

鬼沢 友和 安井 浩之 松山 実
武蔵工業大学

1 はじめに

コンピュータネットワーク上において、構造化された文章やデータを管理、交換、配布するためのマークアップ言語として XML (eXtensible Markup Language) [1]が注目されている。従来の HTML に比べて文書要素の名前が自由に定義できるなど XML 文書を作成する際の自由度が高いという反面、内容が類似するものでも文書の構造は作成者により大きく異なることもある。XML 文書をコンピュータで処理する際、文書型定義 (Document Type Definition) から文書の構造を判断するため、異なる構造の数だけ DTD を定義しなければならない。ここでは無駄に数の多い定義を避けるため、複数の XML インスタンスから統一化された DTD を生成する支援ツールの開発を試みている。

2 手作業による XML 化の問題

本来 XML インスタンスは DTD に基づいて作成されるが、XML には、タグの構造を自由に決定ができるという特徴がある。そのため、既存の文書やデータ (構造化データ) から XML 文書に変換する際、先に XML インスタンスを作成し、インスタンスに合わせて DTD を作成するという場合も多い。しかし、この手順を用いると僅かな違いしかない構造化データに対しても個々に固有の DTD が必要となり効率が悪い (図 1)。また、同じ構造化データでも別の人が XML 化すると異なった構造のインスタンスになってしまう可能性がある (図 2)。このように DTD も XML インスタンスも一意には決まらな

Production of unified DTD from XML Instances

Tomokazu Kizawa, Hiroyuki Yasui, Minoru Matsuyama

Musashi Institute of Technology

いという問題がある。

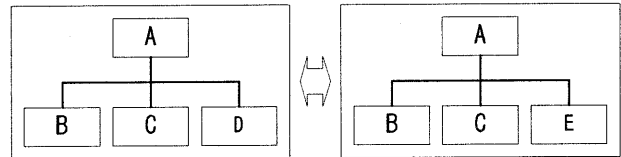


図 1. 固有の DTD

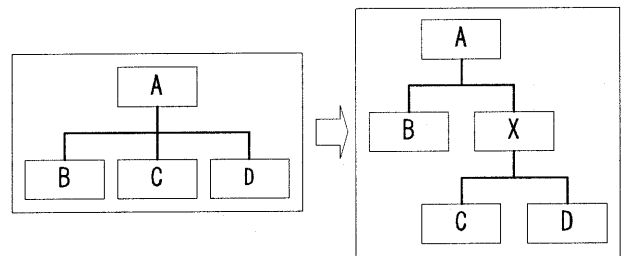


図 2. 同じ意味でも異構造

3 支援ツールの概要

ここに報告する XML 文書作成支援ツールは実際の構造化データから XML 文書化する際に DTD の制約を気にすることなく、インスタンスから作成することを可能にし、既に存在する XML 文書に情報を追加したり、まとめるなどの作業を可能な範囲で自動化するものである。その機能を以下に示す。ただし本ツールは XML エディタの機能を持たないので、テキストエディタや XML エディタで作成された XML インスタンスを別途用意する必要がある。

1)XML インスタンスからの DTD の自動生成
DTD の無いインスタンスの構造を解析して適合する DTD を自動生成する (図 3)。

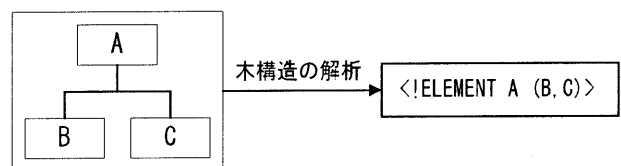


図 3. DTD の自動生成-1

2)複数のインスタンスからの DTD の自動生成

内容は似ているが構造が異なるインスタンス同士を比較して全てに適合する DTD を生成する。ただし<名前>と<氏名>のように一つにまとめる方が良くとユーザが判断した場合には同じものとして処理する。具体的にはユーザが選択、もしくは木構造の幅や深さを比較してユーザが望むものをコンピュータが選択し、それを基準にその他のインスタンスと比較していく (図 4)。

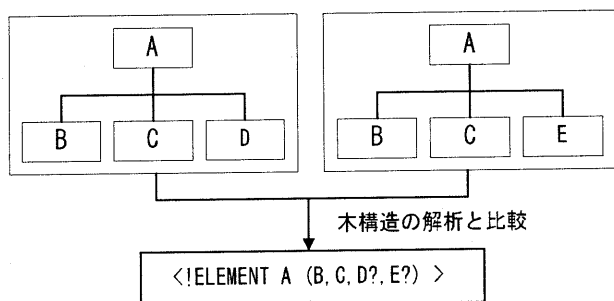


図 4. DTD の自動生成-2

3)インスタンスの構造の統一

2)で生成された DTD に基づいて該当するインスタンス全てを適合するように書き換える。またエンティティの宣言や属性の追加などの作業も行う。これにより対象のインスタンス全てが同じ構造になる。

4 試験的実装

本格的な実装の前段階として試験的実装を行った。

4.1 実装環境

Java2 が動作する Windows98 上で DOM Level1^[2], SAX version1 の各パーサ API を呼び出すことが可能な環境で行った。

4.2 実装例

2 つの構造の異なる XML インスタンスから共通の DTD を生成し、その DTD に適合するようにインスタンスを整形した (図 6, 図 7)。

なお、図 6 の例では instance I は統一 DTD に適合しているので書き換えは行われない。

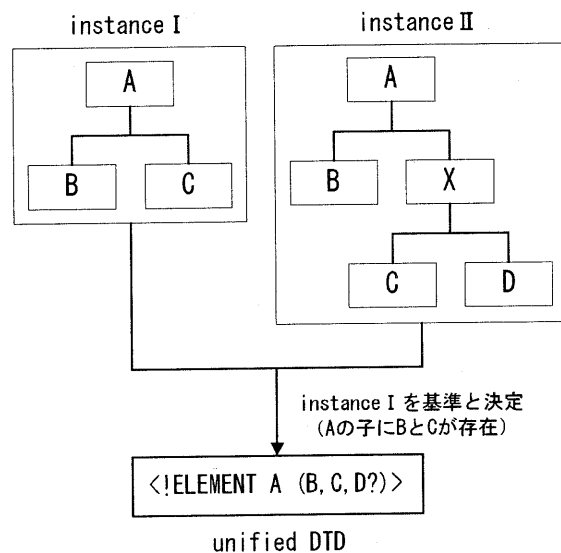


図 6. 実装例 (DTD の統一)

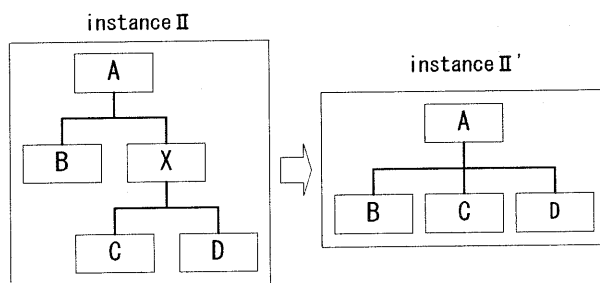


図 7. 実装例 (インスタンスの変更)

5 おわりに

今回の試験的実装は単純なモデルを対象としたが本来は大量の文書を処理することを目的としているため、最終的な実装段階では膨大なデータに対処できるようなシステムにする必要がある。またユーザから指示を与える場面が多いので、統一の自動化アルゴリズムの検討や、容易な操作を目指す必要がある。

参考文献

- [1] W3C, "Extensible Markup Language (XML) 1.0", <http://www.w3.org/TR/REC-xml/>, 1998
- [2] W3C, "Document Object Model (DOM) Level 1 Specification", <http://www.w3.org/TR/REC-DOM-Level-1/>, 1998