

引用格式:

田梦晖, 陈明, 席晓桃. 融合 Albert 模型的珍稀濒危植物知识图谱的构建[J]. 湖南农业大学学报(自然科学版), 2023, 49(5): 616–623.

TIAN M H, CHEN M, XI X T. Construction of the knowledge graph for the rare and endangered plants based on Albert model[J]. Journal of Hunan Agricultural University(Natural Sciences), 2023, 49(5): 616–623.

投稿网址: <http://xb.hunau.edu.cn>



## 融合 Albert 模型的珍稀濒危植物知识图谱的构建

田梦晖<sup>1,2</sup>, 陈明<sup>1,2\*</sup>, 席晓桃<sup>1,2</sup>

(1.上海海洋大学信息学院, 上海 201306; 2.农业农村部渔业信息重点实验室, 上海 201306)

**摘 要:** 针对珍稀濒危植物形态特征、分类等级、濒危系数、保护措施等知识不明确的问题, 设计了文本融合轻量级双向转换编码表示模型(Albert)的知识抽取模型框架, 实现批量抽取珍稀濒危植物知识, 从而构建珍稀濒危植物知识图谱: 1) 在现存一般性植物本体的基础上, 采用自顶向下的方式构建珍稀濒危植物本体, 得到 5 个体系, 即物种分类体系、生长形态特征体系、命名体系、保护现状体系和生态习性体系; 2) 采取 Albert 预训练模型来增强下游任务模型输入向量的珍稀濒危植物属性描述文本语义的表征能力; 3) 利用 BiLSTM-CRF 模型和 BiGRU-Attention 模型分别实现命名实体识别和关系抽取。在珍稀濒危植物数据测试集上对模型的有效性进行验证, 结果表明, 命名实体识别模型和关系抽取模型的召回率和准确率的调和平均值(F1)值分别达到 98.07% 和 93.76%, 将得到的大量的实体和关系所形成的三元组存储在图数据库 Neo4j 中, 完成珍稀濒危植物知识图谱的可视化展示。

**关 键 词:** 珍稀濒危植物; Albert 模型; 知识图谱; 本体; 命名实体识别; 关系抽取

中图分类号: TP391.1

文献标志码: A

文章编号: 1007-1032(2023)05-0616-08

## Construction of the knowledge graph for the rare and endangered plants based on Albert model

TIAN Menghui<sup>1,2</sup>, CHEN Ming<sup>1,2\*</sup>, XI Xiaotao<sup>1,2</sup>

(1.College of Information Technology, Shanghai Ocean University, Shanghai 201306, China; 2.Key Laboratory of Fisheries Information, Ministry of Agriculture and Rural Affairs, Shanghai 201306, China)

**Abstract:** Aiming at the problem of unclear knowledge of morphological characteristics, classification levels, endangerment coefficients, and protection measures in the field of rare and endangered plants, a knowledge extraction model framework based on Albert is designed to realize the batch extraction of rare and endangered plant knowledge and construct the knowledge graph of rare and endangered plants: 1) On the basis of the existing general plant ontology, the rare and endangered plant ontology is constructed in a top-down manner, and five systems are obtained, namely, species classification system, growth morphological characteristic system, and nomenclature system, conservation status system and ecological habit system; 2) The Albert model was adopted to enhance the representation ability of the text semantics of the rare and endangered plant attribute description text input vector of the downstream task model; 3) The BiLSTM CRF model and BiGRU Attention model are used to realize named entity recognition and relation extraction, respectively, and the effectiveness of the model was verified on the rare and endangered plant data test set, and the results showed that the harmonic mean (F1) values of recall and accuracy of the named entity recognition model and the relation extraction model reached 98.07% and 93.76%, respectively, and the triples formed by a large number of entities and relationships were stored in the graph database Neo4j in order to complete the visual display of the knowledge graph of rare and endangered plants.

**Keywords:** rare and endangered plants; Albert model; knowledge graph; ontology; named entity recognition; knowledge extraction

收稿日期: 2022-05-08

修回日期: 2023-06-20

基金项目: 上海市科学技术委员会项目(20dz1203800)

作者简介: 田梦晖(1996—), 女, 湖北孝感人, 硕士研究生, 主要从事知识图谱研究, 1178851697@qq.com; \*通信作者, 陈明, 博士, 教授, 主要从事农业信息技术和知识图谱研究, mchen@shou.edu.cn



现状体系和生态习性体系。其中物种分类体系指按照植物物种间的亲缘关系进行的物种分类，按照界、门、纲、目、属、种的层级结构进行分类；生长形态特征体系主要是对植物生长形态特征概念的描述；命名体系主要是针对植物的命名概念的描述；保护现状体系是对植物的各种濒危系数的概念描述；生态习性体系主要囊括了对植物生长环境、物候期、地理分布等概念的描述。

表 1 和表 2 分别展示了基于本体的珍稀濒危植物数据中实体关系和实体属性的映射。表 1 主要描述本体概念层中实体关系三元组实例，包括植物的别称、植物分类、濒危系数等关系，例如〈窄果脆兰，科类，兰科〉这个三元组表达的语义信息是“窄果脆兰是属于兰科的植物”。表 2 主要描述实体的属性关系，包括植物的生活型、高度、生长形态、共用价值、致危因子、保护措施等。

表 1 珍稀濒危植物本体概念层的实体关系

Table 1 The entity relationships of rare and endangered plant ontologies at the conceptual level			
实体关系	实体1	实体 2	三元组
拉丁名	种	学名	〈百山祖冷杉，拉丁名， <i>Abies beshanzuensis</i> 〉
别称	种	别称	〈巨柏，别称，雅鲁藏布江柏木〉
属的类别	种	属	〈宝岛美冠兰，属类，美冠兰属〉
科的类别	种	科	〈宝岛美冠兰，科类，兰科〉
分布	种	地理分布	〈华西杓兰，分布于，云南西北部(中甸)〉
花期	种	花期	〈裸木果，花期，5—7 月〉
果期	种	果期	〈东北红豆杉，果期，8 月〉
国家保护等级	种	国家保护等级	〈百山祖冷杉，国家保护等级，Ⅰ级〉
ICUN	种	濒危等级	〈百山祖冷杉，濒危等级，极危(CR)〉
CITES	种	附录等级	〈沙冬青，附录等级，Ⅱ级〉

表 2 珍稀濒危植物本体概念层的实体属性

Table 2 The entity properties of rare and endangered plant ontologies at the conceptual level		
实体	属性	属性值
种	生活型	乔木，灌木，多年生草本
种	颜色	叶子颜色，树皮颜色
种	表皮毛	小枝的表皮毛，叶的表皮毛
种	形状	树形，冠形
种	分枝方式	多分枝，几不分枝
种	叶的形态	叶数 5~10，叶序二列互生
种	花的形态	花序聚伞花序
种	果的类型	蒴果，蓇葖果，聚花果，瘦果
种	功用价值	全草民间作药用，果核含脂肪，含油量约 40%，油供工业用
种	致危因子	生境退化或丧失，直接采挖或砍伐，栖息地质量衰退，物种内在因素
种	物种现状	2008 年 7 月，中国江西省安福县林业工作者在该县发现一片天然篦子三尖杉
种	保护措施	目前已由当地林场负责保护，应在分布较集中的地区建立自然保护区，严禁樵采
种	优先保护理由	数量稀少，特有，受威胁严重，科学及文化意义大，经济价值高

2 珍稀濒危植物本体知识抽取模型的建立

2.1 知识抽取任务

为了实现珍稀濒危植物的知识抽取，选取了预

训练模型和神经网络模型结合的方式实现知识抽取任务<sup>[14]</sup>。以《中国珍稀濒危植物图鉴》珍稀濒危植物海南风吹楠(*Horsfieldia hainanensis*)为例，利用人工标注提取的三元组如图 2 所示。



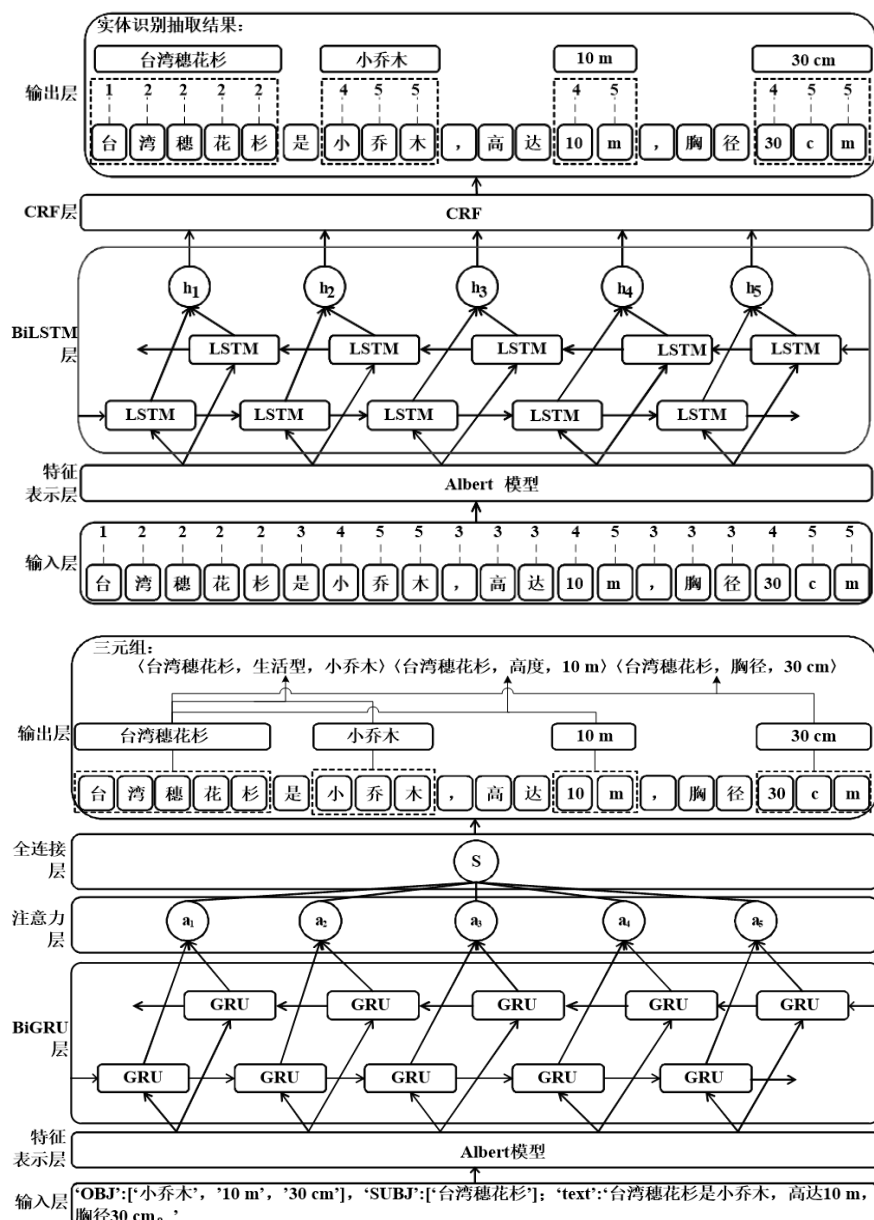


图4 命名实体识别模型和关系抽取模型

Fig.4 Named entities recognition model and relation extraction model

## 2.3 关系抽取模型

设计的关系抽取任务模型如图4所示, 自底向上分为输入层、Albert层、BiGRU层、注意力层、全连接层和输出层。设计训练数据时, 考虑到1条语句中存在1个subject实体与多个object实体相交的情况, 将存在掩码的语句作为1条训练数据输入到Albert模型, 以获得特征向量传输到BiGRU层进行上下文信息语义理解, 最终结合注意力机制层进一步实现对局部信息的抽取, 进行词级别权重分配; 权重高代表对关系抽取的影响较大, 最后全连接层的加入进一步提升关系预测的准确率, 从而实现对实体之间关系的分类。

## 2.4 模型训练与评价

### 2.4.1 数据源

针对数据来源, 首先利用Python脚本, 遵循http协议, 通过解析网页结构请求网站并且获取相应的网页url, 其次采用Xpath和lxml网页解析工具按照设计好的数据规则进行批量爬取, 同时将数据以Json的格式存储在MongoDB数据库当中。最终获取的标注数据总计11688条珍稀濒危植物物种文本。以“猪血木”为例, 收集到的文本数据有: 常绿乔木, 高约15~20 m, 胸径约1.5 m, 全株除顶芽和萼片外均无毛; ……星散分布于广东阳春八甲

村及广西平南思旺村和巴马县灵禄乡；生于海拔 100~400 m 的低丘疏林中或村旁林缘，数量极少；根据实地调查，目前仅在广东阳春八甲村村旁丘陵地田边及八甲小学校园中尚保存有 3 株大树，其他各地似乎均已灭绝。

另外在其结构化数据中，“猪血木”的拉丁学名、科名、国家保护级别、CITES、IUCN 和特有性的实体分别为 *Euryodendron excelsum*、山茶科、I 级、II、CR、中国特有。

2.4.2 数据标注及评价指标

在模型的训练过程中，按照训练需求将数据设

计为特定三元组格式，得到最终的训练数据。例如原文本句子“text”：“海南韶子是常绿乔木，高 5~20 m”。标注的三元组训练数据格式为：“spo\_list”：[{"predicate": “生活型”，“object\_type”: “生活型”，“subject\_type”: “中文名”，“object”: “常绿乔木”，“subject”: “海南韶子”},{“predicate”: “高度”，“object\_type”: “高度”，“subject\_type”: “中文名”，“object”: “5~20 m”，“subject”: “海南韶子”}]。为了避免训练过程出现过拟合现象，将语料按照 7 : 3 的比例将其划分为训练集和测试集，主要对珍稀濒危植物属性描述文本中 28 类关系(表 3)进行分类提取。

表 3 珍稀濒危植物知识抽取试验数据集的分布

Table 3 Distribution of the knowledge extraction experimental datasets for rare and endangered plants					
关系及属性类型	训练集三元组数	测试集三元组数	关系及属性类型	训练集三元组数	测试集三元组数
别称	304	33	花期	996	433
植物分类	2258	968	果期	301	160
表皮毛	342	151	胸径	40	18
颜色	108	114	茎的变态	598	249
生活型	292	263	长度	417	173
高度	543	230	宽度	110	29
形状	1089	479	叶端	82	22
质地	76	55	叶数	33	6
分枝方式	11	3	叶基	68	12
生长方式	257	99	叶序	27	7
生长形态	59	7	果的类型	230	60
染色体	47	20	花序	118	66
地理分布	3073	1379	生长海拔高度	997	427
模式标本地点	936	398	训练语句总数(无重复)	8183	3505
味道	65	27			

2.5 对比结果

按照任务流程，评价模型应充分考虑实体识别和关系抽取的准确率，因此采取精确率(*P*)、召回率(*R*)、精确率和召回率的调和平均值(*F1*)来对模型的性能进行评判。命名实体识别模型结果如表 4 所示。

表 4 珍稀濒危植物命名实体识别模型的精确率和召回率及调和平均值

Table 4 Precision rate, recall rate and F1 value of the named entities recognition model for rare and endangered plants			
模型	<i>P</i>	<i>R</i>	<i>F1</i>
Albert-BiLSTM	0.946 4	0.951 7	0.949 1
所建立的模型	0.982 1	0.979 4	0.980 7
Albert-BiGRU-CRF	0.980 6	0.978 6	0.979 6
BiLSTM-CRF	0.964 7	0.964 8	0.964 7

从表 4 可以看出，所用的 Albert-BiLSTM-CRF 模型的 NER 结果中，在珍稀濒危植物实体识别任务上优于其他模型。Albert-BiLSTM-CRF 模型与 BiLSTM-CRF 发现，加入 Albert 后识别的精确率提高了 1.74%，召回率提高了 1.46%，*F1* 值提高了 1.6%，说明预训练模型的加入能够更好地对输入文本进行语义编码，捕捉到深层次语义信息，从而提高模型识别性能，达到有效识别珍稀濒危植物实体的目的。对比 Albert-BiLSTM-CRF 模型与 Albert-BiLSTM 模型发现，加入 CRF 层后识别的精确率提高了 3.57%，召回率提高了 2.77%，*F1* 值提高了 3.16%，说明加入 CRF 层后对修正最终结果起积极作用，使得识别结果准确率更高；对比 Albert-BiLSTM-CRF 模型与 Albert-BiGRU-CRF 模型发现，精确率提高了 0.15%，召回率提高了 0.08%，

$F1$  值提高了 0.11%，说明针对 NER 任务，BiLSTM 在识别效果上优于 BiGRU。其中，subject 的识别准确度比 object 要高，主要原因在于 subject 是植物中文名，较为单一，而 object 针对的主要是本体中除了中文名的一些概念实体，由于存在人工标注的错误和对某些概念缺乏足够的标注数据的原因，导致 object 的识别率比 subject 要低。

关系抽取模型结果如表 5 所示。

表 5 珍稀濒危植物关系抽取模型的精确率和召回率及调和平均值

**Table 5** Precision rate, recall rate and *F1* value of the relation extraction model for rare and endangered plants

模型	$P$	$R$	$F1$
Albert-BiGRU	0.906 1	0.924 4	0.902 4
所建立的模型	0.935 4	0.943 1	0.937 6
Albert-BiLSTM-Attention	0.928 2	0.958 5	0.933 8
BiGRU-Attention	0.874 9	0.850 8	0.852 5

所建立的关系抽取模型 Albert-BiGRU-Attention 模型效果较好。与 BiGRU-Attention 模型对比,精确率、召回率和  $F1$  值分别提高了 6.05%、9.23%、8.51%,说明 Albert 更好地捕捉到上下文信息,有助于提高模型预测实体之间的关系类别的性能;与 Albert-BiGRU 模型对比,精确率、召回率

和  $F1$  值分别提高了 2.93%、1.87%、3.52%，说明 Attention 层的加入提高了关系抽取的准确度，有效地实现了对珍稀濒危植物实体之间的关系类别的预测；与 Albert-BiLSTM-Attention 模型对比，精确率和  $F1$  值分别提高 0.72% 和 0.38%。

在 28 类关系抽取当中,有一些关系的抽取还存在错误,例如“生长形态”和“分枝方式”,主要是这两类关系在实验语料数据集所占比例较低,以及由于人工标注数据的一些误差,导致一些文本中还存在概念的重叠,因此在训练过程中降低了模型对“生长形态”“生长方式”关系抽取的准确率。

### 3 珍稀濒危植物知识图谱的存储及应用

采用 Neo4j 图数据库存储实体和关系, 采取 2 种方式导入数据: 一是在导入 Json 格式数据集时借助 Neo4j 的 Python 工具包 py2neo, 按照设计好的数据规则直接编码导入 Neo4j; 二是针对 CSV 格式数据, 直接采用 Neo4j 脚本语言 Cypher 语句加载 CSV 至 Neo4j 中, 并实现知识融合, 避免信息冗余, 释放内存压力<sup>[19]</sup>。输入“MATCH (n: Plant) where n.name=“斑叶杓兰” RETURN n”语句后, 在 Neo4j 中可出现如图 5 所示的植物知识图谱可视化结果。

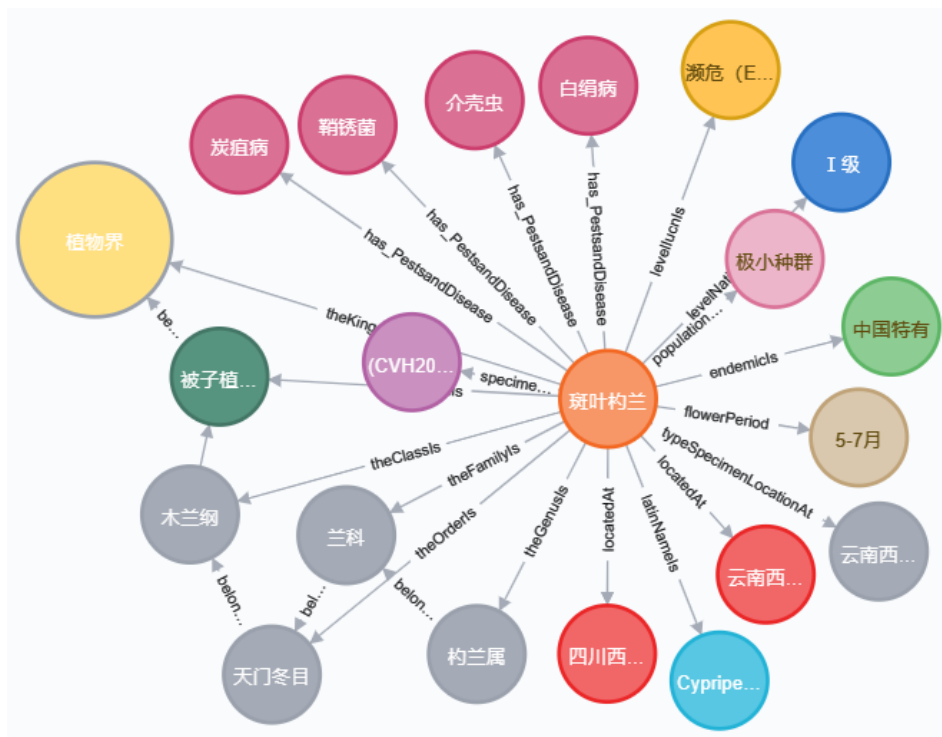


图 5 珍稀濒危植物知识图谱的可视化

**Fig.5 Visualization of the knowledge graph of rare and endangered plants**



## 4 结论

以中国珍稀濒危植物信息系统和植物科学数据中心所提供的《中国植物志》电子版为主要知识来源,构建了珍稀濒危植物本体,设计了一种知识抽取模型流程框架,通过知识图谱构建技术进行知识抽取,获取大量珍稀濒危植物领域三元组,试验结构表明该框架可以实现大批量的知识抽取,有效地提高了准确度,并且满足数据存储要求。

基于知识抽取获取的珍稀濒危植物知识图谱,明确了珍稀濒危植物的物种形态特征、濒危等级、保护现状等信息,为实现植物领域智能系统提供技术支持。后续将围绕基于珍稀濒危植物知识图谱,构建智能问答系统,加强珍稀濒危植物知识关联度。

## 参考文献:

- [1] 陈亚东,鲜国建,寇远涛,等.我国苹果产业知识图谱构建研究[J].中国农业资源与区划,2017,38(11):40-45.
- [2] 于合龙,沈金梦,毕春光,等.基于知识图谱的水稻病虫害智能诊断系统[J].华南农业大学学报,2021,42(5):105-116.
- [3] 张桥英,吴勇.大巴山国家地质公园珍稀濒危植物资源[J].生态环境学报,2018,27(11):2011-2016.
- [4] 王双蕾,韩航,冯金朝,等.基于文献计量学分析沙冬青属植物的研究进展[J].中央民族大学学报(自然科学版),2020,29(1):24-35.
- [5] LAN Z, CHEN M, GOODMAN S, et al. Albert: a lite BERT for self-supervised learning of language representations[C]//ICLR 2020 Area Chairs. ICLR 2020. Addis Ababa: ICLR, 2020.
- [6] 段宇锋,黄思思.基于 BFO 构建中文植物物种多样性领域本体的研究[J].现代图书情报技术,2015(12):72-79.
- [7] 罗贝,吴洁,曹存根,等.从文本中获取植物知识方法的研究[J].计算机科学,2005,32(10):6-13.
- [8] 中国科学院植物研究所.植物智:中国植物志[EB/OL].[2021-10-11].<http://www.iplant.cn/frps>.
- [9] 国家林业局野生动物保护和自然保护区管理司,中国科学院植物研究所.中国珍稀濒危植物图鉴[M].北京:中国林业出版社,2013:249.
- [10] 中国科学院植物研究所.中国珍稀濒危植物信息系统:中国珍稀濒危植物图鉴[EB/OL].[2021-10-11].<https://www.plantplus.cn/rep/protlist>.
- [11] 吴征镒,路安民,汤彦承,等.中国被子植物科属综述[M].北京:科学出版社,2004:6-7.
- [12] 李晓娟,李建秀.山东水龙骨科植物孢粉学研究及其在分类上的意义[J].广西植物,2020,40(4):443-451.
- [13] 刘博,张佳慧,李建强,等.大气污染领域本体的半自动构建及语义推理[J].北京工业大学学报,2021,47(3):246-259.
- [14] 鄂海红,张文静,肖思琪,等.深度学习实体关系抽取研究综述[J].软件学报,2019,30(6):1793-1818.
- [15] 隗昊,周爱,张益嘉,等.深度学习生物医学实体关系抽取研究综述[J].计算机工程与应用,2021,57(21):14-23.
- [16] 马诗语,黄润才.基于 ALBERT 与 BILSTM 的糖尿病命名实体识别[J].中国医学物理学杂志,2021,38(11):1438-1443.
- [17] 张德政,范欣欣,谢永红,等.基于 ALBERT 与双向 GRU 的中医脏腑定位模型[J].工程科学学报,2021,43(9):1182-1189.
- [18] 宋晔璇,陈钊,武刚.基于部分标签数据和经验分布的命名实体识别[J].中文信息学报,2021,35(4):51-57.
- [19] 闫丽华.基于知识图谱的葡萄病虫害自动问答系统[D].杨凌:西北农林科技大学,2021.

责任编辑:罗慧敏

英文编辑:吴志立