

音声認識読み記号および音声関連ソフトウェアに係わる用語の試行標準案

松浦 博¹ 西本 卓也² 金子 宏 磯谷 亮輔³
石川 泰⁴ 西村 雅史⁵ 伊藤 克亘⁶ 新田 恒雄⁷

1 東芝研究開発センター 2 東京大学大学院 情報理工学系研究科 3 NEC マルチメディア研究所
4 三菱電機情報技術総合研究所 5 日本 IBM 東京基礎研究所 6 産業技術総合研究所 7 豊橋技術科学大学大学院工学研究科
〒441-8580 愛知県豊橋市天伯町雲雀ヶ丘 1-1 E-mail: nitta@tutkie.tut.ac.jp⁷

あらまし： 本報告では、平成 14 年に発足した情報処理学会試行標準専門委員会下の WG4 小委員会（音声言語処理インタフェース）が策定した、(A) 音声認識読み記号（エンドユーザ向け）の試行標準案、および (B) 音声認識・合成に係わる用語（開発者向け）の試行標準検討案の内容を紹介する。

キーワード： 音声言語処理, 標準化, 用語, 読み記号, 音声認識, 音声合成

IPJSJ Trial Standard Concerning Spoken Language Interface

*Hiroshi MATSU'URA¹, Takuya NISHIMOTO², Hiroshi KANEKO, Ryosuke ISOTANI³,
Yasushi ISHIKAWA⁴, Masafumi NISHIMURA⁵, Katunobu ITOU⁶, and Tsuneo NITTA⁷*

1 Toshiba Corp., 2 Univ. of Tokyo, 3 NEC Corp.,

4 Mitsubishi Electric Corp., 5 IBM Japan Ltd., 6 AIST, 7 Toyohashi Univ. of Tech.

1-1 Hibariga-oka, Tempaku, Toyohashi, 441-8580 JAPAN E-mail: nitta@tutkie.tut.ac.jp⁷

Abstract: This paper describes the activity of IPJSJ Trial Standard WG4 that has the objectives of standardization for spoken language interface. In this report, we introduce two working drafts: (A) a phonemic system described in Kana used by end-users of ASR, (B) Terminology on ASR and TTS used by application developers.

Key words: Spoken Language Processing, Standardization, Terminology, Phonemic System, Speech Recognition, Speech Synthesis

1. はじめに

平成 14 年より情報処理学会に学会試行標準専門委員会（委員長石崎俊慶^{慶応義塾大学}教授）が発足し、この委員会の下で WG4 小委員会：音声言語処理インタフェースが発足し活動している[1]。学会試行標準は、情報処理学会の Web で公開され、国内外の規格化作業に役立てて頂くことを目指している。WG4 小委員会では、今回、(1) 音声認識読み記号の試行

標準案、(2) 音声認識・合成に係わる用語の試行標準検討案を策定した。本報告は、音声言語処理研究・開発に携わる多くの方々の意見を収集することを目的にまとめたものである。

2. 音声認識読み記号

音声認識の対象単語に用いられる読み記号は、表記の標準がないため、製品での対応は各社まちまちである。また、音

声認識のAPIにおいても、読みまでには踏み込まず、使用記号を規定するとどめているケースが見受けられる。読みの付与が使用者にゆだねられると混乱を招く恐れが生じる。こうした事情から、WG4 は読み記号を早期に情報処理学会試行標準とする作業を開始した。今後、業界団体等において規格化へ向けた議論が高まることを望むと同時に、音声認識用APIなどへも普及することを期待したい。

2.1 標準化の範囲

● 「音声認識読み記号」として表1に示す各記号を標準とする。網掛けされたアルファベット記号は発音の参考であり、標準化の対象ではない。標準に準拠した認識システムは、この範囲で入力された記号を少なくとも受理するか、あるいは受理できない旨のエラーを返すものとする。ただし、音声認識システムによっては受理した結果、必ずしも異なる記号が異なる発音にならないことがあっても、これは標準の範囲内として許容する。

● 音声合成を用いて認識結果をトークバックする(音声で知らせる)場合などを考慮し、JEIDA規格である「日本語テキスト音声合成用記号の規格」(JEIDA-62-2000)との連携を図る。こうした利用と、開発者がより詳細に読みを記述することを目的に表2の「音声認識詳細読み記号」を設ける。

● 音声認識読み記号の並び方に関する規則や制約を細かく規定すると、ルールが複雑になりすぎ、かえって利用することが困難になる。そのため、音声認識読み記号の並び方に関する規定は特に設けない。すなわち、表1に示す「音声認識読み記号」の並び方および表2に示す「音声認識詳細読み記号」の並び方は、各々の範囲内で自由に設定できるものとする。また、音声認識システムによっては、並び方に一部制限を設けることも考えられるが、標準の範囲内として許容する。

2.2 音声認識読み記号に関する説明

エンドユーザが使用できることを念頭においている。

<規定>

- 「ひらがな」をベースに、表1に示す「音声認識読み記号」を定める。
- 鼻濁音化の記号(例えば「^h」)については、エンドユーザは鼻音と鼻濁音を必ずしも区別できないと考えられるため、採用しない。
- 無声化を表す記号も同様の理由で採用しない。
- 表1のアルファベットは発音の参考である。ただし、現代かなづかい[2]に従うと読み記号は発音と一致しないことが

表1 音声認識読み記号

あ	い	う	え	お	や	ゆ	よ	わ	うい	うえ	うお
a	i	u	e	o	ya	yu	yo	wa	wi	we	wo
か	き	く	け	こ	きゃ	きゅ	きょ	くわ	くい	くえ	くお
ka	ki	ku	ke	ko	kye	kyu	kyo	kwa	kwi	kwe	kwo
さ	し	ず	せ	そ	しゃ	しゅ	しょ	つあ	つい	つえ	つお
sa	shi	su	se	so	sha	shu	sho	tso	tsoi	tsoe	tsoo
た	ち	つ	て	と	ちゃ	ちゅ	ちよ	つあ	つい	つえ	つお
ta	chi	tsu	te	to	chy	chu	cho	tso	tsoi	tsoe	tsoo
な	に	ぬ	ね	の	にゃ	にゅ	にょ				
na	ni	nu	ne	no	nye	nyu	nyo				
は	ひ	ふ	へ	ほ	ひゃ	ひゅ	ひょ	ふあ	ふい	ふえ	ふお
ha	hi	fu	he	ho	hya	hyu	hyo	fo	foi	foe	foo
ま	み	む	め	も	みゃ	みゅ	みょ				
ma	mi	mu	me	mo	mya	myu	myo				
ら	り	る	れ	ろ	りゃ	りゅ	りょ				
ra	ri	ru	re	ro	rya	ryu	ryo				
が	ぎ	ぐ	げ	ご	ぎゃ	ぎゅ	ぎょ	ぐあ	ぐい	ぐえ	ぐお
ga	gi	gu	ge	go	gye	gyu	gyo	gwa	gwi	gwe	gwo
ざ	じ	ず	ぜ	ぞ	じゃ	じゅ	じょ				
za	ji	zu	ze	zo	jya	ju	jo				
だ	ぢ	ぢう	ぢで	ぢど							
da	ji	ju	je	jo							
ば	び	ぶ	べ	ぼ	びゃ	びゅ	びょ				
ba	bi	bu	be	bo	bya	byu	byo				
ぱ	ぴ	ぷ	ぺ	ぽ	ぴゃ	ぴゅ	ぴょ				
pa	pi	pu	pe	po	pya	pyu	pyo				
てい	たい				てゅ						
tei	tai				tyu						
すい					すゅ						
su					tyu						
すい					すゃ	すゅ	すょ				
su					sha	shu	sho				
ん	っ	ー	・		ん	づ	を				
n	ts	-	.		n	zu	wo				

促音：っ(JIS 0x2443)

区切り：・(JIS 0x2126)

長音：ー(JIS 0x213c)

[注1] 表記例のコードは全角平仮名(JIS X 0208-1983「情報交換用漢字符号系」)を使用するが、実際に使用するコード系まで規定するものではない。

[注2] この表に定義されていない「ゐ」「ゑ」等の仮名記号は使用しない。

[注3] 長音記号として「- (マイナス JIS 0x215d)」は使用しない。

[注4] 網掛けされたアルファベット表記は発音の参考であり、音声認識読み記号の規定対象ではない。

ある。このような場合は、例えば現代かなづかいの本則で「とうきょう」「ていせい」などと記述される時、音声認識システム側で長音化処理を行い、それぞれ「とーきよー」「てーせい」等の発音がありうると解釈し処理するものとする。

● 一方、「小売」→「こうり」のように長音化しない場合がある。このような場合のため、長音化しないことを示す区切り記号として「・」(中黒あるいは中点)の使用をオプションとして認める。すなわち、「小売」は「こ・うり」と記述できる。なお、音声認識システムによっては「・」の代わりに「/」ほかの採用を宣言して用いても良い。

● 現代かなづかいの例外と位置付けられるもののうち、

「wa」と発音する助詞「は」、 「e」と発音する助詞「へ」は発音どおりにそれぞれ「わ」、「え」と記述することとする。したがって、この記述だけは現代かなづかいの規定と一致しないことになる。

● 他の現代かなづかいについてはこれに従うこととする。この規定は認識システムからみても問題は少ない。例えば、助詞「を」を一般に「o」、場合によって「wo」に近く発声したとしても、ユーザは「を」と記述すれば良い。

<補足説明>

● エンドユーザが利用できる簡易な形式を目指し、より馴染みやすい「ひらがな」を用いた。

● JEIDA 規格 日本語テキスト音声合成用記号の規格 (JEIDA-62-2000) に掲載された表 2-1 をベースとして作成している。また、国語審議会の「外来語の表記」[3] を参考に、一般的な外来語表記をカバーするようにした。

● 音声認識読み記号から発音への変換精度を上げる様々な工夫 (音声合成を用いた読み上げによる確認、発声登録、単語辞書との突合せなど) は、音声認識システムの製作者にゆだねた。入力誤り、ユーザごとの発音の揺れなどへの対処も、音声認識システムの製作者にゆだねた。

● 理解の手助けとして以下に例を示す。ここで () に記した発音表記は、認識システムが読み記号を解釈し内部記号に置き換えた一例である。

(凡例) 画面に表示する文字 音声認識読み記号の表記 (ある認識システムで解釈された発音表記の例)

東京 とうきょう (to-kyo,toukyou)

武蔵新城 むさししんじょう

(musashishinnjo-, musashishinnjou)

丸の内 まるの・うち (marunouchi)

丸の内 まるのうち (maruno-chi,marunouchi)

大阪阿倍野橋 おおさか・あべのばし (o-sakaabenobashi)

経営 けいえい (ke-e,keiei)

経営 け・いえい (keie-,keiei)

ユーザ ゆーざー (yu-za-) …… 外来語もひらがなで記述

秋田 あきた (ak(i)ta) …… 記号列から音声認識システム

が無声化を推測した例

鼻血 はなぢ (hanaji)

続く つづく (tsuzuku)

私は わたしわ (watashiwa)

表 2 音声認識詳細読み記号

ア	イ	ウ	エ	オ	ヤ	ユ	イェ	ヨ	リ	ウィ	ウエ	ウオ
カ	キ	ク	ケ	コ	キヤ	キユ	キエ	キヨ	クア	クイ	クエ	クオ
カ	キ	ク	ケ	コ	キヤ	キユ	キエ	キヨ	クア	クイ	クエ	クオ
サ	シ	ス	セ	ソ	シヤ	シユ	シエ	シヨ	スア	スイ	セエ	ソオ
サ	シ	ス	セ	ソ	シヤ	シユ	シエ	シヨ	スア	スイ	セエ	ソオ
タ	チ	ツ	テ	ト	チャ	チュ	チェ	チヨ	ツア	ツイ	ツエ	ツオ
タ	チ	ツ	テ	ト	チャ	チュ	チェ	チヨ	ツア	ツイ	ツエ	ツオ
ナ	ニ	ヌ	ネ	ノ	ニヤ	ニユ	ニエ	ニヨ	ヌア	ヌイ	ヌエ	ヌオ
ナ	ニ	ヌ	ネ	ノ	ニヤ	ニユ	ニエ	ニヨ	ヌア	ヌイ	ヌエ	ヌオ
ハ	ヒ	フ	ヘ	ホ	ヒヤ	ヒユ	ヒエ	ヒヨ	フア	フイ	フエ	フオ
ハ	ヒ	フ	ヘ	ホ	ヒヤ	ヒユ	ヒエ	ヒヨ	フア	フイ	フエ	フオ
マ	ミ	ム	メ	モ	ミヤ	ミユ	ミエ	ミヨ	ムア	ムイ	ムエ	ムオ
マ	ミ	ム	メ	モ	ミヤ	ミユ	ミエ	ミヨ	ムア	ムイ	ムエ	ムオ
ラ	リ	ル	レ	ロ	リヤ	リユ	リエ	リヨ	ルア	ルイ	ルエ	ルオ
ラ	リ	ル	レ	ロ	リヤ	リユ	リエ	リヨ	ルア	ルイ	ルエ	ルオ
ガ	ギ	グ	ゲ	ゴ	ギヤ	ギユ	ギエ	ギヨ	グア	グイ	グエ	グオ
ガ	ギ	グ	ゲ	ゴ	ギヤ	ギユ	ギエ	ギヨ	グア	グイ	グエ	グオ
ザ	ジ	ズ	ゼ	ゾ	ジヤ	ジユ	ジエ	ジヨ	ズア	ズイ	ズエ	ズオ
ザ	ジ	ズ	ゼ	ゾ	ジヤ	ジユ	ジエ	ジヨ	ズア	ズイ	ズエ	ズオ
ダ	ヂ	ヅ	ヅ	ヅ	ヂヤ	ヂユ	ヂエ	ヂヨ	ヅア	ヅイ	ヅエ	ヅオ
ダ	ヂ	ヅ	ヅ	ヅ	ヂヤ	ヂユ	ヂエ	ヂヨ	ヅア	ヅイ	ヅエ	ヅオ
バ	ビ	ブ	ベ	ボ	ビヤ	ビユ	ビエ	ビヨ	ブア	ブイ	ブエ	ブオ
バ	ビ	ブ	ベ	ボ	ビヤ	ビユ	ビエ	ビヨ	ブア	ブイ	ブエ	ブオ
パ	ピ	プ	ペ	ポ	ピヤ	ピユ	ピエ	ピヨ	プア	プイ	プエ	プオ
パ	ピ	プ	ペ	ポ	ピヤ	ピユ	ピエ	ピヨ	プア	プイ	プエ	プオ
タイ	トク				チャ	チュ	チェ	チヨ	クア	クイ	クエ	クオ
タイ	トク				チャ	チュ	チェ	チヨ	クア	クイ	クエ	クオ
ヴァ	ヴィ	ヴ	ヴェ	ヴォ	ヴィヤ	ヴィユ	ヴィエ	ヴィヨ	ヴァ	ヴィ	ヴェ	ヴォ
ヴァ	ヴィ	ヴ	ヴェ	ヴォ	ヴィヤ	ヴィユ	ヴィエ	ヴィヨ	ヴァ	ヴィ	ヴェ	ヴォ
カ	キ	ク	ケ	コ	キヤ	キユ	キエ	キヨ	クア	クイ	クエ	クオ
カ	キ	ク	ケ	コ	キヤ	キユ	キエ	キヨ	クア	クイ	クエ	クオ
スイ					フヤ	フユ	フエ	フヨ				
スイ					フヤ	フユ	フエ	フヨ				
ン	ッ	ー	()									
ン	ッ	ー	()									

鼻濁音: ° (JIS 0x212c)
 促音: ッ (JIS 0x2543)
 長音: ー (JIS 0x213c)
 無声化・母音脱落: ((JIS 0x214a) および) (JIS 0x214b)
 [注1] 表記例のコードは全角片仮名 (JIS X 0208-1983「情報交換用漢字符号系」) を使用するが、実際に使用するコード系まで規定するものではない。
 [注2] この表に定義されていない「ヂ」「ヅ」「ヲ」「キ」「エ」等の仮名記号は使用しない。
 [注3] 長音記号として「ー (マイナス JIS 0x215d)」は使用しない。
 [注4] 網掛けされたアルファベット表記は発音の参考であり、音声認識詳細読み記号の規定対象ではない。
 [注5] 半角片仮名で表記する場合、濁音では濁音記号「ㇰ」を使用する。
 [注6] 半濁音では半濁音記号「ㇱ」を使用する。

こんにちは こんにちは (konnnichiwa)
 学校へ がっこうえ (gaqko-e, gaqkoue)
 私を わたしを (watashio)
 私が わたしが (watashiga, watashinga)

2.3 音声認識詳細読み記号に関する説明

主に開発者の使用を念頭においている。

<規定>

- 「カタカナ」をベースに、表 2 に示す「音声認識詳細読み記号」を定める。
- 極力、発音どおりに記述すること。

- 鼻濁音の記号として「^h」を採用する。ただし、「カ^h」など表2にあげた形式での使用に限る。
- 長音化しないことを表す区切り記号は使用しない。
- 母音の無声化・脱落を表現するため括弧の記号を規定する。

<補足説明>

- 音声認識詳細読み記号は、音声合成読み記号との共通化を重視し、JEIDA規格 日本語テキスト音声合成用記号の規格 (JEIDA-62-2000) に掲載された表2-1を参考にした。
- 開発者が認識対象の発音をコントロールできることを重視した。ただし、各音声認識システムの性能向上のために、内部変換等の工夫をしてよい。
- 理解の手助けとして以下に例を示す。ここで()に記した発音表記は、認識システムが読み記号を解釈し内部記号に置き換えた一例である。

(凡例) 画面に表示する文字 読み記号の表記 (ある認識システムで解釈された発音表記の例)

東京 トーキョー (to-kyo-)
 丸の内 マルノウチ(marunouchi)
 武蔵新城 ムサシシンジョー (musashishinnjo-)
 大阪阿倍野橋 オーサカアベノバシ (o-sakaabenobashi)
 経営 ケーイー (ke-e-)
 経営 ケイエー (keie-)
 私を ワタシオ (watashio)
 学校へ ガッコウエ (gakko-e)
 こんにちは コンニチワ (konnichiwa)
 私が ワタシガ (watashiga)
 私が ワタシガ^h (watashinga)
 デッドボール デッドボール (deqdoboru)
 デッドボール デットボール (deqtoboru)
 全員 ゼンイン (zennin)
 全員 ゼーイン (ze-inn)
 秋田 ア(キ)タ (ak(i)ta)
 私です ワタシデ(ス) (watashides(u))
 ストップ(ス) トッ(ブ) (s(u)toqp(u))
 多い オーイ (o-i)
 氷 コーリ (ko-ri)
 姉さん ネーサン (ne-sann)
 鼻血 ハナジ (hanaji)
 続く ツズク (tsuzuku)

2.4 全般的補足

- JEIDAの音声合成用読み記号は長音化を長音記号で表記するが、音声認識の場合は、エンドユーザが確実に長音化させて発音するか否かを予め記述することは難しい。これは、他の音韻についても同様であり、読み記号は音声合成用と音声認識用とでは、自ずからその役割・性質が異なる。
- 現行の各社PCソフトウェアの仕様は、一般的な「読み」であり、「とーきょー」「とおきょお」等は推奨されていない。
 (例) ドラゴンスピーチ6 片仮名で読み (トウキョウワ)
 ViaVoiceV9 平仮名で読みと発音 (とうきょうわ)
 SmartVoiceV4.0i 平仮名で読み (とうきょうわ)
 LaLaVoice2001 平仮名で読み (とうきょうわ)
- 「音声認識読み記号」は、音声認識システムのサポートによって認識性能向上が期待できるが万能ではない。一方、「音声認識詳細読み記号」は、より正確に記述できる反面、間違っただけで記述された場合、致命的になる可能性が高いことに留意する必要がある。

3. 音声認識・合成に係わる用語

音声認識や音声合成の開発者向けドキュメントは、特定の動作環境やエンジンの実装に依存して書かれていることが多く、用語とそれらの定義・関連に揺らぎがある。WG4では、日本語の音声認識・音声合成エンジンに係わる用語を対象に、学会試行標準を検討している。標準化によって、システム開発者が既存のエンジンを使用する際に、動作環境やエンジンベンダーが変わっても違和感なく開発できることを目指す。

標準化のプロセスとして、まず日本国内で音声認識技術を製品化しているベンダー数社の現状調査を行った上で、実行可能性のある標準化案を小委員会内で検討している。主な方針は以下のとおりである。

- 既存の標準を尊重する (JEITA,W3C等)
- 関連分野のデファクトを尊重する (かな漢字変換等)
- 専門書や研究論文等の定義や慣用的用法を尊重する
- 現状調査 (既存の製品仕様やドキュメント) を考慮する
 今後メーリングリスト等を通じて、素案に関するアンケートを行う予定である。

3.1 標準化の対象と目標

開発者やユーザにとって用語の不統一が混乱を招きやすいと思われる事項で、かつ、統一化の合意が得られやすい提案に限定する。また、既存の記述の変更を求めるものではなく、

開発者間のコミュニケーションにおいて曖昧性が生じやすい場合に推奨する用語を提案する。

用語の選定に当たっては、用語間の階層関係（依存関係）を考慮するなどにより用語の位置づけを明確にすること、および用語がどのような状況で用いられどのように役立つべきかを明確にすることを目指す。

エンドユーザが製品を通じて目にする用語には、様々な不統一が存在している。しかし既に製品として出荷されていることもあり、エンドユーザが慣れ親しんだ用語を変更することには問題も多い。そこで現時点では、エンドユーザ向け用語の統一案は作成せず各社の判断にゆだねる。

一方、開発者向け用語については、以下のような目標を設定して用語標準化を提案する。

- 音声認識・音声合成エンジンの機能や性能を明確に記述でき、開発者がエンジンを選択する際の補助となること
 - システムに組み込まれた特定のエンジンを、同等の機能・性能を持つ他のエンジンと容易に交換できること
- これらの目標が達成されることで、エンジンの開発や改良、および応用システムの開発が促進されることを期待している。

3.2 音声に関する用語

「音声」の用法には、広く音響信号を示す場合と、言語情報を伴う音に意味を限定する場合がある。特に、音声認識や音声合成が応用システムの一部として用いられる場合には、これらの概念を明確にしておくことが重要である。

入力音声

音声認識システムの入力となる音を総称して「音響信号」と呼ぶ。「音響信号」の内容は「音声信号」と「雑音信号」に分離される。また、音響信号の時系列情報から「休止区間」を取り除いて「発話区間」を得ることを（発話区間の）「切り出し」と呼ぶ。

出力音声

音声合成システムの出力となる音を総称して「音響信号」と呼ぶ。「音響信号」には、「録音音声」および「合成音声」のほかに楽音や効果音などが含まれる。テキスト情報から「合成音声」を得る手法を「テキスト音声合成」と呼ぶ。「テキスト音声合成」において埋め込まれる休止を「ポーズ」と呼ぶ。

なお「プロンプト」は、アプリケーション中で音声を利用する方法のひとつを示した用語であり、上に定義した用語群とは区別して使用すべきである。

3.3 音声認識に関する用語

音声認識手法

「連続音声認識」では、用途を表す用語と音声認識エンジンの手法そのものが混同されがちである。そこで連続音声認識の手法として次の二つの分類を使用する。

- A: 統計的言語モデル(n-gram)を用いた連続音声認識
- B: 構文制約情報を用いた連続音声認識

いずれかの手法を限定したい場合は明示的に記述する。

例: Julius は統計的言語モデルを用いた「連続音声認識」を行っている。

例: Julian は構文制約情報を用いた「連続音声認識」を行っている。

なお、「ディクテーション」「ナビゲーション」「音声コマンド」などは音声認識手法ではなくアプリケーションを示す用語であり、手法を示す際には用いない。

音声認識性能

音声認識エンジンの性能の求め方や結果の表示方法は統一されていない。一般に性能とは認識率を指すが、処理速度が重視される場合もある。

アプリケーション開発者が目的に合った性能のエンジンを選ぶには、性能を記述する用語とその定義を標準化する必要がある。

アラインメント作業が必要な場合の音声認識性能は、「正解精度」(Accuracy) もしくは「誤り率」(Error Rate)を用いて示すことを推奨する。

$$\text{正解精度} = (\text{正解数} - \text{挿入数}) \times 100 / \text{入力総数}$$

$$\text{誤り率} = 100 - \text{正解精度}$$

1発話で1タスクが達成される場合も、挿入・脱落はないものとして、「正解精度」もしくは「誤り率」を用いて示す。

「正解精度」や「誤り率」の値を示す場合は単位を明記して「文字正解精度」「単語正解精度」「音韻正解精度」のように用いる。

音声認識応用システム全体の性能は、上記定義の他にも発話者の声質や利用環境、ユーザインタフェース設計などに大きく依存する。将来的には、エンジン性能の評価尺度をさまざまな観点から明らかにし、必要な情報の整備を促したい。

3.4 音声合成に関する用語

3.4.1 音声合成制御

テキスト音声合成の制御機能を表す用語は、様々なものが用いられている。また音声合成アプリケーションの開発者を対象とした標準やデファクトとして、JEIDA 日本語テキスト

音声合成用記号, W3C Aural CSS, W3C SSML, Microsoft speech API などが策定されている。

一方、個々のエンジンは既存の標準に従いつつ、独自の解釈や拡張が施されている場合も多い。機能ごとに設定される値は、物理量などの客観的な値ではなく相対値を用いる場合が多く、エンジンがAPIによって抽象化されていても実質的には互換性が妨げられている。

以下では標準案の候補となる用語とその定義を示す。

合成音声の速さ

「速さ」「速度」「話速」は合成音声の速さを示す値であり、速いほど大きい値として定義される。「速さ」の例として、1秒あたりのモーラ数や1分当たりの単語数などがある。

英語等では「話速」として1分当たりの単語数[Words per minute, WPM]を用いることが多く、Aural CSS (W3C [4])のspeech-rate 属性もそのように定義されている。マルチリンガルの音声システムを実現する上では、言語によらず統一的に話速を表現できることが望ましい。

声の高さ

「高さ」「ピッチ」「基本周波数」は、合成音声の声の高さをヘルツ[Hz]で示した値である。Aural CSS (W3C)のpitch 属性も周波数[Hz]として定義されている。ただし値の用法は実装依存となる。実際のアプリケーションでは、規定値からの相対指定(例: 75%, 150%など)の方が便利のため多く用いられている。

声の大きさ

「音量」「ボリューム」「大きさ」は声の大きさを表す値である。物理量としては音圧レベル[dBSPL]に対応する。しかし既存の実装においては、0を最小値とし、特定の値を最大値とした相対値で振幅を決めている場合が多い。

Aural CSS (W3C)のvolume 属性は0-100の値を取ると定義しているが、こうした値は環境に応じて制御されるものであり、同じvolume 値であっても騒音レベルによって異なる振幅となるなど、アプリケーションの観点から与えた定義をそのままエンジン制御に適用するのは適さない。

声質

異なる声質の合成音声を使い分け、GUIにおける色やフォントのように利用することが考えられる。「音声フォント」「キャラクタ」「ボイスフォント」「声質」「バリエーション」は声質の識別子の呼称である。一般に「音声フォント」には人名など任意の文字列を使用できる。

Aural CSS (W3C)はvoice-family 属性を提供しており、male, female, child など一般的な呼称と、特定のインスタンスを表す呼称を併用できる。

3.4.2 音声合成の手法および性能

アプリケーション開発者は、音声合成エンジンを選ぶ際に、エンジンが想定している音質レベルを知る必要がある。

音声符号化の分野ではCD相当(約44KHz, 16bit), FMラジオ相当(約22KHz, 16bit), 電話相当(約8KHz, 8bit)など既存メディアを目安にする方法があるが、合成音声の品質は帯域やビット数以外に、手法や実装の違いによっても異なる。また、特に視覚障害者向けアプリケーションでは速い話速が要求されることが多く、実用的な話速の最高値などもエンジン選びの指標となり得る。

一方、テキスト音声合成では、読み付与の精度も重要な評価尺度となる。ここでは2つの性能評価尺度を示す。

「読み付与精度」: 仮名漢字混じり文から読みを生成した際の単語単位の精度(計算方法は音声認識の単語認識精度に準ずる)

「正聴率」: 被験者に合成音声を取らせ書き起こしを行った際に、正しく単語を聞き取れた割合

$$\text{正聴率} = \frac{\text{正解音節数}}{\text{全音節数}}$$

評価方法の詳細はJEITAでの標準化活動を参考にする。

3. おわりに

試行標準は、今後、メーリングリスト等を通して関連分野の研究者・技術者の方々の意見を反映させた後、情報処理学会のWeb上で順次公開し、「役立つ試行標準」の提供を行っていきたい。

参考文献

- [1] 新田ほか: 音声言語情報処理に関する情報処理学会の試行標準策定活動, 情報処理学会研究報告, SLP-40-10, pp.57-60 (2002-02).
- [2] 小泉保「日本語の正書法」, 大修館書店
- [3] 国語審議会「外来語の表記」
<http://www.konan-wu.ac.jp/~kikuchi/kanji/gairai.html>
- [4] Aural style sheets,
<http://www.w3.org/TR/REC-CSS2/aural.html>