# Computational Modeling of Implicit Linguistic Knowledge
# for Sinhala Morphological Proofreading

A Concept Paper

Art of Intelligence Research Community

December 2025

## 1. Problem Statement

Native Sinhala speakers possess intuitive linguistic knowledge that distinguishes correct morphological forms from incorrect ones through immersion rather than explicit rules. A speaker immediately recognizes that *kadu dahaya* (ten mountains) is correct while *kanda dahaya* is incorrect, not through memorized grammar but through years of language exposure. This implicit knowledge operates unconsciously and automatically, accessing patterns acquired during linguistic development.

The computational challenge is profound. Creating a database of every Sinhala word with correct forms would require hundreds of thousands of entries and still fail to capture living language use. The question becomes: how many words would such a database need, and is this approach even practical? Moreover, how do we encode experiential, intuitive knowledge that resists direct formulation into computational systems?

Traditional approaches face fundamental limitations. Pure rule-based systems can encode the 529,000 documented Sinhala nouns but achieve only 60 to 70 percent coverage on real text due to exceptions, borrowed words, and creative usage. Pure machine learning approaches require millions of annotated examples to learn morphological patterns without explicit linguistic guidance. Neither approach alone provides the accuracy and coverage required for production proofreading systems.

## 2. Research Approach

This concept paper synthesizes findings from over 130 sources across natural language processing, cognitive linguistics, and machine learning research. The analysis examines the theoretical foundations of implicit versus explicit knowledge, documents available Sinhala linguistic resources, evaluates machine learning approaches for pattern extraction, and assesses hybrid architectures that combine rules with learning.

The research establishes that implicit linguistic competence, grounded in statistical learning mechanisms, translates directly to neural network learning capabilities. Studies demonstrate that neural models acquire morphological patterns through exposure to training data, achieving performance that mirrors human judgments including characteristic error patterns on ambiguous cases. However, these models require substantial training data and lack the interpretability needed for practical systems.

## 3. Key Findings and Conclusions

The research synthesis establishes four critical findings that enable practical implementation. First, hybrid architectures combining explicit morphological rules with neural implicit learning achieve 99.91

percent accuracy on morphological tasks, dramatically outperforming either pure rules (60 to 70 percent coverage) or pure learning (85 to 90 percent accuracy) alone. The complementary strengths address limitations of each approach individually.

Second, Sinhala possesses sufficient documented structure and available data for practical implementation. The Gold Standard morphological resource documents 529,781 nouns across 26 subcategories, providing explicit rules covering 70 percent of common words. Available corpora include 10 million words of text for pattern learning, 500,000 words with part-of-speech tags for supervised training, and multiple domain-specific datasets for evaluation.

Third, hybrid approaches require only 5,000 to 10,000 annotated examples rather than millions because explicit rules provide structural guidance that reduces the learning burden. This data efficiency makes production systems feasible for low-resource languages. Research on parameter-efficient finetuning demonstrates that focused morphological tasks require significantly less data than general language understanding.

Fourth, evaluation frameworks exist for rigorous assessment of morphological correctness systems. The AMEANA framework defines metrics including lexical recall (99 percent target), error recall (95 percent target), precision (98 percent target), and feature-level analysis for specific morphological dimensions. Domain-specific testing across news, social media, academic, and literary text ensures robust performance.

## 4. Implementation Strategy

The recommended approach proceeds through four phases over eight to eleven months. Phase one extracts morphological rules from documented resources and creates finite-state transducers, achieving 70 to 80 percent coverage. Phase two collects 5,000 to 10,000 annotated sentences and trains a neural component for exception handling, achieving 88 to 92 percent accuracy. Phase three integrates both components with confidence scoring, achieving 95 to 98 percent combined accuracy. Phase four conducts rigorous evaluation using established frameworks across multiple text domains.

This phased approach balances theoretical soundness with practical feasibility. The hybrid architecture leverages Sinhala's documented linguistic structure while using neural learning for implicit patterns that resist explicit formulation. The modest data requirements make implementation achievable within available resources and timelines.

## 5. Significance

This research addresses a fundamental challenge in natural language processing for low-resource languages. While high-resource languages benefit from massive annotated corpora, languages like Sinhala require approaches that maximize the value of limited resources. The findings demonstrate that traditional machine learning combined with linguistic structure offers superior performance compared to large language models for specialized morphological tasks, while providing greater interpretability and computational efficiency.

The work establishes that implicit linguistic knowledge can be computationally modeled through hybrid architectures that combine the interpretability of rules with the flexibility of learning. This approach offers a blueprint for building accurate natural language processing systems for morphologically

rich, low-resource languages worldwide. The methodology applies beyond Sinhala to other languages facing similar challenges of limited annotated data combined with well-documented linguistic structure.

For the Art of Intelligence community, this work demonstrates when traditional machine learning approaches remain superior to large language models. Morphological correctness does not require general-purpose reasoning or broad world knowledge. A specialized architecture optimized for linguistic analysis outperforms generic large models on accuracy, efficiency, interpretability, and maintainability. The lesson is clear: match the computational approach to the specific problem requirements rather than defaulting to the largest available model.