

音声言語コーパスを用いた  
日本語話し言葉の構文解析に関する研究

大野 誠寛



## 概要

近年，音声処理，ならびに，自然言語処理の技術の発展を背景に，音声対話や音声翻訳，音声要約，会話マイニングなど音声言語処理システムに関する研究が盛んに行われている．しかし，現状の音声言語処理システムの多くは，あらかじめ定められた単語や言い回しなどキーワードを処理する方式に留まっている．より高度な処理を実現するための次なるステップとして，音声言語の構文情報の活用が検討されつつあり，音声言語に対する高い性能を備えた構文解析器の開発が望まれている．

一方，構文解析の研究は，これまで，主に書き言葉である新聞記事を対象として行われてきた．言語はその出現形態によって書き言葉と話し言葉に分類でき，さらに話し言葉は，複数の話者が交替で話す「対話」と一人の話者のみが話す「独話」に分類できる．これらはそれぞれ性質が異なるため，書き言葉を対象とした従来の構文解析手法を単に話し言葉に適用しただけでは，様々な問題が生じる．

本論文では，上述した音声言語処理システムの高度な話し言葉処理の実現に必要なとなる話し言葉の高性能な構文解析器を開発することを目的とする．この開発には，音響情報処理，言語情報処理，視覚情報処理などの複数の分野が互いに関連しており，各分野における技術的向上が不可欠であるが，本研究では，そのうち，自然言語処理技術による解決が主に求められる課題に焦点を絞り，解決することを試みる．すなわち，本研究の目標は，話し言葉の高性能な構文解析を実現するための要素技術として，

### 1) 頑健な構文解析手法

話し言葉（特に，自由対話）では，書き言葉にはない倒置やフィラー，言い淀み，言い直し，言い誤りなど，書き言葉の文法に逸脱する言語現象が頻出する．これら非文法的言語現象を含む文に対して，頑健な解析を実現する．

### 2) 効率的な構文解析手法

解説や講演など一人の話者のみにより話される独話では，書き言葉と比べ，文の切れ目に対する意識が低下するため，極端に長い文が頻出する．一般に，文が長くなればなるほどその解析時間は指数関数的に増加するため，解析効率が

低下することになる．このような話し言葉に頻出する極端に長い文に対して，従来の構文解析手法と同程度以上の解析精度を備えた高速な解析を実現する．

### 3) 漸進的な構文解析手法

一般に，書き言葉では，情報が文字列として一度に提示されるのに対して，話し言葉では，情報が音素列として時間軸上で逐次提示される．このため，話し言葉では，入力に追従して情報を処理する応用システムが考えられ（例えば，同時通訳），その要素技術となる構文解析においては，話し言葉の入力に対して漸進的な処理を行うことが求められる．特に，独話では，文が長くなる傾向にあり，文全体の入力を待って解析を開始すると著しく同時性が損なわれることになるため，漸進的な解析が望まれる．そこで，従来の構文解析手法と同程度の解析精度を維持しつつ，話し手の話速に追従できる程度の漸進性を備えた解析を実現する．

を開発することである．本研究では，音声言語コーパスに基づく統計的な手法を用いることにより，これら3つの構文解析手法を開発する．なお，本研究では日本語を構文解析の対象とする．

本論文は全6章から構成される．第1章は本論文の序論であり，話し言葉の構文解析に関する課題及び研究動向を示すとともに，本論文の位置づけとアプローチを述べたものである．

第2章では，話し言葉の構文的特徴を明らかにするための分析データ及び統計的構文解析手法の学習データとして利用することを目的に構築した，対話と独話の2つの構文構造付き音声言語コーパスについて述べる．対話の構文構造付き音声言語コーパスはCIAIR 車内音声対話コーパスに対して，独話の構文構造付き音声言語コーパスはNHKの解説番組「あすを読む」の書き起こしコーパスに対して，それぞれ構文構造を付与することにより構築する．両コーパスとも話し言葉に特有な言語現象に対しては新たな付与基準を設けている．また，構文構造付き音声独話コーパスには，節境界情報や複数の係り先を付与しているという特徴がある．構築した対話と独話の構文構造付きコーパスは，それぞれ，85,870形態素，192,495形態素規模を備えている．

第3章では，大規模音声言語コーパスを用いた話し言葉の頑健な係り受け解析手法を提案する．本章の研究では，フィラーや言い淀み，倒置などの非文法的言語現象が頻出する対話文を構文解析の対象とした．実際に，第2章で構築したCIAIR 構文構造付き音声対話コーパスを分析した結果，従来の係り受け解析手法では，係り

受けの非交差性，後方修飾性，係り先の唯一性の3つの制約が用いられてきたが，対話音声では，後方修飾性を満たさない倒置現象や係り先の唯一性を満たさない文節などを含む発話が頻出することが分かった．そこで本手法では，後方修飾性の制約及び係り先の唯一性に関する制約は統計情報を反映させつつ緩和する．また，本手法では，構築した構文構造付き音声対話コーパスから各文節間の係り受け確率を統計的に獲得し，それを用いて係り受け構造の尤度を計算する．これにより，非文法的な特徴をもつ発話の解析が可能になる．CIAIR 構文構造付き音声対話コーパスに対して係り受け解析実験を行い，その結果，本手法により，自然発話文に対しても，書き言葉を対象とした従来の係り受け解析手法と同等の高い精度で係り受けを抽出できることを確認した．特に，係り先を持たない文節と倒置，発話単位をまたぐ係り受けの解析に対する本手法の頑健性を明らかにした．

第4章では，文の分割に基づく話し言葉の効率的な係り受け解析手法を提案する．本手法では，文の分割単位として節を採用し，節レベルと文レベルの二段階で係り受け解析を実行する．節は，構文的かつ意味的にまとまった単位であるため，文に代わる解析単位として利用できると考えられるためである．まず，節境界解析により文を節に分割し，各節に対して係り受け解析を行うことにより，節内の係り受け関係を同定する．次に，節境界をまたぐ係り受け関係を定め，文全体の係り受け構造を作り上げる．これにより，話し言葉に出現する極端に長い文に対する効率的な解析の実現が期待できる．極端に長い文が頻出する独話データとして，第2章で構築した「あすを読む」構文構造付き音声独話コーパスを用いた係り受け解析実験を行い，その結果，本手法により，従来の係り受け解析手法と比べ，解析精度を改善しつつ解析時間を約  $1/5$  に短縮できることがわかった．

第5章では，話者による音声入力に従って順次，解析を行う漸進的係り受け解析手法を提案する．本手法では，独話音声に対して，節が入力されるたびにその節の内部の係り受け構造を作り上げるとともに，すでに入力されている節の係り先を決定することを試みる．節の係り先となる文節の決定は，後続するいくつかの文節との係り受けの尤度を考慮した動的なタイミングで行う．これにより，独話の入力途中の段階で構造情報を随時出力する漸進的な解析が可能となる．「あすを読む」構文構造付き音声独話コーパスを用いた係り受け解析実験を行い，その結果，本手法により，従来の係り受け解析手法と同程度の解析精度と解析時間を備えつつ，解析の漸進性の向上が可能となることを確認した．

最後に，第6章において本論文を総括し，今後の研究課題ならびに将来の展望について示す．



# 目次

第1章	まえがき	15
1.1	話し言葉の構文解析	15
1.2	話し言葉の構文解析に関する研究動向	17
1.2.1	非文法的言語現象を含む文の構文解析	17
1.2.2	音声認識誤りを含む文の構文解析	18
1.2.3	長い文の構文解析	19
1.2.4	逐次的な情報提示に対する構文解析	20
1.2.5	非言語情報を用いた構文的曖昧性解消を行う構文解析	21
1.3	本論文の目的	21
1.4	本論文の内容	23
1.5	本論文の構成	24
第2章	構文構造付き音声言語コーパスの構築	27
2.1	はじめに	27
2.2	構文構造付き音声対話コーパスの構築	28
2.2.1	CIAIR 車内音声対話コーパス	28
2.2.2	CIAIR 構文構造付き音声対話コーパス	30
2.2.3	CIAIR 構文構造付き音声対話コーパスの構築手順	33
2.2.4	評価実験	38
2.2.5	CIAIR 構文構造付き音声対話コーパスの規模と特徴	40
2.3	構文構造付き音声独話コーパスの構築	41
2.3.1	「あすを読む」書き起こしコーパス	41
2.3.2	「あすを読む」構文構造付き音声独話コーパス	42
2.3.3	「あすを読む」構文構造付き音声独話コーパスの構築手順	47
2.3.4	「あすを読む」構文構造付き音声独話コーパスの規模と特徴	50
2.4	2章のまとめ	51

<b>第3章</b>	<b>日本語話し言葉の頑健な係り受け解析</b>	<b>53</b>
3.1	はじめに	53
3.2	自然発話の言語的分析	54
3.3	統計的係り受け解析	56
3.3.1	係り受けの構文的制約	56
3.3.2	話し言葉の統計的係り受け解析	56
3.3.3	解析例	59
3.4	解析実験	60
3.4.1	実験の概要	60
3.4.2	実験結果	60
3.5	考察	61
3.5.1	話し言葉に特有な現象に対する頑健性	61
3.5.2	データスパースネス問題に対する頑健性	65
3.6	3章のまとめ	65
<b>第4章</b>	<b>日本語話し言葉の効率的な係り受け解析</b>	<b>67</b>
4.1	はじめに	67
4.2	独話文の係り受け解析における処理単位	68
4.2.1	節と係り受け	68
4.2.2	節境界単位	69
4.2.3	節境界単位と係り受けの関係	69
4.3	節境界に基づく係り受け解析	70
4.3.1	節レベルの係り受け解析	71
4.3.2	文レベルの係り受け解析	73
4.4	解析実験	73
4.4.1	実験の概要	73
4.4.2	実験結果	74
4.5	考察	76
4.5.1	節境界解析エラーの影響	76
4.5.2	節境界単位内部の文節に対する解析精度	78
4.5.3	節境界単位の最終文節に対する解析精度	78
4.5.4	節境界をまたぐ係り受け関係	79
4.6	関連研究	82

4.7	4章のまとめ	83
<b>第5章</b>	<b>日本語話し言葉の漸進的な係り受け解析</b>	<b>85</b>
5.1	はじめに	85
5.2	独話の漸進的係り受け解析における処理単位	86
5.2.1	節と節境界単位	86
5.2.2	節境界単位の分析	87
5.2.3	漸進的係り受け解析の処理単位としての節境界単位	87
5.3	節境界に基づく漸進的係り受け解析	88
5.3.1	独話レベルの係り受け解析	89
5.4	漸進的係り受け解析アルゴリズム	90
5.4.1	独話レベルの漸進的解析アルゴリズム	91
5.4.2	解析例	92
5.5	解析実験	94
5.5.1	実験に使用したデータ	94
5.5.2	実験の概要	94
5.5.3	実験結果	95
5.6	文末検出性能	98
5.7	5章のまとめ	99
<b>第6章</b>	<b>あとがき</b>	<b>101</b>
6.1	本論文のまとめ	101
6.2	今後の課題と将来への展望	102



## 図一覽

1.1	言語の分類	16
1.2	本論文の構成	25
2.1	データ収集車	29
2.2	CIAIR 車内音声対話コーパスの書き起こしデータ	29
2.3	フィラー・言い淀みを含む文の係り受け構造	31
2.4	受け文節が省略された文節を含む文の係り受け構造	31
2.5	言い直しを含む文の係り受け構造	31
2.6	発話単位を越える係り受けを含む文の係り受け構造	32
2.7	CIAIR 構文構造付き音声対話コーパスの例	34
2.8	「今日朝パン食べてお昼はおそばを食べたんですよ」の係り受け構造	35
2.9	CIAIR 構文構造付き音声対話コーパス構築の流れ	36
2.10	解析結果修正用インタフェース	38
2.11	統計的文節まとめあげの実験結果	39
2.12	統計的係り受け解析の実験結果	40
2.13	「あすを読む」書き起こしコーパスの例	42
2.14	「今朝東京へ行く途中で考えた」に対する係り受け構造	45
2.15	「小泉さんが二倍近くの差をつけて圧勝した」に対する係り受け構造	45
2.16	「あすを読む」構文構造付き音声独話コーパスの例	46
2.17	「最高裁判所は今日検察側が死刑を求めて上告をしておりました(強盗殺人事件について二審と同じように無期懲役の判決を言い渡しております)」の係り受け構造	47
2.18	修正インタフェース DSMT	48
3.1	文節間距離と係り受け数の関係	55
3.2	「えーと/コンビニ/ないかな/<pause>/そ/そこの/近くに」に対する係り受け構造	59
3.3	最終文節以外で受け文節がない係り受けの内訳	61

3.4	受け文節のない係り受けの例 . . . . .	62
3.5	前方の文節への係り受けの例 . . . . .	63
3.6	前方の文節への係り受けの文節間距離 . . . . .	63
3.7	ポーズをまたぐ係り受けの例 . . . . .	64
3.8	係り受け及びターンに対する正解率と学習データ量の関係 . . . . .	66
4.1	節と係り受けの関係 . . . . .	69
4.2	節境界をまたぐ係り受けの例 . . . . .	70
4.3	文の長さと解析時間の関係 . . . . .	75
4.4	節境界解析の適合率を低下させる節境界検出失敗によって生じた係り 受け解析の失敗例 . . . . .	77
4.5	節境界解析の再現率を低下させる節境界検出失敗によって生じた係り 受け解析の失敗例 . . . . .	77
4.6	係り受けがまたぐ節境界の種類とその割合 . . . . .	80
4.7	節境界「主題八」をまたぐ係り受けの例 . . . . .	80
4.8	節境界「連体節」をまたぐ係り受けの例 . . . . .	81
4.9	節境界「テ節」をまたぐ係り受けの例 . . . . .	82
5.1	節境界単位と係り受けの関係 . . . . .	88
5.2	漸進的係り受け解析の流れ . . . . .	91
5.3	漸進的係り受け解析の例（不変化閾値3の場合） . . . . .	93
5.4	不変化閾値と1番組あたりの解析時間の関係 . . . . .	97
5.5	不変化閾値と平均遅延時間の関係 . . . . .	97

## 表一覧

2.1	話し言葉に特有な言い回し表現 . . . . .	32
2.2	CIAIR 構文構造付き音声対話コーパスの規模 . . . . .	40
2.3	節境界の大分類と小分類 . . . . .	43
2.4	「あすを読む」構文構造付き音声独話コーパスの規模 . . . . .	50
3.1	係り受け分析データ . . . . .	54
3.2	係り文節と係りの種類の例 . . . . .	58
3.3	文節間係り受け確率 . . . . .	59
3.4	実験結果（係り受け正解率） . . . . .	60
3.5	受け文節のない係り受けの解析結果（ポーズの直前でなく、フィラー・ 言い淀み以外） . . . . .	62
3.6	前方の文節への係り受けの解析結果 . . . . .	62
3.7	ポーズをまたぐ係り受けの解析結果 . . . . .	64
4.1	「あすを読む」200 文の基礎統計 . . . . .	70
4.2	実験データ（「あすを読む」構文構造付き音声独話コーパス） . . . . .	74
4.3	実験結果（解析時間） . . . . .	75
4.4	実験結果（係り受け正解率） . . . . .	76
4.5	節境界解析プログラム CBAP の解析結果 . . . . .	76
4.6	本手法と従来手法の比較（節境界単位内部の係り受け解析結果） . . . . .	78
4.7	本手法と従来手法の比較（節境界単位の最終文節に対する解析結果） . . . . .	79
4.8	節境界をまたぐ係り受けに対する解析精度 . . . . .	79
5.1	実験データ（「あすを読む」構文構造付き音声独話コーパス） . . . . .	94
5.2	不変化閾値ごとの係り受け正解率 . . . . .	96
5.3	CBAP の節境界解析結果 . . . . .	96
5.4	文末検出の適合率・再現率・F 値 . . . . .	98



# 第1章 まえがき

## 1.1 話し言葉の構文解析

近年，音声処理，ならびに，自然言語処理の技術の発展を背景に，音声言語処理システムの研究が盛んに行われている．音声は，人間にとって最も自然で手軽な情報伝達手段であり，機械とのインタフェースに利用することが考えられる．そのような応用システムとして，音声対話 [43, 51, 77] や音声検索（音声による情報検索） [16, 17, 63] があり，カーナビゲーションシステムなど，実用化を目指した研究が展開されている．また，社会のグローバル化が急速に進み，異言語間音声コミュニケーションの機会が増加しているものの，言葉の壁は依然として残されている．そのような人間同士のインタラクションを支援することを目的に，音声翻訳システム [64, 93, 98] が研究されている．さらに，会議や講演，解説などで話された言葉から得られる情報を知的資源とみなし，蓄積された音声データへの効率的なアクセスやその効果的な再利用を実現するために，音声データの検索 [31, 71, 79] や音声要約 [12, 26, 100]，会話マイニング [70, 81, 101] などのシステムが開発されている．

しかし，現状の音声言語処理システムの多くは，あらかじめ定められた単語や言い回しなどキーワードを処理する方式に留まっている．より高度な処理を実現するための次なるステップとして，音声言語の構文情報の活用が検討されつつあり [25]，音声言語に対する高い性能を備えた構文解析器の開発が望まれている．

一方，構文解析の研究は，これまで，主に書き言葉である新聞記事を対象として行われてきた．言語はその出現形態によって書き言葉と話し言葉に分類でき，さらに話し言葉は，複数の話者が交替で話す「対話」と一人の話者のみが話す「独話」に分類できる（図 1.1 参照）．これらはそれぞれ性質が異なるため，従来の書き言葉を対象とした構文解析手法を単に話し言葉に適用しただけでは，少なくとも以下のような問題が生じる．

### (1) 非文法的な言語現象の処理

話し言葉（特に，自由対話）では，書き言葉にはない倒置やフィラー，言い淀み，言い直し，言い誤りなど，書き言葉の文法に逸脱する言語現象が頻出す

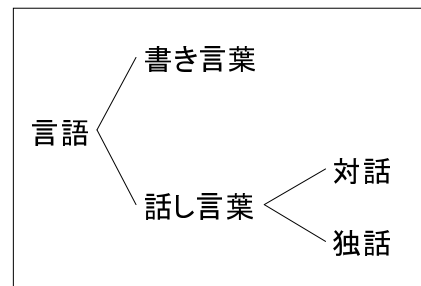


図 1.1: 言語の分類

る<sup>1</sup>．従来の構文解析手法では，このような言語現象の出現を考慮していないため，取り扱うことができない．

## (2) 音声認識誤りに対する処理

話し言葉に対して自動的に構文構造を付与するためには，話し言葉を音声認識し，その結果に対して構文解析する必要がある．近年の音声認識技術の進展は目覚ましいものの，100%の認識精度を達成することは難しく，音声認識誤りに対処する必要がある．音声認識誤りには，単語の変換誤りのほか，文末記号（ピリオド，句点など）の挿入誤りが存在する．話し言葉では，文末表現が明示的に使われるとは限らず，文境界を同定することは難しい．書き言葉を対象とした従来の構文解析手法では，音声認識結果を解析対象としていないため，このような誤りが含まれていることは想定していない．

## (3) 極端に長い文の処理

解説や講演など一人の話者のみにより話される独話では，書き言葉と比べ，文の切れ目に対する意識が低下するため，極端に長い文が頻出する．一般に，文が長くなればなるほどその解析時間は指数関数的に増加するため，解析効率が低下することになる．

## (4) 逐次的な情報提示に対する処理

一般に，書き言葉では，情報が文字列として一度に提示されるのに対して，話し言葉では，情報が音素列として時間軸上で逐次提示される．このため，話し言葉では，入力に追従して情報を処理する応用システムが考えられ（例えば，同時通訳），その要素技術となる構文解析においては，話し言葉の入力に対し

<sup>1</sup>独話では，事前に準備した内容が話されることが多いため，対話と比べ，このような言語現象は少ない傾向にある．

て漸進的な処理を行うことが求められる．特に，独話では，文が長くなる傾向にあり，文全体の入力を待って解析を開始すると著しく同時性が損なわれることになるため，漸進的な解析が望まれる．しかし，従来の書き言葉に対する構文解析では，このような解析性能は要求されていないため，上記の要請を満たすことはできない．

(5) 非言語情報による構文的曖昧性解消に対する処理

話し言葉では，言語情報とともに，動作や音響情報などの非言語情報を利用した情報伝達が行われる．従って，言語情報だけでは構文的曖昧性が存在する言葉を，非言語情報を利用しつつ発話することにより，その曖昧性を解消し正確に情報を伝えることができる．例えば，商品の検索などで，文脈なしに「昔の中国の本」と発話するとき，「昔の」の後にポーズを入れたり，「中国の本」を話速を早めつつ音量を上げて発話することにより，「昔の」が「中国の本」を修飾するということを正しく伝えることができる．従来の構文解析手法は，非言語情報を利用していないため，このような発話を正しく解析できない．

これらの問題を解決するため，これまでいくつかの研究が行われている．次節以降では，まず，これらの問題の観点から話し言葉の構文解析に関する研究動向を概観し，次に，本論文の目的と内容について述べる．

## 1.2 話し言葉の構文解析に関する研究動向

本節では，前節で挙げた各問題ごとに，話し言葉の構文解析に関する研究動向を示す．

### 1.2.1 非文法的言語現象を含む文の構文解析

話し言葉には倒置やフィラー，言い淀み，言い直しなどが頻出する．話し言葉の高精度な構文解析を実現するには，これらの話し言葉に特有な言語現象を含む文を頑健に解析する必要がある．ここでは，非文法的な言語現象を含む文の構文解析を試みた研究を概観する．

話し言葉を分析し，人手で辞書や規則を作成して，非文法的な言語現象に対処した構文解析の研究がいくつかあり，例えば，

- 句構造文法に基づく構文解析において、話し言葉に特有な言語現象を扱うための句構造規則や語の辞書を新たに作成することにより、非文法的な言語現象に対処した手法 [11, 69]
- 句構造文法に基づく構文解析に言い直しを扱うための特別な仕組みを持たせることにより、構文解析と並行して言い直しを処理する手法 [20, 24]（ただし、助詞省略や倒置は、助詞の有無、及び、前方依存の可能性を辞書に記述することにより、処理している。）
- 対話文の分析を通して、助詞省略と倒置に関するヒューリスティックを作成し、それを用いて助詞省略文や倒置文を解析する手法 [99]

などが挙げられる。しかし、これらの手法は、人手で作成した辞書や文法規則、優先規則の網羅性に問題がある<sup>2</sup>。さらに、解析対象を変更したり拡大するたびに、規則の修正や追加を人手により行う必要があるため、保守管理が煩雑になるという問題が生じる。

一方、コーパスに基づく解析手法として、伝 [14] は、言い直しや言い淀みなどを語と語の間の係り受け関係の一種と捉え、従来の係り受け解析を拡張した解析手法を提案している。しかし、扱える話し言葉の言語現象は限られている。

### 1.2.2 音声認識誤りを含む文の構文解析

音声言語処理システムにおいて、話し言葉の構文情報を利用するためには、話し言葉を音声認識し、その結果に対して構文解析する必要があり、音声認識誤りを頑健に処理することが望まれる。

話し言葉の音声認識結果に対して構文解析を試みた研究には、以下のような研究がある。

- 古瀬ら [21] は、多言語音声翻訳の前処理として音声認識結果に対する構文解析を試みている。この手法では、音声認識誤りにより生じる、文法的な逸脱表現を含む文や、句読点なしで連続する複数の文から、構文解析が可能な部分を取り出すことにより、部分的な構文解析を実現している。

<sup>2</sup> 自然言語を解析するための辞書や文法規則、優先規則を完全に網羅することは難しい。しかし、文法情報が付与された大量の実データから自動的に獲得することにより、より網羅性の高い辞書や文法規則、優先規則を作成することができる。

- 船越ら [19] は、前節で挙げた、言い直し、助詞落ち、倒置などの非文法的な言語現象を含む文に対して構文解析を行う手法 [20] を音声認識結果に適用し、この手法の頑健性を評価している。
- Shitaoka ら [85] は、話し言葉、特に独話では、句点が明示的でないため、音声認識により文末記号を挿入することが難しいことに着目し、文境界推定を行いつつ、係り受け解析を行う手法を提案している。この手法では、係り受け情報を利用して文境界推定を行い、その結果を係り受け解析にフィードバックし、再度、係り受け解析を実行することにより、解析精度を改善している。

しかし、いずれの手法も十分な解析精度が得られているとは言い難く、現状では音声認識技術の一層の進展が望まれる。

### 1.2.3 長い文の構文解析

解説や講演など一人の話者のみによる話し言葉では極端に長い文が頻出する。話し言葉の高性能な構文解析を実現するには、話し言葉に出現する極端に長い文に対しても、解析精度を維持しつつ、効率的に解析することが必要となる。

講演や解説などで話される話し言葉データは、CSJ プロジェクト [53]、LDC<sup>3</sup>、ELRA<sup>4</sup>等で、最近になって蓄積され始め、その研究環境がようやく整いつつある段階である。現在のところ、話し言葉で出現する長い文の解析に焦点を当てた構文解析の研究は見当たらない。

一方、書き言葉の構文解析において、新聞記事で見られる長文を対象にした研究がいくつか存在する。しかし、以下に挙げるように、そのほとんどが解析精度の向上を目的とした研究である。

- Lyon と Dickerson [52] は、英語文を pre-subject, subject, predicate の3つの構成素に分割することにより、構文的曖昧性を軽減し、解析精度を改善している。
- Agarwal と Boggess [1] や黒橋ら [48] は、長文解析の精度を高めるため、並列構造を構成する文節列の類似性を動的計画法を用いて検出し、並列関係を解析している。

<sup>3</sup>Linguistic Data Consortium, <http://www.ldc.upenn.edu/>

<sup>4</sup>European Language Resources Association, <http://www.elra.info/>

- 白井ら [84] や宇津呂ら [94] は、構文解析の曖昧性を解消するために、従属節の末尾に現れる表層表現から従属節を階層的に分類し、その順序関係を利用している。
- 市丸ら [28] は、接続助詞などのすべての機能語間の結合順位をコーパスから獲得し、その結合順位を確率的文脈自由文法 (PCFG) の生成規則における適用確率値の計算に埋め込むことにより、長文に対する解析精度を向上させている。

ただし、これらの研究では、解析効率について考慮していない。

一方、解析時間の短縮に焦点をあてた書き言葉の構文解析の研究には以下の研究がある。

- Sekine [83] は、決定的有限状態変換器を用いて線形時間で処理を行う構文解析器を提案している。この手法では、係り先の場所の 97% は文中の 5 つの候補で網羅されるという調査に基づき、考慮する係り先の文節数を制限することにより、解析の効率化を図っている。
- 颯々野 [82] は、文節数に対して線形時間で処理を行う日本語係り受け解析のアルゴリズムを提案している。ただし、この解析アルゴリズムは、文節が必ず後方の文節に係るという構文的制約を仮定して、解析アルゴリズムの単純化を図っている。

しかし、Sekine の手法は、考慮する係り先の文節数を制限することにより、解析精度が大きく損なわれており、書き言葉に対する他の構文解析手法と比べて精度が低い。また、颯々野の手法で仮定している構文制約は、書き言葉では満たされるものの、倒置の場合には満たされないため、この手法を話し言葉の解析に利用することはできない。

#### 1.2.4 逐次的な情報提示に対する構文解析

話し言葉では、入力に追従して情報を処理する応用システムが考えられ、その要素技術として漸進的構文解析が求められる。漸進的構文解析とは、自然言語文をその単語の出現順序に従って順次解析し、文の入力途中の段階でその構文構造を捉える枠組である。

これまでに提案された漸進的構文解析手法として、文脈自由文法に基づく手法 [57] や確率文脈自由文法に基づく手法 [34, 80]、範疇文法に基づく手法 [22, 62]、木接合

文法に基づく手法 [33] などがある．いずれの手法も，入力文中の単語が先頭から順に一単語ずつ入力されるごとに，それまでに入力された文の断片に対する構文構造を生成する．しかし，これらの手法では，英語を対象とした句構造規則に基づく漸進的構文解析手法となっており，英語と比べ文法的制約が弱い日本語への適用可能性は明らかでない．

一部，日本語を対象として文脈自由文法に基づく漸進的な構文解析手法が提案されている [2]．しかし，この研究では，句構造規則を手で作成しており，網羅性の高い規則の実現性については明らかではない．

### 1.2.5 非言語情報を用いた構文的曖昧性解消を行う構文解析

話し言葉では，言語情報とともに伝えられる非言語情報により，構文的曖昧性が解消される場合がある．このような場合，言語情報だけでなく，非言語情報も考慮して構文解析を行う必要がある．

非言語情報を用いて構文的曖昧性解消を行う研究としては，韻律情報を用いて構文解析を行う研究があり，例えば，

- 構文的曖昧性解消の手がかりとして韻律句境界の情報を利用した構文解析手法 [96]
- 韻律情報を用いて隣接句間の結合度をヒューリスティックに定義し，文構造を推定する手法 [44]
- ポーズやピッチ，パワーの韻律情報と係り受け距離の関係を統計的に推定し，この情報を利用した係り受け解析手法 [15]

などが提案されている．これらの研究では，韻律情報を用いることにより解析精度が向上することが報告されているが，さらなる精度改善の余地が残されている．

## 1.3 本論文の目的

1.1 節で挙げた音声言語処理システムの高度な話し言葉処理を実現するためには，話し言葉の高性能な構文解析器を開発する必要がある．この開発には，音響情報処理，言語情報処理，視覚情報処理などの複数の分野が互いに関連しており，各分野における技術的向上が不可欠であるが，本研究では，そのうち，言語情報処理的側面から問題を捉えることとし，自然言語処理技術による解決が主に求められる，(1)，

(3), (4) の課題に焦点をあてる．すなわち，本研究では，話し言葉の高性能な構文解析を実現するための要素技術として，

1. 頑健な構文解析手法

書き言葉には現れない倒置やフィラー，言い淀みなどの文法から逸脱する言語現象を含む文に対して，頑健な解析を実現する．

2. 効率的な構文解析手法

話し言葉に頻出する極端に長い文に対して，従来の構文解析手法と同程度以上の解析精度を備えた高速な解析を実現する．

3. 漸進的な構文解析手法

従来の構文解析手法と同程度の解析精度を維持しつつ，話し手の話速に追従できる程度の漸進性を備えた解析を実現する．

をそれぞれ開発することを具体的な目標とする．なお，本論文では，これらの要素技術を統合した構文解析器の開発には立ち入らない．

本研究では，音声言語コーパスに基づく統計的な手法を用いることにより，これら3つの構文解析手法を開発する．構文解析は，一般に，規則主導型の手法と統計的な手法の2つに分類できる．規則主導型の構文解析手法では，文法規則や優先規則が人手で記述されるのに対して，統計的な構文解析では，コーパスから統計的に推定した確率モデルによって文法規則や優先規則が表現される．最近では，規則主導型の構文解析手法における文法規則の網羅性や保守管理の問題から統計的な構文解析手法に注目が集まっており，書き言葉を対象として多くの成果が得られている[8, 10, 18, 23, 46, 78, 92]．本研究では，非文法的な言語現象が頻出する話し言葉を対象としており，書き言葉以上に文法規則の網羅性の問題が深刻になると考え，統計的な構文解析手法を採用する．また，学習データ及び分析用データとして，対話と独話に対して構文構造を付与した音声言語コーパスを構築した．

特に，頑健な構文解析手法では，音声言語コーパスから獲得した統計情報を反映させつつ構文的制約を緩和することにより，非文法的言語現象を網羅的かつ統一的に扱う．また，効率的な構文解析手法では，文を分割することにより解析時間の改善を図るが，その際，構文解析に適した単位を音声言語コーパスの分析結果に基づき定め，この単位に分割する．これにより，構文的曖昧性が軽減され，解析精度の改善が期待できる．漸進的な構文解析手法では，文の概念を排除し，効率的な構文解析手法で定めた単位を解析の単位とし，この単位ごとに構文解析を実行する．

なお，本研究では，日本語を構文解析の対象とする．日本語には，1) 語順が比較的自由である，2) 格要素の省略が可能である，という特徴があり，この特徴にあった文法を選択することが重要になる．自然言語の構文解析で用いられる文法には，文脈自由文法 (CFG) や主辞駆動句構造文法 (HPSG) [76]，木接合文法 (TAG) [30]，係り受け文法 [60] などが存在する．このうち，CFG や HPSG，TAG など，句構造規則を基本とする文法は，英語のような語順に対する制約が強い屈折言語を扱う場合には有効であるが，日本語に対する有効性は低い．これは，句構造規則によって日本語文の基本的構造を規定しようとする，あらゆる語順，あらゆる省略のパターンに対する規則を用意する必要があるためである．一方，係り受け文法は，依存文法とも呼ばれ，ある文節が他の文節に係る（依存する）という形式で文の構造を表現する．この係り受け文法に基づく構文解析，すなわち，係り受け解析は，文の解析時に各係り受け関係を独立に調べることができるため，日本語の特徴にあった構文解析であると考えられている [67]．そこで，本研究では，構文解析手法として係り受け解析を採用した．

## 1.4 本論文の内容

本論文ではまず第一に，大規模音声言語コーパスを用いた話し言葉の頑健な係り受け解析手法を提案する．本手法では，フィラーや言い淀み，倒置などの非文法的言語現象が頻出する対話文を対象とした．実際に，車内音声対話文に係り受け情報を付与したコーパスを分析した結果，従来の係り受け解析手法では，係り受けの非交差性，後方修飾性，係り先の唯一性の3つの制約が用いられてきたが，対話音声では，後方修飾性を満たさない倒置現象や係り先の唯一性を満たさない文節などを含む発話が頻出することが分かった．そこで本手法では，後方修飾性の制約及び係り先の唯一性に関する制約は統計情報を反映させつつ緩和した．また，本手法では，構築した構文構造付き音声対話コーパスから各文節間の係り受け確率を統計的に獲得し，それを用いて係り受け構造の尤度を計算する．これにより，非文法的な特徴をもつ発話の解析が可能になる．名古屋大学 CIAIR 車内音声対話コーパスに収録された対話文に対して係り受け解析実験を行った結果，本手法の有用性及び頑健性を確認した．

第二に，文の分割に基づく話し言葉の係り受け解析手法を提案する．本研究では，文の分割単位として節に着目した．節は「述語を中心としたまとまり」であり，節の内部で係り受けがまとまりやすいと考えられるためである．本手法では，節レベ

ルと文レベルの二段階で係り受け解析を実行する．まず，節境界解析により文を節に分割し，各節に対して係り受け解析を行うことにより，節内の係り受け関係を同定する．次に，節境界をまたぐ係り受け関係を定め，文全体の係り受け構造を作り上げる．これにより，話し言葉に出現する極端に長い文に対する効率的な解析の実現が期待できる．極端に長い文が頻出する独話データとして，NHKの解説番組「あすを読む」の書き起こしデータを用いて，係り受け解析実験を行い，本手法の有効性を評価した．その結果，本手法により解析精度を改善しつつ解析時間を大幅に短縮できることを確認した．

第三に，話者による音声入力に従って順次，解析を行う漸進的係り受け解析手法を提案する．本手法では，独話音声に対して，節が入力されるたびにその節の内部の係り受け構造を作り上げるとともに，すでに入力されている節の係り先を決定することを試みる．節の係り先となる文節の決定は，後続するいくつかの文節との係り受けの尤度を考慮した動的なタイミングで行う．これにより，独話の入力途中の段階で構造情報を随時出力する漸進的な解析が可能となる．NHKの解説番組「あすを読む」の書き起こしデータを用いた解析実験の結果，本手法により，文を解析単位とした係り受け解析手法と同等の解析精度を維持しつつ，解析の漸進性を実現できることを確認した．

## 1.5 本論文の構成

本論文の構成を図1.2に示す．まず，2章では，構文構造付き音声言語コーパスの構築について述べる．構築したコーパスは，構文解析手法を開発する上で，話し言葉の言語的特徴の分析や統計的解析手法における学習データとして使用する．

3章では，話し言葉に特有な言語現象に対して頑健な構文解析手法を提案する．本手法は，統計に基づく係り受け解析を行うことにより，音声対話文に頻出する倒置やフィラー，言い淀みなどの言語現象を網羅的かつ統一的に扱うことができるという特色をもつ．

4章では，1文の長さが極端に長くなることがある独話文に対して，効率的な構文解析手法を提案する．本手法は，日本語記述文法に従って文を節に近似的に分割し，解析を簡単化することによって，解析精度を維持しつつ解析の高速化を実現する．

5章では，話し言葉の漸進的な構文解析手法を提案する．本手法は，節に相当する単位ごとに解析を実行することにより，解析精度を同程度に維持しつつ，話し手の話速に追従できる程度の解析の漸進性，及び，解析速度を備えた構文解析を実現

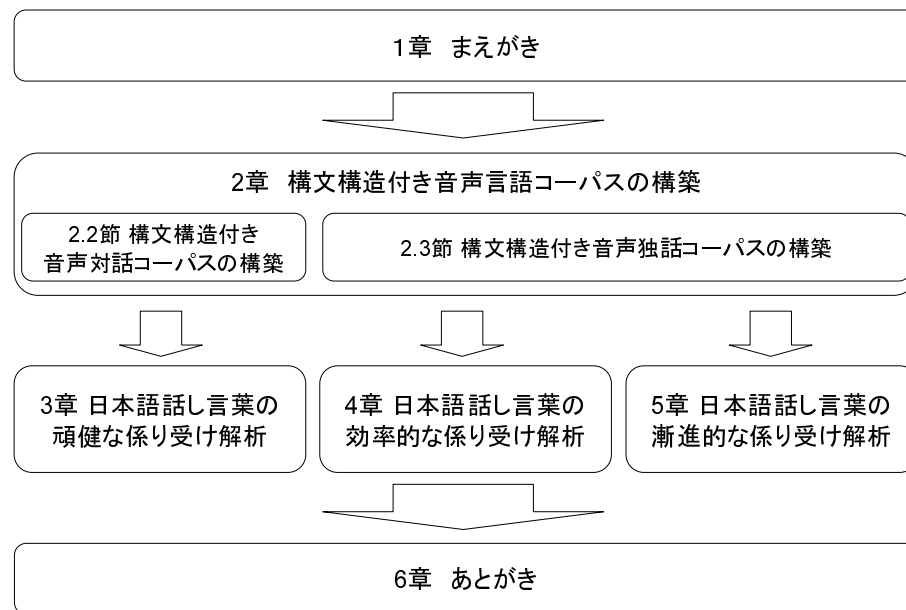


図 1.2: 本論文の構成

する。

最後に 6 章では、本論文のまとめと残された課題、将来の展望について述べる。



## 第2章 構文構造付き音声言語コーパスの構築

### 2.1 はじめに

構文構造が付けられた大規模テキストデータ（以下，構文構造付きコーパス）は，自然言語処理において，重要な役割を果たしている．実際，新聞等，書き言葉の大規模言語データに基づく構文構造付きコーパスが数多く構築されている．例えば，海外では，Penn Treebank[54] や NEGRA Treebank[87]，TIGER Treebank[7]，Prague Czech-English Dependency Treebank[9] などが構築されている．また，日本では，EDR コーパス [29]，京都テキストコーパス [49] などがあり，いずれも情報検索や要約，機械翻訳など広く利用されている．また，構文解析の分野では，これらの構文構造付きコーパスを統計情報として活用することにより，高精度な構文解析が実現されている [8, 10, 18, 46, 78, 92]．

一方，話し言葉の構文構造付きコーパスに目を向けると，Switchboard コーパス [61] や Verbmobil Treebank[97]，Spoken Dutch コーパス [95] などが存在する．しかし，日本語を対象とした話し言葉の構文構造付きコーパスはほとんどない．最近では，CSJ プロジェクト [53] で集められた話し言葉データの一部に構文構造が付与された例があるものの，書き言葉の構文構造付きコーパスに比べて，十分に整備されているとは言い難い．

そこで，本研究では，話し言葉の構文的特徴を明らかにするための分析データ，及び，統計的構文解析手法の学習データとして利用することを目的に，話し言葉の構文構造付きコーパスを構築した．話し言葉は，複数の話者が交替で話す「対話」と一人の話者のみが話す「独話」に分類できる．対話には，倒置やフィラー，言い直し，省略など書き言葉には現れない言語現象が頻出する．一方，独話には，対話文に比べ，1文の長さが長く文の構造が複雑であるといった特徴がある．本研究では，各々の特徴を考慮した構文構造の付与基準を作成し，対話と独話の構文構造付きコーパスを構築した．

以下では，まず，2.2 節で構文構造付き音声対話コーパスの構築について述べ，次

に、2.3 節で構文構造付き音声独話コーパスの構築について述べる．最後に、2.4 節で本章をまとめる．

## 2.2 構文構造付き音声対話コーパスの構築

本節では、構文構造付き音声対話コーパスの構築について述べる．対話には、倒置やフィラー、言い直し、省略など書き言葉には現れない言語現象が頻出する．本研究では、このような話し言葉に特有な言語現象を含む文に対しても構文構造を付与するために、独自の構文構造付与基準を設け、コーパスを構築した．構築したコーパスを統計情報として活用することにより、頑健で精度の高い音声言語係り受け解析の実現が期待できる．

本研究では、音声対話コーパスとして、CIAIR 車内音声対話コーパス [36, 37, 38, 42] を使用し、構文構造を付与する作業を実施した．

### 2.2.1 CIAIR 車内音声対話コーパス

名古屋大学統合音響情報研究拠点（CIAIR）では、1999 年度から 2003 年度まで車内音声対話の収録を実施してきた [36, 37, 38, 42]．図 2.1 に示す専用の収録車を開発し、ドライバ約 800 人に対して各 60 分の対話音声を収録しており、そのデータ量は約 400GB にのぼる．対話は、ドライバとナビゲータとの間で遂行され、タスクとして店検索や道案内などを設定している．収集した音声データの書き起こしは、日本語話し言葉コーパス (CSJ) [53] に準拠しており、作業は人手により行っている．

書き起こしデータの例を図 2.2 に示す．対話音声を 200 ミリ秒以上のポーズで分割し、各々を発話単位として、発話単位ごとに発話音声を書き起こしている．図 2.2 は、4 つの発話単位の書き起こしデータを抜粋したものである．各発話の先頭行には、その ID と開始時間、終了時間、話者の性別 (男性/女性)、話者役割 (ドライバ/ナビゲータ)、対話タスク (道案内/情報検索など)、雑音状況に関する情報を付与している．例えば、1 行目の

0009 - 00:41:960-00:50:228 F:O:I:R:

では、「0009」が発話単位の ID を、「00:41:960」が開始時間を、「00:50:228」が終了時間を、「F:O:I:R:」が「話者が女性 (F) で、話者役割がナビゲータ (O) で、対話タスクが情報検索 (I) で、走行音 (R) の雑音が行っていること」をそれぞれ表している．また、各発話単位の先頭行以降は、発話音声の書き起こしとその読みが文節ごとに改行さ

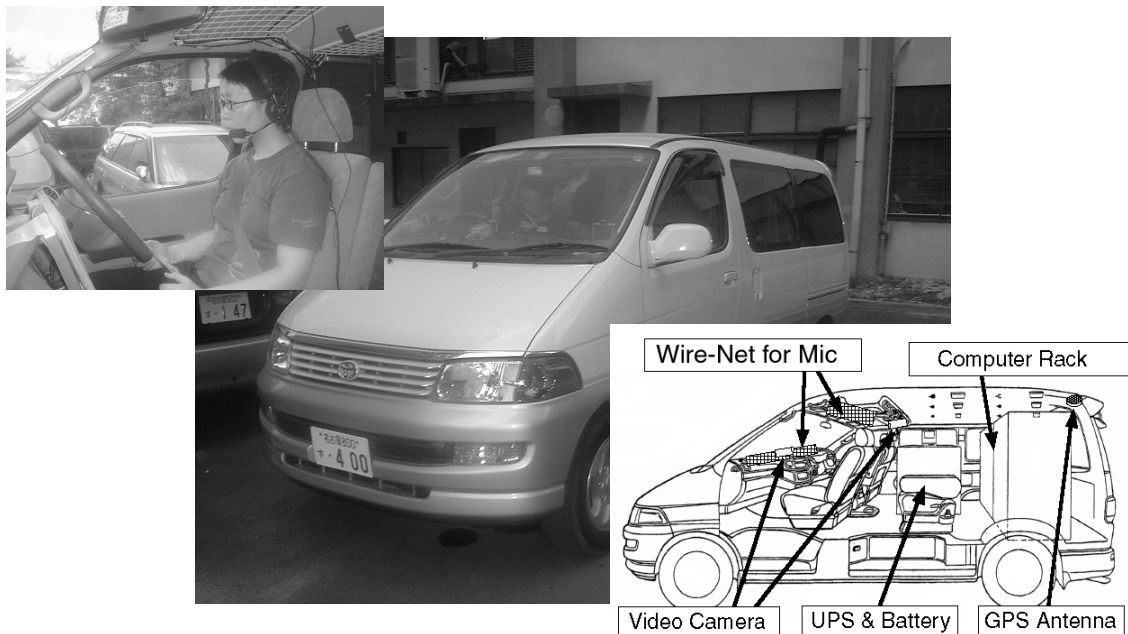


図 2.1: データ収集車

0009 - 00:41:960-00:50:228 F:O:I:R:  
 この先 & コノサキ  
 三百メートルほど先に & サンビャクメートルホドサキニ  
 左手に & ヒダリテニ  
 サンクス & サンクス  
 五百メートルほど先に & ゴビャクメートルホドサキニ  
 セブンイレブンが & セブンイレブンガ  
 ございます<SB> & ゴザイマス<SB>  
 0010 - 00:51:888-00:52:100 F:D:P:C:  
 (F ん) & (F ン)  
 0011 - 00:54:633-00:55:394 F:D:P:R:  
 (? 次左かな) & (? ツギヒダリカナ)  
 0012 - 01:02:750-01:03:993 F:O:I:R:  
 どちらに & ドチラニ  
 なさいますか<SB> & ナサイマスカ<SB>

図 2.2: CIAIR 車内音声対話コーパスの書き起こしデータ

れて記されている．さらに，データの言語学的分析として，文末やフィラー，言い淀み，言い誤りなどにタグを付与している．図 2.2 の例では，8, 15 行目の「<SB>」が文末を，10 行目の「(F ん)」の「F」がフィラーを表すタグである．なお，12 行目の「(? 次左かな)」の「?」は，聞き取りに不安が残る場合に記すタグである．

CIAIR 車内音声対話コーパスの 195 対話の全ドライバ発話を分析した結果，1 発話あたりのフィラーや言い淀み，言い間違いの平均出現回数は，それぞれ 0.34，0.07，0.04 であった．これは，人間と人間の通常行われる会話と同程度の流暢性を示しており [66, 68]，このコーパス中の収録音声が自然な発話であることを意味している．

### 2.2.2 CIAIR 構文構造付き音声対話コーパス

CIAIR 車内音声対話コーパスのドライバの発話に対して，以下の情報を付与することにより，構文構造付き音声対話コーパスを作成した．

- 形態素情報
  - － 形態素区切り
  - － 形態素の読み，基本形，品詞，活用型，活用形
- 係り受け情報
  - － 文節区切り
  - － 文節間の係り受け

ここで，品詞体系は形態素解析器茶筌 [59] の IPA 品詞体系 [3] に，文節区切りは日本語話し言葉コーパスの作成基準 [53] に，係り受け文法は京都テキストコーパスの作成基準 [49] にそれぞれ準拠した．ただし，話し言葉特有の言語現象については，以下の作成基準を新たに設けた．

- フィラーや言い淀みは，係り先が存在しない．すなわち，単独で係り受け構造を形成する．図 2.3 に例を示す．この例では，フィラー「う」「あっ」や言い淀み「おはな」の係り先はない，としている．
- 受け文節が省略された文節の係り先は存在しない．図 2.4 に例を示す．この例では「お願いします」が省略されているため「一つで」と「コーラで」の係り先はない，としている．また，図 2.5 では「予約は」と発話した後，言い直しが起きているために，この「予約は」の係り先が省略されたと判断している．

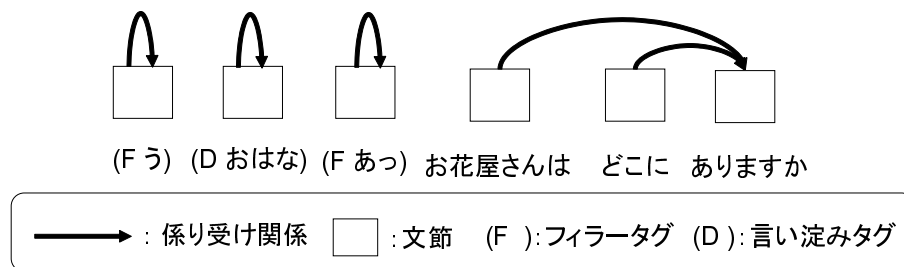


図 2.3: フィラー・言い淀みを含む文の係り受け構造

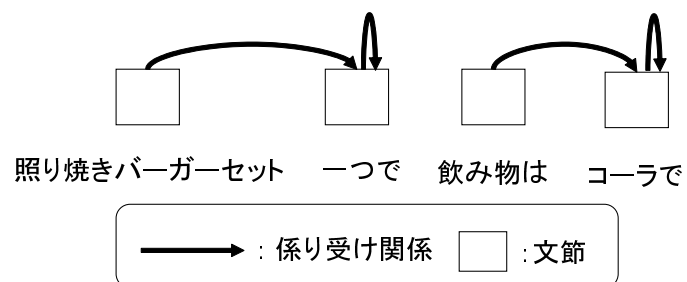


図 2.4: 受け文節が省略された文節を含む文の係り受け構造

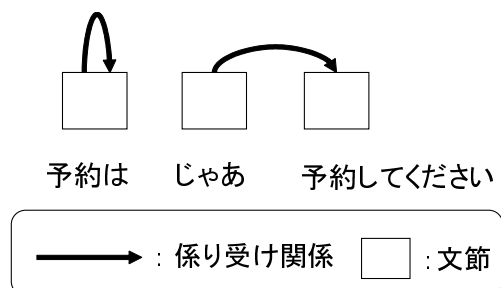


図 2.5: 言い直しを含む文の係り受け構造

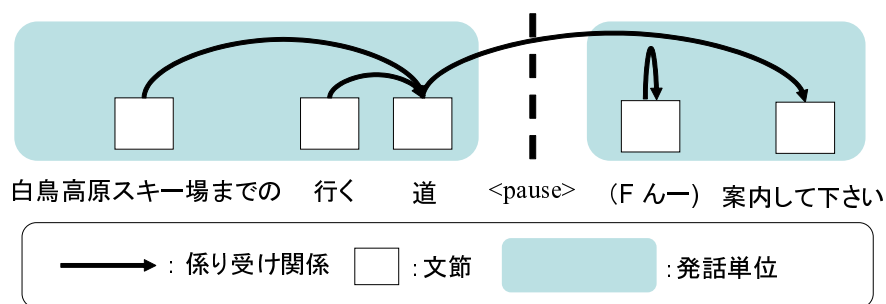


図 2.6: 発話単位を越える係り受けを含む文の係り受け構造

表 2.1: 話し言葉に特有な言い回し表現

話し言葉に特有な表現	品詞
こっ	代名詞
そっこ	代名詞
こっち	代名詞
すげえ	副詞
あんま	副詞
っていうか	接続詞
ほいじゃあ	接続詞
やば	感動詞
そっか	感動詞
よっし	感動詞

- 対話ターンを係り受けの単位とする．すなわち，発話単位を越える係り受けも認めるが，対話ターン間にまたがった係り受けは認めない．発話単位を越える係り受けを含む文の例を図 2.6 に示す．
- 話し言葉特有の言い回し表現（「こっから」、「っていうか」、「ほいじゃあ」など）については，新たな辞書項目を設けて，形態素ごとに品詞を定める．話し言葉特有の言い回し表現として，新たに設けた辞書項目は 124 個あり，その一部を表 2.1 に示す．
- 形式名詞となりえる名詞は，一律に内容語とみなして，その直前に文節境界を入れる<sup>1</sup>．ただし，形式名詞「の」は例外として，文節を切らない．以下に例を示す．なお，「/」は文節境界を意味する．

<sup>1</sup>日本語話し言葉コーパスの作成基準 [53] では，「もとの意味がうすれている場合は，形式名詞とみなし，文節を切らない」としている．しかし，もとの意味がうすれているかどうかを判断することは難しく，判断が揺れてしまうので，本研究では，一律に文節を分けることとした．

- 信州の/ほうへ/行きたい
- 今の/ところ/このまま/やる
- 食べる/ところを/探したい
- ユニクロみたいな/とこで
- 食べた/ことないから
- ケーキの/あるのは/どちらですか

CIAIR 構文構造付き音声対話コーパスの例を図 2.7 に示す．この例は，発話

今日朝パン食べてお昼はおそばを食べたんですよ

に対して係り受け構造を付与したデータである．この図に示すように，係り受け構造は係り受け関係の列で表し，各係り受け関係は，係り文節と受け文節からなる．また，各文節には，文節番号とその文節を構成する形態素を列挙する．図 2.7 が表す係り受け構造を図 2.8 に図示する．

### 2.2.3 CIAIR 構文構造付き音声対話コーパスの構築手順

CIAIR 構文構造付き音声対話コーパスの構築では，音声対話の書き起こしデータを自動的に解析し，その結果を人手により修正することにより，構文構造を付与した．また，統計に基づく構文解析器は，学習データが増加するほど解析の精度が高くなることに着目し，係り受け解析の結果に，人手で修正を加えることにより作成した係り受けデータを，別のテキストデータに構文構造を付与するための統計情報として利用するという増殖的なコーパス構築手法を採用した．構築対象のコーパスデータの全てに係り受け解析を施し，その後一括して人手で修正を行うという方法に比べ，修正に要する労力の軽減が期待できる．本コーパスの構築は，統計に基づく係り受け解析の結果を人手で修正することにより行った．効率的な構築作業を実現するために，修正の労力をできる限り抑えることが重要であるが，修正量の程度は，使用する係り受け解析の性能に大きく依存する．精度の高い係り受け解析を実行することが望まれ，一般にそれは，統計情報を獲得するために用いる言語データの規模を増やすことによって実現可能となる．

そこで本論文では，係り受けデータ付与の対象となるコーパスの全てに対して一括して係り受け解析を行うのではなく，コーパスを複数に分割し，逐次的に構築作業を遂行することによって，増殖的な構文構造付きコーパスの作成を実現する．す

```

((1 ((今日 キョウ きょう 名詞 副詞可能 なし なし)))
-> (7 ((食べ タベ 食べる 動詞 自立 一段 連用形)
      (た タた 助動詞 なし 特殊・タ 基本形)
      (ん ンん 名詞 非自立 なし なし)
      (です デスです 助動詞 なし 特殊・デス 基本形)
      (よ ヨよ 助詞 終助詞 なし なし))))

((2 ((朝 アサ 朝 名詞 副詞可能 なし なし)))
-> (4 ((食べ タベ 食べる 動詞 自立 一段 連用形)
      (て テて 助詞 接続助詞 なし なし))))

((3 ((パン パン パン 名詞 一般 なし なし)))
-> (4 ((食べ タベ 食べる 動詞 自立 一段 連用形)
      (て テて 助詞 接続助詞 なし なし))))

((4 ((食べ タベ 食べる 動詞 自立 一段 連用形)
      (て テて 助詞 接続助詞 なし なし)))
-> (7 ((食べ タベ 食べる 動詞 自立 一段 連用形)
      (た タた 助動詞 なし 特殊・タ 基本形)
      (ん ンん 名詞 非自立 なし なし)
      (です デスです 助動詞 なし 特殊・デス 基本形)
      (よ ヨよ 助詞 終助詞 なし なし))))

((5 ((お昼 オヒル お昼 名詞 副詞可能 なし なし)
      (は ハは 助詞 係助詞 なし なし)))
-> (7 ((食べ タベ 食べる 動詞 自立 一段 連用形)
      (た タた 助動詞 なし 特殊・タ 基本形)
      (ん ンん 名詞 非自立 なし なし)
      (です デスです 助動詞 なし 特殊・デス 基本形)
      (よ ヨよ 助詞 終助詞 なし なし))))

((6 ((おオ お 接頭詞 名詞接続 なし なし)
      (そば ソバそば 名詞 一般 なし なし)
      (をヲ を 助詞 格助詞 なし なし)))
-> (7 ((食べ タベ 食べる 動詞 自立 一段 連用形)
      (た タた 助動詞 なし 特殊・タ 基本形)
      (ん ンん 名詞 非自立 なし なし)
      (です デスです 助動詞 なし 特殊・デス 基本形)
      (よ ヨよ 助詞 終助詞 なし なし))))

((7 ((食べ タベ 食べる 動詞 自立 一段 連用形)
      (た タた 助動詞 なし 特殊・タ 基本形)
      (ん ンん 名詞 非自立 なし なし)
      (です デスです 助動詞 なし 特殊・デス 基本形)
      (よ ヨよ 助詞 終助詞 なし なし)))
-> (NO (なし)))

```

図 2.7: CIAIR 構文構造付き音声対話コーパスの例

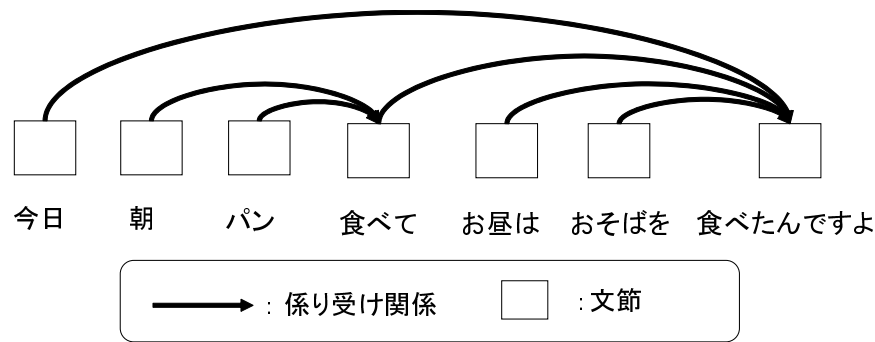


図 2.8: 「今日朝パン食べてお昼はおそばを食べたんですよ」の係り受け構造

なわち，係り受け解析の結果に人手で修正を加えて作成したデータを，すでに構築されている構文構造付きコーパスに追加し，別のテキストコーパスを解析するときの統計情報として利用する．統計情報は修正したデータから自動的に獲得できるので，解析結果の修正以外に人手による作業を必要としない．コーパス構築の過程で，係り受け解析の精度が漸増的に向上し，コーパス修正に要する労力を軽減することが可能となる．

### コーパス構築の流れ

CIAIR 構文構造付き音声対話コーパス構築の流れを図 2.9 に示す．複数セットに分割された CIAIR 車内音声対話コーパスから 1 つを取り出し，以下の作業を順に行う．

- (1) 音声データの書き起こしテキストを，形態素解析システム茶筌 [59] により形態素に区切り，形態素分析を与える．
- (2) 形態素情報を人手により修正する．話し言葉特有の言い回しや固有名詞が新たに出現した場合は，茶筌の辞書に随時追加する．
- (3) 統計情報を利用した文節まとめあげ手法を用いて，文節ごとに区切る．
- (4) 文節区切りの修正を人手により行う．
- (5) 統計的な係り受け解析により，係り受け構造を付与する．
- (6) 係り受け情報を人手により修正する．

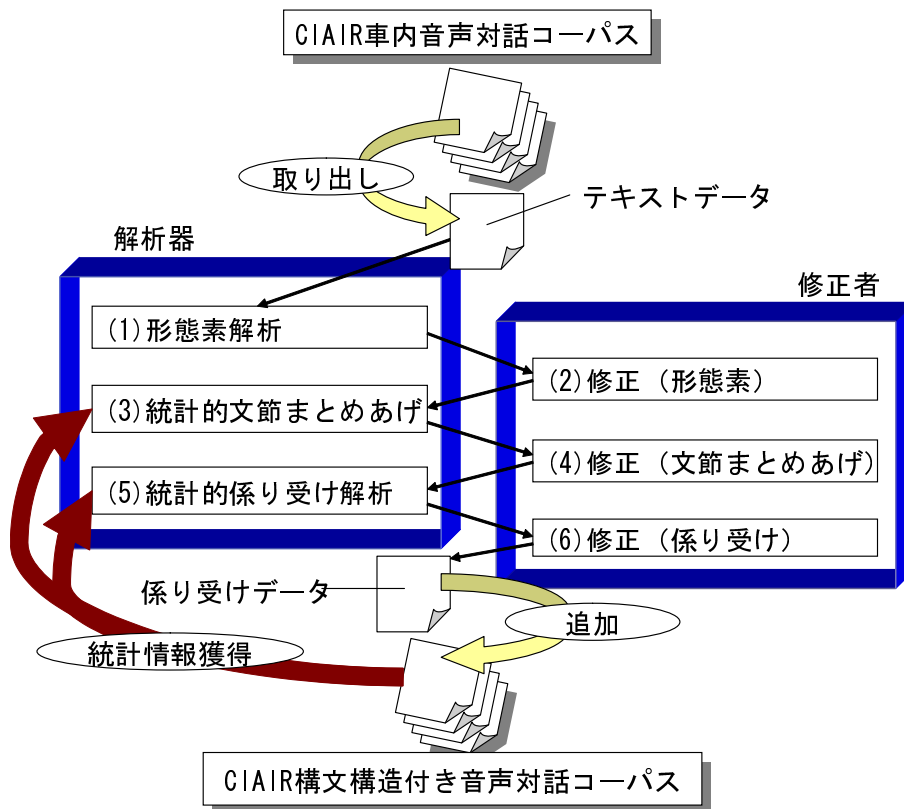


図 2.9: CIAIR 構文構造付き音声対話コーパス構築の流れ

この (6) の結果を CIAIR 構文構造付き音声対話コーパスに追加する。

以下では、本構築手法で用いた統計的な文節まとめあげ手法と人手修正用インタフェースについて述べる。なお、係り受け解析手法は、3.3 節で説明する統計的係り受け解析手法を用いた。

### 統計的文節まとめあげ

文節まとめあげとは、形態素列  $M(= m_1 \cdots m_n)$  の各形態素間に文節区切りを挿入するか否かを判定する問題である。例えば、「今日朝パン食べてお昼はおそばを食べたんですよ」に対する文節まとめあげの効果は次の通りである。

今日/朝/パン/食べ<sub>□</sub>て/お昼<sub>□</sub>は/お<sub>□</sub>そば<sub>□</sub>を/ 食べ<sub>□</sub>た<sub>□</sub>ん<sub>□</sub>です<sub>□</sub>よ

ここで、「/」は文節境界である形態素境界を、「<sub>□</sub>」は文節境界でない形態素境界を示す。

本手法では，文節まとめあげ済みコーパスから，文節区切りに関する統計情報を獲得し，それを用いて文の先頭から順に判定することを試みる．隣接する形態素  $m_i$  と  $m_{i+1}$  の間の文節区切りに注目するとき，統計情報として各形態素の以下に挙げる属性を利用する．

- 形態素の基本形  $h_i, h_{i+1}$
- 形態素の品詞  $t_i, t_{i+1}$
- 形態素の活用形または品詞細分類  $s_i, s_{i+1}$

形態素  $m_i$  と  $m_{i+1}$  の境界で文節区切りとなる確率，すなわち，同一文節に含まれない確率を以下のように計算する．

$$P(m_i/m_{i+1}|m_i, m_{i+1}) = \frac{C(m_i/m_{i+1}, h_i, h_{i+1}, t_i, t_{i+1}, s_i, s_{i+1})}{C(h_i, h_{i+1}, t_i, t_{i+1}, s_i, s_{i+1})}$$

ここで， $m_i/m_{i+1}$  は，連続する形態素  $m_i, m_{i+1}$  の間に文節の区切りが存在することを意味する． $C$  は共起頻度関数である．ここでは， $P(m_i/m_{i+1}|m_i, m_{i+1}) \geq 0.5$  ならば，そこで文節を区切る．

### 解析結果修正用インタフェース

人手による修正の負担を軽減するため，GUI ベースの修正インタフェースを作成した．インタフェース画面を図 2.10 に示す．これを用いて全ての修正を行うことができる．

インタフェースの左側の画面で形態素情報と文節まとめあげの解析結果を修正する．1 つの形態素が 1 行で表現され，各行には，形態素の各種情報を表示したボタンが並んでいる．品詞，品詞細分類，活用型，活用形はそれぞれメニューボタンで表示され，メニューバーから適切なものを選択することにより修正できる．

また，同一の文節に含まれる形態素をすべて同一の色で示すことにより，文節のまとまりを表現する．文節区切りの修正は，画面の左端にある操作ボタンをクリックし，その行の色を変更することにより行う．変更後に係り受け修正ボタンを押すと，係り受け解析プログラムが動作し，修正内容が係り受け構造に反映される．

インタフェースの右側の画面で係り受け構造を修正する．この画面では，文節とその係り先文節を示すとともに，係り受け関係を視覚的に表示している．係り先の文節番号が「NO」となっている場合は，その文節が係り先をもたないことを意味する．



図 2.10: 解析結果修正用インタフェース

## 2.2.4 評価実験

係り受け関係を自動的に与えたときのコーパス全体の修正に要する労力に着目し、増殖的な構築手法の作業効率に関する評価を行った。

### 実験の概要

実験用データとして、CIAIR 車内音声対話コーパスの 221 対話を使用した。データの規模は、45,053 文節からなる 10,995 ターンである (平均ターン長は 4.1 文節)。このうち、10 ターンを基礎学習データとし、残りの 10,985 ターンをテストデータとした。増殖的構築の効果を調べるために、テストデータを 110 個のセット (1 セットは約 100 ターン) に等分割した。便宜上、各セットに 1 から 110 までの番号を与えた。同一のデータに対し、以下の 2 つの手法でコーパスを構築した。

- 増殖的構築手法 2.2.3 節で述べた手順に基づいて係り受けデータを付与した。セット 1 から 110 まで順に構築作業を実施した。各セットへの係り受け付与が完了したデータは、それ以降の解析の学習データに追加した。

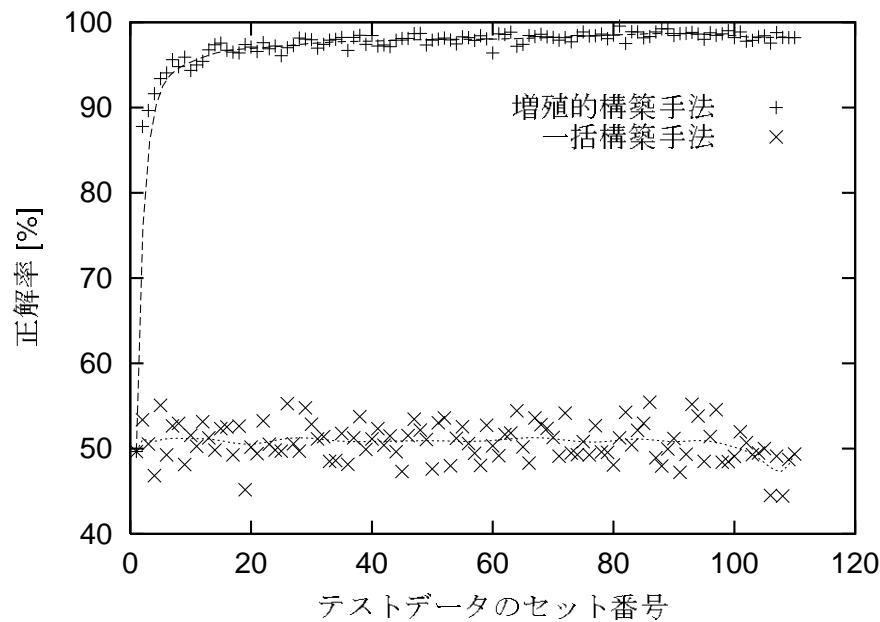


図 2.11: 統計的文節まとめあげの実験結果

- 一括構築手法 全てのテストデータに対し、基礎学習データのみを用いて係り受け解析を行い、その後、一括して人手により修正した。ただし、文節まとめあげと係り受け解析は増殖的構築手法と同一である。

評価は、文節まとめあげと係り受け解析の精度を両構築手法の間で比較することにより行った。

### 実験の結果

両手法による文節まとめあげの結果を図 2.11 に示す。増殖的構築手法の方が 46.4% 高い精度でまとめあげを行うことができた。テストデータに含まれる形態素間の境界は計 96,675 個所あり、そのうち正しく判定できた数は、一括構築手法では 48,969 個所、増殖的構築手法では 93,880 個であった。すなわち、増殖的構築手法を採用することにより、修正個所が 44,911 個削減されており、これは、修正量が一括構築手法の 6% 程度で済むことを意味している。

係り受け解析結果の比較を図 2.12 に示す。増殖的構築手法の正解率は 88.4%、また、一括構築手法の正解率は 61.7% であり、文節まとめあげと同様、増殖的構築手法のほうが高い精度で係り受けデータの付与が行えた。テストデータに含まれる係り受け数 45,012 個のうち、正しくデータを付与できた数は、一括構築手法が 27,767

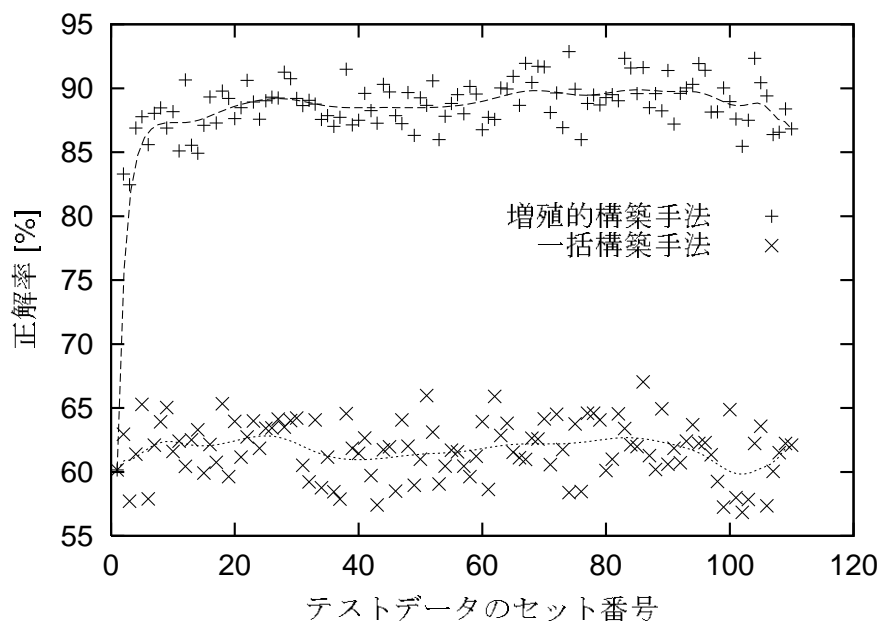


図 2.12: 統計的係り受け解析の実験結果

表 2.2: CIAIR 構文構造付き音声対話コーパスの規模

対話	221
ターン	10,995
発話単位	13,756
文節	45,053
形態素	85,870

個で，増殖的構築手法は 39,804 個であった．すなわち，増殖的構築手法を採用することにより，修正箇所が 12,037 個削減されており，これは，一括構築手法の修正量を約 7 割削減できたことを意味する．

これらのことから，コーパス全体の修正に要する労力を軽減するために増殖的構築手法が有効であることが確認できた．

### 2.2.5 CIAIR 構文構造付き音声対話コーパスの規模と特徴

本研究で構築した構文構造付き音声対話コーパスの規模を表 2.2 に示す．本コーパスでは，話し言葉に特有な言語現象を含む文に対しても構文構造を付与するために，独自の構文構造付与基準を設けている．日本語話し言葉に対して構文構造が付与されている他のコーパスとして，CSJ コーパス [53] がある．しかし，CSJ コーパ

スは講演などの独話に対してのみ構文構造が付与されており，日本語対話文を対象としている本コーパスの意義は大きい．

## 2.3 構文構造付き音声独話コーパスの構築

本節では，構文構造付き音声独話コーパスの構築について述べる．独話の場合には，対話と比べて，非常に長い文が頻出するという傾向がある．これは，対話の場合，その多くが単文であるのに対して，独話は複文や重文が頻出するためである．一般に，文が長くなるとその構造も複雑になるため，文のみを言語的分析や構文解析の絶対的な基本単位として用いることは必ずしも適切とはいえない．

そこで，本コーパスの構築では，係り受け構造情報のほかに，節境界情報も付与した．節は「述語を中心としたまとまり」と定義され，単文に相当し，構文的かつ意味的にまとまった単位であると考えられる[56]．なお，CSJ プロジェクト[53]においても，音声独話に対して係り受け構造情報，及び，節境界情報が付与されたコーパスが構築されているが，係り受け構造は節内部に付与されているのみであるのに対して，本コーパスでは節間の関係も付与されている．

NHK の解説番組「あすを読む」110 番組分の書き起こしデータに対し，上述した独自の基準に基づき係り受け情報，及び，節境界情報を付与することにより，構文構造付き音声独話コーパスを構築した．

### 2.3.1 「あすを読む」書き起こしコーパス

NHK で放映されていた 10 分間のニュース解説番組「あすを読む」<sup>2</sup>327 番組分の発話を書き起こしたコーパスが日本放送協会（NHK）と国際電気通信基礎技術研究所（ATR）の共同研究の中で構築されている．図 2.13 に「あすを読む」書き起こしコーパスの例を示す．独話音声を 200 ミリ秒以上のポーズとフィラーで分割し，各々を発話単位としてその開始時間，終了時間を記録している．図の 1 列目はポーズ，または，フィラー，発話単位の開始時間，3 列目は終了時間，4 列目はポーズ記号，または，フィラーや発話単位の書き起こしをそれぞれ意味している．なお，フィラーは大括弧（〔 〕）で囲んでおり，時間はミリ秒単位で記載している．また，文境界を示す句点も挿入している．例えば，1 行目で 0 ミリ秒から 2,940 ミリ秒までポーズがあり，2 行目で 2,940 ミリ秒から 3,049 ミリ秒まで「え」というフィラーが発話され，

<sup>2</sup> 「あすを読む」は，各回ごとに 1 つの時事問題に関して，異なる解説者が 10 分の解説を行う番組であり，1996 年 4 月から 2006 年 3 月まで放送された．

0	-	2940	PAU
2940	-	3049	[え]
3049	-	3681	今晚は。
3681	-	3981	PAU
3981	-	4097	[え]
4097	-	5716	最高裁判所は今日<kyou>
5716	-	5884	[ま]
5884	-	6937	検察側が
6937	-	7105	[え]
7105	-	9373	死刑を求めて上告をしていました
9373	-	9695	PAU
9695	-	9841	[え]
9841	-	11520	強盗殺人事件について
11520	-	11966	PAU
11966	-	12063	[ま]
12063	-	16011	二審と同じように無期懲役の判決を言い渡しております。
16011	-	17256	PAU
17256	-	17320	[え]
17320	-	18713	裁判で問われましたのは、
(以下省略)			

図 2.13: 「あすを読む」書き起こしコーパスの例

3 行目で 3,049 ミリ秒から 3,661 ミリ秒まで「今晚は。」と発話されたことを示している。

### 2.3.2 「あすを読む」構文構造付き音声独話コーパス

上述した「あすを読む」書き起こしデータ全 327 番組のうち 110 番組分に対して、以下の情報を付与することにより構文構造付き音声独話コーパスを構築した<sup>3</sup>。

- 形態素情報
  - － 形態素区切り
  - － 形態素の読み，原形，品詞，活用型，活用形
- 文節境界情報
- 節境界情報
  - － 節境界区切り

<sup>3</sup>著者は，ATR の研修研究員として「あすを読む」書き起こしコーパスを使用した。

表 2.3: 節境界の大分類と小分類

大分類	小分類
副詞節 (102)	条件・譲歩節 (23), 原因・理由節 (8), 時間節 (21), 様態節 (2), 副詞節その他 (38)
補足節 (10)	補足節 (2), 引用節 (5), 間接疑問節 (3)
連体節 (15)	連体節 (15)
並列節 (12)	並列節 (12)
その他 (8)	文末 (1), 主題八 (1), 談話標識 (1), 体言止 (1), 間投句 (1), 感動詞 (1), 従属文 (1), 従属文その他 (1)

#### － 節境界の種類

##### ● 係り受け情報

#### － 文節間の係り受け

#### － 係り受けの種類

なお，形態素情報は茶筌 [59] の IPA 品詞体系 [3] に，文節境界情報は CSJ コーパス [53] の作成基準に，節境界情報は丸山らの基準 [55] にそれぞれ準拠して付与した．ただし，話し言葉特有の現象については，2.2.2 節と同様の作成基準を新たに設けた．以下では，節境界情報の詳細と係り受け構造情報の付与基準について述べる．

#### 節境界情報

上述したように，節境界情報の付与基準は，丸山ら [55] の基準に準拠している．丸山らの文献 [55] を引用しつつ，節境界情報の詳細について述べる．

丸山らは表 2.3 に示す 147 種類の節境界を定義している [55]．なお，括弧内の数字は，各クラスに登録された節境界の種類の数を示す．節とは「述語を中心としたまとまり」と定義される [56]．節は，文末に配置される「主節」とそれ以外の「従属節」に分けることができる．さらに，従属節は，形態的および機能的な観点から，副詞節，補足節，連体節，並列節の 4 種類に分類される [56]．丸山らは，この 4 種類の従属節を大分類として利用している [55]．なお，大分類「その他」の中の「文末」が「主節」に相当する．また，大分類が「その他」の節境界のうち「文末」以外の要素は，厳密には節境界ではないが，構文的に切れ目になると考えられるため，節境界とみなしている [55]．

例えば、大分類「副詞節」小分類「条件・譲歩節」の中には、「条件節レバ」という節境界名が存在し、「～すれば」という節境界の表現に対応している。このように、それぞれの節境界名は、その形態的特徴、および、文法的な特性の違いを考慮して、「理由節カラ」「連体節ヨウナ」「並列節ケレドモ」などのように表現される [55]。

### 係り受け付与基準

係り受け情報の付与基準は、原則として京都テキストコーパス基準 [49] に準拠しているが、京都テキストコーパス基準が遵守している構文的制約のうち、「係り先の唯一性」を緩和している。これは、日本語の係り受け構造では、ある文節が係り先を1つ以上持つと考えられる次のような場合が存在するためである。

- 主語（または、副詞、接続詞など）と、連体節・主節の述語との関係  
「今朝東京へ行く途中で考えた」に対する係り受け構造を図 2.14 に示す。「今朝」が文全体での係り受け構造としては「考えた」に係るが、「今朝」が「行く」を修飾していると考えても必ずしも誤りであるというわけではない。
- 主語（または、副詞、接続詞など）と、テ節・主節の述語との関係  
「小泉さんが二倍近くの差をつけて圧勝した」に対する係り受け構造を図 2.15 に示す。「小泉さんが」が文全体での係り受け構造としては「圧勝した」に係るが、「つけて」の主語は「小泉さんが」であり、そのような係り受け関係が必ずしも誤りであるというわけではない。

このような場合でも、京都テキストコーパス基準では、係り先として必ずしも誤りでないと考えられる複数の文節から、ただ1つの文節が選択される。

そこで、本コーパスでは、活用形や付属語の文法的働きに従った係り受け関係であり、かつ、意味的に必ずしも誤りでないならば、ある文節が係り先を複数持つことを許した。なお、比較評価のため、京都テキストコーパス基準における係り先とその他の必ずしも誤りでない係り先（以下、意味的係り受け関係）を区別してアノテートしている。

図 2.16 に構築した構文構造付き音声独話コーパスの例を示す。この例は、独話文

最高裁判所は今日検察側が死刑を求めて上告をしておりました（強盗殺人事件について二審と同じように無期懲役の判決を言い渡しております）

に対して係り受け構造を付与したデータである。各文の先頭には、その文の ID や修正日時、修正者名を記載している。この例では、ID が「230-1-2」で、2005 年 6 月

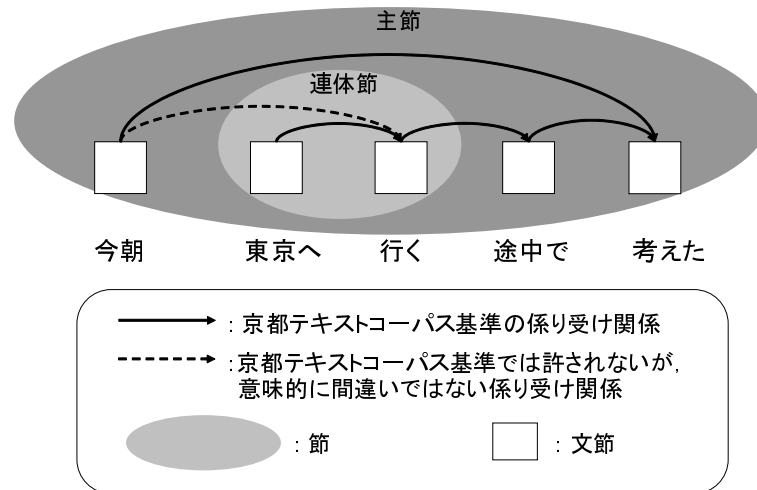


図 2.14: 「今朝東京へ行く途中で考えた」に対する係り受け構造

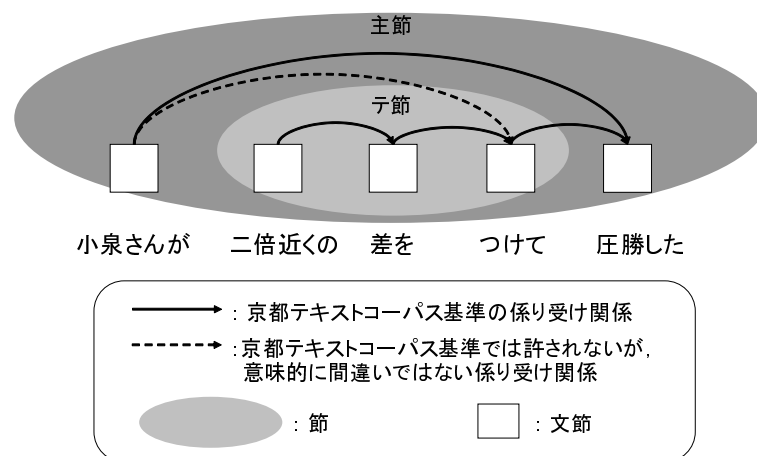


図 2.15: 「小泉さんが二倍近くの差をつけて圧勝した」に対する係り受け構造

ID	230-1-2	2005/6/15	14:18:6	大野誠寛
1	13D	/主題ハ/	LAST	
		最高裁判所	サイコウサイバンショ	最高裁判所 名詞-固有名詞-組織 なし なし
		は	ハ	は 助詞-係助詞 なし なし
2	13D	/テ節/	NONLAST	
		今日	キョウ	今日 名詞-副詞可能 なし なし
3	7D 5S	/テ節/	NONLAST	
		検察	ケンサツ	検察 名詞-サ変接続 なし なし
		側	ガワ	側 名詞-接尾-一般 なし なし
		が	ガ	が 助詞-格助詞-一般 なし なし
4	5D	/テ節/	NONLAST	
		死刑	シケイ	死刑 名詞-一般 なし なし
		を	ヲ	を 助詞-格助詞-一般 なし なし
5	7D	/テ節/	LAST	
		求め	モトメ	求める 動詞-自立 一段 連用形
		て	テ	て 助詞-接続助詞 なし なし
6	7D	/連体節/	NONLAST	
		上告	ジョウコク	上告 名詞-サ変接続 なし なし
		を	ヲ	を 助詞-格助詞-一般 なし なし
7	8D	/連体節/	LAST	
		し	シ	する 動詞-自立 サ変・スル 連用形
		て	テ	て 助詞-接続助詞 なし なし
		おり	オリ	おる 動詞-非自立 五段・ラ行 連用形
		まし	マシ	ます 助動詞 特殊・マス 連用形
		た	タ	た 助動詞 特殊・タ 基本形
		(以下省略)		

図 2.16: 「あすを読む」構文構造付き音声独話コーパスの例

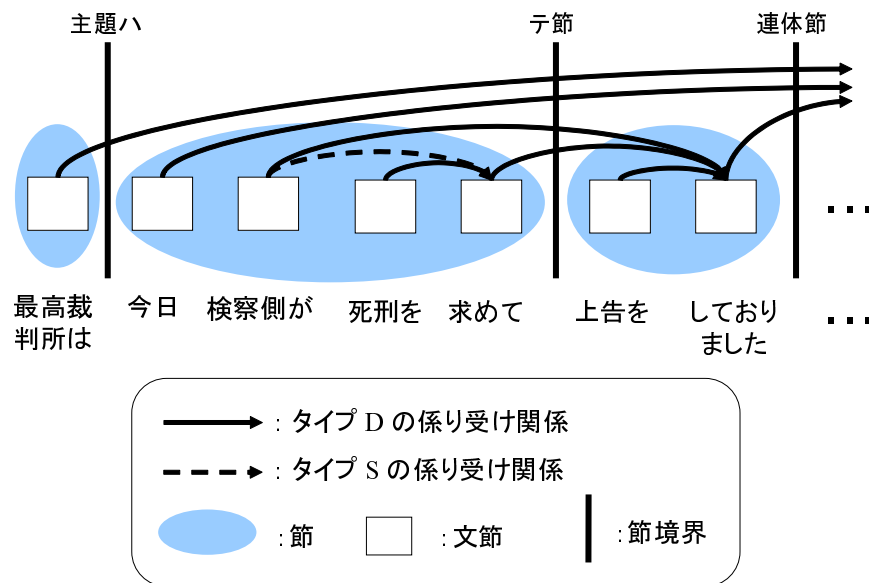


図 2.17: 「最高裁判所は今日検察側が死刑を求めて上告をしております（強盗殺人事件について二審と同じように無期懲役の判決を言い渡しております）」の係り受け構造

15 日 14 時 18 分 6 秒に大野誠寛によって修正されたことが 1 行目に記されている。次に、各文は、その文を構成する文節を列挙する形で記される。各文節の先頭には、文節番号とその受け文節の文節番号（場合によっては複数）、その文節が所属する節の種類、その文節が節末であるか否かの情報を記載している。例えば、3 行目の「1 13D /主題八/ LAST」は、1 番目の文節で、その係り先が 13 番目の文節で、その係り受けのタイプが「通常 (D)」であり、この文節の直後に節境界「主題八」があることを示している。なお、受け文節の番号の横のアルファベットは、係り受けのタイプを示しており、「S」は意味的係り受け関係（京都テキストコーパス基準では受け文節とならないが、必ずしも誤りでない係り先）を表している。図 2.16 では、3 番目の文節が、7 番目の文節だけでなく、意味的係り受け関係として 5 番目の文節に係っていることが記されている。図 2.17 に図 2.16 が表す係り受け構造を図示する。この図では、点線矢印により意味的係り受け関係を示している。

### 2.3.3 「あすを読む」構文構造付き音声独話コーパスの構築手順

「あすを読む」構文構造付き音声独話コーパスの構築では、形態素は茶筌 [59]、文節境界と係り受けは南瓜 [46]、節境界は CBAP [55] によってそれぞれ自動解析し、その結果を手で修正した。人手修正の負担を軽減するため、GUI ベースの修正イン

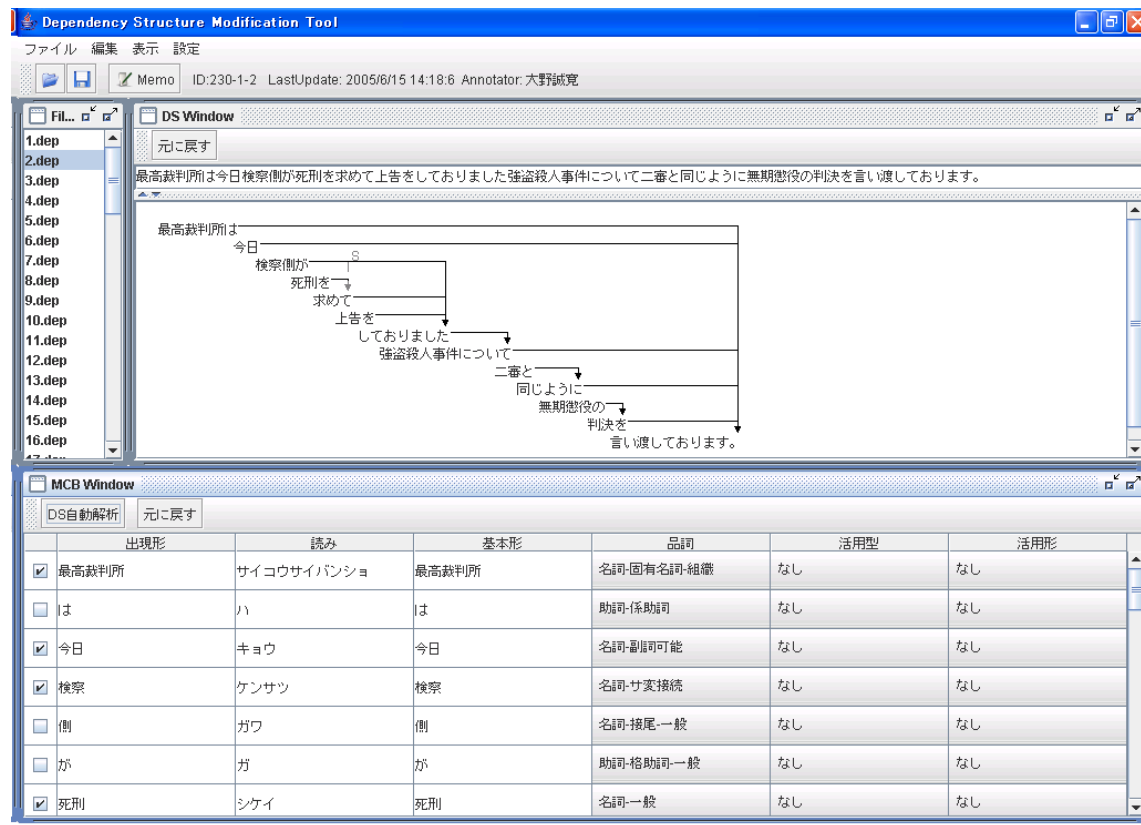


図 2.18: 修正インタフェース DSMT

タフェース DSMT (Dependency Structure Modification Tool) を作成し、これを利用した。以下では、修正インタフェース DSMT について概説する。

### 修正インタフェース DSMT

修正インタフェース DSMT の画面を図 2.18 に示す。DSMT の特徴として、2.3.2 節で述べた付与基準に従って係り受け構造情報を付与できるように、ある文節が係り先を 2 つ以上持つようなアノテーションを可能にしている。また、形態素や文節、文ごとにメモを記録することができる。なお、メモは、形態素情報や文節まとめあげ、係り受けの修正において、判断の迷う箇所や相談を要する箇所に対して付与され、後から参照するために利用される。

DSMT は、以下に示す 5 つの部位から構成される。

- メニューバー：各種コマンドがメニュー表示される。
- ツールバー：「開く」、「保存」、「メモ」のコマンドと文管理情報が表示される。

- File List Window : コーパスファイルのリストを表示する画面 .
- MCB Window : 形態素・文節の修正画面 .
- DS Window : 係り受け構造の修正画面 .

以下では、MCB Window と DS Window について説明する .

形態素と文節の修正は MCB Window で行う . MCB Window では 1 文中の形態素が各行ごとに記述されている . 各列には形態素の文中での表記 ( 出現形 ) , 読み , 基本形 , 品詞 , 活用型 , 活用形が表示されている .

出現形 , 読み , 基本形はテキスト編集によって修正でき , 品詞 , 活用型 , 活用形はメニューによって変更できる . 活用型は品詞に対応するものだけ , 活用形は活用型に対応するものだけがメニュー表示される . また , 以下の操作により , 1 つの形態素を 2 つ以上の形態素に分割したり , 2 つ以上の形態素を 1 つの形態素にまとめることができる .

- 形態素の行の追加

形態素を選択している状態でキーボードの Insert キーを押すと , 選択している形態素の下に行が 1 つ追加される .

- 形態素の行の削除

形態素を選択している状態でキーボードの Delete キーを押すと , 選択している形態素の行が削除される . このとき , 出現形 , 読み , 基本形のテキスト情報は 1 つ下の行にコピーされる .

さらに , 品詞のメニューの末尾には「メモ」という項目があり , この項目を選択すると立ち上がるダイアログボックス上でメモを入力することができる .

文節の区切りは , 文節の先頭の形態素を示すチェックマークを変更することにより修正できる . 形態素や文節の修正の終了時に , 「DS 自動解析」ボタンを押すことにより , 修正した内容に対して係り受け解析が自動的に実行され , その内容が DS Window に反映される .

次に , 係り受け構造の修正は DS Window で行う . この画面では , 1 文中の文節が対角線上に並んでおり , 係り受け関係が 2 つの文節をつなぐ矢印で表される .

係り受け関係の修正は , 係り文節を選択し , その後 , 受け文節を左クリックすることにより行う . また , 係り受け関係の種類 ( 通常 , 並列 , 部分並列 , 同格 ) を , 受け文節上で右クリックすることにより表示されるポップアップメニューにより選択

表 2.4: 「あすを読む」構文構造付き音声独話コーパスの規模

番組	110
文	6,523
節境界単位	30,636
文節	76,573
形態素	192,495

できる．なお、「倒置」は、通常の係り受け関係と同様に前方に位置する受け文節を左クリックすることにより修正できる．同様に、「係り先なし」は、受け文節として係り文節と同一の文節を左クリックすることにより修正できる．また、意味的係り受け関係は、係り文節上で右クリックして、ポップアップメニューを立ち上げ、「追加」を選択することにより追加できる．さらに、係り文節を選択していない状態で、対象文節を右クリックしポップアップメニュー上の「メモ」をクリックすることにより、各文節ごとにメモを入力することができる．

#### 2.3.4 「あすを読む」構文構造付き音声独話コーパスの規模と特徴

本研究で構築したコーパスの基礎統計データを表 2.4 に示す．

本コーパスと同様に、独話データに対して、係り受け情報が付与されたコーパスとして CSJ コーパス [53] がある．CSJ コーパスでは、約 50 万形態素の独話データに対して、形態素情報、文節境界情報、節境界情報、ならびに、係り受け情報が付与されている．本コーパスと CSJ コーパスと大きく異なる点として、本コーパスでは、構文的制約の 1 つ「係り先の唯一性」が緩和され、複数の係り先を持つことが許されていること、また、節間の関係も付与されている点が挙げられる．

次に、新聞記事に対して係り受け情報が付与されている京都テキストコーパス [49] と比較する．京都テキストコーパスでは、5,000 文に対しては、ある述語の省略されている格とその指示対象を示している省略関係等の情報も付与されている．本コーパスにおける「複数の係り先を持つことを許した係り受け情報」と京都テキストコーパスの「省略関係情報」の間には深い関連があると考えられ、今後、比較分析する必要がある．

## 2.4 2章のまとめ

本章では，話し言葉の構文的特徴を明らかにするための分析データ，及び，統計的構文解析手法の学習データとして利用することを目的に，話し言葉の構文構造付きコーパスを構築した．

まず，2.2 節で構文構造付き音声対話コーパスの構築について述べた．本コーパスは，CIAIR 車内音声対話コーパスに収録された対話文に対して，構文構造を付与することにより構築した．話し言葉に特有な言語現象を含む文に対しても構文構造を付与するために，独自の構文構造付与基準を設けている．構築した構文構造付き音声対話コーパスの規模は，45,053 文節，85,870 形態素である．

次に，2.3 節で構文構造付き音声独話コーパスの構築について述べた．本コーパスは，NHK の解説番組「あすを読む」の書き起こしコーパス中の独話文に対して，構文構造を付与することにより構築した．本コーパスでは，形態素情報や係り受け情報のほかに，節境界情報も付与している．構築した構文構造付き音声独話コーパスの規模は，76,573 文節，192,495 形態素に達している．



## 第3章 日本語話し言葉の頑健な係り受け解析

### 3.1 はじめに

近年、音声認識技術の進歩に伴い、音声対話処理システムの実現が大きなテーマとなっている。ユーザフレンドリな音声対話システムの実現のために、実対話環境下において利用可能な頑健な音声理解技術が望まれる。しかし、話し言葉にはフィラーや言い淀み、言い直しなどが頻出するため、それを従来の言語解析手法を用いて処理することは困難である。そのため、話し言葉に対する頑健な解析手法が必要とされている。英語などの比較的文法的制約が強い屈折言語では、句構造文法に基づいた規則主導型の解析手法により話し言葉の解析が実現されている（例えば、[6, 24]）。一方、日本語は、語順の入れ替わりや格助詞の省略が許されるなど文法的制約が弱いいため、これらの解析手法が日本語話し言葉に対しても有効であるとは必ずしもいえない。

本章では、大規模音声対話コーパスを用いた対話音声の係り受けに関する特徴分析について述べ、それに基づく日本語話し言葉の頑健な係り受け解析手法を提案する。従来の係り受け解析手法では、係り受けの非交差性、後方修飾性、係り先の唯一性の3つの制約が用いられてきたが[47]、本研究における分析の結果、対話音声では、倒置現象や係り先がない文節などを含む発話が無視できない頻度で出現することが分かった。そこで本手法では、後方修飾性の制約及び係り先の唯一性に関する制約を緩和する。すなわち、倒置を許し、フィラーなどの文節については、その係り先は存在しないとして解析する。解析結果は、部分的な係り受け構造によって表す。

また、本手法では、音声対話コーパスに付与された係り受け構造から各文節間の係り受け確率を統計的に獲得し、それを用いて係り受け構造の尤度を計算する。近年、大規模テキストコーパスから獲得した統計情報を利用した係り受け解析手法が盛んに研究されている。Collins[10]、藤尾ら[18]は、文節間の共起頻度を用いて係り受け確率を推定し、係り受け解析を行っている。Ratnaparkhi[78]や内元ら[92]、

表 3.1: 係り受け分析データ

対話	81
発話単位	7,781
ターン	6,078
文節	24,250

Charniak[8] は、最大エントロピー法に基づく学習モデルから係り受け確率を求め、その積により 1 文の係り受け構造の尤度を計算している。さらに、春野ら [23]、工藤ら [46] は、係り受け確率の学習器として、それぞれ決定木、SVM を利用する手法を提案している。

しかしながら、これらの手法はいずれも、書き言葉に対する解析手法であるため、話し言葉に適切なモデルとなっているかは明らかではない。話し言葉の係り受け解析手法として、伝 [14] は、言い直しや言い淀みなどを語と語の間の係り受け関係の一種と捉え、従来の係り受け解析を拡張した解析手法を提案しているが、扱える話し言葉の言語現象は限られている。

一方、本手法では、自然発話文に対して、統計情報を活用することにより、最尤の構造を作成する。本手法の有効性を評価するために、名古屋大学 CIAIR 車内音声対話コーパス [36, 37, 38, 42] に収録されたドライバ発話に対して、係り受け解析実験を行った。実験では、車内音声対話コーパスの 81 対話から得られた統計情報を用いて、日本語音声の書き起こしデータに対して係り受け解析を行った。実験の結果、話し言葉の頑健な解析における本手法の有効性を確認した。本章では、特に、係り先を持たない文節と倒置、発話単位をまたぐ係り受けの解析に対する本手法の有効性を報告する。

本章の構成は以下の通りである。次節で日本語話し言葉の言語的分析について述べ、3.3 節で話し言葉の頑健な係り受け解析手法を示す。3.4 節で解析実験について述べる。3.5 節で本手法の頑健性について考察し、3.6 節で本章のまとめと今後の課題について述べる。

## 3.2 自然発話の言語的分析

2.2 節で述べた CIAIR 構文構造付き音声対話コーパスを用いて、実走行車内音声対話の言語分析を行った。分析には、構築した構文構造付き音声言語コーパスの 81 対話を用いた。81 対話の基礎統計を表 3.1 に示す。

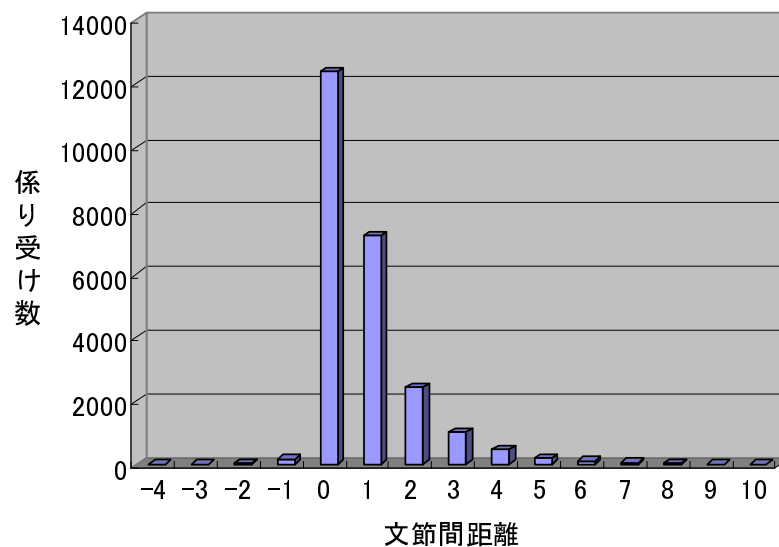


図 3.1: 文節間距離と係り受け数の関係

24,250 文節<sup>1</sup>に対して、11,866 個の係り受けが存在した。1 発話単位あたりの平均係り受け数は 1.52 個であり、1 ターンあたりでは 1.95 個であった。1 文あたりの係り受け数が平均 10 程度である新聞などの書き言葉と比べて、著しく少ない。しかし、このデータは話し言葉の係り受け解析が書き言葉のそれに比べ簡単であることを必ずしも意味しない。なぜなら、話し言葉では、あらゆる文節が受け文節をもっているわけではなく、係り先がない文節の特定も必要であるからである。実際に、全文節数中の約 51.1% が係り先をもたない文節であった<sup>2</sup>。

次に、倒置、及び、発話単位にまたがる係り受けについて調査した。倒置は 256 個存在した。全体の 4.2% のターンに出現しており、話し言葉解析において無視できる頻度ではない。その出現位置について調査したところ、倒置全体のうちの約 85.2% がターンの最後の文節に出現している。一方、発話単位にまたがる係り受けは、92 個存在した。複数発話単位から構成されている 1,362 ターンの 6.8% に出現しており、発話単位を話し言葉解析の処理単位とすることは、必ずしも適当でないことがわかる。

また、係り文節と受け文節の間の距離についても調べた。文節間距離と係り受けの数の関係を図 3.1 に示す。ただし、文節間距離は受け文節の位置と係り文節の位

<sup>1</sup>このうち、フィラー単独の文節は 3,049 個である。

<sup>2</sup>一般に、日本語書き言葉の場合、1 文の最終文節のみが係り先を持たない。一方、本章で対話の解析単位として定義している 1 ターンにも、書き言葉における文末文節と同じ役割をもった係り先を持たない文節がただ 1 つ存在する。書き言葉の文における文末文節に相当する文節を除くと、係り先を持たない文節は、6,306 個存在し、全文節の 26.0% を占める。

置の差で測定する．また，文節間距離が0であることは，受け文節を持たない係り受け，すなわち，自分自身に係る係り受け関係であることを意味する．

### 3.3 統計的係り受け解析

本手法では，従来，係り受け解析で用いられてきた構文的制約を緩和し，また，音声対話コーパスから獲得した統計情報を用いることにより最尤の係り受け構造を作成する．なお，係り受けが発話単位にまたがることもあるという3.2節の分析結果に基づき，本手法では1ターンを解析の単位とした．

#### 3.3.1 係り受けの構文的制約

3.1節で述べたように，従来の日本語係り受け解析では，一般に以下の3つの構文的制約に従う．

係り受けの非交差性 係り受けは互いに交差しない．

係り受けの後方修飾性 文末の文節を除き，必ず後方に位置する文節に係る．

係り先の唯一性 文末の文節を除き，係り先は必ず存在し，かつ，2つ以上存在しない．

しかし，日本語話し言葉には，前節で述べたように話し言葉に固有の現象や係り受けが頻出するため，これらの制約が完全に成り立つことを前提とすることは難しい．そこで本手法では，頑健な係り受け解析を実現するために，上記の構文的制約を緩和した．すなわち，フィラーや言い淀み，言い誤りが多数出現することに着目し，係り先が存在しない文節を認め，その係り先をその文節自身に係るとした（係り先の唯一性の緩和）．また，倒置に対処するため，前方の文節に係ることを認める（係り受けの後方修飾性の緩和）<sup>3</sup>．

#### 3.3.2 話し言葉の統計的係り受け解析

本手法では，形態素解析，及び，文節まとめあげが施された文節列を入力とする．入力文節列  $B (= b_1 \cdots b_n)$  の係り受け構造を  $S$  とするとき， $P(S|B)$  の確率値を最大にする係り受け構造  $S$  を求める．通常，書き言葉に対する係り受け解析手法では，

<sup>3</sup>係り受けが交差することは稀であるため，係り受けの非交差性は遵守する．

係り受けの非交差性，後方修飾性，係り先の唯一性の3つの性質を，絶対的制約として用いるが，本手法では，倒置やフィラー，言い淀み，言い誤りなどが頻出することに注目し，非交差性のみを，満たすべき性質として係り受け構造を求める．ただし，後方修飾性，及び，係り先の唯一性に関する制約は統計情報を反映させつつ緩和した．

それぞれの係り受けは独立であると仮定すると， $P(S|B)$  は以下の式で計算できる．

$$P(S|B) = \prod_{i=1}^n P(b_i \xrightarrow{rel} b_j | B), \quad (3.1)$$

ここで， $P(b_i \xrightarrow{rel} b_j | B)$  は，入力文節列  $B$  が与えられたときに，文節  $b_i$  から  $b_j$  への係り受け関係がある確率を表す．最尤の係り受け構造は，式 (3.1) の確率を最大とする構造であるとして動的計画法を用いて計算する．ここで，本手法では， $P(S|B)$  を計算する時に，従来の書き言葉に対する統計的係り受け解析手法 [18, 92] では利用されない，最終文節の係り受け確率  $P(b_n \xrightarrow{rel} b_j | B)$  も利用していることに注意されたい．

次に， $P(b_i \xrightarrow{rel} b_j | B)$  の計算について述べる．本手法では，統計情報を用いることにより係り受け構造の確からしさを計算する．統計情報として利用する属性は以下の通りである．

- 係り文節  $b_i$ :
  - 係り文節  $b_i$  中の最も遅く発話された自立語の基本形:  $h_i$
  - 係り文節  $b_i$  中の最も遅く発話された自立語の品詞:  $t_i$
  - 係り文節  $b_i$  の係りの種類:  $r_i$
  - 係り文節  $b_i$  の位置:  $l_i$
- 受け文節  $b_j$ :
  - 受け文節  $b_j$  中の最も遅く発話された自立語の基本形:  $h_j$
  - 受け文節  $b_j$  中の最も遅く発話された自立語の品詞:  $t_j$
- 係り文節  $b_i$  と受け文節  $b_j$  の係り受け距離:  $d_{ij}$
- 係り文節  $b_i$  と受け文節  $b_j$  の間に存在するポーズの数:  $p_{ij}$

ここで，係りの種類とは，係り文節が付属語を伴うときはその付属語の語彙，品詞，活用形であり，そうでない場合は一番最後の形態素の品詞，活用形である．表 3.2 に

表 3.2: 係り文節と係りの種類の例

係り文節	係りの種類
大きい	形容詞-連体形
買える	動詞-連体形
ないかな	終助詞「な」
その	助詞-連体化「の」
電話が	格助詞「が」
近くに	格助詞「に」
ちょっと	副詞
コンビニ	名詞
えーと	なし
そ	なし

例を挙げる．係り受け距離  $d_{ij}$  は，倒置現象を取り扱うため，負の値を取ることを許す．また，係り文節の位置  $l_i$  は，その文節が入力ターン内で一番最後の文節か否かを表す．前節で述べたように，倒置の多くは入力文の最後の係り文節に出現する傾向があり，倒置の確率を推定するためにこの属性を設けた．

以上の属性を用いて，各文節間の係り受け確率  $P(b_i \xrightarrow{rel} b_j | B)$  を以下のように計算する．

$$\begin{aligned}
 P(b_i \xrightarrow{rel} b_j | B) &\cong P(b_i \xrightarrow{rel} b_j | h_i, h_j, t_i, t_j, r_i, d_{ij}, p_{ij}, l_i) \\
 &= \frac{C(b_i \xrightarrow{rel} b_j, h_i, h_j, t_i, t_j, r_i, d_{ij}, p_{ij}, l_i)}{C(h_i, h_j, t_i, t_j, r_i, d_{ij}, p_{ij}, l_i)}.
 \end{aligned} \tag{3.2}$$

ここで， $C$  は共起頻度関数である．係り先のない文節はそれ自身に係る（すなわち  $i = j$ ）とみなすことにより，係り先をもたない場合の確率も計算できる．

$P(b_i \xrightarrow{rel} b_j | B)$  を式 3.2 により計算する際に生じるデータスパースネス問題に対処するため，本手法では，藤尾ら [18] によるスムージング手法を利用している．具体的には，式 (3.2) の分子  $C(h_i, h_j, t_i, t_j, r_i, d_{ij}, p_{ij}, l_i)$  が 0 である場合は， $P(b_i \xrightarrow{rel} b_j | B)$  を式 (3.3) により計算した．

$$\begin{aligned}
 P(b_i \xrightarrow{rel} b_j | B) &\cong P(b_i \xrightarrow{rel} b_j | t_i, t_j, r_i, d_{ij}, p_{ij}, l_i) \\
 &= \frac{C(b_i \xrightarrow{rel} b_j, t_i, t_j, r_i, d_{ij}, p_{ij}, l_i)}{C(t_i, t_j, r_i, d_{ij}, p_{ij}, l_i)}.
 \end{aligned} \tag{3.3}$$

表 3.3: 文節間係り受け確率

		受け文節					
		えーと	コンビニ	ないかな	そ	そのの	近くに
係り文節	えーと (フィラー)	1.00	0.00	0.00	0.00	0.00	0.00
	コンビニ	0.00	0.01	0.40	0.00	0.00	0.00
	ないかな	0.00	0.00	0.88	0.00	0.00	0.00
	そ (言い淀み)	0.00	0.00	0.00	1.00	0.00	0.00
	そのの	0.00	0.02	0.00	0.00	0.00	0.75
	近くに	0.00	0.00	0.80	0.00	0.00	0.02

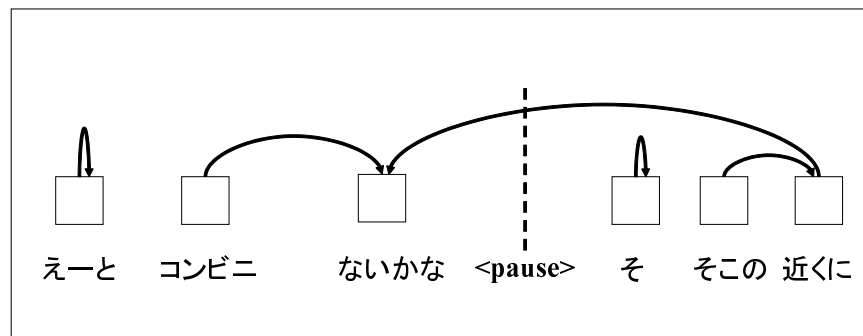


図 3.2: 「えーと/コンビニ/ないかな/&lt;pause&gt;/そ/そのの/近くに」に対する係り受け構造

### 3.3.3 解析例

ポーズを含む日本語発話文「えーとコンビニないかな<pause>そそのの近くに」に対する解析例を示す。

この文を文節に区切ると「えーと/コンビニ/ないかな/<pause>/そ/そのの/近くに」となる。これらの文節の係りの種類、及び、各文節間の係り受け確率をそれぞれ表 3.2, 3.3 に示す。表 3.3 は、例えば、「コンビニ」が「ないかな」に係る確率が 0.40 であることを表している。また、「えーと」が「えーと」に係る確率とは、「えーと」がどの文節にも係らない確率を意味する。表 3.3 から得た各係り受け確率を用いて、最大の確率をもつ構造を作成した結果、図 3.2 のような係り受け構造が得られる。

表 3.4: 実験結果 (係り受け正解率)

	係り受け単位	ターン単位
本手法	87.0% (21,089/24,250)	70.1% (4,260/6,078)
ベースライン 1	53.8% (13,058/24,250)	30.5% (1,853/6,078)
ベースライン 2	57.9% (14,049/24,250)	38.3% (2,329/6,078)

### 3.4 解析実験

本手法の有効性を評価するため、係り受け解析実験を行った。実験には、2.2 節で述べた CIAIR 構文構造付き音声対話コーパスを用いた。

#### 3.4.1 実験の概要

実験データとして 81 対話における運転者の発話を使用した。その規模は、24,250 文節からなる 6,078 ターンである (平均ターン長は 4.0 文節)。

対話ごとに分割し交差検定を行った。すなわち、81 対話におけるある 1 対話をテストセットとし、残りの対話を学習セットとする実験を 81 回繰り返した。なお、上述の分析結果に従って、本手法では 1 ターンを解析の単位とした。

また、比較評価するため、ベースラインとして、以下の 2 つの解析手法によっても係り受け解析を実行した。

- ターンの最終文節は係り先なしとして、その他の文節は直後の文節に係るとして解析する手法 (ベースライン 1)
- 倒置や係り先を持たない文節を許さないほかは、提案手法と同じ手法 (ベースライン 2)

#### 3.4.2 実験結果

表 3.4 に、各手法の係り受け、及び、ターンに対する正解率を示す。本手法は、24,250 文節のうち、正解データと係り先が一致したものが 21,089 個あり、係り受け正解率は 87.0% であった。これは、2 つのベースラインと比較しても非常に高い値であった。本手法により、自然発話文に対しても、書き言葉を対象とした他の係り受け解析手法 [18, 23, 45, 92] と同等の高い精度で係り受けを抽出できることを確認した。

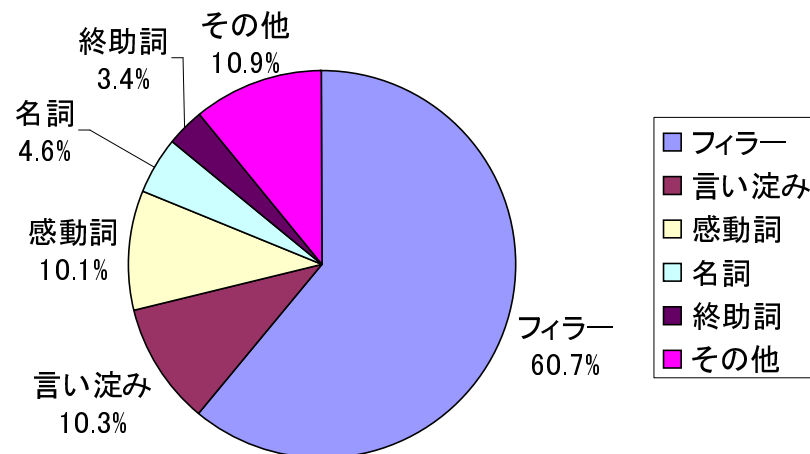


図 3.3: 最終文節以外で受け文節がない係り受けの内訳

## 3.5 考察

本節では、自然な話し言葉に対する本手法の頑健性とデータスパースネス問題について考察する。

### 3.5.1 話し言葉に特有な現象に対する頑健性

以下では、受け文節のない係り受け、前方の文節への係り受け、及び、ポーズをまたぐ係り受け現象に着目し、自然な話し言葉に対する本手法の頑健性について、3.4節の実験結果をもとに考察する。

#### 受け文節のない係り受け

話し言葉では、フィラーや言い淀みなど、必ずしもすべての文節に受け文節が存在するとは限らない。実験で使用したコーパスでは、全文節の51.1%に相当する12,384文節には受け文節がなく、そのうちポーズの直前に位置しないものは4,937文節であった。その内訳を図3.3に示す。その約7割がフィラー及び言い淀みである。これらについては受け文節がないとしてコーパスが作成されており、その特定は難しい。そこで残りの3割の係り受けに対する実験結果を表3.5に示す。また、図3.4に、このような係り受けに対して正しく解析できた例を示す。この3割を構成する

表 3.5: 受け文節のない係り受けの解析結果（ポーズの直前でなく，フィラー・言い淀み以外）

適合率	60.4% (996/1,650)
再現率	69.5% (996/1,434)

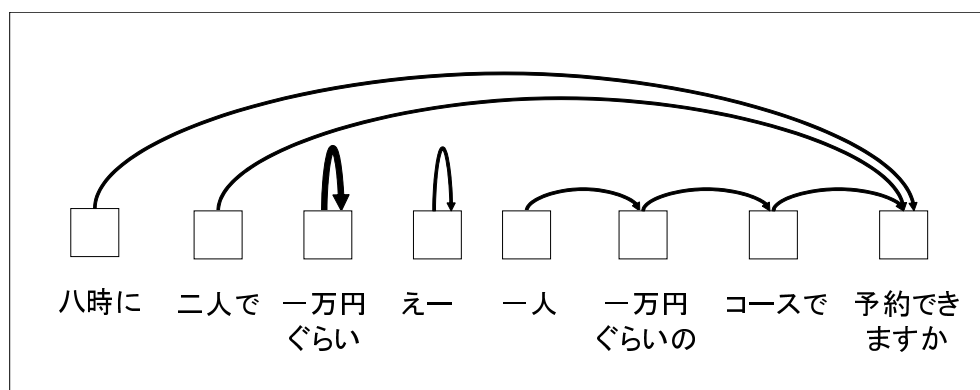


図 3.4: 受け文節のない係り受けの例

表 3.6: 前方の文節への係り受けの解析結果

適合率	60.5% (49/ 81)
再現率	19.1% (49/256)

1,434 文節のうち，996 文節の係り先を正しく特定できており，本手法がこのような係り受けに対しても高い正解率で解析できることを示している．

#### 前方の文節への係り受け

本手法では，倒置関係の同定を可能にするために，後方修飾性を絶対的な性質として定めなかったが，係り受け解析において，倒置，すなわち，前方修飾性を許すと，係り先を特定するための探索空間が約 2 倍に広がるため，正しい解析の実現が困難になる．前方の文節への係り受けは 256 個存在し，必ずしも無視できる数ではないが，割合としては全体のわずか 1% 程度に過ぎず，そのような文節を特定することの意義は必ずしも明らかではない．

前方の文節への係り受けに関する結果を表 3.6 に示す．また，前方の文節への係り受けに対して正しく解析できた例を図 3.5 に示す．再現性については必ずしも高いとはいえないものの，適合率は 60% を超えている．これは，前方への係り受けを

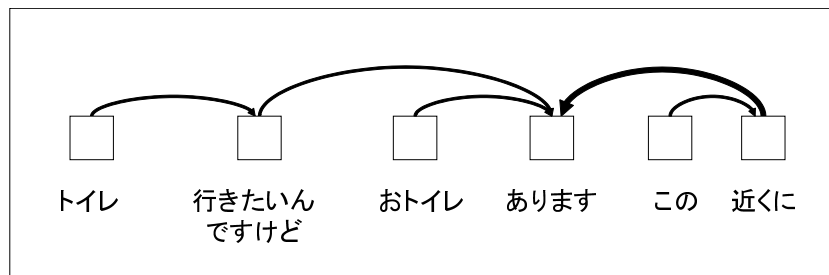


図 3.5: 前方の文節への係り受けの例

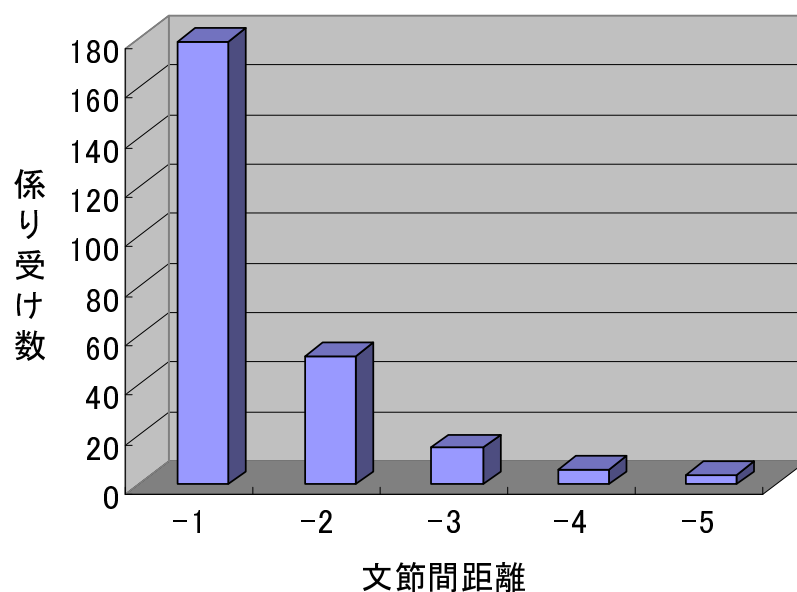


図 3.6: 前方の文節への係り受けの文節間距離

許すことにより、解析精度が上昇することを意味しており、ひいては本手法の倒置現象への頑健性を示している。このような良好な結果が得られた理由として、倒置の出現に関して以下の2つの傾向が存在することが挙げられる。

1つは、文節の位置に関する出現傾向である。すなわち、倒置となる係り受け関係における係り文節の多くは、発話ターンの最終文節に出現する（倒置全体の85.2%）ことを考慮し、係り受け確率を計算する式(3.2)の属性として係り文節の位置を導入したことの効果が現れている。実際、前方への係り受けであるとして同定された81個のうち、係り文節が発話ターンの最後に位置している係り受けの適合率は75.0%であった。

もう1つは、文節間距離に関する出現傾向である。図3.6に示すように、前方の

表 3.7: ポーズをまたぐ係り受けの解析結果

適合率	6.5% (34/521)
再現率	37.0% (34/ 92)

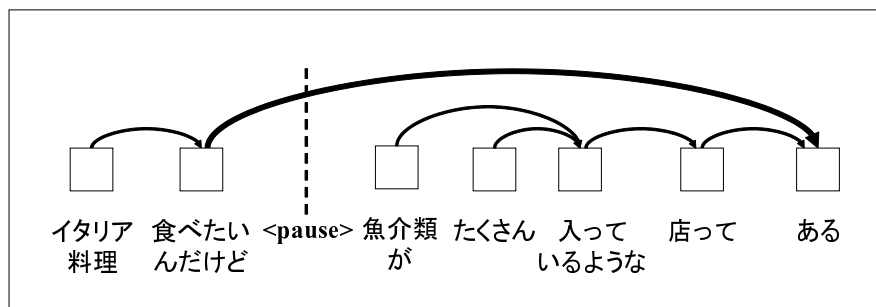


図 3.7: ポーズをまたぐ係り受けの例

文節への係り受けの文節間距離は、全体の 90.2% が  $-1$  あるいは、 $-2$  のいずれかであった。このことは、係り受け確率の計算においても  $d_{ij}$  が負の値をとることを許すことにより、反映されている。実験結果では、倒置のうち、文節間距離が  $-2$  以上の係り受けの適合率は 61.0% であった。

#### ポーズをまたぐ係り受け

話し言葉では、書き言葉でいう文に相当する文法単位を特定するのは容易ではない。ポーズがそのような単位の境界を形成している可能性は高いが、そうでない場合も少なくないため、本実験ではターンを解析単位とした。実験データには 92 個のポーズをまたぐ係り受けが存在した。

表 3.7 にポーズをまたぐ係り受けに関する実験結果を、図 3.7 にこのような係り受けを正しく解析できた例を示す。適合率は著しく低い結果となった。この結果は以下の理由により生じたと考えられる。

- そもそもそのような係り受けの出現頻度が多くない（全体の 0.4%）こと
- ポーズをまたぐ場合の文法的特徴が必ずしも明らかでなく、それを確率の計算式に導入していないこと

その一方で、37.0% の再現率を実現しており、ある程度の効果は見られた。

### 3.5.2 データスパースネス問題に対する頑健性

式(3.2)に示すように、本手法では、話し言葉を頑健に解析するため、多くの属性を利用している。高い解析精度を達成するためには、大規模なコーパスが必要となる。実際、本実験で利用した学習データも、話し言葉に構文構造が付与されたデータとして必ずしも小さいサイズではない。しかしながら、たとえ学習データが十分に大きくななくても、本手法が音声対話の理解に有効であることが期待できる。その理由は以下の通りである。

1. 本手法で利用したスムージング手法がデータスパースネス問題に対して有効であることが、同様のスムージング手法を利用した書き言葉に対する従来研究によって示されている（例えば、[10, 18]）。
2. 現状の音声対話システムの多くは、本実験で対象とした車内音声などの特殊なドメインに限定して利用されている。そして、ユーザーの発話スタイルは、ドメインが限定されている場合、変化が少ない。

データが少なくなることに対する影響を評価するため、新たな実験を行った。実験には、3.4節の実験と全く同じ81対話における6,000ターンを利用した。6,000ターンのうち、500ターン（2,001文節）をテストデータとして利用した。残りの5,500ターンを学習データとして利用し、500ターンずつ学習データを増加させて、計11回の実験を行った。

図3.8に、係り受け及びターンに対する正解率と学習データ量の関係を示す。この図から、学習データの量が多い場合と比べて、少ない場合の正解率が極端に低下していないことがわかる。

## 3.6 3章のまとめ

本章では、日本語話し言葉に対する係り受け解析手法を提案した。本手法は、従来の係り受け解析で用いられてきた構文的制約を緩和し、統計情報を利用することにより、頑健な解析を行うことができる。CIAIR 構文構造付き音声対話コーパスを用いた実験の結果、自然な話し言葉の解析における本手法の有効性を確認した。

本研究では、音声認識率が100%であると仮定して実験を行ったが、音声対話処理システムでは音声認識によって生成されたテキストを処理することになる。そこには、多数の認識誤りが含まれるため、極めて頑健な解析技術が要求される[58]。本

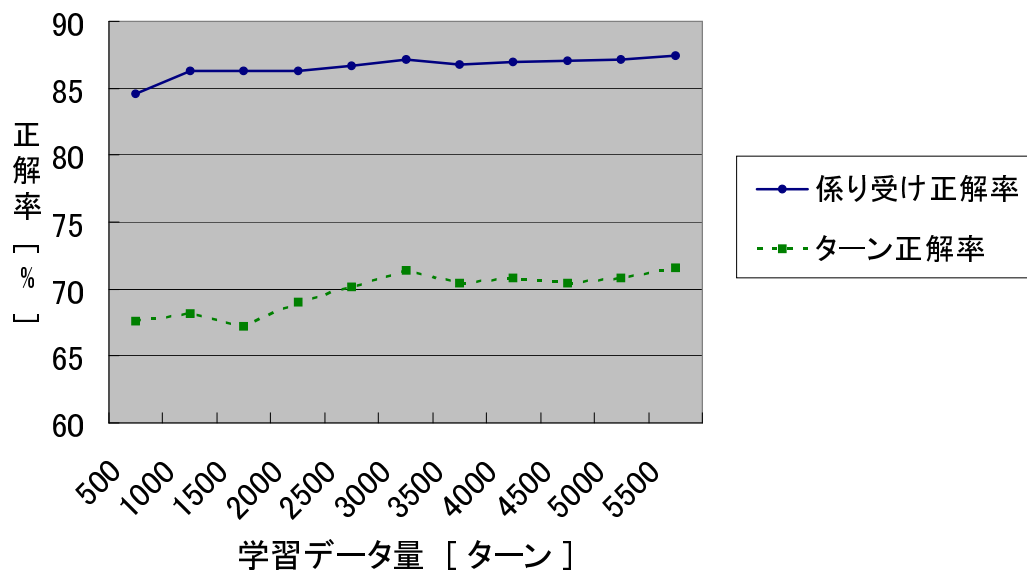


図 3.8: 係り受け及びターンに対する正解率と学習データ量の関係

章で提案した統計的係り受け解析の実用性を検証するために、今後は、音声認識システムを用いた評価を予定している。

## 第4章 日本語話し言葉の効率的な係り受け解析

### 4.1 はじめに

近年，講演や解説などの独話データが人間の貴重な知的財産として注目を集めている．蓄積された独話データを有効かつ効率的に活用するには，単に蓄積するだけでは十分ではない．独話の構造を解析しその情報とともに蓄積することで要約や検索などのシステムがその情報を利用できるようになり，効率的なアクセスや効果的な再利用が可能となる．

しかし，これまで独話を対象とした構文解析の研究はほとんど行われていない．自然に発話された独話には，フィラーや言い淀み，言い直しなどの非文法的な言語現象が頻出する．これらの現象を頑健に扱うため，対話文の構文解析手法が提案されている [6, 11, 13, 24]．一方，独話には，従来の音声言語研究により対象とされてきた対話文と比べて，文の長さが長く，時に極端に長い文が含まれることがあるという特徴もある [32]<sup>1</sup>．一般に，文が長くなればなるほど，その構文構造はより複雑になるため，独話文に対して解析を実行すると，解析時間が長くなるうえ，高い解析精度の達成が難しくなる．高い性能を備えた，より効率的な独話文解析を実現するためには，適切な単位に文を分割し，簡単化することが効果的な方法である．

本章では，文の分割に基づく独話文の係り受け解析手法を提案する．本手法では，節レベルと文レベルの二段階で係り受け解析を実行する．まず，節境界解析により文を節に分割し，各節に対して係り受け解析を行うことにより，節内の係り受け関係を同定する．次に，節境界をまたぐ係り受け関係を定め，文全体の係り受け構造を作り上げる．独話文に対する係り受け解析実験の結果，本手法により解析精度が低下することなく解析時間を大幅に短縮できることを確認した．

本章の構成は以下の通りである．次節で日本語独話文の解析単位について述べる．

<sup>1</sup>本研究で利用した独話コーパス「あすを読む」の平均文長は 29.1 (形態素/文) であるのに対して，対話コーパス SLDB[65] と BTEC[91] の平均文長はそれぞれ 11.7 と 7.9 (形態素/文) である．さらに，独話コーパス「あすを読む」には 100 形態素以上からなる極端に長い文が存在していることが報告されている [32]．

4.3 節で節境界に基づく係り受け解析手法について説明する．4.4 節で解析実験について述べ，4.5 節で考察する．4.6 節で関連研究を示し，4.7 節で本章のまとめと今後の課題について述べる．

## 4.2 独話文の係り受け解析における処理単位

本手法では，文よりも短い単位を解析単位とすることにより，解析を効率化する．1 文が長い独話文の解析では，係り受け関係の探索範囲が狭められ，解析時間を短縮することができる．

### 4.2.1 節と係り受け

節とは，述語を中心としたまとまりであり，複文や重文の場合，文は複数の節から構成される [56]．さらに，節は，構文的かつ意味的にまとまった単位であるため，文に代わる解析単位として利用できる．

そこで本章では，「文は1つ以上の節の接続であり，各節を構成する文節は，節の最終文節を除き，その節の内部の文節に係る」とみなし，それに基づく係り受け解析手法を提案する．

例として，独話文「先日総理府が発表いたしました世論調査によりますと死刑を支持するという人が八十パーセント近くになっております」の係り受け構造を図 4.1 に示す．この文は以下の4つの節から構成される．

- 先日総理府が発表いたしました
- 世論調査によりますと
- 死刑を支持するという
- 人が八十パーセント近くになっております

各節が係り受け構造（図 4.1 の実線矢印）を形成し，それらが節の最終文節からの係り受け関係（図 4.1 の点線矢印）でつながっている．

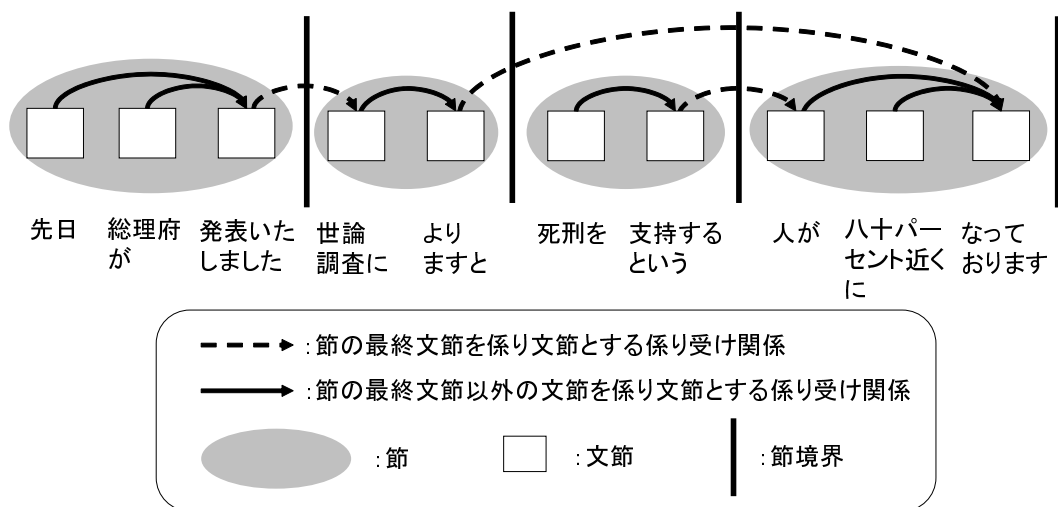


図 4.1: 節と係り受けの関係

### 4.2.2 節境界単位

節を文に代わる解析単位とするためには、係り受け解析の前処理として独話文を節に分割する必要がある。しかし、節には、主節の中に埋め込まれた従属節も存在するため、本来、文を節に一次的に分割することはできない[32]。

ただし、節への分割は、節境界解析プログラム[55]により近似的に実現できる。この節境界解析プログラムでは、節の終端位置を検出することにより節に相当する単位を検出する。さらに、このプログラムは、1～3形態素という局所的な範囲のみから、節の終端位置と種類を特定することができる。日本語形態素解析ツール茶釜[59]で形態素解析した文を入力すると、入力文中に含まれるすべての節境界の位置が特定され、その種類を表す節境界ラベルが挿入される。挿入される節境界ラベルは、「テ節」や「並列節ケレドモ」など、合計147種類である<sup>2</sup>。

本章の研究では、節境界解析により検出された節境界ではさまれた単位を節境界単位と呼び、これを新たな解析単位と考える。

### 4.2.3 節境界単位と係り受けの関係

本手法では、節境界単位は構文的にまとまっていると仮定し、これを係り受け解析の解析単位として利用する。しかし、実際には、図4.2に示すような節境界をまたぐ係り受け関係が存在する。そこで、上述した仮定の妥当性を検証するため、独

<sup>2</sup>節境界の種類の中には「主題八」や「感動詞」「談話標識」など「述語を中心としたまとまり」という節の定義から逸脱しているものも存在する。しかし、これらも、構文的に大きな切れ目になると考えられるため、節境界とみなしている[55]。

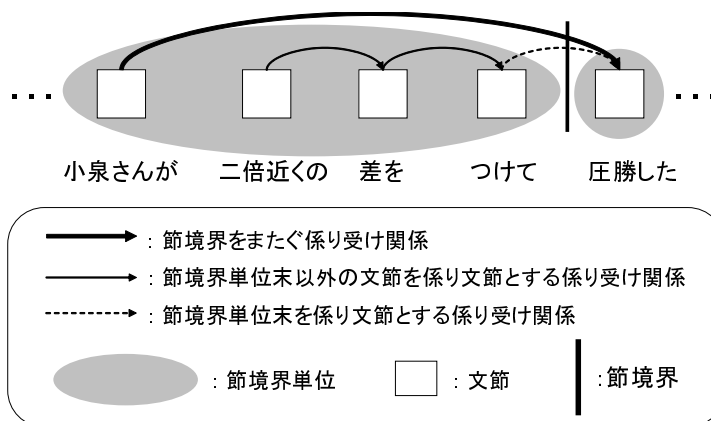


図 4.2: 節境界をまたぐ係り受けの例

表 4.1: 「あすを読む」200 文の基礎統計

文	200
節境界単位	951
文節	2,430
形態素	6,017
節境界をまたぐ係り受け	94

話コーパスを用いて分析した。分析には、2.3 節で構築した「あすを読む」構文構造付き音声独話コーパスの 200 文を用いた。

200 文の基礎統計を表 4.1 に示す。総文節数 2,430 文節のうち、節境界単位の最終文節 (951 文節) を除いた 1,479 文節の中で、94 文節のみが節境界単位の外に位置する文節に属していた。これは、全体の 93.6% (1,385/1,479) の係り受けが節境界単位で閉じていることを意味しており、4.2.1 節で設定した仮定がある程度妥当なものであることを確認した。

### 4.3 節境界に基づく係り受け解析

本手法では、4.2.1 節で設けた仮定に基づき、形態素解析、文節まとめあげ、及び節境界解析が施された文を入力とする<sup>3</sup>。解析の手順は以下の通りである。

<sup>3</sup>独話には明示的な文末標識が存在していないため、事前に独話を文に分割することは容易ではない。しかし、これまでに文境界の検出手法がいくつか提案されており [27, 39, 85, 86]、本手法では、文境界に係り受け解析の前に自動的に検出できることを仮定している。

## 1. 節レベルの係り受け解析

1 文中のすべての節境界単位に対して，各節境界単位ごとにその内部の係り受け構造を解析する．

## 2. 文レベルの係り受け解析

1 文中のすべての節境界単位に対して，その最終文節の係り先を解析する．

なお，本章の以下では，1 文を構成する節境界単位列を  $C_1 \cdots C_m$ ，節境界単位  $C_i$  を構成する文節列を  $b_1^i \cdots b_{n_i}^i$ ，文節  $b_k^i$  を係り文節とする係り受け関係を  $dep(b_k^i)$ ，1 文の係り受け構造を  $\{dep(b_1^1), \dots, dep(b_{n_m}^m)\}$  と記す．

本手法では，まず，あらゆる節境界単位  $C_i$  に対して，節境界単位内の係り受け構造  $\{dep(b_1^i), \dots, dep(b_{n_i-1}^i)\}$  を求める．その次に，節境界単位の最終文節を係り文節とする係り受け関係により構成される係り受け構造  $\{dep(b_{n_1}^1), \dots, dep(b_{n_m-1}^m)\}$  を求める．ここで，本研究では，係り受け構造を求める際，統計的な手法を用いた．これまでに，統計的係り受け解析手法の学習モデルとして，決定木 [23] や最大エントロピー法 [92]，SVM[46] を用いた手法などが提案されているが，本章の研究では，節境界を用いることによる解析性能の向上を目的とするため，学習モデルとしては素朴なモデルを用いることとし，具体的には，文献 [18] で提案された共起確率に基づく統計的係り受け解析手法を採用した．なお，両レベルの解析においても，従来の日本語係り受け解析の制約として用いられている 3 つの構文的制約（係り受けの非交差性，後方修飾性，係り先の唯一性）を満たすことを仮定する<sup>4</sup>．

## 4.3.1 節レベルの係り受け解析

節レベルの係り受け解析は，節境界単位  $C_i$  中の文節列を  $B_i (= b_1^i \cdots b_{n_i}^i)$  とするとき，条件付確率  $P(S_i|B_i)$  を最大にする係り受け構造  $S_i (= \{dep(b_1^i), \dots, dep(b_{n_i-1}^i)\})$  を求める．なお，節レベルの係り受け解析では，節境界単位の最終文節  $b_{n_i}^i$  の受け文節は決定しない．

それぞれの係り受け関係は独立であると仮定すると， $P(S_i|B_i)$  は以下の式で計算

<sup>4</sup>本章の研究では，話し言葉の特徴のうち，長文が頻出するという特徴に着目するため，放送で行われる比較的整った解説における話し言葉を対象としており，このような話し言葉には構文的制約に逸脱する係り受けはほとんど見られないと仮定している．実際，4.2.3 節で分析に用いた 200 文には，非交差性，後方修飾性に逸脱した係り受けは存在しなかった．また，係り先の唯一性に逸脱する係り受けについても，フィラーを除いてほとんど出現しなかった．なお，フィラーは，係り受け解析の前の段階で検出できる [4, 88] と考え，あらかじめ削除している．

できる．

$$P(S_i|B_i) = \prod_{k=1}^{n_i-1} P(b_k^i \xrightarrow{rel} b_l^i|B_i), \quad (4.1)$$

ここで， $P(b_k^i \xrightarrow{rel} b_l^i|B_i)$  は，入力文節列  $B_i$  が与えられたときに，文節  $b_k^i$  が  $b_l^i$  に係る確率を表す．従来の文単位の統計的係り受け解析手法とは異なり，本手法では， $B_i$  は文ではなく節境界単位を構成する文節列である．最尤の係り受け構造は，式 (4.1) の確率  $P(S_i|B_i)$  を最大とする構造であるとして動的計画法を用いて計算する．

次に， $P(b_k^i \xrightarrow{rel} b_l^i|B_i)$  の計算について述べる．まず，係り文節における自立語の基本形を  $h_k^i$ ，その品詞を  $t_k^i$ ，係りの種類を  $r_k^i$  とし，受け文節における自立語の基本形を  $h_l^i$ ，その品詞を  $t_l^i$  とする．また，文節間距離を  $d_{kl}^{ii}$  で記し，文節間距離が1であるかそれ以上かであるかの2値とする．ここで，係りの種類とは，3章で提案した音声対話文に対する頑健な解析手法で用いた属性と同じもので，係り文節が付属語を伴うときはその付属語の語彙，品詞，活用形であり，そうでない場合は一番最後の形態素の品詞，活用形である．なお，これらの属性は，従来の係り受け解析手法 [18, 46, 92] で用いられてきたものと同様である．

さらに，本手法では，受け文節が節境界単位の最終文節であるか否かを示す属性  $e_l^i$  を導入する．文単位で係り受け解析を行う従来手法 [18, 46, 92] では，受け文節が文末であるか否かを示す属性がよく用いられているが，節レベルの係り受け解析では，文に相当する単位として節境界単位を考えているためである．

以上の属性を用いて，確率  $P(b_k^i \xrightarrow{rel} b_l^i|B_i)$  を以下のように計算する．

$$\begin{aligned} P(b_k^i \xrightarrow{rel} b_l^i|B_i) &\cong P(b_k^i \xrightarrow{rel} b_l^i|h_k^i, h_l^i, t_k^i, t_l^i, r_k^i, d_{kl}^{ii}, e_l^i) \\ &= \frac{F(b_k^i \xrightarrow{rel} b_l^i, h_k^i, h_l^i, t_k^i, t_l^i, r_k^i, d_{kl}^{ii}, e_l^i)}{F(h_k^i, h_l^i, t_k^i, t_l^i, r_k^i, d_{kl}^{ii}, e_l^i)}. \end{aligned} \quad (4.2)$$

ただし， $F$  は共起頻度関数である．

なお，本手法では，式 (4.2) により  $P(b_k^i \xrightarrow{rel} b_l^i|B_i)$  を計算するときに起こるデータスパースネスの問題を解決するために，藤尾ら [18] のスムージング手法を用いている．すなわち，式 (4.2) 中の  $F(h_k^i, h_l^i, t_k^i, t_l^i, r_k^i, d_{kl}^{ii}, e_l^i)$  が0である場合は，次式 (4.3) を用いて  $P(b_k^i \xrightarrow{rel} b_l^i|B_i)$  を計算する．

$$\begin{aligned} P(b_k^i \xrightarrow{rel} b_l^i|B_i) &\cong P(b_k^i \xrightarrow{rel} b_l^i|t_k^i, t_l^i, r_k^i, d_{kl}^{ii}, e_l^i) \\ &= \frac{F(b_k^i \xrightarrow{rel} b_l^i, t_k^i, t_l^i, r_k^i, d_{kl}^{ii}, e_l^i)}{F(t_k^i, t_l^i, r_k^i, d_{kl}^{ii}, e_l^i)}. \end{aligned} \quad (4.3)$$

### 4.3.2 文レベルの係り受け解析

節境界単位の最終文節の受け文節を同定する．1文の文節列を  $B(= B_1 \cdots B_m)$  とし，節境界単位 of 最終文節を係り文節とするような係り受け構造  $\{dep(b_{n_1}^1), \dots, dep(b_{n_{m-1}}^{m-1})\}$  を  $S_{last}$  とするとき， $P(S_{last}|B)$  を最大とする  $S_{last}$  を求める． $P(S_{last}|B)$  は以下の式で計算できる．

$$P(S_{last}|B) = \prod_{i=1}^{m-1} P(b_{n_i}^i \xrightarrow{rel} b_l^j | B), \quad (4.4)$$

ここで， $P(b_{n_i}^i \xrightarrow{rel} b_l^j | B)$  は，1文の文節列  $B$  が与えられたときに， $C_i$  の最終文節  $b_{n_i}^i$  が  $b_l^j$  に係る確率を表す．最尤の係り受け構造は，式 (4.4) の確率を最大とする構造であるとして動的計画法を用いて計算する．本手法では，先に解析した節境界単位内部の係り受け構造を前提として決定する．すなわち，後方に位置するすべての文節を受け文節の候補として計算するのではなく，節境界単位内部の係り受け構造から非交差性を満たすものだけを受け文節の候補として計算する．図 4.1 の場合，文節「支持するという」の受け文節は「人が」または「なっております」のいずれかであるとして計算する．

また， $P(b_{n_i}^i \xrightarrow{rel} b_l^j | B)$  は式 (4.5) により計算する．式 (4.5) では，式 (4.2) で利用した全属性を使用し，さらに，受け文節  $b_l^j$  が文末文節であるか否かを示す属性  $s_l^j$  を新たに用いる．ここで，文レベルの係り受け解析においても，属性  $e_l^j$  を利用している．4.2.3 節の 200 文を分析した結果，節境界単位 of 最終文節は，その 70.6%(522/751) が別の節境界単位 of 最終文節に係ることがわかったためである．

$$\begin{aligned} P(b_{n_i}^i \xrightarrow{rel} b_l^j | B) &\cong P(b_{n_i}^i \xrightarrow{rel} b_l^j | h_{n_i}^i, h_l^j, t_{n_i}^i, t_l^j, r_{n_i}^i, d_{n_i l}^{ij}, e_l^j, s_l^j) \\ &= \frac{F(b_{n_i}^i \xrightarrow{rel} b_l^j, h_{n_i}^i, h_l^j, t_{n_i}^i, t_l^j, r_{n_i}^i, d_{n_i l}^{ij}, e_l^j, s_l^j)}{F(h_{n_i}^i, h_l^j, t_{n_i}^i, t_l^j, r_{n_i}^i, d_{n_i l}^{ij}, e_l^j, s_l^j)}. \end{aligned} \quad (4.5)$$

## 4.4 解析実験

独話文の係り受け解析における本手法の有効性を評価するため，解析実験を行った．

### 4.4.1 実験の概要

実験には，2.3 節で構築した「あすを読む」構文構造付き音声独話コーパスを用いた<sup>5</sup>．実験で使ったデータを表 4.2 に示す．テストデータとして 500 文，学習データ

<sup>5</sup>4.2.3 節の分析で使った 200 文とは異なるものを用いた．

表 4.2: 実験データ（「あすを読む」構文構造付き音声独話コーパス）

	テストデータ	学習データ
番組	8	95
文	500	5,532
節境界単位	2,237	26,318
文節	5,298	65,821
形態素	13,342	165,129

として5,532文を用いた．なお，本手法では係り受けが節境界をまたぐことはないことを仮定しているが，この仮定に逸脱している係り受け関係はテストデータ中に152個存在した．これは，本手法の係り受け正解率（文末を除く）が96.8%（4,646/4,798）より高くなることはないことを意味する．なお，本実験では，フィラーが係り受け解析に悪影響を与えること，また，フィラーは音声認識器や形態素解析器により検出できる [4, 88] ことを考慮して，フィラーをあらかじめ「あすを読む」構文構造付き音声独話コーパスから取り除いた<sup>6</sup>．

本手法の有効性を比較評価するために，上述したデータを用いて以下の2つの手法で解析を行い，それぞれの解析時間と解析精度を求めた．

- 本手法 4.2.2 節，4.3 節でそれぞれ述べた，節境界解析，係り受け解析を順に行う手法．
- 従来手法 節境界単位に分割することなく，文の係り受け構造を一度に求める手法．ここで，文節間係り受け確率は式 (4.5) から属性  $e$  を排除した式により求め，1文の係り受け構造の尤度を文節間係り受け確率の積として計算し，最尤の構造を動的計画法により求める．この従来手法は藤尾ら [18] の先行研究に基づいている．

#### 4.4.2 実験結果

両手法の1文あたりの平均解析時間を表 4.3 に示す．節境界単位の解析手法の平均解析時間は，文単位の解析手法に比べて，約5倍向上した．ここで，本手法の解析時間には係り受け解析だけでなく節境界解析の時間も含まれているが，節境界解析プログラム CBAP の1文あたりの平均解析時間は1.2ミリ秒だった．従って，係

<sup>6</sup>フィラー以外の非文法的な言語現象に対しては，通常の文節と同じ扱いをしている．ただし，本実験で利用した「あすを読む」構文構造付き音声独話コーパスには，フィラー以外のこのような現象はほとんど見られなかった．

表 4.3: 実験結果 (解析時間)

	本手法	従来手法
平均解析時間 (msec)	10.9	51.9
使用計算機: Pentium4 2.4 GHz, Linux		
実装言語: LISP		

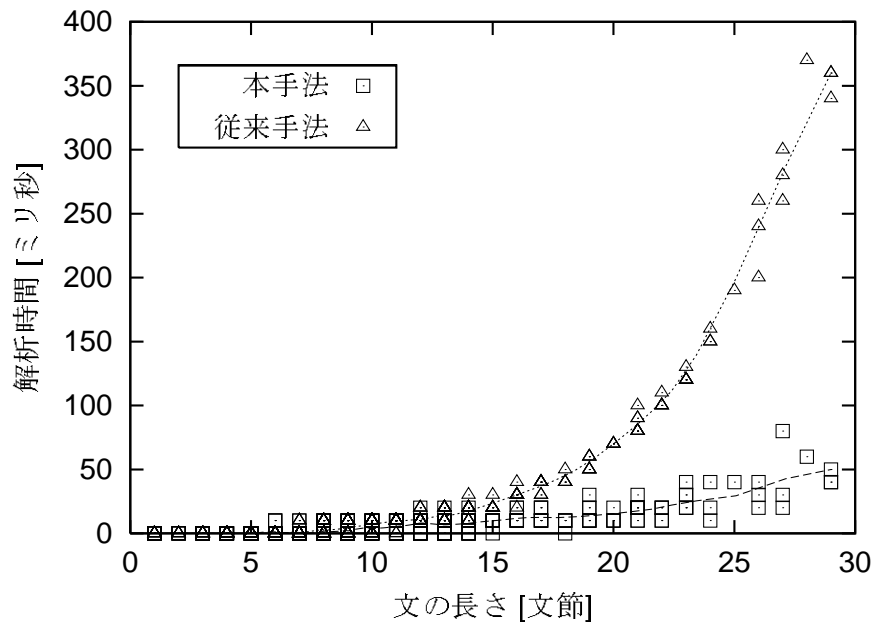


図 4.3: 文の長さ と 解析時間の関係

り受け解析の前処理としてCBAPを用いて節境界解析をする時間的な負担は無視できるほど小さい。文の長さ と 解析時間の関係を図 4.3 に示す。文単位の係り受け解析手法では文の長さが12文節を超えたあたりから、急激に解析時間が上昇するのに対し、節境界に基づく両手法の解析時間はあまり変化していない。実験で使用した6,032文の平均文節数は11.8であり、平均以上の長さをもつ独話文に対する本手法の効果を確認した。

次に、両手法の係り受け正解率を表 4.4 に示す。表 4.4 の第1行は、節境界単位末を除く節境界単位内の全ての文節に対する正解率を、第2行は、文末を除く全ての節境界単位末に対する正解率を示す。節境界単位の内部、最終文節とも、従来手法と比べて、本手法が高い解析精度を備えていることがわかる。

以上の結果から、本手法により、従来手法と比較して解析時間を大幅に短縮でき、解析精度も改善できることを確認した。

表 4.4: 実験結果 (係り受け正解率)

	本手法	従来手法
節境界単位内部	88.2% (2,701/3,061)	84.7% (2,592/3,061)
節境界単位の最終文節	65.6% (1,140/1,737)	63.3% (1,100/1,737)
全体	80.1% (3,841/4,798)	76.9% (3,692/4,798)

表 4.5: 節境界解析プログラム CBAP の解析結果

適合率	99.1% (2,190/2,209)
再現率	97.9% (2,190/2,237)

## 4.5 考察

本手法は、節境界を検出した後、係り受け解析を実行するため、節境界解析が失敗すると、その悪影響が係り受け解析まで及ぶと考えられる。また、本手法は、係り受けは節境界単位で閉じていることを仮定しているため、節境界単位末に係り文節とする係り受けを除き、節境界をまたぐ係り受けを正しく解析することはできない。

以下では、まず、4.5.1 節で節境界解析のエラーによる悪影響について考察する。次に、本手法の解析精度に対する効果を、節境界単位の内部と節境界単位の最終文節の2つに分けて、それぞれ、4.5.2 節と4.5.3 節で考察する。最後に、4.5.4 節で、本手法の問題点である、節境界単位内部の節境界をまたぐ係り受け関係の解析について検討する。

### 4.5.1 節境界解析エラーの影響

本手法で用いた節境界解析プログラム CBAP の解析結果を表 4.5 に示す。この精度は、節ラベルを考慮せず、節境界の位置のみで評価した値である。適合率、再現率ともに非常に高い結果となったが、解析エラーも存在した。まず、正解では節境界がない位置に間違っ検出された節境界が 19 個存在した。これらによって、次に行われる係り受け解析時に 4 個の係り受け関係を正解できなくなった。図 4.4 に、このエラーの例を示す。文節「新しい」と「枠組みを」の間に間違っ節境界を検出している。この誤検出された節境界によって、文節「枠組みを」が文節「貿易の」と異なる節境界単に所属することになり、係り受け解析時に、文節「貿易の」が文節「新しい」に係ると間違っ解析している。

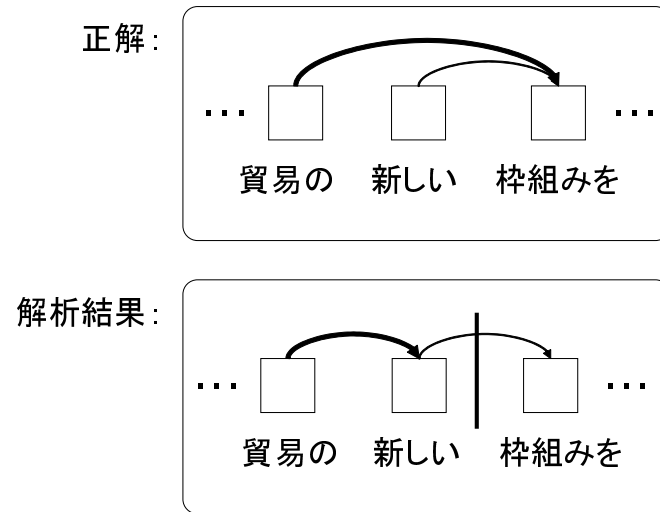


図 4.4: 節境界解析の適合率を低下させる節境界検出失敗によって生じた係り受け解析の失敗例

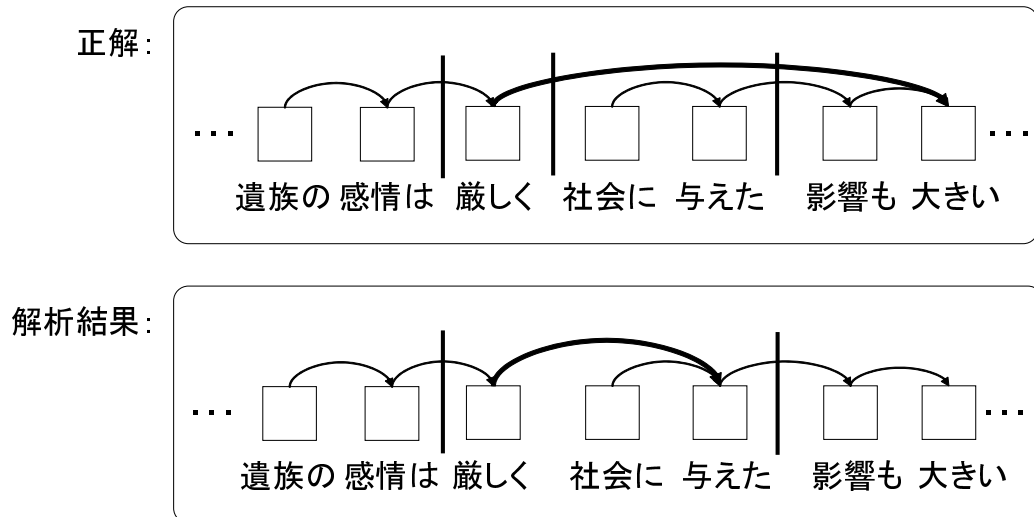


図 4.5: 節境界解析の再現率を低下させる節境界検出失敗によって生じた係り受け解析の失敗例

次に、正解では存在するが、節境界解析では検出できなかった47個の節境界について分析する。この47個の節境界を検出できなかったために、正しく同定することができなかった係り受けが27個あった。このような例を図4.5に示す。文節「厳しく」と「社会に」の間にあるはずの節境界を検出できなかったために、文節「厳しく」の係り先が節レベルの係り受け解析により解析される。そのため、正しい係り先である「大きい」は係り先の候補から外れることになり、係り受け解析が失敗している。

しかし、このような節境界解析のエラーによって、同定できなくなった係り受け

表 4.6: 本手法と従来手法の比較（節境界単位内部の係り受け解析結果）

従来手法 \ 本手法	正解	不正解	合計
	正解	不正解	合計
正解	2,499	93	2,592
不正解	202	267	469
合計	2,701	360	3,061

関係は全体の0.6% (31/4,790) にすぎず、後に行われる係り受け解析に悪影響を与えることはほとんどなかったと考えられる。

#### 4.5.2 節境界単位内部の文節に対する解析精度

表 4.6 は、節境界単位内部の係り受け関係の解析結果における両手法の正解、不正解の関係を示す。節境界単位内の係り受け関係 3,061 個のうち、両手法においてともに正しく解析された係り受け関係は 2,499 個であった。本手法で正解し、従来手法で不正解となったものは 202 個にのぼる。これは、本手法が受け文節の候補を節境界単位内に絞った効果を示している。

一方、従来手法のみで正しく解析できた係り受け関係は 93 個であった。このうち 46 個は、節境界をまたぐ係り受け関係であるため、本手法によりそもそも同定できない。このことは、節境界をまたぐ係り受け関係を除けば、従来手法で正しく解析される係り受け関係のほとんどを本手法によって正しく解析できることを意味する。

#### 4.5.3 節境界単位の最終文節に対する解析精度

表 4.4 が示すように、節境界単位の最終文節（文末を除く）の係り受け正解率は、節境界単位内部の解析と比べて両手法ともかなり低い。これは、節境界単位の最終文節を係り文節とする係り受け関係の同定が難しいことを意味している。

表 4.7 に、節境界単位の最終文節（文末を除く）の解析結果における両手法の正解、不正解の関係を示す。節境界単位の最終文節を係り文節とする係り受け関係 1,737 個のうち、両手法でともに正しく解析された係り受け関係は 1,037 個あった。本手法でのみ正解した係り受け関係は 103 個で、従来手法でのみ正しく解析できた係り受け関係 63 個を上回った。これは、本手法が、先に解析した節境界単位内部の係り受け構造を前提とすることにより、節境界単位の最終文節の受け文節となる候補を効果的に絞った結果であると考えられる。

表 4.7: 本手法と従来手法の比較（節境界単位の最終文節に対する解析結果）

従来手法 \ 本手法	本手法		合計
	正解	不正解	
正解	1,037	63	1,100
不正解	103	534	637
合計	1,140	597	1,737

表 4.8: 節境界をまたぐ係り受けに対する解析精度

	本手法	従来手法
適合率	11.8% (2/ 17)	25.3% (46/182)
再現率	1.3% (2/152)	30.3% (46/152)

#### 4.5.4 節境界をまたぐ係り受け関係

表 4.8 に，節境界をまたぐ係り受け関係に対する両手法の正解率を示す．本手法は，このような係り受け関係は存在しないとして解析を行っているため，1 つも正しく解析できない．なお，実験結果では，2 個の節境界をまたぐ係り受け関係を同定しているが，これは，節境界解析の段階で誤った節境界が付与されたために同定できたものである．一方，従来手法は，テストデータに存在する節境界をまたぐ係り受け関係 152 個のうち 46 個を正しく解析した．従来手法においてもその正解率は 30.3%にとどまっており，このような係り受け関係の同定はそもそも困難であることがわかる．

解析精度の向上のために，節境界をまたぐ係り受け関係を考慮した処理が望まれる．以下では，人手で正解を付与したテストデータ 500 文を用いてその解析可能性について検討する．図 4.6 に，節境界をまたぐ係り受け関係が存在した節境界単位の種類<sup>7</sup>とその割合を示す．「主題八」が最も多く，次いで，「連体節」，「テ節」の順であった．以下では，全体の 70.4%を占めるこの上位 3 つの節についてそれぞれ述べる．

##### 節境界単位の種類「主題八」

節境界をまたぐ係り受け関係 152 個のうち 42 個は，節境界単位「主題八」にその係り文節が存在した．節境界単位「主題八」は「述語を中心としたまとまり」とい

<sup>7</sup> 節境界単位の終端境界に付与された節境界のラベル名をその節境界単位の種類とする．

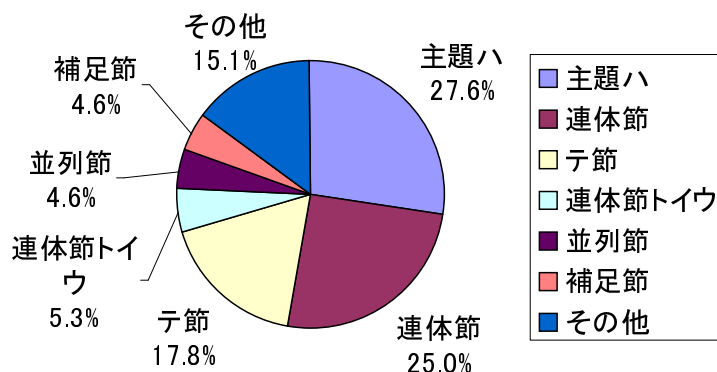


図 4.6: 係り受けがまたぐ節境界の種類とその割合

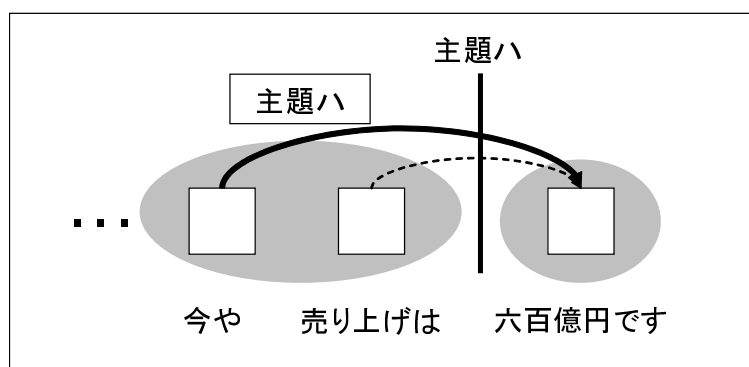


図 4.7: 節境界「主題ハ」をまたぐ係り受けの例

う節の定義に逸脱しているが、構文的に大きな切れ目になると考え [55]，本研究ではこれについても節境界単位としている。

このような節境界単位を調べてみると、『節境界単位「主題ハ」内に述語が存在しないために、述語に係るような文節は節境界単位外に位置する述語に係る現象』が多く見られた。この場合、述語に係る文節については節境界単位外に係り先があるとみなし、そのような規則を作成し検出することが考えられる。

例として、図 4.7 にこのような係り受けを含む文の一部を示す。この例では、副詞句である文節「今や」が、節境界単位「今や売り上げは」の外に位置する述語「六百億円です」に係っている。

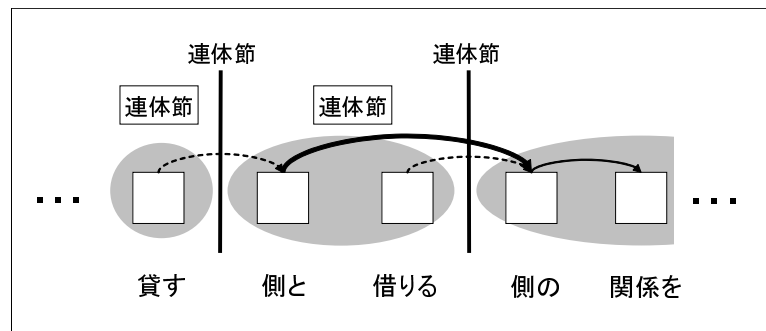


図 4.8: 節境界「連体節」をまたぐ係り受けの例

#### 節境界単位の種類「連体節」

節境界をまたぐ係り受け関係のうち 38 個は、節境界単位「連体節」にその係り文節が存在した。これらを調べてみると、『節境界単位内部の文節がこの連体節が修飾する文節と並列関係や同格関係になっている現象』が多く見られた。一般に、係り受け解析において、並列関係や同格関係の同定は難しいが、これらを検出する手法が報告されており（例えば，[1, 48]），本手法への導入が考えられる。

図 4.8 に例を示す。この例では、「(貸す) 側と」と「(借りる) 側の」が並列関係としての係り受け関係にあり、節境界「連体節」をまたいでいる。

#### 節境界単位の種類「テ節」

節境界をまたぐ係り受け関係のうち 27 個は、節境界単位「テ節」にその係り文節が存在した。これらの中で多く見られたのは、『文全体の係り受け構造としては、節境界をまたぐ係り受け関係になるが、節境界単位内部にも意味的には受けとなる文節が存在する現象』である。この場合、「係り先は唯一である」という制約を緩めて柔軟に評価する、すなわち、節境界単位内部にある受け文節についても正解とすることが考えられる<sup>8</sup>。

この例を図 4.9 に示す。この例では「検察側が」が文全体での係り受け構造としては「上告しておりました」に係るが、同一節境界単位内の「求めて」の主語は「検察側が」であり、そのような係り受け関係が必ずしも誤りであるというわけではない。

<sup>8</sup>2.3.2 節で述べた意味的係り受け関係は、このような評価に利用できると考えられる。

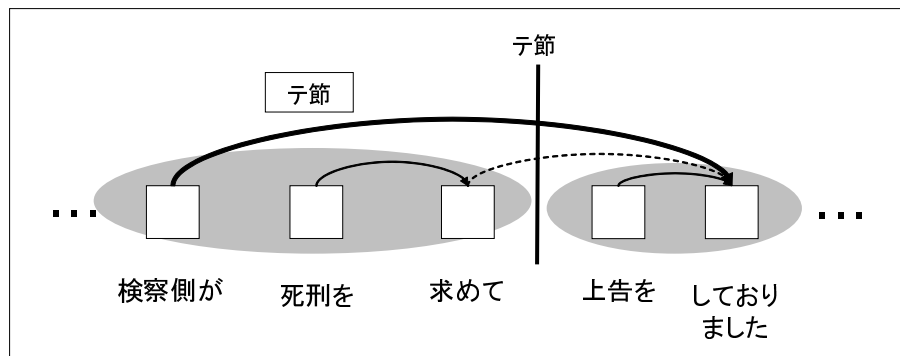


図 4.9: 節境界「テ節」をまたぐ係り受けの例

## 4.6 関連研究

独話文は文の長さが長く文の構造が複雑であるといった特徴があり，このような特徴を考慮することが重要である．長文の構文解析に関する研究は，書き言葉を対象とした研究においていくつか行われており，長文の構文解析の曖昧性を解消するために，「節」に着目した手法がいくつか提案されている．

初めに，並列節に着目した研究として，並列構造を構成する文節列の類似性を動的計画法を用いて検出し，並列関係を明確にすることにより，長文解析の精度を高める試みがある [1, 48]．次に，日本語解析の研究では，複文を構成する従属節のうち，用言に係る連用修飾節は必ずしもすぐ次の節の用言にかかるとは限らないという曖昧性に着目した研究が存在する [84, 94]．これらは，従属節の末尾に現れる表層表現から従属節を階層的に分類し，その順序関係を利用することにより，従属節の最終文節に対する解析精度を改善している．さらに，各従属節の主語を特定し，同一主語を持つ従属節をまとめた単位 “S(ubject)-clauses” に文を分割し，構文解析の曖昧性の解消を目指した研究もある [40]．これらの研究は，特定の節に生じる構文解析の曖昧性を軽減することにより，文全体の解析精度の改善を試みている．

一方，本手法では，特定の種類に限定することなく，全ての節を利用する．文をあらゆる節に網羅的に分割し，その節の中で局所的に構文解析を行うことにより，節内部の文節の係り先の候補を絞ることができ，より高精度な節内部の係り受け解析を実現している．また，これにより，解析時間の短縮も実現している．さらに，本手法は，節ごとに係り受け解析を実行できるので，漸進的な解析への応用が期待できる．なお，漸進的解析への応用については，5章で述べる．

## 4.7 4章のまとめ

本章では，節境界に基づく独話文の係り受け解析手法を提案した．本手法は，文を節に分割することにより，独話文の係り受け解析の効率化と高精度化を実現する．本手法の有効性を評価するために，「あすを読む」構文構造付き音声独話コーパスを用いて係り受け解析実験を行った．実験の結果，本手法が，従来の文単位係り受け解析手法と比べて，解析精度を向上させ，さらに，解析時間を大幅に短縮できることを確認した．

今後の課題として，以下の点が挙げられる．

- 節境界をまたぐ係り受け関係を検出し，その係り先を同定する手法について検討する．
- 4.5.4 節で述べた節境界「テ節」をまたぐ係り受け関係に対して，意味的係り受け関係（2.3.2 節参照）を用いた柔軟な評価を行う．
- 節境界単位の最終文節に対する解析精度の改善を図る．具体的には，節の階層構造における「従属節の選好性」を利用して，節末の係り先の解析精度を改善した先行研究 [84, 94] を本手法に組み込むことを検討する．
- 3章で提案した音声対話文に対する頑健な係り受け解析手法と本手法の結合を図り，フィラーなどの非文法的言語現象も取り扱う．
- 本章で述べた実験では，独話中の長文に対する本手法の有効性を示しているが，本手法は書き言葉の長文に対しても同様に有効であると考えられる．そこで，特許文など書き言葉における長文に対しても実験を行い，その有効性を検証する．



## 第5章 日本語話し言葉の漸進的な係り受け解析

### 5.1 はじめに

音声同時通訳や自動字幕生成のように、音声を入力と同時に処理するような音声言語処理システムでは、話者による音声入力に従って順次、解析することが求められる。特に、独話を対象とする場合、一般に、文が長くなる傾向にあり、従来の言語処理技術と同様に文単位での逐次的な解析を行うと、入力に対する出力の同時性が損なわれることになるため、漸進的な解析が必須となる。実際、これまでに漸進的な構文解析に関する研究がいくつか行われており（例えば、[34, 35, 72, 80]）、ここでは、解析処理の単位、すなわち、どのような言語単位ごとに処理を実行し、結果を出力するのが問題となる。構文解析の漸進性と正確さの双方を満たすために、文より短く、かつ、構文的なまとまりを備えた言語単位を解析処理の単位として採用することが望ましい。

そこで本章では、節を解析単位とする独話の漸進的係り受け解析手法を提案する。節は単文に相当し、構文的にも意味的にもまとまった言語単位であるとともに、文が長くなりがちな複文や重文は、複数の節に分割できることから、節は、漸進的な解析処理の単位として適している。特に、長文の構文的曖昧性を軽減することを目的に、節に着目した研究がいくつか行われ、構文解析の精度 [40, 84, 94] や機械翻訳等の自然言語処理応用システムの性能 [41, 50, 90] の向上が報告されており、節ごとの処理による解析の正確さへの効果が期待できる。

本手法では、独話音声に対して、節が入力されるたびにその節の内部の係り受け構造を作り上げるとともに、すでに入力されている節の係り先を決定することを試みる。節の係り先となる文節の決定は、後続するいくつかの文節との係り受けの尤度を考慮した動的なタイミングで行う。これにより、独話の入力途中の段階で構造情報を随時出力する漸進的な解析が可能となる。

さらに、本手法では、文境界が付与されていない独話データ全体に対してその係り受け構造を解析する。これは、独話には明示的な文末標識がなく、あらかじめ文

単位に区切ることは容易ではないという独話の特徴に対応している．独話データを用いた解析実験の結果，本手法により，文を解析単位とした係り受け解析手法と同等の解析精度を維持しつつ，係り受け解析の漸進性を実現できることを確認した．

本章の構成は以下の通りである．次節で独話の漸進的係り受け解析における処理単位について述べ，5.3節で節境界に基づく漸進的係り受け解析手法を示す．5.4節で漸進的係り受け解析アルゴリズムについて説明し，5.5節で解析実験について述べる．5.6節で提案手法の文末検出性能について考察し，5.7節で本章のまとめと今後の課題について述べる．

## 5.2 独話の漸進的係り受け解析における処理単位

本章の研究では，解析の処理単位として節を採用し，節が入力されるたびにその時点までの係り受け構造を可能な限り決定する漸進的な独話係り受け解析システムを実現する．以下では，4.2.2節で導入した，節の近似的な単位である節境界単位を独話の漸進的係り受け解析における処理単位とすることについて検討を与える．

### 5.2.1 節と節境界単位

節とは，述語を中心としたまとまりであり，複文や重文の場合，文は複数の節から構成される．さらに，節は，構文的，意味的にまとまった単位であるため，文に代わる解析単位として利用できると考えられる．

節を単位とした漸進的係り受け解析とは，節が入力されるたびに，それまでの入力に対する解析結果を出力することを意味し，そのためには，独話の入力と同時的に節への分割を実行できる必要がある．しかし，複文において従属節が主節に埋め込まれる場合など，構文的な解析の前処理として節を漸進的に検出することは必ずしも容易ではない．

そこで，本章の研究では，節への漸進的な分割を近似的に実現するため，4.2.2節で述べた節境界単位を導入する．節境界単位とは，節境界解析 [55] が特定する節の終端境界によりはさまれた単位である．節境界解析 [55] は，形態素の局所的な接続パターンのみを手がかりとしているため，節境界単位への漸進的な分割が可能である．節境界単位は，埋め込み節がある場合には，節と一致しないものの，それ以外の場合であれば，節と完全に一致しており，節の近似的な単位であるといえる．

### 5.2.2 節境界単位の分析

漸進的係り受け解析の漸進性と正確さの双方を実現するために、その処理単位が、節同様、文より短く、かつ、構文的にまとまっていることが望まれる。そこで、本節では、節境界単位がこのような性質を持ち合わせているかについて調査した。具体的には、文と比較した節境界単位の長さ、また、係り受けが節境界単位でどの程度閉じているのかを、実際の独話データを用いて分析した。

分析には、4.2.3 節で用いた「あすを読む」構文構造付き音声独話コーパス 200 文を用いた。まず、節境界単位と文の長さに着目して分析した。文の平均文節長は 12.2 であるのに対して、節境界単位の平均文節長は 2.6 であった。これは、文ごとに解析する場合と比べ、節境界単位ごとに解析することによって大幅に解析の漸進性が改善されることを示唆している。次に、節境界単位と係り受け構造の関係について述べる。総文節数 2,430 文節のうち、節境界単位の最終文節 (951 文節) を除いた 1,479 文節の中で、94 文節のみが節境界単位の外に位置する文節に係っていた。これは、全体の 93.6% (1,385/1,479) の係り受け関係が節境界単位で閉じていることを意味しており、節と同様にある程度構文的にまとまった単位であることを示している。以上の分析結果から、節境界単位が漸進的係り受け解析の処理単位として利用可能であることを確認した。

### 5.2.3 漸進的係り受け解析の処理単位としての節境界単位

前節の分析結果に基づき、本章では、「独話は 1 つ以上の節境界単位の接続であり、各節境界単位を構成する文節は、節境界単位の最終文節を除き、その節境界単位の内部の文節に係る」とみなして、それに基づく漸進的係り受け解析手法を提案する。

例として、図 4.1 の独話文「先日総理府が発表いたしました世論調査によりますと死刑を支持するという人が八十パーセント近くになっております」における節境界単位と係り受けの関係を図 5.1 に示す<sup>1</sup>。この文は 4 つの節境界単位「先日総理府が発表いたしました」、「世論調査によりますと」、「死刑を支持するという」、「人が八十パーセント近くになっております」から構成され、各節境界単位が係り受け構造を形成し、それらが節境界単位の最終文節からの係り受け関係でつながっている。

<sup>1</sup> 図 4.1 は節と係り受けの関係を示しているのに対して、図 5.1 は節境界単位と係り受けの関係を示している。

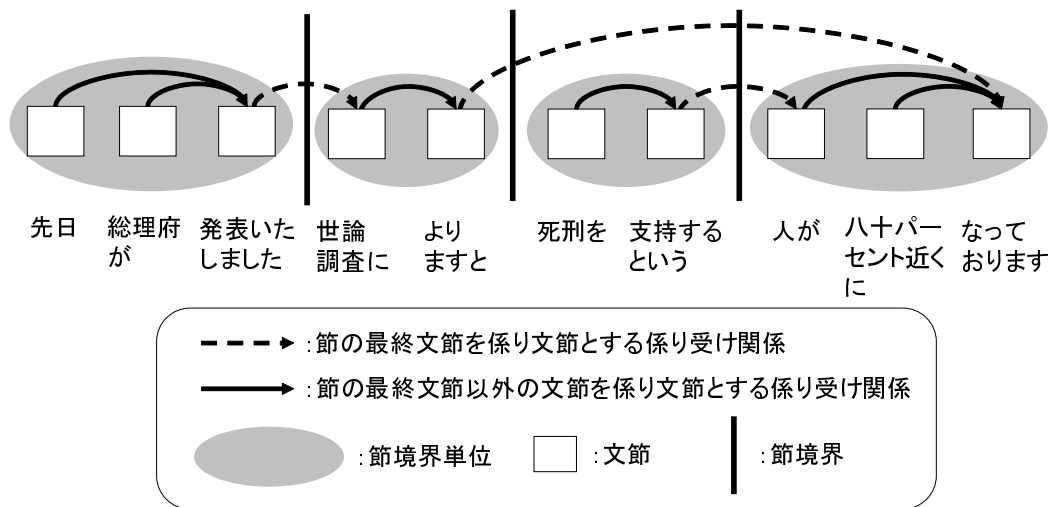


図 5.1: 節境界単位と係り受けの関係

### 5.3 節境界に基づく漸進的係り受け解析

本節では、節境界単位に基づく漸進的係り受け解析について述べる。本手法では、音声入力に対して節境界を随時判定し、節境界単位が同定されると、その時点までの入力に対して係り受け解析を実行する。節境界の判定は節境界解析 [55] により実行する。係り受け解析は、入力された節境界単位の内部の係り受け構造を解析するとともに、すでに入力された節境界単位の最終文節の係り先を可能であれば決定する。

本手法では、形態素解析、文節まとめあげ、及び節境界解析が施された 1 独話を入力とする。ここで、入力データには、文境界が付与されていないことに注意されたい。解析の手順は以下の通りである。なお、具体的なアルゴリズムは 5.4 節で述べる。

#### 1. 節レベルの係り受け解析

1 独話中のすべての節境界単位に対して、その内部の係り受け構造を解析する。

#### 2. 独話レベルの係り受け解析

1 独話中のすべての節境界単位に対して、その最終文節の係り先を解析する。

節レベルの係り受け解析は 4.3.1 節で述べた解析手法を用いた。なお、節レベルの係り受け解析では、係り受けの後方修飾性、係り先の唯一性、非交差性の 3 つの性質を絶対的制約とする<sup>2</sup>。一方、独話レベルの係り受け解析は、4.3.2 節で述べた文レ

<sup>2</sup>4 章の研究と同様に、本章の研究でも、放送で行われる比較的整った解説における話し言葉を対象としており、このような話し言葉には構文的制約に逸脱する係り受けはほとんど見られないと仮定する。

ベルの係り受け解析とは異なるため，次の 5.3.1 節で詳述する．漸進的な係り受け解析を行っている本章の研究では，前章の研究とは異なり，文の概念がなく，独話は節境界単位の列で構成されるとみなす．そのため，4.3.2 節の文レベルの係り受け解析で利用していた，受け文節が文末であるか否かの素性  $s$  を利用していない．また，独話レベルの係り受け解析では，3 つの構文的制約のうち，係り受けの後方修飾性と非交差性の 2 つを絶対的制約とする．すなわち，係り先の唯一性を緩和し，係り先がない文節を許す<sup>3</sup>．

なお，以下では，1 独話を構成する節境界単位列を  $C_1 \cdots C_m$ ，節境界単位  $C_i$  を構成する文節列を  $b_1^i \cdots b_{n_i}^i$ ，文節  $b_k^i$  を係り文節とする係り受け関係を  $dep(b_k^i)$ ，1 独話の係り受け構造を  $\{dep(b_1^1), \dots, dep(b_{n_m-1}^m)\}$  と記す．

### 5.3.1 独話レベルの係り受け解析

節境界単位の最終文節の受け文節を同定する．1 独話の文節列を  $B(= B_1 \cdots B_m)$  とし，節境界単位の最終文節を係り文節とするような係り受け構造  $\{dep(b_{n_1}^1), \dots, dep(b_{n_m-1}^m)\}$  を  $S_{last}$  とするとき， $P(S_{last}|B)$  を最大とする  $S_{last}$  を求める．係り受け関係は互いに独立であると仮定すると， $P(S_{last}|B)$  は以下の式で計算できる．

$$P(S_{last}|B) = \prod_{i=1}^{m-1} P(b_{n_i}^i \xrightarrow{rel} b_l^j | B) \quad (5.1)$$

ここで， $P(b_{n_i}^i \xrightarrow{rel} b_l^j | B)$  は，1 独話の文節列  $B$  が与えられたときに， $C_i$  の最終文節  $b_{n_i}^i$  が  $b_l^j$  に係る確率をあらわす．最尤の係り受け構造は，式 (5.1) の確率を最大とする構造であるとして動的計画法を用いて計算する．また，4.3.2 節と同様に，先に解析した節境界単位内部の係り受け構造を前提として決定する．すなわち，後方に位置するすべての文節を受け文節の候補として計算するのではなく，節境界単位内部の係り受け構造から非交差性を満たすものだけを受け文節の候補とする．図 5.1 の場合，文節「支持するという」の受け文節は「人が」または「なっております」のいずれかであるとして計算する．

次に， $P(b_{n_i}^i \xrightarrow{rel} b_l^j | B)$  の計算について述べる．係り文節における自立語の基本形を  $h_{n_i}^i$ ，その品詞を  $t_{n_i}^i$ ，係りの種類を  $r_{n_i}^i$  とし，受け文節における自立語の基本形を  $h_l^j$ ，その品詞を  $t_l^j$  とする．また，文節間距離を  $d_{n_i l}^{ij}$  とする．ここで，係りの種類とは，係り文節が付属語を伴う場合は文節末の形態素の語彙，品詞，活用形であり，

<sup>3</sup>独話レベルの係り受け解析では，独話全体の節境界単位の最終文節の係り先を同定する．このとき，本来文末であるとされる文節は係り先がないとして解析する．

そうでない場合は文節末の形態素の品詞，活用形である（3.3.2 節参照）。なお，これらの属性は，従来の係り受け解析手法 [18, 46, 92] で用いられてきたものと同様である。

さらに，独話レベルの解析においても，受け文節が節境界単位の最終文節であるか否かを示す属性  $e_l^j$  を導入した。5.2.2 節の 200 文を分析した結果，節境界単位の最終文節も，その 70.6%(522/751) が別の節境界単位の最終文節に係ることがわかったためである。

以上の属性を用いて，確率  $P(b_{n_i}^i \xrightarrow{rel} b_l^j | B)$  を以下のように計算する。

$$\begin{aligned}
 & P(b_{n_i}^i \xrightarrow{rel} b_l^j | B) \\
 & \cong P(b_{n_i}^i \xrightarrow{rel} b_l^j | h_{n_i}^i, h_l^j, t_{n_i}^i, t_l^j, r_{n_i}^i, e_l^j, d_{n_{il}}^{ij}) \\
 & = \frac{F(b_{n_i}^i \xrightarrow{rel} b_l^j, h_{n_i}^i, h_l^j, t_{n_i}^i, t_l^j, r_{n_i}^i, e_l^j, d_{n_{il}}^{ij})}{F(h_{n_i}^i, h_l^j, t_{n_i}^i, t_l^j, r_{n_i}^i, e_l^j, d_{n_{il}}^{ij})} \quad (5.2)
 \end{aligned}$$

ただし， $F$  は共起頻度関数である。

なお，4.3.1 節と同様に，本手法においても，式 (5.2) により  $P(b_{n_i}^i \xrightarrow{rel} b_l^j | B)$  を計算するときに起こるデータスパースネスの問題を解決するために，藤尾ら [18] のスムージング手法を用いている。すなわち，式 (5.2) 中の  $F(h_{n_i}^i, h_l^j, t_{n_i}^i, t_l^j, r_{n_i}^i, e_l^j, d_{n_{il}}^{ij})$  が 0 である場合は，次式 (5.3) を用いて  $P(b_{n_i}^i \xrightarrow{rel} b_l^j | B)$  を計算する。

$$\begin{aligned}
 & P(b_{n_i}^i \xrightarrow{rel} b_l^j | B) \\
 & \cong P(b_{n_i}^i \xrightarrow{rel} b_l^j | t_{n_i}^i, t_l^j, r_{n_i}^i, e_l^j, d_{n_{il}}^{ij}) \\
 & = \frac{F(b_{n_i}^i \xrightarrow{rel} b_l^j, t_{n_i}^i, t_l^j, r_{n_i}^i, e_l^j, d_{n_{il}}^{ij})}{F(t_{n_i}^i, t_l^j, r_{n_i}^i, e_l^j, d_{n_{il}}^{ij})} \quad (5.3)
 \end{aligned}$$

## 5.4 漸進的係り受け解析アルゴリズム

上述した 2 つのレベルの解析のうち，節境界単位内部の係り受け解析は，4.3.1 節で述べた方法により解析すればよい。それに対して節境界単位の最終文節に対する係り受け解析は，その受け文節がいつ入力されるかはあらかじめ明らかであるわけではないため，それを決定するタイミングが問題となる。本章の研究では，節境界単位の最終文節が入力されてからその受け文節が入力されるまでが格段に長くなることは稀であると考え，ある程度解析が進んだ時点でその受け文節を決定することとした。具体的には，節境界単位が入力されるたびにその時点での最尤の係り受け構造を 5.3.1 節で述べた方法により解析し，ある最終文節の受け文節が変化しなかつ

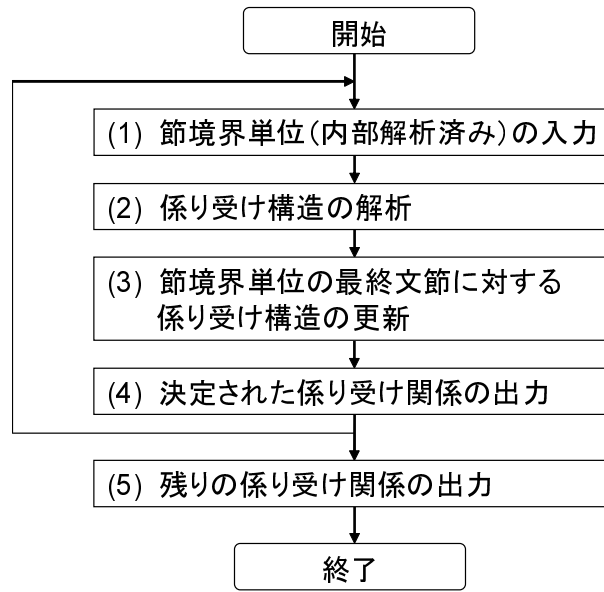


図 5.2: 漸進的係り受け解析の流れ

た回数（不変回数）がある閾値に達したとき，その文節をこの最終文節の受け文節として決定する．以下では，この閾値を不変化閾値と呼ぶ．

本節の以下では，節境界単位の最終文節に対する係り受け解析について説明する．

#### 5.4.1 独話レベルの漸進的解析アルゴリズム

係り受け解析の流れを図 5.2 に示す．解析では，節境界単位  $C_i$  が入力されるごとに，すでに入力された節境界単位  $C_1, \dots, C_{i-1}$  の各最終文節  $b_{n_1}^i, \dots, b_{n_{i-1}}^{i-1}$  に対する係り受け構造  $D = \{(dep(b_{n_j}^j), k) \mid 1 \leq j \leq i-1\}$  を更新することにより実行する．ここで  $k$  は  $dep(b_{n_j}^j)$  の不変回数を示す．以下に係り受け解析アルゴリズムを示す．なお，不変化閾値を  $\lambda$  とする．

- (1) 内部の係り受け構造が決定された節境界単位  $C_i$  を入力する．
- (2) 節境界単位の最終文節のうち，係り先が未決定な文節に対して，それを係り文節とする係り受け関係を 5.3.1 節で説明した方法により求める．
- (3) (2) で生成された係り受け関係  $dep(b_{n_j}^j)$  に基づき，最終文節に対する係り受け関係  $D$  を更新する． $dep(b_{n_j}^j)$  が同一の場合は不変回数を  $k+1$  とし，異なる場合は 1 とする．

- (4)  $k = \lambda$  を満たす係り受け関係  $(dep(b_{n_j}^j), k) \in D$  に対して, 文節  $b_{n_j}^j$  の係り先が決定したとして  $dep(b_{n_j}^j)$  を出力する.
- (5) すべての節境界単位が入力された時点で,  $k < \lambda$  の  $(dep(b_{n_j}^j), k) \in D$  に対して, その係り受け関係  $dep(b_{n_j}^j)$  を出力する.

なお, 本手法では, 文末は係り先がないとして解析する. そのため, 節境界単位末の解析では係り先なしを候補に含める. 具体的には, 式 (5.1) において, 係り先のない文節はそれ自身に係る (すなわち,  $b_{n_i}^i = b_l^j$ ) とし, 係り先なしとなる確率も計算する.

### 5.4.2 解析例

図 5.3 に, 独話「正当な理由がない限り契約期間が切れたといっても明け渡しを請求できない点にあるといわれています」の節境界単位末の文節の係り先を解析する様子を示す. (a) ~ (f) の 6 つの過程から構成され, それぞれ上部に係り受け構造を, 下部に節境界単位の最終文節の係り受け関係を示す.  $(dep(b_{n_j}^j), k) \in D$  の  $dep(b_{n_j}^j)$  が係り文節及び受け文節に,  $k$  が不変回数に, それぞれ相当する. なお, ここでは不変変化閾値が 3 であるとして説明する.

(a) は, 最初の節境界単位 I が入力された状態を, (b) は, 節境界単位 II が入力され, 係り受け構造  $\{dep(限り)\}$  が解析された状態を示す.  $dep(限り)$  は上部の点線矢印に相当し, 「限り」の係り先が「切れた」であり, 不変回数は 1 であることが下部に記録される. 同様にして, (c), (d) は, それぞれ節境界単位 III, IV が入力されたときの最尤の係り受け構造  $\{dep(限り), dep(切れた)\}$ ,  $\{dep(限り), dep(切れた), dep(いっても)\}$  が解析された状態を示す.

(e) は, 節境界単位 V が新たに入力され, 最尤の構造  $\{dep(限り), dep(切れた), dep(いっても), dep(請求できない)\}$  が求まった状態を示している. このとき, 係り受け関係  $dep(切れた)$  の不変回数が不変変化閾値として設定した 3 に達したため, この関係を決定し出力する.

(f) は, 節境界単位 VI が新たに入力され, 最尤の係り受け構造  $\{dep(限り), dep(いっても), dep(請求できない), dep(あると)\}$  が求まった状態を示す. (e) と同様に不変回数が不変変化閾値に達している係り受け関係  $dep(限り), dep(いっても)$  を決定し出力する.

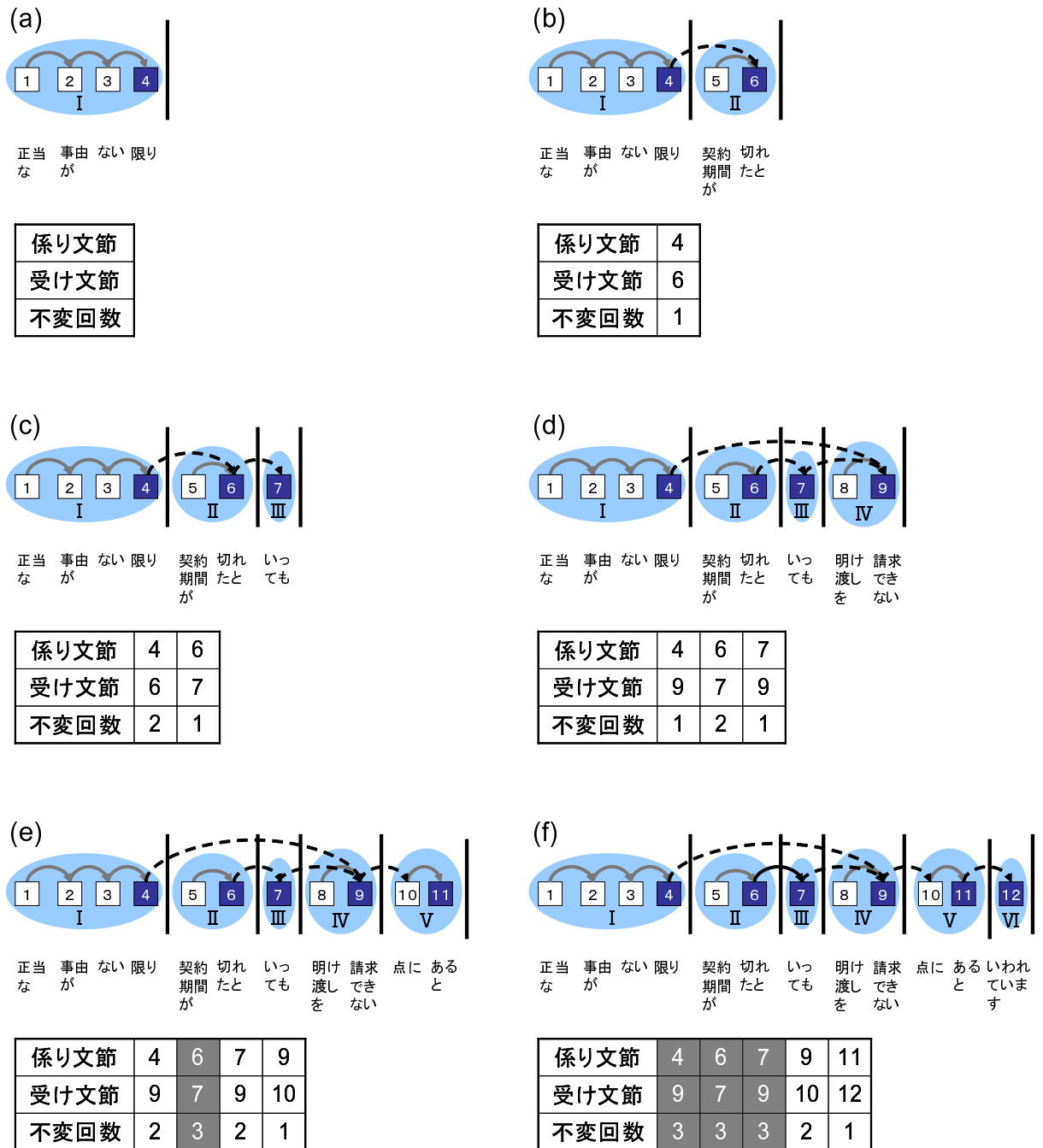


図 5.3: 漸進的係り受け解析の例 (不変化閾値 3 の場合)

表 5.1: 実験データ（「あすを読む」構文構造付き音声独話コーパス）

	テストデータ	学習データ
番組	7	95
文	470	5,532
節境界単位	2,140	26,318
文節	5,054	65,821
形態素	12,753	165,129

## 5.5 解析実験

独話の漸進的係り受け解析における本手法の有効性を評価するため，解析実験を行った．

### 5.5.1 実験に使用したデータ

実験には，2.3 節で構築した「あすを読む」構文構造付き音声独話コーパスを使用した<sup>4</sup>．使用したデータの概要を表 5.1 に示す．テストデータとして 7 番組（470 文），学習データとして 95 番組（5,532 文）をそれぞれ用いた．また，節境界をまたぐ係り受け関係は，テストデータ中に 145 個存在した．これは，本手法の係り受け正解率（番組末を除く）が 97.1% (4,902/5,047) より高くなることはないことを意味する．

### 5.5.2 実験の概要

上述したデータを用いて本手法により解析を行った．5.4 節で説明した不変化閾値を 1 から 12 まで変化させて計 12 回実験し，解析精度，解析時間，解析の漸進性を評価した．

解析精度は係り受け正解率を，解析時間は文節が入力されるごとに実行される解析に要した時間の合計を，それぞれ求めた．また，解析の漸進性を評価するために，遅延時間という尺度を定義し，これを測定した．日本語の係り受けにおける後方修飾性のために，係り受け関係はその受け文節が入力されるまで決定できない．そこで，遅延時間を，受け文節が入力されてからどの程度遅れてその係り受け関係が決定されたかを示す指標として導入することとし，以下の通り定義する．

<sup>4</sup>5.2.2 節の分析で使用した 200 文とは異なるものを用いた．

番組  $p$  の文節列  $b_1 \cdots b_{n_p}$  を解析するとき，ある文節  $b_x$  が入力された時点で解析システムが決定し出力した係り受け関係の集合を  $O_x(p)$  とし，係り受け関係  $\sigma \in O_x(p)$  の受け文節を  $f(\sigma)$  とする．また，ある文節  $b_x$  が入力されたときの時間を  $g(b_x)$  とする<sup>5</sup>．このとき，係り受け関係の決定に致る遅延時間の平均

$$Delay = \frac{\sum_p \sum_{x=1}^{n_p} \sum_{\sigma \in O_x(p)} \{g(b_x) - g(f(\sigma))\}}{\sum_p \sum_{x=1}^{n_p} |O_x(p)|} \quad (5.4)$$

を計算し，これを漸進性の評価値とする．なお，出力された係り受け関係  $\sigma$  の受け文節  $f(\sigma)$  は，文節  $b_x$  が入力された時点ですでに入力されており， $g(b_x) - g(f(\sigma))$  が負の値になることはない．

比較のため，節境界解析を行わず，文ごとに1文の係り受け構造を求める従来手法でも係り受け解析を行った．1文の文節列が与えられたときの文節間の係り受け確率は，式(5.2)，(5.3)の属性  $e_l^i$  を，受け文節  $b_l^i$  が文末であるか否かを示す属性  $s_l^i$  に変更した式を用いて計算する．実験で使用した書き起こしデータには，句点が付与されており，文ごとに解析可能である．なお，本手法により係り受け解析を行うときは，句点を取り除き1番組分の発話を連結した．

これらの解析システムをGNU Common LISPで実装し，CPUがPentium4 2.40GHz，メモリが2GBのLinux PC上で実験した．

### 5.5.3 実験結果

本手法の不変値閾値ごとの係り受け正解率を表5.2に示す．表5.2の第1列は，番組末を除く全ての節境界単位末に対する正解率を，第2列は，番組末を除く全ての文節に対する正解率を示す．不変値閾値が2及び3のときに，節境界単位末に対する正解率が最も高く，全体の正解率は76.2%となった．なお，節境界単位末を除く節境界単位内に対する解析の正解率は87.5% (2,551/2,914)であった<sup>6</sup>．表5.3に，CBAPの節境界解析の精度を示す．これは，CBAPによる節境界の検出性能を適合率，再現率により評価した結果である．適合率，再現率ともに高く，後に行われる解析へ

<sup>5</sup>「あすを読む」の書き起こしデータでは，200ms以上のポーズで区切られた発話ごとに発話時間を付与している．文節の発話時間をモーラ数から近似的に算出し，これを文節の入力時間とした．

<sup>6</sup>このうち，節境界単位末から2つ目の文節に対しては，正解率が97.7% (1,287/1,317)であった．本手法による節境界単位内部の解析では，節境界単位末から2つ目の文節は必ず節境界単位末の文節にかけることになるが，実際，ほとんどの節境界単位で係り受けは閉じているので，高い解析精度が得られたと考えられる．

表 5.2: 不変値閾値ごとの係り受け正解率

不変値閾値	節境界単位末	全体
1	57.6% (1,228/2,133)	74.9% (3,778/5,047)
2	<b>60.8%</b> (1,296/2,133)	<b>76.2%</b> (3,847/5,047)
3	<b>60.8%</b> (1,296/2,133)	<b>76.2%</b> (3,847/5,047)
4	60.4% (1,289/2,133)	76.1% (3,840/5,047)
5	59.8% (1,276/2,133)	75.8% (3,827/5,047)
6	59.4% (1,268/2,133)	75.7% (3,819/5,047)
7	58.8% (1,254/2,133)	75.4% (3,805/5,047)
8	58.6% (1,251/2,133)	75.4% (3,803/5,047)
9	58.7% (1,253/2,133)	75.4% (3,805/5,047)
10	58.4% (1,245/2,133)	75.2% (3,797/5,047)
11	57.6% (1,229/2,133)	74.9% (3,780/5,047)
12	57.9% (1,235/2,133)	75.0% (3,786/5,047)

表 5.3: CBAP の節境界解析結果

適合率	99.1% (2,088/2,106)
再現率	97.6% (2,088/2,140)

の影響は小さい．一方，従来手法の係り受け正解率は，470 の文末文節を除く全ての文節に対して評価し，76.1% (3,490/4,584) であった．このことから，本手法は，文境界情報を利用していないにもかかわらず，従来手法と同程度の解析精度で，係り受け構造を同定できることがわかる．

本手法の不変値閾値と 1 番組あたりの解析時間の関係を図 5.4 に示す．不変値閾値を大きくするにしたがって解析時間が増加している．解析時間が最も短かったのは，不変値閾値が 3 のときで，全 7 番組で 125.3 秒，1 番組あたり 17.8 秒だった．不変値閾値が 2 のときも，解析時間は 1 番組あたり 18.3 秒であり，ほとんど差がなかった．なお，この解析時間には，CBAP による節境界解析の時間も含まれている．節境界解析の解析時間は 1 番組あたり 3 秒程度である．一方，従来手法の解析時間は 1 番組あたり 6.4 秒であった．本手法は，従来手法と違い，文境界解析を行っていない独話全体の文節列を解析の対象にしているにもかかわらず，従来手法の 3 倍程度の時間で 1 番組を解析できていることが分かる．

図 5.5 に，本手法の不変値閾値と係り受け関係決定の平均遅延時間の関係を示す．この図では，係り文節が節境界単位の最終文節である場合と節境界単位内の文節で

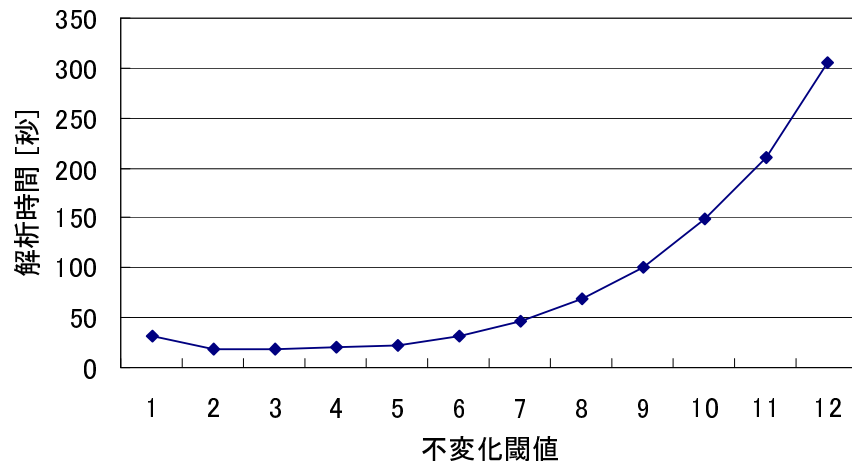


図 5.4: 不変化閾値と1番組あたりの解析時間の関係

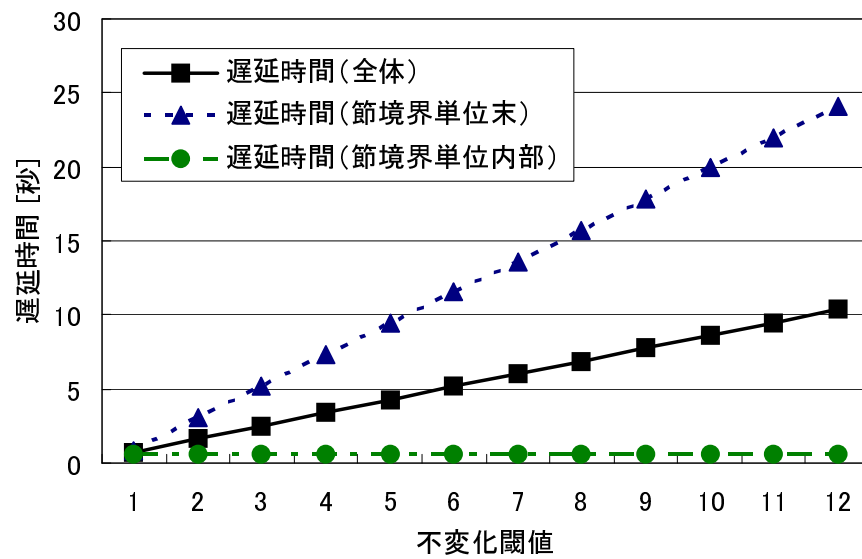


図 5.5: 不変化閾値と平均遅延時間の関係

表 5.4: 文末検出の適合率・再現率・F 値

不変値閾値	適合率	再現率	F 値
1	54.6% (337/617)	<b>72.1%</b> (334/463)	62.1
2	69.8% (312/447)	67.4% (312/463)	68.6
3	74.6% (296/397)	63.9% (296/463)	<b>68.8</b>
4	75.7% (278/367)	60.0% (278/463)	66.9
5	76.8% (271/353)	58.5% (271/463)	66.4
6	76.9% (260/338)	56.2% (260/463)	64.9
7	78.5% (252/321)	54.4% (252/463)	64.3
8	78.7% (247/314)	53.3% (247/463)	63.6
9	<b>81.1%</b> (249/307)	53.8% (249/463)	64.7
10	80.3% (245/305)	52.9% (245/463)	63.8
11	79.9% (238/298)	51.4% (238/463)	62.6
12	79.8% (233/292)	50.3% (233/463)	61.7

ある場合に係り受け関係を分類し、それぞれの平均遅延時間も示している。最も正解率が高かった、不変値閾値が2や3のときの平均遅延時間は、それぞれ1.6秒、2.5秒であった。図から、遅延時間のほとんどが、節境界単位の最終文節を係り文節とする係り受け関係の解析遅延によるものであることがわかる。実際、節境界単位内の文節を係り文節とする係り受け関係の平均解析時間は、不変値閾値の値に関係なく0.5秒程度である。一方、従来手法の平均遅延時間<sup>7</sup>は、3.2秒であった。このことから、本手法は、従来手法とくらべて出力タイミングの漸進性を大幅に改善できることがわかる。

以上の結果から、本実験においては不変値閾値が2もしくは3のとき、最も高い性能を示しており、文単位を入力とする従来手法と比較して、同程度の解析精度を達成し、なおかつ、解析の漸進性を大幅に改善していることを確認した。

## 5.6 文末検出性能

本手法では、文境界が付与されていない独話データを一度に解析する際、本来文末であるとされる文節は係り先がないとして解析を実行している。すなわち、係り先なしと解析された文節を文末であるとみなすことができる。このような観点から、本手法の文末検出性能を評価した。表5.4に文末検出の適合率、再現率、F値を示す。不変値閾値3のときに最も高いF値を示した。独話の文境界検出手法はこれま

<sup>7</sup>1文ごとに解析されるので、1文中のどの係り受け関係も文末の入力時点で決定されるとする。

でもいくつか提案されている [85, 89] . 本手法は, これらと比べ, 精度において劣っているものの, 漸進的係り受け解析と同時的に文末を検出できるという特徴がある. また, 解析精度についても, 従来研究 (例えば, [85]) において, ポーズ長やピッチ, パワーなどの音響情報を利用する効果が報告されており, 今後, これらを利用することにより, 精度向上が期待できる.

## 5.7 5章のまとめ

本章では, 節境界単位での漸進的な独話係り受け解析手法を提案した. 本手法は, 文境界が特定されていない独話発話全体に対して, 漸進的に係り受け関係を同定する. 文末は係り先がないとして解析する. 本手法の有効性を評価するために, 独話コーパスを用いて係り受け解析実験を行った. 実験の結果, 本手法が, 文単位を入力とする従来の係り受け解析手法と同程度の解析精度と解析時間を備えつつ, 解析の漸進性を向上させることができることを確認した.

今後は, 1) 節境界をまたぐ係り受け関係を解析可能にする, 2) 不変化閾値を節境界単位の種類や単語によって動的に変化させる, 3) 節境界情報やポーズ情報を考慮した確率モデルを検討する, 4) より精緻な学習モデルを採用する, などにより, 漸進的係り受け解析のさらなる性能向上を図る予定である.



## 第6章 あとがき

### 6.1 本論文のまとめ

本論文では、音声言語処理システムを実現するための話し言葉処理技術を開発するため、話し言葉の構文解析手法を提案した。話し言葉の特徴から生じる3つの問題を解決するため、話し言葉に特有な非文法的な言語現象を含む文に対して頑健に構文解析する手法、極端に長い文に対して解析精度を維持しつつ効率的に構文解析する手法、及び、話し手の話速に追従できる程度に漸進的に構文解析する手法を提案した。

第1章では、従来の書き言葉に対する構文解析手法を話し言葉に適用したときの問題点を論じ、この問題点ごとに話し言葉の構文解析に関する研究動向を概観した。その中で、話し言葉の構文解析に関する研究はそのほとんどが規則主導型の手法であり、統計的なアプローチによる試みはほとんどなく、話し言葉の特徴から生じる問題はまだ十分に解決されていないことを指摘した。このような現状を踏まえ、本研究では、日本語話し言葉の高性能な構文解析手法を開発することを目的とし、話し言葉の特徴から求められる構文解析器の性能及びその性能を備えるためのアプローチについて議論した。

第2章では、話し言葉の構文特徴を明らかにするための分析データ及び統計的構文解析手法の学習データとして利用することを目的に構築した、対話と独話の2つの構文構造付き音声言語コーパスについて述べた。対話の構文構造付き音声言語コーパスはCIAIR車内音声対話コーパスに対して、独話の構文構造付き音声言語コーパスは「あすを読む」書き起こしコーパスに対して、それぞれ構文構造を付与することにより構築した。対話と独話の構文構造付きコーパスの規模は、それぞれ、85,870形態素、192,495形態素である。両コーパスとも話し言葉に特有な言語現象に対しては新たな付与基準を設けた。また、構文構造付き音声独話コーパスには、節境界情報や複数の係り先を付与した。

第3章では、音声言語コーパスを用いた対話文の頑健な係り受け解析手法を提案した。本手法では、構文構造付き音声対話コーパスから各文節間の係り受け確率を

統計的に獲得し、それを用いて係り受け構造の尤度を計算した。また、本手法では、後方修飾性の制約及び係り先の唯一性に関する制約は統計情報を反映させつつ緩和し、非文法的な特徴をもつ発話の解析を可能にした。CIAIR 構文構造付き音声対話コーパスを用いて係り受け解析実験を行い、自然発話文に対しても、書き言葉を対象とした従来の係り受け解析手法 [18, 23, 45, 92] と同等の高い精度で係り受けを抽出できることを示した。また、特に、係り先を持たない文節と倒置、発話単位をまたぐ係り受けの解析に対する本手法の有効性を報告した。

第4章では、節境界に基づく独話文の係り受け解析手法を提案した。本手法では、文の分割単位として節を採用し、節レベルと文レベルの二段階で係り受け解析を実行する。まず、節境界解析により文を節に分割し、各節に対して係り受け解析を行うことにより、節内の係り受け関係を同定する。次に、節境界をまたぐ係り受け関係を定め、文全体の係り受け構造を作り上げる。これにより、文が長いという特徴をもつ独話文に対して高性能な解析を実現した。本手法の有効性を評価するために、「あすを読む」構文構造付き音声独話コーパスを用いて係り受け解析実験を行い、本手法により解析精度を改善しつつ、分割しない場合の約 1/5 に解析時間を短縮できることを示した。

第5章では、節境界単位での漸進的な独話係り受け解析手法を提案した。本手法では、独話音声に対して、節が入力されるたびにその節の内部の係り受け構造を作り上げるとともに、すでに入力されている節の係り先を決定することを試みる。節の係り先となる文節の決定は、後続するいくつかの文節との係り受けの尤度を考慮した動的なタイミングで行う。本手法は、文境界が特定されていない独話発話全体に対して、漸進的に係り受け関係を同定することができる。「あすを読む」構文構造付き音声独話コーパスを用いた解析実験の結果、本手法により、文単位を入力とする従来の係り受け解析手法と同程度の解析精度と解析時間を備えつつ、解析の漸進性の向上が可能となることを示した。

## 6.2 今後の課題と将来への展望

本論文が残した課題と将来への展望を述べる。

本論文で提案した3つの構文解析手法に共通した研究課題として以下のことが挙げられる。

- 音声認識誤りへの対応

本論文では、音声認識の精度が100%であると仮定して構文解析実験を行った

が、実際の音声言語処理システムでは音声認識によって生成されたテキストを処理することになる。1章で論じた通り、そこには、多数の認識誤りが含まれるため、極めて頑健な解析技術が要求される [58]。本手法の実用性を検証するために、音声認識結果を解析する実験を行うとともに、その結果を分析することにより、音声認識誤りへの対応を検討する。

- 学習モデルの精緻化

本論文では、係り受け構造を求める際、統計的な手法を用いたが、学習モデルとしては素朴なモデル [10, 18] を用いている。これまでに、統計的係り受け解析手法の学習モデルとして、決定木 [23] や最大エントロピー法 [92]、SVM [46] を用いた手法などが提案されており、高い解析精度を達成している。これらの学習モデルを本研究の構文解析手法に組み込むことにより、解析精度の更なる改善が期待できる。本論文では取り上げていないが、本論文で提案した効率的な係り受け解析手法に対して、最大エントロピー法に基づく学習モデル [92] を組み込むことに現在取り組んでおり、解析精度の改善を確認している [74]。他の2つの構文解析手法に関しても、最大エントロピー法に基づく学習モデルを組み込む予定である。

- 音響情報の利用

3章で提案した頑健な係り受け解析手法では、一部、ポーズの情報を使っているものの、4章、5章で提案した節境界に基づく係り受け解析では利用していない。ポーズの時間や話速の変化、パワーやピッチの急激な変化を利用することは、言い直しの検出 [5] や係り受け解析の精度改善 [15] に有効であることが示されており、これらを考慮した係り受け確率モデルを確立する。

節境界に基づく効率的な構文解析と漸進的な構文解析に共通した研究課題としては、以下のことが挙げられる。

- 節境界をまたぐ係り受け関係の検出

4.5.4 節で述べているように、節境界をまたぐ係り受け関係の存在は、本手法の精度低下を招くため、節境界をまたぐ係り受け関係を事前に検出する必要がある。節境界をまたぐ係り受け関係の一部については、すでにその検出を試みており [74]、その他のものについても検討したい。

- 複数の係り先を持つことを許した柔軟な評価

4.5.4 節で述べた節境界「テ節」をまたぐ係り受け関係のように、京都テキス

トコーパスの付与基準における正解以外に、意味的に間違いとはいえない係り先が別に存在する場合がある。そこで、このような複数の係り先を持つ文節があることを認めて、柔軟に評価することが考えられる。本研究では、このような柔軟な評価を目的に、2.3節で構築した「あすを読む」構文構造付き音声独話コーパスに複数の係り先を付与している。今後、2.3.2節で述べた意味的係り受け関係を用いた柔軟な評価方法を確立する。

- 節境界単位の最終文節に対する解析精度の改善

4.5.3節で述べているように、節境界単位の最終文節に対する係り受け正解率は、節境界単位内部の文節に対する係り受け正解率と比べて20%以上も低く、解析精度の改善を図る必要がある。具体的には、節の階層構造における「従属節の選好性」を利用して、節末の係り先の解析精度を改善した先行研究[84, 94]を本手法に組み込む。

- 頑健な構文解析手法との統合

3章で提案した音声対話文に対する頑健な係り受け解析手法と4章や5章で提案した効率的、もしくは、漸進的な係り受け解析手法を統合することにより、フィラーや言い淀み、倒置などの非文法的言語現象に対する頑健さと高い解析効率や解析の漸進性を備えた話し言葉の高性能な構文解析が実現できる。ただし、倒置を解析する枠組みを単純に統合すると、解析の効率性が低下することもあるため、導入については慎重に検討する必要がある。

さらに、本論文で提案した日本語話し言葉の構文解析手法を利用した音声言語処理システムを開発することが挙げられる。現在、本論文で開発した話し言葉の構文解析器を用いて、以下のシステムを開発中である。

- 会話マイニングシステム [73]

音声対話データから情報抽出により話者の知識を獲得するシステムの開発に取り組んでいる。このシステムでは、3章で提案した頑健な係り受け解析手法により対話に含まれる発話を解析し、その結果得られる係り受け構造に対して抽出規則を適用することにより、発話に含まれる情報を抽出する。また、各発話から抽出された情報を統合することにより、対話全体からの情報をまとめ上げることを検討している。

- 字幕生成システム [75]

5章で述べた節境界に基づく漸進的係り受け解析を利用した独話音声のリアル

タイム字幕生成システムを開発中である．本システムでは，節が検出されるたびに同定される係り受け構造に基づいて，順次，要約処理を実行する単位を定め，その単位ごとに要約処理を実行する．これにより，入力音声のリアルタイムな字幕化の実現が期待できる．なお，この単位は，同定された節末の係り受け関係によって連結された節列であるため，構文的，意味的にまとまっており，要約処理に適していると考えられる．また，要約では，係り受け構造を考慮して，冗長と思われる文節単位，節境界単位を順に削除することにより，不自然な日本語の生成を防ぐことができる．

このほかにも，1.1 節で挙げた音声対話システムや音声翻訳システムなどへの応用が期待できる．今後，各応用システムへの統合を検討したい．

将来の展望としては，人間と同程度の発話理解が可能な話し言葉処理技術の開発が挙げられる．このような話し言葉処理が実現できれば，映画やアニメに登場するような人型ロボットの音声インタフェースを開発することも夢ではなくなる．そのためには，本論文で述べたような話し言葉の構文解析や，その後続く意味解析などの処理において，文脈や世界知識などを考慮することが不可欠である．さらに，人間による発話理解には，言語情報や音声情報に劣らず，話者の表情やジェスチャーなどの視覚情報が重要な役割を担っていると言われており，視覚情報を融合した発話理解の枠組みも必要になると考えられる．



## 謝辞

本論文をまとめるにあたり，多大な御教示と御尽力をいただきました，名古屋大学教授の坂部俊樹先生に厚く謝意を申し上げます．また，本論文の詳細について，貴重な御示唆と御指導をいただきました，名古屋大学教授の阿草清滋先生に深く感謝いたします．

本研究を進め，まとめるにあたり，日頃より懇切丁寧な御指導と御鞭撻をいただきました，名古屋大学准教授の松原茂樹先生に厚く感謝いたします．また，松原茂樹先生には，公私共に様々な御相談にのっていただきました．心より感謝の気持ちを申し上げます．

本研究の初期の段階より幅広い角度から御指導と御討論をいただきました，名古屋大学名誉教授で，現在，愛知工業大学教授の稲垣康善先生に深く感謝いたします．

本研究を進めるにあたり，また，研究室生活を送るにあたり，御指導と御支援をいただきました，名古屋大学准教授の河口信夫先生，助教の山口由紀子先生，加藤芳秀先生に深く感謝いたします．

独話を対象とした構文解析に関する研究を進めるにあたり，多大な御協力と御助言をいただきました，ATR 音声言語コミュニケーション研究所の柏岡秀紀氏，NHK 放送技術研究所の田中英輝氏，加藤直人氏，国立国語研究所の丸山岳彦氏に厚く御礼申し上げます．

本研究の初期の段階より有意義な御議論をいただきました，名古屋大学准教授の外山勝彦先生，助教の小川泰弘先生，元助手の杉野花津江先生，名古屋産業大学准教授の福田ムフタル先生に感謝いたします．

研究に関する討論をはじめ，研究以外の面でもいろいろお世話になりました，松原研究室，河口研究室，旧稲垣研究室の皆様心から感謝いたします．特に，笠浩一朗氏には，研究室に配属されて以来これまで，幅広く相談にのっていただきました．ここに，心より感謝の意を表します．

研究活動を行うにあたりいろいろとお世話になりました，名古屋大学松原研究室秘書の土井ひとみさん，上松可奈さん，名古屋大学旧稲垣研究室秘書の伊藤千佳子さんに深く感謝いたします．

本研究では、CIAIR 車内音声対話コーパス及び「あすを読む」書き起こしデータを利用しました。これらの作成に携わった方々に感謝いたします。また、本研究では、名古屋大学大学院国際言語文化研究科の大学院生数名に御協力をいただき、CIAIR 構文構造付き音声対話コーパス及び「あすを読む」構文構造付き音声独話コーパスを作成しました。ここに、謝意を表します。

本研究を進めるにあたり、ATR 音声言語コミュニケーション研究所ならびに情報通信機構（NICT）で各1ヶ月間、研究員として研究活動を行いました。この間お世話になりました皆様に感謝いたします。

最後に、改めまして、本論文をまとめるにあたりご支援をいただいたすべての皆様に心より御礼申し上げます。

## 発表文献リスト

種別	論文名	関連する章
国際会議	Tomohiro Ohno, Shigeki Matsubara, Hideki Kashioka, Naoto Kato, and Yasuyoshi Inagaki. A syntactically annotated corpus of Japanese spoken monologue. In <i>Proceedings of the 5th International Conference on Language Resources and Evaluation</i> , pp. 1590–1595, May 2006.	2 章
国際会議	Tomohiro Ohno, Shigeki Matsubara, Nobuo Kawaguchi, and Yasuyoshi Inagaki. Spiral construction of syntactically annotated spoken language corpus. In <i>Proceedings of the 2003 IEEE International Conference on Natural Language Processing and Knowledge Engineering</i> , pp. 477–483, Oct. 2003.	2 章
論文誌	Tomohiro Ohno, Shigeki Matsubara, Nobuo Kawaguchi, and Yasuyoshi Inagaki. Robust dependency parsing of spontaneous Japanese spoken language. <i>IEICE Transactions on Information and Systems (Special Section on Corpus-based Speech Technologies)</i> , Vol. E88-D, No. 3, pp. 545–552, Mar. 2005.	3 章
国際会議	Tomohiro Ohno, Shigeki Matsubara, Nobuo Kawaguchi, and Yasuyoshi Inagaki. Robust dependency parsing of spontaneous Japanese speech and its evaluation. In <i>Proceedings of the 8th International Conference on Spoken Language Processing</i> , pp. 2173–2176, Oct. 2004.	3 章

- 論文誌 Tomohiro Ohno, Shigeki Matsubara, Hideki Kashioka, 4 章  
Takehiko Maruyama, Hideki Tanaka, and Yasuyoshi Inagaki. Dependency parsing of Japanese spoken monologue based on clause boundaries. *Language Resources and Evaluation (Special Double-Issue on Asian Language Technology: Resources and Processing)*, Springer. (to appear)
- 国際会議 Tomohiro Ohno, Shigeki Matsubara, Hideki Kashioka, 4 章  
Takehiko Maruyama, and Yasuyoshi Inagaki. Dependency parsing of Japanese spoken monologue based on clause boundaries. In *Proceedings of the Joint 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pp. 169–176, Jul. 2006.
- 特許 特開 2005-228060, 2005 年 8 月 25 日 「節境界に基づく 4 章  
係り受け解析装置」
- 論文誌 大野 誠寛, 松原 茂樹, 柏岡 秀紀, 加藤 直人, 稲垣 康 5 章  
善. 節境界に基づく独話の漸進的係り受け解析. 電子情報通信学会論文誌, Vol. J90-D, No. 2, pp. 556–566, Feb. 2007.
- 国際会議 Tomohiro Ohno, Shigeki Matsubara, Hideki Kashioka, 5 章  
Naoto Kato, and Yasuyoshi Inagaki. Incremental dependency parsing of Japanese spoken monologue based on clause boundaries. In *Proceedings of the 9th European Conference on Speech Communication and Technology*, pp. 3449–3452, Sep. 2005.
- 特許 特開 2006-209173, 2006 年 8 月 10 日 「漸進的な独話係 5 章  
り受け解析装置」

## 参考文献

- [1] Rajeev Agarwal and Lois Boggess. A simple but useful approach to conjunct identification. In *Proceedings of the 30th Annual Meeting of the Association for Computational Linguistics*, pp. 15–21, 1992.
- [2] 秋葉友良, 田中穂積. 拡張部分木を用いた漸進的構文解析. 情報処理学会第 45 回全国大会, Vol. 3, pp. 175–176, 1992.
- [3] 浅原正幸, 松本裕治. IPADIC ユーザーズマニュアル, version2.5.1. 奈良先端科学技術大学院大学, 2002.
- [4] Masayuki Asahara and Yuji Matsumoto. Filler and disfluency identification based on morphological analysis and chunking. In *Proceedings of the ISCA & IEEE Workshop on Spontaneous Speech Processing and Recognition*, pp. 163–166, 2003.
- [5] John Bear, John Downing, and Elizabeth Shriberg. Integrating multiple knowledge sources for detection and correction of repairs in human-computer dialog. In *Proceedings of the 30th Annual Meeting of the Association for Computational Linguistics*, pp. 56–63, 1992.
- [6] John Bear and Patti Price. Prosody, syntax, and parsing. In *Proceedings of the 28th Annual Meeting of the Association for Computational Linguistics*, pp. 17–22, 1990.
- [7] Sabine Brants, Stefanie Dipper, Silvia Hansen, Wolfgang Lezius, and George Smith. The TIGER Treebank. In *Proceedings of the Workshop on Treebanks and Linguistic Theories*, pp. 24–41, 2002.
- [8] Eugene Charniak. A maximum-entropy-inspired parser. In *Processings of the 1st Conference of the North American Chapter of the Association for Computational Linguistics*, pp. 132–139, 2000.

- [9] Martin Čmejrek, Jan Cuřin, Jiří Havelka, Jan Hajič, and Vladislav Kuboň. Prague Czech-English Dependency Treebank: Syntactically annotated resources for machine translation. In *Proceedings of the 4th International Conference on Language Resources and Evaluation*, pp. 1597–1600, 2004.
- [10] Michael J. Collins. A new statistical parser based on bigram lexical dependencies. In *Proceedings of the 34th Annual Meeting of the Association for Computational Linguistics*, pp. 184–191, 1996.
- [11] Mark G. Core and Lenhart K. Schubert. A syntactic framework for speech repairs and other disruptions. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*, pp. 413–420, 1999.
- [12] Walter Daelemans, Anja Höthker, and Erik Tjong Kim Sang. Automatic sentence simplification for subtitling in Dutch and English. In *Proceedings of the 4th International Conference on Language Resources and Evaluation*, pp. 1045–1048, 2004.
- [13] Rodolfo Delmonte. Parsing spontaneous speech. In *Proceedings of the 8th European Conference on Speech Communication and Technology*, pp. 1999–2004, 2003.
- [14] 伝康晴. 話し言葉解析のためのコーパスに基づく優先度計算法. 自然言語処理, Vol. 4, No. 1, pp. 41–56, 1997.
- [15] 江口徳博, 尾関和彦. 韻律情報を利用した係り受け解析. 日本音響学会誌, Vol. 52, No. 12, pp. 973–978, 1996.
- [16] Sanda Harabagiu, Dan Moldovan, and Joe Picone. Open-domain voice-activated question answering. In *Proceedings of the 19th International Conference on Computational Linguistics*, pp. 502–508, 2002.
- [17] 藤井敦. 音声による言語バリアフリーな多言語情報アクセス. 情報処理学会研究報告, SLP-44, pp. 195–200, 2002.
- [18] 藤尾正和, 松本裕治. 語の共起確率に基づく係り受け解析とその評価. 情報処理学会論文誌, Vol. 40, No. 12, pp. 4201–4211, 1999.

- [19] 船越孝太郎, 徳永健伸, 田中穂積. 音声対話用構文解析器の頑健性の評価. 情報処理学会研究報告, NL-152, pp. 35–41, 2002.
- [20] 船越孝太郎, 徳永健伸, 田中穂積. 音声対話システムにおける日本語自己修復の処理. 自然言語処理, Vol. 10, No. 4, pp. 33–53, 2003.
- [21] 古瀬蔵, 山田節夫, 山本和英. 頑健な多言語音声翻訳のための不適格入力の分割処理. 情報処理学会論文誌, Vol. 42, No. 5, pp. 1223–1231, 2001.
- [22] Nicholas J. Haddock. Incremental interpretation and combinatory categorial grammar. In *Proceedings of the 10th International Joint Conference on Artificial Intelligence*, pp. 661–663, 1989.
- [23] 春野雅彦, 白井諭, 大山芳史. 決定木を用いた日本語係り受け解析. 情報処理学会論文誌, Vol. 39, No. 12, pp. 3177–3186, 1998.
- [24] Donald Hindle. Deterministic parsing of syntactic nonfluencies. In *Proceedings of the 21th Annual Meeting of the Association for Computational Linguistics*, pp. 123–128, 1983.
- [25] 堀智織, 古井貞熙. 係り受け SCFG に基づく音声自動要約法の改善. 情報処理学会研究報告, SLP-34, pp. 245–250, 2000.
- [26] 堀智織, 古井貞熙. 単語抽出による音声要約文生成法とその評価. 電子情報通信学会論文誌, Vol. J85-DII, No. 2, pp. 200–209, 2002.
- [27] Jim Huang and Geoffrey Zweig. Maximum entropy model for punctuation annotation from speech. In *Proceedings of the 7th International Conference on Spoken Language Processing*, pp. 917–920, 2002.
- [28] 市丸夏樹, 飛松宏征. 接続助詞の結合順位に基づく複文の構文解析. 情報処理学会研究報告, NL-158, pp. 81–86, 2003.
- [29] Japan Electronic Dictionary Research Institute, Ltd. *The EDR electronic dictionary technical guide (second edition)*, Technical Report TR-045. Japan Electronic Dictionary Research Institute, 1995.
- [30] Aravind K. Joshi. Tree adjoining grammar: How much context-sensitivity is required to provide reasonable structural descriptions? In D.R. Dowty,

- L. Karttunen, and A. Zwicky, editors, *Natural Language Parsing*, pp. 206–250. Cambridge University Press, 1985.
- [31] Pierre Jourlin, Sue E. Johnson, Karen Spärck Jones, and Philip C. Woodland. Spoken document representations for probabilistic retrieval. *Speech Communication*, Vol. 32, No. 1-2, pp. 21–36, 2000.
- [32] Hideki Kashioka and Takehiko Maruyama. Segmentation of semantic units in Japanese monologues. In *Proceedings of the International Conference on Speech and Language Technology (ICSLT 2004) and the International Committee for the Co-ordination and Standardization of Speech Databases and Assessment Techniques (Oriental-COCOSDA 2004)*, pp. 87–92, 2004.
- [33] 加藤芳秀, 松原茂樹, 稲垣康善. 確率木接合文法に基づく漸進的構文解析. 情報処理学会研究報告, NL-166, pp. 15–22, 2005.
- [34] 加藤芳秀, 松原茂樹, 外山勝彦, 稲垣康善. 確率文脈自由文法に基づく漸進的構文解析. 電気学会論文誌, Vol. 122-C, No. 1, pp. 2109–2119, 2002.
- [35] 加藤芳秀, 松原茂樹, 外山勝彦, 稲垣康善. 主辞情報付き文脈自由文法に基づく漸進的な依存構造解析. 電子情報通信学会論文誌, Vol. J86-DII, No. 1, pp. 86–97, 2003.
- [36] Nobuo Kawaguchi, Shigeki Matsubara, Kazuya Takeda, and Fumitada Itakura. Multimedia data collection of in-car speech communication. In *Proceedings of the 7th European Conference on Speech Communication and Technology*, pp. 2027–2030, 2001.
- [37] Nobuo Kawaguchi, Shigeki Matsubara, Kazuya Takeda, and Fumitada Itakura. Multi-dimensional data acquisition for integrated acoustic information research. In *Proceedings of the 3rd International Conference on Language Resources and Evaluation*, pp. 2043–2046, 2002.
- [38] 河口信夫, 牛窪誠一, 松原茂樹, 梶田将司, 武田一哉, 板倉文忠. 走行車室内音声対話収録システムの開発. 電子情報通信学会論文誌, Vol. J84-DII, No. 6, pp. 1122–1129, 2001.

- [39] Ji-Hwan. Kim and Philip C. Woodland. The use of prosody in a combined system for punctuation generation and speech recognition. In *Proceedings of the 7th International Conference on Spoken Language Processing*, pp. 2757–2760, 2001.
- [40] Mi-Yong Kim and Jong-Hyeok Lee. Syntactic analysis of long sentences based on s-clauses. In *Proceedings of the 1st International Joint Conference on Natural Language Processing*, pp. 518–526, 2004.
- [41] 金淵培, 江原暉将. 日英機械翻訳のための日本語長文自動短文分割と主語の補完. 情報処理学会論文誌, Vol. 35, No. 6, pp. 1018–1028, 1994.
- [42] Itsuki Kishida, Yuki Irie, Yukiko Yamaguchi, Shigeki Matsubara, Nobuo Kawaguchi, , and Yasuyoshi Inagaki. Construction of an advanced in-car spoken dialogue corpus and its characteristic analysis. In *Proceedings of the 8th European Conference on Speech Communication and Technology*, pp. 1581–1584, 2003.
- [43] 駒谷和範, 上野晋一, 河原達也, 奥乃博. 音声対話システムにおける適応的な応答生成を行うためのユーザモデル. 電子情報通信学会論文誌, Vol. J87-DII, No. 10, pp. 1921–1928, 2004.
- [44] 小松昭男, 大平栄二, 市川薫. 韻律情報を利用した構文推定およびワードスUPPORTにおける会話音声理解方式. 電子情報通信学会論文誌, Vol. J71-D, No. 7, pp. 1218–1228, 1988.
- [45] Taku Kudo and Yuji Matsumoto. Japanese dependency analysis based on support vector machines. In *Proceedings of the 2000 Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, pp. 18–25, 2000.
- [46] 工藤拓, 松本裕治. チャンキングの段階適用による係り受け解析. 情報処理学会論文誌, Vol. 43, No. 6, pp. 1834–1842, 2002.
- [47] Sadao Kurohash and Makoto Nagao. KN Parser: Japanese dependency/case structure analyzer. In *Proceedings of the International Workshop on Sharable Natural Language Resources*, pp. 48–55, 1994.

- [48] Sadao Kurohashi and Makoto Nagao. A syntactic analysis method of long Japanese sentences based on the detection of conjunctive structures. *Computational Linguistics*, Vol. 20, No. 4, pp. 507–534, 1994.
- [49] Sadao Kurohashi and Makoto Nagao. Building a Japanese parsed corpus while improving the parsing system. In *Proceedings of the 1st International Conference on Language Resources and Evaluation*, pp. 719–724, 1998.
- [50] Vilson J. Leffa. Clause processing in complex sentences. In *Proceedings of the 1st International Conference on Language Resources and Evaluation*, pp. 937–943, 1998.
- [51] Esther Levin, Shrikanth Narayanan, Roberto Pieraccini, Konstantin Biatov, E. Bocchieri, Giuseppe Di Fabbrizio, Wieland Eckert, S. Lee, A. Pokrovsky, Mazin Rahim, P. Ruscitti, and M. Walker. The AT&T-DARPA communicator mixed-initiative spoken dialog system. In *Proceedings of the 6th International Conference on Spoken Language Processing*, Vol. 2, pp. 122–125, 2003.
- [52] Caroline Lyon and Bob Dickerson. Reducing the complexity of parsing by a method of decomposition. In *Proceedings of the 6th International Workshop on Parsing Technology*, 1997.
- [53] 前川喜久雄, 籠宮隆之, 小磯花絵, 小椋秀樹, 菊池英明. 日本語話し言葉コーパスの設計. 音声研究, Vol. 4, No. 2, pp. 51–61, 2000.
- [54] Mitchell P. Marcus, Beatrice Santorini, and Mary A. Marcinkiewicz. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, Vol. 19, No. 2, pp. 313–330, 1993.
- [55] 丸山岳彦, 柏岡秀紀, 熊野正, 田中英輝. 日本語節境界プログラム CBAP の開発とその評価. 自然言語処理, Vol. 11, No. 3, pp. 517–520, 2004.
- [56] 益岡隆志, 田窪行則. 基礎日本語文法 改訂版 . くろしお出版, 1992.
- [57] Shigeki Matsubara, Satoru Asai, Katsuhiko Toyama, and Yasuyoshi Inagaki. Chart-based parsing and transfer in incremental spoken language translation. In *Proceedings of the 4th Natural Language Processing Pacific Rim Symposium*, pp. 521–524, 1997.

- [58] 松原茂樹, 佐藤利光, 河口信夫, 稲垣康善. 統計データを用いた話し言葉音声の係り受け解析. 情報処理学会研究報告, NL-143, pp. 63–68, 2001.
- [59] 松本裕治, 北内啓, 山下達雄, 平野善隆, 松田寛, 高岡一馬, 浅原正幸. 形態素解析システム『茶筌』 and version2.2.9 and 使用説明書, 2002.
- [60] Igor' A. Mel'čuk. *Dependency Syntax: Theory and Practice*. SUNY Press, 1987.
- [61] Marie Meteer, Ann Taylor, Robert MacIntyre, and Rukmini Iyer. *Dysfluency annotation stylebook for the Switchboard Corpus*. Linguistic Data Consortium (LDC), 1995.
- [62] David Milward. Incremental interpretation of categorial grammar. In *Proceedings of the 7th Conference of European Chapter of the Association for Computational Linguistics*, pp. 119–126, 1995.
- [63] 翠輝久, 駒谷和範, 清田陽司, 河原達也. 音声対話によるソフトウェアサポートタスクのための効率的な確認戦略. 電子情報通信学会論文誌, Vol. J88-DII, No. 3, pp. 499–508, 2005.
- [64] 水谷研治, 小沼知浩, 遠藤充, 南部太郎, 脇田由実. Pda で動作する旅行会話向け音声翻訳システムのインタフェース評価. 情報処理学会研究報告, HI-103, pp. 1–6, 2003.
- [65] Tsuyoshi Morimoto, Noriyoshi Uratani, Toshiyuki Takezawa, Osamu Furuse, Yasuhiro Sobashima, Hitoshi Iida, Atsushi Nakamura, Yoshinori Sagisaka, Norio Higuchi, and Yasuhiro Yamazaki. A speech and language database for speech translation research. In *Proceedings of the 3rd International Conference on Spoken Language Processing*, pp. 1791–1794, 1994.
- [66] 村上仁一, 嵯峨山茂樹. 自由発話音声における音響的な特徴の検討. 電子情報通信学会論文誌, Vol. J78-DII, No. 12, pp. 1741–1749, 1995.
- [67] 長尾真 (編). 岩波講座ソフトウェア科学 15 自然言語処理. 岩波書店, 1996.
- [68] 中川聖一, 小林聡. 自然な音声対話における間投詞・ポーズ・言い直しの出現パターンと音響的性質. 日本音響学会誌, Vol. 51, No. 3, pp. 202–210, 1995.

- [69] 中野幹生, 島津明. 言い直しを含む発話の解析. 情報処理学会論文誌, Vol. 39, No. 6, pp. 1935–1943, 1998.
- [70] 那須川哲哉, 宅間大介, 竹内広宜, 荻野紫穂. コールセンターにおける会話マイニング. 言語処理学会第13回年次大会発表論文集, pp. 590–593, 2007.
- [71] 西崎博光, 中川聖一. 音声キーワードによるニュース音声データベース検索手法. 情報処理学会論文誌, Vol. 42, No. 12, pp. 3173–3184, 2001.
- [72] Joakim Nivre. Incrementality in deterministic dependency parsing. In *Proceedings of the ACL Workshop Incremental Parsing: Bringing Engineering and Cognition Together*, pp. 50–57, 2004.
- [73] 小野貴博, 大野誠寛, 松原茂樹, 山口由紀子, 河口信夫, 吉川正俊. 話し言葉解析に基づく話者知識の自動獲得. 人工知能学会第19回全国大会講演論文集, 2005.
- [74] 大野誠寛, 松原茂樹, 柏岡秀紀, 稲垣康善. ポーズを考慮した文の分割に基づく独話文の係り受け解析. 言語処理学会第13回年次大会発表論文集, pp. 183–186, 2007.
- [75] 大野誠寛, 松原茂樹, 柏岡秀紀, 加藤直人, 稲垣康善. 同時的な独話音声要約に基づくリアルタイム字幕生成. 情報処理学会研究報告, SLP-62, pp. 51–56, 2006.
- [76] Carl Pollard and Ivan A. Sag. *Head-Driven Phrase Structure Grammar*. Chicago Press, 1994.
- [77] Alexandros Potamianos, Egbert Ammicht, and Hong-Kwang J. Kuo. Dialogue management in the bell labs communicator system. In *Proceedings of the 6th International Conference on Spoken Language Processing*, Vol. 2, pp. 603–606, 2003.
- [78] Adwait Ratnaparkhi. A linear observed time statistical parser based on maximum entropy models. In *Proceedings of the Empirical Method for Natural Language Processing*, pp. 1–10, 1997.
- [79] Steve Renals, Dave Abberley, David Kirby, and Tony Robinson. Indexing and retrieval of broadcast news. *Speech Communication*, Vol. 32, No. 1-2, pp. 5–10, 2000.

- [80] Brian Roark. Probabilistic top-down parsing and language modeling. *Computational Linguistics*, Vol. 27, No. 2, pp. 249–276, 2001.
- [81] Shourya Roy and L Venkata Subramaniam. Automatic generation of domain models for call-centers from noisy transcriptions. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th Annual Meeting of the Association for Computational Linguistics*, pp. 737–744, 2006.
- [82] 颯々野学. 日本語係り受け解析の線形時間アルゴリズム. *自然言語処理*, Vol. 14, No. 1, pp. 3–18, 2007.
- [83] Satoshi Sekine. Japanese dependency analysis using a deterministic finite state transducer. In *Proceedings of the 18th International Conference on Computational Linguistics*, pp. 761–767, 2000.
- [84] 白井諭, 池原悟, 横尾昭男, 木村淳子. 階層的認識構造に着目した日本語従属節間の係り受け解析の方法とその精度. *情報処理学会論文誌*, Vol. 36, No. 10, pp. 2353–2361, 1995.
- [85] Kazuya Shitaoka, Kiyotaka Uchimoto, Tatsuya Kawahara, and Hitoshi Isahara. Dependency structure analysis and sentence boundary detection in spontaneous Japanese. In *Proceedings of the 20th International Conference on Computational Linguistics*, pp. 1107–1113, 2004.
- [86] Elizabeth Shriberg, Andres Stolcke, Dilek Hakkani-Tür, and Gökhan Tür. Prosody-based automatic segmentation of speech into sentences and topics. *Speech Communication*, Vol. 32, No. 1-2, pp. 127–154, 2000.
- [87] Wojciech Skut, Thorsten Brants, Brigitte Krenn, and Hans Uszkoreit. A linguistically interpreted corpus of German newspaper text. In *Proceedings of the 1st International Conference on Language Resources and Evaluation*, pp. 705–711, 1998.
- [88] Andreas Stolcke and Elizabeth Shriberg. Statistical language modeling for speech disfluencies. In *Proceedings of the 1996 IEEE International Conference on Acoustics, Speech and Signal Processing*, Vol. 1, pp. 405–408, 1996.

- [89] 田島幸恵, 難波英嗣, 奥村学. 形態素解析器を利用した講演書き起こしの文境界検出について. 2003FIT 情報科学技術フォーラム講演論文集, pp. 155–156, 2003.
- [90] 武石英二, 林良彦. 接続構造解析に基づく日本語複文の分割. 情報処理学会論文誌, Vol. 33, No. 5, pp. 652–663, 1992.
- [91] Toshiyuki Takezawa, Eiichiro Sumita, Fumiaki Sugaya, Hirofumi Yamamoto, and Seiichi Yamamoto. Toward a broad-coverage bilingual corpus for speech translation of travel conversations in the real world. In *Proceedings of the 3rd International Conference on Language Resources and Evaluation*, pp. 147–152, 2002.
- [92] 内元清貴, 関根聡, 井佐原均. 最大エントロピー法に基づくモデルを用いた日本語係り受け解析. 情報処理学会論文誌, Vol. 40, No. 9, pp. 3397–3407, 1999.
- [93] 浦谷則好, 森元逞, 谷戸文廣. 音声翻訳システム ASURA の開発. 情報処理学会研究報告, CH-18, pp. 49–54, 1993.
- [94] 宇津呂武仁, 西岡山滋之, 藤尾正和, 松本裕治. コーパスからの日本語従属節係り受け選好情報の抽出およびその評価. 自然言語処理, Vol. 6, No. 7, pp. 29–60, 1999.
- [95] Ton van der Wouden, Heleen Hoekstra, Michael Moortgat, Bram Renmans, and Ineke Schuurman. Syntactic analysis in The Spoken Dutch Corpus (CGN). In *Proceedings of the 3rd International Conference on Language Resources and Evaluation*, pp. 768–773, 2002.
- [96] Nanette Veilleux and Mari Ostendorf. Probabilistic parse scoring based on prosodic phrasing. In *Proceedings of the workshop on Speech and Natural Language*, pp. 429–434, 1992.
- [97] Wolfgang Wahlster, editor. *Verbmobil: Foundations of Speech-to-Speech Translation*. Springer, 2000.
- [98] Chao Wang and Stephanie Seneff. High-quality speech translation in the flight domain. In *Proceedings of the 9th International Conference on Spoken Language Processing*, pp. 761–764, 2006.

- [99] 山本幹雄, 小林聡, 中川聖一. 音声対話文における助詞落ち・倒置の分析と解析手法. 情報処理学会論文誌, Vol. 33, No. 11, pp. 1322–1330, 1992.
- [100] Xiaodan Zhu and Gerald Penn. Summarization of spontaneous conversations. In *Proceedings of the 9th International Conference on Spoken Language Processing*, pp. 1531–1534, 2006.
- [101] Geoffrey Zweig, Olivier Siohan, George Saon, Bhuvana Ramabhadran, Daniel Povey, Lidia Mangu, and Brian Kingsbury. Automated quality monitoring in the call center with ASR and maximum entropy. In *Proceedings of the 2006 IEEE International Conference on Acoustics, Speech and Signal Processing*, Vol. 1, pp. 589–592, 2006.