

Transcribing a manuscript with TEI

What does 'digitization' mean?

(Not the same as 'digitalization'!)

- production of digital images of the pages of a manuscript – a *facsimile*
- production of a *transcription* of the content of a manuscript

For manuscripts, the two are often complementary

Facsimile and transcription together

Был тогда день субботний, и, когда солнце взошло, сошлись оба поги. Немножко же и чуть приблизились синий склон поклон, и была тут сечь зла и велела с немедля и чудом; и стоял трек от поминающейся коней и змей от узоров меяй; и казалось, что зефир замерзшее дунулось, и не было видно пада, ибо все было залито кровью.



Бе тогда день скромный,
самим возвращение и
святыниши сии пахали.
Некои же и чудь
прекрасныи склоненіи скозе
плаки, и насть ту сечь зла
и велика немощи и чудь, и
трут, от кнїзъ альбинонъ
здеи от личаго сечинна,
здеи мадо помърдце
длинниту, и не си видети
леду, покоры не есть вѣ
кроме.



Святейшайше Сенату. Уставинопп
объявлено ажною зею, приша
напомицкимъ императоромъ и
покупчимъ токсигономъ, неспытывъ, несто
мъ искорененіемъ. Сенату привѣти
примѣни, азъ азъ именемъ, не въ
коупчимъ токсигономъ, неспытывъ
модицкимъ императоромъ распоряженіемъ
въ окупчимъ токсигономъ, подѣльце. Ф.
А. Г. Челесинскому мѣстечко. Адресовано
губернатору. Въ. М. М. партизанъ
жизниной обороно. Въ привѣтии уважа
ется. Азъ азъ именемъ, азъ зею
и спѣшишь искогонъ, къ систеконъ
артил., въ. А. напоупланутия. А.
напомицкимъ императоромъ. Азъ
жизгогонъ азъ искогонъ и зею, Азъ
прѣжнимъ зею. Азъ искогонъ зею
Хилъ. Сенату привѣти - напомицкимъ

недемоуєтвіа
архієписко^п піво^п
архієписко^п піво^п

В спасах от очевидного
заслуженного, что вади
шьтейн Бомб в воздухе,
пронесли на постаменте
всего лишь минуту. И обнаруж-
или синюю помаду, смешан-
ную с цветами мученика Бориса
Годеба, который из них вредил.
Сама же помада, якобы, яркая, и
все же не красила. Но это не было
импульсом открытия, а было
импульсом на претворение ско-
нов в реальность. И это было
до Судебного здания. В

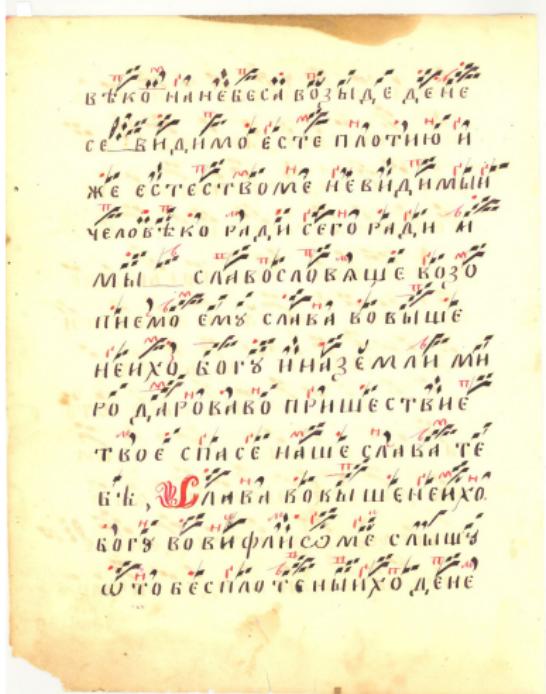
Transcription is not an automatic process

И съяспо́е го́рьзлай здните пе́ска . и пои́де
и ки опоу́гра́ду . и сре́тє иноу́рать .
и не о́усты́дї вѣ́л противопоу́стити .

и сына своего взя из
Витебска, и пошел к
Новгороду, и встретил иную
рать, и не устыдился
противу их стати, и

а сына своего взял из
Витебска и пошел к
Новгороду. И встретил иную
рать, и встал против них,

Transcription: a special kind of reading



What is the goal of your transcription?

- to make a primary source accessible ...
- ... but also comprehensible
- which may imply adding (or suppressing) a lot of information

Because...

- all transcription is selective
- all transcription is imaginative

Transcription

What does a transcription add to a simple facsimile?

Transcribers typically try to make explicit :

- (some) original layout information
- abbreviations and other strange symbols
- 'evident' errors which invite correction or conjecture
- scribal additions, deletions, substitutions, restorations
- non-standard orthography (etc.) which invites normalisation
- irrelevant or non-transcribable material
- passages which are damaged or illegible

What kind of transcription do you want?

- <teiHeader>: provides metadata for the whole thing, at various levels, typically including a <msDesc>
- <text>: contains a structured reading of a document's intellectual content ... its 'text'
- <facsimile>: organizes a set of page images representing a document
- <sourceDoc>: a non-interpretative transcription of a physical document, e.g. for a *dossier génétique*

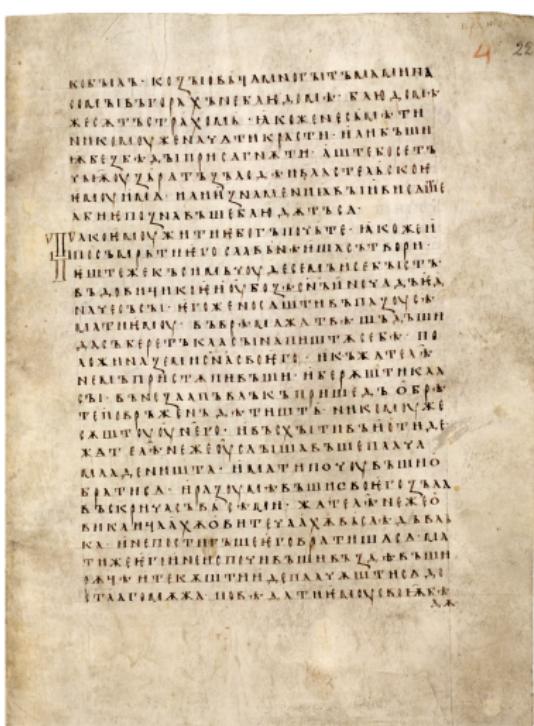
Does your transcription represent a 'text' or a 'document' ?



How is your transcription organised?

- Just pages or folios, composed of blocks or lines
- Sections, paragraphs, verse lines, lists, sentences ...
- Or both?

Layers of transcription



- Paleographic level : identification of characters and graphemes
- Documentary (or diplomatic) level: decide what has been written
- Editorial or semantic level : decide how it should be read

A typical minimal encoding

5

1 demy farbushans de -
soixante sept pieces au marc
a vng felin & demy de remede
contre charme puer po
cinquanti souz tournois.

Sur le pris de quinze livres
tours le marc dargent le roy
fust continu la faburam
des testons a uno denier
de huit graine troye quartz
et fin qui valz vingt deniers
les graine dargent le roy a

<pb n="13" xml:id="B452346101_Ms629_0015"/>
<fw place="top-right" type="pageNum">5</fw>
<lb/>& demy trebuchans de
<lb/>soixante sept pieces au marc
<lb/>a ung felin & demy de remede
<lb/>Cours chacune piece <expan>pour</expan>
<lb/><expan>cinquante</expan> soubz
<expan>tournois</expan><pc>.</pc></p>
<lb/><p><lb/>Sur le pris de quinze livres
<lb/><expan>tournois</expan> le marc dargent
le roy
<lb/>fust <expan>continuée</expan> la
<expan>fabricacion</expan>
<lb/>des testons a dix deniers
<lb/>dix huict grains troyes quartz
<lb/>de fin qui <expan>valent</expan> unze deniers
<lb/>six grains dargent le roy<pc>.</pc> a

Representation of the physical structure

- The physical organisation of a manuscript (its binding, folios, leaves, pages, columns) rarely, if ever, corresponds with its logical organisation (sections, chapters, paragraphs, lines)
- Whichever we choose to represent in our XML structure, we will have to represent the other using empty 'milestone' elements
- For example, in the logical view, we can use `<gb>`, `<pb>`, `<cb>`, or `<lb>` to indicate the start of gatherings, pages, column or lines
- Or in the physical view, we could use a `< milestone >` to indicate the starts of divisions, paragraphs, etc.

Characters and glyphs

- the same character may be represented in many different forms
 - e.g. a A a a a a ... ==> U+0061
 - e.g. S ==> U+0073 l ==> U+017F
- the character or glyph we see may not yet exist in Unicode

The `<g>` element allows us to indicate the presence of a specific glyph, or a non-Unicode character

Using <g>

<lb/>на мѣст<g ref="#ooGlyph">оо</g>
<lb/>гостинницю. гледаше
<lb/>само и шнам<g ref="#ooGlyph">оо</g>
<lb break="no">ти. въсклабив же се рече

Bdinski, fol 7r, detail

There is no Unicode character for the ligatured oo here: we tag it as a <g>

<lb/>на мѣст<g ref="#ooLig">оо</g>. и вышедь въ
<lb/>гостинницю. гледаше
<lb/>само и шнам<g ref="#ooLig">оо</g> хоте видѣ
<lb break="no"/>ти. въсклабив же се рече

#ooGlyph points to a description of the glyph, provided in the TEI header.



Abbreviations &c.

In Western MSS, we commonly distinguish :

Suspensions the first letter or letters of the word are written, generally followed by a point : for example 'e.g.' for 'exempla gratia'

Contractions both first and last letters are written, generally with some mark of abbreviation such as superscript strokes, or points : e.g. 'Mr.' for 'Mister'

Brevigraphs Special signs such as the Tironian *nota* used for 'et', the letter p with a barred tail used for 'per', the letter c with a circumflex used for 'cum' etc.

Superscripts Superscript letters (vowels or consonants) used to indicate various kinds of contraction: e.g. 'w' followed by superscript 'ch' for 'which'.

Most of the symbols needed are available in Unicode, though not necessarily in all fonts.



Abbreviation and Expansion

An abbreviation may be viewed in two different ways:

- as a particular sequence of letters or marks upon the page: thus, a 'p with a bar through the descender', a 'superscript hook', a 'macron'
- as an alternative way of representing a sequence of letters : thus, 'per', 're', 'n'

Two sets of tags

TEI proposes elements for two levels of encoding:

- the whole of an abbreviated word and the whole of its expansion may be marked using `<abbr>` and `<expan>` respectively
- abbreviatory signs or characters and the ‘invisible’ characters they imply may be marked using `<am>` and `<ex>` respectively



A French example

Coures châtrines p[er]s[er] p[er]e.
cinquante. soubz tournois.

We might just note that we have expanded the abbreviations:

```
<p>
<lb/>Cours chacune piece <expan>pour</expan>
<lb/>
<expan>cinquante</expan> soubz <expan>tournois</expan>
<pc>.</pc>
</p>
```

... or we might just record the abbreviated forms

As you noticed, 'pour' was actually written 'po' followed by an 'r' subscript; 'cinquante' as 'cinquāte' with a macron on the 'a' to indicate nasalisation.

```
<p>
  <abbr>po&#xFFFD;</abbr> .... <abbr>cinqu&#x0101;te</abbr>
</p>
```

... or we might look a bit closer

.We can tag the abbreviation markers and the expansion directly

```
<p> po<am> ... or po<ex>u</ex>r </p>
```

... or within the **<abbr>** or **<expan>** as appropriate

```
<abbr>po<am></am>
</abbr>
```

```
<expan>po<ex>u</ex>r</expan>
```

Or we might want to show there's a <choice> ...

The <choice> element wraps alternative mutually exclusive ways of **encoding** the same phenomenon:

- <choice> (groups alternative editorial encodings)
- Abbreviation:
 - <abbr> (abbreviated form)
 - <expan> (expanded form)
- Errors:
 - <sic> (apparent error)
 - <corr> (corrected error)
- Regularization:
 - <orig> (original form)
 - <reg> (regularized form)

Not intended for use with textual variants (for which, use <app>)



Types of abbreviation

The `@type` attribute on `<abbr>` is a useful way of categorising abbreviations, whether for statistical purposes, or to allow for different types to be rendered differently:

```
<choice>
  <abbr type="brevigraphe">po<am>&#xFFD;</am>
  </abbr>
  <expan resp="#LB">po<ex>u</ex>r</expan> en <choice>
    <abbr type="suspension">fin<am>.</am>
    </abbr>
    <expan>fin<ex>ir</ex>
    </expan>
  </choice>
</choice>
```

As elsewhere, the `@resp` and `@cert` attributes can be used to indicate who is responsible for an expansion, and the degree of certainty attached to it.

This encoding might be displayed as : 'po(u)r en finir [LB]'



Corrections and emendations

The `<sic>` element can be used to indicate that the reading of the manuscript is erroneous or nonsensical, while `<corr>` (correction) can be used to provide what in the editor's opinion is the correct reading:

```
<sic>relea</sic>
```

```
<corr>relicta</corr>
```

The two may, of course, be combined within a `<choice>` element:

```
<choice>
  <sic>relea</sic>
  <corr cert="high">relicta</corr>
  <corr cert="low">relatio</corr>
</choice>
```

Normalization

Source texts rarely use modern orthography. For retrieval and other processing reasons, however, the modernized form may be needed.

The `<reg>` (regularized) element is available used to mark a normalized form; the `<orig>` (original) element to indicate a non-standard spelling. These elements can of course be grouped as alternatives using the `<choice>` element

A Russian example

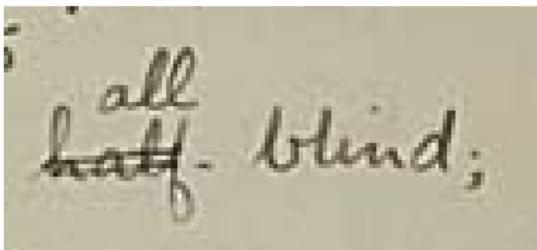
Какъ вѣтеръ мокрый, ты бѣешься въ ставни,
Какъ вѣтеръ черный, поешь: ты мой!
Я зорній часы, я долгъ твой лавній.

```
<!-->
<choice>
  <orig>Какъ</orig>
  <reg>Как</reg>
</choice>
<choice>
  <orig>вѣтеръ</orig>
  <reg>ветер</reg>
</choice> мокрый, ты бѣешься <choice>
  <orig>въ</orig>
  <reg>в</reg>
</choice> ставни,
</l>
<!-->
<choice>
  <orig>Какъ</orig>
  <reg>Как</reg>
</choice>
<choice>
  <orig>вѣтеръ</orig>
  <reg>ветер</reg>
</choice> черный, поешь: ты мой!
</l>
```

Additions, deletions, substitutions

Alterations made to the text, whether by the scribe or in some later hand, can be encoded using `<add>` (addition) or `` (deletion).

Where the addition and deletion are regarded as a single act of *substitution*, they can be grouped together using the `<subst>` (substitution) element



```
<subst>
  <del>half-</del>
  <add>all</add>
</subst> blind
```

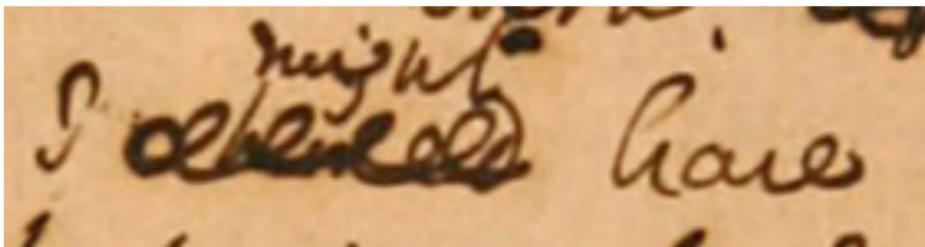
An English example

And towards our distant rest began to trudge,
~~Helping the worst amongst us~~, who'd no boots all
But limped on, blood-shod. All went lame; half-blind;
Drunk with fatigue ; deaf even to the hoots
Of tired, outstripped ~~fif~~ five-nines that dropped behind.

```
<|>And towards our distant rest began to trudge,</|>
<|>
<subst>
  <del>Helping the worst amongst us</del>
  <add>Dragging the worst
        amongst us</add>
</subst>, who'd no boots
</|>
<|>But limped on, blood-shod. All went lame; <subst>
  <del stat us="shortEnd">half-</del>
  <add>all</add>
</subst> blind;</|>
<|>Drunk with fatigue ; deaf even to the hoots</|>
<|>Of tired, outstripped <del>fif</del> five-nines that dropped behind.</|>
```

Semi-legible text

Use `<unclear>` if the text is partly illegible i.e. it can be read but without perfect confidence.



```
I <subst>
  <add place="above">might</add>
  <del>
    <unclear reason="overinking"
      cert="medium" resp="#LDB"> should</unclear>
  </del>
</subst>have
```

Damage to the carrier

Use the `<damage>` element to record the existence of physical damage to the document, whether or not the damaged text is readable :

```
<damage reason="illegible"  
        quantity="3" unit="word"/>Sydney Smith
```

```
<l>The Moving Finger wri<damage agent="water" group="1">es; and</damage>  
having writ,</l>  
<l>Moves  
<damage agent="water" group="1">  
  <supplied>on: nor all  
    your</supplied>  
</damage> Piety nor Wit</l>
```

The `@group` attribute is used to associate together parts of the transcription affected by the same area of damage.

Lacunae

Use the `<gap>` element when something is missing, for example because it is impossible to transcribe, because the carrier is damaged, or because of editorial policy.

In addition to the attributes already mentioned, `@quantity` and `@unit` are available to indicate the size of the gap

```
I am dr Sr yr <gap reason="illegible" quantity="3"  
unit="word"/>Sydney Smith
```

Their arrangement with respect to Jupiter and to each other was as follows:

```
<gap reason="sampling"  
extent="restOfPage">  
  <desc>astrological figure</desc>  
</gap>
```

That is, there were two stars on the easterly side and one to the west; ...

Text made semi-legible through damage

These elements may be used in combination as necessary. In this example, two phrases can be read despite smoke damage, but three lines in between are completely illegible:

```
<damage agent="smoke">  
  <unclear>and the proof of this is</unclear>  
  <gap quantity="3" unit="line"  
    cause="smokeDamage"/>  
  <unclear>margin</unclear>  
</damage>
```

Supplied text

Use the `<supplied>` element if the transcriber has provided a reading not actually visible in the text, perhaps because of scribal error :

...Dragging the worst among`<supplied reason="authorialError" cert="high">s</supplied>`t us...

Attributes can be used to qualify the information further:

- `@reason` why the text has had to be supplied (any word)
- `@source` (if any) from which the text was taken (a pointer)
- `@resp` who is responsible for supplying this markup (a pointer)
- `@cert` the degree of certainty associated with the markup (high, medium, or low)



Text supplied to fill a lacuna

If the transcriber wishes to supply material to fill a lacuna, it should be marked up using `<supplied>` rather than `<gap>`. In this example, 42 lines missing from the bottom of page 5v have been supplied from another source (identified here simply as ed) on the authority of #djb:

расѣдша
се. и ω

```
<lb break="no"/>брѣтени* быти голоубици
<supplied resp="#djb" source="ed"
    reason="missing">
    <pb n="301c"/>
    <lb/>(голоубици) онои въ чре
    <lb break="no"/>вѣ юго. и простъръ роу
    <lb break="no"/>коу блаженныи възѧ
    <lb break="no"/>тъ ю живоу. не имоущу
    <!-- ... -->
    <lb/>днъ и нощь молѧ ба о ню
    <lb break="no"/>и. и по дъвою лѣтоу вѣ
    <lb break="no"/>сть бысть (юмоу. къ-)
</supplied>
<pb n="6r"/>
<lb/>емоу. гдѣ ю<hi rend="sup">с</hi> и како живе
<lb break="no"/>тъ. и оумоливъ нѣкого оу
<lb break="no"/>жикоу своєго посла тамо.
```



Some difficulties

These methods are perfectly adequate where variation is comparatively simple. They rapidly encounter problems when:

- overlap happens (as it always does)
- the sequence of scribal interventions is important
- the layout and the meaning of the writing are not easily separable

How far will the TEI take us ?

In particular, is the TEI scheme adequate for the needs of those transcribing 'modern' manuscripts ?

- surviving medieval or early modern manuscripts generally have a public function, and a more or less conventionalised (if complex) format
- modern manuscripts or authorial drafts however often contain entirely private or idiosyncratic signs, with no clear communicative function



Exercise 1

- Decide on the structure of your encoding
- Decide what you will capture in the transcription and what you will not
- Mark it up!

Look at these files:

- ocsPage.jpg page image
- ocsPage.txt plain text only
- Detailed instructions