

Digitale Edition - Vertiefung und Nutzung













Organisatorisches

- Übungsdateien und Präsentationen
 - Google Drive Link
- Installationen
 - oXygen
 - Anaconda
 - Python-Paket spaCy
 - Python-Paket pydelta:
 - Anaconda prompt: pip install git+https://github.com/cophi-wue/pydelta@next
 - https://github.com/cophi-wue/pydelta
- Wer geht mit in die 10er Marie?
 - Mittwoch ab 16:30 Uhr, auf eigene Kosten

Programmübersicht

https://www.i-d-e.de/aktivitaeten/schools/autumn-school-2018-wien/

Montag	Wiederholung oXygen, XML und TEI Erstellung des Arbeitskorpus XPath, Reguläre Ausdrücke Abendvortrag "Distant Reading Digital Texts" von Jan Rybicki
Dienstag	Basistechnologien für Auswertungsdaten (XSLT, XQuery, CSV, Python) Auswertungen (Geodaten)
Mittwoch	Auswertungen (Netzwerkanalyse, Community Detection, Visualisierung) Heuriger 10er-Marie
Donnerstag	NLP (Vorbereitung von Texten, Part of Speech Tagging, Named Entity Recognition)
Freitag	NLP (Plagiatserkennung, Stilometrie)

Erstellung des Arbeitskorpus I

Wiederholung XML, TEI und oXygen

XML - Extensible Markup Language

- Standard des W3C seit 1998
- XML bezeichnet eine erweiterbare Auszeichnungssprache für die Beschreibung und den Austausch von Daten
- Metasprache zur Ableitung von XML Vokabularen
 - TEI, MEI, SVG, etc.
- XML Dokumente werden nach einem Datenmodell aufgebaut
- XML ist menschen- und maschinenlesbar
- XML trennt Struktur und Darstellung strikt voneinander

Grundregeln von XML

- Jedes XML-Dokument hat genau ein Wurzelelement
- Die Elemente müssen strikt hierarchisch ineinander geschachtelt werden

```
<staedte>
    <stadt>
        <name>Graz</name>
        </stadt>
        <stadt>
        <name>Wien</name>
        </stadt>
        </stadt>
        </stadt>
```

Elemente & Attribute

Elemente

- Korrespondierende Start- und End-Tags
- <element>Inhalt</element>
- o Elemente enthalten Text, Elemente, Text und Elemente oder nichts
- <element />, <element><unterelement>Inhalt</unterelement></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></element></elem

Attribute

- Elemente dürfen beliebige viele Attribute haben
- Elemente dürfen nicht 2x das gleiche Attribute haben
- <element attribut="wert">Inhalt</element>

Dokumentaufbau

Processing Instruction: XML-Deklaration

```
<?xml version="1.0" encoding="UTF-8"?>
```

Processing Instruction: Schemaverweis

```
<?xml-model href="postkarte.rng" type="application/xml"
schematypens="http://relaxng.org/ns/structure/1.0"?>
```

Processing Instruction: Stylesheet

```
<?xml-stylesheet type="text/xsl" href="transformieren.xsl"?>
```

"eigentlichen" Elemente und Attribute

Kommentare

```
<!- Kommentare sind kein Cde -->
```

Namenskonventionen für Elemente

- ... sind frei definierbar
- ... sind case-sensitive <person> ≠ <PERSON> ≠ <Person>
- ... beginnen mit einem Buchstaben oder Unterstrich <person>, <_person>, <1person>
- dürfen nicht mit "xml" beginnen <xmlperson>
- können Buchstaben, Zahlen, Bindestriche, Unterstriche und Punkte enthalten
 <per1son>, <per-son>, <per-son>
- dürfen keine Leerzeichen enthalten <per son>
- Sprechende Namen verwenden, ggf. camelCase-Schreibweise
 <persName>

Wohlgeformtheit & Validität

- XML Dateien werden von einem Parser gelesen
- Ein "wohlgeformtes" XML-Dokument entspricht den Syntaxregeln von XML
- Nur wohlgeformte Dokumente können geparst und verarbeitet werden

Ein "valides" XML Dokument entspricht den Regeln eines assoziierten

Schemas



Datenkontrolle über Schemata

XML Dokumente müssen kein Schema haben, aber ein Schema

- formalisiert die Datenstruktur
- definiert Vokabular und Grammatik einer bestimmten XML Sprache
- definiert, welche Elemente, Attribute und Inhalte in einem XML
 Dokument erlaubt sind
- definiert, wie Elemente und Attribute zu verwenden sind (optional, obligatorisch, Wiederholung, Datentypen, ...)

Entitäten

- Für Zeichen die in XML bereits besetzt sind, gibt es spezielle Entitäten
 - < wird zu <
 - > wird zu >
 - & wird zu & amp;
 - " wird zu "
 - 'wird zu '
- Alternativ können sie als CDATA-Sektionen markiert werden:
 - <![CDATA[Inhalte mit <spitzen> Klammern]]>

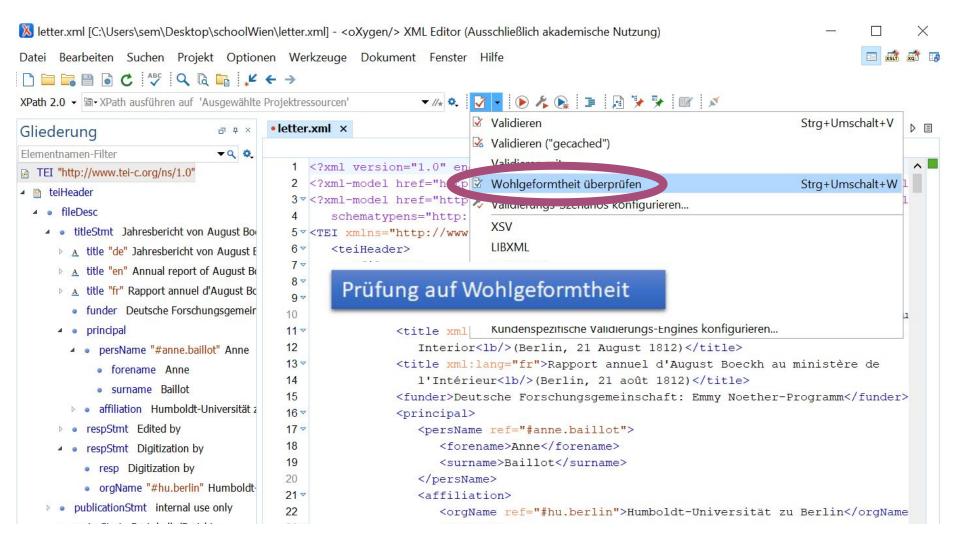
Namensräume

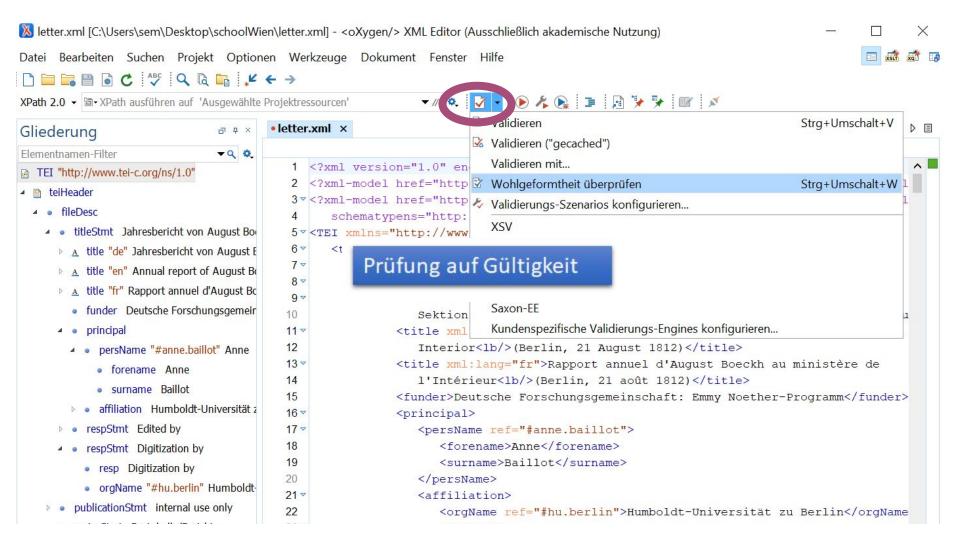
- Element und Attributnamen sind Teil eines Namensraums
- Der Namensraum (Namespace) identifiziert Objekte (Elemente und Attribute) eines verwendeten XML Vokabulars eindeutig
- Damit entstehen keine Namenskonflikte bei gleichlautenden Elementnamen.
- Identifikation des Gegenstandsbereichs über URI

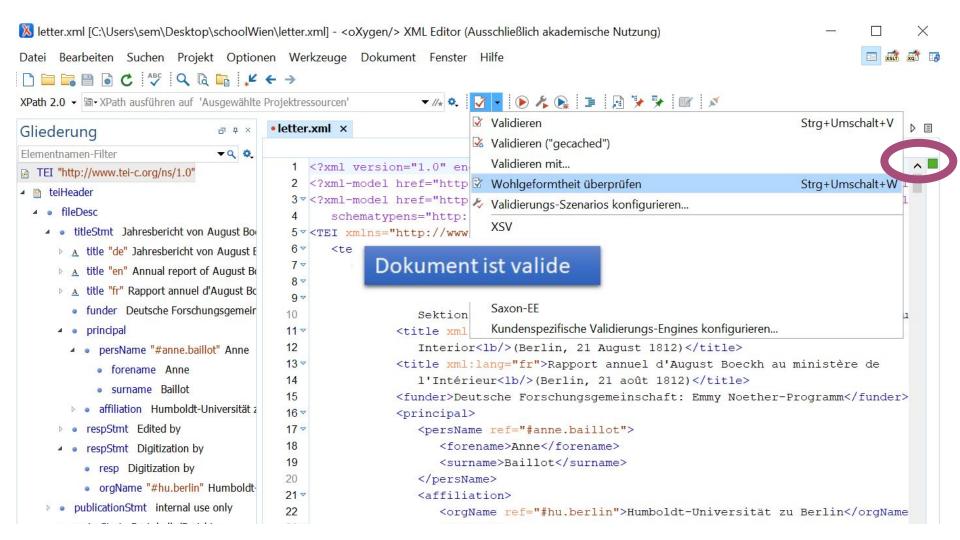
```
<elementname xmlns:Präfix="URI">
```

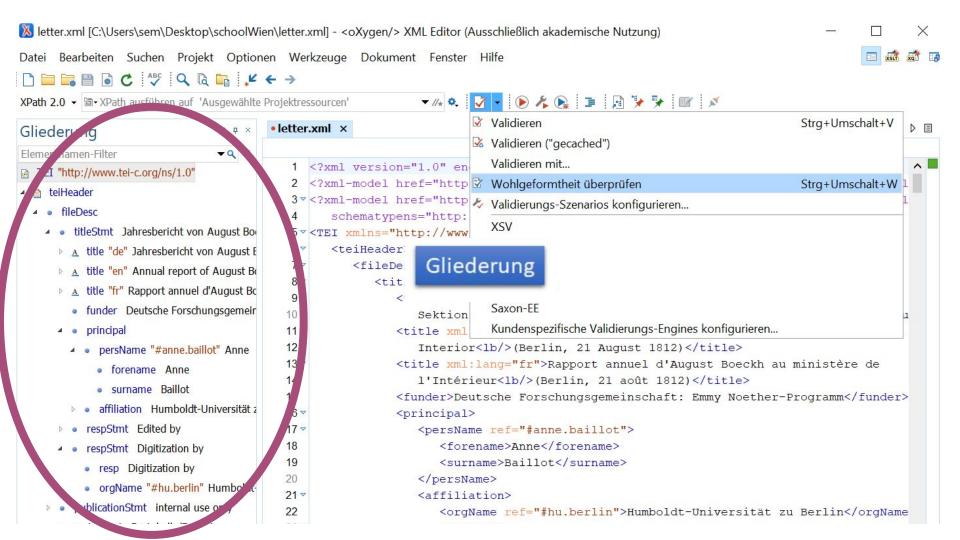
- Die Präfixe können frei vergeben werden Abkürzung
- Namensraum und Schema sind unabhängig voneinander

oXygen XML Editor

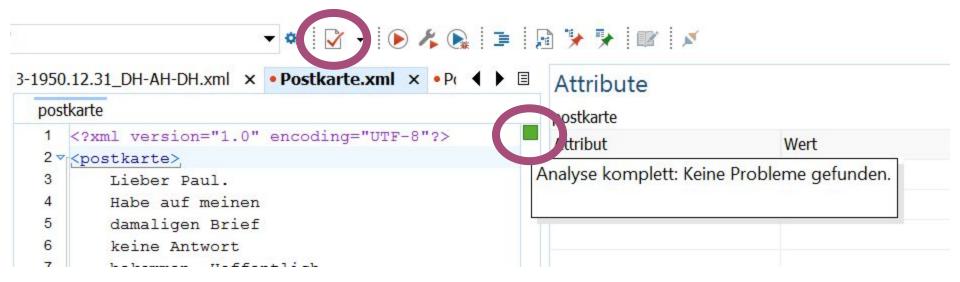






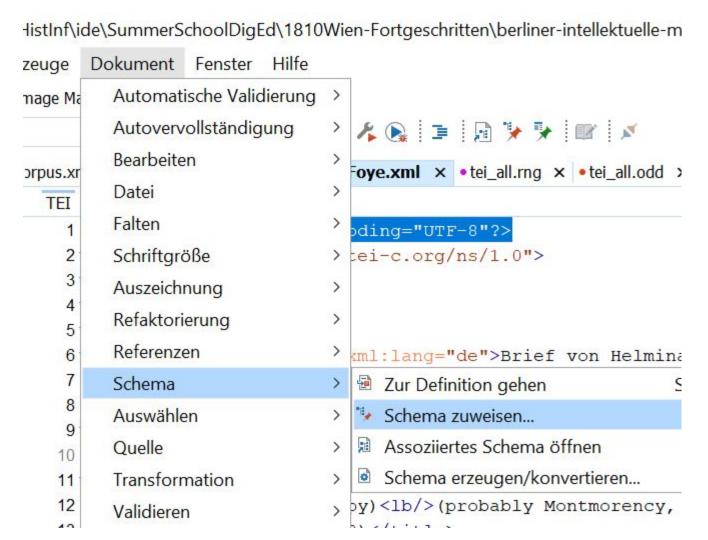


Validität

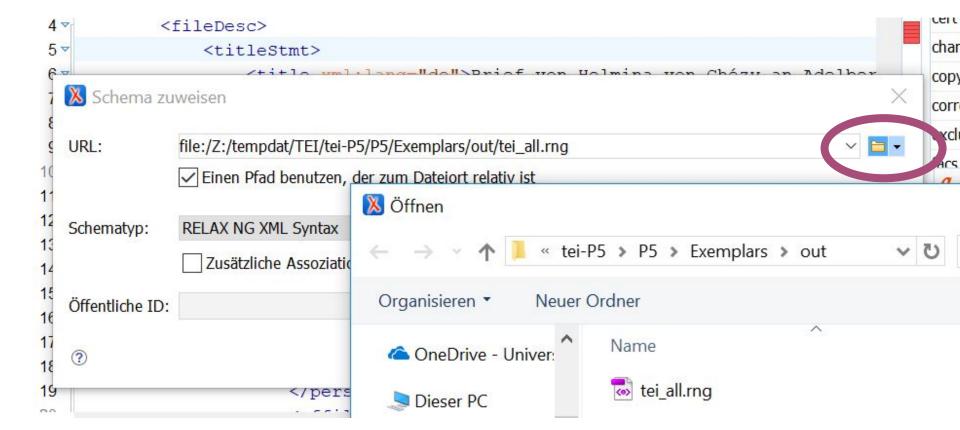


Schema zuweisen

Dokument >
Schema > Schema
zuweisen



Schema zuweisen 2



Integrierte Dokumentation

Integrierte Dokumentation



Wichtige Shortcuts

- Strg/cmd+E: Text mit Element umgeben
- Strg/cmd+Shift+Y: Lange Zeilen umbrechen/nicht umbrechen
- Strg/cmd+Leertaste: Auflisten aller möglichen Elemente an der aktuellen Position
- Strg/cmd+Shift+Beistrich: aktuelle Auswahl in Kommentar umwandeln/Kommentar aufheben.
- Strg/cmd+Shift+P: Dokument formatieren

Übung

- Öffnen Sie Oxygen
- Öffnen Sie die Datei <u>1.10./Brief001ChamissoandeLaFoye.xml</u>
- Laden Sie das Teil All-Schema von 1.10./tei all.rng herunter
- Weisen Sie der XML-Datei dieses Schema zu
- Schauen Sie sich die Fehler an und versuchen Sie, einen davon zu beheben
 - o z.B. wie den letzten: Wie sind die Attribute zu <unclear> geschrieben?

Hilfe zu XML

- W3C Spezifikation: http://www.w3.org/TR/xml/
- W3 Schools Tutorials XML/XSLT/XPath/XQuery: http://www.w3schools.com/xml/



TEI

TEI = Text Encoding Initiative

TEI

Apropos:

Mitalied werden!

- http://www.tei-c.org
- Seit 1986 von Geisteswissenschaftlern entwickelt.
 - Guidelines http://www.tei-c.org/guidelines/p5/
 - darauf aufbauende Schemata http://www.tei-c.org/guidelines/customization/
- derzeit gültige Version: P5, letzte Veröffentlichung: 3.4.0, 23.7.2018
- Vorschläge für 568 Elemente
- Fehler berichten und neue Funktionen vorschlagen:

https://github.com/TEIC/TEI

TEI

- Werkzeuge
 - Stylesheets zur Transformation von TEI-Dokumenten (auch in Oxygen integriert)
 - Roma zur Erstellung eigener Schemata und Dokumentationen
- Aktivitäten
 - TEI Conference and Members Meeting
 - TEI Technical Council, TEI Board of Directors
 - SIGs (Special Interest Groups)
 - Mailingliste http://www.tei-c.org/support/#tei-l
 - Journal of the Text Encoding Initiative: http://journal.tei-c.org

Struktur

23 Module:

- <u>TEI Infrastructure (tei)</u>
- <u>TEI Header (header)</u>
- <u>Elements Available in All TEI Documents (core)</u>
- <u>Default Text Structure (textstructure)</u>
- Für Editionen typisch: representation of primary source (transcr); manuscript description
 (msdescription); names, dates, people and places (namesdates); characters, glyphs, and
 writing modes (gaiji); certainty, precision and responsibility (certainty);

Customization:

- "ODD" (One Document Does it all)
- Erzeugen mit ROMA (http://roma.tei-c.org/, vgl. Folien von 2017)
- oder von Hand

TEI Grundstruktur

Ein TEI Dokument wird durch das Wurzelelement <TEI> repräsentiert, das sowohl Daten als auch Metadaten enthält.

```
<TEI xmlns="http://www.tei-c.org/ns/1.0">
 <teiHeader>
      <!-- obligatorisch: Metadaten zum Dokument -->
 </teiHeader>
 <facsimile>
      <!-- optional: Abbildungen des Dokuments -->
 </facsimile>
 <sourceDoc>
      <!-- optional: dokumentnahe Transkription -->
 </sourceDoc>
 <text>
      <!-- obligatorisch wenn es kein facsimile oder sourceDoc gibt -->
 </text>
</TFI>
```

Das Arbeitscorpus

Berliner Intellektuelle

(ed. A. Baillot et al. 2012ff.)

Inhalt

Bislang unpublizierte Dokumente zum intellektuellen Leben im Berlin des späten 18. und frühen 19. Jahrhunderts

- Briefe
- Dramen/Libretti
- Erzählungen
- Protokolle/Berichte
- Vorlesungsmitschriften/Promotionsschriften
- Rezensionen

Von Adolf Friedrich von Buch, August Boeckh, Adelbert von Chamisso, Helmina von Chézy, Jean Albert Euler, Immanuel Hermann von Fichte, E.T.A. Hoffmann, Adelheid Reinbold, Dorothea Tieck und Ludwig Tieck

Vorgehen im Projekt

- Originalgetreue Wiedergabe, incl. Erhalt unverständlicher Textstellen (Kommentar)
- Wechsel zwischen lateinischer und Kurrentschrift wiedergegeben
- Zuordnung zu Händen
- Kodierungsrichtlinien als PDF: Textbehandlung, Gattungsspezifische Auszeichnung (Briefe, Dramen), Handschriftenbeschreibung, Indizes, Datumskodierung, Template
- Indizes (Orte, Personen, Gruppen, Werke)
- Kein Schema zur Kontrolle angewendet

Das Korpus

Briefe, mit inkonsistentem Briefmarkup

Aufgabe 2: Verbessern Sie das Markup:

- 1. hat jeder Brief eine correspDesc? Sind in dieser auch die Schreiborte der Briefe vermerkt?
- 2. Sind opener, closer und dateline in den Briefen richtig kodiert? Auf die richtigen Elemente angewendet?
- Prüfen Sie, ob in allen Briefen die Transkription in div[@type="transcription"] eingeschlossen ist

Erstellung des Arbeitskorpus II

Praktische Übung

TEI für digitale Editionen

Metadaten

- Die Editoren im teiHeader: <titleStmt><editor>...</editor></titleStmt>
- Die Quelle der Edition im teiHeader: <sourceDesc>
- gedruckte Quellen: <bibl> oder <biblStruct>
- handschriftliche Quellen (oder individuelle Exemplare von Drucken):
 <msDesc>
- Was kann man alles über Handschriften sagen?
 - Identifikation: <msIdentifier> mit <institution> ... </institution> <settlement> ... </settlement>
 <idno> ... </ido> oder <msName> ... <msName>
 - physische Eigenschaften: <physDesc> mit <objectDesc> für Beschreibstoff (<supportDesc>)
 etc., <handDesc> für die Beschriftung, <history> für die Überlieferungsgeschichte <additional> für Hinzufügungen.

Transkription

cert="medium">...</unclear></choice>

```
Tilgungen und Hinzufügungen: <del>, <add>, Ersetzungen: <subst><del>...</del><add>...</add></subst>

Lesbarkeit: <unclear>, <damage>

Abkürzungen und Normalisierungen: <choice>: (<abbr>, <expan>), (<am>, <ex>), (<sic>, <corr>), (<orig>, <reg>), aber auch für alternative unklare Lesarten: <choice><unclear cert="medium">...</unclear><unclear
```

Übung

- Öffnen Sie oXygen
- Suchen Sie sich einen Brieftext aus dem Ordner "1.10./txt" aus
- Erstellen Sie ein neues TEI All Dokument
- Fügen Sie den Briefinhalt ein
- Kodieren Sie den Titel und die Handschriftenbeschreibung zum Brief
- Kodieren Sie die Transkription

TEI für Korrespondenz

Nach einem Foliensatz von Sabine Seifert und Peter Stadler

TEI für Korrespondenzen: Projekte

- Digital Archive of Letters in Flanders (DALF) http://ctb.kantl.be/project/dalf/index.htm
- Van Gogh. The Letters http://vangoghletters.org/
- Alfred Escher Briefedition http://www.briefedition.alfred-escher.ch/
- Carl Maria von Weber Gesamtausgabe (WeGA) http://www.weber-gesamtausgabe.de/
- Berliner Intellektuelle um 1800 http://www.berliner-intellektuelle.eu/?en
- Letters of 1916 http://letters1916.ie/
- Burckhardt Source http://burckhardtsource.org/
- August Wilhelm Schlegel's Korrespondenz
 http://august-wilhelm-schlegel.de/briefedigital/

TEI für Korrespondenzen

Infrastrukturen

- correspSearch http://correspsearch.net
- Early Modern Letters Online (EMLO) http://emlo.bodleian.ox.ac.uk/

Visualisierungen

- Mapping the Republic of Letters http://republicofletters.stanford.edu/
- Visual Correspondence http://letters.nialloleary.ie/
- Nodegoat http://nodegoat.net/

TEI für Korrespondenzen: Theorie

Ein Brief kann als **Objekt** und als **Event** betrachtet werden:

- Objekthaftigkeit
 - Materialität (Papier, Stempel, Wasserzeichen, Schrift etc.)
 - TEI manuscript description <msDesc>
- Eventcharakter
 - Kommunikation (Absender, Empfänger, Absendeort, Datum etc.)
 - TEI header <correspDesc>

Zentrale Fragen zu Korrespondenzen

- Wer hat geschrieben
- An wen wurde geschrieben
- Von wo aus wurde geschrieben
- Wohin wurde geschrieben
- Wo wurde die Nachricht empfangen
- Wann wurde die Nachricht geschrieben
- Wann wurde die Nachricht empfangen
- Was wurde geantwortet
- Wie wurde die Nachricht übermittelt



<person>, <persName>, <org>, <orgName> etc.



<place>, <placeName>, <settlement> etc.

Datierungen

<date>, <time>

Vorangegangene und nachfolgende

Nachrichten

<correspAction>, XXX

Metadaten zu Korrespondenzen

<correspDesc>

Enthält Metadaten, die eine Korrespondenz als Event bzw. Aktion beschreibt.

<msDesc>

Enthält Metainformationen zu den physischen Eigenschaften des Manuskripts.

Eine vollständige Beschreibung der Korrespondenz ergibt sich aus <correspDesc> und <msDesc>

Beispiel zu <correspDesc>

```
cprofileDesc>
  <correspDesc>
    <correspAction type="sent">
      <persName>Carl Maria von Weber</persName>
       <settlement>Dresden</settlement>
       <date when="1817-06-23">23 June 1817</date>
    </correspAction>
    <correspAction type="received">
      <persName>Caroline Brandt</persName>
       <settlement>Prag</settlement>
       </correspAction>
  </correspDesc>
</profileDesc>
```

Beispiel zu <correspContext>

```
cprofileDesc>
  <correspDesc>
     <!-- correspAction -->
     <correspContext>
       <ref type="prev" target="http://www.weber-gesamtausgabe.de/A041209">Previous letter of
          <persName>Carl Maria von Weber</persName> to <persName>Caroline Brandt</persName>:
          <date from="1817-06-19" to="1817-06-20">June 19/20, 1817</date>
       </ref>
       <ref type="next" target="http://www.weber-gesamtausgabe.de/A041217">Next letter of
          <persName>Carl Maria von Weber</persName> to <persName>Caroline Brandt</persName>:
          <date when="1817-06-27">June 27, 1817</date>
       </ref>
     </correspContext>
  </correspDesc>
```

<correspAction>

contains a structured description of the place, the name of a person/organization and the date related to the sending/receiving of a message or any other action related to the correspondence.

- http://www.tei-c.org/release/doc/tei-p5-doc/en/html/ref-correspAction.html
- Über das @type Attribut wird die Natur der Aktion beschrieben
- Empfohlene Werte: sent, received, transmitted, redirected, forwarded

Übung

- Öffnen Sie oXygen
- Suchen Sie sich Ihre Datei von der letzten Übung oder wählen Sie einen Brieftext aus dem Ordner "1.10./txt" aus, den wie in der vorherigen Übung in eine TEI All-Datei übertragen.
- Kodieren Sie die Metadaten des Briefes

Briefinhalte kodieren: Beginn und Ende

<opener>

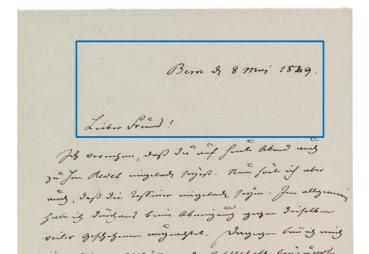
fasst Datumszeile, Verfasserangabe, Anredeformel und ähnliche Phrasen zusammen, die einleitend zu Beginn eines Abschnitts stehen, vor allem bei einem Brief.

<closer>

fasst Grußformeln, Datumszeilen und ähnliche Phrasen zusammen, die am Ende eines Abschnitts stehen, vor allem bei einem Brief.

Beginn eines Briefes

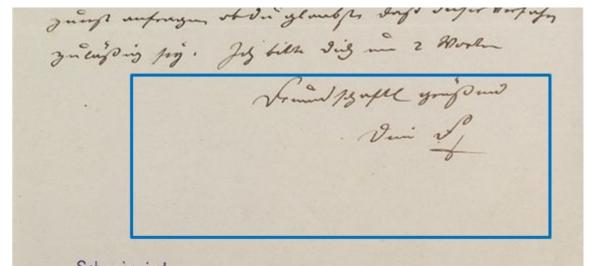
```
<opener>
 <dateline rend="right">
    <placeName>Bern</placeName>
    den <date when="1849-05-08">8 Mai 1849</date>.
 </dateline>
 <salute>Lieber Freund!</salute>
</opener>
Brieftext ..
```



Ende eines Briefes

```
Brieftext ..
<closer rend="right">
  <salute>Freundschaftlich grüßend Dein</salute>
  <signed><persName>F J</persName></signed>
```

</closer>



Postskriptum

```
<postscript>
```

Enthält ein Postskriptum, z.B. bei einem Brief

```
<div type="letter">
 <opener> ... </opener>
  ... 
 <closer> ... </closer>
 <postscript>
     <label>P.S.:</label>
     Leider konnte ich nicht zu deinem Geburtstag kommen ...
 </postscript>
</div>
```

Vorgedruckter Briefkopf

```
<div type="printed_letterhead">
    MAISON DE SANTÉ<1b/>
    de<1b/>
    SAINT-RÉMY<1b/>
    DE PROVENCE<1b/>
    Bouches-du-Rhône
</div>
```



Letter from Théophile Peyron to Vincent van Gogh Saint-Rémy-de-Provence, 1 July 1890 http://vangoghletters.org/vg/letters/let895/letter.html

Briefumschlag

Informationen auf dem Briefumschlag werden im TEI-Header untergebracht:
 <accMat> in <physDesc>

<accMat> (accompanying material): contains details of any significant additional material which may be closely associated with the manuscript being described, such as non-contemporaneous documents or fragments bound in with the manuscript at some earlier historical period.

```
<msDesc>
     <physDesc>
     <accMat>Der Umschlag enthält den Namen des Empfängers ... </accMat>
     </physDesc>
</msDesc>
```

Briefumschlag

Kodierung im <text>-Element; eigenes <div> verwenden

```
<div type="envelope">
                                                          >
<div type="envelope">
                                                             <address>
  >
                                                                <orgName>Università di Bologna</orgName>
     <address>
                                                                <placeName type="country">Italy</placeName>
        <addrLine>To the Dean of St Hugh's College</addrLine>
                                                                <postCode>40126</postCode>
        <addrLine>in Oxford</addrLine>
                                                                <placeName type="city">Bologna</placeName>
     </address>
                                                                <street>via Marsala 24</street>
  </address>
</div>
                                                          </div>
```

Stempel, Siegel und Wasserzeichen

<stamp> contains a word or phrase describing a stamp or similar device. (http://www.tei-c.org/release/doc/tei-p5-doc/en/html/ref-stamp.html)

<seal> contains a description of one seal or similar attachment applied to a manuscript. (http://www.tei-c.org/release/doc/tei-p5-doc/en/html/ref-seal.html)

<watermark> contains a word or phrase describing a watermark or similar device.
(http://www.tei-c.org/release/doc/tei-p5-doc/en/html/ref-watermark.html)

Übung

- Öffnen Sie oXygen
- Suchen Sie sich einen Brieftext aus dem Ordner "1.10./txt" aus
- Öffnen Sie ein neues TEI All Dokument
- Fügen Sie den Briefinhalt ein
- Kodieren Sie die briefspezifischen Metadaten zum Brief

Personen, Orte, Organisationen, Themen

```
<persName> != <person> (desgl. <placeName> | <place>, <orgname> | <org>)
```

Beziehungen zwischen Name und formaler Beschreibung über <persName ref="#myID">...</persName> ... listPerson><person xml:id="myID"><persName>...</persName> ...</persName> ...</person>

Themen / Schlagwörter: <term>...</term>, wenn nicht im Text vorkommend: <index><term>...</term></index>

Übung

- Öffnen Sie oXygen
- Suchen Sie sich einen Brieftext aus dem Ordner "1.10./txt" aus
- Öffnen Sie ein neues TEI All Dokument
- Fügen Sie den Briefinhalt ein
- Kodieren Sie Personen, Orte, Organisationen und Sachthemen

Übung - ODD Customization

- Öffnen Sie oXygen
- Erstellen Sie eine neue ODD Customization (Neu > Framework Vorlagen > ODD Customization)
- Fügen Sie notwendige Module hinzu
- Transformieren Sie die ODD in eine RelaxNG-Datei
- Verknüpfen Sie Ihr TEI Dokument mit dem Schema
- Validieren Sie das Dokument
- Ergänzen Sie fehlende Module



