

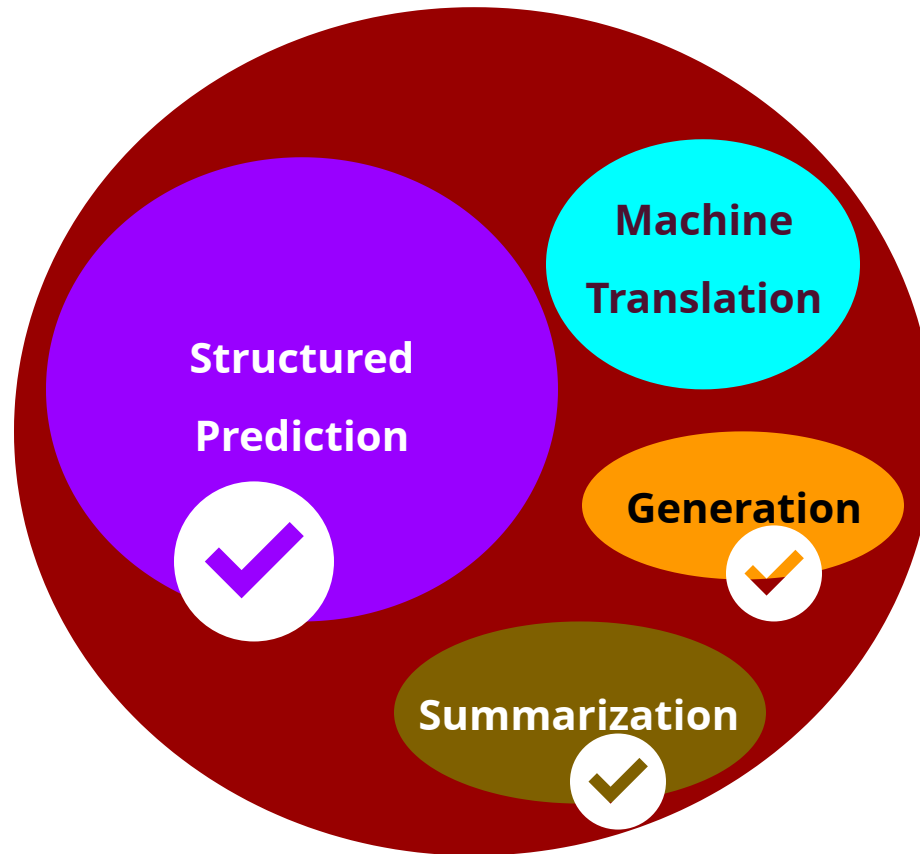
50.040

Natural Language Processing

Lu, Wei

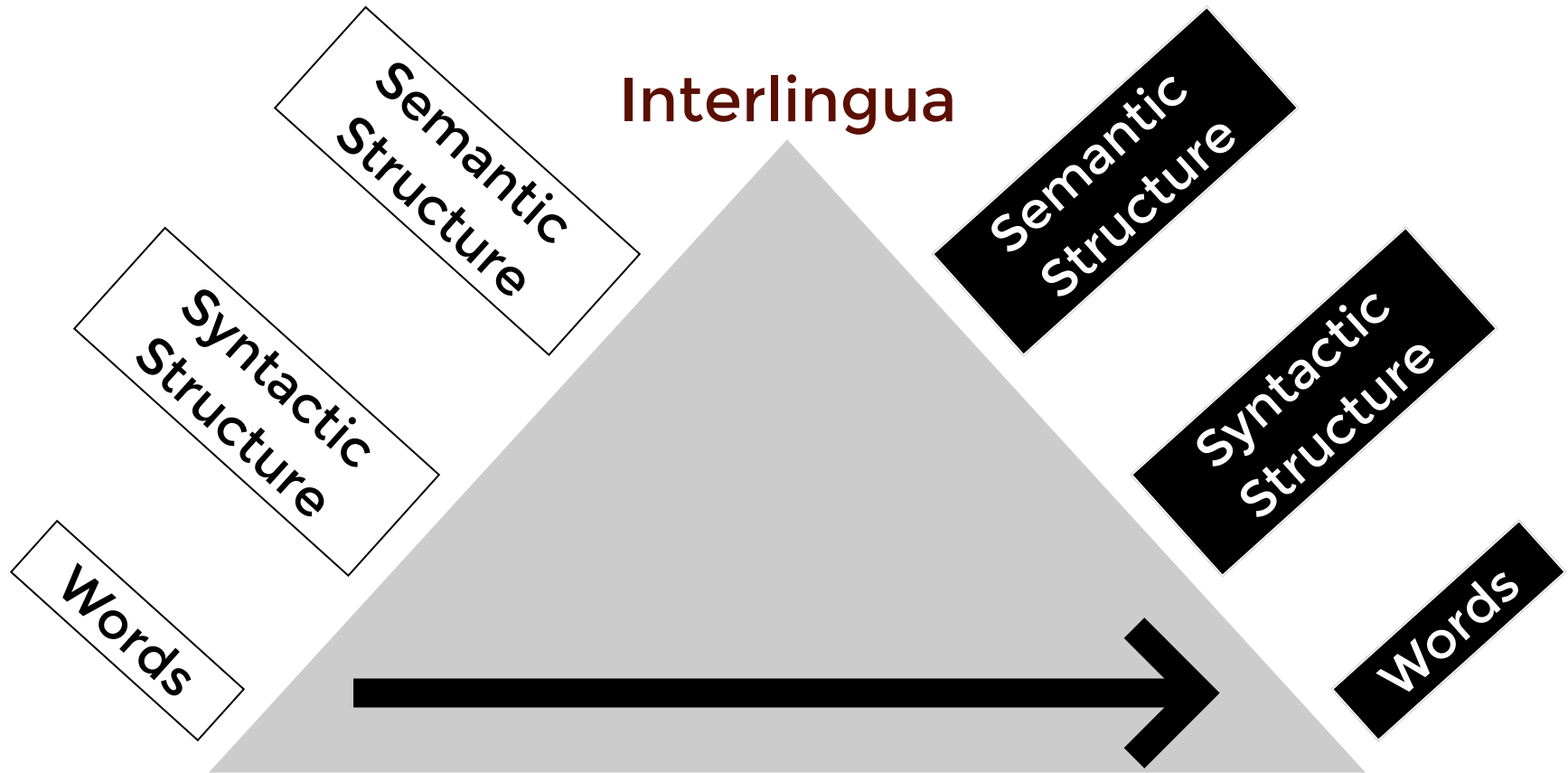


Tasks in NLP



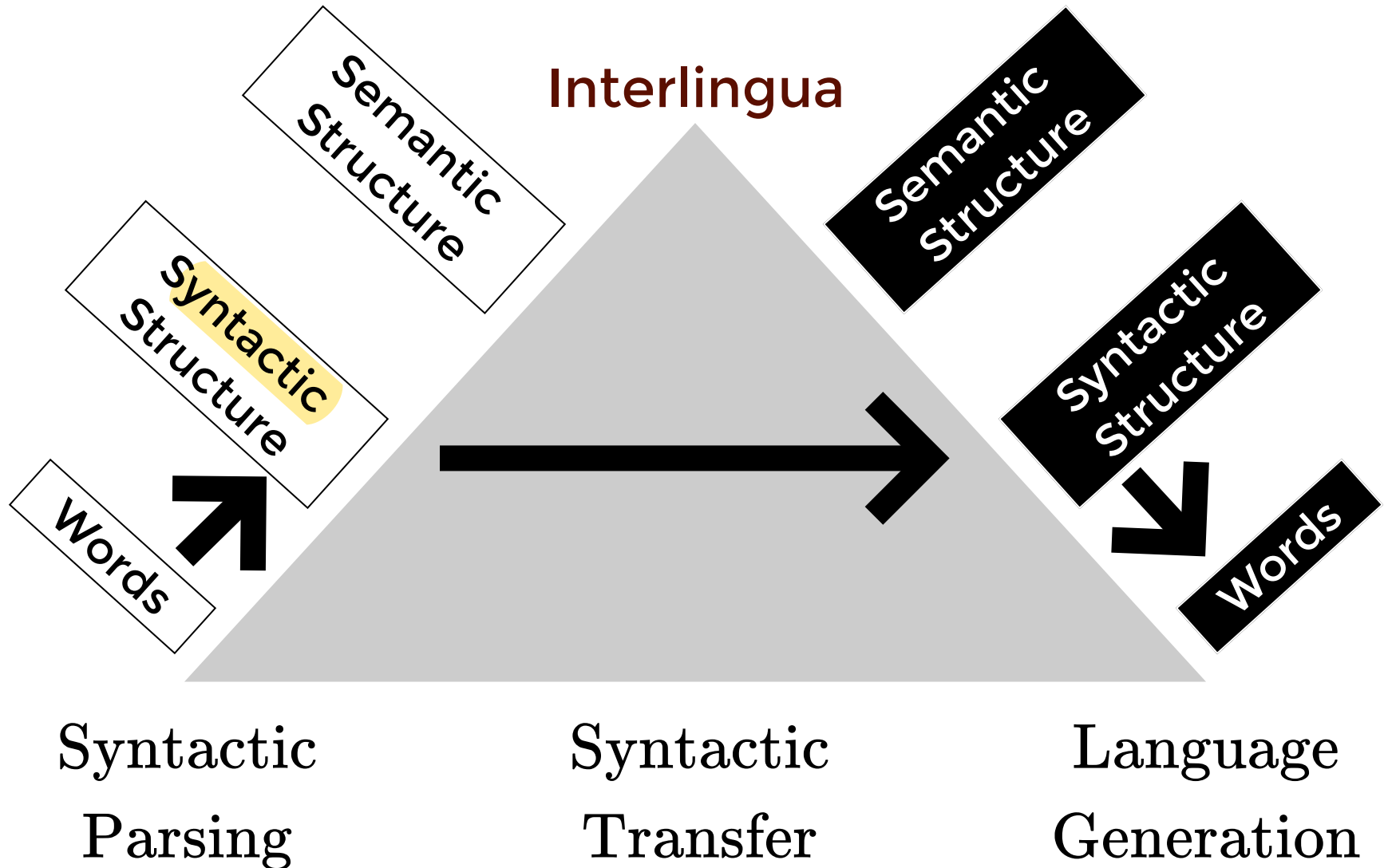
Supervised

Machine Translation

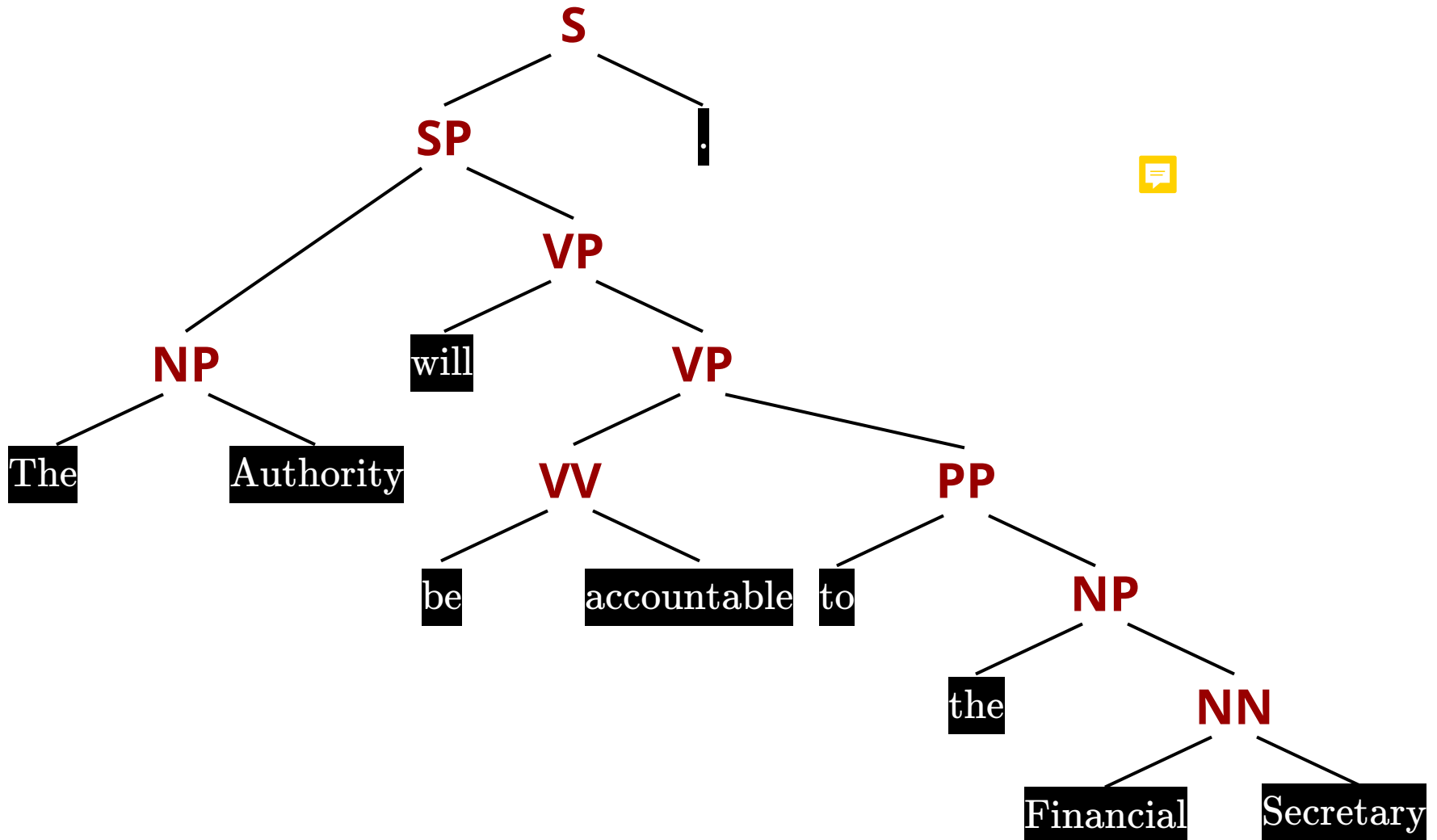


Text-to-text Problem

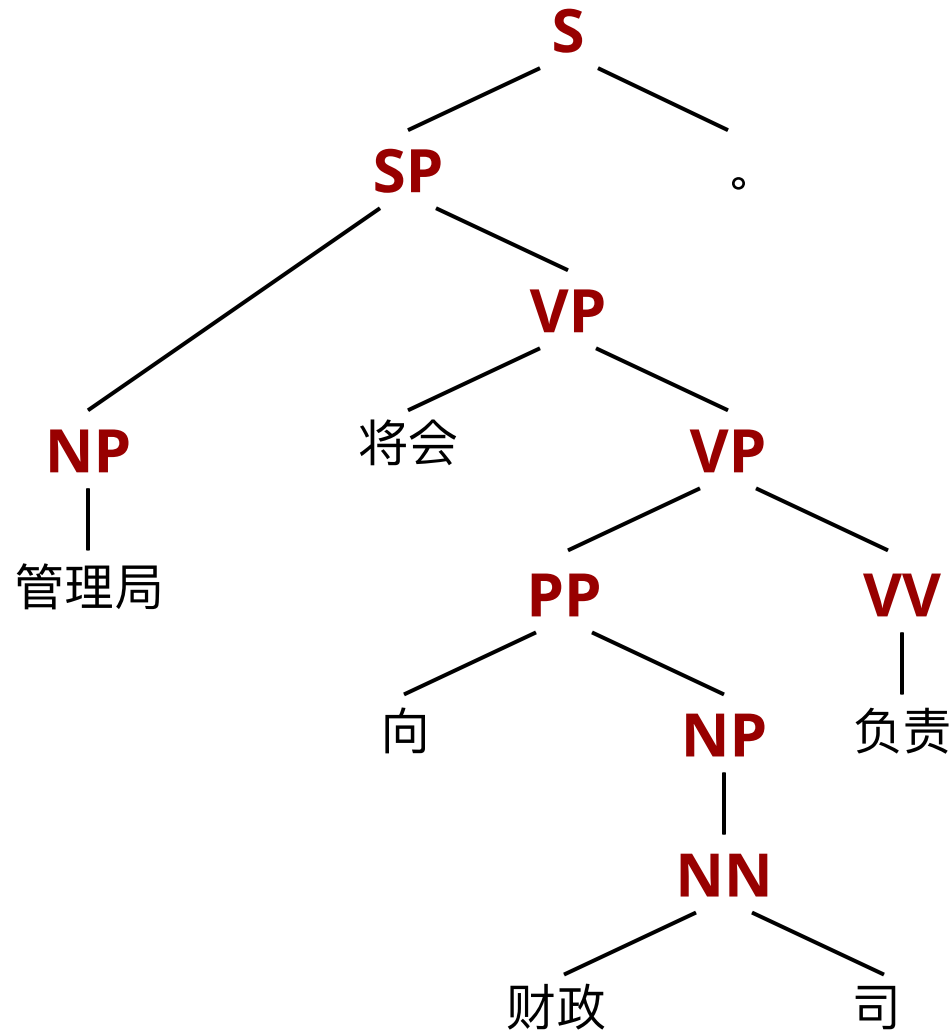
Machine Translation



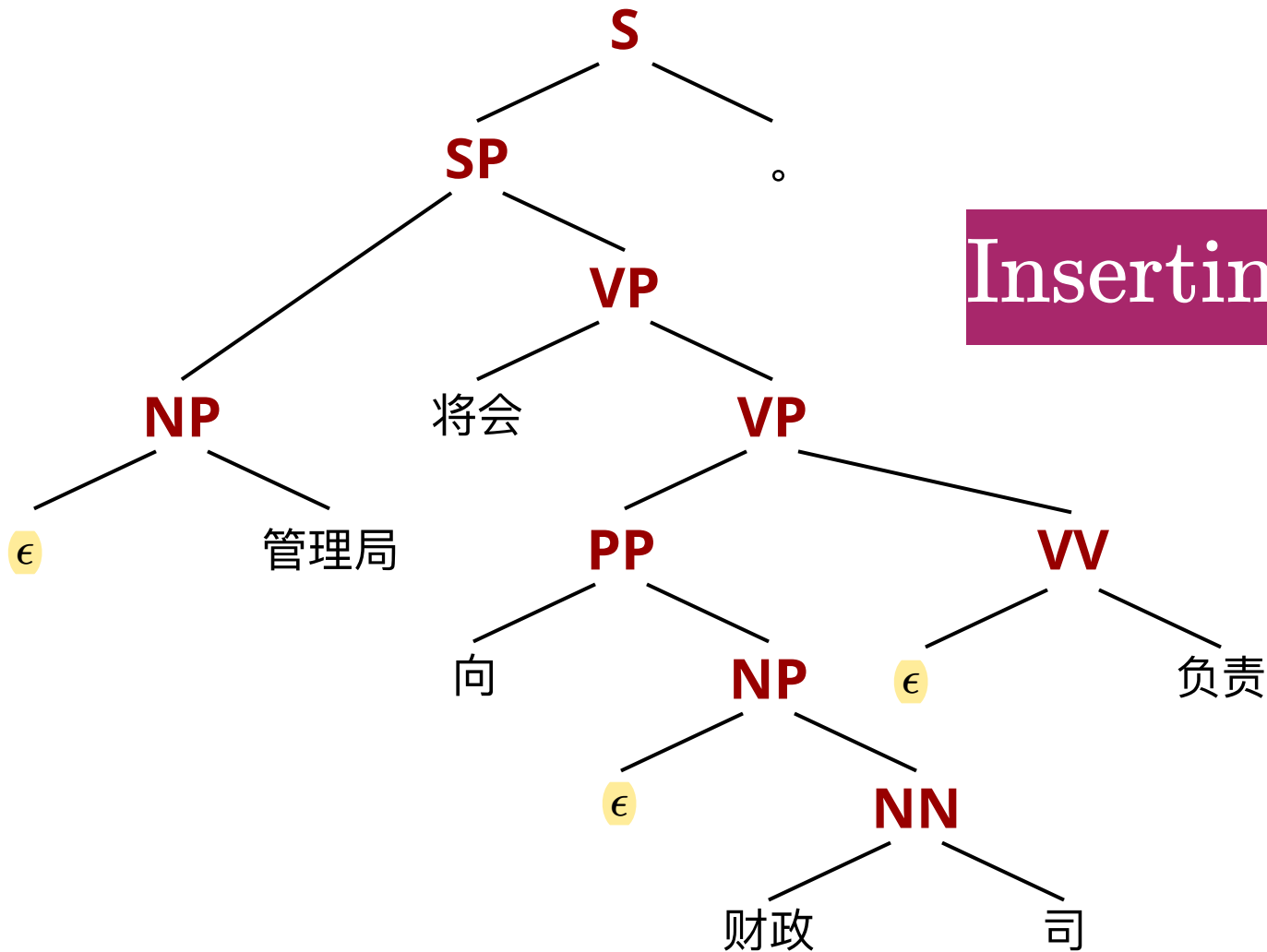
Target Parse Tree



Source Parse Tree

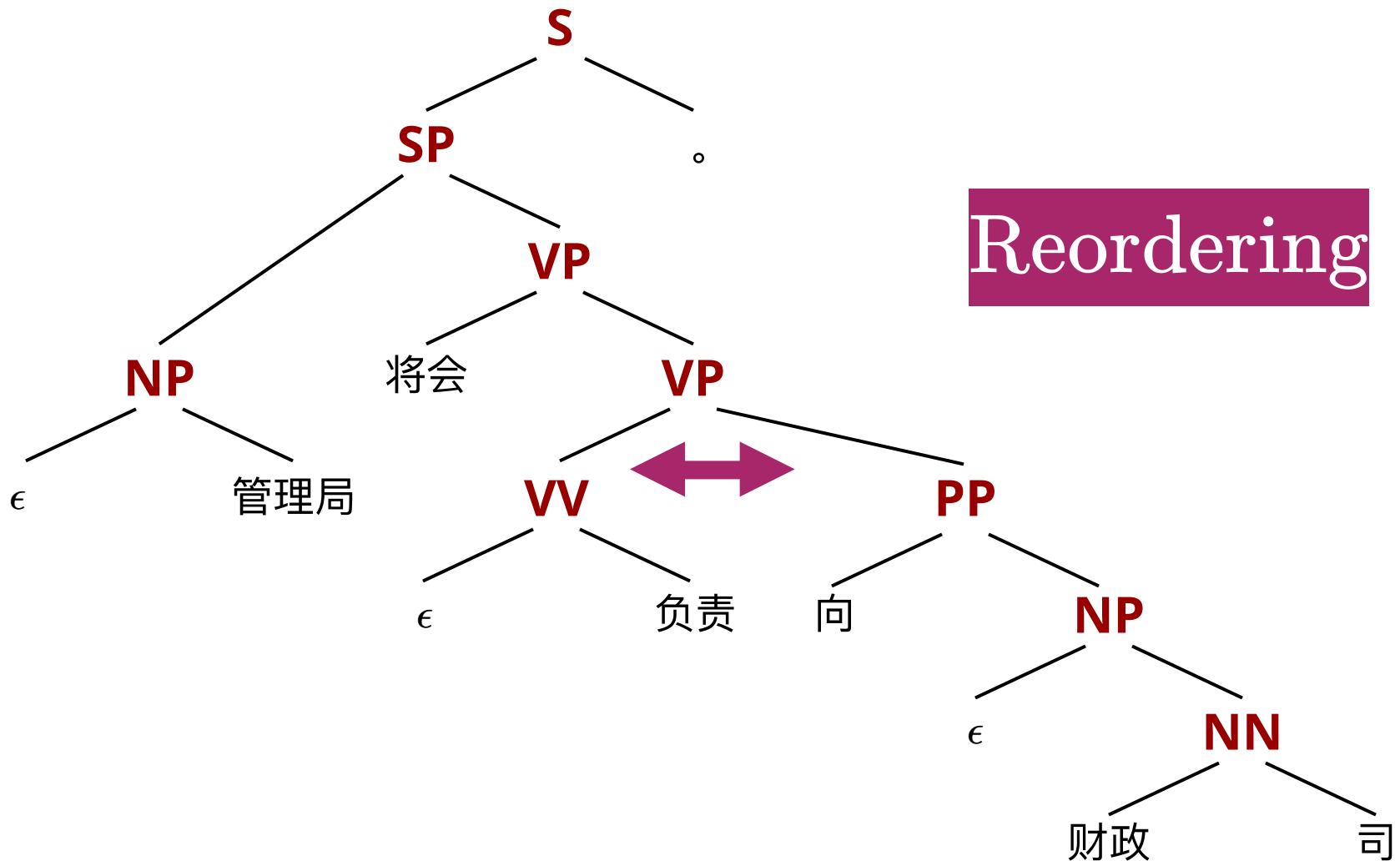


Modified Source Tree

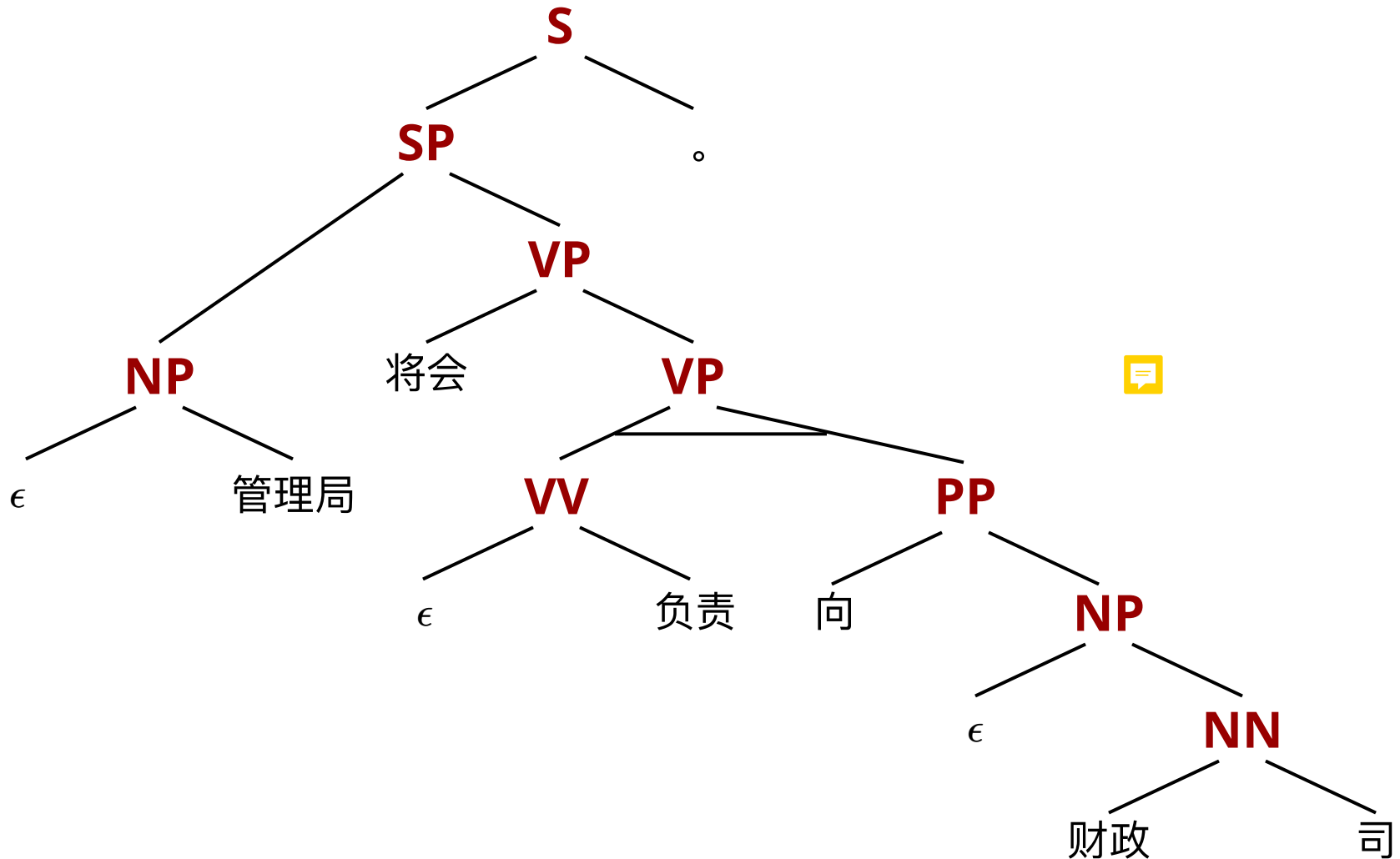


Inserting € nodes

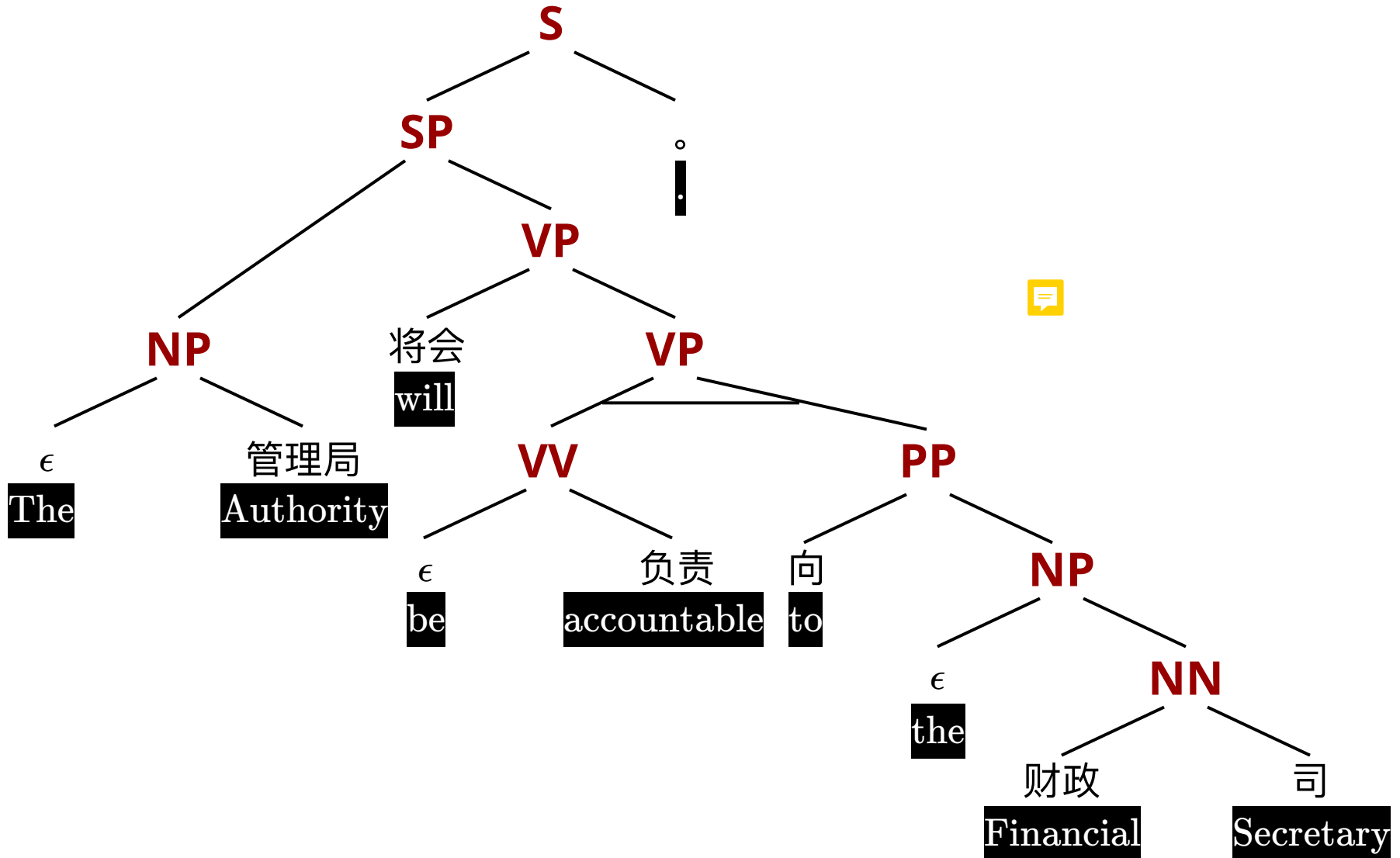
Modified Source Tree



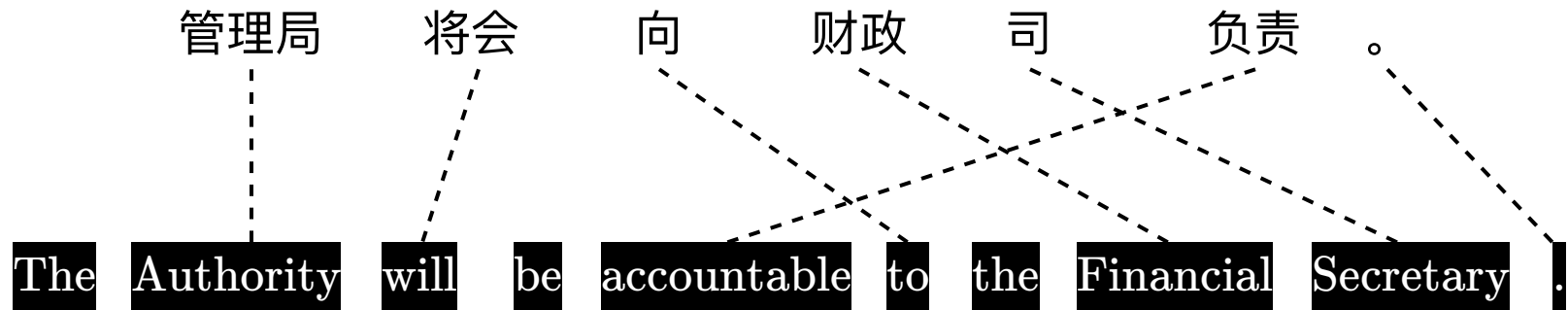
Modified Source Tree



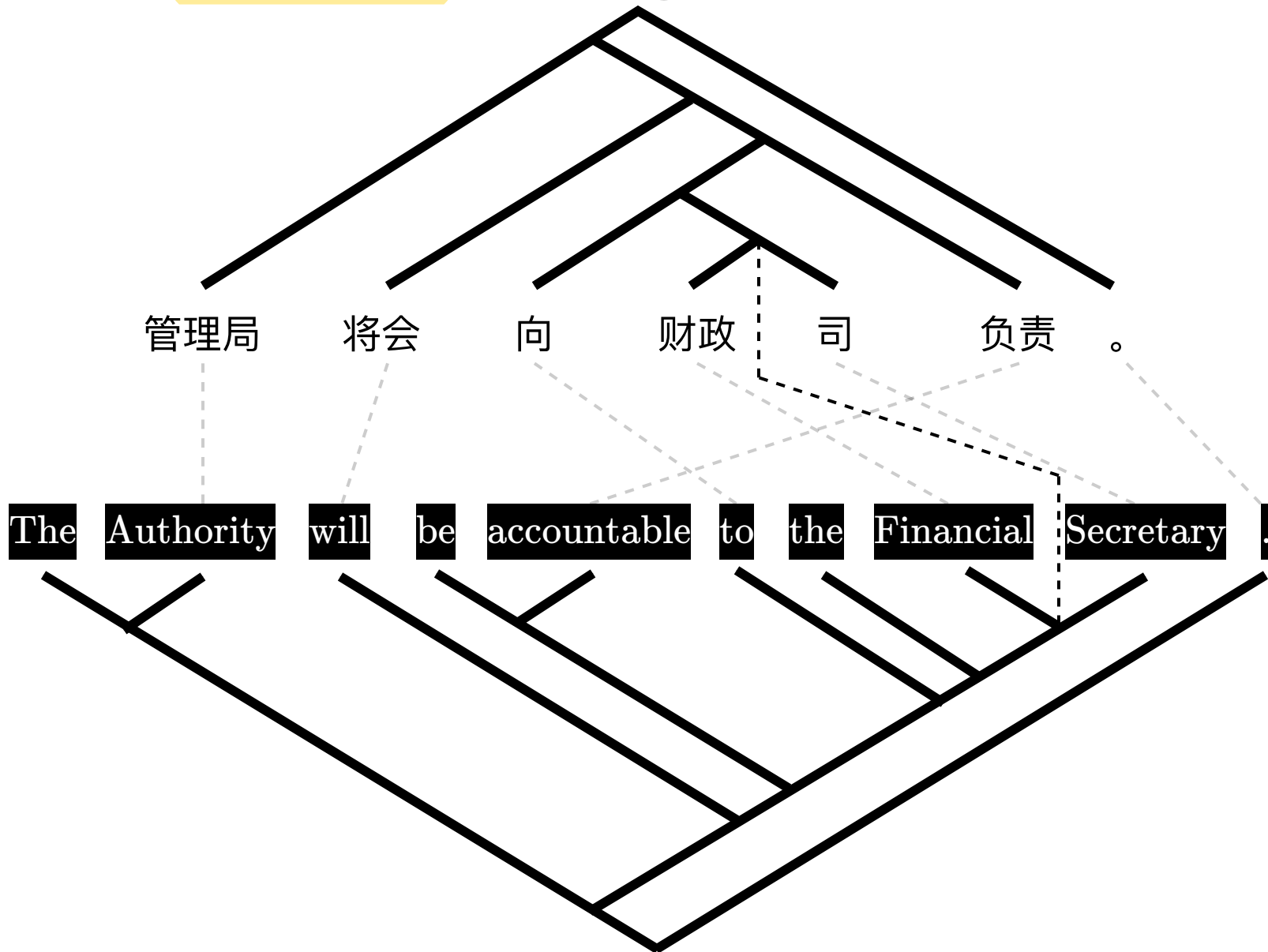
Synchronous Tree



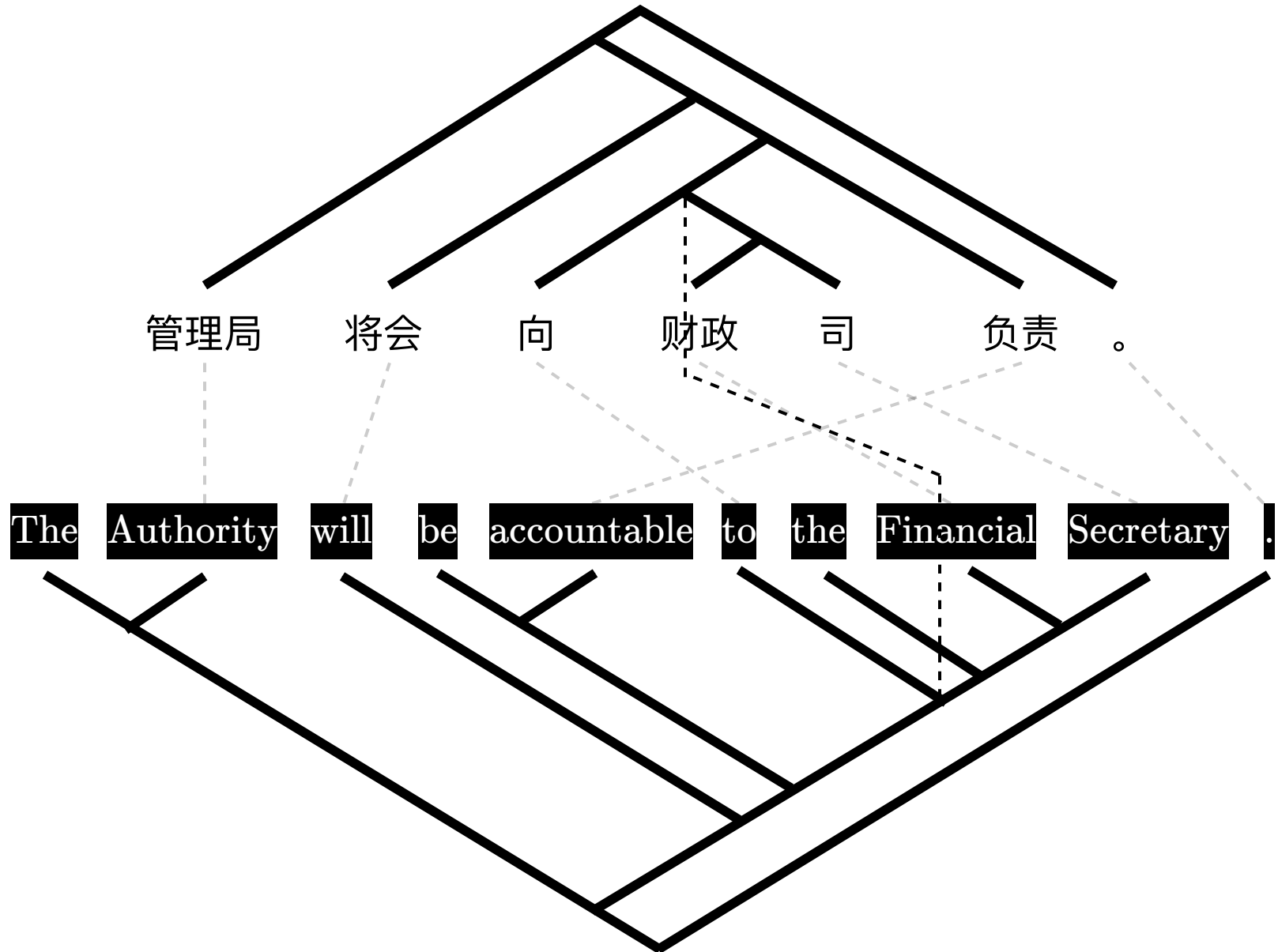
Word Alignment



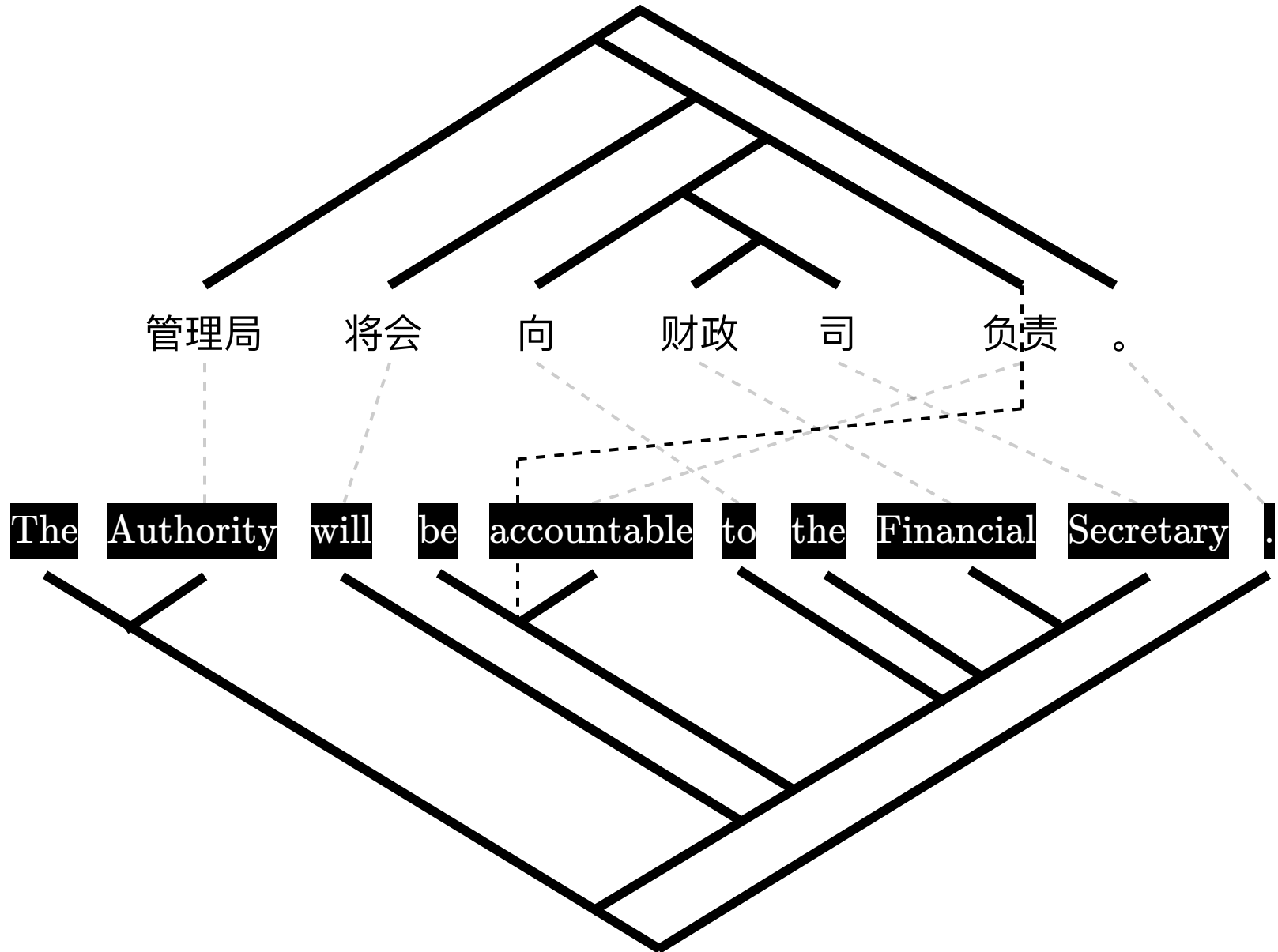
Tree Alignment



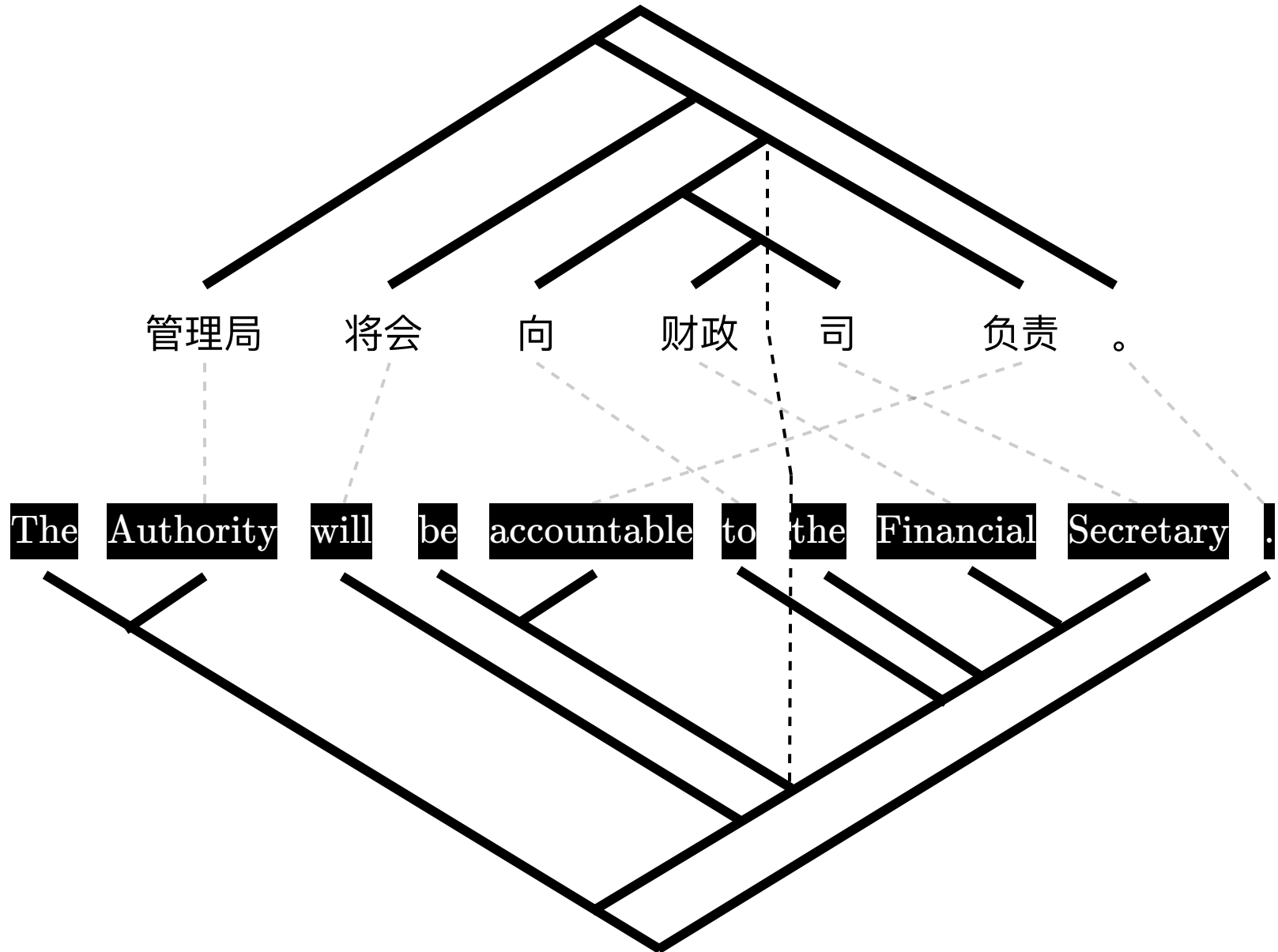
Tree Alignment



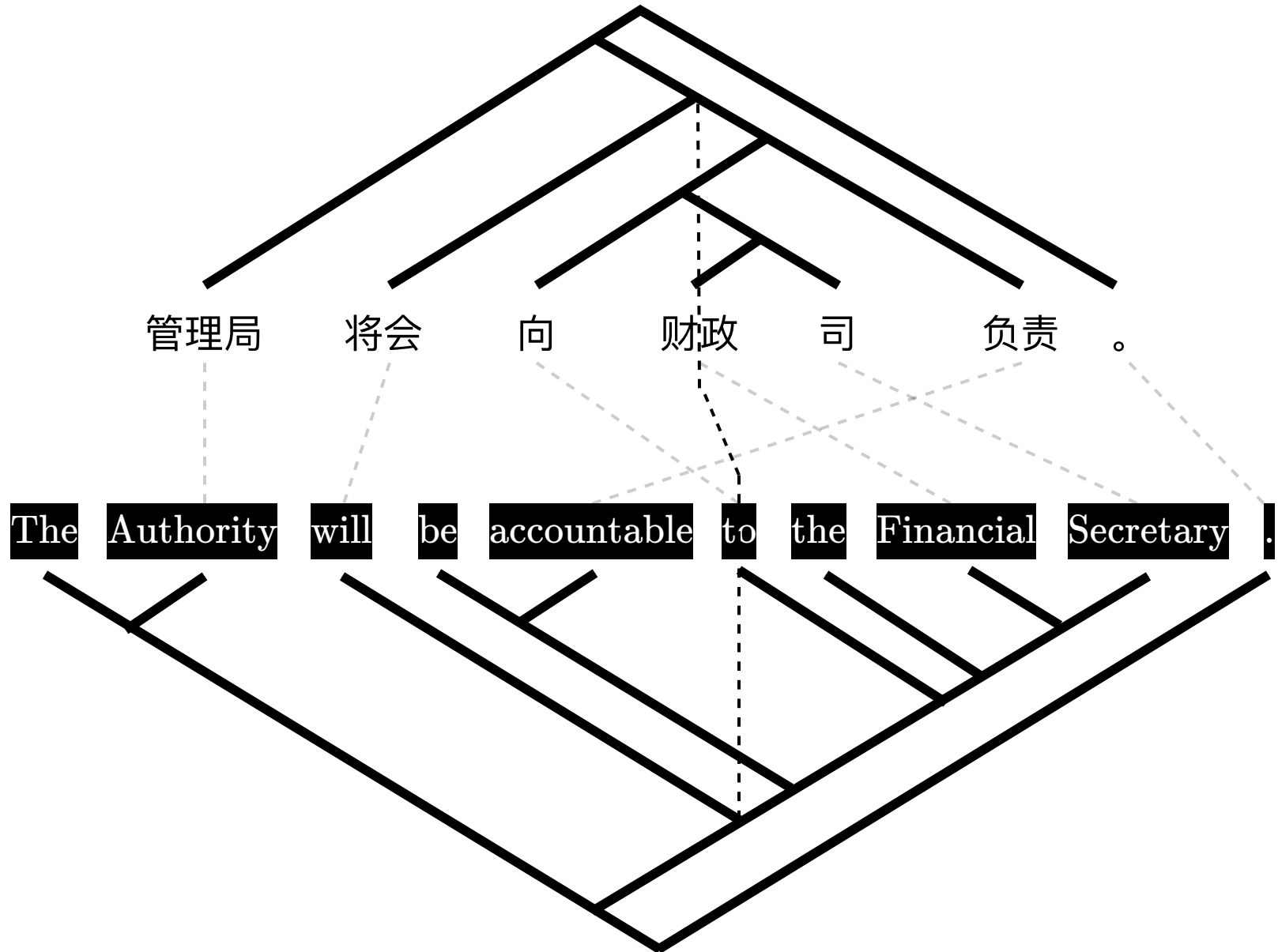
Tree Alignment



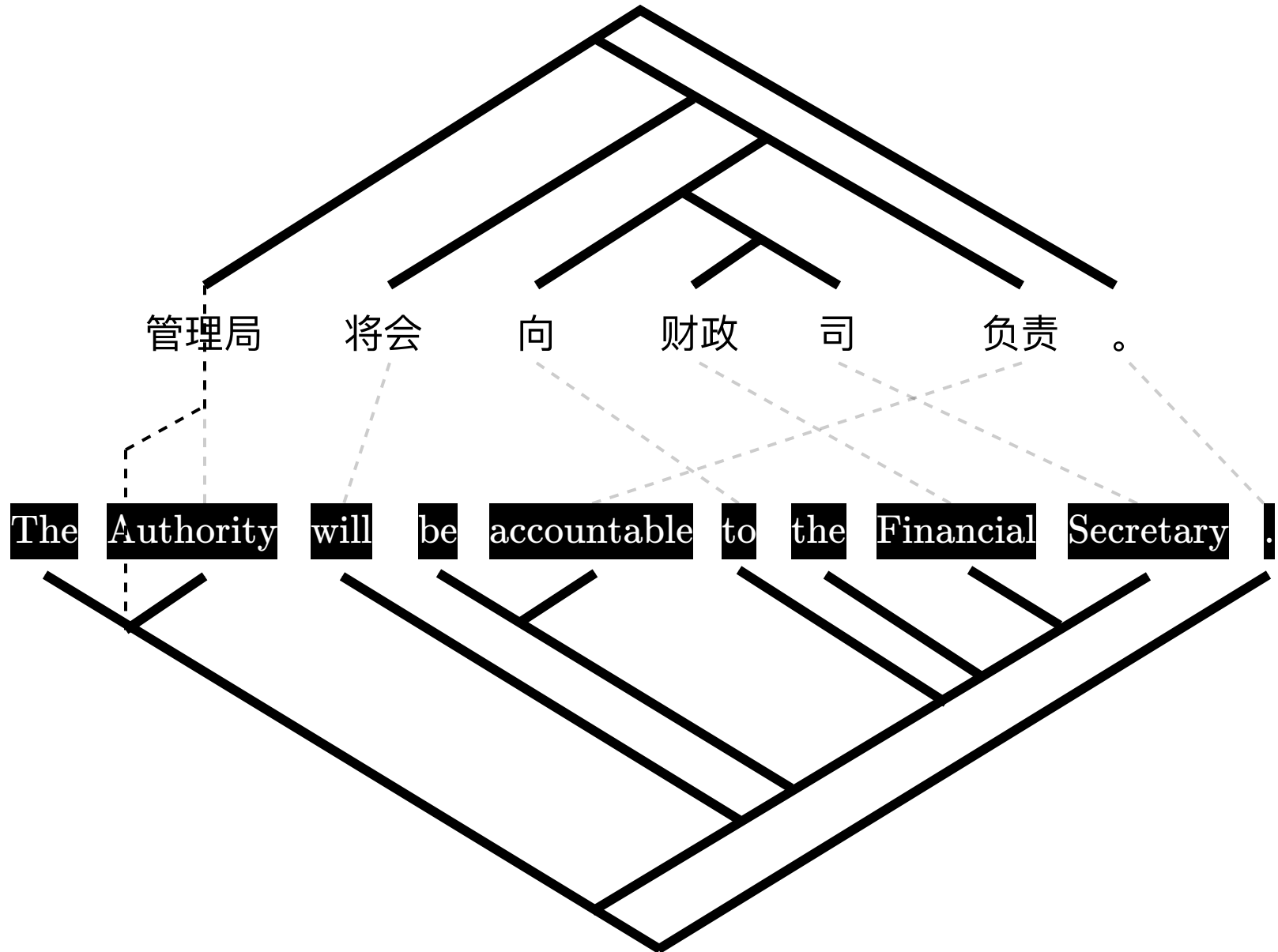
Tree Alignment



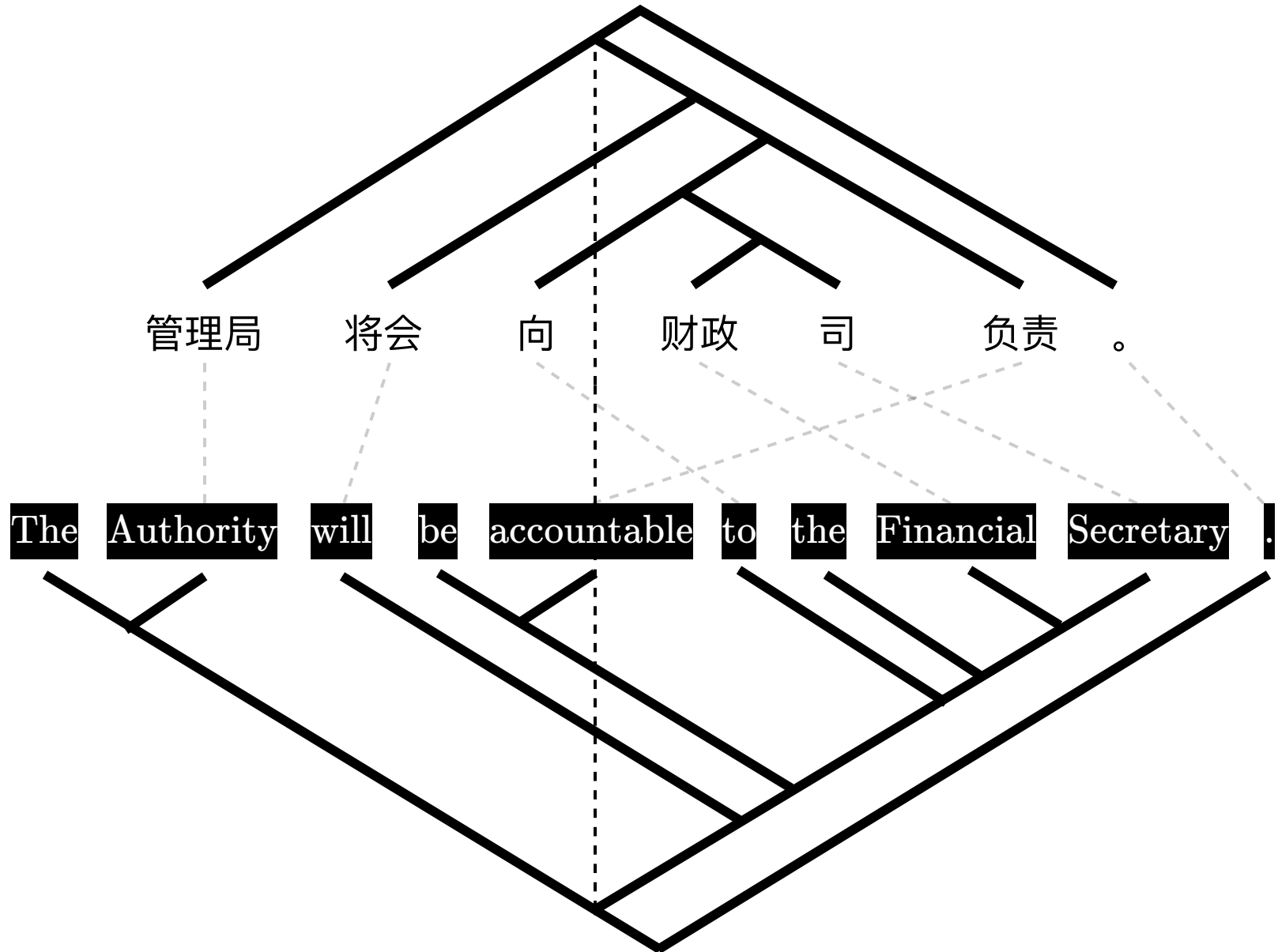
Tree Alignment



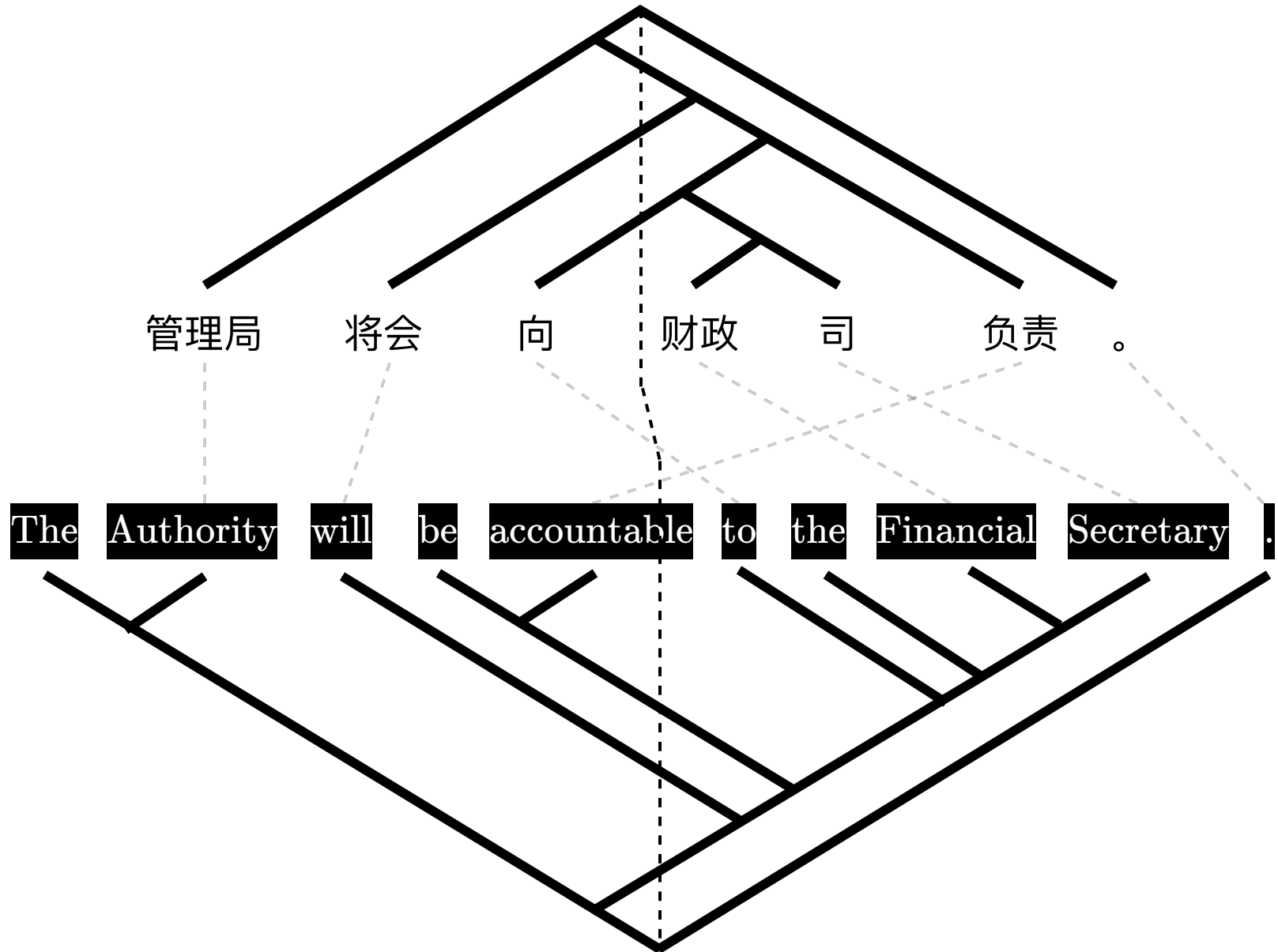
Tree Alignment



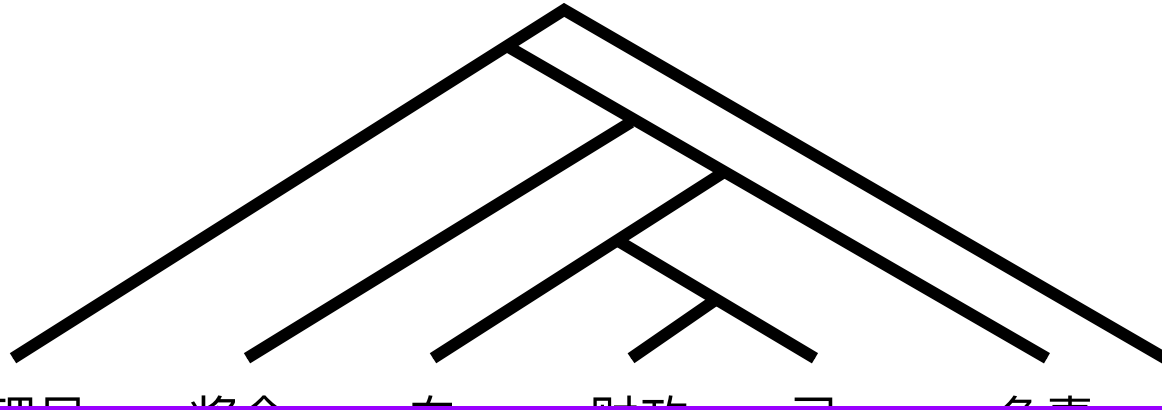
Tree Alignment



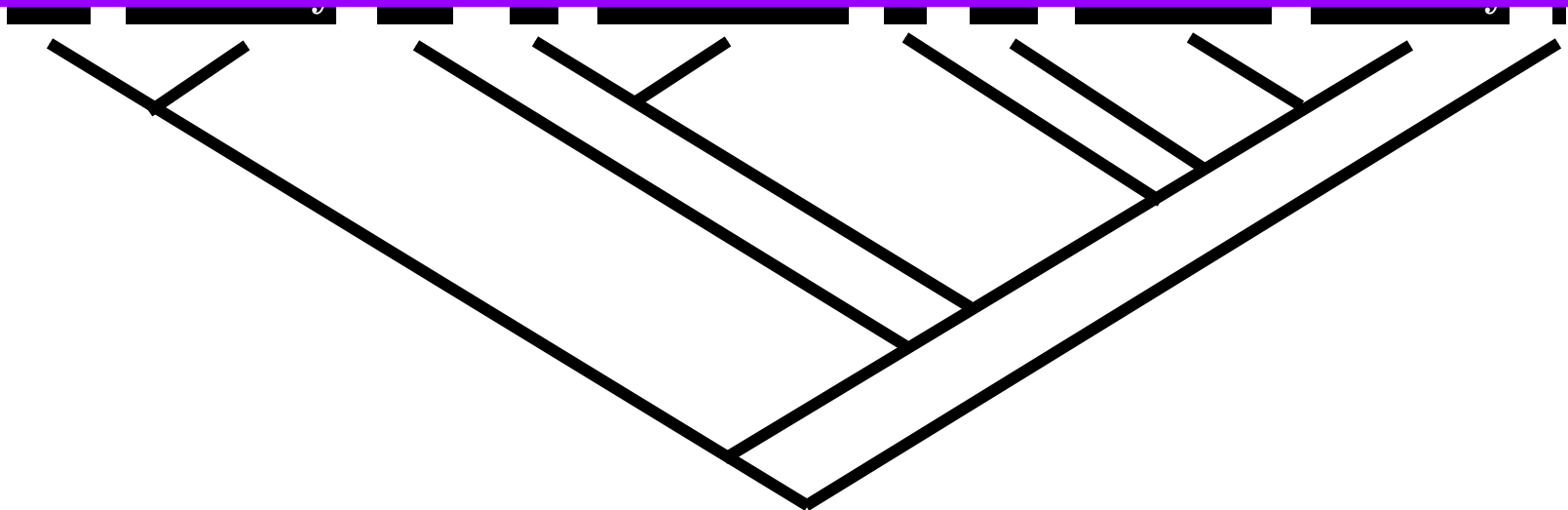
Tree Alignment



Tree Alignment



We need a formal grammar to parse both source and target trees simultaneously.



Synchronous Grammar

Stochastic Inversion Transduction Grammars and Bilingual Parsing of Parallel Corpora

Dekai Wu*
Hong Kong University of Science and Technology

We introduce (1) a novel stochastic inversion transduction grammar formalism for bilingual language modeling of sentence-pairs, and (2) the concept of bilingual parsing with a variety of parallel corpus analysis applications. Aside from the bilingual orientation, three major features distinguish the formalism from the finite-state transducers more traditionally found in computational linguistics: it skips directly to a context-free rather than finite-state base, it permits a minimal extra degree of ordering flexibility, and its probabilistic formulation admits an efficient maximum-likelihood bilingual parsing algorithm. A convenient normal form is shown to exist. Analysis of the formalism's expressiveness suggests that it is particularly well suited to modeling ordering shifts between languages, balancing needed flexibility against complexity constraints. We discuss a number of examples of how stochastic inversion transduction grammars bring bilingual constraints to bear upon problematic corpus analysis tasks such as segmentation, bracketing, phrasal alignment, and parsing.

1. Introduction

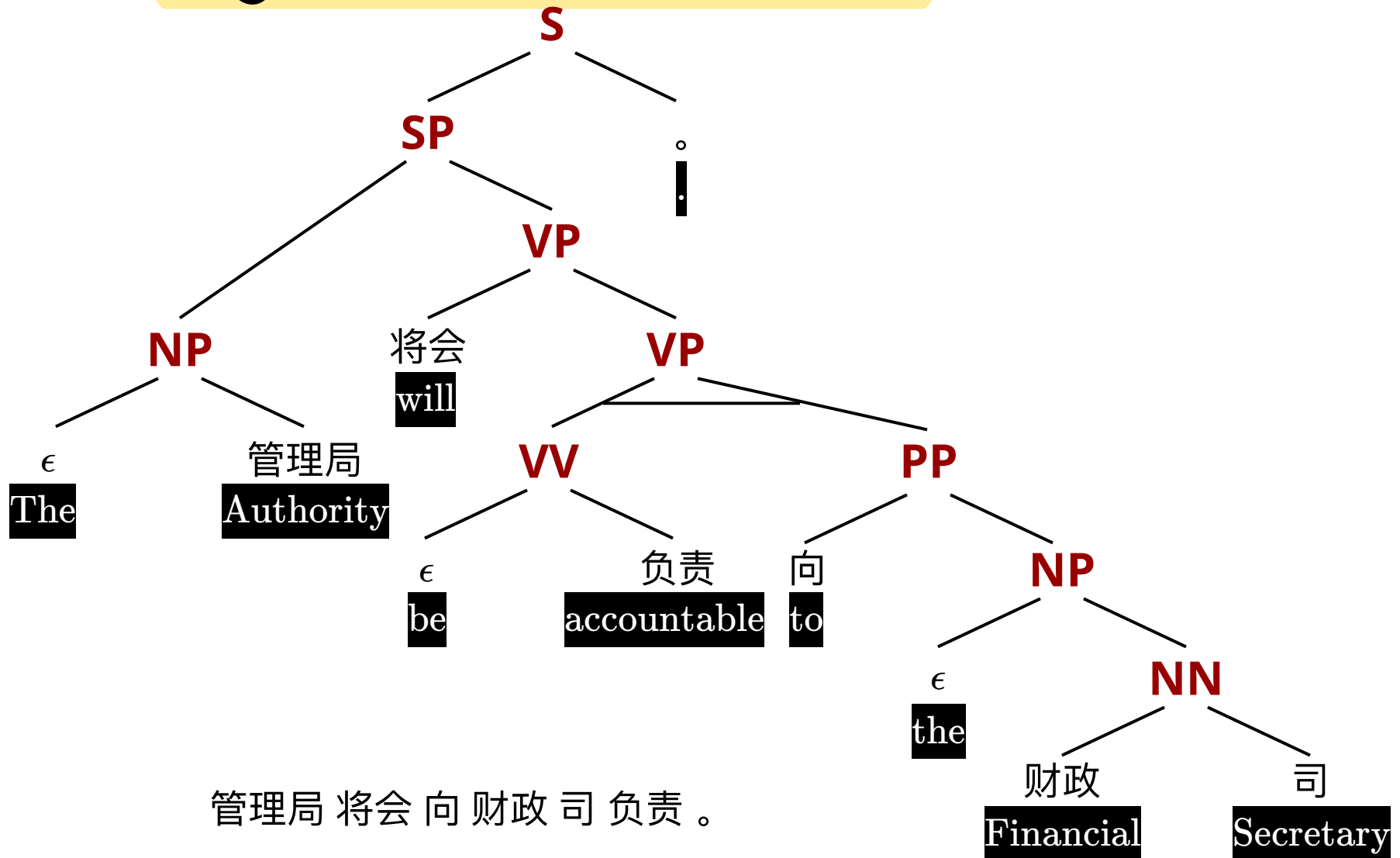
We introduce a general formalism for modeling of bilingual sentence pairs, known as an **inversion transduction grammar**, with potential application in a variety of corpus analysis areas. Transduction grammar models, especially of the finite-state family, have long been known. However, the imposition of identical ordering constraints upon both streams severely restricts their applicability, and thus transduction grammars have received relatively little attention in language-modeling research. The inversion transduction grammar formalism skips directly to a context-free, rather than finite-state, base and permits one extra degree of ordering flexibility, while retaining properties necessary for efficient computation, thereby sidestepping the limitations of traditional transduction grammars.

In tandem with the concept of bilingual language-modeling, we propose the concept of bilingual parsing, where the input is a sentence-pair rather than a sentence. Though inversion transduction grammars remain inadequate as full-fledged translation models, bilingual parsing with simple inversion transduction grammars turns out to be very useful for parallel corpus analysis when the true grammar is not fully known. Parallel bilingual corpora have been shown to provide a rich source of constraints for statistical analysis (Brown et al. 1990; Gale and Church 1991; Gale, Church, and Yarowsky 1992; Church 1993; Brown et al. 1993; Dagan, Church, and Gale 1993;

* Department of Computer Science, University of Science and Technology, Clear Water Bay, Hong Kong
E-mail: dekai@cs.ust.hk

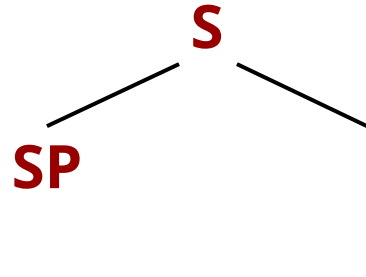


Synchronous Tree



The Authority will be accountable to the Financial Secretary .

Synchronous Derivation



$S \rightarrow (SP \text{ Stop}, SP \text{ Stop})$

$\text{Stop} \rightarrow (., .)$

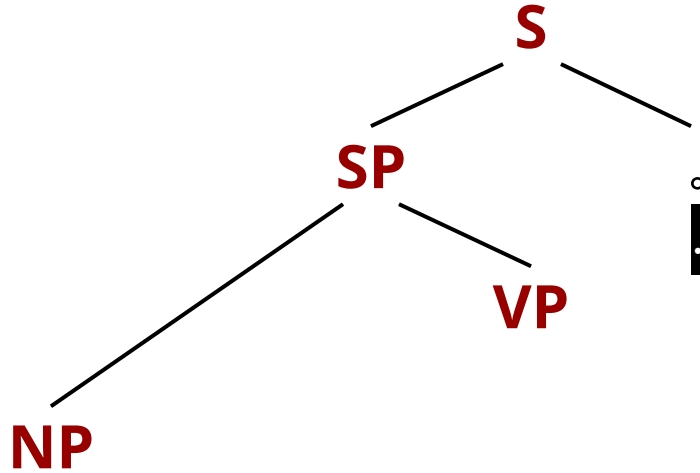
SP
|

o

SP

.

Synchronous Derivation

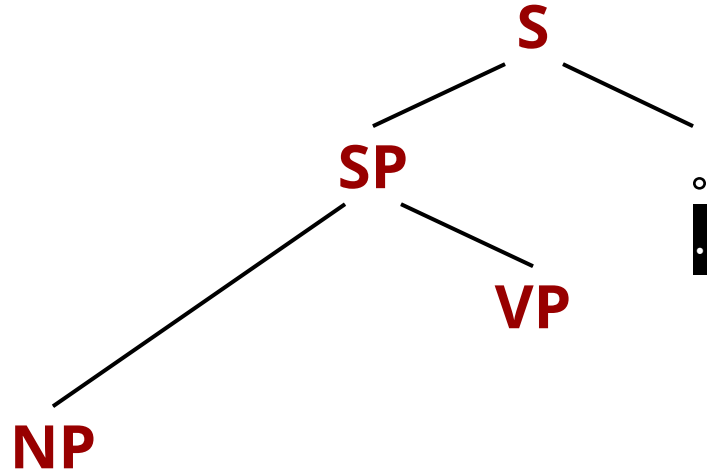


$SP \rightarrow (NP VP, NP VP)$

We will simplify
this rule!



Synchronous Derivation

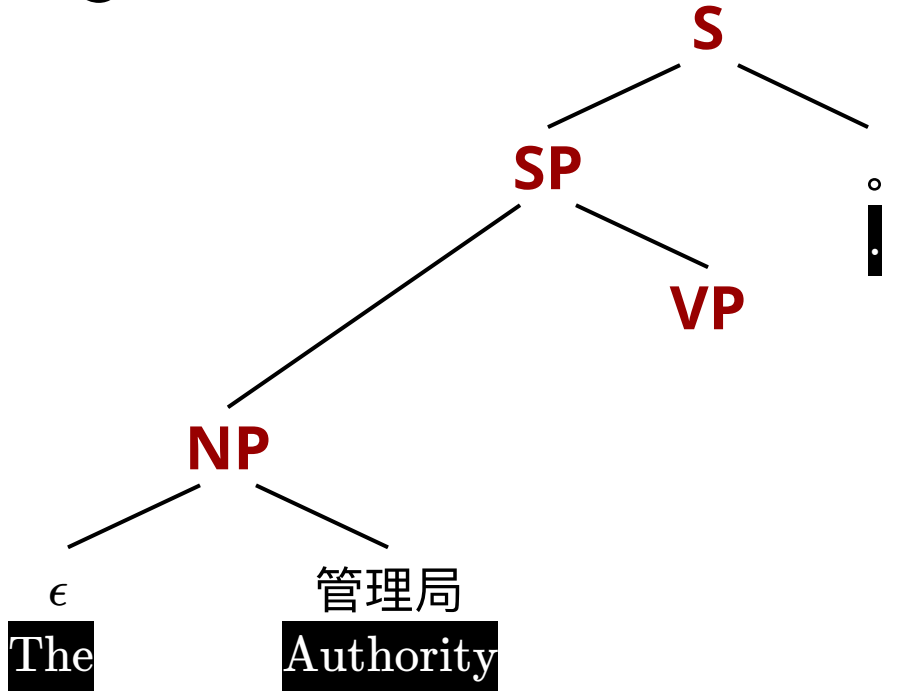


SP \rightarrow (NP VP)

This means both source and target follow the same pattern.



Synchronous Derivation



$\text{NP} \rightarrow (\text{D N})$

$\text{D} \rightarrow (\epsilon, \text{The})$

$\text{N} \rightarrow (\text{管理局}, \text{Authority})$

管理局

VP
/

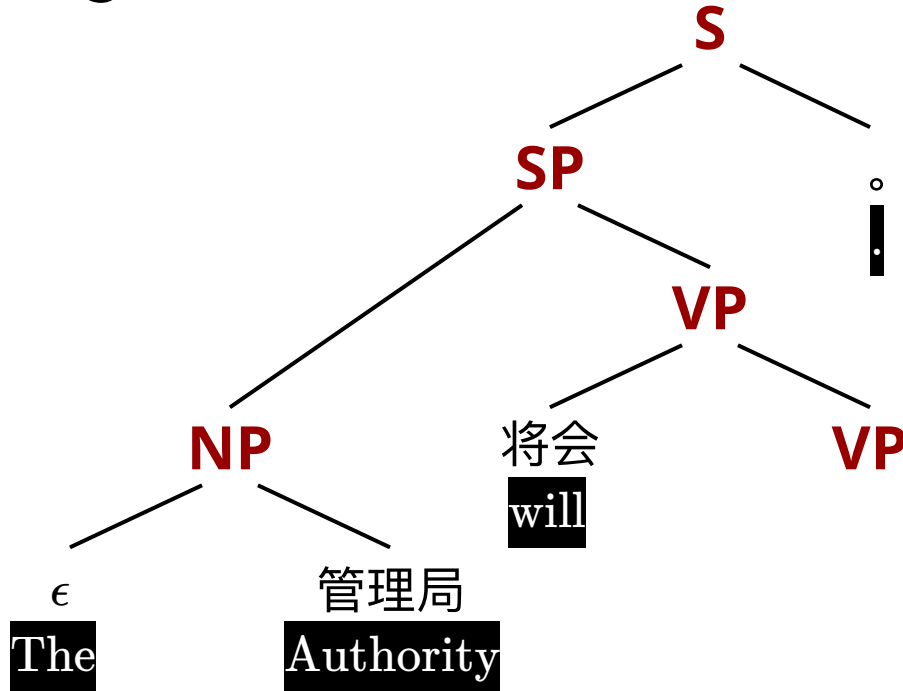
.

The Authority

VP

.

Synchronous Derivation



$VP \rightarrow (V VP)$

$V \rightarrow (\text{将会}, \text{will})$

管理局 将会

VP

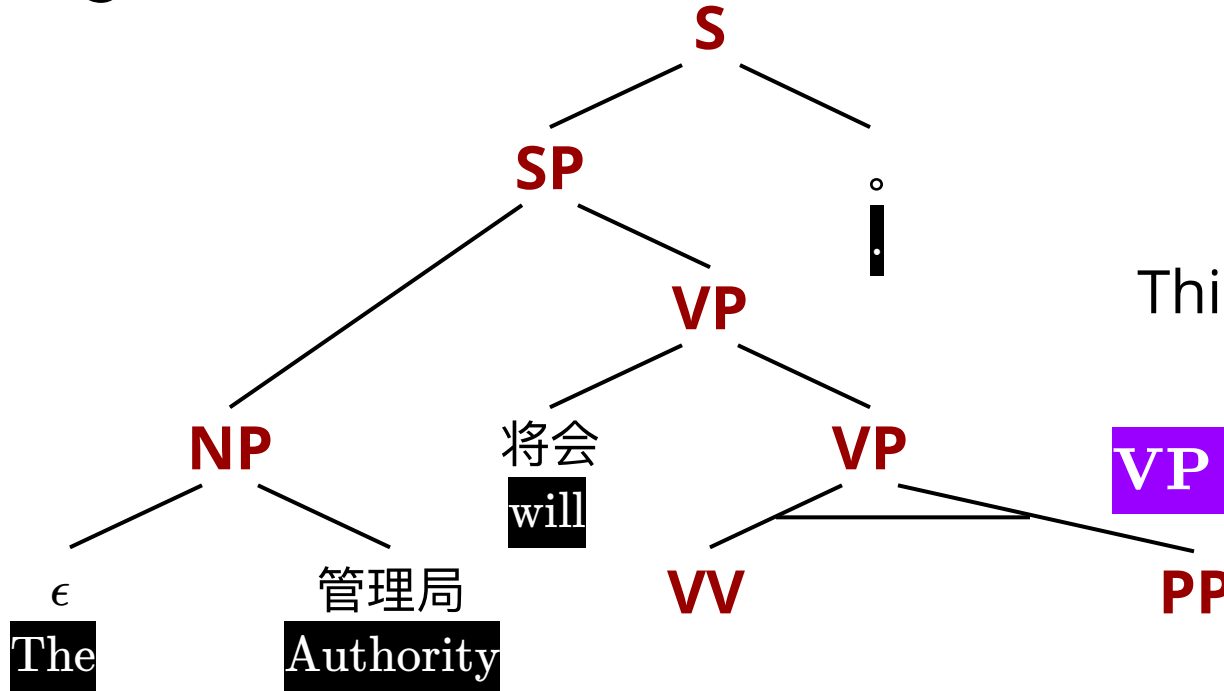
。

The Authority will

VP

。

Synchronous Derivation



$VP \rightarrow \langle VV PP \rangle$

This is equivalent to the following:

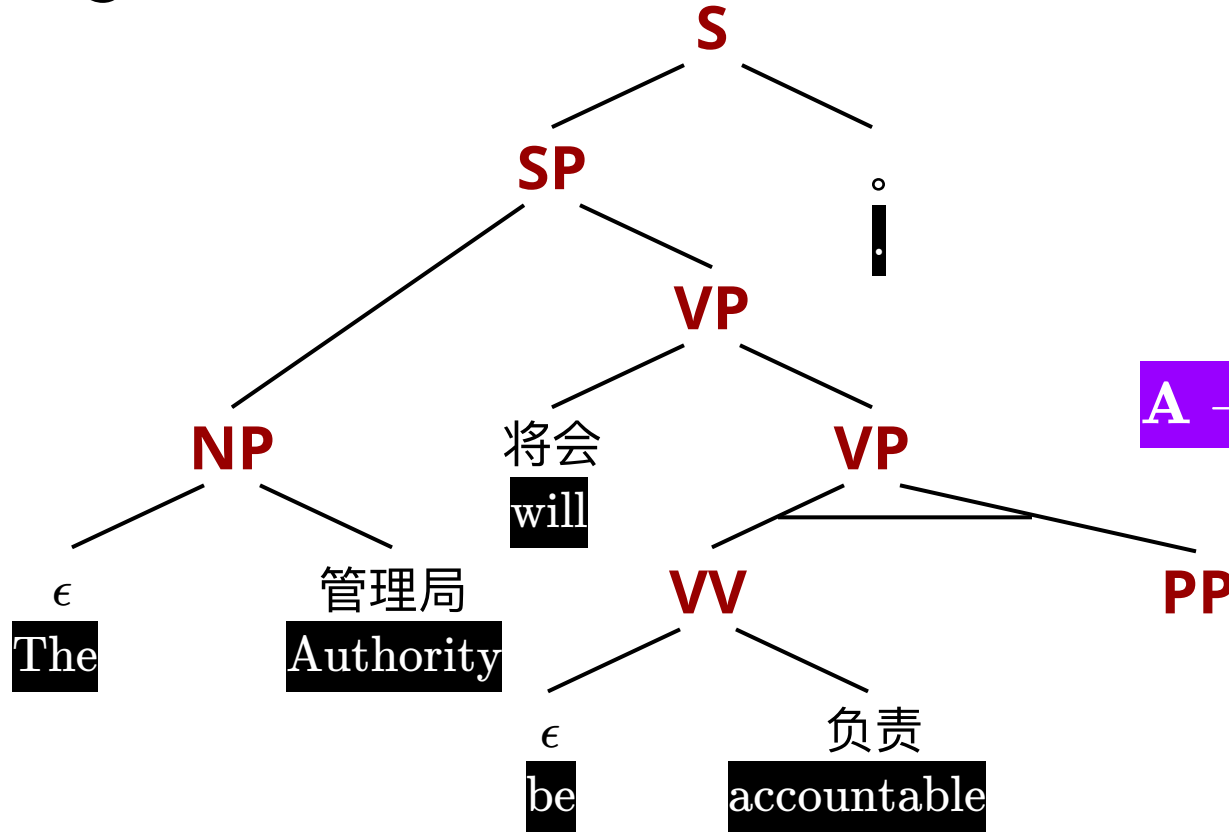
$VP \rightarrow (PP VV, VV PP)$

Note that this definition may be slightly different from what's described in the original paper, but is also valid and is more convenient for our illustration.

管理局 将会 PP VV 。

The Authority will VV PP .

Synchronous Derivation



$VV \rightarrow (V A)$

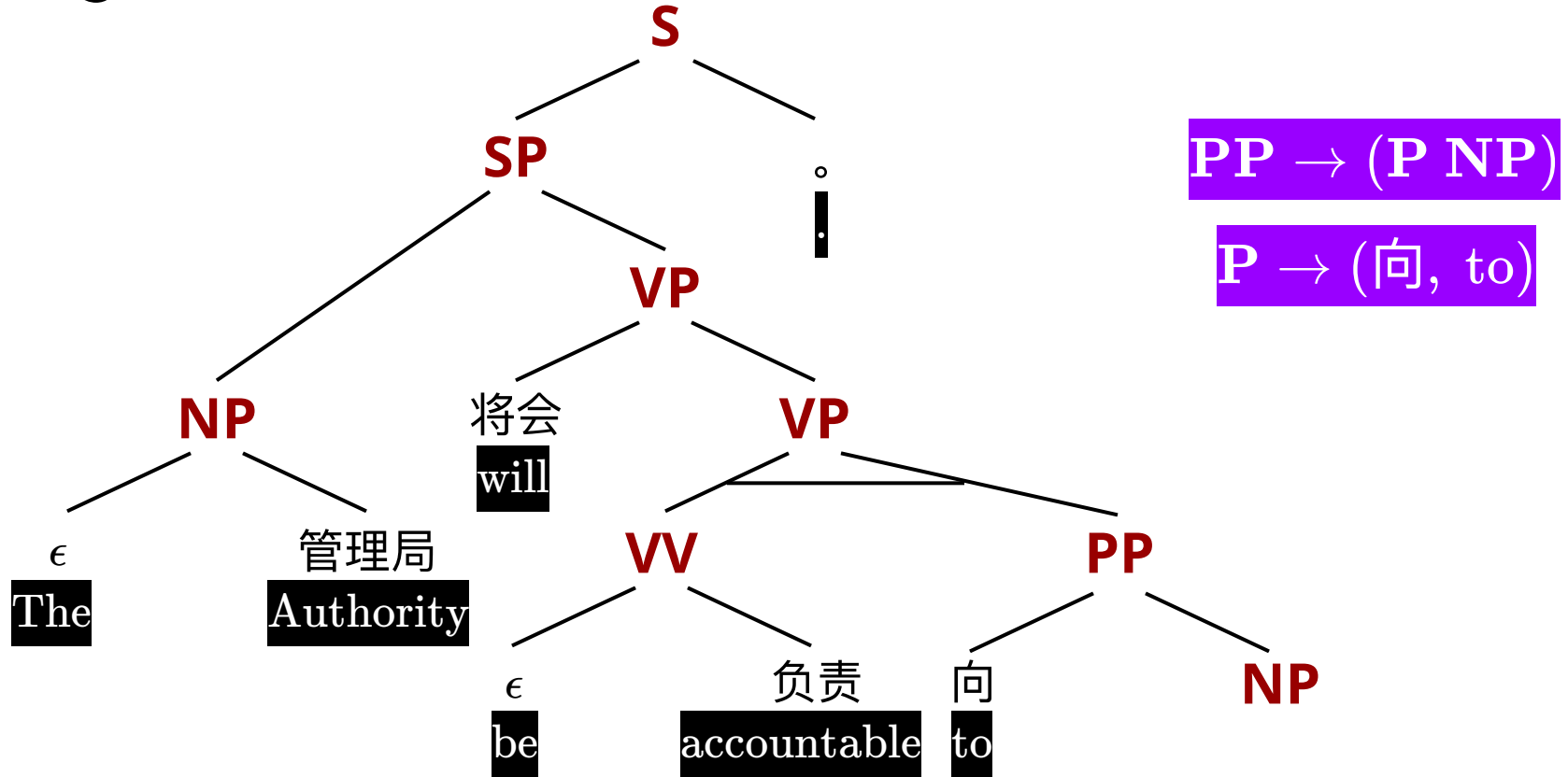
$V \rightarrow (\epsilon, be)$

$A \rightarrow (\text{负责}, accountable)$

管理局 将会 PP 负责。

The Authority will be accountable PP .

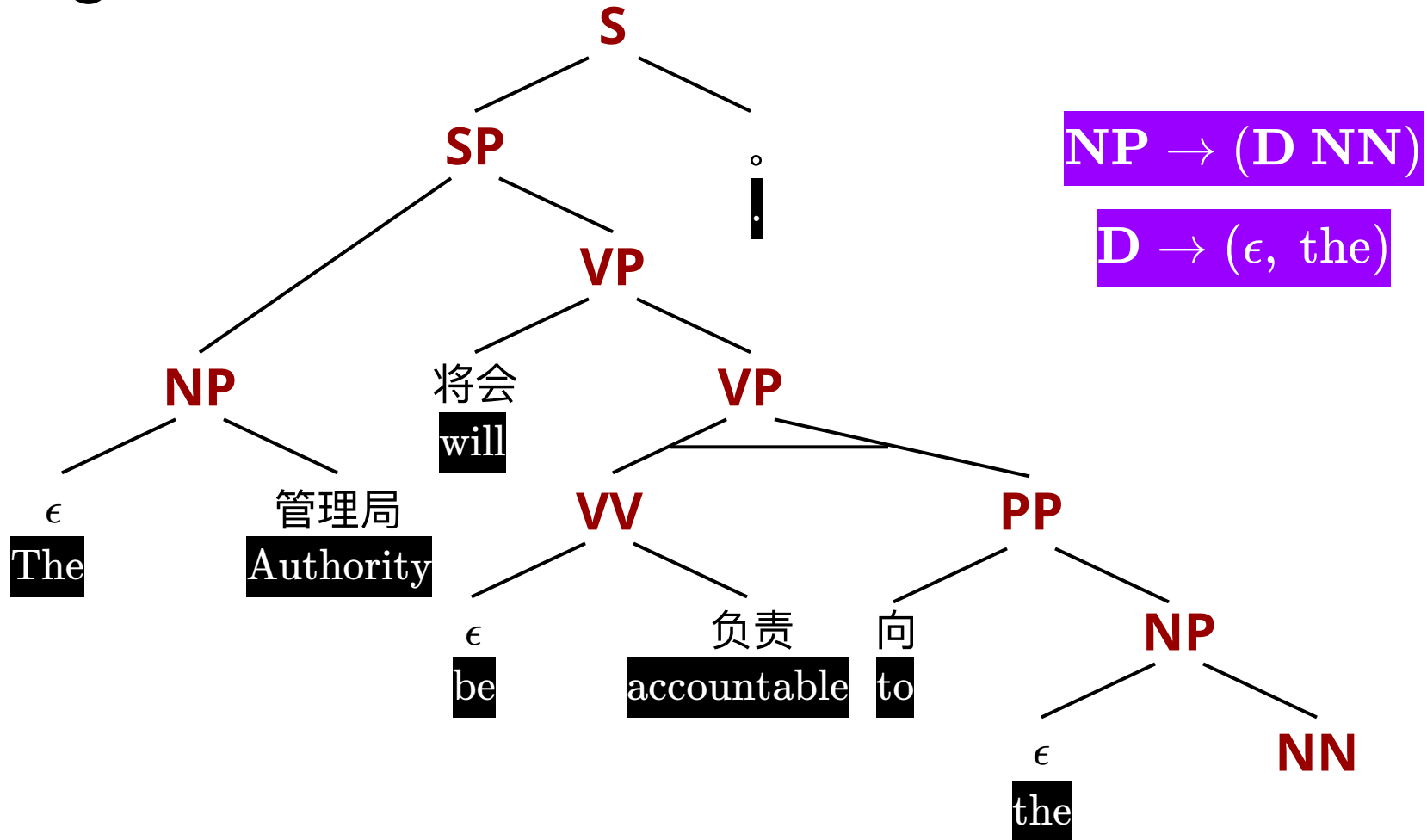
Synchronous Derivation



管理局 将会 向 NP 负责。

The Authority will be accountable to NP .

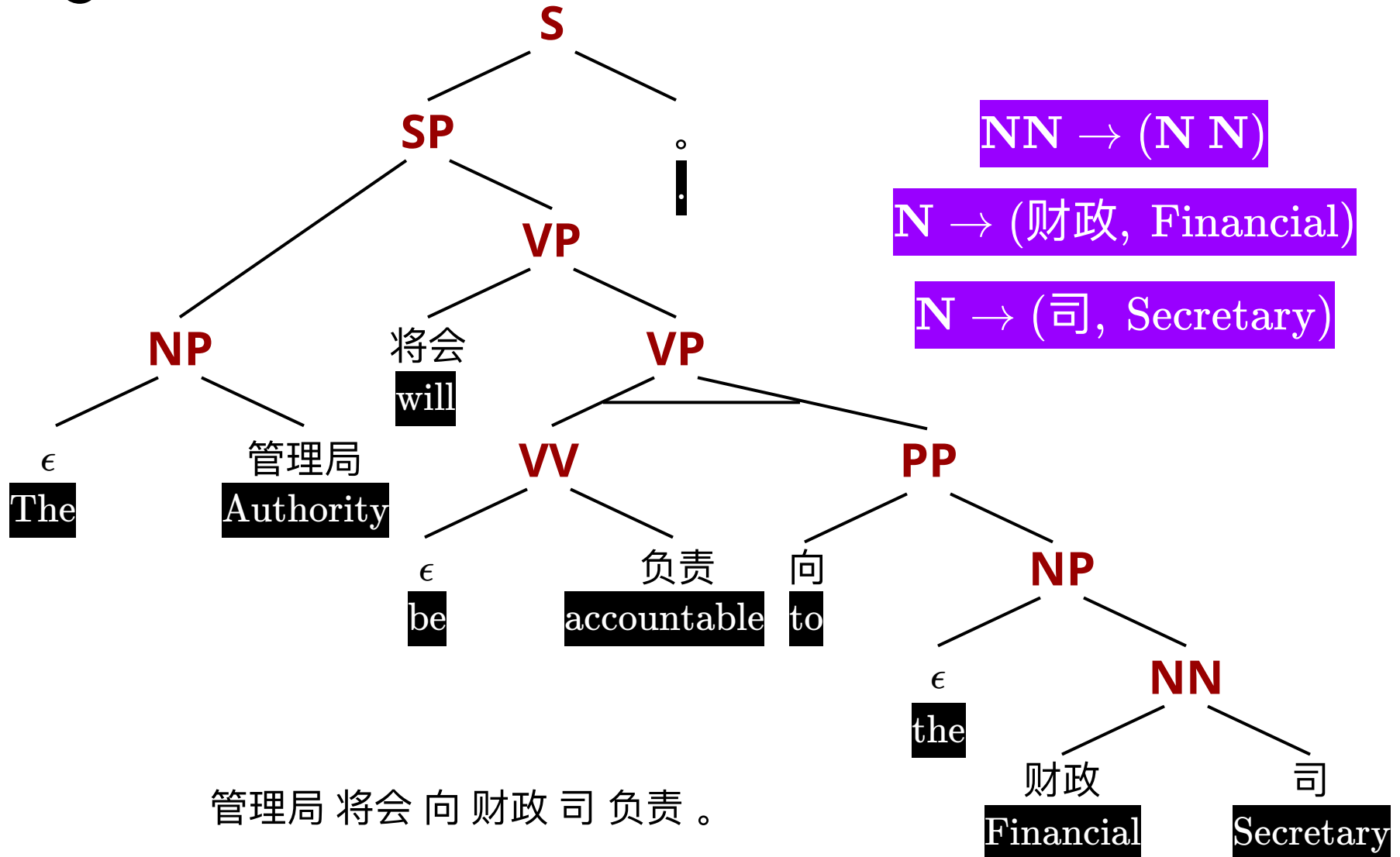
Synchronous Derivation



管理局 将会 向 NN 负责。

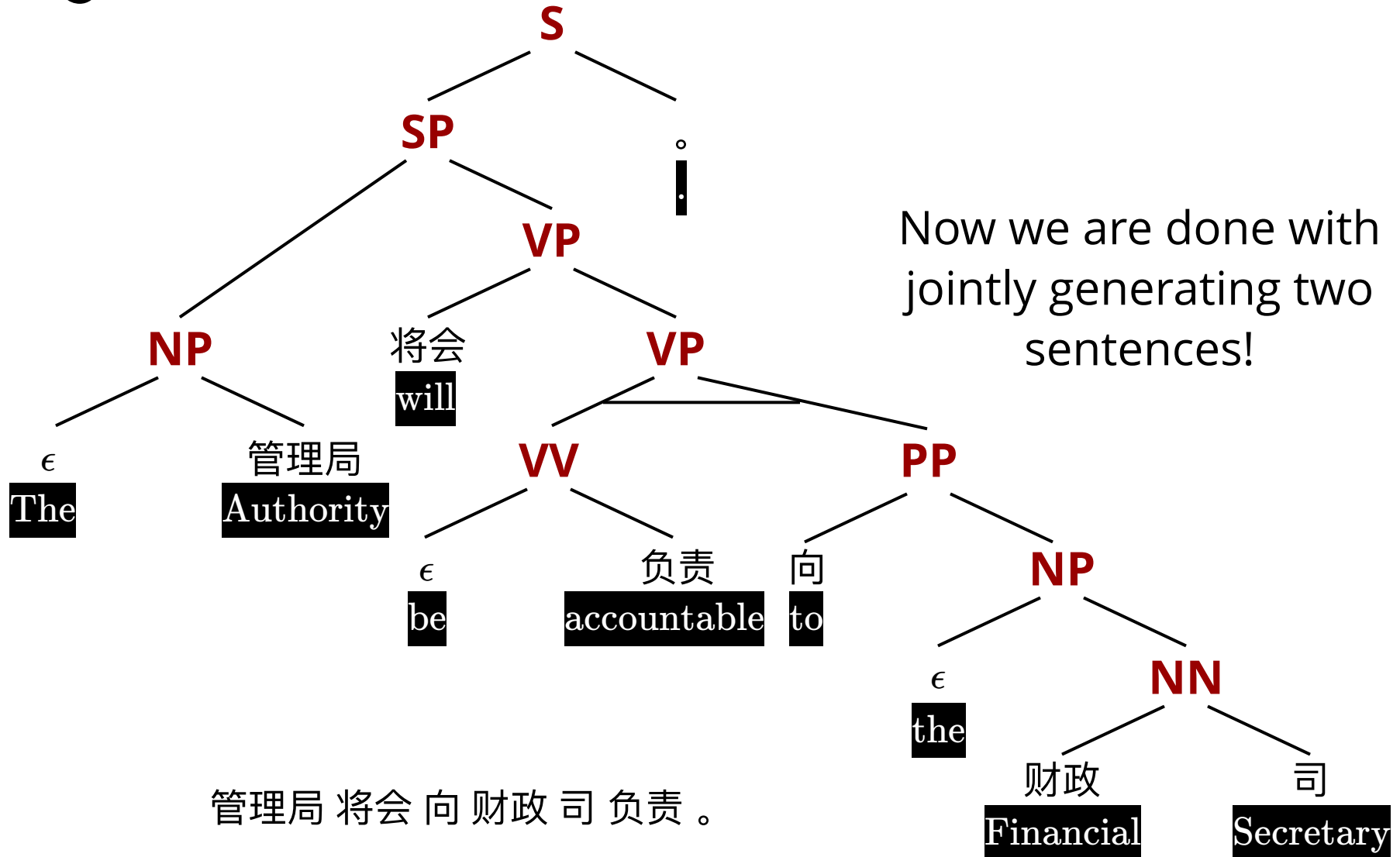
The Authority will be accountable to the NN .

Synchronous Derivation

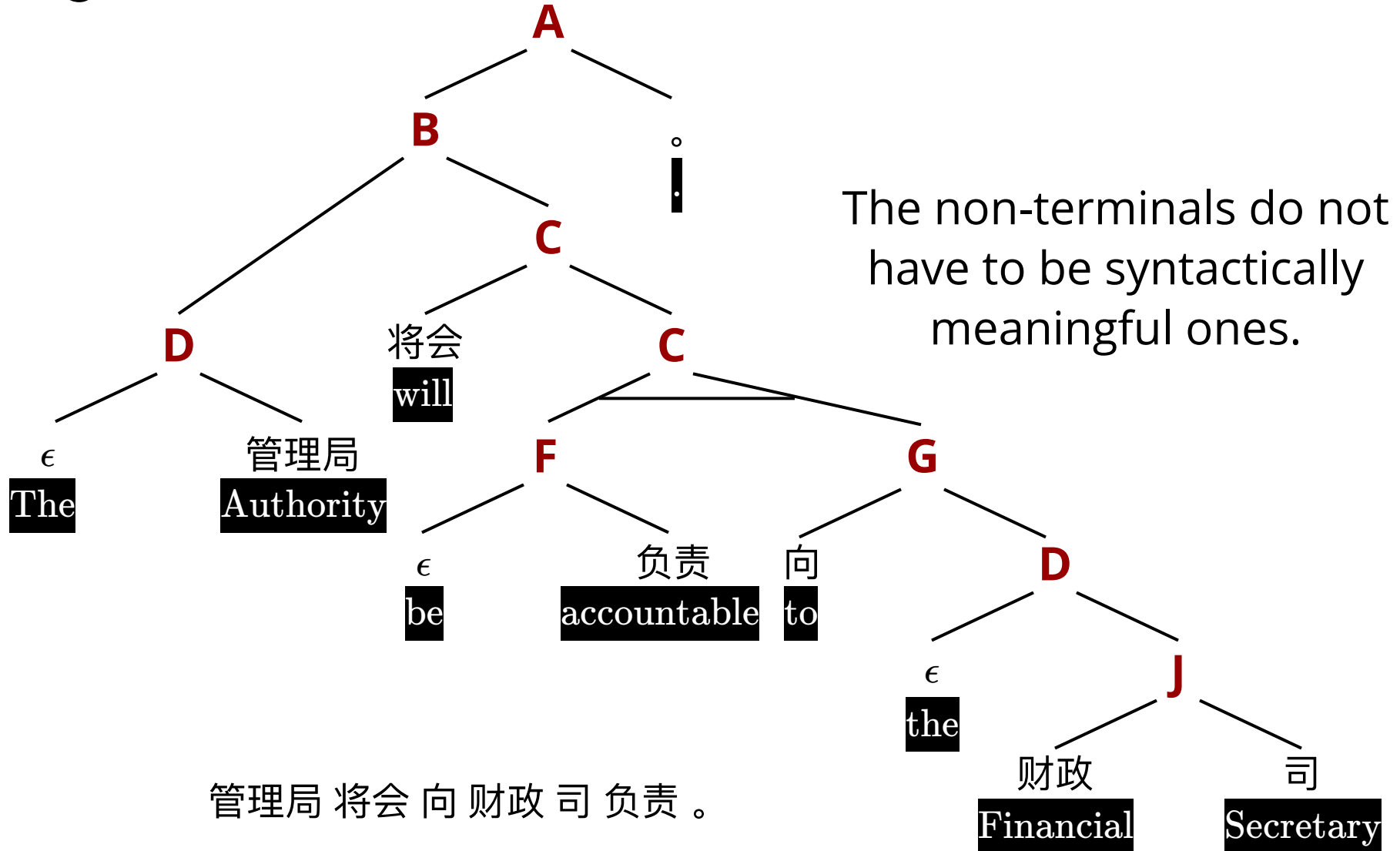


The Authority will be accountable to the Financial Secretary .

Synchronous Derivation

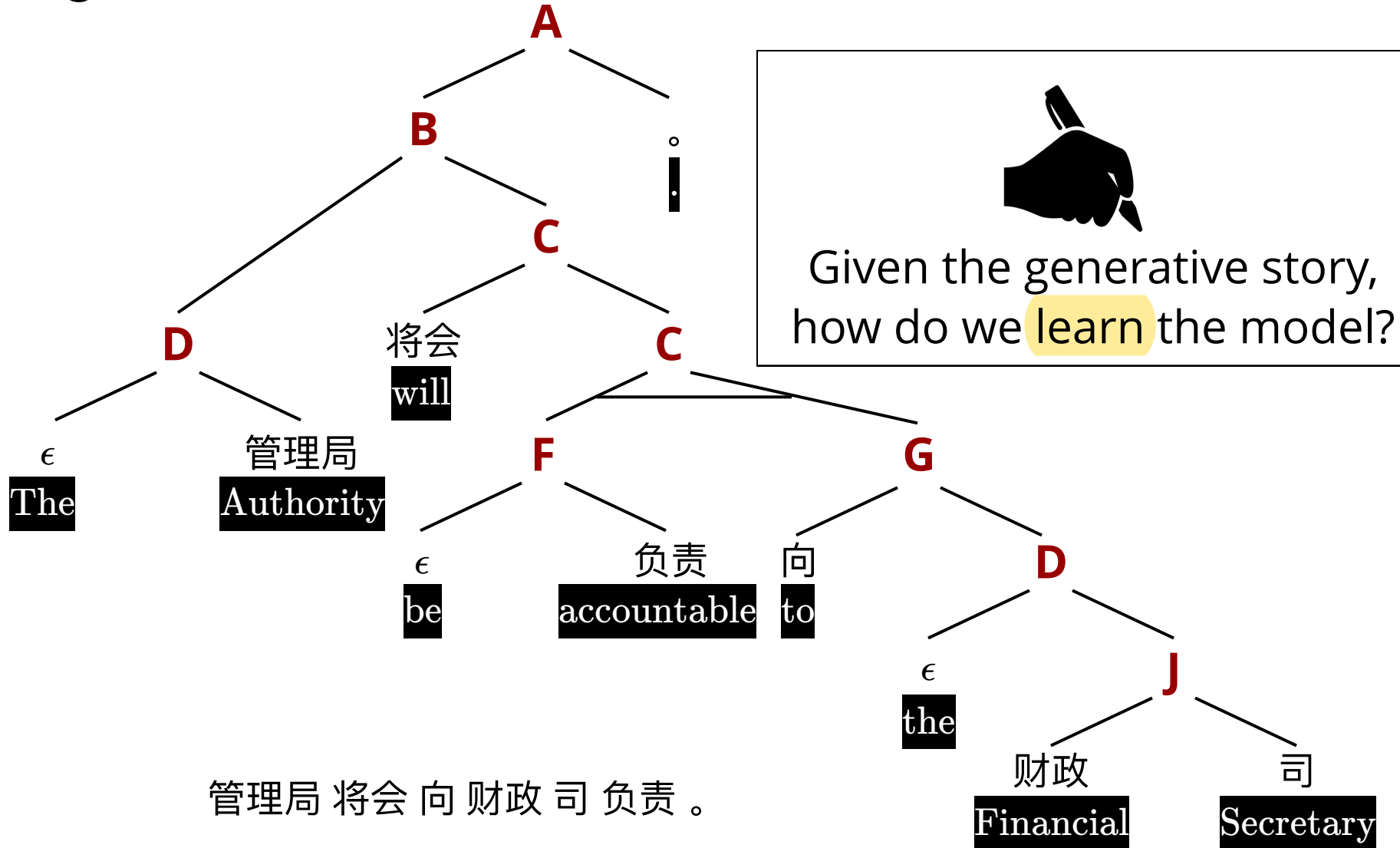


Synchronous Derivation

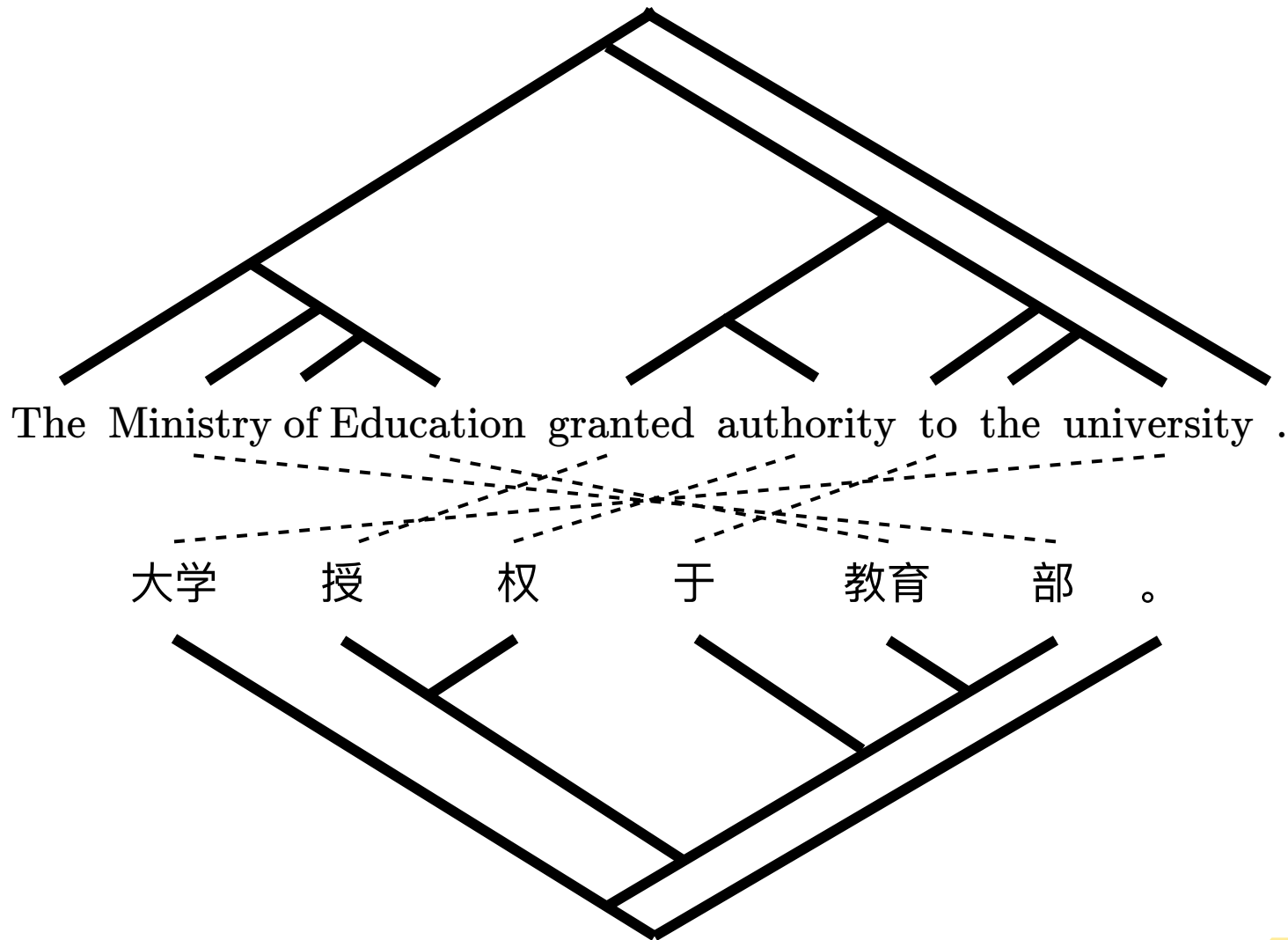


The Authority will be accountable to the Financial Secretary .

Synchronous Derivation

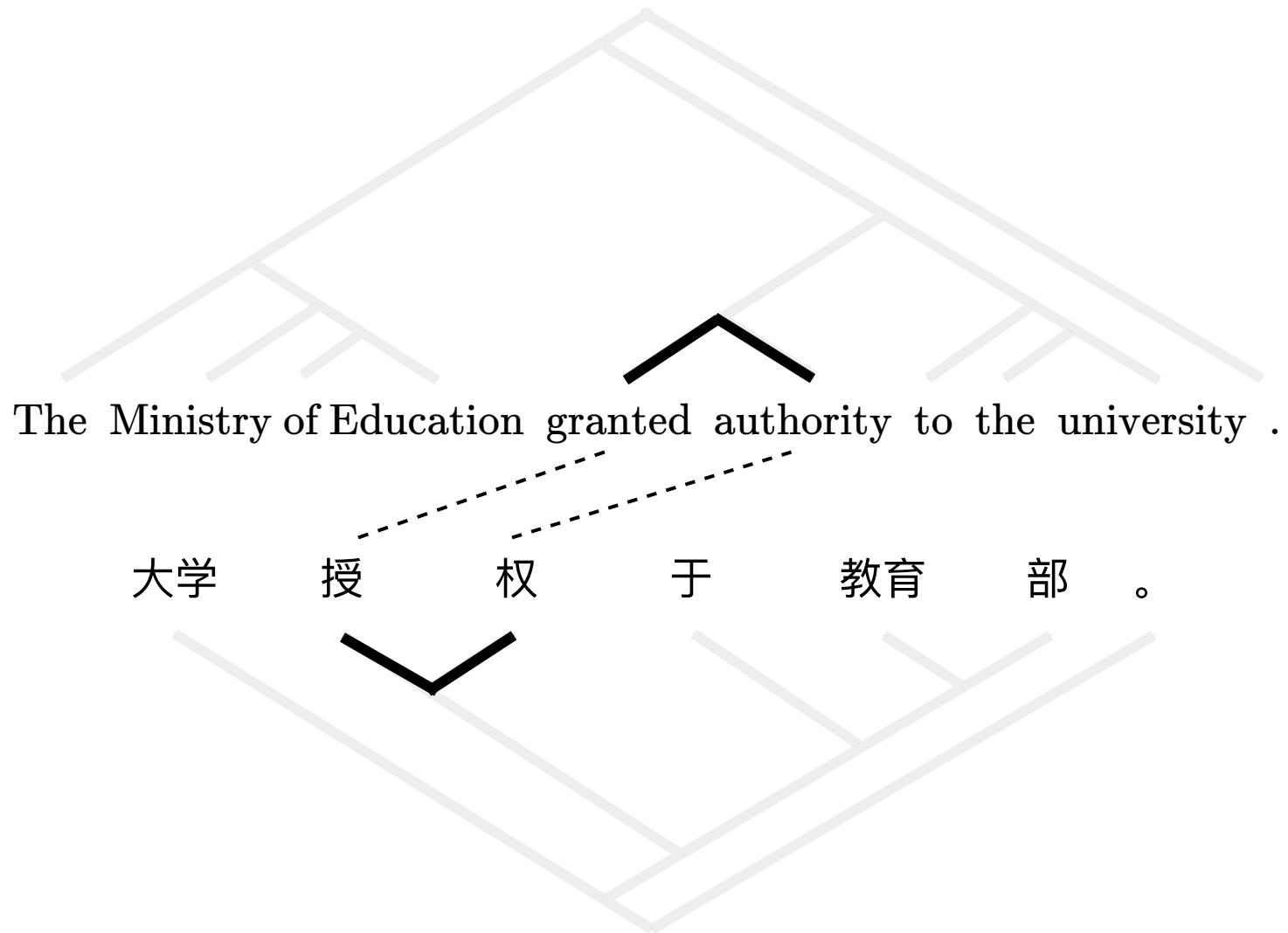


Limitations with ITC

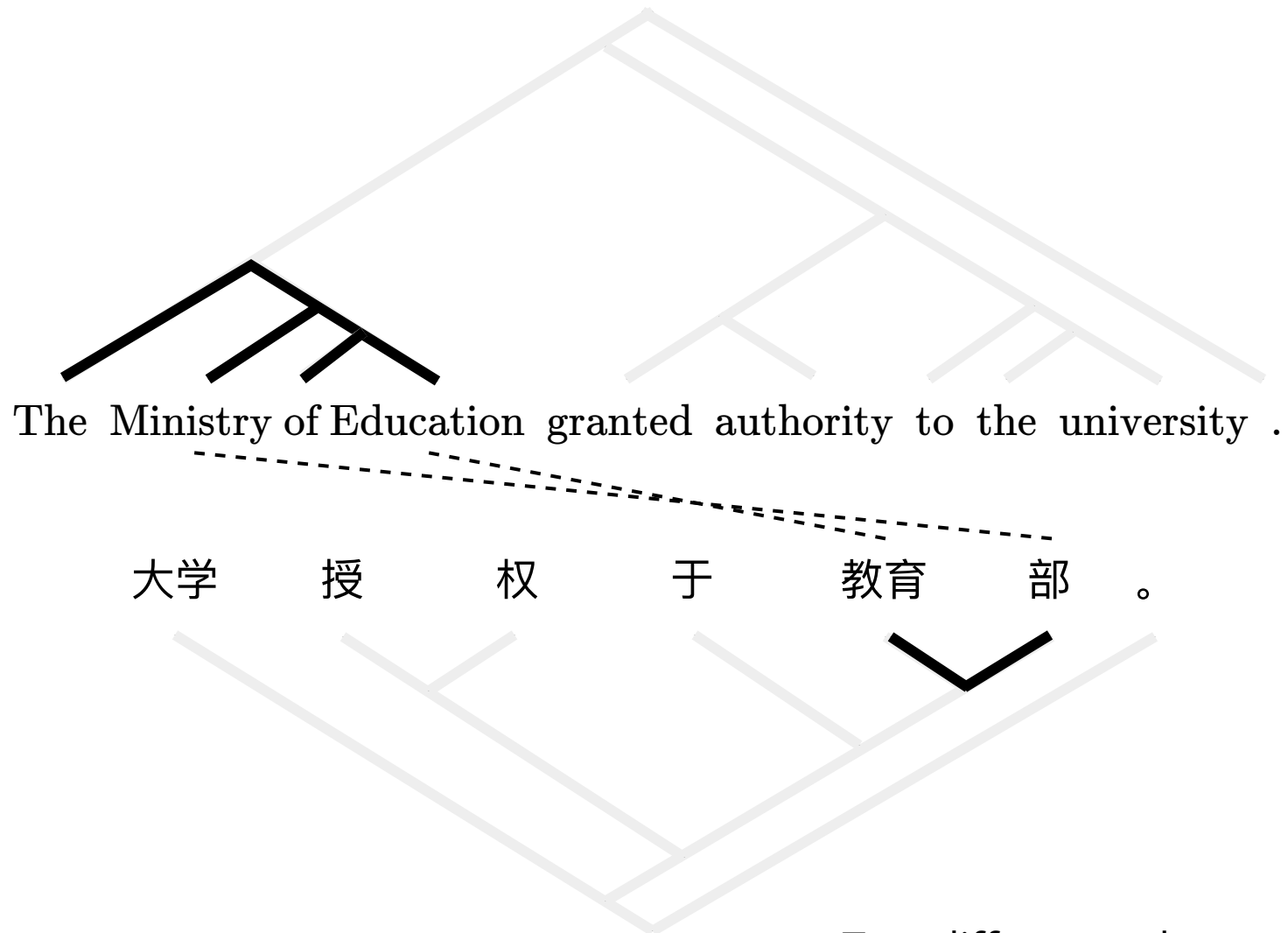


Not possible to find the correct alignment for certain sentence pairs.³⁶

Limitations with ITC

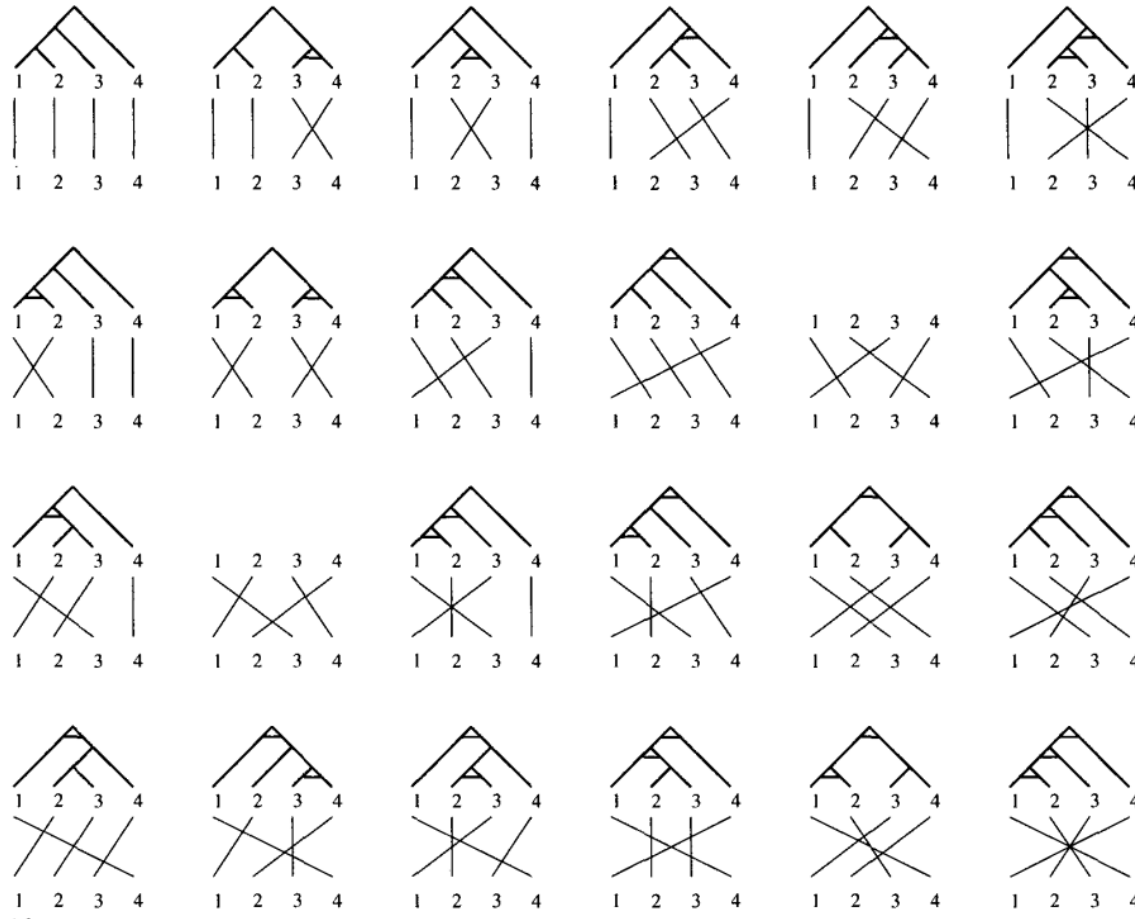


Limitations with ITC



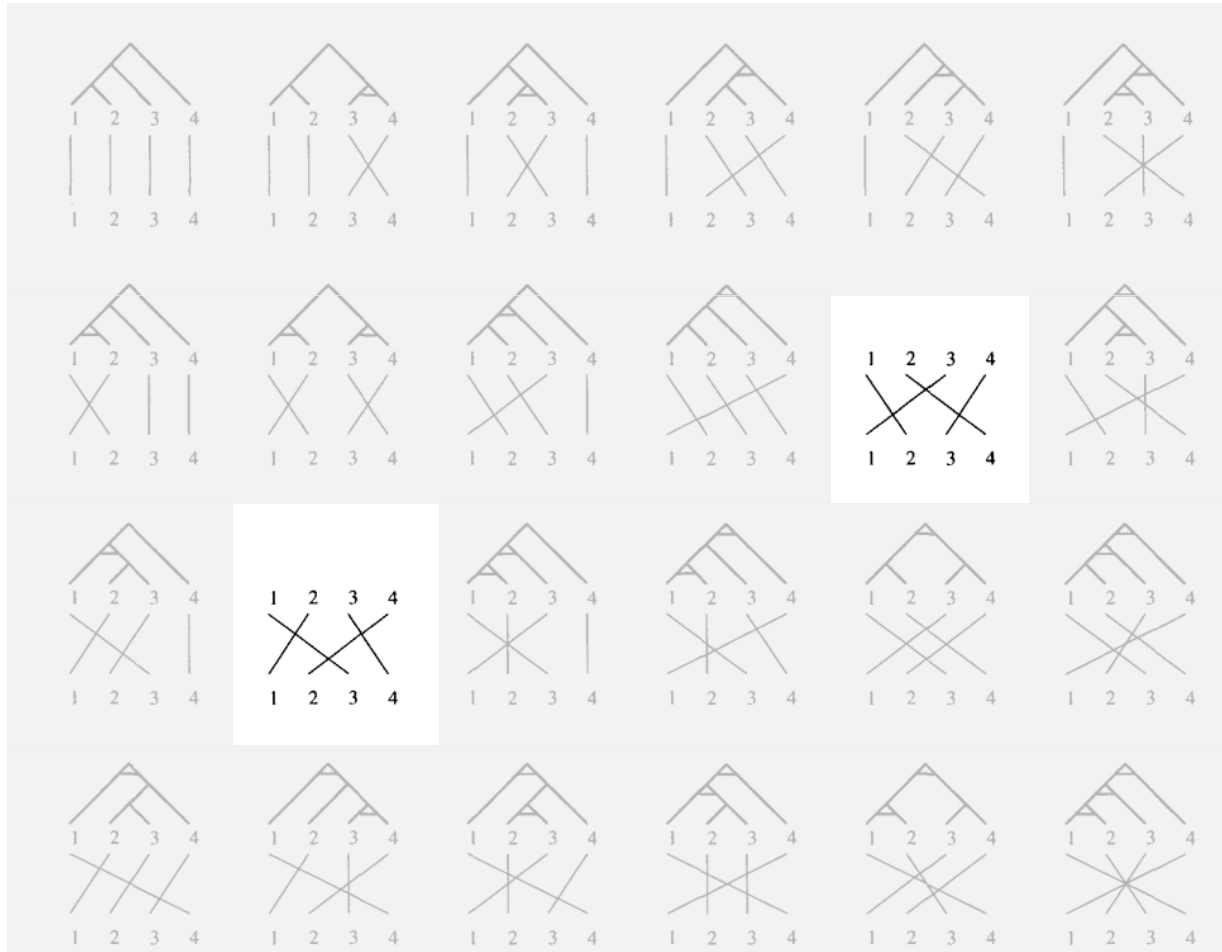
Two different sub-trees are in
different larger sub-trees.

Limitations with ITG



Possible word alignments and the
corresponding ITG parses (if possible)

Limitations with ITG



2 out of 24 cases where the ITG is **unable** to handle.

Limitations with ITG



We will come back to this issue shortly...

2 out of 24 cases where the ITG is unable to handle.

Question

What about parsing the source sentence first?

Tree Transducer

A Syntax-based Statistical Translation Model

Kenji Yamada and Kevin Knight
Information Sciences Institute
University of Southern California
4676 Admiralty Way, Suite 1001
Marina del Rey, CA 90292
{kyamada,knight}@isi.edu

Abstract

We present a syntax-based statistical translation model. Our model trans-
forms a source-language parse tree
into a target-language string by apply-
ing stochastic operations at each node.
These operations capture linguistic dif-
ferences such as word order and case
marking. Model parameters are esti-
mated in polynomial time using an EM
algorithm. The model produces word
alignments that are better than those
produced by IBM Model 5.

1 Introduction

A statistical translation model (TM) is a mathe-
matical model in which the process of human-
language translation is statistically modeled.
Model parameters are automatically estimated us-
ing a corpus of translation pairs. TMs have been
used for statistical machine translation (Berger et
al., 1996), word alignment of a translation cor-
pus (Melamed, 2000), multilingual document cor-
retrieval (Franz et al., 1999), automatic dictionary
construction (Resnik and Melamed, 1997), and
data preparation for word sense disambiguation
programs (Brown et al., 1991). Developing a bet-
ter TM is a fundamental issue for those applica-
tions.

Researchers at IBM first described such a sta-
tistical TM in (Brown et al., 1988). Their mod-
els are based on a string-to-string noisy channel
model. The channel converts a sequence of words
in one language (such as English) into another
(such as French). The channel operations are
movements, duplications, and translations, ap-
plied to each word independently. The movement

is conditioned only on word classes and positions
in the string, and the duplication and transla-
tions are conditioned only on the word identities
(Brown et al., 1993).

One criticism of the IBM-style TM is that it
does not model structural or syntactic aspects of
the language. The TM was only demonstrated on
a structurally similar language pair (English and
French). It has been suspected that a language
pair with very different word order such as En-
glish and Japanese would not be modeled well by
these TMs.

To incorporate structural aspects of the lan-
guage, our channel model accepts a parse tree as
an input, i.e., the input sentence is preprocessed
by a syntactic parser. The channel performs oper-
ations on each node of the parse tree. The oper-
ations are *reordering* child nodes, *inserting* extra
words at each node, and *translating* leaf words.
Figure 1 shows the overview of the operations of
our model. Note that the output of our model is a
string, not a parse tree. Therefore, parsing is only
needed on the channel input side.

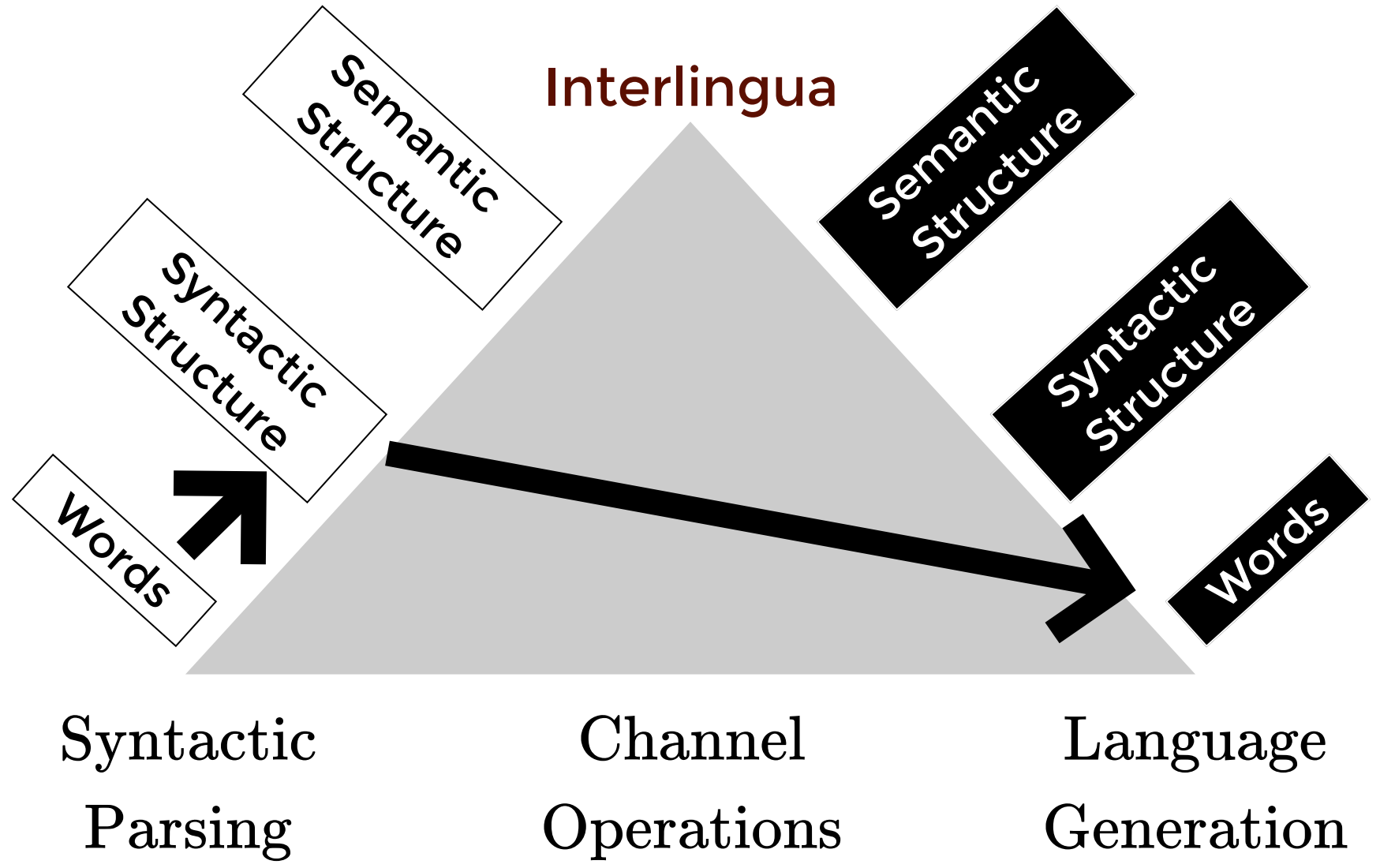
The reorder operation is intended to model
translation between languages with different word
orders, such as SVO-languages (English or Chi-
nese) and SOV-languages (Japanese or Turkish).
The word-insertion operation is intended to cap-
ture linguistic differences in specifying syntactic
cases. E.g., English and French use structural po-
sition to specify case, while Japanese and Korean
use case-marker particles.

Wang (1998) enhanced the IBM models by in-
troducing phrases, and Och et al. (1999) used
templates to capture phrasal sequences in a sen-
tence. Both also tried to incorporate structural as-
pects of the language, however, neither handles

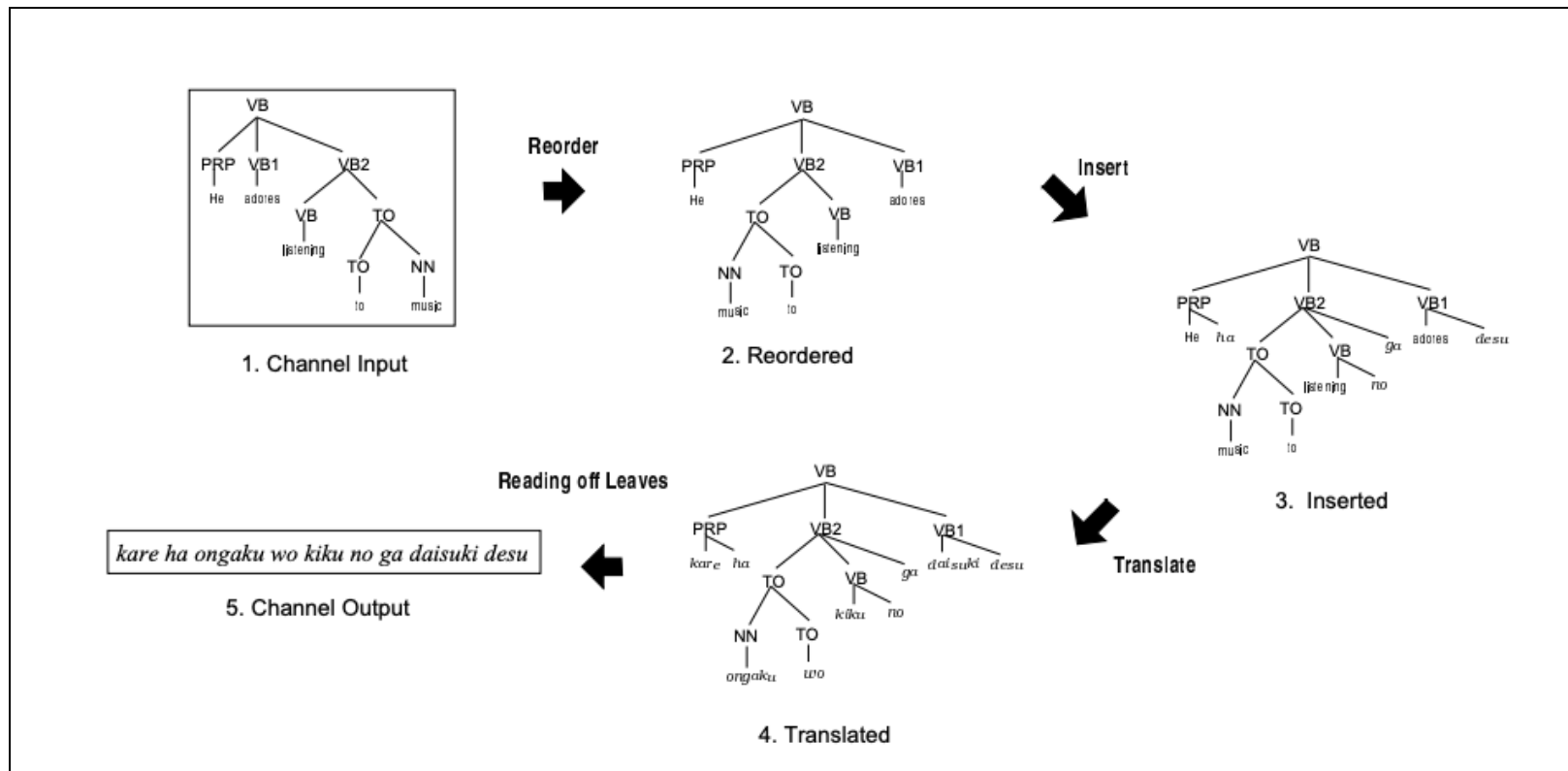



Optional

Tree Transducer



Tree Transducer



Define several operations/actions that can be applied on top of source syntax trees. It incrementally transform the source syntax tree into target sentence. 

Question

Can we introduce **phrases** into the translation process?

Formal Syntax

Hierarchical Phrase-Based Translation

David Chiang*
Information Sciences Institute
University of Southern California

We present a statistical machine translation model that uses hierarchical phrases—phrases that contain subphrases. The model is formally a synchronous context-free grammar but learned from a parallel text without any syntactic annotations. Thus it can be seen as combining fundamental ideas from both syntax-based translation and phrase-based translation. We describe our system's training and decoding methods in detail, and evaluate it for translation speed and translation accuracy. Using BLEU as a metric of translation accuracy, we find that our system performs significantly better than the Alignment Template System, a state-of-the-art phrase-based system.

1. Introduction

The alignment template translation model (Och and Ney 2004) and related phrase-based models advanced the state of the art in machine translation by expanding the basic unit of translation from words to **phrases**, that is, substrings of potentially unlimited size (but not necessarily phrases in any syntactic theory). These phrases allow a model to learn local reorderings, translations of multiword expressions, or insertions and deletions that are sensitive to local context. This makes them a simple and powerful mechanism for translation.

The basic phrase-based model is an instance of the noisy-channel approach (Brown et al. 1993). Following convention, we call the source language “French” and the target language “English”; the translation of a French sentence f into an English sentence e is modeled as:

$$\arg \max_e P(e | f) = \arg \max_e P(e, f)$$
$$= \arg \max_e (P(e) \times P(f | e)) \quad (1)$$

The phrase-based translation model $P(f | e)$ “encodes” e into f by the following steps:

1. segment e into phrases $\tilde{e}_1 \dots \tilde{e}_l$, typically with a uniform distribution over segmentations;

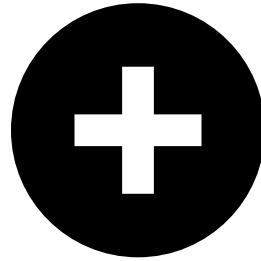
* 4676 Admiralty Way, Suite 1001, Marina del Rey, CA 90292, USA. E-mail: chiang@isi.edu. Much of the research presented here was carried out while the author was at the University of Maryland Institute for Advanced Computer Studies.

Submission received: 1 May 2006; accepted for publication: 3 October 2006.



Formal Syntax

Phrase-based
Translation



Synchronous
Grammar

It combines the idea of
phrase-based translation
and synchronous parsing.

Hierarchical Phrase Pairs

澳洲是与北韩有邦交的少数国家之一

Australia is one of the few countries that have diplomatic relations with North Korea

(北韩, North Korea)

(邦交, diplomatic relations)

(与 1 有 2, have 2 with 1)



Synchronous Grammar

澳洲是与北韩有邦交的少数国家之一

Australia is one of the few countries that have diplomatic relations with North Korea

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

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

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

Subscripts indicate
the alignment

Use some heuristics similar to
what we discussed in **phrase-**
based translation to acquire the
synchronous grammar rules

Grammar Rule

澳洲是与北韩有邦交的少数国家之一

Australia is one of the few countries that have diplomatic relations with North Korea

A non-terminal
symbol


$$\mathbf{X} \rightarrow (\gamma, \alpha)$$


A sequence of *source* terminal
symbols intermixed with non-
terminal symbols

A sequence of *target* terminal
symbols intermixed with
non-terminal symbols

Grammar Rule

澳洲是与北韩有邦交的少数国家之一

Australia is one of the few countries that have diplomatic relations with North Korea

$$S \rightarrow (X_1, X_1) \quad \textcircled{A}$$

$$S \rightarrow (S_1 X_2, S_1 X_2) \quad \textcircled{B}$$

Two special rules considered (auxiliary rules for completion of derivations)

Synchronous Derivation

澳洲是与北韩有邦交的少数国家之一

Australia is one of the few countries that have diplomatic relations with North Korea

$X \rightarrow (\text{与 } X_1 \text{ 有 } X_2, \text{ have } X_2 \text{ with } X_1)$ ①

$X \rightarrow (X_1 \text{ 的 } X_2, \text{ the } X_2 \text{ that } X_1)$ ②

$X \rightarrow (X_1 \text{ 之一, one of } X_1)$ ③

$X \rightarrow (\text{澳洲, Australia})$ ④

$X \rightarrow (\text{北韩, North Korea})$ ⑤

$X \rightarrow (\text{是, is})$ ⑥

$X \rightarrow (\text{邦交, diplomatic relations})$ ⑦

$X \rightarrow (\text{少数 国家, few countries})$ ⑧



Can you work out the derivation for this sentence pair based on the above rules (and the two special rules)?

Synchronous Derivation



澳洲是与北韩有邦交的少数国家之一

Australia is one of the few countries that have diplomatic relations with North Korea

(S_1, S_1)

$\xrightarrow[S_1]{(B)} (S_2 X_3, S_2 X_3)$

Rewrite S_1 with rule (B)

$\xrightarrow[S_2]{(B)} (S_4 X_5 X_3, S_4 X_5 X_3)$

$\xrightarrow[S_4]{(A)} (X_6 X_5 X_3, X_6 X_5 X_3)$

$\xrightarrow[X_6]{(4)} (\text{澳洲 } X_5 X_3, \text{Australia } X_5 X_3)$

$\xrightarrow[X_5]{(6)} (\text{澳洲 是 } X_3, \text{Australia is } X_3)$

$\xrightarrow[X_3]{(3)} (\text{澳洲 是 } X_7 \text{ 之一}, \text{Australia is one of } X_7)$

Synchronous Derivation

$\xrightarrow[\mathbf{X}_7]{(2)}$ (澳洲是 \mathbf{X}_8 的 \mathbf{X}_9 之一, Australia is one of the \mathbf{X}_9 that \mathbf{X}_8)

$\xrightarrow[\mathbf{X}_8]{(1)}$ (澳洲是与 \mathbf{X}_{10} 有 \mathbf{X}_{11} 的 \mathbf{X}_9 之一,

Australia is one of the \mathbf{X}_9 that have \mathbf{X}_{11} with \mathbf{X}_{10})

$\xrightarrow[\mathbf{X}_{10}]{(5)}$ (澳洲是与 北韩 有 \mathbf{X}_{11} 的 \mathbf{X}_9 之一,

Australia is one of the \mathbf{X}_9 that have \mathbf{X}_{11} with North Korea)

$\xrightarrow[\mathbf{X}_{11}]{(7)}$ (澳洲是与 北韩 有 邦交 的 \mathbf{X}_9 之一,

Australia is one of the \mathbf{X}_9 that have diplomatic relations with North Korea)

$\xrightarrow[\mathbf{X}_9]{(8)}$ (澳洲是与 北韩 有 邦交的 少数 国家 之一,

Australia is one of the few countries that have diplomatic relations with North Korea)

Weighted Grammar Rule

澳洲是与北韩有邦交的少数国家之一

Australia is one of the few countries that have diplomatic relations with North Korea

$X \rightarrow (\text{与 } X_1 \text{ 有 } X_2, \text{ have } X_2 \text{ with } X_1) +7.2$

$X \rightarrow (\text{北韩, North Korea}) -0.7$

$X \rightarrow (\text{邦交, diplomatic relations}) +3.9$

We are interested in a weighted
synchronous context-free
grammar (SCFG)



From where can we get the weights?

Score of Derivation

$$\begin{aligned} \textit{score}(D) = & \underbrace{\sum_{\mathbf{X} \rightarrow \langle \gamma, \alpha \rangle \in D} \textit{score}(\mathbf{X} \rightarrow \langle \gamma, \alpha \rangle)}_{\text{Grammar rules}} \\ & + \underbrace{\log q(e)}_{\text{Language model}} \end{aligned}$$



Optional

This decoding algorithm is more complicated (especially the "top-k" decoding algorithm).

Score of Derivation

$$\begin{aligned} \text{score}(D) = & \lambda_{\text{TM}} \times \underbrace{\sum_{\mathbf{X} \rightarrow \langle \gamma, \alpha \rangle \in D} \text{score}(\mathbf{X} \rightarrow \langle \gamma, \alpha \rangle)}_{\text{Grammar rules}} \\ & + \lambda_{\text{LM}} \times \underbrace{\log q(e)}_{\text{Language model}} \end{aligned}$$

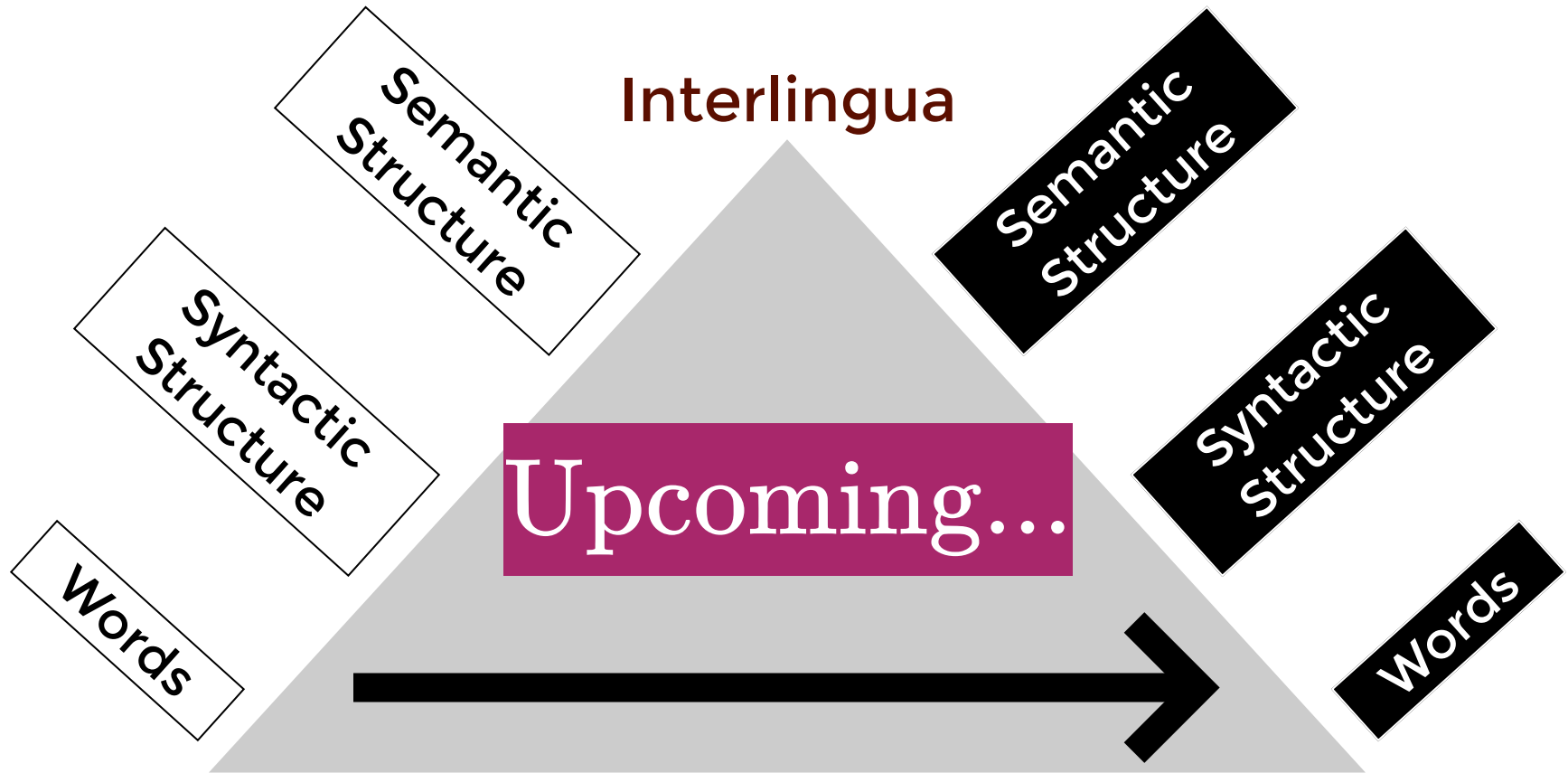
Similarly, we can weight different components differently using **hyperparameters**. Tune these hyperparameters on the development set, using a similar procedure as phrase-based translation.

Question

Can we make use of

word/context embeddings that
we learned earlier?

Machine Translation



Text-to-text problem with **neural networks**