

5J-01 XML をベースとしたタグ格納による検索方法の提案とその効率化

近藤 博[†] 川越恭二[‡]

立命館大学大学院理工学研究科[†]

立命館大学工学部[‡]

1 はじめに

現在, 多くの分野で文書管理に XML が利用され始めており, XML で記述された文書が増大していくことが予想される. しかし, これらを効率良く格納, 検索する技術はまだ十分ではない.

そこで本稿では XML をベースとするタグを利用した XML 文書の効率的な検索方法を提案する. 本方式は従来のように, 検索対象データをシステム側に索引データとして保持するのではなく, タグと文書の場所のみを格納し直接対象データ内を検索するものである. この方法により, 頻繁に更新が行われる XML 文書の検索にも対応可能となる. さらに検索キャッシュと部分インデックスにより効率化した方法を以降で説明する.

2 XML タグ格納方式

提案する方式は物理的に分散されたサーバに分割, 格納された XML 文書ファイル集合を直接検索する方法を実現するものである.

2.1 基本方式

XML 文書ファイル F_i の集合 F を検索対象とする. また, XML 文書ファイルは XML タグとそれに対応するタグ値のペア集合からなる. すなわち, 任意の $F_i \in F$ について, $F_i = \{(T_i(j), V_i(j))\}$ とする. ここで, $T_i(j)$ はタグ集合の要素であり, XML 文書ファイル F_i の中に記述されている j 番目のタグを示す. $V_i(j)$ はこのタグに対応するタグ値である. さらに, XML 文書ファイルは適当に設定された格納アドレス $\text{Address}(F_i)$ が割り当てられ, $\text{Address}(F_i)$ が与え

られれば F_i を取り出すことが出来るものとする.

この文書群に対する検索を行うために指定する質問 Q は以下のように記述されるものとする. $Q = \{(T_k, V_k)\}$ ここで (T_k, V_k) は質問 Q を構成する k 番目のタグとタグ値のペアである. したがって, 質問 Q は k 番目のタグとタグ値のペアを持つ XML 文書ファイル F から選択することを意味する. 質問 Q に対する検索結果 R_Q は K_Q 個の XML 文書ファイル集合 $\{F_k\}$ である. したがって, $\{\text{Address}(F_k)\}$ が求めるべき結果となる.

上記の結果を得るためにシステム内に保持するデータファイル D は (T_s, A_s) の形式をとる. ここで, T_s はタグであり, A_s はタグ T_s を持つ XML 文書ファイル F_{T_s} の格納場所を示す. すなわち, $A_s = \text{Address}(F_{T_s})$ である.

質問 Q が与えられた時, システムはデータファイル D を用いて以下の手順で R_Q を求める.

Step-1

次の条件を満たす格納場所の集合 $A = \{A_i\}$ を求める. “任意の $A_i \in A$ に関して, 質問 Q のすべての T_k について $(T_k, A_i) \in D$ が成り立つ”.

Step-2

任意の $A_i \in A$ について, $A_i = \text{Address}(F_{T_i})$ なる $F_{T_i} \in F$ を取り出し, 次の条件を満たせば A_i を R_Q に含める. “任意の j, k について $T_{T_i}(j) = T_k$ が成り立つならば, $V_{T_i}(j) \in V_k$ が成り立つ”. ここで, $T_{T_i}(j)$, $V_{T_i}(j)$ はタグ T_i を持つ XML 文書ファイル F_{T_i} のタグおよびタグ値である.

2.2 処理の効率化

前節で述べた基本手順をより処理効率化するために, 以下の2つの方法を導入する.

- ・部分インデックスの利用

・検索キャッシュの利用

部分インデックスとは、予め頻繁にアクセスされるXML文書ファイルの部分集合 $F' \subset F$ について作成したインデックスファイル $I = \{T(m), V(m), Address(F_m)\}$ である。ここで $F_m \in F'$ であり、 $T(m)$ はXML文書ファイル F_m の中の適当なタグ、 $V(m)$ はそのタグに対応したタグ値である。

検索キャッシュとは、最近アクセスしたタグおよびタグ値とそれを格納しているアドレスの組からなるキャッシュインデックスファイル $C = \{T[p], V[p], A[p]\}$ である。ここで、 $T[p]$ 、 $V[p]$ 、 $A[p]$ は最近 p 番目にアクセスしたタグ、タグ値、アドレスである。

この2つの方法を用いることで、前節のStep-2をStep-2'に改良する。

Step-2'

(1)もし、“すべての k について $T_k = T[p]$ 、 $V_k = V[p]$ 、 $A_i = A[p]$ 、 $(T[p], V[p], A[p]) \in C$ となる p が存在する”という条件を満たすような $A_i \in A$ があれば、 A_i を R_Q に含める。この A_i を A から削除する。

(2)もし“すべての k について $T_k = T(m)$ 、 $V_k = V(m)$ 、 $A_i = Address(F_m)$ 、 $(T(m), V(m), Address(F_m)) \in I$ となる m が存在する”という条件を満たすような $A_i \in A$ があれば、 A_i を R_Q に含める。この A_i を A から削除する。

(3)もし $A_i = Address(F_{T_i}) \in A$ なる $F_{T_i} \in F$ を取り出し、次の条件を満たせば A_i を R_Q に含める。“任意の j, k について $T_{T_i}(j) = T_k$ が成り立つならば、 $V_{T_i}(j) \in V_k$ が成り立つ”。

2.3 効率化方法の有効性、問題点

部分インデックスを利用した場合、データを参照する時間は大幅に短縮できる。しかし、インデックスファイルを保持するためデータ量の増加によるシステムへの負荷が増加する。また、検索結果の新鮮度が低下するといったデメリットも生じる。

検索キャッシュを利用する場合、ゆるい検索条件を指定して検索を行うとキーワードに適合する確率が下がるため、キャッシュファイルを検索する処

理が無駄になるという問題がある。しかし、特定の分野での検索など、必然的に検索キーワードやタグが限定されるような検索に対して適応すれば検索の効率化が実現できると考えられる。

3. 実装

提案した方式を元にXML文書論文に対する検索システムを試作した。

3.1 システムの構成

システム側は文書中のタグ(タイトル、著者名、発表年度)とその文書の所在場所(アドレス)を保持する。検索要求に対して指定されたタグに対応するアドレスを元に対象文書内を直接検索し、キーワードと合致したもののアドレスを出力する。

3.2 制約事項

事前にXMLのメタデータRDF^[1]を文書ファイル内に埋め込んで置く必要がある。

3.3 アプリケーションの試作

試作システム、実行結果を図1、図2に示す。

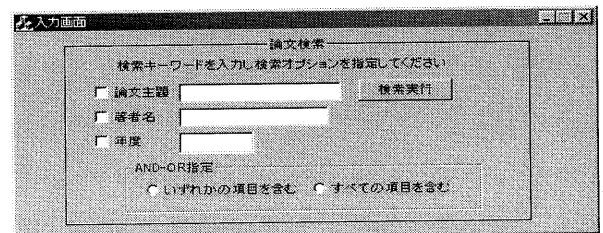


図1.論文検索システム

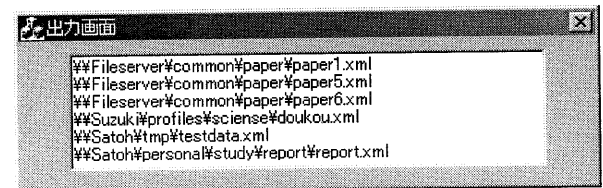


図2.実行結果

4. おわりに

今後は、処理性能の評価や検索の効率化を予定している。

参考文献

[1]“Resource Description Framework(RDF)”

<http://www.w3.org/RDF/>