

Challenges of Machine (Aided) Translation of Historical Texts

Walter V. Hahn

University of Hamburg • Department of Informatics
Natural Language Systems Group

WWW: <http://nls-www.informatik.uni-hamburg.de/view/User/WalterVHahn>
E-Mail: vhahn@informatik.uni-hamburg.de

False Assumptions from Modern Texts

Standard Texts

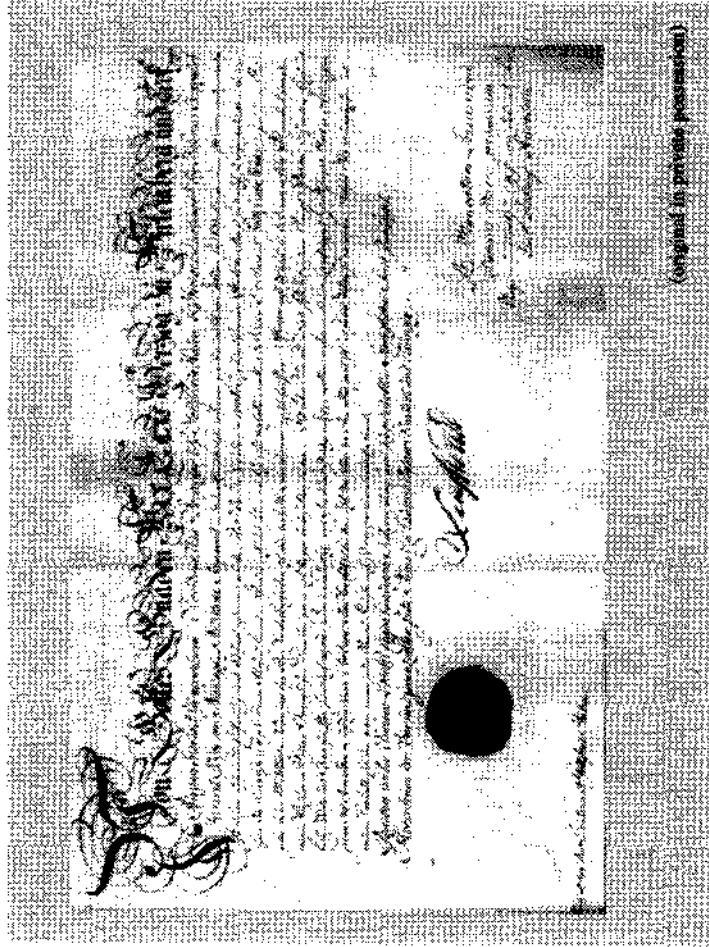
- have an author,
- are written in one language,
- have a well defined orthography,
- have definite grammatical rules,
- belong to a homogeneous linguistic and historical layer
- are readable by every speaker of the language,
- are part of a vast homogeneous corpus.

Why Translation of Historical Texts at all?

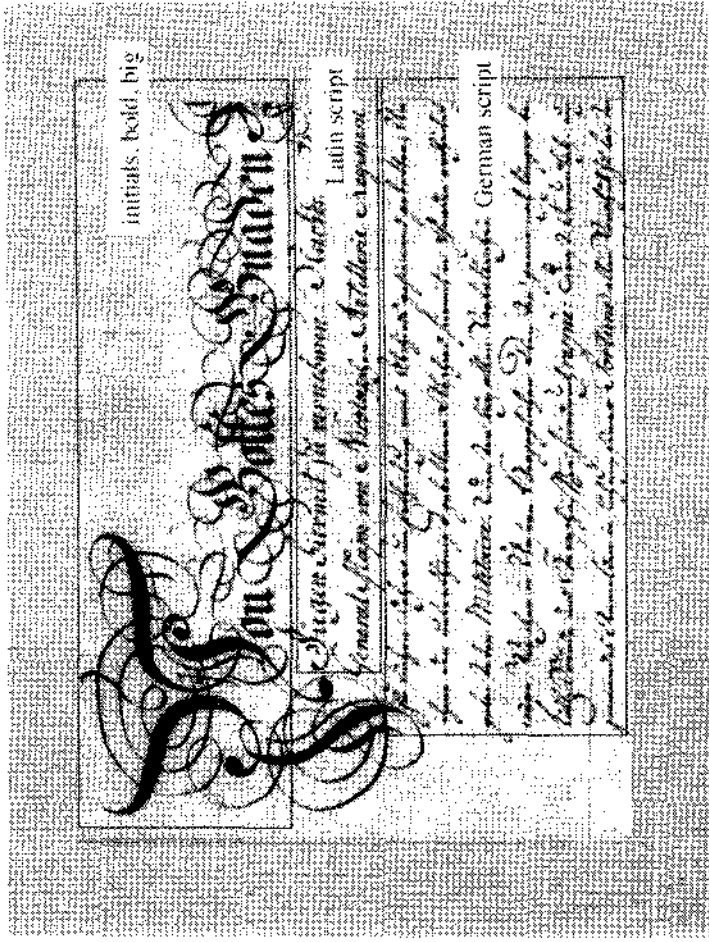
- Availability of resources for researchers in other languages.
— Example: Historical documents of European history in Polish, Russian, French, German, etc.
- Availability of resources for researchers in various fields.
— Example: Linguistics, social science, law, etc. within a cultural or historical area.
- Availability of resources for private foreign readers (e.g. out of genealogical or geographic interest).
- To allow for summarizing and indexing in other languages,
- Multilinguistic comparison of documents in academic education
- Not useful for historical linguists; however, their expertise is difficult to obtain and expensive for other scientists.

Orthography

- there is no orthography, not even writing rules within a document.
- several scripts and languages are mixed in one document.
- arbitrary abbreviations,
- illegible or ambiguous sections.
- All this makes OCR extremely unpromising. Informed transcriptions are mandatory for any sort of translation.
- Normalized orthography is standard for most historical linguistic approaches (there is no middle high German contemporary standard orthography, e.g.)
- → excludes SMT



(original in private possession)



initials, bold, big

Latin script

German script

German script

German script

German script

German script

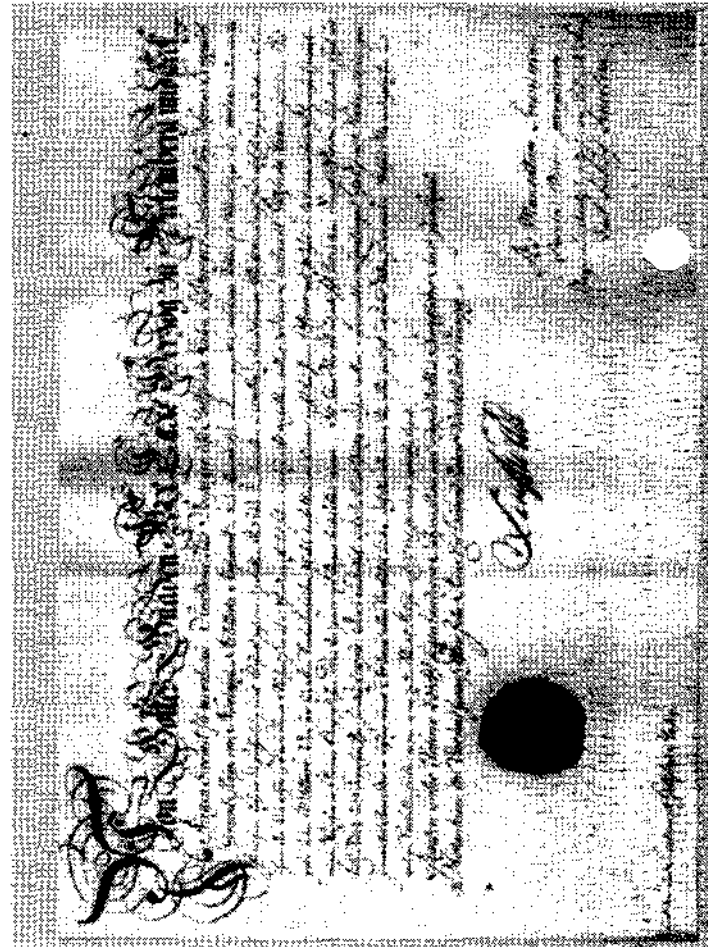
German script

German script

German script

German script

German script



Translation into What?

- A replica of old documents in other old languages is historically impossible, even a "translation" into the same modern language is problematic.
- translation in a modern foreign language distort the document because of
 - changing language
 - changing culture
 - changing domains and facts
 - changing "Zeitgeist"
- modernizing syntax in the same language might be viable.

Modern „Translation“

Verfügung

Der Antragsteller Gottlieb Hahn, derzeit Leutnant beim Artillerieregiment des Generalmajors von Nicolay möchte seine Laufbahn in einer anderen Position fortsetzen und daher kündigen. Ich will der Verbesserung seiner Karriere nicht im Wege stehen und daher der Auflösung des Vertrages zustimmen. Es sei vermerkt, dass Leutnant Hahn während der Tätigkeit bei uns bei allen Gelegenheiten seine Dienstverpflichtungen erfüllt hat.

Karl von Württemberg etc.

Is this the source for an English translation?

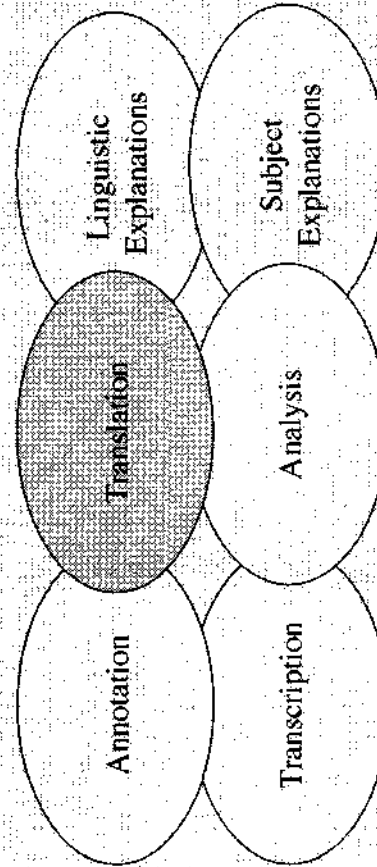
- the text is shorter
- sentences are shorter
- no unknown words, etc.
- completely different text form

Result

So, no translation?

- not in the modern sense
- a translated transcription (word-to-word) does not meet the needs of most users,
- a replica of old texts in other old languages is historically impossible,
- A solution for scientists in other languages and other fields is
 - an integrating tool for machine aided translation, which allows for
 - a modern narrow paraphrase in L2, and
 - a broad documentation of the source text, target text, and translation.

Integration



Functional Sketch of a Translation Tool for Historical Texts

- Displays to the translator
 - Digital source, ✓
 - Domain information,
 - Lexical correspondences via one single ontology, ✓
 - Target text marks for source annotations of languages, font styles, size, script, etc.
- Displays to the user
 - Source and target text, ✓
 - Interactive target text comments
 - Meta data about the translation, ✓
 - Digital source, ✓
- Explanation generator for the translator, ✓
- XML annotation support even for document images, ✓

✓ = facility available in other contexts

Sketch of System Architecture

