

WEB-SCALE FEATURES FOR FULL-SCALE PARSING

Mohit Bansal and Dan Klein, ACL 2011, pp. 693-702

読み手: 岡崎 直観

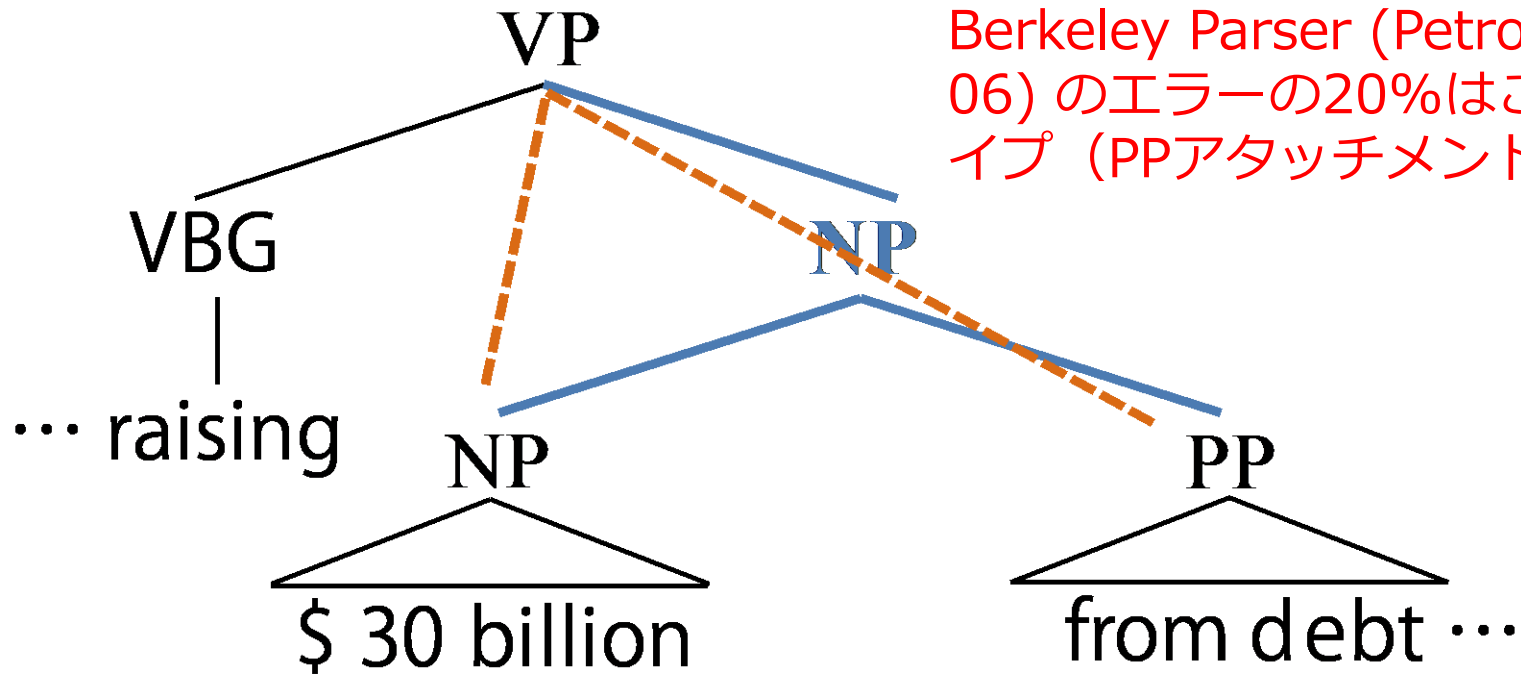
本論文について

- ラベル無しデータ（Web全体の統計量）を使い，常識のような知識を獲得し，言語処理のタスクを改善する例
 - 一見するとパーザの細かい話に見えてしまうが，述語項構造解析や固有表現抽出などのタスクにも使えそう
- 訓練データから文脈をマイニングするのも興味深い

研究の背景

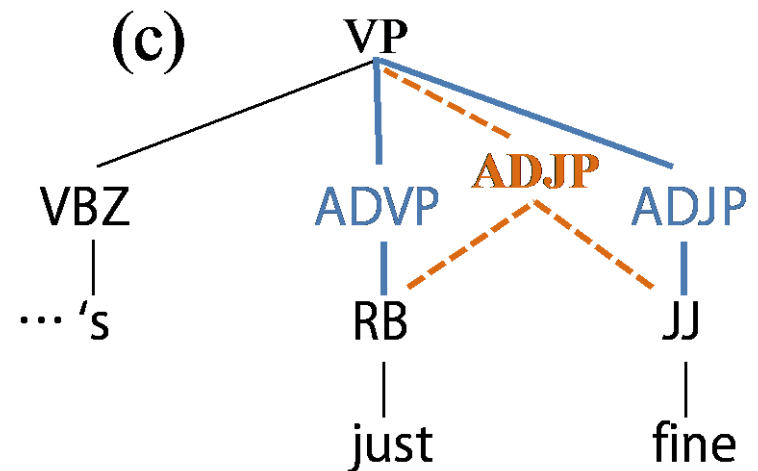
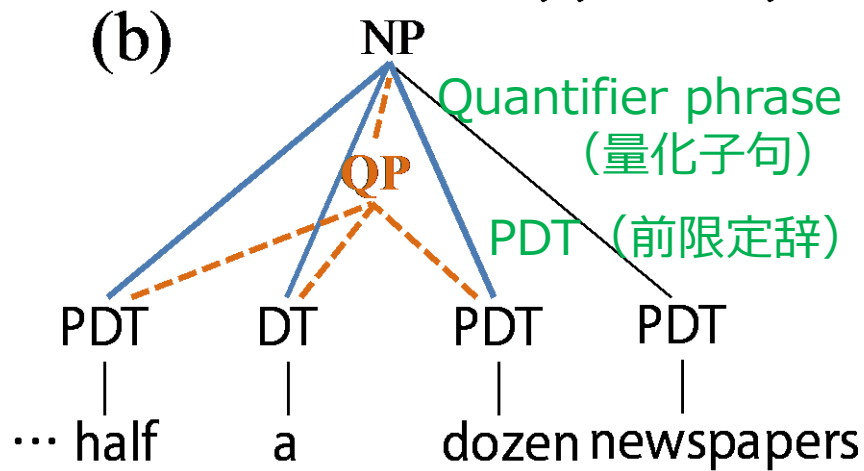
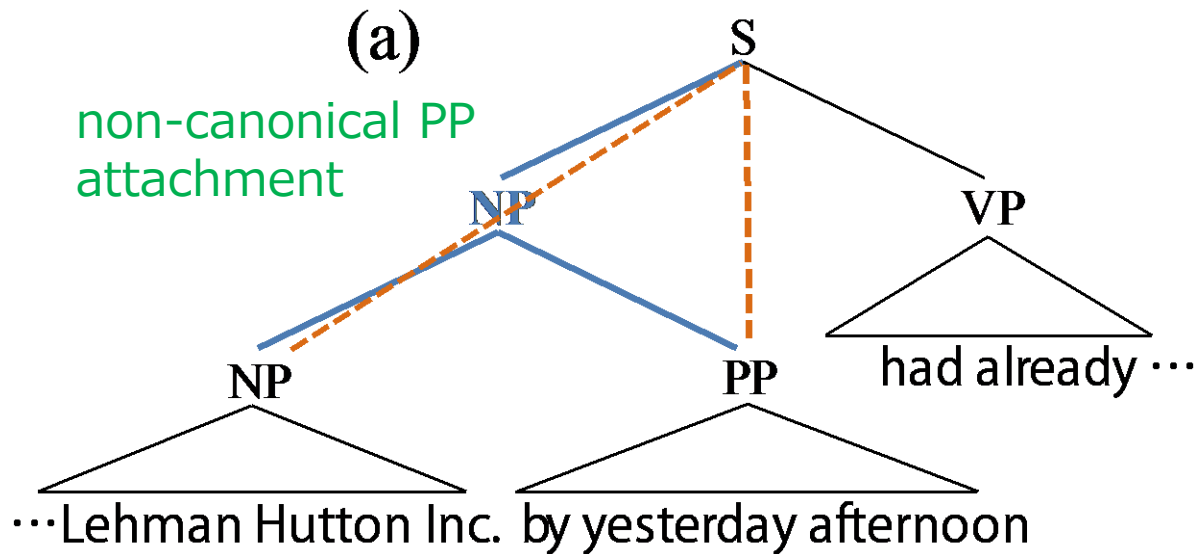
- 最近のパーザはF1で90を達成 (on Penn Treebank)
- アタッチメントなどの構造的なエラーが残されている
 - Berkeley parser (Petrov+ 06) のエラーの20%はPPアタッチメント (次スライド)
 - 他にも色々なタイプのエラーがある (次次スライド)
- **このようなエラーを100万語くらいの学習コーパスだけで改善するのは難しい**

典型的な解析誤り



The complex financing plan in the S&L bailout law includes **raising \$ 30 billion from debt** issued by the newly created RTC.

その他の解析誤り



アタッチメントエラーに対する従来手法

- 親和力 (Lauer 95, Volk 01, Lapata+ 04)
 - 「\$x billion from」よりも「raising from」がよく用いられる
 - 「hydrogen ion exchange」において「hydrogen exchange」よりも「hydrogen ion」の方がよく用いられる
- 言い換え表現を調べる (Nakov+ 05, Pitler+ 10)
 - 「raising it from」が用いられるなら「raising」 → 「from」
 - 「be \$x billion from」が用いられるなら「billion」 → 「from」
 - 「The \$x billion from」が用いられるなら「billion」 → 「from」

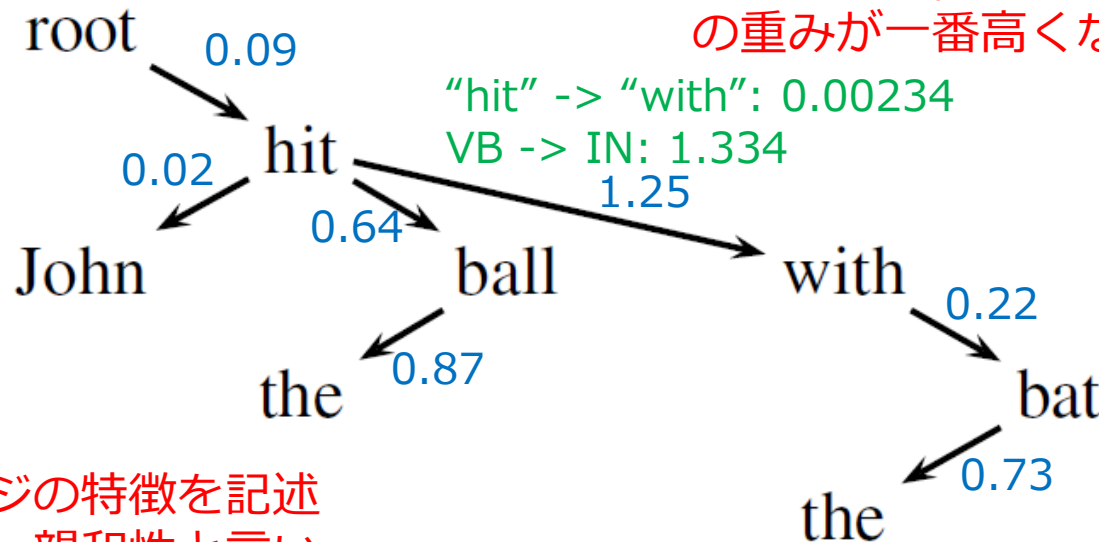
本研究の特色

- 任意のhead-argumentに対応しうる一般的な素性設計
 - 親和力を用いた素性は簡単, 言い換え素性は自明ではない
- 既存のパーザに組み込める実用的な素性設計
 - 係り受け解析器としてMST parser (McDonald+ 06)
 - 句構造解析器としてBerkeley parser (Petrov+ 06)
- Google n-gramコーパスの利用
 - 従来はWebページの検索ヒット件数を用いるものが多かった

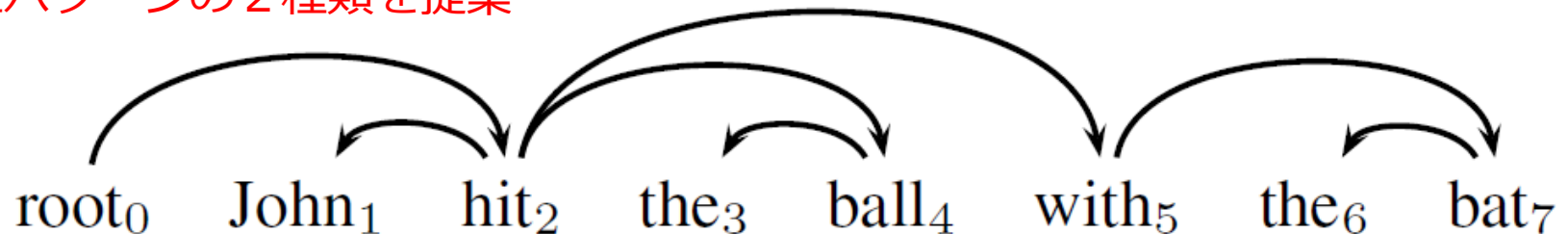
MSTParserの予備知識

- 構文解析: $t = \operatorname{argmax} w \cdot \Phi(x)$

基本的には、各エッジに素性関数を導入し、グラフ全体で素性の重みが一番高くなる木を選ぶ



本研究は、エッジの特徴を記述する素性として、親和性と言い換えパターンの2種類を提案



親和性 (=隣接性) 素性

- 2つの単語 h と a が近くにどのくらい出現するか？
 - $\Phi(h, a)$: h と a の係りやすさ
 - 文中で h の後に a が出現する場合: $q \in \{ha, h*a, h**a, h***a\}$
 - 文中で h の前に a が出現する場合: $q \in \{ah, a*h, a**h, a***h\}$
 - すべての q に対して
 - ADJ: q の出現回数
 - ADJ & POS(h) & POS(a): POS(h) の h と POS(a) の a が隣接する回数
 - ADJ & POS(x) & POS(y) & b : 上の素性の頻度を離散的にしたもの
 - $b = \text{floor}(\log_r(\text{count})/5) * 5$

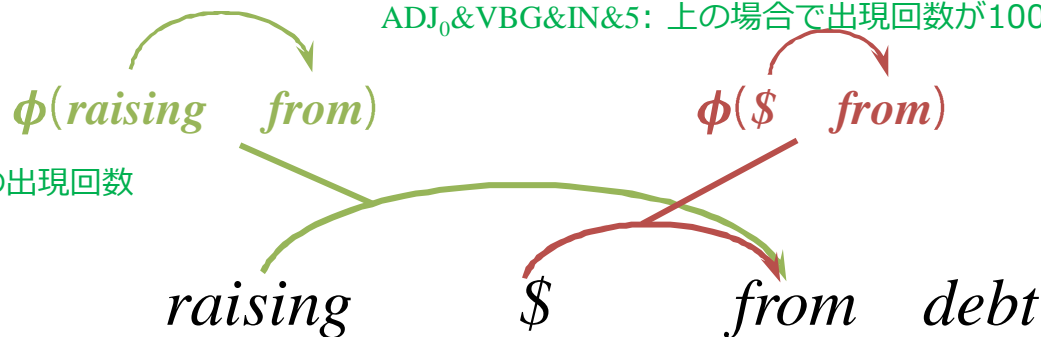
ADJ₀: "raising from" の出現回数

ADJ₀&VBG&IN: "raising" が VBG で from が IN の時の "raising from" の出現回数

ADJ₀&VBG&IN&5: 上の場合で出現回数が 1000 の場合

ADJ₁: "raising * from" の出現回数

...



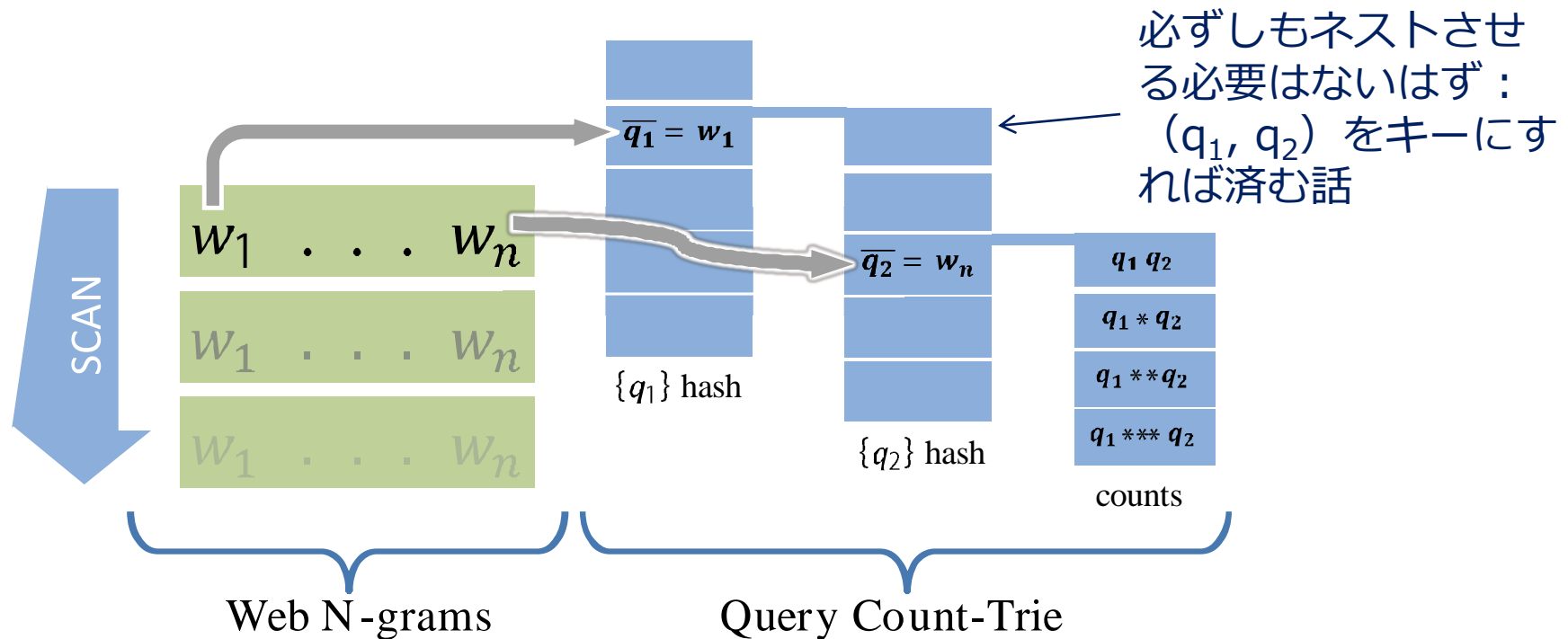
言い換えパターンのマイニング

- 従来手法ではテストに用いるパターンを手動で作成
 - 「VB it IN」が用いられるなら「VB」→「IN」
 - 「be NN IN」が用いられるなら「NN」→「IN」
 - 「The NN IN」が用いられるなら「NN」→「IN」
- 本研究ではパターンを自動生成する
 - 訓練データから正しい係りを収集: 例 (h, a) = (raising, from)
 - 次のパターンで周辺のコンテキスト c を探す: cha, hca, hac
 - 実際には h と a の品詞毎にマイニングをするらしい
 - 例えば“raising c from”で頻出パターンをマイニングすると, $c =$ “it”が見つかる

言い換え素性

- 単語 h と a が特定のパターンでどのくらい出現するか？
 - コンテキストの位置を次のように表現
 - $p = \text{"before"}$ if cha ; $p = \text{"middle"}$ if hca ; $p = \text{"after"}$ if hac
 - 次のようなテンプレートの素性を定義
 - $\text{PARA} \ \& \ \text{POS}(h) \ \& \ \text{POS}(a) \ \& \ c \ \& \ p \ \& \ \text{dir}$ ← h と a の出現順序: $\{\rightarrow, \leftarrow\}$
 - $\text{PARA} \ \& \ \text{POS}(h) \ \& \ \text{POS}(a) \ \& \ \text{POS}(c) \ \& \ p \ \& \ \text{dir}$ ← スペースになるかもしれないので
 - 例えば, hca で $c = \text{"it"}$ の場合
 - $\text{PARA} \ \& \ \text{POS}(h) \ \& \ \text{POS}(a) \ \& \ c \ \& \ p \ \& \ \text{dir}$
 - $\text{PARA} \ \& \ \text{VBG} \ \& \ \text{IN} \ \& \ \text{it} \ \& \ \text{middle} \ \& \ \rightarrow$: 動名詞に続けて it , 前置詞が出現する回数
 - 「raising it from」 や 「lowering it with」 などの表現で素性が共有される
 - $\text{PARA} \ \& \ \text{POS}(h) \ \& \ \text{POS}(a) \ \& \ \text{POS}(c) \ \& \ p \ \& \ \text{dir}$: $\text{PARA} \ \& \ \text{VBG}$
 - $\text{PARA} \ \& \ \text{VBG} \ \& \ \text{IN} \ \& \ \text{PRP} \ \& \ \text{middle} \ \& \ \rightarrow$: 動名詞に続けて代名詞, 前置詞が出現する回数
 - 「raising it from」 や 「providing him with」 などの表現で素性が共有される

データ構造



- 本研究の実験では、Google N-gramコーパスに450万件のクエリを発行したが、検索にかかったトータル時間は115分だった

係り受け解析の実験

- MSTParserに提案した素性を追加
- Penn Treebankで実験
 - 訓練データ: §2-21, 開発データ: §22, 評価データ: §23
 - 品詞タグ付けにはMXPOST (Ratnaparkhi 96) を用いた
 - 素性の次数は1次
- 評価結果
 - 提案手法を用いない2次の素性の場合と比べても改善

	Order 2	+ Web features	% Error Redn.
Dev (sec 22)	92.1	92.7	7.6%
Test (sec 23)	91.4	92.0	7.0%

句構造解析の実験

- Berkeley parserのtop- k 解析木をリランキング
 - リランキングの学習器として平均化パーセプトロンを用いた
 - 素性
 - Berkeley parserが出力した確率値: $\log p(t|w)$
 - 提案手法の素性 (Web features)
 - 非局所情報の素性 (Configurational features)
 - Top- k をとることの可能性 (oracleスコア)

	$k = 1$	$k = 2$	$k = 10$	$k = 25$	$k = 50$	$k = 100$
Dev	90.6	92.3	95.1	95.8	96.2	96.5
Test	90.2	91.8	94.7	95.6	96.1	96.4

← Top-100の中から正しい木を選べばウハウハ

評価結果

Parsing Model	Dev (sec 22)		Test (sec 23)	
	F1	EX	F1	EX
Baseline (1-best)	90.6	39.4	90.2	37.3
$\log p(t w)$	90.4	38.9	89.9	37.3
+ Web features	91.6	42.5	91.1	40.6
+ Configurational features	91.8	43.8	91.1	40.6
+ Web + Configurational	92.1	44.0	91.4	41.4

非局所情報と同程度の貢献

非局所情報と提案手法を組み合わせると、さらに改善

解析誤りの減少

- PPアタッチメント (IN) で改善の効果を確認
- その他のタイプのアタッチメントでも改善が確認できる

Arg Tag	# Attach	Baseline	This Work	% ER
NN	5725	5387	5429	12.4
NNP	4043	3780	3804	9.1
IN	4026	3416	3490	12.1
DT	3511	3424	3429	5.8
NNS	2504	2319	2348	15.7
JJ	2472	2310	2329	11.7
CD	1845	1739	1738	-0.9
VBD	1705	1571	1580	6.7
RB	1308	1097	1100	1.4
CC	1000	855	854	-0.7
VB	983	940	945	11.6
TO	868	761	776	14.0
VBN	850	776	786	13.5
VBZ	705	633	629	-5.6
PRP	612	603	606	33.3

Baselineは2次のMSTParser

重用された親和性素性の例

POS _{head}	POS _{arg}	Example (head, arg)
RB	IN	<i>back</i> → <i>into</i>
NN	IN	<i>review</i> → <i>of</i>
NN	DT	<i>The</i> ← <i>rate</i>
NNP	IN	<i>Regulation</i> → <i>of</i>
VB	NN	<i>limit</i> → <i>access</i>
VBD	NN	<i>government</i> ← <i>cleared</i>
NNP	NNP	<i>Dean</i> ← <i>Inc</i>
NN	TO	<i>ability</i> → <i>to</i>
JJ	IN	<i>active</i> → <i>for</i>
NNS	TO	<i>reasons</i> → <i>to</i>
IN	NN	<i>under</i> → <i>pressure</i>
NNS	IN	<i>reports</i> → <i>on</i>
NN	NNP	<i>Warner</i> ← <i>studio</i>
NNS	JJ	<i>few</i> ← <i>plants</i>

← 名詞が前置詞に係るかどうかは、hとaの親和性から判断しやすい

← 動詞が名詞に係る（目的語をとる）かどうかは、hとaが隣接しやすいかどうかを考慮するとよい

重用された言い換え素性の例

POS _{head}	mid-word	POS _{arg}	Example (head, arg)	bfr-word	POS _{head}	POS _{arg}	Example (head, arg)
VBN	<i>this</i>	IN	<i>leaned, from</i>	<i>second</i>	NN	IN	<i>season, in</i>
VB	<i>this</i>	IN	<i>publish, in</i>	<i>The</i>	NN	IN	<i>role, of</i>
VBG	<i>him</i>	IN	<i>using, as</i>	<i>strong</i>	NN	IN	<i>background, in</i>
VBG	<i>them</i>	IN	<i>joining, in</i>	<i>our</i>	NNS	IN	<i>representatives, in</i>
VBD	<i>directly</i>	IN	<i>converted, into</i>	<i>any</i>	NNS	IN	<i>rights, against</i>
VBD	<i>held</i>	IN	<i>was, in</i>	<i>A</i>	NN	IN	<i>review, of</i>
VBN	<i>jointly</i>	IN	<i>offered, by</i>	<i>:</i>	NNS	IN	<i>Results, in</i>
VBZ	<i>it</i>	IN	<i>passes, in</i>	<i>three</i>	NNS	IN	<i>years, in</i>
VBG	<i>only</i>	IN	<i>consisting, of</i>	<i>In</i>	NN	IN	<i>return, for</i>
VBN	<i>primarily</i>	IN	<i>developed, for</i>	<i>no</i>	NN	IN	<i>argument, about</i>
VB	<i>us</i>	IN	<i>exempt, from</i>	<i>current</i>	NN	IN	<i>head, of</i>
VBG	<i>this</i>	IN	<i>using, as</i>	<i>no</i>	NNS	IN	<i>plans, for</i>
VBD	<i>more</i>	IN	<i>looked, like</i>	<i>public</i>	NN	IN	<i>appearance, at</i>
VB	<i>here</i>	IN	<i>stay, for</i>	<i>from</i>	NNS	IN	<i>sales, of</i>
VBN	<i>themselves</i>	IN	<i>launched, into</i>	<i>net</i>	NN	IN	<i>revenue, of</i>
VBG	<i>down</i>	IN	<i>lying, on</i>	<i>,</i>	NNS	IN	<i>names, of</i>
				<i>you</i>	NN	IN	<i>leave, in</i>
				<i>have</i>	NN	IN	<i>time, for</i>
				<i>some</i>	NN	IN	<i>money, for</i>
				<i>annual</i>	NNS	IN	<i>reports, on</i>

結論

- ラベル付けされていないコーパスの統計量を，教師あり学習で構築される構文解析の枠組みに取り込めた
- 構文解析に効果的なパターンを自動的にマイニングすることができた
- 従来の機械学習アプローチでは解決することが難しかった解析誤りを改善できた