

転記テキストの仕様

Version 1.0

小磯花絵[‡]・間淵洋子^{‡,†}・西川賢哉[‡]・斎藤美紀^{‡,††}・前川喜久雄[‡]
[‡] 国立国語研究所, [†] 東京都立大学大学院, ^{††} 東京大学大学院

目次

0	はじめに	1
1	転記テキストの概要	1
2	転記基本単位	3
3	基本形の表記	4
3.1	字種	4
3.2	漢字と平仮名の使い分け	5
3.3	漢字の使い分け	6
3.4	送り仮名	6
3.5	片仮名	6
3.6	口語表現	6
3.7	用語リスト	6
4	発音形の表記	7
4.1	字種	7
4.2	綴り字における母音連鎖	7
4.3	非語彙的な母音・子音の引き延ばし	7
4.4	発音の怠けや転訛・言い間違い	7
4.5	その他	7
5	各種情報のタグ付け	8
5.1	各種タグの説明	9
5.2	タグ間の共起関係	13
6	再朗読の転記について	15
7	対話の転記について	16

0 はじめに

本文書では、『日本語話し言葉コーパス』における転記テキスト [拡張子: trn] の仕様について概説する¹。転記テキストは 3302 講演全てに対して提供される。基本的に全ての転記テキストは同じ仕様に基づいて作成されているが、再朗読と対話については一部特殊な基準が適用されているため、扱いに注意されたい(6節・7節参照)。またコア(overview.pdf 参照)は、他に比べて発音形の精度が高い。これは、コアに対してのみ提供される分節音ラベリング(segment.pdf 参照)の結果を、転記の仕様の範囲内で発音形に反映したためである。それ以外の点に関しても、コアは他に比べて精度は高いが、両者の違いは精度のみであり仕様自体に変更はない。

1 転記テキストの概要

転記テキストの例を図 1 に示す。この例に基づきながら転記テキストの基本的な構成について説明する。

【講演 ID・講演の開始終了位置の情報】

- ◇ 講演 ID: ファイル冒頭に「%講演 ID:」の形式で記載。
- ◇ 講演の開始位置: 講演開始の直前に「%<SOT>」の形式で記載。
- ◇ 講演の終了位置: 講演終了の直後に「%<EOT>」の形式で記載。

【転記基本単位(転記単位)】

[詳細は 2 節参照](#)

- ◇ 原則、0.2 秒以上のポーズによって挟まれた音声の範囲。
- ◇ 転記単位には以下の 4 種類がある。
 - A: 講演者の言語音
 - B: 講演者のボーカル音(笑い声, 泣き声, 咳, 息)
 - C: 上記 A・B 以外の音で特に目立つ音(聴衆の発話や笑い, 拍手, 発表中のデモの音など)
 - D: その他(朗読間違いの箇所・再朗読ファイルに限定して付与)
- ◇ A(講演者の言語音)の場合, 転記単位情報部と発話部(下記参照)から構成。
 - 図 1 中の発話 ID=0265~0269, 0271, 0273~0277
- ◇ B~D の場合, 単位情報部と<咳>のような音声種別を記したタグから構成。 [タグの詳細は表 6 参照](#)
 - 図 1 中の発話 ID=0270 の<咳>, 0272 の<雑音>

【転記基本単位情報部(単位情報部)】単位情報部は以下 3 つの要素から構成。

- ◇ 発話 ID: 4 桁の通し番号(転記単位を開始時刻の早い順に並べ, 0001 から昇順に ID を付与)。
- ◇ 当該転記単位の開始・終了時刻(秒単位)。
 - 対応する音声ファイルの開始時刻を 0(秒)とした場合の時刻。
 - 発話 ID=0265 の場合「00697.054(秒)」が開始時刻「00701.891」が終了時刻
- ◇ 話者 ID(L/R): モノログファイルは L に固定。
 - 対話(話者 2 名)の場合は L と R(音声ファイルのチャンネルに対応)。

【発話部】

- ◇ 基本形と発音形という 2 つの表記方法を用いて発話内容を記す。
- ◇ & の左側に基本形, 右側に発音形が記される。
- ◇ 「文節」に相当する単位で改行。

[文節の基準は bunsetsu.pdf 参照](#)

【基本形】漢字仮名交じりで表記。可読性が高く, 検索に適する。

[基準の詳細は 3 節参照](#)

【発音形】片仮名表記。聞き取れる範囲で発音を忠実に記す。

[基準の詳細は 4 節参照](#)

【タグ】談話に生じる様々な現象を表現する為のタグ。

[タグの詳細は 5 節参照](#)

発話 ID=0265 コ(?ノ): 音声不明瞭

[タグの一覧は表 6 参照](#)

発話 ID=0266 (w コ;コウ): 発音のなまけ・言い間違い

【コメント】 コメント行は%で開始。コメントには以下の 3 種類がある。

- ◇ 講演全体に対するコメント: 講演 ID と<SOT>の間に記述
- ◇ 転記基本単位に対するコメント: 単位情報部の下に記述
- ◇ 局所的な発話に対するコメント: 当該発話行の下に記述

¹ 転記テキストに含まれる情報は全て XML 形式の文書 [拡張子:xml] でも表現されている。また転記テキストに含まれる情報のうち, コメントを除く全ての情報が形態論情報データ [拡張子:sdb,ldb] にも表現されている(詳細は xml.pdf, wdb.pdf 参照)。

2 転記基本単位

本節では、転記テキストの基本単位である「転記基本単位(以下「転記単位」)」の概要について説明する。

【転記単位の種類】

転記単位A：講演者の言語音。	
転記単位B：講演者のボーカル音。	<笑>，<泣>，<咳>，<息>
転記単位C：A・B以外の音で特に目立つ音。	<フロア発話>，<フロア笑>，<拍手>，<デモ>，<ベル>
上記以外で特に目立つ音。	<雑音> … 音種は特定せず一律「雑音」として扱う
転記単位D：朗読間違いの箇所。再朗読に限定。	<朗読間違い> … 扱いはAと同じ。詳細は6節参照。

【単位認定：転記単位A(D)の場合】

- ◇ 原則：言語音が、0.2秒以上の途切れがなく、連続して生じている区間²。
- ◇ 例外：言語的な文末形式(述語の終止形や終助詞など)が存在している場合には、0.05秒以上0.2秒未満の途切れであっても、転記単位を分割する。
- ◇ 時間的關係：転記単位A同士、及び転記単位AとBは、原則として時間的に重複しない。
対話の場合、二人の話者の言語音は重複し得るが、同一話者内での重複は原則通り存在しない。
転記単位Aのうち、タグ<P>で示されるポーズ範囲内に関しては、転記単位Bと重複し得る(下記例参照)。

0041 00098.407-00100.645 L：
(F えーっと) & (F エーッ<P:00098.768-00100.487>ト)
0042 00098.986-00099.228 L:<咳> … 上記ポーズ区間と時間的に重複

【単位認定：転記単位Bの場合】

- ◇ 原則：講演者のボーカル音が0.2秒以上の途切れなく連続して生じている区間。同種の音種(笑い等)毎に認定。
- ◇ 例外1：転記単位Aと時間的に重複する場合(発話の最中に生じる短い笑い等で、言語音が0.2秒以上途切れのない場合)、独立した転記単位Bとして認定せず、転記単位Aの一部とする(下記例参照)。

転記単位Aのみ認定する場合の例：	0001 00001.114-00001.887 L： それでは & ソレデハ<咳> まず & マズ 0002 00002.104-00002.953 L： 本研究の & ホンケンキューノ								
<table border="0"> <tr> <td>言語音</td> <td>###</td> <td>咳</td> <td>言語音</td> </tr> <tr> <td></td> <td></td> <td>(0.2秒以下)</td> <td></td> </tr> </table>	言語音	###	咳	言語音			(0.2秒以下)		
言語音	###	咳	言語音						
		(0.2秒以下)							
転記単位AとBをそれぞれ認定する場合の例：	0001 00001.114-00001.887 L： それでは & ソレデハ 0002 00001.887-00002.030 L:<咳> 0003 00002.100-00002.900 L： まず & マズ								
<table border="0"> <tr> <td>言語音</td> <td>#####</td> <td>咳</td> <td>言語音</td> </tr> <tr> <td></td> <td></td> <td>(0.2秒以上)</td> <td></td> </tr> </table>	言語音	#####	咳	言語音			(0.2秒以上)		
言語音	#####	咳	言語音						
		(0.2秒以上)							

- ◇ 例外2：<息>は、転記単位A(言語音)の直後に後続しており、言語音と音的に切り離せないものに限定する。
- ◇ 時間的關係：転記単位AとB、転記単位B同士は、原則として時間的に重複しない。

【単位認定：転記単位Cの場合】

- ◇ 原則：転記単位Cの対象音が、0.2秒以上の途切れがなく、連続して生じている区間。同種の音毎に認定。
- ◇ 時間的關係：転記単位Cは、A～Dのいずれの単位(自己を含む)とも時間的に重複し得る。
- ◇ 転記単位Cはあくまで談話の流れを理解する為の補足的な情報であり、A・Bと比べて単位認定の精度は落ちる。例えば、転記単位Cの音は、言語音と重複する等の理由で聞き取りが困難な場合が多い。その為、開始・終了位置が厳密に同定できないこともある。また長いデモンストレーションの音などについては、0.2秒の途切れで細かく分割せずに、1つのまとまったデモンストレーションとして分割するなど、状況に応じて適宜単位の認定を行なうこともある。

² 語中語末の促音、及び語中語末の破裂音・摩擦音の閉鎖区間に相当する途切れは除く。

3 基本形の表記

以下に基本形の表記方針の概要を示す³。

3.1 字種

【基本】基本形の表記には以下の字種を用いる（使用範囲の詳細は表3参照）。

漢字：JIS 第1・2水準 (JIS X 0208-1990) の漢字を使用。

平仮名：表1に挙げる仮名のうち「うぁ」「うい」「う」「うえ」「うぉ」以外の仮名文字。

片仮名：表1に挙げる仮名を、3.5節で定める範囲内で使用。

【オプション】上記字種に併記する形で、以下の字種も使用。表記の際にはタグ (A) を利用 (5節参照)。

算用数字：全角の算用数字 (0 1 2 3 4 5 6 7 8 9)

アルファベット：全角のローマ字・ギリシャ文字 (大文字・小文字)⁴

記号：表2に挙げる記号 (いずれも全角) を「使用条件」の項に示す範囲内で用いる。

表1: 表記に利用する仮名文字のリスト

直音系列	拗音系列	周辺のモーラ A	周辺のモーラ B
アイウエオ	ヤユヨ	イエ	
カキクケコ	キャキュキョ		クワ
ガギグゲゴ	ギャギュギョ		グワ
サシスセソ	シャシュショ	シェ	スイ
ザジズゼゾ	ジャジュジョ	ジェ	ズイ
タチツテト	チャチュチョ	ティトゥチェツァツイツェツォ	テユ
ダヂヅデド		ディドゥデュ	
ナニヌネノ	ニャニユニョ	ニエ	
ハヒフヘホ	ヒャヒュヒョ	ヒエ ファ フィ フェ フォ フュ	
バビブベボ	ビャビュビョ	ブイ	ヴァ ヴィ ヴェ ヴォ
パピプペポ	ピャピュピョ		
マミムメモ	ミャミュミョ	ミエ	
ラリルレロ	リャリュリョ		
ワヲ		ウィ ウェ ウォ	
撥音ン	促音ッ	長音ー	

表2: 記号の使用範囲

記号	コード 区点・JIS	使用条件	例
.	0105・2125	数字の小数点, 節番号等 ▷(A) 内	(A 三 . 一四; 3 . 1 4), (A 二 . 三; 2 . 3) 節
	0191・217B	ゼロの意で発話された「マル」 ▷(A) 内	(A 一 八四; 1 0 8 4) ×一丸八四
・	0106・2126	(1) 片仮名語: a. 姓名の間, b. 固有名詞で慣習的に必要性高いもの c. 前置詞・冠詞が入る場合はその前後 (2) 複合語で語の切れ目が分かりづらい箇所 ▷以上基本形で使用	マルコ・ポーロ, チェコ・スロバキア テキスト・トゥー・スピーチ れる・られる型敬語
&	0185・2175	略語 (アルファベット構成) で慣習的に必要性高いもの ▷(A) 内	(A エムアンドエー; M & A)
-	0161・215D	略語 (アルファベット構成) で慣習的に必要性高いもの, 電話・郵便番号, 住所の区切り ▷(A) 内	(A シーディーアール; C D - R), (A 三の九; 3 - 9) (A 一 一三五四二五; 1 1 3 - 5 4 2 5)
×	0163・215F	講演者名や差別語・誹謗中傷の伏せ字として利用 ▷(R) 内	国語研の (R × ×) と申します

³ 本節で述べる方針に従い、揺れが生じないよう統一的に書き記しているが、固有名詞の表記等については若干例外があるので注意されたい。

⁴ ローマ字は区点コード第3区 (0333-0358, 0365-0390), ギリシャ文字は第6区 (0601-0624, 0633-0656) を使用。

ローマ字は原則大文字, ギリシャ文字は原則小文字表記とし, 一般的に表記が固定しているものは逆も可とする。

表3: 基本形における字種の使用範囲^ア 記号 A: . (ピリオド), 記号 B: , 記号 C: & - , 記号 D: (中点), 記号 E: x

出現環境	文字種	平仮名			片仮名			長音	漢字	算用 数字	アルファ ベット	記号					
		直拗	周A	周B	直拗	周A	周B					A	B	C	D	E ^ク	
(F) フィラー			x	x	x	x	x		x	x	x	x	x	x	x	x	x
(F) 感情表現					x	x	x		x	x	x	x	x	x	x	x	x
(D)									x	x	x	x	x	x	x	x	x
(D2) ^イ			x	x					x	x	x	x	x	x	x	x	x
(A) 左			x	x					x	x				x	x		
(A) 右		x	x	x	x	x	x	x	x					x			
(O) 古語方言					x	x	x	x		x	x	x	x	x	x		
(O) 外国語		x	x	x					カ	x	x	x	x	x	x		
(K) 左			x	x	x	x	x	x	x	x	x	x	x	x	x	x	x
(K) 右		x	x	x	x	x	x	x		x	x			x	x		
(M) 音や文字の引用 ^ウ										x	x	x	x	x	x	x	x
その他 和語漢語 ^エ					x	x	x	x		x	x	x	x	x	x		キ
その他 片仮名語		x	x	x			オ		x	x	x	x	x	x	x		

- ア. (O 言い(F んー)やる) のように複数のタグが重複する場合、字種の使用範囲は一番内側のタグ (この場合は F) の制約に従う。ただし (?) については独自の制約はない。(D (? シェ)) であれば (D) の制約に従う。タグ (R) の例外については「ク」参照。
- イ. 片仮名、長音記号の使用は、数字「ゼロ」の言い直し及び外来語系接辞(「センチ」等)に限定。
- ウ. 音や文字に関するメタ的引用は、原則として平仮名で表記。ただし話者が片仮名語や外国語の音、文字を引用している場合は片仮名で表記。文字の引用も、文字種が文脈から分かる場合は、言及・示唆されている文字種(漢字や片仮名等)で表記。
- エ. 「その他」はタグの付与されていない範囲及び上記以外のタグの範囲(音・文字以外の引用の(M) 含む。ただし(?),(R) は除く)。また和語、漢語の周辺のモーラ A・B については、擬音語、擬態語、一部の俗語的表現(「おとつあん」等)に限定。
- オ. 固有名詞に限定(「ルイ・ヴィトン」など)。
- カ. 中国語など漢字を使用する外国語に限定。
- キ. れる・られる型敬語のように、語の切れ目が分かりづらい複合語の場合に限定。
- ク. タグ (R) 内の伏せ字化でのみ使用。記号 E に関してのみ、外側のタグ (R) の制約が内側のタグ範囲にも影響を与える。(R 国研(D 太)太郎) が発表します (R x x (D x) x x) が発表します

3.2 漢字と平仮名の使い分け

【原則】表記が漢字と平仮名で揺れるもので、両表記とも一般的に使用されるものについては、原則漢字を採用⁵。

例えば / x たとえば, 全て / x すべて

【実質名詞・形式名詞】「こと(事)」「もの(物)」「ところ(所)」については一律平仮名で表記⁶。ただし「事柄」や「物語」のように、単語の構成要素である場合にはその限りでない。

【本動詞とテ形複合動詞】「行く」「来る」「置く」「見る」「貰う」「参る」等は、単独で本動詞として出現する場合漢字で表記。「やっておく」や「食べてみる」のように、テ形複合動詞⁷として用いられる場合には、平仮名で表記。

【「言う/いう」の使い分け】以下の組み合わせで出現した場合に限定して平仮名で表記⁸。ただしこの条件を満たす場合であっても、明らかに動作性を有するものについては漢字で表記。

{ 指示副詞: ああ / こう / そう / どう } + { いう } + { 体言 }
{ 引用の助詞: と / って }

【当て字】常用漢字表の付表に記された熟字訓(「玄人」や「相撲」など)のみ使用可能とし、それ以外(「蕎麦」や「矢張り」など)は用いない。

個々の語の表記については、上記の原則に基づきつつ、関連する語との整合性を検討しながら決定する。例えば、動詞「切る」を漢字で表記するならば、「割り切る」や「逆切れ」「締め切り」のように、この語を構成要素として持つ語も同様に漢字で表記する。ただし、関連語との表記の一致を強く推し進め、無理に表記を統一することはしない⁹。

⁵ これは転記テキストに基づく自動形態素解析の精度向上、及び検索の便宜を図って立てた方針である。

⁶ 実質名詞の場合は漢字で、形式名詞は平仮名で表記される慣習が高いが、その区別は非常に難しく書き分けが困難な為、平仮名に統一。

⁷ 長単位 (pos.pdf 参照) において「助動詞相当句」と認定されたものに限定。

⁸ 「言う」という動作が形骸化された用法では、平仮名書きされることが多いが、形骸化しているか否かの判断は非常に難しく、その書き分けは揺れを招き易い為、このように形式的な判断基準を採用。

⁹ 例えば副詞の「とびきり」を、その語の構成要素である「飛ぶ」と「切る」に表記を合わせ「飛び切り」と表記する、ということまではしない。

3.3 漢字の使い分け

【新字と旧字の揺れ】固有名詞も含め、例外なく新字を採用。

万屋 / × 萬屋, 幕末太陽 伝 / × 幕末太陽 傳

【JIS 第 1 水準と第 2 水準の揺れ】JIS 第 1 水準の漢字を採用。

憧れ / × 憬れ, 一獲千金 / × 一攫千金

ただし、著名人の氏名に関しては例外的に使用可とする。

小 淵 元首相 (通常は「淵」に統一)

【同音異義語】書き分けが困難で表記の揺れが生じ易いものについては、片方の漢字で代用可能である場合に限り表記を統一するが、それ以外の同音異義語については書き分ける。

統一の例: 悲しい / × 哀しい, 会う / × 逢う, 尊ぶ / × 貴ぶ
書き分けの例: 表わす / 現わす, 計る / 測る / 図る

3.4 送り仮名

- ◇ 用言で複数の送り仮名の候補がある場合: 送り仮名の字数の多い方を採用 (行なう / × 行う)
- ◇ 名詞で送り仮名の有無に揺れがある場合: 原則送り仮名を付ける (後ろ / × 後)。慣習等で送り仮名を付与しないものは個別に定義。

合図 / × 合い図, 立場 / × 立ち場, 学割 / × 学割り, 番組 / × 番組み, 関取 / × 関取り, など

3.5 片仮名

片仮名表記の対象は以下の通り。

- ◇ 外来語・外国語¹⁰
- ◇ 専門用語や俗語、固有名詞などで片仮名書きの習慣が強いもの (ダフ屋・タンパク質・ト書き)
- ◇ 動物、植物、魚介、虫の名称 (原則片仮名表記とし、例外は個別に定義)

特に外来語については、「ピオラ / ヴィオラ」のような表記の揺れが多く見られる「ピ」と「ヴィ」、「ウイ」と「ウィ」など、表記の揺れが起き易いパターン毎に表記の方針を整理した上で、コーパスに現われる全ての片仮名語を用語リスト (3.7 節参照) に登録。

3.6 口語表現

CSJ では、(1) 音の転訛を伴い、(2) くだけた場面で (意図的に) 使用される表現で、(3) 一個人に限らず幅広く観察されるもの、という 3 つの条件を満たしたものを口語表現として認め、その形式で基本形に書き記してよいものとする¹¹。それ以外は後述のタグ (w) を利用し、そのまま基本形に記すことはしない。

口語表現: しゃあないっちゅうか & シャーナイッチューカ
非口語表現: しょうがないというか & (w チョー; ショー) ガナイ (w ティユ; トユウ) カ

3.7 用語リスト

上記の作業方針に従い、実際の語の表記を定めた用語リスト (約 10 万語) を作成し、オンライン検索や仮名漢字変換用の辞書として利用した。本公開データに用語リストは含まれていないが、将来的には公開する予定である。

¹⁰ 中国の人名地名、及び韓国の著名人 (「金日成」等) の名前に関しては漢字表記とする。

¹¹ 80 時間のデータを書き起こした段階で、そこに出現した口語調の表現を抽出し、上記 3 つの条件と照らし合わせながら、口語表現として登録する語の選別を行なった。表現を個別に登録するのではなく、ある程度体系的に整理した上で、同じ、あるいは類似した現象は、できるだけ同様の扱いをするように心掛けた。

4 発音形の表記

以下に発音形における表記方針の概要を示す。

4.1 字種

実際に発音された音を，表 1 に挙げる片仮名 (ヂツ以外) を用い，表 4 に示す範囲内のできる限り正確に書き表わす。

表 4: 発音形における字種の使用範囲 (直音等: 直音・拗音・撥音・促音・長音)

	直音等	周辺のモーラ A	周辺のモーラ B
(F) 感情表現, (D), (O), (W) 左項, (M) 音や文字の引用, 片仮名語 (F) フィラー・上記以外の和語漢語		×	×

4.2 綴り字における母音連鎖

「かあさん (kaasaN)」のように，綴り字において母音が連鎖しており，その母音連鎖部の発音が [ka:saN] のように長音化している場合については，以下のように対処する。

タイプ A 長音化が (狭義の) 形態素内で生じるもの¹²: 長音表記・母音表記共に認める (表 5 の ア参照)

カーサン・カアサン (母さん), チーサイ・チイサイ (小さい)・ ケーロ・ケエロ・ケイロ (経路), カナシー・カナシイ (悲しい)

タイプ B 長音化が 2 つの形態素に跨がるもの: 長音表記を認めず，一律母音表記とする (表 5 の イ参照)

ダイイチ・xダイーチ (第 / 一), ケイロ・xケーロ (毛 / 色), オウミ・xオーミ (お / 産み), モノオ・xモノー (もの / 色)

表 5: 母音連鎖パターンにおける発音の扱い

母音連鎖パターン	語彙例	タイプ A: 形態素内の場合			タイプ B: 形態素を跨ぐ場合			
		長音	母音 1 ア	母音 2	語彙	長音	母音 1	母音 2
a a	母さん	カーサン	(カアサン)	-	油揚げ	xアブラーゲ	アブラアゲ	
i i	小さい	チーサイ	(チイサイ)	-	第一	xダイーチ	ダイイチ	-
u u	空気	クーキ	(クウキ)	-	安売り	xヤスーリ	ヤスウリ	-
e e	姉さん	ネーサン	(ネエサン)	ネイサン	影絵	xカゲー	カゲエ	xカゲイ イ
o o	大きい	オーキイ	(オオキイ)	オウキイ	お教え	xオーシエ	オオシエ	xオウシエ イ
e i	経路	ケーロ	(ケエロ)	ケイロ	毛色	xケーロ イ	ケイロ	xケエロ イ
o u	講師	コーシ	(コオシ)	コウシ	子牛	xコーシ イ	コウシ	xコオシ イ

ア: 括弧内に示す母音表記は，一音一音区切って発音する場合など，分節音のレベルで明らかに母音が発音されている場合に限定

イ: この発音が実際に生じた場合には，(w) 対処とする: 毛色 & (w ケエロ;ケイロ), お産み & (w オーミ;オウミ)

4.3 非語彙的な母音・子音の引き延ばし

母音や子音の引き延ばし現象のうち「コレカラー」や「スゴーイ」「サツサガ」「カイツセキ (解析)」のように，語彙的には引き延ばしは存在しないが，パラ言語的意味等が付与されることによって，あるいは言い淀みによって，一時的に引き延ばされているものについては，タグ<H>，<Q>を用い「コレカラ<H>」や「サ<Q>スガ」のように表記する。ただし「ヤハリ」が「ヤツパリ」「ヨホド」が「ヨツポド」となるように，周囲の音が規則的に転訛しているものについては，<Q>ではなく「ッ」で対処し，基本形にも「やっぱり」のように「っ」を含めて表記する。

4.4 発音の怠けや転訛・言い間違い

「国語研」を「コッコケン」「手術」を「シジツ」「あります」を「アリマ」「共分散」を「キョーブサン」と発音するなど，発音の怠けや転訛，言い間違いなどが生じた場合には，実際に発音された音を可能な限り正確に書き表わした上で，5 節で詳述するタグ (w) を利用して，丁寧に発音された場合に生じるであろう音を併記する。

4.5 その他

助詞の「は」「を」「へ」については実際の発音である「ワ」「オ」「エ」を用いて表記「縮む (ちぢむ)」や「続く (つづく)」など，現代仮名遣いの上では「ぢ」「づ」を用いて表記する語であっても，発音形では一律「ジ」「ズ」に統一。

¹² 原則「最小単位」(pos.pdf 参照) と一致 (最小単位における固有名詞等の例外は除く)。また，助動詞の「う」は例外的に本タイプに含め，長音表記を認める。

5 各種情報のタグ付け

転記に用いるタグの一覧を表 6 に示す。

表 6: 転記テキストに用いるタグの一覧

転記単位 A (講演者の言語音) に出現するタグ：文字範囲を指定し、その範囲の性質に言及するタイプ				
タグ	タグの概要	使用例	付与対象 [*1]	範囲 [*2]
(F)	フィラー、感情表出系感動詞	(F あの), (F うわ)	基・発	- . .
(D)	言い直し, 言い淀み等による語断片	(D こ) これ, (D ぼい) 子音の	基・発	- . .
(D2)	助詞, 助動詞, 接辞の言い直し	そこ (D2 が) に, (D2 不) 不自然	基・発	- . .
(?)	聞き取りや語彙の判断に自信がない場合	(? タオングー), (? 堆積, 体積)	基・発 [*3]	文 . .
(M)	音や言葉に関するメタ的な引用	助詞の (M は) は (M わ) と発音	基・発	文・転
(O)	外国語や古語, 方言など	(O ザッツファイン)	基・発	文・転
(R)	講演者の名前, 差別語, 誹謗中傷など	国語研の (R x x) です	基・発	文・転
(X)	非朗読対象発話 (朗読における言い間違い等)	(X 実際は) 実際には,	基・発 / 再朗読	文 . .
(A)	アルファベットや算用数字, 記号の表記	(A シーディーアール; C D - R)	基 . .	- . .
(K)	何らかの原因で漢字表記できなくなった場合	(K たち (F える) ばな橋)	基 . .	- . .
(W)	転訛, 発音の怠けなど, 一時的な発音エラー	(W ギーツ; ギジュツ)	. . 発	- . .
(B)	語の読みに関する知識レベルの言い間違い	(B シブタイ; ジュータイ)	. . 発	- . .
(笑)	笑いながら発話	(笑 ナニガ)	. . 発	文 . .
(泣)	泣きながら発話	(泣 ドンナニ)	. . 発	文 . .
(咳)	咳をしながら発話	シャ(咳 リン) ノ	. . 発	文 . .
(L)	ささやき声や独り言などの小さな声	(L アレコレナンダッケ)	. . 発	文 . .

[*1] 基 / 発: 当該タグが基本形 / 発音形に出現 [*2] 文 / 転: 当該タグの括弧の範囲が複数の文節 / 転記単位を跨ぐことがある
[*3] 状況に応じて, 基本形のみ, 発音形のみ, あるいは両形に出現。詳細は (?) の項参照。

転記単位 A (講演者の言語音) に出現するタグ：音や事象自体を記号で表現するタイプ				
<FV>	ボーカルフライ等で母音が同定できない場合	だから<FV> & タカラ<FV>	基・発	
<VN>	「うん/うーん/ふーん」の音の特定が困難な場合	(F うん) & (F <VN>)	. . 発 / 対話	
<H>	非語彙的な母音の引き延ばし	ソレデ<H>, ス<H>ゴイ	. . 発	
<Q>	非語彙的な子音の引き延ばし	カイ<Q>セキ, スゴ<Q>イ	. . 発	
<笑>	言語音と独立に講演者の笑いが生じている場合	ガクセー<笑> ノ	. . 発	
<咳>	言語音と独立に講演者の咳が生じている場合	デ<咳>	. . 発	
<息>	言語音と独立に講演者の息が生じている場合	ツマリ<息>	. . 発	
<P>	0.2 秒以上のポーズ	オ<P:00453.373-00454.013>モイ	. . 発	

転記単位 B (講演者の発するボーカル音) に出現するタグ：音や事象自体を記号で表現するタイプ				
<笑>	講演者の笑い声		単位情報部	
<泣>	講演者の泣き声		単位情報部	
<咳>	講演者の咳		単位情報部	
<息>	講演者の息 (言語音と連続した息で, 波形上切り離せないものに限定)		単位情報部	

転記単位 C (A・B 以外の音で特に目立つ音) に出現するタグ：音や事象自体を記号で表現するタイプ				
<フロア発話>	聴衆 (司会者等も含む) の発話		単位情報部	
<フロア笑>	聴衆の笑い		単位情報部	
<拍手>	聴衆の拍手		単位情報部	
<デモ>	講演者が発表中に用いたデモンストレーションの音声		単位情報部	
<ベル>	学会講演時に発表時間を知らせる為に鳴らすベルの音		単位情報部	
<雑音>	上記以外の音で特に目立った音		単位情報部	

転記単位 D (その他) に出現するタグ：音や事象自体を記号で表現するタイプ				
<朗読間違い>	転記単位全体が再度読み直された場合		単位情報部 / 再朗読	

5.1 各種タグの説明

本節では表 6 に示したタグのうち、転記単位 A に用いるタグについて概説する¹³。

【タグ (F)】 本タグは、以下に挙げる「フィラー」と「感情表出系感動詞」に対して付与する¹⁴。

フィラー：

- ◇ 場繋ぎ的な機能を持つ表現。語彙を以下のように限定し、その範囲内で場繋ぎの機能を有する場合に付与。
- ◇ 下記以外の表現は、たとえ場繋ぎの機能を有するものであっても (F) は付与しない。
- ◇ 「あの/その」に関して、フィラーと連体詞とで迷う場合には、原則 (F) を付与した上で迷った旨をコメント。

フィラー表現

あ(一)、い(一)、う(一)、え(一)、お(一)、ん(一)、と(一)*、ま(一)*、

う(一)ん、あ(一)(ん)の(一)*、そ(一)(ん)の(一)*、

う(一)ん(一)つと(一)*、あ(一)つと(一)*、え(一)つと(一)*、ん(一)つと(一)*

上記全てのフィラーとの組み合わせ：～ですね(一)、～っすね(一)

[例] あのですね、えーとっすねー

上記「*」印のフィラーとの組み合わせ：～ね(一)、～さ(一)

[例] まーねー、うーんとさー

括弧内は任意

[例] あの、あーの、あーの、あーんのー

感情表出系感動詞：

- ◇ 驚いた時や落胆した時などに発する感動詞。
- ◇ 語彙の限定や表記の統一は行なわず、長音や促音も含め聞こえた通り表記。

例：(F あーあ)、(F あっちゃー)、(F ぎょっ)、(F うおー)

【タグ (D)】

- ◇ 本タグは、以下のケースで生じる「語の断片¹⁵」に対して付与。

言い直しに伴う語断片：「あたら 最新の研究で」の例に見られるように、何かを言い掛け（「あたら」）、それを別の表現（「最新の」）で言い替えた場合の、言い掛けの部分（「あたら」）に付与。

その他の語断片：「その ん こと」など、言い直しに伴う語断片というよりは、発声上の問題で生じたと考えられる断片的な音声に対しても、同様に本タグを付与。

- ◇ 「スライド (F えーと) プロジェクターで」や「それ それについて その問題について考えてみると」のように、言い直された対象が短単位の断片でない場合には、(D) は付与しない。
- ◇ 「ダイガ ノ カイギ (大学の会議)」の「ダイガ」のように、語の断片であっても、言い直されずにそのまま発話されたものについては、「大学の & (W ダイガ; ダイガク) ノ」のように (W) で対処。

(D 形) 形式の、(D だ) (D だいが) 大学の学部の会議、(D じゅ) 精度の上で重要な

(D んみ) 考えもしないし、(D てん) 元気な子が、洗濯機で (D つ) やった

【タグ (D2)】

- ◇ 本タグは、言い直しにおける言い掛け部・訂正部が共に助詞・助動詞・接頭辞・接尾辞から構成される場合、及び言い掛け部・訂正部が共に数字の場合（数字内の数字の言い直し）に付与する。
- ◇ 「ここ か から」のように、機能語の断片が言い直されている場合は (D) を付与する。

評価値 (D2 が) の数値が、組み合わせ (D2 や) (D2 は) については、ここ (D2 です) ですね、広さ (D2 や) たるや

(D2 第) 第一関門を、西洋 (D2 的) (F えー) (D ぶ) 風というか、(A 千 (D2 八百) 九百四十; 1 9 4 0) 年

学習データ (D2 が) (D こん) (F え) (D しゅ) の収集が困難な

¹³ 転記単位 B・C については 2 節を、また転記単位 D については 2 節と 6 節を参照されたい。

¹⁴ 対話 (講演 ID が D から始まるファイル) では応答表現に対しても付与する。詳細は 7 節参照。

¹⁵ ここでいう「語」とは、原則「短単位」(pos.pdf 参照) を指す。ただし外来語に関してのみ語の定義が異なる「カード/ゲーム」のような 2 つの単語から構成される複合語に関して、短単位では 1 つの単位と見なすのに対し、転記では 2 つの単位とする。

【タグ (?)】

- ◇ 本タグは、音の聞き取りや語彙の同定、漢字表記などに自信がない場合に付与する。以下3つの形式を取る。
 1. 値1つ(デフォルト): (? 字数) の & (? ジスー) ノ
 2. 値複数(複数の候補が想定): (? 次数, 実数) & (? ジスー, ジッスー)
 3. 値なし(全く不明な場合): (?) で & (?) で
- ◇ 音の聞き取りが曖昧なのか、それとも語彙や漢字の同定が曖昧なのかにより、本タグを基本形と発音形のどちら(あるいは両方)に付与するか決まる。

(? 字数) の & (? ジスー) ノ ... 音の聞き取り曖昧で語彙も不確定
 それで & (? ソレデ) ... 音の聞き取り曖昧だが文脈から語彙確定
 (? せんさ) 空間の & センサクーカン ... 音は明瞭だが語彙が曖昧。候補があれば漢字、なければ平仮名
 (? 大賞, 対象) の & タイショーノ ... 音は明瞭だが語彙が曖昧。複数の候補あり。

【タグ (M)】

- ◇ 本タグは、音や言葉自体が言及の対象となるようなメタ的引用のうち、特に以下のパターンに限定して付与する¹⁶。
 1. 以下の要素が単独で引用される場合：
 - 単語未満の要素(音、文字、語の断片、接辞)・非文:(M あ)という文字は(M め)とよく似ている
 - 機能語(助詞、助動詞):助詞の(M は)は(M は)と書くが発音は(M わ)
 - 活用する自立語のうち言い切り(終止形、命令形)以外:(M 走ら)は動詞の未然形
 - 活用する自立語の言い切りが引用の「と/って」以外に接続する場合:(M いらっしゃる)が後続すると
 - 引用部の終端が上記1の要素の場合(「と/って」に接続する終助詞・助動詞言いきり形を除く):(M 僕が)が来ると
 - 連体詞・副詞・感動詞・接続詞:(M そして)と接続詞を置くことで
 2. 引用部の始端が機能語または語の断片の場合:(M という)や(M といった)という表現を使って
 3. 引用部の終端が上記1の要素(終助詞・助動詞言い切り形以外)の場合:(M 僕が)が来ること
 4. 記号(句読点・括弧等)の読み上げ:走ります(M 丸)と句点を付ける、平仮名で(M こっかい(M 中括弧閉じる)では)で
 5. 表記や音に関する引用は言及された通り表記した上で(M)を付与:(M カッコ)と片仮名書き,(M えきー)と延ばして発音

【タグ (O)】

- ◇ 本タグは、外国語や古語、方言などCSJが対象とする現代共通日本語から外れている(可能性のある)箇所が付与。
- ◇ 外来語として定着している語であっても、外国語風の発音をしている(通常の日本語の音韻体系から外れている)と考えられるものについては本タグを付与。
- ◇ (O)内の表記は、3節に示す表記の原則や用語リストに定める表記から外れることもある。

【タグ (R)】

- ◇ 本タグは以下の場合に付与する：
 1. 講演者が特定できる情報(講演者・共著者の氏名、学会講演の場合発表タイトル、前後の講演者の氏名¹⁷等)
 2. 一般人が特定できる情報(著名人研究者を除く一般人の氏名・ファーストネームのみの場合は付与せず)
 3. 誹謗中傷・差別語のうち、特に問題になると判断されたもの
 4. 講演者が非公開を希望した箇所
- ◇ 本タグが付与された範囲は伏せ字化される。伏せ字化に際しては、伏せ字化以前の文字の数だけ記号「×(全角)」を繰り返す。(R)内に含まれる記号はそのままとする。また音声ファイルに対しては、本タグの対象となる転記単位の全体を白色雑音で置換する処理を施している。つまり、転記単位の一部でも本タグの範囲であれば、その転記単位全体が白色雑音で置換される。

伏せ字前:(R 話し言葉の(F えー)イントネーション)というタイトルで国語研の(R 佐藤)が発表します
 伏せ字後:(R ×××××(F ××)×××××××)というタイトルで国語研の(R ××)が発表します

【タグ (X)】

- ◇ 本タグは、朗読原稿にない余分な発話に対して付与する。詳細は6節参照。

¹⁶ メタ的引用の前後では、通常の単語の接続パターンから逸脱することもあり、後の自動形態素解析の処理で問題が生じる恐れがある為、特に問題となる可能性の高いパターンに限定して本タグを付与するという方針をとる。

¹⁷ CSJに収録されている場合のみ【例】先程の(R ××)さんの発表にもありましたように

- ◇ 再朗読 (講演 ID が R00 から始まるファイル) に限定して付与。

【タグ (W)】

- ◇ 本タグは「ガクジツ (学術)」や「コッコケン (国語研)」のように、発音の怠けや音の転訛、言い間違いなどが生じた場合に付与する。
- ◇ (W ガクジツ;ガクジュツ) のように、セミコロンの左側に、実際に発音された音を可能な範囲で正確に書き表わすと同時に、セミコロンの右側には、丁寧に発音された場合に生じる (と予想される) 音を併記する。
- ◇ 本タグの範囲は短単位 (9 頁の注 15 参照) とし、その範囲でタグをくくり直す (下記例「そうすると」参照)。ただし、語の融合などが生じて切り離せない場合はその限りでない (「今度は」参照)。

共分散	& キョー (W プサン;ブンサン),	左へ行って	& (W ミダリ;ヒダリ) エイッテ
手術すると	& (W シジツ;シュジュツ) スルト,	けれども	& (W ケード;ケレド) モ
そうすると	& (W ホ;ソー)(W ス;スルト)	今度は	& (W コンダー;コンドウ)

- ◇ タグ (F), <FV> は原則、基本形、発音形の両形に出現するが、(W) の左項 (実際の発音部) に (F), <FV> が生じた場合、下記の通り (W) の右項、及び基本形にこれらのタグは現れない。そのため、(F), <FV> を完全に抽出する為には、基本形 (だけ) ではなく発音形を参照する必要がある。

チェコ・スロバキアの & (W セコスロバ (F アノー) スロバキア;チェコスロバキア) ノ

- ◇ 「アメリカの大統領 エリツイン は」や「これが やります」のように、世界知識や文法のレベルで間違っている、あるいは適格でないものについては修正の対象としない。
- ◇ 本タグと (?) との組み合わせで典型的な例を幾つか示す。

(?手)	& (W 工;(? テ))	... 音は確実。訂正候補は不確実 (W の右項と基本形の両方に?付与)
(?手)	& (W (? 工);(? テ))	... 音も訂正候補も不確実 (W の両項と基本形に?付与)
手	& (W (? 工);テ)	... 音は不確実だが訂正候補は確実 (W の左項のみ?付与)
(?絵,手)	& (? 工,(W 工;テ))	... 音は A と確実だが語彙が A か B かで迷う (?の複数候補の一部に W)
(?毛,手)	& (W 工;(? ケ,テ))	... 訂正候補が不確実で複数の可能性あり (W の右項に?の複数候補)

【タグ (B)】

- ◇ 本タグは、漢字の読みに関する知識レベルの言い間違いに付与する¹⁸。
- ◇ (B シブタイ;ジュータイ) のように、セミコロンの左側に実際に発話された読みを書き表わすと同時に、セミコロンの右側には正しい読みを併記する。
- ◇ 原則として、本タグの範囲は短単位 (9 頁の注 15 参照) とする。

半ば & (B ハンバ;ナカバ), 借家であるか & (B シャクイエ;シャクヤ) デアルカ

- ◇ 以下に該当するものには、(B) を付与しない。
 - 専門的な慣用読み、及び意味の明確化等の理由で本来の (辞書上の) 読みと異なる読みが慣用化している場合

両耳で & リョージデ	... × (B リョージ;リョーミミ) デ
裏面が & ウラメンガ	... × (B ウラメン;リメン) ガ

- 人名における連濁の有無、音転、読み入れ換え

谷川俊太郎 & タニガワトシタロー ... × (W タニガワ;タニカワ) (B トシタロー;シュンタロー)

- 「日本」における「ニホン」「ニッポン」の入れ替え、行政区画名の「町」における音訓入れ替え

大日本帝国憲法 & ダイニホンテーコクケンポー	... ×ダイ (B ニホン;ニッポン) テーコクケンポーガ
櫛形町 & クシガタチョー	... ×クシガタ (B チョー;マチ)

¹⁸ 「メグスリ (目薬)」を「マグスリ」と発音する (母音 /e/ /a/) ; 「ムジン (無人)」を「ブジン」と発音する (子音 /m/ /b/) , といったように、母音 1 つ、ないし子音 1 つのみの交替の場合は、語の読みに関する知識レベルの言い間違いではなく音声レベルの言い間違いの可能性もある為、(B) ではなく (W) 対処とする [例] 無人 & (W ブジン;ムジン) × (B ブジン;ムジン)

【タグ (A)】 本タグは、アルファベット・算用数字・記号 (いずれも全角) を表記する為に用いる。これらの字種は、以下に示すように本タグを利用し漢字仮名に併記する形で記述。

(A 千九百九十五; 1995) 年, (A 二十六.三五; 26.35), (A シーディーアール; CD-R),
(A トリプルエーアイ; AAAI), (A ハンドレッドベースティー; 100BASE-T)

【タグ (K)】 本タグは、本来基本形で漢字や一部の記号 (と.) で表記するはずの語が、(F) や <FV>, 値なしの (?) が途中に生じた為に表記できなくなった場合に用いる。

(K ひ (F いー) だり; 左), (A 三 (K ま<FV>る;) 四; 304)

【タグ (笑), (泣), (咳)】

◇ 本タグは、話者の笑い声・泣き声・咳と発話が、同時もしくは入り混じりながら進行している区間 (笑いながら発話している, など) に付与する。

【タグ (L)】

◇ 本タグは、前後の音声と比べて、かなり小さな声で発話されている区間に付与する。必ずしも独り言であるとは限らない。

【タグ <FV>】

◇ 本タグは、強いきしみ発声 (ボーカルフライ) などによって、母音が明確に同定できない場合に用いる。
◇ きしみ音であっても、母音が同定できる場合には、本タグは使用せずその母音を記す。

【タグ <VN>】

◇ 本タグは、「うん/うーん/ふーん」の音の認定が困難なものに対して用いる。詳細は 7 節参照。
◇ 対話 (講演 ID が D で始まるファイル) に限定して付与。

【タグ <H>, <Q>】

◇ 本タグは、母音 (<H>)・子音 (<Q>) の引き延ばし現象のうち「コレカラー」「スゴイー」や「サッスガ」「カイツセキ (解析)」のように、本来語彙的には引き延ばしは存在しないが、パラ言語的意味等が付与されることによって、あるいは言い淀みによって、一時的に引き延ばされているものに対して付与する (4.3 節参照)。

すごい & ス<Q>ゴ<H>イ ... × すっごーい & スッゴイー

◇ タグ (W), (D), 及び感情表出系感動詞 (F) の内部では、原則本タグは用いず「ー」「っ」で表記 (表 8・10 参照)。

解析が & (W カイツシェーキ; カイセキ) ... × 解析が & (W カイ<Q>シェ<H>キ; カイセキ)
(D ざっつ) 雑音の & (D ザツツ) ザツオンノ ... × (D ざつ) 雑音の & (D ザ<Q>ツ) ザツオンノ
(F うわわっわー) & (F ウワワッワー) ... × (F うわわわ) & (F ウワワ<Q>ワ<H>)

【タグ <笑>, <咳>, <息>】 本タグは、言語音とは独立して生じる話者の笑い・咳・息に付与する。

【タグ <P>】

◇ 本タグは、短単位 (9 頁の注 15 参照) の内部に生じる 0.2 秒以上のポーズに対して付与する¹⁹。
◇ 以下に示す例のように、ポーズの開始・終了時刻の情報 (ミリ秒単位, 形式は時間情報部と同じ) を値に持つ。

0122 00331.802-00334.403 L:
最近 & サイキン
半年 & ハン<P:00333.068-00333.442>トシ<H>
0124 00335.025-00335.170 L:

¹⁹ 本来、0.2 秒以上のポーズで転記単位を分割するが、ポーズが短単位内部に生じた場合には、転記単位を分割せず 1 転記単位とした上で、本タグを利用して 0.2 秒以上のポーズとして時間の情報を残す。

5.2 タグ間の共起関係

本節では、複数のタグが重複して付与される場合の制約について説明する。タグ間の共起関係は、以下に挙げる「同一範囲」と「入れ子」の2種類がある。CSJでは、(1)2つのタグが共起するかしないか、(2)同一範囲の場合どちらのタグが外側になるか、の2点に関して制約を与えている。同一範囲、入れ子毎に表7, 8にまとめる。

【同一範囲】複数のタグが同じ文字範囲に付与：(M (A エス; S)), (D (? かわ)), (F (W トノー; (? ソノー, アノー)))

【入れ子】あるタグの付与範囲の一部に別のタグが付与：(L (? 市外, 市内) 通話), (K たち (F ぃー) ばな橘)

同一範囲のタグの優先順位 (右にある記号が外側)： A, K < ?単 < W < B < F, D, D2 < R < O < M < X < L < 笑

表 7: 【同一範囲】タグ間の共起関係 (: 共起, : 条件付きで共起, x : 共起しない, - : タグ制約上そもそも共起しない)

内\外 タグ	A左	A右	K左	K右	?単	W左	W右	B左	B右	Fフ	F感	F応	D	D2	R	O	M	X	L	笑	?複
A	-	-	x	x	-	-	-	-	-	x	x	x	x	x	-	x	-	-	-	-	-
K	-	-	-	-	-	-	-	-	-	x	x	x	x	-	-	-	-	-	-	-	-
?単	-	-	-	-	-	-	-	-	x	-	-	-	-	-	-	-	-	-	-	-	-
W	-	-	-	-	-	-	-	-	x	-	x	-	x	-	-	-	-	-	-	-	-
B	-	-	-	-	-	-	-	-	-	x	x	x	x	-	-	-	-	-	-	-	-
Fフ	-	-	-	-	-	-	-	-	-	-	-	-	x	x	x	x	-	-	-	-	-
F感	-	-	-	-	-	-	-	-	-	-	-	-	x	x	x	x	-	-	-	-	-
F応	-	-	-	-	-	-	-	-	-	-	-	-	x	x	x	x	-	-	-	-	-
D	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	x	-	-	-
D2	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
R	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
O	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
M	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
X	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
L	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	x
笑	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	x
?複	x	x	x	x	-	-	-	x	x	x	x	x	x	x	x	x	x	x	x	x	-
(?)	x	x	x	x	-	-	-	x	x	x	x	x	x	x	x	x	x	x	x	x	-
<FV>	x	x	x	x	x	-	-	x	x	x	x	x	x	x	x	x	x	x	x	x	x
<VN>	-	-	-	-	-	-	-	x	x	x	x	x	x	x	x	x	-	-	-	-	-
<笑>	-	-	-	-	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x
<H>	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
<Q>	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
<P>	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-

表 8: 【入れ子】タグ間の共起関係 (: 共起, : 条件付きで共起, x : 共起しない, - : タグ制約上そもそも共起しない)

内\外 タグ	A左	A右	K左	K右	?単	W左	W右	B左	B右	Fフ	F感	F応	D	D2	R	O	M	X	L	笑	?複
A	-	-	x	x	-	-	-	-	-	x	x	x	x	x	-	x	-	-	-	-	-
K	-	x	-	-	-	-	-	-	-	x	x	x	x	x	-	-	-	-	-	-	-
?単	-	x	x	x	-	-	-	-	x	-	-	-	-	-	-	-	-	-	-	-	-
W	-	-	-	-	-	-	-	x	x	-	x	x	x	-	-	-	-	-	-	-	-
B	-	-	-	-	-	x	x	-	-	x	x	x	x	x	-	-	-	-	-	-	-
Fフ	-	x	-	x	-	-	-	-	x	-	-	-	x	x	-	-	-	-	-	-	-
F感	-	x	-	x	-	-	x	-	x	-	-	-	x	x	-	-	-	-	-	-	-
F応	-	x	-	x	-	-	-	-	x	-	-	-	x	x	-	-	-	-	-	-	-
D	-	x	x	x	-	x	x	x	x	x	x	x	-	x	-	-	-	-	-	-	-
D2	-	x	x	x	-	x	x	x	x	x	x	x	-	-	-	-	-	-	-	-	-
R	x	x	x	x	-	x	x	x	x	x	x	x	x	x	-	-	-	-	-	-	-
O	x	x	x	x	-	x	x	x	x	x	x	x	x	x	-	-	-	-	-	-	-
M	x	x	x	x	-	x	x	x	x	x	x	x	x	x	-	x	-	-	-	-	-
X	-	x	x	x	-	x	-	x	x	-	x	x	x	-	-	-	-	-	-	-	-
L	-	-	-	-	-	-	x	-	x	-	-	-	-	-	-	-	-	-	-	-	-
笑	-	-	-	-	-	-	x	-	x	-	-	-	-	-	-	-	-	-	-	-	-
?複	x	x	x	x	-	x	-	x	x	-	x	x	x	x	-	-	-	-	-	-	-
(?)	x	x	x	x	-	x	-	x	x	-	x	x	x	x	-	-	-	-	-	-	-
<FV>	-	x	-	x	-	-	-	x	x	-	-	-	-	-	-	-	-	-	-	-	-
<VN>	-	-	-	-	-	-	-	x	x	x	x	x	x	x	-	-	-	-	-	-	-
<笑>	-	-	-	-	-	-	-	x	x	-	-	-	-	-	-	-	-	-	-	-	-
<H>	-	-	-	-	-	-	x	-	x	-	-	-	-	-	-	-	-	-	-	-	-
<Q>	-	-	-	-	-	x	x	-	x	-	x	-	x	-	-	-	-	-	-	-	-
<P>	-	-	-	-	-	-	x	-	x	-	-	-	x	-	-	-	-	-	-	-	x

Fフ...タグ F・フィラー; F感...タグ F・感情表出系感動詞; F応...タグ F・応答詞 (対話の場合のみ); <笑>...<息>, <咳>も含む
 (?)...値なし (?); ?単...値1つ (?); ?複...値複数 (?); 左・右...タグの左項・右項 (例: (A 左項; 右項))

(?) 複数候補に関する補足説明:

- ◇ 同一範囲において、(W)(F)(M)(O)(R)(X)と(?)の複数候補との組合せは、(?)が外に来る場合と内に来る場合の両方を許容しているが、以下の条件で使い分けをしている。(O)と(W)を例に説明する。

- (O)・(W)かそれ以外かで迷う場合は、(?)が外、(O)・(W)が内:

(? (O 寺や), 寺屋) & テラヤ, (? 国保, 国庫) & (? (W コクコ; コクホ), コクコ)

- (O)・(W)であることは確かだが、その中で複数の候補が想定される場合、(O)・(W)が外(?)が内。

(O (? イムニダ, イムニカ)) & (O (? イムニダ, イムニカ))
国保 & (W (? コクコ, コフコ); コクホ), (? 国保, 国庫) & (W コフコ; (? コクホ, コクコ))

- ◇ 同一範囲、入れ子のいずれの条件においても、(?)の複数候補が外に来るケースを概ね許容()しているが、基本的に(?)の複数候補はあまり多用せず、できる範囲で候補を一意に同定する。

表 9: 【同一範囲】 の場合の使用条件・非許容の場合(x)の補足説明

外	内	転記例/許容条件
F	?複	音が曖昧な場合、(F(? うー, えー))のように複数候補も十分あり得るが、多用せずできる限り候補を一つに絞る。
F 応	VN	対話の応答詞のうち、基本形が(F うん), (F うーん), (F ふうん)の場合に限定。
F 応以外	VN	<VN>は必ず(F)と共に出現する。<VN>とそれ以外のタグとの関係は、(F)応と他のタグとの関係に依存する。
D2	?単	(D2)中の(?)は発音形の場合のみ。語が曖昧な場合は(D2)でなく(D)対処:(D(? ご))御本人
D2	W	助詞の「を」を「ウオ(-)」と発音した場合のみ:コレ(D2(Wウオ;オ))オ
?複	R	(? (R x x), 砂糖)のように(R)を複数候補の1つにのみ付与すると伏せ字の意味をなさないため原則認めない。
x AK 左	?複	(A(? エヌ, エム);(? N, M))とはせず、(?)を外に出し(? (A エヌ;N);(A エム;M))とする。
x A 右	?複	(A スリーエヌ;(? 3 N, NNN))など理論上あり得るが、(A)内ではどちらかの候補に倒し複数候補とはしない。
x ?単複	FV,	<FV>は母音が曖昧な場合に使用するタグのため、(?)対処しない。
x ?単複	笑	<笑>等のポーカル音は、そもそも存在が曖昧な場合には表記しない。
x B 右	?短	右項に自信がない場合には(B)対処としない。
x F 感	W	感情表出系感動詞は聞こえた通り記すため(W)対処しない。
x D	?複	(D(? のぼ, のご))などの複数候補は理論上あり得るが、ここまでせず(D(? のぼ))のように単一候補とする。
x ?複	L, 笑	同一範囲の(L, 笑)は(?)複数候補の外に出す:(? (L ソレ), (L コレ)) (L(? ソレ, コレ))
x D	A	(D(A シー;C)) (A シーディー;C D)など、理論的にはあり得るが、(D)内では(A)対処しない。
x D	K	(D(K か<FV>い; 回)) 回復など、理論的にはあり得るが、(D)内では(K)対処しない。
x D	W,B	(D)は発話された音をそのまま表記するため、(D 運) 運動 & (D (W オン;ウン)) ウンドーのように(W), (B)対処せず。

表 10: 【入れ子】 の場合の使用条件・非許容の場合(x)の補足説明

外	内	転記例/許容条件
A 左	D2	数字内の数字の言い直しに限定:(A 千(D2 八百) 九百四十;1 9 4 0)
W 左,F,D	H	母音の引き伸ばしが極めて長い場合に長音記号に追加した形でのみ使用:(W スー<H>;スル), (F アー<H>)
W 右	?F,D2,O,M,X	複数の語が融合して一つの(W)にまとめられた場合に限定:(W コキヤラ;ココ(D2 カラ)) マデ
F フ	W	(F えっとね)や(F あのですね)のように(F)が複数の短単位(例:あの/ですね)から構成される場合に限定。
F フ	?複	(F(? あっと, えっと)ですね)のように複数候補も十分あり得るが、多用せずできる限り候補を一つに絞る。
D2	?単	(D2)中の(?)は発音形の場合のみ許容。語が曖昧な場合は(D2)でなく(D)。
D2	W	助詞の「を」を「ウオ(-)」と発音した場合のみ:ナニ(D2 カ(Wウオ;オ)) カガ
M	M	括弧や句読点などのメタ記号の読み上げに限定:平仮名で(M こっかい (M 中括弧閉じる)では)で
?複	R	Rを複数候補の1つにのみ付与すると伏せ字の意味をなさない場合もあるため、原則認めない。
-	VN	VNは必ず(F)と共に出現するため、FVとの関係は、F応と他のタグとの関係に依存。
x A 左	?複	(A)内では(A(? エヌ, エム)三五;(? N, M) 3)とはせず(?)を外に出す:(? (A エヌ;N 3);(A エム;M 3))
x A 右	?単	理論上あり得るが、右項に自信がない場合には(A)表記としない。
x W 左	D	(W)左項では(W カイ(Dセ) セキユ;カイセキ)のように(D)は使用せず、(W カイセセキユ;カイセキ)とする。
x W 左	?複	(W ホ(? ト,ド);ホントー)のように左項の一部ではなく、(W(? ホト,ホド);ホントー)のように全体に(?)付与。
x B 左	D	(B シブ(D タ) タイジュータイ)のように(D)は使用せず(W)対処:(B (W シブタイ;シブタイ)ジュータイ)
x B 右	?単	右項に自信がない場合には(B)対処としない。
x W 左,F 感,D	Q	聞こえた通り表記する為、<Q>は使用せずに「っ」で表記:(D こっ) 話し言葉
x W 右	R	(W)右項の一部が(R)の対象であっても(R)は(W)全体に付与:(W x x;(R x x) x) (R (W x x;x x x))
x O	A,M	理論上あり得るが、タグ(O)を付与した時点で(A)や(M)は付与しない。

6 再朗読の転記について

本節では、再朗読（講演 ID が R00 から始まるファイル）の転記のうち、通常の講演の転記と異なる点について説明する。朗読（講演 ID が R01,R02,R03 から始まるファイル）については本節で述べる特別な措置は行っていないので注意されたい。

- ◇ 転記単位全体が、朗読原稿にない余分な発話の場合²⁰：<朗読間違い>タグを付与し、その単位の文字化はしない。

0035 00112.070-00114.322 L: それでは & ソレデワ (F あ) & (F ア) 0036 00115.436-00116.199 L: やり直します & ヤリナオシマス 0037 00117.940-00119.938 L: それはですね & ソレワデスネ	右のように対処	0035 00112.070-00114.322 L:<朗読間違い> 0036 00115.436-00116.199 L:<朗読間違い> 0037 00117.940-00119.938 L: それはですね & ソレワデスネ
--	---------	--

- ◇ 転記単位の一部が、朗読原稿にない余分な発話の場合²¹：文字化した上で (X) を付与

0414 01172.918-01176.613 L: (F ま) & (F マ) そういったものを & ソーイッタモノオ 色々 & イロイロ (F えー) & (F エー) (X 見てて) & (X ミテテ) 見てみたんですが & ミテミタンデスガ		転記単位の一部を言い直し
0090 00382.480-00384.785 L: (X (D (? ツ))) & (X (D (? ツ))) 全ての & スベテノ		朗読原稿にない (D)
0360 00788.358-00792.626 L: バリエーションが & バリエーションガ 少ない & スクナイ 為に & タメニ (X (F えー) & (X (F エー) 関西に比べて & カンサイニクラベテ 回答 (D げ)) & カイトー (D ゲ))		転記単位の一部を言い直し
0361 00793.262-00796.283 L:<朗読間違い> 0362 00796.763-00797.614 L:<朗読間違い> 0363 00799.184-00799.758 L:<朗読間違い> 0364 00800.032-00800.486 L:<朗読間違い> 0365 00801.043-00806.870 L: (F えー) & (F エー) 関西に比べて & カンサイニクラベテ 回答 (D げ) 語形の & カイトー (D ゲ) ゴケーノ 異なりというのは<FV> & コトナリトユーノワ<FV> 少なく & スクナク なってるんですけども & ナッテルンデスケレドモ		転記単位全体を言い直し 転記単位全体を言い直し 転記単位全体を言い直し 転記単位全体を言い直し <FV>に関しては (X) を付与せず

- ◇ 上記以外の朗読間違いについては、コメントを付与する。

0120 00289.467-00291.674 L: 仕事を & シゴトオ 持っている & モツテイル 方にのみ & カタニノミ %TYPE=0_ERR 朗読間違い 原稿では「方にのみ」ではなく「方のみに」		
0168 00392.749-00395.306 L: 引き続き & ヒキツズキ 百四ページの & ヒヤクヨンページノ %TYPE=0_ERR 朗読間違い 原稿では「百四」ではなく「百十四」。 グラフーから & グラフィチカラ		

²⁰ 読み間違えたため最初から読み直した場合など。

²¹ 基本形のレベルで朗読原稿と異なるものを対象とする。ただし、<FV>、値なしの (?) のみが異なる場合は本タグの付与対象外とする。

7 対話の転記について

本節では、対話（講演 ID が D から始まるファイル）の転記のうち、通常の講演の転記と異なる点について説明する。対話の転記テキストの例を図 11 に示す。

- ◇ 話者 ID：話者 2 名の ID は L と R。それぞれ音声ファイルのチャンネルに対応。
インタビューの場合、インタビュアーは L、インタビュイーは R に固定。
- ◇ 転記単位の認定方法：単位の認定は L・R 毎に行なう。その際の基準は通常の講演と同じ。
- ◇ 転記単位の順番：L・R に関係なく、開始時間の早い順に並べる。
- ◇ タグの付与範囲：タグは、L・R それぞれの転記単位内で整合的に付与する²²。
- ◇ 応答表現の扱い：
 - ・ 応答表現に対し、通常の講演では特にタグは付与していないが、対話では (F) を付与する²³。
 - ・ 対話に実際に出現し (F) の付与対象となった応答表現は以下の通り：

はい、ええ、うん、うーん、ううん、いえ、いいえ、いや

- ・ 応答表現の表記：基本形における応答表現の表記は、上記リストに示したものに固定する。その為「ハイ」のように母音や子音の引き延ばしが生じた場合には、<H>、<Q>を用いて表記する（発話 ID=0026 参照）。
- ・ 「うん」「うーん」「ふーん」「んー」「ん」の区別は概ね以下に行なう：

下降の音調：(F うん)・(F うーん)， 平坦：(F んー)・(F ん)， 上昇：(F ふーん)
(F ふーん)については、上記の音声的条件と機能的条件（「感嘆」）の両方を満たした場合に限定

- ◇ タグ<VN>： 応答表現「うん」「うーん」「ふーん」の発音において、鼻音が混ざる・母音カテゴリーが曖昧等の理由で、音（分節単位）の認定が困難なものについては、発音形を<VN>で表記する（発話 ID=0028,0034,0037）。
- ◇ 共話の扱い：以下の例に示すように、対話では 2 人で単語や文を共同して発話する場合がある（以下「共話」）。タグの付与や文字の表記は、L・R 毎に行なう為、例えば「共話例 C」の (D ざ) に見られるように、語の途中（この場合は「小遊三」の「ざ」）から発話をしたものについては、例え共話としては語の構成要素であっても、語の断片と見なし (D) が付与される為、扱いに注意されたい。

[共話例 A]
0212 00251.164-00252.195 R:
北野 & (L キタノ)
0213 00253.762-00254.620 L:
大さん & マサルサン

[共話例 B]
0483 00572.810-00573.463 R:
食堂に & ショクドーニ
0484 00573.358-00574.506 L:
に & ニ
行くんだ & イクン (笑ダ<笑>

[共話例 C]
0332 00429.353-00432.664 R:
小遊三さんに & コユーザサンニ
し (D (? て)) & シ (D (? テ))
(F はい) & (F ハイ)
0333 00430.389-00431.390 L:
(D ざ) さん & (D ザ) サン ... 語の断片と見なす
0334 00431.688-00432.236 L:
しましろう & シマシヨー
0335 00437.409-00437.949 R:
(F うーん) & (F <VN>)

²² 例えば、図 11 の発話 ID=0023~0025 に示すように、(R) が、相手話者の発話 (0024) を狭んで次の転記単位に跨がる場合であっても、発話 ID=0024 は (R) の付与範囲とは見なさない。

²³ 質問に対する肯定、依頼に対する承諾、話し手に対する相槌など、様々な機能があるが、機能に関わらず一律 (F) を付与する。

表 11: 対話の転記テキストの例

0021 00024.323-00026.764 L: じゃ 浮かぶ 人 取り敢えず 言ってもらっていいですか	& & & & &	ジャ<H> ウカブ ヒト トリアエズ イッテモラッテイーデスカ	
0022 00026.915-00027.264 R: (F ええ)	&	(F エー)	応答表現「ええ」にもタグ (F) 付与
0023 00028.979-00030.610 R: (F えーとー) (R × × × ×)	& &	(F エートー) (R × × × ×)	タグ (R) の範囲が、話者 L の「はい」を狭んで 0025 まで続く
0024 00030.813-00031.154 L: (F はい)	&	(F ハ<H>イ)	
0025 00030.892-00031.778 R: × × × × ×) さん	&	× × × × ×) サン	
0026 00031.540-00032.076 L: (F はい)	&	(F ハ<H>イ)	応答表現「はい」にもタグ (F) 付与
0027 00033.312-00033.837 R: でしょう	&	デシヨ	
0028 00034.008-00034.573 L: (F うーん)	&	(F <VN>)	「うーん」の発音形 <VN>
0029 00034.520-00036.565 R: それから 水野晴郎さんでしょう	& &	(W (? ホ);ソレ)(W カ;カラ) ミズノハルオサン (? デ) ショ	
0030 00036.053-00036.521 L: (F はい)	&	(F ハ<H>イ)	
0031 00037.377-00038.396 R: この 人	& &	コノ ヒ (? ト)	
誰だっけ	&	ダレダッケ	
0032 00039.707-00042.703 R: (F あー)	&	(F アー)	
あたし	&	アタシ	
結構	&	ケッコ	
名前	&	ナマエ	
知らない	&	シラナイ	
水沢アキさんか	&	ミズサワアキサンカ	
0033 00042.938-00043.613 L: (F はい) (F はい) (F はい) (F はい)	&	(F ハイ) (F ハイ) (F ハイ) (F ハイ)	1 つの応答表現毎に括り直す
0034 00043.636-00044.010 R: (F うーん)	&	(F <VN>)	
0035 00044.577-00045.449 R: それと	&	ソレト<H>	
0036 00046.146-00047.869 R: この 駄洒落	& &	コノ ダジャレ	
言う	&	ユ	
人	&	ヒト	
和田勉さんだ	&	ワダベンサンダ	
0037 00048.011-00048.753 L: 和田勉	&	ワダベン	
(F うん)	&	(F <VN>)	
0038 00048.550-00048.910 R: (F はい)	&	(F ハイ)	
0039 00049.402-00049.652 L:<笑>			