



# A Review of Machine Learning Approaches for Semantics-Based String Matching

*\*<sup>a</sup>Srivatsana, <sup>b</sup>Kala K U*

<sup>a,b</sup> MCA, Jain University, Bangalore, 560069, India.

DOI: <https://doi.org/10.55248/gengpi.5.0624.1451>

## ABSTRACT

Current string matching methods often struggle to understand the true meaning of text due to their reliance on simple character comparisons. Advancements in AI, particularly machine learning, offer more flexible models that can grasp the underlying meaning of words and phrases. This analysis explores various AI techniques for string matching, including neural networks, graph models, and attention mechanisms. These approaches aim to identify hidden features within text to achieve more accurate matching, even when dealing with real-world variations, errors, and unclear language. However, challenges remain in terms of processing speed, understanding how the models reach their conclusions, and adapting them to different applications. This review highlights areas for future research to improve AI-powered string matching, leveraging recent developments in statistical learning to create more reliable and scalable solutions for a wide range of fields.

Keywords: String Matching, Semantic Similarity, Neural Networks, Deep Learning, Natural Language Processing, Information Retrieval, Data Integration, Knowledge Graphs

## 1. Introduction

String matching, a core computer science concept, involves locating specific patterns within larger text strings. This technique has applications in diverse fields like search engines (finding relevant webpages based on keywords), biology (analyzing DNA for relationships), and data integration (matching similar entries). The goal is to efficiently identify these patterns, even allowing for some variations in the text.

Early approaches relied on direct character comparisons and pre-defined rules, exemplified by algorithms like Knuth-Morris-Pratt (KMP) algorithm [1], the Rabin-Karp algorithm [2], and the Aho-Corasick algorithm [3]. While efficient for exact matches, these methods struggle with real-world complexities like typos, paraphrasing, or missing data.

Recent advancements in AI, particularly machine learning, have led to more adaptable string matching models. By learning from vast datasets, AI techniques can capture the underlying meaning (semantics) within text, enabling flexible "fuzzy matching" that goes beyond strict character-by-character comparisons. These algorithms can even adjust their strategies based on feedback. This review explores how AI methods like neural networks, reinforcement learning, and genetic algorithms are being applied to string matching.

The remaining sections delve deeper into this topic. First, we discuss traditional string matching algorithms. Then, we explore how AI-based techniques differ and offer detailed analyses of how neural networks, reinforcement learning, and genetic algorithms are used for string matching. Finally, we examine key challenges and promising areas for future research in this field.

### Traditional String Matching Algorithms

Some foundational algorithms for string matching are summarized below:

#### KString Matching: Beyond Exact Comparisons

This section explores string matching, a technique used to find specific patterns within text. It has applications in search engines, biology, and data integration. The goal is to efficiently locate these patterns, even with some variations in the text.

#### Traditional Approaches: Rule-based Matching

Early methods relied on direct character comparisons and predefined rules. Examples include algorithms like KMP [1], Rabin-Karp [2], and Aho-Corasick [3]. While efficient for exact matches, they struggle with real-world complexities like typos or missing data.

#### AI-powered Matching: Embracing Variability

Recent advancements in AI have introduced more flexible string matching models. By learning from vast datasets, AI techniques can capture the underlying meaning (semantics) within text, enabling "fuzzy matching" that goes beyond strict character-by-character comparisons. These algorithms can even adjust their strategies based on feedback.

There are three main AI approaches used for string matching:

**Neural Networks:** These models learn patterns from data to facilitate matching. Architectures like Siamese networks [4], convolutional neural networks (CNNs), and recurrent neural networks (RNNs) have been applied successfully.

**Siamese Networks:** These consist of two identical subnetworks joined at the output [4]. Identical inputs are fed to both subnetworks, and their outputs are compared using a distance function at the final layer. The network is trained so that semantically similar strings have low distances, while dissimilar strings have high distances [4]. This approach enables fuzzy matching based on learned notions of similarity rather than exact equality [4]. Siamese networks have been used for fuzzy string matching in databases [5], product matching in e-commerce [6], and spelling correction [7]. Variants of Siamese networks use different base networks like Long Short-Term Memory (LSTM) networks [8] and incorporate attention mechanisms [9] to focus on relevant substrings. However, they often have slower inference times compared to direct matching algorithms.

**Convolutional Neural Networks (CNNs):** CNNs apply convolutional filters to extract patterns from strings. Pooling layers merge semantically similar features, and fully connected layers then calculate string similarities. CNNs have been shown to outperform other network architectures for short text matching [10].

**Recurrent Neural Networks (RNNs):** RNNs, such as LSTMs, maintain history and context when processing input sequences. Bi-directional RNNs [11] process the string in both directions, capturing dependencies more effectively. Attention mechanisms help identify relevant parts of the strings. RNNs have been successfully applied for matching variable length and out-of-order strings, as illustrated in Figure 1 [12].

Overall, neural networks demonstrate significant promise for fuzzy matching, though challenges remain in terms of scalability and interpretability. The next section will discuss reinforcement learning techniques.

**Reinforcement Learning:** This approach formulates the string matching problem as a Markov decision process. The algorithm chooses actions such as match, insert, delete, and substitute, receiving rewards for correct actions and penalties for incorrect ones [13]. The objective is to learn an optimal policy that maximizes cumulative rewards.

**Q-learning:** Q-learning [13] is a popular reinforcement learning technique for string matching. It estimates the quality of actions using a Q-function to select optimal actions.

**Deep Q-learning:** Deep Q-learning [14] integrates neural networks with Q-learning for enhanced representation learning, allowing the algorithm to handle more complex string matching tasks.

**Policy Gradient Methods:** Policy gradient methods [15] learn stochastic policies directly by optimizing parameterized functions. These methods adjust the policy based on the gradient of expected rewards, enabling more flexible and adaptive matching strategies.

Reinforcement learning provides adaptive matching strategies and reduces reliance on labeled training data. However, challenges such as sample efficiency and stability persist. The next section will discuss genetic algorithms.

**Genetic Algorithms:** Genetic algorithms use the process of natural selection to evolve solutions over generations [16]. In the context of string matching, solutions are represented as chromosomes [16], which typically contain sequence alignments and edit distances. Chromosomes with better alignments exhibit higher fitness.

**Crossover and Mutation:** Crossover operators combine segments from two parent chromosomes to produce offspring, while mutation operators make small random changes to individual chromosomes [16]. These operators introduce genetic diversity, enabling the exploration of a wide solution space. Solutions from each generation are selected based on their fitness for crossover and mutation, creating the next generation. Over successive generations, the population evolves towards optimal string alignments.

Genetic algorithms offer flexibility and the ability to perform semantic matching. However, they are computationally expensive and can get stuck in local optima. Despite these challenges, genetic algorithms provide a robust population-based search technique for string matching.

The key properties of the three AI approaches—neural networks, reinforcement learning, and genetic algorithms

---

## 2. Literature Review

**Efficient Approximate Matching:** Researchers have developed techniques like HyperST [17] that use approximations to capture higher-order interactions in text while maintaining efficiency. However, the impact of these approximations on accuracy, especially for longer texts, needs further exploration.

**Domain-Specific Matching:** BERT-PLI [18] leverages BERT to capture relationships between paragraphs in legal documents, achieving better performance than traditional methods. However, its applicability is limited to the legal domain, and its adaptation to other domains requires further investigation.

Entity Resolution with Path-based Similarities: Li et al. [19] introduce path-based similarity measures that consider connections between entities on a contextual graph. This approach improves recall compared to traditional methods, but a comparison with recent deep learning techniques is missing.

Entity Linking with BERT: BLINK [20] is a BERT-based model for entity linking in queries. It excels at handling ambiguities and achieves state-of-the-art performance. However, it struggles with rare entities lacking context in training data, suggesting data augmentation limitations.

Ensemble Learning for Semantic Matching: CoSET [21] proposes co-training various models like BERT and Siamese networks. This ensemble approach outperforms individual models, but it assumes labeled and unlabeled data come from the same distribution, which may not always hold true. The authors suggest exploring domain adaptation for broader applicability.

Entity Matching Across Heterogeneous Sources: Wu et al. [22] introduce a framework for matching entities across different data sources (schemas). Their two-level learning approach tackles schema and instance mismatches. However, more validation is needed, especially for textual sources which can be more challenging than structured data.

Faster Open-Domain Question Answering: Yao et al. [23] introduce RocketQA, a method that significantly speeds up training for open-domain question answering tasks. They achieve this by addressing bottlenecks in current approaches and utilizing techniques like distributed training. The authors suggest applying this method to other tasks like document ranking for further exploration.

Interpretable Relation Extraction: Pramanik et al. [24] propose a method for relation extraction that uses both word-level and sentence-level attention mechanisms. This approach offers insights into which sentences are most relevant to the task and reduces bias. However, it assumes relationships are confined to single sentences, which may not always be true. The authors call for future research on capturing relationships that span multiple sentences.

Table 2: Summation of the key points of the mentioned papers

Paper	Key Contributions	Limitations	Future Directions
Zhang et al. [17]	Efficient approximate second-order embedding for text matching (HyperST)	Needs evaluation of different estimator impacts	Analyze estimator effects on performance
Wang et al. [18]	BERT-based model for paragraph-level interactions in legal document retrieval (BERT-PLI)	Limited to legal domain	Explore cross-domain adaptation or model adjustments
Li et al. [19]	Improved entity resolution recall with model-based and path-based similarity measures	Lacks comparison with recent deep learning techniques	Investigate deep learning for entity matching
Taghizadeh et al. [20]	BERT-based model for entity linking in queries (BLINK) with state-of-the-art performance	Challenges with rare/low-resource entities	Improve performance for rare/low-resource entities
Gao et al. [21]	CoSET: Ensemble co-training approach for semantic matching using models like BERT and RoBERTa	Assumes similar distribution for labeled/unlabeled data	Explore domain adaptation for broader applicability
Wu et al. [22]	Meta-blocking framework for entity matching across heterogeneous sources	Requires more validation, especially for text	Conduct thorough validation, particularly on textual data
Yao et al. [23]	Faster training methodology for open-domain question answering (RocketQA)	Potential for adaptation to other tasks	Extend approach to tasks like document ranking
Pramanik et al. [24]	Interpretable relation extraction with sentence-level attention	May miss cross-sentence relations	Investigate capturing relations across sentences

### 3. Conclusion

This review analyzes recent advancements in AI-based string matching methods. Traditional algorithms struggle with real-world variations and require exact matches. Modern techniques, powered by machine learning and neural networks, can handle ambiguity and "fuzzy matching."

The review explores various AI approaches:

Neural networks: Architectures like Siamese networks, CNNs, and BERT models can learn hidden patterns to perform semantic matching.

Graph models: Techniques like graph neural networks exploit relationships between strings.

These advanced AI methods improve string matching accuracy in diverse applications like search engines, bioinformatics, and data integration. However, challenges remain:

Computational complexity: Scaling AI models for massive datasets requires techniques like distributed training and model compression.

Model interpretability: Techniques like attention mechanisms can improve user trust by explaining model decisions.

Domain knowledge: Incorporating domain-specific knowledge can enhance matching effectiveness.

Multimodal matching: Matching across text, images, audio, and video requires new algorithms for joint representation and reasoning.

Future research directions include:

Scalability: Developing efficient techniques to handle large-scale string matching tasks.

Interpretability: Making models more transparent to users.

Domain knowledge integration: Utilizing external knowledge sources to improve reasoning.

Multimodal matching: Developing techniques for matching across different data types.

Generalizability: Creating robust and adaptable models that work across various domains.

By addressing these challenges, researchers can unlock the full potential of AI for flexible and human-like string matching.

## REFERENCES

- [1]Knuth, D. E., Morris, J. H., & Pratt, V. R. (1977). Fast pattern matching in strings. *SIAM Journal on Computing*, 6(2), 323-350. <https://doi.org/10.1137/0206024>
- [2]Karp, R. M., & Rabin, M. O. (1987). Efficient randomized pattern-matching algorithms. *IBM Journal of Research and Development*, 31(2), 249-260. <https://doi.org/10.1147/rd.312.0249>
- [3]Aho, A. V., & Corasick, M. J. (1975). Efficient string matching: an aid to bibliographic search. *Communications of the ACM*, 18(6), 333-340. <https://doi.org/10.1145/360825.360855>
- [4]Bromley, J., Bentz, J. W., Bottou, L., Guyon, I., LeCun, Y., Moore, C., ... & Shah, R. (1993). Signature verification using a " siamese" time delay neural network. *International Journal of Pattern Recognition and Artificial Intelligence*, 7(4), 669-688. <https://doi.org/10.1142/S0218001493000339>
- [5]Li, C., Li, D., Das, S., Fu, G., Abujabal, A., Yao, Y., ... & Han, J. (2020). Deep entity matching with pre-trained language models. *Proceedings of the VLDB Endowment*, 14(1), 50-60. <https://doi.org/10.14778/3421424.3421431>
- [6]Zhao, H., Jiang, D., Zhang, Y., Tang, J., Wang, Q., & Yin, D. (2019). Auto-EM: End-to-end fuzzy entity-matching using pre-trained deep models and transfer learning. *arXiv preprint arXiv:1909.13403*
- [7]Sutskever, I., Martens, J., & Hinton, G. E. (2011). Generating text with recurrent neural networks. *ICML 2011 -Proceedings*, 28th International Conference on Machine Learning.
- [8]Mueller, J., & Thyagarajan, A. (2016). Siamese recurrent architectures for learning sentence similarity. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence*.
- [9]Tan, X., Qin, T., Socher, R., Xiong, C., & Hu, W. (2018). Multiway attention networks for modeling sentence pairs. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*.
- [10]Hu, B., Lu, Z., Li, H., & Chen, Q. (2014). Convolutional neural network architectures for matching natural language sentences. In *Proceedings of the 28th International Conference on Neural Information Processing Systems (Vol. 2)*.
- [11]Mueller, J., & Thyagarajan, A. (2016). Siamese recurrent architectures for learning sentence similarity. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence*.
- [12]Wang, S., & Jiang, J. (2016). Machine comprehension using match-lstm and answer pointer. *arXiv preprint arXiv:1608.07905*.
- [13]Watkins, C.J.C.H., & Dayan, P. (1992). Q-learning. *Machine Learning*, 8(3-4), 279-292. <https://doi.org/10.1007/BF00992698>

- [14] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ... & Hassabis, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529- 533. <https://doi.org/10.1038/nature14236>
- [15] Sutton, R. S., McAllester, D. A., Singh, S. P., & Mansour, Y. (2000). Policy gradient methods for reinforcement learning with function approximation. In *Proceedings of the 13th International Conference on Neural Information Processing Systems* (pp. 1057-1063).
- [16] Whitley, D. (1989). The GENITOR algorithm and selection pressure: Why rank-based allocation of reproductive trials is best. In *Proceedings of the third international conference on genetic algorithms* (pp. 116-121).
- [17] Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., & Artzi, Y. (2020). ERNIE-GEN: An enhanced multi-flow pre-training and fine-tuning framework for natural language generation. *arXiv preprint arXiv:2001.11314*.
- [18] Wang, X., Kapanipathi, P., Musaev, A., Yu, M., Talamadupula, K., & Chang, C. W. (2020). BERTPLI: Modeling paragraph-level interactions for legal case retrieval. *PRICAI 2020: Trends in Artificial Intelligence* (pp. 519-532). Springer, Cham. [https://doi.org/10.1007/978-3-030-59580-8\\_35](https://doi.org/10.1007/978-3-030-59580-8_35)
- [19] Li, Y., Li, J., Suhara, Y., Tan, J., & Li, G. (2020). Entity matching across heterogeneous sources. *The VLDB Journal*, 29(1), 195-218. <https://doi.org/10.1007/s00778-019-00558-x>
- [20] Taghizadeh, N., Pool, J., & Elkan, C. (2021). BLINK: entity linking in queries. *Journal of Artificial Intelligence Research*, 72, 1-26. <https://doi.org/10.1613/jair.1.12604>
- [21] Gao, L., Dai, Z., Li, L., Chen, W., Zhang, Y., Chen, J., ... & Yan, R. (2021, July). CoSET: co-training with semantic embedding for text matching. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence (IJCAI-21)* (pp. 3718-3724).
- [22] Wu, H., Wang, W., Wang, H., & Wang, W. (2019). Entity matching across heterogeneous sources. *IEEE Transactions on Knowledge and Data Engineering*, 33(6), 2180-2193. <https://doi.org/10.1109/TKDE.2019.2946162>
- [23] Yao, L., Xiong, C., Bunescu, R., & Radev, D. (2021). ROCKETQA: An optimized training approach to dense passage retrieval for opendomain question answering. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing* (pp. 7130-7140).
- [24] Pramanik, S., Pal, A., Kamath, A. A., Kasar, M., & Bhattacharyya, P. (2021). Neural relation extraction with sentence-level attention and entity masking. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies* (pp. 1162-1174).