

---

## Agenda for today

---

- Introduction to bi-text parsing
  - Review of finite-state transducers
  - Inversion transduction grammars (ITG)
  - Synchronous context-free grammars (SCFG)
- Machine translation
- Text-Based Language Processing Systems: Winter 2009

---

## Review finite-state string transducers

---

- Simultaneously generate pairs of (related) strings
- Spelling:

fox	f o x
foxes	f o x e s
cat	c a t
cats	c a t s
dog	d o g
dogs	d o g s
donkey	d o n k e y
donkeys	d o n k e y s

---

## Review finite-state string transducers

---

- Simultaneously generate pairs of (related) strings
- Pronunciation:

fox	F AA1 K S
foxes	F AA1 K S AH0 Z
cat	K AE1 T
cats	K AE1 T S
dog	D AO1 G
dogs	D AA1 G Z
donkey	D AA1 NG K IY0
donkeys	D AA1 NG K IY0 Z

---

## Review finite-state string transducers

---

- Simultaneously generate pairs of (related) strings
- ...Translation?:

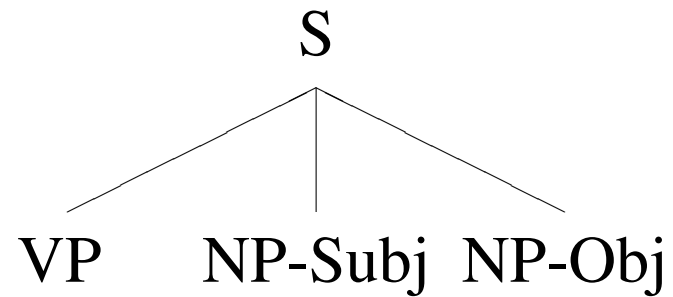
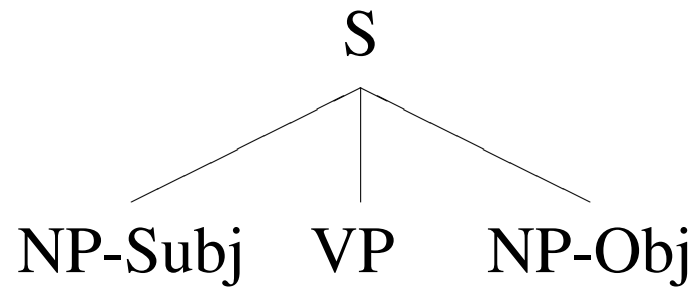
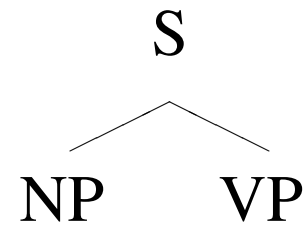
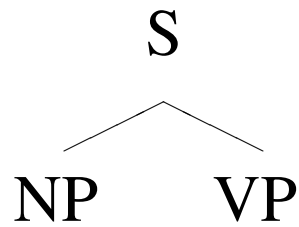
fox	zorro
foxes	zorros
cat	gato
cats	gatos
mouse	ratón
train	formarse (v)
	tren (n)
	ir en tren (v)
to go by train	ir en tren (v)

---

## Tree transducers

---

- Simultaneously generate pairs of (related) trees



---

## Synchronous grammars

---

- Simultaneously generate pairs of recursively related strings
- Will always see *pairs* in the grammar rules

$$A \rightarrow B C / B' C'$$

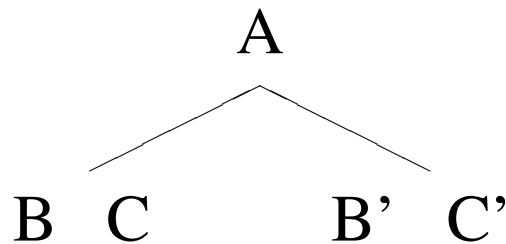
---

## Synchronous grammars

---

- Simultaneously generate pairs of recursively related strings
- Will always see *pairs* in the grammar rules

$$A \rightarrow B C / B' C'$$



---

## Synchronous grammars as transducers

---

$A \rightarrow \text{fox} / \text{f o x}$

$A \rightarrow \text{fox} / \text{F AA1 K S}$

$A \rightarrow \text{fox} / \text{zorro}$

$A \rightarrow \text{train} / \text{formarse}$

$A \rightarrow \text{train} / \text{tren}$



---

## Synchronous grammars as transducers

---

- Drawbacks?
- Ways to address ambiguity?
- Methods similar to automata algorithms for
  - Minimization
  - Determinization
  - Weight-pushing

---

# Bracketing inversion transduction grammars (ITG)

---

[Wu 1995,1997]

- Always binary rules
  - $X \rightarrow [ X_1 X_2 ]$
  - $X \rightarrow < X_1 X_2 >$
  - $X \rightarrow s / t$
  - $X \rightarrow s / \epsilon$
  - $X \rightarrow \epsilon / t$
- Generates both source and target trees
- Requires a common binary tree
  - Derivable from parallel corpora

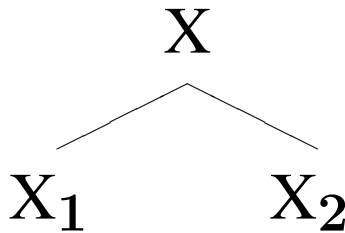
---

## ITG non-terminal rules

---

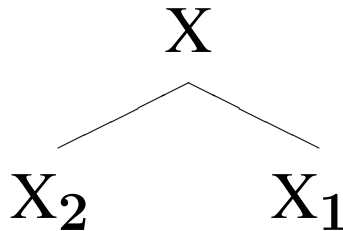
- Simple transduction production:

$$X \rightarrow [ X_1 X_2 ]$$



- **Inverted** transduction production:

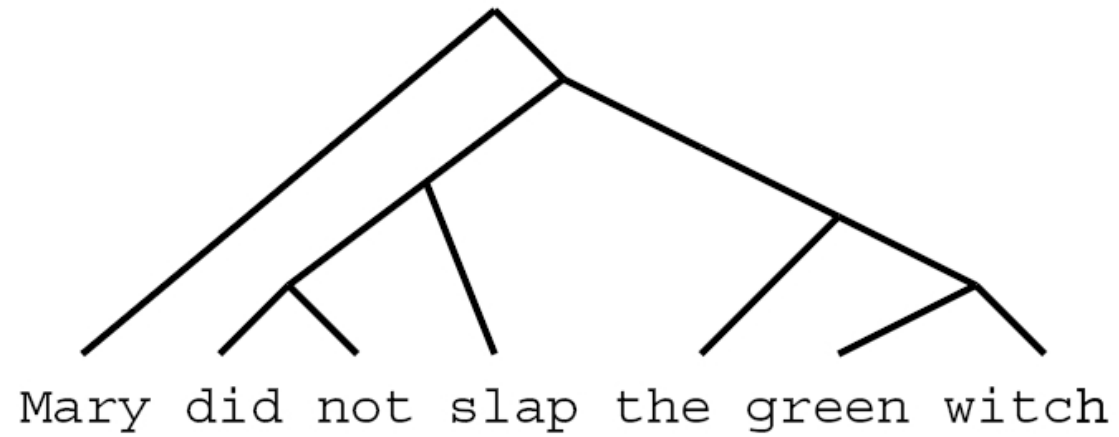
$$X \rightarrow \langle X_1 X_2 \rangle$$



---

## Example: Unlabeled binary tree (English)

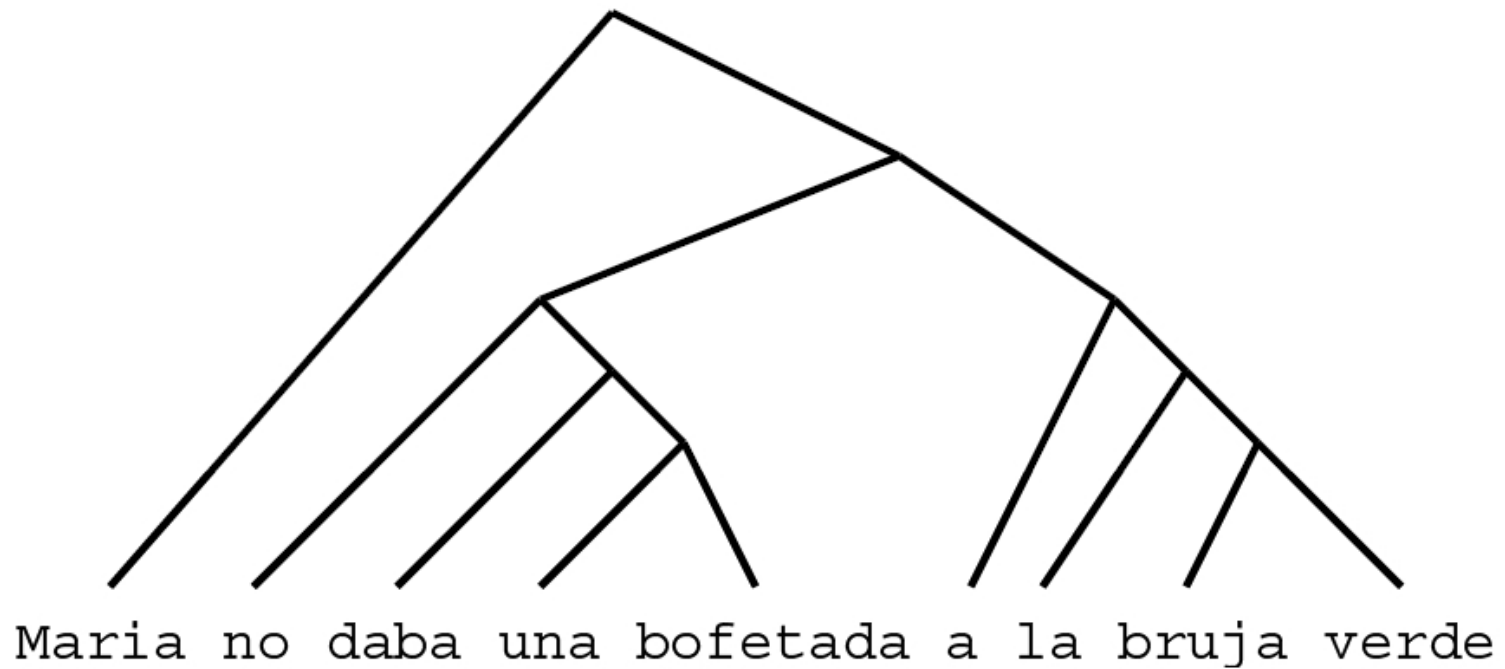
---



---

## Example: Unlabeled binary tree (Spanish)

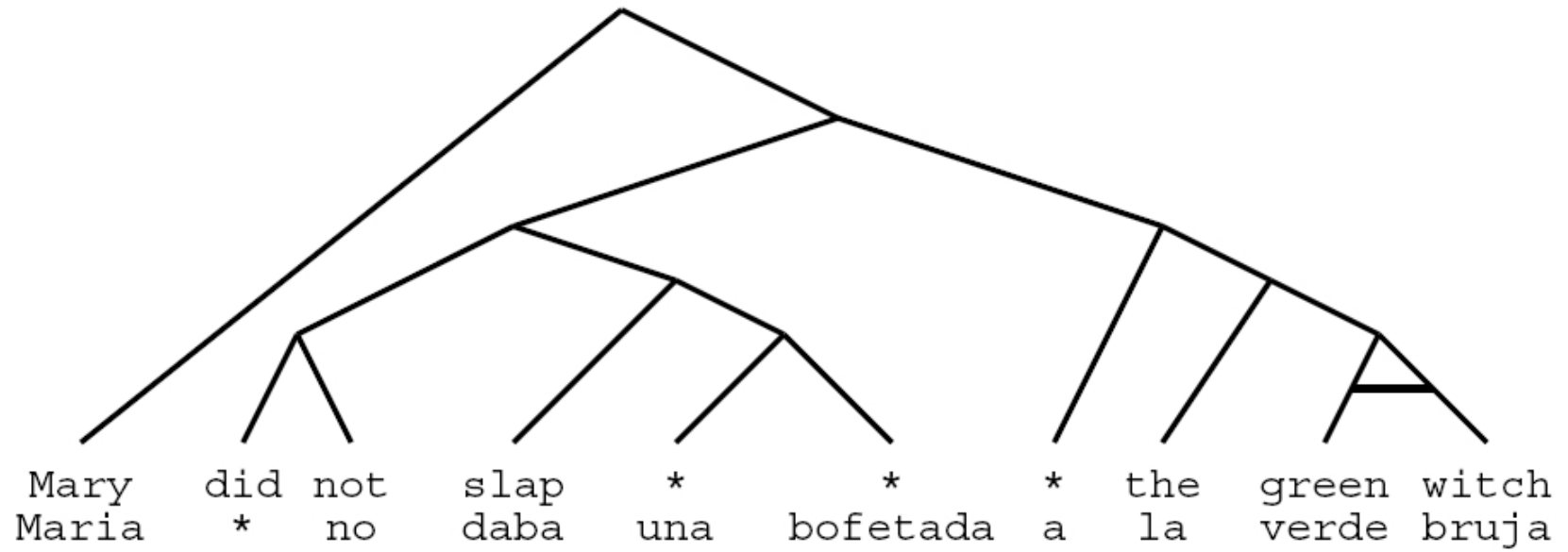
---



---

## Example: ITG tree

---

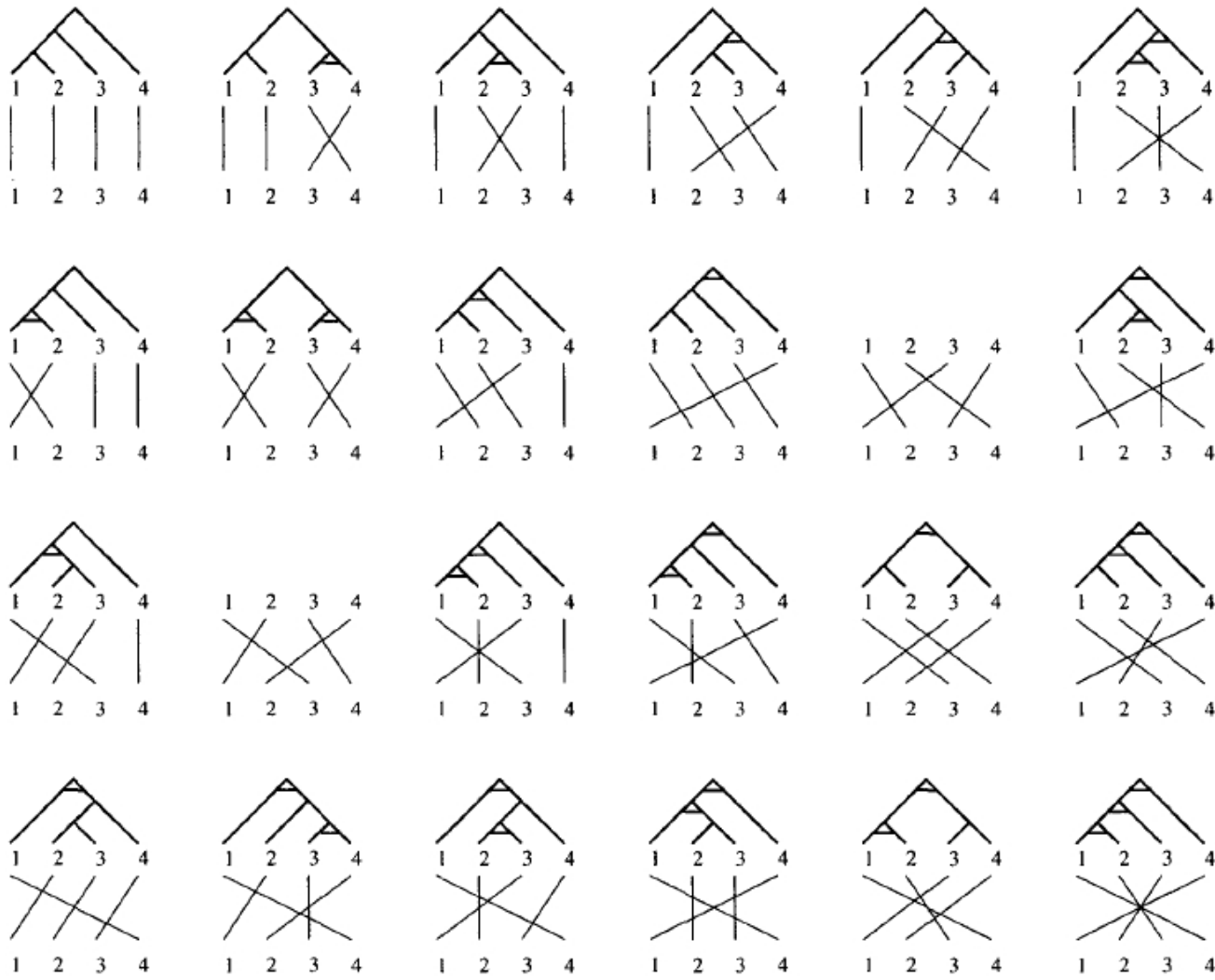




---

## Examples of ITG trees & alignments

---

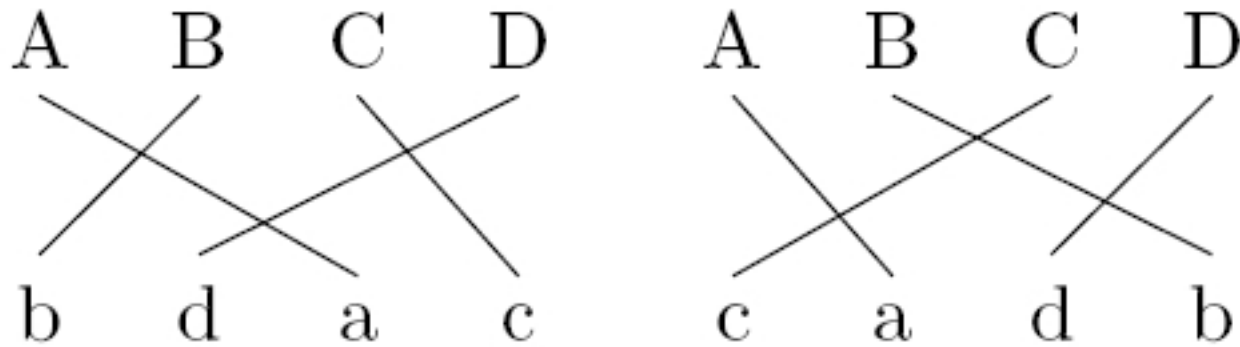




---

## Crossing alignments

---



- Impossible to represent in ITG – no corresponding trees

---

## ITG summary

---

- Drawbacks?
  - Induced structure is not (necessarily) linguistically motivated
  - Induced from parallel corpora  
(no existing human-annotated bracketed treebanks)
  - Similarity of two different languages' grammatical structure
- Benefits?
  - Adds *some* structure to translation

---

## Synchronous context-free grammars (SCFG)

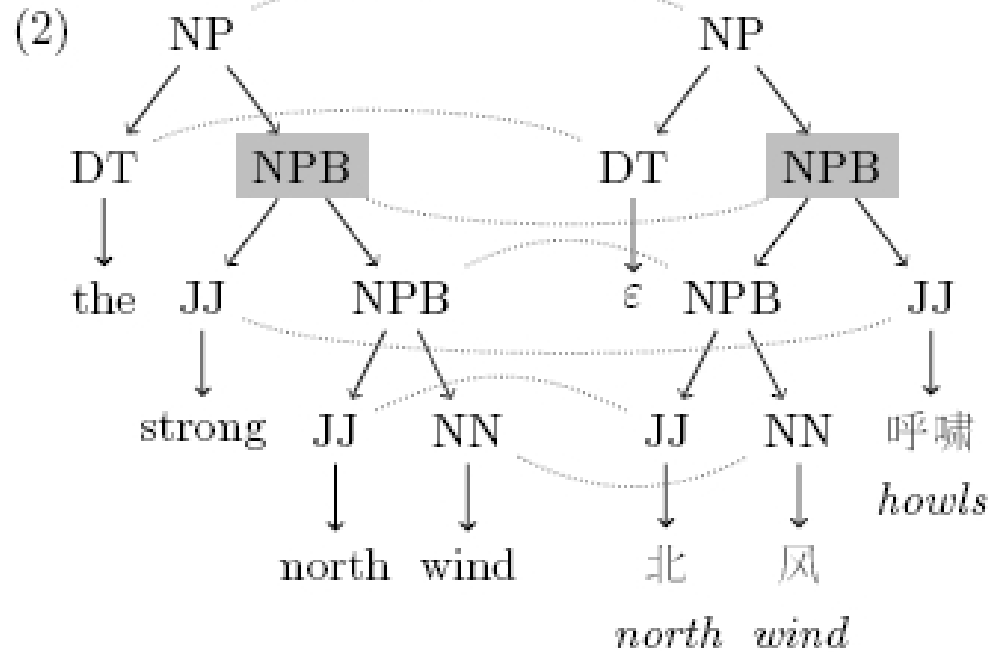
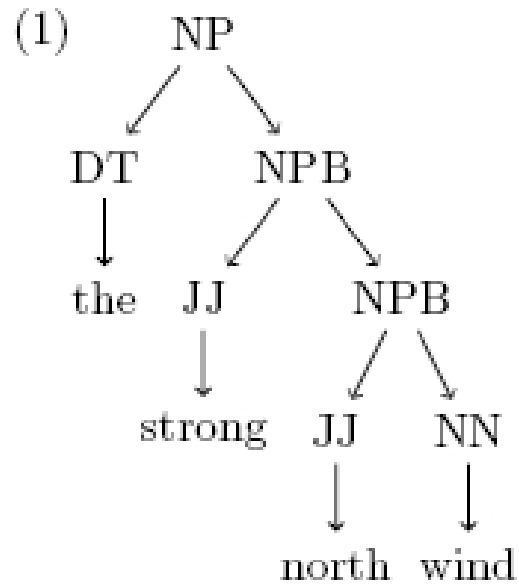
---

- Similar to ITG...
  - Binary rules
  - Inversion productions
- Size of the non-terminal set  $> 1$

a.k.a., labeled nodes on the parse trees

$$\begin{array}{ll} \text{NP} \rightarrow [ \text{DT NPB} ] & \text{NPB} \rightarrow \langle \text{NPB JJ} \rangle \\ \text{NPB} \rightarrow [ \text{JJ NN} ] & \text{DT} \rightarrow \text{the} / \epsilon \end{array}$$

## Example SCFG pair of trees



---

## SCFG summary

---

- Drawbacks?
  - No existing (human-annotated) corpora
  - Relies strongly on similarity of two languages' grammatical structure
- Benefits?
  - Linguistically motivated productions

---

## Hierarchical synchronous context-free grammars

---

[Chiang 05]

- Context-free bi-grammar
- Single non-terminal symbol ( $X$ )
- RHS of rules may include *both* non-terminal and terminals (words)
  - $X \rightarrow X_1 s_1 X_2 / t_1 X_1 X_2$
- Makes translating equivalent to parsing (though  $O(n^6)$ )

---

## Types of hierarchical translation rules

---

Non-terminal rules:

$S \rightarrow S X / S X$

$X \rightarrow X / X$

Terminal rules:

house / casa

blue / bleu

$X \rightarrow \text{slap} / \text{daba una bofetada}$

Mixed non-terminal/terminal:

$X \rightarrow \text{not } X / \text{ne } X \text{ pas}$

$X \rightarrow X_1 \text{ 's } X_2 / X_2 \text{ de } X_1$

$X \rightarrow \text{green } X / X \text{ verde}$

---

## Hierarchical SCFG summary

---

- Drawbacks?
  - Productions are not linguistically-motivated
  - High computational cost
- Benefits?
  - Allows for discontinuous phrases



---

## Machine translation

---

- Word-based
- Phrase-based
- “Syntax”-based

---

## Example translation

---

However , the sky remained clear under the strong north wind .

虽然 北 风 呼啸 ， 但 天空 依然 十分 清澈 。

*Although north wind howls , but sky still extremely limpid .*

---

## Word-based alignment

---

However , the sky remained clear under the strong north wind .

虽然 北 风 呼啸 ， 但 天空 依然 十分 清澈 。

*Although north wind howls , but sky still extremely limpid .*

---

## Word-based translation

---

- English  $\rightarrow$  Chinese:

However  $\rightarrow$  *Although*<sub>1</sub> *but*<sub>6</sub>

,  $\rightarrow$  ,<sub>5</sub>

the  $\rightarrow$   $\epsilon$

sky  $\rightarrow$  *sky*<sub>7</sub>

remained  $\rightarrow$  *still*<sub>8</sub>

clear  $\rightarrow$  *limpid*<sub>10</sub>

under  $\rightarrow$   $\epsilon$

the  $\rightarrow$   $\epsilon$

strong  $\rightarrow$  *howls*<sub>4</sub>

north  $\rightarrow$  *north*<sub>2</sub>

wind  $\rightarrow$  *wind*<sub>3</sub>

.  $\rightarrow$  .<sub>11</sub>

$\epsilon$   $\rightarrow$  *extremely*<sub>9</sub>

- Benefits? Drawbacks?

---

## Phrase-based alignment

---

1. Divide “source” sentence into phrases (how to choose?):

[However] [,] [the sky remained clear] [under the strong north wind] [.]

2. Translate each phrase (from a look-up table):

*[Although] [, but] [sky still extremely limpid] [north wind howls] [.]*

3. Rearrange phrases (using Language Modeling?):

*[Although] [north wind howls] [, but] [sky still extremely limpid] [.]*

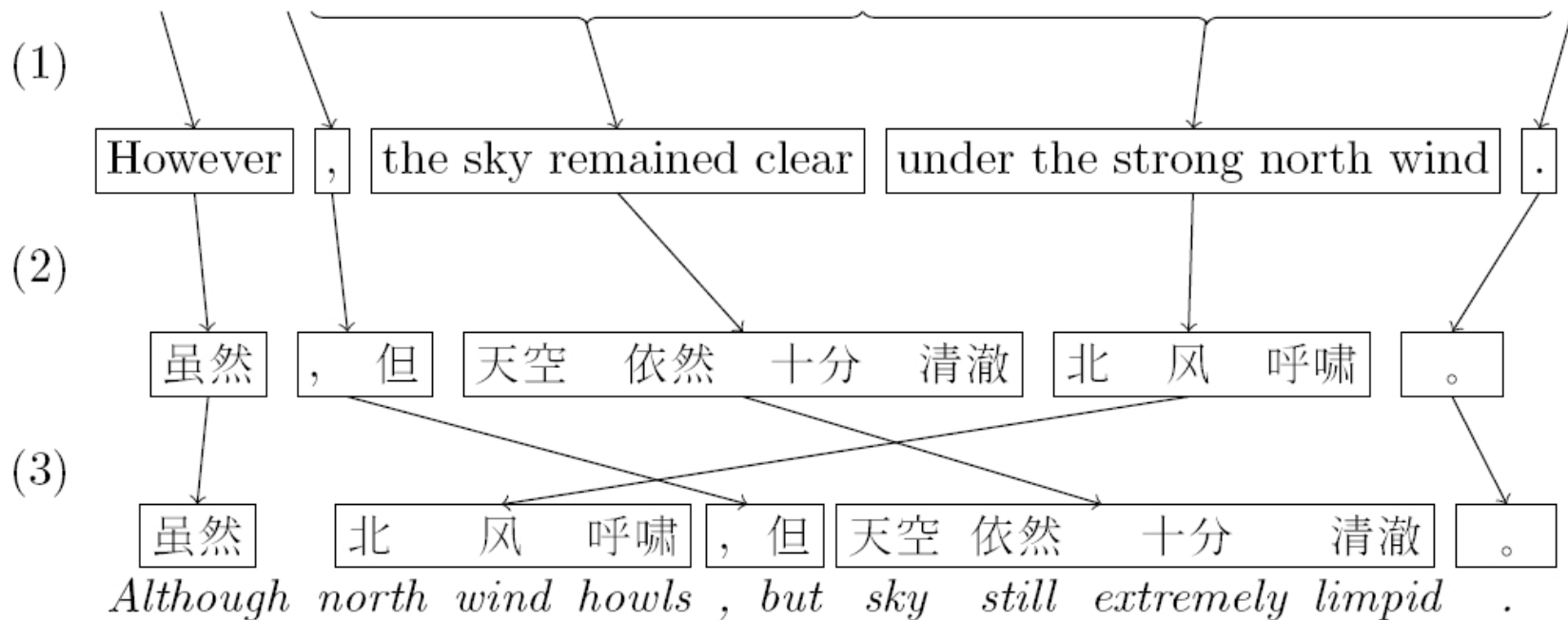
---

## Phrase-based translation

---

- English → Chinese

However , the sky remained clear under the strong north wind .

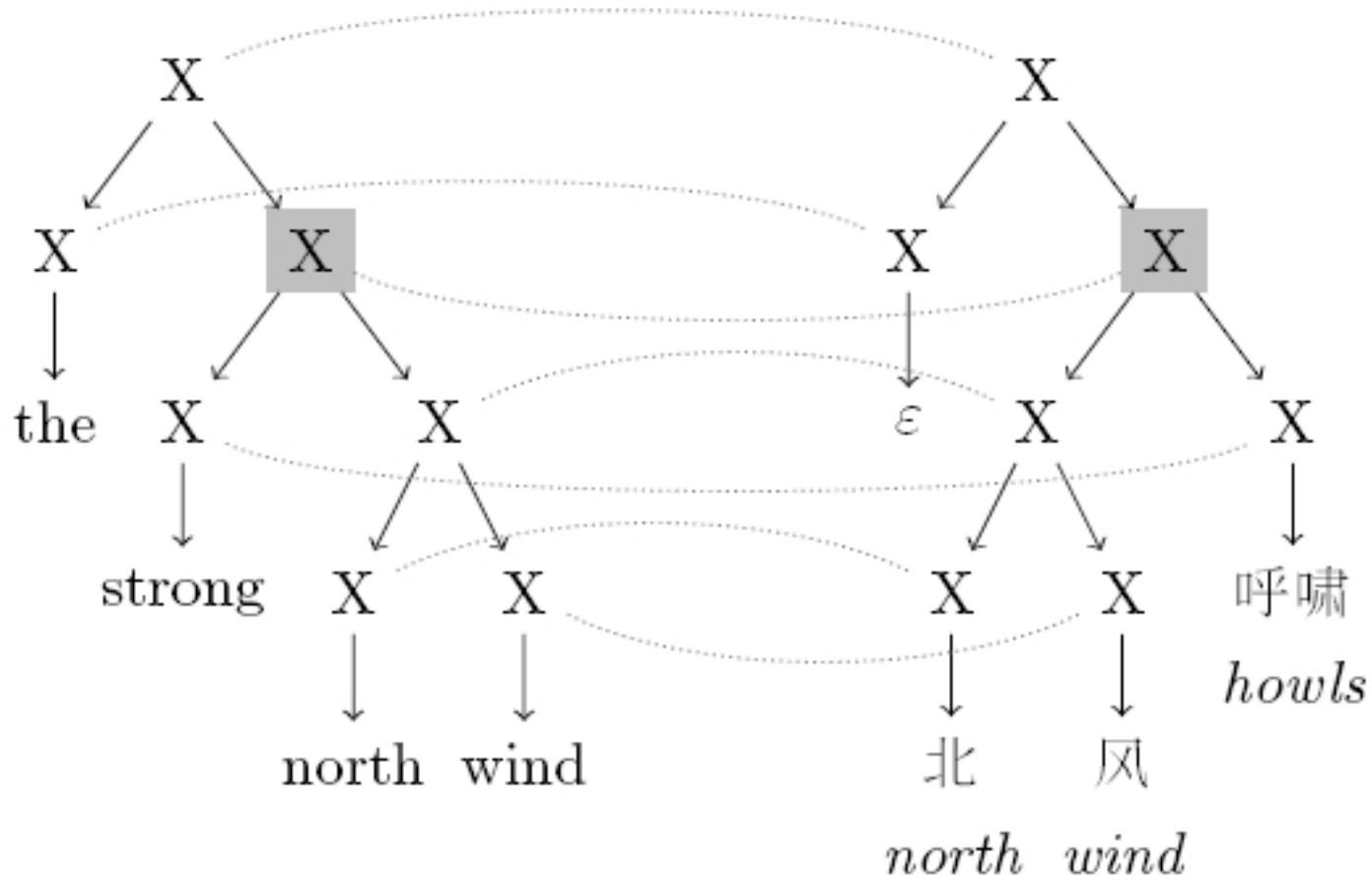


- Benefits? Drawbacks?

---

## ITG alignment

---



---

## ITG translation (derivation)

---

- English  $\rightarrow$  Chinese

$X \rightarrow [ X X ]$

$X \rightarrow \langle X X \rangle$

$X \rightarrow \text{the} / \epsilon$

$X \rightarrow \text{strong} / \textit{howls}$

$X \rightarrow \text{north} / \textit{north}$

$X \rightarrow \text{wind} / \textit{wind}$

- Benefits? Drawbacks?



---

## SCFG grammar

---

$NP \longrightarrow DT_{[1]}NPB_{[2]} / DT_{[1]}NPB_{[2]}$

$NPB \longrightarrow JJ_{[1]}NN_{[2]} / JJ_{[1]}NN_{[2]}$

$NPB \longrightarrow NPB_{[1]}JJ_{[2]} / JJ_{[2]}NPB_{[1]}$

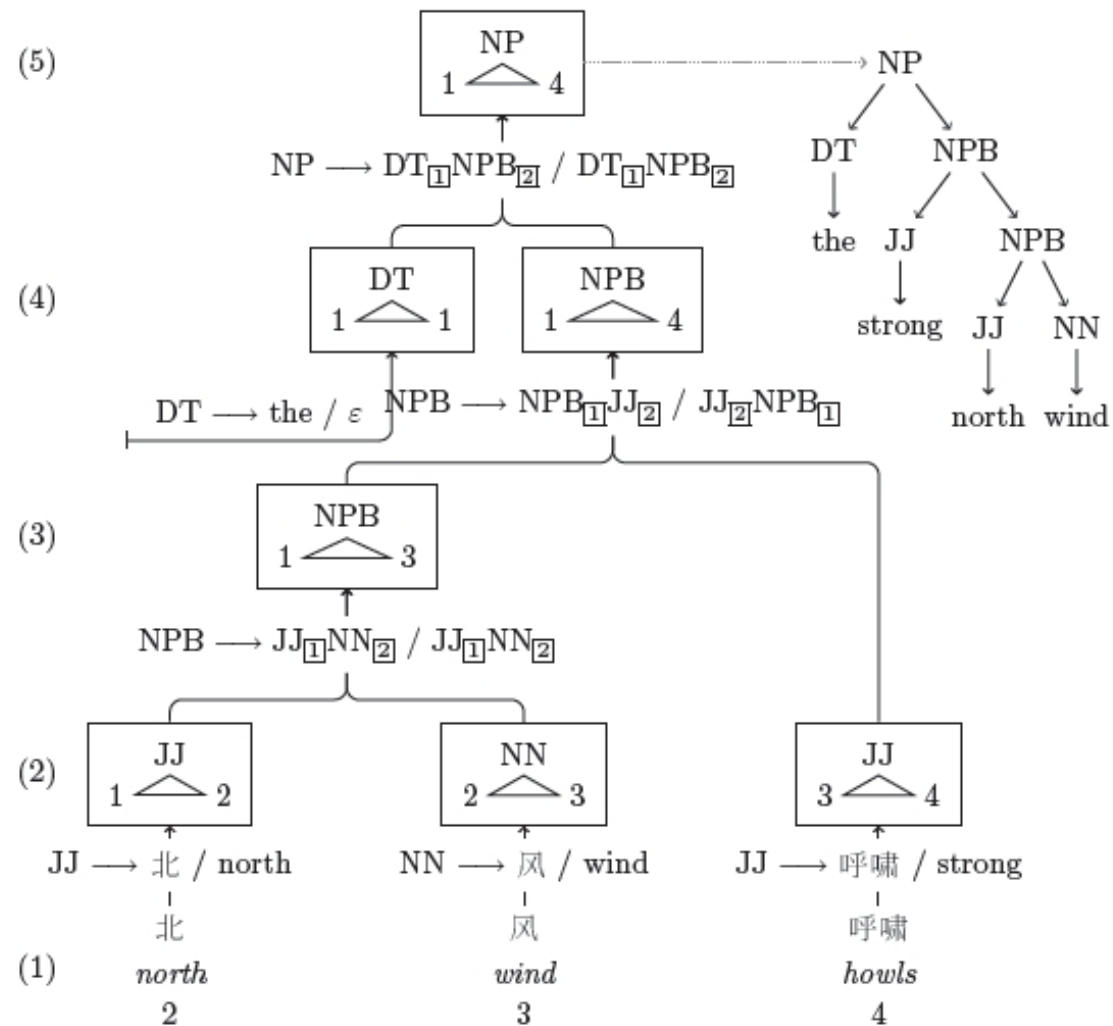
$DT \longrightarrow \text{the} / \varepsilon$

$JJ \longrightarrow \text{strong} / \text{呼啸}$

$JJ \longrightarrow \text{north} / \text{北}$

$NN \longrightarrow \text{wind} / \text{风}$

# SCFG translation (derivation)



---

## Hierarchical SCFG grammar

---

$X \rightarrow$  However ,  $X_{[1]}X_{[2]}$  . / 虽然  $X_{[2]}$  , 但  $X_{[1]}$  。

$X \rightarrow$  under the strong north wind / 北 风 呼 啸

$X \rightarrow$  the sky remained clear / 天 空 依 然 十 分 清 澈

- $X \rightarrow$  However , under the strong north wind the sky remained clear . /

*Although north wind howls , but sky still extremely limp .*

- $X \rightarrow$  However , the sky remained clear under the strong north wind . /

*Although sky still extremely limp , but north wind howls .*

- Benefits? Drawbacks?

---

## Summary

---

- Bi-text parsing:
  - like finite-state transducers for trees
- Grammars provide more powerful transforms
- Higher complexity
- Utility of structure in translation?
  - ...syntax in translation?