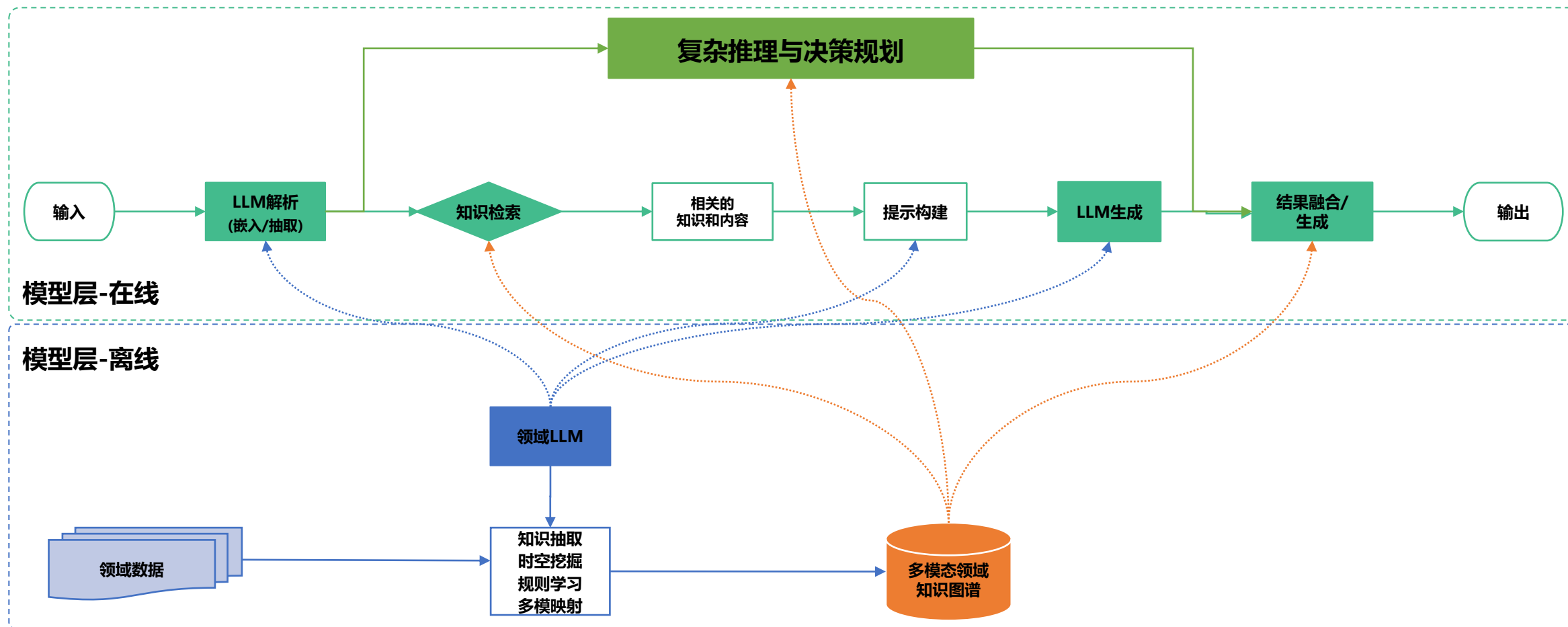
- zero-shot prompting
- Few-shot prompting
- CoT prompting
- Instruct-tuning



KG推理能力

- 图关联分析
- 复杂规则推理
- 业务过程推理
- 时空联合推理
- 案例推理

增强关键技术2：检索/知识增强



增强关键技术2：检索/知识增强

ChatDB: Augmenting LLMs with Databases as Their Symbolic Memory

ChatDB: 用数据库作为符号性记忆模块来增强大语言模型

Chenxu Hu¹, Jie Fu^{2*}, Chenzhuang Du¹, Simian Luo¹, Junbo Zhao³, Hang Zhao^{1*}

¹Tsinghua University ²Beijing Academy of Artificial Intelligence ³Zhejiang University

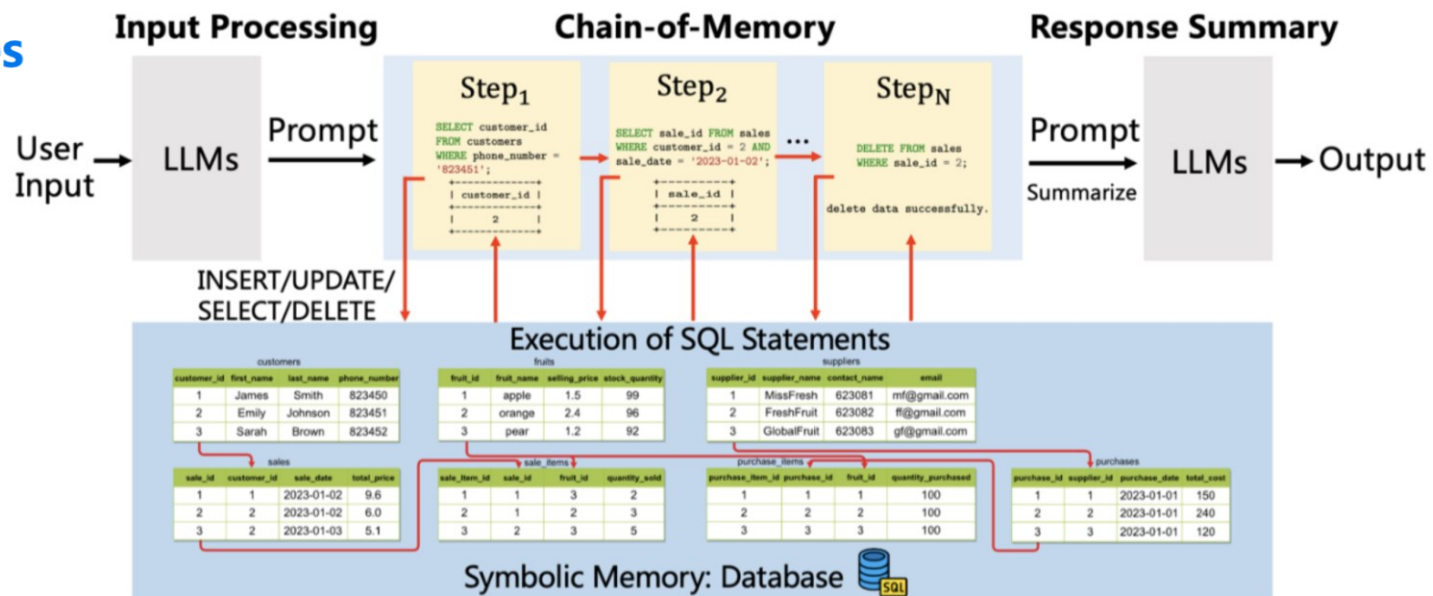
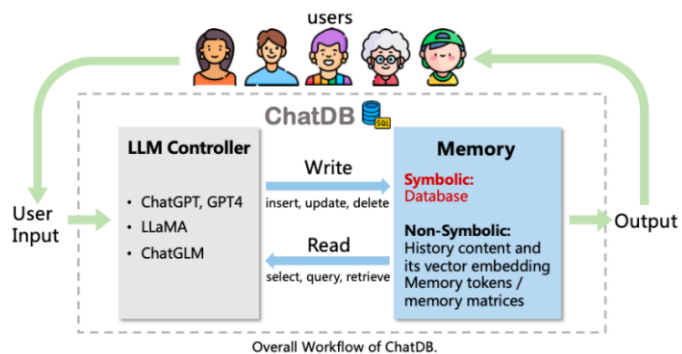


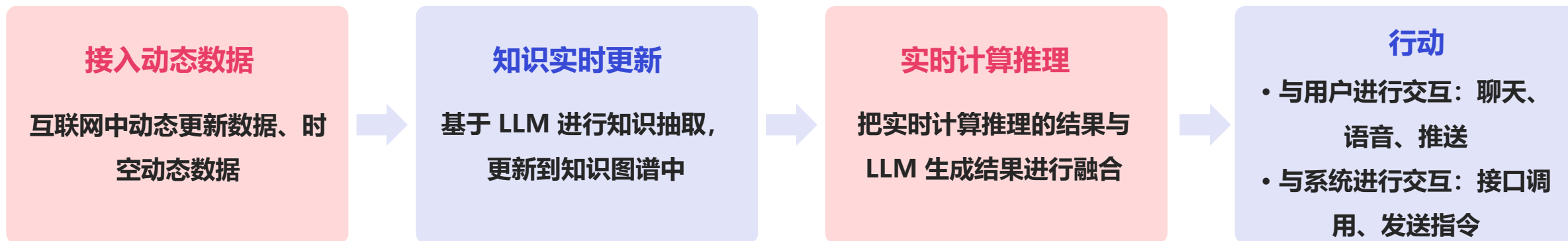
Figure 2: ChatDB framework. The red arrow lines represent the process flow of chain-of-memory, indicating the connection between multiple memory operations. The red arrow lines between database tables represent the reference relationships between primary keys and foreign keys, which start from primary keys to foreign keys. Only the first four columns of each table are shown for brevity. This example showcases the process of returning goods purchased on 2023-01-02 by a customer with the phone number 823451.

增强关键技术3：智能中控

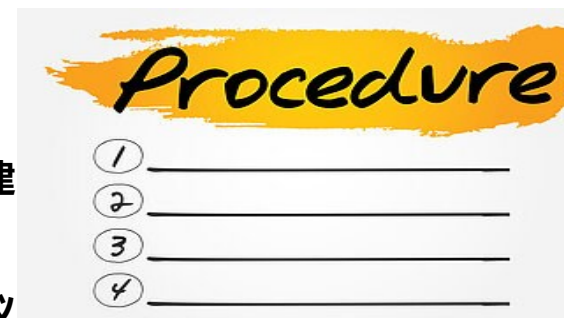
知识图谱以本体实现领域数据、知识及交互的表示，并完成实时数据接入、知识更新到用户交互链路全流程的自动化

数据编排与整合

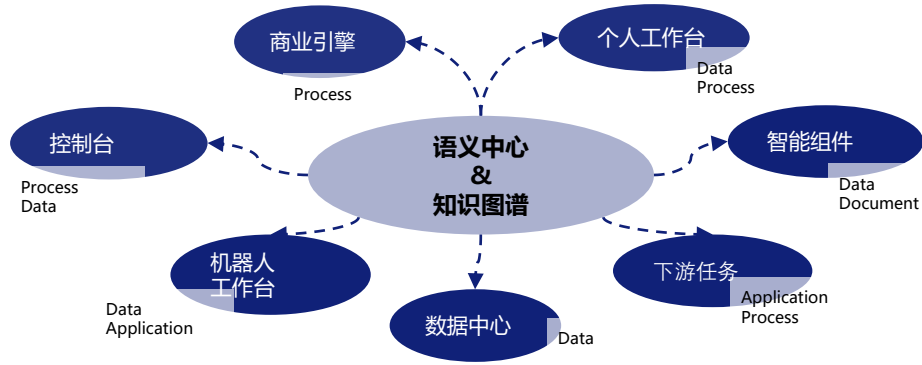
使用编排方式，把以下环节整合形成完整的、端到端的执行流程，实现“超自动化”



- 自动化流程发现
- 一致性检查
- 组织挖掘
- 仿真模型的自动化构建
- 模型扩展与修复
- 案例预测
- 基于历史数据提供建议



增强关键技术3：智能中控 – 跨端主动服务与自适应路由



RPA + Tool use

New Task (Physics QA) Answer this high-school physics question
Input: Hector yanks on the chain with a 72.0 N force at an angle of 35.0° above the horizontal. Determine the horizontal components of the tension force.

TOOL LIBRARY

LLM

OpenAI Codex

python

Input: Hector yanks on the chain with a 72.0 N force at an angle of 35.0° above the horizontal. Determine the horizontal components of the tension force.

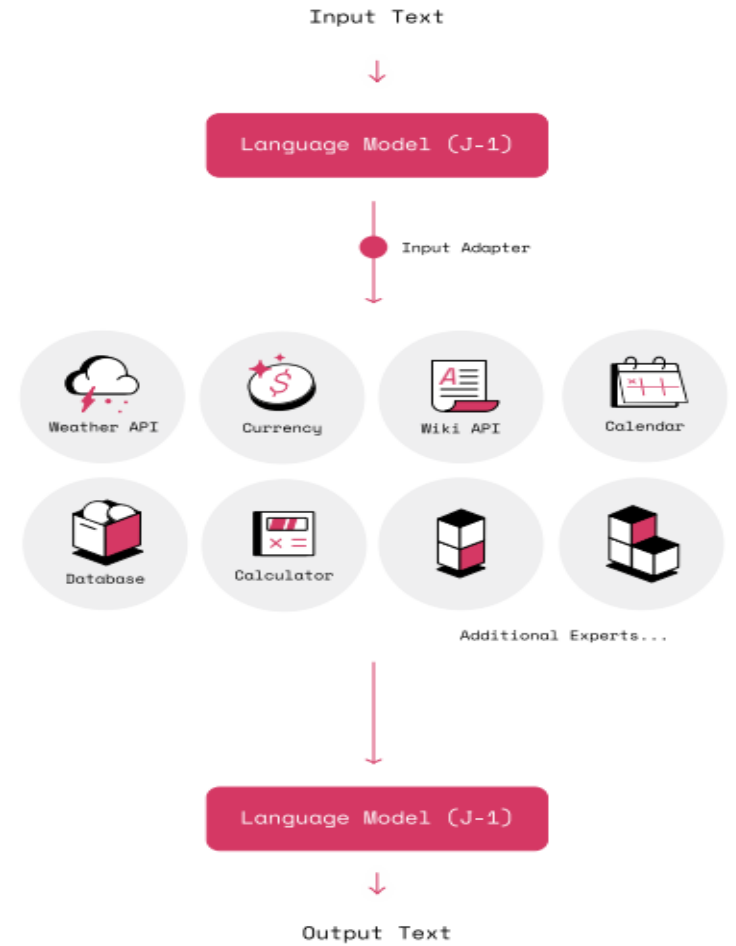
Q1: [search] What is the formula for the horizontal component of the tension force?
#1: The formula for the horizontal component of the tension force is $T \cos \theta$. The horizontal component (F_x) can be calculated as $F \cos(\theta)$ where θ is the angle which the force makes with the horizontal in radians.

Q2: [generate code] Use the formula $F_x = F \cos(\theta)$ to solve: Hank ...
#2: $T = 72.0$, $\theta = 35.0$
 $\text{radians} = \text{math.pi} * \theta / 180$
 $F_x = T * \text{math.cos}(\text{radians})$

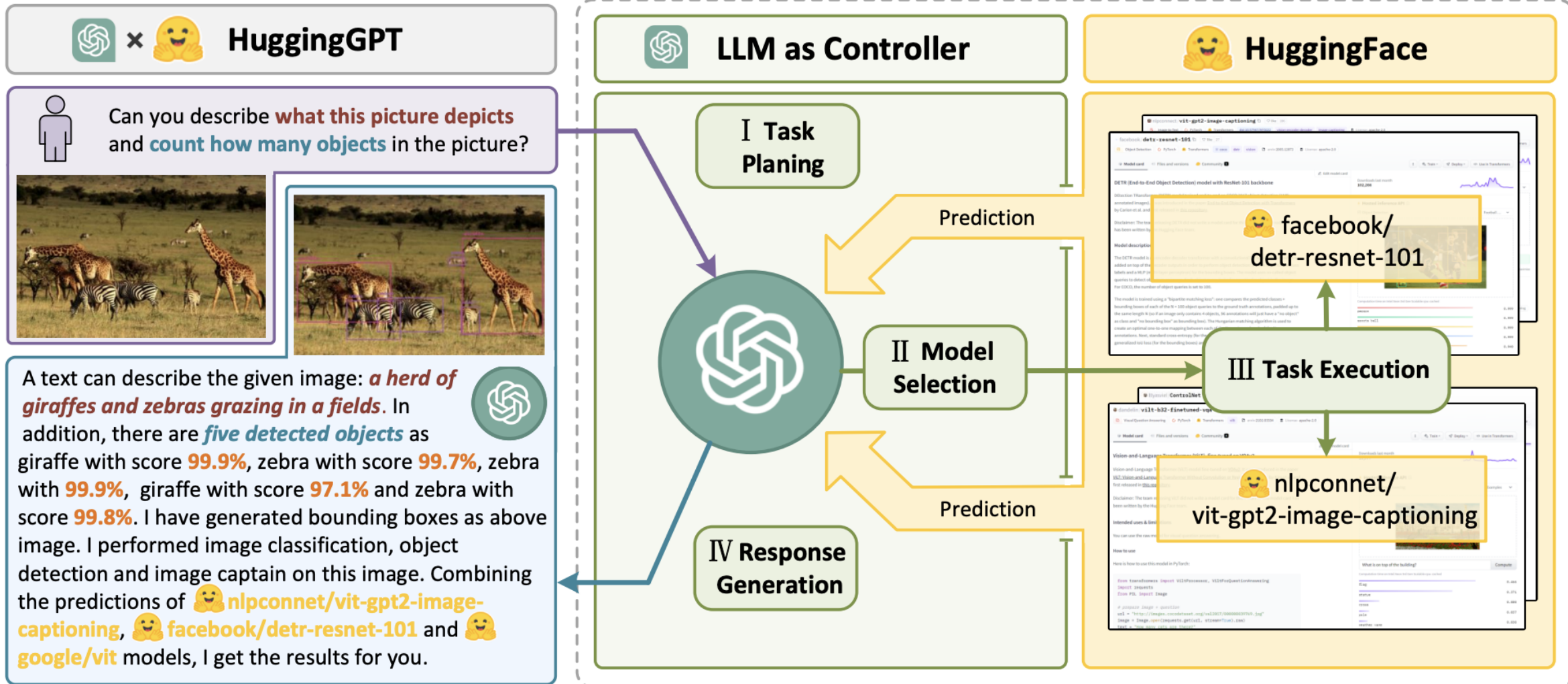
Q3: [code execute] Execute the python code and get the value of "Fx"
#3: 58.9789

Q4: [generate output] Is this an anachronism? #4: No **Q5:** [EOQ] **Ans:** No

自适应路由

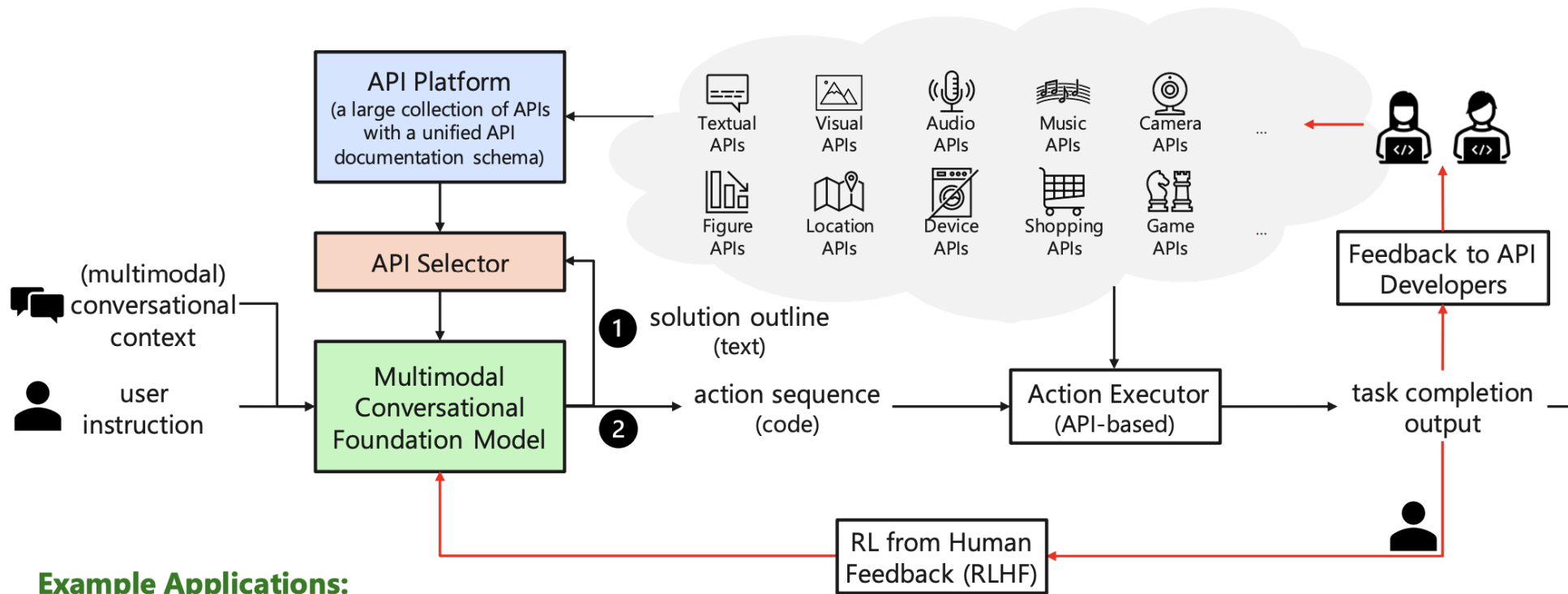


增强关键技术3：智能中控 – 任务规划与模型选择



HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in HuggingFace, 2023

增强关键技术3：智能中控



Example Applications:



AI for Robotics



AI for Office Automation



AI of Things (AIoT)



AI Copilots

TaskMatrix.AI: Completing Tasks by Connecting Foundation Models with Millions of APIs, 2023

增强关键技术3：智能中控 – API链

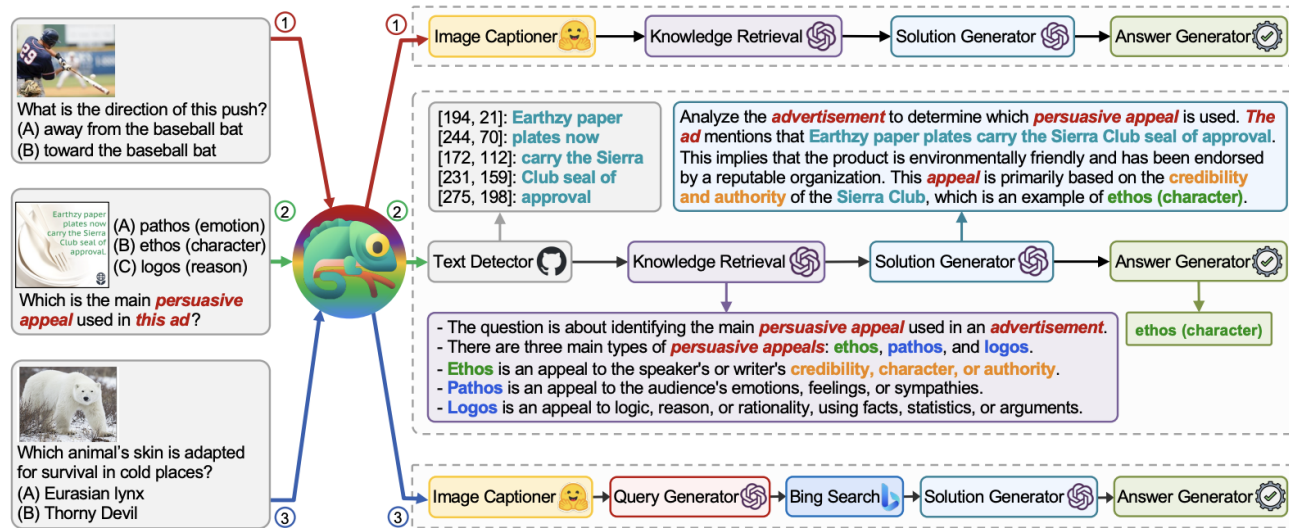


Figure 1: Examples from our Chameleon with GPT-4 on ScienceQA [30], a multi-modal question answering benchmark in scientific domains. Chameleon is adaptive to different queries by synthesizing programs to compose various tools and executing them sequentially to get final answers.

Tool Types	Tools
OpenAI	Knowledge Retrieval, Query Generator, Row Lookup, Column Lookup, Table Verbalizer, Program Generator, Solution Generator
Hugging Face	Image Captioner
Github	Text Detector
Web Search	Bing Search
Python	Program Verifier, Program Executor
Rule-based	Answer Generator

Table 2: Different tools in our module inventory.

•目标：使用LLM与其它工具结合，解决具体领域的问题。在不同类型的数据和各种模型工具之间建立起了桥梁，利用LLM实现了之前需要人工设计的调用顺序和方法。

•当前问题：自然语言大模型LLM由于其自身的限制，无法访问最新信息、无法使用外部工具，无法进行精确的数学推理。

•效果：结合GPT-4，在ScienceQA (86.54%) 和 TabMWP (98.78) 任务中，得到了显著的提升。

•方法：提出chameleon (变色龙)，即插即用的组合推理框架，该框架可以组合多种工具，其中可包含LLM模型、现成的视觉模型、网络搜索引擎、Python 函数和根据用户兴趣定制的基于规则的模块，并将LLM 作为自然语言规划器，将问题拆解成多种工具组合的链条（设计工作流程），然后调用工具协同解决问题，最后通过答案生成器生成回答。

图-1展示了看图回答问题的三个示例，针对第二个问题，展示了从文本识别，信息检索，生成解决方法，最终生成答案的过程。

Chameleon: Plug-and-Play Compositional Reasoning with Large Language Models, 2023

增强关键技术3：智能中控 – 大小模型协同

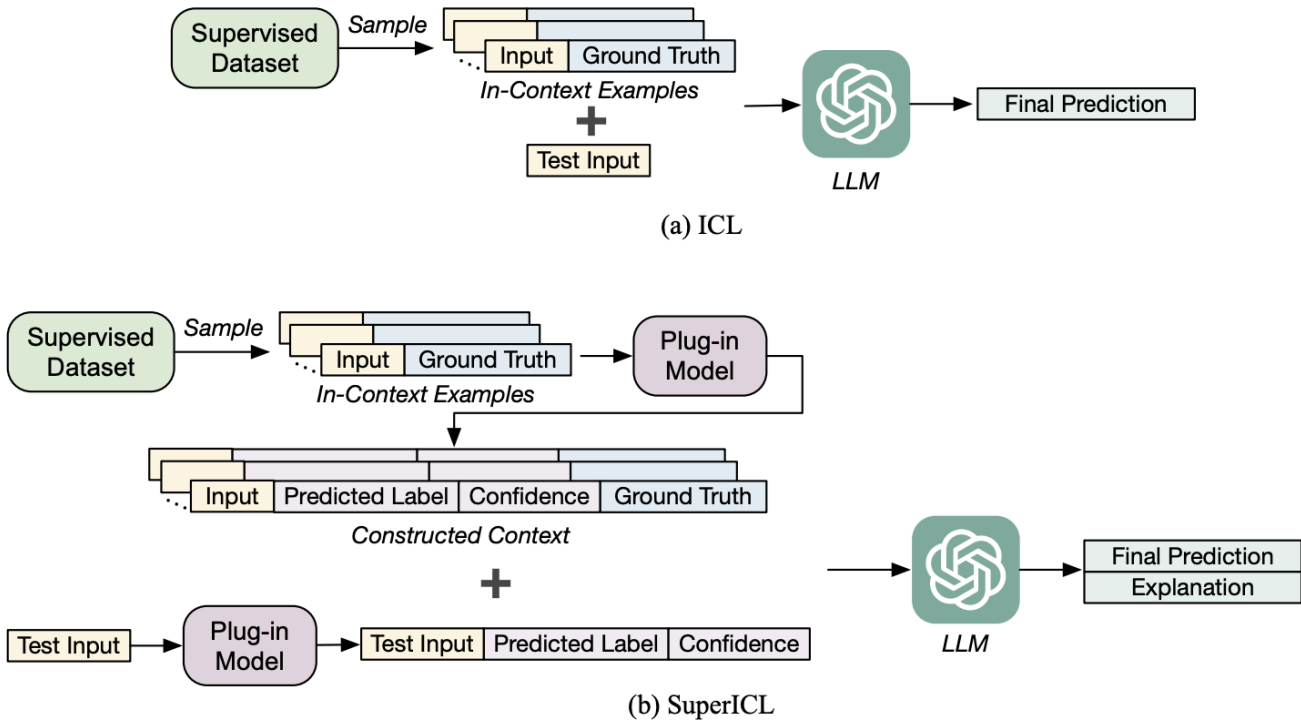


Figure 1: The workflow of ICL and SuperICL. There are three steps in SuperICL: (1) A context is constructed by randomly sampling from the training data and incorporating the plug-in model's predictions, including predicted labels and their corresponding confidence scores. (2) The test input is concatenated after the context, with the plug-in model's prediction attached. (3) Finally, a language model generates the final prediction along with an optional explanation.

•目标：利用自然语言大模型（LLM），提升对大规模的有监督数据的预测效果。

•当前问题：由于上下文长度的限制，只能在对话中给LLM提供有限的上下文提示（In-Context Learning）。

•效果：在效果评测，稳定性，多语言和可解释性方面均表现出其优越性。

•方法
文中提出了SuperICL，将LLM视为黑盒，与本地经过调优的小模型相结合，以提升有监督任务的能力。之前只是将有监督的示例和待预测的测试数据传递给LLM来获得答案。文中提出的方法，首先针对训练集和测试集数据训练了本地模型，预测标签和置信度。然后将这些结果和测试数据一起传递给LLM，从而使LLM不仅学习了推理结果，还学习了决策过程，从而实现了更好的推理和解释能力。

Small Models are Valuable Plug-ins for Large Language Models, 2023

增强关键技术3：智能中控 – 参数化知识增强

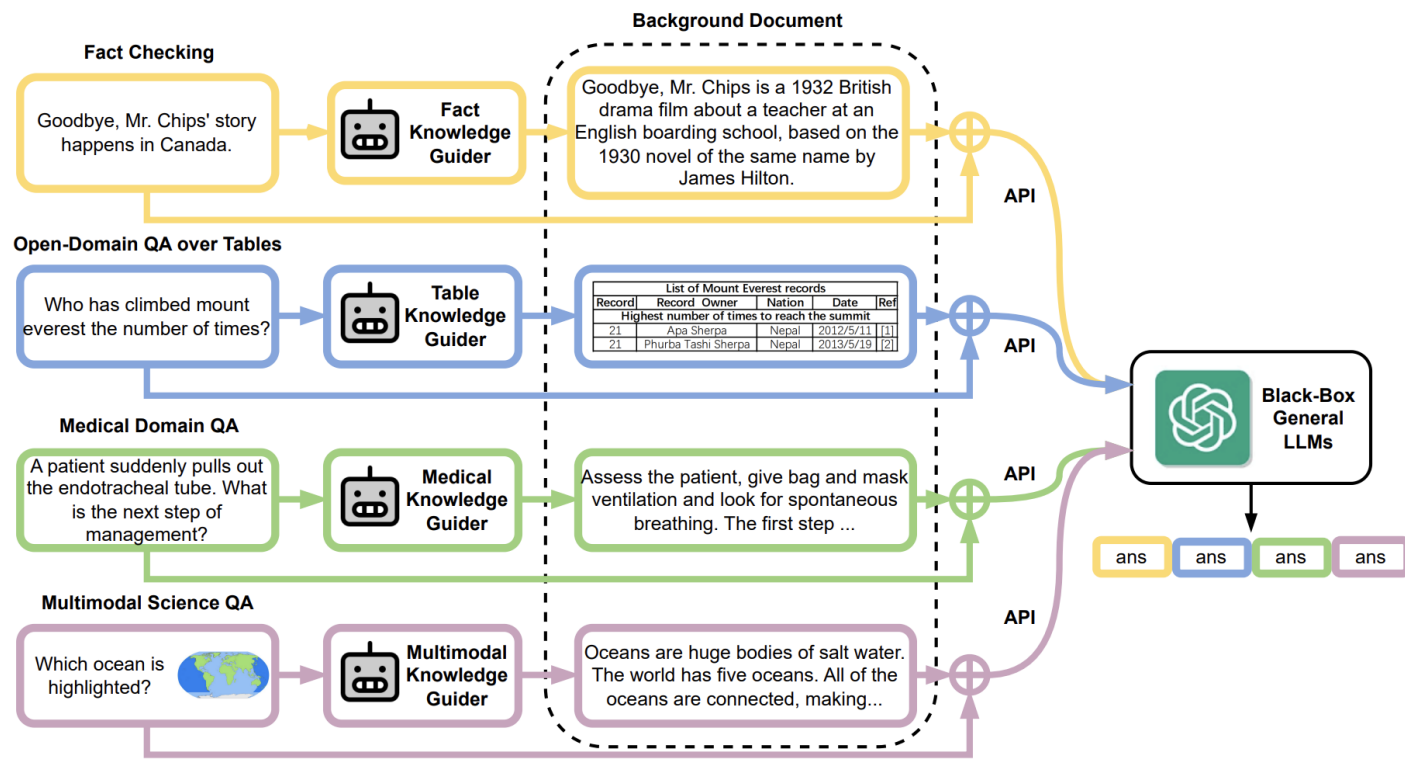


Figure 1: A brief introduction of our parametric knowledge guiding framework (PKG) for augmenting "black box" LLMs on domain-specific tasks.

Augmented Large Language Models with Parametric Knowledge Guiding, 2023

•目标：促进大模型LLM在领域知识密集型任务中的应用

•当前问题：在解决具体问题时，涉及更多领域相关的知识，最新的知识，以及私有数据。

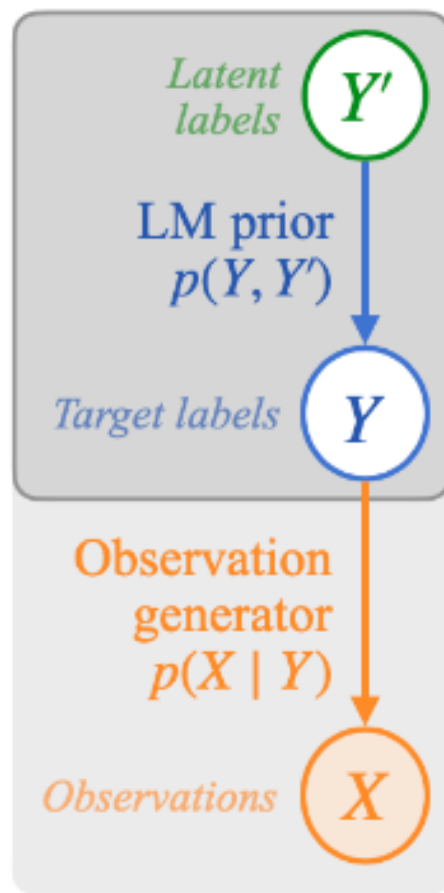
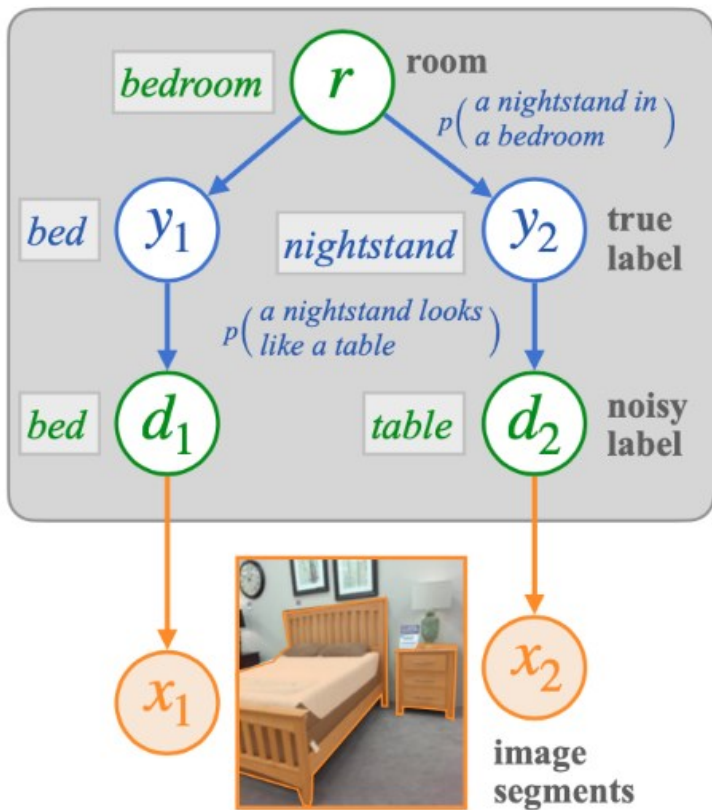
•效果：提升了模型在一系列领域知识密集型任务上的性能，包括事实 (+7.9%)、表格 (+11.9%)、医学 (+3.0%) 和多模态 (+8.1%) 知识。

•方法：提出PKG (Parametric Knowledge Guiding) 参数化知识引导框架，结合本地模型和LLM模型，本地模型基于开源的自然语言模型 (Llama)，它可以存储离线的领域知识，将领域知识转化成参数输出，作为background和问题一起传入大模型。

增强关键技术4：多模态延伸

LAMPP: Language Models as Probabilistic Priors for Perception and Action

Belinda Z. Li¹ William Chen¹ Pratyusha Sharma¹ Jacob Andreas¹

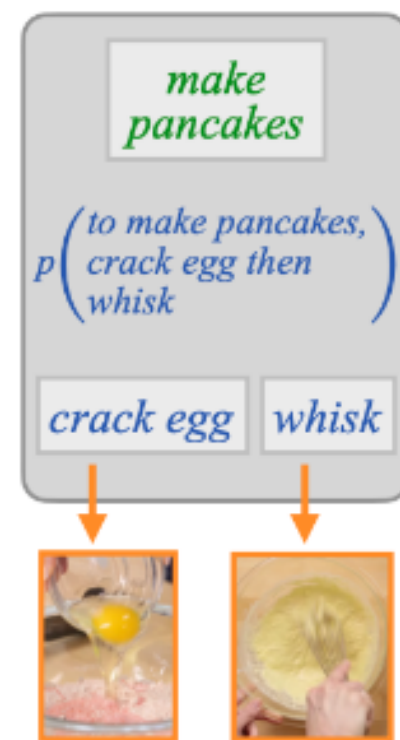


(a)



Image Segmentation

(b)



Action Recognition

(c)

增强关键技术4：多模态与具身延伸

PaLM-E: 具身多模态语言模型

PaLM-E: An Embodied Multimodal Language Model

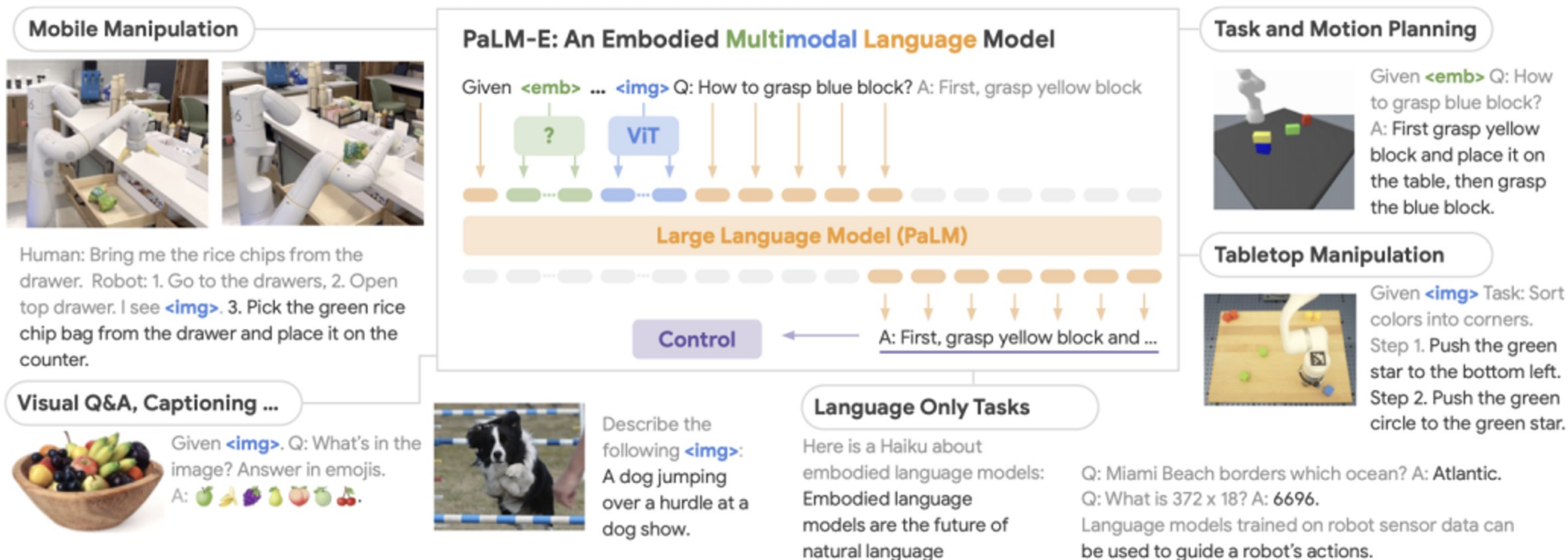
Danny Driess^{1,2} Fei Xia¹ Mehdi S. M. Sajjadi³ Corey Lynch¹ Aakanksha Chowdhery³
Brian Ichter¹ Ayzaan Wahid¹ Jonathan Tompson¹ Quan Vuong¹ Tianhe Yu¹ Wenlong Huang¹
Yevgen Chebotar¹ Pierre Sermanet¹ Daniel Duckworth³ Sergey Levine¹ Vincent Vanhoucke¹
Karol Hausman¹ Marc Toussaint² Klaus Greff³ Andy Zeng¹ Igor Mordatch³ Pete Florence¹

¹Robotics at Google ²TU Berlin ³Google Research

<https://palm-e.github.io>

PaLM-E 将真实世界的传感器信号与文本输入相结合，建立语言和感知的链接。规模最大的模型“PaLM-E-562B”具有562B个参数，将540B的PaLM和22B的ViT集成在一起，这是目前报道的最大的视觉-语言模型。模型输入包括视觉、连续状态估计值和文本输入。作者在多个任务（包括顺序机器人操作规划、视觉问答和字幕生成）中进行了端到端的训练，并通过评估表明，其模型能够有效地解决各种推理任务，并且在不同的观察模态和多个实体上表现出了积极的转移。该模型在进行机器人任务训练的同时，还具有先进的视觉-语言任务表现，并随着规模的增大保持了通用的语言能力。

它支持多模态输入，来自任意模态（例如图像、三维表示或状态，绿色和蓝色）的输入插入文本 token（橙色）旁边作为LLM的输入，进行端到端的训练。



增强关键技术4：多模态与具身延伸

I spilled my drink, can you help?

GPT3

You could try using a vacuum cleaner.

LaMDA

Do you want me to find a cleaner?

FLAN

I'm sorry, I didn't mean to spill it.

I spilled my drink, can you help?

LLM

"find a cleaner"
"find a sponge"
"go to the trash can"
"pick up the sponge"
"try using the vacuum"

Value Functions

"find a cleaner"
"find a sponge"
"go to the trash can"
"pick up the sponge"
"try using the vacuum"



SayCan

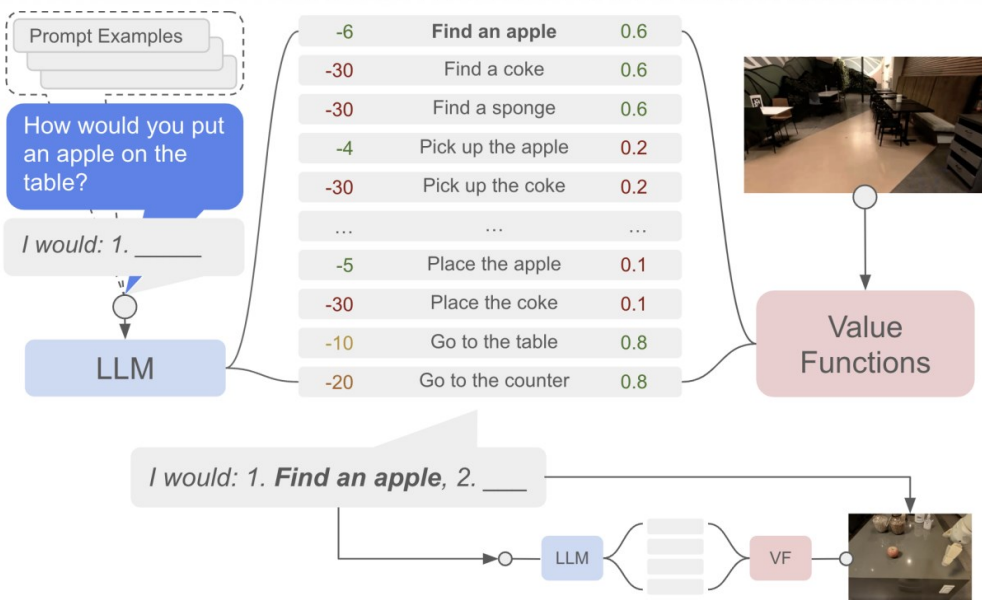
"find a cleaner"
"find a sponge"
"go to the trash can"
"pick up the sponge"
"try using the vacuum"



I would:

1. find a sponge
2. pick up the sponge
3. come to you
4. put down the sponge
5. done

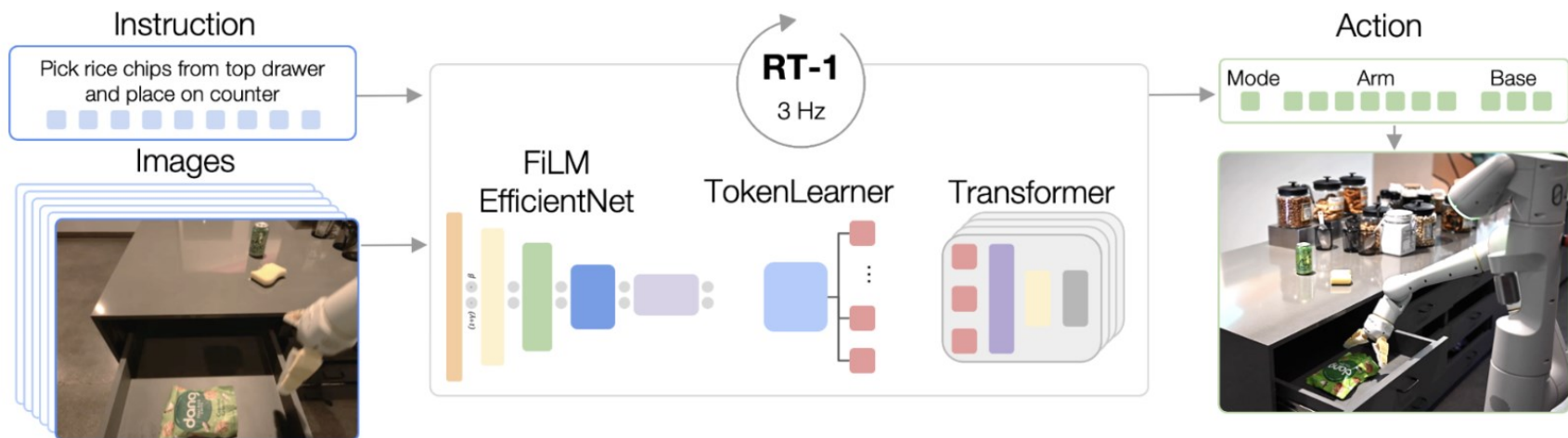
大模型可以把高层级的任务拆分成若干个在语义上符合逻辑的子任务，但是由于LLMs普遍缺乏真实物理世界的经验，无法判断其输出会对环境产生什么样的影响，也不知道真实环境和机器人的状态信息以及机器人是否具备执行这些子任务的能力，**所以其生成的看似逻辑上合理的子任务指令不一定能在某个实际的场景中被机器人顺利执行。**



SayCan的设计逻辑很简单，将决定机器人应该如何执行任务的决策拆成两个部分，**Say**代表大模型LLM，用于输出可用的高层级运动指令，**Can**代表机器人在当前环境下能做的事情，二者通过值函数 (Value Function) 的方式结合起来，共同决定选择哪条指令用于实际执行。

Do As I Can, Not As I Say: Grounding Language in Robotic Affordances, 2023

增强关键技术4：多模态与具身延伸



RT-1使用自然语言和图像作为输入，输出离散的机械臂和底盘动作（位置和姿态）

RT-1研究两个事情，第一**通过Transformer Model学习得到机器人的技能**，**第二如何使用自然语言控制机器人的运动**。

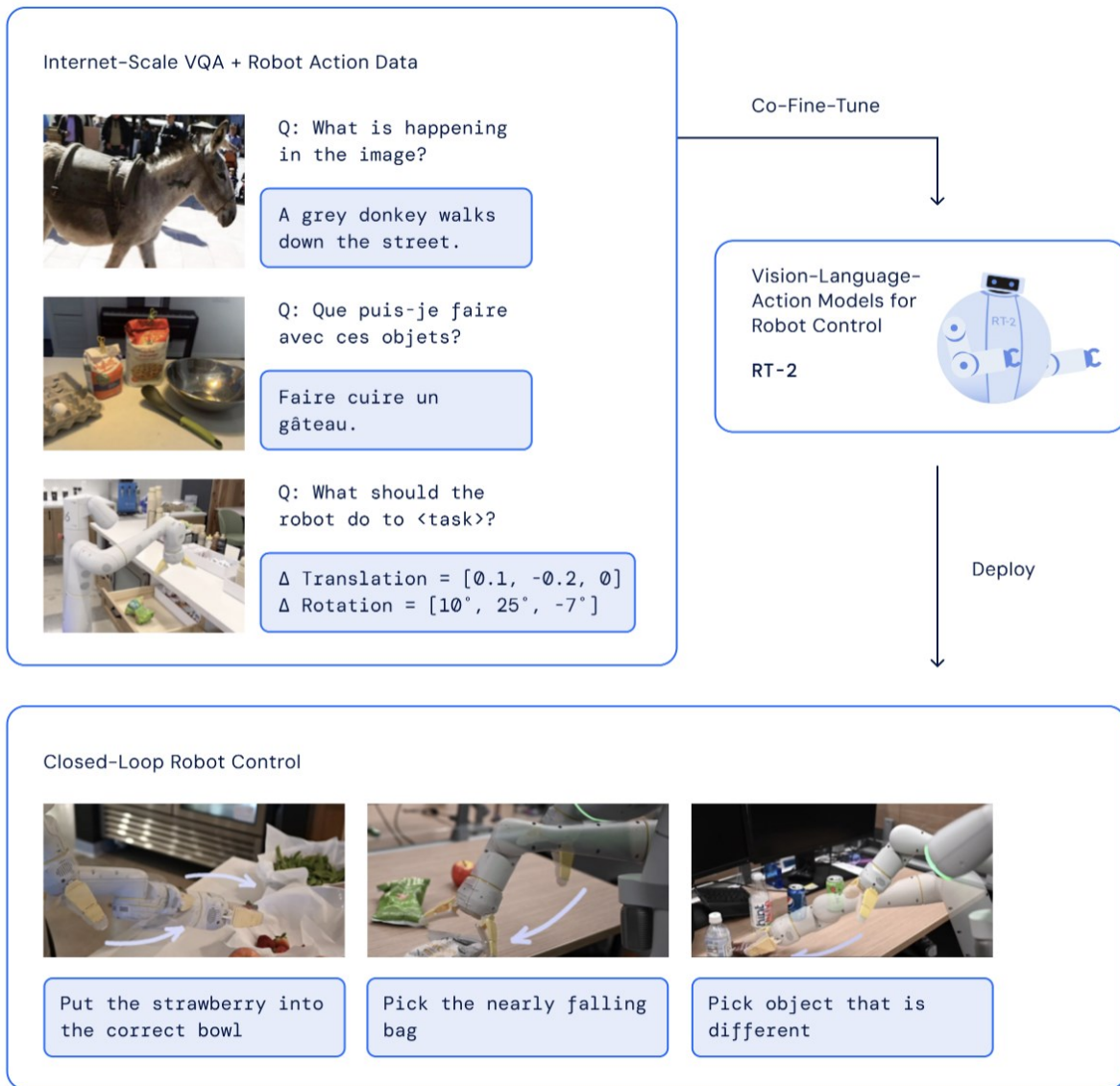
RT-1不是传统意义上的大模型，只是主体使用了Transformer架构，总体参数量只有35M，Transformer部分的参数量只有19M，是面向机器人操作单独训练的一个多任务机器人控制模型。

前面讲的SayCan用大模型做任务理解和任务拆分（谷歌这帮人称之为High-Level）把一个任务拆解为若干个子任务（Sub-task）或者技能（Skill），RT-2用大模型做技能学习（Skill, Low-level），学习到的技能用于支持SayCan的任务执行。



RT-1: Robotics Transformer for real-world control at scale, Dec 2022

增强关键技术4：多模态与具身延伸



RT-2的目的是为了研究将使用互联网规模数据 (Internet-scale data) 训练得到的VLM大模型直接用于端到端的机器人控制，提升机器人操作的泛化能力和语义推理能力。

RT-2抛弃了从头训练Transformer模型的方式，而是直接采用已有的VLM模型作为主模型，再使用更适合机器人任务的fine-tuning方法对其进行微调。简单来讲，RT-2是使用机器人数据集进行fine-tuning且其输出的text被设计成了机器人位置指令形式的VLM (PaLM-E、PaLI-X)。

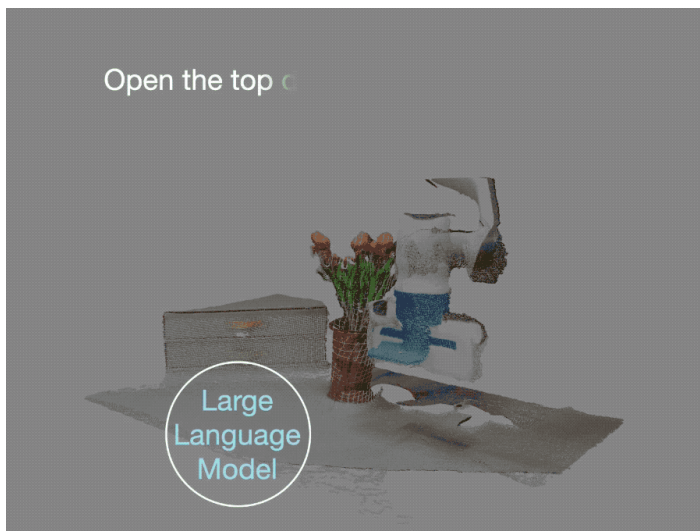
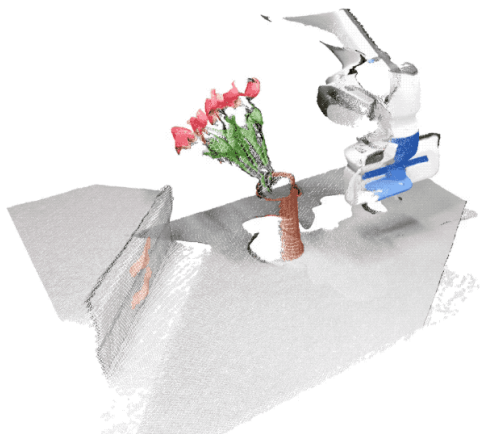
为了更符合机器人的应用，RT-2主要研究的是利用大模型的优势直接生成Low-level的机器人运动指令的内容。

RT-2证明了使用机器人技能数据集对已有的LLMs或者VLMs进行微调，可以快速的利用VLM的海量通识能力，大幅提升机器人的任务执行成功率和泛化能力。

因为RT-2的Backbone模型是PaLM-E 540B，所以无法运行在机器人的硬件上，只能部署在云端，而RT-1可以直接运行在机器人上。

RT-2: New model translates vision and language into action, July 2023

增强关键技术4：多模态与具身延伸



首先，给定环境信息（用相机采集RGB-D图像）和我们要执行的自然语言指令。

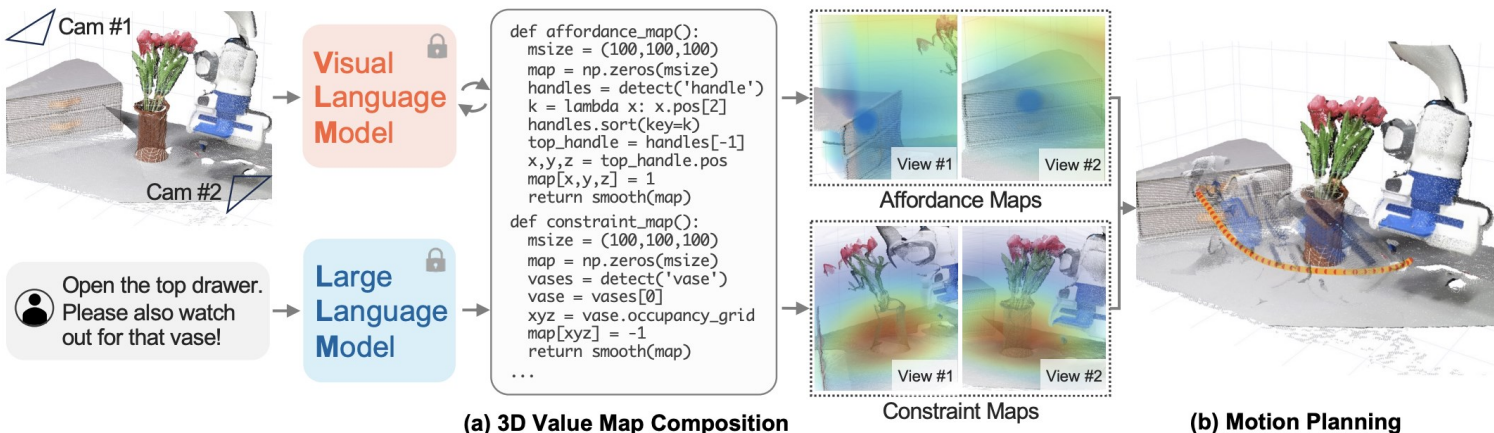
接着，LLM（大语言模型）根据这些内容编写代码，所生成代码与VLM（视觉语言模型）进行交互，指导系统生成相应的操作指示地图，即**3D Value Map**。

所谓3D Value Map，它是Affordance Map和Constraint Map的总称，既标记了“在哪里行动”，也标记了“如何行动”。

无论是哪种环境哪种情况（有无干扰、指令是否可见），它都显著高于基于原语的基线任务。

还惊喜地发现，VoxPoser产生了**4个“涌现能力”**：

- (1) 评估物理特性，比如给定两个质量未知的方块，让机器人使用工具进行物理实验，确定哪个块更重；
- (2) 行为常识推理，比如在摆餐具的任务中，告诉机器人“我是左撇子”，它就能通过上下文理解其含义；
- (3) 细粒度校正，比如执行“给茶壶盖上盖子”这种精度要求较高的任务时，我们可以向机器人发出“你偏离了1厘米”等精确指令来校正它的操作；
- (4) 基于视觉的多步操作，比如叫机器人将抽屉精准地打开一半，由于没有对象模型导致的信息不足可能让机器人无法执行这样的任务，但VoxPoser可以根据视觉反馈提出多步操作策略，即首先完全打开抽屉同时记录手柄位移，然后将其推回至中点就可以满足要求了。



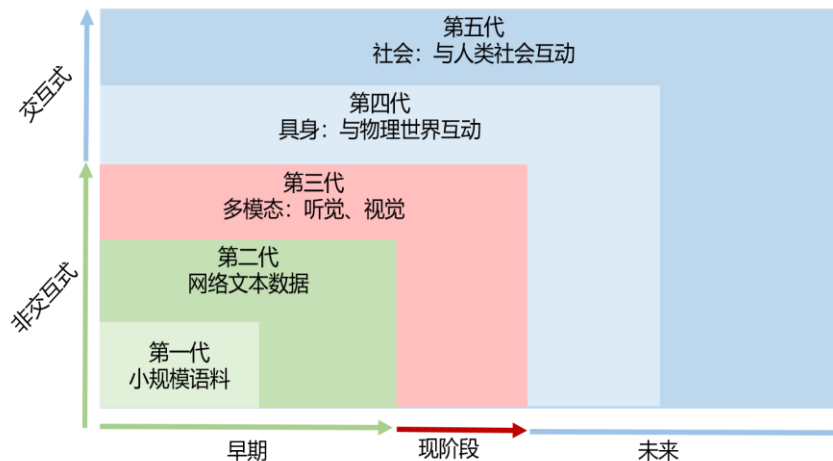
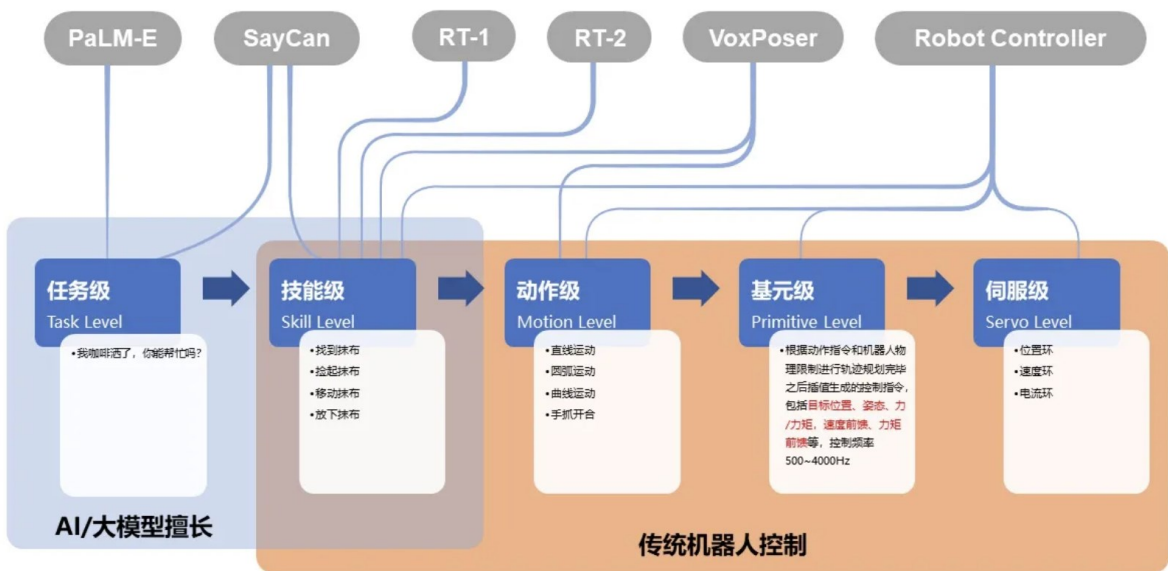
VoxPoser: Composable 3D Value Maps for Robotic Manipulation with Language Models, 2023

增强关键技术4：多模态与具身延伸

□ GPT-X尝试解决与现实世界联结问题

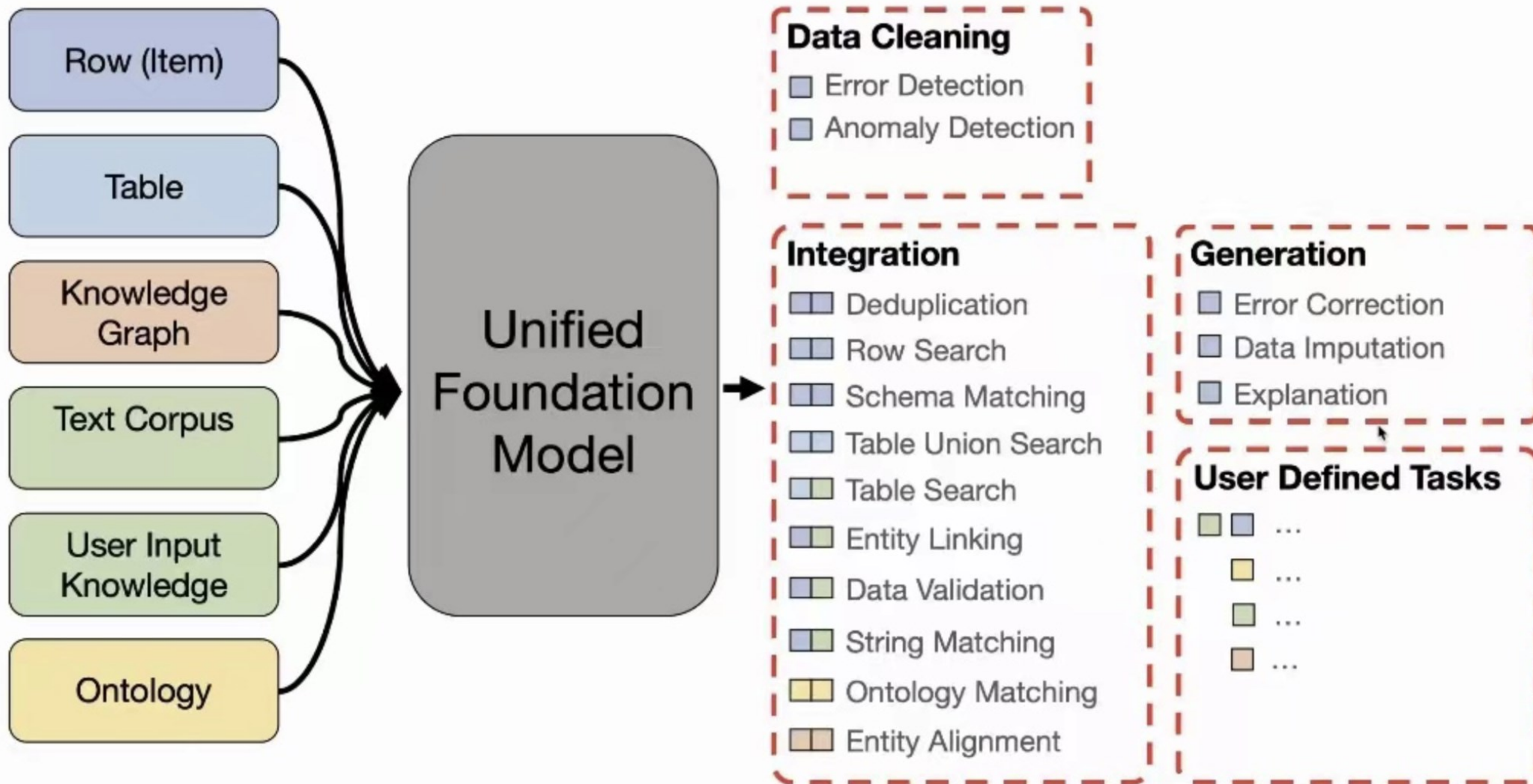
- ChatGPT 解决人与机器的自然语言交互问题；
- GPT-4 突破多模态，从认知世界向融入环境的跨越；
- 近期LLM与人机交互，为机器人生成执行指令[1]
- 未来GPT-X甚至可能产生情感、意识等特性，**并最终通过广义图灵测试**

能力进化速度越来越快



自主智能
人机混合智能
多模态智能
大数据智能

▶ 其他趋势：数据治理也是一种文字生成任务



▶ 其他趋势：数据敏感的场外微调

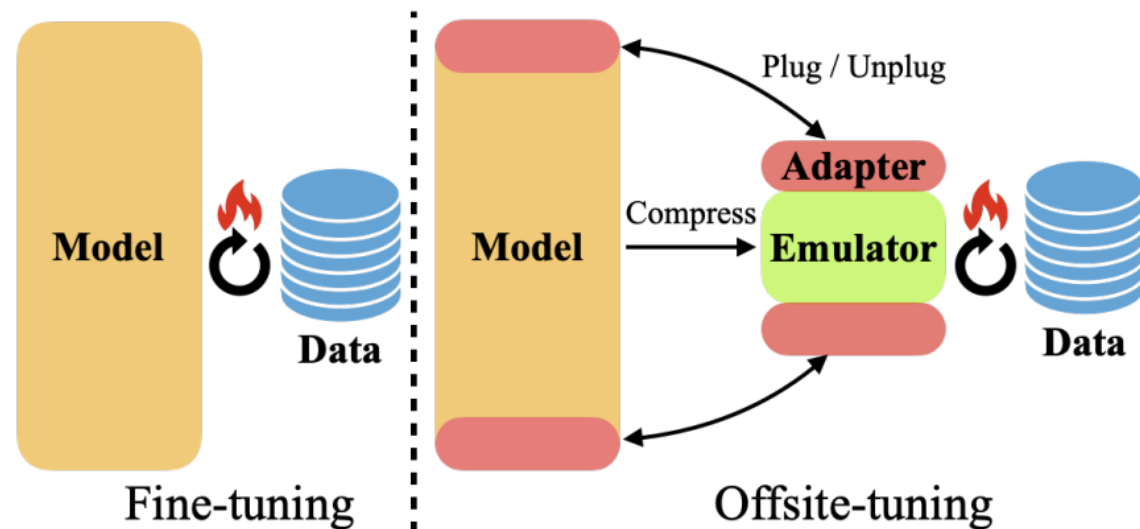
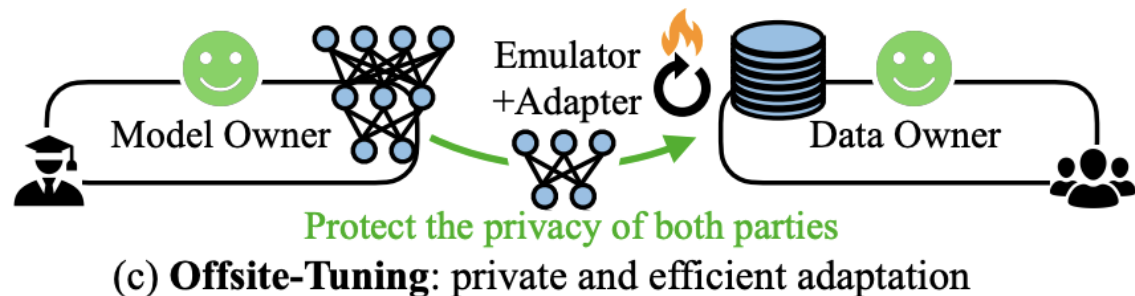
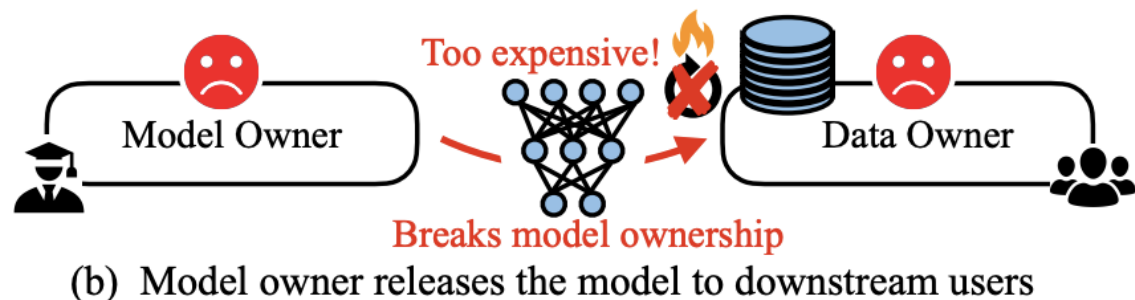
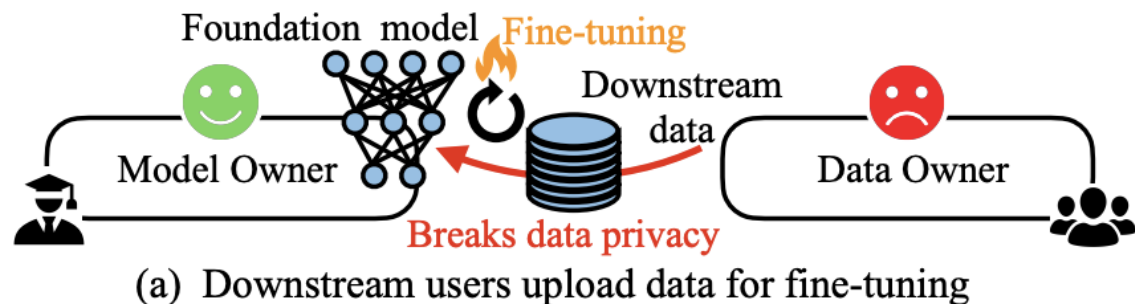


Figure 2. Overview of Offsite-Tuning. Fine-tuning (left) requires access to the full model weights and needs both model and data to be in one location. In Offsite-tuning (right), the model owner sends an adapter and an emulator to the data owner, who fine-tunes the adapter on the downstream data with the emulator's assistance. The fine-tuned adapter is then returned and plugged into the full model to create an adapted foundation model. As neither party needs to share full models or data and the emulator is compressed, offsite-tuning is both privacy-preserving and efficient.

Offsite-Tuning: Transfer Learning without Full Model, 2023

▶ KnowLM: Knowledgeable LLM framework

知识编辑 技术特色1

基于**知识编辑**技术对齐大模型内过时、错误及价值观不正确的知识，解决**知识谬误**问题



知识提示 技术特色2

基于知识图谱等结构化数据的**知识提示生成**和**知识增强约束**技术，解决**知识抽取和推理**问题

知识交互 技术特色3

基于知识动态交互和反馈实现**工具组合学习**及**多智能体协作**，解决大模型**具身认知**问题

.....



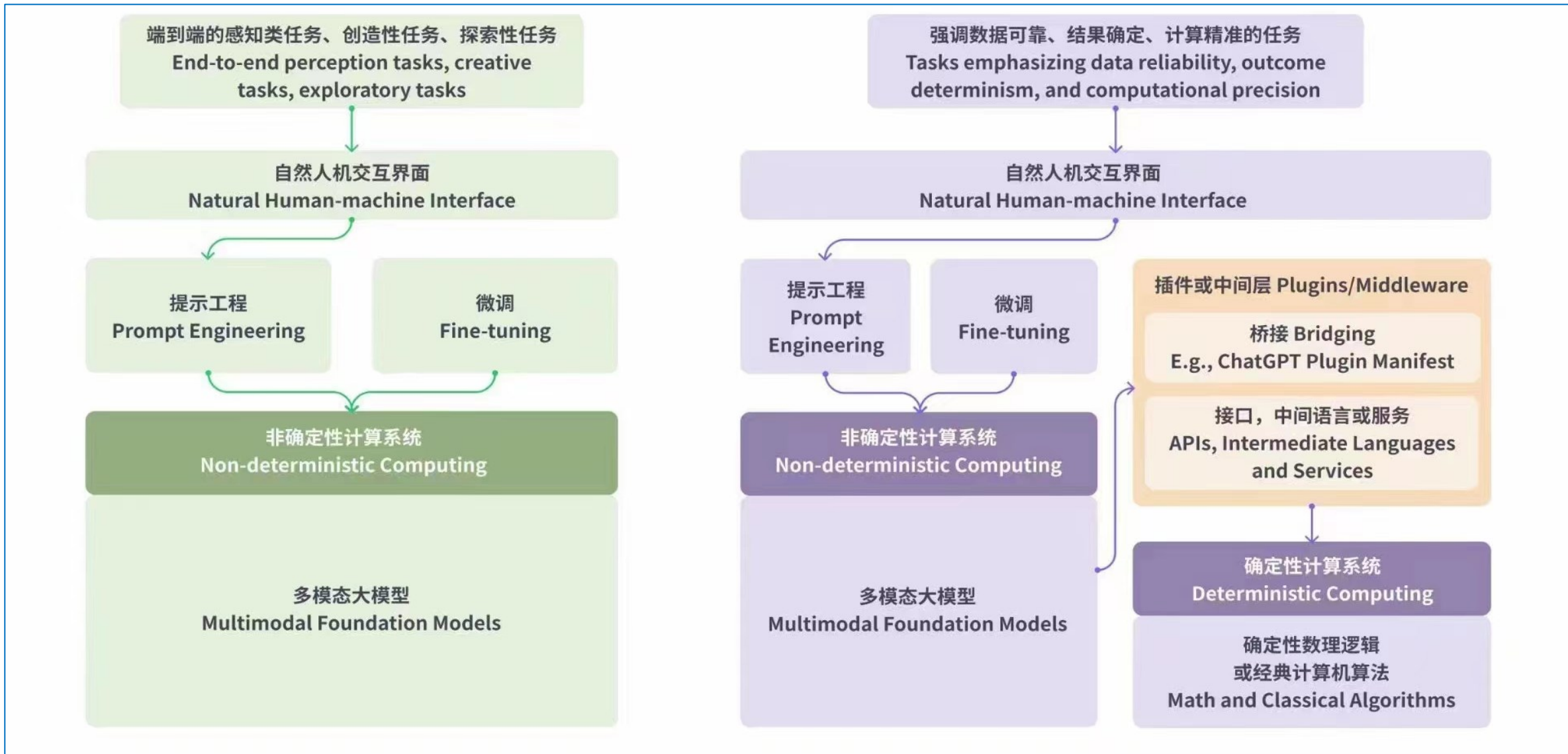
通用知识能力		特色知识能力
结构知识抽取	可信文本生成	公文写作助手
数学/常识推理	高效工具使用

PART 03

应用落地范式



▶ 大语言模型产业落地：基本范式



以大语言模型为基础的生态雏形

以GPT作为基座，融合其他插件



Expedia

Bring your trip plans to life—get there, stay there, find things to see and do.



FiscalNote

Provides and enables access to select market-leading, real-time data sets for legal, political, and regulatory data and information.



Instacart

Order from your favorite local grocery stores.



KAYAK

Search for flights, stays and rental cars. Get recommendations for all the places you can go within your budget.



Klarna Shopping

Search and compare prices from thousands of online shops.



Milo Family AI

Giving parents superpowers to turn the manic to magic, 20 minutes each day. Ask: Hey Milo, what's magic today?



OpenTable

Provides restaurant recommendations, with a direct link to book.



Shop

Search for millions of products from the world's greatest brands.



Speak

Learn how to say anything in another language with Speak, your AI-powered language tutor.



Wolfram

Access computation, math, curated knowledge & real-time data through Wolfram|Alpha and Wolfram Language.



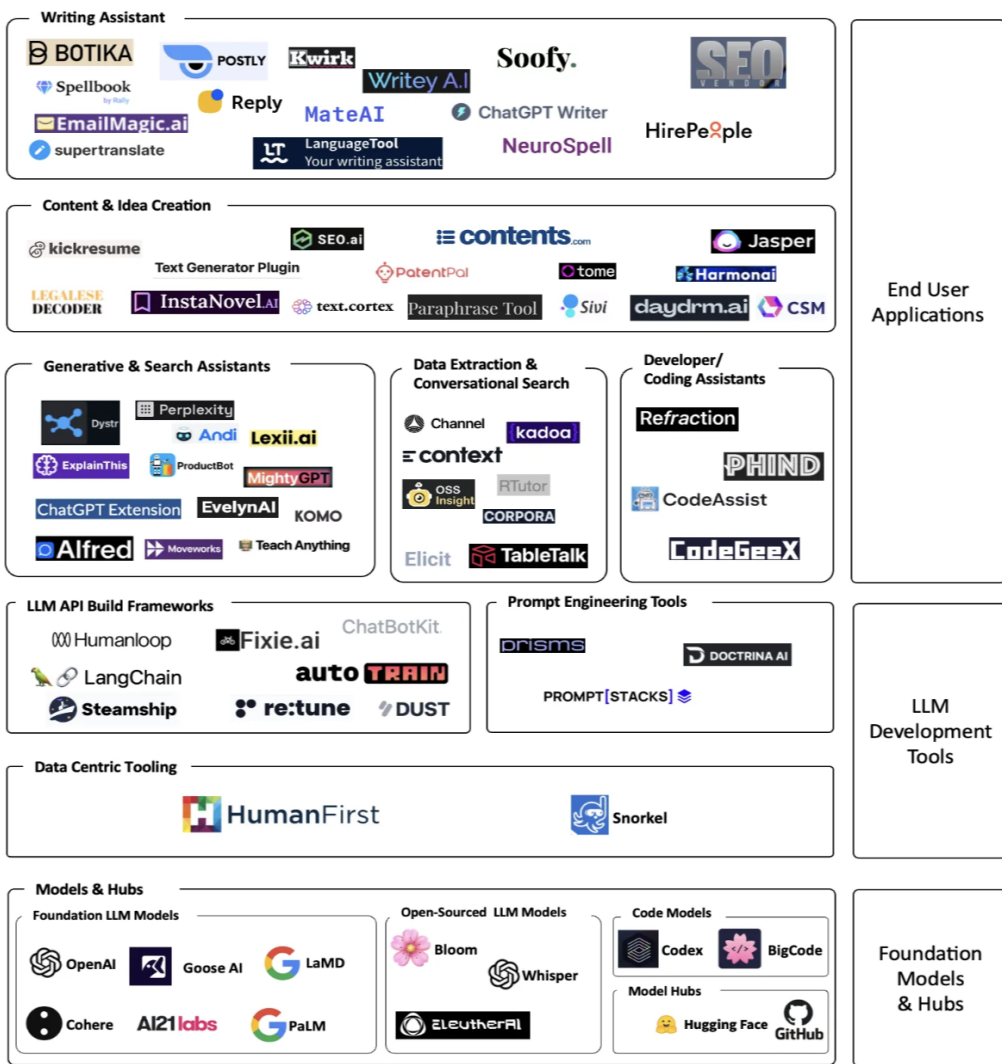
Zapier

Interact with over 5,000+ apps like Google Sheets, Trello, Gmail, HubSpot, Salesforce, and more.

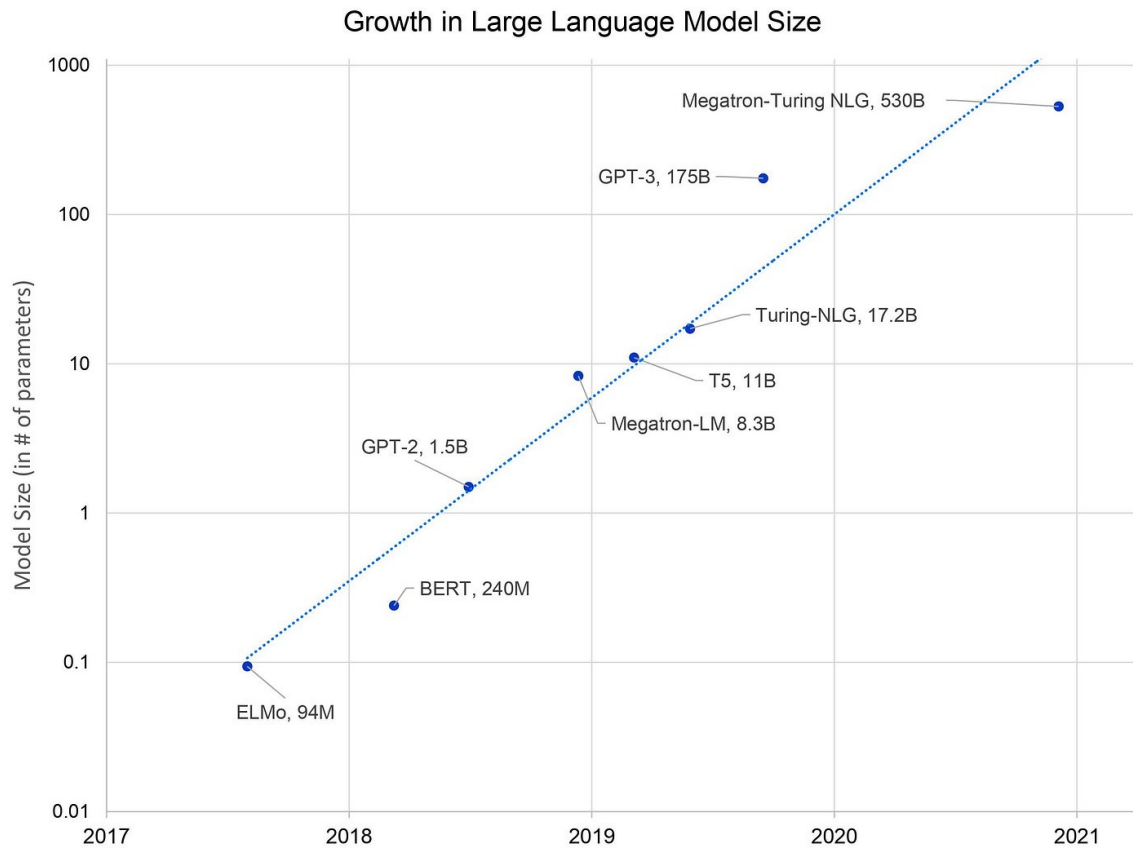
内生性升级，以海纳百川的姿态，旨在形成庞大的生态

大语言模型的技术栈与发展趋势

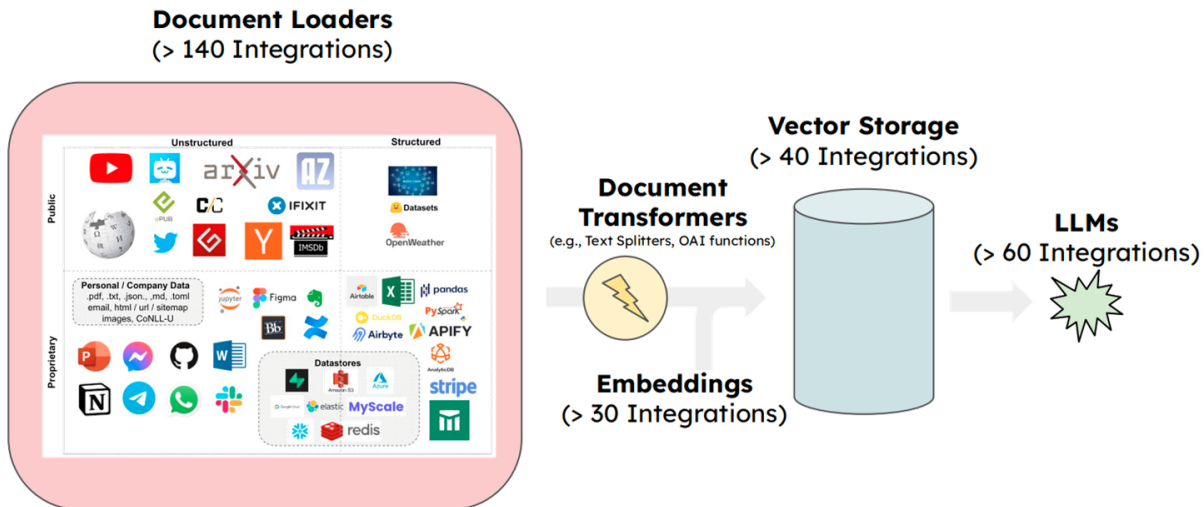
Foundation Large Language Model Stack



SaaS服务+生产力套件

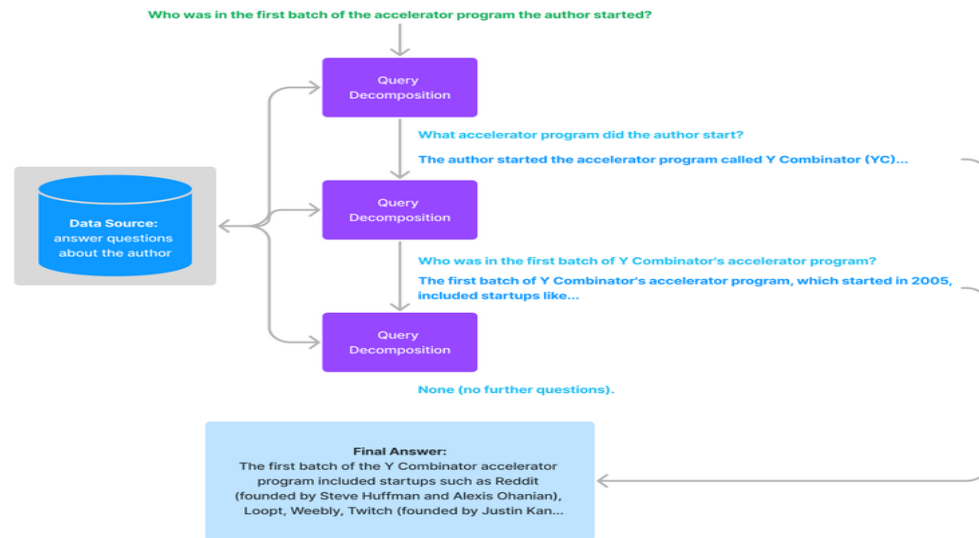


▶ 大语言模型开源工具典型代表



- LLM中间件
- 直接与 OpenAI 的 GPT-3 和 GPT-3.5 模型以及 Hugging Face 的开源替代品进行集成
- 允许用户围绕大型语言模型快速构建应用程序和管道

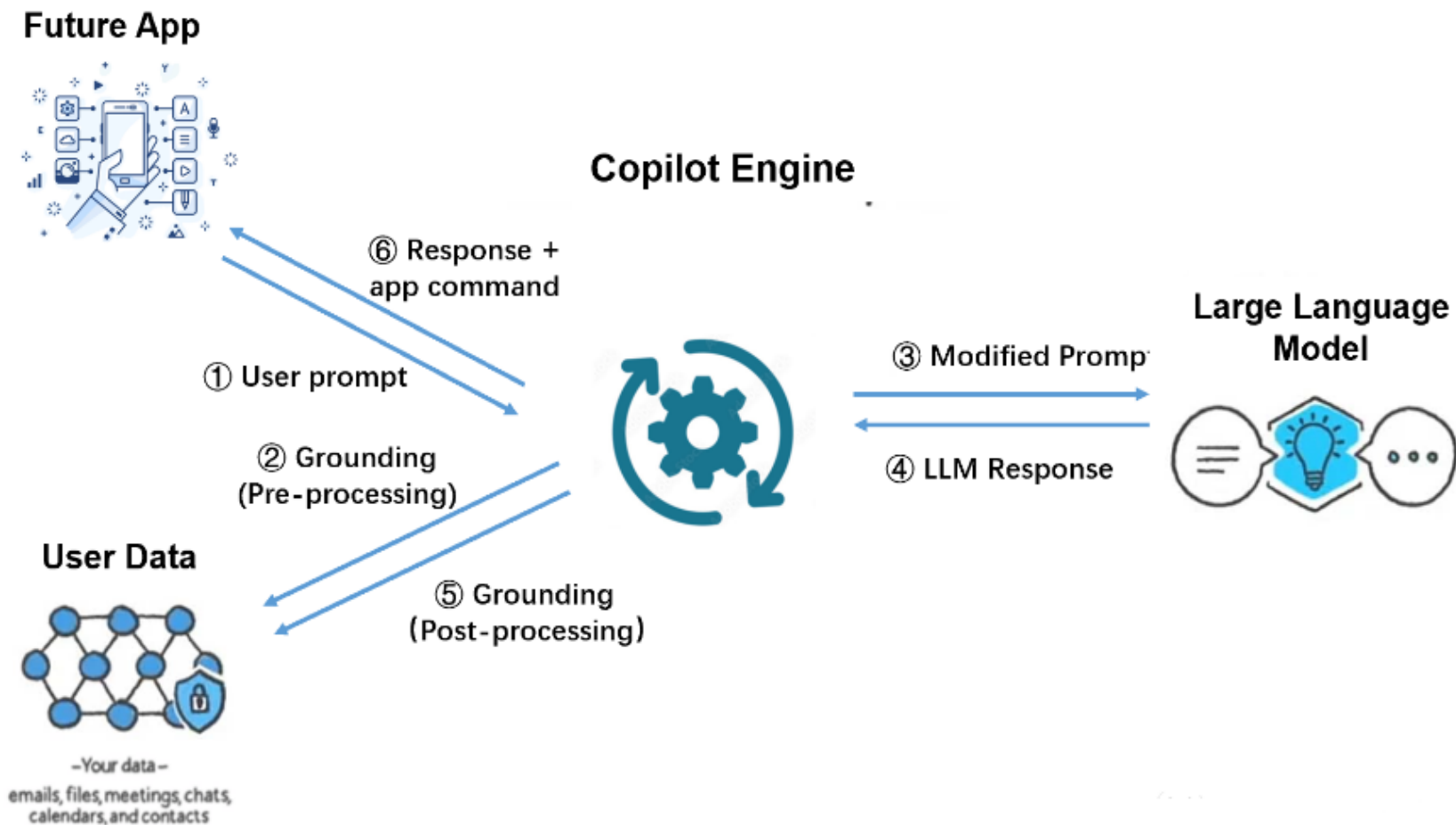
Welcome to LlamaIndex 🐪 (GPT Index)!



- 提供中心接口，将LLM与外部数据连接起来
- 提供了一组数据结构，为各种LLM任务索引大量数据，并消除了提示大小限制和数据摄入的问题
- 支持与LangChain等其他工具或库的集成

▶ 基于大语言模型的应用开发新范式

- ❖ Grounding (pre-processing) : 调用 Microsoft Graph 提取与任务相关的用户数据，作为 prompt 一部分。 (**Prompt 自动构建**)
- ❖ Modified prompt: 结合用户 prompt 和数据，优化 prompt 使产出结果更加稳定高质量，发送给 LLM。
- ❖ Grounding (post-processing) : 对 LLM 返回的结果做 post-processing，并从 User Data 中获取相关数据作为结果的补充。 (**结果处理(补充)**)
- ❖ Response + App command: 把结果变成调用前端应用 (Word、Excel、PPT等) 的命令。比如，Office 支持 VBA，输出的 command 可能是代码。





Generating output...

March 16, 2023, 1:00 PM

Here are some topics to prepare you for your upcoming meeting:

- **Project updates:** The new product line was launched in 3 countries. Project Gamma has been delayed by 2 weeks. ↗
- **Customer Updates:** The sales team finalized a new annual contract with Relecloud ↗ and the meeting with the Relecloud executives has been scheduled for this week.
- **Organization updates:** There were 3 new hires who started in the sales team. The marketing team celebrated Joanne's 25 year anniversary. ↗
- **Team updates:** Matthew returned from paternity leave. The team offsite is scheduled for this week and the agenda was shared with everyone.



Copy



Edit in Word



3 references ↕

Create content with Copilot [Close]

Describe what you'd like to create, and Copilot will generate a draft

[Decorative icons]

1 [Red-bordered box]

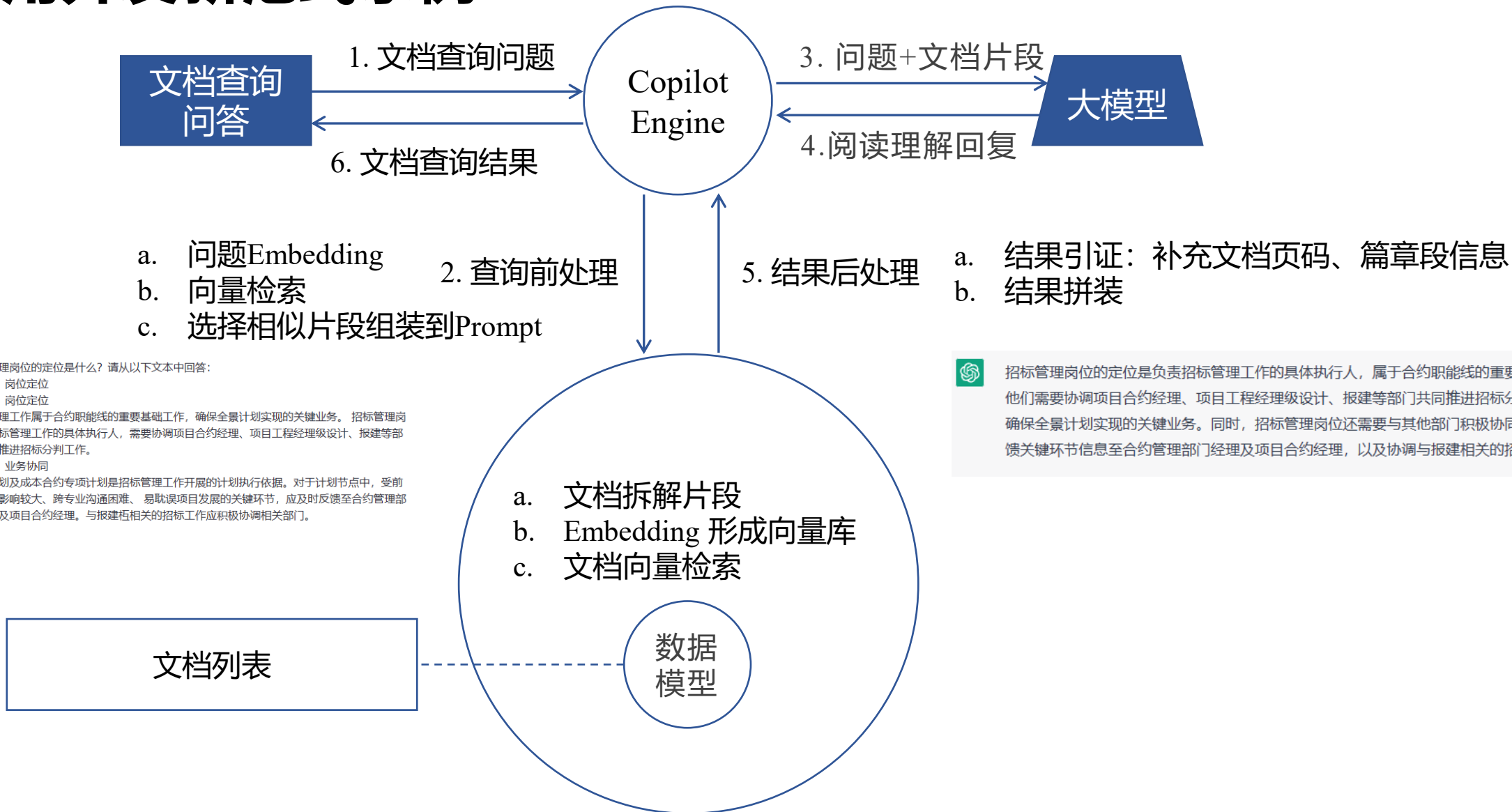
+

Click to add title

Click to add subtitle



应用开发新范式示例：ChatDocument



招标管理岗位的定位是什么? 请从以下文本中回答:

第二章 岗位定位

第四条 岗位定位

招标管理工作属于合约职能线的重要基础工作, 确保全景计划实现的关键业务。 招标管理岗位是招标管理工作的具体执行人, 需要协调项目合约经理、项目工程经理级设计、报建等部门共同推进招标分判工作。

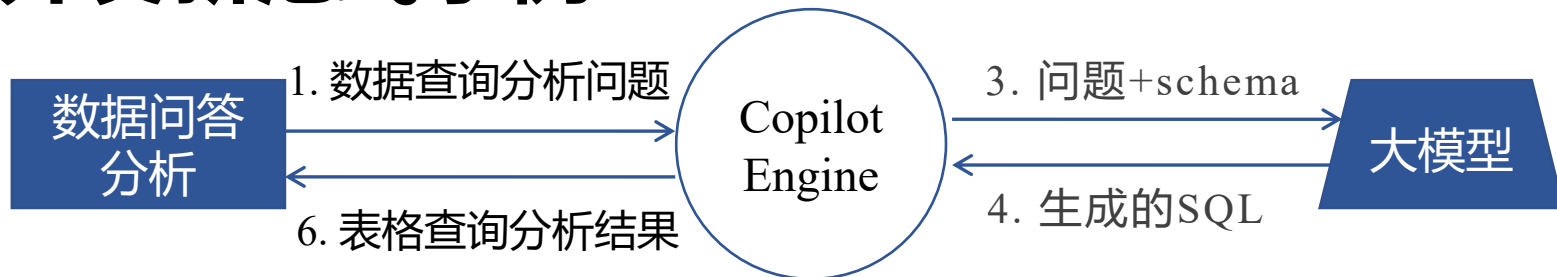
第五条 业务协同

全景计划及成本合约专项计划是招标管理工作开展的计划执行依据。对于计划节点中, 受前置任务影响较大、跨专业沟通困难、易耽误项目发展的关键环节, 应及时反馈至合约管理部门经理及项目合约经理。与报建相关的招标工作应积极协调相关部门。



招标管理岗位的定位是负责招标管理工作的具体执行人, 属于合约职能线的重要基础工作。他们需要协调项目合约经理、项目工程经理级设计、报建等部门共同推进招标分判工作, 以确保全景计划实现的关键业务。同时, 招标管理岗位还需要与其他部门积极协同, 如及时反馈关键环节信息至合约管理部门经理及项目合约经理, 以及协调与报建相关的招标工作。

应用开发新范式示例：ChatTable




- a. 定位查询表 (Locate query table)
- b. 表schema组装到Prompt (Assemble table schema into prompt)
- 2. 查询前处理 (Query preprocessing)
- 5. 结果后处理 (Result post-processing)
- a. 执行SQL查询 (Execute SQL query)

GU 数据库表公民信息表 citizen，表格的列为“唯一标识id”，“姓名name”和“身份证号ID_number”
写一个SQL查询：查找出所有“身份证号相同但姓名不同”的记录

a. 表格 schema 描述

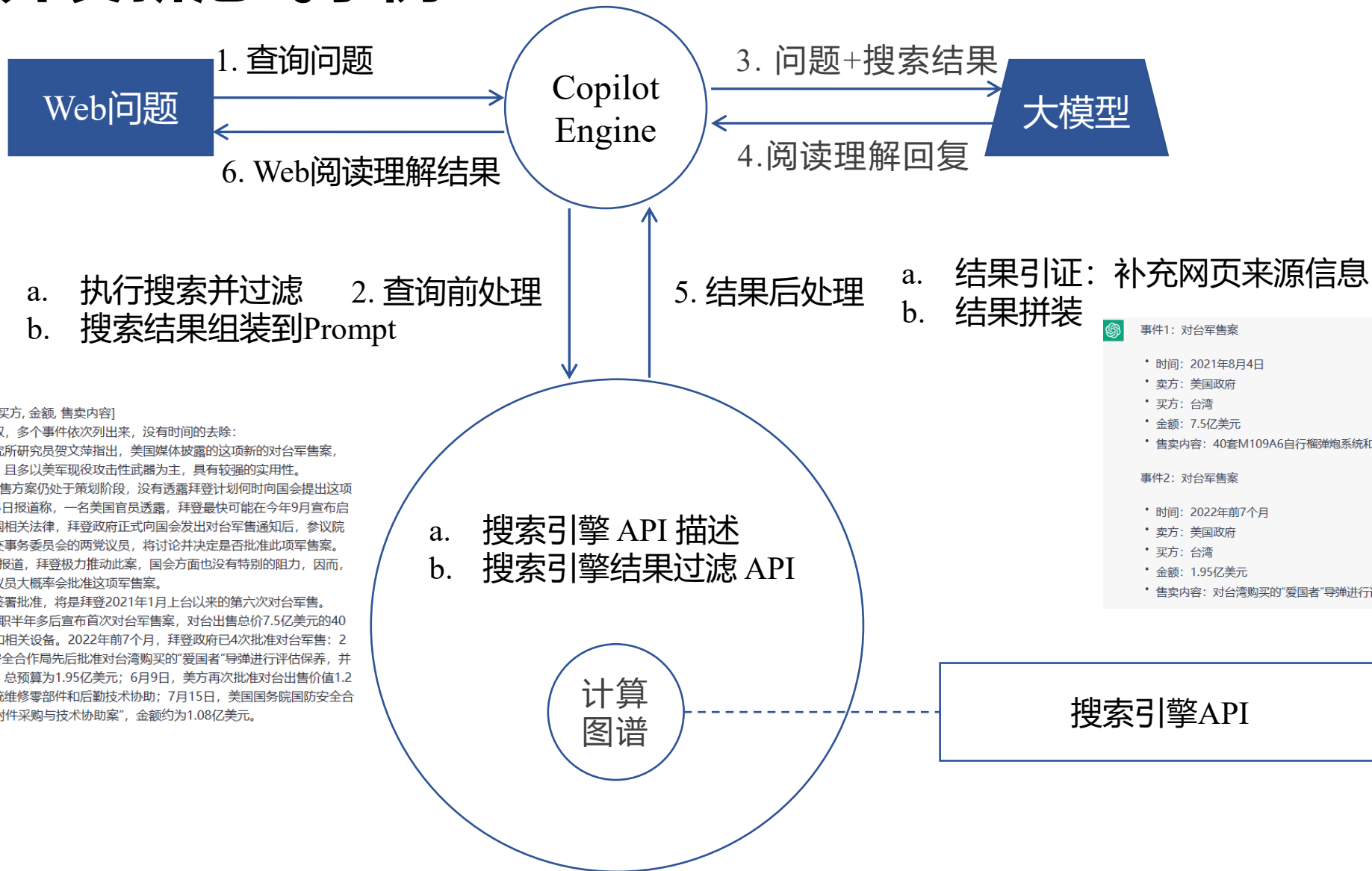
数据库/数据中台的表格

数据模型

 要查询所有身份证号相同但姓名不同的公民信息记录，可以使用子查询来选择具有多个不同名称的唯一身份证号码，然后在唯一身份证号码上将子查询与原始表进行联接：

```
sql
SELECT c.id, c.name, c.ID_number
FROM citizen c
INNER JOIN (
  SELECT ID_number
  FROM citizen
  GROUP BY ID_number
  HAVING COUNT(DISTINCT name) > 1
) c2 ON c.ID_number = c2.ID_number
ORDER BY c.ID_number, c.name
```

应用开发新范式示例：ChatWeb



- a. 执行搜索并过滤
- b. 搜索结果组装到Prompt

5. 结果后处理

- a. 结果引证：补充网页来源信息
- b. 结果拼装

事件1: 对台军售案

- 时间: 2021年8月4日
- 卖方: 美国政府
- 买方: 台湾
- 金额: 7.5亿美元
- 售卖内容: 40套M109A6自行榴弹炮系统和相关设备

事件2: 对台军售案

- 时间: 2022年前7个月
- 卖方: 美国政府
- 买方: 台湾
- 金额: 1.95亿美元
- 售卖内容: 对台湾购买的“爱国者”导弹进行评估保养, 并派遣人员提供直接技术支持

GU

任务类型: 事件抽取

事件类型: 军售

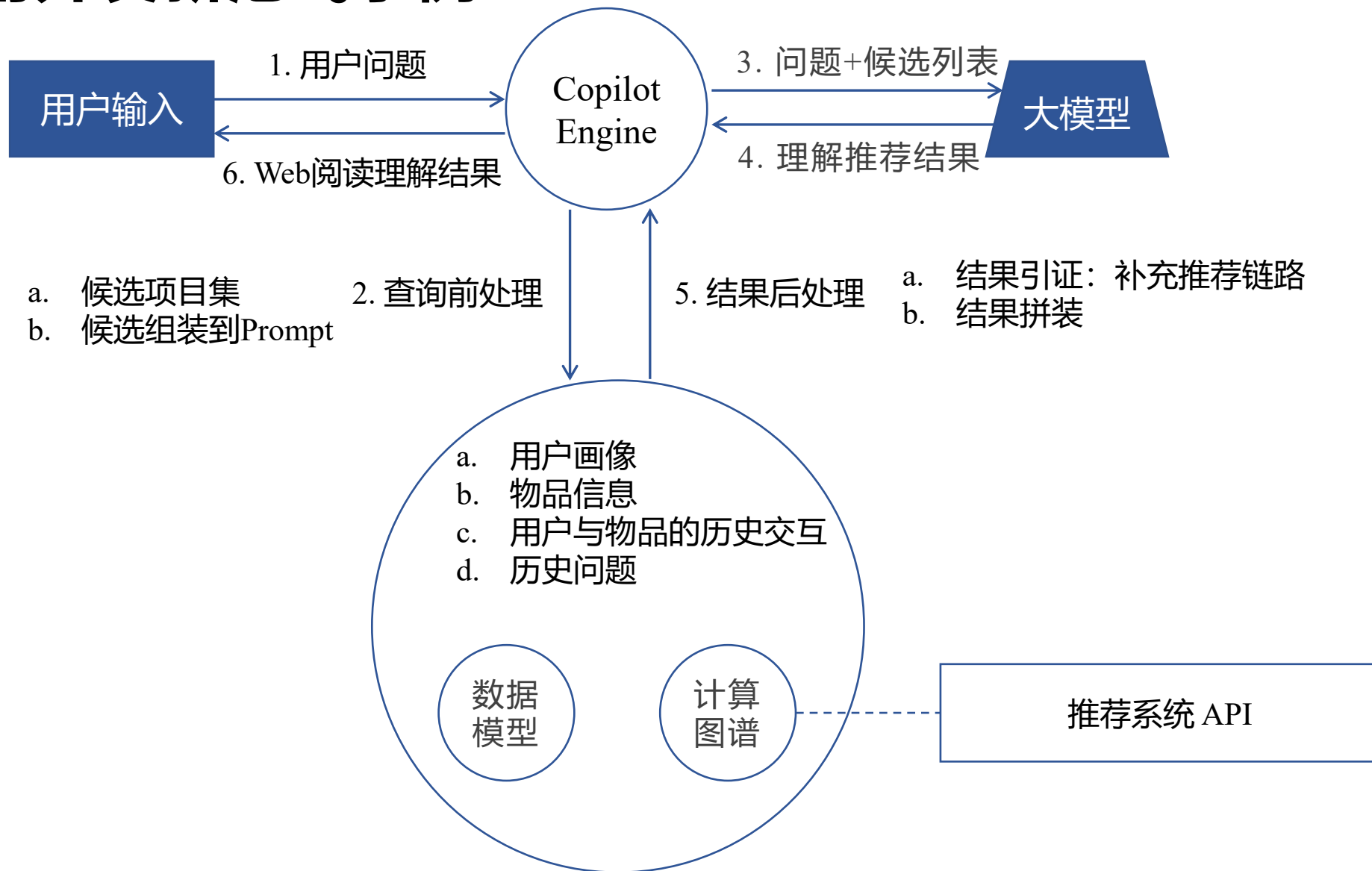
事件论元角色: [时间, 卖方, 买方, 金额, 售卖内容]

请从如下文本中进行事件抽取, 多个事件依次列出来, 没有时间的去除:

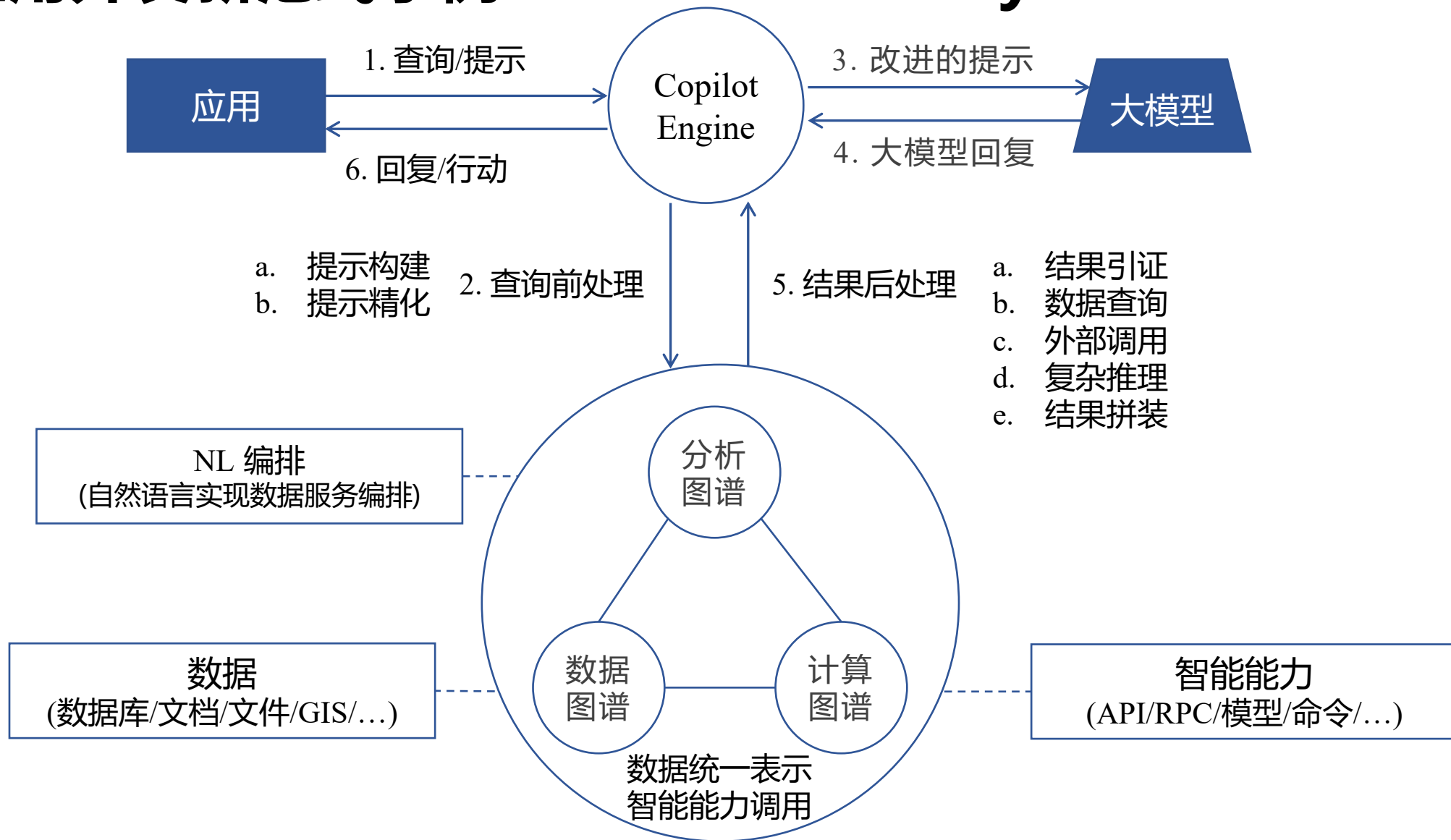
中国社会科学院西亚非洲研究所研究员贺文萍指出, 美国媒体披露的这项新的对台军售案, 价值总额、规模均远超以往, 且多以美军现役攻击性武器为主, 具有较强的实用性。尽管3名美方知情人士都称军售方案仍处于策划阶段, 没有透露拜登计划何时向国会提出这项军售案, 但英国路透社8月25日报道称, 一名美国官员透露, 拜登最快可能在今年9月宣布启动此项新军售项目。按照美国相关法律, 拜登政府正式向国会发出对台军售通知后, 参议院外交关系委员会与众议院外交事务委员会的两党议员, 将讨论并决定是否批准此项军售案。据美国政治新闻网Politico的报道, 拜登极力推动此案, 国会方面也没有特别的阻力, 因而, 美国国会秋季复会后, 两党议员大概率会批准这项军售案。

如果此项军售案获国会议员签署批准, 将是拜登2021年1月上台以来的第六次对台军售。2021年8月4日, 拜登政府就职半年多后宣布首次对台军售案, 对台出售总价7.5亿美元的40套M109A6自行榴弹炮系统和相关设备。2022年前7个月, 拜登政府已4次批准对台军售: 2月7日和4月5日, 美国国防安全合作局先后批准对台湾购买的“爱国者”导弹进行评估保养, 并派遣人员提供直接技术支持, 总预算为1.95亿美元; 6月9日, 美方再次批准对台出售价值1.2亿美元的军用船舶、舰载系统维修零部件和后勤技术协助; 7月15日, 美国国务院国防安全合作局 (DSCA) 批准售台“零附件采购与技术协助案”, 金额约为1.08亿美元。

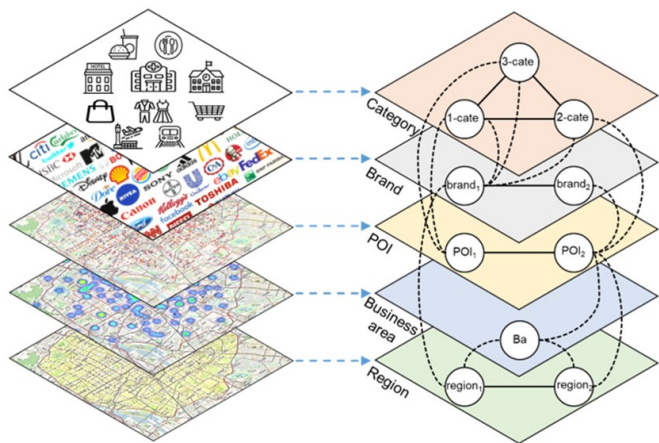
应用开发新范式示例：ChatRec



应用开发新范式示例：Chat and Analyze All Data



领域落地的思考：智慧城市



01 数据异质性：城市领域涉及多种类型的数据，如地理数据、传感器数据、图像、文本等。处理和融合这些异质数据是一个关键挑战。

02 时空数据处理：城市领域涉及大量数据的收集、存储和处理。有效管理和分析这些大数据以提高计算效率和准确性是一个核心挑战。

03 时空特性：城市领域问题需要考虑空间和时间维度。理解和建模这些时空关系，如空间自相关性和时间序列分析，是一个重要挑战。

04 模型泛化：城市领域涉及多种应用场景，如交通、能源、环境等，需要开发具有良好泛化能力的模型。

05 多尺度分析：城市领域通常需要在多个尺度上进行分析。处理和理解这些不同尺度之间的关系是一个关键挑战。

06 可解释性：城市计算模型需要具备一定的可解释性，以便为决策者提供有用的洞察。提高模型的可解释性和可信度是一个关键挑战。

▶ 领域落地的思考：智慧城市

大语言模型（发展相对成熟）

大城市模型（发展初期）

相同点

数据驱动

大语言模型和大城市模型，都依赖于大量数据进行训练。数据量和数据质量对模型性能有很大影响

迁移学习

大语言模型和大城市模型均可以利用迁移学习来提高模型泛化能力和在特定任务上的性能

不同点

主要关注非结构化文本数据

主要处理序列文本数据

主要依赖单一模型结构

主要关注自然语言处理任务

数据类型

需处理多源、多模态数据（如地理数据、传感器数据、轨迹、图像、视频等）

时空特性

需要处理时间和空间维度的关系。通常还需要结合地理知识，如空间自相关性

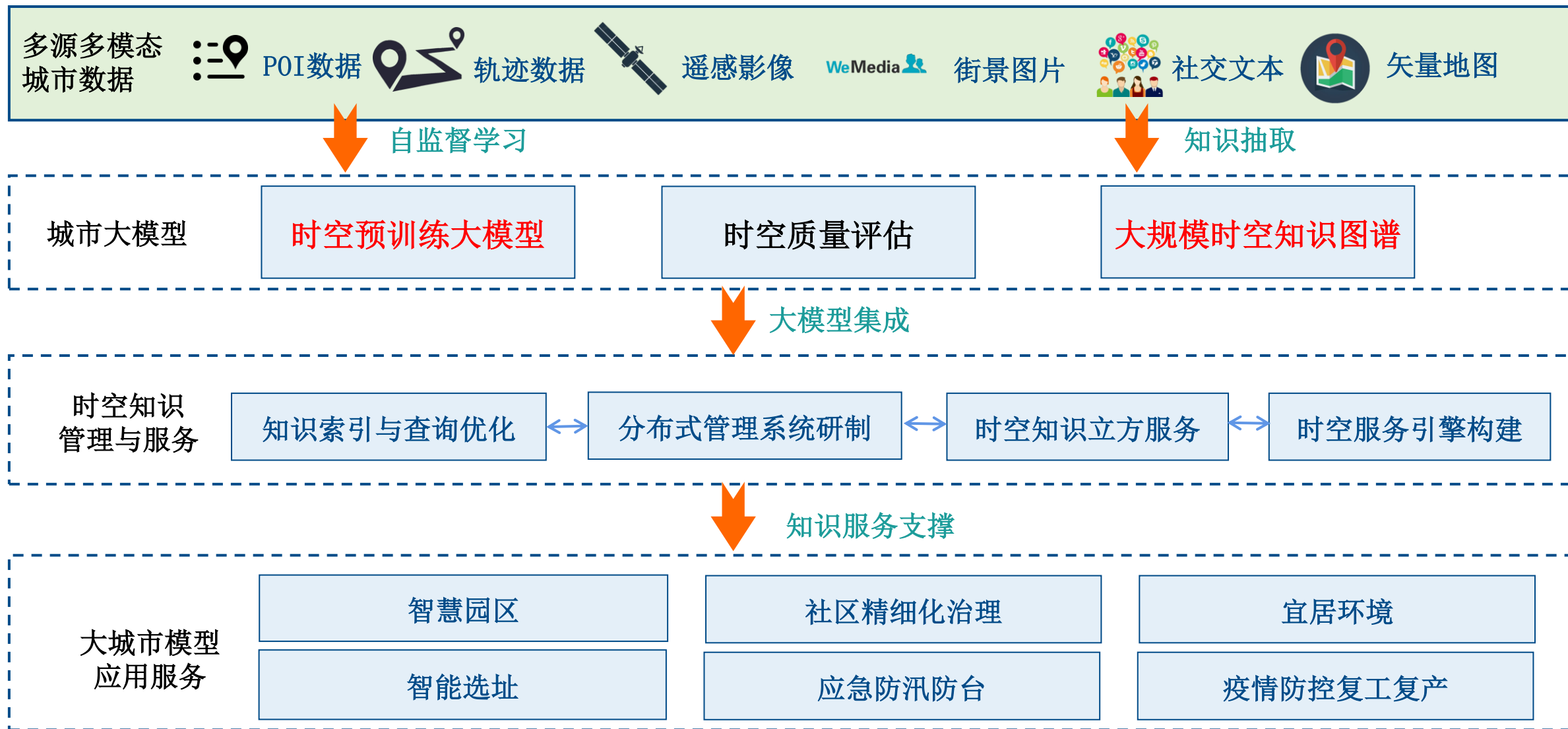
模型结构

需采用更复杂的多模态深度学习结构和网络拓扑结构（如知识图谱、图神经网络等）

应用场景

城市规划、交通、能源、环境等

▶ 领域落地的思考：智慧城市



▶ 领域落地的思考：智慧城市

时空预训练大模型 + 时空知识图谱

大城市模型层

端到端感知，创造性，探索性任务

数据可靠、结果确定、计算精准任务

实体、关系、事件挖掘

辅助构建、更新图谱

知识融合、增强

约束模型输出

时空预训练
大模型

大规模时空
知识图谱

时空数据生成

通用城市表示

城市语义理解

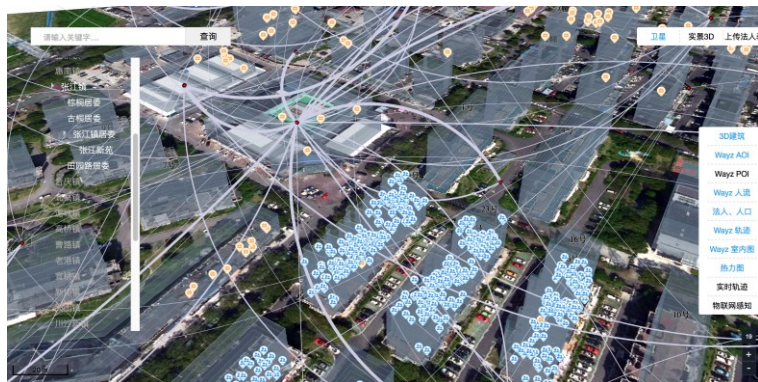
语义搜索

跨域推理

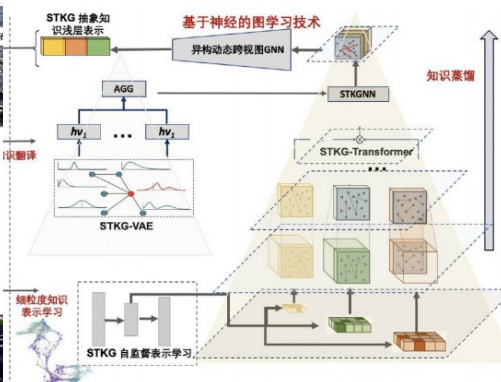
多跳问答

大城市模型协同

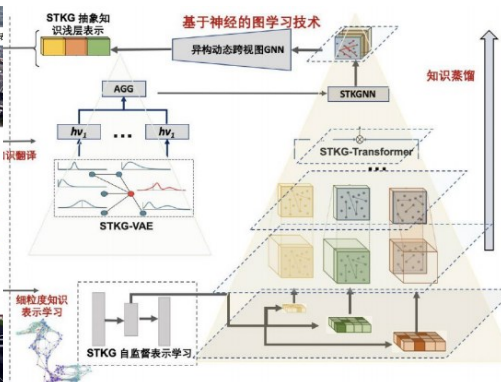
领域落地的思考：智慧城市



园区时空本体



园区时空图谱



园区时空模型

园区时空AI微服务场景



园区体征评价



异常预警



资源配置优化

园区时空AI 微服务接口

园区时空AI定制算法和训练模型

时空知识资产 (动静态知识, 图谱, 预训练模型)

时空AI引擎 (时空大数据, 图谱, AI, 超融合)

业务场景定制开发

智能发现, 智能预测微服务

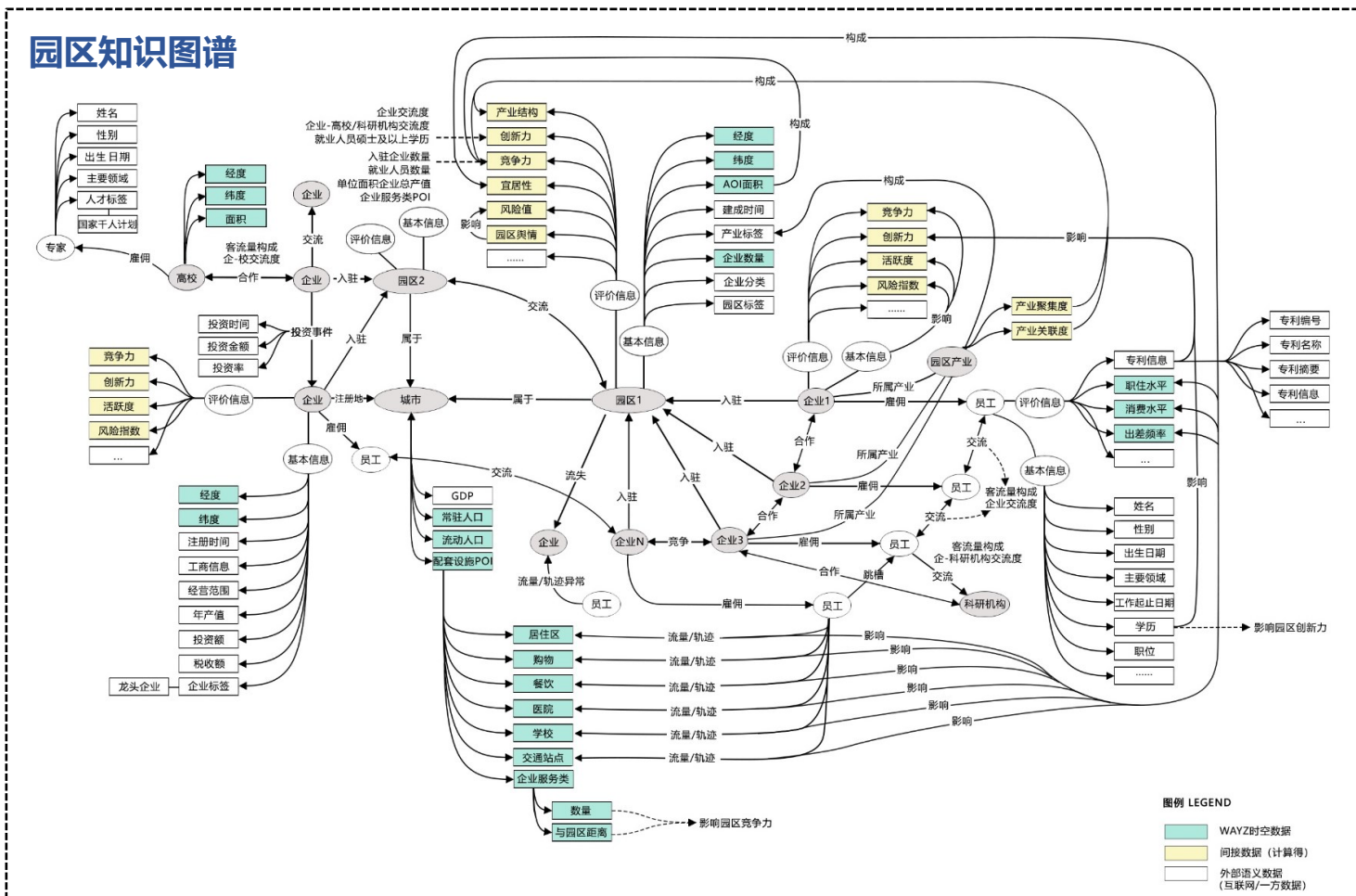
园区业务知识时空融合

开箱即用-园区洞察

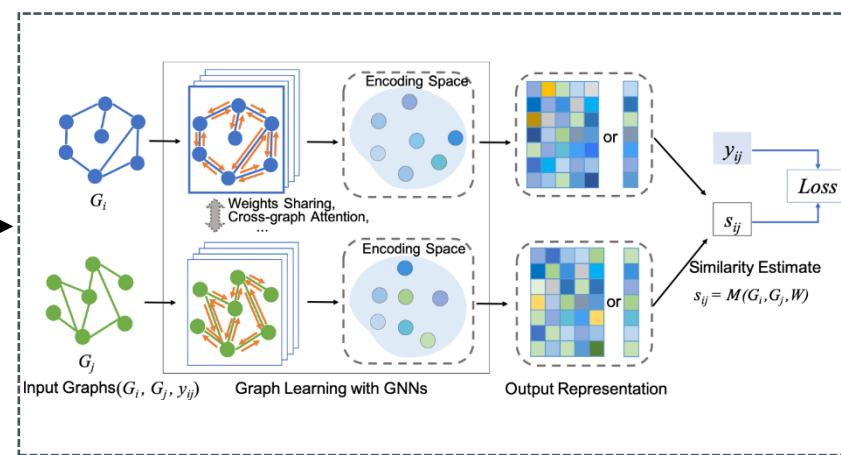


领域落地的思考：智慧城市

- 利用图谱实现园区规划定位、评估、归因、优化全流程。对标榜首产业园区，进行图谱关联可解释分析，找出自身园区的潜在优化点



与头部园区的各项指标的相似度计算



与专业园区对比:

差异点即为
潜在优化点

做专

与综合园区对比:

优化已有园区
产业标签组合

做强

园区知识图谱大屏

- 横向比较张江园区与国内外重点园区
- 触屏交互

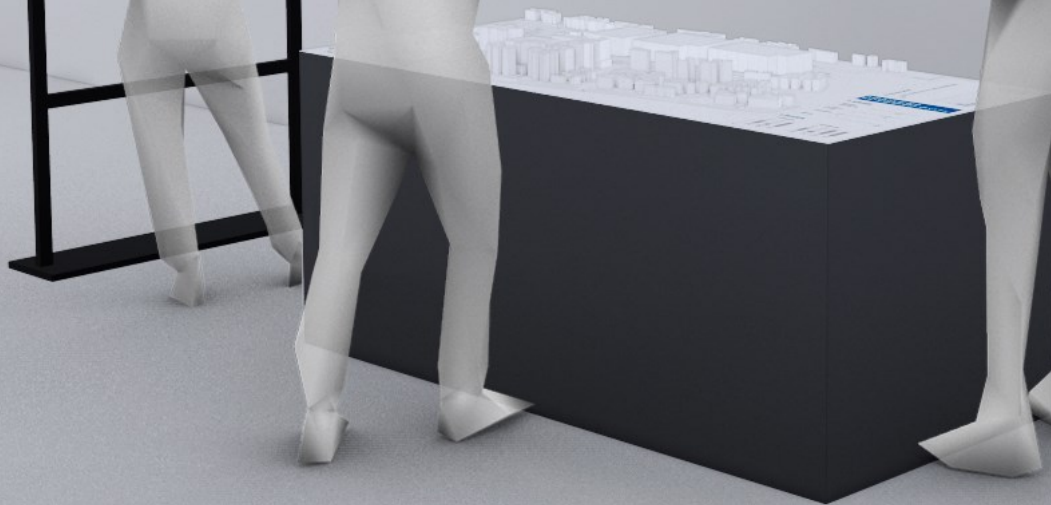


张江园区指标大屏

- 围绕“产、城、创”，为张江园区定制的指标大屏
- 屏幕或投影

园区设计操作台

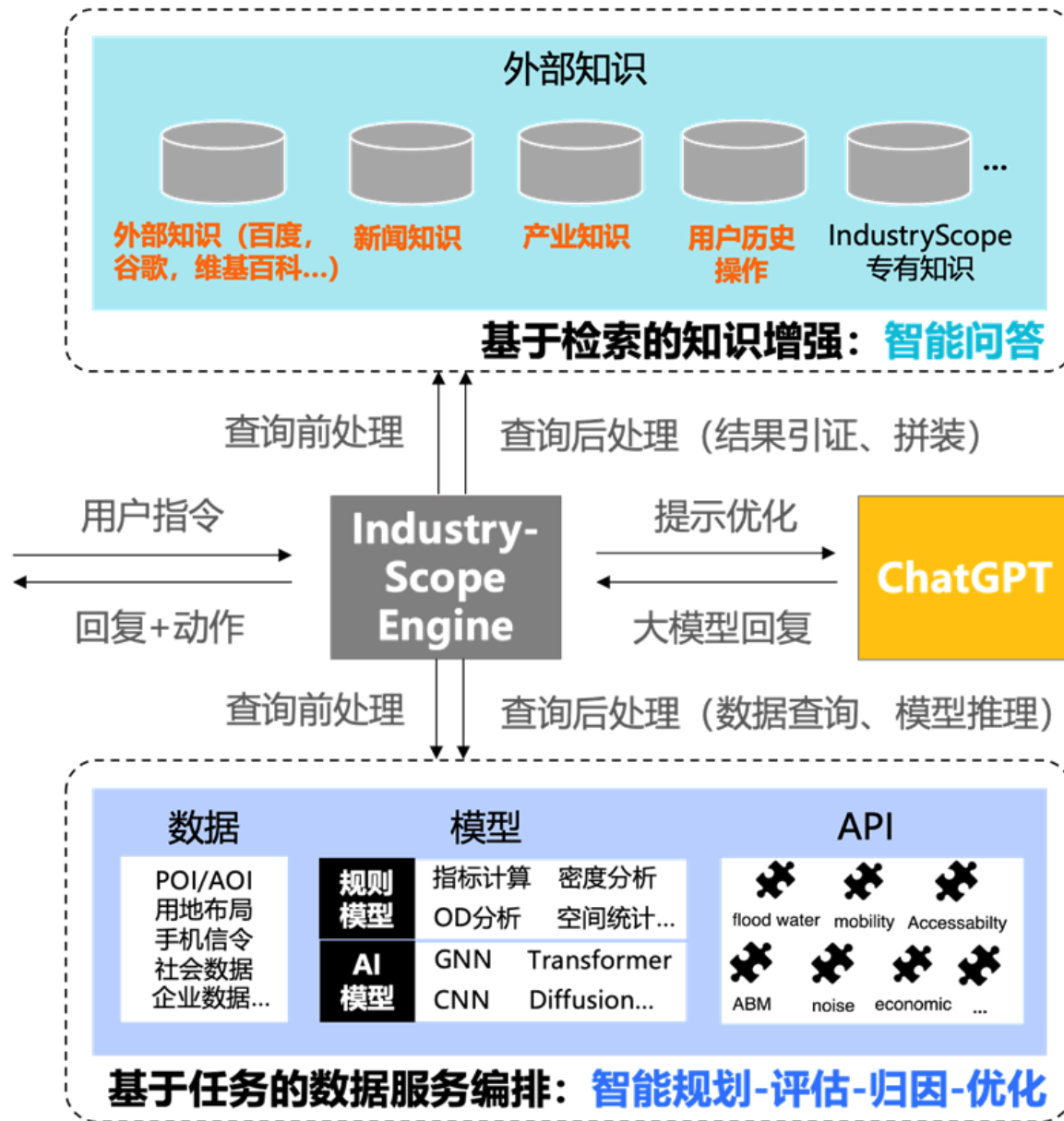
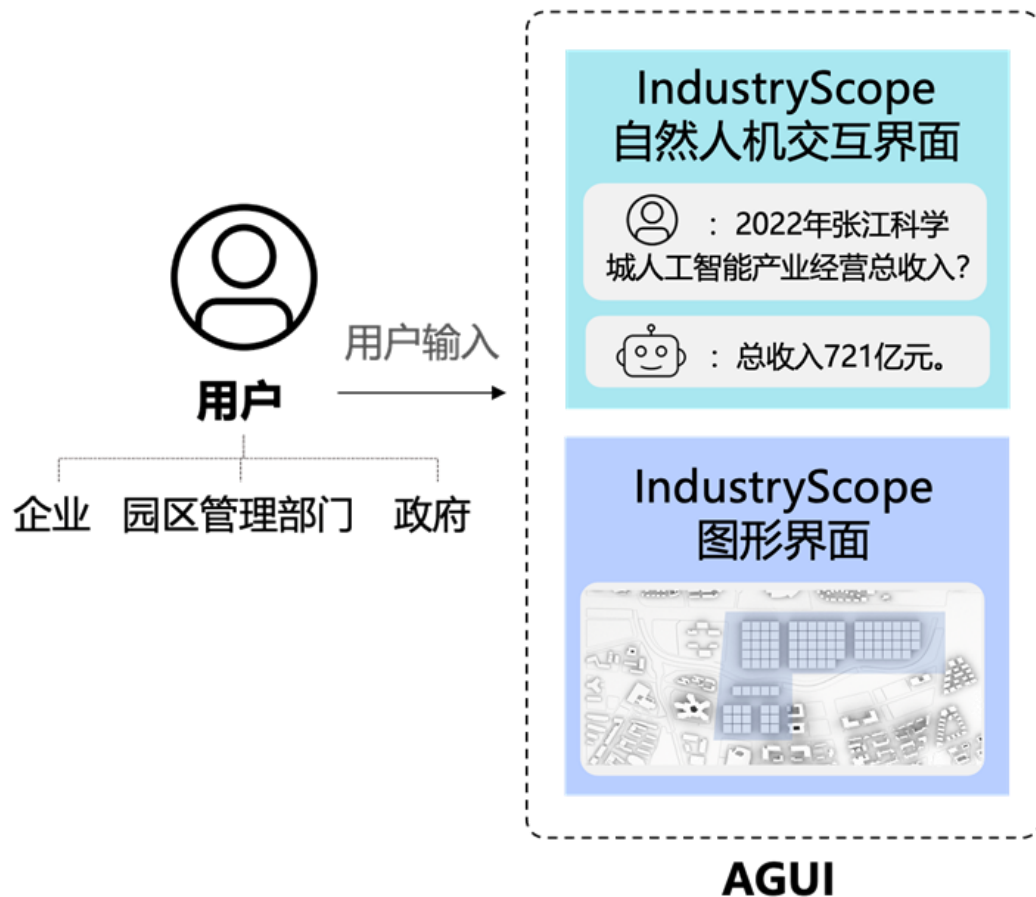
- 基于园区实体模型，采用可触交互技术的城市设计界面
- 以触屏实现的高精度设计参数调整，包括政策力度与部分城市功能属性



▶ 领域落地的思考：智慧城市

IndustryScopeGPT:

整合ChatGPT的Industryscope架构



▶ 领域落地的思考：汽车智能座舱

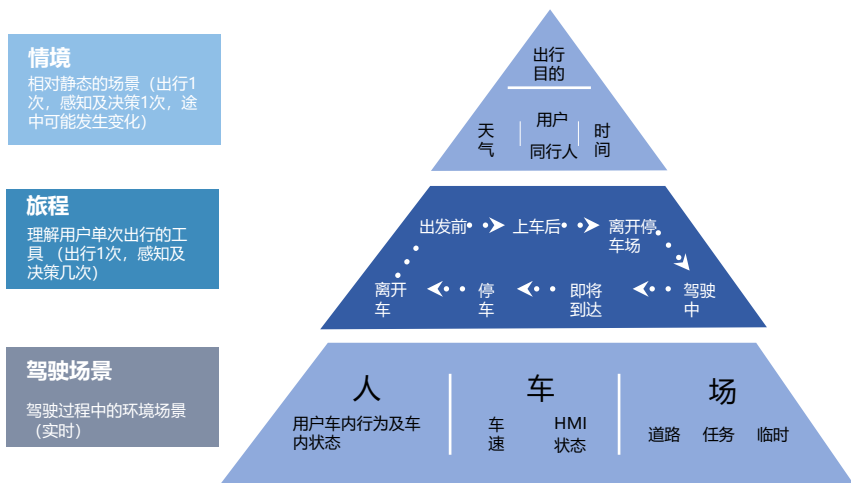
“智能”座舱体验提升亟需场景化、千人千面的解决方案



▶ 领域落地的思考：汽车智能座舱

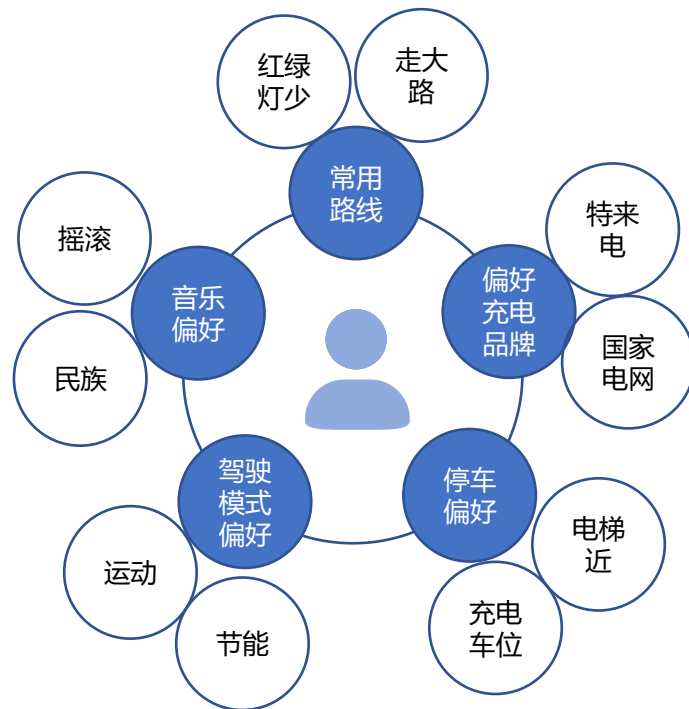
用车场景碎片化

汽车场景要素繁多，各种要素的组合属于乘积关系，造成场景无法穷尽



用户偏好个性化

同样的场景，不同用户的偏好不同，无法用同一个标准进行用户需求的预判



用户需求动态化

受出行目的、交通状况、同行者等场景要素的影响，用户需求随时可能发生改变



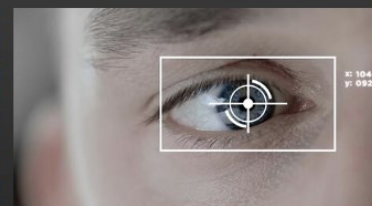
例子：

车内抽烟 -> 开窗
车内抽烟，下雨了 -> 打开外循环
车内抽烟，下小雨，副驾有人 -> 开窗

场景因素不同，用户的决策结果不一样

▶ 领域落地的思考：汽车智能座舱

IOT、大数据、AI等技术为场景化、千人千面的“智能”座舱体验提供了新的可能性



体验+

满足人们在高度动态化、碎片化场景下的个性需求

新体验
体验优化

体验优化

产品定义

智能推送

归因

优化

新方法
体验度量

场景建模

体验度量

知识图谱

情感计算

深度学习

新技术
智能感知

人

车

场

行为

表情

功能

状态

环境

状态

埋点

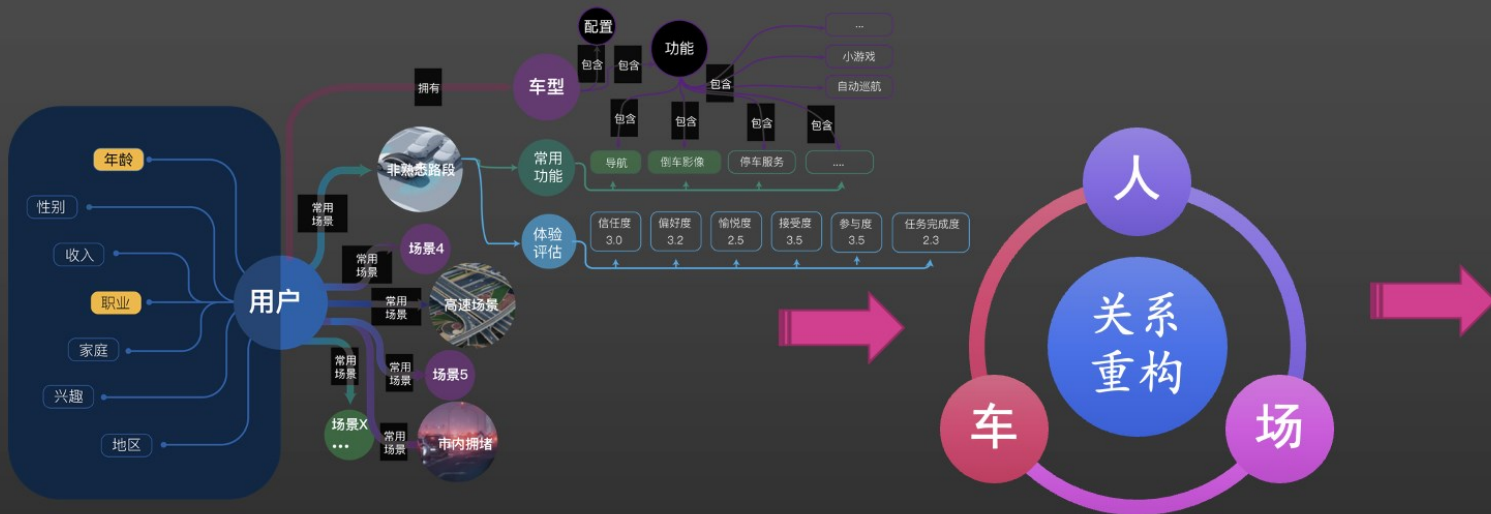
传感

IOT

▶ 领域落地的思考：汽车智能座舱

场景知识图谱辅助优化“智能”座舱体验

场景知识图谱



交互模态 最优路线
POI推荐 信息内容
语音对话 主动交互
配置策略 营销定位

利用知识图谱技术深度融合来自车厂数据和用户的体验数据

用户数据

- 基本profile
- 偏好数据
- 行为数据

场景知识图谱数据类型

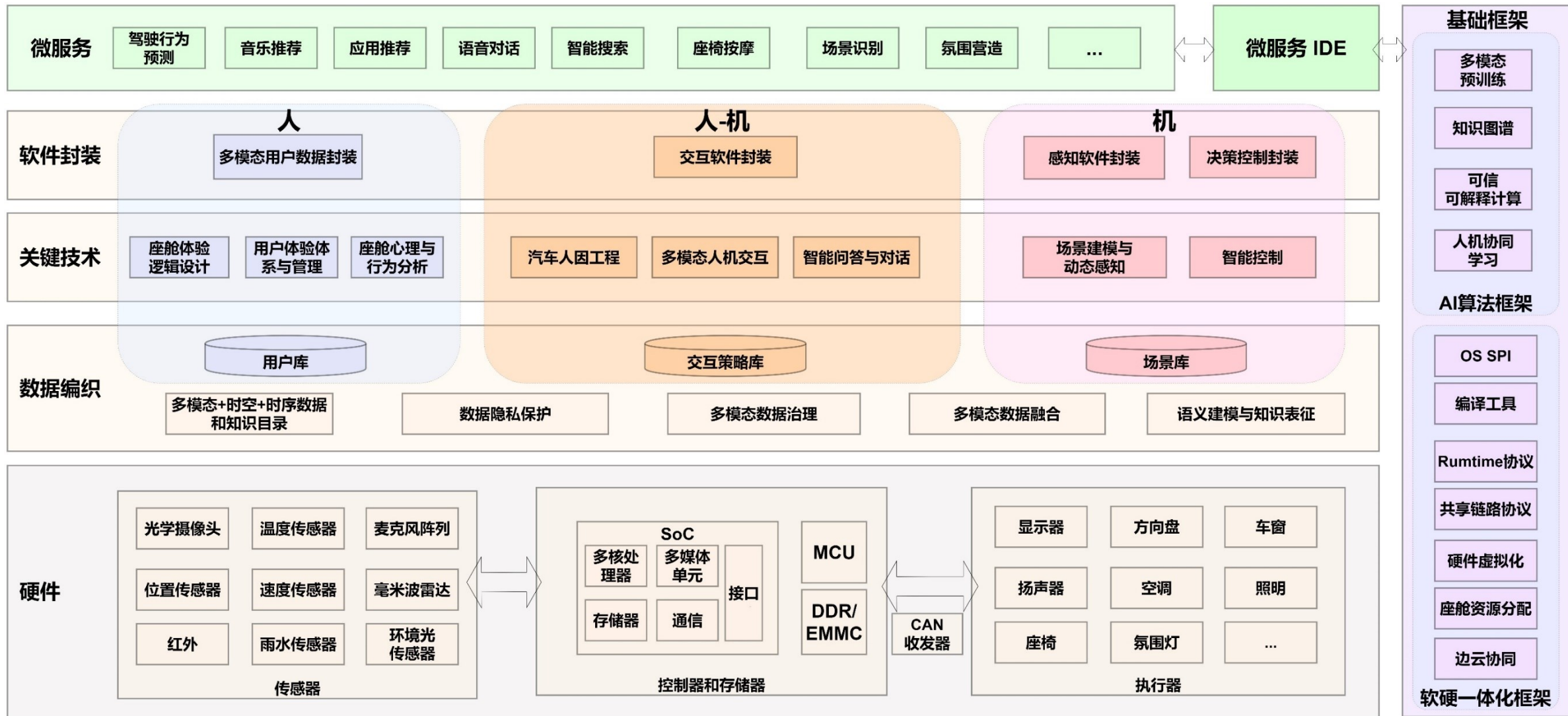
车辆数据

- 位置
- 车内状况
- ...

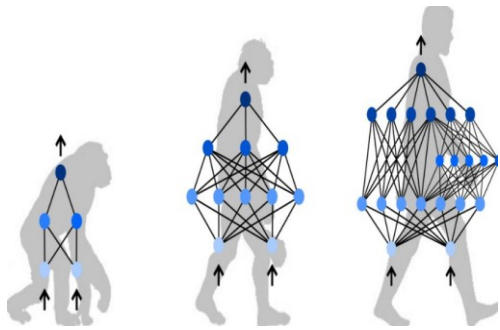
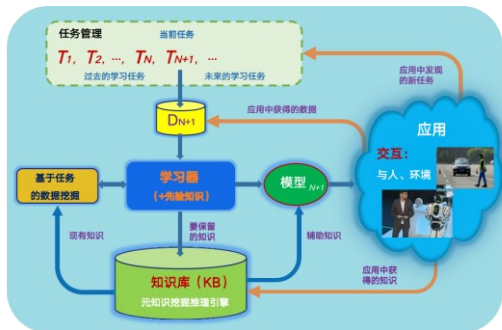
第三方数据

- 天气
- 地图
- ...

▶ 领域落地的思考：汽车智能座舱



▶ 领域落地的思考：汽车智能座舱



- 1) 智能座舱的“小数据，小场景” VS 互联网的“大数据，大场景” --- 人工定义无法覆盖所有“小场景”
- 2) AI算法需要通过多模态人机交互和场景感进行持续学习---开放世界一直在动态变化



01 用户个人知识图谱

- **基础：** 用户标签挖掘、交互本体建模
- **关键技术：** 语义自动化RPA, 分布式端交互
- **目标：** 通用有效的用户行为、交互策略库

02 环境知识图谱

- **基础：** 基于开放世界持续学习的环境上下文感知
- **关键技术：** 模型部署后的持续学习、小样本场景挖掘
- **目标：** 任务驱动式的持续学习，不断积累更新场景库

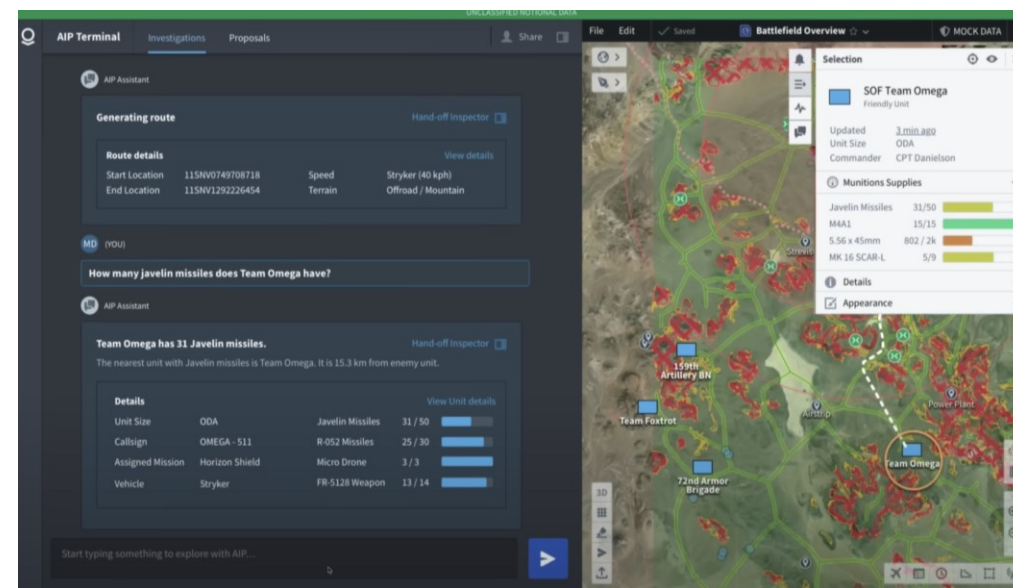
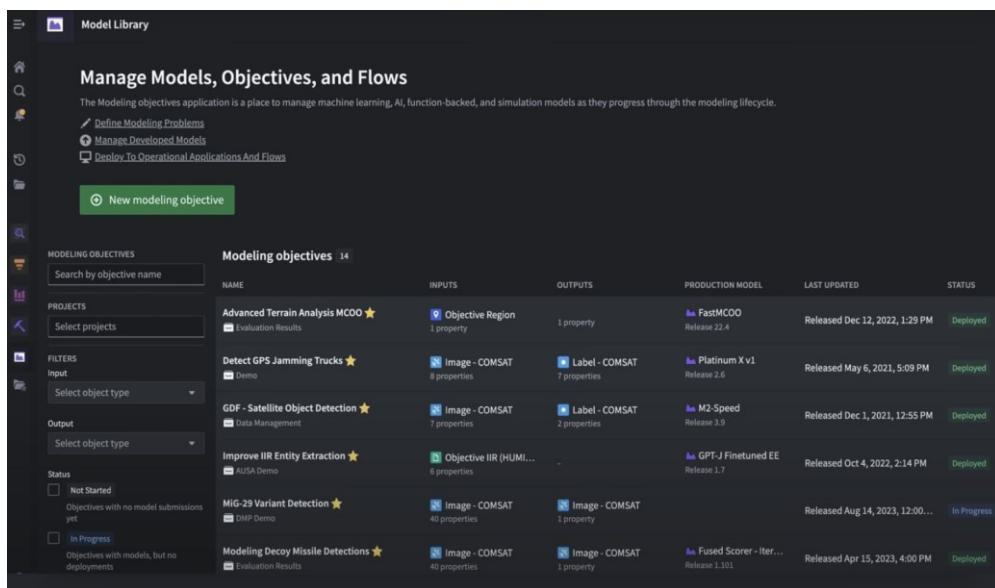
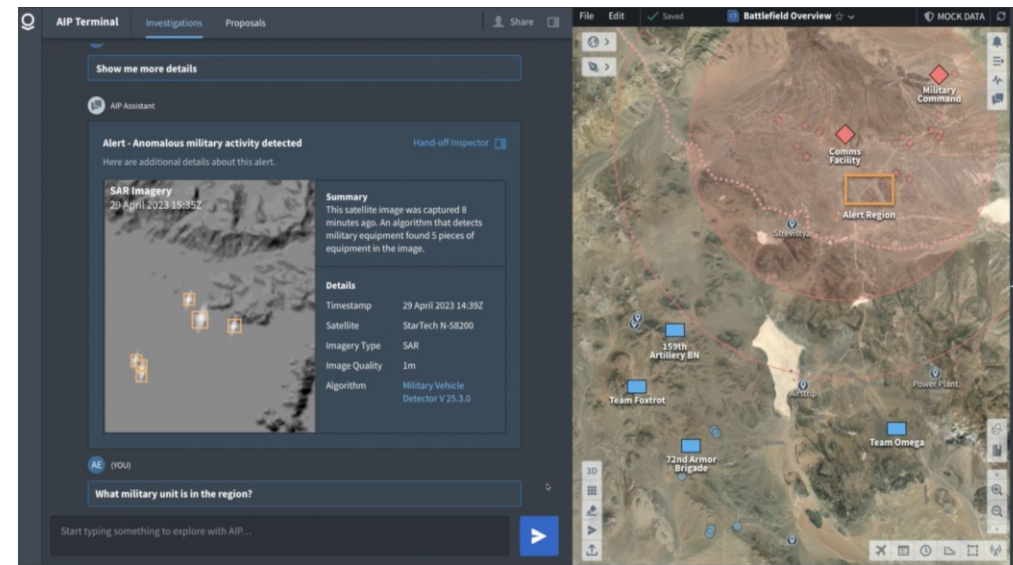
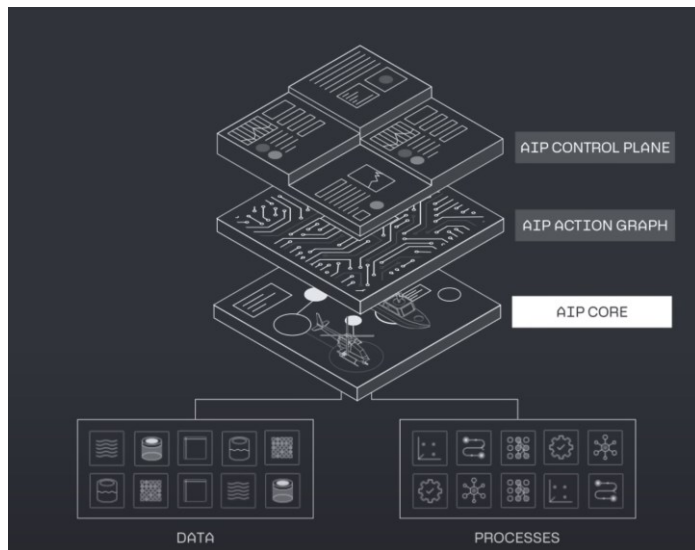


人机物融合的一次尝试 – Palantir AI Platform

AI-Powered Operations, For Every Decision

Activate LLMs and other AI on your private network, safely and securely.

<https://www.palantir.com/platforms/aip/>



AI驱动软件研发全面进入数字化时代

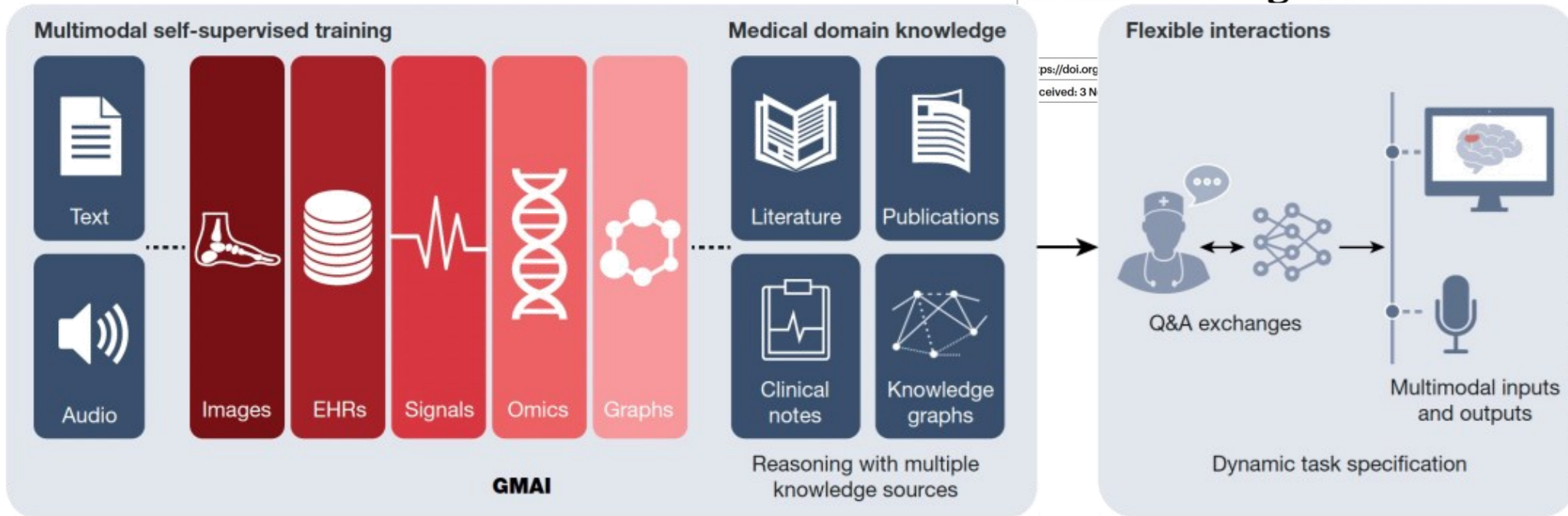
iDDD AI+ 软件研发数字峰会
AI+ software Development Digital summit

领域落地的思考：定制化的领域大模型更关注知识表示与推理

Perspective

Foundation models for generalist medical artificial intelligence

a



b

Applications



Chatbots for patients



Interactive note-taking



Augmented procedures

...



Grounded radiology reports



Text-to-protein generation



Bedside decision support

Regulations: Application approval; validation; audits; community-based challenges; analyses of biases, fairness and diversity

▶ 向领域专用、轻量化小模型发展

领域大模型

■ 金融领域大模型

 BloombergGPT

500亿参数 兼具金融领域和通用领域

■ 安全领域大模型

 Microsoft Security Copilot

微软发布 基于GPT-4与微软安全模型

 360 智脑

360发布 基于领域大模型与360搜索能力

■ 医疗领域大模型

 medGPT

医联发布 基于GPT-4与微软安全模型

■ 自动驾驶大模型

 DriveGPT

毫末智行发布 2000亿参数
4000万公里驾驶数据 5万RLHF数据

轻量化小模型

■ 斯坦福:



Llama-7B模型



Alpaca-7B模型

■ 智谱 AI :



ChatGLM-6B模型

■ 链家Tech:



BELLE-7B 中文对话大模型

■ HuggingFace:



Bloom-7B 中文对话大模型

在通用认知AI基础上，迈向专业人工智能

PART 04

未来展望



▶ 信息技术应用和创新生态新范式

第一阶段



上手使用

+ Prompt工程师

第二阶段



行业应用

+ 专家经验

第三阶段



工具打造

+ 行业积累

第四阶段

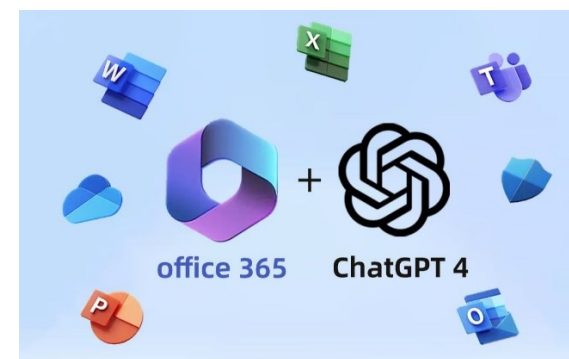
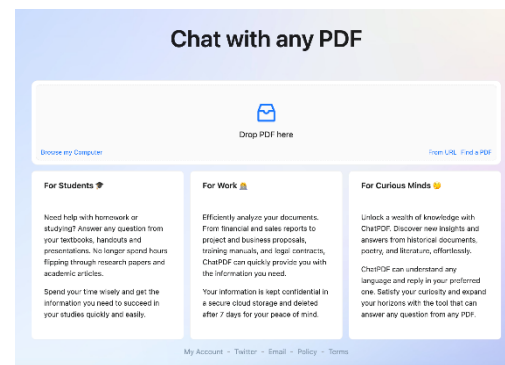
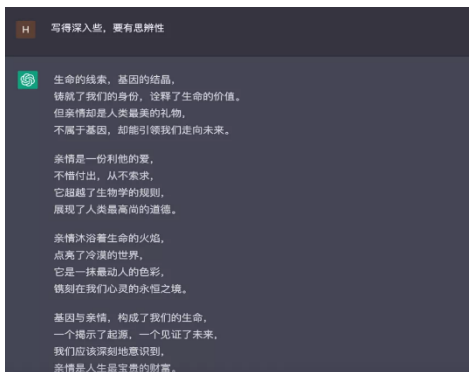


打造生态

ChatGPT
对话、推理、写诗、算术...

X-GPT等
文本分类、情感分析、观点抽取...

ChatPDF等
集成各种工具完成端到端理解

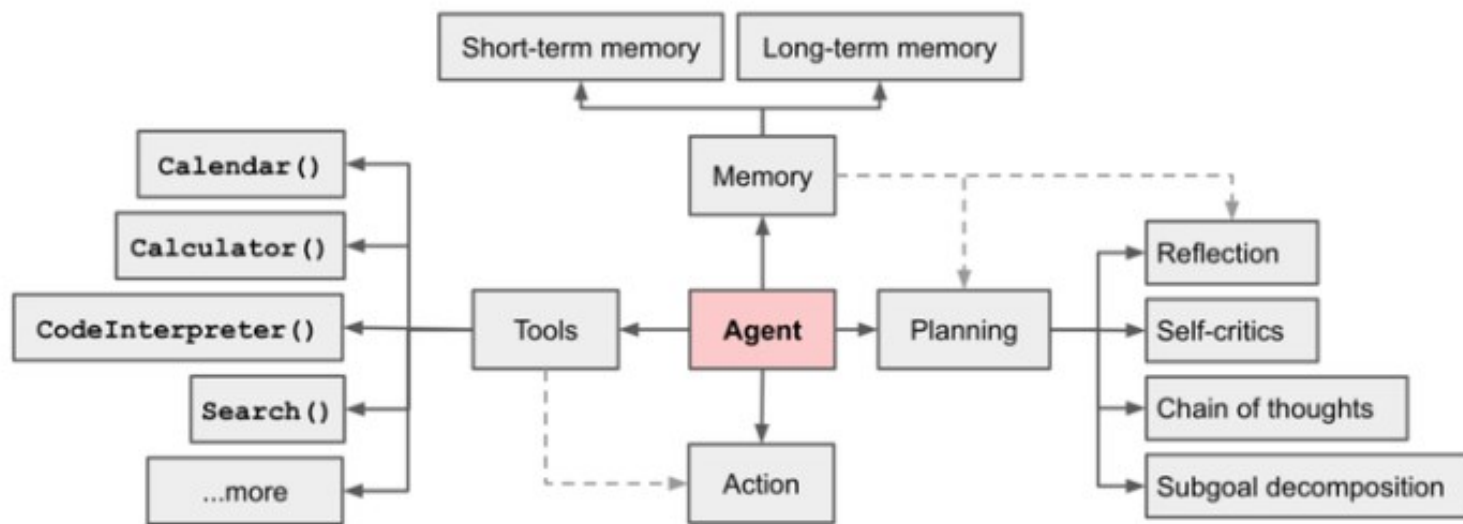


AI驱动软件研发全面进入数字化时代

AI+ 软件研发数字峰会
AI+ software Development Digital summit

Agent时代悄然已至

Yes	Chat Chains <i>ConversationalRetrievalChain</i>	Agents
No	Basic LLM	API / Function Chains <i>APIChain</i>
	No	Yes
	Access to tools	



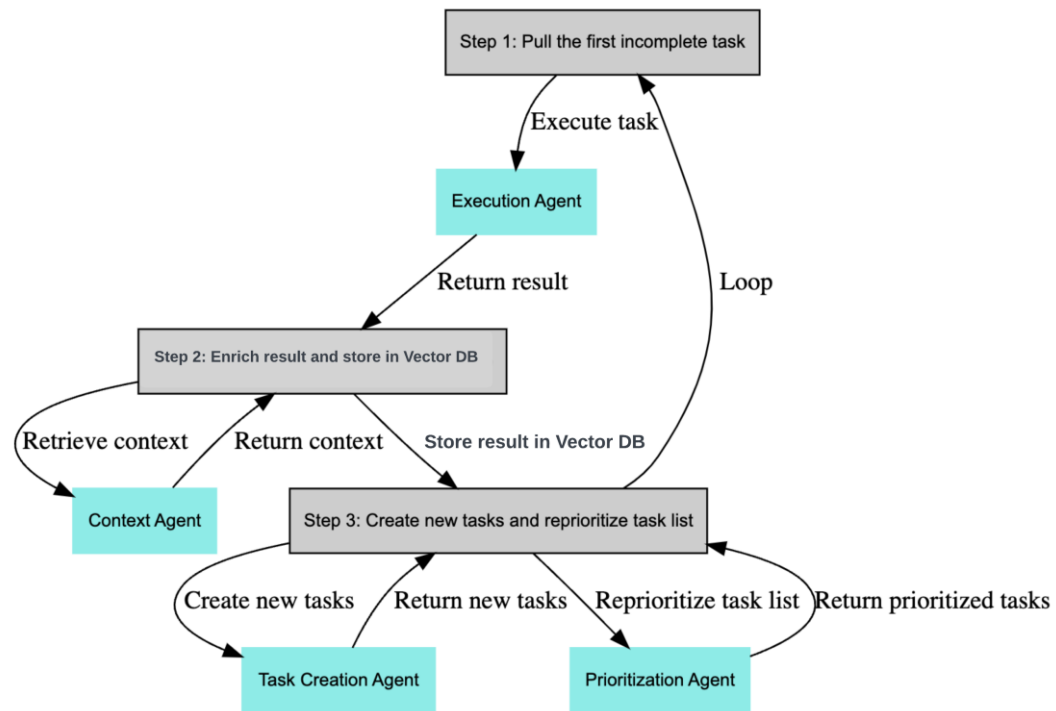
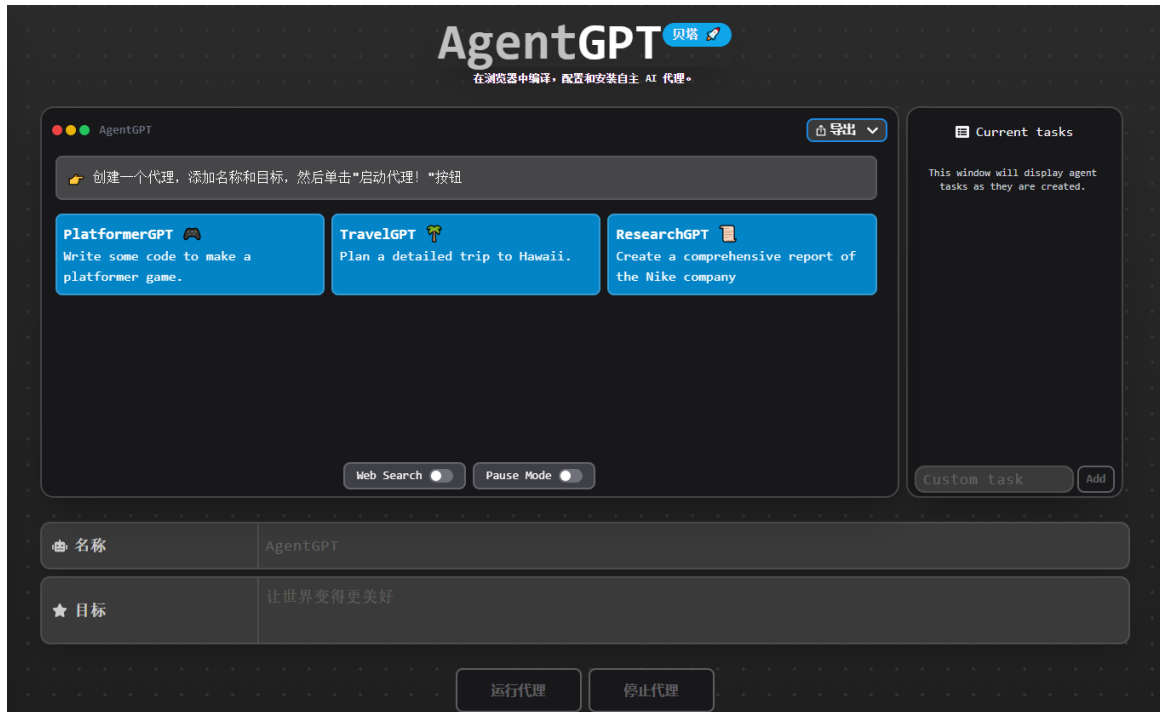
Agent = LLM (大型语言模型) + 记忆 + 规划技能 + 工具使用

自主智能体的崛起

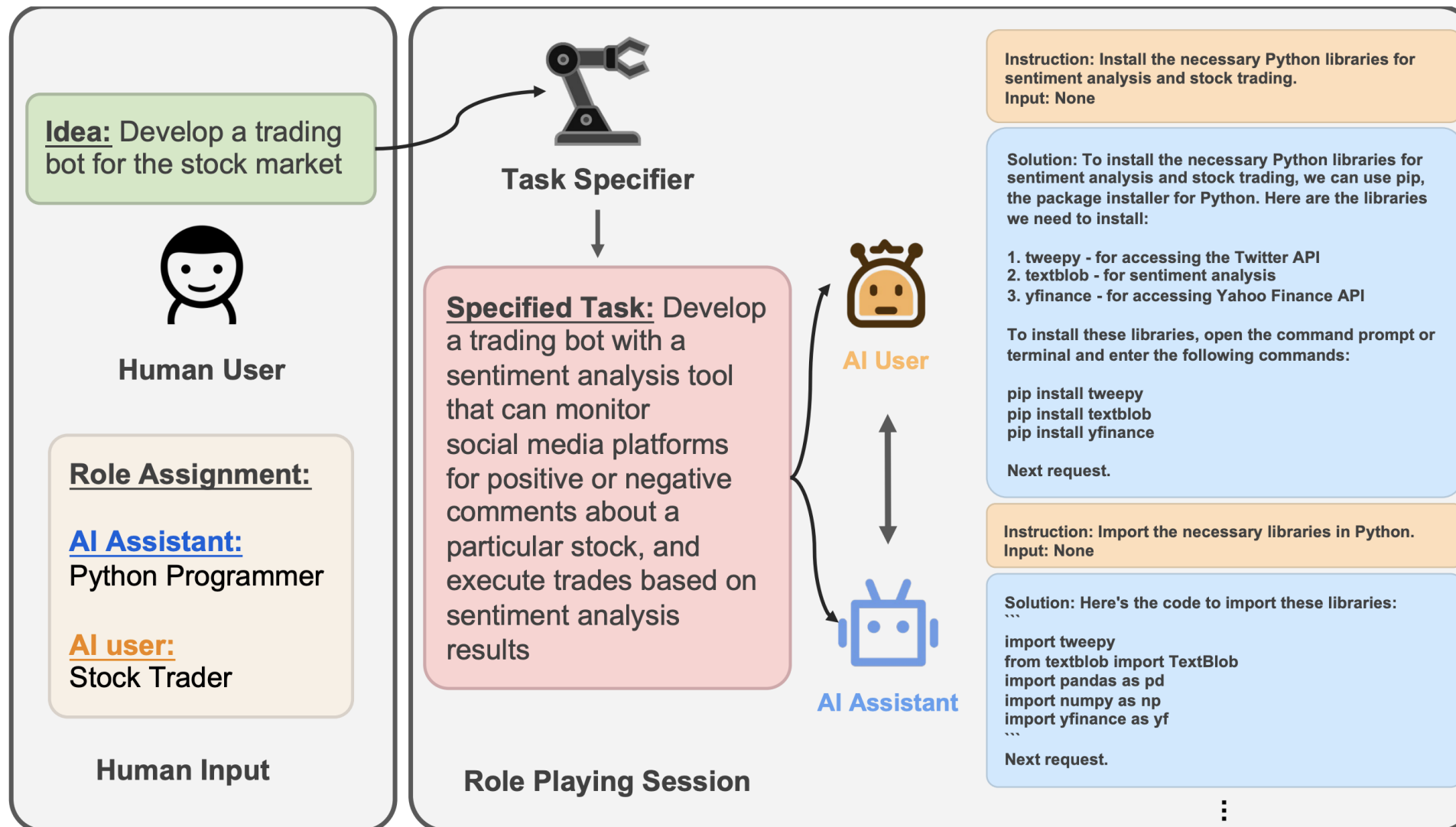
ChatGPT is yesterday's news, now it's all about AI agents

使用ChatGPT等现有LLM时，每步都需要提示，而AI Agent，只需要给它一个总体目标，即可使其自动工作。

如Auto-GPT、AgentGPT、BabyAGI等在通往自主智能体AI Agent的路上



▶ 多智能体协同

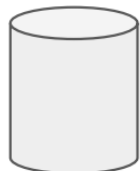


<https://www.camel-ai.org/>

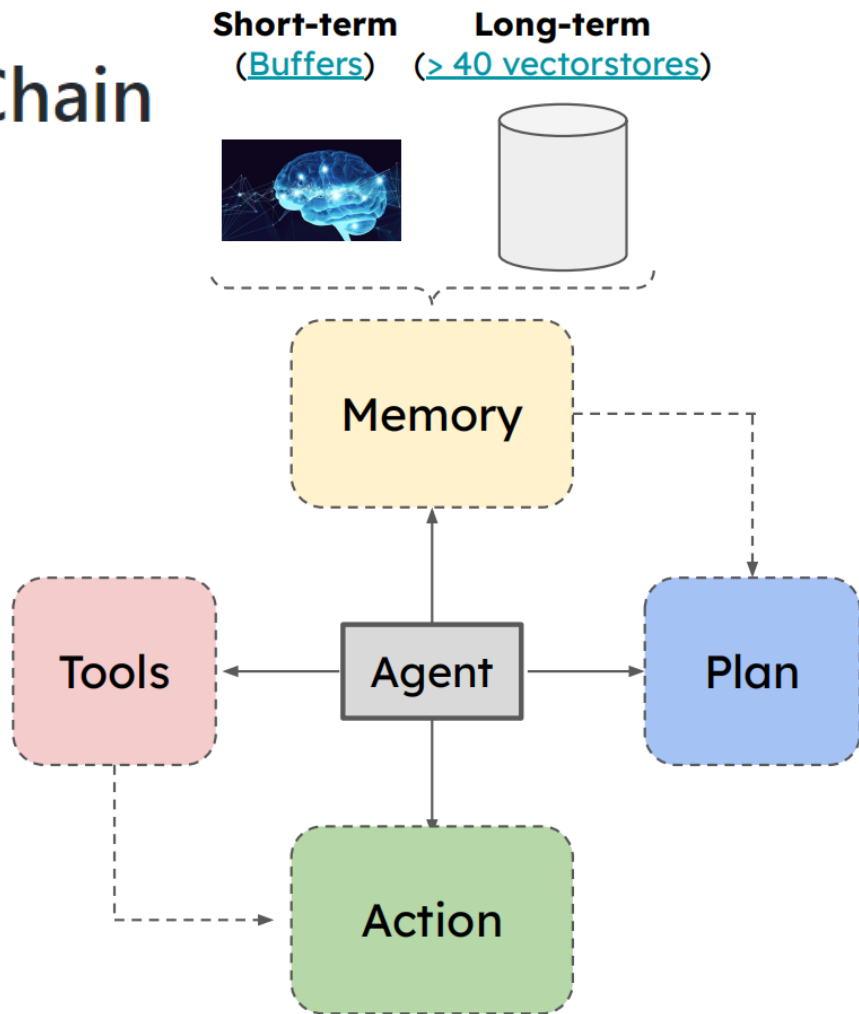
▶ 开源促进Agent生态繁荣



Short-term (Buffers) Long-term (> 40 vectorstores)



Tools (> 60 tools + toolkits)



Agents (> 15 agent types)

Action

```
Thought: ...
Action: ...
Observation: ...
... (Repeated many times)
```

Simulation



Autonomous



LLMs (> 60 integrations)

感谢聆听



OpenKG公众号

知识图谱与大模型技术算法、实战文章、行业案例分享

