

人工智能导论

主讲：王博

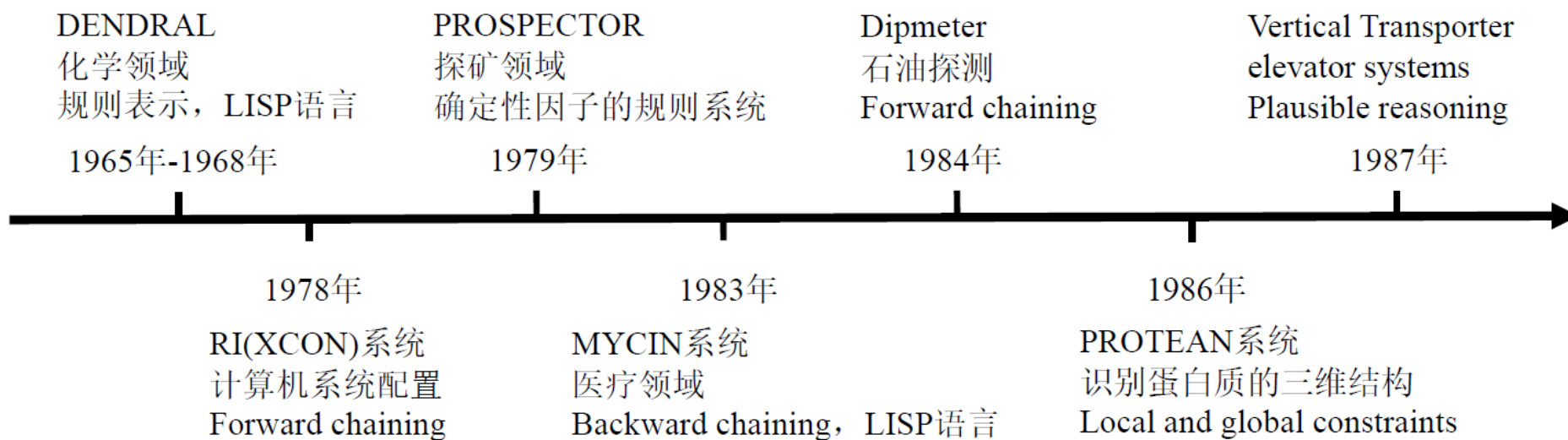
人工智能与自动化学院

符号主义 - 目录

- 3.5 证据理论
- 3.6 模糊推理
- 3.7 知识图谱

传统知识工程代表性系统

- 传统知识工程在规则明确、边界清晰、应用封闭的应用场景取得了巨大成功



传统方法的特点和困难

- 自上而下：严重依赖专家和人的干预
- 知识获取困难：隐性知识、过程知识等难以表达，存在主观性、不一致性，难以完备
- 知识应用困难：超出知识边界、需要常识、处理异常、不确定性推理、知识更新

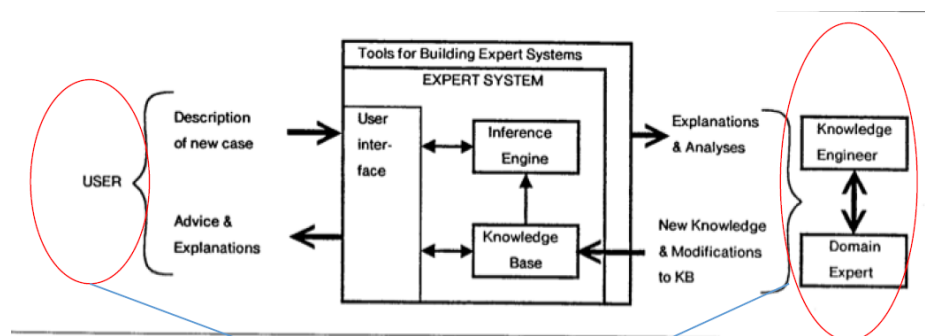
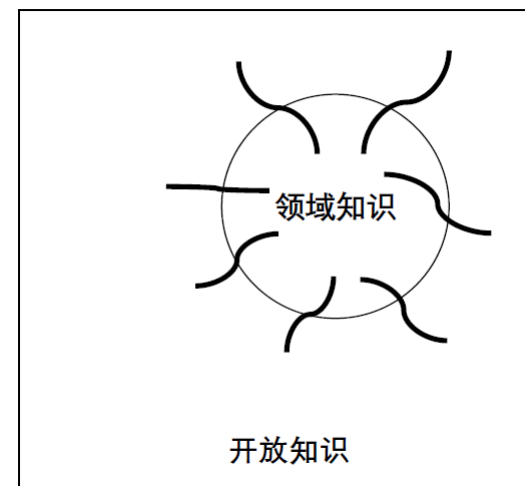


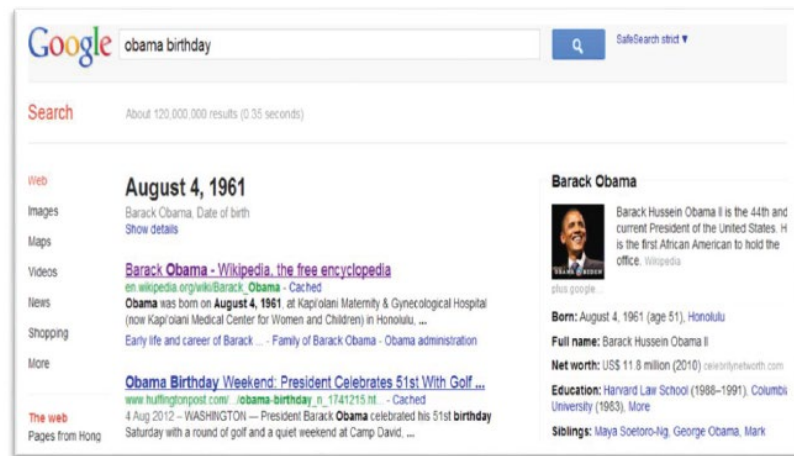
FIGURE 1-2 Interaction of a knowledge engineer and domain expert with software tools that aid in building an expert system. Arrows indicate information flow.

MYCIN专家系统中的人工参与部分



知识图谱的诞生

- 2012年5月，Google收购Metaweb公司，并发布**知识图谱**
- 搜索核心需求： 让搜索通往答案
 - 无法理解搜索关键词
 - 无法精准回答
- 根本问题
 - 缺乏大规模背景知识
 - 传统知识表示难以满足需求

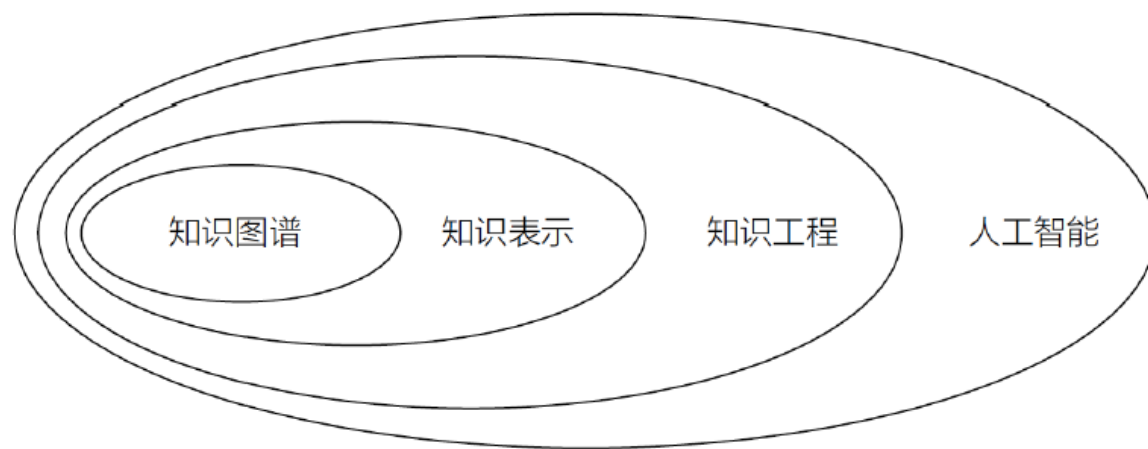


知识图谱的诞生

- Goc

知识图谱与人工智能

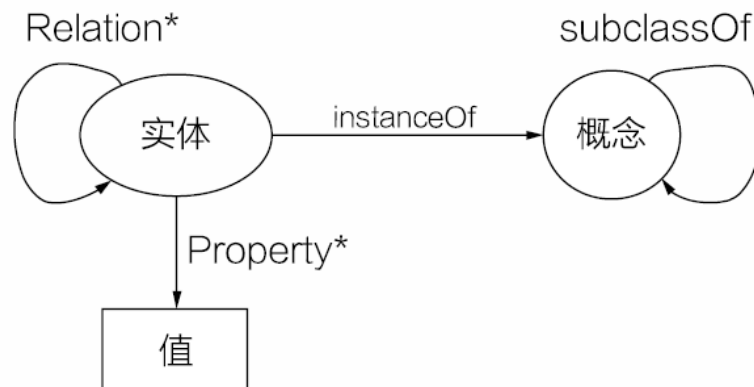
- 作为一种**技术体系**，是大数据时代**知识工程**的代表性进展
- 作为一门**学科**，知识图谱属于**人工智能**范畴
- **知识表示**是发展知识工程最关键的问题之一，而知识表示的一个重要方式就是知识图谱



知识图谱的学科地位

语义网络

- 语义网络是一种以图形化 (Graphic) 的形式、通过点和边表达知识的方式，其基本组成元素是点和边
- 1968年罗斯·奎利恩最先提出
- 节点表示实体、概念和情况等，边表示节点间的关系



语义网络的组成（图中星号表示可以存在多个不同的属性或者关系）

二元语义网络

- 表示一些简单事实，如占有关系和其它情况：以节点表示实体与概念，节点间关系以有向链关联：

- 燕子是一种鸟



- 小燕是一只燕子，燕子是一种鸟

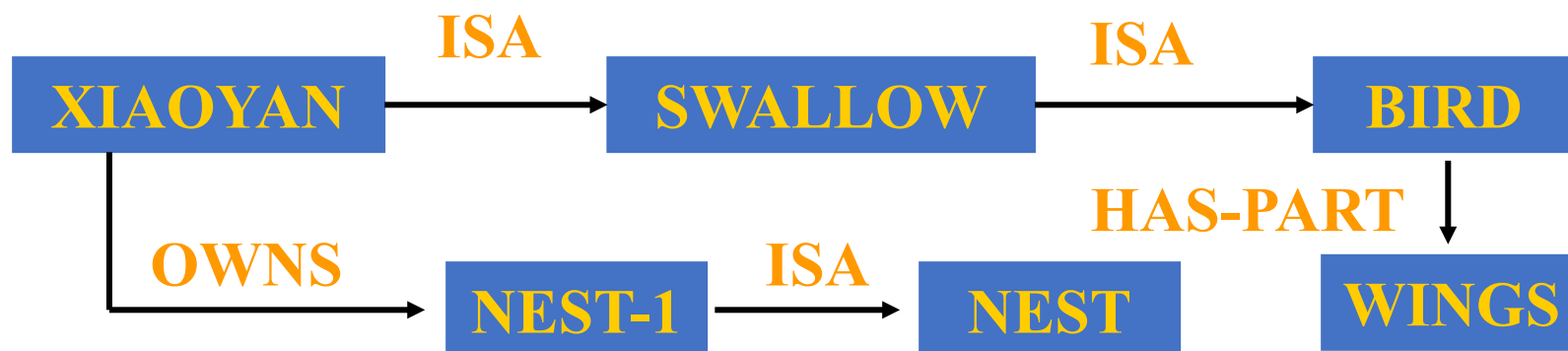


- 小燕是一只燕子，燕子是一种鸟，鸟有翅膀



二元语义网络

- 例：小燕是一只燕子，燕子是一种鸟，鸟有翅膀；巢-1是小燕的巢，巢-1是巢中的一个。



问题：

上述的语义网络为二元关系，无法表示复杂事实，如：小燕从春天到秋天占有巢-1。

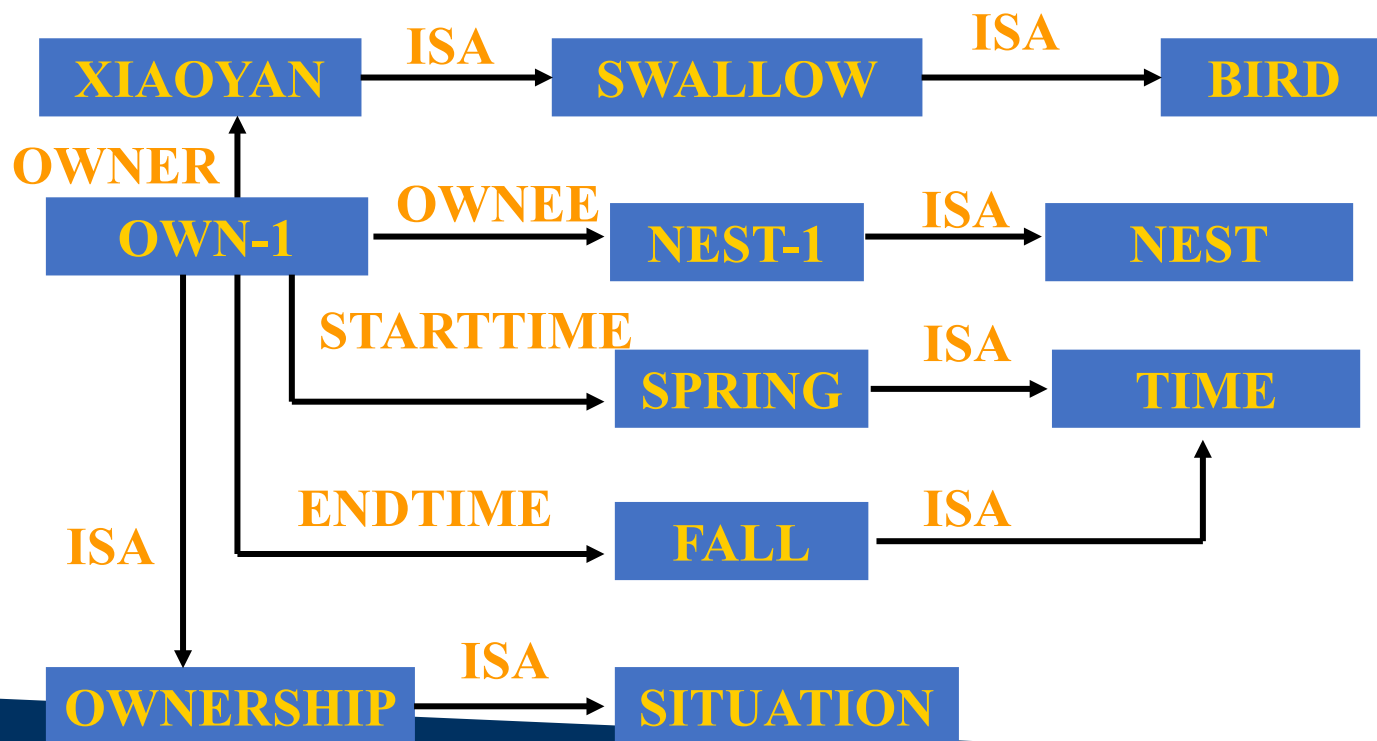
如果采用谓词逻辑表示为一个四元谓词演算：

Owens (XIAOYAN,NET-1,SPRING,FALL)

思考：用语义网络如何表示上述四元谓词公式？

多元语义网络

- Simmons与Slocum扩展了该基本方法：
 - 允许节点既可以表示一个物体或一组物体，也可以表示情况与动作。每一情况节点成为事例框，有一组向外的边，用以说明与该事例有关的各种变量。



多元语义网络

- 多元语义网络表示的实质
 - 把多元关系转化为一组二元关系的组合，或二元关系的合取。

可转换为

$$R(X_1, X_2, \dots, X_n)$$

$$R_{12}(X_1, X_2) \wedge R_{13}(X_1, X_3) \wedge \dots \wedge R_{1n}(X_1, X_n)$$

.....

$$R_{n-1\ n}(X_{n-1}, X_n)$$

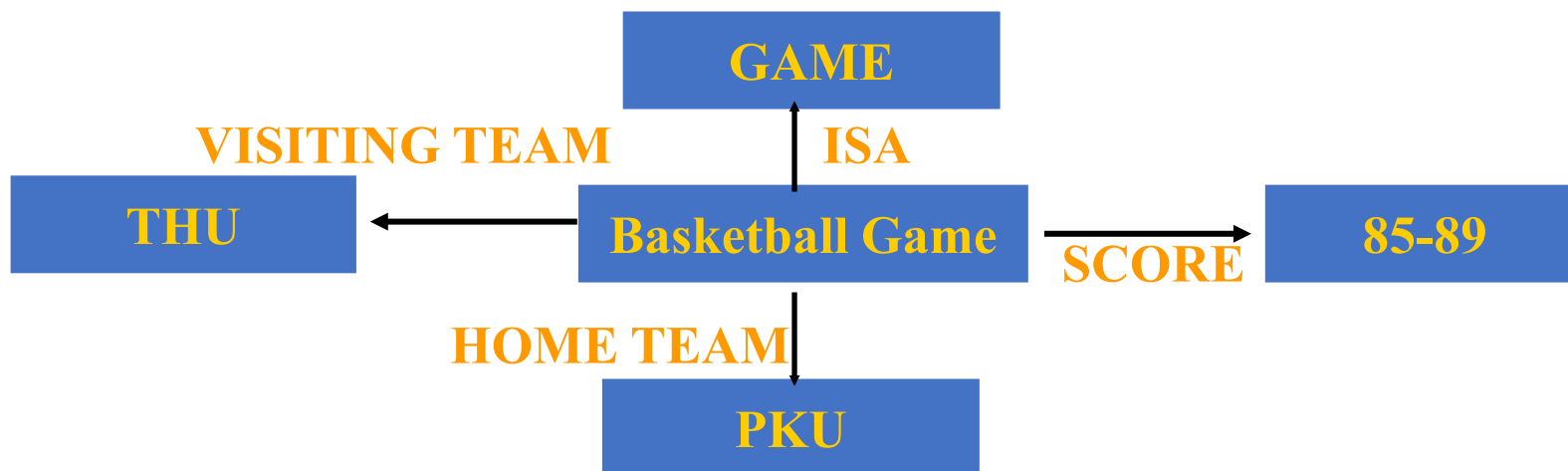
例：三根线a, b, c组成一个三角形

TRIANGLE (a, b, c)

CAT(a, b) \wedge CAT(b, c) \wedge CAT(c, a)

多元语义网络

- 例：北京大学和清华大学两校篮球队在北大进行一场比赛的比分是85比89。
 - 谓词逻辑法：SCORE(pku, thu, (85-89))
 - 语义网络法：



知识图谱的概念

- **知识图谱**(Knowledge Graph)本质上是一种**大规模语义网络**(semantic network)
 - 富含**实体**(entity)、**概念**(concepts)及其之间的各种**语义关系**(semantic relationships)
 - 是大数据时代知识表示的重要方式之一

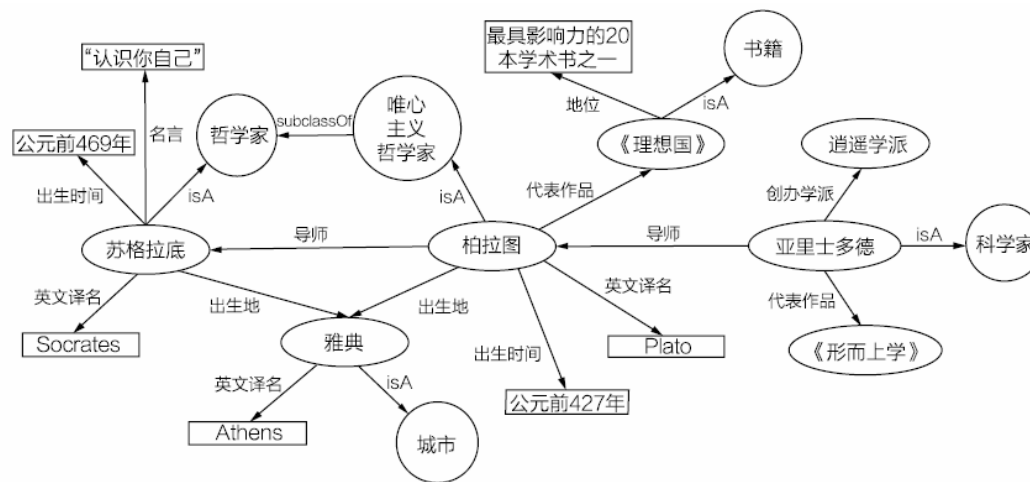


图 1-1 关于古希腊三大哲学家的知识图谱片段

知识图谱的组成：节点

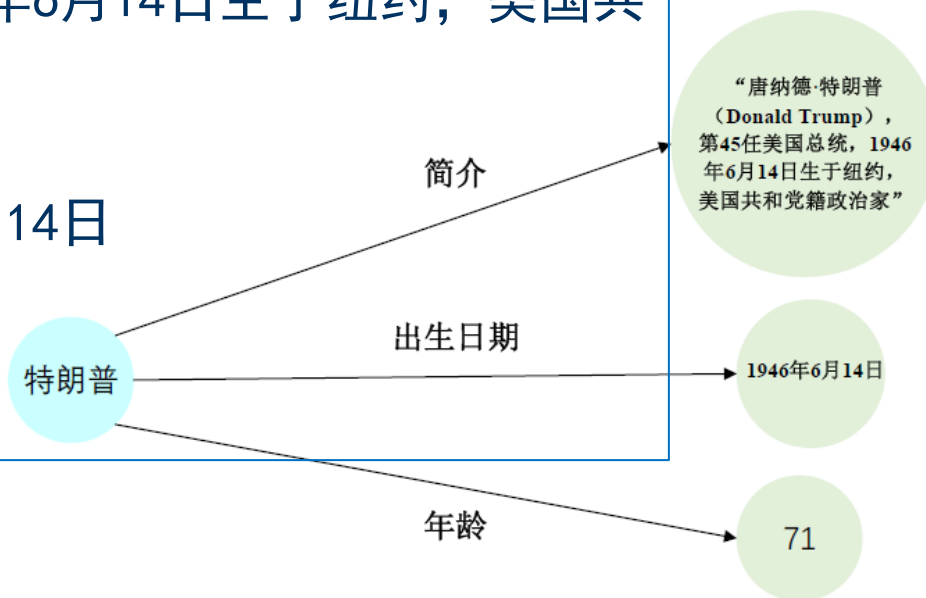
- 实体 Entity/Objects/Instances
 - 具有可区别性且独立存在的某种事物
 - Wikipedia: An entity is something that exists as itself, as a subject or as an object, actually or potentially, concretely or abstractly, physically or not.
 - 黑格尔《小逻辑》：能够独立存在的，作为一切属性的基础和万物本原的东西
- 概念 Concept/Category/Type/Class
 - 具有同种特性的实体构成的集合
 - Concept: In metaphysics, and especially ontology, a concept is a fundamental category of existence.
 - (mental) representations of categories
 - Category: Groups of entities which have something in common
 - Type/Class: A grouping based on shared characteristics.

知识图谱的组成：节点

- 值 Value
 - 对象指定属性的值

例：

- String
 - 特朗普简介：“唐纳德·特朗普（Donald Trump），第45任美国总统，1946年6月14日生于纽约，美国共和党籍政治家”
- Date
 - 特朗普生日：1946年6月14日
- Numeric
 - 特朗普年龄71



知识图谱的组成：边

• 关系 Relation

- 侧重实体之间的关系

- Examples:

Sitting-On: An apple sitting on a table

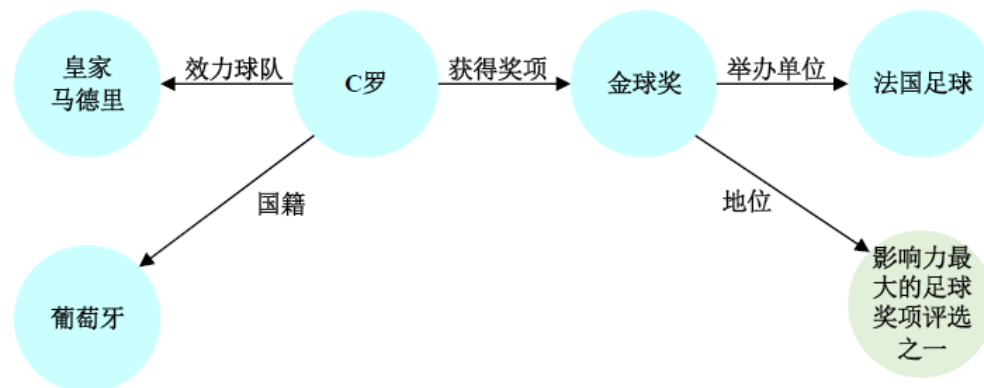
Taller-than: Washington Monument is taller than the White House

• 属性 Property/Attribute/Quality

- 描述一个物体的特性

- Examples:

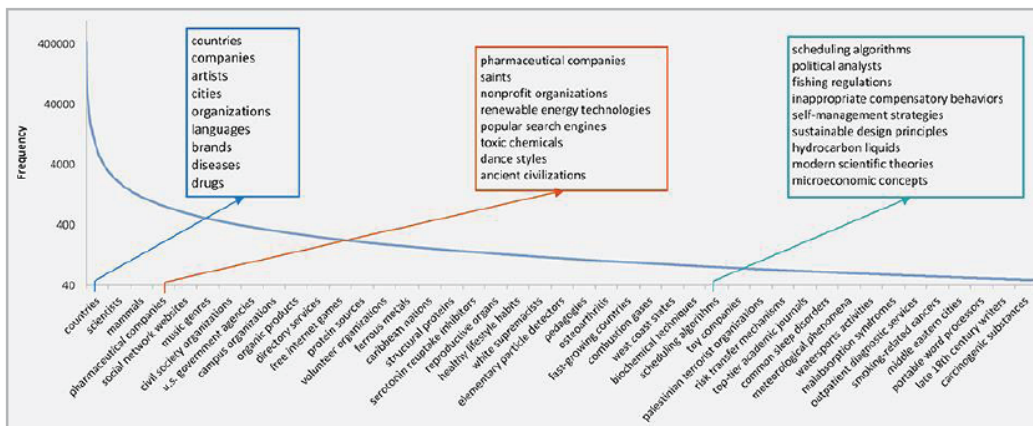
size, color, weight,
composition of an object



知识图谱的优势

- 尺度大 large scale
- Higher coverage over entities and concepts

KGs	# of Entities/Concepts	# of Relations
YAGO	10 Million	120 Million
DBpedia	28 Million	9.5 Billion
Probase	2.7 Million	70 Billion
BabelNet	14 Million	5 Billion
CN-DBpedia	17 Million	200 Million

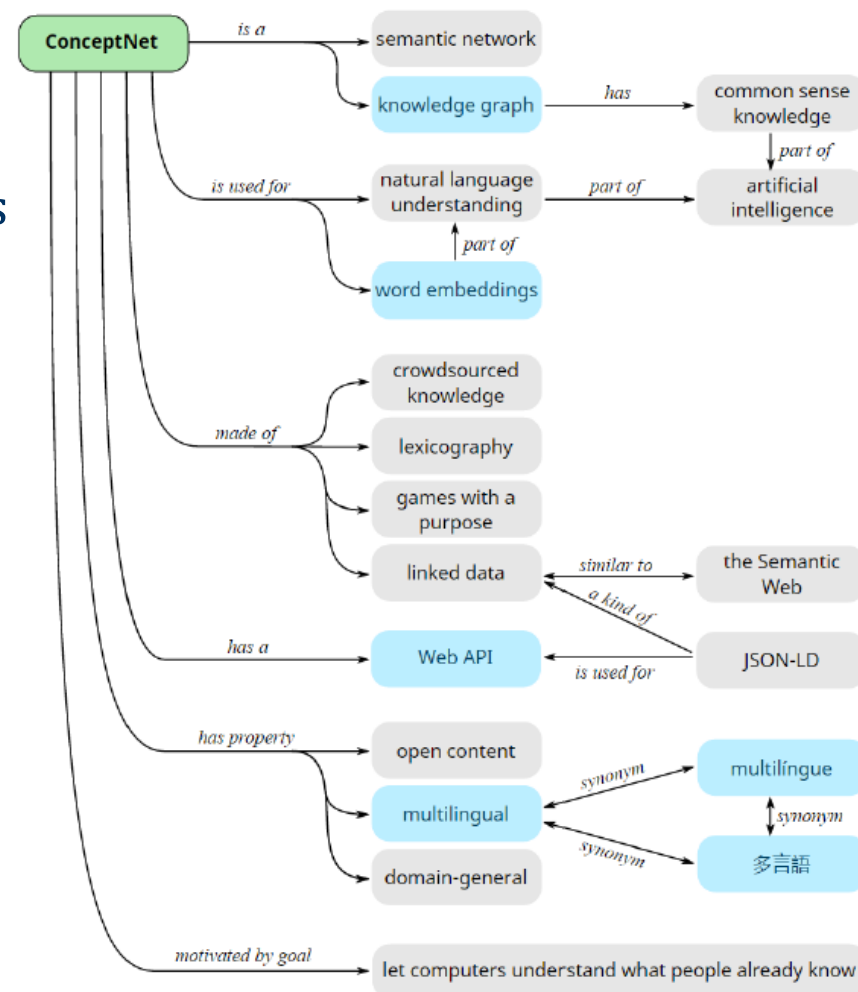


Existing Taxonomies	Number of Concepts
Freebase [5]	1,450
WordNet [13]	25,229
WikiTaxonomy [26]	111,654
YAGO [35]	352,297
DBPedia [1]	259
ResearchCyc [18]	≈ 120,000
KnowItAll [12]	N/A
TextRunner [2]	N/A
OMCS [31]	N/A
NELL [7]	123
Probase	2,653,872

知识图谱的优势

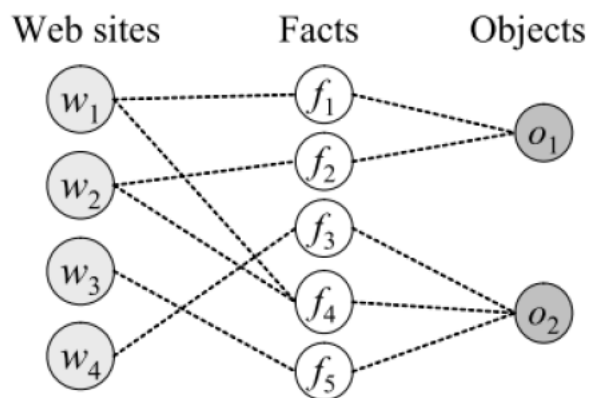
- 语义丰富 semantically rich
 - Higher coverage over numerous semantic relations

KGs	# of Relations
DBpedia	1,650
YAGO1	14
YAGO3	74
CN-DBpedia	100 Thousands



知识图谱的优势

- 质量高 high quality
 - 大数据 Big data: Cross validation by multiple sources
 - 众包 Crowd sourcing: quality guarantee



CN-DBpedia

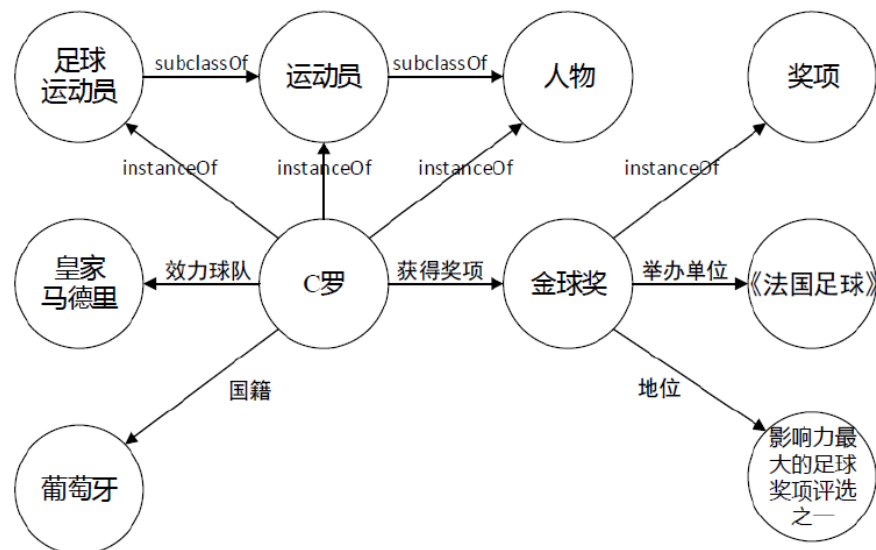
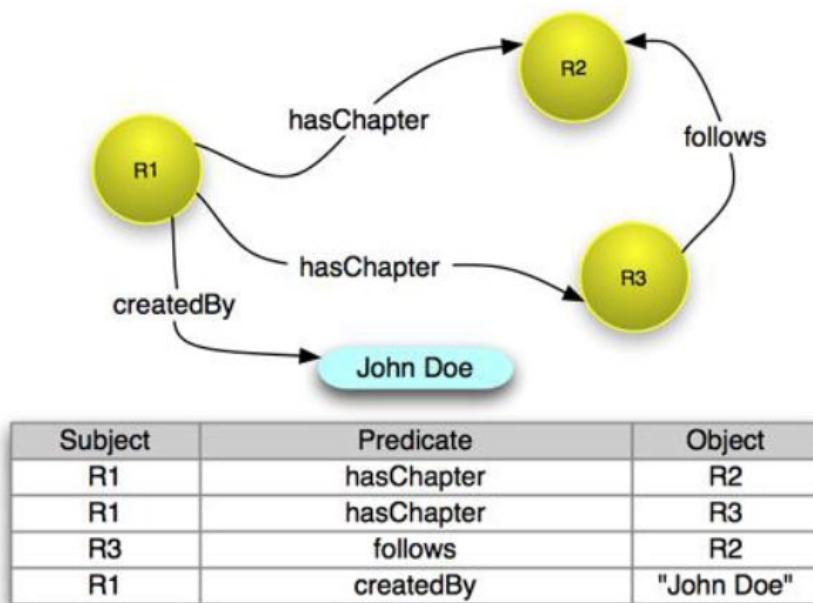
Q InfoBox

专职院士	25人	 
中文名	复旦大学	 
主管部门	中华人民共和国教育部	 
主要奖项	SCI论文单篇被引用次数全国第一	 
主要奖项	诺贝尔奖得主名誉教授10位	 



知识图谱的优势

- 结构易于实现 friendly structure
 - 组织架构
 - By RDF
 - By graph



知识图谱的挑战

- 高质量模式缺失

- 知识图谱在设计模式时通常会采取一种“经济、务实”的做法：也就是允许模式（Schema）定义不完善，甚至缺失
- 模式定义不完善或缺失对知识图谱中的数据语义理解以及数据质量控制提出了挑战

- 封闭世界假设不再成立

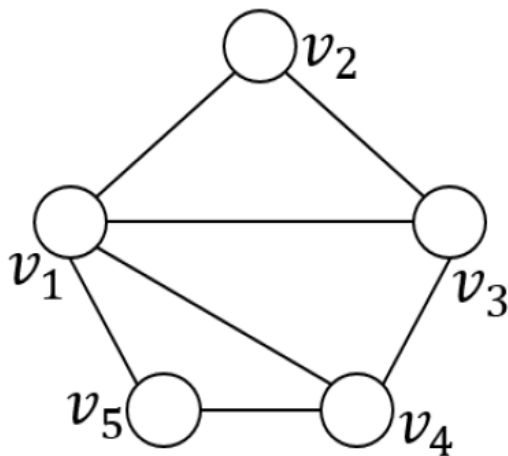
- 传统数据库与知识库的应用通常建立在封闭世界假设（CWA）基础之上，大多数开放性应用不遵守这一假设，在这些应用中缺失的事实或知识未必为假
- 不遵守CWA给知识图谱上的应用带来了巨大的挑战

- 大规模自动化知识获取成为前提

- 大规模自动化知识获取是知识图谱与传统语义网络的根本区别

基于图论的知识图谱表示

- $G=G(V, E)$
 - 其中 V 表示顶点集, $E \subseteq V \times V$ 表示边的集合。
 - 有向图、无向图
 - 邻接表、邻接矩阵
 - 度数、路径、可达 ...



0	1	1	1	1
1	0	1	0	0
1	1	0	1	0
1	0	1	0	1
1	0	0	1	0

基于三元组的知识图谱表示

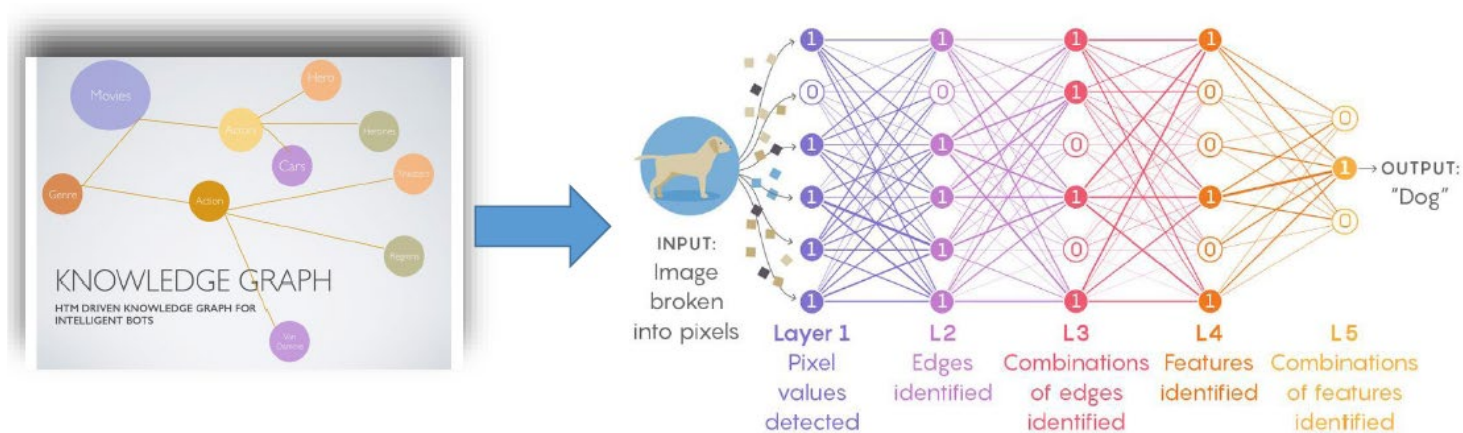
- **RDF** (Resource Description Framework)
 - 是用于描述现实中资源的W3C标准。
 - 现实中任何实体都可以表示成RDF模型中的资源，这些资源是对现实世界中概念、实体和事件的抽象。
- **三元组**包括三个元素：**主体** (subject)、**属性** (property) 及**客体** (object)
 - 也被称为**主体**、**属性**及**属性值** (property value)

e.g. <亚理士多德, 受到影响, 柏拉图>

主体 (Subject)	谓词 (Predicate)	客体 (Object)
<i>Aristotle</i>	<i>influencedBy</i>	<i>Plato</i>
<i>Boethius</i>	<i>placeOfDeath</i>	<i>Pavia</i>
<i>Chalcis</i>	<i>country</i>	<i>Greece</i>
<i>Pavia</i>	<i>postalCode</i>	27100

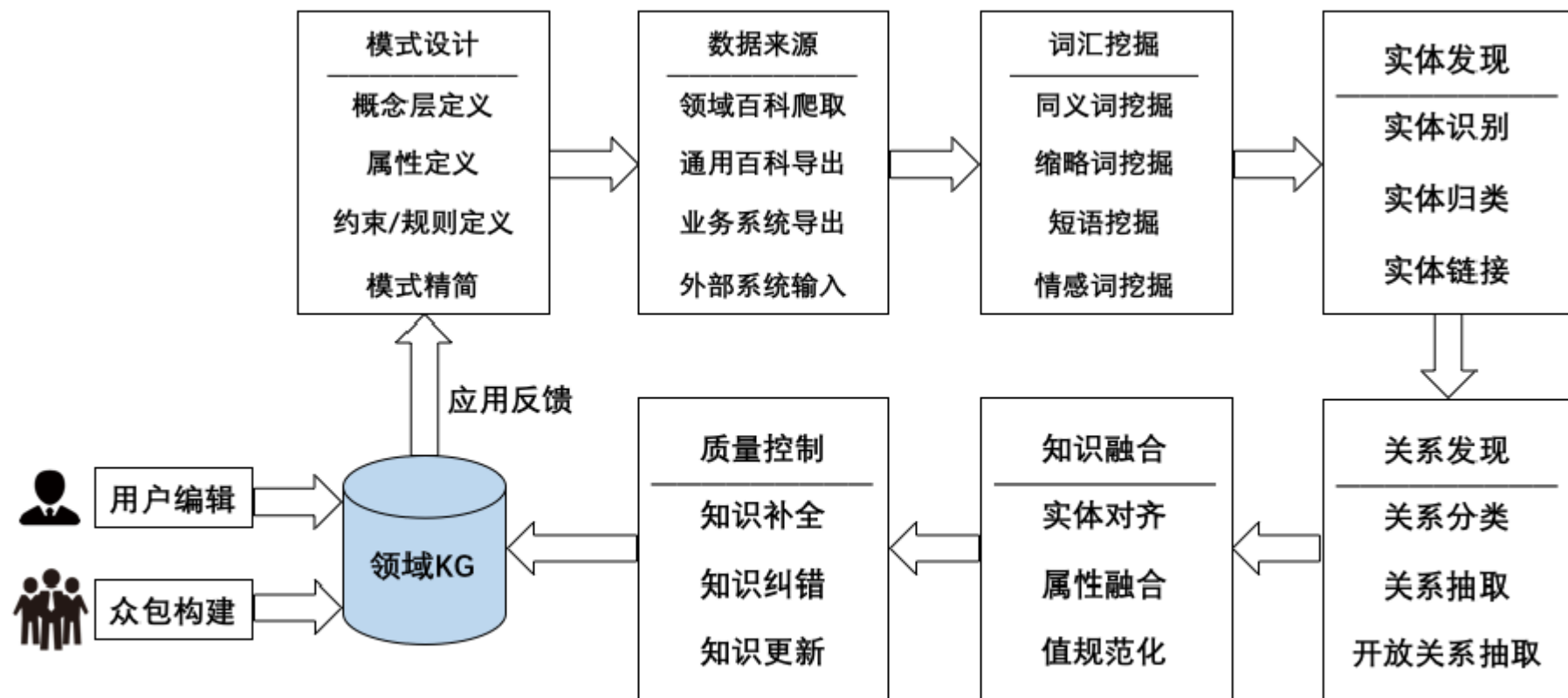
基于数值的知识图谱表示

- 将知识图谱中元素(包括实体、属性概念等)表示为低维稠密实值向量。
- 知识图谱的不同表示各有适用场景：
 - 向量化的表示面向机器处理
 - 符号化表示面向人的理解
 - 符号表示更易于理解、实现符号推理



将知识图谱中的点与边表达成数值化向量

领域知识图谱构建



领域知识图谱构建的基本流程

领域知识图谱构建

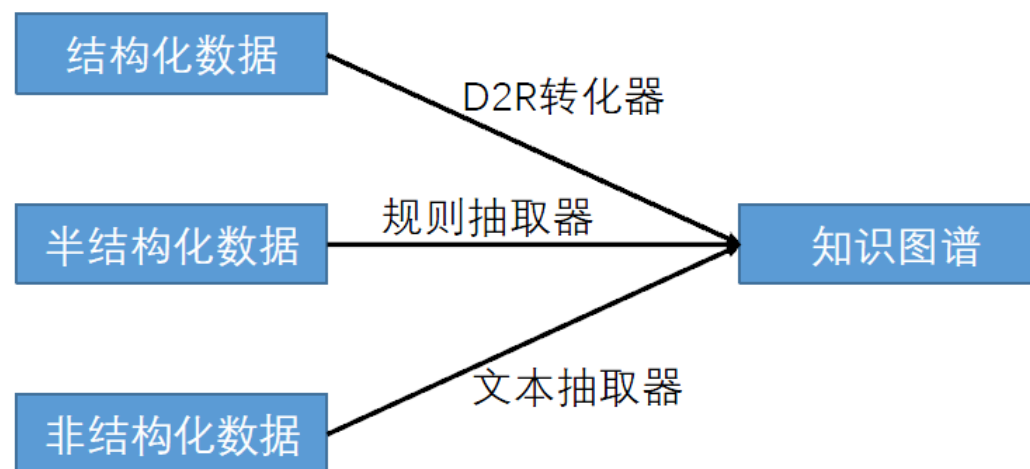
• 1. 模式设计

- 将认知领域的基本框架赋予机器
- 概念层设计
 - 指定领域的基本概念，以及概念之间子类关系
e.g., 足球领域，足球运动员是运动员的子类
- 属性定义
 - 明确领域的基本属性，明确属性的适用概念，属性值的类别或范围
e.g., 效力球队的域为足球运动员，范围为球队
- 约束规则定义
 - 多值属性约束 e.g., 出生日期（单值），获得奖项（多值）
 - 互逆属性约束 e.g., 隶属球员和效力球队为互逆属性

领域知识图谱构建

• 2. 明确数据来源

- 结构化程度较高、质量较好，以尽可能低代价获取数据
 - 互联网上的领域百科爬取
 - 通用百科图谱的导出
 - 内部业务数据的转换
 - 外部业务系统的导入



不同数据来源通过不同的知识获取方式构建知识图谱

领域知识图谱构建

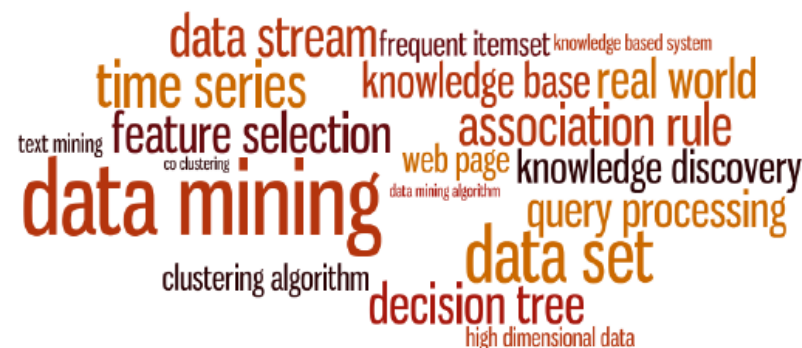
• 3. 词汇挖掘

- 识别出领域中重要短语和词汇
 - 识别领域的高质量词汇
 - 识别同义词
 - 识别缩写词
 - 识别领域常见情感词

Raw Corpus



Quality Phrases

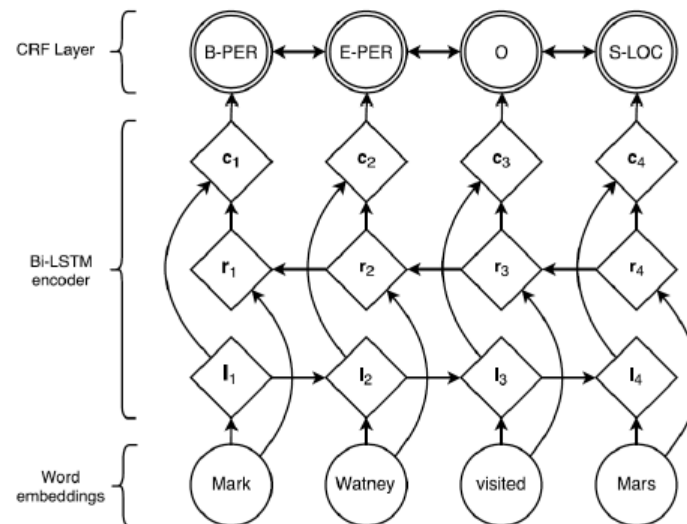


Jiawei Han, etc., Mining Quality Phrases from Massive Text Corpora

领域知识图谱构建

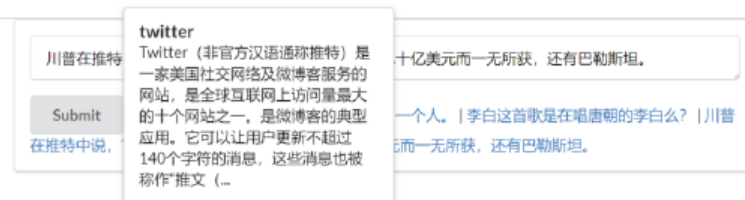
• 4. 实体发现

- 识别出领域中的常见实体
- 理解领域文本和数据的关键一步
 - 实体识别
 - 实体归类
 - 实体链接



Guillaume Lample etc., Neural Architectures for Named Entity Recognition

实体链接



[川普]在[推特]中说，“不仅是[巴基斯坦]让我们支付了几

十亿美元而一无所获，还有[巴勒斯坦]。

知识工场实验室的实体链接DEMO

领域知识图谱构建

• 5. 关系发现

• 填充知识库中的关系实例

- 关系分类：将给定的实体对（entity pairs）分类到某个已知关系

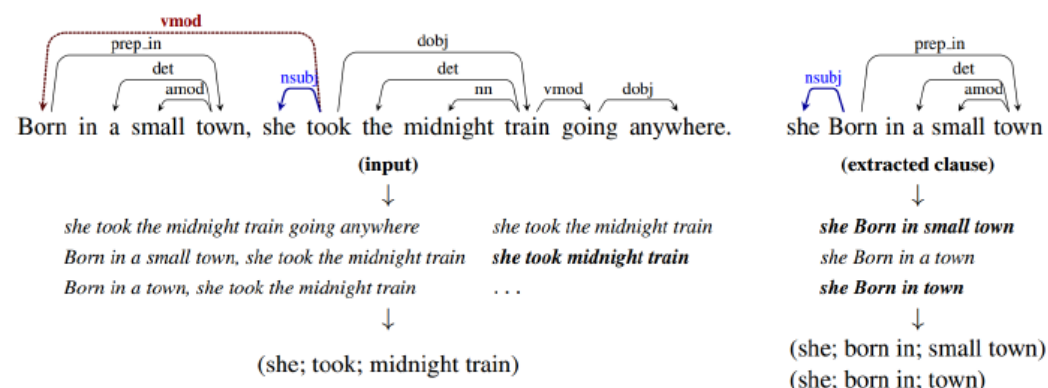
e.g., 李娜-姜山 ->丈夫, 教练

- 关系抽取：从文本中抽取某个实体对的具体关系

e.g., 姜山曾先后两次成为李娜的教练->(李娜, 教练, 姜山)

- 开放关系抽取：从文本中抽取出实体对之间的关系描述

e.g., 上海隔中国东海与日本九州岛相望->(上海, 相望, 日本九州岛)

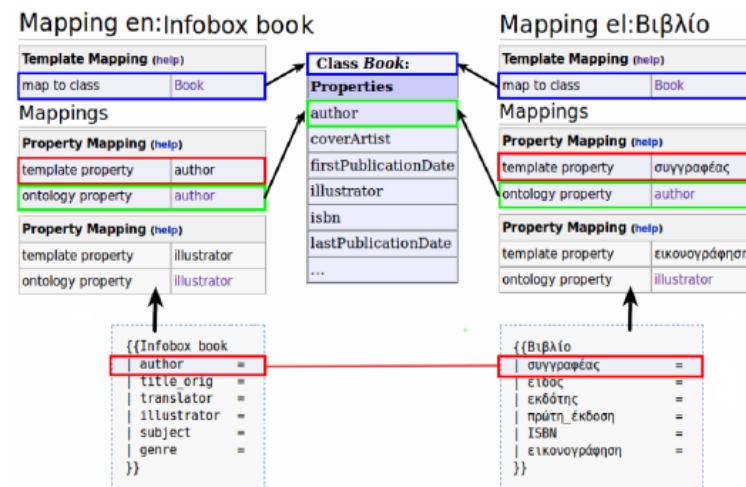


Stanford Open Information Extraction,
<https://nlp.stanford.edu/software/openie.html>

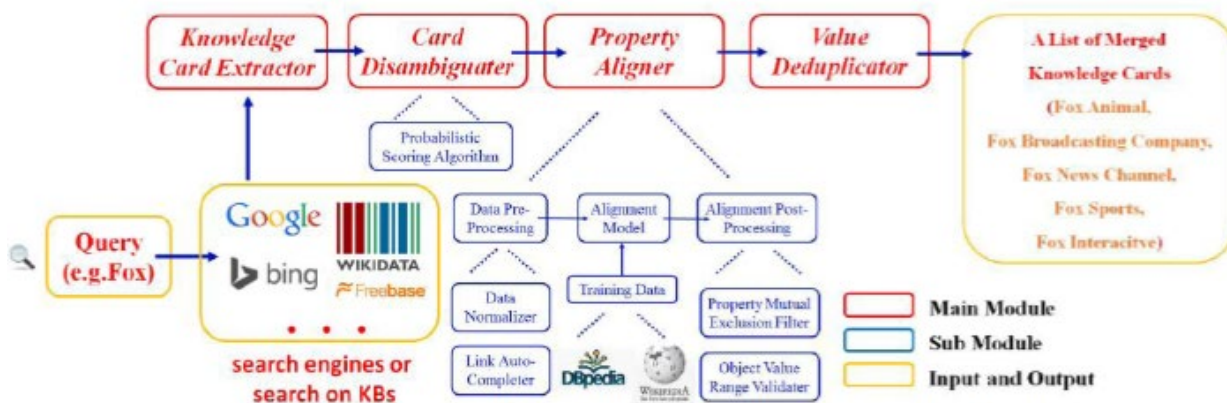
领域知识图谱构建

• 6. 词汇挖掘

- 融合来自不同数据源的知识
 - 实体对齐：识别不同来源的统一实体
e.g., 华中科技大学, HUST
 - 属性融合：识别同一属性的不同描述
e.g., 英文名, 英文名称
 - 值规范化：规范化到统一格式/单位
e.g., 175cm, 1米75



跨语言知识融合



Effective Online Knowledge Graph Fusion

领域知识图谱构建

• 7. 质量控制

- 知识补全

e.g., 如果一个人出生地是中国，推断其国籍也可能是中国

e.g., 从外部互联网文本数据补充知识

- 知识纠错

e.g., 互逆属性纠错：A妻子B，B丈夫C

- 知识更新

• 8. 人工干预

- 人工编辑

- 众包构建

e.g., 利用知识问答验证码来进行知识获取
提升知识图谱的质量



知识工场实验室推出的KADE系统，能够所见即所得的知识图谱编辑

请通过验证

请点击下文中该问题答案的任意部分：毛里西奥·多米齐的出生地在哪里？
太难了，换一个
毛里西奥·多米齐，男，1980年6月28日出生于意大利罗马，是一名出色的足球运动员，曾以后卫效力于拿波里足球队，现效力于乌甸尼斯足球队。

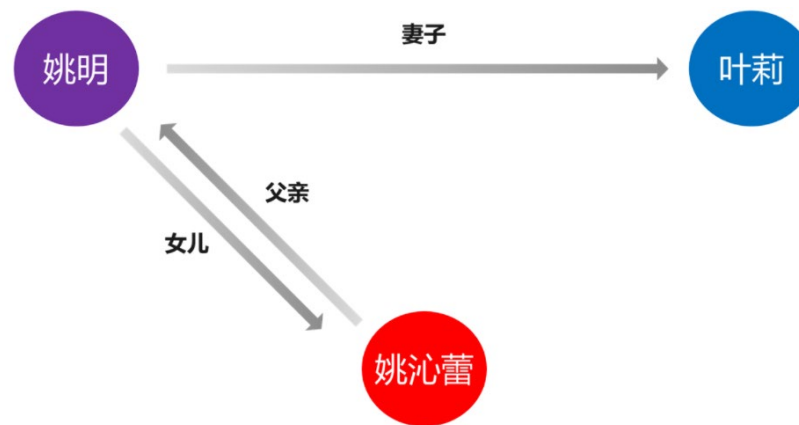
基于文本理解的超级验证码可以实现大规模众包化知识获取

知识图谱的应用



面向知识图谱的推理

- 围绕关系的推理展开
- 已有的事实或关系推断出未知的事实或关系
- 一般着重考察实体、关系和图谱结构三个方面的特征信息
- 辅助推理出新的事实、新的关系、新的公理以及新的规则等

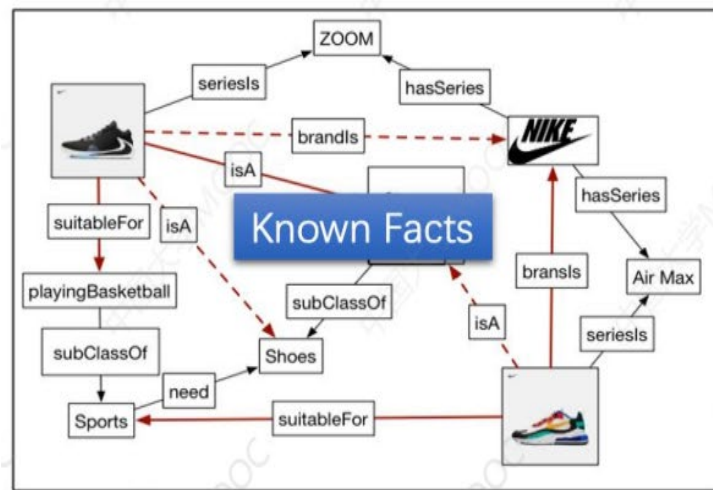


<姚明, 妻子, 叶莉>
<姚明, 女儿, 姚沁蕾>



<叶莉, 女儿, 姚沁蕾>

知识图谱推理的主要任务和作用



Infer

New Facts
New Relations
New Axioms
New Rules

问句扩展

图谱补全

错误检测

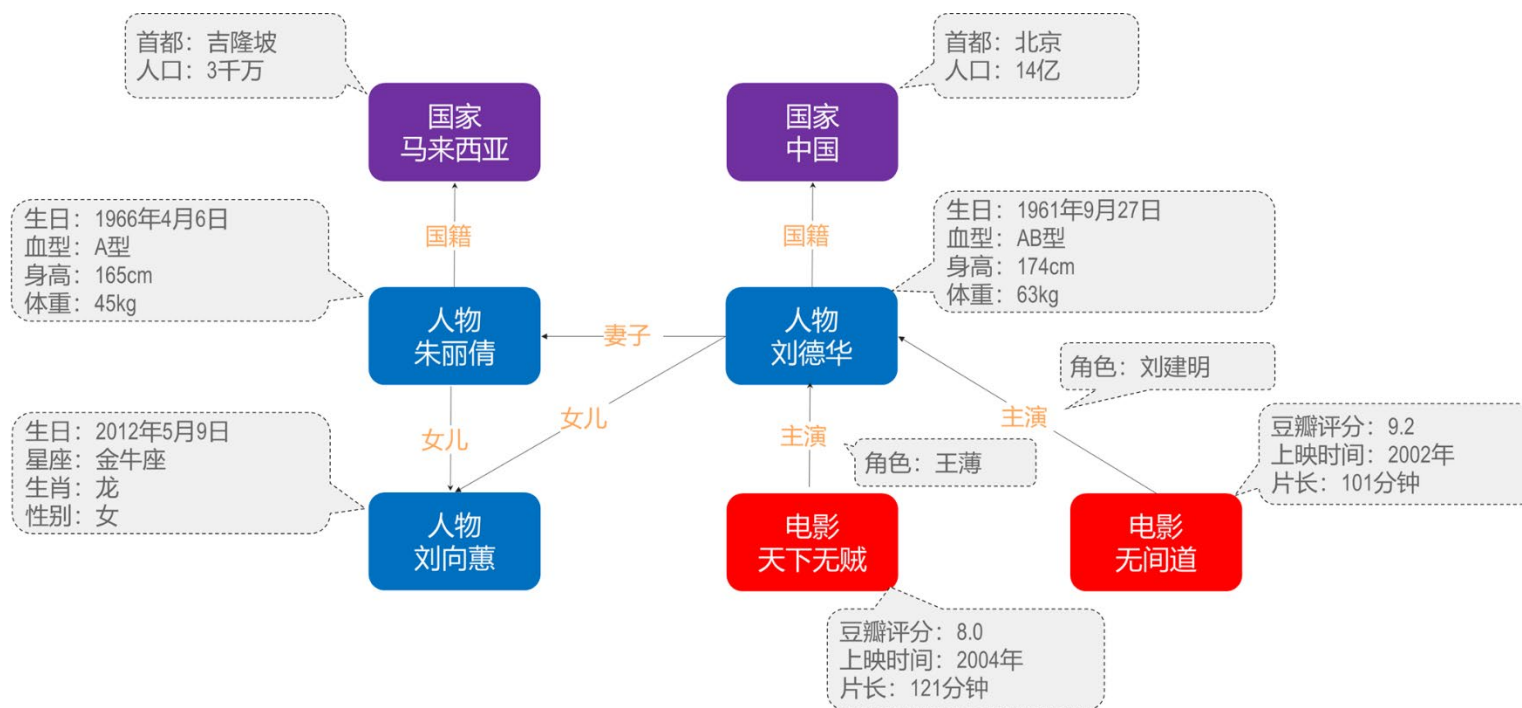
关联关系推理

冲突检测

知识图谱推理的主要任务和作用

• 问句扩展

刘德华主演的电影中豆瓣评分大于8分的有哪些？



知识图谱推理的主要任务和作用

- 图谱补全

- 知识图谱补全 (Knowledge Graph Completion, KGC) 目前主要被抽象成一个预测问题，即预测出三元组中缺失的部分。所以可分成3个子任务：

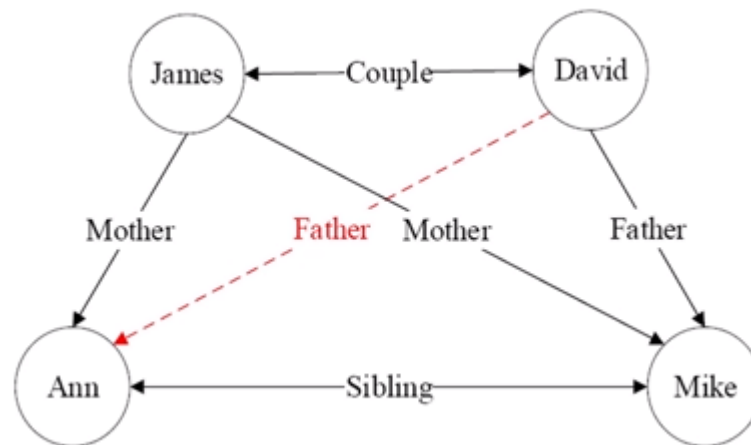
- 头实体预测： $(?, r, t)$  实体预测 (Entity Prediction)
- 关系预测： $(h, ?, t)$  链路预测 (Link Prediction)
- 尾实体预测： $(h, r, ?)$  实体预测 (Entity Prediction)

- 问号表示要预测的部分，而另外两部分是已知的。同时一般按照能否处理新实体或者新关系，可以将知识图谱补全算法分成两类：**静态知识图谱补全**和**动态知识图谱补全**。

Static KGC

- 静态知识图谱补全（Static KGC），该场景的作用是补全已知实体之间的隐含关系。仅能处理实体以及关系都是固定的场景，所以扩展性较差。

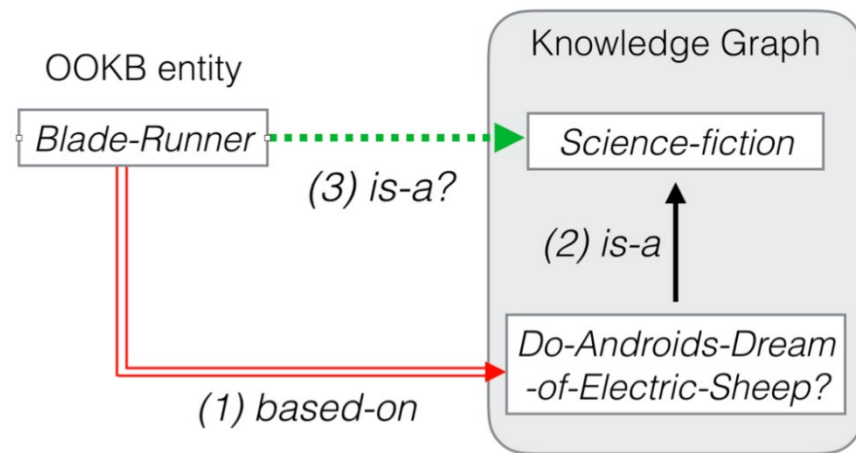
Ann和David的关系是什么？



$$(\forall x)(\forall y)(\forall z)(\text{Mother}(z, y) \wedge \text{Couple}(x, z) \rightarrow \text{Father}(x, y))$$

Dynamic KGC

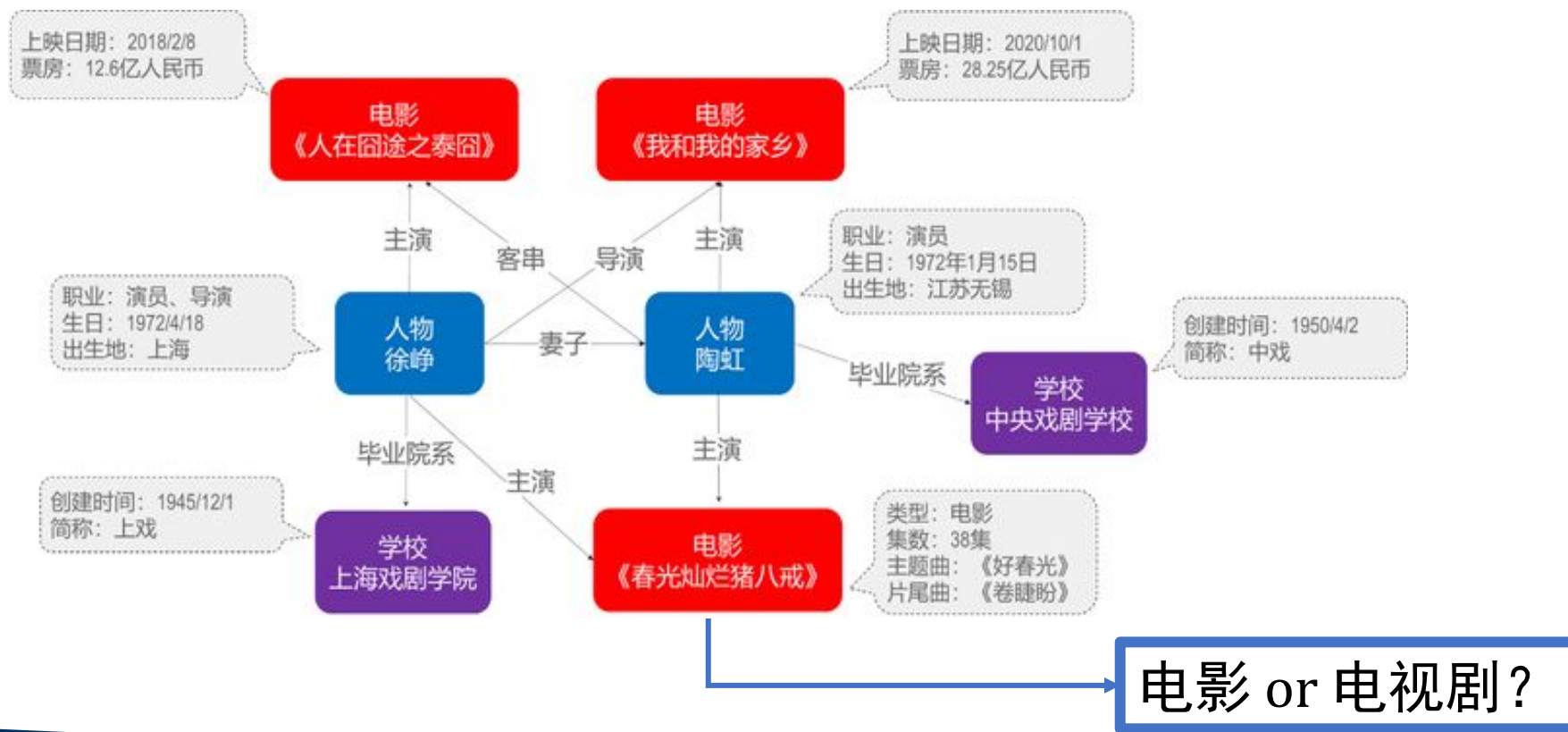
- 动态知识图谱补全（Dynamic KGC），它涉及不止有知识图谱中已有的实体或关系，还涉及一些没有出现的词，或者后期想对知识图谱进行补全的场景下。
- 对于新出现的实体，可以基于该实体与训练好的实体之间的邻居连接状态得到新实体的表示，从而预测新实体与其他实体之间是否存在某关系。



知识图谱推理的主要任务和作用

可以通过推理进行知识图谱纠错。

- 错误检测



知识图谱推理的主要任务和作用

- 关联关系推理
- 虽然目前的知识图谱上已经有了非常多的实体对和关系事实，但是由于数据的更新迭代以及不完整性，注定了知识图谱的不完整，也隐藏着我们难以轻易发现的信息。

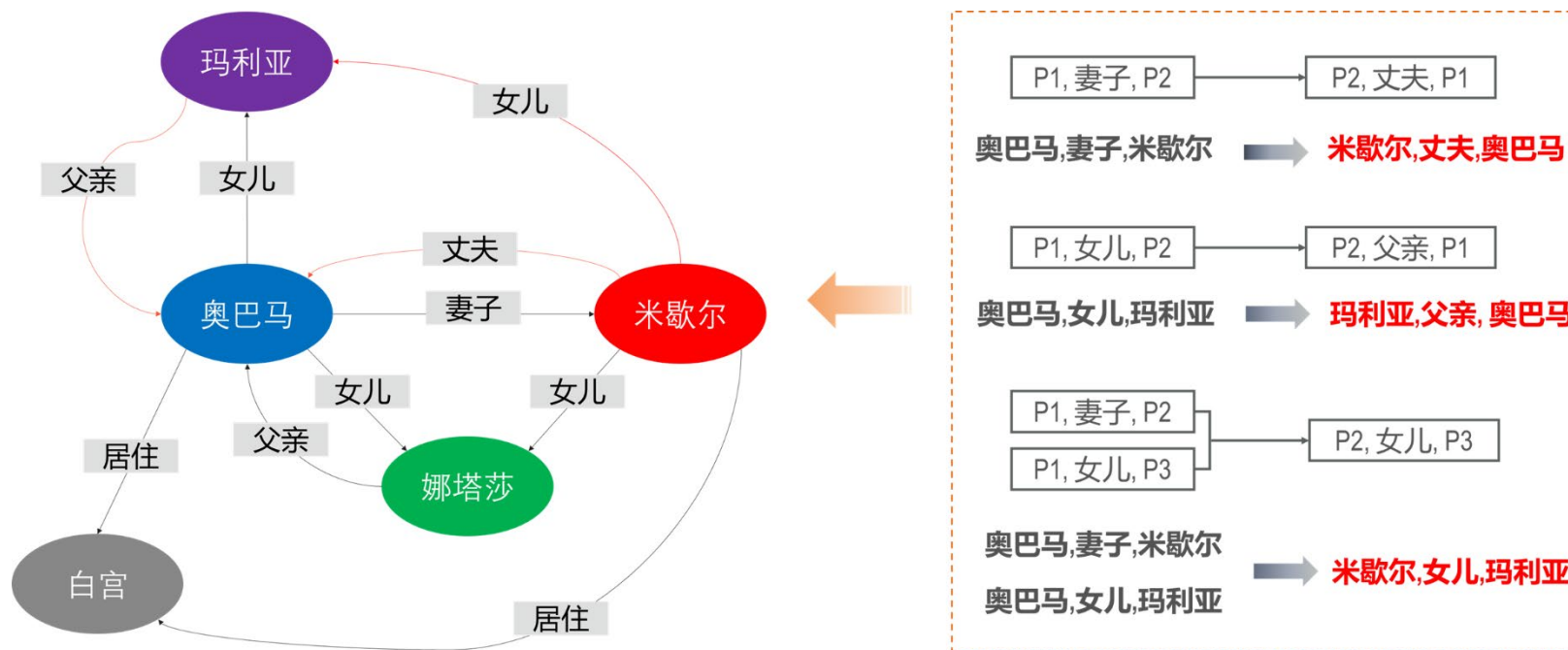
Melinda-spouse-Bill-chairman-Microsoft-HQ-in-Seattle.



Melinda-lives-in- Seattle

基于规则的推理

- 基于规则的推理通过定义或学习知识中存在的规则进行挖掘与推理



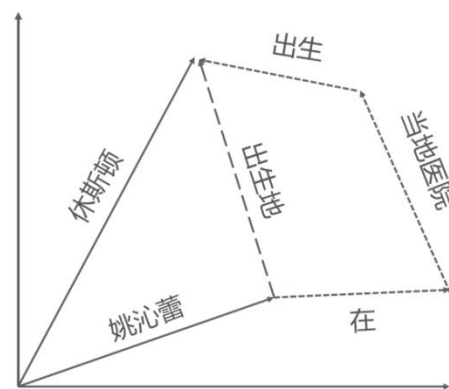
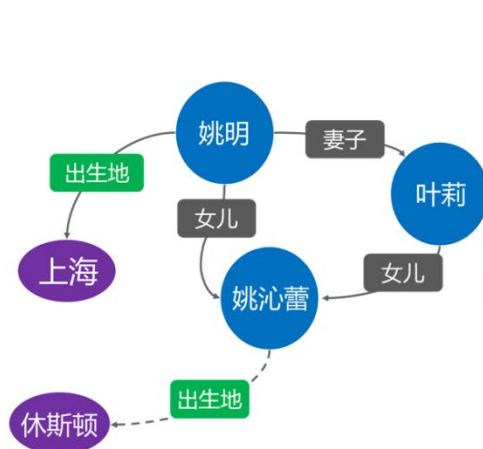
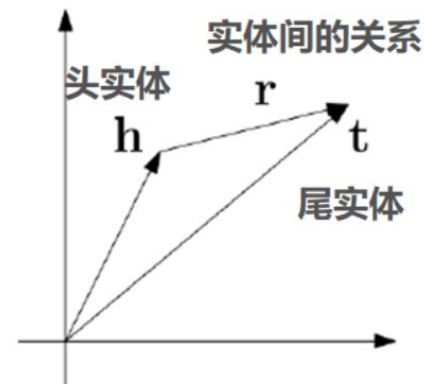
基于本体的推理

- RDFS（一类满足特定规范的用来表述本体的语言）定义了一组用于资源描述的词汇：包括class, domain, range等。其本身就蕴含了简单的语义和逻辑。我们可以利用这些语义和逻辑进行推理。



基于表示学习的推理

“姚沁蕾的出生地是哪儿”



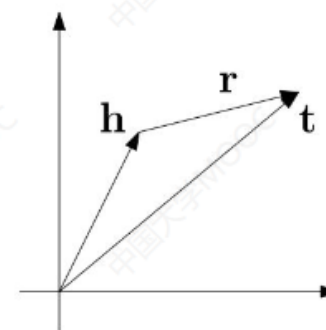
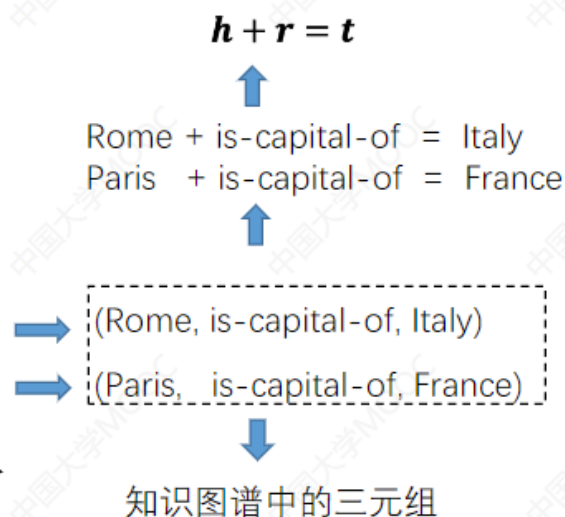
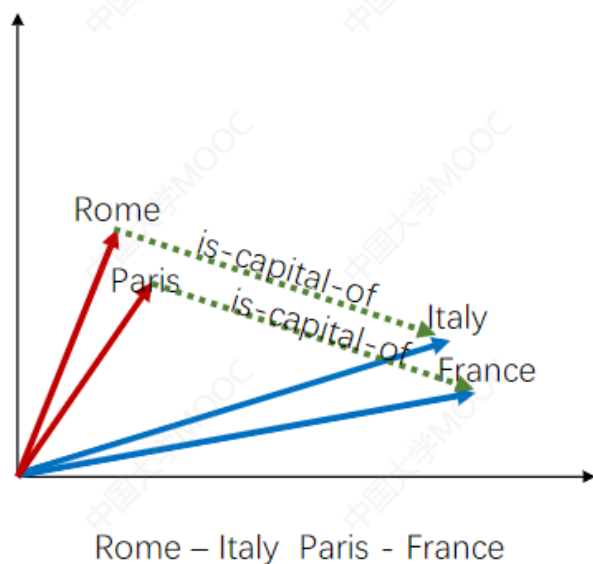
向量表示空间

姚明出生在一个上海篮球世家，父亲和母亲都曾职业篮球运动员。2007年8月6日，姚明与相恋七年的女友叶莉在上海举行婚礼。婚后叶莉随姚明赴美生活，2010年5月21日，叶莉同姚明的女儿在休斯敦当地医院出生。次年姚明在上海公布了女儿的名字：姚沁蕾。

基于算法的推理-例：TransE

- 受词向量在向量空间语义层面的启发
- 转移距离模型的主要思想：知识图谱中向量化后的三元组是否合理→衡量头实体和尾实体之间的距离

$$\text{head} + \text{relation} \approx \text{tail}$$



$$f_r(h, t) = \|h + r - t\|_{L_1/L_2}$$

TransE的训练

$$d(sub, rel, obj) = \|\mathbf{s} + \mathbf{r} - \mathbf{o}\|_2^2$$

To learn such embeddings, we minimize a margin-based ranking criterion over the training set:

$$\mathcal{L} = \sum_{(h, \ell, t) \in S} \sum_{(h', \ell, t') \in S'_{(h, \ell, t)}} [\gamma + d(\mathbf{h} + \mathbf{\ell}, \mathbf{t}) - d(\mathbf{h}' + \mathbf{\ell}, \mathbf{t}')]_+$$

smaller

larger

where $[x]_+$ denotes the positive part of x , $\gamma > 0$ is a margin hyperparameter, and

$$S'_{(h, \ell, t)} = \{(h', \ell, t) \mid h' \in E\} \cup \{(h, \ell, t') \mid t' \in E\}$$



Randomly change an entity

不确定知识图谱推理

- 不确定知识图谱：为每个三元组添加一个置信度来描述三元组的不确定性。
- 将通常知识图谱中三元组(h, r, t)拓展为 $\langle (h, r, t), s \rangle$ ，其中h、t代表头实体尾实体、r代表头尾实体之间的关系，s代表置信度：

1. (choir, relatedto, sing): 1.00
2. (college, synonym, university): 0.99
3. (university, synonym, institute): 0.86
4. (fork, atlocation, kitchen): 0.4

三元组的置信度是如何得到的呢？
不同的知识库计算置信度的策略有所不同：

- 根据众包标注频率计算得到；
- 通过统计三元组的上下文数量计算置信度。

GTransE

- 专注于学习那些置信度更高的三元组，降低那些质量较差、置信度较低的三元组对实体及关系表示的贡献；
- 利用三元组的置信度动态地调整TransE中的Margin；
- 置信度高间隔更大，置信度小间隔更小。

$$L = \sum_{(h,r,t,s) \in Q} \sum_{(h',r,t',s) \in Q'} [f(h,r,t) - f(h',r,t') - s^\alpha M]_+,$$

