

# An Integrated Chinese Grammar Development Environment

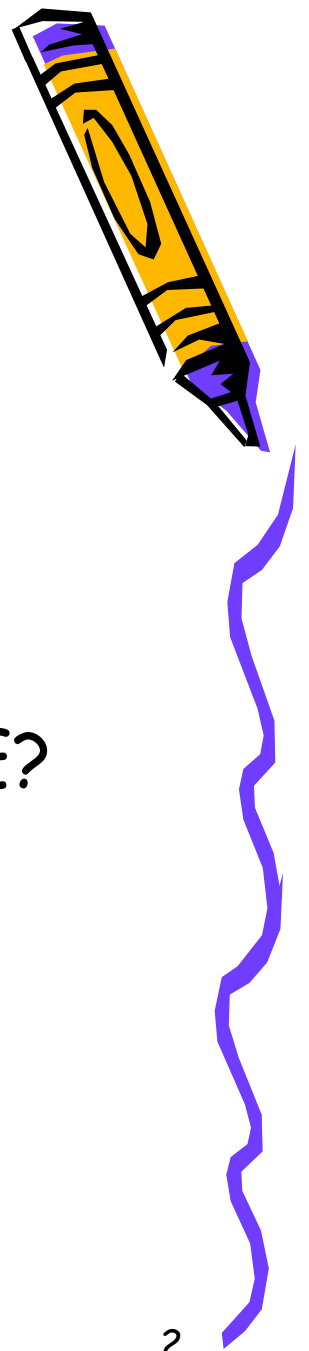
Zhan Weidong\*, Liu Qun, Chang Baobao, Zhang Qinlong, Wu Yonghua

\* Dept. of Chinese Language & Literature, Peking University  
Beijing, 100871, PR.China



Email: [zwd@pku.edu.cn](mailto:zwd@pku.edu.cn)  
Personal Homepage: <http://ccl.pku.edu.cn/doubtfire/>

# Outline of Talk



- 1) Introduction: Why we need ICGDE?
- 2) What does ICGDE consist of?
- 3) What is the characteristics of ICGDE?
- 4) What have been done with ICGDE
- 5) Future works

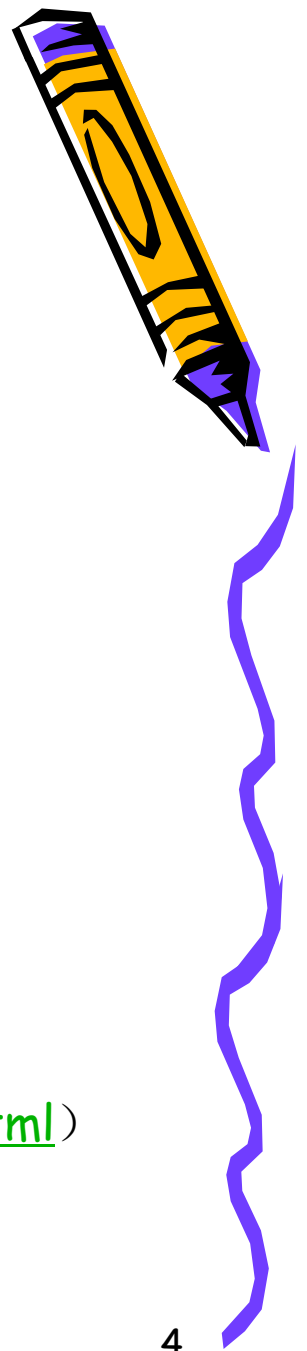


# 1 Introduction

- There are many applications which require a sentence parser as a core or fundamental component. For instance, MT, Human-computer interface, IE, IR, etc..
- Whatever the applications make use of shallow parsing techniques or rely on deep analysis, precise large-coverage grammars, lexicons and annotated corpus of natural languages are always needed to build up.
- Developing such so-called linguistic knowledge base(LKB, for short) is a very time-consuming and difficult activity. The following factors are often needed to taken into consideration: coverage, granularity, efficiency, etc..
- We need a LKB development environment that can process various forms of linguistic knowledge effectively and integrate them into an organized system to support applications of NLP and linguistic research and teaching.

# Overseas Researches on Development of Linguistic Resource

- INTEX: A Linguistic Development Environment  
(<http://www.nyu.edu/pages/linguistics/intex/>)
- The LinGO (Linguistic Grammars Online) project  
(<http://lingo.stanford.edu>)
- Language Resources & Evaluation Conference  
(<http://www.lrec-conf.org/>)
- Xerox Linguistics Environment Project  
(<http://www2.parc.com/istl/groups/nlitt/xle/>)
- XTAG System  
(<http://www.cis.upenn.edu/~xtag/release-8.31.98-html/node12.html>)



# Our Ongoing Projects

Our research work is now supported by several state-funded projects which are listed below.

- 2002.1 -- 2006.12 Construction Rules of Chinese Sentence And Development Environment of Chinese Parsing System, funded by Ministry of Education of China (Project No.200110)
- 2002.1 -- 2004.12 Study on Annotation Specification of Chinese Phrases and Sentence Patterns, funded by Ministry of Education of China (Project No.YB105-49)
- 2002.12 -- 2005.12 Computing Platform for Research on Contemporary Chinese and Information Processing, funded by “211 Project” of The National Tenth Five-year Plan

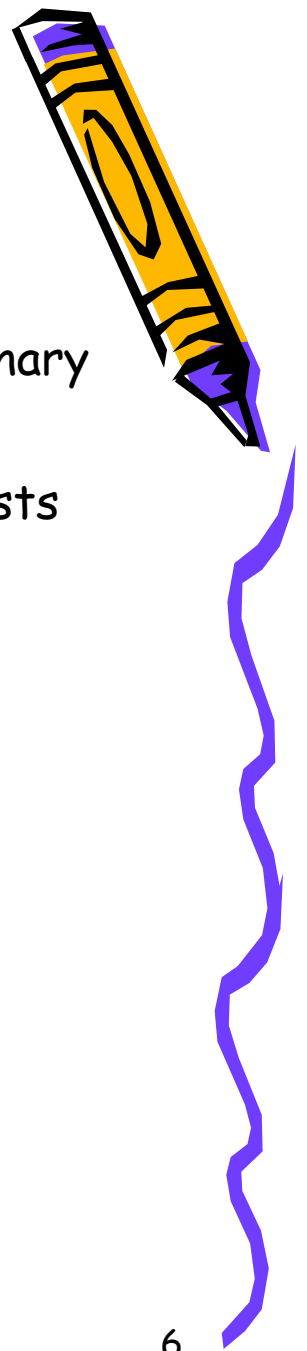


# Research Collaboration

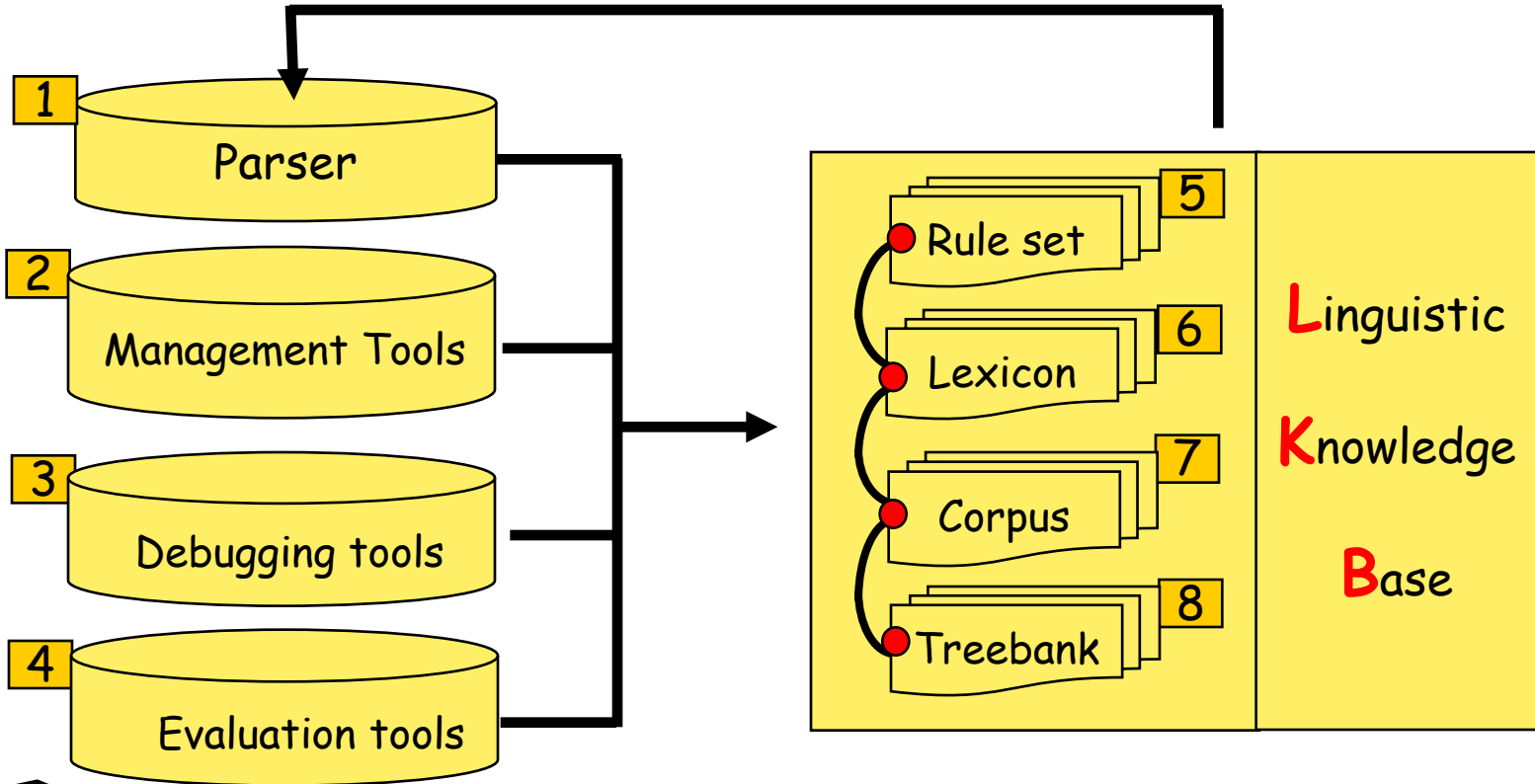
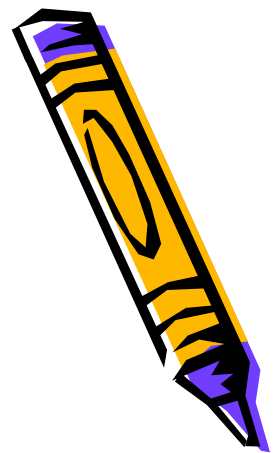
Computing-oriented research on natural language needs interdisciplinary collaboration more than traditional language research works.

Our research team is composed of linguists and computational linguists from four institutes, including:

- Department of Chinese Language and Literature, Peking University
- Center for Chinese Linguistics, Peking University
- Institute of Computing Technology, Chinese Academy of Sciences
- Institute of Computational Linguistics, Peking University



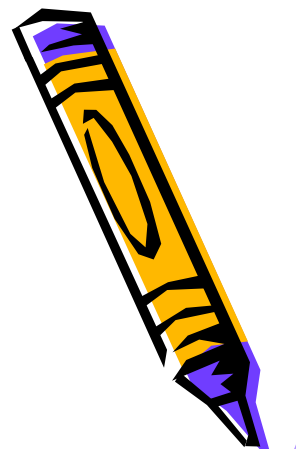
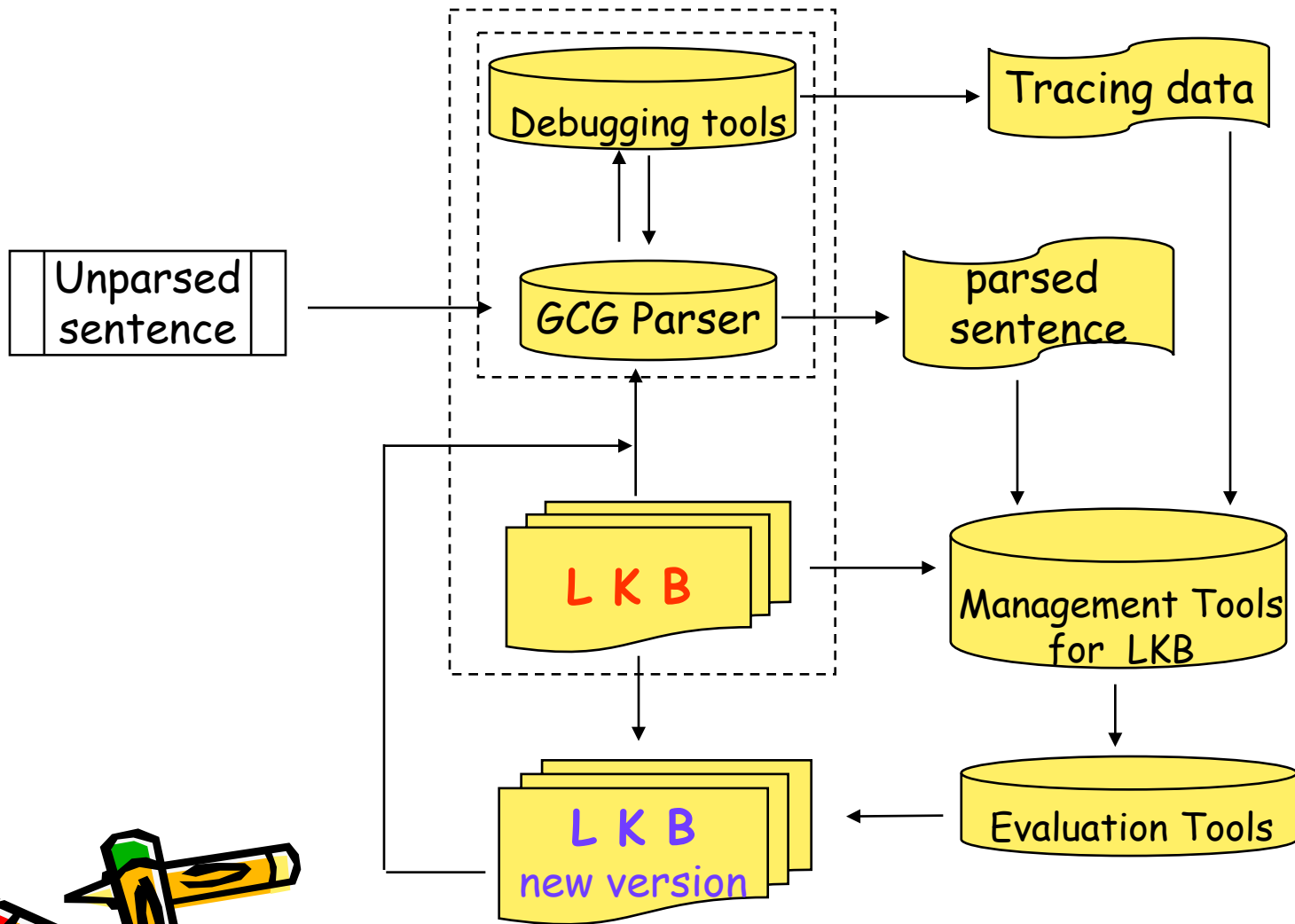
## 2 What does ICGDE consists of



1,2,3,5 will be discussed detailedly in the following sections.  
4 is under construction. 6,7,8 are shown in appendix.



# Workflow of ICGDE



# The main interface of GCG parser

The screenshot shows the main interface of the GCG parser. The window title is "TestParserDoc - [test.trn]". The menu bar includes "文件(F)", "编辑(E)", "查看(V)", "知识库(K)", "分析(T)", "选项(O)", "窗口(W)", and "帮助(H)". The toolbar contains icons for file operations and navigation, along with labels "放大", "缩小", "切换", "清空", and "GOTO".

The left pane displays a list of test sentences and their corresponding GCG annotations:

- 卖给老王的自行车
- 修理老王的自行车
- 研究方法
- 四个
- 四个人
- 八个人
- 八/m 个/q 人/n
- 慢说四个人抬不动，就是八个人也不行。
- 1916年5月21日/t 被/p 定/v
- 为/p 国庆日/n 。 /w
- 1916年5月21日/t 被/p 定为/v
- 国庆日/n
- 1916年
- 5月21日
- 1916年5月
- 1916年5月21日
- 定为国庆日
- 被定为国庆日
- 1916年5月21日被定为国庆日。
- 管理体制
- 北大的管理体制有待改进。
- 日本最大的火山爆发了
- 在房间里打架
- 大家十分关注美军虐萨事件
- 大家十分关注美军虐囚事件
- 大家/r 十分/d 关注/v 美军/n
- 虐囚/v 事件/n
- 大多数人的支持

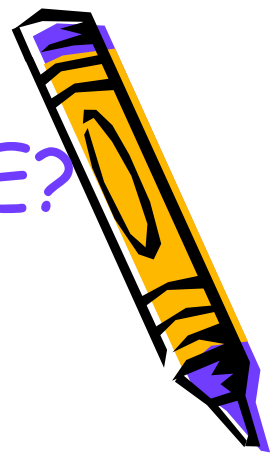
The right pane shows the parse tree for the sentence "修理老王的自行车". The root node is "vp", which branches into "lvp" and "np". "lvp" branches into "lv", which leads to the terminal node "修理". "np" branches into "ap" and "Inp". "ap" branches into "Inp" and "u". "Inp" branches into "In", which leads to the terminal node "老王". "u" leads to the terminal node "的". "Inp" branches into "In", which leads to the terminal node "自行车".

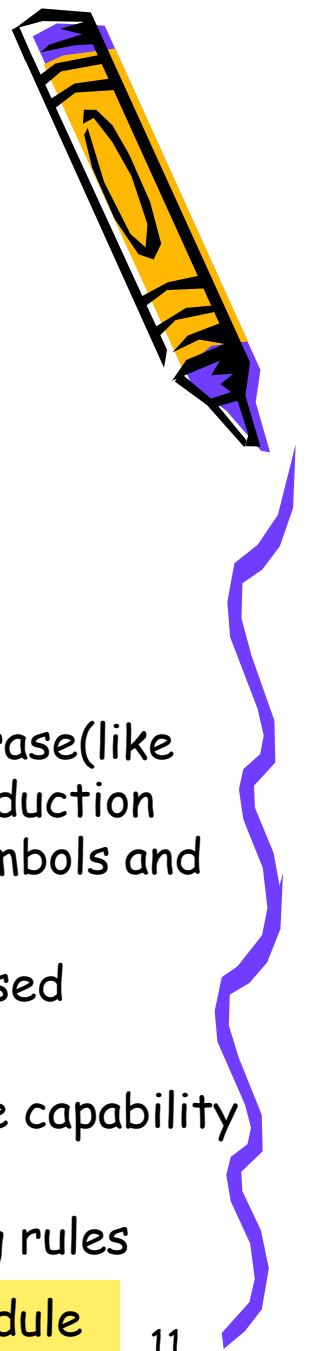
The status bar at the bottom indicates "就绪" (Ready) and "第2行 第5列" (Line 2, Column 5). The execution time is shown as "用时: 20毫秒" (Time: 20 milliseconds).

### 3 What is the characteristics of ICGDE?

In this section, we are going to talk about the highlighted features of the three main modules in ICGDE

- 1) A Chinese Parser with fully constraint-based grammar formalism (GCG parser)
- 2) Tools for linguistic rule writing and debugging
- 3) Tools for building and using Chinese Treebank





# 3.1 A Chinese Parser with fully constraint-based grammar formalism

About the parser①

- Non-deterministic Bottom-up Chart parser
- Towards a multi-engine approach
- Configurability and Customization

About the linguistic knowledge representation

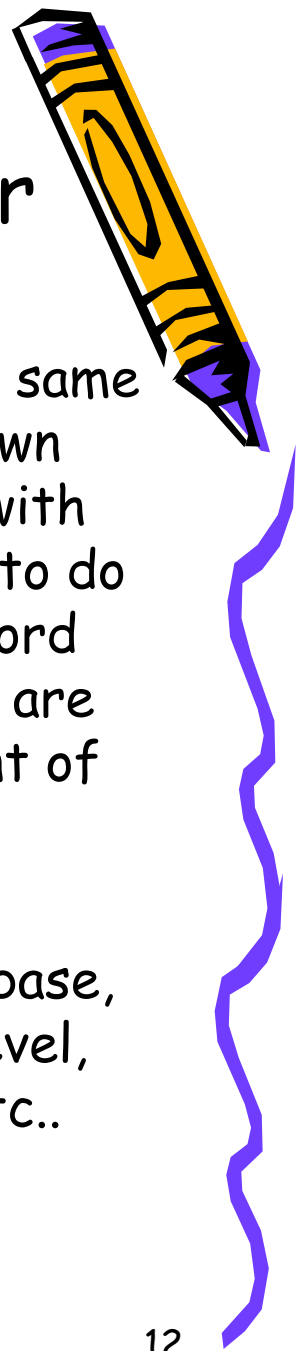
- CFG rule enhanced with marking the head of a phrase (like HPSG formalism) and the right hand side of a production rule is allowed to be composed of non-terminal symbols and trees (like Tree Adjoining Grammar formalism)
- Each rule is decorated with feature structure-based unification
- Built-in Functions and operators for enhancing the capability of construction rules
- Score mechanism for disambiguation of competing rules



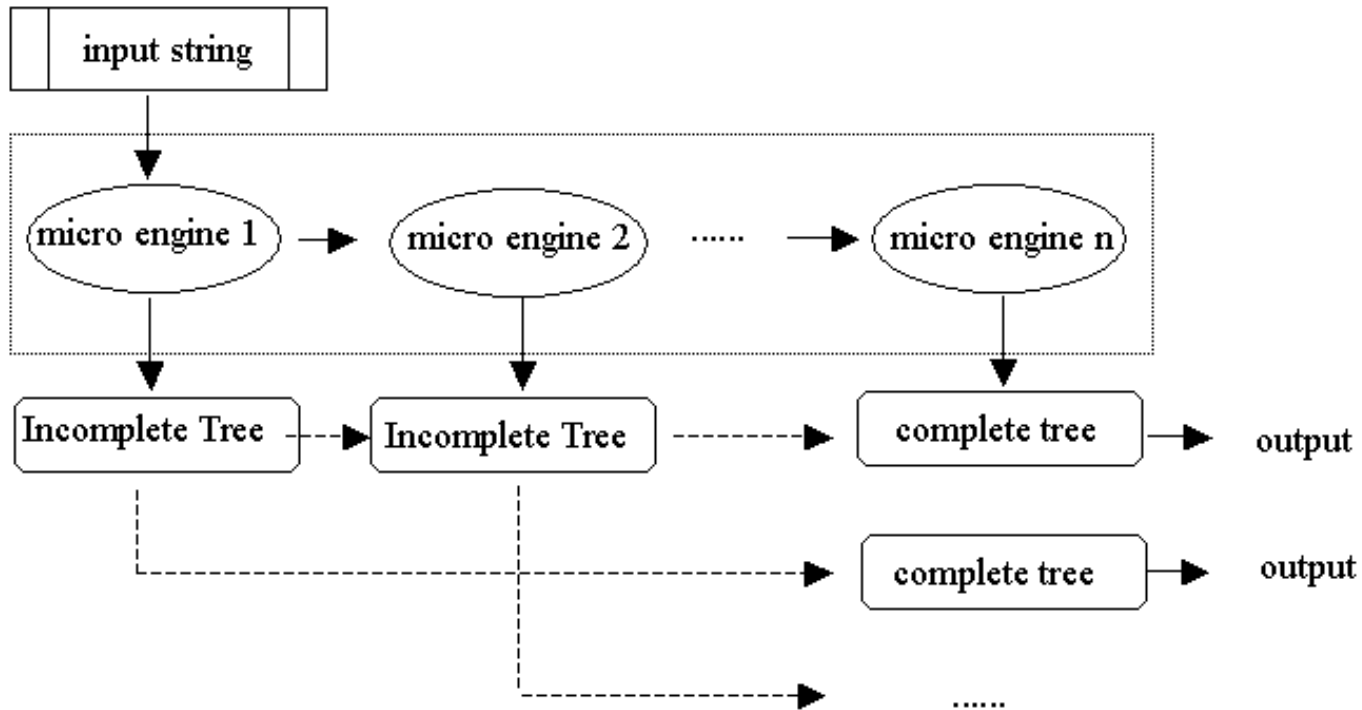
footnote:① The main architecture and the core module of the parser is designed and developed by Dr. Liu Qun

# Multi-engine architecture of the parser

- The parser consists of various micro-engines which share a same **chart data structure** while parsing. Each engine uses its own linguistic resource to make contribution to add new nodes with score into the chart. It is allowed to use different engine to do the same job. For example, we can use different Chinese word segmenter and POS-tagger in our parser. The nodes which are derived from different engines and cover the same segment of input string will be selected according to their score.
- Each micro-engine is bound to its own linguistic knowledge base, which can be various form and on the different linguistic level, such as lexicon, rule set, corpus, probability parameters, etc..

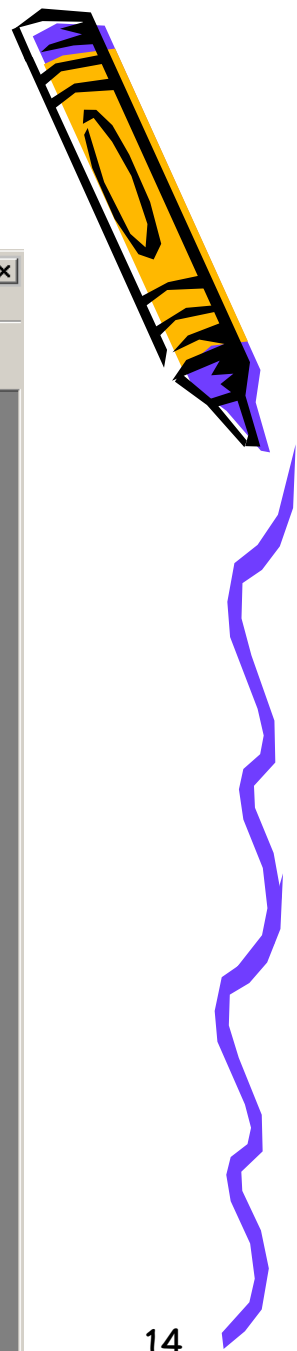


# Multi-engine architecture of the parser



A pipelined multi-engine approach to Chinese parsing

# Configurability and Customization



To set weight for adjusting score of a parsed node depending on different conditions

micro engines are listed here that can be switched on and off

To specify a coefficient to calculate the maximum amount of tree nodes for parsing one sentence

To specify the maximum amount of time for parsing one sentence

TestParserDoc

文件(F) 查看(V) 知识库(K) 选项(O) 帮助(H)

放大 缩小 切换 清空 GOTO

参数设定

内含标点、单侧非标点惩罚因子 (0到1之间): 0.1

内含标点、双侧非标点惩罚因子 (0到1之间): 0.05

同区域、同标记结点惩罚因子 (0到1之间): 0.2

同区域、同规则结点惩罚因子 (0到1之间): 0.2

同区域、同核心结点惩罚因子 (0到1之间): 0.2

分析总结点数约束 ( $a \times N \times \log N$ ) 因子  $a$  (5~50): 20

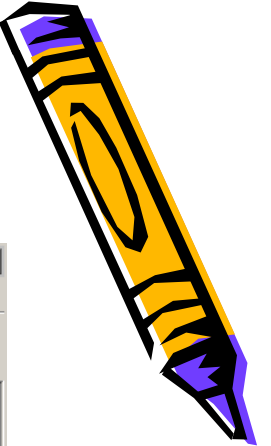
分析时间约束 (10~600秒): 90

屏蔽分析微引擎: 屏蔽生成微引擎:

SegTagUserRecognizer  
SegTagBDRecognizer  
PhraseRecognizer  
DictnRecognizer  
RuleBasedRecognizer  
FailSoftRecognizer

OK  
Cancel

# Configurability and Customization(cont.)



segmentation and pos tagging

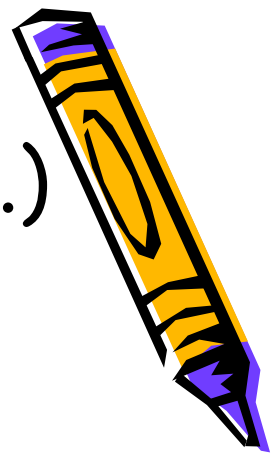
parsing a sentence with word segmentation and pos tag

parsing a sentence unannotated

The screenshot shows the TestParserD application window. On the left is a text editor with a document named 'test.trn' containing several lines of Chinese text. On the right is a '输出选项设定' (Output Options Setting) dialog box. The dialog has several sections: '单语 / 双语' (Monolingual / Bilingual) with radio buttons for '输出单语' (selected) and '输出双语对照'; '切分 / 标注' (Segmentation / Annotation) with radio buttons for '切分' and '切分 + 标注' (selected); '输出格式' (Output Format) with dropdown menus for '切分标注' (set to '格式一 (北京大学)') and '句法分析' (set to '格式一 (计算所)'); and '切分标注输出样例' (Segmentation and Annotation Output Example) showing '我/r 吃/v 了/u 苹果/n 。 /w'. Below this is a '句法分析输出样例' (Syntax Analysis Output Example) showing a complex tree structure. At the bottom of the dialog are 'Cancel' and 'OK' buttons. Below the dialog, the words '修理', '老王', '的', and '自行车' are shown in boxes with lines connecting them to the corresponding words in the text editor. A green text annotation 'User can choose an output format from different styles' points to the '切分标注' dropdown menu.

User can choose an output format from different styles

# Configurability and Customization (cont.)



- Each micro-engine can be switched on and off before parsing a sentence
- User can put a new engine into the parser if necessary, for example, name entity recognizer and shallow parsing engine may be needed to add for information extraction system.
- User can set up the time limitation for parsing a sentence
- User can set up the maximum amount of tree nodes created by parsing one sentence as a limitation
- User can set weight for adjusting score of a parsed node in the different cases. For example, if a node contains punctuation, its score should be reduced slightly
- User can choose different input and output format



# Enhanced CFG production rule



Global rule

-- different coverage--

Local rule

pp->!p np

ap->!ap c ap

np->mp !np

pp -> !p<从> tp v<起|开始>

tp -> mcp !q<点> mcp q<分>

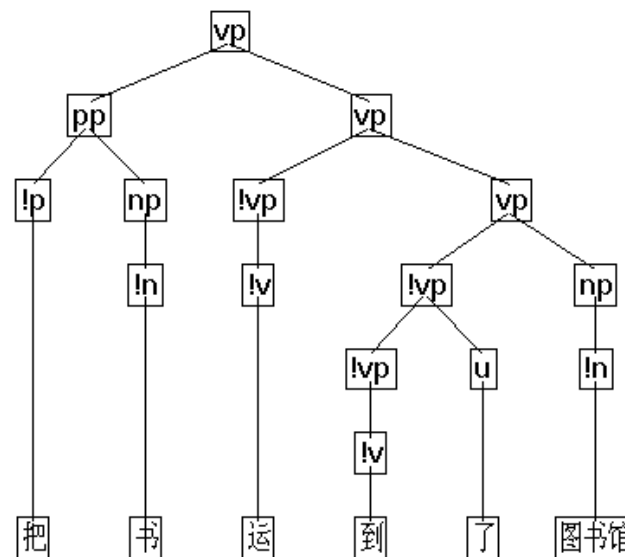
ap -> m<一> q<年|天> p<比> m<一> q<年|天> !ap

... If one or more nodes in the right hand side are specified with terminal symbol, i.e. word, this kind of rule is called "local rule" for its locality.

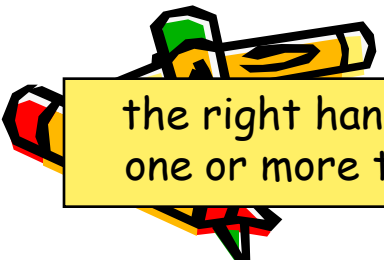
A non terminal symbol in the left hand side of a rule (root node of a tree) can be rewrite as more than two symbols(branches) in the right hand side

vp -> pp( !p<把> np ) !vp(!vp vp(!vp<<到>> np))

so called "mildly context sensitive grammar formalism" like TAG-style grammar

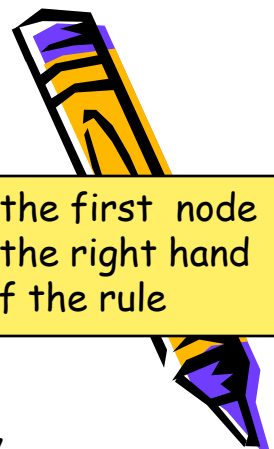


carried the books to library



the right hand side of a rule can contain one or more trees instead of just nodes

# CFG production rule with constraints



start symbol of a rule

stands for root node

stands for the first node  
appears in the right hand  
side of the rule

cfg rewriting rule

&& {npdz11} np->vp !n

Rule-ID

delimiter

```
:: $.内部结构=粘合定中,$.定语=%vp,$.中心语=%n,  
%vp.后名=是,%n.前动=是,%vp.内部结构=单词|联合,  
IF %vp.内部结构=单词 THEN %vp.兼类=~n ENDIF,  
IF %vp.内部结构=单词,%vp.音节=2, %n.音节=2 THEN  
#Score(5) ENDIF,  
IF %vp.内部结构=单词,%vp.音节=1, %n.音节=2 THEN  
#Score(-5) ENDIF,  
IF #GetLength(%vp) > 4 FALSE
```

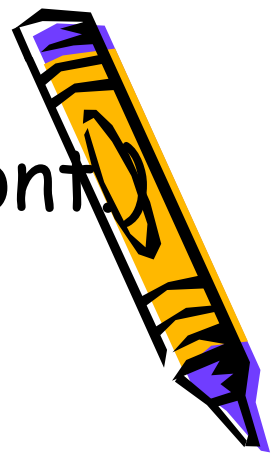
a sequence of  
constraints that  
are separated by  
comma

# symbol is always followed by a built-in function

Basically, each constraint is independent from each other.  
But it will be easier to read if the constraints can be  
written in good order.



# CFG production rule with constraints(cont)

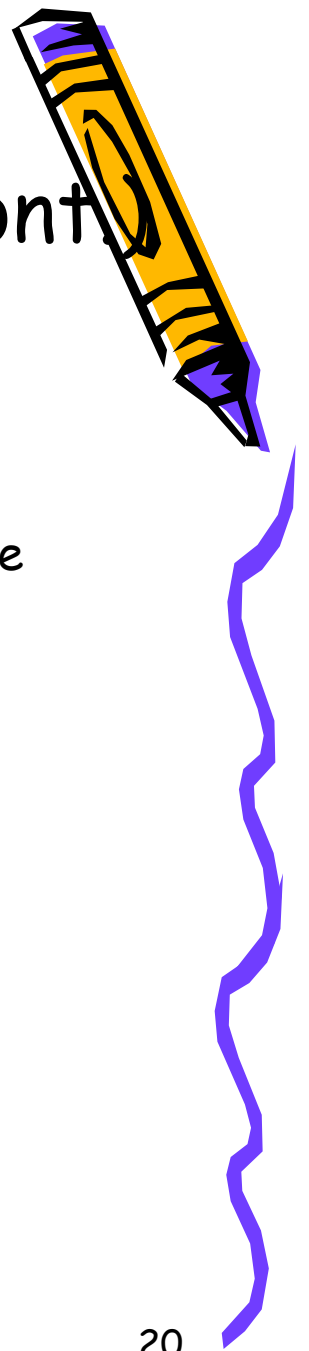


There are four types of constraints listed here one after another.

- 1) The first type of constraint is used to initialize the value of some required features of a rule. For example, the equation “\$.内部结构=粘合定中” means that the internal relation of this construction is Modifier-Head.
- 2) The second type of constraint is used to describe that the constituents of a phrase should satisfy certain conditions. For example, the equation “%vp.后名=是” means that the property “后名” of the first vp in the phrase should have a value “yes”. If the value of a verb is “no”, the condition is unmatched and then the verb is rejected by this rule



# CFG production rule with constraints(cont)



3) The third type of constraint is also used to describe the constraints like the second one. But the form is different. Like programming language designed for computer, we can use the sentence of “IF ... FALSE”, “If ... Then ... Endif”, “IF ... then ... else ... Endif”, etc., to describe compound conditions.

4) The fourth type of constraint contains built-in functions and operators that can be used to enhance the capability of writing constraints for rules. The usage of built-in functions will be illustrated in the next slides.



# Built-in Functions and Operators

All built-in functions can be classified into the following four types according to the type of value returned by a function.

## Bool Function

MatchPattern

Score

## String Function

SubString

## Int Function

NumOfChild

GetLength

FindWord

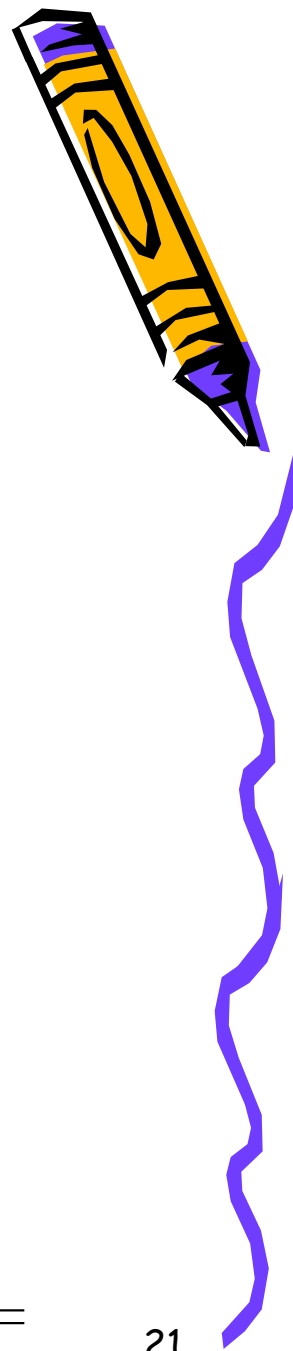
## Node Function

GetMostRightNode

GetMostLeftNode

GetLeafNode

Operator: > < >= <= == !=



# Usage Illustration of Built-in Functions

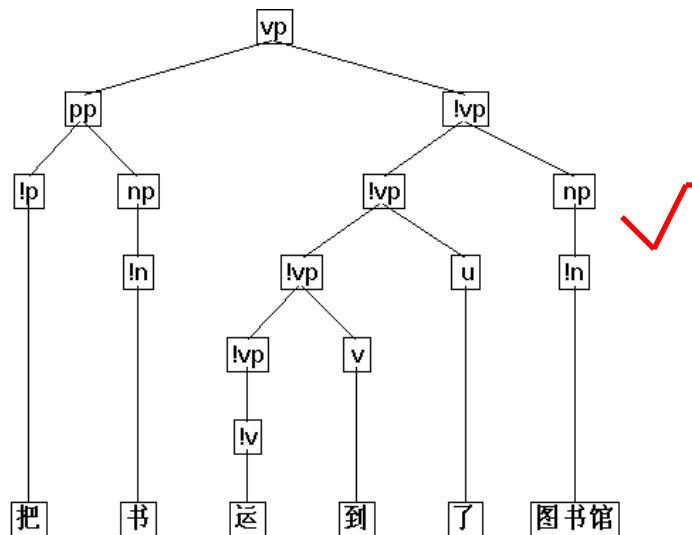
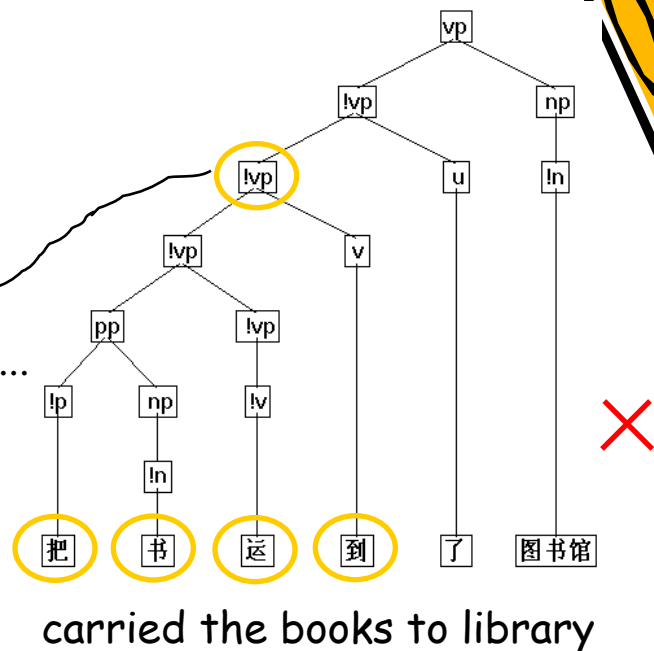
## ex.1 GetLength()

The top right parsing tree is not correct.

In order to avoid it, we can make use of the function *GetLength()* in the following rule:

{vp1} vp -> !vp u :: ..., IF #GetLength(%vp) >=4 FALSE, ...

The verb phrase(vp) “把 + 书 + 运 + 到” consists of four words. So the function *GetLength* took this vp as argument will return the integer value 4, which satisfy the unification condition “#GetLength(%vp) >=4”. The rule vp1, therefore, is rejected while parsing the phrase. And then the expected result is produced. See the bottom right parsing tree.

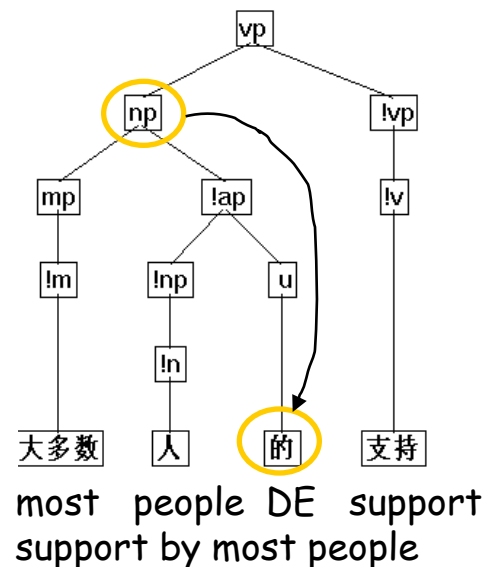


# Usage Illustration of Built-in Functions



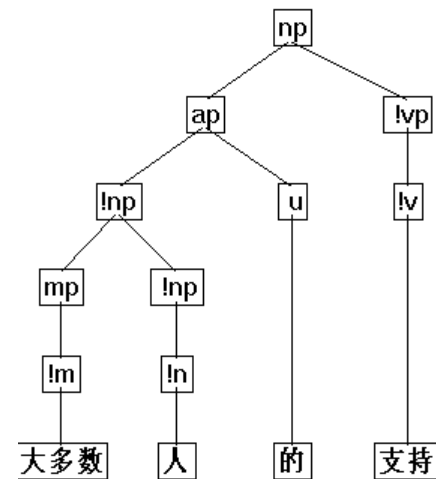
## ex.2 GetLeafNode()

The function `GetLeafNode` takes two arguments. The first argument is a tree node, the second argument is the index of a leaf node which is projected by the tree node that the first argument denotes. For instance, in order to eliminate the incorrect parsing tree (shown top right), we can make use of this function in writing the following rule:



most people DE support  
support by most people

{vp2} vp -> np !vp :: ..., IF %GetLeafNode(%np,-1).原形=的 FALSE, ...



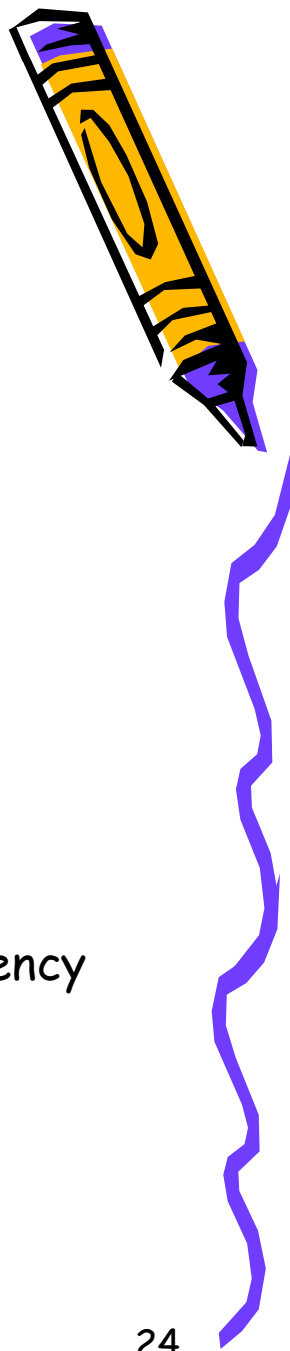
The second argument “-1” is interpreted as the most right leaf node of the first argument, i.e. np. The above constraint means that a np ended with “的” can NOT be positioned before a vp as its adverbial modifier. In the above rule, `GetLeafNode` returns the leaf node “的”, whose feature “原形” (lexmeme) has the value “的”. So the unification will succeed.

# Score mechanism

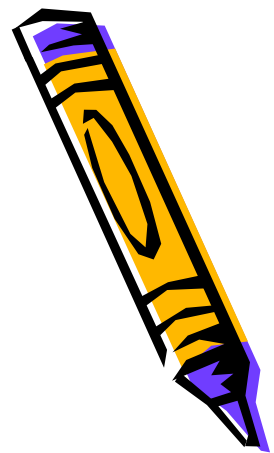
- Using a simple score mechanism to make fine granularity  
set different score thresholds for individual rules or  
different lexicons depending on their credibility

for instance:

- 1) basic lexicon vs. extended lexicon
- 2) global rule vs. local rule
- 3) two global rules with same construction but different frequency



# Score mechanism (cont.)



## Local rule

## Global rule

&& {vpdao1} vp -> pp( !p<把> np ) !vp(!vp vp(!vp<<到>> np))

&& {vpzz1} vp->pp !vp

```
up -2.10694{vpdao1}
====pp -0.694124{pp3}
====p<把> 0{}
====np -0.00097704{np00}
====n<书> 0{}
====up -1.40497{upyundao}
====up -0.00097704{up00}
====u<运> 0{}
====up -0.71085{upsb1}
====up -0.0167254{up1}
====up -0.00097704{up00}
====u<到> 0{}
====u<了> 0{}
====np -0.00097704{np00}
====n<图书馆> 0{}
```

```
up -2.79127{vpzz1}
====pp -0.694124{pp3}
====p<把> 0{}
====np -0.00097704{np00}
====n<书> 0{}
====up -1.404{upsb1}
====up -0.709873{up1}
====up -0.694124{upsbu1}
====up -0.00097704{up00}
====u<运> 0{}
====u<到> 0{}
====u<了> 0{}
====np -0.00097704{np00}
====n<图书馆> 0{}
```

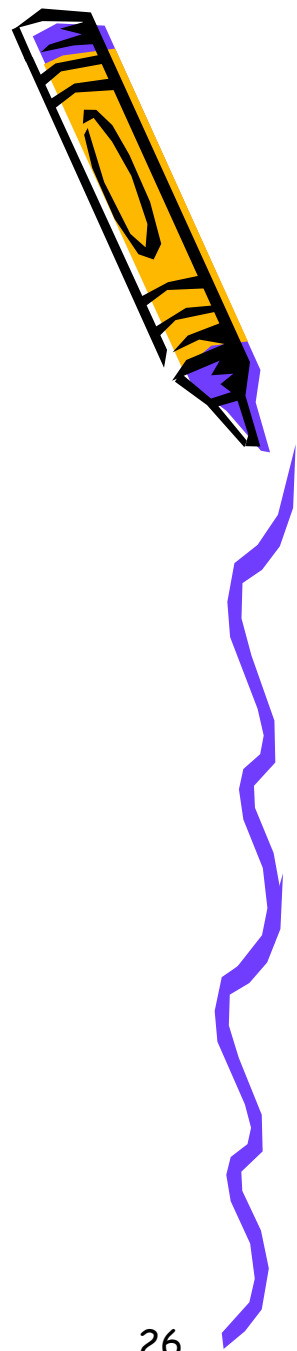
A node generated by local rule is assumed to have higher score than the node with same label that is generated by global rule. In the case of above example, the left tree will be ranked higher than the right one and be selected as the final result to output.



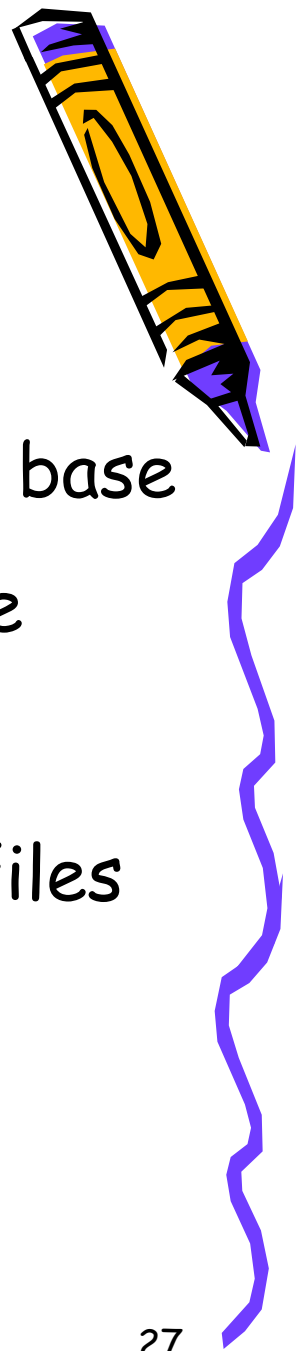
a hybrid formalism for representation of linguistic rule

## 3.2 Tools for linguistic rule writing and debugging

- Management tools for linguistic Knowledge base
- Tracing the parsing process



# Management Tools for Linguistic Knowledge Base



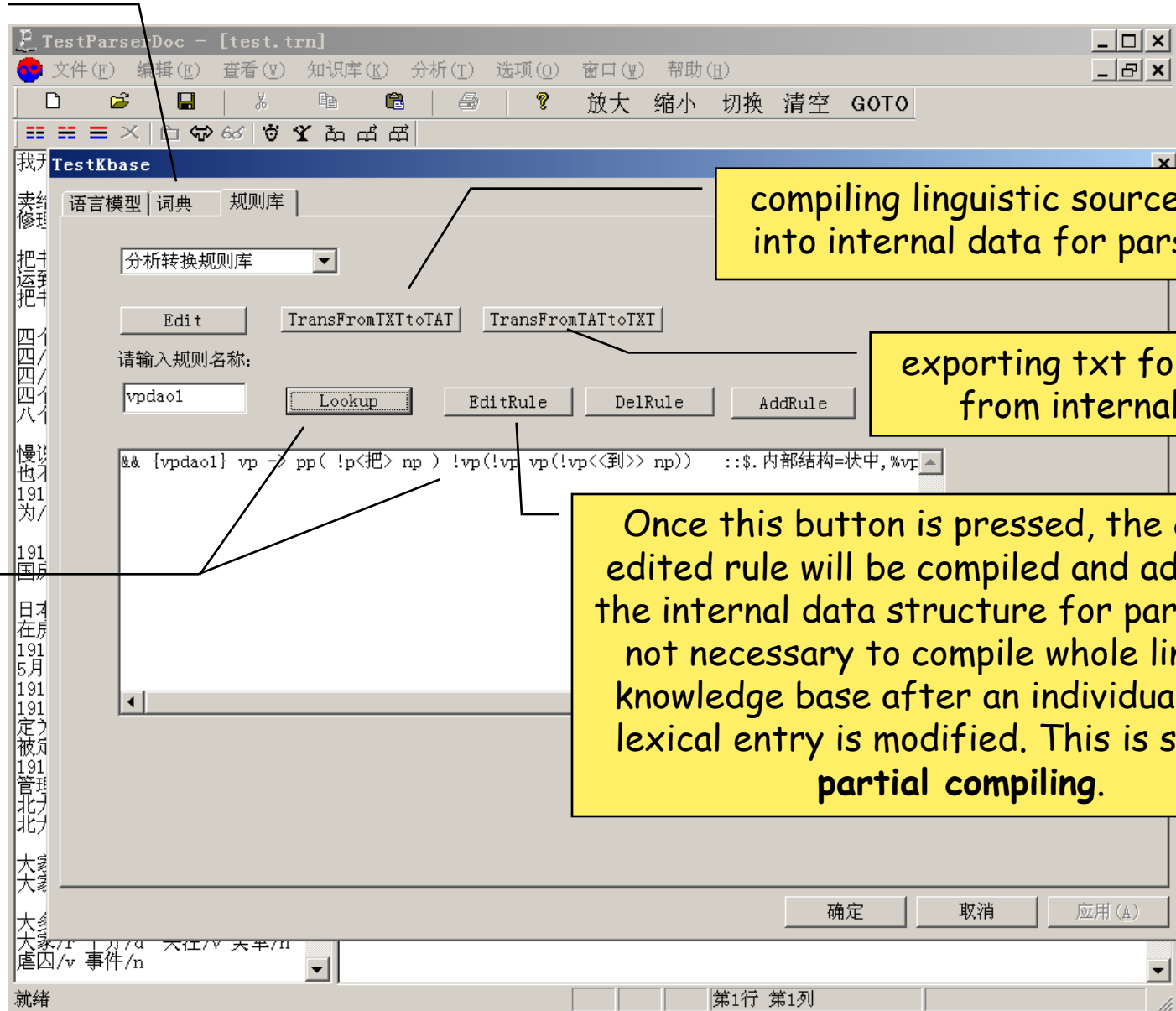
- Grammar checking for linguistic knowledge base
- The interface for linguistic knowledge base  
edit, add, delete, search, etc..
- Compiling linguistic source files to object files  
which are used in parsing indeed.



# Management Interface of LKB

Cascading style menu

Handy editing of a rule searched by it's title



compiling linguistic source file into internal data for parsing,

exporting txt format file from internal data

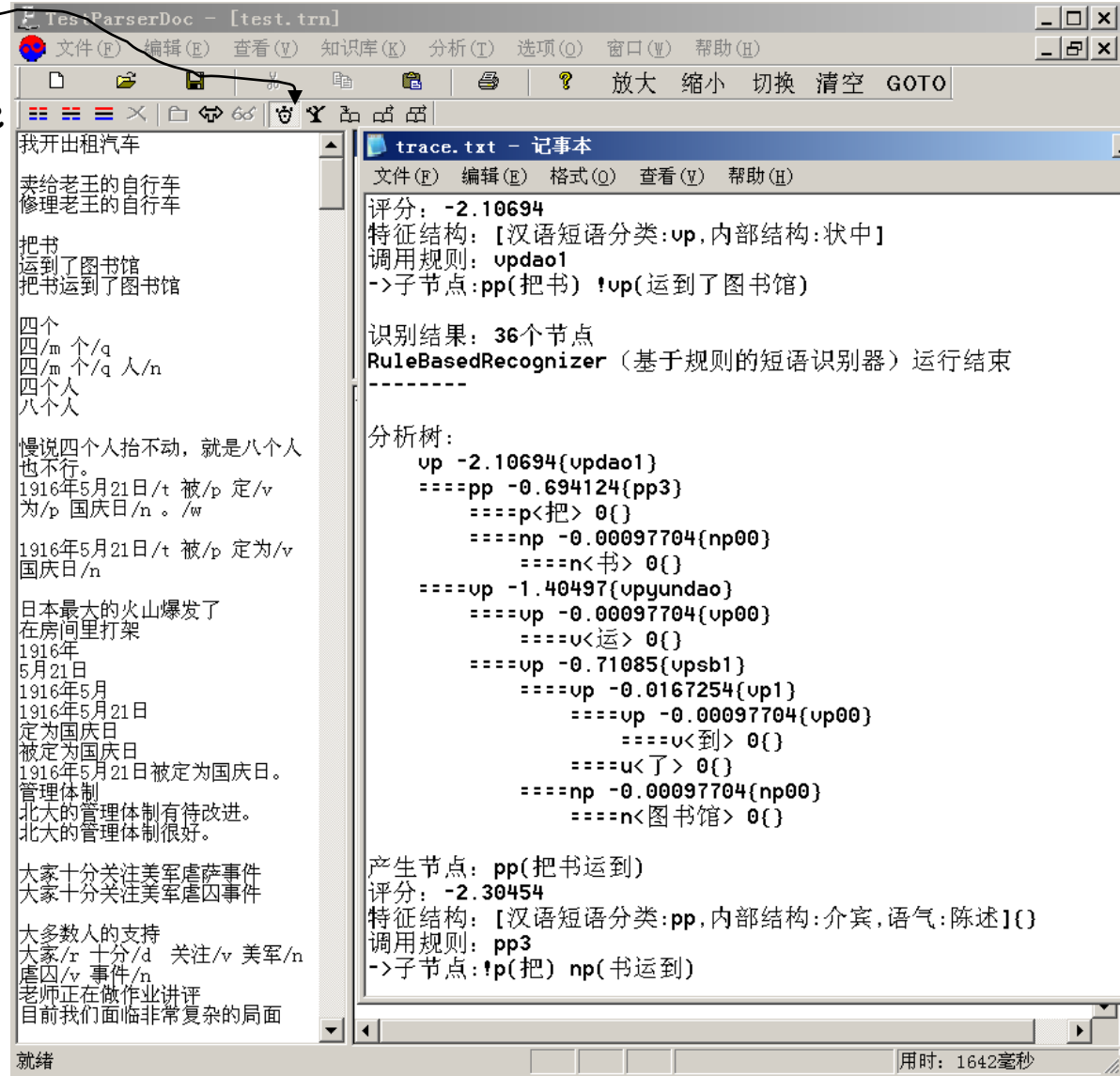
Once this button is pressed, the current edited rule will be compiled and added into the internal data structure for parsing. It's not necessary to compile whole linguistic knowledge base after an individual rule or lexical entry is modified. This is so-called **partial compiling**.

# Tracing

There is a button to trace the execution of the parser. It can be switched on and off before start of parsing. Tracing can give more insight on what the parsing behaves.

All detailed tracing information is saved in a text file which contains the following descriptions.

- The micro-engines are called in parsing
- The feature structure and the score of each node
- The rules are employed in parsing
- The parsed trees



We have planned to add the function of setting breakpoint for enforcing the debugger. But it's not available now.

## 3.3 Tools for building and using Chinese Treebank

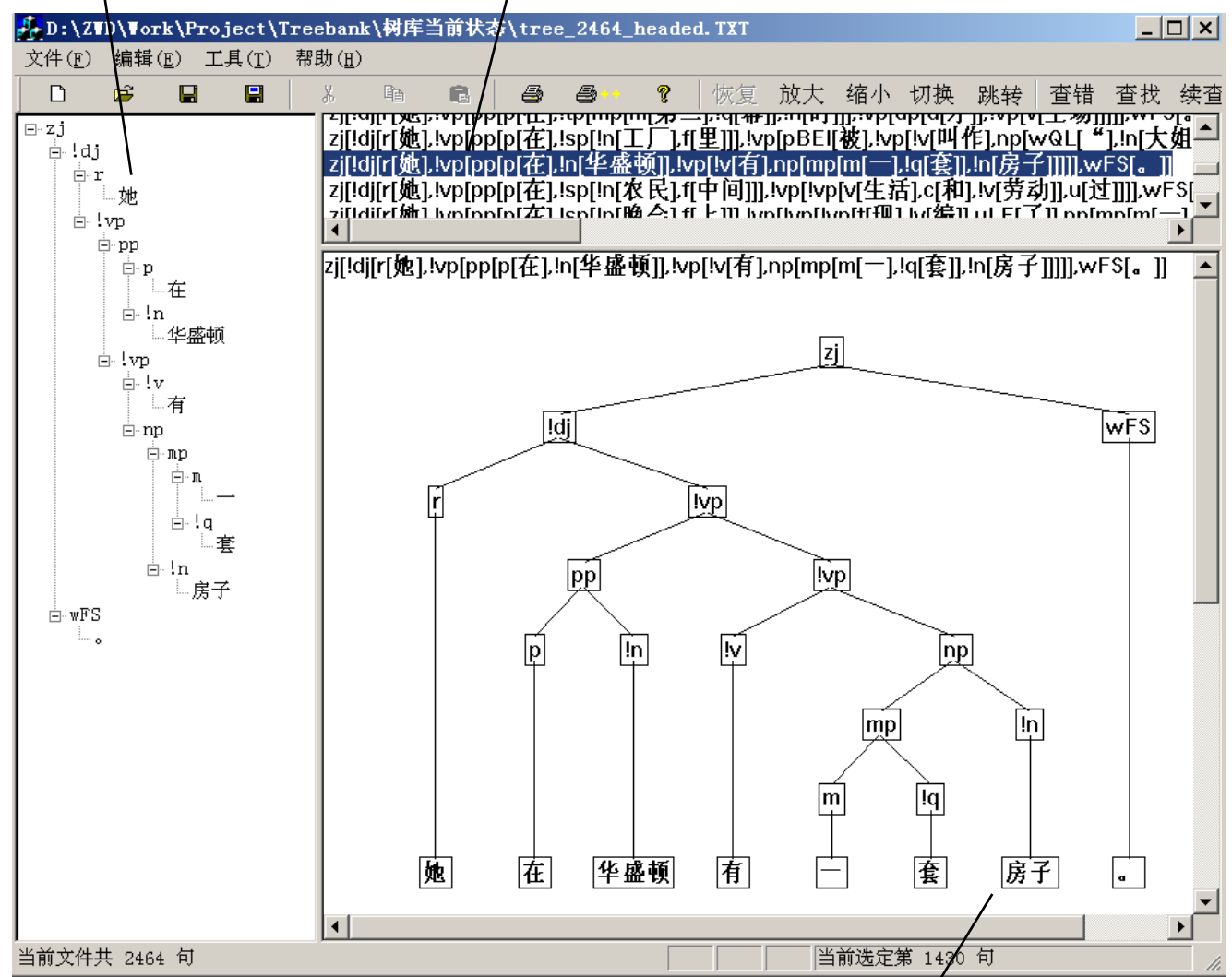
- Tree Editor
- Rule Extraction
- Search in Treebank
- Subtree replacement



# Tree Editor

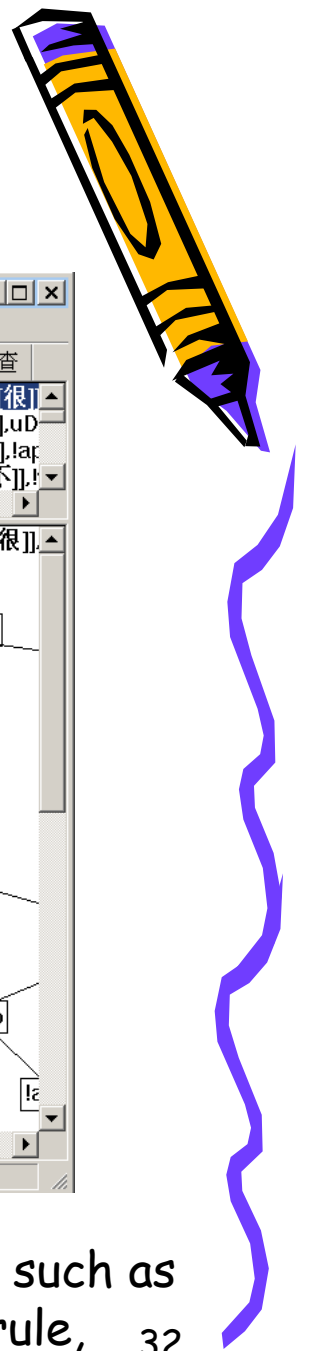
a list of trees inputted from the output of the GCG parser

tree editing area, in which user can edit manually the label of a tree node, add a new child node to a selected tree node, combine two leaf nodes, drag a selected tree node to be the child node of another target tree node, and so on.



tree viewing area, in which a parsed tree is shown as user-friendly graphic inverted tree.

# Rule extraction from a treebank



It can increase the coverage of rule set effectively to extract rules from a manually corrected treebank.

从树库中获取产生式规则，按频度排序

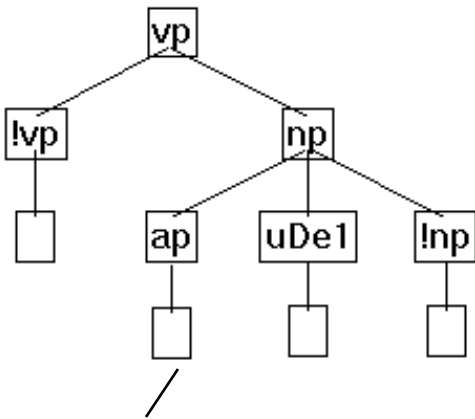
当前选定第 1 句

The rules extracted can be sorted by different criteria, such as the label of root node, the right hand side symbols of a rule, frequency of rule, number of child nodes of a root node



# Search in treebank

The tree you want to search



Note: The white square can match any tree node

树库高级查找

根节点:  OK

子节点:  Cancel

比如: mp[m,q],np[ap[a],np]

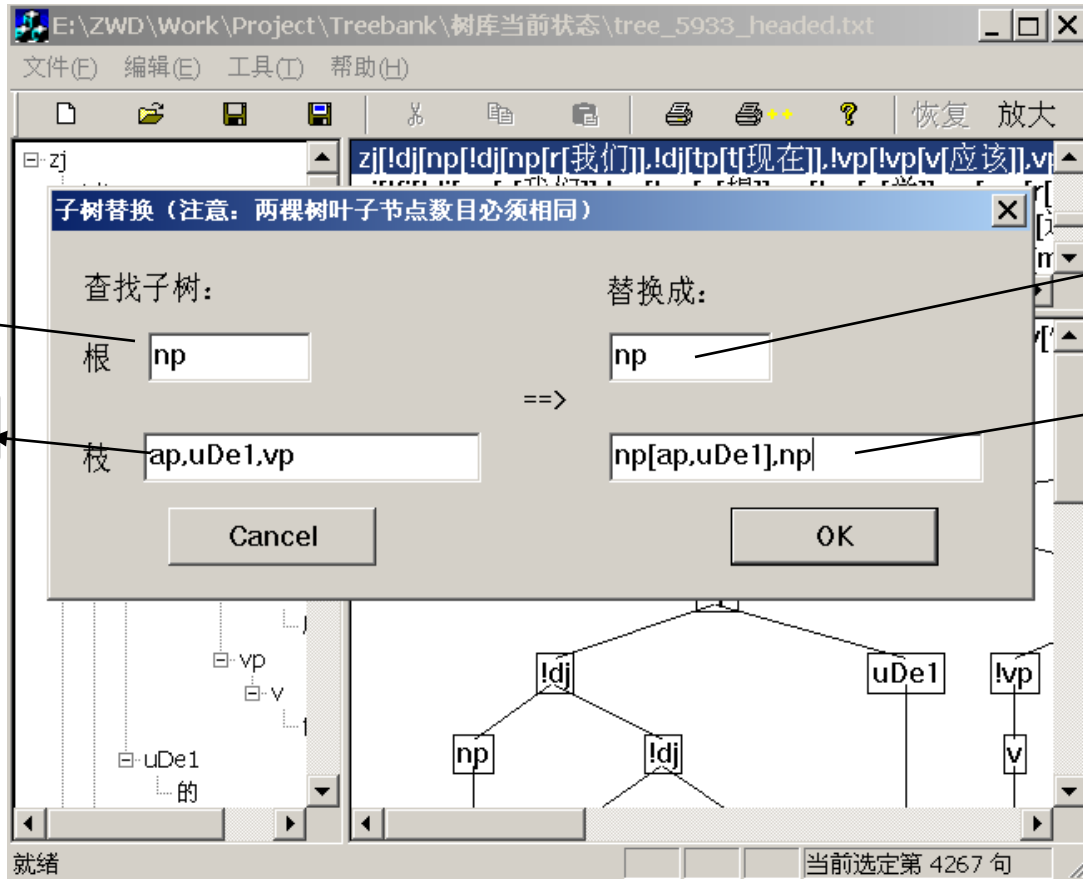
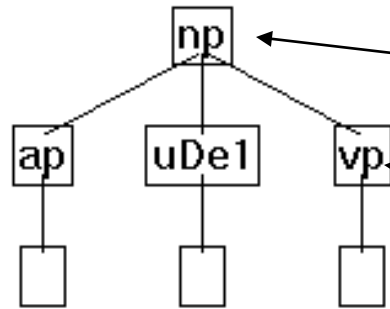
specify the label of root node you want to search

specify the child nodes of root node you want to search

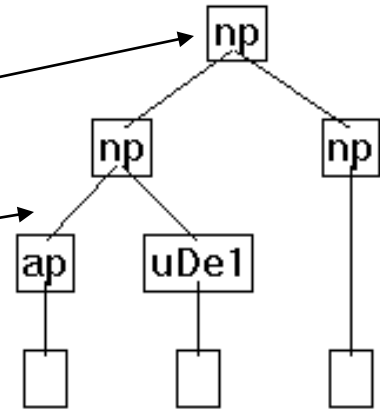
就绪 当前选定第 4267 句

# Subtree replacement

find the tree below that is embedded in larger trees



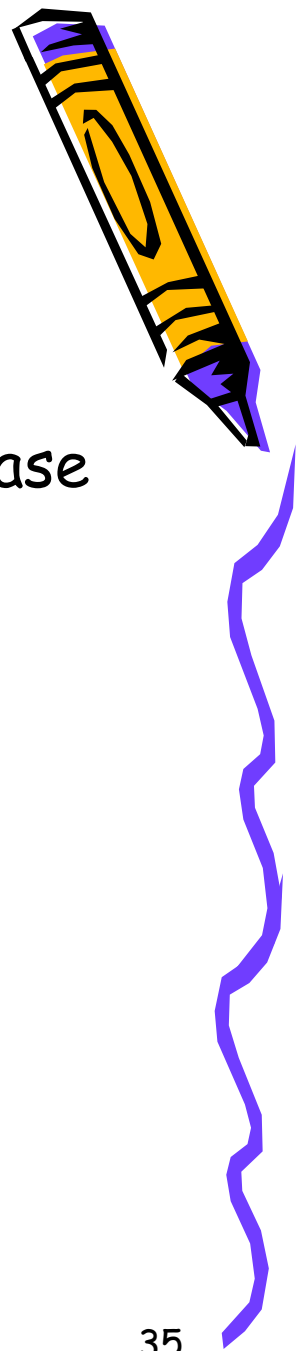
replace the found tree(s) with the tree below



You can type \* as both root node and leaf node that matches any label for a tree node when you input a tree for searching and replacing.

# 4 Experiences with ICGDE

- Statistics of current linguistic knowledge database
- Precision and efficiency of the parser
- Statistics of current PKU Chinese treebank



# 4.1 Statistics of current linguistic knowledge database



- Core lexicon: more than 43,000 entries with rich syntactic and semantic information in attribute-value format
- Appended lexicon: more than 200,000 entries with just part-of-speech tag. Entries in this lexicon are almost compound word or multi-words unit.
- Phrase database: more than 30,000 entries which are annotated with part-of-speech of words that constructs a phrase and the functional category and internal structure of the phrase
- Rule set: more than 900 rules, including 330 global rules and 577 local rules



## 4.2 Precision and efficiency of the parser

- The test set contains 500 sentences, 5000 Chinese characters (4028 words).
- The average sentence length is 10 characters per sentence (or 8 words per sentence).
- The computer used for testing is configured with 2.60GHz Pentium 4 processor, 512MB SDRAM

Large-scale Extended Lexicon engine	Phrase lexicon engine	Rule-based engine	Parsing time (ms)	Bracketing Recall	Bracketing Precision	Complete match	Average crossing	No crossing	2 or less crossing	Tagging accuracy
+	+	+	151678	67.07%	74.18%	6.37%	1.27	49.06%	79.03%	70.37%
+	-	+	114966	67.00%	74.14%	6.39%	1.27	48.87%	78.95%	70.37%
-	+	+	66345	65.06%	73.34%	5.87%	1.27	47.20%	80.80%	69.98%
-	-	+	63842	65.06%	73.34%	5.87%	1.27	47.20%	80.80%	69.98%

The above statistics is produced by Dr. Chang Baobao who used the bracket scoring program developed by Dr. Satoshi Sekine. The program is free for downloading online.

## 4.3 Statistics of current PKU Chinese treebank

- There are 14,736 bracketed sentences that were corrected manually in the treebank.  
9.628 words on average each sentence
- The treebank contains 12,685 word types and 141,879 word tokens now.
- 2477 rules can be extracted from the treebank



# 5 Future works

The ICGDE and the linguistic knowledge base described above can continue to be actively developed.

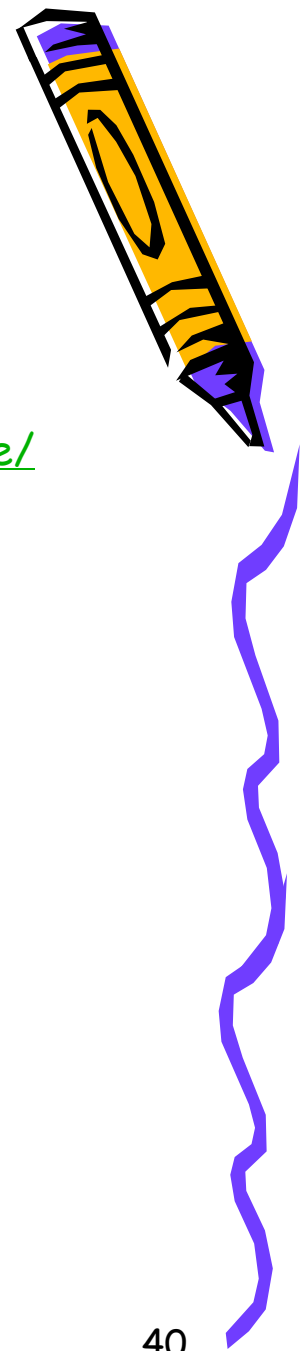
Our planned research is focused in the following areas:

- improving ability of current debugging tool, including to set breakpoint, record more tracing data while parsing, etc..
- extending the current rule set and lexicon through acquisition from corpora
- enhancing functionality of current management tool for linguistic knowledge base, especially focusing on environment of collaborative grammar coding
- developing evaluation technology for integration of existing linguistic resources

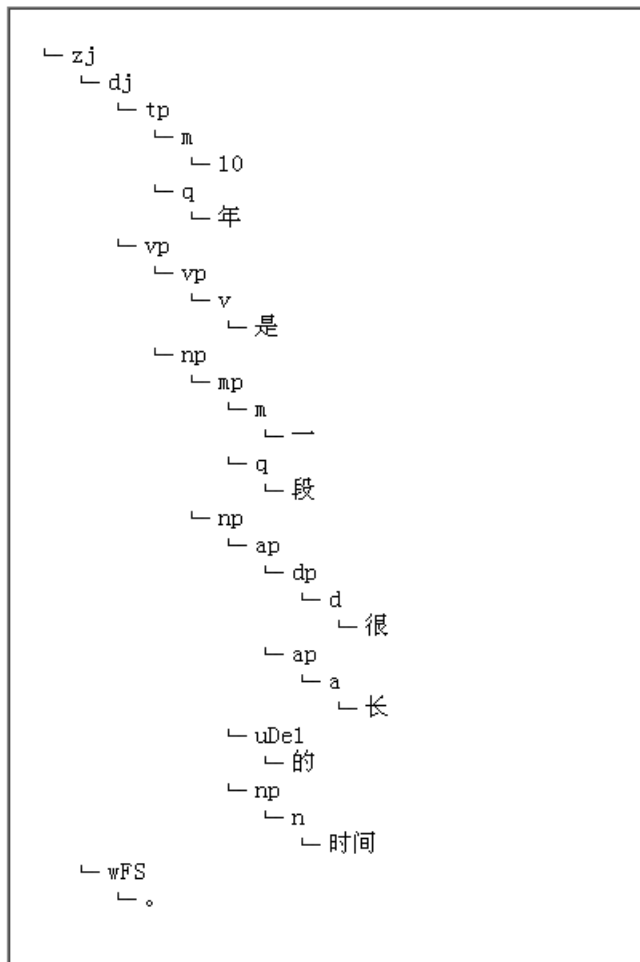


# Appendix: Resources Online

- PKU Chinese Treebank Online  
<http://ccl.pku.edu.cn/doubtfire/projects/chinesesentencestructure/>
- Online Chinese Semantic Lexicon  
[http://ccl.pku.edu.cn/ccl\\_sem\\_dict/](http://ccl.pku.edu.cn/ccl_sem_dict/)
- Online Chinese Balanced Corpus with search engine support  
[http://ccl.pku.edu.cn/ccl\\_corpus/jsearch/](http://ccl.pku.edu.cn/ccl_corpus/jsearch/)



# 北大中文树库 (Chinese Treebank) 工程



浏览第 1 页 / (共17页, 423句)

[后页](#) [尾页](#)

- 1) 10年是一段很长的时间。
- 2) 1991年, 感染霍乱的病人多达3026人, 其中440人死亡。
- 3) 啊, 这是多么美妙的前景!
- 4) 爱国一家, 爱国不分先后。
- 5) 安葬他的地方很美。
- 6) 安装灯的人是我的同学。
- 7) 安娜自己打开了门。
- 8) 八成他不来了。
- 9) 八点时, 他正在吃早饭。
- 10) 八减三得五。
- 11) 八月份物价将上涨。
- 12) 把东西清理干净
- 13) 把花盆搬到外面去。
- 14) 把画挂到墙上!
- 15) 把火拨一拨
- 16) 把课文再读一遍!
- 17) 把龙头打开
- 18) 把那本书扔给我
- 19) 把那本杂志扔给我
- 20) 把你的名字告诉我
- 21) 把你的衣服挂起来
- 22) 把瓶子放在桌上。
- 23) 把其余的吃的留到明天。
- 24) 把全部情况告诉她
- 25) 把身子探出窗外是很危险的。

[后页](#) [尾页](#)

# Online Chinese Semantic Lexicon (CCL, PKU)

Entries can be accessed for browsing, editing, searching via Internet. Authorized users can add new entries and delete entries.

北京汉语语言学研究中心 - Microsoft Internet Explorer

文件(F) 编辑(E) 查看(V) 收藏(A) 工具(T) 帮助(H)

地址(D) http://ccl.pku.edu.cn/ccl%5Fsem%5Fdict/dict\_view.asp?PartOfSpeech=noun&group=2

北京汉语语言学研究中心  
Center for Chinese Linguistics PKU  
2004年10月31日 星期天 欢迎光临CCL

第270575位访客 English | 首页 | 中心概览 | 语言学资源 | 词典 | 语料库 | 课程 | 论坛 | 留言 | 调查 | 搜索 | 收藏

## 浏览名词词典

编辑	词语	拼音	义项编码	释义	语义类	配价数	参照1
编辑	阿弟	a1di4	1		关系	1	个人
编辑	阿爹	a1die1	1		关系		
编辑	阿斗	a1dou3	1	比喻懦弱无能的人	个人		
编辑	阿飞	a1fei1	1		身份		
编辑	阿哥	a1ge1	1		关系	1	个人
编辑	阿訇	a1hong1	1	伊斯兰教主持教仪、讲授经典的人	职业		
编辑	阿姐	a1jie3	1		关系	1	个人
编辑	阿妈	a1ma1	1		关系	1	个人
编辑	阿曼	a1man4	1		处所		
编辑	阿妹	a1mei4	1		关系	1	个人

当前页起始词条: 1 词条总数: 23654 下一页 最后

返回首页 选取分组 高级搜索

本站简介 | 使用帮助 | 免责声明 | 《语言学论丛》 | 北京大学王力语言学奖 | 北京大学中文系 | 北京大学计算语言所 | 全国语言文字标准化技术委员会  
北京大学·北京大学汉语语言学研究中心 版权所有 Copyright  
Powered by 北京大学中文系应用语言学实验室  
建议浏览模式: Internet Explorer 5.0以上版本, 1024\*768分

请输入关键字  
Google

http://ccl.pku.edu.c...

- + 名词
  - + 事物
    - + 具体事物
      - + 生物
        - + 人
          - + 个人
            - 人名
            - 职业
            - 身份
            - 关系
          - + 团体
            - 机构
            - 人群
        - + 动物
          - 兽
          - 鸟
          - 昆虫
          - 鱼
          - 爬行动物
        - + 植物
          - 树
          - 草
          - 花
          - 庄稼
        - 微生物
      - + 非生物
        - + 构件
          - 身体构件
          - 非身体构件
        - + 抽象事物
      - + 过程
      - + 空间
      - + 时间

hierarchical semantic tree for Chinese nouns

# Online Concordance Service for Chinese Corpus

- Both modern and classical Chinese are contained in the balanced corpus that exceeds 100 million Chinese characters now.

- Support both single term query and complex multi-term query for searching which can meet most requirements for language studying.

中文\$3水平-句集搜索 - Microsoft Internet Explorer

文件(F) 编辑(E) 查看(V) 收藏(A) 工具(T) 帮助(H)

后退 搜索 收藏夹

地址(D) http://ccl.pku.edu.cn/ccl\_corpus/jsearch/search?q=%D6%D0%CE%C4%243%CE%AE%C6%BD&start=0&num=50&i 转到 链接

提示: 输入更多关键词可以获得更精确的结果

中文\$3水平 查找 在结果中查找

现代汉语  古代汉语 最多显示字数: 左 30 右 30

语料库统计信息 使用说明

搜索: "中文\$3水平" 共有 21 条结果, 这是第 1 - 21 条。 [搜索用时: 188.0 毫秒]

[下载全部结果](#)

厂长和技术骨干参加了各种形式的乡镇企业管理班学习, 126 名高中文化水平的职工被送到高等院校学习。考中国语文科的B级及格成绩; 开辟单元课程, 以提高现有公务员的中文写作水平; 所有新入职的政务职员人员, 均须在入职的当年参加中文于刻意提高中文播音水平。

考虑到我国8亿农民, 70%以上只具有初中文化水平的现实国情? 而郑培峰是一位盲女, 初中文化水平, 竟在1个月时间内学会了。

农村党支部书记达到大专文化水平, 其他农村基层干部达到中专或高中文化水平。

章的撰写人和编辑者都是台湾省籍人士, 过去主要是学习日文, 所以中文水平不高。

多种措施, 使70%原是农民的职工, 从小学文化水平全部提高到高中文化水平, 一百多人达到大专以上文化程度, 并建立了企业自己的科研只有初中文化水平的曹向东, 深知自己文化浅, 但他更相信勤能补拙。

为了改变这种高投入、低产出的状况, 只有初中文化水平的向大聪啃起了大学课本。

占到全村总人口的10%以上; 而留在村里的中年人也全部达到了初中文化水平; 由于实行农科教相结合, 全村基本劳力每人都掌握了一两项事实上, 目前眼镜店内的验光人员多只有小学或初中文化水平, 有关机构的一次考核表明, 相当一些验光人员不会准确测出几年, 我们村的扫盲工作开展得比较深入, 现在大家普遍都有初、高中文化水平了, 有些还自学大专课程呢。

堂情况, 一个当老师, 一个当学生, 你问我答, 互相校正发音, 提高中文水平。

他这时的中文水平, 已能让他一口气把这封信读完。

时候, 她也犹豫过, 要想进入写戏的这个艺术殿堂之中, 仅仅具有初中文化水平的她, 能否取得成绩, 她没有把握。

中文口语水平相当不错。

两天后周总理再次召开会议, 张与丁均感到维特克的中文水平较低, 对中国近现代史知之甚少, 连中文报纸也看不大懂, 但对只有初中文化水平的小卢哪里听得懂啊!

她具有初中文化水平, 又好看, 深得婆婆疼爱, 尽管小姑子们年龄都比她大, 但都中文口语水平相当不错。

1

提示: 输入更多关键词可以获得更精确的结果

中文\$3水平 查找 在结果中查找

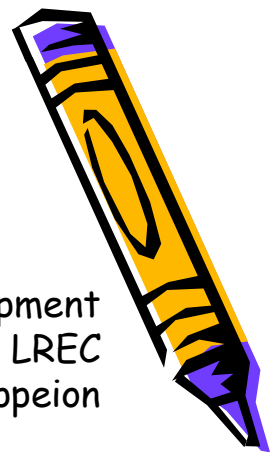
现代汉语  古代汉语 最多显示字数: 左 30 右 30

Copyright (c) 2003 北京大学汉语语言学研究中心  
Powered by [WebLucene](#) on [Lucene](#)

完毕 Internet

# Reference

- Ann, Copestake & Flickinger Dan (2000), An Open Source Grammar Development Environment And Broad-coverage English Grammar Using HPSG, In Proceedings of LREC 2000 (The 2nd International Conference on Language Resource & Evaluation), Zappeion Megaron, Greece, May 31 – June 2, 2000.
- Blevins, James(2003), Feature-based Grammar, In Borsley, R.D. & Borjars, K. eds., Non-transformational Syntax, Oxford: Blackwell, to be published in 2005.
- Blache, Philippe, Marie-Laure Guignon, Tristan van Rullen (2003), A Corpus-based Technique for Grammar Development, In Proceedings of The Shallow Processing of Large Corpora Workshop (SProLaC 2003), Lancaster University (UK), 27 March, 2003.
- Borsley, Robert D., 1996, *Modern Phrase Structure Grammar*, No. 11 in Blackwell textbooks in Linguistics, Blackwell Publishers Inc..
- Chen, Feng-yi, et al. 1999, Sinica Treebank, *Computational Linguistics and Chinese Language Processing*, 4(2):183-204
- Chen, Keh-Jiann & Yu-Ming Hsieh, 2002, Chinese Treebanks and Grammar Extraction, *CJNLP'2002*, Peking University, 2002.10.30-11.2
- Erbach, Gregor (1991), A Flexible Parser for a Linguistic Development Environment, In O. Herzog & C.-R. Rollinger eds., *Text Understanding in LILOG*, Springer, 1991, pp. 74-87



# Reference

- Heinecke, Johannes, Jurgen Kunze, Wolfgang Menzel, and Ingo Schroder (1998), Elimiative parsing with graded constraints. In Proceedings of 17th International Conference on Computational Linguistics, 36th Annual Meeting of the ACL, Coling-ACL '98, Montreal, Canada, 1998.
- Knight, Kevin, 1989, Unification: A Multidisciplinary Survey, *ACM Computing Surveys*, Vol.21, No.1.
- Sag, Ivan A. & Thomas Wasow, 1999, *Syntactic Theory: A Formal Introduction*, CSLI Publications, Stanford, California.
- Schmid, Helmut (1999), YAP: Parsing and Disambiguation With Feature-Based Grammar. PhD thesis, Institute of Maschinelle Sprachverarbeitung, University Stuttgart, Germany, 1999.
- Suzuki, Hisami (2002), A Development Environment for Large-scale Multi-lingual Parsing Systems, In Workshop on Grammar Engineering and Evaluation (Post-conference workshop in conjunction with COLING-2002, Taipei, Sept. 1, 2002).
- Uszkoreit, Hans (2002), New Chances for Deep Linguistic Processing, *Coling2002*, Taipei.
- Volk, Martin & Dirk Richarz (1997), Experiences with the GTU Grammar Development Environment, *ACL workshop on Environments for Grammar Development*, 1997, Madrid, Spain.
- Xu Ruifeng, et al., 2004, The Construction of A Chinese Shallow Treebank, *ACL2004 SACL Workshop*, July 21-16, 2004. Barcelona, Spain.

詹文滢 (2000), 《面向中文信息处理的现代汉语短语结构规则研究》, 清华大学出版社, 广西科学技术出版社。



Thank you  
for your attention

Welcome to  
<http://ccl.pku.edu.cn>

