

# 检索增强生成

《大语言模型》编写团队：李军毅

## ➤ 大模型局限性

➤ 幻觉

➤ 时效性

➤ 缺乏专业领域深度知识

➤ 参数化知识效率低

## ➤ 实际应用的需求

➤ 领域精准问答

➤ 数据更新频繁

➤ 生成内容可溯源

请介绍一下新任美国总统



截至2023年10月，美国总统是乔·拜登（Joe Biden）。他于2021年1月20日就职，成为美国第46任总统。乔·拜登在2020年总统选举中击败了时任总统唐纳德·特朗普（Donald Trump），成功当选。

乔·拜登简介：

- 全名：约瑟夫·罗宾内特·拜登（Joseph Robinette Biden Jr.）
- 出生日期：1942年11月20日
- 党派：民主党
- 政治经历：
  - 曾任美国副总统（2009-2017年，与巴拉克·奥巴马搭档）。
  - 曾任特拉华州联邦参议员（1973-2009年），是美国历史上最年轻的参议员之一。



截至2025年2月22日，美国的新任总统是唐纳德·特朗普（Donald J. Trump）。他在2024年11月的总统选举中击败了民主党候选人、时任副总统贺锦丽（Kamala Harris），成功当选第47任美国总统。这是特朗普的第二个总统任期，他曾于2017年至2021年担任第45任美国总统。

在2024年总统选举中，特朗普与他的竞选搭档、俄亥俄州参议员J·D·万斯（J.D. Vance）获得了312张选举人票，超过了贺锦丽和她的竞选搭档、明尼苏达州州长蒂姆·沃尔兹（Tim Walz）的226张选举人票。

值得注意的是，特朗普在78岁时再次当选，成为美国历史上最年长的总统当选人。

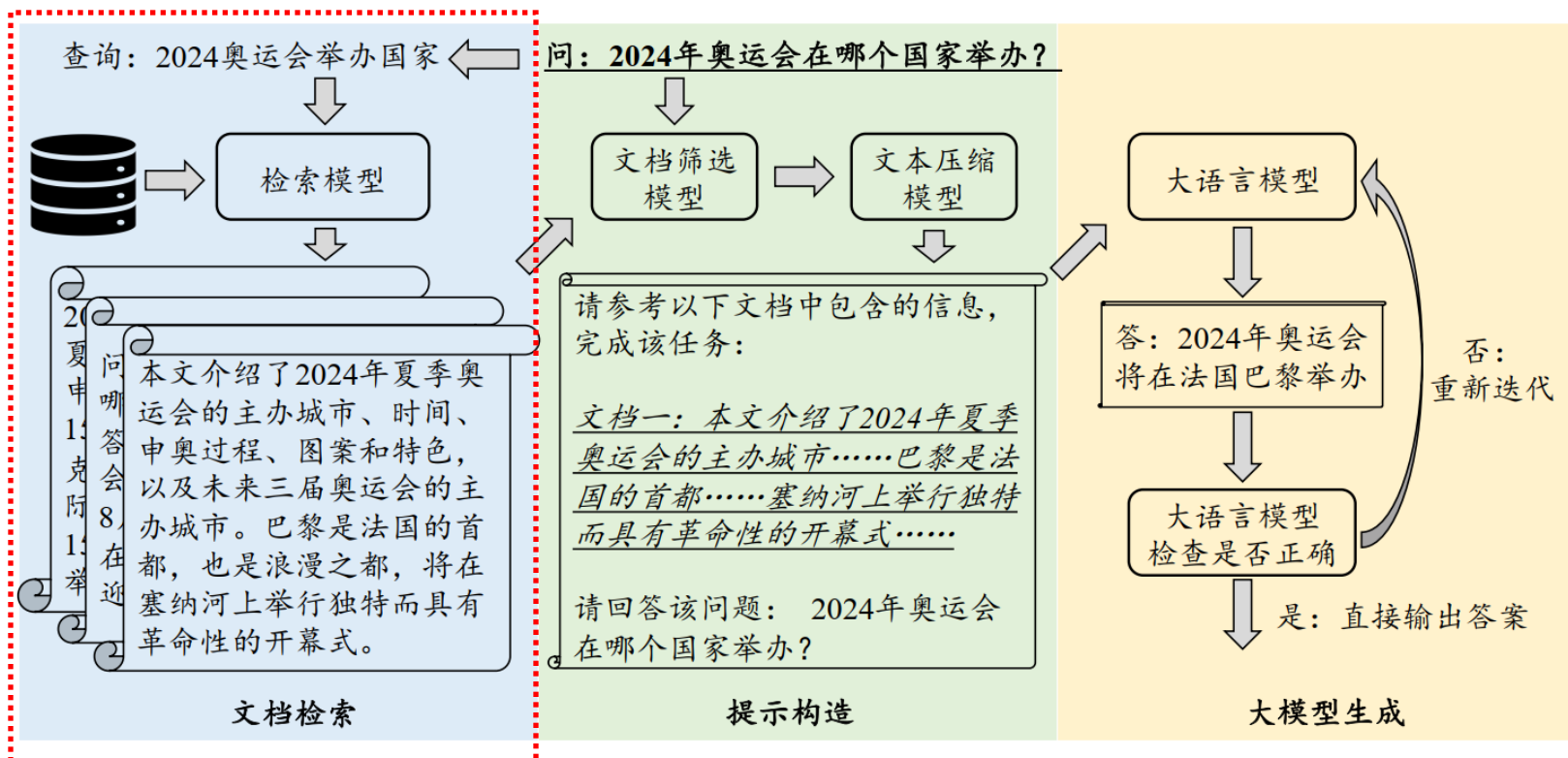
特朗普于2025年1月20日宣誓就职，正式开始他的第二个总统任期。由于当天华盛顿的严寒天气，典礼在国会大厦的圆形大厅内举行。在就职后，特朗普迅速签署了一系列行政命令，包括赦免在四年前国会山骚乱事件中被捕的1500人、宣布南部边境进入紧急状态以加强边境安全、以及宣布美国退出《巴黎协定》和世界卫生组织等。

这是自1893年格罗弗·克利夫兰（Grover Cleveland）以来，首次有美国总统在非连续任期内再次当选。特朗普的副总统J·D·万斯也于同日宣誓就职，成为美国第50任副总统。

# 基本流程

## ➤ 三个阶段

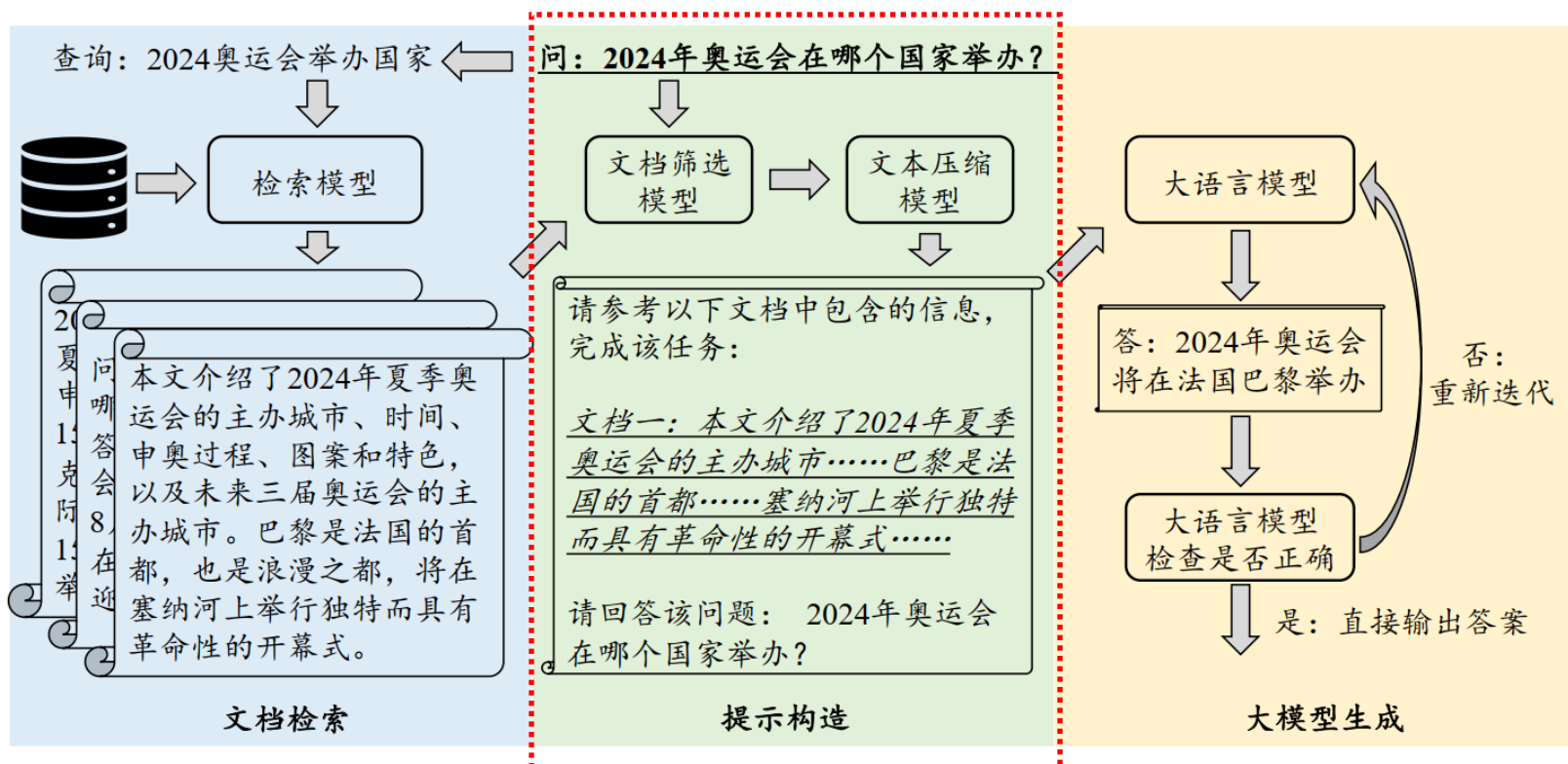
### ➤ 文档检索：从已有信息源中检索与问题相关的文档



# 基本流程

## ➤ 三个阶段

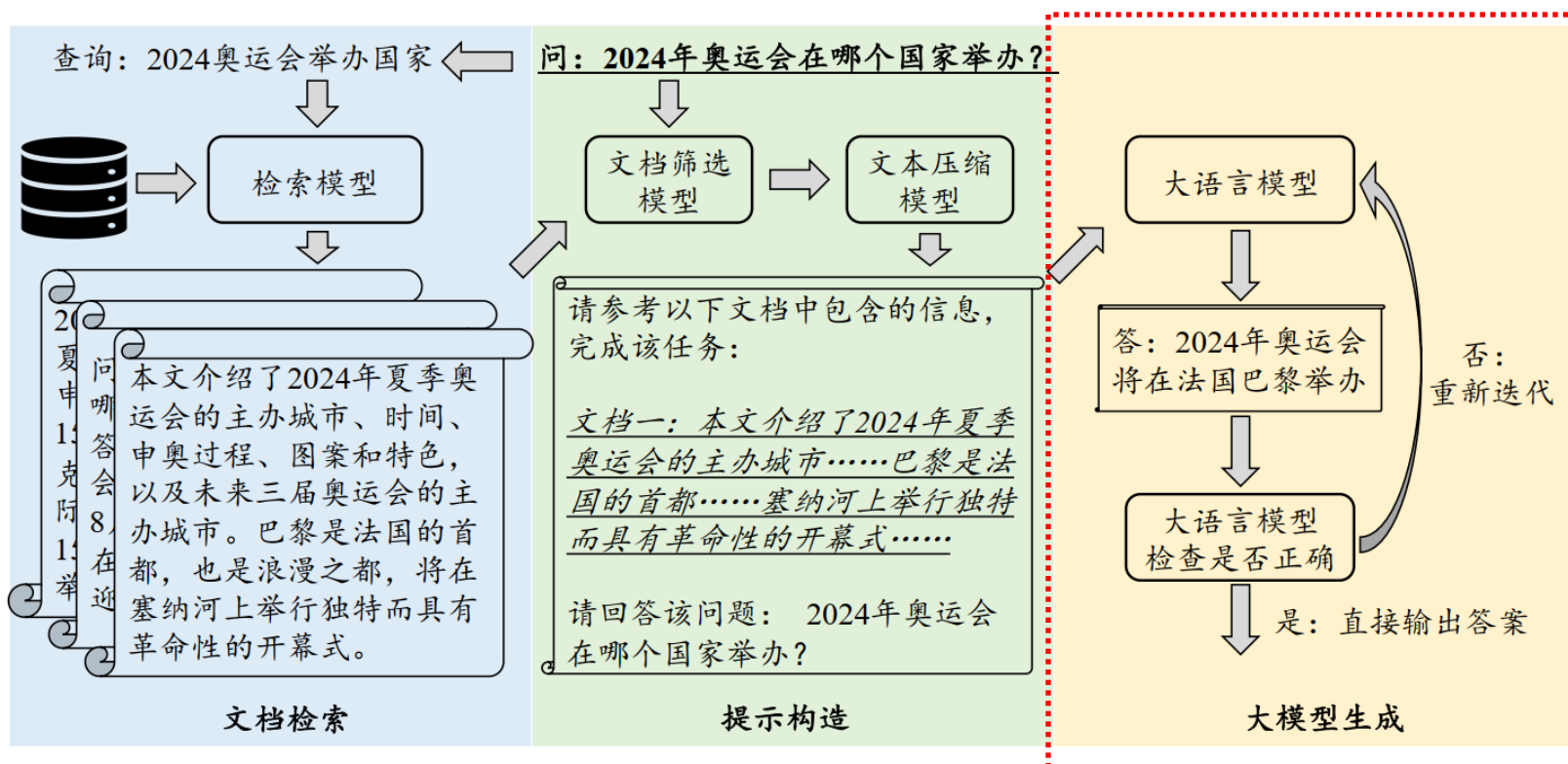
➤ 提示构造：将检索文档加入到输入提示中，并添加对应的任务描述



# 基本流程

## ➤ 三个阶段

### ➤ 大模型生成：模型基于检索增强的提示生成回答

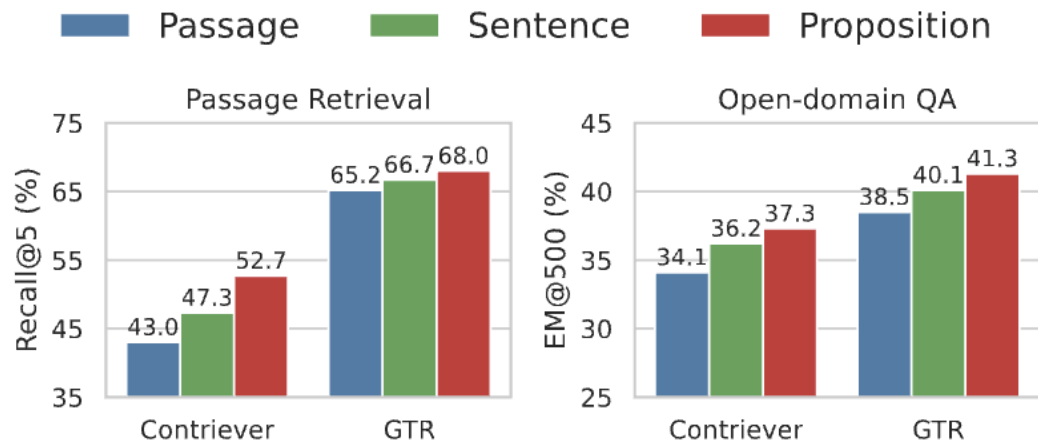


## ➤ 检索优化策略

### ➤ 检索数据源优化

- 检索单元：观点 (语义完整且内容相对独立的文本片段)，需要训练专有模型提取观点
- 对检索内容重组与改写，突出重要信息

文档级检索包含较多的无关信息



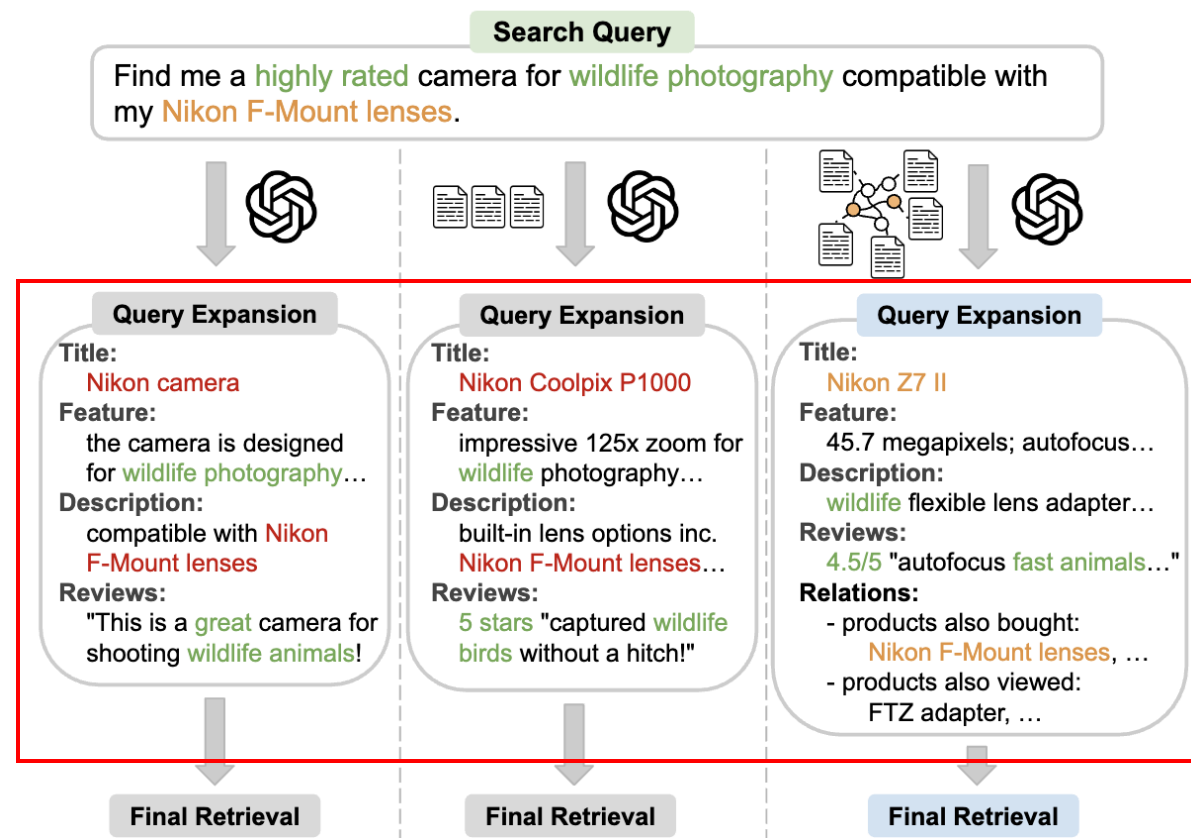
句子级检索召回更多相关信息，但检索时延较高

## ➤ 检索优化策略

### ➤ 查询优化

#### ➤ 查询扩展

- 为原始查询添加补充信息，对于复杂查询则可以分解为若干子查询
- 每个子查询分别进行检索，并将中间结果合并



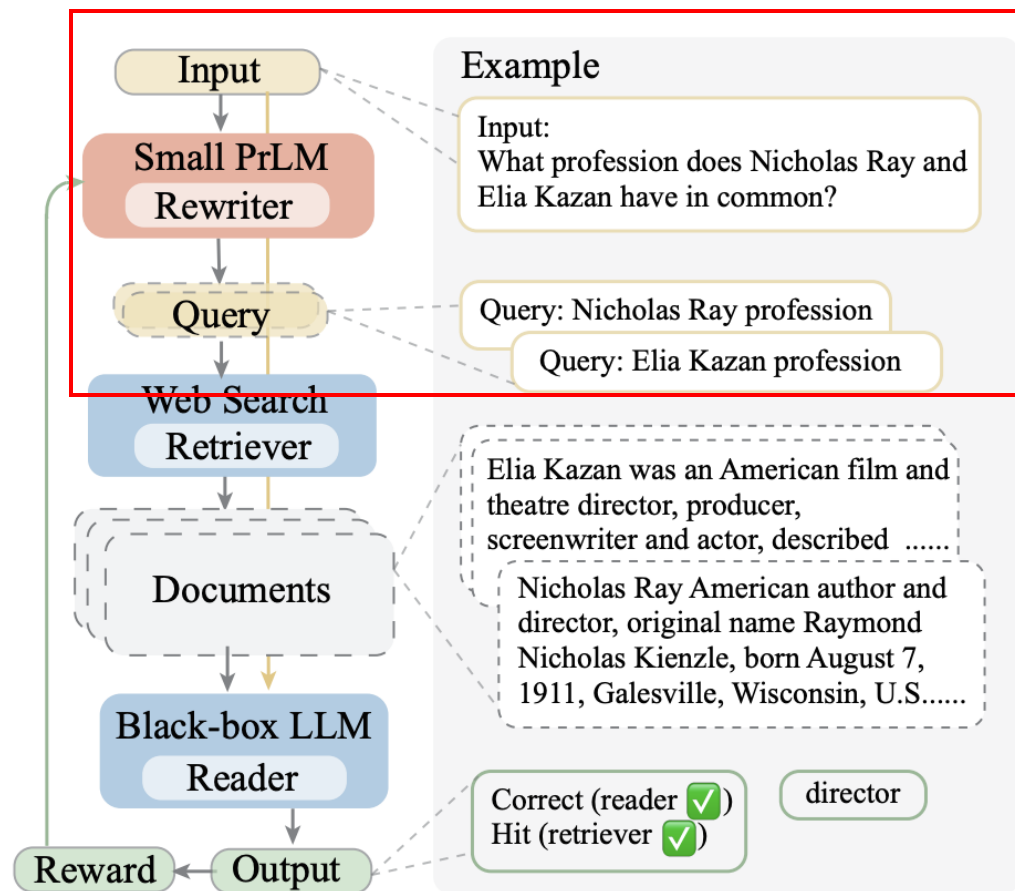


## ➤ 检索优化策略

### ➤ 查询优化

#### ➤ 查询重写

- 使关键信息更为突出，并消除可能存在的歧义
- 可以训练专门的小模型用于查询重写，减少计算开销

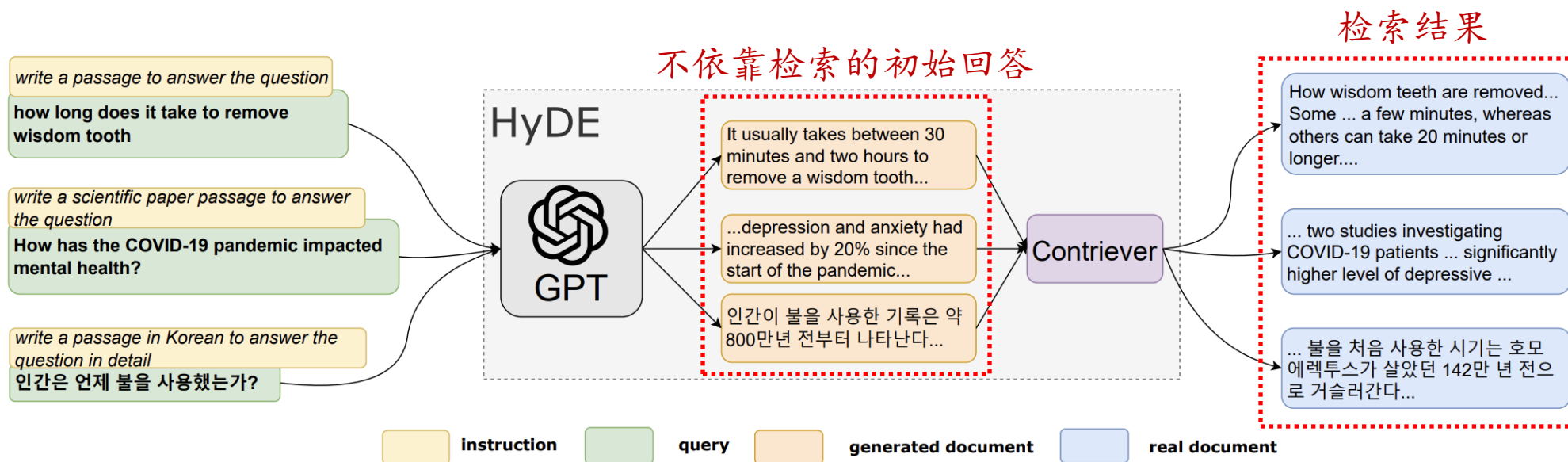




## ➤ 检索优化策略

### ➤ 查询优化

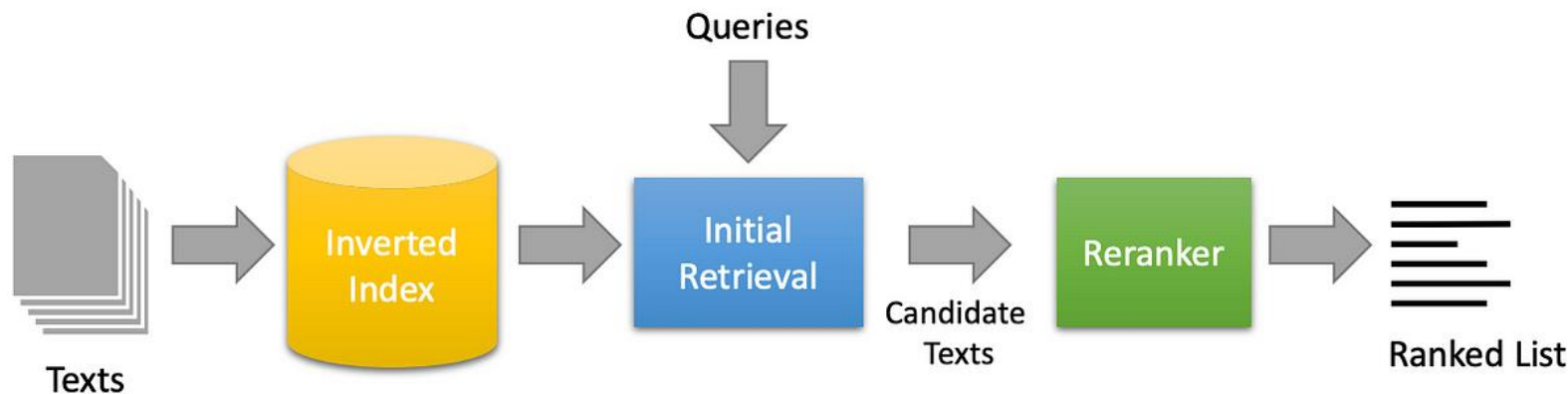
➤ 假设查询：用初始回答作为伪查询，检索得到新内容重新生成



## ➤ 提示优化策略

### ➤ 文档重排序

- 根据文档与输入间的相关程度进行重新排序，过滤掉低质量或不相关的文档，或将相关度较低的文档放置在提示中的非最优位置 (如中间和靠后位置)
- 可将大模型利用检索文档的情况转换为反馈分数用于训练重排序模型，提升重排序模型感知大模型对检索文档的利用能力



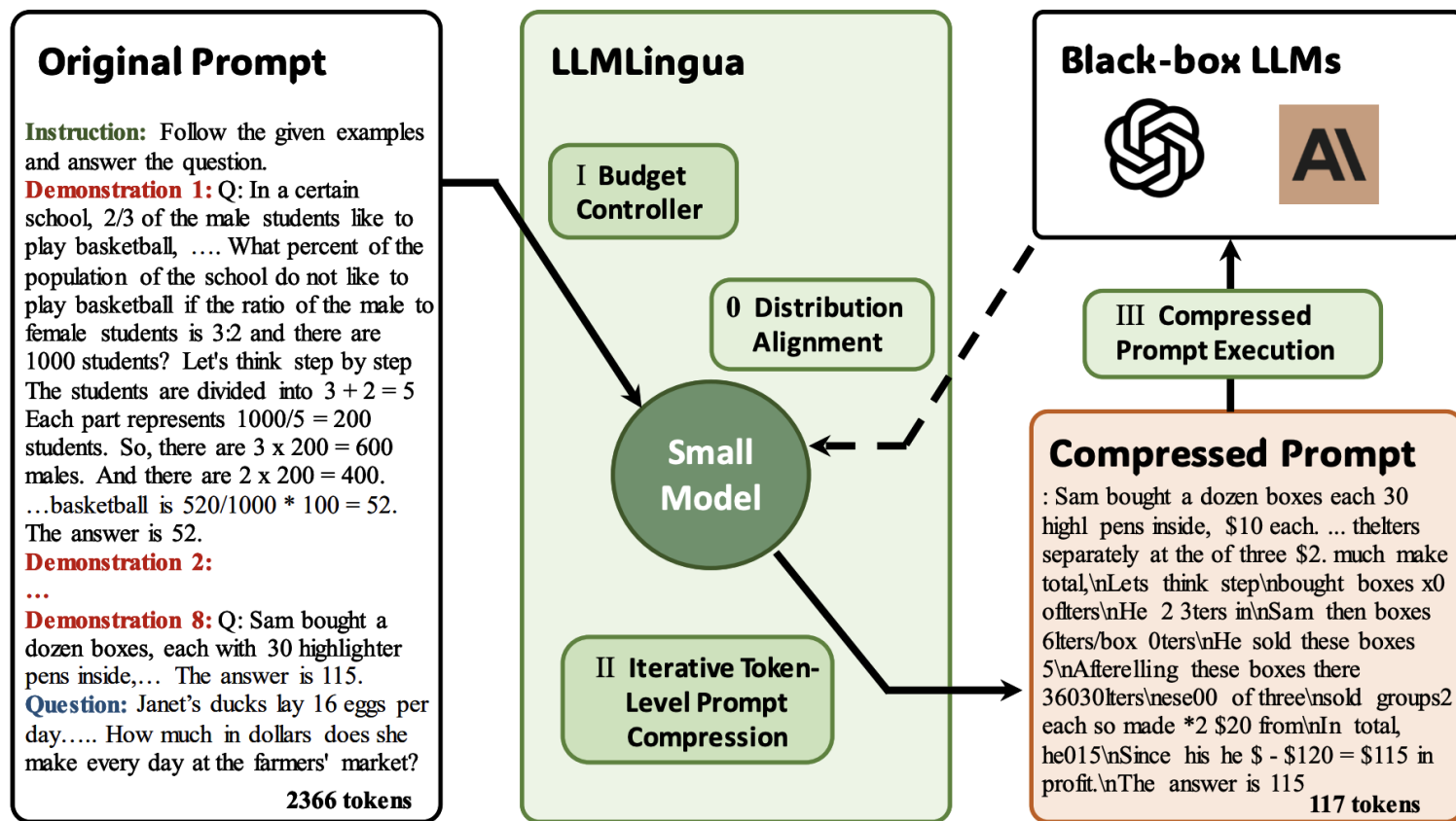
## ➤ 提示优化策略

### ➤ 上下文压缩

#### ➤ 自动摘要或信息抽取技术

#### ➤ 词元级别的压缩策略

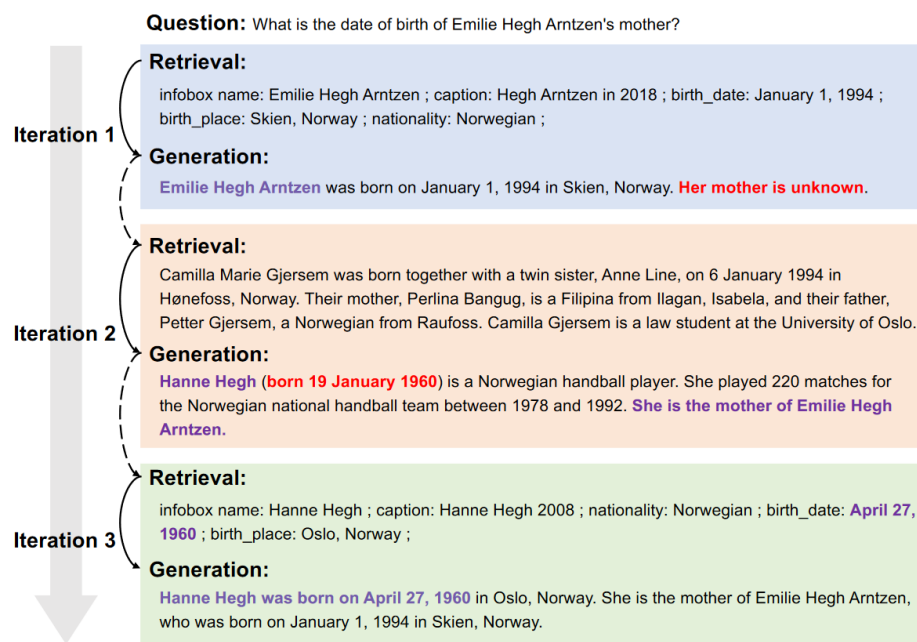
根据大模型对词元的预测  
概率选择出重要的词元



## ➤ 过程优化策略

### ➤ 迭代检索增强：迭代合并模型输出与初始查询，多次检索

### ➤ 检索也可与思维链结合，将思考步骤作为查询进行多步检索



**Identified domains:** factual (Wikidata, Wikipedia)

**Rationale 1:** First, the Argentine actor who directed El Tio Disparate is Fernando Birri.

**Retrieve (Wikidata) 1:** SELECT ?answer WHERE { wd:El Tio Disparate wdt:director ?answer . } -> Palito Ortega

**Retrieve (Wikipedia) 1:** Who directed El Tio Disparate? -> El Tio Disparate is directed by Palito Ortega.

**Corrected rationale 1:** the Argentine actor who directed El Tio Disparate is Palito Ortega.

**Rationale 2:** Second, Palito Ortega was born in 1941.

**Retrieve (Wikidata) 2:** SELECT ?answer WHERE { wd:Palito Ortega wdt:date of birth ?answer . } -> 8 March 1941

**Retrieve (Wikipedia) 2:** When was Palito Ortega born? -> Palito Ortega was born in 8 March 1941.

**Corrected rationale 2:** Palito Ortega was born in 8 March 1941.

**Corrected rationales:** First, the Argentine actor who directed El Tio Disparate is Palito Ortega. Second, Palito Ortega was born in 8 March 1941.

**The answer is 1941.**

## ➤ 过程优化策略

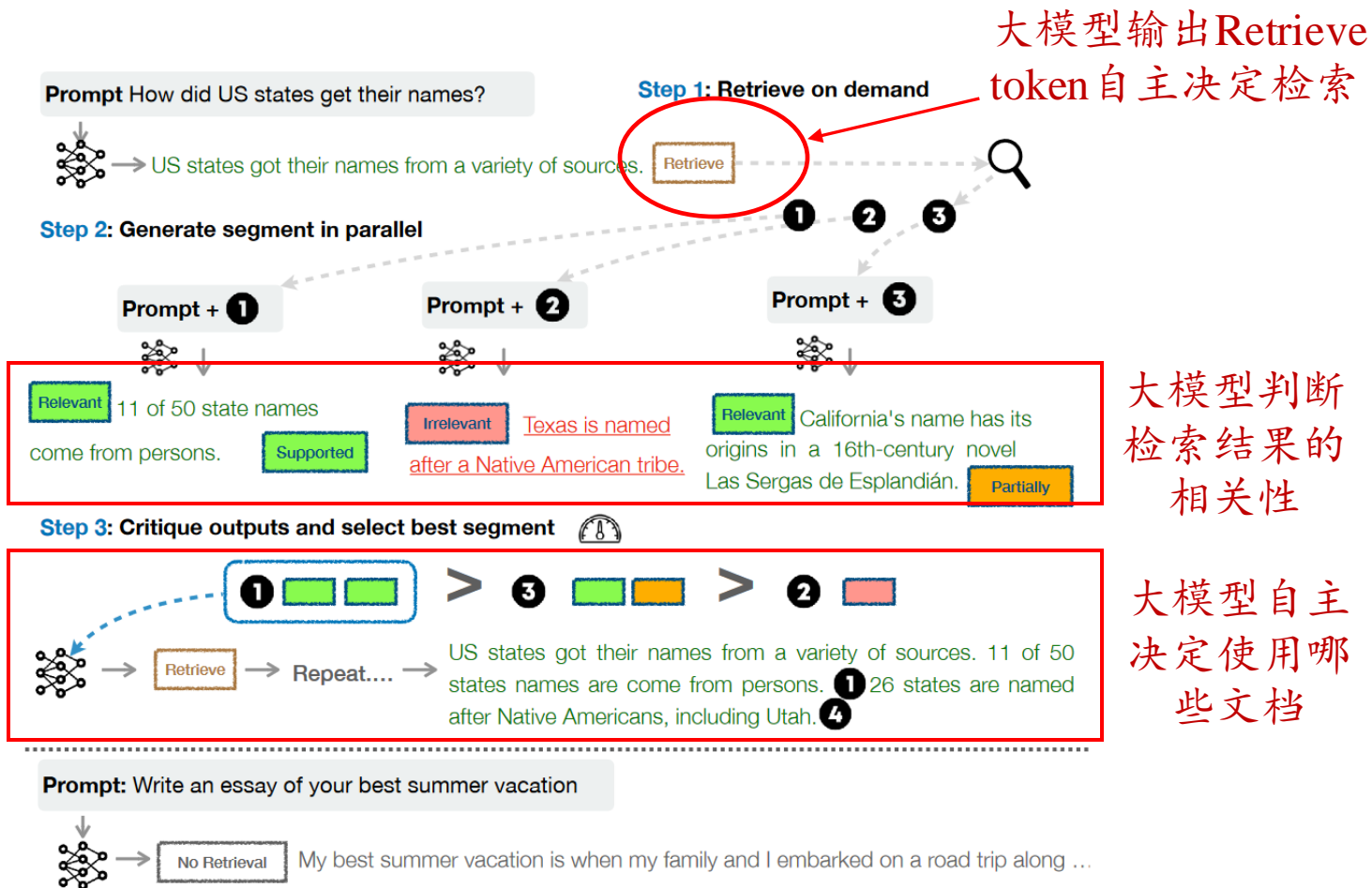
### ➤ 自适应检索增强

#### ➤ 大语言模型自主判断：

➤ 何时需要使用检索器

➤ 如何进行查询

➤ 如何处理检索结果



## ➤ 其他优化策略

### ➤ 指令微调增强策略

- 位置去偏 (相关文档置于不同位置)
- 无关信息过滤 (添加无用文档)

### ➤ 预训练增强策略

- 针对性加强检索生成能力

将文档首段作为查询，训练模型根据检索结果生成文档剩余内容

#### (b) Search-Augmented Prompt

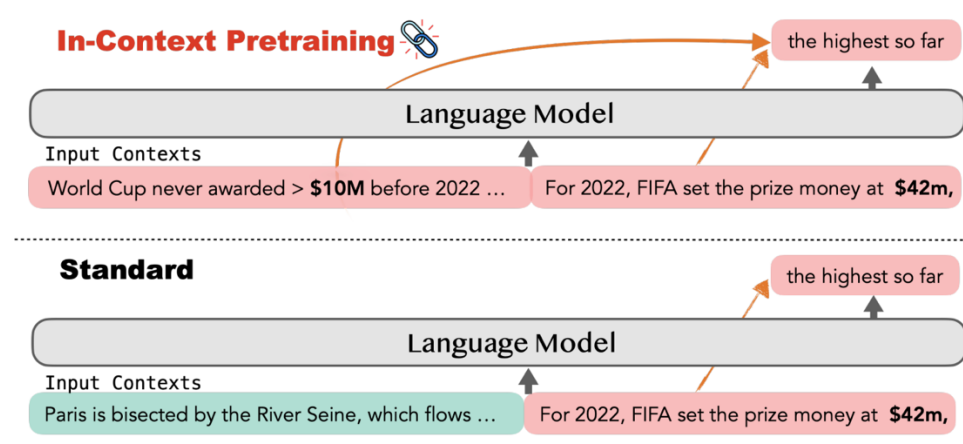
Below is an instruction that describes a task.  
Write a response that appropriately completes the request.

### Related Information:  
[Title 3]\n [Preview 3]  
[Title 2]\n [Preview 2]  
[Title 1]\n [Preview 1]

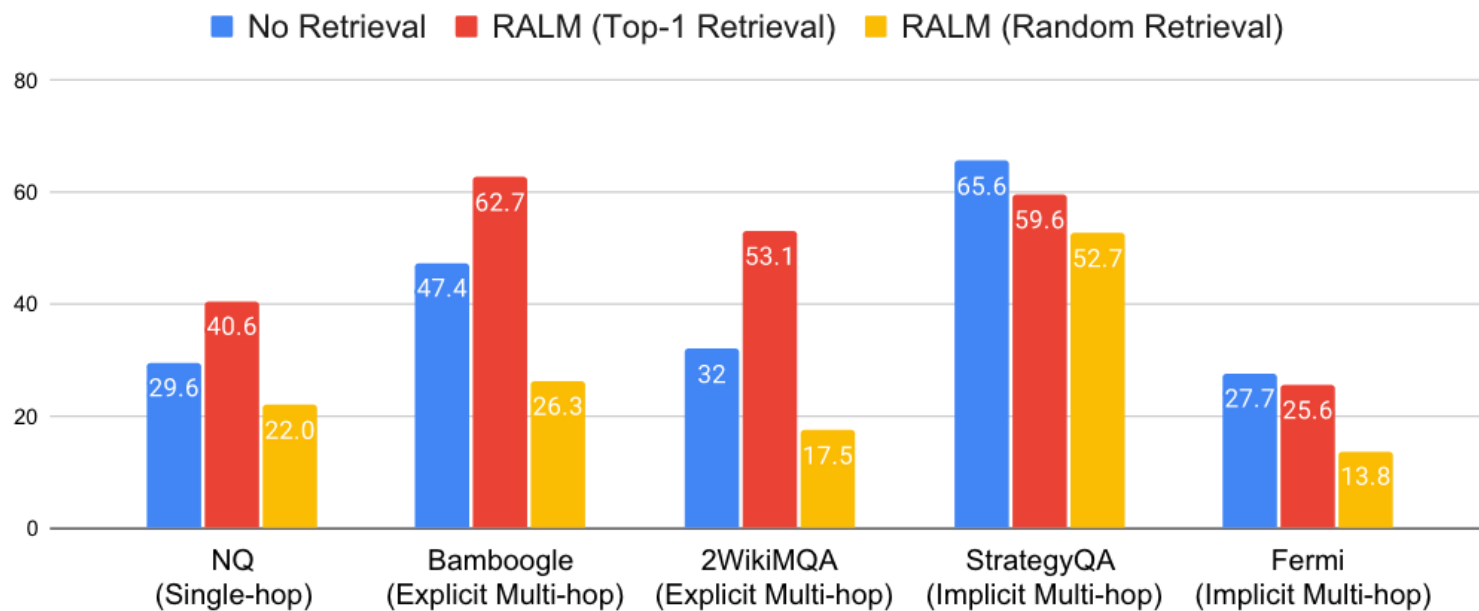
在指令中加入  
外部信息

### Instruction: [Instruction]  
### Input: [Input or None]

### Response:



- 检索一定能增强大模型吗?
- 无关信息可能会带来负面影响
- 大模型已经掌握的知识不需要再检索

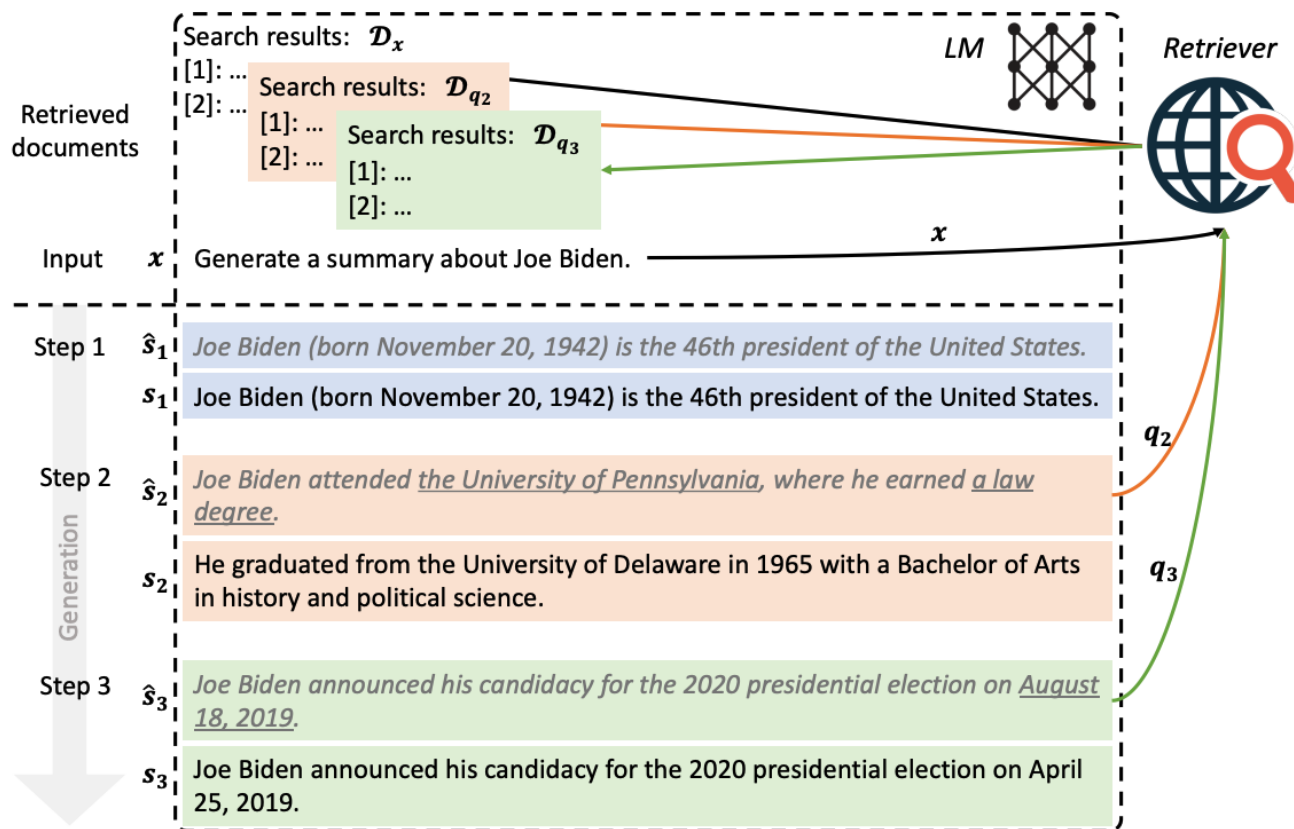




➤ 如何判断是否需要检索?

➤ 先生成后检索?

➤ 可行, 但成本高



➤ 如何判断是否需要检索?

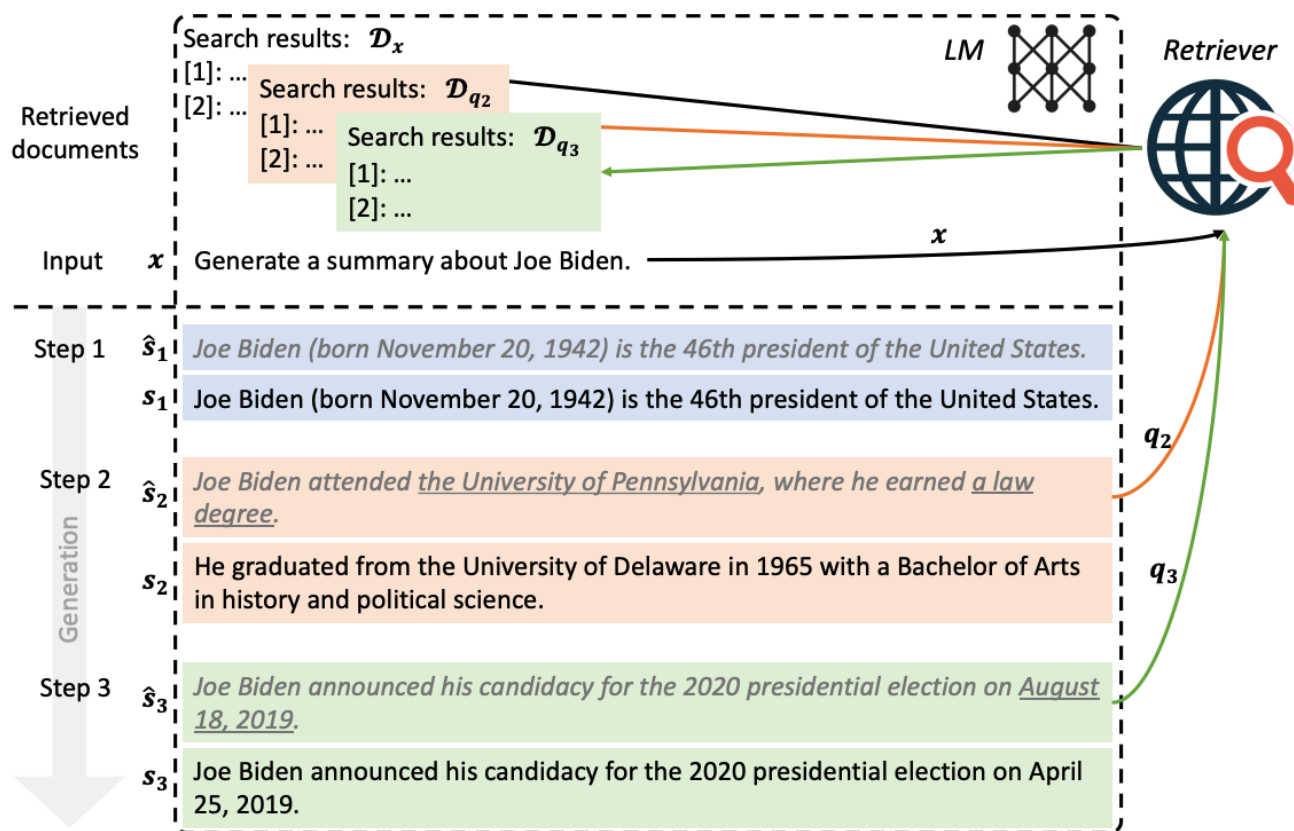
➤ 先生成后检索?

➤ 可行, 但成本高

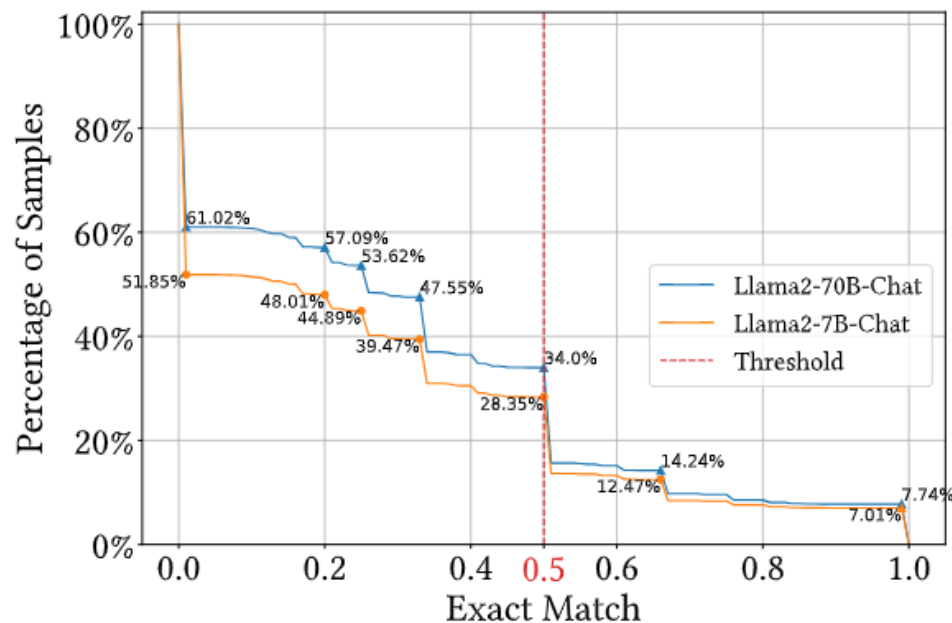
➤ 如何降低成本?

➤ 用小模型进行生成并判断

➤ 大小模型在知识程度上表现出相似性



- 目前的大模型预训练使用了类似的训练语料
- 不同模型在掌握程度高的知识是相当重合的
- 大小模型知识能力的差距主要体现在长尾知识上

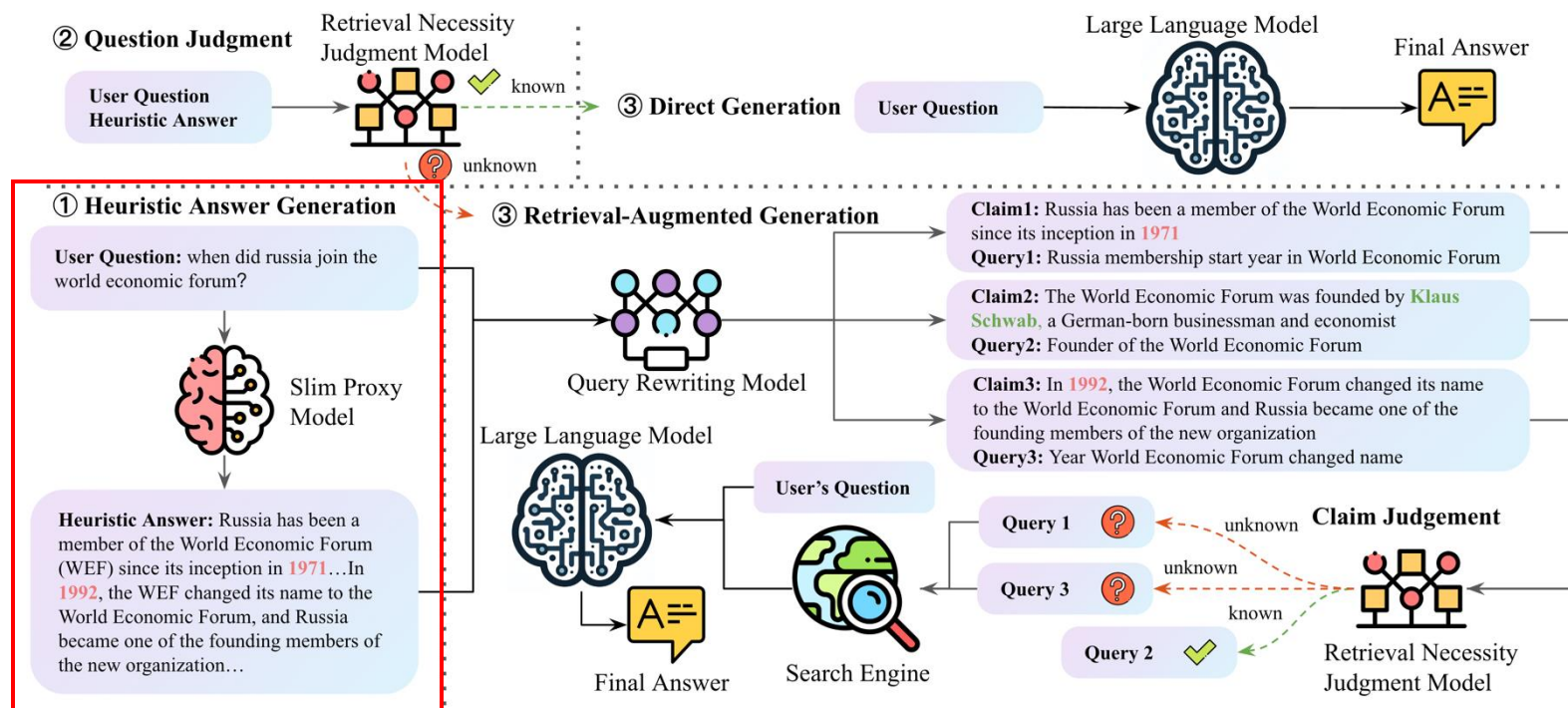


回答质量超过当前EM值的样本占总体样本的比例，越高说明当前模型能力越强

> 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

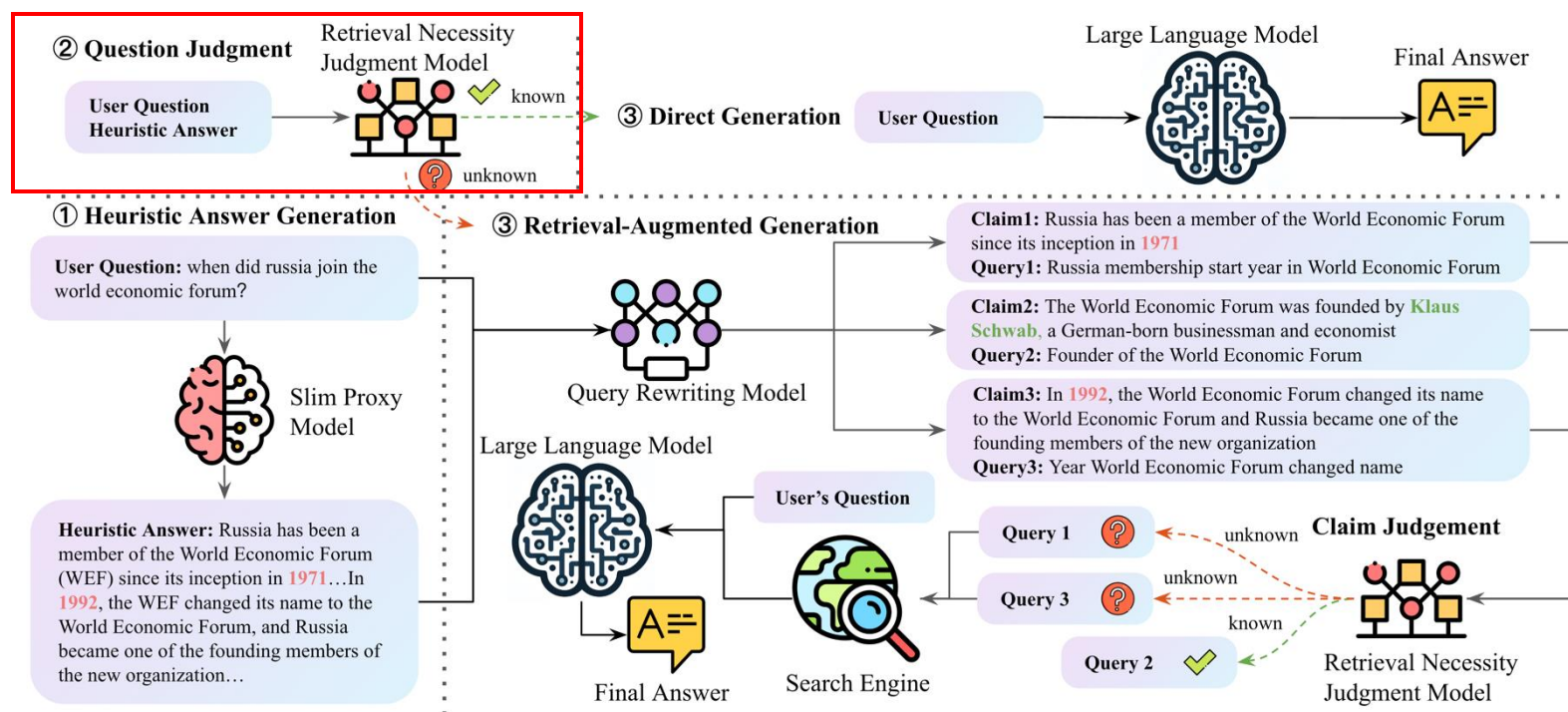
# 基于代理模型的检索时机判定方法

- 代理模型 (不做微调)
- 生成初始回答 (Heuristic Answer)



# 基于代理模型的检索时机判定方法

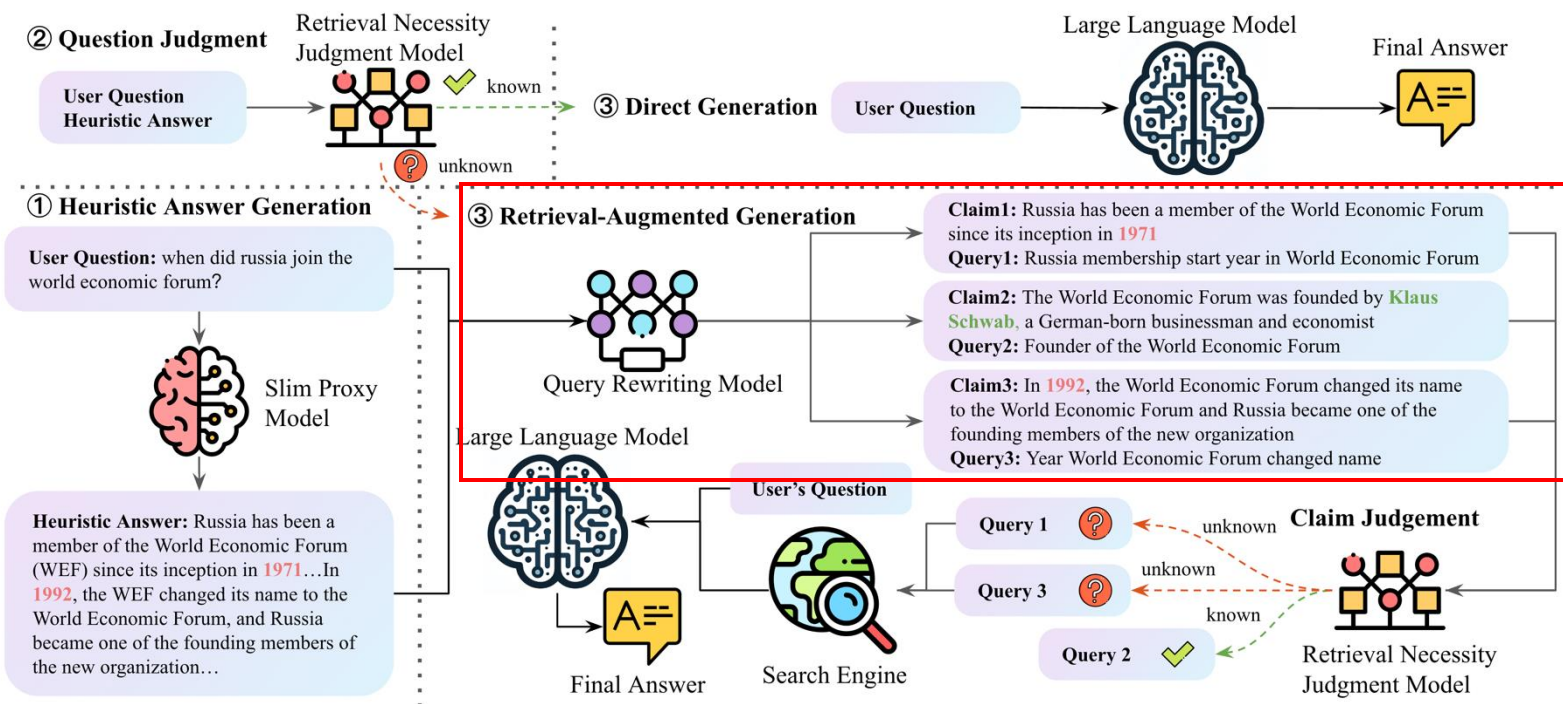
- 检索必要性判别模型 (微调后的 LLaMA-7B)
- 决定当前问题是否需要检索、某个子查询是否检索





# 基于代理模型的检索时机判定方法

- 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，生成若干个子查询



# 基于代理模型的检索时机判定方法



## ➤ 实验结果

- 相比于依赖大模型生成结果或概率的方法，SlimPLM 效果更好
- SlimPLM 只需要大模型推理一次，节省大量开销
- 与仅用用户输入进行判断的方法相比，SlimPLM 具备优势（代理模型解码结果具备指导性）

Method	#API	ASQA		NQ		Trivia-QA		MuSiQue	ELI5		
		EM	Hit@1	EM	Hit@1	EM	Hit@1	EM	ROUGE-1	ROUGE-2	ROUGE-L
Llama2-70B-Chat without Retrieval											
Vanilla Chat	1	29.68	62.50	40.49	55.00	27.44	90.75	11.50	28.66	4.88	14.27
CoT	1	26.21	54.50	35.36	48.75	23.50	79.00	11.50	28.12	4.73	14.06
Llama2-70B-Chat with Retrieval											
Direct RAG	1	27.63	58.00	42.40	56.00	28.07	92.25	10.50	28.61	4.76	<b>15.76</b>
FLARE	2.10	30.08	63.50	41.36	55.75	27.41	89.50	11.25	27.95	4.72	13.91
Self-Eval	2	29.45	60.75	42.15	55.75	27.58	91.50	10.25	28.70	4.83	15.39
Self-Ask	2.67	26.37	60.25	38.56	53.00	26.56	89.50	9.50	-	-	-
ITER-RETGEN	3	30.15	60.50	42.85	55.50	28.31	91.00	13.00	28.44	4.74	15.72
SKR-KNN	1	29.38	61.75	41.90	55.75	28.16	<b>92.25</b>	10.25	28.71	4.80	15.73
SlimPLM (Ours)	1	<b>30.73</b>	<b>65.00</b>	<b>47.43</b>	<b>62.25</b>	<b>28.35</b>	92.00	<b>13.00</b>	<b>29.97</b>	<b>5.61</b>	15.13
Qwen-72B-Chat without Retrieval											
Vanilla Chat	1	26.65	58.50	40.38	53.75	27.82	90.25	11.75	<b>30.61</b>	5.21	15.90
CoT	1	27.74	59.50	40.49	53.75	27.62	91.75	12.75	29.94	4.94	14.75
Qwen-72B-Chat with Retrieval											
Direct RAG	1	25.85	57.00	41.27	52.75	26.39	87.75	7.75	25.93	4.55	<b>16.74</b>
FLARE	2.29	27.68	59.00	40.89	54.50	27.10	88.50	<b>12.75</b>	30.31	5.20	15.77
Self-Eval	2	27.64	60.00	42.43	56.00	27.13	90.50	7.75	29.19	5.14	16.05
Self-Ask	2.76	22.82	52.25	36.16	49.25	25.29	87.50	9.75	-	-	-
ITER-RETGEN	3	<b>29.38</b>	61.50	43.51	57.50	27.16	89.75	12.25	26.15	4.41	16.52
SKR-KNN	1	28.08	61.50	43.08	56.00	26.38	88.50	11.25	27.29	4.75	16.31
SlimPLM (Ours)	1	27.97	<b>62.25</b>	<b>44.07</b>	<b>57.75</b>	<b>28.03</b>	<b>92.25</b>	9.75	29.56	<b>5.91</b>	16.36





谢谢