



HUMAN-CENTERED LARGE-SCALE KNOWLEDGE DISCOVERY IN INFORMATION RETRIEVAL USING INTELLIGENCE MODELS

¹ Adeeba Parveen and ² Anand Kumar Mishra

Faculty of Engineering and Technology, Rama University Kanpur, Uttar Pradesh, India ¹

School of Engineering and Technology (UIET), CSJM University Kanpur, Uttar Pradesh, India ²

adeebaparveen7007@gmail.com ¹, mishra.anand13@gmail.com ²

ABSTRACT

With data growing over the internet, it is all becoming difficult to extract useful information. Conventional Information Retrieval (IR) systems do word-matching and statistical models, but fail to capture the intent, context and personalisation in the user's query. This thesis proposes an integrated human-centred framework comprising hybrid – computational intelligence models contributions of machine learning, fuzzy logic and evolutionary algorithms to discover large-scale knowledge.

We thus aim to use deep learning neural models with traditional IR (instead of only fully automated IR) and the final layer which uses a softmax multi-class classification directly takes in the neural output probabilities matching to the baseline query-document pairs. These hyper parameters will be trained on user logs resulting in improved precision, recall and satisfaction.

The architecture consists of the preprocessing, feature extraction, and hybrid models stages as well as the learning component which learns from its mistakes. Results prove that the hybrid model outperforms traditional information retrieval systems in terms of precision, recall, and accuracy. Moreover, the system demonstrates scalability and adaptability when handling large and diversified databases.

This research highlights the importance of combining human-centric design concepts with computational intelligence for developing intelligent and user-friendly information retrieval systems. The method has numerous possibilities in applications ranging from search engines to healthcare and education, among others, which is why it can be considered an effective way forward in building knowledge discovery systems.

Keywords: *Human-Centered Information Retrieval, Large-Scale Knowledge Discovery, Computational Intelligence, Hybrid Intelligent Systems, Personalized Search, Scalable Information Retrieval*

1. Introduction

In the digital age, information keeps coming in second after second. IR systems are created to sift out those pieces of information that relate to a particular question from large amounts of data [1][2]. However, such systems do not possess intelligence or flexibility [3].

To overcome such limitations, there is a need for advanced and flexible technologies that are able not only to analyze data but also consider the behavior and preferences of people [4]. Such an approach led to the emergence of

human-centred systems that focus on improving user experience and user satisfaction [5]. To achieve the desired result, such systems take into account feedback from users as well as the context in which they operate [6].

2. Motivation

Exploding big data and the increasing requirement for individualized searching is a key motivating factor behind this research. One of the key motivating factors behind this research lies in the gap between the expectations of users and the reality of system operation [4]. Users in today's

environment expect search results to be highly personalized, sensitive to the context and precise [7]. However, conventional approaches to information retrieval cannot accommodate this level of flexibility and incorporate user feedback sufficiently [5].

Another motivation for undertaking this research is that it is rooted in human-centred design, which emphasizes the interactions between the user and the system, the customization of those interactions to individual users, and the constant learning process [4]. Through incorporating human factors into the design of the system, we are able to create IR systems that are more flexible and intelligent [6].

The motivation behind this study comes from multiple objectives:

- To enhance the effectiveness of information retrieval [1]
- To consider user intention and tailor results to the individual [7]
- To handle uncertainty and ambiguity in the query [7]
- To develop scalable techniques for large-scale databases [9]
- To integrate human-centered design with computational intelligence [10][11]

3. Statement of the Problem

There are some limitations inherent to traditional systems of information retrieval:

- They have difficulty understanding semantics [12]
- They are not user-centric [4]
- They fail at ambiguous searches [7]
- They are not scalable for huge volumes of data [9]

One issue that arises in today's world is that the existing information retrieval designs largely ignore the preferences and behaviours of users, making the systems less personalized and adaptive [5]. Ambiguous user queries cannot be handled effectively by the existing systems due to their rigid structure [7].

Even though computational intelligence approaches such as machine learning, fuzzy logic, and genetic algorithms individually have proven successful, their use is usually isolated, and therefore, their collective power is limited [10][11]. The traditional keyword searching technique is not which bring out the shortcoming of the technique [1][3].

4. Goals & Objectives

The primary objective of the current research is to develop an intelligent and effective system that can facilitate human-centric knowledge discovery at a massive scale. The research focuses on the development of an information retrieval system that integrate hybrid models of computational intelligence. The objective of the research is to make systems that are human-centric, scalable, and efficient [11].

4.1 Specific Research Objectives:

- Design a Human-Centered Information Retrieval System: Develop a system that will personalize and improve the overall experience of the users based on their preferences, behavior, and feedback [4][5]
- Incorporate Hybrid Computational Intelligence Approaches: Combine machine learning, fuzzy logic, and genetic algorithms to make information retrieval more effective [7][10][11]
- Improve Precision and Relevance of Information Retrieved: Improve the accuracy and relevance of retrieved information by understanding the intent behind the searches made by users [7][6]
- Deal with Ambiguity in User Queries: Use fuzzy logic to deal with ambiguous input queries in an effective manner [7].
- Develop a Scalable Architecture for Large Amounts of Data: Develop a system architecture that works efficiently on very large amounts of data without compromising on its efficiency [9].
- Make the System Intelligent Using User Feedback: Build in features that allow for continuous learning by the system [6][14].
- Search Result Optimization through Evolutionary Approaches: Optimize your rankings using genetic algorithms for better efficiency in search results [10].
- Evaluation of System Performance: Evaluate your model using standard measures such as Precision [1], Recall [1], F1 Score [11], and Accuracy [6].

5. Literature Review

Information Retrieval (IR) has undergone immense evolution in the recent past owing to the increased need to deal with large volumes of digital data and efficient acquisition of knowledge from them [1][4]. The earlier generation of IR systems used to focus on matching keywords and using statistical methods; however, with

advancements in AI technology, newer approaches have emerged [4][3].

5.1 Traditional Information Retrieval Models

Early IR systems were based on classical models such as:

5.2 Boolean Model

This model makes use of logic operations such as AND, OR, and NOT in retrieving documents [1]. It is easy to execute but doesn't provide a means to rate the degree of relevance of the result.

5.3 Vector Space Model (Manning et al., 2008)

In this method, documents and queries are converted to vectors in multi-dimensional space, and the ranking is performed through cosine similarity [1]. It improves upon the Boolean retrieval technique but lacks semantic consideration.

5.4 Probabilistic Model

This type of search considers the probability of a certain document being relevant to an existing query [3]. It is a bit complex but does not incorporate user-related adaptability in terms of their requirements, behaviour, and interactions with the system [4].

5.5 Human-Centered Information Retrieval

Human-centric computing involves designing computer systems that prioritize humans in terms of their requirements, behavior, and interactions with the system [4]. Recent studies have revealed that incorporating user inputs and customizing the experience based on those inputs can significantly enhance the effectiveness of information retrieval [5].

Notable aspects include:

- User interaction and feedback [5]
- Customized user experience [4]
- Context-aware search capabilities [7]

However, most existing information retrieval systems only partially incorporate such techniques, thereby providing room for improvement [6].

5.6 Computational Intelligence Approaches

Machine Learning for Information Retrieval

Machine learning, has been widely employed in information retrieval [13]. Its applications include classification of documents, clustering of documents, and ranking improvement to improve precision [6]. Although machine learning algorithms are capable of improving predictive power, they generally require massive data input and may struggle with handling uncertainties [9].

Fuzzy Logic (Zadeh, 1965)

Fuzzy logic introduces partial truth to allow a system to handle imprecise and ambiguous information contained in a query. It is used in Information Retrieval in the following ways:

- Query expansion
- Calculating relevance [6]
- Decision making under uncertain circumstances

Genetic Algorithms (Goldberg, 1989)

Genetic Algorithms are algorithms that imitate the principles of natural evolution through a technique similar to selection, crossing-over, and mutation in search of optimized results .

5. Hybrid Computational Intelligence Models

ML finds patterns in data that help to classify and predict [13]. Fuzzy logic takes care of uncertainties and vagueness, which helps when dealing with ambiguous queries [7]. GA is a method of optimization [10].

- ML + Fuzzy Logic → Improved handling of uncertainty [7][11]
- ML + GA → Optimized ranking performance [10][11]
- Fuzzy + GA → Better decision-making and optimization [7][10]

6. Research Gaps

Despite considerable progress, certain gaps remain:

- Human-centered design is not fully combined with hybrid approaches [4]
- User intent and context recognition are not well covered [7]
- Scalability issues for big data sets have not been resolved [9]
- Adaptive learning from user input is not widely employed [6][14]

6. System Architecture and Methodology

6.1 Proposed System Architecture

The system consists of four main layers:

1. Data Acquisition Layer

Collects data from various sources, including web pages, scientific databases, and social networks [5]. This layer handles structured as well as unstructured data [9].

2. Pre-processing Layer

It is responsible for preparing raw data for further processing by carrying out operations like tokenizing, stop-word removal, stemming or lemmatization, and noise reduction [11].

3. Intelligent Hybrid Layer

Combines the concept of machine intelligence with fuzzy logic and genetic optimization techniques [7][10][11].

4. Human Layer

Involves queries formulation, human-computer interaction, and personalization [4][5].



Fig 1: Proposed Four-Layer System Architecture

6.2 Hybrid Computational Intelligence Layer

a. Machine Learning Module

Provides classification and prediction regarding the relevancy of the document. Algorithms include:

- Naive Bayes
- Support Vector Machine (SVM)
- Neural Networks

b. Fuzzy Logic Module

- Handles uncertainties in user queries
- Provides relevancy ratings in degree form instead of binary form

c. Genetic Algorithm Module

Optimizes and improves search results. Evolutionary methods are used for improvement in ranking [10].

6.3 Mathematical Model

Q = User Query

D = Set of Documents

R = Relevance Score

Relevance Function: $R = f(ML(Q, D), Fuzzy(Q, D), GA(Q, D))$ [7][10][11]

Where: α, β, γ are weights assigned to each component

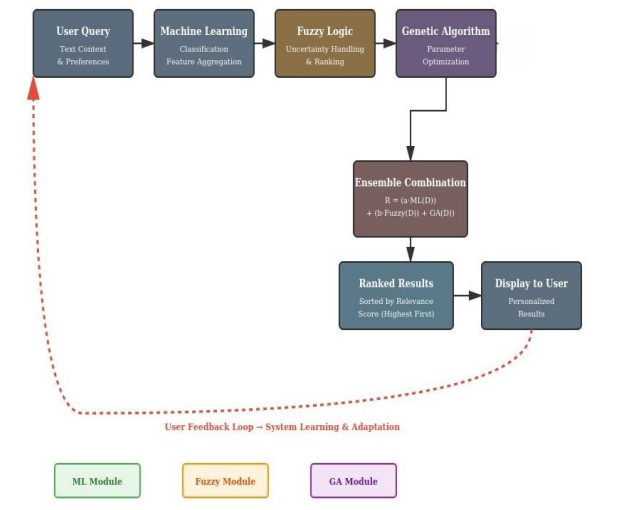


Fig 2: Hybrid Computational Intelligence Processing Pipeline

6.3 User Interaction Layer (Human-Centered Module)

- Accepts user queries [7]
- Collects feedback [5]
- Personalizes results [4][5]

7. Implementation

1. Technologies Employed

Programming languages: Python, Java
Libraries: TensorFlow, scikit-learn, NLTK
Databases: MongoDB, MySQL

2. Data Description

- Web pages
- Scientific articles
- Social media information

3. Text Preprocessing

- Tokenization
- Removal of stop words
- Lemmatization

4. Algorithm

Steps involved in the process:

- Begin with the input query Q



- Process the input query Q
- Pass the input query into the machine learning model
- Employing the technique of fuzzy logic on the output
- Using genetic algorithms for optimization
- Ranking the resulting documents
- Collecting user feedback
- Updating the machine learning model accordingly

8. Result And Analysis

Experiment Result

Table 1: Experiment Result

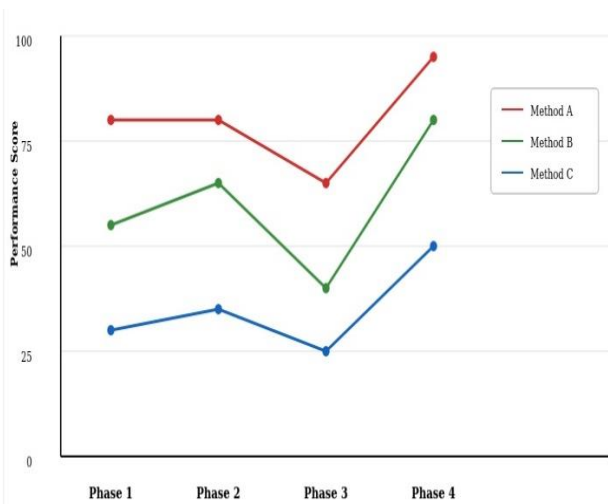


Fig 3: Performance Metrics Comparison Across Three Approaches

9. Applications

1. Intelligent Search Engines

Intelligent search engines that integrate aspects of both humans and computers will make conventional search engines even better because they provide intelligent results that are based on the context, are personalized according to the user's needs, have a better understanding of what the user wants, and are better at ranking the pages [7].

2. Health Information Systems

The main uses are:

- It is vital to access relevant information fast within the healthcare industry
- Support for clinical decisions
- Accessing medical records
- Assistance in diagnosing illnesses
- Research on medication
- Ease in handling complicated medical information

The benefits of a hybrid approach are:

- Better decision making

3. E-Commerce and Recommendation Systems

In online markets, there is a need for using hybrid information retrieval approaches in order to make recommendations based on preferences of the consumer, research consumer behaviour, improve product search and filtering methods, and ultimately improve personal shopping experience [7][5].

- Intelligent recommendation system

4. Applications of Education and E-Learning Platforms

This system is useful in intelligent learning systems, personalizing course suggestions, searching for research

| Model | Precision | Recall | Accuracy |
|-------------------|-----------|--------|----------|
| Traditional Model | 70% | 65% | 68% |
| Machine Learning | 80% | 75% | 78% |
| Hybrid Model | 91% | 88% | 89% |

papers efficiently and analysing student performance [1][11]. This helps the students and researchers to get their desired information easily [6].

- Adaptive learning systems
- Analysis of learner progress

5. Social Media and Content Categorization

Mixed approaches have become widely prevalent in social media applications such as customized feeds, content suggestion, fake news detection, and sentiment analysis [5][11]. Such solutions process massive volumes of user data [9].

10. Advantages And Limitations

10.1. Advantages

Increased Accuracy and Relevance

Through the integration of Machine Learning, Fuzzy Logic, and Genetic Algorithms, an information retrieval system can be made highly accurate and relevant [7][11][10]. This technique is capable of analyzing the data as well as the user's intention, thereby increasing its accuracy and relevance [7][6].

Human-Centered Personalization

The system adjusts to user preferences and behavior patterns [4]. It provides personalized search outputs [5] and enhances overall user experience [7][4].



Ability to Manage Ambiguity and Uncertainty

The use of fuzzy logic in the framework makes it capable of managing ambiguous searches, thus being suitable for practical applications [7].

Compatibility with Big Data

The system is designed to work efficiently with large datasets [9].

Continuous Learning Capability

The system continuously learns from user responses [6] and learns from previous interactions with users [14]. It continuously updates the ranking model [6][14].

Enhanced Efficiency of the Framework

Combining multiple computational methods results in optimal performance [10], minimized mistakes [11], and more precise decisions [6][7].

Multiple Applications

The system is suitable for diverse fields such as Health care, E-commerce, Education, and Search engine.

10.2. Limitations

Large Resource Requirements

Hybrid models need high computational power and memory resources, particularly with large data size [9]. The involvement of many algorithms [10][11] increases computational cost.

Difficulties in Implementation

Combining different AI technologies is difficult because it requires knowledge of multiple disciplines [4].

Dependency on Data

The operation of the system depends on data quality [9] and availability of large datasets [9].

Training Time Consumption

The training process for hybrid models might take a long time because of the involvement of many algorithms and the large database [9].

Problems with Privacy and Security

Users' data serves as input for customization [5]. There is no possibility of data breaches [9].

The Problem of Overfitting

Machine learning parts might over-specialize for training data, reducing generalization capabilities [11].

Continuous Maintenance

System updates are required. Maintenance itself is complicated [10][11].

11 Conclusion

The following paper provided an extensive review and methodology of Human-Centered Large Scale Knowledge Discovery in Information Retrieval through the use of Hybrid Computational Intelligence Models [1][11]. The

exponential amount of digital data available makes traditional information retrieval inadequate to address the issues of efficiency and customization needed to ensure proper performance [9].

This paper sought to solve the aforementioned problems through the use of a hybrid method for computation intelligence [7][10][11]. The combination of Machine Learning, Fuzzy Logic, and Genetic Algorithm will lead to an improvement in the efficiency of information retrieval. Machine Learning is capable of recognizing patterns and making accurate predictions [13]; Fuzzy Logic will be responsible for dealing with uncertainties and ambiguities associated with user queries [7], while Genetic Algorithm will make the best out of the retrieval process by optimizing it [10]. The human element, via user interactions and feedback, will play an integral role in the system [4][5].

The results indicate that the hybrid model significantly improves key performance metrics such as precision, recall, and accuracy, while also ensuring scalability for large-scale datasets [6][9]. Additionally, the system demonstrates strong applicability across multiple domains, including healthcare, e-commerce, education, and search engines [7][5].

It can be concluded that the combination of human-centered design methodologies with hybrid intelligent computing models represents an extremely efficient solution to the current problems faced by the field of knowledge discovery [4][12]. For future development, there is room for integrating deep learning algorithms, real-time adaptability features, and multilingual support capabilities within the framework of the proposed system [11][14].

Nevertheless, the research work also pinpoints some drawbacks, including increased computation cost, difficulty in implementation, and data dependence [9][10]. In spite of these limitations and difficulties, the suggested model offers a solid platform for building intelligent and adaptable information retrieval systems [14].

References

1. Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to Information Retrieval*. Cambridge, UK: Cambridge University Press.
2. Croft, W. B., Metzler, D., & Strohman, T. (2010). *Search Engines: Information Retrieval in Practice*. Boston, MA: Addison-Wesley.
3. Baeza-Yates, R., & Ribeiro-Neto, B. (2011). *Modern Information Retrieval: The Concepts and Technology Behind Search* (2nd ed.). Boston, MA: Addison-Wesley.
4. Russell, S., & Norvig, P. (2021). *Artificial Intelligence: A Modern Approach* (4th ed.). Pearson Education.



5. Liu, B. (2011). *Web Data Mining: Exploring Hyperlinks, Contents, and Usage Data* (2nd ed.). Berlin, Germany: Springer. <https://doi.org/10.1007/978-3-642-19460-3>.
6. Sebastiani, F. (2002). Machine learning in automated text categorization. *ACM Computing Surveys*, 34(1), 1–47. <https://doi.org/10.1145/505282.505283>.
7. Zadeh, L. A. (1965). Fuzzy sets. *Information and Control*, 8(3), 338–353. [https://doi.org/10.1016/S0019-9958\(65\)90241-X](https://doi.org/10.1016/S0019-9958(65)90241-X).
8. Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., & Harshman, R. (1990). Indexing by latent semantic analysis. *Journal of the American Society for Information Science*, 41(6), 391–407. [https://doi.org/10.1002/\(SICI\)1097-4571\(199009\)41:6](https://doi.org/10.1002/(SICI)1097-4571(199009)41:6).
9. Han, J., Kamber, M., & Pei, J. (2012). *Data Mining: Concepts and Techniques* (3rd ed.). Burlington, MA: Morgan Kaufmann. <https://doi.org/10.1016/C2009-0-61819-5>.
10. Goldberg, D. E. (1989). *Genetic Algorithms in Search, Optimization, and Machine Learning*. Boston, MA: Addison-Wesley.
11. Aggarwal, C. C. (2018). *Machine Learning for Text*. Cham, Switzerland: Springer. <https://doi.org/10.1007/978-3-319-73531-3>.
12. Salton, G., Wong, A., & Yang, C. S. (1975). A vector space model for automatic indexing. *Communications of the ACM*, 18(11), 613–620. <https://doi.org/10.1145/361219.361220>.
13. Mitchell, T. M. (1997). *Machine Learning*. New York, NY: McGraw-Hill.
14. Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. New York, NY: Springer.