

# 『日本語歴史コーパス 明治・大正編VI落語 SP 盤』 音声データアラインメントの概要

2023年1月10日

松崎安子・西川賢哉・石本祐一・中村壮範・小木曾智信

『日本語歴史コーパス 明治・大正編VI落語 SP 盤』(短単位データ 1.0)は、2022年3月に公開された『日本語歴史コーパス 明治・大正編VI落語 SP 盤』(短単位データ 0.9)のテキストデータに対し、落語家の音声データを対応付けアップデートしたものである。本文書では音声データアラインメントの工程と、中納言検索結果を通しての音声データ試聴の方法を解説する。

なお、以下『日本語歴史コーパス明治・大正編VI落語 SP 盤』はCHJ「明治・大正編VI落語 SP 盤」と略し表記する。

## 1. 音声アラインメントの工程

### 1. 1 音声データにおけるノイズ除去

CHJ「明治・大正編VI落語 SP 盤」の音源は岡田コレクション「学術研究用デジタル音源集」(日外アソシエーツ株式会社)を用いている。ただし、音源の音声データには口演者の声以外の音が混入しており、そのままでは聞きづらいものとなっていることから、口演者の声の音質を維持しながらその他の雑音成分を減らすため、Wiener フィルタによる雑音抑圧と調波成分の再生成を組み合わせた手法(Plapous et al. 2006)による雑音抑圧処理を行った。なお、この雑音抑圧処理手法の提案者の一人により MATLAB 言語による実装コードが無償で公開されている。

### 1. 2 Praat を使用した音声アラインメント

前節のように用意した音声データとCHJ「明治・大正編VI落語 SP 盤」の書き起こしテキストとの対応づけにあたっては音声分析用フリーソフトウェア Praat (Version6.2.12 2022年4月17日リリース版)を使用した。

Praat 使用下では、音声データに対し既存のテキストを転記単位で割り付けることで TextGrid (Praat アノテーション形式)を作成した(図1、図2参照)。作成にあたって、国立国語研究所が構築した日本語日常会話コーパス(CEJC)における転記テキストについて規定している白田泰如ほか(2018)の基準を参考とし、落語家の発声の切れ目について次のような場合を転記単位の境界とした。

#### 【転記単位の境界】

1. 知覚可能な休止がある場合
2. 異なる音種(言語音・単独の笑い・泣き・歌・その他)が続く場合
3. 発話単位の切れ目がある場合

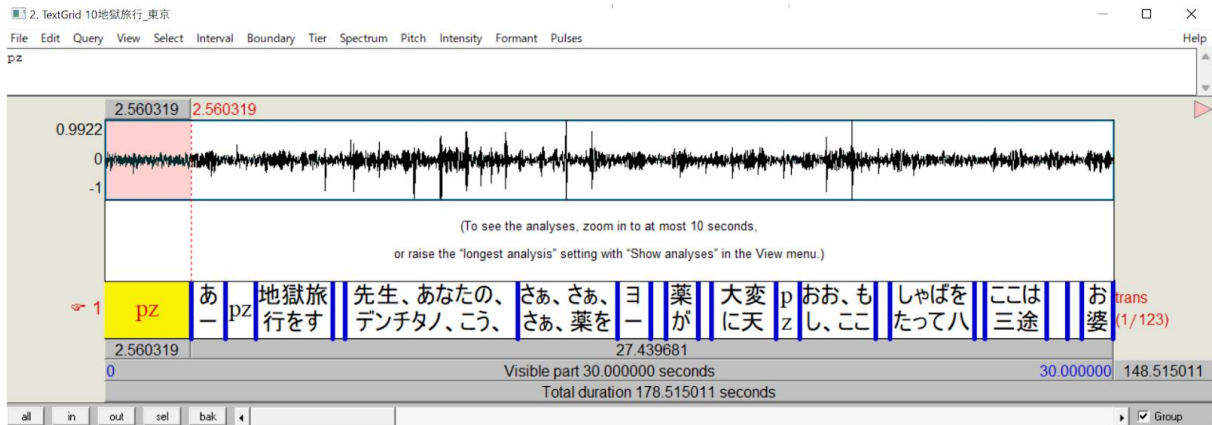


図1 PraatでTextGridを作成

```

1 File type = "ooTextFile"
2 Object class = "TextGrid"
3
4 xmin = 0
5 xmax = 178.5150113378685
6 tiers? <exists>
7 size = 1
8 item []:
9   item [1]:
10     class = "IntervalTier"
11     name = "trans"
12     xmin = 0
13     xmax = 178.5150113378685
14     intervals: size = 123
15     intervals [1]:
16       xmin = 0
17       xmax = 2.5603191509176035
18       text = "pz"
19     intervals [2]:
20       xmin = 2.5603191509176035
21       xmax = 3.5755255765126255
22       text = "あー、お前は、"

```

図2 テキストエディタにTextGridを展開

### 1.3 テキストデータの文字位置と音声データの時刻情報の関連付け

次に、テキストデータの文字位置と音声データの時刻情報とを関連付けを行った。音声データの時刻情報と関連付けるためのテキストの文字位置については、データベースから出力した短単位データを使用した。短単位データには短単位境界とその開始・終了文字位置の情報が付与されており（表1）、前節のPraatでの作業で作成したTextGridと短単位データを照合させることで転記単位での開始・終了の時刻と文字位置の対応表を作成した。

さらに、この対応表に基づいて、形態論情報テーブルに、当該短単位が属する転記単位の開始時刻・終了時刻を付加した（表3）。

表 1 データベースから出力した短単位データ（一部）

file	start	end	boundary	orthToken
10地獄旅行_東京	10	30	B	あー
10地獄旅行_東京	30	40	I	、
10地獄旅行_東京	40	60	I	お前
10地獄旅行_東京	60	70	I	は
10地獄旅行_東京	70	80	I	、
10地獄旅行_東京	80	100	I	地獄
10地獄旅行_東京	100	120	I	旅行
10地獄旅行_東京	120	130	I	を
10地獄旅行_東京	130	150	I	する
10地獄旅行_東京	150	160	I	薬
10地獄旅行_東京	160	170	I	を
10地獄旅行_東京	170	190	I	好む
10地獄旅行_東京	190	200	I	か
10地獄旅行_東京	200	210	I	い
10地獄旅行_東京	210	220	I	。


表 2 開始・終了の時刻と転記単位の文字位置の対応表（一部）

file名	文字 開始位置	文字 終了位置	音声 開始時刻	音声 終了時刻	転記（書字形）	音声ファイル名
10地獄旅行_東京	10	80	2.56	3.576	あー、お前は、	10地獄旅行_東京
10地獄旅行_東京	80	220	4.478	6.786	地獄旅行をする薬を好むかい。	10地獄旅行_東京
10地獄旅行_東京	220	640	7.135	12.269	先生、あなたの、デンタノ、こう、みやげ で評判のお話をするたあを伺いとうがず なあ。	10地獄旅行_東京
10地獄旅行_東京	640	830	12.269	15.238	さあ、さあ、さあ、薬をお上がんさい。	10地獄旅行_東京
10地獄旅行_東京	830	860	15.238	16.225	ヨー。	10地獄旅行_東京
10地獄旅行_東京	860	920	16.687	17.643	薬が早いね、	10地獄旅行_東京
10地獄旅行_東京	920	1100	18.022	19.958	大変に天気は暗くなってきちゃったね。	10地獄旅行_東京
10地獄旅行_東京	1100	1300	20.644	22.909	おお、もし、ここは日本ではありませんよ。	10地獄旅行_東京
10地獄旅行_東京	1300	1500	23.34	25.85	しゃばをたって八幡行きと踏み込んでくる。	10地獄旅行_東京
10地獄旅行_東京	1500	1610	26.169	27.982	ここは三途川の入口だ。	10地獄旅行_東京

表 3 形態論情報と転記単位の開始・終了時刻の対応表（一部）

	subCorpusName	sampleID	start	end	boundary	orthToken	pronToken	reading	lemma	pos	soundFileName	startTime	endTime
1	明治・大正-落語SP盤	60R国遊1903_01010	10	30	B	あー	アー	アー	あー	感動詞-フイラー	10地獄旅行_東京	2.560	3.576
2	明治・大正-落語SP盤	60R国遊1903_01010	30	40	I	、			、	補助記号-読点	10地獄旅行_東京	2.560	3.576
3	明治・大正-落語SP盤	60R国遊1903_01010	40	60	I	お前	オマエ	オマエ	御前	代名詞	10地獄旅行_東京	2.560	3.576
4	明治・大正-落語SP盤	60R国遊1903_01010	60	70	I	は	ワ	ハ	は	助詞-係助詞	10地獄旅行_東京	2.560	3.576
5	明治・大正-落語SP盤	60R国遊1903_01010	70	80	I	、			、	補助記号-読点	10地獄旅行_東京	2.560	3.576
6	明治・大正-落語SP盤	60R国遊1903_01010	80	100	I	地獄	ジゴク	ジゴク	地獄	名詞-普通名詞-一般	10地獄旅行_東京	4.478	6.786
7	明治・大正-落語SP盤	60R国遊1903_01010	100	120	I	旅行	リョウ	リョウ	旅行	名詞-普通名詞-サ変可能	10地獄旅行_東京	4.478	6.786
8	明治・大正-落語SP盤	60R国遊1903_01010	120	130	I	を	オ	ヲ	を	助詞-格助詞	10地獄旅行_東京	4.478	6.786
9	明治・大正-落語SP盤	60R国遊1903_01010	130	150	I	する	スル	スル	為る	動詞-非自立可能	10地獄旅行_東京	4.478	6.786
10	明治・大正-落語SP盤	60R国遊1903_01010	150	160	I	薬	クスリ	クスリ	薬	名詞-普通名詞-一般	10地獄旅行_東京	4.478	6.786
11	明治・大正-落語SP盤	60R国遊1903_01010	160	170	I	を	オ	ヲ	を	助詞-格助詞	10地獄旅行_東京	4.478	6.786
12	明治・大正-落語SP盤	60R国遊1903_01010	170	190	I	好む	コノム	コノム	好む	動詞-一般	10地獄旅行_東京	4.478	6.786
13	明治・大正-落語SP盤	60R国遊1903_01010	190	200	I	か	カ	カ	か	助詞-終助詞	10地獄旅行_東京	4.478	6.786
14	明治・大正-落語SP盤	60R国遊1903_01010	200	210	I	い	イ	イ	い	助詞-終助詞	10地獄旅行_東京	4.478	6.786
15	明治・大正-落語SP盤	60R国遊1903_01010	210	220	I	。			。	補助記号-句点	10地獄旅行_東京	4.478	6.786

## 2. 中納言検索画面からの音声視聴

前節までのように CHJ のテキストと対応付けられた音声情報は、中納言の検索結果画面において直接試聴することができる。図 3 には中納言での検索画面の一例をあげている。このうち、左端のサンプル ID の欄には用例のサンプル ID とともに音声再生用ボタン  が表示されている。これをクリックすると、図 4 のように「(サンプル ID の表示) 再生中」のポップアップが表示され、検索したキーが属する転記単位を含んだ前後 10 秒ほどの音声再生される。



サンプル ID	開始位置	連番	コア	前文脈	キー	後文脈	語彙	語彙	語彙	品詞	活用型	活用形	原文	本文	話者	ジャンル	作品名	成立年	巻名等	作者	生年	底本	ページ番号	底本リンク	参考リンク
60R 圓遊 1903_01010		40	30	1 あー。	お前	「は。 地獄旅行をける 爾を許むか。  #先生。 あななめ。  デシチタル。 」	オマエ	御前	オマエ	代名詞			お前	会話	医者	落語	地獄旅行	1903		三遊亭 圓遊(初代)	1849	地獄旅行 <2651>			
60R 圓右 1911_01049		11830	7220	1 ぼトートー。 だんなんコ ンコ。  #なんだま。  コリヤあ。  #オッホホ。  面白い。  #いや。  あ ー。	お前	「が。 蔵で言ひぬ ら。 おんが。 アー。  太夫で。 は。  待ち やろう。  待ち ー。	オマエ	御前	オマエ	代名詞			お前	会話	三河屋 の旦那	落語	掛取万 歳	1911		三遊亭 圓右(初代)	1860	掛取万歳 <71019>			

図 3 中納言検索結果画面から音声を試聴する (イメージ図)



サンプル ID	開始位置	連番	コア	前文脈	キー	後文脈	語彙	語彙	語彙	品詞	活用型	活用形	原文	本文	話者	ジャンル	作品名	成立年	巻名等	作者	生年	底本	ページ番号	底本リンク	参考リンク
60R 圓遊 1903_01010		40	30	1 あー。	お前	「は。 地獄旅行をける 爾を許むか。  #先生。 あななめ。  デシチタル。 」	オマエ	御前	オマエ	代名詞			お前	会話	医者	落語	地獄旅行	1903		三遊亭 圓遊(初代)	1849	地獄旅行 <2651>			
60R 圓右 1911_01049		11830	7220	1 ぼトートー。 だんなんコ ンコ。  #なんだま。  コリヤあ。  #オッホホ。  面白い。  #いや。  あ ー。	お前	「が。 蔵で言ひぬ ら。 おんが。 アー。  太夫で。 は。  待ち やろう。  待ち ー。	オマエ	御前	オマエ	代名詞			お前	会話	三河屋 の旦那	落語	掛取万 歳	1911		三遊亭 圓右(初代)	1860	掛取万歳 <71019>			

図 4 音声再生中のポップアップ (イメージ図)

### 付記

本コーパスは、国立国語研究所共同研究プロジェクト「開かれた共同構築環境による通時コーパスの拡張」の研究成果によるものである。

### 関連資料

SP 盤貴重音源 岡田コレクション「学術研究用デジタル音源集」日外アソシエーツ株式会社

### 関連 URL

コーパス検索アプリケーション「中納言」 <https://chunagon.ninjal.ac.jp/>

国立国語研究所『日本語歴史コーパス明治・大正編VI落語 SP 盤』

[https://clrd.ninjal.ac.jp/chj/meiji\\_taisho.html#rakugo](https://clrd.ninjal.ac.jp/chj/meiji_taisho.html#rakugo)

Praat 公式サイト <https://www.fon.hum.uva.nl/praat/>

Wiener filter for Noise Reduction and speech enhancement - MATLAB Central

<https://www.mathworks.com/matlabcentral/fileexchange/24462>

#### 参考文献

白田泰如・川端良子・西川賢哉・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

<http://doi.org/10.15084/00001602>

服部紀子・松崎安子・小木曾智信（2022）「『日本語歴史コーパス 明治・大正編 VI 落語 SP 盤』の公開」『日本語学会 2022 年度春季大会予稿集』pp.199-204

松崎安子・西川賢哉・石本祐一・中村壮範・小木曾智信（2022）「『日本語歴史コーパス明治・大正編VI落語 SP 盤』音声アラインメントの公開」pp.115-118

Cyril Plapous, Claude Marro, and Pascal Scalart (2006). "Improved Signal-to-Noise Ratio Estimation for Speech Enhancement." *IEEE Transactions on Audio, Speech and Language Processing*, 14:6, pp. 2098–2108.