

计算语言学

第3讲 词法分析（一）

刘群

中国科学院计算技术研究所

liuqun@ict.ac.cn

中国科学院研究生院2010年春季课程讲义

内容提要



内容提要

词法分析问题 ↗

课程的组织

问题与方法

按问题组织课程

语言的形态

英语词法分析

tokenization

stemming

POS-tagging

汉语词法分析

重叠词

前后缀

离合词

切分歧义

未定义词识别

词性标注

问题与方法

- 计算语言学主要问题：
 - 机器翻译
 - 自动问答
 - 音字转换
 - 自动文摘
 - 信息抽取
 -

问题与方法

- 计算语言学问题的抽象：
 - 序列评估问题
 - 序列标注问题
 - 序列结构化问题
 - 序列转换问题

序列评估问题

- 输入：一个符号序列
- 输出：
 - 合法性评估：是否合法
 - 可能性评估：概率值
- 常见具体问题：
 - 文本校对
 - 汉语词语切分、音字转换等很多问题都可以转化成序列评估问题

序列标注问题

- 输入：一个符号序列
- 输出：给每一个输入符号赋予一个标记
- 常见具体问题：
 - 音字转换：拼音序列 → 汉字序列
 - 词性标注：词语序列 → 词性序列
 - 词义排歧：词语序列 → 词义标记序列

序列结构化

- 输入：一个符号序列
- 输出：一个结构，刻画符号之间的关系
- 常见具体问题：
 - 成分句法分析：词语序列 → 短语结构树
 - 依存句法分析：词语序列 → 依存树
 - 语义分析：词语序列 → 语义网络

问题与方法

- 计算语言学常用方法:

- 规则方法

- 形式语法理论
 - 形式逻辑
 -

- 统计方法

- n元语法模型
 - 隐马尔科夫模型
 - 最大熵模型
 -

本课程的组织

- 按问题组织
 - 词法分析
 - 句法分析
 - 语义篇章分析
 - 机器翻译
 -

本课程的组织

- 问题中穿插方法的介绍
 - 词法分析：
 - 语言模型、HMM模型、最大熵.....
 - 句法分析：
 - 概率语法.....
 - 机器翻译：
 -

内容提要

词法分析问题 ↗

课程的组织

问题与方法

按问题组织课程

语言的形态

英语词法分析

tokenization

stemming

POS-tagging

汉语词法分析

重叠词

前后缀

离合词

切分歧义

未定义词识别

词性标注

语言的形态

- 形态: Morphology
 - The study of the internal structure of words, and of the rules by which words are formed, is called morphology.
(from V. Fromkin & R. Roman: An Introduction to Language)
 - 单词的内部结构的研究, 以及单词形成的规律, 被称之为形态学。
(选自V. Fromkin & R. Roman: 语言介绍)

语言的形态

- 形态：又叫词形变化，同一个词在造句时，因其句法位置的差异而发生的不同变化，是表达语法意义的重要手段。这些不同的变化形成一个聚合。包括词尾，内部曲折，异根等方面。
- 语法范畴：词的变化形式所表示的意义方面的聚合。常见的语法范畴有：性、数、格、体、时态、人称、级等。
- 形态跟语法范畴有对应关系，但不是一回事。

语言的分类

传统语言学根据词的形态把语言分为四大类：

- 分析语：每个词只有一个词素
 - 孤立语（词根语）：词基本上没有专门表示语法意义的附加成分，形态变化很少，语法关系靠词序和虚词来表示。如汉语。
- 综合语：每个词有多个词素
 - 黏着语：词内有专门表示语法意义的附加成分，一个附加成分表达一种语法意义，一种语法意义也基本上由一个附加成分来表达，词根或词干跟附加成分的结合不紧密。如芬兰语、日语、蒙古语等。
 - 屈折语：用词的形态变化表示语法关系，一个形态成分可以表示若干种不同的语法意义，词根或词干跟词的附加成分结合得很紧密，往往不易截然分开。如：英语、德语和法语等。
 - 多式综合语（编插语）：最复杂的综合语，一个词语通常非常长，由很多词素组成。

内容提要

词法分析问题 ↗

课程的组织

问题与方法

按问题组织课程

语言的形态

英语词法分析

tokenization

stemming

POS-tagging

汉语词法分析

重叠词

前后缀

离合词

切分歧义

未定义词识别

词性标注

屈折型语言的词法分析

- **Tokenization:** 把字符串变成词串 (tokens)
I'm a student. → I 'm a student .
- **Lemmatization:** 对词的内部结构进行分析, 并还原到词典形式 (lemma)。主要是对屈折进行还原。
takes → take + ~s
took → take + ~ed
- **Stemming:** 对词的内部结构进行分析, 不仅对曲折进行还原, 而且也对派生进行还原。
tokenization → token + ~ize + ~tion
- Stemming有时等同Lemmatization。
- **POS-Tagging:** 词性标注

Tokenization

- 数字：123,456.78 90.7% 3/8 11/20/2000
- 缩略（包含不同的情况）：
 - 字母一点号一字母一点号组成的序列，比如：U.S. i.e. 等等；
 - 字母开头，最后以点号结束，比如：A. b. Mr. eds.prof. ；
- 包含非字母字符，比如：AT&T Micro\$oft
- 带杠的词串，比如：three-year-old, one-third, so-called
- 带撇号的词串，比如：I'm can't dog's let's
- 带空格的词串，比如："and so on", "ad hoc"
- 其他：如网址（<http://ict.ac.cn>）、公式等

Tokenization问题

- 例外较多，跟文本来源有关
- 歧义现象（如点号的句子边界歧义）

数字的识别

数词的识别一般可以用有限状态自动机来实现

- 识别分数的正则表达式：
 - $[0-9]^+ / [0-9]^+$
 - e.g. 12/21
- 识别百分数的正则表达式：
 - $([+ | -]) ? [0-9]^+ (\cdot [0-9]^*) ? \%$
 - e.g. -5.9% 91%
- 识别十进制数字的正则表达式：
 - $([0-9]^+ (,) ?)^+ (\cdot [0-9]^+) ?$
 - e.g. 12,345

Tokenization算法

- 输入：一段文本
- 输出：单词串
- 算法：（略）

Stemming

屈折型语言的词语变化形式：

- 屈折变化：即由于单词在句子中所起的语法作用的不同而发生的词的形态变化，而单词的词性基本不变的现象，如（take, took, takes）。识别这种变化是词法分析的最基本的任务。
- 派生变化：即一个单词从另外一个不同类单词或词干衍生过来，如morphological ← morphology，英语中派生变化主要通过加前缀或后缀的形式构成；在其他语言中，如德语和俄语中，同时还伴有音的变化。
- 复合变化：两个或更多个单词以一定的方式组合成一个新的单词。这种变化形式比较灵活，如well-formed, 6-year-old等等。

Stemming的目的：将上述变化还原

Stemming常见的问题

- 半规则变化
 - flied → fly + ~ed
 - rebelled → rebel + ~ed
- 不规则变化
 - good, better, best
 - child, children
- 歧义现象
 - better → good + ~er or well + ~er ?
 - works → work + ~s or works ?

Stemming规则示例 (1)

- 名词复数

***s → *, (PLUR)**

***es → *, (PLUR)**

***ies → *y, (PLUR)**

- 动词第三人称单数

***s → * (SINGULAR) (THIRDPERSON) (PRESENT)**

***es → * (SINGULAR) (THIRDPERSON) (PRESENT)**

***ies → *y (SINGULAR) (THIRDPERSON) (PRESENT)**

Stemming规则示例 (2)

- 动词现在分词
 - *ing → * (VING)
 - *ing → *e (VING)
 - *ying → *ie (VING)
 - *??ing → *? (VING)
- 动词过去分词、过去式
 - *ed → * (PAST, VEN)
 - *ed → *e (PAST, VEN)
 - *ied → *y (PAST, VEN)
 - *??ed → *? (PAST, VEN)

Stemming算法

- 输入：一个单词
- 输出：一个或多个单词，其中每个单词还原为原形加前后缀（可以有多个）
- 算法：（略）

基于有限状态自动机的Stemming

- 有限状态自动机是Stemming中的常用算法
- 有限状态自动机的优点是表现形式直观，效率高

Stemming要做到何种程度

- 词干层。如：
impossibilities → impossibility+ies
- 词根层。如：
impossibilities → im+poss+ibil+it+ies
- 分析程度取决于自然语言处理系统的深度：
 - 不解决未定义词，分析到词干层
 - 解决未定义词，要分析到词根层。

内容提要

词法分析问题 ↗

课程的组织

问题与方法

按问题组织课程

语言的形态

英语词法分析

tokenization

stemming

POS-tagging

汉语词法分析

重叠词

前后缀

离合词

切分歧义

未定义词识别

词性标注

汉语词法分析所面临的问题

- 重叠词、离合词、词缀
- 汉语词语的切分歧义
- 汉语未定义词
- 词性标注

汉语双字形容词的重叠形式

形容词(AB)	ABAB式	AABB式	A里AB式
高兴	高兴高兴	高高兴兴	
明白	明白明白	明明白白	
热闹	热闹热闹	热热闹闹	
潇洒	潇洒潇洒	潇潇洒洒	
糊涂		糊糊涂涂	糊里糊涂
流气			流里流气
粘乎	粘乎粘乎	粘粘乎乎	
凉快	凉快凉快	凉凉快快	

汉语单字形容词的重叠形式

形容词（A）	AA式	ABB式	ABCD式
黑	黑黑	黑压压	黑不溜秋
白	白白	白花花	白不毗咧
红	红红	红彤彤	
亮	亮亮	亮晶晶	
恶		恶狠狠	
香	香香	香喷喷	
滑	滑滑	滑溜溜	

汉语双字动词的重叠形式

动词(AB)	ABAB式	AABB式
研究	研究研究	
讨论	讨论讨论	
哆嗦		哆哆嗦嗦
唠叨	唠叨唠叨	唠唠叨叨
嘀咕		嘀嘀咕咕

汉语单字动词的重叠形式

动词（V）	VV式	V—V式	V了V式	V了一V式
听	听听	听一听	听了听	听了一听
想	想想	想一想	想了想	想了一想
玩	玩玩	玩一玩	玩了玩	玩了一玩
醒	醒醒	醒一醒		
试	试试	试一试	试了试	试了一试
笑	笑笑	笑一笑	笑了笑	笑了一笑
讲	讲讲	讲一讲	讲了讲	讲了一讲

汉语其他词类的重叠形式

- 名词
 - 哥哥，人人
 - 山山水水，是是非非，方方面面，头头脑脑
- 数词
 - 一一做了回答，两两结伴而来
- 量词
 - 个个都是好样的，回回考满分
- 副词
 - 常常，仅仅，的确确

汉语重叠词的特点

- 汉语词能否重叠具有很强的个性特点
 - 研究研究 ✓
 - 工作工作 ✕
- 有些词重叠后词性发生了变化
 - 形容词重叠后一般成为状态词
 - 个别量词重叠后可以成为其他词性
 - 回回：副词
 - 个个：名词

汉语词缀

- 前缀
 - 老鹰、老虎、老三、老王
 - 超豪华、超标准、超高速
 - 非党员
- 后缀
 - 骨头、砖头、甜头、苦头、盼头、想头
 - 桌子、椅子、孩子、票子、房子
 - 文学家、指挥家、艺术家
 - 科学性、可能性、学术性
 - 碗儿、花儿、玩儿、份儿、片儿

汉语离合词

- 汉语动词存在离合词现象
 - 游泳：游了一会儿泳
 - 理发：发理了没有
 - 担心：担什么心
 - 洗澡：洗了个热水澡
- 白硕的解释：语义重心偏移
 - 动词虚化（类似英语DO）
 - 语义重心落在后面的名词性语素上
 - 游泳：游了一会儿泳：DO了一会儿游泳
 - 理发：发理了没有：理发DO了没有
 - 担心：担什么心：DO什么担心
 - 洗澡：洗了个热水澡：DO了一个热水洗澡

处理识别词形变化的规则 (1)

- @ @ VV -- v [重叠形式:VV] << V -- v
- @ @ UVUV -- v [重叠形式:UVUV] << UV -- v
- @ @ V了V -- v [重叠形式:V了V] << V -- v
- @ @ V—V -- v [重叠形式:V—V] << V -- v
- @ @ V了N -- v [重叠形式:V了N] << VN -- v [趋向动词:否]
- @ @ VVN -- v [重叠形式:VVN] << VN -- v
- @ @ V过N -- v [重叠形式:V过N] << VN -- v [趋向动词:否]
- @ @ V了一N -- v [重叠形式:V了一N] << VN -- v
- @ @ V过一N -- v [重叠形式:V过一N] << VN -- v
- @ @ V不了N -- v [重叠形式:V不了N] << VN -- v

处理识别词形变化的规则 (2)

@ @ AA -- a [重叠形式:AA] << A -- a

@ @ AABB -- a [重叠形式:AABB] << AB -- a

@ @ ABAB -- a [重叠形式:ABAB] << AB -- a

@ @ DD -- d [重叠形式:DD] << D -- d

@ @ R俩 -- r << R们 -- r

@ @ N儿 -- n [重叠形式:N儿] << N -- n

@ @ MN儿 -- n [重叠形式:MN儿] << MN -- n

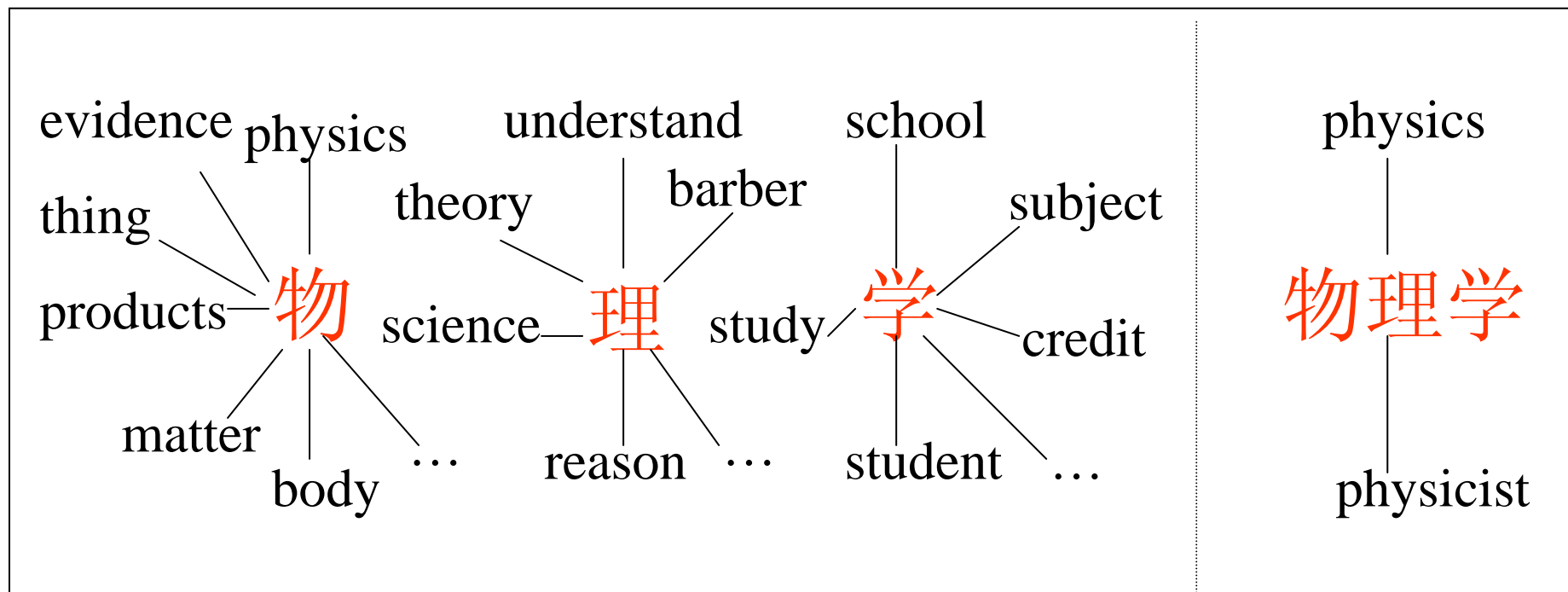
汉语的词语切分

物理学很有趣
(Physics is very interesting)

**Chinese Word
Segmentation**

物理学 / 很 / 有趣
Physics Very Interesting

为什么要做词语切分



为什么要做词语切分

物理学很有趣

(Physics is very interesting)



物理学

/

很

/

有趣

Physics

Very

Interesting



物理

/

学

/

很

/

有趣

Physics

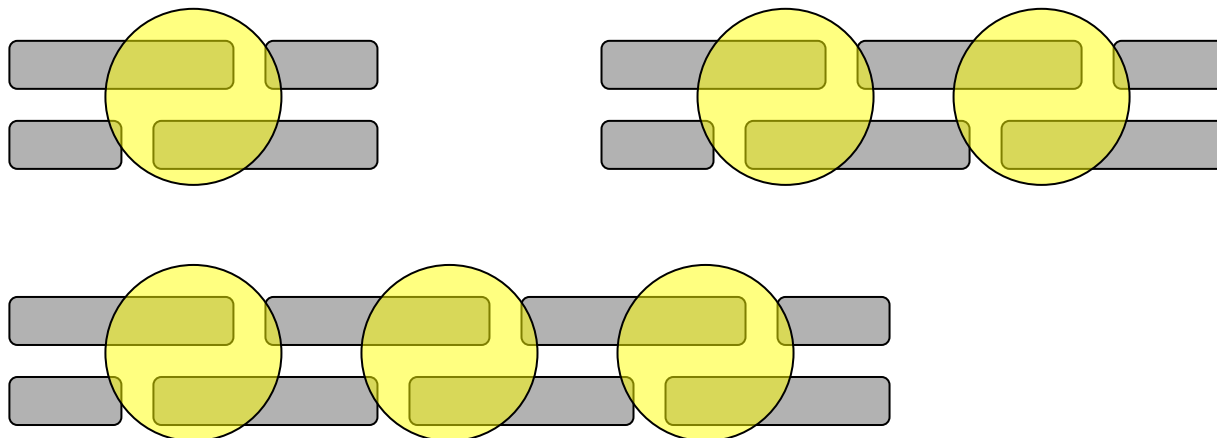
Learn

Very

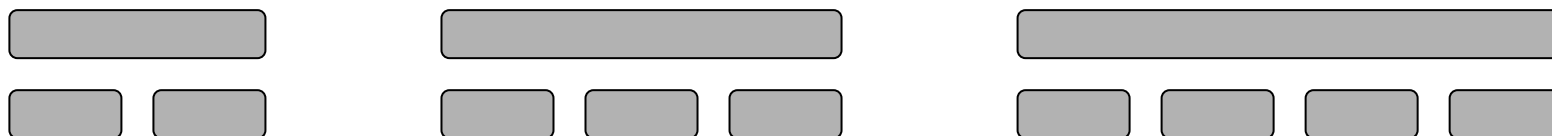
Interesting

切分歧义的类型

交集型歧义（交叉型歧义）



组合型歧义（覆盖型歧义）



汉语的切分歧义

- 交集型歧义（交叉型歧义）：如果字串abc既可切分为ab/c，又可切分为a/bc。其中a，ab，c和bc是词
 - 有意见：我 对 他 有 意见。 总统 有 意 见 他。
- 组合型歧义（覆盖型歧义）：若ab为词，而a和b在句子中又可分别单独成词
 - 马上：我 马上 就 来。 他 从 马 上 下来。
 - 将来：我 将来 要 上 大学。 我 将 来 上海。
- 混合型歧义：由交集型歧义和组合型歧义自身嵌套或两者交叉组合而产生的歧义
 - 人才能：这样 的 人 才 能 经受 住 考验。
 - 人才能：这样 的 人 才能 经受 住 考验。
 - 人才能：这样 的 人 才 能 经受 住 考验。

交集型歧义

乒 乓 球 拍 卖

完 了 。



乒 乓

球拍

卖

完

了

。

ping pong

racket

sell

over

PAST

.

The ping pong rackets are sold out.



乒 乓 球

拍 卖

完

了

。

ping pong ball

sell by auction

over

PAST

.

The ping pong balls are sold out by auction.

交集型歧义



中国	有
China	have



尚未	来
not yet	come



中	国有
in	state-owned



尚	未来
still	future

交集型歧义



结合

combination

成分

ingredient



结

tie

合成

compound

分

divide

交集型歧义



为

人民

工作

for

people

work



为人

民工

作

humanness

migrant worker

do

交集型歧义



中国

产品

质量

China

product

quality



中

国产

品质

量

in

home made

quality

amount

组合型歧义

他	将来	要	上学
He	future	will	go to school

He will go to school in the future

他	将	来	北京
He	will	come to	Beijing

He will come to Beijing.

组合型歧义

机器翻译

很

难

machine translation

very

difficult

Machine translation is very difficult.

用

机器

翻译

文章

很

难

use

machine

translate

article

very

difficult

It is very difficult to use machine to translate an article.

交集型歧义字段的链长

- 链长：交集型歧义字段中含有交集字段的个数，称为链长。
 - 链长为1：和尚未
 - 链长为2：结合成分
 - 链长为3：为人民工作
 - 链长为4：中国产品质量 结合成分分子时
 - 链长为6：努力学习语法规则
 - 链长为8：治理解放大道路面积水

真实语料中歧义字段的分布

刘开瑛，2000，《中文文本自动分词和标注》，
商务印书馆，第65页。

（500万新闻语料的统计结果）

链长	1	2	3	4	5	6	7	8	总计
词次数	47402	28790	1217	608	29	19	2	1	78248
比例	50.58	47.02	1.56	0.78	0.04	0.02	0.00	0.00	100
字段数	12686	10131	743	324	22	5	2	1	23914
比例	53.05	42.36	3.11	1.35	0.09	0.02	0.01	0.01	100

真歧义和伪歧义

- 真歧义
 - 确实能在真实语料中发现多种切分形式
 - 比如“应用于”、“地面积”
- 伪歧义
 - 虽然有多种切分可能性，但在真实语料中往往取其中一种切分形式
 - 比如“挨批评”、“市政府”

真歧义



中国

有

China

have

中国有十几亿人。

There are more than 1 billion people in China.



中

国有

in

state-owned

本地企业中国有企业占30%。

There are 30% of local enterprises are state-owned enterprises.

伪歧义



尚未

来

not yet

come

他尚未来过此地。

He has not been here yet.



尚

未来

still

future

Note: No real sentence contains such word sequence.

词语切分标准 (1)

- 建立汉语词语切分标准的必要性
 - 汉语词语定义不明确
 - 牛肉是词，鸡肉是不是？
 - 打倒是词，打死、打伤、饿死、涂黑是不是？
 - 为操作的方便，必须确定统一的标准或规范
 - 采用“分词单位”的说法
- 问题
 - 取舍理由不够充分，人为色彩过重
 - 过于复杂，难于把握

词语切分标准 (2)

- 相关的标准
 - 《信息处理用汉语分词规范》
GB/T13715-92, 中国标准出版社, 1993
 - 《资讯处理用中文分词规范》台湾中研院
 - 《人民日报》语料库词语切分规范
 -

词语切分标准 (例) (3)

《人民日报》标注语料库词语切分规范（述补结构的切分）

未收入词典的双音节述补结构，若拆开各是一个词，通常作为两个切分单位。

走/v 到/v，撞/v 上/v， 调/v 好/a，坐/v 稳/a

若拆开了，其中至少有一个是语素，通常就不切分，作为一个切分单位。

形成/v， 鼓动/v， 说明/v， 震动/v

双音节的述补结构中间插入“得”或“不”一般应予切分，

走/v 得/u 到/v，走/v 不/d 到/v，安/v 得/u 上/v，安/v 不/d 上/v

但是如果去掉“得”或“不”后，前后两个字不构成一个词的，则作为一个分词单位。

来得及/v，来不及/v，对得起/v，对不起/v，说得过去/l，说不过去/l

有的去掉“得”或“不”后虽然是一个合成词，但其中至少有一个是语素，拆开了却是难以理解的，仍作为一个切分单位

形得成/v，形不成/v

未定义词的类型

- 汉语人名：李素丽 老张 李四 王二麻子
- 汉语地名：定福庄 白沟 三义庙 韩村河 马甸
- 翻译人名：乔治·布什 叶利钦 包法利夫人
- 翻译地名：阿尔卑斯山 新奥尔良 约克郡
- 机构名：方正公司 联想集团 国际卫生组织 外贸部
- 商标字号：非常可乐 乐凯 波导 杉杉 同仁堂
- 专业术语：万维网 主机板 模态逻辑 贝叶斯算法
- 缩略语：三个代表 五讲四美 打假 扫黄打非 计生办
- 新词语：卡拉OK 波波族 美刀 港刀

未定义词识别的困难

- 未定义词没有明确边界
- 未定义词的构成单元（汉字）本身都可以独立成词

未定义词识别的依据

- 内部构成规律（用字规律）
- 外部环境（上下文）
- 重复出现规律

未定义词识别的研究进展

- 较成熟
 - 中国人名、译名
 - 中国地名
- 较困难
 - 商标字号
 - 机构名
- 很困难
 - 专业术语
 - 缩略语
 - 新词语

中国人名的内部构成规律 (1)

- 在汉语的未定义词中，中国人名是规律性最强，也是最容易识别的一类；
- 中国人名一般由以下部分组合而成：
 - 姓：张、王、李、刘、诸葛、西门、范徐丽泰
 - 名：李素丽，张华平，王杰、诸葛亮
 - 前缀：老王，小李
 - 后缀：王老，赵总
- 中国人名各组成部分用字比较有规律

中国人名的内部构成规律 (2)

- 根据统计，汉语姓氏大约有1000多个，姓氏中使用频度最高的是“王”姓，“王，陈，李，张，刘”等5个大姓覆盖率达32%，姓氏频度表中的前14个高频度的姓氏覆盖率为50%，前400个姓氏覆盖率达99%。人名的用字也比较集中。频度最高的前6个字覆盖率达10.35%，前10个字的覆盖率达14.936%，前15个字的覆盖率达19.695%，前400个字的覆盖率达90%。

中国人名的内部构成规律 (3)

- 中国人名各组成部分的组合规律
 - 姓 + 名
 - 姓
 - 名
 - 前缀 + 姓
 - 姓 + 后缀
 - 姓 + 姓 + 名（海外已婚妇女）

中国人名的上下文构成规律

- 身份词：
 - 前：工人、教师、影星、犯人
 - 后：先生、同志
 - 前后：女士、教授、经理、小姐、总理
- 地名或机构名：
 - 前：静海县大丘庄禹作敏
- 的字结构
 - 前：年过七旬的王贵芝
- 动作词
 - 前：批评，逮捕，选举
 - 后：说，表示，吃，结婚
-

中国人名识别的难点

- 一些高频姓名用字在非姓名中也是高频字
 - 姓氏：于，马，黄，张，向，常，高
 - 名字：周鹏和同学，周鹏和同学
- 人名内部相互成词，指姓与名、名与名之间本身就是一个已经被收录的词
 - [王国]维、[高峰]、[汪洋]、张[朝阳]
- 人名与其上下文组合成词
 - 这里[有关]天培的壮烈；
 - 费孝通向人大常委会提交书面报告
- 人名地名冲突
 - 河北省刘庄

中国地名的识别

- 中国地名委员会编写了《中华人民共和国地名录》，收集了全国乡镇以上（含乡镇）各级行政区域的名称，以乡镇人民政府所在地为主的居民聚落名称，山、河、湖、海、岛、高原、盆地、沙溪等自然地理实体名称，名胜古迹、纪念地、古遗址、水库、桥梁、电站等名称。共收录地名**10万**多条。这个地名录中使用的汉字共**2662**个，频度最高的前**65**个汉字占总频度的**50.22%**，前**622**个汉字占总频度的**90.01%**，前**1872**个汉字占总频度的**99%**。
- 与人名的用字情况相比较，地名用字分散得多
- 地名内部也有一定的结构，右边界比左边界更容易识别

音译名的识别 (1)

- 音译名用字非常集中 《英语姓名译名手册》中共收英语姓氏, 教名约4万个, 经计算机统计得出英语姓名译名用字表共476个:

“啊阿埃艾爱昂奥巴白柏拜班邦包保堡鲍北贝倍本比彼边别滨宾玻波博勃伯卜布采蔡藏策查察昌彻陈楚垂茨慈次聪存措达大戴代丹当道德得登邓迪底地蒂第帝丁东杜敦顿多厄恩耳尔法凡范方菲费芬丰冯佛夫福弗辅富盖甘冈高哥戈葛格各根贡古顾瓜圭郭果哈海罕翰汉杭豪赫黑亨洪侯胡华怀惠霍基吉季计嘉佳加贾简姜焦杰捷金津京久居喀卡开凯坎康考柯科可克肯孔扣寇库夸匡奎魁坤昆阔拉腊莱来赖兰朗劳勒乐雷黎理李里礼荔丽历利立莲连廉良列琳林霖龄留刘流柳龙隆卢鲁露路吕略伦萝罗洛玛马麦迈满曼芒茅梅门蒙孟米密敏明名摩莫墨默姆木穆拿娜纳乃奈南内嫩能妮尼年涅宁牛纽农努女诺欧帕派潘庞培佩彭蓬皮匹平泼朴普漆奇齐契恰钱强乔切钦琴青琼丘邱屈让热仁日荣茹儒瑞若撒萨塞赛三缮桑瑟森莎沙珊山尚绍舍申生盛圣施诗石什史士寿舒朔斯思丝松孙索所塔泰坦汤唐陶特藤提惕田铁汀廷亨通透图托脱娃瓦万旺威韦为维伟魏卫温文翁沃乌武伍西锡希悉席霞夏显香向晓肖歇谢欣辛兴幸姓雄休修雪逊雅亚延扬阳尧耀耶叶依易意因英永尤雨约宰赞早泽曾扎詹湛章张哲者珍真芝知智治朱卓兹子宗祖佐丕谟葆薇岑弼娅缪珀璠赉滕斐熙鸠窠艮麟黛”。

音译名的识别 (2)

- 音译名内部很难划分出结构，但有一些常见音节，如“斯基、斯坦”等
- 不同语言的音译规律不尽相同，如法语、俄语、蒙古语译名用字与英语就有较大区别（蒙语人名举例：“那顺乌日图、青格勒图”），如果按不同的语言训练不同的模型可能会比使用统一的模型效果更好
- 音译名可以是人名、地名或其他专名，上下文规律差别较大
- 由于音译名用字比较集中，识别正确率较高

机构名的内部构成规律 (1)

- 机构名一般都是定中结构
- 机构名的后缀一般比较集中，识别相对容易
- 机构名左边界识别非常困难
- 机构名中含有大量的人名、地名、企业字号等专有名称。在这些专有名称中，地名所占的比例最大，其中未登录地名又占了相当一部分的比例。所以机构名识别应在人名、地名等其他专名识别之后进行，其他专名识别的正确率对机构名识别正确率有较大影响

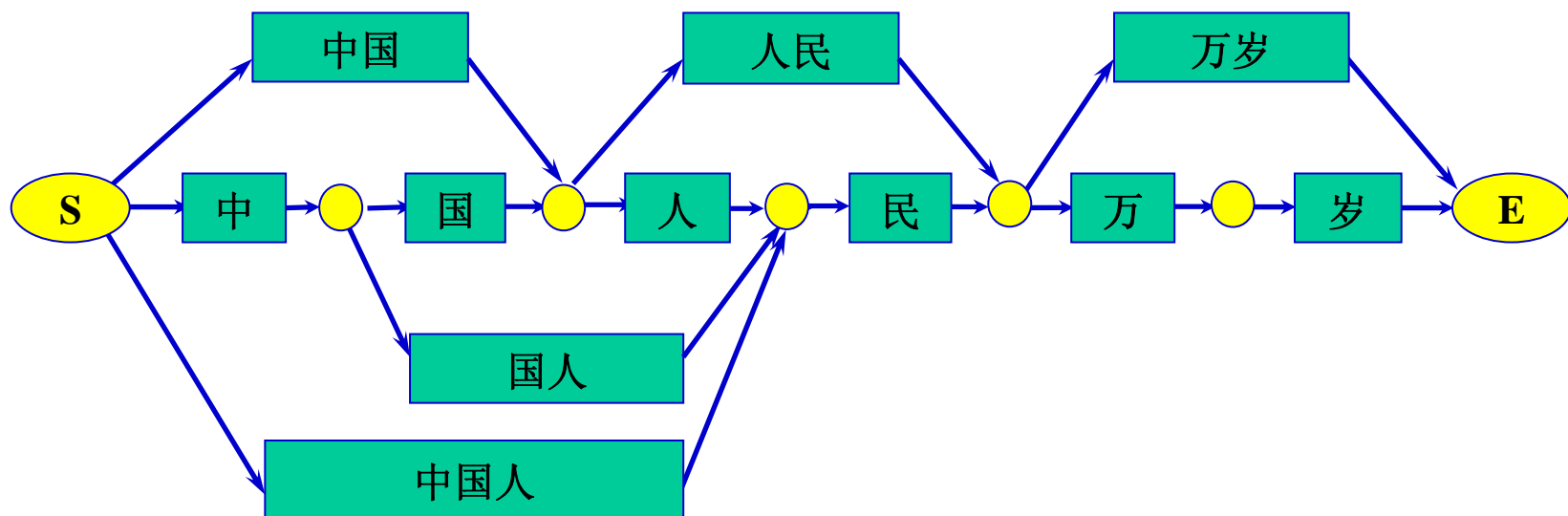
机构名的内部构成规律 (2)

- 中文机构名用词非常广泛。通过对人民日报1998年1月其中的10817个机构名所含的19986个词进行统计，共计27种词，其中名词最多（9941个），地名其次（5023个），依次为简称（1169个）、专有名词（1125个）、动词（848个）以及机构名（714个）等
- 机构名长度极其不固定
- 机构名很不稳定。随着社会不断发展，新机构不断涌现，旧机构不断被淘汰、改组或更名

内容提要



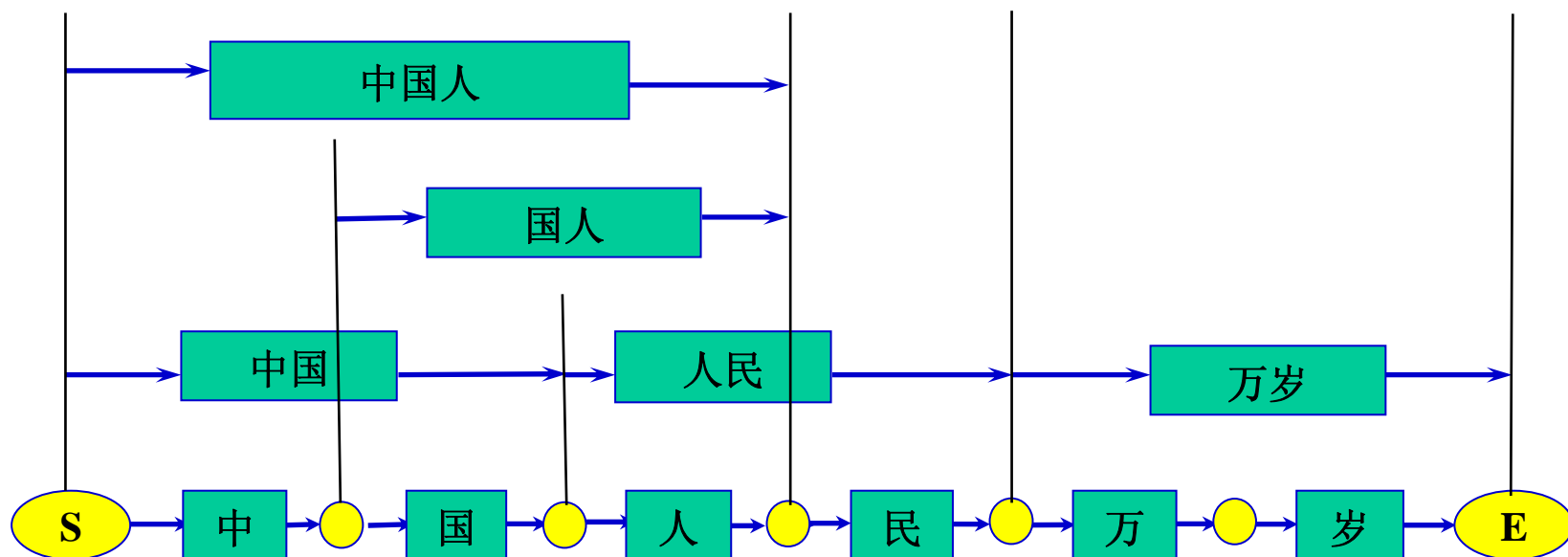
汉语切分的数据结构—词图



在词图（**word graph**）上，我们可以把词法分析中的几种操作转化为：

- 给词图上添加边（查词典，处理重叠词、离合词、前后缀和未定义词）；
- 寻找一条起点**S**到终点**E**的最优路径（切分排歧）；
- 给路径上的边加上标记（词性标注）；

汉语切分的数据结构—词格



- 词格 (**word lattice**) 和词图 (**word graph**) 是等价的表示形式

基于词典的词语机械切分方法

- 输入：
 - 一个词表（词典）
 - 一个待切分句子
- 输出：
 - 一个词语序列

基于词典的词语机械切分方法

- 全切分
- 最大匹配方法
- 最短路径方法
- 基于记忆的交叉歧义排除法
- 基于规则的切分算法

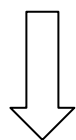
全切分方法

- 给出所有的切分结果
- 算法（略）
- 算法的时间复杂度随着句子长度的增加呈指数增长

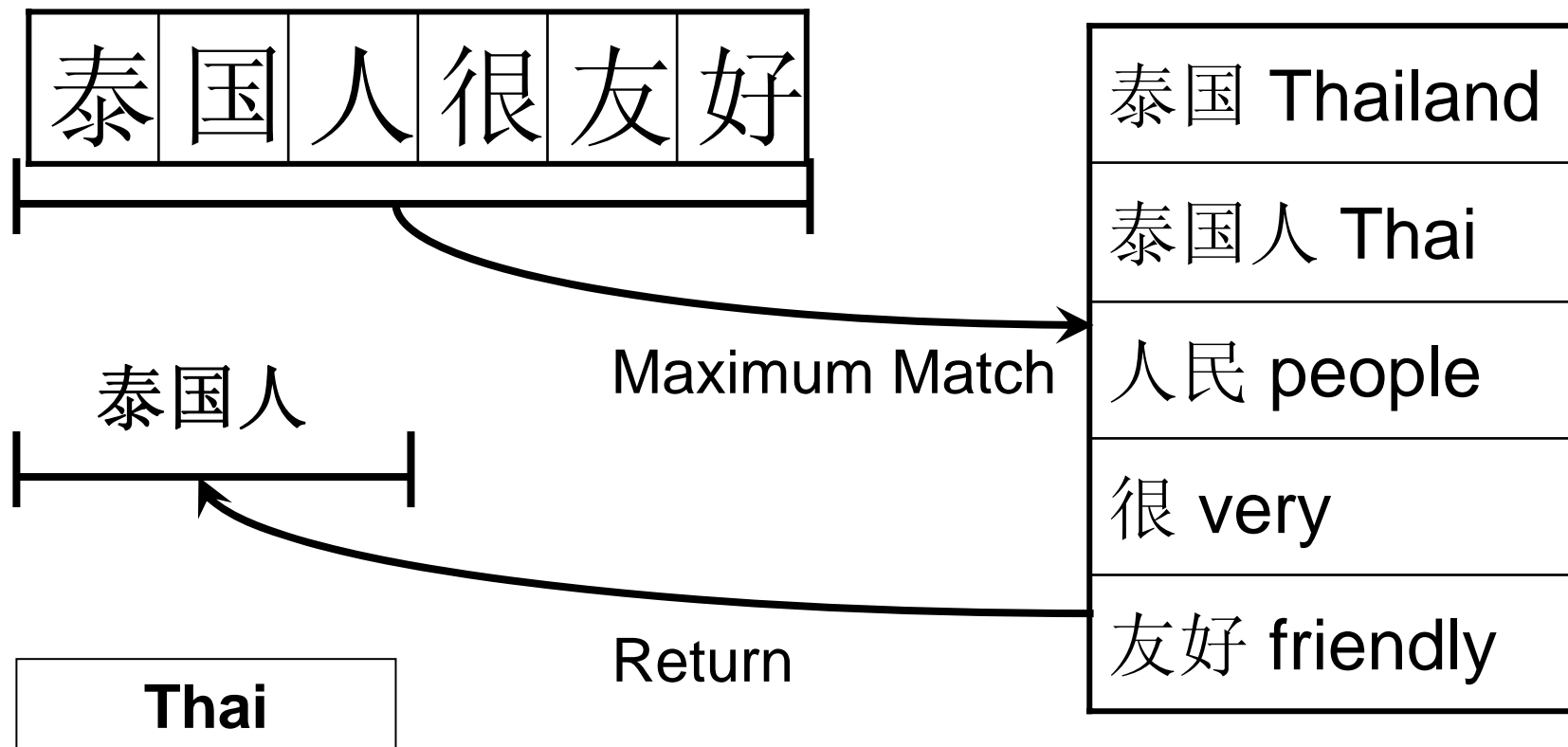
最大匹配方法 (1)

- 正向最大匹配 (MM)
 - 自左往右
 - 每次取最长词
- 逆向最大匹配 (RMM)
 - 自右往左
 - 每次取最长词
- 双向最大匹配
 - 依次采用正向和逆向最大匹配
 - 如果结果一致则输出
 - 如果结果不一致再用其他方法排歧

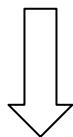
正向最大匹配



词典



正向最大匹配



词典

泰	国	人	很	友	好
---	---	---	---	---	---

Maximum Match

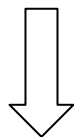
很

Return

Thai	very
------	------

泰国 Thailand
泰国人 Thai
人民 people
很 very
友好 friendly

正向最大匹配



词典

泰	国	人	很	友	好
---	---	---	---	---	---

Maximum Match

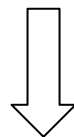
友好

Return

Thai	very	friendly
------	------	----------

泰国 Thailand
泰国人 Thai
人民 people
很 very
友好 friendly

正向最大匹配



词典

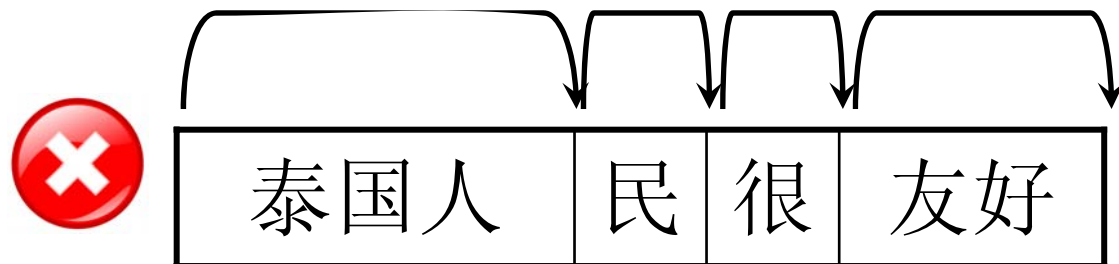
泰国人	很	友好
-----	---	----

Thai are very friendly

泰国 Thailand
泰国人 Thai
人民 people
很 very
友好 friendly

正向最大匹配

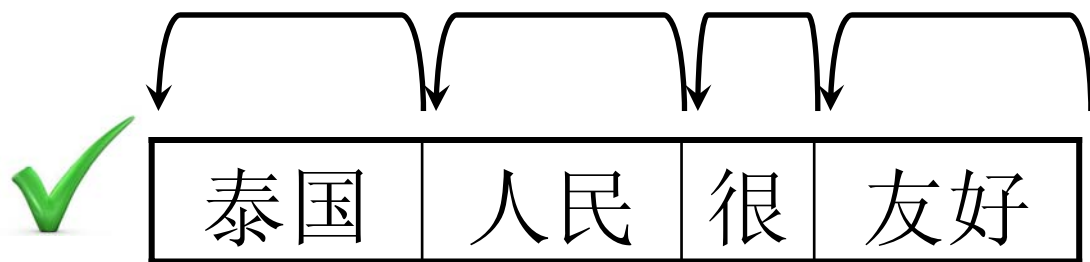
- 正向匹配算法显然很容易导致错误



泰国	Thailand
泰国人	Thai
人民	people
很	very
友好	friendly

逆向最大匹配

- 这个例子用逆向匹配算法切分正确.



Thai people are very friendly

泰国 Thailand
泰国人 Thai
人民 people
很 very
友好 friendly

最大匹配方法 (2)

- 优点
 - 简单、快速
 - 在某些应用场合已经足够
- 缺点
 - 单向最大匹配会忽略交集型歧义和组合型歧义
幼儿园 地 节目 / 独立自主 和平 等 互利的 原则
 - 双向最大匹配会忽略链长为偶数的交集型歧义和组合型歧义
原子 结合 成分 子时 / 他从马上下来

最短路径方法

- 基本思想：
 - 在词图上选择一条词数最少的路径
- 算法：
 - 动态规划算法
- 优点：好于单向的最大匹配方法
 - 最大匹配：独立自主 和平 等 互利 的 原则(6)
 - 最短路径：独立自主 和 平等互利 的 原则(5)
- 缺点：忽略了所有覆盖歧义，也无法解决大部分交叉歧义
 - 结合 成分 子时

基于记忆的交叉歧义排除法

- 孙茂松，左正平，邹嘉彦，1999, 高频最大交集型歧义切分字段在汉语自动分词中的作用，中文信息学报，Vol.13, No.1, 1999
- 该文考察了一亿字的语料，发现交集型歧义字段的分布非常集中。其中在总共的22万多个交集型歧义字段中，高频的4,619个交集型歧义字段占有所有歧义切分字段的59.20%。而这些高频歧义切分字段中，又有4,279个字段是伪歧义字段，也就是说，实际的语料中只可能出现一种切分结果。这样，仅仅通过基于记忆的方法，保存一种伪歧义切分字段表，就可以使交集型歧义切分的正确率达到53%，再加上那些有严重偏向性的真歧义字段，交集型歧义切分的正确率可以达到58.58%。

基于规则的切分方法

- @@ 高峰(A+B, AB)
CONDITION FIND(L,NEXT,X){%X.yx=最|更} SELECT 1
OTHERWISE SELECT 2
- @@ 分成(A+B, AB)
CONDITION FIND(R,NEXT,X){%X.ccat=m} SELECT 1
OTHERWISE SELECT 3
- @@ 是因为(A+B,AB)
CONDITION FIND(L,FAR,X){%X.yx=之所以} SELECT 2
OTHERWISE SELECT 1
- @@ *(A+BC, AB+C)
CONDITION %A.ccat =p, %BC.ccat =n |f |s |t SELECT 1
- @@ *(A+BC, AB+C)
CONDITION %A.ccat=v, %BC.ccat=b|d|t, %AB.ccat=r,
%C.ccat=n SELECT 2