

# Automatisk konvertering til TEI-format



## Introduktion

I denne "tutorial" gennemgår vi hvordan du automatisk kan konvertere dine tekstfiler i txt-format eller RTF-format til xml-formatet TEIP5. Konverteringen følger det TEI-skema der blev defineret i DK-CLARIN-projektet 2008-2011. (Se dokumentation her: <https://info.clarin.dk/clarin-dk-infrastrukturen/tutorials/text-format.pdf> og <https://info.clarin.dk/clarin-dk-infrastrukturen/tutorials/text-header.pdf> )

I xml-filerne er al layout fjernet og teksten gemt med et ord på hver linje, linjerne (dvs. ordene) nummererede og afsnittene markerede. Pga. denne strukturering er det muligt at lave automatisk analyse af teksten; men ikke længere muligt at "læse" teksten som før. Tekstfilen er blevet et forskningsobjekt.

Efter klargøringen skal du selv verificere kvaliteten af den automatiske proces, fx om det umiddelbart ser ud til at konverteringen ikke har ødelagt ordene.

## Trin 1

Gå til NLP Tools (Værktøjskassen) – Klargøring af tekstresurser <https://clarin.dk/clarindk/toolchains-upload.jsp>, vælg den fil du vil have konverteret til TEI-formatet, og vælg filformat txt (flad) eller RTF.

### Værktøjskasse

Klargøring af tekstresurser | Klargøring af dataresurser | Lingvistisk opmærkning

#### Upload filer og konverter til TEI

Vælg en eller flere filer som skal opmærkes. Alle valgte filer skal have samme filformat (pdf, rtf, eller txt).

File(r):  Der er ikke valgt nogen fil  ▼

#### Tilføj metadata

**Kildetekst:**

Værktitel \*

Titel på større værk (hvis værk er en del af dette større værk) \*

Forfatter \*

Udgiver

Udgivelsesdato

Dato for værkets tilblivelse

Sprog \*  ▼

Domæne (English:Subject)

**Elektroniske version:**

Titel på digitalt værk \*

Den ansvarlige for skabelsen af den digitale version \*

Den ansvarlige for distribution af den digitale version \*

Sponsor

Konverteringsdato (automatisk genereret)

Beskrivelse \*

Størrelse (automatisk genereret)

Id \*

## Trin 2

Derefter skal du udfylde de mest basale metadata:

\* betyder at metadataelementet er obligatorisk

Metadata	Betydning	Eksempel
*Værktitel	Titlen på kildeteksten. Evt. e-navn hvis resursen er skabt elektronisk	<i>Foraarets død</i>
Titel på større værk (hvis værk er en del af dette større værk)	Titlen på den del af værket som filen dækker, evt. samme som Værktitel hvis filen ikke hører under et samlet hele	<i>Kongens Fald</i>
*Forfatter	Forfatter af kildeteksten, kan evt. være anonym eller ukendt. Hvis der er flere forfattere, skrives en på hver linje	<i>Johannes V. Jensen</i>
Udgiver	Navnet på den organisation der er ansvarlig for udgivelsen eller distribution af værket. Evt. ens egen organisation.	<i>Gyldendal</i>
Udgivelsesdato	Format yyyy-mm-dd eller yyyy hvis den præcise dato ikke er kendt	<i>1900</i>
Dato for værkets tilblivelse	Bruges hvis man ønsker at angive at kildeteksten er skabt på et andet tidspunkt end udgivelsesdato. Format yyyy-mm-dd eller yyyy hvis den præcise dato ikke er kendt	<i>1898</i>
*Sprog	Sproget i den udgivelse der her uploades udtrykt i ISO639-1 Se: <a href="https://en.wikipedia.org/wiki/List_of_ISO_639-1_codes">https://en.wikipedia.org/wiki/List_of_ISO_639-1_codes</a>	<i>da</i>
Domæne	Emne eller domæne, gerne med henvisning til et anerkendt klassifikationssystem som DK5 eller DDC der bruges af biblioteksvæsner. Emneangivelse er oftest relevant når det drejer sig om fagtekst	<i>Dansk skønlitteratur (DK5-86)</i>
*Titel på digitalt værk	Titlen på den elektroniske version. Kan være identisk med Værktitel hvis resursen er skabt elektronisk eller man ønsker at resursen skal være synlig i CLARIN under Værktitlen	<i>Kongens Fald - TEIP5-version</i>
*Personansvarlig for digital version	Person der er ansvarlig for <i>skabelse</i> af den elektroniske version. Hvis der er flere, skrives en på hver linje.	<i>Peter Petersen</i>
*Ansvarlig for distribution af digital version	Person eller organisation der er ansvarlig for <i>distribution</i> af den elektroniske version. Hvis der er flere, skrives en på hver linje.	<i>NorS, Københavns Universitet</i>
Sponsor	Her kan angives projekt, institution, bevilling hvis man skønner at det er relevant.	<i>CLARIN-DK-projektet</i>
Konverteringsdato	Dato for digital version af værket i formatet yyyy-mm-dd eller yyyy. Udregnes og tilskrives automatisk.	<i>2011</i>
*Beskrivelse	Kommentarer om hvilke forskningsammenhænge teksten kan bruges i, eller andre oplysninger det er relevant at huske eller at dele med andre	<i>Kongens fald i xml TEIP5 format er en del af JVJ online tilvejebragt gennem DK-CLARIN-projektet til brug i tekstanalyse</i>
Størrelse	Værkets størrelse. Udregnes og tilskrives automatisk.	<i>17.563 words</i>
Id	Tilskrives automatisk.	<i>2100106490</i>

Når du har klikket *Submit* genereres tekstfilen i TEIP5-formatet. Der genereres desuden en segmenteret fil som er et mellemtrin og en index-fil der beskriver den tekniske tilblivelsesproces.

### Trin 3

Du kan se indholdet i en browser og du kan redigere filen med en tekst editor som f.eks. *Notepad* eller med en XML-editor som f.eks. 'Oxygen' eller Microsofts 'Visual Studio'. Hvis du åbner filen i fx Word bliver indholdet fortolket, dvs. at alle xml-koderne er usynlige.

Metadataene vil se således ud:

```
<teiHeader type="text">
  <fileDesc>
    <titleStmt>
      <title>Kongens Fald - TEIP5-version</title>
      <sponsor>CLARIN-DK-projektet</sponsor>
      <respStmt>
        <resp>a_annotation</resp>
        <name><name>Peter Petersen</name>
          <note type="method">flat2cbf</note>
          <date when="2016-06-28"/></name>
        </resp>
      </respStmt>
    </titleStmt>
    <extent><num type="words">1100</num></extent>
    <publicationStmt>
      <distributor>Nordisk Forskningsinstitut, Københavns Universitet</distributor>
      <idno type="ctb">20160609-830-step3</idno>
      <availability status="free">
        <ab type="public"/>
      </availability>
    </publicationStmt>
    <notesStmt>
      <note>Kongens fald i xml TEIP5 format er en del af JVJ online tilvejebragt gennem DK-CLARIN-projektet til brug i
      tekstanalyse</note>
    </notesStmt>
    <sourceDesc>
      <biblStruct>
        <analytic>
          <title xml:lang="da">Foraarets død</title>
          <author><name>Johannes V. Jensen</name></author>
        </analytic>
        <monogr>
          <title>Kongens Fald</title>
          <imprint>
            <publisher n="n/a">Gyldendal</publisher>
            <date when="1900"/>
          </imprint>
        </monogr>
      </biblStruct>
    </sourceDesc>
  </fileDesc>
  osv. ....
</teiHeader>
```

Og tekstdelen således:

```
<body>
<p n="1">
```

```
</p>
<p n="2">
    <w xml:id="i1.1">1</w>
    <w xml:id="i2.1">Vejen</w>
    <c xml:id="i2.2" type="s"/>
    <w xml:id="i2.3">bøjede</w>
    <c xml:id="i2.4" type="s"/>
    <w xml:id="i2.5">tilvenstre</w>
    <c xml:id="i2.6" type="s"/>
    <w xml:id="i2.7">over</w>
    <c xml:id="i2.8" type="s"/>
    <w xml:id="i2.9">en</w>
    <c xml:id="i2.10" type="s"/>
    <w xml:id="i2.11">Bro</w>
    osv. ....
</body>
```

Hvis det er klart at processen ikke er lykkedes, kan man gå ind i en af filerne rtf-filen eller txt-filen, rette de ting der er gået galt og uploade filen igen startende ved Trin 1. Metadata skal desværre indtastes igen; men de gamle metadataværdier kommer frem som forslag når man taster det første bogstave.

Det er ikke tilrådeligt at rette i tekstdelen i xml-filen da hvert enkelt ord er unikt nummereret hvilket er påkrævet i den videre lingvistiske processering. Metadataene kan man dog rolig rette i.