



现代汉语语素系统的研发与应用

刘扬

北京大学计算语言学研究所

中文系应用语言学报告 @ 2024/4/18



报告提纲

- **1.** 回望汉语世界中的字词关系
- **2.** 汉语的语素概念表示与语义构词分析
- **3.** 汉语语素系统的应用与展望



报告提纲

- 1. 回望汉语世界中的字词关系
- 2. 汉语的语素概念表示与语义构词分析
- 3. 汉语语素系统的应用与展望



1. 回望汉语世界中的字词关系

- 关于汉语语言事实与特点的讨论
 - 汉语中的形态、句法规则较弱，属于意合型语言
 - 在语言单位形成上，“由字组词、由词造句”大体了遵循同一原则
 - 汉语的字是清晰独立的表义单位，而词的界限并不分明
 - 在词汇化历程中，字作为不同的语素出现，其语素义相对稳定
 - 因此，虽然“词无定形”，但是字和词之间具有很强的推导性
- 汉语语素与构词分析具有突出价值
 - 遵循“从结构形式到意义”的分析方法，并达成对词的理解、计算
- 汉语语素与构词分析相关资源建设
 - 苑春法：汉语语素数据库
 - 周亚民：汉字知识本体（Hantology）
 - 董振东：知网（HowNet）



报告提纲

- **1.** 回望汉语世界中的字词关系
- **2.** 汉语的语素概念表示与语义构词分析
- **3.** 汉语语素系统的应用与展望

2. 汉语的语素概念表示与语义构词分析

■ 以《现代汉语词典》刻画的汉语语素及语素义为依据

■ 语素是语言中最小的语法单位，也就是最小的语音、语义结合体

选 (選) xuǎn ① **动** 挑选：筛~ | ~拔 | ~派 | ~种。② **动** 选举：~民 | 普~ | ~代表。③ 被选中了的(人或物)：入~ | 人~。④ 挑选出来编在一起的作品：文~ | 诗~ | 民歌~。

挑 ¹ tiāo **动** ① 挑选：~心爱的买。② 挑剔：~毛病。

挑 ² tiāo ① **动** 扁担等两头挂上东西，用肩膀支起来搬运：~担 | ~水 | ~着两筐土。② (~儿) **名** 挑子：挑~儿。③ (~儿) **量** 用于成挑儿的东西：一~儿白菜。

另见 1354 页 tiǎo。

材 cái ① 木料，泛指材料①：木~ | 钢~ | 药~ | 就地取~。② **名** 棺材：寿~ | 一口~。③ 资料：教~ | 题~ | 素~。④ 有才能的人。⑤ (cái) **名** 姓。

才 ¹ cái ① **名** 才能：德~兼备 | 多~多艺 | 这人很有~。② 有才能的人：干~ | 奇~。③ (cái) **名** 姓。

才 ² (纔) cái **副** ① 表示以前不久：你怎么~来就要走？② 表示事情发生得晚或结束得晚：他说星期三动身，到星期五~走 | 大风到晚上~住了。③ 表示只有在某种条件下然后怎样(前面常常用“只有、必须”或含有这类意思)：只有依靠群众，~能把工作做好。④ 表示发生新情况，本来并不如此：经他解释之后，我~明白是怎么回事。⑤ 表示数量小，次数少，能力差，程度低等等：这个工厂开办时~几十个工人 | 别人一天干的活儿，他三天~干完。⑥ 表示强调所说的事(句尾常用“呢”字)：麦子长得~好呢 | 我~不信呢！

2. 汉语的语素概念表示与语义构词分析

- 汉语的语素概念表示工程
 - 工程目标 1: 确认“语素义编码”
 - 收集汉字的各个语素义，明确并补齐其所属的语素类（POS），并赋予其唯一的“语素义编码”
 - 工程目标 2: 提取“语素概念”
 - 对语素的释文句子，借助语义相似度计算，形成各个“同义（或同类）语素集”，用以表征“语素概念”，或称“语义基元”
 - 工程目标 3: 构建“语素概念体系”
 - 对现代汉语的全部“语素概念”，构建“语素概念体系”
 - 工程目标 4: 建立“语义构词描述”
 - 通过对构词结构标注及其下的“语素概念”绑定，建立基于“语素概念体系”的“语义构词描述”

2. 汉语的语素概念表示与语义构词分析

- 当前工程进展：目标 1 的阶段性数据成果
 - 获得 8514 个汉字（TYPE）14 个语素类（POS）下的 20855 个语素义（SENSE），并按规则给每个语素义赋予唯一的“语素义编码”
 - 这其中，成词的自由语素（free morpheme）和不成词的黏着语素（bound morpheme）大致各占一半
 - “选”字的各个“语素义编码”：
 - 选1_04_01: 动词 挑选：筛~ | ~拔 | ~派 | ~种
 - 选1_04_02: 动词 选举：~民 | 普~ | ~代表
 - 选1_04_03: 名语素 被选中了的(人或物)：入~ | 人~
 - 选1_04_04: 名语素 挑选出来编在一起的作品：文~ | 诗~ | 民歌~
 - “材”字的各个“语素义编码”：
 - 材1_05_01: 名语素 木料，泛指材料：木~ | 钢~ | 药~ | 就地取~
 - 材1_05_02: 名语 棺材：寿~ | 一口~
 - 材1_05_03: 名语素 资料：教~ | 题~ | 素~
 - 材1_05_04: 名语素 有才能的人
 - 材1_05_05: 名语 (Cái)姓

2. 汉语的语素概念表示与语义构词分析

- 当前工程进展：目标 2 的阶段性的数据成果
 - 作为主体的汉语的名、动、形语素分别聚合形成了2018、1630、549个相对独立、封闭的“语素概念”
 - 这构成了汉语世界中操纵计算的基本意义单元（即“语义基元”）
 - 数据成果具备优良特性：明确、可判的覆盖度与准确性
 - 形成了用于表征、操作汉语词型、词义的新的数据基础
 - 动语素“语素概念”举例，某概念 X：“选择、挑选”
 - { 刷3_01_01, 抡1_01_01, 拔1_08_03, 拣1_01_01, 择1_02_01, 择2_02_01, 挑1_02_01, 擢1_02_02, 调4_02_02, 选1_04_01, 遴1_01_01, 铨1_02_01 }
 - 名语素“语素概念”举例，某概念 Y：“有才能的人”
 - { 匠1_02_02, 哲1_02_02, 器1_05_04, 尖1_09_06, 彦1_02_01, 才1_03_02, 材1_05_04, 杰1_03_01, 模1_04_03, 氏1_05_03, 秀2_04_04, 英1_03_02, 豪1_04_01, 贤1_04_02, 通2_12_07, 骥1_02_02 }

2. 汉语的语素概念表示与语义构词分析

名语素“语素概念”覆盖、分布情况

集合字数	概念个数	概念比例	概念示例	概念说明
1866	1	0.05	丁七万三上下与丐丑专且世丘丙业从东丞两严...	姓氏
261	1	0.05	丽毫囊令任侯侏儂克兹匠单屋厦吓吴喷扮劫坻...	地域的简称、别称
136	1	0.05	匏柘棘稂稊篇篇缩芳艾芴芴芴芴芴芴芴芴...	草本植物
123	1	0.05	们刺妨葵岷建桂汉汜汜汜汜汜汜汜汜汜汜汜汜...	水域、河流名称
118	1	0.05	鸟鳧鳳凰鷓鴣鸂鶒鸂鶒鸂鶒鸂鶒鸂鶒鸂鶒...	鸟类
113	1	0.05	朴杉枹枹枹枹枹枹枹枹枹枹枹枹枹枹枹枹...	木本植物
110	1	0.05	蛇鳗鳗鳗鳗鳗鳗鳗鳗鳗鳗鳗鳗鳗鳗鳗鳗...	鱼类
91	1	0.05	且偃仙佺佺佺佺佺佺佺佺佺佺佺佺佺佺佺...	人名用字
89	1	0.05	蝦拉享莛蚶蚶蚶蚶蚶蚶蚶蚶蚶蚶蚶蚶蚶...	节肢动物
83	1	0.05	金鍮鍮钋钋钋钋钋钋钋钋钋钋钋钋钋钋钋...	金属
71	1	0.05	元兕冉南卫召吴周唐商夏斌宋巴明晋曹朝杞梁...	朝代或国家名
63	1	0.05	吁吗叨吠哞哞哞哞哞哞哞哞哞哞哞哞哞...	制药的有机化合物
53	1	0.05	玉玊瑪瑒玢玢玢玢玢玢玢玢玢玢玢玢玢...	各种玉或玉器
52	2	0.10	卜椿瓜瓠笋芋茨芥芫芫芫芫芫芫芫芫...	植物中的蔬菜类
41	1	0.05	华台圃麥研研研研研研研研研研研研研研...	山名
40	1	0.05	亿佻佻佻佻佻佻佻佻佻佻佻佻佻佻佻佻...	少数民族
39	1	0.05	騾駝駝駝駝駝駝駝駝駝駝駝駝駝駝駝...	马类
37	2	0.10	呢哔布帛彩纛纛纛纛纛纛纛纛纛纛纛纛...	丝织品
36	2	0.10	云京冀台川广晋桂楚沪豫浙渝港湖湘滇澳琼中...	行政区划的省级单位
28	1	0.05	匏李杏杞杞杞杞杞杞杞杞杞杞杞杞杞杞杞...	植物的果实
27	3	0.15	丞令侯倭倭倭倭倭倭倭倭倭倭倭倭倭倭...	古代官名
26	3	0.15	甬庖登尊彝斗罍樽爵甗甗甗甗甗甗甗甗...	与酒有关的器皿
24	1	0.05	冈坂阪场埜埜埜埜埜埜埜埜埜埜埜埜埜...	山体的某一部分
23	1	0.05	禾秫秬稷稷稷稷稷稷稷稷稷稷稷稷稷...	粮食作物
22	1	0.05	治府寺监部馆所厅院科局股处室课段署家司部...	政府机关部门
21	1	0.05	陶席笆笊篋筍筍筍筍筍筍筍筍筍筍筍筍...	竹或木制成的容器
20	3	0.15	酎酎酒醖醖醖醖醖醖醖醖醖醖醖醖醖醖...	酒类
19	5	0.25	帘帟幃幃幃幃幃幃幃幃幃幃幃幃幃幃...	旗帜
18	4	0.20	内心牙肝肠肺肾胃胆胰胱膀胱脾脾脾脾...	动物内脏
17	4	0.20	土地垌垌垌垌垌垌垌垌垌垌垌垌垌垌垌垌...	土地田地
16	7	0.35	匠哲器彦才材杰氏秀英豪贤通骥模尖	有才能的人
15	6	0.30	体例则制宪彝律法矩科桌规贯轨辟	法律、规章制度
14	11	0.55	刀戈戟戣枪架矛鏃鏃鏃鏃鏃鏃鏃鏃...	长矛类兵器
13	9	0.45	丫栝本杈条杪枝柯标株椹椹椹	树枝
12	15	0.74	丘冢圻坟墓塋塋穴窆莹阡陵	坟墓
11	14	0.69	上储后君圣帝庙王皇辟驾	古代帝王
10	24	1.19	位品地档流等级职衔阶	社会地位、等级
9	19	0.94	粇糕酥饼饼饼饼饼饼	糕点
8	37	1.83	劄印戳玺章篆铃鼎	印章、图章
7	50	2.48	唾喇沫津涎酸痰	唾液
6	81	4.01	舟舡航舫船舫	船只
5	99	4.91	志标符记识	标志
4	156	7.73	楊箠鞘鞭	用于鞭打的工具
3	212	10.51	仲孟季	兄弟排行位置
2	405	20.07	伞盖	用于挡雨、遮阳的工具
1	821	40.68	驮	牲口驮着的货物
不限	2018	100	N/A	N/A

2. 汉语的语素概念表示与语义构词分析

- 当前工程进展：目标 3 的阶段性数据成果
 - 以面向对象思想、生成词库理论等为指导，构建汉语的“语素概念体系”（字义系统）
 - 汉语的名、动、形语素的层次结构是大致同构、映射的，形成同语素类内的聚合关系以及跨语素类间的组合关系，以方便认知、计算和推理
 - 在多种特征赋值基础上，借助“特征序列”的设定实现多分类树功能
 - 陈刚，刘扬. 基于特征序列的语义分类体系的自动构建. 中文信息学报, 2015, 29 (3): 52-57



2. 汉语的语素概念表示与语义构词分析

- 当前工程进展：目标 4 的阶段性的数据成果
 - “语义构词描述” 1: “构词结构”标注，以“选材”为例
 - 语法结构标签：述宾，语义结构标签：受事 [已融合了：对象]
 - “语义构词描述” 2: “语素概念”绑定，以“选材”为例
 - “选”字的语素义为“选择、挑选”，语素义编码为“选1_04_01”
 - 绑定了“语素概念” X:
 - {刷3_01_01, 抡1_01_01, 拔1_08_03, 拣1_01_01, 择1_02_01, 择2_02_01, 挑1_02_01, 擢1_02_02, 调4_02_02, 选1_04_01, 遴1_01_01, 铨1_02_01}
 - “材”字的语素义为“有才能的人”，语素义编码为“材1_05_04”
 - 绑定了“语素概念” Y:
 - {匠1_02_02, 哲1_02_02, 器1_05_04, 尖1_09_06, 彦1_02_01, 才1_03_02, 材1_05_04, 杰1_03_01, 模1_04_03, 氏1_05_03, 秀2_04_04, 英1_03_02, 豪1_04_01, 贤1_04_02, 通2_12_07, 骥1_02_02}



报告提纲

- **1. 回望汉语世界中的字词关系**
- **2. 汉语的语素概念表示与语义构词分析**
- **3. 汉语语素系统的应用与展望**

3. 汉语语素系统的应用与展望

- 这些基础性研发工作有望推动 人文领域 和 计算应用 的开展
 - 充分考虑汉语的语言事实与特点，从汉字及其基本意义入手
 - 刘扬, 林子, 康司辰. 汉语的语素概念提取与语义构词分析. 中文信息学报, 2018, 32(2): 12-21
- 一、在 人文领域 方面的应用尝试与探索方向
 - 1. 支持新型汉语电子词典的编撰、出版与浏览使用
 - 字、词意义空间划分、对应的验证与考核，严格的可计算化
 - 不必再依赖字形或拼音了，以电子数据方式，让基于“语素概念”的查询成为可能，并快捷遍历所有（形式或）意义相关的字、词
 - 2. 支持汉语字、词的母语教学以及对外汉语教学
 - 网站查询与浏览服务
 - 手机微信小程序应用
 - 作为插件，对第三方软件（涉及汉语学习理解）的广泛的、无缝的支持

3. 汉语语素系统的应用与展望

■ 一、在人文领域方面的应用尝试与探索方向

■ 3. 支持更加细致、深入的汉语语言认知研究

■ 汉语的字义演变与词汇化过程的解释研究

- 康司辰, 虞梦夏, 刘扬. 基于平行周遍原则的汉语未登录词的知识表示与预测. 中文信息学报, 2020, 34 (8): 23-31

■ 基于语言认知(隐喻、转喻等)理解的实证研究

- 陈龙, 饶琪, 刘扬. 汉语词的非字面义表示与应用. 中国科学: 信息科学, 2019, 49 (8): 1005-1018

■ 汉语基本情感与复合情感的表征、分析与计算

■ 汉语的语义本体构建与应用评估

■ 二、在计算应用方面的应用尝试与探索方向

■ 1. 面向理解的汉语未登录词的词义知识表示及语义预测

- 田元贺, 刘扬. 汉语未登录词的词义知识表示及语义预测. 中文信息学报, 2016, 30 (6): 26-34



3. 汉语语素系统的应用与展望

■ 二、在计算应用方面的尝试方向

■ 2. 基于语素概念与语义构词分析等汉语特性的认知计算

- 康司辰, 刘扬. 基于语义构词的汉语词语语义相似度计算. 中文信息学报, 2017, 31 (1): 94-101
- 康司辰, 虞梦夏, 刘扬. 基于平行周遍原则的汉语未登录词的知识表示与预测. 中文信息学报, 2020, 34 (8): 23-31
- 郑嫫, 刘扬, 殷雅琦, 王悦, 代达励. 基于词信息嵌入的汉语构词结构识别研究. 中文信息学报, 2022, 36 (5): 31-40, 66 (第二十届中国计算语言学大会 CCL2021, 会议最佳论文)
- 王悦, 刘扬, 梁启亮, 王涵思. 汉语语义构词的资源建设与计算评估. 语言文字应用, 2023(4): 105-117

■ 3. 汉语等语言中的语素义、词义的表达、应用与可解释性研究

- Lin Zi, Yang Liu. Implanting Rational Knowledge into Distributed Representation at Morpheme Level. AAAI 2019



3. 汉语语素系统的应用与展望

■ 二、在计算应用方面的尝试方向

- 3. 汉语等语言中的语素义、词义的表达、应用与可解释性研究
 - Hua Zheng, Damai Dai, Lei Li, Tianyu Liu, Zhifang Sui, Baobao Chang, Yang Liu. Decompose, Fuse and Generate A Formation-Informed Method for Chinese Definition Generation. NAACL 2021
 - Hua Zheng, Lei Li, Damai Dai, Deli Chen, Tianyu Liu, Xu Sun, Yang Liu. Leveraging Word-Formation Knowledge for Chinese Word Sense Disambiguation. EMNLP Findings 2021
 - Yaqi Yin, Yue Wang, Yang Liu. Chinese Morpheme-informed Evaluation of Large Language Models. COLING 2024
 - Yue Wang, Hua Zheng, Yaqi Yin, Hansi Wang, Qiliang Liang, Yang Liu. Morpheme Sense Disambiguation: A New Task Aiming for Understanding the Language at Character Level. COLING 2024
 - Yue Wang, Qiliang Liang, Yaqi Yin, Hansi Wang, Yang Liu. Disambiguate Words like Composing Them: A Morpheme-Informed Approach to Enhance Chinese Word Sense Disambiguation. ACL 2024 (to appear)



谢谢各位，请批评指正！

liuyang@pku.edu.cn