

動向情報の要約と可視化とその展開

- MuST (動向情報の要約と可視化に関するワークショップ) 活動報告 -

加藤 恒昭¹ 松下 光範² 神門 典子³
東京大学¹ 関西大学² 国立情報学研究所³

概要

動向情報の要約と可視化に関するワークショップ MuST について、その活動内容を報告する。このワークショップは、動向情報のマルチモーダル要約について、協調的かつ競争的に研究を進めていくことを目的としている。課題の重要性と枠組み、その中で行われた興味深い研究の幾つか、評価課題の設計について説明し、共通の研究資源として作成されたデータセットの概要を述べる。

Multi-modal Summarization for Trend Information: Its Past and Future Development

- An Activity Report on the MuST Workshop -

Tsuneaki Kato¹ Mitsunori Matsushita² Noriko Kando³
The University of Tokyo¹ Kansai University² National Institute of Informatics³

Abstract

The MuST (Multimodal Summarization for Trend Information) workshop was designed to encourage cooperative and competitive studies on summarization and visualization for trend information. In this paper, the framework of its research subject and the reasons why it is worth addressing are explained, some representative themes studied in the workshop are introduced, and the design of the data sets constructed and the evaluation tasks conducted are explained.

1 はじめに

動向とは「今年に入って原油とガソリンの価格はどう動いているのだろうか」「06年からゲーム機業界はどんな感じになったのか」「去年の台風はひどかったのか」等で示される利用者の関心に対する最初の回答となるものであり、その関心に応える情報の全体像を概観できるものである。

筆者らはこのような動向情報に注目して、「MuST: 動向情報の要約と可視化に関するワークショップ」(以下、MuST)を提案し、活動を続けている。このワークショップは、動向情報の要約と可視化に関する技術(動向情報のマルチモーダル要約)について、協調的かつ競争的に研究を進めていくことを目的としている。共通の素材を用いて、緩い意味で共通の課題に取り組むことによる議論と研究の活性化、研究コミュニティの形成、ツールやコーパス類の蓄積を目的として開始したが、その後、多くの参加者が取り組んでいた問題を評価課題として具体化し、評価型ワークショップの側面を合わせて持たせるようにした。

本稿では MuST の活動について報告する。まず、MuST が取り組んでいる課題について説明し、どのようにワークショップが進められたかを述べる。続いて、MuST の枠組みの中で行われた研究、MuST を通じて構築された研究資源についてそれらの概要を説明する。最後にこれまでの意義をまとめ、今後の展望を述べる。

2 課題の位置づけ

動向情報と、要約や可視化によるその生成は、以下ののような興味深い性質を持つ。

- 一定期間にわたる情報を総合的にまとめあげた要約となっていることが必要であるが、そのためには、複数の情報源に分散しかつ重複の多い情報の組織化が必要となる。
- 言語的な情報だけでなく、生産台数やシェアのように時系列データである統計情報への言及を含む場合や、台風や地震のように地理的な情報を含む場合が多い。
- 情報を整理する観点も、時系列、地理的空間に加えて、シェアにおける企業や支持率における政党等の様々な可能性があり、多次的で、視覚的な情報提示の利点が大い。
- 新聞記事や blog のような様々なテキスト情報に加えて、統計量に関する詳細な数値情報が白書等に存在し、それらを横断した情報のジャンルや形式にとられない情報収集が必要となる。
- 統計量への言及等の事実情報は動向情報の一部に過ぎず、それに加えて、その解釈や原因の推測や波及効果の予測等が重要な意味を持つ。

特に情報源と情報提示の両方に言語情報と非言語情報を協調させ活用することの必要性あるいは利点が大いことから、動向情報の要約と可視化のためには、それらを横断

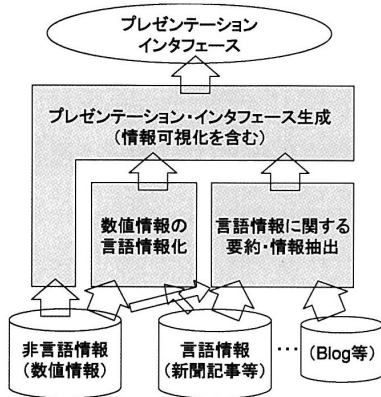


図1 動向情報のマルチモーダル要約

的に扱う技術が重要となる。その研究は、要約や情報抽出の技術をはじめとした言語情報（テキスト）を扱う自然言語処理技術と非言語情報（数値データ）を扱う情報可視化技術との融合と発展を促し、その結果、従来から研究されているような精密な意味表現を前提とするマルチメディアプレゼンテーション生成ではなく、実世界の情報に対応できるマルチモーダル要約^{*1}という新しい技術の確立に繋がることが期待される [6]。

また、動向情報が利用者の蓋然的な関心に対する最初の回答であることから、提示された動向情報は対話的で探索的な情報アクセスの入り口となる。提示された動向情報が対話的でその後の情報アクセスのインタフェースとなっていれば、それとのインタラクションを通じて、関心の詳細化や具体化が進み、利用者が本当に必要とするより詳細な情報への到達を支援することができる。この点で、マルチモーダル要約の研究は対話的探索的な情報アクセスとも深く関係する [9]。^{*2}

動向情報のマルチモーダル要約を行うシステムのイメージを図1に示す。動向情報が提示されるものでありインタフェースの役割も持つことから「プレゼンテーション・インタフェース」という表現を用いている。テキスト要約の技術（言語→言語）、テキストから統計量等に関する情報を抽出する技術（言語→非言語）、統計量の値の変化の様子等、数値情報の列を言語的に説明する文章を生成する技術（非言語→言語）というようなメディア変換を含む広義の要約が要素技術となる。更に、例えば、数値情報を参照したテキスト要約等、これらが相互に関連することで高度化する。これらにより得られた情報要素あるいは情報源から直接得られた数値情報が、情報可視化を含むプレゼンテ

ーション・インタフェース生成によって、有機的にまとめあげられることで、動向情報のプレゼンテーションあるいはインタフェースが構成される。なお、要約や情報抽出については、利用者から与えられた特定の関心に基づくものばかりではなく、テキストマイニングのような上昇的なものを考えることもできる。それは、与えられた情報源が総体として示している動向を捉えることであり、そこでもその提示やその後の対話が重要な役割を持つ。

以上のように、動向情報は対話的探索的な情報アクセスにおいて重要な位置づけを持っており、そのマルチモーダル要約は様々な技術を融合させ発展させる契機となる興味深い研究テーマとなっている。

3 特徴と経緯

MuSTは、動向情報のマルチモーダル要約に関する技術について、協調的かつ競争的に研究を進めていこうというワークショップである。取り組む対象が新しい研究分野であったため、提案当初の2年間（第一、第二サイクル）は、ワークショップの目的を議論と研究の活性化、研究コミュニティの形成、ツールやコーパス類の蓄積とし、この目的のために、共通の研究資源を用いて、緩い意味で共通の課題に取り組むことを行っていった [2]。

MuSTの特徴は共通の研究資源を用いる点にあり、そこに求心力と参加への動機付けを求めた。動向情報のマルチモーダル要約は、様々な研究分野にまたがった様々な要素技術が必要とし、トータルなシステムの構築は必ずしも容易ではない。そこで、MuST データセットと名付けた共通の研究資源を提供することで、研究者各人が関心ある要素技術に取り組むことを可能とした。この MuST データセットは研究対象となる素材であるだけでなく、それへの処理の中間結果を注釈づけたコーパスを中心としている。研究資源を共通化することで、今まで異なる分野に属すると考えられていた研究者達の議論が可能になる。もちろん、同じ分野の研究者は、この共通の素材を使って一定の客観的評価が行えることになる。これらのことを通じて、研究の加速と活性化を図っていた。オーガナイザは、MuST データセットの構築と保守、各種ミーティングの企画に加えて、HP や ML による情報発信、データやツールを参加者で共有するための呼びかけ等を行うことで、活性化を支援した。

第一、第二サイクルで、複数の参加組織が類似した研究課題に取り組んでいたこと、システム構成や要素技術についての意識が共有され整理されたことを受けて、第三サイクルでは、共通的な問題を具体的な課題として提案し、評価型ワークショップの側面を持たせた。これらの課題を評価課題、それまでの緩い意味で共通しているが、参加者独自の着眼と設定による研究を自由課題と呼んでいる。このサイクルでは、オーガナイザは評価課題の提案・具体化と実施・評価を行った。それに加えて、自由課題と評価課題に関する研究で利用できる資源を開発し提供することで、

^{*1} マルチメディア要約は映像情報の抜粋を指すことが多いので、それと区別するためにマルチモーダル要約という用語を使う。

^{*2} 初期には、動向情報を得ること自体を情報アクセスの目的と考えており、情報アクセスの入り口としての動向情報という見方は希薄であった [1, 3]。

表1 MuSTの活動

2004.11	NL・NLC 合同研究会にて提案 [1]
～ 2006.3	第一サイクル (15 組織が参加)
2005.11	ラウンドテーブルミーティング (closed)
2005.12	NTCIR-5 workshop meeting にて概要報告 [4]
2006.2	テキスト情報の要約と提示に関する自然言語処理シンポジウム (信学会 NLC 研究会)
2006.3	成果進捗報告会 (closed, 梗概を HP にて公開)
2006.3	自然言語処理学会第 12 回年次大会ワークショップ「言語処理と情報可視化の接点」
～ 2007.3	第二サイクル (18 組織が参加)
2006.10	知能と情報 (日本知能情報ファジィ学会論文誌) 特集論文「テキストの可視化と要約」
2006.11	ラウンドテーブルミーティング (closed)
2007.3	成果進捗報告会 (closed, 予稿論文を HP にて公開)
2007.6	NTCIR-6 workshop meeting にて概要報告 [8], 一部参加者の研究報告 [11]
～ 2008.12	第三サイクル (13 組織が参加)
2008.1	ラウンドテーブルミーティング (closed)
2008.3	成果進捗報告会 (closed, 予稿論文を HP にて公開)
2008.7	評価課題 formal run 実施 (5 組織が参加)
2008.12	NTCIR-7 workshop meeting にて概要報告 [10], 参加者全員の研究報告 [12]

ワークショップとしての充実を図った。また、数値情報の可視化とそれに言語情報を注釈づけたマルチモーダルプレゼンテーションが容易に作成できる可視化プラットフォームを構築し、共通の研究基盤とすることを進めた。ただし、これについては実際に参加者に利用していただくまでには至らなかった。

MuST は、国立情報学研究所が主催する NTCIR workshop^{*3}の一部であり、第一、第二サイクルは NTCIR-5, NTCIR-6 のパイロットタスクとして運営され、第三サイクルは NTCIR-7 の正式なタスクとなっている。

参加組織数を含めた MuST の活動経緯を表 1 に示す。筆者らが企画に参加し、関係したテーマを持つ会議や特集論文誌を合わせてまとめている。これらに加えて、人工知能学会全国大会の近未来チャレンジにおいて、「情報編纂の基盤技術」を提案し [5]、より広い視野で言語情報処理と非言語情報処理との融合や対話的探索的情報アクセスの支援に取り組んでいるが、それは表に含めていない。

4 取り組まれた研究課題

4.1 自由課題

第一、第二サイクルでの参加者の研究、第三サイクルにおいて自由課題として行われた研究は、2 節の図 1 で示した要素技術のいずれかにほぼ分類できる。紙面の制約で本節ではその一部だけを紹介する。より詳しい情報は MuST の HP^{*4}や文献 [11, 12] を参照されたい。なお、同じく紙面の制約で参考文献の提示が不十分であることをお詫言する。

言語情報の要約・情報抽出として、新聞記事テキストから特定の統計量に関する情報を抽出することが行われている。これは第一、第二サイクルの MuST においても

も盛んに取り組まれたテーマである。統計量の時間変化を示すグラフのプロット点となるような統計量名、時刻、統計量の値の三つ組を抽出する。基本的な手法では、統計量名、時間表現、数値表現を固有表現抽出技術を用いて同定し、出現順序、その回りの表層表現やキーワード、依存構造等を用いてこれらの要素をお互いに関連づける。時刻に関する情報が頻繁に省略されること、複数の統計量に関する情報がひとつの文で表現されることが問題となる。また、統計量名はその表現のバラエティも多いため、ある数値表現がどの統計量の値となっているかを正しく関連づけることにも課題がある。

動向情報ということでは特徴的であるのは、特定の統計量についてできるだけ多くの時点の情報を抽出したいという目的である。このため「前年同月に比べて 5 ドル上昇した」のような間接的に情報を与える比較表現を利用したり、特定の情報の推移を表現する文書を横断した文間関係に着目したり、ある統計量の値は一定の範囲内に収まり大きく変化することはない等のヒューリスティクスを利用したりすることが行われている。また、動向情報の要約で重要となる多様な情報のまとめあげということで、ある新聞を参照して作成された情報抽出のパターンが、異なる新聞でどの程度有効で、共通的に利用できるかの調査や、新聞記事を対象に作成した統計情報に関する情報抽出システムを blog テキストに適用することがどの程度可能であるか等の検討がなされている。

数値情報の言語情報化として、時系列情報からそれを説明する文章を生成することが複数の参加組織によって取り組まれている。数値情報を時間区間毎や部分形状や全体と部分に分割して、それぞれについて適切な言語表現を割り当て、適当な文章になるようにそれらを合成するのが基本的な流れとなる。適切な分割の手法や、言語表現とその割り当てに分野知識をどの程度用いるかが課題となる。

プレゼンテーション・インタフェース生成でも、前述の

^{*3} <http://research.nii.ac.jp/ntcir/index-ja.html>

^{*4} <http://must.c.u-tokyo.ac.jp>

情報抽出結果を折れ線グラフ等で表示するような比較的単純な可視化を含めて、様々な研究がなされている。ひとつは時空間にまたがる情報の可視化で、地震情報を対象に、これを地図上に配置して発生場所に関連づけて表現する、発生の日時に対応した時系列グラフとして表現するという複数の可視化を組み合わせて、それを様々な観点から対話的に眺めるための操作を、LDAP で用いられている概念を可視化に応用することでモデル化している研究がある。

時系列情報とそれに関する言語情報をどのように関連づけるかについては、時系列情報に関連した一連の情報にアクセスするための汎用のインタフェースとして、ズームング等様々な操作が可能な時系列グラフ上にその時点について言及している記事を対応付け、グラフに対する直接操作による概観の理解とより個別の関心に応える詳細情報へのアクセスを縫い目なく繋げるインタフェースが提案されている。また、統計量の変化を表現した折れ線グラフにその変化の要因等を説明する言語的な注釈を与えることも提案されており、アンケート調査によって、値の変化の大きい時点や値が最大や最小となった時点に注釈が欲しいと感じられることが明らかにされている。更に、変化に関する定性的な言語情報を可視化するために様々な変化を表現するパターン化した矢印をグラフに貼付けするという手法も提案されている。

また、言語表現と視覚表現の融合として、要約テキストの生成を、その変化を表現するグラフと対応づけて行うことが試みられている。グラフの粒度が細かく値の詳細な変化を表している場合と粒度の大きいより大局的な変化を表現している場合とで、重要文抽出を行う文書集合の大きさを変化させている。

この他に動向マイニングとも呼べる一群の研究テーマがある。そこでは、利用者から与えられた特定の関心に基づくのではなく、文書集合全体から上昇的に、統計量に関するキーワードの出現傾向等を基に動向を見つけ出したり、統計量の変化とその変化に伴って特徴的に出現したキーワードとの関係、統計量どうしの関係等を理解しようとする。例えば、統計量名を実世界の出来事や現象を特徴づけるものと捉え、それらを文書中から自動抽出し、一種の近傍共起を計算することで、それらの間の因果関係を獲得しようとしている。更にそれをネットワーク表現によって可視化することで複雑で国際的な問題の構造の概観を試みている。

その他に、文書集合全体にどのような情報が含まれているかを概観しようとする試みもある。例えば、与えられた文書集合中のすべての統計量名、日付、値の三つ組を収集し、統計量名と値の類似性に着目してこれらを分類することで、文書集合中に存在するすべての統計量の変化に関する情報を一覧できるようにするという試みがある。統計量だけでなく、文書集合中に示される数値情報間の関係、数値情報と固有表現との関係、更には、数値情報と特定の意味カテゴリの語と関係抽出し、グラフ化することで、そ

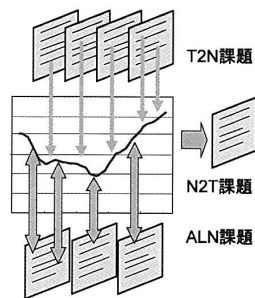


図2 評価課題の位置づけ

こで述べられている様々な関係を概観することも試みられている。

4.2 評価課題

第三サイクルで実施され、評価型ワークショップの側面を構成する評価課題は、動向情報のマルチモーダル要約の基本的な要素技術について客観的定量的評価を行うものとして立案された。言語情報と数値情報の相互変換や対応付けを可能とするための以下の3つで、図2に示すような関係を持つ。

T2N 課題（言語情報に関する情報抽出）文書集合中で言及されている時系列統計情報を一定の形式で抽出する。その文書集合中で、ある統計量のどの時点のどの値が話題や関心となっているかが明らかになる。

N2T 課題（数値情報の言語情報化）数値の列として表現されている時系列統計情報の変化や概要を表現する文章を生成する。得られる文章は数値情報の言語による要約である。

ALN 課題（プレゼンテーション・インタフェース生成の要素、言語情報と数値情報のアラインメント）文書集合中でなされている統計量やその変化への言及を取り出し、数値の列として表現されている時系列統計情報の対応する部分に関連づける。対象とする言及は具体的な数値を述べるものだけでなく、「ピーク」や「なだらかな上昇」等のより定性的な傾向表現を含める。これはメディア・アラインメントと呼べるもので、異なる情報源による異なるメディアの協同的利用を可能にする。

これらのうち、T2N 課題のみが実施された。それ以外については、課題設定の具体化や評価指標の検討が不充分であったこと、共通の評価が行える形でこれらの研究に関心を持つ参加組織が少数であったことから、第三サイクルでの実施は見送られた。

T2N 課題では、ある文書集合から与えられた統計量の値を抽出することになるが、後述の MuST データセットのトピックとなるようなもののみを選び、そこに含まれる複数の統計量を課題とした。同じトピックでは同じセットの記事が抽出の対象として用いられた。統計量は名前と単位

表現（助数詞）で指定している。一部の統計量は複数の名前と単位表現を与えている。8トピック25統計量が課題となり、抽出対象となった記事数は8から20である。これらは1998年から2001年の毎日新聞からとられた。例えば、「ガソリン」トピックでは、レギュラーガソリンの全国平均店頭価格（円）とドバイ原油価格（ドル）が課題となり、12記事から抽出が試みられた。

評価は、一般的な情報抽出と同様に精度と再現率、F値で行った。なお、精度に関連するものとして、この課題が与えられた統計量の値の抽出と、それと時間情報の組み合わせという2つの段階からなることから、誤りをそのどちらに起因するかで分類して分析する、再現率に関連するものとして、抽出された情報のうち、正しいものだけを繋いで実際にグラフを描いてみる等、課題の特徴を考慮した独自の分析が可能と思われるが、実際の評価で用いるには至らなかった。参加したシステムは、自由課題で概説した言語情報の要約・情報抽出の技術を用いている。それらの具体的な処理及び評価については、文献[10, 12]に詳しい。

5 研究資源

前節からも分かるように、動向情報のマルチモーダル要約に関する研究は、様々な広がりを持つ。それらにワークショップとしての求心力を持たせるのが、共通の研究資源である。それにより、参加者各人が関心ある要素技術に取り組むことが可能となり、かつ、それが議論の基盤となり、参加者どうしの議論が深まることが期待される。同時にこれらは、評価課題におけるテストコレクションや参考情報としての役割も持つ。

MuSTで構築された研究資源のうち、MuSTデータセットと変化情報コーパスは、言語情報に関する要約・情報抽出について、研究対象となる素材と処理の中間結果を提供している。可視化プラットフォームは、それを用いて様々なプレゼンテーション・インタフェース生成部が構築できるようなプラットフォームである。統計情報の描画やそれらの拡大・縮小処理など、共通的な機能をライブラリとして、システム構築の負担を軽減する、見た目や操作性に関する統一性を図り、異なるシステムどうしの比較を可能とするというふたつの効果を期待している。以下では、MuSTデータセットと変化情報コーパスについて、その概要を説明する。

5.1 MuSTデータセット

2節の図1における言語情報に関する要約・情報抽出の一般的な流れを図3に示す。ユーザの質問や関心が与えられると、まず、それに関連する文書が収集される。その後、要約の処理であれば、重要文の抽出、抽出された重要文の解析、冗長性除去や整合性確保を考慮した書き換えと文章生成が続く。情報抽出であれば、固有表現抽出と時間表現解析を含む参照表現解析とが行われ、関係抽出が行われる。これらによって得られた情報要素がプレゼンテーション・

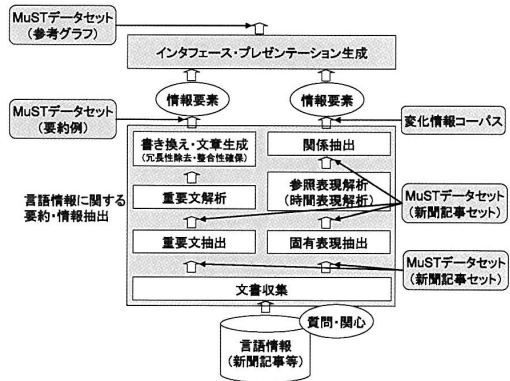


図3 言語情報に関する要約・情報抽出

インタフェース生成に渡される。

MuSTデータセットはユーザの質問や関心に対応するトピックに対して、これらの処理の結果や参考物を提供している。トピックは「ガソリン価格」「自動車生産」等の社会的経済的な分野やアクティビティ、「地震」「台風」等の事件のカテゴリ、「ソニー」「小淵内閣」等の組織、主体等が選ばれている。これらのトピックの動向を構成する情報の基本要素を、それに関する統計量の値についての言及もしくはある出来事の報告であると考えている。例えば、「パソコン業界」の動向は、パソコンの国内出荷台数、国内出荷額、それらのメーカー別シェアという統計量によって説明されるし、「台風」の動向は、個々の台風の発生や上陸に関する出来事の報告とそれによる被害についての統計量等から構成される。このようなデータに加えて、新聞記事等からは状況やその変化についての記述や評価等も得ることができ、それらが総合的にまとめられて動向情報となる。

MuSTデータセットは、言語情報として毎日新聞1998年から2001年の4年分を用い、27のトピックとそれに関連する統計量や出来事のクラス（各トピックについて3つ前後、合計で90の統計量もしくは出来事）について、作成あるいは収集された以下の情報からなる。収集された記事は702記事である。

- それに関する情報を含んだ新聞記事のリスト
- それら記事への注釈を行ったコーパス
- 統計量の変化や出来事に関する200字程度の要約
- それらの統計量のその時期における変化や出来事に関する表やグラフ^{*5}

動向情報のマルチモーダル要約において、新聞記事は情報の収集によって得られた言語情報に相当し、文章による要約は個々の統計量に関する要約の生成の出力として得られる要約結果に相当するものである。グラフや表は収集された情報としてもプレゼンテーションのサンプルとしても捉

^{*5} ただし、要約とグラフ等は第一サイクルでのみ作成したもので、58の統計量もしくは出来事に関してとなっている。

えることができる。

各記事への注釈は、要約における重要文抽出の結果、情報抽出の固有表現抽出と時間表現解析を含む参照表現解析との結果に相当する。統計量や出来事への言及がある文が抽出され、統計量名、時刻、統計量の値、出来事に関する固有表現が注釈づけられている。相対的な時刻表現にはその絶対値が注釈されている。統計量の値や変化に直接関係しない節には、それが出展、原因、評価等への言及であるかの属性が付与されている。具体的な仕様は MuST HP や文献 [8] に詳しい。

5.2 変化表現コーパス

変化表現コーパスは、MuST データセットの一部について、注目している統計量のある時点での変化の情報とその言語表現を抽出したものである。その一部は、注目している統計量のある時点での値という統計量名、時刻、統計量の値の三つ組とそれが抽出された表現で、一般的な関係抽出を行った結果であり、T2N 課題の参考になるデータと言語表現となっている。ただし、このコーパスの特徴は、そのような狭義の定量的な関係に留まらず、定性的なものを含む値と変化の情報と言語表現が抽出されている点にある。

MuST データセットの以下の文章を考えてみる。

原油価格(ドバイ原油)も、昨年10月ごろ1バレル=約20ドルをつけたのをピークに下落が続き、今年1月下旬に同約12ドル50セントまで落ち込んだ。その後、イラク情勢の緊迫化で一時上昇したものの、現在はまた12ドル前後で低迷、… (1998/2/14)

原油価格は指標となるドバイ原油(8月渡し)が15日、18ドル30セント台まで急伸した。2月には10ドルだったので80%以上の上昇になる。(1999/7/17)

ここから「昨年10月ごろ1バレル=約20ドル」等、時点と値の対の情報に加えて、「下落が続き」「～のをピークに」等、ある時間幅や特徴的な時点の変化に関しての定性的な情報を抽出して整理している。「今年1月下旬に同約12ドル50セントまで落ち込んだ」にあるように、時点と値の対は変化の終点や起点としてとらえ、その記事の執筆時点での統計量の値も、「15日、18ドル30セント台まで急伸した」のようにその時点である値となったという変化の終点としてとらえる。つまり、変化の情報を中心的であると考え、そのパラメータとして値を考えている。

この考え方に基づき、「10ドルだった」「10ドルにとどまった」「10ドルになった」「10ドルまで上昇した」では、その時点の統計量の値が10ドルにあることに加えて、それに至る変化の有無や方向が抽出される。また、「10ドルを上回った」では上昇方向の変化があり、10ドルに等しい時点がその時点より過去にあったことが抽出される。つまり、その時点での値に関する情報、T2N 課題で抽出される等値関係だけでなく大小関係が含まれたもの、に加えて、

その時点での値の変化に関する情報(1次微分情報)が抽出されている。この変化の情報にはいつの時点と比較してどれだけ変化したかというような定量的な相対情報も含まれる。更には「上昇を始めた」は上昇方向の変化がその時点で始まったことを示し、「上昇に転じた」はその時点の以前に下降方向の変化があったことを示すというような、変化の変化に関する情報(2次微分情報)まで抽出してコーパスとしている。

MuST コーパス中の9トピック219記事について、1,789件の変化情報が抽出され、変化表現コーパスとなっている。具体的な分析の方針とコーパスの仕様については文献 [7, 10] に詳しい。

6 評価と今後の展開

MuST の場で様々な研究が進められた。それらの研究は、既存の研究分野としては多岐にわたるが、それらに携わる研究者が動向情報という共通の関心で集まり、議論する場を創造できたという点で、MuST は大きい役割を果たした。そこに集まった研究者は、MuST データセットに含まれている新聞記事集合を眺めることで、何が動向情報で、何がそれに係わるかの認識を共有しており、そのことが議論を活性化していると考えられるが、このような認識を作り出したという点で MuST データセットは有益であった。

一方で、データあるいはその仕様として MuST データセットの注釈が期待通りの役割を果たしたかには疑問は残る。MuST データセットの注釈は要約もしくは情報抽出の中間結果という位置づけであったが、情報抽出の研究においては、データセットの注釈ではなく汎用の固有表現抽出システム等による結果を用いるものが多く、注釈結果を前提にその後の処理を研究対象とするものは少なかった。ただ、選ばれたトピックや新聞記事のセットについては利用が多く、共通の入力、処理対象としての役割は果たしたと考えられる。加えて、情報可視化に重点をおく研究では XML パーザを用いて MuST データセットから可視化データの取り出しが行われているし、研究における初期の分析、評価のための参照、正解データの作成としては頻繁かつ有効に利用されているという印象を受ける。その点で注釈付けもその役割を果たしていると言ってよい。

MuST のもうひとつの目的であるツールやデータの共有については、まだ充分とはいえない。多数のデータの提供はあるもののその再利用は残念なことに殆どない。個々の研究の関心が明らかになりその内容も充実し、評価課題も実施されたことで、提供できるデータやツールの質も向上してきていると思われるので、今後の活用期待していきたい。また、可視化プラットフォームもここに加えていく予定である。

評価ワークショップとしては、最初の計画に較べて縮小し、まずはその第一歩を踏み出したという印象である。T2N 課題は、それまでも多くの参加者によって取り上げられていたが、共通の題材を与えることで、それぞれの手

法の特徴を改めて明らかにできた。課題の特徴を反映し、それとシステムの性能の関係をより明確に示す評価指標の提案が今後の課題となる。一方で提案はされたが実施されなかった2つの課題は、評価ワークショップとしての難しさを示している。新しい研究課題に共通の評価を持ち込み、多くの参加者を引きつけるのは必ずしも容易ではない。

現在、NTCIR-8の企画が進められているが、MuSTはそのタスクとしては、提案をおこなっていない。T2N課題を同じ設計で実施することはもちろん可能であったが、評価する技術の明確化、例えば、情報抽出で行われている関係抽出や時間表現処理との関係づけを明確化し、この課題の特徴を反映した評価基準の考案が必要と思われるし、その他の2つの課題についてはいまだに具体化が不十分であった。自由課題については、これまでの参加者は提供された研究資源を用いて研究が続けられる環境にあるし、人工知能学会全国大会「情報編纂の基盤技術」等に発表の場を求めていくことを考えている。

一方、研究課題としてのMuST、動向情報のマルチモーダル要約は、NTCIRが扱うテーマの中で引き続き重要な位置を占めている。特許情報を利用した技術動向の可視化や特許マップの生成をテーマに含めたタスクが、広島市立大学の難波英嗣氏を中心に計画されている。時空間情報を考慮した情報検索のタスクも計画されており、それ自体は対話性や可視化を視野に入れたものではないが、そこで作成されるテストセットはそれらの研究の貴重な素材となると期待される。これらのタスクやそこでの技術課題の展開をにらみつつ、協力させていただきながら、NTCIR-9以降におけるMuSTを考えていく。

研究資源の一般公開も進めている。MuSTデータセット及び変化情報コーパスは近々に研究利用の範囲で一般公開する予定である。国立情報学研究所のご尽力と毎日新聞社のご理解により、注釈づけられた記事については、注釈対象の本文を含めて、研究利用が可能となる予定である。可視化プラットフォームについてもオープンソースでの公開に向けて進めている。

7 おわりに

動向情報の要約と可視化に関するワークショップMuSTについて、その活動内容を報告した。動向情報のマルチモーダル要約という研究分野の重要性、面白さをご理解頂き、そこでの様々な研究に関心を持っていただく契機となれば嬉しい。また、共通の研究資源を活用した非評価型のワークショップから、評価型ワークショップへと展開していくという流れが別の領域での研究活性化においても参考になれば幸いである。

謝辞

MuSTの運営とそれに関連する研究は、NTTと東京大学との産学連携共同研究、ならびに国立情報学研究所のNTTと東京大学との公募型共同研究によって支援されて

います。ご支援をここに感謝いたします。MuSTの活動全般と本稿の執筆はMuSTに参加頂いた皆様によって可能となったものです。あらためて感謝いたします。

参考文献

- [1] 加藤恒昭・松下光範・平尾努 「動向情報の要約と可視化に関するワークショップの提案」 情報処理学会研究報告, 2004-NL-164, pp. 89-94, 2004.
- [2] 加藤恒昭・松下光範・平尾努・神門典子 「評価なきワークショップの試みー「MuST: 動向情報の要約と可視化に関するワークショップ」を例にー」 言語処理学会第11回年次大会併設ワークショップ, 2005.
- [3] 加藤恒昭・松下光範・神門典子 「動向情報の要約と可視化ーその研究課題とワークショップー」 知能と情報 (日本知能情報フェジ学会誌), vol. 17, No. 4, pp. 424-431, 2005.
- [4] Kato, T., Matsushita, M. and Kando, N. MuST: A Workshop on Multimodal Summarization for Trend Information, in Procs of the Fifth NTCIR Workshop Meeting on Evaluation of Information Access Technologies, pp. 556 - 563, 2005.
- [5] 加藤恒昭・松下光範 「情報編纂 (Information compilation) の基盤技術」 第20回人工知能学会全国大会, 1D3-2, 2006.
- [6] 加藤恒昭・松下光範, 神門典子 「動向情報の要約と可視化ー言葉と図で情報をまとめるー」 情報処理, Vol. 47, No. 9, pp. 1013-1020, 2006.
- [7] 加藤恒昭・松下光範 「時系列情報の抽出と可視化に基づく情報アクセスのためのマルチモーダルインタフェースー情報編纂の基盤技術に向けてー」 人工知能学会論文誌, Vol. 22, No. 5, pp. 553 - 562, 2007.
- [8] Kato, T., Matsushita, M. and Kando, N. Expansion of Multimodal Summarization for Trend Information -Report on the First and Second Cycles of the MuST Workshop -. in [11], pp. 235-242, 2007.
- [9] 加藤恒昭・松下光範 「技術展望: 情報アクセスインタフェースとしての要約・可視化と動向情報」 ヒューマンインタフェース学会学会誌, Vol. 9, No. 4, pp. 311-316, 2007.
- [10] Kato, T. and Matsushita, M. Overview of MuST at the NTCIR-7 Workshop Challenges to Multi-modal Summarization for Trend Information, in [12], pp. 475-488, 2008.
- [11] National Institute of Informatics. Proceedings of the 6th NTCIR Workshop Meeting on Evaluation of Information Access Technologies, 2007.
- [12] National Institute of Informatics. Proceedings of the 7th NTCIR Workshop Meeting on Evaluation of Information Access Technologies, 2008.