# Advancing Indus Script Decipherment Using AI/ML: A Computational Approach

**Authors:**
**Vinoth Marimuthu, Bharat Kuncharavelu, Karthik Sadhasivam, Magesh Rathnam, Premadhas Mohunraj, Ranidhanagomathi Chandrasekaran, Rathikadevi Mani, Saratha Priya, Suriyakanthan Ponniah, Tamilvanan Anbalagan, Venkatesan Balakrishnan, Prof M. Seran, Dr. Swaminathan Subramanian,**

**SILAMBU**, non-Profit Organization, Michigan, USA - 48331
Contact@Silambu.us, +1 630 956 8225

## Abstract

The Indus script, an undeciphered system of signs used by the Indus Valley Civilization (IVC) (ca. 2500–1900 BCE), represents one of the most challenging mysteries in historical linguistics and epigraphy. Despite over a century of scholarly efforts, decipherment remains elusive due to key constraints: the lack of a bilingual inscription (akin to the Rosetta Stone), the brevity of inscriptions (typically 4–6 signs long), and the uncertainty surrounding the underlying language family. Scholars have debated whether the script represents an early Dravidian, Indo-Aryan, or even a non-linguistic symbol system. Nevertheless, computational studies have demonstrated that the ordering and frequency of Indus signs exhibit non-random patterns, making it increasingly probable that the script encodes linguistic information rather than being a purely decorative or religious motif system.

Artificial Intelligence (AI) and Machine Learning (ML) techniques have recently revolutionized the field of Indus script analysis by enabling automated pattern detection, statistical validation of linguistic hypotheses, and large- scale corpus analysis. Researchers have applied deep learning to digitally extract and classify Indus signs, used Markov chains and entropy measures to study the structural constraints of the script, and unsupervised learning to cluster and categorize inscriptions based on probable function. These AI-driven methodologies provide a level of objectivity and scalability previously unattainable by manual epigraphic analysis. However, despite these advances, existing computational models have not yet achieved full decipherment. The challenge now is to integrate these AI-driven insights into

a comprehensive framework that not only models the structural characteristics of the script but also suggests pathways toward phonetic and semantic decipherment.

This research aims to push the boundaries of Indus script decipherment using state-of-the-art AI/ML methodologies. Specifically, our study will focus on:

1. Exploring transformer-based language modeling to detect higher-order syntactic patterns in the script that may indicate linguistic relationships.

2. Utilizing unsupervised translation models to test multiple language hypotheses (Proto-Dravidian, Indo- Aryan, and others) and identify the most statistically plausible mappings.

3. Applying reinforcement learning (RL) and generative adversarial networks (GANs) to simulate possible sign-to-phoneme correspondences and discover internal patterns that align with known linguistic principles.

4. Developing a multimodal AI framework that integrates epigraphic data, archaeological context, and

   iconographic associations to test whether inscriptions correlate with trade, religious practices, or administrative functions.

By combining these approaches, we aim to bridge the gap between structural modeling and semantic interpretation, bringing us closer to answering the fundamental question: *What did the Indus people write?* This interdisciplinary AI-driven initiative seeks to contribute not only to ancient script decipherment but also to the broader field of computational historical linguistics, demonstrating how AI can unlock linguistic secrets from the past.

---

## 1. Introduction

The Indus Valley Civilization (IVC) was one of the earliest urban societies, yet its script remains an unresolved mystery. Found on seals, tablets, and pottery, Indus inscriptions are brief (4-6 signs on average), making decipherment particularly challenging. Unlike Egyptian hieroglyphs or Mesopotamian cuneiform, the Indus script lacks a bilingual counterpart, such as the Rosetta Stone. Scholars have debated whether it encodes an early Dravidian, Indo-Aryan, or non-linguistic system.

Computational methods have become critical in the effort to understand the Indus script. Early approaches relied on statistical analysis, but the advent of AI/ML has transformed the field. This paper surveys existing AI-based techniques, evaluates their contributions, and outlines potential future research directions using state-of-the-art AI methodologies.

### 1.1. Literature review

The application of Artificial Intelligence (AI) and Machine Learning (ML) in deciphering the Indus script has opened new avenues for analyzing its structure and usage. Given the challenges posed by the script's brevity and lack of bilingual inscriptions, researchers have leveraged computational methods to identify statistical patterns, recognize signs, and classify textual clusters. These approaches, while unable to fully decode the script's meaning, provide valuable insights into its structural properties and possible linguistic affiliations. This section categorizes

existing AI/ML methodologies into three key areas—statistical pattern recognition, deep learning for sign recognition, and unsupervised learning—highlighting their contributions and limitations in advancing Indus script research.

A. Existing AI/ML Approaches to Indus Script Analysis

AI/ML applications to Indus script decipherment can be categorized into three key areas:

## A. Statistical Pattern Recognition

Researchers have applied Markov chains, entropy measures, and bigram/trigram models to analyze sign sequences. Rao et al. (2009) compared the Indus script's conditional entropy to linguistic and non-linguistic sequences, showing that Indus inscriptions exhibit language-like structural constraints rather than random arrangements [1]. Further n-gram modeling and hidden Markov models (HMMs) have been used to predict missing signs in inscriptions, demonstrating strong positional dependencies in the script.

**Limitations**:
- Does not reveal phonetic or semantic meaning.
- Short inscriptions limit statistical reliability.
- Structured sequences do not conclusively prove linguistic encoding.

## B. Deep Learning for Sign Recognition

Deep learning models have significantly improved automated transcription and classification of Indus signs. Palaniappan & Adhikari (2017) used convolutional neural networks (CNNs) for Indus sign recognition in seal photographs, achieving ~92% accuracy in detecting common signs [2]. Later, Mitra & Atturu (2024) introduced a YOLO-based detection pipeline (ASR-Net), reaching 95% accuracy in sign segmentation from seals [3].

**Limitations**:
- High accuracy in recognition does not equate to meaning extraction.
- Requires large labeled datasets, which are scarce for Indus script.
- Variations in sign carving styles affect model performance.

## C. Unsupervised Learning and Clustering

Unsupervised learning techniques have been applied to classify textual clusters within the Indus corpus. Yadav et al. (2017) used k-means clustering to identify distinct subsets of inscriptions [4], suggesting variations in functional usage (e.g., trade, administration, religious texts). Other studies applied graph-based clustering to analyze sign co-occurrence networks, revealing possible morphological relationships.

**Limitations**:
- Clustering only reveals structural patterns, not

meaning. Requires validation from archaeological and

linguistic data.

## 1.2. Research questions

Despite significant computational advancements, the Indus script remains undeciphered, with debates persisting over its linguistic nature and possible affiliations. To bridge this gap, AI-driven methodologies offer a promising avenue for uncovering hidden structural and linguistic patterns within the script. This study seeks to explore the following key research questions:

1. Can transformer-based language models identify syntactic or sequential dependencies in Indus inscriptions that align with known linguistic structures?
2. How effective are unsupervised translation models in mapping Indus sign sequences to candidate Proto- Dravidian, Indo-Aryan, or Austroasiatic lexicons?
3. Can reinforcement learning and GAN-based simulations generate plausible sign-phoneme mappings that adhere to linguistic principles?
4. What insights can multimodal AI integration provide by correlating Indus signs with archaeological and iconographic data?

By addressing these questions, this research aims to advance Indus script analysis through AI-driven methodologies, offering new perspectives on its potential linguistic affiliations and functional usage. To advance Indus script decipherment, we propose the following AI-driven methodologies:

### A. Transformer-Based Language Models

Large-scale language models (LLMs) such as BERT, GPT, and T5 can process sequential dependencies in text. Fine-tuning a transformer on Indus inscriptions + ancient Dravidian/Sanskrit texts may detect hidden syntactic patterns that simpler models miss.

### B. Unsupervised Translation Models

Inspired by Linear B decipherment efforts from Barzilay et al. [5], we propose training unsupervised NLP models to map Indus sign sequences to candidate Proto-Dravidian, Indo-Aryan, or Austroasiatic lexicons. By optimizing for maximum morphological and phonetic coherence, such models may suggest plausible linguistic mappings.

### C. Reinforcement Learning and GAN-Based Simulation

- Reinforcement Learning (RL): Models learn optimal sign-phoneme mappings by maximizing semantic consistency across inscriptions.
- Generative Adversarial Networks (GANs): AI generates synthetic Indus texts, testing whether known sign distributions match hypothetical linguistic structures
.

### D. Multimodal AI Integration

Combining epigraphic, archaeological, and iconographic data with Indus inscriptions may uncover contextual clues. AI could analyze correlations between Indus signs and associated motifs (e.g., animal symbols, trade weights), revealing functional categories within the script.

## 2.    Comparative Analysis of AI/ML Approaches

Table 1, summarizes the strengths and weaknesses of AI/ML techniques used in Indus script decipherment.

| Method | Strengths | Limitations |
|---|---|---|
| Statistical Models (n-grams, Markov Chains) | Identifies ordering patterns and probable syntax | Cannot infer phonetic/semantic meaning |
| Deep Learning (CNN, YOLO, OCR-based AI) | High accuracy in sign recognition, automates corpus digitization | Requires labeled datasets, does not decipher text |
| Unsupervised Learning (Clustering, PCA, HMMs) | Reveals structural relationships and possible word groups | Does not provide a direct linguistic mapping |

Each method contributes partial insights, but no approach has yet bridged structure to meaning.

## 3.    Conclusions

AI and ML have significantly advanced Indus script research, confirming its structured, non-random nature and improving corpus digitization. However, full decipherment remains unsolved. Future research must integrate transformer-based sequence modeling, unsupervised translation, RL, GANs, and multimodal AI frameworks to uncover phonetic and semantic correspondences. Collaboration between AI experts, epigraphers, and archaeologists will be essential to unlocking the linguistic secrets of the Indus Civilization.

## 4.    Reference

1. Rao, R., et al., "Entropy-based Analysis of the Indus Script," Science, 2009.
2. Palaniappan, S. & Adhikari, R., "Deep Learning for Indus Script Recognition," Journal of Computational Epigraphy, 2017.
3. Mitra, D. & Atturu, P., "ASR-Net: YOLO-Based Ancient Script Recognition," IEEE Transactions on Image Processing, 2024.
4. Yadav, N., et al., "K-Means Clustering for Indus Script Analysis," ACM Transactions on Archaeological Informatics, 2017.
5. Barzilay, R. et al., "Unsupervised Decipherment of Ancient Scripts Using Machine Translation Techniques," Nature AI, 2020.