

中文信息处理前沿技术进展

Recent Advances in Chinese Information Processing

刘 挺 (Ting Liu)

哈尔滨工业大学社会计算与信息检索研究中心

Research Center of Social Computing and Information Retrieval

Harbin Institute of Technology (HIT-SCIR)

2018年11月7日

提纲 (Outline)

- 概述
- 阶段性进展
- 技术挑战及应对之道
- 行业应用

从感知到认知 (From Perception to Cognition)



运算智能
能存会算

完胜人类



感知智能
能听会说，能看会认

与人类媲美



认知智能
能理解会思考

与人类有一定差距

认知智能是人工智能的高级阶段，是目前制约人工智能取得更大突破和更广泛应用的关键瓶颈

自然语言处理是认知智能的核心

(NLP is the core of cognitive intelligence)



“深度学习的下一个大的进展应该是让神经网络真正理解文档的内容”

深度网络之父: Geoffrey Hinton



“如果给我10亿美金，我会用这10亿美金建造一个NASA级别的自然语言处理研究项目。”

机器学习专家、美国双院院士
Michael I. Jordan



“深度学习的下一个前沿课题是自然语言理解。”

Facebook人工智能负责人: Yann LeCun



“下一个十年，懂语言者得天下”

微软全球执行副总裁: 沈向洋

自然语言处理的研究内容 (Research Fields of NLP)

应用系统 (NLP+)

- 教育, 医疗, 司法, 金融, 机器人

应用技术研究

- 信息抽取, 机器翻译, 问答系统, 文本挖掘

基础研究

- 分词, 词性标注, 句法分析, 语义分析

资源建设

- 语言学知识库建设, 语料库资源建设

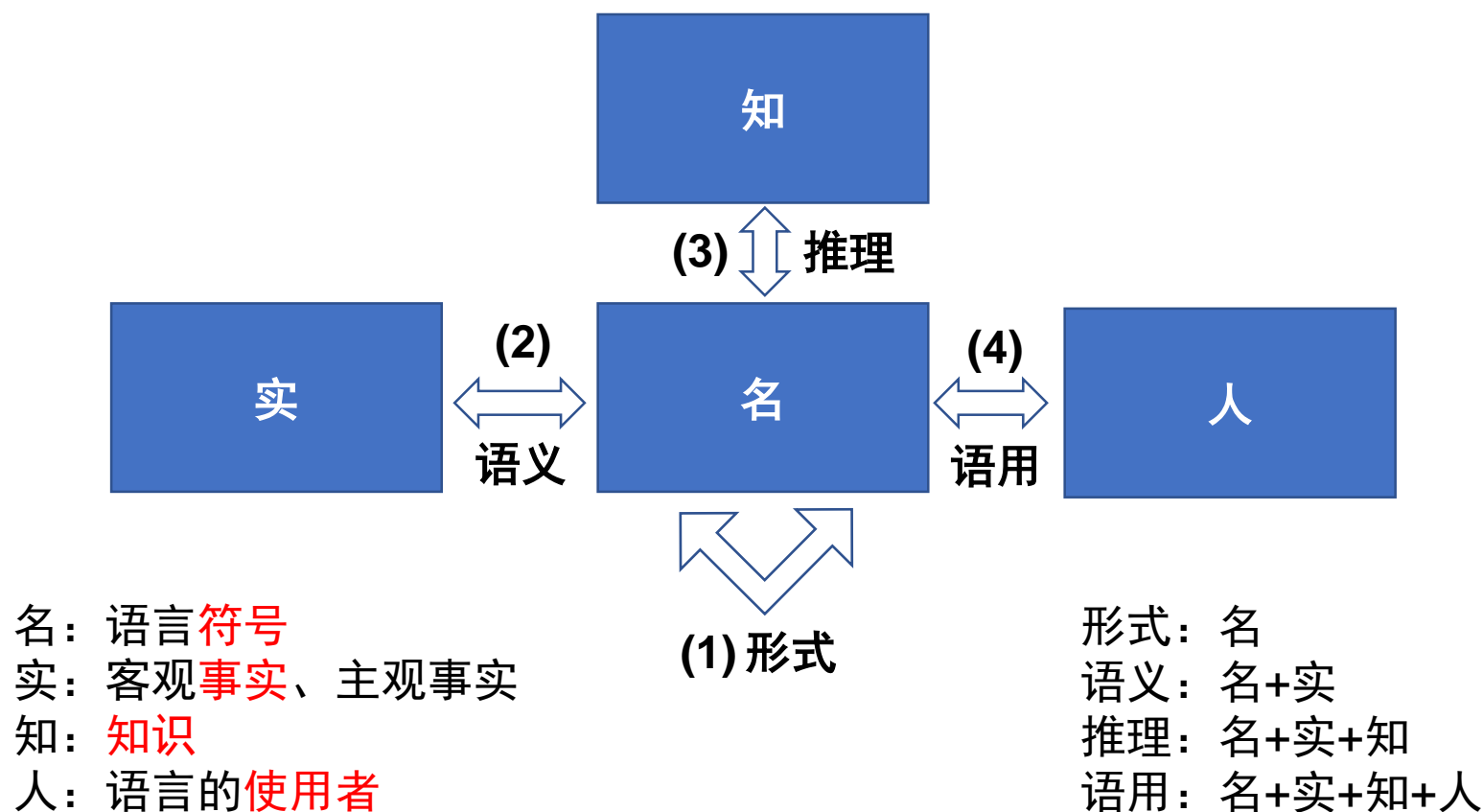
自然语言处理的特点

(Features of NLP problems)

- 优势
 - 存在大量可以利用的**先验知识**
- 难点
 - 研究问题**纷繁复杂**，难以被单一模型处理
 - 难以获得**大量**标注数据
 - 难度大，触及**常识、推理**等认知能力
 - 部分课题**评测**难度高
 - **通用性弱**，与行业关联性强

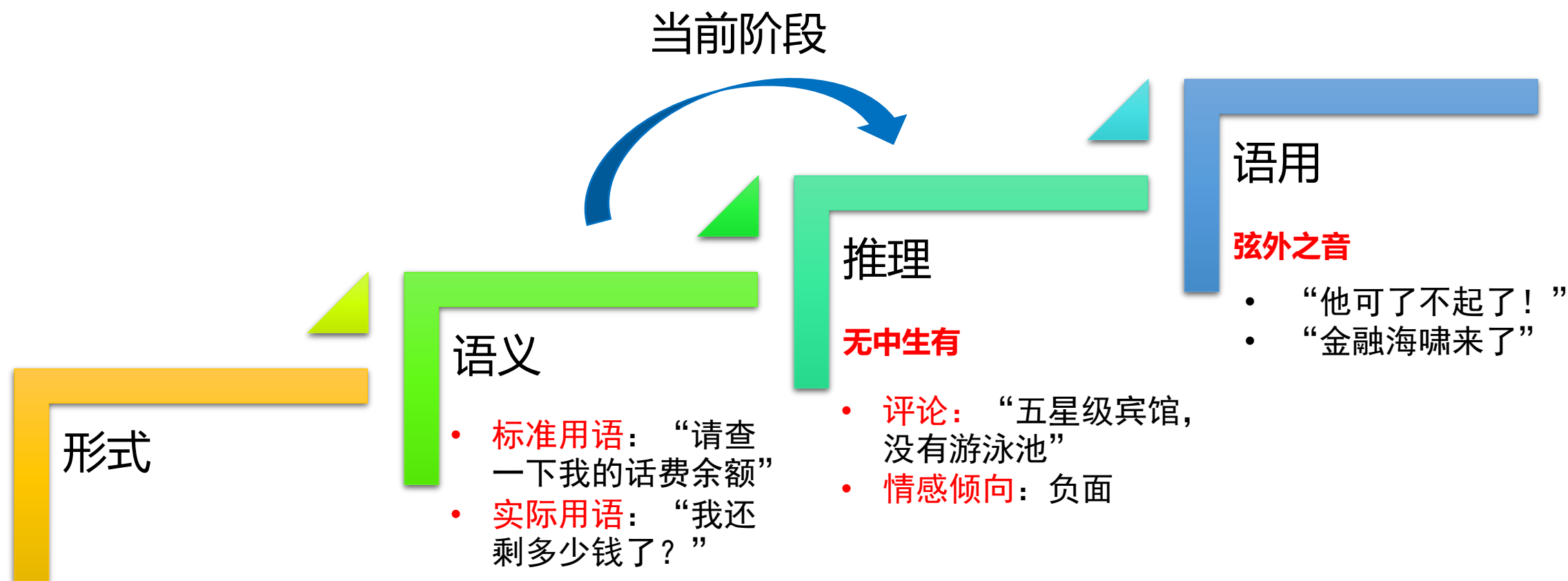
语言理解的四个空间

(Four Spaces of Language Understanding)



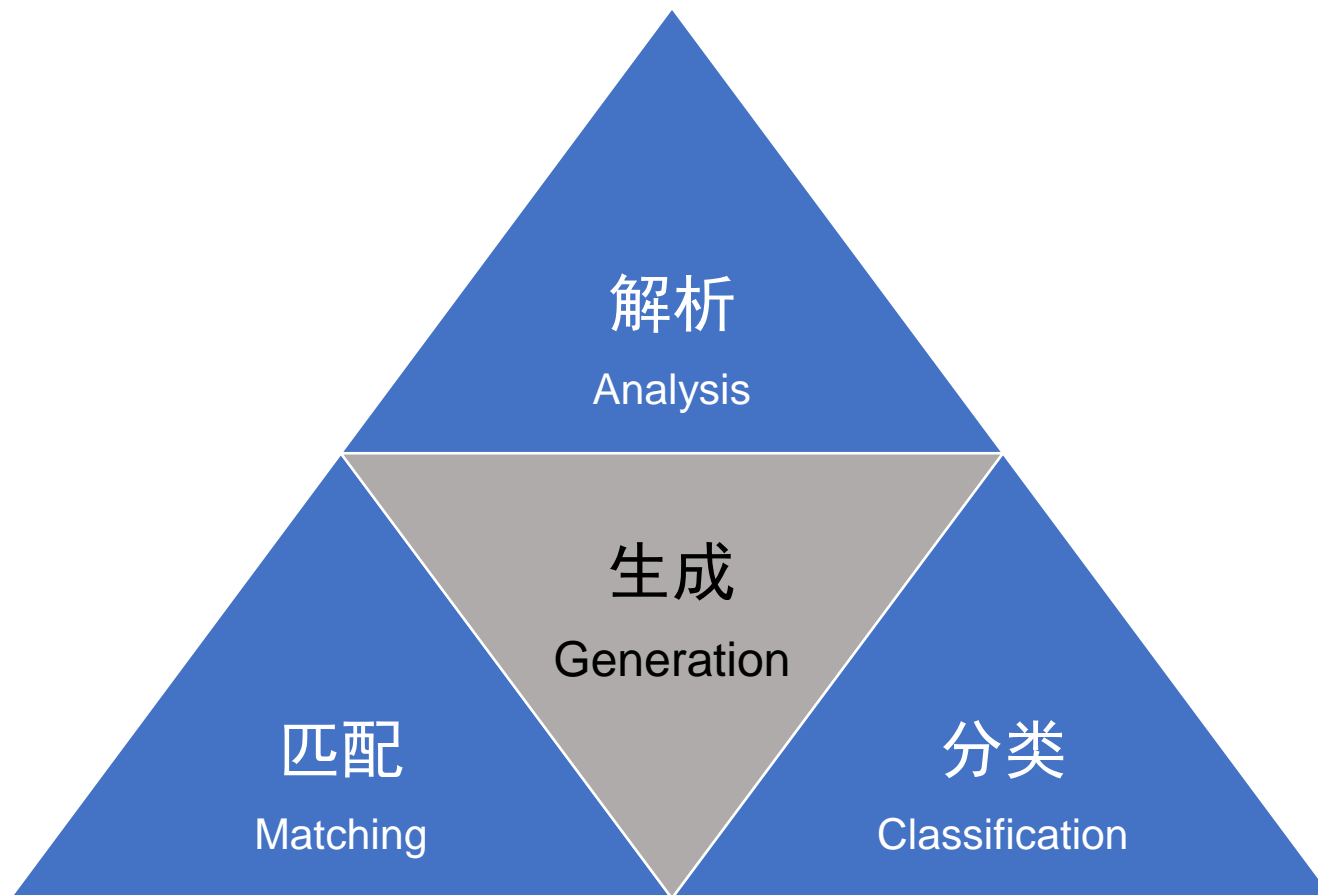
自然语言处理由浅入深的四个层面

(Four Stages of NLP)



自然语言处理的四类问题

(Four Mainstream Tasks in NLP)



NLP，层面×问题二维表 (Stage × Task)

代表性 课题	匹配	分类	解析	生成
形式	搜索	文本分类	词性标注	机械式文摘
语义	问答	情感分析	语义角色标注	机器翻译
推理	文本蕴含	隐式情感分析		写故事结尾
语用		反语		聊天

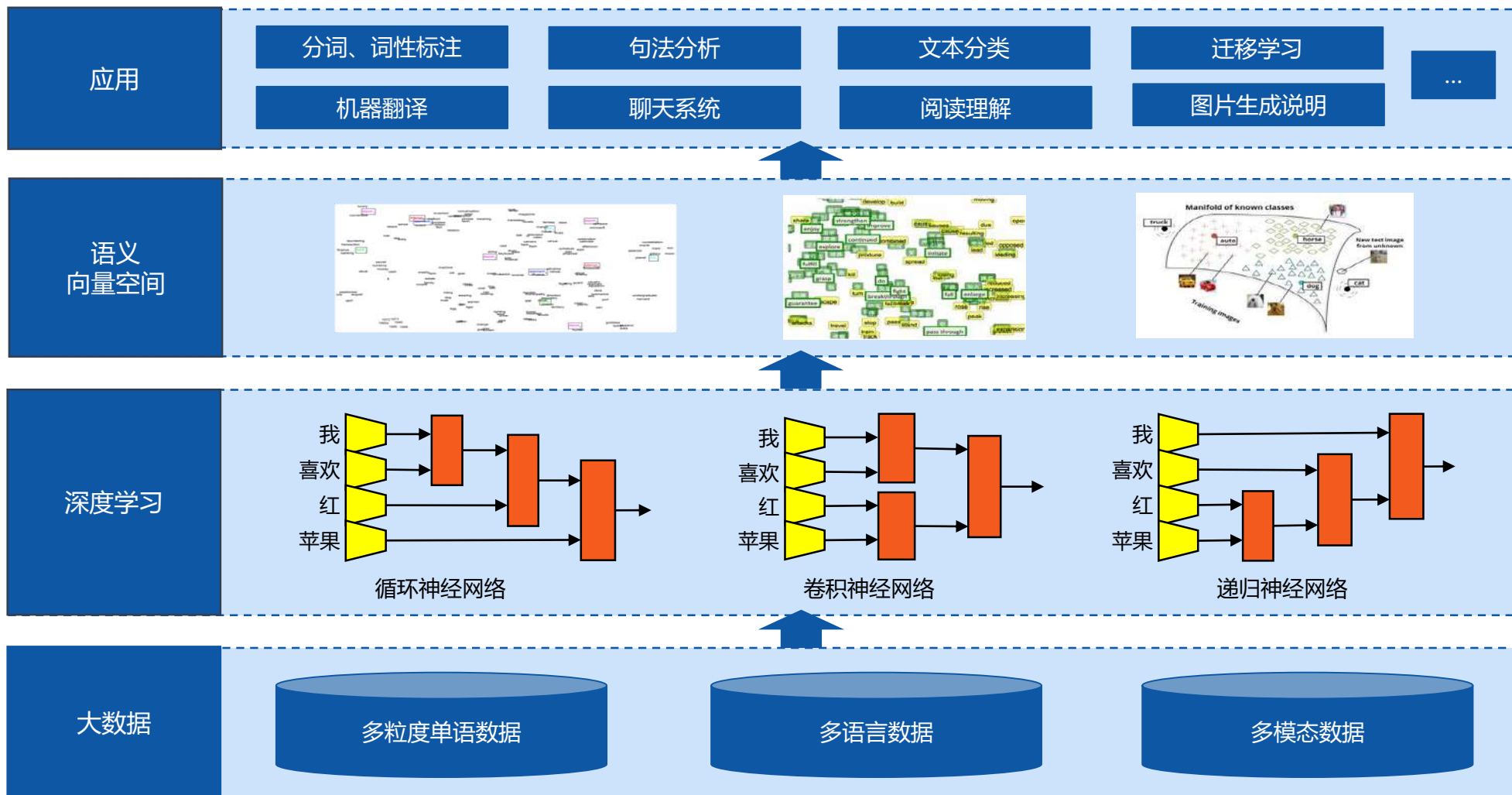
提纲 (Outline)

- 概述
- 阶段性进展
- 技术挑战及应对之道
- 行业应用

阶段性进展 (Progresses in NLP)

- 进展1：广泛采用分布式语义表示
- 进展2：知识图谱开始发挥实际作用
- 进展3：模型预训练显著提升技术指标
- 进展4：机器阅读理解技术在某些数据集上超过人类平均水平
- 进展5：文本情感分析进展明显
- 进展6：文本生成从研究到实用
- 进展7：自然语言处理平台陆续开放
- 进展8：对话系统从应用到平台化

深度学习：目前自然语言处理所采用的主要技术手段 (Deep Learning: the Core Technology in NLP)



进展1：广泛采用分布式语义表示 (Distributed Representation)

• 符号表示

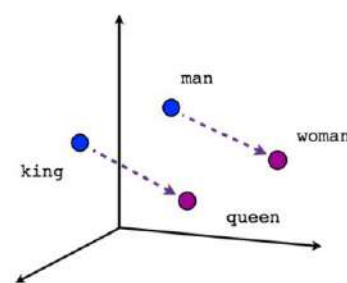
- 离散、高维、稀疏
- One-Hot表示, 词袋表示等



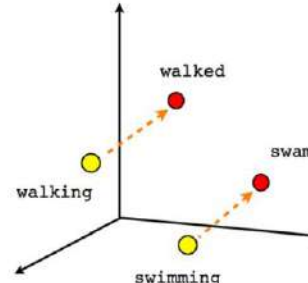
• 分布表示

- 连续、低维、稠密
- 词、短语、句子、篇章
- 便于计算语言单元之间的距离和关系

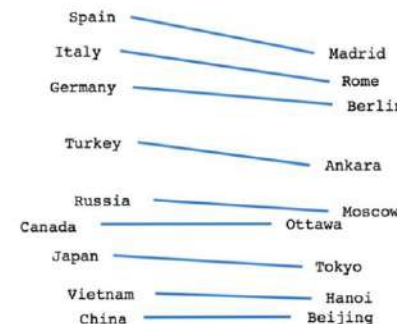
小学 -> [1, 0, 0, 0, 0]
 中学 -> [0, 1, 0, 0, 0]
 大学 -> [0, 0, 1, 0, 0]
 硕士 -> [0, 0, 0, 1, 0]
 博士 -> [0, 0, 0, 0, 1]



Male-Female



Verb tense



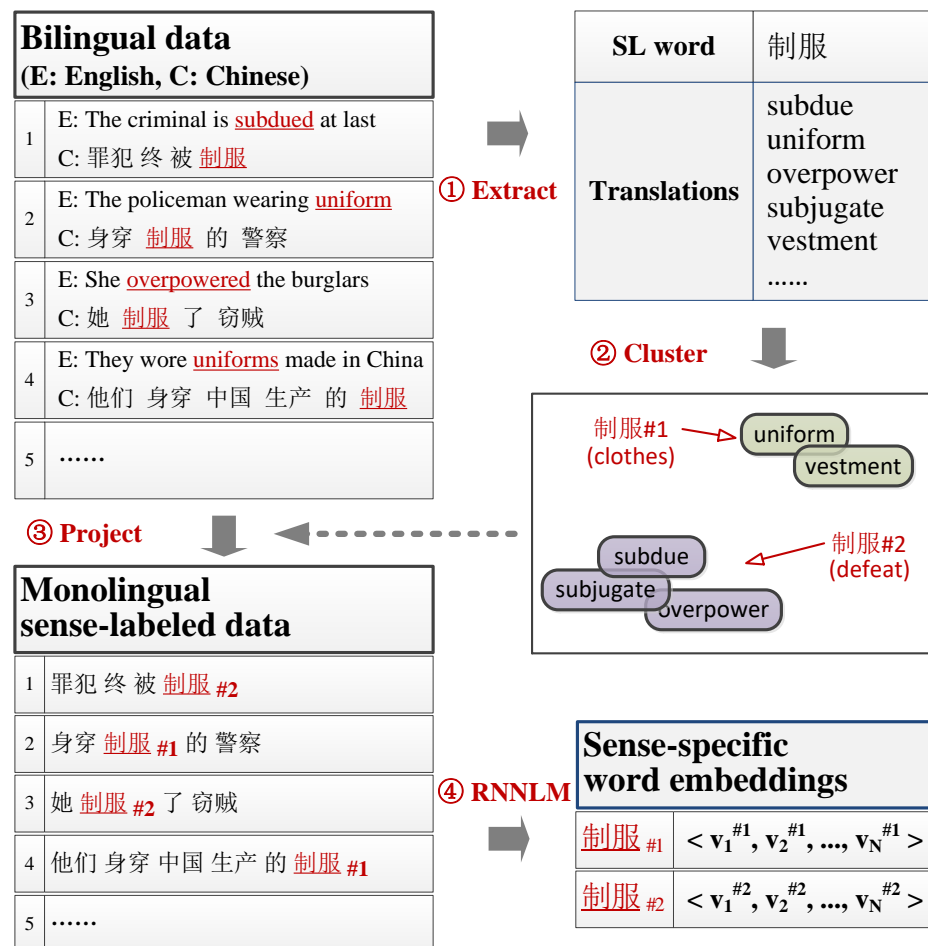
Country-Capital

基于双语的多义词词向量学习

(Learning Sense-specific Word Embeddings via Bilingual Resources)

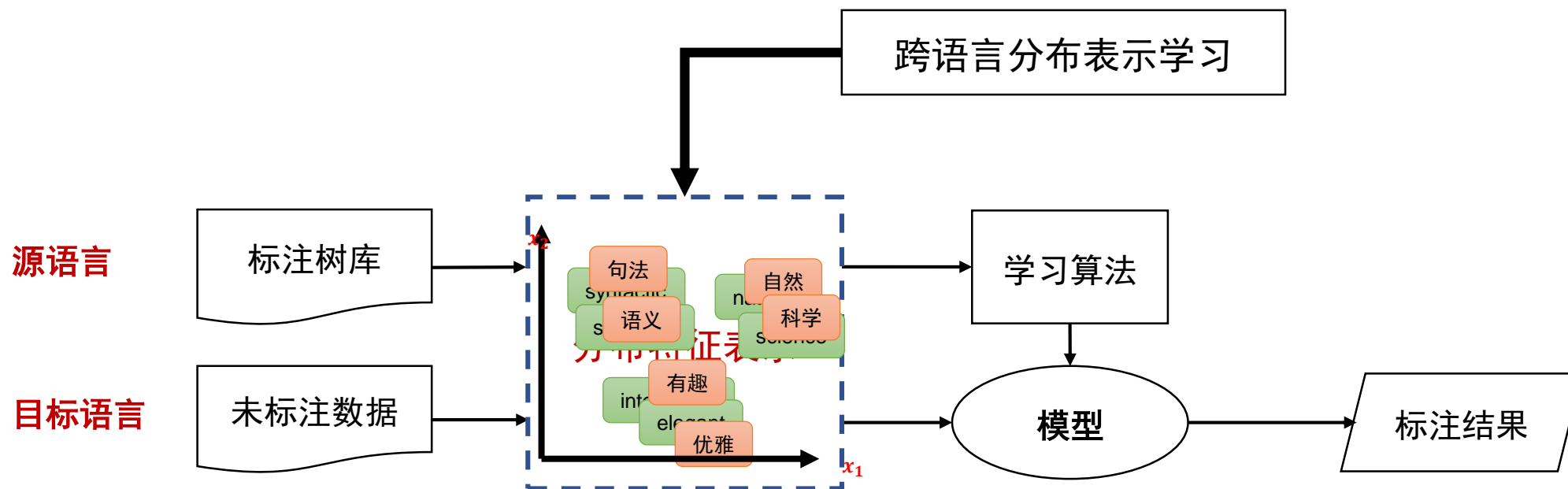
- 传统每个词一个词向量的方法忽略了多义词情况
- 使用双语数据，为多义词的每个词义学习一个词向量
- 评价结果
 - 词相似度任务相关系数绝对值提高超过10%
 - 命名实体识别任务提升1.5%

Jiang Guo, Wanxiang Che, Haifeng Wang, Ting Liu
Learning Sense-specific Word Embeddings By Exploiting Bilingual Resources, COLING 2014



基于跨语言分布表示的句法模型迁移

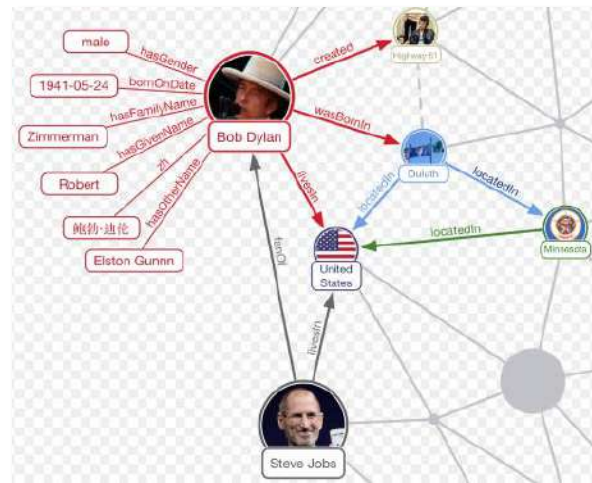
(Cross-Lingual Parsing based on Distributed Representations)



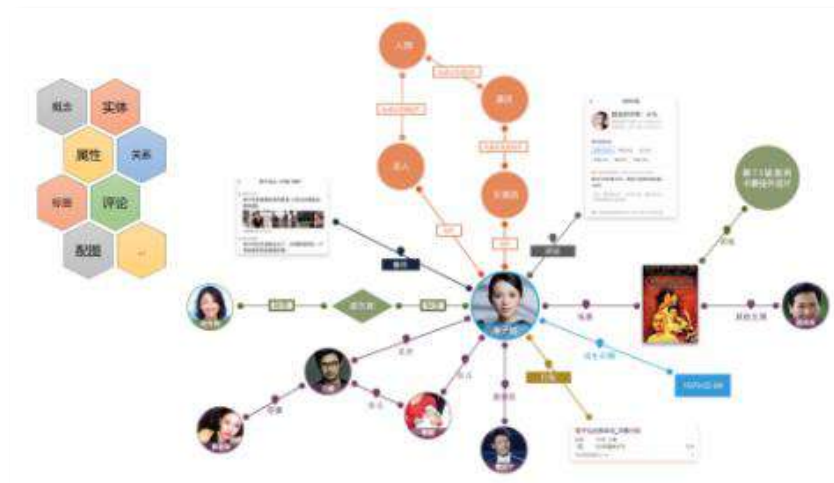
Jiang Guo, Wanxiang Che, David Yarowsky, Haifeng Wang and Ting Liu. Cross-Lingual Dependency Parsing Based on Distributed Representations. ACL 2015.

进展2：知识图谱开始发挥实际作用 (Knowledge Graph)

- 从强逻辑到轻语义
- 知识类型：从概念型到事实型
- 节点：从体词性到谓词性（事理图谱出现）
- 应用目标：从推理到搜索、问答、推理、分析等



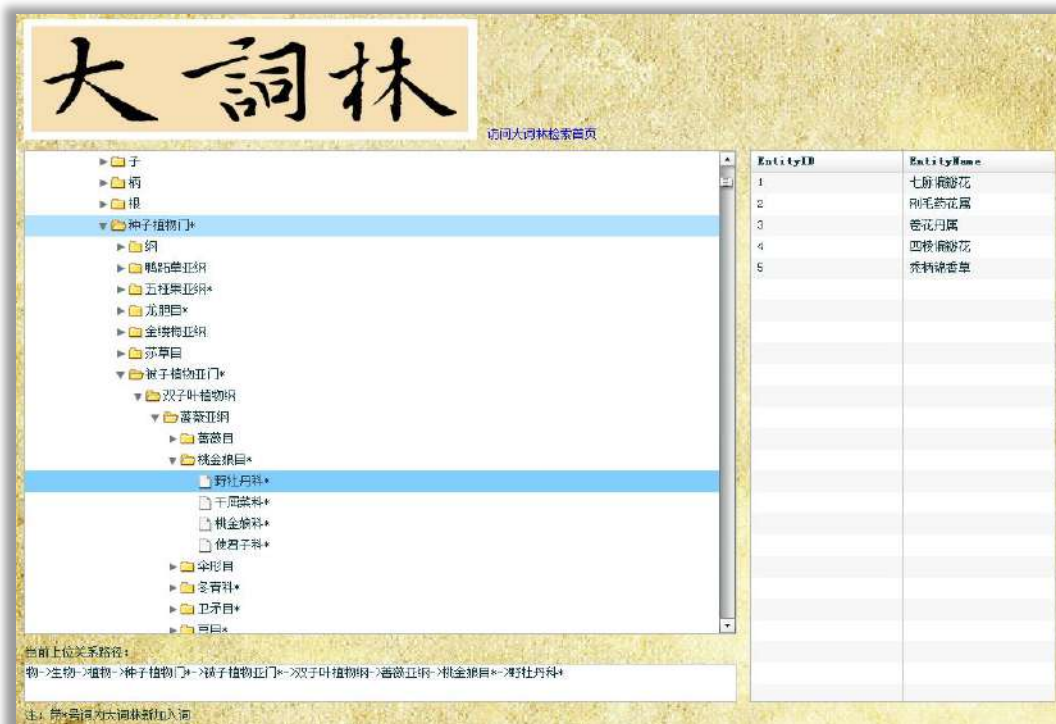
Google
知识图谱



百度
知识图谱

大词林 (BigCilin)

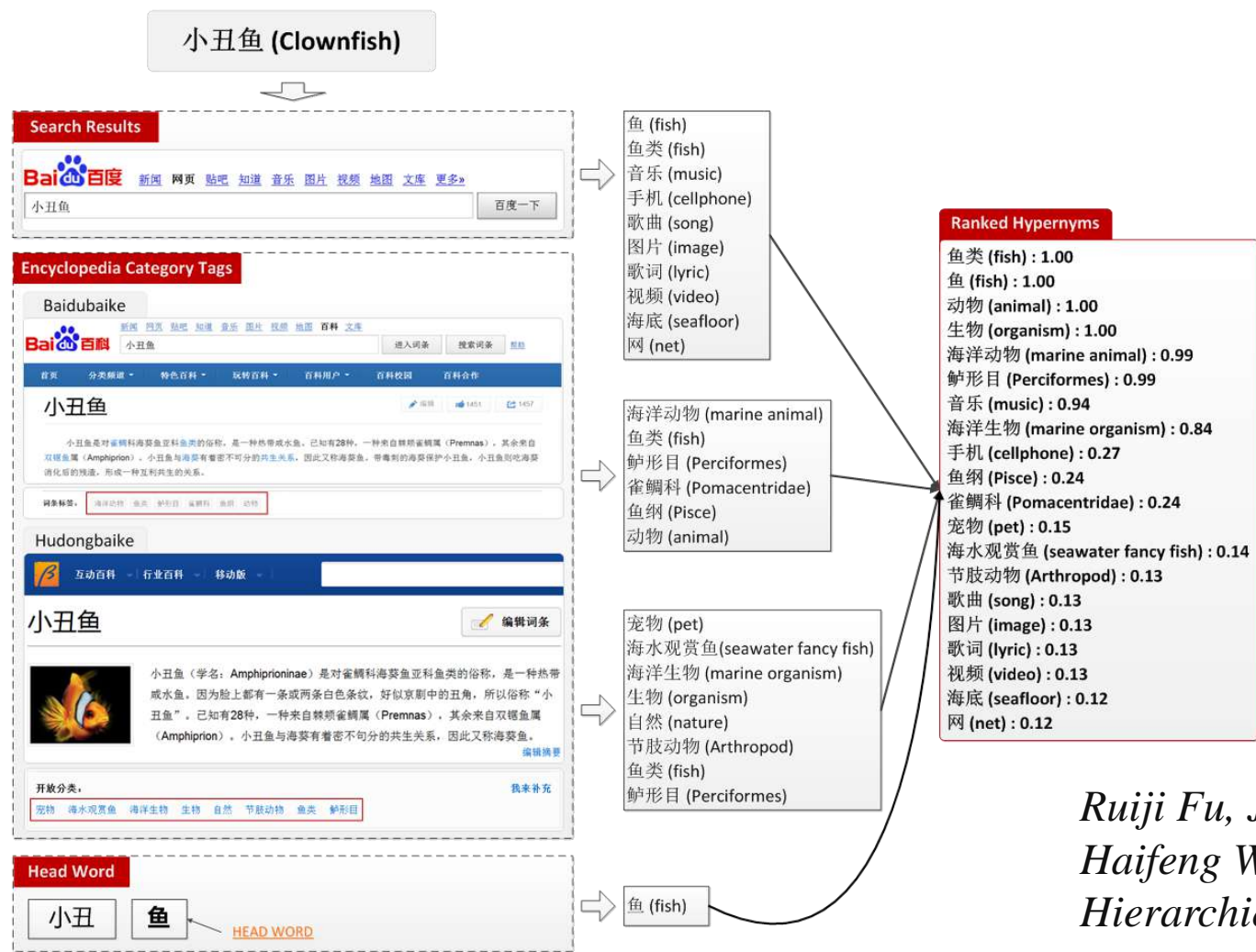
- 大词林以《同义词词林（扩展版）》为骨架，包含细粒度的、多层次的类别体系
- 类别体系结构的层数不固定，依据实体词的不同而动态变化



哈工大SCIR研制，2014年11月27日发布

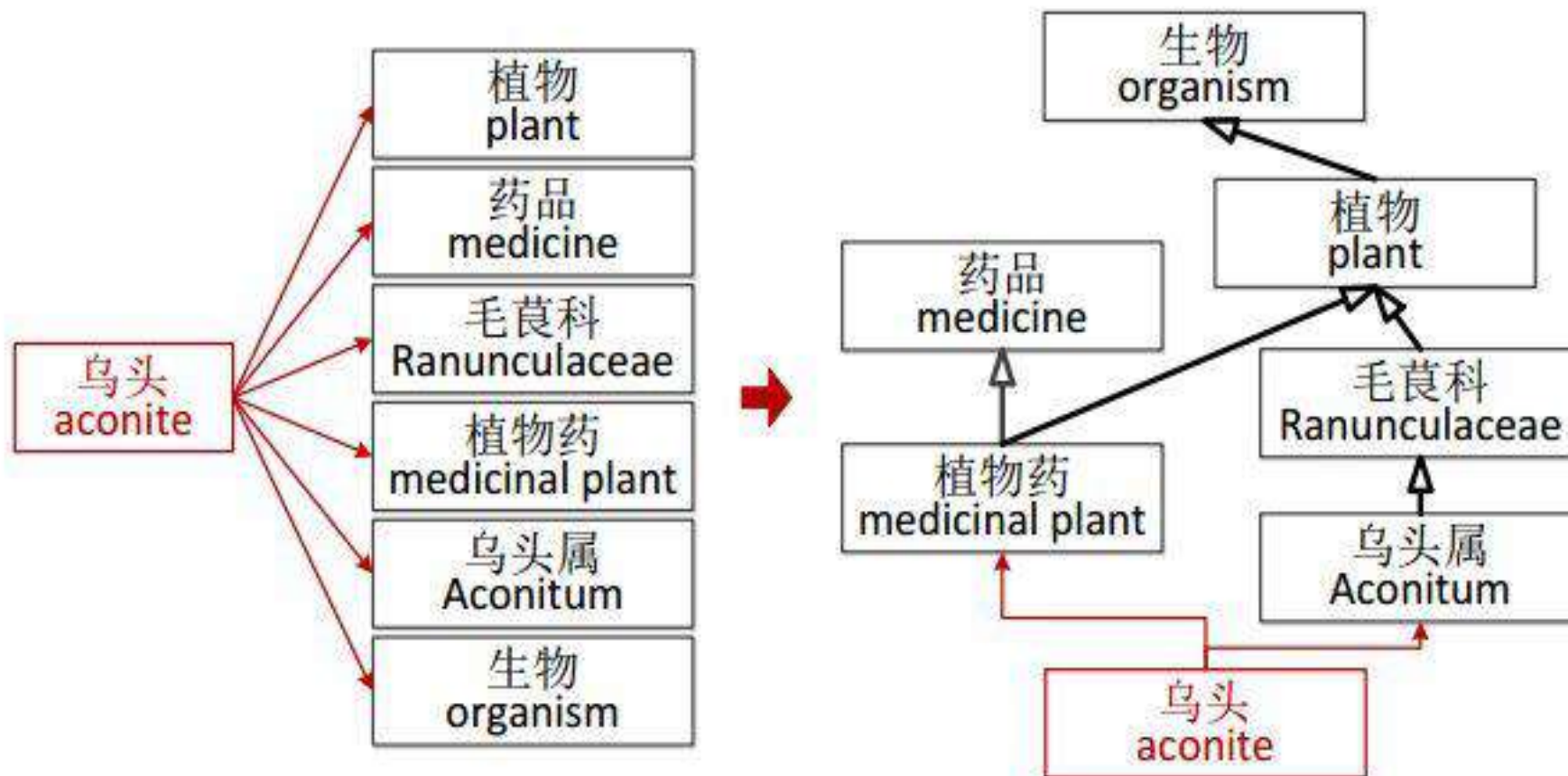


候选类别词获取 (Candidates Selection)



Ruiji Fu, Jiang Guo, Bing Qin, Wanxiang Che, Haifeng Wang and Ting Liu. Learning Semantic Hierarchies via Word Embeddings, ACL 2014

类别词层次划分 (The Hierarchy of Category)



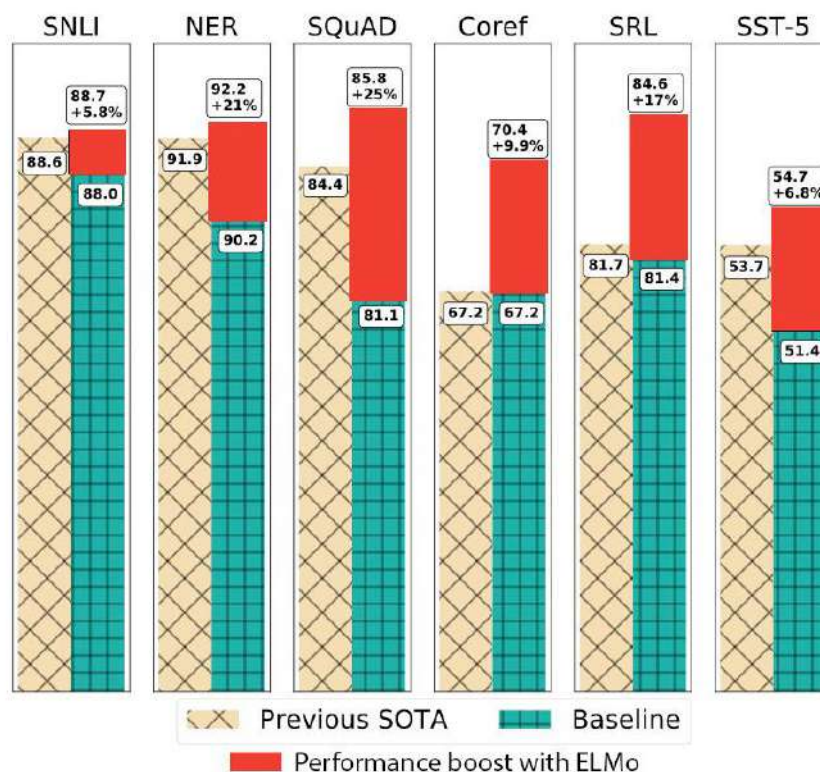
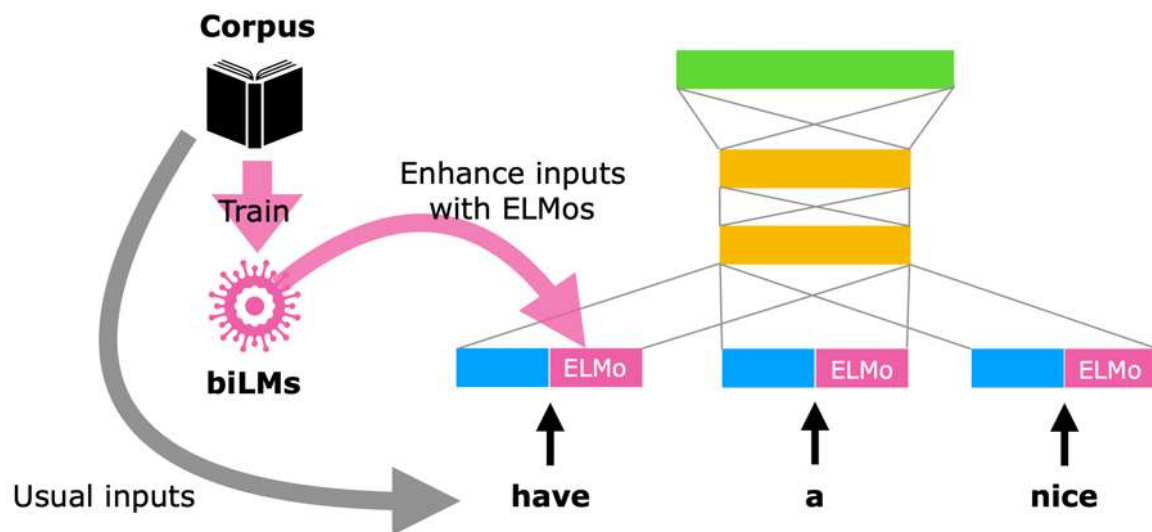
多家大企业付费使用《大词林》

A number of First-tier Companies paid to use BigCilin



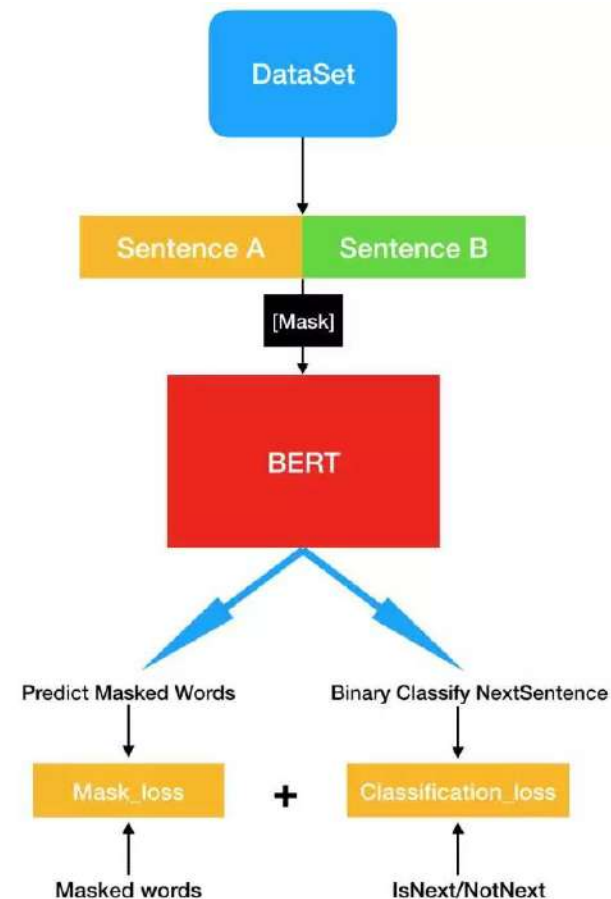
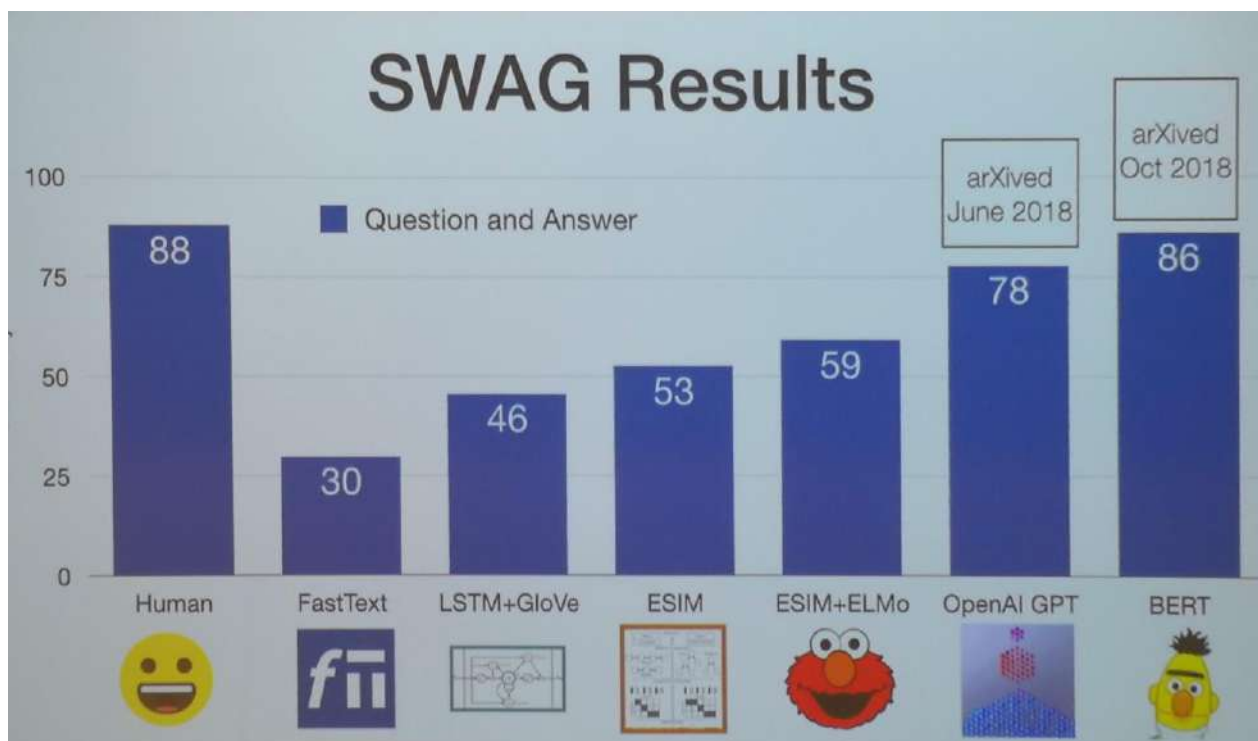
进展3：模型预训练显著提升技术指标 (Pre-training)

- 使用辅助任务预训练，如语言模型
 - ELMo, BERT 等
 - 类似图像处理中基于ImageNet的预训练
 - **大幅提高**多种任务的准确率



最近Google发布了BERT

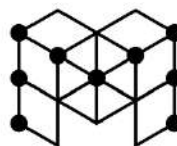
- *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. 2018.*
 - 刷新了 11 项 NLP 任务的当前最优性能记录



进展4：机器阅读理解技术在某些数据集上超过人类平均水平(Machine Reading Comprehension)

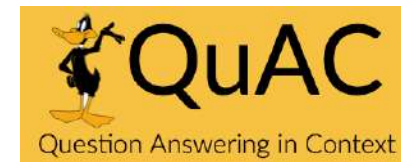


MCTest: A Challenge Dataset for the Open-Domain Machine Comprehension of Text



MS MARCO

Microsoft MACHiNE Reading COMprehension Dataset



Simple RC (2013 ~ 2015)

- bAbI
- MCTest
- 填空类: CNN/Daily Mail, CBT

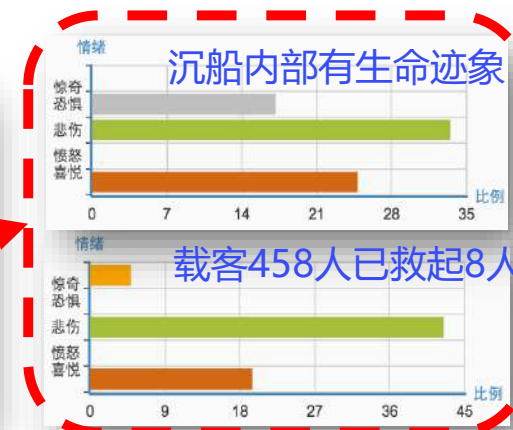
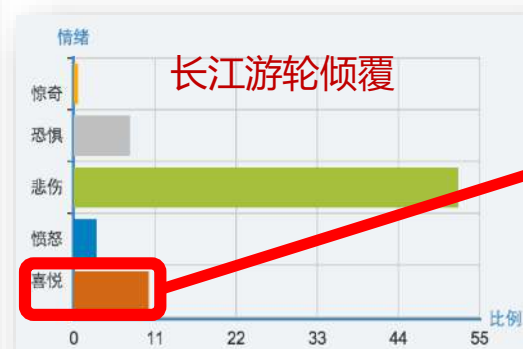
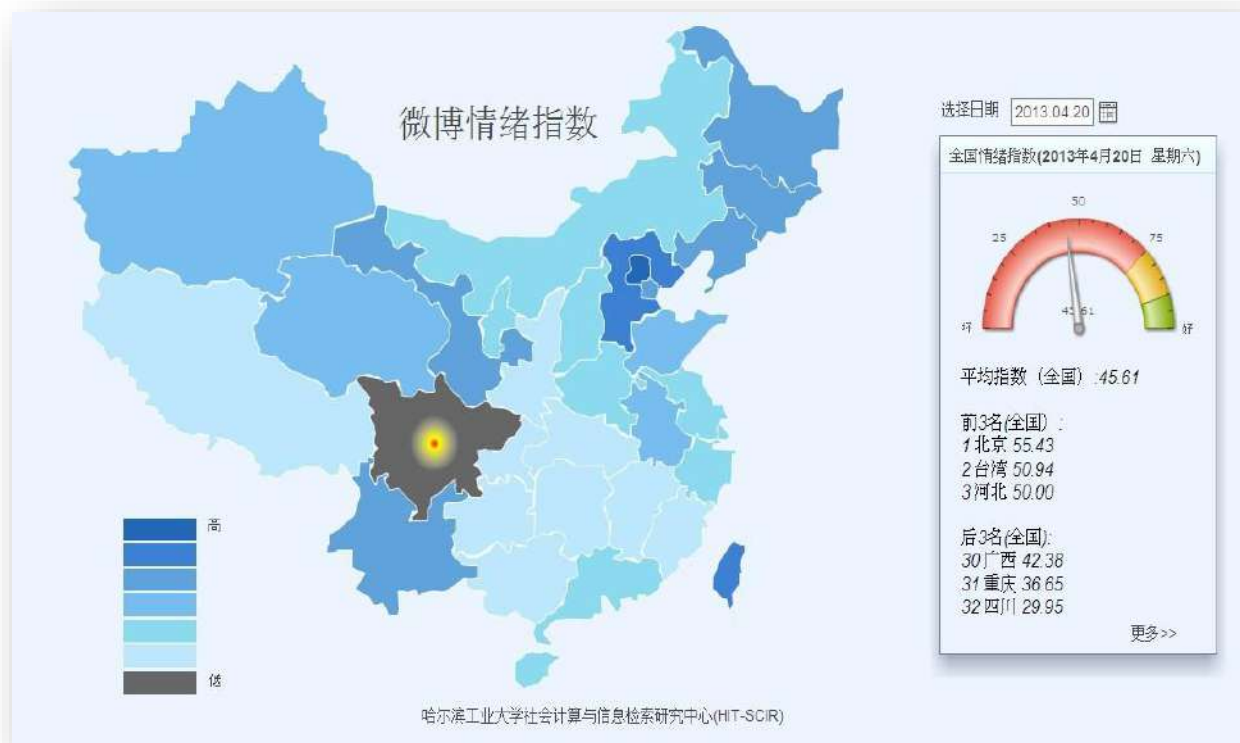
Complex RC (2016 ~)


- 片段抽取: SQuAD
- 多选题: RACE
- 常识阅读理解: SemEval 2018 Task 11
- 开放域: MS MARCO


General RC (2018 ~)

- 拒答问题
- 对话式阅读理解
- 多跳阅读理解

进展5：文本情感分析进展明显 (Sentiment Analysis)



正义网 
#长江400人游轮倾覆# 【沉船船底已露出水面 记者近距离拍摄救援照片】据航道部门扫测，沉船位置已确定，事故水域水深约15米，目前沉船船底已露出水面，沉船处已设沉船标。从事故现场反馈回消息，潜水员潜入后敲击船体，船内有应答，已发现生命迹象！❤️ (@湖北日报 记者张朋)

江苏新闻 
【一位85岁老太太被成功救援出水！❤️】12点55分，一位85岁老太太被成功救援出水，并立即安排救护。目前健康状况良好。相信奇迹，期待奇迹！❤️ @湖北日报

进展6：文本生成从研究到实用 (Text Generation)



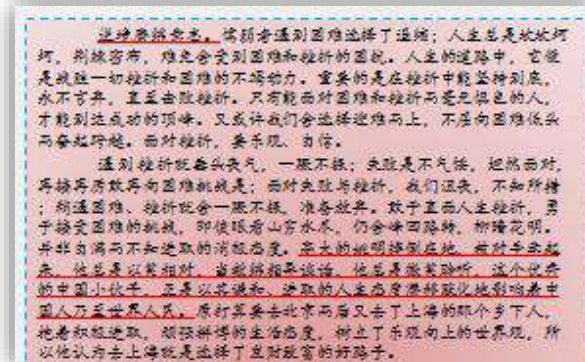
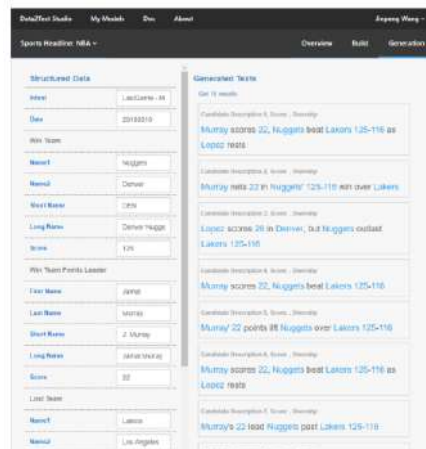
橄榄球赛自动报道
(美联社+Automated Insights公司)

写古体诗
(清华九歌)

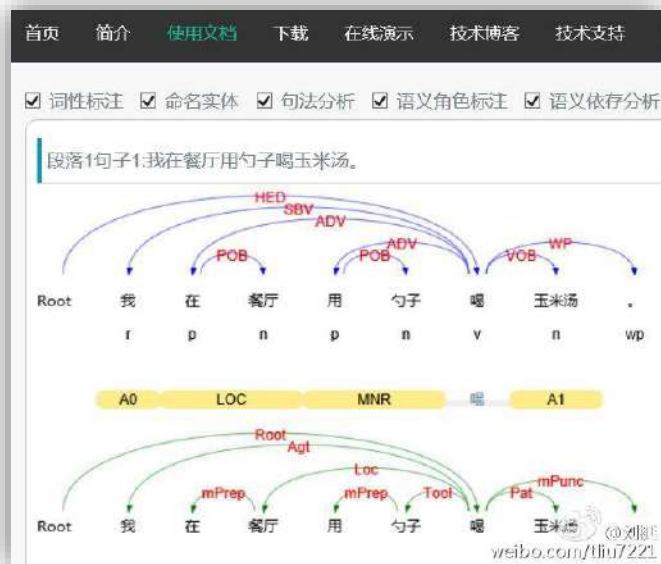
写现代诗
(微软小冰)

Data2Text
(微软)

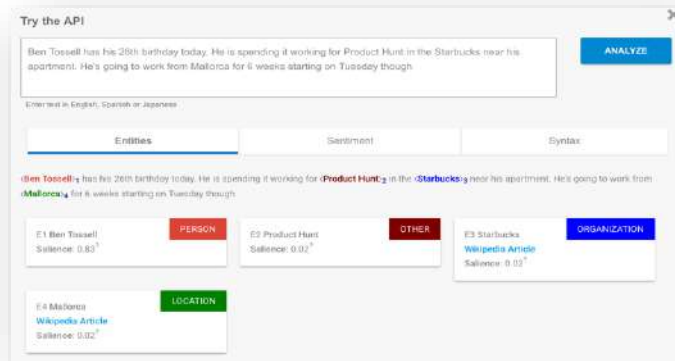
写作文 (哈工大SCIR)



进展7：自然语言处理平台陆续开放 (NLP Open Platform)



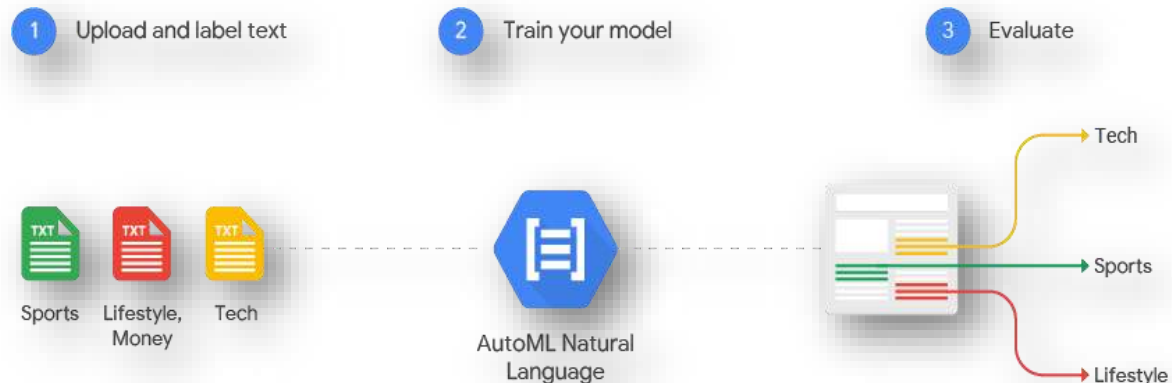
HIT-SCIR LTP (2011)



Google NLP Cloud (2016)



Baidu AI (2017)



Google AutoML NLP (2018)

“语言技术平台”的应用

(Applications of Language Technology Platform)

- 2003年开始研制
- 国内外500余家学术机构签署协议，使用该平台
 - 卡耐基梅隆大学、俄亥俄州立大学、德克萨斯州立大学、伊利诺斯州大学、哈佛商学院、新加坡国立大学、日本德岛大学，……
- 百度、腾讯、讯飞、华为、金山等企业付费使用



哈工大SCIR在CoNLL 2018评测中排名第一

(HIT-SCIR Ranked Top-1 in CoNLL 2018)



- CoNLL 2018多语言通用依存句法分析评测
 - 包括分句、分词、词性标注、依存句法分析任务
 - 数据：57种语言、82个树库
 - 统一的词性和句法标注体系
- 参赛队伍包括：斯坦福大学、IBM等著名大学和企业
- 哈工大获得**第1名**，高出第二名**2.5%**

LAS Ranking

1. HIT-SCIR (Harbin)	75.84 ± 0.14 [OK]	(p<0.001)
2. TurkuNLP (Turku)	73.28 ± 0.14 [OK]	(p=0.039)
3-5. UDPipe Future (Praha)	73.11 ± 0.13 [OK]	(p=0.221)
3-5. LATTICE (Paris)	73.02 ± 0.14 [OK]	(p=0.461)
3-5. ICS PAS (Warszawa)	73.02 ± 0.14 [OK]	(p<0.001)
6. CEA LIST (Paris)	72.56 ± 0.14 [OK]	(p=0.036)
7-8. Uppsala (Uppsala)	72.37 ± 0.15 [OK]	(p=0.191)
7-8. Stanford (Stanford)	72.29 ± 0.14 [OK]	(p<0.001)
9-10. AntNLP (Shanghai)	70.90 ± 0.15 [OK]	(p=0.242)
9-10. NLP-Cube (București)	70.82 ± 0.14 [OK]	(p=0.032)
11. ParisNLP (Paris)	70.64 ± 0.14 [OK]	(p<0.001)
12. SLT-Interactions (Bengaluru)	69.98 ± 0.14 [OK]	(p<0.001)
13. IBM NY (Yorktown Heights)	69.11 ± 0.16 [OK]	(p<0.001)
14. UniMelb (Melbourne)	68.66 ± 0.15 [OK]	(p=0.002)
15. LeisureX (Shanghai)	68.31 ± 0.16 [OK]	(p<0.001)
16. KParse (İstanbul)	66.58 ± 0.16 [OK]	(p=0.015)
17. Fudan (Shanghai)	66.34 ± 0.15 [OK]	(p<0.001)
18. BASELINE UDPipe 1.2 (Praha)	65.80 ± 0.15 [OK]	(p=0.048)
19. Phoenix (Shanghai)	65.61 ± 0.16 [OK]	(p<0.001)
20. CUNI x-ling (Praha)	64.87 ± 0.16 [OK]	(p<0.001)
21. BOUN (İstanbul)	63.54 ± 0.15 [OK]	(p<0.001)
22. ONLP lab (Ra'anana)	58.35 ± 0.15 [81]	(p<0.001)
23. iParse (Pittsburgh)	55.83 ± 0.11 [65]	(p<0.001)
24. HUJI (Yerushalayim)	53.69 ± 0.15 [80]	(p<0.001)
25. ArmParser (Yerevan)	47.02 ± 0.11 [66]	(p<0.001)
26. SParse (İstanbul)	1.95 ± 0.00 [2]	

汉语与英语存在的显著差异

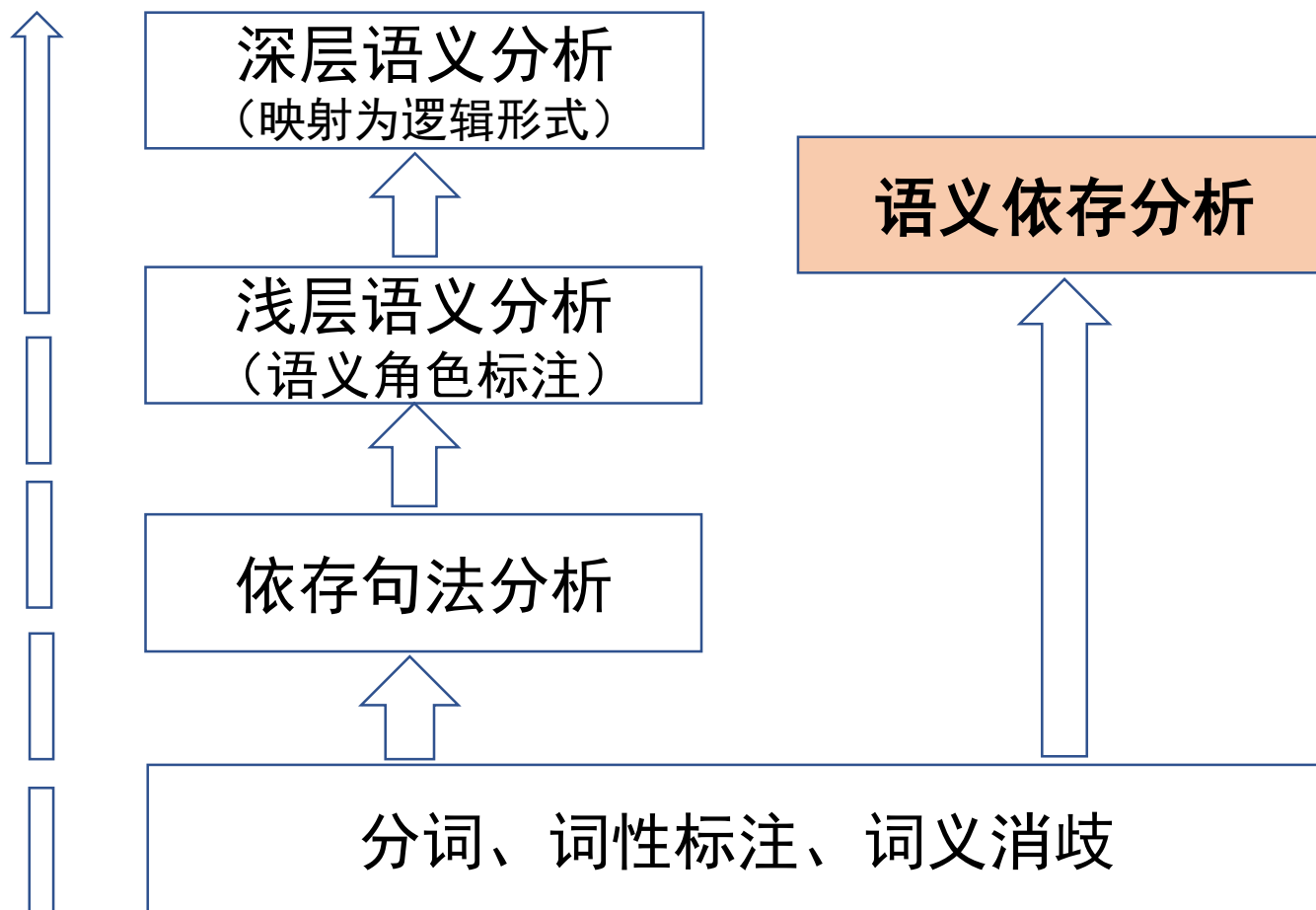
(Significant Differences between Chinese and English)

- 英语重形合，汉语重意合，缺乏形态变化，语法灵活

	中文	英文	例子
形态变化	无	有	机器翻译 - Machine Translation 翻译人员 - Translator 翻译小说 - Translated Novel 翻译经验 - Translating experiences
时态变化	无	有	我回家了 I will go home. / I went home.
句法形式	灵活	规范	你买 票 了么？ / 票 你买了么？ Have you bought the ticket?
复杂名词短语	多	少	中国北京红十字芦山抢险救援队 “五一” 节期间工作掠影

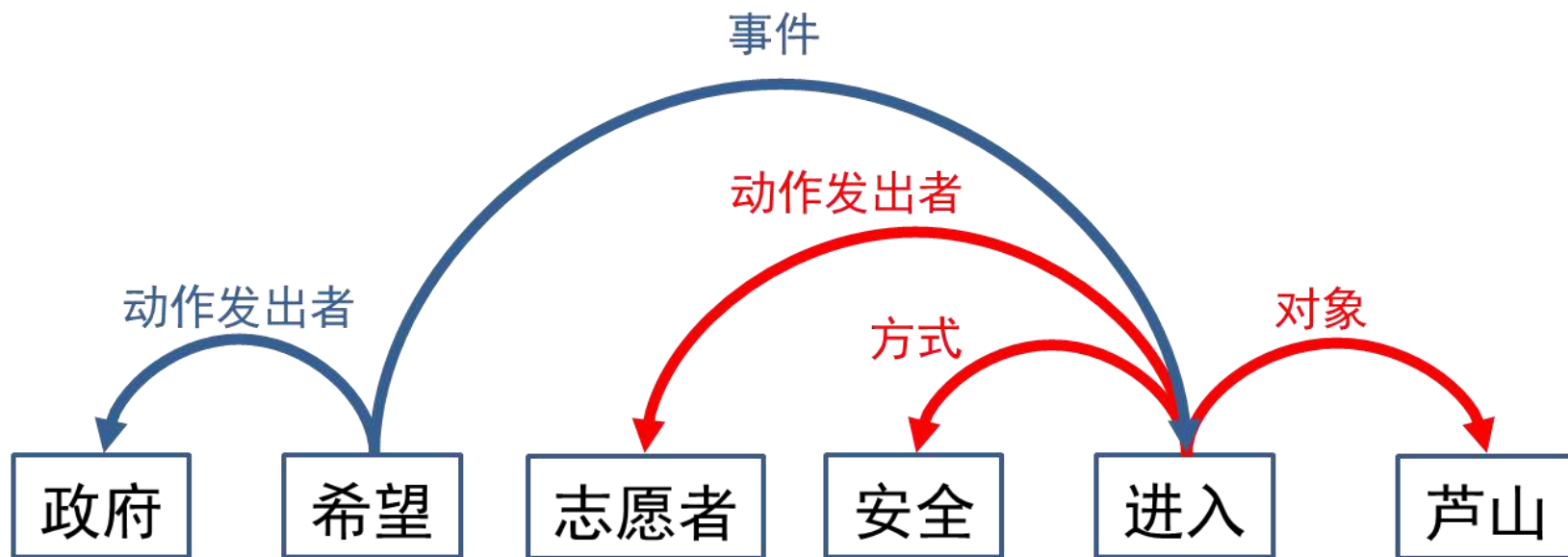
新型的句子级中文分析技术架构 (Sentence-level Chinese Technology Platform)

错误级联效应导致最终结果的精度很低



中文句子语义依存树表示

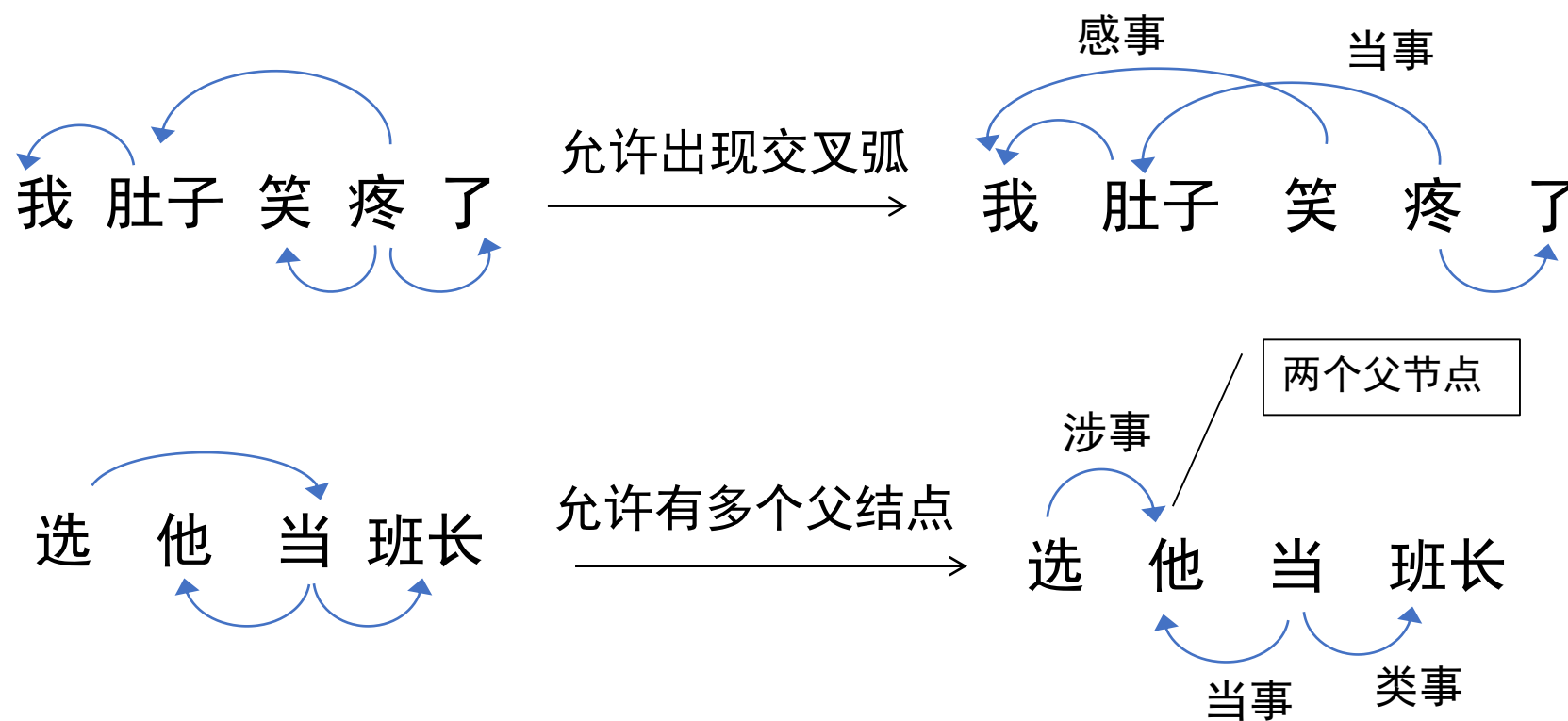
(Chinese Sentence Semantic Dependency Tree)



形式上类似于依存语法，但弧上标注的是“语义”关系
本质上是以语义与句法的结合体

从“树”到“图” (From "Tree" to "Graph")

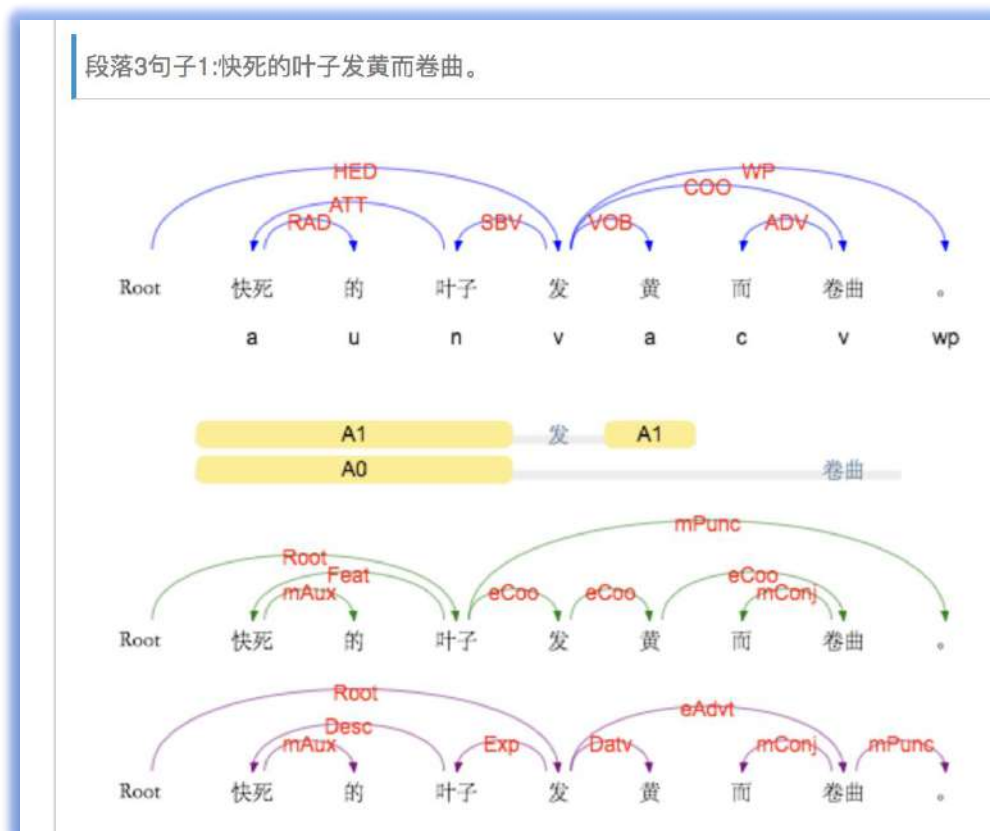
- 将句子自动分析的语义表示形式从“树”推进到“**图**”



语义依存分析系统

(The System of Semantic Dependency Graph Parsing)

- <http://ltp.ai/demo.html>



进展8：对话系统从应用到平台化 (Dialog System)



2011年，苹果



聊天机器人 (2014年)

- 微软: 小冰
- 百度: 小度
- **哈工大SCIR: 笨笨**



对话平台 (2015年)

- 微软: LUIS
- FaceBook: wit.ai
- Google: api.ai
- 百度: UNIT
- **哈工大SCIR: DTP**

提纲 (Outline)

- 概述
- 阶段性进展
- 技术挑战及应对之道
- 行业应用

技术挑战(Technical Challenges)

- 挑战1：带标数据不足
- 挑战2：常识知识不足
- 挑战3：可解释性问题
- 挑战4：知识工程与统计方法的融合问题
- 挑战5：文本领域迁移问题
- 挑战6：文本推理问题
- 挑战7：对话语(Discourse)、语用(Pragmatics)的研究
- 挑战8：基于多模态融合的文本理解

自然语言处理中的数据、知识

大类	细分类	特点	举例
数据 (自动、隐性)	有标注	专家标注、众包	Penn TreeBank
	无标注	原始语料	《人民日报》、微博
	弱标注	量大	情感分析中对表情符的利用
知识 (人工、显性)	元知识	关于知识的知识	人工定义的表示，特征工程
	语言知识	词典、规则库	WordNet、大词林(BigCilin)
	常识知识	很难从文本中挖到	CYC
	世界知识	可以从文本中挖到	知识图谱

挑战1：带标数据不足 (Unavailability of Large-scale labeled Data)

- 深度学习需要大量的带标记数据

	简单问题	复杂问题
小数据	无明显优势 (词性标注)	有劣势 (语义分析)
大数据	优势最明显 (语言模型)	较明显 (机器翻译)

解决之道：弱标注数据 (Weakly Annotated Data)

- 弱标注数据
 - 是带标签的训练数据
 - 不曾面向所研究的任务进行人工标注
 - 标签是样本的近似答案，而不是精确答案
- 弱标注数据的类型
 - 自然标注大数据 (Naturally Annotated Big Data)
 - 自动产生的数据 (Automatically Generated Big Data)
- 弱标注数据的优点与不足
 - 优点：数据量大，制作成本低
 - 不足：因为是近似标注，所以噪声大

获取“弱标注数据”的两种方法

(Two Methods of Constructing Weakly Annotated Data)

获取弱标注数据的方法	要点	样本与真实样本相比	标签与真实标签相比
寻“找”	同一语义有多种表现形式，其中一种形式的语义歧义性小	一致	近似
制“造”	对样本进行编辑操作形成伪数据	近似	一致

制“造”弱标注数据 (Construction of Weakly Annotated Data)

		任务	方法
制造	修改（换）	词义消歧	等价伪词
	删除（挖）	阅读理解	基于挖词模型
		零指代	基于挖词模型
	增加（插）	文本顺滑	序列标注

弱标注数据样本与真实样本近似，弱标注数据标签与真实标签一致

制造出来的大规模有标注阅读理解题

(Constructing Weakly Annotated Data in Reading Comprehension)

Original Version	Anonymised Version
Context The BBC producer allegedly struck by Jeremy Clarkson will not press charges against the “Top Gear” host, his lawyer said Friday. Clarkson, who hosted one of the most-watched television shows in the world, was dropped by the BBC Wednesday after an internal investigation by the British broadcaster found he had subjected producer Oisin Tymon “to an unprovoked physical and verbal attack.” ...	the <i>ent381</i> producer allegedly struck by <i>ent212</i> will not press charges against the “ <i>ent153</i> ” host , his lawyer said friday . <i>ent212</i> , who hosted one of the most - watched television shows in the world , was dropped by the <i>ent381</i> wednesday after an internal investigation by the <i>ent180</i> broadcaster found he had subjected producer <i>ent193</i> “ to an unprovoked physical and verbal attack . ” ...
Query Producer X will not press charges against Jeremy Clarkson, his lawyer says.	Producer X will not press charges against <i>ent212</i> , his lawyer says.
Answer Oisin Tymon	<i>ent193</i>

零代词问题消解 (Zero Pronoun Resolution)

- 零代词问题举例：
 - “我有一个苹果，非常甜。”
 - 第一步“可消解零代词识别”：“我有一个苹果，[] 非常甜。”
 - 第二步“零代词的消解”：“我有一个苹果，[苹果]非常甜。”
- 存在问题：严重缺乏训练语料
- 通过挖词构造零指代样本，成功将挖词模型应用在零指代消解任务上，方法纯净，准确率比当前最好方法提高5%

Ting Liu, Yiming Cui, Qingyu Yin, etc. ,Generating and Exploiting Large-scale Pseudo Training Data for Zero Pronoun Resolution, ACL 2017

零代词问题消解示例

(Constructing Zero Pronoun Instance)

Document:

- 1 || welcome both of you to the studio to participate in our program ,
欢迎 两位 呢 来 演播室 参与 我们的 节目 ,
- 2 || it happened that i was going to have lunch with a friend at noon .
正好 因为 我 也 和 朋友 这个 , 这个 中午 一起 吃饭 .
- 3 || after that , i received an sms from 1860 .
然后 我 就 收到 1860 的 短信 .
- 4 || uh-huh , it was by sms .
嗯 , 是 通过 短信 的 方式 ,
- 5 || uh-huh , that means , er , you knew about the accident through the source of radio station .
嗯 , 就是 说 呃 你 是 通过 台 里面 的 一个 信息 的 渠道 知道 这儿 出 了 这样 的 事故 .
- 6 || although we live in the west instead of the east part , and it did not affect us that much ,
虽然 我们 生活 在 西部 不 是 在 东部 , 对 我们 影响 不 是 很 大 ,
- 7 || but i think it is very useful to inform people using sms .
但是 呢 , 我 觉得 有 这样 一个 短信 告诉 大家 呢 是 非常 有 用 的 啊 .

Query:

- 8 || some car owners said that <blank> was very good .
有 车主 表示 , 说 这 <blank> 非常 的 好 .

Answer:

sms

短信

构造弱标注零代词问题训练数据的方法：将连续出现的两个名词中的后一个变成空槽，答案就是该词本身

两步训练过程

(A Two-step Training Mechanism)

- 两步训练
 - 大规模弱标注数据打底
 - 领域专用数据做自适应
- 首次将挖词模型应用于实际专用领域

	Kong and Zhou			Chen and Ng(2014)			Chen and Ng(2015)			Our Approach [†]		
	R	P	F	R	P	F	R	P	F	R	P	F
Overall	44.9	44.9	44.9	48.4	48.9	48.7	50.0	50.4	50.2	55.3	55.3	55.3
NW (84)	34.5	34.5	34.5	38.1	38.1	38.1	46.4	46.4	46.4	59.2	59.2	59.2
MZ (162)	32.7	32.7	32.7	30.9	31.1	31.0	38.9	39.1	39.0	51.3	51.3	51.3
WB (284)	45.4	45.4	45.4	50.4	50.4	50.4	51.8	51.8	51.8	60.5	60.5	60.5
BN (390)	51.0	51.0	51.0	45.9	45.9	45.9	53.8	53.8	53.8	53.9	53.9	53.9
BC (510)	43.5	43.5	43.5	53.8	53.8	53.8	49.4	49.4	49.4	55.5	55.5	55.5
TC (283)	48.4	48.4	48.4	53.7	56.1	54.9	52.7	52.7	52.7	52.9	52.9	52.9

Table 3: Experimental result on the test data. The strongest F-score in each row is in boldface. [†] indicates that our approach is statistical significant over the baselines (using t-test, with $p < 0.05$).

挑战2：常识知识不足

(Lack of Common Sense)

数学题：“有若干只鸡兔同在一个笼子里，从上面数，有35个头，从下面数，有94只脚。问笼中各有多少只鸡和兔？”

鸡（家禽种类） 编辑			
本词条由“科普中国”百科科学词条编写与应用工作项目 审核。			
鸡是一种家禽，家鸡源出于野生的原鸡，其驯化历史至少约4000年，但直到1800年前后鸡肉和鸡蛋才成为大量生产的商品。鸡的种类有火鸡、乌鸡、野鸡等。而且鸡也是12生肖中的一属。			
中文学名	鸡	亚 纲	今鸟亚纲
拉丁学名	Gallus gallus domesticus	目	鸡形目
界	动物界	科	雉科
门	脊索动物门	族	雉族
亚 门	脊椎动物亚门	属	原鸡属
纲	鸟纲	种	红原鸡
		亚 种	家鸡

百度百科里也无法查到“鸡有2条腿”

人：你多大了？
机：8岁
人：你结婚了？
机：结婚了，我儿子20岁了

人机对话片段

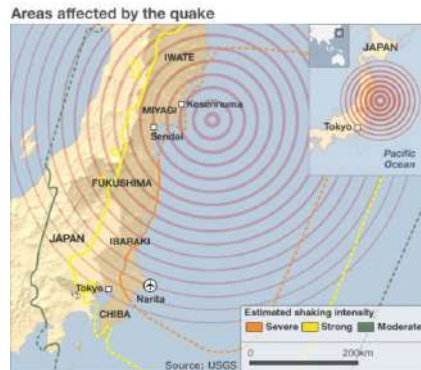
需要常识：

1. 8岁不能结婚
2. 父母比子女年龄大

事件关系可归纳：从事实到认知

(Event Relationships: from Facts to Cognition)

- 例子：A massive 8.9-magnitude *earthquake* hit northeast Japan on Friday, which cause a large amount of *houses collapsed*.
(星期五，日本东北部发生8.9级大地震，造成大量房屋倒塌。)



常识性知识:

earthquake → house collapse

事理图谱：刻画动态常识知识

(Event Evolutionary Graph: Modeling Dynamic Common Sense)

- 事理图谱：Event Evolutionary Graph (EEG)
- 事理图谱是一个事理逻辑知识库，描述事件之间的演化规律和模式
- 事理图谱是一个有向有环图，节点代表事件，有向边代表事件之间的顺承、因果关系。



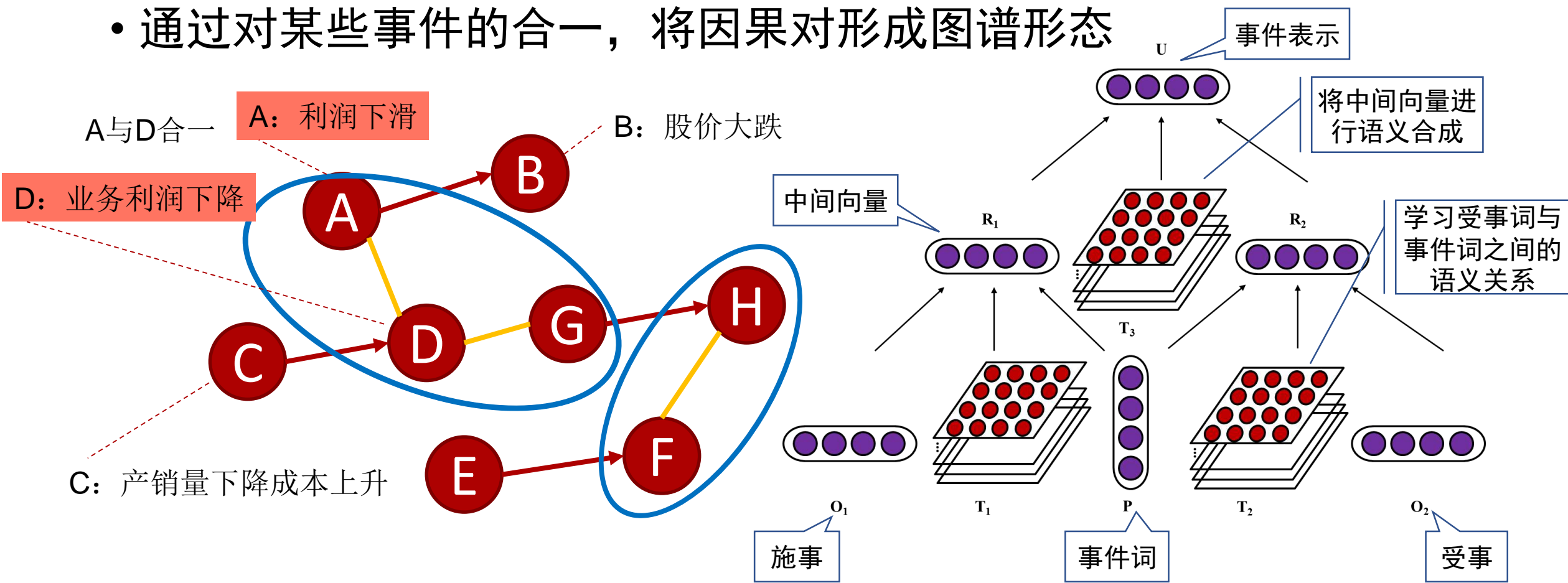
事理图谱与知识图谱的区别与联系

Differences between Event Evolutionary Graph and Knowledge Graph

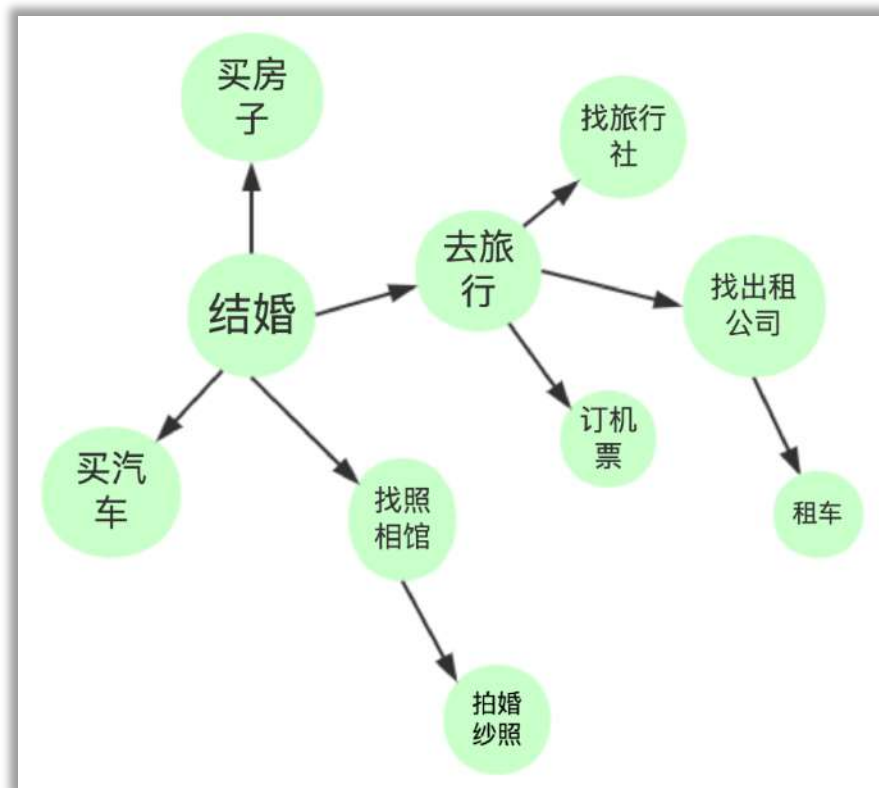
	事理图谱	知识图谱
研究对象	谓词性事件及其关系	名词性实体及其关系
组织形式	有向图	有向图
主要知识形式	事理逻辑关系，以及概率转移信息	实体属性和关系
知识的确定性	事件间的演化关系多数是不确定的	多数实体关系是确定性的

事件泛化 (Event Generalization)

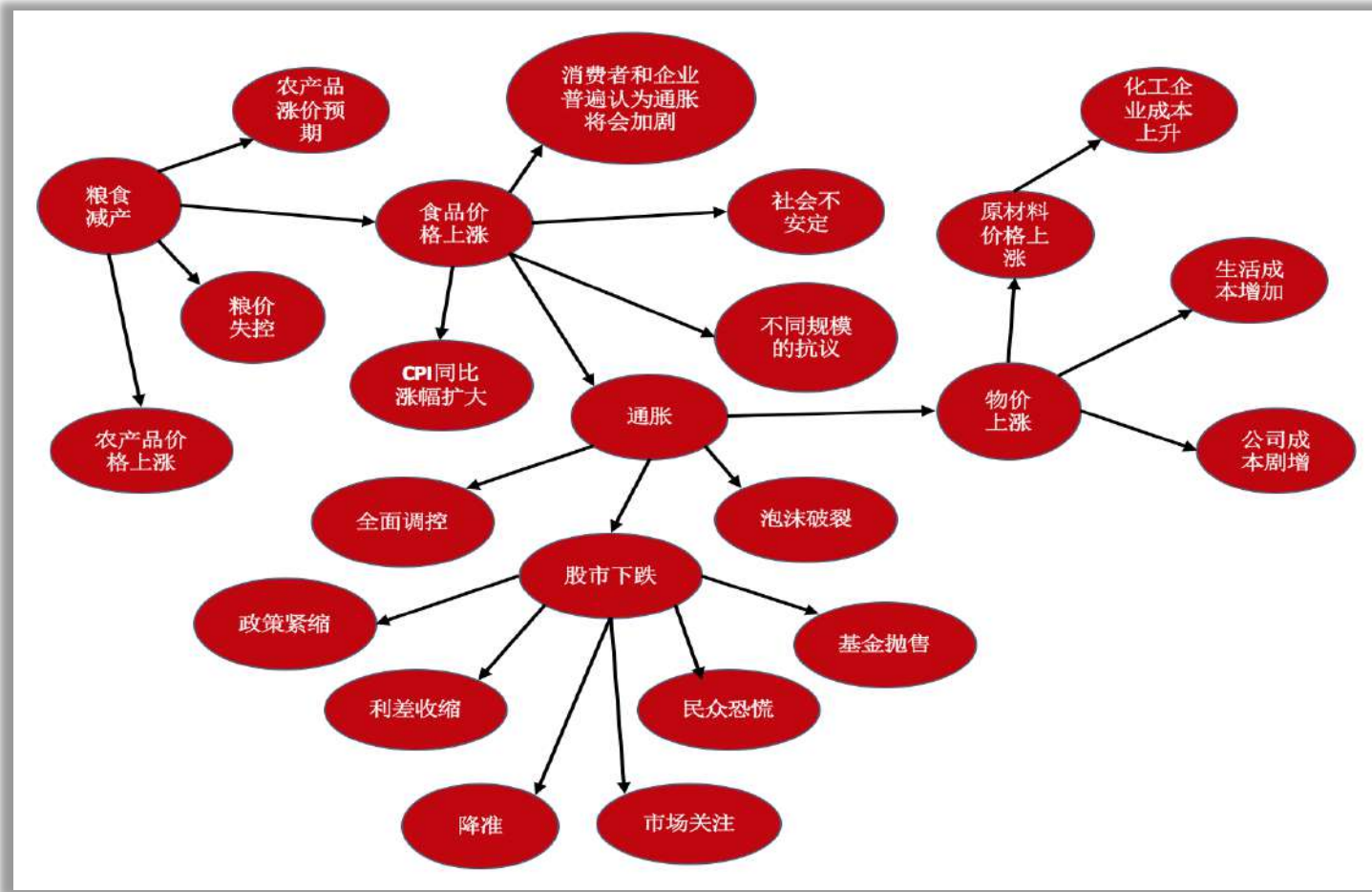
- 通过对某些事件的合一，将因果对形成图谱形态



事理图谱样例(Examples of EEG)



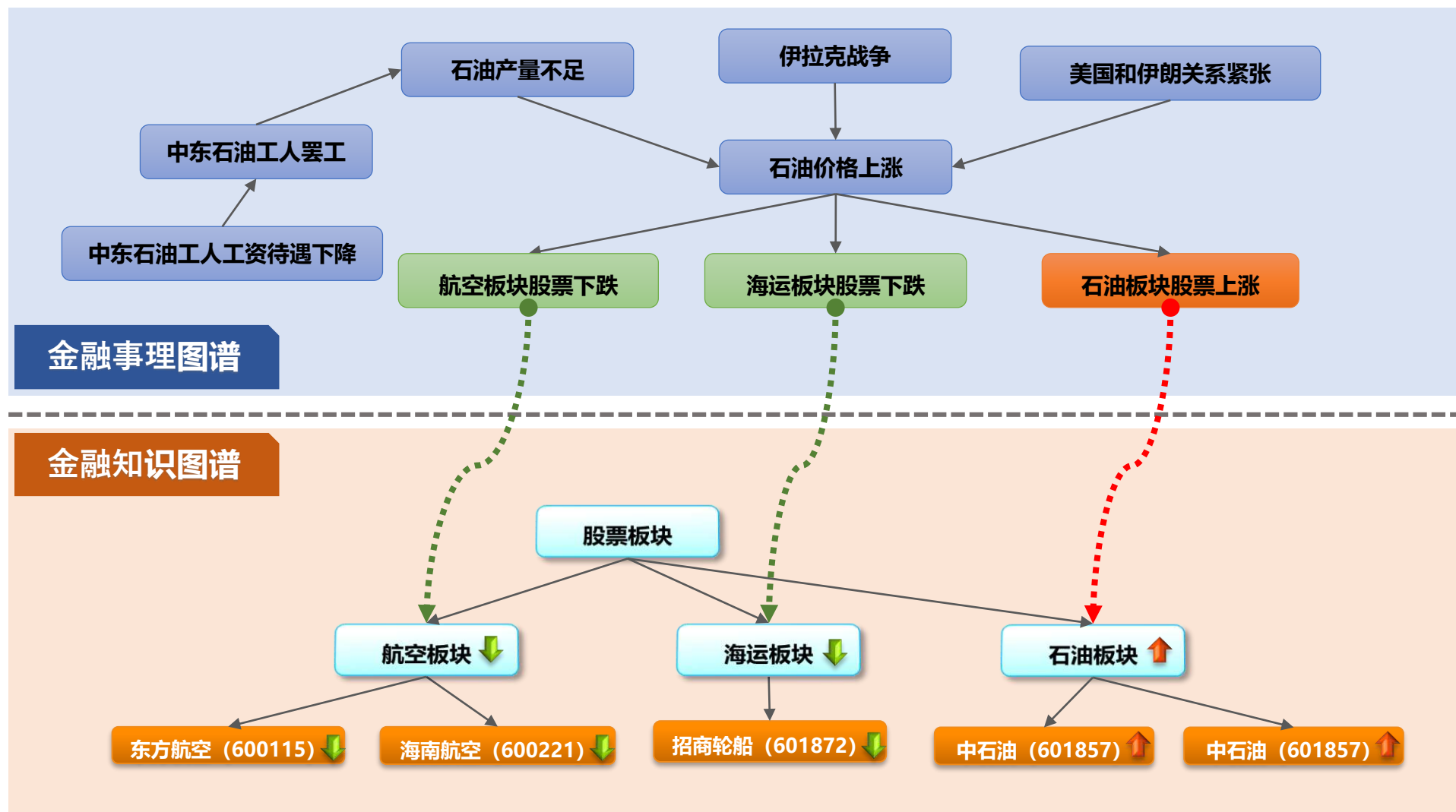
出行领域事理图谱



金融事理图谱演示地址: <http://eeg.8wss.com>

知识图谱与事理图谱的融合

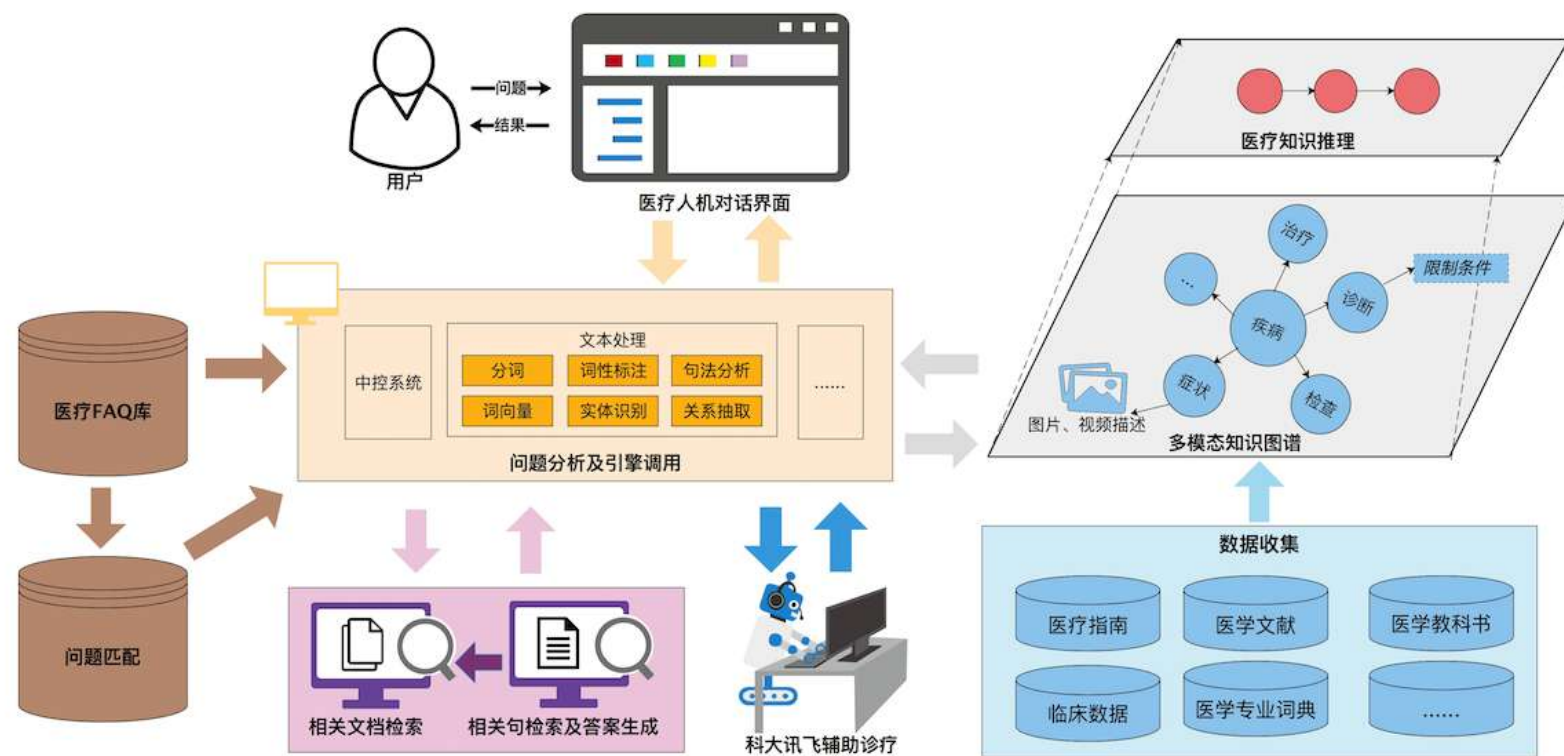
(Integration of KG and EEG)



提纲 (Outline)

- 概述
- 阶段性进展
- 技术挑战及应对之道
- 行业应用

智慧医疗 (Intelligent Healthcare)



哈工大、北大、鹏城实验室联合研制，“云知医”公益性智慧医疗开放平台

智能教育 (Smart Education)

- 哈工大讯飞联合实验室 (HFL) 研制“中文作文评阅技术”
- 作文评阅在大规模考试场景达到了人工专家水平，教学场景通过人机耦合方式可提高批改效率。

18. 作文。

我懂了

乱世烽火，焰花飞舞。乱世硝烟漫，肆意豪情。我懂了这乱世英魂，懂了你如诗似画的情义。

阿房炬·繁华逝

十八铜人仍矗立于前，阿房宫金碧辉煌，闪烁丝丝尊贵。你披长袍，骑乌骓，手上是熊熊燃烧的火把！你策转马头，仰望天空，长叹一息。那熊熊烈火如千军万马奔涌。顿时，硝烟弥漫，阿房宫在烈火中低吟，低吟繁华已逝，盛世不覆。你眼里是升腾的焰火。那焰火是阿房宫的繁华，是你的肆意豪情。

阿房炬，繁华逝，笙歌落，斗志燃。阿房宫前，我懂你，懂你乱世中的肆意豪情。

鸿门宴·天机变

帐外楚旗飞扬，帐内却是剑拔弩张的气氛。你身披黑袍，眼眸似眼凌厉的剑般犀利。你脸上已无少年时的稚嫩，征战天下战未休。你紧握酒杯，酒泉甘酿萦绕鼻间，清酿中映射项庄舞剑的意图。你眼睛环视四周，凌厉的气息弥漫。心中轻轻流动的思绪拨动心弦，心间的年少轻狂轻轻叩击着心扉。酒过三巡，刘邦望风而逃。

优美句子

评分

38

【一类卷】

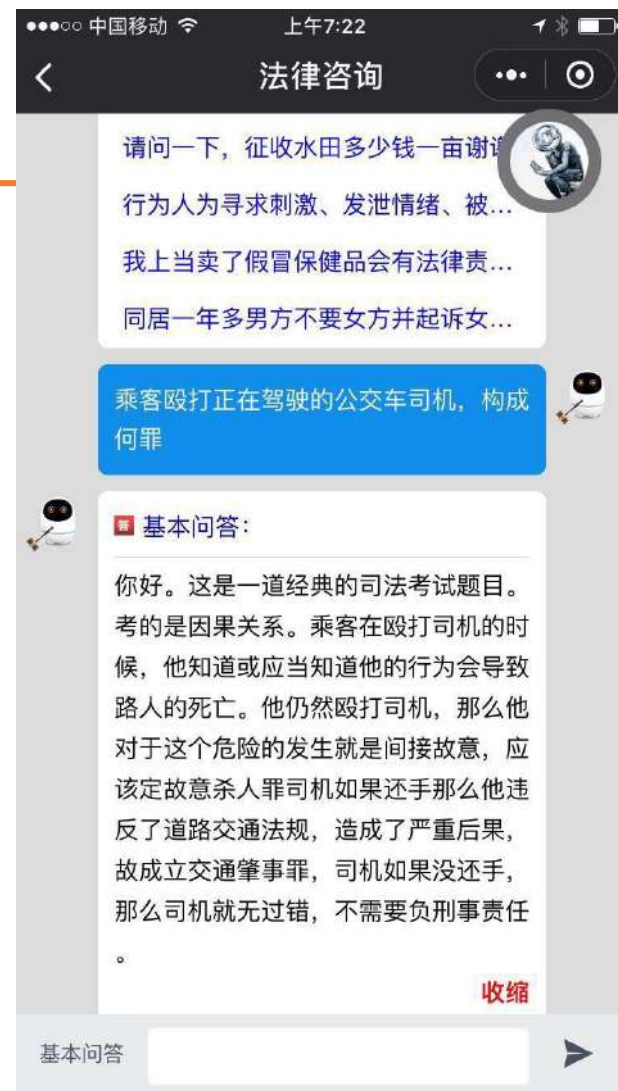
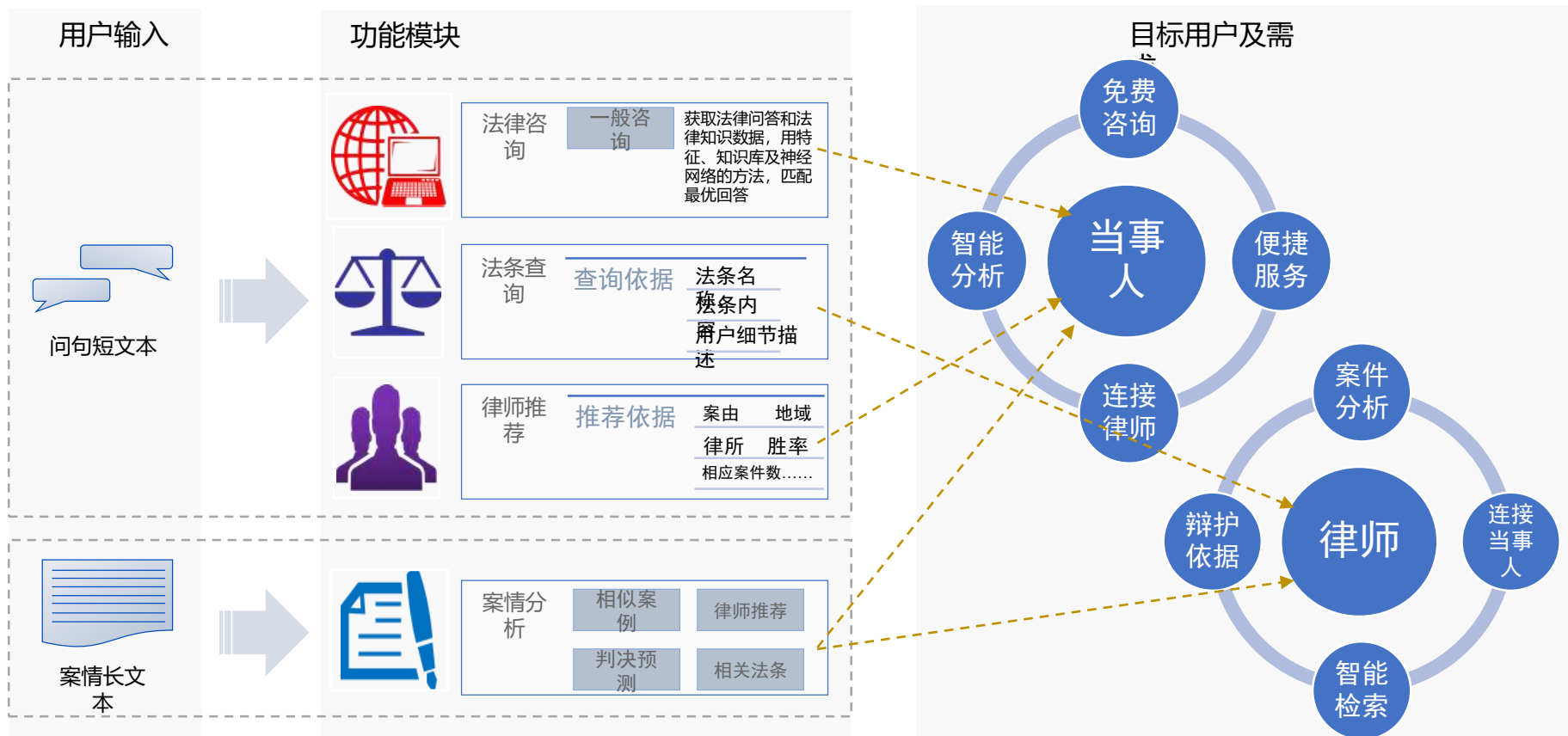
评语

内容方面：本文内容丰富、引经据典、紧扣题意。

表达方面：本文以总-分-总的结构行文，运用了小标题的形式使全文的层次分明；长短句结合运用，语言优美。

综合来看，本文是一篇结构清晰、语言优美、文采飞扬的好文章。

智慧司法 (Smart Judiciary)



“法小飞”
(哈工大讯飞联合实验室研制)

人工智能+知识服务类应用的共性问题

Common Problems of AI in Different Knowledge Services

	教育（作文）	司法	医疗
分类	几类文	什么罪	什么病
分级/评分	多少分	判几年	病情严重程度
证据	——	证词、物证	化验单
依据	评分准则	法条	医学知识、经验
案例	范文	案例	病例
交互式咨询	答疑	分诊、医疗咨询、 健康咨询	法律咨询
推荐	读什么文章	律师	医生 吃什么药
抽样质检	抽查评分结果	抽检判决结果	诊断书质量抽查

总结(Summary)

- 自然语言处理正在从语义阶段向推理阶段发展
- NLP在大数据和深度学习的有力推动下，取得了一批阶段性成果
- 带标数据不足、常识知识不足等问题困扰着NLP的发展
- NLP在向各个行业快速渗透，尤其是知识密集型行业
- 自然语言处理是人工智能皇冠上的明珠，前途无限光明！

Thanks

致谢合作者：
车万翔教授
秦兵教授等



理解语言，认知社会
以中文技术，助民族复兴



长按二维码，关注哈工大SCIR
微信号：HIT_SCIR

