

전자문서의 XML 문서로의 변환 및 저장 시스템

주원균^o 양명석 김태현 이민호 최기석
한국과학기술정보연구원

{joo^o, msyang, heemang, cokeman, choi}@kisti.re.kr

Rule Based Document Conversion and Information Extraction on the Word Document

WonKyun Joo^o, MyungSeok Yang, TaeHyun Kim, MinHo Lee, KiSeok Choi
Korea Institute of Science and Technology Information

요약

본 논문은 HWP, DOC와 같은 전자 문서에서 사용자가 제공한 구조적인 규칙과 XML 기반 전자 문서 변환 기법을 이용함으로써, 사용자의 관심 영역에 해당하는 다양한 형태(표, 리스트 등)의 정보를 효과적으로 추출(변환)하여 저장하기 위한 방법에 관한 것이다. 본 논문에서 제시한 시스템은 3가지의 중요한 요소들로 구성되어 있는데, 1)전자문서의 원시 XML 문서로의 변환 방법, 2)XML 기반 구조적인 규칙과 작성된 규칙을 이용하여 원시 XML 문서에서 정보를 추출(변환)하는 방법, 3)추출된 정보에서 최종 XML을 생성하거나 DB에 저장하는 방법이 그것이다. 전자문서의 변환을 위해서 독립적으로 동작하는OCX 기반의 전자문서 변환 데몬(Daemon)을 개발하였고, 사용자의 정보 추출(변환)과정을 돕기 위해서 XSLT를 확장한 형태의 스크립트 언어를 개발하였다. 스크립트 언어는 비교적 간단한 문법 구조를 가지고 있고, 데이터 처리를 위한 자체 정의 함수와 변수를 사용한다. 추출된 정보는 원하는 형태의 데이터 포맷으로 생성하거나 DB에 저장할 수 있다. 본 시스템은 전자 문서 원문 정보에 대한 데이터베이스 구축 및 서비스의 제공, 혹은 구축된 데이터베이스를 이용하여 다양한 현황 통계를 제공하는 분야에서 유용하게 사용할 수 있다. 실제로 연구과제관리시스템과 성과정보시스템에 적용하여 그 성과를 입증하였다.

1. 서론

XML 문서가 인터넷을 비롯한 다양한 분야에서 정보 교환을 위한 표준으로 널리 사용되면서 XML 문서의 변환에 대한 필요성이 널리 인식되고 있다. 더욱이 정보교환과정에서 동적으로 생성되는 XML 문서뿐만 아니라, 오프라인에서 작성되는 워드 프로세스와 같은 전자문서 편집기의 결과물인 전자문서들도 고유 DTD 혹은 스키마를 이용하여 새로운 XML 형태로 변환하는 방법에 대한 필요성이 날로 증가하고 있다.

XML의 출현과 HWP의 인기와 더불어 이와 관련하여 많은 연구들이 진행되었다. 관련연구로서 [1]은 HWP 문서를 전자책 표준 중의 하나인 EBKS 문서로 변환하기 위한 기법에 관한 연구이다. 문서 작성과정에서 HWP에서 제공하는 고유기능인 스타일을 이용하는데, 스타일은 문서의 제목과 수준을 제공하는 역할을 한다. HWP의 스타일을 이용하여 변환하고자 하는 영역에 미리 정의된 스타일을 지정하고, HWP 문서를 HWP의 XML형태인 HWPXML로 저장한다. 저장된 XML 문서에 대상으로 XSLT 변환[2]을 수행하여 전자책(EBKS)으로 최종 변환하는 방법을 제안하였다. 이전 연구가 HWP의 일부 구조적인 특징을 이용한다면, 원시 XML 문서가 완전한 계층적인 형태를 가지고 있는 도메인에 대한 연구들도 수행되었다. [3][4][5]에서는 원시 XML 문서와 대상 XML 문서 사이에 템플릿을 이용하여 DTD 매핑(mapping)을 수행하고, XML 문서간의 자동변환을 처리하는 방법을 연구하였다. 이 연구들의 주안점은 템플릿을 이용한 변환을 위한 XSLT 스크립트를 생성해내는 방법에 관한 것이다. 한글과 컴퓨터에서는 보다 근본적인 방법으로 전자문서에서 사용자가 원하는 XML을 생성하는 방법을 제안하

였다[6]. XML 스키마와 HWP 문서를 접목한 한글 XML 서식을 이용하여 원시 문서를 XML 형태로 생성해냄으로써, 사용자가 보다 쉽게 전자문서에 접근할 수 있는 방법을 제공하였다. 그러나 HWP에 특화된 방법이라는 점에서 다른 형태의 전자 문서로의 적용은 불가능하다.

위의 연구들은 전자문서 편집기의 특정 지원 기능을 이용하거나 원시 XML 문서의 내용적인 구조성을 요구하고 있다. 본 논문에서는 물리적인 문서 구조상으로는 계층구조를 가지고 있지만 내용적 평면성을 보이는 범용 전자문서 편집기의 결과 XML을 대상으로 한다. 본 논문에서는 전자문서를 원시 XML로 변환하고, 이 중에서 사용자가 선택한 부분을 사용자의 의도에 맞추어 구성된 XML 문서로 변환하여 DB에 저장하는 방법을 제안한다. 특히 이번 논문에서는 전반적인 시스템의 설계 및 구현에 초점을 맞추었다. 2장에서는 시스템의 기본적인 설계 구조에 대해서 설명하고, 3장에서는 본 논문에서 사용한 XML 문서들 간의 변환 방법에 대해서, 4장에서는 결과 XML 문서를 저장하는 방법에 대해서 설명하고, 5장에서 결론을 맺는다.

2. 전자문서 변환 및 저장 시스템

본 논문에서 제안하는 전자문서 변환 및 저장 시스템은 [그림 1]과 같이 크게 3부분으로 구성되어 있다.

□ HWP2HML 변환

HWP2HML 변환 부분은 전자문서의 한 형태인 HWP를 원시 XML 문서로 저장하기 위한 하위 모듈이다. HWP는 현재 원시

XML 문서 저장과정에 HWPML 2.1을 사용하여 확장자 HML을 가지는 원시 XML 문서를 생성한다. HWP 변환을 위해 한글과 컴퓨터에서 제공하는 OCX 라이브러리를 사용하는 HWP2HML 변환 데몬을 개발하였고, 데몬은 윈도우즈 환경에서 쓰레드 기반 서비스 형태로 동작 하도록 설계하였다([그림 2] 참조).

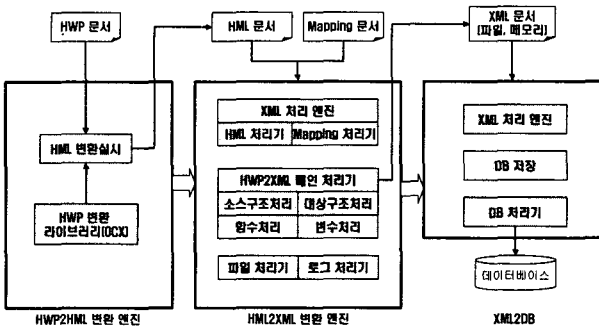


그림 1 전자문서 변환 및 저장 시스템 구조

본 연구에서 제안하는 방법은 범용 전자 문서에 적용 가능한데, 우선 HWP 문서를 대상으로 하여 설계/구현하였다. 향후에 다른 전자문서에 적용하기 위해서는 HWP2HML 변환 엔진 부분을 확장하는 것으로 충분하다.

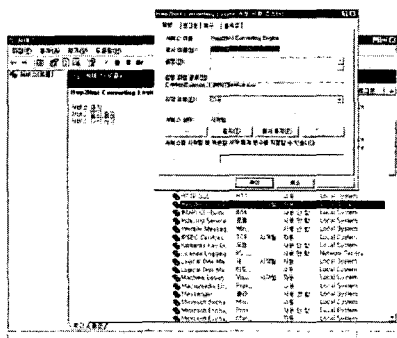


그림 2 HWP2HML 변환 데몬 서비스 적용 화면

□ HML2XML 변환

이 부분에서는 원시 XML 문서를 사용자 위주의 XML 문서 형태로 변환한다. 원시 XML 문서는 HWPML을 이용하여 구성된 문서로서 [그림 3]과 같은 형태를 가지고 있다. 원시 XML 문서는 물리적으로 구조적인 형태로 구성되었지만, 내용적 계층성은 포함하고 있지 않다. 또한 전자문서의 특성상 폰트를 명시한 부분과 같은 불필요한 태그를 많이 포함하고 있다. 전자 문서의 변환에 관한 부분은 3장에서 자세히 설명한다.

□ XML2DB 저장

최종 XML 문서는 이 모듈을 이용하여 DB에 저장한다. 현재 DB에 저장할 수 있는 XML 문서는 관계형 DBMS에 적합한 형태를 가지는 XML 문서라는 제약점을 가지고 있다.

```
<?xml version="1.0" encoding="UTF-8" standalone="no" ?>
- <HWPML Style="embed" SubVersion="5.0.0.0" Version="2.5">
- <HEAD SecCnt="1">
- <DOC SUMMARY>
  <TITLE>자재구매의뢰서</TITLE>
  <DATE>2001년 6월 26일 화요일, 13시 46분</DATE>
</DOC SUMMARY>
+ <DOC SETTING>
<INSIDEMARGIN Bottom="140" Left="140" Right="140"
Top="140" />
- <ROW>
- <CELL BorderFill="2" ColAddr="0" ColSpan="1" Dirty="false"
Editable="false" HasMargin="false" Header="false" Height="5720"
Protect="false" RowAddr="0" RowSpan="2" Width="2196">
  - <PARALIST LineWrap="Break" TextDirection="0"
VertAlign="Center">
  - <P ParaShape="12" Style="0">
  - <TEXT CharShape="0">
  <CHAR>결</CHAR>
</TEXT>
</PARALIST>
</CELL>
</ROW>
...

```

그림 3 HWPML을 이용한 문서 구성 예제(원시 XML, HML)

3. XML 문서 간의 변환 방법

본 논문에서 제안하는 XML 문서 간의 변환 방법은 일종의 정보 추출과 유사한 내용으로서 전체 문서를 대상으로 하지 않고, 전체 문서 중의 일부 관심 부분으로 한정한다.

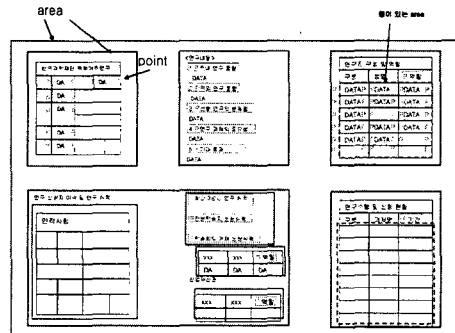


그림 4 정보 추출 및 변환 대상 예시

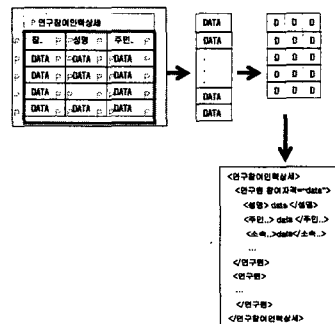


그림 5 정보 추출 및 결과 XML 문서 생성

또한 적용도메인의 특성상 많은 정보들이 [그림 4]와 같은 표 형태로 구성되어 있다. 그러나 표 이외의 정보에 대한 추출도 지원하도록 설계하였다. [그림 5]는 특정 표를 대상으로 하여 정보를 추출하는 과정을 설명하는 것으로써, 먼저 표 영역을 검색하여 단일 배열로 만들고, 이중 배열 처리 후에 결과 XML을 생성해낸다.

이 과정에서 [그림 6]과 같은 규칙을 사용한다. 많은 부분이 XSLT를 이용한 방법과 유사한데, 대상으로 한 정보의 특수성으로 인해 XSLT를 확장한 형태로 설계하였다. 그 특수성으로는 1)대상으로 하는 정보가 많은 부분 표, 리스트 등과 관련이 있고, 2)특정 영역에 있는 정보만을 추출해야하기 때문에 강화된 검색 기능(Point 개념)이 필요하고, 3)세밀한 검증(validation)과 함수를 필요로 한다는 점이다.

```

...
<area>
  <sPoint><![CDATA[<BOX TYPE=TABLE]]></sPoint>
  <ePoint type="joint"/>
  <apply action="TableToXml" columns="10"
includeHeader="false">
    <![CDATA[ResultOf()]]>
    <select pattern="RepeatedArea" breakCount="10">
      <cSText><![CDATA[<CELL]]></cSText>
      <cEText><![CDATA[</CELL]]></cEText>
    </select>
    <event name="AfterEveryRow">
      <apply action="AppendPreparedNode" path="." id="
구성원"/>
      <apply action="SetAttr" path="참여자격"
name="code"><![CDATA[<column0]]></apply>
      <apply action="AddText" path="성명"
"><![CDATA[<column1]]></apply>
      <apply action="AddText" path="주민등록번호"
"><![CDATA[<string.PickNum(<column2>)]></apply>
      <apply action="AddText" path="참여기간/시작일"
"><![CDATA[<date.SimForm(<column5>,"yyymmdd","y")]></appl
y>
      <apply action="AddText" path="참여율"
"><![CDATA[<string.PickNum(<column6>)]></apply>
      </event>
    </apply>
  </area>
...

```

그림 6 원시 XML 변환을 위한 규칙

4. XML 저장 방법

[그림 5]의 하단 부분에 묘사된 것과 같은 형태의 결과 XML 문서는 관계형 DB에 적합한 형태로 변형되어 데이터베이스에 저장된다.

```

<?xml version="1.0" encoding="EUC-KR"?>
<Database ip="203.250.200.93" port="1521" dbms="oracle"
user="radis" password="radis123" name="HWPXML">
<DefSection>
<Table name="tab01" type="single">
<Column name="fd_01_01" source="/자재구매/결재/담당" />
<Column name="fd_01_02" source="/자재구매/결재/검토" />
<Column name="fd_01_03" source="/자재구매/결재/결재" />
<Column name="fd_01_04" source="/자재구매/의뢰사업장" />
<Column name="fd_01_05" source="/자재구매/의뢰일자" />

```

```

</Table>
...
</DefSection>
<ApplySection overwrite="no">
  <Apply name="tab01" onError="stop" />
  <Apply name="tab02" onError="skip" />
</ApplySection>
</Database>

```

그림 7 관계형 DB에 저장하기 위한 저장 규칙

이 과정에서 데이터베이스 저장을 위한 규칙을 사용하는데, 규칙은 [그림 7]에 설명되어 있다.

5. 결론 및 향후연구

본 논문에서는 XML 기반 변환 기법을 이용하여 전자문서 형식에서 필요한 정보만을 추출하여 DB에 저장하는 방법에 대해 제안하였다. 특히 일반 사용자들이 널리 사용하는 HWP문서를 대상으로 하는 방법을 중심으로 언급하였으나, 일반적인 전자문서들(MS-Word, 기타)에서도 적용가능하다. 또한 저작도구의 특수 기능을 이용하지 않고 범용의 전자문서 저작도구 환경에서도 적용가능 하도록 설계/구현하였다. 대상 문서를 HWP를 대상으로 한정 할 경우, 스타일이나 누름틀 기능, HWP 고유 XML 저장 기법을 이용하여 보다 효과적인 처리가 가능할 것으로 보인다.

본 시스템은 전자 문서 원문 정보의 데이터화를 제공해야 하는 분야에서 유용하게 사용할 수 있다. 특히 연구과제관리시스템, 연구개발성과정보 시스템과 같은 분야에서는 과제, 논문, 지적재산권, 세미나, 인력, 예산에 이르기까지 방대한 규모의 데이터 구축이 필요하다. 지금까지는 데이터 구축에 걸리는 인력 및 시간의 소요가 상당할 뿐 아니라 이로 인해 해당 시스템에서 제공하는 데이터의 정밀도와 서비스 성능에 제약을 받아왔다. 이런 정보의 상당 부분은 이미 사업계획서, 과제요약서와 같은 전자 문서에 자세하게 기록되어 있어, 본 시스템을 접목함으로써 데이터 구축 오버헤드를 줄이고 보다 자세하고 정확한 정보를 구축하여 서비스의 질을 향상시킬 수 있다.

참고문헌

- [1] 고승규, 정병희, 손원성, 이경호, 임순범, 최원철, "HWP 문서와 EBKS 문서간의 변환 기법에 관한 연구", 한국멀티미디어학회 추계학술발표, 553-557, 2001.
- [2] W3C, "XSL Transformations (XSLT) Version 1.0", <http://www.w3.org/TR/xslt>, 1999.
- [3] 신동훈, 이경호, "XML 문서의 자동 변환을 위한 XSLT 스크립트 생성", 한국정보과학회 불 학술발표논문집, 31권, 1호, 160-162, 2004.
- [4] 이준승, 신동훈, 이경호, "XML 문서의 자동변환", 한국멀티미디어학회 추계학술발표대회논문집, 822-826, 2004.
- [5] 곽동규, 박호병, 유재우, "XML 스키마의 의미 구조 분석을 이용한 XML 문서의 변환", 한국정보과학회 추계학술발표 논문집, 32권 2호, 592-594, 2005.
- [6] 한글과컴퓨터, <http://www.hncxml.com/>