

# Talk 13: Customising and Documenting Your TEI project

James Cummings

28 January 2014

# Customising the TEI

- How the TEI is constructed
- Making a TEI schema
- Specifying your profile of the TEI
- Generating your own documentation

**Every** use of the TEI involves making use of a customisation of the TEI.

## Some terminology

- The TEI encoding scheme consists of a number of *modules*
- Each module contains a number of *element specifications*
- Each element specification contains:
  - a canonical name (`<gi>`) for the element, and optionally other names in other languages
  - a canonical description (also possibly translated) of its function
  - a declaration of the *classes* to which it belongs
  - a definition for each of its *attributes*
  - a definition of its *content model*
  - usage examples and notes
- a TEI *schema* specification (`<schemaSpec>`) is made by selecting modules or elements and (optionally) modifying their contents
- a TEI document containing a schema specification is called an *ODD* (One Document Does it all)

## What is a module?

- A convenient way of grouping together a number of element declarations
- These are usually on a related topic or specific application
- Most chapters of P5 focus on elements drawn from a single module, which that chapter then defines
- A TEI Schema is created by selecting modules and adding or removing elements from them as needed

## Which modules exist?

Module name	Chapter
analysis	Simple Analytic Mechanisms
certainty	Certainty and Responsibility
core	Elements Available in All TEI Documents
corpus	Language Corpora
dictionaries	Dictionaries
drama	Performance Texts
figures	Tables, Formulae, and Graphics
gaiji	Representation of Non-standard Characters and Glyphs
header	The TEI Header
iso-fs	Feature Structures
linking	Linking, Segmentation, and Alignment
msdescription	Manuscript Description
namesdates	Names, Dates, People, and Places
nets	Graphs, Networks, and Trees
spoken	Transcriptions of Speech
tagdocs	Documentation Elements
tei	The TEI Infrastructure
textcrit	Critical Apparatus
textstructure	Default Text Structure
transcr	Representation of Primary Sources
verse	Verse

## How do you choose?

- Just choose everything (not really a good idea)
- The TEI provides a small set of predefined combinations (TEI Lite, TEI Bare...)
- Or you could roll your own (but then you need to know what you're choosing)

**Roma** a command line script, with a web front end,  
designed to make this process much easier

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

# Roma: New



## Roma: generating customizations for the TEI

These pages will help you design your own TEI customization, as a DTD, RELAX NG or W3C Schema.

### Create a new or upload existing customization

- Build up: create a new customisation by adding elements and modules to the smallest recommended schema
- Reduce: create a new customization by removing elements and modules from the largest possible schema
- Create customization from template
- Open existing customization

**Start**

Roma was written by Arno Mittelbach and is maintained by Sebastian Rahtz. Sanity check written by Ioan Bernevig. Please direct queries to the [TEI@Oxford](mailto:TEI@Oxford) project.

# Roma: Customize



## Roma: generating validators for the TEI

You are currently working on **TEI Absolutely Bare**

### Set your parameters

[New](#) [Customize](#) [Language](#) [Modules](#) [Add Elements](#) [Change Classes](#) [Schema](#) [Documentation](#) [Save Customization](#) [Sanity Checker](#)

#### Set your parameters

Title	<input type="text" value="TEI Absolutely Bare"/>
Filename	<input type="text" value="tei_bare"/>
Namespace for new elements	<input type="text" value="http://www.example.org/ns/nonTEI"/>
Prefix for TEI pattern names in schema	<input type="text"/>
Language	<input type="radio"/> English, <input type="radio"/> Deutsch, <input type="radio"/> Italiano, <input type="radio"/> Español, <input type="radio"/> Français, <input type="radio"/> Portugues, <input type="radio"/> Russian, <input type="radio"/> Svenska, <input type="radio"/> 日本語, <input type="radio"/> 中文
Author name	<input type="text" value="Laurent Romary"/>
Description	<pre>This customization creates a TEI schema with the bare minimum of tags to make a recognizable document. It combines the four basic modules, but removes most of the available elements and changes several attribute classes.</pre>

[Save](#)

Roma was written by Arno Mittelbach and is maintained by Sebastian Rahtz. Sanity check written by Ioan Bernevig. Documentation language en. Please direct queries to the [TEI @ Oxford](#) project. This is Roma version 4.9, last updated 2012-06-16. Using TEI P5 version 2.1.0



# Roma: Schema



Roma:

## Production des validateurs TEI

Vous travaillez actuellement sur **TEI Absolutely Bare**

Enfin je vous rends votre schéma...

[Nouvelle](#) [Personnaliser](#) [Langage](#) [Modules](#) [Ajouter des éléments](#) [Modifier les classes](#) [Schéma](#)  
[Documentation](#) [Enregistrer](#) [Contrôleur de validité](#)

Creation du schéma en cours

Quel format  
préférez-vous?

RELAX NG schema (compact syntax) ·  
RELAX NG schema (compact syntax)  
RELAX NG schema (XML syntax)  
ISO Schematron  
Schematron  
W3C schema (in ZIP archive)

**Generate**

We use RELAX NG DTD

[http://en.wikipedia.org/wiki/XML\\_Schema\\_Language\\_comparison](http://en.wikipedia.org/wiki/XML_Schema_Language_comparison).

Roma was written by Arno Mittelbach and is maintained by Sebastian Rahtz. Sanity check written by Ioan Bernevig.  
Documentation language en. Please direct queries to the [TEI@Oxford](#) project.  
This is Roma version 4.9, last updated 2012-06-16. Using TEI P5 version 2.1.0

# Roma: Documentation



Roma:

## Production des validateurs TEI

Vous travaillez actuellement sur **TEI Absolutely Bare**

### Documentation?

[Nouvelle Documentation](#) [Personnaliser](#) [Langage](#) [Modules](#) [Ajouter des éléments](#) [Modifier les classes](#) [Schéma](#)  
[Enregistrer](#) [Contrôleur de validité](#)

Création de la documentation en cours

Quel format  
préférez-vous?

**Generate**

HTML web page ·  
HTML web page  
PDF  
TEI Lite  
TEI ODD

Roma was written by Arno Mittelbach and is maintained by Sebastian Rahtz. Sanity check written by Ioan Bernevig.  
Documentation language en. Please direct queries to the [TEI @ Oxford](#) project.  
This is Roma version 4.9, last updated 2012-06-16. Using TEI P5 version 2.1.0

## What did we just do?

We processed a pre-existing ODD file which contained (as well as some discursive prose) the following schema specification:

```
<schemaSpec ident="tei_bare" start="TEI">
  <moduleRef key="core"/>
  <moduleRef key="tei"/>
  <moduleRef key="header"/>
  <moduleRef key="textstructure"/>
  <elementSpec ident="abbr" mode="delete" module="core"/>
  <elementSpec ident="add" mode="delete" module="core"/>
<!-- ... -->
  <elementSpec ident="trailer" mode="delete" module="textstructure"/>
  <elementSpec ident="title" mode="change" module="core">
    <attList>
      <attDef ident="level" mode="delete"/>
    </attList>
  </elementSpec>
<!-- ... -->
</schemaSpec>
```

We selected four modules, deleted loads of elements, and also deleted an attribute.

## Roma provides an interface to the detail

- The [Modules] tab shows the modules available
- Selecting a module from it shows the elements within that module, and gives you the choice to
  - include all of them (and then remove some)
  - exclude all of them (and then put back the ones you want)
- You can also change an element's attribute list, and the values they permit

# Roma: Modules

Liste des modules TEI			
	Titre du module	Description brève	Modification
ajouter	<a href="#">analysis</a>	🔗 Mécanismes analytiques simples	
ajouter	<a href="#">certainty</a>	🔗 Degré de certitude et responsabilité	
ajouter	<a href="#">core</a>	🔗 Éléments disponibles pour tous les documents TEI	modifié
ajouter	<a href="#">corpus</a>	🔗 Corpus linguistiques	
ajouter	<a href="#">dictionaries</a>	🔗 Dictionnaires	
ajouter	<a href="#">drama</a>	🔗 Théâtre	
ajouter	<a href="#">figures</a>	🔗 Tableaux, formules et graphiques	
ajouter	<a href="#">gaiji</a>	🔗 Représentation des caractères et des glyphes non standard	
ajouter	<a href="#">header</a>	🔗 En-tête TEI	modifié
ajouter	<a href="#">iso-fs</a>	🔗 Structures de traits	
ajouter	<a href="#">linking</a>	🔗 Liens, segmentation et alignement	
ajouter	<a href="#">msdescription</a>	🔗 Description de manuscrits	
ajouter	<a href="#">namesdates</a>	🔗 Noms, dates, personnes et lieux	
ajouter	<a href="#">nets</a>	🔗 Graphes, réseaux et arbres	
ajouter	<a href="#">spoken</a>	🔗 Transcriptions de la parole	
ajouter	<a href="#">tagdocs</a>	🔗 Éléments de déclaration d'un modèle	
ajouter	<a href="#">textcrit</a>	🔗 Apparat critique	
ajouter	<a href="#">textstructure</a>	🔗 Structure textuelle par défaut	modifié
ajouter	<a href="#">transcr</a>	🔗 Représentation de sources primaires	
ajouter	<a href="#">verse</a>	🔗 Poésie	

## Liste des Modules sélectionnés

supprimer [core](#)  
 tei  
 supprimer [header](#)  
 supprimer [textstructure](#)

# Roma: Change Module

## Liste des éléments compris dans ce modulecore

	Inclure	Exclure	Nom	Description	Attributs
abbr	○	*	<input type="text" value="abbr"/>	? (abréviation) contient une abréviation quelconque.	<a href="#">Changer les attributs</a>
add	○	*	<input type="text" value="add"/>	? ( ajout) contient des lettres, des mots ou des phrases insérées dans le texte par un auteur, un copiste, un annotateur ou un correcteur.	<a href="#">Changer les attributs</a>
addrLine	○	*	<input type="text" value="addrLine"/>	? (ligne d'adresse) contient une ligne d'adresse postale.	<a href="#">Changer les attributs</a>
address	○	*	<input type="text" value="address"/>	? contient une adresse postale ou d'un autre type, par exemple l'adresse d'un éditeur, d'un organisme ou d'une personne.	<a href="#">Changer les attributs</a>
analytic	○	*	<input type="text" value="analytic"/>	? (niveau analytique) contient des éléments descriptifs qui décrivent la bibliographie d'une ressource (par exemple un poème ou un article de revue) publiée à l'intérieur d'une monographie ou d'une ressource et non publiée de façon indépendante.	<a href="#">Changer les attributs</a>
author	*	○	<input type="text" value="author"/>	? (auteur) dans une référence bibliographique contient le nom de la (des) personne(s) physique(s) ou du collectif, auteur(s) d'une oeuvre ; la première mention de responsabilité comme seul élément bibliographique.	<a href="#">Changer les attributs</a>
bibl	○	*	<input type="text" value="bibl"/>	? (référence bibliographique.) contient une référence bibliographique faiblement structurée dans laquelle les	<a href="#">Changer les attributs</a>

# What does our REED corpus need?

A simple selection of elements, but also

- we want to allow only certain values for *@type* on `<div>`
- we may want to create a new element (can you think of something?)

Other constraints are possible — we might want to insist that a `<div type="record">` contains a date, for example.

## The ODD advantage

We can express these constraints in our ODD meta-schema, and then generate a formal schema to enforce them using whichever schema language we like.

- TEI schemas can be generated in
  - ISO RELAX NG language
  - W3C Schema Language
  - XML DTD language
- ODD itself defines an element's content models using a subset of RELAX NG syntax
- Datatypes are defined in terms of W3C datatypes
- Some facilities (e.g. alternation, namespaces) cannot be expressed in DTDs — RELAX NG schema is recommended
- Additional constraints can be expressed in Schematron



# Roma: selecting attributes

List of attributes: div

Add new attributes

Change attribute	Include	Exclude	Name	Description	Delete
<u>ana</u>	<input type="radio"/>	<input type="radio"/>	<input type="text" value="ana"/>	indicates one or more elements containing interpretations of the element on which the ana attribute appears.	
<u>change</u>	<input type="radio"/>	<input type="radio"/>	<input type="text" value="change"/>	points to one or more change elements documenting a state or revision campaign to which the element bearing this attribute and its children have been assigned by the encoder.	
<u>copyOf</u>	<input type="radio"/>	<input type="radio"/>	<input type="text" value="copyOf"/>	points to an element of which the current element is a copy.	
<u>corresp</u>	<input type="radio"/>	<input type="radio"/>	<input type="text" value="corresp"/>	points to elements that correspond to the current element in some way.	
<u>decls</u>	<input type="radio"/>	<input type="radio"/>	<input type="text" value="decls"/>	identifies one or more declarable elements within the header, which are understood to apply to the element bearing this attribute and its content.	
<u>exclude</u>	<input type="radio"/>	<input type="radio"/>	<input type="text" value="exclude"/>	points to elements that are in exclusive alternation with the current element.	
<u>fac</u> s	<input type="radio"/>	<input type="radio"/>	<input type="text" value="fac"/>	points to all or part of an image which corresponds with the content of the element.	
<u>met</u>	<input type="radio"/>	<input type="radio"/>	<input type="text" value="met"/>	contains a user-specified encoding for the conventional metrical structure of the element.	
<u>n</u>	<input type="radio"/>	<input type="radio"/>	<input type="text" value="n"/>	gives a number (or other label) for an element, which is not necessarily unique within the document.	

# Roma: constraining attribute values

## Add a new attribute

**Attribute name** type

**Class name**

**Is it optional?**

- yes
- no

**Contents**

Text

**Default value**

**Closed list?**

- yes
- no

**List of values**

**Description**

Save

## What did we just do?

Our ODD now includes something like this:

```
<elementSpec ident="div" module="textstructure" mode="change">
  <attList>
    <attDef ident="type" mode="change" usage="req">
      <valList type="closed" mode="replace">
        <valItem ident="prose"/>
        <valItem ident="verse"/>
        <valItem ident="drama"/>
      <!-- ... -->
    </valList>
  </attDef>
</attList>
</elementSpec>
```

Note that we can also add documentation to the ODD:

```
<valItem ident="verse">
  <gloss>contains (parts of ) a poem</gloss>
</valItem>
```

## Defining a new element

When defining a new element, we need to consider

- its name and description
- what attributes it can carry
- what it can contain
- where it can appear in a document

The TEI class system helps us answer all these questions (except the first).

# The TEI Class System

- The TEI distinguishes over 500 elements,
- Having these organised into classes aids comprehension, modularity, and modification.
- *Attribute class*: the members share common attributes
- *Model class*: they can appear in the same locations (and are often semantically related)
- Classes may contain other classes
- An element can be a member of any number of classes, irrespective of the module it belongs to.

## Attribute Classes

- Attribute classes are given (usually adjectival) names beginning with **att.**; e.g. *att.naming*, *att.typed*
- all members of **att.naming** inherit from it attributes *@key* and *@ref*; all members of **att.typed** inherit from it *@type* and *@subtype*
- If we want an element to carry the *@type* attribute, therefore, we add the element to the **att.typed** class, rather than define those attributes explicitly.

## A very important attribute class: `att.global`

All elements are usually members of `att.global`; this class provides, among others:

`@xml:id` a unique identifier

`@xml:lang` the language of the element content

`@n` a number or name for an element

`@rend` how the element in question was rendered or presented in the source text.

## Model Classes

- Model classes contain groups of elements which are allowed in the same place. e.g. if you are adding an element which is wanted wherever the `<bibl>` is allowed, add it to the `model.biblLike` class
- Model classes are usually named with a **Like** or **Part** suffix:
  - members of `model.pLike` are all things that 'behave like' paragraphs, and are permitted in the same places as paragraphs
  - members of `model.pPart` are all things which can appear *within* paragraphs. This class is subdivided into
    - `model.pPart.edit` elements for simple editorial intervention such as `<corr>`, `<del>` etc.
    - `model.pPart.data` 'data-like' elements such as `<name>`, `<num>`, `<date>` etc.
    - `model.pPart.msdesc` extra elements for manuscript description such as `<seal>` or `<origPlace>`



## Basic Model Class Structure

Simplifying wildly, one may say that the TEI recognises three kinds of element:

**divisions** high level major divisions of texts

**chunks** elements such as paragraphs appearing within texts or divisions, but not other chunks

**phrase-level elements** elements such as highlighted phrases which can occur only within chunks

There are 'base model classes' corresponding with each of these, and also with the following groupings:

**inter-level elements** elements such as lists which can appear either in or between chunks

**components** elements which can appear directly within texts or text divisions

And yes, there is a class **model.global** for elements that can appear *anywhere* inside a text — at any hierarchic level.

# Defining a new element

What other elements is it like?

What other elements can contain it?

What can it contain?

Conclusions:

- What model do we select?
- What content model do we select?

# Roma: Defining a new element

## Defining a new element:

**Name**

**Namespace**

**Description**

**Model classes**

- model.addrPart
- model.addressLike
- model.applicationLike
- model.availabilityPart
- model.biblLike
- model.biblPart
- model.castItemPart
- model.catDescPart
- model.certLike
- model.choicePart
- model.common
- model.dateLike
- model.dimLike
- model.div1Like
- model.div2Like
- model.div3Like
- model.div4Like
- model.div5Like
- model.div6Like
- model.div7Like
- model.divBottom
- model.divBottomPart
- model.divGenLike
- model.divLike
- model.divPart
- model.divPart.spoken
- model.divTop
- model.divTopPart
- model.divWrapper
- model.editorialDeclPart
- model.egLike
- model.emphLike
- model.encodingDescPart
- model.entryLike

## Defining a content model

- A typical TEI element defines its content by referencing *classes* of element which it can contain, rather than using specific elements.
- Content models are defined using the RELAXNG vocabulary
- Here are some very common predefined content models:
  - `macro.paraContent` content of paragraphs and similar elements
  - `macro.limitedContent` content of prose elements that are not used for transcription of extant materials
  - `macro.phraseSeq` a sequence of character data and phrase-level elements
  - `macro.phraseSeq.limited` a sequence of character data and those phrase-level elements that are not typically used for transcribing extant documents
  - `macro.specialPara` the content model of elements which either contain a series of component-level elements or else contain a series of phrase-level and

# Roma: Defining a new element 2

## Attribute classes

- att.ascribed
- att.combinable
- att.datable
- att.datable.w3c
- att.declaring
- att.docStatus
- att.duration.w3c
- att.entryLike
- att.global.change
- att.handFeatures
- att.interpLike
- att.metrical
- att.personal
- att.pointing.group
- att.readFrom
- att.segLike
- att.spanning
- att.timed
- att.typed
- att.breaking
- att.coordinated
- att.datable.custom
- att.datcat
- att.dimensions
- att.duration
- att.editLike
- att.global
- att.global.facs
- att.identified
- att.lexicographic
- att.msExcerpt
- att.placement
- att.ranging
- att.responsibility
- att.sortable
- att.tableDecoration
- att.transcriptional
- att.witnessed
- att.canonical
- att.damaged
- att.datable.iso
- att.declarable
- att.divLike
- att.duration.iso
- att.enjamb
- att.global.analytic
- att.global.linking
- att.internetMedia
- att.measurement
- att.naming
- att.pointing
- att.rdgPart
- att.scoping
- att.sourced
- att.textCritical
- att.translatable

## Contents

User content ▾

```
<content xmlns:rng="http://relaxng.org/ns/structure/1.0">
</content>
```

Save

## What did we just do?

We added a new element specification to our ODD, like this:

```
<elementSpec
  ident="something"
  ns="http://www.example.org/ns/nonTEI"
  mode="add">
  <desc> contains something division like.</desc>
  <classes>
    <memberOf key="model.divPart"/>
    <memberOf key="att.typed"/>
  </classes>
  <content>
    <rng:ref name="something"/>
    <rng:oneOrMore>
      <rng:ref name="model.pLike"/>
    </rng:oneOrMore>
  </content>
</elementSpec>
```

Note that this new element is *not* in the TEI namespace. It belongs to this specific project only!

## Other kinds of constraints

- You can also constrain the content of an element or the value of an attribute to be of a particular *datatype* (for example, to insist that the *@when* attribute of the element `<date>` contains only a date)
- This can be done by using one of a set of predefined *macros* to define the content. Examples include
  - `data.word` a single word or token
  - `data.name` an XML Name
  - `data.enumerated` a single XML name taken from a documented list
  - `data.temporal.w3c` a W3C date
  - `data.truthValue` a truth value (true/false)
  - `data.language` a human language
  - `data.sex` human or animal sex
- Or you can define a more complex constraint, e.g. using Schematron

## Schematron constraints

- An element specification can also contain a `<constraintSpec>` element which contains rules about its content expressed as ISO Schematron *constraints*

```
<elementSpec ident="div" mod-  
ule="teiststructure" mode="change"  
  xmlns:s="http://purl.oclc.org/dsdl/schematron">  
  <constraintSpec ident="div" scheme="isoschematron">  
    <constraint>  
      <s:assert test="@type='record' and ../tei:date">prose  
must include a paragraph  
      </s:assert>  
    </constraint>  
  </constraintSpec>  
</elementSpec>
```

However...

- You can only add such rules by editing your ODD file: Roma doesn't know about them.
- Not all schema languages can implement these constraints.



## TEI customization for fun and profit

With the Stationers' Register project the keying company would:

- do double-keying data entry matching **any** XML schema
- and charged per kilobyte of output!

So we:

- used TEI ODD to create a highly compressed / byte-reduced renaming schema
- and up-converted the result back to pure TEI!

## tei\_corset example

```
<nm r="rm">iijli vjs viijd</nm> from  
<n r="b">Thomas barthelett</n>
```

was up-converted to

```
<seg type="notFee" rend="roman-numerals">  
  <num type="totalPence" value="800">  
    <num type="poundsAsPence" value="720"> iij<hi rend="superscript">li</hi>  
    </num>  
    <num type="shillingsAsPence" value="72"> vj<hi rend="superscript">s</hi>  
    </num>  
    <num type="pence" value="8"> viij<hi rend="superscript">d</hi>  
    </num>  
  </num>  
</seg> from  
<persName role="stationer" rend="bold">  
  <forename>Thomas</forename>  
  <surname>barthelett</surname>  
</persName>
```

## Next

Any Questions? That is a quick look at some of the basic things one can do with the TEI ODD language, and the Roma web tool. Now let's do an exercise where we try out customising the TEI!