

<b>Title</b>	Overview of historical corpora
<b>Version</b>	1.0
<b>Author(s)</b>	Darja Fišer, Jakob Lenardič
<b>Date</b>	06-06-2018
<b>Status</b>	For distribution
<b>Distribution</b>	BoD, NCF, UI
<b>ID</b>	CE-2018-1237



## Table of contents

1. Introduction .....	1
2. Historical corpora in the CLARIN infrastructure .....	2
2.1. Summary .....	11
2.1.1. Identification .....	11
2.1.2. Availability.....	11
2.1.3. Metadata.....	13
3. Non-CLARIN historical corpora .....	16
3.1. Summary .....	19
3.1.1. Availability.....	19
3.1.2. Metadata.....	20
4. Linguistic tools for the annotation of historical corpora .....	21

## 1. Introduction

In the following report, we present an overview of historical corpora, primarily focusing on those that are part of the CLARIN infrastructure (i.e., they are either listed in the VLO or in the repositories of the national consortia).

The report was conducted in two steps:

- (i) manually searching the VLO and the national consortia with keywords like “historic\* corpus”, “medieval corpus”, etc.
- (ii) input provided by CLARIN UI and NC coordinators

The full results are available in a Google Docs Spreadsheet.<sup>1</sup> In total, around 90 historical corpora were identified. Information on most of the corpora was provided by UI and NC coordinators, whom we would like to thank for their invaluable input. In Section 2, we provide a comprehensive list of the historical corpora that are part of the CLARIN infrastructure, describing their identification (i.e., listed in the VLO or not), their availability (download or through a concordancer), and their metadata (language, size, annotation, license). In section 3, we provide a list of historical corpora are available outside the CLARIN infrastructure. In

<sup>1</sup> [https://docs.google.com/spreadsheets/d/1orDix-eyoHE3K2H44HVNs3ttA\\_SmdDwzCrFKtcJPV8E/edit](https://docs.google.com/spreadsheets/d/1orDix-eyoHE3K2H44HVNs3ttA_SmdDwzCrFKtcJPV8E/edit)

section 4, we provide a short list of tools and services that were provided by participants in the Google forms survey.

## 2. Historical corpora in the CLARIN infrastructure

Table 1 lists 68 historical corpora that are part of the CLARIN infrastructure

Table 1: historical corpora in the CLARIN infrastructure

Corpus	Language	Description
<a href="#">Open Richly Annotated Cuneiform Corpus, Korp Version</a> <b>Size:</b> 741,100 tokens <b>Annotation:</b> tokenised <b>Licence:</b> CC-BY-SA	Akkadian	This corpus contains cuneiform texts from Ancient history.  The corpus is available through the concordancer Korp.
<a href="#">Greek Medieval Texts</a> <b>Size:</b> 3.4 million words <b>Licence:</b> CC-BY	Ancient Greek	This corpus contains texts from the 4 <sup>th</sup> to the 16 <sup>th</sup> century.  The corpus is available for download from the clarin:el repository.
<a href="#">Sheffield Corpus of Chinese</a> <b>Licence:</b> CC-BY-NC-SA 3.0	Chinese	This corpus contains fictional and non-fictional texts from the Medieval and Modern Chinese periods.  The corpus is available for download from the Oxford Text Archive.
<a href="#">"PoDiLemma" Middle Polish Diachrone Lemmatised Corpus</a> <b>Annotation:</b> tokenised, lemmatised	Czech, German, Latin, Polish	This corpus contains political, religious and scientific texts from the 16 <sup>th</sup> to the 18 <sup>th</sup> century.  The corpus is available for download from the CLARIN-D repository.
<a href="#">Medieval Charter Sections Corpus</a> <b>Size:</b> 57 chapters <b>Annotation:</b> manually-tagged, named entities <b>Licence:</b> CC-BY-NC-SA 4.0	Czech, Latin	This corpus contains Latin charters created in the era of John the Bling, King of Bohemia.  The corpus is available for download from LINDAT.
<a href="#">Brieven als buit (Letters as loot)</a> <b>Size:</b> 460,000 words <b>Annotation:</b> lemmatised, PoS-tagged, grammatically tagged	Dutch	This corpus contains 1,000 letters from the 17 <sup>th</sup> to the 18 <sup>th</sup> century.  The corpus is available through a dedicated concordancer.
<a href="#">Corpus Gysseling</a> <b>Size:</b> 1.5 million words <b>Annotation:</b> PoS-tagged, lemmatised <b>Licence:</b> INT Licence for researchers	Dutch	This corpus contains texts from the 13 <sup>th</sup> century.  The corpus is available for download from the Instituut voor de Nederlandse Taal and through a dedicated concordancer.

<p><a href="#">A Corpus of English Dialogues 1560-1760 (CED)</a></p> <p><b>Size:</b> 1.2 million words  <b>Licence:</b> Oxford Text Archive licence</p>	English	<p>This corpus contains dialogues from literary and didactic works from 1560 to 1760.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">Corpus of Early English Correspondence Sampler (CEECS)</a></p> <p><b>Size:</b> 450,000 words  <b>Annotation:</b> none  <b>Licence:</b> Oxford Text Archive licence</p>	English	<p>This corpus contains 1147 letters from 1418 to 1680.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">Corpus of Late Modern English prose / David Denison</a></p> <p><b>Annotation:</b> none  <b>Licence:</b> Oxford Text Archive licence</p>	English	<p>This corpus contains fictional texts from 1837 to 1926.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">Hansard Corpus</a></p> <p><b>Size:</b> 1.6 billion tokens  <b>Annotation:</b> tokenised, PoS-tagged, lemmatised, semantic tags</p>	English	<p>This corpus contains parliamentary debates from 1803 to 2005.</p> <p>The corpus is available through a dedicated concordancer.</p>
<p><a href="#">Helsinki Corpus of Scottish Correspondence (1540-1750)</a></p> <p><b>Size:</b> 0.5 million tokens  <b>Annotation:</b> tokenised  <b>Licence:</b> CLARIN ACA</p>	English	<p>This corpus contains personal correspondence from 1540 to 1750.</p> <p>The corpus is available through the concordancer Korp.</p>
<p><a href="#">Older Scottish texts : the Edinburgh DOST corpus / A.J. Aitken, Paul Bratley and Neil Hamilton-Smith</a></p> <p><b>Size:</b> 877,000 tokens  <b>Annotation:</b> tokenised  <b>Licence:</b> CC-BY-NC-SA 3.0</p>	English	<p>This corpus contains texts from 1450 to 1600.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">Pamphlets of the American Revolution : [selections] / edited by Bernard Bailyn</a></p> <p><b>Annotation:</b> none  <b>Licence:</b> CC-BY-NC-SA 3.0</p>	English	<p>This corpus contains pamphlets of the American Revolution from 1750 to 1776.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>

<p><a href="#">Parsed Corpus of Early English Correspondence (PCEEC)</a></p> <p><b>Size:</b> 2.2 million words  <b>Annotation:</b> tokenised, PoS-tagged, syntactically parsed  <b>Licence:</b> Oxford Text Archive licence</p>	English	<p>This corpus contains correspondence from around 1410 to 1681.</p> <p>This corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">The English language of the north-west in the late Modern English period: a Corpus of late 18c Prose</a></p> <p><b>Size:</b> 300,000 words  <b>Annotation:</b> COCOA-style  <b>Licence:</b> Oxford Text Archive licence</p>	English	<p>This corpus contains texts from 1761 to 190.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">The Lampeter Corpus of Early Modern English Tracts</a></p> <p><b>Annotation:</b> none  <b>Licence:</b> CC-BY-NC-SA 3.0</p>	English	<p>This corpus contains tracts from 1640 to 1740.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">The Lancaster Newsbooks Corpus</a></p> <p><b>Licence:</b> CC-BY-NC-SA 3.0</p>	English	<p>This corpus contains two collections of English newsbooks from 1654 to 1655.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">Corpus of Historical American English - Kielipankki Korp version 2017H1</a></p> <p><b>Size:</b> 385 million tokens  <b>Annotation:</b> tokenised  <b>Licence:</b> CLARN ACA</p>	English (American)	<p>This corpus contains texts from 1810 to 2009.</p> <p>The corpus is available through the concordancer <i>Korp</i>.</p>
<p><a href="#">The Old Bailey Corpus</a></p> <p><b>Size:</b> 134 million words  <b>Annotation:</b> detailed sociobiographical, pragmatic and textual annotation  <b>Licence:</b> CC-BY-NC-SA 4.0</p>	English (Late Modern)	<p>This corpus contains proceedings of the Old Bailey (i.e., legal documents) from 1674 to 1913.</p> <p>The corpus is available for download from the CLARIN-D repository and through the CQPConcordancer.</p>
<p><a href="#">Anthology of Middle English texts / Santiago Gonzalez y Fernandez-Corugedo</a></p> <p><b>Size:</b> 4000 words  <b>Annotation:</b> none  <b>Licence:</b> Oxford Text Archive licence</p>	English (Middle), Hebrew	<p>This corpus contains literary texts from 1100 to 1400.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>

<p><a href="#">Helsinki corpus of English texts</a></p> <p><b>Size:</b> 240,000 words  <b>Annotation:</b> none  <b>Licence:</b> Oxford Text Archive licence</p>	<p>English (Old and Middle)</p>	<p>This corpus contains Biblical and fictional texts from 730 to 1710.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">The York-Helsinki parsed corpus of Old English poetry (YCOEP)</a></p> <p><b>Size:</b> 71,500 words  <b>Annotation:</b> syntactically-parsed  <b>Licence:</b> Oxford Text Archive licence</p>	<p>English (Old)</p>	<p>This corpus contains poems from 730 to 1710.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">Dictionary of Old English Corpus in Electronic Form (DOEC)</a></p> <p><b>Licence:</b> Oxford Text Archive licence</p>	<p>English (Old), Latin</p>	<p>This corpus contains 3037 texts from 600 to 1150.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">The York-Toronto-Helsinki Parsed Corpus of Old English prose (YCOE)</a></p> <p><b>Size:</b> 1.5 million words  <b>Annotation:</b> syntactically-parsed  <b>Licence:</b> Oxford Text Archive licence</p>	<p>English (Old), Latin</p>	<p>This corpus contains fictional texts from 600 to 1150.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">Hamburg Corpus of Old Swedish with Syntactic Annotations (HaCOSSA)</a></p> <p><b>Size:</b> 128,000 words  <b>Annotation:</b> MSD-tagged, syntactically parsed  <b>Licence:</b> CLARIN RES</p>	<p>English, German, Latin, Old Norse, Swedish</p>	<p>This corpus contains texts written in the Late Old Swedish period (from 1375 to 1550).</p> <p>The corpus is available for download from the repository of the University of Hamburg.</p>
<p><a href="#">The Electronic Text Corpus of Sumerian Literature. Revised edition.</a></p> <p><b>Annotation:</b> Each word form in the composite transliterations has been assigned to a lexeme which is specified by a citation form, word class information and basic English translation.  <b>Licence:</b> CC-BY-NC-SA 3.0</p>	<p>English, Sumerian</p>	<p>This corpus contains transliterations and English translations of 394 Sumerian compositions from approximately 2100 to 1700 BCE.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>

<p><a href="#">Corpus of Old Written Estonian</a></p> <p><b>Size:</b> 2 million tokens  <b>Annotation:</b> tokenised, 16.-18. century texts have been tagged with contemporary Estonian, morphological and language information. 19. century texts are unannotated.  <b>Licence:</b> CC-BY</p>	<p>Estonian</p>	<p>This corpus covers secular and religious texts from the 16<sup>th</sup> to the 18<sup>th</sup> century.</p> <p>The corpus is available through a dedicated concordancer.</p>
<p><a href="#">Classics of Finnish Literature, Kielipankki Version</a></p> <p><b>Size:</b> 1.5 million words  <b>Licence:</b> EUPL v.1.1 SA</p>	<p>Finnish</p>	<p>This corpus contains literary texts from 1880 to 1949.</p> <p>The corpus is available through the concordancer Korp.</p>
<p><a href="#">Corpus of Old Literary Finnish</a></p> <p><b>Size:</b> 4.1 million words  <b>Licence:</b> EUPL v.1.1 SA</p>	<p>Finnish</p>	<p>This corpus contains literary texts from 1543 to 1810.</p> <p>The corpus is available through the concordancer Korp.</p>
<p><a href="#">The Finnish Gutenberg Corpus</a></p> <p><b>Size:</b> 34.5 million words  <b>Annotation:</b> no linguistic annotation  <b>Licence:</b> CC-BY</p>	<p>Finnish</p>	<p>This corpus contains books published up to 1925 that are made available through the Gutenberg project.</p> <p>The corpus is available through the concordancer Korp.</p>
<p><a href="#">The Finnish Sub-corpus of the Newspaper and Periodical Corpus of the National Library of Finland, Kielipankki Version</a></p> <p><b>Size:</b> 5.2 billion tokens  <b>Annotation:</b> tokenised  <b>Licence:</b> CC-BY-SA</p>	<p>Finnish</p>	<p>This corpus contains newspaper articles from 1840 to 2011.</p> <p>The corpus is available through the concordancer Korp.</p>
<p><a href="#">The Morpho-Syntactic Database of Mikael Agricola's Works</a></p> <p><b>Size:</b> 428,300 tokens  <b>Annotation:</b> tokenised, PoS-tagged, morphological components and syntactic function  <b>Licence:</b> CC-BY-ND</p>	<p>Finnish</p>	<p>This corpus contains texts from 1544 to 1551 written by the clergyman Mikael Agricola.</p> <p>The corpus is available through the concordancer Korp.</p>
<p><a href="#">Virtual Old Literary Finnish (VVKS) - Kielipankki Korp version</a></p> <p><b>Size:</b> 48 texts  <b>Licence:</b> CC-BY-NC-ND</p>	<p>Finnish</p>	<p>This corpus will contain literary texts from 1543 to 1791.</p>

<a href="#">Finnish Folk Poetry</a> <b>Size:</b> 7.1 million words <b>Licence:</b> CC-BY-NC	Finnish, Karelian, Ludian, Latin, Swedish, Olonets, Izhorian, Votic	This corpus contains poems from 1564 to 1939. The corpus is available through the concordancer Korp.
<a href="#">Corpus of Early Modern Finnish, Kielipankki Version</a> <b>Size:</b> 8.6 million words <b>Licence:</b> EUPL v.1.1 SA	Finnish, Russian, German, Latin	This corpus contains texts from 1809 to 1899. The corpus is available through the concordancer Korp.
<a href="#">Aleksis Kivi Corpus (SKS)</a> <b>Size:</b> 413,700 words <b>Licence:</b> CC-BY-NC	Finnish, Swedish	This corpus contains the works by Finnish author Aleksis Kivi from 1855 to 1871. The corpus is available through the concordancer Korp.
<a href="#">Classics Library of the National Library of Finland - Kielipankki version</a> <b>Licence:</b> CC-BY	Finnish, Swedish	This corpus will contain literary texts from 1549 to 1944.
<a href="#">The Letters of Paul Sinebrychoff, Kielipankki Version</a> <b>Size:</b> 8.6 million words <b>Licence:</b> CC-BY	Finnish, Swedish	This corpus contains letters from 1895 to 1909. The corpus is available through the concordancer Korp.
<a href="#">The Newspaper and Periodical Corpus of the National Library of Finland, Kielipankki Version</a> <b>Size:</b> 8.7 billion words <b>Licence:</b> CC-BY	Finnish, Swedish	This corpus contains newspaper articles from 1770 to 2011. The corpus is available through the concordancer Korp.
<a href="#">The Newspaper and Periodical OCR Corpus of the National Library of Finland (1771-1874)</a> <b>Licence:</b> CC-BY	Finnish, Swedish	This corpus contains newspaper articles from 1771 to 1874. The corpus is available for download from the Language Bank of Finland.
<a href="#">The Newspaper and Periodical OCR Corpus of the National Library of Finland (1875-1920)</a> <b>Size:</b> 8.7 billion tokens <b>Annotation:</b> tokenised <b>Licence:</b> CLARIN ACA	Finnish, Swedish	This corpus contains newspaper articles from 1875 to 1920. The corpus is available for download from the Language Bank of Finland.
<a href="#">Partonopeus de Blois: transcriptions of all manuscripts and fragments</a> <b>Licence:</b> CC-BY-NC-SA 3.0	French (Old)	This corpus contains transcriptions of the manuscripts and fragments of the romance Partonopeus de Blois. The corpus is available for download from the Oxford Text Archive.

<a href="#">Syntactic Reference Corpus of Medieval French</a>  <b>Size:</b> 245,000 tokens <b>Annotation:</b> tokenised, syntactically-parsed <b>Licence:</b> CLARIN ACA	French (Old)	<p>This corpus contains texts from the 9<sup>th</sup> to the 13<sup>th</sup> century.</p> <p>The corpus is available for download from a dedicated webpage.</p>
<a href="#">Austrian Baroque Corpus</a>  <b>Size:</b> 200,000 tokens <b>Annotation:</b> tokenised, PoS-tagged, lemmatised, named entities	German	<p>This corpus contains sermons from 1650 to 1750.</p> <p>The corpus is available through a dedicated concordancer.</p>
<a href="#">DDR-Preseportal (GDR press portal)</a>  <b>Size:</b> 1.1 billion tokens <b>Annotation:</b> tokenised, lemmatised, PoS-tagged, normalised orthography <b>Licence:</b> CLARIN ACA	German	<p>This corpus contains newspaper texts from 1945 to 1994.</p> <p>The corpus is available through a concordancer provided by CLARIN-D.</p>
<a href="#">Deutsches Textarchiv (DTA)</a>  <b>Licence:</b> CLARIN PUB	German	<p>This corpus contains texts from the 17<sup>th</sup> to the 20<sup>th</sup> century.</p>
<a href="#">Die Grenzboten (journal)</a>  <b>Size:</b> 89 million tokens <b>Annotation:</b> tokenised, lemmatised, PoS-tagged, normalised orthography	German	<p>This corpus contains texts from 1842 to 1921.</p> <p>The corpus is available for download from the Deutsches Text Archiv and through a concordancer.</p>
<a href="#">Dinglers Polytechnisches Journal (Polytechnical Journal of Dingler)</a>  <b>Size:</b> 77.5 million tokens <b>Annotation:</b> tokenised, PoS-tagged, lemmatised, normalized orthography <b>Licence:</b> CC-BY-NC-SA 3.0	German	<p>This corpus contains academic texts from 1820 to 1931.</p> <p>The corpus is available for download from the Deutsches Text Archiv and through a concordancer.</p>
<a href="#">GeMi Corpus</a>  <b>Size:</b> 120,000 tokens <b>Annotation:</b> TEI Lite markup, no linguistic annotation <b>Licence:</b> CC-BY-NC-SA 3.0	German	<p>This corpus contains medical writing from 1500 to 1700.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<a href="#">GerManC. A Historical Corpus of German Newspapers 1650-1800</a>  <b>Size:</b> 700,000 words <b>Licence:</b> CC-BY-NC-SA 3.0	German	<p>This corpus contains personal letters, sermons and fictional, scholarly (i.e., humanities), scientific and legal texts from 1650 to 1800.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>



<a href="#">Mannheimer Korpus Historischer Zeitungen und Zeitschriften</a> <b>Size:</b> 3532 pages	German	<p>This corpus contains texts from the 18th and 19th centuries.</p> <p>The corpus is available for download directly through the VLO</p>
<a href="#">Referenzkorpus Mittelhochdeutsch (Middle High German Reference Corpus)</a> <b>Size:</b> 2.5 million tokens <b>Annotation:</b> tokenised, PoS-tagged, lemmatised, normalised, morphosyntactic description <b>Licence:</b> CC-BY-SA 4.0	German	<p>This corpus contains texts from 1050 to 1350.</p> <p>The corpus is available for download from the Deutsches Text Archiv and through a concordancer.</p>
<a href="#">B4 Historisches Predigtenkorpus zum Nachfeld</a> <b>Size:</b> 92,500 tokens <b>Annotation:</b> tokenised, syntactic and discursive annotation <b>Licence:</b> CLARIN ACA	German (Middle High)	<p>This corpus contains sermons from an Upper German (Balvarian-Alemannic) dialect area.</p> <p>The corpus is available for download from the repository of the University of Hamburg.</p>
<a href="#">B4 Ludolf</a> <b>Size:</b> 6,690 tokens <b>Annotation:</b> tokenised, tagged for clause type and grammatical function <b>Licence:</b> CLARIN ACA	German (Middle Low)	<p>This corpus contains texts from a journey diary from 1350.</p> <p>The corpus is available for download from the repository of the University of Hamburg.</p>
<a href="#">Reference Corpus Middle Low German/Low Rhenish (1200-1650)</a> <b>Size:</b> 200,700 tokens <b>Annotation:</b> tokenised, MSD-tagged <b>Licence:</b> CC-BY	German (Middle Low)	<p>This corpus contains texts from the 13<sup>th</sup> century to the middle of the 17<sup>th</sup> century.</p> <p>The corpus is available for download from the repository of the University of Hamburg.</p>
<a href="#">OROSSIMO Corpus – History</a> <b>Size:</b> 553,000 tokens <b>Annotation:</b> structural annotation (paragraph) <b>Licence:</b> CC-BY	Greek	<p>This corpus contains historic academic texts.</p> <p>The corpus is available for download from the clarin:el repository.</p>
<a href="#">Hungarian Historical Corpus</a> <b>Size:</b> 30 million words	Hungarian	<p>This corpus contains historical texts from the 18th century to the 2000s.</p> <p>The corpus is available through a dedicated concordancer.</p>

<p><a href="#">B4 Tatian Corpus of Deviating Examples 2.1</a></p> <p><b>Size:</b> 11,300 tokens  <b>Annotation:</b> tokenised, MSD-tagged  <b>Licence:</b> CC-BY</p>	<p>Latin, German (Old High)</p>	<p>This corpus contains the OHG Tatian, which is one of the largest prose texts from the Old High German period.</p> <p>The corpus is available for download and through a concordancer from the repository of the University of Hamburg.</p>
<p><a href="#">Menota</a></p> <p><b>Size:</b> 1.6 million tokens  <b>Annotation:</b> tokenised, MSD-tagged, lemmatised  <b>Licence:</b> CC-BY</p>	<p>Old Norse</p>	<p>This corpus contains Medieval Nordic texts.</p> <p>The corpus is available through the concordancer <i>Corpuscle</i>.</p>
<p><a href="#">Polish language of the 1960s</a></p> <p><b>Size:</b> 500,000 words  <b>Annotation:</b> MSD-tagged  <b>Licence:</b> CC-BY-NC-SA 3.0</p>	<p>Polish</p>	<p>This corpus contains essays, news articles, and scientific and literary texts from 1963 to 1967.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">Corpus of biblical text in Scots / John Kirk</a></p> <p><b>Annotation:</b> none  <b>Licence:</b> Oxford Text Archive licence</p>	<p>Scots</p>	<p>This corpus contains Biblical texts.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">The Helsinki corpus of Older Scots : [1450-1700]</a></p> <p><b>Annotation:</b> none  <b>Licence:</b> CC-BY-NC-SA 3.0</p>	<p>Scots</p>	<p>This corpus contains texts of different domains and genres (e.g., burgh records, diaries, pamphlets, scientific treatises, sermons) from 1450 to 1700.</p> <p>The corpus is available for download from the Oxford Text Archive.</p>
<p><a href="#">Digital library and corpus of historical Slovene IMP 1.1</a></p> <p><b>Size:</b> 17.7 million tokens  <b>Annotation:</b> tokenised, lemmatised, PoS-tagged  <b>Licence:</b> CC-BY-SA 4.0</p>	<p>Slovenian</p>	<p>This corpus contains 658 unique texts from 1584 to 1919.</p> <p>The corpus is available for download from the CLARIN.SI repository and through the concordancer KonText.</p>
<p><a href="#">Reference corpus of historical Slovene goo300k 1.2</a></p> <p><b>Size:</b> 300,000 tokens  <b>Annotation:</b> manually tokenised, lemmatised, PoS-tagged, modern synonyms for archaic words  <b>Licence:</b> CC-BY 4.0</p>	<p>Slovenian</p>	<p>This corpus contains 89 unique texts from 1584 to 1899.</p> <p>The corpus is available for download from the CLARIN.SI repository and through the concordancer KonText.</p>

<a href="#">The Swedish Sub-corpus of the Newspaper and Periodical Corpus of the National Library of Finland, Kielipankki Version</a>  <b>Size:</b> 3.5 billion tokens <b>Annotation:</b> tokenised <b>Licence:</b> CC-BY-SA.	Swedish	This corpus contains newspaper articles from 1770 to 1950.  The corpus is available through the concordancer Korp.
<a href="#">Språkbanken's historical corpora</a>  <b>Size:</b> 1.34 billion tokens <b>Annotation:</b> tokenised, PoS-tagged, lemmatised, syntactically parsed, word sense (for materials more recent than 1800)	Swedish, German, French and others	This collection of corpora contains – among others – diachronic legal texts, Bible translations, medieval letters, digitized newspapers from the Swedish National Library and 19 <sup>th</sup> century fiction from the Swedish Literature Bank.  The corpora are available through the concordancer Korp.
<a href="#">Historical Corpus of the Welsh Language 1500-1850</a>  <b>Size:</b> 420,000 words	Welsh	This corpus contains 30 texts from 1500 to 1850.  The corpus is available for download from a dedicated website and through a dedicated concordancer.

## 2.1. Summary

### 2.1.1. Identification

56 out of the 68 CLARIN corpora can be found in the VLO. The following 12 corpora are not yet included in the VLO:

1. [Menota](#)
2. [Greek Medieval Texts](#)
3. [Austrian Baroque Corpus](#)
4. [OROSSIMO Corpus - History](#)
5. [DDR-Presseportal \(GDR press portal\)](#)
6. [The Morpho-Syntactic Database of Mikael Agricola's Works](#)
7. [Classics Library of the National Library of Finland - Kielipankki version](#)
8. [Virtual Old Literary Finnish \(VVKS\) - Kielipankki Korp version](#)
9. [The Newspaper and Periodical OCR Corpus of the National Library of Finland \(1875-1920\)](#)
10. [The Finnish Sub-corpus of the Newspaper and Periodical Corpus of the National Library of Finland, Kielipankki Version](#)
11. [The Swedish Sub-corpus of the Newspaper and Periodical Corpus of the National Library of Finland, Kielipankki Version](#)
12. [Språkbanken's historical corpora](#)

### 2.1.2. Availability

The following 9 corpora are available for download and through a concordancer:

1. [Reference corpus of historical Slovene goo300k 1.2](#)

2. [Digital library and corpus of historical Slovene IMP 1.1](#)
3. [Historical Corpus of the Welsh Language 1500-1850](#)
4. [The Old Bailey Corpus](#)
5. [B4 Tatian Corpus of Deviating Examples 2.1](#)
6. [Dinglers Polytechnisches Journal \(Polytechnical Journal of Dingler\)](#)
7. [Referenzkorpus Mittelhochdeutsch \(Middle High German Reference Corpus\)](#)
8. [Die Grenzboten \(journal\)](#)
9. [Corpus Gysseling](#)

The following 22 corpora are available through a concordancer:

1. [Hungarian Historical Corpus](#)
2. [Corpus of Historical American English - Kielipankki Korp version 2017H1](#)
3. [Helsinki Corpus of Scottish Correspondence \(1540-1750\)](#)
4. [Menota](#)
5. [Austrian Baroque Corpus](#)
6. [Hansard Corpus](#)
7. [DDR-Presseportal \(GDR press portal\)](#)
8. [Brieven als buit \(Letters as loot\)](#)
9. [The Morpho-Syntactic Database of Mikael Agricola's Works](#)
10. [The Finnish Gutenberg Corpus](#)
11. [Aleksis Kivi Corpus \(SKS\)](#)
12. [Finnish Folk Poetry](#)
13. [Classics of Finnish Literature, Kielipankki Version](#)
14. [Corpus of Old Literary Finnish](#)
15. [Corpus of Early Modern Finnish, Kielipankki Version](#)
16. [The Letters of Paul Sinebrychoff, Kielipankki Version](#)
17. [The Newspaper and Periodical Corpus of the National Library of Finland, Kielipankki Version](#)
18. [Open Richly Annotated Cuneiform Corpus, Korp Version](#)
19. [The Finnish Sub-corpus of the Newspaper and Periodical Corpus of the National Library of Finland, Kielipankki Version](#)
20. [The Swedish Sub-corpus of the Newspaper and Periodical Corpus of the National Library of Finland, Kielipankki Version](#)
21. [Språkbanken's historical corpora](#)
22. [Corpus of Old Written Estonian](#)

The following 35 corpora are available for download:

1. [Medieval Charter Sections Corpus](#)
2. [Sheffield Corpus of Chinese](#)
3. [GerManC. A Historical Corpus of German Newspapers 1650-1800](#)
4. ["PoDiLemma" Middle Polish Diachrone Lemmatised Corpus](#)
5. ["PoDiLemma" Middle Polish Diachrone Lemmatised Corpus](#)
6. [Parsed Corpus of Early English Correspondence \(PCEEC\)](#)
7. [Syntactic Reference Corpus of Medieval French](#)
8. [Hamburg Corpus of Old Swedish with Syntactic Annotations \(HaCOSSA\)](#)
9. [Reference Corpus Middle Low German/Low Rhenish \(1200-1650\)](#)

10. [B4 Ludolf](#)
11. [B4 Historisches Predigtenkorpus zum Nachfeld](#)
12. [Mannheimer Korpus Historischer Zeitungen und Zeitschriften](#)
13. [Greek Medieval Texts](#)
14. [OROSSIMO Corpus – History](#)
15. [Older Scottish texts : the Edinburgh DOST corpus / A.J. Aitken, Paul Bratley and Neil Hamilton-Smith](#)
16. [Anthology of Middle English texts / Santiago Gonzalez y Fernandez-Corugedo](#)
17. [Helsinki corpus of English texts](#)
18. [Corpus of biblical text in Scots / John Kirk](#)
19. [Pamphlets of the American Revolution : \[selections\] / edited by Bernard Bailyn](#)
20. [Corpus of Late Modern English prose / David Denison](#)
21. [The Helsinki corpus of Older Scots : \[1450-1700\]](#)
22. [The Lampeter Corpus of Early Modern English Tracts](#)
23. [The York-Helsinki parsed corpus of Old English poetry \(YCOEP\)](#)
24. [Corpus of Early English Correspondence Sampler \(CEECS\)](#)
25. [The York-Toronto-Helsinki Parsed Corpus of Old English prose \(YCOE\)](#)
26. [The English language of the north-west in the late Modern English period: a Corpus of late 18c Prose](#)
27. [Polish language of the 1960s](#)
28. [Dictionary of Old English Corpus in Electronic Form \(DOEC\)](#)
29. [Partonopeus de Blois: transcriptions of all manuscripts and fragments](#)
30. [A Corpus of English Dialogues 1560-1760 \(CED\)](#)
31. [The Electronic Text Corpus of Sumerian Literature. Revised edition.](#)
32. [The Lancaster Newsbooks Corpus](#)
33. [GeMi Corpus](#)
34. [The Newspaper and Periodical OCR Corpus of the National Library of Finland \(1771-1874\)](#)
35. [The Newspaper and Periodical OCR Corpus of the National Library of Finland \(1875-1920\)](#)

### **2.1.3. Metadata**

#### **2.1.3.1. Language**

There are 51 monolingual historical corpora in the CLARIN infrastructure for the following 16 languages:

1. English (15 corpora)
2. German (12 corpora)
3. Finnish (6 corpora)
4. French (2 corpora)
5. Dutch (2 corpora)
6. Greek (2 corpora)
7. Slovenian (2 corpora)
8. Scots (2 corpora)

9. Akkadian (1 corpus)
10. Chinese (1 corpus)
11. Estonian (1 corpus)
12. Hungarian (1 corpus)
13. Old Norse (1 corpus)
14. Polish (1 corpus)
15. Swedish (1 corpus)
16. Welsh (1 corpus)

The other 17 corpora are multilingual, for the following language combinations:

1. Finnish, Swedish (6 corpora)
2. English, Latin (2 corpora)
3. Czech, German, Latin, Polish (1 corpus)
4. Czech, Latin (1 corpus)
5. English, Hebrew (1 corpus)
6. English, German, Latin, Old Norse, Swedish (1 corpus)
7. English, Sumerian (1 corpus)
8. Finnish, Karelian, Ludian, Latin, Swedish, Olonets, Izhorian, Votic (1 corpus)
9. Finnish, Russian, German, Latin (1 corpus)
10. Latin, German (1 corpus)
11. Swedish, German, French, etc. (1 corpus)

#### 2.1.3.2. Size and temporal span

Information on word/token size is available for 53 corpora. The following 15 corpora lack such information:

1. [Sheffield Corpus of Chinese](#)
2. ["PoDiLemma" Middle Polish Diachrone Lemmatised Corpus](#)
3. [Deutsches Textarchiv \(DTA\)](#)
4. [Corpus of biblical text in Scots / John Kirk](#)
5. [Pamphlets of the American Revolution : \[selections\] / edited by Bernard Bailyn](#)
6. [Corpus of Late Modern English prose / David Denison](#)
7. [The Helsinki corpus of Older Scots : \[1450-1700\]](#)
8. [The Lampeter Corpus of Early Modern English Tracts](#)
9. [Dictionary of Old English Corpus in Electronic Form \(DOEC\)](#)
10. [Partonopeus de Blois: transcriptions of all manuscripts and fragments](#)
11. [The Electronic Text Corpus of Sumerian Literature. Revised edition.](#)
12. [The Lancaster Newsbooks Corpus](#)
13. [Classics Library of the National Library of Finland - Kielipankki version](#)
14. [Virtual Old Literary Finnish \(VVKS\) - Kielipankki Korp version](#)
15. [The Newspaper and Periodical OCR Corpus of the National Library of Finland \(1771-1874\)](#)

The largest corpus is [The Newspaper and Periodical Corpus of the National Library of Finland, Kielipankki Version](#), containing 8.7 billion tokens.

In terms of temporal span, there is no common denominator, but generally the oldest corpora reach back to the 5<sup>th</sup> century (e.g., there are 3 Old English corpora). The exception is [Open Richly Annotated Cuneiform Corpus, Korp Version](#), which contains Sumerian texts and is thus the oldest corpus by far.

#### 2.1.3.3. Annotation

Information on linguistic annotation is available for 47 corpora. The following 21 corpora lack such information:

1. [Hungarian Historical Corpus](#)
2. [Sheffield Corpus of Chinese](#)
3. [Historical Corpus of the Welsh Language 1500-1850](#)
4. [GerManC. A Historical Corpus of German Newspapers 1650-1800](#)
5. [Deutsches Textarchiv \(DTA\)](#)
6. [Mannheimer Korpus Historischer Zeitungen und Zeitschriften](#)
7. [Greek Medieval Texts](#)
8. [Dictionary of Old English Corpus in Electronic Form \(DOEC\)](#)
9. [Partonopeus de Blois: transcriptions of all manuscripts and fragments](#)
10. [A Corpus of English Dialogues 1560-1760 \(CED\)](#)
11. [The Lancaster Newsbooks Corpus](#)
12. [Aleksis Kivi Corpus \(SKS\)](#)
13. [Finnish Folk Poetry](#)
14. [Classics of Finnish Literature, Kielipankki Version](#)
15. [Corpus of Old Literary Finnish](#)
16. [Corpus of Early Modern Finnish, Kielipankki Version](#)
17. [The Letters of Paul Sinebrychoff, Kielipankki Version](#)
18. [The Newspaper and Periodical Corpus of the National Library of Finland, Kielipankki Version](#)
19. [Classics Library of the National Library of Finland - Kielipankki version](#)
20. [Virtual Old Literary Finnish \(VVKS\) - Kielipankki Korp version](#)
21. [The Newspaper and Periodical OCR Corpus of the National Library of Finland \(1771-1874\)](#)

18 corpora are PoS/MSD-tagged, 13 corpora are lemmatised and 6 corpora are syntactically parsed.

#### 2.1.3.4. Licence

Licence information is available for 60 corpora. The following 8 corpora lack such information:

1. [Hungarian Historical Corpus](#)
2. [Historical Corpus of the Welsh Language 1500-1850](#)
3. ["PoIDiLemma" Middle Polish Diachrone Lemmatised Corpus](#)
4. [Mannheimer Korpus Historischer Zeitungen und Zeitschriften](#)
5. [Austrian Baroque Corpus](#)
6. [Hansard Corpus](#)
7. [Die Grenzboten \(journal\)](#)
8. [Brieven als buit \(Letters as loot\)](#)

35 corpora are available under CC-BY, 11 corpora under the Oxford Text Archive licence, 6 under CLARIN ACA, 3 under EUPL v.1.1 SA, the rest have miscellaneous licences.

### 3. Non-CLARIN historical corpora

Table 2 lists 25 corpora that are not part of the CLARIN infrastructure.

Table 2: non-CLARIN historical corpora

Corpus	Language	Description
<a href="#">DIAKORP v6</a>  <b>Size:</b> 4 million tokens <b>Annotation:</b> basic structural markup <b>Licence:</b> CC-BY-NC-SA	Czech	This corpus contains texts from the 14 <sup>th</sup> to the 20 <sup>th</sup> century.  The corpus is available through a dedicated concordancer.
<a href="#">Brieven als buit (Letters as loot)</a>  <b>Size:</b> 460,000 words <b>Annotation:</b> lemmatised, PoS-tagged, grammatically tagged	Dutch	This corpus contains 1,000 letters from the 17 <sup>th</sup> to the 18 <sup>th</sup> century.  The corpus is available through a dedicated concordancer.
<a href="#">ARCHER Corpus</a>	English	The corpus contains texts from 1600 to 1999.  Availability of the corpus is restricted.
<a href="#">ECCO-TCP</a>  <b>Size:</b> 74 million tokens <b>Annotation:</b> no linguistic annotation <b>Licence:</b> CC-1	English	This corpus contains texts (literature, philosophy, politics, religion, geography, science and all other areas of human endeavour) from 1700 to 1800.  The corpus is available for download from the University of Oxford Text Archive and through a dedicated concordancer.
<a href="#">EEBO-TCP</a>  <b>Size:</b> 766 million tokens <b>Annotation:</b> no linguistic annotation <b>Licence:</b> CC-0	English	This corpus contains texts (literature, philosophy, politics, religion, geography, science and all other areas of human endeavour) from 1450 to 1750.  The corpus is available for download from the University of Oxford Text Archive and through a dedicated concordancer.
<a href="#">EVANS-TCP</a>  <b>Size:</b> 766 million tokens <b>Annotation:</b> no linguistic annotation <b>Licence:</b> CC-0	English	This corpus contains American texts from 1640 to 1821.  The corpus is available for download from the University of Oxford Text Archive and through a dedicated concordancer.
<a href="#">Historical Corpora at Lancaster University</a>  <b>Annotation:</b> tokenised, PoS-tagged, partial semantic tagging (USAS system)	English	The corpus contains texts in various domains (e.g., fiction, newspaper texts, religious texts) from 1500 on.  The corpus is available through the CQPConcordancer.



<p><a href="#">Frantext</a></p> <p><b>Size:</b> 300 million words  <b>Annotation:</b> PoS-tagged, lemmatised</p>	French	<p>This corpus contains texts from the 10<sup>th</sup> to the 21<sup>st</sup> century.</p> <p>The corpus is available through a dedicated concordancer.</p>
<p><a href="#">Bundesblatt/Feuille fédérale/Foglio federale</a></p> <p><b>Size:</b> 203,585,806 tokens (German), 239,125,036 tokens (French), 85,223,085 tokens (Italian)  <b>Annotation:</b> tokenised, syntactically-parsed</p>	German, French, Italian	<p>This corpus contains texts from 1849 to 2014.</p> <p>The corpus is available through the CQPWeb concordancer.</p>
<p><a href="#">Corpus of Old and Middle Hungarian court records and private correspondence</a></p> <p><b>Size:</b> 850,000 words  <b>Annotation:</b> tokenised, MSD-tagged, lemmatised, sociolinguistic metadata</p>	Hungarian	<p>This corpus contains private letters and testimonies from the 16<sup>th</sup> to the 18<sup>th</sup> century.</p> <p>The corpus is available through a dedicated concordancer.</p>
<p><a href="#">Old Hungarian Corpus</a></p> <p><b>Size:</b> 3 million tokens  <b>Annotation:</b> tokenised, partially normalized, partially MSD-tagged</p>	Hungarian	<p>This corpus contains texts (codices, letters) from the 12<sup>th</sup> to the 17<sup>th</sup> century.</p> <p>The corpus is available for download from a dedicated webpage and through a dedicated concordancer.</p>
<p><a href="#">Archivio Datini</a></p> <p><b>Annotation:</b> lemmatised</p>	Italian	<p>This corpus contains Italian texts.</p> <p>The corpus is available through a dedicated concordancer.</p>
<p><a href="#">Corpus testuale del Tesoro della Lingua Italiana delle Origini</a></p> <p><b>Size:</b> 23 million tokens  <b>Annotation:</b> tokenised, lemmatised</p>	Italian	<p>This corpus contains early Italian texts before 1375.</p> <p>The corpus is available through a dedicated concordancer.</p>

<a href="#">DiaCORIS</a>	Italian	This corpus contains texts from 1861 to 1945.  The corpus is available through a dedicated concordancer.
<a href="#">M.I.DIA. (Morfologia dell'Italiano in DIACronia)</a>  <b>Size:</b> 7.5 million tokens <b>Annotation:</b> tokenised <b>Licence:</b> CC-BY-NC 4.0	Italian	This corpus contains texts from the 13 <sup>th</sup> to the 20 <sup>th</sup> century.  The corpus is available through a dedicated concordancer
<a href="#">Chronopress</a>  <b>Size:</b> 16 million tokens <b>Licence:</b> CC-BY-SA	Polish	This corpus contains newspaper articles from 1945 to 1954.  The corpus is available through a dedicated concordancer.
<a href="#">Corpus of the 19. century Polish (Korpus polszczyzny XIX-wiecznej)</a>  <b>Size:</b> 625,000 tokens <b>Annotation:</b> tokenised, PoS-tagged, lemmatised, transliteration, transcripton	Polish	This corpus contains texts from 1830 to 1918.  The corpus is available through a dedicated concordancer.
<a href="#">IMPACT GT corpus (Korpus GT projektu IMPACT)</a>  <b>Size:</b> 1.5 million tokens <b>Annotation:</b> transcription	Polish	This corpus contains texts from 1570 to 1756.  The corpus is available through a dedicated concordancer.
<a href="#">Corpus of old Polish texts until 1500 (Korpus tekstów staropolskich do roku 1500)</a>  <b>Size:</b> 620,000 tokens <b>Annotation:</b> tokenised	Polish, Latin	This corpus contains texts until 1500.  The corpus is available for download from a dedicated webpage.
<a href="#">Corpus of the 16. century Polish (Korpus polszczyzny XVI wieku)</a>  <b>Annotation:</b> lemmatised, transliteration	Polish, Latin	This corpus contains texts from the 16 <sup>th</sup> century.  The corpus is available through a dedicated concordancer.
<a href="#">eFontes Mediae et Infimae Latinitatis Polonorum (Elektroniczny korpus polskiej łaciny średniowiecznej)</a>  <b>Size:</b> 5 million tokens <b>Annotation:</b> tokenised, lemmatised	Polish, Latin	This corpus contains texts from the 11 <sup>th</sup> to the middle of the 16 <sup>th</sup> century.  The corpus is available through a dedicated concordancer.

<a href="#">The Electronic Corpus of the 17th and 18th century Polish (Korpus tekstów polskich z XVII i XVIII w.)</a>  <b>Size:</b> 12 million tokens <b>Annotation:</b> tokenised, partially PoS-tagged, structural annotation	Polish, Latin	<p>This corpus contains texts from 1601 to 1772.</p> <p>The corpus is available through a dedicated concordancer.</p>
<a href="#">XV century New Testament translations (Piętnastowieczne przekłady Nowego Testamentu – elektroniczna konkordancja staropolska)</a>  <b>Size:</b> 400,000 tokens <b>Annotation:</b> tokenised	Polish, Latin	<p>This corpus contains Biblical texts from 1380 to 1500.</p> <p>This corpus is available for download from a dedicated webpage and through a dedicated concordancer.</p>
<a href="#">Corpus Informatizado do Português Medieval</a>  <b>Size:</b> 2 million tokens <b>Annotation:</b> tokenised, PoS-tagged	Portuguese	<p>This corpus contains texts from the 9<sup>th</sup> to the 16<sup>th</sup> century.</p> <p>The corpus is available through a dedicated concordancer.</p>
<a href="#">Parsed Corpus of Historical Portuguese</a>  <b>Size:</b> 3.3 million <b>Annotation:</b> tokenised, PoS-tagged (2 million), treebanked (1.2 million)	Portuguese	<p>This corpus contains 76 texts written by authors born between 1380 and 1881.</p> <p>The corpus is available through a dedicated concordancer.</p>

### 3.1. Summary

#### 3.1.1. Availability

The following 5 corpora are available for download and through a concordancer:

1. [XV century New Testament translations \(Piętnastowieczne przekłady Nowego Testamentu – elektroniczna konkordancja staropolska\)](#)
2. [Old Hungarian Corpus](#)
3. [ECCO-TCP](#)
4. [EEBO-TCP](#)
5. [EVANS-TCP](#)

The following 18 corpora are available through a concordancer.

1. [Corpus Informatizado do Português Medieval](#)
2. [Parsed Corpus of Historical Portuguese](#)
3. [Historical Corpora at Lancaster University](#)
4. [Corpus testuale del Tesoro della Lingua Italiana delle Origini](#)
5. [DiaCORIS](#)
6. [M.I.DIA. \(Morfologia dell'Italiano in DIACronia\)](#)
7. [Archivio Datini](#)
8. [Frantext](#)
9. [eFontes Mediae et Infimae Latinitatis Polonorum \(Elektroniczny korpus polskiej łaciny średniowiecznej\)](#)
10. [Corpus of the 16. century Polish \(Korpus polszczyzny XVI wieku\)](#)
11. [The Electronic Corpus of the 17th and 18th century Polish \(Korpus tekstów polskich z XVII i XVIII w.\)](#)
12. [Corpus of the 19. century Polish \(Korpus polszczyzny XIX-wiecznej\)](#)
13. [IMPACT GT corpus \(Korpus GT projektu IMPACT\)](#)
14. [Chronopress](#)
15. [Bundesblatt/Feuille fédérale/Foglio federale](#)
16. [DIAKORP v6](#)
17. [Corpus of Old and Middle Hungarian court records and private correspondence](#)
18. [Brieven als buit \(Letters as loot\)](#)

One corpus is available for download:

1. [Corpus of old Polish texts until 1500 \(Korpus tekstów staropolskich do roku 1500\)](#)

### 3.1.2. Metadata

#### 3.1.2.1. Language

There are 19 monolingual corpora in Table 2 for the 8 following languages:

1. English (5 corpora)
2. Italian (4 corpora)
3. Polish (3 corpora)
4. Hungarian (2 corpora)
5. Portuguese (2 corpora)
6. Czech (1 corpus)
7. Dutch (1 corpus)
8. French (1 corpus)

There are 6 multilingual corpora for the following 2 language combinations:

1. Polish, Latin (5 corpora)
2. German, French, Italian (1 corpus)

#### 3.1.2.2. Size and temporal span

Information on size is available for the majority of the corpora – 17 out of 22. The following corpora 5 corpora lack this information:

1. [ARCHER Corpus](#)
2. [Historical Corpora at Lancaster University](#)
3. [DiaCORIS](#)
4. [Archivio Datini](#)
5. [Corpus of the 16. century Polish \(Korpus polszczyzny XVI wieku\)](#)

The oldest corpus reaches back to the 9<sup>th</sup> century.

#### 3.1.2.3. Annotation

Information on annotation is likewise available for the majority of the corpora – 19 out of 22. The following 3 corpora lack this information:

1. [ARCHER Corpus](#)
2. [DiaCORIS](#)
3. [Corpus of the 19. century Polish \(Korpus polszczyzny XIX-wiecznej\)](#)

#### 3.1.2.4. Licence

Licence information is available for 4 corpora:

1. [Chronopress](#)
2. [ECCO-TCP](#)
3. [EEBO-TCP](#)
4. [EVANS-TCP](#)

## 4. Linguistic tools for the annotation of historical corpora

The following tools/services were contributed by a NC/UI coordinator in the Google Spreadsheet. All are available through a CLARIN-D node:

1. [Part-of-speech tagging: mixed text](#)  
**Description:** A webapplication that reads in mixed texts written in Latin and Middle English in txt-Format
2. [TreeTagger -- Middle High German parameter file](#)  
**Description:** A trained parameter file for the TreeTagger (part-of-speech tagger) to tag Middle High German text.
3. [OCR Post-Correction](#)  
**Description:** A web application that reads in scans of German books from the 18th to 20th century in Antiqua and Fraktur and runs a pipeline of OCR with tesseract and a post-correction system using different techniques.