JOST GIPPERT,
The TITUS Project. 25 years of corpus building in ancient languages,

in:
*Perspektiven einer corpusbasierten historischen Linguistik und Philologie. Internationale Tagung des Akademienvorhabens „Altägyptisches Wörterbuch" an der Berlin-Brandenburgischen Akademie der Wissenschaften, 12. – 13. Dezember 2011*, herausgegeben von Ingelore Hafemann, Berlin 2013, S. 169-192.

BERLIN-BRANDENBURGISCHE AKADEMIE DER WISSENSCHAFTEN

**Thesaurus Linguae Aegyptiae**

**4**

BERLIN 2013

BERLIN-BRANDENBURGISCHE AKADEMIE DER WISSENSCHAFTEN

# Perspektiven einer corpusbasierten historischen Linguistik und Philologie

Internationale Tagung des Akademienvorhabens „Altägyptisches Wörterbuch" an der Berlin-Brandenburgischen Akademie der Wissenschaften, 12. – 13. Dezember 2011

herausgegeben von Ingelore Hafemann

BERLIN

2013

# INHALTSVERZEICHNIS

VORWORT

Die internationale Tagung „Perspektiven einer corpusbasierten historischen Linguistik und Philologie" vom 12. – 13. Dezember 2011 am Akademienvorhaben „Altägyptisches Wörterbuch" der Berlin-Brandenburgischen Akademie der Wissenschaften (BBAW) war dem Thema des Aufbaus und der Nutzungsperspektiven elektronischer Textcorpora und Wörterbücher in den historischen Sprachen gewidmet. Die Teilnehmer, Vertreter der Ägyptologie, der Hethitologie, Indogermanistik sowie Referenten aus der historischen Lexikographie des Mittel- und Frühneuhochdeutschen und des Altfranzösischen diskutierten vor allem über die Veränderungen, die mit dem Einsatz elektronischer Erfassungs- und Verarbeitungsprozeduren einhergehen. Vertreter der Computerlinguistik vom „Zentrum Sprache" der BBAW wurden in die Diskussionen einbezogen. Dort beschäftigt man sich seit Jahren mit dem Aufbau großer elektronischer Textcorpora (DWDS), darunter auch solcher, die historische Texte (DTA) für die elektronische Nutzung ermöglichen.

Die größte Herausforderung dieser neuen elektronischen Corpora und Wörterbücher ist es, sowohl den Methoden und damit den wissenschaftlichen Ansprüchen der traditionellen Philologie und Lexikographie unbedingt verpflichtet zu bleiben als auch neue Gebiete wie die Corpus- und Computerlinguistik für die historischen Sprachen zu öffnen. Die Teilnehmer haben gemeinsam und disziplinenübergreifend die Möglichkeiten und Grenzen der Datenerfassung, ihrer Präsentation und den Nutzen neuer Auswertungsprozeduren diskutiert.

Unter dem ersten Thema „Historische Corpusprojekte – synchron und diachron" wurden elektronische Corpora vorgestellt und ein intensiver Austausch darüber geführt, welche Datenstrukturen die linguistischen Inhalte in adäquater Weise abbilden. Wichtig war die Frage, auf welche Resonanz diese elektronischen Corpora bei den Nutzern gestoßen sind und welche Erwartungen und Anforderungen aus den verschiedenen Fachdisziplinen an die Projekte herangetragen werden. Der Austausch über Nutzungsperspektiven elektronischer Corpora schloss auch die Diskussion über die Erarbeitung projektübergreifend einsetzbarer Standards der Codierung und Strukturierung historischer Textdaten mit ein. Hinsichtlich einer mittel- und langfristigen Nutzbarkeit sowie einer langfristigen Datensicherheit stehen solche Fragen zunehmend im Focus und einige aktuelle Initiativen dazu wurden vorgestellt. Spezielle technische Aspekte

elektronischer Datenerfassung und automatischer Analyse- und Speicherungsverfahren elektronischer Textdaten konnten am letzten Tag als ein Themenschwerpunkt mit den Programmierern diskutiert werden.

Ein zweiter Schwerpunkt waren konkrete Fragestellungen aus der historischen Lexikographie und diachronen Textanalyse. Für das Ägyptische ist der diachrone Ansatz auf Grund der über viertausendjährigen Textüberlieferung von großer Relevanz. Themen wie historischer und/oder textgattungsspezifischer Wortgebrauch, die Erarbeitung diachroner Wortlisten und Aspekte des kontaktindizierten Sprachwandels konnten disziplinübergreifend zwischen den Ägyptologen und den Kollegen der historischen Lexikographie des Mittel- und Frühneuhochdeutschen und des Altfranzösischen behandelt werden.

Mit dem Abendreferenten Gregory Crane, dem Begründer der „Perseus Digital Library", wurde ein breites Publikum angesprochen. In seinem Vortrag hat er noch einmal die hohe Relevanz und die neuen Möglichkeiten der Einbeziehung zahlreicher Wissenschaftler und einer interessierten Öffentlichkeit in die Projektarbeit demonstriert, die das Internet auf völlig neue Weise eröffnet hat. Die Herausgeberin ist sehr froh, seinen programmatischen Beitrag zu diesem Thema, dessen schriftliche Form er gemeinsam mit Alison Babeu erarbeitet hat, ebenfalls in diesem Band präsentieren zu können.

Allen Autoren dankt die Herausgeberin für ihre anregenden Diskussionen und die qualitätvollen Beiträge in diesem Band.

THE TITUS PROJECT

25 YEARS OF CORPUS BUILDING IN ANCIENT LANGUAGES

JOST GIPPERT

The article summarizes the contents and the structurtal premises of the "Thesaurus Indogermanischer Text- und Sprachmaterialien" (TITUS), focussing on search functions and facilities and questions of the encoding of ancient languages written in various scripts. Examples are taken from Tocharian, Greek, Vedic Sanskrit, and other ancient Indo-European languages covered by TITUS.

In September 1987, a group of Indo-Europeanists decided to join efforts in the digitization of primary sources that are essential for their research, by creating a common pool of the electronic texts to be prepared. Eversince,[1] the text pool has developed into a comprehensive retrieval system covering a large amount of relevant materials. The scope, the contents and the structural premises of the "Thesaurus Indogermanischer Text- und Sprachmaterialien" (TITUS) are summarized in the following pages.[2]

1. Since its foundation, the primary goal of the TITUS project consisted in the compilation of a comprehensive text database of ancient Indo-European languages that were not covered by concurrent projects such as the *Thesaurus Linguae Graecae*.[3] To reach this aim, a practical way of cooperation was decided upon: everybody who was able to contribute to the database was granted, as a member of the TITUS team, access to the complete database. In the 1980ies, this still presupposed data exchange via floppy or, later, compact disks, as internet facilities were not yet available in a sufficient way. Nevertheless, as early as 1988 the complete text of the Old Indic Rigveda Samhita, which had been electronically prepared as a text file of ca. 1.5 MB by H.S. Ananthanarayana under the supervision of W.P. Lehmann in Austin/Texas, was successfully transferred via a data line from the USA to the Berlin Free University, which hosted

---

[1] The project was announced under the title "Thesaurus altindogermanischer Textcorpora auf Datenträgern" in: *Die Sprache* 32/2, 1987 [1988], 151t.

[2] For previous accounts of the TITUS project cf. GIPPERT (1995a; 1995b; 1996; 2001; 2010).

[3] Project of the University of California at Irvine; cf. http://stephanus.tlg.uci.edu/.

the data pool then. By 1994, when the facilities of the internet emerged, the exchange of data was put on an online basis by establishing an FTP server at the University of Frankfurt, and soon after, the first web pages of the project were launched under the new name of "TITUS" which had meanwhile been agreed upon by the members.[4] Since 1996, the TITUS project has been promoting the use of Unicode to ensure a reliable encoding of its data, and the independent web server of the project established then[5] was one of the first sites world-wide to make a considerable amount of textual data available in this way of encoding. Thanks to a generous grant of the WordCruncher Company, the project was able in 1997 to install, along with its web site, a special "WordCruncher" server for the search and retrieval of data from the database. This service has been maintained until recently but has now been given up as most of the facilities it provides have meanwhile been implemented in an SQL-based online retrieval engine that has been publicly accessible since 2000.[6] Today, the TITUS database comprises not only corpora of ancient Indo-European languages such as Avestan, Vedic Sanskrit, Phrygian, or Umbrian, many of them covering the complete textual heritage of the languages involved, but also materials in more recent Indo-European as well as neighbouring languages, among them the largest corpus of Old and Middle Georgian available world-wide.[7] Many of the TITUS corpora have been the basis for more specialized corpus projects such as, e.g., the Referenzkorpus Altdeutsch project,[8] which aims at a full annotation of all textual materials in Old High German and Old Saxon; the Sanskrit Library project at Brown University, which aims at providing grammatical and other information pertinent to Sanskrit texts;[9] or the National Corpus of the Georgian Language, an international project aiming to cover the complete

---

[4] The clumsy URL was http://www.rz.uni-frankfurt.de/home/ftp/pub/titus/public_html/.

[5] URL: http://titus.uni-frankfurt.de/.

[6] URL: http://titus.fkidg1.uni-frankfurt.de/search/query.htm.

[7] Cf. http://titus.uni-frankfurt.de/texte/texte.htm for a full account of available corpora and texts.

[8] A common project of the universities of Berlin (Humboldt), Frankfurt and Jena, financed by the Deutsche Forschungsgemeinschaft since 2009 and part of the initiative "Deutsch-diachron-digital"; cf. http://www.deutschdiachrondigital.de.

[9] Cf. http://sanskritlibrary.org.

written history of Georgian from the 5th century A.D. up to the present day.[10]

2. With the establishment of the WordCruncher server in 1997, the TITUS project has moved far away from its original concept of being a mere exchange base of text files. Instead, the focus has shifted towards providing sophisticated search facilities within and across the text corpora, thus supporting online research into the languages and literatures in question. A few examples may suffice to illustrate the facilities that have been developed meanwhile.

2.1 One of the Indo-European languages for which TITUS may claim to cover the complete textual heritage in its corpus, is Tocharian, a language that was spoken in two different varieties in East Turkestan in the first millennium of our era. The textual remnants of the two Tocharian varieties (East- or A- and West- or B-Tocharian) are contained in ca. 5,000 manuscripts written in a "Northern" type of Brahmi script that were found in a region extending from Kucha to Turfan and Dunhuang along the Silk Road.[11] The largest part of these manuscripts is preserved in the Turfan Collection of the Berlin-Brandenburg Academy of Sciences (BBAW) today (ca. 4,000 manuscripts);[12] other major collections are hosted in London, Paris, and St. Petersburg. Within the TITUS project, work on the Tocharian manuscripts started in 1996 with the digitization of the printed editions of A- and B-Tocharian texts of the Berlin collection (in Romanized transcription), which formed the foundation of the emerging corpus. In the same year still, TITUS and the BBAW agreed upon preparing a complete set of digital images of the Tocharian manuscripts from Berlin to provide them online along the transcribed texts;[13] this endeavour, which was kindly supported by T. Tamai, resulted, in 2000, in one of the first frame-based online editions providing images, transcribed texts, and metadata as to each manuscript side by side. Today, this online-edition comprises the complete set of Berlin manuscripts including the ca. 3,000 hitherto unpublished

---

[10] Cf. http://georgiannationalcorpus.ac.ge.

[11] Cf. http://titus.uni-frankfurt.de/didact/karten/turkstan/turkst.htm for a map showing the locations.

[12] Cf. http://www.bbaw.de/en/research/turfanforschung as to the Turfan Studies project of the BBAW.

[13] Cf. GIPPERT (1997 and 1998) as to the technical foundations of the digitization project.

fragments, all manually transliterated by T. Tamai (cf. Fig. 1 showing a screen-shot of the site).[14]

2.2 In parallel to the online edition of the Berlin collection, which provides access to the corpus only via the catalogue number of a given manuscript,[15] the Tocharian data of all major collections have been prepared for a word-form based retrieval via the TITUS search engine.[16] This is built upon a more fine-grained referencing system where every single line of a manuscript (page) can be addressed directly (cf. Fig. 2 showing line 3 of the recto of the A-Tocharian fragment THT 634, which is part no. 1a in the edition by SIEG & SIEGLING 1921). To facilitate investigations into the paleography of the Brahmi script used for Tocharian, each line is further provided with an *akṣara*-based transliteration alongside the "normal" word-based transcription (as visible in Fig. 2). On the basis of a preindex-ation of the complete corpus, this allows for searching for both word-forms and individual *akṣaras*, either by clicking upon an item as displayed in the text or by using a query form. E.g., clicking upon the word-form *kumseñc* 'they come' in the given text line invokes (via a javascript underlying the word) a query for all (eight) occurrences of the same word-form throughout the A-Tocharian corpus,[17] which is output as a list of keyword-in-context entries with full referentiation of the text passages in question (cf. Figs. 3 and 4). Each entry is linked to the corresponding text passage so that the wider context can be accessed at will (cf. Fig. 5 showing the context of THT 935 = 302a, line 5). Note that in the transliteration, "Ä\" stands for the combination of the diaeresis-like vowel mark (which usually stands

---

[14] Cf.       http://titus.fkidg1.uni-frankfurt.de/texte/tocharic/thtframe.htm;       most elements of the edition were prepared in cooperation by K. Kupfer and T. Tamai.

[15] In the TITUS edition, the Berlin manuscripts are referenced according to their catalogue number in the Turfan Archive ("THT"). Of the 4074 manuscripts listed there, nos. 1–633 are B-Tocharian, and nos. 634–1099 are A-Tocharian (numbered 1–467 in the printed edition by SIEG & SIEGLING 1921). Several manu-scripts have been missing since the Second World War; in some cases, digitized images could be provided from existing photographs.

[16] Cf. http://titus.uni-frankfurt.de/texte/texte2.htm#toch; for the time being, access to the B-Tocharian corpus, which is still under construction, is restricted to TITUS members and other registered users (cf. http://titus.uni-frankfurt.de/titusstd.htm for a form to apply for registration).

[17] The javascript causes an SQL-query to be sent to the following ASP script: http://titus.fkidg1.uni-frankfurt.de/database/titusinx/titusinx.asp?LXLANG=58285&LXWORD=kumseF100c&LCPL=0&TCPL=0&C=H&PF=26.

for the shewa vowel rendered as *ä* in transcriptions of Tocharian) with a *virāma* in word-final position (i.e., *kumseñc$^{ä}$* in the transcription system applied in the editions by Sieg and Siegling); in the word-based search, it is ignored. In a similar way, the so-called "Fremdzeichen" are represented by capital letters, with "A" standing for their inherent vowel (i.e., KA etc. stands for *ḵa* etc. in Sieg/Siegling's transcription); in the word-based search, *A* is treated as an equivalent of *ä* and the difference between "Fremdzeichen" and "Indian" akṣaras is ignored. The unsyllabic *u* vowel indicated by subscript *u* with a bent line above in the editions is represented by *ù* in the corpus; in the word-search, this is treated as being equivalent with plain *u*. In the *akṣara*-based transliteration, * stands for a (missing or unreadable) complete *akṣara* and +, for a (vocalic or consonantal) element of an *akṣara*; ˆ stands for a word boundary within an *akṣara* (ignored in the search). All this is warranted by a specific structure of the underlying relational database, which contains "normalised" variants of the word-forms wherever applicable (cf. Fig. 6).[18]

2.3 A more flexible and powerful query method than the hyperlink-based retrieval is provided via special input forms. In general, the TITUS "Search Engine" comprises two different methods of input-based access to its data, one yielding a list of word-forms matching the query input, and one, the keyword-in-context output of occurrences as shown above (cf. Fig. 7).[19] In both input forms, the language of the search must be determined first, either specifically (e.g., "Tocharian A" as in Fig. 8) or more generally (e.g., "Tocharian" as in Fig. 9). The word-form to be searched can then be entered either in toto or partially, in exact Unicode encoding or in a substitutional plain-ASCII representation (or in a mixed representation), and with two types of "wild cards" replacing explicit characters: the question mark, "?", stands for one single character, and the asterisk, "*", for any sequence of characters (including zero). E.g., the word-form *kumseñc* can be entered as such or in the form *kumsen~c* (with the diacritic adscripted to match ASCII-based keyboards, cf. Fig. 10), and *kumseñc* will also be found if the query string is reduced to

---

[18] The database used at present is IBM DB-2 Express version 9.1, a powerful and yet free SQL-based system with full Unicode support.

[19] Both query types are accessible via http://titus.fkidg1.uni-frankfurt.de/search/query.htm. – A third query type, styled "unspecified", consists of a mere link to a Google search over the TITUS site.

*k?ms\*c*, i.e., with one character between *k* and *m* and any sequence of characters between *s* and *c* (cf. Fig. 11), or kums\*, i.e., with any sequence of characters following *kums*, as the resulting word-list shows (cf. Fig. 12). Similarly, the word-form *NAmseñc* 'they revere', which occurs in the same line of THT 634 as *kumseñc*, can be retrieved by entering *NAmseñc, nämseñc, n\*ms??c, näm\**, etc. (cf. Fig. 13 showing the variant spellings in parentheses). In the word-list output, every word-form is provided with a hyperlink to the relevant context query; this means that by clicking upon *NAmseñc* or *nämseñc* in the list, all 3 occurrences of the word-form (the number of occurrences is indicated in square brackets for each list entry) will be listed in an extra window (cf. Fig. 14).[20] Of course, the same result can also be achieved with the "context output" form, entering, e.g., *n?mse?c* (cf. Fig. 15). If the query string has more than one match as, e.g., in the case of *kums\** (cf. above), the occurrences of the different matching forms will be output in alphabetic order (cf. Fig. 16 showing first an occurrence of *kumsanträ*, 3rd pers.pl.pres.ind.med., then one of *kumse*, 3rd pers.pl.pres.ind.act., shortened form). In addition, in the case of verb forms, the header of the output list indicates the underlying root (if determinable), again provided with a hyperlink (cf. Fig. 17); this leads to a special table which illustrates for every Tocharian verb which of its paradigm forms are attested in the two dialects (cf. Fig. 18).

2.4 As was stated above, the input of search strings in the query forms can be exactly "as is", i.e., in Unicode encoding, or in a substitutional plain-ASCII format with adscript diacritics. This is true not only for the input of Latin-based scripts (or transcriptions) but also for other scripts. Thus, e.g., to search for the attestations of Greek ἄνδρα (acc.sg. of ἀνήρ 'man'), both the Greek spelling and Latin *andra* can be entered (cf. Figs. 19 to 21). Note that the entry of Greek diacritics is not necessary as unaccented variants are stored in the database for all word-forms; this means that all occurrences of ἄνδρα will also be found in a search for (less specified) ανδρα. This is even true for the same word spelt with an initial capital (Ἄνδρα) or with an acute on the word-final vowel (ἄνδρά, to be expected in the

---

[20] This is again invoked by a javascript which sends the SQL-query to the following URL: http://titus.fkidg1.uni-frankfurt.de/database/titusinx/titusinx.asp?lxlang=941&lxword=N4100mseF100c&LCPL=1&TCPL=1&C=H.

position preceding a clitic), which are matched by ανδρα (and *andra*) but not by ἄνδρα (cf. Figs. 21 and 22).

2.5 A special feature of the context-related search is the "combined search" function. Up to four query patterns (word-forms, stems, etc.) can be entered in parallel for a search of co-occurrences in a given context; cf., e.g., Figs. 24 and 25 showing a combined search for *thaz* 'the' and *uuort* 'word' in the Old High German corpus. The amount of context envisaged here can be determined by the user. If the "distance" is set to "0" (the default setting), the context in question is the lowest reference level of a given text (usually a sentence, a verse or a line); in the given example, this yields 111 co-occurrences of *thaz* and *uuort,* irrespective of the order of the two words (and including spelling variants such as *tház* and *uuórt*). Setting the distance to "1 - exact" (cf. Fig. 26) yields but 33 co-occurrences, with *uuort* immediately following *thaz*. (cf. Fig. 27).

2.6 A feature that has only been implemented for Old Indic (Sanskrit) and Avestan so far is the "thesaurus search" function. Different from the word-form based queries illustrated above, this function admits of searching for complete paradigms of words irrespective of a common ("matching") string structure of the individual word-forms. Starting, e.g., from *r̥tvíjo* as the genitive singular case form of the Vedic noun *r̥tvíj-* 'priest' (cf. Fig. 28), the output displays all occurrences of all case forms of this word as met with in the corpus[21] beginning with the nominative singular *r̥tvík*, provided with a grammatical analysis of each form[22] and a German translation of the respective lemmata[23] (cf. Fig. 29).

3. It is obvious that the latter type of retrieval presupposes a thorough modelling of the morphology of the language concerned. To implement similar facilities for all languages covered by TITUS

---

[21] Including spelling variants (caused by sandhi and accentuation), the list comprises the following forms: *r̥tvik, r̥tvík, r̥tvig, r̥tvíg, r̥tvijam, r̥tvíjam, r̥tvijam̐, r̥tvíjam̐, r̥tvijā, r̥tvíjā, r̥tvijah, r̥tvíjah̠, r̥tvijas, r̥tvíjas, r̥tvijaś, r̥tvíjaś, r̥tvijo, r̥tvíjo, r̥tvija, r̥tvíja, (r̥tvijah̠, r̥tvijah̠), r̥tvijām, r̥tvíjām, r̥tvijām̐, r̥tvíjām̐.*

[22] The analysis was worked out by R. Gehrke in the course of the AUREA project ("Avesta und Rigveda: Elektronische Analyse") financed by the Deutsche Forschungsgemeinschaft in 1997-1999; cf. http://titus.uni-frankfurt.de/curric/aurea/aurea.htm.

[23] Based upon the dictionary by K. MYLIUS (1992).

therefore means an immense task for the project that still has to be undertaken. Another task for the future that can be envisaged today concerns improvements in the rendering of non-Latin scripts as in the case of the cuneiform inscriptions of Old Persian for which a Unicode-based encoding in the original script has recently been provided by A. Sarhadi and M. Esnaashari (cf. Fig. 30).[24] In some cases, this is still hampered by the fact that the corresponding code-points of the Unicode standard are not yet available; e.g., it would be possible now to encode the Avestan texts in the original script[25] but for Middle Persian (Pahlavi) passages that are often met with in Avestan contexts, a Unicode rendering is not yet available. As in former cases, the members of the TITUS project are ready to support the standardisation process with their expertise.

---

[24]  Cf. http://titus.uni-frankfurt.de/texte/etcs/iran/airan/apers/apers.htm.

[25]  Cf. GIPPERT (forthcoming) as to details.

BIBLIOGRAPHY

GIPPERT, J., 1995a: TITUS. Das Projekt eines indogermanistischen Thesaurus, in: *LDV-Forum* 12/1, 35-47.

GIPPERT, J., 1995b: TITUS. Von der Keilschrifttafel zur Textdatenbank, in: *Forschung Frankfurt* 4/1995, 46-56.

GIPPERT, J., 1996: TITUS – Alte und neue Perspektiven eines indogermanistischen Thesaurus, in: *Studia Iranica, Mesopotamica et Anatolica* 2, 1996 [1997], 49-76.

GIPPERT, J., 1997: Digitization of Tocharian Manuscripts. Short notice about a new project, in: *Tocharian and Indo-European Studies* 7, 265-266.

GIPPERT, J., 1998: Digitization of Tocharian Manuscripts from the Berlin Turfan Collection, in: *Manuscripta Orientalia. International Journal for Oriental Manuscript Research* 4/1, 49-57.

GIPPERT, J., 2001: Der TITUS-Server: Grundlagen eines multilingualen Online-Retrieval-Systems, in: WILLÉE, G. *et al.* (eds.), *Computerlinguistik. Was geht, was kommt? / Computational Linguistics. Achievements and Perspectives. Festschrift für Wilhelm Lenders,* Bonn 2002, 81-85.

GIPPERT, J., 2010: Manuscript Related Data in the TITUS Project, in: *Comparative Oriental Manuscript Studies Newsletter* 1, 2011, 7-8.

GIPPERT, J., forthcoming: The Encoding of Avestan: Problems and Solutions, to appear in: *Journal for Language Technology and Computational Linguistics,* 2012.

MYLIUS, K., 1992: *Wörterbuch Sanskrit-Deutsch,* 4. Auflage, Leipzig [u.a.].

SIEG, E. & W. SIEGLING, 1921: *Tocharische Sprachreste,* Band I: *Die Texte,* A: *Transcription,* B: *Tafeln,* Berlin [u.a.].

FIGURES



*Figure 1*



*Figure 2*

*Figure 3*



*Figure 4*

*Figure 5*

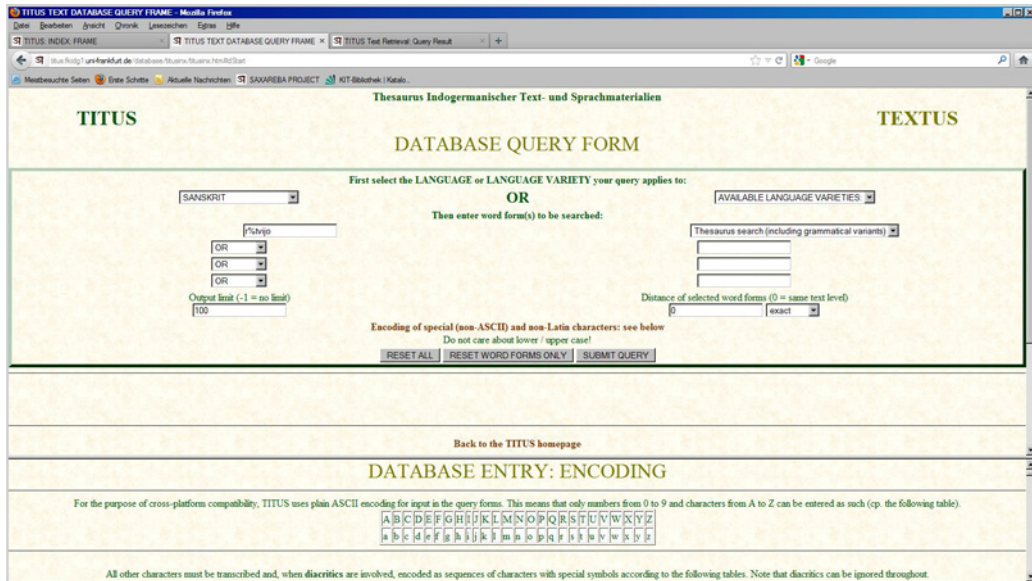| 941 | 57 | tlu | | 26 | 6 | 232 |
|---|---|---|---|---|---|---|
| 941 | 57 | ne | | 26 | 6 | 233 |
| 941 | 57 | tsra | | 26 | 6 | 236 |
| 941 | 57 | ṣi | | 26 | 6 | 237 |
| 941 | 57 | śśi | | 26 | 6 | 238 |
| 941 | 57 | mā | | 26 | 6 | 239 |
| 941 | 57 | kni | | 26 | 6 | 240 |
| 941 | 57 | ṣpa | | 26 | 6 | 241 |
| 941 | 57 | lntu | | 26 | 6 | 242 |
| 941 | 57 | tsra | | 26 | 6 | 243 |
| 941 | 57 | ṣi | | 26 | 6 | 244 |
| 941 | 57 | śśi | | 26 | 6 | 245 |
| 941 | 57 | mā | | 26 | 6 | 246 |
| 941 | 57 | K\ | k\ | 26 | 6 | 247 |
| 941 | 57 | * | | 26 | 6 | 248 |
| 941 | 57 | * | | 26 | 6 | 249 |
| 941 | 56 | ṣṣeñ | | 26 | 7 | 250 |
| 941 | 56 | NAmseñc | nämseñc | 26 | 7 | 252 |
| 941 | 56 | yäsluS | yäsluṣ | 26 | 7 | 253 |
| 941 | 56 | tsraṣisac | | 26 | 7 | 254 |
| 941 | 56 | kumseñc | | 26 | 7 | 255 |
| 941 | 56 | yärkanT | yärkant | 26 | 7 | 256 |
| 941 | 56 | tsraṣisac | | 26 | 7 | 257 |
| 941 | 56 | tsraṣiñ | | 26 | 7 | 259 |
| 941 | 56 | waste | | 26 | 7 | 260 |
| 941 | 56 | wrasaśśi | | 26 | 7 | 261 |
| 941 | 57 | ṣṣe | | 26 | 7 | 262 |

*Figure 6*

*Figure 7*



*Figure 8*

*Figure 9*



*Figure 10*

*Figure 11*



*Figure 12*

*Figure 13*



*Figure 14*

*Figure 15*



*Figure 16*

*Fugure 17*



*Figure 18*

*Figure 19*



*Figure 20*

*Figure 21*



*Figure 22*

*Figure 23*



*Figure 24*

*Figure 25*



*Figure 26*

*Figure 27*



*Figure 28*

*Figure 29*



*Figure 30*