

アノテーションの概念を用いた情報共有・処理フレームワークの開発 —自然言語処理と視覚化によるコンテンツ生成・獲得支援—

1. 背景

アノテーション技術の発展により任意の種類データやエンティティを対象とした、機械翻訳、情報検索、自動要約、質問応答、知識発見システムなどの実用化や、より高度なデータの加工提示や情報共有が可能となり、その有用性は計り知れない。またアノテーションとはその情報を他人が利用できることでその価値が大幅に増幅されるものであり、アノテーション情報共有のためのフレームワークの構築が望まれる。

2. 目的

本プロジェクトでは、アノテーションの概念に基づく高度情報共有・処理環境の実現を目的とした応用システム構築の上で有用となるライブラリ及びコンポーネントを実装する。

3. 開発の内容

データ主体のシステム構築を行う上で、それに関する要件として大きく“データ構造の決定”、“処理系の実装”、“ユーザインタフェースの構築”が挙げられる。

近年のマルチメディア情報アクセスアプリケーションでは高度のレベルのインタフェースやツールを提供するが、ばらばらの記述仕様、処理系、ユーザインタフェースに頼っているのが現状である。

我々は、平成 14 年度「未踏ユース事業」の中で策定した汎用アノテーション記述言語 MAML (Multimedia Annotation Markup Language) [1] 及び平成 15 年度「未踏ソフトウェア創造事業」にて様々な応用アプリケーションを開発してきた。

つまり、共通のデータ構造を用いることによりソフトウェア開発にかかるコストの省力化を目指した。しかしながら処理系やユーザインタフェースに関しては既存のシステム同様、それぞれ独自に開発する必要があった。

上記問題点を解決するために、今年度は主に以下の3点について実装を行った。

第一にアノテーションデータのための汎用性の高い処理系の構築である。

RDF を例に挙げれば、Java アプリケーション構築のためのフレームワーク Jena に代表される。Jena では利用される領域を想定した上で、データセットの特徴を最大限に活用できるような幅広い枠組みでの機能提供がなされている。本処理系においても、ジェネレータ及びパーザの機能に加え、自然言語中心の構造であることを鑑み、自然言語処理技術を応用した検索、条件無し(大域)要約、条件付要約、関連情報抽出等の処理プロセスを提供する。

第二にアノテーションデータの生成に加え、処理系で提供されている検索、(条件付き、条件無し)要約、情報抽出処理を同一の簡易な操作で実現できる汎用性の高いインタラクティブ型情報可視化インタフェースを実装した。可視化技術の利点としては、マウス操作などによって、キーボード入力とは異なる直接的なインタフェースを与えられることや、情報空間での自分の位置の把握ができることなどが考えられる。逆に、問題点としては、構造がほぼ一意にしか与えられないことと、システム側に主導権があり、ユーザ側には視覚表現を変える手段が十分に与えられていないことが指摘されていた。つま

りはインタラクティブな情報視覚化手法が特に重要である。そこで、可視化技術を用いてアノテーションデータを視覚的にわかりやすい形で提示し、メディア（アノテーションの対象）へのアクセスも実現する。さらにユーザの要求に応じた多様な表示形式を実現する。このインタフェース上でユーザ群が協調してアノテーションデータを生成していく、つまり協調型オーサリングツールとしても機能する。これにより情報の共有や創造、ユーザ主導のマイニング処理をサポートする。スクリーンショットを図1に示す。

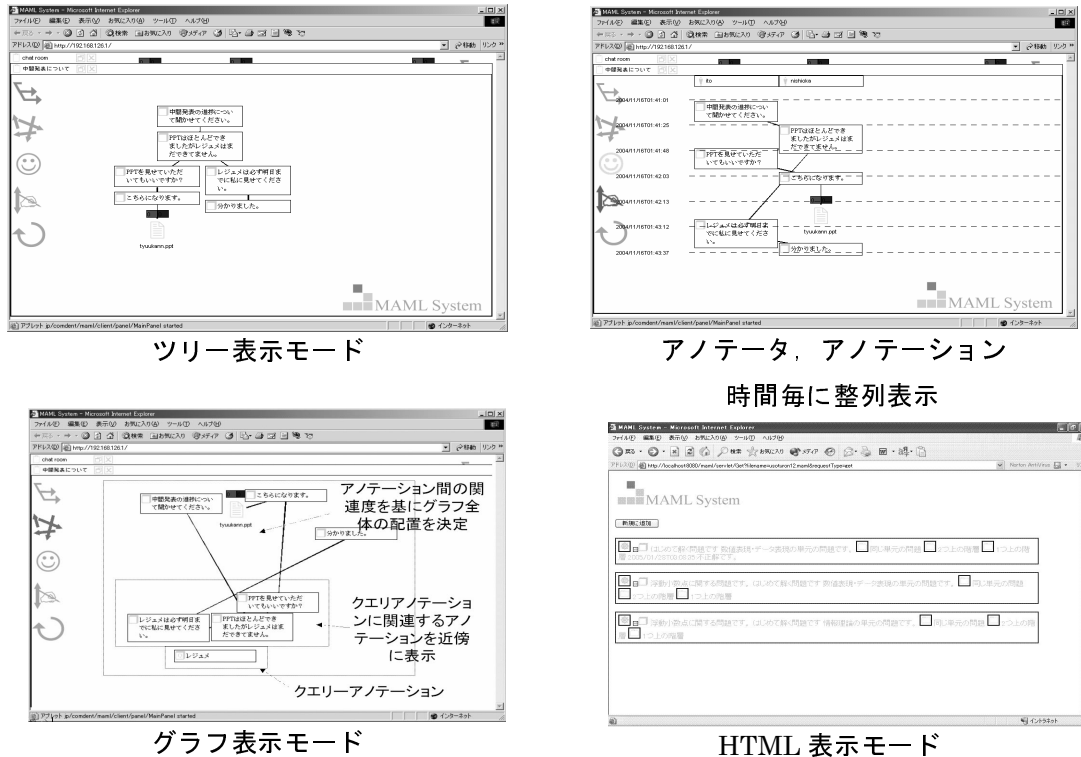


図1 スクリーンショット

第三に、アノテーションサーバの構築である。多人数アノテーションを考える上でサーバは重要である。アノテーションサーバを中心とする全体のシステム構成図を図2に示す。本サーバの特徴として、単にアノテーション情報を入出力するだけではなく、アノテーションの追加を例にあげると、さらに構文解析器によって言語情報を GDA[2]によりタギングする。次に、追加されたアノテーションに関する他のアノテーションとの関連度（類似度）、重要度（スコア）、アノテーションの表示における最適配置座標を計算し、データオブジェクトファイルとして保持し、これらの情報が更新されたことを各クライアントに通知する。これにより、クライアントは必要に応じて、アノテーションデータだけではなく、その類似度、類似度、配置座標に関する情報も獲得することができる。さらには、クライアントからのアノテーションの追加などに際して、その情報を瞬時に他のクライアント群に通知できる。これによりチャットのようなリアルタイム性が要求される仕組みも実現する。

本システムを利用した応用事例の一部を図3に示す。これら多くの機能を横断的に利用することが可能である。

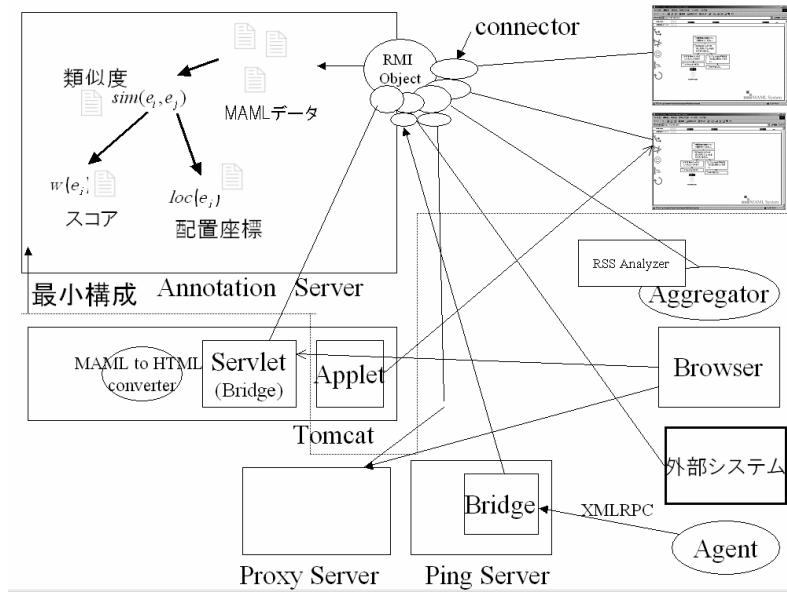
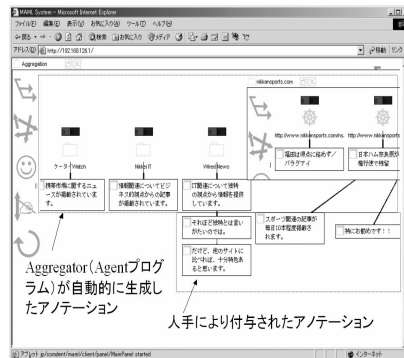


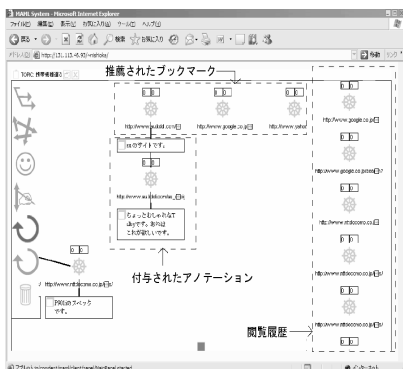
図2 システム構成図



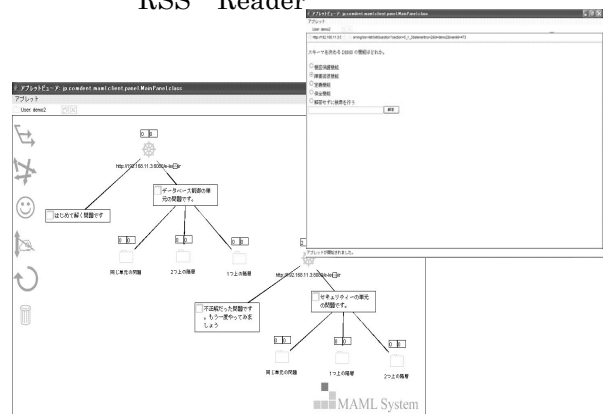
Google 検索ラッパ



RSS Reader



協調型ブラウジング支援



e-Learning システム

図3 応用事例

4. 従来の技術(または機能)との相違

これまでに、特定のドメイン、例えば Web マイニングや動画解析などの分野で可視化コンポーネントが開発されてきたが、可視化することが主目的でそれらに対するユーザの多様な要求をうまく反映させることが出来なかった。本フレームワークは、任意のドメインにそのまま適用でき、またユーザの能動的操作により対話的やり取りを行いながら情報を生成および獲得できる。

5. 期待される効果

汎用のコンポーネントを開発することはデータマイニング、アクティブマイニング研究やそれに関わるシステム開発領域において有用性が高いと考えられる。また、将来的に Semantic Web 分野への応用を視野に入れている。Semantic Web で想定している高度に構造化されたデータは、本研究で題材とした自然言語中心のデータに比べて、言うまでもなくその応用を考える上での技術的障害は少なく、またその活用領域は比較にならない。しかし、それに係る構築のためのコストの問題が解決されず、実際にデータが生成されていなければ、卓上の理論に終わってしまうであろう。MAML 程度のデータの構造化が現状、自然な形で構築されうる妥当な水準と考えている。当然、MAML のデータから意味解析レベルの自然言語処理や機械的学習論等の要素技術を組み合わせ、さらにはその過程で生じた精度の低下や過りは人手による修正作業により、段階的かつ部分的であれ、計算機が自立的に解析可能な、オントロジーが構築していくというスタンスも必要であろう。つまり一般の多くの人間が間接的であれ、オントロジー生成に干渉できるような方策というのが望ましいと考える。Semantic Web が掲げる目的は、人工知能に関わる研究者共通の夢である。そういう意味で、将来的に Semantic Web 構築の一助を担うことができるのでないかと考えている。

6. 普及(または活用)の見通し

実際に何らかのシステムを構築する上で、独自の追加実装はほとんど何も必要ないという点が本提案の最大の利点である。用途はほぼアノテーションデータの内容にのみ依存することとなる。つまりはプログラミングの能力が全くない多くの一般ユーザでも様々な応用システムを構築できる。これにより有用なアノテーション(コンテンツ)が多く創造されることが期待される。

また、プロジェクト終了後より積極的に研究会・ワークショップでのデモや論文誌掲載など精力的に啓蒙活動を行っている[3]-[6]。

7. 開発者名(所属)

*伊藤 一成(慶應義塾大学 大学院理工学研究科後期博士課程)

2005年4月より青山学院大学理工学部情報テクノロジー学科助手

[参考文献]

[1]伊藤一成, 斎藤博昭: 汎用アノテーション記述言語 MAML の提案とその生成・処理プロセス, 情報処理学会論文誌 TOD, Vol. 45 No. SIG7 (TOD22),

pp. 137-150, (2004).

- [2] 橋田浩一：GDA 意味的修飾に基づく多用途の知的コンテンツ，人工知能学会論文誌，Vol. 13, No. 4, pp. 528-535, (1999).
- [3] 伊藤一成：アノテーションの概念に基づく情報可視化インタフェース，日本データベース学会論文誌 DBSJ Letters, Vol4 No. 1 (2005)
- [4] 伊藤一成，斎藤博昭：アノテーションの概念に基づく情報可視化インタフェースの提案，電子情報通信学会・データベース学会合同ワークショップ DEWS, (2005)
- [5] 山根木 浩平，伊藤 一成，斎藤 博昭：汎用アノテーションシステム(MAML System)の e-Learning への適用，電子情報通信学会・データベース学会合同ワークショップ DEWS, (2005)
- [6] 滝本湖，伊藤一成，斎藤博昭：汎用アノテーションシステム (MAML System) を利用した Web 検索結果のグラフ表示，データベースワークショップ DBWS2005 , (2005)