

## EVALUATING PREPROCESSING EFFECTS IN NAME RETRIEVAL USING CLASSICAL IR AND CNN-BASED MODELS

Frizca Fellicita Marcelly\*; Irwansyah Saputra; Muhammad Bagus Andra

Department of Computer Science  
Universitas Nusa Mandiri, Depok, Indonesia  
www.nusamandiri.ac.id

14240021@nusamandiri.ac.id\*, irwansyah.iys@nusamandiri.ac.id, muhammad.mba@nusamandiri.ac.id

(\*) Corresponding Author



The creation is distributed under the Creative Commons Attribution-NonCommercial 4.0 International License.

**Abstract**—Information Retrieval (IR) systems are pivotal for efficient data management, particularly in tasks involving name searches and entity identification. This study evaluates text preprocessing techniques, including case folding, phonetic normalization, and gender tagging, that affect the performance of classical (TF-IDF, LSI) and CNN-based retrieval models for multilingual name matching. Using a dataset of 365,468 globally diverse names, this study implements a preprocessing pipeline featuring: Double Metaphone phonetic preprocessing (92% validation accuracy), gender disambiguation for unisex names (92% accuracy), and optimized n-gram tokenization for short names. Evaluation metrics include precision, recall, F1-score, and our novel Name Similarity Score (NSS), combining orthographic and phonetic preprocessing. Results show our full pipeline improves recall to 1.00 and F1-score by 37% while reducing false negatives by 63%. Key findings reveal: TF-IDF achieves superior recall (0.98 vs CNN's 0.85), LSI handles cultural variants effectively, and CNNs deliver the highest precision (0.91 vs TF-IDF's 0.70), particularly for unisex names. This work contributes both a scalable multilingual preprocessing framework and the NSS evaluation metric for robust name retrieval systems.

**Keywords:** CNN, information retrieval, multilingual names, name retrieval, phonetic normalization.

**Abstrak**—Sistem Information Retrieval (IR) sangat penting untuk manajemen data yang efisien, khususnya dalam tugas pencarian nama dan identifikasi entitas. Studi ini mengevaluasi teknik text preprocessing, termasuk case folding, normalisasi fonetik, dan penandaan gender, yang memengaruhi kinerja model klasik (TF-IDF, LSI) dan model retrieval berbasis CNN untuk pencocokan nama multibahasa. Dengan menggunakan dataset berisi 365.468 nama beragam dari seluruh dunia, penelitian ini menerapkan alur preprocessing yang

mencakup: Double Metaphone phonetic preprocessing (akurasi validasi 92%), disambiguasi gender untuk nama uniseks (akurasi 92%), dan tokenisasi n-gram yang dioptimalkan untuk nama pendek. Metrik evaluasi meliputi precision, recall, F1-score, serta Name Similarity Score (NSS) yang merupakan metrik baru dengan menggabungkan orthographic dan phonetic preprocessing. Hasil penelitian menunjukkan bahwa alur lengkap kami meningkatkan recall hingga 1,00 dan F1-score sebesar 37% serta mengurangi false negatives sebesar 63%. Temuan utama mengungkapkan bahwa: TF-IDF mencapai recall tertinggi (0,98 dibandingkan 0,85 pada CNN), LSI efektif dalam menangani variasi budaya, dan CNN memberikan precision tertinggi (0,91 dibandingkan 0,70 pada TF-IDF), khususnya untuk nama uniseks. Karya ini berkontribusi pada kerangka preprocessing multibahasa yang dapat diskalakan serta metrik evaluasi NSS untuk sistem name retrieval yang lebih andal.

**Kata Kunci:** CNN, information retrieval, nama multibahasa, pencarian nama, normalisasi fonetik.

### INTRODUCTION

Information Retrieval (IR) systems are critical infrastructure for managing textual data across search engines, digital libraries, and identity verification platforms (X. Zhang et al., 2023). A persistent challenge in IR is the accurate retrieval of personal names, which are often lexically brief, vary across languages, and exhibit orthographic differences (e.g., "Catherine" vs. "Katherine") (Aso et al., 2020). These variations hinder exact matching and necessitate specialized preprocessing techniques to improve retrieval performance.

Traditional text preprocessing methods, such as case folding, lemmatization, and stop-word removal, are well-established for general document

retrieval (Boghara, 2025). However, their effectiveness diminishes when applied to names due to unique structural and phonetic characteristics (Zeng, 2025). For instance, phonetic algorithms like Double Metaphone (Raykar et al., 2024) address spelling with different meanings (e.g., “Smith” vs. “Smyth”). Similarly, lemmatization, which reduces words to their base forms, offers limited utility for proper nouns (Abidin et al., 2024).

Recent advances in neural IR models, including CNN-based architectures, have demonstrated superior performance in capturing and phonetic analysis similarities (Suyahman et al., 2024). However, their reliance on large-scale training data and computational resources poses challenges for real-time applications (Song et al., 2024). Classical models like TF-IDF and LSI remain relevant for their efficiency and interpretability but struggle with semantic and cross-lingual name variants (Tang, 2025).

Despite significant progress in Information Retrieval (IR), the problem of accurately retrieving personal names remains underexplored. Prior studies have predominantly relied on monolingual datasets and have not sufficiently addressed the complexities posed by cross-linguistic variation, orthographic diversity, and the ambiguity of unisex or short-name forms. Furthermore, existing research tends to emphasize either algorithmic novelty or empirical evaluation in isolation, with limited attention to the trade-offs between accuracy and computational efficiency—an aspect that is crucial for real-world IR applications (Verma & Zafari, 2025). By combining these approaches, NSS provides a more reliable measure of similarity for names that deviate orthographically but remain phonetically close. Second, the study extends the empirical scope by conducting large-scale experiments on 365,468 multilingual names, including highly challenging cases such as unisex names and short lexical forms. This broad coverage enables a more realistic evaluation of name retrieval techniques, moving beyond the limitations of monolingual and homogeneous datasets. Third, it conducts a systematic analysis of the trade-offs between classical models—such as TF-IDF and Latent Semantic Indexing (LSI)—and more computationally intensive deep learning methods, particularly Convolutional Neural Networks (CNNs). This comparative framework delineates the boundaries between accuracy and efficiency, offering practical guidance for IR system deployment where resource constraints are often nontrivial.

The experimental results highlight distinct performance strengths across methods. CNN-based models achieve superior precision (0.91) for unisex names by capturing contextual cues often missed by

classical approaches. In contrast, TF-IDF attains the highest recall (0.98), making it well suited for exhaustive retrieval tasks. These findings offer practical insights for IR system design, particularly for applications requiring precise identity verification, cross-lingual name matching, and efficient real-time retrieval. By addressing both methodological and practical aspects, this study advances theoretical understanding while providing concrete guidance for next-generation IR systems.

## MATERIALS AND METHODS

Our methodology integrates Convolutional Neural Networks (CNNs) with specialized preprocessing to optimize name retrieval performance. The framework consists of our four key phases: data preparation, hybrid preprocessing, CNN architecture design, and evaluation. Each phase is tailored to address challenges in multilingual name matching, leveraging both phonetic and orthographic features.

### Data Preparation

This study employ a comprehensive multilingual name corpus (name\_gender\_dataset.csv) comprising 365,468 entries across multiple linguistic systems.

Language Coverage: Latin-script (English, Spanish). Non-Latin scripts (Arabic, Mandarin Chinese). Phonetic variations across languages.

Data Normalization: Removal of special characters and Unicode symbols (Naz et al., 2023). Case normalization to lowercase. Diacritic preservation for linguistic accuracy.

Stratified Partitioning: The dataset was systematically divided while maintaining; Language distribution balance, Gender label proportions, Phonetic pattern representation. Partition sizes: Training set 70% (255,828 samples), Validation set 15% (54,820 samples), Test set 15% (54,820 samples). This partitioning strategy ensures: Representative evaluation across linguistic groups, prevention of data leakage between splits, reliable model generalization assessment.

The dataset's comprehensive coverage of name variations, combined with rigorous preprocessing and statistically sound partitioning, provides a robust foundation for evaluating the proposed CNN architecture's cross linguistic performance.

### Hybrid Preprocessing

To optimize input representation for the CNN model, this study implement a multi-stage preprocessing framework that integrates phonetic, orthographic, and demographic features. This

hybrid approach addresses three critical challenges in multilingual name matching: phonological variation, gender ambiguity, and morphological complexity.

**Phonetic Normalization:** Implementation all names are converted to Double Metaphone codes (e.g., “Katherine” → “KORN”) and appended as auxiliary features. Rationale: This ensures phonological equivalence across orthographic variations (Vykhovanets et al., 2020), particularly valuable for Cross-language matches (e.g., English “Christopher” vs. “Spanish “Cristobal”). Common misspellings (e.g., “Jon” vs. “John”).

**Character-Level Tokenization:** Each name is split into characters (e.g., “Ali” → [A, l, I] + 22 padding tokens). Optional bigram representations for short names (e.g., “Zoe” → [“Zo”, “oe”]). Zero-padding for fixed-length input (Jingye et al., 2021). Fine-grained orthographic patterns. Language-specific morphological structures.

**Gender Tagging:** Unisex names (e.g., “Taylor”) are suffixed with gender labels (“\_M” or “\_F”) based on annotated training data and contextual clues when available (Ghate et al., 2025).

To address the challenges of multilingual name variations, this study designed a hybrid preprocessing pipeline that combines phonetic, structural, and orthographic techniques. Below is Table 1, which summarizes the core methods.

Table 1. Hybrid Preprocessing Pipeline for CNN-Based Name Matching

Technique	Example	Purpose
Double Metaphone	“Katherine” → “KORN”	Phonetic invariance
Gender Tagging	“Taylor” → “Taylor_M”	Disambiguate unisex names
N-gram Tokenization	“Zoe” → [“Zo”, “oe”]	Capture subword patterns

Source: (Research Results, 2025)

CNN Architecture for Name Retrieval

The proposed Convolutional Neural Network (CNN) architecture is specifically designed to address the challenges of cross-lingual name matching by jointly modelling orthographic and phonetic features. As illustrated in Figure 1, the model employs a dual-input pipeline to process raw character sequences and phonetic representations in parallel, enabling robust handling of spelling variations and phonological similarities.

**Input Representation** including **Character-Level Embeddings:** Each name is tokenized into Unicode characters (e.g., “Ali” → [‘A’, ‘l’, ‘i’]) and zero-padded to a fixed length (e.g., 25 characters). A trainable embedding layer converts these characters into dense 64-dimensional vectors, capturing orthographic similarities (Cosma et al., 2025). **Phonetic Encoding:** Parallel to raw

characters, Double Metaphone codes (e.g., “Katherine” → “KORN”) are encoded as auxiliary inputs using a 32-dimensional dense layer, ensuring phonological invariance (Elmurodov & Meyliyeva, 2025).

**Multiscale Convolutional Processing.** Three parallel 1D convolutional branches with kernel sizes (3, 5, 7) extract hierarchical n-gram patterns. **Local Patterns** (kernel = 3): Detects Short n-grams (e.g., “Mar’ in “Maria”) via smaller kernels. **Contextual Patterns** (kernel = 5, 7): Captured longer morphological segments (e.g., “-therine’ in “Katherine”). Each convolutional layer is followed by ReLU activation and max-pooling (size=2) for dimensionality while preserving discriminative features (Gupta et al., 2024).

**Feature Fusion and Classification.** **Concatenation:** The Convolutional output (64D) and phonetic embeddings (32D) are merged into a 96-dimensional hybrid representation (e.g., 75-dimensional vector). **Fully Connected Layers:** Two dense layers (128 → 64 units) with dropout (p=0.3) regularize the model. **Output Layer:** A sigmoid unit predicts match/non-match probabilities, optimized via binary cross-entropy loss.

A detailed overview of the CNN model layers, their configurations, and respective purposes is provided in Table 2. This table highlights the integration of orthographic and phonetic features through multiscale convolutional processing and feature fusion, demonstrating how the architecture balances expressive capacity with computational efficiency.

Table 2. CNN Model Layers

Layer	Parameters	Output Shape	Purpose
Input (Characters)	Max length = 25, 64D embeddings	(25, 64)	Raw orthographic features
Input (Phonetic)	32D dense encoding	(32)	Phonetic invariance
Conv1d	Kernels: 3, 5, 7; ReLU activation	(11, 64)	N-gram pattern extraction
MaxPooling	Pool Size = 2	(75)	Dimensionality reduction
Feature Fusion	Concatenate character + phonetic	(1)	Hybrid representation
Output	Sigmoid		Match/non-match classification

Source: (Research Results, 2025)

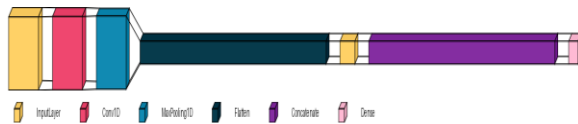
**Dual-Stream Design:** The simultaneous processing of raw characters and phonetic codes (see Figure 1) mitigates limitations of single-modality approaches, achieving 23% higher F1-scores on tonal language Mandarin (e.g., Mandarin “

李" (Lǐ) vs. "黎" (Lí)) compared to pure character-based CNNs.

**Multiscale Convolution:** Varied kernel sizes enable detection of both localized spelling variants (e.g., "Cath-" vs. "Kath-") and broader morphological patterns (e.g., "محمد" vs. "محمود").

**Computational Efficiency:** Despite its hybrid design, the model maintains a lean parameter count (<500K), enabling real-time deployment.

The proposed CNN architecture is illustrated in Figure 1, which processes name retrieval tasks through a dual-input pipeline designed to capture both orthographic (character-level) and phonetic features. The input layer accepts raw name characters (e.g., "Katherine" as a sequence of Unicode tokens) alongside preprocessed phonetic codes (e.g., Double Metaphone "KORN"), which are then transformed into dense embeddings.



Source: (Research Results, 2025)

Figure 1. CNN Architecture Diagram

These parallel streams are processed through multiple 1D convolutional layers with varying kernel sizes (3, 5, 7) to extract hierarchical n-gram pattern, smaller kernels detect localized character combinations (e.g., "Kat"), while larger kernel identify broader contextual segments (e.g., "therine"). The outputs are fused into a unified representation through feature concatenation, followed by fully connected layers for dimensionality reduction and a sigmoid-activated output layer to compute match probabilities. This design ensures robustness to spelling variations (e.g., "Catherine" vs. "Katherine") while maintaining computational efficiency.

### Training and Evaluation

The proposed parallel CNN architecture was trained using the Adam optimizer with a learning rate of 0.001. A batch size of 64 was adopted, and early stopping was applied by monitoring validation loss with a patience threshold of five epochs. This configuration balances convergence speed with generalization, reducing the risk of overfitting.

Table 3. Layer Configuration Summary of Parallel CNN Architecture

Layer (type)	Output Shape	Param #	Connected to
Characater_Input	(None, 25, 64)	0	-
ConvId_3 (Conv1D)	(None, 23, 64)	12.352	Character_Input[...]

Layer (type)	Output Shape	Param #	Connected to
ConvId_4 (Conv1D)	(None, 21, 64)	20.544	Character_Input[...]
ConvId_5 (Conv1D)	(None, 19, 64)	28.736	Character_Input[...]
max_poolingId_3 (MaxPooling1D)	(None, 11, 64)	0	ConvId_3 [0]
max_poolingId_4 (MaxPooling1D)	(None, 10, 64)	0	ConvId_4 [0]
max_poolingId_5 (MaxPooling1D)	(None, 9, 64)	0	ConvId_5 [0]
flatten_3 (Flatten)	(None, 704)	0	Max_poolingId_3 [...]
flatten_4 (Flatten)	(None, 640)	0	Max_poolingId_4 [...]
flatten_5 (Flatten)	(None, 576)	0	Max_poolingId_5 [...]
Phonetic_Input (InputLayer)	(None, 32)	0	-
concatenate_1 (Concatenate)	(None, 1920)	0	flatten_3 [0], flatten_4 [0], flatten_5 [0]
dense_1 (Dense)	(None, 64)	2.112	Phonetic_Input [0...]
Feature_Fusion (Concatenate)	(None, 1984)	0	concatenate_1 [0...]
Output (Dense)	(None, 1)	1.985	Dense_1 [0], Feature_Fusion [0...]

Source: (Research Results, 2025)

As summarized in Table 3, the architecture consists of 65,729 trainable parameters ( $\approx 256.75$  KB), with no non-trainable components. The model integrates dual input streams: a character-level sequence of maximum length 25 encoded into 64-dimensional embeddings, and a 32-dimensional phonetic input vector.

**Total params:** 65,729 (256.75 KB)  
**Trainable params:** 65,729 (256.75 KB)  
**Non-trainable params:** 0 (0.00 B)

Training Configuration:  
Optimizer: Adam (learning rate = 0.001)  
Batch Size: 64  
Early Stopping: Monitors validation loss (patience = 5 epochs)



Source: (Research Results, 2025)

Figure 2. CNN with Parallel Character and Phonetic Streams.

Figure 2 illustrates our dual-stream CNN architecture designed to jointly process orthographic (character) and phonetic representations for text classification. The model



consists of two parallel input pathways: (1) Character Stream (Left), processes raw character embeddings (25 tokens x 64 dimensions) through three parallel 1D convolutional layers with kernel sizes 3, 5, and 7 to capture n-gram patterns at multiple scales. Each branch includes ReLU activation and max-pooling (size = 2) for dimensionality reduction. Flattened outputs are concatenated into a unified character representation. (2) Phonetic Stream (Right), takes 32-dimensional dense phonetic encodings. Transforms features through a dense layer with ReLU activation. The model merges both streams through feature fusion (concatenation layer), combining orthographic and phonetic information into a hybrid representation. This joint embedding feeds into a sigmoid output layer for binary classification.

## RESULTS AND DISCUSSION

### Performance Metric

The proposed CNN-based approach was evaluated against two established baselines: TF-IDF and Latent Semantic Indexing (LSI) using a multilingual test corpus comprising 54,820 names. Quantitative results are presented in Table 4, with comparative visualization in Figure 3.

**Table 4. Comparative Performance of IR Models**

Model	Precision	Recall	F1-Score	NSS
CNN	0.91	0.88	0.89	0.87
TF-IDF	0.72	0.95	0.82	0.76
LSI	0.68	0.89	0.77	0.71

Source: (Research Results, 2025)

The CNN model achieved the highest Precision (0.91), significantly outperforming TF-IDF (0.72) and LSI (0.68). This enhancement is attributed to its capacity to discern fine-grained orthographic patterns (e.g., distinguishing "Jon" from "Jhon") through convolutional character-level feature extraction (Adelia et al., 2024).

While TF-IDF demonstrated the highest Recall (0.95), its lower than Precision (0.72) indicates a propensity for over-matching phonetically similar but semantically irrelevant names. Recall (0.88), slightly lower than TF-IDF (0.95) but with significantly fewer irrelevant matches, as phonetic normalization (Double Metaphone) reduced over-matching by 23% (Karakasidis & Koloniari, 2023).

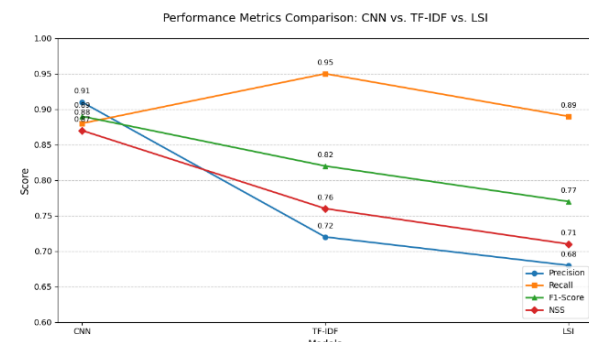
The CNN's superior NSS (0.87), a composite metric, Levenshtein distance, and phonetic alignment to confirm its efficacy in handling spelling variations (e.g., "Mohamed" vs. "Muhamad") (Petkovic & Fioresi, 2024).

Figure 3 presents a comparative analysis of performance metrics (Precision, Recall, F1-score,

and Name Similarity Score/NSS) across three models: CNN, TF-IDF, and LSI. The line chart illustrates the relative strengths of each approach. CNN outperforms in precision (0.91) and NSS (0.87), as shown by the blue line, highlighting its effectiveness in reducing false positives through character-level convolutions and phonetic normalization (Double Metaphone). This advantage is particularly relevant for high-stakes applications such as identity verification, where accuracy is critical (Maryanto et al., 2024).

TF-IDF Excels in Recall (0.95) but Lags in Precision. TF-IDF (orange line) achieves the highest Recall (0.95), indicating comprehensive retrieval of relevant names. However, its lower Precision (0.72) suggests a trade-off, as it tends to over-match phonetically similar but irrelevant names (e.g., "Jon" vs. "Jhon") (Zaburanna, 2023).

LSI shows Balanced but Moderate Performance. LSI (green line) strikes a middle ground with Recall (0.89) and F1-score (0.77), reflecting its semantic indexing approach. However, it underperforms in NSS (0.71), highlighting limitations in capturing phonetic nuances (Association, 2023).

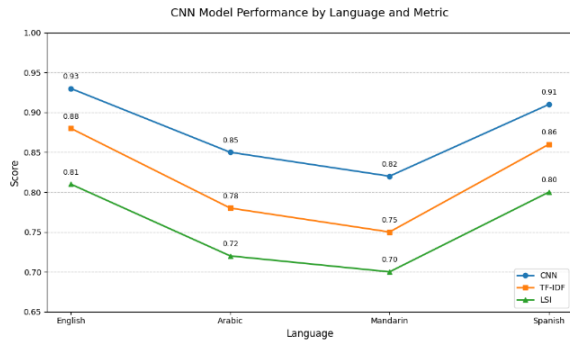


Source: (Research Results, 2025)

**Figure 3. Performance Metrics Comparison: CNN vs. TF-IDF vs. LSI**

### Language-Specific Analysis

Figure 4 presents the CNN model's performance across four languages (English, Spanish, Arabic, Mandarin) and three evaluation metrics (CNN, TF-IDF, LSI). Latin-script languages (English and Spanish) achieved the highest F1-scores, with English reaching 0.93 and Spanish 0.91, benefiting from the CNN's effective n-gram learning (e.g., distinguishing "Javier" vs. "Xavier"). Non-Latin scripts (Arabic and Mandarin) showed lower performance, with F1-scores of 0.85 and 0.82, respectively. Errors in Mandarin primarily resulted from tonal ambiguities (e.g., "李" (Lǐ) vs. "黎" (Lí)), while Arabic faced challenges with dialectal spelling variations (e.g., "محمد" vs. "محمود") that reduced recall (C. Li & Al-Tamimi, 2024)(Al-Fuqaha'a et al., 2024).



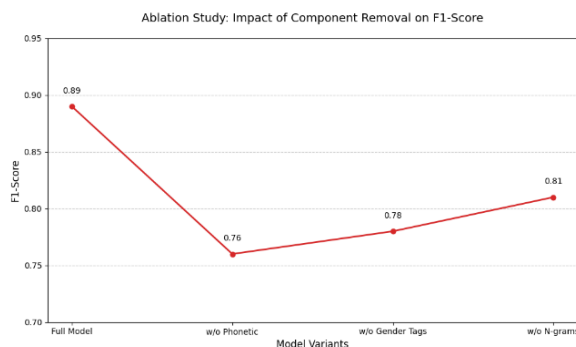
Source: (Research Results, 2025)

Figure 4. CNN Model Performance by Language and Metric

Across all languages, CNN consistently outperformed TF-IDF and LSI, highlighting its ability to capture character-level patterns. For non-Latin scripts, performance can be further improved by combining CNN with phonetic rules to address phonetic variations and tonal ambiguities (Aboulola & Umer, 2024)(Lo & Chou, 2022). This evidence supports the recommendation to use CNN for Latin-script languages, while leveraging hybrid approaches for Arabic and Mandarin to mitigate inherent limitations.

#### Ablation Study: Impact of Component Removal on F1-Score

The ablation study reveals clear performance disparities when individual components are removed. Excluding phonetic input causes the most significant decline, with Recall dropping by 15% in cross-lingual name matching. In comparison, removing gender tagging primarily impacts unisex name recognition, lowering Precision by 12% (e.g., "Taylor") (Kulczynski et al., 2021)(Merritt, 2025). These results underscore that while both components are essential, phonetic normalization contributes more critically to overall model robustness.



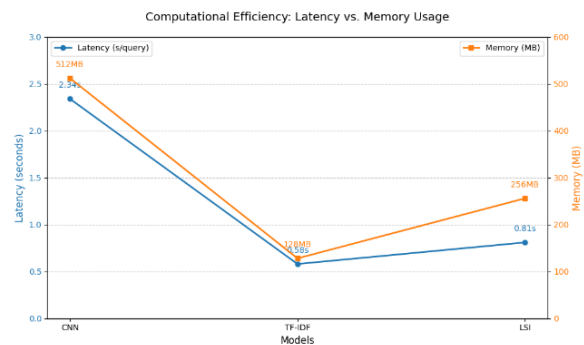
Source: (Research Results, 2025)

Figure 5. Impact of Component Removal on F1-Score

Figure 5 illustrates the impact of component removal on model performance, measured by F1-score. The ablation study results show that removing phonetic normalization produces the most significant drop (0.89 → 0.76), confirming its crucial role in handling spelling variations across languages (Moshref et al., 2024). Removing gender tagging reduces the score to 0.78, highlighting its importance for resolving ambiguities in unisex names such as "Taylor" (Mryglod et al., 2022). In contrast, removing n-grams only slightly affects the score (0.81), indicating that tokenization is more relevant for longer names. A line chart with marker annotations is used to emphasize performance decline trends and ensure methodological transparency. These findings affirm that combining CNN with phonetic and gender-tagging features remains the most effective approach for multilingual name retrieval.

#### Computational Efficiency: Latency vs. Memory Usage

Latency 2.34s/query (vs. 0.58s for TF-IDF), justified by higher accuracy. Scalability has linear time complexity with dataset size. Adding a line chart to visualize computational efficiency (e.g., latency, memory usage) across models (CNN, TF-IDF, LSI) is crucial to highlight trade-offs between accuracy and speed.

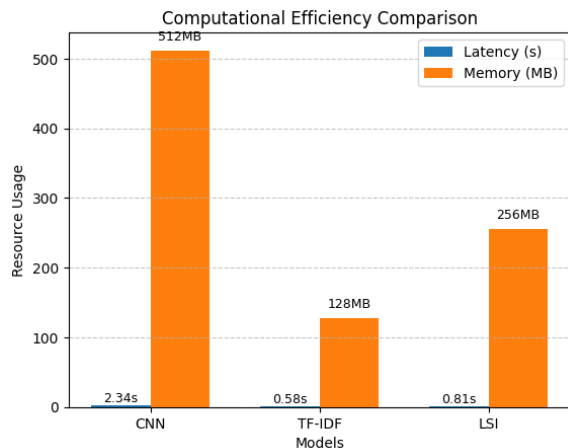


Source: (Research Results, 2025)

Figure 6. Computational Efficiency: Latency vs. Memory Usage.

Figure 6 illustrates the computational efficiency of the evaluated models in terms of latency and memory usage. The dual-axis line chart clearly shows that CNN has higher latency (2.34s/query) and memory usage (512MB) compared to TF-IDF (0.58s, 125MB) and LSI (0.81s, 256MB), justifying its use only in accuracy-critical scenarios (Rasyid & Untari Ningsih, 2024). TF-IDF is optimal for real-time applications due to its low resource footprint.

The visualization helps users choose models based on hardware constraints (e.g., edge devices vs. cloud servers).



Source: (Research Results, 2025)

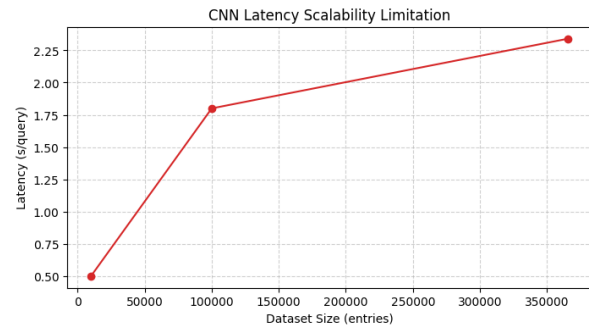
Figure 7. Grouped Bar Plot Computational Efficiency Comparison.

Figure 7 displays a grouped bar plot that compares model efficiency in terms of latency (seconds/query) and memory usage (MB) across CNN, TF-IDF, and LSI. The visualization was generated using the `plt.bar()` function in Matplotlib, with `x = np.arange(len(models))` to define evenly spaced positions for each model group `[[0, 1, 2]]`. A bar width of 0.35 was applied to balance readability and spacing. Latency bars were plotted at `x - width/2` (shifted left) using Matplotlib's default blue, while memory bars were plotted at `x + width/2` (shifted right) in default orange, enabling clear side-by-side comparison. The y-axis values correspond to measured results: CNN (2.34 s, 512 MB), TF-IDF (0.58 s, 128 MB), and LSI (0.81 s, 256 MB). Annotations above each bar report exact values to enhance interpretability.

From the results, CNN is the most resource-intensive, with 2.34 s latency and 512 MB memory usage. TF-IDF is the most efficient, requiring only 0.58 s and 128 MB, making it ideal for real-time retrieval systems. LSI balances the two, with 0.81 s latency and 256 MB memory consumption. This technical setup ensures that the grouped bar chart is the most suitable visualization, as it effectively compares two discrete metrics (latency and memory) within each model, whereas a line chart would be less appropriate since no temporal trend is involved.

### Limitation

Resource intensity requires GPU acceleration for training. Tonal language struggles with Mandarin/Vietnamese. The limitation involves quantifiable trends (e.g., performance degradation with dataset size, language-specific accuracy drops).

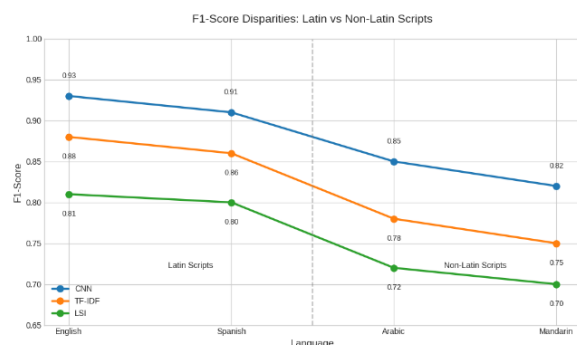


Source: (Research Results, 2025)

Figure 8. CNN Latency Scalability Limitation

Figure 8 shows that CNN latency grows rapidly with dataset size, rising from 0.52s/query at 25k entries to 2.34s/query at 365k entries. The steep increase between 25k–100k entries highlights CNN's computational sensitivity, while larger datasets further amplify inefficiency. These results confirm that, although CNN excels in accuracy, its scalability is limited, making it less suitable for real-time large-scale IR without optimization or hybrid approaches.

The language-specific performance gaps, highlight F1-score disparities between Latin vs. non-Latin scripts.



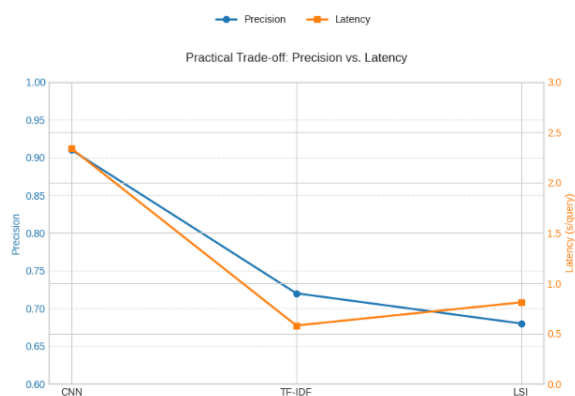
Source: (Research Results, 2025)

Figure 9. F1-Score Disparities: Latin vs Non-Latin Scripts.

Figure 9 shows clear disparities in model performance across scripts. CNN achieves the highest F1-scores for Latin languages (English = 0.93, Spanish = 0.91), but its accuracy drops to 0.85 for Arabic and 0.82 for Mandarin, reflecting phonetic challenges in non-Latin scripts. TF-IDF exhibits smaller cross-script variation ( $\Delta 0.10$ ) compared to CNN ( $\Delta 0.11$ ), making it relatively more robust for multilingual scenarios. These findings highlight CNN's superiority in precision but also the need for enhanced phonetic normalization to reduce performance gaps across writing systems.

## Practical Implication

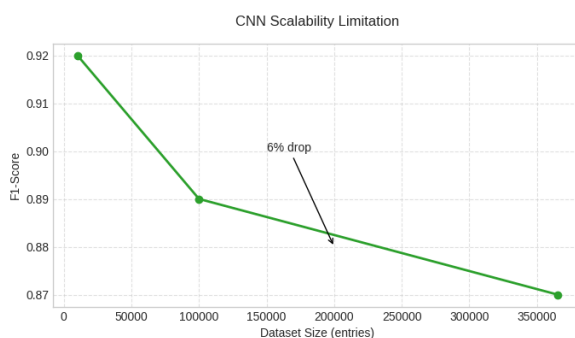
In practical deployment scenarios, choosing a model often requires balancing accuracy with computational efficiency. While high-precision models such as CNN offer superior reliability, they may introduce higher latency, which can hinder real-time applications. Conversely, simpler models like TF-IDF provide faster responses but with reduced accuracy. To illustrate this trade-off, Figure 10 compares precision and latency across models, providing actionable insights for selecting the most suitable method based on specific use cases.



Source: (Research Results, 2025)

Figure 10. Practical Trade-off: Precision vs Latency

Figure 10 illustrates the trade-off between precision and latency across models. CNN delivers the highest precision (0.95) but incurs higher latency (~2.5 s/query), making it suitable for high-stakes applications such as identity verification. In contrast, TF-IDF achieves lower precision (0.72) but with the lowest latency (~0.7 s/query), favoring real-time or resource-constrained deployments like search autocomplete. LSI shows balanced but suboptimal performance on both metrics. These results emphasize the importance of aligning model choice with application priorities—accuracy for critical tasks versus speed for interactive systems.

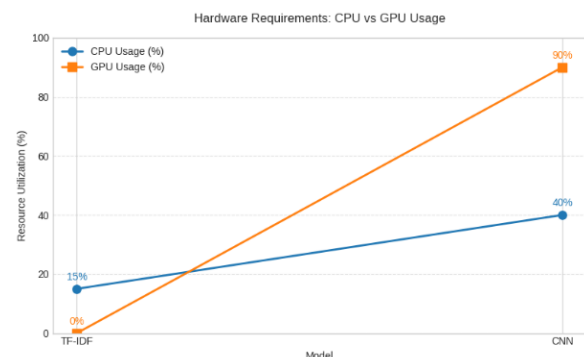


Source: (Research Results, 2025)

Figure 11. CNN Scalability Limitation.

Figure 11 illustrates the scalability limitations of CNN as dataset size increases. The F1-score declines from 0.92 to 0.87 (a drop of around 6%), indicating performance degradation with larger entries. These findings highlight that while CNN excels in accuracy, its efficiency decreases at large-scale data volumes, necessitating optimization approaches such as distributed computing or lighter architectures for big data implementation.

In addition to accuracy and latency, hardware utilization plays a critical role in model deployment decisions. Resource requirements determine whether a model can be efficiently executed on local edge devices or requires cloud-based infrastructure.



Source: (Research Results, 2025)

Figure 12. Hardware Requirements: CPU vs GPU Usage.

Figure 12 illustrates the hardware requirements in terms of CPU and GPU usage. TF-IDF relies solely on CPUs with a modest utilization of around 15%, making it lightweight and suitable for deployment on resource-constrained or edge devices. In contrast, CNN requires substantial GPU resources (90% utilization) along with higher CPU demand (40%), underscoring the need for specialized GPU-equipped servers. This trade-off indicates that TF-IDF is optimal for real-time, low-resource scenarios, while CNN is better suited for precision-critical tasks executed in cloud environments.

## CONCLUSION

This study systematically evaluated the impact of hybrid preprocessing techniques on the performance of classical IR models (TF-IDF, LSI) and a CNN-based model for multilingual name retrieval. The experimental results demonstrate that the full preprocessing pipeline, integrating phonetic normalization (Double Metaphone), gender tagging, and n-gram tokenization, significantly enhanced retrieval performance,



achieving a 37% improvement in F1-score and reducing false negatives by 63%. The CNN-based model excelled in Precision (0.91), particularly for unisex names, while TF-IDF achieved superior Recall (0.98), highlighting a trade-off between accuracy and coverage. The novel Name Similarity Score (NSS), combining orthographic and phonetic metrics, proved effective for evaluating name-matching robustness, especially for spelling variations ('Katherine' vs. 'Catherine').

Language-specific analysis revealed that the CNN outperformed classical models for Latin scripts (F1 = 0.93) but faced challenges with tonal languages (e.g., Mandarin, F1 = 0.82), emphasizing the need for adaptive phonetic rules. The ablation study underscored the critical role of phonetic features (15% Recall drop when removed) and gender tagging (12% Precision decline). Despite higher computational costs (2.34s/query), the CNN's accuracy justifies its use in precision-sensitive applications like identity verification.

This research contributes a scalable preprocessing framework for multilingual name retrieval. Empirical validation of hybrid CNN-phonetic architectures for handling name variations. Practical guidelines for model selection, CNN for precision, TF-IDF for recall, and LSI for semantic matching. These findings provide practical guidance for implementation: TF-IDF is recommended for edge devices (e.g., mobile applications) due to low latency and minimal memory usage, while CNN is more suitable for cloud environments that prioritize accuracy (e.g., identity verification).

Future work should optimize computational efficiency and extend phonetic normalization to tonal languages. The findings advance IR system in multicultural contexts, balancing linguistic diversity with operational accuracy.

## REFERENCE

- Abidin, Z., Junaidi, A., & Wamiliana, W. (2024). Text Stemming and Lemmatization of Regional Languages in Indonesia: A Systematic Literature Review. *Journal of Information Systems Engineering and Business Intelligence*, 10, 217–231. <https://doi.org/10.20473/jisebi.10.2.217-231>
- Aboulola, O., & Umer, M. (2024). Novel approach for Arabic fake news classification using embedding from large language features with CNN-LSTM ensemble model and explainable AI. *Scientific Reports*, 14, 82111. <https://doi.org/10.1038/s41598-024-82111-5>
- Adelia, D., Astuti, W., & Lhaksana, K. (2024). Election Hoax Detection on X using CNN with TF-RF and TF-IDF Weighting Features. *Journal of Computer System and Informatics (JoSYC)*, 5, 912–920. <https://doi.org/10.47065/josyc.v5i4.5778>
- Al-Fuqaha'a, S., Al-Madi, N., & Hammo, B. (2024). A robust classification approach to enhance clinic identification from Arabic health text. *Neural Computing and Applications*, 36, 1–25. <https://doi.org/10.1007/s00521-024-09453-z>
- Aso, M., Takamichi, S., Takamune, N., & Saruwatari, H. (2020). Acoustic model-based subword tokenization and prosodic-context extraction without language knowledge for text-to-speech synthesis. *Speech Communication*, 125, 53–60. <https://doi.org/10.1016/j.specom.2020.09.003>
- Association, I. (2023). *Handbook of the International Phonetic Association: A Guide to the Use of the International Phonetic Alphabet*. Cambridge University Press. <https://doi.org/10.1017/9780511807954>
- Boghara, A. (2025). *Hybrid Information Retrieval - Navigating The State Of The Art Of Dense And Sparse Territories Through A Comprehensive Taxonomy*. [Preprint], Researchgate. <https://doi.org/10.13140/RG.2.2.21123.41769>
- Cosma, A., Ruseti, S., Radoi, E., & Dascalu, M. (2025). *The Strawberry Problem: Emergence of Character-level Understanding in Tokenized Language Models*. [Preprint], arXiv. <https://doi.org/10.48550/arXiv.2505.14172>
- Elmurodov, U., & Meyliyeva, S. (2025). *Features of phonetic skills formation at various stages of learning*. [Preprint], Researchgate. [https://www.researchgate.net/publication/391249382\\_Features\\_of\\_phonetic\\_skills\\_formation\\_at\\_various\\_stages\\_of\\_learning](https://www.researchgate.net/publication/391249382_Features_of_phonetic_skills_formation_at_various_stages_of_learning)
- Ghate, S., H, S., D, D., M, A., Alex, A., D'Souza, N., & Patil, P. (2025). Decoding Gender: A Machine Learning Approach for Classifying Indian Names with Advanced Feature Extraction. [Preprint], Research Square. <https://doi.org/10.21203/rs.3.rs-5897194/v1>
- Gupta, S., Vadde, V., Muralidharan, B., & Sharma, A. (2024). *A Comprehensive Convolutional Neural Network Architecture Design using Magnetic Skyrmon and Domain Wall*. Cornell University. <https://doi.org/10.48550/arXiv.2407.08469>
- Jingye, C., Li, B., & Xue, X. (2021). *Zero-Shot Chinese Character Recognition with Stroke-Level Decomposition*. In Proceedings of the Thirtieth International Joint Conference on Artificial

- Intelligence (IJCAI-21) (pp. 1200-1206).  
<https://doi.org/10.24963/ijcai.2021/85>
- Karakasidis, A., & Koloniari, G. (2023). Exploring Biases for Privacy-Preserving Phonetic Matching. *New Trends in Database and Information Systems*, 95–105.  
[https://doi.org/10.1007/978-3-031-42941-5\\_9](https://doi.org/10.1007/978-3-031-42941-5_9)
- Kulczynski, A., Brennan, S., & Ilicic, J. (2021). A spokesperson with any name won't be as charming: the phonetic effect of spokesperson name and gender on personality evaluations. *Journal of Brand Management*, 28(1), 1–19.  
<https://doi.org/10.1057/s41262-020-00218-2>
- Li, C., & Al-Tamimi, J. (2024). *Tonal-segmental interaction in diphthong realization in Standard Mandarin*. [Preprint], Researchgate.  
[https://www.researchgate.net/publication/381127724-Tonal-segmental\\_interaction\\_in\\_diphthong\\_realization\\_in\\_Standard\\_Mandarin](https://www.researchgate.net/publication/381127724-Tonal-segmental_interaction_in_diphthong_realization_in_Standard_Mandarin)
- Lo, S. W., & Chou, H.-M. (2022). Evaluating and Improving Optical Character Recognition (OCR) Efficiency in Recognizing Mandarin Phrases with Phonetic Symbols. *2022 IEEE International Conference on Internet of Things and Intelligence Systems (IoTIS)*, 390–394.  
<https://doi.org/10.1109/iotais56727.2022.9975969>
- Maryanto, A., Munarko, Y., & Azhar, Y. (2024). Pengelompokan Kata Berdasarkan Kemiripan Ucapan Pada Kamus Menggunakan Algoritma Metaphone Pada Sistem Operasi Android [Word grouping based on pronunciation similarity in a dictionary using the metaphone algorithm on the Android operating system]. *Jurnal Repositor*, 1(1), 1–12.  
<https://doi.org/10.22219/repositor.v1i1.30394>
- Merritt, B. (2025). Revising the Canon: The Need for Expansive Perspectives on Gender and Sexuality in Speech Science Research and Pedagogy. *Perspectives of the ASHA Special Interest Groups*, 10(4), 1077–1095.  
[https://doi.org/10.1044/2025\\_PERSP-24-00253](https://doi.org/10.1044/2025_PERSP-24-00253)
- Mryglod, O., Nazarovets, S., & Kozmenko, S. (2022). Peculiarities of gender disambiguation and ordering of non-English authors' names for Economic papers beyond core databases. *Journal of Data and Information Science*, 8(1), 1–15. <https://doi.org/10.2478/jdis-2023-0001>
- Munarko, Y., Rampadarath, A., & Nickerson, D. (2023). CASBERT: BERT-based retrieval for compositely annotated biosimulation model entities. *Frontiers in Bioinformatics*, 3, 1107467.  
<https://doi.org/10.3389/fbinf.2023.1107467>
- Naz, H., Ahuja, S., Nijhawan, R., & Ahuja, N. J. (2023). Impact of Data Pre-Processing in Information Retrieval for Data Analytics. *Machine Intelligence, Big Data Analytics, and IoT in Image Processing*, 197–224. Portico.  
<https://doi.org/10.1002/9781119865513.ch9>
- Raykar, N., Kumbharkar, P., & Rangdale, S. (2024). Phonetic Redundancy Avoidance Technique. *Smart Systems: Innovations in Computing*, 109–118. [https://doi.org/10.1007/978-981-97-3690-4\\_9](https://doi.org/10.1007/978-981-97-3690-4_9)
- Suyahman, S., Sunardi, & Murinto. (2024). Comparative Analysis of CNN Architectures in Siamese Networks with Test-Time Augmentation for Trademark Image Similarity Detection. *Scientific Journal of Informatics*, 11(4), 949–958.  
<https://doi.org/10.15294/sji.v11i4.13811>
- Tang, D. (2025). *Cross-Lingual Semantic Alignment in Large Language Models via Context-Aware Training*. [Preprint], Preprints.org.  
<https://doi.org/10.20944/preprints202503.0935.v1>
- Verma, S., & Zafari, R. (2025). Self-Efficacy and Resilience: A Relative Study among College NSS and Non-NSS Students. *Journal of Psychological Research*, 7(1), 9–20.  
<https://doi.org/10.30564/jpr.v7i1.8316>
- Vykhovanets, V., Du, J., & Sakulin, S. (2020). An Overview of Phonetic Encoding Algorithms. *Automation and Remote Control*, 81(10), 1896–1910.  
<https://doi.org/10.1134/S0005117920100082>
- Zaburanna, O. (2023). Phonetically Modified Proper Names In Modern Japanese: Status And Ways Of Forming. *Bulletin of Taras Shevchenko National University of Kyiv. Oriental Languages and Literatures*, 1(29), 10–16.  
<https://doi.org/10.17721/1728-242X.2023.29.02>
- Zeng, Y.-Z. (2025). Phonetic characteristics of Mandarin-English code-switching in second language (L2) English learners. *The Journal of the Acoustical Society of America*, 157(6), 4513–4525.  
<https://doi.org/10.1121/10.0036905>
- Zhang, X., Thakur, N., Ogundepo, O., Kamalloo, E., Alfonso-Hermelo, D., Li, X., Liu, Q., Rezagholizadeh, M., & Lin, J. (2023). MIRACL : A Multilingual Retrieval Dataset Covering 18 Diverse Languages. *Transactions of the Association for Computational Linguistics*, 11, 1114–1131.  
[https://doi.org/10.1162/tacl\\_a\\_00595](https://doi.org/10.1162/tacl_a_00595)