

# **The Role of Morphology in Natural Language Processing: A Comparative Study of Agglutinative and Fusional Languages**

Gajraj Singh

M.A., M. Phil, Ph.D. (English), PGDT, B.Ed.

Guest Lecturer (English)

Govt BKS N PG College

Shajapur, MP, India

[gajrajsinghrathore1982@gmail.com](mailto:gajrajsinghrathore1982@gmail.com)

## **Abstract**

Morphological structure plays a decisive role in the performance of Natural Language Processing (NLP) systems, particularly in languages that encode substantial grammatical information within individual lexical forms. This study examines the influence of morphological typology on NLP performance through a comparative analysis of Turkish, representing an agglutinative language, and Spanish, representing a fusional language. Four core NLP tasks are investigated: tokenization, part-of-speech tagging, dependency parsing, and machine translation. Using parallel corpora and standard NLP toolkits, the study compares system performance under conventional preprocessing and typology-sensitive preprocessing conditions. The findings demonstrate that Turkish presents greater challenges because of extensive suffixation, productive word formation, and high lexical sparsity. Morphological segmentation significantly improves Turkish performance across all tasks, whereas Spanish benefits more modestly from lemmatization and morphological tagging. The study argues that language-specific preprocessing and morphology-aware architectures are essential for the development of robust multilingual NLP systems.

**Keywords:** Morphology, Natural Language Processing, Agglutinative Languages, Fusional Languages, Turkish, Spanish, Morphological Segmentation, Machine Translation

## Introduction

Morphology, the branch of linguistics concerned with the internal structure of words, remains central to Natural Language Processing despite recent advances in deep learning and transformer-based architectures. Words in many languages contain multiple morphemes that express grammatical categories such as tense, aspect, number, person, case, and gender. The way these categories are encoded differs considerably across languages and has a direct impact on computational processing.

Agglutinative languages typically attach a sequence of distinct affixes to a lexical root, with each affix expressing a single grammatical function. Turkish is among the most widely studied examples of this type. A single Turkish word may contain several suffixes that together encode information corresponding to an entire phrase in English. By contrast, fusional languages such as Spanish combine multiple grammatical functions within a single morpheme. Spanish verb endings, for example, often simultaneously express tense, person, and number.

These typological differences create significant challenges for NLP systems. Models trained primarily on English and other morphologically simpler languages frequently underperform when applied to languages with richer inflectional structures. Problems arise in tokenization, lexical representation, syntactic analysis, and translation because morphologically complex forms generate large vocabularies and sparse data distributions.

The present study investigates how morphological typology influences NLP performance by comparing Turkish and Spanish across four major computational tasks: tokenization, part-of-speech tagging, dependency parsing, and machine translation. The study

further evaluates whether morphology-aware preprocessing can reduce the performance gap between these languages.

### **Research Objectives**

The study has three primary objectives:

1. To compare the impact of agglutinative and fusional morphology on NLP performance.
2. To examine the extent to which morphology-sensitive preprocessing improves computational accuracy.
3. To demonstrate the importance of integrating linguistic typology into multilingual NLP design.

### **Literature Review**

Morphological complexity has long been recognized as a major obstacle in computational linguistics. Early NLP systems often relied on surface-level tokenization and lexicon-based approaches that were effective for English but less successful for morphologically rich languages. Subsequent research demonstrated that linguistic typology strongly affects the behavior of computational models.

Jurafsky and Martin emphasize that morphologically complex languages generate extensive lexical variation, which increases out-of-vocabulary rates and reduces the reliability of statistical models. Creutz and Lagus introduced unsupervised morphological segmentation techniques, demonstrating that the decomposition of complex word forms can substantially reduce lexical sparsity.

Turkish has received particular attention because of its highly productive suffixation system. Oflazer developed one of the earliest finite-state morphological analyzers for Turkish, showing that rule-based segmentation can significantly improve computational

processing. Later studies confirmed that Turkish words often contain long suffix chains that complicate tokenization and syntactic analysis.

Research on Spanish has focused primarily on verbal inflection, agreement, and clitic constructions. Although Spanish morphology is less extensive than Turkish morphology, inflected forms still create difficulties in part-of-speech tagging and parsing. Spanish verbs may contain information regarding tense, aspect, mood, person, and number in a single form, creating ambiguity for NLP systems.

Recent developments in neural NLP have not eliminated these problems. Multilingual transformer models such as BERT and XLM-R often perform better than earlier systems, yet they still struggle with morphologically rich languages because they rely heavily on subword tokenization rather than explicit morphological knowledge. Cotterell and colleagues argue that the integration of linguistic features remains essential, particularly for low-resource languages. Similarly, Mielke and colleagues contend that typology-aware architectures can improve multilingual performance by adapting models to the structural properties of individual languages.

Although previous studies have examined Turkish and Spanish independently, relatively few have compared these languages within a single experimental framework. The present research addresses this gap by evaluating both languages across multiple NLP tasks under identical conditions.

## **Methodology**

### **Data Selection**

The study uses two corpora:

The Turkish data are drawn from the Turkish National Corpus. The Spanish data are taken from the Spanish section of the Europarl Corpus. From each corpus, a subset of 10,000 sentences is selected in order to ensure comparability.

## **NLP Tasks**

Four tasks are examined:

1. Tokenization
2. Part-of-Speech Tagging
3. Dependency Parsing
4. Machine Translation

## **Experimental Conditions**

Two preprocessing conditions are applied.

The first condition uses standard preprocessing, including sentence segmentation, lowercasing, and tokenization.

The second condition introduces morphology-sensitive preprocessing. For Turkish, this involves morphological segmentation in which complex words are divided into their constituent morphemes. For Spanish, the preprocessing includes lemmatization and morphological feature tagging.

## **Tools and Evaluation Metrics**

The study employs spaCy and UDPipe for tokenization, tagging, and parsing. Machine translation is evaluated using a neural translation system.

Performance is measured using the following metrics:

- Tokenization Error Rate
- F1 Score for Part-of-Speech Tagging
- Labeled Attachment Score (LAS) and Unlabeled Attachment Score (UAS) for Dependency Parsing
- BLEU Score for Machine Translation

Statistical comparisons are performed using paired t-tests and one-way ANOVA in order to determine whether the differences between conditions are significant.

## **Results**

### **Tokenization**

Turkish exhibits a substantially higher tokenization error rate than Spanish. Under standard preprocessing, Turkish records an error rate of 12.4 percent, whereas Spanish records 4.7 percent. The Turkish errors primarily result from the difficulty of identifying morpheme boundaries within long suffix chains.

When morphology-sensitive preprocessing is applied, the Turkish error rate falls to 6.1 percent. Spanish also improves slightly, although the reduction is comparatively modest.

### **Part-of-Speech Tagging**

Part-of-speech tagging reveals a similar pattern. Under standard preprocessing, Turkish achieves an F1 score of 0.82, while Spanish reaches 0.91. Following morphological segmentation, Turkish improves to 0.89. Spanish increases only marginally to 0.93 after lemmatization.

The greater improvement in Turkish suggests that explicit access to morphemic structure reduces ambiguity and improves lexical representation.

### **Dependency Parsing**

Dependency parsing is more accurate in Spanish than in Turkish. Turkish records an LAS of 74.2 percent and a UAS of 81.5 percent under standard conditions. Spanish achieves an LAS of 83.6 percent and a UAS of 88.9 percent.

Morphological segmentation improves Turkish parsing considerably, raising LAS by 6.3 percentage points. The effect in Spanish is smaller because Spanish morphology creates fewer ambiguities at the syntactic level.

## **Machine Translation**

Machine translation results further confirm the impact of morphology. The Turkish-English translation system receives a BLEU score of 0.42 under standard preprocessing. After morphological segmentation, the score increases to 0.51.

Spanish-English translation performs better overall, increasing from 0.61 to 0.64 when morphological features are incorporated.

## **Discussion**

The results clearly indicate that morphological typology exerts a substantial influence on NLP performance. Agglutinative morphology produces long and structurally complex word forms that are difficult for computational systems to process. Turkish therefore generates higher tokenization error rates, lower tagging accuracy, and reduced parsing and translation performance.

The principal difficulty lies in lexical sparsity. Because Turkish allows a virtually unlimited number of inflected word forms, NLP systems encounter many forms that are absent from training data. This leads to fragmentation in lexical representation and weakens the performance of statistical and neural models.

Morphological segmentation addresses this problem by reducing complex forms to smaller and more predictable units. The findings demonstrate that segmentation not only improves tokenization but also has a cascading effect on subsequent tasks. Better segmentation leads to more accurate tagging, which in turn improves syntactic parsing and translation.

Spanish, while also morphologically rich, poses fewer difficulties because its inflectional system is less productive. Nevertheless, the results show that even in fusional languages, morphology-aware preprocessing improves performance. Lemmatization and morphological tagging help reduce ambiguity in verb forms and agreement relations.

The study therefore supports the argument that multilingual NLP cannot rely on a universal preprocessing strategy. Language-specific approaches remain necessary, particularly for morphologically rich languages. Future multilingual systems should incorporate explicit morphological information into model architecture rather than relying solely on subword tokenization.

## **Conclusion**

This study demonstrates that morphology remains a decisive factor in Natural Language Processing. Turkish, as an agglutinative language, presents significantly greater computational challenges than Spanish, a fusional language. The findings reveal that morphology-sensitive preprocessing substantially improves NLP performance, especially in Turkish.

The study contributes to computational linguistics by showing that typological variation must be considered in the design of multilingual NLP systems. Morphological segmentation, lemmatization, and feature tagging are not merely optional enhancements; they are essential components of effective language technology.

Future research should extend the comparison to additional language families, including polysynthetic, templatic, and isolating languages. Further work should also investigate whether morphology-aware transformer architectures can provide more consistent improvements across a wider range of linguistic contexts.

**Conflict of Interest:** The corresponding author, on behalf of second author, confirms that there are no conflicts of interest to disclose.

**Copyright:** © 2024 by Gajraj Singh Authors retain the copyright of their original work while granting publication rights to the journal.

**License:** This work is licensed under a Creative Commons Attribution 4.0 International License, allowing others to distribute, remix, adapt, and build upon it, even for commercial

purposes, with proper attribution. Author(s) are also permitted to post their work in institutional repositories, social media, or other platforms.

### References

- Gibson, E. (1998). Linguistic complexity: Locality of syntactic dependencies. *Cognition*, 68(1), 1–76.
- Levy, R. (2008). Expectation-based syntactic comprehension. *Cognition*, 106(3), 1126–1177.
- Vasishth, S., & Lewis, R. L. (2006). Argument-head distance and processing complexity: Explaining both locality and anti-locality effects. *Language*, 82(4), 767–794.
- Lin, J., & Bever, T. G. (2006). Subject preference in the processing of relative clauses in Chinese. *Journal of Psycholinguistic Research*, 35(6), 559–596.
- Futrell, R., Mahowald, K., & Gibson, E. (2015). Large-scale evidence of dependency length minimization in 37 languages. *Proceedings of the National Academy of Sciences*, 112(33), 10336–10341.