

D-29 XML 文書の事前形式変換によるデータ処理性能改善の検討 Improving Performance of Processing XML Documents by Format Transformation in Advance

吉田 茂† Shigeru Yoshida 矢作裕紀† Hironori Yahagi 井谷宣子† Noriko Itani

1. まえがき

近年、伝票や顧客名簿等、レコード構成のデータへの XML 適用が BtoB 用途で広がっている。XML はデータ表現力、柔軟性、拡張性が高い反面、データ処理時にはその冗長性が負担になる。現状は、数千レコード(数 MB)以上の XML 文書は、通常、RDB に格納してデータ処理する。XML 文書をランダムに操作する場合、W3C 標準の API である DOM [1]を用いるが、DOM によるデータ処理は、CPU 負荷が高く、XML ファイル容量の 5~10 倍もの大量のメモリを消費し、メモリ上で大容量 XML 文書を扱うには課題があった。本稿では、レコード構成の XML 文書に対して DOM 適用ソフトでのデータ処理性能を改善するため、XML 文書を別な形式の XML 文書に事前に変換する方法を検討・評価したので、結果を報告する。

2. 従来技術と狙い

負荷を軽減する XML 形式として、携帯機器の WAP [2]で使われているバイナリ XML が知られている。バイナリ XML は特定用途には有効だが、全ての用途向けに一般化したバイナリ符号を作るのは困難であり [3]、新たに専用 DOM ソフトも必要であり、W3C XML 規格と整合が難しい。本検討では、実用性を重視し、適用ソフトにわずかな修正で適用できて、変換後に基本的に元の XML 文書と同様(トランスパレント)に扱える形式を目指すことにした。

XML の導入現場では、既存データを XML 化することで処理能力不足が生じる場合も多く、これを補うため、ユーザーによって個々に XML 文書の形式を修正する幾つかの方策が採られている [4]。基本的な方策は、(a) XML の定義情報の長さを短くする、(b) 冗長な情報を一つにまとめる、(c) 業務に不要な XML データを取り除く、である。

XML 文書の要素名を省いてデータ量を圧縮するため、Trotter [5]は、レコード毎に全要素内容を CSV(Comma Separated Values)形式で繋いで 1 個の要素にまとめた XML 文書に変換する方法を「ZXML」と称して提案した。ZXML では、XML 文書の変換は XSL 変換 [6]として実行し、XML 文書がレコード内 1 階層の場合の簡単な例を示した。変換 XML 文書のデータ処理には専用ソフトを用いることを想定した。

以降では、これらの従来技術をデータ処理性能改善のツールとして使える汎用手法に高め、XML 文書の事前変換後の形式と、その変換の実行法について述べる。

3. 技術開発

3.1 XML 文書の事前変換形式

DOM 処理は XML 文書の全要素をメモリ上に展開するため負荷が重くなるが、その適用ソフトでは全種の要素をレコード横断のデータ処理の対象にすることはまれである。冗長な情報を一つにまとめて要素数を減らすことが有効な

ので、レコード毎に、適用ソフトで処理対象とする要素はそのままにして、処理対象としない複数の要素の要素内容を CSV 形式で一つの要素にまとめた XML 文書に変換する(図 1 (a), (b))。これにより、冗長な要素名が省かれ、要素数が削減されて、DOM 展開時の負荷が軽減される。本方式を「CSV 圧縮変換」と呼ぶ。

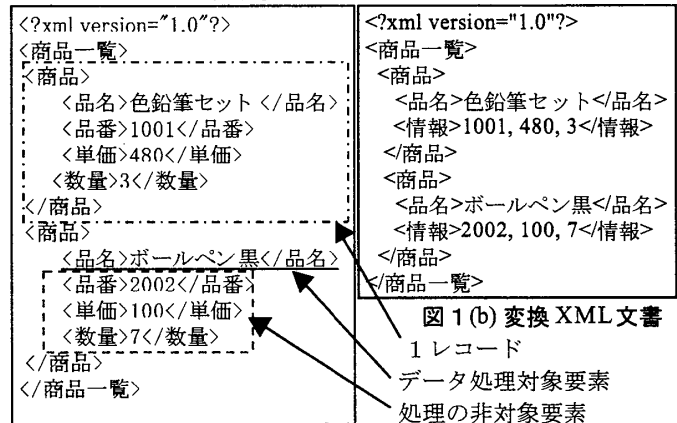


図 1 (a) 変換前 XML 文書の例

3.2 CSV 圧縮変換の実行

XML→XML 文書形式変換である CSV 圧縮変換を、「XSL 変換」として実行する(図 2)。特別なソフトを作らず変換が実行できるように、処理対象要素名と非対象要素名を列挙した「変換仕様 XML 文書」(図 3)をユーザが作成すれば、XSL 変換/逆変換に用いる XSL シートを自動的に生成できるようにする。この自動生成も、XSL シートを用いて、XSL 変換として実行する。

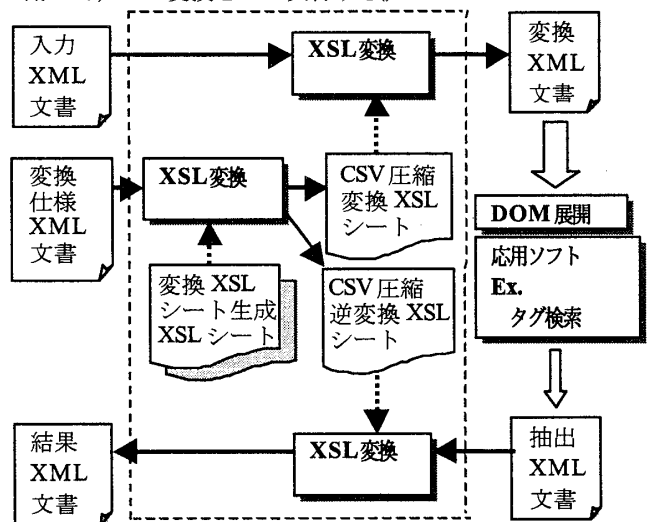


図 2. 変換 XSL シート生成、CSV 圧縮変換、及び、応用ソフトへの適用の全体フロー

† (株) 富士通研究所 パリフェラルシステム研究所

```

<?xml version="1.0"?>
<transform>
<structure>
  <root>商品一覧</root>
  <record>商品</record>
</structure>
<items>
  <merging_tag>情報</merging_tag>
  <key>品名</key>
  <nonkey>品番</nonkey>
  <nonkey>単価</nonkey>
  <nonkey>数量</nonkey>
</items>
</transform>

```

各タグによる指定

<root> : ルート要素名

<record> : レコード要素名

<merging_tag> : 変換先 CSV 形式の要素名

<key> : 処理対象要素名

<nonkey> : 処理非対象要素名

図 3. 図 1(a) から図 1(b) への変換仕様 XML 文書の例

図 1(b) のような表形式データでは、処理非対象要素名は変換/逆変換 XSL シートの中に吸収され、変換 XML 文書中から省かれる。この作用によって、変換 XML 文書のデータ量は、レコード毎に全要素を CSV 形式で 1 要素にまとめた場合、元の最大約 1/3 に圧縮される。

3.3 複雑な XML 文書の扱い

XML 文書の特徴は、非表形式や、深い階層構造のデータが表現できることである。複雑な XML 文書は次のようにして扱うことができる。これらは変換仕様 XML 文書で指定して図 2 と同様に変換 XSL シートを生成すればいい。

(1) 非表形式のデータの変換形式

レコード内に出現する要素が不定であるため、変換 XML 文書上に、処理非対象の要素名を保持して、要素内容と対応付ける必要がある。CSV 形式の要素内容と同じ並び順で、処理非対象要素名を CSV 形式で繋いで、変換先 CSV 形式の要素の属性として置く。

CSV 圧縮変換によって非表形式データの要素数を最大限に減らしたいときは、レコード内の要素を、固定して現れる要素と、不定の要素との 2 グループに分け、それぞれを変換先の表形式の CSV 要素と、CSV 形式の要素名を属性に置いた非表形式の CSV 要素、の 2 個に格納すればいい。

(2) レコード内 2 階層以上、属性を持つ場合の変換形式

表形式データの場合は前述と同様に、処理非対象要素の要素名、属性名は変換 XSL シートに吸収される。非表形式データの場合は(1)で述べた変換先 CSV 形式の要素の属性に、処理非対象要素や属性のレコード内の階層位置を XPath で表して、CSV 形式で繋いだものを置く。その XPath の表現が長くなる場合は、XPath を短縮名に置換えれば、データ量を増やさずに一意に識別することができる。

4. 入出力性能の改善効果

レコード内 2 階層の XML 文書が扱える変換/逆変換 XSL シート生成 XSL シートを作成して、方法論の確認と、変換 XML 文書の性能評価を行なった。CSV 形式にまとめて減らした要素数の割合にほぼ比例して、DOM のメモリ使用量、メモリへの展開時間、DOM 木を辿る時間が削減された (図 4 (a),(b),(c))。CSV 圧縮変換は、同じ評価環境で、レコード毎に全要素を 1 個の CSV 要素にまとめたとき 10MB の XML 文書の変換が 20 秒、逆変換が 8 秒であった。

5. むすび

本技術の特長は、アプリケーションソフトでデータ処理の対象にする要素の種類が予め分かれば XML 文書を必要最小限の DOM 木で扱えること、XSL 変換として実行するため種々の言語で利用できること、特別なソフトを書かずに変換が実行できること、である。種々の実際の応用で本技術の効果を確認することが今後の課題である。XML 技術は普及途上であり、様々な便利ツールが必要とされるが、本技術もその一端を担うものと考えられる。

参考文献

[1] DOM <http://www.w3.org/DOM>

[2] WAP(The Wireless Application Protocol), <http://www.atmarkit.co.jp/aig/01xml/wap.html>

[3] "Intuition and Binary XML", <http://www.xml.com/pub/a/2001/04/18/binaryXML.html>, 2001.04.16

[4] 「見えてきた万能幻想の真実 XML の”常識”を覆す」 pp. 52-71, 日経コンピュータ 2001.2.26 号

[5] Alain Trotter, "Building an XML Bloat Buster using ZXML -- XML Compression Mehod", 有料情報配信 Web サイト <http://www.asptoday.com> 2001.1.3 の記事 概要は <http://www.xml.com/pub/r/904>

[6] XSLT <http://www.w3.org/TR/xslt>

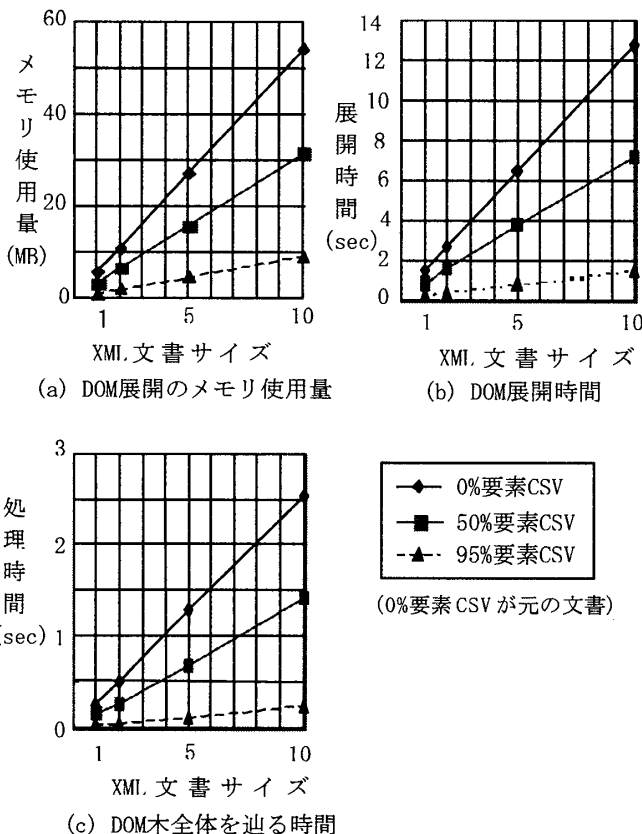


図 4. CSV 圧縮変換の性能評価

評価環境: PentiumII 450MHz Windows2000 マシン
 JDK1.3, DOM ソフト Crimson 1.1 で測定
 試料: レコード内 1 階層の XML 文書, 50 要素 (2KB) / レコード