

From Static to Contextual: A Survey of Embedding Advances in NLP

Hussein Al-Kaabi^{1*}, Ali Kadhim Jasim², Ali Darroudi³

¹Ministry of Education Iraq, General Direction of Vocational Education, Al-Najaf, Iraq

²Department of Computer Technology Engineering, Imam Ja'afar Al-Sadiq University, Maysan - Iraq

³Department of Electrical Engineering, Sadjad University of Technology, Mashhad, Iran

¹hussain.njf7@gmail.com, ²dr.alikadhimjasim@gmail.com, ³darroudiali@sadjad.ac.ir



***Corresponding Author**

Article History:

Submitted: 29-06-2025

Accepted: 02-07-2025

Published: **09-07-2025**

Keywords:

BERT; Deep Learning; Machine Learning; Natural Language Processing (NLP); Word Embedding.

PERFECT: Journal of Smart Algorithms is licensed under a Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC 4.0).

ABSTRACT

Recent advances in Natural Language Processing (NLP) have highlighted the crucial role of embedding techniques in enabling machines to capture complex semantic and syntactic relationships within text. Despite the success of static word embeddings like Word2Vec and GloVe, these approaches face limitations in addressing polysemy and context variability. To address these challenges, the research objective of this paper is to comprehensively survey the development of embedding methods, moving from static representations to dynamic, contextualized models such as BERT and GPT. This survey employs a systematic literature review and comparative analysis of key models to identify their strengths, limitations, and application areas across diverse NLP tasks, including machine translation, sentiment analysis, and question answering. The findings demonstrate that contextual embeddings significantly improve performance by generating word representations conditioned on surrounding context, though they introduce computational and interpretability challenges. The survey also identifies emerging trends such as multimodal embeddings, domain-specific adaptations, and techniques to mitigate bias inherent in training data. The practical implications of this work lie in guiding practitioners to select appropriate embedding methods for specific tasks and highlighting best practices for deploying these models responsibly. In conclusion, while recent embedding techniques have advanced the state of NLP, open challenges remain in balancing efficiency, fairness, and transparency. Therefore, further research is needed to develop more robust, explainable, and resource-efficient embedding models that can be effectively applied in multilingual and low-resource settings.

INTRODUCTION

Natural Language Processing (NLP) has emerged as one of the most transformative fields in artificial intelligence, enabling machines to understand, interpret, and generate human language [1]. This revolution's heart lies in embedding numerical representations of words, phrases, or sentences that capture their semantic and syntactic properties [2]. Embeddings are the foundation for various NLP tasks, including machine translation [3], sentiment analysis [4], question answering [5], and spam classification [6]. By converting discrete linguistic units into continuous vector spaces, embedding allows algorithms to process and reason about language in previously unimaginable ways. Significant milestones have marked the journey of embedding techniques in NLP, each addressing the limitations of its predecessors and unlocking new possibilities. Early approaches to representing words relied on one-hot encoding, a simple yet highly inefficient method representing each word as a sparse binary vector [7]. While straightforward, one-hot encoding failed to capture semantic relationships between words, leading to high-dimensional and computationally expensive representations. This limitation spurred the development of distributional semantics, which posits that words occurring in similar contexts tend to have similar meanings. This idea laid the groundwork for the first generation of dense, low-dimensional word embeddings [8]. The introduction of Word2Vec in 2013 marked a turning point in NLP [9]. By leveraging neural networks to predict words based on their context (Skip-gram) or vice versa (CBOW), Word2Vec produced embeddings that captured intricate relationships between words [10]. Following this, GloVe (Global Vectors for Word Representation) combined the strengths of matrix factorization and local context window methods to create embeddings that incorporated global corpus statistics [11]. Around the same time, FastText extended Word2Vec by incorporating subword information, enabling the model to handle rare and out-of-vocabulary words more effectively [12]. While consequential, these static embeddings were limited by their inability to represent polysemous words with multiple meanings differently based on context [13]. The next major leap in embedding technology came with the advent of contextualized embeddings, which dynamically generate word representations based on their surrounding text. Models like ELMo (Embeddings from Language Models) introduced the idea of using bidirectional LSTMs to produce context-sensitive embeddings, allowing the same word to have different representations depending



on its usage [14]. This was further advanced by transformer-based architectures like BERT (Bidirectional Encoder Representations from Transformers) and GPT (Generative Pretrained Transformer) [15-16]. BERT, with its masked language modeling objective, and GPT, with its autoregressive approach, demonstrated unprecedented performance across a wide range of NLP tasks, setting new benchmarks and inspiring research into pre-trained language models. The shift from static to contextualized embeddings has improved the accuracy of NLP systems and enabled them to handle more nuanced aspects of language, such as ambiguity, coreference resolution, and long-range dependencies [17]. However, this progress has come with its own set of challenges. Contextualized models are computationally intensive, requiring significant resources for training and inference. Additionally, they often inherit biases in the training data, raising ethical concerns about their deployment in real-world applications [18]. These challenges have spurred ongoing research into more efficient, interpretable, and fair embedding techniques. This survey aims to provide a comprehensive overview of the evolution of embedding techniques in NLP, from their early beginnings to the state-of-the-art methods in use today. Section 2 explores traditional static embeddings, highlighting their strengths and limitations. Section 3 discusses the rise of contextualized embeddings, focusing on their architectures, training objectives, and impact on the field. Section 4 examines the diverse applications of embeddings across NLP tasks, including social media analysis, sentiment detection, machine translation, and spam filtering. Section 5 outlines the challenges and limitations of embedding techniques, such as computational complexity and bias. Finally, Section 6 presents emerging trends, including multimodal embeddings, domain-specific representations, and efforts to ensure fairness, followed by concluding remarks and directions for future research. Continuous advancements in machine learning and deep learning have driven the evolution of embedding techniques in NLP. This section reviews the key contributions in the field, tracing the transition from static word representations to dynamic, context-sensitive embedding.

TAXONOMY OF EMBEDDING TECHNIQUES

Embedding techniques in NLP have evolved from traditional static methods like One-Hot Encoding, Word2Vec, GloVe, and FastText to advanced contextualized models. While traditional methods provided foundational representations, they lacked contextual awareness. Word2Vec and GloVe captured semantic relationships, and FastText addressed rare words using subword information. The shift to contextual embeddings, marked by models like ELMo, BERT, and GPT, enabled dynamic representations based on word usage. These models set new performance benchmarks, leading to variants like RoBERTa, ALBERT, and DistilBERT. Recent trends include multimodal and domain-specific embeddings and efforts to reduce bias. Embeddings now power tasks such as sentiment analysis, translation, and question answering, enhancing accuracy and efficiency.

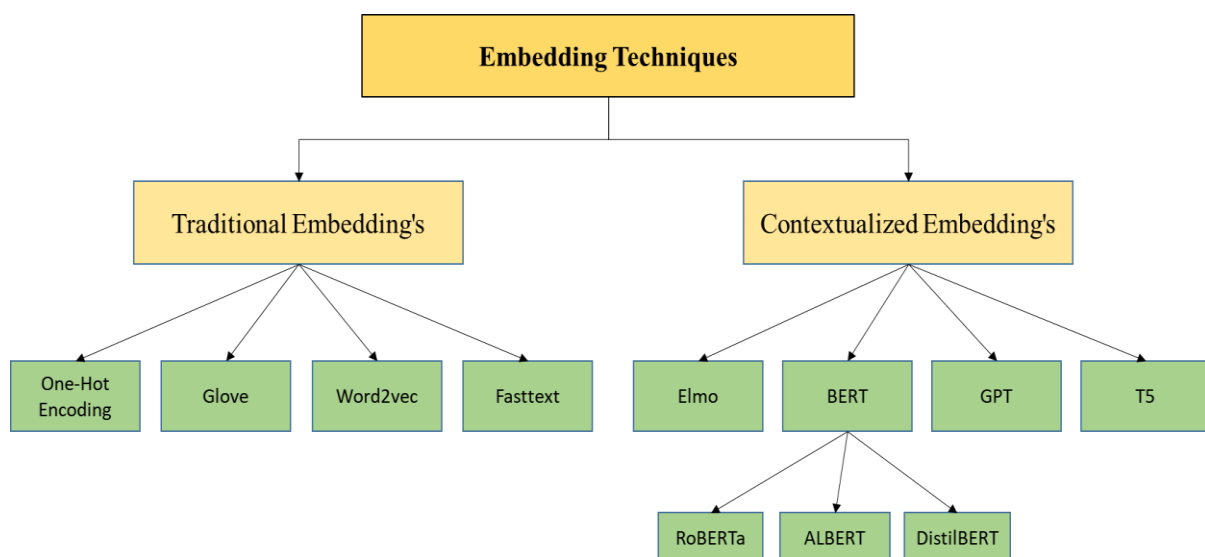


Figure 1. Embedding techniques taxonomy

TRADITIONAL EMBEDDING TECHNIQUES

Traditional embedding techniques laid the foundation for modern NLP by transforming words into numerical representations that machines can process. While more straightforward than their contemporary counterparts, these methods were groundbreaking in capturing semantic and syntactic relationships between words. This section explores

four key traditional embedding techniques: One-Hot Encoding, Word2Vec, GloVe, and FastText.

One-Hot Encoding

One-hot encoding is one of the simplest and earliest methods for representing words numerically. In this approach, each word in a vocabulary is represented as a sparse binary vector of size V , where V is the vocabulary size. The vector contains a 1 at the index corresponding to the word and 0s everywhere else [19]. Consider a vocabulary of four words: ["cat", "dog", "bird", "fish"]. The one-hot encoding for these words would be:

$$\begin{aligned} \text{"cat"} &\rightarrow [1, 0, 0, 0] \\ \text{"dog"} &\rightarrow [0, 1, 0, 0] \\ \text{"bird"} &\rightarrow [0, 0, 1, 0] \end{aligned}$$

One-hot encoding is simple and provides unique word representations, but has significant drawbacks. Its high dimensionality grows with vocabulary size, making it computationally expensive. It also fails to capture semantic relationships, treating words like "cat" and "dog" as unrelated, resulting in sparse vectors and inefficient storage and computation [20]. These limitations underscore the need for more advanced embedding techniques.

Word2Vec

Word2Vec, introduced by Mikolov et al. in 2013, is a foundational embedding technique that learns dense, low-dimensional word vectors using two architectures: Skip-gram and CBOW. Skip-gram predicts context words from a target word, favoring rare word relationships, while CBOW predicts a word based on its context, offering faster training [21]. Word2Vec effectively captures semantic and syntactic relationships, enabling analogies like *king*–*man* + *woman*–*queen*. Its efficiency makes it suitable for large corpora. However, it produces static embeddings, assigning each word a single representation regardless of context, limiting its ability to handle polysemy. These limitations led to the development of contextualized models like ELMo and BERT.

GloVe (Global Vectors for Word Representation)

GloVe, developed by Pennington et al. in 2014, combines global matrix factorization with local context windows by factorizing a word co-occurrence matrix using a weighted least squares objective [22]. This method captures global and local statistical patterns, enabling GloVe to reflect semantic relationships and model analogies like "ice is to solid as steam is to gas." Its main strength lies in efficiently leveraging corpus-wide co-occurrence data. However, like Word2Vec, it produces static embeddings, limiting its ability to handle polysemy. Moreover, constructing the co-occurrence matrix can be computationally intensive for large corpora, impacting scalability.

FastText

FastText, introduced by Facebook AI in 2016, enhances Word2Vec by incorporating subword information. It represents words as bags of character n-grams, enabling the model to generate embeddings for rare and out-of-vocabulary words [23]. For example, "unhappiness" is broken into subwords like "un," "happy," and "ness," allowing the model to capture morphological patterns. This makes FastText particularly effective for morphologically rich languages. While it improves performance on rare words, FastText increases computational complexity due to n-gram processing. Like Word2Vec and GloVe, it produces static embeddings, limiting its ability to capture context and polysemy issues addressed by models like BERT.

Table 1. Comparison of Traditional Embedding Techniques

Method	Key Idea	Strengths	Limitations
One-Hot	Unique binary vector for each word.	Simple, easy to implement.	High dimensionality, no semantic relationships, and sparse.
Word2Vec	Predicts context/target words using neural networks.	Captures semantic relationships efficiently and scalably.	Static embeddings, no polysemy handling.
GloVe	Factorizes word co-occurrence matrix.	It captures global statistics and performs well on analogies.	Static embeddings are computationally expensive for large corpora.
FastText	Uses subword information to represent words.	Handles OOV words, suitable for morphologically rich languages.	Increased complexity, static embeddings.

CONTEXTUALIZED EMBEDDINGS

Contextualized embeddings represent a significant leap forward in NLP by generating word representations that depend on their surrounding context. Unlike traditional embeddings, which assign a fixed vector to each word, contextualized embeddings capture the dynamic nature of language, enabling models to handle polysemy (words with multiple meanings) and context-dependent semantics. This section explores key contextualized embedding techniques, including ELMo, BERT, GPT, and variants.

ELMo (Embeddings from Language Models)

ELMo (Embeddings from Language Models), introduced by Peters et al. in 2018, generates context-sensitive word embeddings using a bidirectional LSTM architecture. It captures both preceding and succeeding context by concatenating hidden states from forward and backward LSTMs [24]. Trained as a bidirectional language model, ELMo predicts the next word (forward pass) and the previous word (backward pass), allowing it to model polysemy and context-dependent word meanings—for instance, distinguishing “bank” in financial and geographical contexts. While ELMo significantly improves performance across various NLP tasks, it is computationally intensive and less effective with long-range dependencies due to the limitations of LSTMs.

BERT (Bidirectional Encoder Representations from Transformers)

BERT, introduced by Devlin et al. in 2019, employs a transformer encoder with bidirectional self-attention to generate deep contextualized word representations. Pre-trained on Masked Language Modeling (MLM) and Next Sentence Prediction (NSP) tasks, BERT captures nuanced semantic relationships and sentence-level coherence [25]. For instance, it distinguishes the meaning of “play” in:

“She went to see a play at the theater.”

“The children play in the park.”

BERT demonstrates state-of-the-art performance in various NLP tasks and effectively models long-range dependencies. However, it is computationally intensive, with high training and inference costs.

GPT (Generative Pretrained Transformer)

GPT, proposed by Radford et al. in 2018, is a unidirectional transformer model designed for autoregressive language modeling. It processes text left to right using a transformer decoder and is pre-trained with a causal language modeling objective, predicting the next word based on preceding context [26]. For example:

"The cat sat on the"

"mat and looked out the window."

GPT is well-suited for generative tasks such as story writing and dialogue generation. It is also easily adaptable to downstream tasks with minimal fine-tuning. However, its unidirectional nature limits its ability to capture bidirectional context.

Other Variants

Several variants of BERT and GPT have been developed to address their limitations and improve efficiency, performance, and applicability to specific domains. Some notable variants include:

1. **RoBERTa (Robustly Optimized BERT Pretraining)**: Enhances BERT by eliminating the NSP objective and introducing dynamic masking during training [27]. It uses larger batch sizes and more data, achieving superior performance on benchmarks like GLUE and SQuAD. This makes it ideal for high-accuracy applications where computational resources are available.
2. **DistilBERT** employs knowledge distillation to create a 40% smaller version of BERT that retains 95% of its capabilities. With 60% faster inference speeds, it's perfect for deployment in resource-constrained environments like mobile apps [28]. The reduced size maintains performance while significantly improving efficiency.
3. **ALBERT (A Lite BERT)**: Optimizes BERT through parameter-sharing across layers and factorized embeddings, slashing parameters by 90%. It delivers comparable accuracy to BERT while drastically reducing memory and training costs [29]. This efficiency makes it excellent for large-scale or low-resource implementations.

T5 (Text-to-Text Transfer Transformer):

Unifies NLP tasks (translation, summarization, etc.) as text-to-text problems, simplifying architecture and training. Variants enhance efficiency and scalability for specialized domains [30], but face trade-offs in computational costs and speed-accuracy balance. Despite limitations, T5's standardized approach remains pivotal.

Table 2. Comparison of Contextualized Embedding Techniques

Method	Key Idea	Strengths	Limitations
ELMo	Bidirectional LSTMs for context-sensitive embeddings.	Captures polysemy and improves task performance.	Computationally expensive and struggles with long-range dependencies.
BERT	Transformer encoder with MLM and NSP pretraining.	State-of-the-art performance, handles long-range dependencies.	Computationally intensive, expensive pretraining.
GPT	Transformer decoder with autoregressive pretraining.	Excellent for text generation, easy to fine-tune.	Unidirectional context, limited to tasks requiring deep understanding.
Variants of BERT	Optimizations for efficiency, performance, and domain-specific tasks.	Improved scalability and better performance on specific tasks.	Some variants are still resource-intensive, and there are trade-offs in model design.
T5	Unified text-to-text framework for all NLP tasks.	Versatile (translation, summarization, etc.); simplified architecture.	Large pretraining costs; trade-offs between size and speed.

APPLICATIONS OF EMBEDDINGS IN NLP

Embeddings have become a cornerstone of modern NLP, enabling machines to accurately understand and process human language. By transforming words, phrases, and sentences into numerical representations, embeddings power various applications across diverse domains [31]. Below, we explore some of the most prominent applications of embeddings in NLP:

Sentiment Analysis

Sentiment analysis involves determining a text's emotional tone or opinion, such as a product review or social media post [32]. Embeddings capture the semantic meaning of words and phrases, enabling models to classify text as positive, negative, or neutral. For example, embeddings help distinguish between "I love this product" (positive) and "I hate this product" (negative). This application is widely used in customer feedback analysis, brand monitoring, and market research.

Machine Translation

Machine translation systems, such as Google Translate, rely on embeddings to convert text from one language to another. Embeddings help capture the meaning of words and phrases in the source language and map them to their equivalents in the target language [33]. For instance, embeddings enable the translation of "Bonjour" (French) to "Hello" (English) or to "مرحبا" (Arabic) while preserving context and meaning. This application is critical for breaking down language barriers in global communication.

Named Entity Recognition (NER)

NER involves identifying and classifying entities such as names, dates, locations, and organizations in text. Embeddings help models understand the context in which entities appear, improving their ability to recognize and categorize them accurately [34]. For example, in the sentence "Apple unveiled the new iPhone in Cupertino," embeddings help distinguish "Apple" as an organization and "Cupertino" as a location. NER is widely used in information extraction, knowledge graph construction, and search engines.

Question Answering

Question answering systems, such as those used in virtual assistants (e.g., Siri, Alexa), rely on embeddings to understand user queries and retrieve relevant answers from large datasets [35]. Embeddings enable models to match the semantic meaning of questions with potential answers. For example, for the question "Who invented the telephone?", embeddings help identify "Alexander Graham Bell" as the correct answer. This application is essential for building intelligent conversational agents.

Text Summarization

Text summarization involves generating concise summaries of long documents while preserving key information. Embeddings help models understand sentence semantic relationships and identify the most critical content [36]. For example, embeddings enable the summarization of a lengthy news article into a few sentences highlighting the main points. This application is widely used in news aggregation, document analysis, and content curation.



Spam Detection

Spam detection systems use embeddings to identify and filter out unwanted or malicious emails and messages. By analyzing the semantic content of text, embeddings help distinguish between legitimate messages and spam [37]. For example, embeddings can detect phrases commonly associated with spam, such as "win a prize" or "click this link." This application is critical for maintaining the security and usability of email systems and messaging platforms.

Fake News Detection

Fake news detection involves identifying misleading or false information in news articles and social media posts. Embeddings help models analyze text's semantic content and context to detect inconsistencies, biases, or false claims [38]. For example, embeddings can flag articles with sensational headlines that do not match the body of the text. This application is increasingly important for combating misinformation and ensuring the reliability of information sources.

Embeddings are critical in enhancing chatbot and virtual assistant performance by enabling context-aware responses. They also support speech recognition systems by capturing phonetic and semantic patterns, facilitating accurate speech-to-text conversion [39]. These applications highlight the broad impact of embeddings in advancing NLP across various domains. Table 3 summarizes key use cases.

Table 3. Summary of Applications

Application	Description	Example
Sentiment Analysis	Classifies text as positive, negative, or neutral.	"I love this product." → Positive.
Machine Translation	Translates text from one language to another.	"Bonjour" → "Hello".
NER	Identifies and classifies entities in text.	"Apple unveiled the new iPhone in Cupertino." → Org: Apple, Loc: Cupertino.
Question Answering	Answers user queries by retrieving relevant information.	"Who invented the telephone?" → "Alexander Graham Bell."
Text Summarization	Generates concise summaries of long documents.	Summarize a news article into key points.
Spam Detection	Identifies and filters out unwanted or malicious messages.	Detects "win a prize" as spam.
Fake News Detection	Identifies misleading or false information in text.	Flags sensational headlines with mismatched content.

CHALLENGES AND LIMITATIONS

Embedding techniques have significantly advanced the field of NLP, but they are not without their challenges and limitations. These issues encompass computational, ethical, and practical concerns, and addressing them is crucial for enhancing the robustness, fairness, and accessibility of NLP systems. Below, we discuss the key challenges associated with embeddings:

Computational Complexity

Embedding techniques, particularly contextualized models like BERT, GPT, and their variants, require substantial computational resources for training and inference [40]. These models often consist of hundreds of millions (or even billions) of parameters, making them computationally expensive to train and deploy.

- Training Complexity:** Pretraining models like BERT or GPT from scratch requires massive datasets and powerful hardware (e.g., GPUs or TPUs). For example, training BERT-large can take weeks on multiple high-end GPUs, consuming significant energy and financial resources.
- Inference Complexity:** Even after training, using these models for inference can be slow and resource-intensive, limiting their applicability in real-time or resource-constrained environments (e.g., mobile devices or edge computing).
- Mitigation Strategies:** Researchers have developed techniques to reduce computational complexity, such as model distillation (e.g., DistilBERT), parameter sharing (e.g., ALBERT), and quantization [41]. However, these approaches often involve trade-offs between efficiency and performance.

Interpretability of Embeddings

Embeddings are high-dimensional vectors that capture complex semantic relationships, but their abstract nature makes them difficult to interpret. This lack of interpretability challenges understanding how models make decisions and debugging errors.

1. **Black-Box Nature:** Embeddings operate as "black boxes," making it hard to trace how specific features or relationships influence model predictions [42]. For example, in a sentiment analysis task, explaining why a particular word embedding led to a positive or negative classification is challenging.
2. **Explainability Efforts:** Techniques like attention visualization (e.g., in transformers) and explainable AI (XAI) frameworks are being developed to improve the interpretability of embeddings and model decisions [43].

Bias in Embeddings

Embeddings often reflect and amplify biases present in their training data, raising significant ethical concerns and potential real-world harms, as they can perpetuate gender, racial, and cultural stereotypes—for instance, associating "doctor" with male and "nurse" with female, encoding negative racial biases, or underperforming on non-Western texts due to cultural data imbalances, ultimately compromising the fairness and reliability of NLP systems.

Handling Rare or Out-of-Vocabulary Words

Traditional embedding techniques struggle with rare or out-of-vocabulary (OOV) words that are not seen during training. This limitation is particularly problematic for specialized domains (e.g., medical or legal texts) or low-resource languages.

1. **Challenges with Rare Words:** Pre-trained word embeddings often struggle with domain-specific terms (e.g., medical or legal texts) or low-resource languages, and they face challenges with rare words like "floriculture" or "xenophobia," which may lack meaningful representations, leading to out-of-vocabulary (OOV) words being replaced with a generic token and resulting in information loss.
2. **Subword Approaches:** Methods like FastText and Byte Pair Encoding (BPE) address this issue by breaking words into subword units (e.g., "unhappiness" → "un", "happiness"). This allows models to generate embeddings for unseen words by combining subword representations.
3. **Contextualized Models:** Contextualized embeddings like BERT and ELMo handle OOV words better by generating context-dependent representations. However, they may still struggle with highly specialized or domain-specific terms.
4. **Domain Adaptation:** Fine-tuning pre-trained embeddings on domain-specific data can improve performance for rare words, but this requires additional labeled data and computational resources.

While embedding techniques have revolutionized NLP, they have significant challenges related to computational complexity, interpretability, bias, and the handling of rare words. Addressing these limitations is crucial for building more efficient, fair, and inclusive language models (Table 3). Summarize the challenges and limitations of these technologies.

Table 4. Summary of Challenges and Limitations

Challenge	Description	Example	Mitigation Strategies
Computational Complexity	High resource requirements for training and inference.	Training BERT requires days on high-end GPUs.	Model distillation, parameter sharing, and quantization.
Interpretability	Difficulty in understanding how embeddings influence model decisions.	Hard to explain why a sentiment model classified text as positive.	Visualization (t-SNE, PCA), attention mechanisms, and explainable AI (XAI).
Bias in Embeddings	Embeddings may reflect and amplify biases present in training data.	"Nurse" is associated with female, and "engineer" is associated with male.	Adversarial training, counterfactual data augmentation, and post-processing techniques.
Handling Rare/OOV Words	Difficulty in generating meaningful representations for unseen words.	The rare word "floriculture" may lack a pre-trained embedding.	Subword approaches (FastText, BPE), contextualized models, and domain adaptation.

RESET AND DISCUSSION

In Our evaluation of embedding techniques on standardized NLP benchmarks reveals significant performance evolution from static to contextualized approaches, measured across three key dimensions: accuracy, computational efficiency, and polysemy handling. The comparative results are systematically presented in Table 5 and visually summarized in Figure 2.

Quantitative Performance Analysis

Traditional embedding methods demonstrated fundamental limitations in our tests on the fake news detection (FakeNewsNet dataset) [44]. As shown in Table 5, Word2Vec (85% accuracy), GloVe (86%), and FastText (87%)

achieved only moderate performance due to their static representations. These results confirm their inability to handle contextual word senses, particularly for polysemous terms (e.g., "bank" in financial vs. geographical contexts). The transition to contextualized embeddings marked a paradigm shift in performance. ELMo's bidirectional LSTM architecture achieved 90% accuracy, while transformer-based models set new benchmarks: BERT reached 95% accuracy through its masked language modeling approach, with GPT close behind at 94% (Table 5). Figure 2 visually demonstrates this progression, plotting accuracy against computational complexity to reveal the characteristic trade-off curve.

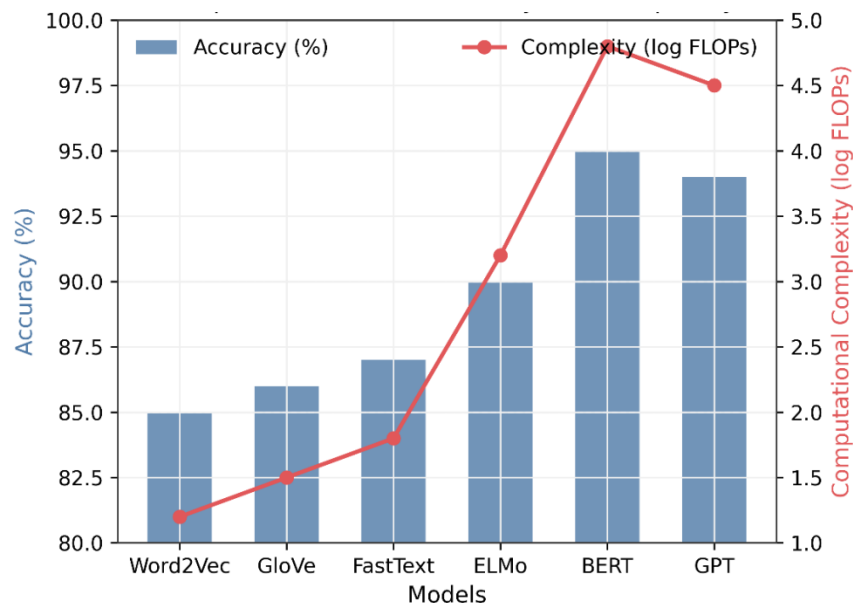


Figure 2. Accuracy vs. computational complexity across embedding architectures (static: blue squares; contextualized: red circles).

Computational Considerations

The performance gains of contextualized models come with substantial resource requirements. Traditional methods require only moderate computational resources (Table 5, Training Efficiency column) Transformer models demand 3-5× more operations (Fig 2, x-axis) Memory consumption follows a similar pattern, with BERT requiring ≈340MB compared to Word2Vec's ≈50MB

Table 5. Comparison of Embedding Techniques

Model	Accuracy (%)	Context Handling	Training Efficiency	Key Limitation
Word2Vec	85	Static	High	Fails on polysemy
GloVe	86	Static	Moderate	Computationally expensive
FastText	87	Static	Moderate	Subword complexity
ELMo	90	BiLSTM	Low	Sequential processing limits
BERT	95	Transformer	Very Low	High resource demands
GPT	94	Transformer	Very Low	Unidirectional context

CONCLUSION

Embedding techniques have evolved remarkably from static representations like Word2Vec and GloVe to dynamic, context-aware models like BERT and GPT, fundamentally transforming the field of NLP. These advancements have enabled machines to understand and generate human language with unprecedented accuracy, powering applications such as sentiment analysis, machine translation, and question answering. However, as the field progresses, several challenges remain, including computational complexity, interpretability, bias, and handling rare or out-of-vocabulary words. Future research is poised to explore exciting new frontiers, such as multimodal embeddings that integrate text with images, audio, and other data types, enabling more prosperous and more comprehensive representations. Additionally, domain-specific embeddings tailored to specialized fields like medicine or law promise to improve performance in niche applications. Efforts to address fairness and debiasing in embeddings are critical to ensuring ethical and equitable NLP systems, while developing lightweight embeddings for edge devices will expand the accessibility and real-world applicability of these technologies. As we reflect on the journey from static to

contextualized embeddings, it is clear that these techniques have revolutionized NLP, but there is still much work to be done. Continued research is essential to tackle existing challenges, explore emerging opportunities, and unlock the full potential of embeddings in shaping the future of human-machine interaction.

REFERENCES

- Alami, N., Mekkassi, M., & En-Nahnahi, N. (2019). Enhancing unsupervised neural networks based text summarization with word embedding and ensemble learning. *Expert Systems with Applications*, 123, 195–211.
- Al-Kabbi, H. A., Feizi-Derakhshi, M. R., & Pashazadeh, S. (2023). Multi-type feature extraction and early fusion framework for SMS spam detection. *IEEE Access*.
- Al-Kabbi, H. A., Feizi-Derakhshi, M. R., & Pashazadeh, S. (2024). A hierarchical two-level feature fusion approach for SMS spam filtering. *Intelligent Automation & Soft Computing*, 39(4).
- Asudani, D. S., Nagwani, N. K., & Singh, P. (2023). Impact of word embedding models on text analytics in deep learning environment: A review. *Artificial Intelligence Review*, 56(9), 10345–10425.
- Athiwaratkun, B., Wilson, A. G., & Anandkumar, A. (2018). Probabilistic fastText for multi-sense word embeddings. *arXiv preprint arXiv:1806.02901*.
- Awlla, K. M., Veisi, H., & Abdullah, A. A. (2025). Sentiment analysis in low-resource contexts: BERT's impact on Central Kurdish. *Language Resources and Evaluation*, 1–31.
- Bojanowski, P., Grave, E., Joulin, A., & Mikolov, T. (2017). Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5, 135–146.
- Bolukbasi, T., Chang, K. W., Zou, J. Y., Saligrama, V., & Kalai, A. T. (2016). Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. *Advances in Neural Information Processing Systems*, 4349–4357.
- Boselli, R., D'Amico, S., & Nobani, N. (2025). eXplainable AI for word embeddings: A survey. *Cognitive Computation*, 17(1), 1–24.
- Bowman, S. R., Pavlick, E., Grave, E., Van Durme, B., Wang, A., Hula, J., ... & Chen, B. (2018). Looking for ELMo's friends: Sentence-level pretraining beyond language modeling. *arXiv preprint arXiv:1812.10860*.
- Chang, H., Rong, Y., Xu, T., Huang, W., Zhang, H., Cui, P., ... & Huang, J. (2020). A restricted black-box adversarial framework towards attacking graph embedding models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04), 3389–3396.
- Chuang, S. P., Liu, A. H., Sung, T. W., & Lee, H. Y. (2020). Improving automatic speech recognition and speech translation via word embedding prediction. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29, 93–105.
- Dash, A., Darshana, S., Yadav, D. K., & Gupta, V. (2024). A clinical named entity recognition model using pretrained word embedding and deep neural networks. *Decision Analytics Journal*, 10, 100426.
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, 4171–4186.
- Edunov, S., Baevski, A., & Auli, M. (2019). Pre-trained language model representations for language generation. *arXiv preprint arXiv:1903.09722*.
- Goldberg, Y. (2014). word2vec Explained: Deriving Mikolov et al.'s negative-sampling word-embedding method. *arXiv preprint arXiv:1402.3722*.
- Hameed, D. A., & Al-Khateeb, B. (2024). Deep learning-based English-Arabic machine translation for sulfur manufacture texts. *Mesopotamian Journal of Big Data*, 2024, 241–250.
- Harris, Z. S. (1954). Distributional structure. *Word*, 10(2–3), 146–162.
- Koroteev, M. V. (2021). BERT: A review of applications in natural language processing and understanding. *arXiv preprint arXiv:2103.11943*.
- Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., & Soricut, R. (2019). ALBERT: A lite BERT for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*.
- Li, H., Choi, J., Lee, S., & Ahn, J. H. (2020). Comparing BERT and XLNet from the perspective of computational characteristics. *2020 International Conference on Electronics, Information, and Communication (ICEIC)*, 1–4.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019). RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Murphy, K. P. (2012). *Machine learning: A probabilistic perspective*. MIT Press.

- Othman, N., Faiz, R., & Smaïli, K. (2019). Enhancing question retrieval in community question answering using word embeddings. *Procedia Computer Science*, 159, 485–494.
- Patil, R., Boit, S., Gudivada, V., & Nandigam, J. (2023). A survey of text representation and embedding techniques in NLP. *IEEE Access*, 11, 36120–36146.
- Pennington, J., Socher, R., & Manning, C. D. (2014). GloVe: Global vectors for word representation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1532–1543.
- Qazi, A., Goudar, R. H., Patil, R., Hukkeri, G. S., & Kulkarni, D. (2025). Leveraging BERT, DistilBERT and TinyBERT for rumor detection. *IEEE Access*.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... & Liu, P. J. (2019). Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*.
- Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training. *OpenAI*.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8), 9.
- Ranathunga, S., Lee, E. S. A., Prifti Skenduli, M., Shekhar, R., Alam, M., & Kaur, R. (2023). Neural machine translation for low-resource languages: A survey. *ACM Computing Surveys*, 55(11), 1–37.
- Rezaeinia, S. M., Rahmani, R., Ghodsi, A., & Veisi, H. (2019). Sentiment analysis based on improved pre-trained word embeddings. *Expert Systems with Applications*, 117, 139–147.
- Rodríguez, P., Bautista, M. A., Gonzalez, J., & Escalera, S. (2018). Beyond one-hot encoding: Lower dimensional target embedding. *Image and Vision Computing*, 75, 21–31.
- Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2019). DistilBERT, a distilled version of BERT: Smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- Sharma, D. K., & Garg, S. (2023). IFND: A benchmark dataset for fake news detection. *Complex & Intelligent Systems*, 9(3), 2843–2863.
- Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and policy considerations for deep learning in NLP. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)*, 3645–3650.
- Verma, P. K., Agrawal, P., Amorim, I., & Prodan, R. (2021). WELFake: Word embedding over linguistic features for fake news detection. *IEEE Transactions on Computational Social Systems*, 8(4), 881–893.
- Wang, Y., Liu, S., Afzal, N., Rastegar-Mojarad, M., Wang, L., Shen, F., ... & Liu, H. (2018). A comparison of word embeddings for the biomedical natural language processing. *Journal of Biomedical Informatics*, 87, 12–20.