# Harmonizing word alignments and syntactic structures for extracting phrasal translation equivalents

Dun Deng, **Nianwen Xue**, Shiman Guo

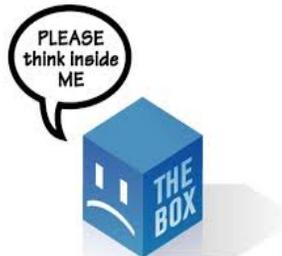Brandeis University

2015-06-04

Denver, Colorado

# What's an ideal language resource for statistical machine translation?

- Lots and lots of parallel texts?
- A large word and phrase aligned parallel treebank?
- A large semantically annotated parallel corpus?
  - Parallel Abstract Meaning Representation Corpus
- Other

# What's an ideal language resource for statistical machine translation?

- Lots and lots of parallel texts?
- ✓ **A large word and phrase aligned parallel treebank?**



- A large semantically annotated parallel corpus?
  - Parallel Abstract Meaning Representation Corpus
- Other

# Issues with treebanks and word alignments

- Status quo:
  - Manually annotated treebanks and word alignments independently conceived
    - Inevitable incompatibilities between word alignment and the syntactic structure
  - Treebanks on both sides are independently conceived
    - Incompatibilities between the parse trees of a sentence pair
  - Treebanks not optimized for MT
    - Trees too shallow, too deep
- Frustrated MT users
  - Forget about parse trees!
  - Or use automatically induced parse trees!

# Addressing this problem

- Create a hierarchically aligned Chinese English parallel treebank that
  - harmonizes word alignment and the syntactic structure of a sentence,
  - synchronizes the parse trees of the sentence pair
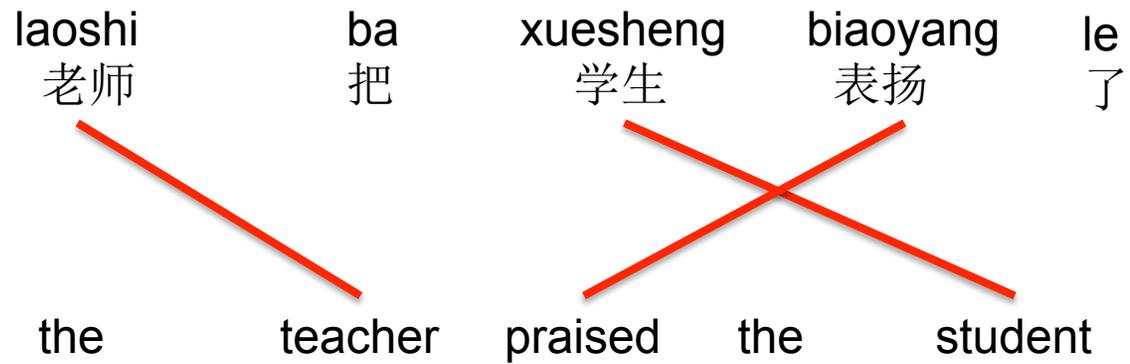  - Empirically determine the amount of structure that is needed during alignment
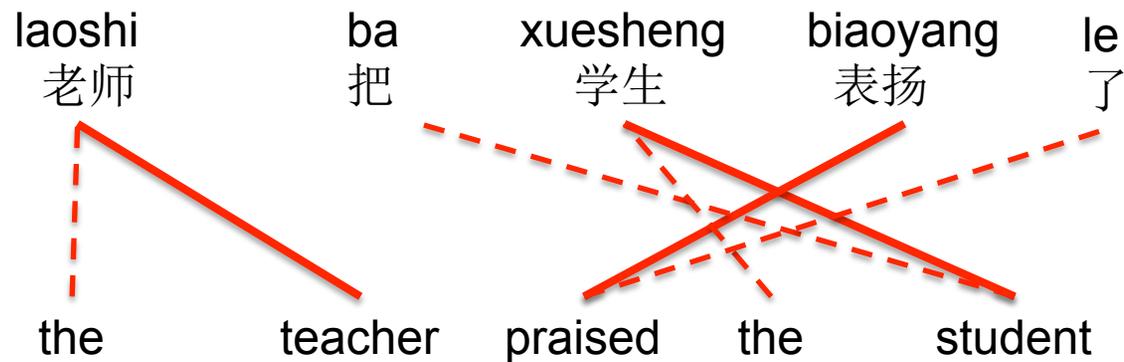
# Word alignment is deceptively complicated

✓ An equivalent exists in the target-language sentence, which matches the word in both lexical meaning and grammatical function

**?** There is a candidate in the target-language sentence, which does not have the same lexical meaning and/or grammatical function as the word but could be aligned in the given context

**?** The word has no translation counterpart in the target-language sentence at all.
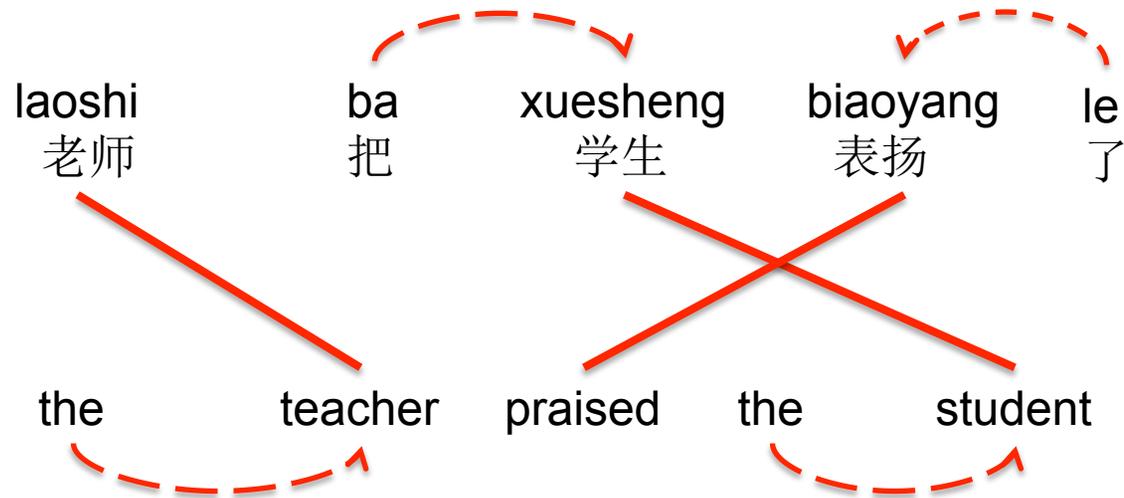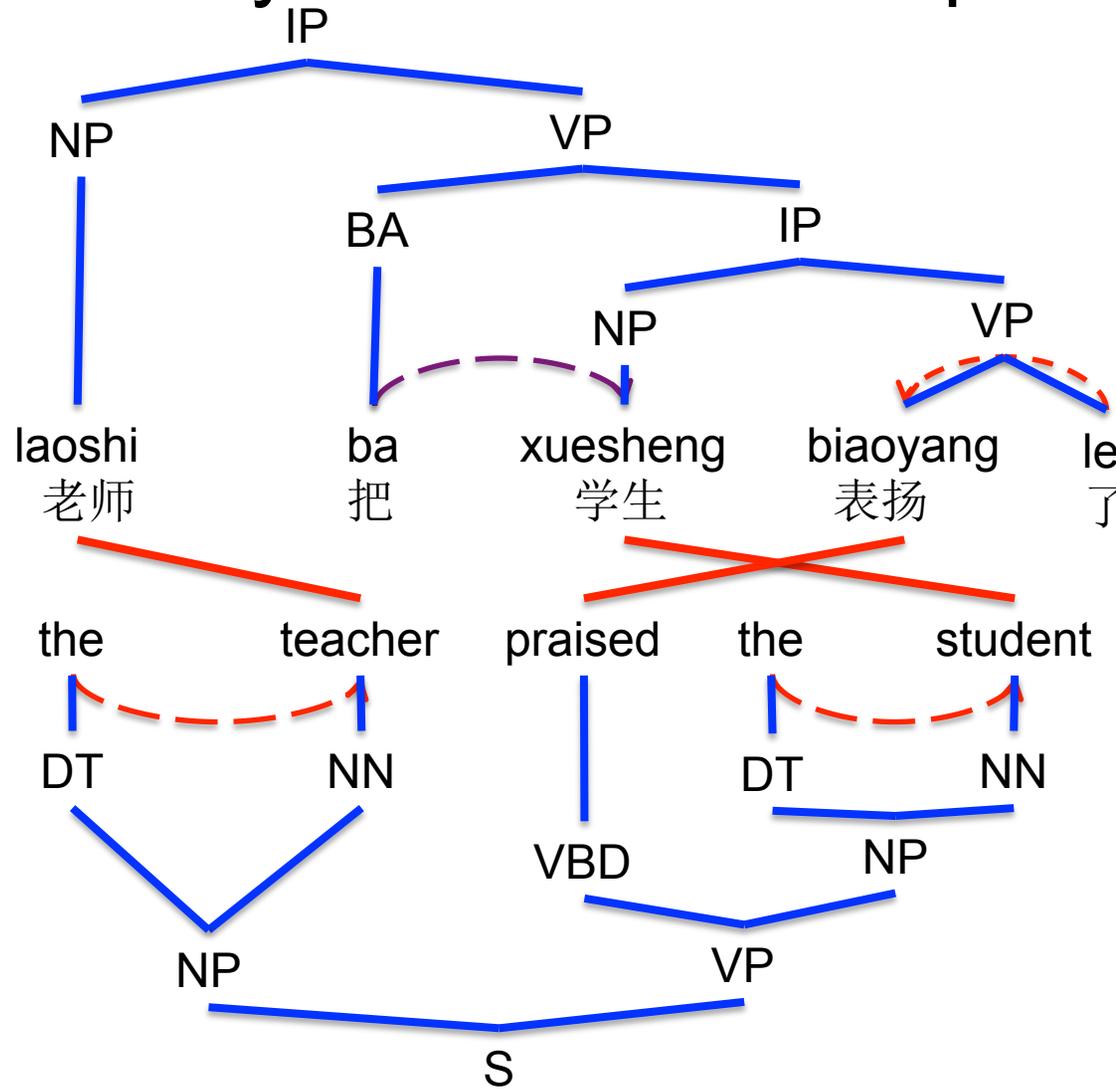
# Word alignment

laoshi        ba      xuesheng     biaoyang     le

老师        把      学生      表扬     了

the      teacher     praised     the      student

# Word alignment



laoshi     ba     xuesheng     biaoyang     le
老师     把     学生     表扬     了

the     teacher     praised     the     student

老师⇔teacher
老师⇔the teacher

学生⇔student
学生⇔the student
把学生⇔**student**
把学生⇔**the student**

# "Syntactic annotation" in word alignment
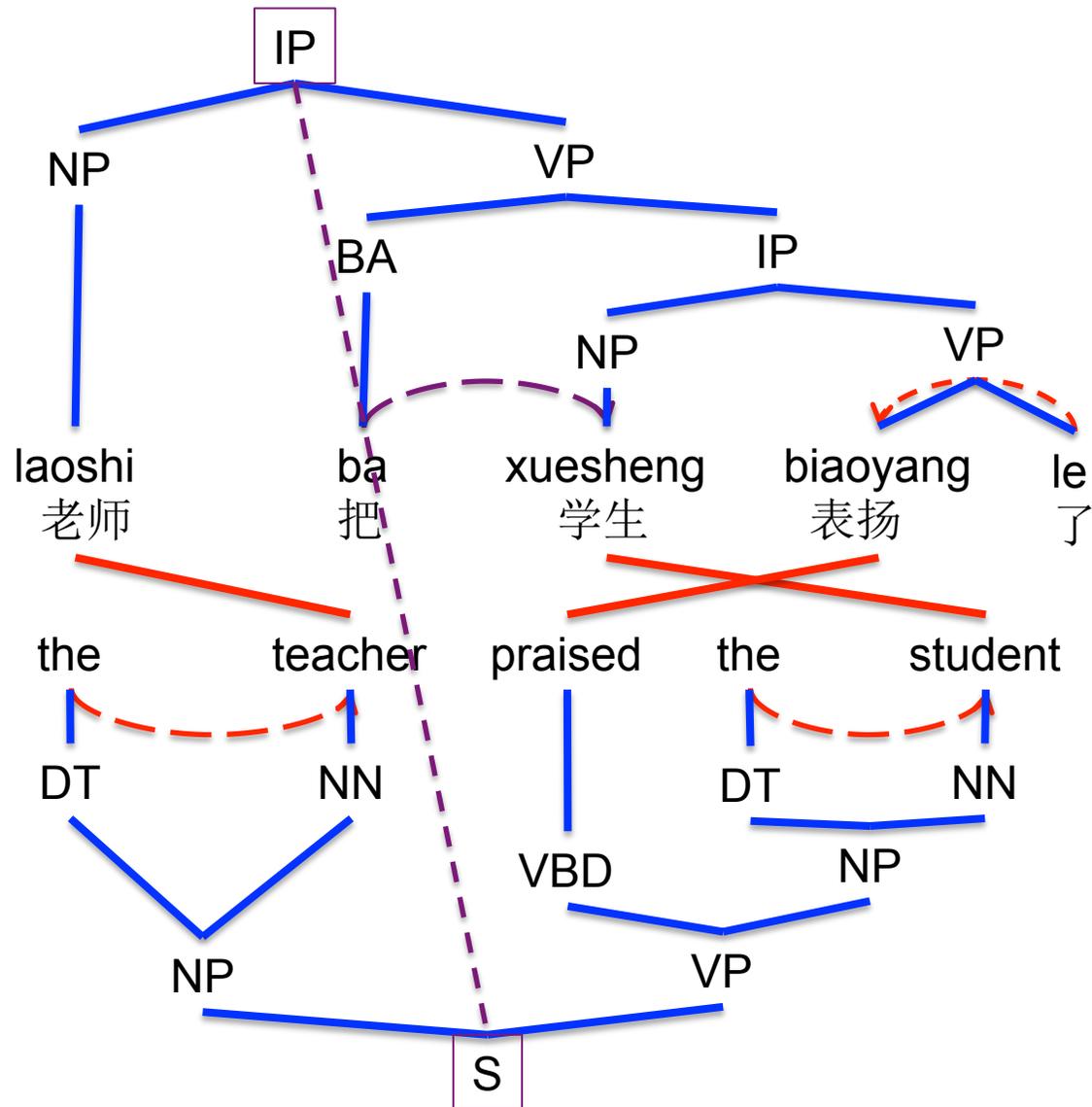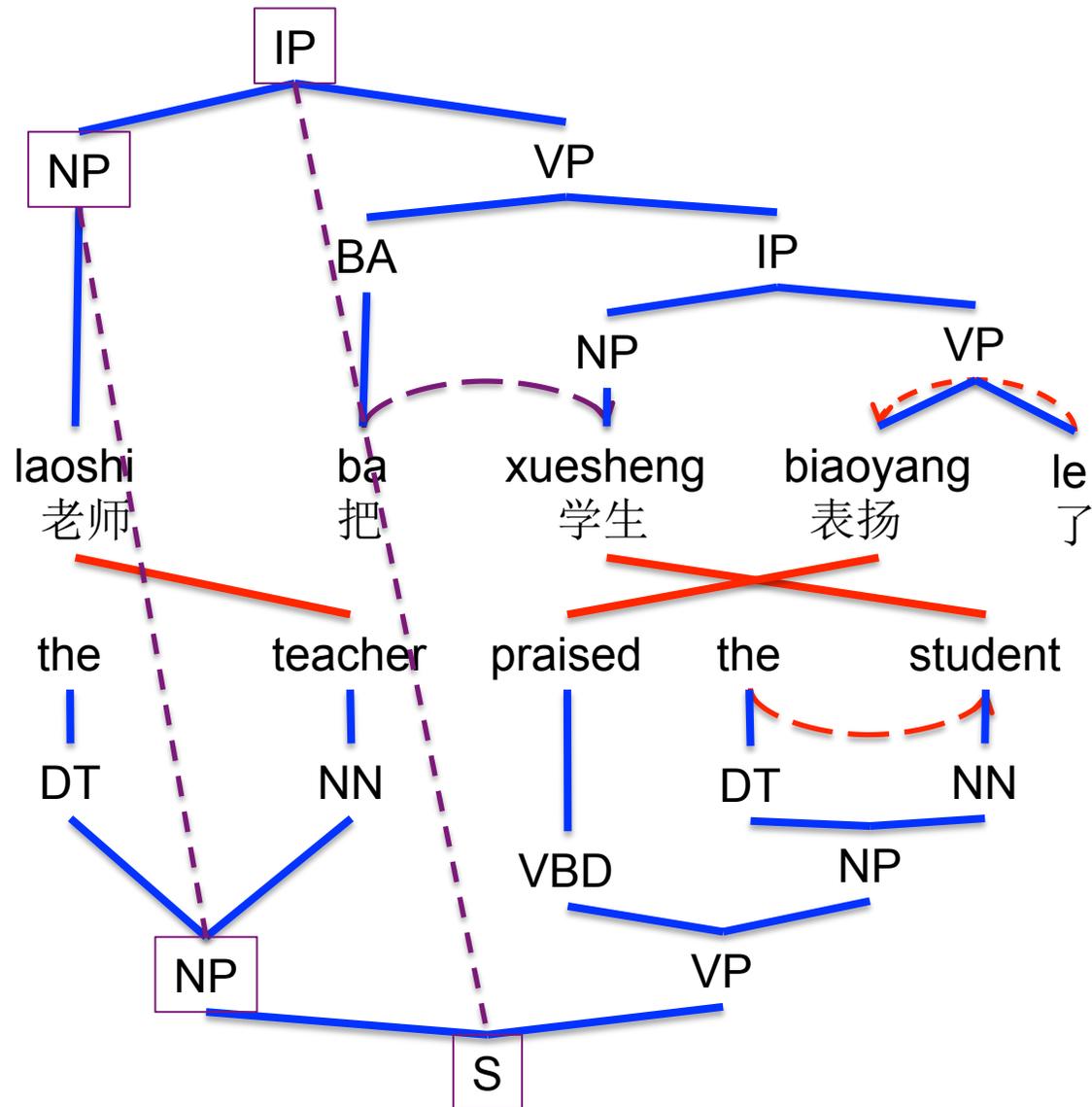
# Redundancy and conflict with parse trees

# Causes of conflict

- Word alignment standards and treebanking standards are independently established and are meant to be used separately
- Need an approach that systematically considers the interaction between word alignment and syntactic structure to maximize the utility of aligned parse trees
- Solution: hierarchical alignment
  - "sure" alignments only at word level, others aligned indirectly via node alignment in their context
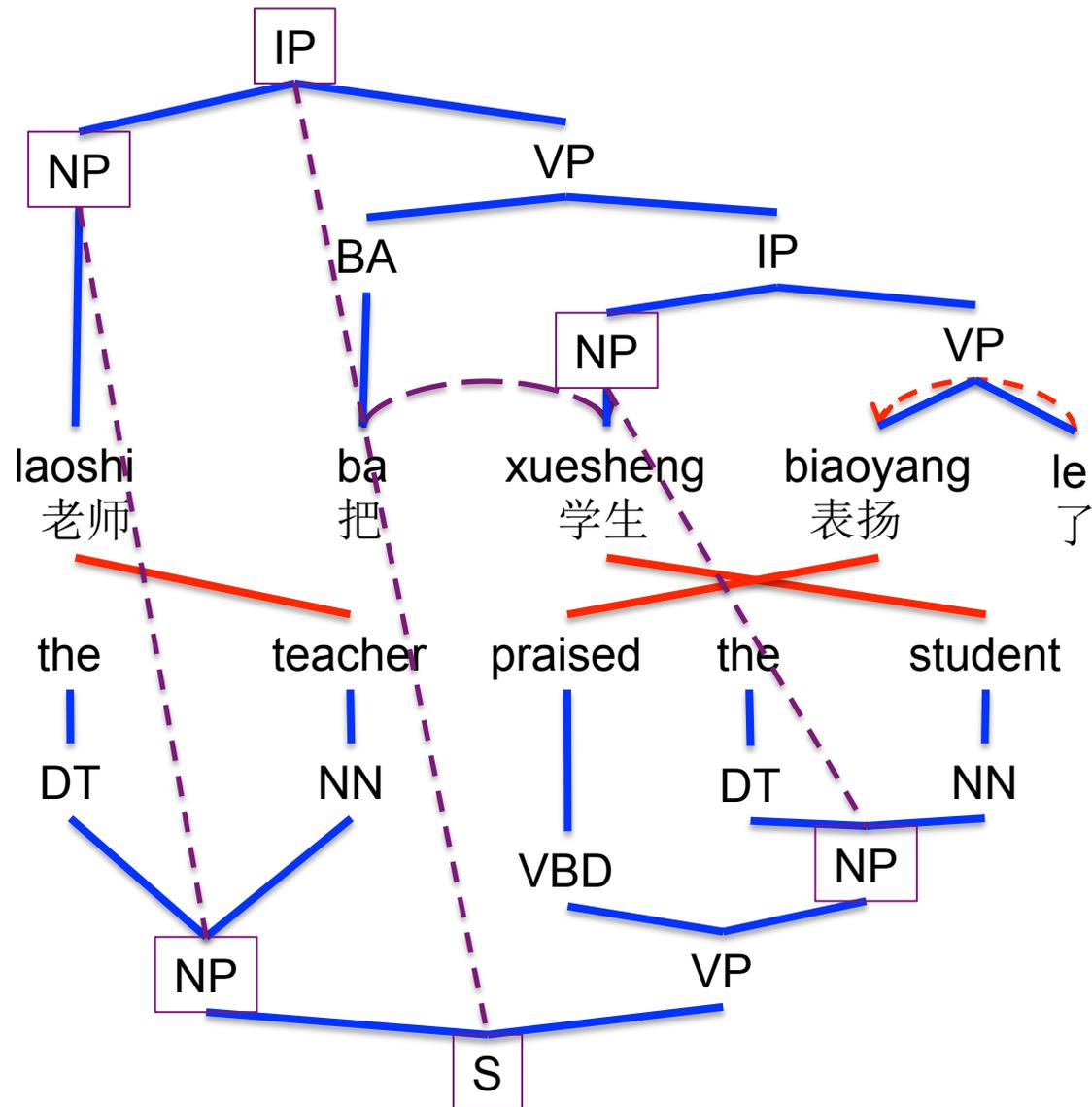
# Hierarchical alignment

# Hierarchical alignment



13

# Hierarchical alignment
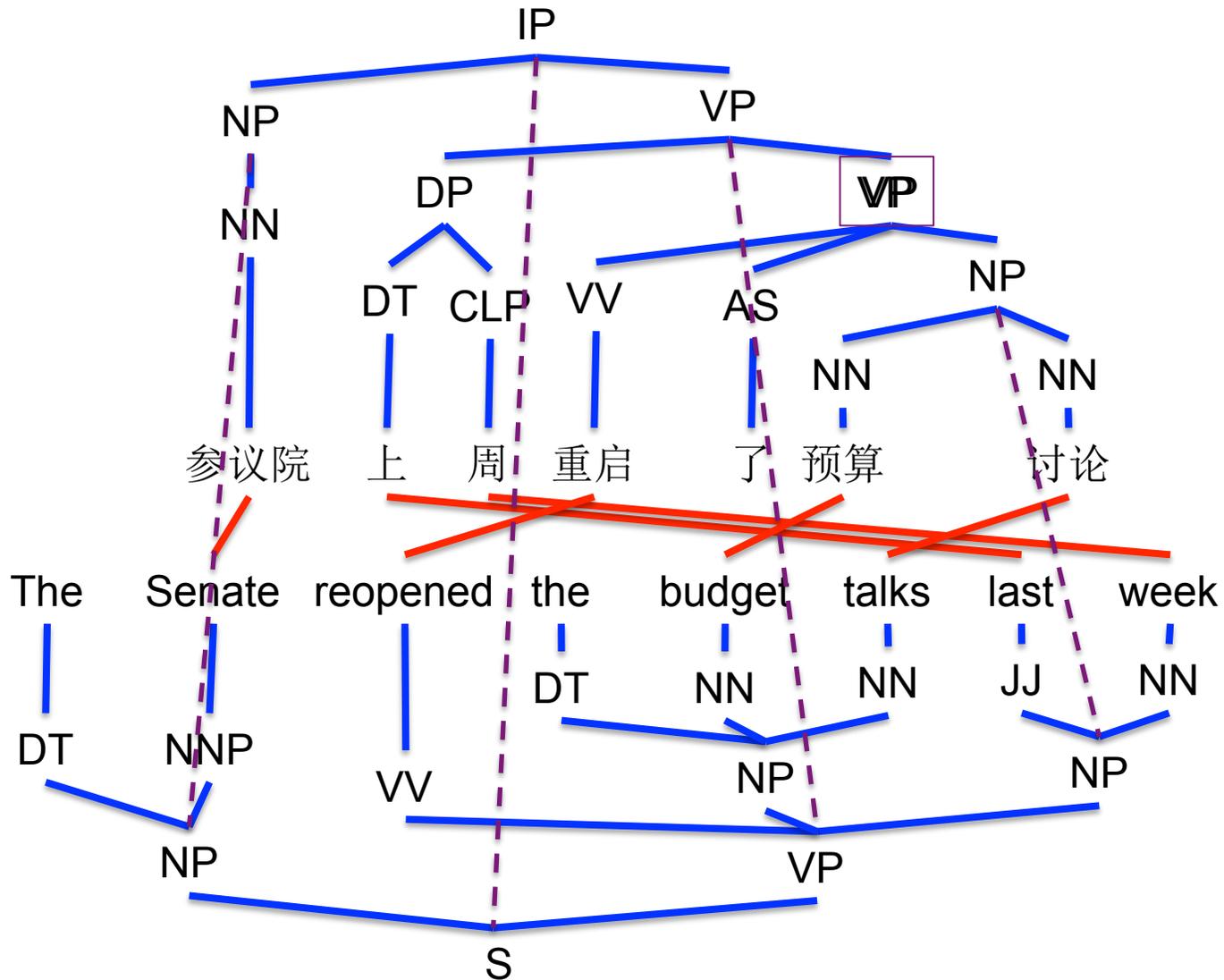
# Hierarchical alignment

# Issues in aligning parse trees

- ## Too much structure
  - Not all nodes in a parse tree are aligned or alignable.
  - This might suggest that not all the syntactic structures annotated in existing treebanks are necessary for MT purposes

- ## Too little structure
  - Flat structures in the parse trees can prevent legitimate alignments
  - Needs to revise trees in order to get proper alignments

# Too little structure

# Potential benefits of hierarchical alignment

- Captures unconstrained long-distance lexical dependencies

- Allows the extraction of linguistically interpretable rules

- Smaller but equally potent translation table

- Better MT accuracy?

# Fewer rules



学生⇔the student
把学生⇔the stuent
把学生⇔student
把X表扬了⇔praised X

19

# Long distance dependency: Passives

S66



[可 禁止 $X_1$ $X_2$ <--> $X_1$ is prohibited from $X_2$]

# Long distance dependency: 把 construction

S66



| TOP─IP─ADVP─────────AD | 1 | g | 最终 | | Eventually | g | 1 | RB─────────────ADVP_TMP─S─TOP |
|---|---|---|---|---|---|---|---|---|
| ─────────────────PU | 2 | extra | ， | | we | g | 2 | PRP──────────────NP_SBJ |
| ─NP_SBJ──────────PN | 3 | g | 我们 | | will | g | 3 | MD─────────────────VP |
| VP─VP─ADVP─────────AD | 4 | g | 将 | | outlaw | g | 4 | VB─────────VP/REU─VP |
| └VP─────────────BA | 5 | extra | 把 | | gravity | g | 5 | NN────────────NP |
| └IP_OBJ─NP_SBJ─NN | 6 | g | 重力 | | so | extra | 6 | IN────────────SBAR_PRP |
| └VP──────────VV | 7 | g | 定为 | | that | extra | 7 | IN |
| └NP_OBJ─ADJP─JJ | 8 | g | 非法 | | sludge | g | 8 | NN──────NP_SBJ_1─S/REO |
| └NP──────────NN | 9 | g | 因素 | | is | extra | 9 | VB────────────VP |
| ─────────────PU | 10 | extra | ， | | prohibited | g | 10 | VBN────────────VP |
| └VP─ADVP──────────AD | 11 | extra | 这样 | | from | extra | 11 | IN──────────PP_CLR |
| └VP────────────VV | 12 | extra | 可 | | running | g | 12 | VBG────────VP─S_NOM |
| └VP────────────VV | 13 | g | 禁止 | | downhill | g | 13 | RB─ADVP_DIR |
| ─NP_OBJ──────────NN | 14 | g | 淤泥 | | . | g | 14 | . |
| └IP─VP──────────VV | 15 | g | 流到 | | | | | |
| └NP_OBJ────────NN | 16 | g | 山下 | | | | | |
| ─────────────PU | 17 | g | 。 | | | | | |

[把 X 定为非法因素 <--> outlaw X]

21

# Long distance dependency: discourse Connectives



| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| TOP—IP—ADVP————————AD | 1 | g | 最终 | | Eventually | g | 1 | RB————————ADVP_TMP—S—TOP |
| ————————————————PU | 2 | extra | ， | | we | g | 2 | PRP————————————NP_SBJ |
| —NP_SBJ————————PN | 3 | g | 我们 | | will | g | 3 | MD————————————VP |
| —VP—VP—ADVP————AD | 4 | g | 将 | | outlaw | g | 4 | VB————————VP/REU—VP |
| ⌐VP————————BA | 5 | extra | 把 | | gravity | g | 5 | NN————————————NP |
| ⌐IP_OBJ—NP_SBJ—NN | 6 | g | 重力 | | so | extra | 6 | IN————————SBAR_PRP |
| ⌐VP————————VV | 7 | g | 定为 | | that | extra | 7 | IN |
| ⌐NP_OBJ—ADJP—JJ | 8 | g | 非法 | | sludge | g | 8 | NN————NP_SBJ_1—S/REO |
| ⌐NP————————NN | 9 | g | 因素 | | is | extra | 9 | VB————————————VP |
| ————————————PU | 10 | extra | ， | | prohibited | g | 10 | VBN————————————VP |
| ⌐VP—ADVP————AD | 11 | extra | 这样 | | from | extra | 11 | IN————————PP_CLR |
| ⌐VP————————VV | 12 | extra | 可 | | running | g | 12 | VBG————VP—S_NOM |
| ⌐VP————————VV | 13 | g | 禁止 | | downhill | g | 13 | RB—ADVP_DIR |
| —NP_OBJ————————NN | 14 | g | 淤泥 | | . | g | 14 | .————————————— |
| ⌐IP—VP————————VV | 15 | g | 流到 | | | | | |
| ⌐NP_OBJ————————NN | 16 | g | 山下 | | | | | |
| ————————————PU | 17 | g | 。 | | | | | |

[将 $X_1$ ， 这样 $X_2$ <--> will $X_1$ so that $X_2$]

# Long distance dependency: Questions

S1550



[ $X_1$ 为什么 不 $X_2$ <--> Why do n't $X_1$ $X_2$]

# Alignment in syntactic context: Conjunction

S3850



[X , Y <> X and Y]

# Conclusions and future work

- Described a hierarchical alignment approach that systematically considers the interaction between word alignment and the syntactic structure of a sentence.

- Showed that such alignments can be used to extract translation rules that cover long-distance dependencies.

- Aligned 10K sentence pairs annotated with PTB and CTB trees

- To do: revising treebanking guidelines and synchronizing the parse trees

- More info can be found at:

  – http://www.cs.brandeis.edu/~clp/ctb/hacept/

# Brandeis University

# Acknowledgements

- Libin Shen for inspirations
- Tool support by the IBM Multilingual NLP Technologies Group: Niyu Ge, Abe Ittycheriah, Salim Roukos
- Annotators: Gao Hui, Guo Shiman, Wang Tse-Ming, Zhou Linya
- Funded by DARPA BOLT Program via IBM