

# Historisches Korpus zum Schlagwort "Utraquismus"

- Über die Datensammlung und Dokumentation
- Umfang und Aufbau des Korpus
  - Konkretere Informationen zur Korpusaufbau und Korpusgröße in Bezug auf die Originalquellen
  - v 0.1
  - v 0.2
- Metadaten
- Orthographische und morphologische Normalisierung
  - Abweichtungstypen und Abweichungen

# Über die Datensammlung und Dokumentation

## Inhalte der Datensammlung

Das Korpus enthält Textausschnitte rund um das **Schlagwort**

**"utraquistisch"/"Utraquismus"/"Utraquist" (Abfrage: "utra??is\*" aus historischen Zeitungen und Zeitschriften, die im [ANNO \(AustriaN Newspapers Online\)](#) der Österreichischen Nationalbibliothek verfügbar und durchsuchbar sind. Die dort verfügbaren OCR-Scans wurden manuell überprüft und korrigiert sowie in Bezug auf Orthographie und Morphologie normalisiert, um bestmögliche automatische Verarbeitung mit auf die Gegenwartssprache des Deutschen ausgelegten Taggern und Lemmatizern zu ermöglichen.**

Es wurden keine kompletten Zeitungsausgaben oder Zeitungsartikel sondern nur Textausschnitte bearbeitet, die wie folgt definiert wurden:

**ein maximaler Textausschnitt** = der Satz mit dem Schlagwort plus bis zu fünf Sätze vor und/oder nach diesem innerhalb eines Absatzes

Ein Textausschnitt (kurz: "text") enthält demnach maximal 11 Sätze innerhalb eines Absatzes.

Jeder Textausschnitt (XML-Element `text`) ist in Bezug auf die folgenden **Metadaten** (als Attribute) annotiert:

- Zeitung/Zeitschrift (`source`)
- eventuell: Beilage der Zeitung/Zeitschrift (`source2`)
- Veröffentlichungsort (`place`)
- Veröffentlichungsdatum (`date`)
- Periode (für einfachere Durchsuchbarkeit, als Zuordnung zu einem Jahrzehnt, `period`)
- Genrezuordnung (aktuell noch provisorisch, `genre`)

- Thema des Artikels (aktuell noch provisorisch, `topic`)
- Bezugsraum des Artikels (`space`)

Jede **orthographische/morphologische Normalisierung** (XML-Element `edit`) ist durch folgende Attribute näher spezifiziert:

- originale Variante (`original`)
- Klassifikation des Abweichungstyps (aktuell noch provisorisch, `type`)

## Arbeit mit der Datensammlung

Das Korpus wird einerseits als **.xml-Dokument** zur Verfügung gestellt. Andererseits ist es auf **SketchEngine** verfügbar und kann dort auf Anfrage freigeschalten werden.

**ACHTUNG:** Der Lemmatizer erkennt die verschiedenen Formen der Schlagwörter nicht als jeweils ein Lemma! Um alle Formen (aller Schlagwörter zu finden) muss mit der **CQL-Abfrage** `[lemma="utraquis.*"]|[lemma="Utraquis.*"]` gearbeitet werden.

## Entstehung

Die Datensammlung ist das **Ergebnis einer Hausaufgabe** in einem von Agnes Kim am Institut für Germanistik der Universität Wien im Sommersemester 2021 angebotenen [Proseminar zur "Historischen Soziolinguistik"](#). Ziel der Aufgabe war, zwischen acht und zehn Texte aus einem vorgegebenen Paket rund um Suchergebnisse für die **Abfrage "utra??is\*" aus dem ANNO**, also aus historischen Zeitschriften aus dem Bestand der Österreichischen Nationalbibliothek in ein XML-Dokument zu übertragen und so zu normalisieren, dass die Texte auch mit den auf die Gegenwartssprache trainierten, auf SketchEngine zur Verfügung stehenden Taggern für das Deutsche bearbeitet werden können. Die Daten wurden von den Studierenden eingegeben und von Agnes Kim überprüft und angeglichen.

## Verfügbarkeit

Da es sich aktuell um eine erste Version dieser Datensammlung handelt, ist sie noch ausschließlich als **Download aus der [dioecloud](#)** verfügbar. Außerdem ist das Korpus auf SketchEngine verfügbar und kann auf Anfrage (per Mail an [agnes.kim@univie.ac.at](mailto:agnes.kim@univie.ac.at)) freigeschalten werden!

Version/Dateiname	Veröffentlichungsdatum	ungefähre Größe	Link
<b>v 0.1</b> / Korpus_Utraquis_v0.1.CSV	12.05.2021	156 Textausschnitte 31.176 Tokens	<a href="https://dioecloud.trans">https://dioecloud.trans</a>
<b>v 0.2</b> / Korpus_Utraquis_v0.2.CSV	06.06.2021	156 Textausschnitte 31.176 Tokens	<a href="https://dioecloud.trans">https://dioecloud.trans</a>

**Änderungsnotizen** können den [Seiten zum Korpusaufbau](#) entnommen werden.

Nach weiteren Überprüfungen und eventuellen Ergänzungen ist geplant, das Korpus in einer ersten stabilen Version zu archivieren und auf SketchEngine öffentlich zur Verfügung zu stellen.

## Zitation

### Zitation der Datensammlung

Kim, Agnes (JAHRESZAHL DES VERÖFFENTLICHUNGSDATUMS DER VERSION): Historisches Korpus zum Schlagwort "Utraquismus". Datensammlung. VERSIONSNUMMER ANGEBEN In: Dies. (Hg.): Datensammlungen aus der bzw. für die Forschung und Lehre. Online verfügbar unter: LINK ZUR VERSION.

### Zitation der Dokumentation

Kim, Agnes (2021): Historisches Korpus zum Schlagwort "Utraquismus". Dokumentation. In: Dies. (Hg.): Datensammlungen aus der bzw. für die Forschung und Lehre. Online verfügbar unter: <https://wiki.dioe.at/books/historisches-korpus-zum-schlagwort-%22utraquismus%22>.

## Fragen? Fehler gefunden?

Sie haben einen Fehler gefunden, haben Ergänzungen oder Nachfragen? Bitte wenden Sie sich per Mail an mich: [agnes.kim@univie.ac.at](mailto:agnes.kim@univie.ac.at)

# Umfang und Aufbau des Korpus

# Konkretere Informationen zur Korpusaufbau und Korpusgröße in Bezug auf die Originalquellen

In der folgenden Tabelle werden die bearbeiteten Textausschnitte in Relation zu einem Schlagwortkorpus, das das gesamte 19. und das frühe 20. Jahrhundert (bis 1918) grob abdeckt und nachvollziehbar macht, gesetzt. Insgesamt können (am 17. 05. 2021) mit der Suchabfrage **"utra??ist\*" im ANNO bis inklusive 1918 6.338 Treffer** (= Ausgaben von Zeitungen und Zeitschriften, in denen die Suchabfrage mindestens ein Mal vorkommt) gefunden werden, von denen **6.179** auf Zeitungen und 159 auf Zeitschriften entfallen. Letztere wurden bei den Überlegungen zu einem möglichen Gesamtkorpus außer Acht gelassen.

## Periode 1800–1829

Publikationsort	Link zu ANNO	Treffer	davon im Korpus
Brno	<a href="#">Link</a>	15	13 (2 Duplikate)
Wien	<a href="#">Link</a>	11	11
L'viv	<a href="#">Link</a>	1	0 (weil Duplikat)

## Periode 1850–1859

Publikationsort	Link zu ANNO	Treffer	davon im Korpus
-----------------	--------------	---------	-----------------

Wien	<a href="#">Link</a>	<b>29</b>	<b>8</b>
Olomouc	<a href="#">Link</a>	3	<b>3</b>
Salzburg	<a href="#">Link</a>	3	<b>3</b>
Brno	<a href="#">Link</a>	1	<b>1</b>
Opava	<a href="#">Link</a>	1	<b>1</b>
Litoměřice	<a href="#">Link</a>	1	<b>0</b>

## Periode 1867–1868 (stellvertretend für die 1860er)

<b>Jahr</b>	<b>Publikationsort</b>	<b>Link zu ANNO</b>	<b>Treffer</b>	<b>davon im Korpus</b>
1867	Wien	<a href="#">Link</a>	<b>13</b>	<b>11</b>
1867	Plzeň	<a href="#">Link</a>	4	<b>0</b>
1867	Graz	<a href="#">Link</a>	1	<b>0</b>
1867	Klagenfurt	<a href="#">Link</a>	1	<b>0</b>
1867	Znojmo	<a href="#">Link</a>	1	<b>0</b>
1867	Olomouc	<a href="#">Link</a>	4	<b>0</b>
1867	Praha	<a href="#">Link</a>	2	<b>0</b>
1868	Wien	<a href="#">Link</a>	<b>22</b>	<b>8</b>
1868	Znojmo	<a href="#">Link</a>	7	<b>0</b>
1868	Praha	<a href="#">Link</a>	6	<b>0</b>
1868	Graz	<a href="#">Link</a>	1	<b>0</b>
1868	Plzeň	<a href="#">Link</a>	1	<b>0</b>



## Jahr 1890

Publikationsort	Zeitung	Link zu ANNO	Treffer	davon im Korpus
Wien	Neue Freie Presse	<a href="#">Link</a>	<b>11</b>	<b>8</b>
Wien	Die Presse	<a href="#">Link</a>	8	<b>8</b>
Wien	Das Vaterland	<a href="#">Link</a>	<b>12</b>	<b>10</b>
Wien	Wiener Montags-Journal	<a href="#">Link</a>	1	<b>0</b>
Wien	Wiener Zeitung	<a href="#">Link</a>	6	<b>6</b>
Wien	Deutsches Volksblatt	<a href="#">Link</a>	4	<b>4</b>
Praha	alle	<a href="#">Link</a>	<b>13</b>	<b>0</b>
Maribor	alle	<a href="#">Link</a>	9	<b>6</b>
Salzburg	alle	<a href="#">Link</a>	4	<b>4</b>
Litoměřice	alle	<a href="#">Link</a>	8	<b>0</b>
Linz	alle	<a href="#">Link</a>	7	<b>0</b>
Znojmo	alle	<a href="#">Link</a>	3	<b>0</b>

## Jahr 1910

Publikationsort	Zeitung	Link zu ANNO	Treffer	davon im Korpus
Wien	Neue Freie Presse	<a href="#">Link</a>	<b>13</b>	<b>10</b>
Wien	Wiener Zeitung	<a href="#">Link</a>	<b>13</b>	<b>10</b>
Wien	Neues Wiener Tagblatt	<a href="#">Link</a>	6	<b>6</b>

Wien	Zollämter- und Finanzwachzeitung	<a href="#">Link</a>	4	<b>4</b>
Wien	Deutsches Volksblatt	<a href="#">Link</a>	5	<b>5</b>
Wien	Arbeiterzeitung	<a href="#">Link</a>	3	<b>3</b>
Praha	alle	<a href="#">Link</a>	<b>30</b>	<b>0</b>
Maribor	alle	<a href="#">Link</a>	<b>15</b>	<b>6</b>
Klagenfurt	alle	<a href="#">Link</a>	<b>10</b>	<b>9</b>
Graz	alle	<a href="#">Link</a>	<b>24</b>	<b>10</b>
Czernowitz	alle	<a href="#">Link</a>	<b>13</b>	<b>0</b>

## Mitarbeit?

Sie haben Interesse, einzelne Ausgaben zu einem der Datenpakete zu ergänzen? Gerne sende ich Ihnen eine genaue Liste der noch fehlenden Ausgaben zu! Melden Sie sich bei:  
 agnes.kim@univie.ac.at

# v 0.1

## Veröffentlichung

**Datum:** 12.05.2021

**Link:** <https://dioecloud.trans.univie.ac.at/index.php/s/36iS7xZ4f8dD6fK>

## Größe des Korpus

**Textausschnitte:** 156

**Tokens:** 31,176

**Wörter:** 26,806

## Überblick über den Korpusaufbau

Die Grafiken beziehen sich auf die Gesamtzahl der Token. Sie wurden mit SketchEngine erstellt.

## Zeitliche Dimension: Korpusgröße nach Jahrzehnten

Jahrzehnt	Tokenzahl
1800	382
1820	1.115
1850	2.343
1860	3.419
1890	10.024

1900	2.650
1910	11.241

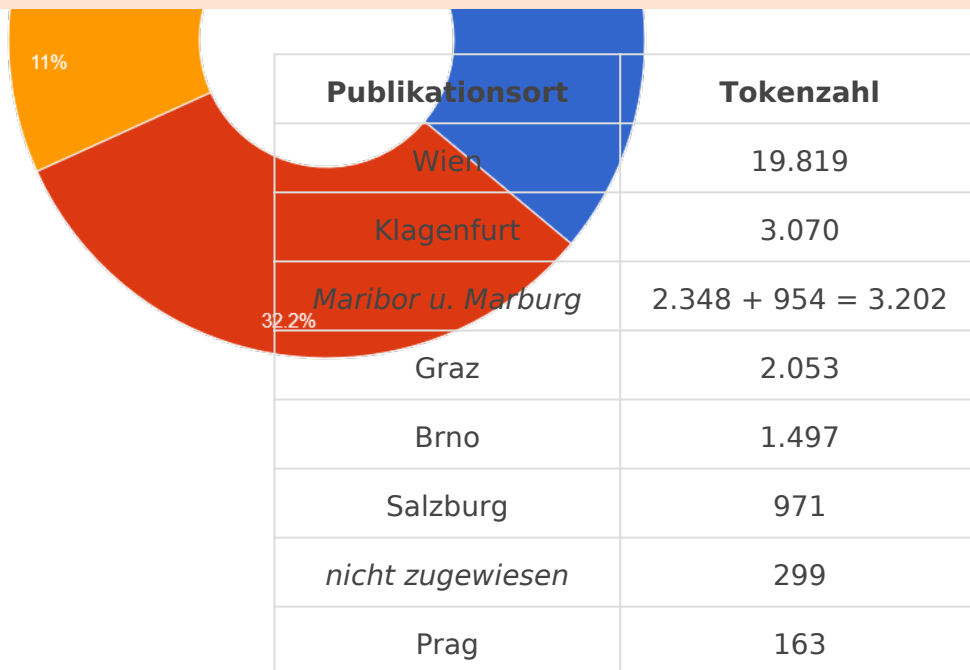
text - period

## Räumliche Dimension: Korpusgröße nach Publikationsorten

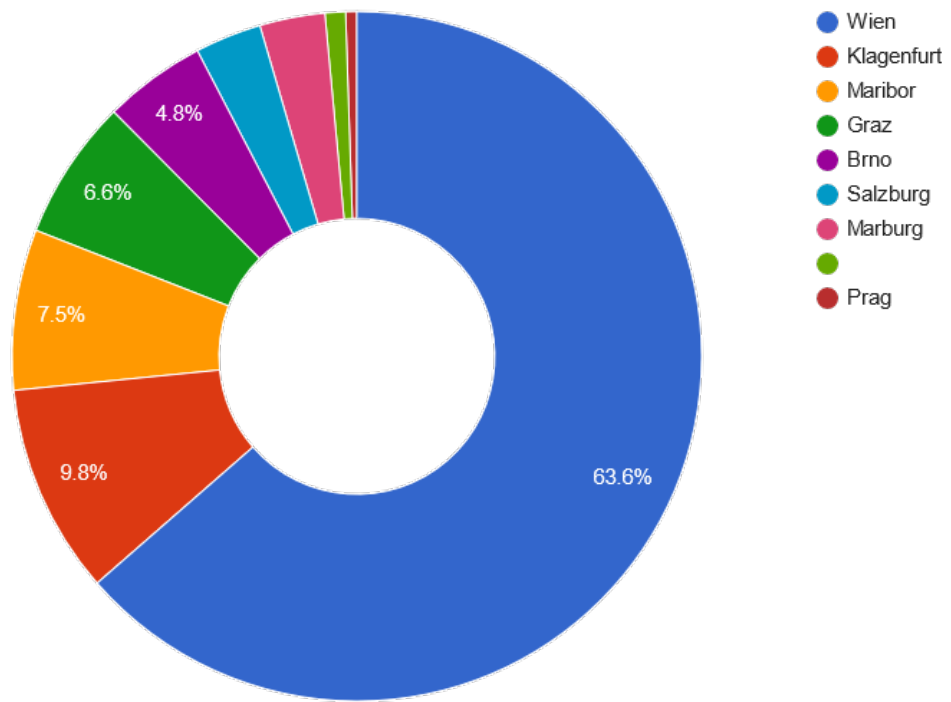


**Achtung:** In v 0.1 bestehen noch Probleme in Bezug auf die Zuweisung bzw.

Vereinheitlichung der Publikationsorte, die in der Tabelle durch Kursivsatz hervorgehoben werden!



text - place



# v 0.2

## Veröffentlichung

**Datum:** 06.06.2021

**Link:** <https://dioecloud.trans.univie.ac.at/index.php/s/P7p2jMcoFMgWFbx>

## Änderungsnotizen

- Vereinheitlichung der Periodenzuordnung ( `period` )
- Vereinheitlichung der Publikationsorte ( `place` )
- Vereinheitlichung der Annotation der orthographischen und morphologischen Normalisierungen ( `edit type` )

## Größe des Korpus

**Textausschnitte:** 156

**Tokens:** 31,176

**Wörter:** 26,806

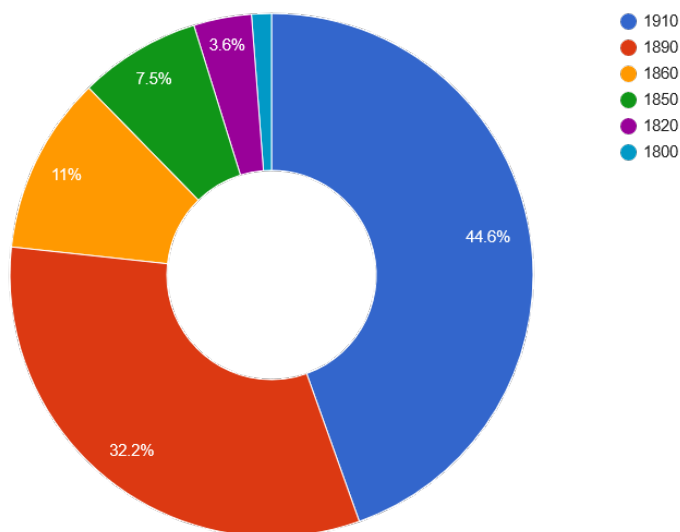
## Überblick über den Korpusaufbau

Die Grafiken beziehen sich auf die Gesamtzahl der Token. Sie wurden mit

## Zeitliche Dimension: Korpusgröße nach Jahrzehnten

Jahrzehnt	Tokenzahl
1800	382
1820	1.115
1850	2.343
1860	3.419
1890	10.024
1910	13.888

text - period



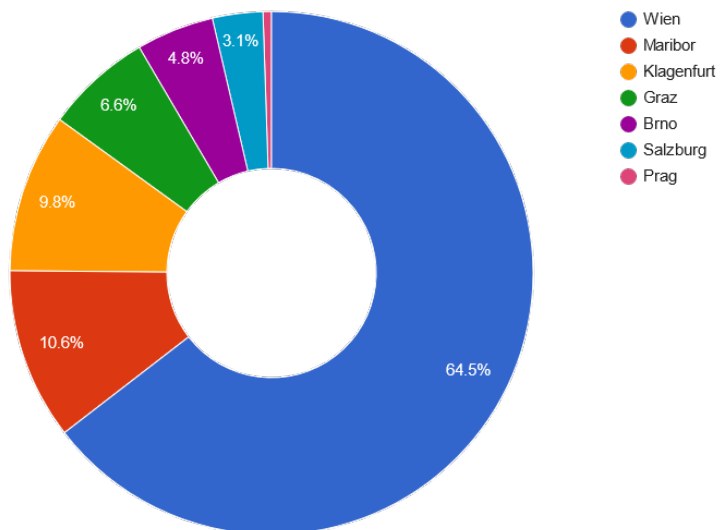
## Räumliche Dimension: Korpusgröße nach Publikationsorten

**Achtung:** In v 0.1 bestehen noch Probleme in Bezug auf die Zuweisung bzw.

Vereinheitlichung der Publikationsorte, die in der Tabelle durch Kursivsatz hervorgehoben werden!

Publikationsort	Tokenzahl
Wien	20.119
Maribor	3.302
Klagenfurt	3.070
Graz	2.053
Brno	1.497
Salzburg	971
Prag	163

text - place







# Metadaten

# Orthographische und morphologische Normalisierung

# Abweichungstypen und Abweichungen

Abkürzung	Erklärung	Beispiele
<b>A</b>	<a>-Schreibung	<i>Kommissär</i>
<b>ABK</b>	Abkürzung	<i>Abg.</i>
<b>AE</b>	<Ä>-Schreibung	<i>Aenderung</i>
<b>DAT</b>	Dativ-e	
<b>EI</b>	<ei>-Schreibung	<i>ey</i>
<b>ERR</b>	Abweichungen, die als Fehler (morphologisch oder orthographisch) eingestuft werden	<i>Agitarion</i>
<b>IE</b>	<i> statt <ie> in {ieren}/{ierung}	<i>finanziren</i>
<b>K</b>	<k>-Schreibung	<i>Collission</i>
<b>KOMP</b>	Schreibung von Komposita	
<b>LEX</b>	lexikalische Abweichungen	
<b>M</b>	<m>-Schreibung	<i>gesammt</i>
<b>MOR</b>	morphologische Abweichungen	
<b>OE</b>	<Ö>-Schreibung	<i>Oesterreich</i>
<b>ORTHKON</b>	andere Abweichungen in der Schreibung von Konsonanten	

<b>ORTHVOK</b>	andere Abweichungen in der Schreibung von Vokalen	
<b>T</b>	<th>-Schreibung	<i>Muth</i>
<b>TSCH</b>	<tsch>-Schreibung	<i>czechisch</i>
<b>UE</b>	<UE>-Schreibung	<i>Uebung</i>
<b>VOKL</b>	Vokallänge	<i>Abendmahl, Schooß</i>
<b>W</b>	<w>-Schreibung	<i>slovenisch</i>
<b>Z</b>	<z>-Schreibung	<i>Civilisation</i>