# Exploring Word Embedding Tools: GloVe, FastText, Word2Vec and BERT

**Pinaki Sahu[1]**

1 IIPP Research Intern Asia University, Taichung, Taiwan

(e-mail:0000pinaki1234.kv@gmail.com).

⋮ ABSTRACT Encoding tools have transformed natural language processing by improving the understanding and Use of textual content by machines. This article examines four well-known word processing models: Glove, FastText, Word2Vec, and BERT. Each of these tools has its own unique benefits and applications, making them necessary for NLP's various responsibilities. We delve into their fundamental principles, identify their core abilities, and explain their practical applications. These models not only made it possible for machines to understand sophisticated language, but they also revolutionized language-driven engineering. As NLP continues to evolve, these models continue to be at the forefront, enabling deeper human-machine interactions.

⋮ KEYWORDS Word embeddings, GloVe, FastText, Word2Vec, BERT, Natural Language Processing (NLP)

## 1. INTRODUCTION

In a time when texts on digital platforms from social media reports and academic research have the capacity to analyze both the data and the time, the ability to analyze both of these things is essential. [1]. In addition to this, it aims at reducing the gap in terms of computer capability. To be more specific, it aims to reduce the gap between natural language and artificial language. Word input plays a key role in natural language processing (NLP), a new technology that changes the way machines learn and understand human language Human language is a complex set of symbols, namely context, culture internal effects are closely related to cognitive aspects In order for machines to recognize language and process it correctly to represent words and phrases that can be understood. Natural language processing has come a long way quickly because of to word embeddings. These models not only helped computers understand language better, but they also made it possible for language-driven AI to do new things.

Here are four of the most common word embedding models that we will talk about: These are GloVe, FastText, Word2Vec, and BERT.
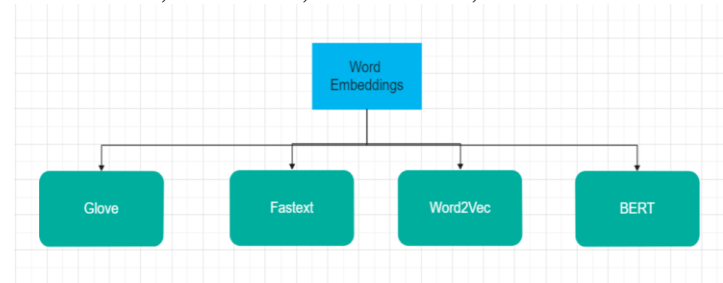


Fig-1 Word Embeddings Types

## 2. Related Works

**1. Semantic analysis:** Word entry is often used for semantic analysis, which includes finding synonyms, evaluating synonyms, and looking for semantic links between words, among other things. Researchers are always finding better ways to get important semantic information from word vectors [2]. Deeper semantic analysis, and the ability to capture and quantify subtleties, paves the way for more accurate and sophisticated understanding of language, pushing the development of many NLP applications

2. Multilingual and multilingual encoding: Multilingual encoding enables the representation of words from different languages in a common vector space, facilitating multilingual NLP tasks.

Current research [3] focuses on methods for generating effective multilingual embeddings and cross-language transfer learning. The search of a universal language representation is not only a matter of linguistic balance but also a means of facilitating improved global communication.

**3.Bias Mitigation**: Word embeddings typically inherit biases from training data. To ensure that AI systems are fair and neutral in their language understanding and generation, Bias mitigation is not only a technical challenge, but also a moral imperative that highlights the responsibility of AI in a world that is becoming increasingly diverse and interconnected researchers are developing techniques to reduce these biases [4].

**4.Zero-Shot and Few-Shot Learning:** Emerging areas of research include zero-shot and few-shot learning with word embeddings. These techniques allow AI systems to perform tasks with minimal training data, thereby creating new opportunities for adaptability and versatility. The utilization of word embeddings, in particular, contextual embeddings like BERT, is opening new opportunities for zero-shot and few-shot learning, allowing AI systems to promptly adapt to new domains and scenarios [5].

## 3.Word Embeddings

### 3.1 Word2Vec

Google's Word2Vec has been an essential model in the field of word embeddings, changing how we perceive and represent words. Continuous Bag of Words (CBOW) and Skip-gram, each with their own unique principles and applications, are used by Word2Vec.

**1.CBOW: Continuous Bag of Words**

CBOW's task is to guess a target word by looking at the words that are around it. Based on the words within the target word, it guesses how likely it is to appear.

The model's goal is to find the word vector that increases the conditional chance of the target word to the highest level given the situation.

**2.Skip Grams**
Whereas the skip-gram model reverses the objective. This tool provides contextual terms that are associated with a given target word. The primary aim of the model is to optimize the probability of context words in relation to the target word..

The utilization of Word2Vec is widespread in several natural language processing (NLP) applications due to its ability to represent words inside a continuous vector space, where words with similar meanings are situated in close proximity. Here are a few uses where Word2Vec excels:

**1. Document Grouping:** Utilizing Word2Vec improves the process of clustering together similar documents. Document content can be seen as a vector in this context by aggregating the word vectors. Document clustering is a form of unsupervised learning in which documents with similar information are grouped together.

**2. Response to a Question:** Word2Vec is essential in enhancing the efficacy of question-answering systems. It gives these systems the ability to comprehend the semantic relationships between sentences and paragraphs [6].

### 3.2 Glove

The creation of Glove (Global Vectors for Word Representation) at Stanford University is a significant advancement in the field of word embedding models. Its primary objective is to discover complex global word-word relationships by analyzing large text corpora.

**Co-occurrence Statistics:**

GloVe computes the probability of co-occurrence of words within a particular context window. The co-occurrence of words shows crucial information about their semantic relationships.

Glove's flexibility extends to various natural language processing (NLP) applications, where its

ability to understand the semantic relationships between words is of the greatest importance. Here are a few applications where GloVe really excels:

**1. Sentiment Analysis**: Sentiment analysis, a crucial endeavor in NLP, relies on understanding the emotional tone of the text. GloVe supports sentiment analysis models in determining the underlying sentiment by supplying them with semantically dense word embeddings. Glove's word vectors encode contextual information that aids in recognizing nuances in sentiment, making them useful for responsibilities such as product reviews, social media sentiment analysis, and customer feedback analysis.

**2. Machine Translation:** Machine translation systems heavily depend on word embeddings to convert text from one language to another. The embeddings of GloVe allow these systems to map source-language words to their semantically equivalent target-language counterparts. GloVe helps enhance the quality and accuracy of machine translation models by identifying the global semantic relationships between words.

GloVe improves the precision and relevance of information retrieval by aligning the semantic space of queries and documents [6].

## 3.3 FastText

The innovative word embedding model created by Facebook AI Research is FastText. FastText surpasses conventional word embeddings by providing a novel approach that includes subworld information to improve its capabilities [6].

**Subword Representation:**

FastText stands out by its representation of words as the sum of their character n-grams. This sub word representation enables FastText to capture linguistic information at a granular level within words. The model treats words as a "bag of character n-grams," which means it considers the constituent elements at the character level. This

method is particularly useful for managing non-vocabulary words and morphological variants.

FastText is an essential tool in the field of NLP due to its innovative approach to sub word information, which enables a variety of applications[7].

**1.Multilingual NLP:** FastText's capabilities are truly beneficial to NLP systems that support multiple languages. The ability to adapt to languages with complex grammar enables the development of multilingual models capable of handling a variety of linguistic structures.

**2. Grammar Correction:** The sub word information provided by FastText is invaluable for grammar correction applications. It can suggest word replacements based on similarity at the character level, making it an effective tool for automatically correcting typographical errors in text.

## 3.4 BERT

An new model based on transformers was made by Google AI called BERT (Bidirectional Encoder Representations from Transformers). BERT changed the field of natural language processing (NLP) forever when it came up with contextual embeddings, which look at the whole sentence context of a word [8].

**Pre-training and Fine-Tuning:**

BERT undergoes pre-training using extensive quantities of textual material, which facilitates the acquisition of a profound understanding of linguistic nuances. During this phase, the system is taught to predict absent words in sentences, a process known as hidden language modeling.

BERT can subsequently be fine-tuned for specific following tasks by adapting its pre-trained knowledge to various NLP applications.

**1.Sentiment Analysis and Text Classification:** BERT is highly effective in sentiment analysis and text classification tasks due to its ability to

understand word context. It is capable of understanding the subtle nuances of text, which is necessary for distinguishing sentiment and categorizing text into relevant categories**.**

**2. Response to a Question:** The contextual awareness of BERT excels in question-answering duties. It can analyze the context and relationships between words in order to provide precise answers to natural language queries.

**3.Language Translation:** The contextual embeddings developed by BERT have improved machine translation systems. By more accurately capturing sentence structures and idiomatic expressions, BERT has contributed to improvements in the quality of language translation, thus improving cross-lingual communication.

**4. Healthcare Information Technology**: In healthcare, the contextual awareness of BERT is crucial for medical record analysis. It assists in clinical decision support, patient history extraction, and disease identification by understanding the complex medical language and context, ultimately contributing to enhanced patient care and health outcomes.

## 4.Challenges in word embeddings

**1.Size of Data and Pretraining:** Models such as BERT require huge amounts of data and extensive pretraining, which can be computationally costly. Smaller organizations and initiatives may struggle for the necessary resources.

**2.Interpretability:** It can be difficult to interpret the inner workings of deep learning models such as BERT, raising concerns in applications where clarity is crucial, such as medical and legal fields.

**3.Overfitting**: There is a risk of overfitting when fine-tuning pre-trained models like BERT on limited datasets for certain tasks. Maintaining a balance between model complexity and generalization remains crucial.

**4.Support for Multiple Languages:** While these models perform well for many languages, achieving the same level of performance across all languages remains a challenge, particularly for languages with limited training data.

## 5.Conclusion

The area of natural language processing has been changed by word embedding technologies like Glove, FastText, Word2Vec, and BERT. Each model has its own assets and applications, allowing NLP professionals to more effectively solve a wide variety of problems.

The choice of word embedding tool depends on the specific mission, available language, and data. When selecting the most appropriate tool for their NLP projects, researchers and engineers must carefully consider the model's principles, strengths, and limitations.

As NLP moves forward, these embedding models will no doubt evolve, helping machines understand content and communicate more intelligently and effectively, and ultimately it will change how we interact with technology. These developments will affect how we interact with information technology in our daily lives, and develop a deeper understanding between humans and machines.

## 5.References

[1] Lebret, R. P. (2016). Word embeddings for natural language processing (No. THESIS). EPFL.

[2]Yu, L. C., Wang, J., Lai, K. R., & Zhang, X. (2017, September). Refining word embeddings for sentiment analysis. In Proceedings of the 2017 conference on empirical methods in natural language processing (pp. 534-539).

[3]Ruder, S., Vulić, I., & Søgaard, A. (2019). A survey of cross-lingual word embedding models. Journal of Artificial Intelligence Research, 65, 569-631.

[4] Rahman, S., Khan, S., & Porikli, F. (2018). A unified approach for conventional zero-shot,

generalized zero-shot, and few-shot learning. IEEE Transactions on Image Processing, 27(11), 5652-5667.

[5]Papakyriakopoulos, O., Hegelich, S., Serrano, J. C. M., & Marco, F. (2020, January). Bias in word embeddings. In Proceedings of the 2020 conference on fairness, accountability, and transparency (pp. 446-457).

[6] Dharma, E. M., Gaol, F. L., Warnars, H. L. H. S., & Soewito, B. (2022). The accuracy comparison among word2vec, glove, and fasttext towards convolution neural network (cnn) text classification. J Theor Appl Inf Technol, 100(2), 31.

[7] I. Santos, N. Nedjah and L. de Macedo Mourelle, "Sentiment analysis using convolutional neural network with fastText embeddings," 2017 IEEE Latin American Conference on Computational Intelligence (LA-CCI), Arequipa, Peru, 2017, pp. 1-5, doi: 10.1109/LA-CCI.2017.8285683.

[8]Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.

[9] Zhang, J., et al. (2021). A secure decentralized spatial crowdsourcing scheme for 6G-enabled network in box. IEEE Transactions on Industrial Informatics, 18(9), 6160-6170.

[10]Shankar, K., et al.(2021). Synergic deep learning for smart health diagnosis of COVID-19 for connected living and smart cities. ACM Transactions on Internet Technology (TOIT), 22(3), 1-14.

[11]Prathiba, S. B.,et al.. (2021). SDN-assisted safety message dissemination framework for vehicular critical energy infrastructure. IEEE Transactions on Industrial Informatics, 18(5), 3510-3518.

[12]Gaurav, A., et al.(2022). A comprehensive survey on machine learning approaches for malware detection in IoT-based enterprise information system. Enterprise Information Systems, 1-25.