

Tagged PDF

Hans Hagen

1. Introduction

Occasionally users asked me if ConT_EXt can produce tagged pdf and the answer to that has been: I'll implement it when I need it. However, users tell me that publishers show an increasing demand for tagged pdf files, although one might wonder what for, except maybe for accessibility. Another reason for not having spent too much time on it before is that the specification was not that inviting.

At any rate, when I saw Ross Moore¹ presenting tagged math at TUG 2010, I decided to look up the spec once more and see if I could get into the mood to implement tagging. Before I started it was already clear that there were a couple of boundary conditions:

- Tagging should not put a burden on the user but users should be able to do the tagging themselves.
- Tagging should not slow down a run too much; this is no big deal as one can postpone tagging till the last run.
- Tagging should in no way interfere with typesetting, so no funny nodes should be injected.
- Tagging should not make the code look worse, neither the document source, nor the low level ConT_EXt code.

And of course implementing it should not take more than a few days' work, certainly not during an exceptionally hot summer.

You can 'google' for one of Ross's documents (like `DML_002-2009-1_12.pdf`) to see how a document source looks at his end using a special version of pdfT_EX. However, the version on my machine didn't support the primitives shown, so I could not see what was happening under the hood. Unfortunately it is quite hard to find a properly tagged document so we have only the reference manual as a starting point. As the pdfT_EX approach didn't look that pleasing anyway, I just started from scratch.

Tags can help Acrobat Reader when it is asked to read out the text aloud. But you cannot browse the structure in the no-cost version of Acrobat and, as not all users have the professional version of Acrobat, the fact that a document has structure can go unnoticed. Add to that the fact that the overhead in terms of bytes is quite large due to a significant increase in generated objects and you will understand why this feature is not enabled by default.

2. Implementation

So, what does tagging boil down to? We can best look at how tagged information is shown in Acrobat. Figure 1 shows the content tree that has been added (automatically) to a document while figure 2 shows a different view.

¹ He is often exploring the boundaries of pdf, Unicode and evolving techniques related to math publishing, so you'd best not miss his presentations when you are around.

contextgroup > context meeting 2011



Figure 1: A tag list in Acrobat.

In order to get that far, we have to do the following:

- Carry information with {typeset} text.
- Analyse this information when shipping out pages.
- Add a structure tree to the page.
- Add relevant information to the document.

That first activity is rather independent of the other three and we can use that information for other purposes as well, like identifying where we are in the document. We carry the information around using attributes. The last three activities took a bit of experimenting mostly using the “Example of Logical Structure” from the pdf standard 32000-1:2008.

This resulted in a tagging framework that uses explicit tags, meaning the user is responsible for the tagging:

```
\setupstructure[state=start,method=none]  
  
\starttext  
  
\startelement[document]
```

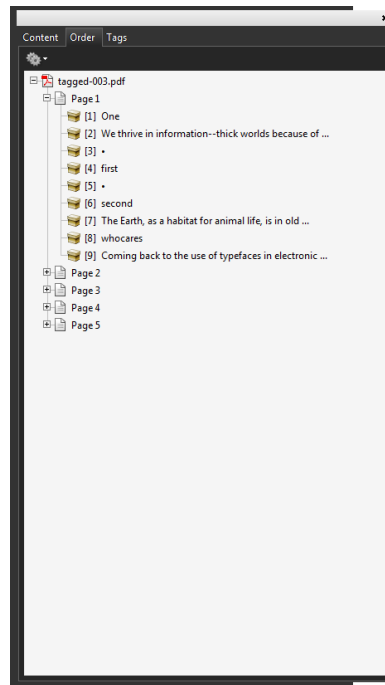
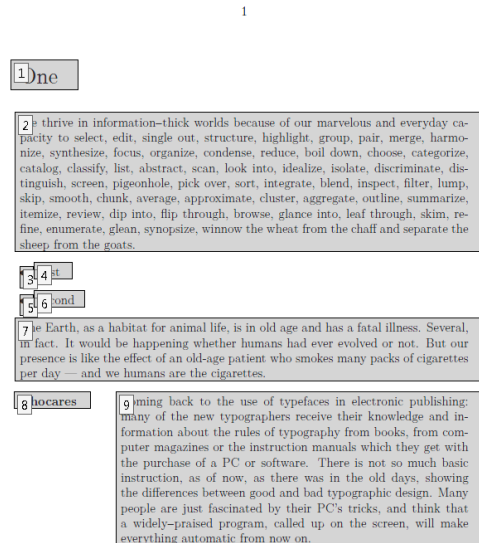


Figure 2: Acrobat showing the tag order.



```

\startelement[chapter]
  \startelement[p] \input davis \stopelement \par
\stopelement

\startelement[chapter]
  \startelement[p] \input zapf \stopelement \par
  \startelement[whatever]
    \startelement[p] \input tuftte \stopelement \par
    \startelement[p] \input knuth \stopelement \par
  \stopelement
\stopelement

\startelement[chapter]
  oeps
  \startelement[p] \input ward \stopelement \par
\stopelement

\stopelement

\stoptext

```

Since this is not much fun, we also provide an automated variant. In the previous example we explicitly turned off automated tagging by setting `method` to `none`. By default it has the value `auto`.

contextgroup > context meeting 2011

```
\setupstructure[state=start] % default is method=auto
\definedescription[whatever]
\starttext
\startfrontmatter
  \startchapter[title=One]
    \startparagraph \input tufte \stopparagraph
    \startitemize
      \startitem first \stopitem
      \startitem second \stopitem
    \stopitemize
    \startparagraph \input ward \stopparagraph
    \startwhatever {Herman Zapf} \input zapf \stopwhatever
  \stopchapter
\stopfrontmatter
\startbodymatter
.....
```

If you use commands like `\chapter` you will not get the desired results. Of course these can be supported but there is no real reason for it, as in MkIV we advise using the start-stop variant.

It will be clear that this kind of automated tagging brings with it a couple of extra commands deep down in ConT_EXt and there (of course) we use symbolic names for tags, so that one can overload the built-in mapping.

```
\setuptaglabeltext[en][document=text]
```

As with other features inspired by viewer functionality, the implementation of tagging is independent of the backend. For instance, we can tag a document and access the tagging information at the T_EX end. The backend driver code maps tags to relevant pdf constructs. First of all, we just map the tags used at the ConT_EXt end onto themselves. But, as validators expect certain names, we use the pdf rolemap feature to map them to [less interesting] names. The next list shows the currently used internal names, with the pdf ones between parentheses.

combination (Span), combinationcaption (Span), combinationcontent (Span), combinationpair (Span), construct (Span), delimited (Quote), delimitedblock (BlockQuote), description (Div), descriptioncontent (Div), descriptionsymbol (Span), descriptiontag (Div), division (Div), document (Div), float (Div), floatcaption (Caption), floatcontent (P), floatlabel (Span), floatnumber (Span), floattext (Span), formula (Div), formulacaption (Span), formulacontent (P), formulalabel (Span), formulanumber (Span), formulaset (Div), highlight (Span), ignore (Span), image (P), item (Li), itemcontent (LBody), itemgroup (L), itemtag (Lbl), label (Span), line (Code), lines

[Code], link [Link], list [TOC], listcontent [P], listdata [P], listitem [TOCI], listpage [Reference], listtag [Lbl], maction [Span], margintext [Span], margintextblock [Span], math [Div], merror [Span], metadata [Div], metavariable [Span], mfenced [Span], mfrac [Span], mi [Span], mid [Span], mn [Span], mo [Span], mover [Span], mpgraphic [P], mroot [Span], mrow [Span], ms [Span], msqrt [Span], msub [Span], msubsup [Span], msup [Span], mtable [Table], mtd [TD], mtext [Span], mtr [TR], munder [Span], munderover [Span], number [Span], p [P], paragraph [P], register [Div], registerentries [Div], registerentry [Span], registerpage [Span], registerpagerange [Span], registerpages [Span], registersection [Div], registersee [Span], registertag [Span], section [Sect], sectioncontent [Div], sectionnumber [H], sectiontitle [H], sorting [Span], sub [Span], subformula [Div], subsentence [Span], subsup [Span], sup [Span], synonym [Span], table [Table], tablecell [TD], tablerow [TR], tabulate [Table], tabulatecell [TD], tabulaterow [TR], verbatim [Code], verbatimblock [Code], verbatimline [Code], verbatimlines [Code].

So, the internal ones show up in the tag trees as shown in the examples but applications might use the rolemap which normally has less detail.

Since we keep track of where we are, we can also use that information for making decisions.

```
\doifinelementelse{structure:section}           {yes} {no}
\doifinelementelse{structure:chapter}          {yes} {no}
\doifinelementelse{division:*-structure:chapter} {yes} {no}
\doifinelementelse{division:*-structure:*}     {yes} {no}
```

As shown, you can use `*` as a wildcard. The elements are separated by `-`. If you don't know what tags are used, you can always enable the tag related tracker:

```
\enabletrackers[structure.tags]
```

This tracker reports the identified element chains to the console and log.

3. Special care

Of course there are a few complications. First of all the tagging model sort of contradicts the concept of a nicely typeset document where structure and outcome are not always related. Most \TeX users are aware of the fact that \TeX does not have space characters and does a great job on kerning and hyphenation. The tagging machinery on the other hand uses a rather dumb model of strings separated by spaces.² But we can trick \TeX into providing the right information to the backend so that words get nicely separated. The non-optimized function that does this looks as follows:

² The search engine on the other hand is rather clever on recognizing words.

contextgroup > context meeting 2011

```
function injectspaces(head)
  local p
  for n in node.traverse(head) do
    local id = n.id
    if id == node.id("glue") then
      if p and p.id == node.id("glyph") then
        local g = node.copy(p)
        local s = node.copy(n.spec)
        g.char, n.spec = 32, s
        p.next, g.prev = g, p
        g.next, n.prev = n, g
        s.width = s.width - g.width
      end
      elseif id == node.id("hlist") or id == node.id("vlist") then
        injectspaces(n.list,attribute)
      end
      p = n
    end
  end
end
```

Here we squeeze in a space (given that it is in the font, which it normally is when you use ConT_EXt) and make a compensation in the glue. Given that your page sits in box 255, you can do this just before shipping the page out:

```
injectspaces(tex.box[255].list)
```

Then there are the so-called suspects: things on the page that are not related to structure at all. One is supposed to tag these in a special way to prevent the built-in reading equipment from getting confused. So far we could get around them simply because they don't get tagged at all and therefore are not seen anyway. This might well be enough of a precaution.

Of course we need to deal with mathematics. Fortunately the presentation MathML model is rather close to T_EX and so we can map onto that. After all we don't need to care too much about back-mapping here. The currently present code is rather experimental and might get extended or thrown out in favour of inline MathML. Figure 3 demonstrates that a first approach does not even look that bad. In future versions we might deal with table-like math constructs, like matrices.

This is a typical case where more energy has to be spent on driving the voice of Acrobat but I will do that when we find a good reason.

As mentioned, it will take a while before all relevant constructs in ConT_EXt support tagging, but support is already quite complete. Some screen dumps are included as examples at the end.

contextgroup > context meeting 2011

1 chapter

test oeps test **whow** test

test

`\whatever[goes]{here}`

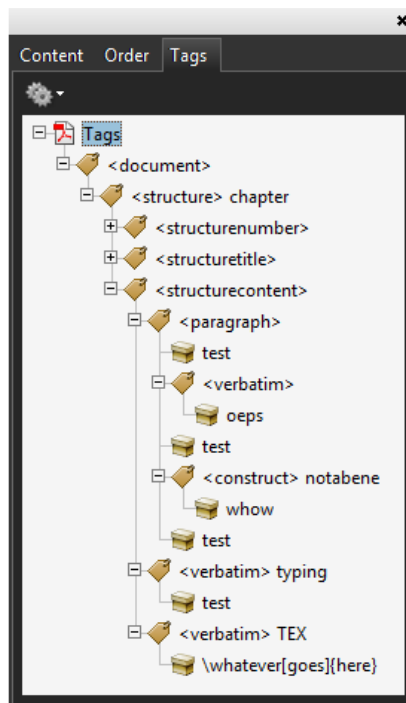


Figure 4: Verbatim, including dedicated instances.

test 11	test 12
test 21	test 22
test 33	

test Coming back

new typograp

typography fi

which they ge

sic instructio

between good

by their PC's

the screen, w

test Coming back

new typograp

typography fi

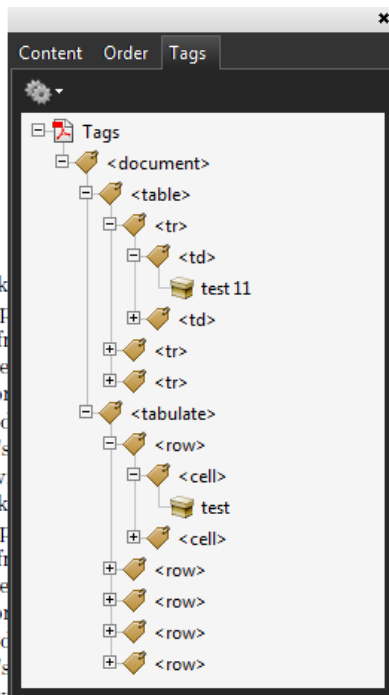
which they ge

sic instructio

between good

by their PC's

the screen, w



ublishing: many of the

ation about the rules of

the instruction manuals

There is not so much ba

showing the differences

ople are just fascinated

l program, called up on

on.

ublishing: many of the

ation about the rules of

the instruction manuals

There is not so much ba

showing the differences

ople are just fascinated

l program, called up on

on.

Figure 5: Natural tables as well as the tabulate mechanism is supported.

Contents

1	One	
1.1	alpha	
1.2	beta	
1.3	gamma	
1.4	delta	

		2
		2
		2
		2
		2

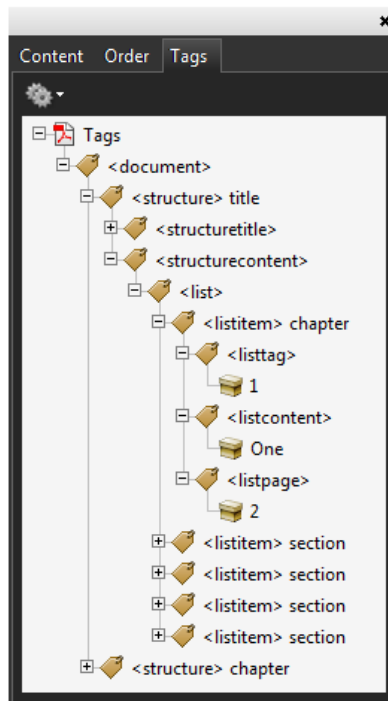


Figure 6: Tables of content with specific entries tagged.

Index

o		
one	1, 2	
t		
two	1, 2	

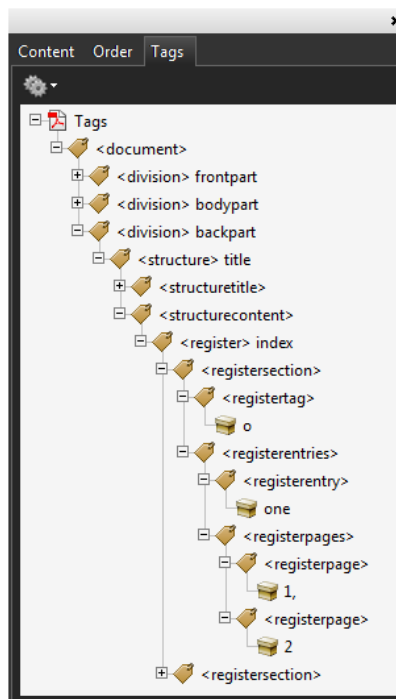
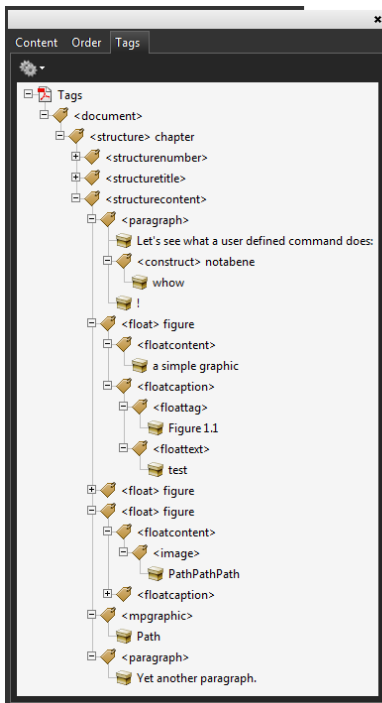


Figure 7: A detailed view of registered is provided.



1 chapter

Let's see what a user defined command does: **whow!**

a simple graphic

Figure 1.1 test

a simple graphic

Figure 1.2 test

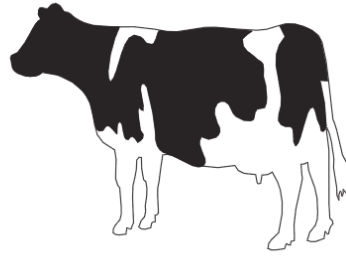


Figure 1.3 test

Yet another paragraph.

Figure 8: Floats tags end up in text stream. Watch the user defined construct.

sheep from the

- first!
- second

The Earth, as
in fact. It wou
presence is like
per day — and

whocares

1 test

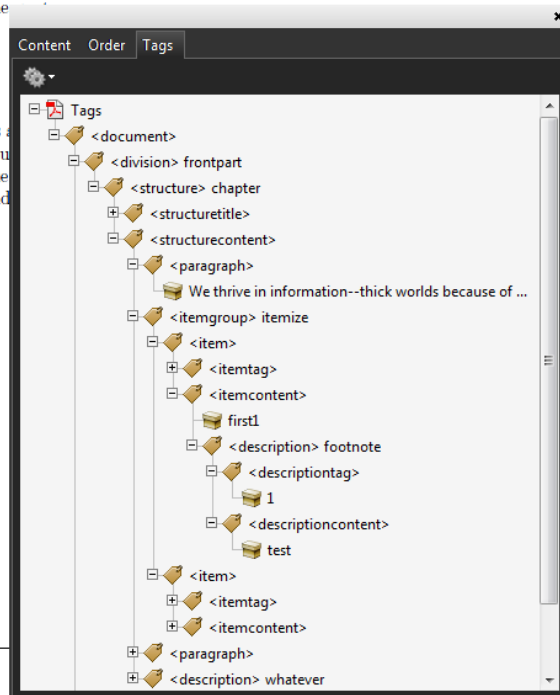


Figure 9: Footnotes are shown at the place in the input (flow).