

A complete schema definition language for the Text Encoding Initiative

Lou Burnard and Sebastian Rahtz

XML London, June 16th 2013



Reminder: what is the TEI?

A 25 year old project to define Guidelines for text encoding:

- mainly targetted at digital editions of existing texts
- covers manuscripts, dictionaries, transcribed text, spoken corpora, and facsimiles, as well as simple books
- governed by an international membership consortium
- defines a very rich language, with about 550 elements managed in 22 modules and an infrastructure of model and attributes classes
- Specialist vocabularies such as XInclude, MathML and SVG are used where appropriate.

<http://www.tei-c.org/>

The domain of the TEI

AN ANSWER TO WILLIAMS'S ALABLAST! his MOTIVES.

By ROGER FENTON F. Grayes Inns.



Imprinted W.A.F. ya

Agreed upon and Published at the Month of September, 1755.



HIST. 6. 19
The Treasurers Almanack,
for a brief View of the present state
OR
The Money-Master:
of the **Money-Master:**

Wherein we show the Loss of so much Money by Interest, the Loss of so much Money by Long...

Being most...

Days	Interest	Factor	A Table of Good Money
1	3/4	12	12
2	1 1/2	24	24
3	2 1/4	36	36
4	3 1/4	48	48
5	4 1/2	60	60
6	5 1/2	72	72
7	6 1/2	84	84
8	7 1/2	96	96
9	8 1/2	108	108
10	9 1/2	120	120
11	10 1/2	132	132
12	11 1/2	144	144
13	12 1/2	156	156
14	13 1/2	168	168
15	14 1/2	180	180
16	15 1/2	192	192
17	16 1/2	204	204
18	17 1/2	216	216
19	18 1/2	228	228
20	19 1/2	240	240
21	20 1/2	252	252
22	21 1/2	264	264
23	22 1/2	276	276
24	23 1/2	288	288
25	24 1/2	300	300
26	25 1/2	312	312
27	26 1/2	324	324
28	27 1/2	336	336
29	28 1/2	348	348
30	29 1/2	360	360
31	30 1/2	372	372
32	31 1/2	384	384
33	32 1/2	396	396
34	33 1/2	408	408
35	34 1/2	420	420
36	35 1/2	432	432
37	36 1/2	444	444
38	37 1/2	456	456
39	38 1/2	468	468
40	39 1/2	480	480
41	40 1/2	492	492
42	41 1/2	504	504
43	42 1/2	516	516
44	43 1/2	528	528
45	44 1/2	540	540
46	45 1/2	552	552
47	46 1/2	564	564
48	47 1/2	576	576
49	48 1/2	588	588
50	49 1/2	600	600
51	50 1/2	612	612
52	51 1/2	624	624
53	52 1/2	636	636
54	53 1/2	648	648
55	54 1/2	660	660
56	55 1/2	672	672
57	56 1/2	684	684
58	57 1/2	696	696
59	58 1/2	708	708
60	59 1/2	720	720
61	60 1/2	732	732
62	61 1/2	744	744
63	62 1/2	756	756
64	63 1/2	768	768
65	64 1/2	780	780
66	65 1/2	792	792
67	66 1/2	804	804
68	67 1/2	816	816
69	68 1/2	828	828
70	69 1/2	840	840
71	70 1/2	852	852
72	71 1/2	864	864
73	72 1/2	876	876
74	73 1/2	888	888
75	74 1/2	900	900
76	75 1/2	912	912
77	76 1/2	924	924
78	77 1/2	936	936
79	78 1/2	948	948
80	79 1/2	960	960
81	80 1/2	972	972
82	81 1/2	984	984
83	82 1/2	996	996
84	83 1/2	1008	1008
85	84 1/2	1020	1020
86	85 1/2	1032	1032
87	86 1/2	1044	1044
88	87 1/2	1056	1056
89	88 1/2	1068	1068
90	89 1/2	1080	1080
91	90 1/2	1092	1092
92	91 1/2	1104	1104
93	92 1/2	1116	1116
94	93 1/2	1128	1128
95	94 1/2	1140	1140
96	95 1/2	1152	1152
97	96 1/2	1164	1164
98	97 1/2	1176	1176
99	98 1/2	1188	1188
100	99 1/2	1200	1200

Account from 1. 10. to 31. 12. by the Way, and so on by the Day. For if you come by the way, 10. a Day is 250 1/2. a Year, and so on: a Day is 250 1/2. a Year, and so on: a Day is 1. 10. to 1000.

THE FIRST MOTIVE.

and some thereby, for by these motives, these are... (text continues with dense, partially obscured text)

28) The count is that: I believe not matter... (text continues with handwritten notes and calculations)

29) A count of the count is that... (text continues with handwritten notes)

30) I have been... (text continues with handwritten notes)

31) ... (text continues with handwritten notes)

32) ... (text continues with handwritten notes)

33) ... (text continues with handwritten notes)

34) ... (text continues with handwritten notes)

PETIT L'EGARDE

... (text continues with dense, partially obscured text)

... (text continues with dense, partially obscured text)

... (text continues with dense, partially obscured text)



B&M PRIVATE
 H&M HOLLDS
 W&P W&P W&P
 JULY 1916 PAGE 21

... (text continues with dense, partially obscured text)

The TEI manifesto

- 1 The Guidelines are **descriptive** of many different ways and levels of encoding a digital text, not **prescriptive**
- 2 The Guidelines should be **technology-agnostic**. They currently use XML, but are prepared to change
- 3 The schema is modelled as **independently** as possible, though it currently uses RELAX NG to describe content models
- 4 A project is **actively** encouraged to develop an **appropriate subset** of the Guidelines, and apply domain-appropriate constraints

The TEI is built using a literate programming system: ODD (one language does it all)

A set of TEI elements which describe

- elements and attributes
- descriptions (in multiple languages)
- examples
- content models and datatypes
- information about how it can be used
- constraints
- equivalences (eg to formal ontologies like FRBR or CIDOC CRM)

Original tagdoc for <resp> element in TEI P2 (20 years ago)

```
ktagdoc usage=rwa id="resp"><gi>resp</gi>
<name>statement of responsibility</name>
<desc>supplies information about someone other than an author,
sponsor, funder or principal researcher responsible for the
intellectual content of a text, edition, recording, or
series.</desc>
<attlist></attlist>
<exemplum><eg><![CDATA [
  <resp><role>transcribed from original ms</role>
  <name>Claus Huitfeldt</name>
</resp>
]]>
</eg></exemplum>
<exemplum><eg><![CDATA [
  <resp><role>converted to SGML encoding</role>
  <name>Alan Morrison</name>
</resp>
]]>
</eg></exemplum>
<remarks></remarks>
<part>auxiliary tag set for TEI headers</part>
<classes>
<files names='teihdr2'>
<datadesc></datadesc>
<parents>editionStm recording seriesStm titleStm</parents>
<children>role name</children>
<elemdecl>
<![CDATA [
<!ELEMENT resp          - o ((role & name)+)          >
]]>
</elemdecl>
<attldecl>
<![CDATA [
<!ATTLIST resp          %a.global;                    >
]]>
</attldecl>
<xref target='hd21'>
</tagdoc>
```

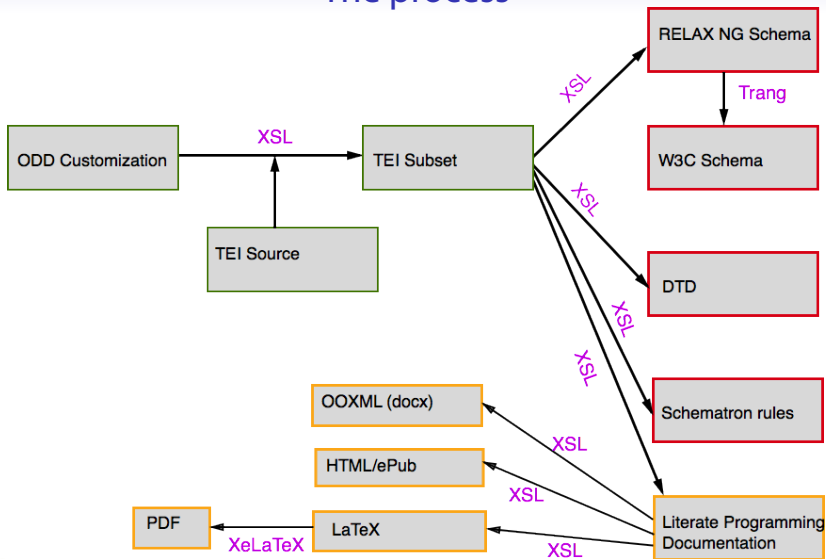
How we do ODD now

```
<elementSpec module="core" ident="respStmt">
  <gloss>statement of responsibility</gloss>
  <desc versionDate="2007-01-21" xml:lang="it">fornisce una dichiarazione di
    responsabilità per qualcuno responsabile del contenuto intellettuale di un
    testo, curatela, registrazione o collana, nel caso in cui gli elementi
    specifici per autore, curatore ecc. non sono sufficienti o non
    applicabili.</desc>
  <classes>
    <memberOf key="att.global"/>
    <memberOf key="model.respLike"/>
    <memberOf key="model.recordingPart"/>
  </classes>
  <content>
    <rng:group>
      <rng:oneOrMore>
        <rng:ref name="resp"/>
      </rng:oneOrMore>
      <rng:oneOrMore>
        <rng:ref name="model.nameLike.agent"/>
      </rng:oneOrMore>
    </rng:group>
  </content>
  <exemplum versionDate="2008-04-06" xml:lang="fr">
    <egXML><respStmt>
      <resp>Nouvelle édition originale</resp>
      <persName>Geneviève Hasenohr</persName>
    </respStmt>
  </egXML>
</exemplum>
</elementSpec>
```


We use the same language to define a customization

```
<schemaSpec
  ident="myschema"
  source="http://www.tei-c.org/release/xml/tei/odd/p5subset.xml">
  <moduleRef key="tei"/>
  <moduleRef key="core"/>
  <moduleRef key="header"/>
  <moduleRef key="textstructure"/>
  <moduleRef key="namesdates" include="persName placeName"/>
  <moduleRef key="figures" except="formula"/>
  <elementSpec ident="title" mode="change">
    <attList>
      <attDef ident="type" mode="change">
        <datatype minOccurs="1" maxOccurs="unbounded">
          <rng:text/>
        </datatype>
        <valList mode="replace" type="closed">
          <valItem ident="biography"/>
          <valItem ident="chronology"/>
          <valItem ident="introduction"/>
          <valItem ident="project"/>
        </valList>
      </attDef>
    </attList>
  </elementSpec>
</schemaSpec>
```

The process



What's the problem?

We're neither one thing nor the other.

Currently in P5:

- Element content models are expressed using a subset of RNG
- Attribute datatypes are expressed using RNG references to W3C datatypes
- Semantic constraints are expressed using ISO Schematron rules

Why don't we just write a huge RELAX NG schema and embed TEI documentation in it?



Choices

- 1 Keep on as we are
 - 2 Rewrite everything in pure RELAX NG
 - 3 Define the whole schema language in TEI
- 1 We have two ways to do things. This is a recipe for confusion
 - 2 We would tie ourselves to one technology
 - 3 We need to show added value from doing so

Looking at element content models

ODD must be intended to support (as far as possible) the **intersection** of what is possible using the current three different schema languages.

In practice, this reduces our modelling requirements quite significantly.

(It also reduces the scope of what we can model)

Requirements for our content modelling system

- 1 It must support alternation, repetition, and sequencing of individual elements, element classes, or sub-models (groups of elements)
- 2 Only one kind of mixed content model — the classic `(#PCDATA | foo | bar)*` — is permitted
- 3 The SGML ampersand connector — `(a & b)` as a shortcut for `((a, b) | (b, a))` is not permitted
- 4 A parser or validator is not required to do look ahead and consequently the model must be deterministic, that is, when applying the model to a document instance, there must be only one possible matching label in the model for each point in the document

Change 1: Define new ODD elements to represent syntax of content models

Specifically:

- `<sequence>` to indicate that its children form a sequence within a content model
- `<alternate>` to indicate that its children can be alternated within a content model

Change 2: provide new att.repeatable class of attributes

- Attributes *@minOccurs* and *@maxOccurs* are currently defined locally on the `<datatype>` element
- Instead provide them via a new class, to which existing elements `<elementRef>`, `<classRef>` and `<macroRef>` elements are added
- Default value for both *@minOccurs* and *@maxOccurs* is 1.

Change 3: re-express generic `<rng:ref>` elements as appropriate XML ODD elements

For example,

```
<rng:ref name="model.pLike"/>
```

becomes

```
<classRef key="model.pLike"/>
```

Example 1 — repeated alternation

((a, (b|c)*, d+), e?) is expressed as follows:

```
<sequence>
  <sequence>
    <elementRef key="a"/>
    <alternate minOccurs="0" maxOccurs="unlimited">
      <elementRef key="b"/>
      <elementRef key="c"/>
    </alternate>
    <elementRef key="d" maxOccurs="unlimited"/>
  </sequence>
  <elementRef key="e" minOccurs="0"/>
</sequence>
```

Example 2 — repeated sequence

$((a, (b^* | c^*)))^+$ is expressed as follows:

```
<sequence maxOccurs="unlimited">  
  <elementRef key="a"/>  
  <alternate>  
    <elementRef key="b" minOccurs="0" maxOccurs="unlimited"/>  
    <elementRef key="c" minOccurs="0" maxOccurs="unlimited"/>  
  </alternate>  
</sequence>
```

Example 3 — treatment of class references

Each class reference is understood to mean any one member of the class:

```
<sequence>
  <classRef key="model.a"/>
  <classRef key="model.b" maxOccurs="unlimited"/>
  <alternate minOccurs="0" maxOccurs="unlimited">
    <classRef key="model.c"/>
    <classRef key="model.d"/>
  </alternate>
</sequence>
```

The *@expand* attribute is used to vary this behaviour in the same way as the existing *@generate* on `<classSpec>`

Examples using *@expand*

Supposing that elements a and b constitute the members of class model.ab:

`<classRef key="model.ab" expand="sequence"/>` is interpreted as a, b

`<classRef key="model.ab" expand="sequenceOptional"/>` is interpreted as a?, b?

`<classRef key="model.ab" expand="sequenceRepeatable"/>` is interpreted as a+, b+

`<classRef key="model.ab" expand="sequenceOptionalRepeatable"/>` is interpreted as a*, b*

Example 4 — mixed content

A mixed content model such as $(\#PCDATA \mid a \mid \text{model.b})^*$ would be expressed as follows, borrowed the *@mixed* attribute from XSD:

```
<alternate minOccurs="0" maxOccurs="unlimited" mixed="true">  
  <elementRef key="a"/>  
  <classRef key="model.a"/>  
</alternate>
```

New and old

```
<alternate>
  <sequence>
    <elementRef key="resp" maxOccurs="unbounded"/>
    <classRef key="model.nameLike.agent" maxOccurs="unbounded"/>
  </sequence>
  <sequence>
    <classRef key="model.nameLike.agent" maxOccurs="unbounded"/>
    <elementRef key="resp" maxOccurs="unbounded"/>
  </sequence>
</alternate>
```

```
<rng:choice>
  <rng:group>
    <rng:oneOrMore>
      <rng:ref name="resp"/>
    </rng:oneOrMore>
    <rng:oneOrMore>
      <rng:ref name="model.nameLike.agent"/>
    </rng:oneOrMore>
  </rng:group>
  <rng:group>
    <rng:oneOrMore>
      <rng:ref name="model.nameLike.agent"/>
    </rng:oneOrMore>
    <rng:oneOrMore>
      <rng:ref name="resp"/>
    </rng:oneOrMore>
  </rng:group>
</rng:choice>
```

The start of added value

The ability to specify repetition at the individual class level gives a further level of control not currently possible.

'no more than two consecutive sequences of all members of the class model.nameLike'

```
<classRef key="model.nameLike" maxOccurs="2" expand="sequence"/>
```


Progress so far

- 1 New and changed elements defined as a TEI customization
- 2 Processing tools enhanced to cover the new elements
- 3 Conversion from old content models written and tested

<https://github.com/TEIC/pureodd>

A few practical issues to look at

- Investigate how to model embedded MathML and SVG (just use NVDL?)
- Develop native conversion to W3C Schema (remove dependency on trang)

Going large (1)

Suppose we forget about supporting only the intersection of current schema language facilities?

- features which are present in one schema language (but not all) are probably there because someone wanted them!
- can we rethink ODD to cater for (potentially) **all** schema language features, rather than their intersection?
- one possible implementation strategy: use an additional constraint language such as ISO Schematron to mop up the parts that a specific schema language cannot support.
- We would like to recast some constraints which are currently in raw Schematron into pure TEI

Going large (2)

Example 1 we want a content model like (a&b&c&d) but only RELAX NG provides `interleave`

Example 2 we want different content models for `teiHeader//p` and for `text//p` but only W3C Schema has concept of `base` types

Going large (3)

Example 1 add an attribute `@preserveOrder` with values `true` or `false` to the `<sequence>` element

Example 2 add an attribute `@context` with an XPath expression as value to `<elementRef>` and friends

```
<elementRef key="s" context="ancestor::text" minOccurs="1"/>  
<macroRef key="macro.limitedContent" context="ancestor::teiHeader"/>
```

In the absence of an exact equivalent in the target schema language, an ODD processor can choose to

- flag the construct as illegal
- overgenerate, by producing code which validates the target construct plus others
- compensate, by over-generating but also producing Schematron code to catch 'false positives'

Autocritique

- 1 TEI, schmei. Just use HTML5 and stop being obfuscatory
 - 2 you're just re-expressing RELAX in a similar language
 - 3 who cares? validation is so 20th century
 - 4 you're imposing a bottleneck in processing, limited by a single implementation of an under-specified idea
- 1 is `` better than `<unclear>`?
 - 2 yes, at the start. but now we can extend
 - 3 if you want interoperability, you need a language in which to express 'business rules'
 - 4 a fair cop, sort of. but the old system already the bottleneck, now we are simplifying it

Conclusions

Was this worth it? **Yes:**

- we expect a lot more human reading and changing of constraints than most schemas
- a single language to express as many constraints as possible helps our users
- we have a coherent platform on which to express more of our semantic rules
- the TEI is positioning itself to be free of XML, if alternatives appear

An extensible independent notation for expressing text encoding Guidelines takes the TEI back to its roots