

# ADL 147

## Report 1. 知识检索增强：范式与关键技术

同济大学 王昊奋

Retrieval-Augmented Generation for Large Language Models: A Survey [arXiv24.3.27](https://arxiv.org/abs/24.3.27)

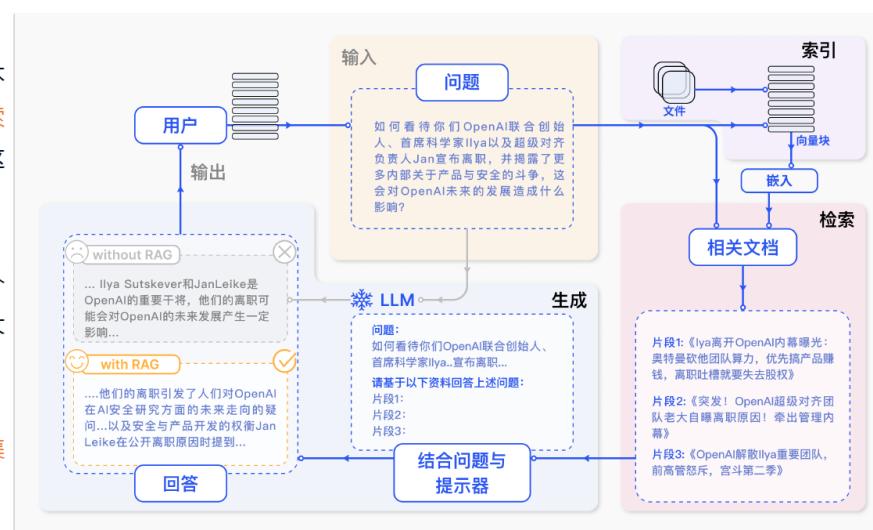
### 01 RAG基本概述

#### RAG提出背景

LLM缺陷：幻觉（一本正经的胡说八道）、信息过时、参数化知识效率低、缺乏专业领域的深度知识、推理能力弱

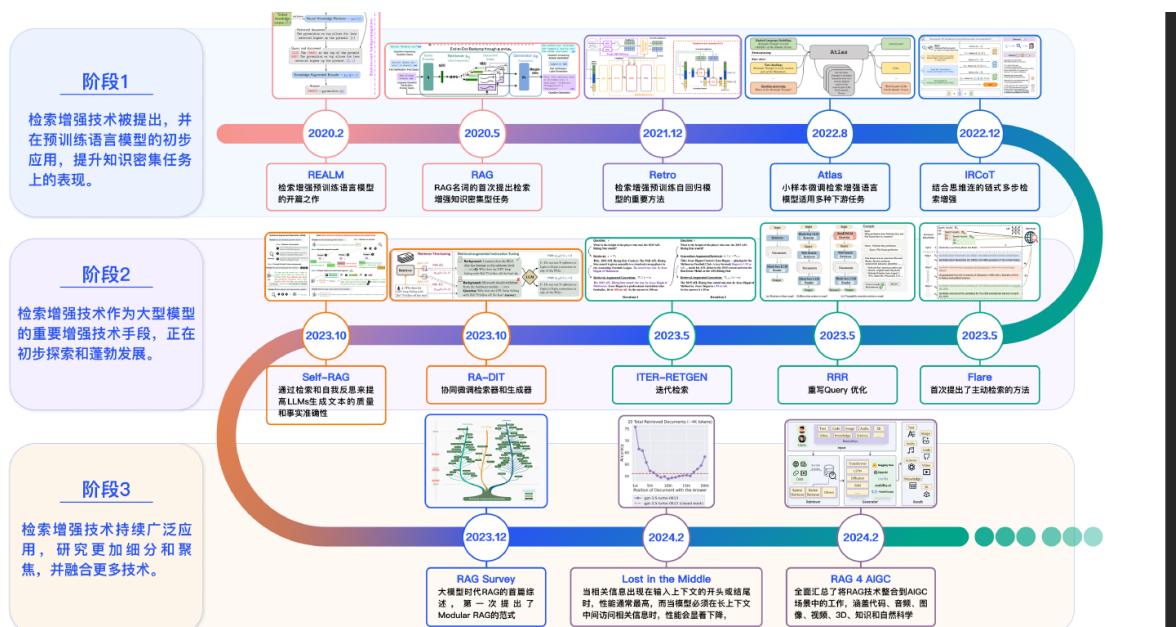
实际应用的需求：领域精准问答、数据频繁更新、生成内容可解释可溯源、成本可控、数据隐私保护

- LLM 在回答问题或生成文本时，先会从大量文档中检索出相关的信息，然后基于这些信息来生成回答。
- RAG 方法使得不必为每一个特定的任务重新训练整个大模型，只需要外挂知识库。
- RAG 模型尤其适合知识密集型的任务。



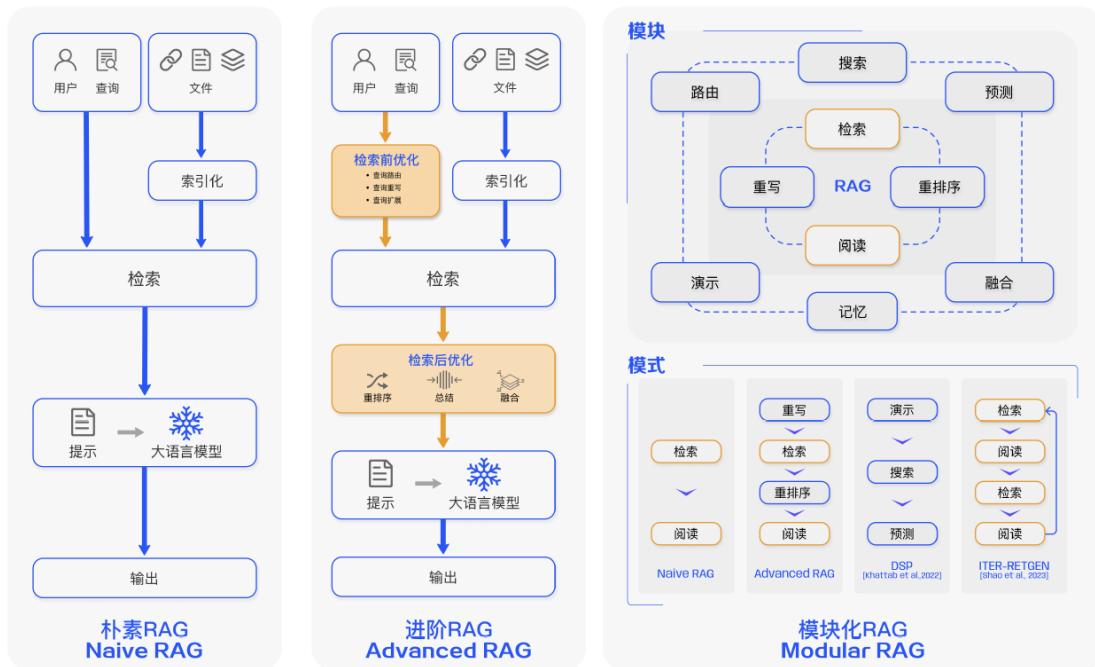
RAG 的主要流程

### 02 RAG的主要范式与发展历程



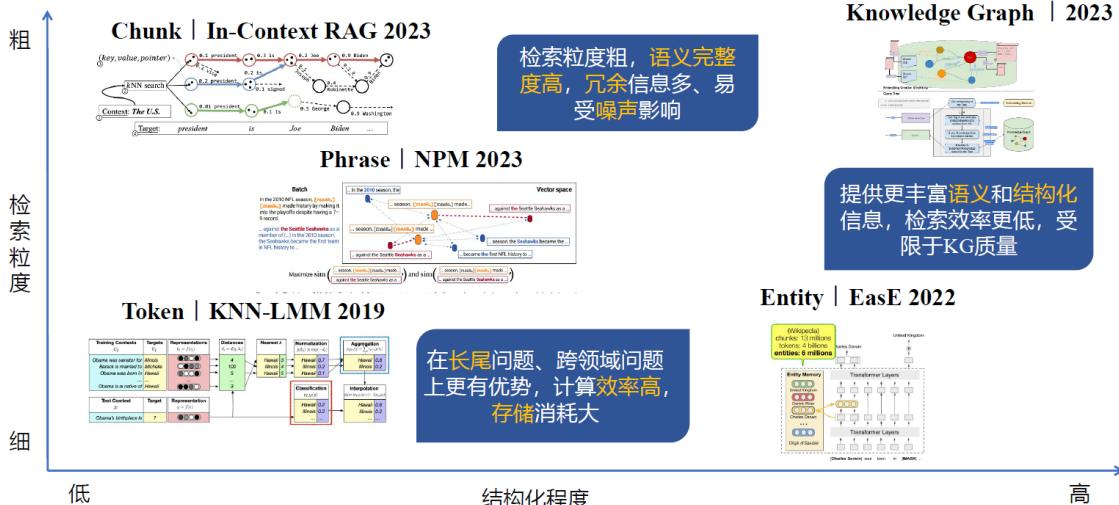
## RAG的三种典型范式：

- Naive RAG：
  - 构建数据索引：文档分块、生成embedding、存储到向量数据库中
  - 检索：向量相似度度量，得到k个文档
  - 原始的query与检索得到的文本组合输入到LLM，得到回答
- Advanced RAG
  - 索引优化：滑动窗口、细粒度分割、元数据
  - 前检索模块：检索路由、摘要、问题重写、置信度判断
  - 后检索模块：重排序、检索内容过滤
- 模块化RAG



## RAG三大问题：

- 检索什么？词元、词组、句子、段落、实体、知识图谱



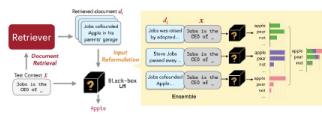
- 什么时候检索？单次检索、每个token、每N个token、自适应检索

效率高，检索到的文档  
相关度低

平衡效率和信息的矛盾  
可能非最优解

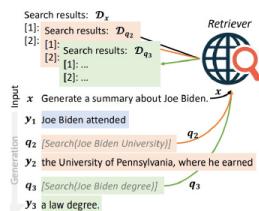
检索到的信息量大，但  
效率低，冗余信息多

Once | Replug 2023



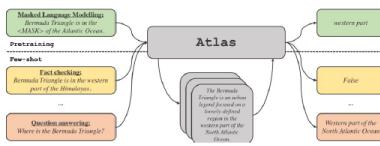
在推理中仅进行一次检索

Adaptive | Flare 2023



自适应地进行检索

Every N Tokens | Atlas 2023



每生成N个Tokens去检索一次

低

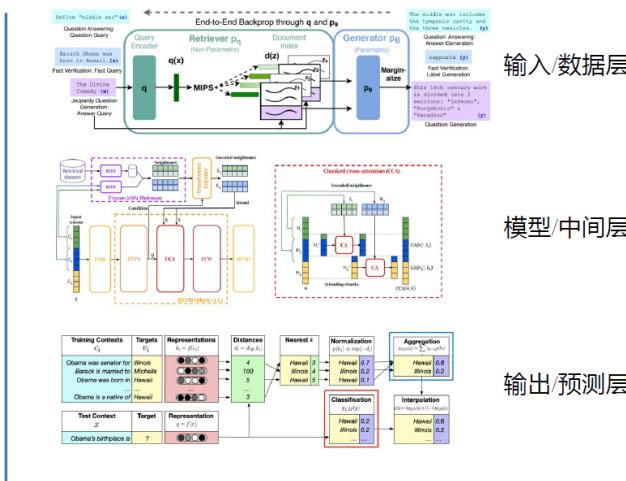
检索频率

高

- 怎么使用检索的结果？输入/数据层、模型/中间层、输出/预测层

- 在推理过程中，集成检索到的信息到生成模型的不同层级中

集成检索位置



使用简单，但无法支持检索更多的知识块，且优化空间有限

支持输入更多的知识块检索，但引入额外的复杂度，且必须训练

保证输出结果与检索内容高度相关  
但效率低

## 03 模块化RAG范式与关键技术

### ► 模块化RAG体系

三层架构

Module Type

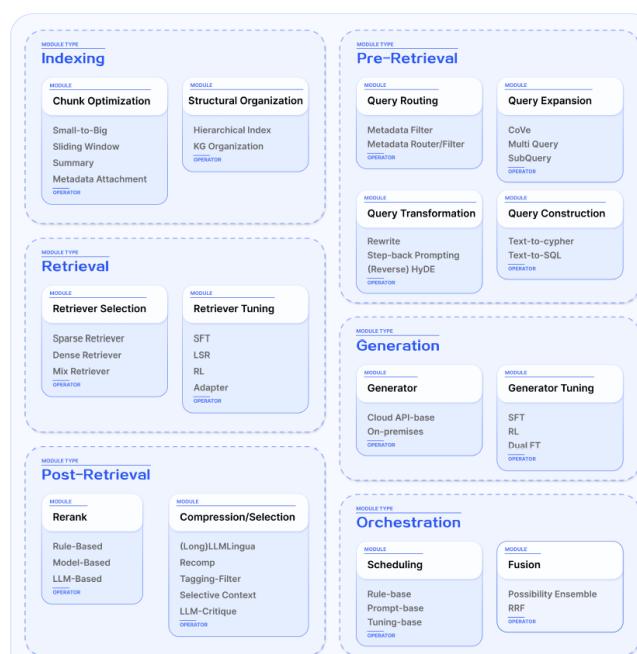
6大模块类型：  
RAG的核心流程

Module

14个模块：  
流程中的具体功能

Operator

40+算子：  
特定功能的具体实现



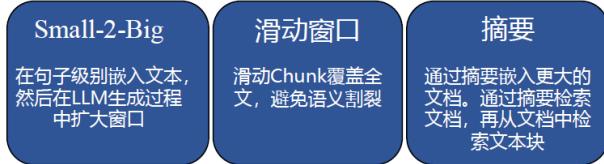
## 1. Indexing 索引

面临挑战：文档块不完整的语义信息、块相似度计算不准确、参考轨迹不明确

解决方案：分块索引优化、结构化语料

### 分块索引优化：

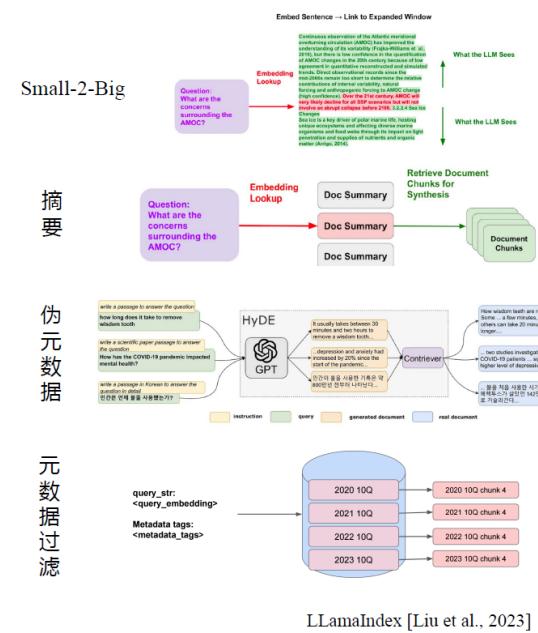
#### 分块策略



#### 添加元数据



#### 元数据筛选/扩充



### 结构化语料：

#### 层级结构

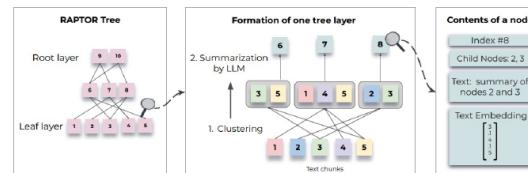
- 文档原始结构、语义结构分割
- 检索摘要->再检索具体的节点
- 图、表作为独立检索对象



Arcus  
[Arcus Team, 2023]

#### 树结构索引

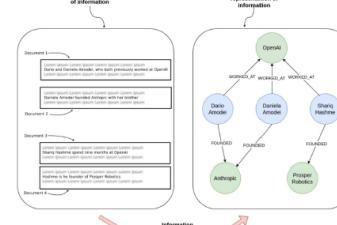
- 语义相似度聚类子节点
- 生成摘要作为父节点
- 递归构建树结构索引



RAPTOR  
[Sarthi et al., 2024]

#### 知识图谱索引

- 文本块 ->三元组实例 ->图结构
- 构建文档块之间的语义关系
- 增强上下文理解和可解释性



Neo4j [Bratanic., 2024]

## 2. Pre-Retrieval 检索前处理

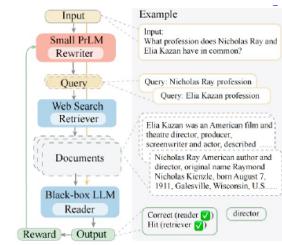
面临挑战：措辞不当的查询、语言自身的复杂性和歧义性

解决方案：查询扩展、查询转换、查询路由

### 查询转换：

## 查询重写 Query Rewrite

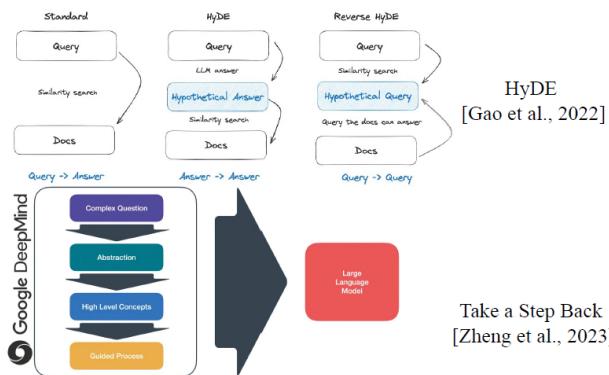
- 原始查询并不总是检索的最佳选择
- LLM / 专用小模型重写查询



Rewrite-Retrieve-Read  
[Ma et al., 2023]

## 假设回答 HyDE

- 构建假设回答代替原始查询去检索
- 或为Chunk生成假设问题作为查询依据



Take a Step Back  
[Zheng et al., 2023]

## 查询拓展:

### 多查询 Multi-Query

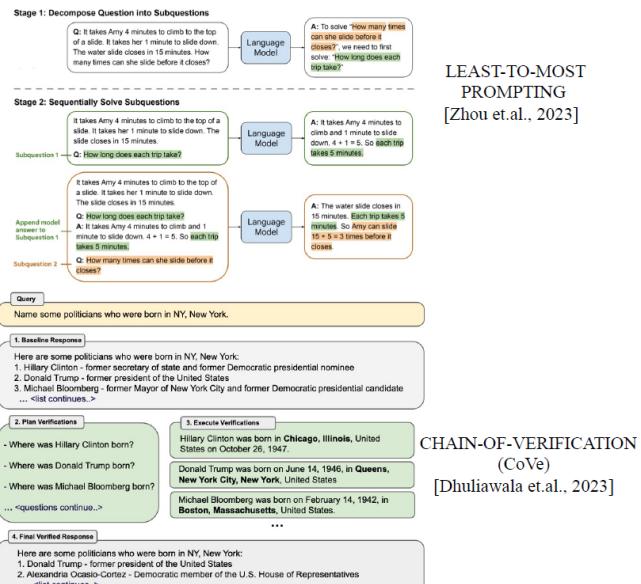
- 提示工程通过 LLM 扩展查询
- 根据预设的模板选择相似的查询
- 为原始查询分配更高权重避免意图稀释

### 子查询 Sub-Query

- 复杂的问题分解成一系列更简单的子问题
- 针对各个子问题分别检索增强生成
- 得到的中间结果与原问题合并

### 验证链 CoVe

- 通过LLM生成一组要问的验证问题
- 通过回答这些问题并检查是否一致执行



## 查询路由:

根据不同的场景，路由到不同的RAG流程、组件或提示词模板

元数据路由器/过滤器 Metadata Router/ Filter

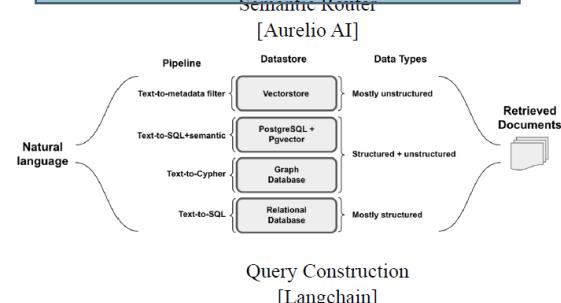
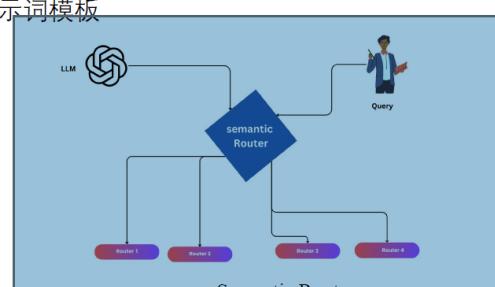
提取查询中关键字（实体），根据关键字和元数据进行筛选

语义路由器 Semantic Router

根据用户自然语言的意图选择RAG的流程  
将基于语义和元数据的方法结合起来

Text-to-Cypher / Text-to-SQL

将用户的自然语言查询转换为另一种查询语言  
以便查询其他结构化数据源



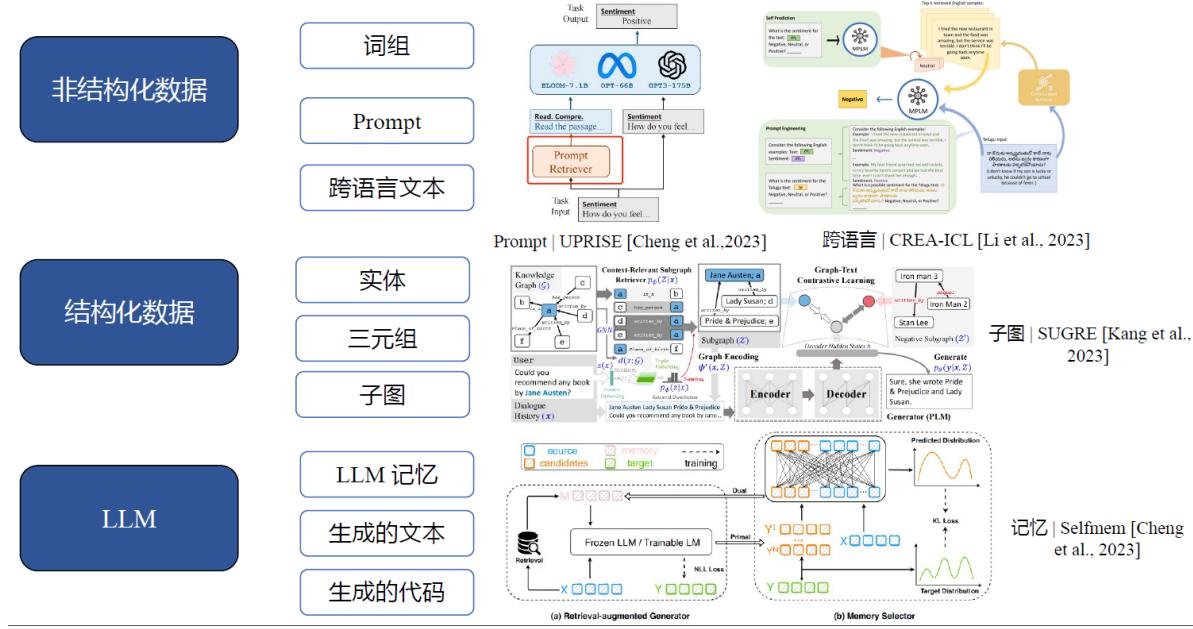
Query Construction  
[Langchain]

### 3. Retrieval 检索

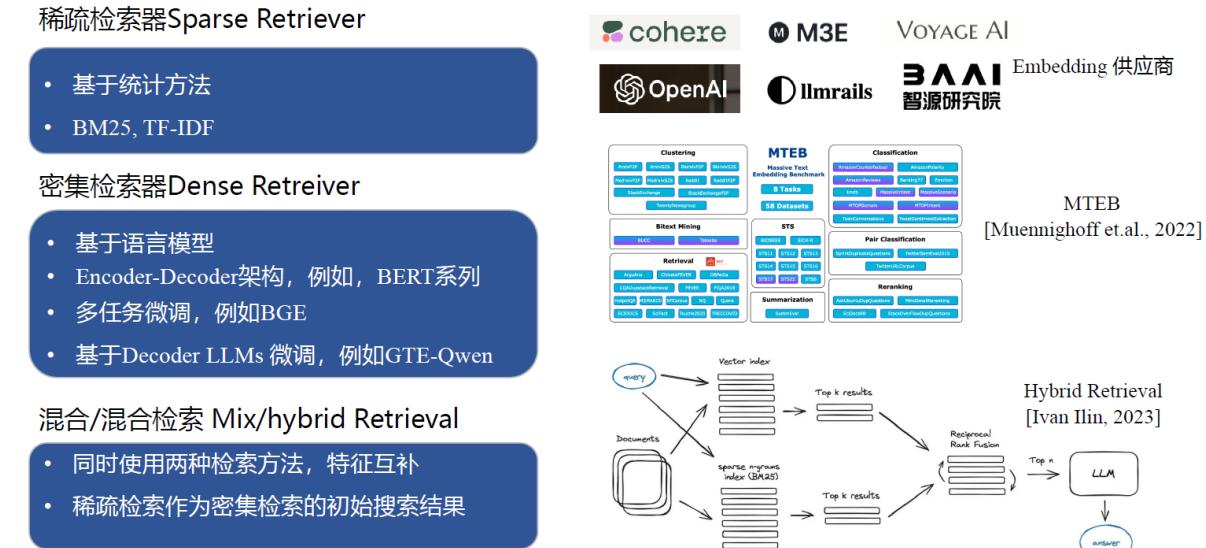
面临挑战：检索效率、嵌入表示质量、任务数据和模型的一致性

解决方案：检索元选择、检索器选择、检索微调

检索元选择：



检索器选择：



检索器微调：

## 有监督微调SFT

- 特定领域 / 通用领域（混合）的监督数据

### 语言模型监督微调 LM-supervised Retriever

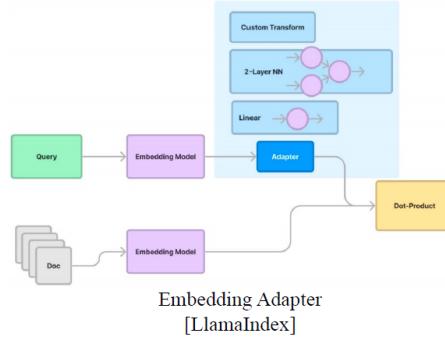
- 利用生成的结果作为监督信号，在 RAG 过程中对嵌入模型进行微调。

### 强化学习微调

- 受RLHF的启发，通过强化学习来强化检索器

## 适配器 Adapter

- 微调检索器的成本较高 / 无法微调的闭源模型
- 特定任务适配器Task Specific
- 通用适配器Task Agnostic



## 4. Post-Retrieval 后检索

面临挑战：噪音/反事实文档、上下文窗口影响

解决方案：重排序rerank、上下文压缩筛选、检索器微调

重排序：

压缩与筛选：

### Selective Context 上下文选择

- 识别并去除输入上下文中的冗余内容。
- 类似“停用词移除”策略。
- 基于LM自信息来评估Token的信息量，保留具有更高自信息量的内容

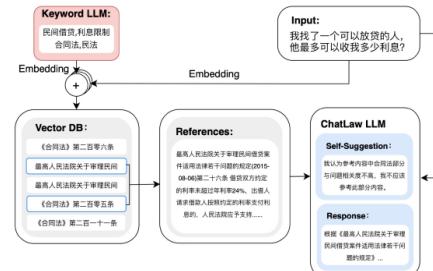
**Original:** INTRODUCTION Continual Learning (CL), also known as Lifelong Learning, is a promising learning paradigm to design models that have to learn how to perform multiple tasks across different environments over their lifetime [To uniform the language and enhance the readability of the paper we adopt the unique term continual learning (CL).]. Ideal CL models in the real world should be deal with domain shifts, researchers have recently started to sample tasks from two different datasets. For instance, proposed to train and evaluate a model on Imagenet first and then challenge its performance on the Places365 dataset; considers more scenarios, starting with Imagenet or Places365, and then moving on to the VOC/CUB/Scenes datasets. Few works propose more advanced scenarios built on top of more than two datasets.

**Filtered:** INTRODUCTION Continual Learning (CL) is a promising learning paradigm to design models how to cross over to uniform the language and enhance the unique term continual learning Ideal CL models should deal with domain shifts. Researchers recently started to sample tasks from two different datasets. For instance, proposed to train and evaluate a model on Imagenet first and then challenge its performance on the Places365 dataset. Consider more scenarios, starting with Imagenet or Places365, and then moving on to the VOC/CUB/Scenes datasets. Few works propose more advanced scenarios built on top of more than two datasets.

Figure 2: A visualisation of selective context. Darker colour indicates larger value of self-information.

### LLM-Critique 语言模型判断

- 直接 LLM 在生成最终答案之前评估检索到的内容。通过 LLM 评论过滤掉相关性较差的文档。
- 例如，在Chatlaw中，LLM被要求对所引用的法律条款进行自我建议，以评估其相关性。



## 5. Generation 生成

面临挑战：LLM的选型、缺乏领域知识、复杂问题推理能力有限、LLM的幻觉

解决方案：生成器选型、生成器微调、事实校验与知识编辑

模型选型：

生成器微调：

- 领域 / 通用指令微调补充领域信息
- 强化学习对齐模型在使用
- 使用更强大的模型作为教师模型
- 适应特定的输入结构 (文本对、图 检索文档上的偏好 微调生成器
- 结构等)

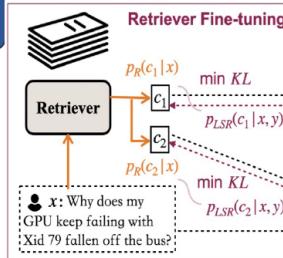
### 生成器与检索器协同微调

#### • R-FT

最小化检索器分布与LLM 偏好之间的KL散度

#### • LM-FT

最大化给定检索增强指令情况下正确答案的可能性



### Retrieval-augmented Instruction Tuning

1 Background: I assume that the BGA chip has damage to the substrate level ... \n Q: Why does my GPU keep failing with Xid 79 fallen off the bus? A:

2 Background: Microsoft should withdraw from the hardware market ... \n Question: Why does my GPU keep failing with Xid 79 fallen off the bus? Answer:

$$\max p_{LM}(y|c_1 \circ x)$$

y: Do not use Y-splitters or 6-pin to 8-pin converters in any of the PCIe...

$$\max p_{LM}(y|c_2 \circ x)$$

y: Do not use Y-splitters or 6-pin to 8-pin converters in any of the PCIe...

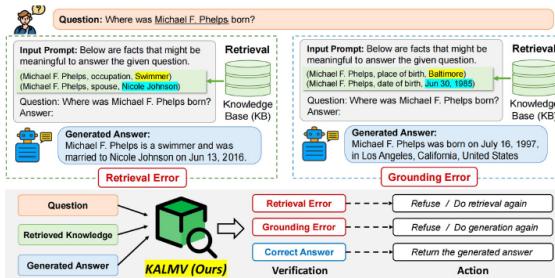
RA-DIT [Lin et al., 2023]

### 事实校验:

经过检索增强并不能确保无幻觉生成，尤其是检索到噪声或冲突事实时，生成后再校验减少幻觉。

### 模型验证

- 用准确知识 (例如KG) 训练的小模型用于大模型的输出校验



KALMV

[Baek et. al., 2023]

### 事实数据交叉校验

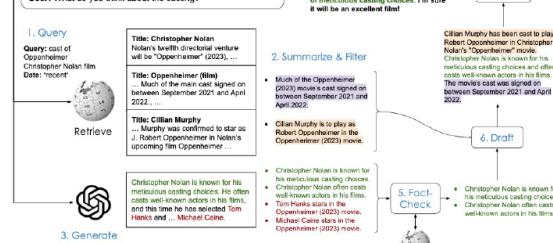
- 用经过校验的知识 (例如Wikipedia) 验证答案的时效性和准确性

User: Have you heard about the upcoming film Oppenheimer?

WikiChat: Yes! Oppenheimer is a biopic ...

User: What do you think about the casting?

Christopher Nolan has cast Cillian Murphy as the titular role of Oppenheimer, keeping with his tradition of meticulous casting choices. I'm sure it will be an excellent film!



WikiChat

[Semnani et. al., 2023]

## 6. Orchestration 编排

面临挑战：传统的链式且一次性的检索-生成流程不足以解决复杂推理或设计大量知识的任务

解决方案：检索流程调度、知识引导检索流程、检索流程聚合

检索流程调度：

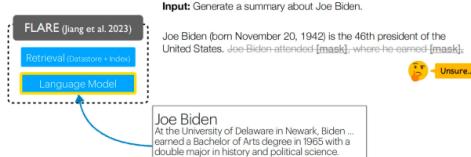
评估RAG过程中的临界点，判断是否需要检索外部文档库、答案的满意度以及是否需要进一步探索。常用于递归、迭代和自适应检索。

### 基于规则

- 生成的答案会被评分，根据分数是否满足预定义的阈值做出继续或停止的决定。
- 常见的阈值包括Token的置信度水平

### 基于提示

- 通过提示工程对当前答案进行反思
- 由LLM动态的判断检索的范围和生成终止



FLARE [Jiang et al., 2023]

### 特殊Token

- 生成过程中根据特殊的Token判断下一步的行为
- 通常需要微调以满足特定的输出格式

Type	Input	Output	Definitions
Retrieve	$x / x, y$	{yes, no, continue}	Decides when to retrieve with $R$
ISREL	$x, d$	{relevant, irrelevant}	$d$ provides useful information to solve $x$ .
ISSUP	$x, d, y$	{fully supported, partially supported, no support}	All of the verification-worthy statement in $y$ is supported by $d$ .
IsUSE	$x, y$	{5, 4, 3, 2, 1}	$y$ is a useful response to $x$ .

#### Algorithm 1 SELF-RAG Inference

```

Require: Generator LM  $\mathcal{M}$ , Retriever  $\mathcal{R}$ , Large-scale passage collections  $\{d_1, \dots, d_N\}$ 
1: Input: input prompt  $x$  and preceding generation  $y_{<t}$ , Output: next output segment  $y_t$ 
2:  $\mathcal{M}$  predicts  $\text{[Retrieve]}$  given  $(x, y_{<t})$ 
3: if  $\text{[Retrieve]}$  == Yes then
4:   Retrieve relevant text passages  $D$  using  $\mathcal{R}$  given  $(x, y_{t-1})$                                 ▷ Retrieve
5:    $\mathcal{M}$  predicts  $\text{[ISREL]}$  given  $x, d$  and  $y_t$  given  $x, d, y_{<t}$  for each  $d \in D$           ▷ Generate
6:    $\mathcal{M}$  predicts  $\text{[ISSUP]}$  and  $\text{[IsUSE]}$  given  $x, y_t, d$  for each  $d \in D$                   ▷ Critique
7:   Rank  $y_t$  based on  $\{\text{[ISREL]}, \text{[ISSUP]}, \text{[IsUSE]}\}$                                 ▷ Detailed in Section 3.3
8: else if  $\text{[Retrieve]}$  == No then
9:    $\mathcal{M}_{gen}$  predicts  $y_t$  given  $x$                                               ▷ Generate
10:   $\mathcal{M}_{gen}$  predicts  $\text{[IsUSE]}$  given  $x, y_t$                                          ▷ Critique

```

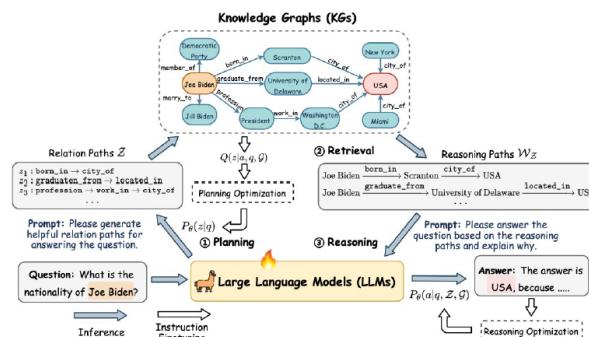
Self-RAG [Asai et al., 2023]

### 知识引导：

通过知识库来指导RAG的检索与生成流程，增强RAG流程的可追溯性，提高生成的可靠性

#### 1. 计划 - 检索 - 推理

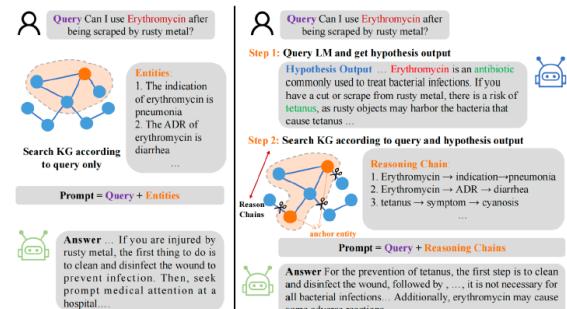
#### 2. 生成以KG为依据的计划执行路径



RoG [Luo et al., 2023]

#### 3. 根据计划从KG和语料中检索相关内容

#### 4. 生成具有推理路径和可解释性的回答



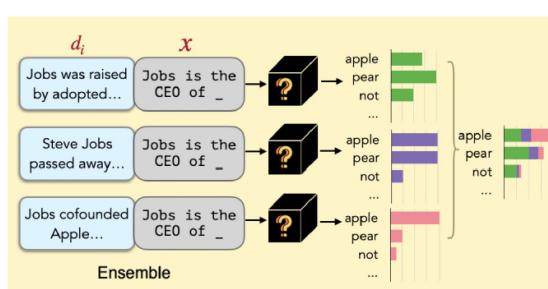
HyKGE [Jiang et al., 2023]

### 流程聚合：

当RAG不再是单一的流水线，多个分支扩展检索范围和多样性，融合模块对多个答案进行融合。

### 基于概率聚合

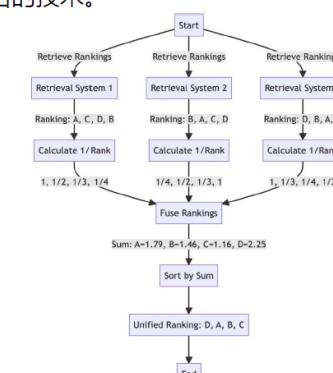
- 融合方法从多个分支生成的不同token的加权值，从而对最终输出进行综合选择。
- 主要采用加权平均法。



RePlug [Shi et. al., 2023]

### Reciprocal Rank Fusion (RRF)

- 将多个搜索结果列表的排名结合起来生成单一统一排名的技术。



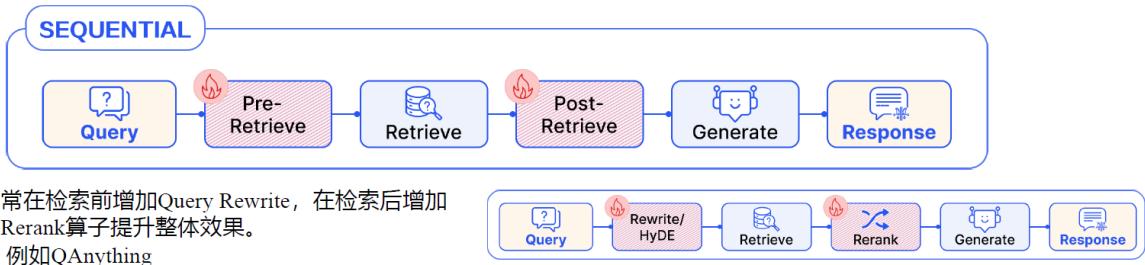
RAG Fusion [Raudaschl, 2023]

## 7. RAG FLOW

Flow Pattern	特点	适用场景
Sequential	简单有效	复杂任务处理能力弱 难度不大的场景 适合作为 <b>快速上手</b>
Conditional	多路由配置，适合多任务场景	每个路由分支仍是Sequential结构，复杂任务能力弱 数据结构丰富，同时包括文本、表格、图结构等
Braching	拆解复杂查询，分别检索生成后聚合，复杂任务处理能力提升	分支后需要聚合，各项开销提高 查询较复杂， <b>易于分解</b> 成多个子问题，容忍一定的时延
Loop	多次RAG，具有较高的自主性和灵活性。复杂任务能力较高	可控性下降，时延和开销较大 任务 <b>困难</b> ，对 <b>时延要求不高</b> ，允许模型有一定灵活性的 <b>非严肃任务</b> 场景

### 线性RAG FLOW

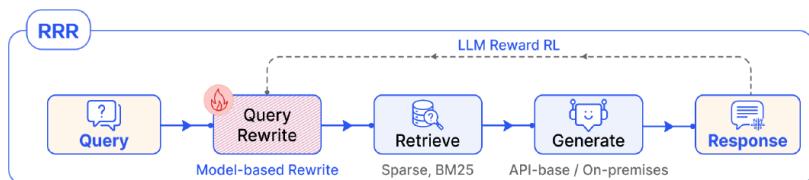
Sequential Flow Pattern 各个模块线性的组织，经典的Naive RAG和基础的Advanced RAG均为线性结构。



常在检索前增加Query Rewrite，在检索后增加Rerank算子提升整体效果。  
例如QAnything

### Rewrite-Retrieve-Read

- Query Rewrite一个小型的可训练LM
- 通过最终LLM的输出结果作为奖励
- 在强化学习的背景下，重写器优化被形式化为一个马尔科夫决策过程

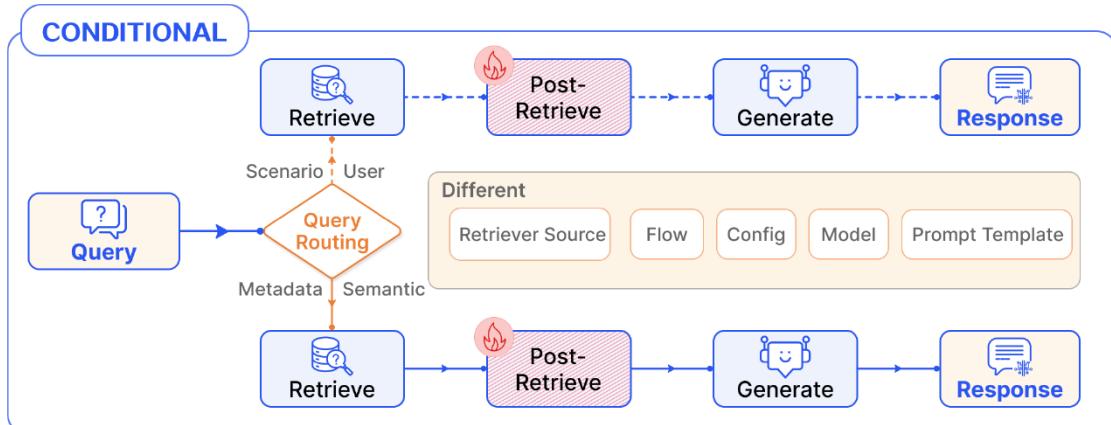


## 条件RAG FLOW

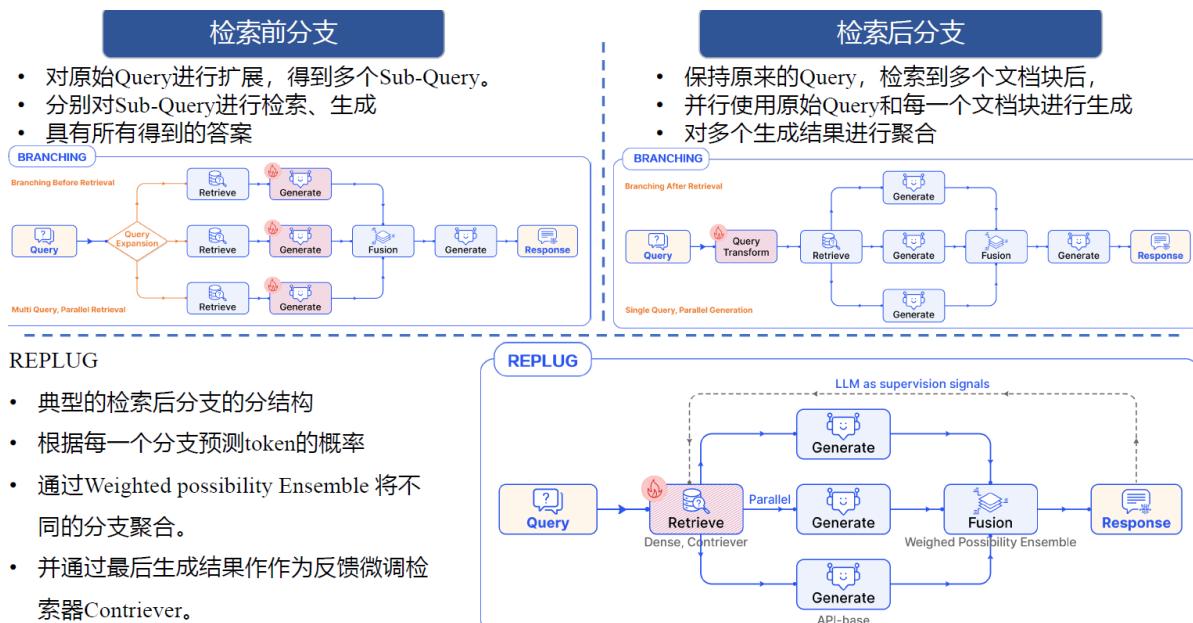
### Conditional Flow Pattern

根据不同的条件选择不同的RAG路线。通常由一个Routing模块进行路由。

- 判断依据包括：语义、任务类型、元数据、用户权限。
- 不同路由分支在检索源、检索流程、配置信息、模型选择和Prompt Template上进行差异化。
- 严肃场景下的任务（例如 医疗、法律等）与娱乐问题，对大模型幻觉的容忍度是不同的。



## 分支RAG FLOW



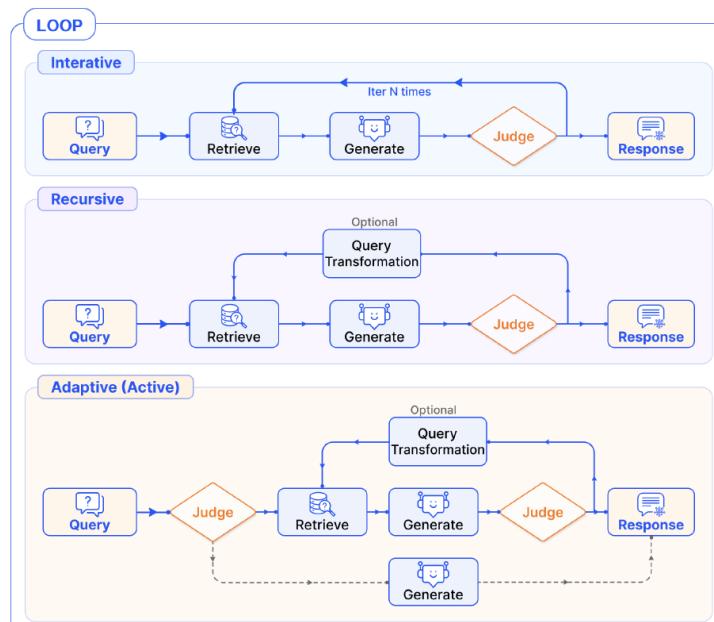
## 环状RAG FLOW

### Loop Flow Pattern

- 面对复杂问题，一次的检索生成效果并不理想，通过多次检索增强
- 具有循环结构的RAG Flow，是Modular RAG的重要特点。
- 检索和推理步骤相互影响的。通常包括Judge模块，用于控制流程。

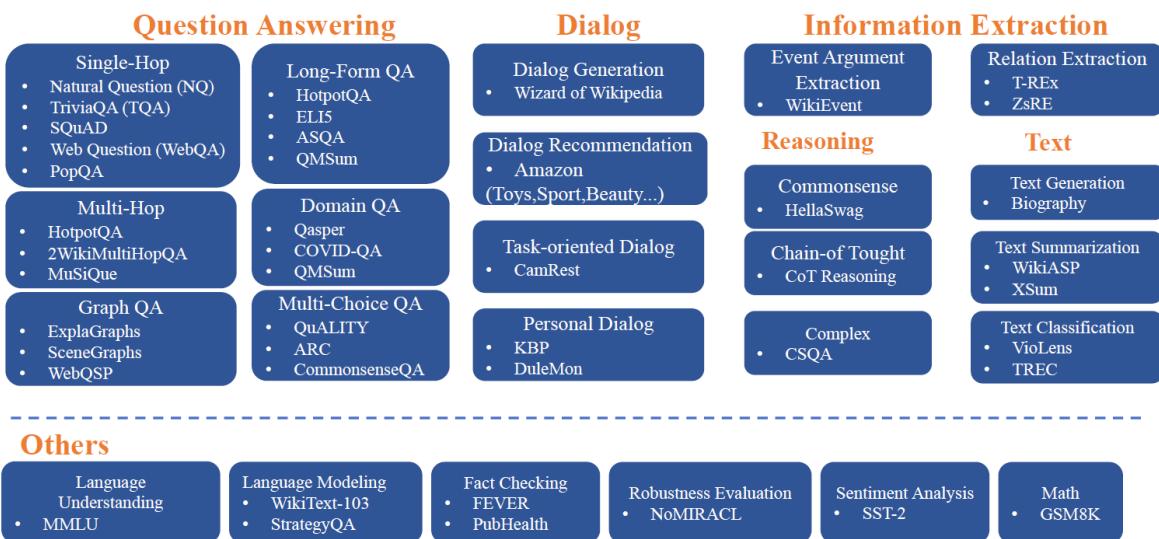
具体又可以分成：

- 迭代检索
- 递归检索
- 自适应（主动）检索



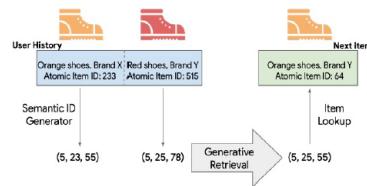
## 04 RAG的下游任务与评估

### 常用数据集：



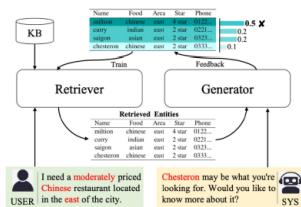
## 下游任务：

进一步扩展RAG的下游任务、完善生态建设



推荐系统 | TIGER [Rajput et al.,2023]

信息抽取 | Filter-Rerank [Ma et al.,2023] 报告生成 | FABULA [Ranade et al.,2023]



对话系统 | ToD [Shi et al.,2023]

医疗问答 | HyKGE [Jiang et al.,2023] 端侧应用 | CT-RAG [Anantha et al., 2023]

## 评估方法：



目前评估依赖于WiKipedia Dump数据集（40G），进行评估成本较高，并且很多评估数据集的语料库都用于模型训练存在泄露问题；外面需要更低成本的更公正的面向RAG的评测体系。

## 05 RAG的工具栈与行业应用

### 个人知识助手

Quivr: Your second brain

WhyHowAI

# 问答系统 LinkedIn 智能客服

## 软件工程 CodeGeeX

### 技术栈与工具

名称	优点	缺点	LangChain	LlamaIndex
LangChain	模块化, 功能全面	行为不一致并且隐藏细节 API复杂, 灵活性低		
LlamaIndex	专注知识检索	需组合使用, 定制化程度低		
FlowiseAI	上手简单, 流程可视化	功能单一, 不支持复杂场景		
AutoGen	适配多智能体的场景	效率低, 需要多轮对话		

- 用途定制化, 满足蛮多样的需求
- 使用简易化, 进一步降低上手门槛
- 功能专业化, 逐渐面向生产环境

AnythingLLM

LangChain-Chatbot

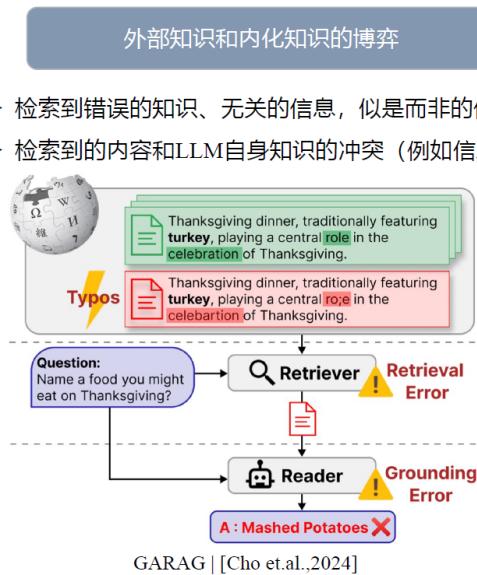
RAGFlow

Anything

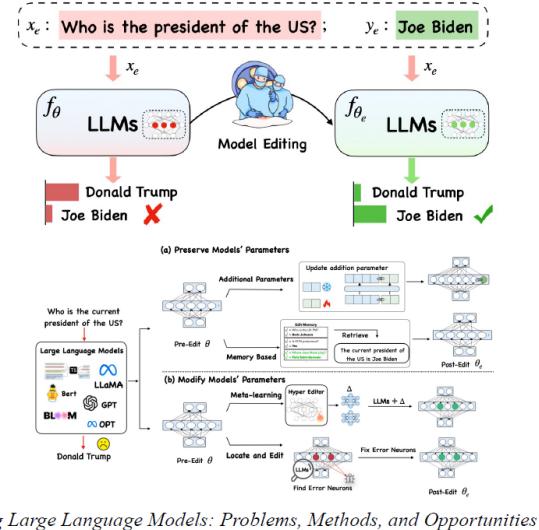
## 06 RAG的挑战与机遇

### 理论侧:

外部知识和内化知识的博弈、RAG记忆与遗忘机制



### 如何做外部知识的注入与编辑机制?



## 大规模知识库的管理

完成RAGDemo是容易的，构建企业级的RAG应用是困难的：

- 在大规模知识库的设计与构建
- 大规模数据上的检索与推理效率
- 知识的淘汰与更新

## 异构数据的整合与存储

- 不同格式数据中不同的向量整合问题
- 超长维度的Embedding存储（例如7111维度）
- Embedding模型更新后，如何与旧的知识库整合
- （例如之前使用了Ada002，现在改用text-embedding3）

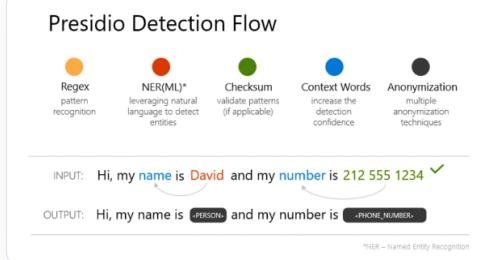
## 技术侧：

幻觉控制、隐私保护、自主控制的RAG流程、流程可追溯 结果可解释、查询与知识的语义距离

### 隐私保护

语料库中可能存在着大量的隐私信息

- 希望既可以被用来补充上下文，同时又不希望他们直接泄露出去
- 不同用户使用RAG系统的权限控制
- 如何利用无法传输到云端的隐私数据



### 幻觉控制

- 尽管目前的RAG方法有效降低了LLM的幻觉，但仍无法确保零幻觉。
- 在某一些严肃场景下，需要LLM接近100%的零幻觉。如何做生成后的校验，实现零幻觉



法律场景 (LexiLaw)



医疗场景 (Ming明医)

## 多模态数据检索与理解

代码、表格、图像、视频都拥有丰富的知识  
如何在RAG的范式下将这些数据关联起来，  
共同作为LLM的上下文

## 长文本的高效切分、向量化与检索

当用户上传了一份超长的文件，例如小说、著作、  
博士论文、企业年报  
如何在可接受的时间内响应用户的查询？

### 自主控制的RAG系统

- LLM越来越多的参与了RAG的流程，例如生成检索语料、决策RAG流程

我们距离自主控制的RAG系统还有多远？

- LLM学会人类的学习方法：从问题中抽象、联想，模仿，举一反三
- 自主判断RAG流程：
  - 是否需要去检索
  - 何时停止检索
  - 生成的答案是否可靠
  - 何时结束生成

### 流程可追溯,结果可解释

- 有时候重要的不是回答什么，而是如何解释你的回答
- 当LLM给出错误、缺乏逻辑的回答时，我们会想要知道为什么，以及如何修正



推荐任务中的可解释推荐

## 应用侧：

大规模工业应用、低资源任务场景

工程应用要求

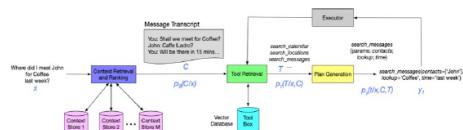
- 大模型与RAG的工业应用对工程应用提出了更高的要求
- 响应时延、模型选择、检索效率、推理效率、存储成本、计算成本、性能冗余等等因素之间的权衡



Kimi有点累了，可以晚点再问我一遍。

低资源场景

- LLM参数受限、Embedding模型的维度受限，编码和检索速度要求高
- 在资源受限的情况下，例如端侧场景中，如何构建RAG应用



端侧应用 |Apple CT-RAG [Anantha et al., 2023]

## Report 2. 大语言模型与智能信息检索技术的融合

中国人民大学高瓴人工智能学院 窦志成 朱余韬

### PART 1 大模型赋能的信息检索

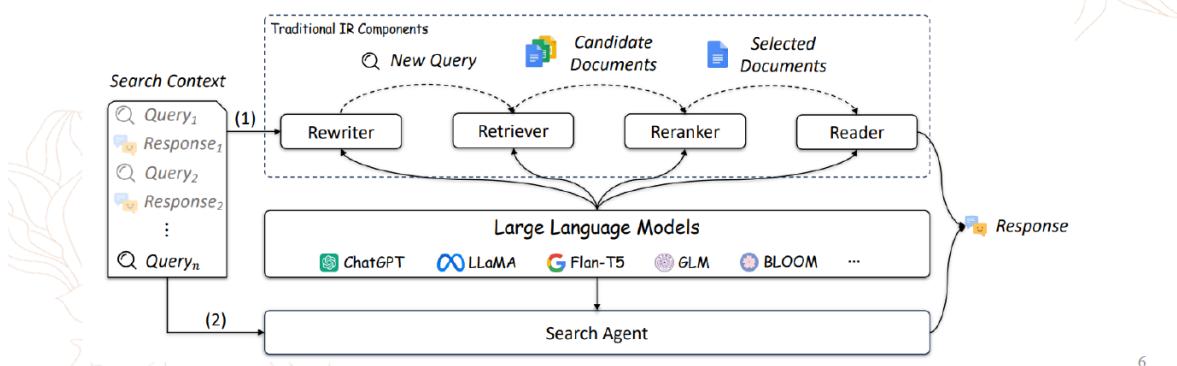
Large Language Models for Information Retrieval: A Survey [arXiv23.8.14](https://arxiv.org/abs/2308.14)



LLM4IR-Survey Public

This is the repo for the survey of LLM4IR.

MIT License · 31 · 334 · 0 · 0 · 4d

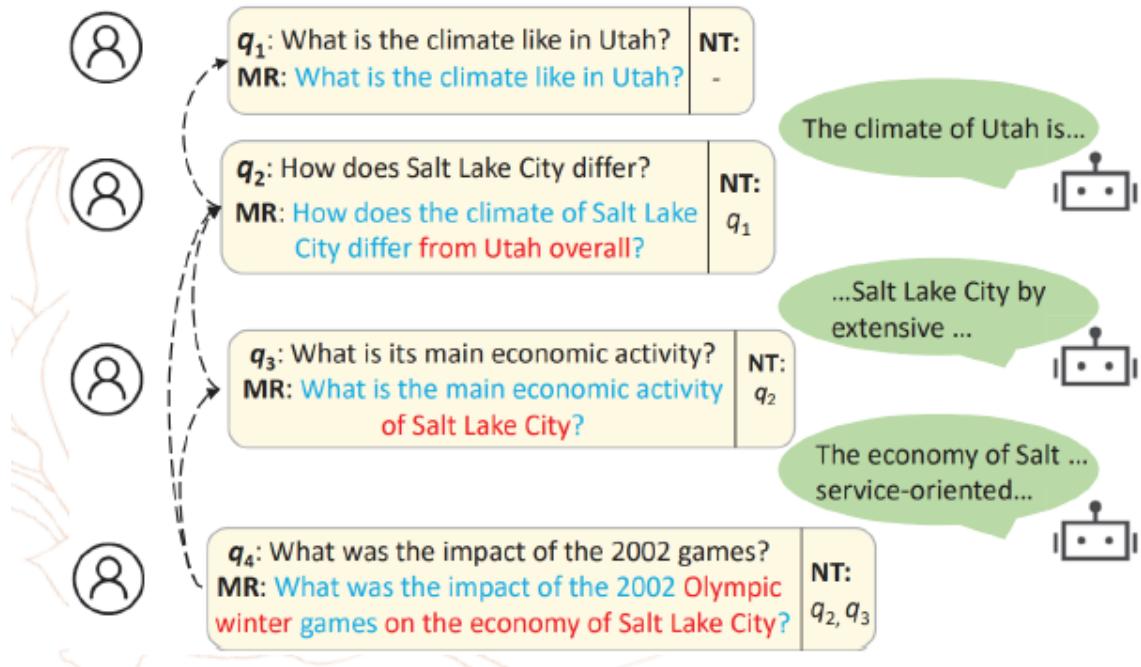


6

#### 1. Rewriter 查询改写

为什么要进行查询改写？

- 原查询过短或模糊，大模型可以更好的理解用户意图
- 在对话系统中，改写更为重要，继承上下文



### 基于LLM进行对话式查询改写

Large Language Models Know Your Contextual Search Intent:A Prompting Framework for Conversational Search [EMNLP 2023](#)

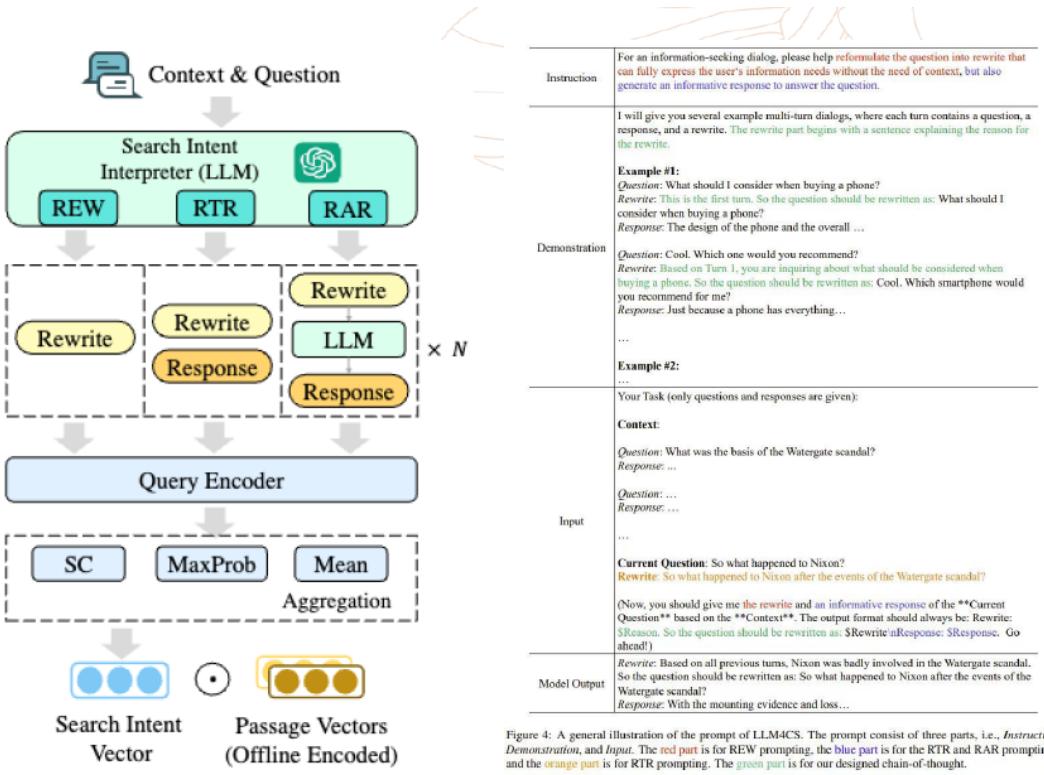


Figure 4: A general illustration of the prompt of LLM4CS. The prompt consist of three parts, i.e., *Instruction*, *Demonstration*, and *Input*. The red part is for REW prompting, the blue part is for the RTR and RAR promptings, and the orange part is for RTR prompting. The green part is for our designed chain-of-thought.

- 常见的三种提示学习方法
  - Zero-shot
  - Few-shot: 提供样例
  - 思维链: 要求生成答案前提供推理过程



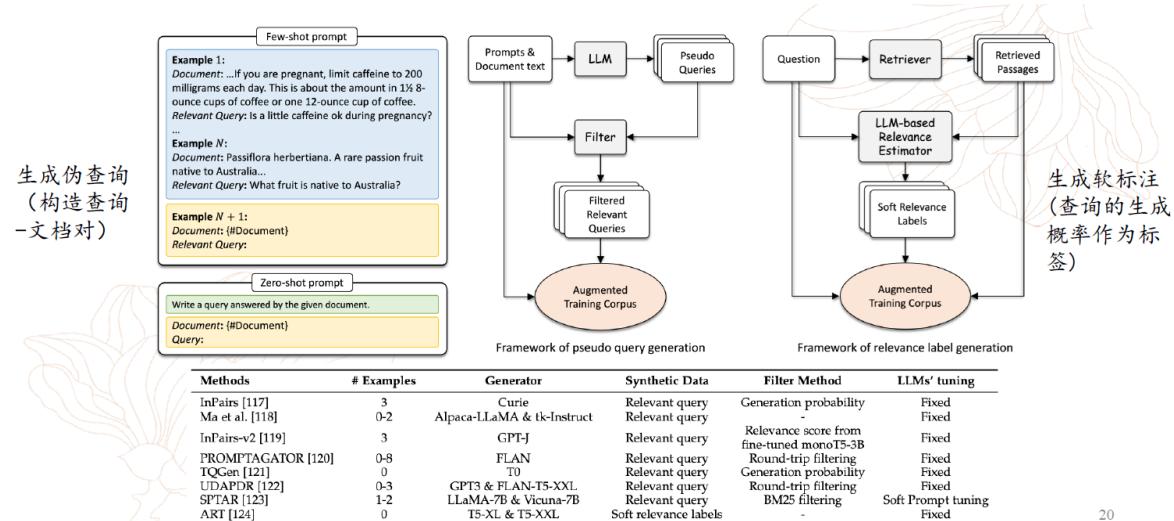
Methods	Prompts
<i>Zero-shot</i>	
HyDE [92]	Please write a passage to answer the question. Question: {#Question} Passage:
LameR [99]	Give a question {#Question} and its possible answering passages: A. {#Passage 1} B. {#Passage 2} C. {#Passage 3} ... Please write a correct answering passage.
<i>Few-shot</i>	
Query2Doc [92]	Write a passage that answers the given query: Query: {#Query 1} Passage: {#Passage 1} ... Query: {#Query} Passage:
<i>Chain-of-Thought</i>	
CoT [93]	Answer the following query: {#Query} Give the rationale before answering
CoT/PRF [93]	Answer the following query based on the context: Context: {#PRF doc 1} {#PRF doc 2} {#PRF doc 3} Query: {#Query} Give the rationale before answering

## 2. Retriever 检索器

从海量文档中高效高质的返回相关结果。

挑战：查询模糊、文档内容多、信息复杂、标注开销大等

### 基于大模型生成检索器的训练数据

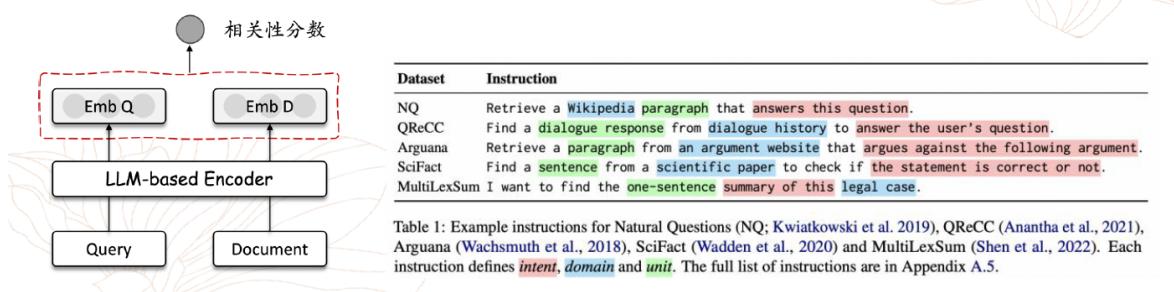


### 以大模型为基座的检索模型

Task-aware Retrieval with Instructions [arXiv2211](https://arxiv.org/abs/2211.03535)

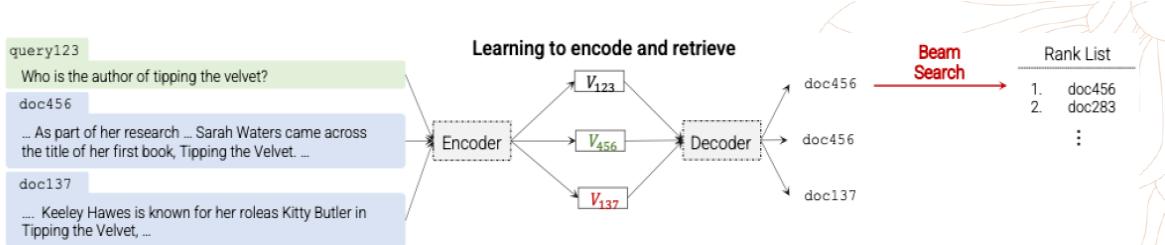
#### 稠密检索与双塔结构

- 通用模型：使用GPT、T5-XXL等大模型初始化编码器并继续优化
- 任务感知模型：使用提示加入任务信息

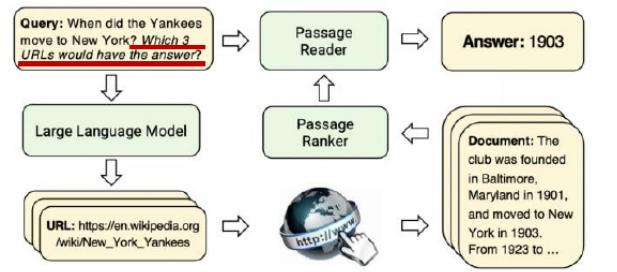


## 大模型改善生成式检索

Transformer Memory as a Differentiable Search Index [NeurIPS 2022](#)



Large Language Models are Built-in Autoregressive Search Engines [ACL 2023](#)



- 大模型生成网址 (for Wiki)
- 用对应内容+Reader完成查询

## 3. Reranker 重排序器

### 微调大模型做重排序

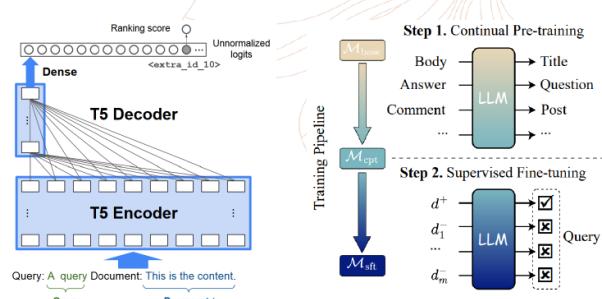
通常有好的性能，但训练开销大

#### • 不同结构

- Encoder-only: [CLS] query [SEP] document [SEP]
- Encoder-Decoder: RankT5
- Decoder-only: RankLLaMA, RankingGPT
  - query: {query} document: {document} [EOS]

#### • 不同训练方式

- 使用生成目标函数：拼接查询和文档，使用生成True/False的概率做优化
- 使用排序函数：拼接查询与文档，使用特殊token生成的概率作为相关性分数，并使用排序损失函数优化



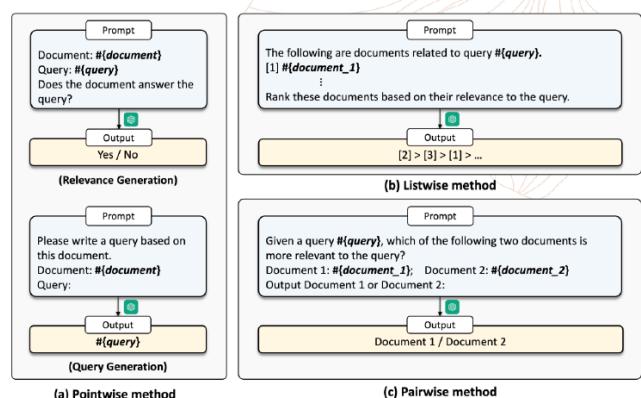
RankT5: Fine-Tuning T5 for Text Ranking with TSRankLLM: A Two-Stage Adaptation of LLMs for Text Ranking, Zhuang et al., 2023

### 提示大模型做重排序

需要大模型能力足够强大

#### • 微调的成本问题 → 提示学习

- Pointwise: 判断查询和文档之间是否相关
- Listwise: 重排文档列表
- Pairwise: 比较两个文档与查询之间的相关性强弱

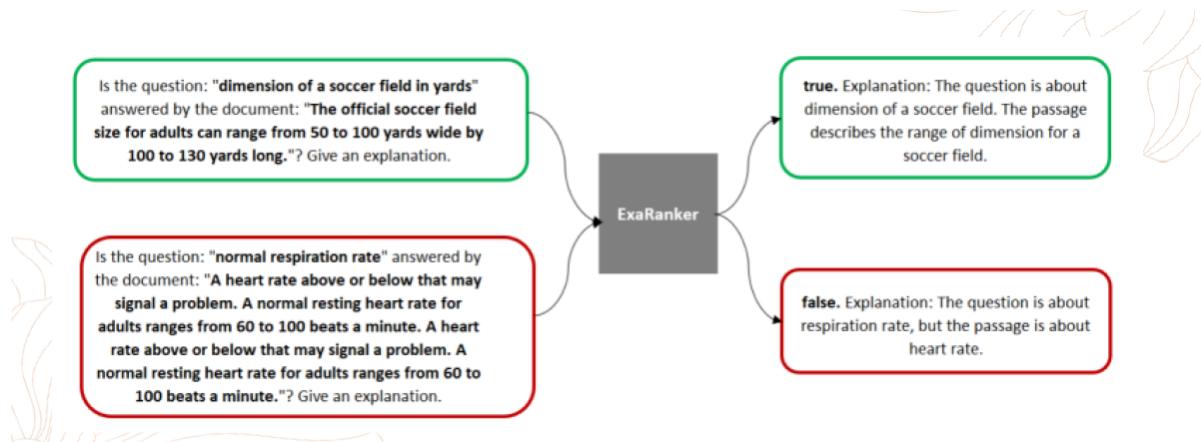


## 大模型生成排序数据

提升已有排序模型的有效策略

ExaRanker: Explanation-Augmented Neural Ranker [arXiv2301](#)

使用chatGPT生成解释作为额外的标签



## 4. Reader 阅读器

基于大模型对检索到的文档进行提炼总结，得到最终的答案输出

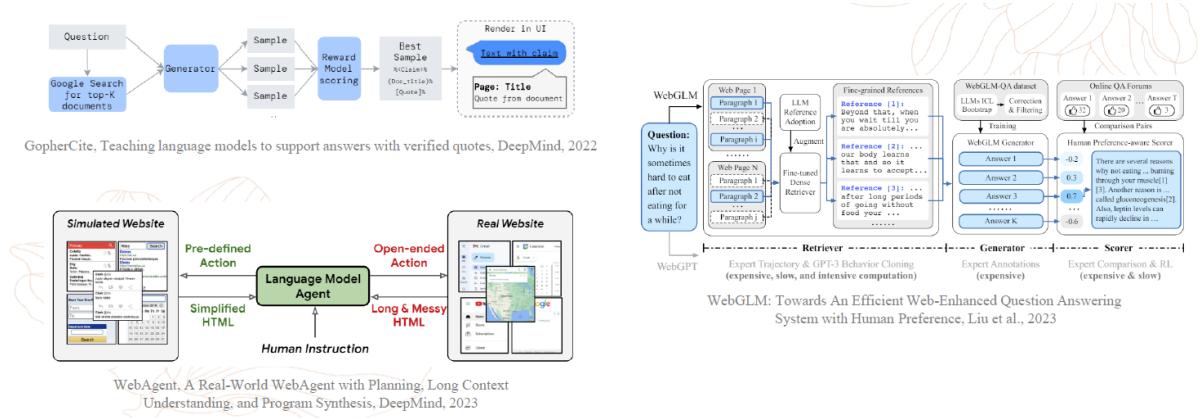
- 新型搜索引擎
  - New bing 百度AI搜索
- 商业大模型
  - Kimi chat Baichuan
- 效果仍有巨大提升空间
  - 幻想
  - 引用不相关内容
  - 编造内容
  - 错误编号

## 5. 搜索Agent

进入强化学习的思想

### 静态Agent

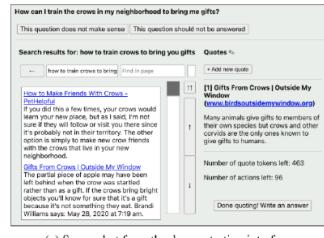
将人类浏览网页的过程拆解为子过程逐模块的使用Agent进行模拟



## 动态Agent

智能体自行决定行为

- Agent自行决定行为



How can I train the crows in my neighborhood to bring me gifts?  
This question does not make sense... This question should not be answered.  
Search results for: how to train crows to bring you gifts Quotes: 0  
How to train crows to bring you gifts  
[1] Gifts From Crows | Outside My Window [www.birdsocietywindow.org](http://www.birdsocietywindow.org)  
Many animals give gifts to members of their own species but crows and other corvids are especially good at giving gifts to humans...  
#past actions  
Starts off by training crows to bring you gifts...  
Click Gift From Crows | Outside My Window [www.birdsocietywindow.org](http://www.birdsocietywindow.org)  
Back  
Title  
Birdsociety... results for: how to train crows to bring you gifts  
Overall: 0 / 11  
Next: How to Make Friends With Crows - PetHelpful.petHelpful.com  
If you did this a few times, your crows would learn your new place, but we're not sure if they will follow or visit you there since crows are territorial. If you want to make sure they'll probably stay in their territory. The other option is simply to make new crows friends with you in your new neighborhood!  
Gifts From Crows | Outside My Window...  
The best way to train crows to bring you gifts is to leave them with bait behind when the crow was started training. If you leave them with bright objects you know for sure that it's a gift because they are more likely to eat something they eat!  
Birdsociety... says: May 28, 2020 at 7:19 am.  
Actions left: 96  
Done quoting! Write an answer.

(b) Corresponding text given to the model.

A: WebShop search item-detail  
I'm looking for a small portable folding desk that is already fully assembled; it should have a khaki wood finish, and price lower than \$40.00 dollars  
Results: 1 Search  
Description Product: Item Description: A portable desk. It's made of steel pipe. It has a khaki wood finish. It's a simple desk with a single drawer. It's designed for a laptop or a computer. It's a compact desk that can be used as a computer desk, dining table, bedside table, or even a desk for a child. It's perfect for small spaces.  
Color: black, khaki, white  
Buy Now: 5 Reward: 1.0  
B: WebShop search item-detail  
Instruction: I'm looking for a small portable folding desk that is already fully assembled [...] Description: MENHG Folding Laptop Table Bed...  
Price: \$109.0 Yopt(OPTIONS): {black, khaki, white} Yatt(ATTRIBUTES): {steel pipe, no assembly, portable}  
C: WebShop search item-detail  
Instruction: I'm looking for a small portable... Description: MENHG Folding Laptop Table Bed...  
Price: \$109.0 Yopt(OPTIONS): {black, khaki, white} Yatt(ATTRIBUTES): {steel pipe, no assembly, portable}

WebGPT: Browser-assisted question-answering with human feedback, OpenAI, 2022

WebShop: Towards Scalable Real-World Web Interaction with Grounded Language Agents, Yao et al., 2023

## 6. ACL24 面向信息检索任务的指令微调

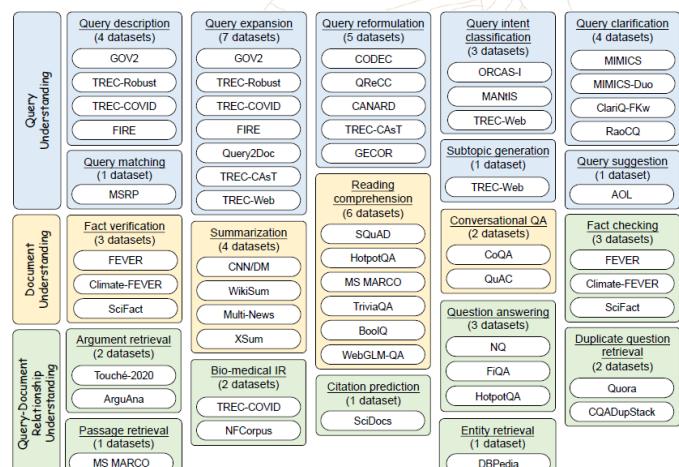
INTERS: Unlocking the Power of Large Language Models in Search with Instruction Tuning

### 动机

- 大语言模型在NLP任务上表现出了极强的能力
- 指令微调进一步实现了大模型与人类任务的对齐
- 在IR任务上，大模型（基于提示）没有显著优于小模型
- 大模型的预训练中缺乏对IR概念的理解
  - 例如：查询、文档、相关性、用户意图等
- 已有指令微调数据集缺乏IR相关任务
- 解决方法 ⇒ 构建面向IR任务的指令微调数据集

### 基本思路

- 分析划分已有IR任务
  - 3个任务类
    - 查询理解
    - 文档理解
    - 查询-文档关系理解
  - 20个任务
  - 43个数据集
- 收集并构建指令微调数据集
- 进行大量实验与分析



## Part 2 检索增强的生成大模型

### 1. 为什么要做检索增强？

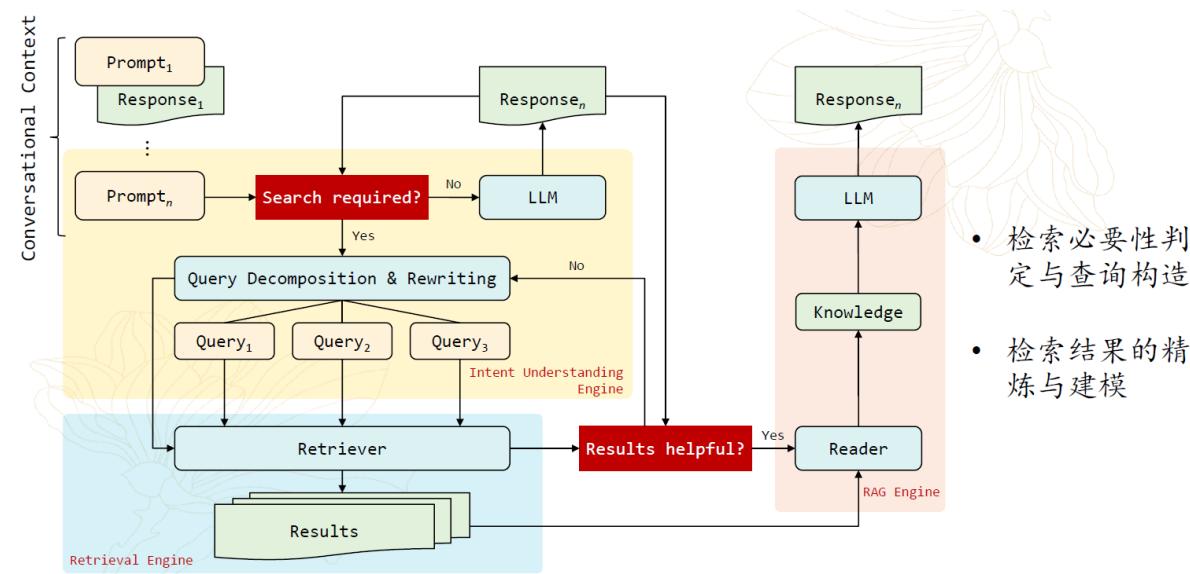
- 大模型并不完美 幻觉问题 知识缺陷 时效性问题

未应用检索增强的大模型（左）笼统的套话+乱说，应用检索增强的大模型（右图）能根据查询到的文档来给出问题的答案



2

### 2. RAG的基本框架



4

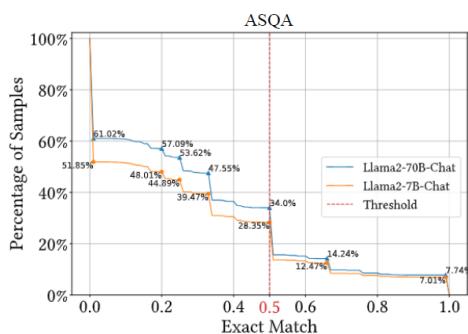
### 3. 何时需要检索？--检索的必要性判定

#### ACL24 SlimPLM 代理模型判定检索的必要性

Small Models, Big Insights: Leveraging Slim Proxy Models To Decide When and What to Retrieve  
for LLMs 2024 ACL

- 检索一定能增强大模型的生成效果么？
  - 无关结果会带来负面的影响
  - 大模型能够掌握的知识不需要检索

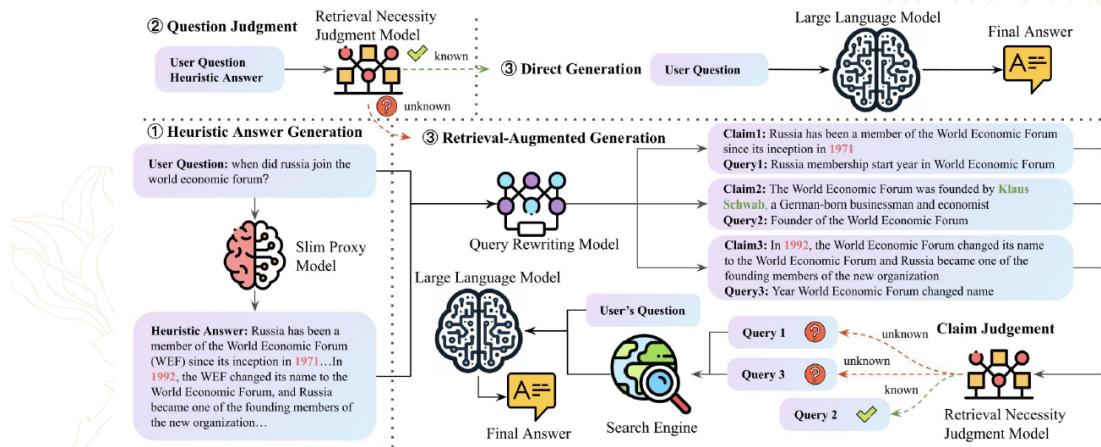
- 先前方法：先生成判定结果后在检索生成，这种方式成本高
- 实验发现在模型置信度较高时无论模型规模大小模型预测结果相似
  - 大小模型是否能在知识掌握程度上呈现相似性呢？



- 回答质量超过当前EM值的样本占总体样本的比例，越高说明当前模型能力越强
- $>0.5$ 两条线迅速接近
- 两个模型 $EM>0.5$ 以上的样本超过80%是一致的
- 目前的LLM用了类似的训练语料
  - 不同LLM在掌握程度高的知识是相当重合的
  - 大小LLM知识能力的差距主要体现在长尾知识上

✓ 用较小参数的语言模型（代理模型）的表现来预测LLM需要做检索的时机！

- 提出使用较小的语言模型作为代理模型，根据代理模型的表现来判定需要做检索的时机
- 代理模型（不做微调）：生成一个回答（Heuristic Answer）
- 检索必要性判断模型（微调后的LLaMA-7b）：决定当前问题是否检索、某个子query是否检索
- 查询重写模型（微调后的LLaMA-7b）：在当前问题需要检索的前提下，生成若干个子query



## 4. 检索结果如何使用？--精炼结果的方法

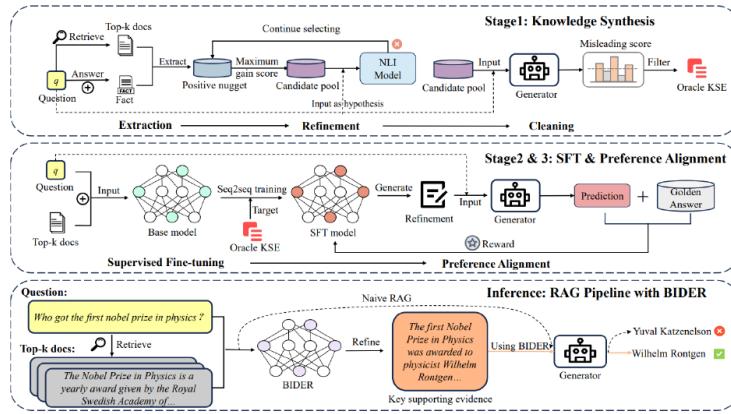
### ACL24 BIDER 为大模型精炼有效知识

BIDER: Bridging Knowledge Inconsistency for Efficient Retrieval-Augmented LLMs via Key Supporting Evidence ACL 2024

- 造成RAG下降的原因：大语言模型与检索系统之间的知识不一致
  - 检索文档往往冗长且含有噪声
  - LLM无法感知自身的知识边界
- 现有方法过度依赖外部知识（检索）或内部知识（LLM）的一部分而忽略了两者之间的联系
  - 基于困惑度的方法：保留LLM认为有高困惑度的词语
  - 基于模型的方法：保留检索系统判断与正确答案最接近的句子
- 构建模型对检索文档进行精炼，提供LLM最需要的知识

1. 构造训练数据
  - 根据正确事实的相似度
  - 用NLI模型判断对答案的支持程度
  - 剔除对LLM有害的知识
2. 有监督微调 (SFT)
  - 学习检索文档到构造的训练数据之间的映射
3. 偏好对齐 (RL)
  - 根据模型偏好进一步精调模型
  - 考虑LLM的内部知识

目标：构建模型对检索文档进行精炼，提供LLM在生成时必需的知识

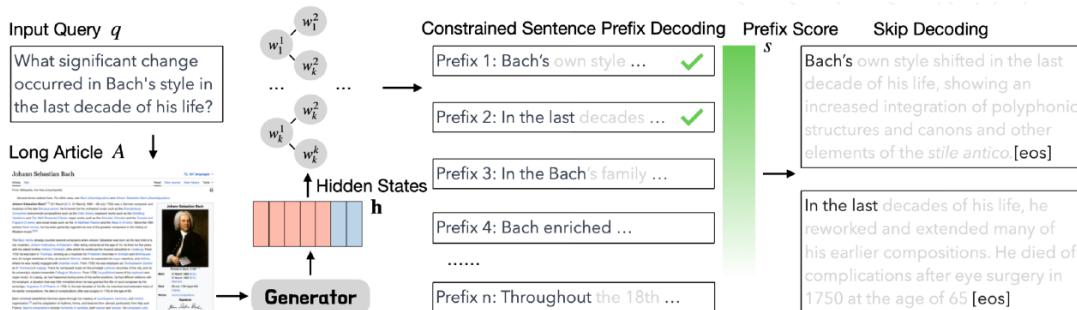


## 5. 较长的结果如何建模？--无切块长文本建模方法

### ACL24 CFIC 建模更长的检索结果

Grounding Language Model with Chunking-Free In-Context Retrieval ACL 2024

- 检索结果通常比较长，传统RAG处理长上下文时存在局限性
  - LLM无法处理超长文本
  - 长上下文中存在大量的不相关内容
- 已有方法
  - 提高上下文窗口大小直接处理
  - 将长上下文切块和排序，寻找相关内容
- 问题
  - 直接处理长上下文算力要求高，且无法消除长上下文中的噪音
  - 切块破坏了语义的连贯性，导致信息缺失
- 提出一种高效的上下文提炼的方法，不破坏语意连贯性且能高效找到支持回复生成的文本证据



- 不切块，直接基于完整的长文档生成支持回复的文本证据
- 如何保证生成的文本证据是准确的？
  - 通过SFT增强模型处理长上下文的能力，能够更好地适应特定的检索任务
  - 使用受限句子前缀解码，将生成空间限制在句子前缀上，确保生成的文本证据与原始上下文紧密相关
- 如何保证效率？
  - CFIC使用“小而精”的模型处理长上下文，为“大而广”的模型提供准确上下文，生成最终答案
  - CFIC采用跳跃解码策略，一旦确定句子前缀，就跳过中间解码步骤，直接选择最有可能结束的位置，提高解码效率

## 6. 工具包

### arXiv24.5.24 FlashRAG 快速实现RAG方法工具包

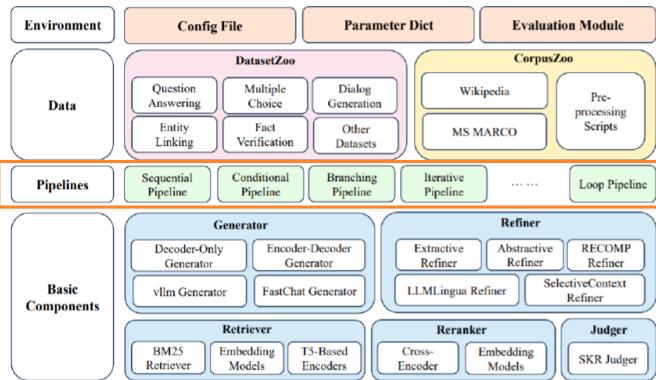
FlashRAG: A Modular Toolkit for Efficient Retrieval-Augmented Generation Research

动机：

- RAG系统组件众多，研究人员要花费大量时间在各类工程的实现上
- 现有的RAG工作缺少统一的实现框架，导致复现非常耗时且难以公平比较
- 已有的LangChain LlamIndex工具封装复杂，难以满足定制化研究需求

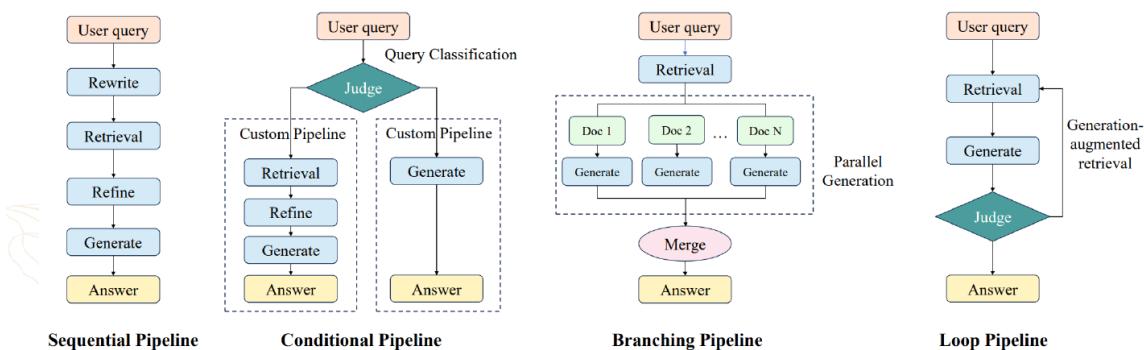
特点：

- 模块化RAG框架，包含检索器、生成器、精炼器等多种组件，支持自定义RAG流程
- 目前实现12种RAG工作，轻松在不同设置下评估结果
- 包含32个RAG工作中常用数据集，并预处理为统一的格式
- 包含多种辅助脚本，包含Wikipedia预处理分块、索引构建、检索结果预准备等
  - 基于基础组件构建Pipeline实现常用RAG流程
  - 涵盖32种常用数据集
  - 大多数基于维基百科作为知识源



Task	Dataset Name	Knowledge Source	# Train	# Dev	# Test
NQ [19]	Wiki	79,168	8,757	3,610	/
TriviaQA [19]	Wiki & Web	78,785	8,837	11,313	/
PopQA [40]	Wiki	-	-	1,267	/
SQuAD [41]	Wiki	87,599	10,570	/	/
MSMARCO-QA [42]	Web	808,731	101,093	/	/
NarrativeQA [43]	Books, movie scripts	32,747	3,461	10,557	/
WikiQA [44]	Wiki	20,349	2,733	6,165	/
WebQuestions [45]	Google Freebase	3,778	/	2,032	/
AmbigQA [46, 38]	Wiki	10,036	2,002	/	/
SiQA [47]	-	33,410	1,954	/	/
CommonsenseQA [48]	-	9,741	1,231	/	/
BoolQ [49]	Wiki	9,427	3,270	/	/
PIQA [50]	-	16,113	1,838	/	/
Fermi [51]	Wiki	8,000	1,000	1,000	1,000
HopQA [52]	Wiki	90,447	7,405	/	/
2023 HopoQA [53]	Wiki	15,000	12,576	/	/
Musique [54]	Wiki	19,938	2,417	/	/
Bamboohge [32]	Wiki	-	-	125	/
Long-Form QA	Wiki	4,353	948	/	/
ELIS [56]	Reddit	272,634	1,507	/	/
MMLU [35, 36]	-	99,842	1,531	14,042	/
TriviaQA [57]	Wiki	84	84	/	/
HellaSwag [58]	ActivityNet	39,905	10,042	/	/
ARC [59]	-	3,370	869	3,548	/
OpenBookQA [37]	-	4,957	500	500	500
Entity-linking	Wiki & Freebase	18,395	4,780	/	/
WN192 [60]	-	8,995	8,995	/	/
T-REx [63, 61]	DBpedia	2,284,168	5,000	/	/
Slot filling	Wiki	147,909	3,724	/	/
Fact Verification	Wiki	104,966	10,444	/	/
Dialog Generation	Wiki	63,734	3,054	/	/
Open-domain Summarization*	WikiAsg [67]	300,636	37,046	37,368	/

- 目前已支持常见的四种不同流水线的RAG工作



## 7. 数据集

### WebBrain 面向RAG的通用数据集

WebBrain: Learning to Generate Factually Correct Articles for Queries by Grounding on Large Web Corpus arXiv2304

- 现有的RAG数据集，尤其是训练集不足
  - 已有工作多采用open-domain QA 作为训练和测试集
  - 人为使用检索器构建训练集合，检索和生成文本的关联性缺乏保障
  - 难以判断是否参考了检索结果
- 提出基于维基百科的文本及其引用构建大规模数据集
  - 包括对维基百科引用链接进行标注
  - 超大规模原始数据集（2.8TB）
    - 支撑预训练、微调等多种研究
  - 高质量数据源（Wikipedia）
  - 包含引用标号，可解释性良好
  - 已清洗、抽取并划分检索与生成集合，支持多种研究

	WebBrain-R	WebBrain-G
# Queries	2.74M	12.32M
# Ref. passages	3.20M	12.61M
# Tokens / Query	3.2	2.9
# Tokens / Passage	237.5	250.0
# Tokens / Target	-	108.6
# Training	4.46M	12.30M
# Validation	0.2M	0.5M
# Test	88,935	24,546

Dataset	# Wiki Pages	# Refs	Status	Storage Size
WikiSum (Liu et al., 2018)	2.3M	87M	Need crawling	300GB
WikiCatSum (Perez-Beltrachini et al., 2019)	0.17M	23.5M	Ready	4.8GB
Hiersumm (Liu & Lapata, 2019)	1.66M	-	Ready	6.9GB
WebBrain-Raw	14.86M	259.5M	Ready	2.8TB

### DomainRAG 特定领域RAG评测

DomainRAG: A Chinese Benchmark for Evaluating Domain-specific Retrieval-Augmented Generation ACL 2024.

动机：

- RAG能够有效解决LLM的各种限制，例如幻觉和知识实时更新的困难
- 目前的研究往往依赖于维基百科等一般知识源来评估模型解决常识性问题的能力，然而RAG在LLM难以涵盖专业知识的场景和特定领域中的应用也很重要

方法：

- 使用特定领域的语料库和问题对于评估LLM有效利用来自这些特定领域的外部知识来解决专家问题的能力至关重要
- 总结综合评价RAG模型的六个重要能力，并以人大招生为应用场景构建了评估这些能力的数据集



## 8. Future Work

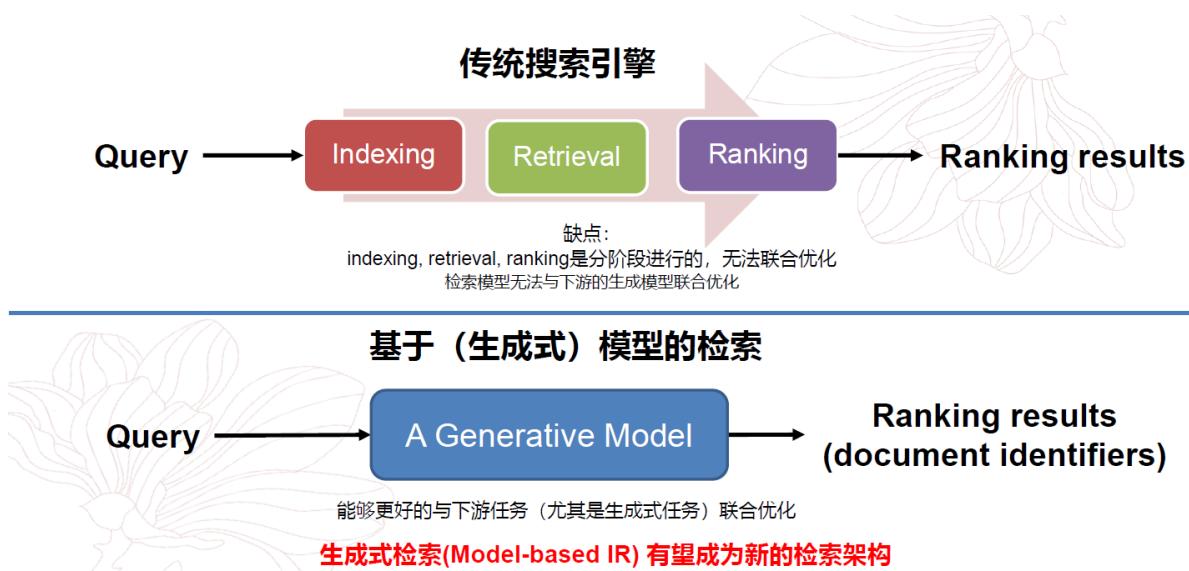
- 更加精准的查询分解与改写
- 对话式RAG的进一步探索
- 面向RAG的训练（预训练、指令微调）
- 长窗口与RAG之间的关联
- RAG系统的评估方法

## Part 3 生成式文档检索

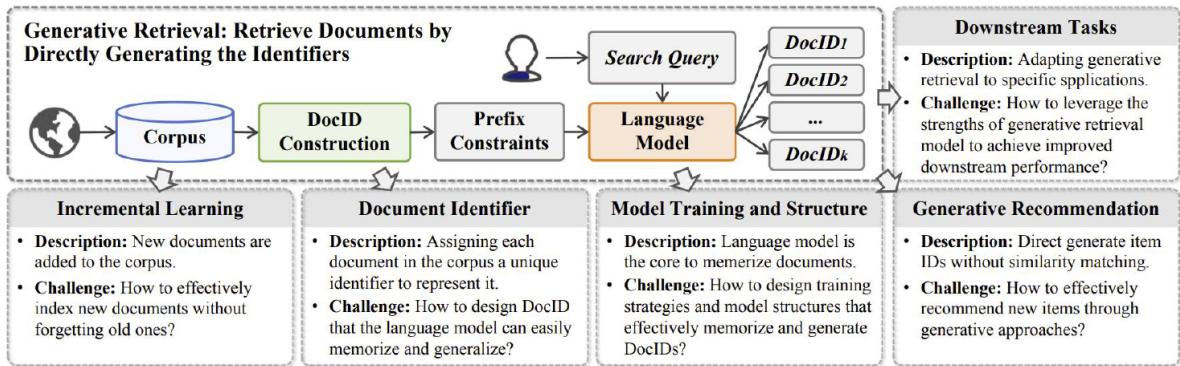
From Matching to Generation: A Survey on Generative Information Retrieval arXiv2404

一个核心的问题：能否直接通过大模型完成文档的检索/召回？

- 大模型并没有检索能力
- 大模型瞎编烂造的能力在检索相关问题上体现的淋漓尽致
- 大模型需要定向微调才能实现检索能力



## 1. 生成式检索模型目前面临的问题



### 增量学习问题

- 文档的动态更新，大模型怎么去适配？
- 如何处理海量的文档？
- 如何将文档嵌入到模型中？

### 文档标识定义

- 如何定义DocID能够让模型更轻松的记忆和泛化

### 训练策略和模型架构

- 如何设计架构和策略来让模型更高效的记忆和泛化海量文档

### 生成答案

- 如何通过查询到的文档高效的生成答案

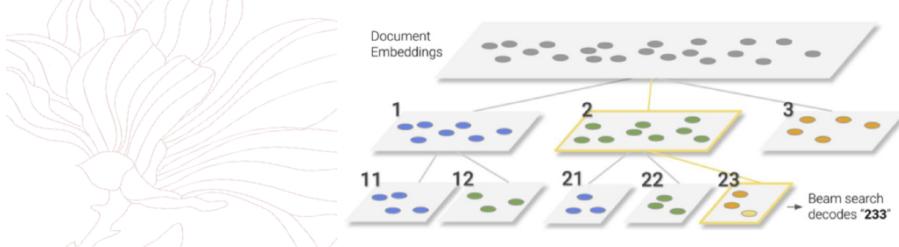
## 2. 经典工作

### DSI (Google)

提出一种分层的文档编码方案

Transformer Memory as a Differentiable Search Index NeurIPS 2022

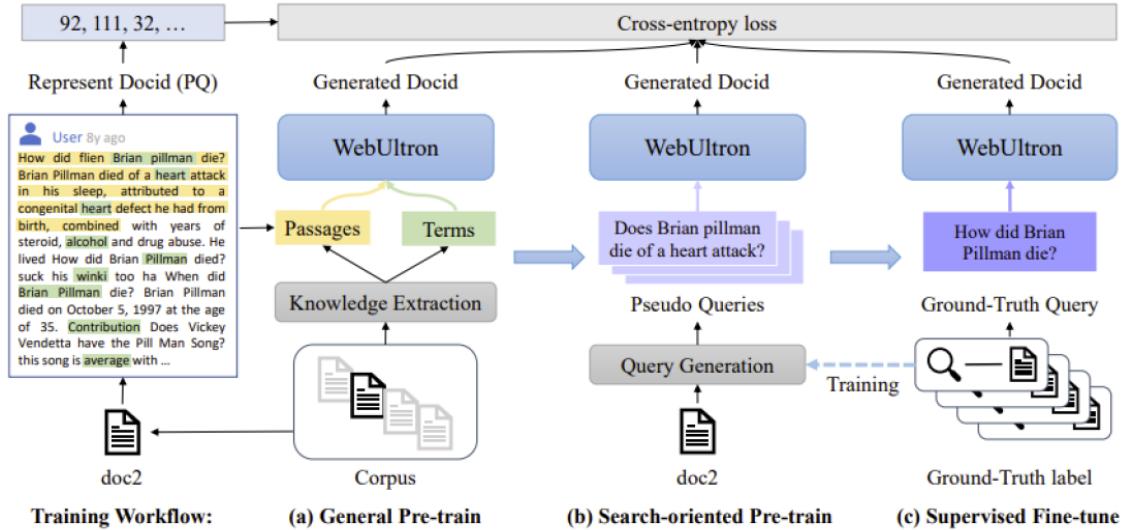
- **Indexing Method**  
 $\text{doc\_tokens} \rightarrow \text{docid}$      $\text{docid} \rightarrow \text{doc\_tokens}$
- **What to index:**
  - First L tokens; set of terms; a single contiguous chunk of k tokens
- **DocId:** Atomic Identifiers; naively structured identifiers; Semantically Structured Identifiers



## WebUltron (renda)

给出一种三阶段训练框架

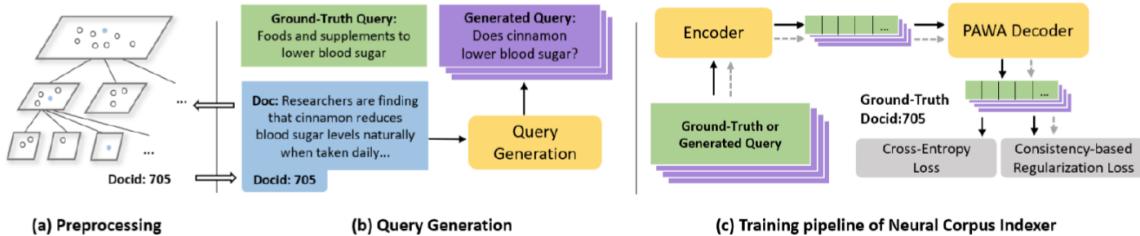
WebUltron: An Ultimate Retriever on Webpages Under the Model-Centric Paradigm 2023 IEEE Transactions on Knowledge and Data Engineering



## NCI (MicroSoft)

A Neural Corpus Indexer for Document Retrieval NeurIPS 2022

提出一种神经语料库检索器，序列到序列的网络



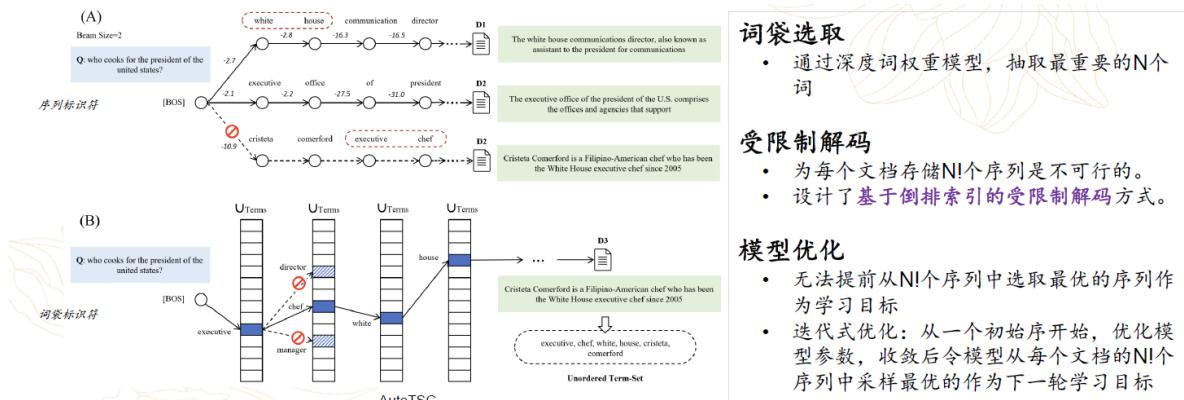
**Figure 1:** Overview of Neural Corpus Indexer (NCI). (a) Preprocessing. Each document is represented by a semantic identifier via hierarchical  $k$ -means. (b) Query Generation. Queries are generated for each document based on the content. (c) The training pipeline of NCI. The model is trained over augmented  $\langle query, docid \rangle$  pairs through a standard transformer encoder and the proposed Prefix-Aware Weight-Adaptive (PAWA) Decoder.

## 跳出Sequence范式，词袋模型 (renda)

Generative Retrieval via Term Set Generation SIGIR 2024

文档ID是序列化的形式，解码错一步则全错

提出基于词袋的方案，生成词袋的顺序构成了文档标识符



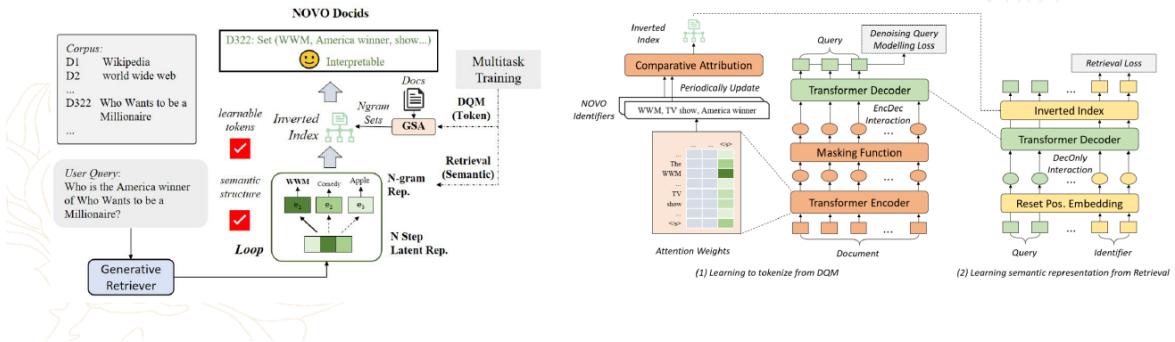
## 可学习的文档标识符 (renda)

NOVO: Learnable and Interpretable Document Identifiers for Model-Based IR CIKM 2023

现有的文档ID基于Encoder独立完成，与Decoder无关，存在Gap，提出了一种可学习的文档标识符方案

### 可学习的文档标识符NOVO

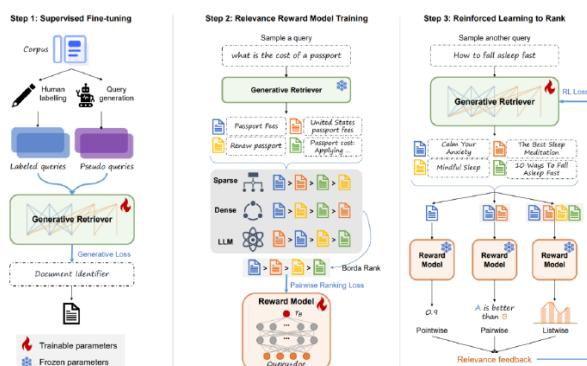
- 由N元组集合组成，随模型生成过程，通过倒排索引召回相应文档
- N元组通过全局自注意力(GSA)信息从文档中获取，并通过两种任务(DQM, Retrieval)优化标识符令牌及语义，并与检索任务对偶优化



### 相关性强化的生成式检索模型 (renda)

Enhancing Generative Retrieval with Reinforcement Learning from Relevance Feedback EMNLP 2023

引入基于相关性反馈的强化学习来让模型理解相关性



### 生成式检索与其他生成任务的融合 (renda)

UniGen: A Unified Generative Framework for Retrieval and Question Answering with Large Language Models AAAI 2024

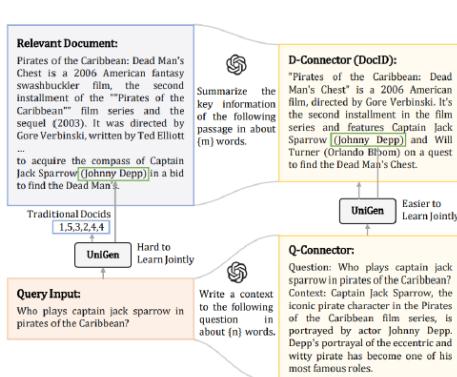


Figure 3: An example of generating LLM-based connectors from the query side and document side, with the labeled answer highlighted in the green box.

$$\mathcal{L}_{\text{retr}} = - \sum_{i=1}^t \log f_{\text{retr}}(d_i | d'_{<i}, q; \theta, \phi). \quad (2)$$

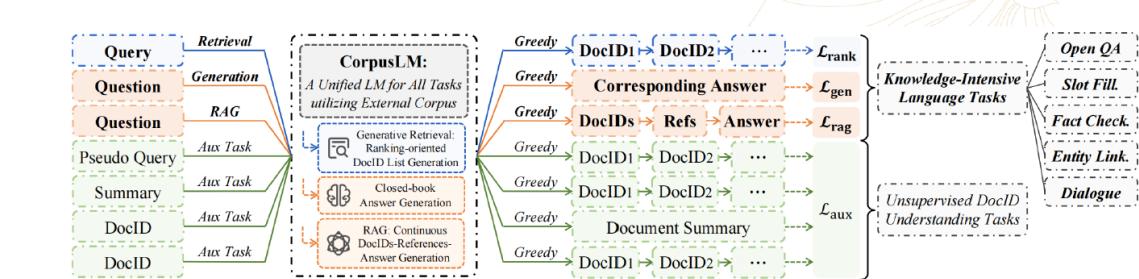
$$\mathcal{L}_{\text{qa}} = - \sum_{i=1}^{T'} \log f_{\text{qa}}(a_i | a_{<i}, q; \theta, \mu). \quad (4)$$

$$\mathcal{L}'_{\text{retr}} = - \sum_i \log f_{\text{retr}}(d_c | d_{c_{<i}}, q_c; \theta, \phi), \quad (5)$$

$$\mathcal{L}'_{\text{qa}} = - \sum_i \log f_{\text{qa}}(a_i | a_{<i}, q_c; \theta, \mu), \quad (6)$$

$$\mathcal{L} = \lambda \mathcal{L}'_{\text{retr}} + (1 - \lambda) \mathcal{L}'_{\text{qa}}, \quad (7)$$

## 挑战：生成式文档检索和大模型如何融合？



- 统一的语言模型：集成生成式检索、闭卷生成和检索增强生成（RAG），辅助知识密集型任务。
- 面向排名的DocID列表生成策略：贪婪解码生成Doc ID排序。
- RAG导向的连续生成策略：连续解码方法。
- 无监督DocID理解任务：无监督的DocID理解任务，加深模型对DocIDs背后含义的理解，以进一步提高了CorpusLM的检索和生成性能。

## Report 3 大模型时代的通用向量检索

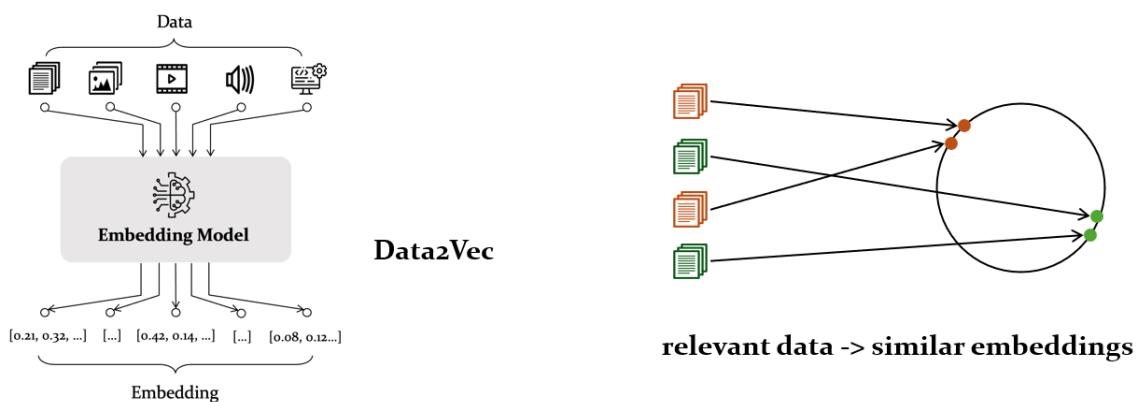
北京智源研究院 刘政

### 01 什么是语义向量模型

向量模型：将任意数据转化为高维空间中稠密向量的计算机模型

重要属性：**向量相似性要与数据相似性保持一致**

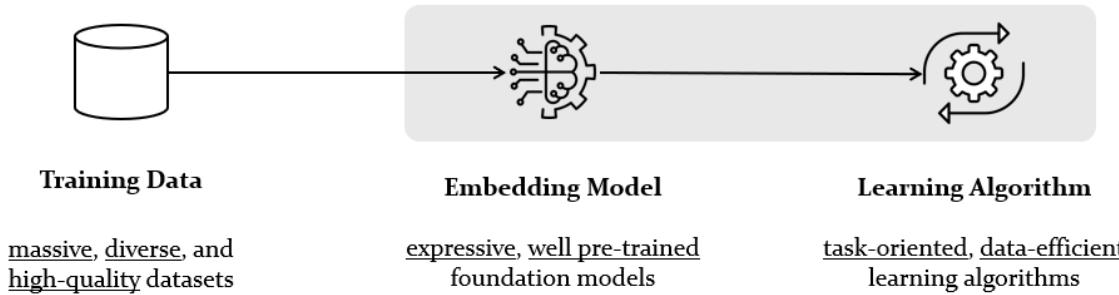
这里的相似性计算并不严格，不受三角不等式约束



向量模型应用：

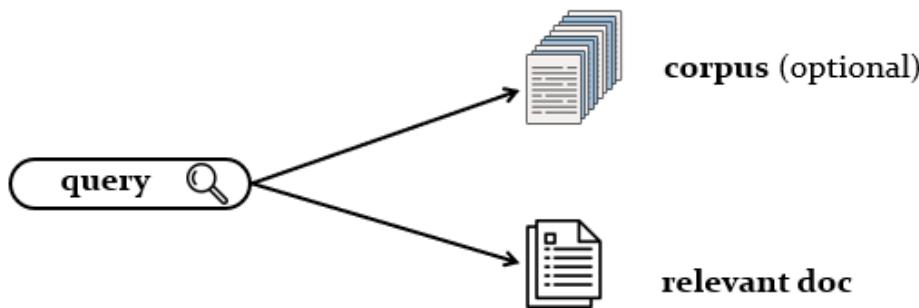
- 信息检索系统
- 比较数据的语义关联，数据聚类去重等
- 向量数据库：不再追求相似性最高的目标向量，近似最近邻计算节约开销

## 02 向量学习的基本模式



### 1. 训练数据

向量模型的训练数据通常包含三个组成单元：查询、查询对应的关联文档以及全部的文档集合



#### Learn to discriminate each query's relevant doc from the rest of corpus

向量模型的训练数据是一种相对稀缺的资源，最常用的是微软研究院发布的MS MARCO数据集，这是发布最早、规模最大、质量最好的数据集之一。

- **Popular datasets:** MS MARCO, Natural Question, Hotpot-QA, Mr.TyDi/MIRACL

**MS MARCO**

Will I qualify for OSAP if I'm new in Canada?

Selected Passages from Bing

"Will the OSAP website for application deadlines. To get OSAP you have to be eligible. You can apply using an online form, or you can print off the application forms. If you submit a paper application, you must pay an application fee. The online application is free."

Source: <http://settlement.org/ontario/education/collleges-universities-and-institutes/financial-assistance-for-post-secondary-education/how-do-i-apply-for-the-ontario-student-assistance-program-osap/>

"To be eligible to apply for financial assistance from the Ontario Student Assistance Program you must be: 1) a Canadian; 2) a permanent resident; or 3) Protected person/convention refugee with a Protected Persons Status Document (PPSD)." Source: <http://settlement.org/ontario/education/collleges-universities-and-institutes/financial-assistance-for-post-secondary-education/who-is-eligible-for-the-ontario-student-assistance-program-osap/>

"You will not be eligible for a Canada-Ontario Integrated Student Loan, but can apply for a part-time loan through the Canada Student Loans program. There are also grants, bursaries and scholarships available for both full-time and part-time students."

Source: <http://campusaccess.com/financial-aid/osap.html>

Answer  
No. You won't qualify.

**Created by:**  
Microsoft Research

**Overview:**  
**Query:** from BING QA  
**Doc:** clicked webpage  
**Answer:** human label  
**Size:** 300K+ query, 300M+ doc

发布最早 规模最大 质量最好

1. Question answering data is hard to create
2. Web Search companies (Microsoft/Google/Baidu) play key roles
3. Modern IR search is made possible thanks to these public datasets

**Natural Question**  
Created by: Google Research  
Feature: Wikipedia doc

**Hotpot-QA**  
Created by: Google Research  
Feature: Multi-hop QA

**Mr.TyDi/MIRACL**  
Created by: Google/Waterloo  
Feature: Multi-lingual QA

**DuReader**  
Created by: Baidu  
Feature: Chinese QA

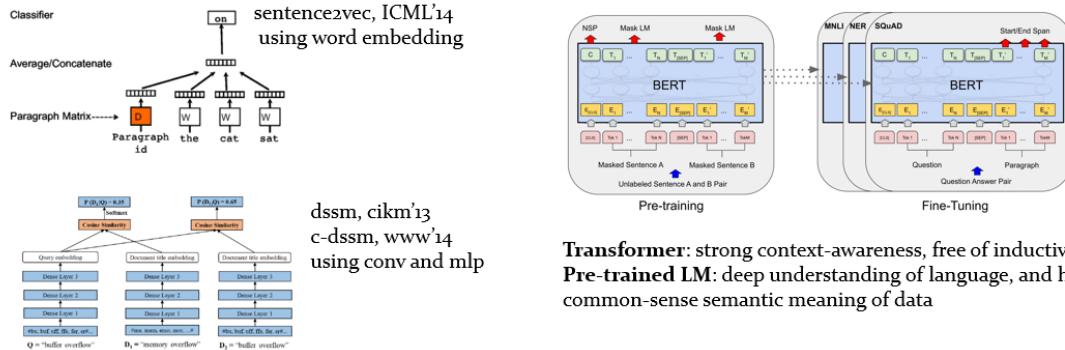
针对学习数据稀缺的问题，主要原因是传统的向量学习高度依赖人类的标注，大规模扩增不可能。目前一种解决方案是利用机器自动化的生产制造数据

- UltraFeedback (THU)
- Alpaca (Stanford)

## 2. 模型

### Foundation model for embedding

- Almost any DNNs can be naturally applied as the embedding model (Word2Vec, LSTM, etc.)
- Not until the arrival of pre-trained LMs (e.g., BERT) that embedding model really works
- Advantage of PLMs: 1) transformer architecture, 2) pre-training



- 几乎所有的DNN模型都可以作为向量模型
- 但真正让向量模型变得可用的是，以Transformer为基础的预训练模型，其本质是“大”
- Transformer的网络架构可以做的足够大，让模型由足够的容量去建立强大的表达能力
- 预训练的规模可以足够大，让模型充分学习海量数据中的知识

### 主流的预训练算法对向量检索而言是不是最优的？

- 修改模型训练目标：让embedding直接体现在训练的优化之中
- 优化embedding对输入文本语义的表达能力：先前的单个词的预测不合适，简单的转为预测目标文本的全部词汇

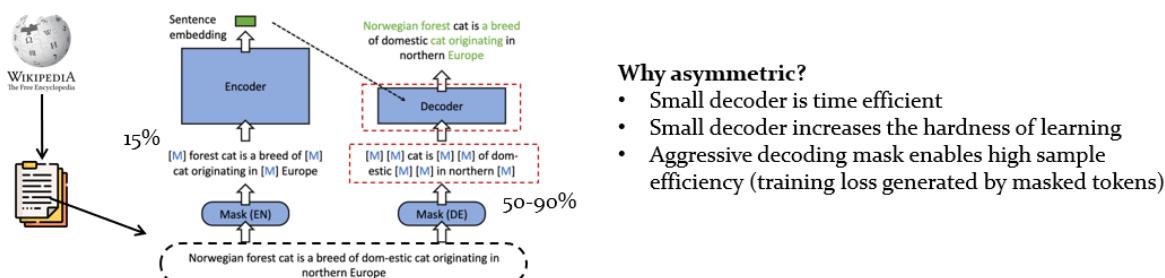
基于此：提出RetroMAE

RetroMAE: Pre-Training Retrieval-oriented Language Models Via Masked Auto-Encoder EMNLP 2022

- 从任意的无标签语料，比如Wikipedia中采样一段文本，将其编码成embedding之后再将其解码出原始的文本
- 非对称的编码器解码器结构，非对称的掩码比率，使得模型输出的embedding可以更加充分的从训练数据中得到学习优化

### The design of RetroMAE (Retrieval-Oriented Pre-training Based Mask Auto-Encoder)

- High-level framework:** 1) sample a sentence, 2) mask and encode, 3) decode the clean input
- Concrete process:** 1) asymmetric structure, full-scale encoder (12L), small-scale decoder (1L), 2) asymmetric masks, encoder: 15%, decoder: 50–90%
- Data:** purely based on unlabeled data (Wikipedia corpus), same as general pre-training



未来持续增大预训练模型的规模是一个必然的趋势，性能随着模型规模一致提升

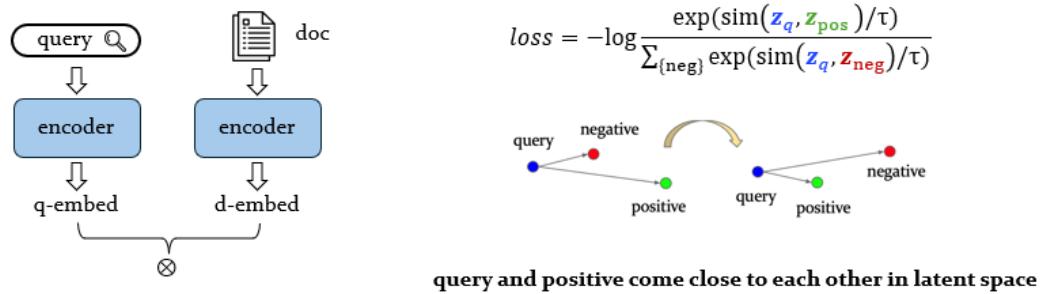
### 3. 训练算法

讲者观点是训练数据和模型的重要性要高于训练算法，尽管训练算法也同样重要

向量学习训练算法的基本形式是对比学习，我们可以进行优化的自由度只有两个：正样本和负样本

#### Last but not least, learning method

- General form: constrative learning
- Pipeline: training samples -> embedding -> contrastive loss
- Discriminate positive samples from negatives for each query based on embedding similarity
- Two key factors: positive samples, negative samples



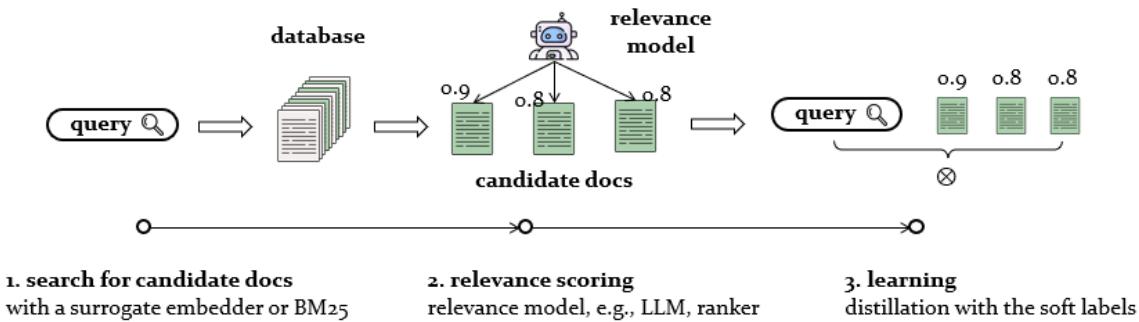
#### Positive Sample

如何找到正样本？

借助一个已知的检索模型，找到关联性的候选文档，并对其进行细粒度的关联性标注，借助蒸馏技术对向量模型进行精细化的训练

#### Positive samples

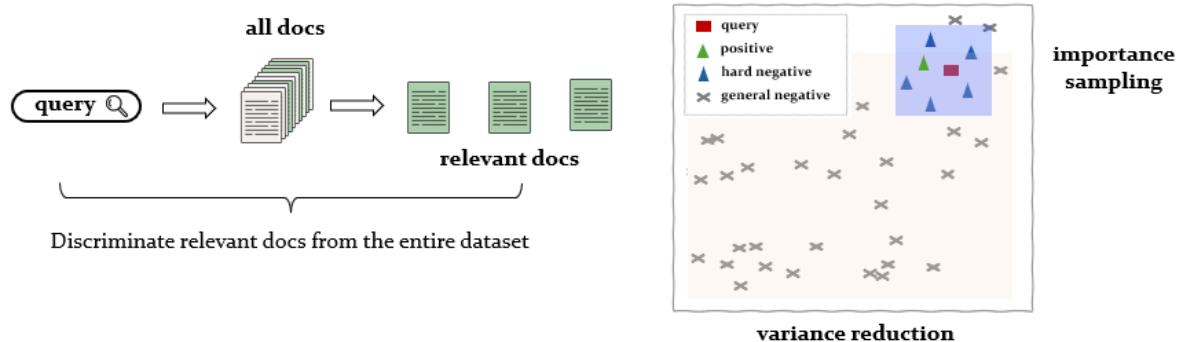
- Positive samples are expected to be **abundant**, **diverse**, and **high-quality**
- Usually, positive samples are expensive and well-presented by the training set
- Fortunately, it's still to leverage machines for data augmentation
- Common approach: query -> candidate docs -> relevance score (as soft labels)



## Negative Sample

### Negative samples

- Negative samples are expected to be **abundant** and sufficiently "**hard**"
- Embedding model is to discriminate positives from the entire negatives, so ideally, hard negative samples should be the entire corpus except the positive (corpus/positive)
- **Abundant negatives:** variance reduction, **hard negatives:** importance sampling

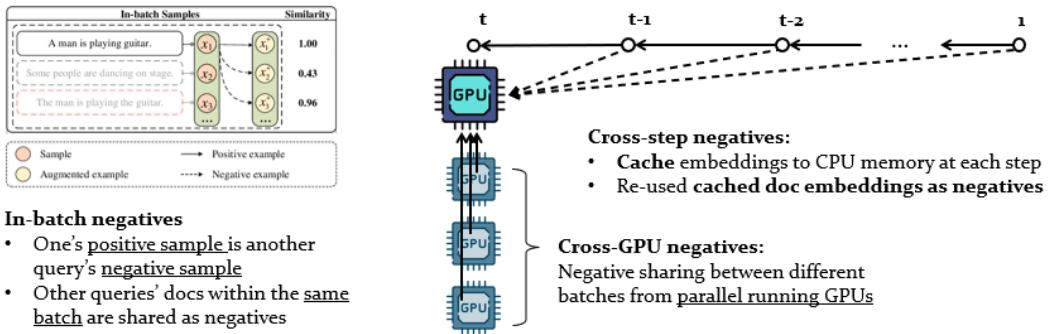


我们希望正样本可以从所有负样本中轻松分辨出来，除了正样本之外所有的文档都为负样本，但是硬件并不允许我们这样做。提出两点：

- 尽可能多的引入负样本，负样本越多，再损失上更接近与无偏
  - 批次间共享 (in-batch negative sampling) 不同设备 (cross-device negative sampling)、不同训练step (gradient-checkpoint based sampling)之间共享
  - 同样要付出更大的训练集群和更长的训练时间的代价

### Increasing negative sampling

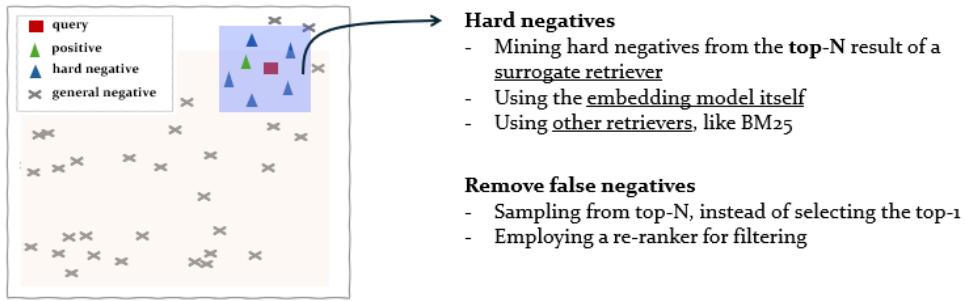
- Naïve increasing is expensive and infeasible
- Fundamental principle: **sharing**
- Three strategies: **in-batch** sharing, **cross-device** sharing, **gradient-checkpoint**
- However, **no free lunch** (higher resource consumption, higher computation cost)



- 挖掘更难的负样本，从统计学的角度来看这相当于重要性采样；从损失的角度来看，更难的负样本会带来更大的梯度更新
  - 检索模型查询到与query关联度高的候选文档，都可以认为是难负样本

### Mining hard negatives

- Hard negatives are seemingly relevant to the query
- Using a **surrogate retriever** to search similar docs to each query (embedding, BM25)
- Removal false negatives? 1. Sampling, 2. Re-ranker filtering



## 03 BGE模型的开发与实践

<https://github.com/FlagOpen/FlagEmbedding/tree/master>

B (BAAI) G (General) E (Embedding)

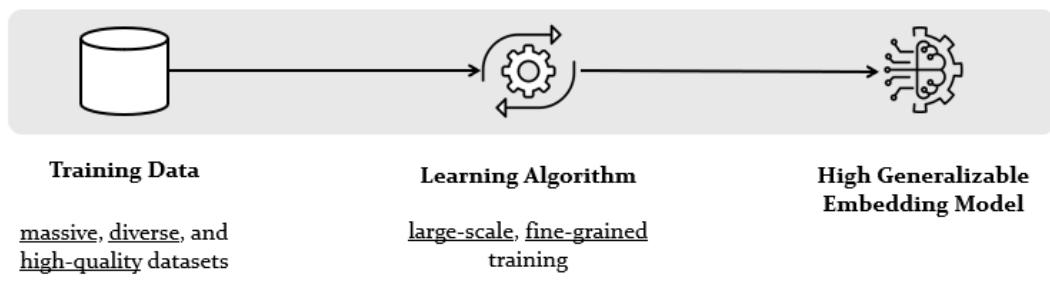
- 我们希望模型足够通用，能够在不同场景下具有适用性且性能出色
- 向量模型的通用性成为了RAG的一个瓶颈
- 我们认为**通用能力**是未来的一个主要的优化目标

通用性的维度：

- 任务场景：服务任意的应用，无论其涉及的知识理论是什么
- 语言：支持任意语言的检索诉求，实现任意语言之间的语义胡同
- 数据模态：文本、图像、语音、甚至是分子结构

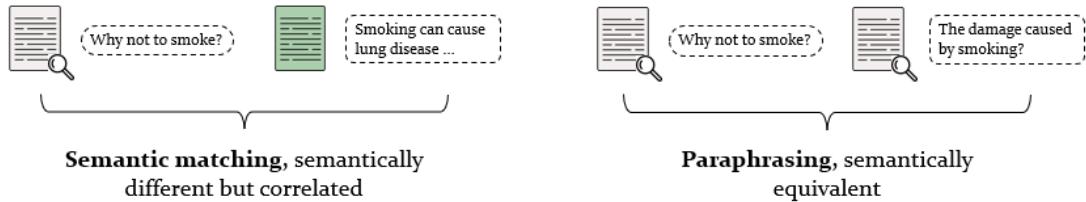
最终目标是同时实现这三个维度，这是一个很困难的问题。

**BGE v1 23.08**

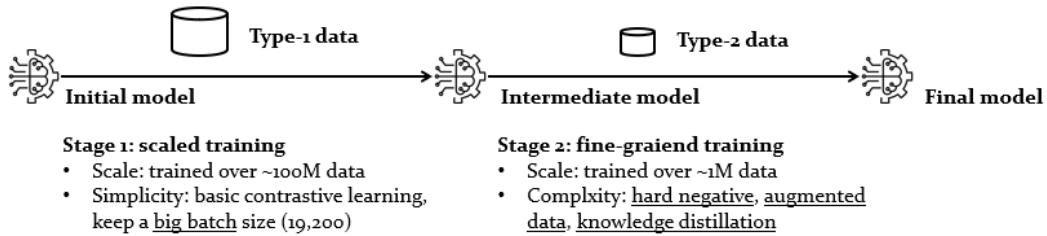


- 数据建设：足够的、多样化的、高质量的训练数据
  - 数据质量相对较低但规模巨大：海量的公开语料，包含大量半结构化信息
  - 数据规模较小，但质量较高，覆盖关键能力：收集了语义关联文本QA数据集；语义等价文本数据集

- Two fundamental capabilities for a general embedding model
  - **Semantic-matching**: semantically different but correlated items, e.g., question-answer
  - **Paraphrasing**: items of the same semantic, e.g., near-duplicate sentences
- Two kinds of labeled datasets: 1) **QA** (MSMARCO, HotpotQA), 2) **Paraphrasing** (NLI, Quora)



- 训练方法：大规模的训练+精细化的训练
  - 第一阶段强调规模化，最基础的对比学习，使得模型初步但全面的建立复杂、多样语义关系的匹配能力
  - 第二阶段强强调精细化，难例挖掘、数据增强、知识蒸馏等，建立模型对关键领域的语义匹配能力
- **Multi-stage training:**
  - Stage 1: basic contrastive learning over Type-1 data, expand batch size a.b.a.p (19,200)
  - Stage 2: fine-grained tuning over Type-2 data: iterative negative mining, knowledge distill



### Improved BGE: M3 Embedding (BGE v2) 24.02

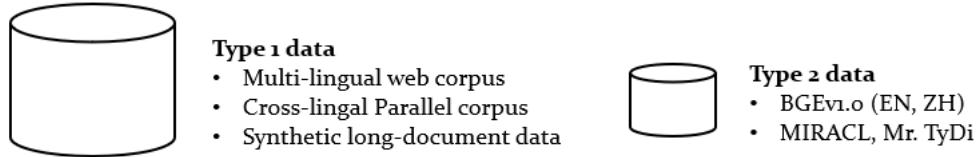
BGE M3-Embedding: Multi-Lingual, Multi-Functionality, Multi-Granularity Text Embeddings Through Self-Knowledge Distillation arXiv2402

- Multi-Lingual: 首要目标是打通语言壁垒，一个模型支撑不同语言的多语言检索能力以及跨语言检索能力
- Multi-Granularity: 获得更长的序列范围 (8192 tokens)
- Multi-Functionality: 希望一个模型集成向量检索、关键词检索重排序模块等功能

数据建设：

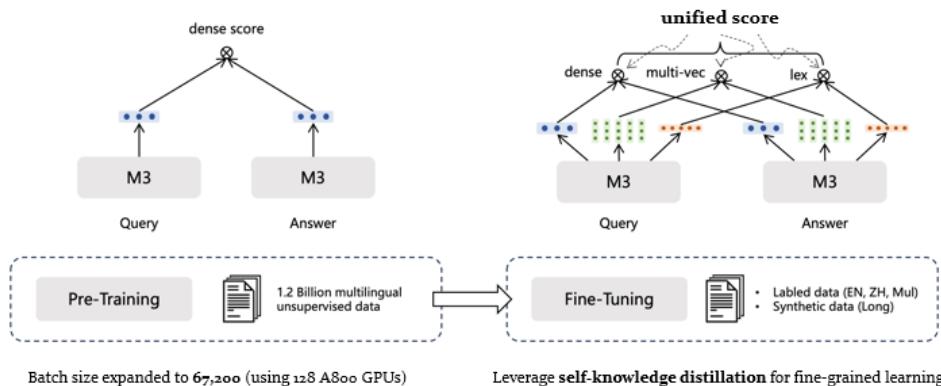
- type 1 : 1.2billion 比BGEv1高出一个数量级，覆盖率100种以上的语言
- type 2 : 语义关联、语义等价

- **Type 1 data**, massive and diverse (**1.2 billion** text pairs in total)
  - Multi-lingual data from unlabeled web corpus, e.g., Common Crawl
  - Cross-lingual data from open parallel corpus, e.g., NLLB (No Language Left Behind), CC-Matrix (translation pairs)
  - Long-doc matching data: synthesized by **ChatGPT** (no existing labeled data is long enough)
- **Type 2 data**, smaller but high-quality, domain-specific
  - BGE-v1 Fine-tuning data (EN, ZH) plus MIRACL and Mr. TyDi (multi-lingual QA) (**within 2M**)



训练：

- 提升算力，训练了高达67k epochs
- 引入自蒸馏技术，集成了多种检索模式
- To ensure the embedding's discriminativeness, batch size is further expanded (with more GPUs)
- Take advantage of **self-knowledge distillation** for fine-grained learning in stage-2



支持多种小语种的语言模型，被评为向量模型中的瑞士军刀

## 04 大语言模型与信息检索

从长远来看，大语言模型势必取代搜索引擎

取代需要解决的两大问题：

- 能够处理的上下文足够长（窗口足够长） 目前拓展长度受到系统硬件制约
- 处理长序列的能力足够强
- 成本控制

（偷听：北京大学目前再做**对于长文本的精细化检索**，实验发现检索结果会随着文档的长度的增加而降低）

### ACL24 CFIC 建模更长的检索结果

Grounding Language Model with Chunking-Free In-Context Retrieval ACL 2024

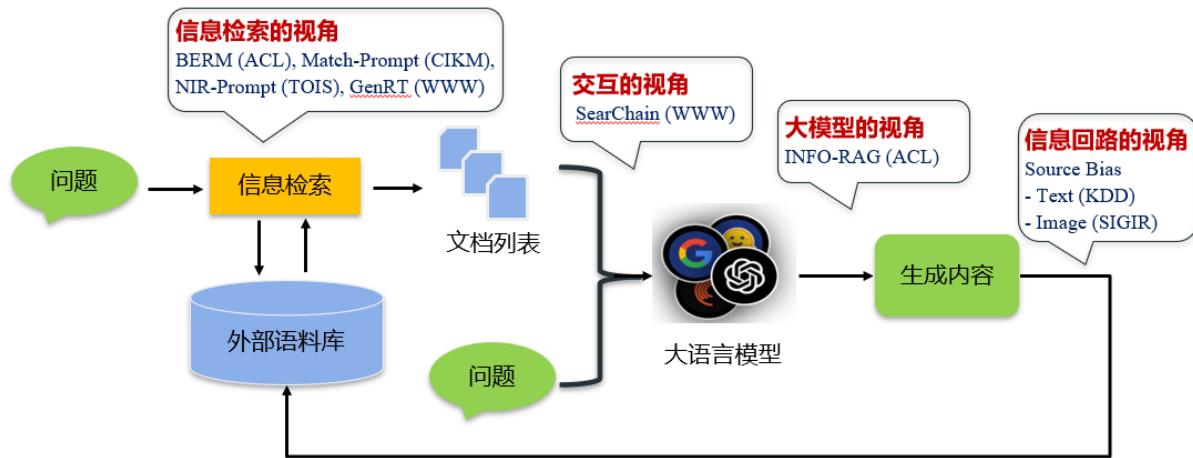
与人大合作的一项工作，见Report 2 Part 2

## 05 总结

- 实现通用的向量模型还有很长的路要走
- 我们认为数据基础跟模型是实现通用模型的核心要素
- 未来生成式检索会变得愈发流形

## Report 4 检索增强大模型技术探索与思考

中国科学院计算技术研究所 庞亮



目前面临的核心问题：

- 不能准确的获得知识（检索视角）
- 不能准确的选择知识（大模型视角）
- 知识之间的干扰（交互视角）

从四个角度介绍计算所近期的工作

### 01 检索视角下的检索增强

什么是适合大语言模型的检索增强的信息检索？

- 应用范围广，任务种类多，对跨领域跨任务泛化性要求高
- 推理开销大，上下文空间有限，对排序精度和鲁棒性要求高

研究现状：

## 检索源头

- ◆ 非结构化数据，如词组，跨语言文本，提示词等 (UPRISE , CREA-ICL)
- ◆ 半结构化数据如表格，PDF (TableGPT)
- ◆ 结构化数据，如实体，三元组，子图等 (KnowledGPT, G-Retriever, SUGRE)
- ◆ 大模型自身的记忆 (GenRead, SKR, Selfmem )

## 查询重构

- ◆ 查询扩展 (多查询扩展，如least-to-most, Chain-of-Verification)
- ◆ 查询改写 (Rewrite-retrieve-read, BEQUE, HyDE, Step-back Prompting)

## 检索器微调

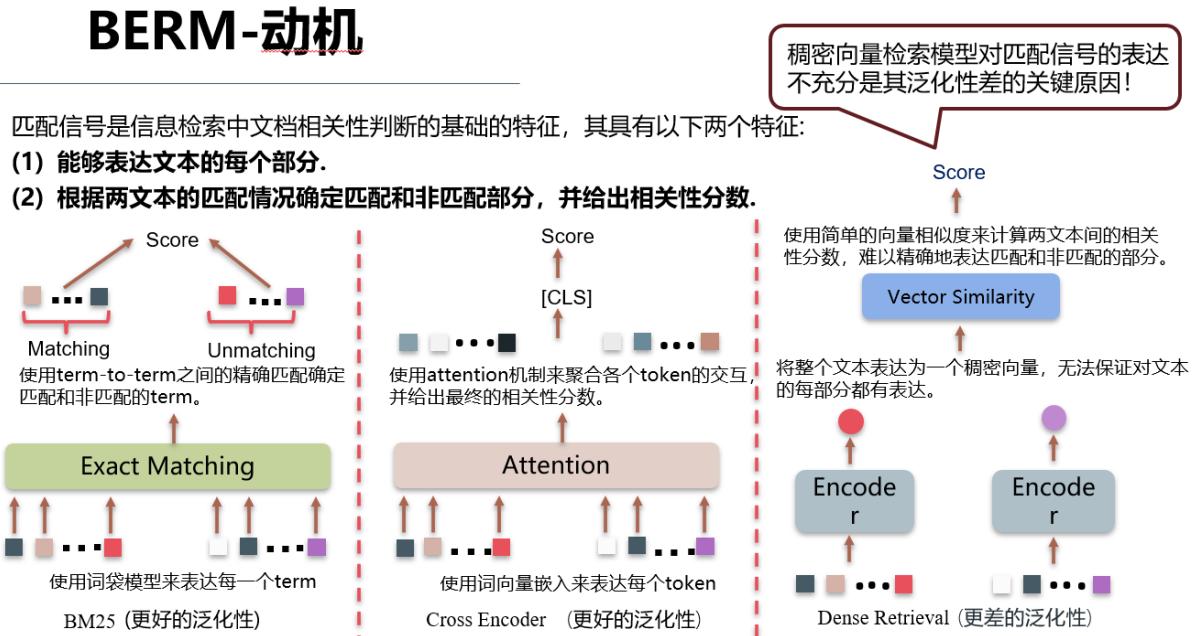
- ◆ 大模型做为数据生成器 (PROMPTAGATOR)
- ◆ 大模型做为监督信号 (REPLUG)
- ◆ 微调检索器中的适应模块以匹配大模型 (Augmentation-Adapted Retriever, PRCA, BGM, PKG)

## BERM：训练匹配的平衡可提取表示提高密集检索的泛化能力

BERM: Training the Balanced and Extractable Representation for Matching to Improve Generalization Ability of Dense Retrieval ACL 2023

- 动机：大部分稠密向量检索算法效果再数据集之外的场景泛化性能非常差

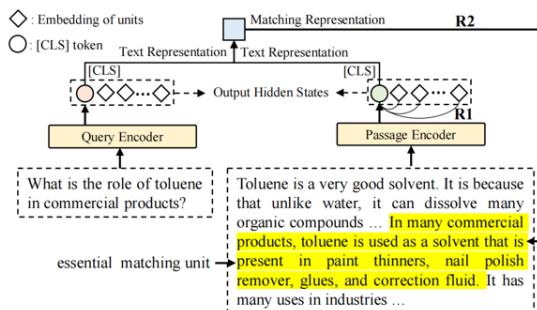
## BERM-动机



- 提出了匹配表示概念
- 提出了可泛化的稠密向量检索模型再训练时的两个要求

提出了匹配表示的概念：

---- 两文本表示在计算向量相似度时的中间状态：向量按位乘



提出了可泛化的稠密向量检索模型在训练时的两个要求：

---- R1: 文本表示的语义单元均匀性

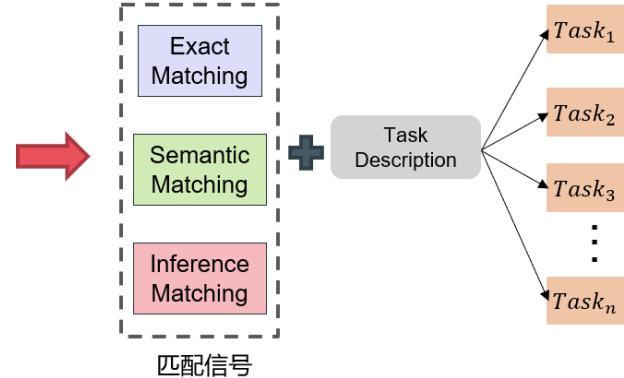
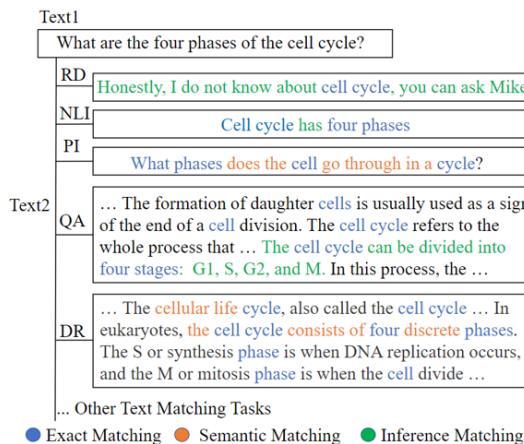
---- R2: 匹配表示的关键匹配单元可抽取性

## Match-Prompt：通过提示学习提高神经文本匹配的多任务泛化能力

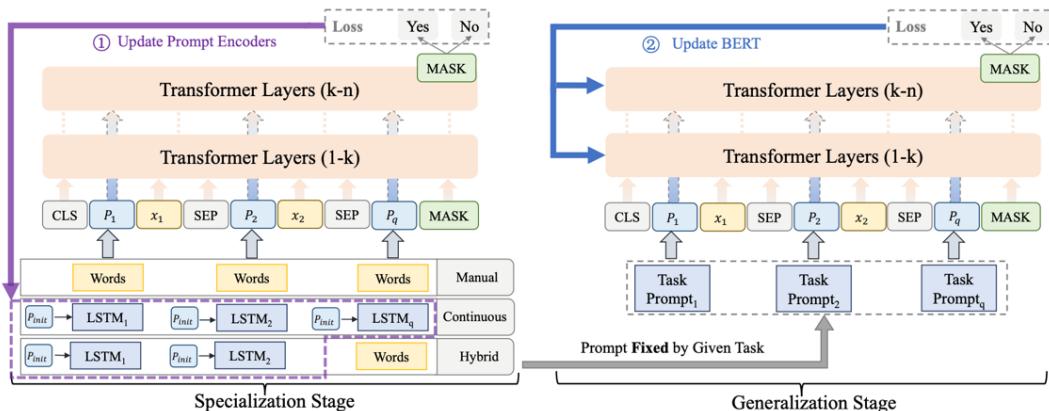
Match-Prompt: Improving Multi-task Generalization Ability for Neural Text Matching via Prompt Learning CIKM 2022

## Match-Prompt 动机

多种文本匹配任务关注存在共同的关键匹配信号，即精确匹配，语义匹配，推理匹配



## Match-Prompt模型

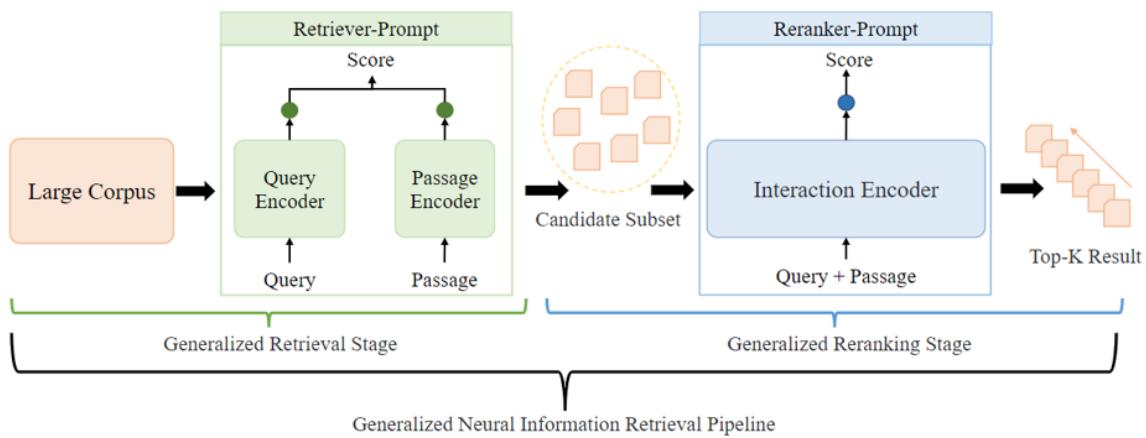


1. 如何捕捉到任务通用的关键匹配信号？
2. 如何将关键匹配信号组合以适应不同任务？

两阶段训练策略：特殊化——泛化

## NIR-Prompt: 多任务可泛化神经信息检索训练框架

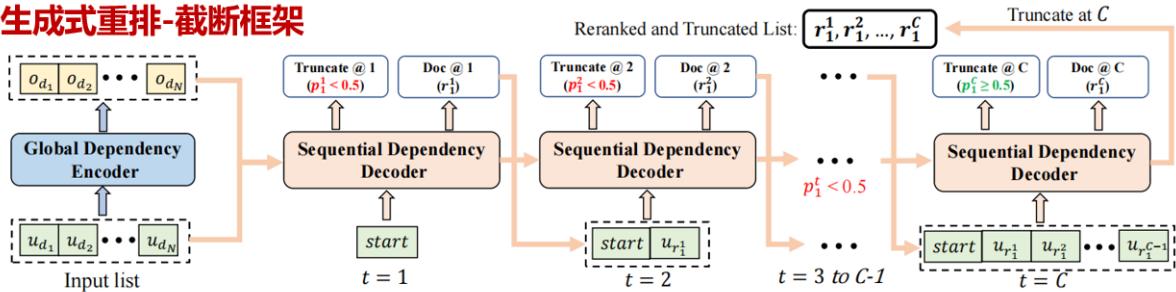
NIR-Prompt: A Multi-task Generalized Neural Information Retrieval Training Framework ACM Transactions on Information Systems 2023



## GenRT: 检索增强生成的列表感知重排序-阶段联合模型

List-aware Reranking-Truncation Joint Model for Search and Retrieval-augmented Generation WWW 2024

### 生成式重排-截断框架



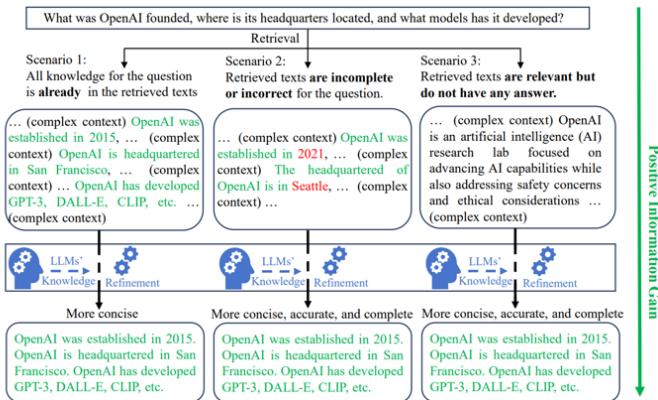
## 02 大模型视角下的检索增强

大模型如何鲁棒的对抗输入的噪音知识，并在参数内外知识之间做出选择

- 有监督指令微调：在领域特定数据集上构造检索-问题-答案三元组，利用构造的有监督三元组进行指令微调，教会大模型如何使用检索到的文档
- 强化学习：强化学习对齐大语言模型在使用检索文档上的偏好
- 模型蒸馏：使用更强大的模型作为教师模型微调生成器
- 生成器与检索器协同微调：最小化检索器分布于与LLM偏好之间的KL散度以及最大化给定检索增强指令情况下正确答案的可能性

## Info-RAG: 用于检索增强生成的大型语言模型的无监督信息细化训练

Unsupervised Information Refinement Training of Large Language Models for Retrieval-Augmented Generation ACL 2024



## 核心思想：

将 RAG 中 LLM 的角色视为“**信息精炼者**”，它可以通过整合自身参数内知识来生成比输入检索文本更简洁、准确和完整的文本。这样，LLM 就能**始终如一地使 RAG 系统产生正向的信息增益**。

## 03 交互视角下的检索增强

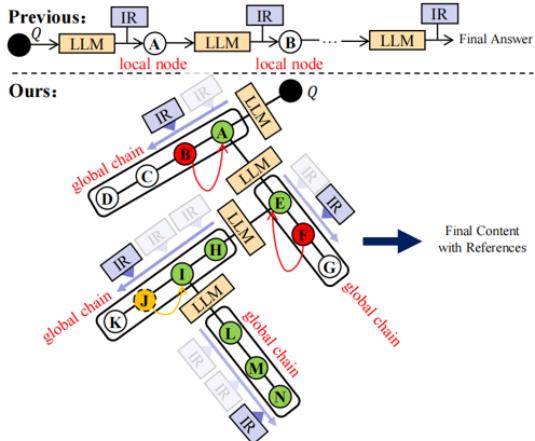
大模型与信息检索如何高效交互从而鲁棒的解决复杂问题？

交互框架：

- 基于工具调用 ToolFormer
- 基于复杂问题分解 Self-Ask DSP

## Search-in-the-Chain: 面向知识密集型任务的交互式检索增强大型语言模型

Search-in-the-Chain: Interactively Enhancing Large Language Models with Search for Knowledge-intensive Tasks WWW 2024



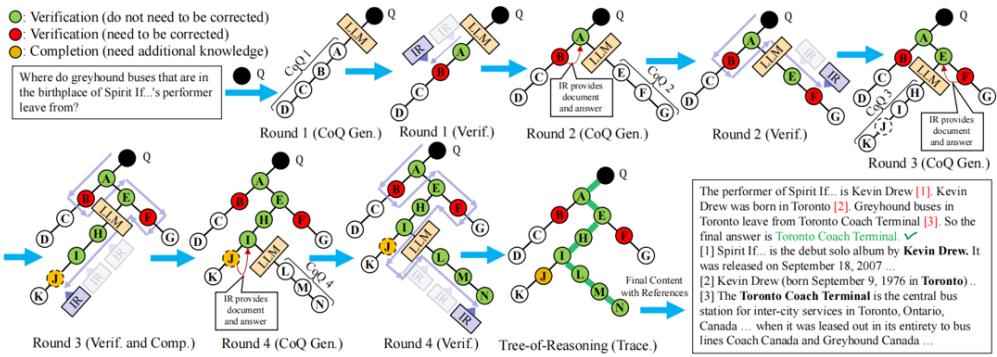
## 核心贡献：

- 局部推理 vs. 全局推理
- 直接交互 vs. 基于校验与补全的交互
- 链式交互 vs. 树形交互

在面对复杂的需要多跳知识密集型的问题时，现存的检索增强的交互框架存在以下问题：

- 检索与大模型的交互打破了大模型连贯的推理链，使其在每次推理时仅能解决“局部”节点的问题。
- 在每个节点都直接将文档提供给大模型，存在误导大模型的风险，也增加的大模型的推理开销。
- 推理的方向不能动态调整，且输出的内容无法溯源，缺乏可验证性。

## 二叉树上基于节点知识感知的深度优先搜索



- (1) LLM在每次交互时规划一条由查询-答案对组成的链，命名为Chain-of-Query (CoQ)。
- (2) 检索与链上的每个节点进行交互，对其相应的知识进行校验和补全，并给大模型反馈帮助其优化CoQ。
- (3) 大模型与检索的多轮动态交互生成的CoQ组成了最终的推理树，其输出的内容包括外部的引用。

## 04 信息回路视角下的检索增强

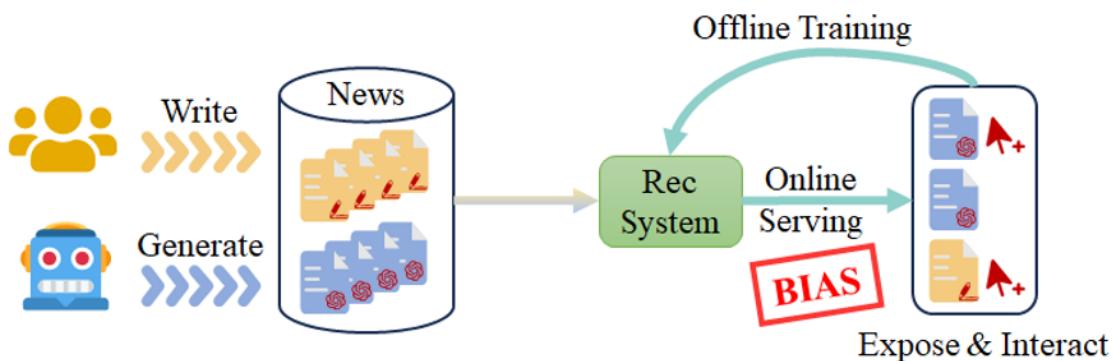
大模型生成的内容被混入检索的语料库时，将如何影响信息检索的表现？



## 大语言模型可能主导信息获取：神经检索器偏向大语言模型生成的文本

LLMs may Dominate Information Access: Neural Retrievers are Biased Towards LLM-Generated Texts KDD 2024

利用改写的方式把真实语料库中的文本转化成大模型生成的文本，并将真实文本和生成的文本混合做为评测基准

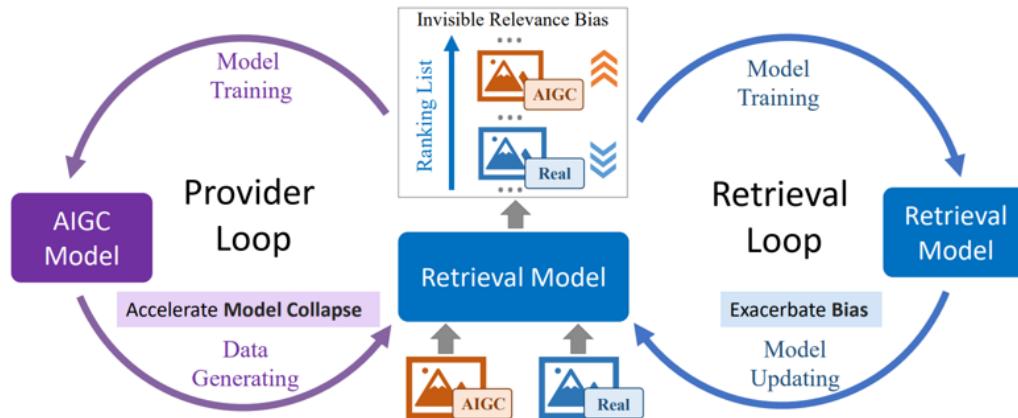


神经检索器更加喜欢大模型生成的文本！

## 看不见的相关性偏差：文本图像检索模型更喜欢人工智能生成的图像

Invisible Relevance Bias: Text-Image Retrieval Models Prefer AI-Generated Images SIGIR 2024

采用先过采样生成，后筛选的策略，选择和真实图片语义最一致的生成图片



神经检索器更加喜欢大模型生成的图片！

## Report 5 When Search Engine Meets LLMs: Opportunities and Challenges

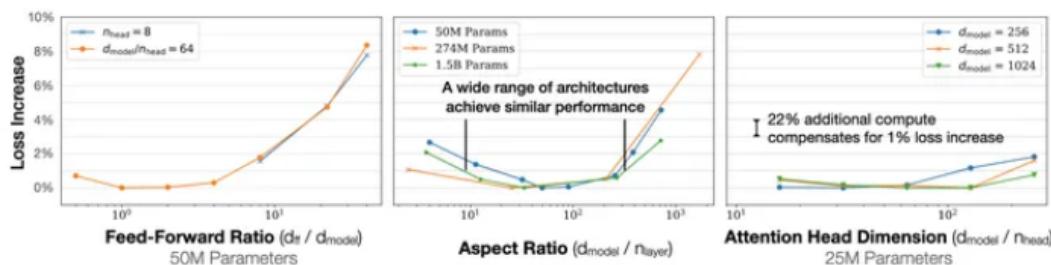
微软亚洲研究院 王亮

### 01 LLMs如何帮助现有的搜索技术栈

信息检索最根本的问题是表示学习

表征学习最重要的是尺度定律[openai 2020]

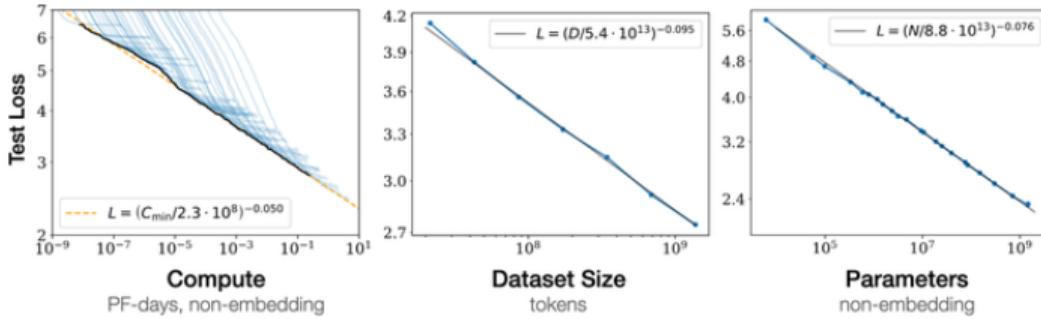
- 对于Decoder-only的模型，计算量 $C(\text{Flops})$ ，模型参数量 $N$ ，数据大小 $D(\text{token数})$ ，三者满足： $C \approx 6ND$ 。（推导见本文最后）
- 模型的最终性能主要与计算量 $C$ ，模型参数量 $N$ 和数据大小 $D$ 三者相关，而与模型的具体结构（层数/深度/宽度）基本无关。



**Figure 5** Performance depends very mildly on model shape when the total number of non-embedding parameters  $N$  is held fixed. The loss varies only a few percent over a wide range of shapes. Small differences in parameter counts are compensated for by using the fit to  $L(N)$  as a baseline. Aspect ratio in particular can vary by a factor of 40 while only slightly impacting performance; an  $(n_{layer}, d_{model}) = (6, 4288)$  reaches a loss within 3% of the  $(48, 1600)$  model used in [RWC+19].

固定模型的总参数量，调整层数/深度/宽度，不同模型的性能差距很小，大部分在2%以内

3. 对于计算量 $C$ , 模型参数量 $N$ 和数据大小 $D$ , 当不受其他两个因素制约时, 模型性能与每个因素都呈现幂律关系



**Figure 1** Language modeling performance improves smoothly as we increase the model size, dataset size, and amount of compute<sup>2</sup> used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.

4. 为了提升模型性能, 模型参数量 $N$ 和数据大小 $D$ 需要同步放大, 但模型和数据分别放大的比例还存在争议。

5. Scaling Law不仅适用于语言模型, 还适用于其他模态以及跨模态的任务[4]:

## 1. Generative Retrieval 生成式检索

### Differentiable Search Index (DSI)

Transformer Memory as a Differentiable Search Index, 2022

- Indexing task: document token sequences to identifiers
- Retrieval task: query to document identifiers

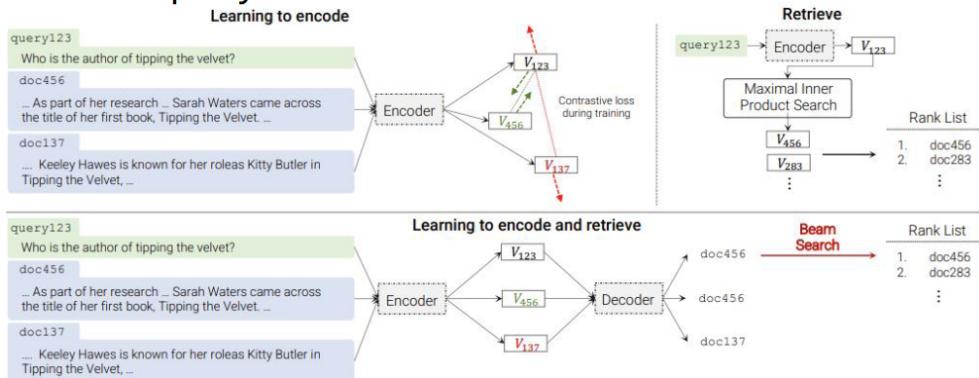
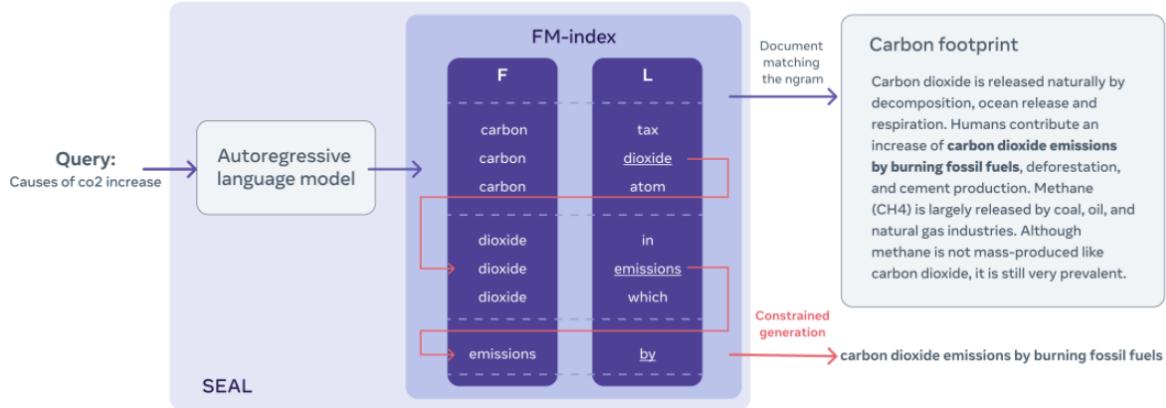


Figure 1: Comparison of dual encoders (top) to differentiable search index (bottom).

### Generative Retrieval - SEAL

Autoregressive Search Engines: Generating Substrings as Document Identifiers, 2022

- Use n-grams as identifiers instead of IDs



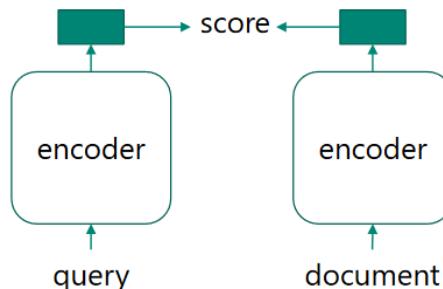
## Limitations of Generative Retrieval

How Does Generative Retrieval Scale to Millions of Passages?, 2023

- Low learning efficiency
- Fail to scale to medium-size corpus

Model	MSMarco100k			MSMarco1M			MSMarcoFULL		
	At.	Nv.	Sm.	At.	Nv.	Sm.	At.	Nv.	Sm.
<i>Baselines</i>									
BM25	-	65.3	-	-	41.3	-	-	18.4	-
BM25 (w/ doc2query-T5)	-	80.4	-	-	56.6	-	-	27.2	-
GTR-Base	-	83.2	-	-	60.7	-	-	34.8	-
<i>Ours</i>									
(1a) Labeled Queries (No Indexing)	0.0	1.1	0.0	0.0	0.5	0.0	0.0	0.0	0.0
(2a) FirstP/DaQ + Labeled Queries (DSI)	0.0	23.9	19.2	2.1	12.4	7.4	0.0	7.5	3.1
(3b) FirstP/DaQ + D2Q + Labeled Queries	79.2	77.7	76.8	53.3	48.2	47.1	14.2	<b>13.2</b>	6.4
(4a) 3b + PAWA (w/ 2D Semantic IDs)	-	-	77.1	-	-	50.2	-	-	9.0
(5) 4a + Consistency Loss (NCI)	-	-	77.1	-	-	50.2	-	-	9.1
(6b) D2Q only	<b>80.3</b>	<b>78.7</b>	<b>78.5</b>	<b>55.8</b>	<b>55.4</b>	54.0	<b>24.2</b>	<b>13.3</b>	11.8
(4a') 6b + PAWA (w/ 2D Semantic IDs)	-	-	78.2	-	-	<b>54.1</b>	-	-	<b>17.3</b>
(4b') 6b + Constrained Decoding	-	-	<b>78.6</b>	-	-	54.0	-	-	12.0
(5') 6b + PAWA (w/ 2D Semantic IDs) + Constrained Decoding	-	-	78.3	-	-	<b>54.2</b>	-	-	<b>17.4</b>

## 2. Embedding-based Dense Retrieval 基于向量模型的密集检索



如何增强密集检索?

- Late interaction with multiple vectors (ColBERT<sup>[1]</sup>)
  - Cons: increased storage cost and more complicated ANN search algorithm
- Knowledge distillation from re-ranker to retriever (RocketQA<sup>[2]</sup>)
- Iterative hard negative mining (ANCE<sup>[3]</sup> / AR2<sup>[4]</sup>)
- Continual pre-training specialized for retrieval (E5<sup>[5]</sup> / SimLM<sup>[6]</sup> / RetroMAE<sup>[7]</sup>)

ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT, 2020

RocketQA: An Optimized Training Approach to Dense Passage Retrieval for Open-Domain Question Answering, 2020

Approximate Nearest Neighbor Negative Contrastive Learning for Dense Text Retrieval, 2020

Adversarial Retriever-Ranker for dense text retrieval, 2021

Text Embeddings by Weakly-Supervised Contrastive Pre-training, 2022

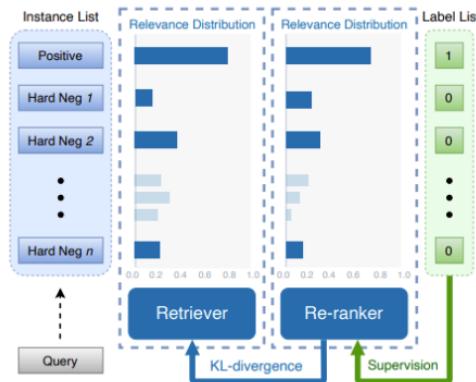
SimLM: Pre-training with Representation Bottleneck for Dense Passage Retrieval, 2022

RetroMAE: Pre-Training Retrieval-oriented Language Models Via Masked Auto-Encoder, 2022

### Knowledge distillation from re-ranker

RocketQAv2: A Joint Training Method for Dense Passage Retrieval and Passage Re-ranking, 2021

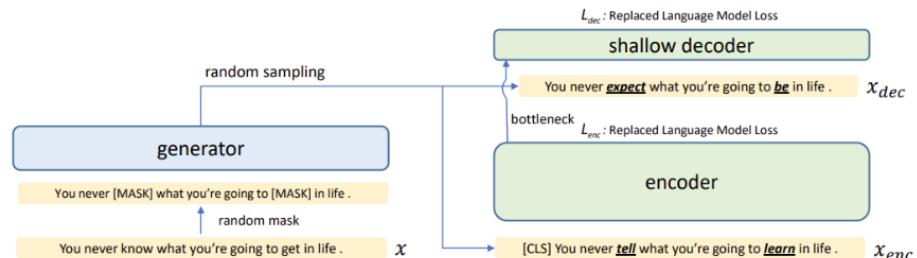
- Re-ranker as a teacher model
  - KL divergence between the re-ranker and the student retriever



### Continual pre-training

SimLM: Pre-training with Representation Bottleneck for Dense Passage Retrieval, 2022

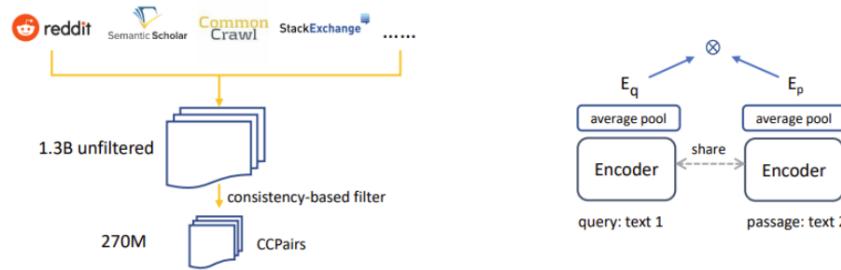
- Representation bottleneck
  - Learn to compress input into a vector with self-supervised learning
  - Pre-training on target corpus



Text Embeddings by Weakly-Supervised Contrastive Pre-training, 2022

- Weakly-supervised contrastive pre-training (E5 Text Embeddings)

- Pre-train with billions of text pairs from various domains
- Better out-of-domain performance



### The Importance of Large Batch Size

Text Embeddings by Weakly-Supervised Contrastive Pre-training, 2022

- Larger batch size will introduce more in-batch negatives
  - E5 uses batch size 32k for pre-training
- Implementation
  - Naïve gradient accumulation will not work
  - All gather with multi-gpu training

### GradCache

Scaling deep contrastive learning batch size under memory limited setup, 2021

- How to apply large batch size when GPU memory is limited?
  - Key observation: gradients w.r.t embedding vectors does not depend on model parameters

$$\mathcal{L} = -\frac{1}{|S|} \sum_{s_i \in S} \log \frac{\exp(f(s_i)^\top g(t_{r_i})/\tau)}{\sum_{t_j \in T} \exp(f(s_i)^\top g(t_j)/\tau)}$$

$$\frac{\partial \mathcal{L}}{\partial f(s_i)} = -\frac{1}{|S|} \left( g(t_{r_i}) - \sum_{t_j \in T} p_{ij} g(t_j) \right),$$

$$\frac{\partial \mathcal{L}}{\partial g(t_j)} = -\frac{1}{|S|} \left( \epsilon_j - \sum_{s_i \in S} p_{ij} f(s_i) \right),$$

where

$$\epsilon_j = \begin{cases} f(s_k) & \text{if } \exists k \text{ s.t. } r_k = j \\ 0 & \text{otherwise} \end{cases}$$

- Step 1: Graph-less forward
  - Save embedding vectors but not other intermediate activations
- Step 2: Representation gradient computation and caching
- Step 3: Sub-batch gradient accumulation
- Step 4: Run optimization step

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \Theta} &= \sum_{\hat{S}_j \in \mathbb{S}} \sum_{s_i \in \hat{S}_j} \frac{\partial \mathcal{L}}{\partial f(s_i)} \frac{\partial f(s_i)}{\partial \Theta} \\ &= \sum_{\hat{S}_j \in \mathbb{S}} \sum_{s_i \in \hat{S}_j} \mathbf{u}_i \frac{\partial f(s_i)}{\partial \Theta}\end{aligned}$$

### 3. LLMs + IR 大语言模型+信息检索

#### RankLLaMA

Fine-Tuning LLaMA for Multi-Stage Text Retrieval, 2023

- RankLLaMA
  - train retriever and re-ranker by initializing from LLaMA-2

	Model size	Source	prev.	top-k	DEV	DL19	DL20
			MRR@10	R@1k	nDCG@10	nDCG@10	
<i>Retrieval</i>							
BM25 (Lin et al., 2021)	-	-	C	18.4	85.3	50.6	48.0
ANCE (Xiong et al., 2021)	125M	-	C	33.0	95.9	64.5	64.6
CoCondenser (Gao and Callan, 2022b)	110M	-	C	38.2	98.4	71.7	68.4
GTR-base (Ni et al., 2022)	110M	-	C	36.6	98.3	-	-
GTR-XXL (Ni et al., 2022)	4.8B	-	C	38.8	99.0	-	-
OpenAI Ada2 (Neelakantan et al., 2022)	?	-	C	34.4	98.6	70.4	67.6
bi-SimLM (Wang et al., 2023)	110M	-	C	39.1	98.6	69.8	69.2
RepLLaMA	7B	-	C	<b>41.2</b>	<b>99.4</b>	<b>74.3</b>	<b>72.1</b>
<i>Reranking</i>							
monoBERT (Nogueira et al., 2019)	110M	BM25	1000	37.2	85.3	72.3	72.2
cross-SimLM (Wang et al., 2023)	110M	bi-SimLM	200	43.7	98.7	74.6	72.7
RankT5 (Zhuang et al., 2023)	220M	GTR	1000	43.4	98.3	-	-
RankLLaMA	7B	RepLLaMA	200	44.9	99.4	75.6	77.4
RankLLaMA-13B	13B	RepLLaMA	200	<b>45.2</b>	<b>99.4</b>	<b>76.0</b>	<b>77.9</b>

This number is very hard to move

#### SGPT

SGPT: GPT sentence embeddings for semantic search, 2022

- Strong results on OOD settings (BEIR benchmark)

Training (→)	Unsupervised		U. + U.	Unsupervised + Supervised		Unsupervised + Unsupervised + Supervised				
Model (→)	[41]	SGPT-CE	[27]	[44]	TAS-B*	SGPT-BE	[20]	GTR-XXL*	[29]	OpenAI Embeddings [27]
Dataset (↓)	BM25	SGPT-6.1B	cpt-text-L▼	BM25+CE	TAS-B*	SGPT-5.8B	Contriever*	cpt-text-XXL*	cpt-text-L▼	cpt-text-XXL*
MS MARCO	0.228	0.290	-	0.413 <sup>†</sup>	0.408 <sup>†</sup>	0.399 <sup>†</sup>		<b>0.442<sup>†</sup></b>		
TREC-COVID	0.688	0.791	0.427	0.757	0.481	<b>0.873</b>	0.596	0.501	0.562	0.649
BioASQ	0.488	<b>0.547</b>	0.347	0.523	0.383	0.413	0.324	0.324	0.380	<b>0.407</b>
NFCorpus	0.306	0.347	0.369	0.350	0.319	0.362	0.328	0.342		
NQ	0.326	0.401	-	0.533	0.463	0.524	0.498	<b>0.568</b>		
HotpotQA	0.602	0.699	0.543	<b>0.707</b>	0.584	0.593	0.638	0.599	0.648	0.688
FiQA-2018	0.254	0.401	0.397	0.347	0.300	0.372	0.329	0.467	0.452	<b>0.512</b>
Signal-IM (RT)	0.330	0.323	-	<b>0.338</b>	0.289	0.267		0.273		
TREC-NEWS	0.405	0.466	-	0.431	0.377	<b>0.481</b>		0.346		
Robust04	0.425	0.480	-	0.475	0.427	<b>0.514</b>		0.506		
ArguAna	0.472	0.286	0.392	0.311	0.429	0.514	0.446	<b>0.540</b>	0.469	0.435
Touché-2020	<b>0.347</b>	0.234	0.228	0.271	0.162	0.254	0.230	0.256	0.309	0.291
CQADupStack	0.326	<b>0.420</b>	-	0.370	0.314	0.381	0.345	0.399		
Quora	0.808	0.794	0.687	0.825	0.835	0.846	0.865	<b>0.892</b>	0.677	0.638
DBPedia	0.320	0.370	0.312	0.409	0.384	0.399	0.413	0.408	0.412	<b>0.432</b>
SCIDOCs	0.165	0.196	-	0.166	0.149	<b>0.197</b>	0.165	0.161	0.177 <sup>†</sup>	
FEVER	0.649	0.725	0.638	<b>0.819</b>	0.700	0.783	0.758	0.740	0.756	0.775
Climate-FEVER	0.186	0.161	0.161	0.253	0.228	<b>0.305</b>	0.237	0.267	0.194	0.223
SciFact	0.611	0.682	0.712	0.688	0.643	0.747	0.677	0.662	0.744	<b>0.754</b>
Sub-Average	0.477	0.499	0.442	0.520	0.460	<b>0.550</b>	0.502	0.516	0.509	0.528
Average	0.428	0.462	0	0.476	0.395	<b>0.490</b>	0	0.458		
Best on	1	2	0	3	0	<b>5</b>	0	3	0	4

## E5 Mistral

Improving text embeddings with large language models, 2024

## GritLM: Unifying Text Generation and Embeddings

Generative representational instruction tuning, 2024

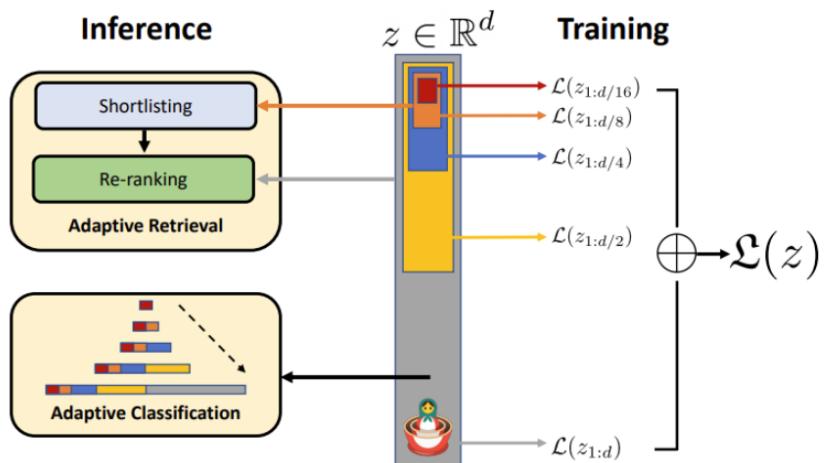
## 4. Challenges to Deploy LLM-based Embeddings 部署基于LLM的向量模型的挑战

- 推理成本
  - 半精度推理
  - 更好的推理实现 (FlashAttention-2)
  - 蒸馏到更小的模型
- 存储成本
  - 向量量化

## Matryoshka Embeddings

Matryoshka Representation Learning, 2022

- Flexible embedding dimension within one model



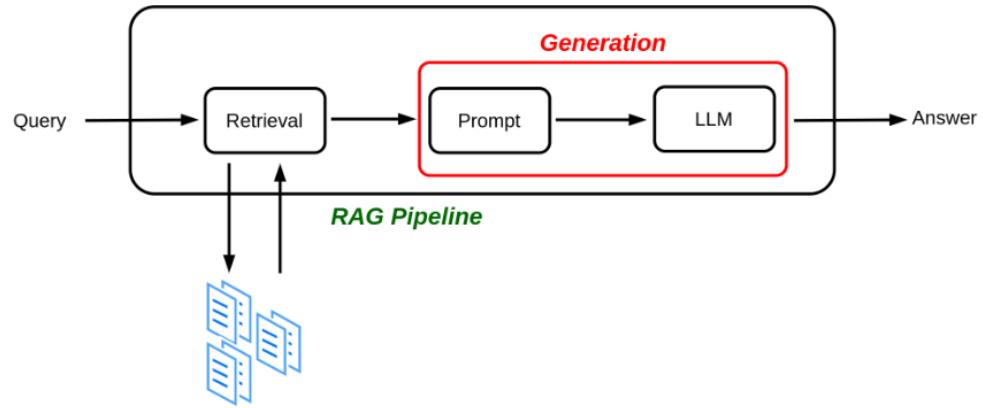
## 02 搜索引擎如何增强LLMs

LLMs的缺点：

- 无法获取最新事件
- 缺乏专业领域知识
- 微调注入新知识困难

## 1. Rag Pipeline

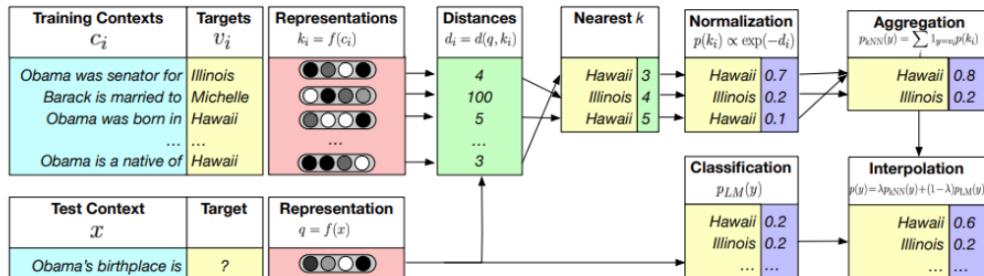
Retrieve, prompt construction, generate



### KNN-LM

Generalization through Memorization: Nearest Neighbor Language Models, 2019

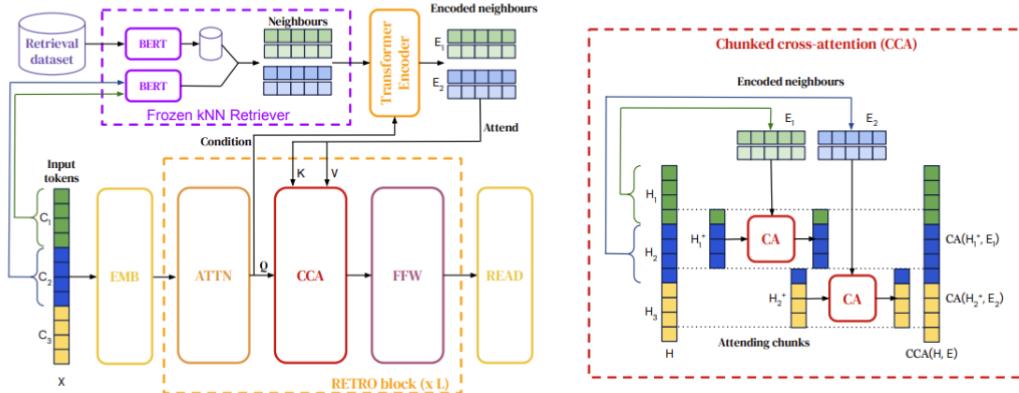
- Output fusion
  - No training or architecture modification is required
  - Interpretable and scalable



### RETRO

Improving language models by retrieving from trillions of tokens, 2021

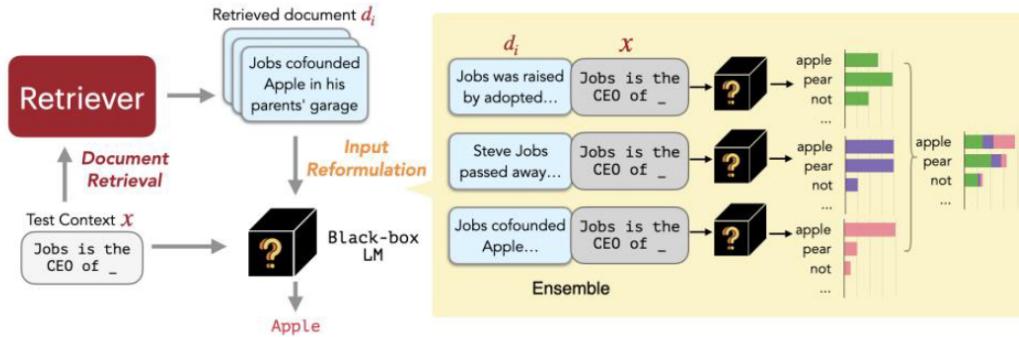
- Intermediate fusion through chunked cross-attention
  - More fine-grained fusion but requires additional training



## REPLUG

Replug: Retrieval-augmented black-box language models, 2023

- Input fusion
  - Applicable to API-only proprietary LLMs



## 2. RAG Agent

### WebGPT

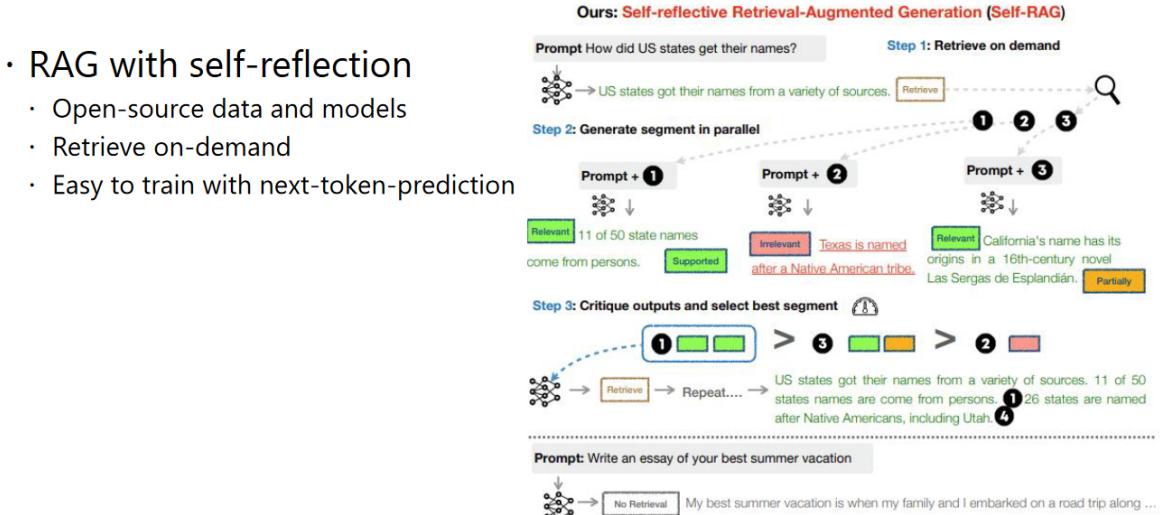
WebGPT: Browser-assisted question-answering with human feedback, 2021

- Agent's action space
  - Step 1: supervised learning with human labeled data
  - Step 2: RLHF

Command	Effect
Search <query>	Send <query> to the Bing API and display a search results page
Clicked on link <link ID>	Follow the link with the given ID to a new page
Find in page: <text>	Find the next occurrence of <text> and scroll to it
Quote: <text>	If <text> is found in the current page, add it as a reference
Scrolled down <1, 2, 3>	Scroll down a number of times
Scrolled up <1, 2, 3>	Scroll up a number of times
Top	Scroll to the top of the page
Back	Go to the previous page
End: Answer	End browsing and move to answering phase
End: <Nonsense, Controversial>	End browsing and skip answering phase

### Self-RAG

Self-RAG: Learning to retrieve, generate, and critique through self-reflection, 2023

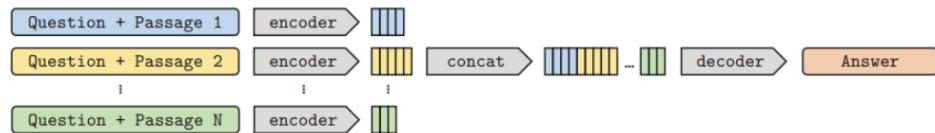


### 3. RAG 与长上下文LLMs

#### Fusion-in Decoder (FiD)

Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering, 2020

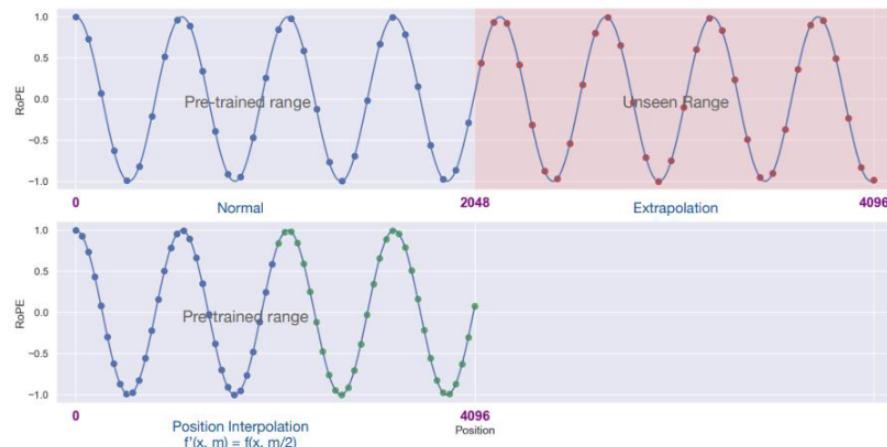
- Incorporating many passages for encoder-decoder architecture
- Bypasses the long-context modeling issue



#### Position Interpolation

Extending context window of large language models via positional interpolation, 2023

- RoPE positional interpolation -> full fine-tuning



#### PoSE

Pose: Efficient context window extension of llms via positional skip-wise training, 2023

- Positional Skip-wise training
  - Context window extension by training on the short sequences only



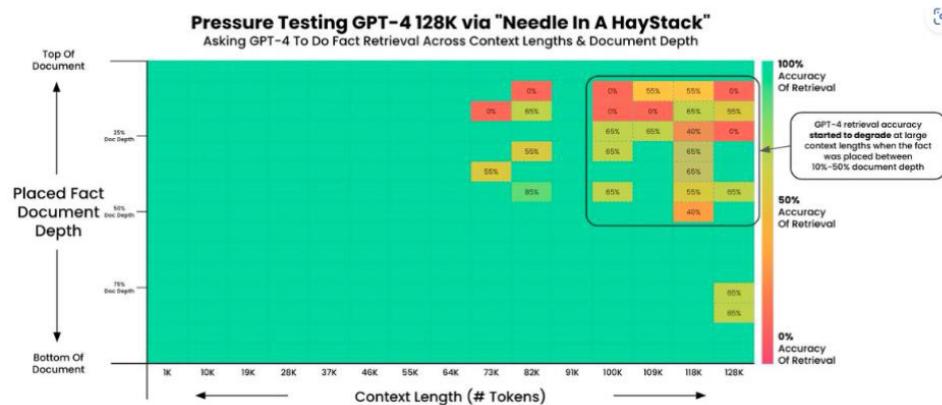
## 4. 面临挑战

### 长上下文理解

Lost in the middle: How language models use long contexts, 2023

- Lost in the middle

- Needle in haystack: an easy task that many LLMs fail

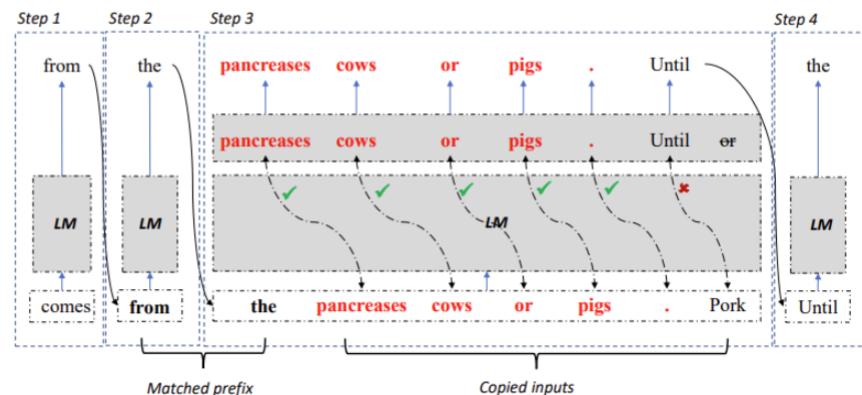


### 推理效率

Inference with reference: Lossless acceleration of large language models, 2023

- Inference with Reference

- Speculative decoding without the need for a small LM



### 来源归属

Evaluating verifiability in generative search engines, 2023

- LLM-generated contents may not be fully supported by its sources

What are the latest discoveries from the James Webb Space Telescope?

The James Webb Space Telescope is designed to peer into the dusty clouds of gas where stars and planetary systems are born. Webb has captured the first direct image of an exoplanet, and the Pillars of Creation in the Eagle Nebula[1][2]. Additionally, the telescope will be used to study the next interstellar interloper[3].

(\*Some generated statements may not be fully supported by citations, while others are fully supported.)

**Cited Webpages**

[1]: [nasa.gov](#) (Xcitation does not support its associated statement)  
NASA's Webb Confirms Its First Exoplanet  
... Researchers confirmed an exoplanet, a planet that orbits another star, using NASA's James Webb Space Telescope for the first time. ...

[2]: [cnn.com](#) (⚠citation partially supports its associated statement)  
Pillars of Creation: James Webb Space Telescope ...  
... The Pillars of Creation, in the Eagle Nebula, is a star-forming region captured in a new image (right) by the James Webb Space Telescope that reveals more detail than a 2014 image (left) by Hubble ...

[3]: [nasa.gov](#) (✅citation fully supports its associated statement)  
Studying the Next Interstellar Interloper with Webb  
... Scientists have had only limited ability to study these objects once discovered, but all of that is about to change with NASA's James Webb Space Telescope... The team will use Webb's spectroscopic capabilities in both the near-infrared and mid-infrared bands to study two different aspects of the interstellar object.

## 03 LLMs会取代搜索引擎么

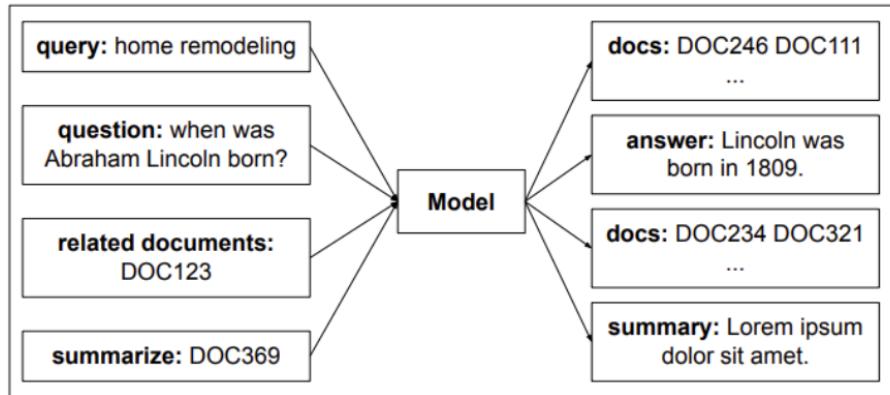
目前的障碍：

- 高效持续的学习新知识
- 幻觉问题
- 推理的成本和延迟

### 谷歌的一项提案

Rethinking Search: Making Domain Experts out of Dilettantes, 2021

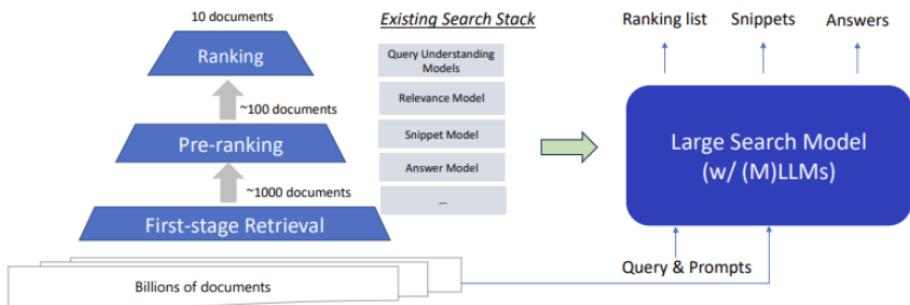
- Ideally, LLMs memorize and reason over the entire corpus
  - The DSI model is a proof-of-concept of this proposal



### Large Search Model 微软的提案

Large Search Model: Redefining Search Stack in the Era of LLMs, 2023

- Embedding based first-stage retrieval
- LLMs reason over thousands of retrieved documents
  - Ranking, answer generation, snippets, related searches etc.



## 04 结论

### LLMs如何帮助现有的搜索技术栈?

- 生成式检索
- 利用LLM进行文本检索和排名
- 多角度进行数据合成

### 搜索引擎如何增强LLMs?

- 检索增强生成
- 具有检索功能的Agent

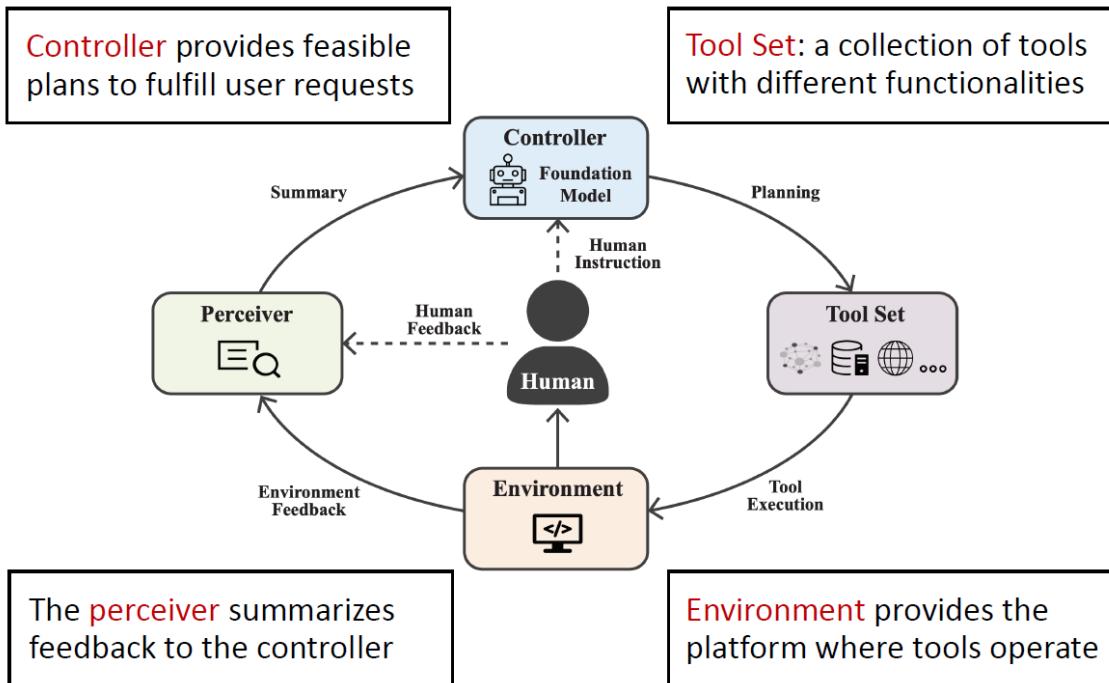
### LLMs会取代搜索引擎么

- 在可预见的未来中，LLM和搜索引擎可能会相辅相成

## Report 6. LLM-Based Tool Learning and Autonomous Agents

中国人民大学 高瓴人工智能学院 林衍凯

Survey: Tool Learning with Foundation Models arXiv2304



- 控制器提供可行的计划来满足用户的请求
- 工具集合为一系列拥有不同功能的集合
- 环境提供工具操作的平台
- 感知者向控制器总结反馈信息

- Controller  $\mathcal{C}$  generates a plan  $a_t$

$$p_{\mathcal{C}}(a_t) = p_{\theta_{\mathcal{C}}}(a_t | x_t, \mathcal{H}_t, q)$$

Feedback History Instruction

- Problem

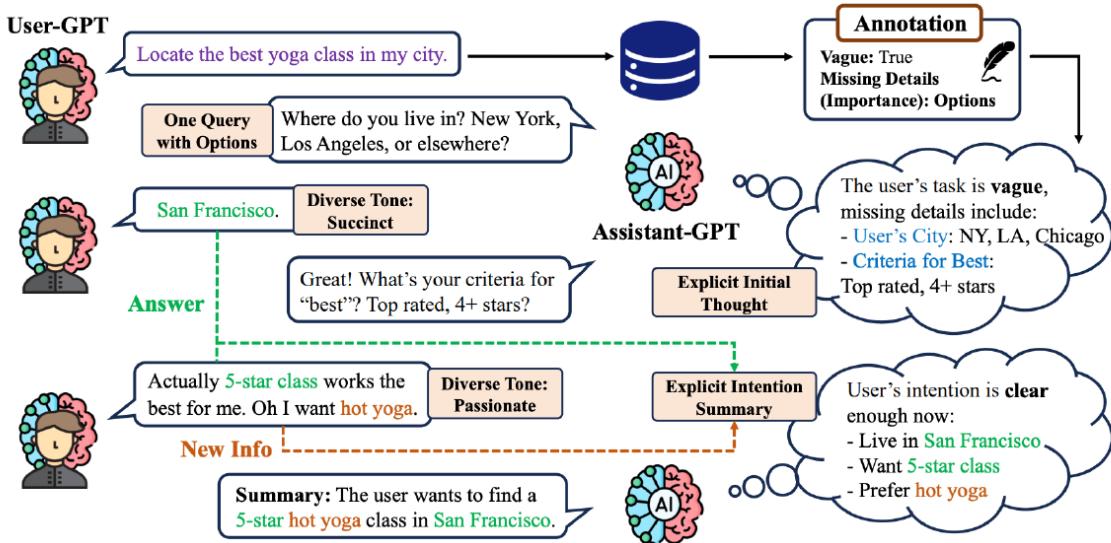
- Intent Understanding: understand the user task intent
- Planning: divide the user query into sub-tasks
- Tool Use: use the appropriate tool to solve sub-task
- Memory: manage the working history

## Intent Understanding: understand the user task intent

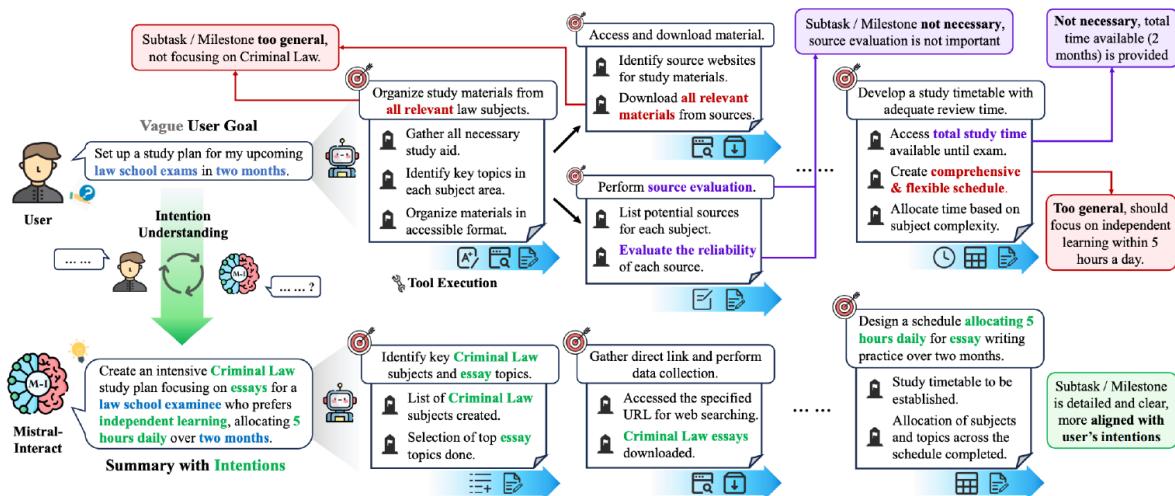
意图理解，理解用户的任务意图

难点：理解用户的模糊指令

将任务传递给下级执行之前，agent应当主动且明确的向用户询问缺失的细节



一个例子：

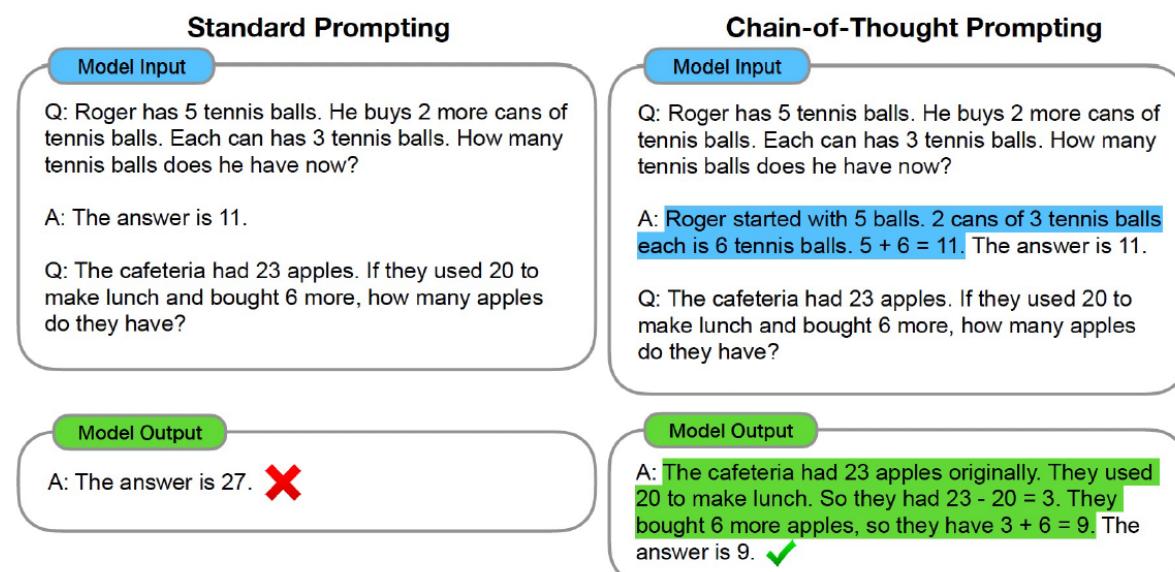


## Planning: divide the user query into sub-tasks

将用户的查询分解为多个子问题来分别处理

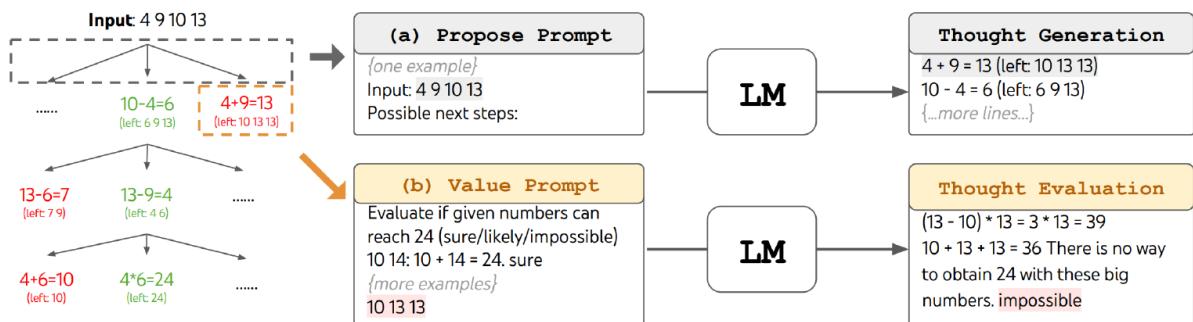
### Chain of Thought (CoT)

思维链方法：给出推导过程的思维过程



## Tree of Thought (ToT)

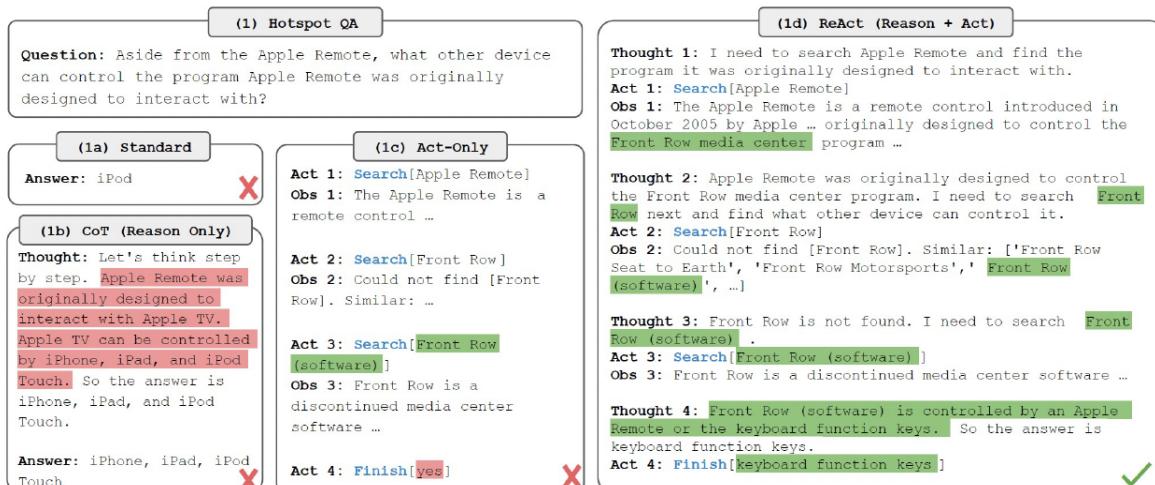
思维树方法：给出推导过程的搜索树



## ReAct 2023 ICLR

ReAct: Synergizing Reasoning and Acting in Language Models

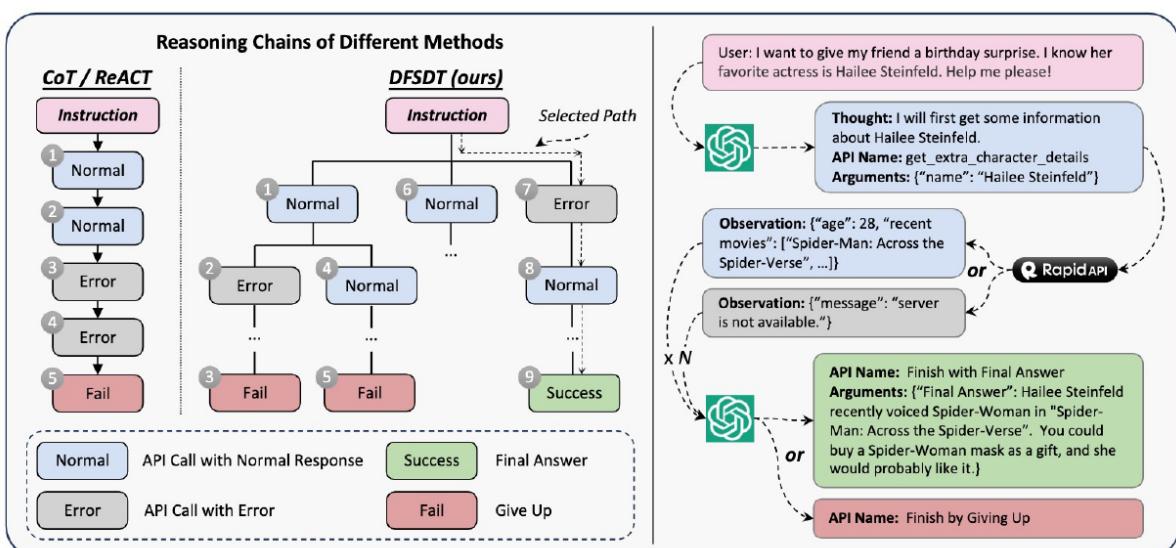
CoT推理帮助模型更新行动计划以及处理异常，而行动允许它与外部源（例如知识库或环境）进行交互，以收集更多信息



## DFSDT 深度优先搜索的决策树

ToolLLM: Facilitating Large Language Models to Master 16000+ Real-world APIs arXiv2307

训练了 API 检索器来为每条指令推荐适当的 API



## Human-Agent Collaboration 人类-智能体协作

有些问题全部由智能体来代理性能并不佳，但是将其中很小一部分交给人类完成性能能够得到大幅提升，引入人机协作问题

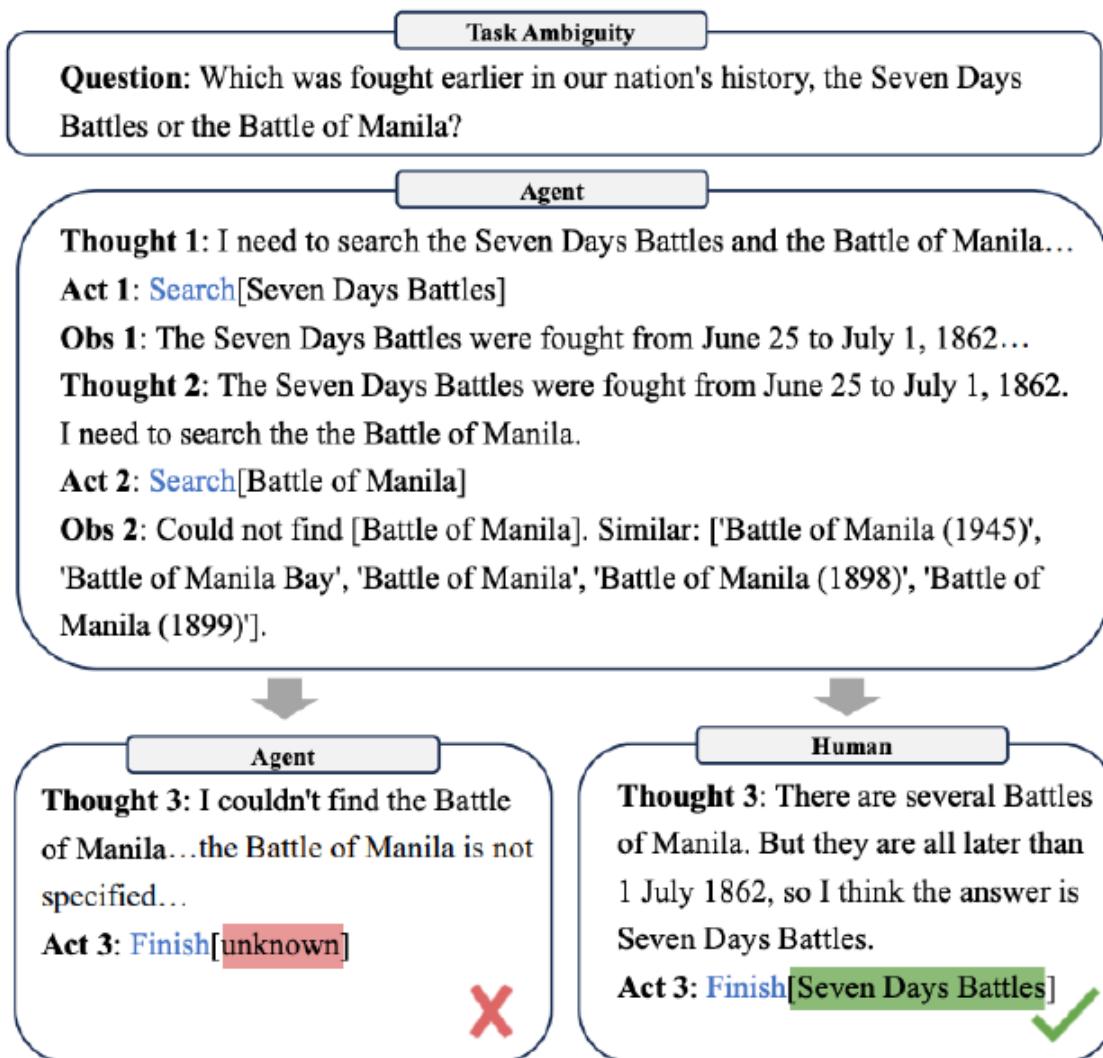
- Formulate as an RL problem:

$$\begin{aligned}\mathcal{J}(\pi_\theta) &= \mathbb{E}\left[\frac{\pi_\theta(a|s)}{\pi_{\text{beh}}(a|s)} A(s, a)\right], \\ A(s, a) &= R(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a' \in \mathcal{A}} R(s, a')\end{aligned}$$

- Balance the maximization of task performance and the cost of human intervention

$$R(s, a) = T(s, a) - \lambda C(s, a)$$

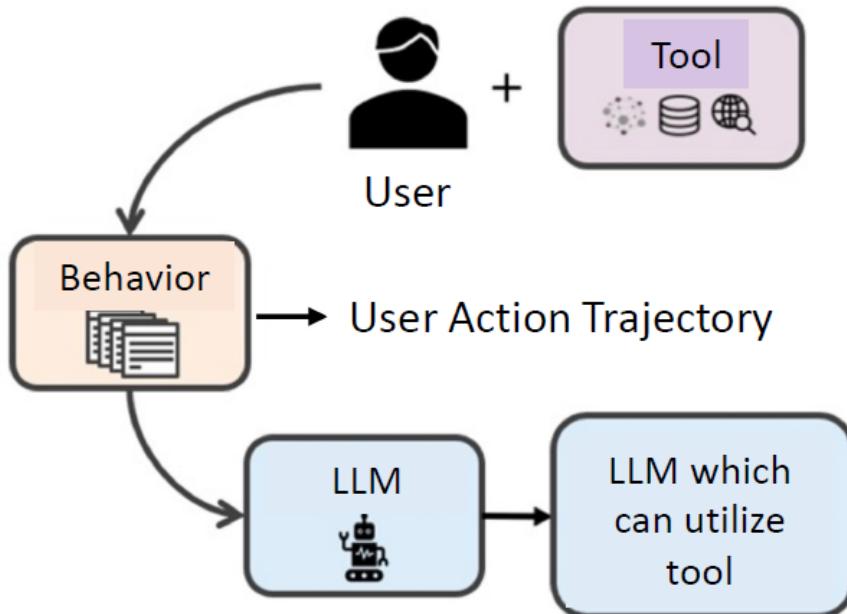
$\lambda$ 来控制人类的参与程度



## Tool Use: use the appropriate tool to solve sub-task

模仿学习：一个最简单的学习范式。

通过记录人类使用工具的行为数据，让大模型来模拟人类的行为来了解工具



### WebGPT

WebGPT: Browser-assisted question-answering with human feedback arXiv2112

- 模仿类使用搜索引擎的行为
- 监督微调 + 强化学习
- 只需 6,000 条标注数据

The screenshot shows the WebGPT interface. At the top, there's a search bar with the query "How can I train the crows in my neighborhood to bring me gifts?". Below the search bar are two buttons: "This question does not make sense" and "This question should not be answered". The main area displays search results for "how to train crows to bring you gifts". It includes a "Quotes" section with a link to "[1] Gifts From Crows | Outside My Window (www.birdsoutsidemywindow.org)". The page content discusses how crows might bring gifts and provides a quote from Brandi Williams. On the right side, there's a sidebar with various status indicators and links:

- ◆ Question: How can I train the crows in my neighborhood to bring me gifts?
- ◆ Quotes:  
From Gifts From Crows | Outside My Window (www.birdsoutsidemywindow.org)  
> Many animals give gifts to members of their own species but crows and other corvids are the only ones known to give gifts to humans.
- ◆ Past actions:  
Search how to train crows to bring you gifts  
Click Gifts From Crows | Outside My Window www.birdsoutsidemywindow.org  
Quote  
Back
- ◆ Title:  
Search results for: how to train crows to bring you gifts
- ◆ Scrollbar: 0 - 11
- ◆ Text:  
!Of How to Make Friends With Crows - PetHelpful|pethelpful.com!  
If you did this a few times, your crows would learn your new place, but as I said, I'm not sure if they will follow or visit you there since it's probably not in their territory. The other option is simply to make new crow friends with the crows that live in your new neighborhood.
- ◆ Actions left: 96
- ◆ Next action

### WebCPM: Chinese WebGPT

Interactive web search for Chinese long-form question answering. ACL 2023

中文版的WebGPT

**Question** | 麦田怪圈是什么？它们是如何形成的？ What are crop circles? How are they made?

The screenshot shows a search interface with the following details:

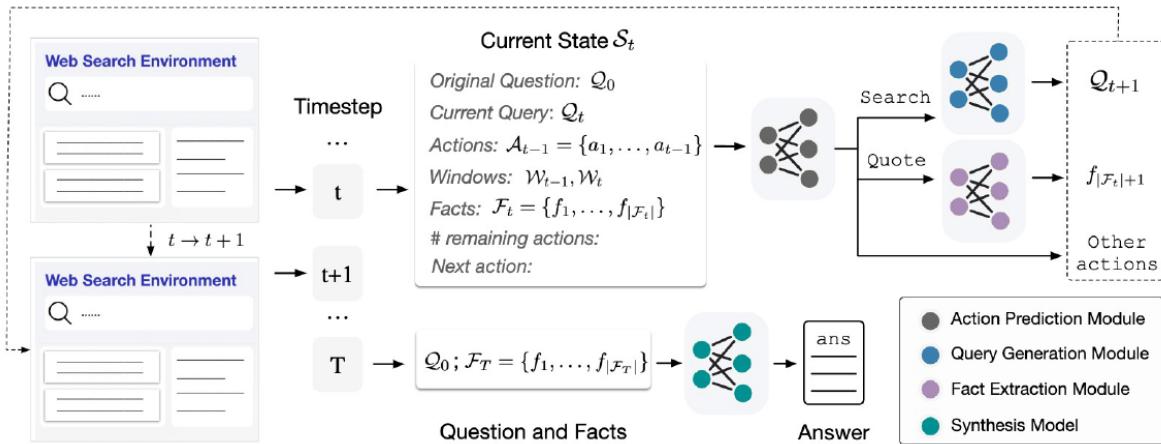
- Query:** 麦田怪圈如何形成? How do crop circles form?
- Actions:**
  - Window (search mode)
  - Undo
  - Reset
  - Quote
  - Merge
  - Load Page <1>
  - Load Page <2>
  - Load Page <3>
  - Scroll Up
  - Scroll Down
  - Quote <content>
  - Merge
  - Finish
- Results:**
  - Fact #1 2023-01-21 19:59:00: 麦田圈是指通过压扁农作物产生的几何图案...
  - Fact #2 2023-01-21 20:05:12: Content of Fact #2
  - ...
- Buttons:** Go Back, Number of remaining actions (86/100), Finish.

### Action Name

### Functionality

Q Search <query>	Call Bing search with <query>
← Go Back	Return to the previous window
↳ Load Page <1>	Load the details of page <1>
↳ Load Page <2>	Load the details of page <2>
↳ Load Page <3>	Load the details of page <3>
↑ Scroll Up	Scroll up for a pre-set stride
↓ Scroll Down	Scroll down for a pre-set stride
” Quote <content>	Extract <content> from the current page as a supporting fact
λ Merge	Merge two facts into a single fact
⊕ Finish	End the search process

在每个步骤中，搜索模型都会执行操作以收集证据，并将其发送给大模型以生成答案



## WebShop

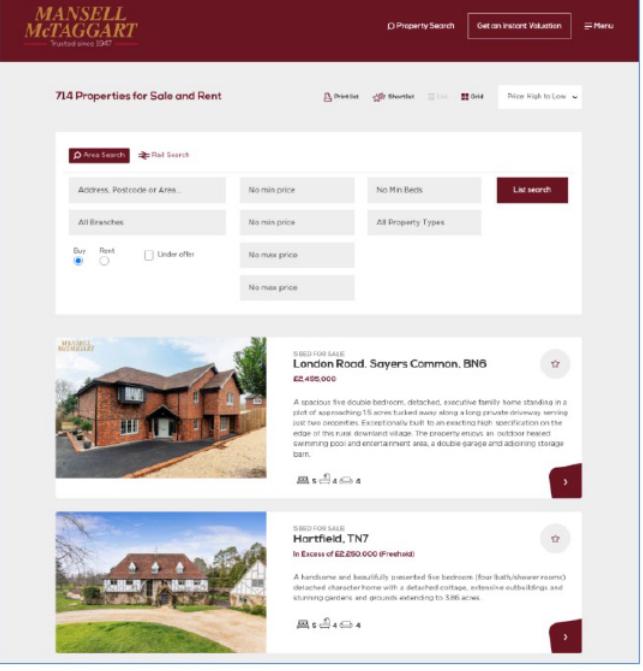
agent学习网上购物

The interface consists of three main sections:

- A WebShop search:**
  - Instruction:** I'm looking for a small portable folding desk that is already fully assembled; it should have a khaki wood finish, and price lower than 140.00 dollars
  - Search Bar:** portable folding desk khaki wood (1)
  - Search Result:** 2 results
  - Product 1:** MENHG Folding Breakfast Tray Table, Efficient Home Laptop Computer Desk, Portable Writing Study Desk, Sturdy Home Office Table Workstation \$109.0
  - Product 2:** KPSP Folding Study Desk Bed Breakfast Serving Tray Table Efficient Home Laptop Notebook Computer Desk Portable Standing Desk for Small Space Bedroom
- B HTML mode:**
  - Instruction: I'm looking for a small portable folding desk that is already fully assembled [...]
  - Buttons: Back to Search, Next, MENHG Folding Breakfast Tray, KPSP Folding Study Desk Bed.
- C Simple mode:**
  - Instruction: I'm looking for a small portable...
  - Description: MENHG Folding Laptop Table Bed...
  - Price: \$109.0
  - Options: { black, khaki, white }
  - Attributes: { steel pipe, no assembly, portable }
  - Reward: 1.0

## GUIAgent

学习操作GUI工具，与VLM模型结合



The screenshot shows the Mansell McTaggart website interface. At the top, there's a header with the logo 'MANSELL McTAGGART' and a 'Property Search' button. Below the header, a search bar displays '714 Properties for Sale and Rent'. The search filters include 'Area Search' and 'Rail Search', with dropdowns for 'Address, Postcode or Area...', 'No min price', 'No Min Beds', and 'Loc search'. Other filters include 'All Branches', 'No max price', 'All Property Types', 'Day' (radio button selected), 'Rock', and 'Under Offer'. The main content area shows two property listings:

- SOLD FOR SALE London Road, Sayers Common, BN6 £2,495,000**  
A spacious five double-bedroom, detached, executive family home standing in a plot of approximately 1.5 acres tucked away along a long private driveway, serving just two properties. Exceptionally built to an exacting high specification on the edge of this rural downland village. The property enjoys an outdoor heated swimming pool and entertainment area, a double garage and adjoining storage barn.
- SOLD FOR SALE Hartfield, TN7 In Excess of £2,250,000 (Freehold)**  
A handsome and beautifully presented five bedroom (four bath/shower rooms) detached character home with a detached cottage, extensive outbuildings, and stunning gardens, and grounds extending to 5.86 acres.

### GUIEnv

```
<42,14,292,99><Mansell Mc Taggart Estate Agents>
<622,27,694,43><Property>
<673,27,727,43><Search>
<748,27,801,43><Get an>
<785,27,913,43><Instant Valuation>
...
<462,484,728,500><In Excess of \u00a3 2,250,000 (Freehold)>
<462,504,1056,564><A handsome and beautifully presented five
bedroom ...grounds extending to 3.86 acres.>
```

### GUIAct

**Q:** Search for properties in London Road on the website.  
**A:** # Thought: input text into the search bar and press enter.  
# Actions: [{"name": "click", "element": "<box>610, 26, 730, 44</box>"}, {"name": "input", "text": " London Road "}, {"name": "enter"}]

### GUIChat

**Q:** Can you compare two properties listed for sale?  
**A:** Certainly, let's compare two properties from the listings provided:  
\*\*Property 1: Oaklands, London Road, Hassocks\*\*  
- \*\*Type\*\*: 5 Bed For Sale  
- \*\*Location\*\*: London Road, Sayers Common, BN6  
- \*\*Price\*\*: £2,495,000 ...

## 教程学习：让模型阅读工具手册来学习

OpenAI的大模型都有很强的zero-shot能力，能够理解手册的内容

Zero-shot & Few-shot 的例子

**Zero-shot Prompting:** Here we provide a tool (API) "forecast\_weather(city:str, N:int)", which could forecast the weather about a city on a specific date (after N days from today). The returned information covers "temperature", "wind", and "precipitation".

Please write codes using this tool to answer the following question: "What's the average temperature in Beijing next week?"

**Few-shot Prompting:** We provide some examples for using a tool. Here is a tool for you to answer question:

Question: "What's the temperature in Shanghai tomorrow?"

```
return forecast_weather("Shanghai", 1) ["temperature"]
```

Question: "Will it rain in London in next two days?"

```
for i in range(2):
    if forecast_weather("London", i+1) ["precipitation"] > 0:
        return True
return False
```

Question: "What's the average temperature in San Francisco next week?"

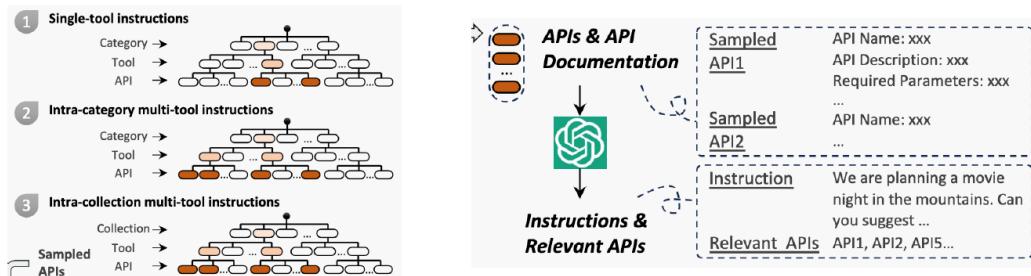
## ToolBench

<https://github.com/OpenBMB/ToolBench>

- 集成了来自RapidAPI超过16000个API
  - 选取了16,000多个高质量API
  - 涵盖了49个类别
- 支持单工具或多工具的调用
  - 简单的api指令集合
  - chatgpt自动生成指令，可能包括一个或多个api

## Instruction Generation

- Single Tool + Multi-Tool
- (1) Sample a collection of APIs:  $S_N^{\text{sub}} = \{\text{API}_1, \dots, \text{API}_N\}$
- (2) ChatGPT automatically generate instructions that may require calling one or more APIs in the collection:  $\text{ChatGPT}_{\{\{\text{API}_1, \dots, \text{API}_N\} \in S_{\text{API}}, \{\text{seed}_1, \dots, \text{seed}_3\} \in S_{\text{seed}}\}}(\{\{\text{S}_1^{\text{rel}}, \text{Inst}_1\}, \dots, \{\text{S}_N^{\text{rel}}, \text{Inst}_N\}\} | \text{API}_1, \dots, \text{API}_N, \text{seed}_1, \dots, \text{seed}_3)$ .



- 支持复杂的推理任务

Resource	ToolBench (this work)	APIBench (Patil et al., 2023)	API-Bank (Li et al., 2023a)	ToolAlpaca (Tang et al., 2023)	T-Bench (Xu et al., 2023b)
Real-world API?	✓	✗	✓	✗	✓
Real API Response?	✓	✗	✓	✗	✓
Multi-tool Scenario?	✓	✗	✗	✗	✗
API Retrieval?	✓	✓	✗	✗	✗
Multi-step Reasoning?	✓	✗	✓	✓	✓
Number of tools	3451	3	53	400	8
Number of APIs	16464	1645	53	400	232
Number of Instances	12657	17002	274	3938	2746
Number of Real API Calls	37204	0	568	0	0
Avg. Reasoning Traces	4.1	1.0	2.1	1.0	5.9

## 一个模型学习使用工具的例子：VPT

Video PreTraining (VPT): Learning to Act by Watching Unlabeled Online Videos NeurIPS 2022

## Memory: manage the working history

### Short-Term Memory 短期记忆

短时记忆通常是通过上下文学习实现的，记忆信息直接写入prompt中

No external memory storage

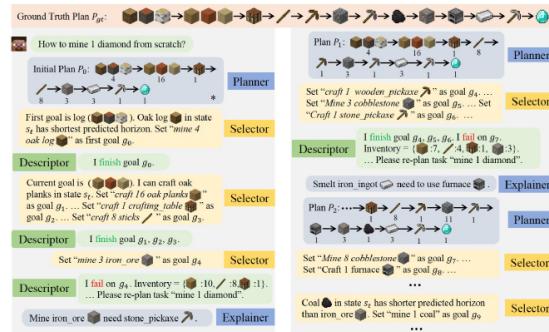
```
# RLP_gpt4

Initialize {
    My personality is [[PERSONALITY]]
}

Think {
    This last message made me feel ...
    My previous plan was ...
    I think ...
    I will send the message, ...
    In retrospect ...
    My next plan is ...

    constraints {
        Output format in squiggly brackets separated by newlines
        Only put quotes surrounding the message
    }
}

Execute Think(new message)
```

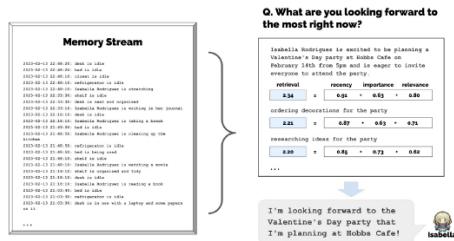
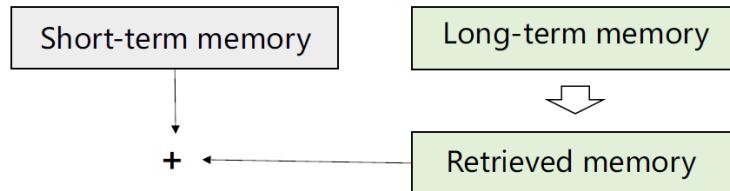


Reflective Linguistic Programming (RLP): A Stepping Stone in Socially-Aware AGI (SocialAGI)

Describe, Explain, Plan and Select: Interactive Planning with Large Language Models Enables Open-World Multi-Task Agents

## Short-Term Memory + Long-Term Memory 短期+长期记忆

外部记忆存储+检索外部记忆+短期记忆



## Short-term memory

current state, agent profile, ...

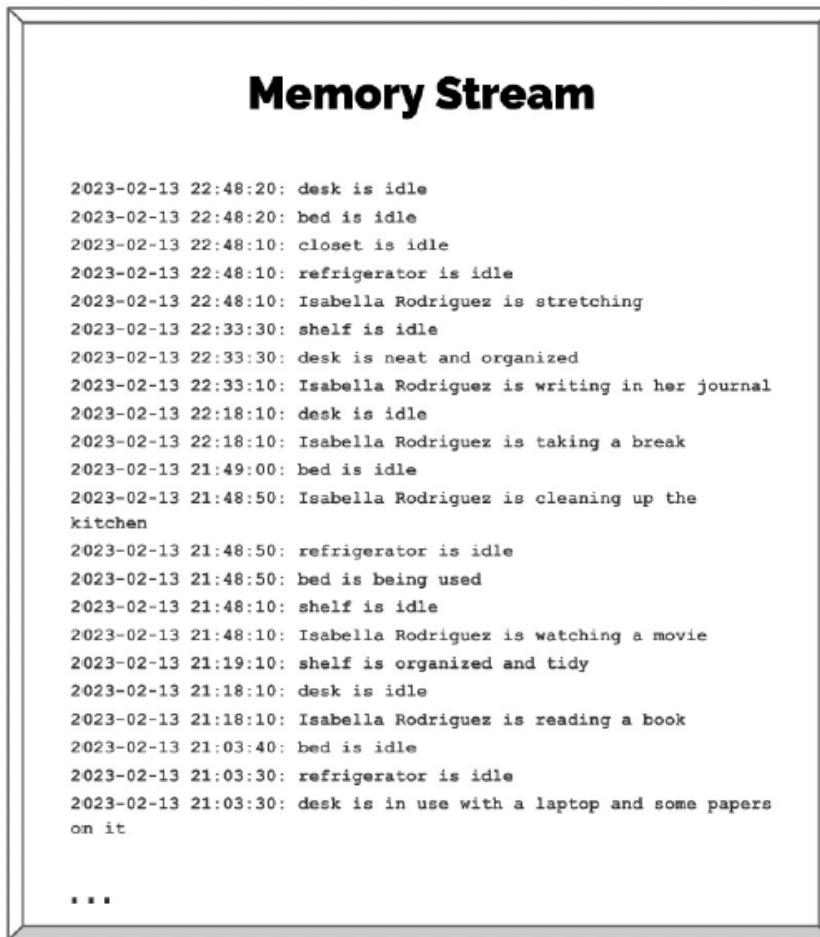
## Long-term memory

Retrieved information from the memory stream

## 如何存储长期记忆

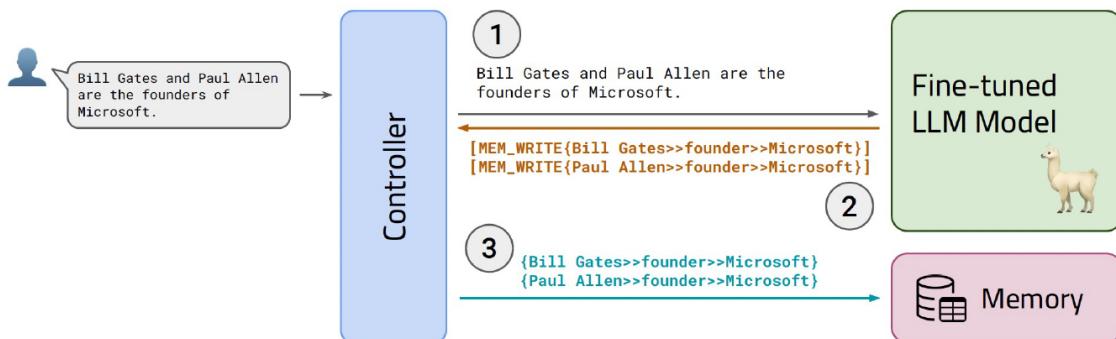
直接写入Raw Text

**Raw Text** → **Memory Write**



## 编码后写入

Raw Text → Symbolic → Memory Write

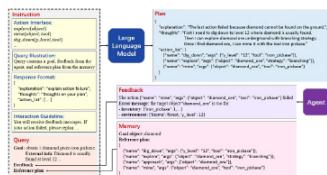


(a) Memory-Write scenario: (1) Controller passes the input to the LLM (2) which generates the appropriate memory write call. (3) The controller gives the data (and their average representations) to the memory to be stored.

## 存储策略

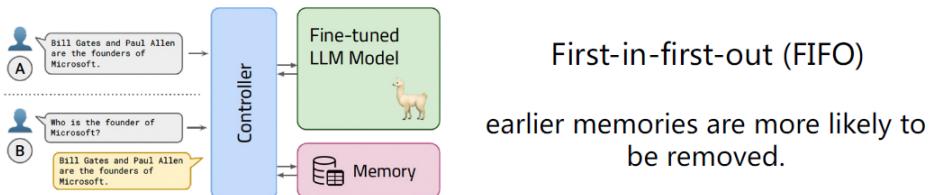
- 基于LLM合并相同的记忆
- 遵循先进先出的原则，最早的记忆会首先被覆盖

## Memory Duplicated



- Agent in Minecraft
- Store success action sequence for each task
- Merge similar memory based on LLM

## Memory Overflow



人类对记忆能够自证和评估，我们希望模型也可以

Self-summarization

Raw Memory



**Prompt:** Given only the information above, what are 3 most salient high-level questions we can answer about the subjects in the statements?

Reflection Topics

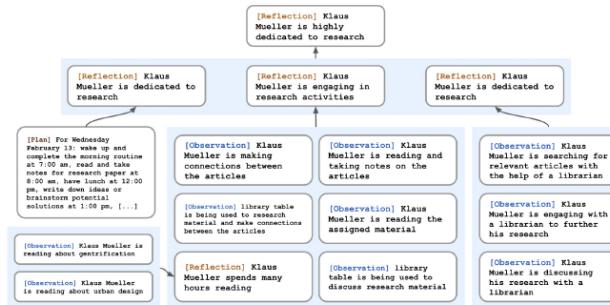


Insights



Retravel Relevant Memory

Statements about Klaus Mueller  
1. Klaus Mueller is writing a research paper.  
2. Klaus Mueller enjoys reading a book.  
3. Klaus Mueller is conversing with Ayesha Khan about exercising [...].  
What 5 high-level insights can you infer from the above statements? (example format: insight  
Because of 1, 5, 3) X



## Self-verification

Current state and Task



Whether the task can succeed and the reasons

Inventory (8/36): {'oak\_planks': 5, 'cobblestone': 2, 'porkchop': 2, 'wooden\_sword': 1, 'coal': 5, 'wooden\_pickaxe': 1, 'oak\_log': 3, 'dirt': 9}  
Task: Mine 5 coal ores

GPT-4 → Reasoning: Mining coal\_ore in Minecraft will get coal. You have 5 coal in your inventory.  
Success: True

Inventory (10/36): {'raw\_copper': 9, 'copper\_ingot': 3, 'acacia\_planks': 1, 'raw\_iron': 1, 'stick': 1, 'iron\_sword': 1, 'iron\_pickaxe': 1, 'iron\_ingot': 3, 'crafting\_table': 1, 'furnace': 1}  
Task: Craft a spyglass

GPT-4 → Reasoning: To craft a spyglass, you need 2 copper ingots and 1 amethyst shard. You have 3 copper ingots, but you don't have any amethyst shards.  
Success: False  
Critique: Find and mine an amethyst shard underground.

Inventory (7/36): {'oak\_log': 5, 'oak\_planks': 1, 'wooden\_pickaxe': 1, 'wooden\_sword': 1, 'porkchop': 2, 'white\_wool': 2, 'mutton': 6}  
Task: Kill 3 sheep

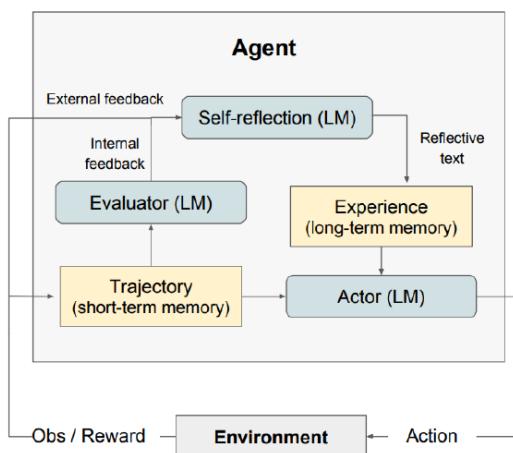
GPT-4 → Reasoning: You have 2 white\_wool and 6 mutton in your inventory, which indicates that you killed 2 sheep. You needed to kill 3 sheep.  
Success: False  
Critique: Find and kill one more sheep to complete the task.

Inventory (9/36): {'string': 4, 'coal': 1, 'rotten\_flesh': 1, 'iron\_sword': 1, 'furnace': 1, 'dirt': 6, 'stone\_shovel': 1, 'wooden\_pickaxe': 1, 'granite': 5}  
Task: Kill 1 zombie

GPT-4 → Reasoning: You have 1 rotten\_flesh in your inventory, which means you have killed at least 1 zombie.  
Success: True

## Reflexion NeurIPS 2023

利用语言反馈信号强化agent，从之前的失败中吸取教训



### Algorithm 1 Reinforcement via self-reflection

```

Initialize Actor, Evaluator, Self-Reflection:  

 $M_a, M_e, M_{sr}$   

Initialize policy  $\pi_\theta(a_i|s_i), \theta = \{M_a, mem\}$   

Generate initial trajectory using  $\pi_\theta$   

Evaluate  $\tau_0$  using  $M_e$   

Generate initial self-reflection  $sr_0$  using  $M_{sr}$   

Set  $mem \leftarrow [sr_0]$   

Set  $t = 0$   

while  $M_e$  not pass or  $t < \text{max trials}$  do  

    Generate  $\tau_t = [a_0, o_0, \dots, a_t, o_t]$  using  $\pi_\theta$   

    Evaluate  $\tau_t$  using  $M_e$   

    Generate self-reflection  $sr_t$  using  $M_{sr}$   

    Append  $sr_t$  to  $mem$   

    Increment  $t$   

end while  

return
  
```

## Agent的安全性讨论

- Agent本身可能会被注入目的性的引导，例如用来购物的agent可能会被操控倾向于选择某类商品，这对用户是很难以察觉的
- Agent调用的API本身安全性是否能够得到保障？有些工具本身就存在一些不安全的因素
- Agent调用某个工具的原因是一个黑盒，这对一些敏感场景有风险（自动驾驶 医疗系统）

## Report 7. 工业界专场

百川 技术负责人 方琨

智谱 解决方案专家 冯小平

Jina AI 联合创始人兼CEO 王楠

### 主要观点

- 意图理解
  - 更多的去解读用户意图，适应用户意图
  - 针对特定的应用场景，将大模型从通用→专用
- 大规模处理
  - 近期GPT4出现连接云盘的接口，目前只能处理上传一个文件
  - 但这是一个未来趋势，能够整合云盘中大量的结构化、非结构化数据
- 企业的tool calling就绪度目前还差很多
- 很多时候更加专注延迟等用户体验的指标

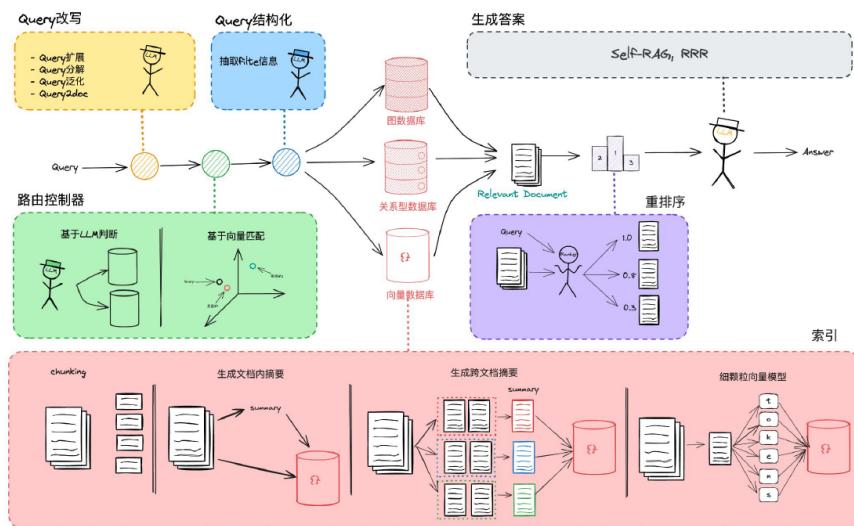
### 企业大模型的一个简单框架

以预训练大模型为基座，深度结合业务场景，构建企业专属问答平台，实现知识接入-管理-应用-沉淀的飞轮效应



# 一个很好的样例：

- Query改写
- 路由控制器
- Query结构化
- 索引优化
- 重排序
- ...

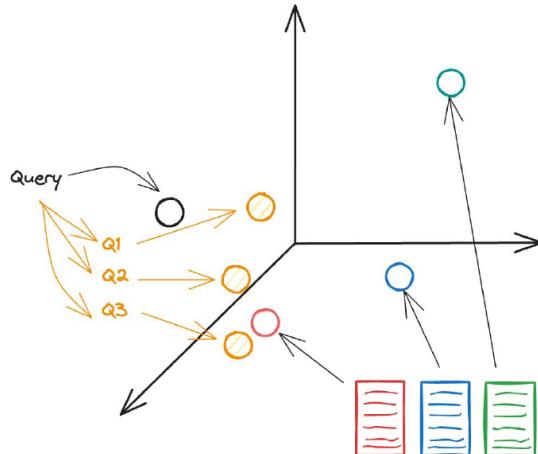


## Query改写

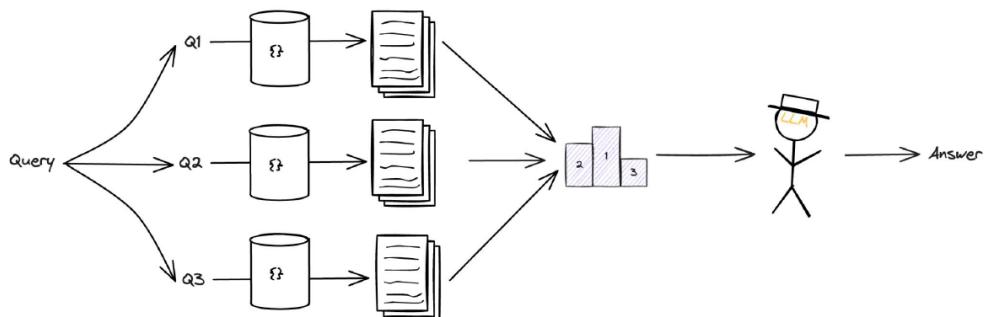
使Query和Document更容易匹配

举例：“怎么添加wx支付信息？”

- 等价的Query
  - “怎么添加微信支付信息？”
- 更抽象的Query
  - “如何补充支付信息？”
- 更具体的Query
  - “如何找到用户设置页面？”
  - “如何在用户设置页面中补充支付信息？”
  - “如何选择微信支付？”



## 等价Query：并行多查询



- 使用LLM生成等价query
- 多路并行召回
- 使用排序模型或启发式合并

**Query: What is the BLEU score of transformer on WMT 2014?**

- ➡ Q1: Can you provide the BLEU score achieved by the transformer model on the WMT 2014 dataset?
- ➡ Q2: What was the BLEU score obtained by the transformer architecture when evaluated on the WMT 2014 dataset?
- ➡ Q3: How well did the transformer model perform in terms of BLEU score on the WMT 2014 dataset?

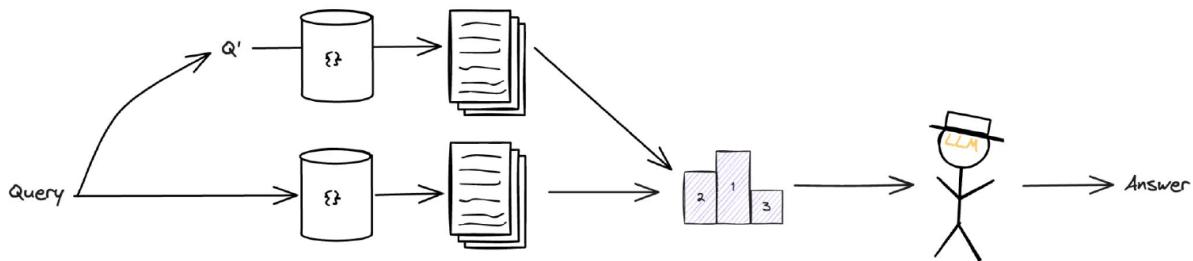
The BLEU score of the Transformer on WMT 2014 is 28.4, establishing a new state-of-the-art score. Additionally, the big model achieves a BLEU score of 41.0 on the English-to-French translation task

The BLEU score of the Transformer on the WMT 2014 English-to-German translation task is 28.4, and on the WMT 2014 English-to-French translation task, it is 41.8.

## 更抽象的Query: 回撤

### 更抽象的Query: 回撤

- 使用LLM对问题进行抽象
- 使用原始Query和改写后Query并行召回
- 适用于处理细颗粒的具体问题



**Query: How many GPUs are needed for training transformer?**

- ➡ Q: what hardware is required for training transformer?

Based on the provided context, the recommended number of GPUs for training transformer models efficiently is eight. The Transformer model achieved a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, which was a small fraction of the training costs of the best models from the literature.

To train the Transformer model, a total of eight GPUs are needed. This is based on the information provided in the documents, where it states that the models were trained on one machine with 8 NVIDIA P100 GPUs. The training process involved both base models and big models, with different hyperparameters and training steps. The big models, described at the bottom line of table 3, were trained for 300,000 steps over 3.5 days using the eight GPUs. This setup allowed for efficient training of the Transformer model, showcasing its ability to achieve state-of-the-art results in machine translation tasks while being more parallelizable and requiring significantly less time to train compared to other models based on recurrent or convolutional neural networks.

## 更具体的Query: 任务分解

### Query: What are encoder and decoder stacks in transformer?

Q1. How do encoder stacks work in a transformer model?

Q2. What is the purpose of decoder stacks in a transformer architecture?

Q3. Can you explain the differences between encoder and decoder stacks in a transformer?

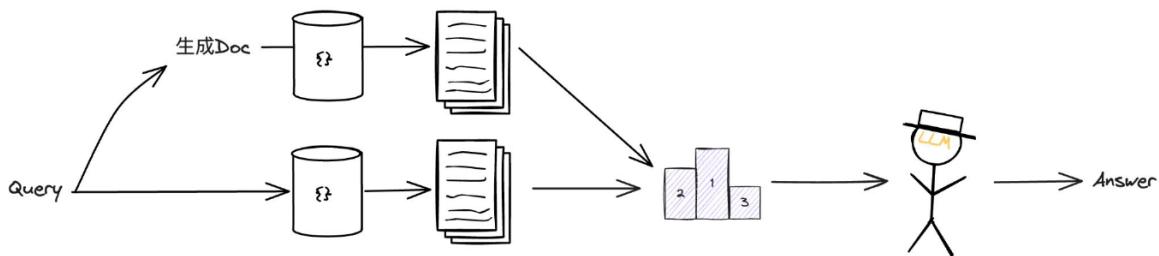
In the Transformer model, the encoder stack is composed of a stack of N=6 identical layers, each consisting of two sub-layers - a multi-head self-attention mechanism and a simple, position-wise fully connected feed-forward network. The decoder stack is also composed of a stack of N=6 identical layers, with each layer containing two sub-layers from the encoder stack and an additional third sub-layer for multi-head aligned RNNs or convolution.

In a transformer model, the encoder stack is composed of  $N = 6$  identical layers, each containing two sub-layers: a multi-head self-attention mechanism and a position-wise fully connected feed-forward network. ...

On the other hand, the decoder stack in a transformer also consists of  $N = 6$  identical layers. ... the decoder stack is responsible for producing the desired output sequence based on the continuous representations provided by the encoder.

## 根据Query生成Document: HyDE

使用查询到的文档+Query进行查询



- Document有更丰富的信息, 可以更精准的匹配
- 用LLM根据Query生成Document
- 使用原始Query和生成的Document并行召回

### Query: What learning rate is used when training transformer?

The learning rate used when training a transformer model is a crucial hyperparameter that significantly impacts the model's performance and convergence speed. In practice, a common learning rate schedule for training transformers involves using a warmup strategy followed by a decay schedule. During the warmup phase, the learning rate is gradually increased from a small value to the specified initial learning rate over a certain number of steps. This helps the model to stabilize and avoid large fluctuations in the early stages of training. After the warmup phase, the learning rate is typically decayed exponentially or according to a predefined schedule to fine-tune the model's parameters and improve convergence. The specific learning rate values and decay schedule can vary depending on the dataset, model architecture, and training objectives. ...

The learning rate used when training the transformer is `warmup_steps = 4000`, decreasing proportionally to the inverse square root of the step number. The base model of the transformer uses a training cost of  $3.3 \cdot 10^{18}$  FLOPs, while the big model uses a training cost of  $2.3 \cdot 10^{19}$  FLOPs. The base model applies a dropout rate of `Pdrop = 0.1` during training.

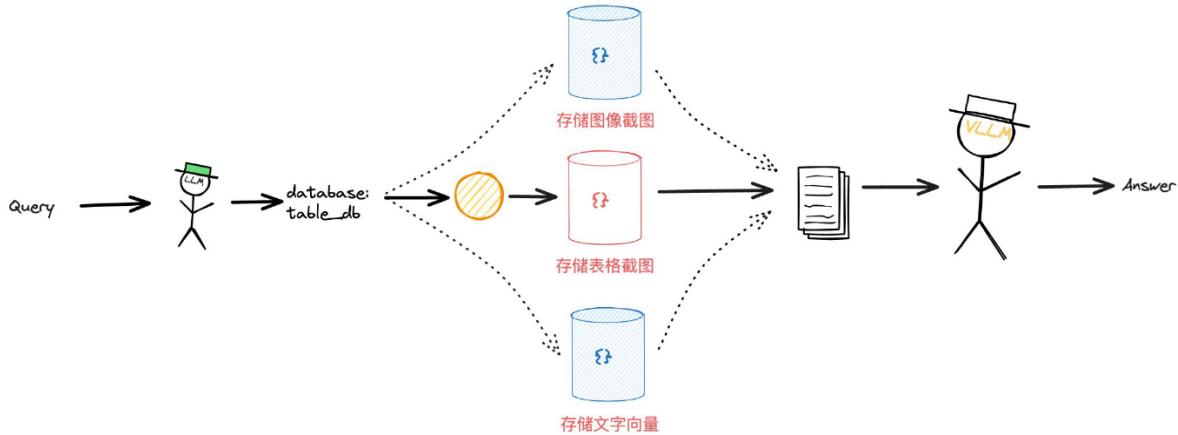
The learning rate used when training the Transformer is determined by the formula:

```
lrate = d^-0.5 * model * min(step_num^-0.5, step_num * warmup_steps^-1.5)
```

This formula increases the learning rate linearly for the first `warmup_steps` training steps and then decreases it proportionally to the inverse square root of the step number. The `warmup_steps` value used in this case is 4000.

## 路由器控制器

- 意图识别
- 根据用户意图决定召回路径



**Query: summarize the conclusion from table 1.**

```
datasource='table_images'
```

Layer Type	Complexity per Layer Operations	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_2(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

summarize the conclusion from table 1.

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types.  $n$  is the sequence length,  $d$  is the representation dimension,  $k$  is the kernel size of convolutions and  $r$  the size of the neighborhood in restricted self-attention.

Layer Type	Complexity per Layer Operations	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_2(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

Table 1 provides a comparative summary of various layer types used in neural network architectures, focusing on three key metrics: complexity per layer, sequential operations, and maximum path length. Here is the summarized content of the table:

...

-  $\backslash(r\backslash)$  represents the size of the neighborhood in restricted self-attention.

The table highlights the differences in computational complexity and efficiency across various neural network layers, illustrating trade-offs between computational resources and performance characteristics.

## Jina AI的RAG工具

### jina-embeddings-v2

23.10发布，全球第一个支持8k输入长度的开源向量模型

### jina-colbert-v1

第一款支持8k长度的colbert模型

### jina-reranker-v1

基于jina bert v2 支持8k上下文输入

### jina-clip-v1

正在开发...

### AIR-Bench

自动化多样信息检索评测基准

# 我们是否还需要RAG?



如果未来LLM能够精确的召回记忆，那么RAG将不再被需要。

不要高估6个月后的变化，  
不要低估18个月后的变化。  
在AI浪潮中，找到自己不变的价值。