

# 计算语言学

## 第 11 讲 机器翻译

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 2012 年春季课程讲义

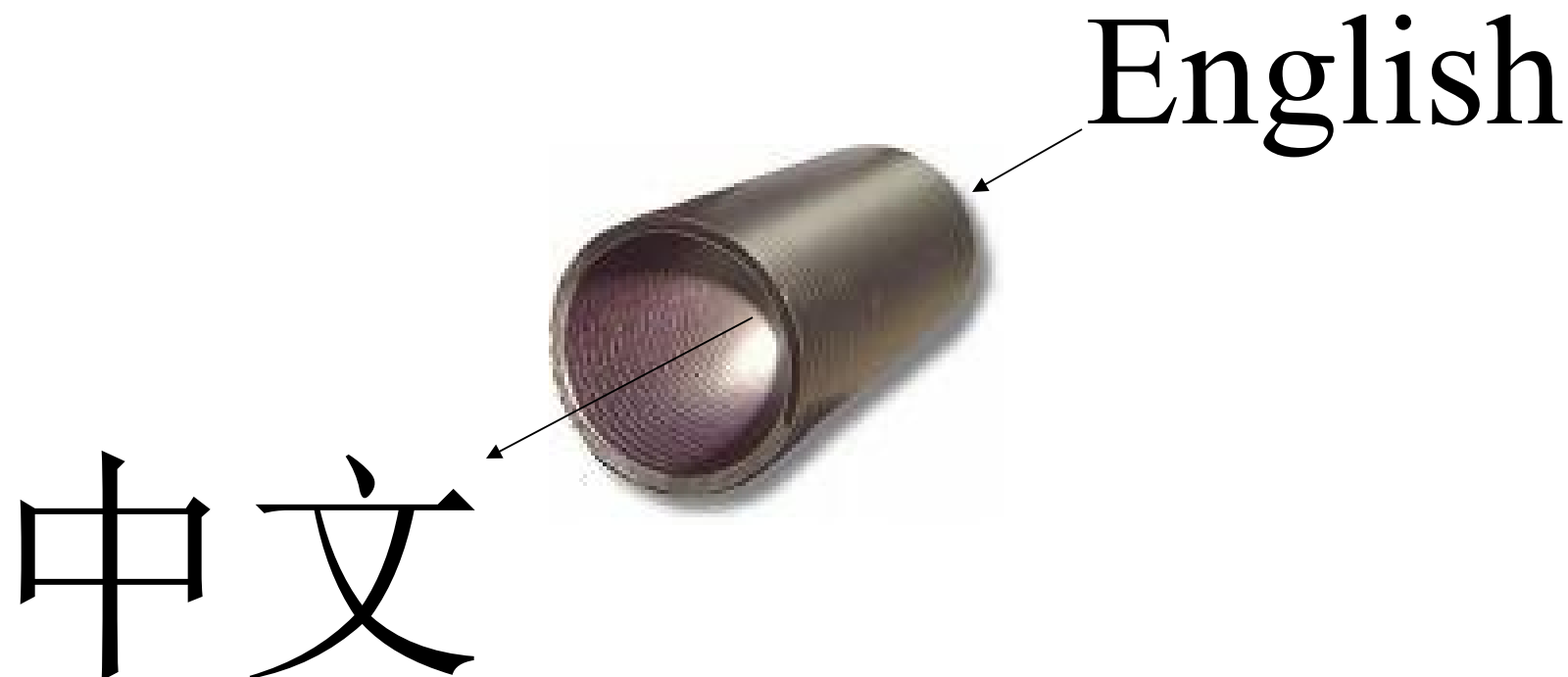
# 内容提要

机器翻译概述

机器翻译方法

机器翻译评价

# 什么是机器翻译



# 什么是机器翻译

- 机器翻译 (machine translation) 是使用电子计算机把一种自然语言 ( 源语言 ,source language) 翻译成另外一种自然语言 ( 目标语言 ,target language) 的一门学科
- 这门学科同时也是一种技术 . 它涉及到语言学、计算机科学、数学等许多部门 , 是非常典型的多边缘的交叉学科
  - 在语言学中 , 机器翻译是计算语言学的一个研究领域
  - 在计算机科学中 , 机器翻译是人工智能的一个研究领域
  - 在数学中 , 机器翻译是数理逻辑和形式化方法的一个研究领域 .

以上定义引自冯志伟《**澄清对机器翻译的一些误解（论文提要）**》，现代语文（语言研究）， 2005.1，做了个别修改

# 机器翻译面临的问题

- 从语言学角度看
  - 自然语言理解的所有问题都会对机器翻译产生影响
  - 要做到完美的机器翻译翻译，必须对自然语言处理所有层面的歧义进行消解
  - 机器翻译的难度等同于自然语言理解的难度
- 从操作角度看，机器翻译可以简单理解为只有两个问题：
  - 为每个词寻找合适的译文词
  - 给译文词排序

# 机器翻译的历史

- W. J. Hutchens, latest Development in MT Technology: Beginning a New Era in MT Research. In : Proceedings of Machine Translation Summit-IV, Kobe, Japan, 1993.
- 冯志伟, 自动翻译, 上海知识出版社, 1987 年
- 冯志伟, 自然语言机器翻译新论, 语文出版社, 1994 年
- 冯志伟, 自然语言的计算机处理, 上海外语教育出版社, 1996 年

# 机器翻译的历史

- 萌芽期（17 世纪 -1930 年代）
- 草创期（1946-1964）
- 萧条期（1964-1960 年代后期）
- 复苏期（1970 年代初期）
- 繁荣期（1970 年代后期 -1980 年代初期）
- 平台期（1980 年代后期 -1999 年）
- 再度繁荣期（1999- 现在）统计方法！

# 内容提要

机器翻译概述

机器翻译方法

机器翻译评价

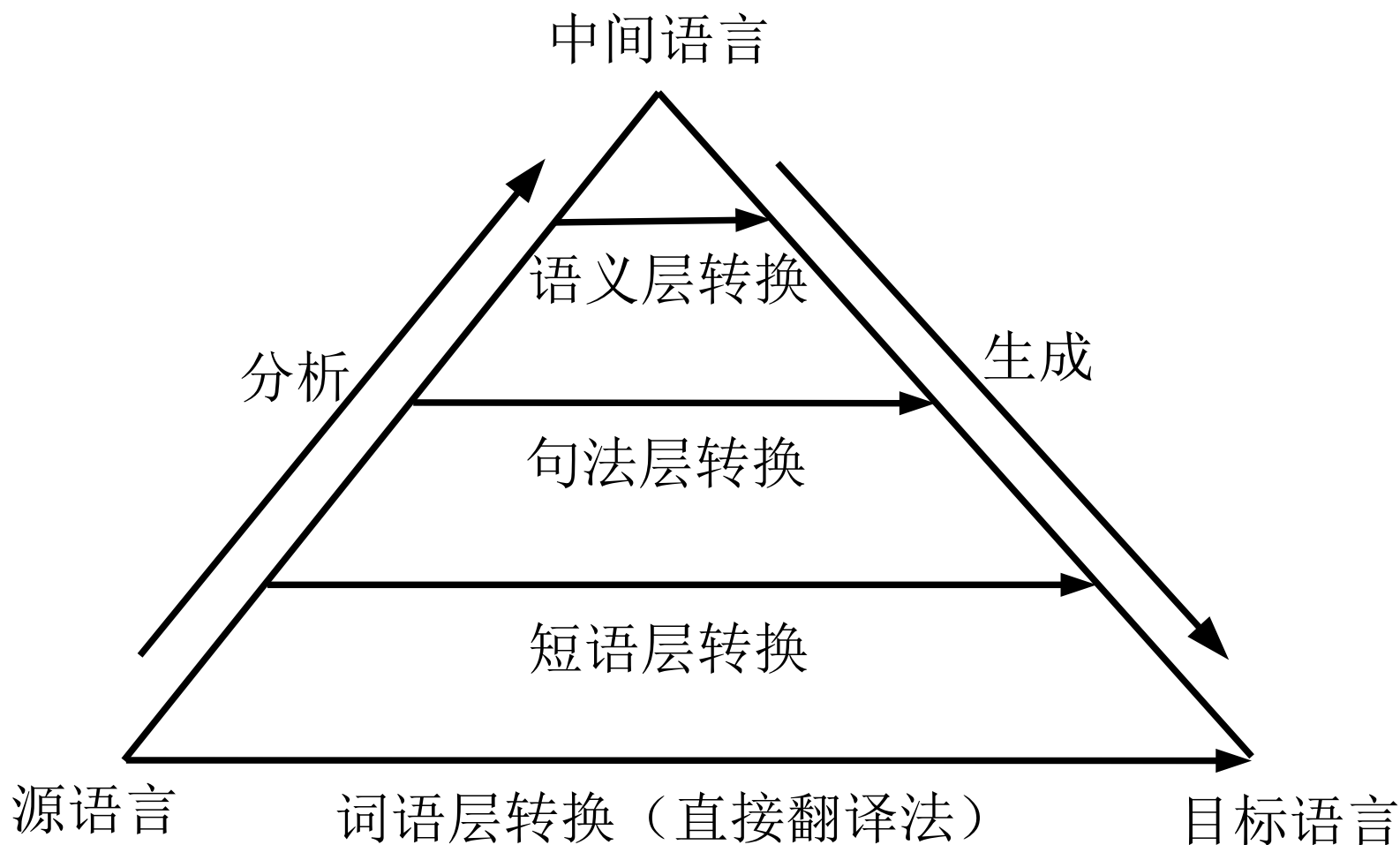


# 机器翻译方法

按转换层面划分

按知识表示划分

# 机器翻译的转换层面



# 直接翻译方法

- 通过词语翻译、插入、删除和局部的词序调整来实现翻译，不进行深层次的句法和语义的分析，但可以采用一些统计方法对词语和词类序列进行分析
- 早期机器翻译系统常用的方法，近期 **IBM** 提出的统计机器翻译模型也可以认为是采用了这一范式
- 著名的机器翻译系统 **Systran** 早期也是采用这种方法，后来逐步引入了一些句法和语义分析

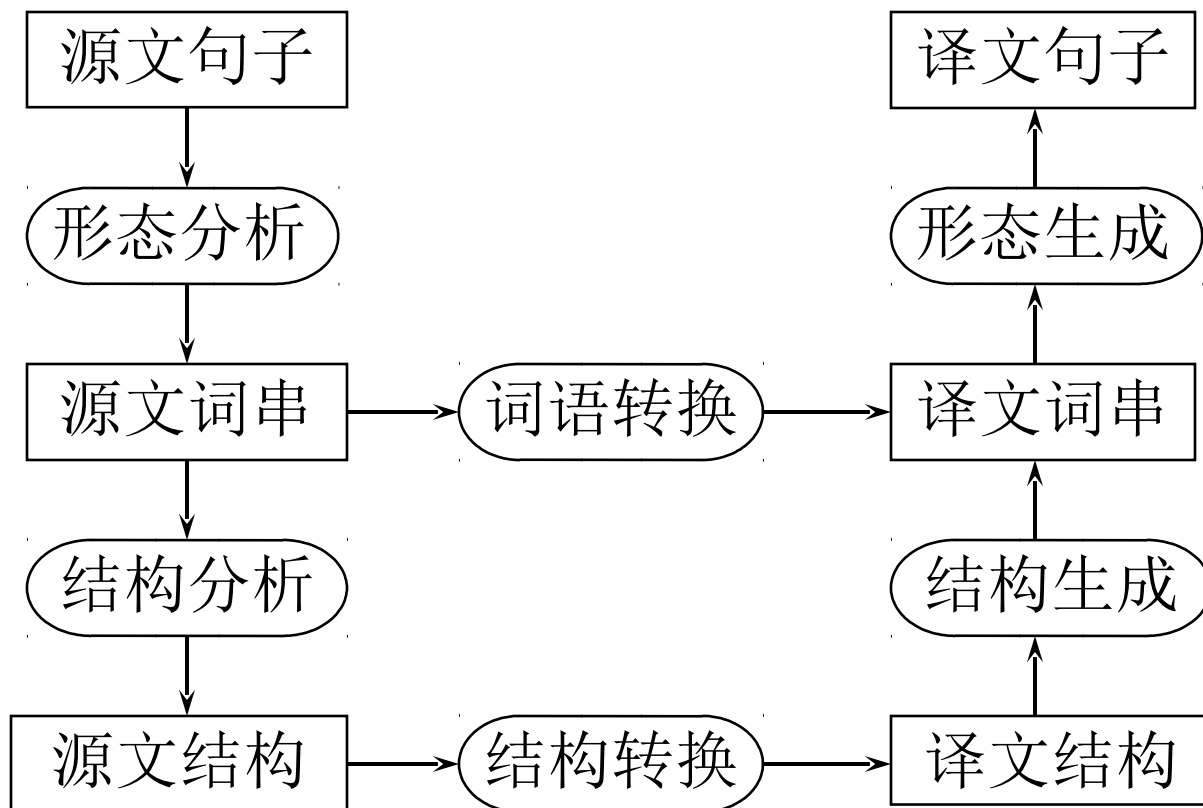
# 转换方法 (1)

- 整个翻译过程分为“分析”、“转换”、“生成”三个阶段；
- 分析：源语言句子→源语言深层结构
  - 相关分析：分析时考虑目标语言的特点
  - 独立分析：分析过程与目标语言无关
- 转换：源语言深层结构→目标语言深层结构
- 生成：目标语言深层结构→目标语言句子
  - 相关生成：生成时考虑源语言的特点
  - 独立生成：生成过程与源语言无关

# 转换方法 (2)

- 理想的转换方法应该做到独立分析和独立生成，这样在进行多语言机器翻译的时候可以大大减少分析和生成的工作量；
- 转换方法根据深层结构所处的层面可分为：
  - 句法层转换：深层结构主要是句法信息
  - 语义层转换：深层结构主要是语义信息
- 分析深度的权衡
  - 分析的层次越深，歧义排除就越充分
  - 分析的层次越深，错误率也越高

# 转换方法 (3)



基于转换方法的翻译流程

# 例：句法层面的转换方法 (1)

她把一束花放在桌上。  $\Longrightarrow$  She put a bunch of flowers on the table.

切分 / 标注

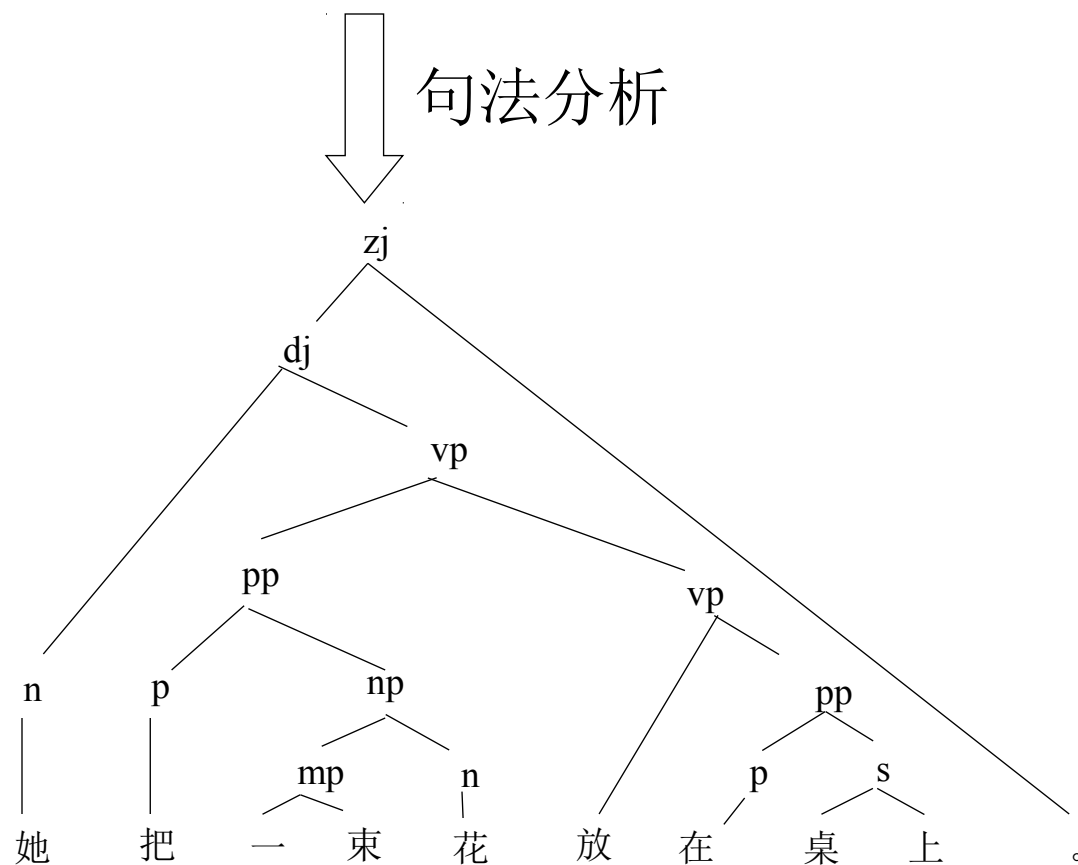
她 /r 把 /p-q-v-n 一 /m-d 束 /q 花 /n-v-a 放 /v 在 /p-d-v 桌 /n 上 /f-v 。 /w

标注排歧

她 /r 把 /p 一 /m-d 束 /q 花 /n 放 /v 在 /p-v 桌 /n 上 /f-v 。 /w

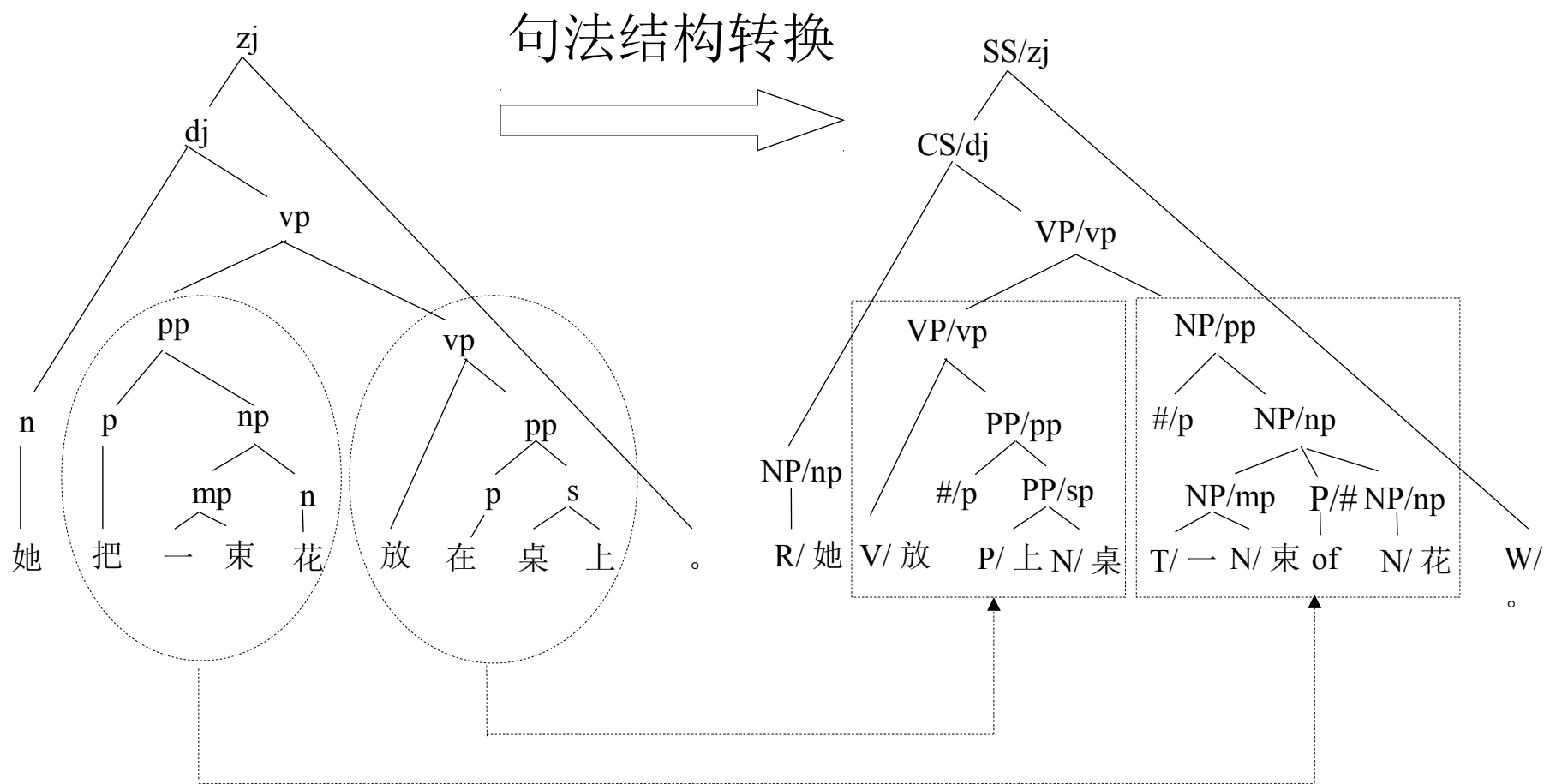
# 例：句法层面的转换方法 (2)

她 /r 把 /p 一 /m-d 束 /q 花 /n 放 /v 在 /p-v 桌 /n 上 /f-v 。 /w

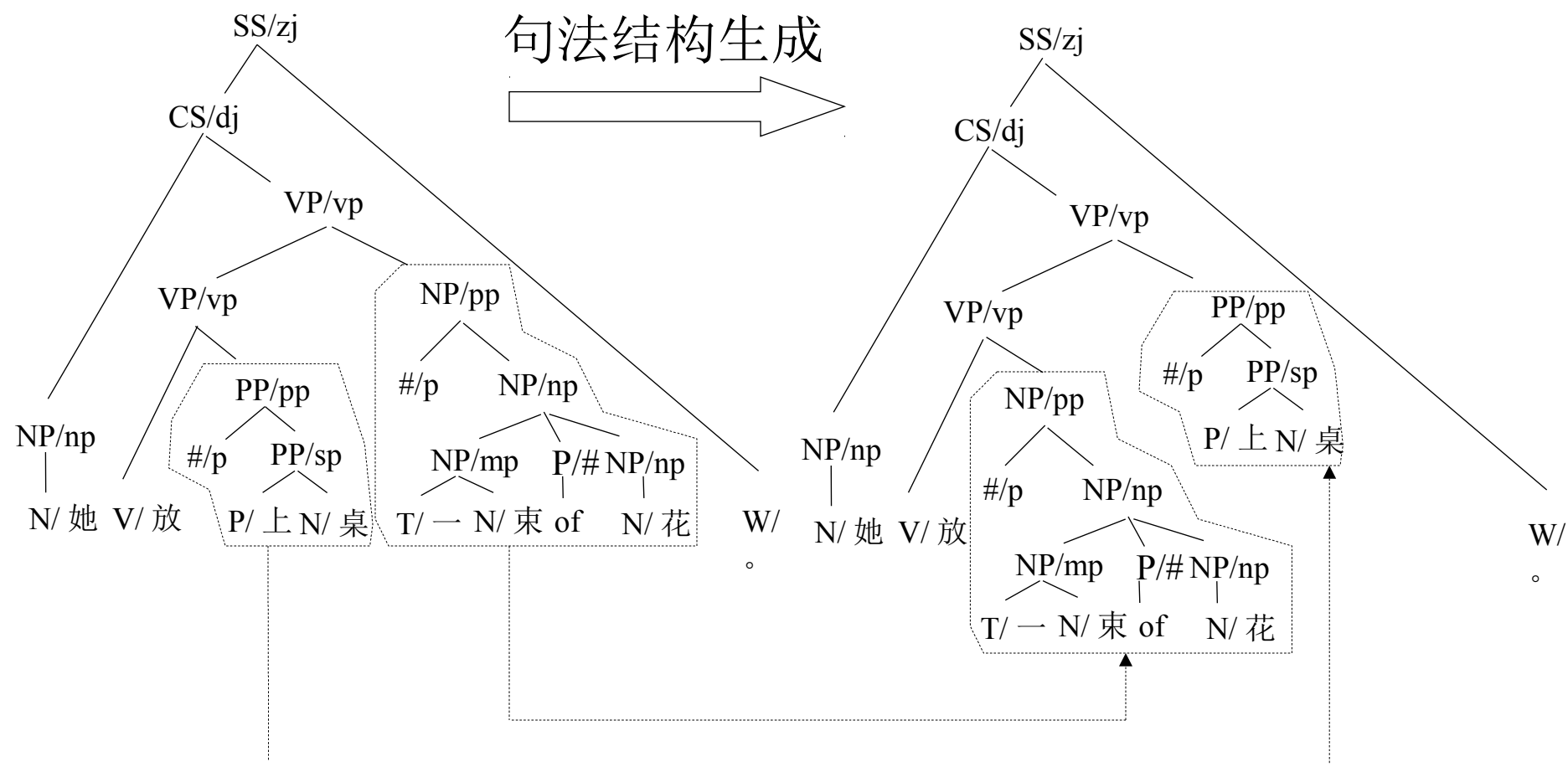




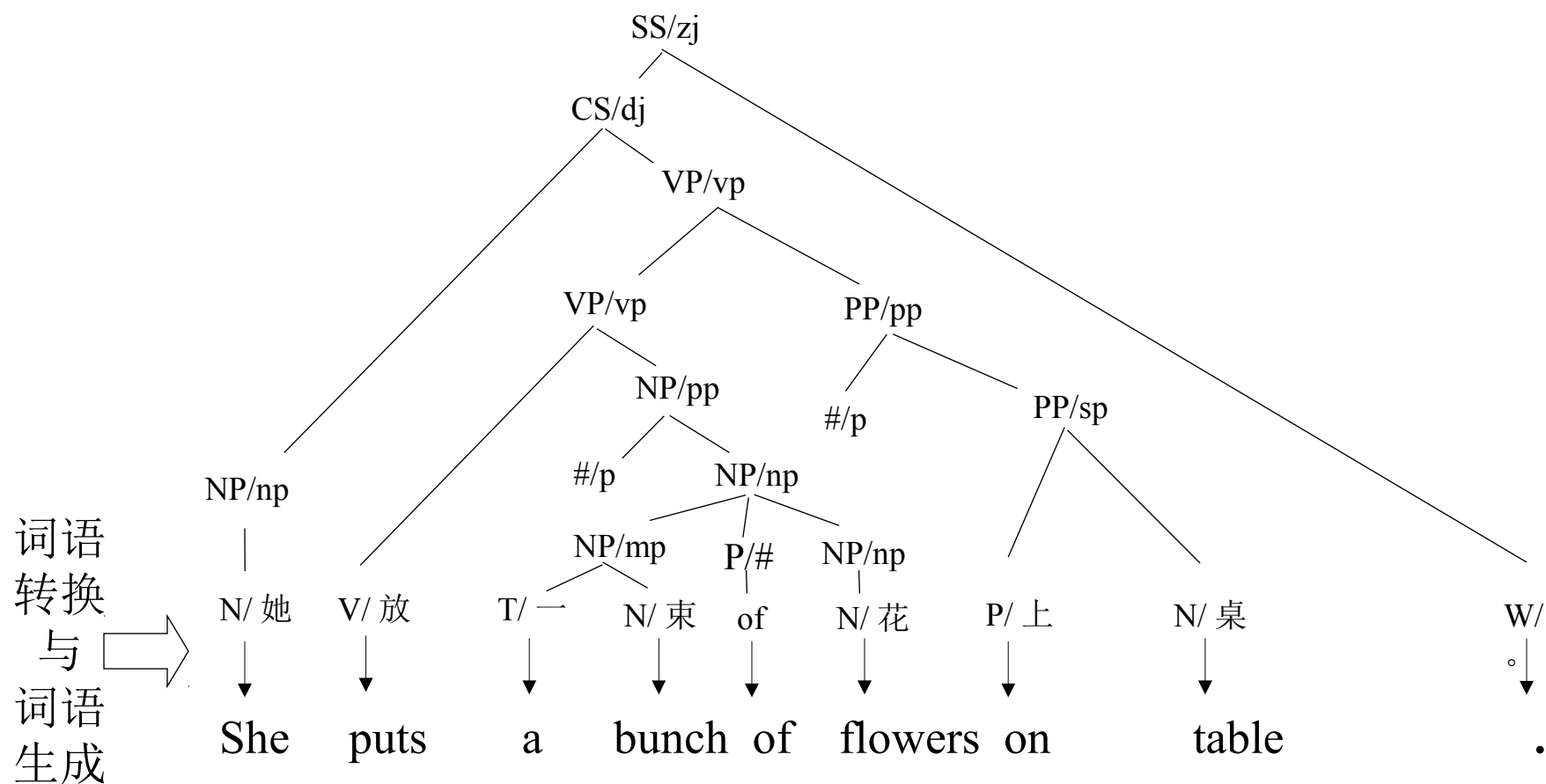
# 例：句法层面的转换方法 (3)



# 例：句法层面的转换方法 (4)



# 例：句法层面的转换方法 (5)



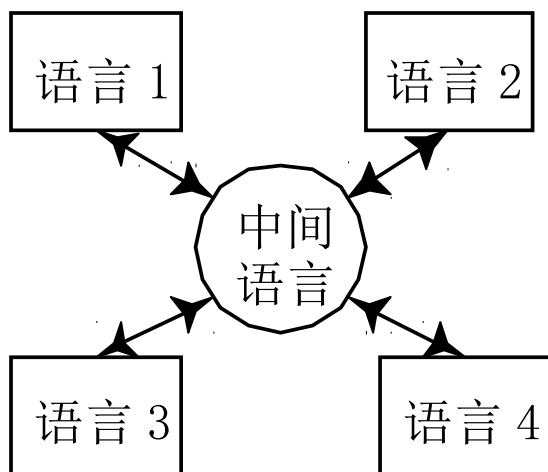
# 中间语言方法 (1)

- 利用一种中间语言（ interlingua ）作为翻译的中介表示形式；
- 整个翻译的过程分为“分析”和“生成”两个阶段
- 分析：源语言→中间语言
- 生成：中间语言→目标语言
- 分析过程只与源语言有关，与目标语言无关
- 生成过程只与目标语言有关，与源语言无关

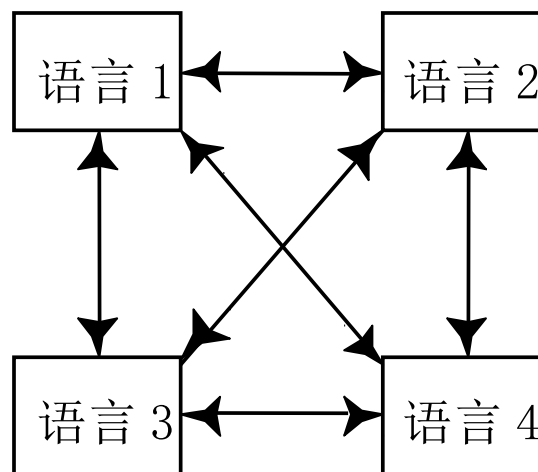
## 中间语言方法 (2)

- 中间语言方法的优点在于进行多语种翻译的时候，只需要对每种语言分别开发一个分析模块和一个生成模块，模块总数为  $2*n$ ，相比之下，如果采用转换方法就需要对每两种语言之间都开发一个转换模块，模块总数为  $n*(n-1)$

# 中间语言方法 (3)



中间语言方法



转换方法

# 中间语言方法 (4)

- 中间语言通常并不是一种真正的语言，而是某种知识表示形式，比如语义网络、逻辑表达式等等
- 基于中间语言的机器翻译方法有时也称为基于知识的机器翻译方法

# 中间语言方法 (5)

- Makoto Nagao (Kyoto University) said: “.. when the pivot language [i.e. interlingua] is used, the results of the analytic stage must be in a form which can be utilized by all of the different languages into which translation is to take place. This level of subtlety is a practical impossibility.” (Machine Translation, Oxford, 1989)
- Patel-Schneider (METAL system) said: “METAL employs a modified transfer approach rather than an interlingua. If a meta-language [an interlingua] were to be used for translation purposes, it would need to incorporate all possible features of many languages. That would not only be an endless task but probably a fruitless one as well. Such a system would soon become unmanageable and perhaps collapse under its own weight.” (A four-valued semantics for terminological reasoning, Artificial Intelligence, 38, 1989)



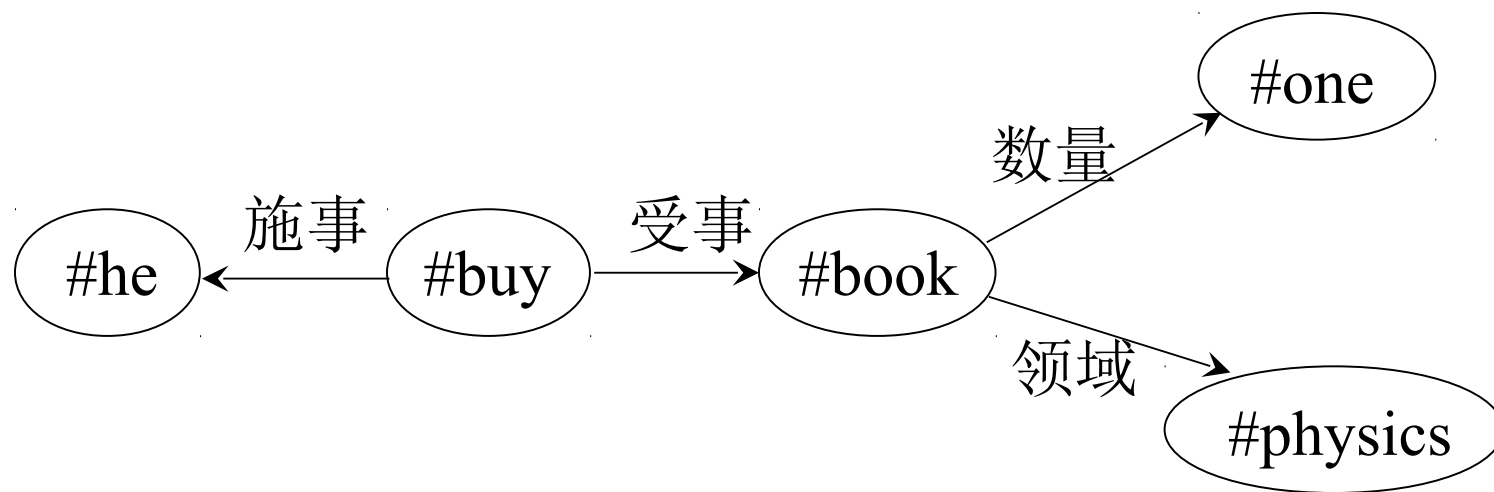
# 中间语言方法 (6)

- 基于中间语言方法一般都用于多语言的机器翻译系统中。
- 从实践看，采用某种人工定义的知识表示形式作为中间语言进行多语言机器翻译都不太成功，如日本主持的亚洲五国语言机器翻译系统，总体上是失败的。
- 在 **CSTAR** 多国语口语机器翻译系统中，曾经采用了一种中间语言方法，其中间语言是一种带话语信息的语义表示形式，由于语音翻译都限制在非常狭窄的领域中（如旅游领域或机票预定），语义描述可以做到比较精确，因此采用中间语言方法有一定的合理性。但该方法最终也不成功。
- 实际上，领域特别窄的场合可以采用中间语言方法。一个适合于中间语言方法的例子是数词的翻译，采用阿拉伯数字作为中间语言显然是非常合理的。

# 中间语言示例一语义网络

英语： He bought a book on physics.

汉语： 他买了一本关于物理学的书。



说明：这里 # 后面表示的是概念，而不是英语词。

# 中间语言示例一框架

英语： He bought a book on physics.

汉语： 他买了一本关于物理学的书。

谓词	概念	#buy	
	施事	概念	#he
	受事	概念	#book
		数量	#one
		领域	#physics

说明： 这里 # 后面表示的是概念，而不是英语词。

# 中间语言示例一概念词典

概念	语义类	中文词	英文	格框架
#he	指代词	他	he	
#buy	获得	买	buy	施事，受事
#book	出版物	书	book	
#physics	学科	物理	physics	
#one	数量	一	one	

# 中介语言方法

- 在多语言机器翻译中，很多研究人员开始采用某种自然语言作为中介语言（这时又称“枢纽语言”或“桥接语言”，英文是 **Pivot Language**）。中介语言也可以是人造语言，如世界语。
- 这种方法不同于前述的中间语言方法，因为这种方法中经过了两个独立的翻译过程，而这两个过程可能有各自独立的分析、转换、生成模块。而中间语言方法只有一个分析过程和一个转换过程。
- 也有文献把中介语言方法归入到中间语言方法，阅读文献的时候请注意区分这些概念的具体含义。
- 多语言统计机器翻译（如 **Google** 翻译）比较适合采用中介语言方法。主要原因是英语到其他语言的双语语料库比较容易获得，而其他语言之间的双语语料库很难获得。

# 机器翻译方法

按转换层面划分

按知识表示划分

# 机器翻译的知识表示

程序代码：直接机器翻译方法

翻译规则：基于规则的机器翻译方法

翻译实例：基于实例的机器翻译方法

模型参数：统计机器翻译方法

# 直接机器翻译方法

- 在早期的直接机器翻译方法中，翻译程序代码和翻译知识表示没有明显的区分，所有翻译知识直接以代码形式固化在翻译程序代码中，无法实现复杂的知识表示，翻译方法也相对比较简单



# 机器翻译的知识表示

程序代码：直接机器翻译方法

翻译规则：基于规则的机器翻译方法

翻译实例：基于实例的机器翻译方法

模型参数：统计机器翻译方法

# 基于规则的方法 (1)

- 采用规则作为知识表示形式
  - 重叠词规则
  - 切分规则
  - 标注规则
  - 句法分析规则
  - 语义分析规则
  - 结构转换规则（产生译文句法语义结构）
  - 词语转换规则（译词选择）
  - 结构生成规则（译文结构调整）
  - 词语生成规则（译文词形生成）

# 基于规则的方法 (2)

- 优点
  - 直观，能够直接表达语言学家的知识
  - 规则的颗粒度具有很大的可伸缩性
    - 大颗粒度的规则具有很强的概括能力
    - 小颗粒度的规则具有精细的描述能力
  - 便于处理复杂的结构和进行深层次的理解，如解决长距离依赖问题
  - 系统适应性强，不依赖于具体的训练语料

# 基于规则的方法 (3)

- 缺点

- 规则主观因素重，有时与客观事实有一定差距
- 规则的覆盖性差，特别是细颗粒度的规则很难总结得比较全面
- 规则之间的冲突没有好的解决办法（翘翘板现象）
- 规则一般只局限于某一个具体的系统，规则库开发成本太高
- 规则库的调试极其枯燥乏味

# 基于规则的方法—译词选择

\$\$ 开

**\*\*{v} v \$=[...]**

|| \$. 主体 = 是, \$. 主体 . 语义类 = 植物

→ V<bloom> \$=[...]

|| \$. 客体 = 是, \$. 客体 . 汉字 = 灯 | 机 | 器

→ V( !V<turn> D<on> ) \$=[...]

|| \$. 客体 = 是, \$. 客体 . 语义类 = 交通工具

=> V<drive> \$=[...]

|| OTHERWISE

=> V<open> \$=[...]

# 基于规则的方法—结构转换

&& {mp7} mp->r !mp :: \$. 内部结构 = 组合定中 ,...

|| %mp. 定语 . 内部结构 = 单词 , %mp. 定语 .yx= 一 , %mp. 量词子类 = 集体 | 种类 | 容量 | 时量 | 度量 | 成形

=> NP(T/r !NP/mp) %T.TNNUM=%NP.NNUM /\* 这一年 \*/

|| %mp. 定语 . 内部结构 = 单词 , , %mp. 定语 .yx= 一 , %mp. 量词子类 = 个体

=> T(T/r M<one>) /\* 这一个 哪一个 \*/

|| %r.yx= 这 | 那 , IF %mp. 定语 . 内部结构 = 单词 , %mp. 定语 .yx= 一 FALSE

=> NP(T/r !M/mp) %T.TNNUM=PLUR,\$.NNUM=PLUR /\* 这两张 \*/

=> NP(T/r !NP/mp) %T.TNNUM=PLUR,\$.NNUM=PLUR

|| %r.yx=~ 这 ~ 那 , IF %mp. 定语 . 内部结构 = 单词 , %mp. 定语 .yx= 一 FALSE

=> NP(T/r !M/mp) \$.NNUM=%M.NNUM

=> NP(T/r !NP/mp) %T.TNSUB=%NP.NSUBC,...

# 基于规则的方法 (1)

- 采用规则作为知识表示形式
  - 重叠词规则
  - 切分规则
  - 标注规则
  - 句法分析规则
  - 语义分析规则
  - 结构转换规则（产生译文句法语义结构）
  - 词语转换规则（译词选择）
  - 结构生成规则（译文结构调整）
  - 词语生成规则（译文词形生成）

# 基于规则的方法 (2)

- 优点
  - 直观，能够直接表达语言学家的知识
  - 规则的颗粒度具有很大的可伸缩性
    - 大颗粒度的规则具有很强的概括能力
    - 小颗粒度的规则具有精细的描述能力
  - 便于处理复杂的结构和进行深层次的理解，如解决长距离依赖问题
  - 系统适应性强，不依赖于具体的训练语料



# 基于规则的方法 (3)

- 缺点

- 规则主观因素重，有时与客观事实有一定差距
- 规则的覆盖性差，特别是细颗粒度的规则很难总结得比较全面
- 规则之间的冲突没有好的解决办法（翘翘板现象）
- 规则一般只局限于某一个具体的系统，规则库开发成本太高
- 规则库的调试极其枯燥乏味

# 基于规则的方法—译词选择

\$\$ 开

**\*\*{v} v \$=[...]**

|| \$. 主体 = 是, \$. 主体 . 语义类 = 植物

→ V<bloom> \$=[...]

|| \$. 客体 = 是, \$. 客体 . 汉字 = 灯 | 机 | 器

→ V( !V<turn> D<on> ) \$=[...]

|| \$. 客体 = 是, \$. 客体 . 语义类 = 交通工具

=> V<drive> \$=[...]

|| OTHERWISE

=> V<open> \$=[...]

# 基于规则的方法—结构转换

&& {mp7} mp->r !mp :: \$. 内部结构 = 组合定中 ,...

|| %mp. 定语 . 内部结构 = 单词 , %mp. 定语 .yx= 一 , %mp. 量词子类 = 集体 | 种类 | 容量 | 时量 | 度量 | 成形

=> NP(T/r !NP/mp) %T.TNNUM=%NP.NNUM /\* 这一年 \*/

|| %mp. 定语 . 内部结构 = 单词 , , %mp. 定语 .yx= 一 , %mp. 量词子类 = 个体

=> T(T/r M<one>) /\* 这一个 哪一个 \*/

|| %r.yx= 这 | 那 , IF %mp. 定语 . 内部结构 = 单词 , %mp. 定语 .yx= 一 FALSE

=> NP(T/r !M/mp) %T.TNNUM=PLUR,\$.NNUM=PLUR /\* 这两张 \*/

=> NP(T/r !NP/mp) %T.TNNUM=PLUR,\$.NNUM=PLUR

|| %r.yx=~ 这 ~ 那 , IF %mp. 定语 . 内部结构 = 单词 , %mp. 定语 .yx= 一 FALSE

=> NP(T/r !M/mp) \$.NNUM=%M.NNUM

=> NP(T/r !NP/mp) %T.TNSUB=%NP.NSUBC,...

# 机器翻译的知识表示

程序代码：直接机器翻译方法

翻译规则：基于规则的机器翻译方法

翻译实例：基于实例的机器翻译方法

模型参数：统计机器翻译方法

# 基于实例的机器翻译方法

- 基于实例的机器翻译方法简介
- 平行语料库对齐技术
- 翻译记忆方法
- 基于短语实例的机器翻译方法
- 基于模板（模式）的机器翻译方法
- 基于依存实例的机器翻译方法

# 基于语料库的机器翻译方法

- 机器翻译的实例方法和统计方法都是基于语料库的机器翻译方法
- 优点
  - 使用语料库作为翻译知识来源，无需人工编写规则，系统开发成本低，速度快
  - 从语料库中学习到的知识比较客观
  - 从语料库中学习到的知识覆盖性比较好
- 缺点
  - 系统性能依赖于语料库
  - 数据稀疏问题严重
  - 语料库中不容易获得大颗粒度的高概括性知识

# 基于实例的机器翻译 (1)

- 长尾真 (Makoto Nagao) 在 1984 年发表了《采用类比原则进行日 - 英机器翻译的一个框架》一文，探讨日本人初学英语时翻译句子的基本过程，长尾真认为，初学英语的日本人总是记住一些最基本的英语句子以及一些相对应的日语句子，他们要对比不同的英语句子和相对应的日语句子，并由此推论出句子的结构。参照这个学习过程，在机器翻译中，如果我们给出一些英语句子的实例以及相对应的日语句子，机器翻译系统来识别和比较这些实例及其译文的相似之处和相差之处，从而挑选出正确的译文。
- 长尾真指出，人类并不通过做深层的语言学分析来进行翻译，人类的翻译过程是：首先把输入的句子正确地分解为一些短语碎片，接着把这些短语碎片翻译成其它语言的短语碎片，最后再把这些短语碎片构成完整的句子，每个短语碎片的翻译是通过类比的原则来实现的。
- 因此，我们应该在计算机中存储一些实例，并建立由给定的句子找寻类似例句的机制，这是一种由实例引导推理的机器翻译方法，也就是基于实例的机器翻译。

# 基于实例的机器翻译 (2)

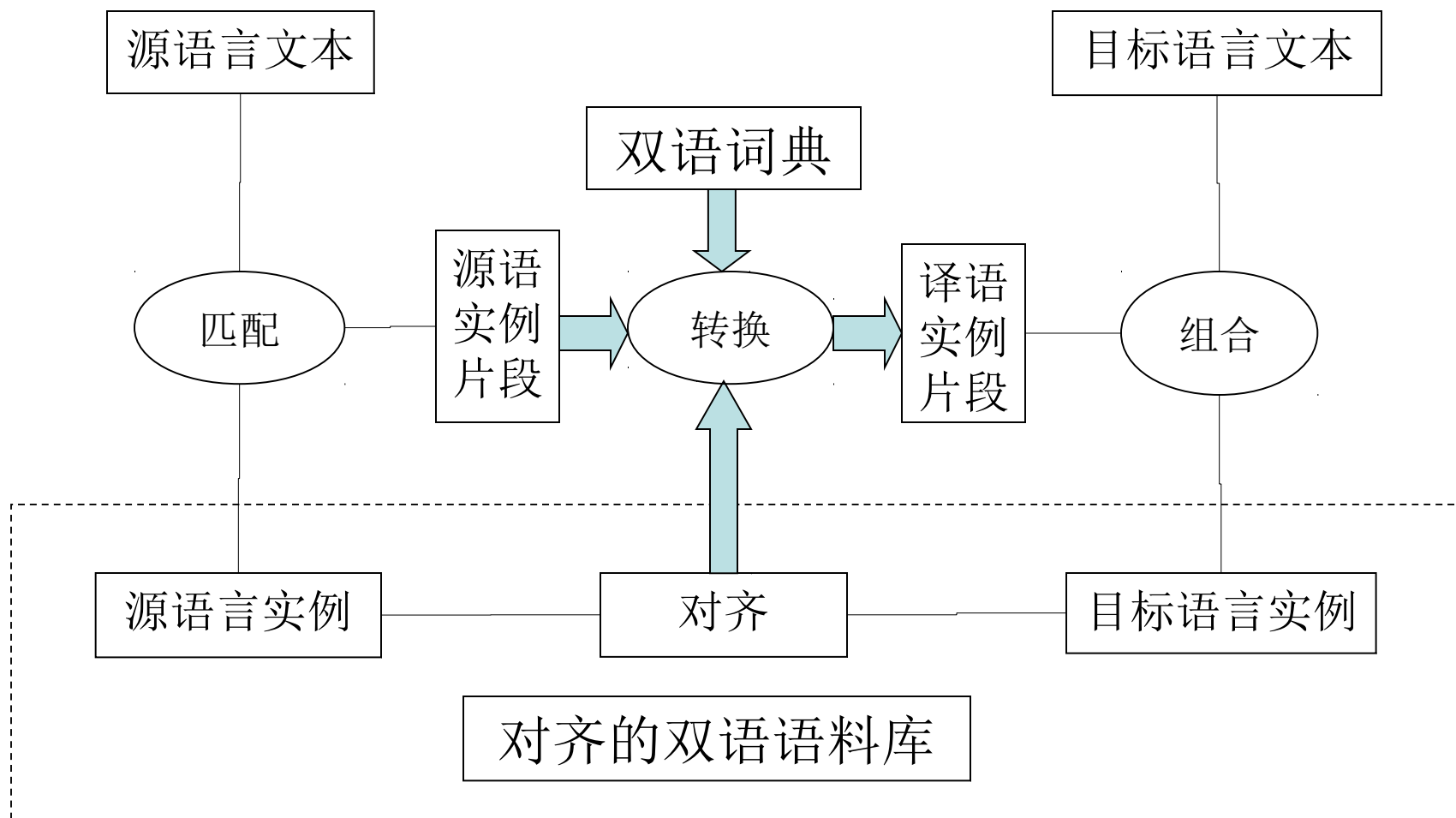
- 在基于实例的机器翻译系统中，系统的主要知识源是双语对照的翻译实例库，实例库主要有两个字段，一个字段保存源语言句子，另一个字段保存与之对应的译文，每输入一个源语言的句子时，系统把这个句子同实例库中的源语言句子字段进行比较，找出与这个句子最为相似的句子，并模拟与这个句子相对应的译文，最后输出译文。
- 基于实例的机器翻译系统中，翻译知识以实例和义类词典的形式来表示，易于增加或删除，系统的维护简单易行，如果利用了较大的翻译实例库并进行精确的对比，有可能产生高质量译文，而且避免了基于规则的那些传统的机器翻译方法必须进行深层语言学分析的难点。在翻译策略上是很有吸引力的。



# 基于实例的机器翻译 (3)

- 优点
  - 直接使用对齐的语料库作为知识表示形式，知识库的扩充非常简单
  - 不需要进行深层次的语言分析，也可以产生高质量的译文
- 缺点
  - 覆盖率低，实用的系统需要的语料库规模极大（百万句对以上）

# 基于实例的机器翻译系统结构



# 基于实例的机器翻译一举例

要翻译句子：

(E1) He bought a book on physics.

在语料库中查到相似英语句子及其汉语译文是：

(E2) He wrote a book on history.

(C2) 他写了一本关于历史的书。

比较 (E1) 和 (E2) 两个句子，我们得到变换式：

(T1) replace(wrote, bought) and replace(history, physics)

将这个变换式中的单词都换成汉语就变成：

(T2) replace( 写, 买 ) and replace( 历史, 物理 )

将 (T2) 作用于 (C2)

(C1) 他买了一本关于物理学的书。

# 基于实例的机器翻译方法

- 基于实例的机器翻译方法简介
- 平行语料库对齐技术
- 翻译记忆方法
- 基于短语实例的机器翻译方法
- 基于模板（模式）的机器翻译方法
- 基于依存实例的机器翻译方法

# 平行语料库齐技术

- 双语语料库（ **Bilingual Corpus** ）或平行语料库（ **Parallel Corpus** ），在 **EBMT** 中又称为实例库
- 双语语料库对齐的级别
  - 篇章对齐
  - 段落对齐
  - 句子对齐
  - 词语对齐
  - 短语块对齐
  - 句法结构对齐
- 基于实例的机器翻译中实例库必须至少做到句子级别的对齐

# 不同对齐级别的差异

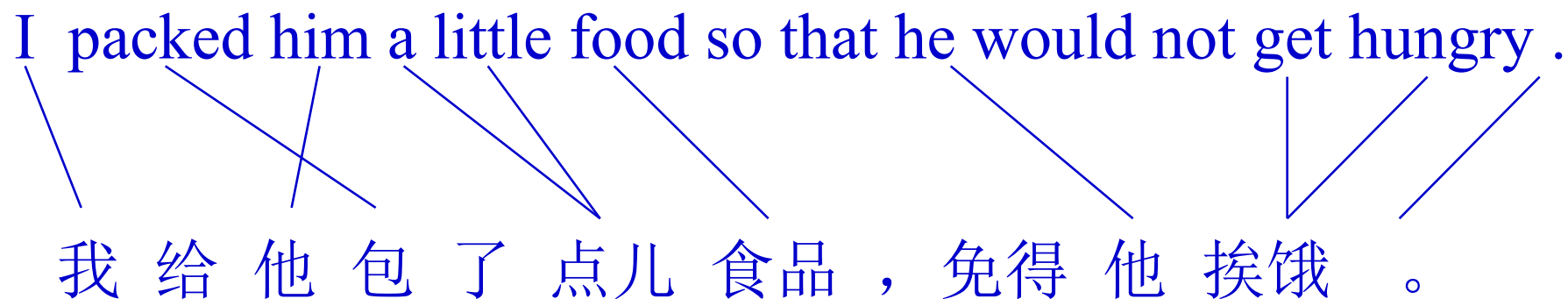
- 段落对齐和句子对齐
  - 要求保持顺序（允许局部顺序的调整）
  - 只有一个层次
- 词语对齐和短语块对齐
  - 不要求保持顺序
  - 只有一个层次
- 句法结构对齐
  - 不要求保持顺序
  - 多层次对齐

# 句子对齐

汉语	英语	模式
1995 年初我来成都的那天，没想到会是在一个冬季的漆黑的日子。	I little thought when I arrived in Chengdu in the dark, dark days of winter, early in 1995, that I would still be here more than five years later.	1:1
那时我也根本没有想到会在这儿呆上五年，也不知道我会遇到一位成都的女儿，并且后来还娶她为妻。  一个完全陌生的家庭接纳了我，我也因此成为成都的一部分。	I little knew that I would meet one of Chengdu's daughters, and later marry her, thus acquiring a whole new family who embraced me as one of them, and thus I became part of this place.	2:1

# 词语对齐 (1)

I packed him a little food so that he would not get hungry .  
我 给 他 包 了 点 儿 食 品 ， 免 得 他 挨 饿 。



- 特点：
  - 保序性不再满足
  - 对齐模式复杂：一对多、多对一、多对多都非常普遍



# 词语对齐 (2)

- 困难：
  - 翻译歧义：一个词出现两个以上的译词
  - 双语词典覆盖率有限：非常普遍的现象
  - 位置歧义：出现两个以上相同的词
  - 汉语词语切分问题
  - 虚词问题：虚词的翻译非常灵活，或没有对译词
  - 意译问题：根本找不到对译的词

## 词语对齐 (3)

- 一般而言，一个单词对齐的模型可以表述为两个模型的乘积：
  - 词语相似度模型 (word similarity model)
  - 位置扭曲模型 (word distortion model)用公式表示如下：

$$Score(e_i, c_j) = S(e_i, c_j) \times D(i, j)$$

# 基于实例的机器翻译方法

- 基于实例的机器翻译方法简介
- 平行语料库对齐技术
- 翻译记忆方法
- 基于短语实例的机器翻译方法
- 基于模板（模式）的机器翻译方法
- 基于依存实例的机器翻译方法

# 翻译记忆方法 (1)

- 翻译记忆方法（ Translation Memory ）是基于实例方法的特例；
- 也可以把基于实例的方法理解为广义的翻译记忆方法；
- 翻译记忆的基本思想：
  - 把已经翻译过的句子保存起来
  - 翻译新句子时，直接到语料库中去查找
    - 如果发现相同的句子，直接输出译文
    - 否则交给人去翻译，但可以提供相似的句子的参考译文

## 翻译记忆方法 (2)

- 翻译记忆方法主要被应用于计算机辅助翻译（**CAT**）软件中
- 翻译记忆方法的优缺点
  - 翻译质量有保证
  - 随着使用时间的增加匹配成功率逐步提高
  - 特别适用于重复率高的文本翻译，例如公司的产品说明书的新版本翻译
  - 与语言无关，适用于各种语言对
  - 缺点是匹配成功率不高，特别是刚开始使用时

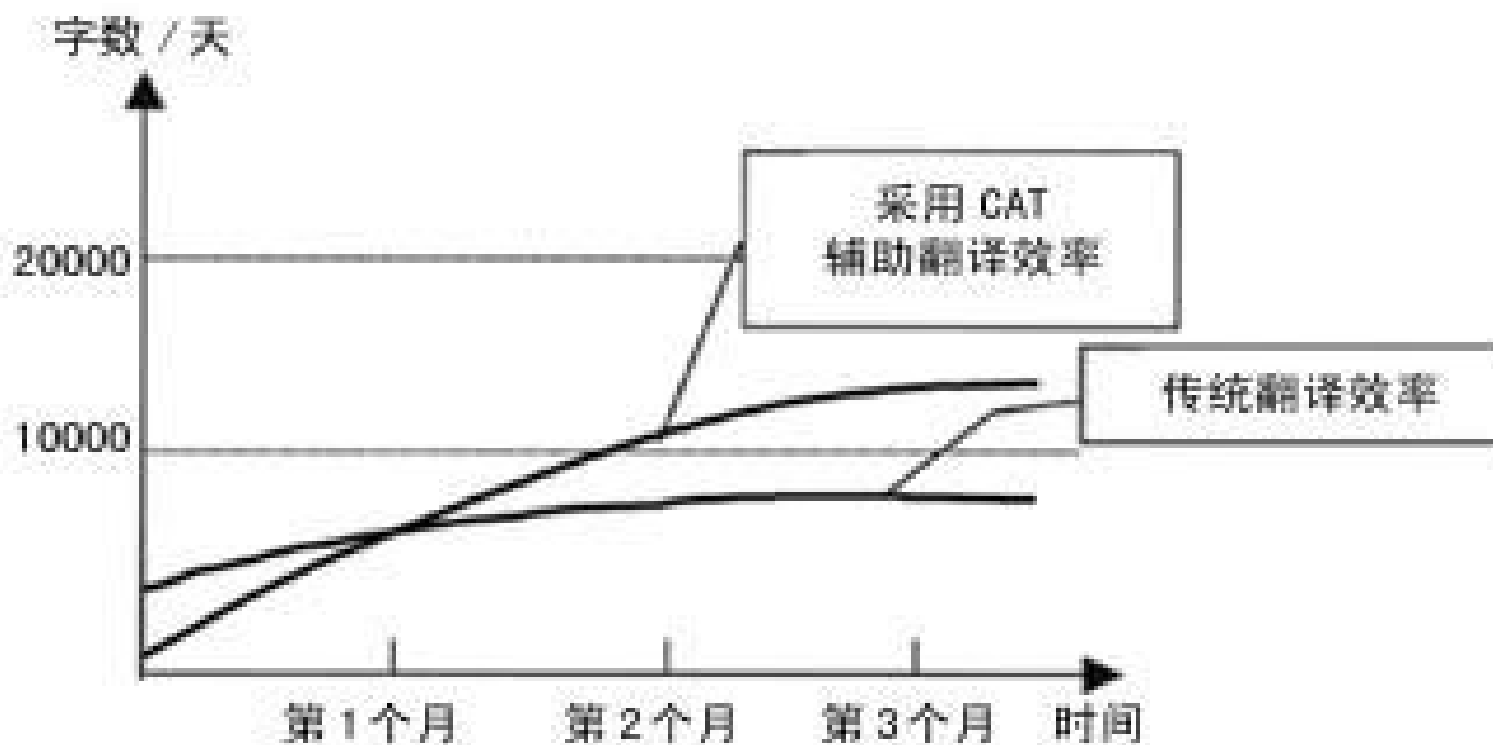
# 翻译记忆方法 (3)

- 计算机辅助翻译（**CAT**）软件已经形成了比较成熟的产业
  - **TRADOS**
    - 号称占有国际 **CAT** 市场的 70%
    - **Microsoft**、**Siemens**、**SAP** 等国际大公司和一些著名的国际组织都是其用户
  - 雅信 **CAT**
    - 适合中国人的习惯
    - 产品已比较成熟
  - 国际组织：**LISA**（**Localisation Industry Standards Association**）
- 面向用户：专业翻译人员
- 数据交换：**LISA** 制定了 **TMX**（**Translation Memory eXchange**）标准。

# 翻译记忆方法 (4)

- 完整的计算机辅助翻译软件除了包括翻译记忆功能以外，还应该包括以下功能
  - 多种文件格式的分解与合成
  - 术语库管理功能
  - 语料库的句子对齐（历史资料的重复利用）
  - 项目管理：
    - 翻译任务的分解与合并
    - 翻译工作量的估计
  - 数据共享和数据交换

# 翻译记忆方法 (5)





# 基于实例的机器翻译方法

- 基于实例的机器翻译方法简介
- 平行语料库对齐技术
- 翻译记忆方法
- **基于短语实例的机器翻译方法**
- 基于模板（模式）的机器翻译方法
- 基于依存实例的机器翻译方法

# 基于短语实例的机器翻译方法

- 把整个句子切分成若干短语或单词，使得每个短语都可以在实例库中找到相似度很高的短语，而每个单词都可以在词典中找到译文
- 短语的翻译需要利用实例的词语对齐信息，查找相应的译文片段
- 问题
  - 源语言短语的划分
  - 源语言短语和目标语言短语的对齐
  - 原语言片段的匹配（相似度）和翻译
  - 目标语言短语的组合

# 基于实例的机器翻译方法

- 基于实例的机器翻译方法简介
- 翻译记忆方法
- 基于短语实例的机器翻译方法
- **基于模板（模式）的机器翻译方法**
- 基于依存实例的机器翻译方法

# 基于模板（模式）的机器翻译方法 (1)

- 基于模板（**Template**）或者模式（**Pattern**）的机器翻译方法通常也被看做基于实例的机器翻译方法的一种延伸
- 所谓“翻译模板”或者“翻译模式”可以认为是一种颗粒度介于“翻译规则”和“翻译实例”之间的翻译知识表示形式
  - 翻译规则：颗粒度大，匹配可能性大，但过于抽象，容易出错
  - 翻译实例：颗粒度小，不易出错，但过于具体，匹配可能性小
  - 翻译模板（模式）：介于二者之间，是一种比较合适的知识表示形式
- 一般而言，单语模板（或模式）是一个常量和变量组成的字符串，翻译模板（或模式）是两个对应的单语模板（或模式），两个模板之间的变量存在意义对应关系

# 基于模板（模式）的机器翻译方法 (2)

- 模板举例：
  - 这个 X 比 Y 更 Z。
  - The X is more Z than Y.
- 模板方法的主要问题
  - 对模板中变量的约束
  - 模板抽取
  - 模板的冲突消解

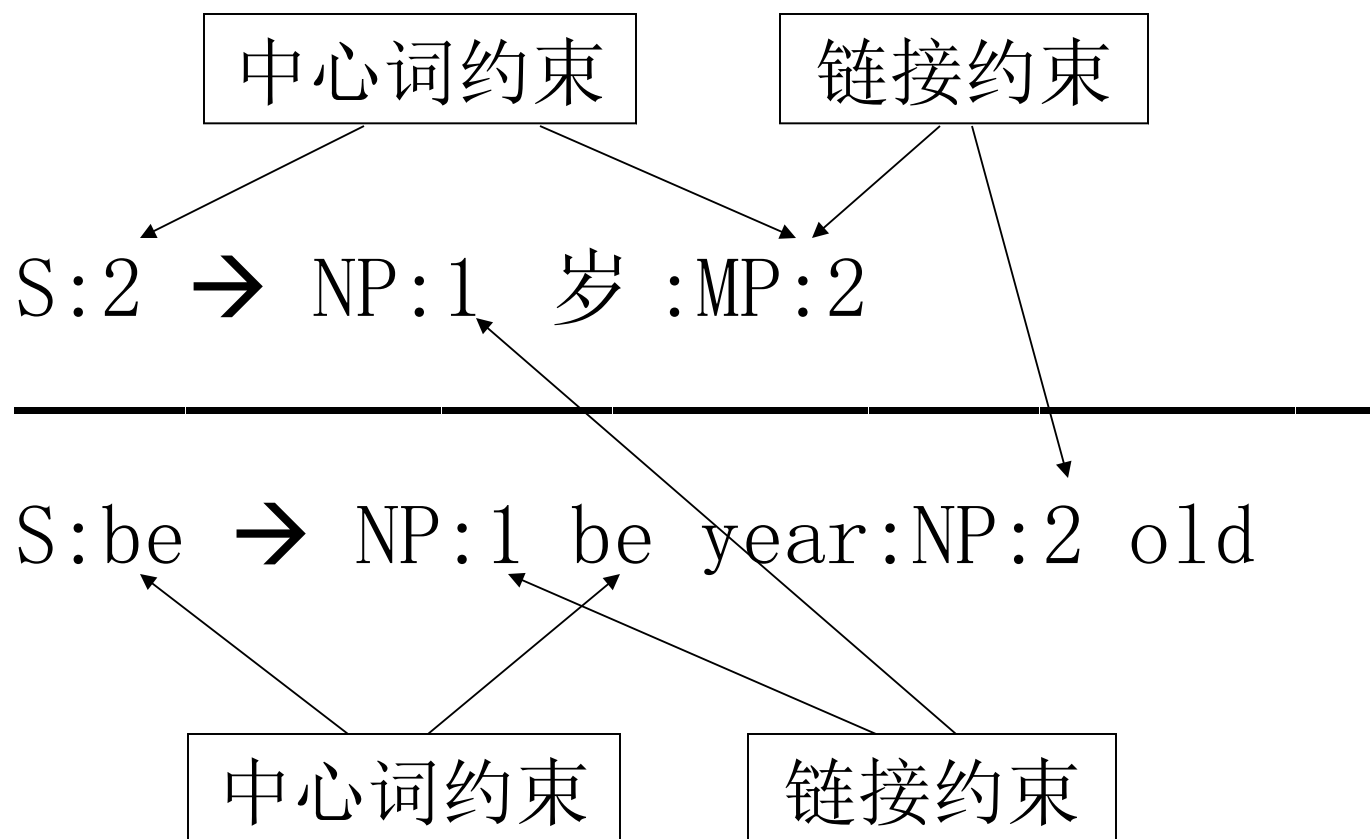
# Pattern-Based CFG for MT(1)

- Koichi Takeda, Pattern-Based Context-Free Grammars for Machine Translation, Proc. of 34th ACL, pp. 144-- 151, June 1996
- 给出了翻译模式的一种形式化定义，并给出了相应的翻译算法以及算法复杂性的理论证明

# Pattern-Based CFG for MT(2)

- 每个翻译模板由一个源语言上下文无关规则和一个目标语言上下文无关规则（这两个规则称为翻译模板的骨架），以及对这两个规则的中心词约束和链接约束构成；
- 中心词约束：对于上下文无关语法规则中右部（子结点）的每个非终结符，可以指定其中心词；对于规则左部（父结点）的非终结符，可以直接指定其中心词，也可以通过使用相同的序号规定其中心词等于其右部的某个非终结符的中心词；
- 链接约束：源语言骨架和目标语言骨架的非终结符子结点通过使用相同的序号建立对应关系，具有对应关系的非终结符互为翻译。

# Pattern-Based CFG for MT(3)





# Pattern-Based CFG for MT(4)

- 翻译的过程分为三步：
  - 使用源语言 **CFG** 骨架分析输入句子 **s**
  - 应用源语言到目标语言的 **CFG** 骨架的链接约束，生成一个译文 **CFG** 推导序列
  - 根据译文 **CFG** 推导序列产生译文
- 模板排序的启发式原则：
  - 对于源文 **CFG** 骨架相同的模板，有中心词约束的模板优先于没有中心词约束的模板；
  - 对于同一跨度上的两个结点，比较其对应的模板的源文 **CFG** 骨架，非终结符少的模板优先于非终结符多的模板；
  - 中心词约束被满足的结点优先于中心词约束不被满足的结点；
  - 对于一个输入串而言，分析步骤越短（推导序列越短）越优先。

# Pattern-Based CFG for MT(5)

- 模板库的获取：假设  $T$  是一组翻译模板， $B$  是双语语料库， $\langle s, t \rangle$  是一对互为翻译的句子
  - 如果  $T$  能够翻译句子  $s$  为  $t$ ，那么 do nothing；
  - 如果  $T$  将  $s$  译为  $t'$ （不等于  $t$ ），那么：
    - 如果  $T$  中存在  $\langle s, t \rangle$  的推导  $Q$ ，但这个推导不是最优解，那么给  $Q$  中的模板进行实例化；
    - 如果不存在这种推导，那么加入适当的模板，使得推导成立；
  - 如果根本无法翻译  $s$ （分析失败），那么将  $\langle s, t \rangle$  直接加入到模板库中。

# 模板的自动提取

- 利用一对实例进行泛化
  - Jaime G. Carbonell, Ralf D. Brown,  
Generalized Example-Based Machine Translation  
<http://www.lti.cs.cmu.edu/Research/GEBMT/>
- 利用两对实例进行比较
  - H. Altay Guvenir, Ilyas Cicekli, Learning Translation Templates from Examples  
Information Systems, 1998
  - 张健, 基于实例的机器翻译的泛化方法研究, 中科院计算所硕士论文, 2001

# 通过泛化实例得到翻译模板

- 已有实例：
  - Karl Marx was born in Trier, Germany in May 5, 1818.
  - 卡尔·马克思于 1818 年 5 月 5 日出生在德国特里尔城。
- 泛化：
  - <Person> was born in <City> in <Date>
  - <Person> 于 <Date> 出生在 <City>
- 对齐
  - <Person>  $\Leftrightarrow$  <Person>
  - <City>  $\Leftrightarrow$  <City>
  - <Date>  $\Leftrightarrow$  <City>

# 通过比较实例得到翻译模板

- 已有两对翻译实例：
  - 我给玛丽一支笔  $\leftrightarrow$  I gave Mary a pen.
  - 我给汤姆一本书  $\leftrightarrow$  I gave Tom a book.
- 双侧单语句子分别比较，得到：
  - 我 给 #X 一 #Y #Z  $\leftrightarrow$  I give #W a #U.
- 查找变量的对应关系：
  - #X  $\leftrightarrow$  #W
  - #Y  $\leftrightarrow \phi$
  - #Z  $\leftrightarrow$  #U

# 基于实例的机器翻译方法

- 基于实例的机器翻译方法简介
- 翻译记忆方法
- 基于短语实例的机器翻译方法
- 基于模板（模式）的机器翻译方法
- 基于依存实例的机器翻译方法

# 基于依存实例的机器翻译

- [Sato & Nagao 1990] 的方法:
  - 将实例按照**词语依存树**配对的形式进行存储, 同时保存结点对应关系链接的集合

He eats vegetables.

e([e1,[eat,v],  
[e2,[he,[pron]],  
[e3,[vegetable,n]]])).

**clinks**([[e1, j1], [e2, j3], [e3, j5]]).

Kare ha yasai wo taberu.

e([j1,[taberu,v],  
[j2,[ha,p],  
[j3,[kare,pron]]],  
[j4,[wo,p],  
[j5,[yasai,n]]])).

# 基于依存实例的机器翻译

- [Sato & Nagao 1990] 的方法：
  - 在翻译的过程中，每一个输入句子都被表示为一个或多个匹配表达式。
  - 每一个匹配表达式表示在实例库中找到的某个依存子树的特定结点上所进行的某种操作（即插入、删除和替换）。
  - 利用这些操作，可以通过数据库中找到的实例片段来组合得到输入的句子。

输入英语句子：“He eats mashed potatoes.”

匹配表达式为：[ e1, [ r, e3, [ e<sup>x</sup> ] ] ]

这里 r 表示替换，整个表达式的意思是

“在实例 e1 中，用结点 e<sup>x</sup> 替换结点 e3”



# 机器翻译的知识表示

程序代码：直接机器翻译方法

翻译规则：基于规则的机器翻译方法

翻译实例：基于实例的机器翻译方法

模型参数：统计机器翻译方法

# 统计机器翻译

- 统计机器翻译：一种新的研究范式
- 经典的统计机器翻译方法  
—基于词的 **IBM** 模型
- 最成熟的统计机器翻译方法  
—基于短语的模型
- 目前统计机器翻译研究的热点  
—基于句法的模型

# 统计机器翻译：一种新的研究范式

- 统计机器翻译的成功在于采用了一种新的研究范式（ **paradigm** ）
- 这种研究范式已在语音识别等领域中被证明是一种成功的翻译，但在机器翻译中是首次使用
- 这种范式的特点：
  - 公开的大规模的训练数据
  - 周期性的公开评测和研讨
  - 开放源码的工具

# 统计机器翻译

- 统计机器翻译：一种新的研究范式
- 经典的统计机器翻译方法  
—基于词的 IBM 模型
- 最成熟的统计机器翻译方法  
—基于短语的模型
- 目前统计机器翻译研究的热点  
—基于句法的模型

# 基于词的统计机器翻译方法

- 统计机器翻译—为翻译建立概率模型
- IBM 的信源信道模型
- 语言模型— $n$  元语法模型
- 翻译模型—IBM 模型 1-5
- Candide 系统

# 为翻译建立概率模型

- 假设任意一个英语句子  $e$  和一个法语句子  $f$  , 我们定义  $f$  翻译成  $e$  的概率为:

$$\Pr(e | f)$$

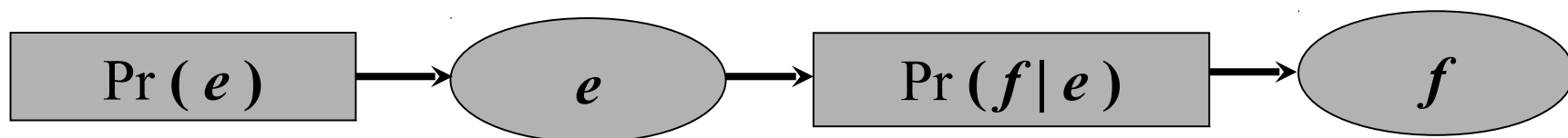
其归一化条件为:

$$\sum_e \Pr(e | f) = 1$$

- 于是将  $f$  翻译成  $e$  的问题就变成求解问题:

$$\hat{e} = \operatorname{argmax}_e \Pr(e | f)$$

# 信源信道模型 (1)



- 假设我们看到的源语言文本  $f$  是由一段目标语言文本  $e$  经过某种奇怪的编码得到的，那么翻译的目标就是要将  $f$  还原成  $e$ ，这也就是就是一个解码的过程。
- 注意，在信源信道模型中：
  - 噪声信道的源语言是翻译的目标语言
  - 噪声信道的目标语言是翻译的源语言

这与整个机器翻译系统翻译方向的刚好相反

# 信源信道模型 (2)

$$\hat{e} = \arg \max_e \Pr(e) \Pr(f | e)$$

- P.Brown 称上式为统计机器翻译基本方程式
  - 语言模型:  $P(E)$
  - 翻译模型:  $P(F|E)$
- 语言模型反映 “ E 像一个句子” 的程度: 流利度
- 翻译模型反映 “ F 像 E” 的程度: 忠实度
- 联合使用两个模型效果好于单独使用翻译模型, 因为后者容易导致一些不好的译文。



# 信源信道模型 (3)

- 统计机器翻译分解为以下三个问题：
  - 语言模型的定义和参数估计
  - 翻译模型的定义和参数估计
  - 解码

# 语言模型 — n 元语法模型

- 语言模型在机器翻译中具有极为重要的作用
- 到目前位置，统计机器翻译中最常用、而且最有效的模型仍然是 n 元语法模型
- 模型的阶数越来越高：3 元、4 元、5 元
- 模型的训练语料越来越大：
  - Google 提供了公开的 Web 1T 语料库，其中的 n 元共现词频数据是从 web 中得到的 1T 英文词的语料库中统计得到的（剪切掉了低频组合）
  - Google 号称使用了 2T 英文词训练的语言模型
  - 大规模的数据为系统实现带来很大的困难

# 翻译模型

- 翻译模型  $P(F|E)$  反映的是一个源语言句子  $E$  翻译成一个目标语言句子  $F$  的概率
- 由于源语言句子和目标语言句子几乎不可能在语料库中出现过，因此这个概率无法直接从语料库统计得到，必须分解成词语翻译的概率和句子结构（或者顺序）翻译的概率

# 翻译模型与对齐

- 翻译模型的计算，需要引入隐含变量：  
对齐  $A$ ：

$$P(F | E) = \sum_A P(F, A | E)$$

- 翻译概率  $P(F|E)$  的计算转化为对齐概率  $P(F,A|E)$  的估计
- 对齐：建立源语言句子和目标语言句子的词与词之间的对应关系和句子结构之间的对应关系

# 词语对齐的表示 (1)

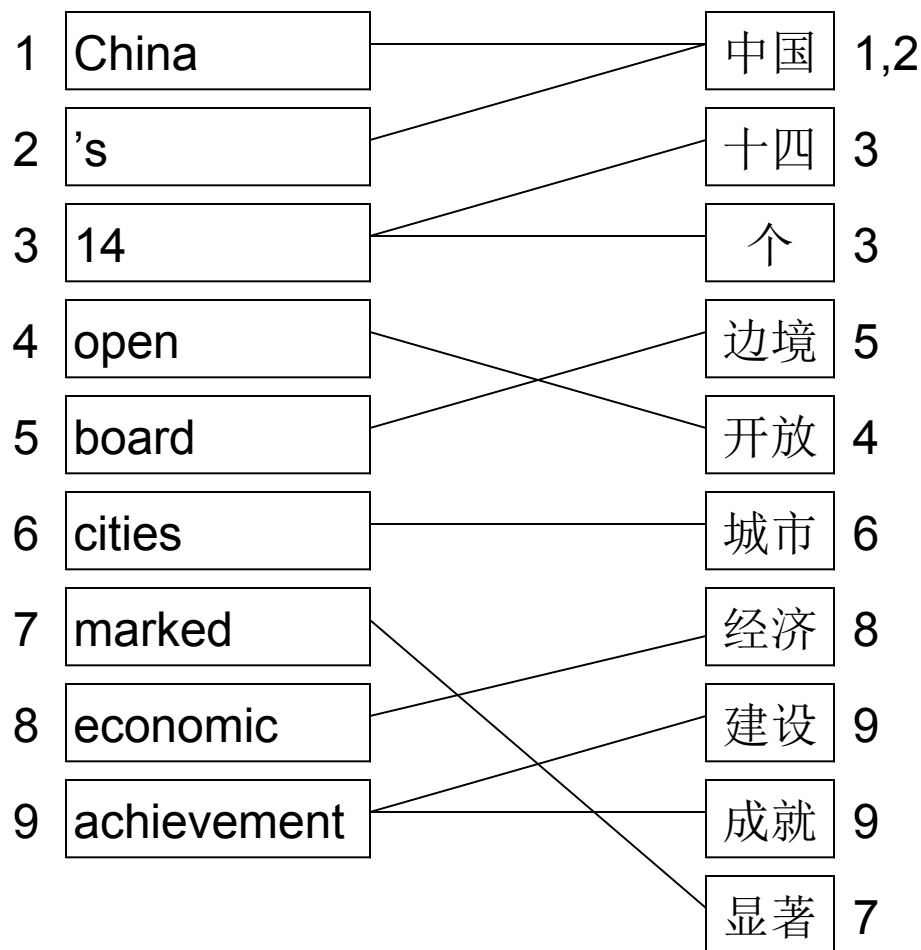
- 图形表示

- ✓ 连线

- ✓ 矩阵（见下页）

- 数字表示

- ✓ 给每个目标语言单词标记其所有对应的源语言单词



# 词语对齐的表示 (2)

achievement										
economic										
marked										
cities										
board										
open										
14										
's										
China										
	中国	十四	个	边境	开放	城市	经济	建设	成就	显著

# IBM Model 1

- 最简单的理解，可以句子 **e** 翻译成 **f** 的概率，就是 **e** 中每一个词语翻译成 **f** 中对应词语的概率的乘积
- 这就是 IBM Model 1 的基本思想
- IBM Model 1 可以利用句子对齐的双语语料库，通过 EM 无指导学习算法获得
- IBM 提出了复杂度递增的 5 个统计翻译模型， IBM Model 1 是最简单的模型

# IBM Model 1-5

- IBM Model 1 仅考虑词对词的互译概率
- IBM Model 2 加入了词的位置变化的概率
- IBM Model 3 加入了一个词翻译成多个词的概率
- IBM Model 4
- IBM Model 5



# IBM 公司的 Candide 系统 (1)

- 基于统计的机器翻译方法
- 分析—转换—生成
  - 中间表示是线性的
  - 分析和生成都是可逆的
- 分析（预处理）：
  1. 短语切分
  2. 专名与数词检测
  3. 大小写与拼写校正
  4. 形态分析
  5. 语言的归一化

# IBM 公司的 Candide 系统 (2)

- 转换（解码）：基于统计的机器翻译
- 解码分为两个阶段：
  - 第一阶段：使用粗糙模型的堆栈搜索
    - 输出 140 个评分最高的译文
    - 语言模型：三元语法
    - 翻译模型：EM Trained IBM Model 5
  - 第二阶段：使用精细模型的扰动搜索
    - 对第一阶段的输出结果先扩充，再重新评分
    - 语言模型：链语法
    - 翻译模型：最大熵翻译模型（选择译文词）

# IBM 公司的 Candide 系统 (3)

- ARPA 的测试结果 :

	Fluency		Adequacy		Time Ratio	
	1992	1993	1992	1993	1992	1993
Systran	.466	.540	.686	.743		
Candide	.511	.580	.575	.670		
Transman	.819	.838	.837	.850	.688	.625
Manual		.833		.840		

# 统计机器翻译

- 统计机器翻译：一种新的研究范式
- 经典的统计机器翻译方法  
—基于词的 IBM 模型
- 最成熟的统计机器翻译方法  
—基于短语的模型
- 目前统计机器翻译研究的热点  
—基于句法的模型

# 基于短语的统计机器翻译方法

- 从信源信道模型到对数线性模型
- 翻译模型的发展—基于短语的模型
- 短语的自动抽取
- 短语语序的调整

# 统计机器翻译的对数线性模型 (1)

- Och 于 ACL2002 提出，思想来源于 Papineni 提出的基于特征的自然语言理解方法，该论文获得 ACL2002 的最佳论文称号
- 是一个比信源—信道模型更具一般性的模型，信源—信道模型是其一个特例
- 原始论文的提法是“最大熵”模型，现在通常使用“对数线性（Log-Linear）模型”这个概念。“对数线性模型”的含义比“最大熵模型”更宽泛，而且现在这个模型通常都不再使用最大熵的方法进行参数训练，因此“对数线性”模型的提法更为准确。
- 与 NLP 中通常使用的最大熵方法的区别：使用连续量（实数）作为特征，而不是使用离散的布尔量（只取 0 和 1 值）作为特征

# 统计机器翻译的对数线性模型 (2)

假设  $e$ 、 $f$  是机器翻译的目标语言和源语言句子,  $h_1(e, f)$ ,  $\dots, h_M(e, f)$  分别是  $e$ 、 $f$  上的  $M$  个特征,  $\lambda_1, \dots, \lambda_M$  是与这些特征分别对应的  $M$  个参数, 那么翻译概率可以用以下公式模拟:

$$\begin{aligned}\Pr(e | f) &\approx p_{\lambda_1 \dots \lambda_M}(e | f) \\ &= \frac{\exp[\sum_{m=1}^M \lambda_m h_m(e, f)]}{\sum_{e'} \exp[\sum_{m=1}^M \lambda_m h_m(e', f)]}\end{aligned}$$

# 统计机器翻译的对数线性模型 (3)

对于给定的  $f$ ，其最佳译文  $e$  可以用以下公式表示：

$$\begin{aligned}\hat{e} &= \arg \max_e \{\Pr(e | f)\} \\ &\approx \arg \max_e \left\{ \sum_{m=1}^M \lambda_m h_m(e, f) \right\}\end{aligned}$$



# 对数线性模型 vs. 噪声信道模型

- 取以下特征和参数时，对数线性模型等价于噪声信道模型：
  - 仅使用两个特征
  - $h_1(e, f) = \log p(e)$
  - $h_2(e, f) = \log p(f|e)$
  - $\lambda_1 = \lambda_2 = 1$

# 对数线性模型的优点

- 噪声模型只有在理想的情况下才能达到最优，对于简化的语言模型和翻译模型，取不同的参数值实际效果更好；
- 对数线性模型大大扩充了统计机器翻译的思路；
- 特征的选择更加灵活，可以引入任何可能有用的特征。

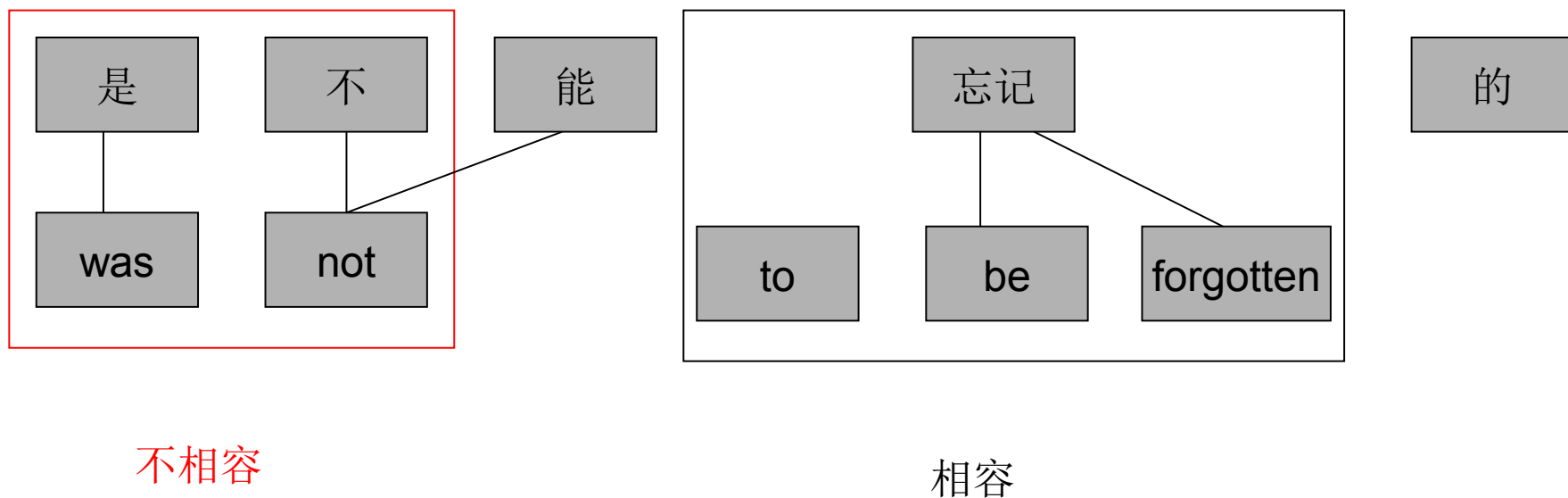
# 翻译模型的发展—基于短语的模型

- 基于词的 **IBM** 翻译模型有明显的缺陷：一个词在翻译的时候基本上不考虑上下文，孤立地进行翻译，导致了大量的错误；词序调整模型近乎无礼，很难准确调整词序，对词序差别较大的语言之间的翻译效果太差。
- 人们很容易想到，将一个短语捆绑起来进行翻译，可以大大提高翻译的准确率
- 很多不同的研究人员尝试了各种各样的基于短语的翻译模型，最终形成了目前比较成熟的基于短语的翻译模型

# 基于短语的翻译模型

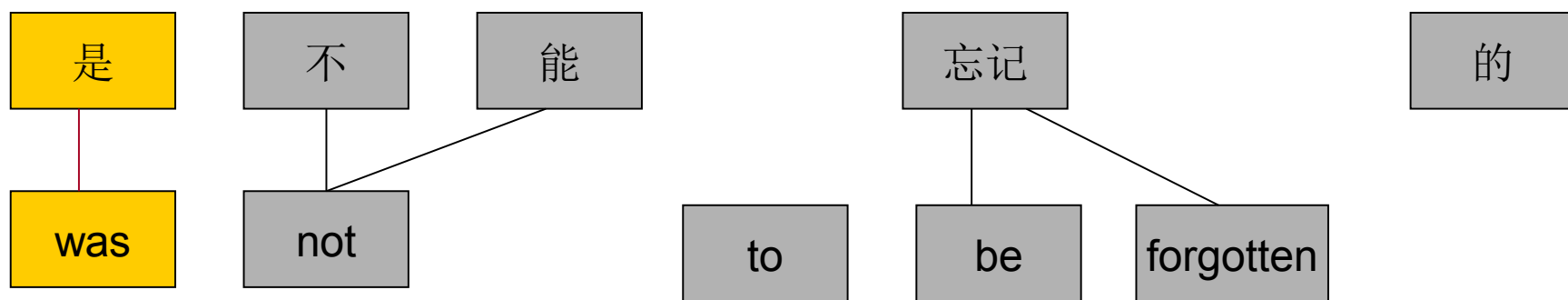
- 基本思想
  - 把训练语料库中所有对齐的短语及其翻译概率存储起来，作为一部带概率的短语词典
  - 这里所说的短语是任意连续的词串，不一定是一个独立的语言单位
  - 翻译的时候将输入的句子与短语词典进行匹配，选择最好的短语划分，将得到的短语译文重新排序，得到最优的译文
- 问题：
  - 短语如何抽取？
  - 短语概率如何计算？

# 基于词语对齐的短语自动抽取



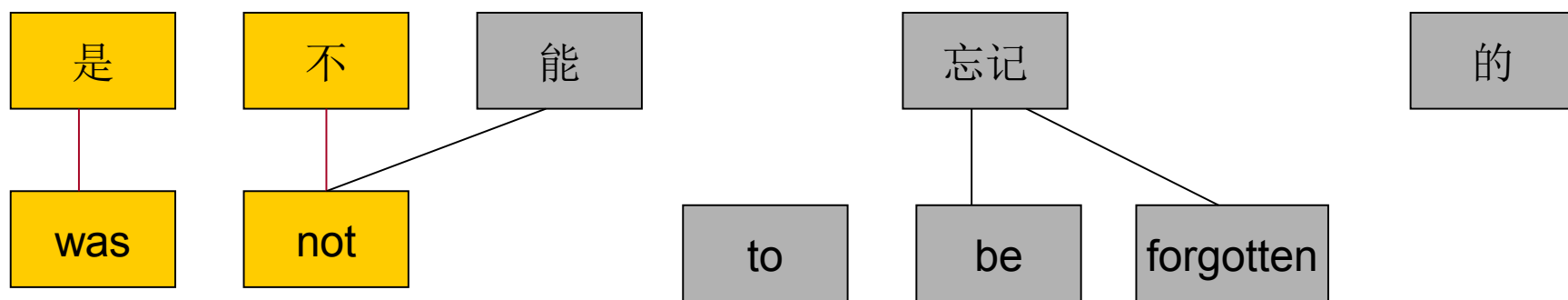
# 短语自动抽取算法运行示例 (1)

- 列举源语言所有可能的短语，  
根据对齐检查相容性



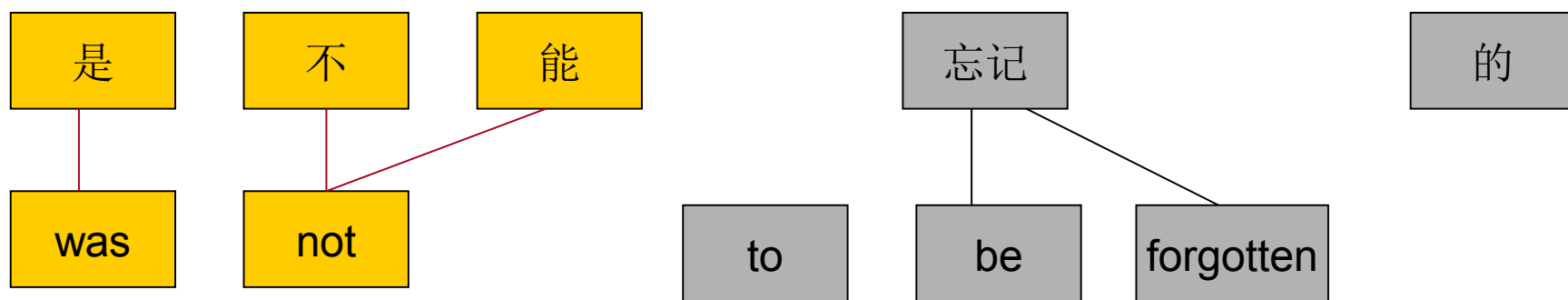
( 是, was )

# 短语自动抽取算法运行示例 (2)



不相容

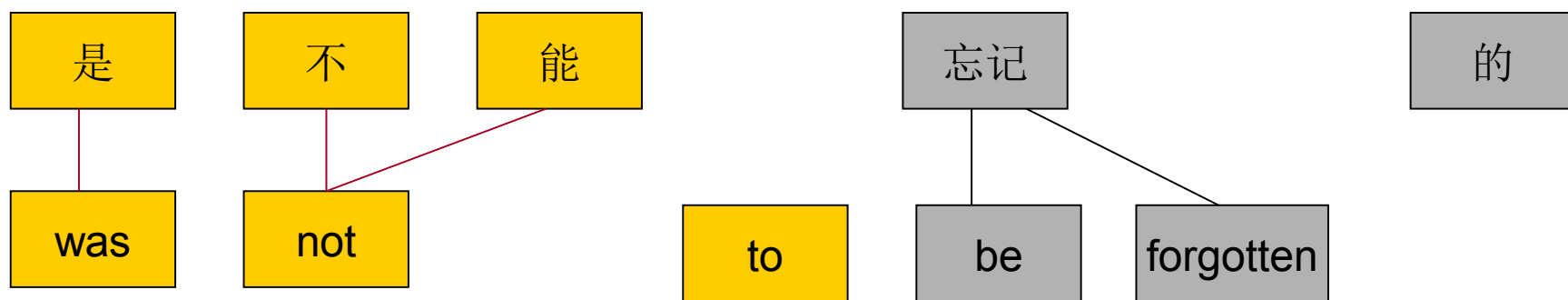
# 短语自动抽取算法运行示例 (3)



(是不能, was not)

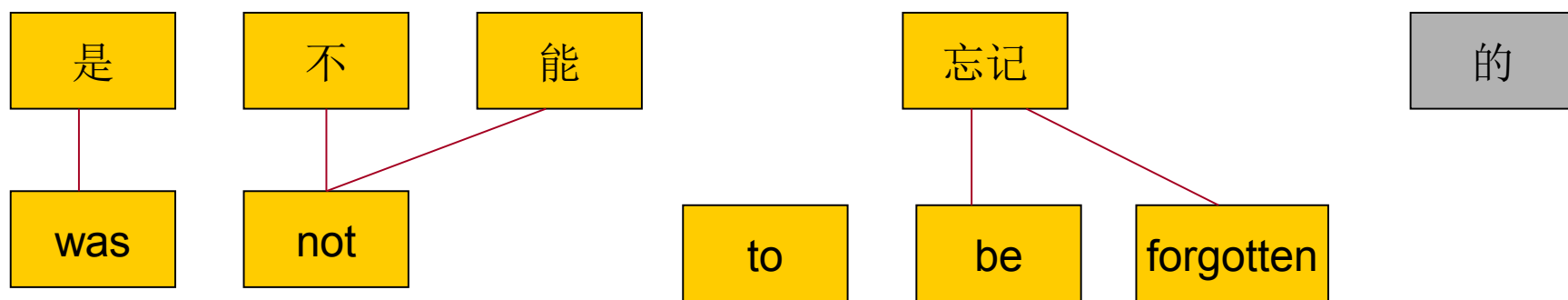


# 短语自动抽取算法运行示例 (4)



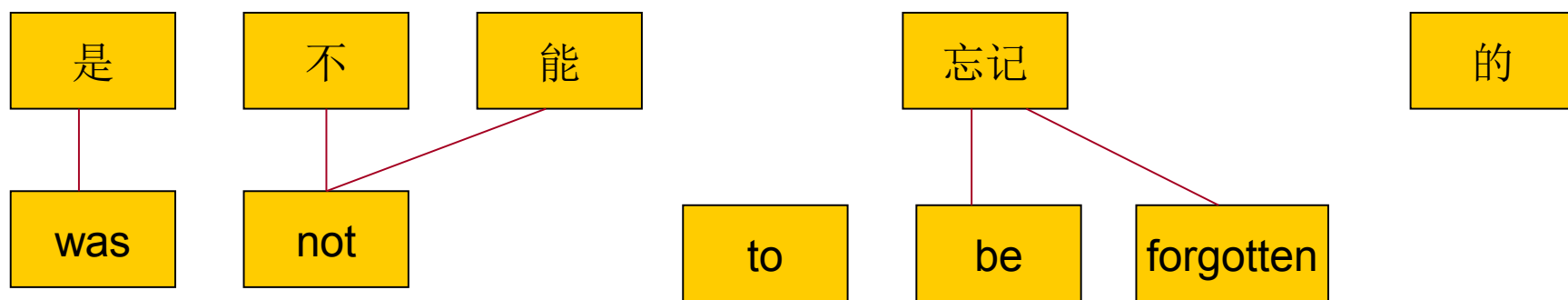
(是不能, was not to)

# 短语自动抽取算法运行示例 (5)



(是不能忘记, was not to be forgotten)

# 短语自动抽取算法运行示例 (6)



(是不能忘记的, was not to be forgotten)

# 短语表

• 是	was
• 是不能	was not
• 是不能	was not to
• 是不能忘记	was not to be forgotten
• 是不能忘记的	was not to be forgotten
• 不 能	not
• 不 能	not to
• 不 能 忘记	not to be forgotten
• 不 能 忘记 的	not to be forgotten
• 忘记	be forgotten
• 忘记	to be forgotten
• 忘记 的	be forgotten
• 忘记 的	to be forgotten

# 短语语序的调整

- 在基于短语的模型中，短语内部的顺序无需调整，只需要调整短语之间的顺序
- 短语的调序模型类似于基于词的模型，允许任意的语序调整
- 为了避免搜索空间的过于膨胀，通常限制语序调整的距离

# 统计机器翻译

- 统计机器翻译：一种新的研究范式
- 经典的统计机器翻译方法  
—基于词的 **IBM** 模型
- 最成熟的统计机器翻译方法  
—基于短语的模型
- 目前统计机器翻译研究的热点  
—基于句法的模型

# 基于句法的模型

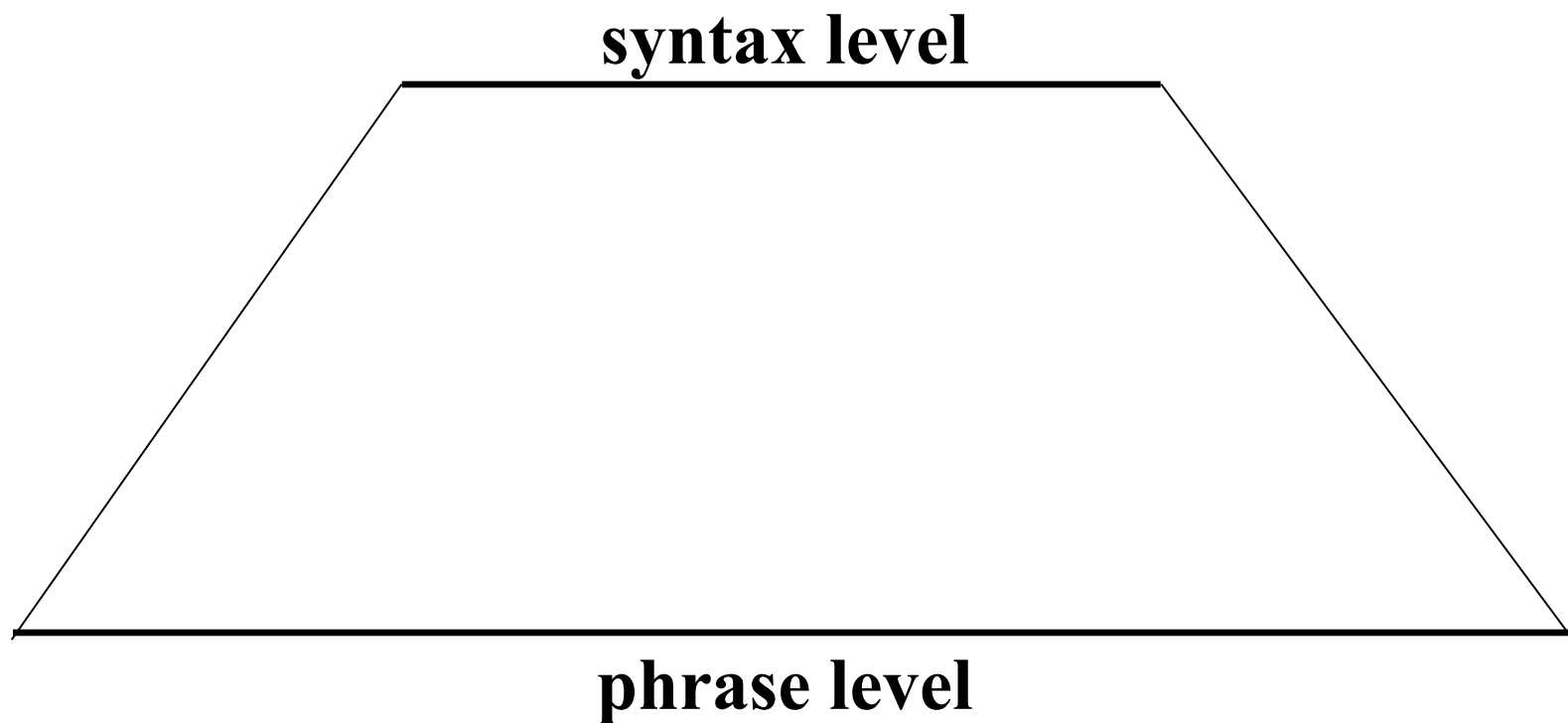
- 翻译模型的发展—基于句法的模型
- 基于句法的模型概述
- 形式上基于句法的模型
  - ITG 和 BTG
  - 最大熵 BTG 模型
  - 层次短语模型
- 语言学上基于句法的模型
  - 树到串模型
  - 串到树模型

# 翻译模型的发展—基于句法的模型

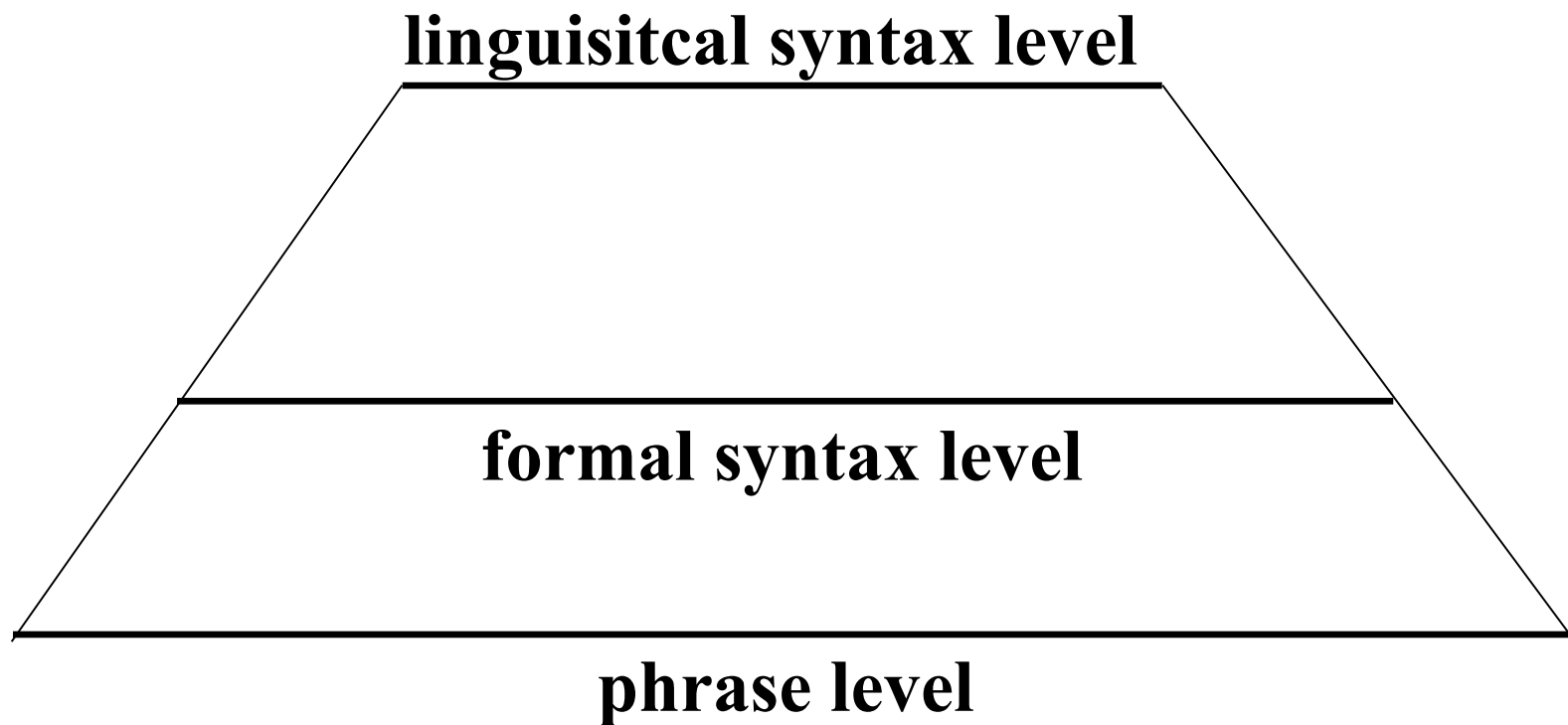
- 基于短语的模型比基于词的模型性能有了较大提高，但对于短语之间的语序调整，仍然没有提供合理的解决方案
- 经验表明，在基于短语的统计机器翻译系统中，绝大多数匹配的短语长度都是 **2-3** 个词，**1** 个词的短语也占相当大的比例
- 要解决长距离语序的调整，引入句法信息是个必然的选择



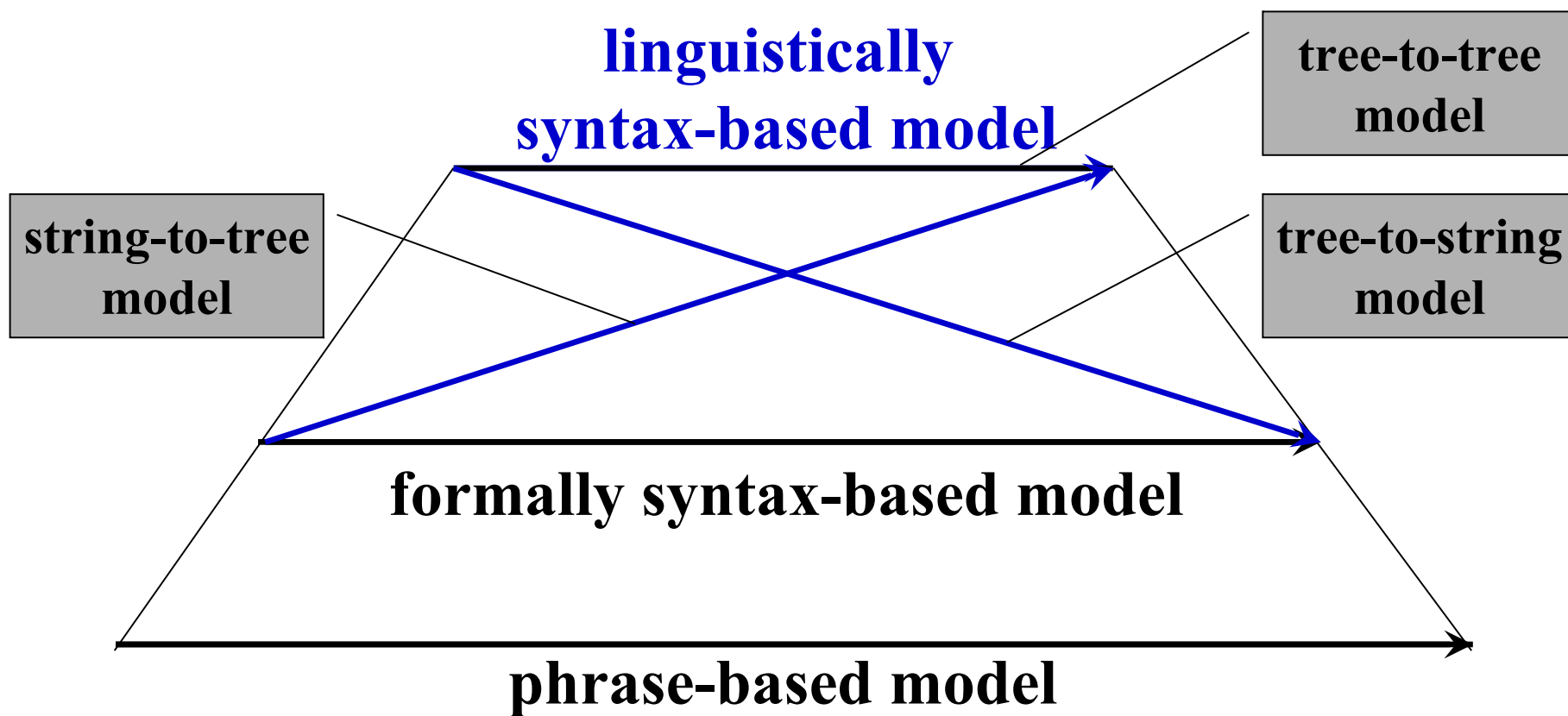
# 基于句法的统计翻译模型 (1)



# 基于句法的统计翻译模型 (1)



# 基于句法的统计翻译模型 (1)



# 基于句法的统计翻译模型 (2)

- 基于句法的统计翻译模型，通常的做法都是分别为源语言和目标语言句子建立某种句法结构，并在这两种句法结构之间建立某种对应关系
- 基于句法的统计翻译模型有两种不同的做法
  - 形式上基于句法的统计翻译模型：并不采用语言学上的句法分析，而是从词语对齐的双语语料库中自动获取某种双语平行的句法结构
  - 语言学上基于句法的统计翻译模型：利用语言学上的句法分析，为源语言句子和目标语言句子建立句法结构，并借助词语对齐建立句法结构的对应关系

# 基于句法的统计翻译模型 (3)

- 语言学上基于句法的统计翻译模型又有三种不同的做法
  - 树到串模型：在源语言端进行句法分析并得到源语言句法结构，然后根据词语对齐建立对应的目标语言句法结构（可称为伪句法结构）
  - 串到树模型：在目标语言端进行句法分析并得到目标语言句法结构，然后根据词语对齐建立对应的源语言句法结构（也是伪句法结构）
  - 树到树模型：在源语言端和目标语言端分别进行句法分析并得到双语的句法结构，然后根据词语对齐建立这两种句法结构之间的对应关系

# 形式上基于句法的模型

- 反向转录语法（ ITG ）和括号转录语法（ BTG ）  
Inversion (Bracketing) Transduction Grammar (ITG,BTG), Dekai Wu 1997
- 有限状态中心词转录机  
Finite-State Head Transducer, Alshawhi 2000
- 基于层次短语的翻译模型  
Hierarchical Phrase-based Model, David Chiang 2005
- 最大熵括号匹配语法的翻译模型  
Maximal Entropy Bracket Transduction Grammar (ME-BTG), Deyi Xiong 2006

# 语言学上基于句法的模型

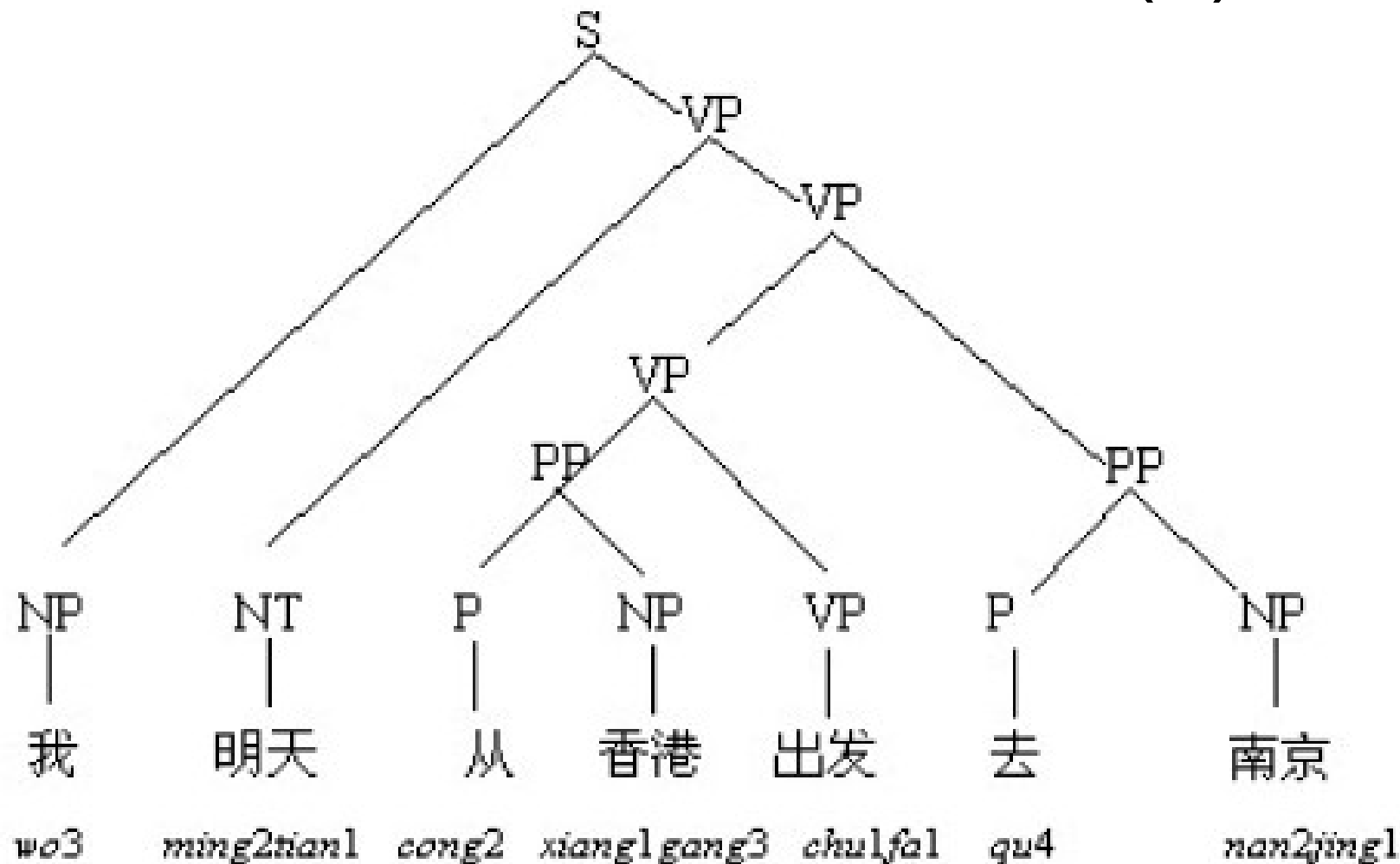
- 串到树模型 **String-to-Tree Model**
  - 美国南加州大学信息科学研究所（**ISI/CSU**）的工作  
**Yamada 2001, Galley 2006, Marcu 2006**
- 树到串模型 **Tree-to-String Model**
  - 中科院计算所的工作  
**Tree-to-string Alignment Template Model (TAT), Liu Yang 2006**
  - 微软研究院的工作（依存模型）  
**Dependency Treelet Translation, Quirk 2005**
- 树到树的模型 **Tree-to-Tree Model**

# 基于 ITG 的机器翻译 (1)

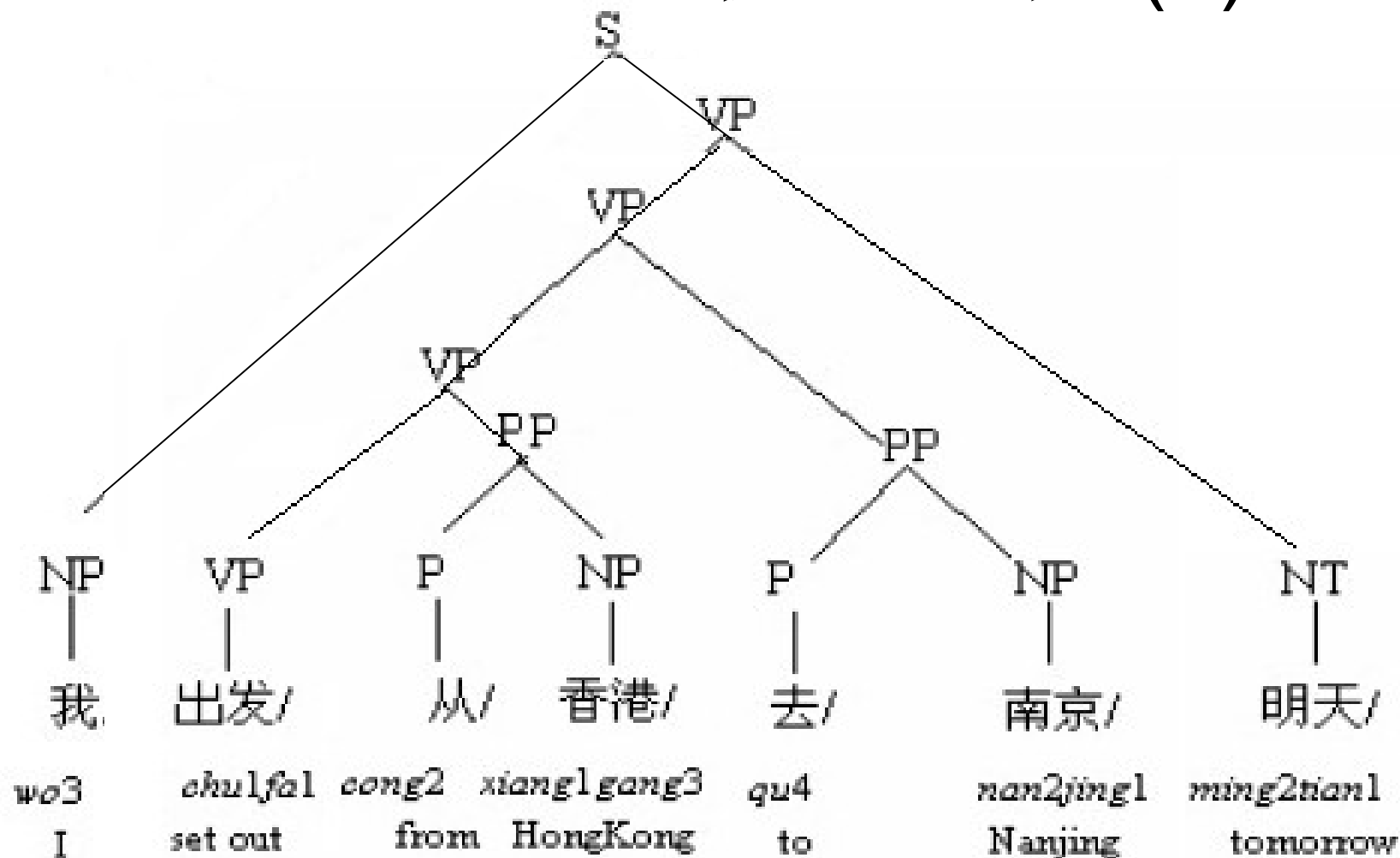
- 训练：  
从词语对齐的双语语料库中自动获得 **ITG** 规则
- 解码：  
类似于传统的基于规则的机器翻译方法
  - 先用 **ITG** 的源语言端规则对源语言进行句法分析
  - 根据 **ITG** 规则的映射关系，确定源语言句法树中每条源语言句法规则对应的目标语言句法规则
  - 生成目标语言句法树



## 基于 ITG 的机器翻译 (2)



## 基于 ITG 的机器翻译 (3)



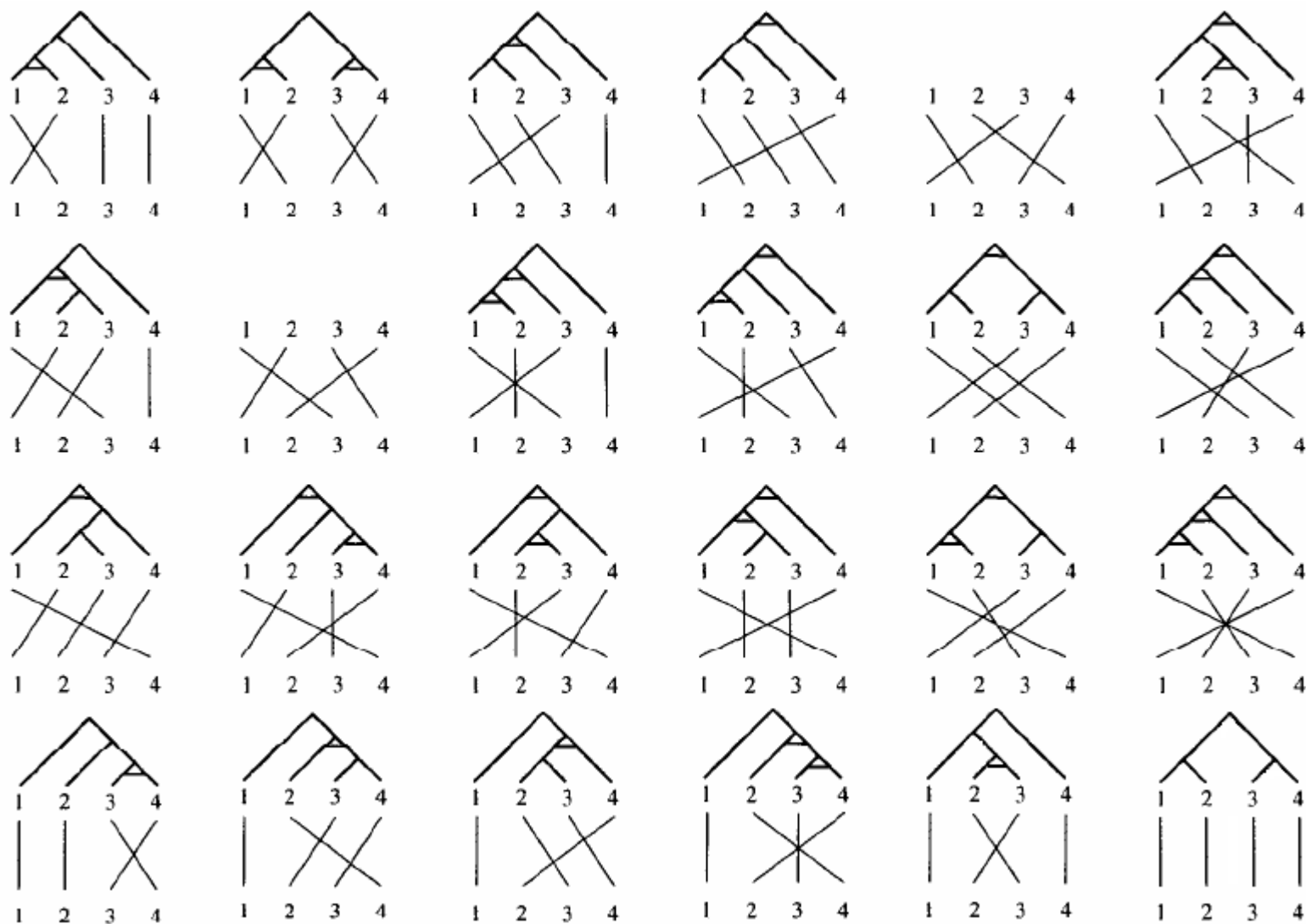
## 基于 ITG 的机器翻译 (4)

- 在 ITG 中，仍然使用了 NP、VP 之类的句法标记，这对于训练语料库提出了比较高的要求
- 如果我们不考虑标记，也就是说，认为所有的标记都是相同的，只有一个非终结符标记 X，那么 ITG 就退化成 BTG

# ITG 约束

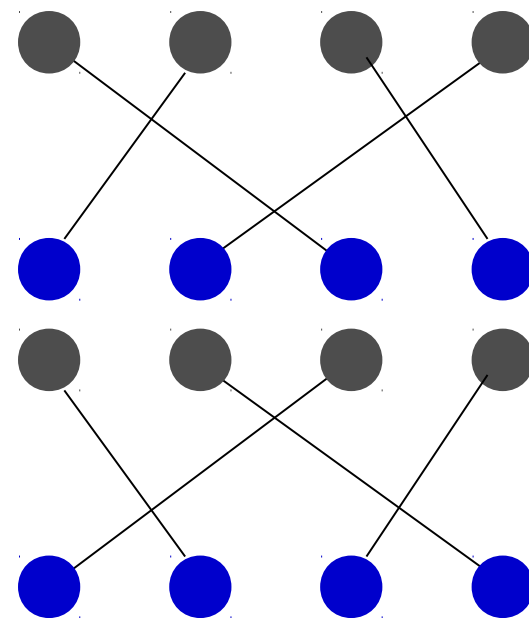
- **BTG 约束 ( BTG constraint )**
  - 只有满足某种 **BTG** 对应关系的目标语言词序才是允许的，否则排除在搜索空间之外
  - 解码的时候，采用类似于 **CYK** 句法分析的方式进行解码，就可以穷尽所有可能的 **BTG** 约束下的词序
  - 在 **BTG** 约束下，可能的对齐方式是多项式级的
  - 无需限制长距离的词序调整

这里给出了四个词在 **BTG** 约束下所有可能的词序调整方案  
其中有两种方案在 **BTG** 约束下是不允许的



# ITG 约束 (3)

$f$	BTG	all matchings	ratio
0	1	1	1.000
1	1	1	1.000
2	2	2	1.000
3	6	6	1.000
4	22	24	0.917
5	90	120	0.750
6	394	720	0.547
7	1806	5040	0.358
8	8558	40320	0.212
9	41586	362880	0.115
10	206098	3628800	0.057
11	1037718	39916800	0.026
12	5293446	479001600	0.011
13	27297738	6227020800	0.004
14	142078746	87178291200	0.002
15	745387038	1307674368000	0.001
16	3937603038	20922789888000	0.000



**word reordering  
which are not  
permitted in BTG**

# ITG 约束 (4) 一个反例

- For Chinese and English, almost true.

– an exception:



- For some other languages with free order, not true.

# 层次短语模型 (1)

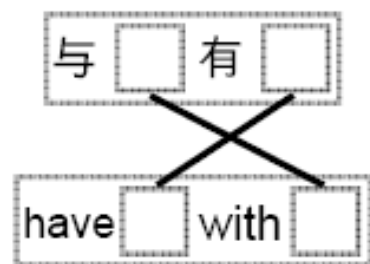
- David Chiang. A Hierarchical Phrase-Based Model for Statistical Machine Translation. ACL2005. (Best Paper Award)
- 本讲义这一部分内容直接引用了以下讲义的部分内容，特此说明并向原作者表示感谢：
  - David Chiang, Hiero: Finding Structure in Statistical Machine Translation, in National University of Singapore



## 层次短语模型 (2)

- 传统的基于短语的翻译模型中，短语是平面的，不能嵌套
- 在层次短语模型中，引入了嵌套的层次短语
- 采用平行上下文无关语法作为理论基础，但只使用唯一的非终结符标记
- 效果比传统的短语模型有很大提高

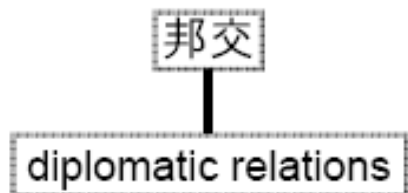
# 用同步语法表示层次短语 (1)



$(X \rightarrow \text{与 } X_1 \text{ 有 } X_2, X \rightarrow \text{have } X_2 \text{ with } X_1)$

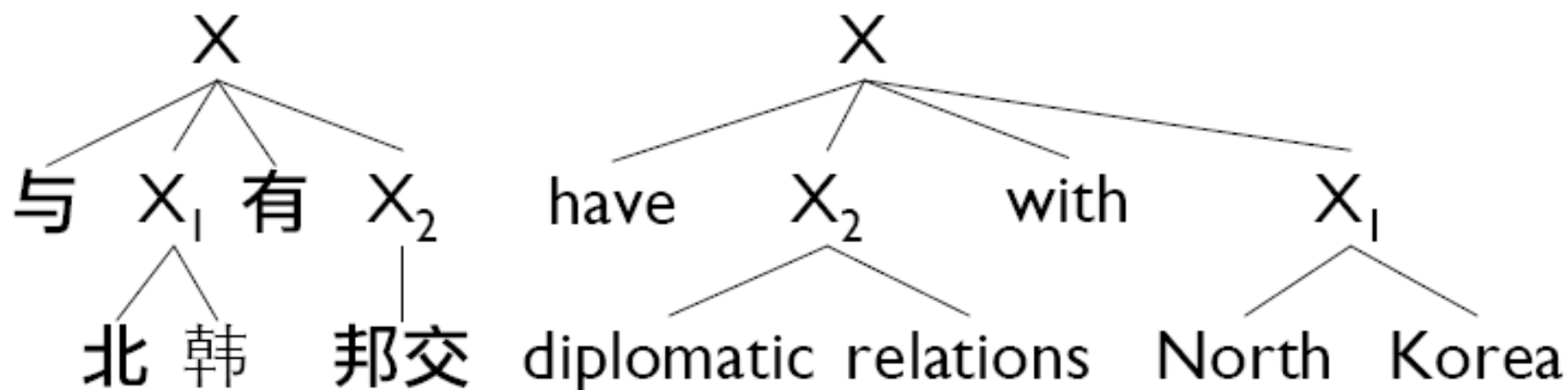


$(X \rightarrow \text{北 韩}, X \rightarrow \text{North Korea})$



$(X \rightarrow \text{邦交}, X \rightarrow \text{diplomatic relations})$

## 用同步语法表示层次短语 (2)



# 规则举例

$X \rightarrow \text{的}$	$X \rightarrow \text{'s}$
$X \rightarrow X_1 \text{ 的 } X_2$	$X \rightarrow \text{the } X_2 \text{ of } X_1$
$X \rightarrow X_1 \text{ 的 } X_2$	$X \rightarrow \text{the } X_2 \text{ that } X_1$
<hr/>	
$X \rightarrow \text{在}$	$X \rightarrow \text{in}$
$X \rightarrow \text{在 } X_1 \text{ 下}$	$X \rightarrow \text{under } X_1$
$X \rightarrow \text{在 } X_1 \text{ 前}$	$X \rightarrow \text{before } X_1$
<hr/>	
$X \rightarrow \text{今年 } X_1$	$X \rightarrow X_1 \text{ this year}$
$X \rightarrow X_1 \text{ 之一}$	$X \rightarrow \text{one of } X_1$
$X \rightarrow X_1 \text{ 总统}$	$X \rightarrow \text{president } X_1$

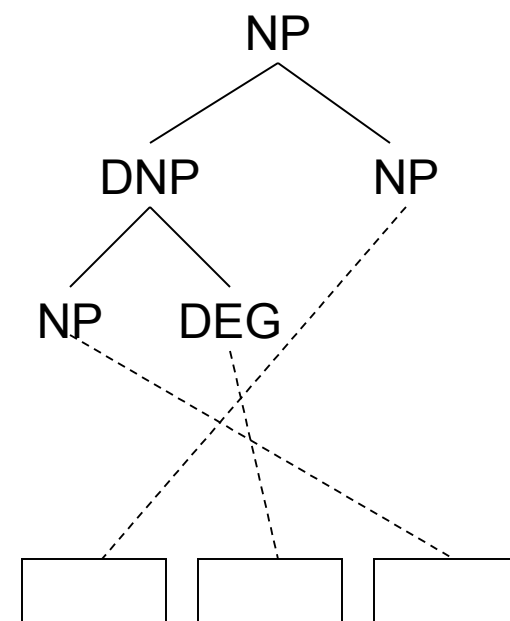
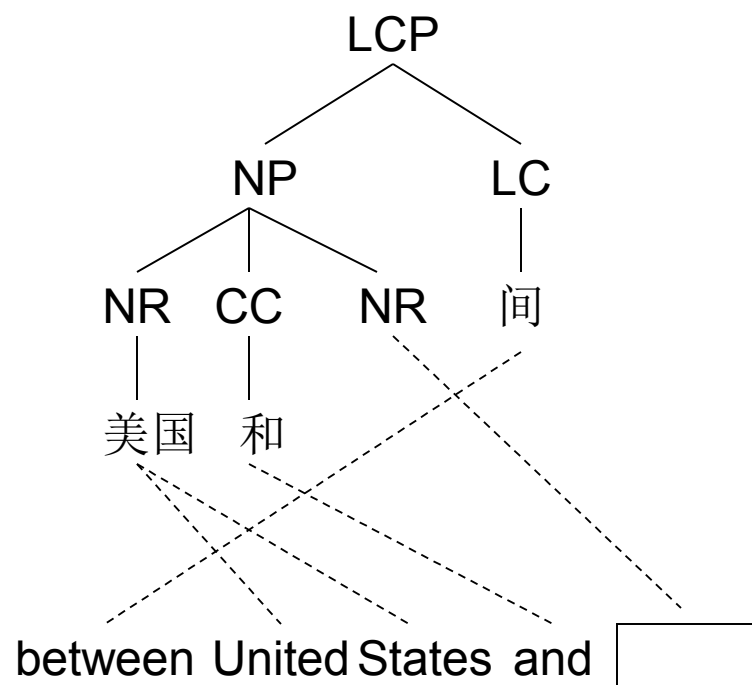
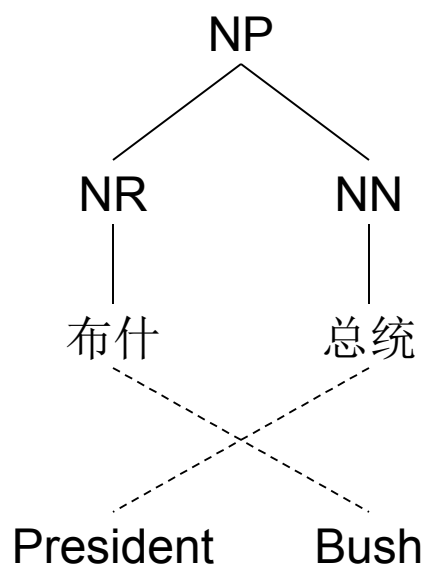
# 树到串翻译模型

- 树到串翻译模型  
Tree-to-String Translation Model
- Yang Liu, Qun Liu, and Shouxun Lin. 2006.  
Tree-to-String Alignment Template for  
Statistical Machine Translation. COLING-ACL  
2006, Sydney, Australia, July 17-21.
- Yang Liu, Yun Huang, Qun Liu and Shouxun  
Lin, Forest-to-String Statistical Translation  
Rules, ACL2007, Prague, Czech , June 2007

# 树到串翻译模型

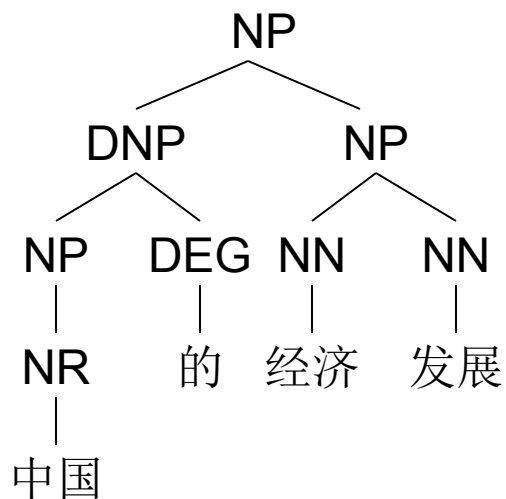
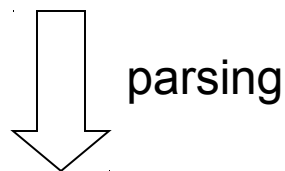
- 树到串统计翻译模型是一种在源语言进行句法分析的基于语言学句法结构的统计翻译模型
- 树到串翻译规则既可以生成终结符也可以生成非终结符，既可以执行局部重排序也可以执行全局重排序
- 从经过词语对齐和源语言句法分析的双语语料库上自底向上自动抽取 TAT
- 自底向上的柱搜索算法

# 树到串翻译规则



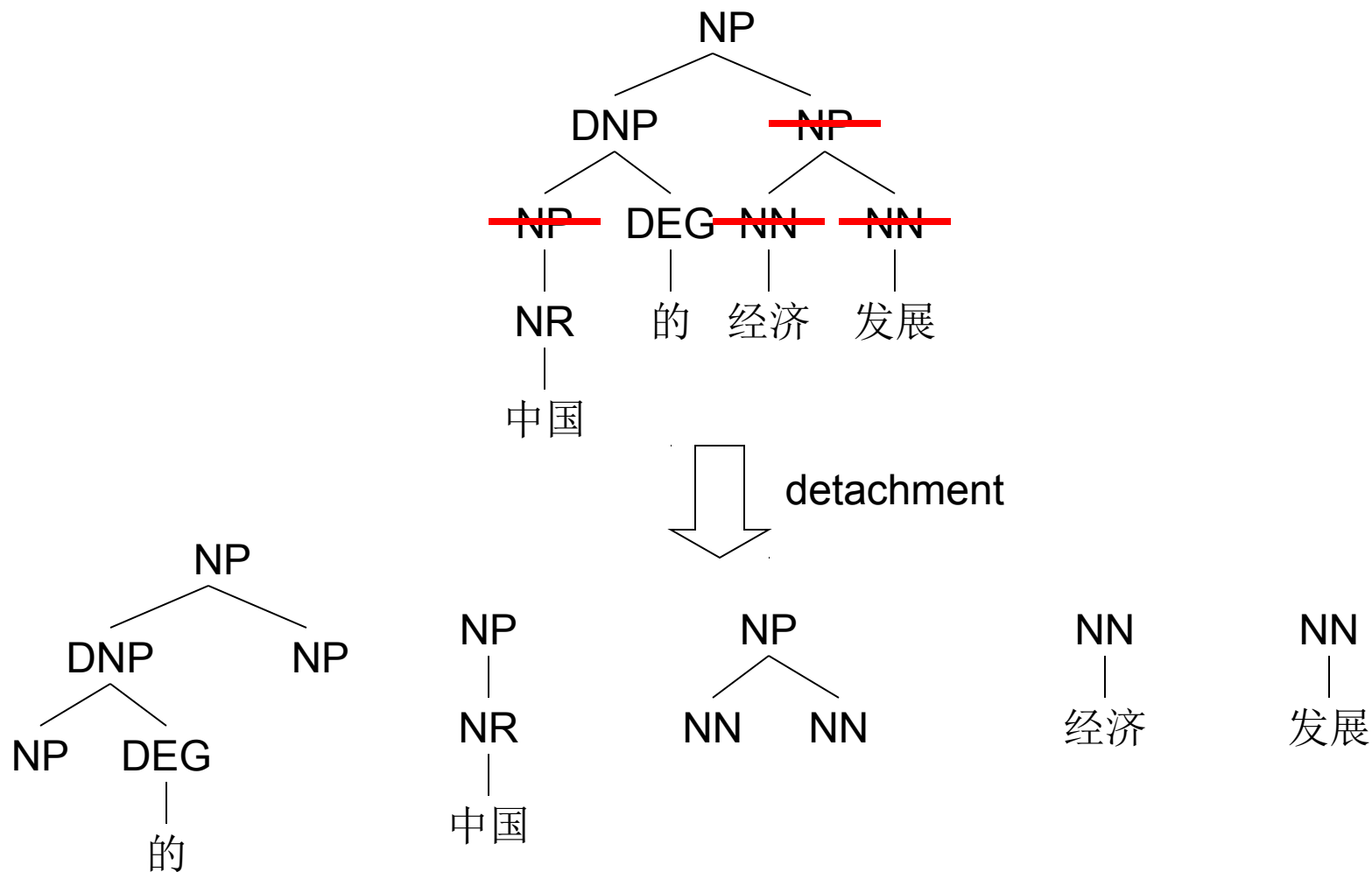
# 翻译过程 : Parsing

中国 的 经济 发展

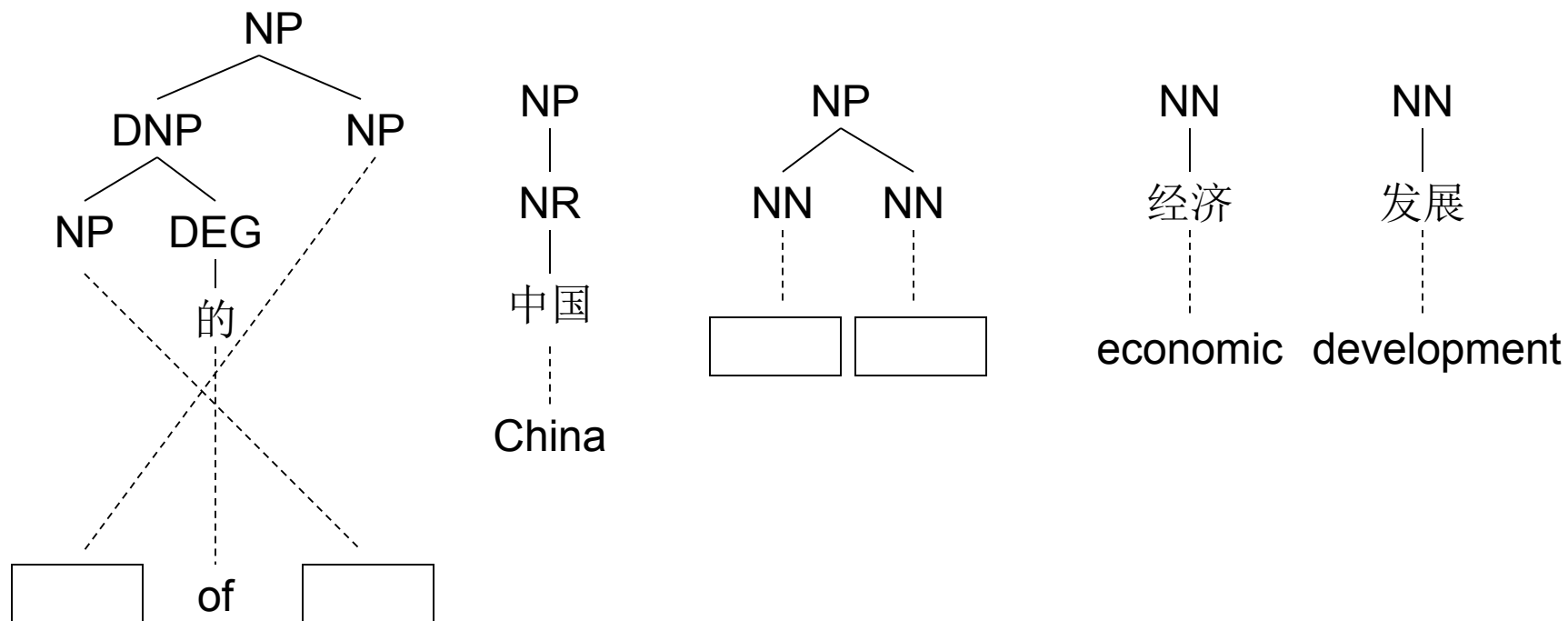




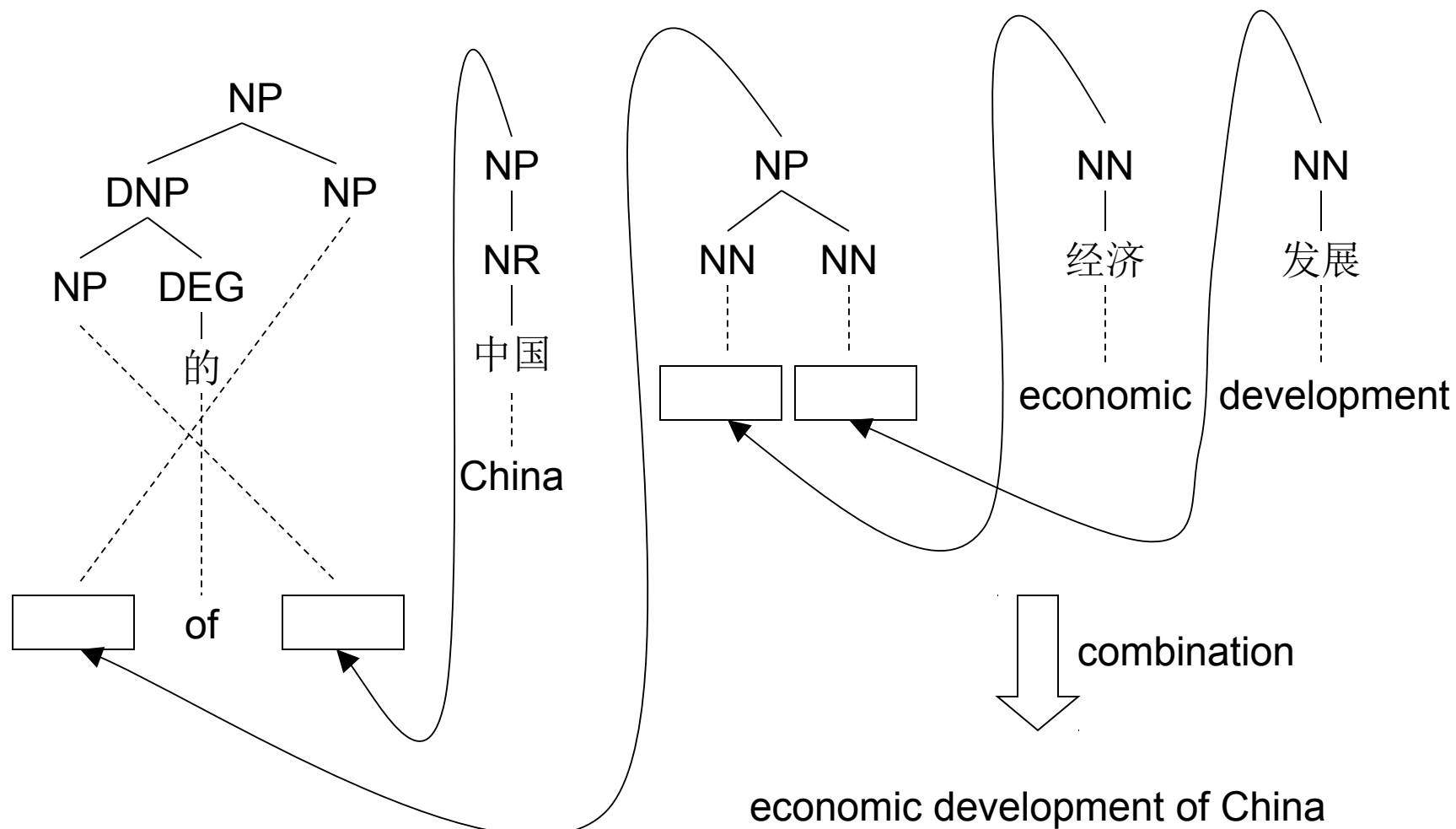
# 翻译过程：Detachment



# 翻译过程 : Production



# 翻译过程：Combination



# 串到树的统计翻译模型 (1)

- **USC-ISI** 的系列工作
- 发表了大量论文，但还没有一个完整的论述
- 性能优异，在 **NIST2006** 汉英项目平常中超过了 **Google** （**Google** 使用的语言模型规模比 **ISI** 大得多）

# 串到树的统计翻译模型 (2)

- 基本思想
  - 在目标语言端进行句法分析
  - 根据目标语言端的句法结构，和词语对齐，建立源语言端的句法结构（伪树）
  - 利用两个句法结构自动抽取带概率的平行上下文无关语法
  - 对平行上下文无关语法进行二叉化
  - 解码时类似规则方法，复杂度等价于句法分析
    - 源文分析
    - 规则映射
    - 译文生成

# 串到树翻译示例 (1)

枪手

被

警方

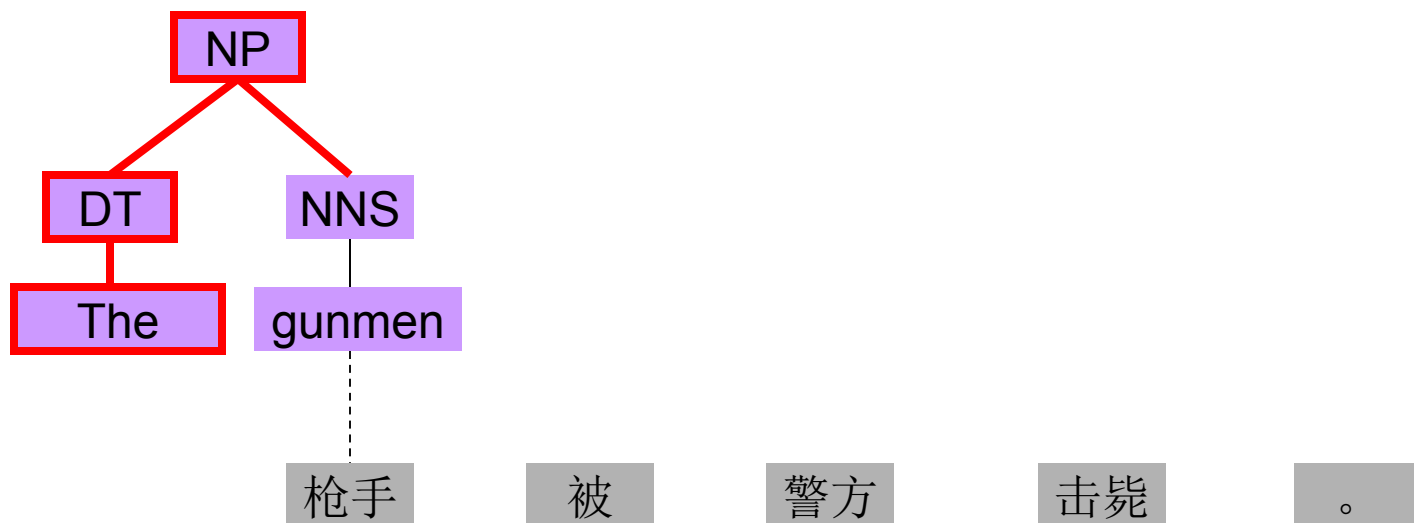
击毙

。

# 串到树翻译示例 (2)

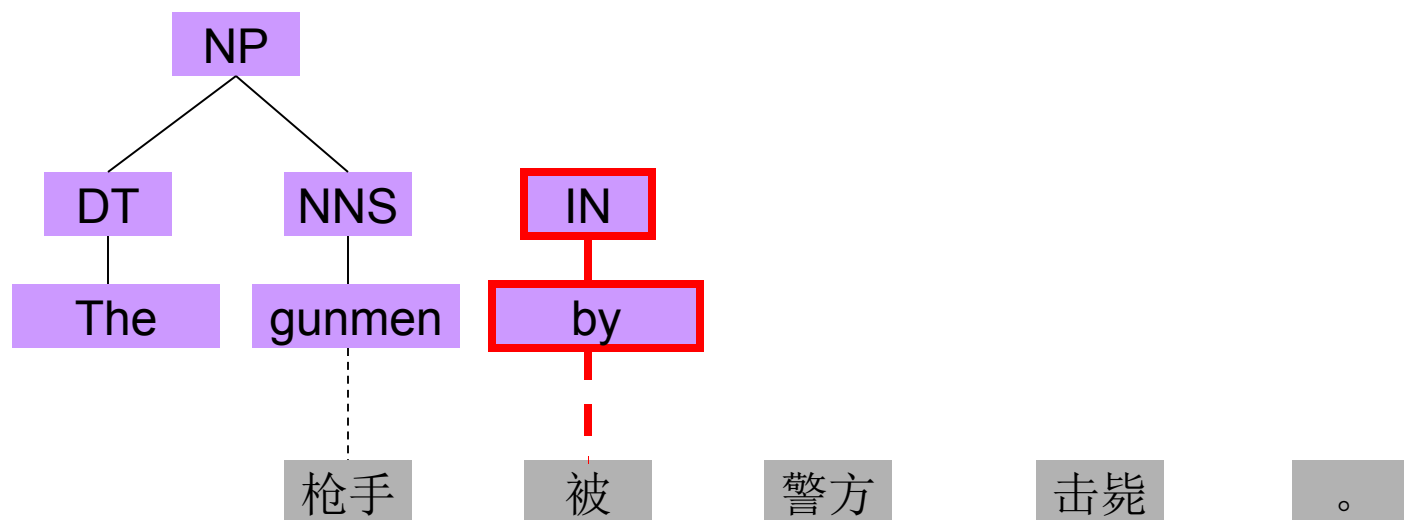


# 串到树翻译示例 (3)

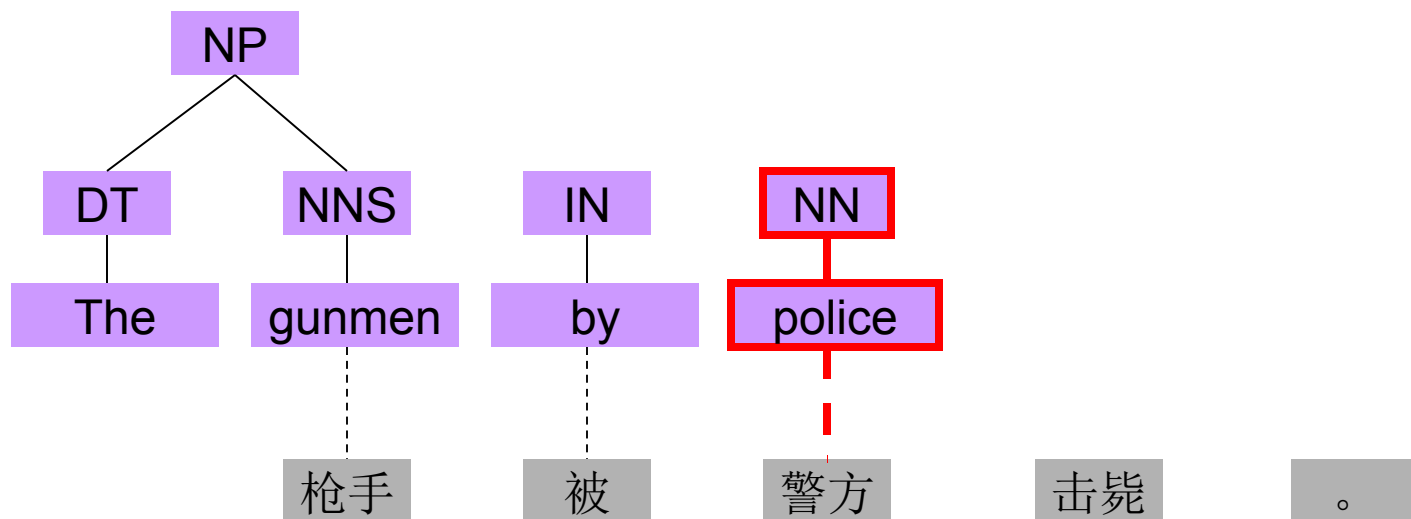




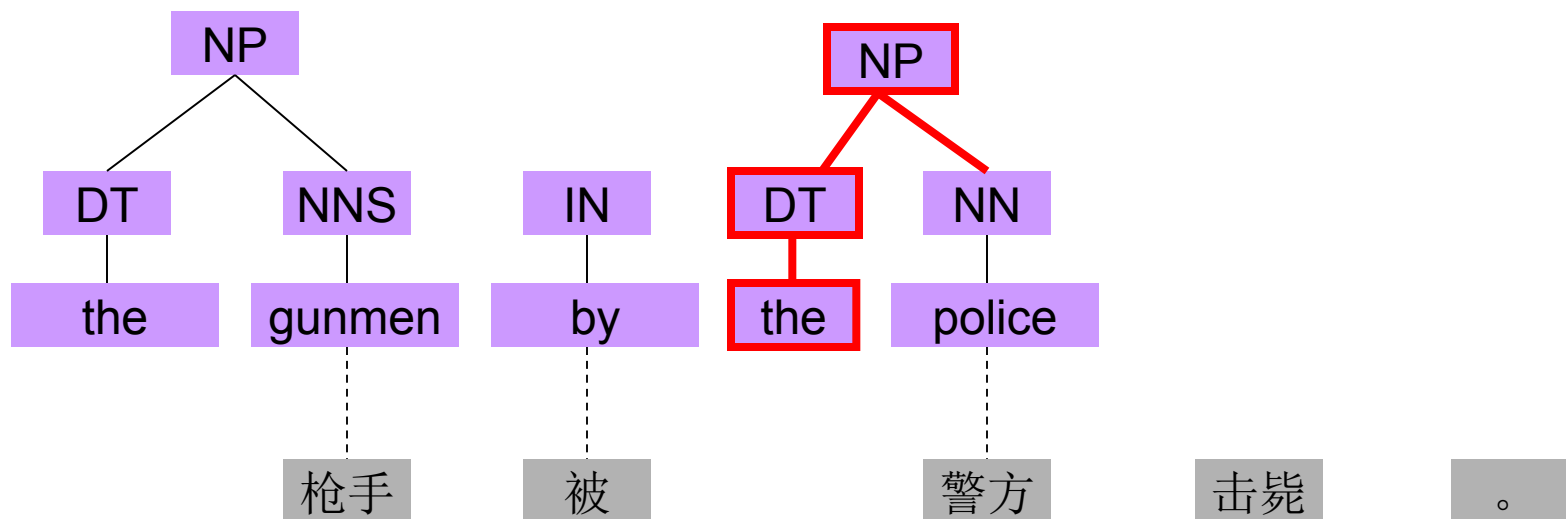
# 串到树翻译示例 (4)



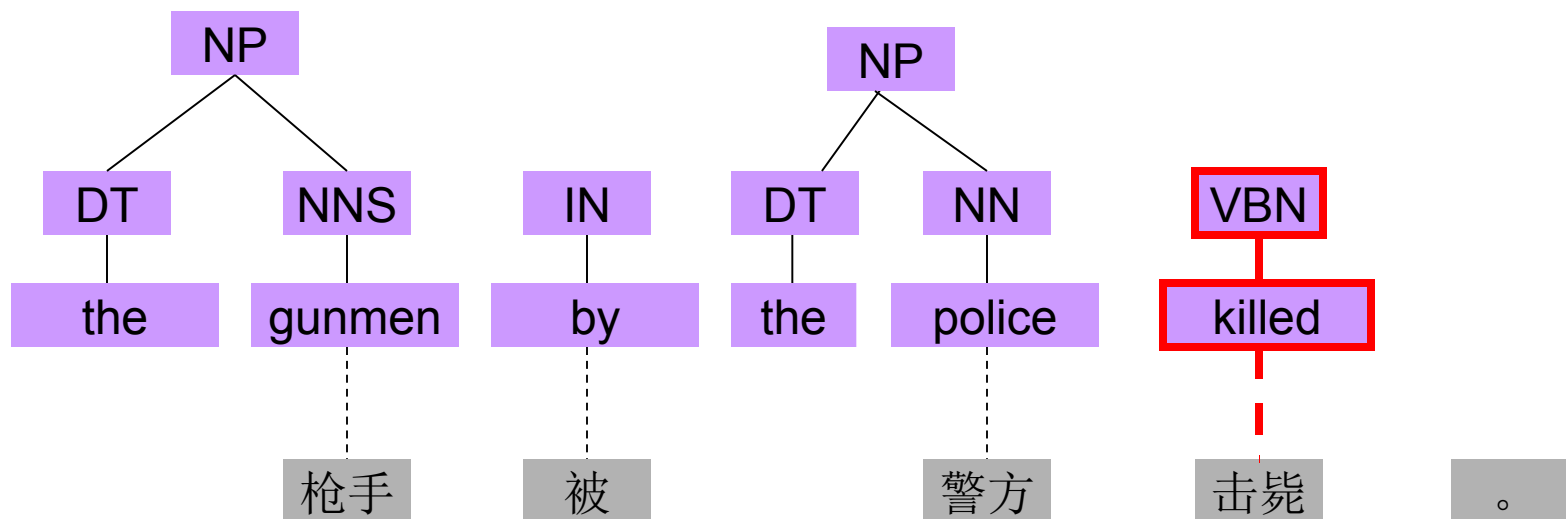
# 串到树翻译示例 (5)



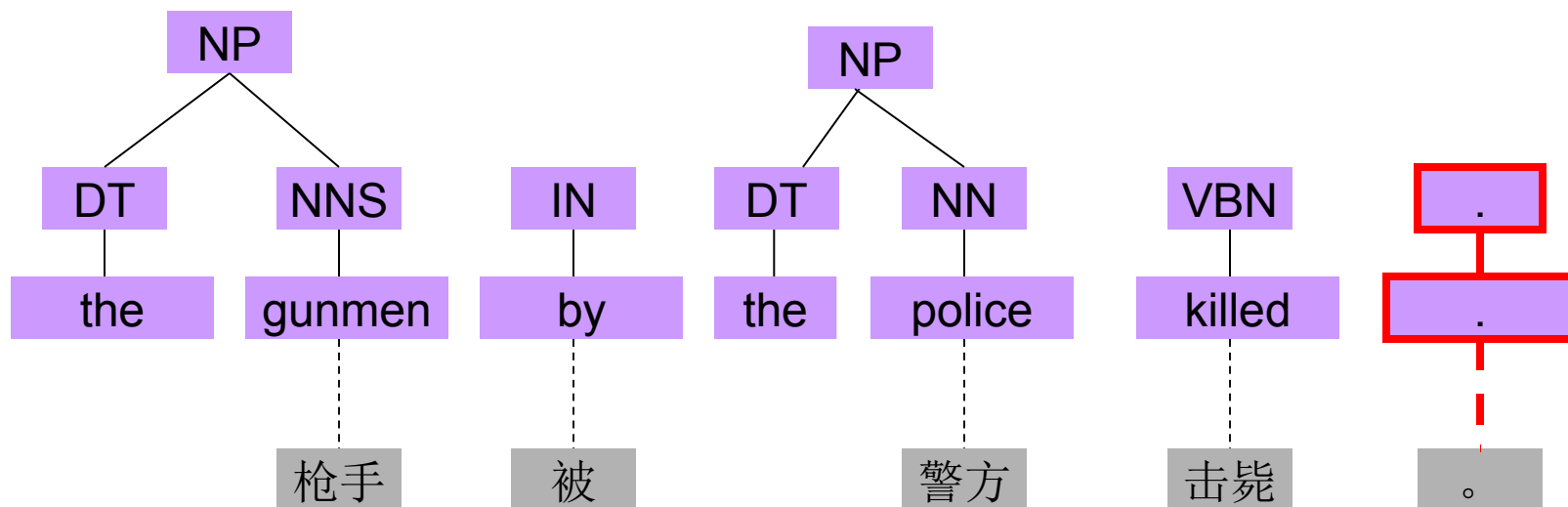
# 串到树翻译示例 (6)



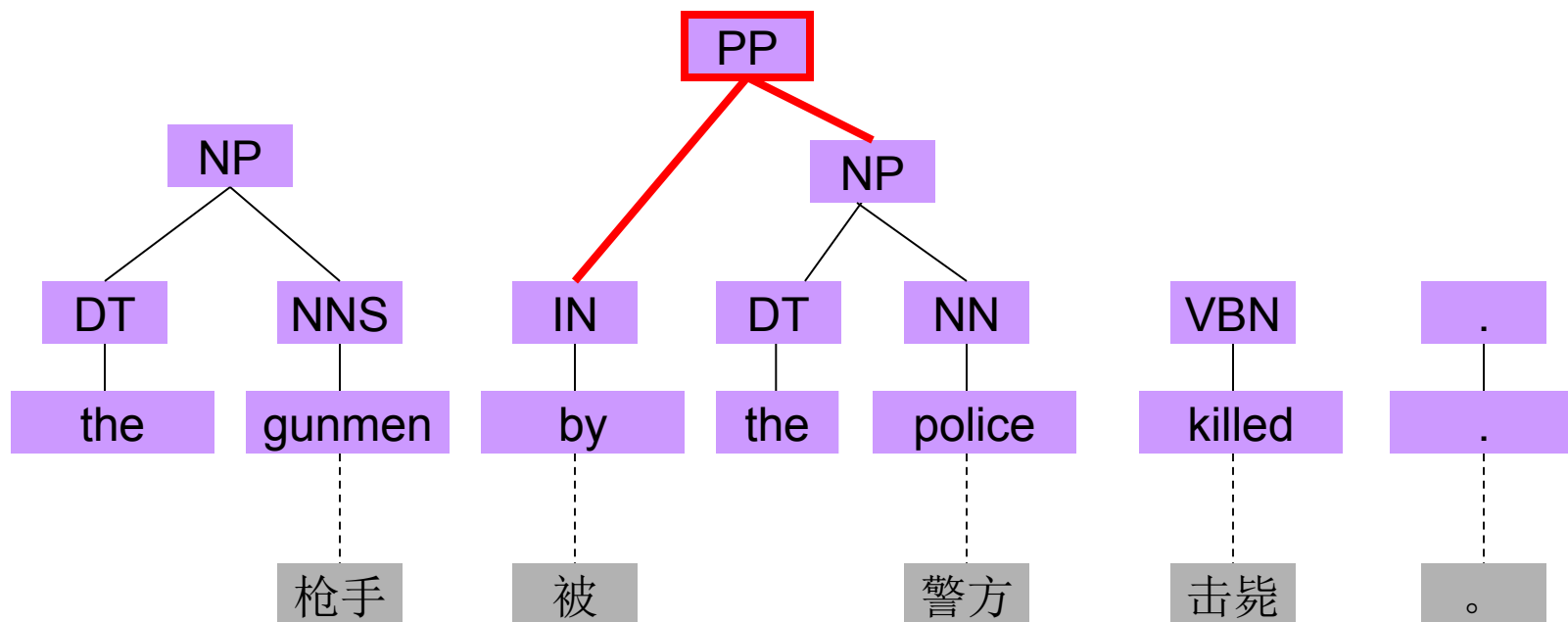
# 串到树翻译示例 (7)



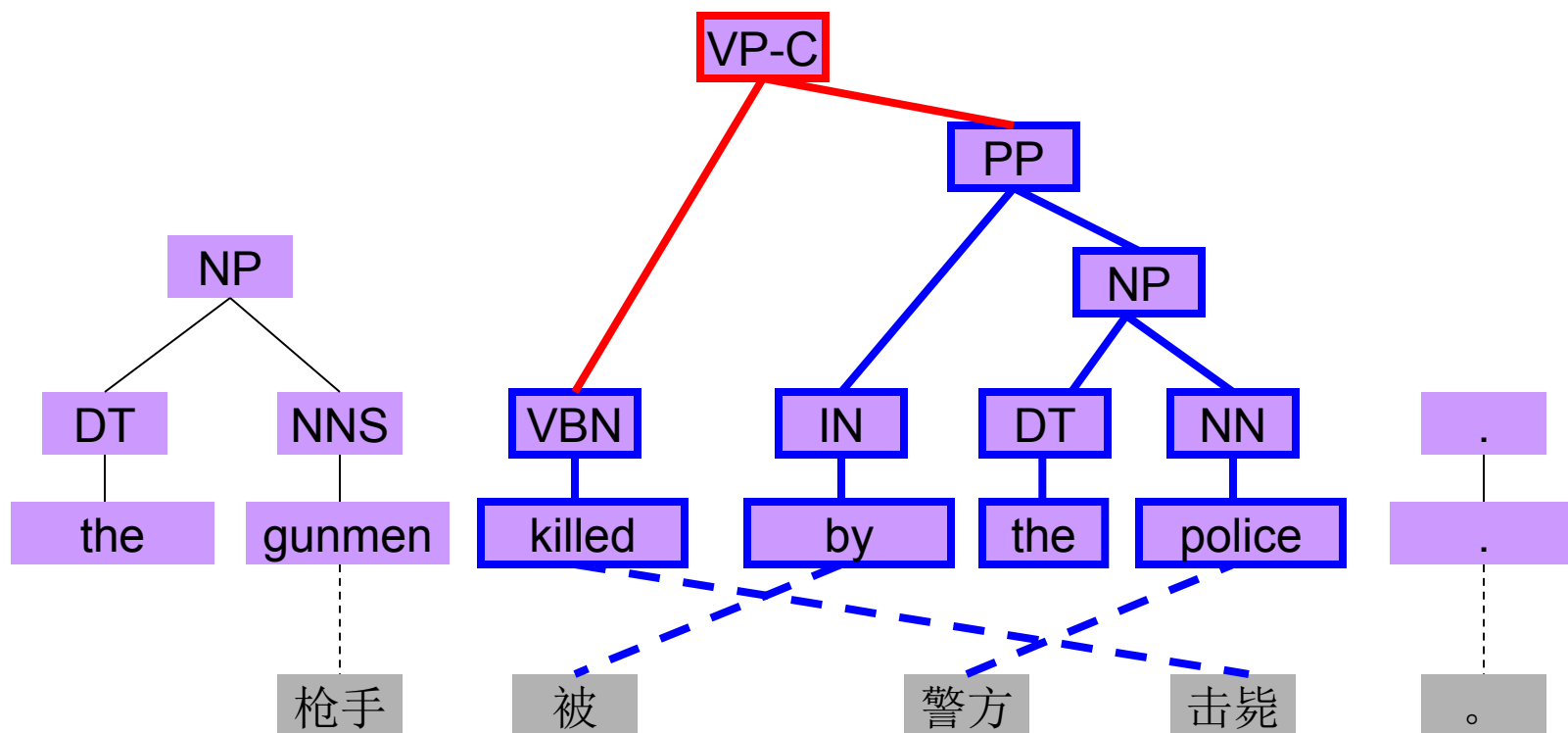
# 串到树翻译示例 (8)



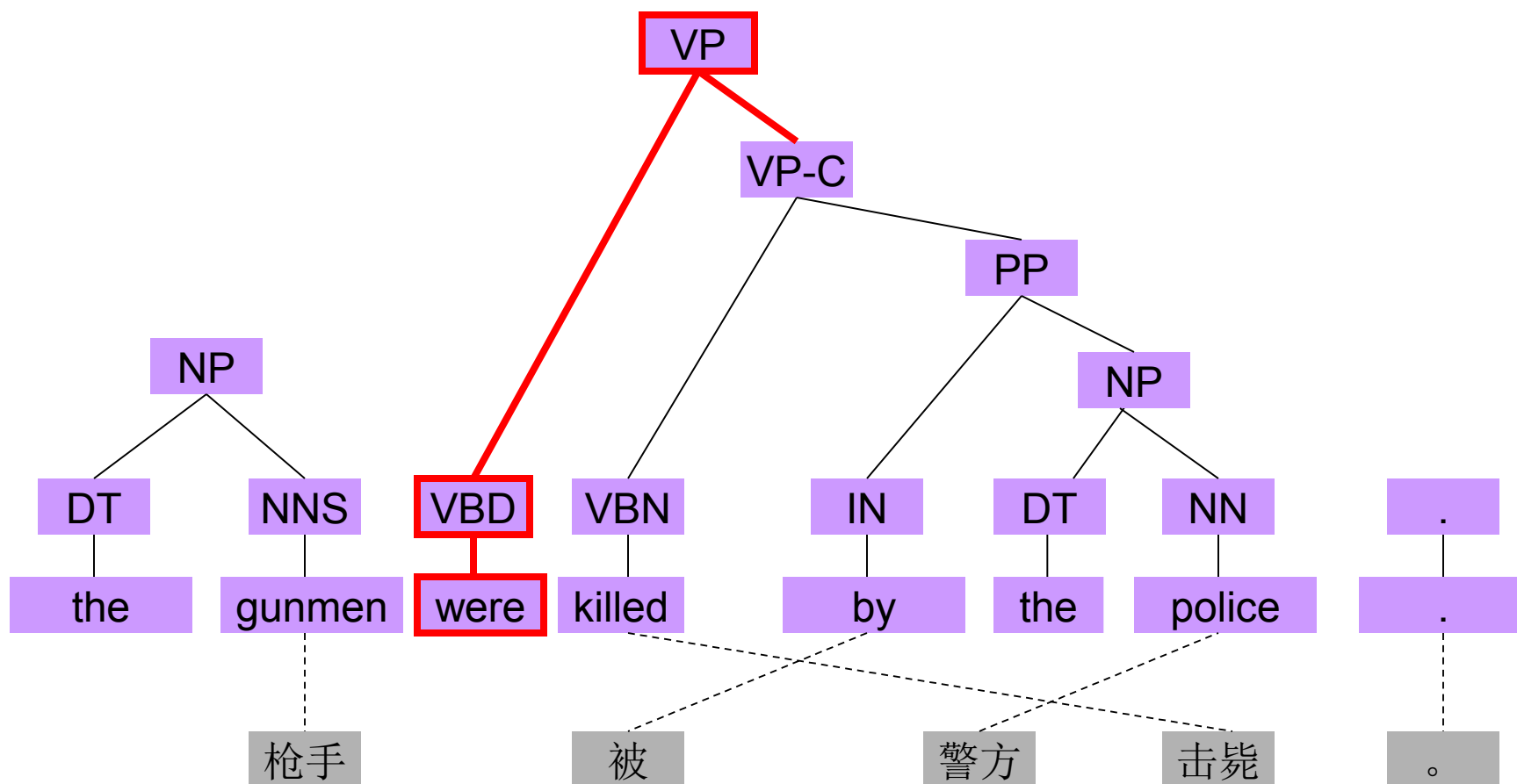
# 串到树翻译示例 (9)



# 串到树翻译示例 (10)

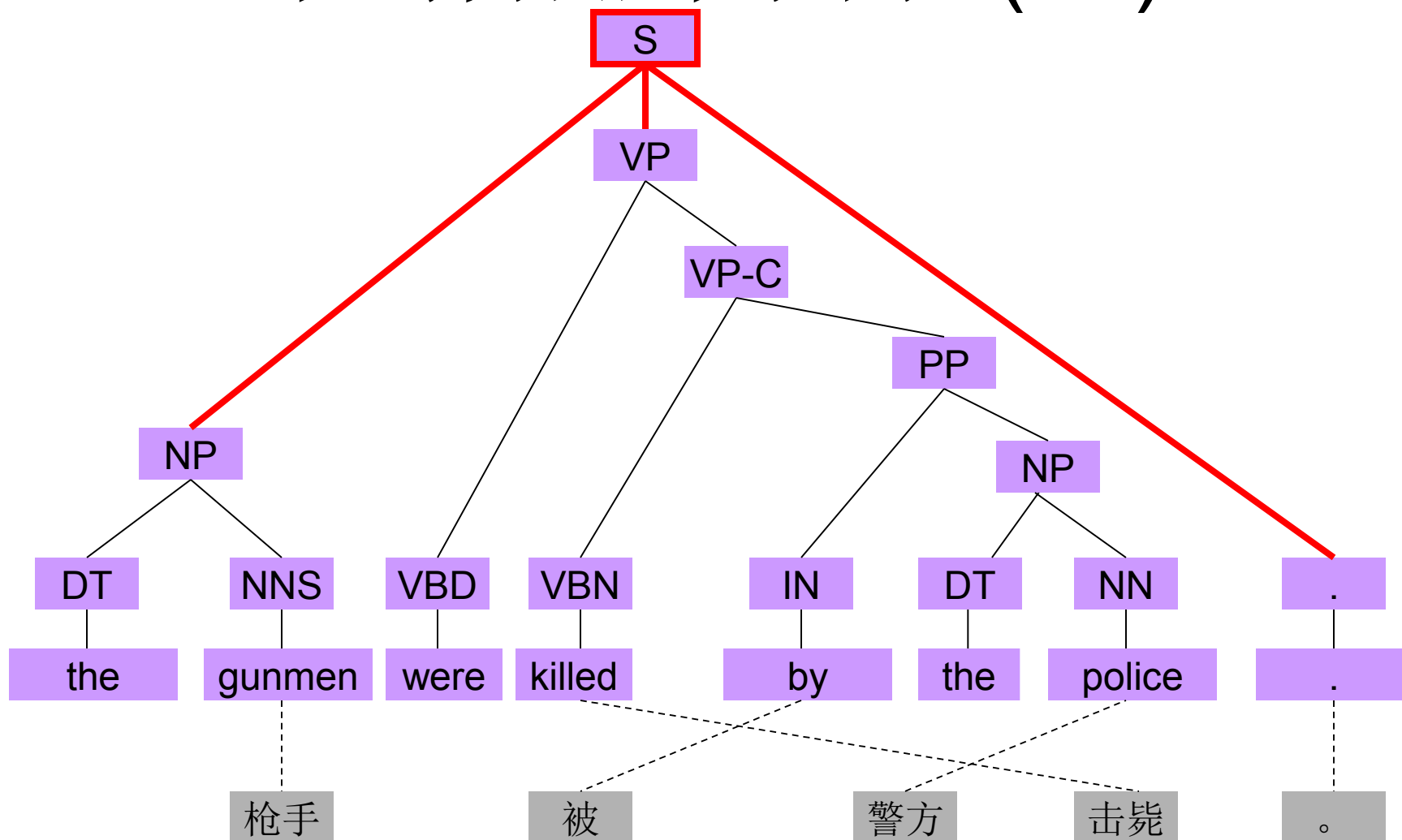


# 串到树翻译示例 (11)





# 串到树翻译示例 (12)



# 内容提要

机器翻译概述

机器翻译方法

机器翻译评价

# 机器翻译评价 (1)

- 最早的机器翻译评价： **ALPAC** 报告
- 机器翻译评价的常用指标
  - 忠实度（ **Adequacy** ）：译文在多大程度上传递了源文的内容；
  - 流利度（ **Fluency** ）：译文是否符合目标语言的语法和表达习惯；
  - 信息度（ **Informative** ）：用户可以从译文中获得信息的程度（通过选择题评分）
- 绝对评价和相对评价

# 机器翻译评价 (2)

- 人工评价
  - 准确
  - 成本极高
  - 不能反复使用
- 自动评价
  - 准确率低
  - 成本低
  - 可以反复使用

# 机器翻译评价 (3)

- 机器翻译的评价一直是机器翻译研究领域中的一个备受关注的问题；
- 机器翻译的自动评价越来越引起重视
  - “评测驱动”成为自然语言处理研究的一个主要动力
  - 大规模语料库的出现、各种机器翻译算法的提出，使得开发过程中频繁的评测成为必需
  - 开发过程中频繁的评测只能通过采用自动评测方法

# 机器翻译的自动评测

- 完全匹配方法
  - 与参考译文完全相同的译文才被认为是正确的
  - 显然该标准过于严格，不适用
- 编辑距离方法
- 基于测试点的方法
- 基于 **N** 元语法的方法

# 基于编辑距离的机器翻译评测 (1)

- 编辑距离定义：

从候选译文到参考译文，所需要进行的插入、删除、替换操作的次数
- 举例说明：
  - 原文： She is a star with the theatre company.
  - 机器译文：她 是 与 剧院 公司 的 一 颗 星 。
  - 参考译文：她 是 剧团 的 明星 。
  - 编辑距离： 6
    - 删除：与 公司 一 颗
    - 替换：剧院\*剧团 星\*明星

# 基于编辑距离的机器翻译评测 (2)

- 单词错误率：编辑距离除以参考译文中单词数
  - 这个指标是从语音识别中借鉴过来的。
  - 由于语音识别的结果语序是不可变的，而机器翻译的结果语序是可变的，显然这个指标存在一定的缺陷。
- 与位置无关的单词错误率：计算编辑距离时，不考虑插入、删除、替换操作的顺序
  - 也就是说，候选译文与参考译文相比，多出或不够的词进行删除或插入操作，其余不同的词进行替换操作。
  - 这个指标与单词错误率相比，允许语序的变化，不过又过于灵活。



# 基于测试点的机器翻译评测 (1)

- 俞士汶等，机器翻译译文质量自动评估系统，中国中文信息学会 1991 年论文集， pp. 314 ~ 319
- 基本思想
  - 对于每一个句子，孤立测试点，简化测试目标（模拟人类标准化考试的办法）
  - 对于每一个句子，采用一种 TDL 语言描述的 BNF 去与译文匹配，匹配成功则正确，否则错误
  - 大批量出题，全面评价机器翻译译文质量

# 基于测试点的机器翻译评测 (2)

- 测试点分组：
  - 单词、词组、词法、语法（初、中、高级）
- 测试点示例：
  - 原文： I am a student.
  - 测试： 译文中出现“学生 / 大学生”为正确
  - 原文： I bought a table with three dollars.
  - 测试： “买”出现在“美元”之后为正确
  - 原文： I bought a table with three legs.
  - 测试： “买”出现在“腿”之前为正确

# 基于测试点的机器翻译评测 (3)

- 优点：
  - 全自动
  - 实验证明，评价结果是可信的
  - 可以按照人类专家的要求进行单项评测
- 缺点
  - 题库的构造需要具有专门知识的专家，并且成本较高

# 基于 N 元语法的机器翻译评测 (1)

- 基本思想
  - 用译文中出现的 N 元组和参考译文中出现的 N 元组相比，计算匹配的 N 元组个数与候选译文的所有 N 元组总个数的比例
  - 允许一个源文有多个参考译文，综合评分
- 参考文献
  - Kishore Papineni, Salim Roukos, Todd Ward, Wei-Jing Zhu, Bleu: a Method for Automatic Evaluation of Machine Translation, IBM Research, RC22176 (W0109-022) September 17, 2001

# 基于 N 元语法的机器翻译评测 (2)

原文：党指挥枪是我党的行动指南。

候选译文：

- It is a guide to action which ensures that the military always obeys the command of the party
- It is to insure the troops forever hearing the activity guidebook that party direct

参考译文：

- It is a guide to action that ensures that the military will forever heed party commands
- It is the guiding principle which guarantees the military forces always being under the command of the party
- It is the practical guide for the army to heed the directions of the party

# 基于 N 元语法的机器翻译评测 (3)

- 两个改进：
  - 对于候选译文中某个  $n$  元接续组出现的次数，如果比参考译文中出现的最大次数还多，要把多出的次数“剪掉”（不作为正确的匹配）。
  - 为了避免“召回率”过低的问题，**BLEU** 的评价标准又对比参考译文更短的句子设计了“惩罚因子”。

# 基于 N 元语法的机器翻译评测 (4)

- BLEU 的总体评价公式如下：

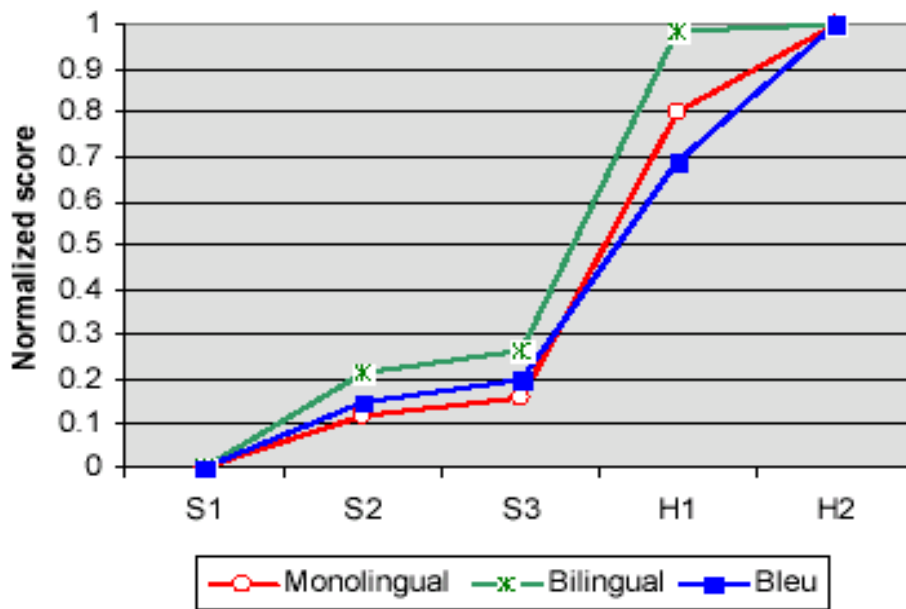
$$\text{BLEU} = \text{BP} \cdot \exp \left( \sum_{n=1}^N w_n \log p_n \right).$$

其中， $p_n$  是出现在参考译文中的  $n$  元词语接续组占候选译文中  $n$  元词语接续组总数的比例， $w_n = 1/N$ ， $N$  为最大的  $n$  元语法阶数（实际取 4）。

$$\text{BP} = \begin{cases} 1 & \text{if } c > r \\ e^{(1-r/c)} & \text{if } c \leq r \end{cases}.$$

其中  $c$  为候选译文中单词的个数， $r$  为参考译文中与  $c$  最接近的译文单词个数。

# 基于 N 元语法的机器翻译评测 (5)



其中 S1、S2、S3 分别是三个不同的机器翻译系统提供的译文，H1 和 H2 是两个人类翻译者提供的译文。蓝线是 BLEU 系统评测的结果，红线是只懂目标语言的人类专家提供的评测结果，绿线是同时懂源语言和目标语言的人类专家提供的评测结果。



# 基于 N 元语法的机器翻译评测 (6)

- 这种方法比较好地模拟了人对机器翻译结果的评价
  - 对于低质量译文比高质量译文的评价更准确;
  - 评价结果与只懂目标语言的人的评价结果更接近  
(相对于懂双语的人而言)
- 优点
  - 全自动
  - 可以提供多种参考译文综合考虑, 结果更全面
  - 容易构造测试集, 不需要专门知识

# 复习思考题

- 访问一些知名的网上翻译网站，直观了解机器翻译
  - SYSTRAN Homepage
  - WordLingo
  - 看世界
- 尝试写一些规则，将英语句子 “**He wrote a book on history.**” 翻译成汉语句子 “他写了一本关于历史的书。”
- 写一个程序实现英语数字、汉语数字和阿拉伯数字之间的互译
- 写一个程序实现英语和汉语之间时间表达式的互译
- 实现一个基于实例的机器翻译中的实例匹配模块，也就是说，将一个输入的句子分解为实例库中的句子片段的组合，并使得这种组合尽可能简单