

# Leveraging AI for Historical Linguistics

**Ms. Aditi Patil**

Student, B.Tech Computer Engineering, Ramrao Adik Institute of Technology, D Y Patil Deemed to be University, Nerul, Navi Mumbai, Maharashtra.

## Abstract

Decoding ancient languages plays a crucial role in comprehending historical occurrences, safeguarding cultural legacy, and advancing linguistic expertise. Nevertheless, this undertaking is rife with obstacles stemming from inadequate data accessibility, intricate linguistic structures, and the gradual loss of linguistic knowledge. The utilization of artificial intelligence (AI) methodologies presents promising avenues for surmounting these challenges and expediting the decipherment procedure. This manuscript delves into the application of AI techniques, such as machine learning, natural language processing, and pattern recognition, to interpret ancient scripts and unveil linguistic phenomena. Through an examination of the strategies utilized by linguists and adept individuals in deciphering ancient languages, the manuscript underscores the potential of AI-driven methodologies in this realm. Moreover, it deliberates on pragmatic challenges and puts forth recommendations for the integration of AI technologies into decipherment endeavors. By fostering interdisciplinary cooperation among AI experts, linguists, archaeologists, and historians, this manuscript seeks to pave the way for the creation of AI systems tailored for unraveling ancient languages and augmenting our comprehension of human history and cultural legacy.

**Keywords:** Artificial Intelligence, Machine Learning, Historical Linguistics

## 1. Introduction

Understanding the past, safeguarding cultural heritage, expanding linguistic knowledge, and interpreting archaeological discoveries are all contingent upon our capacity to decode ancient languages. Researchers can acquire deeper insights into past events, societies, and ancient civilizations through the interpretation of ancient texts. This process not only progresses linguistic theories and methodologies but also aids in the conservation of invaluable cultural heritage. Moreover, comprehending antiquated languages offers a glimpse into human communication, cognition, and the development of written language systems. It promotes intercultural communication and global comprehension, elucidates historical enigmas, and enriches archaeological analyses. In essence, the comprehension of our shared human history and legacy greatly relies on our proficiency in deciphering ancient languages. The examination of these texts equips society with a historical vantage point that influences present-day decision-making, policy formation, and societal discussions, fostering empathy, admiration, and acceptance of cultural diversity, thereby contributing to a more inclusive and cohesive society.

Numerous impediments exist in the realm of deciphering historical languages, such as limited resources, a shortage of bilingual texts for comparative analysis, and the gradual erosion of linguistic expertise over time. An illustration of this challenge is evident in the case of Linear A from the Minoan civilization, which remains unintelligible to researchers due to the absence of a bilingual inscription or a

comprehensible cipher. Similarly, the enigmatic Indus Valley Script, associated with one of the most ancient urban societies on the planet, continues to baffle scholars without a counterpart akin to the Rosetta Stone for decipherment. Likewise, the Etruscan language utilized by the ancient Italian civilization poses difficulties owing to the scant availability of surviving texts and the absence of a direct linguistic successor. These instances underscore the complexities inherent in the comprehension of historical languages, often necessitating interdisciplinary collaboration, innovative methodologies, and persistent endeavors to unravel their enigmas.

## 2. Objective

The essential objective of the term paper is to demonstrate the application of manufactured insights (AI) strategies within the translation of old dialects and the disclosure of related phonetic wonders. The inquiry about endeavors to show how AI methods such as machine learning, normal dialect preparing, and design acknowledgment can be utilized to decode antiquated scripts by conducting a comprehensive examination of these innovations. AI calculations inside the domain of fake insights can distinguish designs, phonetic characteristics, and syntactic configurations inside broad collections of engravings, tablets, and original copies. This encourages the method of unraveling phonetic elucidations, semantic significances, and linguistic standards implanted within the antiquated literary sources.

The term paper points to address viable challenges in old dialect translating in expansion to encouraging scholarly information. Analysts can decrease human translation predispositions, overcome information fracture, and speed up the decipherment handle by utilizing AI advances. AI-driven techniques also encourage the integration of assorted information sources, counting hereditary, chronicled, and archeological records, permitting for a more comprehensive understanding of antiquated dialects and their social settings.

The overall goal of the term paper is to appear how fake insights (AI) can revolutionize our understanding of human history, etymological differences, and social legacy by helping to translate old dialects. The paper points to fortify participation between AI analysts, language specialists, archeologists, and students of history in arrange to open the riddles of the past and protect our shared social bequest for future eras. It does this by combining cutting-edge AI techniques with thorough academic examination.

## 3. Proposed Methodology

Studying ancient Textbooks has consistently been a protracted and arduous undertaking in the scholarly realm. Academics are tasked with meticulously examining archaic inscriptions such as eulogies or calligraphies with utmost precision. The deciphering process demands a comprehensive comprehension of the symbols and terminologies employed, which can prove to be intricate given the antiquity of these educational materials originating from civilizations that have long vanished. To ascertain the meanings of the terms and their contextual coherence, researchers undertake comparative analyses with languages they are familiar with, scrutinizing for parallels in linguistic structures. This meticulous process is time-consuming and necessitates collaboration among experts from diverse disciplines like history and linguistics. It can be analogized to unraveling a profound enigma, requiring a profound depth of knowledge and patience. While artificial intelligence can potentially utilize comparable techniques to traditional decryption methods, its capacity for unparalleled speed and efficacy is unprecedented in this domain.

### 3.1 Pre-existing translation tools

Utilizing artificial intelligence, as previously conducted by [1] Ronojoy Adhikari and his research team,

have developed a “deep-learning” algorithm that can read the Indus script from images of artefacts such as a seal or pottery that contain Indus writing.

Scanning the image, the algorithm smartly “recognises” the region of the image that contains the script, breaks it up into individual graphemes (the term in linguistics for the smallest unit of the script) and finally identifies these using data from a standard corpus. In linguistics the term corpus is used to describe a large collection of texts which, among other things, are used to carry out statistical analyses of languages. The software they created underwent a comprehensive analysis of a wide array of linguistic and non-linguistic data, including but not limited to computer programming code, DNA sequences, established languages, which in turn revealed the contrasting levels of entropy between natural languages and structured codes. Upon application to the Indus script, the program illustrated similarities to established languages such as Tamil and Sanskrit, thus indicating the existence of a linguistic framework within the script itself. This significant finding not only validated their approach but also showcased the potential of AI-driven methodologies in deciphering ancient languages.

### **3.2. Syntax and Semantic Analysis**

After the completion of this particular stage, Artificial Intelligence has the capability to be utilized in the comprehension of both the syntax and semantics encapsulated within the script. Diverse Natural Language Processing methodologies, for example, part-of-speech tagging, dependencies parsing, and syntactic parsing, can be implemented for the purpose of scrutinizing the syntax of sentences by disassembling them into distinct grammatical constituents and subsequently identifying the interrelations existing among the words. The utilization of AI algorithms becomes instrumental in deducing the connotations of words, phrases, and sentences through the application of semantic analysis. This process entails the utilization of methodologies such as word embeddings, which function by portraying words in the form of numerical vectors derived from their contextual utilization within a given corpus of text.

### **3.3 Leveraging Simulated Data Generation**

When no data is available for a particular ancient language, AI methodologies are capable of being employed by means of a technique known as simulated data generation, a process in which experts in linguistics collaborate closely with researchers in the field of artificial intelligence to produce simulated datasets that are constructed based on the existing knowledge pertaining to related languages, historical contexts, and recognized linguistic characteristics. This particular method entails the creation of synthetic texts, lexicons, and grammatical principles that are designed to closely resemble the presumed structure and attributes of the ancient language under investigation. The AI algorithms are subsequently subjected to training using these simulated datasets, employing methodologies such as generative adversarial networks (GANs) or reinforcement learning to discern and internalize the fundamental patterns and configurations of the language in question. Over time, with a series of iterative adjustments and validation procedures carried out by linguistic experts, the AI model incrementally enhances its capacity to scrutinize, construe, and produce text in the simulated ancient language. Despite the reliance of this approach on expert knowledge and presumptions, it presents a valuable initial phase for further investigation and exploration of the target language, particularly in situations where authentic data is scarce or unattainable.

### **3.4 Detection of concealed patterns and linguistic attributes**

Moreover, through the utilization of unsupervised learning techniques, artificial intelligence (AI) algorithms are assigned the task of detecting patterns, linguistic attributes, and stylistic norms directly from the textual content itself, eliminating the necessity for pre-existing datasets. Various methodologies like clustering algorithms, topic modeling, and word embedding can be utilized to reveal concealed structures

and connections within the text. Through the analysis of disparities in writing styles, regional linguistic variations, and errors in transcription, AI has the capability to assist scholars in gaining a deeper understanding of the intricacies of language and cultural contexts embedded in ancient manuscripts. This continuous process of scrutinizing and investigating enables researchers to gradually unveil the latent patterns and interpretations encoded within the textual material, ultimately contributing to enhancing translation precision and comprehension.

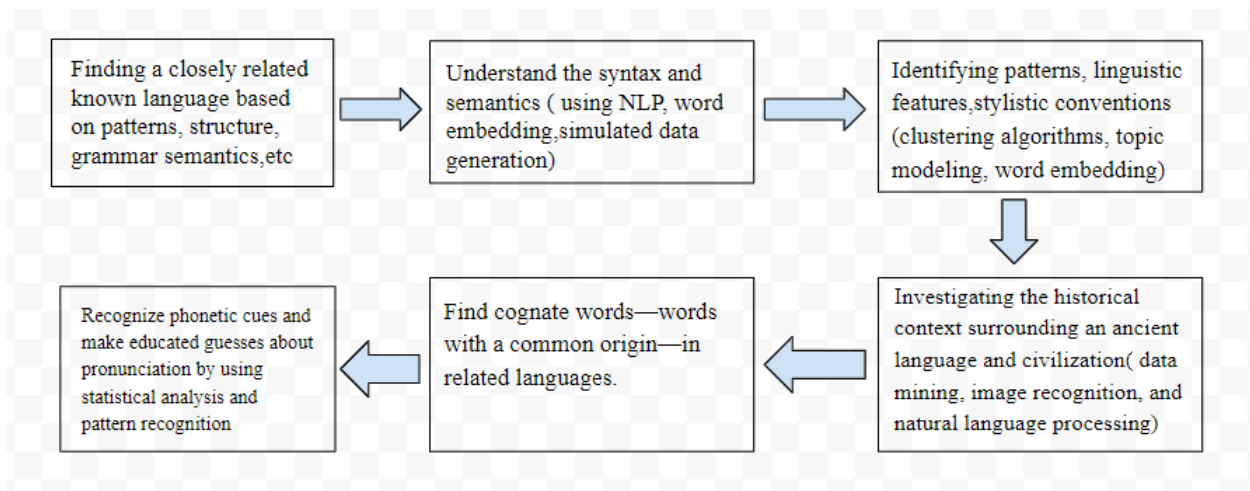
AI can then aid in investigating the historical context surrounding an ancient language and civilization by employing techniques such as data mining, image recognition, and natural language processing. Firstly, AI algorithms can analyze vast repositories of historical texts, scholarly articles, and archaeological reports to extract relevant information about the time period, geographical region, cultural practices, and socio-political events associated with the civilization in question. Additionally, AI-powered image recognition can analyze artifacts, inscriptions, and archaeological sites to identify patterns, symbols, and motifs, providing insights into the material culture and socio-economic organization of the civilization. By processing textual and visual data, AI enables researchers to uncover previously unrecognized connections and correlations, facilitating a deeper understanding of the historical context and cultural dynamics of the ancient civilization.

### **3.5 Finding Cognate Words in related languages**

Once done, AI can be used find cognate words—words with a common origin—in related languages. Cognate analysis aids in reconstructing linguistic changes over time and etymological relationships. Artificial Intelligence (AI) can identify words with similar meanings and shared linguistic features by analyzing textual data from different languages and utilizing computational algorithms like word embeddings and similarity measures. AI can also use morphological and phonetic analysis to find structural and acoustic similarities between words in different languages. Artificial intelligence (AI) algorithms can identify possible cognates and etymological relationships by clustering words based on their phonetic and semantic properties. Even in the lack of previous data, this procedure enables researchers to reconstruct linguistic changes over time and infer historical connections between languages. Further, AI can be used to search for bilingual inscriptions or texts containing both the ancient language and a known language. Bilingual texts provide valuable clues for deciphering and translating the ancient language by allowing AI to make direct comparisons and identify correspondences between the two languages.

### **3.6 Recognizing Phonetic cues and grammatical markers**

Statistical analysis and pattern recognition enable artificial intelligence to recognize phonetic cues and make educated guesses about pronunciation. Morphological analysis can be performed by AI algorithms that use unsupervised learning techniques to identify patterns of inflection, derivation, and compounding. These algorithms look at the internal structure of words to determine morphemes and grammatical markers. AI models trained on universal grammar principles can be used to achieve syntactic parsing, which deduces meaning from sentence structures, word order, and grammatical constructions. Artificial intelligence (AI) systems with natural language understanding capabilities can help with semantic interpretation. These systems can interpret words, phrases, and sentences in their respective cultures based on idiomatic usage, metaphorical expressions, and contextual cues.



**Fig 3.6: Process of using AI for language translation**

Overall, AI offers valuable tools for deciphering ancient scripts and understanding their linguistic and cultural significance, even in the absence of prior data. But regardless of the approach used, AI is dependent on high-quality data being available in a machine-readable format. This remains a key challenge when it comes to ancient texts, given that they often come to chipped, eroded, or incomplete in some other form. Scholars can spend decades debating the uniqueness of symbols: Is that a scratch next to a known character, for instance, or a new character altogether? Given how little there is to work with when it comes to long-lost languages, noisy or incomplete data can seriously curtail decipherment efforts.

#### 4. Conclusion

In summary, the current manuscript serves as a cornerstone in laying the groundwork for the development of sophisticated artificial intelligence systems specifically tailored for the purpose of deciphering ancient languages. The manuscript provides valuable and significant perspectives into the challenges and opportunities within this domain, elucidating the methodologies and approaches employed by linguists and experts in the interpretation of archaic texts. Integration of artificial intelligence (AI) methodologies with advancements in machine learning, natural language processing, and pattern recognition presents a promising avenue for expediting the decryption procedures and unveiling the linguistic and cultural heritage embedded within aged manuscripts. Moreover, the research underscores the critical importance of fostering interdisciplinary collaboration among linguists, historians, archaeologists, and AI scholars to enhance the efficacy of AI-facilitated decipherment methodologies. This manuscript serves as a catalyst propelling the continuous advancement of artificial intelligence technologies dedicated to unraveling the enigmatic languages of antiquity, thereby enriching our comprehension of historical narratives and cultural pluralism.

#### 5. Acknowledgement

I extend our heartfelt gratitude to all those who contributed to the completion of this research endeavor. I am indebted to Ramrao Adik Institute of Technology for providing the necessary resources, facilities, and funding that made this research possible. Your support has been indispensable and greatly appreciated. I am also thankful to the participants and volunteers who generously contributed their time, knowledge, and expertise to this study. Without your cooperation and participation, this research would not have been

feasible. Additionally, I acknowledge the contributions of colleagues, friends, and family members who provided encouragement, support, and understanding throughout the duration of this project.

Finally, I express our gratitude to the academic community, whose collective wisdom, insights, and scholarly contributions have informed and inspired our research endeavors.

## 6. References

1. Satish Palaniappan, Ronojoy Adhikari, “Deep Learning the Indus Script”, PLOS submission, Feb 2017, rXiv:1702.00523v1 [cs.CV]
2. Koskeniemi S, Parpola A, Parpola S. A method to classify characters of unknown ancient scripts. *Linguistics*. 1970;8(61):65–91.
3. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature*. 2015;521:436–444.
4. Dr. Shikha Chadha, Ms. Neha Gupta, Dr. Anil B, and Ms. Rosey Chauhan, “A Novel Framework for Ancient Text Translation Using Artificial Intelligence”, *ADCAIJ: Advances in Distributed Computing and Artificial Intelligence Journal Regular Issue*, Vol. 11 N. 4 (2022), 411-425 eISSN: 2255-2863.