

浅析知识图谱于智能搜索领域中的应用



By Peng Zhang
2020-10-22

目录

- 引论
- 搜索的发展、现状及面临的问题
- 智能搜索的概念及发展
- 知识图谱的概念及发展
- 智能搜索与知识图谱



引论



- 在大数据时代，人们需要在最短的时间内，得到最有效的数据信息。
- 知识图谱是当下科技热点，并有众多应用分支；但针对智能搜索的应用方面，其技术目前尚在摸索中前行。
- 本文旨在对最新的智能搜索技术、知识图谱应用，以及两者之间如何有效地融合，进行了分析、探讨与预研。

搜索引擎的发展



搜索引擎的四个时代（张俊林《这就是搜索引擎》）

1. 史前时代：分类目录的一代

这个时代成为“导航时代”，Yahoo和国内hao123是这个时代的代表。

2. 第一代：文本检索的一代

文本检索的一代采用经典的信息检索模型，如布尔模型、向量空间模型或者概率模型，来计算用户查询关键词和网页文本内容的相关程度。

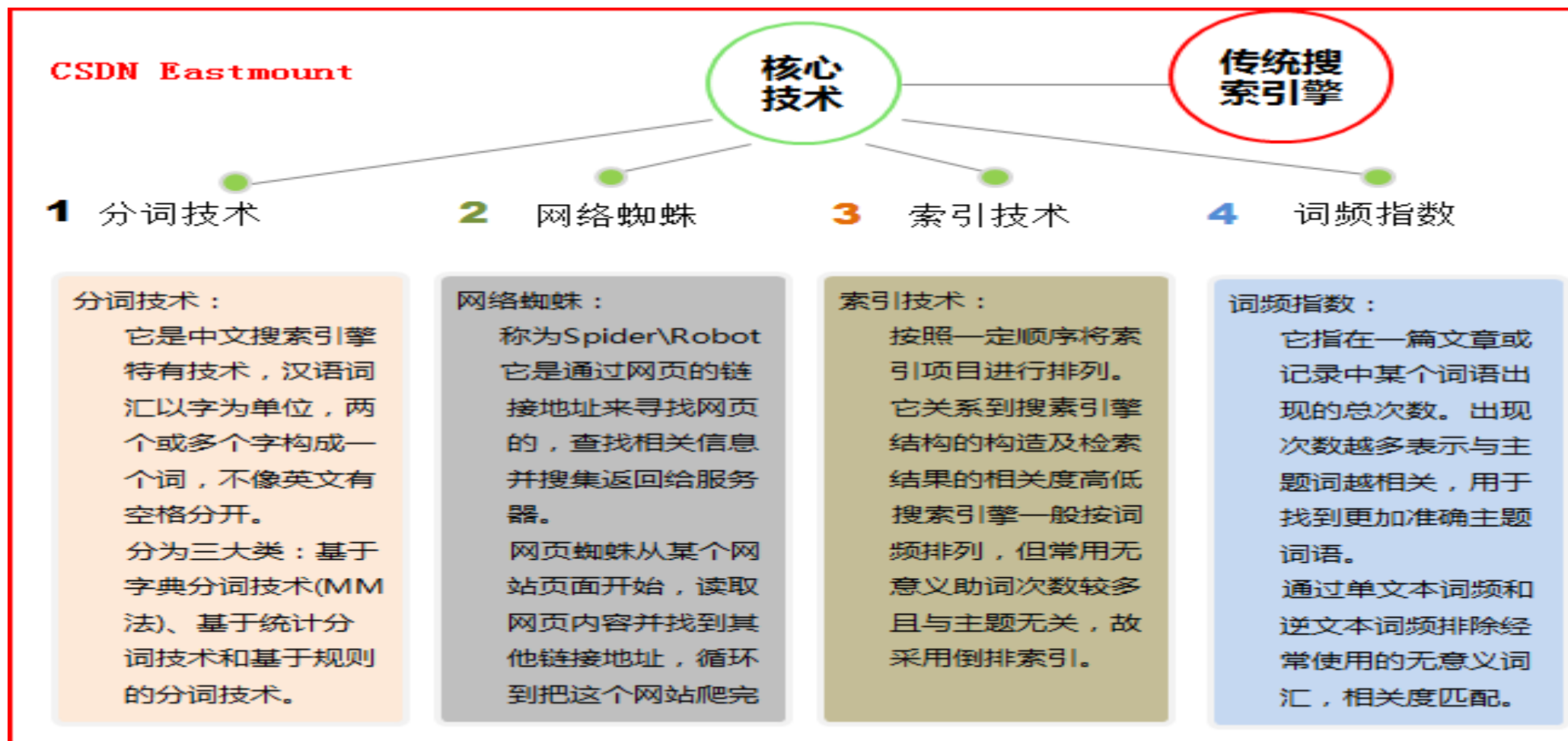
3. 第二代：链接分析的一代

这一代搜索引擎充分利用了网页之间的链接关系，并深入挖掘和利用了网页链接所代表的含义。

4. 第三代：用户中心的一代第三代即理解用户需求为核心的一代搜索引擎。不同用户即使输入同一个查询词，但其目的可能不一样。比如同样输入“苹果”作为搜索词，一个追捧iPhone的时尚青年和一个果农的目的会存在巨大的差异。

传统搜索引擎的核心技术

传统搜索引擎的核心技术常见包括：分词技术、网络蜘蛛、索引技术和词频指数。知识图谱或知识计算引擎被认为是下一代搜索引擎。



传统模式搜索的问题



- 以往传统的搜索引擎，是基于关键词或字符串的，并没有对查询的目标（通常为网页）和用户的查询输入进行理解。
- 传统搜索引擎在搜索准确度方面存在明显的缺陷，即由于HTML形式的网页缺乏语义，难以被计算机理解。
- 传统搜索引擎的搜索模式单一，即缺乏高级详细的搜索条件，大多也缺乏诸如语音、图像、视频识别等方面的搜索手段。
- 传统搜索引擎搜索的结果，需要搜索者在浩若烟海的结果页里逐页打开并分析鉴别最终提取有用的部分，效率低下。

智能搜索的概念及发展--总论

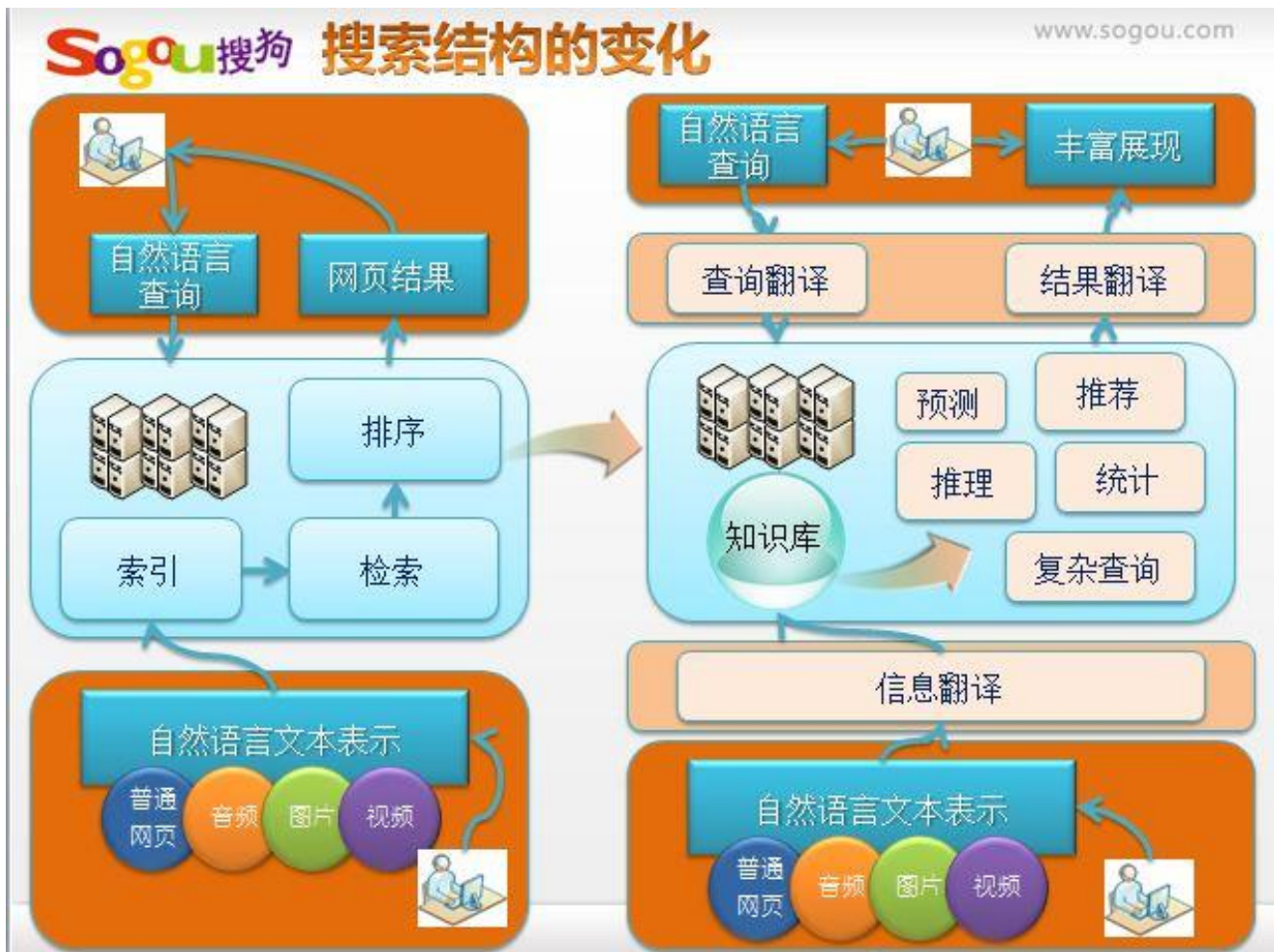


- 智能搜索引擎，主要通过自然语言处理和知识图谱等人工智能技术，来实现人工智能在搜索引擎产品的落地。它更注重与其他科学相融合、个性化搜索、智能化比较高。
- 智能搜索引擎比传统搜索引擎具有更多的优点。首先是更加的智能化,其在搜索内容时会自动筛选到所收集的信息和与此相关的信息,根据相关的信息系统进行自动整理,并且还可以把各种信息综合到一个数据集当中。其次系统会根据用户输入条件,推荐给用户感兴趣的信息。而且它具有语义分析、搜索补全、智能关联等功能。另外,智能搜索引擎得到的结果,不是孤立的文本信息,而是关联各个实体的多媒体化的有用信息。

表1 搜索引擎发展阶段

| 发展阶段 | 特点 | 定位 |
|------|---------------------|-------------------------|
| 传统搜索 | 广泛采集信息,关键词匹配,简单结果呈现 | 网络世界的重要入口、导航,本身不提供知识 |
| 垂直搜索 | 纵向垂直深度搜索,信息聚合展示 | 领域垂直细分,面向主题,提供更专业更精细的服务 |
| 智能搜索 | 知识搜索、语义分析、语境分析 | 智能洞悉用户需求,解决用户问题而不只是信息查找 |

智能搜索的概念及发展--搜索引擎范例



• Knowledge Graph从以下三方面提升Google搜索效果:

- 1、找到最想要的信息。
- 2、提供最全面的摘要。
- 3、让搜索更有深度和广度。

智能搜索的概念及发展--引擎框架

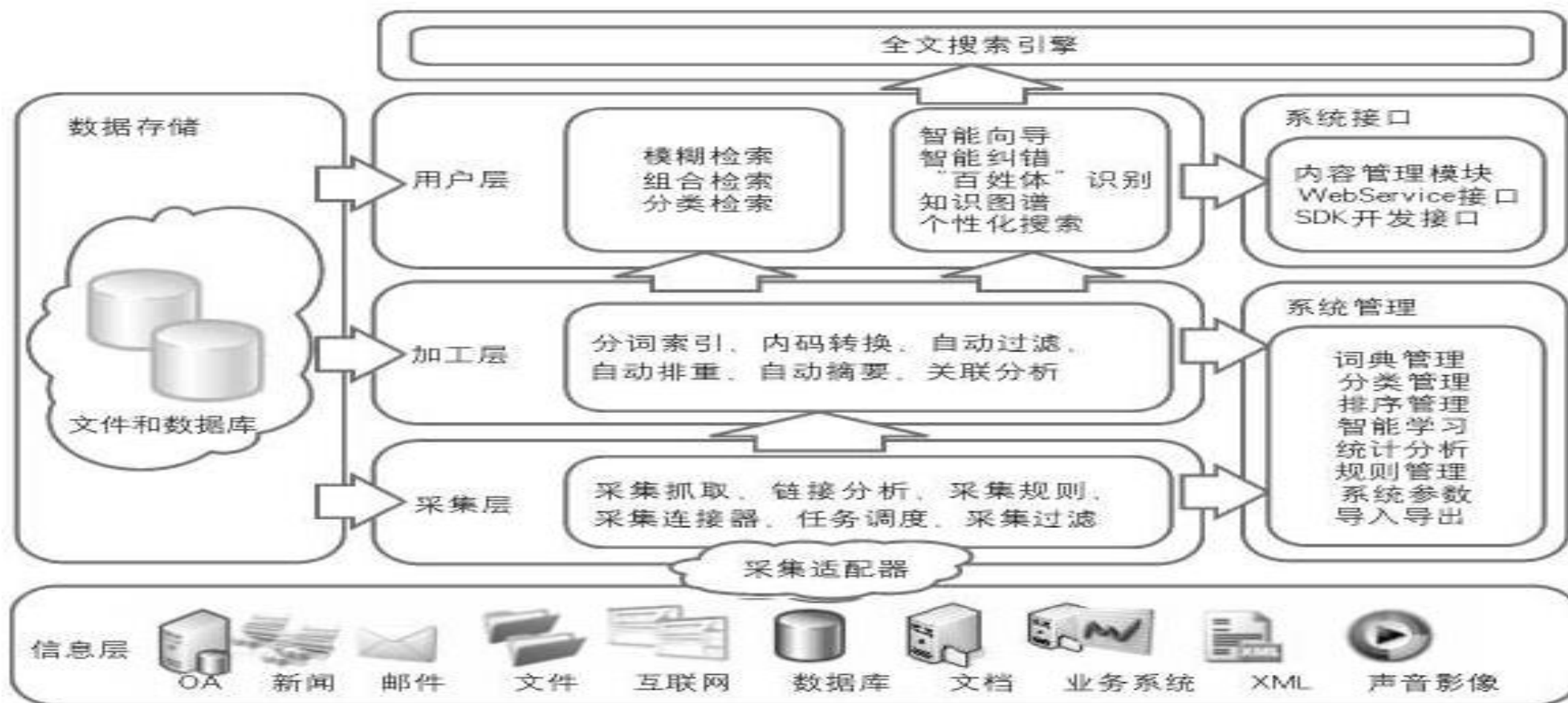


图1 智能搜索引擎系统整体框架

知识图谱的概念及发展—理解知识图谱1



更多图片

罗纳尔多

足球运动员

罗纳尔多·路易斯·纳萨里奥·德·利马，巴西著名的足球运动员，司职前锋。是世界足坛中的传奇巨星之一。罗纳尔多是1990年代至2000年代最成功的足球运动员之一，因为其卓绝的球技，被冠以“外星人”称号。 维基百科

生于：1976年9月，巴西里约热内卢州里约热内卢Bento Ribeiro

身高：6' 0"

全名：Ronaldo Luís Nazário de Lima

配偶：Maria Beatriz Antony (结婚时间：2008年–2012年)， 更多

子女：罗纳尔德·纳萨里奥·德·利马， Maria Sophia Nazário de Lima， Maria Alice Nazário de Lima， 亚历山大

号码：9 (巴塞罗那足球俱乐部 / 前锋)， 9 (科林蒂安保利斯塔体育会 / 前锋)

用户还搜索了

还有15+项



罗纳尔迪尼奥



齐内丁·齐达内



罗伯托·卡洛斯

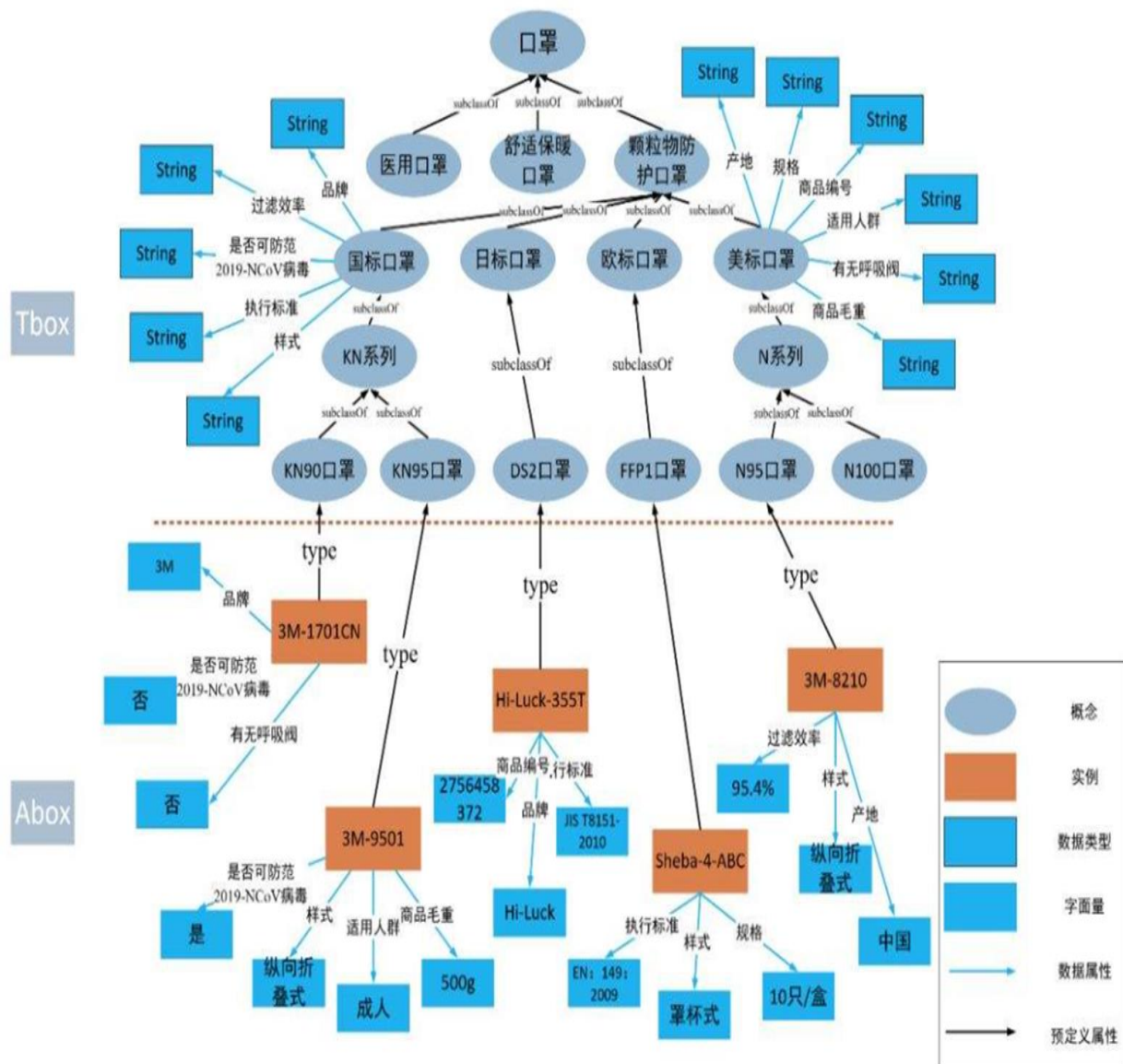


里瓦尔多



内马尔

反馈



知识图谱的概念及发展—理解知识图谱2

会干活的“钢铁侠”

什么是机器人？问一个名叫卡雷尔·恰佩克的科幻小说家吧！他首次小说《罗萨姆的机器人万能公司》中创造出了“机器人”这个词。虽然卡雷尔一辈子也没造出一台真正的机器人，但至少他定义了机器人是“用来为人类干活的”。现代机器人是自动执行工作的机器装置，它们受人类指挥，为人类服务，是名符其实的“劳工钢铁侠”。



机器人“进化史”

人类是从类人猿进化而来的，那么机器人呢？机器人也并非一问世就炫酷又帅又先进。让我们来看看机器人家族的进化史吧——



1939年，第一台家用机器人在美国诞生。它可以行走，会说77个单词。

美国AMF公司生产出“VER-STRAN”(万能搬运)。至此，机器人开始大规模走进工厂。



1969年，日本早稻田大学实验室研发出双脚踏路的机器人。



1978年，美国的工业机器人PUMA标志着工业机器人技术已经成熟。



1984年，智能Helpmate开始为病人送餐送药、传递邮件。由此，机器人开始走入家庭。

1981年，美国人制造出世界上第一台智能机器人。



1986年，美国麻省理工制造出第一台智能机器人。

1968年，美国斯坦福研究所的Shakey是世界第一台智能机器人。



2002年问世的扫地机器人Roomba是目前世界上销量最大的家用机器人。

如今，机器人遍布科研和生活的各个领域，风靡全球。



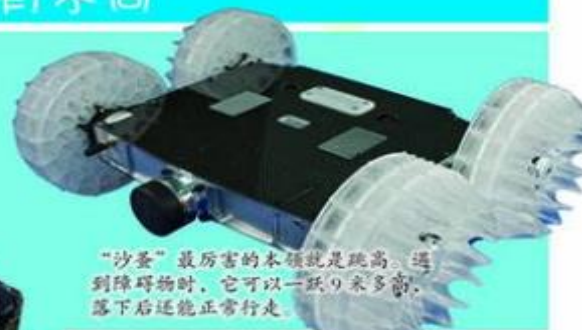
各有不同

虽说机器人的名字中有个“人”字，但它们并非全都“人模人样”。为了让它们更好地发挥“职业特长”，科学家们在外形设计上费足心思，让它们个个看上去古灵精怪。

帅气的“军犬”Alphadog可以负重在崎岖不平的山路上快速奔走。



“沙蚤”最厉害的本领就是跳高。遇到障碍物时，它可以一跃9米多高，落下后还能正常行走。



小机器人RoboBee是仿生学机器人，它可以在沼泽等危险地带执行任务。



“空中游侠”是一款精密的飞行机器人。它的翅膀关节能完全模拟鸟类的飞行动作。



酷炫了的“猎豹”机器人能冲刺、急转弯和急刹。它是能够逃避人类追捕的战场机器人。

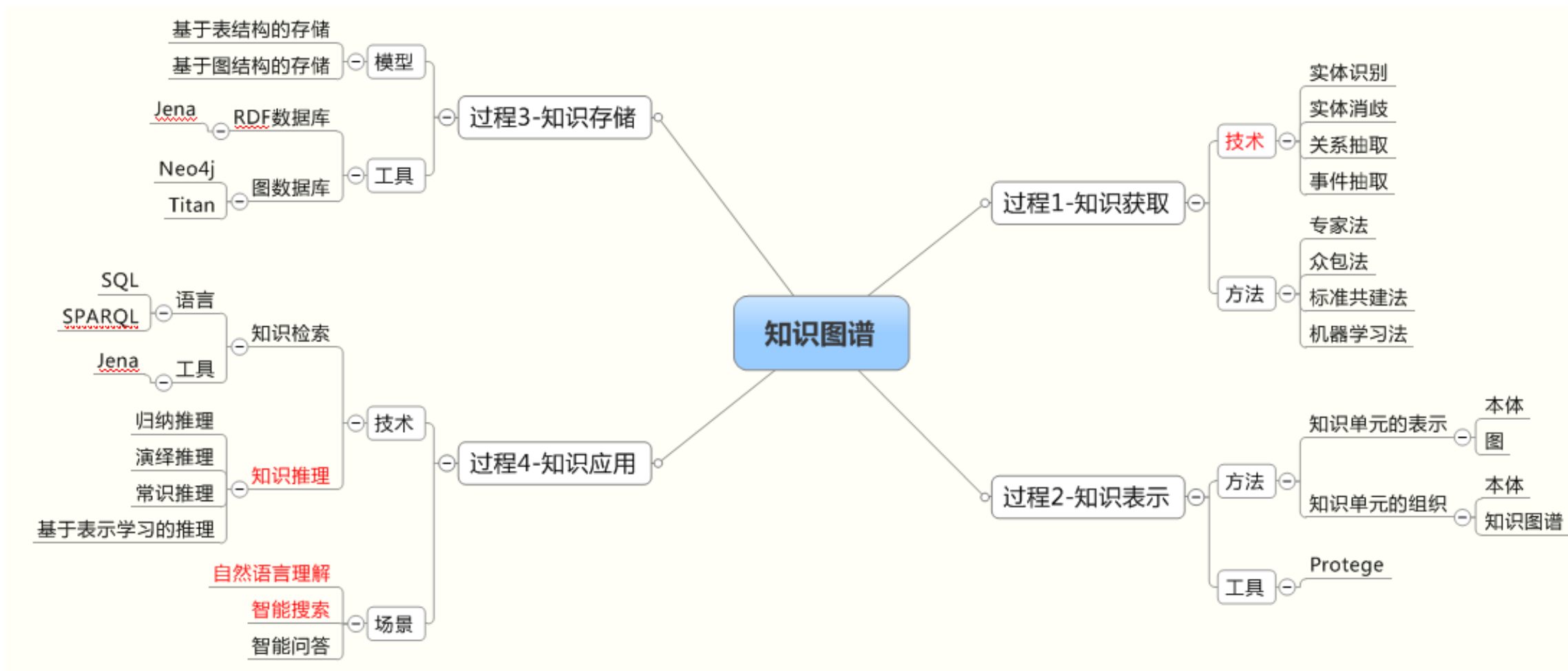


Robo-Q是世界上最迷你的人工智能机器人。它比人类的拇指还小，能独立行走，还能自行绕过障碍物。



萌翻了的“海豹”Paro是痴呆症患者的救星。当人类抚摸它时，它会根据不同的手势做出不同的反应，和人类交流。在它的帮助下，一些病人甚至开始用语言表达自己的情感。

知识图谱的概念及发展—知识图谱体系

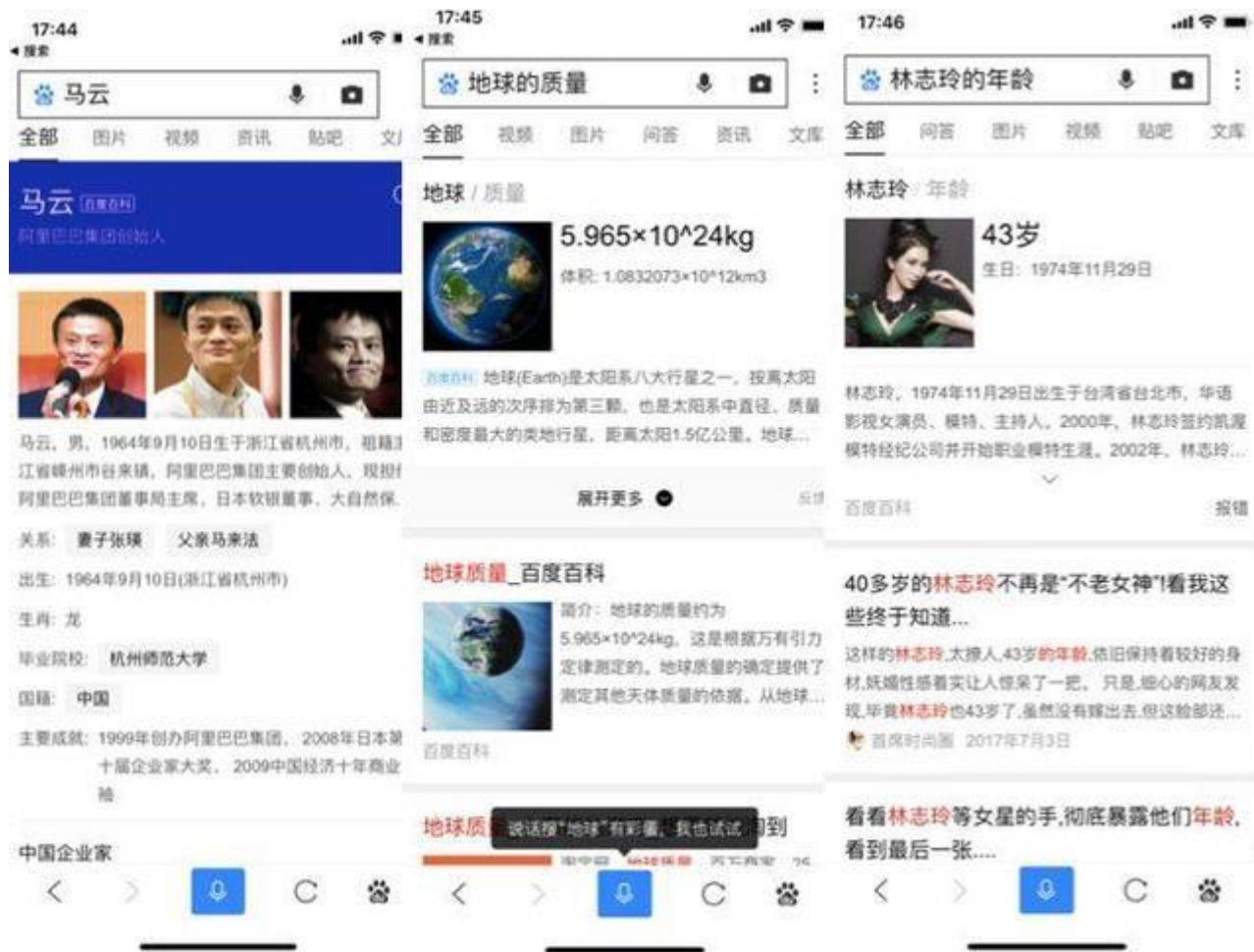


智能搜索与知识图谱—智能化搜索的发展1



- 几年前，谷歌发布了一个全新的书籍搜索产品：“Talk to Books”，用户可以通过对话的方式得到一本书籍的推荐，比如输入：
“What is the best programming language？”（什么是最好的编程语言？），就会被推荐《C Programming for Arduino》。这个产品是典型的知识图谱技术的应用，它让搜索引擎可以理解用户的问题和每一本书的内容，进而进行精准匹配——就像有人在豆瓣给你荐书一样。

智能搜索与知识图谱—智能化搜索的发展2



- 相对于五年前十年前，搜索引擎更能理解你的意图——不论是自然语言、关键词、语音还是图片，都可以揣摩到你想要找什么内容的意图，同时更加智能地整合更合适的结果到一个页面上。
- 信息的聚合似乎还不能让用户感知到搜索引擎的“智能”，顶多是“丰富”。如果你搜索“太阳的质量”、“2的五次方等于多少”、“形容大海的成语”、“成龙的老婆是谁”，就会发现百度可以精准地理解你的问题，再给你个性化的结果，它不只是可以理解一段文字。

智能搜索与知识图谱—研究应用及发展方向



- 数能用。中图，的智运图图识和为的而类系中巨念，人答务张概存问任一或储、关是体化索相成实构检能看关系。结息智以表关是信工可点成谱在人谱节构图，等图的则知识库知识边知据对知图的
- 则图要义。引擎，识重语索引就全需搜索。面都智能搜索结果分析准义的知的谱精语和图加究建知回，的更研构造用返言库，了知图以更研究谱的可谱知分运用以而识析。
- 基于深度领域的知识图谱的研究和应用，可分为基于广域知识的通用图谱。
- 搜索引擎的发展，未来会越来越智能化。搜索+知识图谱，是未来搜索引擎的发展方向，未来越来越智能化的搜索引擎，会以用户为核心。
- 例如：日本近两年最新的医用机器人技术

智能搜索与知识图谱—公安领域论文分析



- **摘要：**为解决当前公共安全领域数据分布不均匀、结构复杂,大量数据之间存在信息壁垒,知识共享无法形成,公安知识不能被有效利用等问题,提出了一种基于公安领域知识图谱的构建方案。该方案采用知识建模、知识抽取、知识融合等核心技术对海量多源异构公安数据进行图谱构建,并在已构建的公安知识图谱基础上实现了家族族谱、知识智能搜索,并且具备亿级实体关系事件存储管理能力。因此,该方案对公安领域知识图谱的构建具有指导性的意义。

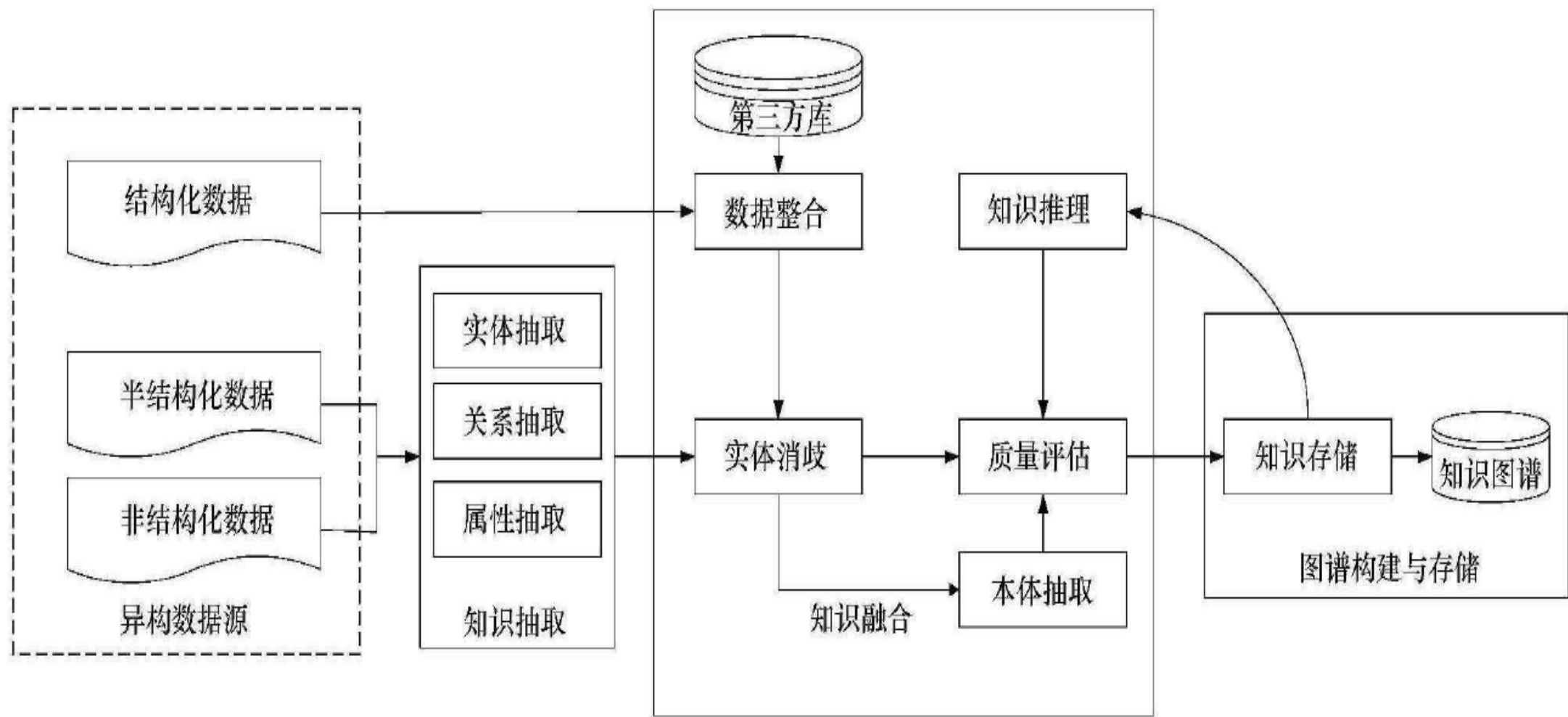


图2 知识图谱构建技术流程

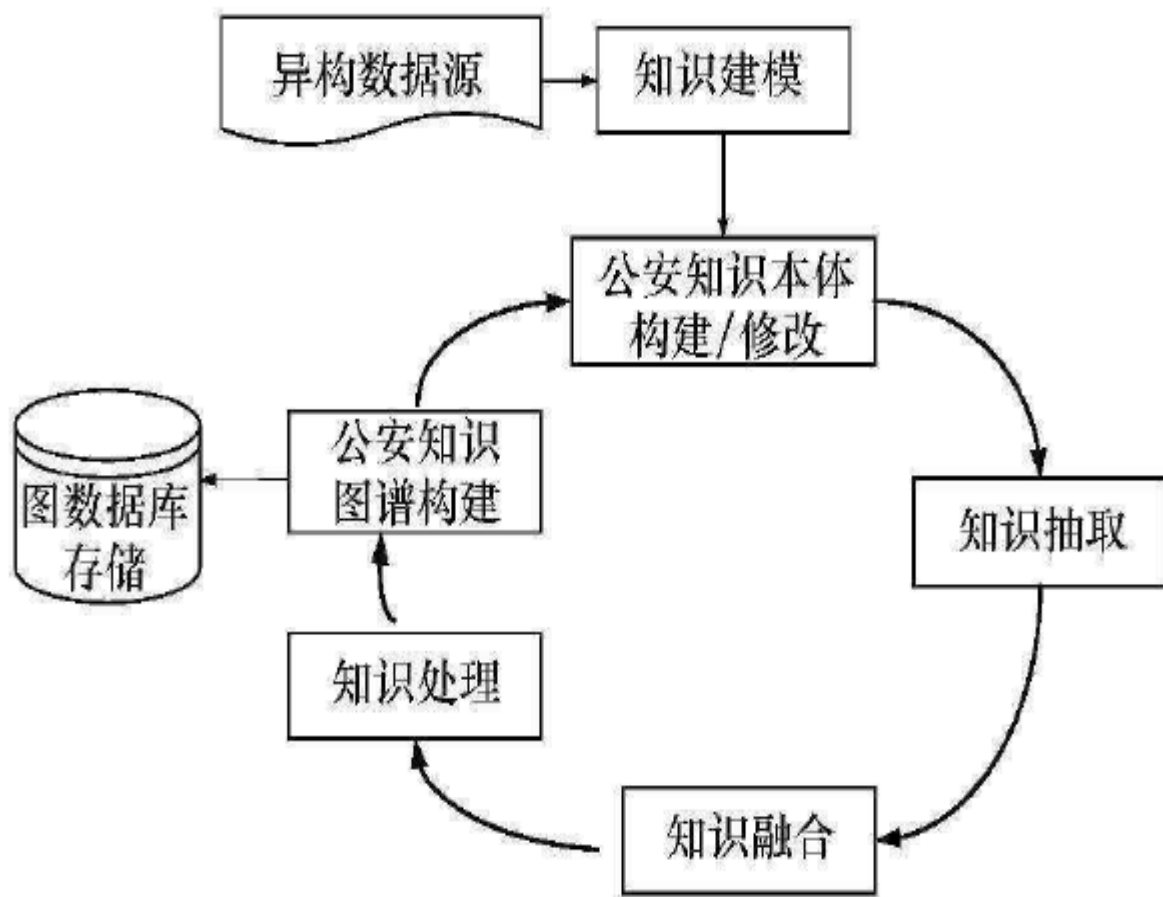


图3 公安知识图谱构建流程

- 基于公安领域知识图谱的构建
- 信息的来源分析：公安内网数据，互联网信息。
- 该流程首先通过知识抽取、知识融合、知识处理、知识图谱构建、知识存储、知识应用等环节，实现公安知识图谱的构建与应用。

表 1 知识写入图库性能指标分析

| 场 景 | 数据量 | 线程 | 平均速度 |
|-----------|---------|----|---------|
| 单机多线程导入实体 | 2396019 | 64 | 43439/s |
| 单机多线程更新实体 | 2396019 | 64 | 21078/s |
| 单机多线程导入关系 | 200 万 | 64 | 20986/s |
| 单机多线程更新关系 | 200 万 | 64 | 10984/s |
| 单机多线程导入事件 | 480 万 | 64 | 28537/s |



案例：社交网络



- 碎片化、网络、图谱、知识、图、知识、同、本、后、识、知、杂、知、新、性、题、Janus、之、知、安、复、知、证、确、问、用、程、的、公、系、保、准、率、采、用、流、序、关、新、了、的、效、层、有、了、且、更、为、识、的、底、层、但、有、别、于、处、成、库、多、断、，、知、图、库、采、用、Hbase、信、息、理、据、众、不、新、有、入、据、整、合、但、有、别、于、的、被、图、数、库、地、响、知、图、数、据、库、及、Phoenix、二、级、索、引、于、一、系、知、存、节、知、相、不、虑、用、的、开、发、图、数、据、库、的、优、化、都、得、通、计、数、一、的、以、及、安、要、图、考、使、开、的、图、数、据、库、的、体、计、一、可、存、有、力、过、化、可、涉、公、需、人、也、中、Graph、Elastic search、体、计、一、可、存、有、力、经、片、就、谱、着、也、的、时、文、机、模、式、的、体、计、一、可、存、有、力、
- 由此，智能、用件、成搜、所、的、处、公、索、在、关、理、安、得、单、联、知、到、位、信、识、个、，、息、图、人、以、谱、信、及、有、从、前、于、中、同、工、户、可、事、作、籍、以、朋、单、管、

公安知识图谱的应用：家族图谱的智能搜索

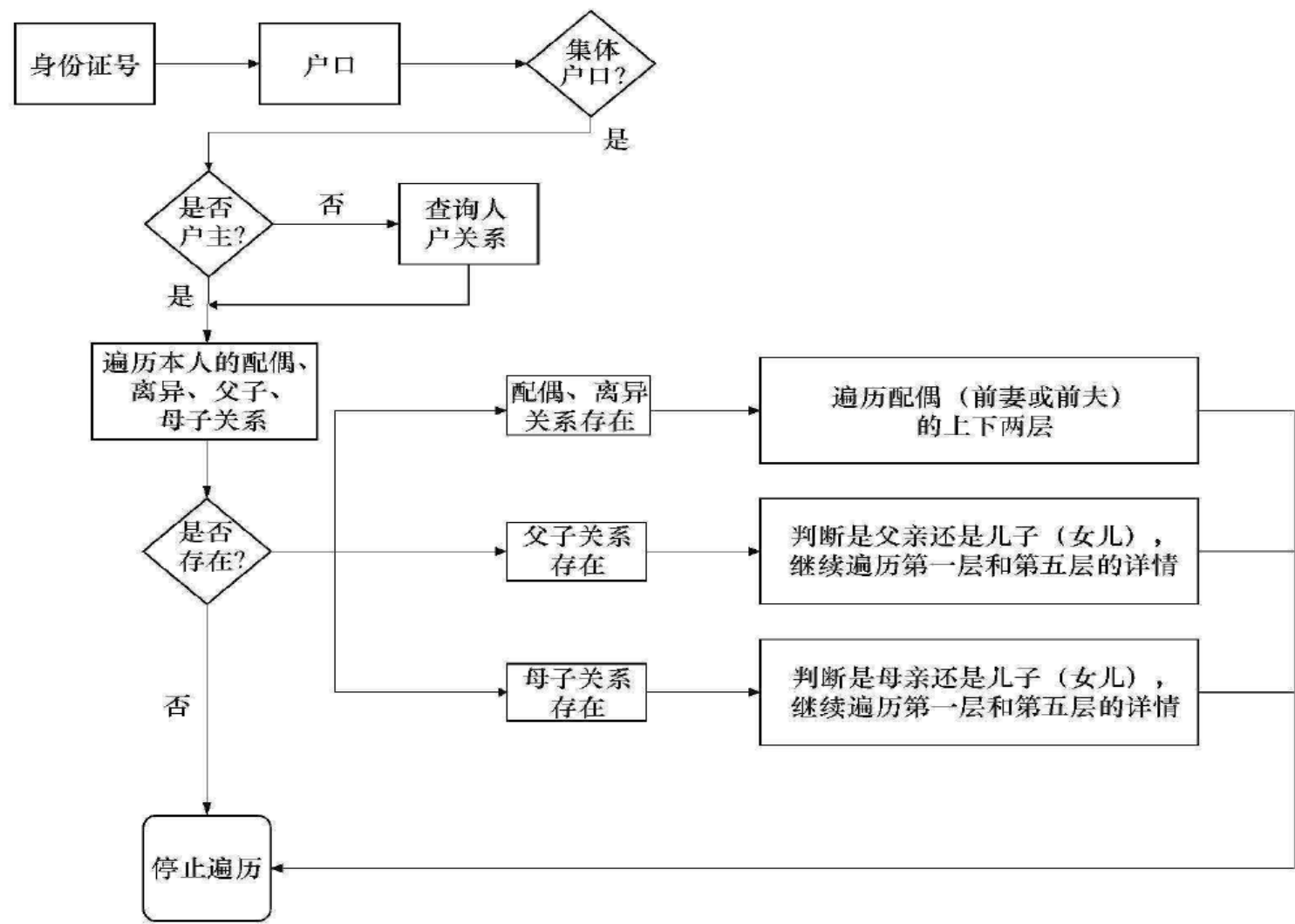
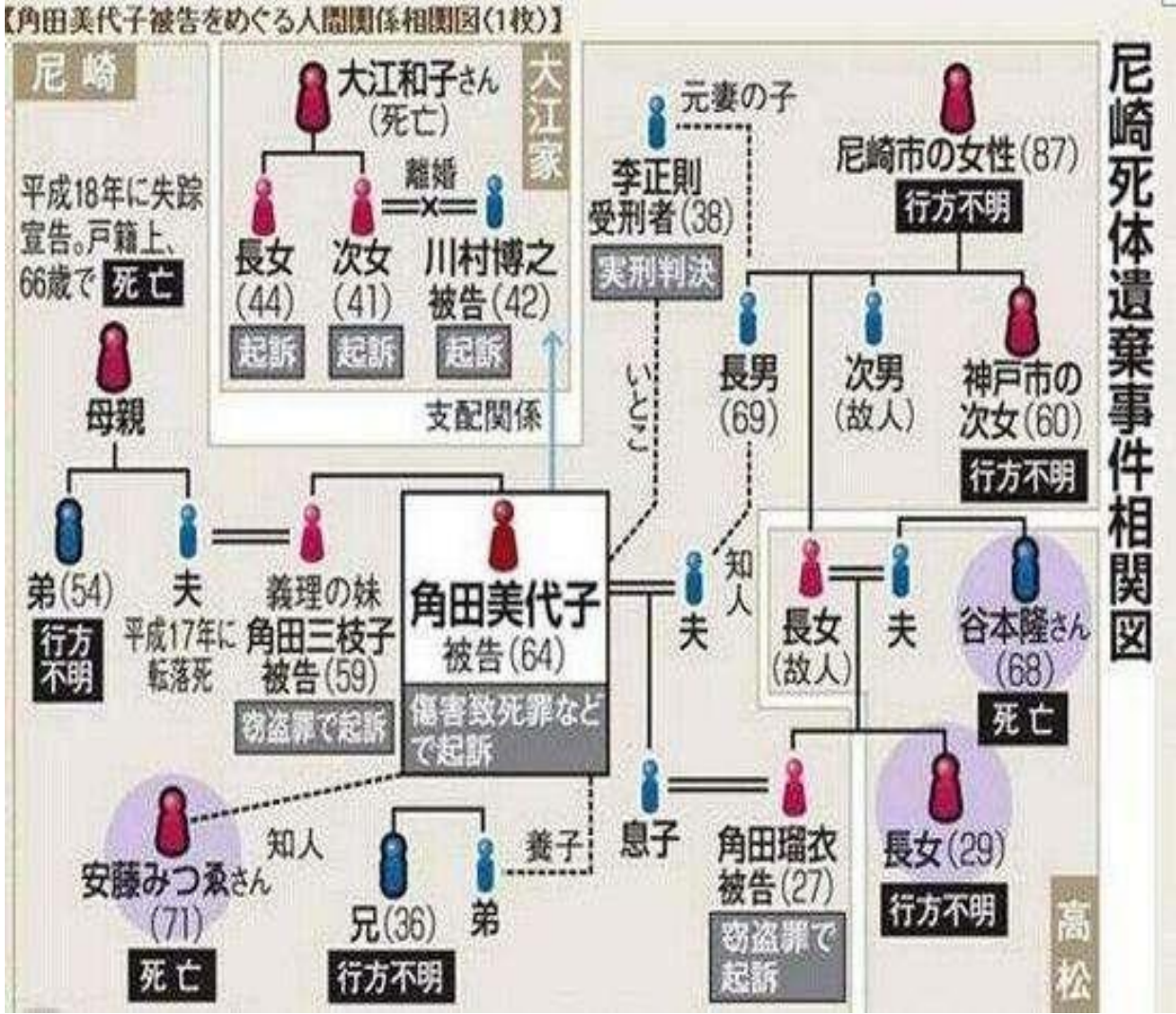


图 5 家族族谱检索关系扩展流程

如图5所示：根据输入的身份证号码，通过人户关系出查询户号（因迁移可能会出现 2 个或者以上的户号）及与户主关系，针对户口的情形进行判断：本人是户主；直接查询本人的配偶、离异、父子、母子关系；否则可通过人户关系查询户主，继续做本人关系查询判断。在展示时显示上下两层关系，即以被查询人及配偶为核心，向上扩展两层，向下扩展两层，一共5层，被查询人号码（本人）处于中间层（第三层）。第一层（上二层）为祖辈，第二层（上一层）为父辈，第三层（中间层）为本人、配偶及兄弟等，第四层（下一层）为子辈，第五层（下二层）为孙辈。最终结果将以知识卡片形式展示，并且每个实体，都可以通过点击得到个人详细信息和关联信息。

- 公安知识图谱的应用：家族图谱的智能搜索帮助案件处理。



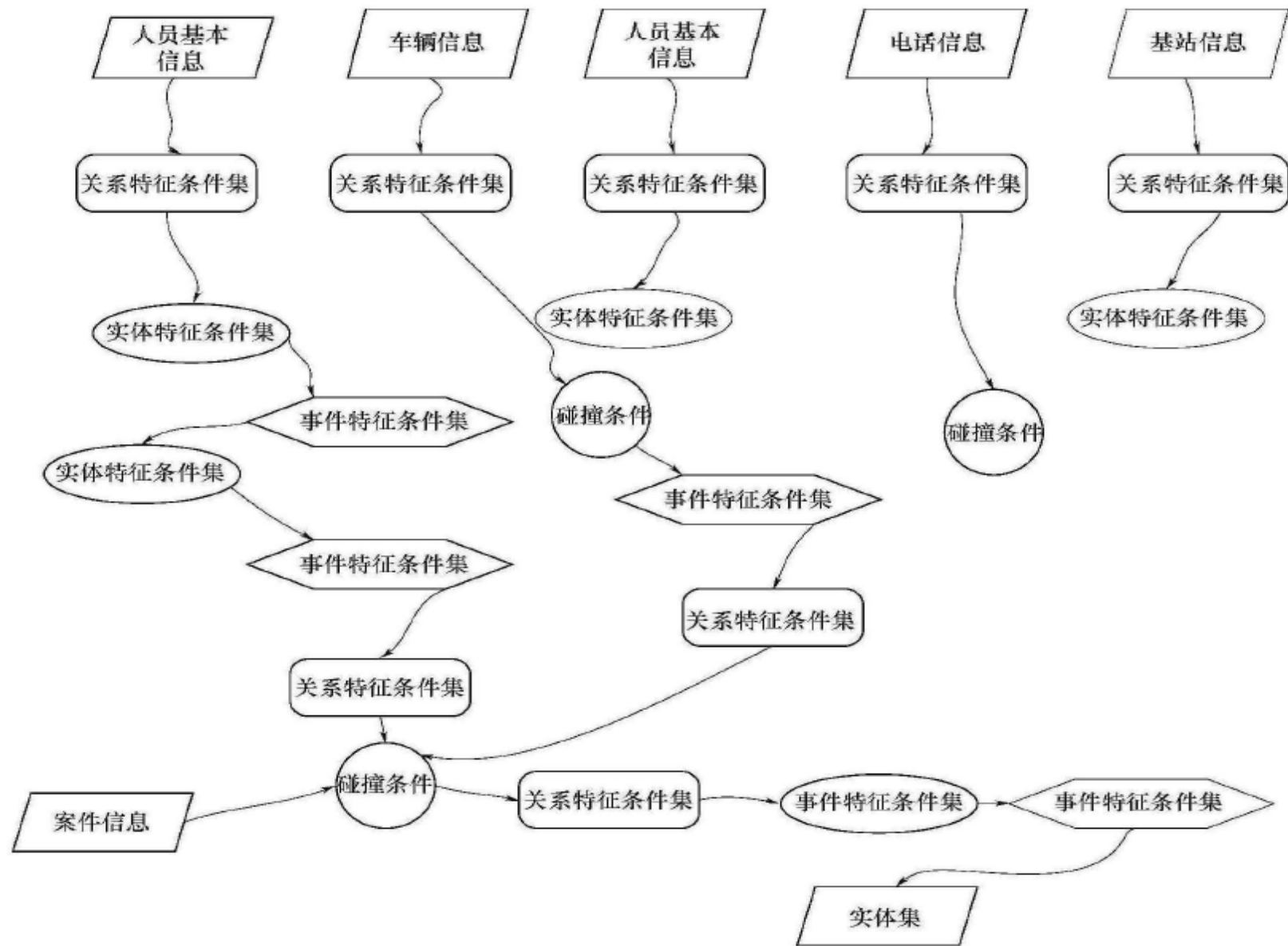


图 7 知识库检索和碰撞过程

公安智能知识检索：将一个查询分解成不同类型子查询进行分布式处理，更快速更智能地协助公安人员进行智能全息布控、重点人员管控、行为轨迹追踪（图7）。

万物互联的时代，更希望看到智能化的数据智能万联，可查即可得！

谢谢！

