

新冠开放知识图谱

020.11.15

2020.11.15 CCKS工业论坛



Z \dashv Ш Z

- 新冠开放知识图谱构建及关键技术
- 跨数据集关联与融合
- 数据规范与统一访问接口
- 新冠开放知识图谱潜在应用与发展方向
- 新冠开放知识图谱相关数据竞赛

▶疫情袭来,信息过载



截至9月5日24时,全球新冠肺炎确诊病例累计达**26,953,120**例,累计死亡**878,256**例。疫情态势依然严峻复杂。

每天大量关于新冠肺炎的信息分布在各种媒体网站、研究刊物、官方文件等,需要将过载的信息整合,提高信息利用价值,有效助力抗疫行动。





AI 正在向"认知智能"演进



国务院新一代人工智能发展规划

2017年7月8日 国务院明确提出了

建立新一代人工智能关键共性技术体系重点任务,

特别强调了研究跨媒体统一表征、关联理解与知识挖掘、知识图谱构建与学习、知识演化与推理等技术。

新一代人工智能重大项目2020年度第一批项目申报指南

"**认知**"出现 16 次,"知识图谱"出现 13 次 鼓励在金融、客服、教育、工业、医疗等领域构建行业知识图谱。













新冠开放知识图谱构建及关键技术







我想对新型冠状病毒和新冠肺炎有一个基本的认识...<百科>



得了新冠肺炎会有哪些症状,现在都怎么治疗的...<\(\text{\mathbb{k}}\) >
<a href="http



个人防护、场所防控都权威的指导 手册...<^{防控>}



疫情中,物资现在是什么状态...

<物资>



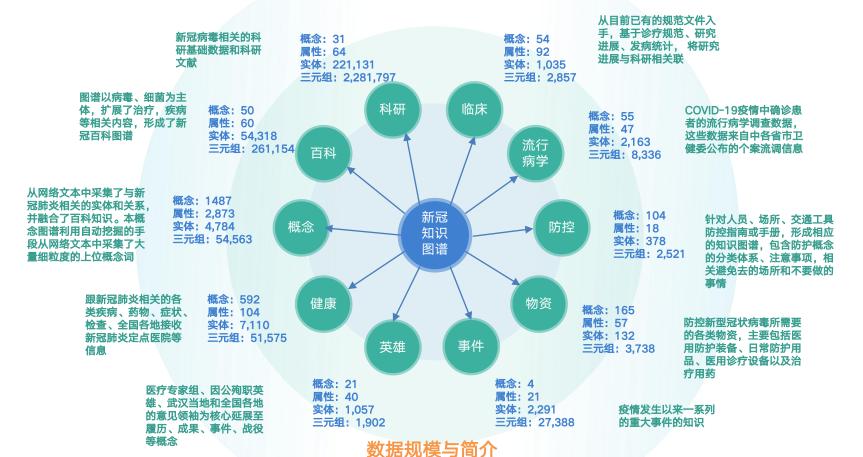


新冠开放 知识图谱

- 5 现在有没有研发出疫苗,有没有特效药呢...<科研>
- 有没有技术可以追溯传染源, 探索被接触者…<流行病学>
- び场"战疫"中,有哪些伟大的英雄...<英雄>
- 有关疫情的热点事件有哪些... <_{事件}>

新冠开放知识图谱助力疫情防控





新冠开放知识图谱助力疫情防控



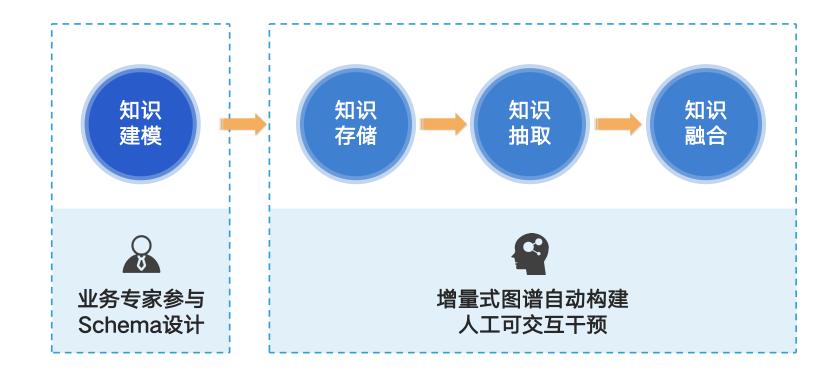
预测新病毒的生物学分类,新物种发 《新型冠状病毒感染的肺 基干该图谱讲行知识问答 现、变异性、热稳定性、易感群体、 NCBI(美国国家生物 炎诊疗方案》、Wiki百科、 宿主、致病部位、可导致的症状,可 技术信息中心网站) 中医药知识服务平台、 缓解症状的药物、潜在治疗的药物. 医疗器械分类目 如老药新用、传播途径、传播种类 临床 基于流行病学知识图谱的 科研 百度百科、 面向新冠相关术语的语义检索、 分析推理技术为医护人员 《流行病学-第7 智能问答,并可用于新冠相关 互动百科、 和疾病防控人员提供技术 版》、《流行病学复 中文维基百科 文档的智能搜索和推荐 流行 支撑,加速流调研究 习考试指导》、各地 病学 卫健委公开信息 已应用于深睿医疗开发的 人民日报、丁香 提供基本防控知识问题, 集成 新冠 新冠肺炎小睿医生助手中 新型冠状病毒肺炎 医生、腾讯、新 概念 防控 于流程化信息处理平台, 用于 用于计算问句之间的相似 知识 的各种防护手册 浪微博 各场所检查防控措施是否正确 度以及辅助解答用户提问 图谱 《新型冠状病毒感染的肺 物资 健康 《新型冠状病毒感染的肺 炎诊疗方案》、百度百科、 炎诊疗方案》、《国家基 多轮人丁智能问答 流行病调查研究,基于 卫生健康委、北京妙医佳 本药物目录》等 图谱的新冠肺炎健康防 健康科技集团有限公司 事件 英雄 护间答 百度百科、微信 人民日报、丁香 公众号、知网、 支持对新型冠状病毒的事件在时间 医生、腾讯、新 澎湃新闻 上的正向和反向索引。并提供系列 浪微博 基于图谱进行英雄 事件发展脉络的枚举。支持热点事 人物信息动态展示 件的查证溯源。和区块链技术结合

来源与潜在应用

可具备对事件的存真鉴伪的功能

▮新冠开放知识图谱─构建流程

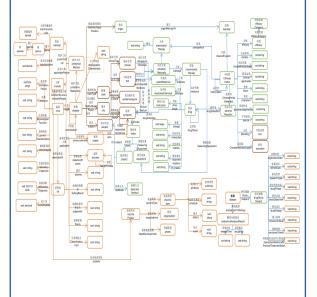






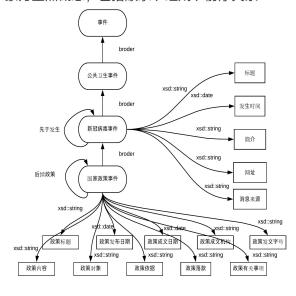
新冠临床图谱

以新冠肺炎疾病的患者病例、治疗方案、药物 研究、中西药物作用等为重点



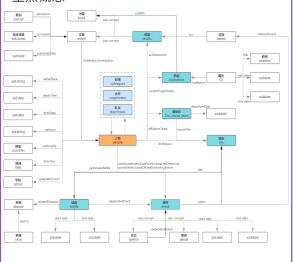
新冠事件图谱

以政策、发文机构、发布日期、依据、标题、对 象为重点概念,包括顺承、组成、前序关系



新冠英雄图谱

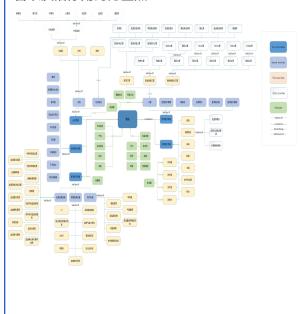
以人物、战役、事件、成果、文章、履历等为 重点概念





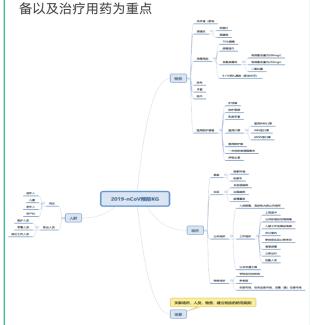
新冠物资图谱

以医用防护装备、日常防护用品、医用诊疗设 备以及治疗用药为重点



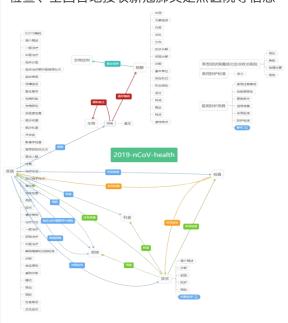
新冠防控图谱

以医用防护装备、日常防护用品、医用诊疗设 备以及治疗用药为重点



新冠健康图谱

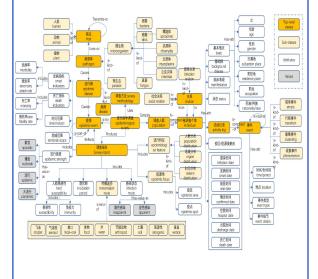
包含和新冠肺炎相关的各类疾病、药物、症状、检查、全国各地接收新冠肺炎定点医院等信息





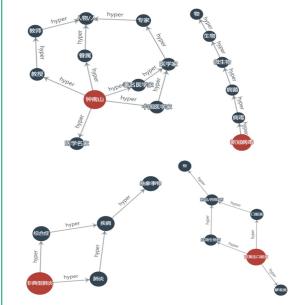
新冠流行病学图谱

重点刻画流行病学的基本概念、流行病学调查 等内容



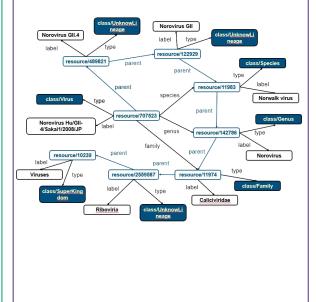
新冠概念图谱

包括括疾病、人物、症状等实体和关系,以及大量细粒度的上位概念

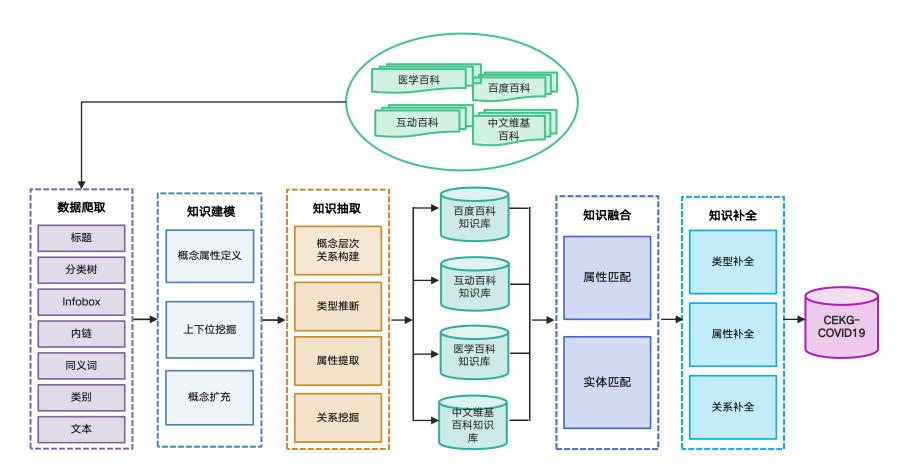


新冠科研图谱

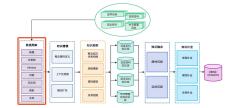
以病毒分类、新冠基本信息、抗病毒药物、病毒亲缘关系为重点



新冠开放知识图谱一构建流程(以新冠百科图谱构建过程为例)



新冠百科知识图谱—数据爬取





1 内部链接(InnerLink)

在本病的诊断中病史极为重要,往往曾反复发生。 <u>鱼家炎、支气管炎</u>或肺炎,或曾患麻疹、百日咳、流行性感冒或腺病毒肺炎。确定诊断需要结合病史、症状和丝检查。

② 分类树

百科分类材 知识地图 搜分类 病毒 - 病毒 收起 - DNA病毒 收起 + 双链DNA病毒 展开 + 单链DNA病毒 展开 - RNA病毒 收起 + 双链RNA病毒 展开 + 正链RNA病毒 展开 + 负链RNA病毒 展开 + 负链RNA病毒 展开 + 负链RNA病毒 展开 + 逆转录病毒 展开 + 逆转录病毒 展开

3 类别

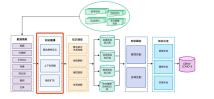
异烟肼磺酸钠

开放分类: 医学 | 处方药 | 药品 | 西药

异烟肼磺酸钠,别名 甲磺烟肼钠,异烟肼磺酸钠,异烟肼甲烷磺酸钠,用于各型结核病的治疗。口服或肌注: 0.1~0.3g/次,2~3次/日。

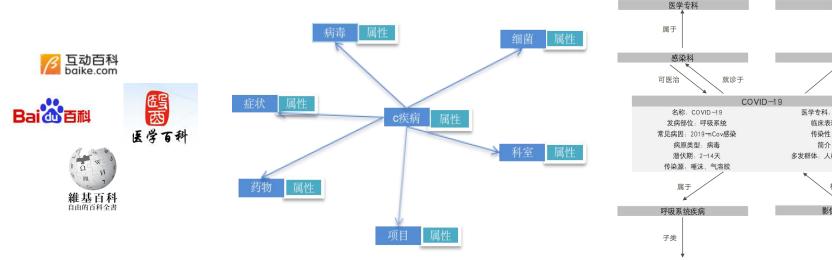
辆铒那

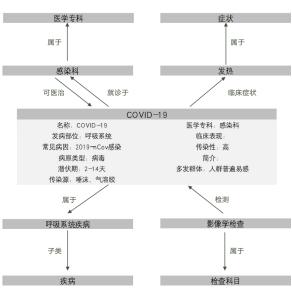
▋新冠百科知识图谱─知识建模





从COVID-19的基本信息、症状、致病原因、治疗药物、检查项目、诊治科室六个切入点,选取 疾病、症状、细菌等七个顶级概念,获取四大百科category system中这七大概念的所有子概念, 上下位关系取四大百科的交集部分,对齐BabelNet,并以此扩充概念,构建新冠百科知识图谱 Schema.





新冠百科知识图谱一知识抽取



OpenKG

类型推断(Type Inference)

- 句法分析
- Infobox 模板
- 外部知识库
- 观察特征



▮新冠百科知识图谱─知识抽取





属性与属性值提取





非结构化数据

- 1.网页抽取特定属性和属性值同时出现的句子,"该疾病的**传播途径**包括 近距离飞沫传播、呼吸道分泌物传播 等。"
- 2.把1.中抽取的句子作为正例,其他作为负例,训练最大熵模型
- 3.通过模型判断出潜在正例句子,使用CRF标注属性值,如"COVID主要传播途径有<u>唾液</u>,<u>呼吸道</u>,<u>密切接</u>触。"

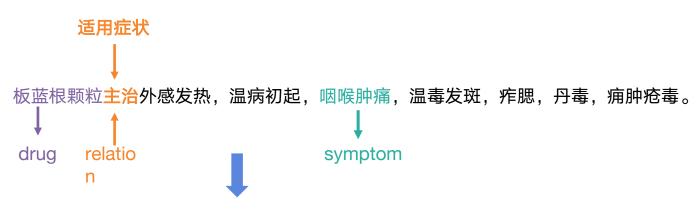
▮新冠百科知识图谱─知识抽取





关系提取

- 1. 抽取潜在relation语义的句子(Regex Matchig)
- 2. Bert+BiLSTM+CRF进行实体识别
- 3. 三元组构建



<板蓝根颗粒,适用症状,咽喉肿痛>

▮新冠百科知识图谱─知识融合



实例: 重症急性呼吸综合征



属性融合

- 属性的同义词库
- 谓语频数统计与归纳

实体融合

● 挖掘等价规则

首先从种子实例中挖掘出等价属性,如"别名"和"别称"、"常见症状"和"临床表现"等是等价属性,然后根据已知的等价实体的等价属性的出现频率情况,挖掘匹配规则。

百度百科

<mark>別称</mark>: 严重急性呼吸综合征、SARS、非典、非典型肺炎 英文别名;;;;||;;;severe acute respiratory syndrome,

SARS

就诊科室;;;||;;;感染科

常见病因;;;;ll;;;;SARS冠状病毒感染

常见症状;;;;|1;;;;发热、头痛、肌肉酸痛、呼吸衰竭

传染性;;;;||;;;;有

传播途径;;;;ll;;;;近距离飞沫传播或接触患者呼吸道分泌物

传播

互动百科

中文名: 重症急性呼吸综合征

外文名: severeacuterespiratorysyndrome 别名: 严重急性呼吸综合征、SARS、非典、

季节分布:四季传染病:是

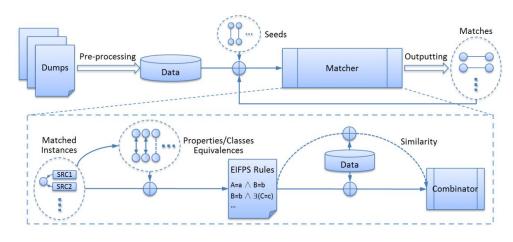
传播途径: 近距离飞沫传播或接触患者呼吸道分泌物传播

临床表现:发热、头痛、肌肉酸痛、呼吸衰竭

就诊科室: 传染科

常见病因: SARS冠状病毒感染

英文别名: severeacuterespiratorysyndrome, SARS



新冠百科知识图谱一知识补全





属性值补全

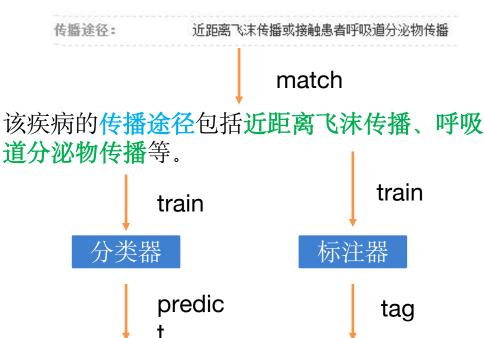
- 从Infobox中抽取特定属性与属性值
- 在页面中抽取特定属性和属性值同时 出现的句子
- 这些句子作为正例,其他句子作为负 例训练分类器,以预测某一句子是否 可能包含特定属性的属性值。
- 对于预测包含属性值的句子,使用标 注方法对属性值进行预测

类型补全

特定关系所链接的两个实体类型可预测。

关系补全

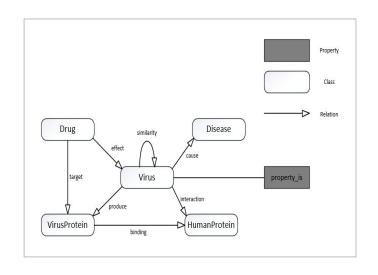
基于title和inner link(内链)之间的context。



COVID-19主要传播途径有唾液,呼吸道,密切接

其他典型图谱知识抽取举例(科研图谱)

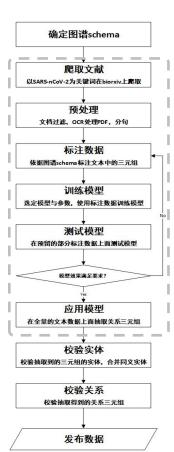




数据源是期刊论文,从病毒及与病毒相关的多种类型生物分子/化学分子出发,梳理其关系,构建新冠科研-基本信息图谱。

```
"text": "Binding interactions of 2019-nCoV S-RBD
with human ACE2 With the complex model in Fig.",
"spo_list": [
   "predicate": "binding",
   "subject_type": "VirusProtein",
   "object_type": "HumanProtein",
   "subject": "2019-nCoV S-RBD".
   "subject position": 24,
   "object": "ACE2",
   "object_position": 51
```

标注数据样例

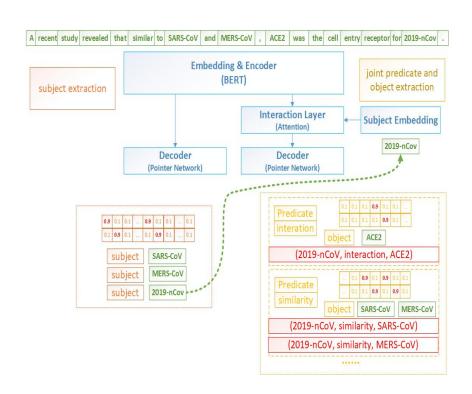


其他典型图谱知识抽取举例(科研图谱)



基于《A Unified MRC Framework for Named Entity Recognition》,提出了一个从医疗文献中抽取实体关系的模型。这个模型同样由两部分组成: 主语抽取部分(subject extract)和谓词宾语联合抽取部分(joint predicate and object extraction),如右图所示。

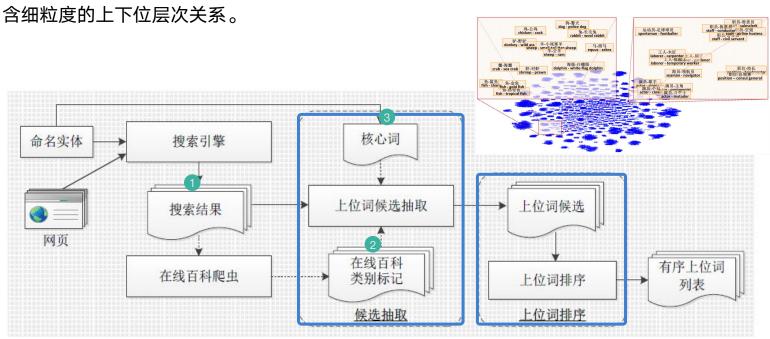
- 在主语抽取部分,模型将会抽取能够作为三元组主语的 实体,输出部分是一个双层指针网络,标识出实体的头 尾位置;
- 在谓词宾语联合抽取部分,模型将基于前一部分抽取到的主语,抽取出对应的宾语并判断关系。一个注意力层,一边的输入是句子的BERT输出,另一半的输入是这个主语在句子中的BERT输出。最后输出层是2*N层的指针网络(N是所需抽取的关系数量)。每一类关系与2层指针网络相对应,标识出宾语的头尾位置。



| 其他典型图谱知识抽取举例(概念图谱)



从网络文本中自动挖掘大量细粒度的概念词,自动构建概念的is-a层次结构(schema),富







对新冠开放知识图谱每个数据集进行了数据质量评测,针对每种关系随机采样一定数目的三元组,由对数据熟悉的人员或者专业人员评测。为了将采样评测结果扩展到整个数据集,我们计算了置信度为95%的威尔逊区间,评测结果如下:

图谱名称	数据规模	采样比例	采样数目	正确数目	准确率预估
百科	261,154	1.91%	5,000	4,778	95.52%±0.58%
临床	2,857	22.82%	652	620	94.81%±1.71%
科研	2,281,797	0.37%	8556	8555	99.96%±0.03%
事件	27,388	0.73%	200	198	96.35%±1.69%
英雄	1,902	29.97%	570	570	99.65%±0.35%
防控	2,521	17.29%	436	434	99.09%±0.79%
物资	3,738	9.76%	365	359	97.83%±1.42%
健康	51,575	0.94%	487	483	98.78%±0.91%
概念	54,563	1.20%	100	96	92.31%±4.96%
流行病学	8,336	2.4%	200	200	98.08%±1.92%





跨数据集关联与融合



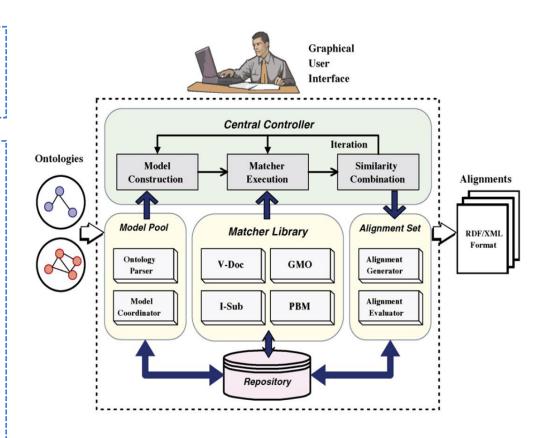
▮新冠开放知识图谱─本体匹配



由于没有中心本体,我们对各个图谱的本体成对匹配。采用了南京大学研发的Falcon-AO工具,利用文本和结构特征进行自动化本体匹配

匹配器:

- V-Doc: 采用了语言学的方法来做本体匹配, 其最大的亮点在于, 构建了虚拟文档
- I-Sub: 将字符串的共性和差异结合考虑,使用字符串距 离的度量方法比较本体的类名和属性名
- GMO是一种迭代式的结构的匹配方法,它使用RDF二部图 来表示本体,计算类或属性以及三元组之间的结构相似度
- PBM使用分而治之的方法来寻找大规模本体之间的块映射



新冠开放知识图谱一本体匹配结果



表中结果斜线左侧是类的个数,右侧是属性的个数,结果显示各个新冠知识图谱本体有重叠,但是重叠数量较少。

类/属性	百科	防控	概念	健康	科研	临床	流行病	事件	物资	英雄
百科		0/0	37/0	9/9	2/4	4/6	4/1	0/0	0/0	0/0
防控			7/0	0/0	0/0	2/2	0/0	5/16	1/17	2/1
概念				41/0	6/0	25/ 0	18/0	1/0	18/0	11/0
健康					2/3	4/13	3/3	0/1	6/1	0/0
科研						4/3	0/0	0/1	0/0	5/0
临床							3/6	0/3	0/0	5/1
流行病								1/0	0/0	3/4
事件									2/16	1/0
物资										0/0
英雄										

新冠开放知识图谱一实体对齐



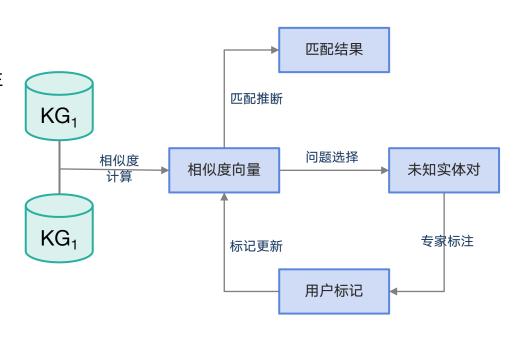
基于实体对相似度向量的主动学习方法

动机

- 实例数据属性以Datatype Property为主
- 缺少训练数据

方法

- 基于属性对齐的相似度计算
- 基于偏序的问题选择
- 基于偏序的匹配推断



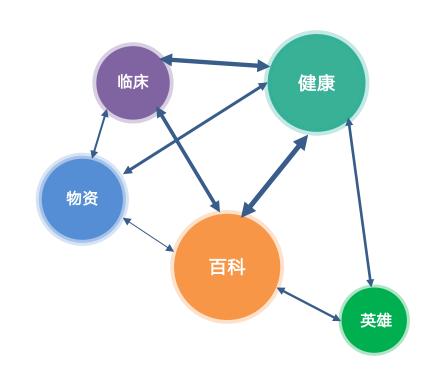
新冠开放知识图谱一实体对齐结果



实体对齐结果:

- 总计1055个匹配
- 绝大多数匹配的类是药物
- 百科数据集自身存在匹配(已有owl:sameAS)
- 百科和其他数据集之间的匹配较多
- 百科和物资图谱匹配较少
- 健康和英雄图谱存在匹配(医疗方向)

	健康	临床	英雄	物资
百科	836	55	11	2
健康		95	19	28
临床				9







数据规范与质量评估



新冠开放知识图谱一命名规范



资源URI格式:

- 概念: http://{namespace}/{graphname}/class/{localname}
- 实体: http://{namespace}/{graphname}/resource/{localname}
- 属性: http://{namespace}/{graphname}/property/{localname}
- Namespace统一采用: http://www.openkg.cn/2019-nCoV/

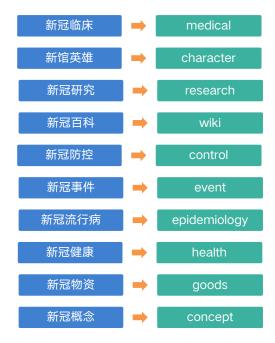
LocalName采用类Wikidata编码方式:

- 概念: C + ID (Integer)
- 实体: R + ID (Integer)
- 属性: P + ID (Integer)
- 概念: http://www.openkg.cn/2019-nCoV/character/class/C10
- 实体: http://www.openkg.cn/2019-nCoV/character/resource/R101
- 属性: http://www.openkg.cn/2019-nCoV/character/property/P5

RDF文件格式:

支持RDF/XML、N-TRIPLE、TURTLE、N3、 JSON-LD

新冠开放知识图谱命名:



新冠开放知识图谱一三元组声明规范



01 声明声称

谓语: http://www.w3.org/2000/01/rdf-schema#label

4 声明实体

谓语: http://www.w3.org/1999/02/22-rdf-syntax-ns#type

07 声明单实体同义词

谓语: http://www.openkg.cn/2019nCoV/ontology/alias 02 声明概念

谓语: http://www.w3.org/1999/02/22-rdf-syntax-ns#type

宾语: http://www.w3.org/2000/01/rdf-

schema#Class

05 声明属性

谓语: http://www.w3.org/1999/02/22-rdf-syntax-ns#type

08 声明属性domain

谓语: http://www.w3.org/2000/01/rdf-schema#domain

3 概念继承

谓语: http://www.w3.org/2000/01/rdf-schema#subClassOf

06 声明实体间同义

谓语: http://www.w3.org/2002/07/owl#sameAs

09 声明属性range

谓语: http://www.w3.org/2000/01/rdf-schema#range

▮新冠开放知识图谱─三元组声明规范



对于对象属性,宾语使用概念URI 对于数值属性,宾语使用如下:

文本	http://www.w3.org/2001/XMLSchema#string	•
整数	http://www.w3.org/2001/XMLSchema#integer	
浮点数	http://www.w3.org/2001/XMLSchema#double	
日期	http://www.w3.org/2001/XMLSchema#date	
时间	http://www.w3.org/2001/XMLSchema#time	Ø
日期时间	http://www.w3.org/2001/XMLSchema#datetime	•





新冠开放知识图谱统一访问接口



▶新冠开放知识图谱─OpenKG



OpenKG 发布的所有新冠知识图谱采用 CC-by SA 相似署名开放许可协议,供大家免费下载使用。

欢迎大家访问新冠图谱专题链接,获取更多新冠知识图谱:

http://openkg.cn/group/coronavirus

同时可以通过访问下面链接,得到新冠图谱数据规范以及任何新冠图谱schema:

http://openkg.cn/dataset/covid-19-schema



▮新冠开放知识图谱─OpenBase



在OpenBase可以搜索任何粒度的实体, 搜索结果会展示出该实体的属性和关系

欢迎大家使用OpenBase:

http://openbase.openkg.cn/





新冠开放知识图谱一OpenBase

OpenKG

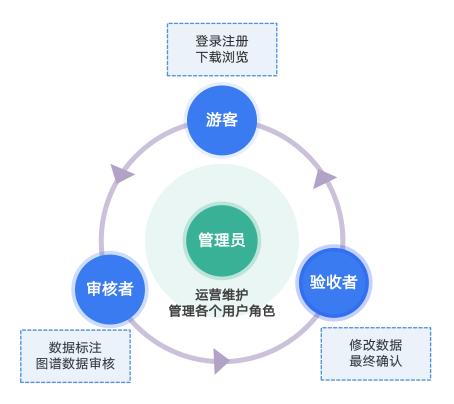
新冠的数据还接入了OpenBase的审核功能,你可以联系我们成为OpenBase的数据审核员,可以在Openbase的网站上查看数据的审核,我们有关系类审核和实体审核两种众包任务。



新冠开放知识图谱一OpenBase



OpenBase整体的众包,我们会抽样选择出新冠数据可能不准确的部分,然后切割成最简单的任务供志愿者来一起审核,分为游客、审核者、验收者、管理员四种角色。





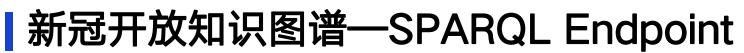


为了方便大家碎片化时间也可以做众包任务,我们还做了微信小程序版本,在这里大家可以随时在手机上进行标注。







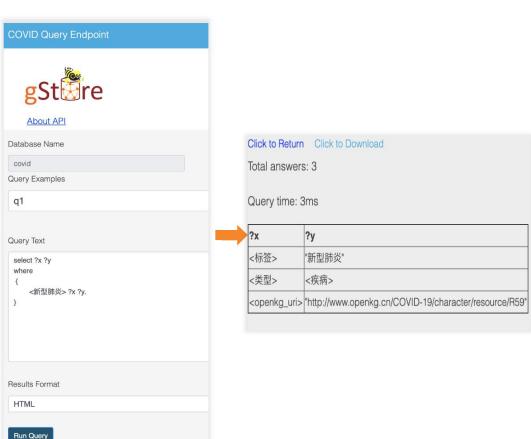




SPARQL Endpoint提供在线查询,同时也可以点击"Click to Download"下载json格式的查询结果。

访问链接:

http://covid.gstore-pku.com/







新冠知识图谱潜在应用与发展方向



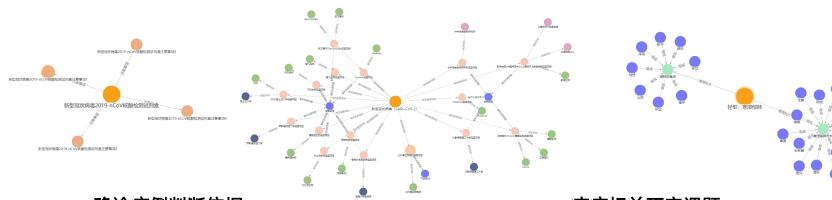


新冠开放知识图谱一主题库可视化



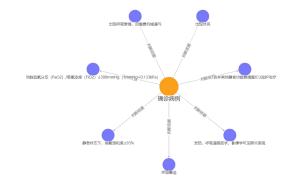
新冠肺炎疫苗项目进展情况

中医证型推荐处方及组成中药材



确诊病例判断依据

疾病相关研究课题





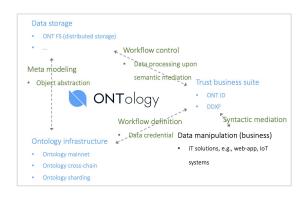
工具和平台建设



知识众包



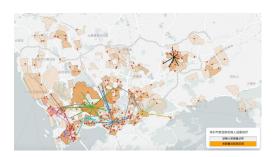
知识上链



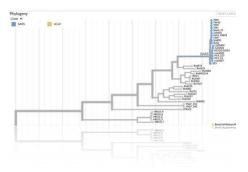
交互问答



辅助决策



预测分析







新冠开放图谱相关数据竞赛



■新冠开放图谱相关数据竞赛



Task1 新冠百科知识图谱类型推断: 围绕新冠百科知识图谱构建中的实体类型推断 (Entity Type Inference)展开,从实体百科页面出发,从给定的数据推断相关实体的类型。

Task2 新冠概念图谱的上下位关系预测:利用自动挖掘的手段从网络文本中采集的细粒度上位概念词,实体和上位词之间以及不同上位词之间复杂的层次关系,自动准确的构建细粒度的上下位层次关系

Task3 新冠科研抗病毒药物图谱的链接预测:依据抗病毒药物图谱Schema及知识图谱的实体、实体属性、实体之间的关系,预测新的两个实体的关系,以进行关系预测,如药物和病毒的靶向作用、蛋白间的交互作用等。

Task4 新冠知识图谱问答评测:面向新型冠状病毒构造了针对健康、医药、疾病防控等特定主旨的问答数据,输入中文问题后,期望问答系统从给定知识库中选择若干实体或属性值作为该问题的答案,并且既能处理百科类的浅层问题,也能处理具备一定领域知识(如流行疾病等)的较深层问题。

相关报名网站: https://www.biendata.com/competition/





THANK YOU



请扫码关注OpenKG,查看更多行业内容

