XML Forth Glossary

Here is a rough description of the DocBook Forth Glossary extension I have been considering. This is subject to review.

 $\langle \text{glossary} \rangle := \langle \text{wordset} [\langle \text{id} \rangle] [\langle \text{number} \rangle] [\langle \text{name} \rangle] [\langle \text{core} \rangle] [\langle \text{system} \rangle] > \\ (\langle \text{worddef} \rangle^* | \langle \text{rationale} \rangle^* | \langle \text{testing} \rangle^* | \langle \text{implement} \rangle^*) \\ \langle \text{/wordset} \rangle \\ \langle \text{id} \rangle := \text{id} = " \langle \text{label} \rangle " \\ \langle \text{number} \rangle := \text{number} = " \langle \text{sec-no} \rangle " \\ \langle \text{name} \rangle := \text{name} = " \langle \text{string} \rangle " \\ \langle \text{core} \rangle := \text{core} = " \langle \text{boolean} \rangle " \\ \langle \text{system} \rangle := \text{system} = " \langle \text{label} \rangle " \\ \langle \text{boolean} \rangle := \text{true} | \text{false} \\ \langle \text{sec-no} \rangle := (\langle \text{upper} \rangle | \langle \text{digit} \rangle) (\langle \text{digit} \rangle | .)^* \\ \langle \text{label} \rangle := (\langle \text{upper} \rangle | \langle \text{lower} \rangle | \langle \text{digit} \rangle | ! | ' | (|) | * | + | , | - | . | / | : | ; | = | @ | [|])^+$

The $\langle wordset \rangle$ tag defines a number of defaults as well as encapsulating all of the word definitions $(\langle worddef \rangle)$, rationale $(\langle rationale \rangle)$, test cases $(\langle testing \rangle)$ or implementation $(\langle implement \rangle)$ for the word set.

id	Label for cross referencing. Note that characters that will cause a problem in one of the intended rendering notations have been forbidden.
number	The location of the word set within the document. The current section number will be by default.
name	The name of the word set (normally a single upper case word)
core	Indicate that the word set describes the core word set (true, the default) or the extended word set (false).
system	The name of the Forth system (or standard) on which this wordset is defined.

<wordset id="memory" number="14.6.1" name="MEMORY" system="Forth2012">

1 Word definition

 $\begin{array}{ll} \langle worddef \rangle := & < worddef \langle id \rangle [\langle number \rangle] \langle name \rangle [\langle wordset \rangle] [\langle english \rangle] > \\ & & \langle description \rangle \\ & & [\langle word-rat \rangle] \\ & & [\langle word-test \rangle] \\ & & [\langle word-te$

 $\langle english \rangle := english = "\langle string \rangle "$

A glossary definition is enclosed in <worddef> tags. The tag takes four optional arguments:

id The label used to identify the word in the cross referencing system.

number The word's number (and sub-number) within the document. The next sequential number will be used by default. This number is appended to the number attribute of the enclosing $\langle wordset \rangle$ (if given).

name The word name, normally in upper case with special characters escaped.

wordset The word set this definition appears in. The wordset attribute of the enclosing $\langle wordset \rangle$ will be used by default.

english The pronunciation for the word.

<worddef id="core:DOES" number="1250" name="DOES>" english="does">

Question:

In the $L^{A}T_{E}X$, *id* is optional, the default label is made by separating the *id* attribute of the enclosing *<wordset>* and *name* attributes with a colon. The *id* attribute is only provided to override the *name* attribute when that uses a character that can not be types easily. Thus, if the *id* attribute of the enclosing *<wordset>* is "core" we would have:

```
\begin{array}{ll} \textit{name}{=}"IF" & \Rightarrow \textit{core:}IF \\ \textit{name}{=}"\>" \textit{id}{=}"more" & \Rightarrow \textit{core:more} \end{array}
```

Should this behaviour be reproduced in the XSLT?. Or, to put it another way, should the *id* attribute be optional?

1.1 Description

The main description of the operation of the word. The description can be broken down into different $\langle section \rangle s$.

```
 \langle \text{description} \rangle := \langle \text{description} \rangle \\ ( [\langle \text{interpret} \rangle] [\langle \text{compile} \rangle] [\langle \text{init} \rangle] [\langle \text{execute} \rangle] [\langle \text{runtime} \rangle] \langle \text{para-text} \rangle^* ) | \\ [\langle \text{sec-text} \rangle] \\ \langle \text{description} \rangle \\ \langle \text{interpret} \rangle := \langle \text{interpret} \rangle \langle \text{sec-text} \rangle \langle /\text{interpret} \rangle \\ \langle \text{compile} \rangle := \langle \text{compile} \rangle \langle \text{sec-text} \rangle \langle /\text{compile} \rangle \\ \langle \text{init} \rangle := \langle \text{init} \rangle \langle \text{sec-text} \rangle \langle /\text{init} \rangle \\ \langle \text{execute} \rangle := \langle \text{execute} [\langle \text{type} \rangle] \rangle \langle \text{sec-text} \rangle \langle /\text{runtime} \rangle \\ \langle \text{runtime} \rangle := \langle \text{runtime} [\langle \text{type} \rangle] \rangle \langle \text{sec-text} \rangle \langle /\text{runtime} \rangle \\ \langle \text{sec-text} \rangle := \langle \text{para-text} \rangle^*
```

A <description> should include a general description of the word which may be broken down into further sections.

1.2 Rationale

If the word requires additional discussion, that is not part of its definition, this discussion should be included in the rationale, this should include the reasoning for including the word in the word set. The rationale for a word can be provided in two places. It can be included with the word definition (inside a <worddef>) or the rationale for all of the words in a word set can be collected together (inside a <worddet>).

 $\langle word\text{-}rat \rangle := \langle rationale \rangle \langle sec\text{-}text \rangle \langle rationale \rangle$

 $\langle att-def \rangle := \langle id \rangle [\langle number \rangle] \langle name \rangle [\langle wordset \rangle]$

 $\langle rationale \rangle := \langle rationale \langle att-def \rangle > \langle sec-text \rangle \langle rationale \rangle$

When defined within a word definition (<worddef>) the \langle word-rat \rangle element is used. This requires no additional information as the word definition provides a context for the rationale.

When defined within a separate collection of rationale's the $\langle rationale \rangle$ element is used. The definition context is not available, thus a number of additional attributes are used to provide this context. The attributes have the same function as for the $\langle worddef \rangle$.

Question:

In the $L^{A}T_{E}X$, the label for a rationale is constructed by adding the text "rat:" to the start of the word's label. Thus if the word's label is "core:IF" the label for the rationale would be "rat:core:IF". The label is constructed in the same way as for (worddef), thus the *id* attribute has the same function.

Therefore the same question applies, should this behaviour be replicated in XSLT?

The element $\langle rref \rangle$ is used to cross-reference the rationale section of a word's definition.

When processing the file the first pass extracts the $\langle word-rat \rangle$ elements, collating them in a $\langle wordset \rangle$ for inclusion in the appendix. Subsequent processing ignores the $\langle word-rat \rangle$ element.

1.3 Testing

The unit testing for a word should be given in the testing section. As with $\langle rationale \rangle$, this can be included with the word's definition ($\langle word\text{-test} \rangle$) or collated into a list of test cases ($\langle testing \rangle$).

 $\langle word\text{-}test \rangle := \langle testing \rangle (\langle sec\text{-}text \rangle | \langle test \rangle)^* \langle testing \rangle$

 $\langle \text{testing} \rangle := \langle \text{testing} \langle \text{att-def} \rangle > (\langle \text{sec-text} \rangle | \langle \text{test} \rangle)^* \langle \text{testing} \rangle$

The "test:" prefix is added to the word's label to generate the label for the reference implementation. The $\langle tref \rangle$ element is used to reference an implementation.

In addition to normal text ($\langle sec-text \rangle$) the $\langle test \rangle$ tag can be used to describe an individual test case.

$$\begin{array}{ll} \langle test \rangle := < test \ [\ type="(\ X \ | \ R \)*" \]> \\ \langle inline-source \rangle \\ < post> \langle inline-source \rangle < /post> \\ < /test> \end{array}$$

For example,

```
<test type="XX">0 DUP<post>0 0</post></test>
```

describes the following test:

```
T{ 0 DUP -> 0 0 XX}T
```

The type argument is used to provide the type of the parameters required for the floating-point test harness.

1.4 Implement

Provides a reference implementation. As with < rationale >, this can be included with the word's definition ($\langle word\text{-}imp \rangle$) or collated into a list of implementations ($\langle implement \rangle$).

 $\langle word\text{-}imp \rangle := \langle implement \rangle \langle multi-line\text{-}source \rangle \langle /implement \rangle$

 $\langle implement \rangle := \langle implement \langle att-def \rangle > \langle mulit-line-source \rangle </implement>$

The "imp:" prefix is added to the word's label to generate the label for the reference implementation. The $\langle iref \rangle$ element is used to reference an implementation.

1.5 Cross referencing

Cross references to other words should be included in the $\langle see \rangle$ section. This consists of a collection of $\langle wref \rangle$, $\langle rref \rangle$, $\langle tref \rangle$, $\langle iref \rangle$ and $\langle xref \rangle$ tags.

$$\langle see \rangle := \langle see \rangle \langle cross-ref \rangle^* \langle see \rangle \\ \langle cross-ref \rangle := \langle wref \rangle | \langle rref \rangle | \langle tref \rangle | \langle iref \rangle | \langle xref \rangle \\ \langle att-ref \rangle := word="\langle string \rangle" [\langle wordset \rangle] \\ \langle wref \rangle := \langle wref \langle att-ref \rangle / \rangle \\ \langle rref \rangle := \langle rref \langle att-ref \rangle / \rangle \\ \langle tref \rangle := \langle tref \langle att-ref \rangle / \rangle \\ \langle iref \rangle := \langle iref \langle att-ref \rangle / \rangle \\ \langle xref \rangle := \langle xref | abel = "\langle string \rangle" > \langle text \rangle \langle xref \rangle$$

The cross referencing functions <wref>, <rref>, <tref> and <iref> all follow the same pattern. The word attribute is the label of the word being referenced. They look in the current word set for the word, if not found they will look for the word in the core word set. The optional wordset attribute can be used to override the current word set context. When the word is found it will output the reference number and the name, making them both link to the definition being referenced. Thus:

<wref word="Sq" />

will look in the current word set context for the word with the id "Sq", looking for the label " $\langle wordset \rangle$:Sq". If this is not found, it will look in the core word set (the label "core:Sq"). It will then output the following hyper link:

6.1.2165 S"

The function <rref> references the rationale for the word (prefixes the label with "rat:") effectively providing a separate label name space. The <tref> functions references the testing ("test:") name space while <iref> references the reference implementation ("imp:") name space.

Question:

This describes the operation of the existing $E^{A}T_{E}X$ macros. The name space could be passed to the generate cross reference function:

 $\langle cross-ref \rangle := \langle ref \langle att-ref \rangle [\langle namespace="(rat | test | imp) " \rangle] /> | \langle xref \rangle$

where the label is made up from " $\langle namespace \rangle$: $\langle wordset \rangle$: $\langle word \rangle$ ".

An alternative would be to insist the user provides a full label.

Finally, the $\langle xref \rangle$ tag allows for the cross referencing between sections of the document. If the section identified by the label attribute can be found, a full reference to the section is produced, complete with a link to the section. However, if the section cannot be found the $\langle text \rangle$ is used.

Question:

The *<xref>* tag is already used by DocBook so we may need to use a different tag for this function. Suggestions as to the name?

2 Block / Paragraph text

The (sec-text) element consists of any number of block or paragraph level elements described by the (para-text) element.

 $\langle \text{para-text} \rangle := \langle \text{stack} \rangle | \langle \text{note} \rangle | \langle \text{item} \rangle | \langle \text{code} \rangle | \langle \text{DBpara} \rangle$

 $\langle type \rangle := type = "\langle string \rangle "$

 $\langle note \rangle := \langle note | \langle type \rangle | \rangle \langle text \rangle \langle note \rangle$

 $\langle item \rangle := \langle item \langle type \rangle \rangle \langle text \rangle \langle item \rangle$

 $\langle DBpara \rangle := \dots DocBook \ paragraph \ level \ elements \ \dots$

Some sections take an optional type attribute. This is normally displayed in the label for that section of the definition, with the exception of the *<item>* section where the type is the whole of the label.

The list of paragraph level DocBook elements includes the *<para>* (paragraph) element.

2.1 Stack description

Many sections will require a stack description. The *<stack>* element is considered a paragraph (block) level element:

 $\begin{array}{l} \langle stack \rangle := \langle stack \ [\ type="\langle upper \rangle" \] > \\ & [\ \langle pre \rangle \ \langle stack\text{-}item \rangle^* \ \langle /pre \rangle \] \\ & [\ \langle post \rangle \ \langle stack\text{-}item \rangle^* \ \langle /post \rangle \] \\ & \langle /stack \rangle \end{array} \\ \\ \langle stack\text{-}item \rangle := \ \langle char \rangle^* \ | \ I \ | \ _ \langle char \rangle \ | \ _ \{ \ \langle char \rangle^* \ \} \ | \ * \ | \ \langle parsed\text{-}text \rangle \\ \\ \langle parsed\text{-}text \rangle := \ \langle \ (char \ | \ chars \ | \ space \ | \ spaces \ | \ quote \ | \ paren \ | \ eol \) \ / > \end{array}$

A stack description gives both the before $(\langle pre \rangle)$ and after $(\langle post \rangle)$ items. The type attribute is used to indicate the which stack the stack description is being being applied. (I.e., type="R" for the return stack).

The $\langle char \rangle$ element of the $\langle stack-item \rangle$ is rather misleading as a number of characters are treated as "special". The special meaning has been taken from the IAT_{EX} macros:

x_{i}	x subscript i	x_i
x_i	x subscript i	x_i
x*i	x multiplied by i	x * i
n u	$n ext{ or } u$	n u
"name <spa< td=""><td>ace/>"</td><td></td></spa<>	ace/>"	
-	the syntax element $name$ followed by a space	" $name \langle space \rangle$ "

The intention was to make the stack item appear as similar to the output text as possible.

Table 2.1 of the '94 document defines the seven syntactic elements (in $\langle parsed-text \rangle$) that can be used when parsing text. Defining these as empty elements allows there use in this way.

Question:

Should we consider using XML for the stack notation?

If so, then what notation,	something along the lines of:
x _i	$\Rightarrow x_i$
x <times></times> y	$\Rightarrow x * y$
n∣u	$\Rightarrow n u$
n❘u	$\Rightarrow n u$

3 Inline text

A number of additional functions are provided for text level. These functions can be used within any paragraph / block level element.

This is normal text with the following additional in-line tags:

 $\begin{array}{l} \langle text \rangle := \langle DBinline \rangle \mid \langle cross\text{-ref} \rangle \mid \langle word \rangle \mid \langle c \rangle \\ \mid & < \text{param} > \langle stack \; item \rangle^* < / \text{param} > \\ \mid & < \text{arg} > \langle string \rangle < / \text{arg} > \\ \langle word \rangle := < \text{word} \; \langle att\text{-ref} \rangle \; / > \end{array}$

5

The following additional formatting tags are provided:

$\langle DBinline \rangle$	DocBook inline elements, including <emphasis></emphasis> .		
$\langle cross-ref \rangle$	The cross reference functions ($\langle wref \rangle$, $\langle rref \rangle$, $\langle tref \rangle$, $\langle iref \rangle$, $\langle xref \rangle$).		
<word></word>	Represents a forth word. The word and wordset attributes work in the same way as for the cross reference tags. This will only output the name of the referenced word. If this is the current word, the name is output in bold, otherwise the word is a link to the word's definition. $<$ word word="IF"/> \Rightarrow IF $<$ word word="Sq" wordset="file" /> \Rightarrow S" $<$ word word="DOES" /> \Rightarrow DOES>		
$\langle c angle$	Represents program code fragment within the body of the text.		
<param/>	Represents a stack item within the the $c-add len$	$e \text{ text.} \\ \Rightarrow c \text{-} a ddr \ len$	
<arg></arg>	Represents an argument <arg>text</arg>	$\Rightarrow \langle text \rangle$	

4 Source code

Source code is output in one of two environments: the paragraph level < code > tag; the inline text level <c > tag.

 $\langle code \rangle := \langle code \rangle \langle multi-line-source \rangle \langle /code \rangle$ $\langle c \rangle := \langle c \rangle \langle inline-source \rangle \langle /c \rangle$

 $\langle multi-line-source \rangle := (\langle para-text \rangle | \langle inline-source \rangle)^*$

 $\langle \text{inline-source} \rangle := (\langle \text{char} \rangle | \langle \text{word} \rangle)^*$

In both environments the white space is preserved. All Forth words are cross referenced back to their definition using the <word> tag.

Question:

We could dispense with the *<word>* tag and assume every sequence of non-white spaces is a label for a word. Any label not found in the current word set or the core word set can be output directly.

The the current, rather user hostile, example:

<word word=":" /> TEST <word word="POSTPONE" /> <word word="Sq" wordset="file" />
... <word word=";" />

would become the considerable more friendly:

: TEST POSTPONE file:Sq \ldots ;

Note the use of the word's label rather than the word name, i.e., Sq and not S''. Also note the use of the word set prefix.

While this would be considerable more user friendly, I have no idea how to go about implementing this is XSLT.