

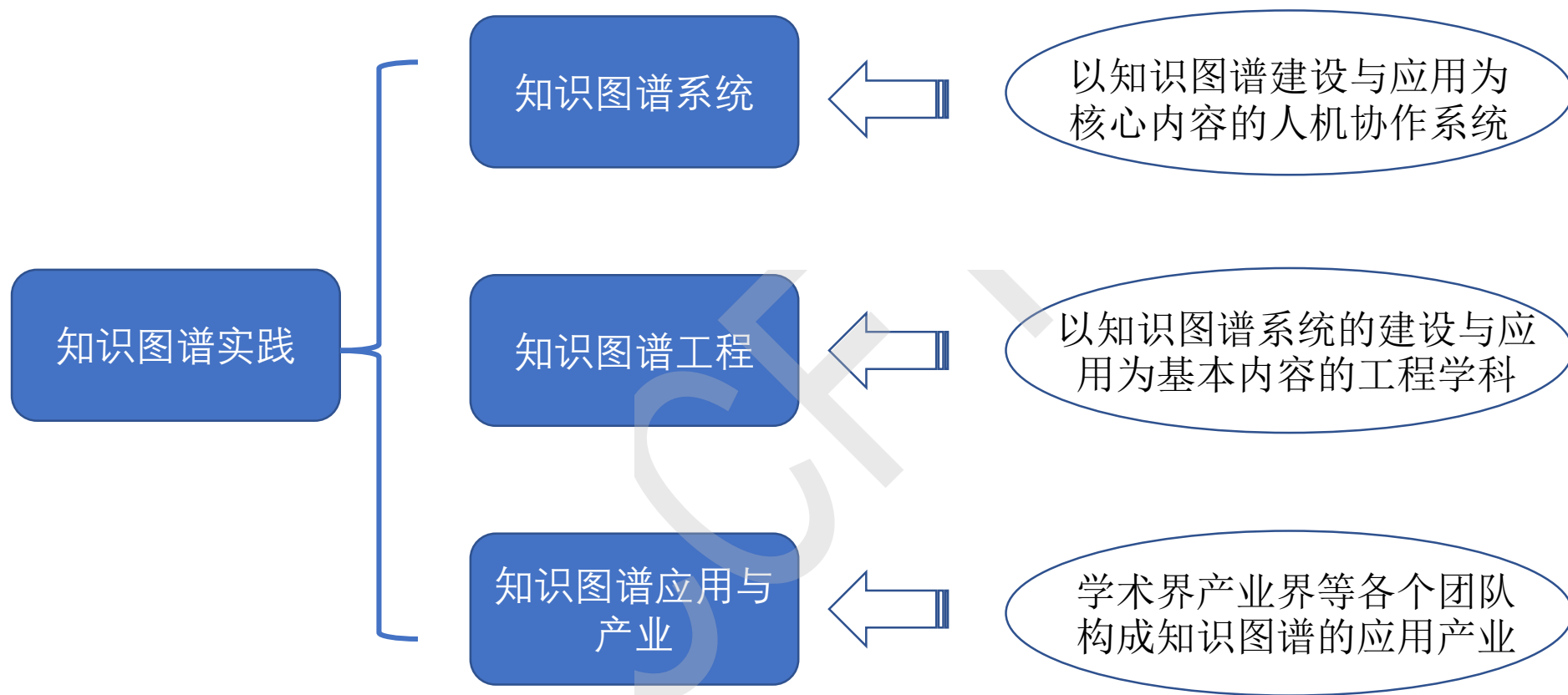
# 知识图谱工程实践的基本 原则与“最佳”实践

肖仰华

复旦大学

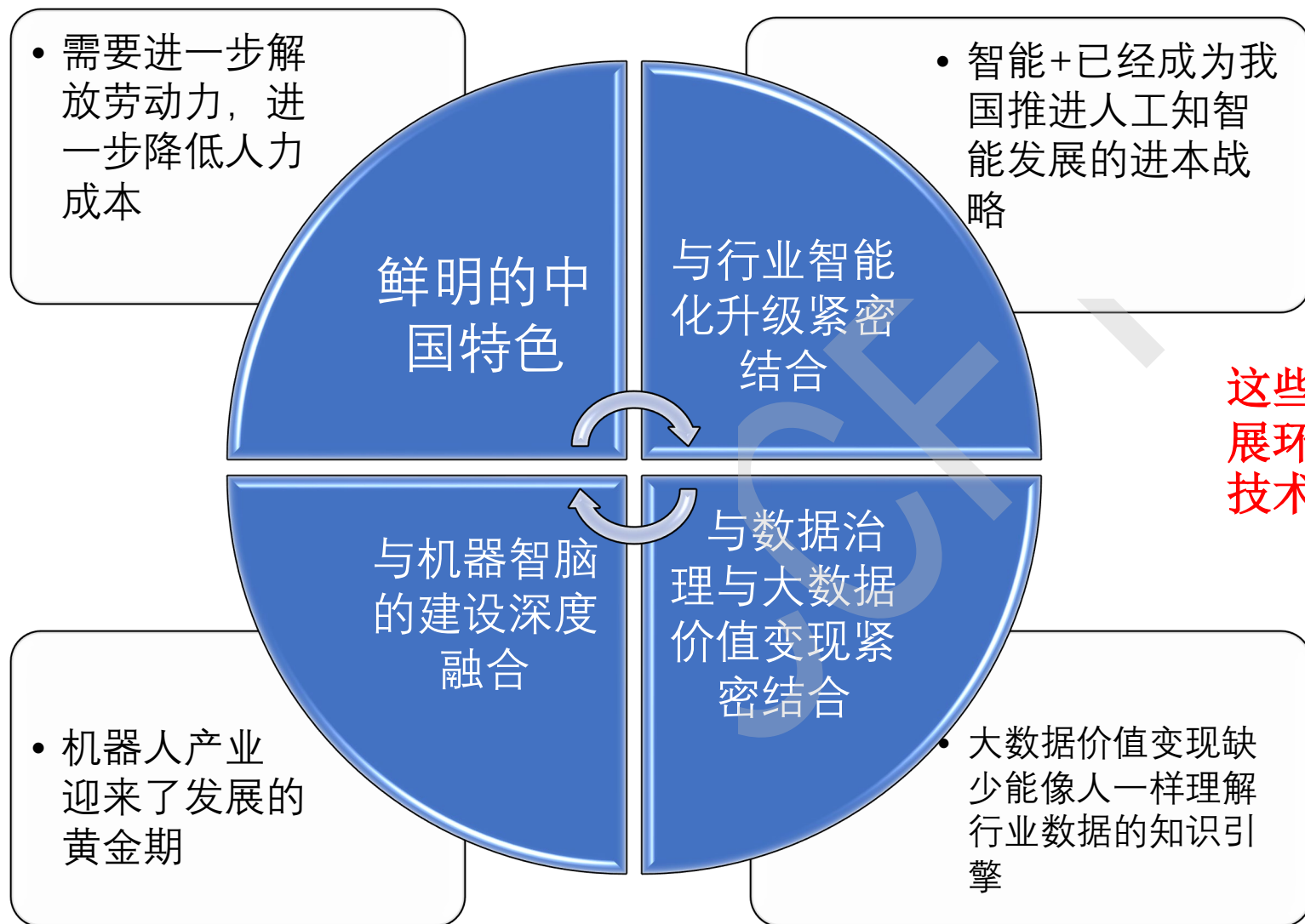
知识工场实验室

# 概览



知识图谱实践的相关问题

# 知识图谱应用的推动力



这些特点背后体现了当下的宏观发展环境以及技术生态对于知识图谱技术需求的迫切性。

# 我国知识图谱应用现状

数据与  
服务

IBM Watson、微软认知服务、百度大脑平台等

产品与  
系统

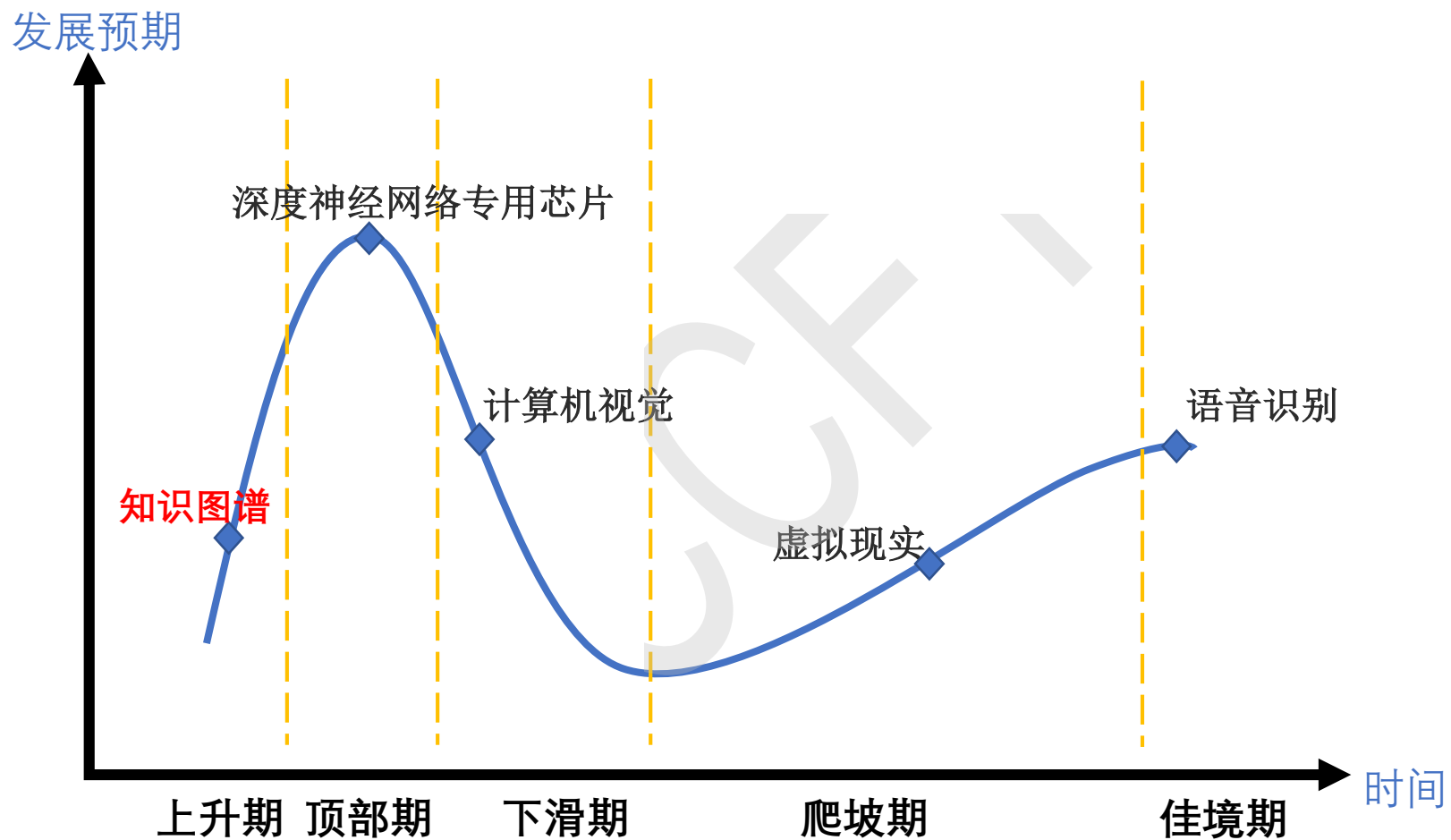
大规模的分布式爬虫系统、图数据库系统等

咨询与解  
决方案

企业知识资源严格保护，知识图谱标品化的服务与产品仍然稀缺

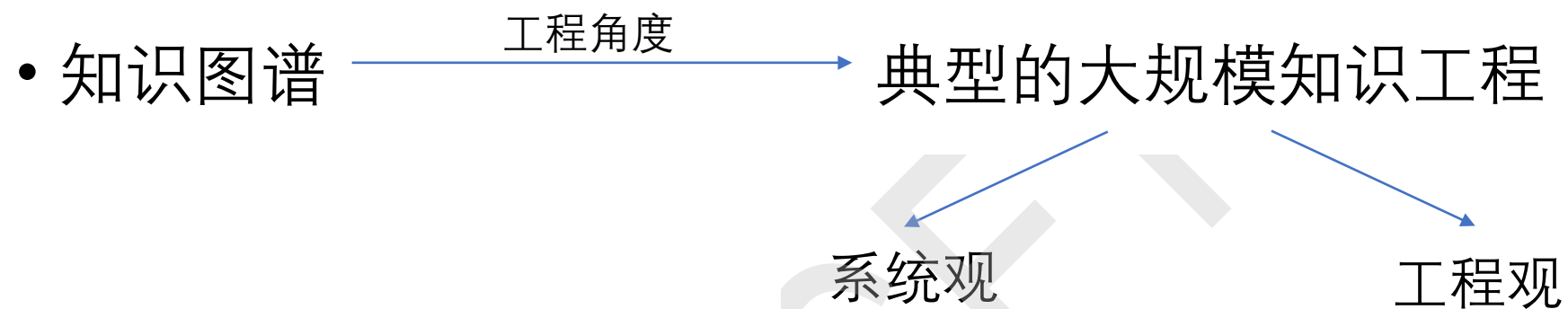
知识图谱的三种典型产业形态

# 知识图谱技术成熟度



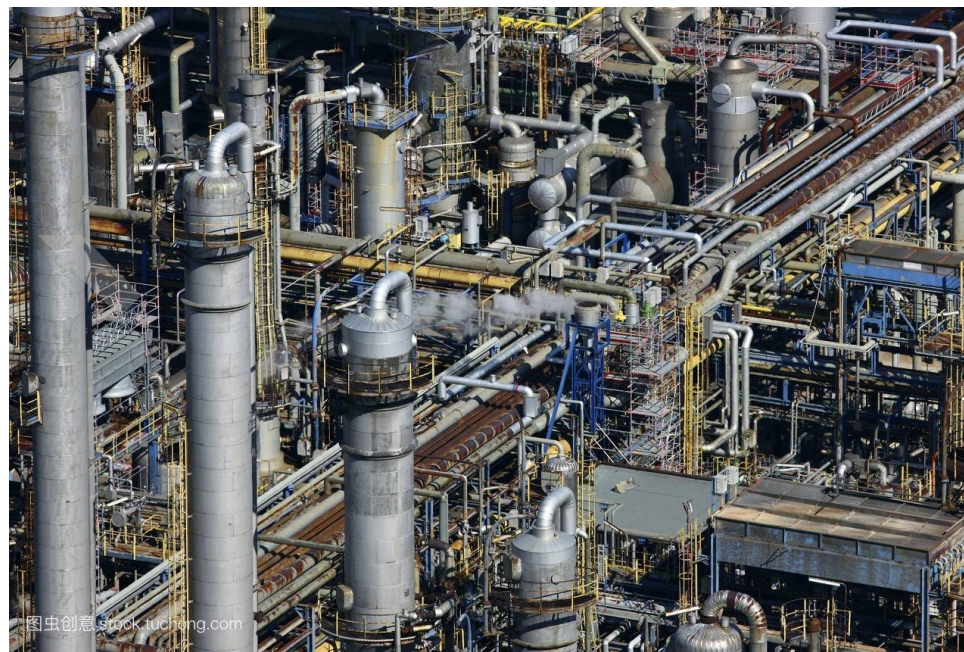
人工智能成熟度曲线

# 知识图谱的系统工程观念



# 知识图谱的工程观

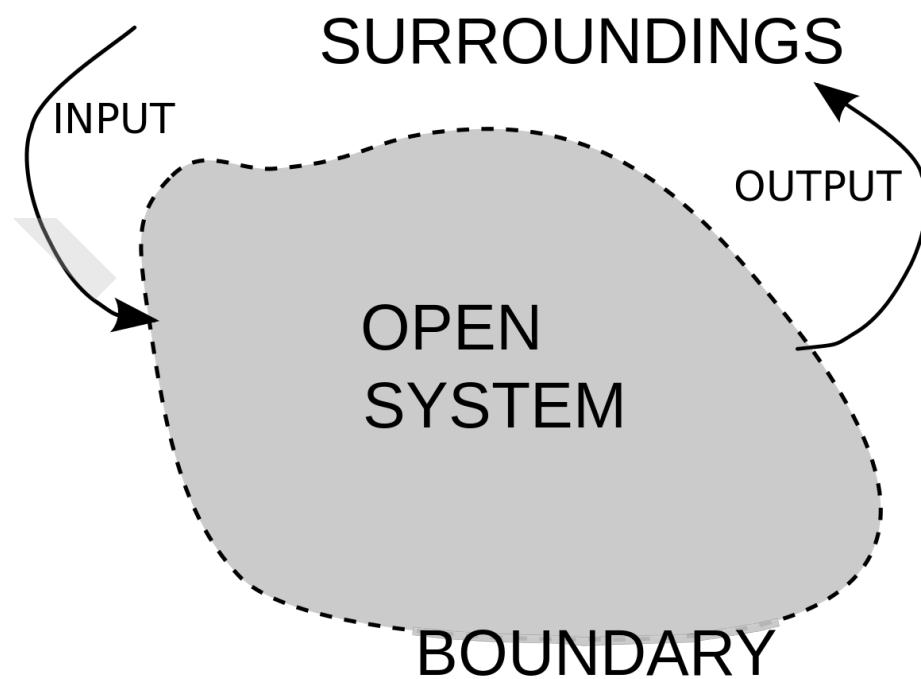
- 知识图谱的工程观
  - 利用数学和科学原理提出解决实际问题的有效方案的观念
  - 带约束的最优化问题求解思路
  - 实践的重要性



知识加工与石油萃取

# 知识图谱的系统观

- 知识图谱系统是“由相互作用相互依赖的若干组成部分结合而成的，具有特定功能的有机整体”
- 与传统知识工程的差别：知识的规模化表示与应用
- 特征
  - 涌现性
  - 交互性
  - 演化性



知识图谱系统本质上是一类开放系统



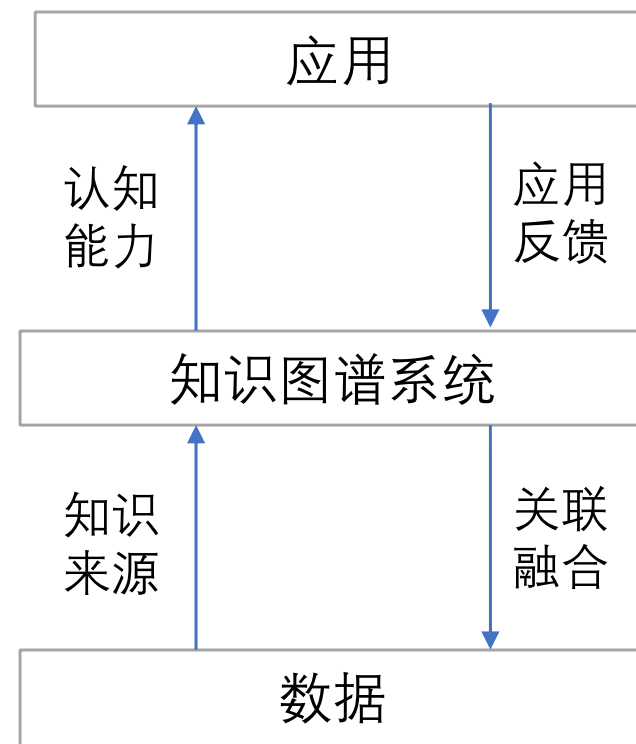
# 知识图谱系统的外部环境

- 向下统摄数据

- 各业务系统数据是知识图谱构建的来源
- 知识图谱中为各业务数据的关联与融合提供了支撑
  - “Id”与“身份证”

- 向上支撑应用

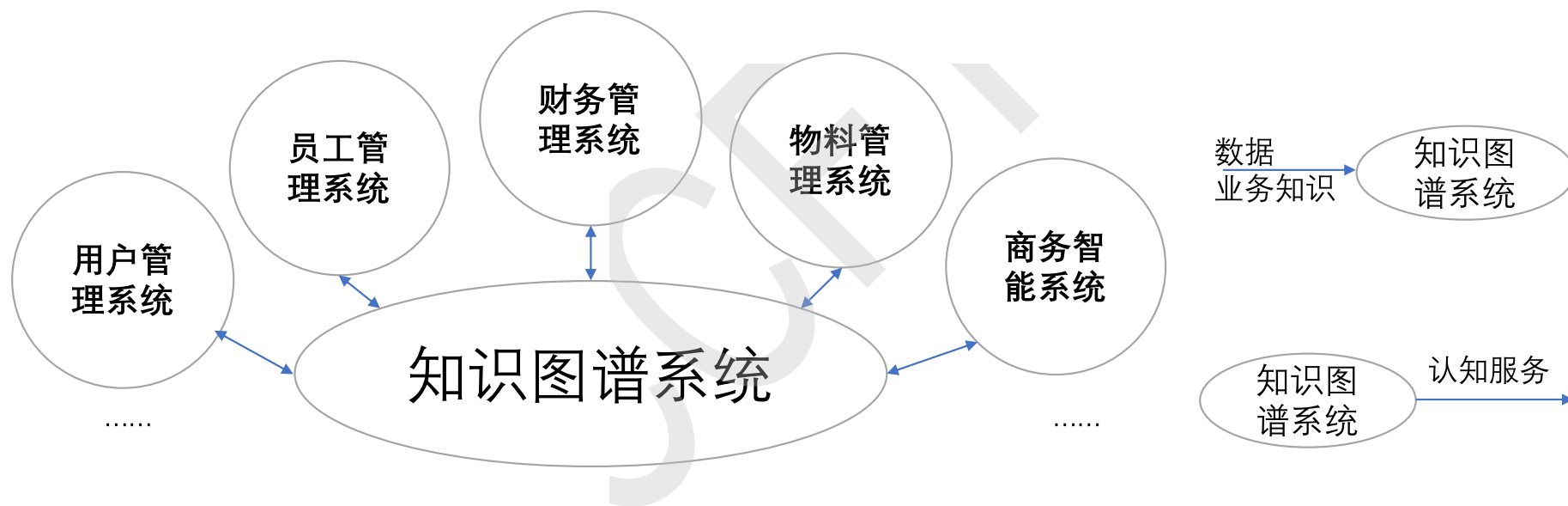
- 认知服务支撑行业智能化升级
- 各类应用为知识图谱系统提供应用反馈



知识图谱系统向上支撑应用向下统摄数据的核心地位

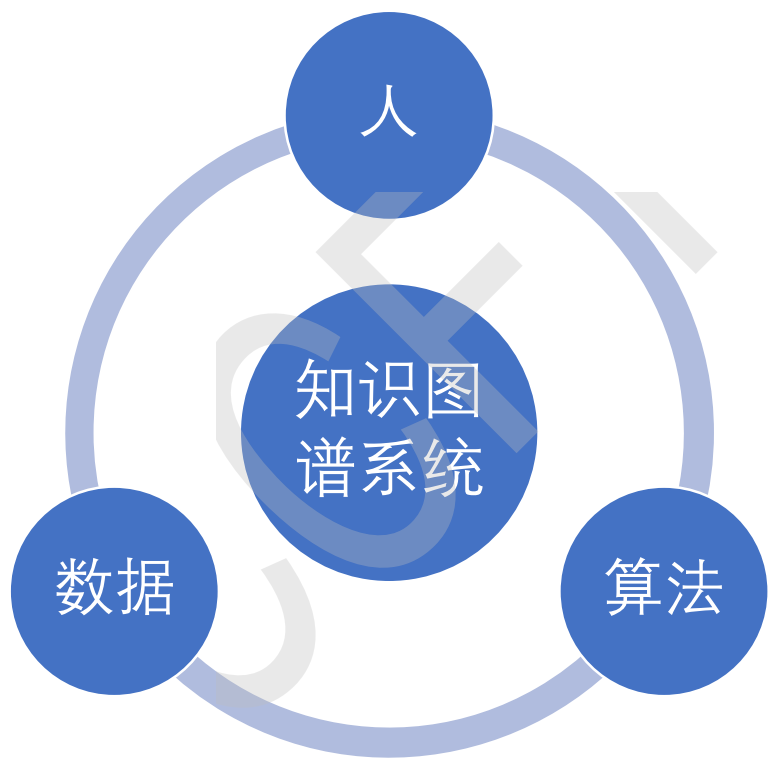
# 知识图谱系统的外部环境

- 在企业环境下，知识图谱系统与其他业务系统关系密切



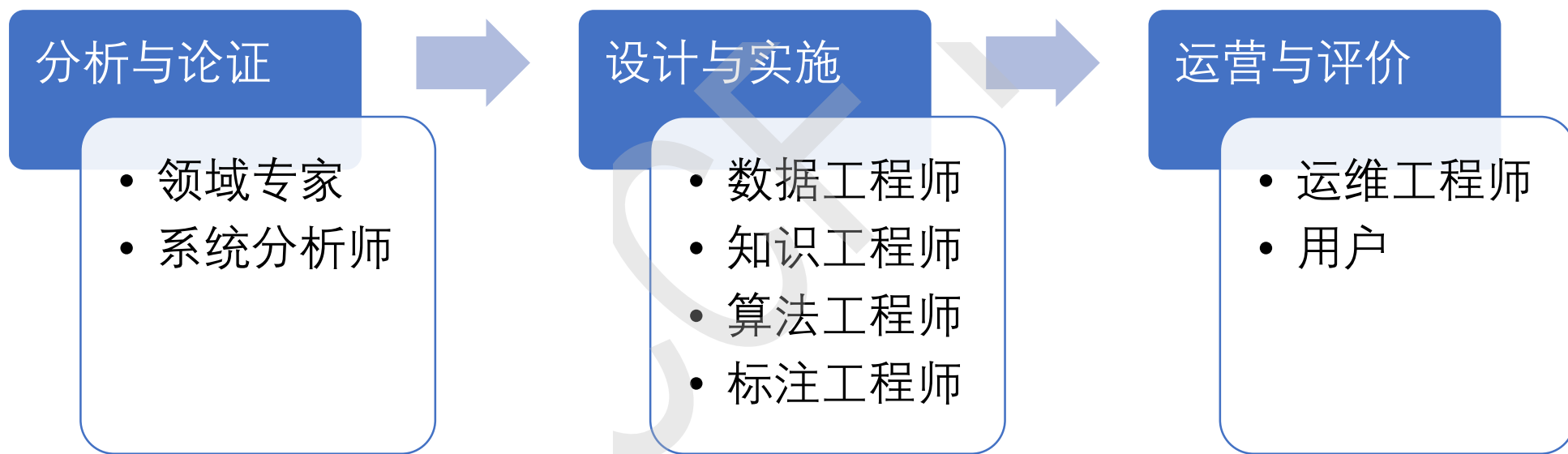
知识图谱系统与其他业务系统之间的关系

# 知识图谱系统关键要素



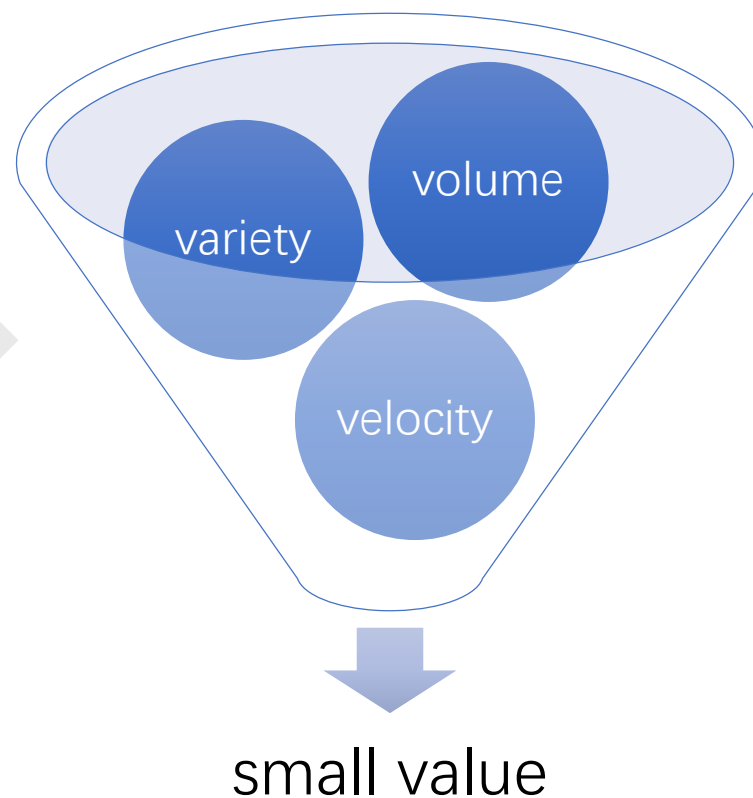
# 知识图谱系统关键要素-人

- 人是知识图谱系统的发起者、设计者、建造者与评估者



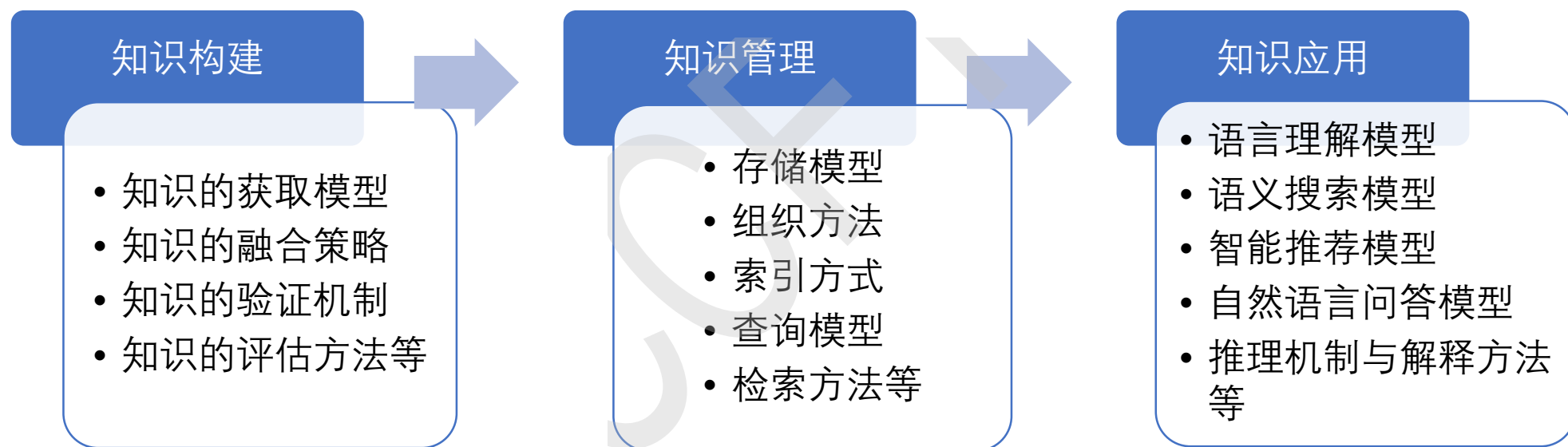
# 知识图谱系统关键要素-数据

- 作为知识图谱来源的各类数据
- 复杂多样
  - 类型众多
  - 来源众多
- 大规模自动化知识获取是当前大数据知识工程的巨大挑战

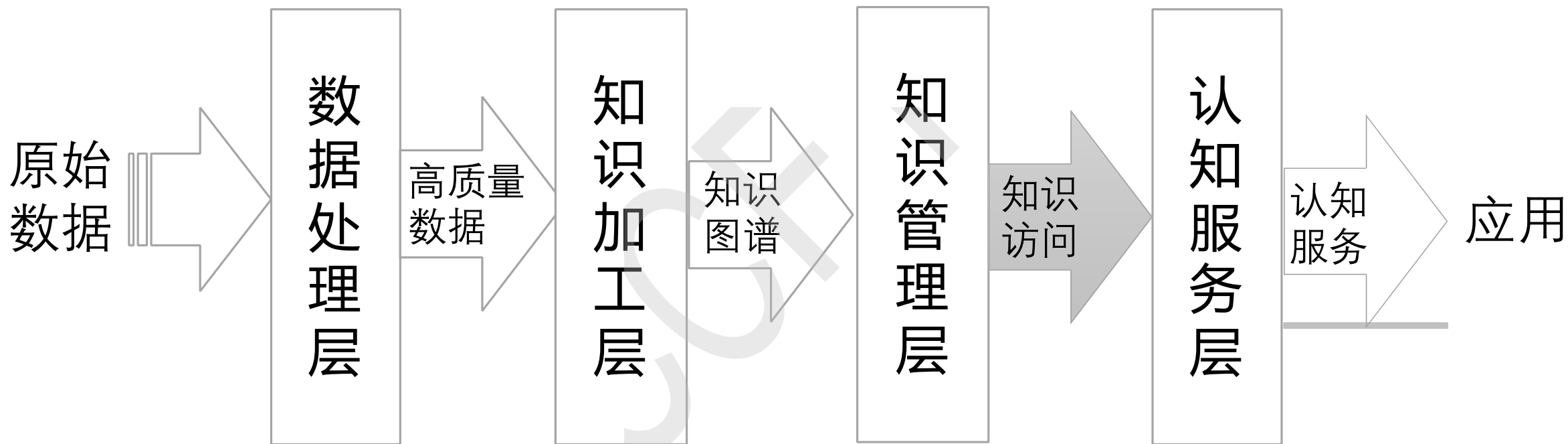


# 知识图谱系统关键要素

- 知识图谱系统整个生命周期涉及的自动化计算过程、模型、策略

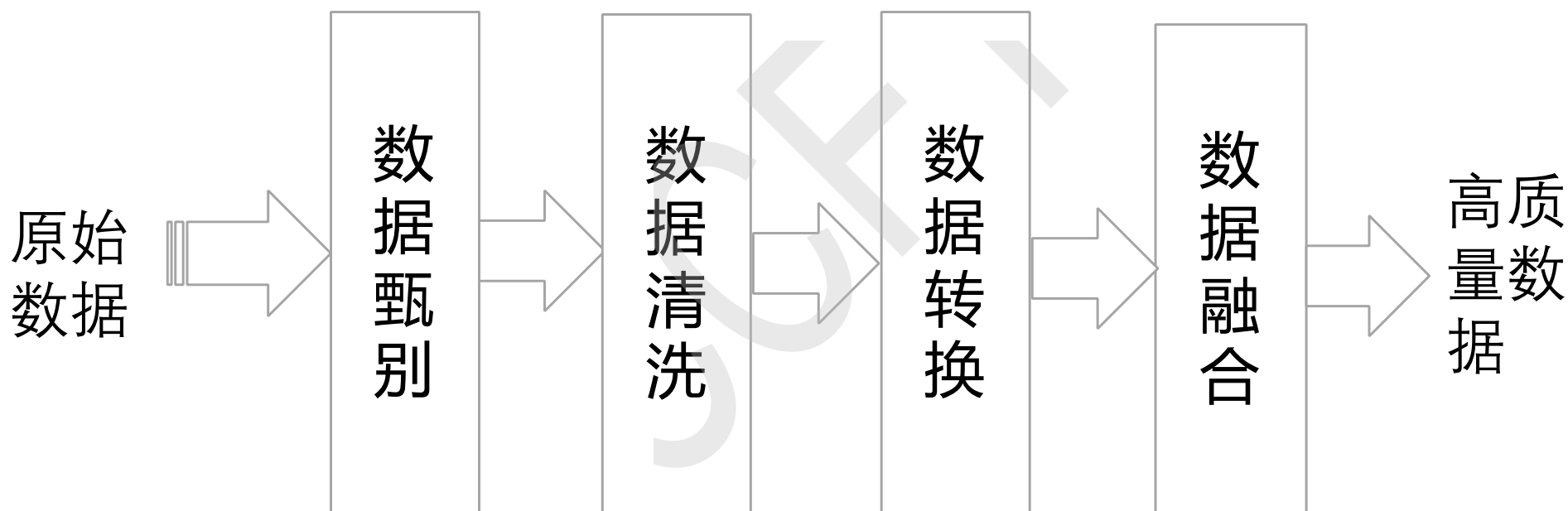


# 知识图谱系统的典型架构



# 知识图谱系统的典型架构-数据处理层

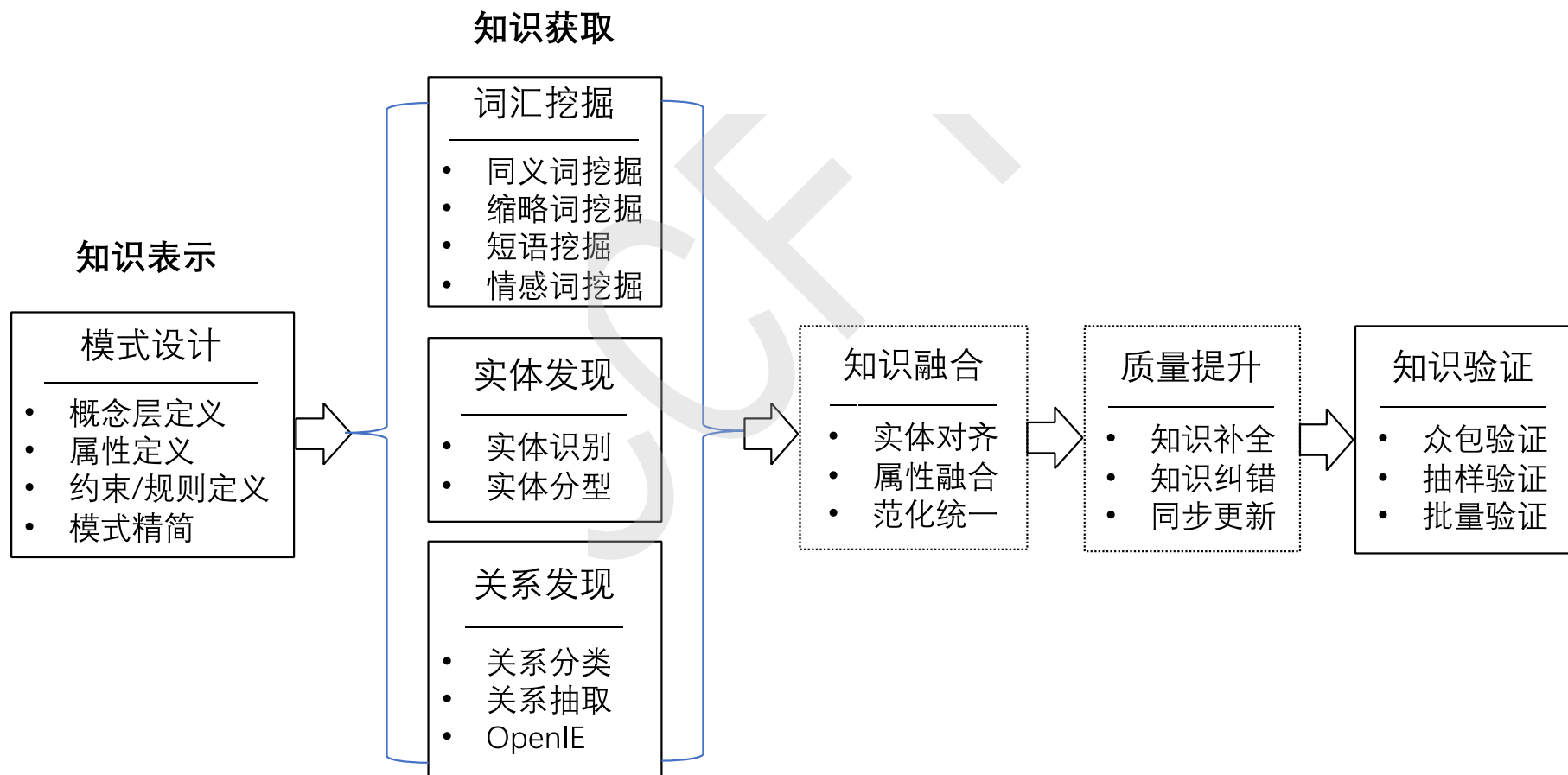
- 接受原始数据作为输入，经过数据处理形成高质量的数据





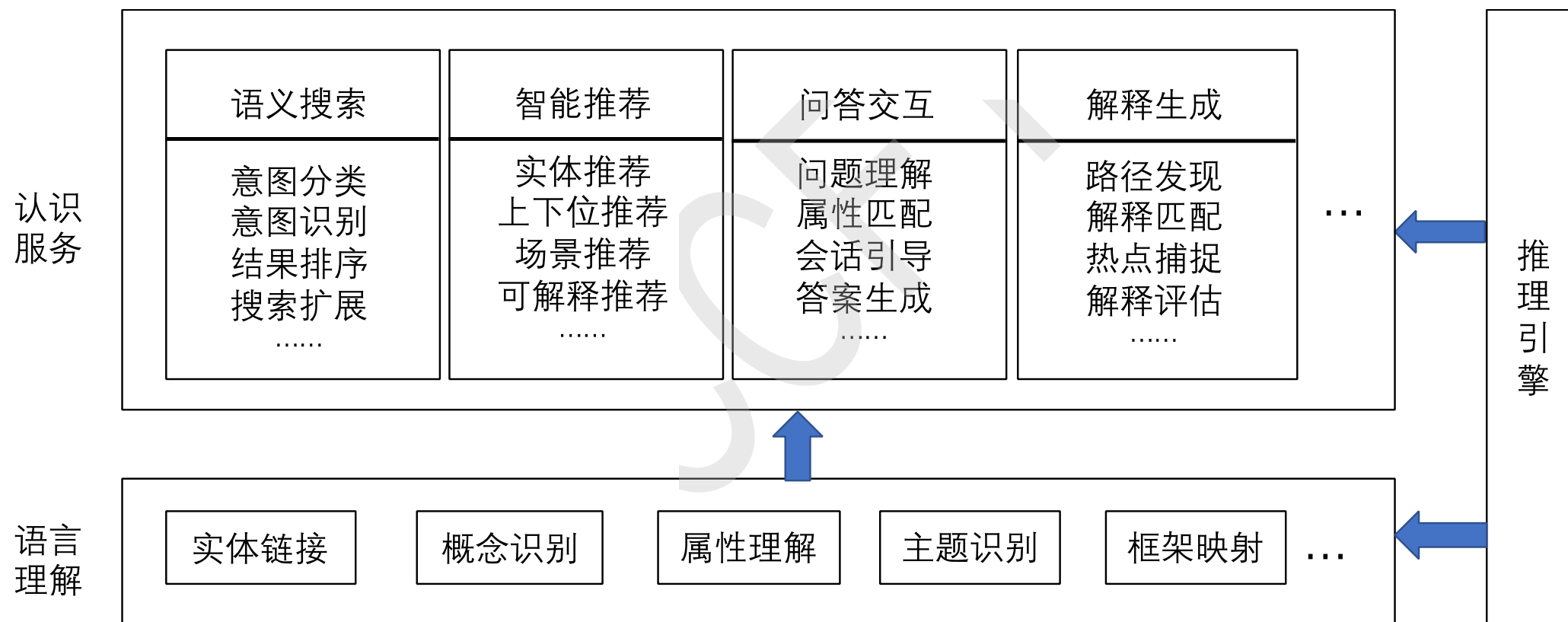
# 知识图谱系统的典型架构-知识加工层

- 知识图谱系统的核心：接受数据处理层形成高质量的数据作为输入，输出高质量的知识图谱



# 知识图谱系统的典型架构-认知服务层

- 提供认知能力，包括语言理解、认知服务、推理引擎



# 推理引擎

- 弥补知识的缺失，提升系统的智能程度

## 知识图谱+推理规则

- 定义“父亲的父亲是爷爷”这样的规则，从知识图谱的基础事实进行推理

## 基于知识图谱的分布式推理

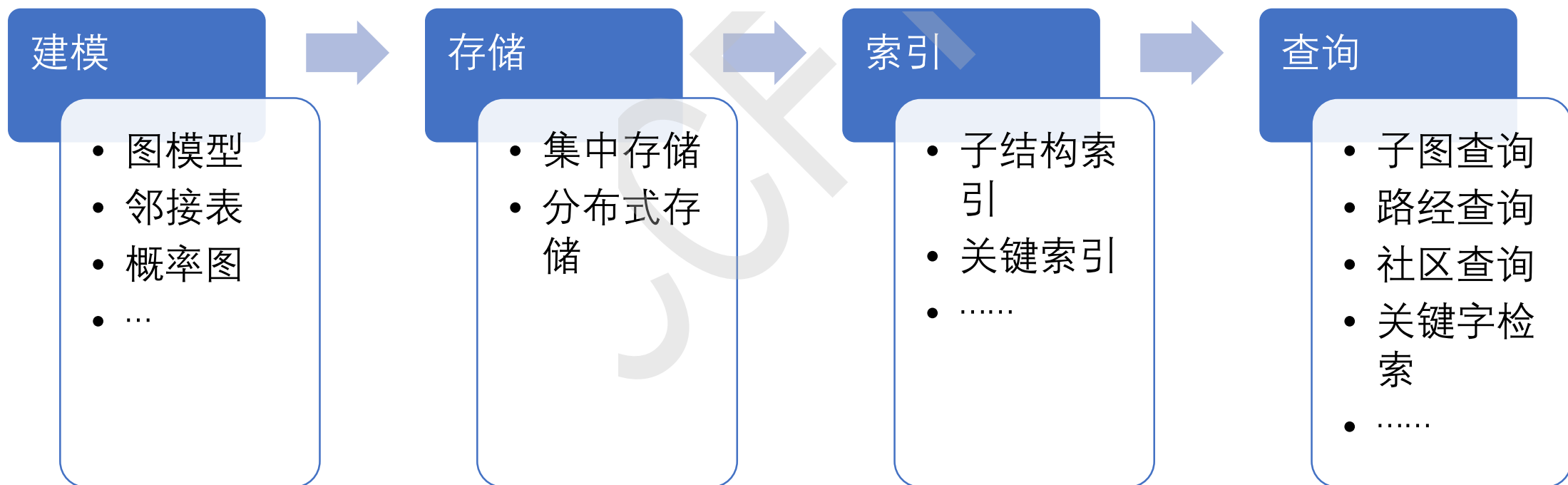
- 给定两个实体 $h$ 与 $t$ 的向量表示(比如 $h$ ,  $t$ )，如果 $h$ ,  $t$ 的向量距离足够相近，则推断 $h$ 与 $t$ 语义相近，甚至其间存在一定的语义关系。

## 基于知识图谱上的显式推理

- 当两个实体 $h$ 与 $t$ 在知识图谱之间存在多条可达路径，且路径上的语义关联强度足够大，则推断 $h$ 与 $t$ 语义相近。

# 知识图谱系统的典型架构-知识管理层

- 实现知识图谱数据的有效管理和高效访问



# 知识图谱工程的基本原则

- 合理定位

- 期望过高，或者期望明显高于当前技术水平会带来不良后果
- 必须心怀敬畏

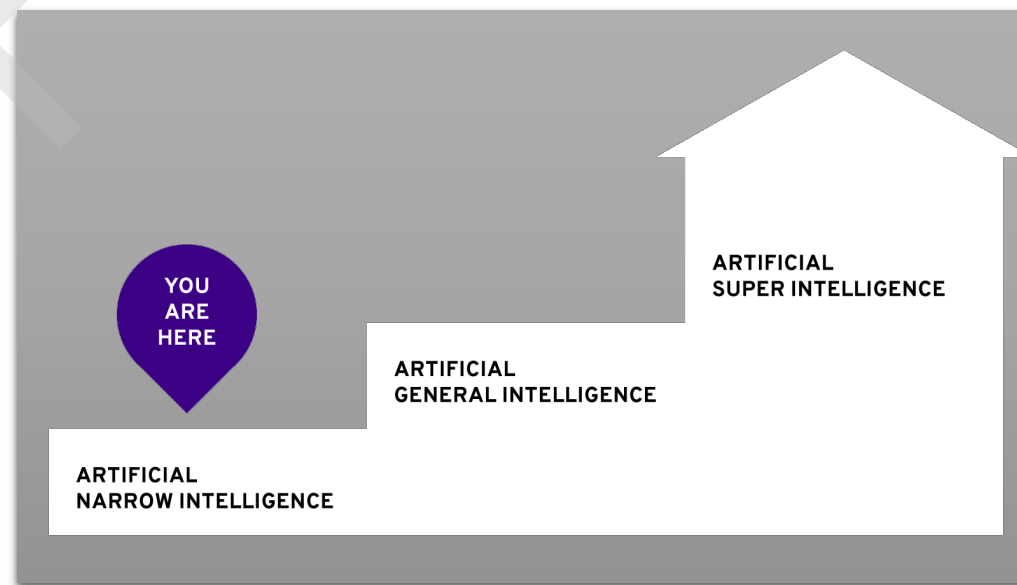
- 任何一个普通人在知识方面所具有的智能，都是当前机器所无法企及的
- 代替专家助理的工作是个合适的目标
- 代替领域专家的工作仍然十分苦难
  - 专家的很多知识是隐性的，难以言明的，难以外化的。
  - 专家的思维方式、知识适配、异常处理均是机器无法企及的

# 知识图谱工程的基本原则

- 应用牵引

基于知识图谱的认知智能还没发展到普适、通用智能的阶段  
知识图谱技术平台化发展的异常艰难

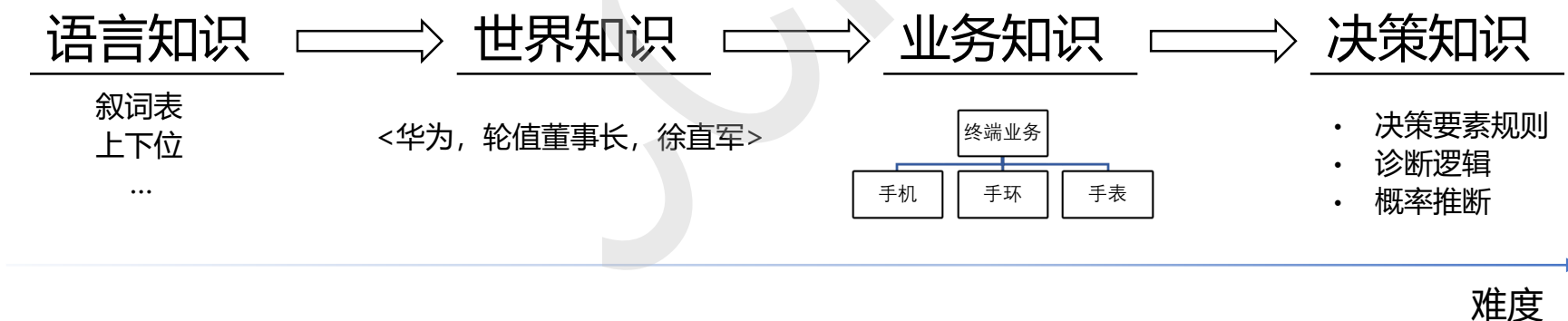
- 互联网时代：业务模式简单且具有同质化，技术与业务的平台化受到追捧
- 当下人工智能的发展多以场景化应用为主



# 知识图谱工程的基本原则

- 循序渐进

先从结构化程度高的数据中抽取出易于获得的语言知识，再从半结构化数据中抽取出世界知识，进而总结出业务知识，最后再触及决策知识



# 知识图谱工程的基本原则

- 先简后难

知识图谱各项技术成熟程度不均衡是当前知识图谱产业实践的基本情形，决定了先简后难的推进策略

- 知识图谱技术体系复杂多样，包括知识表示、知识抽取、知识融合、知识推理、知识存储和知识检索等。
- 每类关键技术的成熟度不同，有的已投入实用化，有的仍处于学术研究阶段。



先摘成熟的果子，坐等其他苹果自然成熟再行采摘



# 知识图谱工程的基本原则

- 由粗到精

粒度越细表达越精准，知识获取的难度也越大，越容易存在模糊性与不确定性  
由粗到精，逐步求精

“机动车变道，应打开相应的变道指示灯”



IF (机动车变道) THEN (打开相应的变道指示灯)



IF (机动车变道) THEN (打开相应的变道指示灯)

条款级

规则级

实体、动作级

汉堡是食物？



汉堡是健康食物？



# 知识图谱工程的基本原则

- 求同存异

认知的模糊性导致认知主观性，细微的主观性差异。  
暂且搁置争议，先解决无争议的问题

- 不同人对于“高个子”认识不同，但是没有人会认为2.2M还不是高个子。
- 有争议的事实，可使用数据驱动的方法来自定义
  - 对于“矮个子NBA球星”，如果大部分用户在这一搜索关键词下，点击的球星都在1.8m以下，那么1.8M以下对与NBA球星而言或许就是矮个子。

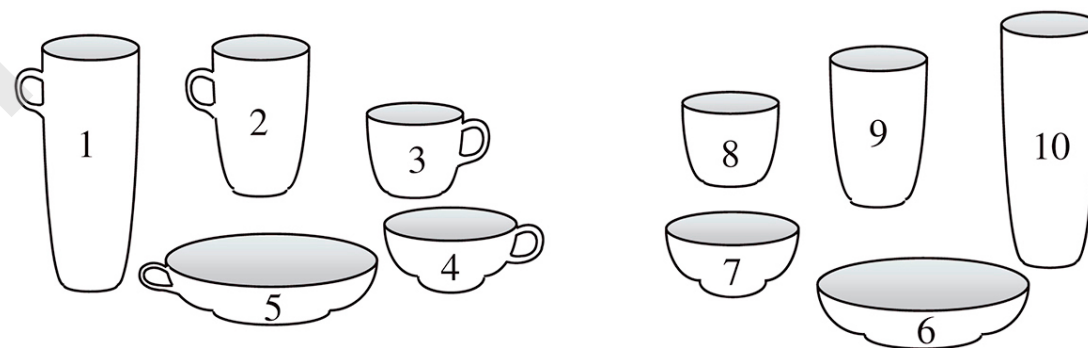


Figure 11.2 Cups, bowls, vases

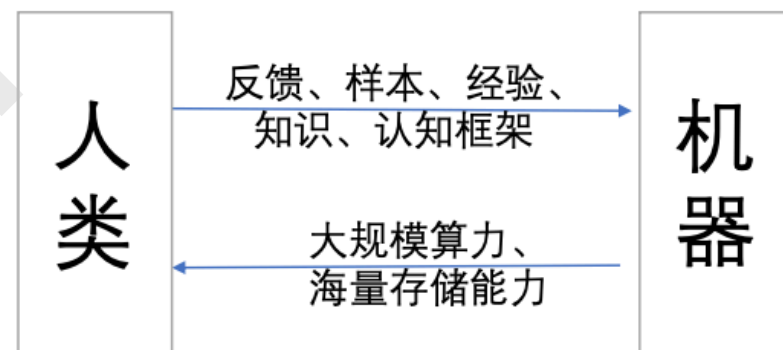
© Sebastian Löbner 2013

# 知识图谱工程的基本原则

- 人机协同

- 知识图谱落地，需要机器和人，二者不可或缺
- 人在环中是当前AI的基本形态

- 认知世界的框架、样本、知识均需要人类给予机器



Human-in-the-loop

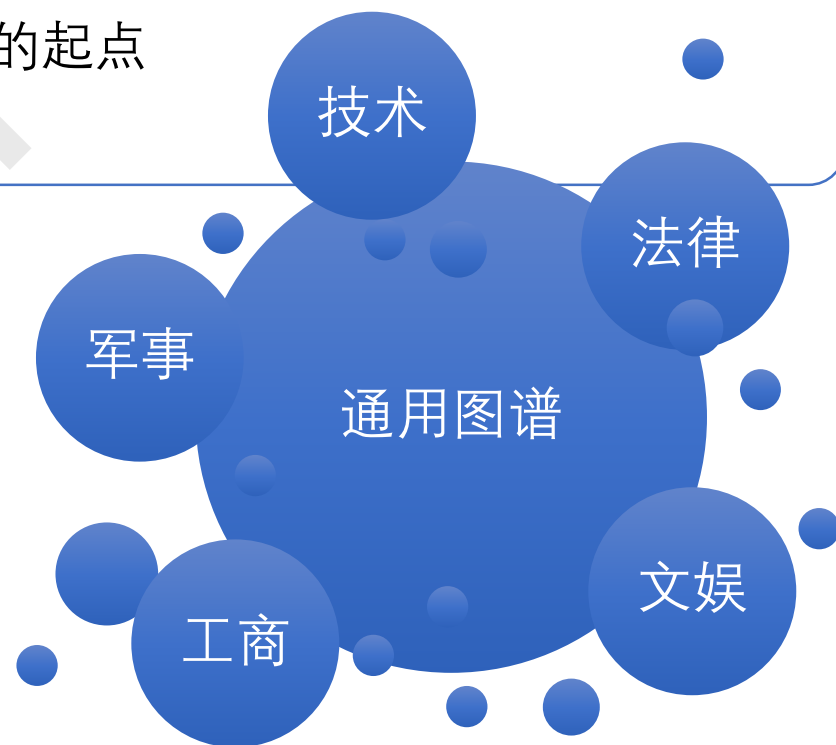
图 15.12 人在环中的人工智能发展模式

# 知识图谱工程的基本原则

- 快速启动

- 充分利用已有资源，提高领域知识图谱构建的起点
- 尽可能复用是知识资源建设的重要策略之一

- 从通用领域导出
- 从临近领域迁移



# 知识图谱工程过程模型



## 分析与论证

明确知识图谱的应用目标，分析知识图谱的业务价值，论证知识图谱项目上线的必要性；对所设定目标所涉及的数据资源、人员投入等、技术可行性等进行评估，作出可行性评估；对于整个知识图谱工程项目的进行计划。

## 设计与实施

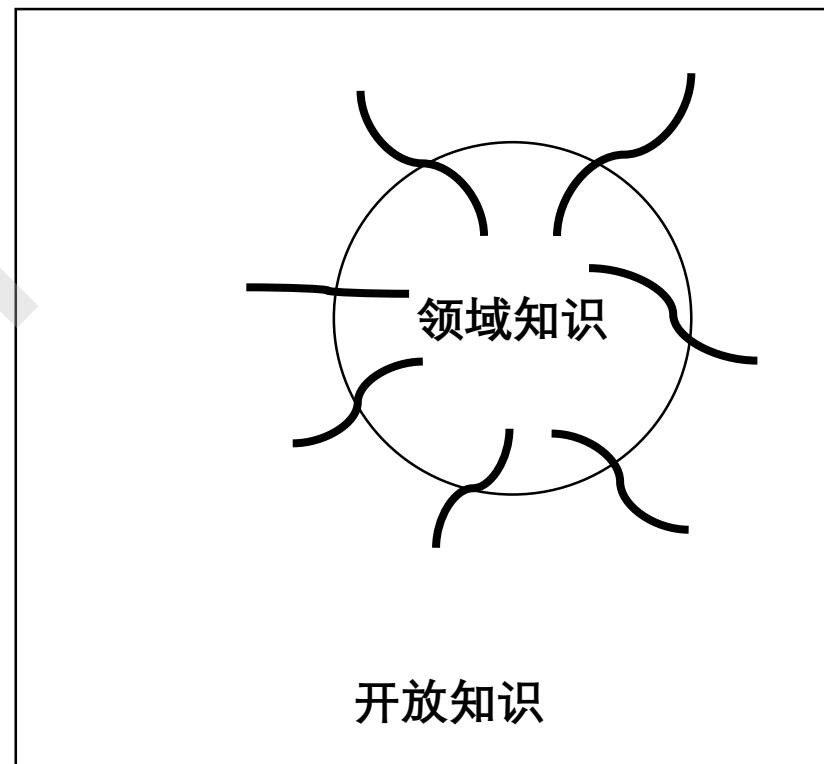
对知识图谱系统相关的数据库、数据流程、系统架构、关键算法、系统选型等等进行设计，制定详细的设计方案。进行代码开发，实现相关算法，集成相关系统，完成系统上线。

## 运营与评价

获取用户的使用日志、评估反馈是十分关键的。这些反馈与日志是知识图谱工程持续演进的重要依据，是下一轮建设周期的发动机。

# 知识图谱工程可行性分析

- 是否是封闭应用
  - 一切实体都身处在一个复杂的因果网络中，世界是普遍关联的
  - 很多领域的应用不是封闭的。
- 基于金融知识图谱的关联分析往往会牵扯出几乎万事万物。



行业应用中的知识需求难以封闭于  
预设的领域知识边界内

# 知识图谱工程可行性分析

- 是否涉及常识
  - 常识难以建模
  - 常识难以获取
  - 机制尚未探明

## 例一：对话机器人 人工“智障”

问题与痛点

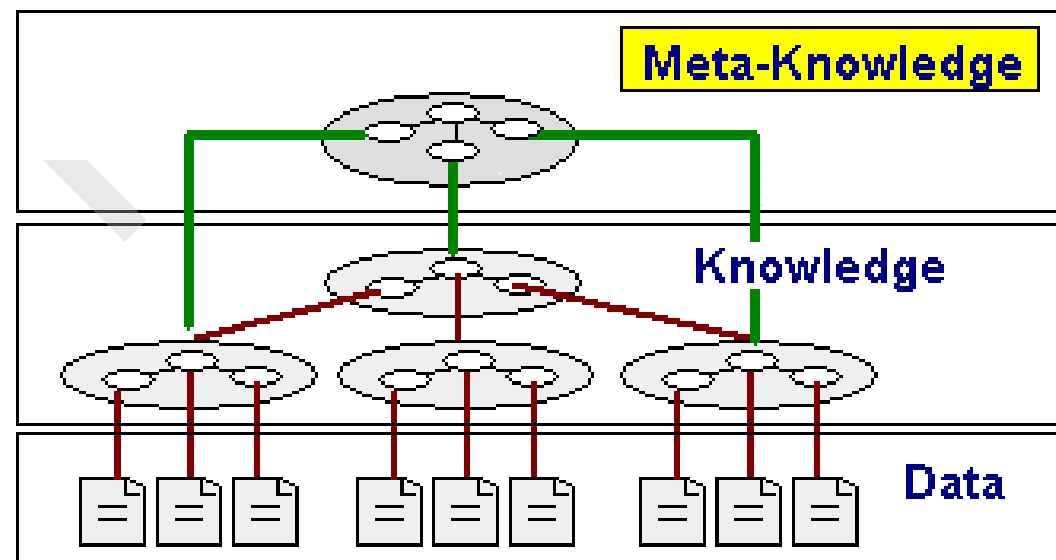
## 例二：生活安全

失火了，猫和人应该先救  
哪一个？



# 知识图谱工程可行性分析

- 是否涉及元知识
  - 元知识是指有关知识的知识。
    - 包括属性的领域（Domain）与范围（Range）
      - 比如“父亲”作为属性发生在人物这个类别的实体上（这是在制定Domain），取值也只能是个人物。
    - 领域内的约束
      - 比如父亲都必须比子女年龄大
    - 如何使用知识的知识，
      - 比如吃了不洁净的物品呕吐了，我们立即就会判断有可能是因为不洁饮食导致的食物中毒。





# 可行性检查列表-1

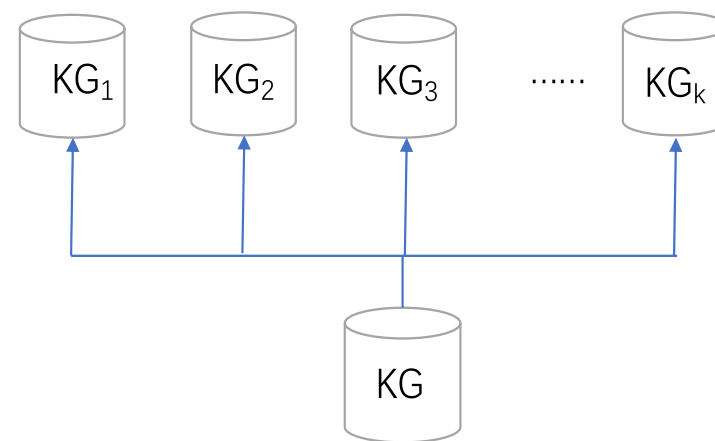
	考察角度	问题列表
数据禀赋	数据完整	Q1:知识图谱所需数据是否完整（或存在缺失）？
	数据质量	Q2:数据格式是否统一？ Q3:数据来源是否多样？ Q4:数据库类型是否多样？ Q5:数据缺失、完整性、一致性等情况如何？
	数据结构化程度	Q6:数据的结构化程度？ Q7:数据是否已经电子化？是否需要OCR？ Q8:数据是否存在明显的结构化标记？
	数据可用性	Q9:是否存在行业壁垒？ Q10:是否存在国家安全红线？ Q11:是否涉及个人隐私？

# 可行性检查列表-2

	问题列表
应用复杂性	<p>Q12:是否用到常识?</p> <p>Q13:是否用到元知识?</p> <p>Q14:是否单一问题模型即可建模（比如分类或者回归）?</p> <p>Q15:是否涉及长程推理?</p> <p>Q16:用到的知识类型是否多样?</p> <p>Q17:领域专家的学习周期是否很长?</p> <p>Q18:是否简单的岗位培训就能胜任应用需求?</p> <p>Q19:应用是否封闭?</p>
知识复杂度	<p>Q20:知识是否容易发生变化?</p> <p>Q21:是否涉及复杂过程的描述?</p> <p>Q22:是否涉及分支繁复的推理决策?</p>
知识资源积累	<p>Q23:是否存在领域本体?</p> <p>Q24:是否存在叙词表?</p> <p>Q25:是否存在领域词典?</p>

# 知识图谱工程实践建议

视图



- 合理控制不同视角下的不同图谱

- 知识的不同视角往往是程度之分。
  - 比如“龙”，在东方人的视角下是吉祥的，在西方人的视角下往往是凶恶的有贬义
  - 搜索“5G 文档”，市场部门的人更愿意看到售前方案，而技术研发人员可能关心研发资料。“物美价廉的水果”这个品类对于不同人理解完全不同。
- 要针对不同的角色，定制相应的图谱。
- 考虑到图模型的普适性，可以定制不同的权重，以体现不同角色对与不同知识的认知程度。

# 知识图谱工程实践建议

- 区别对待冷启动与热运营两阶段的不同策略

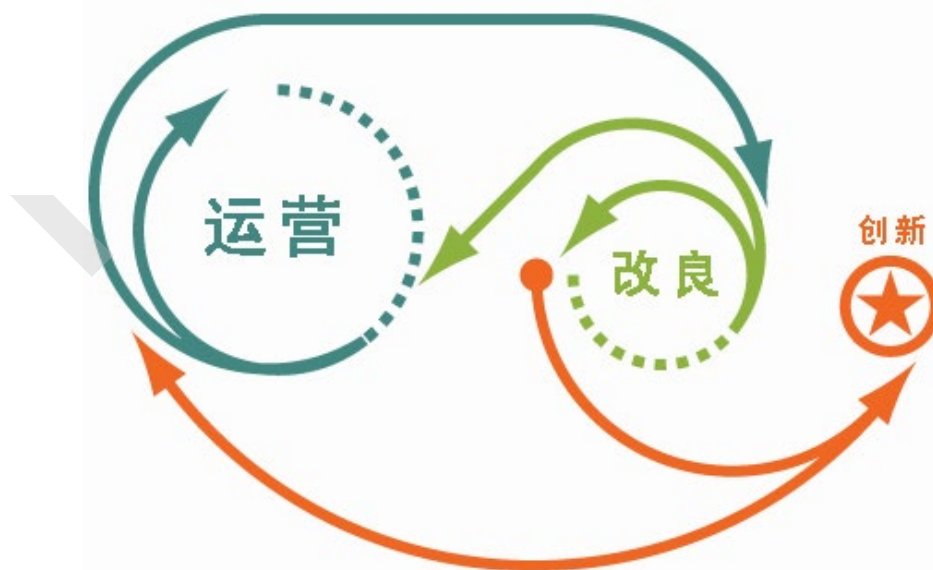
表 15.3 冷启动与热运营的不同策略

	方法	数据	思路
冷启动	基于规则	无用户行为数据	依赖专家知识
热运营	基于学习	丰富的用户行为	依赖数据驱动

# 知识图谱工程实践建议

- **知识图谱应用建设与运营并重**

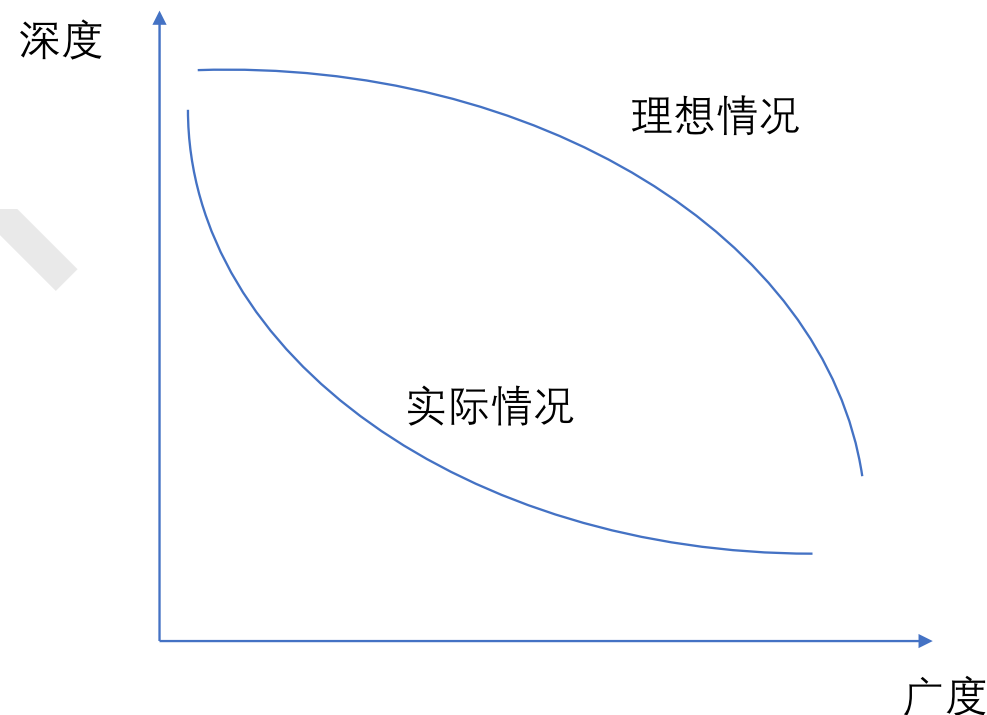
- 智能系统需要经历多轮迭代方能成熟
- 持续运营保持系统处于最佳状态。
- 智能系统必须随着用户的演进而演进，否则容易失效。



# 知识图谱工程实践建议

- **合理处理知识的扁平化与纵深化矛盾**

- 深度知识在安全相关等领域十分重要
  - 恶意意图的行为往往具有隐蔽性等特点，难以简单语义关联所发现。
- 通用知识图谱有着广度，但是缺乏深度，体现在平均关系数小于相应的领域知识图谱。
- 广度与深度需要相互平衡



# 知识图谱工程实践建议

- **坚持迭代式演进路径**

- 螺旋迭代式发展路径是知识图谱工程实践的基本形态。
- 知识库积累与知识抽取模型的迭代发展
- 螺旋上升是在技术集群尚未全面成熟时的典型不成熟时达到目标的典型发展路径。



# 知识图谱工程实践建议

- **区别对待静态知识与动态知识**

- 知识的动态变化是绝对的，静止不变是相对的。
- 绝大部分知识在有限时间内变化的可能性是极低的
  - 比如地球是圆的，在很长一段时间人们对于这个事实的信念是不会发生改变的。
- 一般而言事实是相对易变的，而模式是相对不变的

变与不变



# 版权声明

- 本ppt用到了来自互联网的一些图片，特此说明。
- 本ppt大部分内容与观点来自即将出版的《知识图谱：概念与技术》一书“第三部分：实践与问题篇”第十五章《知识图谱实践》

