

音声言語処理のための構文解析ツールキット

田代敏久 森元 還
ATR 音声翻訳通信研究所

概要

頑健かつ高精度で、処理効率が良い構文・意味解析機構を目指して、多種多様な言語知識を有効に利用でき、外部モジュールと容易にリンクできる構文解析ツールキットの開発を進めている。本稿では、構文解析ツールキットの概要、および予備的な言語解析実験の結果を報告する。

A Parsinig Toolkit for Spoken Language Processing

Toshihisa Tashiro, Tsuyoshi Morimoto
ATR Interpreting Telecommunications Research Laboratories

Abstract

To obtain a robust, accurate and efficient parser, we are developing a parsing toolkit, which can be linked with external modules easily and can make good use of various knowledge sources. In this paper, we report the outline of the toolkit and the results of preliminary tests.

1はじめに

音声翻訳システムのような音声言語処理システムの構築のためには、頑健で解析精度が高く、処理効率も良い構文・意味解析機構の研究が必要である。また、単に解析能力の優劣だけではなく、他の処理モジュールとの協調的な動作が可能かどうかということも、重要な構文・意味解析機構の評価尺度である。しかし、従来の構文・意味解析の研究の多くは、以下の3つの型に分類することができ、それぞれに問題を抱えている。

• 文法理論重視型:

構文・意味解析の研究に、文法理論や他の言語学上の研究成果を取り入れるのは当然である。しかし、少なくとも現状の文法理論は、現実世界の多種多様な言語現象を矛盾なく説明できるほど緻密ではない。よって過度に特定の文法理論に依存していくは、頑健な解析機構の開発は不可能である。

• 計算メカニズム重視型:

構文・意味解析の研究に、ソフトウェア科学上の成果を取り入れることも重要かつ必要である。しかし、構文・意味解析機構は必然的に大量のデータ処理を必要とするのに対し、ソフトウェア科学

上の新しい計算メカニズムは、しばしば現実には不可能な資源(メモリおよび計算量)を要求するために、実際に動作する構文・意味解析機構の研究には向かないことが多い。また、新しい計算メカニズムは、記号処理の世界に閉じていることが多く、音声言語システムの開発に必要となる記号処理と非記号的情報処理との協調作業を困難にしている。

• 開発重視型:

上記の2つの型が理論重視なのに対し、もっぱら実践を旨とする研究も存在する。この種の研究では、かなり解析精度が高く、処理効率も良い結果を得ていることが多い。しかし、この種の研究は、ある特定のシステムに依存した処理機構やデータを前提としていたり、ドメインに依存したヒューリスティックを無批判に利用していたりする場合もある。また言語学的基盤や計算メカニズムが明確でないため、実験結果を客観的に評価することが困難な場合が多い。

我々は従来、文法理論としてはJPSGのような制約に基づく句構造文法、計算メカニズムとして単一化演算を基礎に、音声言語解析の研究を進めてきた[4]。しかし、特定の文法理論や計算メカニズムに深く依

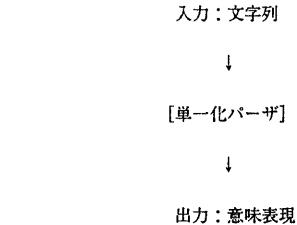


図 1: 単一化文法に基づく構文解析のイメージ

存した研究は、前述のような問題を持ち、我々が目標とする“自発的に発声された話し言葉 (spontaneous speech) の処理”を行なうためには不十分であることがわかつた。そこで、我々は構文解析機構の研究・開発を、以下のような方針で行なうこととした。

- 制約に基づく句構造文法の枠組を守りながらも、他の文法理論や言語学上の知見、コーパスから自動的に学習された統計的・量的な知識等を積極的に採り入れ、話し言葉に出現する幅広い言語現象を処理できる文法・言語モデルを開発する。
- 単一化演算をそのまま実装、利用するのは計算コスト、多様な言語的知識の柔軟な利用、等の点で問題があるので、より処理効率が良く、改良・改造や他のモジュールとのリンクが容易な計算機構を用意する。

本稿では、上記の方針に従って整備を進めている構文解析ツールキットの概要、および予備的な言語解析実験の結果を報告する。

2 構文解析ツールキットの概要

2.1 目標とする解析機構のイメージ

前述のように、单一化文法に基づく構文解析は、多様な言語的知識を柔軟に利用したり、他のモジュールとリンクしたりすることが困難である。これは、单一化文法に基づく構文解析機構は、基本的に図 1 で示すような硬直化した設計思想に基づいて作成されているためである。

このような設計思想で作成された構文解析機構には、以下のような問題がある。

- 文法・辞書には、統語的な知識や意味的・運用論的な知識を混在させて記述する必要があるので、大規模な語彙や言語現象をカバーすることが困難である。
- 出力が固定されているため、出力構造と相性の悪い外部モジュールは、パーザの出力を利用できない。

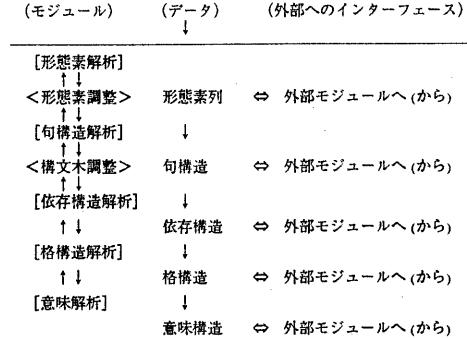


図 2: 望ましい構文解析機構のイメージ

そこで、我々は図 2 で示すような、1) 処理機構および知識がモジュール化されており、2) さまざまなレベルの出力を外部モジュールに提供可能な、構文・意味解析機構を目指す。

2.2 開発にあたっての留意点

我々は、以下のようないくつかの留意点に留意して構文・意味解析機構を開発している。

• 処理機構と知識を分離する

单一化に基づく構文解析のように、計算機構および知識の記述形式が統一されている場合には、特に注意しなくとも処理機構と知識は形式的に分離される。しかし、我々の枠組は、計算機構や知識の記述形式の多様性を認めるので、不注意により処理機構の中に知識が埋め込まれてしまう恐れがある。

• 制約と選好を区別する

单一化に基づく構文解析では、すべての言語的知識は制約として扱われる。しかし、長尾ら [3] が述べているように、制約としての知識と選好としての知識を区別して考慮しなくては、頑健で高精度な解析機構は開発できない。

• データ(構造)の書き換え(変換)のための手続き・知識と、曖昧性解消のための手続き・知識を区別する

单一化文法に基づく構文解析に限らず、従来の構文解析の研究の多くは、曖昧性解消の重要性を強調するあまり、解析の基本的な手続き・知識と曖昧性解消のための手続き・知識を区別することを怠ってきた。本当は、曖昧性解消の研究を効率的に進めるためにこそ、両者を区別することが重要なのである。我々が開発している構文解析ツールキットのすべてのモジュールは、以下の原則を遵

守しているので、曖昧性解消の研究をより柔軟に行なうことができる。

- 各モジュールは極めて単純なデータの変換・書き換え機能 (=基本機能) のみを有する。
- 曖昧性の解消は、各モジュール固有のヒューリスティックや統計情報、他のモジュールの処理結果等を利用して、基本機能とは分離可能な計算機構により行なう。

2.3 各モジュールの概要

前述の方針や原則に基づき開発を進めている各モジュールの概要を説明する。なお図3に、各モジュールの基本機能、基本機能に必要な知識、曖昧性の解消に必要な知識についてまとめる。

2.3.1 形態素解析モジュール

形態素解析モジュールとは、入力文字列を単語単位に分割し、品詞・活用形等の形態素情報ラベルを付与するモジュールである。辞書を入れ換えることにより、任意の形態素情報体系のもとで動作する必要がある。また、形態素解析時に発生する曖昧性の解消には、様々な知識が提案されているので、それらの知識を有効に利用できるメカニズムを用意する必要がある。図4に形態素解析モジュールの入出力例を示す。

2.3.2 形態素調整モジュール

形態素調整モジュールとは、ある形態素情報体系(単語の分割および品詞情報)に基づく形態素列を、別な形態素情報体系に基づく形態素列に書き換えるモジュールである。このモジュールは本来、異なる形態素情報体系で作成された形態素情報コーパスを有効利用するために開発されたが[5]、形態素解析と句構造解析の間での形態素情報体系の調整にも利用できる。このモジュールも、形態素解析モジュールと同様、任意の形態素情報体系で動作し、曖昧性解消のための様々な知識を利用できなくてはならない。図5に形態素調整モジュールの入出力例を示す。

2.3.3 句構造解析モジュール

句構造解析モジュールとは、形態素列を入力にとり、句構造規則に従い、構成要素(consituent)をノードとする木構造(句構造)を作成するモジュールである。

我々は、句構造解析に文脈自由文法を用いている。自然言語の持つ文脈依存性は、句構造解析モジュールとは独立した計算機構および知識を利用して対応する。図6に句構造解析モジュールの入出力例を示す。

2.3.4 構文木調整モジュール

構文木調整モジュールとは、ある構文解析文法に基づく句構造(構文木)を、別な句構造(構文木)に書き換

えるモジュールである。このモジュールは本来、異なる構文解析文法で作成された構文木コーパスを有効利用するために開発されたが[6]、句構造解析と依存構造解析間での構文解析文法の調整にも利用できる。図7に構文木調整モジュールの入出力例を示す。

2.3.5 依存構造解析モジュール

依存構造解析モジュールとは、句構造(構文木)を入力とし、構成要素(consituent)の統語的な主従関係を判定し、“語”¹をノードとするラベル無し有向グラフ²に変換するモジュールである。

一般的には、このモジュールは、次に述べる格構造モジュールと一体をなすものとして考えられているが、我々はできる限りモジュールを細分化する方針なので、あえて1モジュールとして独立させている。図8に依存構造解析モジュールの入出力例を示す。

2.3.6 格構造解析モジュール

格構造解析モジュールとは、依存構造(ラベル無し有効グラフ)を入力とし、格解析辞書を利用して“語”的主従関係の役割を決定し、格構造(ラベル付き有向グラフ)に変換するモジュールである。

なお、このモジュールはあくまでノード間の関係を決定するだけで、主従関係の逆転や、表層にない要素の生成等は行わない。格構造のノードは必ず表層の語と直接的な対応関係を保っている。図9に格構造解析モジュールの入出力例を示す。

2.3.7 意味解析モジュール

意味解析モジュールとは、ラベル付き有向グラフを、グラフ書き換え規則を利用して変形するモジュールである。

このモジュールは、主従関係の逆転、ノードの挿入・削除等、任意の変形操作を行うことができる。結果として、意味解析結果と表層の語とは直接的な対応関係は失われてしまう。図10に意味解析モジュールの入出力例を示す。

2.4 自由発話への対応

我々の解析対象は話し言葉なので、必ずしも文法的に適格な文が入力されるわけではない。そこで、不適格な入力を処理するために、Jensen[1]により提案された Fitted Parse の手法を用いることにした。Fitted Parse とは、

- 適格文を対象とする核文法(core grammar)を用いてボトムアップに統語解析する。

¹ 通常は単語と考えてよい。動詞等の活用する語は、複数の形態素で一つの単語を構成すると考える。

² ほとんどの場合、木構造で十分表現できる。

形態素解析モジュール	
基本機能	文字列から形態素列への変換
データ書き換え知識	辞書(形態素辞書)
曖昧性解消知識	
	単語の N-gram、品詞(ラベル)の N-gram、接続テーブル、最長一致等のヒューリスティック等
形態素調整モジュール	
基本機能	形態素列のデータ書き換え
データ書き換え知識	形態素列書き換え規則
曖昧性解消知識	単語の N-gram、品詞(ラベル)の N-gram、接続テーブル、最長一致等のヒューリスティック等
句構造解析モジュール	
基本機能	形態素列から句構造への変換
データ構造変換知識	句構造規則(文脈自由文法)
曖昧性解消知識	確率文法、規則の連鎖統計情報、語や句の共起関係、Mental OS 等のヒューリスティック等
機文木調整モジュール	
基本機能	句構造のデータ書き換え
データ書き換え知識	句構造書き換え規則
曖昧性解消知識	確率文法、規則の連鎖統計情報、語や句の共起関係、Mental OS 等のヒューリスティック等
依存構造解析モジュール	
基本機能	句構造から依存構造への変換
データ構造変換知識	注釈付き句構造規則
曖昧性解消知識	規則の適用確率、語や句の共起関係、統語的な制約、枝分かれ等に関するヒューリスティック等
格構造解析モジュール	
基本機能	依存構造(ラベルなしグラフ)から格構造(ラベル付きグラフ)への変換
データ構造変換知識	格情報辞書
曖昧性解消知識	意味素性、規則の適用確率、語や句の共起関係、シーケンス、ドメインに依存するヒューリスティック等
意味解析モジュール	
基本機能	ラベル付きグラフの書き換え
データ構造変換知識	書き換え規則
曖昧性解消知識	規則の適用確率、メタ規則等

図 3: 構文解析ツールキットの各モジュールの概要

- 統語解析に失敗した場合、保存されている途中結果(部分木)を出力し、後続の処理に委ねる。

という手法である。

この手法を実現するために、我々の構文解析ツールキットは以下のようないくつかの特徴を持っている。

- すべてのモジュールは、入力に対し、機械的に定義できないような制約(“入力は文でなくてはいけない”、等の制約)を科さない。
- すべてのモジュールは、処理に失敗した場合でも、部分的な解を出力することができる。
- 文字列、形態素列、句構造等、各レベルにおいてデータの分割・併合ツールを用意し、処理単位を再構成できるようにする。

3 解析実験

前節で概説した構文解析ツールキットは、まだ開発途上にある。しかし、句構造解析モジュール等は既に

ほとんど完成しているので、極めて簡単な解析実験を行なってみた。

解析対象は、ATRで作成している音声言語データベース[8]に含まれている2352文(64会話)である。これらの文は既に形態素解析および句構造解析が済んでおり、形態素解析結果を入力とし、句構造解析結果を正解ファイルとすることにより、句構造解析実験を容易に行なうことができる。なお、これらの文は必ずしも文法的に適格ではないので、図11で示すような、部分木(の集合)としてしか解釈できない文もかなりある。

使用した文法は、適格な日本語文を想定して作成された229規則からなる純粋な文脈自由文法(核文法)である。前節で述べたように、我々の構文解析ツールキットは部分的な解を出力することができるので、適格文のみを想定した文法を利用しても、必ず何らかの解析結果を得ることができる。

純粋な文脈自由文法を用いた構文解析では、当然のことながら大量の曖昧性が生じる。我々は、曖昧性解消のための知識として文法規則の統計情報を用いることにし、解析対象とは別の4237文(104会話)から、1)通常の確率文法と、2)文脈依存の確率文法の一種である北[2]により提案された言語モデルの2つを学習した。表1に実験条件を示す。

解析の戦略も極めて単純なものにした。句構造解析モジュールはボトムアップ探索を行なうチャートパーザなので、とりあえずチャートの弧の数が上限(現在は20000)に達するまで全解探索を行ない、入力が適格な文として解釈された場合には、文としての全ての解釈を結果として出力した。入力が不適格な文であったり、適格な文でもメモリ不足になった場合には、チャートに保持されている部分木を、左最長優先のヒューリスティックを利用して探索し、最大50通りの部分木の組合せを結果として出力した。こうして出力した木(または部分木の集合)のすべての中から、

1. たまたま最初に見つかった結果(FIRST-HIT)
2. 通常の確率文法を用いてスコアリングし、もっとも高いスコアを得た結果(PCFG)
3. 北により提案された言語モデルを用いてスコアリングし、もっとも高いスコアを得た結果(RULE-BIGRAM)

の3つの解を求め、評価した。なお解の評価は、Black[7]により提案された手法で行なった。表2に実験結果を示す。

今回の実験では、あらかじめ正しく形態素解析された入力という、現実にはあり得ない入力を用いているので、結果の数値の絶対値には意味がない。しかし、ある程度まとまった量の解析実験なので、北の言語モデル(RULE-BIGRAM)は曖昧性解消のため

文法規則(核文法)	229 規則
テスト集合	2352 文(64 会話)
最長	49 語
最短	2 語
平均	11.6 語
訓練集合	4237 文(104 会話)

表 1: 実験条件

	recall	precision	crossing
FIRST HIT	88.2%	88.8%	9.2%
PCFG	89.4%	90.5%	8.2%
RULE-BIGRAM	92.1%	92.9%	6.1%

表 2: 実験結果

の知識としてかなり優れている、と判断していいだろう。

4 おわりに

本稿では、現在開発を進めている構文解析ツールキットの概要と、ツールキットを利用した簡単な解析実験の結果を報告した。複雑で難解な構文・意味解析機構をできる限り単純なモジュールの組合せで実現することにより、様々な言語知識を有效地に利用したり、外部モジュールとのリンクが容易になることが期待される。

今後は、開発途中の各モジュールを完成させるとともに、

- 分割された各モジュールをどのように統合して利用するか。単純な階層型インターフェースでよいのか、あるいは別のインターフェースを研究する必要があるのか。
- 分散して管理されることになる言語知識の一貫性をどう保つか。

- 音声認識部とのインターフェースをどうするか。

等の検討を行ない、より頑健かつ高精度で処理効率が良い構文・意味解析機構を目指していく。

入力:
ニューワシントンホテルでございます。

出力:
((ニューワシントンホテル 固有名詞)(で 助動詞)
(ございま 補助動詞)(す 語尾)(。 記号))

図 4: 形態素解析モジュールの入出力例

入力:
((ニューワシントンホテル 固有名詞)(で 助動詞)
(ございま 補助動詞)(す 語尾)(。 記号))

出力:
((ニューワシントンホテル <固有名詞>)(で <助動詞>
(ございま <補助動詞>)(す <語尾>)(。 <記号>))

図 5: 形態素調整モジュールの入出力例

入力:
((ニューワシントンホテル <固有名詞>)(で <助動詞>
(ございま <補助動詞>)(す <語尾>)(。 <記号>))

出力:
(<文>
(<筋>
(<動詞句>
(<動詞>
(<名詞句>
(<固有名詞> ニューワシントンホテル))
(<助動詞> で))
(<補助動詞>
(<補助動詞語幹> ございま)
(<語尾> す)))
(<句点> 。))

図 6: 句構造解析モジュールの入出力例

入力:
(文
(主語文節 ((ニューワシントンホテル <固有名詞>
(で <助動詞>)))
(補語文節 ((ございま <補助動詞>)(す <語尾>
(。 <記号>))))

出力:
(<文>
(<筋>
(<動詞句>
(<動詞>
(<名詞句>
(<固有名詞> ニューワシントンホテル))
(<助動詞> で))
(<補助動詞>
(<補助動詞語幹> ございま)
(<語尾> す)))
(<句点> 。))

※文節文法から一般の句構造規則への調整結果を例として示す。

図 7: 構文木調整モジュールの入出力例

入力:
(<文>
(<筋>
(<動詞句>
(<動詞>
(<名詞句>
(<固有名詞> ニューワシントンホテル))
(<助動詞> で))
(<補助動詞>
(<補助動詞語幹> ございま)
(<語尾> す)))
(<句点> 。))

出力:
(。):(<句点>)
└ (ございま す):(<補助動詞語幹> <語尾>)
└ (で):(<助動詞>)
└ (ニューワシントンホテル):(<固有名詞>)

図 8: 依存構造解析モジュールの入出力例

入力:
 (。)<句点>
 ↘(ございます)<補助動詞語幹><語尾>
 ↘(て)<助動詞>
 ↘(ニューワシントンホテル)<固有名詞>

出力:
 [[語義見出し。]]
 [[カテゴリ <句点>]]
 [[任意 [[語義見出し *ございます*]]
 [[カテゴリ <補助動詞語幹>]]
 [[任意 [[語義見出し *です*]]
 [[カテゴリ <助動詞>]]
 [[OBJ *未定義*]]
 [[IDEN
 [[[[語義見出し
 ニューワシントンホテル]]
 [[カテゴリ <固有名詞>]]]]]]]]]

図 9: 格構造解析モジュールの入出力例

入力:
 [[語義見出し。]]
 [[カテゴリ <句点>]]
 [[任意 [[語義見出し *ございます*]]
 [[カテゴリ <補助動詞語幹>]]
 [[任意 [[語義見出し *です*]]
 [[カテゴリ <助動詞>]]
 [[OBJ *未定義*]]
 [[IDEN
 [[[[語義見出し
 ニューワシントンホテル]]
 [[カテゴリ <固有名詞>]]]]]]]]]

出力:
 [[RELN *INFORM*]
 [AGEN *SPEAKER*]
 [RECP *HEARER*]
 [OBJE [[RELN *COPULA*]
 [OBJE *UNSPECIFIED-COMPLEX*]
 [[IDEN
 [[RELN *ニューワシントンホテル*]]]]]]]]]

図 10: 意味解析モジュールの入出力例

(
 (<感動詞>
 (<副詞句>
 (<副詞> 大変))
 (<感動詞>
 (<感動詞> 申し訳ございません)
 (<読点> 、)))
 (<名詞句>
 (<人名> 鈴木)
 (<接尾辞> 様))
 (<句点> 。)
)

図 11: 不適格文の部分木による表現

参考文献

- [1] Jensen, K. and Heidorn, G.E: "The Fitted Parsing: 100% Parsing Capability in a Syntactic Grammar of English," ANLP83, 1983.
- [2] Kita, K et al.: "Continuously Spoken Sentence Recognition by HMM-LR," ICSLP-92, pp.305-308, 1992.
- [3] 長尾 碓, 丸山 宏: "自然言語処理における曖昧さとその解消、情報処理," Vol.33, No. 7, 1992.
- [4] Nagata, M. and Morimoto, T.: "A Unification-Based Japanese Parser for Speech-to-Speech Translation," IEICE Trans. Inf. & Syst., Vol.E76-D, No.1, pp.51-61 1993.
- [5] Tashiro, T., Uratani, N., Morimoto, T, "Restructuring Tagged Corpora with Morpheme Adjustment Rules", COLING94, 1994.
- [6] 田代敏久, 柏岡 秀紀, Ezra W.Black, "構文木コーパスの再構成手法", 情報処理学会第49回全国大会, 1994.
- [7] Black, E., et al.: "A Procedure for Quantitatively Comparing the Syntactic Coverage of English Grammars", DARPA Speech and Natural Language Workshop, 1991.
- [8] Morimoto, T. et al. : "A Speech and Language Database for Speech Translation Research", IC-SLP94, 1994.