

# 談話研究のツールとしての転記エディターと談話データベース

亀山 真一

(ハンブルク大学言語学部)

## 1. はじめに

話し言葉のコーパスにもとづく言語学的分析を試みるとき、収録した音声データやビデオ画像データをどのように文字化し、分析しやすい形態へ移行するかといった問題に直面することになる。そのとき、次のような点について考慮し、データ作成に取り組んでいく必要があるように思う。

- ・ 現象の捉えやすい文字化データ：文字化データは、談話における言語行動や非言語行動を視覚的に捉えやすいかたちで書き起こしたものであることが望ましい。また、文字化データのもととなった画像・音声データに立ち戻り、関心のある現象を直接見直したり聞き直したりすることができることよい。
- ・ 加工処理のしやすい文字化データ：大規模なコーパスの分析に取り組む場合は、文字化データを計算機で処理できるような形態へ移行していく必要がある。談話の文字化を行うときからそのことを念頭において、データが加工しやすいような形式で転記しておくことよい。
- ・ 検索処理のしやすいコーパス・システム：コーパス・システムは、検索・整理などの作業がしやすく、全体像が把握しやすいものであることが望ましい。そのためには、文字化データをデータベース化していくなどの加工処理が必要となる。
- ・ データの関連性の捉えやすいコーパス・システム：談話そのものと談話の展開に影響を及ぼす諸要因との関連が捉えやすいように、コーパス・システムは文字化データと文字化した談話に関する諸情報との関連をモデル化したものであることが望ましい。

本稿では、このような条件を満たすコーパス作成システムの一例として、ハンブルク大学で開発された転記エディターを用いた談話データベースの構築方法について報告する<sup>1</sup>。本稿の構成としては、談話データ文字化のための転記法、転記エディター、転記テキストをもとにした談話データベースの構築方法、談話データベースの構造、談話データベースの活用方法、ハンブルク談話データベース、という順で実例を交えながら紹介する。

## 2. 談話データの文字化

### 2.1. 文字化のための転記法

談話データを文字化する場合、何らかの転記法の規定に従って作業を進めていく必要がある。ここでは数多くある転記法のなかから、Edwards&Lampert(1993)などで紹介され、筆者も現在使用しているHIAT転記法を取り上げる<sup>2</sup>。HIAT転記法は、Ehlich&Rehbeinが1976年に開発、発表したもので、以来ドイツでの談話分析では標準的な転記法とみなされている。

以下、図1にHIAT転記法による転記テキストの例を示す。

40	MUR	/150 /151 /152	[Soo]. ((warai)) Bodikon. (Sore de moo) mā [bodii konshasu] na n da kedo, bodikon
	KRÜ	/153 /154	Naruhodo]. ((----warai----)) (Mā)...
			<i>[warainagara]</i> <i>[body conscious]</i>
41	MUR	/155 /156 /157	[suutsu] toka ne, bodikon to iu to, moo sugoi hora, • karada ni [fitto] shita - ichiji
	KRÜ	Hài Ä:: Hai	
			<i>[suits]</i> <i>[fit]</i>
42	MUR	/158 /159 /160	hayatta deshoo soo iu [suutsu] ga - <u>pita::</u> to shita, soo iu no o bodikon tte iu no.
	KRÜ	Ä:: Hài Ä	
			<i>[suits]</i>
43	MUR	/161/162 /163 /164	Uñ ((1s)) Dakara itsumo soo iu koto bakari ne, ano ki ni shiteru • • hitotachi no
	KRÜ	soo (Taito)...	Hài
44	MUR	/165 /166 /167	koto o • bodikon... <u>Anmari</u> otoko de bodikon no hito inai wa ne.
	KRÜ	Otoko mo onna • mo so iu?	( )

図1 HIATスコア

図1は、HIATによる転記テキスト（スコア形式）の一例として、筆者の作成した日本語談話資料のなかから取り上げたものである。談話タイプは、日本語の個人授業の授業談話で、話者は、MURが教師で母語が日本語、KRÜは生徒で母語はドイツ語、いずれも女性である。転記テキストの構造としては、一つ一つの譜面（40-44）の中の発話（150-167）を平行して読むようになっており、枠外にあるコメントは角括弧内 [ ] の語について言及するものである。その他の特殊記号としては、問投詞の場合などモーラ上のアクセントの上昇、下降、上昇下降、下降上昇、平坦アクセントといったパターンを持つものについてはアクセントの表記をするほか、聞き取り困難な箇所は丸括弧（）、笑いなどの言葉でない音声は二重の丸括弧（（））、下線は強調、音の引き伸ばしはコロンの、発話内の自己修正はスラッシュ/、発話の中断は...、続けて話すことによる発話の結合は、休止（ポーズ）には中丸・一つが途切れる程度、二つが0.5秒以下、三つが0.5秒以上1秒未満、それ以上長いものでは二重の丸括弧（（））に長さを明示したもの、といったものを採用している。

HIAT転記法には、データの捉えやすさを考慮した3つの特徴があげられる。

楽譜に倣ったスコア形式を採用していること。

発話を単位として談話をセグメントに区切っていること<sup>3</sup>。

書き言葉の表記法を話し言葉の用途に適用した表記を用いていること。

スコア形式とは、いわゆるオーケストラの楽譜の場合のように話者一人一人の言語行動を同一の譜面に並記したものであり、談話中の発話の時間的關係が転記テキストにおいて捉えやすくなるという利点を持っている。これを応用すれば、話者の非言語行動を体位、身振り、表情、視線の動向などに細分化した情報や、発話に附随するコメント、行間翻訳、或いは様々な範疇によるタグのような二次情報などを行別に加えていくことが可能となる<sup>4</sup>。基本的な考え方は、談話中にみられる様々な言語行動、非言語行動情報の時間的關係をできるだけ正確に書き表そうということである。発話と発話の時間的關係を相対的に捉えることは、たとえば話者交代やあいづちの研究をするときなどは必要不可欠である。

発話という単位に談話を区切るということは、話者が聴者に理解できるよう自分の言語行動を発話単位で分割しながら遂行していることを記述することである<sup>5</sup>。また発話は、統語構造や韻律が談話においてどのような機能を果たしているか(たとえば発話モード、発話のテーマ・レーマ構造など)ということを見ていく上で、必要な単位である。

書き言葉の表記法を話し言葉に適用した表記とは、転記テキストにおいて、話し言葉の特徴を最大限に活かしながらも、テキストとしての読みやすさを保つため、書き言葉の規定、つまり句読点の用法、言葉の綴りに関するルールなどにできるだけ従い、それをうまく応用しながら、そして必要最小限の特殊記号しか使わず、話し言葉を転記していくということである。話し言葉にしかないような修正やポーズなどについては特殊記号を使ったり、間投詞の場合など必要に応じてアクセントをつけたり、発音や言葉遣いの特殊なものはそれなりに忠実に転記したりする必要もあるが、書き言葉の規定が無理なく適応できる範囲では、視覚的に捉えやすいように、従来の見慣れた表記を使用するというのが基本的な考え方である<sup>6</sup>。

転記テキストの精度を高め、転記の標準化を図るため、Rehbein, Griebhaber, Löning, Hartung & Bührig(1993)は長期にわたる実際の転記経験をもとに転記のためのマニュアルを作成した。ハンブルク大学ではそれを参照しながら文字化をするようにしている。

## 2.2. 文字化のための転記エディター

転記エディターsyncWRITERは1990年に談話研究や手話研究の資料作成ツールとしてハンブルク大学で開発されたものである<sup>7</sup>。ここではsyncWRITERでHIATに準じた転記を行うというふうに両方併せて紹介するが、厳密に言えばsyncWRITERは純粋に転記のためのツール、HIATはそれとは別に転記法というふうに分けて考えなくてはならない。つまりsyncWRITERで転記を行う場合、上で紹介したHIAT転記法に限らず、ほかのどのような転記法で転記を行ってもよいし、HIATで転記を行う場合、勿論別のエディターを使っても構わない、ということである<sup>8</sup>。

従来のテキスト・エディターでは、テキスト入力を上から下へ行ごとに行う。談話を図1のようにスコアとして入力する場合、様々な編集作業が必要となる。そのため、一度入力したものを訂正していくとなると、いろいろと面倒な作業が伴う。談話を転記する目的で開発されたsyncWRITERの場合は、一度設定したトラックに改行なしでどんどん文字入力が行なえるようになっている。また、タブを使ってトラック別になっているテキストを揃えることができる。入力を済

ませてから自動改行を行うため、改行や譜面の枠（図1参照）のことを考える必要がない上、入力と自動改行の工程が幾度でもくり返せるので、テキストの訂正が容易に行なえる。

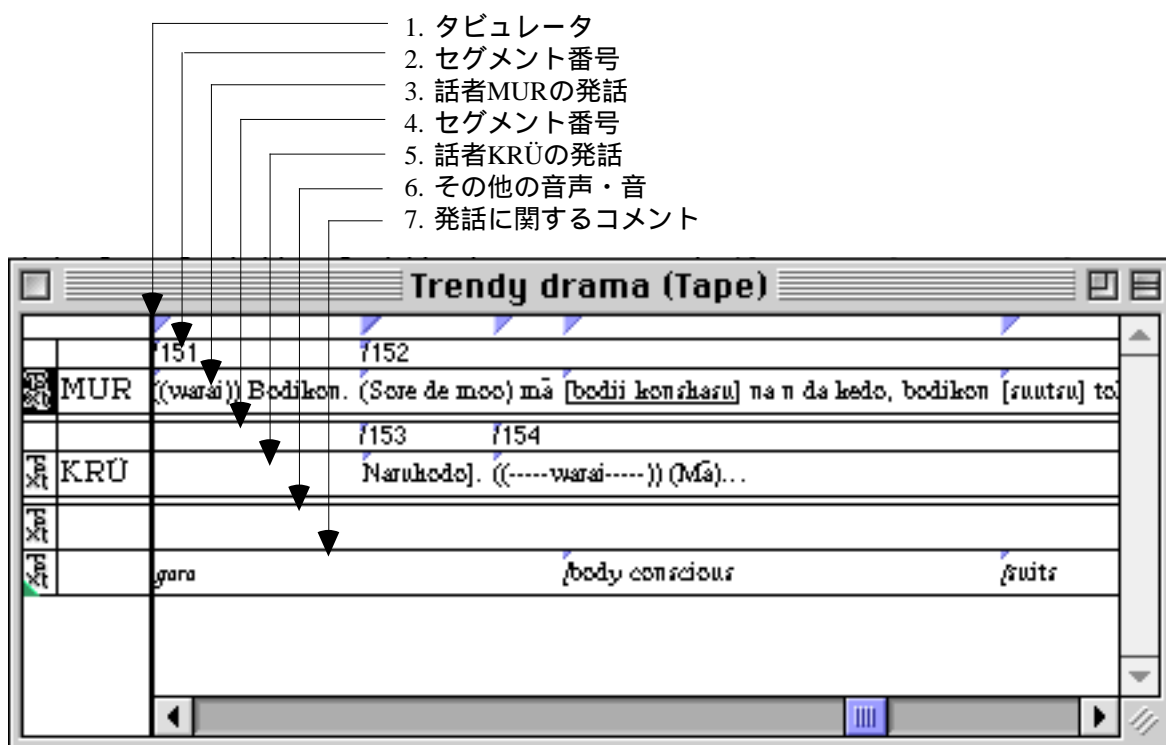


図2 syncWRITERの入力画面

図2はsyncWRITERを使った転記テキスト作成の一例である。入力の後、スコア設定をして自動改行を行うと、図1のようなスコアが出力できる。最上列にある小さな三角形のアイコンが発話を同期化するためのタブで、ここではたとえば発話セグメント152と発話セグメント153を同期化するため、それから発話にセグメント番号やコメントを添付するためなどに使われている。その下の列が第一トラックで、そこから話者MURの発話のセグメント番号、話者MURの発話、話者KRÜの発話のセグメント番号、話者KRÜの発話、話者が判別できない音声その他の音、コメントといった順になっている。その他行間に見える細いトラック（行間翻訳）は現在閉じてあるもので、必要に応じて開くことができる。閉じているトラックは、テープを改行してスコア形式に直す場合とテキストをリスト形式で出力する場合、除外される。

SyncWRITERでは、テキストのもととなった音声や動画像信号を添付し、再生することもできる。また、音声信号の一部をくり返し再生したり、ビデオ静止画像の連続を作成したりすることも可能である。こういったマルチメディア機能は、プレゼンテーションに利用できるだけでなく、分析のとき、もととなった音声をすぐに再生したり、ビデオ静止画像の連続を使って非言語行動を詳しく見たりする際にも便利である。

SyncWRITERで作成された転記テキストは、上でも述べたように、スコア形式のほかに、セグメント・リスト（発話ごとに改行した単純なリスト形式のテキスト）としても出力できる。そのためには、「ピリオド+スペース」、「疑問符+スペース」などの文字列を幾つか設定しておいて、

転記テキストの該当箇所に挿入しておく。つまり、転記テキストの中でそのような文字列があるところは、発話という単位の区切りになっているということである。

発話セグメントを分割する文字列



図3 syncWRITERによるビデオ静止画像のシーケンス

セグメント・リストを転記テキストとして使うためには、更に発話順に番号をふったり、編集し直したりする必要がある。スクリプトによる自動変換処理を施すと、図2の入力テープから以下のような発話リストができる。

- (152) MUR (Sore de moo) mā [bodii konshasu] na n da kedo, bodikon [suutsu] toka ne, bodikon to iu to, moo sugoi hora, • karada ni [fitto] shita - ichiji hayatta deshoo soo iu suutsu ga - pita:: to shita, soo iu no o bodikon tte iu no.  
 /body conscious /suits fit
- (153) KRÜ Naruhodo]
- (154) KRÜ ((warai)) (Mā)...
- (155) KRÜ Hài
- (156) KRÜ A::
- (157) KRÜ Hai
- (158) KRÜ A::
- (159) KRÜ Hài
- (160) KRÜ A soo
- (161) MUR Uñ
- (162) MUR ((1s)) Dakara itsumo soo iu koto bakari ne, ano ki ni shiteru •• hitotachi no koto o • bodikon...

図4 HIAT発話リスト

発話リストは、その簡潔さゆえに、個々の発話の構造の特色や言い換えなどの場合のように関連する発話ペアの構造の違いをつかむのには適しているが、談話中の発話の重なりについては時間的に正確な記述をしていないので、ある発話の連なりを見ていく場合、注意する必要がある。従って、分析を行うときは、発話リスト(図4)とスコア形式の転記テキスト(図1)とを併せて見ていくようにするとよい。

### 3. 談話データベース

#### 3.1. 談話データベースの構築方法

本節では、談話コーパスを作成する方法の一つとして、既存のデータベース・プログラムを利用した談話データベースの構築方法を紹介する<sup>9</sup>。この談話データベースは、概念や基本構造、そして構築方法などはEissenhauer&Benecke(1995)が構築したハンブルク談話データベースに倣い、筆者が研究目的で構築したものである。規模やシステムなどに多少の違いがあるものの、どちらの談話データベースも、文字化したデータの検索やタグ付け、データの関連付けの多様性を追求するところから出発しており、談話の転記テキストをもとにしたいわゆるテキスト・データベースであるという点では変わらない。

コーパスを構築していくとき注意しなければならないのは、システム自体がどれだけ汎用性を持つ便利なものでも、データの加工処理や入力処理に手間がかかりすぎれば、データの追加、メンテナンスなどの作業が困難になってしまう点である。従って、文字化データを使った談話データベースを構築するとき、もととなる転記テキストからコーパスに至る処理をできるだけ簡単にすることが必要である。ここで紹介する談話データベースの場合、データの加工・入力処理を図化すると、図5に示したようなものになる。

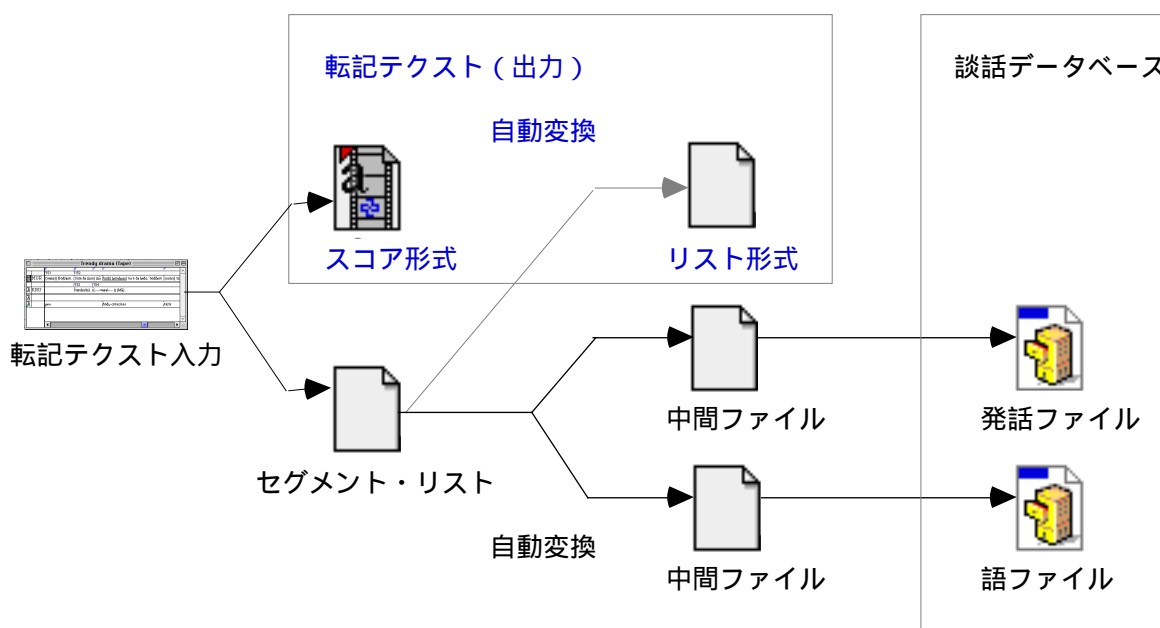


図5 転記テキストから談話データベースまで

上で紹介した転記テキストは、発話や語の単位に区切られているため、少し手を加えただけでデータベースに取り込めるような形態になる。SyncWRITERで出力したセグメント・リストを一旦中間ファイル(データセット・リスト)に自動変換し、文字化データを発話単位、語単位で取り込めば、談話データベースの中核となる発話ファイルと語ファイルがほぼできあがる。文字化された転記テキストそのものはそれぞれ発話単位・語単位といった形態で発話ファイル・語ファイ

ルの発話・語のフィールドに取り込まれる（図6右参照）。その時、発話や語以外に、発話の含まれている談話のID、転記テキストのID、転記テキスト中の発話のID番号、発話の中の語のID番号といった発話や語の出所の情報とともに、その発話や語を発した話者のID（略称）などが一緒に発話ファイル・語ファイルに取り込まれる。

談話データベースを構築するには、そのほかに談話に関するデータ（談話のタイプ、使用言語、発話状況など）、話者に関するデータ（母語、使える外国語など）や収録に関するデータ（録音機器やマイクの種類など）といった二次データが必要となる。また、談話の発話データを文字化するときには、転記テキストのほかに転記に関するデータ（使用した転記法、転記した箇所など）の記録をとっておく必要がある。そのような二次データについては、談話を収録および転記した後に、直接データベースへ入力するか、または事前に必要項目の用紙を作っておいて、その場でデータを記録しておくようにすれば、それほどの手間はかからない。

### 3.2. 談話データベースの構造

談話コーパスを作成する場合、上述のように談話、話者や収録、転記に関するデータなどの二次データが文字化データそのものとともに保存、管理できるシステムを作っていくことが望ましい。それは、収集した談話の発話がどのような状況で発せられたものなのか、手元にある文字化データはどのようなプロセスを経て作成されたものなのか、といった情報を正確に記録しておかないと、分析の段階で発話データの位置付けができなくなるからである。談話データベースというのは、そういった諸々のデータをただ保存、管理できればよいというだけではない。当然、分析の段階でこういった多様な情報が十分に活用できるかということが問題になってくる。どれだけ優れた談話データを集めたコーパスでも、データをうまく検索・整理できるシステムが伴わないかぎり、研究に役立てることは困難である。そのような構造の談話データベースを作っていくためには、全体としてどのようなデータがあるか、そしてその様々な項目のデータが互いにどのような関連を持っているかということを考え、整理していかなければならない。

筆者の構築した談話データベースの場合、文字化データとして談話を発話・語単位で含んだ発話ファイルと語ファイルのほかに、談話ファイル（談話に関する情報を含む）、話者ファイル、収録ファイル、転記ファイル（転記に関する情報を含む）がある。ファイルの項目が互いにどのような関連をもっているかということを考えて関連が組んであるので、関連付けされたファイルの項目による複合条件検索が可能である。談話データベースの構造を各ファイルの項目、ファイル間の関連について図化すると、図6のようなものになる<sup>10</sup>。



図6 談話データベースの構造

たとえば文字化データが含まれている発話ファイル・語ファイルは話者フィールドをキーにして話者ファイルと、また談話フィールドをキーにして談話ファイルと関連付けてある。そうすることによって、発話ファイル・語ファイルで発話や語を検索する際に、話者ファイルや談話ファイルの項目を使った複合条件検索ができるようになる。

### 3.3. 談話データベースの活用方法

上の図6のような関連付けのなされた構造の談話データベースの活用方法例として、ここでは二次情報を使った複合条件検索を二つ紹介する。

最初の例(図7)は、日本語の談話(談話データの「言語」の項目が「日本語」)で、ドイツ人日本語学習者(話者データの「母語」の項目が「ドイツ語」、「外国語」の項目が「日本語」)の発話で、助詞の「は」(発話フィールドが"wa")が使われているもの、といった複数の条件で談話コーパスにある発話を検索したものである。





図7 発話ファイルの検索例（ドイツ語母語話者の発話で「は」を含むもの）

このように、ある特定の現象（たとえば、ドイツ語母語話者の「は」の使用例）を調べたいとき、コーパスにはそれについてどのような用例があるかが瞬時に把握できる。

また、この発話ファイルには、タグとして言語行動が入力できるようになっているが、言語行動によるタグ付けをしていくことによって、ある特定のタイプの言語行動の発話だけを引き出して分析の対象としていくことが可能である。Kameyama(1997)では実際にそのような方法を取り、《時間の打ち合わせ》の言語行動を分析した。まず言語行動のタグ付けを行い、それにもとづいて《提案要求》や《提案》などの言語行動の命題内容を実現する言語要素の集合を求めた<sup>11</sup>。

図8は、ドイツ語の"gut"（「（都合が）良い」）という《評価の表現》、"können"（「できる」）という《可能性の表現》を使った言語行動《提案》の発話をリスト化したものだが、このように様々な角度からある言語行動を実現した発話で使われている表現を整理・分類することによって、規則性やバリエーションを見い出すことが容易となる<sup>12</sup>。

転記テキスト	ターン 番号	ターン発話 番号	テキスト発話 番号	発話
g072a.trl	1	6	6	und zwar würde mir sehr <i>gut</i> passen Anfang Mai gleich <P> am vierten , das ist ein Mittwoch .
g072a.trl	15	4	59	mir würd's sehr <i>gut</i> passen<Z> am Dienstag zum Beispiel , dem vierzehnten Juni .
G142A.TRL	7	1	13	<äh> der vierte Mai w<Z>ürde mir<Z> <P> <i>gut</i> passen .
G145A.TRL	9	2	14	und zwar <A> würd' es mir sehr <i>gut</i> an einem Wochenende passen ,
...	...	...	...	...
M105D.TRL	19	2	99	<ähm> <P> oder +/m=/+ meinetwegen <i>können</i> wir <!könn ma> auch Donnerstag vormittag machen .
M093D.TRL	21	1	53	dann <i>können</i> wir uns gleich morgen treffen , wenn Ihnen das recht ist ?
N002K.TRL	38	1	70	wenn Ihnen <!Ihnn> der Nachmittag etwas ungeschickt ist dann <i>können</i> wir 's ruhig auch früher machen .
N017K.TRL	7	3	15	ich würde Ihnen den fünfzehnten dann vielleicht anbieten <i>können</i> <A> .
...	...	...	...	...

図 8 活用用例 1 "gut", "können" という表現を含む言語行動《提案》の発話

談話データベースを分析に活用する場合に注意しなければならないのは、ある現象に的を絞って追及しているうちに、検索した発話や語が文脈からみて談話中のどのようなところに位置しているかを見失いがちになる点である。たとえば、上の提案の発話を見ていく場合、《提案》が常に一つの発話によって実現されているものではなく、複数の発話の連なりによって実現されている可能性もあるということを忘れてはならない。従って、一度談話データベースで検索した現象でも、再度転記テキストへ立ち戻って見直すことが必要である。

### 3.4. ハンブルク談話データベース

談話データベースのもう一つの魅力的なところは、インターネット技術などを使って、多数の研究者がアクセスできるようなサーバ・クライアント・システムが比較的容易に構築できるところである。ハンブルク大学ではベルブモビルで収集したデータをもとに<sup>13</sup>、標準のウェブ・ブラウザを使ってインターネット・アクセスすることができる談話データベースが実現された。ハンブルク談話データベース略称 HaDageS というもので、現在のところドイツ語や英語の「時間の打ち合わせ」の談話が自然談話 50、ロールプレイ談話 650、計 700 の談話（語数にして約 35 万語）が含まれている。実験的に音声データの一部も圧縮された形でアクセスできるようになっており、音声ファイルを聞くことが可能である。今後は、スコア形式の転記テキストも直接アクセスできるようにしたり、日本語、ポルトガル語、トルコ語などのデータも処理できるようにする。外部からのデータのアクセスに関しては、一般公開可能なデータをオープン・アクセスにするほか、データ交換の条件で一部の談話データについてはパスワード・アクセスによる使用を可能にする予定である。

#### 4. おわりに

本稿では、コーパス作成システムの一例として、ハンブルク大学で開発された転記エディターを用いた談話データベースの構築方法について報告し、活用の可能性について紹介した。紹介したコーパス作成システムは、未だ完成されたものではない。これから解決していかなければならない問題が幾つか残されている。なかでも特に、今後、多言語の談話コーパスを作っていく上で、一つ一つの言語の特質を考慮した処置が必要となるものと思われる。

談話データは従来、研究者が自分の研究目的に応じて独自の転記システムで転記テキストを作成し、分析の終了後は、とりあえずそれを何らかのかたちでただ保存しておく、ということが多かった。しかし、同じデータをほかの目的の研究に利用したり、多数の研究者でデータを共有したり、または大量データによる分析を行ったりする場合、後から対処の方法を考えるのでは結構面倒なことになる。どうしたら多数の研究に応用できるようなコーパスが実現できるか、データをどうやって計算機処理できるような形態にしていってよいか、どのような形で多数の研究者が同時にアクセスできるようなコーパスを作ったら良いのか、どのようにしたら大量の談話データを用いた精密な分析ができるのか、といった問題について事前に考え、一貫性のあるコーパスやそれを構築する簡潔なルーチンを作っていくようにすれば、そういった諸問題によりよく対処していくことができるのではないだろうか。ここで紹介した談話研究のツールがそのような対策を講じる上で何かの役に立つことがあれば幸いである。

(本稿を執筆するにあたって、国立国語研究所の佐々木倫子氏、小磯花絵氏、熊谷智子氏をはじめ、多くのの方々にお世話になった。)

#### 注

- 1 ここでの「データベース」とは、上村(1997)の定義による狭義のデータベースのことであって、「大容量の記憶装置に電子化され貯蓄されコンピュータで検索可能な形に整理・統合されたデータの集合」のことを指す。それに対して「コーパス」という場合は、松本・小磯(1996)がいうように一般的に「何らかの視点や対象を明確にして集められた言語データのこと」を指す。つまり「言語データベース」は加工度の高い「コーパス」であるともいえる。
- 2 HIAT転記法についてはEhlich&Rehbein(1976)やEhlich(1993)を、HIATにおけるイントネーションの転記についてはEhlich&Rehbein(1979)、非言語行動の転記についてはEhlich&Rehbein(1981)を参照。
- 3 セグメントとは、談話中のターンより小さい音節、形態素、語、句、発話など、分析の結果得られることのできる言語単位のことをいうが、ここでは発話単位という意味のみで使っている。
- 4 タグ付けなどによる簡易分析の問題点についてはRehbein&Mazeland(1981)を参照。タグセットなどコーパスの標準化動向については竹沢&末松(1995)を参照。
- 5 発話は、単純にイントネーション・ユニットや文のような統語単位と一致するものではないため、発話の区切りを判断するのは容易ではない場合もある。発話には、間投詞によるあいづちなど、完成された発話内行為 (illocutionary act) も命題内容 (propositional content) も含まないものもあり、その形態は多様である。また、発話と言語行動の関係も従来言語行為論で考えられてきたように、一つの発話によって一つの言語行動が具現化されるという単純な関係ではなくて、一つの発話によって幾つもの言語行動が具現化されたり、一つの言語行動が一連の発話によって具現化されるといったこともありうる。発話については、Rehbein(1994)を参照。

- 6 従って、筆者としてはできれば日本語の談話の場合、ローマ字だけではなくて、ほかにカナ漢字まじりの表記も採用したいところなのだが、今のところ、転記エディターが2バイトコード対応になっていないので、残念ながらまだそうすることができない。
- 7 SyncWRITERには、現在アップル・マッキントッシュ対応の英語版とドイツ語版のものがある。SyncWRITER開発の背景についてはGrießhaber(1992)を参照。
- 8 HIATによる転記が行えるエディターとしてほかにミュンヘン大学で開発されたHIATDOSがある。
- 9 計算機処理可能な談話資料を作成するとき、転記テキストにそのまま、SGMLやXMLなどの文章の理論構造を記述する言語でタグ付けをしていく方法もある。このような記述言語によるタグ付けは単位を自由に設定できるという利点がある。
- 10 一つのファイルですべてのデータを管理するより、幾つかのファイルに分けて関連付けしたほうがデータの重複部分が少ないため、効率良くデータ入力ができ、データ容量が少なくて済む。また、談話と話者の関係など一対一、一対多の関係でないものを無理なくモデル化できる。
- 11 言語行動パタン《時間の打ち合わせ》については、Rehbein,Kameyama&Maleck(1994/1998)を参照。
- 12 図8の転記テキストはベルブモビルの転記法に準じて文字化されたデータであるため、HIAT転記法のものとは違った特殊記号が使われている(<A>は息継ぎ、<P>はポーズ、<Z>は引き伸ばし、+/m=/+ は修正、<!hnn>は発音が標準語でないものなど)。ベルブモビルの転記テキストは、もともとターン単位での区切りしかしておらず、発話ごとのタグ付けができるように発話単位で切り直す必要があった。
- 13 ベルブモビルはドイツ科学技術省の総合プロジェクトで、日独語の自動通訳機の開発を目的としたもの。