

# 道案内WOZシステムとの対話における言い淀み表現の分析

千田 恭子 ((財)電力中央研究所), 伊藤 克亘, 後藤 真孝 (電子技術総合研究所)

senda@denken.or.jp, kito@etl.go.jp, goto@etl.go.jp

## 1 はじめに

発話の際、次の言葉の選択に迷うことを言い淀みという。本論での言い淀みとは、とりあえずつなぎの言葉をしゃべったり(以後、つなぎ語と呼ぶ)、音節をひきのばして時間かせぎをすることである(以後、長音化と呼ぶ)。具体例を挙げると、以下の話者Bの「えーと」がつなぎ語、「それから」の語末のひきのばし部分が長音化である。(以降の例文は、断りがない限りは[1]の音声対話データからの引用。)

A) どんなものが売ってますか

B) えーと、ハンバーガーとか、それから、フライドポテト

他に、言葉につまんで黙ってしまったたり、言い直したりすることを、言い淀みに含む場合もあるが、本論では含まないこととする。

この言い淀みは話し言葉に非常に多く、例えばつなぎ語は1/3から2/3の発話に出現すると報告されている[2,3]。このように頻度高く出現するため、話し言葉を扱う音声認識や対話システムでは、言い淀みの現象は無視できないものである。そこで筆者らはこれまでに、言い淀み現象の的確な認識、活用を目的として、自動検出システムを構築し[4]、言い淀み表現のタグ付けを行った対話コーパスを用いて評価した[5]。本論では、同コーパスを用いて、言い淀みの言語的特徴を分析したので報告する。

以降の節では、先行研究の概要と課題についてまず説明し、次に対話コーパスの分析結果を説明し、最後にまとめを述べる。

## 2 つなぎ語と長音化の関係

### 2.1 先行研究の概要と残された課題

言い淀み現象を計算機で言語的に認識するためには、つなぎ語に関しては辞書に見出しとして登録する語形や接続、長音化に関しては長音化が起こる語の特徴等を明らかにする必要がある。しかし、先行研究では、その点に関して述べたものはほとんどなかった。

たとえばつなぎ語に関しては、社会言語学的もしくは語学教育的見地から述べたもの、また談話機能について分析したものがほとんどである[6-9]。これらの研究の一部は、コーパスから抽出したつなぎ語の一覧を

示しているので、辞書の見出しには、そこに挙がるつなぎ語をそのまま登録すればよいと思われるかもしれない。しかしつなぎ語には、語形が互いに類似した語のグループ(例:「えと」「えーと」「えーとー」)が幾つかあるため、核となる基本形と派生形に整理しないと、辞書登録の効率性・整合性が悪くなる。

またつなぎ語は、ドメインや聞き手・話し手の関係によって、異なったものが使用される[6]。そのため、あるコーパスから抽出しても、別のコーパスのつなぎ語に適合するとは限らず、解析に役立たない可能性がある。先に述べた基本形と派生形の整理は、ドメイン等によって異なるつなぎ語の変化/不変部分を見極め、語形変化の予測をある程度たてることにも役立つ。

言い淀みの際の長音化に関しては、これまで語末で(一部の文献では文節末で)起きる傾向が指摘されてきた[6,9,10]。だが語末ではなく、語の途中で長音化が起きているように見える場合もある。たとえば、先に類似語形のつなぎ語の例として、「えと」「えーと」を挙げたが、見方を変えれば、「えーと」は「えと」の語中が長音化したように見える。このように言い淀みの長音化が語中で起こることについては、従来研究では明確には説明されていない。

また、長音化は言い淀んでいない時でも起きるため、言い淀みとそれ以外(以後、非言い淀みと呼ぶ)の長音化が起きる語の違いを明らかにして、両者を識別できるようにすることが、言い淀み現象の的確な認識には必要である。たとえば次の二つの文は、前者は語調を柔らげるもしくは強める表現として、後者は近年の若者に多い一種の口ぐせとして、長音化が起きている例である。

C) わかんない

D) とりあえず手作りの一、ものを作れるように一、道具や一、小物が一、揃っています

筆者らが長音の継続時間について分析したところ、言い淀みの長音は非言い淀みの長音より長い傾向にあることがわかった[5]。但しその分布の重なりは大きいため、両者の識別には継続時間のような物理的なパラメータだけでなく、言語的な情報も必要である。例文Cのような一種の強調の長音については、主に形容詞、副詞で起きると指摘する研究はあるが[11]、他の品詞で起きる場合もあり、どの品詞でどれぐらいの割合で起きるか、実データに基づくさらに詳しい分析が必要である。例文Dのような口ぐせの長音に関しては、長音

化が起きる語の特徴についての研究はほとんどない。

## 2.2 本論の着眼点と分析方法

前節では、言い淀みの認識に重要でありながら、先行研究でほとんど扱われてこなかった課題として、以下の三つをあげた。

1. 類似語形が多いつなぎ語の整理
2. 語末でなく語中で起きる長音化の説明
3. 言い淀みと非言い淀みの長音化の語の性質の違い

これらは、一見別々の課題に見えるかもしれない。しかし本論では、これらを相互に密接に関連した課題としてとらえる。なぜならば、つなぎ語は複数の要素に分解でき、長音化はその各要素末で起き得るために、少しずつ異なる語形のつなぎ語が存在したり(課題1)、長音化がつなぎ語中で起きるように見えたりしているが(課題2)、そのことは、言い淀みの長音化が起こる語の性質(課題3)とつなぎ語の各構成要素の性質との類似性から、検証できると考えるからである。

つなぎ語が複数の要素に分解でき、長音化が各要素末で起きている例を示そう。たとえば、2節で類似語形のつなぎ語例とした挙げた「えと」「えーと」等で、本論の分析コーパスに出現したものは、以下の三つのグループに整理できた。(括弧内の「えとね」はなかったが、異なるコーパスならば存在し得ると考えて挙げたである)。

- 1) え／えー
- 2) えと／えーと／えーとー
- 3) (えとね)／えーとね／えーとねー

これを見ると、類似した語形の主な差異は、1)から3)のグループの違いである語末の接辞的要素「と」「ね」の有無と、各グループ内での各要素末の長音の有無にあることに気づく。つまりつなぎ語には、1)の頭の「え」のような語の基本形と、そこに接辞「と」「ね」がつく派生形があり(グループ1)~3)の違いに対応)、さらにその各構成要素が長音化することで別の派生形があるように見える。(グループ内の各語の違いに対応)。

## 3 コーパスの分析

前節で述べた考えを裏づけるため、本節では言い淀みの長音化が起こる語の性質と、つなぎ語の各構成要素の性質との類似性を分析する。

## 3.1 コーパスの概要

本研究では、渋谷の道案内をタスクとして、システムになりすました人と被験者が対話を行う、Wizard of Oz方式によって収録した音声対話コーパス[1]を分析対象として用いた。このコーパスには、書き起こした対話文のつなぎ語、長音化個所にタグが付けられている[5]。

表1が、付与したタグの定義である。長音化である

表1 言い淀みのタグと個数

現象の種類		タグ	個数	
長音化	言い淀み	明らか	@	942
		微妙	:	224
	非言い淀み	明らか	-	260
		微妙	-	322
つなぎ語		[ ]	771	

かの判断は微妙であるため、明らかに長音化と判断できる場合と、微妙と判断した場合とでタグを変えている。なお、つなぎ語のタグの[ ]は、つなぎ語をはさむように付与し、その他は、長音化で音節が引きのばされた音の直後に付与している。以下は、タグを付与した対話文の例で、これには三種類のタグ(@, -, [ ])がふられている。

E [と、]ファーストフード系の、お店が@、いいんですけど\_

ちなみに、このタグ付け作業では、2名の作業者が同一個所を相互にチェックしている。

## 3.2 分析結果

### 3.2.1 長音化の語の性質

まず、言い淀みと非言い淀みの長音化の性質の違いについて分析する。

[語末] 言い淀み、非言い淀みの各長音化が語のどの個所で起こるか調べた。@の言い淀みの長音化が語末で起きる割合は全体の95%と高く、語末が長音化するという通説は、ほぼ裏づけられた。また残り5%の語中の@について、その内訳を調べた結果、その約半分(48%)が固有名詞であった。本コーパスは渋谷の道案内をタスクとするため、「ダルゴロニューゼ」「中華めん菜」といった渋谷地域内の地名や店名が多いが、被験者には不慣れな語であるため、語中で言い淀み長音化することが多かったと考えられる。

これに対し、~がつく非言い淀みの長音化は語中が11%と、言い淀みの倍以上の割合であった。これは、非

言い淀みの長音化のうち、語調を強めたり柔らげる表現は、2節の例文Cのように語中で起きやすいためである。

**[文節末]** 語末の長音化について、文節中での位置を調べた。@の長音化は文節末の割合は59%で、~の同割合は98%であった。@の非文節末は41%と、通説から予測されるよりその値は大きい。内訳を調べたところ、@の非文節末の96%は、名詞(代名詞、固有名詞含む)または名詞性接尾辞、名詞性の強い副詞の語末だった。また逆に名詞、名詞性接尾辞、名詞性副詞の語末の@を調べてみたところ、うち82%は非文節末だった。この理由については、名詞は後接する助詞の選択が文全体の構成に関わり難いため、助詞の前で言い淀むのではないかと考えられこのような名詞の例外を除けば、その他の語は通説通りほとんど文節末に起きている。

**[品詞]** 言い淀みと非言い淀みの長音化が起こる語の品詞について調べた。つなぎ語以外で長音化(@と~)が生じた主な語について、品詞別にまとめたものが図1であるが、品詞的にも、言い淀みと非言い淀みの長音化との語種は異なる傾向にある。(但し、長音化の個所は語末語中に関わらずカウントしている)。図1に見るよ

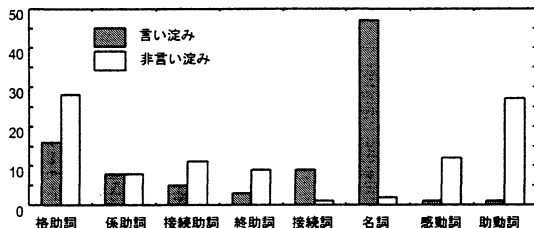


図1 長音化する語の品詞 (%)

うに、名詞(代名詞、固有名詞、名詞接尾辞、名詞性副詞を含む)、接続詞で起きる長音化は、ほとんどが言い淀みの長音化(@)であるのに対して、助動詞や感動詞は、非言い淀みの長音化(~)がほとんどである。助動詞や感動詞に非言い淀みが多いのは、2節の例Cで示したように、語調を強めたり柔らげて感情を表現する長音が多いためであると考えられる。また格助詞と係助詞は、言い淀み/非言い淀みのどちらに対しても比較的高い割合で現れ、文節末の語末でおこるという点でも似ているため、音声以外の言語的な情報だけで二つを識別することは難しい可能性がある。

### 3.2.2 つなぎ語の語構成

本コーパスに現われるつなぎ語の語形を、正規表現を用いてまとめたものが、表2である。本コーパスのつなぎ語は、表2の上から順に、副詞系列、連体詞系、人の呻吟の際に発せられる言語化されていない音(以後、呻吟音と呼ぶ)に基づく呻吟音系に分けられる。下から二つは、それぞれ頻度数1のみで、つなぎ語として聞き覚えのない音でもあり、他の音が欠落して理解不能に見える可能性もあるため、本論では特には扱わない。

表2 つなぎ語の語構成

つなぎ語の種類	つなぎ語の形	出現個数
副詞系	まあ?	3
連体詞系	あの@?	13
呻吟音系	あ([@:]ん?と(@ね@?)?)?)?	136
	う@ん?と(@ね@?)?)?	48
	ん([@:]?つ?と(@ね?)?)?	72
	え([@:]?つ?と(@ね@?)?)?)?	456
	つ?と@?	49
	えの@	1
	に@	1

?は、直前の記号があってもなくてもよいことを表わす

**[副詞系、連体詞系]** 表2のつなぎ語を構成する各要素について、3.2.1節で述べた言い淀みの長音化の性質と比較しながら分析する。まず表中の上の二つの「副詞系」「連体詞系」について、表の語形では語末が長音化することを示している。(副詞系の語末には@がついていないが、語末の「あ」が長音の一種の可能性もある)。この性質は勿論、3.2.1節で確認した、言い淀みの長音は語末で起こりやすいという性質と一致している。

また、つなぎ語以外で長音化が起こった語の品詞について調べた結果、連体詞の言い淀みの長音化(@:)は8例、非言い淀み(~)は1例、副詞の言い淀みの長音化(@:)は20例、非言い淀み(~)は3例のみだった。連体詞、副詞の長音化に関しては、言い淀みの長音化である可能性が高いことを示している。

よって、つなぎ語の構成要素のうち、副詞系、連体詞系の語は、言い淀みの語の性質と一致する傾向を示したといえる。

**[呻吟音系]** 次に、呻吟音系のつなぎ語について、長音化によって分割される各要素の性質が、つなぎ語以外の長音化の語の性質と類似性が高いかどうか分析する。

呻吟音系のつなぎ語では、語頭にある「あ」「う」「ん」「え」の要素は各グループの類似語形に共通して含まれるものであり、これが核の基本形を成している。この語は、その語形やニュアンスからいって、人の呻吟の

際に発せられる言語化されていない音に基づくものと本論ではとらえる。このような語は、つなぎ語以外の語では、感動詞にしか表われないであろう。感動詞の長音化は、言い淀みでは6例、非言い淀みでは57例あるが、言い淀みはそのうち50%が「あ」「ん」といった呻吟音的な語であるのに対し、非言い淀みは91%が「こんにちは」といった呼びかけの感動詞か「はい」といった応答詞だった。よって、語種の点からみると、つなぎ語の基本形である呻吟音(の語)末の長音化は、言い淀みの長音化の性質と一致する。(呻吟音系の語中の「ん」や促音「っ」は、上記呻吟音の長音化部分の一種と本論ではみなす。)

呻吟音の語に後接する助詞トについては、つなぎ語以外で長音化する例は、4例しかないため、このトの性質について論じることは難しい。但し、このトはつなぎ語末(つまり文節末)の時は長音化が起きるが、終助詞ネが後接する時は長音化しないという性質を示した。これは、言い淀みの長音化の性質である、文節末の語が長音化しやすいという性質と一致する。また後述する終助詞ネは、この助詞トを介しないと、上記の呻吟音には後接しない。よってこのトは、非言語的な呻吟音に付くことで、呻吟音を言語的なレベルに引き上げる役割を担っているといえる。

助詞トの次に接続するネは、そのニュアンスや、副詞系・連体詞系のつなぎ語にも後節し得ることから、終助詞ネと解釈できる。つなぎ語以外での終助詞ネの長音化は、言い淀み(㊦)が9例、非言い淀みが2例あった。よってつなぎ語のネで起きる長音化は、言い淀みの長音化である可能性が高いと言える。

以上、説明してきたように、長音化可能個所を手がかりに分解可能な、つなぎ語の各構成要素は、言い淀みの長音化が起こる他の語の性質となんらかの共通点があり、二つの語の類似性は高いと考えられる。これは、2.2節で述べた説を支持するものである。

なおこの節では、つなぎ語の構成要素となる助詞を幾つか挙げたが、そのうち終助詞ネは、場面や話し手と聞き手の関係に左右されやすい表現である。つなぎ語につくこのような表現として他には、「です」「だ」の助動詞、「な」「さ」といった終助詞があげられる。(但し、ドメインや聞き手・話し手の関係が固定されている今回のデータにはない)。たとえば、公式な場面や目上の相手に対しては、「えーとですね」のように「です」が用いられたり、くだけた場面や発話者と同等以下の地位の相手に対しては、「えーとだな」「あのさ」というように、「だ」「さ」が用いられると考えられる。よって、今回コーパスから抽出したつなぎ語は、ドメインや話し手と聞き手の関係の異なる対話では、助動詞・終助詞部分が異なる形で現われるだろう。

## 4 まとめ

前節までの分析で、長音化する語と、つなぎ語の構成要素の語との共通点を検証し、その類似性の高さを確認した。また、この分析を通して、つなぎ語は長音化可能個所をてがかりに整理して基本形と派生形に分けられることを示せた。これによって、つなぎ語の検出処理のためには、計算機用辞書にどのような語を基本形として登録するべきか、またその派生形としてはどのような形があり得るか、それはドメインや話し手・聞き手の関係によってどう変わり得るかを、ある程度明らかにできた。

また同分析を通し、言い淀みの長音化が起きる語の特徴について、通説通り語末や(一部の例外を除けば)文節末に起きやすいこと、非言い淀みの長音とは品詞的に異なる傾向があることを明らかにした。この言語的特徴は、継続時間のような物理的パラメータで検出しにくい、言い淀みの長音化表現を検出する際に役立つ可能性がある。

今後は、音声認識システムに本論で整理した情報を提供し、そのシステムのつなぎ語の認識率の向上にどれほど貢献するかを実験して、この研究の評価としたと考えている。

## 参考文献

- [1] Katunobu Itou et al. A Japanese spontaneous speech collected using automatically inferencing wizard of oz system. *J. Acoust. Soc. Jpn.(E)*, Vol. 20, No. 3, pp. 207-214, 1999.
- [2] 上條俊一, 他. 音声対話データの分析と発話理解への応用. 情処研報, 94-SLP-3, pp. 31-36, 1994.
- [3] 村上他. 自由発話音声における音響的な特徴の検討. 信学論, Vol. J 78-D-II, No. 12, pp. 1741-1749, 1995.
- [4] Masataka Goto et al. A real-time filled pause detection system for spontaneous speech recognition. In *Proc. of Eurospeech '99*, pp. 227-230, 1999.
- [5] 後藤真孝, 他. 有声休止個所のリアルタイム検出システムの評価. 日本音響学会講演論文集 春季 3-8-8, 2000.
- [6] 塩沢孝子. 日本語の hesitation に関する一考察. 社会言語学シリーズ 2 ことばの諸相, pp. 151-166. 1979.
- [7] 田吹昌俊. 日本語コーパスによる繋ぎ語とその配列・衝突についての分析. 九州工業大学研究報告(人文・社会科学), No. 4, pp. 17-24, 1999.
- [8] 李麗燕. 日本語母語話者の会話における「情報伝達行動の持続」. 世界の日本語教育, Vol. 7, pp. 61-75, 1997.
- [9] 田窪他. 応答詞・感動詞の談話的機能. 音声文法研究会(編), 文法と音声, pp. 257-279. くろしお出版, 1997.
- [10] 小出慶一. 言いよどみ. 講座日本語の表現 3 話しことばの表現, pp. 81-88. 1983.
- [11] 相沢佳子. 強度強調の長音. 音声学会会報, Vol. 167, No. 4, pp. 5-8, 1981.