

**Datenbank:** WBIS Online

**Provider:** De Gruyter

		WBIS Online
Access	Web address, API, Dumps, offline back up copy	<ul style="list-style-type: none"> <li>• <a href="https://wbis.degruyter.com">https://wbis.degruyter.com</a></li> <li>• API access is not available</li> </ul>
Documentation	Web address	<ul style="list-style-type: none"> <li>• n/a</li> </ul>
Distribution		<ul style="list-style-type: none"> <li>• Continuously</li> </ul>
Scope	Content Purpose Field of use	<ul style="list-style-type: none"> <li>• The most comprehensive biographical database available</li> <li>• Biographical information on over 6 million people from the 8th century B.C. to the present</li> <li>• Included are 8.5 Million digital facsimile articles from biographical reference works</li> </ul>
Time, Place, Language	temporal, local reference	<ul style="list-style-type: none"> <li>• From 8th century B.C. to the present</li> <li>• International</li> <li>• Multiple languages</li> </ul>
Data type	What are the basic data types?	<ul style="list-style-type: none"> <li>• XML files</li> <li>• TIF Images</li> </ul>
Provenance, dependencies, accompanying material	original data source, manufacturer, data collection procedure, dependencies / links to other data sets / online resources, old versions	<ul style="list-style-type: none"> <li>• The definitive dataset for biographical information is kept in a proprietary XML format, and from this its online biographical datasets are derived</li> <li>• The digital facsimile articles are provided in TIF format and accompanied by metadata in the related biographical XML and a source XML entry</li> </ul>
Description Structured text data	Text markup or data structure e.g. TXT, XML, ALTO, TEI, versions	<ul style="list-style-type: none"> <li>• Each XML entry contains biographic information for a specified person and bibliographical information for the specified sources</li> <li>• Also included are XML files for the sources themselves</li> </ul>
Description of databases, tabular data	data tables, existing / recommended data splits (e.g. training / test set)	<ul style="list-style-type: none"> <li>• n/a</li> </ul>
Description of image formats	as precisely as possible (e.g. resolution, greyscale / bitonal)	<ul style="list-style-type: none"> <li>• TIF</li> <li>• 400 dpi, greyscale</li> </ul>
Standards, vocabularies	as precisely as possible: standards and vocabularies used	<ul style="list-style-type: none"> <li>• No standard vocabularies were used but WBIS Online has its own classification system, e.g. for occupations, that accompanies each entry and can be provided separately and standardized keywords are also included</li> </ul>

Data quality: OCR; missing, incorrect, redundant data, noise	For example. OCR error rate, OCR process; different raw data available? Used software?	<ul style="list-style-type: none"> <li>n/a</li> </ul>
Administration, cleanups,	e.g. handling of missing data, cutting, rescaling, NLP preprocessing, used software	<ul style="list-style-type: none"> <li>Missing data is added when identified, happens very rarely</li> </ul>
Scope /Size	size of data records	<ul style="list-style-type: none"> <li>Over 6 million biographical entries with 8.5 million facsimiles</li> </ul>
Metadata	Format/ Standards,	<ul style="list-style-type: none"> <li>Proprietary, see above</li> </ul>
Rights	licenses for metadata, full texts (TDM), rights / use (e.g. on-site, groups, scientific use)	<ul style="list-style-type: none"> <li>Institutions have rights for non-commercial use by Authorised Users of the institutions only.</li> </ul>
Ethical Issues	Personal and / or Confidential Information; Bias / representation; offensive / insulting / sensitive content	<ul style="list-style-type: none"> <li>n/a</li> </ul>
Use	Recommendations for use/ not recommended use	<ul style="list-style-type: none"> <li>All purposes of TDM</li> </ul>
Text and Data Mining	Additional costs? If so, how much? Trial possible?	<ul style="list-style-type: none"> <li>Small cost for delivering the data on hard drives.</li> </ul>

Stand 30.3.2022