

# AI 算法进阶 (Advanced AI System)

@Tyler

3. 符号主义：知识库与自动推理 (Symbolism: Knowledge Bases and Automated Reasoning)

# 目录

- **模块一：符号表示和知识库**
- 模块二：规则引擎与自动推理



# 检索技术的演进：从搜索引擎到RAG系统

- 我们将一同回顾“检索”这项核心技术在不同时代的演化路径：
  - 从早期的**经典搜索引擎**，以关键词匹配和倒排索引为基础；
  - 到**现代推荐系统**，引入用户建模与内容理解，提升相关性与个性化；
  - 再到**大语言模型驱动的RAG系统**，实现语义理解与生成能力的深度融合。
- 本节的核心目标，是理解**技术如何持续回答同一个根本问题**：
  - 如何在海量信息中，**高效、精准地定位用户真正需要的内容**？
- 不同技术时代，答案虽异，**本质未变**。
- 正是对这一问题的持续追问，推动了AI系统从“**信息检索**”到“**知识注入**”的跨越。

# 回到原点：什么支撑了搜索引擎的速度与规模？



- 当你在谷歌或百度搜索“大语言模型”，系统是如何在**0.1 秒内**，从**数十亿网页**中准确返回结果的？
- 这个问题的答案，指向了一个古老而强大的核心机制：
- **倒排索引 (Inverted Index)** —— 自 20 世纪 60 年代以来，一直统治文本检索领域的关键数据结构。

# 正排 vs 倒排：如何重构“查找”的方向？

- 我们熟悉的书本目录是一种**正排索引**：  
*页码 → 内容*
- 而倒排索引的原理，则更接近新华字典的**部首查字法**：  
*字词 → 所在页码集合*
- 它彻底**颠倒了“文档-词语”之间的索引关系**，建立了一张映射表：
- **Token（词元） → Posting List（包含该词的文档ID集合）**

# 倒排索引的本质：将内容拆解为检索入口

- 每个词元（token）成为一个“入口”，指向所有包含它的文档。
- 这让搜索系统无需遍历全文，只需根据用户输入的关键词，快速定位候选文档集合，大幅提升检索效率。

类型	示例	类比说明
正排索引	文档ID → 内容	书本目录（按页找内容）
倒排索引	Token → 文档ID列表	字典部首查字（按字找页）
查询流程	用户输入关键词 → 快速定位包含词的文档集合	精确过滤

# 倒排索引的第一步：将文本变为“可索引”的单元

极客时间

训练营

- 我们以三个简短文档为例：
- **原始文档 (Documents)**
- 文档1: “大语言模型是人工智能的未来。”
- 文档2: “语言是人类沟通的桥梁。”
- 文档3: “深度学习模型驱动了人工智能的发展。”

# 倒排索引的第一步：将文本变为“可索引”的单元

- **Step 1: 分词 (Tokenization)**
- 将自然语言句子拆分为词元 (Token) 序列。  
(这里使用粒度较粗的中文分词作为示意)
- 文档1 → ["大", "语言", "模型", "是", "人工智能", "的", "未来"]
- 文档2 → ["语言", "是", "人类", "沟通", "的", "桥梁"]
- 文档3 → ["深度学习", "模型", "驱动", "了", "人工智能", "的", "发展"]



# 倒排索引的构建过程

- 完成分词后，我们可以建立“词元 → 文档ID列表”的映射关系，也就是**倒排索引表（Inverted Index Table）**：
- 倒排索引示意表：**

Token（词元）	Posting List（文档ID列表）
语言	[文档1, 文档2]
模型	[文档1, 文档3]
人工智能	[文档1, 文档3]
深度学习	[文档3]
沟通	[文档2]
桥梁	[文档2]
未来	[文档1]
驱动	[文档3]
发展	[文档3]

# 倒排索引的构建过程

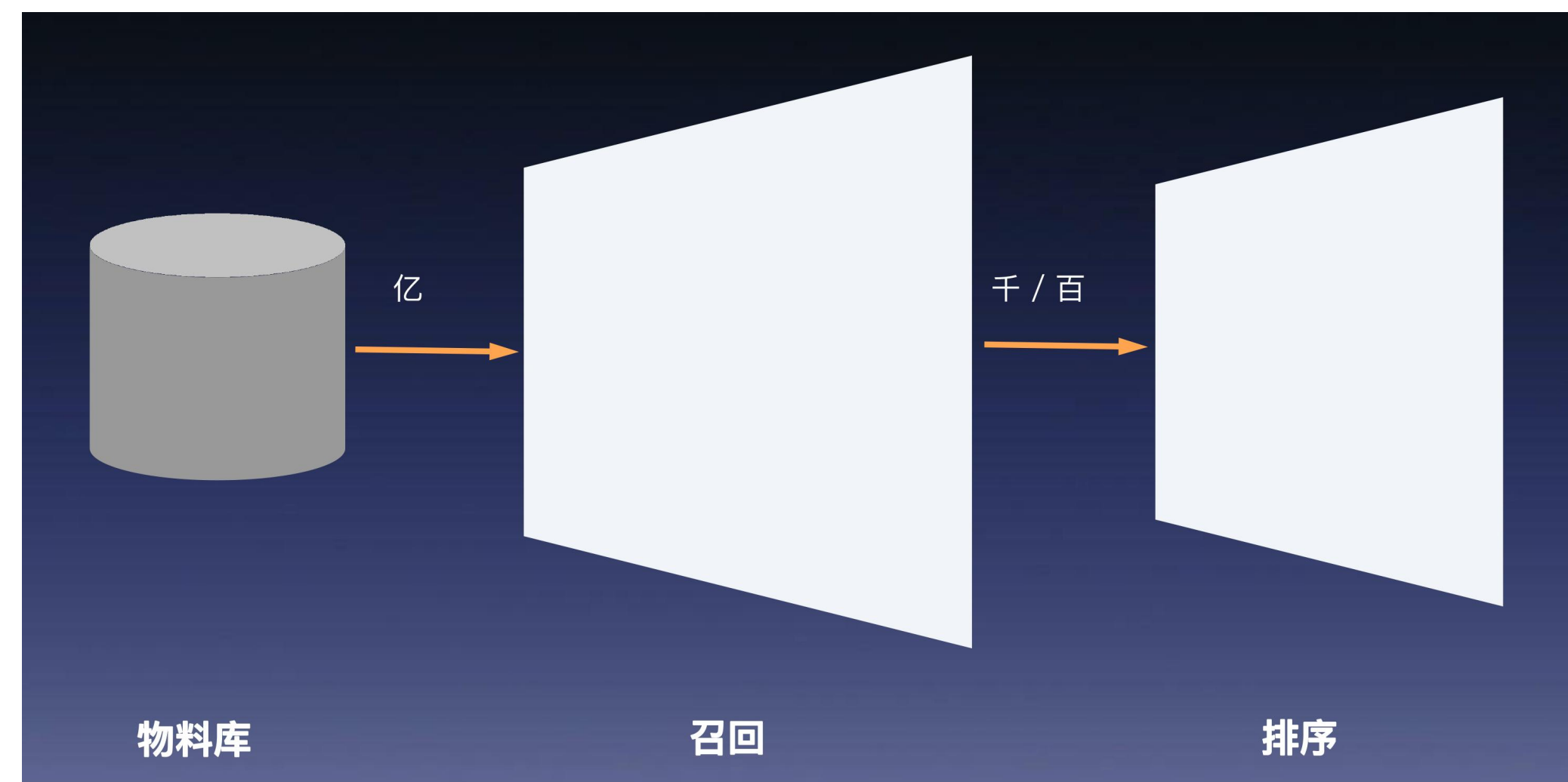
- **示例：搜索“语言 模型”时，系统如何响应？**
- 用户输入的查询词为：“语言 模型”。系统的处理流程如下：
- 根据倒排索引快速查找两个词的文档列表：
  - “语言” → [文档1, 文档2]
  - “模型” → [文档1, 文档3]
- 执行集合运算（取交集）：
  - $[文档1, 文档2] \cap [文档1, 文档3] = [文档1]$
- 最终结果：**文档1是唯一同时包含“语言”和“模型”的文档。**
- **核心优势：**无需遍历全文，只需查两个列表 → 运算复杂度大幅降低；

# 倒排索引：找到了，但还不够

- 倒排索引**完美解决了“找得到”的问题** —— 给定关键词，系统可以迅速定位包含它的所有文档。但这只是第一步。
- **新问题：哪个更相关？**
- 以用户搜索“语言”为例，命中以下两个文档：
- 文档1：“**大语言模型**是人工智能的未来。”
- 文档2：“**语言**是人类沟通的桥梁。”
- 现在的问题是：
- **当包含关键词的文档成千上万时，哪个文档应该排在前面？哪个才更符合用户的真实意图？**
- 我们需要的不仅是“包含”，而是“更相关”。

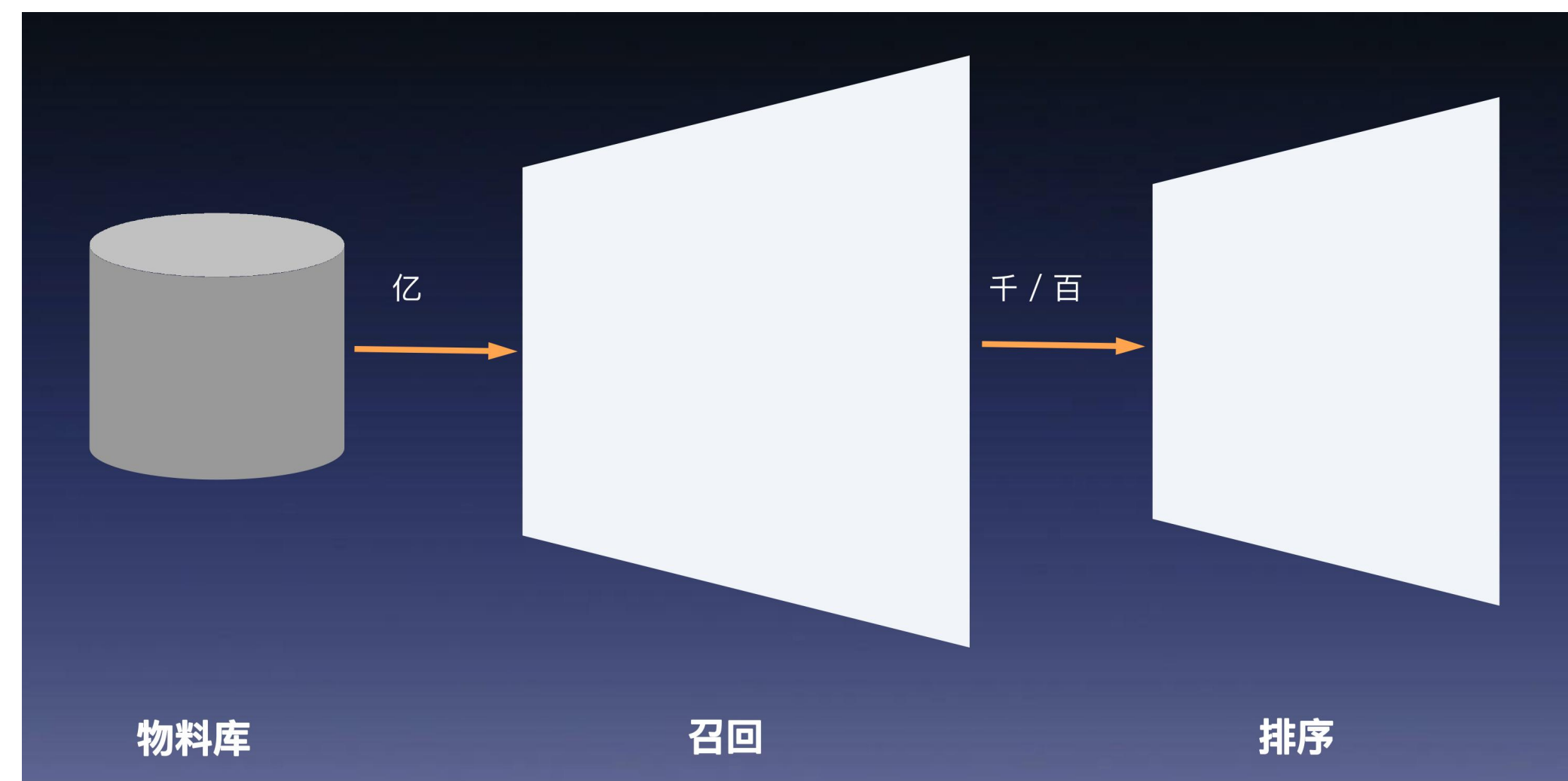
# 现代检索系统的主流范式：召回 + 排序架构

- 为兼顾**效率**与**精度**，当代大规模检索系统普遍采用两阶段结构：
- **1 召回阶段 (Recall Stage)**
- 使用高性能的检索算法（如倒排索引、向量召回等）
- 从数百万甚至数亿候选文档中**快速筛选出一个 Top-K 的初选集合**（通常几百到几千条）
- 目标：高覆盖、高召回率、低延迟
- 特点：粗粒度过滤，牺牲准确率换取速度和范围



# 现代检索系统的主流范式：召回 + 排序架构

- **2 排序阶段 (Ranking Stage)**
- 对召回得到的候选文档进行更精细的特征计算与相关性打分，使用排序模型进行最终打分排序
- 目标：精准判断用户意图匹配度
- 特点：精模型、慢计算、小规模、高准确率
- **这种分层结构将“找得全”和“排得准”解耦，**
- **是搜索、推荐、问答系统中的通用工程架构模板。**





# 倒排索引：找到了，但还不够

- 引出关键机制：TF-IDF 排序 (Term Frequency — Inverse Document Frequency)
- 搜索引擎时代，正是为了回答“谁更相关”这个问题，提出了经典的排序算法：

$$\text{TF-IDF} = \text{词频} \times \text{逆文档频率}$$

- 它不仅考虑某个词在一篇文档中出现了几次 (TF)，还考虑这个词在整个文档库中有多“稀有” (IDF)：
  - 出现频率高 → 说明对该文档重要；
  - 出现得“越稀有” → 区分度越强，信息量越大。
- 倒排索引解决的是“有没有”，TF-IDF 解决的是“哪一个更好”。

# TF-IDF：搜索引擎时代的相关性评分公式

- 我们已经能“找出包含关键词的文档”，但如何判断**哪个文档更相关**？
- 核心思想只有一句话：**一个词的重要性，取决于它在这篇文档中的“特殊性”。**
- **TF: Term Frequency (词频)**
- **目的：**为关键词在文档中的出现次数增加权重。
- **直觉：**一个词在**某篇文档中出现得越多**，这篇文档就**越有可能和它相关**。
- **例子：**如果一篇文章反复出现“乔布斯”，它很可能是在谈论苹果公司或乔布斯的传记。
- **IDF: Inverse Document Frequency (逆文档频率)**
- **目的：**当用户查询被分词成多个关键词时，**哪个查询词更重要、更有“区分力”**，对排序贡献更大。
- **直觉：**如果一个词在**所有文档中都很常见**，那它信息量低（如“的”、“一个”）；而如果一个词在全库中**很少出现**，那它更能体现一篇文档的独特性。
- **例子：**“人工智能”比“技术”更稀有，因此在一篇文章中出现时，其区分度更高。

# TF-IDF：搜索引擎时代的相关性评分公式

- **排序分数 (查询Q, 文档D) =  $\Sigma \{TF-IDF(\text{词}t, \text{文档}D)\}$**
- 当一个词在某篇文档中：查询词出现得多（高 TF）且在整个语料库中很少见（高 IDF）
- 那么这个词对该文档来说就**高度相关**，这篇文档也应在搜索结果中排得更靠前。
- TF-IDF 是搜索引擎时代衡量文本相关性的基础算法。
- 它帮助系统不仅找到文档，更学会了判断：**哪一篇才是用户真正想要的？**

# TF-IDF解决了“相关性”，却没能回答“为什么点它”



- TF-IDF 帮助我们衡量**文本之间的相关性**，但现实世界中的用户行为，远比“文本匹配”复杂得多。
- **用户点击的内容，并不总是“文本最相关”的那一篇。**
- **思考几个真实场景：**
- **时效性 (Timeliness)**
- 搜索：“苹果发布会”
- TF-IDF 找到了 2017、2019、2024 年的文章.....
- 但用户**只想看昨天那场发布会的最新报道。**

# TF-IDF解决了“相关性”，却没能回答“为什么点它” 极客时间

- **权威性 (Authority)**
- 搜索：“感冒怎么办”
- TF-IDF 可能把关键词最多的博客文章排在最前
- 但用户更信任**来自协和医院或CDC官网的权威内容**。
  
- **用户偏好 (Personalization)**
- 搜索：“机器学习入门”
- 初学者倾向选择**通俗易懂的博客教程**
- 专业人士更偏好**一篇系统综述或学术论文**



# TF-IDF解决了“相关性”，却没能回答“为什么点它”



- **结论：相关性只是冰山一角**
- 真正的“好”结果，不只是“词对上了”，而是：
- **综合考虑了文本匹配、时效性、权威性、点击反馈、用户历史偏好等上百种信号**的结果。
- TF-IDF 等基于规则的算法，只能在一个维度上打分，在复杂、多维的用户需求面前，**力不从心**。
- 我们需要一个**能够自动融合多种信号、学习判断标准**的“智能裁判”：
- **机器学习模型，替代手工规则，成为排序系统的核心引擎**。
- 它不再死记规则，而是**从点击数据中学“哪些结果更可能被点击”**，最终演化出以用户反馈为监督信号的“排序学习系统”。

# 核心转变：排序变成监督学习问题

- 在前面的传统搜索中，“排序”曾经是靠规则打分、靠权重组合。
- 但在后来，广泛采用了我们在上一节课学习的方法：**将排序问题转化为一个“二分类”的监督学习问题**。也就是我们熟悉的：**点击率预估（CTR Prediction）问题**。
- **通过学习一个模型，去预测：某个文档被用户点击的概率有多大，并用这个概率排序。**

# 核心转变： 排序变成监督学习问题

- 输入特征 (X)：来自多个维度的信号融合
- 每个候选文档会被编码为**特征**，通常包含上百个维度：

特征类型	示例说明
权威性	是否来自权威域名
时效性	文档发布时间与当前的时间差
用户行为	历史点击率 (CTR)
页面结构	URL 长度、是否含关键词
设备匹配	查询设备与页面兼容性等
...	...

- 这些特征共同刻画了文档是否“值得点击”的可能性。

# 核心转变：排序变成监督学习问题

- **监督信号 (y)：用户行为就是标签**
- 若用户看了，并且点击了该文档  $\rightarrow y=1$
- 若用户看了，但是没有点击它  $\rightarrow y=0$
- 成千上万的搜索浏览行为，构成了一个**真实世界的训练集**。

# 核心转变：排序变成监督学习问题

- **训练模型：逻辑回归的自然契合**
- 我们可以直接使用上节课学习的**逻辑回归（Logistic Regression）**，对上述特征向量进行建模。
- 这个模型学习到了：**在什么条件下，用户更可能点击某个搜索结果。**
- 传统：手工打分（TF-IDF）
- 现代：用户点击行为监督 + 机器学习预测点击概率。
- **排序系统不再由工程师指定“什么重要”，而是由用户的行为告诉模型“什么重要”。**



# 排序系统的现代工作流程：从规则到学习

- 我们来看，一个现代搜索引擎如何使用机器学习进行排序：
- **Step-by-Step：机器学习排序的三步流程**
- **1 特征提取**
- 对召回阶段找到的**Top 1000 个候选网页**，系统为每个网页构造一组多维特征：
  - 权威性（来源网站）
  - 时效性（发布时间）
  - 用户行为（历史点击率）
  - 页面结构（标题长度、URL深度）
  - .....
- 每一个网页变成一个 **高维特征向量**  $x$

# 排序系统的现代工作流程：从规则到学习

- **2 CTR 预估（点击率预测）**
- 将每个网页的特征向量输入到训练好的**逻辑回归模型**中：
- 模型输出的是一个点击概率预测值（CTR）：越高 → 用户越可能点击
  
- **3 排序并展示**
- 对所有候选网页按 CTR 值从高到低排序：
- Top 1 排在搜索结果的最上面
- 依次向下排列
- 最终展示给用户的，是 **点击概率最高的一批内容**

# 排序系统的现代工作流程：从规则到学习

- **从这一步开始，排序任务被机器学习正式“接管”**
  - 这标志着搜索技术的核心范式从此发生跃迁：
  - **从“人工规则排序” → “用户行为驱动模型排序”**
  - 而这一套流程，也成为后续所有复杂排序模型的雏形：
  - 模型更强 → 表达能力增强
  - 特征更丰富 → 信息融合更全面
  - 用户体验更优 → 持续闭环学习
- 
- **排序 = 特征工程 + 监督学习 + 实时决策**
  - 这就是推荐系统、广告系统和搜索引擎的共同排序底座。

# 范式转变：从搜索引擎到推荐系统

- 在搜索引擎时代，人与系统的关系是“我问，你答”：
- 用户通过关键词表达明确意图，系统负责尽快、准确地满足这个意图。
- 这是一个**基于需求响应的匹配系统**。
  
- **而现在，我们进入了一个全新的时代：推荐系统（Recommender System）**
- 这就是抖音、淘宝、小红书们的时代。
- **用户不再输入搜索词**，只是在“逛”——刷视频、滑商品、浏览内容。
- 系统不再等用户提出请求，而是需要**主动出击、实时推荐**。

# 范式转变：从搜索引擎到推荐系统

- 搜索 vs 推荐：核心差异

维度	搜索引擎	推荐系统
用户输入	明确的搜索词	无输入或仅有隐式行为（浏览、点击）
系统角色	被动响应，满足明确查询	主动预测，引导兴趣探索
排序目标	精确匹配用户的当前意图	理解用户的长期偏好和短期兴趣
技术动因	关键词匹配 → 排序优化	用户建模 → 意图预测 → 实时召回与排序

- 技术结论：从“理解内容”到“理解人”
- 搜索系统的召回逻辑是：输入这个**关键词**的背后想要获取什么信息？  
推荐系统的召回逻辑则是：此刻使用我的App的**这个人**可能对什么感兴趣？
- 召回的单位从“匹配关键词”升级为“匹配人心”。
- 这意味着我们不再围绕**内容本身做索引**，而是围绕**用户兴趣建模**。



# 推荐系统的信条：理解用户，才能推荐得准

- 推荐系统的核心思想非常直接：**如果我们能深刻理解一个用户是谁，就能预测他会喜欢什么。**
- 这就是推荐系统与搜索系统的根本区别：不是理解查询，而是理解人。
- 查询从**我要什么（路人）**，变成了**我是谁（会员制VIP）**。
- **用户画像（User Profile）：数字世界中的“用户建模”**
- 所谓用户画像，就是为每一位用户贴上一系列结构化的“标签”或“特征”，它是对用户的**属性、兴趣、行为模式**的综合刻画。
- 用户画像 = 内容召回的线索来源 + 排序系统的决策依据。

# 推荐系统的信条：理解用户，才能推荐得准

- 画像信息来源：静态 + 动态，全方位建模用户
- ♀ 静态属性 (Static Attributes)
- 长期不变或变化缓慢，描绘“你是谁”
- 性别：女
- 年龄段：25–30 岁
- 所在城市：北京
- 手机型号：iPhone 15 Pro

# 推荐系统的信条：理解用户，才能推荐得准

- **动态属性 (Behavioral Attributes)**
- 由用户行为实时沉淀，刻画“你现在可能在想什么”
- **长期兴趣 (长期行为统计)**
  - 典型标签：#健身 #美食 #旅游
  - 来源：过去半年浏览 / 点击 / 购买数据
- **短期意图 (即时兴趣波动)**
  - 典型标签：#徒步鞋 #防晒霜
  - 来源：最近三天的搜索 / 加购 / 收藏行为
- **活跃时段**：工作日晚间高活跃
- **消费能力**：高客单价，付费转化率高

# 推荐系统的信条：理解用户，才能推荐得准

- 为什么用户画像如此关键？
  - 因为推荐系统不再是“谁来谁说话”，而是：
  - **系统要在用户还没说话之前，就“说对话”**
  - 这就要求模型必须**先理解用户**，再决定展示什么内容。
- 
- **理解人 → 匹配兴趣 → 实时推荐**
  - 用户画像是推荐系统的“语境”输入
  - 精准的兴趣建模，是高点击率、高留存的基石

# 用户画像机制的优点

- **可解释性强**
- 我们能够清楚解释推荐的理由：
  - “因为你关注美食话题，所以我们推荐了这家人气餐厅。”
  - 这使得推荐结果具备**透明性与可控性**，利于用户信任与系统调试。
- **覆盖面广，支持多兴趣建模**
- 通过标签体系，我们不仅能捕捉用户当前的主兴趣，还能识别其**潜在的、长期演化出的多样化偏好**：
  - #健身、#母婴、#家装、#旅游.....
  - 为多样性推荐、多品类电商、跨域兴趣挖掘等提供支持。

# 用户画像机制的局限性

- **冷启动问题严重**
- 依赖标签的完整性与内容库的标注质量：
- 新用户 → 行为数据稀少，难以贴上准确标签
- 新内容 → 缺乏明确标签，无法被正确召回
- 系统在初期对人和物都“不了解”，导致推荐能力受限。
  
- **难以发现“惊喜”兴趣**
- 用户画像只反映“已经表达过的兴趣”，  
**对那些用户**未曾点击、但可能喜欢的内容捕捉能力不足：
- 用户从未看过“滑雪”视频，但他可能会一见钟情。
- 传统标签匹配系统容易陷入“兴趣窄化”与“回音室效应”。



# 用户画像机制的局限性

- **标签匹配的尽头：机器只能认“字”，却不懂“意”**
- 尽管标签匹配直观、易实现，但它在实际推荐中暴露出一个**致命缺陷**：
- 它看得见“标签”，却看不懂“语义”。
- **案例思考：系统看不懂的“兴趣相关”**
- 假设某位用户最近频繁浏览“徒步鞋”，系统据此打上用户画像标签：#徒步鞋
- 而物品库中有一款高评分的“登山杖”，标签是：#登山杖
- 在**标签匹配的系统**中，这两个词是完全不同的字符串，没有任何交集。
- 于是： 登山杖很优质    用户潜在感兴趣    但系统无法召回

# 用户画像机制的局限性

- 问题的本质：标签 ≠ 语义

问题类型	描述
语义鸿沟 (Semantic Gap)	系统只会比对“#徒步鞋”和“#登山杖”这两个标签字面是否一致，而无法理解它们在语义空间中高度相关
泛化能力弱 (Poor Generalization)	用户对一个物品感兴趣，系统无法泛化到同一兴趣簇的其他物品（如登山、户外装备）
稀疏性问题 (Sparsity)	长尾物品、新内容、新用户标签稀缺，导致系统无法将其纳入有效召回路径

- 结果：错失了本应推荐的优质内容，用户兴趣“失联”
- 转折：推荐系统迫切需要“语义理解”能力
- 我们需要一种技术，能够突破标签的表层限制：
- 理解“徒步鞋”与“登山杖”的兴趣相关性
- 泛化“用户行为”到语义邻近的内容区域
- 补全长尾物品和冷启动用户的语义连接

# 引向下一个技术范式：表示学习

- 与其依赖离散的标签，不如让系统“学会”如何将物品和用户嵌入到同一个**连续语义空间**中。  
在这个空间里：
- **“兴趣相近”的物品距离更近，“行为相似”的用户聚集在一起**
- 为了让机器真正理解用户与物品之间的**语义联系**，现代推荐系统引入了最重要的表示学习技术基石：**Embedding（向量嵌入）**。
- 我们不再用孤立的离散标签（如 #美食、#徒步鞋）描述用户或物品，
- 而是将**每个用户（User）和每个物品（Item）**都映射为同一个高维空间中的一个**向量坐标**：

实体	向量表示（示例）
"北京"	[0.32, -0.15, ..., -0.41]
"上海"	[0.35, -0.11, ..., -0.38]
"徒步鞋"	[-0.50, 0.81, ..., 0.12]
"登山杖"	[-0.48, 0.79, ..., 0.15]

# 跨越语义鸿沟的关键武器：Embedding（向量嵌入）



- **语义距离的定义：相近 = 相似**
- 在这个空间中：
- “北京”和“上海”向量靠得很近 → 表示它们在地理/文化维度上相似
- “徒步鞋”和“登山杖”几乎重叠 → 表示它们在功能/兴趣维度上密切相关
- **语义相似性 → 空间上的接近程度**  
推荐系统可以据此实现“语义级的召回”
- **为什么 Embedding 是一次范式飞跃？**

标签匹配	向量嵌入匹配
离散、精确匹配	连续、模糊匹配（容忍表达差异）
无法理解相似标签	可以自动“感知”相似物品
强依赖人工标注	可通过大规模行为数据端到端学习
冷启动/稀疏性突出	向量可迁移泛化，弱化稀疏性

# “语义召回”：让系统真正理解兴趣

- 引入 Embedding 技术之后，推荐系统终于具备了一个全新的、强大的能力：
- **不再匹配标签，而是匹配兴趣向量**
- 这使得推荐系统具备了“语义级”的召回能力。
- 进一步推广，我们将：
- 每个用户的画像、行为、上下文信息，压缩为一个**用户向量**（User Embedding）
- 每个物品的属性、内容、类型信息，压缩为一个**物品向量**（Item Embedding）
- 向量之间的**空间距离 / 相似度（如内积、余弦相似度）**，就代表了推荐强度：距离越近 → 推荐概率越高

# “语义召回”：让系统真正理解兴趣

- 向量召回的三大核心阶段
- **1 离线训练 (Offline Embedding Training)**
- 利用大规模用户行为日志（点击、点赞、加购等）进行监督或自监督训练，构建表示模型。
- 输入：用户行为序列（用户看过/点击过的物品对）
- 模型结构：如 DSSM、Two-Tower、Item2Vec 等
- 输出：
  - 用户编码模型：将用户特征映射为向量（空间坐标）
  - 物品编码模型：将物品特征映射为向量（空间坐标）
  - 得到每个用户、每个物品的高维向量表示



# “语义召回”：让系统真正理解兴趣

- **2 向量建库与索引构建 (Vector Indexing)**
- 将所有物品向量写入向量数据库（如 Faiss、Milvus），构建近似最近邻（ANN）索引结构。
- 支持大规模快速检索（百万级、亿级）
- 常用索引结构：HNSW、IVF、PQ 等
- 支持 Top-K 相似向量查找，查得快、精度高

# “语义召回”：让系统真正理解兴趣

- **3 在线召回 (Online Vector Retrieval)**
- 用户进入系统，实时计算用户向量，与向量库中所有物品进行相似度计算，取 Top-K 作为召回候选集。
- 用户向量：基于长期兴趣、短期行为、上下文等实时生成。
- 检索方式：以用户向量为查询，在物品向量库中找相似向量。
- 返回结果：语义空间中最“接近”的 K 个物品，即为该用户的候选推荐集。

# “语义召回”：让系统真正理解兴趣

- 系统视角总结：向量召回的三大优势

对比维度	传统召回	向量召回（Embedding）
匹配方式	标签/规则匹配	语义空间距离（向量相似度）
泛化能力	差，冷启动严重	强，即使没看过也能找到“语义近邻”
排序可解释性	高（可解释标签）	较低（需配合可解释建模）
工程扩展性	规则体系复杂、难维护	统一建模 + 向量库易扩展

- 它将推荐系统从规则驱动的标签匹配，推向了真正意义上的**语义驱动召回**。  
是从“知道你做过什么” → “理解你为什么做” → “预测你还可能喜欢什么”的进化跳板。

# 向量召回：让推荐系统真正“理解兴趣”的技术突破



- 在现代推荐系统中，向量召回已经成为**最核心、最有效**的召回机制之一。
- 它通过嵌入技术（Embedding），将人和物映射到统一的语义空间中，实现了前所未有的**泛化能力与语义匹配能力**。
- **核心优势一：强大的泛化能力**
- 向量召回不依赖于精确标签或历史点击记录，而是依赖于用户与物品之间在语义空间中的**邻近关系**。
- 这使推荐系统具备了“**猜你喜欢**”的能力，大幅提升了推荐的**多样性与惊喜感**。

# 向量召回：让推荐系统真正“理解兴趣”的技术突破



- **核心优势二：缓解冷启动问题**
- 对**新上线的物品**，即使没有用户行为数据，我们也可以通过其内容（如标题、描述、图片）生成初始向量表示：
- “零行为” ≠ “零推荐”
- 只要它与某类已有物品相似，系统就能自动将其**推荐给兴趣相近的用户**，让新物品迅速获得曝光与冷启动机会。

# 向量召回：让推荐系统真正“理解兴趣”的技术突破



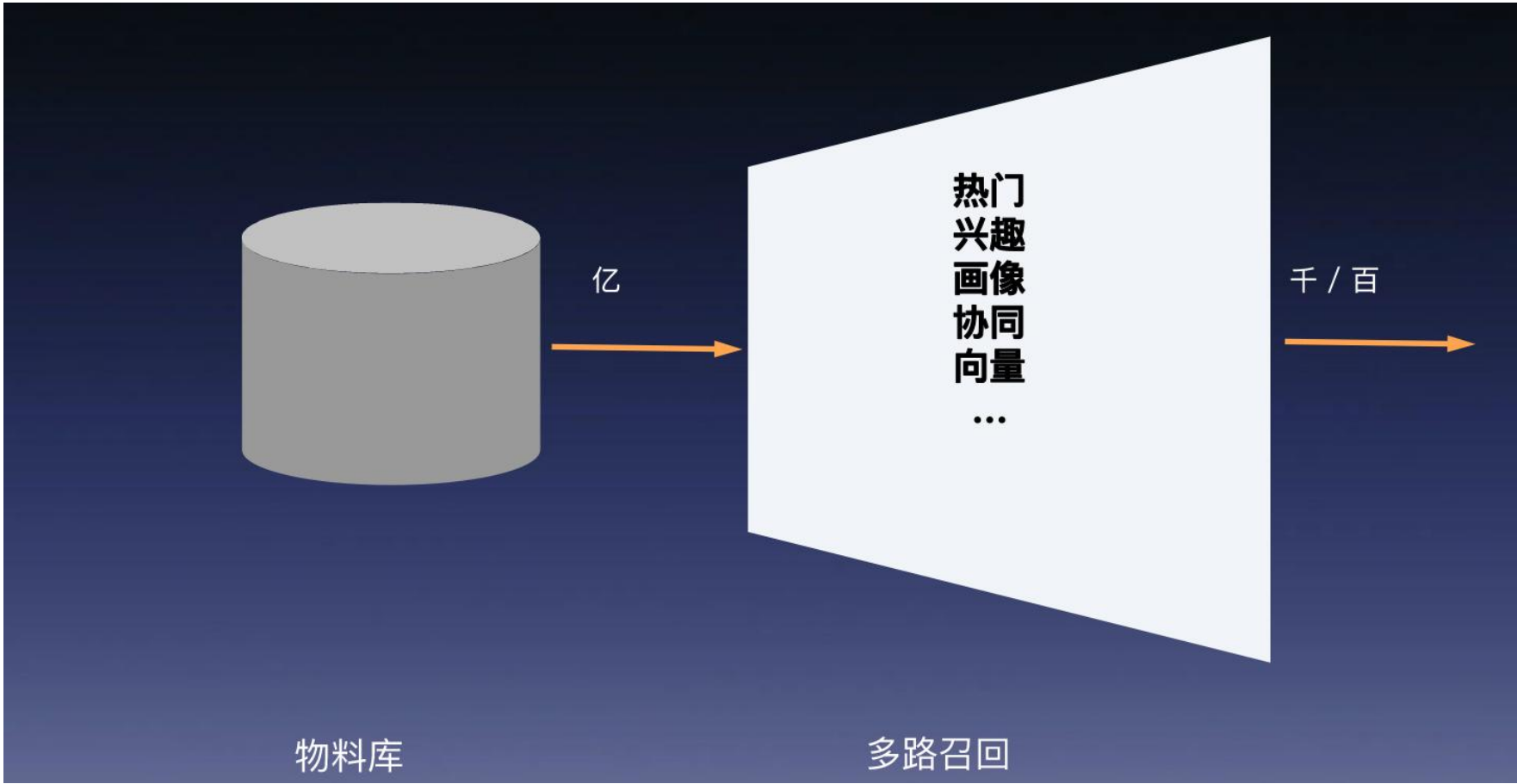
- **核心优势三：实现真正的语义匹配**
- 与传统“标签是否一致”的硬匹配不同，向量召回关注的是：“**兴趣是否相投**”，“**语义是否接近**”。
  - 标签世界是离散、片段的。
  - 向量世界是连续、可泛化的。
- 推荐系统因此从“规则驱动”进化为“兴趣建模驱动”，更接近人类认知方式。



# 多路召回：让不同策略各展所长的系统化协同机制

- 既然每种召回方式都有其独特优势与局限，最有效的做法并不是“二选一”，而是：
- 取长补短，兼收并蓄构建一个“多通道协同”的召回体系
- 工业级推荐系统的标准架构：多路召回（Multi-channel Recall）
- 在真实的大规模系统中，召回层通常由**多个策略通道并行工作**，形成一个**稳定、高覆盖、强泛化**的组合机制。
- 常见的召回通道及其功能分工

通道类型	功能定位与技术特点
向量召回	主力召回方式，负责兴趣泛化与语义匹配，发现长尾好物
用户画像召回	基于静态标签与兴趣分面，提供可解释性与多样性保障
热门召回	基于全站热榜，作为新用户冷启动时的“兜底通道”，确保有内容可看
地理召回	针对本地场景（如外卖、商圈、出行）推送附近相关内容，提升相关性与转化率
规则策略召回	由运营策略驱动（如活动、手动推品），具备强可控性和精准定向能力

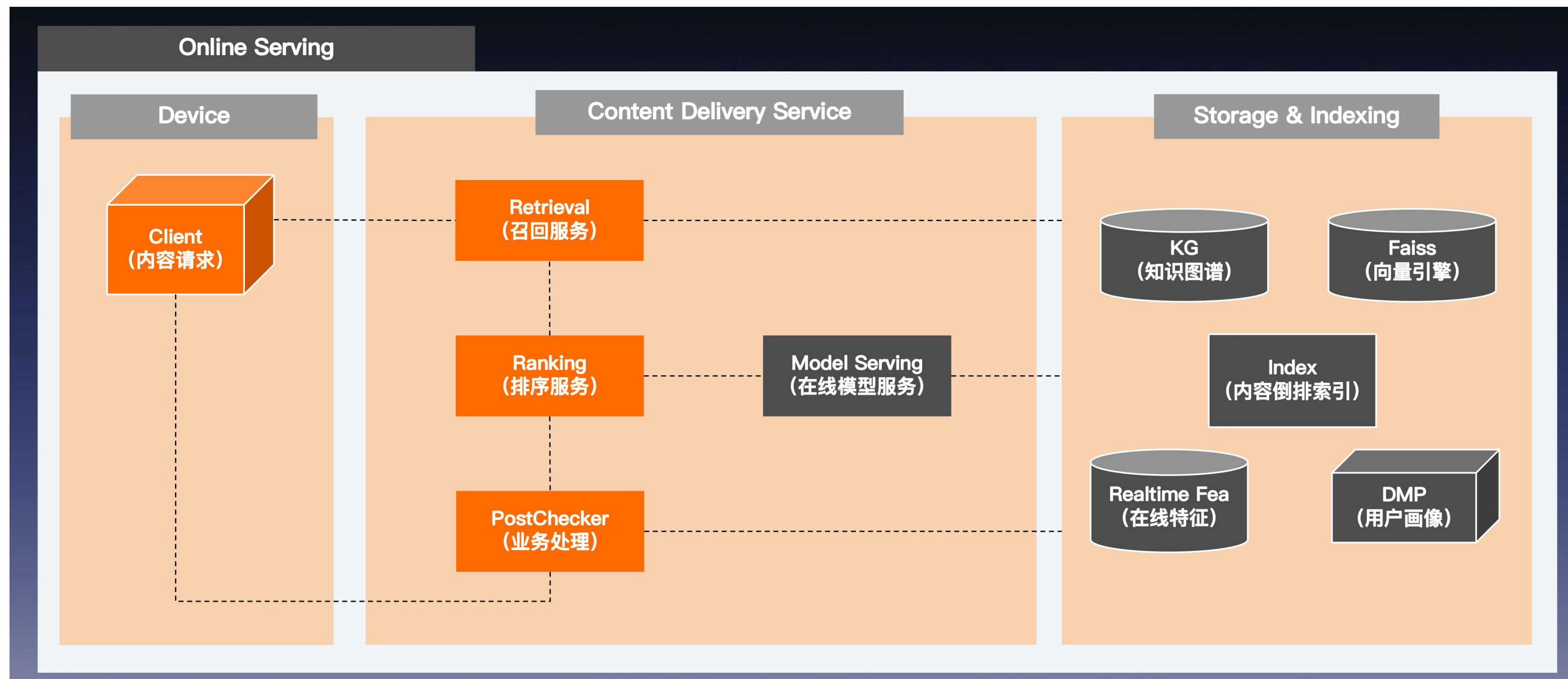


# RAG系统与结构化推理：从语言到知识的融合 极客时间

- **协同机制：召回结果的融合与精炼**
- 多通道召回的结果将被**统一汇总**，系统进行**去重、去异常、打分合并**等预处理，构成一个高质量的候选池。
- 随后，这些候选内容会被送入**排序系统（Ranking）**，由CTR模型进行最终的兴趣打分与精选排序。
- **总结：召回多通道 + 排序深建模 = 推荐系统的黄金双层结构**
- 多路召回确保**覆盖广、响应快、稳定性强**
- 精排模型确保**匹配深、选择准、转化率高**
- 推荐系统正是靠这套“召回—排序”协同架构，实现了从理解内容 → 理解人 → 匹配需求的完整闭环。
- 随着 ChatGPT、文心一言、Claude 等大语言模型（LLM）的崛起，我们迎来了一个全新的时代 —— **生成式 AI（Generative AI）时代**。而 **RAG** 则带来了检索系统的另一次飞跃。

# 代码演示

- 需提前预习项目代码，具体讲解将课程内容更新后的，代码演示视频中进行。



# 从推荐到生成：语言模型时代的新挑战

- 大模型展现出令人惊叹的**语言理解与生成能力**，能撰写文案、回答问题、生成代码、模拟对话……几乎无所不能。
- 但强大并不意味着完美，LLM 也存在一个众所周知的**核心缺陷**：
- **幻觉（Hallucination）：语言模型的软肋**
- 语言模型有时会**自信满满地给出看似合理、实则虚假的内容**，例如编造引用文献、虚构数据、捏造事件，甚至杜撰不存在的人名和机构。
- 它看起来像在回答问题，实际上却是在“编造”一个概率上最合理的答案。

# 大模型与幻觉问题：从生成力到事实实力的缺口



- 根本原因：本质是“模式学习”，而非“事实记忆”
- LLM 并非知识库，而是一个**大规模概率建模器**
  - 它学习的是**语言分布中的“模式”**，而非事实本身
  - 它的知识全部封存在训练数据中，并且**一旦训练完成，就被冻结不变**
  - 无法访问外部世界，也无法保证事**实时效**与准确性

# 大模型与幻觉问题：从生成力到事实力的缺口



- **灵魂拷问：语言强者，是否注定“言之无据”？**
- 我们已经拥有了：能写能说、上下文理解强、生成质量高的语言模型
- 但却缺少：对事实敏感、对数据可信、对知识保持最新的“连接能力”
- 有没有一种方式，能把 LLM 的语言能力，和**外部、可信、实时的知识源结合**起来？
- **这正是 Retrieval-Augmented Generation (RAG) 被提出的根本动因：**
- 让语言模型“会说话”的同时，也“**靠谱地说**”；
- 不仅生成合理，更生成真实；不仅理解语言，更连接知识



# 标准 RAG 流程

## 核心流程：先检索，后生成

- 在传统 LLM 的调用方式中，我们直接将用户问题输入模型：“你认为 Transformer 是怎么工作的？” → 模型凭训练记忆作答。
- 而在 RAG 中，我们引入了一个明确的**两阶段机制**。

# 标准 RAG 流程

## 检索 (Retrieve)

- 将用户问题视为查询 (Query)，利用向量搜索等手段，从**外部知识库**中找出相关文档片段：
- 例如，从技术文档、百科资料、API手册中找出与“Transformer原理”最相关的段落。

## 生成 (Generate)

- 将检索到的文档片段，作为“辅助信息”拼接到用户原始问题之前，形成一个包含背景知识的**扩展提示语**：背景材料 + 问题 → **输入给 LLM**。
- 此时的 LLM 不再是“凭记忆答题”，而是被要求：
- “请基于提供的材料来回答问题”，“只在你看到的上下文内进行合理生成”。
- 这大幅减少了幻觉现象，提升了输出内容的可验证性、事实性与可追溯性。

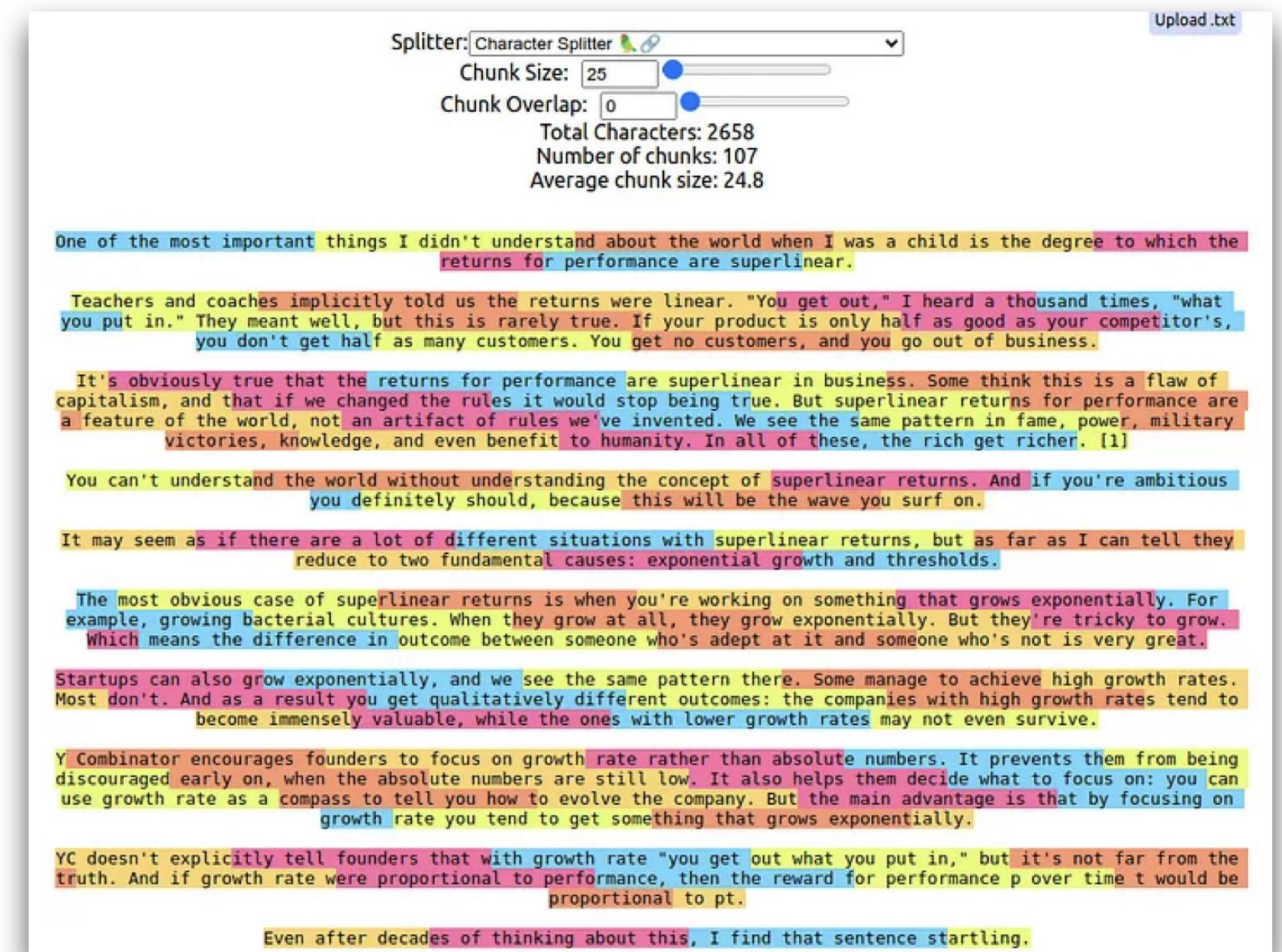
# 标准 RAG 流程

- RAG 的核心创新之一，并不是“重新发明检索”，而是巧妙地复用了推荐系统时代打磨成熟的向量检索技术。

- 标准 RAG 工作流程：语义检索驱动生成

## 1 知识库切分 (Chunking)

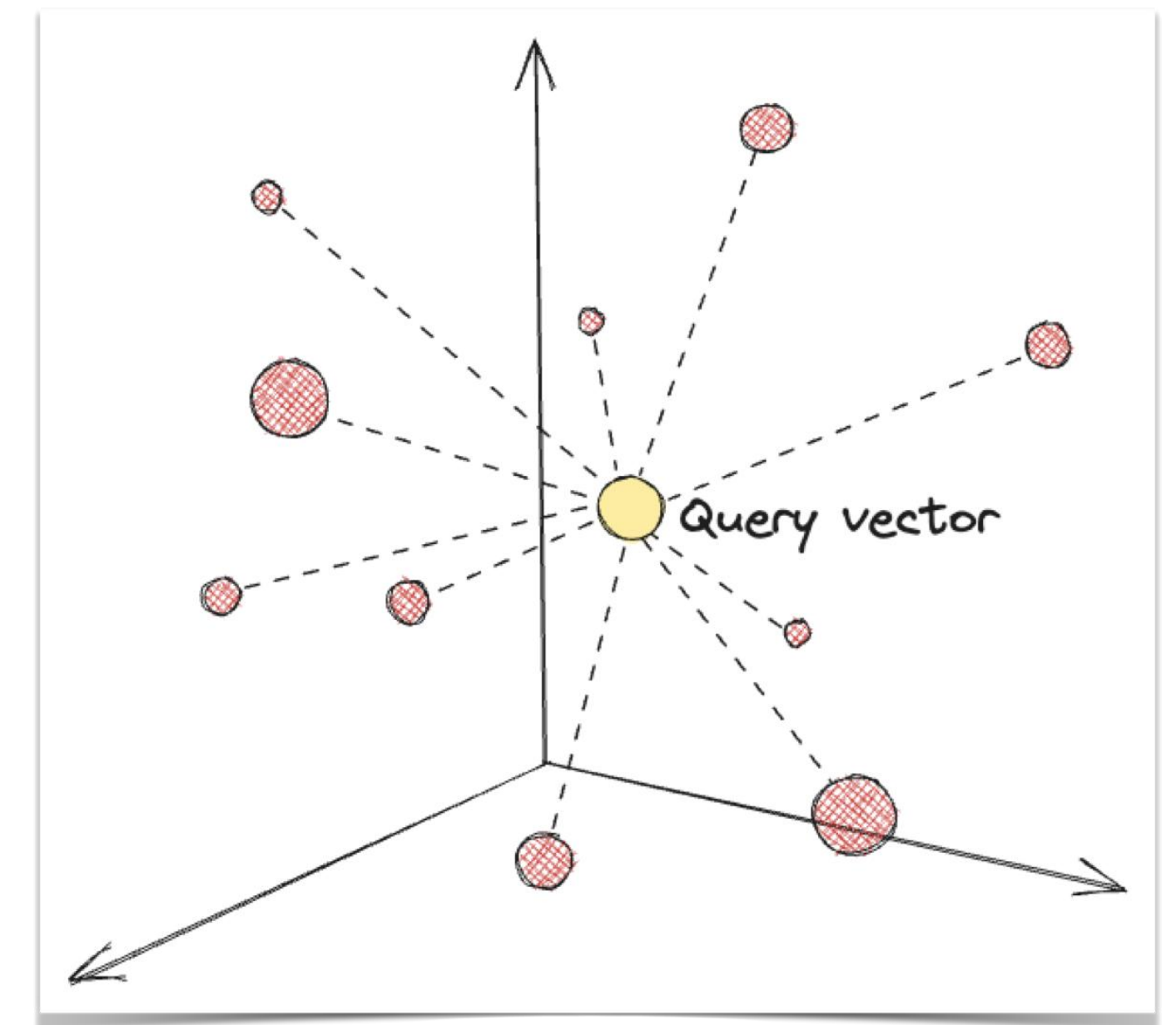
- 将原始知识源（如 PDF、Word、网页等）进行预处理，按照语义边界或固定长度切分为独立的文本块（Chunk）：
- 一段产品介绍、一节FAQ、一段法规说明.....
- 这些 Chunk 是 RAG 检索的基本单位。





# 标准 RAG 流程

- **2 嵌入建库 (Embedding & Indexing)**
- 使用一个强大的嵌入模型（如 BGE、OpenAI Embedding、Cohere Embed）将每个文本块转换为向量表示，然后将这些向量写入向量数据库（如 Faiss、Milvus、Weaviate），并建立索引结构（如 HNSW、IVF 等）以支持高效检索。
- **3 查询向量化 (Query Embedding)**
- 当用户提问时（如：“如何提升转化率？”），系统使用同一个嵌入模型将其转换为查询向量。
- **4 向量召回 (Vector Recall)**
- 以查询向量  $q$  为检索关键，系统在向量库中执行一次“向量搜向量”的 ANN 检索，找出与之**语义最接近的 Top-K 个文本块**作为上下文注入 LLM。



# 推荐与 RAG 的技术继承关系：向量召回骨架的延续



- 技术观察：RAG 与推荐系统是**同一个“骨架”**
- 我们惊讶地发现，RAG 的这套流程与推荐系统的“向量召回”高度相似：

推荐系统	RAG 系统
用户向量 (User Vector)	问题向量 (Query Vector)
物品向量 (Item Embedding)	知识块向量 (Chunk Embedding)
召回候选物品	召回候选文本块
用户-物品相似度召回	问题-知识块语义相似度召回

- 技术栈高度复用
- 检索范式跨领域迁移
- 嵌入学习的工程成熟度得以延续
- **RAG 没有推翻推荐系统的旧体系，而是站在其肩膀上向“语言理解 + 知识融合”迈进**
- 推荐系统解决“**人与物品**”的匹配，RAG 系统解决“**问题与知识**”的匹配
- 底层技术几乎同构，目标从“推荐喜欢的”转为“回答可靠的”

# RAG 的结构性瓶颈：无法多跳推理，语义链断裂

极客时间  
训练营

- 我们已经构建起了一个强大的 RAG 系统，
- 它能够准确处理大量“事实类”、“描述类”的提问，例如：
  - “Transformer 有哪些核心结构？”
  - “BERT 是哪年提出的？”
- 这些属于**点状知识**的语义匹配任务，RAG 表现优异。
- 但当问题需要**多步推理**、**多实体关联**、**复杂关系挖掘**时，RAG 往往力不从心。
  - **思考两个真实的复杂查询**
  - “《流浪地球》导演的妻子主演过哪些电影？”
  - “开发了 BGE 模型的机构的 CEO 是谁？”
- 这类问题不再是“查一个块能回答”的事实类问题，而是需要**跨块拼接 + 关系链推理**。

# RAG 的结构性瓶颈：无法多跳推理，语义链断裂

极客时间  
训练营

- 标准 RAG 的结构性困境
- 1. 知识割裂
- RAG 的知识库是由非结构化文档切分后的 **碎片化文本块（Chunk）** 构成：
  - Chunk A：“郭帆是《流浪地球》的导演。”
  - Chunk B：“吴京主演了《战狼2》。”
  - Chunk C：“吴京和谢楠结婚了。”
- 这些文本块之间是**彼此独立的**，缺乏显式结构与连接方式。同时向量检索机制无法保证我们找到**所有回答问题所必要的片段**。



# RAG 的结构性瓶颈：无法多跳推理，语义链断裂

极客时间  
训练营

- 2. 跨块推理难
- 标准 RAG 的检索器是“以查询向量找语义近邻”
- 它可能找到了 Chunk A 或 Chunk B，甚至找到了 C，但即便如此，也**无法明确**三者之间存在**一条可组合的关系链**。
- 缺乏“结构意识”的检索器，无法对“隐式的多步逻辑路径”做出有效拼接。
- 最终生成结果容易：忽略关键信息，逻辑跳步，幻觉补全（猜测关系）。

# RAG 的结构性瓶颈：无法多跳推理，语义链断裂

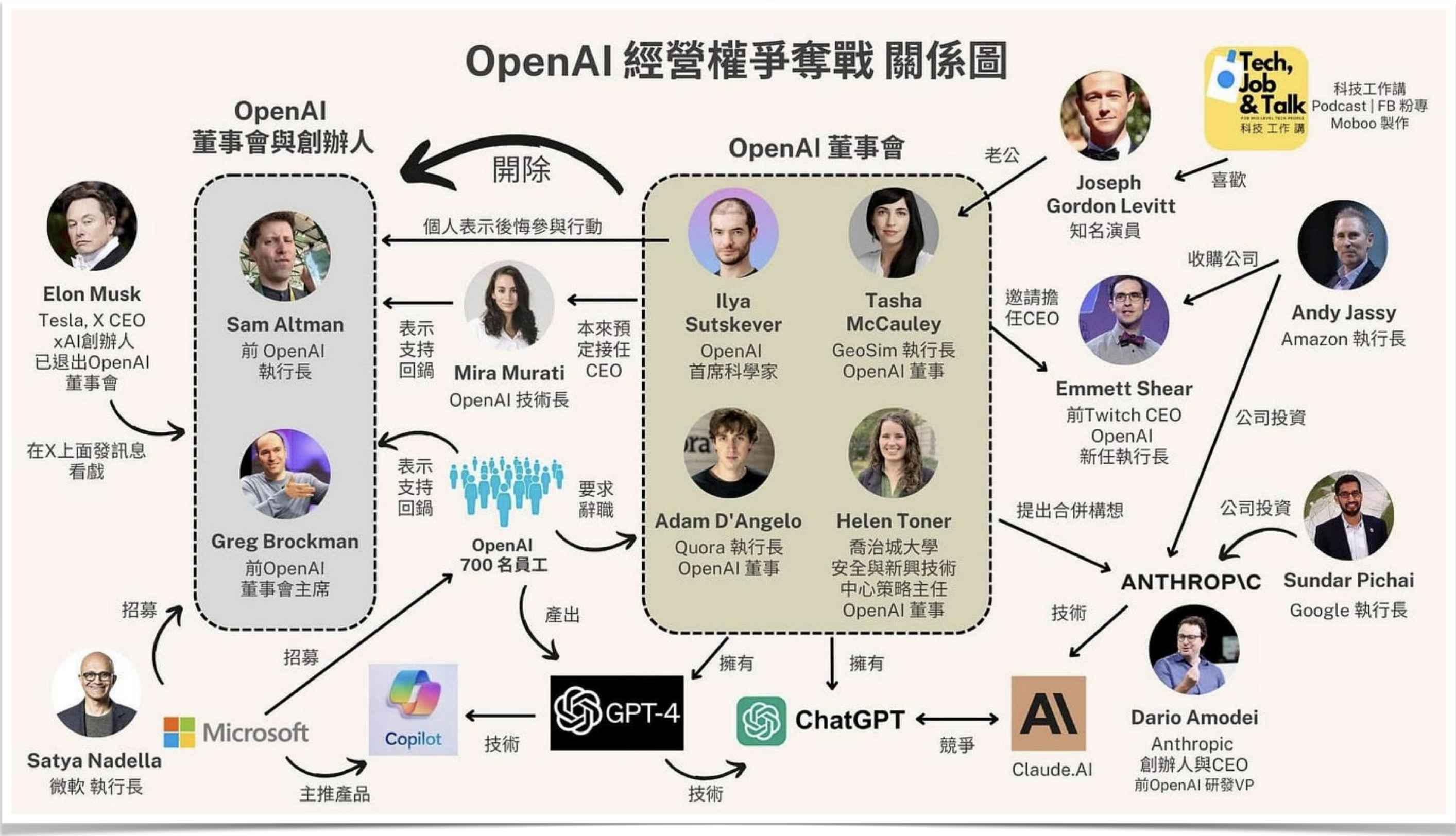
极客时间  
训练营

- 本质问题：在**结构化任务**中使用**非结构化（分块数据）检索**，是一种“低维错配”
- 我们希望模型能处理“图状知识”或“多跳路径”，但实际提供的只是“独立文本块”的向量近似检索：
- 结构性问题，被非结构化手段处理，必然牺牲推理能力。
- 因此，面对结构化知识的查询场景，标准 RAG 不足以胜任复杂推理任务。
- 这正引出了下一阶段的技术趋势：
  - 如何将**结构化知识**引入 RAG？
  - 如何构建“**图结构**”、“**三元组**”、“**知识库**”驱动的增强检索？
  - 如何让 LLM 从“**阅读段落**”进化为“**操控知识图谱**”？



# 极客时间

- 往往，用户想知道的是那些**文字中看不到**的信息。



# 知识图谱增强 RAG：实体关系图结构的检索机制

极客时间  
训练营

- **关系推理的突破口：请回符号主义的“老朋友”——知识图谱（Knowledge Graph）**
- 面对标准 RAG 无法处理的跨块推理、实体关系理解等问题，我们必须重拾**符号主义学派**最经典、最强大的知识表示工具：
- **知识图谱（Knowledge Graph）**
- 它不是一堆文本段落的拼接，而是一个**结构化、可计算、便于推理的“事实网络”**。
- **什么是知识图谱？——知识的“结构版本”**
- 知识图谱是一种将真实世界知识表示为三元组（Triple）的结构：  
**（实体1, 关系, 实体2）**
- 它构成了一个由“实体”为节点、“关系”为边的**多跳有向图结构**。



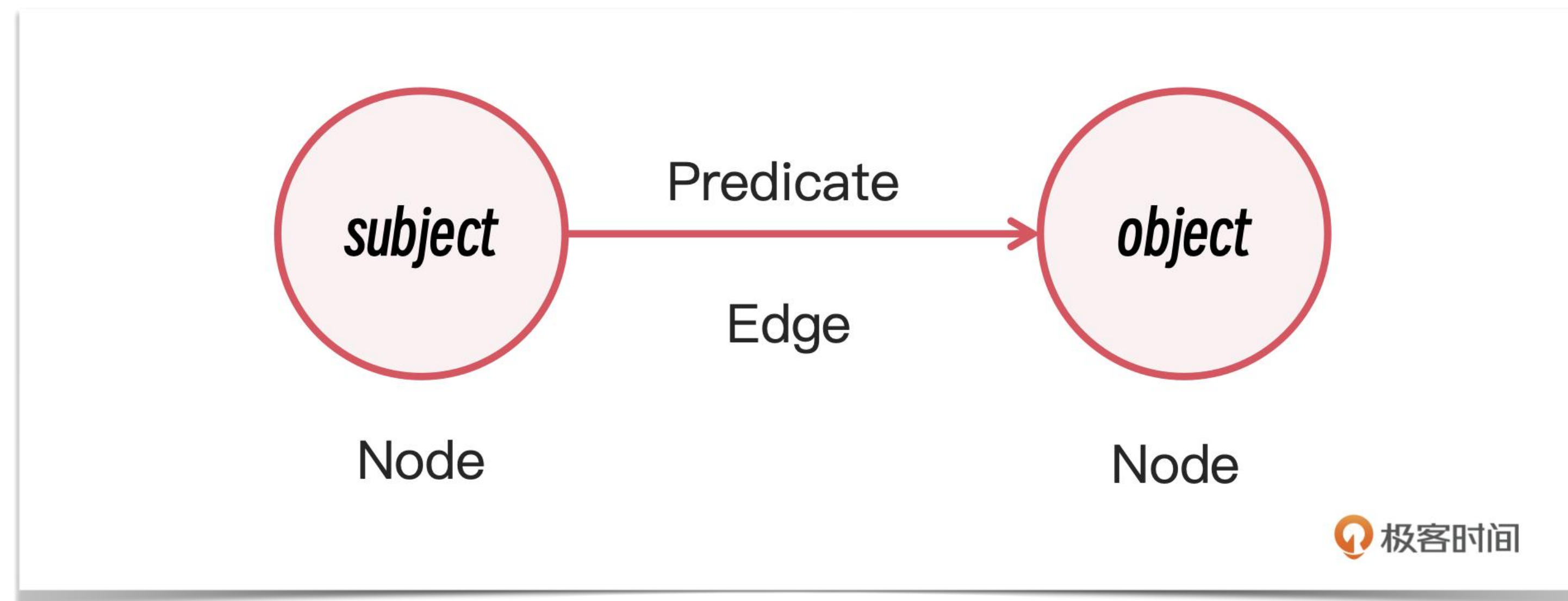
# 流程全解析（实体识别 → 子图转述 → 增强生成）



- **实体 (Entity)**
  - 代表现实世界中的具体对象 —— 图谱中的“节点 (Node) ”
  - 如：吴京/谢楠/《流浪地球》/《战狼2》
  - 类别：导演、电影、人物、机构等
- **关系 (Relation)**
  - 表示实体之间的语义联系 —— 图谱中的“边 (Edge) ”
  - 类别：导演/主演/配偶是/CEO 是/所属机构为.....

# 流程全解析（实体识别 → 子图转述 → 增强生成）

- 三元组的拓扑结构：
- Entity 1: Albert Einstein
- Relation: 提出
- Entity 2: 相对论



# 流程全解析（实体识别 → 子图转述 → 增强生成） 极客时间

- **-任务目标-** 对于给定的中文文本，请根据指定的实体类型列表 [{entity\_types}], 识别并提取文中的所有实体，并进一步找出这些实体之间存在的明确关联关系。
- **-执行步骤-** （注：输出时请使用 {{tuple\_delimiter}} 作为元组内的分隔符）
  - **1. 识别所有实体** 从文本中识别所有符合指定类型的实体。对于每一个被识别的实体，提取以下信息：
    - entity\_name (实体名称): 从文本中提取的实体标准名称。
    - entity\_type (实体类型): 必须是 [{entity\_types}] 列表中的一种。
    - entity\_description (实体描述): 依据文本，对该实体的关键属性、特点或行为进行全面、客观的描述。
- **输出格式应为：** ("entity"{{tuple\_delimiter}}<实体名称>{{tuple\_delimiter}}<实体类型>{{tuple\_delimiter}}<实体描述>)



# 流程全解析（实体识别 → 子图转述 → 增强生成） 极客时间

- **2. 识别实体间的关系** 基于步骤1的结果，找出所有存在明确关联的（源实体，目标实体）配对。
- 对于每一对相关实体，提取以下信息：
  - source\_entity (源实体): 源实体的名称，必须与步骤1中识别的名称完全一致。
  - target\_entity (目标实体): 目标实体的名称，必须与步骤1中识别的名称完全一致。
  - relationship\_description (关系描述): 依据文本内容，用一句话清晰地阐述源实体与目标实体之间的具体关系（例如：A是B的创始人，A在B公司工作，A与B共同出席了C事件）。
  - relationship\_strength (关系强度): 一个从1到10的整数。1表示关联最弱（如仅在同一段落提及），10表示关联最强（如存在直接、明确的互动或定义性关系）。
- **输出格式应为：** ("relationship"{{tuple\_delimiter}}<源实体>{{tuple\_delimiter}}<目标实体>{{tuple\_delimiter}}<关系描述>{{tuple\_delimiter}}<关系强度>)

# 流程全解析（实体识别 → 子图转述 → 增强生成）

- 真实场景示例如下：
- 1. 给定的**知识库文本**：马化腾是腾讯公司的联合创始人和首席执行官。腾讯公司成立于1998年，总部位于深圳，是中国领先的互联网增值服务提供商。该公司开发了著名的即时通讯软件QQ和多功能社交应用微信。微信不仅是一款聊天工具，还提供了支付和公共服务等功能，深刻地改变了中国人的生活方式。
- 2. 任务参数：
- entity\_types: ['人物', '组织机构', '产品', '地点']
- tuple\_delimiter: |

# 流程全解析（实体识别 → 子图转述 → 增强生成）

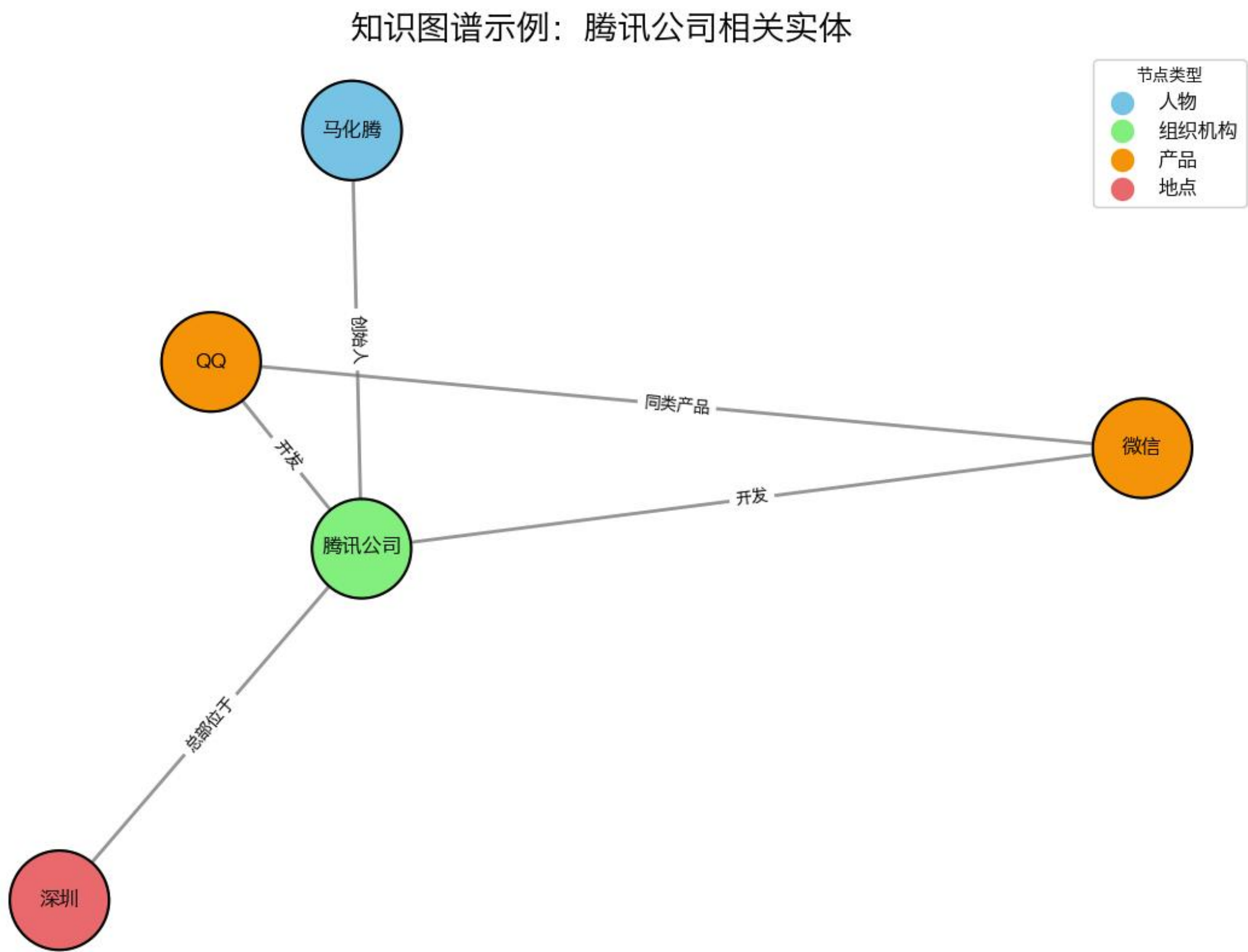
- **步骤1：识别所有实体**
- ("entity"|马化腾|人物|腾讯公司的联合创始人和首席执行官。)
- ("entity"|腾讯公司|组织机构|成立于1998年，总部位于深圳，开发了QQ和微信的中国互联网服务提供商。)
- ("entity"|深圳|地点|腾讯公司的总部所在地。)
- ("entity"|QQ|产品|由腾讯公司开发的即时通讯软件。)
- ("entity"|微信|产品|由腾讯公司开发的多功能社交应用，提供聊天、支付和公共服务功能。)

# 流程全解析（实体识别 → 子图转述 → 增强生成）

- **步骤2：识别实体间的关系**
- ("relationship"|马化腾|腾讯公司|马化腾是腾讯公司的联合创始人和首席执行官。|10)
- ("relationship"|腾讯公司|深圳|腾讯公司的总部位于深圳。|9)
- ("relationship"|腾讯公司|QQ|腾讯公司开发了即时通讯软件QQ。|10)
- ("relationship"|腾讯公司|微信|腾讯公司开发了多功能社交应用微信。|10)
- ("relationship"|QQ|微信|同为腾讯公司开发的产品。|6)

# 流程全解析（实体识别 → 子图转述 → 增强生成）

- 步骤3：编织三元组为知识图谱关系



# 流程全解析（实体识别 → 子图转述 → 增强生成）

- **知识图谱的力量：结构化推理的基石**
- 知识图谱的优势不仅在于可读性，更在于：
  - 关系明确，适合路径遍历与逻辑组合
  - 易于构建**推理链**（Multi-hop Reasoning）
  - 支持**实体链接、消歧、溯源**等关键能力
  - 可作为 RAG 检索的高级形式，实现**更精细、更可靠**的内容召回

# 流程全解析（实体识别 → 子图转述 → 增强生成）

- 标准 RAG 系统的局限，在于它的检索对象是**非结构化文本块**，而复杂推理问题，恰恰需要在**结构化知识网络**中进行路径级理解。
- 将 RAG 的“检索单位”，从**文本块**升级为结构化的“**子图 (Subgraph)**”。



# 流程全解析（实体识别 → 子图转述 → 增强生成）

## 1. 实体与关系识别（Entity and Relation Recognition）

- 从用户提出的问题中识别出核心实体 **和** 它们之间的关系。
- **示例问题：**“《流浪地球》主演的妻子主演过哪些电影？”**一个更完整的识别结果应该是：**
  - **链接初始实体：** 系统将文本“《流浪地球》”成功链接到知识图谱中代表这部电影的唯一节点（例如：Entity ID: 123）。
  - **遍历第一层关系：** 从《流浪地球》节点出发，沿着“主演”（has\_actor）的关系边，找到链接到的实体节点，例如 吴京。
  - **遍历第二层关系：** 从 吴京 节点出发，沿着“妻子”（has\_spouse）的关系边，找到链接到的实体节点 谢楠。
  - **锁定目标实体：** 此时，系统已将“《流浪地球》主演的妻子”这个复杂的描述成功解析并链接到了知识图谱中的具体实体：谢楠。
  - **执行最终查询：** 以 谢楠 为主语，查询她“主演过”的所有“电影”，并返回结果列表。

# 流程全解析（实体识别 → 子图转述 → 增强生成）

## 2. 子图转文本（Subgraph Verbalization）

- 为了准确回答最初的问题（“《流浪地球》主演的妻子主演过哪些电影？”），系统在第二步遍历图谱后，构建的子图应该围绕核心人物“吴京”和“谢楠”。因此，转译成的文本也应反映这个正确的子图。
- 正确的子图结构应为：**
  - （《流浪地球》，有主演，吴京）
  - （吴京，配偶是，谢楠）
  - （谢楠，客串，《战狼2》）
  - （谢楠，参演，《幸福额度》）
  - ...（以及谢楠参演或主演的其他作品）
- 基于此正确子图，转化得到的自然语言“背景资料”应为：**
- “根据知识库：吴京是《流浪地球》的主演。吴京的配偶是谢楠。谢楠曾客串出演电影《战狼2》，并参演了电影《幸福额度》。”

# 流程全解析（实体识别 → 子图转述 → 增强生成）

## 3. 增强生成 (Augmented Generation)

- **构建正确的 LLM 提示输入 (Prompt):**
  - 根据知识库：吴京是《流浪地球》的主演。吴京的配偶是谢楠。谢楠曾客串出演电影《战狼2》，并参演了电影《幸福额度》。
  - **问题：** 《流浪地球》主演的妻子主演过哪些电影？
- **LLM输出：** 根据提供的资料，《流浪地球》的主演吴京的妻子是谢楠。她参演过的电影包括《战狼2》（作为客串）和《幸福额度》。

# 流程全解析（实体识别 → 子图转述 → 增强生成）



- 知识图谱增强 RAG 的本质：语言理解 × 知识结构 = 多跳推理能力

标准 RAG	知识图谱增强 RAG
检索单位：非结构化文本块	检索单位：结构化知识子图
匹配逻辑：语义向量相似度	匹配逻辑：图结构路径 + 实体关系链遍历
优势：覆盖广、泛化强	优势：逻辑清晰、关系显式、支持推理链构建

- 让大模型拥有“图谱思维”，不仅查得准，还能形成语义通路。

# 三大优势：精准性、推理性、可解释性

- 将结构化的知识图谱（Knowledge Graph）引入 RAG 系统，不仅弥补了标准 RAG 的语义盲区，还带来了以下三大关键能力提升：
  - **1. 精准性（Precision）——从“猜测”到“确认”**
- 传统向量召回依赖语义相似度，容易“以词会意”，造成事实偏差；
- 图谱路径是**明确的、结构化的实体-关系对**，可直接作为**强事实支撑**；
- 不再是“模型编造了谁是谁的配偶”，而是“图谱中明确指出了配偶关系”。
- 这极大地降低了语言模型幻觉（Hallucination）产生的概率。

# 三大优势：精准性、推理性、可解释性

- 将结构化的知识图谱（Knowledge Graph）引入 RAG 系统，不仅弥补了标准 RAG 的语义盲区，还带来了以下三大关键能力提升：
  - **2. 推理能力（Reasoning）——支持多跳关系链的复杂理解**
- 图谱提供了天然的**有向边结构**，可实现多步**灵活的关系跳转**；
- 通过**图遍历**，系统可以构建一条**从问题实体出发 → 多跳扩展 → 回答所需信息**的逻辑路径；
- 像“谁的妻子主演了什么”这类复杂链式问题，LLM 在结构支持下也能**有依据、有步骤地完成**。



# 三大优势：精准性、推理性、可解释性

- 将结构化的知识图谱（Knowledge Graph）引入 RAG 系统，不仅弥补了标准 RAG 的语义盲区，还带来了以下三大关键能力提升：
- **3. 可解释性（Explainability）——让推理路径变得透明可信**
- 图谱路径是**可视、可追溯、可验证**的；
- 每一条答案都能明确指出其**推理来源与语义链条**；
- 不再是“AI 说它知道”，而是“AI 告诉你它从哪条路径知道的”。
- 这对于在医疗、金融、法律等**高可信度领域的应用**尤为关键。



# 检索的进化之路：从关键词匹配到语言理解的智能跃迁

- 在本模块中，我们完整追溯了一条技术演进的长链：一条关于“检索”如何不断变得更精准、更智能、更理解人类意图”的道路。
- 它穿越了搜索引擎、推荐系统，到今天的大语言模型，从最早的“关键词匹配”，一路走到“结构化语义推理”。

# 检索的进化之路：从关键词匹配到语言理解的智能跃迁



- **第一阶段：搜索引擎时代（符号匹配）**
- **倒排索引 + TF-IDF**  
构建了文本检索的基本能力，让系统能在**海量信息中快速匹配关键词**。
- 但它仍停留在“字符级理解”——系统知道你在找什么词，却不理解你**为什么找**。

# 检索的进化之路：从关键词匹配到语言理解的智能跃迁

- **第二阶段：推荐系统时代（从符号到语义）**
- **排序模型**  
将“用户点击行为”引入排序决策，实现从“文档为中心”到“用户为中心”的转变。
- **用户画像召回**  
通过静态标签与长期兴趣建模，实现了第一代**个性化推荐**。
- **向量召回（Embedding Matching）**  
一次革命性跃迁！  
将用户与物品映射到同一语义空间，从此**系统开始真正理解“兴趣相近”而非“标签一致”**。
- **多路召回架构**  
引入多个召回通道（向量、热度、地理、画像等），实现**系统级协同与覆盖率优化**。

# 检索的进化之路：从关键词匹配到语言理解的智能跃迁

- **第三阶段：大语言模型时代（检索 × 生成的融合）**
- **标准 RAG (Retrieval-Augmented Generation)**  
借助向量召回，为大模型提供“开卷考试”的能力，让生成内容更真实、可控、有支撑。
- **RAG + Reranker（重排序机制）**  
引入更强大的排序器，优化上下文质量，提升 LLM 的生成准确性与相关性。
- **知识图谱增强 RAG**
- 回归结构化符号主义，注入**实体-关系路径**，让语言模型具备“多跳逻辑推理”的能力，实现了从“匹配内容”到“理解关系”的能力跃升。

# 检索的进化之路：从关键词匹配到语言理解的智能跃迁



- 这一部分**最需要大家记住**的是：无论是传统搜索还是现代RAG，**单一模型都无法构成工业级系统**。AI系统的演进，是一部多技术协同、不断形成**复杂专家系统**的历史。

技术阶段	核心任务	关键技术	解决的问题
第一阶段：搜索引擎	关键词匹配	倒排索引 + TF-IDF	“找得到”：解决了海量文本的快速定位问题。
第二阶段：推荐系统	兴趣匹配	CTR模型 + 用户画像 + Embedding	“猜你喜欢”：从“理解内容”进化到“理解人”，通过向量化实现了语义级的兴趣匹配与泛化。
第三阶段：RAG时代	知识增强生成	检索器 (向量/图谱) + 大语言模型 (LLM)	“有据可查”：将LLM的生成能力与外部知识库结合，让模型“开卷考试”，大幅缓解了幻觉问题。

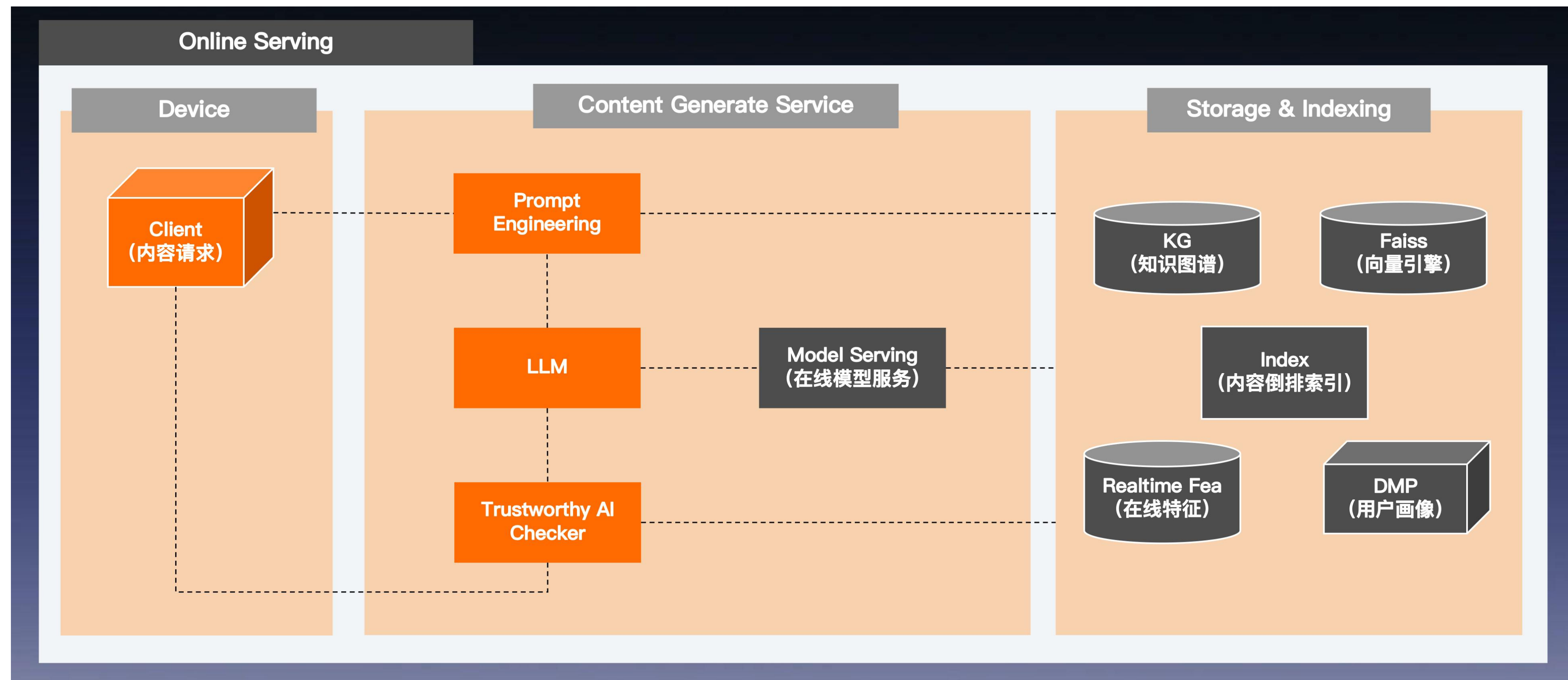
# 当“检索”走向极限，“推理”必须登场

- 在最后，我们引入**知识图谱**，回归符号主义，发现了**知识推理**对AI系统的重要意义。
- 我们发现，为了让AI系统真正像专家一样思考，仅有强大的信息检索与关联能力（如向量召回）是远远不够的。当任务的复杂度超越了简单的信息匹配，我们就必须赋予系统一种更古老而强大的能力——**自动推理**。
- 这正是**符号主义学派**的核心价值所在。**知识图谱**增强的RAG，本质上就是一次“**检索系统**”与“**轻量级推理系统**”的成功融合。它证明了，只有将不同AI范式协同，才能构建出真正强大的专家系统。
- **那么，如何系统性地构建这种推理能力？**
- 如何定义规则，让机器能够像人类专家一样，在复杂的条件下做出判断？
- 如何构建一个高效的“推理引擎”，自动执行这些规则，并解决更通用的问题？
- **这正是我们下一模块将要深入探讨的核心。**



# 代码演示

- 需提前预习项目代码，具体讲解将课程内容更新后的，代码演示视频中进行。



# 目录

- 模块一：符号表示和知识库
- **模块二：规则引擎与自动推理**

# 从知识获取到逻辑决策：为何AI需要推理能力？



- 在上一模块中，我们触及了现代检索系统的前沿能力。
  - **能力**：通过向量检索和知识图谱，我们让AI具备了前所未有的**信息获取与关联**能力。
  - **成就**：它能精准、快速地从海量数据中定位事实，为决策提供“知识原料”。
  - **边界**：然而，它回答了如何**高效地获取事实**，但是在**复杂推理任务**面前显得力不从心。面对需要严密逻辑、多步推演的复杂任务，它依然存在结构性短板。
- 但是，仅有强大的**记忆**和**知识**，不足以构成真正的智能。

# 从知识获取到逻辑决策：为何AI需要推理能力?



- AI系统必须掌握一种更根本的能力——**推理 (Reasoning)**。这是从“信息处理器”到“问题解决者”的关键一跃。
- 在接下来的模块中，我们将共同探索机器推理的奥秘，核心路径包括：
  - **1. 经典的路径规划与演绎**：回归本源，学习 **A\* 搜索算法**如何以最优效率在复杂空间中找到解决方案。
  - **2. 现代的语言模型驱动推理**：探索大语言模型（LLM）如何通过**思维链 (Chain-of-Thought)** 等模式，模拟人类进行多步规划与逻辑推演。

# 从知识获取到逻辑决策：为何AI需要推理能力?



- 最终目标：两大能力合二为一，实现真正的智能决策
- 检索（知道什么） + 推理（知道怎么办） = 综合智能
- 我们将致力于将符号主义学派的这两大能力深度融合，构建不仅能**引经据典**，更能**深思熟虑**、步步为营的下一代AI系统。
- 欢迎进入AI的核心腹地——机器推理的世界。

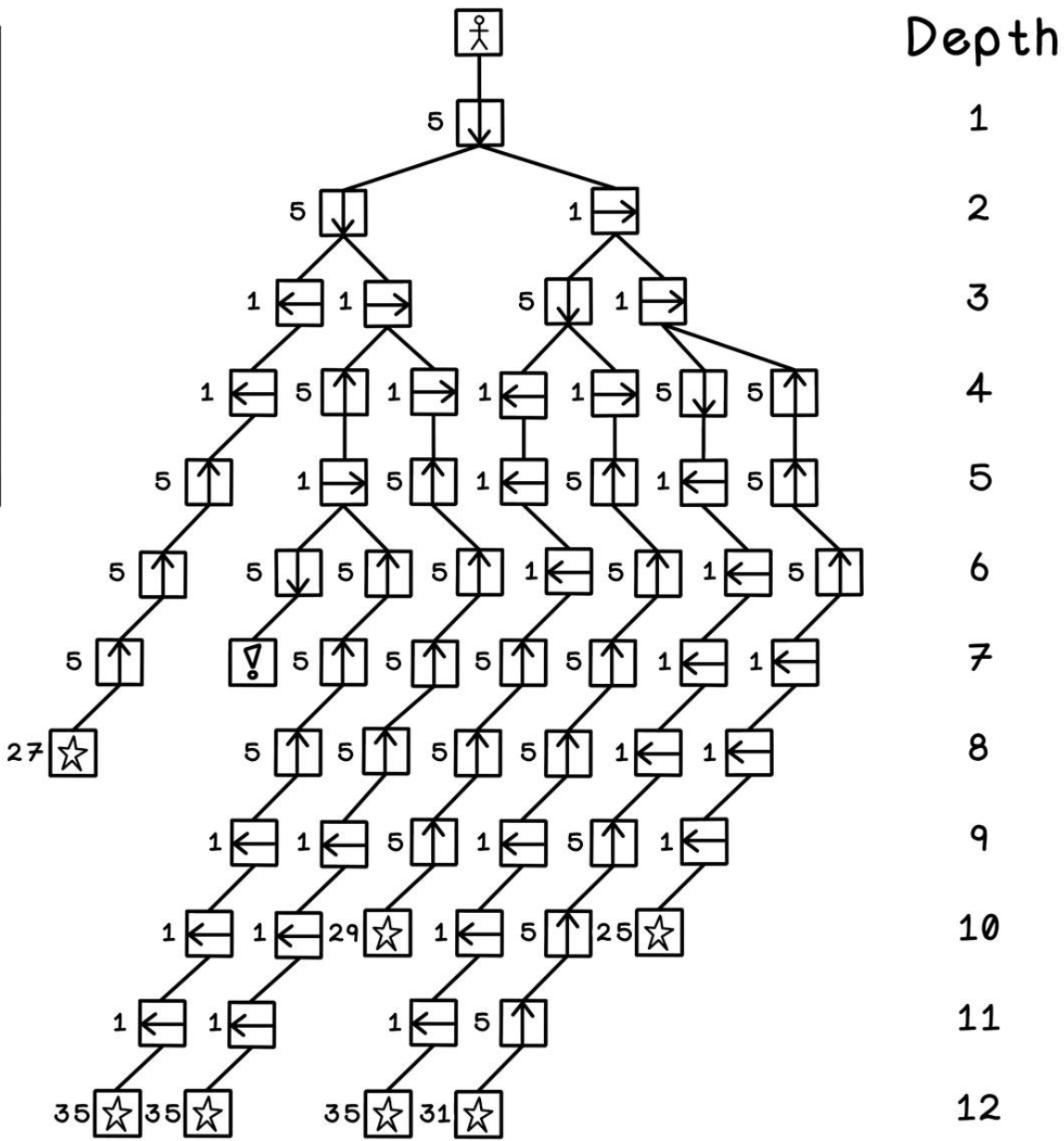
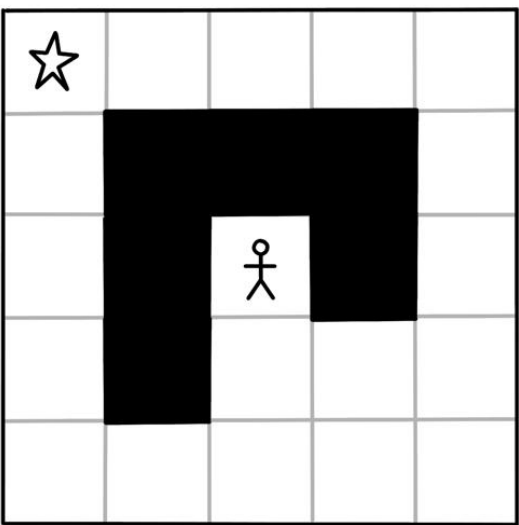


# 通用推理建模：将问题抽象为“状态空间搜索”



- 在AI眼中，所有需要“**规划**”和“**决策**”的复杂问题，都可以被简化为一个在巨大地图上**寻路**的游戏。

概念	解释：以“解迷宫”为例
状态空间 (State Space)	地图本身：包含了所有可能的情况（迷宫中的每一个格子）。
初始状态 (Initial State)	你的起点：迷宫的入口 A。
目标状态 (Goal State)	你的终点：迷宫的出口 B。
操作 (Action)	你可以做的动作：向上、下、左、右移动一格（前提是不能穿墙）。
解 (Solution)	你要找的答案：一条从起点到终点的有效路径。

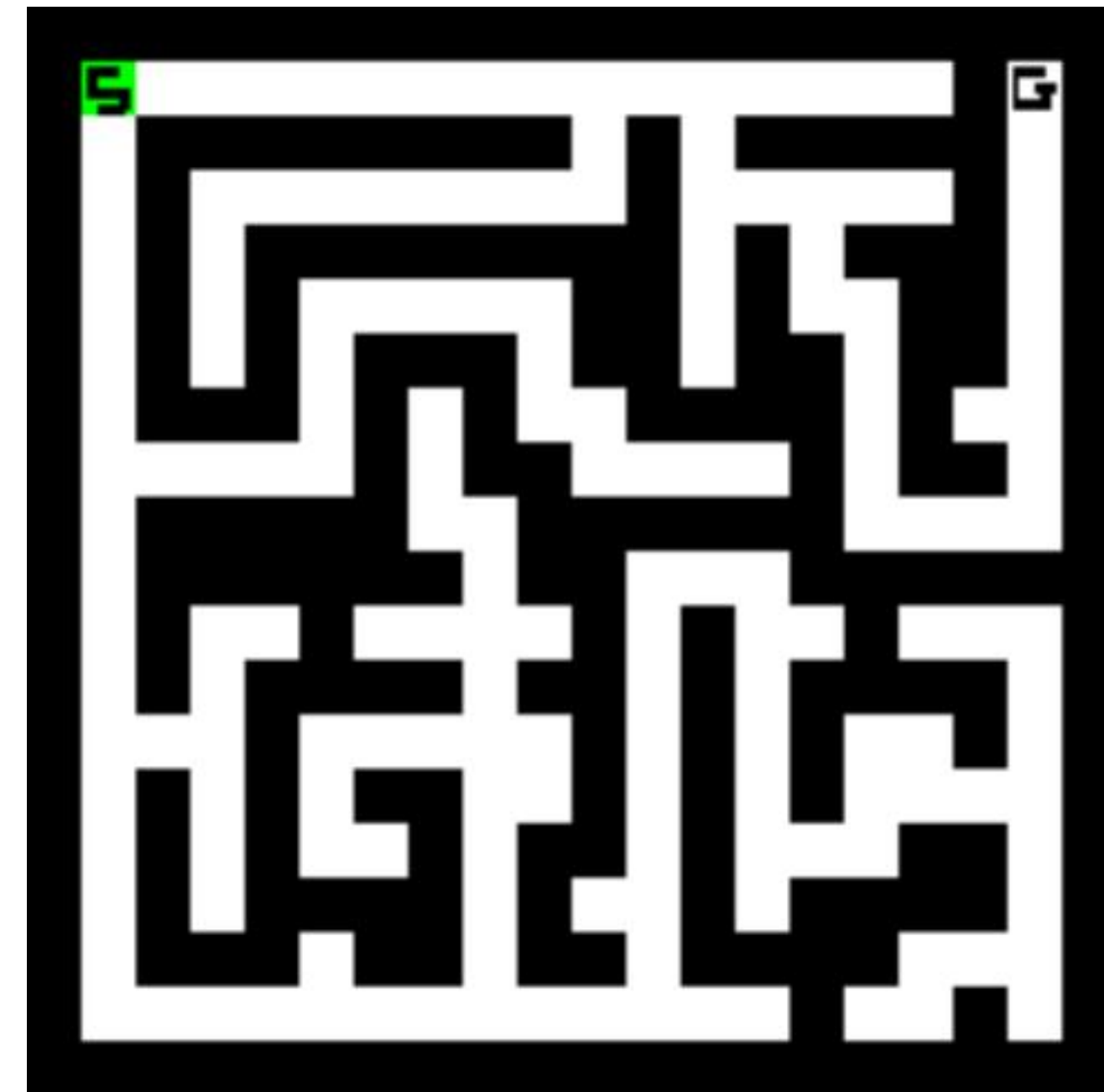


- 通过这套语言，我们可以用统一的框架来描述下棋、解谜、路径规划、甚至更复杂的决策问题。



# 盲目搜索：AI的最初尝试

- 当机器开始解决这个问题时，它并不知道捷径。它会从起点开始，系统性地探索所有可能的下一步。**探索过程：**
- **从 A 出发**，发现可以向 **上** 或向 **右** 走。
- **选择向上**，发现下一步可以向 **左** 或向 **右** 走。
- ...不断重复这个过程。



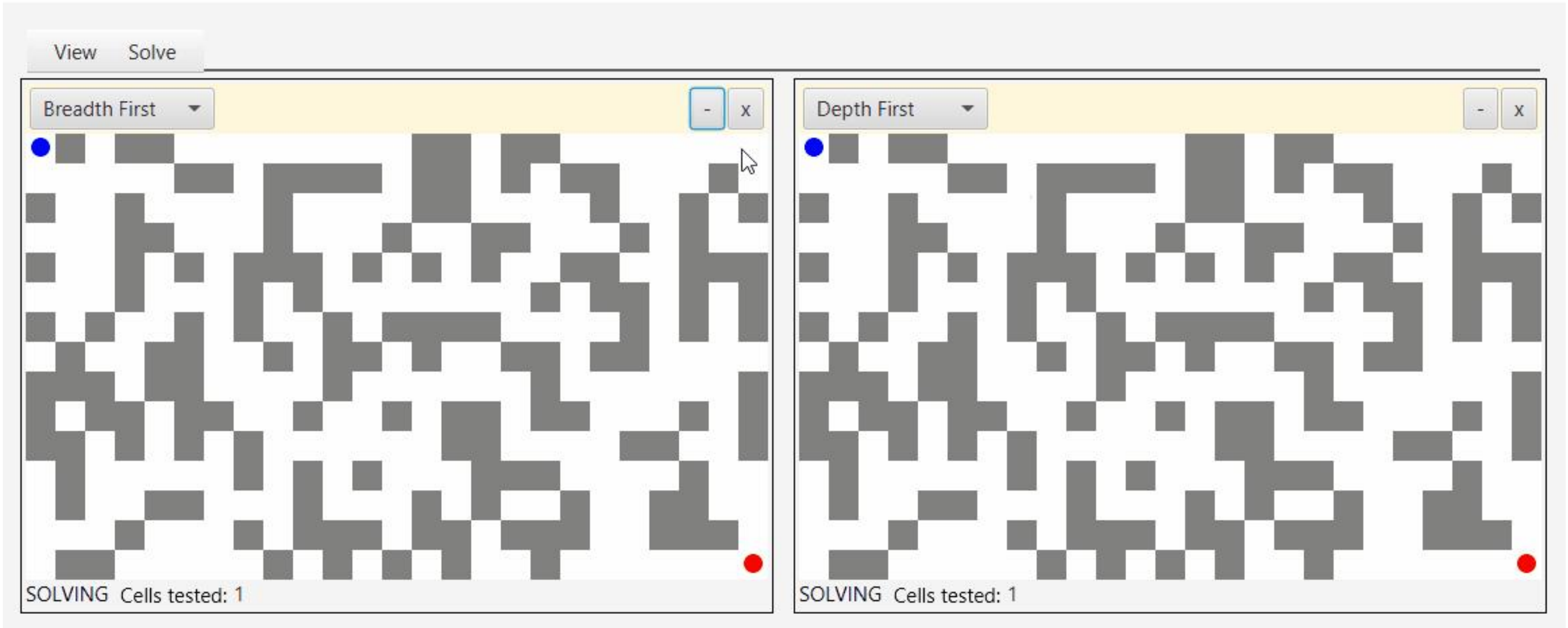
# 盲目搜索：AI的最初尝试

- **决策结构 (右侧生成的搜索树):** 这个探索所有可能路径的过程，会自然地形成一棵庞大的“**搜索树 (Search Tree)**”。
  - **树的根节点:** 是我们的 **初始状态 (A)**。
  - **树的每个分支:** 代表一次 **操作 (移动)**。
  - **树的每个节点:** 代表一个 **中间状态 (迷宫中的某个位置)**。
  - **我们的目标:** 是在这棵树的某个叶子节点上，找到 **目标状态 (B)**。
- 推理问题，被成功转化为了一个“**图搜索**”问题。接下来AI算法的优劣，就体现在 **遍历这棵搜索树的效率** 上。是广度优先？深度优先？还是更智能的 **A\* 算法**？这便是我们即将深入的核心。

# 盲目搜索：AI的最初尝试

- 在没有任何额外信息（无地图、无启发）时，我们只能采用“笨办法”——系统性地遍历所有可能性。

广度优先搜索 (Breadth-First Search, BFS)	深度优先搜索 (Depth-First Search, DFS)
核心策略：地毯式搜索	核心策略：一条路走到黑
“稳扎稳打，一层一层地向外探索。” 先看离起点1步远的所有点，再看2步远的所有点...	“不撞南墙不回头，随机选一个方向深入到底。” 碰壁后，才回溯到上一个路口，换条路继续。
池中涟漪投入石子后，波纹逐圈均匀扩散。	靠墙寻路在迷宫里，始终用右手摸着墙走。
优点：最优保证只要有解，保证能找到最短路径。	优点：空间轻量内存开销极小，只需记住当前走过的路径。
缺点：内存爆炸需要存储所有待探索的“边界”节点，空间复杂度极高。	缺点：质量随机找到的第一个解不一定是最优的，且可能在无限深的分支中“迷路”而耗尽时间。



# 盲目搜索：AI的最初尝试

- **BFS 用空间换时间，保证了结果的“最优性”；而 DFS 用时间换空间，追求了过程的“经济性”。**
- 这两种看似简单的“笨办法”，构成了所有复杂搜索算法的基石。它们的鲜明缺点也引出了一个至关重要的问题：
- **有没有一种方法，既能像BFS一样“有远见”，又能像DFS一样“轻装上阵”？**

# A\* 搜索：最优性与效率的结合

- 要打破BFS和DFS的困境，我们需要引入额外的信息，一种“启发式”的指引。

什么是“启发 (Heuristic)”？	一个直观的例子
定义：一种“经验法则”或“有根据的猜测”。它不是100%精确的，但能为我们在决策的十字路口提供一个明智的指引，告诉我们哪个选项“看起来更有希望”。	场景：你在城市中找路 起点：悉尼中央车站 终点：邦迪海滩 你的启发式知识：“邦迪海滩在东边。” 决策：在每个路口，你虽然不知道确切的最佳路线，但你会优先选择朝东走的路。这个“朝东走”的倾向，就是一种强大的启发，它能帮你排除大量明显错误的方向。
在AI搜索中，我们用一个函数 $h(n)$ 来量化这种启发： $h(n)$ = 从当前节点 $n$ 到终点的 预估 距离/成本 核心价值：启发信息让我们能评估未来的可能性，从而让搜索变得“有目的”，不再是盲目的。	最常用的启发函数：直线距离 (欧几里得距离) 虽然实际路线（蓝线）是弯曲的，但“两点之间直线最短”（红虚线）这个常识，为我们提供了一个绝佳的、永远不会高估实际距离的乐观估计。

- “启发”就是我们的“智能指南针”。它将一个盲目的搜索者，变成了一个有策略、有方向感的探索者。



# A\* 搜索：最优性与效率的结合

- 核心思想：A\* 算法将BFS的严谨性与启发式搜索的智能性完美结合。
- A\* 的决策公式： 在每一步，A\* 都会评估所有待探索的节点，并选择评估分最低的那个。评估函数  $f(n)$  定义如下：

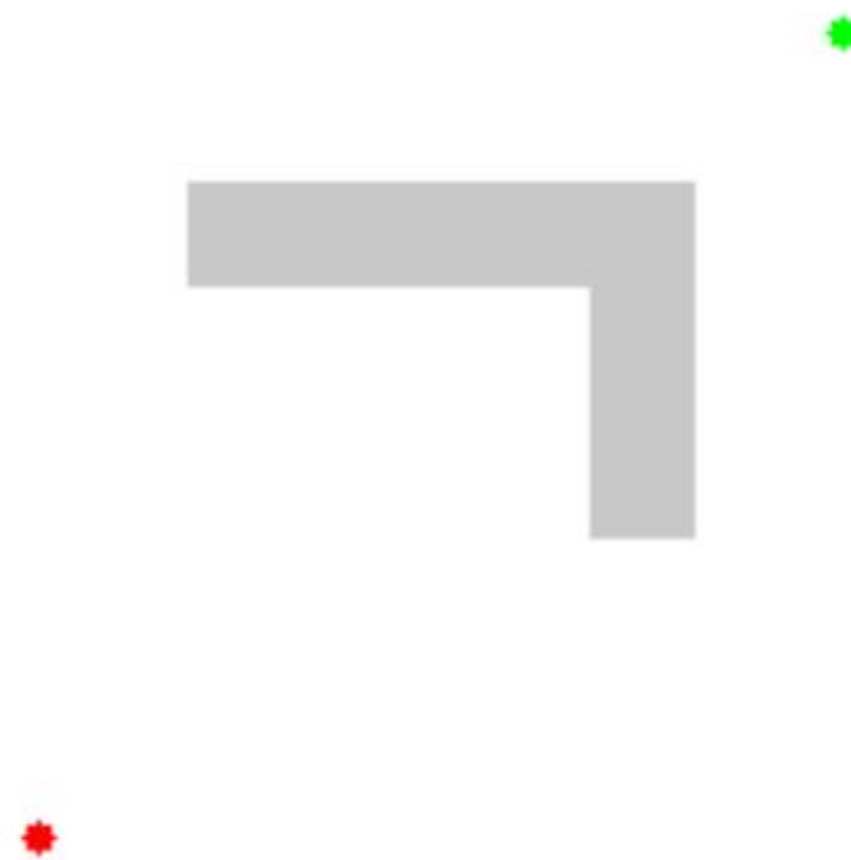
$$f(n) = g(n) + h(n)$$

公式组成	含义	作用
$g(n)$	已付出的代价 (Actual Cost)	我从起点走到当前节点 $n$ 已经花了多少成本。——这是对过去的总结，体现了BFS的“公平性”。
$h(n)$	对未来的预估 (Heuristic Cost)	我猜测从当前节点 $n$ 走到终点 大概还要花多少成本。——这是对未来的展望，体现了启发式的“前瞻性”。
$f(n)$	综合评估值 (Total Estimated Cost)	经过节点 $n$ 到达终点的总成本预估。——这是我们最终的决策依据，平衡了现实与理想。



# A\* 搜索：最优性与效率的结合

- **工作方式：** A\* 算法始终探索“**看起来总路程最短**”的路径。它巧妙地避免了：
- **BFS的盲目扩张：** 如果一条路已经走得太远 ( $g(n)$ 很大)，即使它离终点看起来很近 ( $h(n)$ 很小)，A\*也可能放弃它。
- **DFS的执迷不悟：** 如果一条路虽然很短 ( $g(n)$ 很小)，但它指向了离终点很远的方向 ( $h(n)$ 很大)，A\*也会优先探索其他更有希望的路。
- A\* 算法是在拥有良好启发函数的前提下，能够保证找到最优解，并且效率远超盲目搜索的“黄金标准”算法。它是AI寻路、机器人导航、物流规划等无数领域的绝对基石。



# A\*的边界：理想建模的前提难以满足

- A\* 算法如此强大，为何我们还需要新的方法？因为它依赖于两个苛刻的前提。

前提一：形式化的“世界模型”	前提二：有效的“启发函数”
A* 要求问题被精确地定义为一个“状态空间图”。这意味着，我们必须提前： 1. 定义所有可能的状态 (Nodes) 2. 定义所有合法的操作 (Edges) 3. 量化每次操作的成本 (Weights)	A* 的效率，完全取决于一个好的启发函数 $h(n)$ 。设计一个好的 $h(n)$ 需要： 1. 深刻的领域知识：你需要对问题有洞察，才能做出“有根据的猜测”。 2. 可受理性： $h(n)$ 必须是“乐观的”，即永远不能高估实际成本，否则A*无法保证找到最优解。
现实世界的挑战： 很多问题是模糊的、开放的，并且用自然语言描述的。例如：“帮我策划一场为期三天的北京家庭游。”“我的网站流量下降了，分析一下可能的原因并给出解决方案。”我们几乎不可能为这类问题手动构建一个完整的状态空间图。	现实世界的挑战： 对于许多抽象问题，我们根本无法设计出有效的启发函数。例如：写一首关于秋天的诗，它的“状态”是什么？“成本”又是什么？证明一个复杂的数学定理，如何“预估”离最终证明还有多远？

# A\*的边界：理想建模的前提难以满足

- 经典搜索算法是强大的“解题工具”，但它们需要一个被完美数学化的“理想世界”。
- 当问题无法被轻易“建模”时，我们需要一种能够直接理解和处理**非结构化信息**（如自然语言），并在一个**隐式、开放的世界**中进行推理的全新范式。
- 这，正是大语言模型（LLM）登场的舞台。

# 新范式：大语言模型驱动的推理机制

- 核心思想：大语言模型（LLM）不“搜索”状态空间，它直接“生成”通往答案的思考路径。
- LLM如何“思考”？—— 思维链 (Chain-of-Thought, CoT)
- LLM突破在于，我们可以通过特定的提示（Prompting），引导它模仿人类的逐步推理过程。

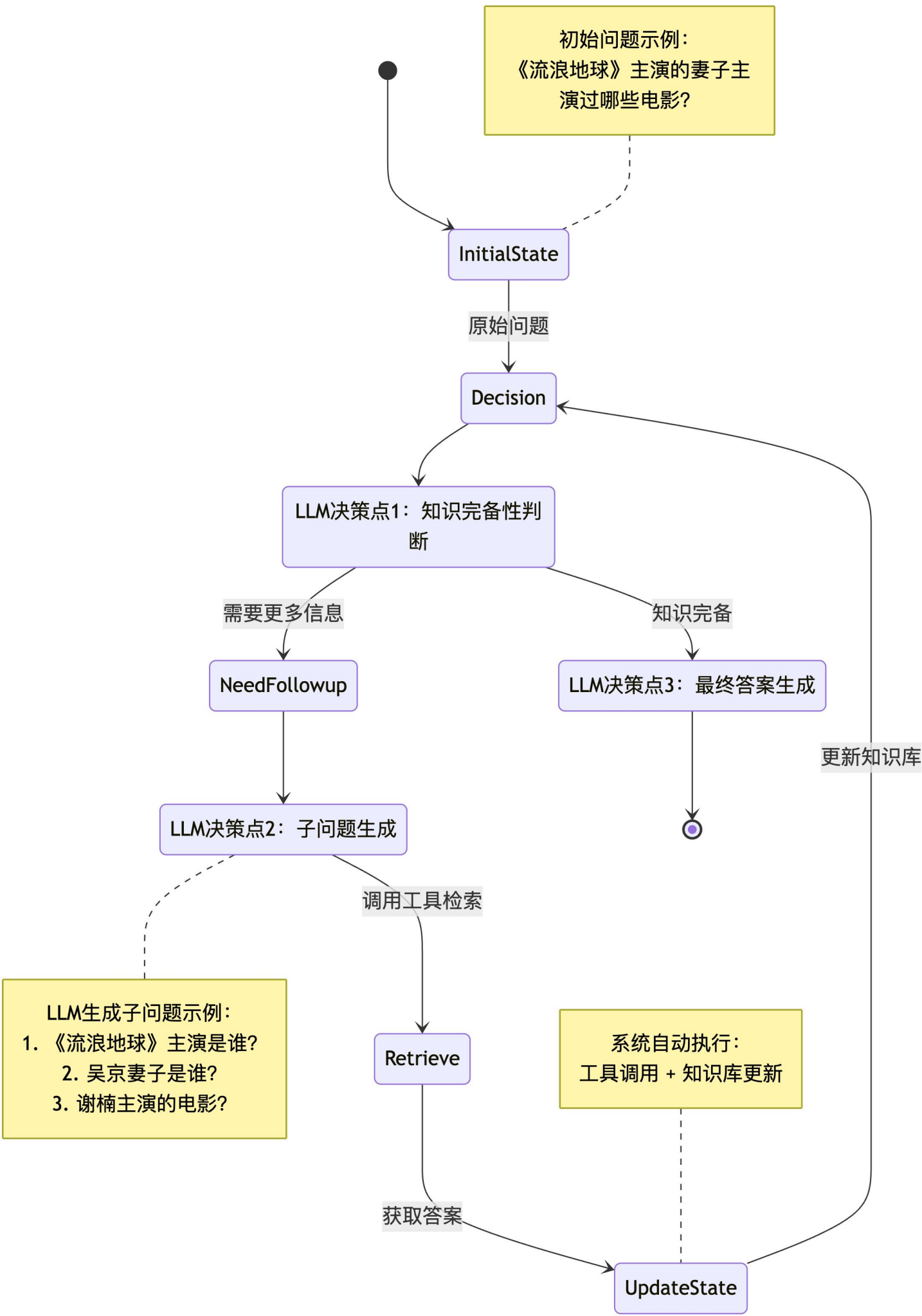
传统提问方式 (Standard Prompting)	思维链提问方式 (Chain-of-Thought Prompting)
输入： 问：一个杂货店有23个苹果。如果他们用了20个来做午餐，又买了6个，现在他们有多少个苹果？ 答：	输入： 问：一个杂货店有23个苹果。如果他们用了20个来做午餐，又买了6个，现在他们有多少个苹果？ 答：让我们一步一步地思考。
LLM的输出 (可能错误): 9个 过程： 模型试图直接“猜”出答案，容易在多步计算中出错。它是一个“黑箱”，我们不知道它是如何得出结论的。	LLM的输出 (更可靠): 1. 一开始有23个苹果。 2. 他们用了20个，所以剩下 $23 - 20 = 3$ 个。 3. 然后他们又买了6个，所以现在有 $3 + 6 = 9$ 个。 最终答案是9个。 过程： 模型被引导着生成了中间的思考步骤。这个过程本身，就是一种动态的、自洽的推理链。



# LLM 作为“智能启发函数”

- 思维链（CoT）解决了“如何思考”，但没解决“知识从哪里来”的问题。当LLM自身的知识不足或过时，它依然会“一本正经地胡说八道”。
- 解决方案：赋予LLM“提问”和“检索知识”的能力。**

Self-Ask	本质：一个由LLM驱动的状态机 (State Machine)
与其直接猜测答案，不如让LLM将一个复杂问题，分解成若干个可以通过工具（如搜索引擎）回答的子问题。LLM会自己问自己，然后调用工具查找答案，再基于新知识继续推理。	我们可以将整个任务过程，看作一个在不同“知识状态”之间跳转的流程：（此处放置一个清晰的状态机流程图）
例子：用户问题：“《流浪地球》主演的妻子主演过哪些电影？”	<div><div>1. 初始状态：仅包含原始问题。</div><div>2. LLM决策 (Action): 分析后发现，需要先知道“主演是谁”。于是生成子问题：“谁是《流浪地球》的主演？”</div><div>3. 工具调用 (Tool Use): 执行搜索 Search(“谁是《流浪地球》的主演?”)，返回“吴京”。</div><div>4. 状态更新：现有知识更新为：“《流浪地球》的主演是吴京。他的妻子主演过哪些电影？”</div><div>5. LLM决策 (Action): 分析新状态，生成下一个子问题：“吴京的妻子是谁？”</div><div>6. ...循环此过程，直到...</div><div>7. 目标状态：所有信息集齐，LLM综合所有知识，生成最终答案。</div></div>



# LLM 作为“智能启发函数”

- **核心洞察：**在Self-Ask这个状态机中，LLM扮演的角色，与A\*算法中的“启发函数”惊人地相似。
  - **让我们再次回到核心的搜索与决策问题上：**LLM 本质上就是一个极其强大的、动态的启发函数。
  - 它不是给出一个简单的代价值，而是凭借其强大的语言理解和世界知识，直接生成一个最有价值的“下一步动作”。
- 
- 因此，Self-Ask 这类由LLM驱动的代理（Agent）系统，其工作流程可以被精炼地概括为：
  - 一个由LLM作为“启发式导航员”，来智能驱动和规划路径的状态机。它通过（分析 -> 提问 -> 检索 -> 整合）的循环，代理我们完成了复杂的知识获取与推理任务。



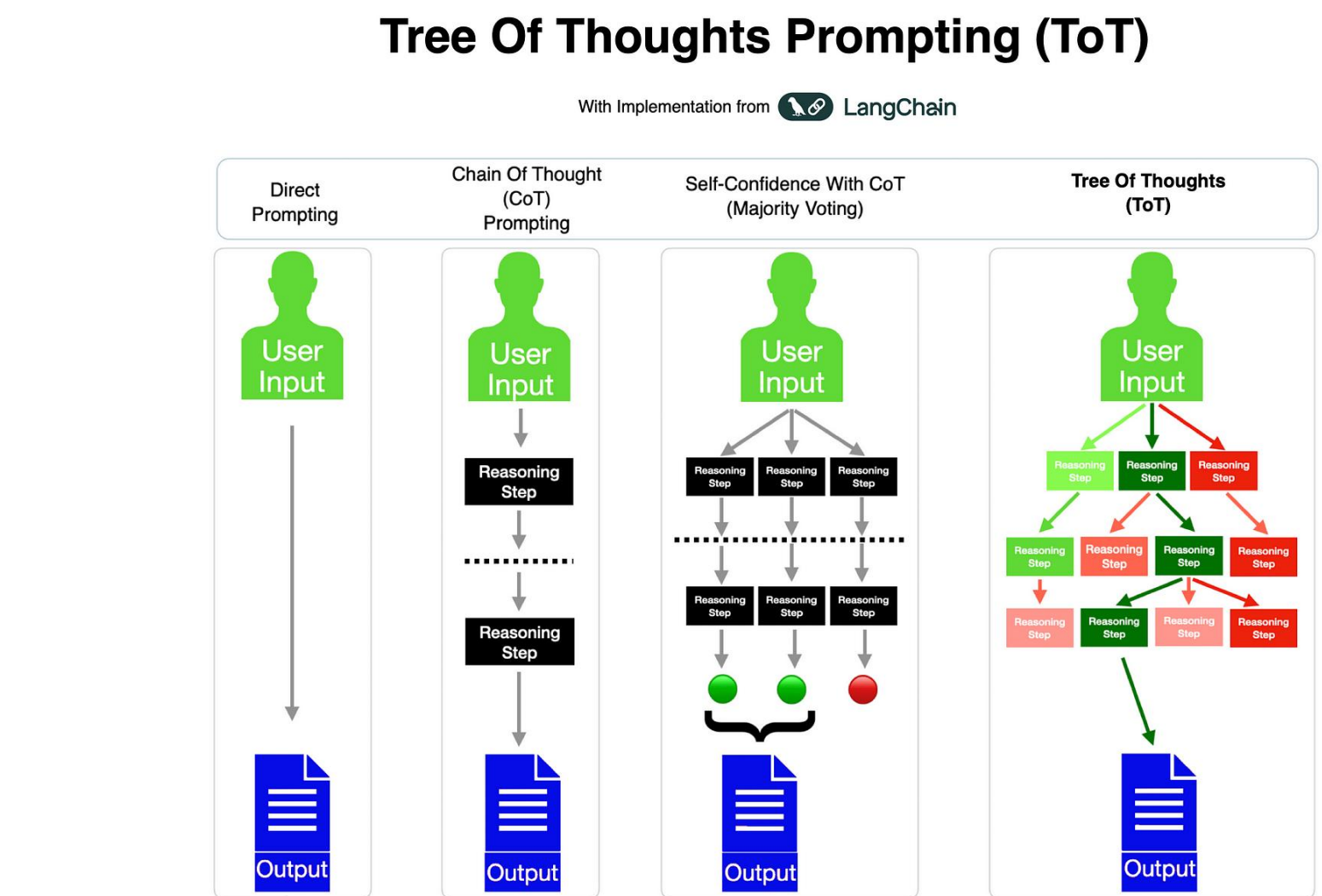
# 双路径融合：从“搜索路径”到“生成思路”

- 在本模块中，我们探索了AI解决复杂问题的两种核心范式。它们都试图回答同一个根本问题：如何在巨大的可能性空间中，智能地找到解决方案？

路径一：经典符号主义推理	路径二：现代生成式推理
核心思想：形式化与最优性	核心思想：理解与生成
将世界抽象成一个精确、结构化的状态空间图。	将世界看作一个可通过自然语言来交互和理解的开放场域。
代表算法：A* 搜索	代表模式：LLM 代理 (Self-Ask, ReAct等)

# 双路径融合：从“搜索路径”到“生成思路”

- 我们看到，经典推理像是一位严谨的逻辑学家，而现代LLM推理则更像一位拥有渊博知识和强大直觉的创造者。
- 真正的智能系统，不应是“二选一”。未来的趋势必然是**两者的强强联合（如右图 ToT）**：
- 用 **LLM 的理解能力**去解析和分解现实世界中模糊、复杂的用户需求。
- 用 **经典算法的严密性**去执行那些需要高可靠性、最优性保证的核心规划与推理步骤。
- 本模块的探索至此结束。我们不仅学会了机器如何“搜索”路径，更理解了它们如何“生成”思路。
- **知识库+专家系统+自动推理**这几种符号主义学派能力的结合，是通往更强大人工智能的必由之路。



# 推荐阅读

- 《信息检索导论》 (Introduction to Information Retrieval) 作者: Christopher D. Manning, Prabhakar Raghavan, Hinrich Schütze
- 《计算广告学》 作者: 刘鹏
- 《自动机理论、语言和计算导论》 (Introduction to Automata Theory, Languages, and Computation) 作者: John E. Hopcroft, Rajeev Motwani, Jeffrey D. Ullman



# THANKS

 极客时间 | 训练营