# From Inscription to Semi-automatic Annotation of Maya Hieroglyphic Texts

## Christian M. Prager, Cristina Vertan

Rheinische Friedrich-Wilhelms-Universität Bonn, Berlin-Brandenburgische Akademie der Wissenschaften
Oxfordstrasse 15, 53111 Bonn, Unter den Linden 8, 10117 Berlin
cprager@uni-bonn.de, vertan@uni-hamburg.de

## Abstract

The Maya script is the only readable autochthonous writing system of the Americas and consists of more than 1000 word signs and syllables. It is only partially deciphered and is the subject of the project "Text Database and Dictionary of the Classic Maya"[1]. Texts are recorded in TEI XML and on the basis of a digital sign and graph catalog, which are stored in the TextGrid virtual repository. Due to the state of decipherment, it is not possible to record hieroglyphic texts directly in phonemically transliterated values. The texts are therefore documented numerically using numeric sign codes based on Eric Thompson's catalog of the Maya script. The workflow for converting numerical transliteration into textual form involves several steps, with variable solutions possible at each step. For this purpose, the authors have developed ALMAH "Annotator for the Linguistic Analysis of Maya Hieroglyphs". The tool is a client application and allows semi-automatic generation of phonemic transliteration from numerical transliteration and enables multi-step linguistic annotation. Alternative readings can be entered, and two or more decipherment proposals can be processed in parallel. ALMAH is implemented in JAVA, is based on a graph-data model, and has a user-friendly interface.

**Keywords:** Digital Epigraphy, Linguistic Annotation, Maya Hieroglyphic Writing

## 1. The Maya and Their Writing System



Figure 1: Detail of hieroglyphic inscription carved on Stela 2 from Dos Pilas, Guatemala. Karl Herbert Mayer, 1978 (CC BY 4.0).

This paper addresses the semi-deciphered written language of the Classic Maya, whose cultural area extended over territories of the present-day nation states of Mexico, Guatemala, Belize and Honduras. Maya hieroglyphic writing was used between between 300 BC and AD 1500. It is a mixed, morphographic and syllabic writing system comparable to Egyptian hieroglyphs or cuneiform of Mesopotamia. As a visual language, Classic Mayan survived in more than ten thousand texts (Houston and Martin, 2016). Most sources exhibit biographical information on political elites and provide written evidence for political relations between the more than sixty ruling dynasties (Martin, 2020). The inscription's focus lies on religious and political events that marked elite daily life (Stuart, 1998). Maya kings made their public claim to power through writing and iconography. In this context, written and pictorial records, especially those on stone

(Figure 1), wood, ceramics, bone and fig-bark paper, not only served as vehicles for cultural memory at the time, but today form the most important material basis for reconstructing elite history and culture. Furthermore, most texts display calendar dates that record exact sequences of events, providing not only historical insights, but also unique data on the history of Maya writing and language.

The Maya writing system is considered a hieroglyphic script because of the iconic character of its more than 1,000 graphs depicting figurative and abstract objects from the natural environment, flora, fauna, material culture, human and animal body parts, or portraits of supernaturals. Typologically, it is a logographic-syllabic writing system with two basic, functional sign types: syllabic signs and logographs (Grube, 1994). The latter denote concrete words and bound morphemes, whereas the former represent vowels and open syllables and thus permit syllabic spellings of lexical and grammatical morphemes. In addition, syllabic signs were used as phonetic complements that were pre- or post-fixed to morphographs. Thus, it was possible to write words entirely with syllabic signs, by using morphographs alone or by combining the two sign types. To create hieroglyphic text, graphs were squeezed and stacked into quadratic or rectangular blocks (Figure 1). It is the basic structural unit of a Classic Mayan text that usually corresponded to the emic concept of a word. The blocks were usually arranged in double columns to be read from left to right and from top to bottom. Researchers identified a range of calligraphic principles with which not only individual graphemes, but also Classic Mayan words could be realized in a variety of ways (Zender, 1999). The high aesthetic quality of an overall work was meant to catch the eye, monotony, conformity and repetition, it seems to today's viewer of the hieroglyphs, were to be avoided by applying a common set

of graphetic and graphemic principles described by Prager and Gronemeyer (2018)

## 2. The Digital Exploration of Classic Mayan

Maya writing and language forms the subject of the long term research project "Text Database and Dictionary of Classic Mayan"[2] (Prager et al., 2018). The project's goal is to compile a text database and a dictionary of Classic Mayan. Such efforts would permit a detailed and precise investigation of the Classic Mayan literary language, for instance by comparing text passages using co-text and co-occurrence analysis. Until now, such systematic and cross-linked work with text, image, and information carriers was impossible, because the necessary technology did not yet exist in this field of research. This undertaking can only be initiated using methods and technologies from the digital humanities, whereby the project is drawing upon tools and technologies that are already available in the virtual research environment TextGrid or that are being developed and implemented in the context of the project, e.g. an annotator for the linguistic analysis of the Maya hieroglyphs (Grube et al., 2014).

For this purpose, the inscribed artefacts and their illustrations are currently being researched in the literature, in archives and photo collections and are made accessible with the help of digital methods and technologies in the virtual research environment TextGrid (Prager, 2015). At the present time, about one third of the known text carriers including their metadata have been recorded, and the relevant literature has been documented. Images of the texts are continuously added to the project's online "Maya Image Archive"[3] (Diederichs et al., 2020). In the long term, research data will be published in the TextGrid repository, including persistent identifiers, and made freely available through a research portal[4]. In cooperation with the Bonn University and State Library the project is also publishing selected content from the TextGrid repository in the "Archive of Maya Hieroglyphic Texts", as part of ULB's Digital Collections[5]. In the past years the project started to transfer the hieroglyphic texts into an XML/TEI-based machine-readable format[6]. For this purpose, the project has simultaneously implemented a digital inventory for the signs in Maya script, which currently comprises almost 1000 elements (Diehr et al. 2018, 2018). Due to the vague state of decipherment of the Maya script, it is not possible to record hieroglyphic texts in phonemically transliterated values, in contrast to comparable projects in Egyptology or cuneiform research (Diehr et al., 2019). Therefore, Maya texts are numerically transcribed using sign codes adapted from Eric Thompson's catalog of Maya hieroglyphs (1962). Since the start of the project this catalog has been critically scrutinized and supplemented with signs that were not included in the original work (Prager and Gronemeyer 2018). Thompson's inventory is still regarded as the standard work for Maya epigraphers, which is why the project has been adopting his nomenclature while removing misclassifications and duplicates, merging graph variants under a common nomenclature, and adding new signs or

allographs to the sign index in sequence (Diehr et al., 2018).

In order to generate linguistic documents from these numerically encoded hieroglyphic texts the project in cooperation with Cristina Vertan has developed an annotation tool for the linguistic analysis of the machine-readable texts, which takes into account the vague decipherment status of the Maya script and the current state of research on Classic Maya language (Gronemeyer, 2014; Law and Stuart, 2017). In order to generate a readable text from the text corpus encoded in TEI XML, the tool, called ALMAH "Annotator for the Linguistic Analysis of Maya Hieroglyphs", queries the linguistic transliteration values stored in the digital sign catalog and, on this basis, semi-automatically generates a phonemic transliteration of the texts, which are further processed manually. Based on this workflow, the corpus-based Dictionary of Classic Maya is generated, which digitally maps the dictionary of Classic Mayan and its use in writing and forms the prerequisites for a deeper understanding of Maya culture, history, religion and society.

The digital-based epigraphic analysis of an inscription according to digital methods begins with the topographical description of the hieroglyphic writing (Iglesia et al., 2021). Thereby the individual graphs of the inscribed monuments are classified numerically. Based on these annotations, the linguistic analysis consisting of transliteration, transcription, morphological segmentation, linguistic interpretation and translation is performed using the annotation tool ALMAH (see chapter 4), and the results are finally published in the text database.

## 3. Encoding of Maya Hieroglyphic Texts

To document the arrangement of signs in the hieroglyphic block, the project applies Thompson's annotation convention to the XML/TEI scheme (Iglesia et al., 2021), according to which adjacent signs are separated by a period (.), superposed ones by a colon (:). Block segments within the hieroglyphic block are enclosed with square brackets [ ]. If a sign is inserted into another sign, it is marked with a degree sign (°) and the merging of two signs is indicated with a plus sign (+) (Prager and Gronemeyer, 2018). Definitions and editorial conventions, such as annotation of text structure, reading direction, topographic text arrangement, unreadable or reconstructed text passages, and text carrier design (shape, relief depth, framing, coloring, etc.) are predefined in the TEI schema and specified in the editorial guidelines. In the TEI annotation, the signs are referenced to the sign catalog using a TextGrid URI. For this purpose, the TextGrid URI to a graph must be retrieved in order to specify it in the TEI document. This is done using a TEI parser developed by Maximilian Behnert-Brodhun, which searches for the references from a numeric transcription code and generates the corresponding TEI structure automatically. Subsequently, the TEI document is parsed from an XML file in which only information about the text-carrying surfaces is given and the text fields and the individual hieroglyphic blocks are defined with the help of alphanumeric IDs. For each

block, the numerical transliteration of the graph entered in the digital sign catalog is entered using the conventions defined in the projects editorial guidelines, e.g. a hieroglyphic block transcribed using sign codes based on the catalog of Maya hieroglyphs : 1br.[501st:25st]. With the increasing number of encoded inscriptions, the sign catalog, which currently counts more than 1000 signs, is also being completed. With the help of the TEI parser, the XML files with the previously created numerical transliterations of the hieroglyphic texts are transformed into TEI documents and saved in TextGrid. At the same time, the TEI file is displayed online and can be viewed on screen and checked for errors.

A special feature for the quality control of our epigraphic work is the display of the original spelling. For this purpose, the parser, as well as the annotator ALMAH, retrieve the image of the graph from the digital sign catalog using the numeric character codes and displays it next to the numeric transliteration in the parser's result window. This visual validation allows the numeric transliteration to be checked and, if necessary, corrected before processing the TEI document. If the transliterations are correct, the generated TEI document can be checked and validated in TextGrid. For the annotation of unreadable and reconstructed text passages, for example, the project uses a specific TEI-P5 application profile and follows the EpiDoc Guidelines[7] to document classic or ancient texts in TEI XML. Damages, reconstructions, explanations for reconstructed text passages as well as the layout of the text carrier are not created into the XML by the parser, but have to be edited manually in the document according to our editorial guidelines. In the further course of the project, the parser will be extended to include these editorial functions so that these areas can also be created automatically in the future.

## 4. Annotator for the Linguistic Analysis of Maya Hieroglyphs (ALMAH)

1. Alphanumeric transliteration according to graphic variants

[512st:25st].181br

2. Numeric transliteration according to sign number

[512:25].181

3. Graphematic transliteration (broad transliteration)

[chu:ka].ja

4. Graphemic transliteration

chu-ka-ja

5. Phonemic transliteration indicating morphemes

chu-ka=ja

6. Morphological transcription according to morphemic units

chu[h]k-aj-Ø

7. Morphophonemic transcription (free and bound morphemes)

chuhk-aj-Ø

8. Morphosyntactic glossing (ling. description)

chuhk-    aj-    Ø
capture-PASS-V.INTR.MOD-3s.ABS

9. Consolidated transcription

chuhkaj

10. Literal translation

was captured

11. Free translation of the inscription

"… on the day 7 Imix 14 Tzec he was captured …"

Figure 2: Eleven annotation levels of Maya texts used in ALMAH (concept and terminologies by S. Gronemeyer, layout by Prager)

Linguistic transliterations and transcriptions of the inscriptions are generated automatically with the help of the analysis or annotation tool ALMAH in the next step. The linguistic backbone model is developed and extensively described in (Gronemeyer 2014). It processes a total of eleven epigraphic annotation levels (Figure 2), which are dynamically generated from the annotation of the previous level. The analysis and annotation typically proceeds as follows: The annotation tool is accessing the data in TextGrid or locally via an OAI-PMH interface. Once a file is selected, the TEI document is loaded and the automatic analysis process begins. Analysis levels 1 - 4 are first generated automatically: 1) and 2) Numeric transliteration 1 and 2 with graph and character numbers. 3) and 4) Graphemic transliteration 1 and 2 with possible manual rearrangement of the reading order of the signs. Here the results of automatic transliteration are displayed block by block. In addition to the numerical transliteration, the images of the individual graphemes are imported from TextGrid into ALMAH and displayed with analysis level 1-4. From the third annotation level on, manual corrections, additions and multiple analytical variants are possible, so that we can, for example, operate simultaneously with several decipherment suggestions. For example, if several linguistic readings are available for a sign, the analysis in graphemic transliteration allows selection of a particular reading or readings stored in the digital sign catalog via a selection window. If two or more readings are selected, ALMAH generates a corresponding number of graphemic transliteration variants that can be analyzed in parallel by the editors. However, if no reading is entered in the sign catalog, ALMAH takes the sign number and inserts it into the transliteration. On the level of graphemic transliteration 2, the reading order of the signs can also be rearranged as well as the morpheme boundaries can be changed with the help of a graphical interface. The conversion of the reading order becomes necessary when it does not correspond to the original writing order. From the graphemic transliteration of level 4, the phonemic transliteration of level 5 is created in the following step. Here, the morpheme boundaries between the phonemes are defined with the help of a graphical interface in order to distinguish free and bound morphemes. At level 6, the morphologically segmented transcription, the lexical and grammatical morphemes, such as inflections, derivations, proclitics or enclitics are segmented, reconstructed or superfluous sounds or sound loss are marked. For this purpose, transcriptions are dissected into phonetic chains, whereby superfluous sounds are removed, needed ones are inserted, morpheme boundaries are set, or null morphemes are used. At level 7, the morphophonemically consolidated transcription is created. At level 8, the consolidated morphosyntactic glossing is done. In this process, the brackets and special characters inserted at level 7 are removed and only the cleaned transcription is displayed, on which the interlinear morpheme glossing of the lexical and grammatical morphemes is performed. Interlinear morpheme glosses indicate the meanings and grammatical properties of individual words and parts of words. The morpheme glossing used in ALMAH is based on the Leipzig glossing rules, which have been extended and adapted by Frauke Sachse and Michael Dürr (2016) for the analysis of Mayan languages. The glosses are assigned in the tool to the lexical classes nouns, verbs, adjectives, adverbs, particles, pronouns, articles, classifiers, conjunctions, demonstratives, numerals, and prepositions, and are searchable and selectable via a matrix of language examples. If a definite assignment is not possible, several

---

glosses can be assigned to one morpheme. Based on these analysis steps, the consolidated transcription of the inscription (without special characters and brackets) is automatically generated on level 9. On annotation level 10, the editors can create the literal translation of the inscription, and finally, on level 11, the free translation. Free annotations of the hieroglyphic blocks also allow scholars to annotate calendrical information, nominal phrases, place names, or events and to ontologically link them to datasets from TextGrid in order to interpret the text and vocabulary of Classic Maya embedded in their historical and sociocultural context. In this way, over one hundred and fifty years of epigraphic research history and findings can be linked to our current analyses in an ontology.
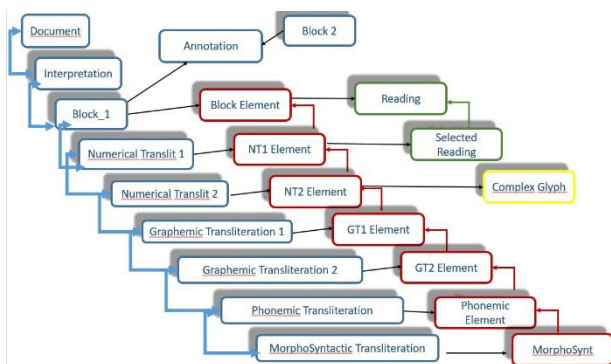
## 5. Architecture and Functionalities
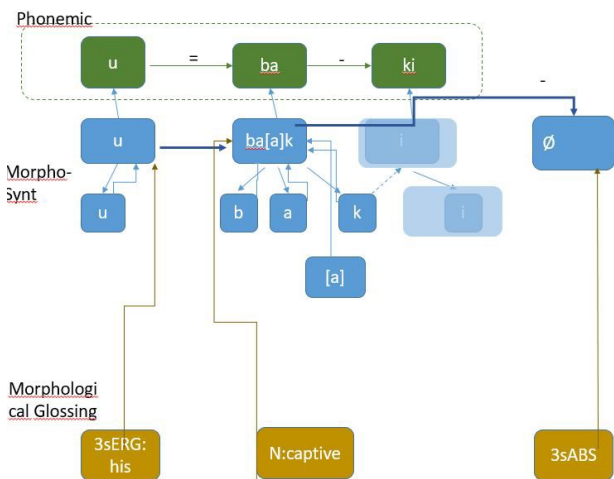


Figure 3: ALMAH Data Structure



Figure 4: Data Model Example (interconnection of Elements among several transliteration levels)

This complex linguistic model is mapped on a graph-based data model. Each transcription level represents a node in a tree structure. Each node contains information about the current transliteration level (Id, label) and a nested graph representing the structure of the transliteration. A transliteration is represented by a succession of elements (nodes of the structure graph) and operators (labelled edges in the structure graph). Elements of each transliteration know their ancestors. In this way we have the possibility at every moment to reconstructs the analysis path. The data Structure is presented in Figure 3 and an example in figure 4.

The structure gives also the possibility to operate dynamically changes on the graph label. Each transliteration level can generate several variants at the next level (working hypothesis). The first for levels are automatized: readings of the elements are extracted from the RDF-Database. If an element has several readings, the user is asked to select the possible ones for the current block. If more than one alternative reading is selected, the tool generates all possible combinations. At the linguistic level we give the possibility of linking the semantic annotation with English Wordnet-Sysets (only when the meaning of the word truly corresponds with a wordnet sysnset). In Figure 5 we present an example of processing done with the ALMAH Tool:
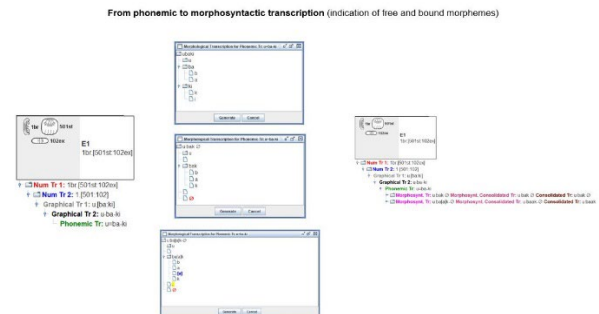


Figure 5: Interface of the ALMAH Tool

At this moment the linguistic information at phonetic and morphological level is done manually. Given the fact that the deciphering process is not completed it is quite common, that the user works at each level with more than one working hypothesis. A rule-based linguistic annotation approach, as known from the state-of-the art computational linguistics is in this case not possible. A supervised machine learning approach is in absence of a large annotated training corpus (given the number of featured to be learned) not realistic at this stage. However we are planning to use the manually annotations for building such a corpus, and introduce in a further version of the system a translation-memory –like approach. At each step, the system will search in the database for existent solution and will present the user possible annotation hypothesis, from which one or more will be manually selected. A fully unsupervised machine-learning algorithm is at this moment not appropriate, as long as the grammar of the language is not completely researched. In a third step, we envisage the possibility of exporting ALMAH –output in an ANNIS[8]-compatible format, which will allow corpus-linguistics specific queries.

## 6. Conclusions and further work

The newly developed tool ALMAH supports the epigraphic annotation and linguistic analysis of Maya hieroglyphic texts by standardising the decipherment process through semi-automatic processes and improving the epigraphic

---

workflow through machine learning. The tool provides the necessary flexibility to operate with alternative readings where a unique identification of characters in a block is not possible or multiple reading variations exist for a character or hieroglyph. ALMAH combines the linguistic annotation of hieroglyphs, including morphoglossification, with the creation of lemmas, which form the basis for the dictionary of Classic Mayan.

The tool is written in Java 8 as client application. An Internet connection is for the data reading and save necessary. Although it relies on a complex data-structure the interface is user friendly and transparent. The graph data structure is represented as such (through usage of graph libraries) and users can change edges, order of the graph nodes , i.e. realise permutation od elements, rename edges of the graph). Data is stored in an instance of OrientDB[9], which is the only database allowing graph and document data structures. Further work concerns the (semi) automatisation of the annotation steps (through a learning mechanism) as well as the generation of entries for a lexicon of Classic Mayan, the language of the hieroglyphs.

## 7. Bibliographical References

Diederichs K., Prager C.M., Brodhun M., and Tamignaux C. (2020), *„Ich brauch' mal ein Foto …": der Umgang mit Bildern im Projekt Textdatenbank und Wörterbuch des Klassischen Maya* [in:] "Bilddaten in den Digitalen Geisteswissenschaften.," C. Hastik, P. Hegel (eds.), Wiesbaden: Harrassowitz, pp. 175–197.

Diehr F., Brodhun M., Gronemeyer S., Diederichs K., Prager C.M., Wagner E., and Grube N. (2018), *Ein digitaler Zeichenkatalog als Organisationssystem für die noch nicht entzifferte Schrift der Klassischen Maya*, [in:] "Knowledge Organization for Digital Humanities: Proceedings of the 15th Conference on Knowledge Organization WissOrg'17 of the German Chapter of the International Society for Knowledge Organization (ISKO) [30th November - 1st December 2017, Freie Universität Berlin]," C. Wartena, M. Franke-Maier, E. de Luca (eds.), pp. 37–43, Berlin: Freie Universität Berlin.

Diehr F., Gronemeyer S., Prager C.M., Diederichs K., Grube N., and Sikora U. (2019), *Modelling Vagueness – A Criteria-based System for the Qualitative Assessment of Reading Proposals for the Deciphering of Classic Mayan Hieroglyphs* [in:] "Proceedings of the Workshop on Computational Methods in the Humanities 2018," Lausanne: Université de Lausanne, pp. 33–44.

Gronemeyer S. (2014), *The Orthographic Conventions of Maya Hieroglyphic Writing: Being a Contribution to the Phonemic Reconstruction of Classic Mayan*, Ph.D. Dissertation, Department of Archaeology, La Trobe University, Melbourne. http://hdl.handle.net/1959.9/321048

Grube N. (1994), *Mittelamerikanische Schriften* [in:] "Schrift und Schriftlichkeit: ein interdisziplinäres Handbuch internationaler Forschung = Writing and its Use: an Interdisciplinary Handbook of International Research," H. Günther, O. Ludwig (eds.), Berlin: Walter de Gruyter, Vol. 1, pp. 405–415.

Grube N., Prager C.M., Diederichs K., Gronemeyer S., Wagner E., Brodhun M., Diehr F., Maier P. (2014), *Jahresabschlussbericht 2014* [Electronic Document]. http://mayawoerterbuch.de/?p=4477

Houston S.D. and Martin S. (2016), *Through Seeing Stones: Maya Epigraphy as a Mature Discipline*, Antiquity 90(350):443–455.

Iglesia M. de la, Diehr F., Sikora U., Gronemeyer S., Behnert-Brodhun M., Prager C.M., and Grube N. (2021), *The Code of Maya Kings and Queens: Encoding and Markup of Maya Hieroglyphic Writing*, Journal of the Text Encoding Initiative Issue 14. Retrieved from https://journals.openedition.org/jtei/3336

Law D., and Stuart D. (2017), *Classic Mayan: An Overview of Language in Ancient Hieroglyphic Script* [in:] "The Mayan Languages," J. Aissen, N.C. England, R. Zavala (eds.), London; New York: Routledge / Taylor & Francis Group, pp. 128–172.

Martin S. (2020), *Ancient Maya Politics: A Political Anthropology of the Classic Period 150–900 CE*, Cambridge: Cambridge University Press.

Prager C.M. (2015), *Das Textdatenbank- und Wörterbuchprojekt des Klassischen Maya: Möglichkeiten und Herausforderungen digitaler Epigraphik* [in:] "TextGrid: Von der Community - für die Community: Eine Virtuelle Forschungsumgebung für die Geisteswissenschaften," H. Neuroth, A. Rapp, S. Söring (eds.), Glückstadt: Werner Hülsbusch, pp. 105–124.

Prager C.M., and Gronemeyer S. (2018), *Neue Ergebnisse in der Erforschung der Graphemik und Graphetik des Klassischen Maya* [in:] "Ägyptologische 'Binsen'-Weisheiten III: Formen und Funktionen von Zeichenliste und Paläographie," S.A. Gülden, K.V.J. van der Moezel, U. Verhoeven-van Elsbergen (eds.), Stuttgart: Franz Steiner Verlag, pp. 135–181.

Prager C.M., Grube N., Brodhun M., Diederichs K., Diehr F., Gronemeyer S., and Wagner E. (2018), *The Digital Exploration of Maya Hieroglyphic Writing and Language* [in:] "Crossing Experiences in Digital Epigraphy: From Practice to Discipline," A. De Santis, I. Rossi (eds.), Berlin: De Gruyter, pp. 65–83.

Sachse F., and Dürr M. (2016), *Morphological Glossing of Mayan Languages under XML: Preliminary Results* [Electronic Document]. http://mayawoerterbuch.de/?p=2122

Stuart D. (1998), *Dynastic History and Politics of the Classic Maya* [in:] "Maya Civilization," P. Schmidt, M. de la Garza, E. Nalda (eds.), London: Thames and Hudson, pp. 320–335.

Thompson J.E.S. (1962), *A Catalog of Maya Hieroglyphs*, Norman, OK: University of Oklahoma Press.

Zender M. (1999), *Diacritical Marks and Underspelling in the Classic Maya Script: Implications for Decipherment*, M.A. Thesis, Department of Archaeology, University of Calgary, Calgary.

[9]https://orientdb.org/