

# Report NLP Assignment-1

Akash Manna, 2024801008

February 2026

## 1 Task-1

Table 1: Analysis of Tokenization Methods: English

Lang	Tokenizer	Category	Example Snippet	Justification
En	Whitespace	Sensible	[‘easily’, ‘obtained’, ‘online’, ...], [‘The’, ‘following’, ‘steps’, ... ‘success’, ‘’], [‘Let’, ‘us’, ‘look’,... ‘afair’, ‘’]	Clean splitting on spaces; preserves whole words and semantic integrity.
	Whitespace	Not Sensible	[‘Saturday’, ‘Dec’, ‘1’, ‘-’, ‘10’, ...], [“”, ‘Fisherman’,... ‘ass’, ‘-’, ‘-’, ‘stand’, ‘back’, ‘!’, ...], [‘L’, ‘’, ‘is’,..., ‘is’]	Poor handling of dates and time ranges; symbols like ‘-‘ are isolated excessively.
En	Regex	Sensible	[“I’d”, ‘write’, ‘them’, ...], [‘Close’, ‘to’, ‘this’, ‘horn’, ...] [‘The’, ‘card’, ‘is’, ‘to’, ...]	Semantics preserved; punctuation like apostrophes handled within context effectively.
	Regex	Not Sensible	[‘A’, ‘testamur’, ‘(’], [‘This’, ‘feature’, ... ‘\$', ‘500’, ‘to’, ‘\$', ‘1’, ‘,’, ‘500’, ‘’], [‘@’, ‘proseimprint’, ‘@’, ...]	Over-segmentation of currency (“\$, ‘500’) and punctuation clusters; breaks logical entities.
En	BPE	Sensible	[‘e’, ‘at’, ‘en</w>’, ‘h’, ...], [‘N’, ‘a’, ..., ‘the</w>’] [‘un’, ... ‘been</w>’, ‘pri’, ‘v’, ‘il’, ‘e’, ‘g’, ‘ed</w>’, ‘to</w>’].	(Per input label): Functional, but exhibits extreme over-segmentation bordering on character-level tokens. ”privileged” being 5 pieces
	BPE	Not Sensible	[‘str’, ‘u’, ‘gg’, ‘le,</w>’, ...], [ ‘j’, ‘us’, ‘ti’, ‘fi’, ‘ab’, ‘le,</w>’,...], [‘re’, ‘li’, ‘gi’, ‘ous’, ...]	Inconsistent fragmentation; common roots (‘struggle’) are broken into non-semantic chunks.

Figure 1: Analysis of Tokenization Methods: Mongolian

## 2 Task-2: Language models

## 2.1 Perplexity Score

Table 2: Perplexity score of 9 Languagu models

Tokenizer	Smoothing	Perplexity Score
Whitespace	None	66771.900
	Witten-Bell	3911.900
	Kneser-Ney	37846.629
Regex	None	72964.946
	Witten-Bell	4619.654
	Kneser-Ney	43083.070
BPE	None	33893.713
	Witten-Bell	278.490
	Kneser-Ney	9463.389

## 2.2 Comparative study

- Smoothing Impact: None → MLE (crashes on unseen) → Witten–Bell (light discount) → Kneser–Ney (continuation-aware, best).
  - Tokenization Impact: Whitespace (coarse, fast) → Regex (syntax-aware) → BPE (morphological, generalizable).
  - Best Overall: BPE + Kneser–Ney (lowest perplexity, fluent, robust). Fastest Baseline: Whitespace + None (but poor quality).
  - Syntax-Critical: Regex + Witten–Bell (dialogue, punctuated text).

### 2.3 Examples:

Table 3: Comparison of Language Models: Tokenization, Smoothing, and Performance

Tokenization	Smoothing	Sparsity Handling	Key Strength	Key Weakness
Whitespace	None (MLE)	Zero-probability crashes	Simplicity, speed	Repetition on unseen n-grams
	Witten–Bell	Reserves mass for novel events	Better generalization	Still favors seen n-grams
	Kneser–Ney	Continuation probability	Semantic coherence, lowest perplexity	Computational cost
Regex	None (MLE)	Zero-probability crashes	Syntax preservation	Extreme sparsity, broken structure
	Witten–Bell	Reserves mass for punct. pairs	Syntax-aware smoothing	Moderate perplexity
	Kneser–Ney	Continuation-aware	Best syntax handling	Slowest training
BPE	None (MLE)	Morphological fallback	Handles novel words	Overfits subword patterns
	Witten–Bell	Subword generalization + smoothing	Morphological robustness	Moderate for unseen forms
	Kneser–Ney	Continuation + subword	Best overall: low perplexity, fluent	Highest complexity

Table 4: Analysis of Language Model Failure Modes and Root Causes

Model	Primary Failure Mode	Root Cause
WS + None	Repetition loops on rare contexts	Zero-probability MLE crash → argmax locks on high-freq token
WS + Witten-Bell	Still vulnerable to sparse n-grams	Discount $T/(count + T)$ is too weak ( $\sim 10\text{--}20\%$ reserved mass)
WS + Kneser-Ney	Only fails on unseen context diversity	Continuation probability cannot invent contexts not present in training data
Regex + None	CATASTROPHIC punctuation repetition	Explicit tokenization → sparse (word, punct.) pairs → MLE collapse
Regex + Witten-Bell	Punctuation context bias remains	Witten-Bell discount treats word-word and word-punct. pairs equally
Regex + Kneser-Ney	Nearly never fails (except ultra-rare cases)	Continuation probability deprioritizes punctuation; semantic recovery via backoff
BPE + None	Rare subword n-gram crashes	BPE reduces OOV (Out-Of-Vocabulary) but does not fix zero-probability n-grams
BPE + Witten-Bell	Subword fragmentation drift	Merges can vary; context learned on different fragmentation leads to mismatch
BPE + Kneser-Ney	Minimal failure (only cross-domain/adversarial)	Morphology + continuation probability + recursive backoff = trifecta defense

Table 5: Scenarios where the model performs well

Model	Best Scenario
WS + None	Exact phrase match present in training data
WS + Witten-Bell	Moderate-frequency phrases in a balanced domain
WS + Kneser-Ney	Any standard input, handling both rare and common contexts effectively
Regex + None	Unpunctuated simple phrases (behaves like whitespace)
Regex + Witten-Bell	Semi-rare phrases without heavy punctuation
Regex + Kneser-Ney	Punctuation-heavy phrases (e.g., titles, lists)
BPE + None	Morphologically common words where subwords are frequent
BPE + Witten-Bell	Rare words with clear morphological structure (prefixes/suffixes)
BPE + Kneser-Ney	Any input (rare, punctuation-heavy, or Out-Of-Vocabulary)