

CZECH TECHNICAL UNIVERSITY IN PRAGUE
FACULTY OF ELECTRICAL ENGINEERING
DEPARTMENT OF CYBERNETICS
CZECH INSTITUTE OF INFORMATICS, ROBOTICS AND CYBERNETICS



From FastText to Transformer Models, and their Application in Retrieval-Augmented Generation

Bachelor's Thesis

Yauheni Zviazdou

Prague, May 2024

Study programme: Open Informatics
Branch of study: Artificial Intelligence and Computer Science

Supervisor: Ing. Jan Šedivý, CSc.

Acknowledgments

Firstly, I would like to express my gratitude to my supervisor.

I. Personal and study details

Student's name: **Zviazdou Yauheni** Personal ID number: **507333**
Faculty / Institute: **Faculty of Electrical Engineering**
Department / Institute: **Department of Cybernetics**
Study program: **Open Informatics**
Specialisation: **Artificial Intelligence and Computer Science**

II. Bachelor's thesis details

Bachelor's thesis title in English:

From FastText to Transformer Models, and their Application in Retrieval-Augmented Generation

Bachelor's thesis title in Czech:

Od FastText k Transformer model m a jejich aplikace v Retrieval-Augmented generování

Guidelines:

Review the representation of words, sentences, and paragraphs, progressing from traditional methods like FastText to advanced transformer-based models such as BERT. The primary focus is to evaluate selected representations using analogy tests and confusion matrix. Use the UPV corpus set for evaluation.

In the second part of the study, emphasis will shift towards selecting optimal representations for Retrieval-Augmented Generation (RAG) algorithms. The investigation will determine the most efficient embeddings and optimal text chunk size for question-answering tasks, particularly in the context of natural language answers generation from technical manuals. Conduct a comprehensive evaluation with a particular focus on suggesting an optimal representation model that balances factuality and CPU requirements.

Bibliography / sources:

- [1] FastText documentation and tutorials, <https://fasttext.cc/>
- [2] LangChain documentation, <https://js.langchain.com/docs>
- [3] Word2Vec tutorials, <https://www.tensorflow.org/text/tutorials/word2vec>

Name and workplace of bachelor's thesis supervisor:

Ing. Jan Šedivý, CSc. Big Data and Cloud Computing CIIRC

Name and workplace of second bachelor's thesis supervisor or consultant:

Date of bachelor's thesis assignment: **02.01.2024** Deadline for bachelor thesis submission: **24.05.2024**

Assignment valid until: **21.09.2025**

Ing. Jan Šedivý, CSc.
Supervisor's signature

prof. Dr. Ing. Jan Kybic
Head of department's signature

prof. Mgr. Petr Páta, Ph.D.
Dean's signature

III. Assignment receipt

The student acknowledges that the bachelor's thesis is an individual work. The student must produce his thesis without the assistance of others, with the exception of provided consultations. Within the bachelor's thesis, the author must state the names of consultants and include a list of references.

Date of assignment receipt

Student's signature

Declaration

I declare that presented work was developed independently, and that I have listed all sources of information used within, in accordance with the Methodical instructions for observing ethical principles in preparation of university theses.

Date
.....

Abstract

This thesis investigates the application of word and sentence embeddings in Retrieval-Augmented Generation (RAG) for factual Question Answering (QA) tasks using technical manuals. The study explores the effectiveness of traditional FastText embeddings and advanced transformer-based models like Bidirectional Encoder Representations from Transformers (BERT) in capturing semantic relationships within text. We evaluate the quality of these representations using analogy tests and confusion matrix analysis on the UPV corpus set.

Subsequently, we will select optimal representations for RAG algorithms and assess their impact on factual accuracy and computational efficiency during QA. By analyzing the performance with different text chunk sizes, we aim to identify the optimal configuration for factual RAG in technical domains. This research contributes to the field of Natural Language Processing (NLP) by providing insights into selecting effective representations that balance factual accuracy and computational efficiency for QA systems.

Keywords NLP, Word Embedding, Transformers, RAG, QA, Semantic Textual Similarity (STS)

Abstrakt

Tato práce zkoumá aplikaci embeddingů slov a vět v modelu Retrieval-Augmented Generation (RAG) pro úlohy Question Answering (QA) zaměřené na fakta s využitím technických příruček. Studie se zabývá účinností tradičních embeddingů FastText a pokročilých modelů založených na transformátoru, jako je Bidirectional Encoder Representations from Transformers (BERT), při zachycování sémantických vztahů v textu. Kvalitu těchto reprezentací hodnotíme pomocí analogových testů a analýzy confusion matrix na korpusu UPV. Následně vybereme optimální reprezentace pro algoritmy RAG a posoudíme jejich vliv na faktickou přesnost a výpočetní efektivitu během QA. Analýzou výkonu s různými velikostmi textových fragmentů se snažíme identifikovat optimální konfiguraci pro faktické RAG v technických oborech. Výzkum přispívá k oblasti zpracování přirozeného jazyka (Natural Language Processing (NLP)) tím, že poskytuje poznatky o výběru efektivních reprezentací, které vyvažují faktickou přesnost a výpočetní efektivitu pro systémy QA.

Klíčová slova NLP, Word Embedding, Transformátory, RAG, QA, Sémantická podobnost textu

Abbreviations

RAG Retrieval-Augmented Generation

NLP Natural Language Procession

STS Semantic Textual Similarity

QA Question Answering

BERT Bidirectional Encoder Representations from Transformers

CBOW Continuous Bag-of-Words

GloVe Global vectors

ML Machine Learning

NN Neural Network

mBERT Multilingual Bidirectional Encoder Representations from Transformers

MTEB Massive Text Embedding Benchmark

ALBERT A Lite BERT

SimCSE Simple Contrastive Learning of Sentence Embeddings

RetroMAE Retrieval-oriented Language Models Via Masked Auto-Encoder

mE5 Multilingual E5

LaBSE Language-Agnostic BERT Sentence Embedding

MLM Masked Language Modeling

XLM-R XLM-Roberta

NSP Next Sentence Prediction

BGE BAAI General Embeddings

GTE General Text Embedding

Contents

1	Introduction	1
1.1	Text representation	1
1.2	Evolution of text representation methods	1
1.3	Research objective	1
2	Literature Review	2
2.1	Traditional word embedding methods	2
2.1.1	Word2Vec	2
2.1.2	GloVe	2
2.1.3	FastText	3
2.2	Transformer-based models	3
2.3	TODO: REMOVE, STRUCTURE	4
3	Methodology	5
3.1	Evaluation process	5
3.1.1	Text Data Preparation	5
3.1.2	Datasets	5
3.1.3	Baseline	6
3.1.4	Chosen transformer models	7
3.2	Optimizing Text Representations for RAG in Technical QA	11
3.2.1	RAG Algorithm and Text Representation Reliance	11
3.2.2	Key factors	12
3.3	REMOVE: Methodology structure	13
4	Experiments and Results	14
4.0.1	Balanced models	15
5	Discussion	16
6	Conclusion	17
7	References	18
A	Appendix A	20

1 Introduction

1.1 Text representation

The human language, with its nuances and complexities, presents a significant challenge for machines to understand. Natural Language Processing (NLP) bridges this gap, and at its core lies the critical concept of text representation. This process acts as a translator, bridging the gap between the richness of text and the numerical language that machines understand. By effectively capturing the meaning within words and their relationships, text representation empowers NLP models to leverage machine learning's capabilities. From sentiment analysis to machine translation, this ability to represent meaning fuels the advancements in NLP, enabling machines to interact with and decipher human language with ever-increasing accuracy.

1.2 Evolution of text representation methods

NLP has undergone a significant transformation in its approach to text representation. Early methods, such as one-hot encoding, while simple to implement, suffered from limitations in efficiency due to dimensionality and sparsity issues.

Word embedding techniques (e.g., Word2Vec, Global vectors (GloVe), FastText) offered a significant improvement by capturing semantic relationships between words through high-dimensional word vectors. However, these techniques primarily focused on local context within a limited window, hindering their ability to capture complex relationships within sentences or documents.

The emergence of deep learning architectures, particularly transformer-based models like Bidirectional Encoder Representations from Transformers (BERT), revolutionized the field of text representation. These models allow to not only understand the meaning of individual words but also consider their interaction and context within a sentence or document.

1.3 Research objective

This research aims to evaluate the effectiveness of various word, sentence, and paragraph representations for their subsequent application in Retrieval-Augmented Generation (RAG) algorithms, with a specific focus on the domain of technical Question Answering (QA).

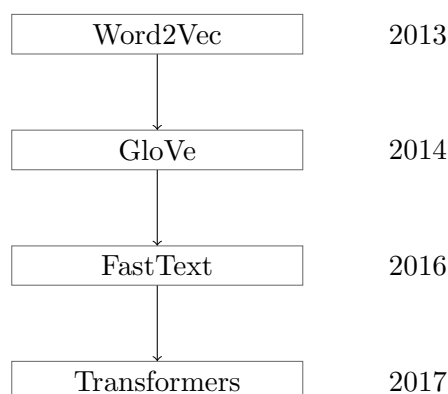


Figure 1.1: Evolution of the text representation methods.

■ 2 Literature Review

■ 2.1 Traditional word embedding methods

■ Word2Vec

Word2Vec [27] is algorithm that generates word embedding using information about target word (context). Word2Vec uses Neural Network (NN) and Machine Learning (ML) techniques to generate word embedding for every word in vocabulary during training. As NN architecture are used Continuous Bag-of-Words (CBOW) and Skip-gram, Fig. 2.1.

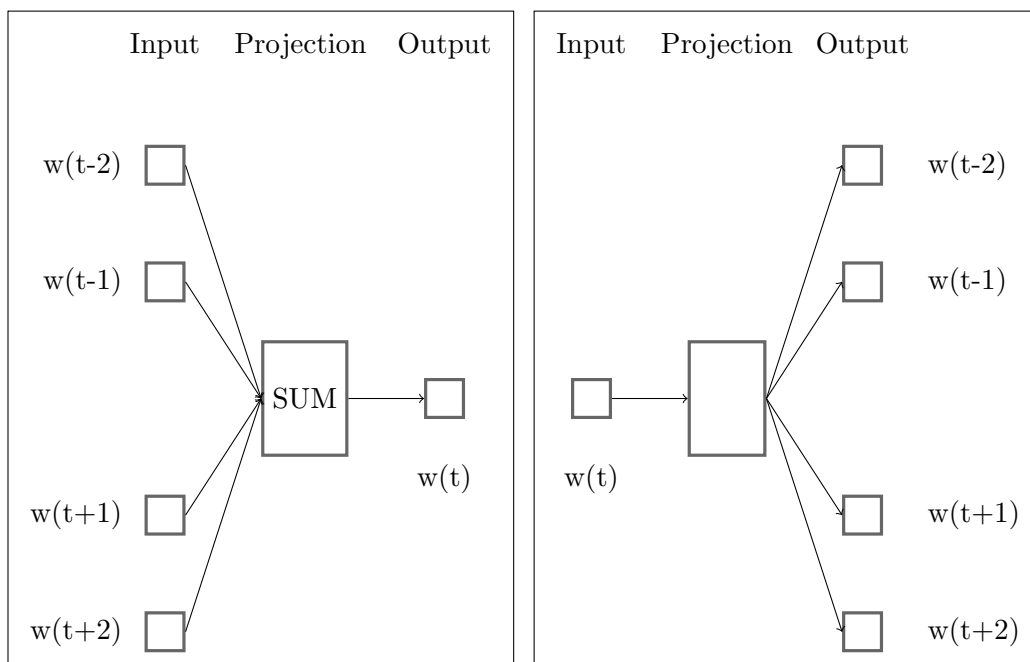


Figure 2.1: CBOW and Skip-gram schemes respectively

Due to its algorithmic simplicity and efficiency, Word2Vec has established itself as a strong baseline for numerous NLP tasks. Compared to more recent and complex models, Word2Vec requires minimal hyperparameter tuning, making it a relatively straightforward approach.

However, it is important to acknowledge that Word2Vec has limitations. These include its inability to capture **global information** within a document, its challenges in effectively handling **morphologically rich languages** (languages with many word variations), and its lack of awareness of the **broader context** beyond a limited window of surrounding words.

■ GloVe

GloVe [26] leverages the co-occurrence statistics of words within a corpus to learn vector representations. This approach involves constructing a co-occurrence matrix, where each entry reflects the frequency of two words appearing together within a predefined window size. This matrix essentially captures the relative importance of various word pairings.

A core principle of GloVe lies in the notion that word vectors should effectively encode the ratios between co-occurrence probabilities of words. By analyzing these ratios, GloVe can identify semantic relationships between words. This is achieved by factorizing the co-occurrence matrix into a lower-dimensional space, allowing for efficient representation and manipulation of word meanings.

To optimize the learned word embeddings, GloVe employs a weighted least squares objective function. This function aims to minimize the discrepancy between the dot product of two word vectors and the logarithm of their co-occurrence probability. Through iterative adjustments of the word vectors, GloVe converges on a solution that yields the desired word embeddings.

■ FastText

FastText [24] utilizes similar NN architectures as word2vec, namely CBOW and Skip-gram, but applies them to character n-grams (subwords) instead of entire words. This decomposition allows FastText to represent a word's meaning by considering its constituent subword components. Consequently, FastText offers advantages in two key areas:

- **Rare Word Embeddings:** Unlike word2vec, which struggles with words appearing infrequently in the training data, FastText can construct meaningful representations for rare words. By leveraging known subwords, FastText can represent unseen words, making it particularly valuable for working with large and diverse datasets.
- **Handling Morphologically Rich Languages:** Languages with complex morphology, where words are formed through prefixes and suffixes, often pose challenges for word2vec. FastText overcomes this limitation by capturing the shared subwords between derived words and their root forms. This allows FastText to represent the inherent relationships between words in these languages, leading to more accurate NLP tasks.

However, it's important to acknowledge that FastText also has limitations:

- **Context Insensitivity:** Similar to word2vec, FastText embeddings do not inherently capture the order or context in which words appear within a sentence. This can be a drawback for tasks like sentiment analysis or machine translation, where word order and context are crucial for accurate interpretation.
- **Limited Long-Range Dependency Capture:** While subwords allow FastText to capture local context, it might not effectively capture long-range dependencies within sentences. This can be a disadvantage for tasks requiring analysis of complex sentence structures, where understanding the relationships between words across larger distances is important.

■ 2.2 Transformer-based models

Transformer models [7] underpin powerful NLP models like BERT. A key advantage is their self-attention mechanism, which assigns importance to words based on context, not just position. This enables efficient parallel processing of entire sentences. Architecture of transformers visualized in Fig. 2.2

BERT builds on transformers with pre-training on a massive text corpus. Masked Language Modeling (MLM) and Next Sentence Prediction (NSP) further enhance BERT's capabilities, fostering deep contextual understanding and grasp of sentence relationships.

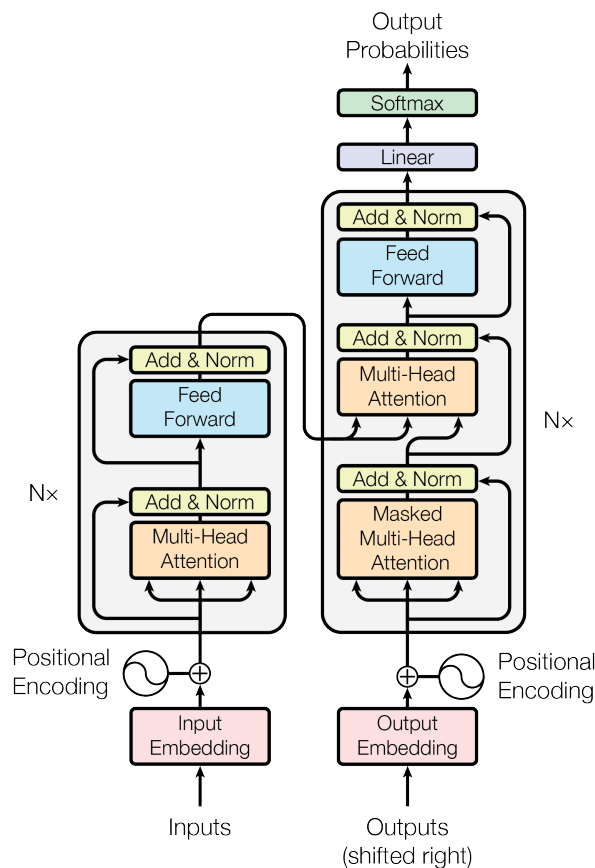


Figure 2.2: The Transformer - model architecture.

These strengths make transformers, particularly BERT, well-suited for NLP tasks. Their advantage lies in capturing contextual understanding, leading to richer text representations and superior comprehension of semantic relationships.

Furthermore, BERT excels in transfer learning, readily adapting to various tasks (sentiment analysis, QA) with minimal modifications. Additionally, efficiency and speed are benefits due to parallel processing and pre-training.

Transformer models effectiveness is validated by state-of-the-art performance across NLP benchmarks. Finally, BERT's robustness allows it to handle nuances in text without significant performance degradation.

■ 2.3 TODO: REMOVE, STRUCTURE

- Discuss traditional word embedding methods like FastText and their limitations.
- Explain the concept of transformer-based models like BERT and their advantages for text representation.
- Review related work on RAG algorithms and their dependence on effective text representations. Discuss existing research on evaluating text representations using analogy tests and confusion matrices.
- Briefly mention the UPV corpus set as the chosen evaluation benchmark.

■ 3 Methodology

■ 3.1 Evaluation process

This work specifically targets the evaluation of word, sentence, and paragraph representation methods on datasets in the Czech language. This focus on Czech allows for a deeper understanding of how these methods perform in a language with specific characteristics, such as a rich inflectional morphology and the presence of diacritics.

■ Text Data Preparation

In certain foreign languages, a common issue arises when individuals incorrectly write words by omitting diacritics or altering letters, Fig. 3.1. This problem is prevalent in social media, chatbots, and other informal written communications. As a result, embedding models face challenges in comprehending text without diacritics (hereinafter diacriticless). A potential solution involves adapting data representation to accommodate both formal and informal styles of writing.

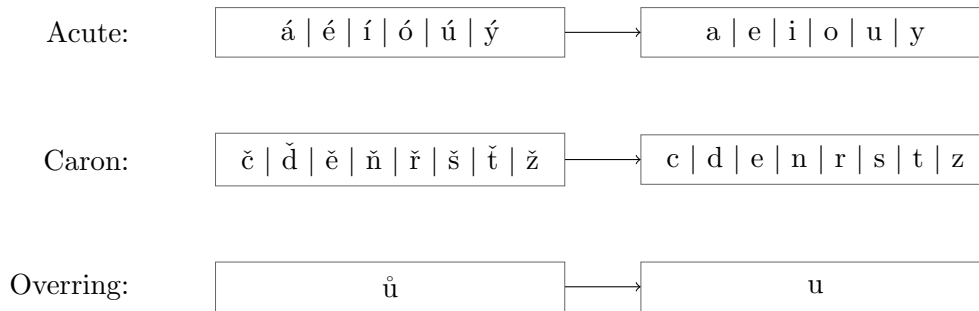


Figure 3.1: Usual changes in informal Czech texts.

This study will employ two distinct text representations: text with diacritics and text without diacritics. To ensure optimal evaluation, the diacritic text will be assessed using datasets that preserve these diacritics, while the diacriticless text will be evaluated using datasets that lack diacritics.

As detailed in Lst. 3.1, this script is used for creating diacriticless versions of the datasets.

```
sed 's/./\L&/' "$1" | iconv -f utf-8 -t ascii//TRANSLIT > diacriticless/"$1"
```

Listing 3.1: Script for removing diacritics using Unix utilities

■ Datasets

UPV FAQ

Dataset is comprised of frequently asked questions (FAQs) and their answers from the Industrial Property Office of the Czech Republic website, organized into four distinct groups, as detailed in Table 3.1.

This work will evaluate two key metrics using UPV FAQ dataset

- **Question matching accuracy:** This metric involves calculating the cosine similarity between all possible question pairs within a dataset. A question is considered successfully matched if its second-highest cosine similarity score corresponds to another question belonging to the same class (i.e., the question with the highest similarity is likely the same question itself). The overall question matching accuracy is then computed as the ratio of successfully matched questions to the total number of question pairs evaluated.
- **Answer matching accuracy:** This metric assesses the system’s ability to identify the correct answer for a given question. The system accomplishes this by directly comparing the question with pre-generated answer embeddings. By evaluating the similarity between the question and each answer embedding, the system classifies the question as corresponding to the answer with the highest similarity score. The overall answer matching accuracy is then calculated as the proportion of questions for which the system correctly identifies the corresponding answer.

Tests set	Number of tests
FAQv5	2054
FAQ50	561
FAQ76	2025
FAQ76v2	1965
All	6605

Table 3.1: Information about UPV FAQ tests

■ Baseline

Due to the inherent morphological richness of the Czech language, this study adopts **FastText** as the baseline word embedding method. This decision is motivated by FastText’s ability to effectively capture morphological variations within words, a characteristic that has been shown to be advantageous for languages like Czech. While other techniques like Word2Vec and GloVe have been explored for word embedding generation, they have demonstrated lower performance in this context. The FastText word embedding model will be trained using the fastText.cc library ¹

Training parameters

- Architecture: CBOW
- Vector dimensionality: 300
- Loss function: Negative sampling loss
- Dictionary threshold (the frequency of the word to be included in the dictionary): 130

Training data

FastText model is trained on a segment of a preprocessed Common Crawl repository, which encompasses raw data from web pages, as well as metadata and text extractions. Given the nature of this dataset, it is expected to include misspellings and text lacking diacritics.

¹fastText.cc library website: <https://fasttext.cc/>

- (i) Eliminating duplicate entries
- (ii) Filtering out lines containing fewer than 9 characters
- (iii) Excluding URLs
- (iv) Breaking lines into individual sentences
- (v) Converting all text to lowercase
- (vi) Removing lines containing words exceeding 30 characters
- (vii) Excluding HTML tags with less than 100 characters
- (viii) Removing lines surpassing 500 characters
- (ix) Omitting words longer than 21 characters

Processed corpus contains 3.87 billion words.

Model summary

The FastText models trained in this study exhibit differences in vocabulary size. The diacritic model possesses a dictionary of approximately 800,000 words, while the diacriticless model contains roughly 762,000 words. This discrepancy reflects the inherent reduction in word count due to the removal of diacritics in the diacriticless dataset.

These vocabulary sizes directly influence the number of trainable parameters within each model. The number of parameters (N) can be calculated using the formula (3.1)

$$N = V_{\text{size}}d + d \tag{3.1}$$

where V_{size} is model vocabulary size
 d is chosen dimensionality of word embeddings

Applying this formula to the vocabulary sizes of our models:

- The diacritics model possesses approximately 240 million parameters.
- The diacriticless model has approximately 229 million parameters.

■ Chosen transformer models

To ensure a comprehensive evaluation, a curated selection of embedding models will be utilized. This selection encompasses three distinct categories:

- (i) **Existing Czech Embedding Models:** This category incorporates established Czech embedding models developed within the Czech NLP community. Their inclusion allows for a focused analysis of how these models perform specifically for the Czech language.
- (ii) **Multilingual Models from the Massive Text Embedding Benchmark (MTEB):** The evaluation will leverage highly regarded multilingual models readily available through the MTEB. This inclusion enables an assessment of how these models generalize to the Czech language, providing insights into their adaptability across languages.
- (iii) **Popular Monolingual Models from the MTEB (rank is lower than 100):** In addition to multilingual models, this selection will also include well-regarded monolingual models (models trained on a single language) from the MTEB. This allows for a comparative analysis of how these models, potentially trained on English or other high-resource languages, perform on the Czech dataset.

All models should be open-source. And sentence-bert. 1B. With information about training process. If exist two version of models - monolingual (no czech lang) and multilingual, we evaluate only multilingual

Czert-B

The Czert model [13] is a set of Czech BERT-like language representation models developed specifically to enhance performance in processing the Czech language. These models leverage the BERT and A Lite BERT (ALBERT) [15] architectures and are designed to outperform multilingual models by training exclusively on Czech data. The training set includes a comprehensive corpus of Czech texts, such as Wikipedia articles, news, and other texts, accumulating to around 36GB of data.

There are 2 variants of the Czert model, Czert-A and Czert-B. Unfortunately Czert-A model is not available, so we will test only Czert-B model. Czert-B model is based on the traditional BERT architecture (110M parameters). Model are pre-trained from scratch using MLM and NSP tasks. However, a slight modification is made to the NSP task to adapt it better to the Czech language corpus structure.

Seznam's models

This study leverages a group of compact word embedding models specifically designed for the Czech language. These models were developed by the Seznam research group with a focus on efficient word representation generation [5].

- RetroMAE-Small: BERT-small model trained with the Retrieval-oriented Language Models Via Masked Auto-Encoder (RetroMAE) objective [10] on a custom Czech corpus.
- Dist-MPNet-ParaCrawl: BERT-small model distilled from the *sentence-transformers/all-mpnet-base-v2*² model, using parallel cs-en dataset [18] for training.
- Dist-MPNet-CzEng: BERT-small model distilled from the *sentence-transformers/all-mpnet-base-v2*² model, using parallel cs-en dataset Czeng [14] for training.
- SimCSE-RetroMAE-Small: The RetroMAE-Small model fine-tuned with the Simple Contrastive Learning of Sentence Embeddings (SimCSE) [11].
- SimCSE-Dist-MPNet-ParaCrawl: The Dist-MPNet-ParaCrawl model fine-tuned with the SimCSE.
- SimCSE-Dist-MPNet-CzEng: The Dist-MPNet-CzEng fine-tuned with the SimCSE.
- SimCSE-Small-E-Czech: Czech ELECTRA model [12] fine-tuned with the SimCSE objective to enhance sentence embeddings.

The developed models are about eight times smaller and five times faster than conventional base-sized models, making them suitable for real-time applications where computational efficiency is critical. The models are trained using techniques like pre-training, knowledge distillation, and unsupervised contrastive fine-tuning to adapt to the limited availability of labeled Czech data.

² Model on the huggingface website: <https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

mBERT

Multilingual Bidirectional Encoder Representations from Transformers (mBERT) [17] leverages the same transformer-based architecture as the original BERT model but boasts an increased number of parameters (178M) to enhance its capabilities. A key distinction of mBERT lies in its ability to understand and process text data across multiple languages. To achieve this multilingual proficiency, mBERT is trained on a massive dataset sourced from Wikipedia entries in 104 different languages. This corpus is meticulously constructed to ensure balanced representation, meaning each language is included regardless of the size or depth of its corresponding Wikipedia. This approach ensures that even languages with limited resources are adequately represented within the training data, fostering better performance for these languages.

Multilingual E5 (mE5)

mE5 models by Microsoft [3] are advanced text embedding models designed to operate across multiple languages based on English-only E5 models [4]. These models are available in three variants — small, base, and large — catering to different computational efficiency and performance needs. The mE5 models are trained using a two-phase approach. The first phase involves weakly-supervised contrastive pre-training on about 1 billion text pairs sourced from diverse multilingual corpora (Wikipedia, mC4, Multilingual CC News, Reddit, etc.). The second phase is supervised fine-tuning on approximately 1.6 million data points from high-quality labeled datasets (MS MARCO [22], Natural Questions (NQ) [19], TriviaQA [23], SQuAD [25], etc.).

LaBSE

Language-Agnostic BERT Sentence Embedding (LaBSE) model [9], developed by Google, is a state-of-the-art model for generating sentence embeddings that are effective across 109 languages. It leverages the transformer architecture and is trained on both monolingual (from sources like CommonCrawl and Wikipedia) and bilingual data (mined from web pages). LaBSE utilizes a dual-encoder structure with BERT-based encoding modules. This setup enables the efficient processing of text pairs in multiple languages.

XLM-R

The XLM-Roberta (XLM-R) model [16] is a significant advancement in unsupervised cross-lingual representation learning, introduced by Facebook AI. It is specifically designed to improve performance across a wide range of cross-lingual tasks. XLM-R is pre-trained on a dataset dubbed 'CC-100', derived from Common Crawl, covering about 2.5 terabytes of text across 100 languages. This dataset is significantly larger than the ones used by its predecessors, offering a broader and more diverse linguistic foundation.

SentenceTransformers multilingual models [20]

- `distiluse-base-multilingual-cased-v2`: Multilingual knowledge distilled version of multilingual Universal Sentence Encoder [21]. This version supports 50+ languages including Czech.

- paraphrase-multilingual-MiniLM-L12-v2: Multilingual version of paraphrase-MiniLM-L12-v2, trained on parallel data for 50+ languages including Czech.
- paraphrase-multilingual-mpnet-base-v2: Multilingual version of paraphrase-mpnet-base-v2, trained on parallel data for 50+ languages including Czech.

UAE-Large-V1

The UAE-Large-V1 model [6] focuses on enhancing short and long Semantic Textual Similarity (STS) tasks through a novel angle-optimized text embedding approach that works by dividing text embeddings into real and imaginary components in a complex space. This model is designed to address the challenges posed by the saturation zones of the cosine function, which can impede learning by causing vanishing gradients.

The model is trained using a hybrid objective that combines cosine similarity, in-batch negatives, and angle differences in complex space. This approach helps overcome the limitations of traditional cosine similarity measures by ensuring better gradient flow during training. Training dataset includes around 21K samples and is specifically designed to evaluate STS performance on long texts, which are common in real-world applications.

Nomic Embed V1 and Nomic Embed V1.5

The Nomic model [2], designed for reproducible long-context text embedding, employs a modified BERT architecture optimized for 8192 token sequences. It includes innovations such as rotary positional embeddings and SwiGLU activations, enhancing its ability to process longer texts. The model undergoes a comprehensive training regime starting with unsupervised contrastive pretraining on large-scale datasets, followed by supervised fine-tuning using human-annotated data. This process leverages contrastive learning to improve text embedding capabilities, significantly enhancing performance on both short and long-context tasks.

Nomic Embed V1.5 is an improvement upon Nomic Embed V1 that utilizes Matryoshka Representation Learning [1] which gives developers the flexibility to trade off the embedding size for a negligible reduction in performance.

snowflake-arctic-embed

This study investigates the use of snowflake-arctic-embed, a suite of text embedding models optimized for retrieval tasks. The suite offers multiple model sizes (xs, s, m, m-long, l) for flexibility.

Snowflake-arctic-embed leverages a two-stage training pipeline

- (i) Pretraining on large batches of query-document pairs with dynamically generated negative examples using dataset with 400M samples combining public data and proprietary web search data.
- (ii) Fine-tuning the model on a smaller dataset (1M samples) consisting of query-positive document-negative document triplets to improve retrieval accuracy.

GTE and GTE-v1.5

The General Text Embedding (GTE) model is a pivotal development in NLP, utilizing a deep Transformer encoder based on a BERT-like architecture for generating dense text embeddings. Initially, the GTE model is unsupervisedly pre-trained on approximately 800 million text pairs from diverse web sources, enabling broad semantic coverage. It then undergoes supervised fine-tuning with 3 million annotated text triples from varied datasets, including MS MARCO and Natural Questions, applying contrastive learning to enhance text relevance detection and similarity assessments. This dual-stage training strategy equips the GTE model to excel in complex NLP tasks, demonstrating significant versatility and robust performance across multiple applications.

There is also updated version GTE-v1.5.

BGE-v1.5

BAAI General Embeddings (BGE) Models [8] are built on a BERT-like architecture, available in three sizes: small, base, and large. They are trained using a sophisticated multi-stage process involving pre-training on large unlabeled data, contrastive learning for fine-tuning on text pairs, and multi-task learning with high-quality labeled datasets. This training regimen equips BGE models to handle a wide range of text embedding tasks with high efficiency and accuracy.

TaylorAI tiny models

This study incorporates two distilled transformer models from TaylorAI for evaluation:

- **TaylorAI/bge-micro-v2** ³: This model is a 2-step distilled version of the small BGE-v1.5 model
- **TaylorAI/gte-tiny** ⁴: This model is a distilled version of the small version of GTE model

■ 3.2 Optimizing Text Representations for RAG in Technical QA

This section explores the application of RAG algorithms for generating natural language answers to questions posed on technical documents. We investigate the crucial role of text representations in RAG's effectiveness and outline the evaluation approach to identify optimal representation models for this specific task.

■ RAG Algorithm and Text Representation Reliance

The RAG algorithm leverages a two-stage approach for answer generation: retrieval and generation. Both stages rely heavily on the chosen text representation technique.

- **Retrieval:** The algorithm retrieves relevant passages from a technical document (e.g., manual) that address the user's question. To achieve this, RAG represents both the question and potential answer passages using the same text representation model (e.g., word

³TaylorAI/bge-micro-v2 model on the huggingface website: <https://huggingface.co/TaylorAI/bge-micro-v2>

⁴TaylorAI/gte-tiny model on the huggingface website: <https://huggingface.co/TaylorAI/gte-tiny>

embeddings, sentence embeddings). The similarity between the question and passage representations determines the relevance of a passage for answering the question.

- **Generation:** Once retrieved, the passages serve as input for the generation model, which constructs a natural language answer. The quality of the text representation in these retrieved passages is crucial. Rich representations that capture semantic relationships and context allow the generation model to create informative and well-structured answers.

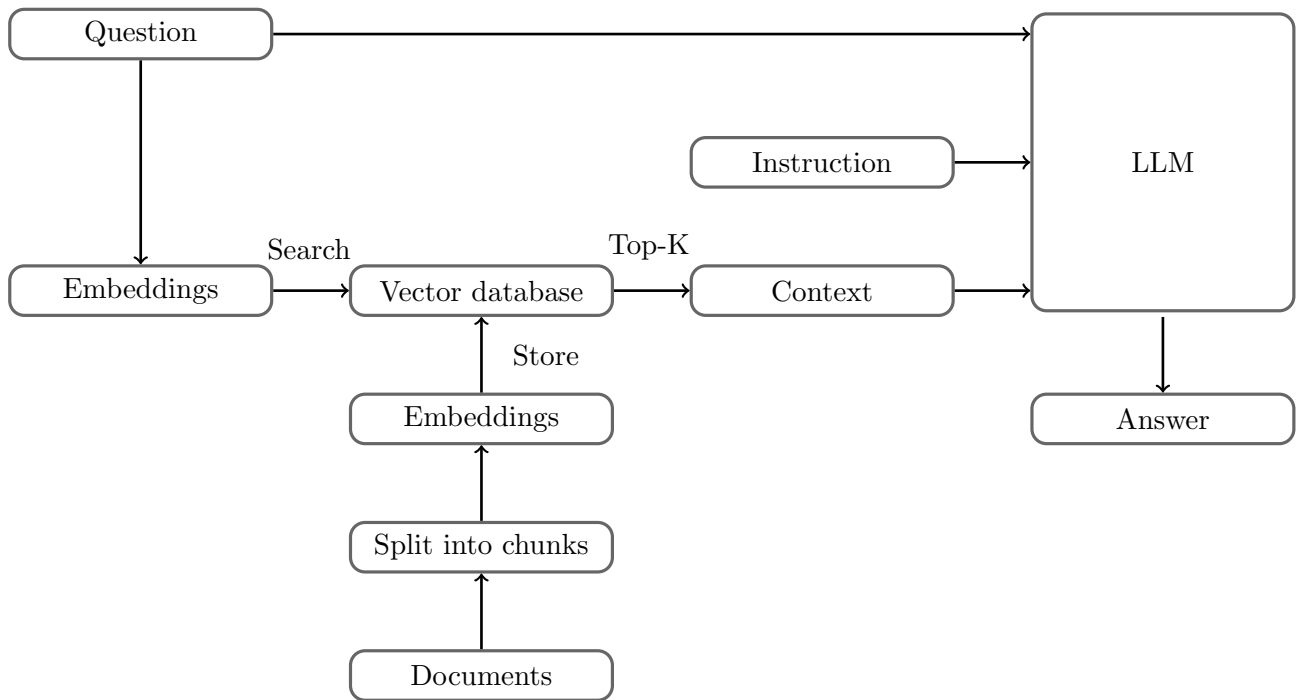


Figure 3.2: RAG scheme

■ Key factors

To identify the most suitable text representation model for RAG in technical QA, we propose a comprehensive evaluation approach that considers the following key factors:

- **Embedding Efficiency:** Computational efficiency is crucial for real-world applications. We will evaluate the processing time required for different representation models, considering factors like embedding dimensionality and model complexity. This ensures the chosen model can handle real-time QA tasks within reasonable processing time constraints.
- **Text Chunk Size:** The size of text chunks used by the RAG algorithm (e.g., words, sentences, or paragraphs) can impact performance. We will investigate the optimal chunk size for technical QA tasks. Here, we will balance the granularity of information retrieved by the model with computational efficiency. Smaller chunks (words) might capture finer details but require more processing, while larger chunks (paragraphs) might be faster to process but might miss relevant details.
- **Factuality of Generated Answers:** The primary objective is to generate answers that are factually accurate and consistent with the technical document. We will evaluate

the models based on metrics that assess the factual correctness and coherence of the generated answers in response to technical questions. This ensures the generated answers are reliable and trustworthy for the user.

■ 3.3 REMOVE: Methodology structure

- Describe the evaluation process for different text representations.
 - Specify the chosen word, sentence, and paragraph representation models (e.g., Fast-Text, BERT variants).
 - Explain the usage of analogy tests and confusion matrices for evaluation.
 - Detail the selection process for the UPV corpus set and its suitability for technical QA tasks.
- Outline the second part of the study focusing on RAG for technical QA.
 - Explain the RAG algorithm and its reliance on text representations.
 - Describe the evaluation approach for selecting optimal representations for RAG.
 - Mention the factors considered during evaluation, such as embedding efficiency, text chunk size, and factuality of generated answers.

■ 4 Experiments and Results

Model	\mathbf{QMA}_d	\mathbf{AMA}_d	\mathbf{QMA}_{dl}	\mathbf{AMA}_{dl}	#params
BASELINE					
fasttext	0.8304	0.2899	0.832	0.302	240M
CZECH MODELS					
Czert-B	0.8759	0.2469	0.8388	0.0977	110M
RetroMAE-Small	0.8651	0.2893	0.8634	0.2437	24M
Dist-MPNet-ParaCrawl	0.8540	0.1089	0.8344	0.0808	24M
Dist-MPNet-CzEng	0.8705	0.0487	0.8322	0.0426	24M
SimCSE-RetroMAE-Small	0.8682	0.3647	0.8649	0.3316	24M
SimCSE-Dist-MPNet-ParaCrawl	0.8817	0.2552	0.8602	0.2322	24M
SimCSE-Dist-MPNet-CzEng	0.8833	0.2278	0.8556	0.1450	24M
SimCSE-Small-E-Czech	0.8094	0.1074	0.8160	0.0878	13M
MULTILINGUAL MODELS					
mBERT	0.8584	0.2012	0.8361	0.1306	178M
mE5 _{Small}	0.8952	0.6078	0.8564	0.4446	118M
mE5 _{Base}	0.8961	0.6019	0.8726	0.5134	278M
mE5 _{Large}	0.9084	0.6559	0.8944	0.5593	560M
LaBSE	0.8875	0.3525	0.8594	0.3264	471M
XLM-R _{Base}	0.8011	0.0098	0.7701	0.0198	279M
XLM-R _{Large}	0.7884	0.0460	0.7411	0.0298	560M
distiluse-base-multilingual-cased-v2	0.8335	0.2978	0.7784	0.2369	135M
paraphrase-multilingual-MiniLM-L12-v2	0.8502	0.4062	0.8029	0.2576	118M
paraphrase-multilingual-mpnet-base-v2	0.8752	0.4538	0.8354	0.3174	278M
MONOLINGUAL MODELS					
UAE _{Large} v1	0.8241	0.2913	0.8237	0.2931	335M
GTE _{Small}	0.8549	0.3632	0.8543	0.3634	33M
GTE _{Base}	0.8443	0.3645	0.8437	0.3643	109M
GTE _{Large}	0.8376	0.3345	0.8370	0.3352	335M
GTE _{Base} v1.5	0.8501	0.3336	0.8499	0.3305	137M
GTE _{Large} v1.5	0.8592	0.3294	0.8586	0.3289	434M
BGE _{Small} v1.5	0.8479	0.3816	0.8474	0.3798	33M
BGE _{Base} v1.5	0.8368	0.3246	0.8362	0.3240	109M
BGE _{Large} v1.5	0.8244	0.2938	0.8238	0.2955	335M
TaylorAI/BGE _{Micro} v2	0.8476	0.3616	0.8475	0.3616	17M
TaylorAI/GTE _{Tiny}	0.8492	0.3343	0.8488	0.3342	23M

Table 4.1: **Evaluation of models.** We show evaluation results where: \mathbf{QMA}_d (\mathbf{QMA}_{dl}) are Question Match Accuracy for diacritics (diacriticless) model. \mathbf{AMA}_d (\mathbf{AMA}_{dl}) are Question Match Accuracy for diacritics (diacriticless) model. #params is total number of parameters.

■ Balanced models

To ensure the effectiveness of the evaluation process, a selection criterion was applied to the initial set of candidate models. This criterion focused on Question Matching Accuracy and Answer Matching Accuracy for both diacritic and diacriticless models. Models that exhibited performance below the established baseline for their respective category (diacritic or diacriticless) were excluded from further evaluation.

Additionally, models with lower performance metrics were removed if a smaller, more efficient model demonstrated comparable or superior accuracy. This approach ensures that the final selection of models for evaluation represents a balance between effectiveness and efficiency.

Model	$\mathbf{QMA_d}$	$\mathbf{AMA_d}$	$\mathbf{QMA_{dl}}$	$\mathbf{AMA_{dl}}$	$\mathbf{\#params}$
simcse-retromae-small-cs	0.8682	0.3647	0.8649	0.3316	24M
GTE _{Small}	0.8549	0.3632	0.8543	0.3634	33M
mE5 _{Small}	0.8952	0.6078	0.8564	0.4446	118M
mE5 _{Base}	0.8961	0.6019	0.8726	0.5134	278M
mE5 _{Large}	0.9084	0.6559	0.8944	0.5593	560M

Table 4.2: **Balanced models.** We show most factual models according to their efficiency, where: $\mathbf{QMA_d}$ ($\mathbf{QMA_{dl}}$) are Question Match Accuracy for diacritics (diacriticless) model. $\mathbf{AMA_d}$ ($\mathbf{AMA_{dl}}$) are Question Match Accuracy for diacritics (diacriticless) model. $\mathbf{\#params}$ is total number of parameters.

- Present the results of the evaluation for different text representations using analogy tests and confusion matrices.
- Discuss the findings regarding the effectiveness of each representation model for capturing semantic relationships in technical text.
- Analyze the results from the RAG evaluation, highlighting the impact of different representations and text chunk sizes on answer generation quality and CPU efficiency.
- Identify the representation model that achieves a balance between factuality of answers and computational demands.

■ 5 Discussion

- Interpret the overall findings and their implications for choosing suitable text representations for RAG in technical QA tasks.
- Discuss the strengths and limitations of the chosen evaluation methods.
- Address potential challenges encountered during the study and suggest improvements for future research.

■ 6 Conclusion

Summarize the achieved results. Can be similar as an abstract or an introduction, however, it should be written in past tense.

- Summarize the key takeaways from the research, emphasizing the most effective text representation model for RAG in technical QA based on the evaluation criteria.
- Briefly mention the trade-offs between factuality, CPU usage, and other factors in selecting representations for RAG.
- Suggest potential future research directions, such as exploring other text representation methods or evaluating RAG performance on different datasets.

7 References

- [1] A. Kusupati, G. Bhatt, A. Rege, *et al.*, *Matryoshka representation learning*, 2024. arXiv: [2205.13147 \[cs.LG\]](#).
- [2] Z. Nussbaum, J. X. Morris, B. Duderstadt, and A. Mulyar, *Nomic embed: Training a reproducible long context text embedder*, 2024. arXiv: [2402.01613 \[cs.CL\]](#).
- [3] L. Wang, N. Yang, X. Huang, L. Yang, R. Majumder, and F. Wei, *Multilingual e5 text embeddings: A technical report*, 2024. arXiv: [2402.05672 \[cs.CL\]](#).
- [4] L. Wang, N. Yang, X. Huang, *et al.*, *Text embeddings by weakly-supervised contrastive pre-training*, 2024. arXiv: [2212.03533 \[cs.CL\]](#).
- [5] J. Bednář, J. Náplava, P. Barančíková, and O. Lisický, *Some like it small: Czech semantic embedding models for industry applications*, 2023. arXiv: [2311.13921 \[cs.CL\]](#).
- [6] X. Li and J. Li, “Angle-optimized text embeddings,” *arXiv preprint arXiv:2309.12871*, 2023.
- [7] A. Vaswani, N. Shazeer, N. Parmar, *et al.*, *Attention is all you need*, 2023. arXiv: [1706.03762 \[cs.CL\]](#).
- [8] S. Xiao, Z. Liu, P. Zhang, and N. Muennighoff, *C-pack: Packaged resources to advance general chinese embedding*, 2023. arXiv: [2309.07597 \[cs.CL\]](#).
- [9] F. Feng, Y. Yang, D. Cer, N. Arivazhagan, and W. Wang, *Language-agnostic bert sentence embedding*, 2022. arXiv: [2007.01852 \[cs.CL\]](#).
- [10] S. Xiao, Z. Liu, Y. Shao, and Z. Cao, *Retromae: Pre-training retrieval-oriented language models via masked auto-encoder*, 2022. arXiv: [2205.12035 \[cs.CL\]](#).
- [11] T. Gao, X. Yao, and D. Chen, “Simcse: Simple contrastive learning of sentence embeddings,” in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021, pp. 6894–6910.
- [12] M. Kocián, J. Náplava, D. Štancl, and V. Kadlec, *Siamese bert-based model for web search relevance ranking evaluated on a new czech dataset*, 2021. arXiv: [2112.01810 \[cs.IR\]](#).
- [13] J. Sido, O. Prazák, P. Pribán, J. Pasek, M. Seják, and M. Konopík, “Czert - czech bert-like model for language representation,” *CoRR*, vol. abs/2103.13031, 2021. arXiv: [2103.13031](#). [Online]. Available: <https://arxiv.org/abs/2103.13031>.
- [14] T. Kocmi, M. Popel, and O. Bojar, *Announcing czeng 2.0 parallel corpus with over 2 gigawords*, 2020. arXiv: [2007.03006 \[cs.CL\]](#).
- [15] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, *Albert: A lite bert for self-supervised learning of language representations*, 2020. arXiv: [1909.11942 \[cs.CL\]](#).
- [16] A. Conneau, K. Khandelwal, N. Goyal, *et al.*, “Unsupervised cross-lingual representation learning at scale,” *CoRR*, vol. abs/1911.02116, 2019. arXiv: [1911.02116](#). [Online]. Available: <http://arxiv.org/abs/1911.02116>.
- [17] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, *Bert: Pre-training of deep bidirectional transformers for language understanding*, 2019. arXiv: [1810.04805 \[cs.CL\]](#).
- [18] M. Esplà, M. Forcada, G. Ramírez-Sánchez, and H. Hoang, “Paracrawl: Web-scale parallel corpora for the languages of the eu,” in *Proceedings of Machine Translation Summit XVII: Translator, Project and User Tracks*, European Association for Machine Translation, 2019, pp. 118–119.
- [19] T. Kwiakowski, J. Palomaki, O. Redfield, *et al.*, “Natural questions: A benchmark for question answering research,” *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3877–3886, 2019.

- [20] N. Reimers and I. Gurevych, “Sentence-bert: Sentence embeddings using siamese bert-networks,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Nov. 2019. [Online]. Available: <http://arxiv.org/abs/1908.10084>.
- [21] Y. Yang, D. Cer, A. Ahmad, *et al.*, *Multilingual universal sentence encoder for semantic retrieval*, 2019. arXiv: [1907.04307](https://arxiv.org/abs/1907.04307) [cs.CL].
- [22] P. Bajaj, D. Campos, N. Craswell, *et al.*, *Ms marco: A human generated machine reading comprehension dataset*, 2018. arXiv: [1611.09268](https://arxiv.org/abs/1611.09268) [cs.CL].
- [23] M. Joshi, E. Choi, D. Weld, and L. Zettlemoyer, “Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension,” *arXiv preprint arXiv:1705.03551*, 2017.
- [24] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, “Enriching word vectors with subword information,” *arXiv preprint arXiv:1607.04606*, 2016.
- [25] P. Rajpurkar, R. Jia, and I. Polosukhin, “Squad: A question answering dataset for readers,” *arXiv preprint arXiv:1606.05250*, 2016.
- [26] J. Pennington, R. Socher, and C. D. Manning, “Glove: Global vectors for word representation,” in *Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1532–1543. [Online]. Available: <http://www.aclweb.org/anthology/D14-1162>.
- [27] T. Mikolov, K. Chen, G. Corrado, and J. Dean, *Efficient estimation of word representations in vector space*, 2013. arXiv: [1301.3781](https://arxiv.org/abs/1301.3781) [cs.CL].

■ A Appendix A

Model	L	H_m	H_{ff}	A	D	V	#params
Czert-B	12	768	3072	12	768	31K	110M
retromae-small-cs	12	256	1024	4	256	58K	24M
dist-mpnet-paracrawl-cs-en	12	256	1024	4	256	58K	24M
dist-mpnet-czeng-cs-en	12	256	1024	4	256	58K	24M
simcse-retromae-small-cs	12	256	1024	4	256	58K	24M
simcse-dist-mpnet-paracrawl-cs-en	12	256	1024	4	256	58K	24M
simcse-dist-mpnet-czeng-cs-en	12	256	1024	4	256	58K	24M
simcse-small-e-czech	12	256	1024	4	256	31K	13M
mBERT	12	768	3072	12	768	120K	178M
mE5 _{Small}	12	384	1536	12	384	250K	118M
mE5 _{Base}	12	768	3072	12	768	250K	278M
mE5 _{Large}	24	1024	4096	16	1024	250K	560M
LaBSE	12	768	3072	12	768	502K	471M
XLM-R _{Base}	12	768	3072	12	768	250K	279M
XLM-R _{Large}	24	1024	4096	16	1024	250K	560M
distiluse-base-multilingual-cased-v2	6	768	3072	12	512	120K	135M
paraphrase-multilingual-MiniLM-L12-v2	12	384	1536	12	512	250K	118M
paraphrase-multilingual-mpnet-base-v2	12	768	3072	12	514	250K	278M
UAE _{Large} v1	24	1024	4096	16	1024	31K	335M
GTE _{Small}	12	384	1536	12	384	31K	33M
GTE _{Base}	12	768	3072	12	768	31K	109M
GTE _{Large}	24	1024	4096	16	1024	31K	335M
GTE _{Base} v1.5	12	768	3072	12	768	31K	137M
GTE _{Large} v1.5	24	1024	4096	16	1024	31K	434M
BGE _{Small} v1.5	12	384	1536	12	384	31K	33M
BGE _{Base} v1.5	12	768	3072	12	768	31K	109M
BGE _{Large} v1.5	24	1024	4096	16	1024	31K	335M
TaylorAI/BGE _{Micro} v2	3	384	1536	12	512	31K	17M
TaylorAI/GTE _{Tiny}	6	384	1536	12	384	31K	23M

Table A.1: **Details on model sizes.** We show the number of layers L , the number of hidden states of the model H_m , the dimension of the feed-forward layer H_{ff} , the number of attention heads A , the dimension of output embedding D , the size of the vocabulary V and the total number of parameters #params. For Transformer encoders, the number of parameters can be approximated by $4LH_m^2 + 2LH_mH_{ff} + VH_m$. While this table gives more hindsight on the difference of capacity of each model, note it does not highlight other critical differences between the models.