

Can a Machine Translate Without Knowing What a Verb Is?

Kevin Knight

USC/Information Sciences Institute



Joint work with ISI/LW folks...

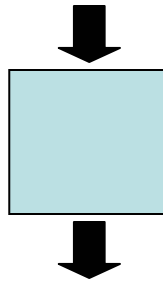
*Daniel Marcu, Wei Wang, Jonathan Graehl,
Michael Pust, Jens Voeckler, Ignacio Thayer,
Radu Soricut, Dragos Munteanu, Alex Fraser,
Steven DeNeefe, Jonathan May*

... and exceptional summer visitors!

*Michel Galley, Mark Hopkins,
Liang Huang, Hao Zhang
Victoria Fossum, David Kauchak*

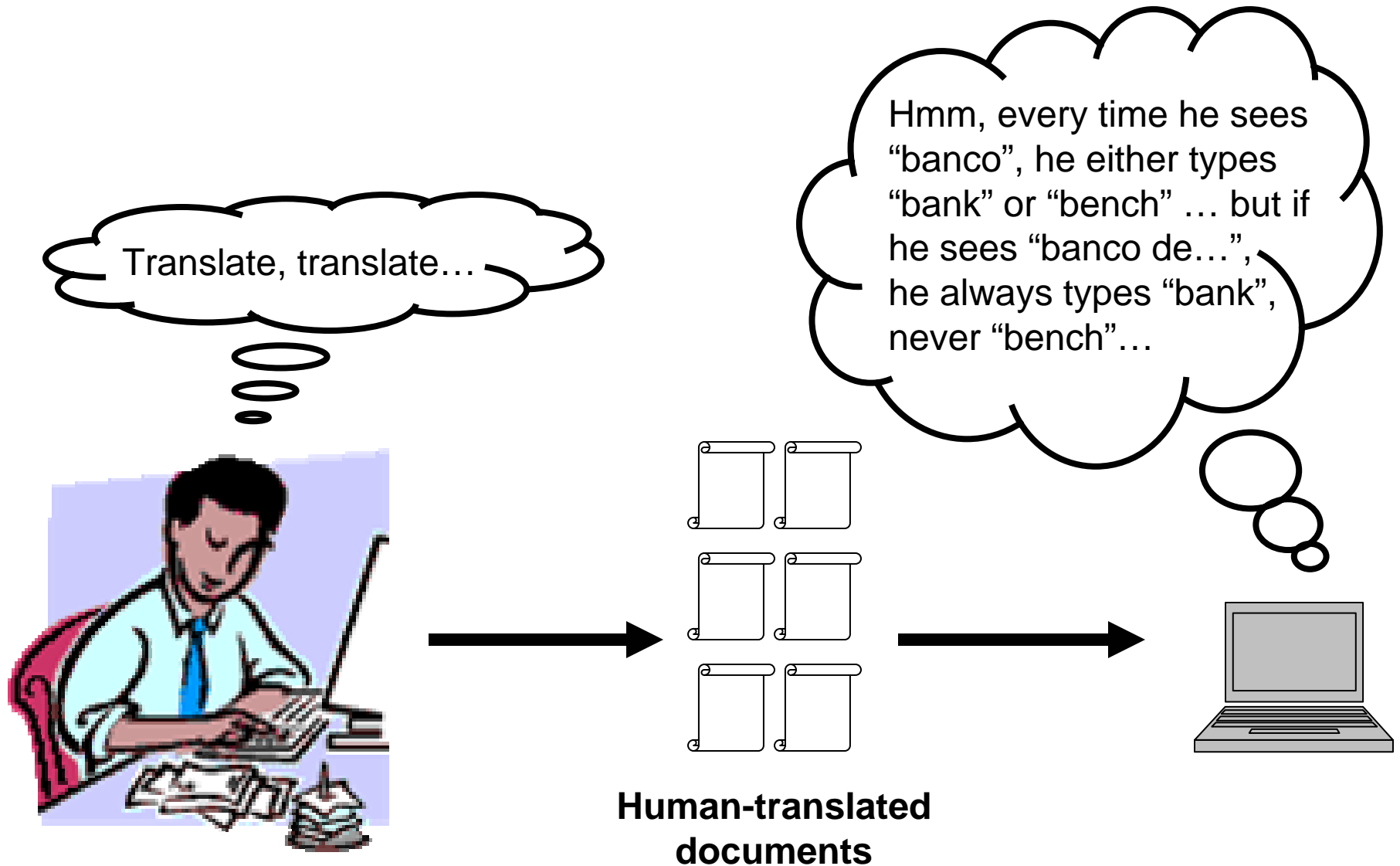
Machine Translation

美国关岛国际机场及其办公室均接获一名自称沙地阿拉伯富商拉登等发出的电子邮件，威胁将会向机场等公众地方发动生化袭击後，关岛经保持高度戒备。

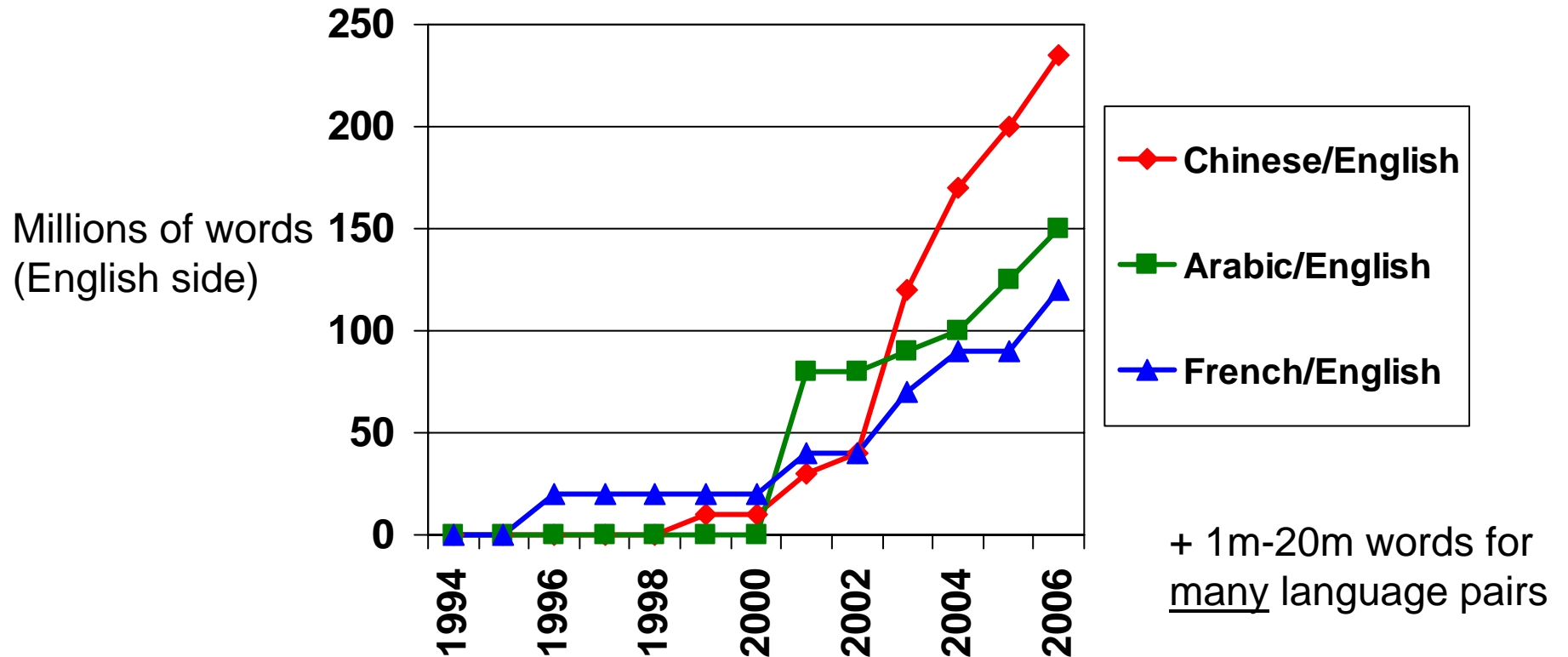


The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport.

Statistical Machine Translation



Ready-to-Use Online Bilingual Data



Bilingual Text (200m words)

English
strings

[illegible]

Bilingual text

Chinese strings



Bilingual Text (200m words)

English
strings

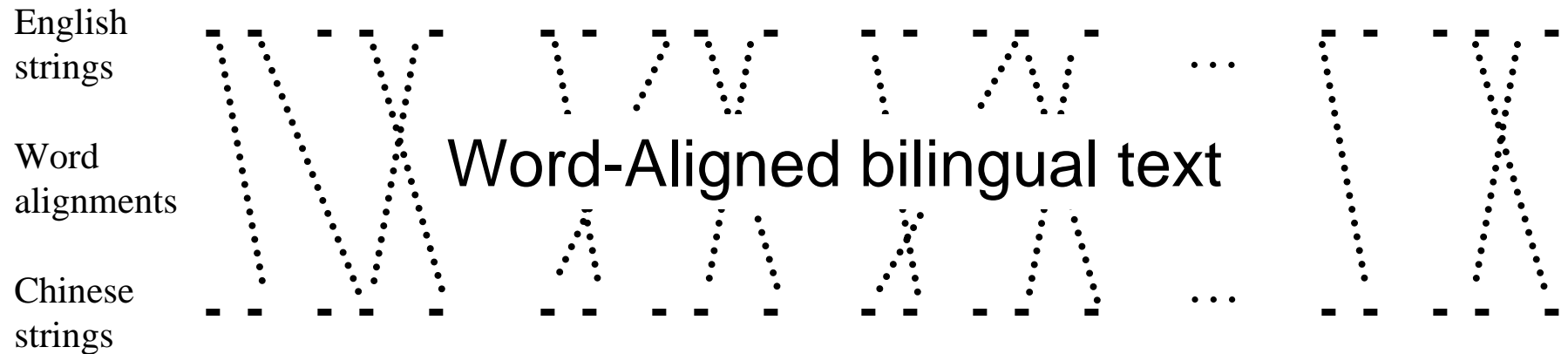
Word alignments

Chinese strings

Word-Aligned bilingual text



Bilingual Text (200m words)



Phrase Pair Extraction [Och & Ney, 2004]

Vast Database of Phrase Pairs



Phrase-Based Translation

这 7人 中包括 来自 法国 和 俄罗斯 的 宇航 员 .

| | | | | | | | | |
|-------|-----------------------|--------------------------------|-------------------|--------------------|-----------------|---------------|-----------------------------|---------------------|
| the | 7 people | including | by some | and | the russian | the | the astronauts | , |
| it | 7 people included | | by france | and the | the russian | | international astronautical | of rapporteur . |
| this | 7 out | including the | from | the french | and the russian | the fifth | | . |
| these | 7 among | including from | | the french and | of the russian | of | space | members . |
| that | 7 persons | including from the | | of france | and to | russian | of the | aerospace members . |
| | 7 include | | from the | of france and | russian | | astronauts | . the |
| | 7 numbers include | | from france | | and russian | | of astronauts who | . ” |
| | 7 populations include | | those from france | | and russian | | astronauts . | |
| | 7 deportees included | | come from | france | and russia | in | astronautical | personnel ; |
| | 7 philtrum | including those from | | france and | russia | a space | | member |
| | | including representatives from | | france and the | russia | | astronaut | |
| | | include | came from | france and russia | | by cosmonauts | | |
| | | include representatives from | | french | and russia | | cosmonauts | |
| | | include | came from france | | and russia 's | | cosmonauts . | |
| | | includes | coming from | french and | russia 's | | cosmonaut | |
| | | | | french and russian | | 's | astronavigation | member . |
| | | | | french | and russia | | astronauts | |
| | | | | | and russia 's | | | special rapporteur |
| | | | | | , and russia | | | rapporteur |
| | | | | | , and russia | | | rapporteur . |
| | | | | | , and russia | | | |
| | | | | | or | russia 's | | |

Table 1: #11# the seven - member crew includes astronauts from france and russia .

Scoring: Try to use phrase pairs that have been frequently observed.
Try to output a sentence with frequent English word sequences.

Phrase-Based Translation

| | | | | | | | | | | |
|-------|-----------------------|--------------------------------|------------------|--------------------|-----------------------------|-----------------|--------------------|---------|---|-------|
| 这 | 7人 | 中包括 | 来自 | 法国 | 和 | 俄罗斯 | 的 | 宇航 | 员 | . |
| the | 7 people | including | by some | and | the russian | the | the astronauts | | | , |
| it | 7 people included | by france | and the | the russian | international astronautical | of rapporteur | | | | . |
| this | 7 out | including the | from | the french | and the russian | the fifth | | | | . |
| these | 7 among | including from | the french and | of the russian | of | space | members | | | . |
| that | 7 persons | including from the | of france | and to | russian | of the | aerospace | members | | . |
| | 7 include | from the | of france and | russian | astronauts | | | | | . the |
| | 7 numbers include | from france | and russian | of astronauts who | | | | | | ." |
| | 7 populations include | those from france | and russian | astronauts | | | | | | . |
| | 7 deportees included | come from | france | and russia | in | astronautical | personnel | | | ; |
| | 7 philtrum | including those from | france and | russia | a space | member | | | | |
| | | including representatives from | france and the | russia | astronaut | | | | | |
| | | include | came from | france and russia | by cosmonauts | | | | | |
| | | include representatives from | french | and russia | cosmonauts | | | | | |
| | | include | came from france | and russia 's | cosmonauts | | | | | . |
| | | includes | coming from | french and | russia 's | cosmonaut | | | | |
| | | | | french and russian | 's | astronavigation | member | | | . |
| | | | | french | and russia | astronauts | | | | |
| | | | | and russia 's | | | special rapporteur | | | |
| | | | | , and | russia | | rapporteur | | | |
| | | | | , and russia | | | rapporteur | | | . |
| | | | | , and russia | | | | | | |
| | | | | or | russia 's | | | | | |

Table 1: #11# the seven - member crew includes astronauts from france and russia .

Scoring: Try to use phrase pairs that have been frequently observed.
Try to output a sentence with frequent English word sequences.

Phrase-Based Translation

| | | | | | | | | | | |
|-------|-----------------------|--------------------------------|------------------|--------------------|-----------------|-------------------|-----------------------------|--------------------|---------|-------|
| 这 | 7人 | 中包括 | 来自 | 法国 | 和 | 俄罗斯 | 的 | 宇航 | 员 | . |
| the | 7 people | including | by some | and | the russian | the | the astronauts | | | , |
| it | 7 people included | by france | | and the | the russian | | international astronautical | of rapporteur | | . |
| this | 7 out | including the | from | the french | and the russian | the fifth | | | | . |
| these | 7 among | including from | | the french and | of the russian | of | space | members | | . |
| that | 7 persons | including from the | | of france | and to | russian | of the | aerospace | members | |
| | 7 include | from the | | of france and | russian | | astronauts | | | . the |
| | 7 numbers include | from france | | and russian | | of astronauts who | | | | . |
| | 7 populations include | those from france | | and russian | | astronauts | | | | . |
| | 7 deportees included | come from | france | and russia | | in | astronautical | personnel | | ; |
| | 7 philtrum | including those from | france and | russia | | a space | | member | | |
| | | including representatives from | france and the | russia | | astronaut | | | | |
| | | include | came from | france and russia | | by cosmonauts | | | | |
| | | include representatives from | french | and russia | | cosmonauts | | | | |
| | | include | came from france | and russia 's | | cosmonauts | | | | . |
| | | includes | coming from | french and | russia 's | | cosmonaut | | | |
| | | | | french and russian | 's | astronavigation | member | | | . |
| | | | | french | and russia | astronauts | | | | |
| | | | | and russia 's | | | | special rapporteur | | |
| | | | | , and | russia | | | rapporteur | | |
| | | | | , and russia | | | | rapporteur | | . |
| | | | | , and russia | | | | | | |
| | | | | or | russia 's | | | | | |

Table 1: #11# the seven - member crew includes astronauts from france and russia .

Scoring: Try to use phrase pairs that have been frequently observed.
Try to output a sentence with frequent English word sequences.

Phrase-Based Translation

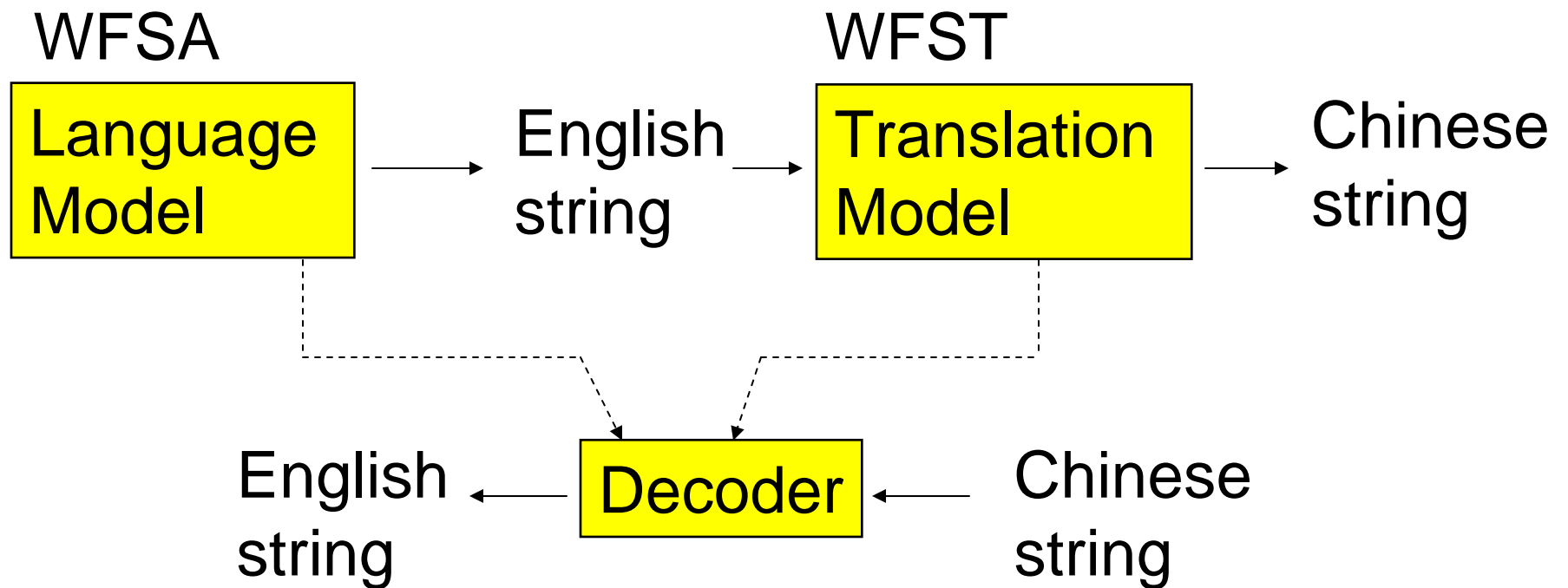
| | | | | | | | | | | |
|---|----|-----|----|----|---|-----|---|----|---|---|
| 这 | 7人 | 中包括 | 来自 | 法国 | 和 | 俄罗斯 | 的 | 宇航 | 员 | . |
|---|----|-----|----|----|---|-----|---|----|---|---|

| | | | | | | | | |
|-------|-----------------------|--------------------------------|------------------|--------------------|-----------------|-------------------|-----------------------------|-----------------|
| the | 7 people | including | by some | and | the russian | the | the astronauts | , |
| it | 7 people included | by france | | and the | the russian | | international astronautical | of rapporteur . |
| this | 7 out | including the | from | the french | and the russian | the fifth | | . |
| these | 7 among | including from | | the french | and | of the russian | of | space |
| that | 7 persons | including from the | | of france | and to | russian | of the | aerospace |
| | 7 include | from the | | of france and | russian | | astronauts | . |
| | 7 numbers include | from france | | and russian | | of astronauts who | | . |
| | 7 populations include | those from france | | and russian | | astronauts . | | |
| | 7 deportees included | come from | france | and russia | | in | astronautical | personnel |
| | 7 philtrum | including those from | france and | russia | | a space | | member |
| | | including representatives from | france and the | russia | | astronaut | | |
| | | include | came from | france and russia | | by cosmonauts | | |
| | | include representatives from | french | and russia | | cosmonauts | | |
| | | include | came from france | and russia 's | | cosmonauts . | | |
| | | includes | coming from | french and | russia 's | cosmonaut | | |
| | | | | french and russian | 's | astronavigation | member . | |
| | | | | french | and russia | astronauts | | |
| | | | | and russia 's | | | special rapporteur | |
| | | | | , and | russia | | rapporteur | |
| | | | | , and russia | | | rapporteur . | |
| | | | | , and russia | | | | |
| | | | | or | russia 's | | | |

Table 1: #11# the seven - member crew includes astronauts from france and russia .

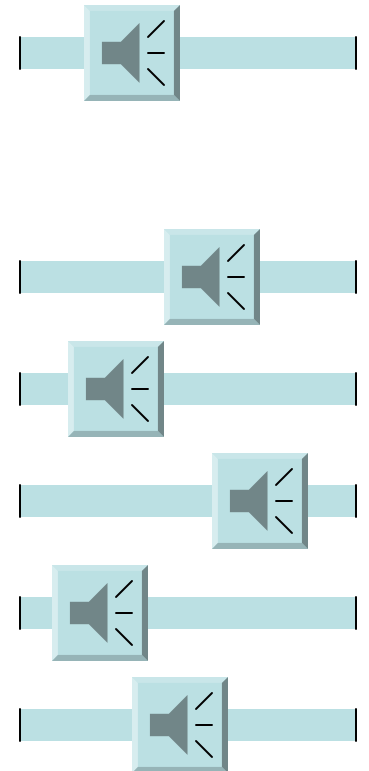
Scoring: Try to use phrase pairs that have been frequently observed.
 Try to output a sentence with frequent English word sequences.

Phrase-Based Noisy Channel



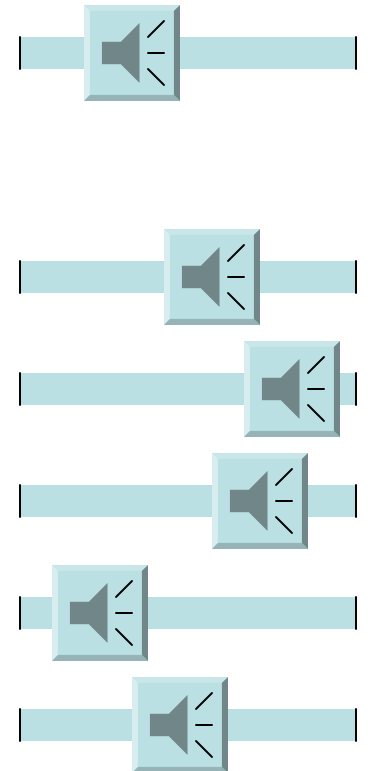
Features and Tuning

- English trigram language model
- Phrase pairs
 - Conditional probability
 - Bad-phrase spotter
 - Word-drop spotter
 - “Move Me” preference
- English output length



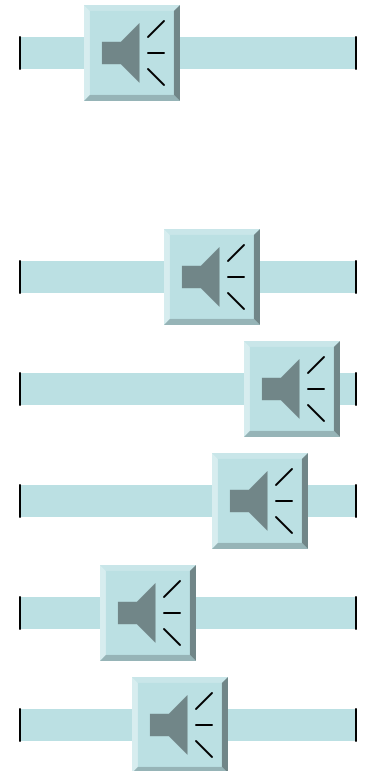
Features and Tuning

- English trigram language model
- Phrase pairs
 - Conditional probability
 - Bad-phrase spotter
 - Word-drop spotter
 - “Move Me” preference
- English output length



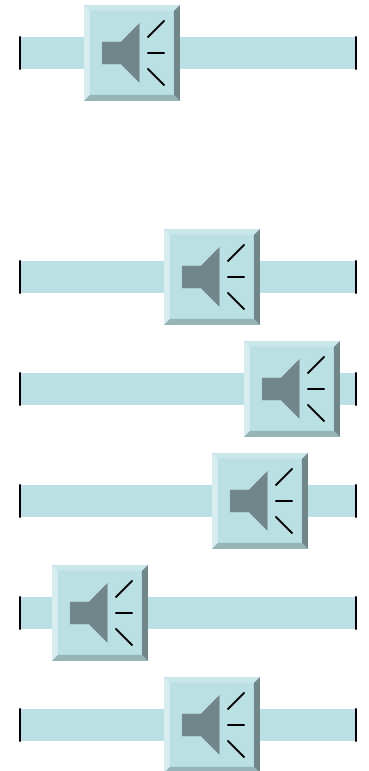
Features and Tuning

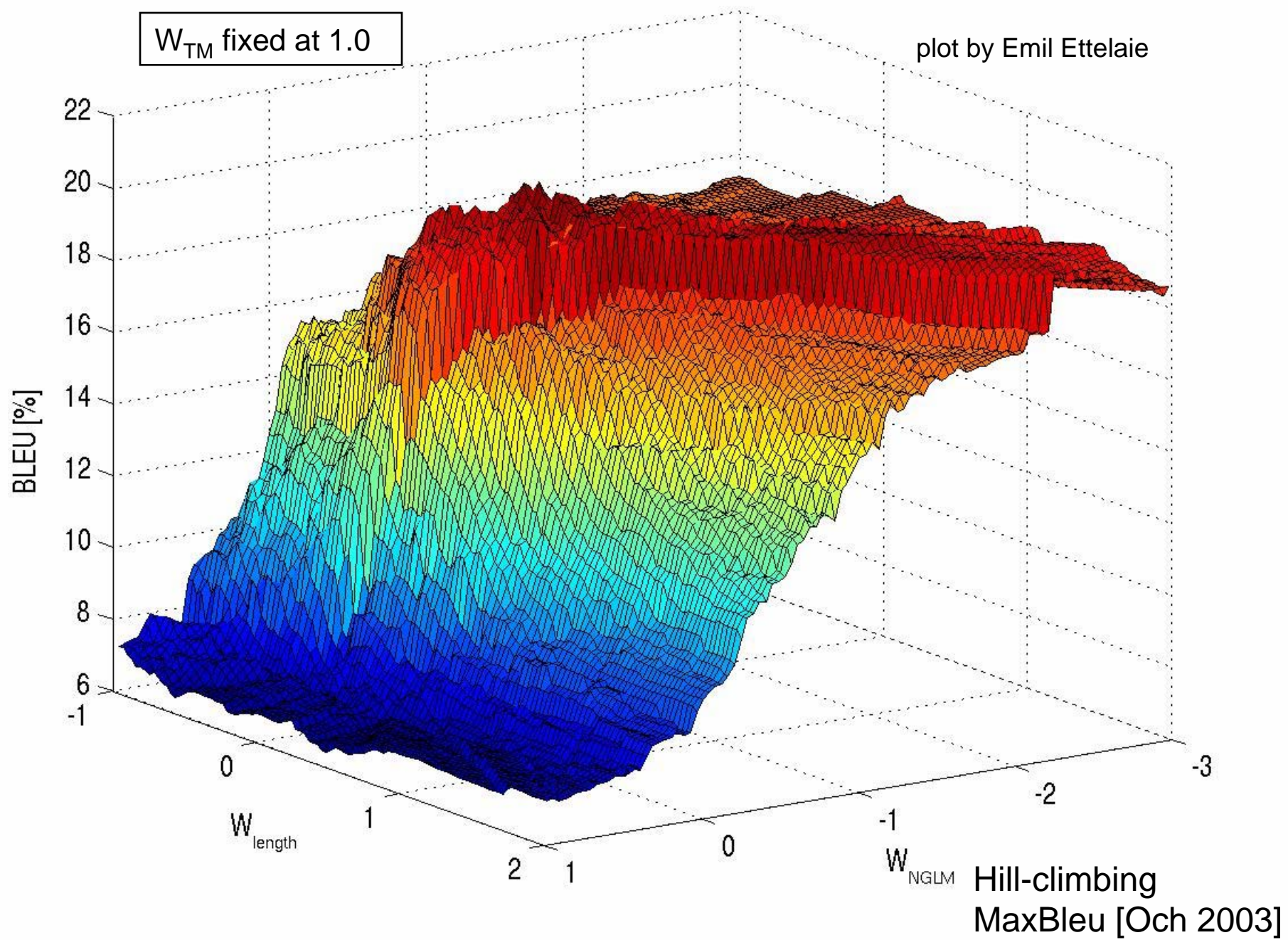
- English trigram language model
- Phrase pairs
 - Conditional probability
 - Bad-phrase spotter
 - Word-drop spotter
 - “Move Me” preference
- English output length



Features and Tuning

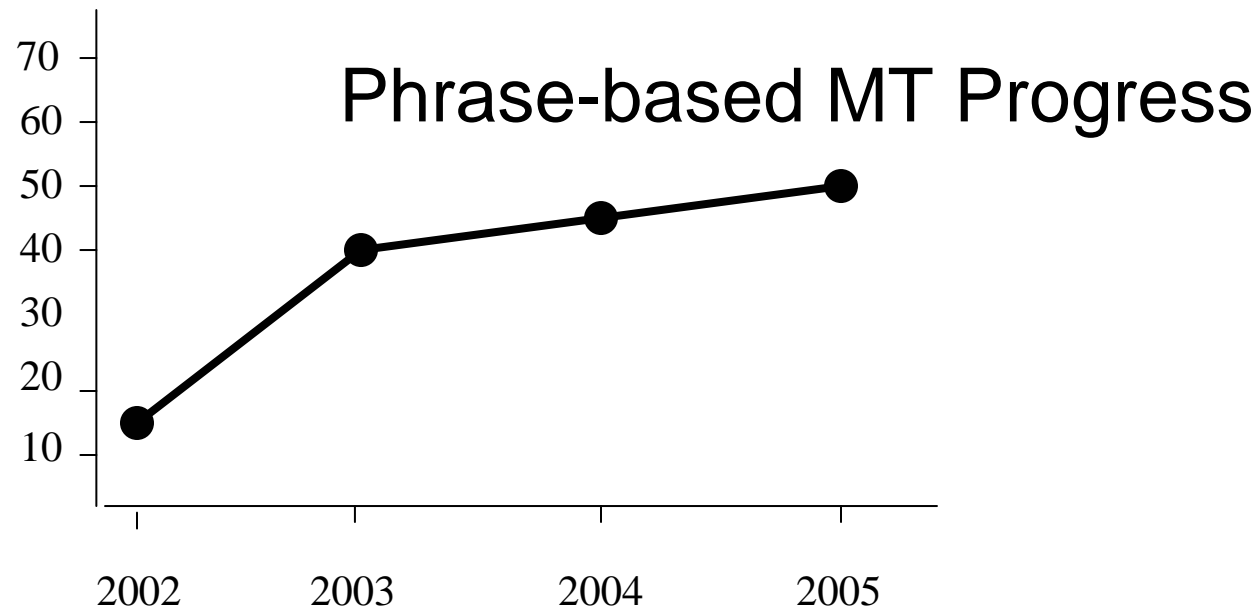
- English trigram language model
- Phrase pairs
 - Conditional probability
 - Bad-phrase spotter
 - Word-drop spotter
 - “Move Me” preference
- English output length





These Ideas Work!

Translation Quality
(BLEU)



NIST Common Evaluations
(Arabic/English)

Can a machine translate between Chinese and English without knowing what a verb is?

- Of course
- But the output is frequently bad

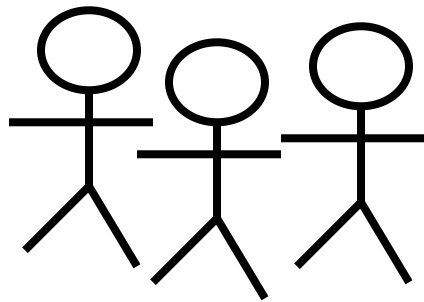
“Frequent high-tech exports are bright spots for foreign trade growth of Guangdong has made important contributions.”

- This phrase-based story is a little bit crazy

Syntax

Maybe we need some grammar?

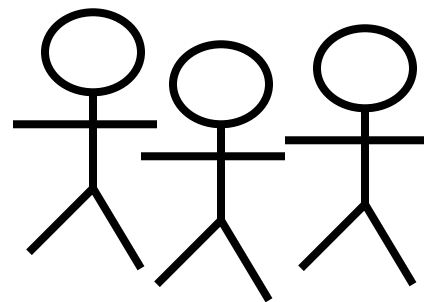
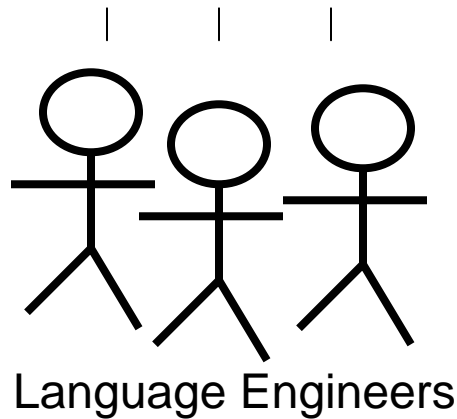
MT Research Landscape



Working on syntax-based approach
to translation (nouns, verbs,
prepositional phrases...)

MT Research Landscape

Syntax will never work!
We're better off without syntax!
Syntax has been *shown* to make things worse!
It has never worked in speech recognition!
You are crazy!

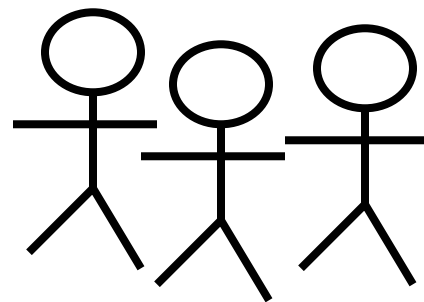
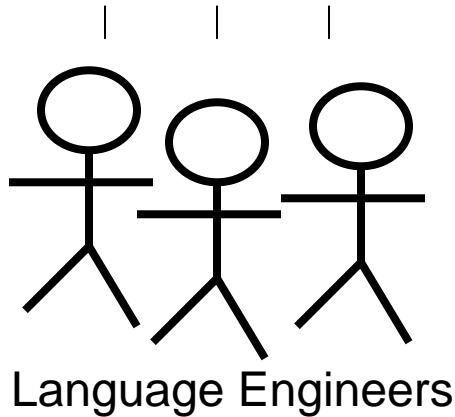


Working on syntax-based approach
to translation (nouns, verbs,
prepositional phrases...)

MT Research Landscape

Syntax will never work!
We're better off without syntax!
Syntax has been *shown* to make things worse!
It has never worked in speech recognition!
You are crazy!

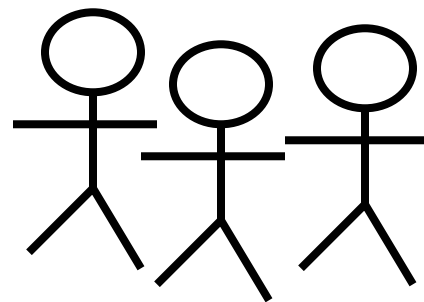
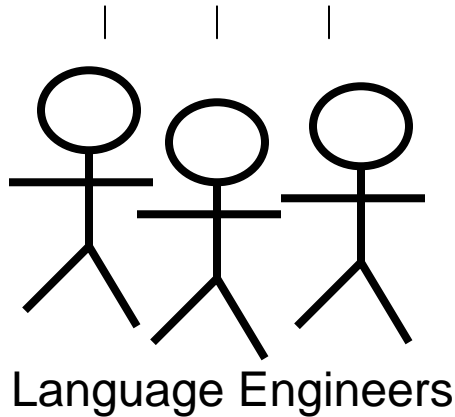
[Koehn et al 03]



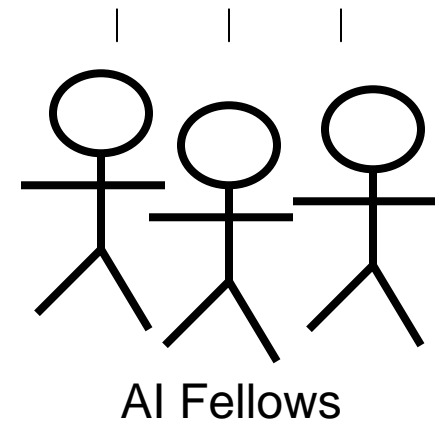
Working on syntax-based approach
to translation (nouns, verbs,
prepositional phrases...)

MT Research Landscape

Syntax will never work!
We're better off without syntax!
Syntax has been *shown* to make things worse!
It has never worked in speech recognition!
You are crazy!



Syntax will never work!
You need *semantics*!
Language is about the world!
You are crazy!



MT Research Landscape

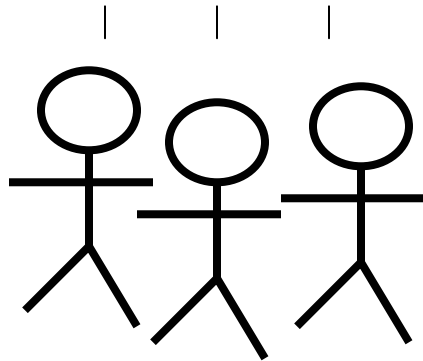
Syntax will never work!

We're better off without syntax!

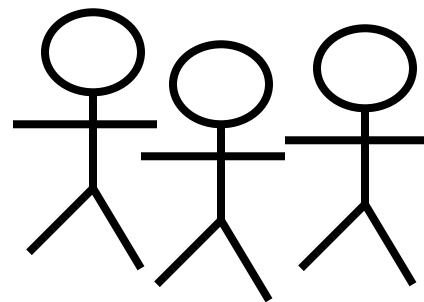
Syntax has been *shown* to make things worse!

It has never worked in speech recognition!

You are crazy!



Language Engineers



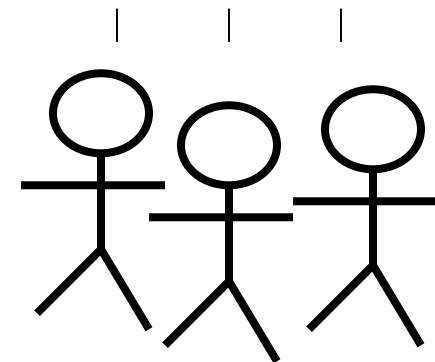
Working on syntax-based approach
to translation (nouns, verbs,
prepositional phrases...)

Syntax will never work!

You need *semantics*!

Language is about the world!

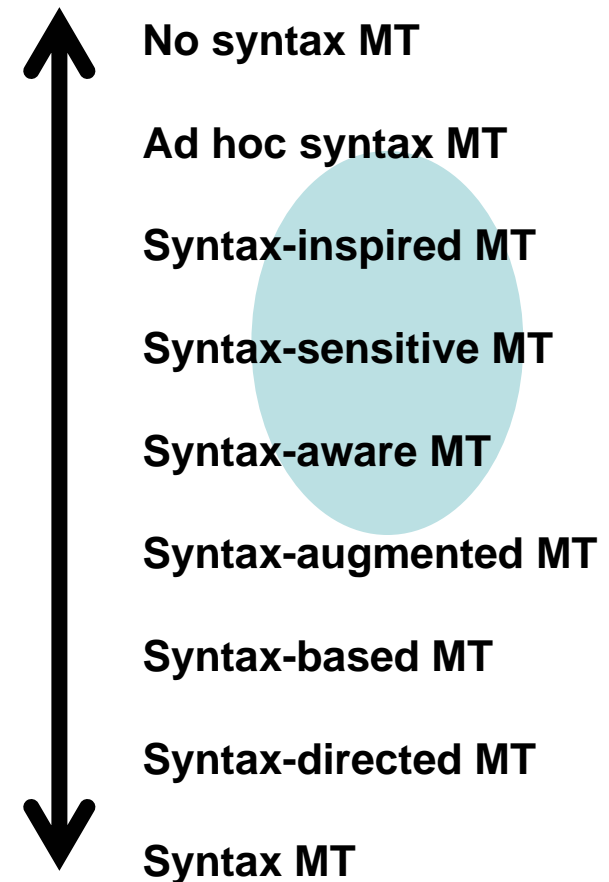
You are crazy!



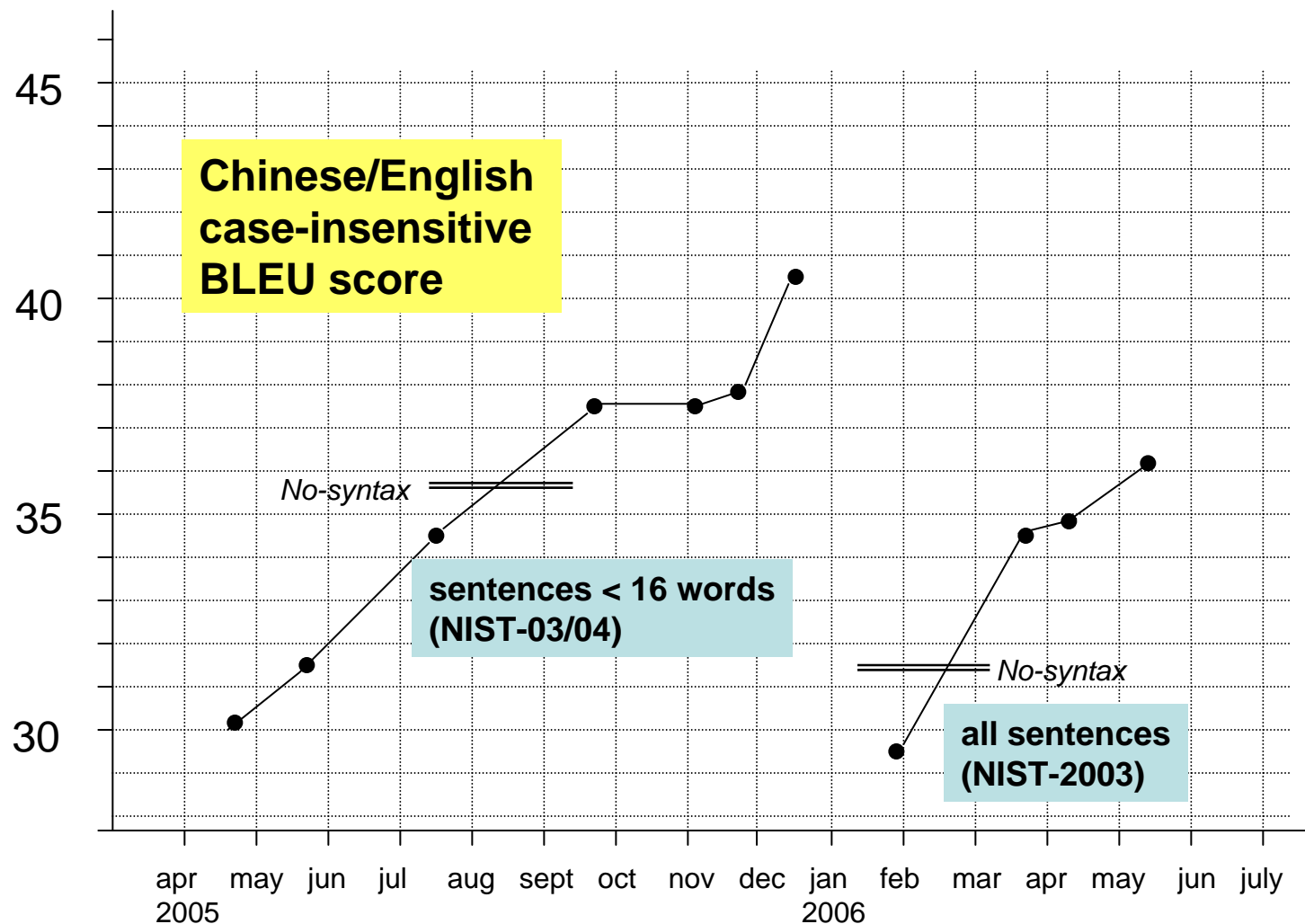
AI Fellows

NIST-2005 Common Evaluation of Machine Translation Systems

- Chinese/English
 - ISI No-Syntax system: 30.7
 - ISI Syntax system: 24.3
 - Google system: ~35
- Higher is better (not like golf)

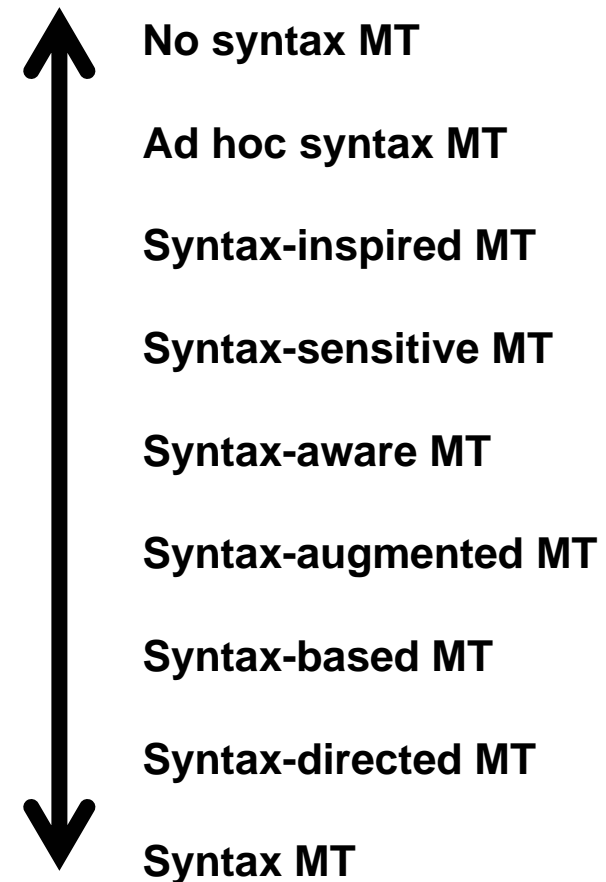


Syntax Started to Work in 2006...



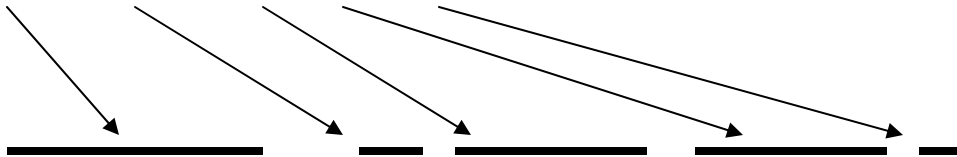
NIST-2006 Common Evaluation of Machine Translation Systems

- Chinese/English
 - ISI No-Syntax system: ~30
 - ISI Syntax system: 33.9
- Similar results for BN genre
- No-syntax = Syntax for BC genre and for Arabic/English
- Detailed testing with ASR just beginning



Phrase-Based Output

枪手 被 警方 击毙 .



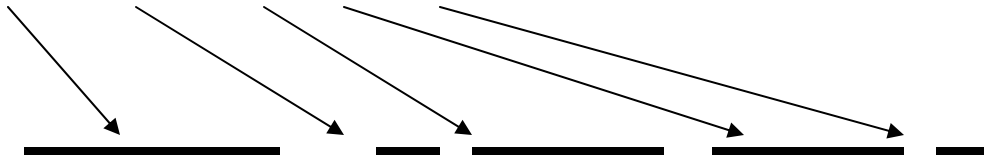
Gunman of police killed .

The diagram illustrates word alignment between the Chinese sentence "枪手 被 警方 击毙 ." and the English sentence "Gunman of police killed .". Arrows point from the Chinese words to the English words: "枪手" to "Gunman", "被" to "of", "警方" to "police", and "击毙" to "killed". The period "." in both sentences is aligned with the period in the English sentence.

*Decoder
Hypothesis #1*

Phrase-Based Output

枪手 被 警方 击毙 .



The diagram shows four arrows pointing from the Chinese characters to the English words in the hypothesis below:

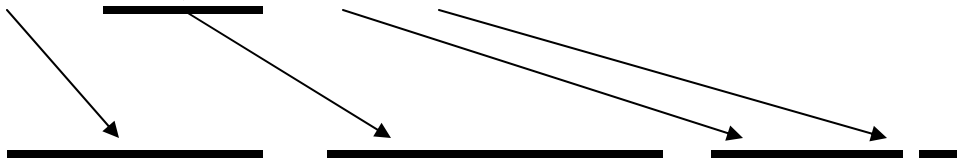
- 枪手 (Gunman) → 枪手
- 被 (of) → of
- 警方 (police) → police
- 击毙 (attack) → attack

Gunman of police attack .

*Decoder
Hypothesis #7*

Phrase-Based Output

枪手 被 警方 击毙 .

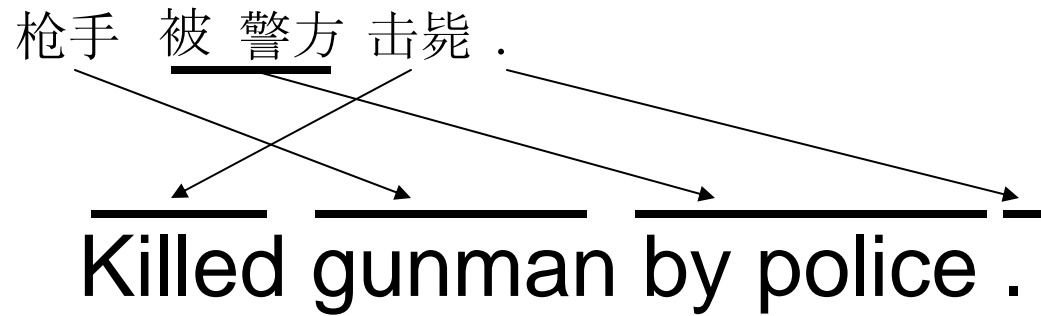


Gunman by police killed .

The diagram illustrates word alignment between the Chinese sentence "枪手 被 警方 击毙 ." and the English sentence "Gunman by police killed .". Arrows point from the Chinese words to the English words: "枪手" to "Gunman", "被" to "by", "警方" to "police", and "击毙" to "killed". The Chinese characters "被" and "警方" are underlined in the original image. The English words "Gunman", "by", "police", and "killed" are each underlined in the original image.

*Decoder
Hypothesis #12*

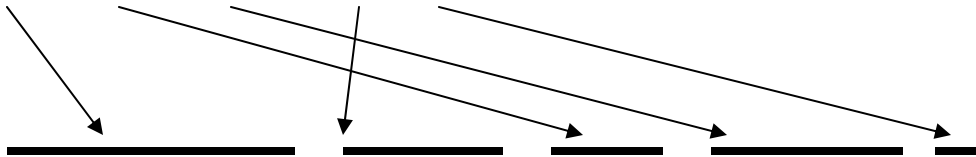
Phrase-Based Output



*Decoder
Hypothesis #134*

Phrase-Based Output

枪手 被 警方 击毙 .

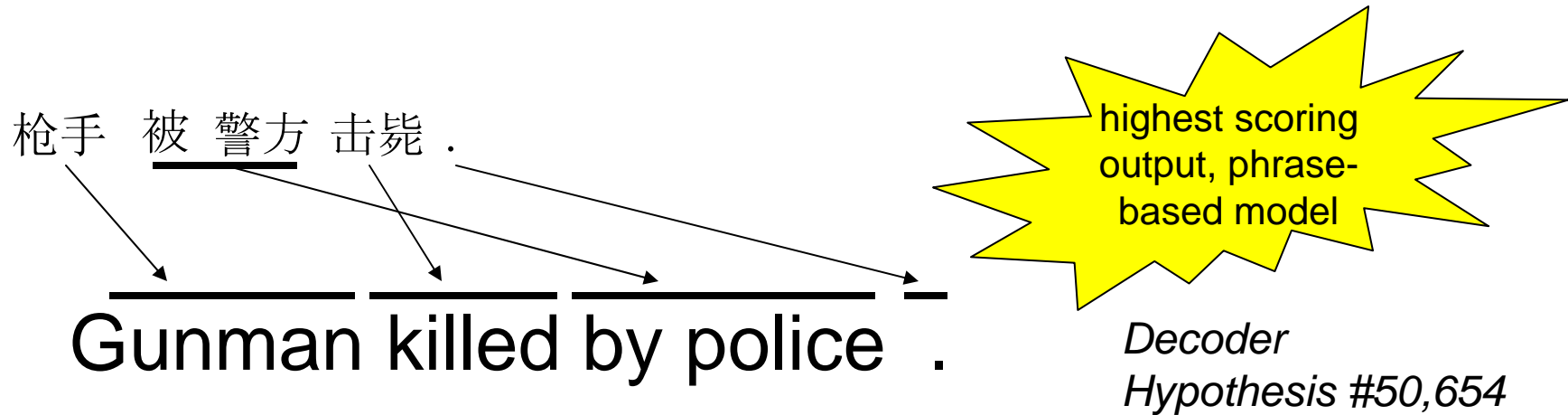


Gunman killed the police .

The diagram illustrates word alignment between the Chinese sentence "枪手 被 警方 击毙 ." and the English sentence "Gunman killed the police .". Arrows point from the Chinese words to the English words: "枪手" to "Gunman", "被" to "killed", "警方" to "the", and "击毙" to "police".

Decoder
Hypothesis #9,329

Phrase-Based Output



Problematic:

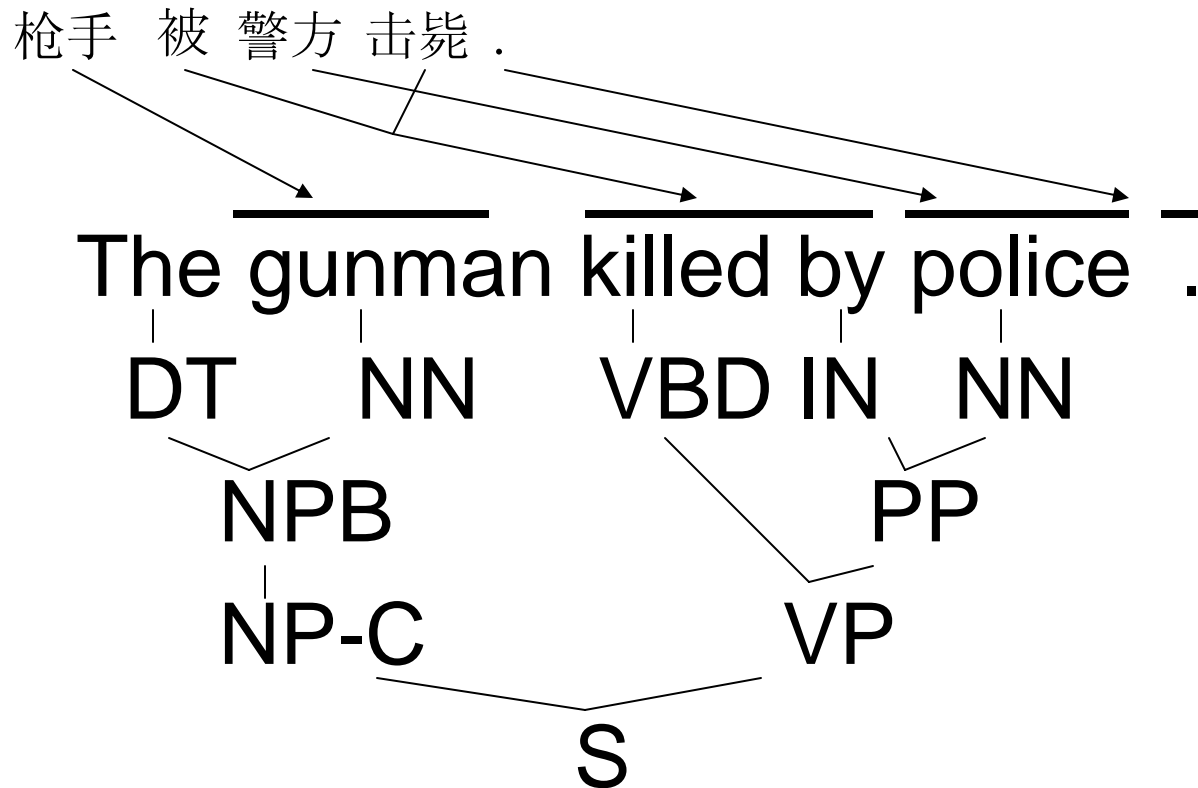
- VBD “killed” needs a direct object
- VBN “killed” needs an auxiliary verb (“was”)
- countable “gunman” needs an article (“a”, “the”)
- “passive marker” in Chinese controls re-ordering

Can't enforce/encourage any of this!

How to Get Grammar into the Statistical MT Picture?

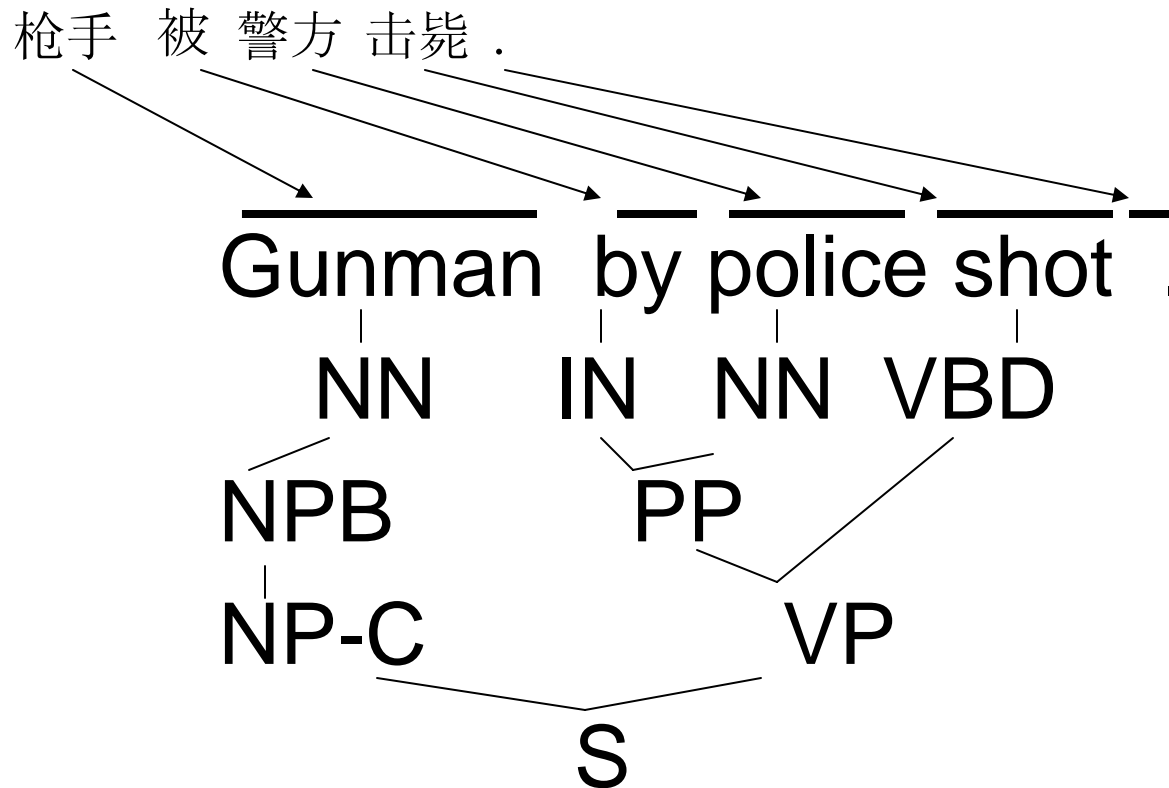
- Original work by Dekai Wu
- Yamada & Knight (2001, 2002)
- Galley, Hopkins, Marcu & Knight (2004)

Syntax-Based Output



*Decoder
Hypothesis #1*

Syntax-Based Output



*Decoder
Hypothesis #16*

Syntax-Based Output

枪手 被 警方 击毙 .

The gunman was killed by police .

*Decoder
Hypothesis #1923*

DT NN AUX VBN IN NN

NPB

PP

NP-C

VP

S

highest scoring
output, syntax-
based model

Syntax-Based Output

- Better modeling of target language structure
 - Always a verb
 - Verb is always in the right place
- Better handling of function words
 - They often don't translate
 - They control translation
- Better generalization in translation patterns

Why Target Trees Instead of Source Trees?

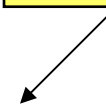
- Human translators need to know a lot more about the target language.
- MT system seems to know Chinese just fine.
 - Any evidence to the contrary?
- But the system does not know English!
 - Lots of evidence
- Speech input to MT
 - We don't have to parse source speech recognition
 - We can generalize to source lattices instead of strings

这 7人 中包括 来自 法国 和 俄罗斯 的 宇航 员 。

cstring



estring



These 7 people include astronauts coming from France and Russia .

这 7人 中包括 来自 法国 和 俄罗斯 的 宇航 员 .



cstring

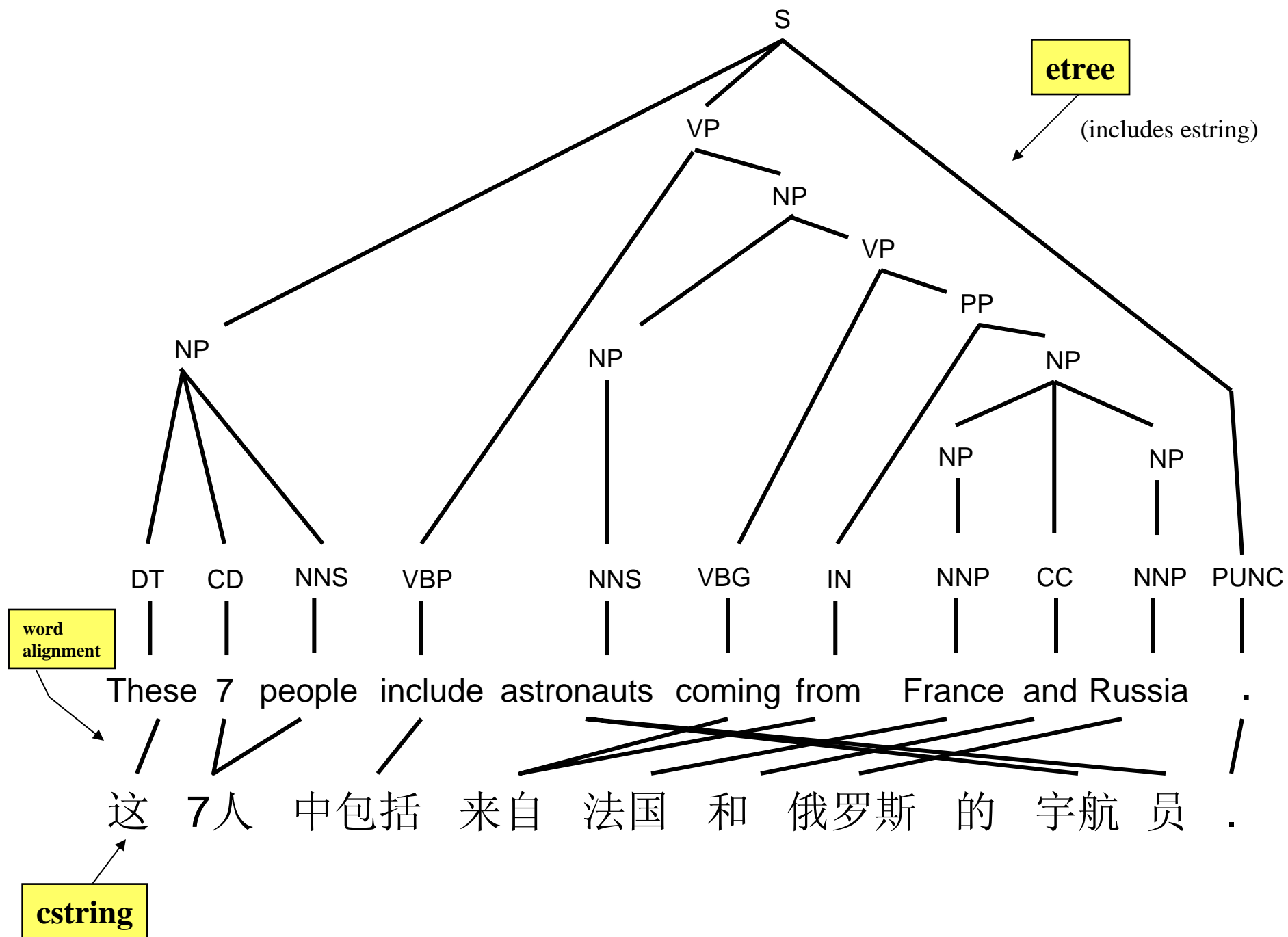
word alignment

estring

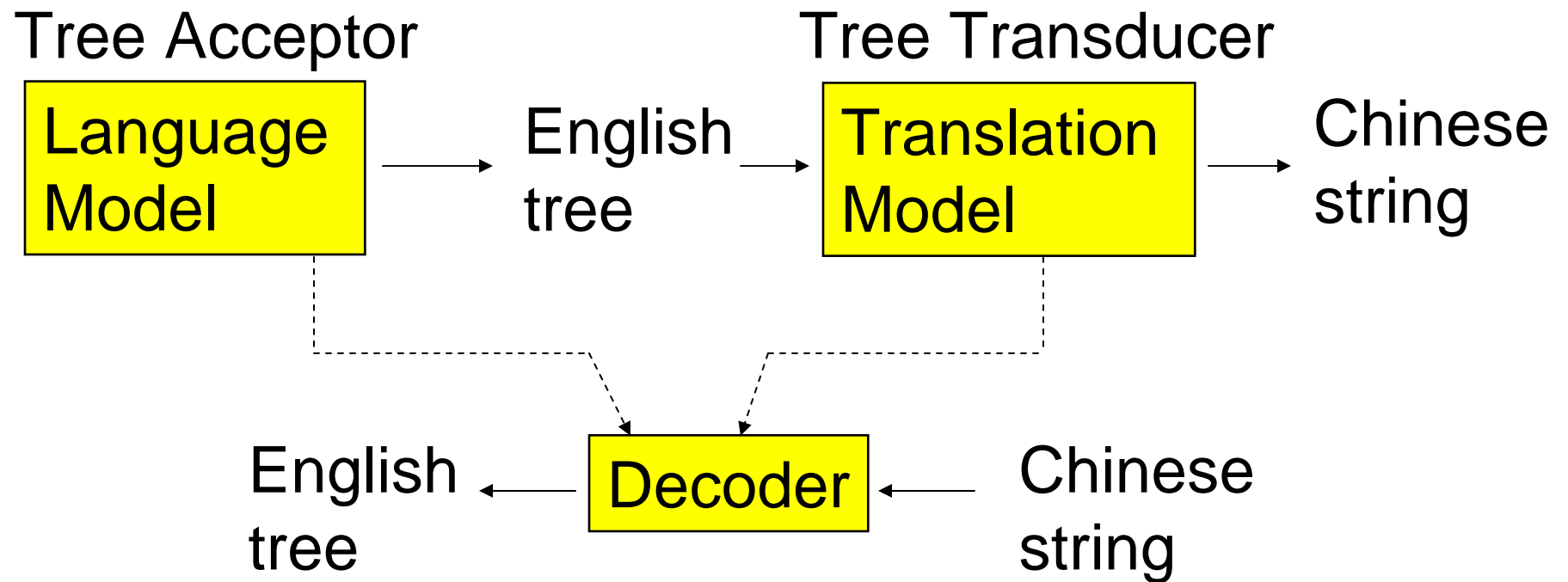
These 7 people include astronauts coming from France and Russia .

cstring

这 7人 中包括 来自 法国 和 俄罗斯 的 宇航员 .

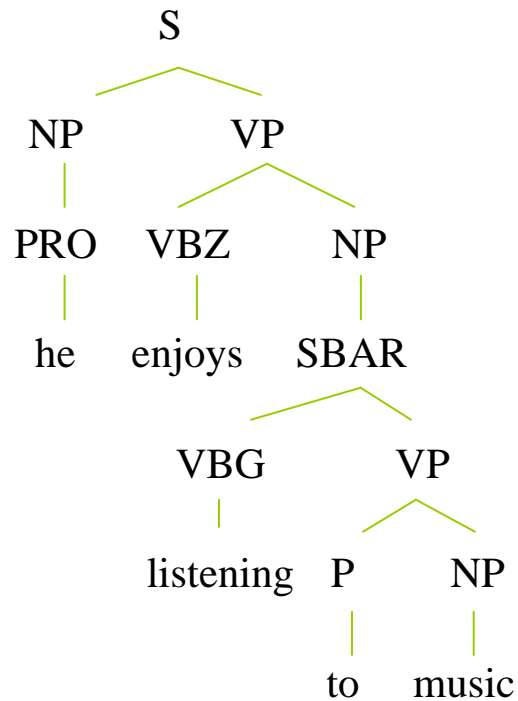


Syntax-Based Noisy Channel

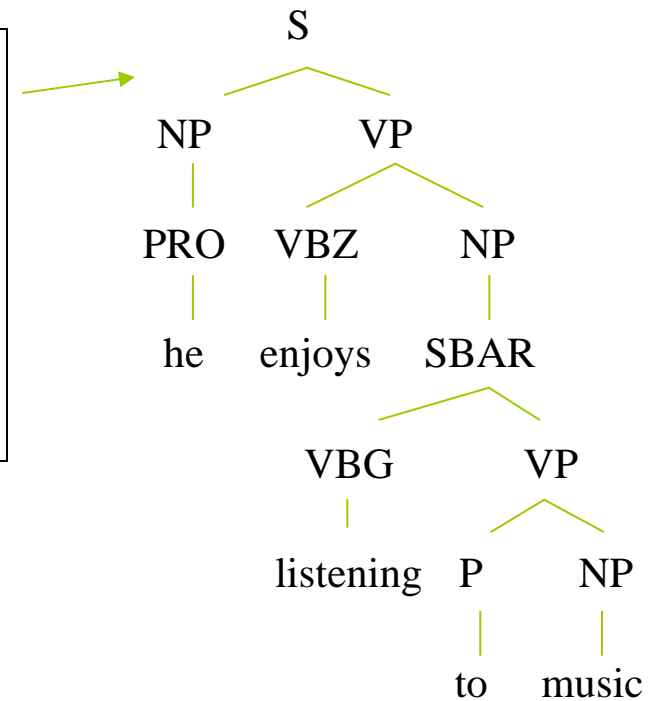
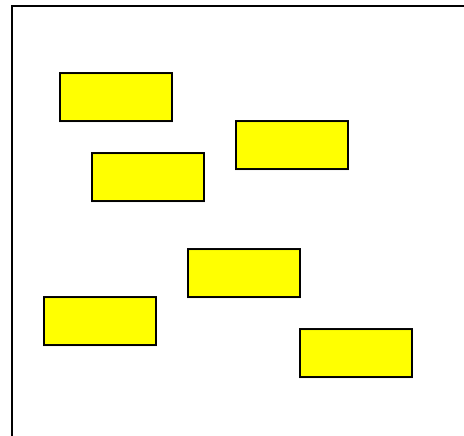


Top-Down Tree-to-String Transducer

Original input:

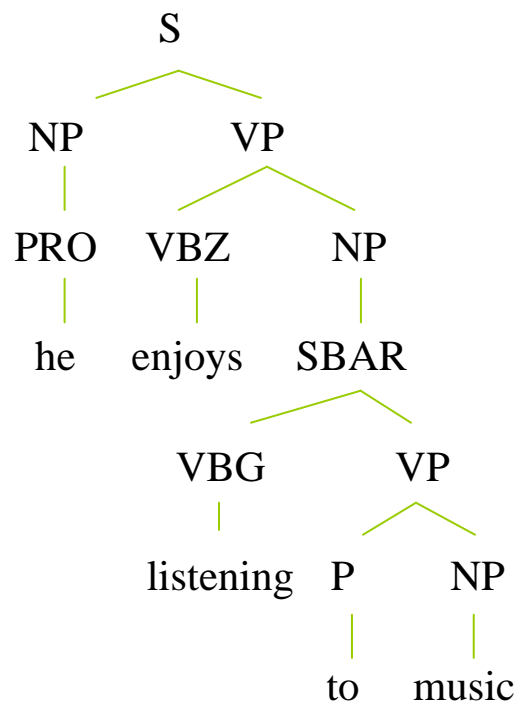


Transformation:

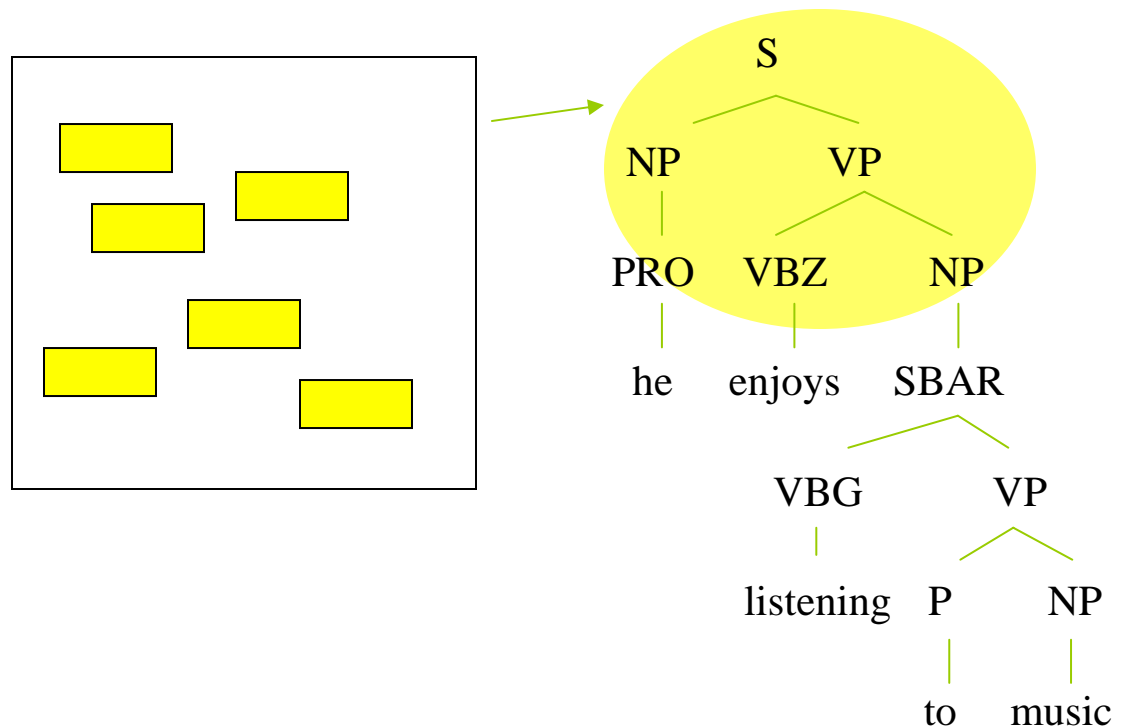


Top-Down Tree-to-String Transducer

Original input:

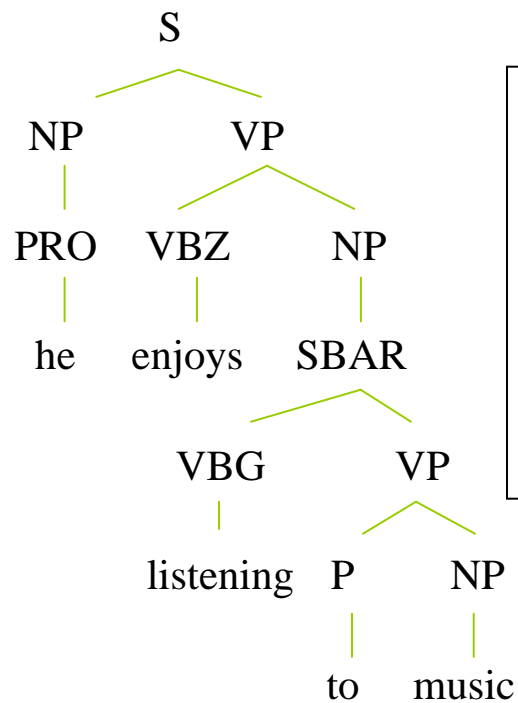


Transformation:

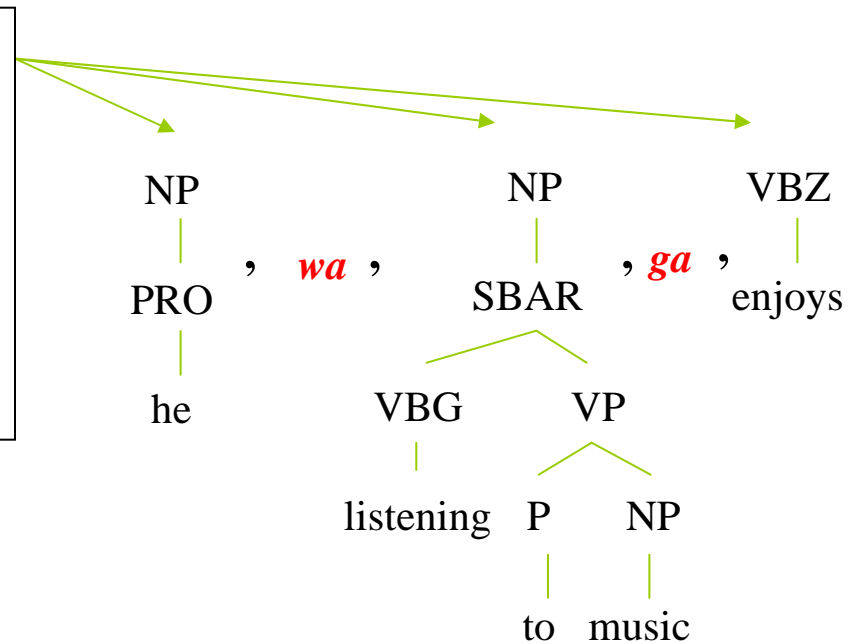


Top-Down Tree-to-String Transducer

Original input:

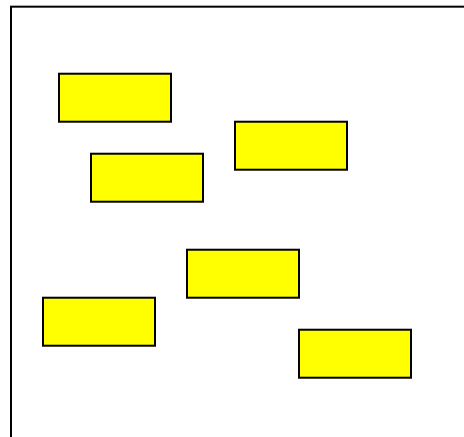
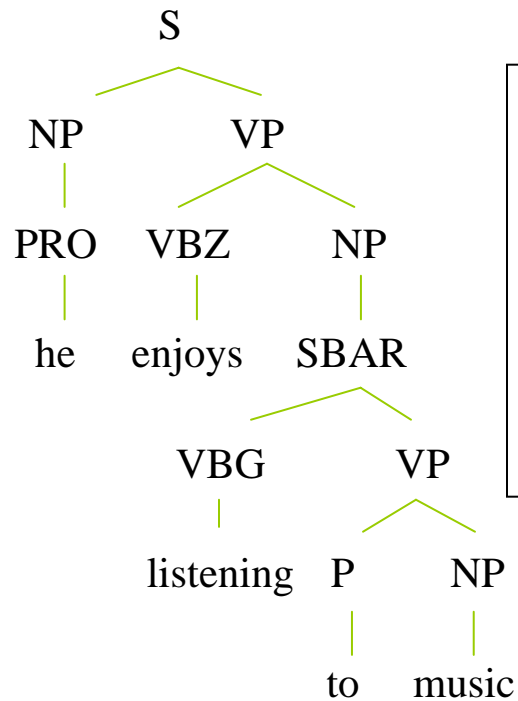


Transformation:

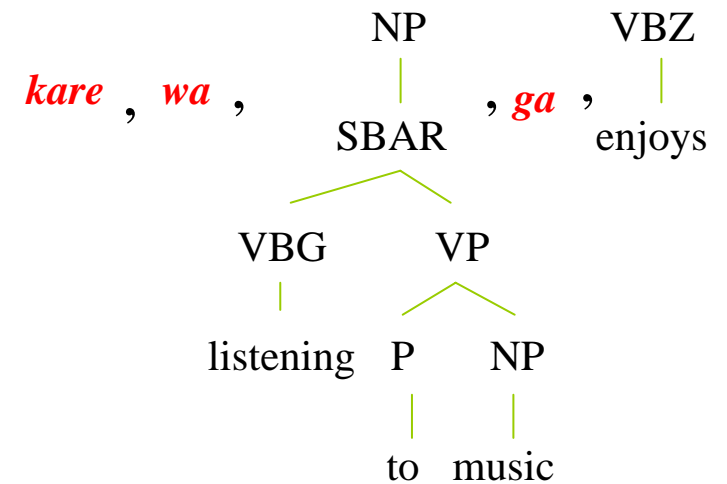


Top-Down Tree-to-String Transducer

Original input:



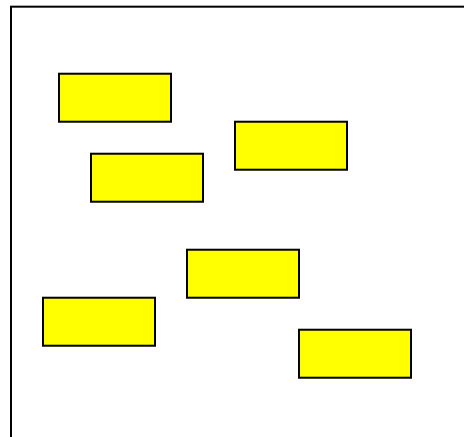
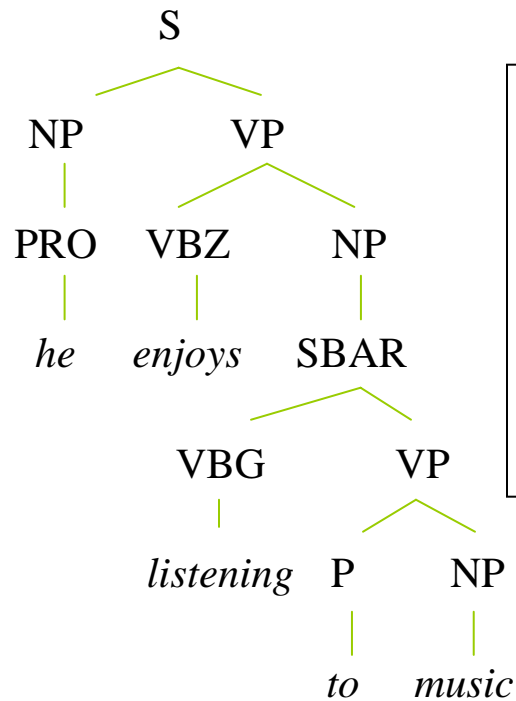
Transformation:



Top-Down Tree-to-String Transducer

Original input:

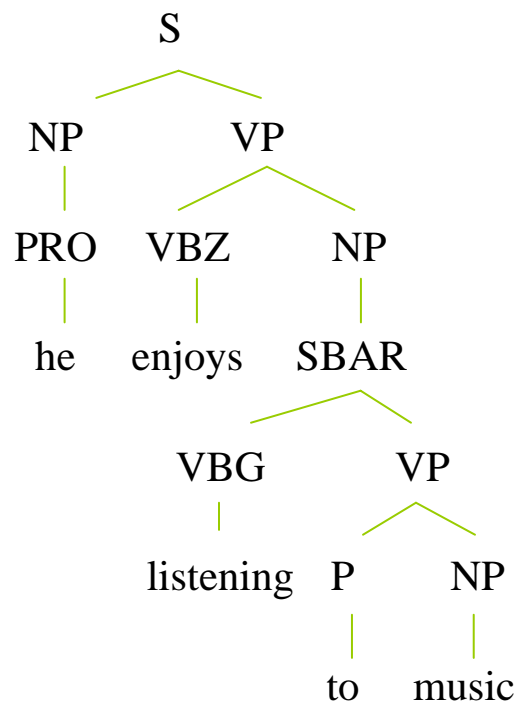
Final output:



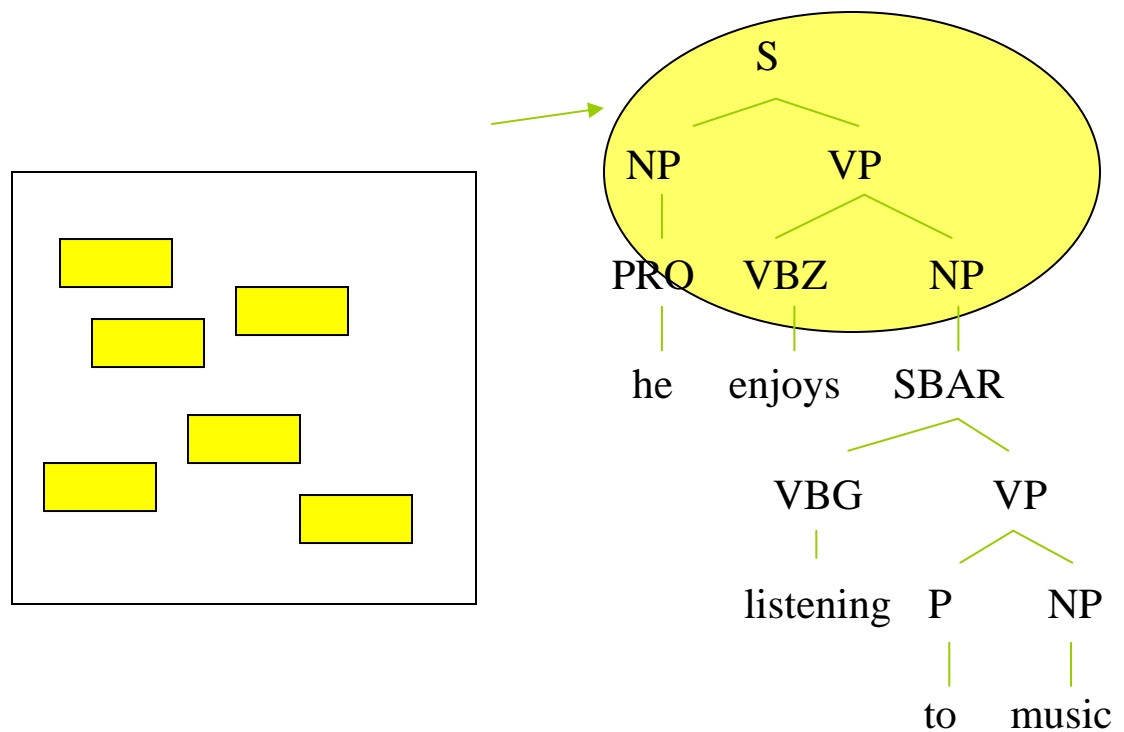
kare , wa , ongaku , o , kiku , no , ga , daisuki , desu

Top-Down Tree-to-String Transducer

Original input:

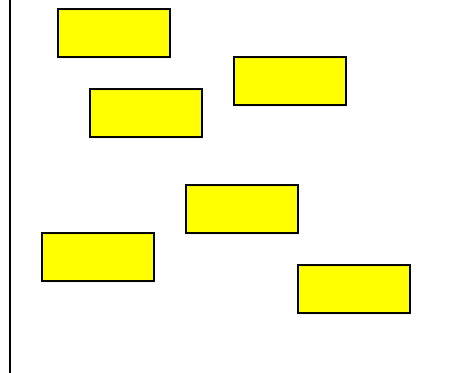
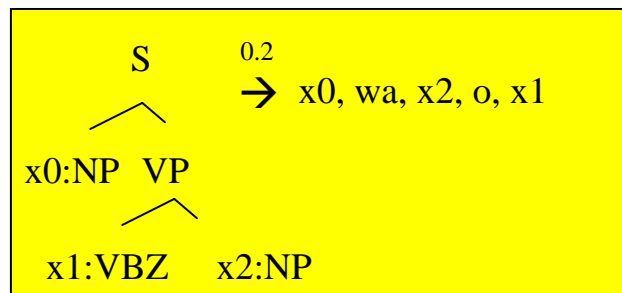
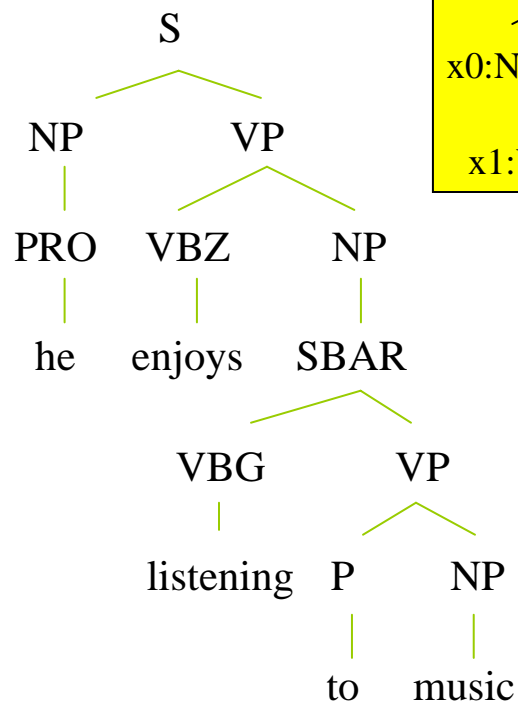


Transformation:

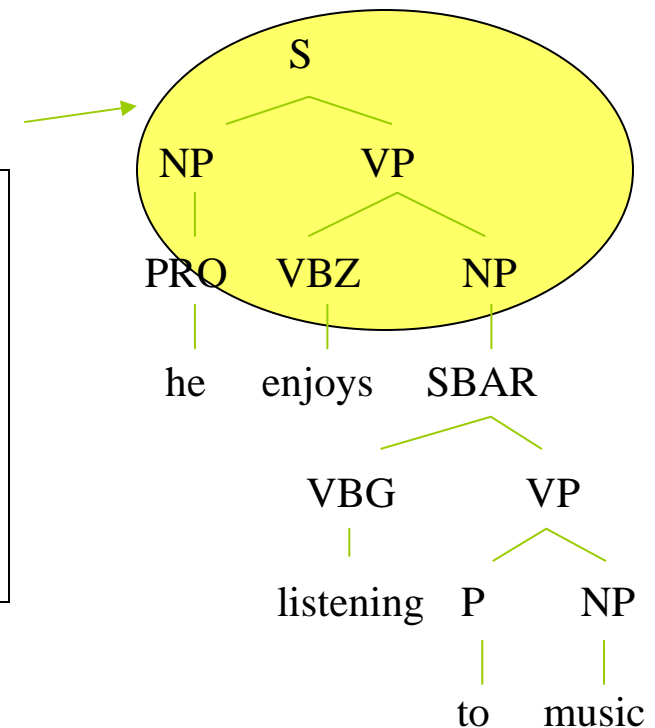


Top-Down Tree-to-String Transducer

Original input:



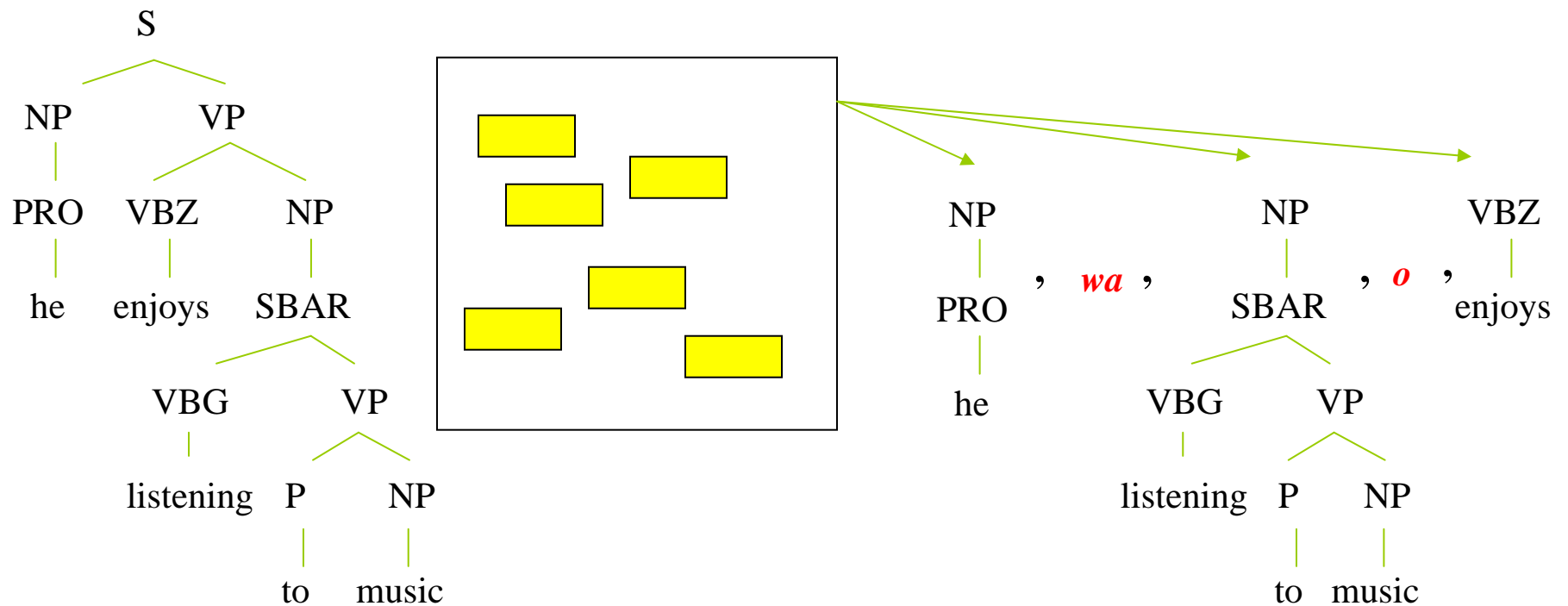
Transformation:



Top-Down Tree-to-String Transducer

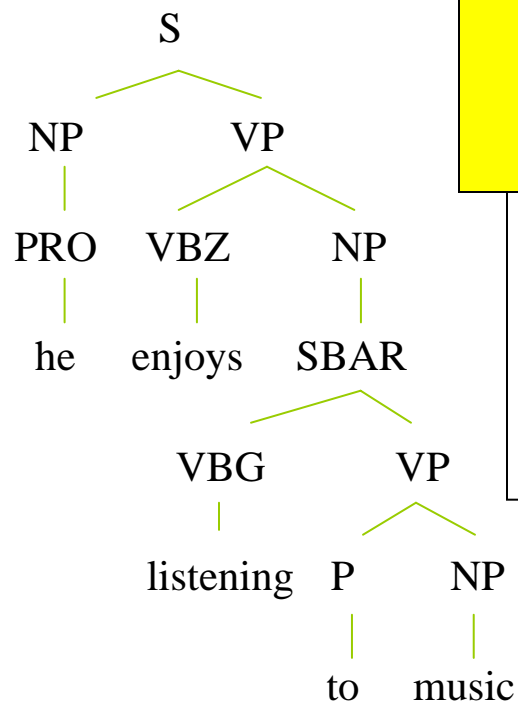
Original input:

Transformation:

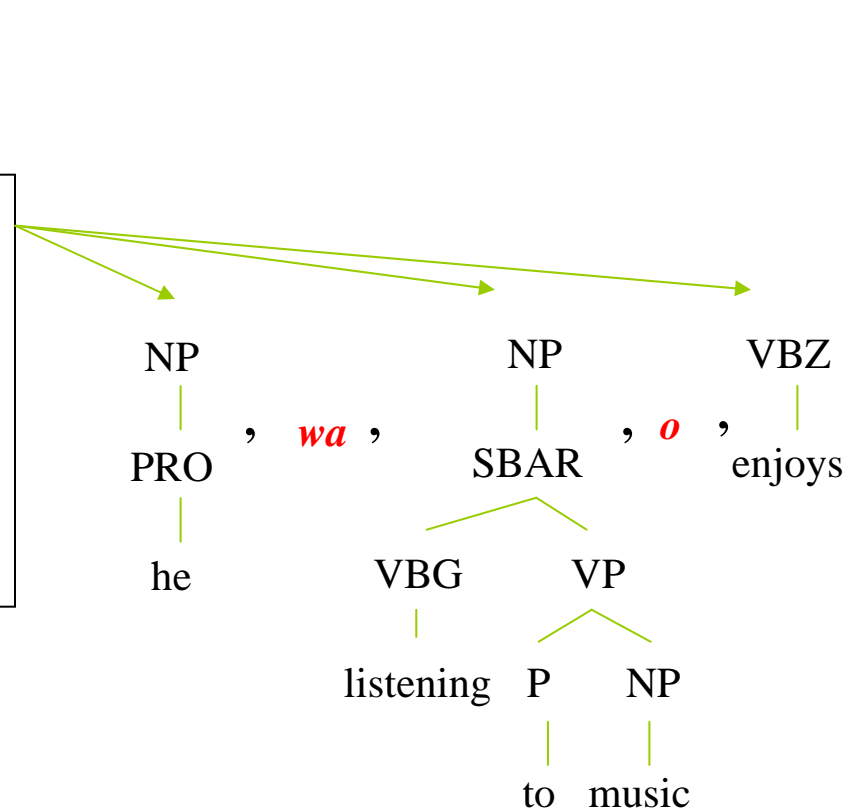


Top-Down Tree-to-String Transducer

Original input:

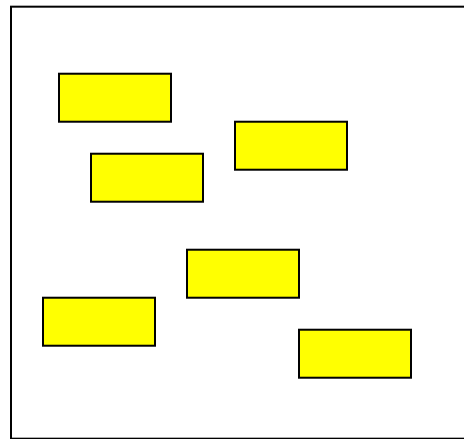
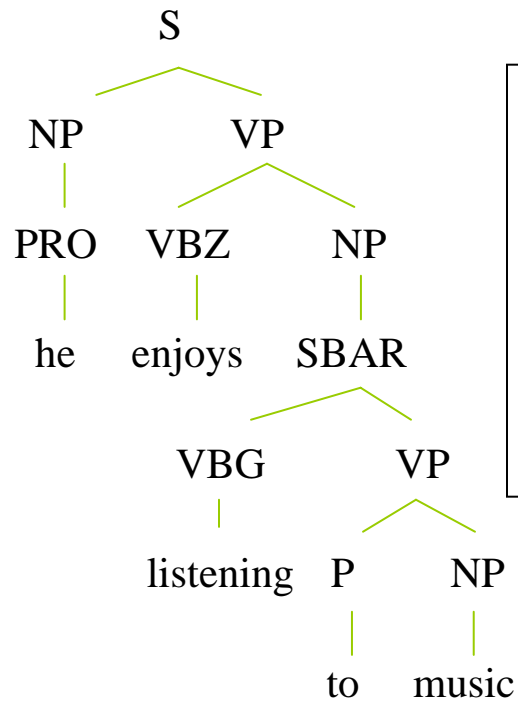


Transformation:

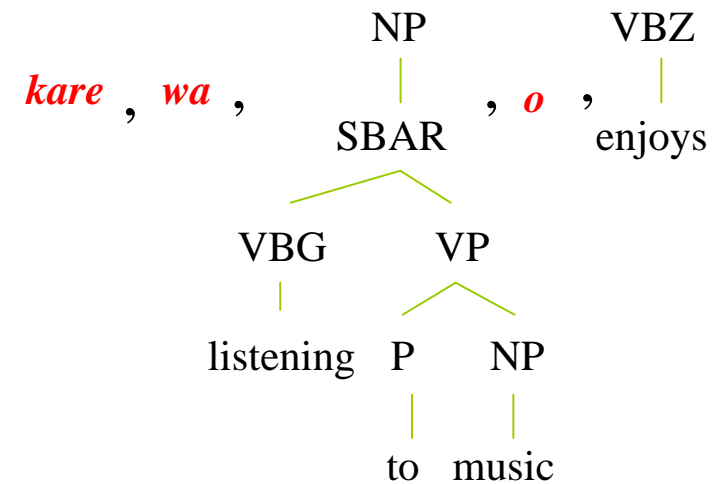


Top-Down Tree-to-String Transducer

Original input:

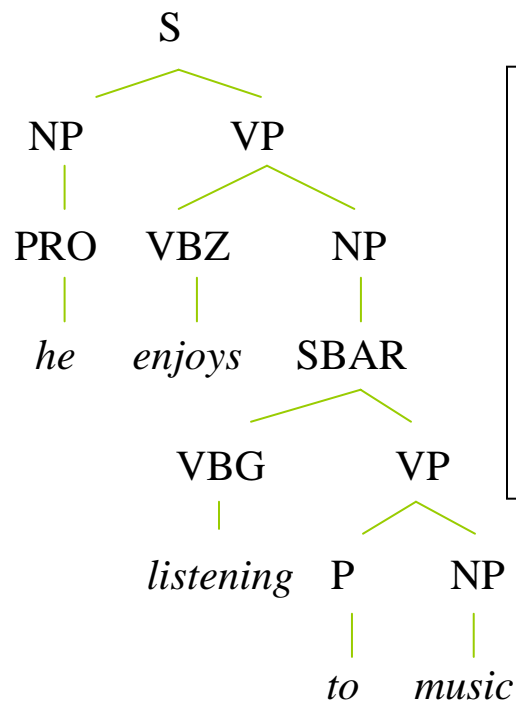


Transformation:

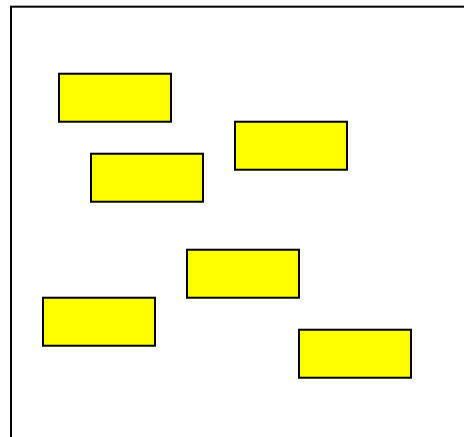


Top-Down Tree-to-String Transducer

Original input:



Final output:

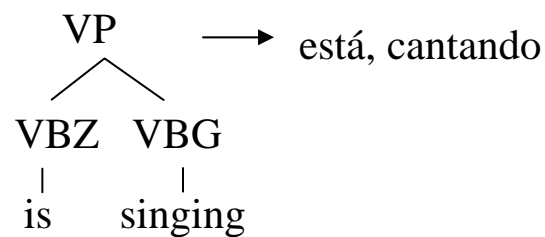


kare , wa , ongaku , o , kiku , no , ga , daisuki , desu

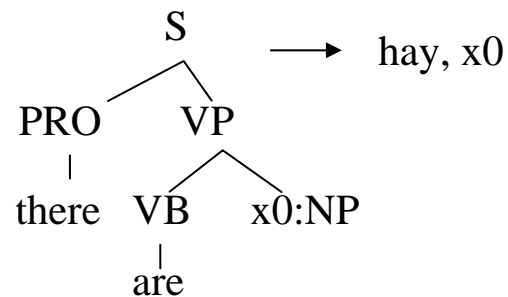
To get total probability,
multiply probabilities of the
individual steps.

Transducer Format is Expressive

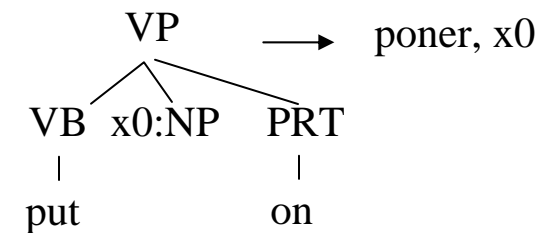
Phrasal Translation



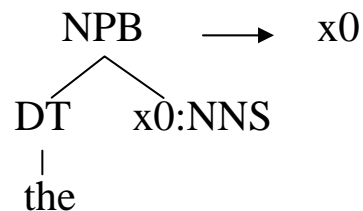
Non-constituent Phrases



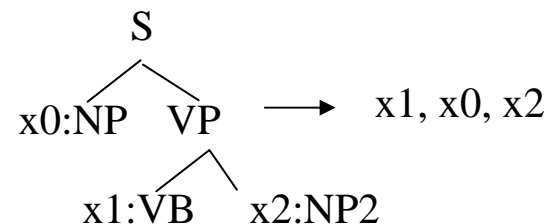
Non-contiguous Phrases



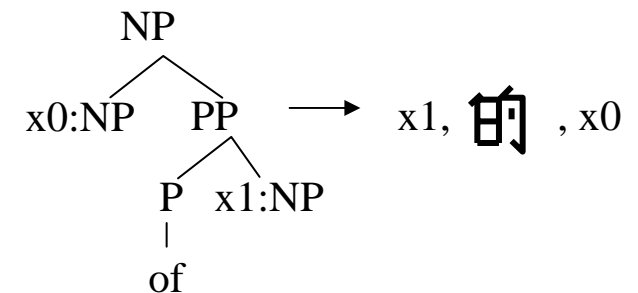
Context-Sensitive Word Insertion



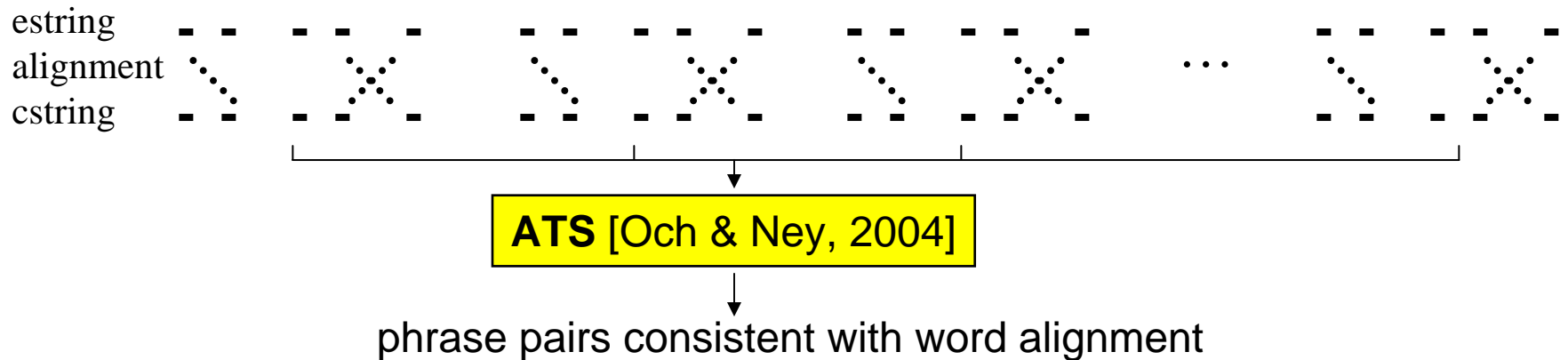
Multilevel Re-Ordering



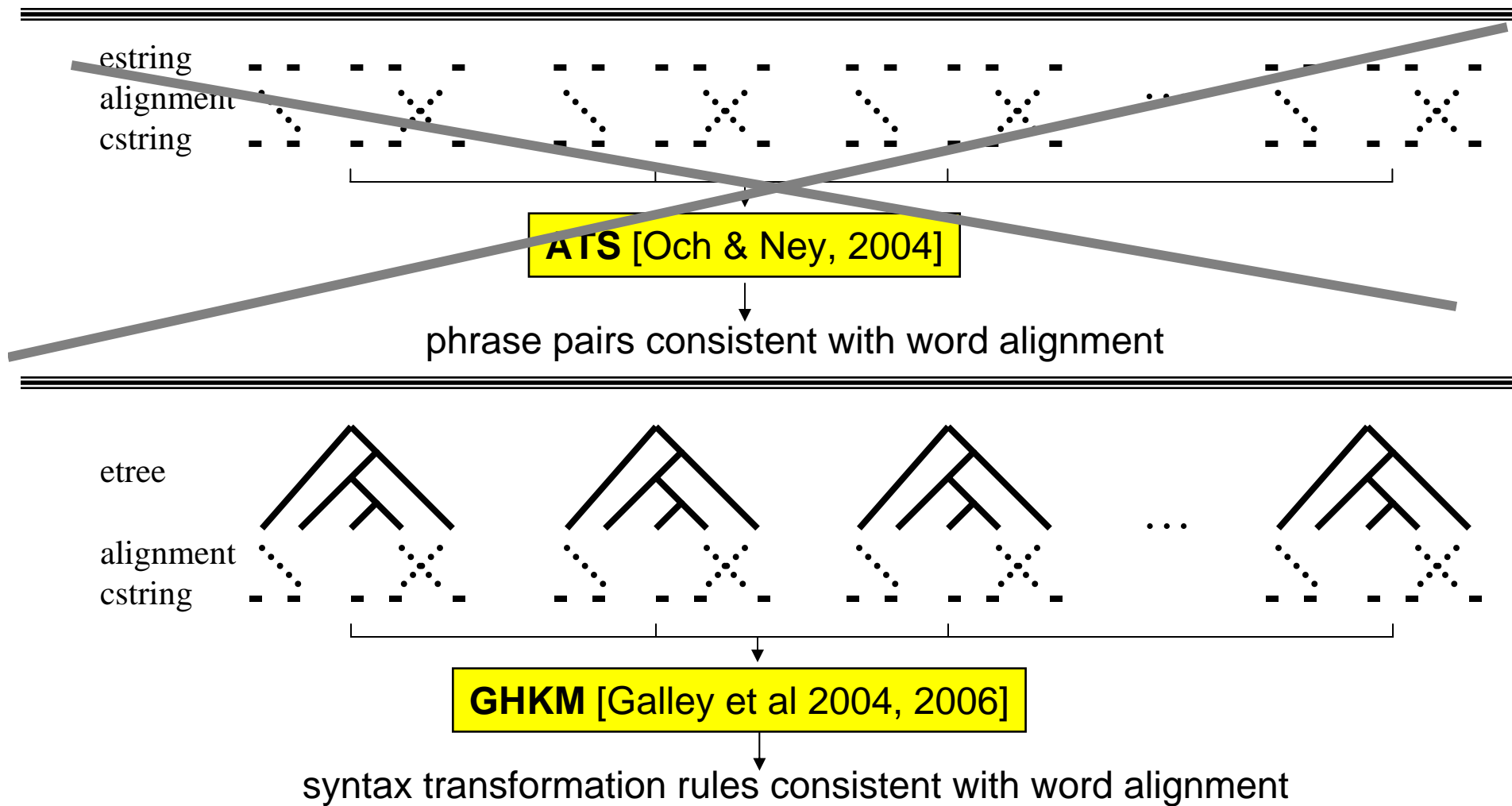
Lexicalized Re-Ordering



Phrase-Based and Syntax-Based Pattern Extraction

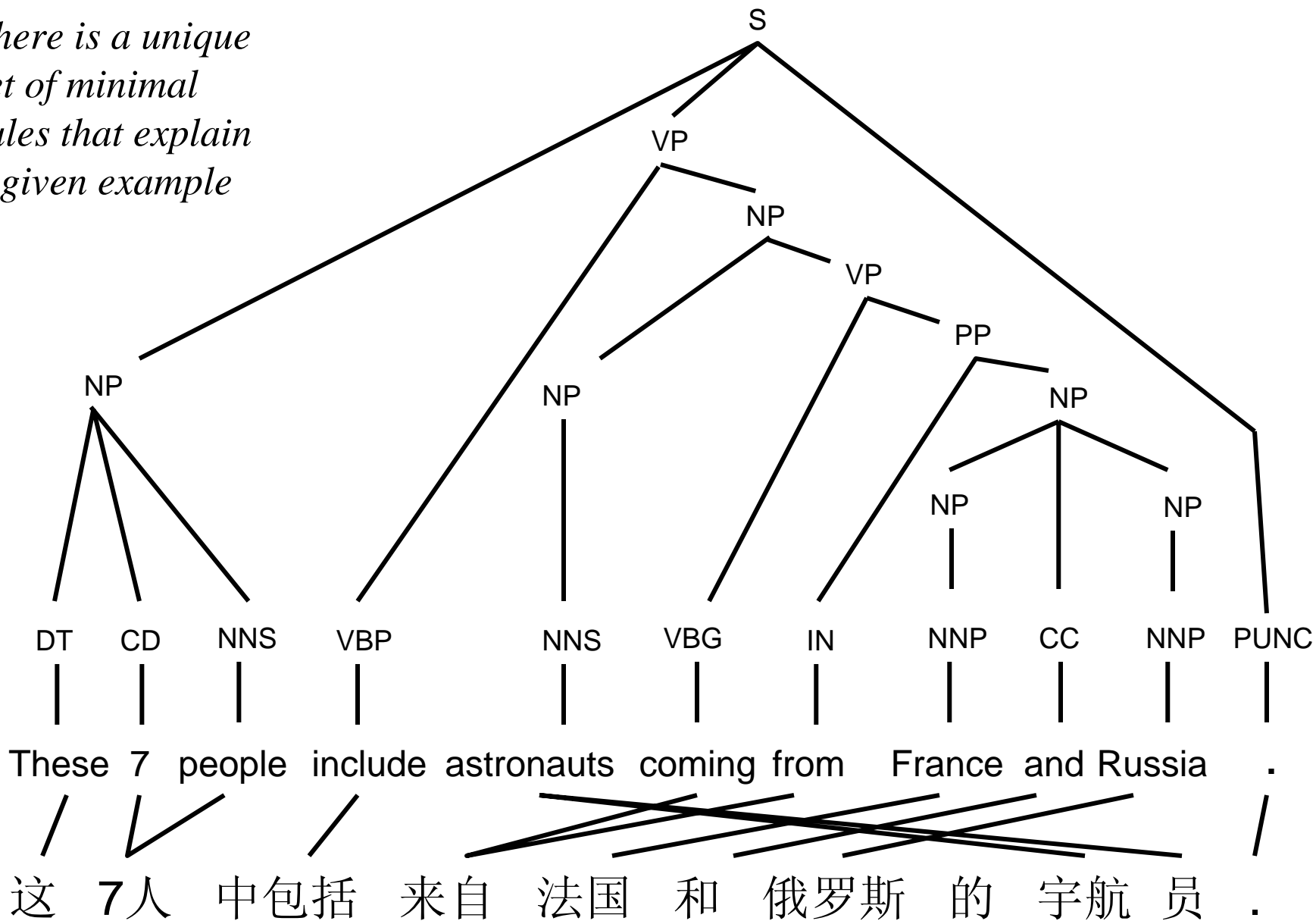


Phrase-Based and Syntax-Based Pattern Extraction



GHKM

*There is a unique
set of minimal
rules that explain
a given example*

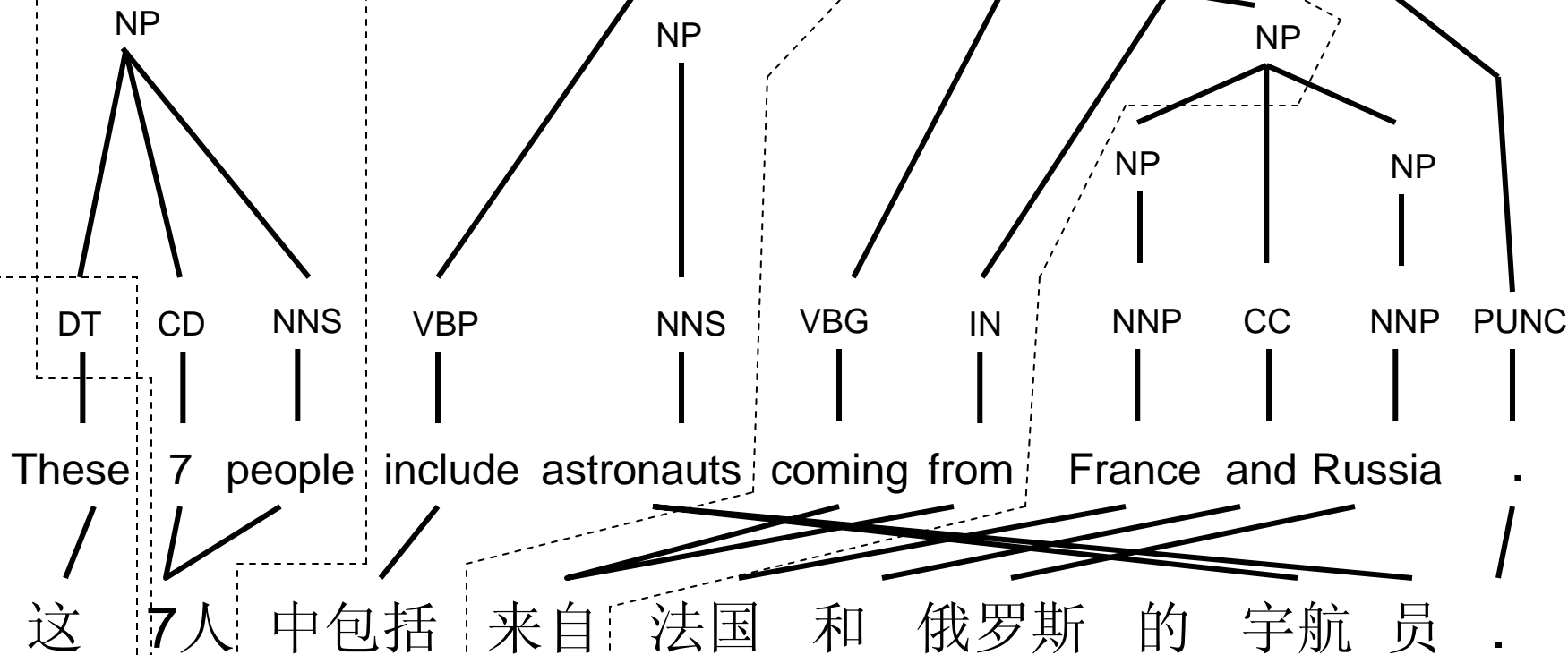


GHKM

*There is a unique
set of minimal
rules that explain
a given example*

VP(VBG(coming), PP(IN(from), x0:NP)) → 来自, x0

NP(x0:DT, CD(7), NNS(people)) → x0, 7人

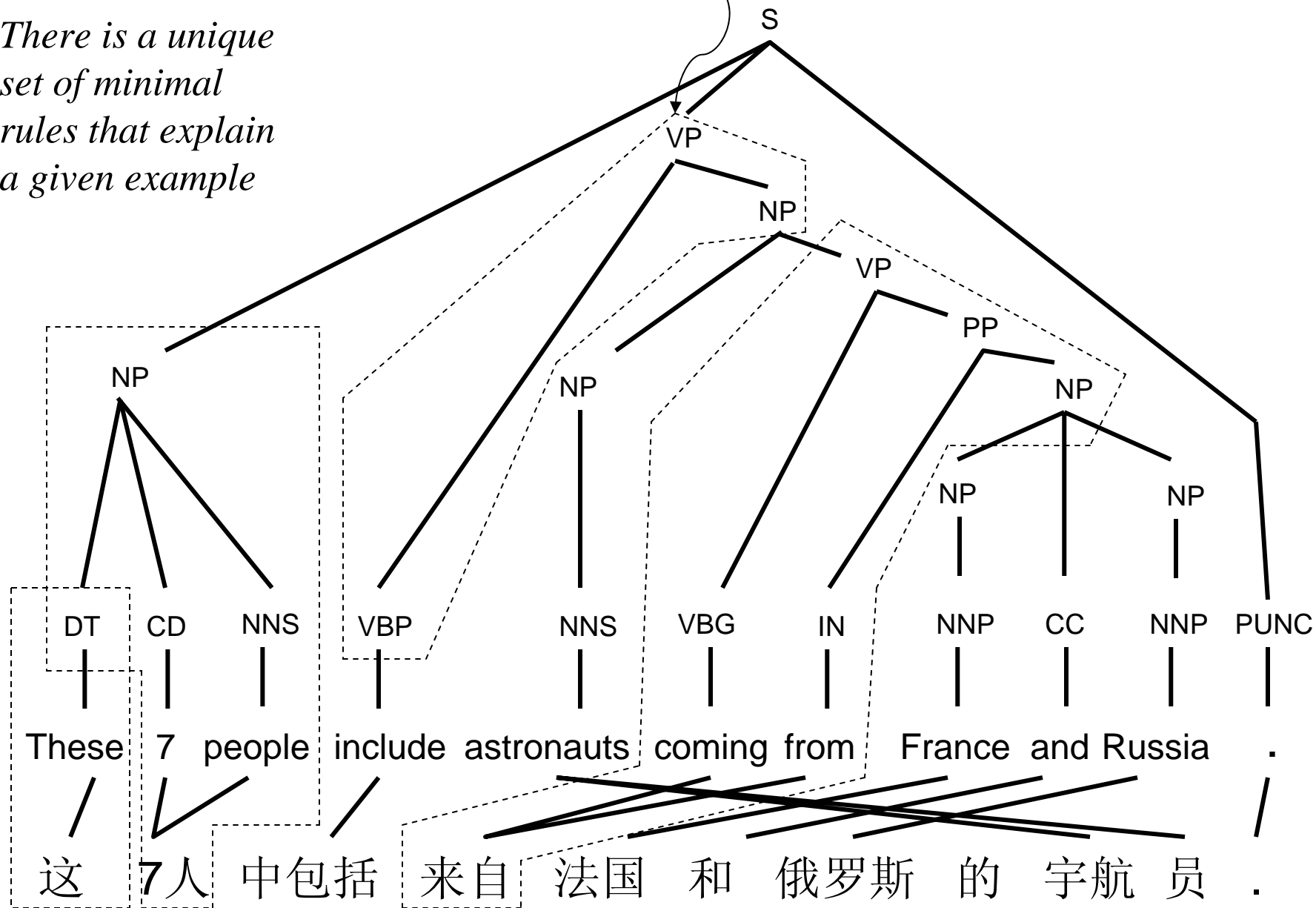


DT(these) → 这

GHKM

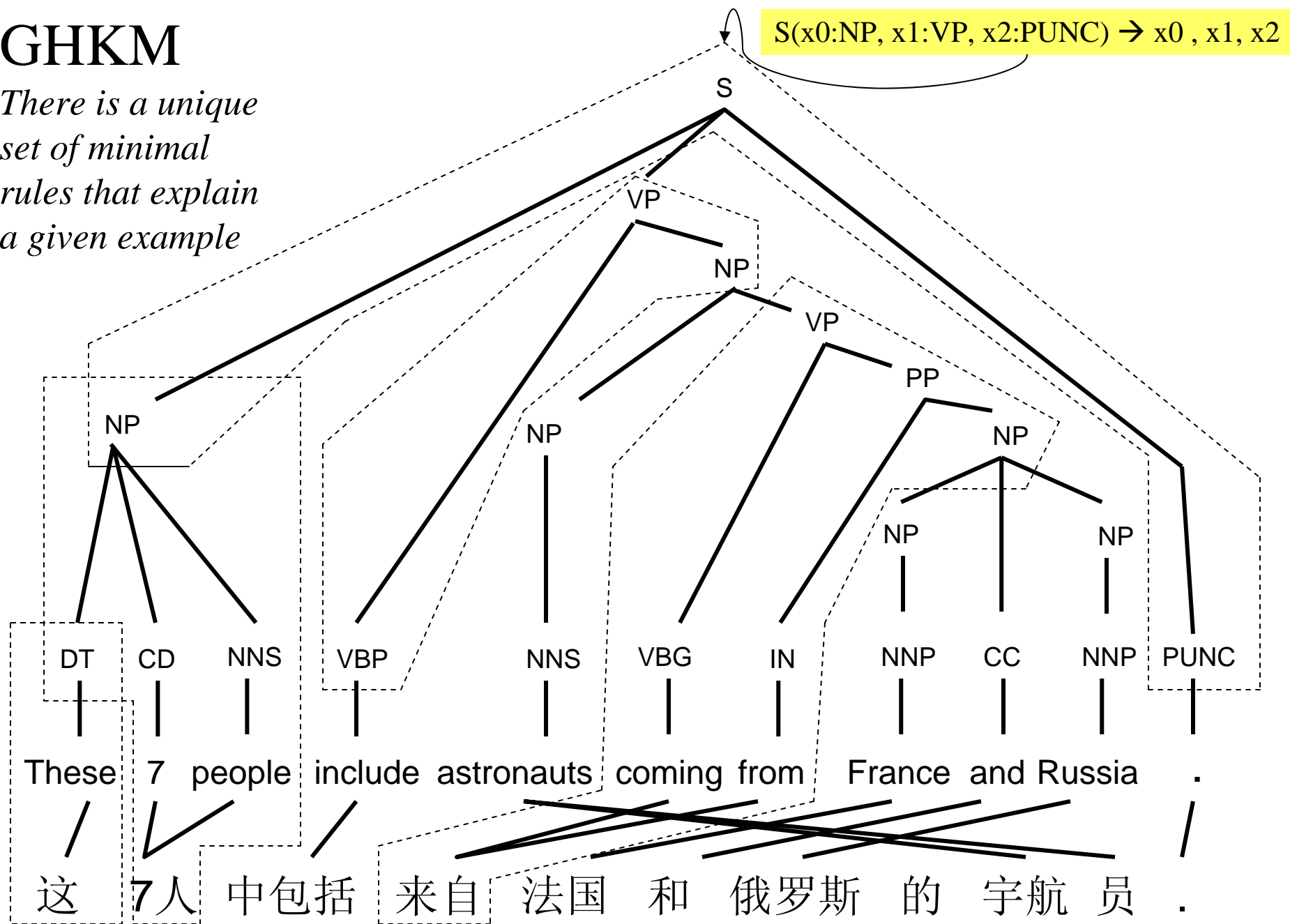
$VP(x_0:VBP, x_1:NP) \rightarrow x_0, x_1$

*There is a unique
set of minimal
rules that explain
a given example*

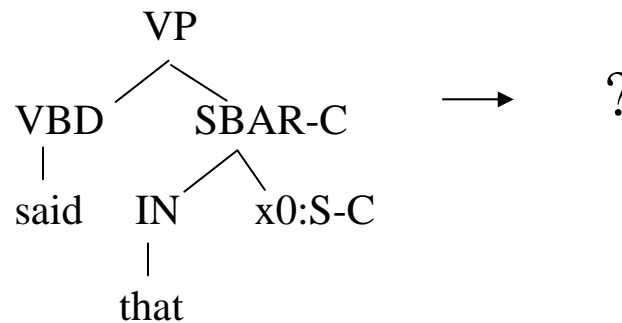


GHKM

*There is a unique
set of minimal
rules that explain
a given example*

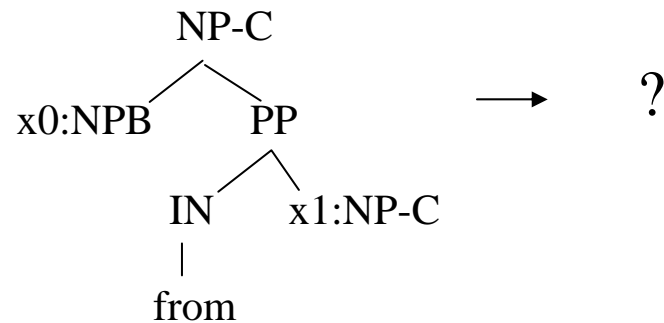


Sample “said that” rules



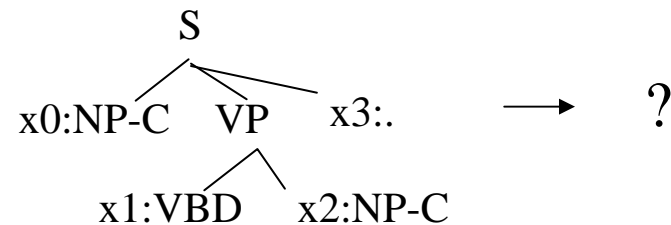
- 0.57 VP(VBD("said") SBAR-C(IN("that") x0:S-C)) -> "说" ", " x0
- 0.09 VP(VBD("said") SBAR-C(IN("that") x0:S-C)) -> "说" x0
- 0.02 VP(VBD("said") SBAR-C(IN("that") x0:S-C)) -> "他" "说" ", " x0
- 0.02 VP(VBD("said") SBAR-C(IN("that") x0:S-C)) -> "指出" ", " x0
- 0.02 VP(VBD("said") SBAR-C(IN("that") x0:S-C)) -> x0
- 0.01 VP(VBD("said") SBAR-C(IN("that") x0:S-C)) -> "表示" x0
- 0.01 VP(VBD("said") SBAR-C(IN("that") x0:S-C)) -> "说" ", " x0 "的"

Sample “NP-from-NP” rules



- 0.27 NP-C(x0:NPB PP(IN("from") x1:NP-C)) -> x1 x0
- 0.15 NP-C(x0:NPB PP(IN("from") x1:NP-C)) -> "来自" x1 x0
- 0.06 NP-C(x0:NPB PP(IN("from") x1:NP-C)) -> x1 "的" x0
- 0.06 NP-C(x0:NPB PP(IN("from") x1:NP-C)) -> "从" x1 x0
- 0.06 NP-C(x0:NPB PP(IN("from") x1:NP-C)) -> "来自" x1 "的" x0
- 0.02 NP-C(x0:NPB PP(IN("from") x1:NP-C)) -> x0 "从" x1
- 0.01 NP-C(x0:NPB PP(IN("from") x1:NP-C)) -> "自" x1 x0
- 0.01 NP-C(x0:NPB PP(IN("from") x1:NP-C)) -> x1 x0 ", "

Sample SVO rules



CHINESE / ENGLISH

- 0.82 S(x0:NP-C VP(x1:VBD x2:NP-C) x3:.) -> x0 x1 x2 x3
- 0.02 S(x0:NP-C VP(x1:VBD x2:NP-C) x3:.) -> x0 x1 "," x2 x3
- 0.01 S(x0:NP-C VP(x1:VBD x2:NP-C) x3:.) -> x0 ", " x1 x2 x3

ARABIC / ENGLISH

- 0.54 S(x0:NP-C VP(x1:VBD x2:NP-C) x3:.) -> x0 x1 x2 x3
- 0.44 S(x0:NP-C VP(x1:VBD x2:NP-C) x3:.) -> x1 x0 x2 x3

Language Models

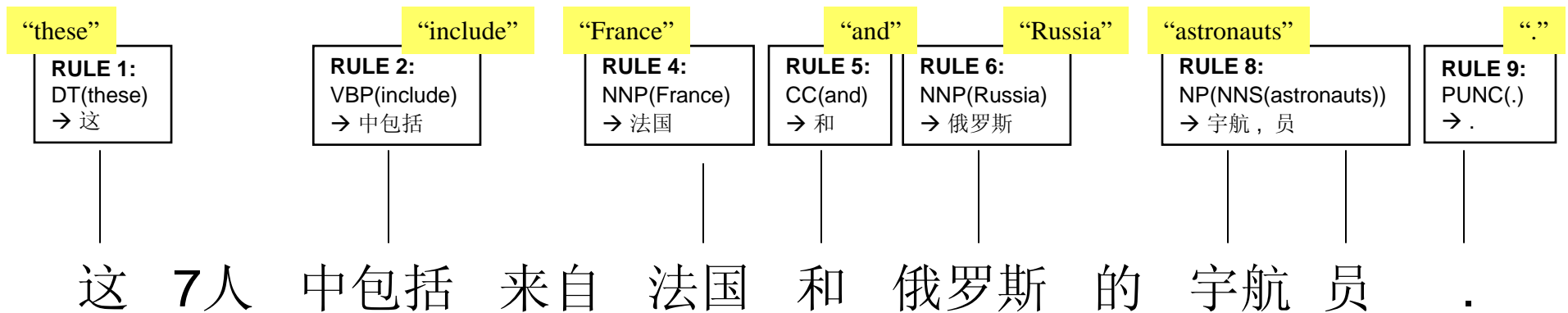
- Syntax-based Language Model
 - Assigns $P(\text{tree})$ [Collins, 1997; Charniak, 2001]
 - Unlike parser, must be trained on domain data
 - Still unproven!
- N-gram Language Model
 - Standard trigram model
 - “Only judge a tree by its leaves”
 - Used in current syntax-based MT systems

Decoder

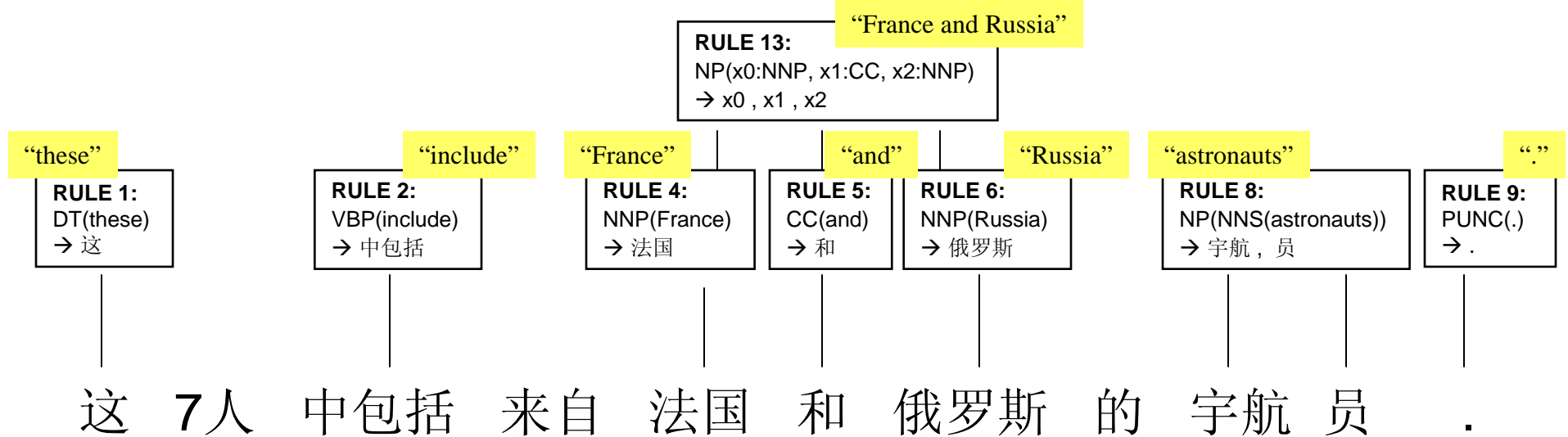
- Bottom-up CKY parser
- Builds English constituents on top of Chinese spans
- Record of rule applications (the derivation) provides information to construct English tree
- Returns k-best trees

Syntax-Based Decoding

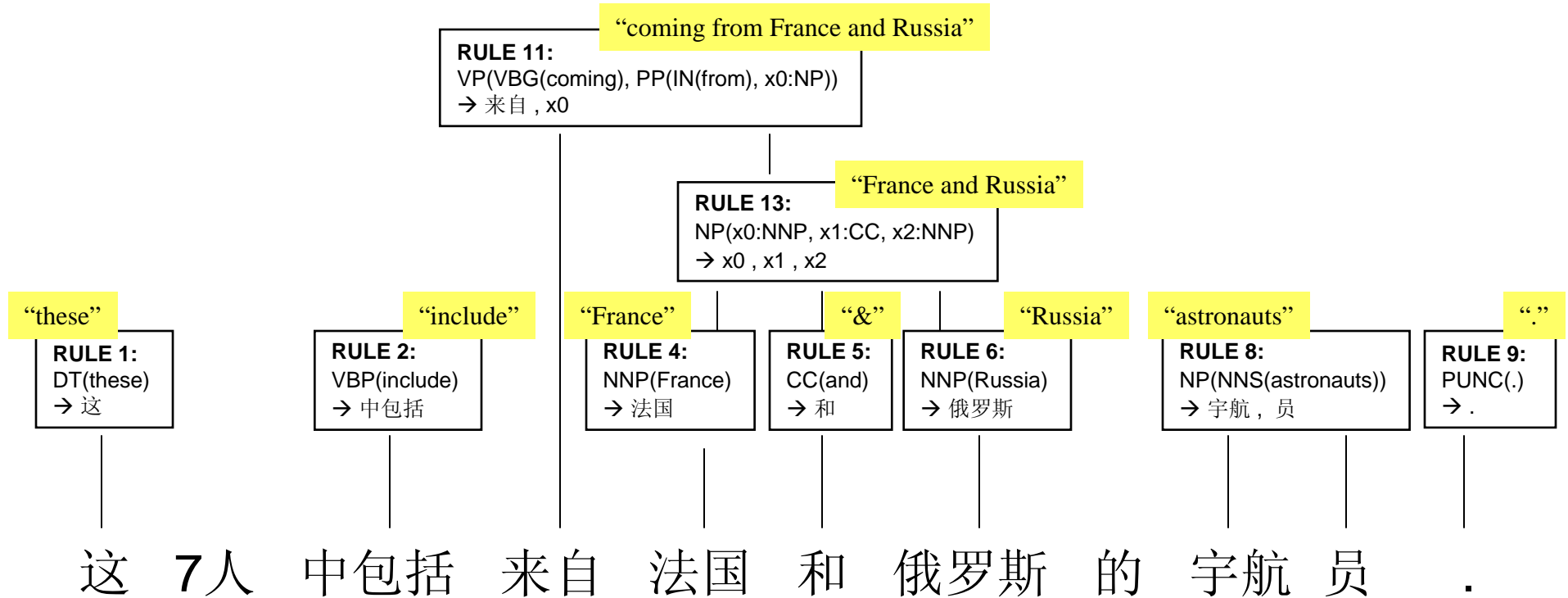
Rules apply when their right-hand sides (RHS) match some portion of the input.



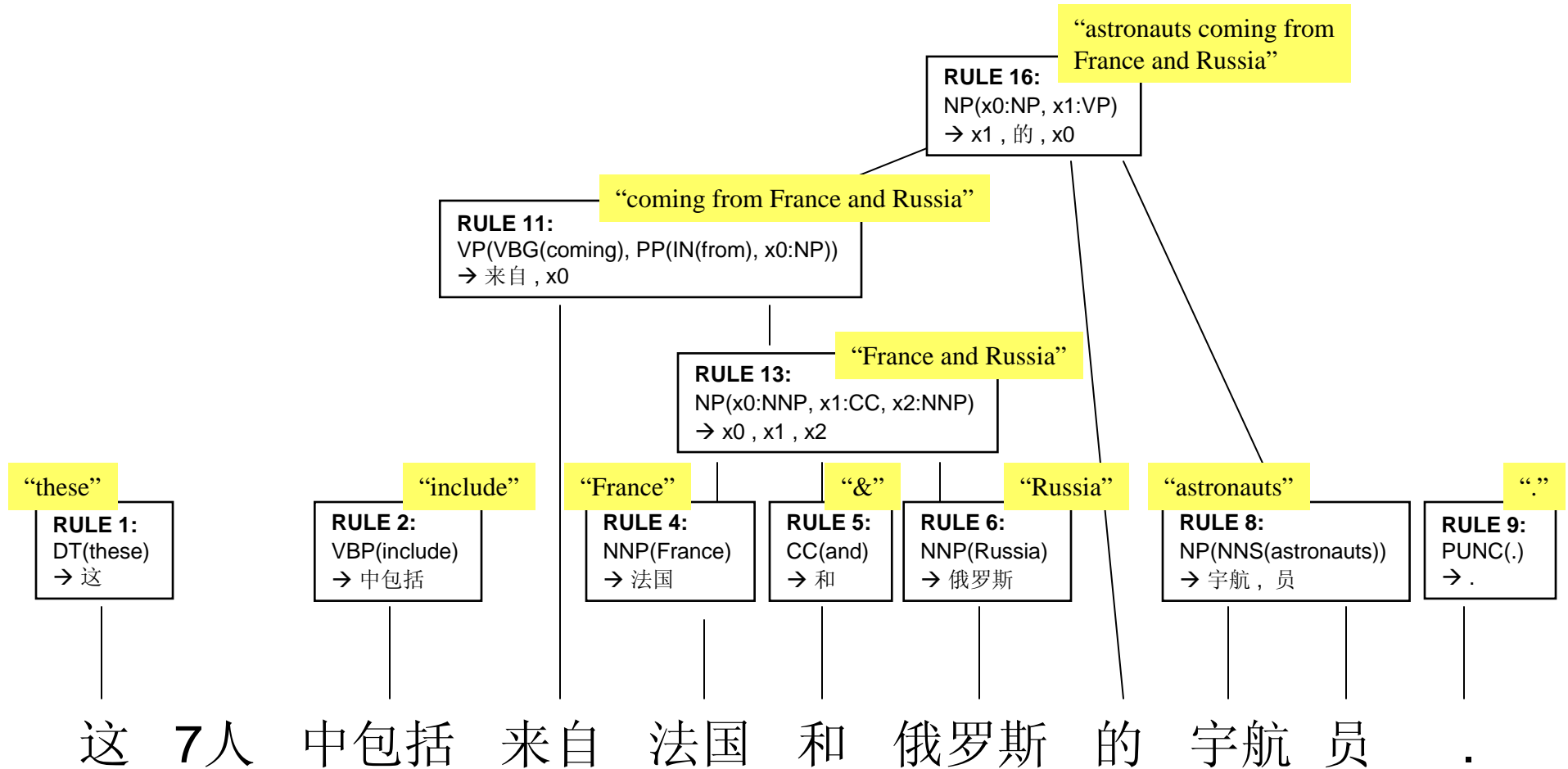
Syntax-Based Decoding

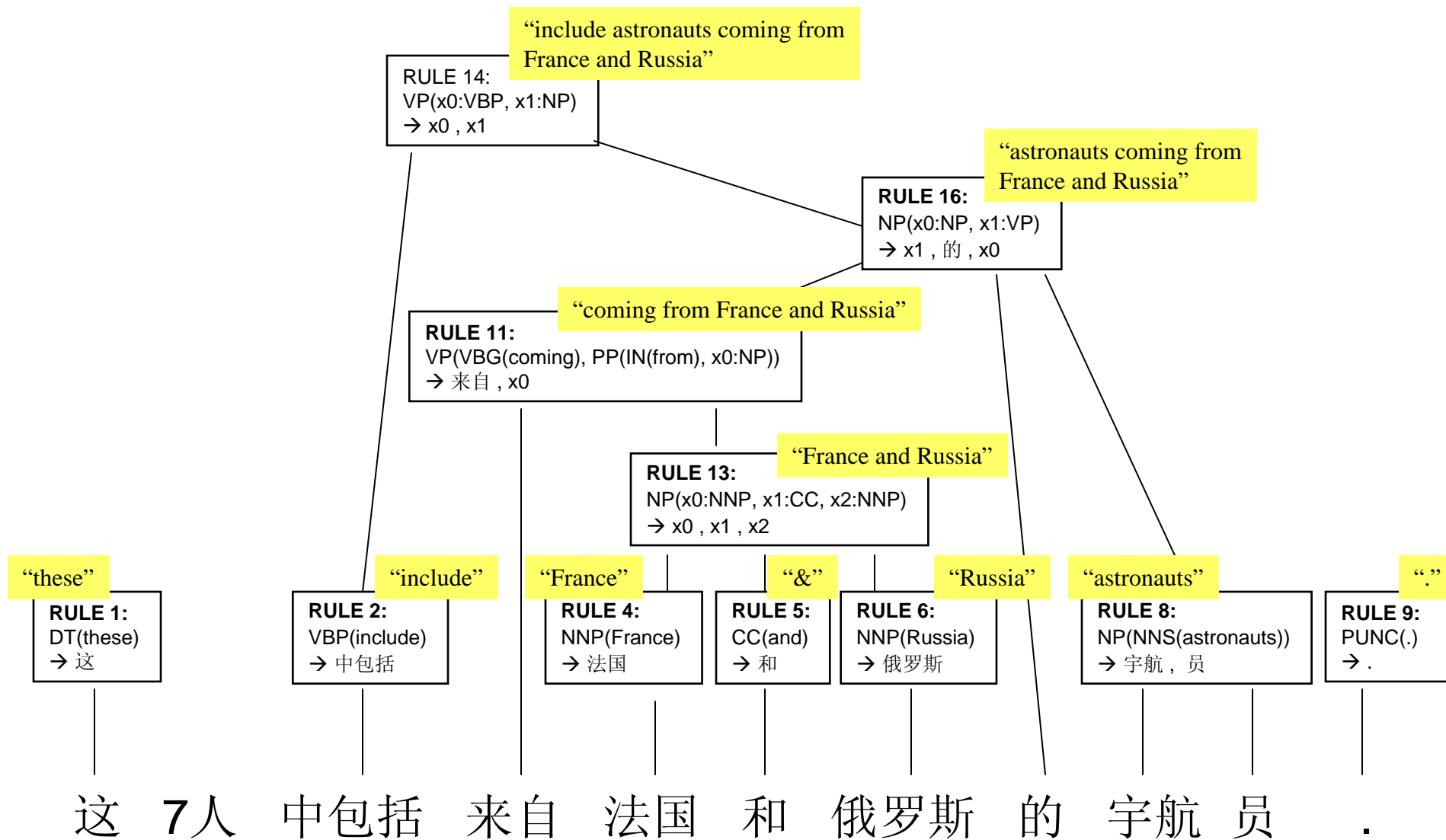


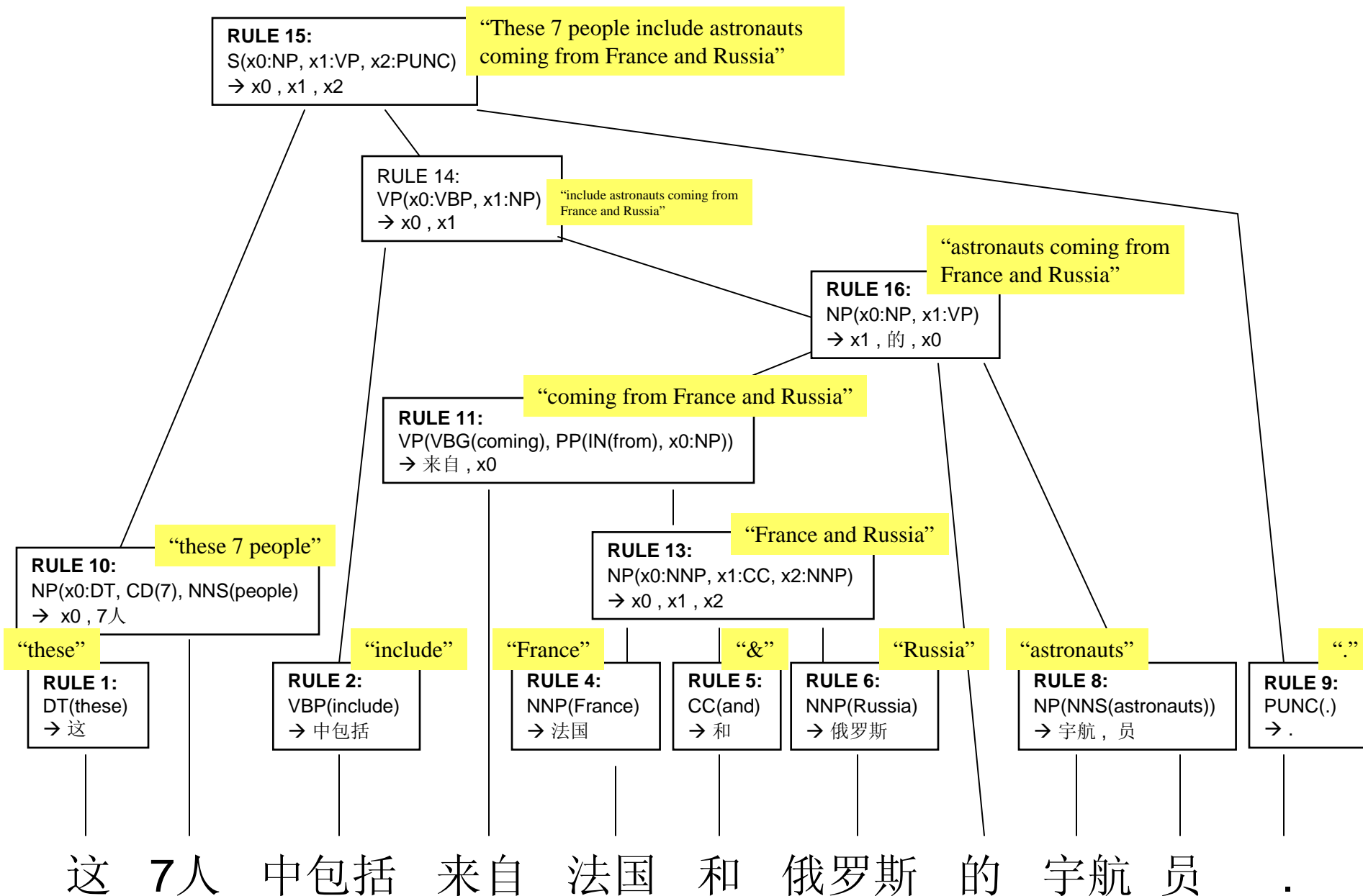
Syntax-Based Decoding



Syntax-Based Decoding

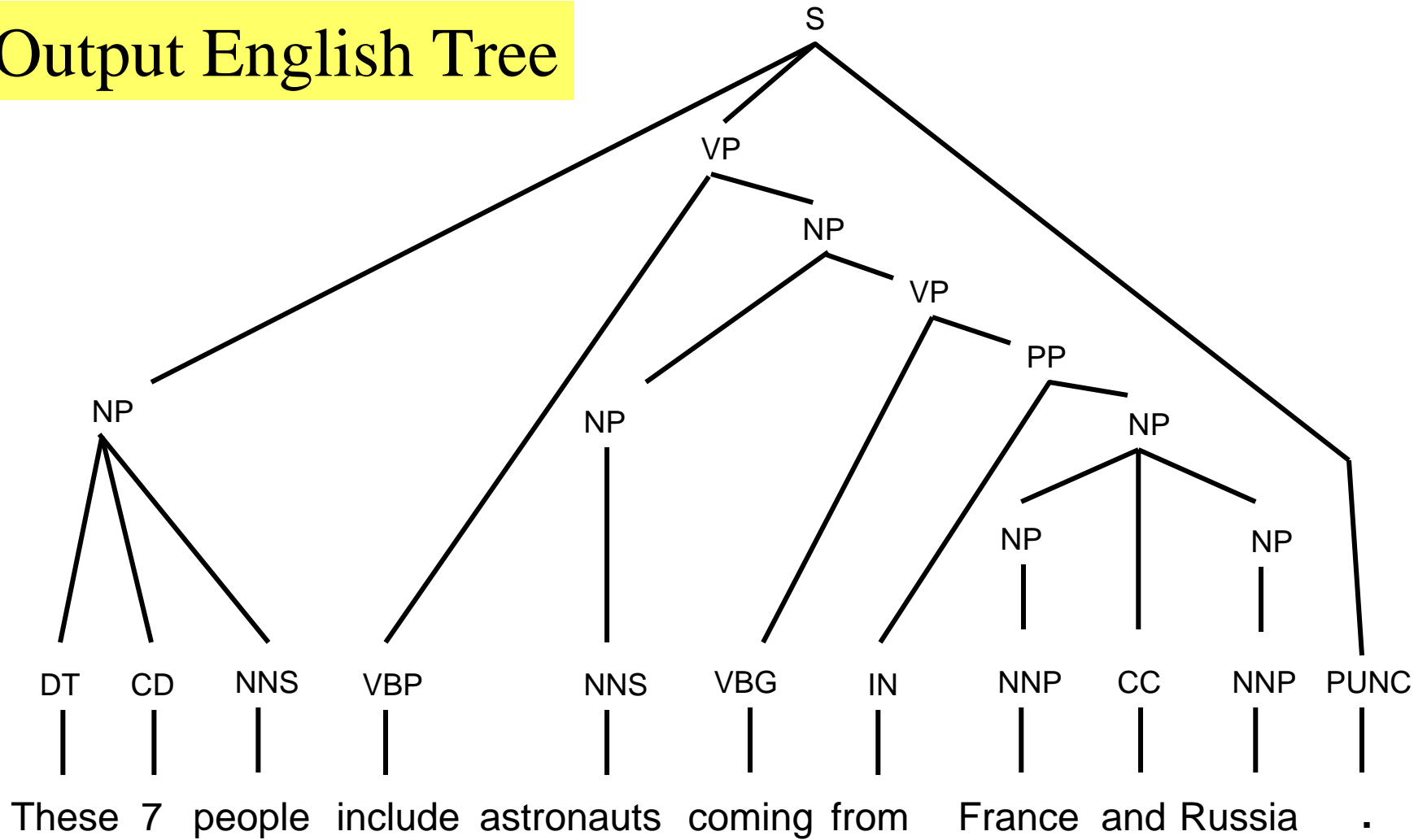






Derivation Tree

Output English Tree



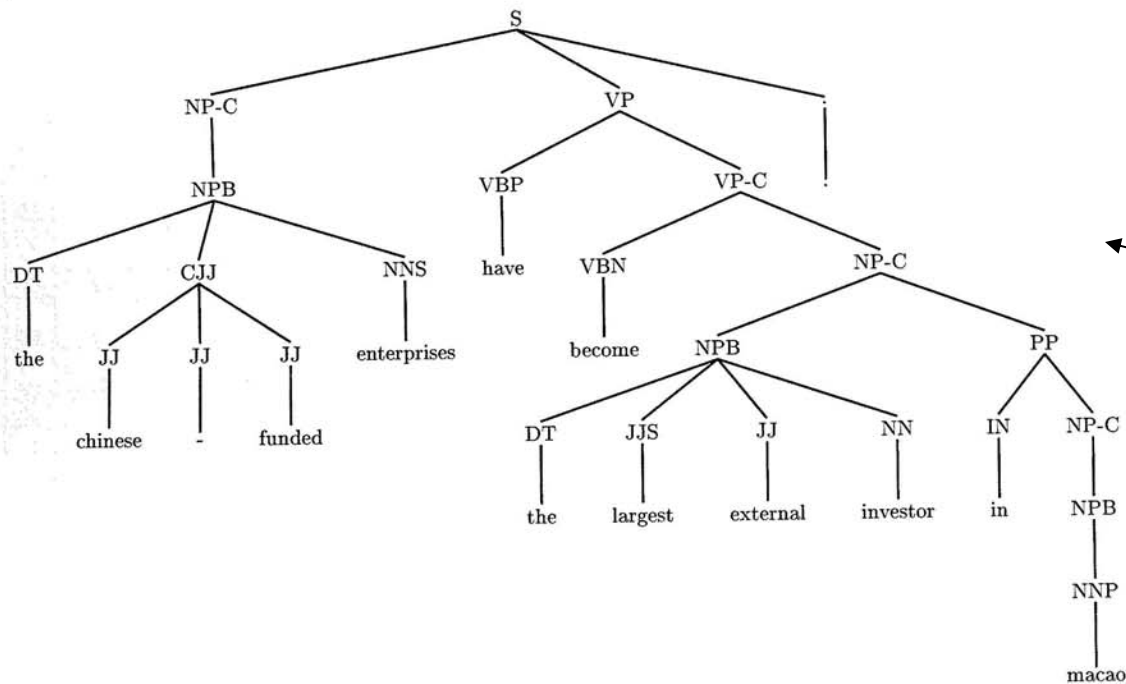
Empirical Questions

- With the acquired rules:
 - Can you always reach a top S for new Chinese sentences?
 - Does reaching a top S result in overall grammaticality?
 - Is it only possible to reach a top S by changing the meaning?
 - Is the overall translation accuracy good?

dev-little (line 47) - dev-little

Input: 中资已成为澳门最大的外来投资者。
Reference: the chinese enterprises have become the biggest outside investors in macao .
AlTemp-e: investment₀ | in₁ | macao₂ | has₃ | become₄ | the largest₅ | foreign₆ | investors₇ | .₈
AlTemp-f: 中₁ | 资₀ | 已₃ | 成为₄ | 澳门₂ | 最大的₅ | 外来₆ | 投资者₇ | .₈
[dev-little] 1-Best: the chinese - funded enterprises have become the largest external investor in macao .

[dev-little] 1-Best Tree



input

phrase-based
system output

syntax-based
system output

dev-little (line 59) - dev-little

Input: 基纳纳对中国过去向坦桑提供的大量援助表示感谢。

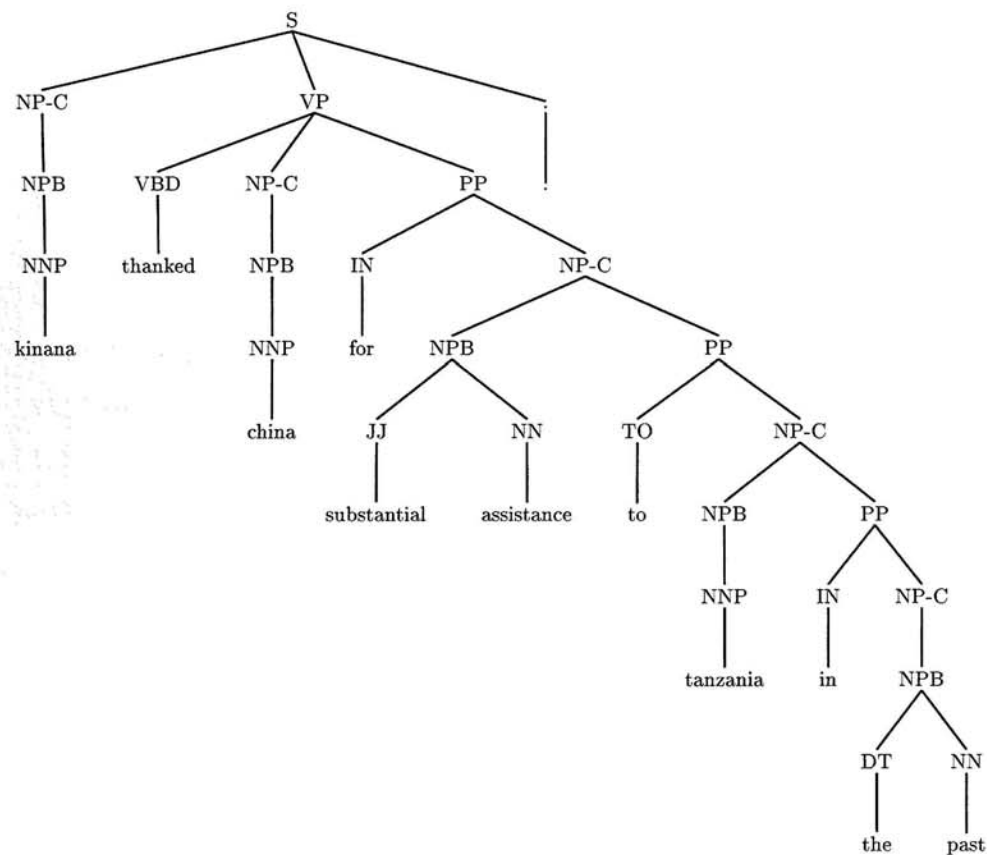
Reference: keenana expressed gratitude to china for its great assistance to tanzania in the past .

AlTemp-e: kinana₀ | to₁ | china₂ | in the past to₃ | tanzania₄ | expressed appreciation₅ | for the substantial₆ | assistance₇ | .₈

AlTemp-f: 基纳纳₀ | 对₁ | 中国₂ | 过去向₃ | 坦桑₄ | 提供的大量₆ | 援助₇ | 表示感谢₅ | .₈

[dev-little] 1-Best: kinana thanked china for substantial assistance to tanzania in the past .

[dev-little] 1-Best Tree



input

phrase-based
system output

syntax-based
system output

dev-little (line 38) - dev-little

Input: 此次 为期 两天 的 研讨会 , 由 世界贸易组织 上海 研究中心 与 上海市 对外 服务 有限公司 联合 举办 。

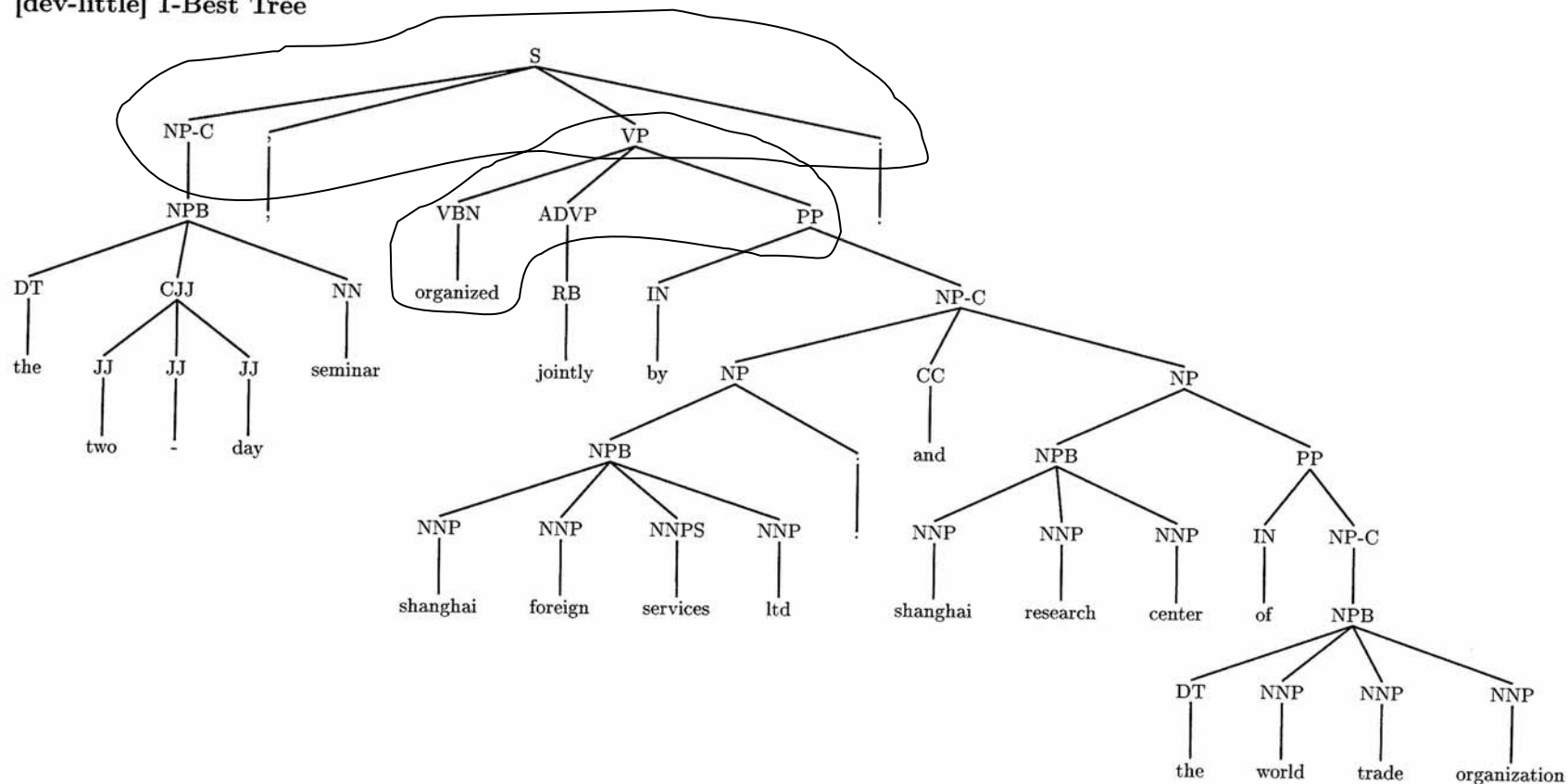
Reference: the two - day seminar is jointly sponsored by the wto shanghai research center and shanghai foreign service company limited .

AlTemp-e: the 0 | two - day 1 | seminar 2 | by the world trade organization 3 | , 4 | shanghai research center 5 | and 6 | shanghai foreign service 7 | co . , ltd . 8 | jointly 9 | . 10

AlTemp-f: 此次 0 | 为期 两天 的 1 | 研讨会 2 | , 4 | 由 世界贸易组织 3 | 上海 研究中心 5 | 与 6 | 上海市 对外 服务 7 | 有限公司 8 | 联合 举办 9 | . 10

[dev-little] 1-Best: the two - day seminar , organized jointly by shanghai foreign services ltd . and shanghai research center of the world trade organization .

[dev-little] 1-Best Tree



lev-little (line 53) - dev-little

Input: 丁豪在儿童福利院读完小学，随后进入附近乡里一所学校上初中。

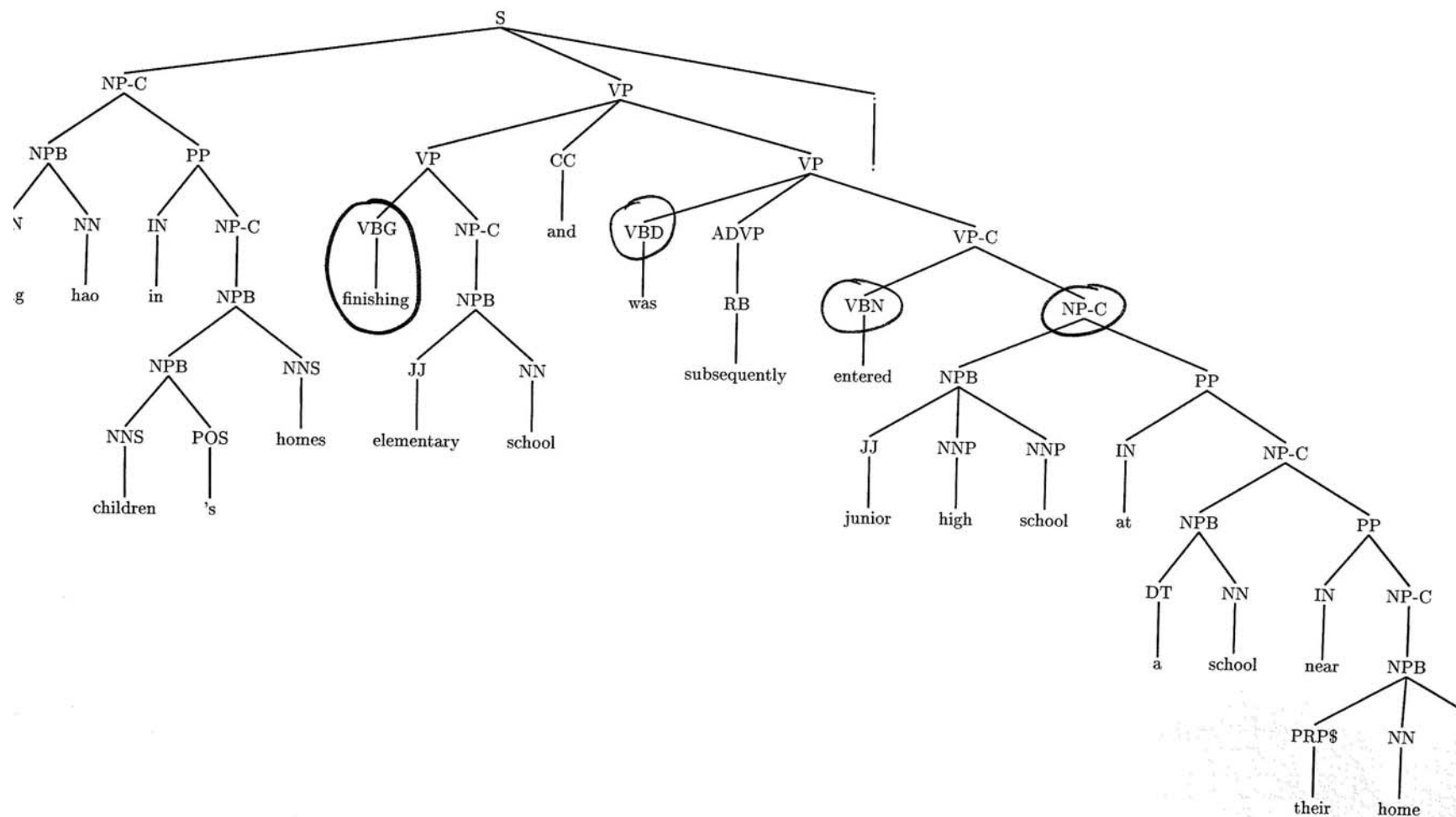
Reference: ding hao completed his primary school at the children welfare school , and then went to a nearby township middle school .

AlTemp-e: ding hao ₀ | in ₁ | children 's welfare institute ₂ | complete primary school ₃ | , ₄ | then entered ₅ | commune ₆ | near ₇ | a school ₈ | , ₉ | junior high school level . ₁₀

AlTemp-f: 丁豪₀ | 在₁ | 儿童福利院₂ | 读完小学₃ | , ₄ | 随后进入₅ | 附近₇ | 乡里₆ | 一所学校₈ | 上₉ | 初中 。 ₁₀

dev-little] 1-Best: ding hao in children 's homes finishing elementary school and was subsequently entered junior high school at a school near their home towns .

lev-little] 1-Best Tree



dev-little (line 64) - dev-little

can become very good partners

Input: 他确信，加、中两国可以成为很好的合作伙伴。

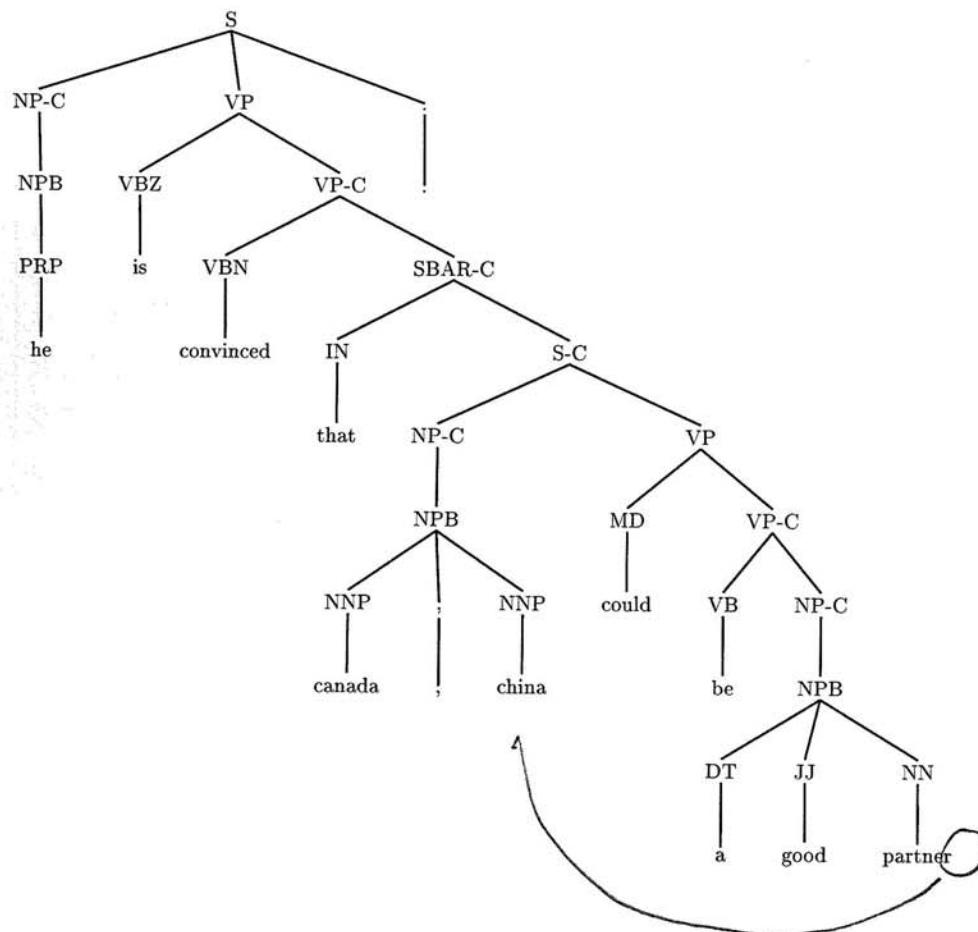
Reference: he assured that canada and china can become very good partners .

AlTemp-e: he was convinced that ₀ | the ₁ | two countries ₂ | , ₃ | can ₄ | become good ₅ | partners ₆ | . ₇

AlTemp-f: 他确信，₀ | 加 ₁ | 、中 ₃ | 两国 ₂ | 可以 ₄ | 成为很 好的 ₅ | 合作 伙伴 ₆ | 。 ₇

[dev-little] 1-Best: he is convinced that canada , china could be a good partner .

[dev-little] 1-Best Tree



Subj - obj
number agreement.

dev-little (line 51) - dev-little

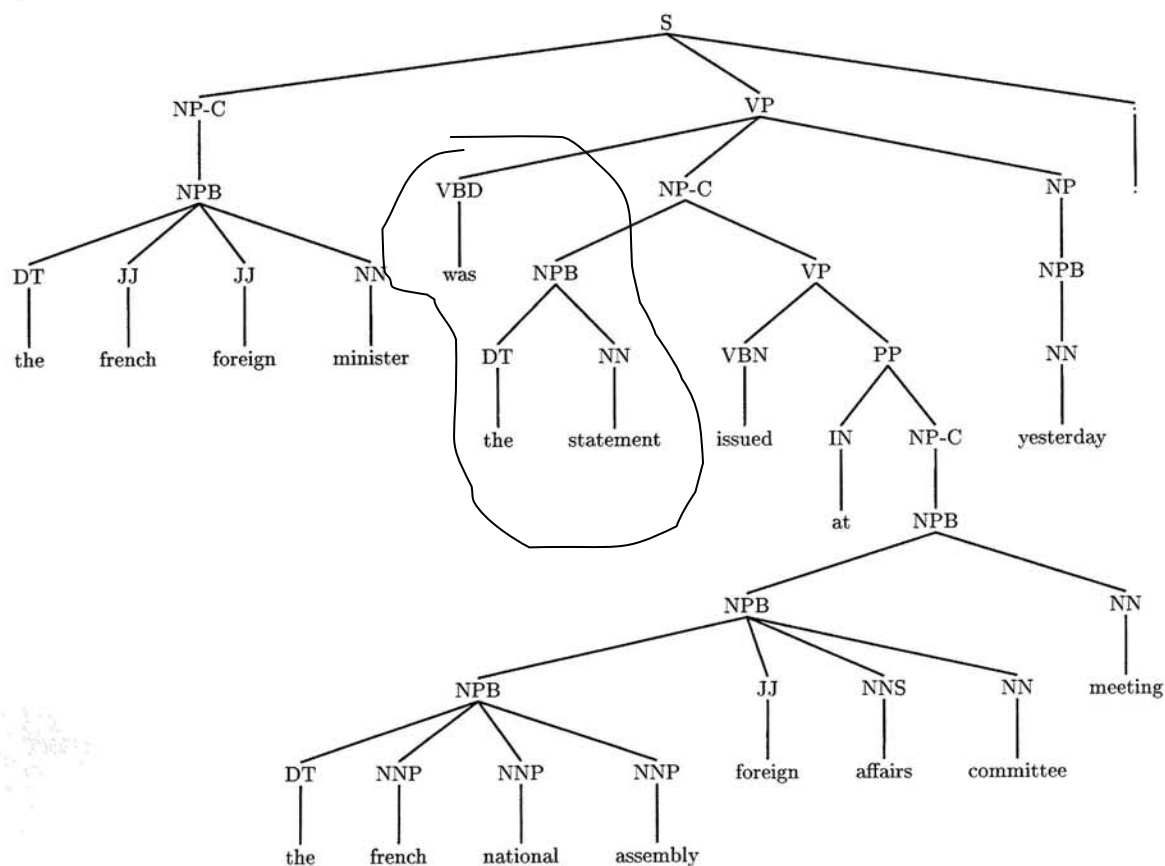
Input: 法国 外长 昨天 是 在 法国 国民议会 外事 委员会 会议 上 发表 上述 声明 的 。

Reference: the french foreign minister made the above statement in a meeting of the foreign affairs commission of the french national congress .

AlTemp-e: french₀ | foreign minister₁ | in the french national assembly₂ | yesterday₃ | the statement delivered by₄ | foreign affairs₅ | committee meeting₆ | .₇

AlTemp-f: 法国₀ | 外长₁ | 昨天 是₃ | 在 法国 国民议会₂ | 外事₅ | 委员会 会议 上₆ | 发表 上述 声明₄ | 的 。₇

[dev-little] 1-Best: the french foreign minister was the statement issued at the french national assembly foreign affairs committee meeting yesterday .



dev-little (line 125) - dev-little

Input: 今年在加利福尼亚州和南部地区^{rain}的豪雨都归咎于厄尔尼诺作宠。

Reference: the torrential rain this year in california and its southern part is attributed to the el nino .

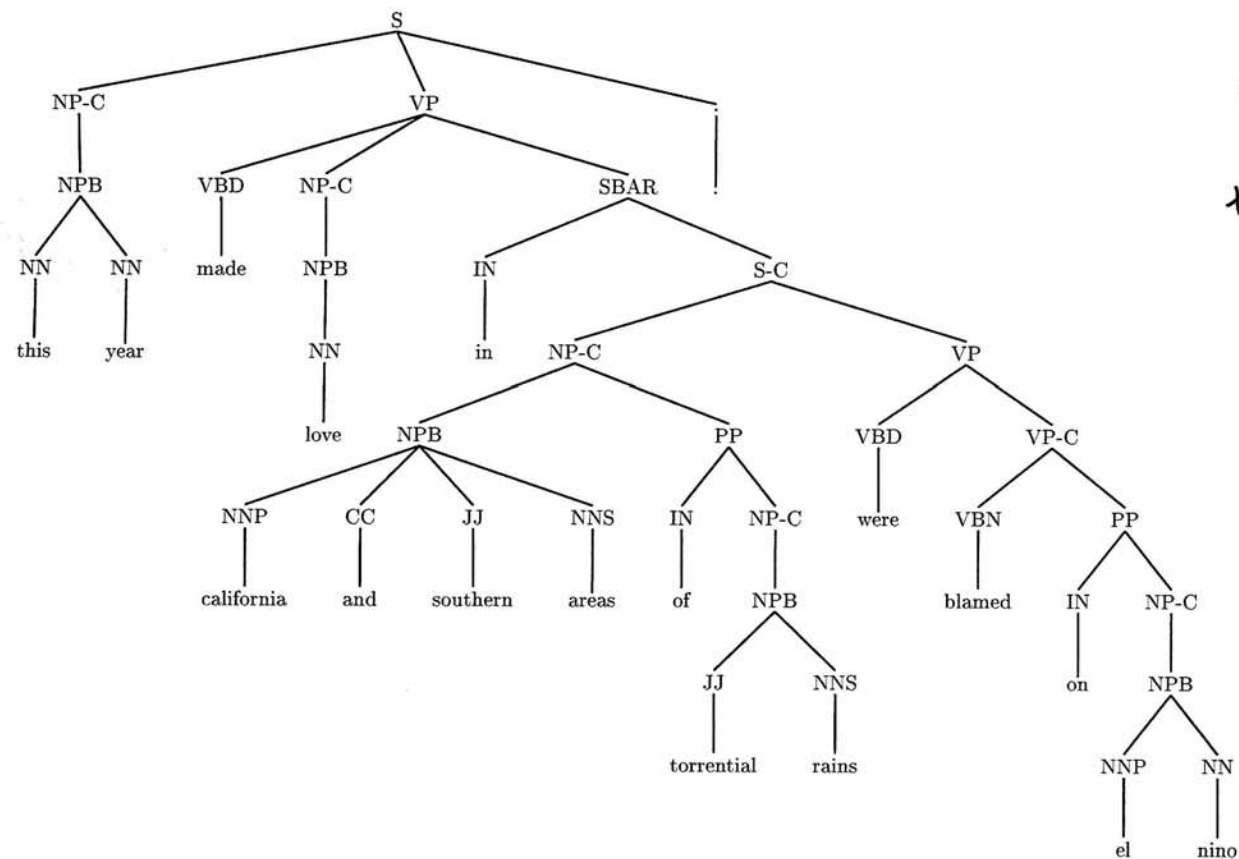
AlTemp-e: this year₀ | in california₁ | and southern₂ | areas₃ | of heavy rains₄ | attributed to₅ | a₆ | favorite₇ | el nio₈ | .₉

AlTemp-f: 今年₀ | 在加利福尼亚州₁ | 和南部₂ | 地区₃ | 的豪雨₄ | 都归咎于₅ | 厄尔尼诺₈ | 作₆ | 宠₇ | 。₉

[dev-little] 1-Best: this year made love in california and southern areas of torrential rains were blamed on el nino .

[dev-little] 1-Best Tree

Scope.



funny.

this year really NP-C

Lots of Open Problems

- Specific to MT:
 - Choosing syntactic categories that are appropriate for translation
 - Decoder search errors
 - More context for rule choice
 - Syntax-based language models
- For of NLP and beyond:
 - Modeling with tree transducers
 - Algorithms for tree transducers
 - Generic software toolkits for tree transducers

Tiburon: A Tree Automata Toolkit

- Developed by Jonathan May, USC/ISI
- First version distributed in April (www.isi.edu...)
- You cast your problem in terms of tree acceptors and transducers
- You get implemented algorithms for free
 - Kumar/Byrne'03 do this for phrase-based MT
 - Pereira/Riley'96 do this for ASR
- Wealth of tree automata literature to draw on
- Still lots of open problems in tree automata and in choosing formalisms for modeling NLP

Tiburon: A Tree Automata Toolkit

| | String World | Tree World |
|--------------------------|---------------------------|------------------|
| Weighted Sets | String acceptors (WFSA) | Tree acceptors |
| Weighted Transformations | String transducers (WFST) | Tree transducers |

Tiburon: A Tree Automata Toolkit

| | String World | Tree World |
|------------------------|--|--|
| N-best ... | ... paths through a lattice (Viterbi, 1967; Eppstein, 1998) | ... trees in a forest (Huang & Chiang, 2005) |
| EM training | Forward-backward EM (Baum & Welch, 1971) | Tree transducer EM training (Graehl & Knight, 2004) |
| Determinization ... | ... of weighted string acceptors (Mohri, 1997) | ... of weighted tree acceptors (May & Knight, 2005) |
| Intersection | WFSA intersection | Tree acceptor intersection (despite CFG not closed) |
| Applying transducers | string \rightarrow WFST \rightarrow WFSA | tree \rightarrow TT \rightarrow weighted tree acceptor |
| Transducer composition | WFST composition (Pereira & Riley, 1996) | Many tree transducers are not closed under composition! (Rounds, 1970; Engelfriet, 1975; Graehl, Hopkins, Knight) |

Classes of Tree Transducers

copying

non-copying

deleting

non-deleting

R

RL

RLN



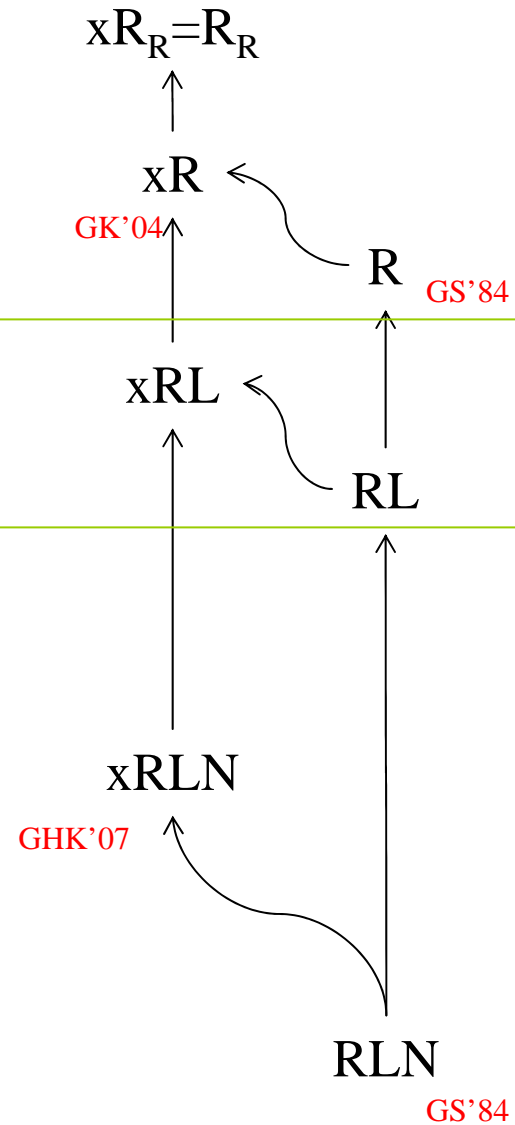
Classes of Tree Transducers

copying

non-copying

deleting

non-deleting



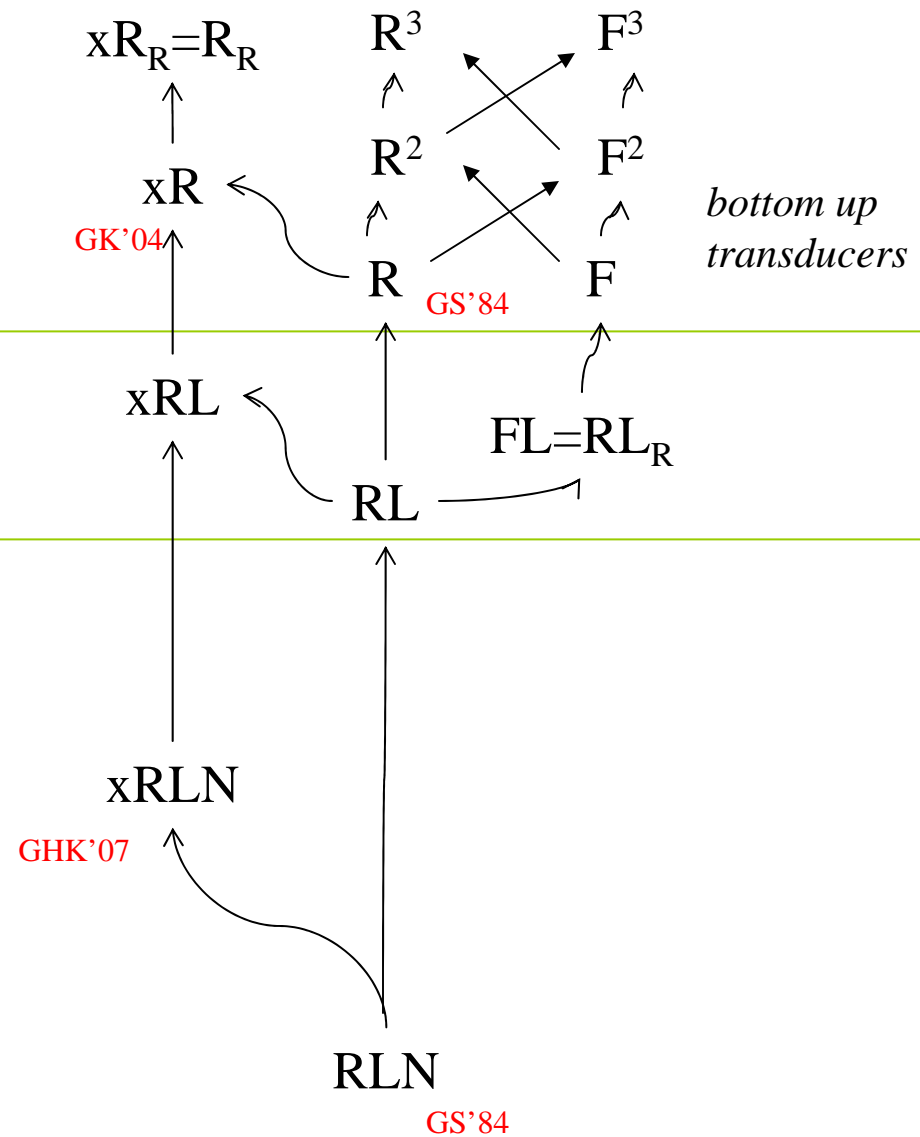
Classes of Tree Transducers

copying

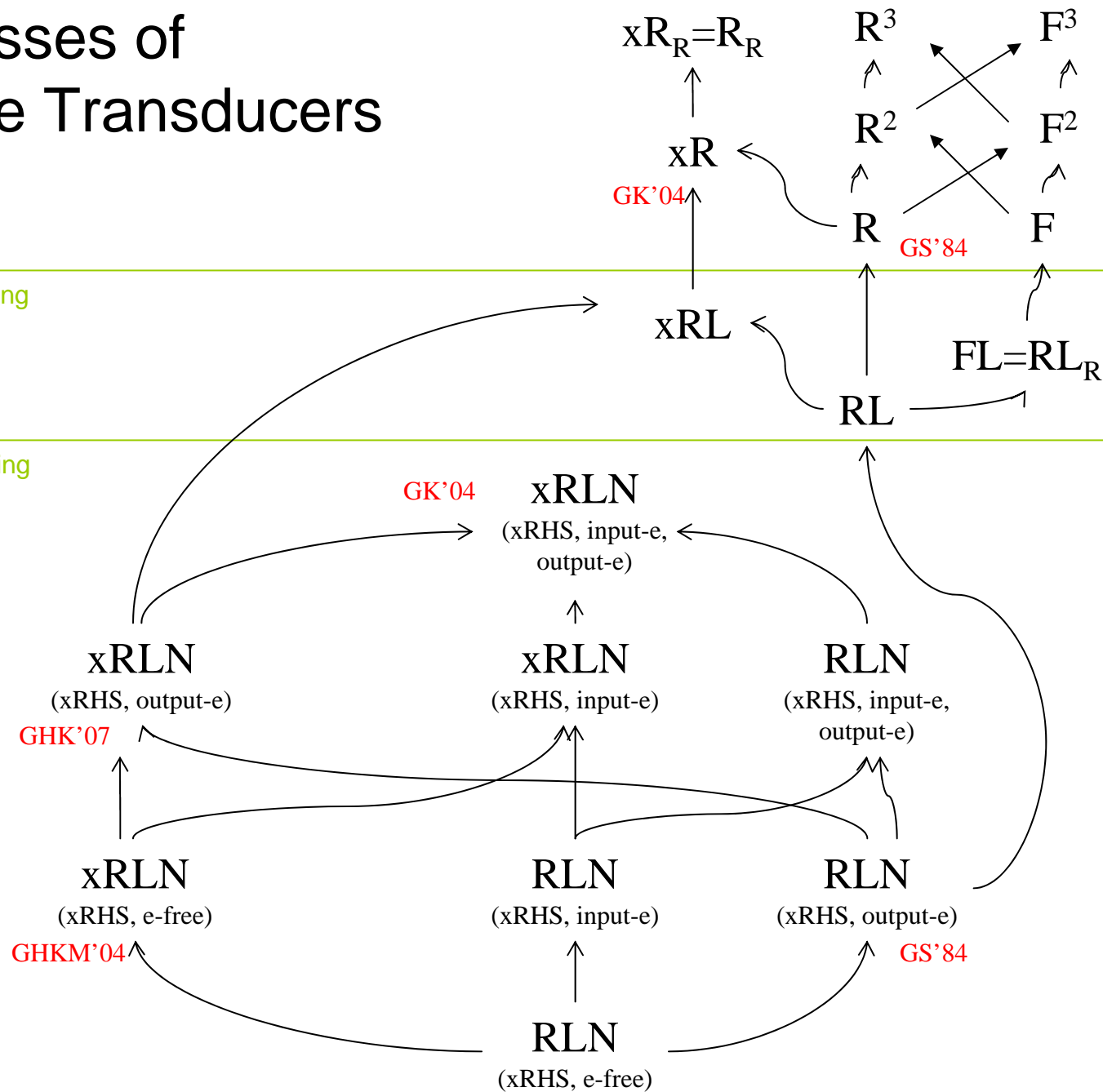
non-copying

deleting

non-deleting



Classes of Tree Transducers



Conclusion

- Making progress on machine translation
- Opening up field of tree automata to NLP
- Interdisciplinary Research
 - Machine Learning
 - Engineering
 - Linguistics
 - Efficient search algorithms
 - Automata theory
 - Grid computing

~~the end~~
beginning!