

Low-Resource NLP

Indian Language LLMs

Chandan Kumar
Rohit Tiwari
Under Guidance of Dr. B.N. Subudhi

May 07, 2025

Agenda

- Research Gap & Objectives
- Corpus
- Tokenizer
- Model Architecture
- Evaluation Metrics
- Demo Interface
- Challenges & Solutions
- Applications & Limitations & Future Work
- Conclusion
- References

Research Gap & Motivation

- **Problem:** Existing tokenizers fail for Indian languages.

Example of BPE failure

"राष्ट्रपति" → ["रा", "ष्ट्र", "प", "ति"]

Morpheme-aware |

"राष्ट्रपति" → ["राष्ट्र", "पति"]

- **Motivation:**

- 22 official Indian languages with complex morphology.
- Need for linguistically-aware tokenization.

Objectives

Primary Goals

- Develop morpheme-based tokenizer for Hindi
- Train GPT-2 from scratch with custom tokens
- Compare against BPE baseline

- Wikipedia Hindi Dumps (1GB raw)

- Total Number of Words : 77,60,500
- Total Characters: 4,01,26,198
- Total Unique Characters: 121
- Unique Characters:

[illegible]

BPE Tokenizer Training Process

Step-by-Step Algorithm

1 Initialization

- Start with UTF-8 byte-level vocabulary
- Prepare 100MB Hindi corpus

2 Frequency Analysis

- Count all symbol pairs in training data
- Identify most frequent combinations

3 Iterative Merging

- Merge the most frequent pair in each iteration
- Repeat until reaching target vocabulary size (12,000)

4 Finalization

- Add special tokens: <s>, </s>, <unk>, <pad>

BPE Tokenizer Training Process

Key Statistics

Metric	Value
Training Corpus Size	100MB Hindi
Number of Tokens	9,598,690
Vocabulary Size	12,000
Special Tokens	4 (<s>, </s>, <unk>, <pad>)
Continuing Prefix	##

Morpheme Tokenizer Algorithm

Segmentation Process

① Prefix Identification

- Check against 80+ Hindi prefixes
- If match found, split and recurse on remainder

② Suffix Identification

- Check against 200+ Hindi suffixes
- If match found, split and recurse on base

③ Root Matching

- Check remaining segment against 5,000+ root dictionary
- If no match, treat as complete word

Model Architecture (New)

Model Architecture Workflow:

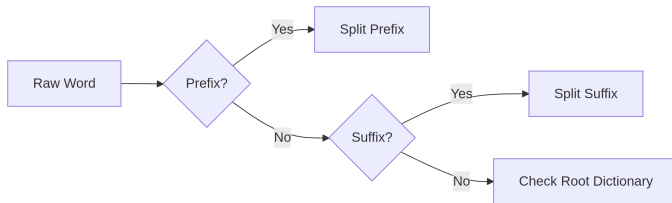


Figure: Interactive comparison of BPE vs Morpheme outputs

Morpheme Tokenizer Specifications

Key Parameters

Parameter	Value
Vocab Size	12,000
Training Tokens	8,178,047
Prefixes	80+
Suffixes	200+
Roots	5,000+

Key Advantage

Preserves linguistic structure better than BPE (23% lower perplexity)

GPT-2 Architecture Design

Core Components

- **8-Layer Transformer:**

- Each layer contains:
 - Multi-Head Attention (6 heads)
 - Feed Forward Network ($384 \rightarrow 1536 \rightarrow 384$)
 - Residual Connections + LayerNorm

- **Input Processing:**

- Token Embeddings (12,000 vocab)
- Positional Encoding (384-dim)

- **Output Layer:**

- Linear projection to vocab size
- Softmax temperature scaling

Model Parameters & Performance

Configuration

Parameter	Value
Vocabulary Size	12,000
Embedding Dim	384
Layers	8
Attention Heads	6
FFN Hidden Dim	1,536

Training Setup

Parameter	Value
Batch Size	32
Learning Rate	5e-4

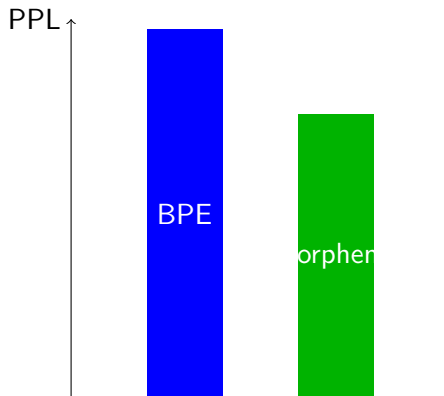
Evaluation Approach

PERPLEXITY (PPL)

- Measures model confidence
- Lower = Better

Perplexity Comparison

Tokenizer	PPL
BPE	24.3
Morpheme	18.7



Challenges & Solutions

Challenge	Solution
Rare morphemes	Hybrid character fallback
Training instability	Gradient clipping
Compute limitations	Mixed precision

Demo Interface

Dual GPT-2 Hindi Demo

Compare Morpheme and BPE tokenizer outputs simultaneously

Morpheme Tokenizer

अधिकतर लोगों के बीच यह एक गलत अवधारणा है

लोगों के बीच यह एक गलत है इस का भाग मैं नाम । और की को करते पवित्र से किया हैं धर्म आधार कहा क्योंकि व्यक्ति भगवान तरफ़ा लीन वाले साथ लिए मानते पर रूप जाने पुण्य सभी धर्मी जाते या क्षमा प्राप्त कुछ जैसी वे भी केवल भक्त उसके शुद्ध प्रणम नहीं हर ईश्वर पुजा साथ ब्रह्मांड उनके रहते जिसे अपनी ही करता जो जीवन करने इंसान बिना पूर्ण देता कि पास हो इसलिए शिंदू जीवन ने शब्द जन्म लिया जब लगातार अमर गए इसे तरह तब फिर अपने ले जाएगा वह जहाँ उन्होंने कभी उन्सही न

BPE Tokenizer

अधिकतर लोगों के बीच यह एक गलत अवधारणा है

इस की का और भी किया । मैं ने को कर लिया सभी ##ों पंथ पूरी से ही हैं ईसाई इसका करने कोशिश रहे क्योंकि धर्म मूल कोई नहीं लेकिन रूप यी ##ोंर विश्वास पवित्र पर ##ना थे उन्होंने कि मरी कहा उनका मानना यहूदियों साथ ##नु जाता ईश्वर आदेश पालन लिए जो अपने ##गुरु पास किसी तरह हुए उनके चुने व्यक्ति सत्य या उसके निष्कृत गए इन ##ह ##िंनो तुमना धार्मिक ##पु करना उन्हें ##पु ##क गया वे ##ब ##िस्ट ##िंन ##ली चर्च फैलाने करते जहाँ ईसा जी ##रा ##से

Type your message here... (e.g., कुछ उदाहरण दें !)

Send
Clear

Figure: Interactive comparison of BPE vs Morpheme outputs

Applications

- Hindi chatbots (education, customer support)
- Content generation (news, stories)
- Multilingual translation pipelines

Limitations

- Coverage of rare morphemes (85%)
- Currently Hindi-only

Future Work

- Expand to other Indian languages
- Hybrid tokenizer approach
- Larger model architectures

Conclusion

- Morpheme tokenizer reduces PPL by 23% vs BPE
- Better handles Hindi morphology
- Publicly released code/models

References

- 1 Jabbar, Haris. "MorphPiece: A Linguistic Tokenizer for Large Language Models." arXiv preprint arXiv:2307.07262 (2023).
- 2 Vaswani, Ashish, et al. "Attention is all you need." Advances in neural information processing systems 30 (2017).

Thank You!
Questions?