

# 言語学的に妥当な CCGツリーバンク構築に向けて

2023年3月17日 言語処理学会第29回年次大会併設ワークショップJLR2023

富田朝<sup>1</sup> 谷中瞳<sup>2</sup> 戸次大介<sup>1</sup>

<sup>1</sup>お茶の水女子大学

<sup>2</sup>東京大学

# 目次

## 研究概要

1. 組合せ範疇文法
2. ツリーバンク
3. 背景と目的

## 先行研究

1. 日本語CCGbank
2. ABCツリーバンク
3. lightblue

## 提案手法

1. リフォーミングの流れ
2. リフォーミングの詳細
3. Left-corner chart parsing

## エラー分析

1. エラーケース1
2. エラーケース2

## まとめ

1. まとめと展望

01

# 研究概要

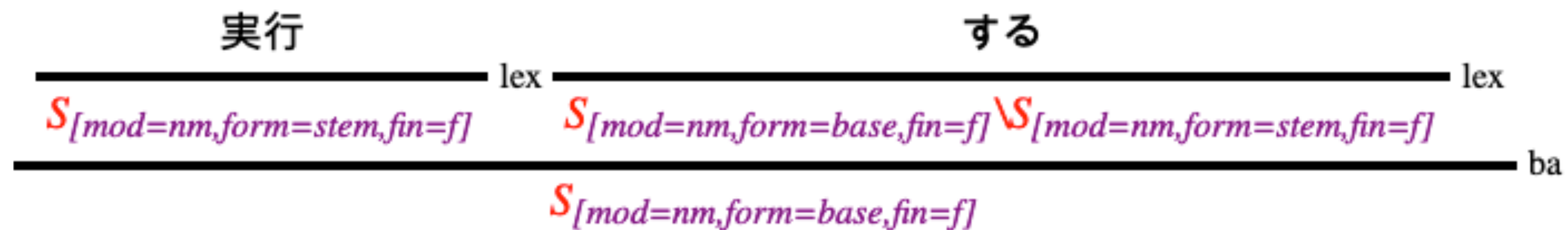
## ツリーバンク

各文に**統語構造**が付与されている**コーパス**

※コーパス：新聞や本などの内容をテキスト化したデータを大量に集めたデータベース

※統語構造：語同士がどのような順番や規則で結合しているのかを表す構造

例：日本語CCGbank [Uematsu+; 2013]



# 組合せ範疇文法 (Combinatory Categorical Grammar; CCG)

語彙化文法の一種

**辞書** (統語構造や意味に関する情報の記述)

例: Keats     $\vdash NP$   
      eats     $\vdash (S \setminus NP) / NP$   
      apples  $\vdash NP$

+

**組み合わせ規則**

例: 関数適用規則  
      関数合成規則  
      型繰り上げ規則

## 背景と目的

CCGパーザは学習・評価データとしてCCGツリーバンクを利用する

CCGパーザの妥当性はCCGツリーバンクの妥当性に依存する

言語学的に妥当なCCGツリーバンクの構築が必要である

→Bekki (2010) に基づいている

02

# 先行研究

# 先行研究

01

日本語CCGbank  
[Uematsu+ 2013]

02

ABCツリーバンク  
[Kubota+ 2019]

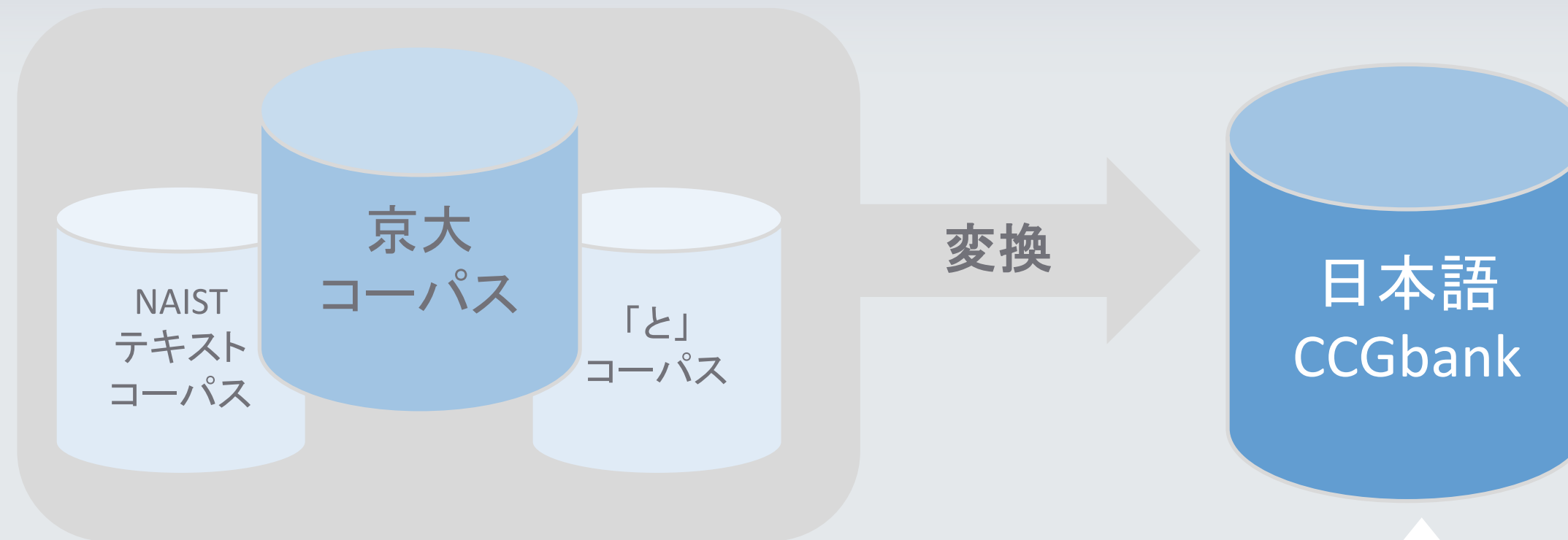
03

lightblue  
[Bekki+ 2016]



## 日本語CCGbank [Uematsu+, 2013]

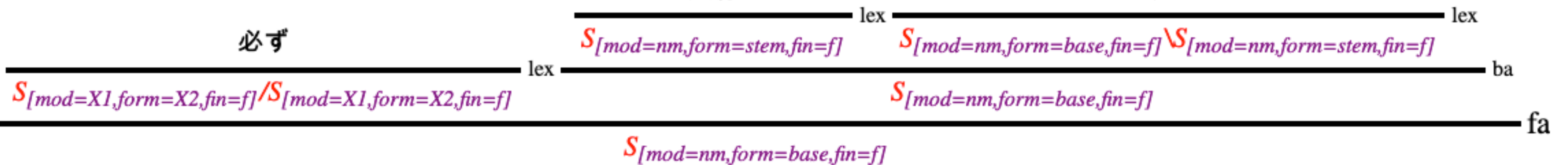
1/3



実行

する

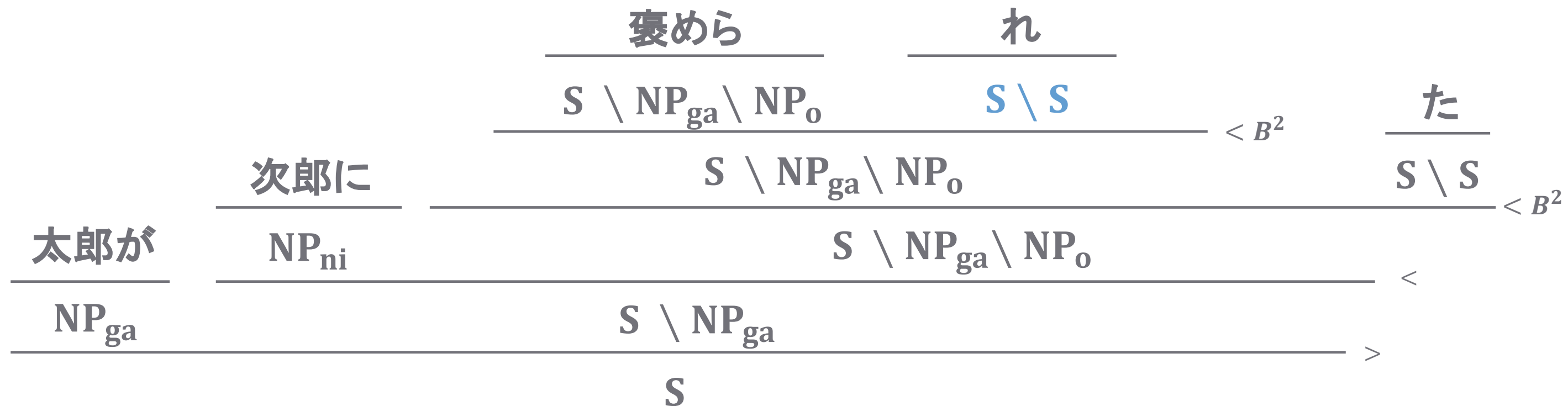
必ず



## 課題

受身・使役の構文に対して誤った分析がなされている。  
[Bekki & Yanaka, 2023]

## 日本語CCGbankの誤ったCCG木



## 課題

受身・使役の構文に対して誤った分析がなされている。  
[Bekki & Yanaka, 2023]

## 日本語CCGbankの分析

れ  
—  
S \ S

## 言語学的に正しいCCG木

褒めら	れ	
$S \setminus NP_{ga} \setminus NP_o$	$S \setminus NP_{ga} \setminus NP_{ni} \setminus (S \setminus NP_{ga} \setminus NP_{ni o})$	$<$
次郎に	$S \setminus NP_{ga} \setminus NP_{ni}$	$< B^2$
太郎が	$S \setminus NP_{ga} \setminus NP_{ni}$	$>$
$T / (T \setminus NP_{ni})$	$S \setminus NP_{ga}$	$>$
$T / (T \setminus NP_{ga})$	S	

# 先行研究

01

日本語CCGbank  
[Uematsu+ 2013]

02

ABCツリーバンク  
[Kubota+ 2019]

03

lightblue  
[Bekki+ 2016]

- ABC文法を用いたツリーバンク

## ABC文法

関数適用のみからなるAB文法 + 関数合成(Function Composition)規則

### 関数適用

$$A / B \quad B \Rightarrow A$$

$$B \quad B \setminus A \Rightarrow A$$

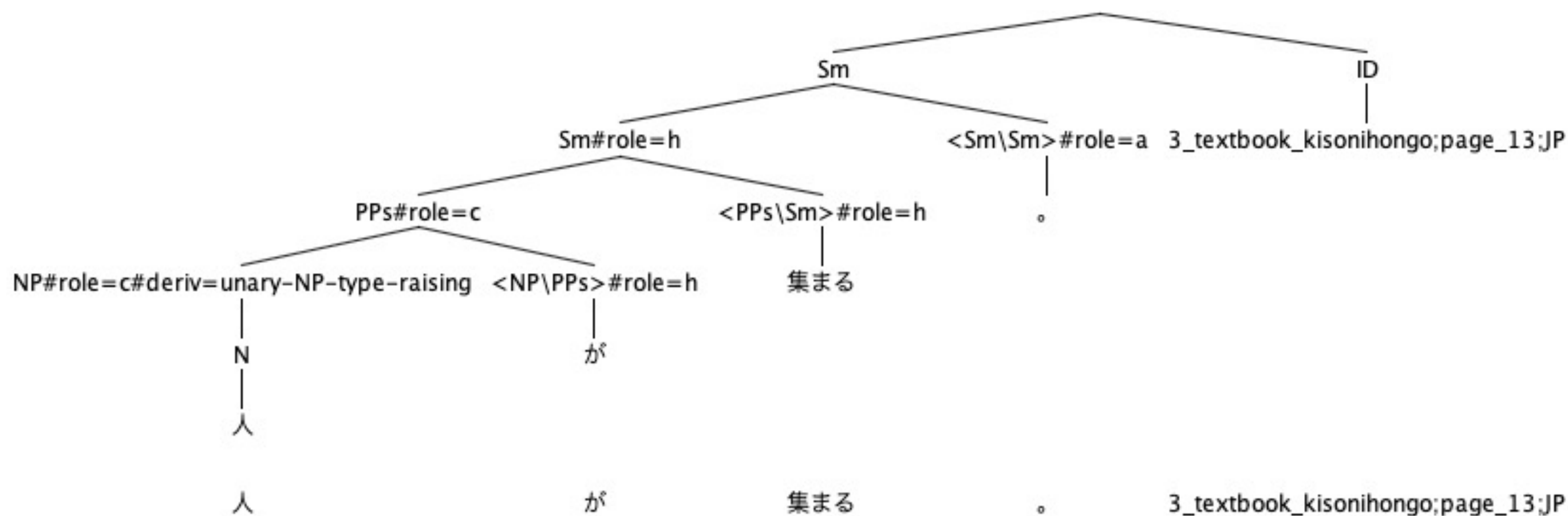
### 関数合成

$$A / B \quad B / C \Rightarrow A / C$$

$$C \setminus B \quad B \setminus A \Rightarrow C \setminus A$$

## ABCツリーバンク [Kubota+, 2019]

2/5

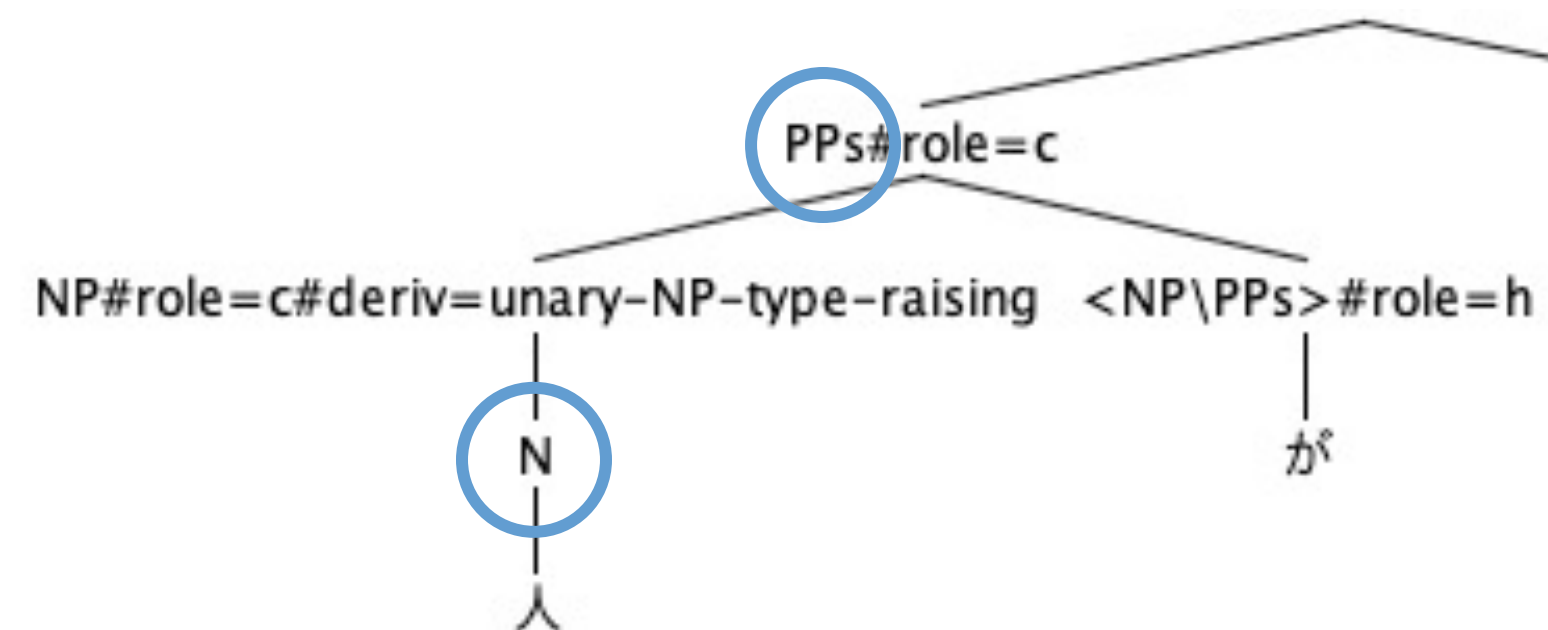


## ABCツリーバンク [Kubota+, 2019]

3/5

## ABC文法の原子的カテゴリ

統語範疇	名称	素性
CP	節	f, t, t-sbj, q, q-sbj, x
FRAG	断片	
<b>N</b>	<b>名詞</b>	<b>s</b>
NP	名詞句	q
NPR	固有名詞句	
NUM	数量詞	
NUMCLP	数量詞句	
<b>PP</b>	<b>後置詞句</b>	<b>s, s2, o1, o2</b>
PRO	代名詞	
Q	量化詞	
S	節 (IP相当)	a, e, imp, m, nml, rel, smc, sub
WNUM	疑問数量詞	
WPRO	疑問代名詞	



## 特徴

日本語CCGbankの項構造などの課題点の多くが改善された

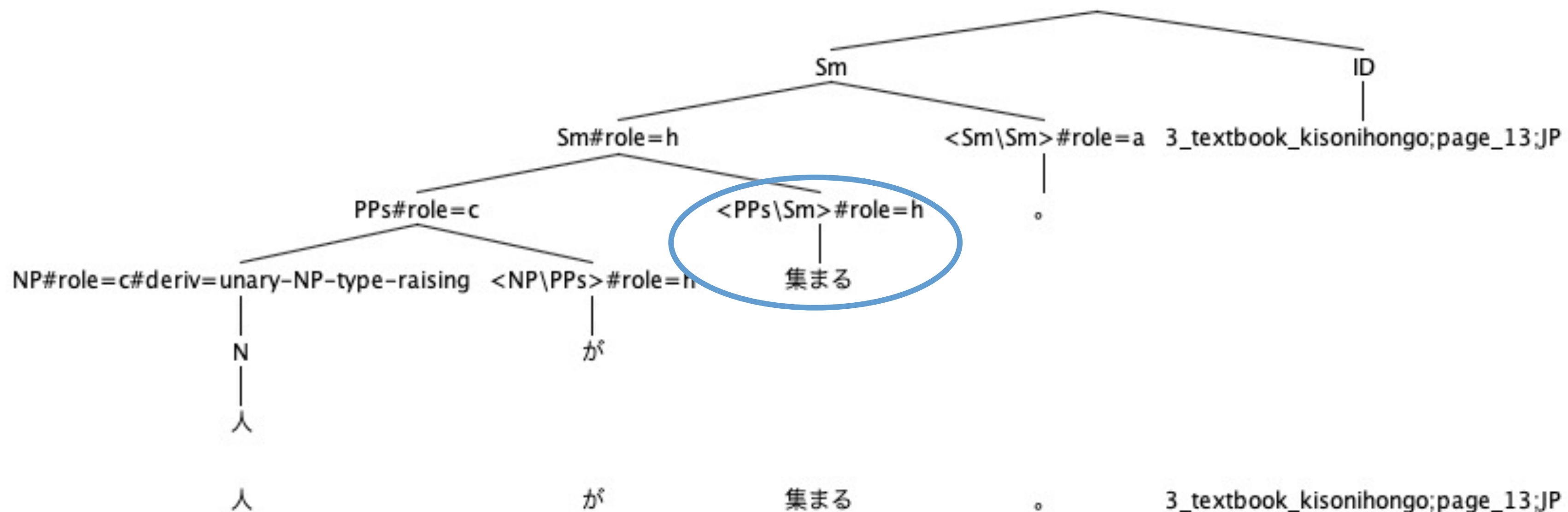
## 課題

活用の種別などの統語情報が不十分



## ABCツリーバンク [Kubota+, 2019]

5/5



# 先行研究

01

日本語CCGbank  
[Uematsu+ 2013]

02

ABCツリーバンク  
[Kubota+ 2019]

03

lightblue  
[Bekki+ 2016]

## lightblue [Bekki+, 2016]

1/5

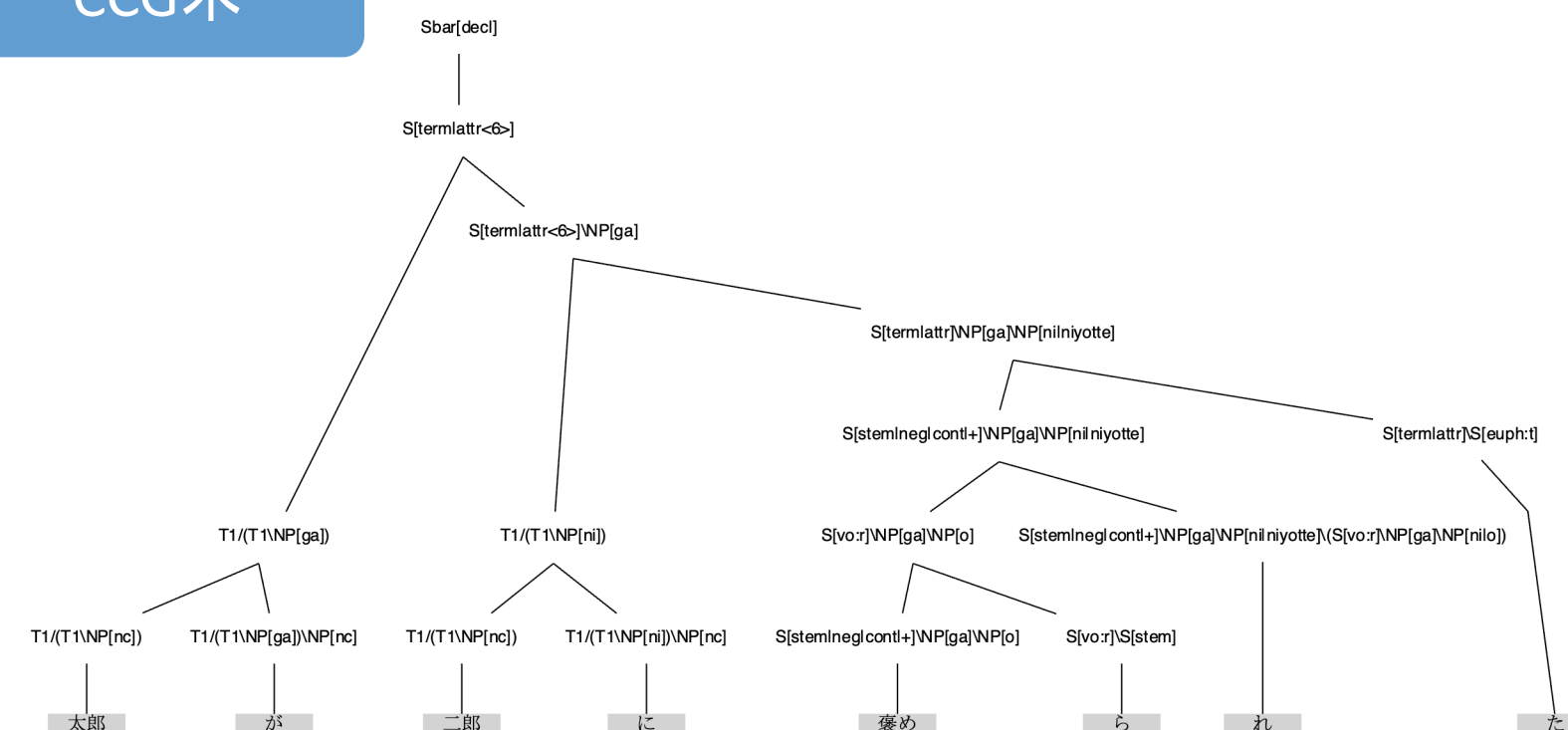
- CCG統語解析器
- 格フレームを元に作られた語彙辞書とCCGの組合せ規則に基づいて CCG統語解析を行う

自然言語

「太郎が二郎に  
褒められた」

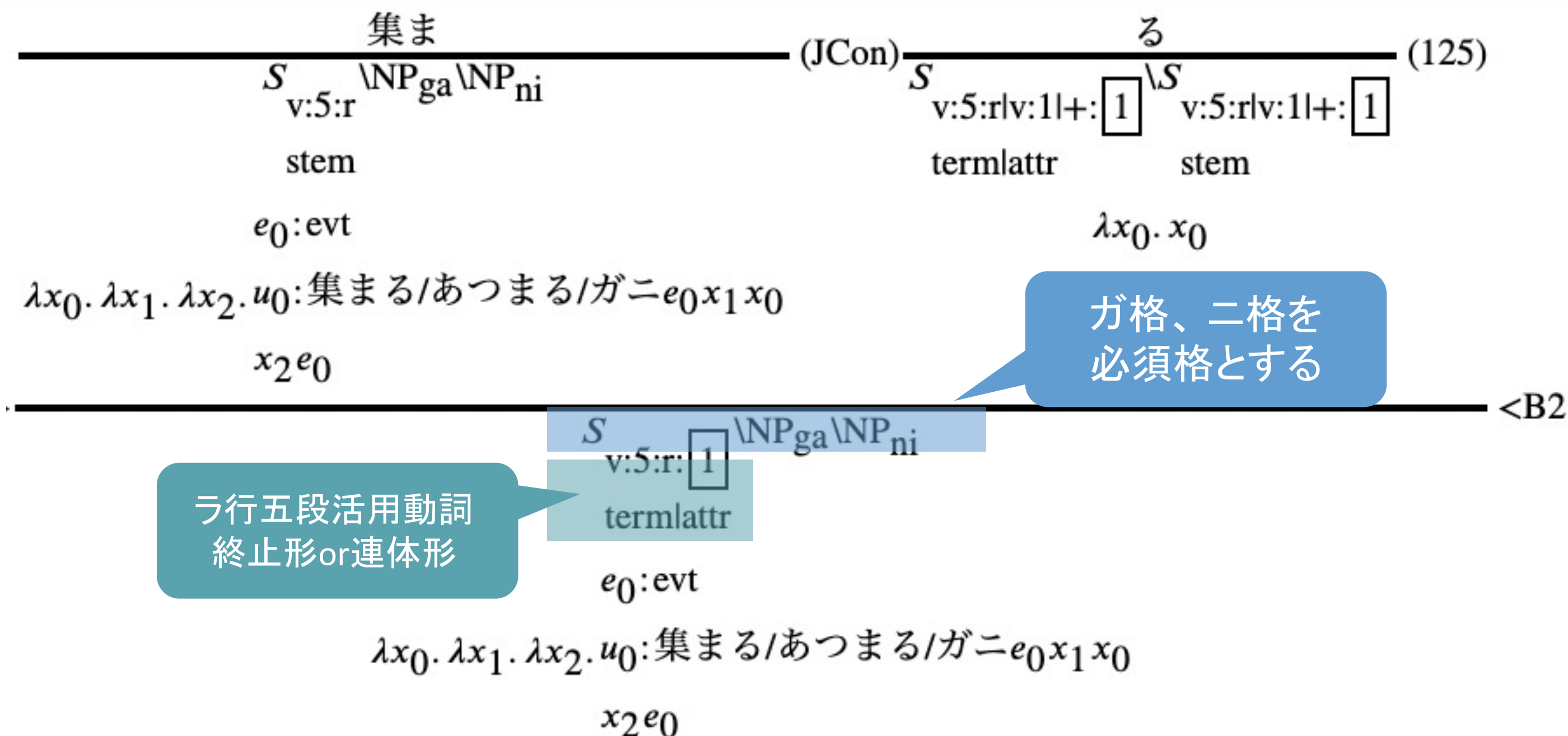
lightblue

CCG木



## 特徴

詳細な統語素性を含むCCG木を出力する

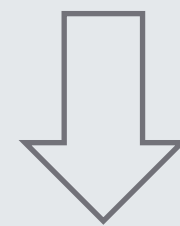


## 特徴

詳細な統語素性を含むCCG木を出力する

## 課題

項構造に関する誤りが多い



- 本来自然ではない用言が選択される
- 用言の格フレームに間違いが含まれる

03

# 提案手法

## 提案手法

### ABCツリーバンクの特徴

項構造が人手によって  
正確に記述されている

+

### lightblueの特徴

出力に詳細な統語素性を  
与えることができる

言語学的に妥当で詳細な情報を持った  
ツリーバンクを構築する



# ABCツリーバンクのリフォーミングの流れ

リフォーミング：ツリーバンクを分解し再構築する手法

1. ABCツリーバンクのパーズ

2. lightblueの辞書の書き換え

3. ツリーバンクの再構築

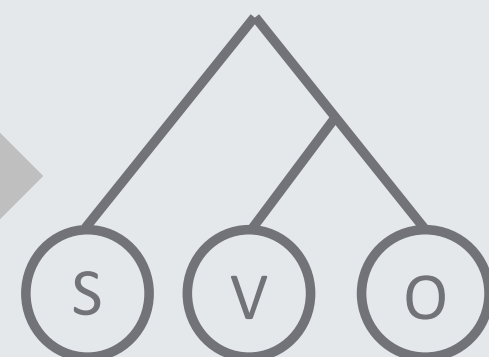
# ABCツリーバンクのリフォーミングの詳細

## 1. ABCツリーバンクのパーズ

入力

ABC  
ツリーバンク

ABCパーズ



用言を抽出



各文を抽出

フィルタリング

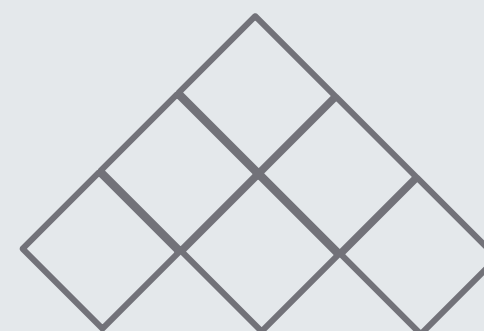
## 2. lightblueの辞書の書き換え

格  
フレーム

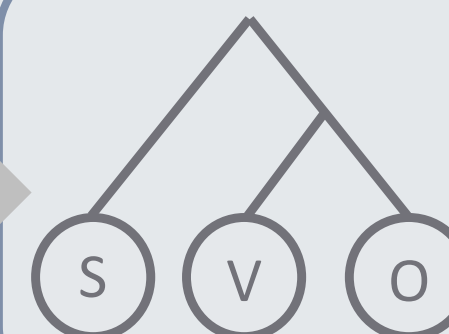
lightblue用  
辞書

辞書から  
選択された  
語彙項目

チャート  
パーズ



## 3. ツリーバンクの再構築



出力

日本語  
CCGツリーバンク

1. ABCツリーバンクのパーズ

2. lightblueの辞書の書き換え

3. ツリーバンクの再構築

## リフォーミング -ABCツリーバンクのパーズ

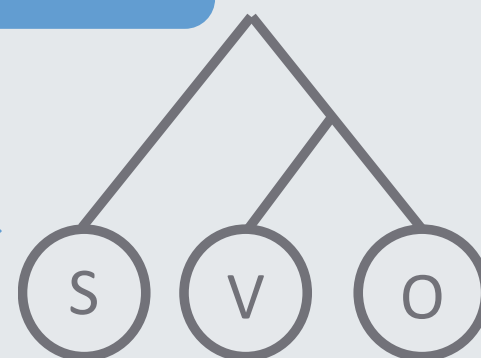
2/3

## 1. ABCツリーバンクのパーズ

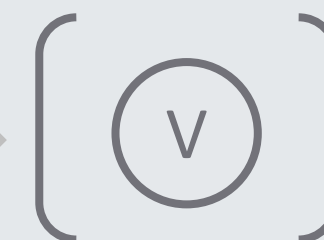
入力

ABC  
ツリーバンク

パーズ



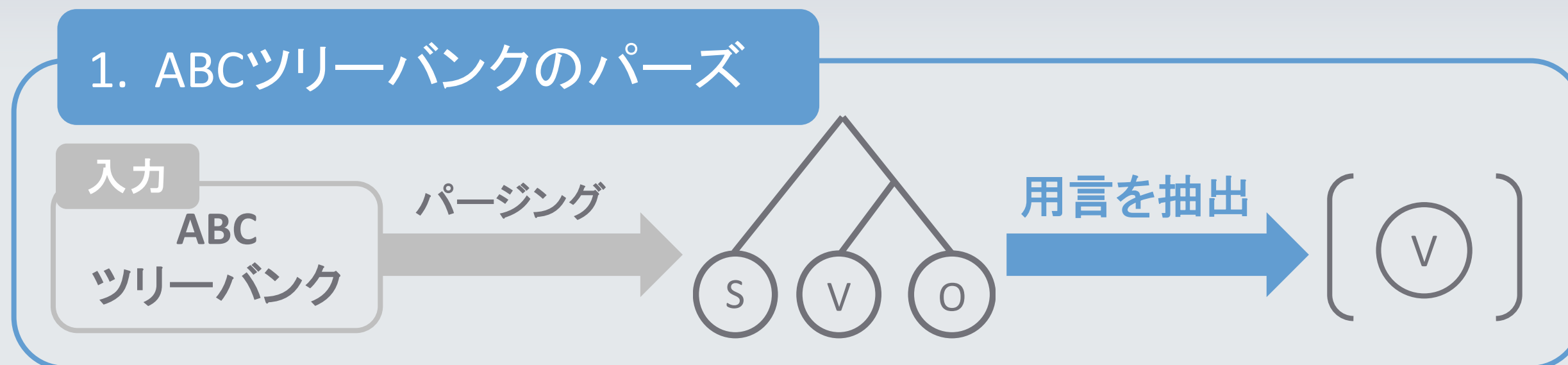
用言を抽出



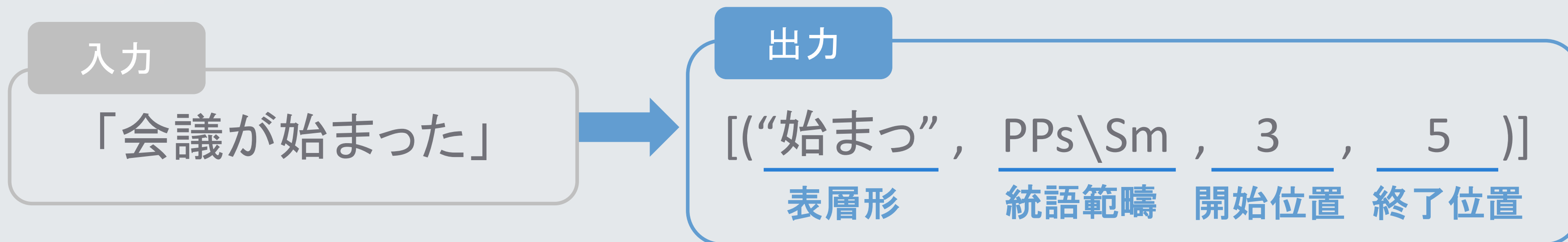
ABCツリーバンクのパーズを作成し木構造データ  
で扱えるようにする

## リフォーミング -ABCツリーバンクのパーズ

3/3



例



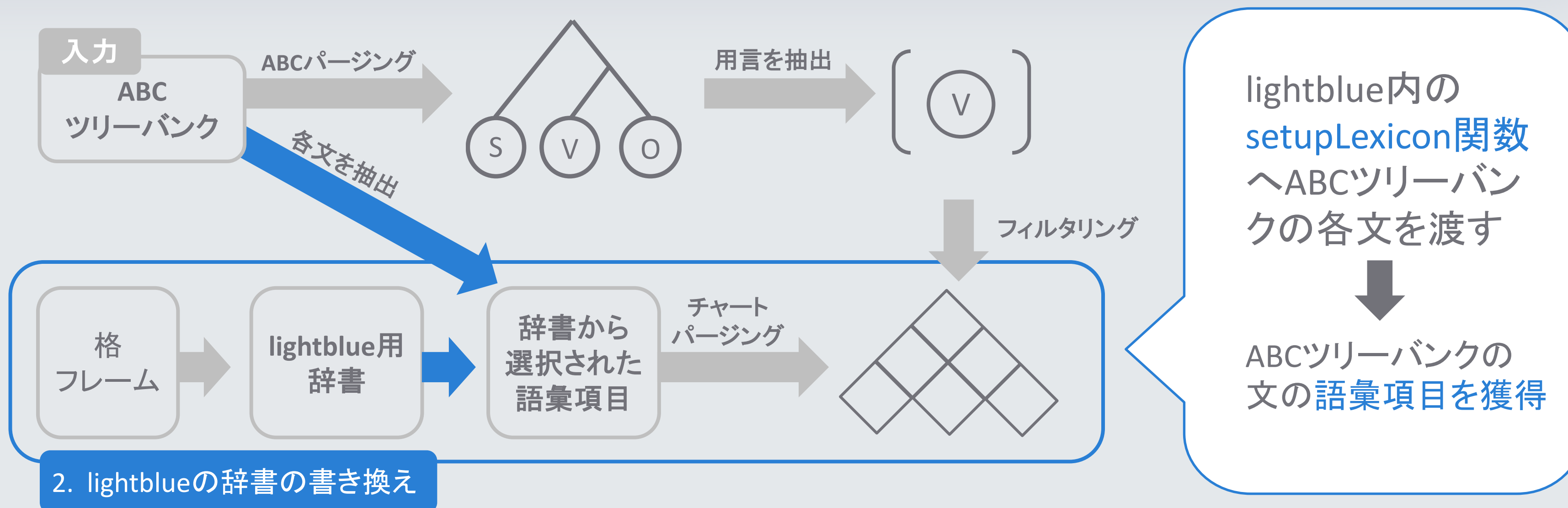
1. ABCツリーバンクのパーズ

2. lightblueの辞書の書き換え

3. ツリーバンクの再構築

## リフォーミング – lightblueの辞書の書き換え

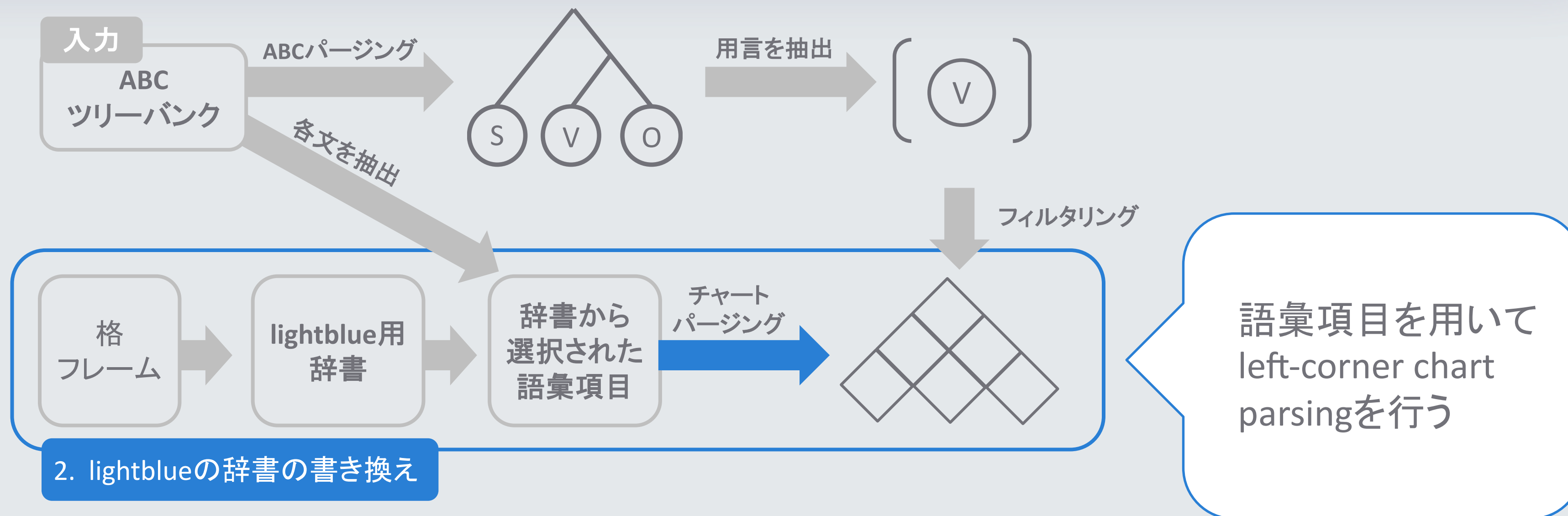
2/9



setupLexicon関数: 文を受け取ると、受け取った文を解析するために必要な語彙項目のリストを返す関数

# リフォーミング – lightblueの辞書の書き換え

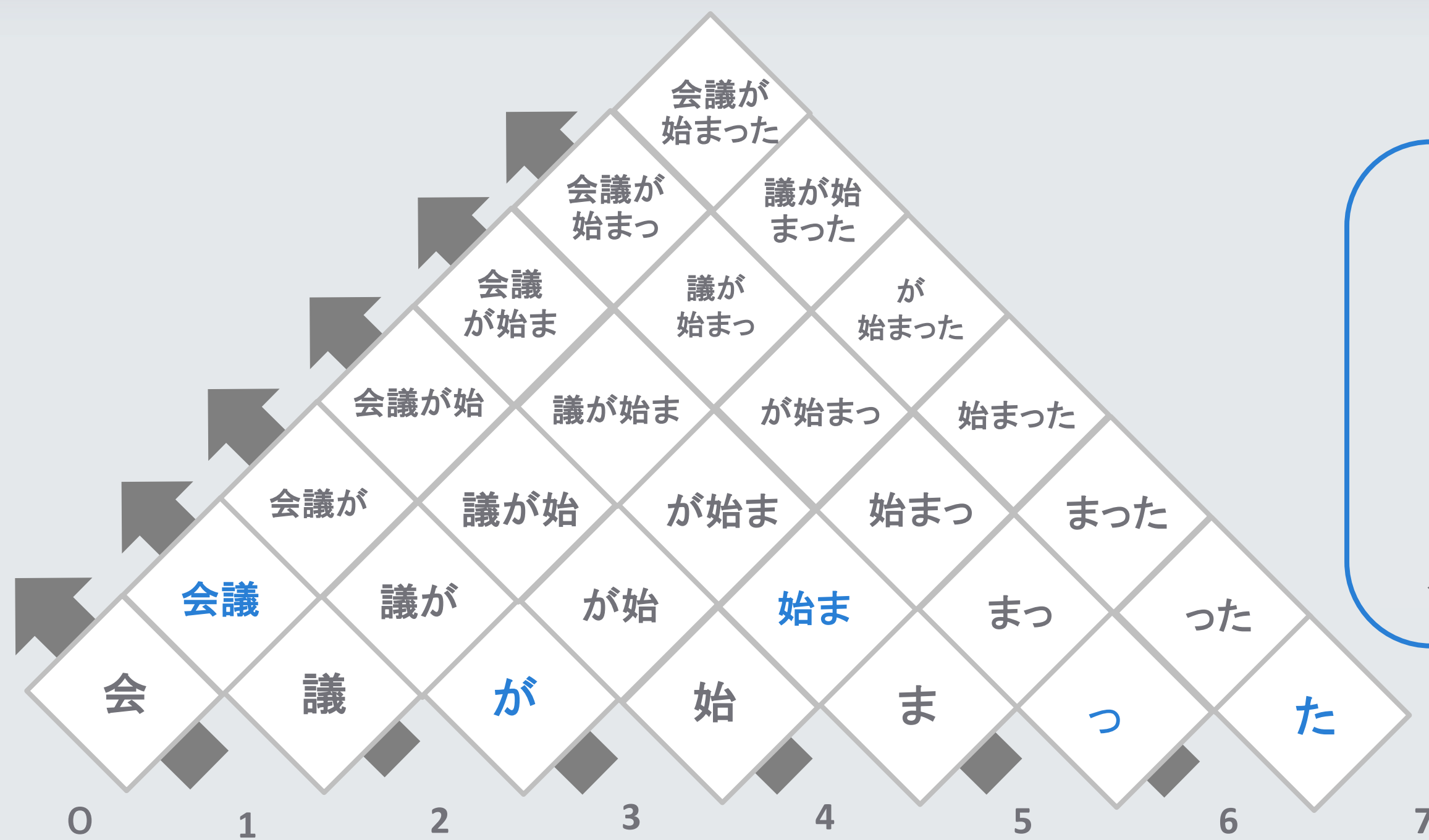
3/9





# Left-corner chart parsing

4/9



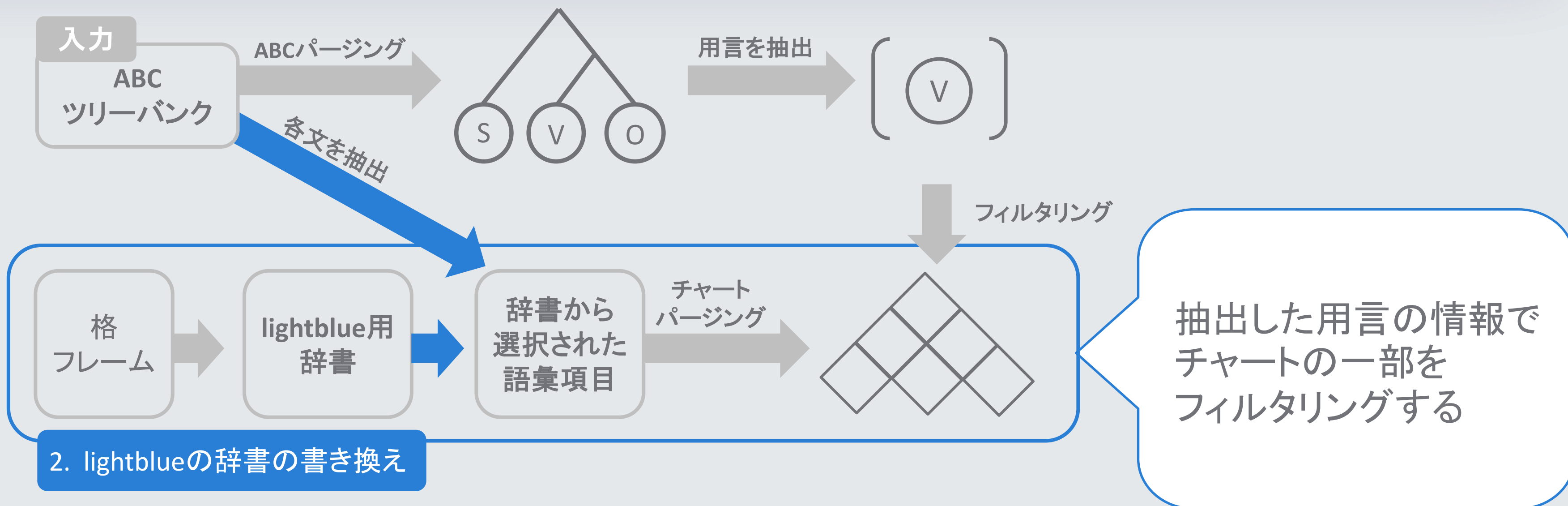
抽出した文の部分文字列  
全ての組み合わせにスコアをつける



スコアの高い組み合わせを出力する

# リフォーミング – lightblueの辞書の書き換え

5/9



## リフォーミング – lightblueの辞書の書き換え

6/9

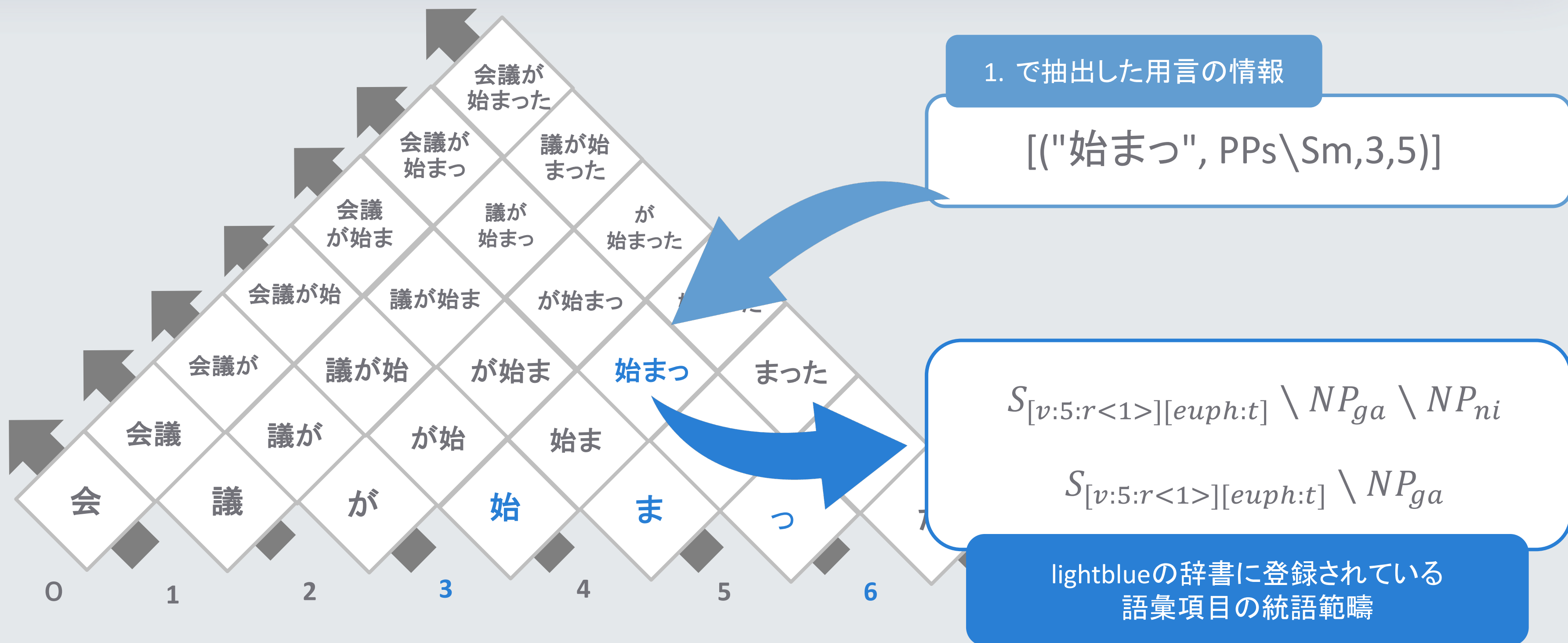


1. で抽出した用言の情報

[("始まっ", PPs\Sm,3,5)]

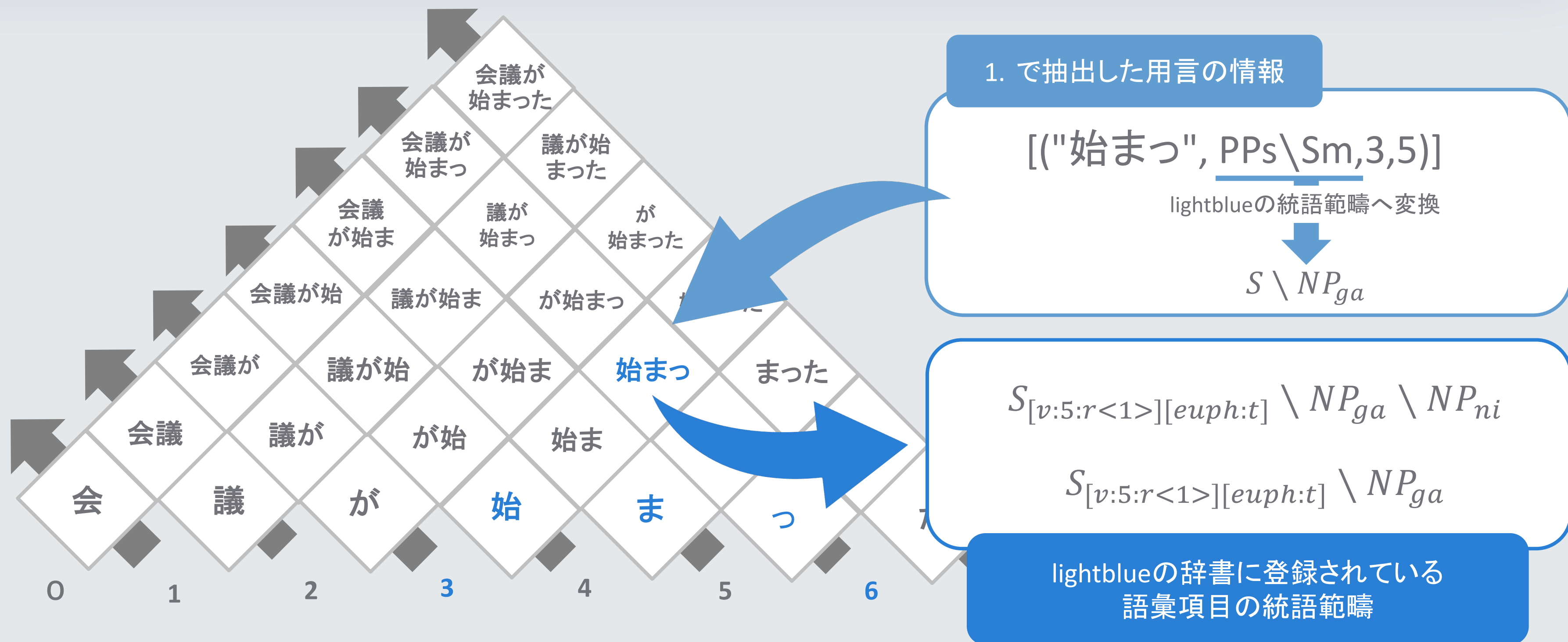
## リフォーミング – lightblueの辞書の書き換え

7/9



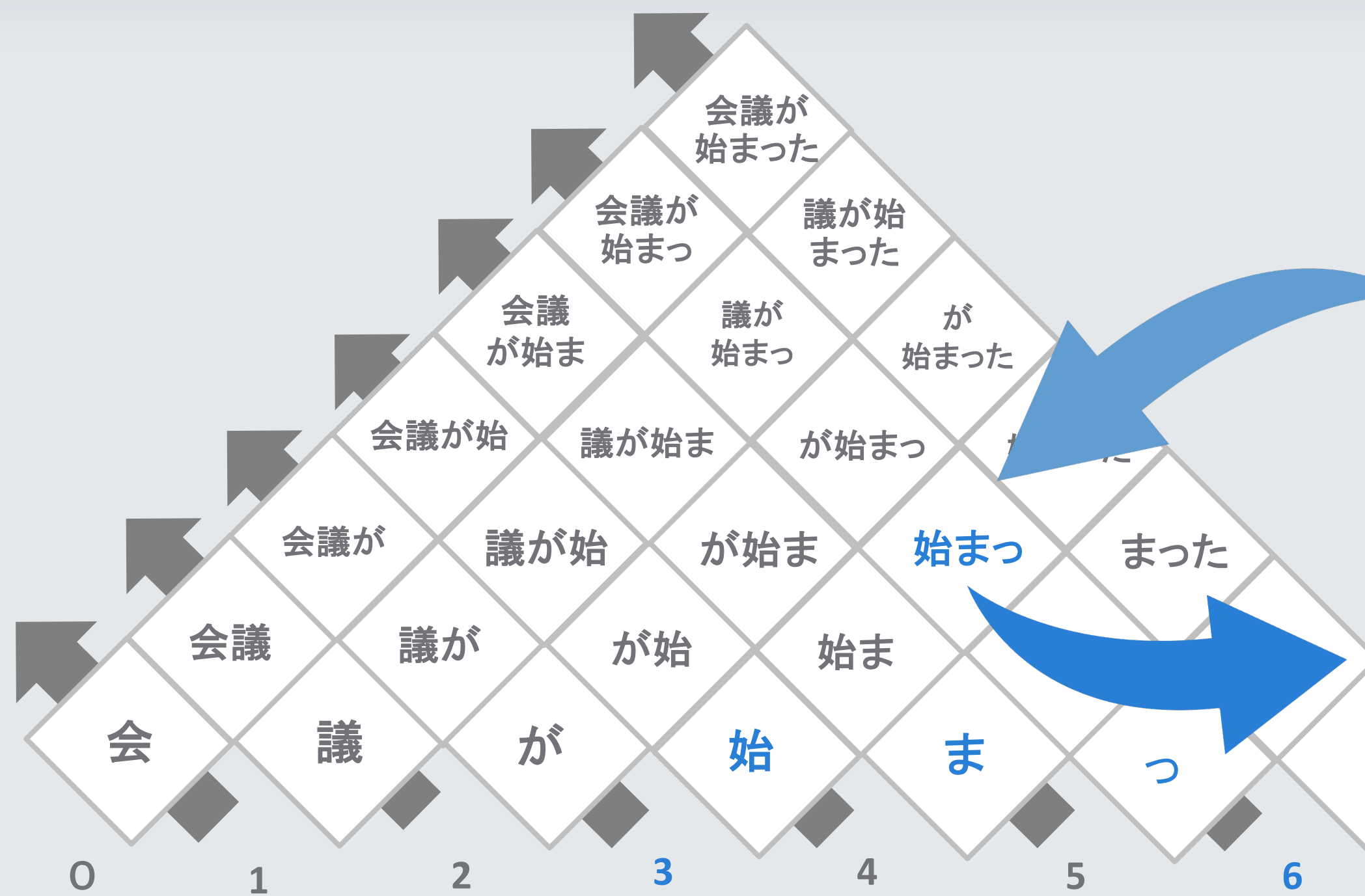
## リフォーミング – lightblueの辞書の書き換え

8/9



## リフォーミング – lightblueの辞書の書き換え

9/9



## 1. で抽出した用言の情報

$$[("始まっ", \underline{PPs \setminus Sm, 3, 5})]$$

lightblueの統語範疇へ変換

$$S \setminus NP_{ga}$$

$$\underline{S_{[v:5:r<1>][euph:t]} \setminus NP_{ga} \setminus NP_{ni}}$$

$$S_{[v:5:r<1>][euph:t]} \setminus NP_{ga}$$
lightblueの辞書に登録されている  
語彙項目の統語範疇

1. ABCツリーバンクのパーズ

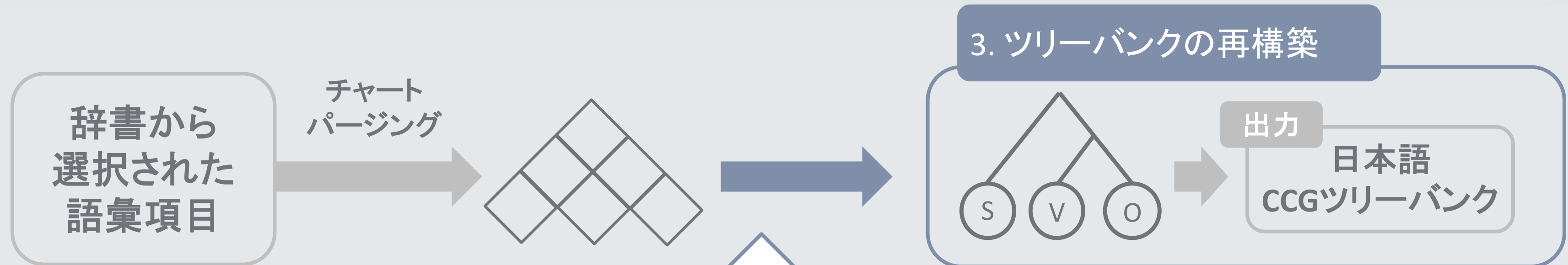
2. lightblueの辞書の書き換え

3. ツリーバンクの再構築

## リフォーミング - ツリーバンクの再構築

2/3

## 3. ツリーバンクの再構築



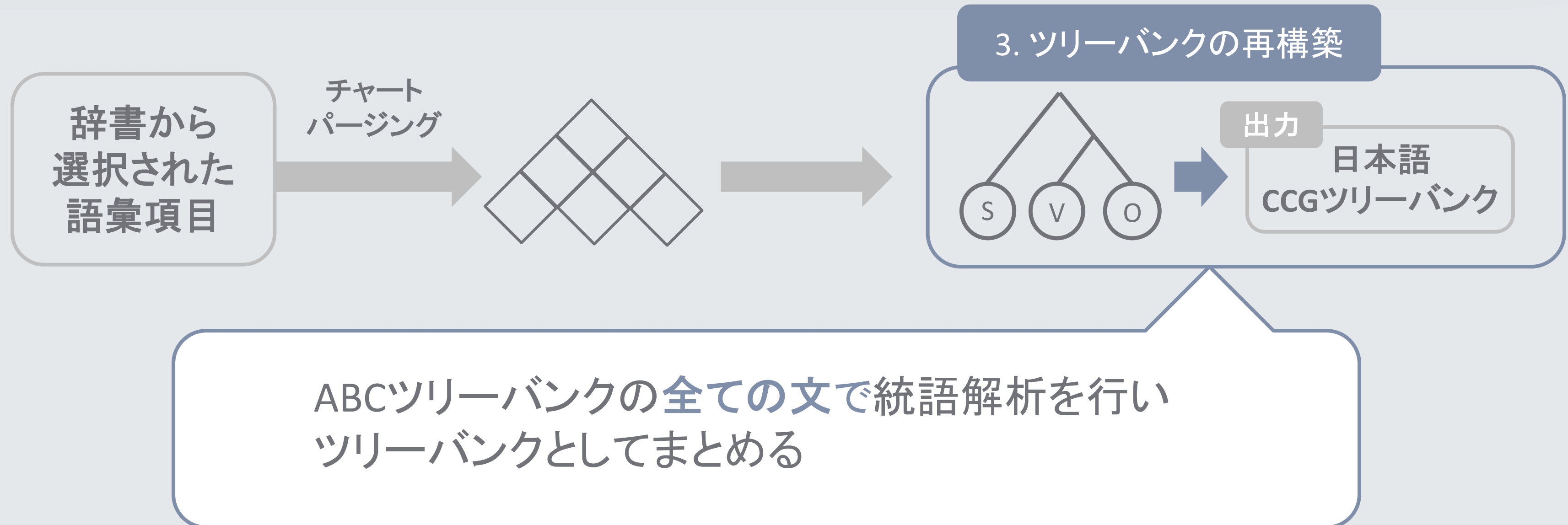
フィルターされたチャートを用いてABCツリーバンク  
の文の統語解析を行う

成功 → CCG統語構造が出力  
失敗 → エラーが出力



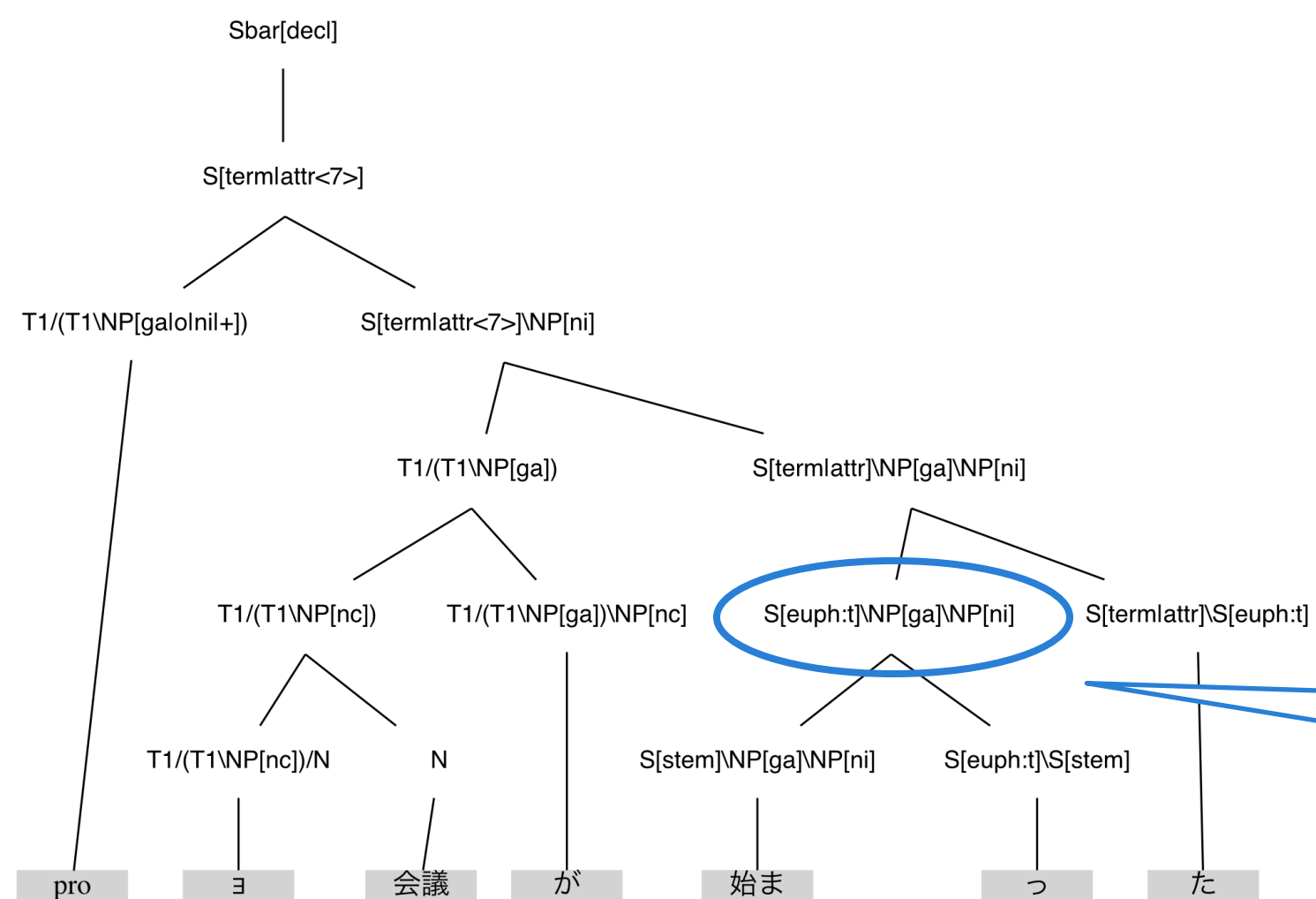
## リフォーミング - ツリーバンクの再構築

3/3



# リフォーミングの成功例

## フィルタリング前



lightblueの辞書に登録されている  
語彙項目の統語範疇

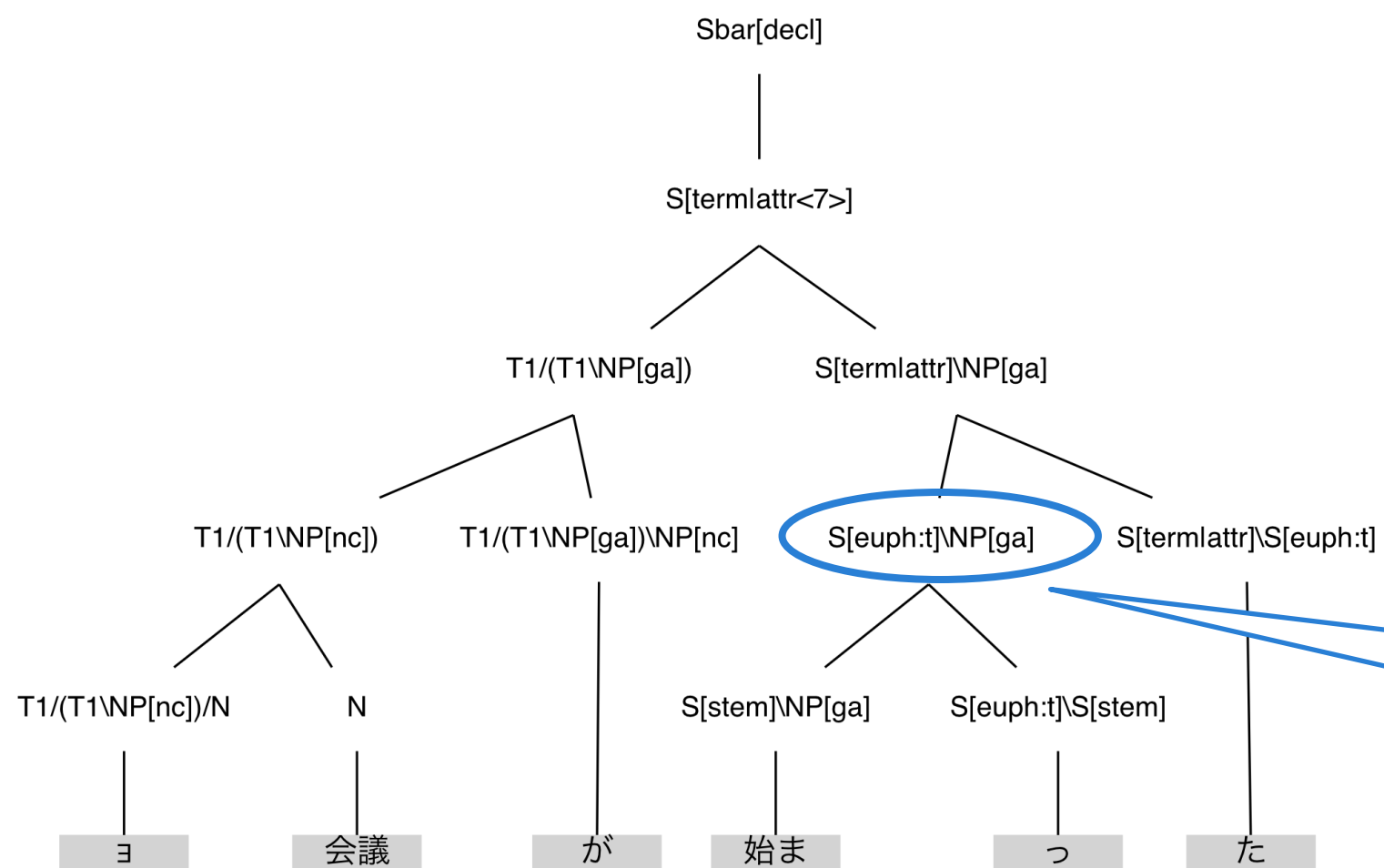
$$S_{[v:5:r<1>][euph:t]} \setminus NP_{ga} \setminus NP_{ni}$$

$$S_{[v:5:r<1>][euph:t]} \setminus NP_{ga}$$

必須格はガ格のみだが、二格も必須格としている

# リフォーミングの成功例

フィルタリング後



lightblueの辞書に登録されている  
語彙項目の統語範疇

~~$S_{[v:5:r<1>][euph:t]} \setminus NP_{ga} \setminus NP_{nt}$~~

$S_{[v:5:r<1>][euph:t]} \setminus NP_{ga}$

必須格が**ガ格**になっている

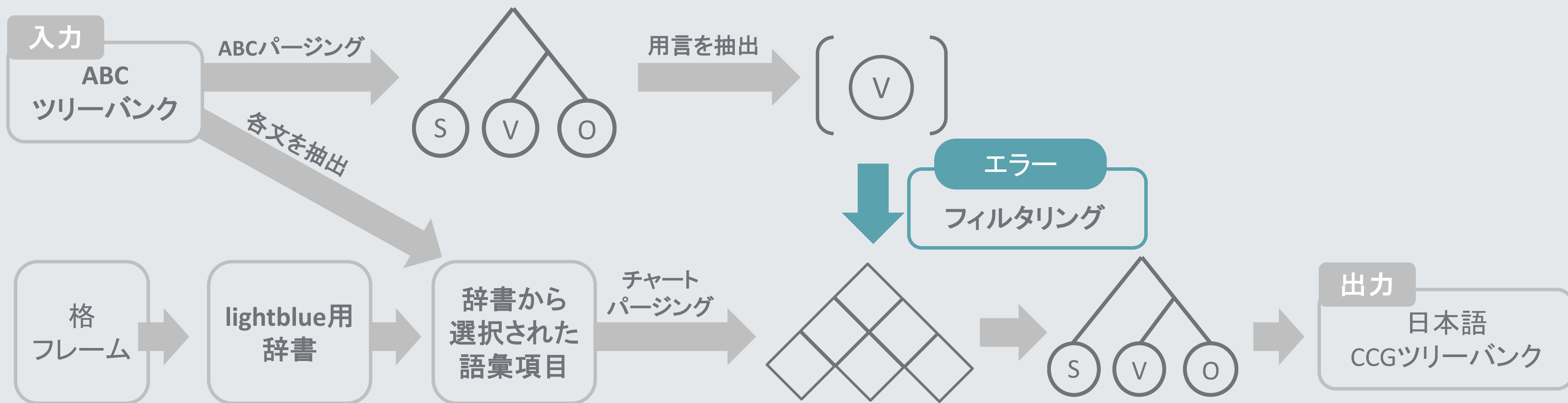
04

# エラー分析

# エラー分析

## 課題

lightblueの用言の語彙項目を ABC ツリーバンクから得られた用言の情報によってフィルタリングする際エラーが発生



# エラー分析

## 課題

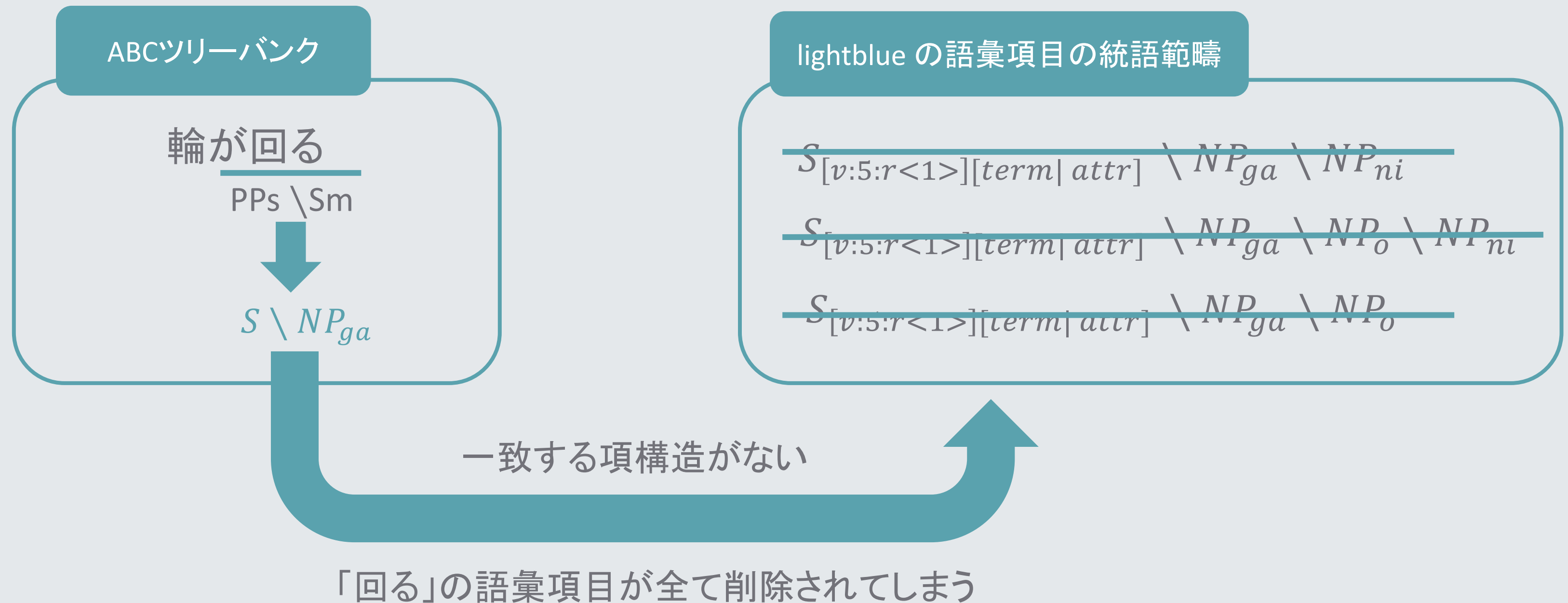
lightblueの用言の語彙項目をABC ツリーバンクから得られた用言の情報によってフィルタリングする際エラーが発生

01 lightblue の語彙項目に正しい用言が含まれないケース

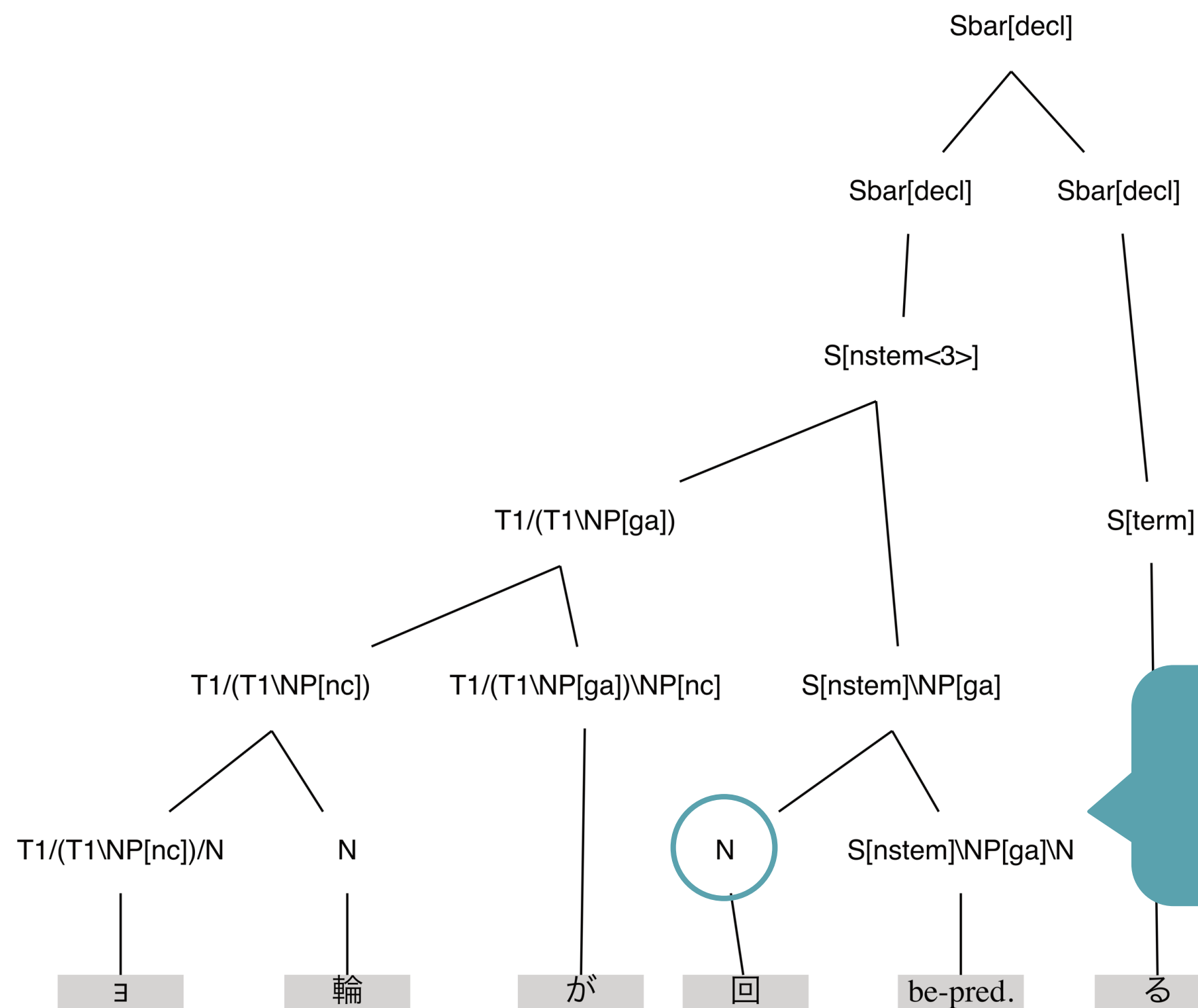
---

02 ABC ツリーバンクの項構造が誤っているケース

# エラー分析 - lightblueの語彙項目に正しい用言が含まれないケース



# エラー分析 ABCツリーバンクの項構造が誤っているケース



「回」は 名詞 として  
分析されている



# エラー分析 – ABCツリーバンクの項構造が誤っているケース

## ABCツリーバンク

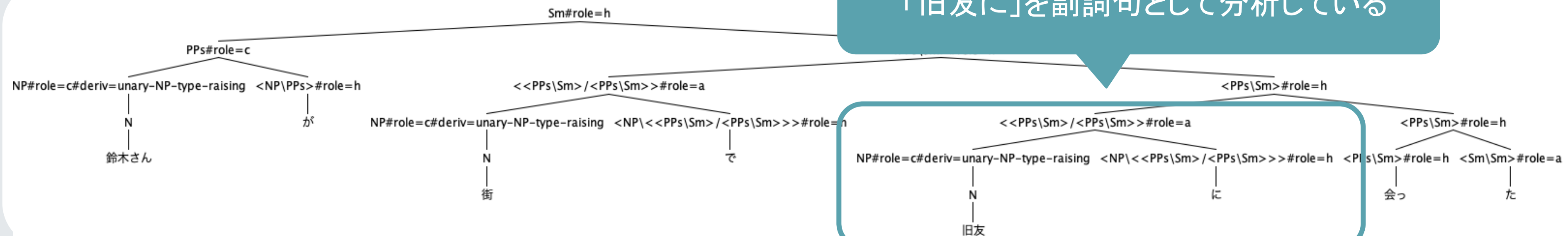
鈴木さんが街で旧友に会った。

PPs \ Sm

S \ NP<sub>ga</sub>

「会う」の必須格はガ格、二格になるべきである

「旧友に」を副詞句として分析している



# エラー分析 – ABCツリーバンクの項構造が誤っているケース

## ABCツリーバンク

鈴木さんが街で旧友に会った。

PPs \ Sm

$S \setminus NP_{ga}$

## lightblue の語彙項目の項構造

~~$S_{[v:5:w<1>][euph:t]} \setminus NP_{ga} \setminus NP_{ni} \setminus Sbar_{ga}$~~

~~$S_{[v:5:w<1>][euph:t]} \setminus NP_{ga} \setminus Sbar_{ga}$~~

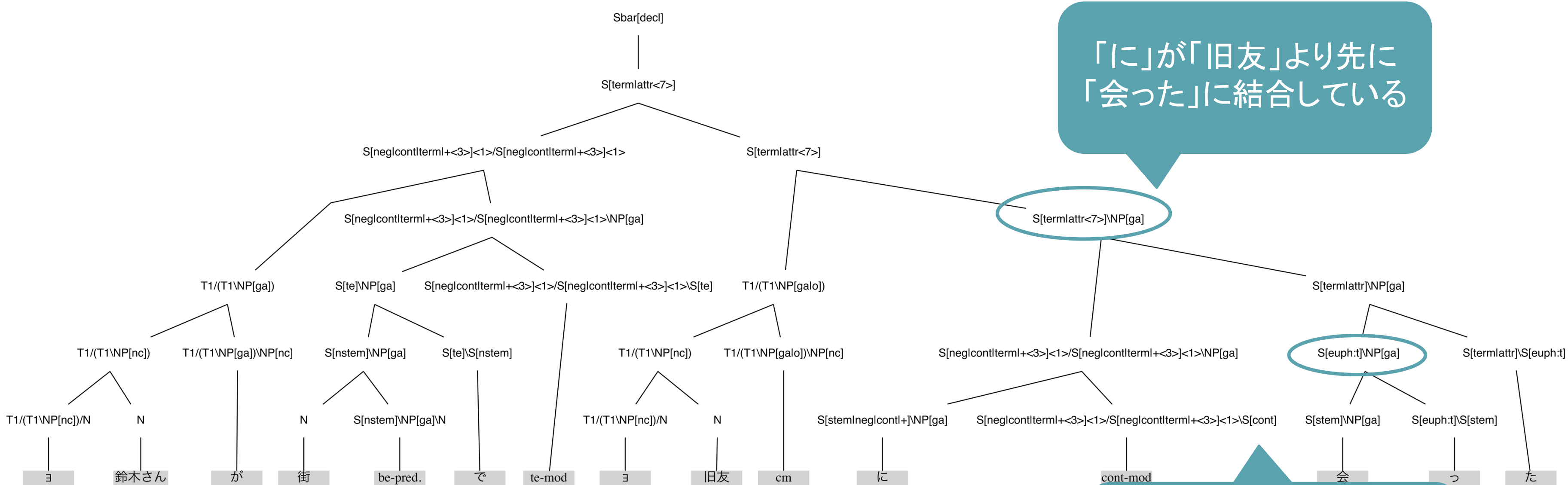
~~$S_{[v:5:w<1>][euph:t]} \setminus NP_{ga} \setminus NP_{ni}$~~

$S_{[v:5:w<1>][euph:t]} \setminus NP_{ga}$

本来は  $S \setminus NP_{ga} \setminus NP_{ni}$  が残るべき

正しい語彙項目が削除されてしまっている

# エラー分析 ABCツリーバンクの項構造が誤っているケース



05

まとめ

## まとめ

- 01 ABCツリーバンクと lightblue が持つ利点を組み合わせることで言語学的に妥当で詳細な統語情報を持った日本語CCGツリーバンクを構築する手法(リフォーミング)を提案した。
- 02 リフォーミングによって正しい日本語CCG統語構造が部分的に得られたが、誤りが含まれるケースも残されている

### 今後の展望

より正しい CCG 統語構造を出力できるようにフィルタリングのアルゴリズムを改良する

## 参考文献

1. Daisuke Bekki and Ai Kawazoe. Implementing variable vectors in a CCG parser. In Logical Aspects of Computational Linguistics. Celebrating 20 Years of LACL (1996–2016), pp. 52–67, Berlin, Heidelberg, 12 2016. Springer Berlin Heidelberg.
2. Daisuke Bekki and Hitomi Yanaka. Is Japanese CCGBank empirically correct? A case study of passive and causative constructions. In Proceedings of Treebanks and Linguistic Theories (TLT) 2023 (to appear), the workshop in the Georgetown University Round Table on Linguistics 2023 (GURT2023), forthcoming.
3. Daisuke Kawahara and Sadao Kurohashi. A fullylexicalized probabilistic model for Japanese syntactic and case structure analysis. In Proc. of the Human Language Technology Conference of the NAACL, Main Conference, June 2006.
4. Yoshikawa Masashi, Noji Hiroshi, and Matsumoto Yuji. A\* CCG parsing with a supertag and dependency factored model. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, pp. 277–287, Vancouver, Canada, 2017. Association for Computational Linguistics.
5. Mark Steedman. The Syntactic Process. MIT Press, 2000.
6. Mark J. Steedman. Surface Structure and Interpretation. The MIT Press, Cambridge, 1996.
7. 植松すみれ, 松崎拓也, 花岡洋輝, 宮尾祐介, 美馬秀樹. 統語・意味コーパスの統合と再解釈による大規模な日本語CCG文法の開発. 人工知能学会全国大会論文集, Vol. JSAI2013, pp. 4B11–4B11, 2013.
8. 戸次大介. 日本語文法の形式理論. くろしお出版, 東京, 2010.
9. 窪田悠介, 峯島宏次, 林則序, 岡野伸哉. 汎用的な範疇文法 ツリーバンクの構築. 言語処理学会 第 25 回年次大会 発表論文集 (2019 年 3 月), pp. 143–146. 一般社団法人 言語処理学会, 2019.
10. 窪田悠介, 峯島宏次, 林則序, 岡野伸哉. ABC ツリーバンク: 学際的な言語研究のための基盤資源. 言語処理学会 第 27 回年次大会 発表論文集 (2021 年 3 月), pp. 1529–1534. 一般社団法人 言語処理学会, 2021.
11. 花岡洋輝, 増田勝也, 植松すみれ, 美馬秀樹. 日本語助詞「と」コーパスの構築. 言語処理学会 第 18 回年次大会 発表論文集 (2012 年 3 月), pp. 247–250. 一般社団法人 言語処理学会, 2012.

## 補足 – 語彙項目、統語蘇生

### 語彙項目

辞書に記述される語と統語範疇の対応付け

例: Keats       $\vdash NP$   
      eats  $\vdash (S \setminus NP) / NP$   
      apples     $\vdash NP$

### 統語素性

統語範疇についてのより詳しい指定

例: eats  $\vdash (S \setminus NP_{nom}) / NP_{acc}$

## 補足 – 格フレーム

### 格フレーム

私がカメラで走っている少女を見た。

計算機には

{  
カメラで見た  
カメラで走っている

のどちらが妥当かを判定するのが難しい



## 補足 – 格フレーム

### 格フレーム

用言とそれに関係する名詞を用言の各用法ごとに整理したもの

例:「とまる」

1. 「停止する」の意味の「とまる」

{ 車, バス, タクシー, ... } が { 駐車場, 横, 路肩, ... } に とまる

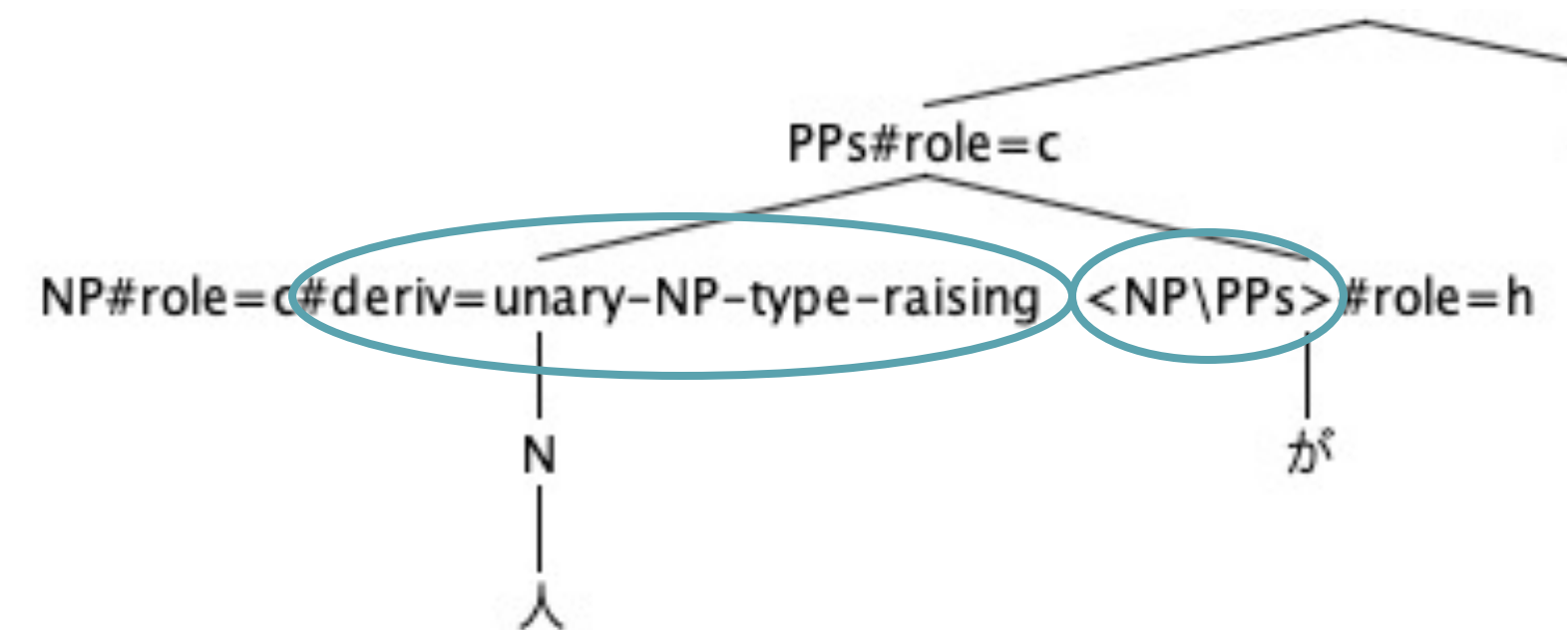
2. 「吊るす」の意味の「かける」

{ 人, 3人, ... } が { ホテル, 家, ... } に { 3人, 家族, ... } で とまる

# 補足 – ABC文法の代表的な規則

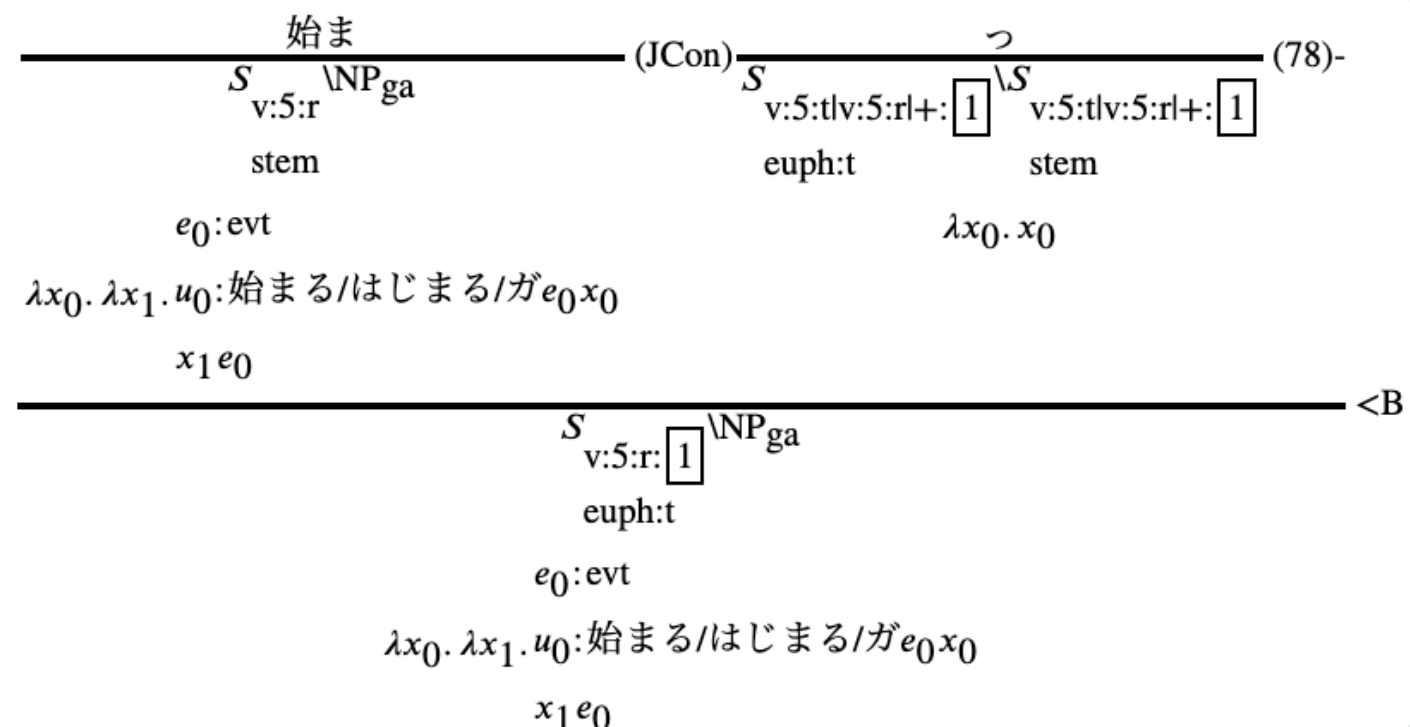
## ABC文法の代表的な規則

名称	規則	
関数適用(順方向)	$A/B \ B \Rightarrow A$	ABC 文法
関数適用(逆方向)	$B \ B \backslash A \Rightarrow A$	
関数合成(順方向)	$A/B \ (B/C)/\$ \Rightarrow (A/C)/\$$	
関数合成(逆方向)	$\ (C \ B)C \ A \Rightarrow \ (C \ A)$	
名詞の名詞句化	$N \Rightarrow NP$	unary 規則
名詞句の述語化	$NP \Rightarrow PPs \backslash S$	
名詞句の後置詞句化	$NP \Rightarrow PPs; \ NP \Rightarrow PPo1; \text{ etc.}$	
関係節による名詞修飾	$PPs \backslash Srel \Rightarrow N/N; \ PPo1 \backslash Srel \Rightarrow N/N; \ Srel \Rightarrow N/N; \text{ etc.}$	
連用節化	$Sa \Rightarrow S/S$	
名詞の副詞化	$NP \Rightarrow S/S; \ NP \Rightarrow (PPs \backslash S)/(PPs \backslash S); \text{ etc.}$	
スクランブリング	$PPo1 \backslash (PPs \backslash S) \Rightarrow PPs \backslash (PPo1 \backslash S); \text{ etc.}$	
疑問節化・感嘆節化	$Ssub \Rightarrow CPq; \ Ssub \Rightarrow CPx$	



# 補足 – lightblueの辞書の形式

## lightblueの辞書の形式



項目	用言「始まっ」の辞書	備考
規則名	BFC1	関数合成(逆方向)
表層形	始まっ	
CCG統語範疇	$S[v:5:r<1>][\text{euph:t}][\backslash \text{NP}[ga]$	ラ行五段活用動詞、 必須格ガ格、
DTS意味表示	$\lambda x_0. \lambda x_1. (e_0: \text{evt}) \times (u_0: \text{始まる/はじまる/ガ}(e_0, x_0)) \times x_1(e_0)$	
シグネチャー	$[["\text{始まる/はじまる/ガ}", (x_0: \text{entity}) \rightarrow (e_0: \text{evt}) \rightarrow \text{type}]]$	
子ノード	$[LEX \text{ 始ま } S[v:5:r][\text{stem}][\backslash \text{NP}[ga]$ $\lambda x_0. \lambda x_1. (e_0: \text{evt}) \times (u_0: \text{始まる/はじまる/ガ}$ $(e_0, x_0)) \times x_1(e_0) \text{ (JCon) [0.99] Sig. [始まる/はじま}$ $\text{る/ガ}: (x_0: \text{entity}) \rightarrow (e_0: \text{evt}) \rightarrow \text{type}]$ $, LEX \text{ っ } S[v:5:t v:5:r v:5:w +<1>][\text{euph:t}][]$ $\backslash S[v:5:t v:5:r v:5:w +<1>][\text{stem}][] \lambda x_0. x_0 \text{ (78)-}$ $[1.00] \text{ Sig. []}$	「始ま」と「っ」の辞書
スコア	99 / 100	0.99
ソース	""	