

From RTF to XML to L^AT_EX

Andre Dierker Arne Jans

Stephan Lehmke

QuinScape GmbH, Thomasstraße 1, 44135 Dortmund, Germany

{André.Dierker,Arne.Jans,Stephan.Lehmke}@QuinScape.de

<http://www.QuinScape.de>

Februar 26, 2005

Abstract

This paper shows how the widely used Rich Text Format (originally specified by Microsoft) can be processed to produce XML and how the resulting XML-file can be processed by L^AT_EX to produce print-quality PDF-files. To convert the RTF-files to a specific XML-format we use the open-source-tool Majix which can be found on Sourceforge.

We then take XMLEX to process the generated XML-file. We had to create an output that is as near as possible to Word's. An effort was made to reach this goal even with constructs such as lists, tabstops and especially tables.

Using XML as an intermediary format in typesetting RTF has the advantage that structural transformations are much easier based on XML even if the XML 'only' reproduces the RTF as faithfully as possible.

Filtering or transforming certain objects or attributes and even correcting typesetting errors can be done by appropriate transformations of the XML files.

1 Introduction

For a commercial project it was necessary to typeset larger documents from automatically generated XML data with embedded references to external RTF files. It was decided to first transform the RTF to XML and include that into the existing XML-structure, all together then being typeset by XMLEX.

For this, the open-source-tool Majix was extended to achieve the RTF-XML translation. For typesetting the resulting XML an appropriate implementation using XMLEX was created which will be uploaded to CTAN eventually, providing another open-source way of handling RTF with T_EX.

2 Results

We begin with a RTF-File. As you can see in figure 2 OpenOffice isn't able to handle recorded changes to the file correctly. For example in 'Betrag für 20042 Tsd EUR' (first row, second column) the '20042' is meant to be a year. Originally is was 2002 but the last '2' was deleted and replaced with '4'. However the '2' remained as old version in the RTF and is tagged as 'deleted'. In opposite to Majix OpenOffice doesn't recognise the responsible control word. The corresponding RTF-Code is shown in the follwing listing. We have selected a quite readable portion. When it comes to font management RTF isn't readable at all.

The RTF-Code looks like this

```

1 \s16\qj\sa40\widctlpar\adjustright \f18\fs16\lang1031\cgrid {\expnd0\expndtw-2 Anteil, der
2 aufgrund Artikel 9 des Verwaltungsabkommens vom 5.\~9.\~1957 i. d. F. vom 28.\~2.\~1991 zwi
3 schen Bund und L'e4ndern \fc
4 ber die Err
5 ichtung eines Wissenschaftsrats im Haushaltsjahr 1994 voraussichtlich entf\`e4llt.
6 \par } \pard
7 \plain \s18\qc\sa40\widctlpar\adjustright \f18\fs16\lang1031\cgrid {\dcbericht \fcber di
8 e Einnahmen und Ausgaben\line des Wissenschaftsrates
9 \par } \trow
10 \clvertalt\clbrdrt\brdrs\brdrw20 \clbrdrb\brdrs\brdrw20 \cltxlrb \cellx1021\clvertalt\cl
11 brdr\brdrs\brdrw20 \clbrdrb\brdrs\brdrw20 \clbrdrr\brdrs\brdrw20 \cltxlrb \cellx2042\clve

```

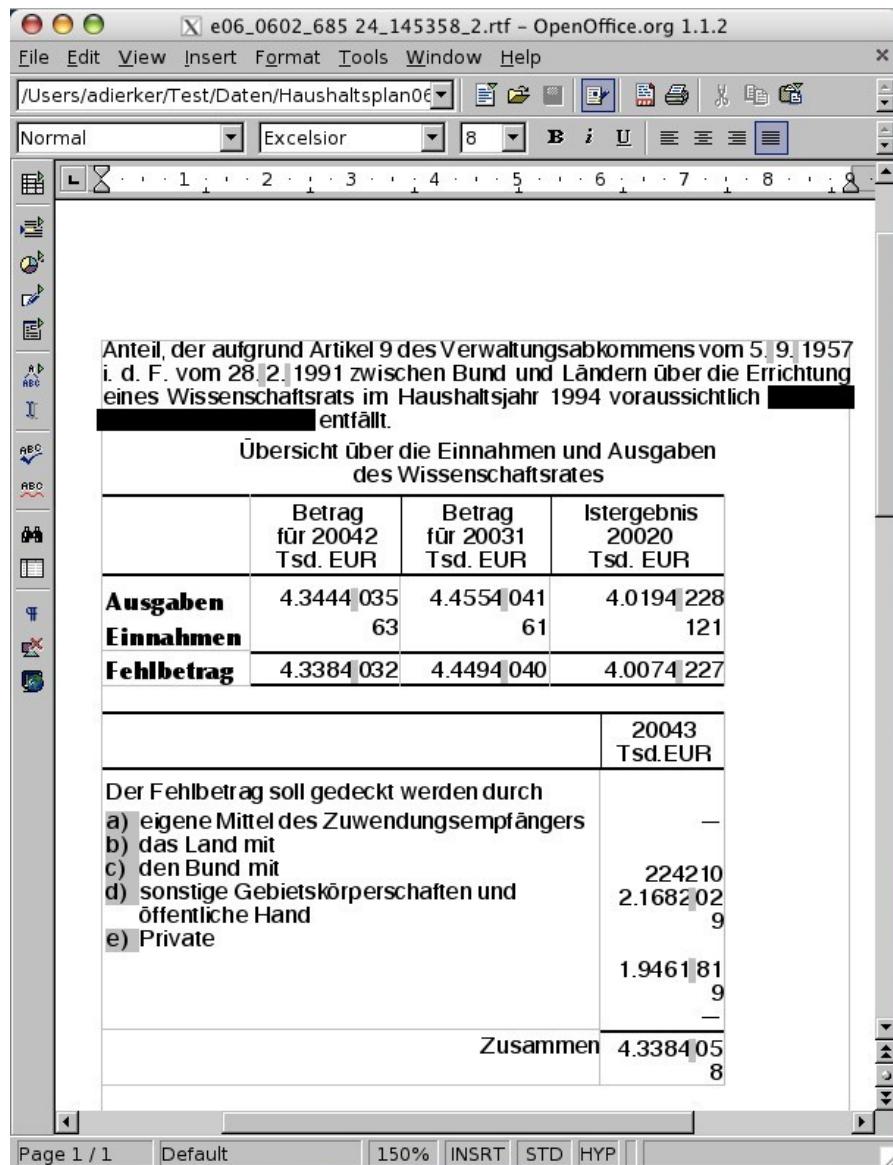


Figure 1: We first have a RTF file. Note the mixed up old and new version of the years and amounts. (Please ignore the font type. The used one wasn't installed and had to be substituted by OpenOffice)

```

12 rta|t\clbrdrt\brdrs\brdrw20 \clbrdrb\brdrs\brdrw20 \cltxlrb
13 \cellx3063\clvertalt\clbrdrt\brdrs\brdrw20 \clbrdrb\brdrs\brdrw20 \c
14 ltxlrb \cellx4253\pard\plain \qc\sb40\sa40\widctlpar\intbl\adjustright \f18\fs16\lang1031\
15 cgrid {\cell Betrag\line f\`fcr 200}{\revised\revauth1\revdttm1182249745
16 }{\dele
17 \revauthdel1\revdttmde1182249745 2}{\line Tsd. EUR\cell }{\pard \qc\f1\li-1\sb40\widctlpar
18 \intbl\adjustright {Betrag\line f\`fcr 200}{\revised\revauth1\revdttm1182249745 3}{\dele
19 \revauthdel1\revdttmde1182249745 1}{\line Tsd. EUR\cell
20 Istergebnis
21 \line 200}{\revised\revauth1\revdttm1182249745 2}{\dele\revauthdel1\revdttmde1182249745
22 0}{\line Tsd. EUR\cell }{\pard \widctlpar\intbl\adjustright {\row }\trowd \clvertalt\cltxlr
23 tb \cellx1021\clvertalt\cltxlrb \cellx2042
24 \clvertalt\
25 cltxlrb \cellx3063\clvertalt\cltxlrb \cellx4253\pard\plain \s1\sb80\sa40\keepn\widctlpar\
26 intbl\outlinelevel0\adjustright \b\fs20\lang1031\cgrid {\f18\fs16 Ausgaben}

```

After the conversion by Majix we get the XML-Code shown in the next listing. Majix did a great job in giving the data a meaningful structure.

The generated XML is much more readable

```

1 <par align="justified">
2   <tabdeflist>
3     <tabdef type="default" align="left" position="12.49mm"/>
4   </tabdeflist>
5   <parcontent>
6     Anteil, der aufgrund Artikel 9 des Verwaltungsabkommens vom 5. 9. 1957 i. d. F.
7     vom 28. 2. 1991 zwischen Bund und Ländern über die Errichtung eines
8       Wissenschaftsrats im Haushaltsjahr 1994 voraussichtlich entfällt.
9   </parcontent>
10  </par>
11  <par align="center">
12    <tabdeflist>
13      <tabdef type="default" align="left" position="12.49mm"/>
14    </tabdeflist>
15    <parcontent>
16      Übersicht über die Einnahmen und Ausgaben<linebreak/>des Wissenschaftsrates
17    </parcontent>
18  </par>
19  <table>
20    <tbody>
21      <tr>
22        <td width="18.0093mm" valign="top" border-top="0.0pt"
23          border-bottom="1.0pt" border-left="0.0pt">
24        </td>
25        <td width="18.0093mm" valign="top" border-top="0.0pt"
26          border-bottom="1.0pt" border-right="1.0pt">
27          <par align="center">
28            <tabdeflist>
29              <tabdef type="default" align="left" position="12.49mm"/>
30            </tabdeflist>
31            <parcontent>
32              Betrag<linebreak/>für 2004<linebreak/>Tsd. EUR
33            </parcontent>
34          </par>
35        </td>
36        <td width="18.0093mm" valign="top" border-top="0.0pt"
37          border-bottom="1.0pt">
38          <par align="center">
39            <tabdeflist>
40              <tabdef type="default" align="left" position="12.49mm"/>
41            </tabdeflist>
42            <parcontent>
43              Betrag<linebreak/>für 2003<linebreak/>Tsd. EUR
44            </parcontent>
45          </par>
46        </td>
47        ...
48      </tr>

```

The resulting PDF is shown in figure 2. As you can see there are some differences between the OpenOffice and the L^AT_EX version. These are the result of some filtering on the XML data demanded by our client.

3 Filtering

Having a well-formed XML one can easily filter the data. For example harmonizing the indentation of unordered lists can be done by deleting the necessary attributes in the XML and providing corresponding defaults.

Another class of transformations deals with the enrichment of semantics by converting visual markup to logical markup: one could search for special XML-constructs perhaps with specific attributes and/or contents and replace them with a ‘meaningful’ structure.

A good example for this are the so called VEGrids. These are special tables which always have the same layout. A typical VEGrid is shown in figure 3. The headline is always the same as is the tablehead. In the first column there are years, the bottom row is a sum. Usually VEGrids are already given in a XML-structure but sometimes it seems the person dealing with the case doesn’t use the right program to input the data but uses Word to create a RTF file with a VE-lookalike. After the XML conversion we get a noname-table as shown in

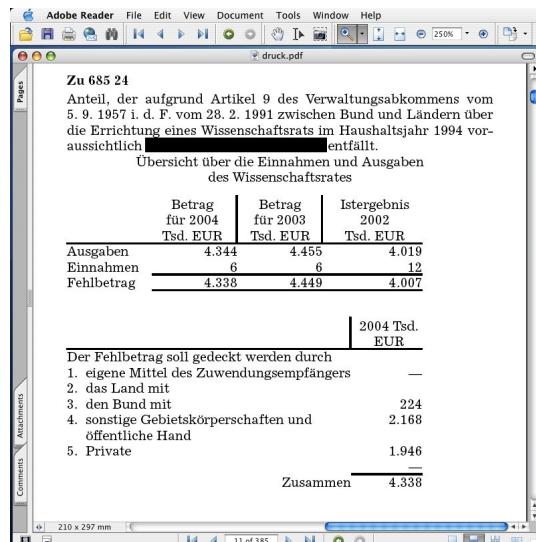


Figure 2: After processing with XMLTEX

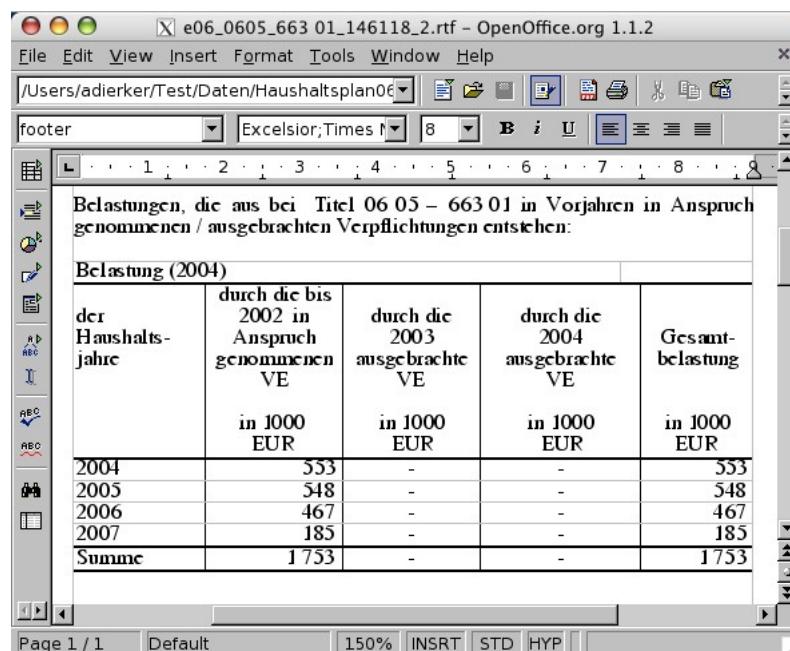


Figure 3: A special table called VEGrid in a RTF file

the next listing. Because of the constant layout we could identify these lookalikes and exchange them by ‘real’ VEGrids. This is shown in the last listing on the next page.

The generated XML (note the numerous omissions)

```

1 <table>
2   <tbody>
3     <tr border-bottom="0.5pt" border-left="0.5pt" border-right="0.5pt"
4       border-horizontal="0.5pt" border-vertical="0.5pt">
5       <td width="72.5135mm" valign="top" border-top="0.0pt"
6         border-left="0.0pt" colspan="4">
7         <par align="left">
8           <tabdeflist>
9             <tabdef type="default" align="left" position="12.49mm"/>
10            </tabdeflist>
11            <parcontent>Belastung (2004) </parcontent></par></td>
```

```

12   <td width="17.4978mm" valign="top" border-top="0.0pt"
13     border-right="0.0pt" colspan="2"></td>
14 </tr>
15 <tr border-top="0.5pt" border-bottom="0.5pt" border-left="0.5pt" ...>
16   <td width="17.5154mm" valign="top" border-top="1.0pt" ...>
17   ...
18   <parcontent>
19     <linebreak/>der<linebreak/>Haushalts-<linebreak/>jahre
20   </parcontent></par></td>
21 ...
22 </tr>
23 <tr border-top="0.5pt" border-bottom="0.5pt" border-left="0.5pt" ...>
24   <td width="17.5154mm" valign="top" border-top="1.0pt" ...>
25   ...
26   <parcontent>2004</parcontent></par></td>
27   <td width="18.1151mm" valign="top" border-top="1.0pt" ...>
28   ...
29   <parcontent>553 </parcontent></par></td>
30   <td width="18.1328mm" valign="top" border-top="1.0pt" ...>
31   ...
32   <parcontent>-</parcontent></par></td>
33 ...
34 </tr>
35 <tr border-top="0.5pt" border-bottom="0.5pt" border-left="0.5pt" ...>
36   <td width="17.5154mm" valign="top" border-top="1.0pt" ...>
37   ...
38   <parcontent>Summe</parcontent></par></td>
39   <td width="18.1151mm" valign="top" border-top="1.0pt" ...>
40   ...
41   <parcontent>1 753 </parcontent></par></td>
42 ...
43 </tr>
44 </tbody>
45 </table>

```

The enriched and filtered XML

```

1 <VEGRID>
2   <VEROW>
3     <VECOLUMN1>2005</VECOLUMN1>
4     <VECOLUMN2>553</VECOLUMN2>
5     <VECOLUMN3>null</VECOLUMN3>
6     <VECOLUMN4>null</VECOLUMN4>
7     <VECOLUMN5>553</VECOLUMN5>
8   </VEROW>
9   <VEROW>
10    <VECOLUMN1>2004</VECOLUMN1>
11    <VECOLUMN2>548</VECOLUMN2>
12    <VECOLUMN3>null</VECOLUMN3>
13    <VECOLUMN4>null</VECOLUMN4>
14    <VECOLUMN5>548</VECOLUMN5>
15  </VEROW>
16 ...
17  <VESUMROW>
18    <VECOLUMN1>Summe</VECOLUMN1>
19    <VECOLUMN2>1753</VECOLUMN2>
20    <VECOLUMN3>null</VECOLUMN3>
21    <VECOLUMN4>null</VECOLUMN4>
22    <VECOLUMN5>1753</VECOLUMN5>
23  </VESUMROW>
24 </VEGRID>

```

4 Further Development

Because there was no need up to now we ignore changes of font size and type. Besides that the management of colours wasn't implemented yet. Perhaps these features will be implemented in future.

As said before it is planned to release the code of the package on CTAN. The extended version of the RTF converter Majix is already available via SourceForge.