

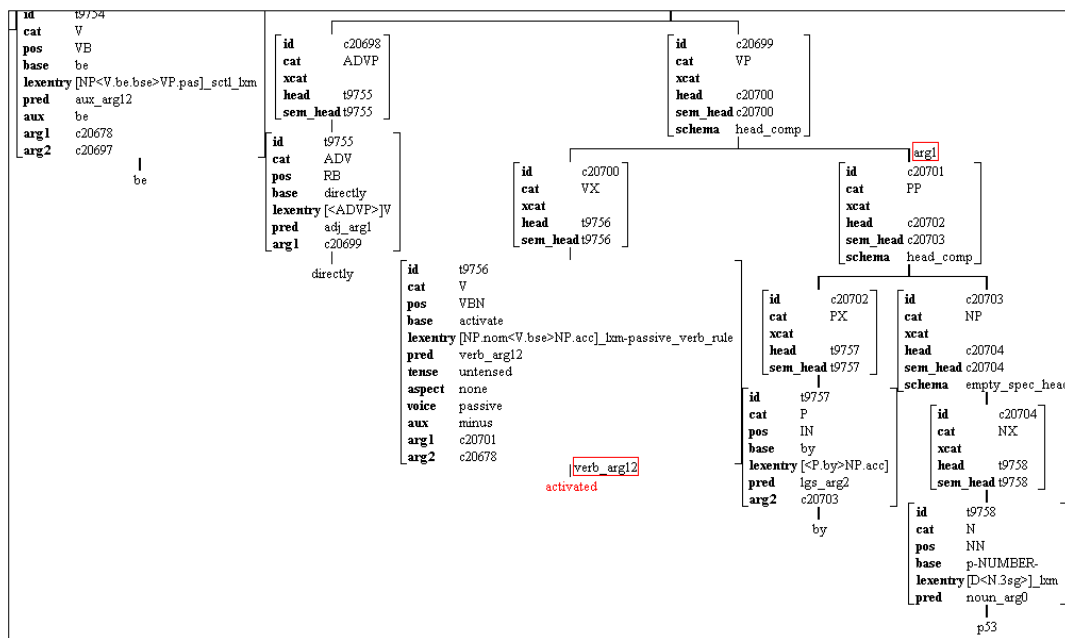
自然言語処理における 構文解析と言語理論の関係

宮尾祐介

国立情報学研究所

研究紹介

- 言語理論に基づく構文解析
- 英語HPSGパーザ Enju
 - HPSG理論に基づき統語構造(構文木)と意味構造(述語項構造)を計算
- 中国語、日本語も



Enju の出力

- 述語項構造

The Bcl-2 protein has been isolated and shown to be directly activated by p53.



- 述語論理式

- Stately, plump Buck Mulligan came from the stairhead, bearing a bowl of lather on which a mirror and a razor lay across. A yellow dressinggown ungirdled was sustained gently behind him on the mild morning air.
- $\exists v x_1 x_2 x_3$ (plump_buck_mulligan(z) \wedge lather(x) \wedge bowl+of(x_2, x) \wedge mirror(v) \wedge lay_across+on(v, x_2) \wedge razor(x_1) \wedge lay_across+on(x_1, x_2) \wedge bear_{prog}(z, x_2) \wedge stairhead(y) \wedge stately_come_{past}+from(z, y) \wedge yellow_dressinggown_ungirdled (x_3) \wedge mild_morning_air(u) \wedge $\exists x_4$ ($x_4 = z$ \wedge sustain_{past,passive}—gently+behind+on (x_3, x_4, u)))

自然言語の構文解析

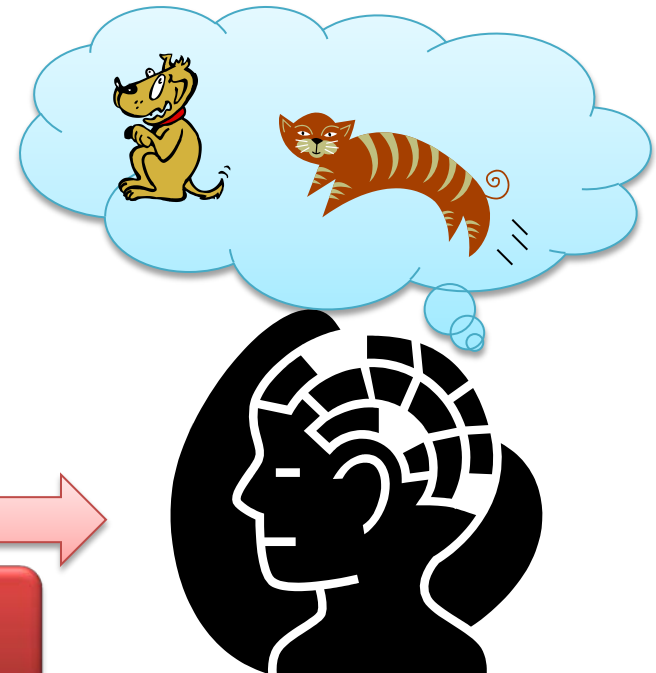
- テキストを入力として、その意味表現を計算する

入力: テキスト

A cat chased a dog.

構文解析

出力: 意味表現



アプローチ

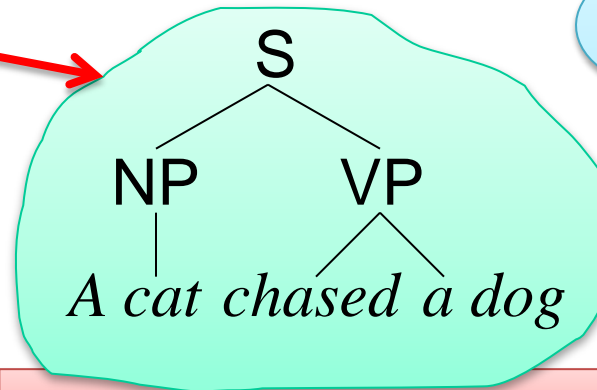
- 統語構造: 文と意味表現の橋渡し
- 文法: 構文構造を導出する規則
→ 文法を実装すれば、構文解析ができる!

出力: 意味表現

文法(←言語理論)

S → NP VP
VP → V NP
NP → NP PP
NP → D N
...

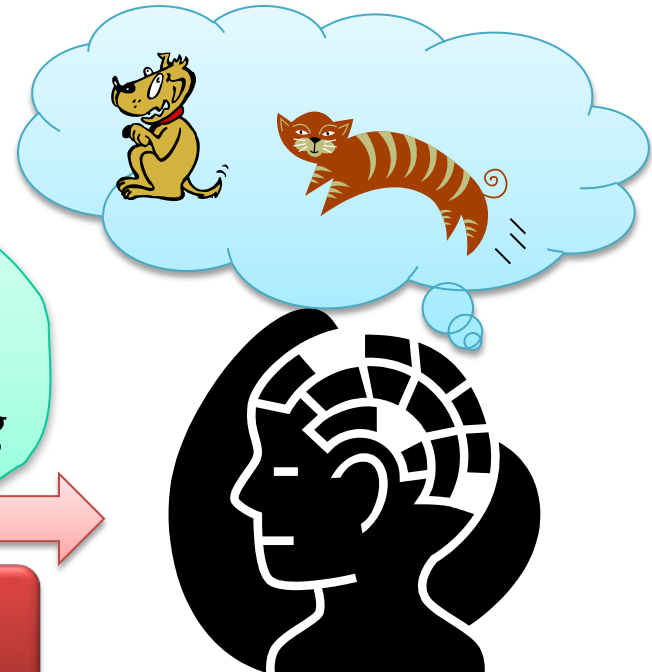
統語構造



入力: テキスト

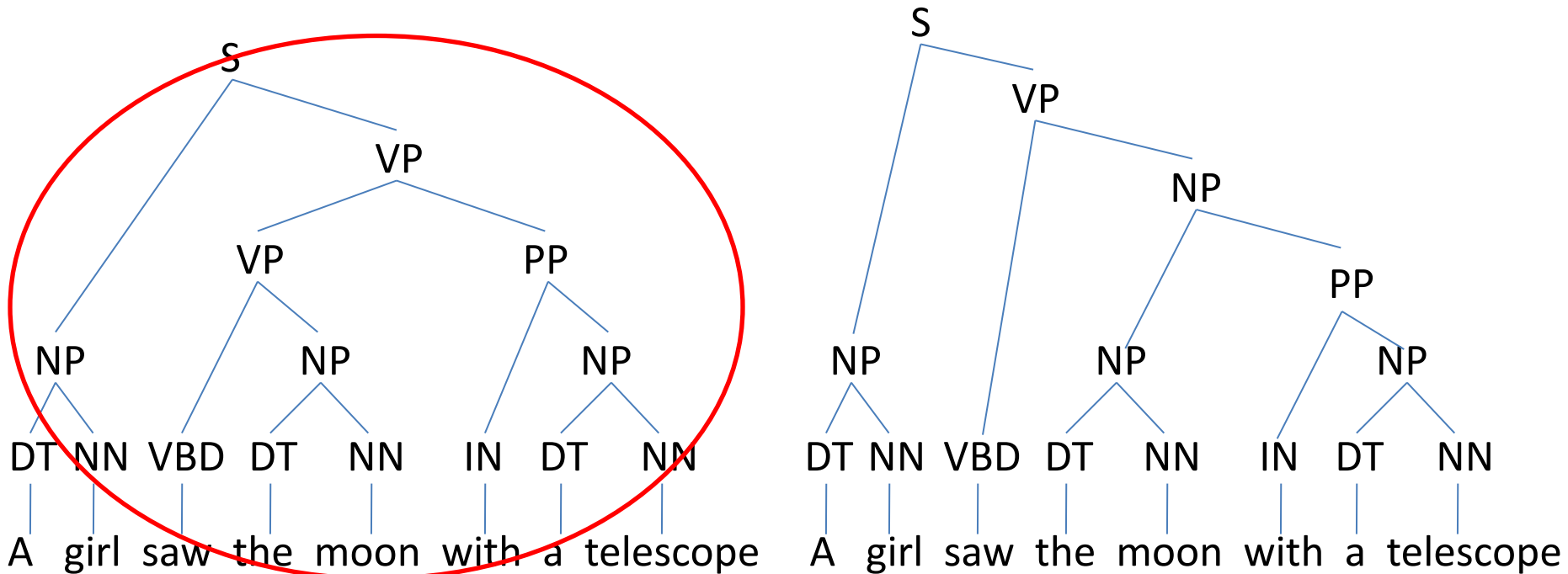
A cat chased a dog.

構文解析



構文解析と文法

- 文法を実装すれば、構文解析ができる → **間違い**
- 文法は**文法的に適切な構造**を規定する
- 文法的に適切な構造の集合 ≫ 人間の解釈
- 構文解析では、**人間の解釈に相当する構造**が欲しい



構文解析と曖昧性解消

- 曖昧性解消がない構文解析は意味がない
 - 構文解析 = 探索問題
 - 高い構文解析精度 = よい曖昧性解消
- 自然言語処理における構文解析研究のほとんどは曖昧性解消(≒統計モデル)について
 - 逆に、文法の役割は忘れられている
- 今日は、構文解析において文法や言語理論が果たす役割について議論

範疇文法(CG)

- カテゴリを組み合わせることで統語構造を計算
- 基本カテゴリ: S, NP, N
- 複合カテゴリ
 - カテゴリを “/” や “\” でつなげたもの
 - X/Y は、カテゴリ Y が右側に来るとくっついて X になる
 - $X\backslash Y$ は、カテゴリ Y が左側に来るとくっついて X になる
- 例:
 - 自動詞: $S\backslash NP$
 - 他動詞: $S\backslash NP/NP$

NP	S\NP	→	S
<i>John</i>	<i>walked</i>		

規則

- 関数適用規則

– $X/Y \quad Y \Rightarrow X \quad (>)$

– $Y \quad X \backslash Y \Rightarrow X \quad (<)$

- 例:

– $S/NP \quad NP \Rightarrow S$

– $NP \quad S \backslash NP \Rightarrow S$

NP	S \ NP	→	S
<i>John</i>	<i>walked</i>		

構文解析

- 各単語にカテゴリ(語彙カテゴリ)を割り当て

NP

John

S\NP/NP

loves

NP

Mary

構文解析

- 規則を適用し、カテゴリをくっつける

$$X/Y \quad Y \Rightarrow X \quad (>)$$

$$X = S \setminus NP, Y = NP$$

NP

John

S \ NP / NP

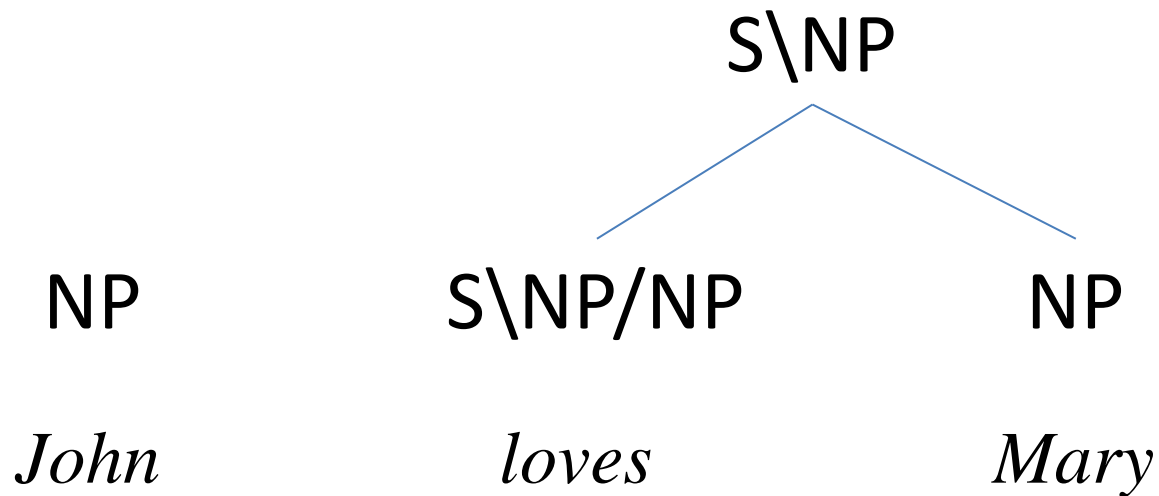
loves

NP

Mary

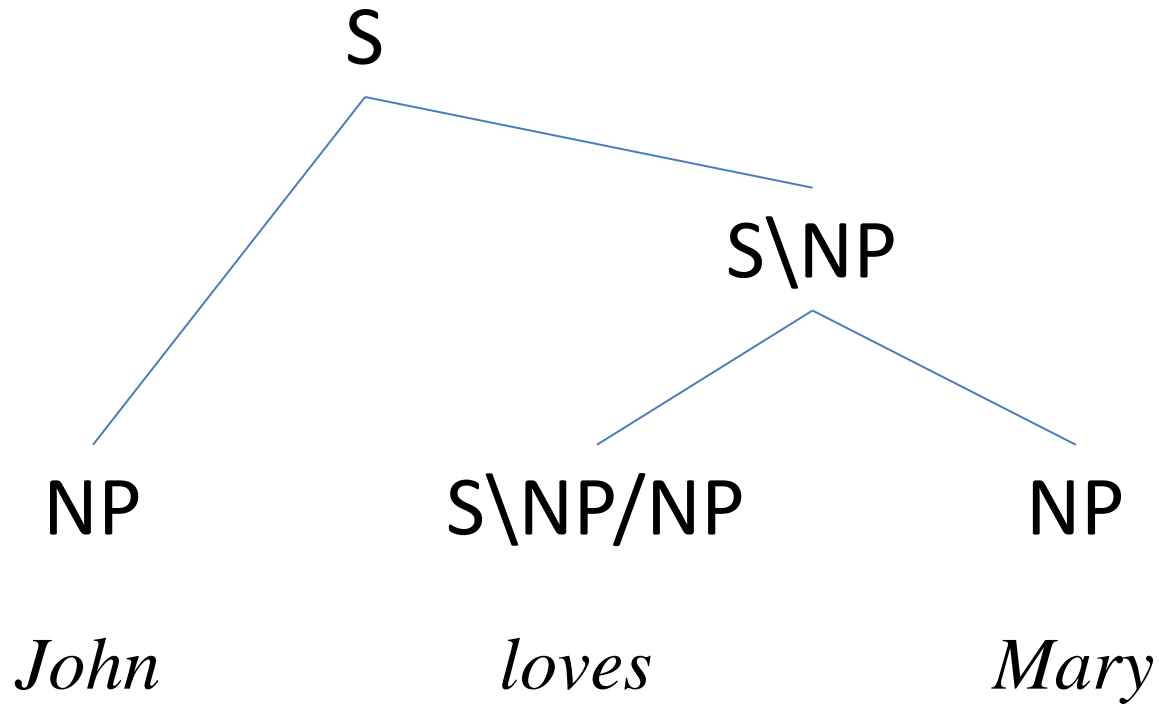
構文解析

- 規則を適用し、カテゴリをくっつける

$$Y \quad X \backslash Y \quad \Rightarrow \quad X \quad (<)$$
$$X = S, Y = NP$$


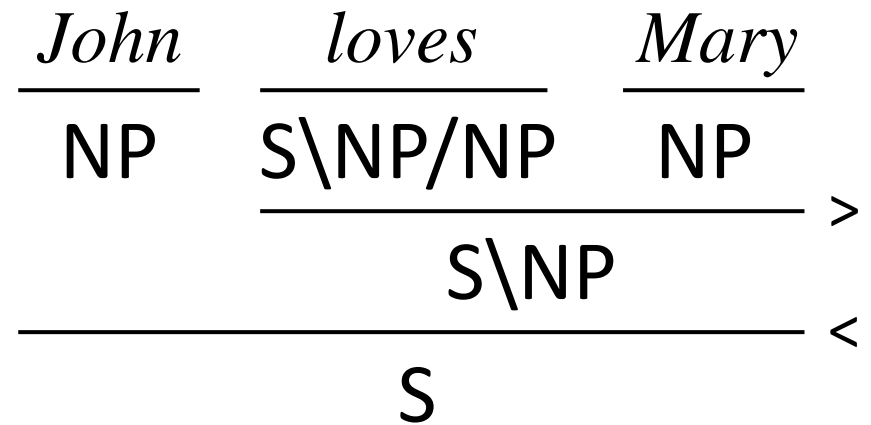
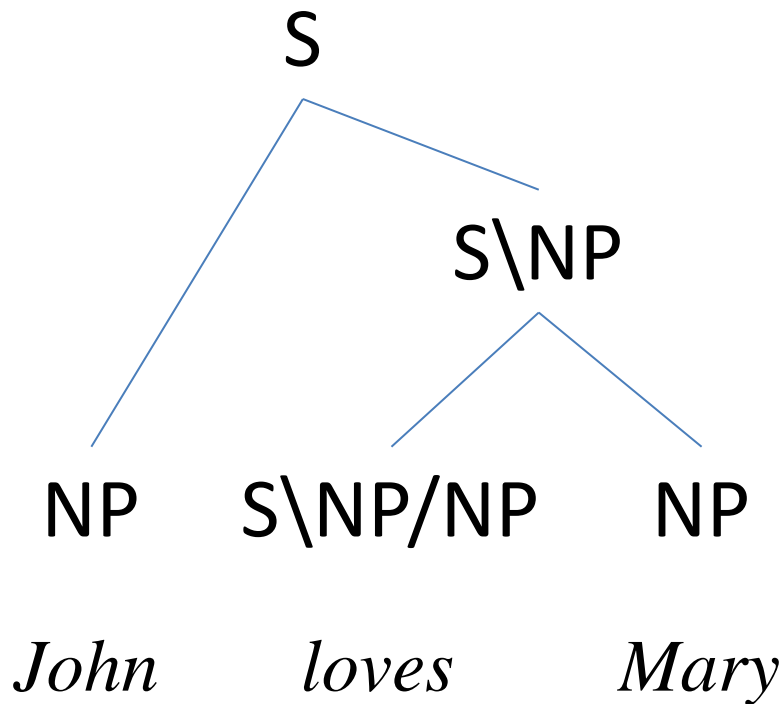
構文解析

- 規則を適用し、カテゴリをくっつける



構文解析 = 証明

- 範疇文法分野では、伝統的に構文解析を証明として表す



意味構造の計算

- 各カテゴリには、意味構造の λ 式が割り当てられる
- 統語構造に沿って λ 計算を行うと、意味構造が計算できる

<u>John</u>	<u>loves</u>	<u>Mary</u>
NP: john	$S \backslash NP / NP: \lambda x. \lambda y. \text{love}(y, x)$	NP: mary
	$S \backslash NP: \lambda y. \text{love}(y, \text{mary})$	>
	S: love(john, mary)	<

組合せ範疇文法(CCG)

- 範疇文法に数個の規則を追加することで、自然言語の様々な構文がすっきり解析できる

CG

カテゴリ
関数適用規則

CCG

+ 拡張規則
(組合せ規則)

CCG の例

<i>John</i>	<i>loved</i>	<i>and</i>	<i>Bob</i>	<i>hated</i>	<i>Mary</i>
NP	S\NP/NP	CONJ	NP	S\NP/NP	NP
S/(S\NP)			S/(S\NP)		
S/NP			S/NP		
S/NP					
S					

- 関数合成規則

- $X/Y \quad Y/Z \Rightarrow X/Z \quad (>B)$
- $Y\Z \quad X\Y \Rightarrow X\Z \quad (<B)$

- 型繰上げ規則

- $X \Rightarrow T/(T\X) \quad (>T)$

主辞駆動句構造文法(HPSG)

- 文法 = 語彙項目 + 構文規則

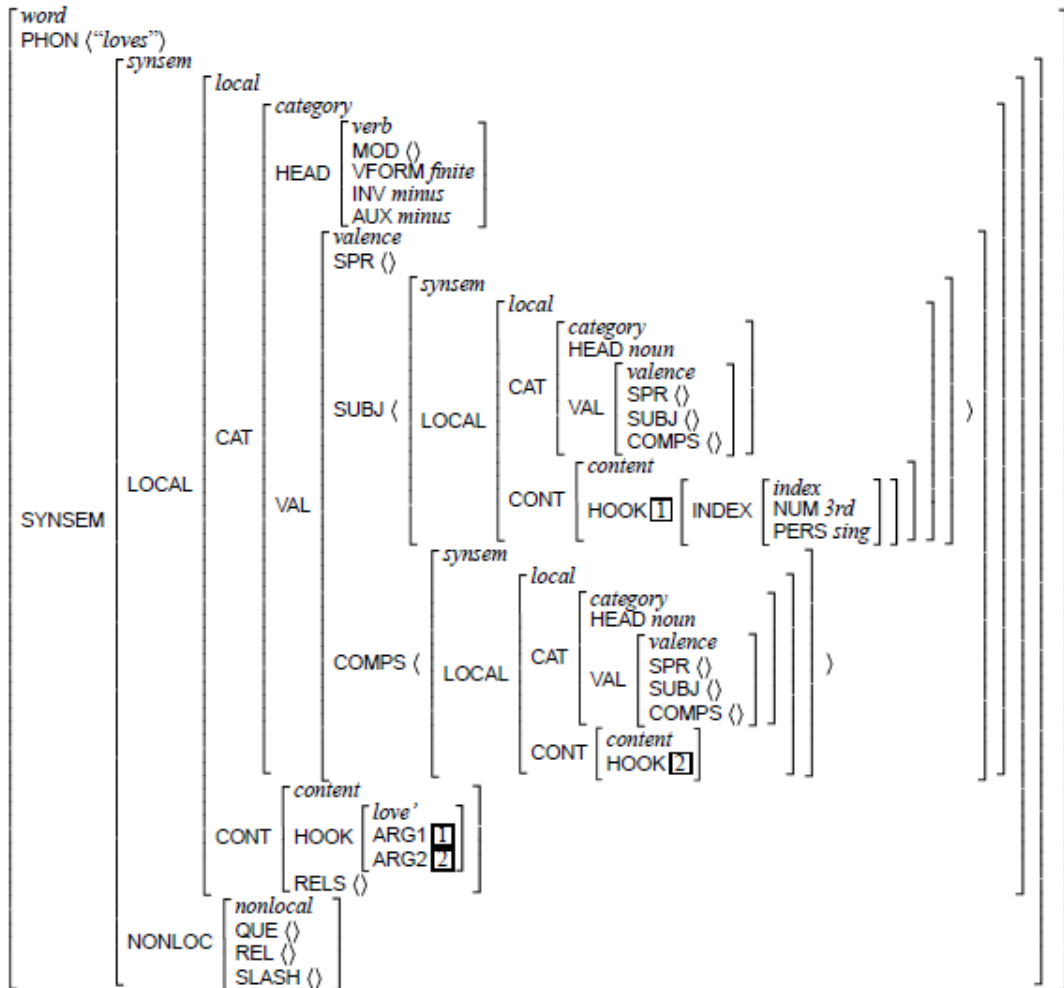
- 語彙項目 (≡ 語彙カテゴリ)

$$\left[\begin{array}{l} \text{HEAD } \textit{verb} \\ \text{SUBJ } \langle \text{NP} \rangle \\ \text{COMPS } \langle \text{NP} \rangle \end{array} \right] \quad (\equiv S \backslash \text{NP} / \text{NP})$$

- 構文規則 (≡ 組合せ規則)

$$\boxed{2} \left[\begin{array}{l} \text{HEAD } \boxed{1} \\ \text{SUBJ } \langle \boxed{2} \rangle \\ \text{COMPS } \langle \rangle \end{array} \right] \Rightarrow \left[\begin{array}{l} \text{HEAD } \boxed{1} \\ \text{SUBJ } \langle \rangle \\ \text{COMPS } \langle \rangle \end{array} \right] \quad (\equiv Y \ X \backslash Y \Rightarrow X)$$

語彙項目



- 単語の文法的性質を表す
- 単純化した記法で十分

[PHON "loves"
 HEAD verb
 SUBJ <NP₁>
 COMPS <NP₂>
 CONT love([1], [2])]

[HEAD noun
 SUBJ <>
 COMPS <>
 CONT [1]]

NP を主語に取り
 NP を目的語に取る
 動詞

HPSG 構文解析

- 各単語に語彙項目を割り当て

語彙項目

[HEAD *noun*
SUBJ <>
COMPS <>]

John

[HEAD *verb*
SUBJ <NP>
COMPS <NP>]

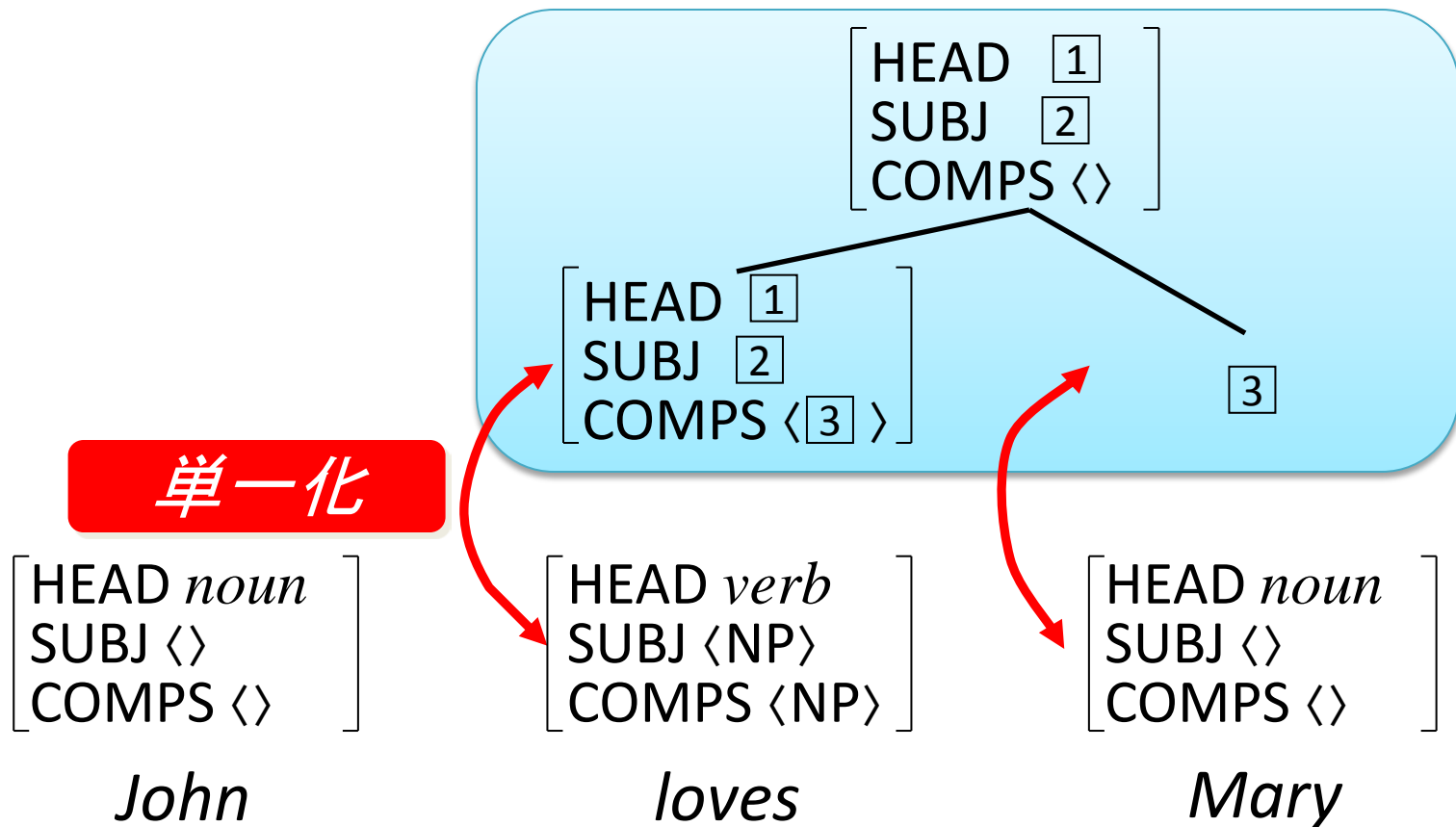
loves

[HEAD *noun*
SUBJ <>
COMPS <>]

Mary

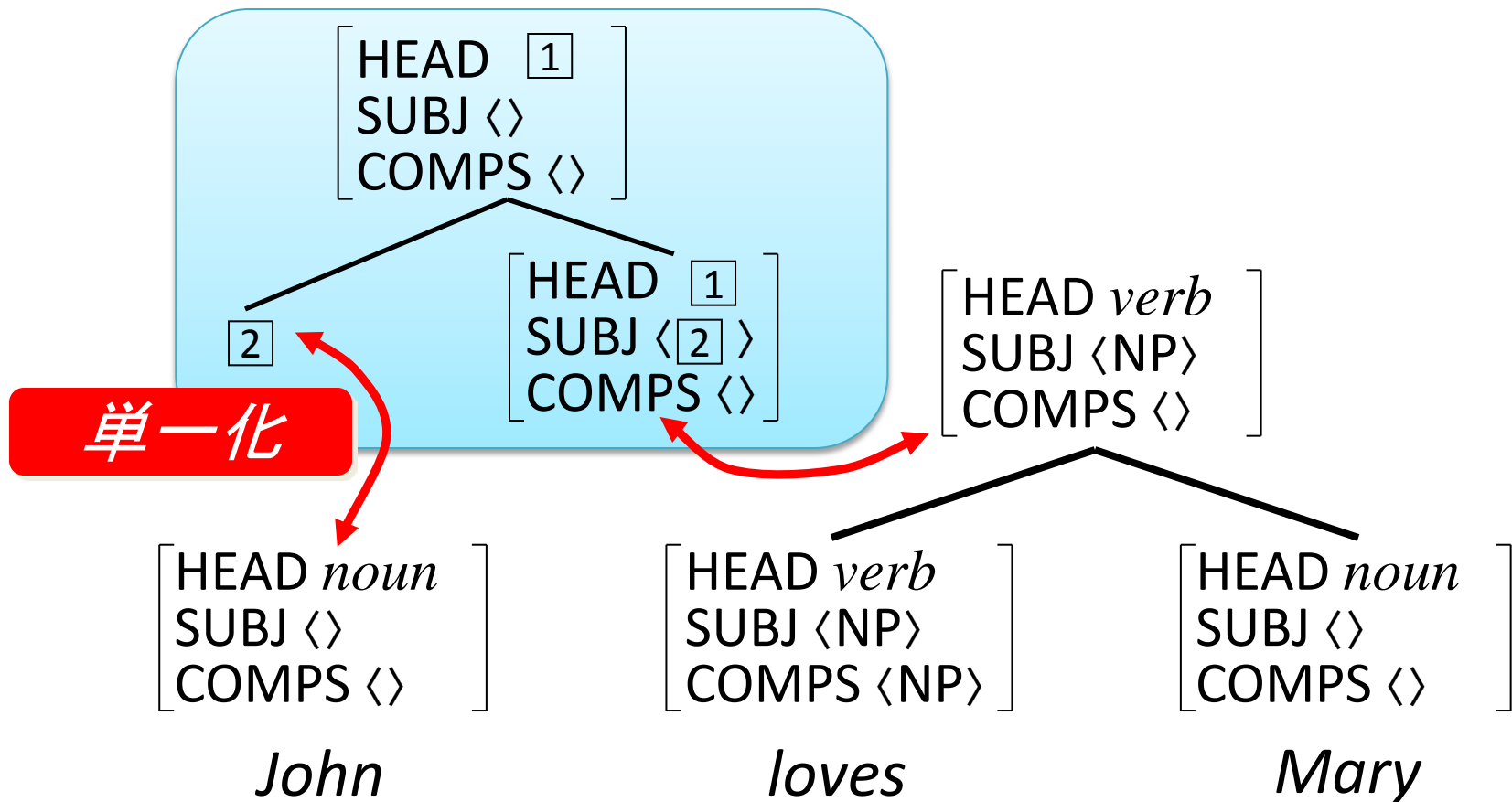
HPSG 構文解析

- 構文規則で句を作る



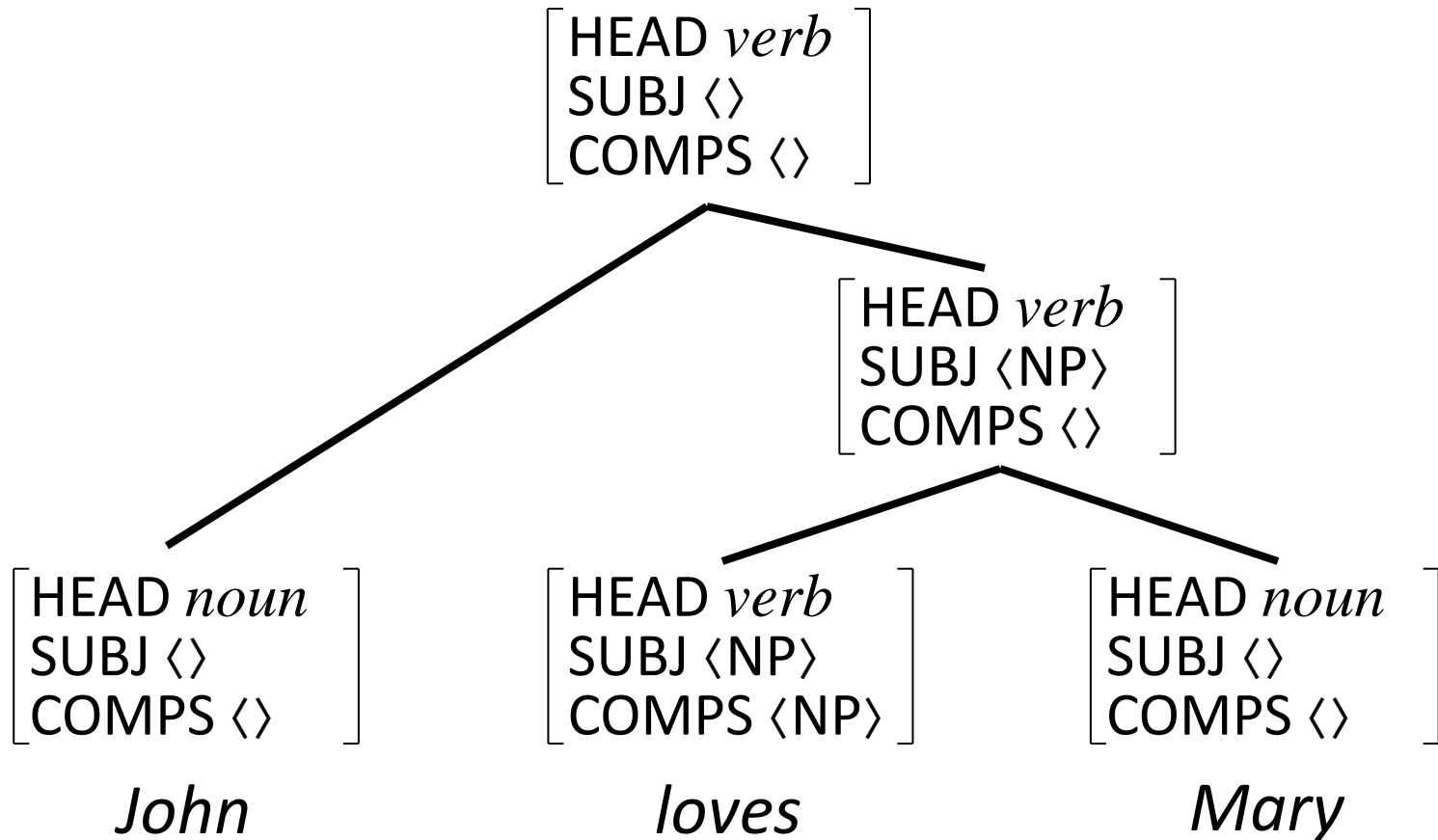
HPSG 構文解析

- 構文規則で句を作る



HPSG 構文解析

- 構文規則で句を作る

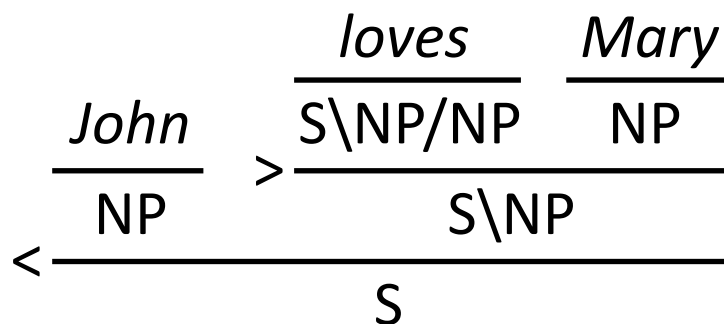
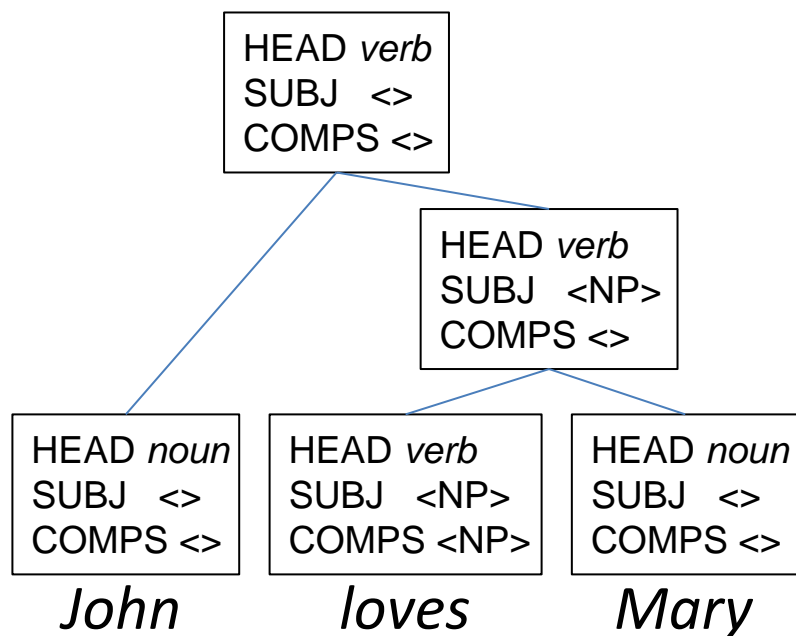


HPSG \doteq CCG

- 基本的な構文の分析はだいたい同じ

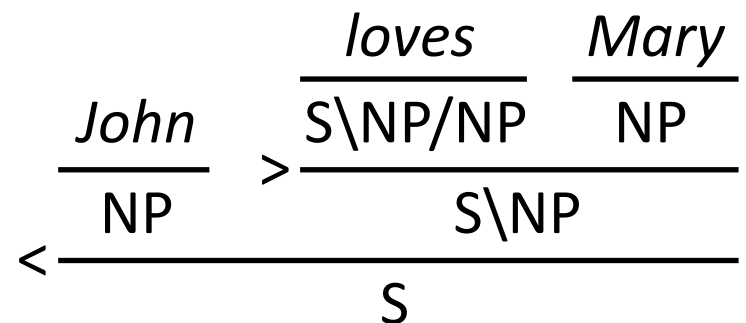
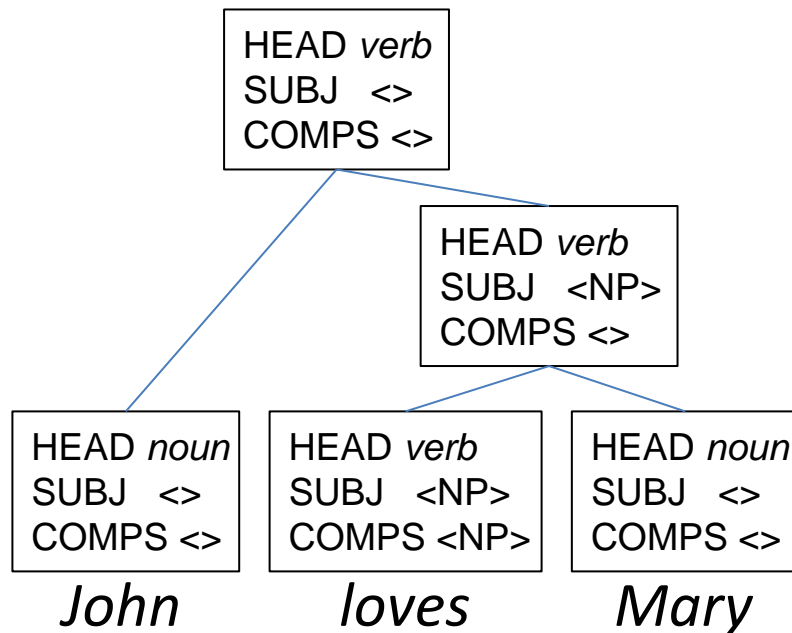
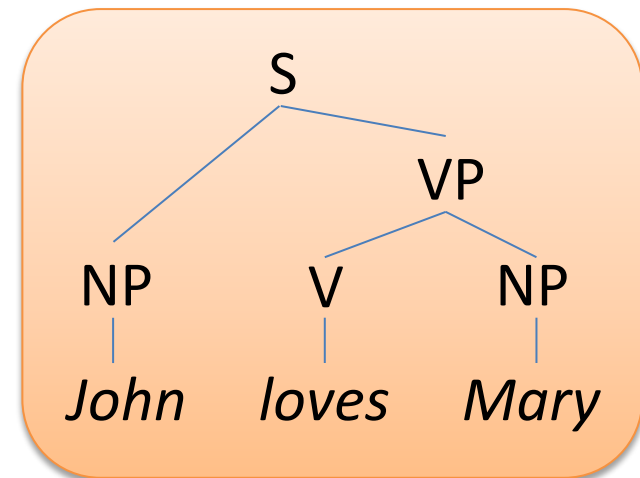
$$\left[\begin{array}{l} \text{HEAD } \textit{verb} \\ \text{SUBJ } \langle \textit{NP} \rangle \\ \text{COMPS } \langle \textit{NP} \rangle \end{array} \right] \doteq S \backslash \textit{NP} / \textit{NP}$$

– もちろん、分析のしかたが異なる言語現象も



HPSG/CCG 構文解析 ≡ CFG 構文解析

- CFG の解析アルゴリズムがそのまま利用できる
 - CYK法、確率モデル、探索手法、...
- (非)終端記号、生成規則が違う
 - 終端記号 = 語彙項目、語彙カテゴリ
 - 生成規則 = 構文規則、組合せ規則



文法と曖昧性解消

- 曖昧性解消を前提とすると、最適解を一つ求めることが目標
 - 全解探索は必要ない
 - 全解探索を前提とした理論やアルゴリズムはあまり重要でない
 - 解析アルゴリズム、計算量、学習可能性、etc.
- 曖昧性解消を前提として、文法と構文解析の関係を考えるべき

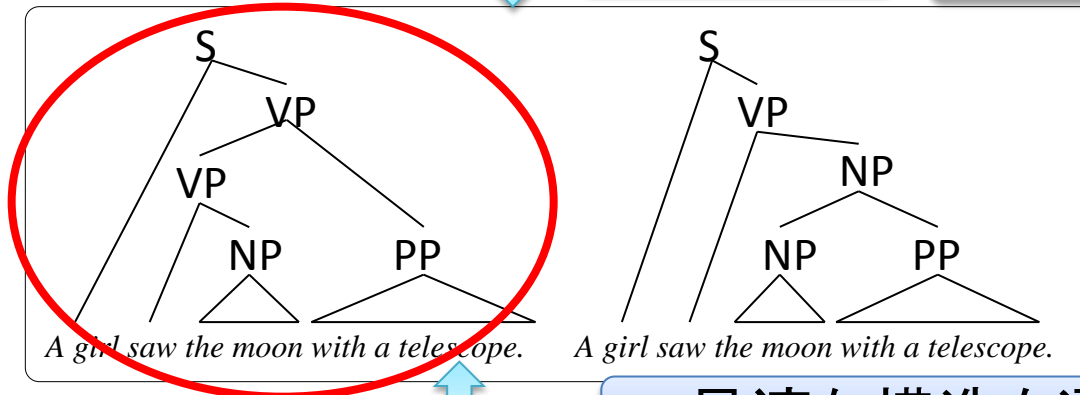
文法と曖昧性解消

- 曖昧性解消を前提としたときの文法の役割
 - 解候補の集合を決める
 - 曖昧性解消に有用な構造を提供

A girl saw the moon with a telescope

文法

文法的に適格な構造の集合



最適な構造を選択

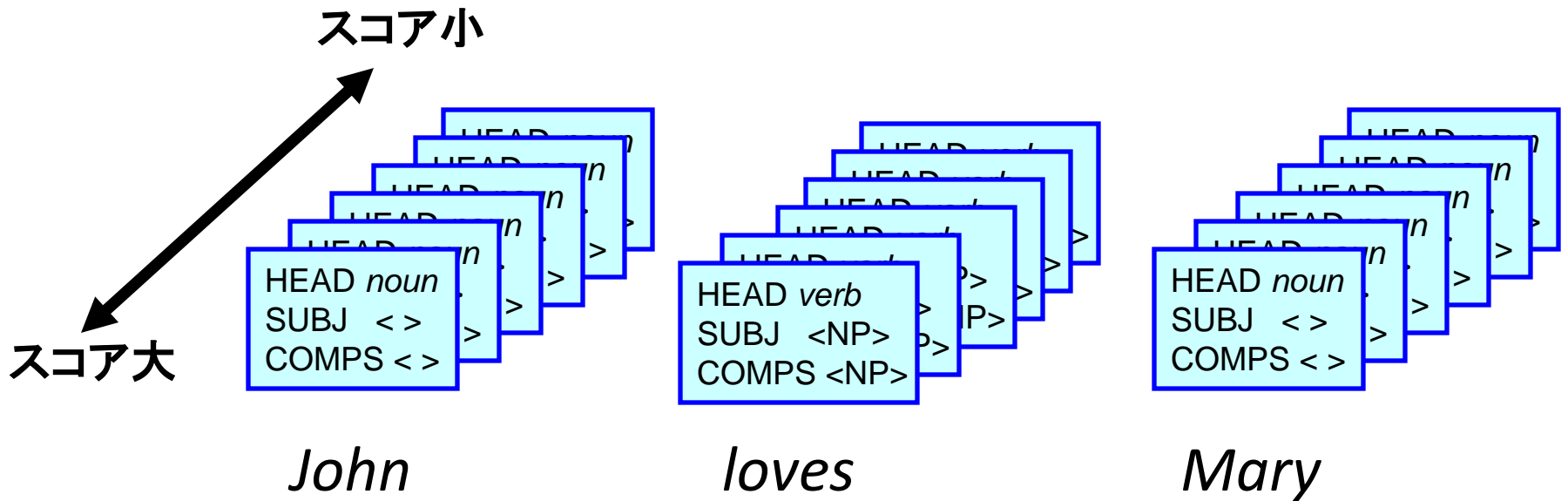
曖昧性解消

HPSG/CCG と CFG の違い

- 一般論: HPSG や CCG は難しそうに見える(?)
- 曖昧性解消を考えない場合
 - 文(統語構造)の集合が違う
- 曖昧性解消を考える場合
 - 曖昧性解消の単位が違う(後述)
- 計算量
 - 全解探索の計算量: $\text{HPSG} > \text{CCG} > \text{CFG}$
 - CCG: mildly context-sensitive
 - HPSG: type 0 (?)
 - 実際の解析時間: $\text{CFG} \gg \text{HPSG/CCG}$

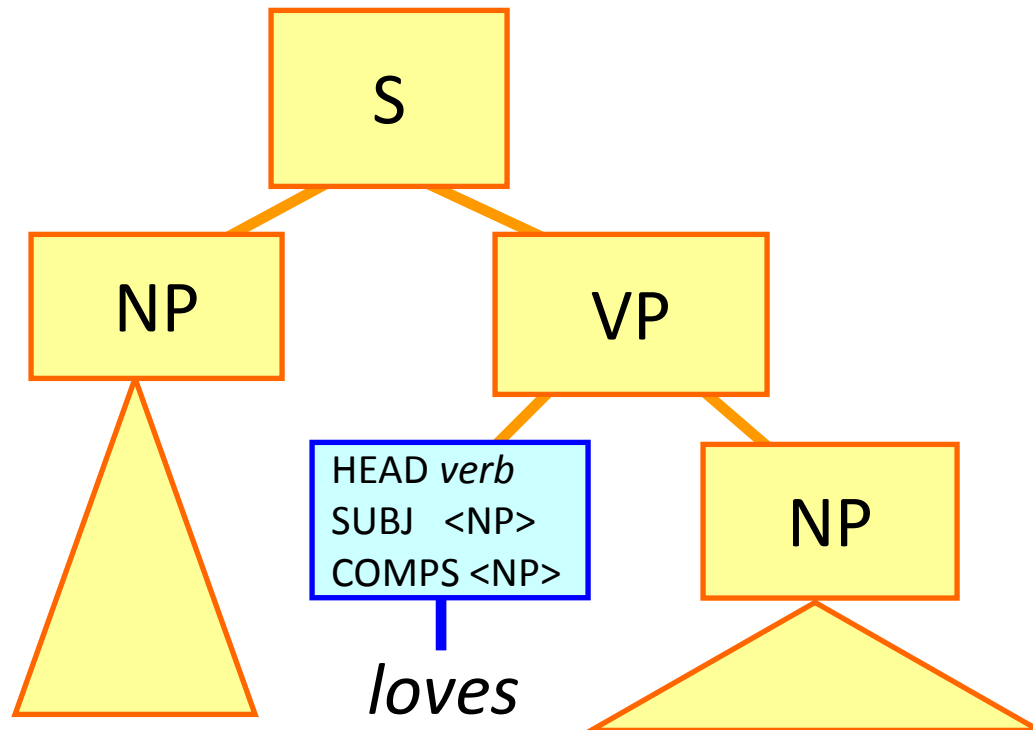
Supertagging

- Supertag = 語彙項目、語彙カテゴリ
- Supertagging = 単語に supertag を割り当てる



Supertagging = “almost parsing”

- Supertag が決まると、その上に作られる統語構造がほぼ決まる
- 曖昧性のほとんどが supertagging で解消



Supertagging は意外と簡単

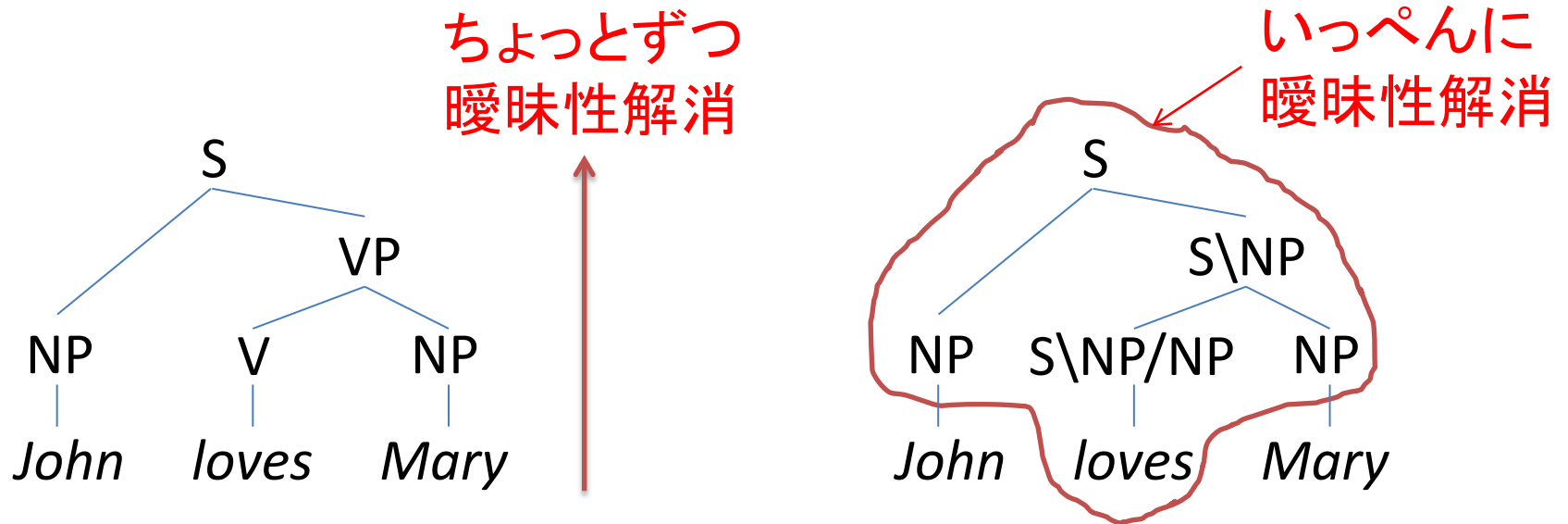
- Supertagging は系列ラベリングの一種
 - $O(n)$
 - 簡単な分類器でうまくいく
 - 多くの場合、supertag は局所的情報で決まる
- 構文解析のほとんどが、単純・高速な supertagging で終わってしまう

おそらく目的語コントロール動詞

... man **forced** his friend to ...
... NN VBD PRP\$ NN TO ...

文法と曖昧性解消

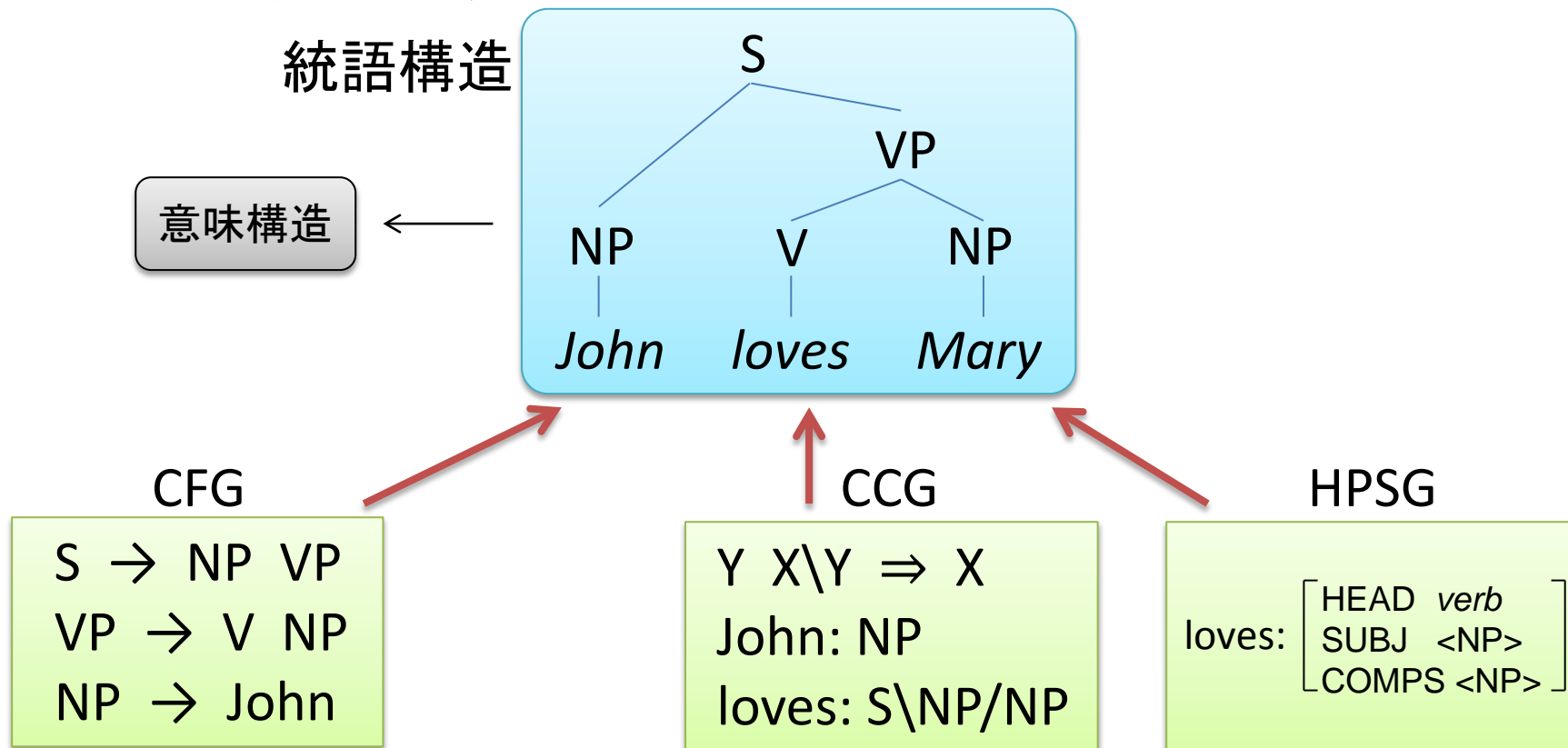
- CFG構文解析とは、曖昧性解消の単位が違う



- 現在の構文解析技術の多くが同様のアイデアに基づいている
 - Lexicalization, head percolation, symbol annotation...

文法と計算プロセス

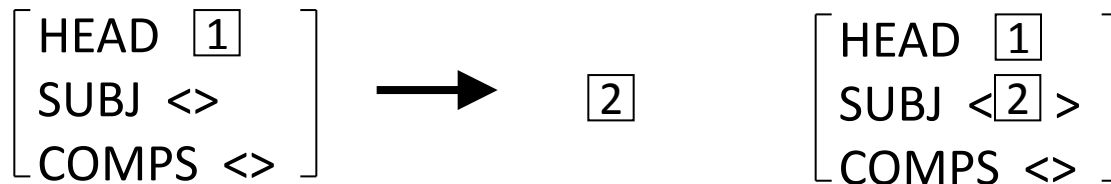
- 統語構造 = 構文解析の**計算結果**
- CFG と HPSG/CCG の違いは、計算方法(プロセス)の違いと見ることもできる



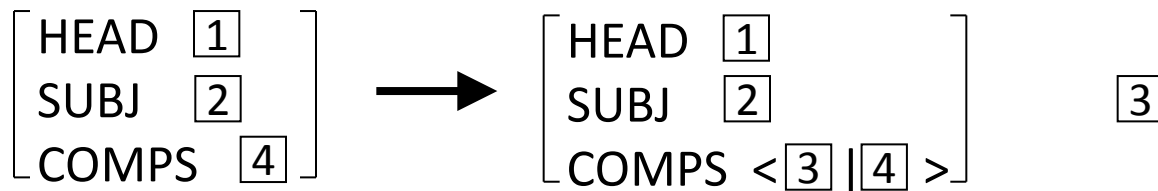
HPSG 構文規則

- HPSG の構文規則は、類似したCFG規則を一般化したもの

主語 + 主辞



主辞 + 目的語・補語



プリンシプル＝文法の文法

- 構文規則や語彙項目は、それ自身が**プリンシプル**から導出される
- 主辞素性プリンシプル
 - HEADの値は主辞から受け継がれる

[HEAD 1] → ^{主辞}[HEAD 1]

- 下位範疇化プリンシプル
- 意味合成プリンシプル
- 語彙規則
- ...

プリンシプル＝文法の文法

- プリンシプルから文法を導出する

主辞素性プリンシプル

[HEAD 1] → [HEAD 1]
head daughter

下位範疇化プリンシプル

[SUBJ 2] → [SUBJ 2]
head daughter

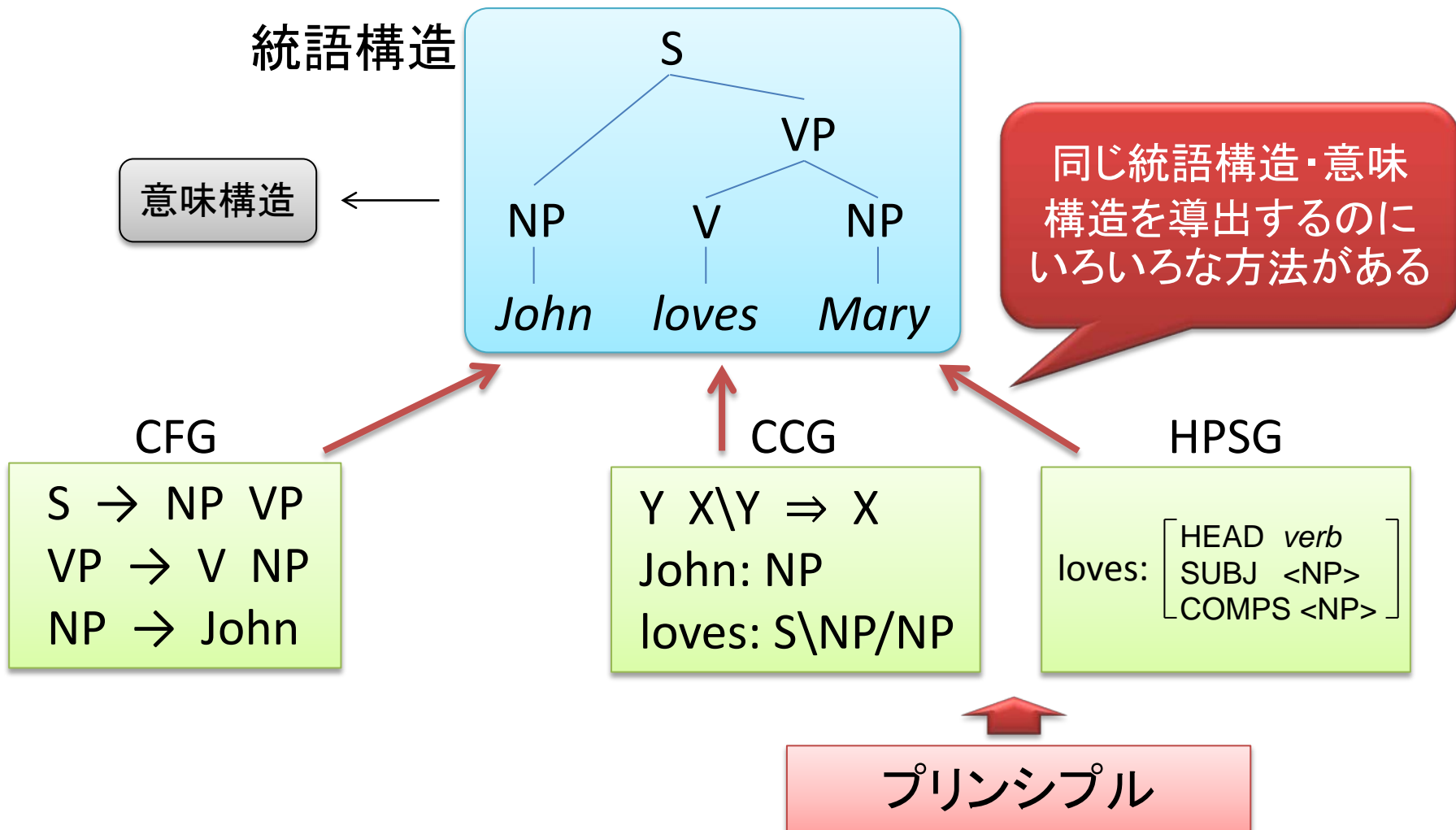
意味合成プリンシプル、...

導出

構文規則

[HEAD 1]
[SUBJ 2]
[COMPS 4] → [HEAD 1]
[SUBJ 2] [3]
[COMPS < 3 | 4 >]

いろいろな計算プロセス (=分解・一般化)の可能性



文法と計算プロセス

- CFG と HPSG/CCG の違いは、統語構造に至る計算方法(プロセス)の違い
 - HPSG と CCG の間にも重要な違いがある(割愛)
- 学習の観点から見ると、統語構造の分解(一般化)のしかたが違う
- 文法に分解(一般化)した上で、さらに一般化した構造を探求することもできる **＝言語理論**
 - cf. Grammar Matrix

おわりに

- 構文解析と文法、言語理論、曖昧性解消の関係について
- 自然言語処理における現状
 - 現在の構文解析手法は、言語理論が提示する一般化・構造(=モデル)まで達していない
 - 構文解析に本当に有効な一般化・構造は、言語理論のそれと一致するとは限らない
- 統語構造・意味構造を高精度、高速かつシンプルに計算できる「原理」を探したい