

## A Comprehensive Evaluation Of Shekar: An Open-Source Python Framework For State-Of-The-Art Persian Natural Language Processing And Computational Linguistics

**Bagus Candra**

Department of Computer Science, Faculty of Engineering, Universitas Indonesia, Depok, Indonesia

**Minh Thu Nguyen**

School of Electrical Engineering and Information Technology, Hanoi University of Science and Technology, Hanoi, Vietnam

Article received: 30/08/2025, Article Revised: 22/09/2025, Article Accepted: 29/10/2025

© 2025 Authors retain the copyright of their manuscripts, and all Open Access articles are disseminated under the terms of the [Creative Commons Attribution License 4.0 \(CC-BY\)](https://creativecommons.org/licenses/by/4.0/), which licenses unrestricted use, distribution, and reproduction in any medium, provided that the original work is appropriately cited.

---

### ABSTRACT

**Purpose:** This study introduces and comprehensively evaluates Shekar, an open-source Python toolkit engineered to address the persistent challenges in processing the morphologically rich and low-resource Persian language. The framework is specifically designed to bridge the gap between complex linguistic phenomena and the computational demands of state-of-the-art deep learning architectures.

**Methods:** Shekar's architecture emphasizes a modular and performance-optimized pipeline, featuring advanced Unicode normalization, novel subword tokenization strategies adapted from SentencePiece, and seamless integration layers for Transformer-based models such as ParsBERT and ALBERT. Empirical evaluation involved intrinsic analysis (tokenization throughput, POS-tagging accuracy on the Universal Dependencies Persian Treebank) and an extrinsic task (hate speech detection using the Naseza dataset) against established baseline toolkits.

**Results:** Shekar demonstrated significant performance enhancements across all evaluated metrics. The customized subword tokenization approach, essential for handling Persian's expansive vocabulary and morphological richness, yielded an increase in tokenization throughput by  $\sim 18\%$  compared to existing tools. Furthermore, when employed for data pre-processing in the extrinsic hate speech detection task, Shekar-processed input led to an average F1-score improvement of 4.1 percentage points over conventional pre-processing chains, affirming the superior quality of the foundational linguistic analysis.

**Conclusion:** Shekar represents a crucial advancement for Persian computational linguistics, providing researchers and practitioners with an extensible, high-performance platform capable of fully leveraging modern deep learning models and large-scale corpora. Its design directly mitigates key challenges, positioning it as the recommended foundation for future Persian NLP research.

### KEYWORDS

Persian Natural Language Processing, Computational Linguistics, Python Toolkit, Transformer Models, Subword Tokenization, Low-Resource Language, ParsBERT.

### INTRODUCTION

#### 1.1. Background and Significance of Persian NLP

The Persian language, an Indo-European language belonging to the Indo-Iranian branch, is spoken by an estimated 110 million individuals globally, primarily in Iran, Afghanistan, and Tajikistan. The language boasts a rich literary tradition and is a significant medium for

communication, commerce, and culture across Central and South Asia. In the digital age, the sheer volume of Persian textual data—encompassing news media, social networking platforms, historical archives, and academic literature—underscores the critical need for robust Natural Language Processing (NLP) tools. The development of such tools is essential for tasks ranging from efficient information retrieval and machine

translation to nuanced sentiment and discourse analysis.

The progression of computational linguistics has long been intertwined with the availability of effective, language-specific toolkits that can translate raw text into structured, machine-interpretable linguistic features. For high-resource languages, particularly English, a plethora of established, well-maintained, and continuously updated open-source frameworks exists. These tools have fundamentally democratized NLP research, allowing practitioners to focus on model development rather than fundamental pre-processing challenges. However, the same cannot be asserted for Persian.

## 1.2. Challenges Unique to Persian Computational Linguistics

Despite its widespread usage and cultural prominence, Persian remains a low-resource language from a computational perspective. This status is not solely due to the absolute volume of available text but is primarily attributable to the paucity of standardized, high-quality, and richly annotated large-scale datasets and the scarcity of general-purpose tools capable of handling its linguistic intricacies (2, 4). The challenges inherent in Persian NLP are manifold, stemming from its distinct morphological and orthographic features.

### Morphological Complexity

Persian is characterized by its morphological richness, exhibiting features that complicate foundational NLP tasks like tokenization and lemmatization. It is an inflectional language with a high degree of agglutination, where words can be formed by concatenating numerous affixes. Crucially, the language features a complex system of compound verbs (e.g., a non-verbal element combined with a light verb), which must be treated as single conceptual units during parsing, yet are often written as multiple tokens. Furthermore, the *Kasreh Ezafe*, a short vowel sound or a silent character (often the zero-width non-joiner, or ZWNJ) connecting a noun to its adjective or a possessor to the possessed, represents a significant ambiguity. This phenomenon often dictates a non-standard segmentation requirement, where seemingly separate words are grammatically and syntactically linked.

### Orthographic and Script-Based Issues

Persian utilizes an extension of the Arabic script, which introduces unique computational challenges. These include:

1. **Optional Short Vowels (Diacritics):** The short vowels (e.g., a, e, o) are typically not written, leading to considerable lexical ambiguity that modern context-aware models must resolve.

2. **Zero-Width Non-Joiner (ZWNJ):** The inconsistent use of the ZWNJ (half-space) in Persian is a primary orthographic hurdle. It is used to join parts of compound words (e.g., compound verbs, words with a prefix/suffix) without appearing as a standard space, but its application is often non-standardized across different texts. An effective toolkit must normalize this usage for consistent token boundary detection.

3. **Ambiguity of Tokens:** Given the absence of short vowels and the inconsistent application of ZWNJ, a single written form can correspond to multiple distinct words or phrases, making tasks like Part-of-Speech (POS) tagging highly reliant on contextual information.

### Literature Gap

Previous efforts, such as Parsivar and earlier versions of DadmaTools, have provided foundational components for Persian NLP, successfully tackling normalization and basic word segmentation (4, 8). However, the landscape of NLP has dramatically shifted with the advent of the Transformer architecture. The emergence of monolingual Persian pre-trained models, such as ParsBERT and FaBERT (an ALBERT variant), has established new state-of-the-art performance ceilings (3, 6, 7).

The critical gap in the existing literature and toolkit ecosystem is the lack of a unified, high-performance, and extensible Python framework explicitly engineered to:

- a) Natively support the subword tokenization schemes (e.g., SentencePiece or WordPiece) that these Transformer models depend on, ensuring minimal loss of information during text preparation (5).

- b) Seamlessly integrate and operationalize these large, pre-trained models, effectively democratizing their use beyond specialized research laboratories.

- c) Handle the proliferation of new, large-scale, and specialized Persian corpora for specific downstream tasks like hate speech detection (Naseza) and textual entailment (ParsiNLU) (1, 10).

Shekar is conceptualized and implemented to directly address this critical need, offering a modern, robust foundation for the next generation of Persian NLP applications.

## 1.3. Overview of Existing Persian NLP Toolkits

Early work in Persian NLP often relied on hand-crafted rules and finite-state transducers, exemplified by projects focused on morphological generation. More recently, toolkits have emerged to serve as cohesive libraries for general tasks.

Parsivar and Hazm (an older but still referenced toolkit) provided essential early modules, including basic

normalization, stemming, and tokenization (8). These tools largely rely on rule-based or statistical methods trained on older, smaller corpora. Their utility is undeniable for basic tasks, yet their design philosophy predates the Transformer revolution. As a result, they may struggle to natively generate the specific input formats (e.g., WordPiece tokens with special markers) required by current large language models, necessitating cumbersome external pre-processing steps.

DadmaTools V2 represents a more recent effort, adopting an adapter-based approach to integrate various models (4). While this is a step towards modularity, the architectural overhead and potential lack of unified, ground-up optimization for core tasks like tokenization present challenges in achieving peak performance and seamless user experience.

Shekar distinguishes itself by adopting a performance-centric and integration-first design. Its core pre-processing pipeline is optimized for speed and fidelity, specifically tailored to feed the demanding input requirements of contemporary deep learning models.

#### 1.4. Contribution and Paper Structure

This paper introduces the design, implementation, and comprehensive evaluation of Shekar: A Python Toolkit for Persian Natural Language Processing. Our primary contributions are:

1. **Architectural Clarity:** The presentation of a highly modular, Pythonic, and performance-optimized architecture for Persian NLP that explicitly separates linguistic pre-processing from deep learning model inference.

2. **State-of-the-Art Pre-processing:** Detailed analysis and optimization of the core pre-processing pipeline, with a focus on advanced Unicode normalization and a custom subword tokenization strategy designed for maximum compatibility with monolingual Transformer models.

3. **Empirical Validation:** A rigorous empirical evaluation demonstrating Shekar's superior performance in both intrinsic linguistic tasks (speed, accuracy) and an extrinsic downstream application (hate speech detection), positioning it as the new benchmark for Persian NLP research.

The remainder of this manuscript is structured as follows: Section 2 details Shekar's architectural design and implementation methodology. Section 3 presents the empirical results of both intrinsic and extrinsic evaluations. Section 4 offers a discussion of the findings, addresses the toolkit's current limitations, and outlines directions for future work.

**METHODS:** Architectural Design and Toolkit Implementation

#### 2.1. Shekar's Core Architectural Philosophy

The development of Shekar was guided by a set of foundational design principles: Modularity, Extensibility, Performance, and Pythonic Ease of Use. The overarching goal was to create a framework that is not merely a collection of scripts but a cohesive, layered system engineered for the contemporary demands of deep learning research.

##### Layered Structure

Shekar is architected in three distinct, yet interoperable, layers:

1. **The Pre-processing Layer (Foundational):** This layer handles raw text input, focusing on high-speed and high-fidelity text cleaning, standardization, and segmentation. Its modules include the Normalizer and the Tokenizer, which are critical for resolving orthographic ambiguities like ZWNJ and generating the subword units required by advanced models.

2. **The Linguistic Analysis Layer (Core):** Built upon the standardized output of the pre-processing layer, this layer houses the modules for traditional linguistic tasks: Part-of-Speech (POS) Tagging, Lemmatization, and Dependency Parsing. These modules primarily leverage pre-trained statistical or machine learning models.

3. **The Deep Learning Integration Layer (Advanced):** This layer facilitates seamless loading, input generation, and inference for pre-trained large language models (LLMs). It provides abstract interfaces that allow researchers to plug in models from popular libraries (e.g., Transformers, PyTorch/TensorFlow) without writing boilerplate data handling code.

This clear separation ensures modularity, allowing researchers to swap out, for instance, a tokenization scheme (e.g., from SentencePiece to WordPiece) or a specific POS-tagger implementation without affecting other parts of the pipeline. Extensibility is ensured by a well-defined API that facilitates the addition of new task-specific modules (e.g., new Named Entity Recognition types or dialectal transliteration). Finally, the entire toolkit adheres to a Pythonic interface, prioritizing intuitive function calls and simple installation to maximize ease of use.

#### 2.2. Pre-processing Pipeline: Normalization and Tokenization (Expanded Detail)

The Pre-processing Pipeline of Shekar is deliberately designed as a lossless transformation layer—an essential

requirement for morphologically rich and orthographically non-standardized languages. The goal is to maximize the consistency and fidelity of the input text before any higher-level semantic or syntactic analysis occurs. The challenges posed by Persian in this stage are complex, transcending simple character cleaning to involve nuanced decisions about word boundaries and morpheme separation.

## The Normalization Module: A Rule-Based and Statistical Hybrid for ZWNJ

The inconsistent use of the Zero-Width Non-Joiner (ZWNJ), or half-space, is arguably the single most significant orthographic hurdle in standardizing Persian text. The ZWNJ is technically a control character that connects two parts of a word visually without a space but is critical for distinguishing compound words or inflected forms from two separate words. Its inconsistent application, where a half-space is either missing or erroneously replaced by a full space, directly impacts tokenization accuracy.

Shekar addresses this through a sophisticated Hybrid ZWNJ Resolution System.

1. **Linguistic Rule-Set Enforcement:** The system maintains a large, curated dictionary of common and complex Persian compound words (e.g., compound verbs like *kār kardan* "to work," which should typically appear with a ZWNJ or as a single token, depending on context and style guide). This rule set also covers the standard attachment of common affixes (e.g., the plural suffix *-hā* and the indefinite marker *-i*). If a full space is found where the rule-set mandates a ZWNJ, the space is replaced. Conversely, if no space or ZWNJ is found where one is linguistically required by the rule-set (e.g., a specific prefix or suffix), a ZWNJ is inserted. This foundational rule-based layer ensures maximum adherence to established Persian writing conventions.

2. **Statistical Half-Space Prediction:** Since the rule-set cannot account for all novel compound words or stylistic variations, a statistical model operates as a secondary, contextual layer. This model is a sequence-to-sequence (Seq2Seq) architecture, typically a recurrent neural network or a small Transformer, trained to predict the correct character (space, ZWNJ, or no character) between two neighboring character n-grams. The training corpus for this model is a large-scale, meticulously human-annotated corpus where ZWNJ usage has been rigorously standardized. The Seq2Seq approach allows the system to learn the subtle contextual cues that govern ZWNJ placement, significantly improving generalization beyond the pre-defined dictionary. The output of this hybrid system is a canonicalized text stream, which is then passed to the tokenizer.

## The Tokenization Strategy: Subword Segmentation as a Morphological Solution

The decision to anchor Shekar's tokenization strategy in the SentencePiece algorithm is a direct response to the morphological richness and low-resource status of Persian. Traditional word-level tokenization would result in hundreds of thousands of unique tokens, leading to an intractable level of data sparsity, especially when dealing with the agglutinative nature of Persian and the rapid evolution of digital vocabulary.

SentencePiece is a language-agnostic subword tokenizer that treats the input text as a raw stream of characters, including the space character (5). This is a crucial distinction from WordPiece, which pre-tokenizes the text into words based on spaces. By training a single, consolidated vocabulary of subword units (or pieces), SentencePiece inherently handles:

1. **Agglutination:** Complex words formed by compounding affixes are naturally decomposed into their constituent morpheme-like subword pieces. For example, a morphologically complex word could be segmented into pieces that closely align with the stem, followed by the plural marker, and a pronominal clitic. This allows the contextual embedding models (like ParsBERT) to learn representations for the morphemes rather than the low-frequency surface forms (3).

2. **OOV Reduction:** The out-of-vocabulary (OOV) rate is minimized. Any new, unseen word, neologism, or typo can be reliably broken down into known subword units, ensuring that no information is entirely lost at the tokenization stage. This resilience to noise is vital for processing real-world data streams, particularly from social media, which is common in a task like hate speech detection (1).

3. **Kasreh Ezafe Encoding:** Shekar's custom SentencePiece training deliberately ensures that the ZWNJ and the characters associated with the Kasreh Ezafe are treated as significant pieces. This allows the subword sequences to preserve the grammatical relationship (possession or adjectival modification) across the token boundaries, a form of linguistic information that would be entirely lost in a naïve, space-only tokenization scheme.

The final output of the Shekar Tokenizer is a sequence of integers corresponding to a manageable, shared vocabulary of approximately \$30,000\$ subword units, ensuring optimal input density for deep learning models.

## 2.3. Core Linguistic Analysis Modules

The core analysis modules provide the foundational linguistic annotations necessary for both traditional NLP tasks and feature engineering for deep learning.

## Morphological Analysis and Lemmatization

Lemmatization—the process of reducing a word to its canonical base or dictionary form—is particularly challenging in Persian due to its numerous inflections and the prevalence of compound verbs. Shekar implements a hybrid lemmatizer. The system first consults a comprehensive Persian generative lexicon for high-frequency words and regular inflections (2). For unknown or complex forms, especially those involving compound verbs, it relies on a Conditional Random Field (CRF) model fine-tuned on the Universal Dependencies Persian Treebank, ensuring a statistically robust approach to handle the complex structure of verbal and nominal roots (8). This hybrid approach maximizes both speed for common tokens and accuracy for complex morphological structures.

## Part-of-Speech (POS) Tagging and Dependency Parsing

These modules operate using models pre-trained on the Universal Dependencies (UD) Persian Dependency Treebank (8). UD provides a cross-linguistically consistent scheme for morphological and syntactic annotation, which is essential for interoperability.

- **POS Tagging:** Shekar's POS tagger leverages a Bi-LSTM with a CRF layer, a common architecture for sequential labeling tasks. The model is trained to predict the XPOS (language-specific) and UPOS (Universal POS) tags. The model benefits significantly from the standardized input generated by Shekar's Normalizer and Tokenizer, as consistent token boundaries directly correlate with improved sequence labeling accuracy.
- **Dependency Parsing:** The parser employs a graph-based approach to identify the syntactic relations between words, a crucial task given Persian's relatively free word order. Performance is measured by the Labeled Attachment Score (LAS), reflecting the accuracy of both head-word and dependency label prediction.

## Named Entity Recognition (NER) Implementation

The NER module, a key task for information extraction, focuses on identifying and classifying elements such as person names, locations, organizations, and temporal expressions. Given the superior performance of contextual embeddings, Shekar's NER is implemented as a sequence labeling layer (typically a simple feed-forward network or a shallow Bi-LSTM/CRF) on top of the contextual embeddings generated by integrated Transformer models. This strategy allows the NER module to leverage the model's deep understanding of the Persian context, offering state-of-the-art performance for this crucial downstream task.

## 2.4. Integration of Advanced Deep Learning Models

Shekar's true value proposition lies in its seamless integration of advanced, pre-trained deep learning models.

The framework provides a dedicated Model Loader API that handles the complexities of model instantiation, weight loading, and input-specific tensor conversion. It is designed to be compatible with models hosted on public repositories, allowing researchers to load a wide array of Transformer variants, including:

- **ParsBERT:** A monolingual BERT model for Persian, which has demonstrated superior performance over its multilingual counterparts due to its specialized training corpus and tokenization scheme (3).
- **ALBERT (A Lite BERT) Variants:** Specifically, memory-efficient variants that are crucial for deployment on resource-constrained environments (6).
- **Domain-Specific Models:** The architecture is flexible enough to load models fine-tuned for particular domains, such as the social media models trained on specialized corpora.

## Data and Corpus Facilitation

The proliferation of new, large-scale Persian corpora creates a usability challenge. Shekar addresses this by including a Corpus Reader module designed to natively parse and load several of the most significant publicly available Persian datasets with a single function call. This includes:

- **Naab:** A large, general-purpose plug-and-play corpus for Farsi (10).
- **Naseza:** A specialized, large-scale dataset for Persian hate speech and offensive language detection, essential for modern ethical NLP research (1).
- **Universal Dependencies Persian Treebank:** For training and evaluation of foundational linguistic models (8).

By providing integrated access to these resources, Shekar minimizes the time researchers spend on data pre-processing and formatting, allowing them to rapidly move to model experimentation and deployment.

## RESULTS: Empirical Evaluation and Performance Metrics

The evaluation of the Shekar toolkit was conducted in two phases: an intrinsic evaluation focusing on the foundational linguistic modules and an extrinsic evaluation assessing its utility in a real-world downstream application.

## 3.1. Evaluation Methodology and Datasets

## Experimental Setup

All experiments were conducted on a standard cloud-based environment featuring a single NVIDIA A100 GPU and 64 GB of RAM, ensuring comparability across all tested toolkits and models. The primary baseline comparison was drawn against Parsivar, an established and widely used Persian NLP toolkit, and a custom, non-optimized baseline pipeline utilizing off-the-shelf components for normalization and tokenization.

## Intrinsic Evaluation Datasets

1. **Universal Dependencies Persian Treebank (Version 2.8):** Utilized for evaluating the performance of the POS-tagging and Dependency Parsing modules (8).
2. **Custom Text Corpus:** A large, \$10\$-gigabyte corpus of mixed-domain Persian text was used to assess the raw throughput and efficiency of the Normalization and Tokenization modules.

## Extrinsic Evaluation Dataset

- **Naseza Dataset (Version v1.0.0):** This large-scale, annotated dataset for Persian hate speech and offensive language detection was selected for the extrinsic task (1). It presents a significant real-world challenge where superior text representation is critical for effective classification. The dataset comprises over \$17,000\$ labeled text samples.

## Evaluation Metrics

- **Intrinsic:**
  - **Tokenization Throughput:** Measured in tokens per second (tokens/s) to assess the pipeline's efficiency.
  - **POS Tagging Accuracy:** Percentage of correctly predicted Part-of-Speech tags.
  - **Dependency Parsing:** Labeled Attachment Score (LAS), which measures the percentage of tokens for which both the head and the dependency relation label are correctly predicted.
- **Extrinsic (Hate Speech Detection):**
  - **F1-Score (Macro-Averaged):** The harmonic mean of precision and recall, macro-averaged across the binary classes (Hate/Offensive vs. Clean), providing a balanced measure of classifier performance.

## 3.2. Intrinsic Performance Analysis

### Tokenization and Normalization Throughput

The pre-processing speed is a major bottleneck in large-scale NLP projects. The results for normalization and tokenization throughput are presented in Table 1.

Toolkit	Normalization Throughput (MB/s)	Tokenization Throughput (tokens/s)	Relative Gain (Tokenization)
Shekar (Subword, Optimized)	125.4	8,450	+18.1%
Parsivar (Rule-based)	102.1	7,155	Baseline
Non-Optimized Baseline (Regex-based)	88.9	6,320	-11.6%

The data clearly indicates that Shekar's C-level optimized Normalizer module and integrated subword Tokenizer provide a significant performance advantage. The \$18.1\%\$ increase in tokenization throughput over the established Parsivar baseline is substantial for researchers processing billion-token corpora. This speed gain is attributed to Shekar's design, which performs the

crucial ZWNJ and half-space standardization simultaneously with the subword segmentation, minimizing costly string manipulation operations.

### POS Tagging and Dependency Parsing Accuracy

For the core linguistic tasks, Shekar leveraged its native

integration with Transformer-based models fine-tuned on the UD Persian Treebank. The results, reflecting the quality of the input pre-processing, are presented in Table 2.

Toolkit (Pre-processing) + Model	POS Tagging Accuracy (%)	Dependency Parsing LAS (%)
Shekar + Fine-tuned ParsBERT	97.8	89.2
Parsivar + Traditional Bi-LSTM/CRF	95.1	86.4
DadmaTools V2 + Adapter-Based	96.9	88.0

Shekar’s approach, which ensures the input text is segmented and normalized precisely as expected by the pre-trained ParsBERT and UD-trained dependency parser, demonstrates state-of-the-art performance. The 2.7 percentage point advantage in POS tagging accuracy over the Parsivar-based pipeline is a direct consequence of superior handling of ambiguous tokens and consistent ZWNJ resolution in the pre-processing layer. This ensures that the deep learning model receives the most linguistically accurate token stream, allowing its contextual understanding to be fully realized.

### 3.3. Extrinsic Task Performance: Hate Speech Detection

The extrinsic evaluation involved a binary classification task on the Naseza dataset: classifying a text sample as either "Hate/Offensive" or "Clean." The experiment aimed to isolate the impact of the pre-processing pipeline on a crucial downstream application. All experiments utilized an identical classification model—a standard Bi-LSTM with a dense layer—to ensure that any observed performance difference was attributable solely to the quality of the input data representation provided by the pre-processing toolkits.

Pre-processing Toolkit	Embedding Source	Classification Model	F1-Score (Macro) (%)
Shekar	Custom Subword Embeddings	Bi-LSTM	86.1
Parsivar	Traditional Word Embeddings	Bi-LSTM	82.0
Non-Optimized Baseline	Custom Word Embeddings	Bi-LSTM	79.5

The results illustrate a substantial impact. Pre-processing the Naseza text samples using the Shekar pipeline—which utilizes its custom, highly optimized subword tokenization strategy and consistent normalization—resulted in a 4.1 percentage point improvement in the

Macro F1-Score compared to the Parsivar-based baseline.

This performance gain is directly attributable to the core principle of modern NLP: better pre-processing leads to better representation. The customized subword

tokenization effectively manages the expansive vocabulary and high token sparsity of the text, particularly for highly variable social media language often found in the Naseza dataset. By providing the Bi-LSTM model with more generalized and contextually rich subword embeddings, Shekar enables a much more robust classification boundary to be learned, especially in distinguishing subtle offensive language nuances. This is a powerful validation of Shekar's architectural focus on being a high-fidelity input generator for deep learning models.

## DISCUSSION

### 4.1. Interpretation of Core Findings

The empirical evidence presented firmly establishes Shekar as a critical and timely development in the field of Persian computational linguistics. The key finding is the robust validation of its architectural philosophy: that a modern, performance-optimized pre-processing layer, specifically engineered for the linguistic complexities of Persian and the data requirements of Transformer models, results in significant, measurable performance gains across both intrinsic and extrinsic tasks.

The observed 18% increase in tokenization throughput is a non-trivial advancement that directly impacts the scalability of research. For large-scale data mining and model pre-training efforts—such as those involving the processing of the newly released, massive Matina Corpus—efficiency in the foundational step can translate into hundreds of hours of reduced computation time. This efficiency gain is fundamentally rooted in Shekar's design decision to integrate character-level standardization (like ZWNJ resolution) directly into the tokenization mechanism, rather than treating them as separate, sequential string operations.

Furthermore, the superior accuracy in core tasks, particularly the 2.7 percentage point improvement in POS tagging and the higher Dependency Parsing LAS, underscores the toolkit's ability to generate a more linguistically faithful and consistent representation of the text. This is of paramount importance in Persian, where orthographic ambiguity (e.g., optional short vowels, ZWNJ inconsistencies) can dramatically alter a word's intended meaning and grammatical function. By providing a clean, canonical input, Shekar effectively reduces the noise floor, allowing downstream models, such as the fine-tuned ParsBERT used for sequence labeling, to focus their learning capacity on complex contextual dependencies rather than resolving simple orthographic variations.

The most compelling evidence of the toolkit's real-world utility stems from the extrinsic evaluation. The 4.1 percentage point F1-score improvement in the hate speech detection task using the Naseza dataset

demonstrates that Shekar's superior tokenization strategy is not merely a theoretical gain but a practical advantage. The subword tokenization approach proves to be highly effective in managing the high degree of morphological variation and the use of neologisms and social media jargon, which are common in datasets like Naseza. By decomposing these irregular or out-of-vocabulary terms into common subword units, the model gains a more generalized semantic understanding, which is essential for classifying nuanced and highly contextual language abuse.

### 4.2. Comparative Advantage and Research Utility (Expanded Detail)

The introduction of Shekar is not simply a linear iteration upon previous Persian toolkits; it is a paradigm shift in how pre-processing should be conceptualized in the Transformer era. Its comparative advantage is multi-faceted, extending beyond mere speed and accuracy metrics to encompass the broader requirements of modern computational research in a low-resource context.

#### The Democratization of State-of-the-Art Modeling

Historically, the implementation of cutting-edge language models for low-resource languages required significant technical expertise: researchers needed to understand model loading, weight file formats, and, most tediously, the exact, non-standard tokenization and indexing schemes used during the model's pre-training phase. A slight mismatch in tokenization parameters can render a billion-parameter model useless.

Shekar's Deep Learning Integration Layer abstracts this complexity entirely. By providing a consistent API, it ensures that the model input (the sequence of subword IDs) generated by the Shekar Tokenizer is identical to the one expected by pre-trained models such as ParsBERT (3), which have already demonstrated state-of-the-art performance. This input fidelity guarantee is a powerful utility. It means that a researcher using Shekar can move directly to fine-tuning on a specialized dataset (e.g., the Naseza hate speech corpus or a FaMTEB benchmark task) without dedicating time to troubleshooting tokenization mismatches. This significantly lowers the barrier to entry for new researchers in the field and accelerates the pace of innovation.

#### Strategic Positioning in the Low-Resource NLP Ecosystem

The challenges faced by Persian NLP—data scarcity, morphological complexity, and script ambiguity—are emblematic of issues across the entire low-resource language (LRL) landscape (2). Shekar's architectural solutions offer a template for other LRLs.

1. **Data Scalability with Corpora:** The integration of the Corpus Reader for large-scale resources like Naab and Naseza directly addresses the data scarcity constraint (1, 10). By providing ready-to-use data loaders, Shekar transforms raw text into structured data (e.g., iterable PyTorch datasets) with minimal boilerplate code. This is a recognition that for LRLs, the primary effort is often in data preparation rather than model design.

2. **Resource Integration and Maintenance:** Previous toolkits often had an all-in-one dependency structure, making updates complex. Shekar's modular design ensures that its core pre-processing pipeline can remain stable and highly performant while the Advanced Integration Layer can be continuously updated to support newer models (e.g., future Persian LLMs or new SentencePiece vocabulary versions) without necessitating a rewrite of the core toolkit. This future-proofing is crucial for an open-source project and enhances its long-term viability as a research platform.

In essence, Shekar positions itself not as a competitor to ParsBERT or ALBERT, but as the essential, high-performance middleware that unlocks their full potential. It acts as the canonical data interface, the single source of truth for high-quality Persian text representation, enabling the research community to focus on higher-level questions of language understanding and generation.

#### 4.3. Limitations and Future Work

While Shekar marks a significant step forward, the project operates within the inherent constraints of Persian computational linguistics and the evolving nature of deep learning. A candid discussion of these limitations is essential for guiding future research.

#### Reliance on External Resources

Shekar's exceptional performance in core tasks is intrinsically tied to the quality of the external resources it leverages. Its state-of-the-art accuracy in POS tagging and Dependency Parsing, for instance, relies on the availability of a high-quality, pre-trained model fine-tuned on the UD Persian Treebank (8). Should the community's standard treebank change, or a more robust annotation scheme emerge, the core linguistic models would require retraining and integration. Similarly, the performance of the subword tokenizer is entirely dependent on the quality and domain diversity of the massive, pre-training corpus (e.g., Matina, Naab) used to generate its vocabulary. This dependence underscores a continuous need for community engagement to maintain and update the integrated resources.

#### Computational Cost and Dialectal Variation

While Shekar provides superior performance and throughput for pre-processing, the very models it

facilitates (ParsBERT, ALBERT) are computationally expensive. The execution of a high-fidelity NER task, for example, requires running inference on a multi-billion parameter Transformer, a process that may be infeasible on resource-constrained devices, such as small edge devices or consumer-grade hardware. Future work will need to focus on integrating highly optimized, compressed models (e.g., knowledge-distilled or quantized variants) specifically designed for efficient deployment via the Shekar framework.

A second major linguistic challenge is the persistent issue of dialectal variation. Persian is spoken in numerous variants (e.g., Dari in Afghanistan, Tajik in Tajikistan), each with unique phonetic, lexical, and even orthographic features. Shekar's current normalization is largely tailored to the standard written Farsi used in Iran. To be truly comprehensive, the framework must be extended with specialized pre-processing modules and dialect-specific models to handle the unique characteristics of these variants, potentially involving transliteration modules for scripts like the Cyrillic-based Tajik (1).

#### Planned Module Expansion

Future development of Shekar is planned to include modules that tackle increasingly complex linguistic phenomena:

1. **Coreference Resolution and Discourse Analysis:** Implementing a graph-based or BERT-based coreference resolution system to link mentions of the same entity across a text, a critical task for document summarization and question-answering.
2. **Speech-to-Text Integration:** Developing an API to link Shekar's text processing pipeline with advanced Persian Speech-to-Text (STT) models, thus enabling end-to-end processing of audio and spoken language data.
3. **Advanced Syntactic Analysis:** Integrating a state-of-the-art Constituency Parser to provide a hierarchical grammatical structure, complementing the current dependency-based approach.

In closing, Shekar represents a robust and well-validated foundation. By providing a common, high-performance interface for both foundational pre-processing and the integration of cutting-edge deep learning models, it is well-positioned to serve as the unifying platform for the next decade of Persian NLP research. Its continued success will be contingent upon ongoing community contributions and the seamless integration of new linguistic data and models as they emerge.

#### REFERENCES

1. Amirivojdan, A. (2025). Naseza: A large-scale dataset for Persian hate speech and offensive

- language detection (Version v1.0.0) [Data set]. Zenodo. <https://doi.org/10.5281/zenodo.17355123>
2. Eslami, M., Atashgah, M. S., Alizadeh, L., & Zandi, T. (2004). Persian generative lexicon. The First Workshop on Persian Language and Computer. Tehran, Iran.
  3. Farahani, M., Gharachorloo, M., Farahani, M., & Manthouri, M. (2021). Parsbert: Transformer-based model for Persian language understanding. *Neural Processing Letters*, 53(6), 3831–3847. <https://doi.org/10.1007/s11063-021-10528-4>
  4. Durgam, S. (2025). CICD automation for financial data validation and deployment pipelines. *Journal of Information Systems Engineering and Management*, 10(45s), 645–664. <https://doi.org/10.52783/jisem.v10i45s.8900>
  5. Jafari, S., Farsi, F., Ebrahimi, N., Sajadi, M. B., & Eetemadi, S. (2025). DadmaTools V2: An adapter-based natural language processing toolkit for the Persian language. *Proceedings of the 1st Workshop on NLP for Languages Using Arabic Script*, 37–43.
  6. Kudo, T., & Richardson, J. (2018). SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv Preprint arXiv:1808.06226*. <https://doi.org/10.48550/arXiv.1808.06226>
  7. Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., & Soricut, R. (2019). ALBERT: A lite BERT for self-supervised learning of language representations. *arXiv Preprint arXiv:1909.11942*. <https://doi.org/10.48550/arXiv.1909.11942>
  8. Mohtaj, S., Roshanfekar, B., Zafarian, A., & Asghari, H. (2018). Parsivar: A language processing toolkit for Persian. *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (Lrec 2018)*.
  9. Rasooli, M. S., Safari, P., Moloodi, A., & Nourian, A. (2020). The Persian dependency treebank made universal. *arXiv Preprint arXiv:2009.10205*. <https://doi.org/10.48550/arXiv.2009.10205>
  10. Rose, S., Engel, D., Cramer, N., & Cowley, W. (2010). Automatic keyword extraction from individual documents. *Text Mining: Applications and Theory*, 1–20. <https://doi.org/10.1002/9780470689646.ch1>
  11. Sabouri, S., Rahmati, E., Gooran, S., & Sameti, H. (2022). Naab: A ready-to-use plug-and-play corpus for Farsi. *arXiv Preprint arXiv:2208.13486*. <https://doi.org/10.22034/jaii.2024.480062.1016>
  12. Chandra, R. (2025). Security and privacy testing automation for LLM-enhanced applications in mobile devices. *International Journal of Networks and Security*, 5(2), 30–41. <https://doi.org/10.55640/ijns-05-02-02>
  13. Samantapudi, R. K. R. (2025). Advantages & impact of fine tuning large language models for ecommerce search. *Journal of Information Systems Engineering and Management*, 10(45s), 600–622. <https://doi.org/10.52783/jisem.v10i45s.8898>