



# 百度基础技术研发

赵世奇 百度

2012-08-16

定位

前瞻

通用

基础

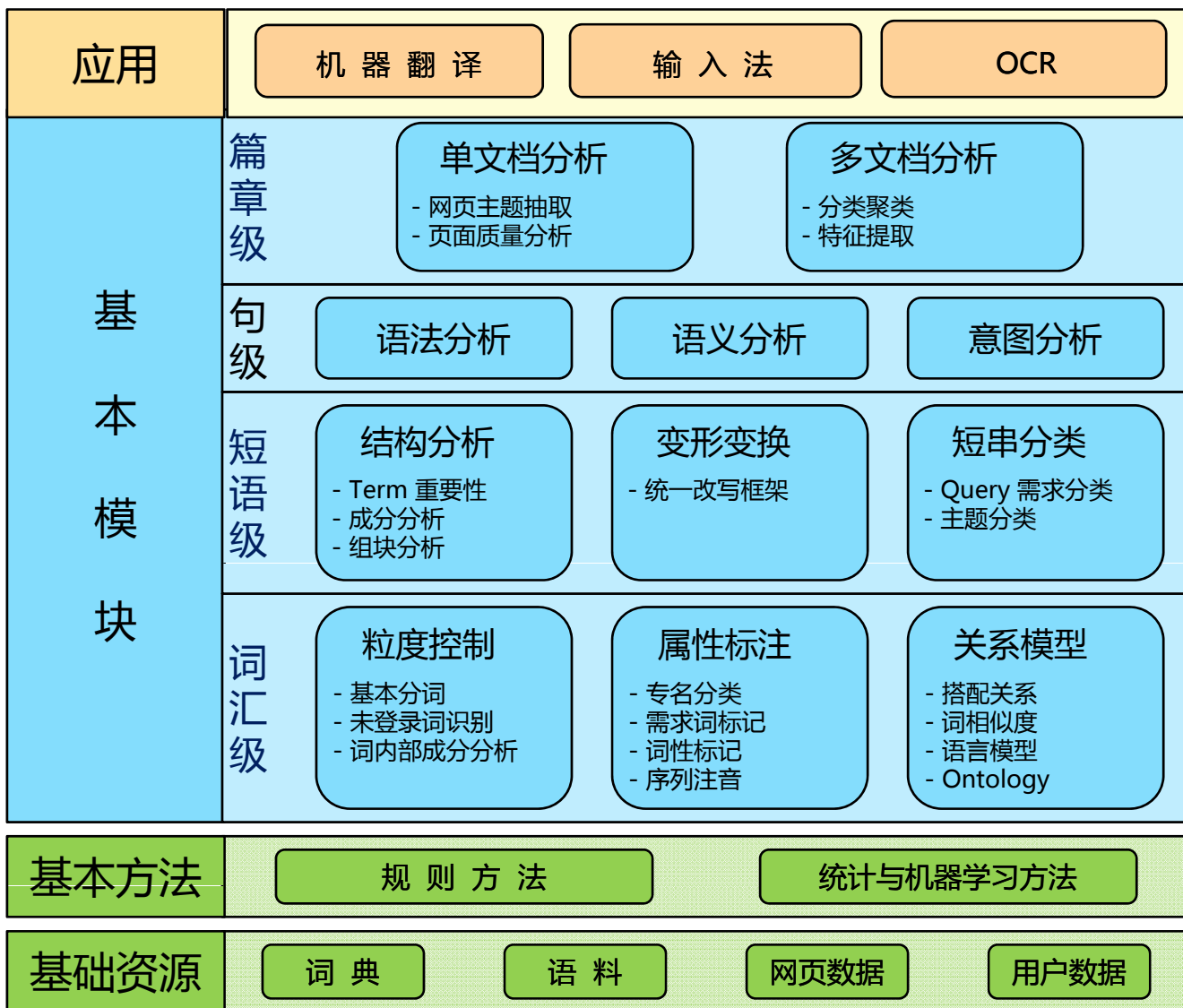
研究内容

声 图 文

# 研发内容

- 自然语言处理
- 互联网数据
  - Spider
  - 互联网数据挖掘
- 多媒体
  - 语音
  - 视觉

# 百度自然语言处理全景图



## 全产品线

基础数据

通用模型

前瞻研发

解决方案

应用订制

知识共享



# NLP业务

- 基础模型与算法
  - 统计机器学习模型
- 基础数据与技术
  - 基础语言资源
  - 用户数据
  - 词法、句法、语义及篇章分析
- 应用技术
  - Query处理技术
  - 用户理解和通用表示空间构建
- 独立产品
  - 机器翻译
  - 输入法
  - OCR



# Ontology建设与应用

# 提纲



Ontology应用价值



Ontology构建 — 新词挖掘



Ontology构建 — 关系挖掘



Ontology应用实例



# Ontology应用意义

- 互联网应用
  - 基础应用：分词与专名识别
    - 新词/专名的挖掘、分类与属性计算
  - 互联网搜索：从传统搜索到语义搜索
    - 丰富的用户引导与知识展现
  - 推荐与个性化：用户建模的底层框架
    - 将用户模型映射到Ontology的知识结构
  - 自动问答：问答系统的知识基础
    - 提供事实性问题的精准答案
  - .....

# Ontology应用意义

- 应用中体现的新问题
  - 新词层出不穷
    - 各类专名（电影、网络红人等）、新词等不断涌现
  - 知识动态变化
    - 如“价格”等类别的属性、“相关实体”等类别的关系
  - 涉及领域繁多
    - 应用的方方面面：母婴、医疗、旅游、美食、娱乐……
  - 数据规模巨大
    - 亿级别数量级
  - 挖掘方法多样
    - 往往是将多种方法的挖掘结果融合在一起

# 提纲



Ontology应用价值



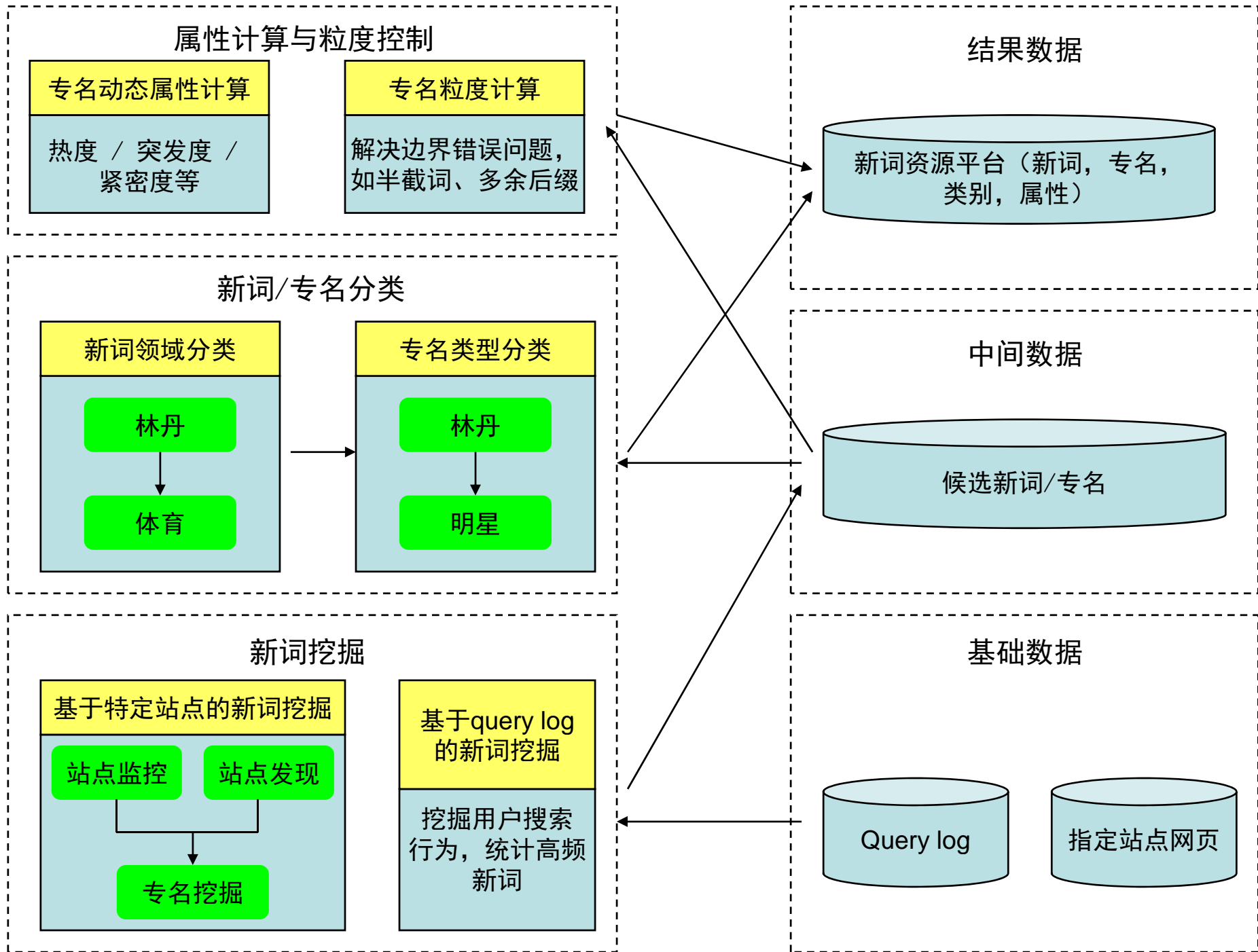
Ontology构建 — 新词挖掘



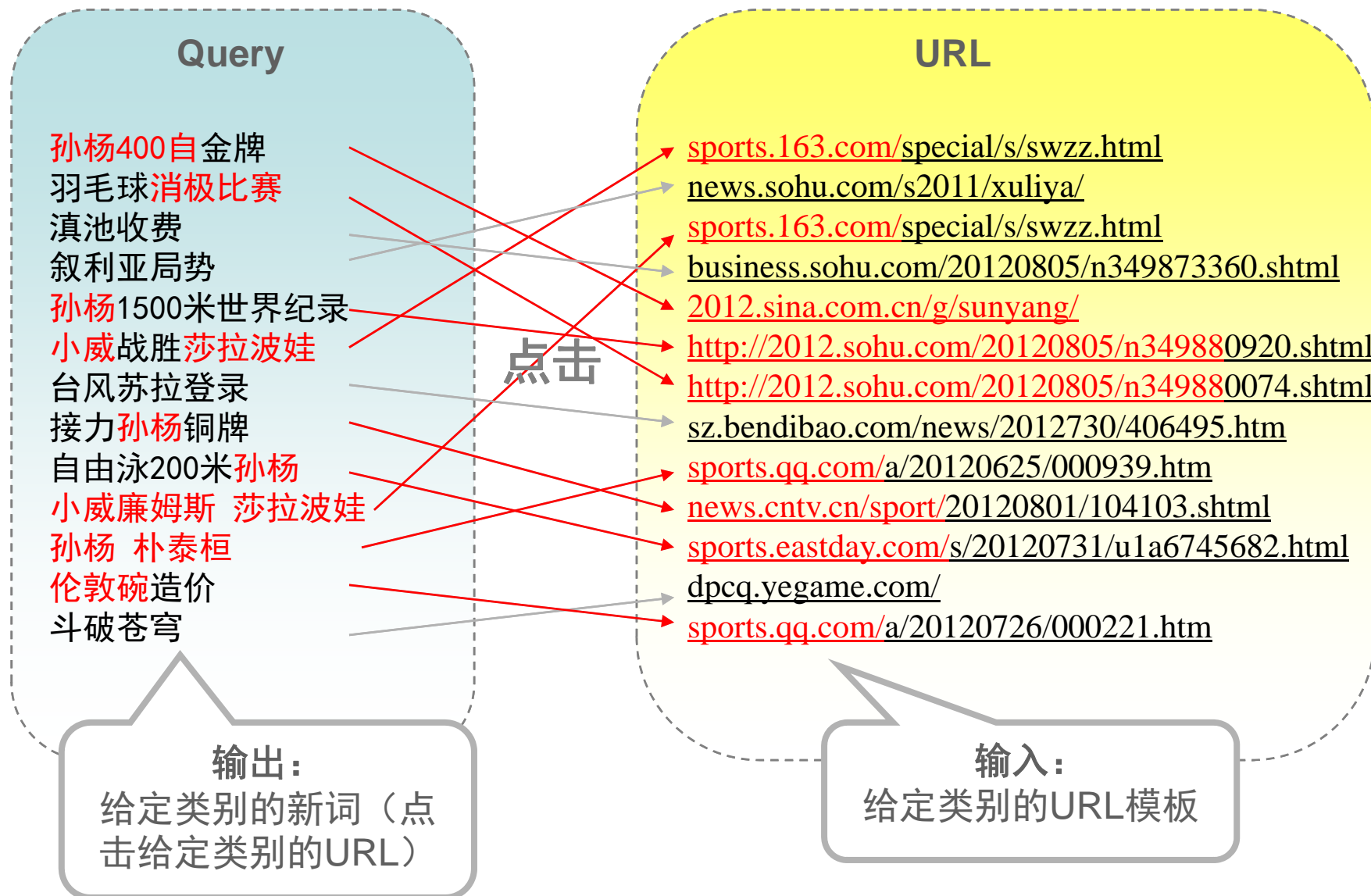
Ontology构建 — 关系挖掘



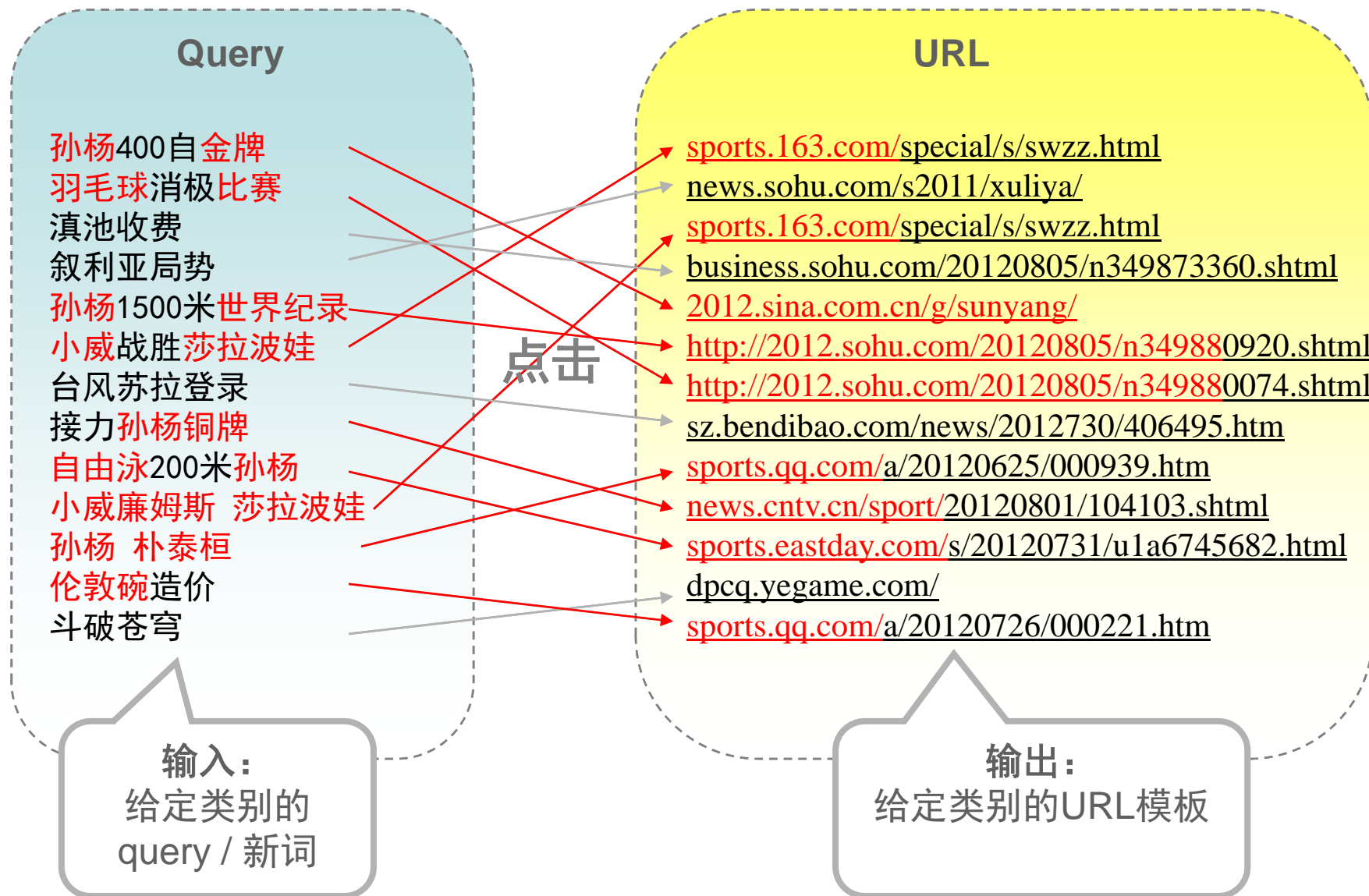
Ontology应用实例



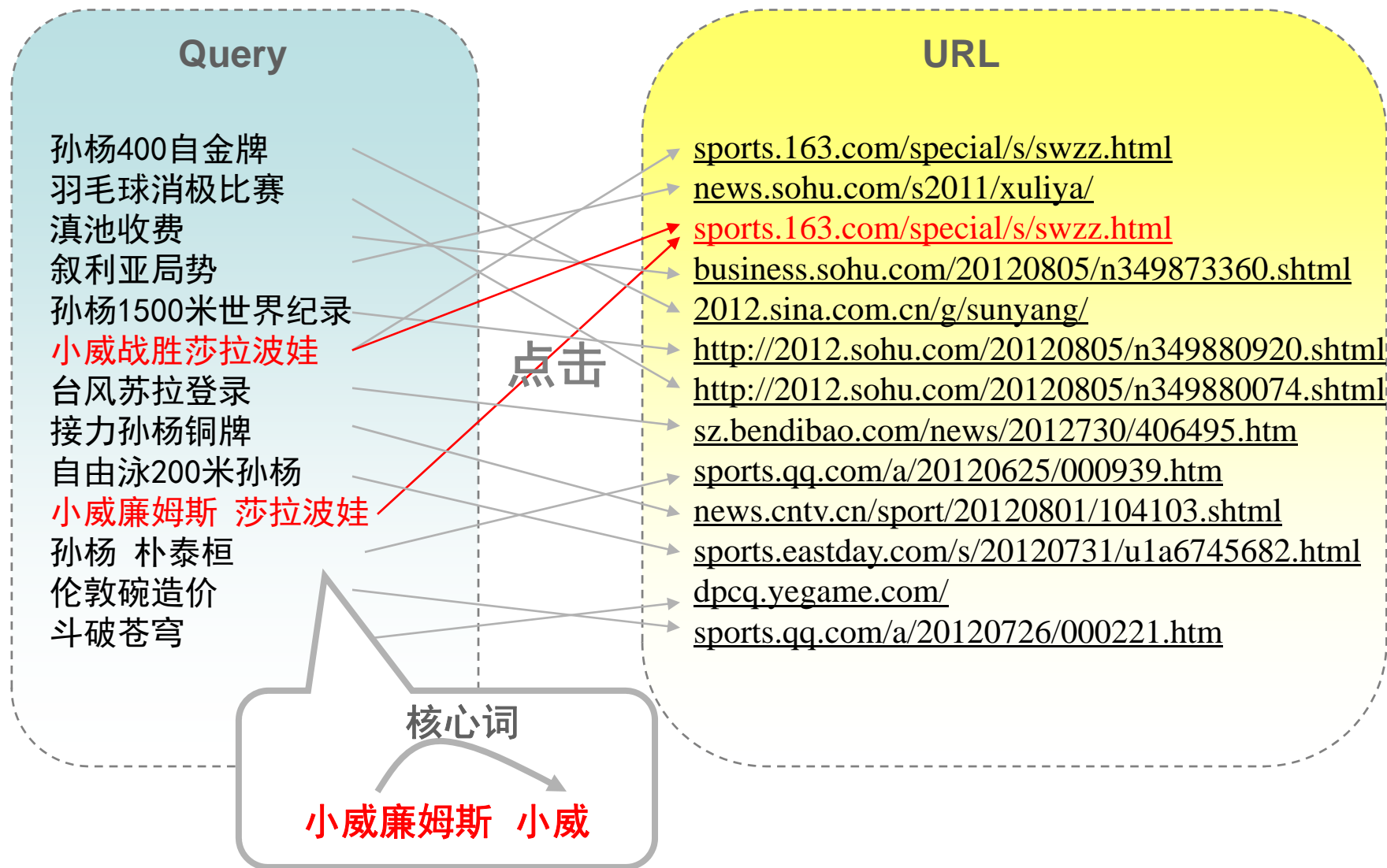
# 基于用户日志的中文新词挖掘 - 新词分类



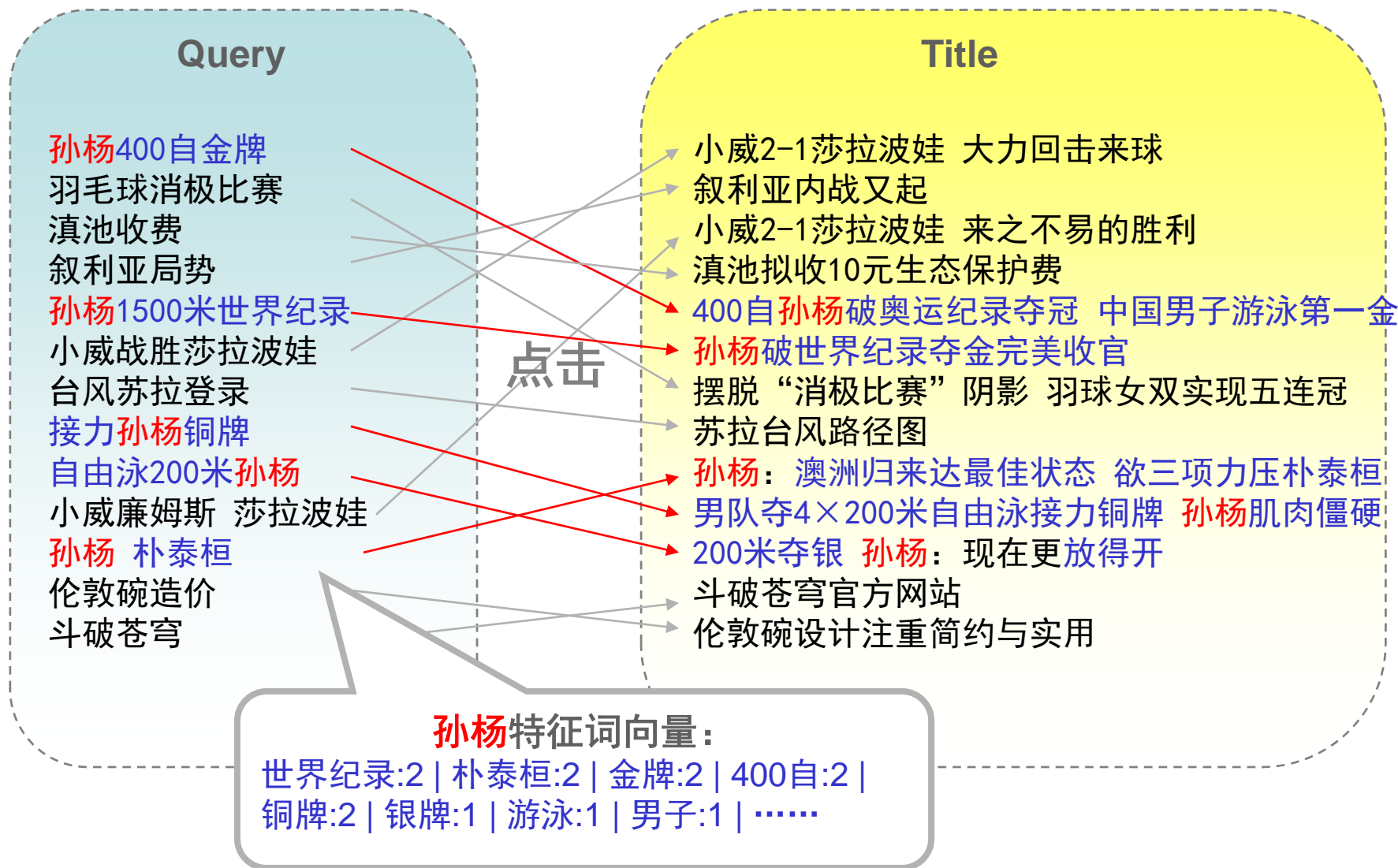
# 基于用户日志的中文新词挖掘 - URL模板挖掘



# 基于用户日志的中文新词挖掘 - 核心词挖掘



# 基于用户日志的中文新词挖掘 - 特征词挖掘





# 提纲



Ontology应用价值



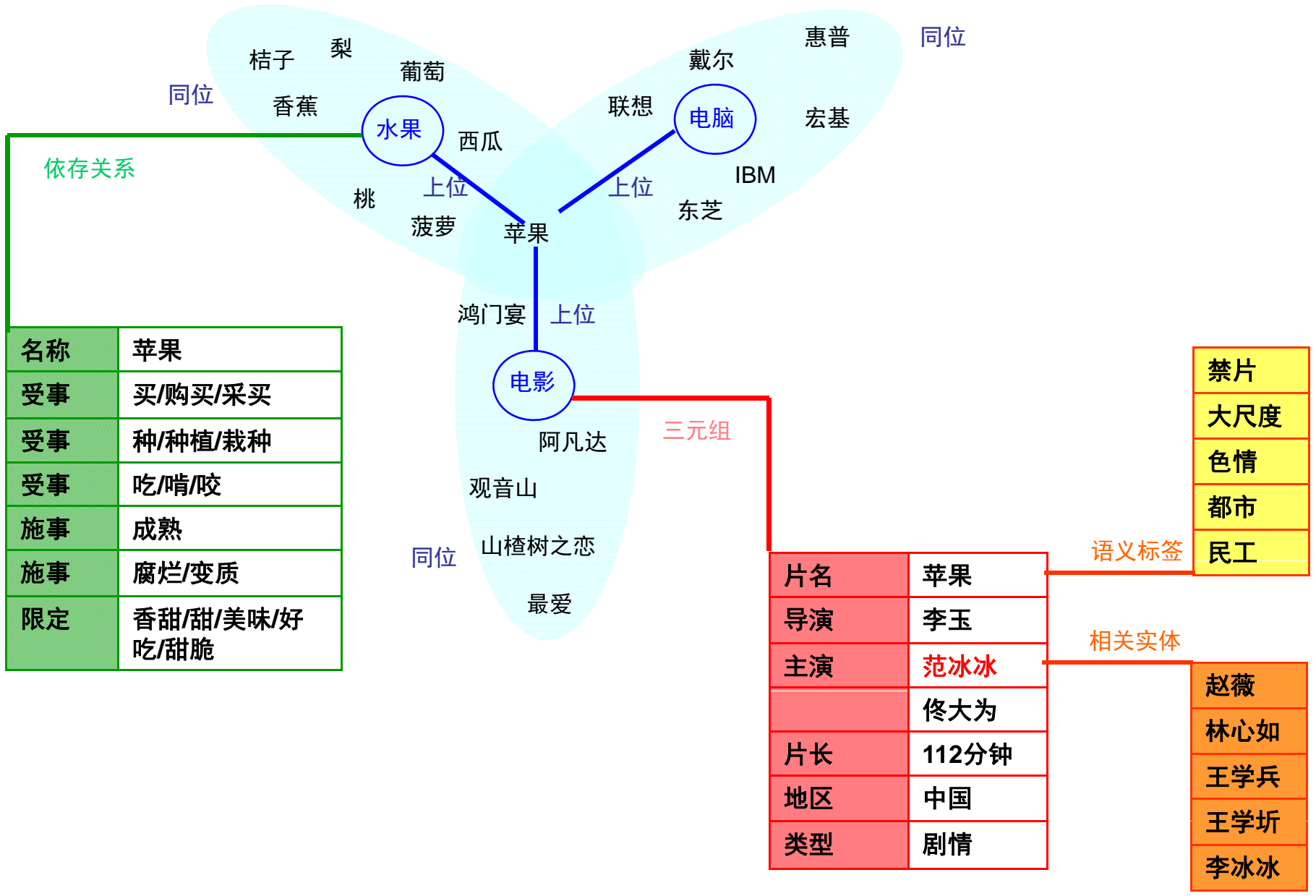
Ontology构建 — 新词挖掘



Ontology构建 — 关系挖掘



Ontology应用实例



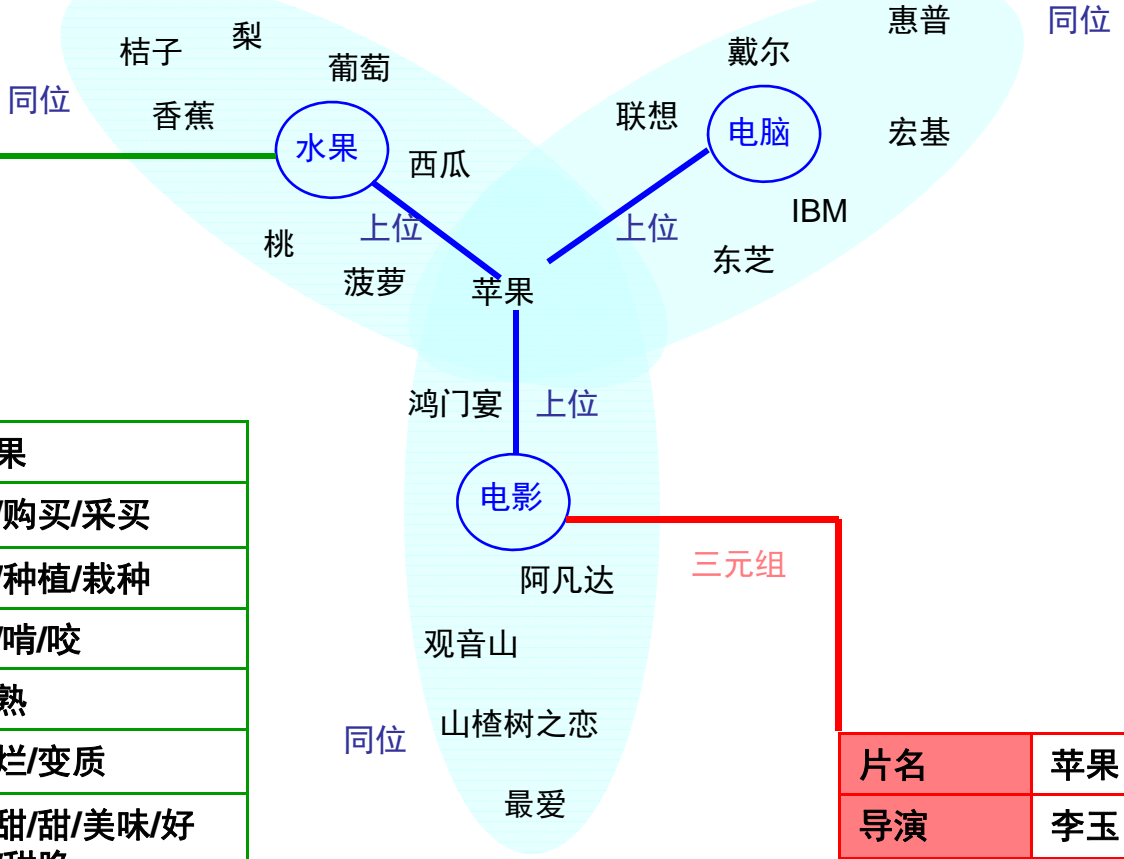
依存关系

名称	苹果
受事	买/购买/采买
受事	种/种植/栽种
受事	吃/啃/咬
施事	成熟
施事	腐烂/变质
限定	香甜/甜/美味/好吃/甜脆

片名	苹果
导演	李玉
主演	范冰冰
	佟大为
片长	112分钟
地区	中国
类型	剧情

禁片
大尺度
色情
都市
民工

赵薇
林心如
王学兵
王学圻
李冰冰



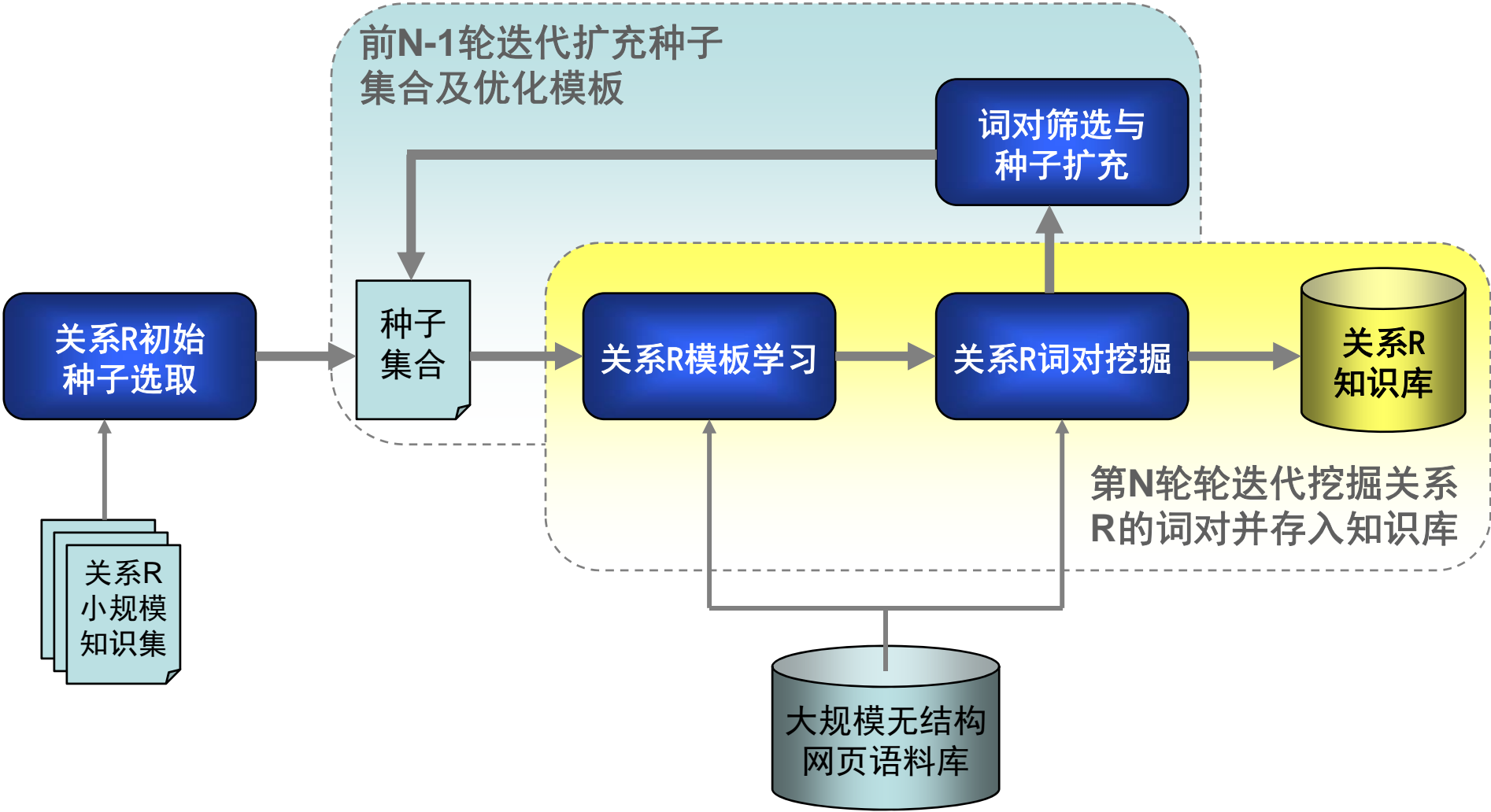
名称	苹果
受事	买/购买/采买
受事	种/种植/栽种
受事	吃/啃/咬
施事	成熟
施事	腐烂/变质
限定	香甜/甜/美味/好吃/甜脆

片名	苹果
导演	李玉
主演	范冰冰
	佟大为
片长	112分钟
地区	中国
类型	剧情

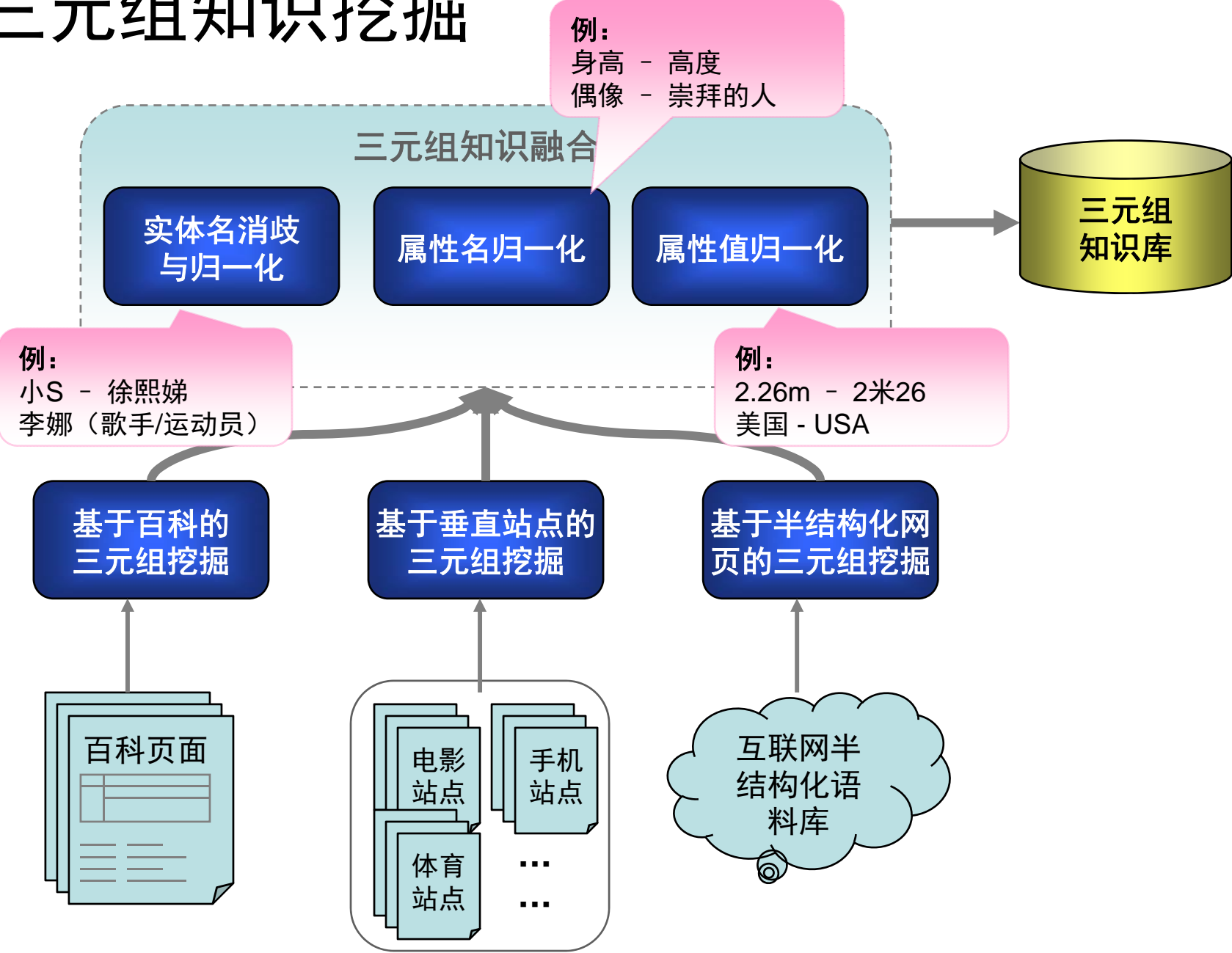
禁片
大尺度
色情
都市
民工

赵薇
林心如
王学兵
王学圻
李冰冰

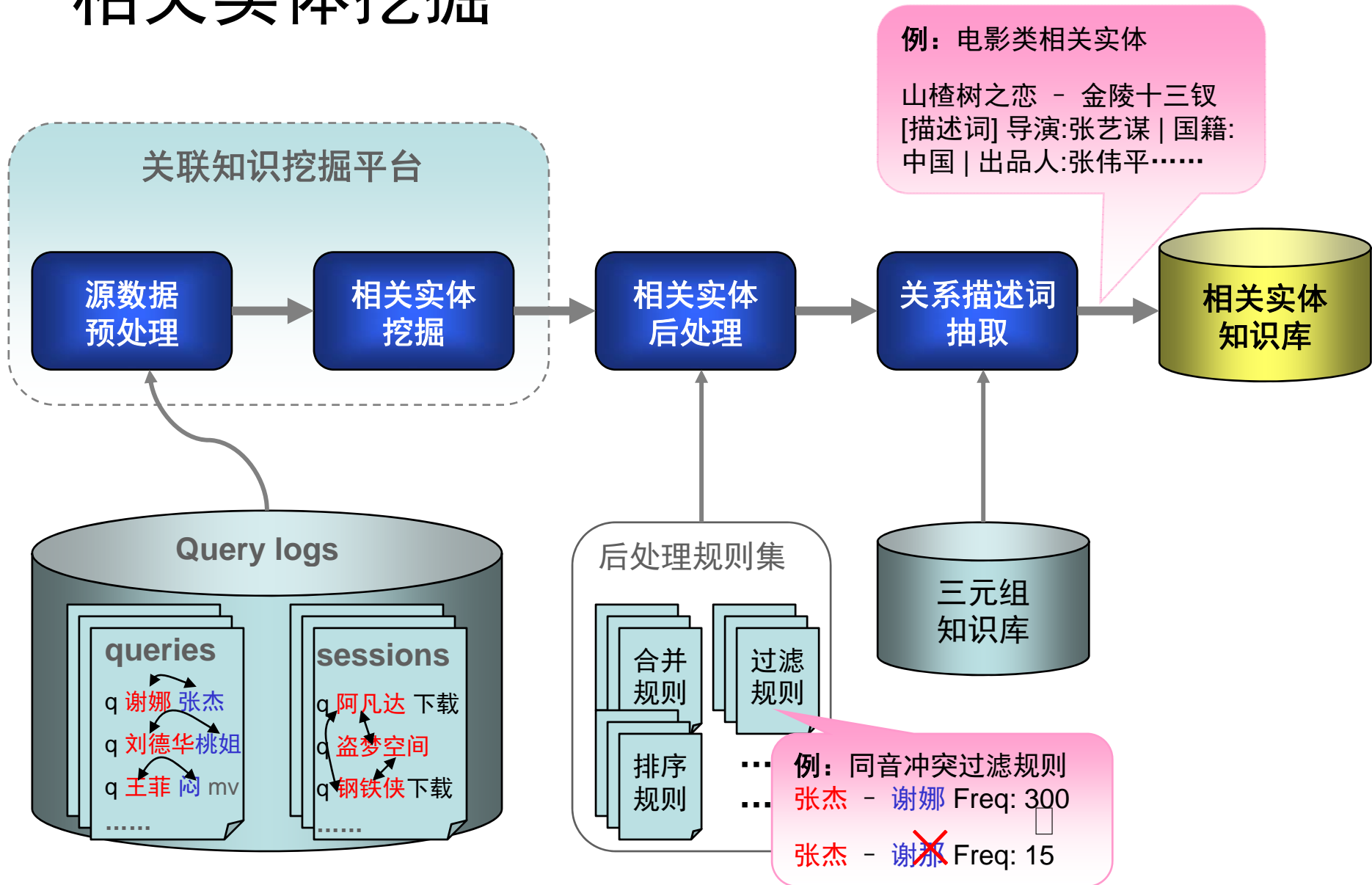
# 上下位/同位知识挖掘



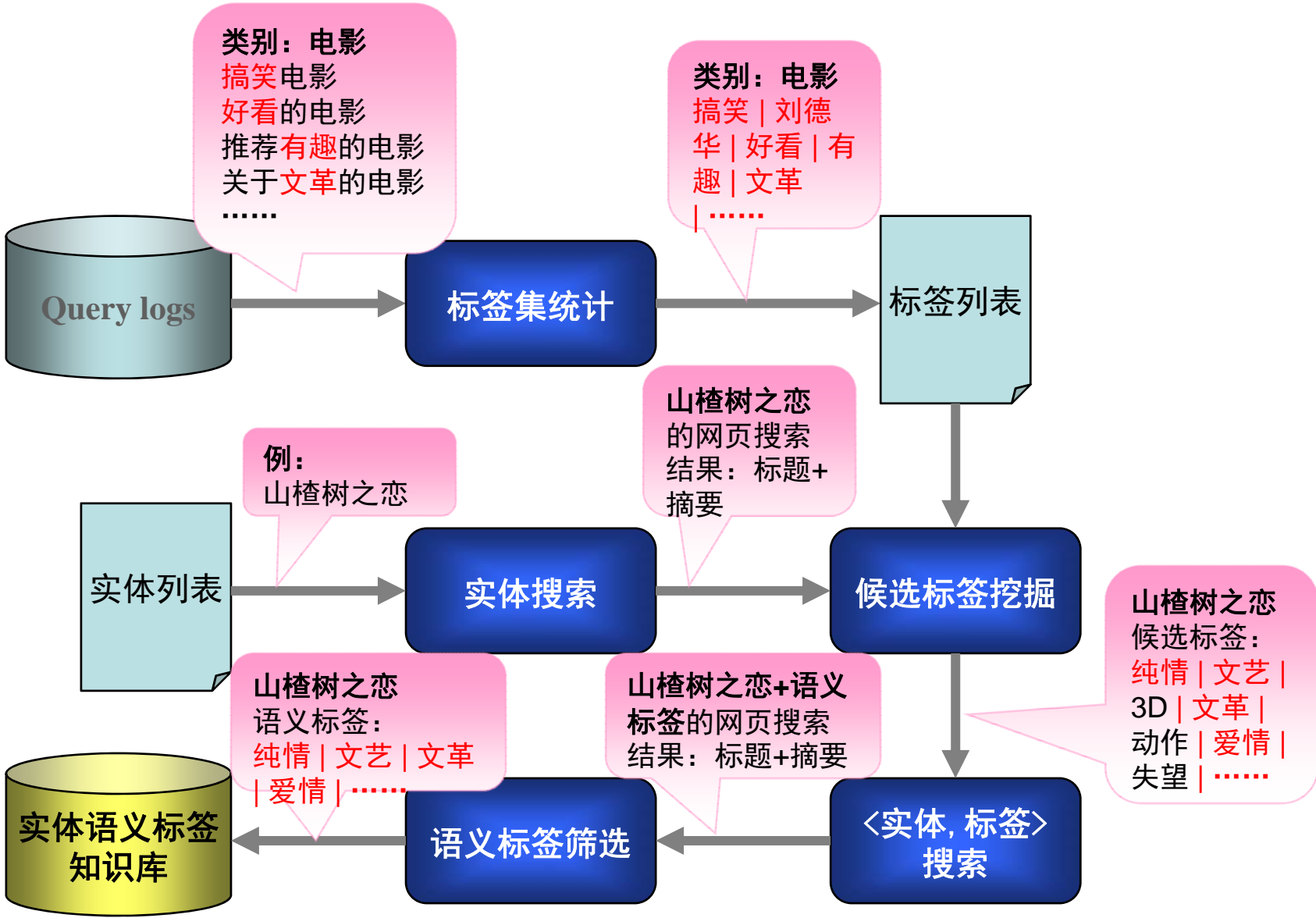
# 三元组知识挖掘



# 相关实体挖掘



# 实体语义标签挖掘



# 提纲



Ontology应用价值



Ontology构建 — 新词挖掘



Ontology构建 — 关系挖掘



Ontology应用实例

# Ontology应用实例

- 网页搜索

- 基于新词挖掘资源改善查询词紧密度

[醇酸磁漆](#) [百度百科](#)  
醇酸磁漆，是由醇酸树脂、颜料、助剂、溶剂等经研磨调配而成的油漆涂料，广泛用作遭受化工大气、  
[基本信息](#)  
baike.baik

[醇酸磁漆](#)  
2个回答  
最佳答案：  
呈磁光色：  
zhidao.ba

[旧家具](#)  
[醇酸磁](#)  
[什么是](#)  
[更多知](#)

[美林君渡](#) [楼盘详情](#) [户型图](#) [业主论坛](#) [北京新浪乐居](#)  
小区均价：均价8600元/平方米...  
位置：燕郊燕顺路西侧...  
开发商：美林地产  
[价格走势](#) [业主论坛](#) [最新图片](#)  
[house.sina.com.cn](#)

[美林君渡](#) [楼盘详情-北京搜房网](#)  
搜房网(SouFun.com)提供美林君渡售楼电话(400-813-0000 转 50448)、最新房价、地址、交通和周边配套、开盘动态、户型图、实景图等楼盘信息。搜房网(SouFun.com...  
[meilinjundu.soufun.com/ 2012-7-13 - 百度快照](#)

[美林·君渡](#) [业主论坛-美林·君渡](#) [业主社区-北京搜狐焦点](#)  
57条回复 - 发帖时间: 2012年8月13日



# Ontology应用实例

- 输入法
  - 基于新词资源更新输入法词库

The diagram illustrates the process of updating an input method dictionary using an ontology. It shows a comparison between a traditional input method's suggestions for 'zhuo'si'tian'cheng' and an updated version for 'zhuo si tian cheng'.

**Traditional Input Method:** Input: zhuo'si'tian'cheng | Suggestions: 1. 捉死天成 2. 卓思天成 3. 桌 4. 卓 5. 捉

**Updated Input Method:** Input: zhuo si tian cheng | Suggestions: 1. 卓思甜橙 2. 卓 3. 桌 4. 着 5. 捉

The comparison is indicated by a starburst labeled "对比" (Comparison) and a downward arrow pointing from the traditional suggestions to the updated ones.

**Search Results for "北京卓思天成国际市场研究咨询有限公司":**

[北京卓思天成国际市场研究咨询有限公司](#)  
北京卓思天成国际市场研究咨询有限公司卓思是一家定位于为主流汽车厂商提供信息咨询服务的  
的企业，目前主要为各大汽车厂家提供新产品上市策略、品牌研究、竞争分析、用户...  
[www.sunjob.com.cn/guanggao/adView.asp?Adl ... 2012-7-29 - 百度快照](#)

[北京卓思天成国际市场研究咨询有限公司北京卓思天成国际市场研究...](#)  
北京卓思天成国际市场研究咨询有限公司简介 卓思是一家新近成立的中小型信息咨询企业，  
以汽车行业市场研究咨询业务为主。卓思定位于为少数高端客户提供高品质产品和...  
[www.01hr.com/company/c-342334696357.html 2012-8-3 - 百度快照](#)

# Ontology应用实例

- 自动问答
  - 基于三元组知识提供确定性问题的答案
    - <http://zhidao.baidu.com>



# Ontology应用实例

- APP推荐
  - 基于相关实体知识推荐相关的app

The screenshot displays a music application interface. On the left, a playlist titled "段林希合集" (Duan Linxi Collection) is shown, containing 30 songs. The main area features a video player for "音悦Tai" (Yinyue Tai) with the song "正在播放: 段林希 - 我的快乐" (Playing: Duan Linxi - My Happiness) from "女声学院" (Girls' College). On the right, a "相关应用" (Related Applications) section lists four artists: 苏妙玲 (Su Miaoling), 洪辰 (Hong Chen), 黄英 (Huang Ying), and 谭杰希 (Tan Jieshi). Red circles highlight the playlist title and the artist list.

百度应用 我的应用

音悦Tai 正在播放: 段林希 - 我的快乐 女声学院 点击下载

段林希合集 - 音悦  
共有歌曲30首  
火柴天堂 感谢你用心  
点击听歌

相关应用

苏妙玲

洪辰

黄英

谭杰希

谢谢!

Q&A

