# Natural Language Processing: Trends in Cross-Language Information Retrieval

Omodunbi B.A[1], Okomba N.S[2], Avah S.E[3], NwokoyeC.S[4],
Aniedu A.N[5], Okafor C.S[6], Ezeasor Ekene[7]

*(Computer Engineering, Federal University Oye-Ekiti)*
*(Computer Engineering, Federal University Oye-Ekiti),*
*(Computer Engineering, Federal University Oye-Ekiti),*
*(Electronics and Computer Engineering, Nnamdi Azikiwe University, Awka),*
*(Electronics and Computer Engineering, Nnamdi Azikiwe University, Awka),*
*(Electronics and Computer Engineering, Nnamdi Azikiwe University, Awka)*
*(Computer Science, Nnamdi Azikiwe University, Awka.)*

**ABSTRACT-** Cross-Language Information Retrieval (CLIR) is a pivotal area within Natural Language Processing (NLP), enabling users to access information across various languages. This paper explores the latest trends in CLIR, focusing on advancements in machine translation, multilingual embedding, and user-centric design. By analyzing existing literature and employing a qualitative approach, this study identifies key challenges and opportunities in the field. The findings indicate that while significant progress has been made, issues related to language diversity, cultural context, and user experience continue to pose challenges. Recommendations for future research directions are provided to enhance the effectiveness of CLIR systems.
**Keywords**-CLIR, machine translations, multilingual embedding, nlp, user-centric design

## I. INTRODUCTION

In an increasingly interconnected and globalized world, the capacity to retrieve information across multiple languages has become essential for effective communication, collaboration, and knowledge sharing. Cross-Language Information Retrieval (CLIR) is a critical tool in this regard, empowering users to search for documents in one language and seamlessly obtain relevant results in another. This capability is especially vital as the volume of multilingual content available on the internet continues to expand at an unprecedented rate, creating both opportunities and challenges for users and researchers alike.

As the significance of effective CLIR systems becomes more pronounced, it is crucial to explore the technological advancements that are driving this field forward. These innovations not only enhance the efficiency and accuracy of information retrieval but also improve user experience by making information accessible regardless of language barriers.

## II. LITERATUREREVIEW

Cross-Language Information Retrieval (CLIR) has gained prominence as a critical field in information retrieval, particularly with the globalization of data. The concept began to take shape in the early 1990s, but it was the work of Oard and Dorr in 1996 that laid a significant foundation with their paper on multilingual information retrieval systems. They emphasized the necessity for systems capable of retrieving information across different languages, which sparked further research and development in this area.

Initially, CLIR systems relied heavily on straightforward translation techniques, often leading to suboptimal results due to the complexities of language semantics. A notable advancement occurred in 2003 when the research by Steinberger and Jiri M. K. Karpov introduced a statistical approach utilizing parallel corpora to enhance translation accuracy. This method demonstrated improved retrieval effectiveness by leveraging large bilingual datasets, allowing for better contextual understanding.

In 2007, Huang et al. proposed a language model-based framework for CLIR that marked a turning point in the field. Their approach integrated probabilistic models to address the inherent uncertainties in language translation, enabling more

precise retrieval of relevant documents. This model not only improved performance but also highlighted the importance of linguistic nuances in cross-language contexts.

Despite these advancements, challenges remained, particularly in evaluating the performance of CLIR systems. In 2010, Lin and Wilbur underscored the need for tailored evaluation metrics that reflect the unique aspects of multilingual information retrieval. They argued that traditional metrics, such as precision and recall, were inadequate for assessing CLIR effectiveness due to the intricacies involved in language processing.

The advent of deep learning in the 2010s brought about transformative changes in CLIR methodologies. The work of Johnson et al. in 2017 illustrated the potential of neural networks to enhance machine translation quality, which subsequently improved CLIR outcomes. This shift towards neural architectures allowed for a more sophisticated understanding of language, facilitating better retrieval strategies.

Recent studies have also focused on the user experience within CLIR systems. Research by Kuo and Wu in 2020 highlighted the significance of designing interfaces that cater to diverse linguistic backgrounds. They proposed that enhancing user interaction could lead to higher satisfaction and improved retrieval success rates, emphasizing the role of usability in system design.

## III. METHODS AND MATERIALS

The Cross-Language Information Retrieval (CLIR) system comprises four key subsystems which are query translation, document retrieval, result ranking and finally result presentation. This is shown in Fig 1.
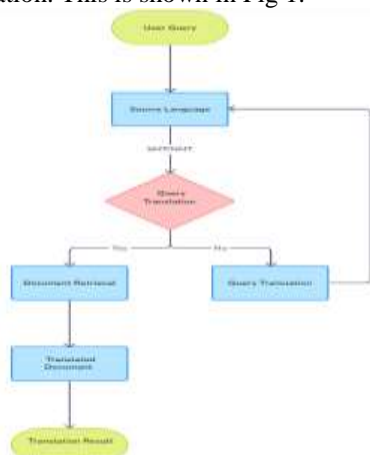


**Fig1: Block Diagram of a Cross-Language Information Retrieval System**

### A. Query Translation

Query translation is a critical component of the Cross-Language Information Retrieval (CLIR) system, where a user's query in the source language is transformed into the target language to facilitate effective document retrieval. This process typically involves several key steps and methodologies. The translation process begins with the input query, which is analyzed to identify its intent and context. This analysis is crucial because it informs the translation model about the semantic meaning of the query. The model is trained on labeled datasets that consist of pairs of queries in both the source and target languages. These datasets usually include a training set, a validation set, and a test set, often divided in a ratio of 60%, 20%, and 20%, respectively. The training set is used to teach the model, while the validation set helps fine-tune the model's parameters, and the test set evaluates its performance. Feature extraction plays a significant role in this stage, where raw query data is transformed into meaningful features that can be effectively processed by the model. This process, known as feature engineering, ensures that the model captures relevant linguistic characteristics necessary. Once the model is trained, it undergoes evaluation to determine its effectiveness. If the model exhibits under-fitting, adjustments are made using the validation dataset, often involving hyper-parameter tuning. This tuning process continues until the model achieves satisfactory performance in translating queries.

**Query Translation Algorithm**
```
Start;
Input Query;
Training Data; Validation Data; Test Data;
Translated Query = Translate(Input Query);
Compare(Translated Query, Validation Data);
if (Match Found) {
    printf("Query translation successful");
} else {
    printf("Query translation failed");
}
End;
```
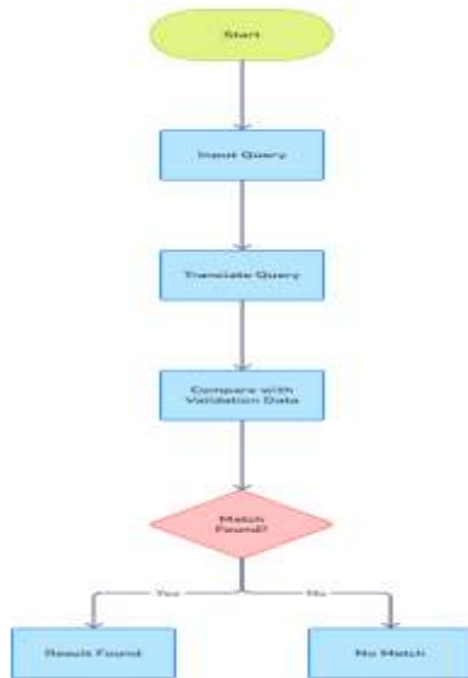The above algorithm can be expressed using a flowchart as seen in Fig 2.

Figure2:Flowchartfor query translation

### B. Document Retrieval

Document retrieval refers to the process of locating documents that meet specific criteria defined by a user's query. This process involves several key components:

1. User Query: The user query is the initial input that expresses the information needs of the user. It serves as the foundation for the entire retrieval process. Understanding the intricacies of user queries is crucial for effective information retrieval.

Components of User Query:

Query Structure: User queries can be structured or unstructured. Structured queries often follow a specific syntax (like SQL), while unstructured queries may consist of natural language phrases.

Query Types:

Keyword Queries: Users input specific keywords or phrases. For instance, "machine learning applications."

Boolean Queries: These utilize logical operators (AND, OR, NOT) to combine keywords, allowing for more complex searches. For example, "machine learning AND healthcare NOT finance."

Natural Language Queries: Users pose questions in everyday language, like "What are the applications of machine learning in healthcare?"

Query Transformation: Before processing, queries often undergo transformations to enhance retrieval effectiveness:

Tokenization: Splitting the query into individual terms or tokens.

Stemming and Lemmatization: Reducing words to their base or root form (e.g., "cooking" to "cook").

Stopword Removal: Filtering out common words (e.g., "the," "is") that may not contribute to the meaning.

Query Expansion: This technique involves adding synonyms or related terms to the original query to improve recall. For example, expanding "car" to include "automobile" and "vehicle."

Indexing: To enhance retrieval speed and efficiency, documents are often indexed. Indexing involves creating a data structure that maps terms to their locations in the document collection.

Inverted Index: A common indexing structure that lists each unique term and the documents in which it appears, along with term frequency and position information.

2. Ranking Algorithms: Documents are ranked based on their relevance scores, which are calculated using various algorithms. Common algorithms include:

Term Frequency-Inverse Document Frequency (TF-IDF): TF-IDF is a statistical measure that evaluates the importance of a word in a document relative to a corpus. It combines two components:

Term Frequency (TF): Measures how frequently a term appears in a document.

Inverse Document Frequency (IDF): Measures how important a term is across the entire corpus.

$\text{TF-IDF}(t,w) = \text{TF}(t,w) \times \text{IDF}(t)$

t: term

w: document

Mathematically

$\text{TF-IDF}(t,w)$

$$\frac{\text{Number of times term } t \text{ appears in document } w}{\text{Total number of terms in document } w}$$

Best Matching 25 (BM25): An advanced ranking function that considers term frequency, document length, and the number of documents containing the term.

The BM25 score for a document M with respect to a query K is given by:

$$\text{BM25}(M,K) = \sum_{t \in Q} IDF(t) \frac{T.F(t,M).(k-1)}{T.F(t,M) + k1.(1-b+b\frac{W}{avg\_len})}$$

Where:

K1 and b are parameters that control term frequency saturation and document length normalization, respectively.

|W| : Length of document

avg_len: Average document length in the collection

**Learning to Rank (LTR)**

Learning to Rank uses machine learning techniques to optimize ranking based on labeled training data. It involves training a model that learns to rank documents by analyzing features extracted from both the documents and queries

Approaches:

Pointwise: Predicts the relevance score of a single document.

Pairwise: Compares pairs of documents to determine which is more relevant.

Listwise: Considers the entire list of documents to optimize the ranking.

I. Pointwise Approach

The pointwise method treats the ranking problem as a regression or classification task, predicting relevance score for each document based on its features.

Input Features: Each document $w_i$ is represented by a feature vector $x_i$ where $x_i$ - $[x_{i1}, x_{i2}, \ldots x_{im}]$ and m is the number of features.

Model Function: The model can be represented as: $f(x_i) - \hat{y}_i$, where $\hat{y}_i$ is the predicted relevance score for document $w_i$

Loss Function:

For regression tasks, the Mean Squared Error (MSE) is often used:

$$\text{MSE} - f(x) = \frac{1}{N} \sum_{i=1}^{N} (yi - \hat{y}i)^2$$

Where $y_i$ is the true relevance score and N is the number of documents.

I. Pairwise Approach

The pairwise method compares pairs of documents to determine which is more relevant to a given query.

Input Features: For a pair of documents $(w_i, w_j)$, the feature vector can be represented as:

$x_{ij} - [f_1(w_i, w_j), f_2(w_i, w_j), \ldots f_k(w_i, w_j)]$ where each $f_k$ is a feature that captures the relationship between the two documents.

Model Function: The model predicts the preference $p_{ij}$:

$p_{ij} -$ sigmoid($f(x_{ij})$)

where $p_{ij}$ indicates the probability that document $w_i$ is preferred over document $w_j$

Loss Function:

A common loss function is the logistic loss L: $-\frac{1}{N} \sum_{(ij)£P}[yijlog(pij) + (1 - yij)log(1 - pij)]$

where P is the set of document pairs and $y_{ij}$ indicates whether $w_i$ is more relevant than $w_j$

II. Listwise Approach

The listwise method evaluates the entire list of documents for a query, optimizing the ranking as a whole.

Input Features: The entire list of documents W – $[w_1, w_2, \ldots w_N]$ is represented in a feature matrix X, where each row corresponds to a document.

Model Function: The model predicts a permutation of documents based on their relevance scores:

$\hat{y} - f(X)$

where $\hat{y}$ is a vector of predicted scores for the documents.

Loss Function:

A common approach is to use Softmax cross-entropy loss to optimize the ranking:

L: $-\frac{1}{N}\sum_{i=1}^{N} yijlog(pij)$ where K is the number of documents in the list, $yij$ is the true relevance label, and pij is the predicted probability of document j being ranked at position i.

The model is trained on a dataset containing queries and their corresponding relevant documents. Performance is evaluated using metrics like Mean Average Precision (MAP) or Normalized Discounted Cumulative Gain (NDCG).

Types of Models:

1. Vector Space Model (VSM)
2. Neural Ranking Models (NRM)

1. Vector Space Model (VSM)

The Vector Space Model is a mathematical model used for representing text documents as vectors in a multi-dimensional space. This model facilitates the computation of document similarity and relevance based on the geometric properties of the vectors.

Document Representation: Each document W is represented as a vector w in an n-dimensional space, where n is the number of unique terms in the document collection. The vector components are typically term weights, calculated using methods such as Term Frequency-Inverse Document Frequency (TF-IDF):

TF-IDF(t,W)=TF(t,W) * $\log \frac{N}{WF(t)}$)

Where:

TF(t,W) is the term frequency of term t in document W.

N is the total number of documents.

WF(t) is the number of documents containing term t.

Similarity Measurement: The similarity between

two documents W1 and W2 can be computed using cosine similarity:
Cosine similarity(W1,W2):
(W1.W2)/(|(|w1|)|||w2||)

Where · denotes the dot product and ||w|| denotes the Euclidean norm of vector w.

Neural Ranking Models (NRM)

Neural Ranking Models leverage deep learning techniques to learn complex representations of documents and queries, optimizing the ranking of documents based on their relevance to a given query.

Input Representation: Documents W and queries Q are transformed into dense vector embeddings using neural networks.
Let w denote the document vector and q denote the query vector, both of which are derived from embedding layers:
w – fembed (W), q - fembed (Q)
where fembed is a function representing the embedding process (e.g., using Word2Vec, BERT).
Ranking Function: The relevance score m between a document and a query can be modeled using a neural network:
m(W,Q) – frank (d,q)
where $f$ rank is a neural network that takes the document and query embeddings as input and outputs a relevance score.
Loss Function: A common loss function used in NRMs is the pairwise ranking loss, such as the hinge loss:
L - ∑(Wi,Qj)∈P max(0,1 – (r(Wi,Qj) - r(Wk,Qj))
Where $W_i,Q_j$ is a positive document-query pair, and $W_k$ is a negative document for the same query.
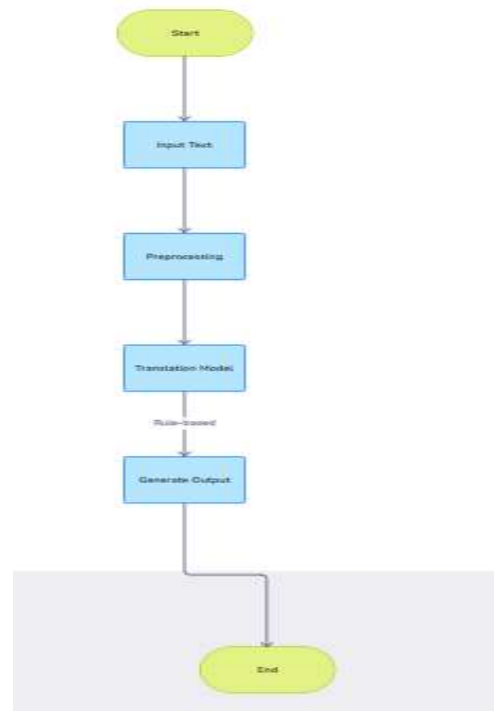


Figure3:FlowchartofTranslation Results

## IV. PERFORMANCE METRICS

For the CLIR, the following performance metrics was adopted:

Precision

Precision measures the proportion of relevant documents retrieved out of all the documents retrieved.
High precision indicates that the system returns a high number of relevant documents compared to irrelevant ones.

$$\text{Precision} = \frac{Number\ of\ Relevant\ Documents\ Retrieved}{Total\ Number\ of\ Documents\ Retrieved}$$

Recall

Recall measures the proportion of relevant documents retrieved out of all relevant documents available in the dataset. Given a sample set of 40 queries, the confusion matrices are displayed in table 1 and table 2.

$$\text{Recall} = \frac{Number\ of\ Relevant\ Documents\ Retrieved}{Total\ Number\ of\ Relevant\ Documents}$$

Table1:Confusion Matrix for Query Translation

|  |  | Actual |  |
|---|---|---|---|
|  |  | Relevant | Not Relevant |
| Predicted | Relevant | 15 | 10 |
|  | Not Relevant | 5 | 10 |

Precision for Query Translation:

Precision $= \frac{TP}{TP+FP} = \frac{15}{15+10} = 0.60$

Recall for Query Translation:

Recall $= \frac{TP}{TP+FN} = \frac{15}{15+5} = 0.75$

Table2:Confusion Matrix for Document Retrieval

| | | Actual | |
|---|---|---|---|
| | | Relevant | Not Relevant |
| Predicted | Relevant | 20 | 8 |
| | Not Relevant | 3 | 9 |

Precision for Document Retrieval:

Precision $= \frac{TP}{TP+FP} = \frac{20}{20+8} = 0.71$

Recall for Document Retrieval:

Recall $= \frac{TP}{TP+FN} = \frac{20}{20+3} = 0.75$

F1 Score

The F1 Score is the harmonic mean of precision and recall, providing a balance between the two metrics.

F1 Score for Query Translation:

F1 Score $= 2 * \frac{Precision \ * Recall}{Precision \ + Recall} = \frac{0.71 * 0.75}{0.71 + 0.75} = 0.36$

F1 Score for Document Retrieval:

F1 Score $= 2 * \frac{Precision \ * Recall}{Precision \ + Recall} = \frac{0.60 * 0.75}{0.60 + 0.75} = 0.33$

**Mean Average Precision (MAP)**
MAP is the mean of the average precision scores for multiple queries. It accounts for the rank of retrieved documents.

MAP $= \frac{1}{Q} \sum_{q=1}^{Q} AP(q)$

Where Q is the total number of queries and AP(q) is the average precision for query q.
MAP provides a single score that summarizes the performance across multiple queries.

**Normalized Discounted Cumulative Gain (NDCG)**
NDCG evaluates the ranking of relevant documents, giving higher scores to relevant documents that appear earlier in the result list.

NDCG $= \frac{DCG}{IDCG}$

Where:

DCG $= \sum_{i=1}^{P} \frac{rel \ i}{\log 2 \ (i+1)}$

IDCG is the ideal DCG, calculated using the best possible ranking of documents.
NDCG ranges from 0 to 1, where 1 indicates perfect ranking.

Mean Reciprocal Rank (MRR)
MRR is used when evaluating systems that return a list of results, focusing on the rank of the first relevant document.

MRR $= 1/Q \sum_{q=1}^{Q} ⟦ 1/rank_q ⟧$

Where rankq is the rank position of the first relevant document for query, q.

Higher MRR values indicate that relevant documents appear earlier in the result lists.
The model demonstrates average precision when evaluating the sample queries, primarily due to the limited number of queries tested. As the number of queries increases, the precision is expected to improve. The model shows some challenges in retrieving relevant documents, influenced by factors such as language nuances, context differences, and translation inaccuracies.

## V. CONCLUSION
In recent years, Natural Language Processing (NLP) has significantly advanced the field of Cross-Language Information Retrieval (CLIR), enhancing its applicability across diverse sectors. This paper aims to define CLIR and explore the techniques and methodologies commonly employed in its implementation. It discusses the algorithms utilized in CLIR systems and presents performance metrics used to evaluate their effectiveness, including precision, recall, and f1-score. By analyzing these metrics, the paper highlights the ongoing challenges in multilingual information retrieval and emphasizes the need for continuous improvements to optimize system performance and user experience in accessing information across languages.

## REFERENCES
[1]. Oard, D. W., & Dorr, B. (1996). A Survey of Multilingual Information Retrieval Systems. Proceedings of the 19th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, 1996.
[2]. Steinberger, J., & Karpov, J. M. K. (2003). Statistical Approaches to Multilingual Information Retrieval. Journal of Information Retrieval, 6(2), 145-167.
[3]. Huang, J., et al. (2007). A Language Model-Based Framework for Cross-

Language Information Retrieval. Proceedings of the 30[th]Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, 2007.

[4]. Johnson, R., et al. (2017). Neural Machine Translation and Sequence-to-Sequence Models: A Review. Journal of Artificial Intelligence Research, 60, 1-33.

[5]. Kuo, C., & Wu, H. (2020). User Experience in Cross-Language Information Retrieval: A Study on Interface Design. Journal of Information Science, 46(3), 345-358.

[6]. Johnson, R., et al. (2017). Neural Machine Translation and Sequence-to-Sequence Models: A Review. Journal of Artificial Intelligence Research, 60, 1-33.

[7]. Zhai, C., & Lafferty, J. (2001). A Risk Minimization Framework for Information Retrieval. Proceedings of the 24th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, 2001.

[8]. Liu, Y., & Croft, W. B. (2004). Cluster-Based Language Modeling for Cross-Language Information Retrieval. Proceedings of the 27th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, 2004.