

ロボットとの対話によるマルチメディアブログ創作システム

奥村 明俊 池田 崇博 西沢 俊広 安藤 真一 安達 史博
NEC メディア情報研究所

1. はじめに

ブログは、近年、消費者生成メディアの代表的な存在として利用者数が急増し、個人利用だけでなく、企業など組織内のコミュニケーションツールとしても活用されている。それに伴いブログ市場は拡大し、関連市場も含めると約 1377 億円に達し、事業者から様々なサービスが提供されている[1]。このように、ブログの目的やサービス内容が多様化する中、最近では、テキストだけでなく、映像や音声、動画などのマルチメディアコンテンツを取り込んだブログが急増している。ブログのマルチメディア化には、1) 録音した音声や録画した動画に見出しとなるキーワードをつけるもの、2) テキストのブログに対して関連する映像や音声、イラストなどのコンテンツを盛り込むものがある。マルチメディアブログは、臨場感豊かにメッセージを伝えることができ、注目されているが、WEB やインターネット、コンテンツ検索などに慣れていない、いわゆる情報リテラシの低いユーザにとって、魅力的なものを作成することは容易ではない。いわゆるデジタルデバイドの課題が顕著となる。また、ブログ作成経験のあるユーザであっても、1) の場合、携帯電話などを用いて簡単にメッセージを入力することはできるが、他のユーザから検索・閲覧されるための見出しとなるキーワードの付与に手間がかかる。2) の場合、メッセージに対して関連コンテンツを検索したり、それらを編集・コーディネートする手間がかかるなど、いずれの場合も作成コストの課題がある。さらに、せっかく作成したブログであっても他者からのコメントやトラックバックがないと持続しにくい。最近、自動的にブログにコメントやトラックバックするエージェント的なペットが登場しているが、ユーザの感情や気分に合わせてコメントが期待されている。

これらの課題を解決するために、ユーザがロボットに話しかけ対話することで、ロボットがユーザの発話内容を解析し、関連するコンテンツを探してコメントとともにマルチメディアブログを作成するシステムを提案する。本システムは、ノート PC と同程度のハードをベースに開発された対話型ロボット PaPeRo(パペロ) [2][3] 上に、旅行日記作成のためのプロトタイプとして構築される。以下に、まずマルチメディアブログ創作手法について述べ、次に全体システム構成を説明し、最後に動作例と評価結果を示す。

2. マルチメディアブログ創作

2.1. 創作処理の概要

本システムは、ロボットとの対話によって入力されたビデオメッセージから、見出し用キーワードやメッセージの内容に関連するコンテンツ(イラスト、音楽など)と、ロボットからのコメントを含むマルチメディアブログを作成する。ロボットとの対話によって、情報リテラシの低いユーザでも簡単にマルチメディアブログを作成することができる。

提案システムの処理の流れを図 1 に示す。システムは、入力されたビデオメッセージを蓄積し、発話音声抽出し、映像発話音声認識によって発話テキストに変換して見出しとなるキーワードを抽出する。次に、自然言語文検索により発話テキストと関連するマルチメディアコンテンツを、予め指定された WEB やディレクトリから検索する。これらの機能によってユーザのマルチメディアブログ作成コストを軽減する。さらに、発話テキストの表現やモダリティからユーザの心的状態を推定し、それに合わせたコメントを生成する。そして、入力ビデオメッセージ、見出し用キーワード、検索されたマルチメディアコンテンツおよびコメントをコーディネートしてブログを創出する。

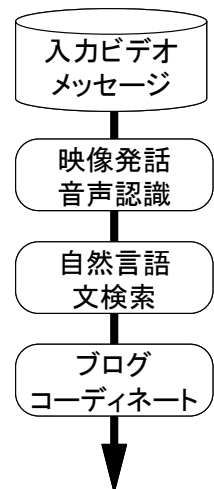


図 1: 処理の流れ

2.2. 映像発話音声認識

ブログ映像の中では、さまざまな単語や表現を含む話し言葉を扱う大語彙連続音声認識が必要である。大語彙連続音声認識は、一般に多くの処理能力やメモリ等のリソースを必要とするが、対話や動作など様々な処理を行う小型ロボットにおいて話し言葉を認識するためには、省メモリでコンパクトに動作する大語彙連続音声認識が必要である。そこで、今回、PDA 程度の端末でも動作する、コンパクトでかつスケーラブルな大語彙連続音声認識フレームワークを映像発話音声認識に用いる[4]。本フレームワークは、以下の手法により

コンパクトかつ高速な処理を実現している。

(1) 音響モデル・距離計算

音響モデルの使用メモリ量と距離計算の計算量を削減するために以下の3つの手法を用いる。

- 記述長最小基準に基づく効率的な分布数削減[5]
- 対角共分散行列の共有化による分布の簡易化
- 分布の木構造化に基づく出力確率の高速計算[6]

(2) 単語辞書・言語モデル

言語モデルとしては、単語の n 個組の連鎖確率である n -gram モデルを用いている。利用可能なメモリ量や処理能力に応じて単語 2-gram, クラス 2-gram, 単語 3-gram を組み合わせている。クラスは品詞をベースに、対象分野に応じて意味的なクラスや自動クラスターリングにより細分化して用いている。

(3) 最適単語列探索

最適単語列探索は距離計算結果と言語モデルを用いて、入力音声に同期して可能性の高い候補のみに絞り込みながら、辞書中の単語との照合を行う。辞書は先頭の共通部分を束ねることで木構造化して圧縮した表現で保持し、動的に展開して用いる。一定間隔ごとのワークメモリのガベージコレクションや、言語モデルの計算結果の再利用などにより、メモリ量、計算量を削減している。

今回、旅行日記作成をタスクとして、旅行会話向け日英音声翻訳システム[7]の日本語認識エンジンをベースに約 5 万語の大語彙連続音声認識を構築し、発話音声テキストに変換する。

2.3. 自然言語文検索

自然言語文検索は、ユーザによる発話の音声認識結果を検索要求文としてマルチメディアコンテンツの検索を行う。一般に、マルチメディアコンテンツには、十分なインデックスが与えられていることは少ない。そのため、コンテンツが含まれている WEB ページやドキュメントのテキストを手がかりとして検索する必要がある。本システムでは、検索要求文中の自立語をキーワードとしてテキストの検索を行い、その上位結果の近傍に含まれるマルチメディアコンテンツをブログに使用する。

また、検索結果に関するユーザとのインタラクションを無しにブログを創出するためには、自然言語文検索に高い適合率が要求される。本システムでは、Okapi BM25 式[8]による検索モデルをベースとして、適合率を高めるために以下の3つの拡張を行い、上位候補からより適切なコンテンツを抽出する[9]。

(1) 係り受け関係にある単語ペアの利用

例えば、同じ「魚」という単語でも、レストランで魚を食べた場合と、水族館で魚を見た場合とでは表示すべきイラストは異なる。そこで、予め単語間の係り受け

関係を求めておき、係り受け関係にある2つの単語を組にした単語ペアを重みの計算に利用する。具体的には、検索要求文中に含まれる単語ペアがテキスト中に含まれる数に応じてテキストの重みを加算する。

(2) 否定表現と肯定表現の区別

「楽しかった」という表現と「楽しなかった」という表現のように、否定表現と肯定表現とでは表す意味が正反対になるため、両者を区別して扱わなければならない。そこで、各単語が肯定的に使用されているか否定的に使用されているかを、その単語に付随する付属語によって判別し、検索においては、肯定の単語と否定の単語を異なる単語として扱う。

(3) 同義語の同一視

ユーザの発話内に含まれる単語から確実にマルチメディアコンテンツを見つけるために、同義語辞書を用意し、同義語を事前に特定の単語に統一した上で、テキストの重みを計算する。本システムでは、事前に数百語からなる同義語辞書を用意して、利用している。

2.4. ブログコーディネート

ブログコーディネートは、発話テキストからユーザの心的状態を推定してユーザに共感するコメントを PaPeRo のひとこととして生成する。そして、そのコメントを、入力ビデオメッセージ、見出し用キーワード、検索されたマルチメディアコンテンツとともに創出する。

発話中の表現は、人間の感情を推定する重要な手がかりである[10]。そのような感情を直接表現する語彙に注目して心的状態を推定する。ここではモダリティも含めて表層的な発話表現から推定可能な心的状態を定義し、心的状態と発話表現の対応データベースを構築する。このデータベースを2.3.の自然言語文検索によって、係り受け関係、否定・肯定表現、同義語辞書を用いて検索する。今回、ユーザの心的状態として、喜び、怒り、悲しみ、願望などデフォルト状態も含めて10種類を定義し、それぞれに100種類ほどの発話表現を対応づけた。そして、心的状態ごとに生成すべき PaPeRo のコメントを用意し、検索結果に合わせてコメントを生成する。コメントは、「やったね」や「残念」といった言語的なものに、PaPeRo の動作を示す GIF アニメをリンクしている。生成されたコメントは、作成されるブログ中に PaPeRo の動作とともに表示される。

3. 全体システム構成

PaPeRo は、図2に示すモジュールから構成される。全体制御モジュールが、入出力デバイス制御や認識・検知・合成など各種モジュールを対話動作シナリオに基づいて実行する。入力デバイス制御は、マイク、カメラ、圧力センサーを制御して、音声、映像、タッチを入力情報として全体制御プラットフォームに伝達する。全

体制御モジュールは、入力された情報を対話動作シナリオに基づいて、音声認識や顔認識などのモジュールに伝達してモジュールの機能を実行する。モジュールの実行結果は、全体制御モジュールに伝達され、対話動作シナリオに基づいて、さらに各種モジュールや出力デバイス制御部に伝達される。出力デバイス制御部は、各モジュールの出力結果を対象となるデバイスから出力する。

今回、入力デバイスとして、ビデオメッセージ入力用にマイクとカメラ、メッセージの開始と終了のスイッチとしてタッチセンサを利用する。また、出力用デバイスとして、動作のためのモータ、PaPeRo の口や耳の動きを示すライト、声を出すためのスピーカを利用する。これらのデバイスを用いてユーザと対話するシナリオを対話・動作シナリオに追加した。さらにマルチメディアブログを WEB に掲載するための WEB 掲載モジュールを実装した。2 節で述べたマルチメディアブログ創作の機能は、PaPeRo のひとつのモジュールとして実装され、その出力結果が出力デバイス制御を介して WEB 上に掲載される。

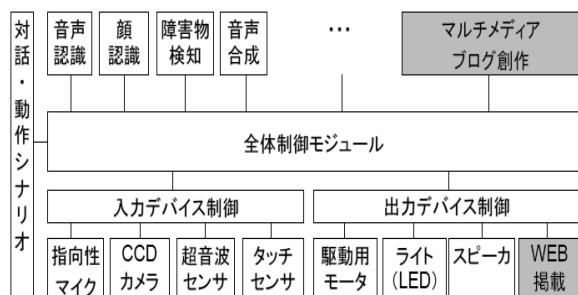


図 2：全体システム構成図

4. システムの動作例

旅行日記作成をタスクとし、イラスト 1500 と音楽 550 を含む約二千ページを検索対象としてプロトタイプシステムを構築した。その動作例として、ユーザと PaPeRo の対話例を表 1、対話の様子を図 3、作成されたブログ画面を図 4に示す。図 4A には、PaPeRo の内蔵カメラとスピーカから撮影されたビデオメッセージが貼り付けられ、Play ボタンで再生可能である。図 4B に、ユーザの発話から見出し用キーワードとして抽出されたヨセミテ、ラスベガス、ルーレットが掲載されている。図 4C は、ヨセミテ、ラスベガス、ルーレットのイラストが検索結果として表示されている。図 4D には、検索された音楽が貼り付けられ Play ボタンで再生可能である。図 4E に、PaPeRo のひとつこととして「そうなのいいね」というコメントが GIF アニメとともに掲載されている。これは、「行きたい」という表現から願望という心的状態が推定された結果である。

表 1: ユーザと PaPeRo の対話例

ユーザ：ブログ作って
PaPeRo：はいブログ作っちゃうんですね。
頭のボタンを押すとビデオの撮影を始めます。
ユーザ：（ボタンを押す）
PaPeRo：終わったらもう一回押してくださいね。
では、用意スタート。
ユーザ：ヨセミテに行きました。そのあとラスベガスでルーレットをしました。また行きたいです。
（ボタンを押す）
PaPeRo：ばっちり撮影できました。これからブログページを作成します。しばらくお待ちください。その間踊らせていただきます。（PaPeRo 踊る）
PaPeRo：ブログページができました。パソコンや携帯電話でご覧ください。自信作ですよ。



図 3：ユーザと PaPeRo の対話の様子

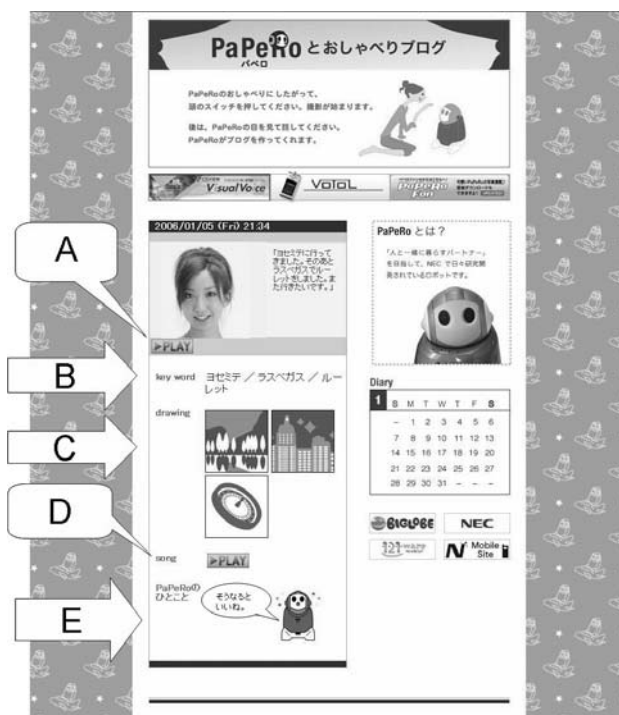


図 4：作成されたブログ画面

5. システム評価

5.1. 評価方法

システム評価のためには、映像発話音声認識、自然言語文検索など個別機能の性能評価と、全体としての出来栄や使い勝手などトータル性能評価が必要である。それぞれの評価項目が多岐にわたり、今後どのような観点に注目して評価を進めるかも課題である。個別機能の性能評価は、旅行会話向け日英音声翻訳システム[7]や音声入力によるテキスト検索システム[9]など他のアプリケーションで行ってきた。今回は、ユーザにブログを作成してもらい、出来上がったブログや使い勝手などや期待などのヒアリング調査を行う。その結果を、音声認識精度とブログ作成経験の有無から分析し今後の課題を抽出する。

5.2. 結果と考察

9名のユーザが、旅行や旅行以外の日記文を一人平均20文程度入力してブログを作成し、その結果に対してヒアリング調査を行った。高い音声認識精度(単語正解率85%以上95%以下)の結果作成されたブログに対する意見と、低い認識精度(単語正解率50%以上85%未満)で作成されたブログに対する意見に分け、さらにユーザのブログ作成経験の有無により、主な意見を表2にまとめた。考察結果を以下に記す。

(1) ロボットのブログ作成とコンテンツ拡充への期待

認識精度が高い場合、手間をかけずに簡単にブログが作成でき、ブログ作成経験の有無に関わらずユーザの満足度は高い。現在、十分な検索インデックスが付与されたマルチメディアコンテンツは少ないが、提案システムによりインデックスが付与されたマルチメディアブログが普及することになる。その結果、マルチメディアブログと検索可能となるマルチメディアコンテンツが相乗的に普及・拡充していくことが期待される。

(2) 音声認識誤りに対する頑健性と寛容性

認識精度が低い場合でも、自立語が正しければ、まずまずの検索結果を得る。また、PaPeRoがブログを作るという本システムは、誤りも面白さとなりユーザに許容されることもある。ただし、必ずしも面白い誤りとなるわけではないので、見出しキーワードの修正や全体レイアウトの編集インタフェースは必要である。

(3) PaPeRoとの対話への期待

認識精度に関係なく、システムとしての面白さや完成度を高めるために、PaPeRoとのやりとり、対話バリエーションや機能の充実が求められている。PaPeRoは、ユーザの発話を解析して内部的に保持しており、以前のブログ内容を反映したコメントやトラックバックが可能である。また、対話シナリオを追加することで、ビデオ

撮影中にユーザの発話を解析して、PaPeRoからユーザに質問してブログとすることも可能である。

表2:ヒアリング調査結果

		ブログ作成経験		
		有り(4名)	無し(5名)	共通
音声認識精度	高	・検索できるマルチメディアコンテンツの拡充と普及に期待	・PaPeRoとならブログを作りたい	・話すだけでブログが作れるのは楽 ・PaPeRoのひとことが面白い
	低	・見出しキーワードや検索結果、全体レイアウトを編集したい	・検索結果はまずまず。 ・誤認識の予想外のイラストも面白い	・誤りもPaPeRoなら許せる ・誤認識テキストが表示されているとビデオ再生したくなる
	共通	・PaPeRoのトラックバックが欲しい ・対話バリエーションの充実	・PaPeRoに曲の選択理由を聞きたい	・PaPeRoのコメントに以前のブログの解析結果を反映

6. おわりに

誰もが簡単に魅力的で楽しいマルチメディアブログを作成できるよう、ロボットとの対話によるマルチメディアブログ創作システムを提案し、PaPeRo上にプロトタイプを構築して評価・分析を行った。その結果、音声認識誤りも含めてユーザから好印象を得るシステムの基盤となることが判った。今後、対話機能の充実やインタフェースの改良を進め、ブログを創る楽しみをより多くの人に与えるシステムとして改良を進める予定である。

参考文献

- [1] 総務省: ブログ・SNSの現状分析および将来予測, 2005年5月
- [2] 藤田善弘, “パーソナルロボットPaPeRoの開発,” 計測と制御, Vol.42, No.6 (2003.6)
- [3] <http://www.incx.nec.co.jp/robot/>
- [4] 磯谷亮輔 他: “話し言葉認識に向けた基本技術と応用”, 情処研報, 2005-NL-169, pp.109-116, 2005.
- [5] 篠田他: 「音声認識のためのMDL基準を用いた効果的なガウス数削減」, 信学技報, SP2001-83, 2001-10.
- [6] Watanabe et al.: “High Speed Speech Recognition Using Tree-Structured Probability Density Function”, ICASSP-95, pp.556-559, 1995.
- [7] 山端潔他: “PDAで動作する旅行会話向け日英双方向音声翻訳システム”, 情処研報, 2002-NL-150-9, 2002.
- [8] S. E. Robertson et al.: Okapi at TREC-3, TREC-3, pp.109-126, 1995.
- [9] 池田崇博他: “自由発話音声入力による携帯電話向けテキスト検索システム”, 言語処理学会第10回年次大会, pp.109-112, 2004.
- [10] 中村明 編者: “感情表現辞典”, 東京堂出版