# Fast Extraction of Word Embedding from Q-contexts

Junsheng Kong*†
School of Software Engineering,
South China University of Technology
China
sescut_kongjunsheng@mail.scut.edu.cn

Weizhao Li*†
School of Software Engineering,
South China University of Technology
China
se_weizhao.li@mail.scut.edu.cn

Zeyi Liu
University of Cambridge
United Kingdom
zl411@cam.ac.uk

Ben Liao‡
Tencent Quantum Lab
China
bliao@tencent.com

Jiezhong Qiu
Tsinghua University
China
qiujz16@mails.tsinghua.edu.cn

Chang-Yu Hsieh
Tencent Quantum Lab
China
kimhsieh@tencent.com

Yi Cai†‡
School of Software Engineering,
South China University of Technology
China
ycai@scut.edu.cn

Shengyu Zhang‡
Tencent Quantum Lab
China
shengyzhang@tencent.com

## ABSTRACT

The notion of word embedding plays a fundamental role in natural language processing (NLP). However, pre-training word embedding for very large-scale vocabulary is computationally challenging for most existing methods. In this work, we show that with merely *a small fraction of contexts (Q-contexts)* which are *typical* in the whole corpus (and their mutual information with words), one can construct high-quality word embedding with negligible errors. Mutual information between contexts and words can be encoded canonically as a sampling state, thus, Q-contexts can be fast constructed. Furthermore, we present an efficient and effective WEQ method, which is capable of extracting word embedding *directly* from these typical contexts. In practical scenarios, our algorithm runs 11∼ 13 times faster than well-established methods. By comparing with well-known methods such as matrix factorization, word2vec, GloVe and fasttext, we demonstrate that our method achieves comparable performance on a variety of downstream NLP tasks, and in the meanwhile maintains run-time and resource advantages over all these baselines.

## CCS CONCEPTS

• **Computing methodologies** → **Lexical semantics**.

*Equal contribution. Work was done during internship at Tencent Quantum Lab.
†Also with The Key Laboratory of Big Data and Intelligent Robot (South China University of Technology), Ministry of Education.
‡Corresponding Author.

## KEYWORDS

word embedding, Q-contexts, fast extraction, large-scale

## 1 INTRODUCTION

Word embedding plays a fundamental role in the development and real-world applications of natural language processing (NLP). It efficiently provides meaningful representations of individual words in a continuous space, allowing smooth integration with machine learning models in various downstream NLP tasks [17, 22]. The notion of word embedding is also the predecessor of follow-up deep contextualization models, including the recently discovered powerful pre-trained contextual embedding models such as ELMo [28] and BERT [12].

High-quality embedding of words can help boost the performance of many machine learning models in NLP tasks. Recent work about word embedding can be categorized into two genres, i.e., neural network based methods [4, 24, 27] and global matrix factorization based methods [2, 10, 19]. Word2Vec, GloVe and fasttext are the most popular neural network based methods. Most of the existing methods focus on improving the performance of word embedding. However, it is computationally expensive to obtain such word embedding — it takes several days and, typically, around a hundred CPU cores to attain decent quality representation of words [23–25, 27]. Global matrix factorization based methods for generating word embedding have roots stretching as far back as LSA[10]. These methods utilize low-rank approximations to decompose large matrices that capture statistical information about a

corpus. Previous work has shown that both the word2vec and GloVe methods can be viewed as implicit factorization of special information matrices[2, 19]. Although global matrix factorization based method is more efficient than the neural network based methods, it still needs to factorize a large $n \times n$ information matrix for large-scale vocabulary, where $n$ is the size of vocabulary. This makes it highly expensive to directly factorize and calculate for large-scale word embedding learning.

To address the efficiency limitations of current work, we propose to study word embedding learning for large-scale vocabulary with the goal of efficiency and theoretical guarantees. Recent literature has shown that quantum perspective can thus provide advantages for classical machine learning [29, 38]. Coecke et al. have previously demonstrated a potential quantum advantage for NLP in various ways including by algorithmic speed-up for search-related or classification tasks [7]. By mimicking how word-meanings are encoded in quantum states, we design our algorithms implemented on classical computers to speed up the word embedding learning problem. The main idea is to construct a small and typical information matrix which is a good approximation of the original information matrix. Both the construction and the factorization of the small matrix require a low cost. With this design, we are able to demonstrate running-time supremacy for solving a large-scale word embedding problem and maintain accuracy for various downstream NLP tasks.

We reveal a simple relation between a word vector $e_w$ for a target word $w$ and what we call *Q-contexts*. Q-contexts are a small fraction of contexts that are typical in the whole corpus, capable of capturing the most important information encoded in the information matrix $M$ – they are certain rows of $M$ chosen to represent the original information matrix. The word vector of $w$ is shown to be a combination of its interaction with these contextual environments

$$e_w \approx \sum_{c \in \text{Q-contexts}} \lambda_c M_{c,w} \tag{1}$$

where $M_{c,w}$ is the entry in the information matrix $M$ indexed by $c$ and $w$, $\lambda_c$ is a constant vector for a context $c$ to be determined (for a more detailed description see Section 3.1). Information matrix $M$ can be naturally encoded as a sampling state, enabling a fast construction of Q-contexts. To the best of our knowledge, it has not been studied to extract meaningful word embedding from its mutual information with a few contexts.

Based on these, we develop a WEQ method that substantially accelerates the word embedding learning — in fact, our method is at least $11 \sim 13$ times faster than well-established methods, and has fewer resource requirements on CC corpus. We show empirically that WEQ achieves comparable performance in comparison to well-known methods, such as the direct matrix factorization, word2vec, GloVe and fasttext [4, 19, 24, 27] and maintains high accuracy in various downstream NLP tasks. Further, WEQ's efficiency and effectiveness are theoretically backed up. The small Q-contexts matrix is a good approximation of the original information matrix with negligible error, maintaining the representation power of its learned embedding.

The organization of the rest of this article is as follows. In Section 2, we recall a general matrix factorization perspective on well-known word embedding methods. In Section 3, we introduce the WEQ method. In Section 4, we analyze the approximation error of

the Q-contexts with theoretical proof. In Section 5, we conduct a comprehensive set of experiments demonstrating the accuracy and efficiency of our method. In Section 6, we review the related work of word embedding. Finally, we give a conclusion in Section 7.

### Table 1: Notation.

| Notation | Description |
|---|---|
| $P$ | the multiset of context-word pairs |
| $w$ | the target word |
| $c$ | context word around the target word |
| $\#(c, w)$ | the number of co-occurrences of $c$ and $w$ in $P$ |
| $\#(c)$ | the number of times of $c$ appears in $P$ |
| $\#(w)$ | the number of times of $w$ appears in $P$ |
| $\|P\|$ | $\sum_c \sum_w \#(c, w)$ |
| $M_{c,w}$ | entry in information matrix indexed by $c$ and $w$ |
| $e_w$ | word (row) vector for the target word $w$ |
| $E_w$ | $E_w = \begin{pmatrix} e_{w_1} \\ \vdots \\ e_{w_n} \end{pmatrix}$ |
| $e_c$ | context (row) vector for context $c$ |
| $\|\cdot\|$ | the $\ell^2$-norm of a vector |
| $A$ | matrix |
| $\|A\|_F$ | the Hilbert-Schmidt norm of $A$ |
| $\|A\|_{op}$ | the operator norm of matrix $A$ |
| $A_{i,*}$ | row $i$ of $A$ |
| $A_{*,j}$ | column $j$ of $A$ |
| $A^\top$ | the transpose of $A$ |
| $\text{nnz}(A)$ | the number of nonzeros in $A$ |
| $R$ | Q-contexts matrix |
| $\tilde{R}$ | the normalized version of $R$ |

## 2 PRELIMINARIES

Commonly, the problem of word embedding is learned by capturing the semantic relationship between word-context pairs $(w, c)$. For a target word $w$, its context word c is obtained from the neighborhood centering around the locations where $w$ appears in the corpora. Previously established results show that the factorization of information matrices provides a united framework for many important existing word embedding algorithms, including word2vec, GloVe, PMI, and NCE [2, 19, 25, 27]. The main difference between these methods lies in different choices of mutual information matrix $M_{c,w}$ between contexts and words.

As indicated in [2, 6, 8, 37], factorizing the following point-wise mutual information matrix (PMI) yields effective word representations $M_{c,w} = \log \frac{\#(c,w)|P|}{\#(c)\#(w)}$, where $P$ is the multiset of context-word pairs, #(c,w) is the number of co-occurrences of context word $c$ and target word $w$ in the $P$, $\#(c)$ and $\#(w)$ are the number of times $c$ and $w$ appear in $P$ respectively. It is shown in [19] that word vectors from word2vec can be obtained from factorization of a shifted version of PMI: $M_{c,w} = \log \frac{\#(c,w)|P|}{\#(c)\#(w)\cdot\kappa}$, where $\kappa$ denotes the number of negative samples. They also show that NCE model [25] is in fact factorizing $M_{c,w} = \log \frac{\#(c,w)}{\#(c)\cdot\kappa}$. To improve performance, a positive

version of PMI (PPMI) $M_{c,w} = \log_+ \frac{\#(c,w)|P|}{\#(c)\#(w)}$ and a shifted version of PPMI (SPPMI with shift parameter $\kappa$) $M_{c,w} = \log_+ \frac{\#(c,w)|P|}{\#(c)\#(w)\cdot\kappa}$ are proposed, where $\log_+(x) = \max(\log x, 0)$.

It is shown [2] that GloVe objective is in fact optimizing (modulo some error term)

$$\sum \#(c,w)\Big(\log\#(c,w) - e_c \cdot e_w - \|e_c\|^2 - \|e_w\|^2\Big)^2,$$

where $e_c$ is the context vector for context $c$ and $e_w$ is the word vector for the target word $w$. In their theory, the authors also show that for some constant $Z$:

$$\log p(c,w) \approx \|e_c + e_w\|^2/2d - 2\log Z$$
$$\log p(w) \approx \|e_w\|^2/2d - \log Z.$$

Since $p(w) \approx \frac{\#(w)}{|P|}$, we conclude

$$e_c \cdot e_w \approx \log\left[(|P|Z)^{4d} \cdot \frac{\#(c,w)}{(\#(c)\#(w))^{2d}}\right].$$

Factorization of mutual information matrices constitutes a unified framework of these word embedding algorithms: $M_{c,w} = e_c \cdot e_w$. We list necessary notations and their descriptions in Table 1.

## 3 METHOD

In this section, we present WEQ method which is an efficient and effective method for large-scale word embedding learning problem. We develop the WEQ method to construct and factorize a small typical information matrix that approximates the original information matrix. The WEQ method is composed of three steps, as illustrated in Fig. 1. First, it calculates the information matrix $M$ from co-occurrence matrix $X$. Secondly, it constructs the small typical information matrix (Q-contexts) from the original information matrix through $\ell^2$-norm sampling. Third, it conducts the singular value decomposition of Q-contexts matrix to obtain the word embedding.

### 3.1 Q-contexts Definition

We first introduce the $\ell^2$-norm sampling, then describe the definition of Q-contexts.

$\ell^2$-**norm sampling**: The $\ell^2$-norm sampling technique has well exhibited its effectiveness in machine learning [16, 31] and randomized linear algebra [13]. In fact, the work by Frieze, Kannan, and Vempala [14] shows that with certain $\ell^2$-norm sampling assumptions, a form of singular value estimation can be achieved in time independent of the size of input matrix. Further, the work by Tang [34] shows that sampling from the projection of a vector onto a subspace is not outside the realm of feasibility.

Inspired by these, we leverage the $\ell^2$-norm sampling to construct a small typical information matrix to solve the large-scale word embedding learning problem. We now elaborate on our proposed relation between Q-contexts and words in Equation (1).

Given an information matrix $M$, we first encode the information matrix $M$ into an $\ell^2$-norm state which will be described in detail in Section 3.2. Now that the mutual information $M$ is prepared to be a state, each $\ell^2$-norm sampling yields a context $c_i$ with probability $p_{c_i}$. The collection of the corresponding row vectors $r_{c_i} = M_{c_i,*}$ is what we call Q-contexts. In our proposed scheme, $p_c$ is designed to
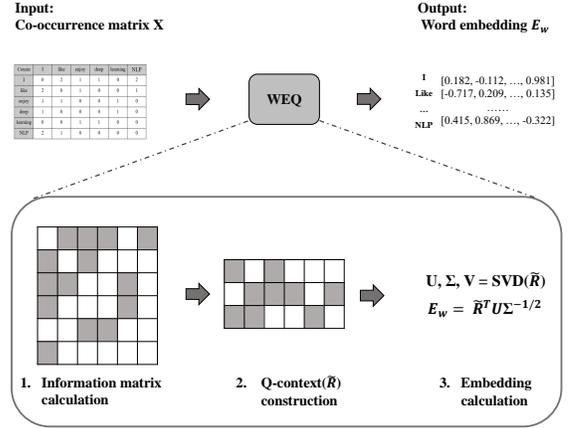
correctly reflect the amount of information carried by the context $c$, so that contexts with more information are more likely to be sampled.

*Definition 3.1.* A Q-contexts matrix $R \in \mathbb{R}^{k\times n}$ is the collection of $k$ rows $r_{c_i}$ in the information matrix $M$:

$$R = \begin{pmatrix} r_{c_1} \\ \vdots \\ r_{c_k} \end{pmatrix}$$

where context $c_i$ is the $i$-th sampling outcome from the information matrix $M$.

The central idea in Equation (1) is that a few $\ell^2$-norm sampling of the information matrix provides sufficient information such that a good word embedding can be obtained from a linear combination of its mutual information with Q-contexts

$$e_w \approx \sum_{c \text{ runs over row indices in } R} r_{c,w}\lambda_c$$

where the *context coefficients vectors* (i.e., $\lambda_c$'s) are to be determined in Section 3.2.

### 3.2 Word Embedding from Q-contexts (WEQ)

We now give the full algorithm (Algorithm 2) for computing word embedding based on Q-contexts matrix $R$. WEQ method consists of three steps: information matrix calculation, Q-contexts construction, and calculating word embedding from Q-contexts.

**Step 1: Information matrix calculation.** As we have discussed in Section 2, there are different kinds of information matrices. Here, we choose PPMI and SPPMI matrices proposed in [19] which shows that exact factorizing (PPMI / SPPMI) matrix with SVD is at least as good as SGNS's solutions. Given a co-occurrence matrix, we construct the PPMI information matrix as follows:

$$M_{c,w}(PPMI) = \log_+ \frac{\#(c,w)|P|}{\#(c)\#(w)}$$



**Figure 1: WEQ method.**

and a shifted version of PPMI (SPPMI with shift parameter $\kappa$)

$$M_{c,w}(SPPMI) \ = \ \log_+ \frac{\#(c,w)|P|}{\#(c)\#(w) \cdot \kappa}.$$

**Step 2: Q-contexts construction.** As in Definition 3.1, we will encode the information matrix $M$ into an $\ell^2$-norm state, which admits our fast construction of Q-contexts matrix $R$.

In our method, we propose a data structure that achieves fast $\ell^2$-norm sampling in practice. For the information matrix $M \in \mathbb{R}^{n \times n}$, we store the cumulative summation of its row-squared ($\|r_1\|^2, \|r_1\|^2 + \|r_2\|^2, \cdots, \|M\|_F^2$). To perform the $\ell^2$-norm sampling, we first generate a random number $a$ from the uniform distribution $\mathcal{U}(0, \|M\|_F^2)$ and perform an efficient binary search algorithm to find the leftmost index such that $a$ is less than or equal to the corresponding cumulative sum. Through repeating the $\ell^2$-norm sampling for $k$ times, we get the small Q-contexts $R$. We normalize each row to get $\tilde{R}$.

In addition, we can use column sampling to reduce the matrix $\tilde{R}$ again. Applying the theorem and similar analysis to $\tilde{R}^\top$, we get a smaller matrix $\tilde{C}$ for which with high probability

$$\tilde{R}\tilde{R}^\top \approx \tilde{C}\tilde{C}^\top.$$

Then

$$U\Sigma^2 U^\top \approx \tilde{R}\tilde{R}^\top \approx \tilde{C}\tilde{C}^\top.$$

We obtain good approximations of the right singular vectors $U$ and singular values $\Sigma$ of $\tilde{R}$, by simply doing the same calculations for the much smaller matrix $\tilde{C}$.

---

**Algorithm 1:** Q-contexts construction

**input** : information matrix $M \in \mathbb{R}^{n \times n}$;
number of samples $k$
**output** : Q-contexts matrix $\tilde{R} \in \mathbb{R}^{k \times n}$

1 /* Prepare the state: compute the cumulative sum of M*/
$S(M) =$
$(\|r_1\|^2, \|r_1\|^2 + \|r_2\|^2, \|r_1\|^2 + \|r_2\|^2 + \|r_3\|^2, \cdots, \|M\|_F^2)$;
2 **for** $i \leftarrow 1$ **to** $k$ **do**
3    /* sample a row $i$ from the state of $M$ */
4    Generate $s \sim \mathcal{U}(0, \|M\|_F^2)$;
5    Search $i$ such that $S(M)_{i-1} \leq s < S(M)_i$;
6    /* normalization */
7    Form $\tilde{R}_{i,*} = \frac{\|M\|_F}{\sqrt{k}\|M_{r_i,*}\|} M_{r_i,*}$;
8 **end**
9 Return $\tilde{R}$;

---

**Step 3: Calculating word embedding from Q-contexts.** It is known [25] that embedding matrix $E_w$ taken to be the form $E_w = V_w \sqrt{\Sigma}$ can be beneficial in predictive performance, where $V_w$ and $\Sigma$ are right singular vectors and singular values of information matrix $M$. In view of Theorem (4.2) later, with high probability $\tilde{R}^\top \tilde{R} \approx M^\top M$, which implies that $\tilde{R} \approx U \Sigma V_w^\top$ where $U$ is the left singular vectors of $\tilde{R}$. Since $U$ is a (partial-)isometry, we obtain the matrix form of Equation (1):

$$E_w = V_w \sqrt{\Sigma} \approx \tilde{R}^\top U \Sigma^{-1/2} = R^\top \Lambda, \tag{2}$$

where $\Lambda = D^{-1} U \Sigma^{-1/2}$ is the matrix of context coefficients.

**Complexity Analysis.** We get Algorithm 2 by putting the above procedures together. As for line 1, it requires $O(\text{nnz}(X))$ time to perform point-by-point operations on the nonzero elements of the co-occurrence matrix $X$. As for line 2 and line 3, time complexity for state preparation is $O(n)$ and for $\ell^2$-norm sampling is $O(k \log n)$. As for line 4, $O(\text{poly}(k))$ time is spent in singular vectors computation. In MF method, it needs to take highly expensive $O(\text{poly}(n))$ time to compute singular vectors of original information matrix $M$. As for line 5, $O(kdn)$ time is spent in matrix multiplication.

---

**Algorithm 2:** WEQ method

**input** : sparse co-occurrence matrix $X \in \mathbb{R}^{n \times n}$;
        number of samples $k$, and embedding dimension $d$
**output** : embedding matrix $E_w$
1 Compute information matrix $M$ (e.g. PPMI, SPPMI) from $X$ ;
2 Construct the Q-contexts matrix $\tilde{R}$ according to Algorithm 1;
3 (Optional) Construct the much smaller Q-contexts matrix $\tilde{C}^\top$ from $\tilde{R}^\top$ according to Algorithm 1 again;
4 Compute the top $d$ left singular vectors $U$ and singular values $\Sigma$ for the Q-contexts matrix $\tilde{R}$ or $\tilde{C}$;
5 Return $E_w = \tilde{R}^\top U \Sigma^{-1/2}$;

---

## 4 THEORETICAL PROOF

In this section, we demonstrate that very few $\ell^2$-norm samplings on the information matrix $M$ suffice to extract high-quality word embeddings.

The idea that word embedding hidden in mutual information matrix $M$ can be extracted from a few $\ell^2$-norm sampling on the state of $M$ might seem surprising at first glance. The root of this phenomenon stems from the fact that Q-context $r_c$ is a randomized object, its probabilistic behavior has an almost deterministic nature in the sense of Central Limit Theorem.

In more precise terms, word representation hides in the symmetric form $M^\top M$. In view of the random nature of context $r_c$, this form can be written as the expectation of the rank-one matrix $S = \frac{1}{p_c} r_c^\top r_c$ relating to context $c$,

$$M^\top M = r_1^\top r_1 + \cdots + r_n^\top r_n = \sum_c p_c \times \frac{1}{p_c} r_c^\top r_c, \tag{3}$$

where $p_c$ is the probability of getting context $c$ from a measurement, and $c$ runs over the set of contexts. Precisely, $S$ is the random matrix taking value $\frac{1}{p_c} r_c^\top r_c$ with probability $p_c$.

The idea of Central Limit Theorem and its general form of probabilistic measure concentration is also valid in a functional analytic (matrix) context, where a scalar-valued random variable is generalized to a matrix-valued one. One of the most celebrated theorems is the operator/matrix Bernstein concentration inequality [36].

THEOREM 4.1. *Let $S_i \in \mathbb{R}^{n \times n}$, $i = 1, \cdots, k$ be a sequence of independent identically distributed symmetric norm-bounded matrices with mean $\mu$. Then for sufficiently small $\epsilon$, we have*

$$\mathbb{P}\left[ \left\| \frac{1}{k} \sum S_i - \mu \right\|_{op} \geq \epsilon \right] \leq 4 \cdot sr(\varsigma) \cdot \exp\left( -\frac{k\epsilon^2}{2\|\varsigma\|_{op}^2} \right),$$

where $\varsigma^2$ is the covariance matrix of $S_1 : \varsigma^2 = \mathbb{E}(S_1 - \mu)^2$, and $sr(\cdot) = \frac{\|\cdot\|_F^2}{\|\cdot\|_{op}^2}$ denotes the stable rank of a matrix.

Thanks to the matrix concentration Theorem 4.1 , we see that, with high probability, the expectation $\mu = \mathbb{E}S = M^\top M$ can be well approximated by the average of $k \geq \frac{3\|\varsigma\|_{op}^2}{\epsilon^2} \log sr(\varsigma)$ samples $S_1, \cdots, S_k$. The average of these samples is

$$M^\top M \approx \frac{1}{k}(S_1 + \cdots + S_k) = \tilde{R}^\top \tilde{R},$$

where $\tilde{R}$ is the normalized version of Q-contexts matrix $R$:

$$\tilde{R} = D^{-1}R, D = \text{diag}(\sqrt{kp_{i_1}}, \cdots, \sqrt{kp_{i_k}}).$$

Precisely, we get the following theorem:

THEOREM 4.2. *Let $p_c = \frac{\|r_c\|^2}{\|M\|_F^2}$, then we have*

$$\mathbb{P}\left[\left\|\tilde{R}^\top \tilde{R} - M^\top M\right\|_{op} \geq \epsilon\right] \leq 4 \cdot sr(\varsigma) \cdot \exp\left(-\frac{k\epsilon^2}{2\|\varsigma\|_{op}^2}\right),$$

*where $\varsigma = \sqrt{\mathbb{E}(S - \mu)^2}$ with stable rank $sr(\varsigma)$ being bounded by the rank of $M$, and*

$$\|\varsigma\|_{op} \leq \min\left(\|M\|_F\|M\|_{op}, \sqrt{\|M\|_F^4 - \|M^\top M\|_F^2}\right). \quad (4)$$

PROOF. The probability inequality in the theorem is a direct consequence of Theorem 1. We only need to prove inequality (4) and the bound on $sr(\varsigma)$ in the theorem.

By direct calculation, we have

$$\varsigma^2 = \mathbb{E}S^2 - \mu^2 = \sum_c \frac{1}{p_c}\|r_c\|^2 r_c^\top r_c - (M^\top M)^2.$$

Therefore,

$$\|\varsigma\|_F^2 = Tr(\varsigma^2) = \sum_c \frac{\|r_c\|^4}{p_c} - \|M^\top M\|_F^2$$

$$\sum_c p_c = 1.$$

The Lagrange multiplier Theorem implies that $\|\varsigma\|_F^2$ achieves minimum only when $p_c^2$ is proportional to $\|r_c\|^4$, namely $p_c = \frac{\|r_c\|^2}{\|M\|_F^2}$. In this case,

$$\|\varsigma\|_{op}^2 \leq \|\varsigma\|_F^2 = \|M\|_F^4 - \|M^\top M\|_F^2$$

which is the first part of inequality (4).

For the second part, when $p_c = \frac{\|r_c\|^2}{\|M\|_F^2}$ we see that as matrices

$$\varsigma^2 \leq \sum_c \frac{1}{p_c}\|r_c\|^2 r_c^\top r_c = \|M\|_F^2 M^\top M. \quad (5)$$

Inequality (5) implies

$$\|\varsigma\|_{op} \leq \|M\|_F\|M\|_{op}.$$

The stable rank of $\varsigma$ is bounded as a consequence of inequality (5):

$$sr(\varsigma) \leq rank(\varsigma) = rank(\varsigma^2)$$
$$\leq rank(M^\top M) = rank(M).$$

$\square$

From the proof of Theorem 4.2 (especially the proof of the first part of inequality (4)), we can see that the choice of $p_c$ is motivated by the variance minimization scheme. Therefore such a sampling strategy can lead to the matrix mean at a faster speed.

## 4.1 Error Analysis

We remark that that approximation in Theorem 4.2 is of high quality, both in theory and practice.

From a mathematical point of view, approximation error in Theorem 4.2 is small provided that $k \geq \frac{3\|\varsigma\|_{op}^2}{\epsilon^2}\log sr(\varsigma)$ rows are sampled from the information matrix $M$. Our analysis is tighter than [14] since we make use of a stronger bound on matrix concentration (Theorem 4.1) concerning stable rank and spectral norm of a matrix, instead of dimension and Frobenius norm. Stable rank is upper bounded by the rank of a matrix, which is more effective in a low-rank matrix regime such as word embedding problem. Spectral norm is bounded above by Frobenius norm as well. Moreover, the Frobenius norm can fail to capture the genuine behavior of random matrices in even very simple examples [36]. Therefore, our analysis shows that with very few samples of rows, $\tilde{R}^\top \tilde{R}$ can be a good approximation of $M^\top M$ with negligible error.

These theoretical observations are also consistent with our experiments. Indeed, in Section 5.5, we will see that with information matrix dimensions varying from 50k to 400k, the number of row samples required to achieve a good approximation is relatively stable — in our test it ranges from 40k to 50k, while less than only 2% of accuracy is lost compared to direct calculation with the original information matrix.

## 5 EXPERIMENTS

In this section, we evaluate the proposed WEQ method on three different evaluation tasks: Word Similarity, Text Classification and Named Entity Recognition (NER), which have been commonly used to evaluate previous word embedding methods [19, 39]. We introduce our training corpora and baselines in Section 5.1. We describe the details of implementation and evaluation tasks in Section 5.2 and Section 5.3. We report experimental results and ablation study in Section 5.4 and Section 5.5, respectively.

**Table 2: Data statistics of enwik9, WebBase and CC. '#tokens' indicates the number of total tokens in the whole corpora. '#words' indicates the number of unique words. '#pairs' indicates the number of word-context pairs.**

|  | enwik9 | WebBase | CC |
|---|---|---|---|
| #tokens | 124,301,826 | 2,976,897,565 | 6,065,531,635 |
| #words | 833,184 | 3,107,950 | 10,558,748 |
| #pairs | 42,234,884 | 930,896,198 | 1,861,679,208 |

## 5.1 Training Corpora and Baselines

**Training corpora.** We conduct our experiments on three large scale English corpora: *enwik9* with about 0.1 billion tokens, *Web-Base* [15] with about 3 billion tokens, and the filtering of *Common*

*Crawl corpus (CC)* [26] with about 6 billion tokens. We first lowercase text and then remove the noisy text like HTML span. We mainly implement data pre-processing on the basis of the script[1] used in word2vec. After pre-processed, the vocabulary sizes of three corpora are 56,466 (on enwik9), 277,704 (on WebBase) and 400,000 (on CC) respectively. We prepare co-occurrence matrices with a context window size of 10. The co-occurrence matrices are square symmetric matrices whose row sizes are equal to vocabulary sizes. Our training corpora (enwik9[2],WebBase[3],CC[4]) are publicly available. Note that the original CC is super-scale corpora, which obtains over 418B tokens. To facilitate experiments, we take 6B tokens. The statistics of the three raw corpora are listed in Table 2.

**Baselines.** In our experiments, we mainly compare our method WEQ with popular benchmarks listed below.

- **MF** [19] is a global matrix factorization method that uncovers the semantic information implied in the matrix, such as PPMI and SPPMI matrices.
- **word2vec**[5] [24] is a two-layer neural network that is trained to reconstruct linguistic contexts of words. We choose the SGNS architecture in our experiments.
- **GloVe**[6] [27] is a typical co-occurrence count based neural method, which captures global information from word co-occurrence matrix in the training corpus.
- **fasttext**[7] [4] is a character-based method which represents each word as an n-gram of characters.

### 5.2 Implementation Details

In our experiments, we use code[8] from GloVe to obtain co-occurrence matrix. For MF and WEQ, matrices are stored in a list of lists (LIL) sparse matrix format to allow efficient row operations. The shift parameter of SPPMI is set to $\kappa$=5 on all training corpora. Sample size $k$ for our method WEQ are chosen from [10, 100000] as shown in Fig. 2. We factorize the Q-contexts matrix $\tilde{C}$ to get singular vectors and singular values. According to the official guide, we use 15 iterations to train GloVe and word2vec methods and 5 iterations to train fasttext method. The number of dimensions of embedding is set to 300. Note that we set the dimension of word embedding to sample size $k$ when $k$ is smaller than 300 in ablation study of sample size. All experiments are carried out on a cloud server with Intel(R) Xeon(R) Gold 6231C CPU. We set the number of CPU cores to 20, 20, 10, 1 and 1 for fasttext, word2vec, GloVe, MF and WEQ, respectively.

### 5.3 Evaluation Tasks

To measure the quality of embeddings, we conduct experiments on three different evaluation tasks: Word Similarity, Text Classification and Named Entity Recognition (NER).

**Word similarity.** The word similarity task is to measure how well the semantic relationship between words is captured by word embeddings. Experiments are conducted on MEN [5] and WS353 [1] datasets. To estimate the quality of word embeddings, we compute the Spearman Rank Correlation score between the word embedding cosine similarities and human-annotated scores.

**Text classification.** In the text classification task, the classifier predicts predefined labels for the given texts. We conduct experiments on two datasets, i.e., 5AbstractsGroup (5AG) [20] and Stanford Sentiment Treebank (SST-2) [30].We choose TextCNN [18] as our classifier, and measure the performance by weighted F1 metrics. The sizes of the kernels are 2, 3 and 4 respectively. We use the Adam optimizer with a mini-batch size of 128 and learning rate of 0.001.

**Named entity recognition.** NER is the task of identifying and categorizing the entities in the given text. We conduct experiments on two benchmark datasets: Conll [35] and BTC [11]. We choose wordLSTM+charCNN+CRF [21] as our NER model and measure the performance by entity-level F1 metrics. The wordLSTM is built upon a one-layer bidirectional LSTM structure with 200 hidden size. We employ the SGD algorithm to optimize model parameters on mini-batches of size 100, with a learning rate of 0.001.

### 5.4 Results

Table 3 illustrates the results on various evaluation tasks and training times of all word embedding methods. It can be observed that our method WEQ outperforms all baselines in three corpora on training time (and resources) metric while the performance metrics are fairly close.

In terms of evaluation tasks, we notice that MF and WEQ are superior to GloVe, word2vec and fasttext in most word similarity tasks, whereas performing slightly worse on text classification and NER tasks. For both NLP tasks (Text classification and NER), models using pre-trained word embeddings have remarkable improvement compared to the one using random embeddings. Comparing with MF, results obtained with WEQ are fairly close to the ones obtained with MF and sometimes even better, despite the fact we only use a fraction of the original information matrix.

We also observe that the released GloVe embedding is about 3% better than word embeddings trained by us on Conll task. One possible reason is the mismatch of vocabulary [21]. We reuse the data pre-process script of word2vec excluding punctuations and digits. And the punctuations and digits are important in the Conll task.

To verify the stability of our method WEQ, we run 5 times for each embedding. We only compute the mean and variance of WEQ because other methods need a long time to train word embedding. The small variance shown in Table 3 indicate that the word embedding obtained by WEQ are stable.

With regard to training time, MF method innately has an advantage over neural network based methods (word2vec, GloVe and fasttext) given the same computational settings. Our method WEQ improves significantly on top of that. As can be seen from the Table 3, our method WEQ-PPMI only requires 54.38% (on enwik9), 15.32% (on WebBase) and 7.72% (on CC) of the training time of MF-PPMI respectively. With the CC corpus, while GloVe needs 270.5 minutes with 10 CPU cores to obtain word embeddings, it takes only 31.4 minutes with 1 CPU core with WEQ-PPMI. Although, as the vocabulary size increases, the training time of WEQ increases

---

Table 3: The performance of WEQ and all baselines on word similarity, text classification and NER tasks. $k$ = the sample size in WEQ. By the 'random' method, we randomly initialize the embedding of each word. GloVe[†] is publicly released word embedding pre-trained with 6B tokens. Bold scores are best within groups of baselines and WEQ.

| Dataset | Method | Time / mins ↓ (Cc: CPU core) | Word similarity ↑ | | Text classification ↑ | | NER ↑ | |
|---|---|---|---|---|---|---|---|---|
| | | | MEN | WS353 | SST-2 | 5AG | Conll | BTC |
| enwik9 (0.1B) | MF-PPMI | 27.4 × 1 Cc | **74.24** | **70.31** | 78.61 | 86.63 | 87.72 | 71.06 |
| | WEQ-PPMI ($k$=40000) | 14.9 × 1 Cc | 73.19±0.28 | 69.23±0.41 | 78.73±0.57 | **87.25**±0.29 | **87.74**±0.11 | **71.16**±0.49 |
| | MF-SPPMI | 13.2 × 1 Cc | 72.26 | 67.14 | 78.54 | 86.19 | 86.13 | 69.88 |
| | WEQ-SPPMI ($k$=40000) | **3.9 × 1 Cc** | 71.92±0.58 | 66.52±1.60 | 78.15±0.20 | 87.18±0.30 | 86.34±0.12 | 69.70±0.12 |
| | GloVe | 44.6 × 10 Cc | 69.45 | 66.31 | 79.26 | 86.51 | 86.78 | 71.00 |
| | word2vec | 59.9 × 20 Cc | 74.09 | 68.96 | 79.13 | 85.62 | 86.51 | 70.19 |
| | fasttext | 67.3 × 20 Cc | 73.13 | 65.79 | **80.59** | 86.85 | 87.63 | 70.65 |
| WebBase (3B) | MF-PPMI | 185.4 × 1 Cc | **77.97** | **68.77** | 80.94 | 86.32 | 88.21 | 71.31 |
| | WEQ-PPMI ($k$=50000) | 28.4 × 1 Cc | 77.01±0.44 | 68.74±0.84 | 81.59±0.14 | 87.17±0.22 | 88.48±0.17 | 71.37±0.49 |
| | MF-SPPMI | 67.6 × 1 Cc | 76.84 | 67.62 | 79.13 | 85.34 | 87.07 | 69.31 |
| | WEQ-SPPMI ($k$=50000) | **7.3 × 1 Cc** | 76.31±0.55 | 67.66±0.55 | 78.98±0.28 | 86.86±0.23 | 87.68±0.12 | 70.44±0.24 |
| | GloVe | 182.7×10 Cc | 75.20 | 63.43 | **82.66** | 86.90 | 87.78 | 71.02 |
| | word2vec | 1211.4 × 20 Cc | 74.73 | 63.98 | 80.93 | 87.64 | 88.58 | 69.89 |
| | fasttext | 1843.7 × 20 Cc | 73.10 | 58.80 | 81.07 | **87.74** | **88.93** | **72.83** |
| CC (6B) | MF-PPMI | 406.5 × 1 Cc | 77.72 | **70.87** | **83.99** | 86.85 | 87.48 | 72.55 |
| | WEQ-PPMI ($k$=50000) | 31.4 × 1 Cc | 76.39±0.54 | 69.54±0.54 | 82.90±0.74 | **87.69**±0.15 | 88.59±0.25 | 72.36±0.36 |
| | MF-SPPMI | 153.9 × 1 Cc | 75.30 | 69.35 | 80.25 | 85.92 | 87.95 | 70.37 |
| | WEQ-SPPMI ($k$=50000) | **13.3 × 1 Cc** | 73.92±0.35 | 66.15±0.96 | 80.79±0.66 | 86.36±0.17 | 87.78±0.18 | 70.86±0.56 |
| | GloVe | 270.5 × 10 Cc | 76.37 | 65.41 | 81.70 | 86.47 | 88.56 | 71.23 |
| | word2vec | 2369.7 × 20 Cc | 79.29 | 69.90 | 81.57 | 87.14 | 88.65 | 72.93 |
| | fasttext | 3419.4 × 20 Cc | **79.64** | 68.97 | 82.70 | 87.53 | **89.46** | **73.92** |
| | random | / | 0.31 | 2.38 | 73.82 | 82.51 | 81.63 | 59.22 |
| | GloVe[†] | / | 73.75 | 57.26 | 82.79 | 87.34 | 91.15 | 72.26 |

accordingly, the time taken by MF increases by a significantly larger margin.

Note that in the case closest to practice (6B tokens and 400K vocabularies), we achieve at least 11 ∼ 13 times speed-up compared with other baselines. In more realistic scenarios where the number of tokens is trillions and the size of vocabulary is several million, one can infer from the growing trend that *the speed advantage of our algorithm could potentially reach several orders of magnitude.*

The analysis presented above confirms that WEQ method can substantially minimize the training time while maintaining benchmark performance.

## 5.5 Discussion and Analysis

**Time analysis.** Recall three main steps of the WEQ method: information matrix computation, Q-contexts construction (including state preparation and sampling) and calculating word embedding from Q-contexts (including computation of singular values and left singular vectors, and embedding computation). The breakdown of computational time is displayed in Table.4. Note that the Q-contexts construction step takes up only no more than 11% of the whole time, which shows the efficiency of Q-contexts construction step.

Table 4: The training time decomposition of the different part of WEQ methods on the CC corpus. Total training time of WEQ-PPMI($k$=50000) and WEQ-SPPMI($k$=50000) are 31.4 minutes and 13.3 minutes respectively.

| | Info. matrix computation | State preparation | Sampling | Sing. val. vec. | Embedding calculation |
|---|---|---|---|---|---|
| WEQ-PPMI | 26.37% | 3.74% | 6.93% | 46.40% | 16.56% |
| WEQ-SPPMI | 52.93% | 3.76% | 3.73% | 30.92% | 8.66% |

**Sample size.** In our method, Q-contexts matrix is obtained by sampling from the information matrix. To investigate how the sample size affects the performance of WEQ, we carry out experiments with different sample sizes on all three corpora as shown in Fig.2. We present a representative benchmark score, similarity task MEN, against the training times under corresponding sample sizes. The results of other evaluation tasks are shown in Fig.3. As the sample size increases, it is obvious that the performance of WEQ on MEN approaches that of MF. The convergence is fairly fast, especially when the corpus size is large. We can see that to reach at least 98% of the performance of MF, the sample size of WEQ method
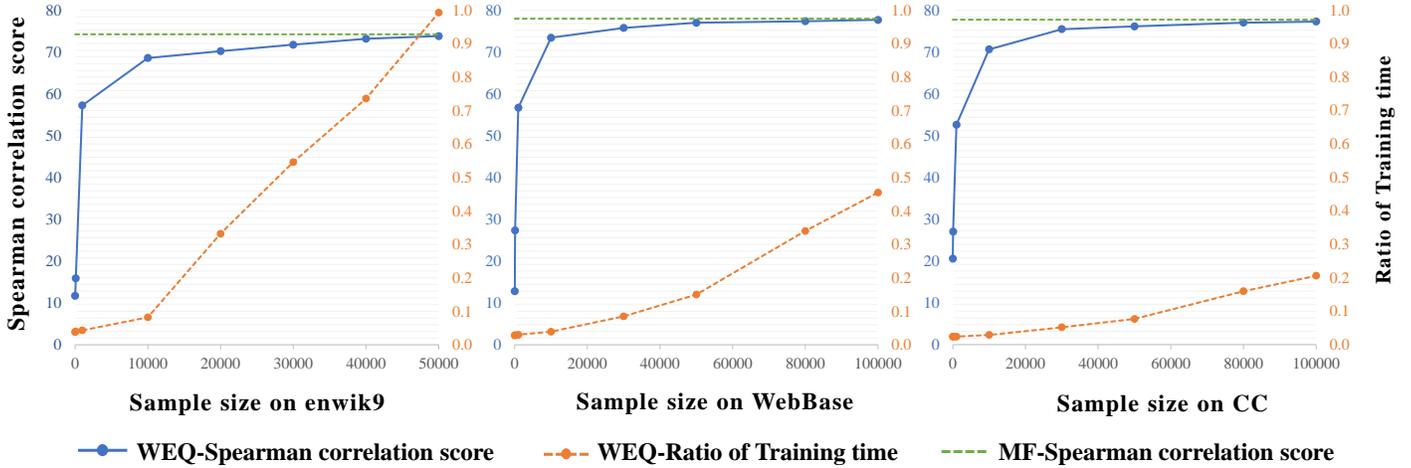
**Figure 2: The scores of similar task (MEN) and the corresponding training times of WEQ-PPMI method under different sample sizes. The original dimensions of the information matrices are about 56K (on enwik9), 278K (on WebBase) and 400K (on CC) respectively. Ratio of Training time = training time of WEQ / training time of MF.**

only needs to be 40k (on enwik9), 50k (on WebBase), 50k (on CC), which are about 70.83%, 18.00%, 12.50% of the original size of the information matrices, respectively. This validates the theoretical observations in Section 4.1. As the size of the original information matrix increases, the percentage of samples needed decreases. In terms of training time, it rises slowly at small sample sizes then trends upward sharply at larger sample sizes. This indicates that the training time can be effectively reduced when sampling a small matrix to get word embedding.
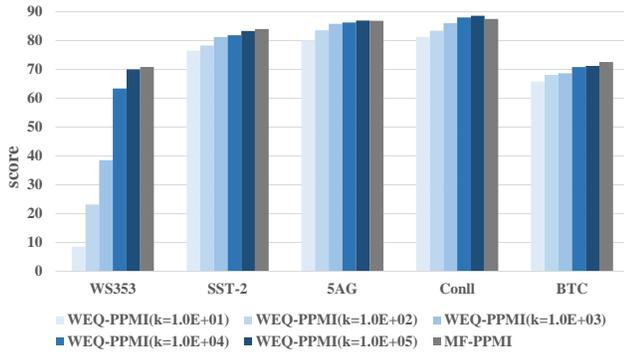


**Figure 3: The scores of all evaluation tasks except MEN under different sample sizes on CC corpus.**

$\ell^2$**-norm sampling vs uniform sampling.** In WEQ method, we use $\ell^2$-norm sampling to capture the typical information of origin information matrix. To verify the effectiveness of the $\ell^2$-norm sampling, we replace the $\ell^2$-norm sampling with uniform sampling in the Q-contexts construction step. Uniform sampling means that each row of information matrix $M$ will be sampled with equal probability. We compare the performance of word embeddings trained by these two different sampling algorithms. The evaluation

results are shown in Fig. 4. $\ell^2$-norm sampling has great advantages over uniform sampling in terms of word similarity scores. In terms of downstream tasks (text classification, NER), $\ell^2$-norm sampling has 1.5% ~ 4% increase compared with uniform sampling. These experimental results are consistent with our theoretical proof in Section 4.
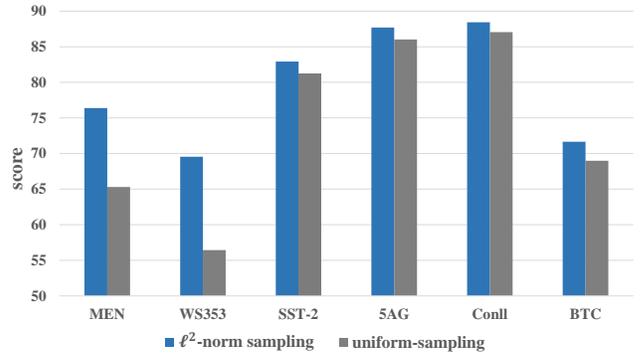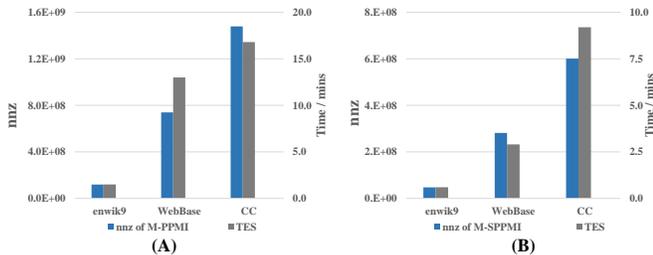


**Figure 4: Comparison between $\ell^2$-norm sampling and uniform sampling.**

**Nonzeros/Sparsity analysis.** To investigate the correlation between matrix sparsity and algorithm efficiency, we show a comparison between time and number of nonzeros in an information matrix $M$ in Fig. 5. By TES, we mean total time excluding SVD computation, which corresponds to Algorithm 3 excluding line 3. Since the time for SVD computation depends on the SVD engine used in our method, we do not provide a detailed analysis. The Pearson Correlation Score between $nnz(M)$ and TES is 0.9583. It demonstrates that our algorithm efficiency is approximately linearly and positively correlated with the number of nonzero elements of $M$. We also observe that number of nonzero elements in SPPMI matrix

is significantly less than that in PPMI. That is because only the positive values higher than $\kappa$ are retained after the shifted operation. The decrease of number of nonzero elements leads to a more sparse matrix and shorter TES.



Figure 5: Comparison between nnz($M$) and TES on three different training corpora. nnz($M$) means number of nonzero elements in information matrix $M$. TES means algorithm training Time Excluding SVD. The Pearson Correlation Score is 0.9583. Panel A is the result for the PPMI information matrix, and Panel B is the result for the SPPMI matrix.

## 6 RELATED WORK

In this section, we review the related work of word embedding and low-rank approximation.

### 6.1 Word Embedding

Constructing the word embedding which could express the semantic features is a fundamental task in Natural Language Processing. It brings benefits to many NLP tasks [17, 22]. Briefly, recent work about word embedding can be categorized into two genres, i.e., neural network based methods [4, 24, 27] and global matrix factorization based methods [2, 10, 19]. It is found that we can bridge these two categories by unifying neural network based methods into a matrix factorization framework [2, 19].

More specifically, most of the neural network based methods model the relation between the target word and its contextual words. Among them, the SGNS (Skip-Gram with Negative Sampling) model is the popular implementation of word2vec [24]. It separates local context windows on the whole corpus, and focuses on maximizing the likelihood of contextual words based on the given word. Fasttext [4] is an extension of the word2vec model, which represents each word as an n-gram of characters by a sliding window. GloVe [27] is trained based on the global word-word co-occurrence counts from a corpus. As for one kind of global matrix factorization based method, MF [19] factorizes the SPPMI matrix explicitly, rather than iteratively tuning the network parameters.

More recently, ELMo [28] and BERT [12] achieve state-of-the-art results in a variety of NLP tasks. Despite the state-of-the-art results achieved by deep contextualization models, pre-trained word embedding should not be neglected. Indeed, ELMo is dependent on GloVe embedding as input during training [28]. Moreover, both ELMo and BERT are extremely large deep networks that require huge computing resources. The pre-training of BERT$_{large}$ takes 4 days to complete on 64 TPUs [12]. And it is also difficult to apply

deep contextualization models directly on low-resource PCs or mobile phones [32, 33].

In the studies described above, researchers paid less attention to the computational cost of word embedding training. It is true that large-scale word embedding training requires significant computational costs [23–25, 27]. In this work, we aim to present a novel efficient method to obtain word embeddings of a large-scale vocabulary, while maintaining its superiority in terms of effectiveness.

### 6.2 Low-Rank Approximation

Low-rank approximation problems in mathematical modeling have been studied for decades [9, 14], which is approximating a matrix by one whose rank is less than that of the original matrix. The aim is to obtain more compact information representation and less complex data modeling with limited losses. It arises in many applications such as recommendation system [3, 34]. Unlike their implementations of random matrices and matrices of small sizes [3], our WEQ investigates the actual large matrices associated with natural language and reveals new relations for word embeddings. Our work is also different from [14, 34] since we provide more detailed theoretical and empirical analysis on matrix concentration, and direct analysis of singular vectors. We also have novel treatments for sparse matrices and better algorithmic designs for state preparation.

## 7 CONCLUSION

In this work, we introduce the notion of Q-contexts (matrix), which can be constructed efficiently. These are only a small fraction (less than 12.5% in the practical scenario) of all the contexts in the entire corpus. We also present a novel relation between word vectors and Q-contexts, and provide a theoretical foundation and rigorous analysis. Based on this relation, we design a novel and efficient WEQ method for fast computation of large-scale word embedding. Empirical experiments show that our algorithm WEQ method runs at least $11 \sim 13$ times faster than well-established matrix factorization methods. Resource and time advantages over word2vec, GloVe and fasttext are even more pronounced in our empirical study. We have also shown that WEQ enjoys decent accuracy performance in a variety of NLP tasks compared to the other methods tested in this study. In the future, we would like to efficiently learn word embedding of million-level vocabulary. It is also an interesting topic to apply our method to other domains such as graph embedding.

## REFERENCES

[1] Eneko Agirre, Enrique Alfonseca, Keith Hall, Jana Kravalova, Marius Paşca, and Aitor Soroa. 2009. A Study on Similarity and Relatedness Using Distributional and WordNet-based Approaches. In *NAACL '09*.
[2] Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. 2016. A latent variable model approach to pmi-based word embeddings. *Transactions of the Association for Computational Linguistics* 4 (2016), 385–399.
[3] Juan Miguel Arrazola, Alain Delgado, Bhaskar Roy Bardhan, and Seth Lloyd. 2020. Quantum-inspired algorithms in practice. *Quantum* (2020).

[4] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics* 5 (2017), 135–146.

[5] Elia Bruni, Nam-Khanh Tran, and Marco Baroni. 2014. Multimodal distributional semantics. *Journal of artificial intelligence research* 49 (2014), 1–47.

[6] Kenneth Ward Church and Patrick Hanks. 1989. Word Association Norms, Mutual Information, and Lexicography. In *ACL '89*.

[7] Bob Coecke, Giovanni de Felice, Konstantinos Meichanetzidis, and Alexis Toumi. 2020. Foundations for Near-Term Quantum Natural Language Processing. *arXiv preprint arXiv:2012.03755* (2020).

[8] Ido Dagan, Fernando Pereira, and Lillian Lee. 1994. Similarity-Based Estimation of Word Cooccurrence Probabilities. In *ACL '94*.

[9] Yogesh Dahiya, Dimitris Konomis, and David P. Woodruff. 2018. An Empirical Evaluation of Sketching for Numerical Linear Algebra. In *KDD '18*.

[10] Scott Deerwester, Susan T Dumais, George W Furnas, Thomas K Landauer, and Richard Harshman. 1990. Indexing by latent semantic analysis. *Journal of the American society for information science* 41, 6 (1990), 391–407.

[11] Leon Derczynski, Kalina Bontcheva, and Ian Roberts. 2016. Broad twitter corpus: A diverse named entity recognition resource. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*. 1169–1179.

[12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL '19*. 4171–4186.

[13] Petros Drineas, Michael W. Mahoney, and S. Muthukrishnan. 2008. Relative-Error CUR Matrix Decompositions. *SIAM J. Matrix Anal. Appl.* 30, 2 (2008), 844–881. https://doi.org/10.1137/07070471X

[14] Alan Frieze, Ravi Kannan, and Santosh Vempala. 2004. Fast Monte-Carlo Algorithms for Finding Low-Rank Approximations. *J. ACM* (2004).

[15] Lushan Han, Abhay L. Kashyap, Tim Finin, James Mayfield, and Jonathan Weese. 2013. UMBC_EBIQUITY-CORE: Semantic Textual Similarity Systems. In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity*.

[16] Elad Hazan, Tomer Koren, and Nati Srebro. 2011. Beating SGD: Learning SVMs in Sublinear Time. In *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain*, John Shawe-Taylor, Richard S. Zemel, Peter L. Bartlett, Fernando C. N. Pereira, and Kilian Q. Weinberger (Eds.). 1233–1241.

[17] Elham Khabiri, Wesley M Gifford, Bhanukiran Vinzamuri, Dhaval Patel, and Pietro Mazzoleni. 2019. Industry specific word embedding and its application in log classification. In *CIKM '19*. 2713–2721.

[18] Yoon Kim. 2014. Convolutional Neural Networks for Sentence Classification. In *EMNLP '14*.

[19] Omer Levy and Yoav Goldberg. 2014. Neural Word Embedding as Implicit Matrix Factorization. In *NeurIPS '14*.

[20] Qian Liu, Heyan Huang, Yang Gao, Xiaochi Wei, Yuxin Tian, and Luyang Liu. 2018. Task-oriented Word Embedding for Text Classification. In *Proceedings of the 27th International Conference on Computational Linguistics*, Santa Fe, New Mexico, USA, 2023–2032.

[21] Xuezhe Ma and Eduard Hovy. 2016. End-to-end Sequence Labeling via Bidirectional LSTM-CNNs-CRF. In *ACL '16*.

[22] Sameen Maruf and Gholamreza Haffari. 2018. Document Context Neural Machine Translation with Memory Networks. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Melbourne, Australia, 1275–1284.

[23] Tomas Mikolov, Kai Chen, Greg S. Corrado, and Jeffrey Dean. 2013. Efficient Estimation of Word Representations in Vector Space. *ICLR Workshop'13* (2013).

[24] Tomás Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. Distributed Representations of Words and Phrases and their Compositionality. In *NeurIPS '13*, Christopher J. C. Burges, Léon Bottou, Zoubin Ghahramani, and Kilian Q. Weinberger (Eds.).

[25] Andriy Mnih and Koray Kavukcuoglu. 2013. Learning word embeddings efficiently with noise-contrastive estimation. In *NeurIPS '13*, C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger (Eds.).

[26] Pedro Javier Ortiz Suárez, Laurent Romary, and Benoît Sagot. 2020. A Monolingual Approach to Contextualized Word Embeddings for Mid-Resource Languages. In *ACL '20*.

[27] Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *EMNLP '14*. 1532–1543.

[28] Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep Contextualized Word Representations. In *Proceedings of NAACL-HLT*. Association for Computational Linguistics, New Orleans, Louisiana, 2227–2237.

[29] Patrick Rebentrost, Masoud Mohseni, and Seth Lloyd. 2014. Quantum support vector machine for big data classification. *Physical review letters* 113, 13 (2014), 130503.

[30] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *EMNLP '13*. 1631–1642.

[31] Zhao Song, David P. Woodruff, and Huan Zhang. 2016. Sublinear Time Orthogonal Tensor Decomposition. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (Eds.). 793–801.

[32] Emma Strubell, Ananya Ganesh, and Andrew McCallum. 2019. Energy and Policy Considerations for Deep Learning in NLP. In *ACL '19*. Association for Computational Linguistics, Florence, Italy, 3645–3650.

[33] Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. 2020. MobileBERT: a Compact Task-Agnostic BERT for Resource-Limited Devices. In *ACL '20*. Association for Computational Linguistics, Online, 2158–2170.

[34] Ewin Tang. 2019. A quantum-inspired classical algorithm for recommendation systems. *STOC '19* (2019).

[35] Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the CoNLL-2003 Shared Task: Language-Independent Named Entity Recognition. In *NAACL '03*.

[36] Joel A. Tropp. 2015. An Introduction to Matrix Concentration Inequalities. (2015).

[37] Peter D. Turney and Patrick Pantel. 2010. From Frequency to Meaning: Vector Space Models of Semantics. *J. Artif. Int. Res.* (2010).

[38] William Zeng and Bob Coecke. 2016. Quantum algorithms for compositional natural language processing. *arXiv preprint arXiv:1608.01406* (2016).

[39] Yun Zhang, Yongguo Liu, Jiajing Zhu, Ziqiang Zheng, Xiaofeng Liu, Weiguang Wang, Zijie Chen, and Shuangqing Zhai. 2019. Learning Chinese word embeddings from stroke, structure and pinyin of characters. In *CIKM '19*. 1011–1020.