



知识图谱在小米的落地与挑战



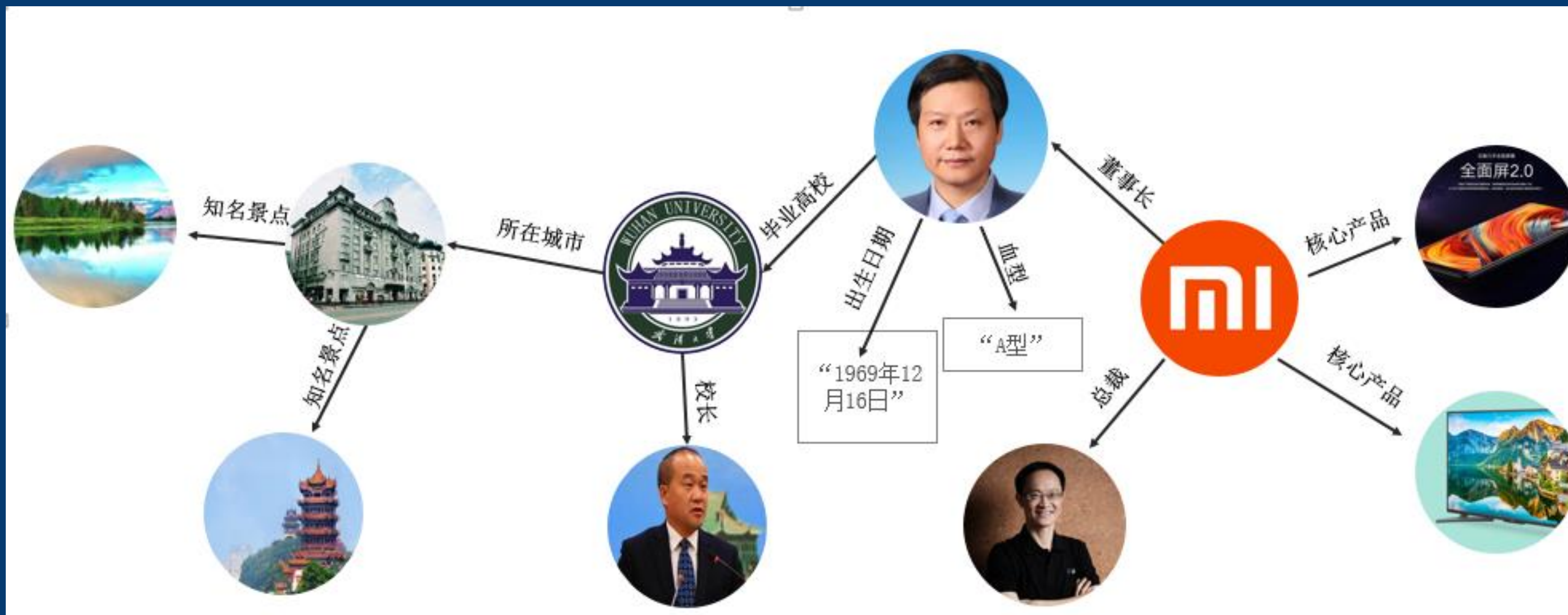
刘作鹏

小米人工智能实验室知识图谱总监

TABLE OF CONTENTS 大纲

- 知识图谱概述
- 知识图谱构建
- 基于图谱的问答
- 图谱的其他典型应用

知识图谱简介



- 本质上，知识图谱是大规模的语义网络（Semantic Web）；
- 语义网络是知识表示的重要方式之一，富含实体、概念和多种语义关系；

语义网络

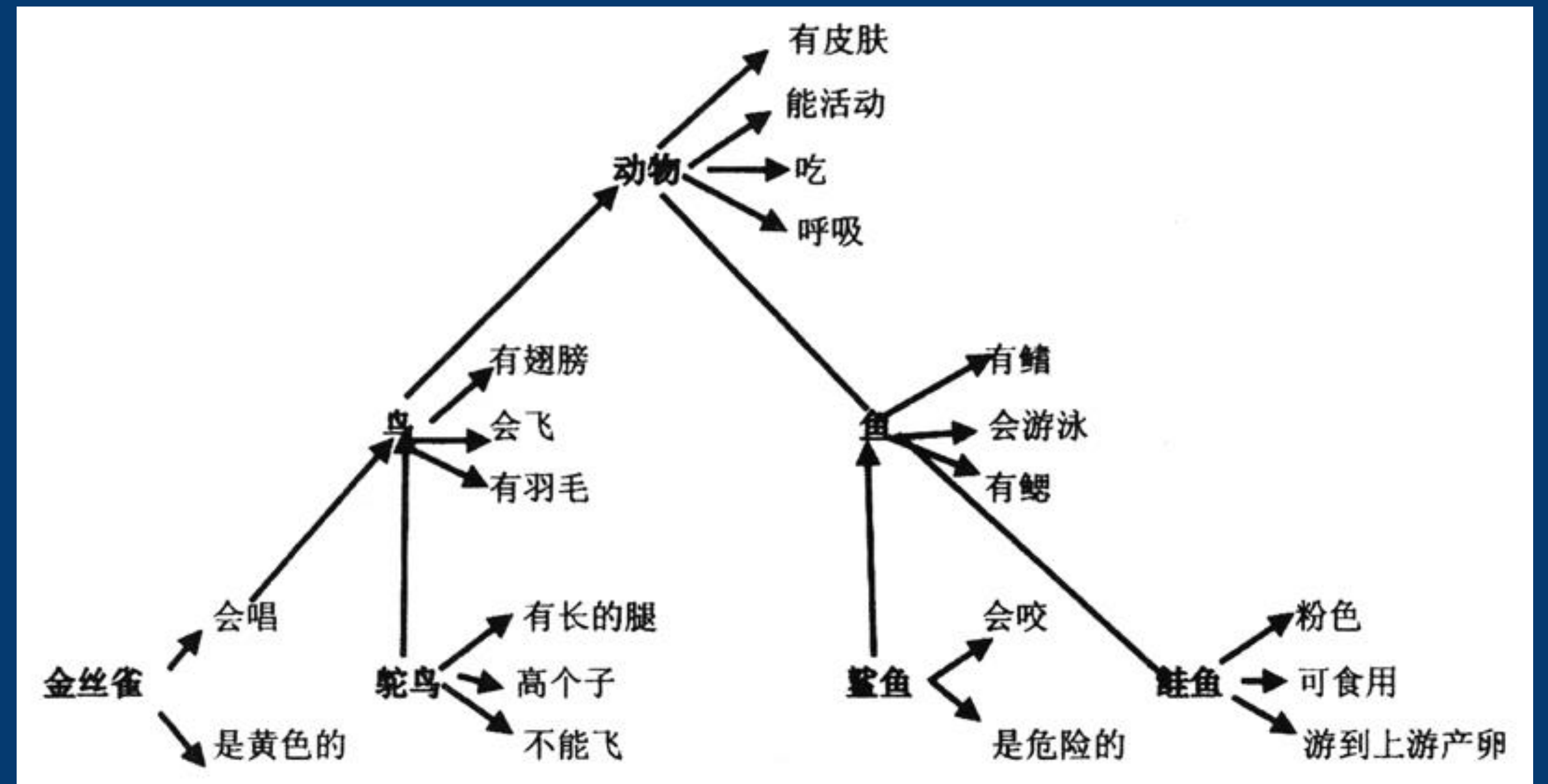
语义网络

知识表示

知识工程

人工智能

语义网络的位置



典型的语义网络

图谱的价值

规模巨大

结构精良

语义丰富

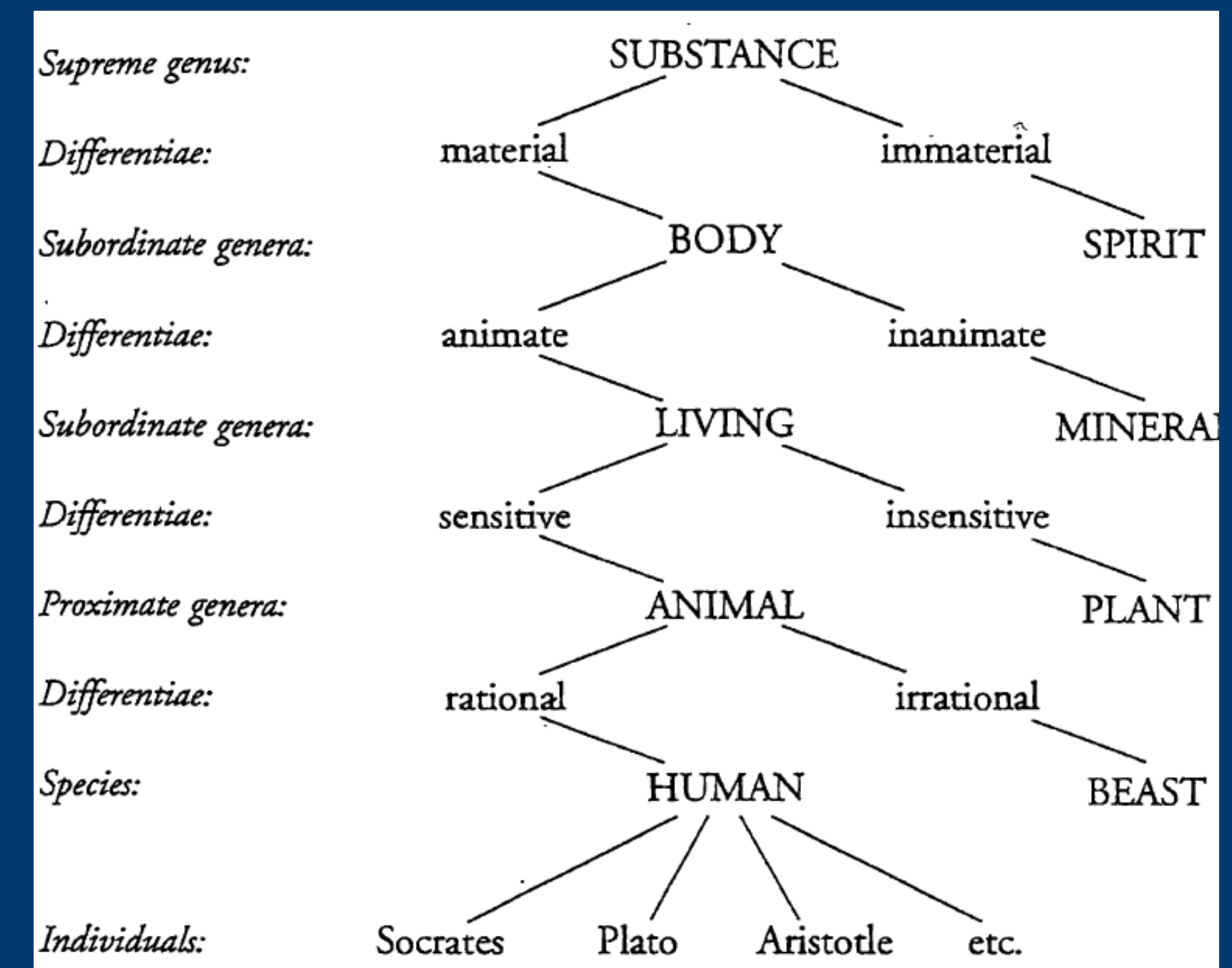
质量上乘

- 知识图谱是让机器具备认知能力的关键技术；
- 同深度学习相比，它在推理、理解、解释等问题的处理上，有独特的优势；

简单回顾

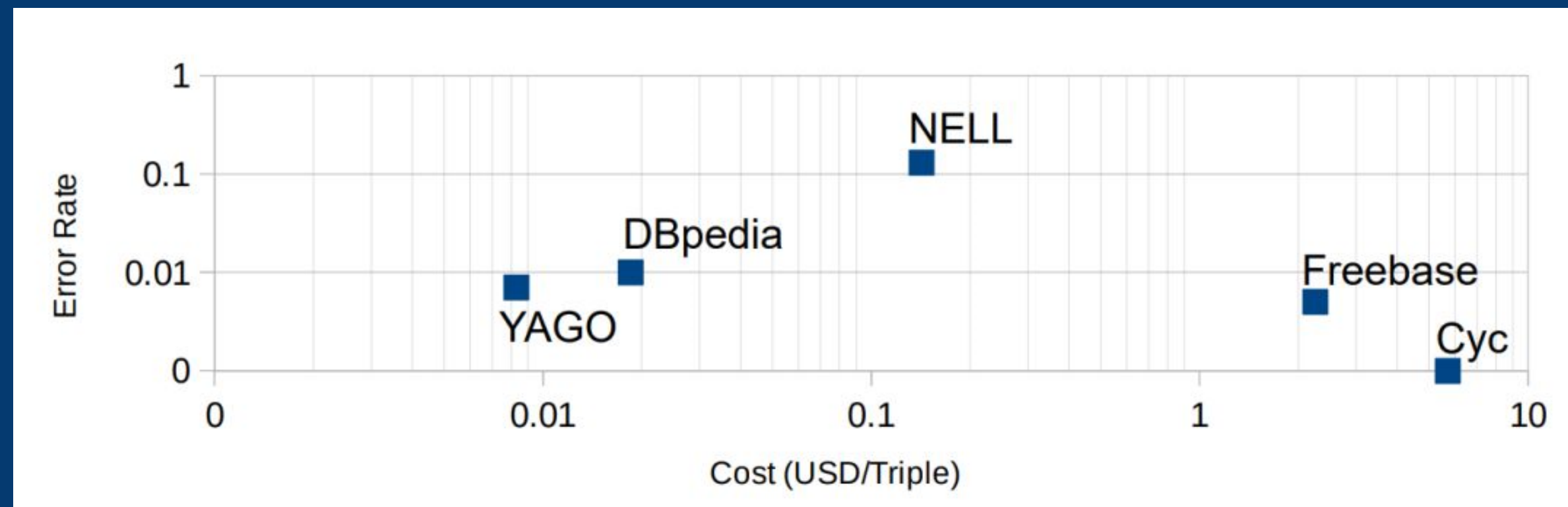
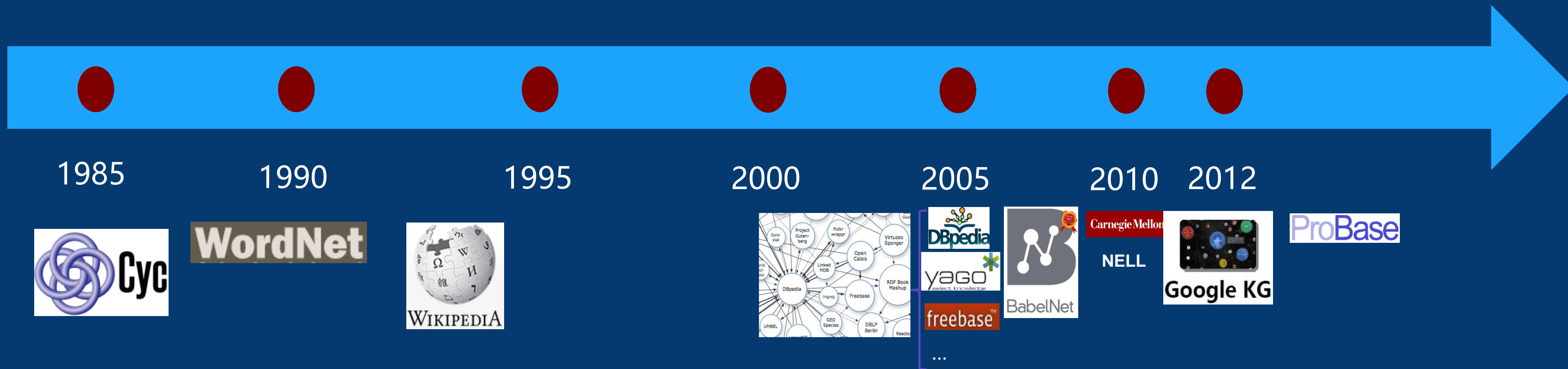


亚里士多德(“百科全书之父”): 三段论 & 形而上学
莱布尼茨: 齿轮计算器 & 符号逻辑



波菲利之树(公元270年)

近代发展



- Cyc 最早通用知识图谱之一，耗时1000人年，成本1.2 亿美金；
- Freebase 有30 亿条英文关系，基于众包，总成本约为 67.5 亿美元；

OpenBase项目



知识的抽取与编辑



知识的挖掘与融合



知识更新



知识众包

联合发起机构



OpenKG



小米集团



清华大学



浙江大学



东南大学



海知智能



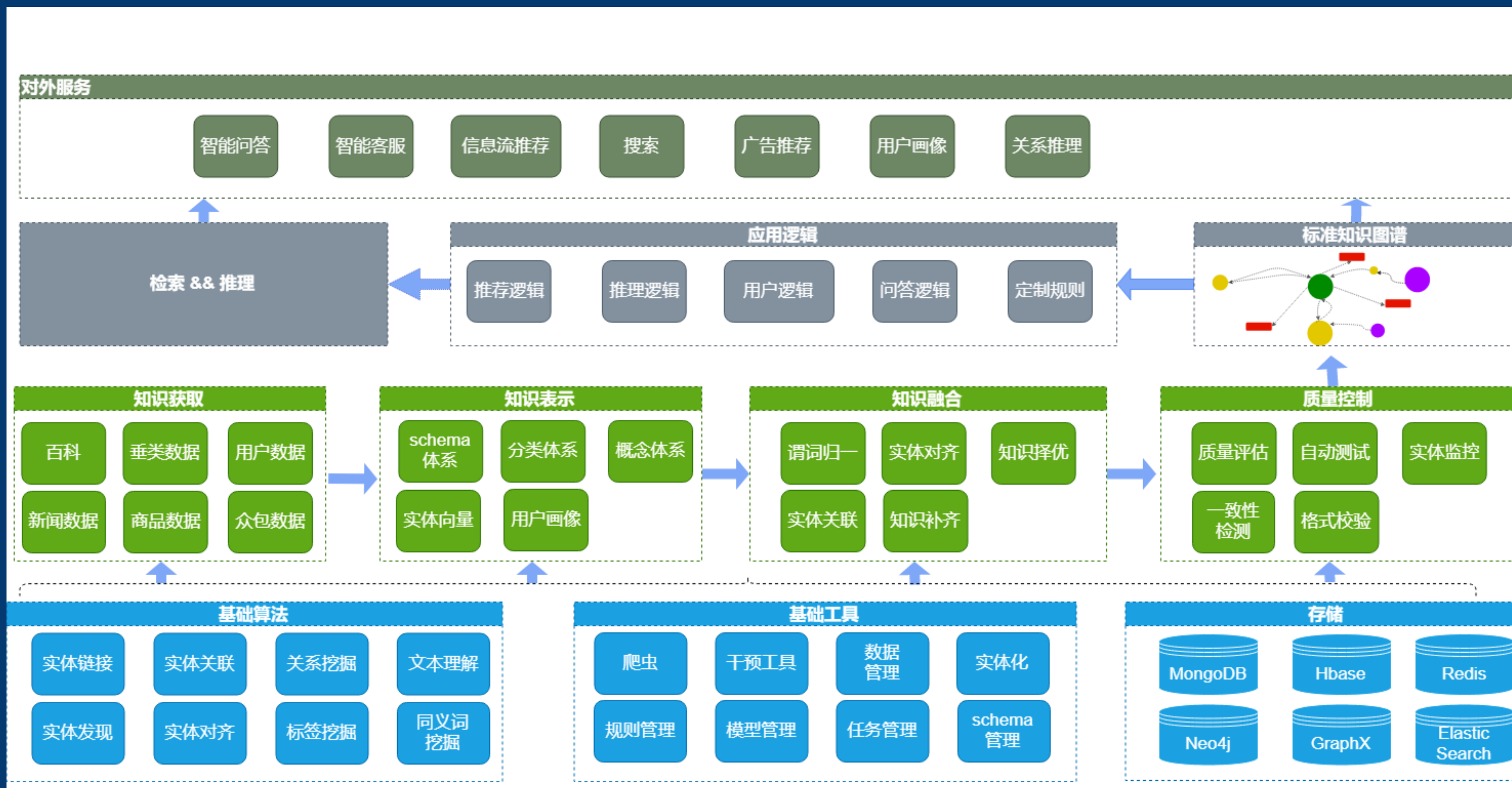
狗尾草

<https://openbase.ai.xiaomi.com/>

TABLE OF CONTENTS 大纲

- 知识图谱概述
- 知识图谱构建
- 基于图谱的问答
- 图谱的其他典型应用

关键构建技术

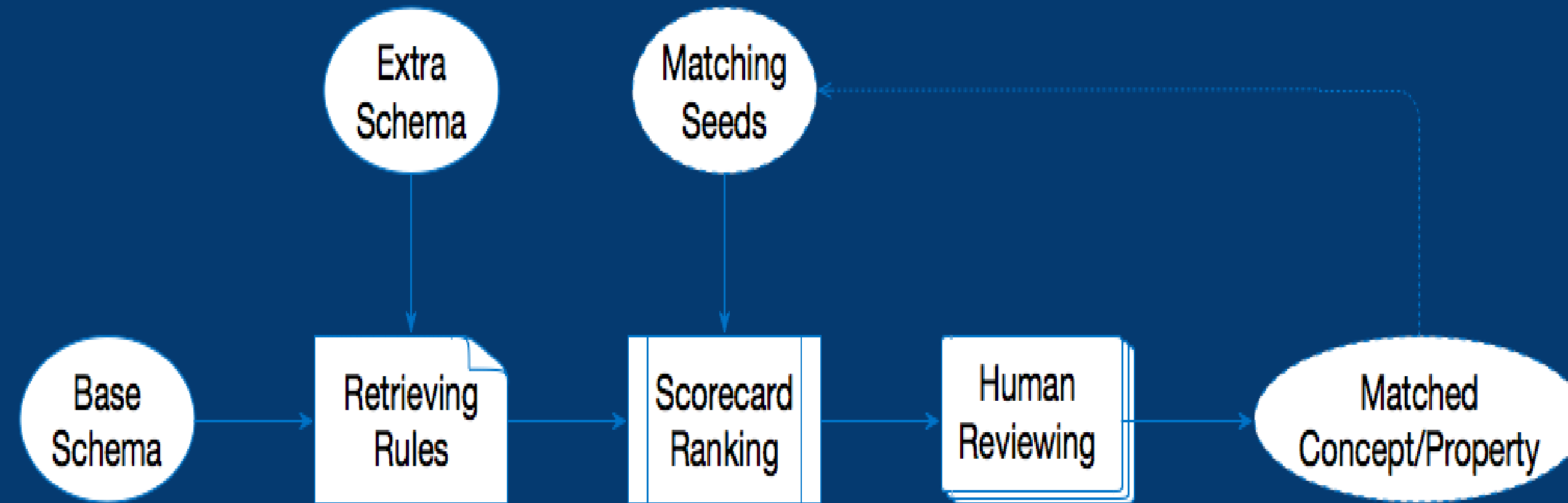


百科知识图谱的构建

知识抽取 → 实体发现 → 实体分类 → 知识补全 → 知识更新 → 知识融合

- 从多个百科数据源开始的、自下而上的图谱构建方式;
- 缺乏完整、一致的本体知识;
- 大规模百科图谱的对齐是一个难题;

百科图谱的对齐：Schema层对齐



- 以cnSchema为基础，半自动地抽取等价概念和上下位关系；
- 利用等价概念发现等价谓词，必要时利用人工干预提升结果质量；
- 等价谓词又可以帮助发现等价概念，然后迭代循环；

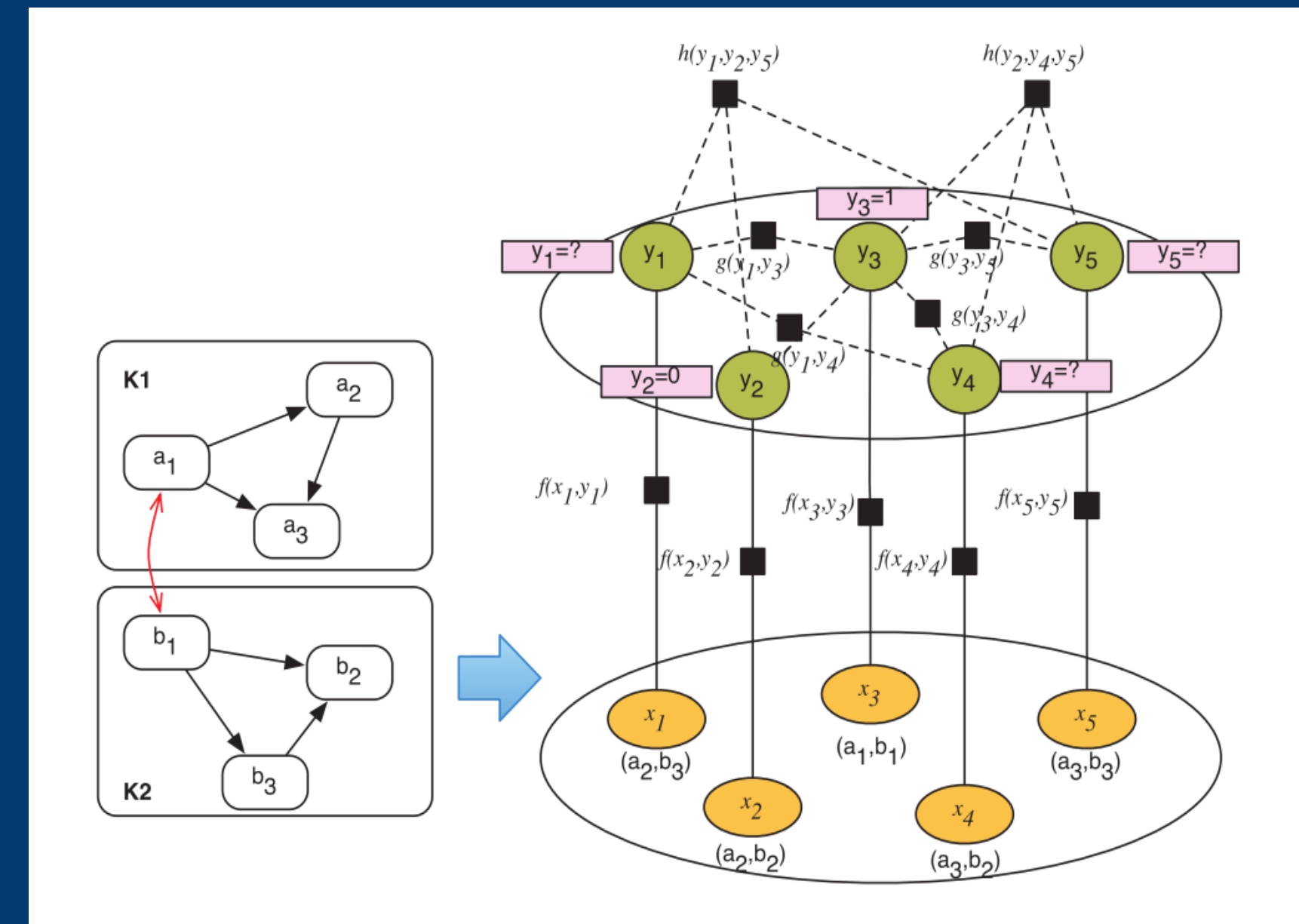
百科图谱的对齐：实例层对齐

- 方法一：使用基于EM算法的半监督学习策略寻找等价实例：
 - 等价实例和属性做为种子集合，并从种子集合中挖掘匹配规则；
 - 根据匹配规则寻找置信度高的等价实例；
 - 将新的等价实例加入到种子集合中，然后迭代循环；

方法二： Linkage Factor Graph Model在跨语言链接上，有着很好的表现

$$p(Y) = \prod_i f(y_i, x_i) g(y_i, G(y_i)) h(y_i, H(y_i))$$

《Cross-lingual Knowledge Linking Across Wiki Knowledge Bases》
(Zhichun Wang... 2012)

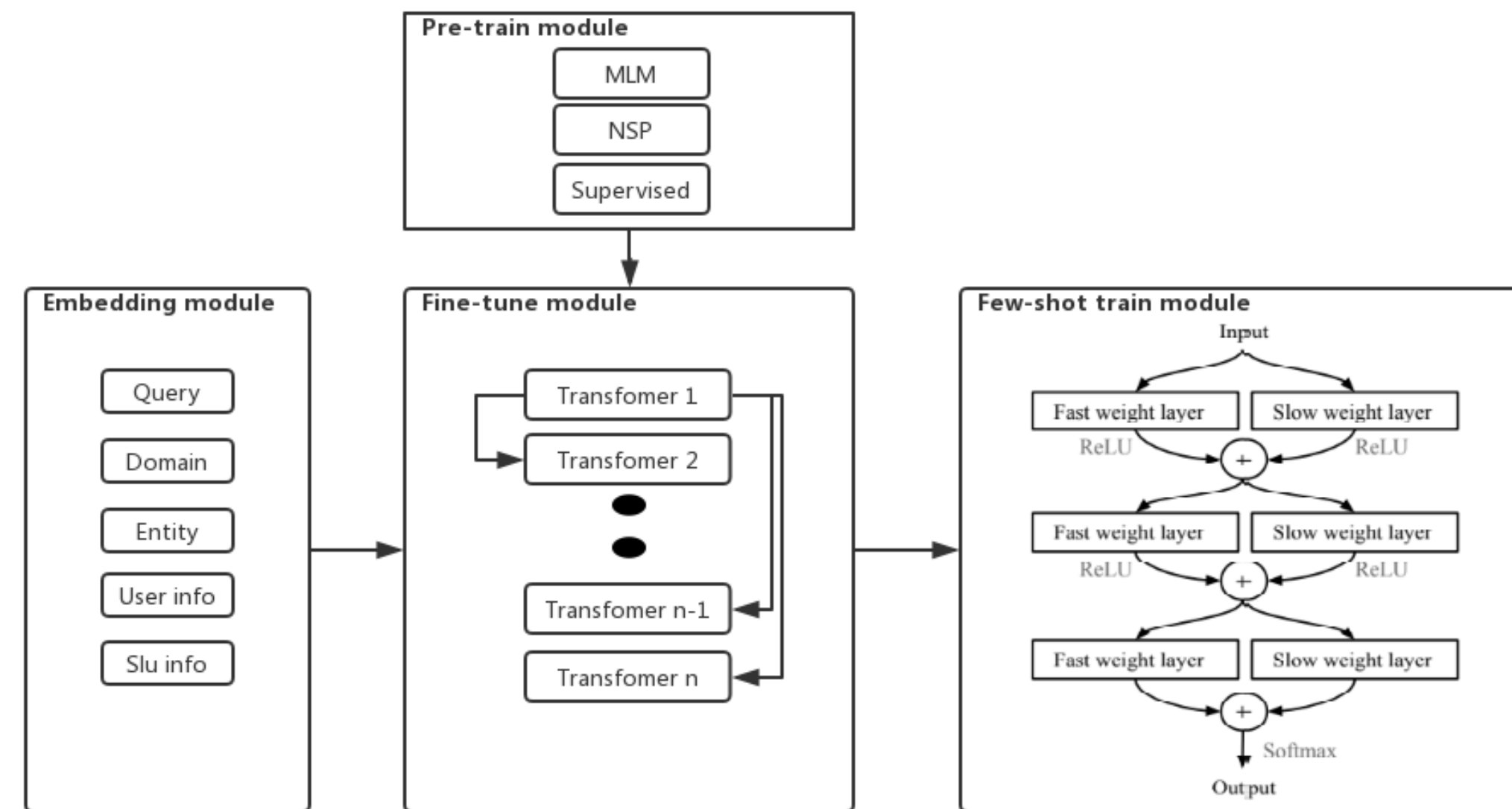


垂域知识图谱的构建

本体定义 → 实体发现 → 知识抽取 → 实体对齐 → 知识择优 → 知识更新

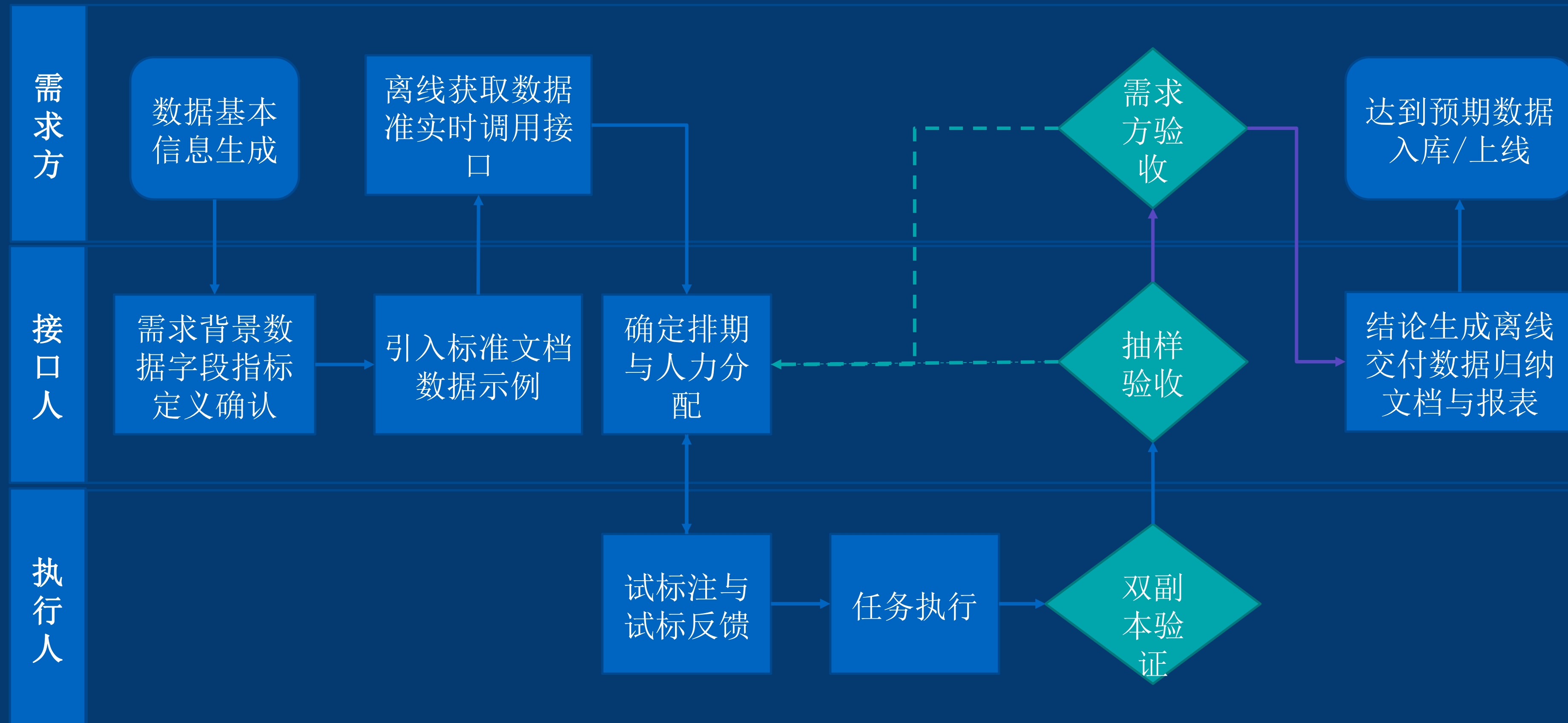
- 垂域图谱的构建从本体定义开始，采取自上而下的方式；
- 数据源的质量不同，知识择优异异常重要；
- 在构建全新的垂域图谱时，基于小样本的模型学习非常重要；

小样本条件下的文本分类模型

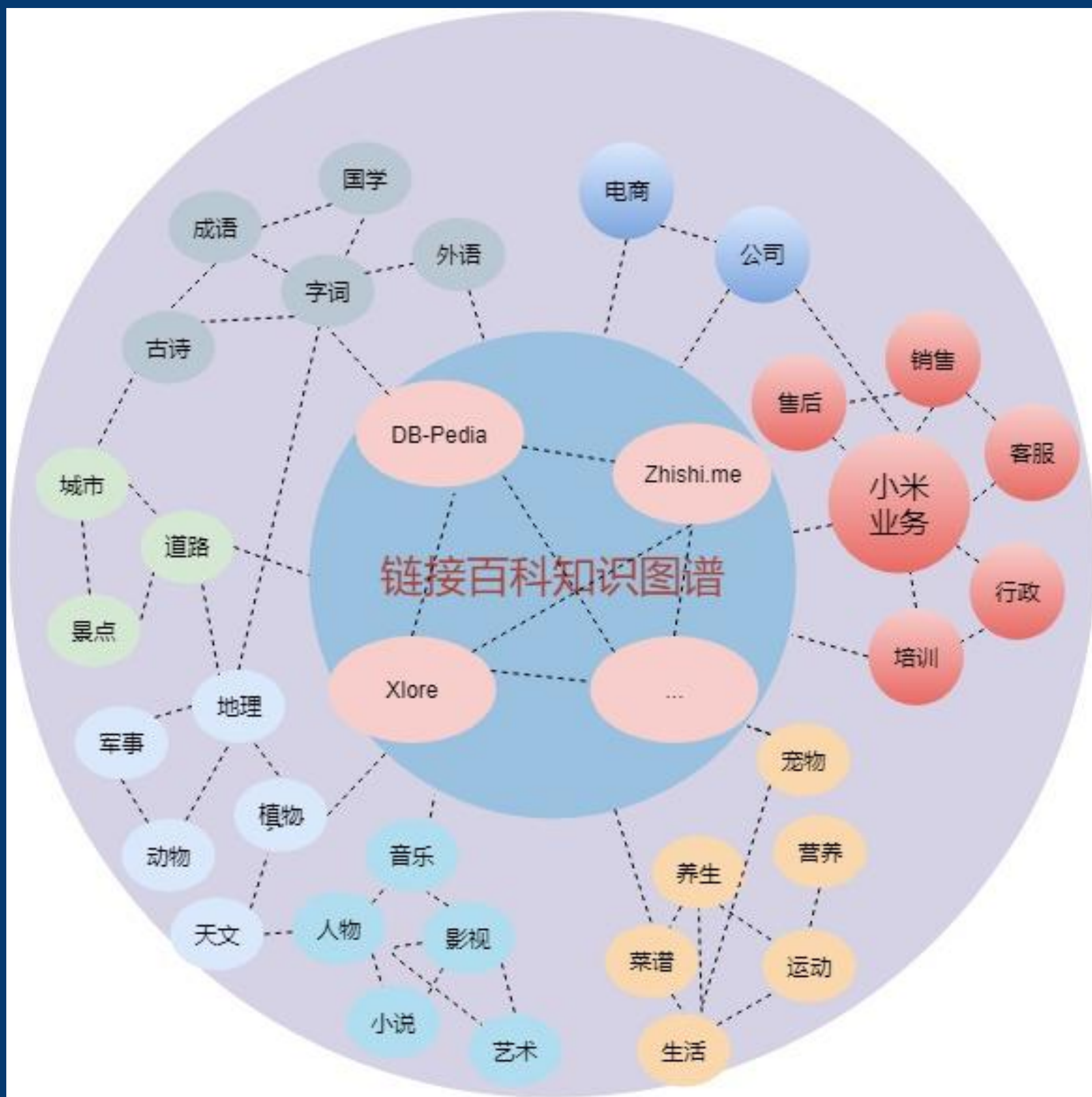


- 联合BERT和Meta Network的网络模型;
- Transformer层之间用残差网络连接, 提升了模型训练速度和分类效果

图谱的质检流程

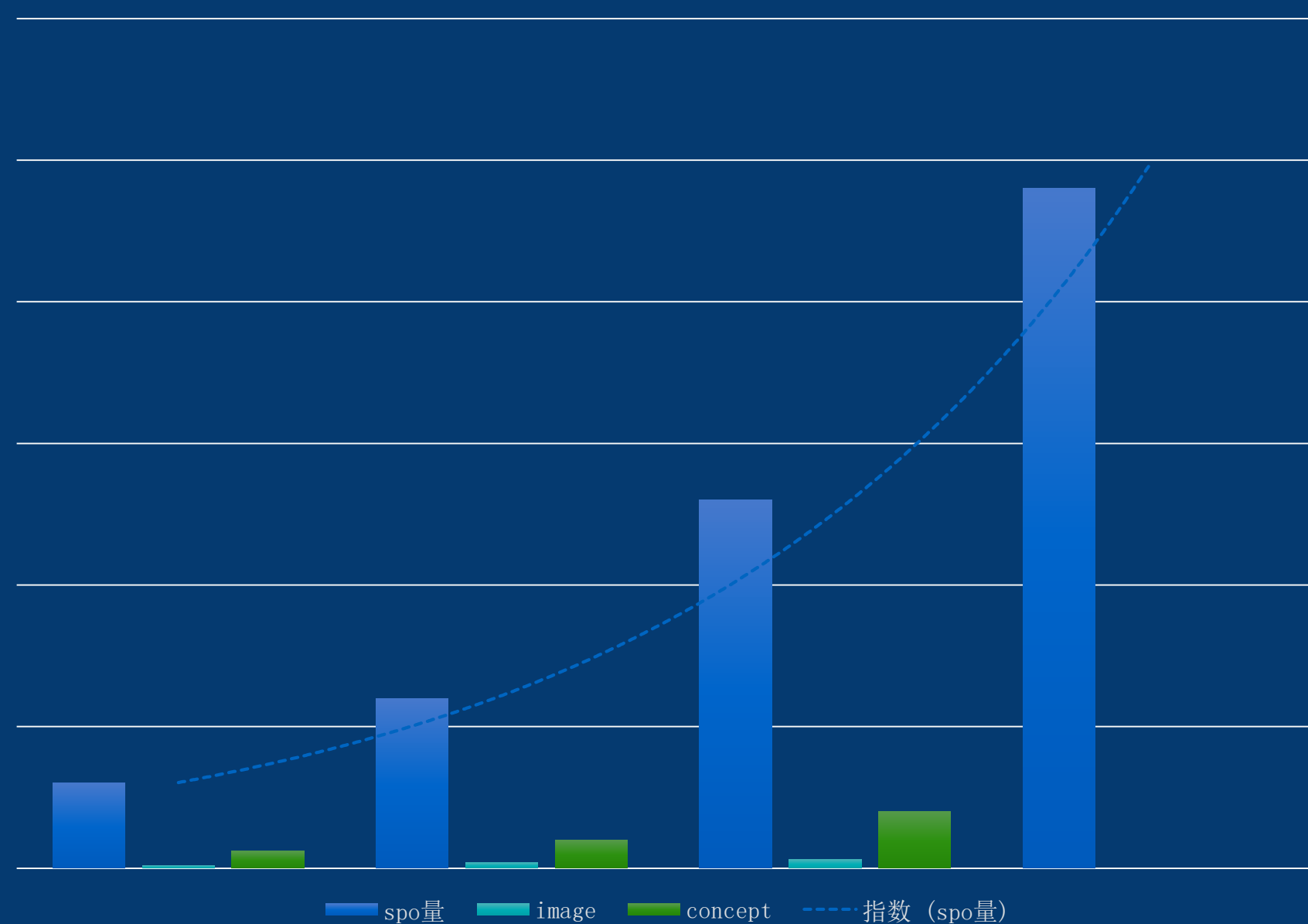


小米图谱



- 以百科图谱为中心，链接了各垂类图谱；
- 在小米智能客服和行政助手场景下，小米业务图谱显著地提升了问答准确率；

小米图谱规模



- 截止到19年Q3，SP0数超过百亿；
- 高质量关系数目每季度翻倍增长；

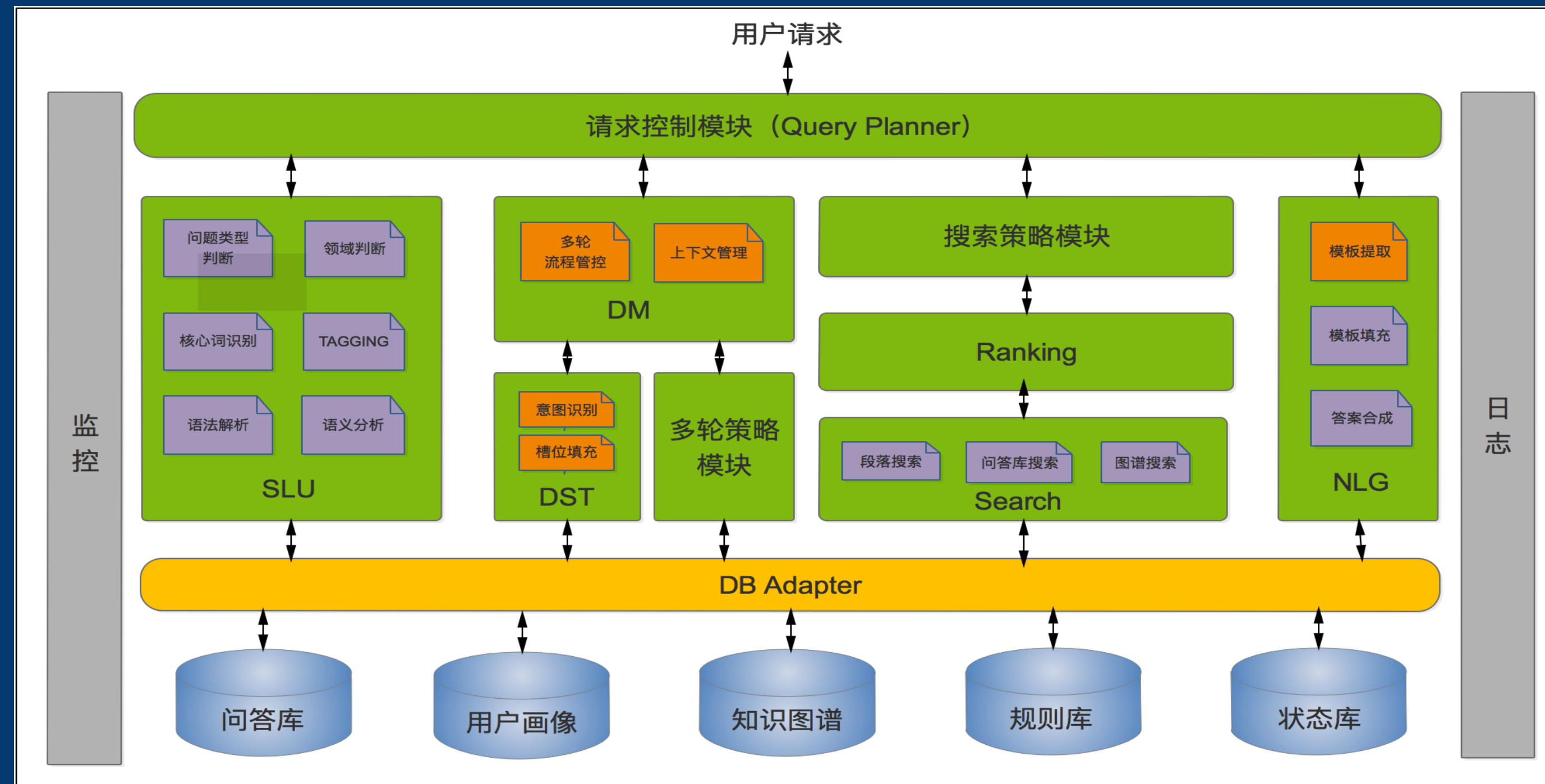
TABLE OF CONTENTS 大纲

- 知识图谱概述
- 知识图谱构建
- 基于图谱的问答
- 图谱的其他典型应用

小爱同学生态架构



小爱开放域问答系统



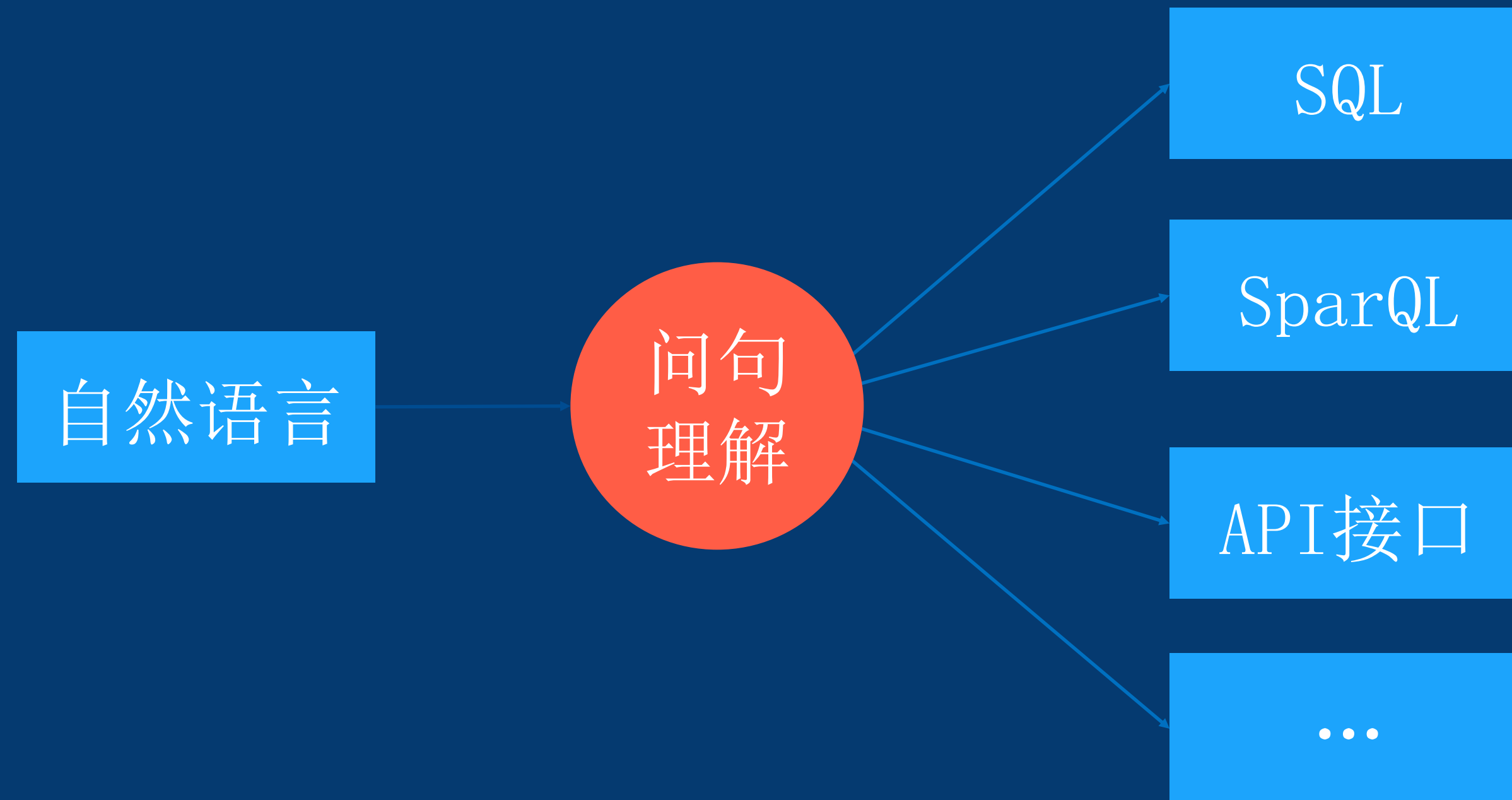
- KBQA是小爱开放域问答系统的重要组成部分；
- 从小爱的线上表现来看，KBQA的准确率，要显著高于基于FAQ的问答和基于阅读理解的问答子系统；

问句理解是当前的难点

问句纠错 → 问句改写 → 意图识别 → 指代消解 → 实体链接 → 谓词判断

- 小爱的KBQA分成许多步，部分借鉴了搜索和NLP的技术
- 需要在80毫秒内完成全部处理，每天处理6千万次以上的请求

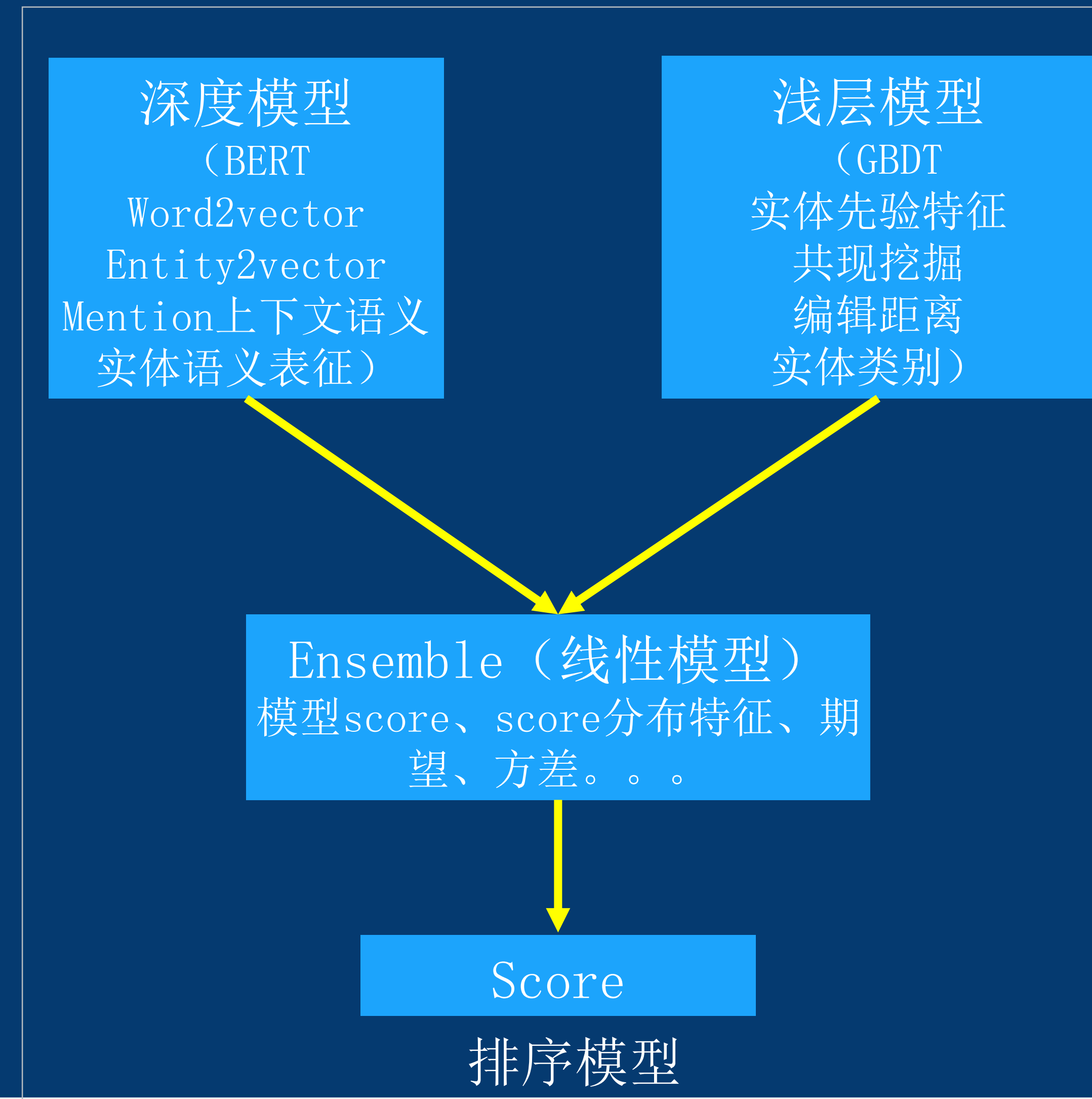
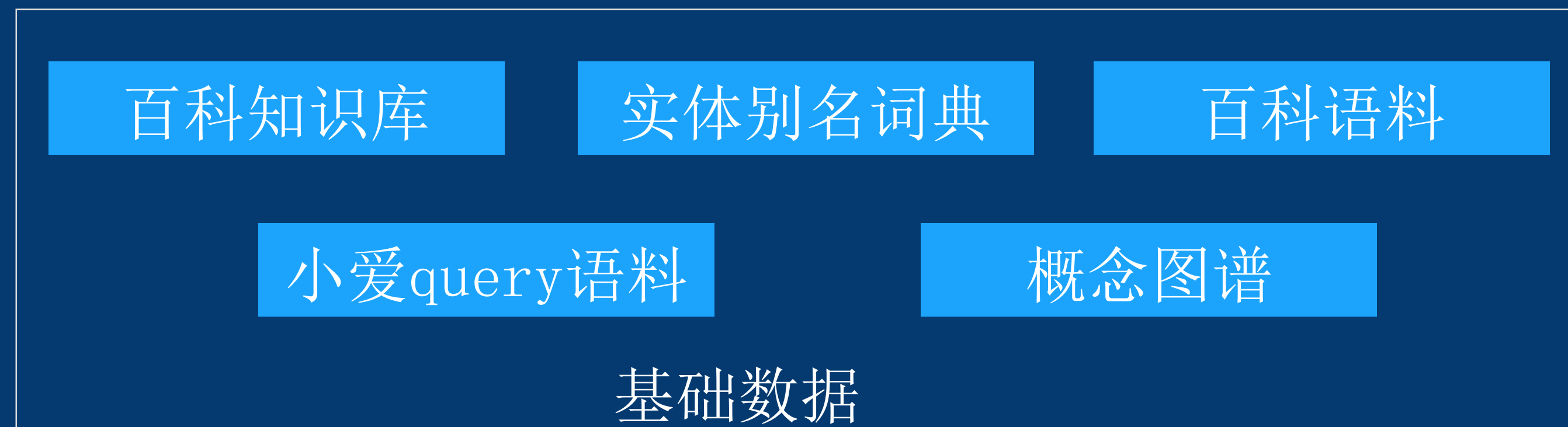
从翻译的角度看用户问句理解



将自然语言翻译成SQL/SparQL等机器语言，方便后续的查询，如：

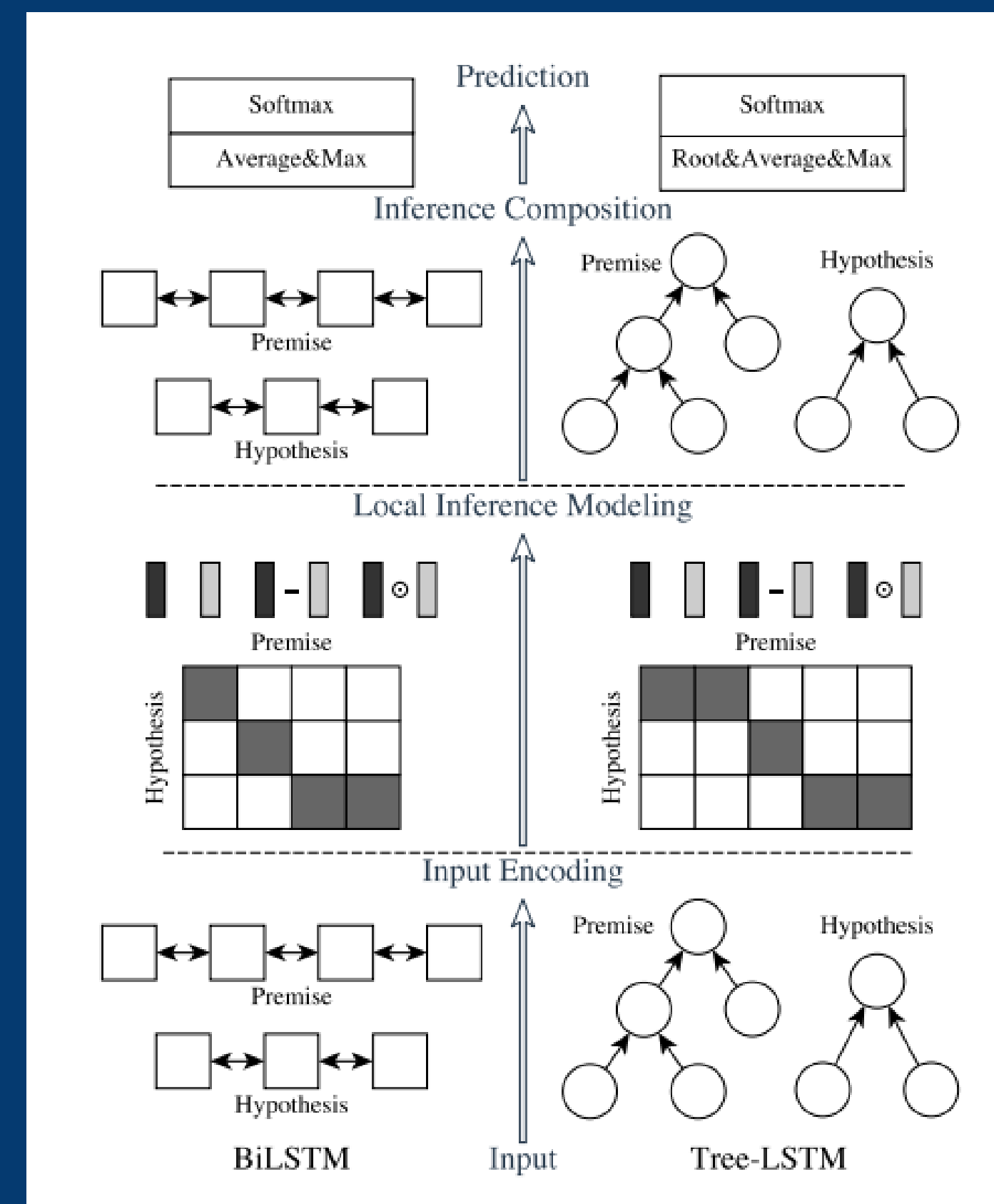
“你知道张三的电话是多少么？” => SELECT Phone_number FROM Employees WHERE Name = “张三”

KBQA关键技术：实体链接



KBQA关键技术：谓词判断

- 头部问题用模板和规则的方式处理，模板的挖掘，可以采用BootStraping方法；
- 结合CFG-CKY Parser来解析成Lambda 算式，
例：“唐家三少的家乡有什么机场” =>
 $\lambda x_y. \text{Birthplace}(\text{唐家三少}, x) \text{ And Airport}(x, y)$;
- 长尾问题采用模型来处理，如利用ESIM模型：
将谓词关联到口语短语，利用短语与用户问句做匹配；

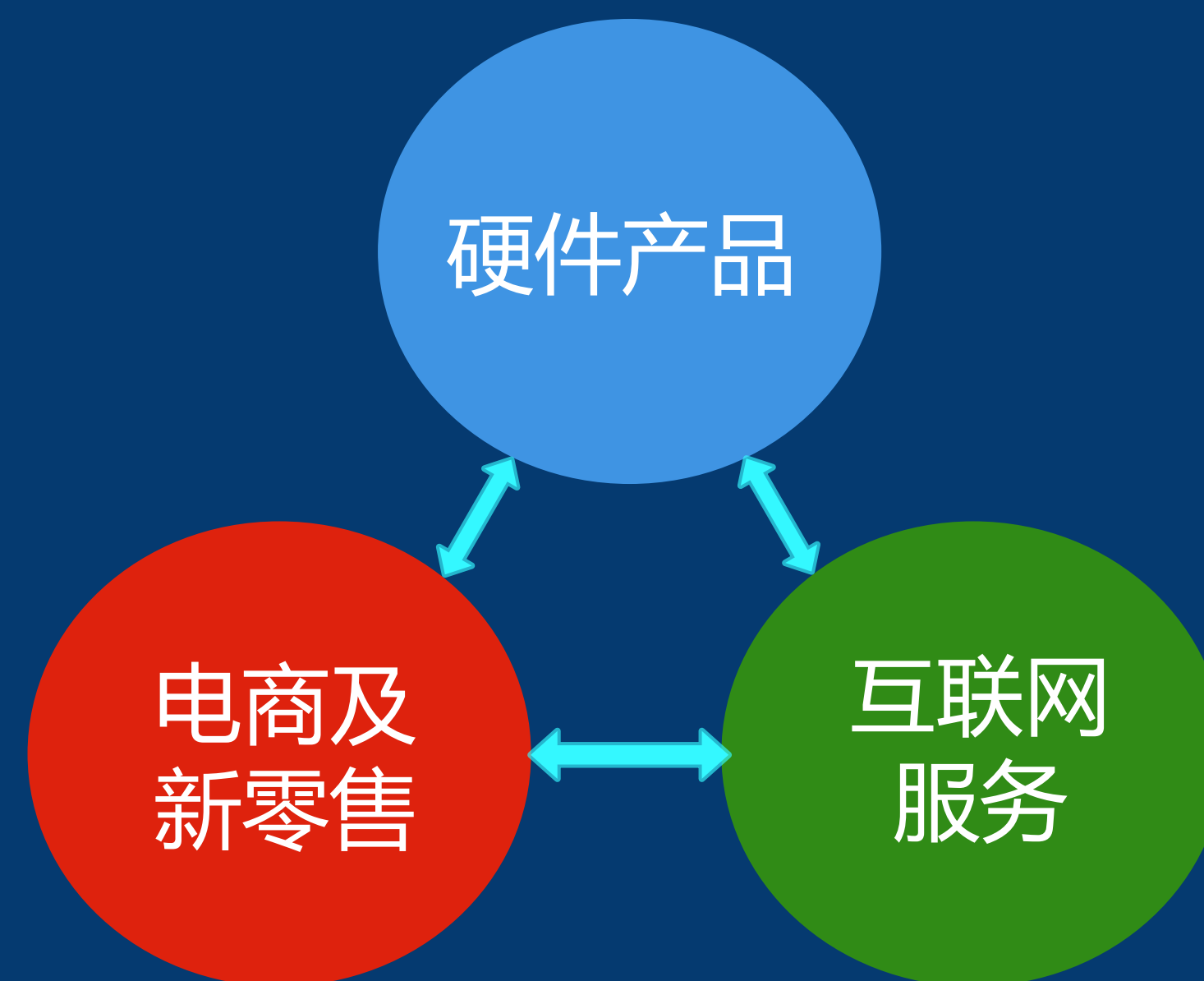


ESIM论文：《Enhanced LSTM for Natural Language Inference》

TABLE OF CONTENTS 大纲

- 知识图谱概述
- 知识图谱构建
- 基于图谱的问答
- 图谱的其他典型应用

众多的应用场景

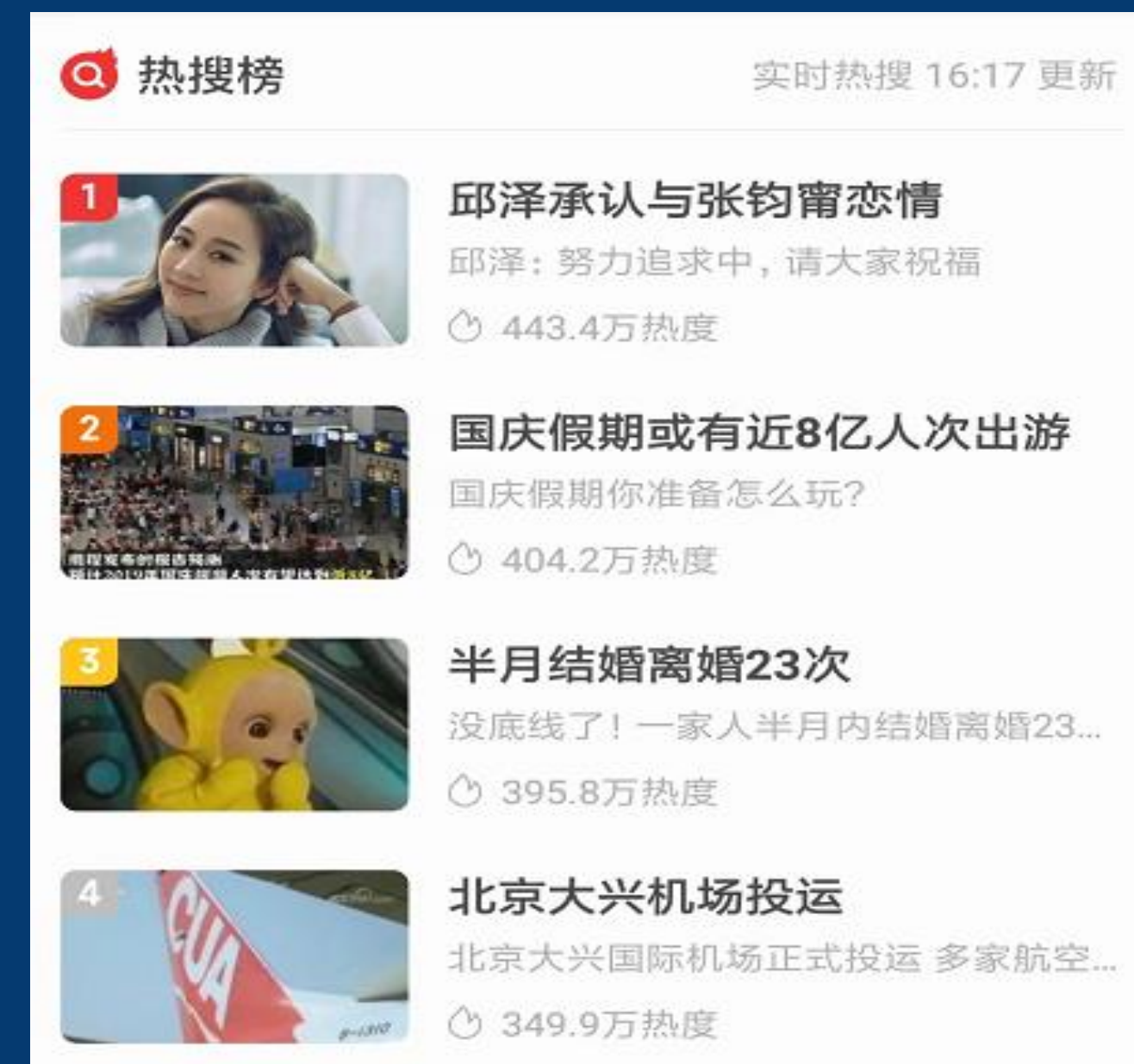


铁人三项

典型场景（一）



有品电商：构建电商图谱，提升推荐精度



新闻推荐：新闻去重

典型场景（二）

事件追踪

07.25

宋慧乔方起诉造谣者

宋慧乔方表示将对传播者采取法律措施

07.22

宋慧乔宋仲基正式离婚

双宋离婚调解成立，正式在法律意义上离婚

07.21

宋慧乔下半年或停工休息

宋慧乔透露自己下半年会休假不接新戏

展开更多

相关人物

宋慧乔

韩国女演员

宋仲基

前夫

宋允儿

好友

金喜善

好友

玉子妍

好友

全局搜索：补充结构化内容

消费购物

消费偏好领域	住宅家具:1、童装/童
偏好品类	传统糕点:1、简易衣
偏好价格区间	0-100:2、500以上:1
消费频率	3.0
近期消费次数	3
近期消费总额	725.0
使用手机类型	
信用等级	

基本属性

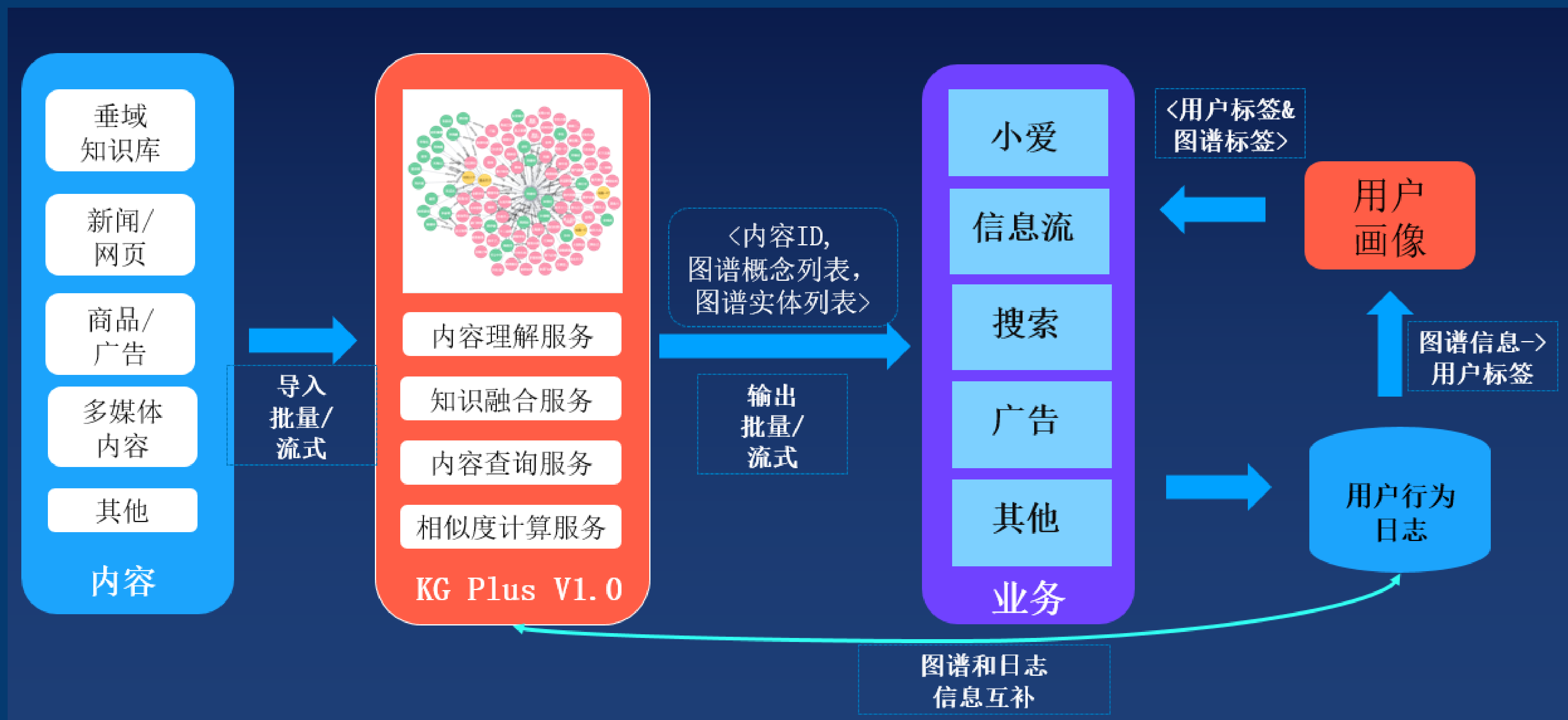
年龄	33
性别	女
生日	1981-6-25
所在国家	01
所在省份	51
所在城市	01
所在县区	00
故乡国家	00
故乡省份	00
故乡城市	00
故乡县区	00
星座	巨蟹座
血型	未知
学校	-中医药大学

交际圈

交际偏好领域	
微博粉丝数	246
微博关注数	270
微博互粉数	85
微博认证类型	无认证
微博认证原因	
微博个人标签	
qq群偏好特征	

用户画像：补充标签数量

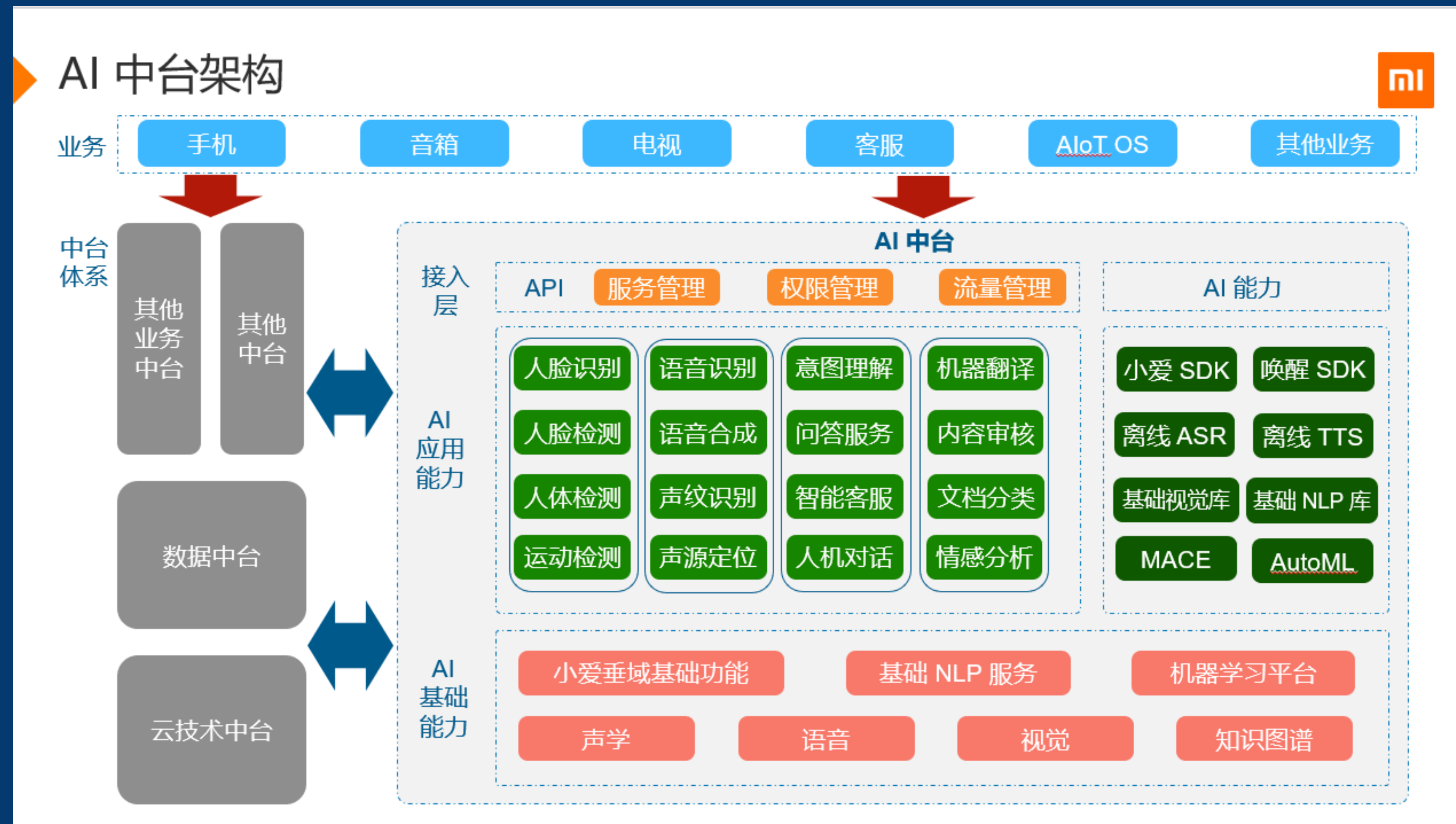
从KG到KG Plus



KG Plus

- 将结构化的图谱数据与非结构化的文本数据和多媒体数据相关联；
- 通过KG Plus平台，统一为各上游产品提供服务；
- KG Plus 平台提供的核心能力如下：
 - 1、KG Plus查询：根据内容名称或者内容ID，从KG Plus里查询内容；
 - 2、理解服务：根据内容名称或者内容ID，查询补充了各类信息的内容数据；
 - 3、融合服务：将结构化的和非结构化的数据融合到KG Plus中；
 - 4、距离计算：查询实体之间或者概念之间的图谱距离 ；
 - 5、实体推荐：根据给定实体，推荐关联实体；

KG Plus在小米AI中台的位置



THANKS

AiCon
全球人工智能与机器学习技术大会