

国立国語研究所学術情報リポジトリ

Overview of the Monitor Version of the Corpus of Everyday Japanese Conversation

メタデータ	言語: jpn 出版者: 公開日: 2019-02-14 キーワード (Ja): キーワード (En): 作成者: 小磯, 花絵, 天谷, 晴香, 居關, 友里子, 臼田, 泰如, 柏野, 和佳子, 川端, 良子, 田中, 弥生, 西川, 賢哉, 伝, 康晴, AMATANI, Haruka メールアドレス: 所属:
URL	https://doi.org/10.15084/00001682

『日本語日常会話コーパス』モニター公開版の概要

小磯花絵*・天谷晴香・居關友里子・白田泰如・柏野和佳子・川端良子・田中弥生

(国立国語研究所音声言語研究領域)

西川賢哉(国立国語研究所コーパス開発センター)

伝康晴(千葉大学人文科学研究院/国立国語研究所音声言語研究領域)

Overview of the monitor version of the *Corpus of Everyday Japanese Conversation*

Hanae Koiso, Haruka Amatani, Yuriko Iseki, Yasuyuki Usuda, Wakako Kashino,

Yoshiko Kawabata, Yayoi Tanaka, Ken'ya Nishikawa

(National Institute for Japanese Language and Linguistics)

Yasuharu Den (Graduate School of Humanities, Chiba University/
National Institute for Japanese Language and Linguistics)

National Institute for Japanese Language and Linguistics)

要旨

国立国語研究所共同研究プロジェクト「大規模日常会話コーパスに基づく話し言葉の多角的な研究」では、200時間規模の日常会話を収めた『日本語日常会話コーパス』の構築を進めている。このコーパスは、多様な日常場面の会話を、映像まで含めて収録・公開するものであり、世界的に見ても新しい試みである。『日本語日常会話コーパス』の本公開は、プロジェクトの最終年度にあたる2021年度を予定しているが、コーパスの利用可能性や問題などを把握し今後の構築に活かすために、50時間のデータについて2018年12月にモニター公開することを予定している。そこで本稿では、モニター公開データの仕様や特徴について報告する。

1. はじめに

国立国語研究所では、2016年度から開始した共同研究プロジェクト「大規模日常会話コーパスに基づく話し言葉の多角的な研究」において、さまざまなタイプの日常会話200時間をバランス良く収めた大規模なコーパス『日本語日常会話コーパス』(*Corpus of Everyday Japanese Conversation*, 以下CEJC)の構築を進めている(小磯2017)。CEJCの特徴は、(1)収録のために集められた状況での会話ではなく、日常場面の中で当事者たち自身の動機や目的によって自然に生じる会話を対象とすること、(2)多様な場面の会話をバランスよく集めること、(3)音声だけでなく映像まで含めて収録・公開することである(小磯ほか2017)。特に、日常生活の中で生じる会話を200時間の規模で映像まで含めて公開するというのは、世界的に見ても新しい取り組みである。CEJCの本公開は、プロジェクトの最終年度にあたる2021年度に行う予定であるが、コーパスの利用可能性や問題などを把握し今後の構築に活かすために、50時間分の会話を2018年12月にモニター公開することを予定しており、現在、その準備を進めている。モ

* koiso@ninjal.ac.jp

モニター公開では、(1) 50 時間分の会話の映像・音声データ、転記テキスト、形態論情報（短単位情報）、ツールなどを収めたハードディスクでの公開と、(2) 形態論情報（短単位情報）をオンラインで検索できる「中納言」⁽¹⁾での公開を予定している。本稿では、両者に共通するものとして、会話の収録法（2 節）、コーパス格納データの選定方針（3 節）、及びモニター公開対象データの特徴として調査協力者や会話参加者の属性、会話の種類の内訳（4 節）について報告する。その上で、ハードディスクに同梱するデータとして、映像・音声データ、転記テキスト、短単位情報の仕様についてまとめる（5 節）。

2. 会話の収録法

日常場面の中で当事者たち自身の動機や目的によって自然に生じる会話をバランスよく収録するために、主として個人密着法と呼ぶ収録法で会話を集めている。個人密着法は、日常生活の中で生じる会話を、一般の調査協力者（以下、協力者）自身に収録してもらうという方法である。性別・年代の点から均衡性を考慮して選別された 40 名の協力者（男女×20 代・30 代・40 代・50 代・60 代以上×各 4 名）に収録機材を 3 ヶ月ほど貸し出し、日常生活における多様な場面の会話を 15 時間程度収録してもらう。この中から、データの質や倫理的・法的な問題、バランス、会話参加者（以下、会話者）の希望などを考慮し、コーパスに格納するデータとして 4~5 時間の会話を選別する⁽²⁾。モニター公開で対象とするのは、このうち 20 名の協力者が収録したデータである。協力者の内訳については 4.1 節で述べる。

個人密着法では、調査者は収録に介入しない。そのため、協力者自身に、会話の映像・音声の収録、会話者への調査内容及びデータ公開方法の説明、同意書への署名の依頼、フェイスシート（性別、出身地などの会話者の属性）記入の依頼、会話の収録状況等の記録など、実に多くのことを担当してもらう必要がある。このように収録調査には各種個人情報と扱うなど重い責任が生じることから、協力者は 20 歳以上の成人に限定している。収録法の詳細については田中ほか (2018) を参照されたい。

3. コーパス格納データの選定方針

本節では、コーパスに格納するデータをどのように選定しているかについて述べる。CEJC は、多様な会話をバランスよく集めることを目標に掲げている。そこで、普段われわれがどのような種類の会話をどの程度行っているかの指標を得るために、約 250 人の成人を対象に、起床から就寝までの間に行った全ての会話について、いつ、どこで、誰と、何をしながら、どのような種類の会話を行ったか、などを問う会話行動調査を実施した (小磯ほか 2016)。この調査結果を一つの目安として格納データの選定を進めている (小磯ほか 2017, Koiso et al. 2018)。モニター公開対象についても、構築状況を見ながらできるだけ多様な会話が含まれるように選定した (4.3 節参照)。

個人密着法では、収録を始める前に、機材の設定や会話者への説明、書類の記入などが必要

⁽¹⁾ <https://chunagon.ninjal.ac.jp/useraccount/register>

⁽²⁾ 個人密着法では収録が難しいと思われる場面を調査者が主体となり収録する方法として、特定場面法を採用する。約 20 時間をこれに当てる予定である。

表1 モニター公開対象データの調査協力者の属性、対象とする収録数・会話数、会話時間

年代	男性				女性			
	職業・職種等	収録数	会話数	時間	職業・職種等	収録数	会話数	時間
20代	大学生	5	5	2.2h	大学生	7	7	2.6h
	大学院生	5	5	2.5h	大学生	5	10	2.6h
30代	自営業・自由業	4	4	2.8h	会社員・公務員等	5	6	2.7h
	会社員・公務員等	6	6	2.2h	専業主婦	7	7	2.8h
40代	会社員・公務員等	4	5	2.1h	会社員・公務員等	5	5	2.6h
	自営業・自由業	6	6	2.4h	パート・アルバイト	6	6	2.6h
50代	会社員・公務員等	7	7	2.4h	パート・アルバイト	6	6	2.6h
	会社員・公務員等	4	4	2.6h	会社員・公務員等	7	7	2.2h
60代以上	その他（非常勤講師）	9	9	2.1h	自営業・自由業	6	6	2.7h
	定年退職	6	8	3.0h	専業主婦	6	7	2.7h

となるため、話が少し進んだところから収録が開始されることもある。また1回の収録は最大でも1時間程度としており⁽³⁾、会話の途中で収録が切れることもある。そのため、協力者が収録したものから、ある程度のまとまりをもった範囲を「会話」として切り出し、コーパスに格納するデータを決めている。倫理的・法的な問題や会話者の希望などを考慮し、問題のある部分をカットした結果、一つの収録データが複数の会話に分かれることもある。

4. モニター公開対象データの特徴

4.1 調査協力者の属性

2018年3月末の時点で、収録調査、コーパス格納データ選定、転記1次作業、フォローアップインタビューを全て終了した協力者の中から、バランスを考慮して、モニター公開対象とする協力者を20名選んだ。協力者の属性、対象とする収録・会話の数、会話時間の合計を表1に示す。収録スケジュールの都合で40代の女性が3名、60代以上の女性が1名となっているが、それ以外は性別・年代をバランスさせて各層2名ずつとした。収録数は全体で116回、126会話、約50時間（平均2.5時間/1人）である。

4.2 会話者の属性

モニター公開対象となる116の収録に含まれる会話者は、延べ392名、異なり237名である⁽⁴⁾。性別・年代ごとの数を図1に示す。

男性・女性ともに未成年者の数が少ないが、2節に記したように、個人密着法に基づく収録調査では重い負担を伴うことから、協力者は成人に限定している。そのため、未成年者の数は他と比べ必然的に少なくなる。また、延べ、異なりともに、女性に関しては40代・50代が多く60代・70代が少ない傾向が見られる。40代女性の協力者が3名と多かったために、同性・

⁽³⁾ コーパスに格納するのは1協力者あたり4～5時間と限られており、バリエーションを確保するために、このような上限を設けている。

⁽⁴⁾ データには店員との注文等のやりとりなども含まれるが、多くの場合、店員はメインの会話者ではないため、数には含めていない。店員であっても、長く会話を続ける場合で、収録・公開の同意を得たものについては、その限りでない。そのほか、配偶者との会話の途中で妹と電話で短い会話をしているものがあるが、この場合の妹も数に含めていない。

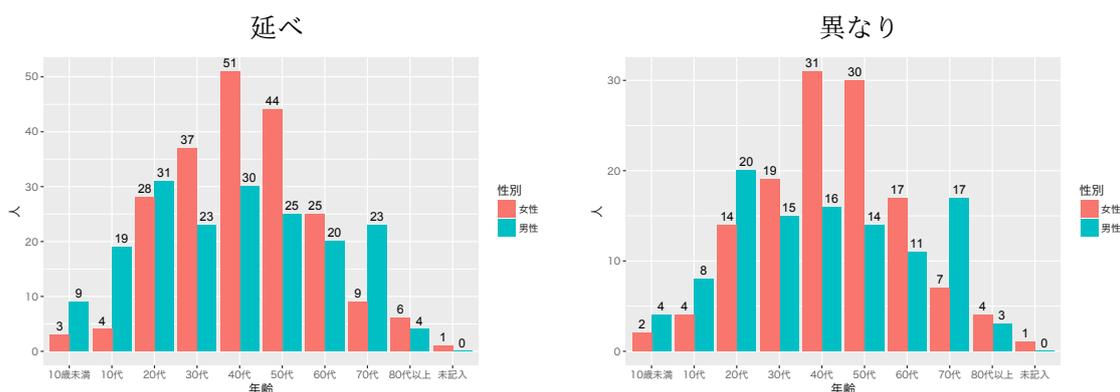


図1 会話者の性別・年代の内訳（人）

表2 会話者の職業の内訳（人）

職業	延べ	異なり	職業	延べ	異なり
会社員・公務員等	127	80	高校生	1	1
自営業・自由業	42	28	中学生	9	4
パート・アルバイト	34	19	小学生	15	6
専業主婦	61	42	就学前	6	4
無職・定年退職	27	19	その他	14	4
大学生・大学院生	54	28	未記入	2	2

同世代の人との会話が他と比べて多くなったためと考えられる。CEJC 全体では、協力者の年齢・性別のバランスをとるようにしているが、こうした対応がコーパス全体の質を保証する上で重要と言える。

会話者の職業の内訳を表2に示す。未成年者が少ないことと関係するが、高校生・中学生・小学生・就学前の人数が少ない。特に高校生についてはモニター公開対象データでは1名のみである。成人については、会社員・公務員等が一番多いものの、それ以外の職業も含めて多様な職業の会話者が含まれている。

協力者から見た相手の会話者との関係を表3に示す。家族や友人知人との会話が多く、仕事関係者や学校等の関係者は少ない。公共商業サービス関係の会話者も少ないが、先に注記したように、注文等で会話した店員などは会話者数には含めていない。そのような店員などを含めると、「サービスを提供する人」は34名となる。

4.3 会話の種類

3節で述べたように、CEJCでは多様な会話をバランスよく集めるために、会話行動調査を実施した。そこで、モニター公開データを対象に「形式」「会話者数」「活動」「場所」の内訳を求め、行動調査の結果と比較する。両者を合わせて図2に示す。図の上段は会話の件数で見た場合の、下段は時間で見た場合の割合の比較である。

「形式」については雑談が約7割を占めており、行動調査より若干多いものの、概ねバランスよくデータが選定できていることが分かる。会議・会合は件数で見ると行動調査より多いが、時間で見ると少ない傾向が見られる。CEJCでは収録の上限を1時間に設定しているのに

表3 調査協力者から見た会話者との関係の内訳（人）

関係性 1	関係性 2	延べ	異なり	関係性 1	関係性 2	延べ	異なり	
家族親戚	配偶者	30	13	仕事	職場の上司	1	1	
	子供	36	17		同僚	7	7	
	父母	17	12		部下	3	3	
	義父母	9	7		取引先など他社の人	6	6	
	自分の兄弟姉妹	7	7		その他	9	9	
	配偶者の兄弟姉妹	3	3		計	26	26	
	自分の祖父母	4	3		先生生徒	学校の先生	4	4
	おい・めい	3	3			学校の生徒・学生	6	3
	その他	5	4			習い事などの生徒	1	1
計		114	69	計	11	8		
友人知人		118	108	公共商業サービス	サービスを提供する人	4	4	
	計	118	108		サービスを受ける人	2	2	
				計	6	6		

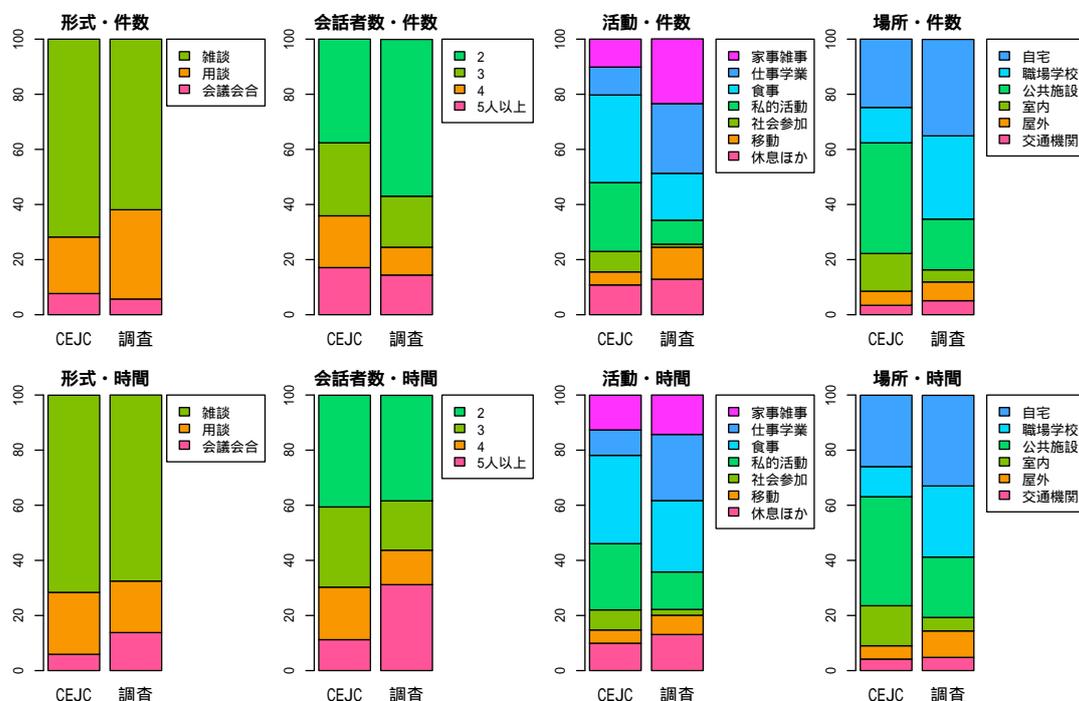


図2 モニター公開対象データと会話行動調査における「形式」「会話者数」「活動」「場所」の内訳

対し、実際の会議・会合は1時間を越えるものが少なからずあることが影響していると考えられる。「会話者数」において、5人以上の場合に、件数では行動調査と同程度だが時間でみると少ない傾向が見られるのも、収録の上限を定めた影響と考えられる。

「活動」と「場所」については、料理や棚の組み立てなどの家事雑事、ボランティアなどの社会参加、屋外・交通機関での移動など、多様な場面の会話が収録できているが、行動調査と比べると、職場や学校における仕事・学業中の会話が少ない。個人密着法ではこの種の会話の収録が難しいためであり、今後、特定場面法で補強する予定である。



図3 基本収録の機材セットで記録した映像の例。左の映像は Kodak PIXPRO SP360 で、右の上下二つの映像は GoPro Hero3+ で録画したもの。論文掲載用に会話者の顔にボカシの処理を加えている。

5. 同梱データの仕様

5.1 映像・音声データ

個人密着法では以下に示す機材を用いて会話を収録している。詳細については田中ほか(2018)を参照されたい。

【基本収録】屋内等での対面会話の収録で主として使用

映像 以下のカメラ2種、計最大3台(最低1台⁽⁵⁾)を使用して会話を録画(図3参照)

1. Kodak PIXPRO SP360 4k (1440×1440, 59.94fps), 最大1台。360度撮影可能なカメラで会話の場の中央に配置。
2. GoPro Hero3+ (1920×1080, 59.94fps), 最大2台。170度の視野角を持つカメラで会話者や会話の状況を俯瞰的に記録。

音声 各会話者の音声と会話全体の収録のために以下2種のICレコーダーを使用

1. 個人用ICレコーダー Sony ICD-SX734 (リニアPCM, 16bit, 44.1kHz), 会話者数に応じて2~6台使用⁽⁶⁾。レコーダーをフォルダーに入れ首から下げて収録。マイク設定は「ズームマイク・音声用・ズーム1」(単一指向性マイクでモノラル録音), 感度は事前調査に基づき定めたレベルに固定。
2. 全体用ICレコーダー Sony ICD-SX1000 (リニアPCM, 16bit, 44.1kHz), 1台。会話の場の中央に配置し会話全体を収録。マイク設定は「ステレオマイク」(ICレコーダー先端左右両端のマイクでステレオ録音), 感度は auto に設定。

⁽⁵⁾ 電話等, 非対面での会話の場合はカメラを用いないことが多い。

⁽⁶⁾ 会話者が6名を越える場合や収録の失敗などの理由で, 全ての会話者の音声が個人用レコーダーで収録できていないこともある。原則として貸出は6台としているが, 協力者の要請により最大15台まで使用したことがある。

表4 公開の映像・音声データのファイル形式

映 像	PIXPRO SP360	mp4, H264, 1440×1440, 29.97fps
	GoPro	mp4, H264, 1280×720, 29.97fps
	HX-A500	mp4, H264, 1280×720, 29.97fps
	合成	mp4, H264, 1360×720, 29.97fps
音 声	ICD-SX734	リニア PCM, 16bit, 16kHz, モノラル
	ICD-SX1000	リニア PCM, 16bit, 16kHz, ステレオ
	合成	リニア PCM, 16bit, 16kHz, ステレオ

【移動時収録】 移動時の会話の収録で主として使用

映像 Panasonic HX-A500 (1920×1080, 29.97fps), 1台。散歩などの移動の際に会話者のうち1名が頭に装着して映像を収録。

音声 個人用 IC レコーダー Sony ICD-SX734 (基本収録と同設定)

複数のカメラによる映像がある場合、図3に示すような合成した映像も作成し、個々の映像データと合わせて公開する。また、全体用 IC レコーダーで収録した音声は何らかの理由で公開できない場合⁽⁷⁾、あるいは、公開はできるが質に問題がある場合、個人用 IC レコーダーで収録した各会話者の音声を合成した音源を作成して公開する。同梱する映像・音声データのファイル形式を表4に示す⁽⁸⁾。

映像データの音声については次の通りとする。HX-A500と合成の映像の場合、全体用 IC レコーダーの音声あるいは合成音声(後者を優先)を用いる。PIXPRO SP360とGoProの場合、原則としてそれぞれのカメラで記録した音声を用いるが、音質などに問題がある場合、全体用 IC レコーダーの音声あるいは合成音声(後者を優先)を用いる。

5.2 転記テキスト

図4に、モニター公開で提供する転記テキストの例を示す。映像分析ソフトウェア ELAN や音声分析ソフトウェア Praat などを用い、映像・音声を参照しながら人手で作成している。原則として漢字仮名まじりで表記するが、母音の延伸や発音エラーなど会話で生じる現象については、表5に示す各種タグを用いて表現する。転記テキストの1行は転記単位に相当する。転記単位とは、知覚可能な休止、発話単位の境界、あるいは相互に異なる音種(言語音と笑い、泣き、歌)の境界のいずれかによって区切られる単位である。転記単位ごとに、発話の開始時間と終了時間が割り当てられており、転記テキストから映像・音声データが容易に参照できるようになっている。句点「。」は発話単位の境界を示す。発話単位とは、Japanese Discourse Research Initiative によって策定された「長い発話単位」に相当する(JDRI 2017)。転記テキストの詳細については白田ほか(2018)を参照されたい。

転記テキストは、2種類の単位(転記単位・発話単位)ごとに、CSV形式・EAF形式(ELAN用)・TextGrid形式(Praat用)で提供する。

⁽⁷⁾ 基本収録において録音に失敗した場合、飲食店などでの収録において第三者の会話音声が大きく写り込んでしまうなどの理由により公開すべきではないと判断した場合、及び移動時の収録のように全体用 IC レコーダーでの収録がもともとなされていない場合などが該当する。

⁽⁸⁾ 映像については、何らかの事情で収録時の設定が変更されてしまったために、ここに示す値と異なることがある。

開始時間	終了時間	会話者の ID	テキスト
2502.617	2503.920	IC01_一ノ宮	(U この前) 飲み会どこで飲んだの。
2504.661	2505.651	IC03_さとし	えっと 赤坂。
2507.718	2508.495	N10A_酒井	赤坂の
2508.791	2509.744	IC03_さとし	(L)
2509.287	2510.202	N10A_酒井	料亭。
2510.912	2511.480	IC03_さとし	(L いやいや)。
2511.432	2512.185	IC01_一ノ宮	違う違う。
2512.749	2513.451	IC01_一ノ宮	居酒屋。
2513.641	2514.236	IC03_さとし	(W イサカヤ 居酒屋)。
2515.464	2516.201	IC03_さとし	(X フタヘルモ)。
2516.999	2519.648	IC03_さとし	同期の (D ヒ)(D フ) 同期と二人で飲んだぐらいで。
2519.670	2521.473	Z10A_酒井	芸能人もいっぱい歩いてるんじゃないの。
2521.473	2522.070	Z10A_酒井	そうすっと。
2522.235	2522.864	IC03_さとし	んな見る余裕。
2522.869	2526.534	IC03_さとし	もう 仕事終わったら家帰ることしか頭に (L ないです)。
2523.585	2524.039	Z10A_酒井	ね:。
2526.541	2527.636	IC03_さとし	(L)
2530.214	2531.759	IC01_一ノ宮	前TBSの地下で:
2532.456	2533.398	IC01_一ノ宮	(R (U たか))(W (D サ) さん) ジュリー見た。@ジュリーを見たのは発話者本人

図 4 転記テキストの例

表 5 転記テキストに用いるタグの一覧

■ 非語彙的な発音の変化や言いよどみに関わるもの

タグ	概要	使用例
:	非語彙的な母音の引き伸ばし	すご:い, けれども:
%	非語彙的な音の詰まり	す%ごい, 解%析
(W)	言い誤り・発音の怠け等の一時的な発音エラー	(W コエ これ), (W ギーツ 技術)
(D)	語の言いさし	(D コ) 明日から

■ 韻律・パラ言語的情報に関わるもの

?	疑問上昇調	行きます?, コップ?
(T)	小さい声で発話している箇所	(T これじゃないのか)
(L)	笑いが生じている箇所, あるいは単独の笑い	これ (L なんですけど), (L)
(C)	泣きながら発話している, あるいは単独の泣き	(C なにが), (C)
(S)	歌いながら発話している, あるいは歌詞を伴わない歌	(S ヘイヘイホー), (S)
<>	発音に類する行為のうち会話の流れに関わるもの	<舌打ち>, <咳>, <口笛>

■ 聞き取り等の判断の信頼性に関わるもの

(U)	聞き取りや語の判断に自信がない箇所	(U ジャック) に, (U 国産/特産)
(X)	語が不明な箇所	(X フンジン) 中に, (X ## #)

■ 転記テキストの可読性や内容理解の補助等に関わるもの

(K)	タグ等のために漢字表記できず可読性が落ちる箇所	(K シ:ツ 質) 間, (K ナ%シ 梨)
(M)	音や言葉が言及されており (W) などで対応すると把握しづらい箇所	(M すごい) を (M すっごい) と発音
(O)	一般的に理解が難しい外国語・方言が用いられる箇所	(O ボッソワー), (O # # #)
(Y)	漢字表記の一般的な読みと発音が異なる箇所	(Y ゼツ 舌), (Y センゲン 浅間)
(G)	可読性が低い口語表現	(G 嫌 や), (G もう も)
(F)	「あの」「その」類が連体詞ではなくフィラーとして用いられる場合	(F あの), (F そーの:)

■ 発話単位・転記単位に関わるもの

。	発話単位末	食べます。 , やったけど。 , うん。
+	1 短単位内の知覚可能な休止により転記単位が分割される場合	す+ごい, 神田+川

■ その他

(R)	個人情報などに関わる仮名・伏字処理を行った箇所	(R 国語研究所) の (R 佐藤) さん
@	発話に対するコメント	お願いします:す。@店員への応答

5.3 短単位情報

モニター公開では、長短2種類の形態論情報のうち短単位情報(小椋 2014)を提供する。転記テキストを対象に、形態素解析器 MeCab と形態素解析用辞書 UniDic を用いて自動解析した上で、人手による修正を加えた。自動解析で得られた語形・発音形のうち、語形・発音形が一意に同定できない語(例: 一日「イチニチ/ツイタチ」、日本「ニホン/ニッポン」)については、音を聴取しながら語形・発音形の確認・修正を行った。

形態素解析が対象とするのは、転記テキストにおける母音の引き述べしや音の詰まりのタグをとった形、また言い誤り・発音の怠け等の一時的な発音エラーがある場合はそれを丁寧に発音した場合に生じると想定される語の形である。例えば「けれども:」「く%さい」「(W ギーツ|技術)」であれば、「けれども」「くさい」「技術」が解析の対象となり、発音形にはこうした母音の引き述べしや音の詰まり、言い誤りの情報は記録されない。そこで、転記テキストから得られる実際の「発音」(この例でいえば「ケレドモー」「クッサイ」「ギーツ」)についても「発音形」とは別に公開する。

短単位情報は、CSV 形式で提供するほか、同梱する全文検索システム「ひまわり」パッケージで検索することができる(山口 2018)。またオンライン検索システム「中納言」でも公開する。

5.4 映像・音声・転記テキストのマスキング処理

会話を収録する際、会話者と交わした同意書では、会話者の名前、所属組織名、自宅・所属組織の住所・電話番号の情報、及び会話者が公開を望まない箇所について、それらが分からないよう会話の音声と文字化資料(転記テキスト)を加工することを定めている。音声は当該箇所をピープ音で置換する処理を施す。転記テキストについては、タグ(R)を付けた上で、仮名あるいは伏字(全角アスタリスク)に置換する処理を施す(図4に示す転記テキストの例の最終行を参照)。

同意書において、映像の加工についての条件は付けておらず、原則として会話者の顔にボカシなどの処理を加えずに公開する。ただし、名札や名刺など個人情報を含むものや、収録・公開の同意を得ていない第三者の容貌が映像に写り込むことも少なからずある。そこで、実際に収録された映像・音声データにもとづき法的・倫理的な観点から問題を整理した上で、公開方針を定めた。この方針に従い、必要と判断される箇所について映像にボカシ処理を加える。公開方針の詳細については小磯・伝(2018)を参照を参照されたい。

6. おわりに

本稿では、『日本語日常会話コーパス』のうち、2018年度12月にモニター公開を予定している50時間分のデータの仕様や特徴について報告した⁽⁹⁾。ハードディスク版については、ハードディスク代・郵送代などの実費相当での提供を予定している。モニター公開の最新情報は、以下のページを参照されたい。

<http://pj.ninjal.ac.jp/conversation/cejc-monitor.html>

⁽⁹⁾ 現在、公開に向けて最終的な確認・修正を進めているところであり、ここでの報告と若干異なる可能性がある。

謝 辞

本研究は、国立国語研究所共同研究プロジェクト「大規模日常会話コーパスに基づく話し言葉の多角的研究」によるものである。コーパスの収録にご協力・ご参加くださった皆さまに感謝します。

文 献

- 小磯花絵 (2017). 「『日常会話コーパス』プロジェクト—コーパスに基づく話し言葉の多角的研究を目指して—」 言語資源活用ワークショップ 2016 発表論文集, pp. 114–119.
- 小磯花絵・居關友里子・白田泰如・柏野和佳子・川端良子・田中弥生・伝康晴・西川賢哉 (2017). 「『日本語日常会話コーパス』の構築」 言語処理学会第 23 回年次大会, pp. 775–778.
- 田中弥生・柏野和佳子・角田ゆかり・伝康晴・小磯花絵 (2018). 「『日本語日常会話コーパス』の構築：会話収録法に着目して」 国立国語研究所論集, 14, pp. 275–292.
- 小磯花絵・土屋智行・渡部涼子・横森大輔・相沢正夫・伝康晴 (2016). 「均衡会話コーパス設計のための一日の会話行動に関する基礎調査」 国立国語研究所論集, 10, pp. 85–106.
- Hanae Koiso, Yasuharu Den, Yuriko Iseki, Wakako Kashino, Yoshiko Kawabata, Kenya Nishikawa, Yayoi Tanaka, and Yasuyuki Usuda (2018). “Construction of the Corpus of Everyday Japanese Conversation: An interim report.” *Proceedings of the 11th edition of Language Resources and Evaluation Conference*, pp. 4259–4264.
- JDRI (2017). 『『発話単位ラベリングマニュアル』 version 2.1』. <http://www.jdri.org/open-data/>
- 白田泰如・川端良子・西川賢也・石本祐一・小磯花絵 (2018). 「『日本語日常会話コーパス』における転記の基準と作成手法」 国立国語研究所論集, 15, pp. 177–193.
- 小椋秀樹 (2014). 「形態論情報」 山崎誠 (編) 『書き言葉コーパス 設計と構築』 2 巻講座 日本語コーパス, 第 4 章 pp. 68–88.
- 山口昌也 (2018). 「『日常会話コーパス』活用環境の構築」 言語資源活用ワークショップ 2018 発表論文集.
- 小磯花絵・伝康晴 (2018). 「『日本語日常会話コーパス』データ公開方針：法的・倫理的な観点からの検討を踏まえて」 国立国語研究所論集, 15, pp. 75–89.