



Ambiguity Resolution for Vt-N Structures in Chinese

Yu-Ming Hsieh^{1,2} Jason S. Chang² Keh-Jiann Chen¹

¹ Institute of Information Science, Academia Sinica, Taipei

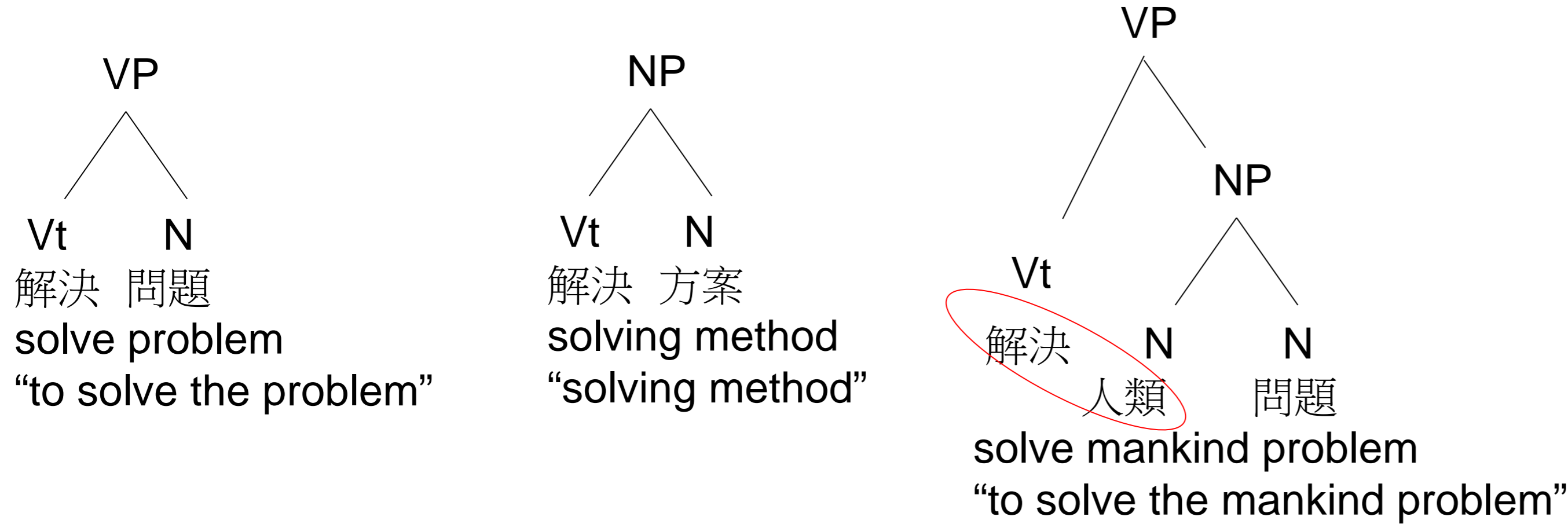
² Department of Computer Science, National Tsing-Hua University, Taiwan

morris@iis.sinica.edu.tw jason.jschang@gmail.com kchen@iis.sinica.edu.tw



Problem Statement and Analysis

- In Chinese, the structure of a transitive verb (Vt) immediately followed by a noun (N) may form a verb phrase (VP), a noun phrase (NP), or there may not be a dependent relation, as shown below.



- Analysis of NP(Vt-N) structures in the Sinica Treebank reveals the following four types of semantic structures.

Types	Examples
Telic(Vt) + Host(N)	研究/工具 research tool 探測/機 detective machine
Host-Event(Vt) + Attribute(N)	攻擊/策略 attacking strategy 書寫/內容 written context
Agentive(Vt) + Host(N)	叫/聲 shouting sound 炸/雞 fried chicken
Apposition(Vt) + Affair(N)	追撞/事故 collision accident 破壞/運動 destructive movement

Table 2. Semantic types of NP(Vt-N) and examples.

- Parsers generally prefer to the VP reading due to the statistical majority.

VP	NP	Other Relation
58%	16%	26%

Table 1. Statistical data from the Sinica Treebank.

- Unlike inflectional languages, Chinese verbs modify nouns without morphological inflection, e.g., 養殖/farming 池/pond.
- Linguistically motivated features to build a Vt-N classifier include lexical words, semantic knowledge, the morphological structure of verbs, neighboring parts-of-speech, and the syllabic length of words.

- Difficulties:

- Feature selection. Decision on which features to adopt and their combination is a difficult task in classification.
- Unknown word issue. Unknown word processing has technical problems that affect the prediction of the semantic types and morph-structures of unknown words.
- Data sparseness problem. Due to the limited size of the current Treebank, we should mine useful information from all available resources.

Proposed Models

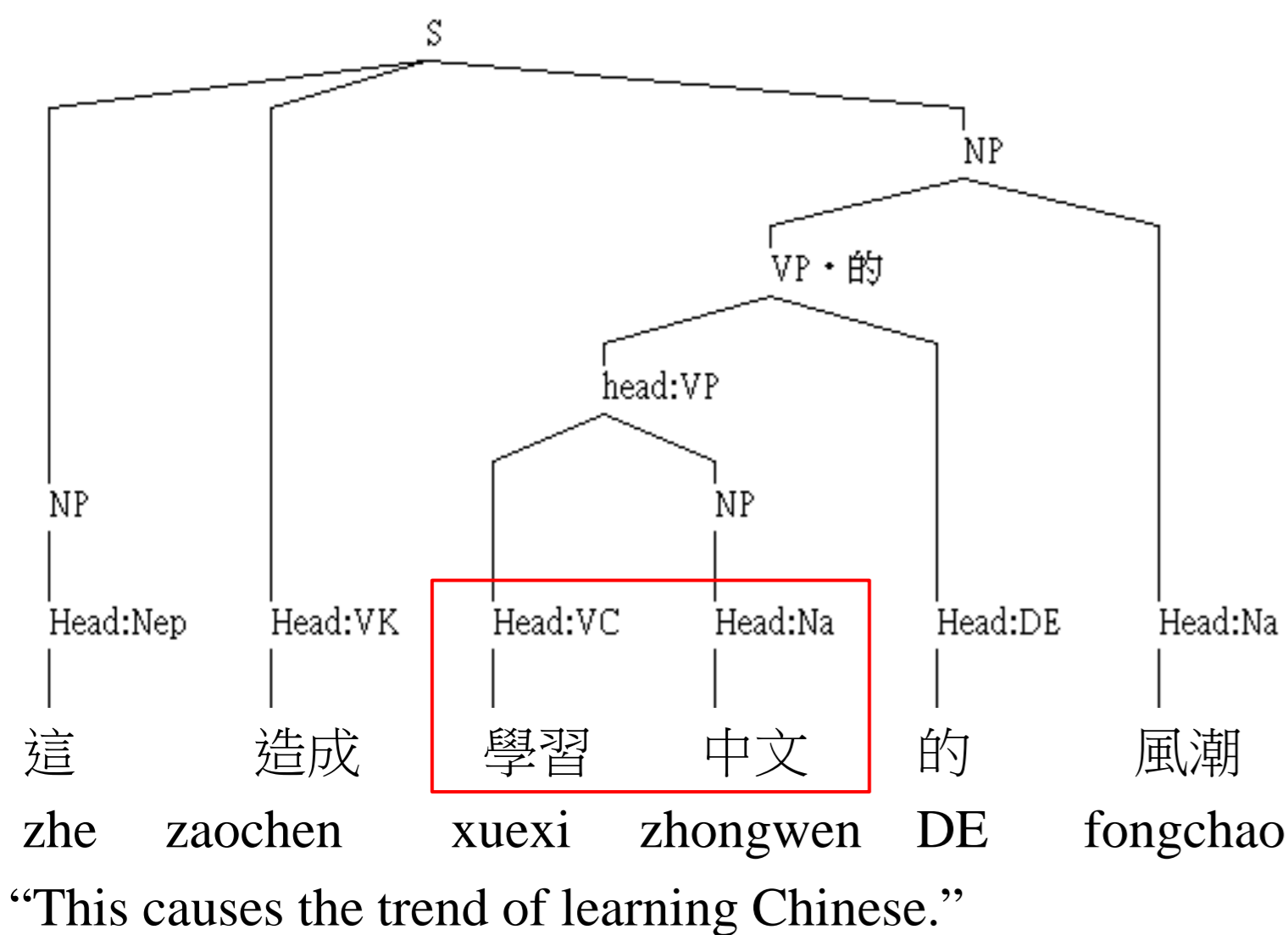
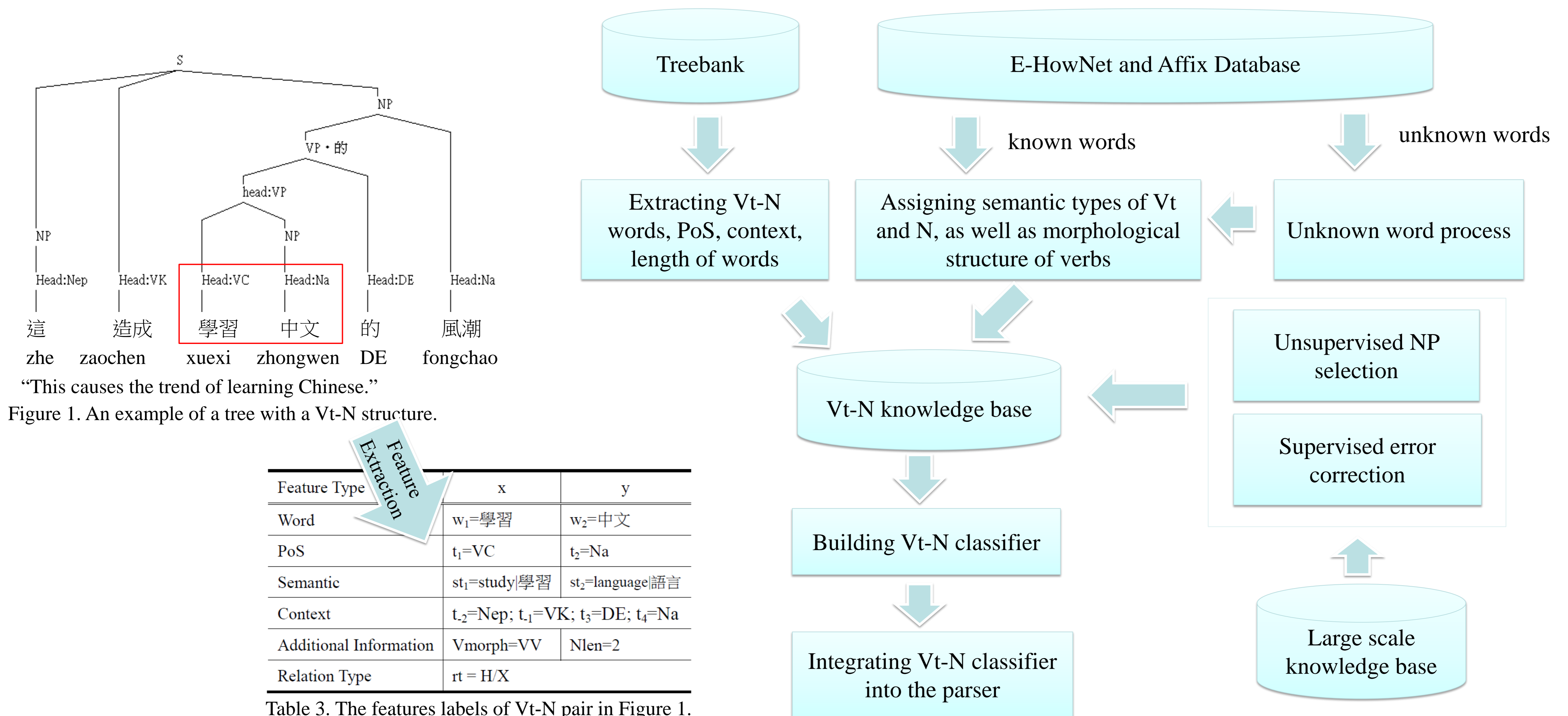


Figure 1. An example of a tree with a Vt-N structure.

Feature Type	x	y
Word	w ₁ =學習	w ₂ =中文
PoS	t ₁ =VC	t ₂ =Na
Semantic	st ₁ =study 學習	st ₂ =language 語言
Context	t ₂ =Nep; t ₁ =VK; t ₃ =DE; t ₄ =Na	
Additional Information	Vmorph=VV	Nlen=2
Relation Type	rt = H/X	

Table 3. The features labels of Vt-N pair in Figure 1.

Experiment Results

Evaluation of the Vt-N Classifier

- The results of Table 4 show that contextual information (M2) and lexical feature (M3) are the most important feature used to disambiguate VP, NP and independent structures. Results also demonstrate the benefits of using the semantic type (M4), verb morph-structure and noun length features (M5).

Models	Feature for Vt-N	P1(%)	P2(%)
M1	(t ₁ , t ₂)	61.94	59.10
M2	+(t ₁ , t ₁) (t ₂ , t ₃) (t ₂ , t ₁ , t ₁) (t ₂ , t ₃ , t ₄) (t ₁ , t ₃)	76.59	72.30
M3	+(w ₁ , t ₁ , w ₂ , t ₂) (w ₁ , w ₂) (w ₂) (w ₁)	83.55	80.20
M4	+(st ₁ , t ₁ , st ₂ , t ₂) (st ₁ , t ₁) (st ₂ , t ₂)	84.63	81.90
M5	+(Vmorph) (Nlen)	85.01	83.00

Table 4. The results of using different feature combinations. P1(%) is the 10-fold cross validation accuracy of the training data; P2(%) is the accuracy of the test data.

Knowledge from Large-scale Unlabeled Data

- We applied the data selection method (*distance=1*, with an intransitive verb followed by an object noun) and data correction method to learn more useful knowledge.
- The results in Table 5 show that the proposed methods can improve the accuracy.

	Treebank+ ASBC+Correction	Treebank+ ASBC-Vi-N	Treebank
size of training instances	56,521	46,258	8,017
M5 - P2(%)	88.40	83.90	83.00

Table 5. Experiment results of classifiers with different training data.

Vt-N classifier and PCFG Parser Integration

- We integrate the Vt-N models into the PCFG parser. The formula of the integrated structural evaluation model is as follows:

$$Score(T, S) = \sum_{i=1}^n (w_1 \times RP_i + w_2 \times VtNP_i)$$

	PCFG + M5 (Treebank+ASBC+Correction)	PCFG + M5 (Treebank)	PCFG
P2(%)	87.88	80.68	77.09
BF(%)	84.68	83.64	82.80

Table 6. The performance of the PCFG parser with and without model M5. The BF (*bracketed f-score*) is the parsing performance metric.

